
Seq2seq, Attention, Self attention, Transformer, BERT

— Tuan Nguyen - AI4E —

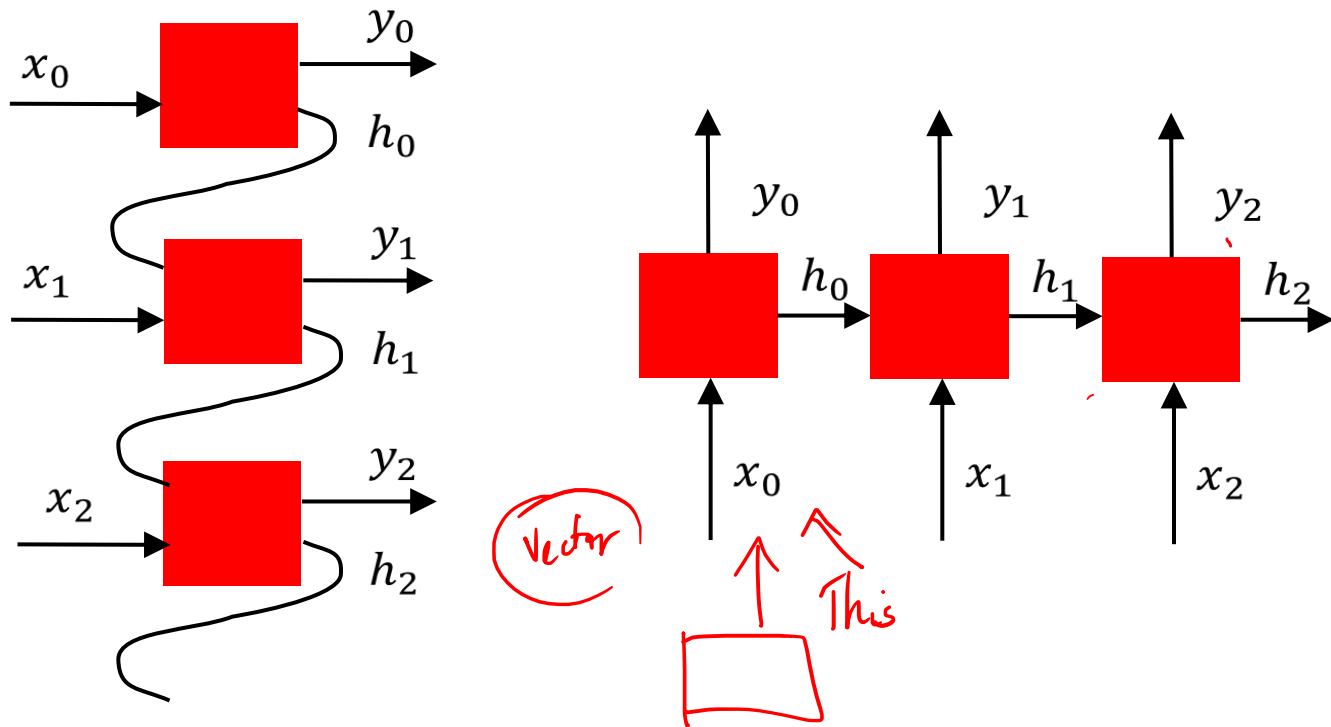
Outline

- RNN review
- Seq2seq
- Beam search
- Attention
- Self-attention
- Transformer
- BERT

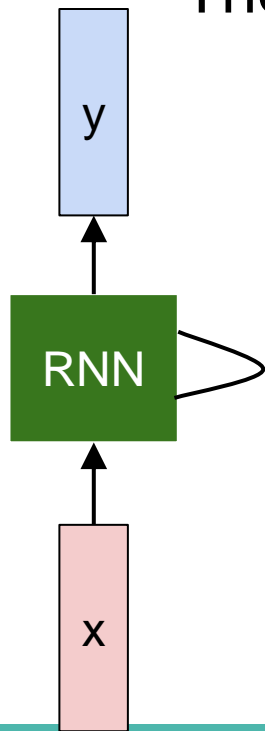
Recurrent Neural Network

*sequence
time-series.*

Usually drawn as:



RNN Formula



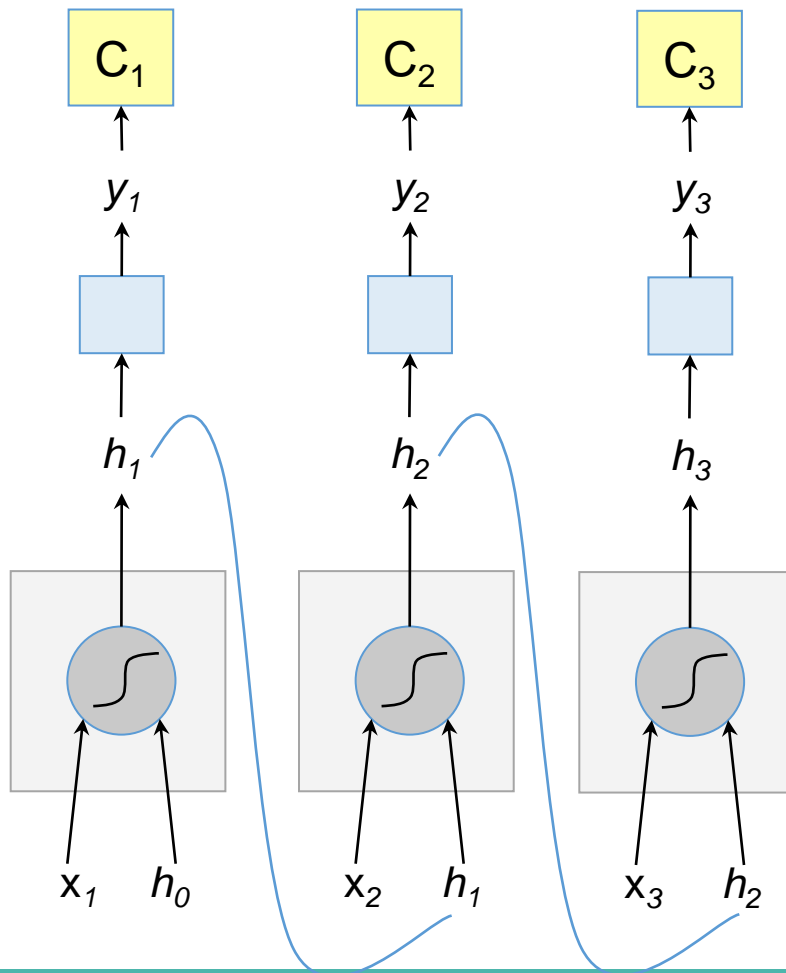
The state consists of a single “*hidden*” vector **h**:

$$h_t = f_{\underline{W}}(h_{t-1}, x_t)$$

↓

$$\left\{ \begin{array}{l} h_t = \tanh(\underline{W}_{hh}h_{t-1} + \underline{W}_{xh}x_t) \\ y_t = \underline{W}_{hy}h_t \end{array} \right.$$

Forward

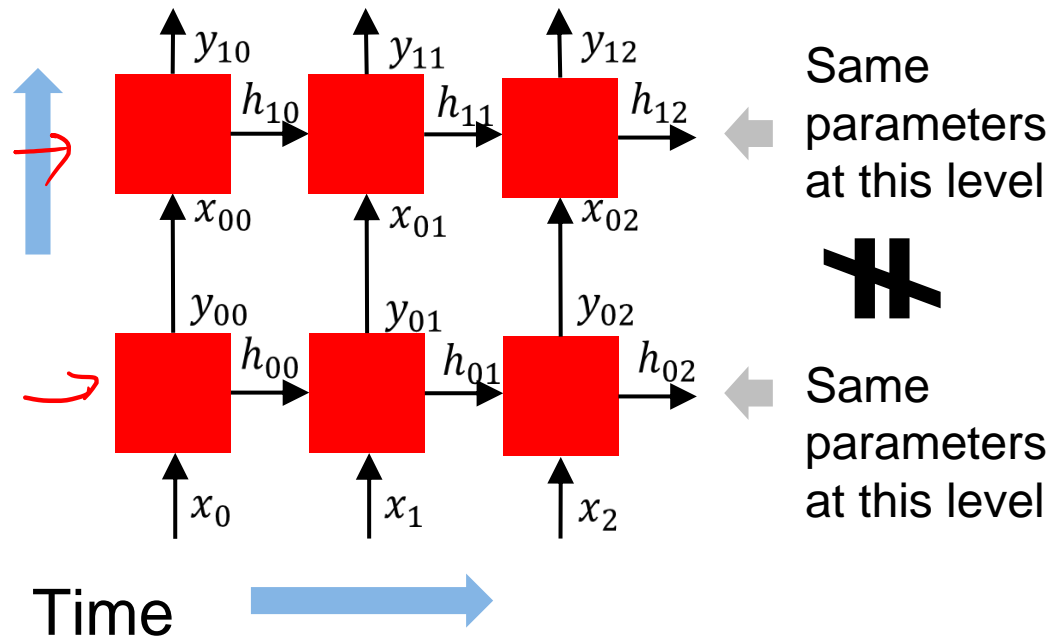
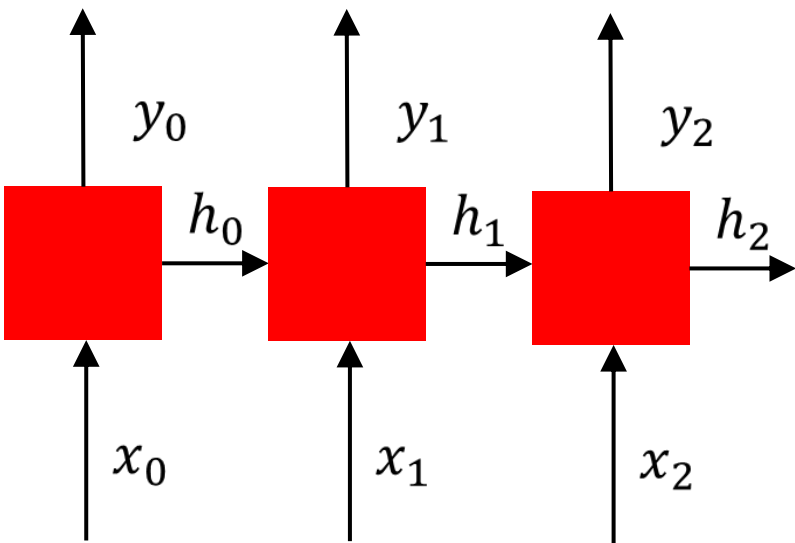


