

# Attention Is All You Need

## 2017



Rak Arina  
HSE Research Seminar  
2022



# Plan

---

- Intro
- Recap: Autoregressive models and attention
- Architecture:
  - Encoder
  - Self-attention
  - Decoder
  - Positional embeddings
- Questions



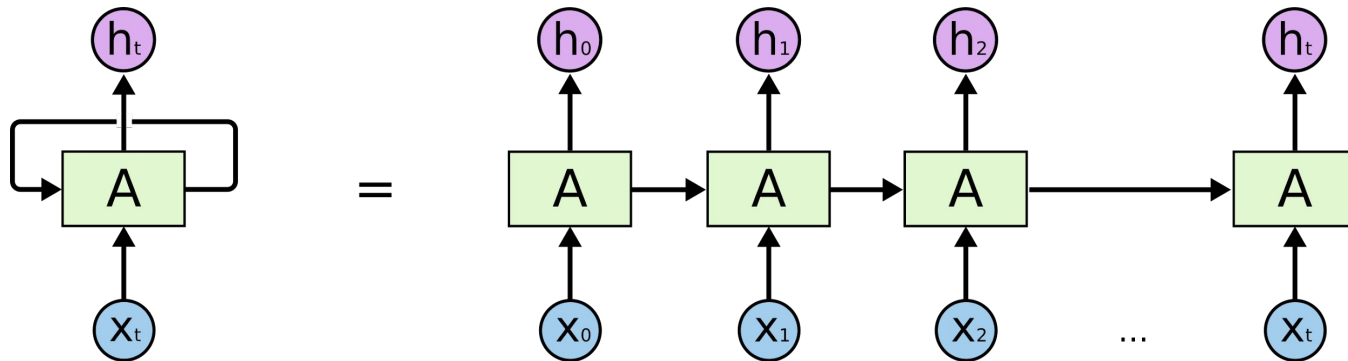
# Task

---

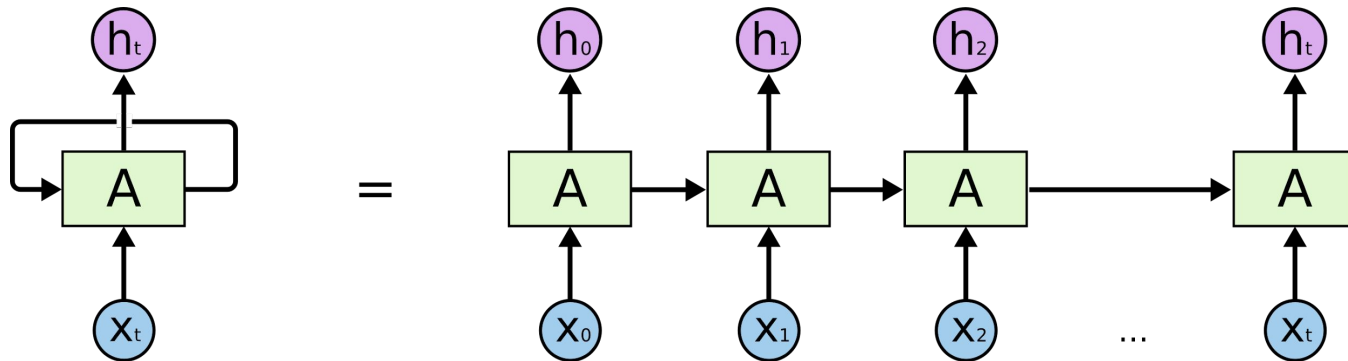
**Sequence to sequence** — transformation of input sequences into output sequences

- Machine translation
- Spelling correction
- Part of speech tagging
- Speech recognition

