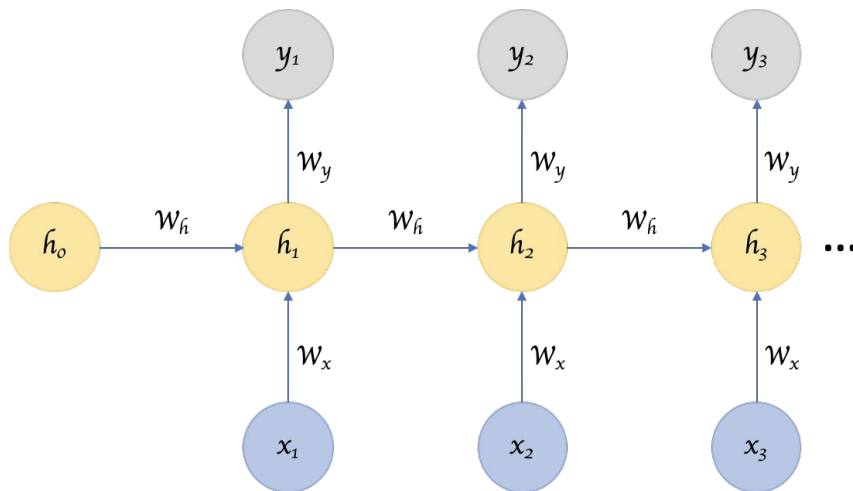


# Deep Learning

## Lecture 5

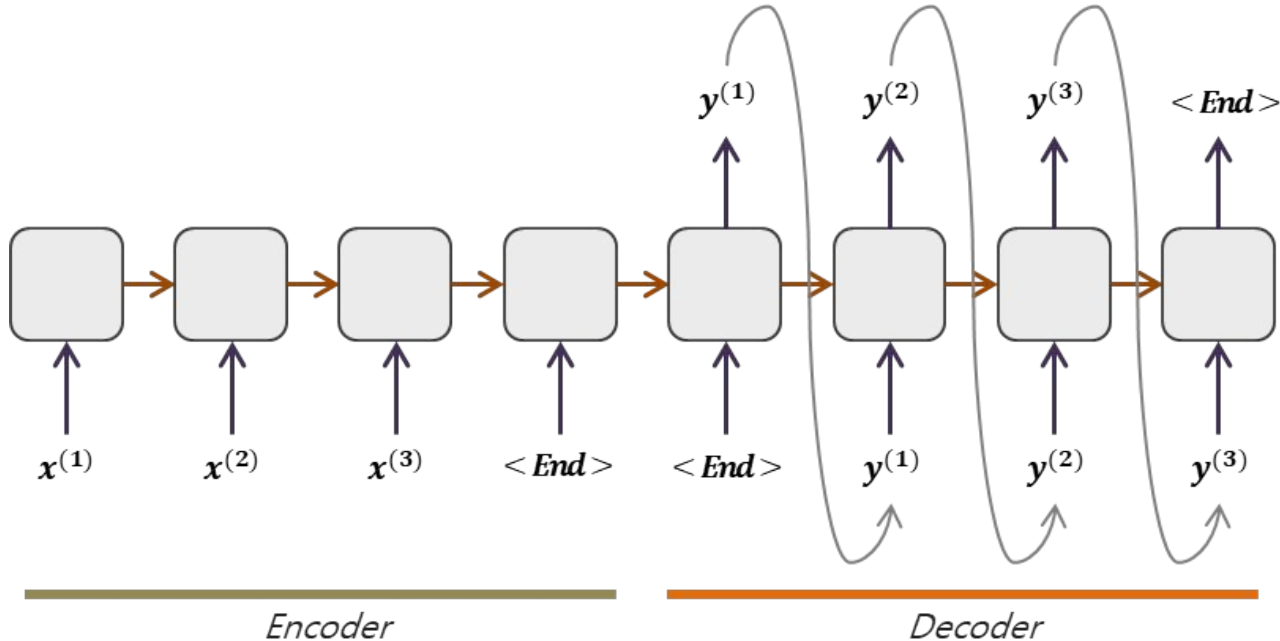
# Recap



- RNN
  - Backward
  - Vanishing gradients
  - Exploding gradients
- LSTM
- GRU
- Applications
- Dropout and BN

# Problem: copy of sequence

aeebdhf#  $\longrightarrow$  aeebdhf aeebdhf (aeebdhf \* 2)



## Attention mechanism

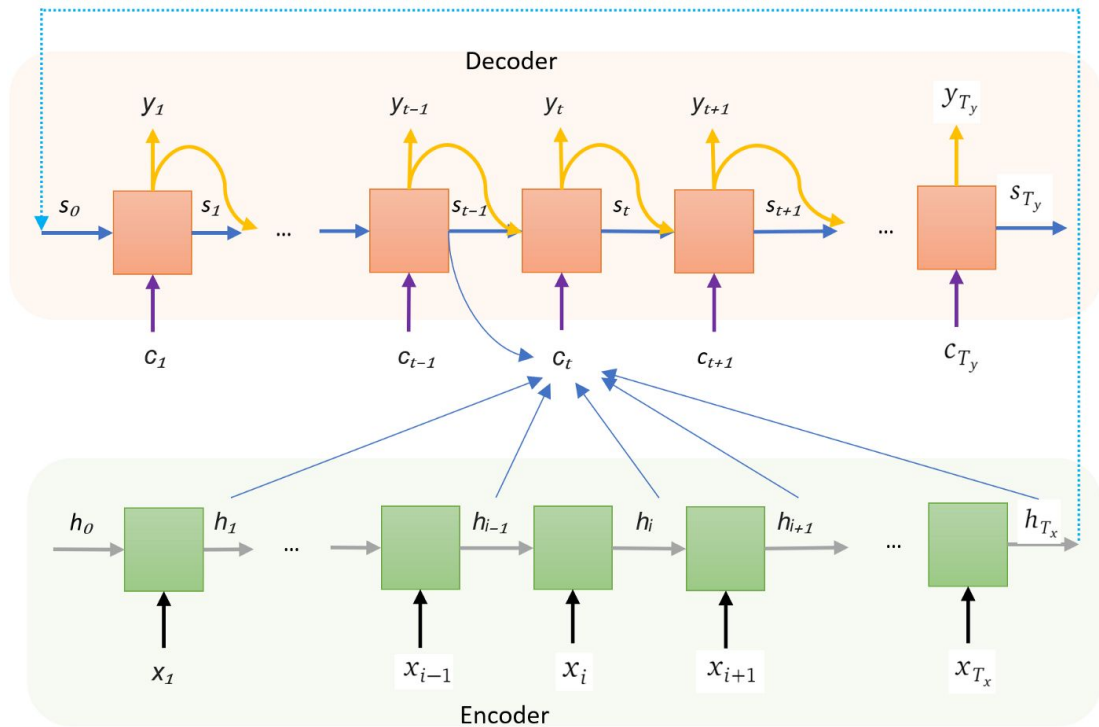
$$s_t \quad h_i$$

$$\text{sim}(s_t, h_i) \in \mathbb{R}$$

$$\alpha_{t,j} = \frac{\exp(\text{sim}(s_{t-1}, h_j))}{\sum_{i=1}^{T_x} \exp(\text{sim}(s_{t-1}, h_i))}$$

$$c_t = \sum_{i=1}^{T_x} \alpha_{t,i} h_i$$

# Attention mechanism



$$\alpha_{t,j} = \frac{\exp(\text{sim}(s_{t-1}, h_j))}{\sum_{i=1}^{T_x} \exp(\text{sim}(s_{t-1}, h_i))}$$

$$c_t = \sum_{i=1}^{T_x} \alpha_{t,i} h_i$$

## Similarity function

$$\textit{sim}(h_1, h_2) - ?$$

$$1) \quad h_1^T h_2 \quad \dim(h_1) = \dim(h_2)$$

$$2) \quad \frac{h_1^T h_2}{\sqrt{\dim(h_1)}} \quad \dim(h_1) = \dim(h_2)$$

Similarity function

$$\textit{sim}(h_1, h_2) - ?$$

$$h_{1i} \sim \mathcal{N}(0, 1)$$

$$\mathbb{E}[h_1^T h_2] = 0$$

$$h_{2i} \sim \mathcal{N}(0, 1)$$

$$\mathbb{D}[h_1^T h_2] = \textit{dim}(h_1)$$

## Similarity function

$$\textit{sim}(h_1, h_2) - ?$$

$$1) \quad h_1^T h_2 \quad \dim(h_1) = \dim(h_2)$$

$$2) \quad \frac{h_1^T h_2}{\sqrt{\dim(h_1)}} \quad \dim(h_1) = \dim(h_2)$$

$$3) \quad w_3 \tanh(w_1 h_1 + w_2 h_2) \quad \dim(h_1) \neq \dim(h_2)$$



# Attention

$$c = \text{Attention}(k : \text{"keys"}, q : \text{"query"}, v : \text{"values"})$$

$$\alpha_t = \frac{\exp(\text{sim}(q, k_t))}{\sum_{i=1}^T \exp(\text{sim}(q, k_i))}$$

$$c = \sum_{i=1}^T \alpha_i v_i$$

# Handwriting generation:

handwriting  $\rightarrow$  handwriting

Next pen position (we predict parameters):