$$h_t = \tanh W \begin{pmatrix} x_t \\ h_{t-1} \end{pmatrix}$$

$$y_t = F(h_t)$$

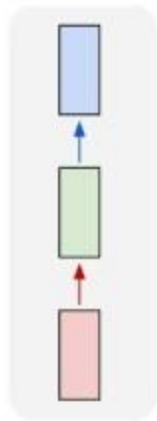
$$C_t = \text{Loss}(y_t, \text{GT}_t)$$

Deep RNN



Recurrent neural network problem

one to one



nn
cnn

one to many

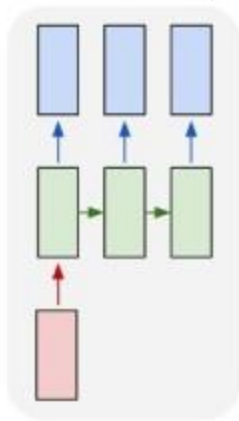
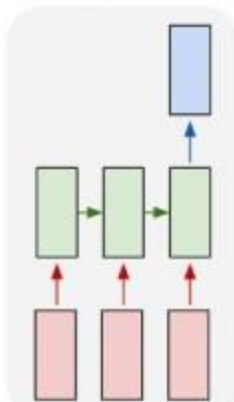


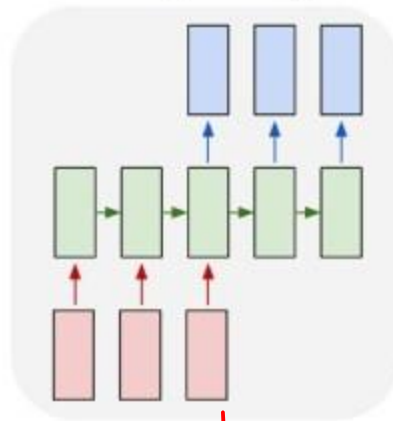
image
captioning

many to one



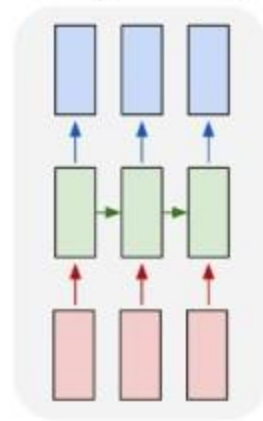
video classification

many to many



translation

many to many

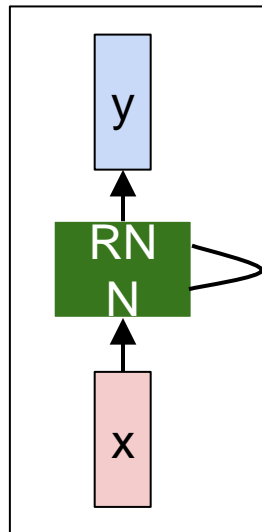


character generation

Character-level language model example

Vocabulary:
[h,e,l,o]

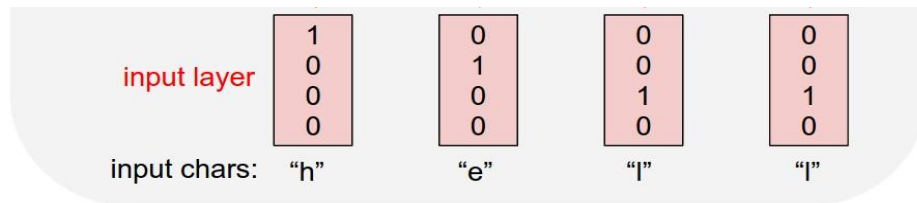
Example training
sequence:
“hello”



Character-level language model example

Vocabulary:
[h,e,l,o]

Example training
sequence:
“hello”



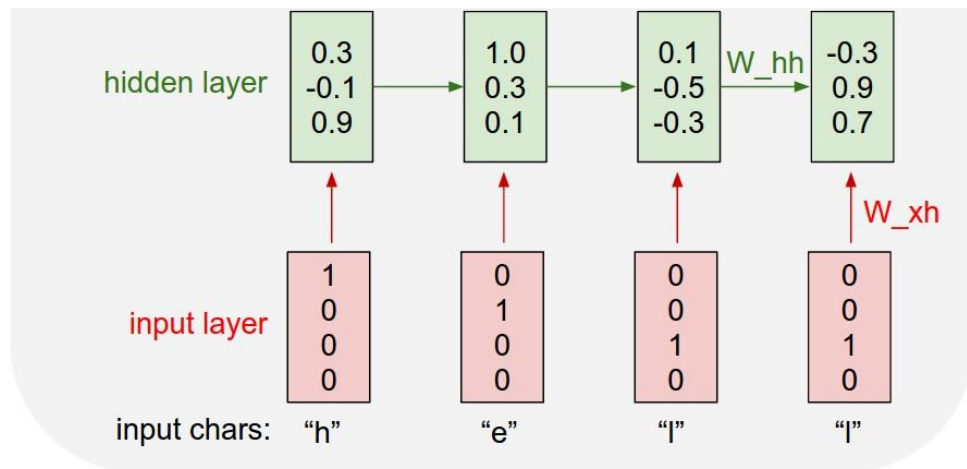
Character-level language model example

Vocabulary:

[h,e,l,o]

Example training
sequence:
“hello”

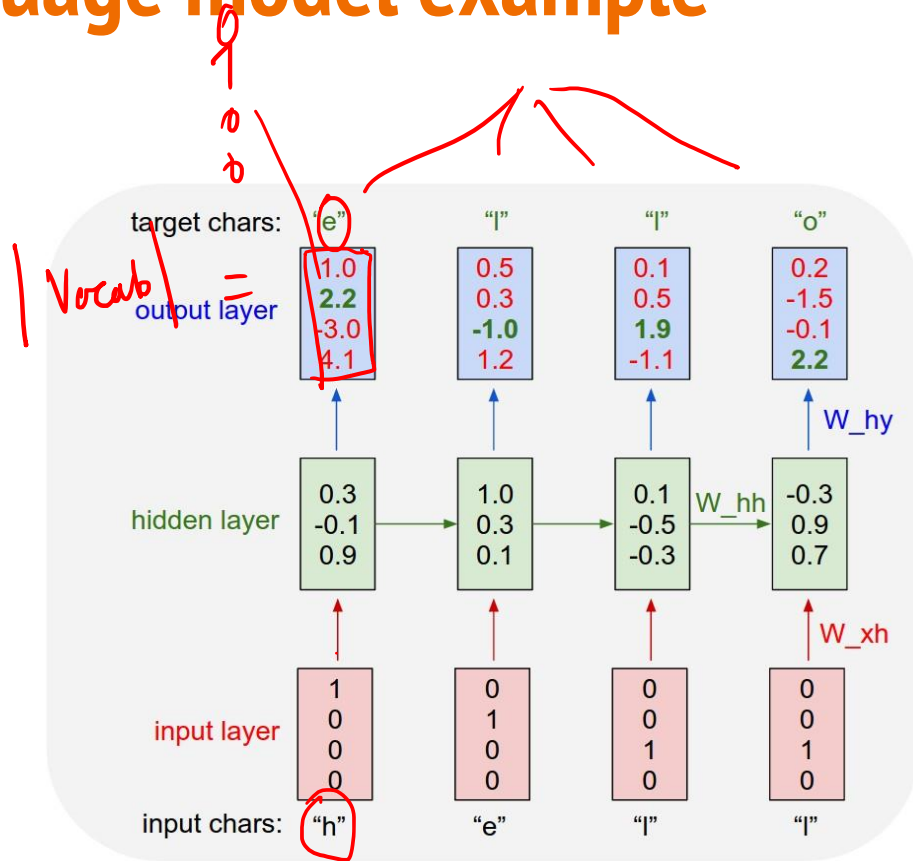
$$h_t = \tanh(W_{hh}h_{t-1} + W_{xh}x_t)$$