# Recap: RNN



## Recap: RNN

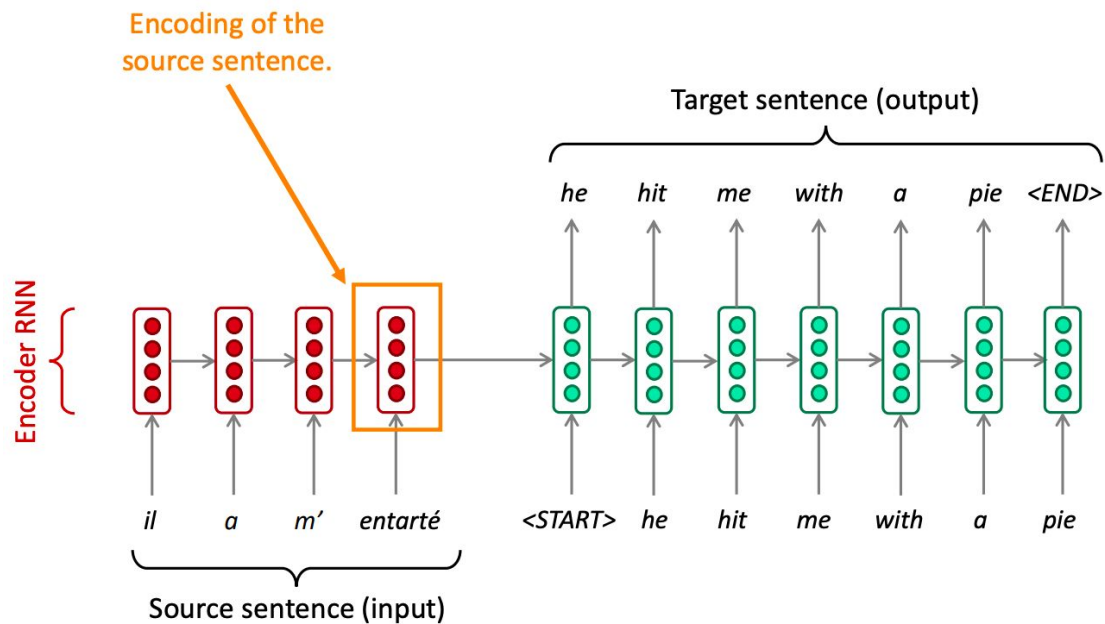


### Problems:

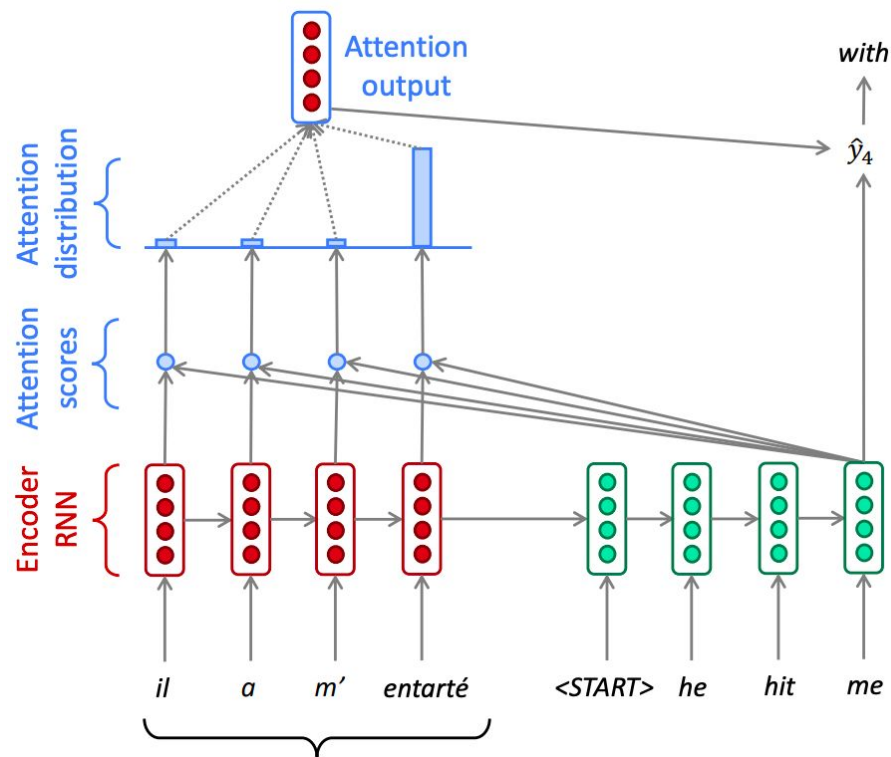
- Learns slow
- Vanishing / exploding gradients
- Catastrophic forgetting



# Recap: RNN



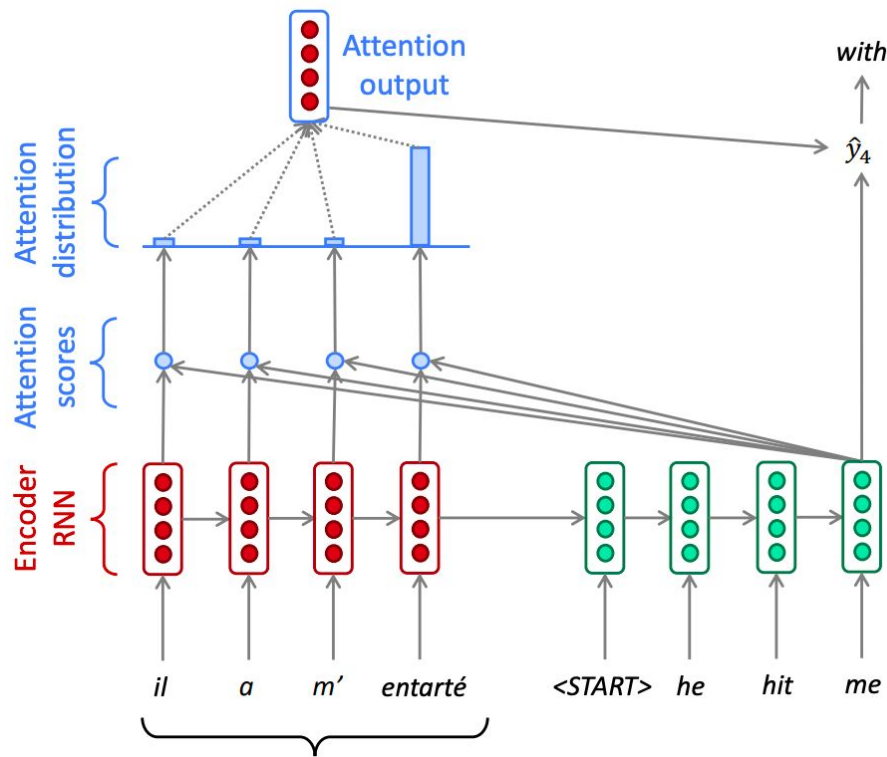
# Recap: Attention



# Recap: Attention

## Attention:

- Improves performance
- Helps with vanishing gradients
- Solves the bottleneck problem
- Helps with interpretability