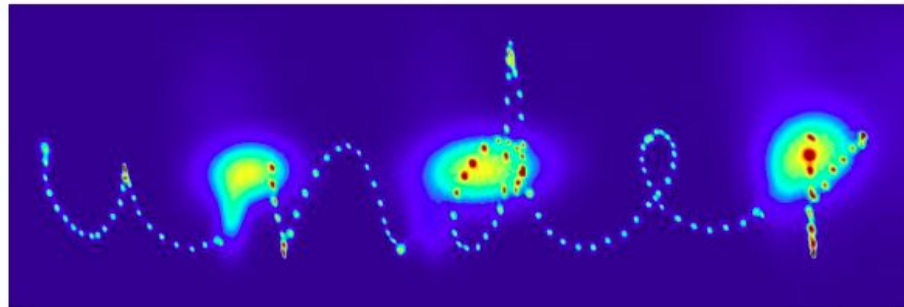
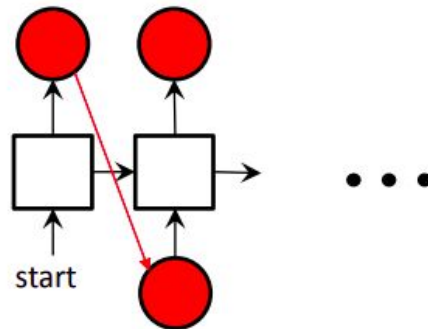
$x_1, x_2$  - mixture of bivariate Gaussians

$x_3$  - Bernoulli distribution

Current pen position:

$x_1, x_2$  - pen offset

$x_3$  - is it end of the stroke



# Handwriting generation:

example

when my under grow cage there will

pegged and the 'bepestures the the

maine Ceneh to of high vraditro'

see Boring a. the acorattures ka

pure n mist taken so learned

bopes & cold mine's wine case

heist. Y Ceehs the gaster in

style satet Boring In doing Te a

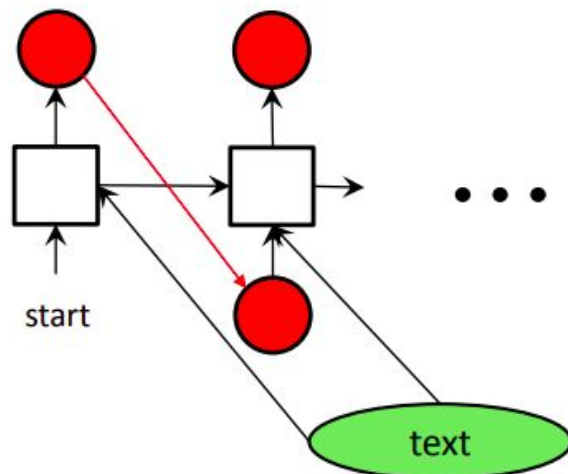
# Handwriting synthesis:

text -> handwriting

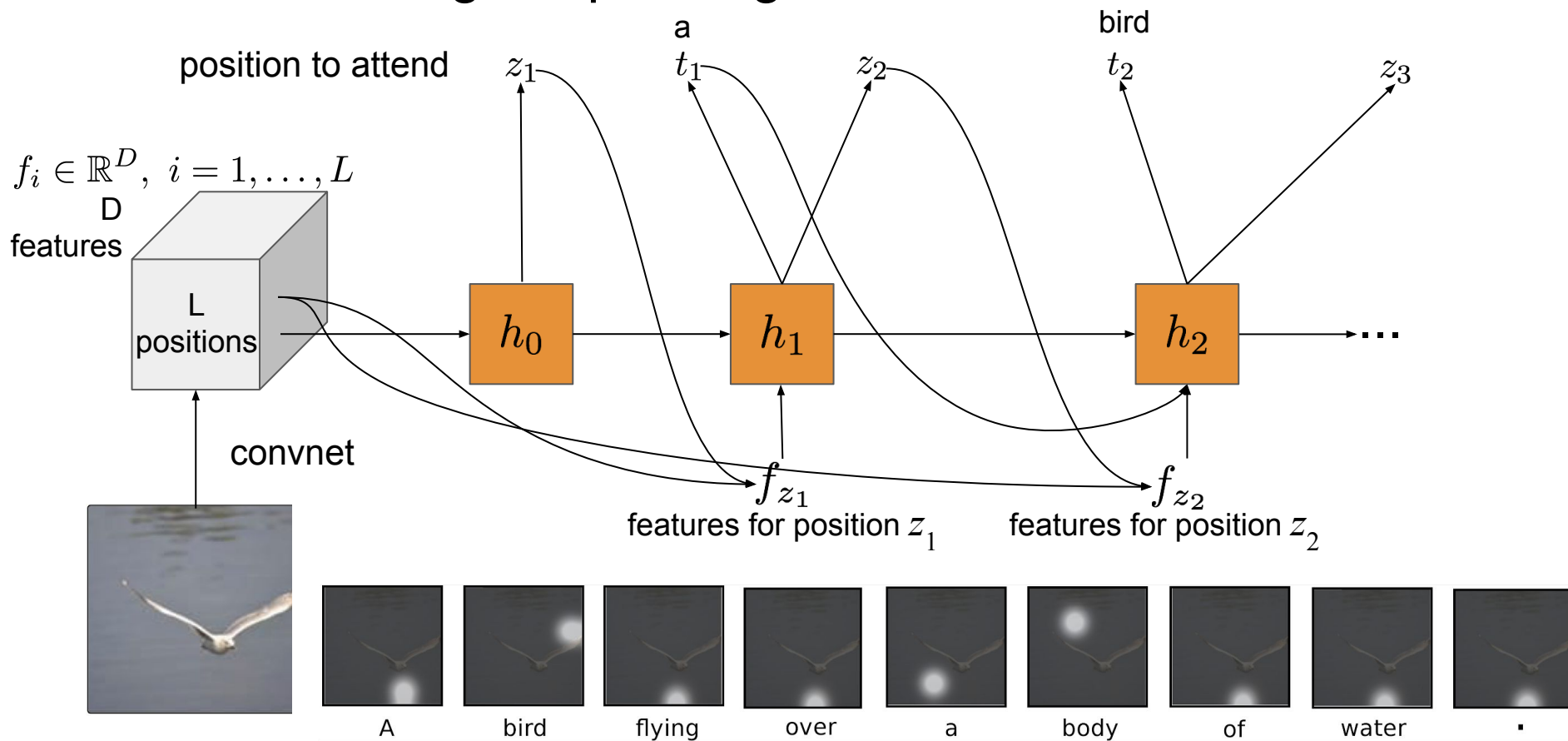
Next pen position

Current pen position

Which letter we write now



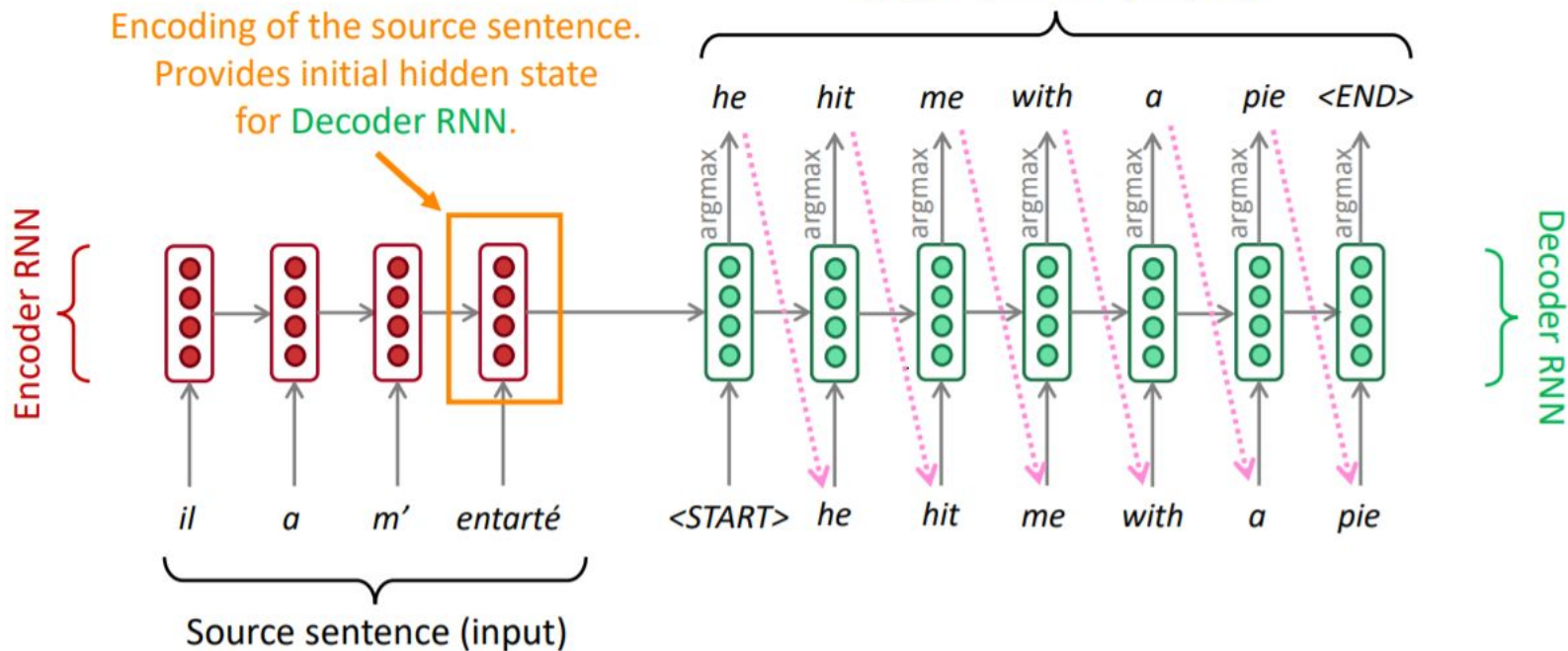
# Attention for image captioning



# Translation with attention

The sequence-to-sequence model

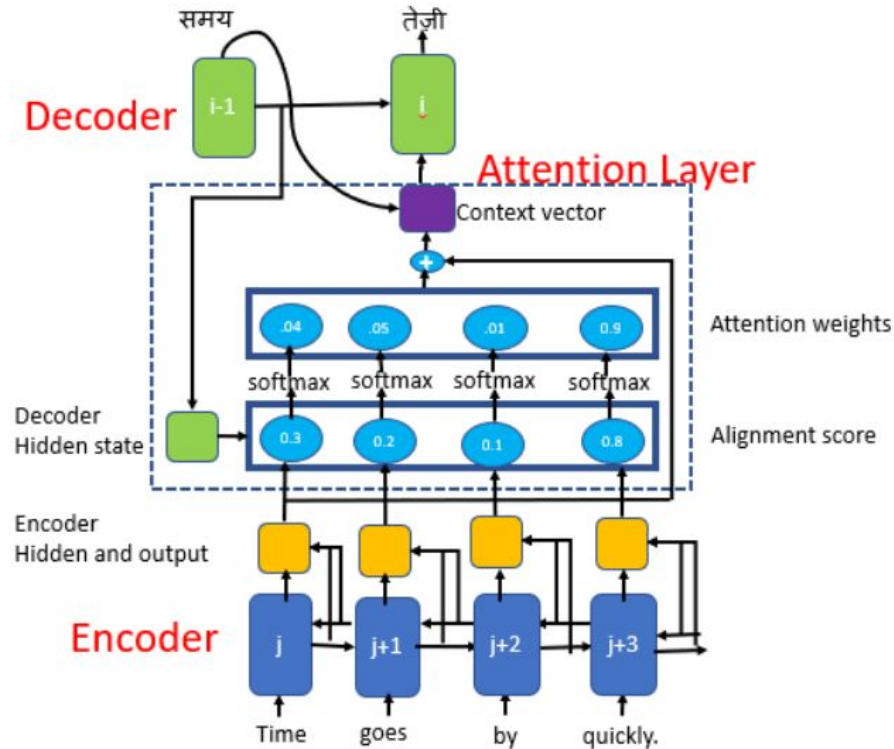
Encoding of the source sentence.  
Provides initial hidden state  
for Decoder RNN.



# Translation with attention

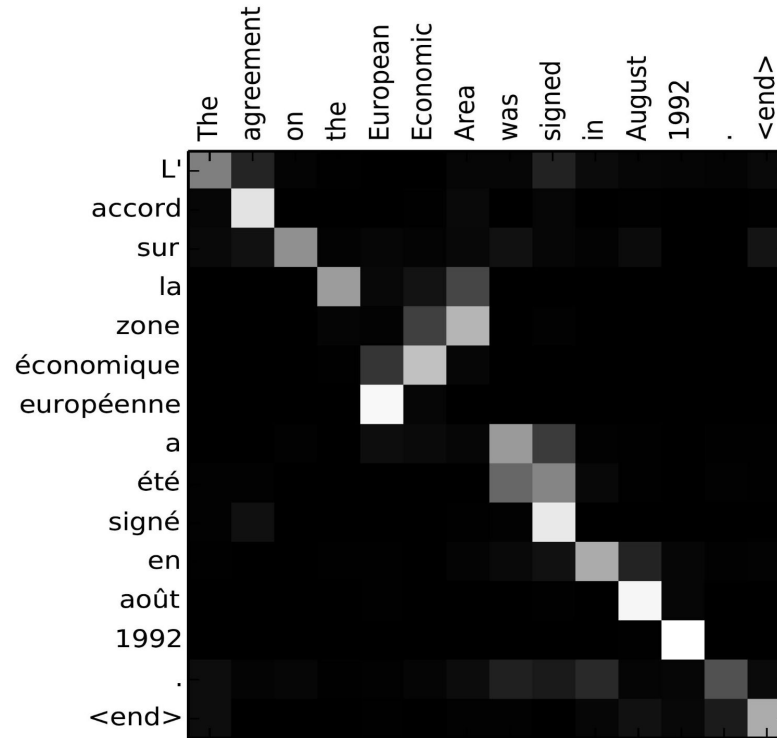


# Translation with attention





# Translation with attention



# Attention types

# Self-attention

The FBI is chasing a criminal on the run .

The FBI is chasing a criminal on the run .

The FBI is chasing a criminal on the run .

The FBI is chasing a criminal on the run .

The FBI is chasing a criminal on the run .

The FBI is chasing a criminal on the run .

The FBI is chasing a criminal on the run .

The FBI is chasing a criminal on the run .

The FBI is chasing a criminal on the run .

The FBI is chasing a criminal on the run .

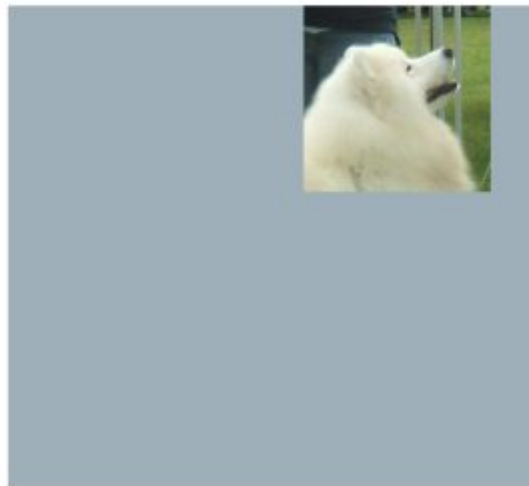
The current word is in red and the size of the blue shade indicates the activation level.