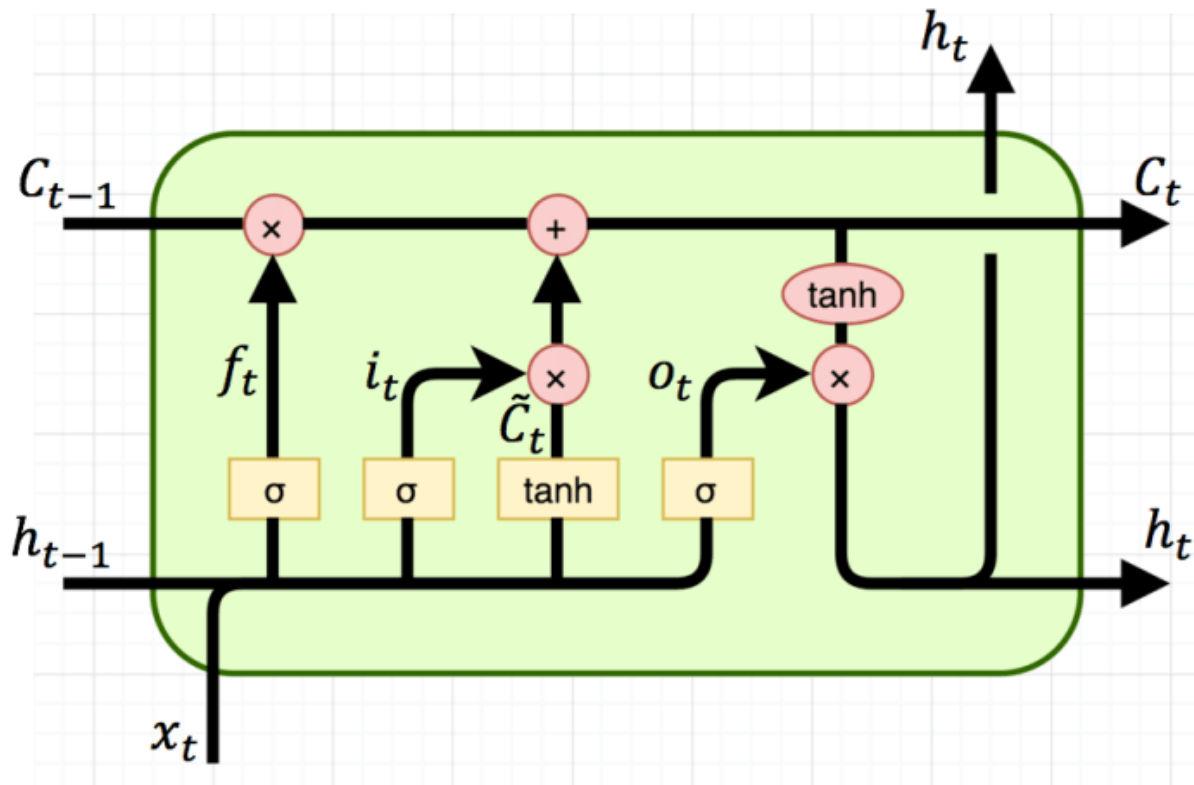
Character-level language model example

Vocabulary:
[h,e,l,o]

Example training
sequence:
“hello”



Long short term memory (LSTM) / GRU.



Translation

Dịch tiếng Anh
Translate English-Vietnamese



English

Translate



A



V

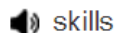
I would like to improve my English **skills**

Dictionary



[Basic] [Technical]