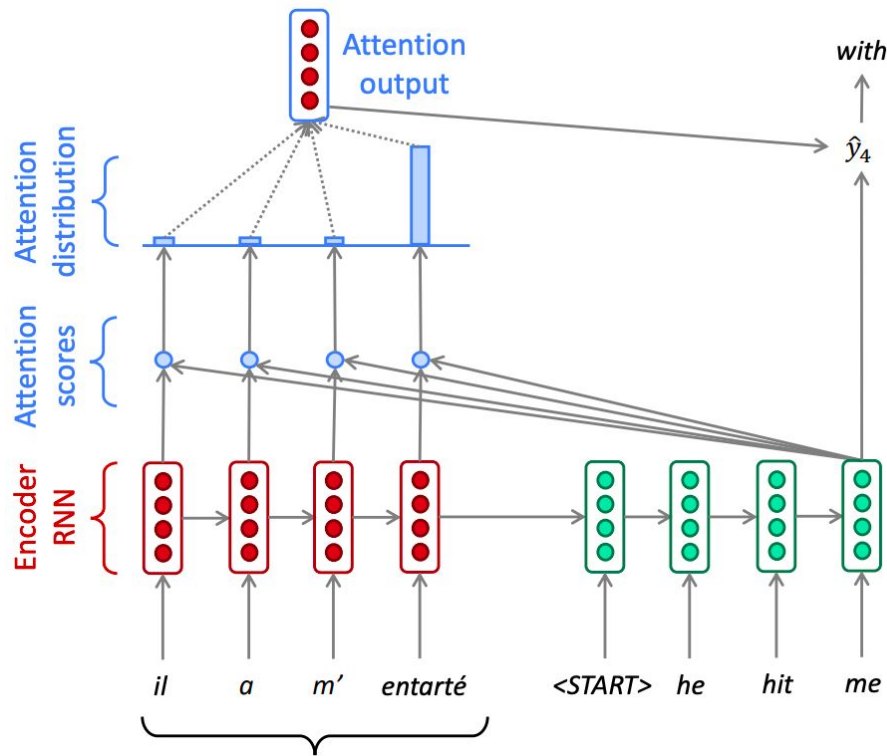
# Recap: Attention

## Attention:

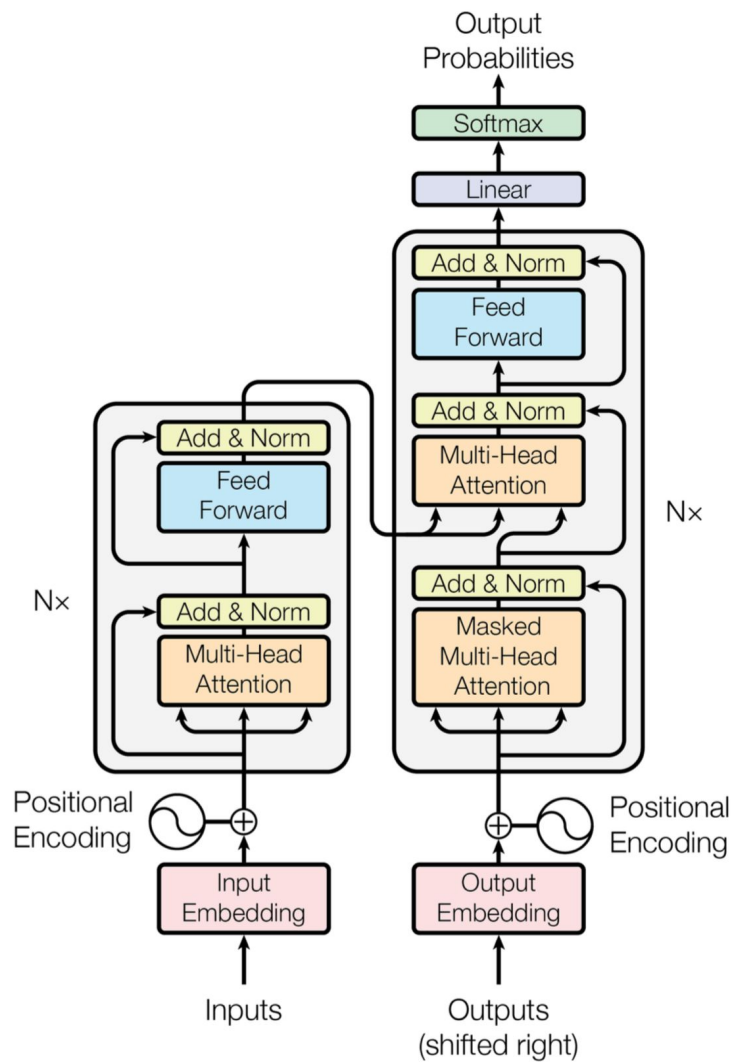
- Improves performance
- Helps with vanishing gradients
- Solves the bottleneck problem
- Helps with interpretability

## BUT:

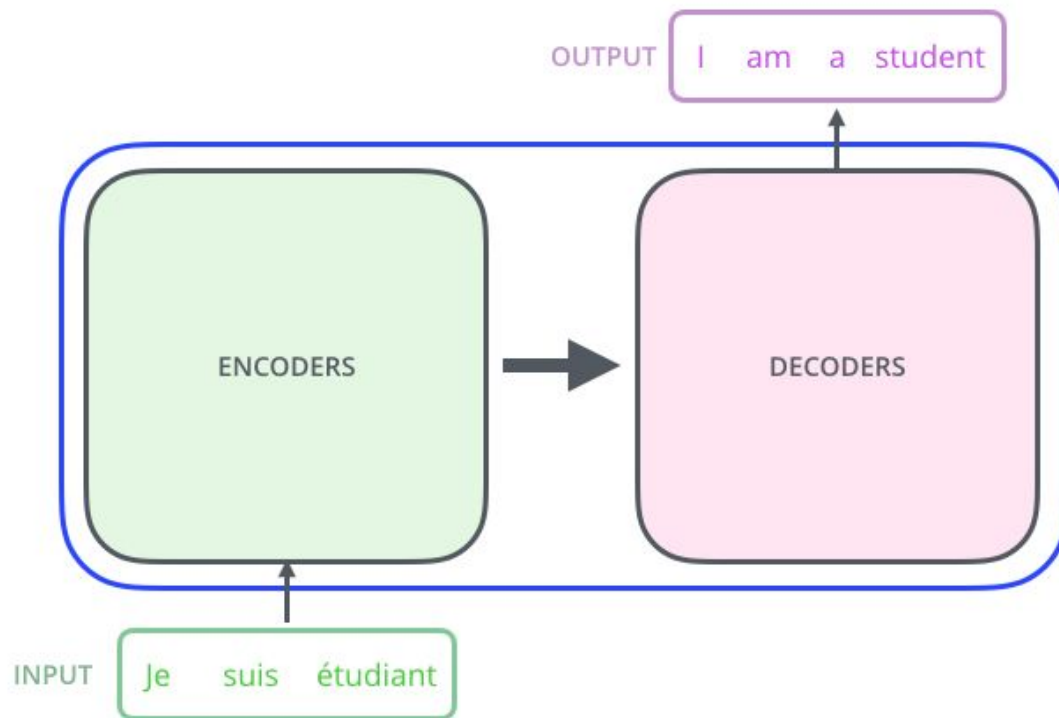
Models get more and more complex and the computations still can not be done in parallel therefore **SLOW**