# Soft vs Hard Attention

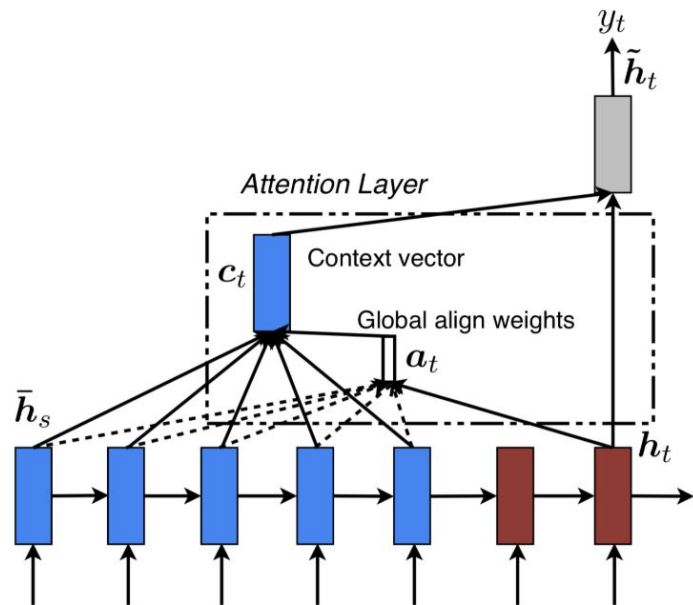
Soft Attention



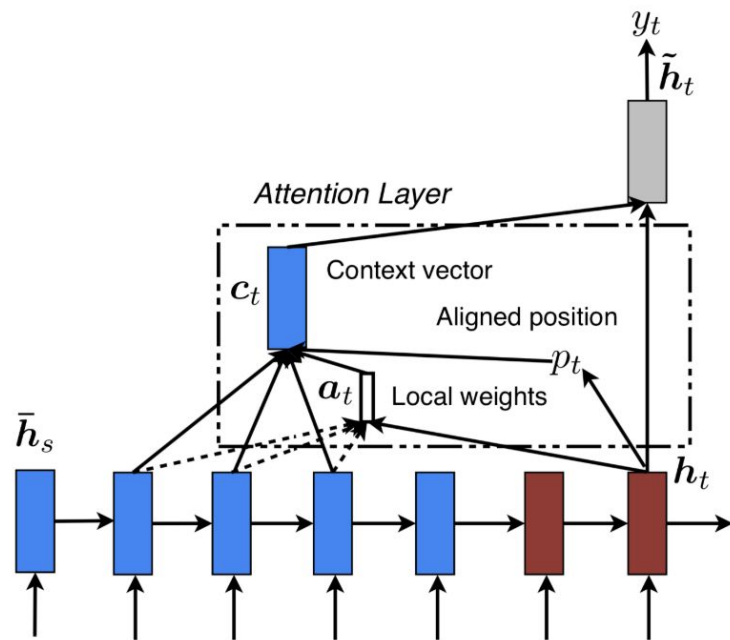
Hard Attention



# Global vs Local Attention



Global Attention Model



Local Attention Model

# What's next?



# Attention is all you need

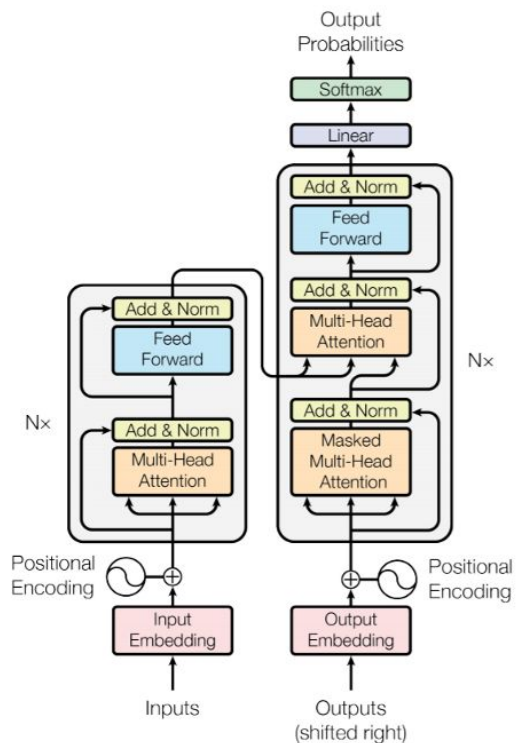
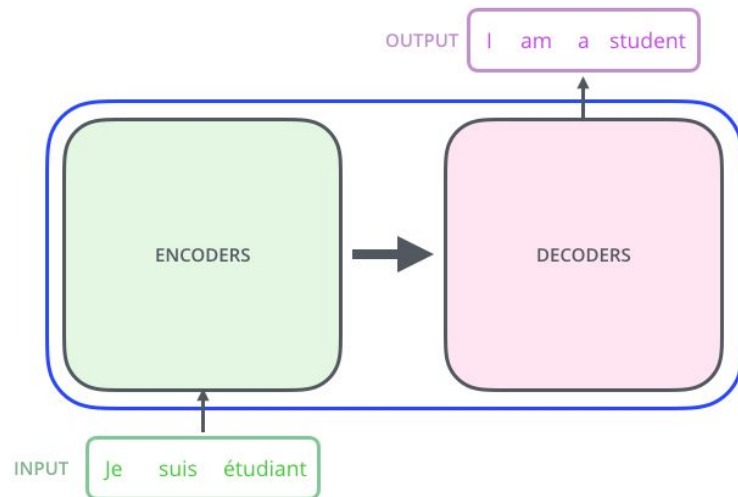
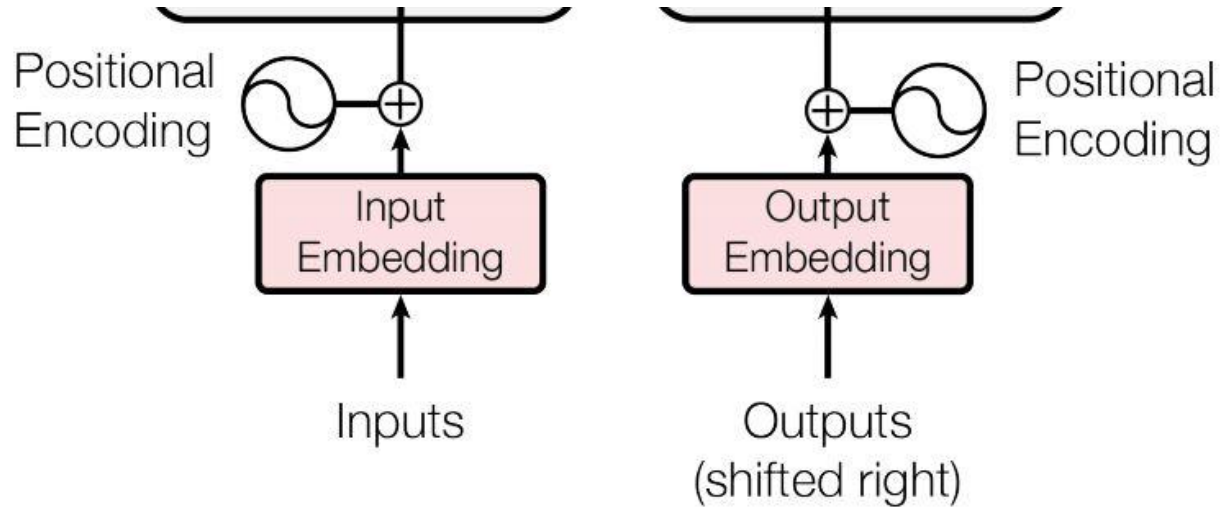


Figure 1: The Transformer - model architecture.



# Input layer





# Positional encoding

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

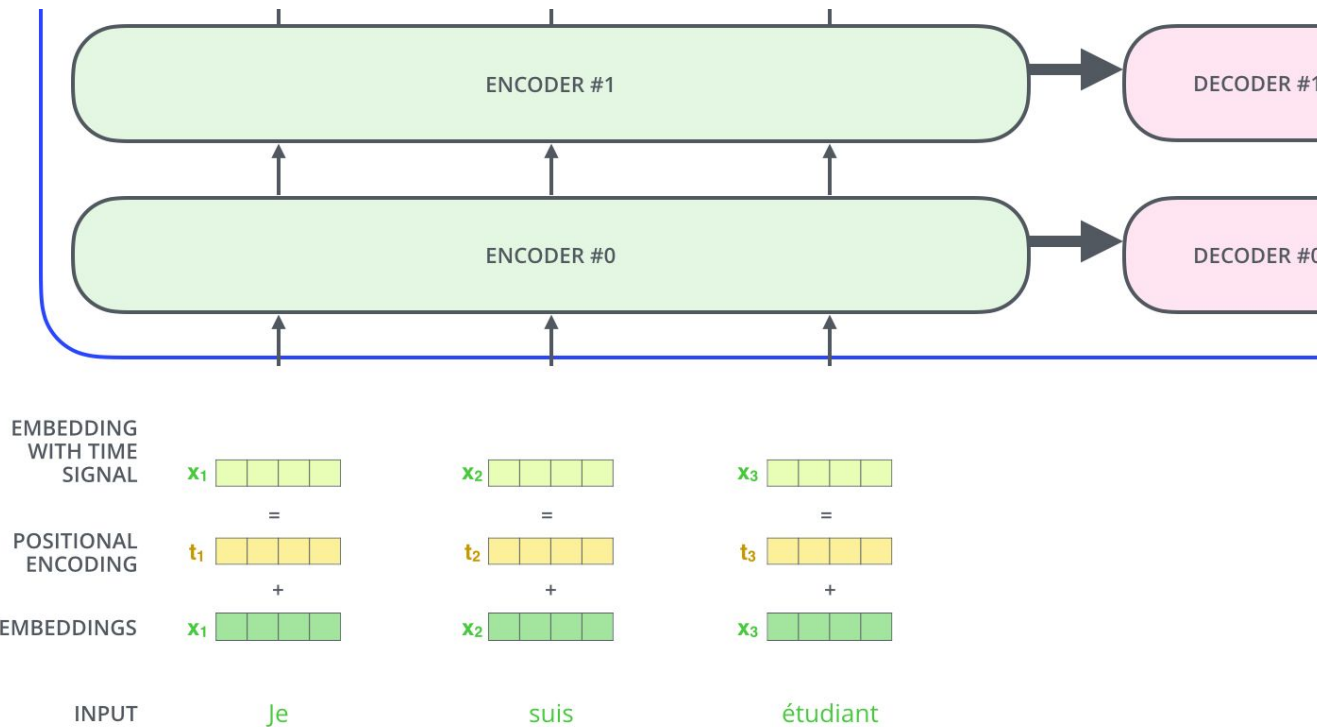
$$PE_{\text{pos},2i+1} = \cos\left(\frac{\text{pos}}{10000^{\frac{2i}{d_{\text{model}}}}}\right)$$

$$PE_{\text{pos}+k,2i} = \sin\left(\frac{\text{pos} + k}{10000^{\frac{2i}{d_{\text{model}}}}}\right) = \sin\left(\frac{\text{pos}}{10000^{\frac{2i}{d_{\text{model}}}}} + \frac{k}{10000^{\frac{2i}{d_{\text{model}}}}}\right) =$$