☒ On/Off



skills

skill /skil/

• danh từ

- sự khéo léo, sự khéo tay, sự tinh xảo; kỹ năng, kỹ xảo

Tiếng Việt

Tôi muốn cải thiện kỹ năng tiếng anh của mình

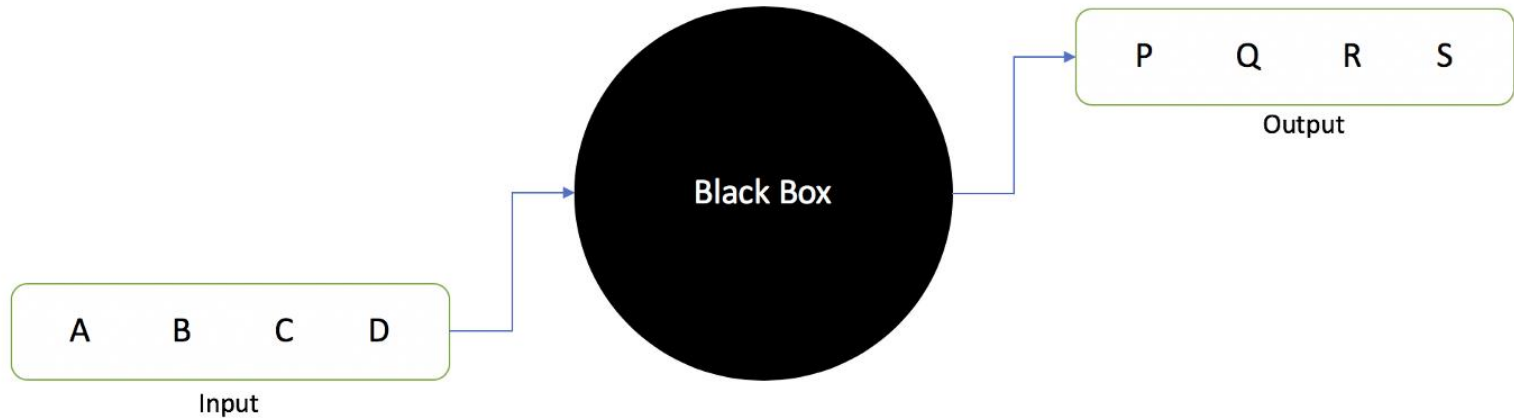
Example

☒ On/Off

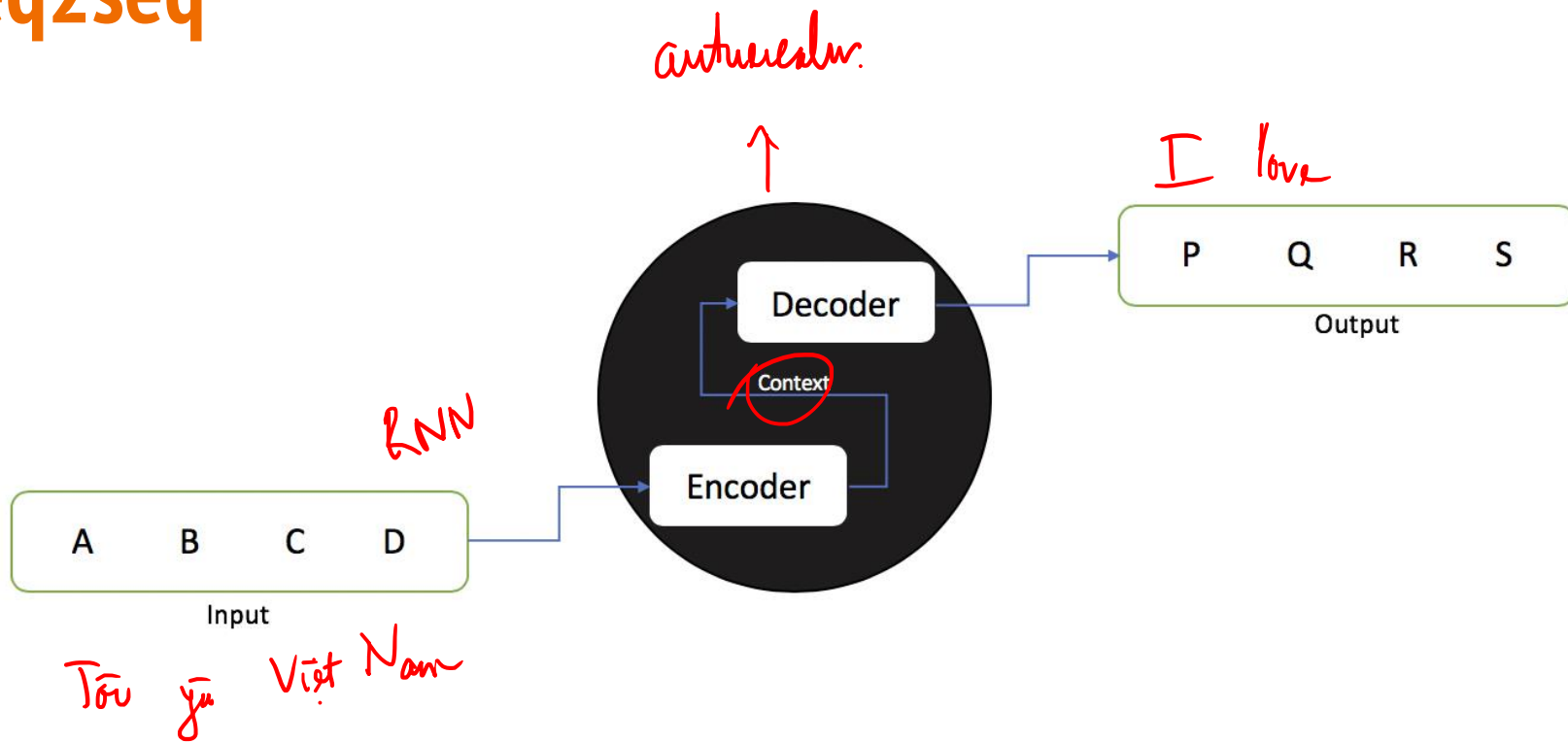
Next >>

- I'm expecting to gain more **skills** and knowledge, Sir.
- I don't have any **skills** at fishing.
- I humbly offer my dedication to practice the medical **skills** and knowledge I gained from 10 long years of study.