# Transformer



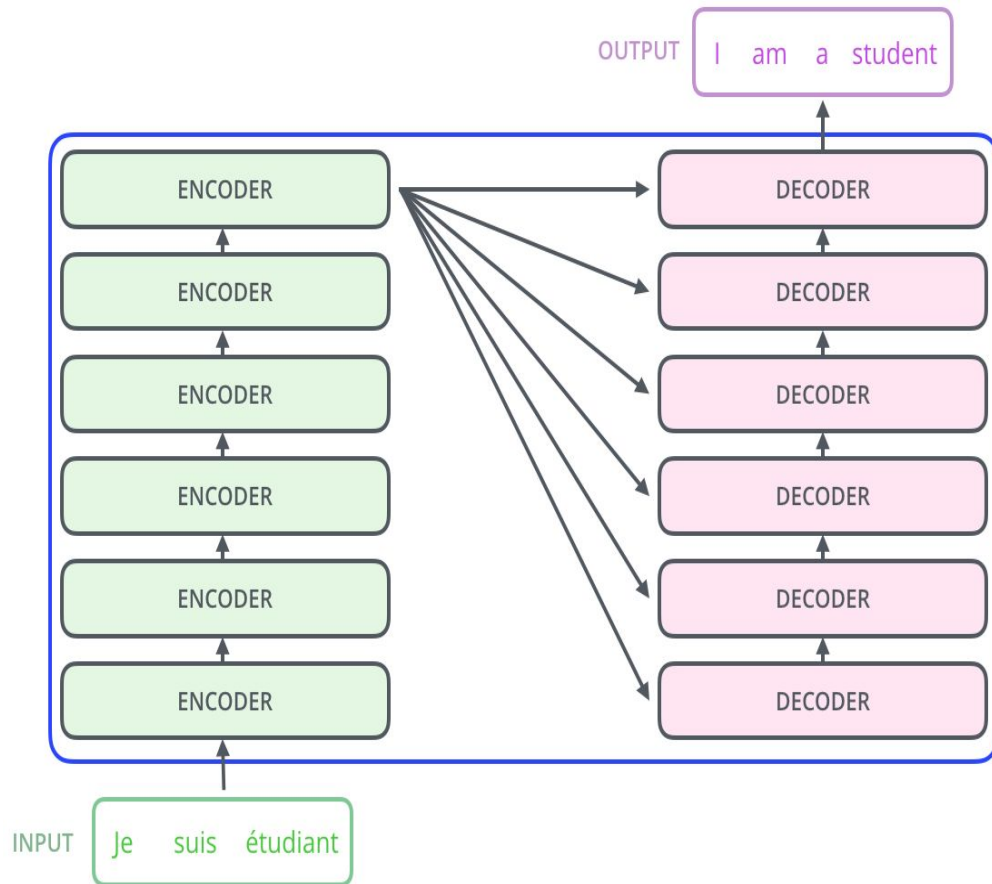
# Architecture



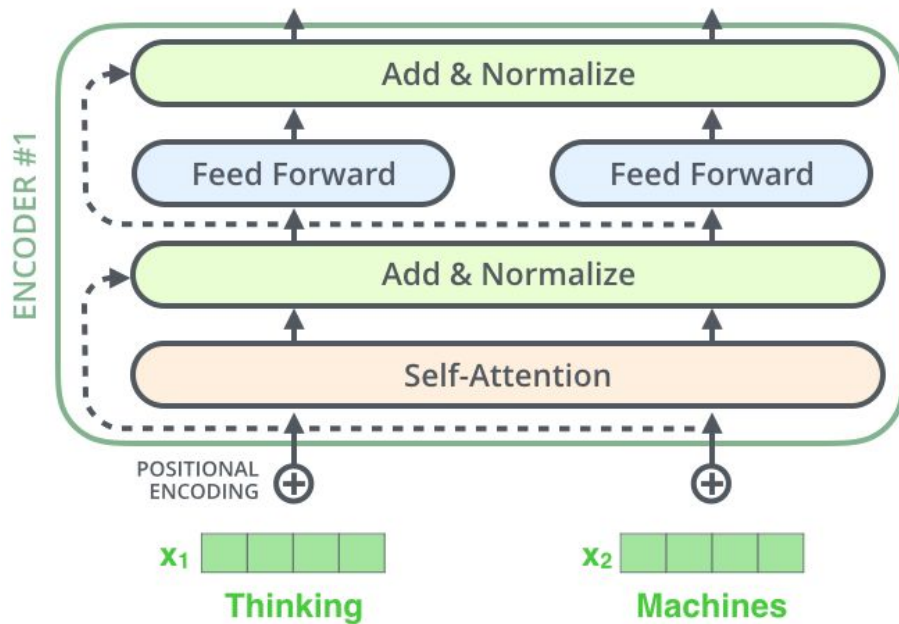
# ● Architecture

Encoder and Decoder:

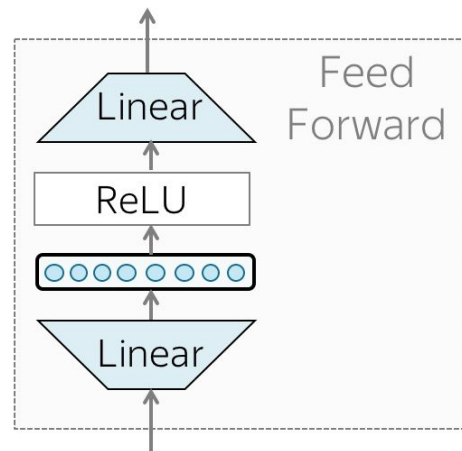
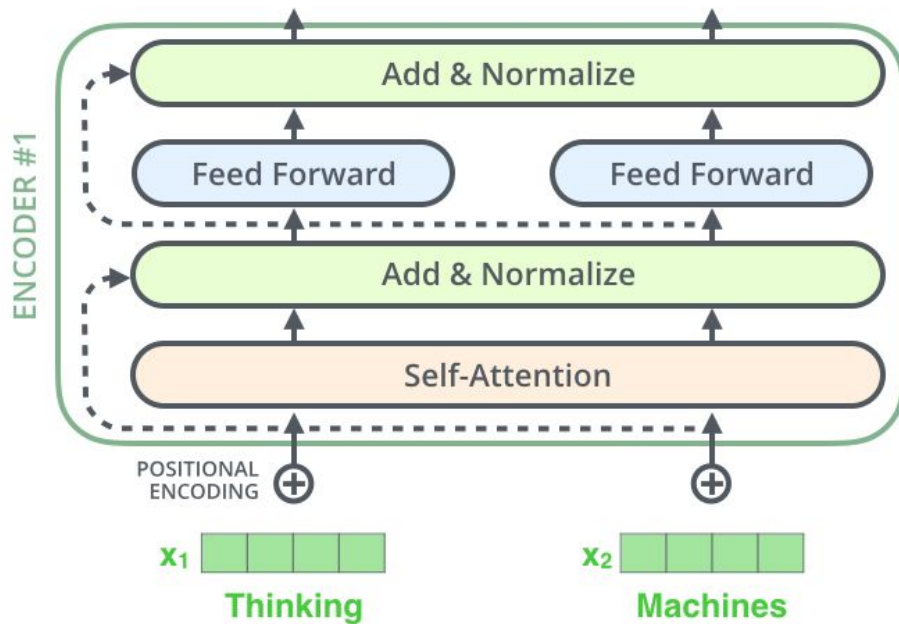
- Consist of blocks
- The blocks have the same architecture
- The blocks do not share weights



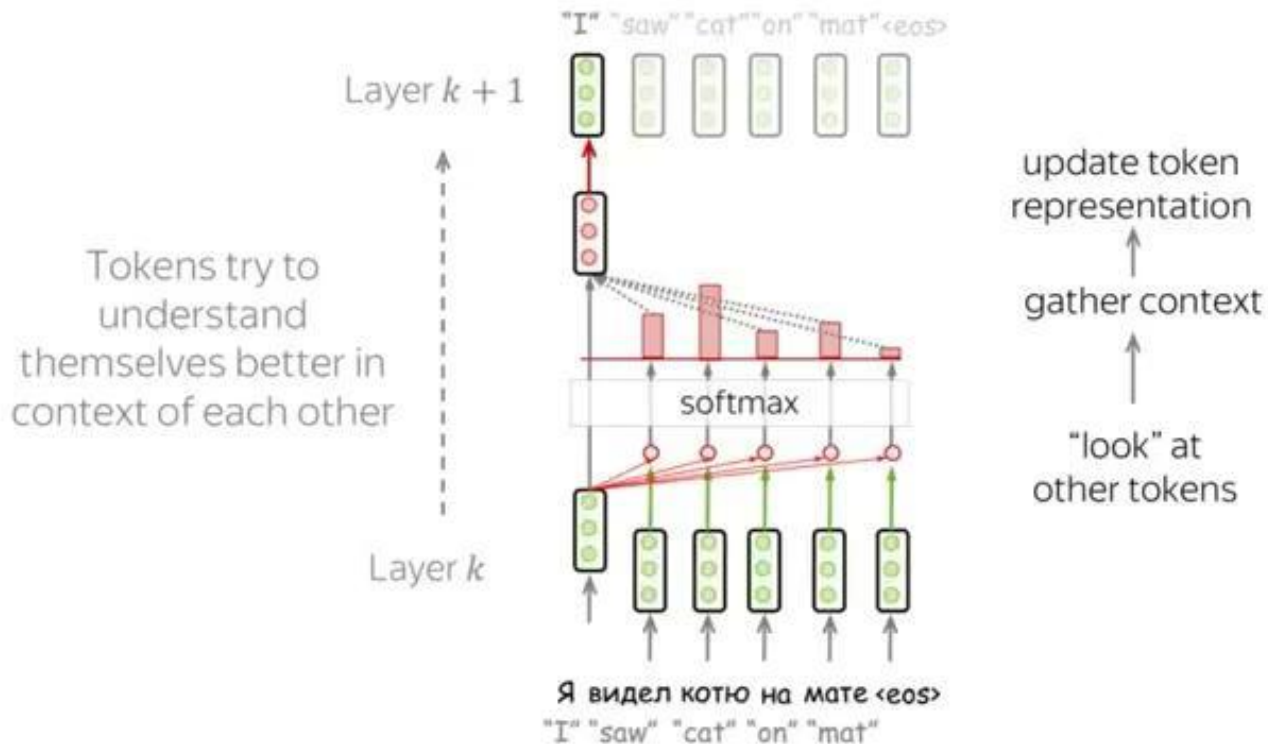
# Encoder



# Encoder



# Self-attention





# Self-attention

Each vector receives three representations (“roles”)

$$\begin{bmatrix} W_Q \end{bmatrix} \times \begin{bmatrix} \text{green} \\ \text{green} \\ \text{green} \end{bmatrix} = \begin{bmatrix} \text{blue} \\ \text{blue} \\ \text{blue} \end{bmatrix}$$

**Query:** vector from which the attention is looking

“Hey there, do you have this information?”

$$\begin{bmatrix} W_K \end{bmatrix} \times \begin{bmatrix} \text{green} \\ \text{green} \\ \text{green} \end{bmatrix} = \begin{bmatrix} \text{yellow} \\ \text{yellow} \\ \text{yellow} \end{bmatrix}$$