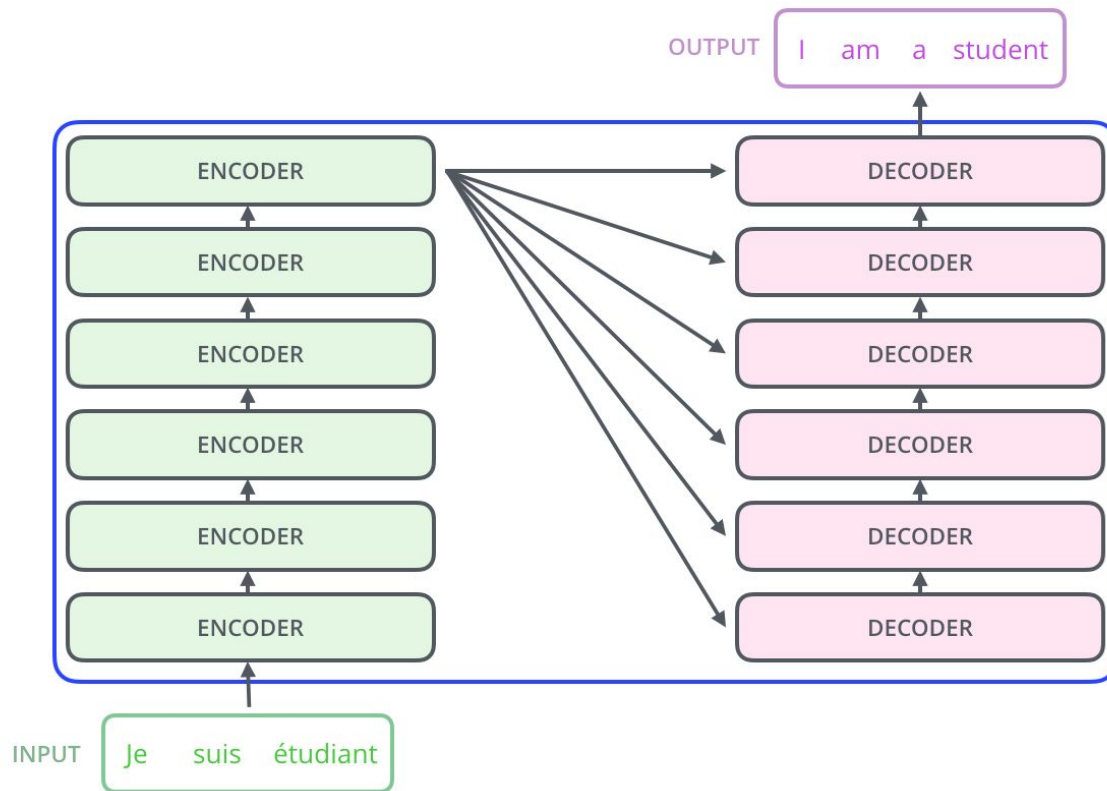
$$= \sin(x + y) = \sin(x) \cos(y) + \cos(x) \sin(y)$$

$$PE_{\text{pos}+k,2i+1} = \cos(x) \cos(y) + \sin(x) \sin(y)$$

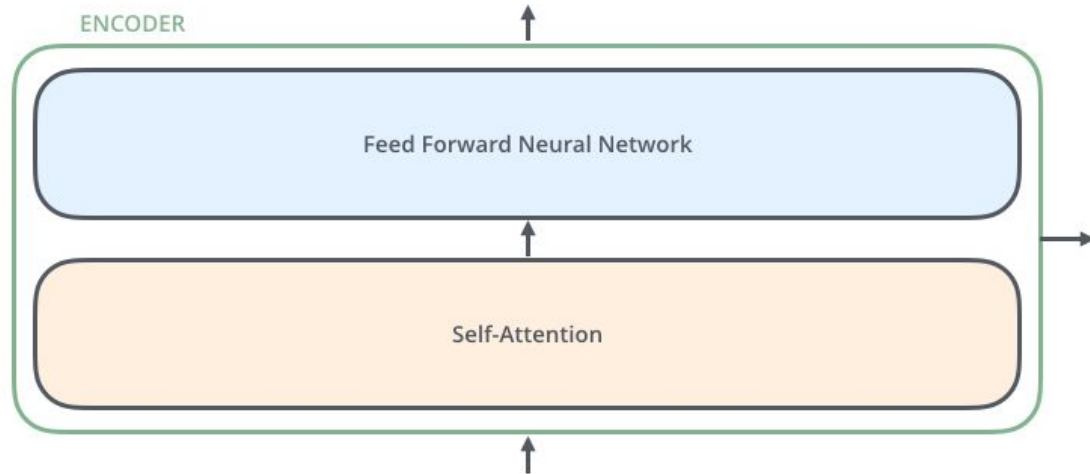
# Positional encoding



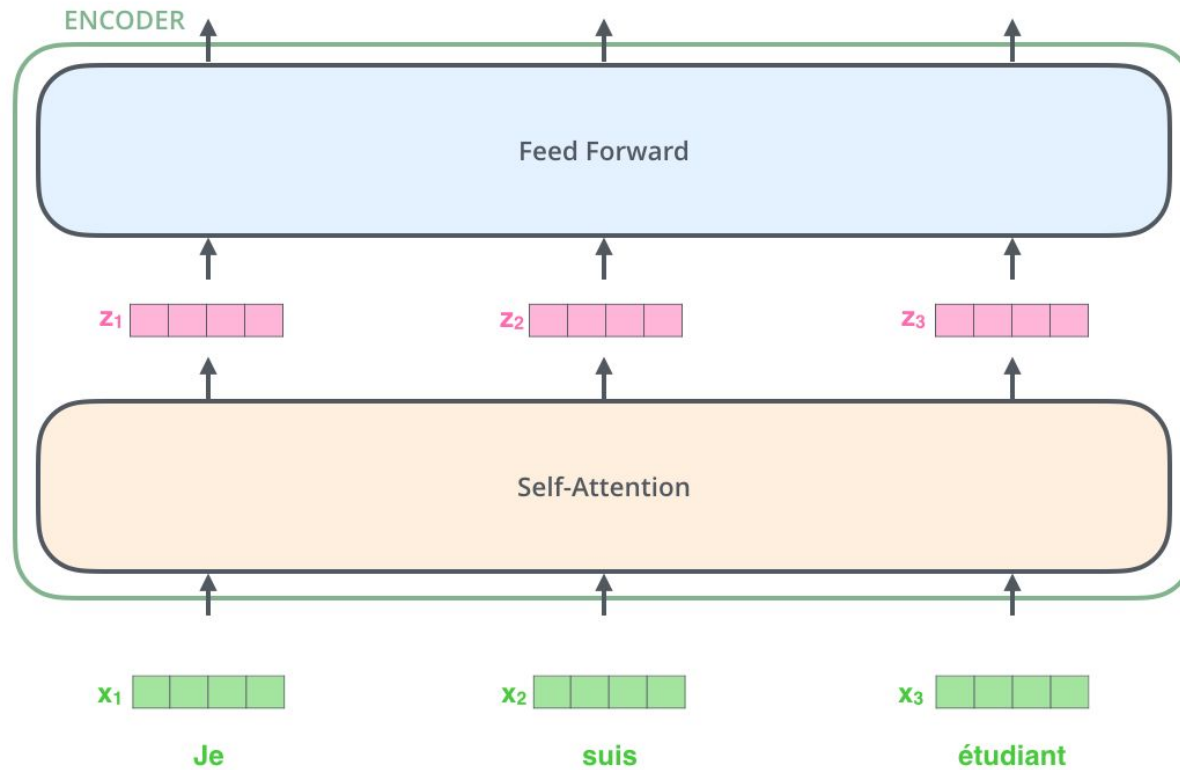
# Attention is all you need



# Encoder

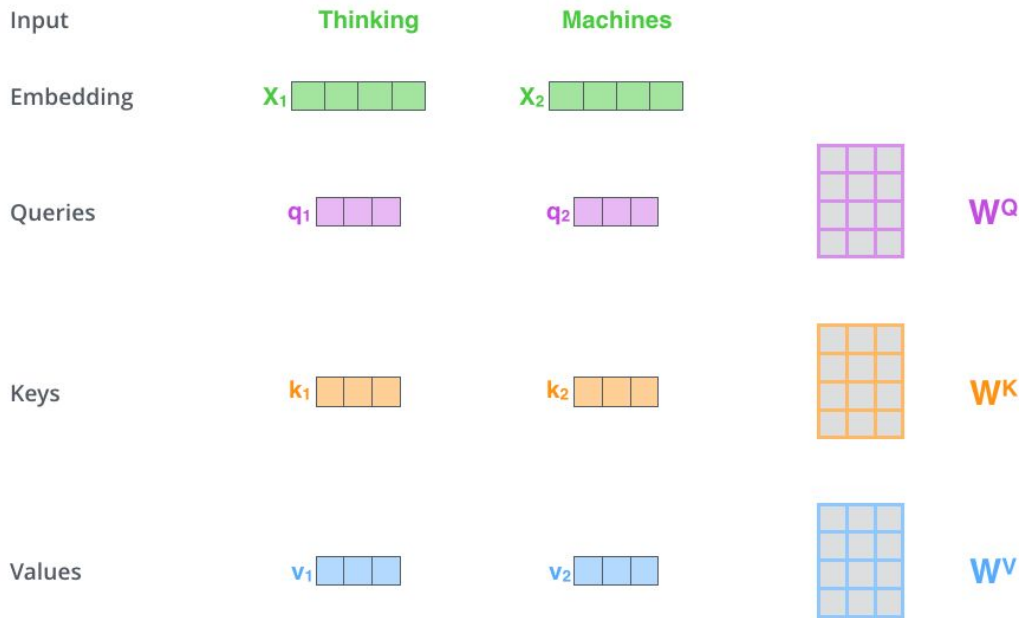


# Encoder



# Self attention

$$\text{Attention}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \text{softmax}\left(\frac{\mathbf{Q}\mathbf{K}^\top}{\sqrt{\dim}}\right)\mathbf{V}$$



# Self attention

Input

Embedding

Queries

Keys

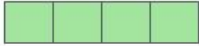
Values

Score

Divide by 8 (  $\sqrt{d_k}$  )

Softmax

Thinking

$x_1$  

$q_1$  

$k_1$  

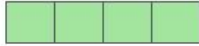
$v_1$  

$q_1 \cdot k_1 = 112$

14

0.88

Machines

$x_2$  

$q_2$  

$k_2$  

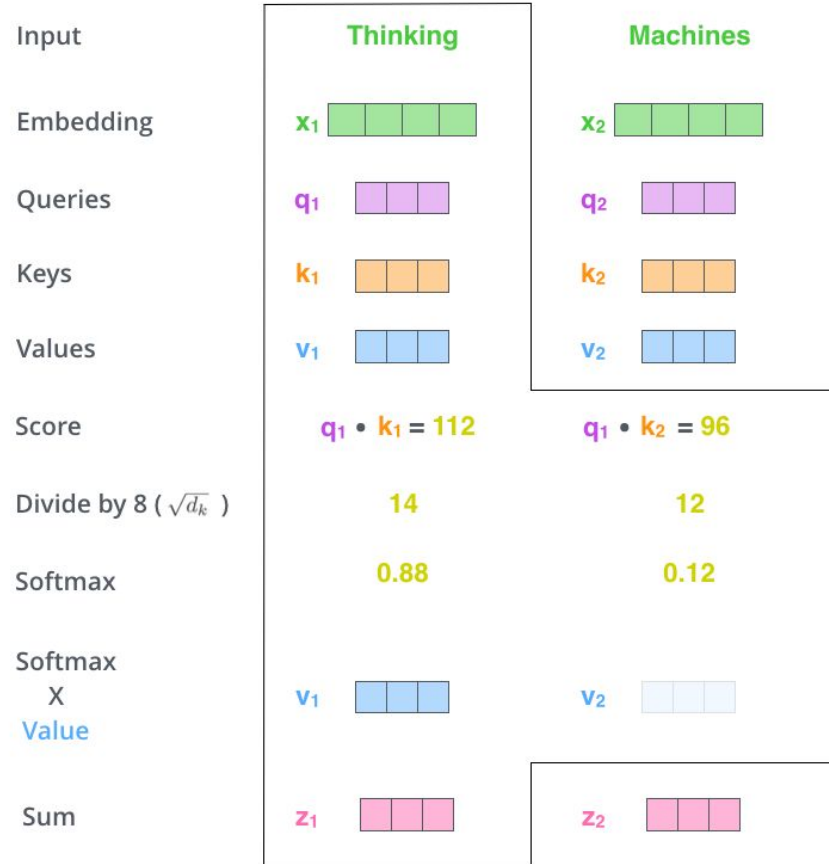
$v_2$  

$q_1 \cdot k_2 = 96$

12

0.12

# Self attention





# Self attention

$$\begin{matrix} \text{X} \\ \begin{array}{|c|c|c|c|} \hline & & & \\ \hline & & & \\ \hline & & & \\ \hline \end{array} \end{matrix} \times \begin{matrix} \text{W}^{\text{Q}} \\ \begin{array}{|c|c|c|} \hline & & \\ \hline & & \\ \hline & & \\ \hline & & \\ \hline \end{array} \end{matrix} = \begin{matrix} \text{Q} \\ \begin{array}{|c|c|c|} \hline & & \\ \hline & & \\ \hline & & \\ \hline \end{array} \end{matrix}$$

$$\begin{matrix} \text{X} \\ \begin{array}{|c|c|c|c|} \hline & & & \\ \hline & & & \\ \hline & & & \\ \hline \end{array} \end{matrix} \times \begin{matrix} \text{W}^{\text{K}} \\ \begin{array}{|c|c|c|} \hline & & \\ \hline & & \\ \hline & & \\ \hline & & \\ \hline \end{array} \end{matrix} = \begin{matrix} \text{K} \\ \begin{array}{|c|c|c|} \hline & & \\ \hline & & \\ \hline & & \\ \hline \end{array} \end{matrix}$$

$$\begin{matrix} \text{X} \\ \begin{array}{|c|c|c|c|} \hline & & & \\ \hline & & & \\ \hline & & & \\ \hline \end{array} \end{matrix} \times \begin{matrix} \text{W}^{\text{V}} \\ \begin{array}{|c|c|c|} \hline & & \\ \hline & & \\ \hline & & \\ \hline & & \\ \hline \end{array} \end{matrix} = \begin{matrix} \text{V} \\ \begin{array}{|c|c|c|} \hline & & \\ \hline & & \\ \hline & & \\ \hline \end{array} \end{matrix}$$

# Self attention

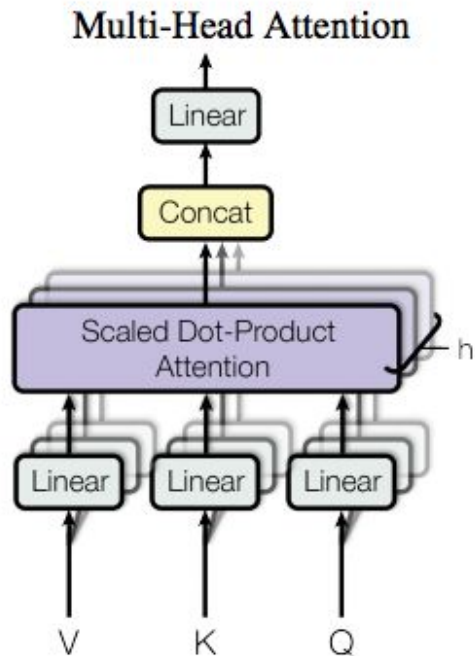
$$\text{softmax} \left( \frac{\begin{matrix} \text{Q} \\ \begin{array}{|c|c|c|} \hline \square & \square & \square \\ \hline \square & \square & \square \\ \hline \end{array} \end{matrix} \times \begin{matrix} \text{K}^T \\ \begin{array}{|c|c|} \hline \square & \square \\ \hline \square & \square \\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

=

$\text{Z}$

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

# Multi-Head Attention



$$\text{MultiHead}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = [\text{head}_1; \dots; \text{head}_h] \mathbf{W}^O$$

$$\text{head}_i = \text{Attention}(\mathbf{QW}_i^Q, \mathbf{KW}_i^K, \mathbf{VW}_i^V)$$

# Multi-Head Attention

1) This is our  
input sentence\*

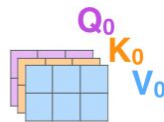
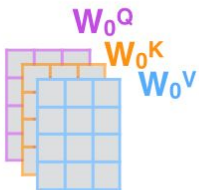
2) We embed  
each word\*

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

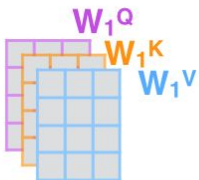
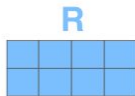
4) Calculate attention  
using the resulting  
 $Q/K/V$  matrices

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

Thinking  
Machines



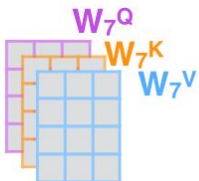
\* In all encoders other than #0,  
we don't need embedding.  
We start directly with the output  
of the encoder right below this one