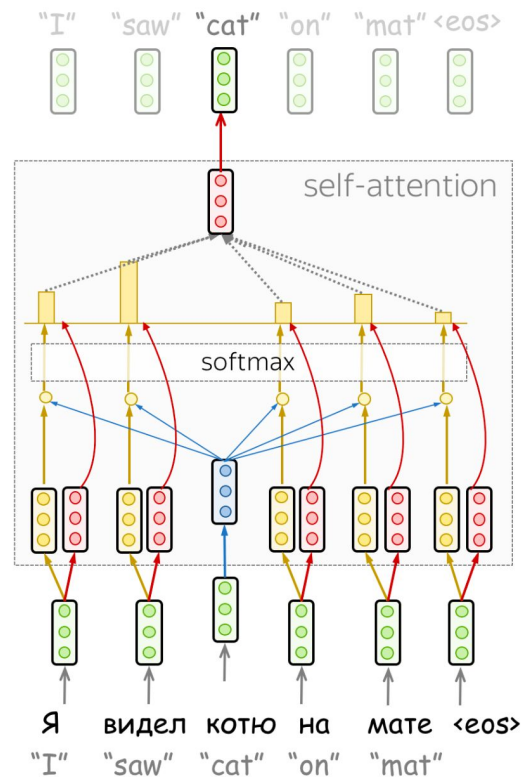
**Key:** vector at which the query looks to compute weights

“Hi, I have this information – give me a large weight!”

$$\begin{bmatrix} W_V \end{bmatrix} \times \begin{bmatrix} \text{green} \\ \text{green} \\ \text{green} \end{bmatrix} = \begin{bmatrix} \text{red} \\ \text{red} \\ \text{red} \end{bmatrix}$$

**Value:** their weighted sum is attention output

“Here’s the information I have!”



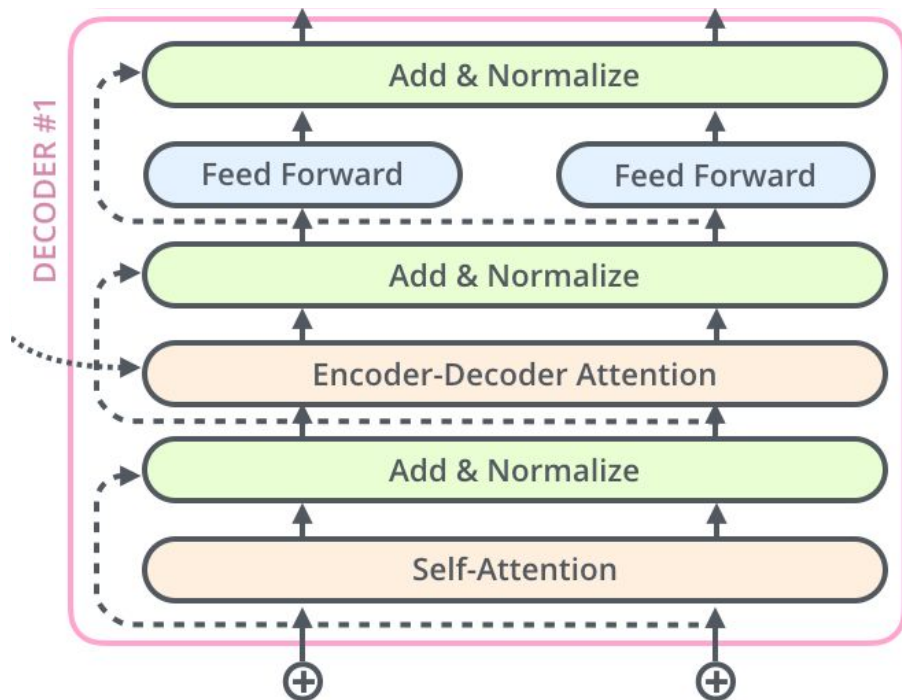


## ● Self-attention

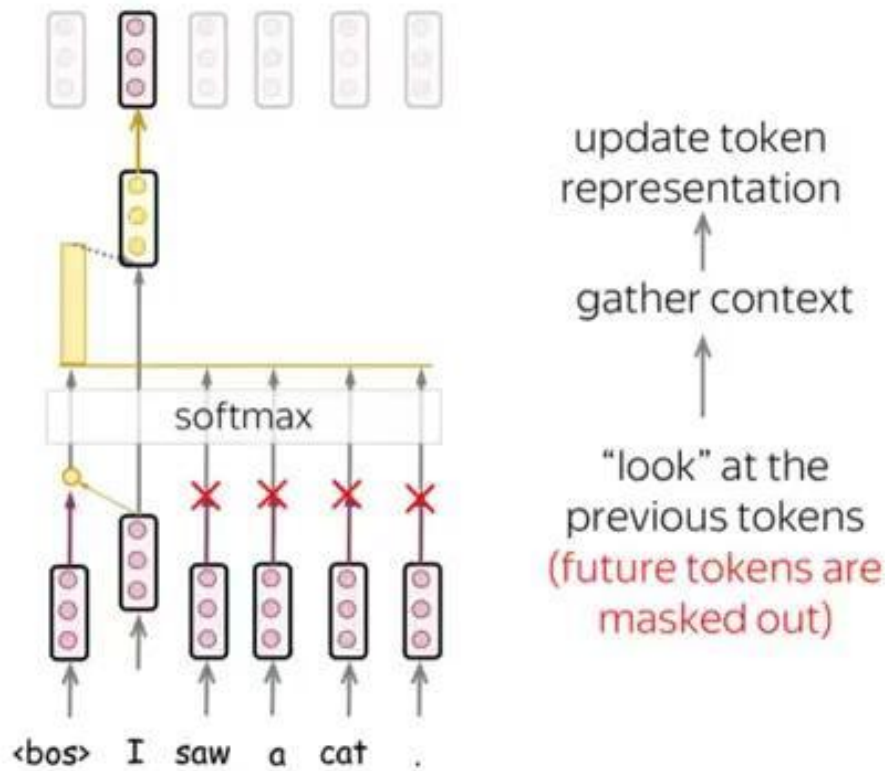
$$\text{Attention}(\underset{\substack{\text{from}}}{q}, \underset{\substack{\text{to}}}{k}, v) = \overbrace{\text{softmax}\left(\frac{qk^T}{\sqrt{d_k}}\right)}^{\text{Attention weights}} v$$