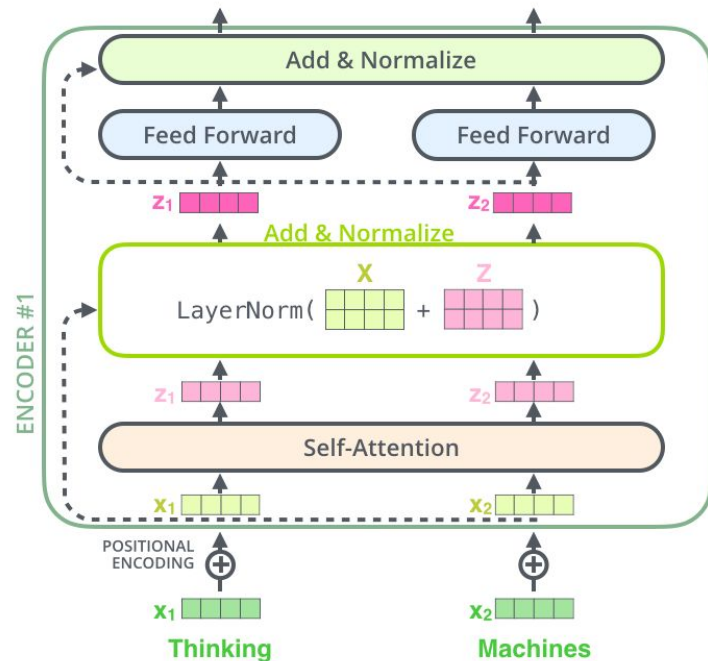
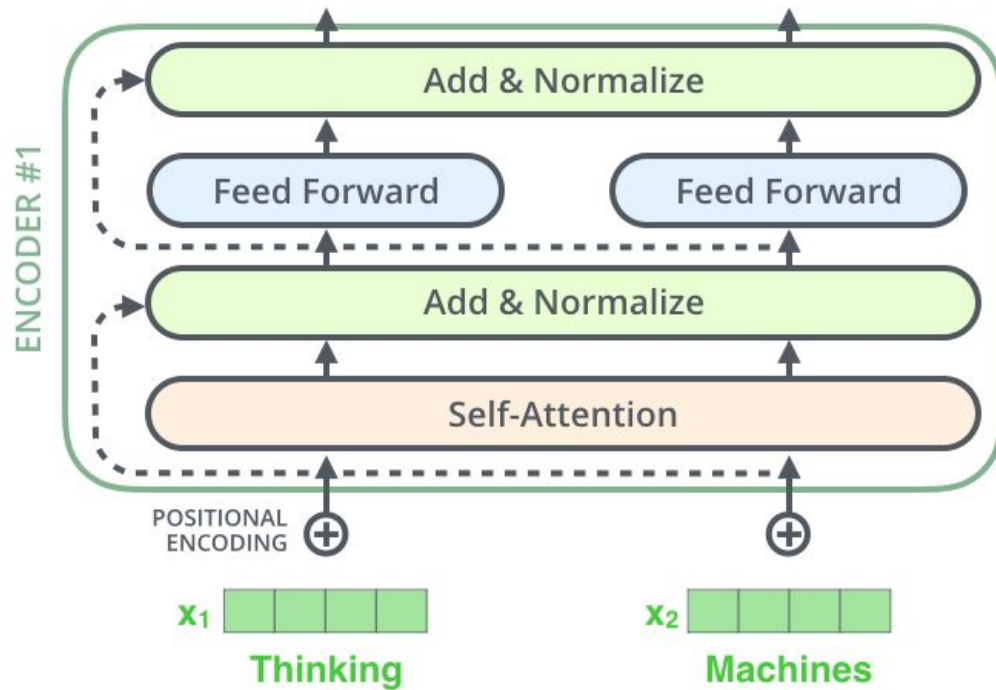
...

...

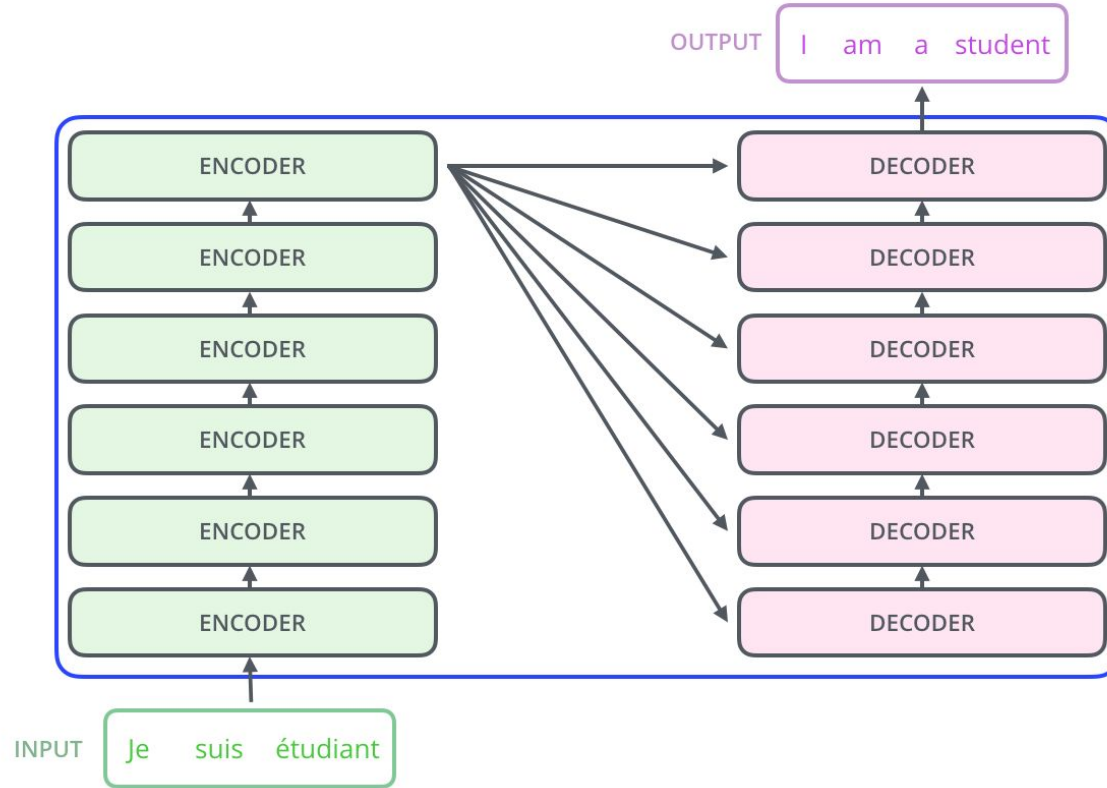
...



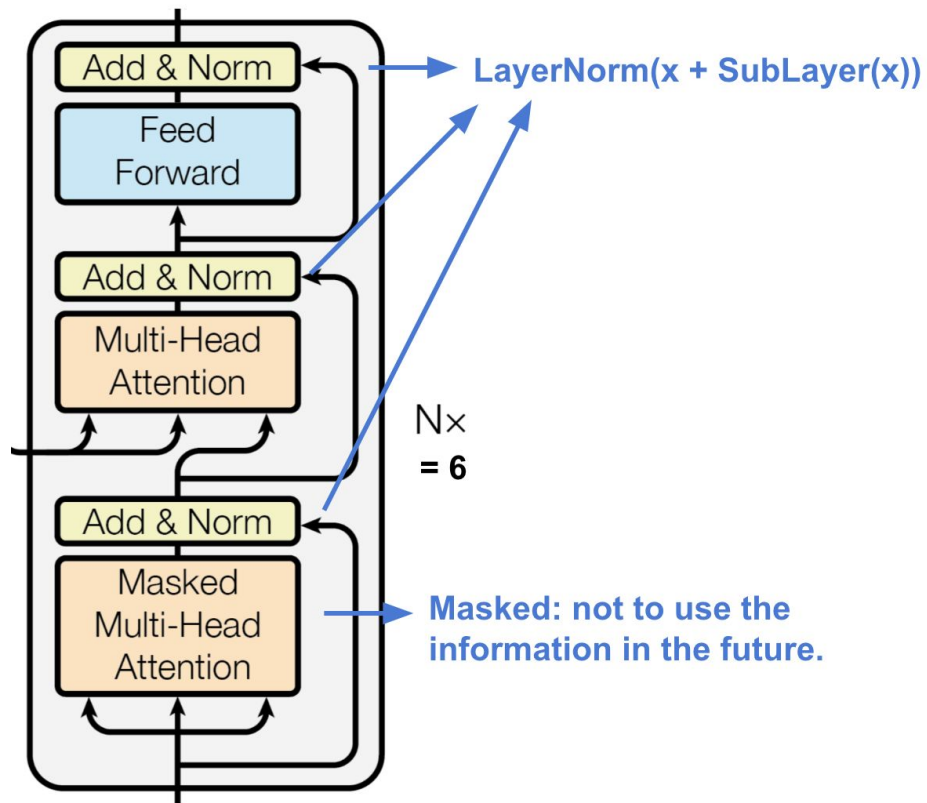
# Encoder



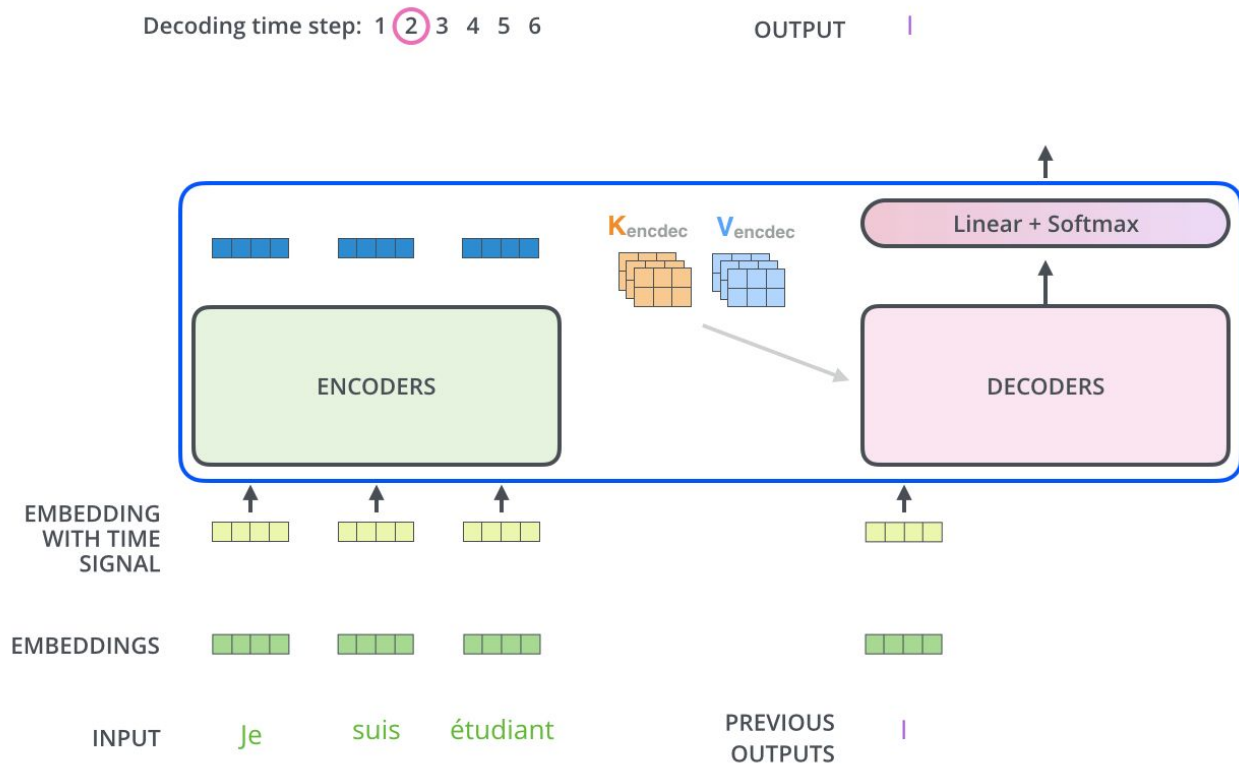
# Encoder\decoder



# Decoder



# Decoder





# Last layer

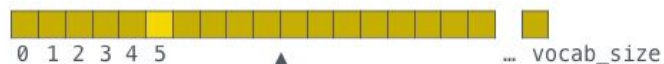
Which word in our vocabulary  
is associated with this index?

Get the index of the cell  
with the highest value  
(argmax)

am

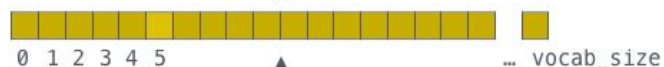
5

log\_probs



Softmax

logits



Linear

Decoder stack output



# Attention is all you need

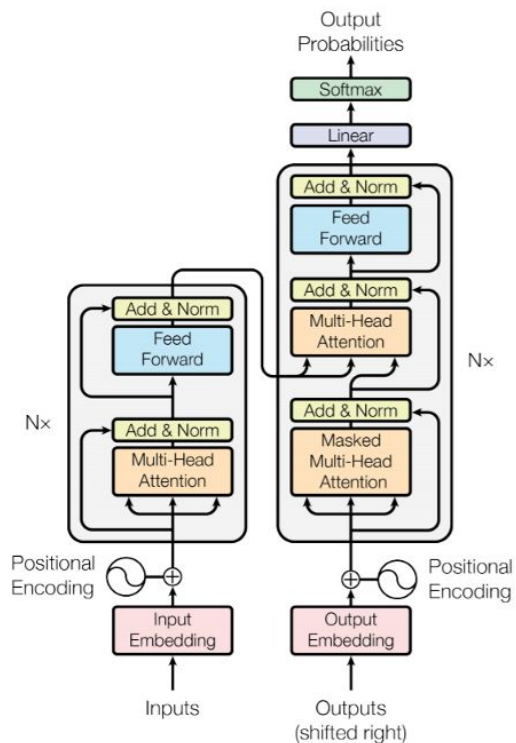
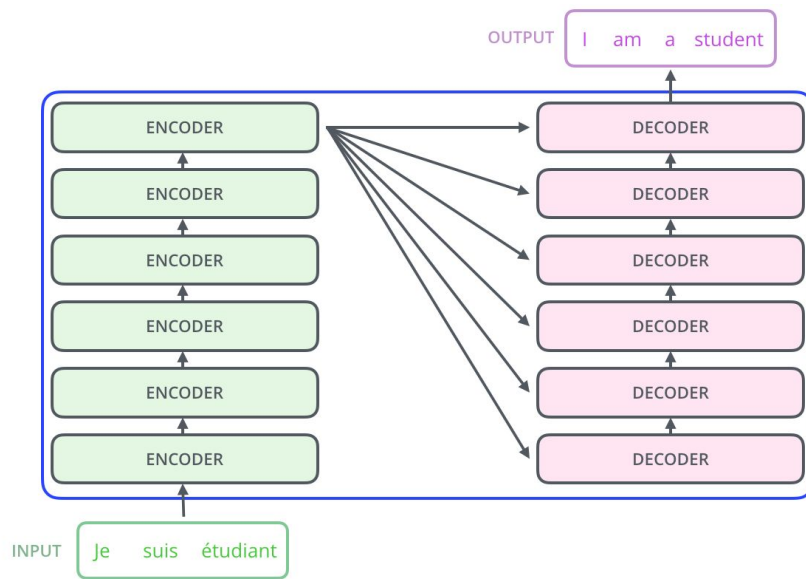
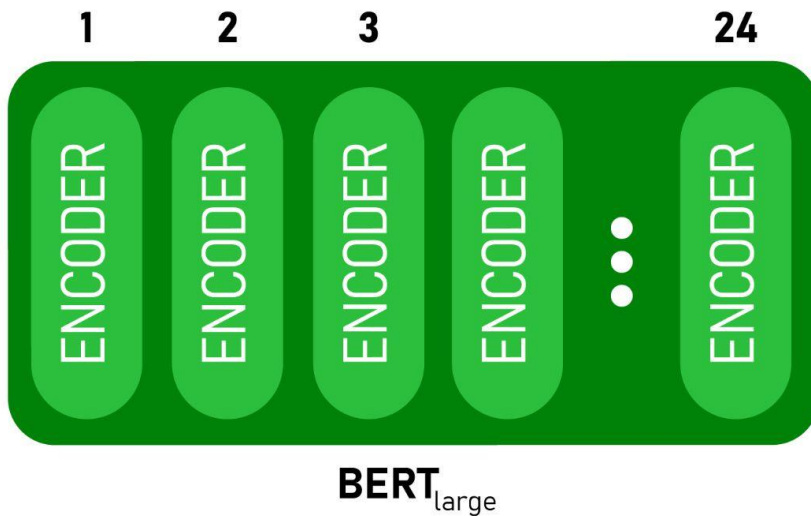
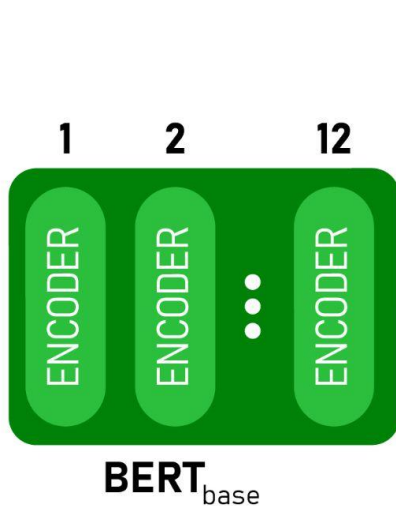


Figure 1: The Transformer - model architecture.

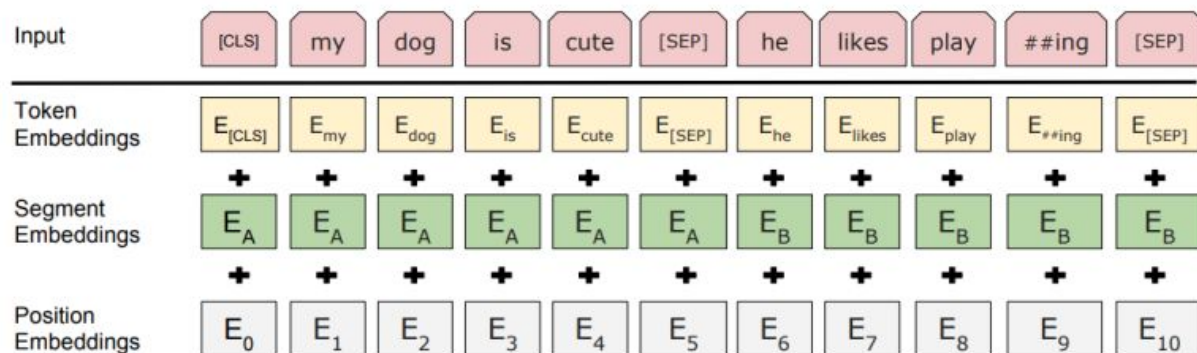


- no time loops
- can learn really long dependencies
- too many parameters

# BERT - Bidirectional Encoding Representation from Transformers



# Masked Language Modeling (MLM) and Next Sentence Prediction (NSP)



- 80% of the time the words were replaced with the masked token [MASK]
- 10% of the time the words were replaced with random words
- 10% of the time the words were left unchanged
- 50% of the time the second sentence comes after the first one.
- 50% of the time it is a random sentence from the full corpus.

# Masked Language Modeling (MLM) and Next Sentence Prediction (NSP)

Input = [CLS] the man went to [MASK] store [SEP]  
he bought a gallon [MASK] milk [SEP]

Label = IsNext

Input = [CLS] the man [MASK] to the store [SEP]  
penguin [MASK] are flight ##less birds [SEP]

Label = NotNext

# Recap

- Attention
- Applications
- Types of attention
- Transformer
  - Positional encoding
  - Self-attention
  - Multi-head attention
- BERT model (MLM)