vector dimensionality of K, V

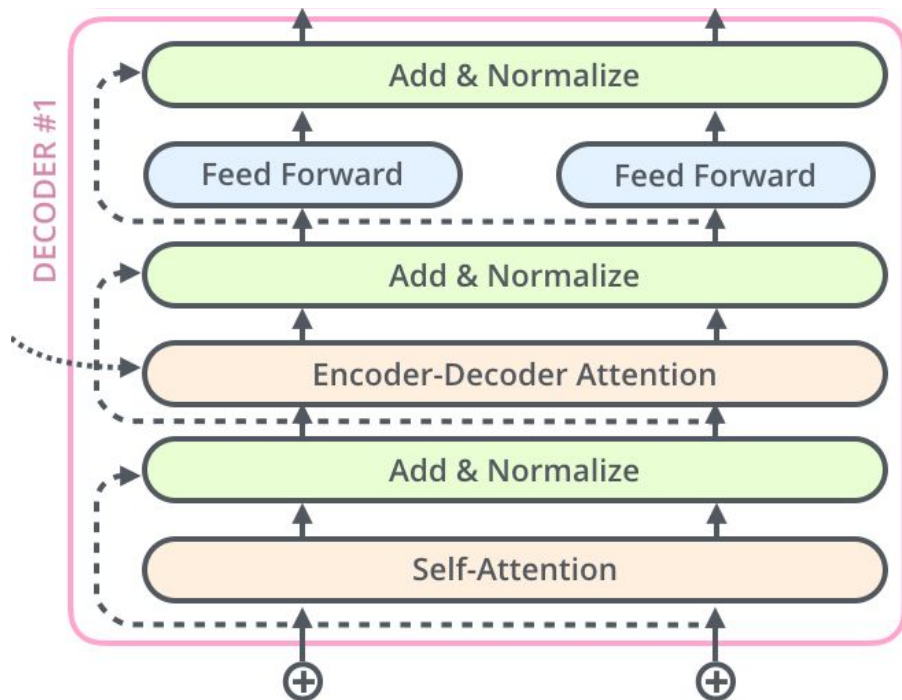
# Decoder



# Decoder self-attention

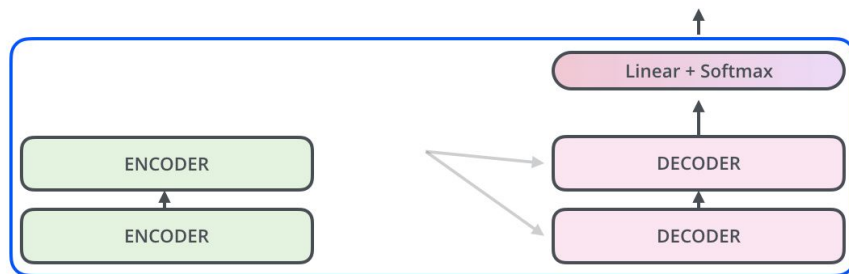


# Decoder



Decoding time step: 1 2 3 4 5 6

OUTPUT



EMBEDDING WITH TIME SIGNAL



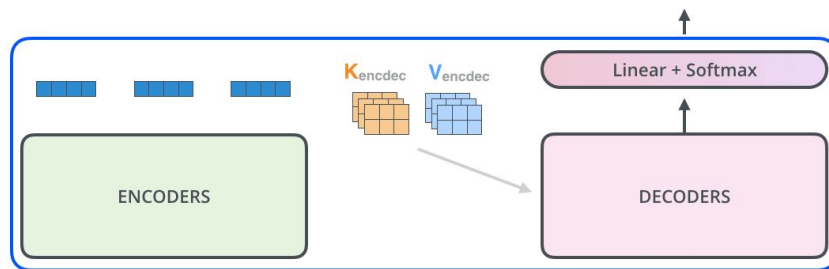
EMBEDDINGS



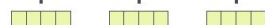
INPUT Je suis étudiant

Decoding time step: 1 2 3 4 5 6

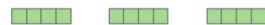
OUTPUT |



EMBEDDING WITH TIME SIGNAL



EMBEDDINGS

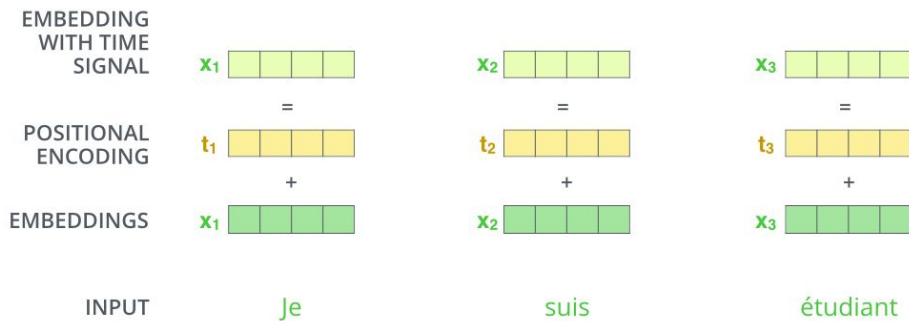


INPUT Je suis étudiant

PREVIOUS OUTPUTS |

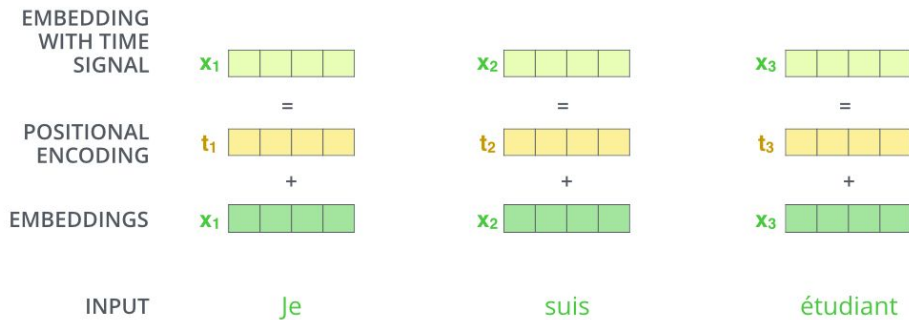


# Positional information





# Positional information

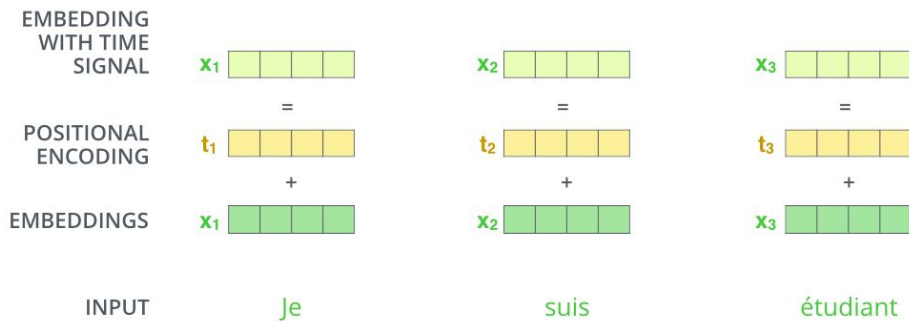


$$PE_{(pos, 2i)} = \sin(pos/10000^{2i/d_{\text{model}}})$$

$$PE_{(pos, 2i+1)} = \cos(pos/10000^{2i/d_{\text{model}}})$$

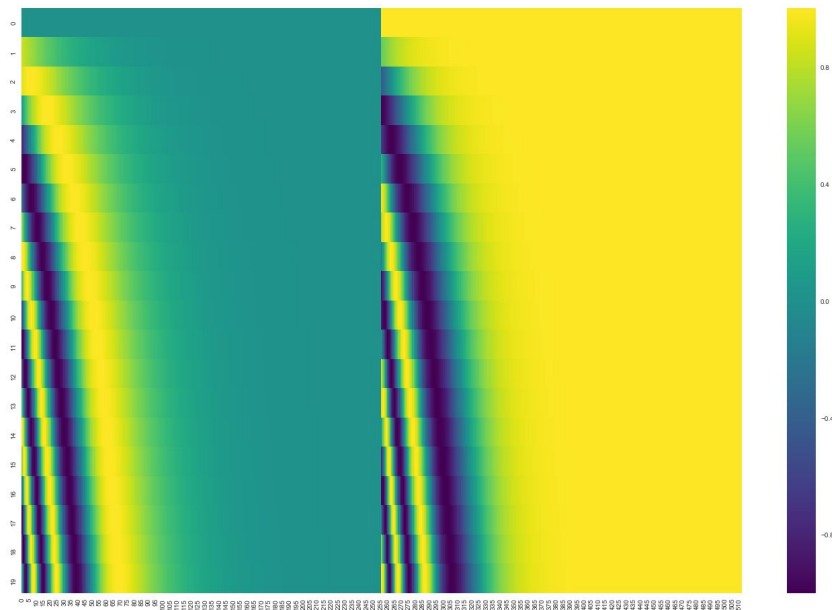


# Positional information



$$PE_{(pos, 2i)} = \sin(pos/10000^{2i/d_{\text{model}}})$$

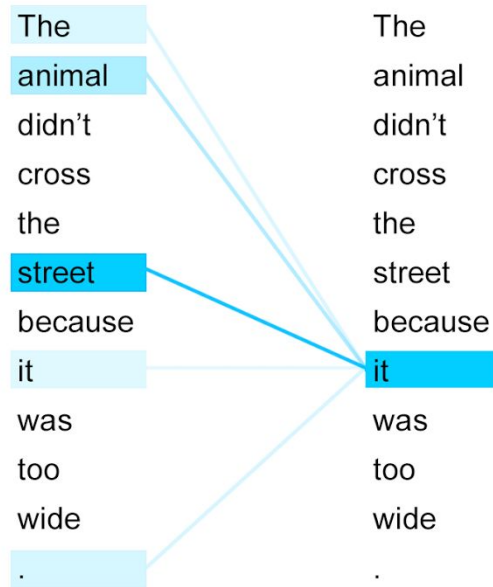
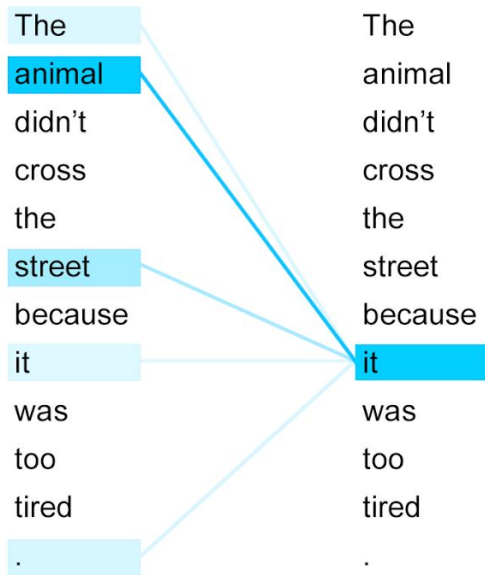
$$PE_{(pos, 2i+1)} = \cos(pos/10000^{2i/d_{\text{model}}})$$

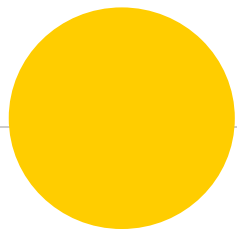






# Self-attention illustration





# Вопросы

- Что подается на вход Encoder-у трансформера?
- Что такое  $q$ ,  $k$ ,  $v$  в слое self-attention?
- Чем отличаются слои self-attention у Encoder-а и Decoder-а?



## Resources

- <https://arxiv.org/pdf/1706.03762.pdf> (original paper)
- [https://lena-voita.github.io/nlp\\_course/seq2seq\\_and\\_attention.html](https://lena-voita.github.io/nlp_course/seq2seq_and_attention.html) (blog post / NLP textbook )
- <http://jalammar.github.io/illustrated-transformer/> (classic detailed explanation)
- [https://www.youtube.com/watch?v=S0KakHcj\\_rs&t=1132s](https://www.youtube.com/watch?v=S0KakHcj_rs&t=1132s) (video on the paper)
- <https://www.youtube.com/watch?v=QEw0qEa0E50&feature=youtu.be> (CS224n, Stanford's course on NLP)
- <https://ai.googleblog.com/2017/08/transformer-novel-neural-network.html> (paper blog post)