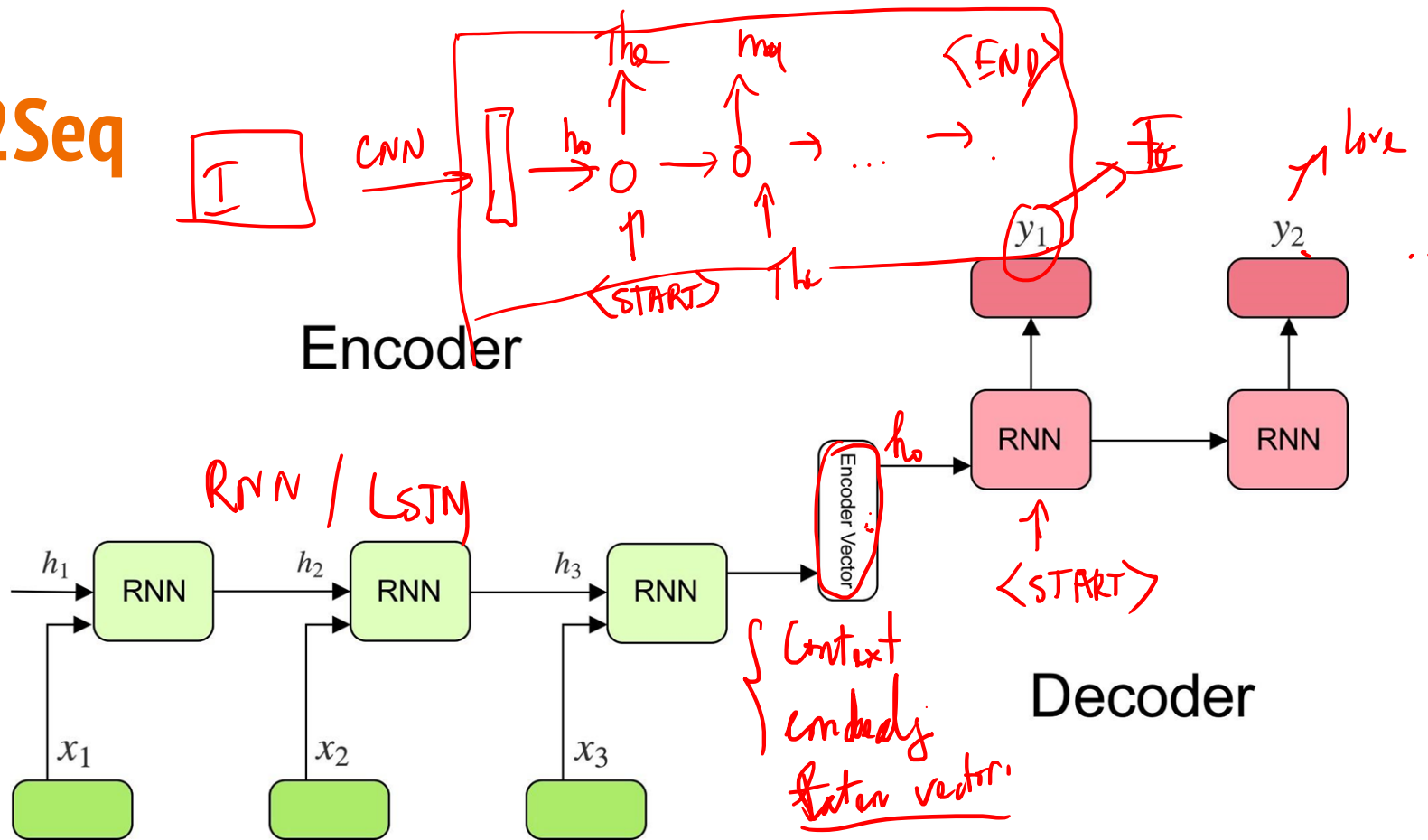
Model



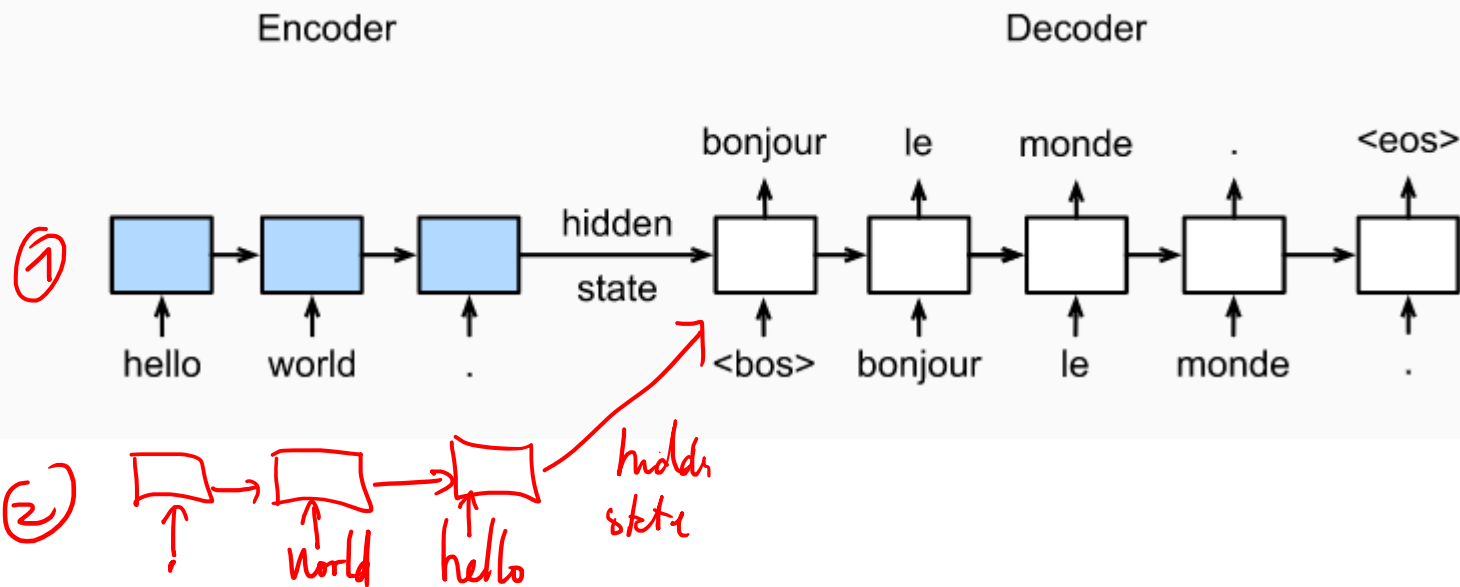
Seq2seq



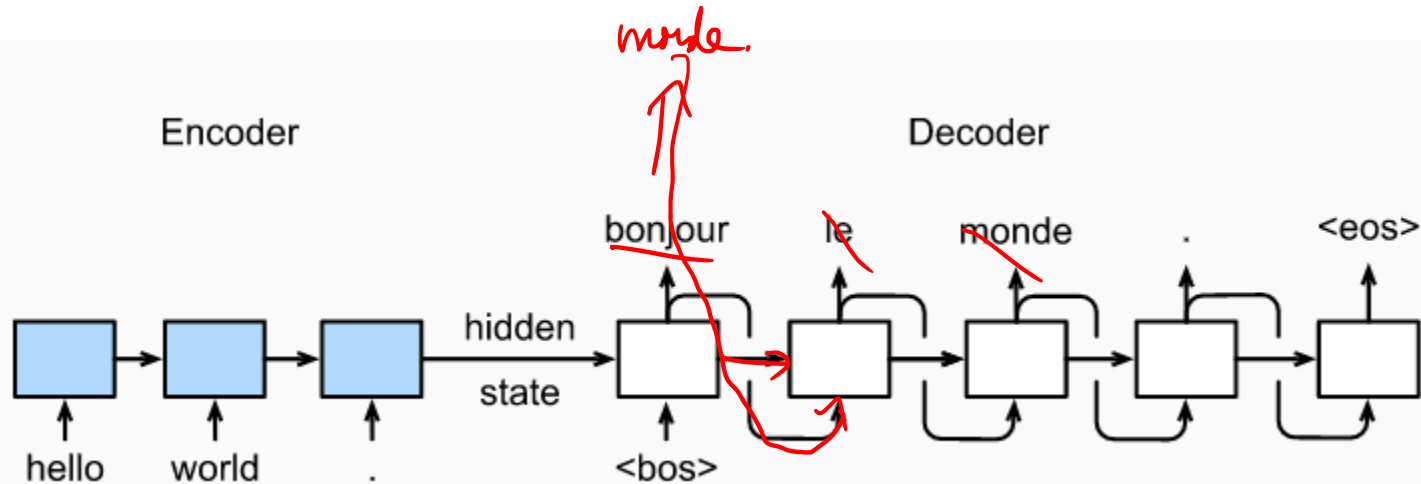
Seq2Seq



Seq2seq



Seq2seq - prediction



Greedy search

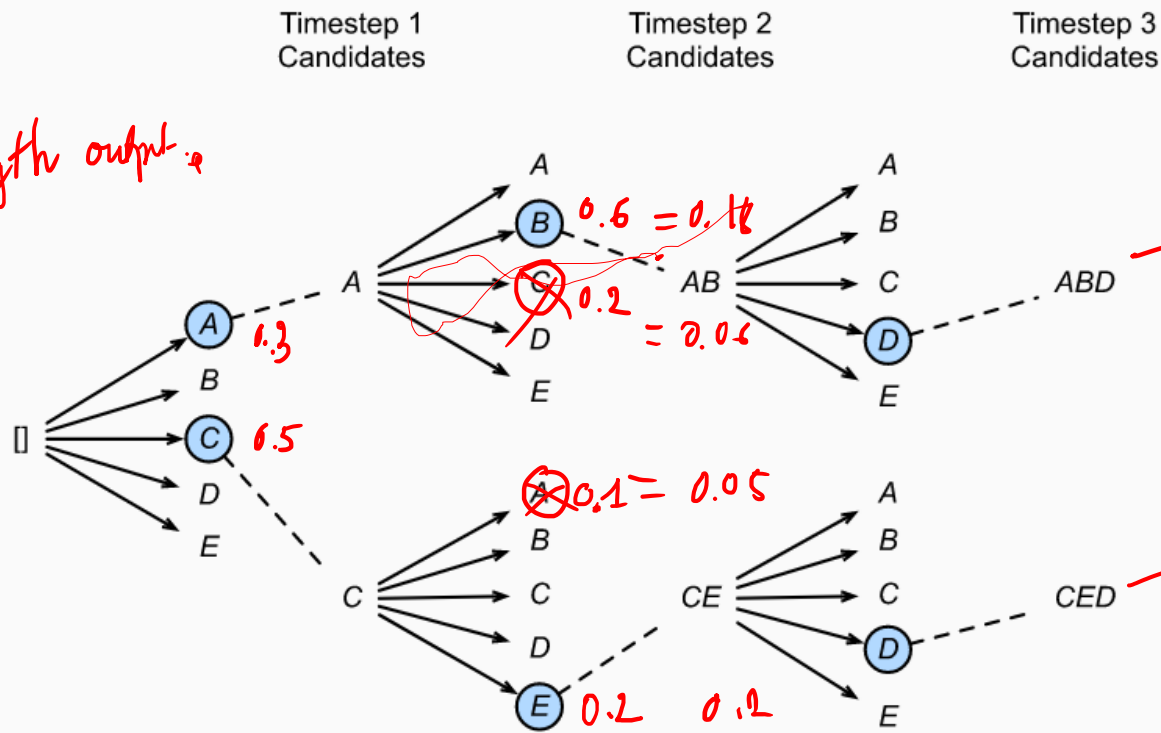
Timestep	1	2	3	4
A	0.5	0.1	0.2	0.0
B	0.2	0.4	0.2	0.2
C	0.2	0.3	0.4	0.2
<eos>	0.1	0.2	0.2	0.6

Beam search, $k=2$, hyper-parameter k giữ k nhánh tốt nhất

2^n , n là length output

$$\sqrt[n]{x_1 \cdot x_2 \cdot x_3 \cdot x_4 \cdot x_5}$$

$$\sqrt[n]{y_1 \cdot y_2 \cdot y_3}$$



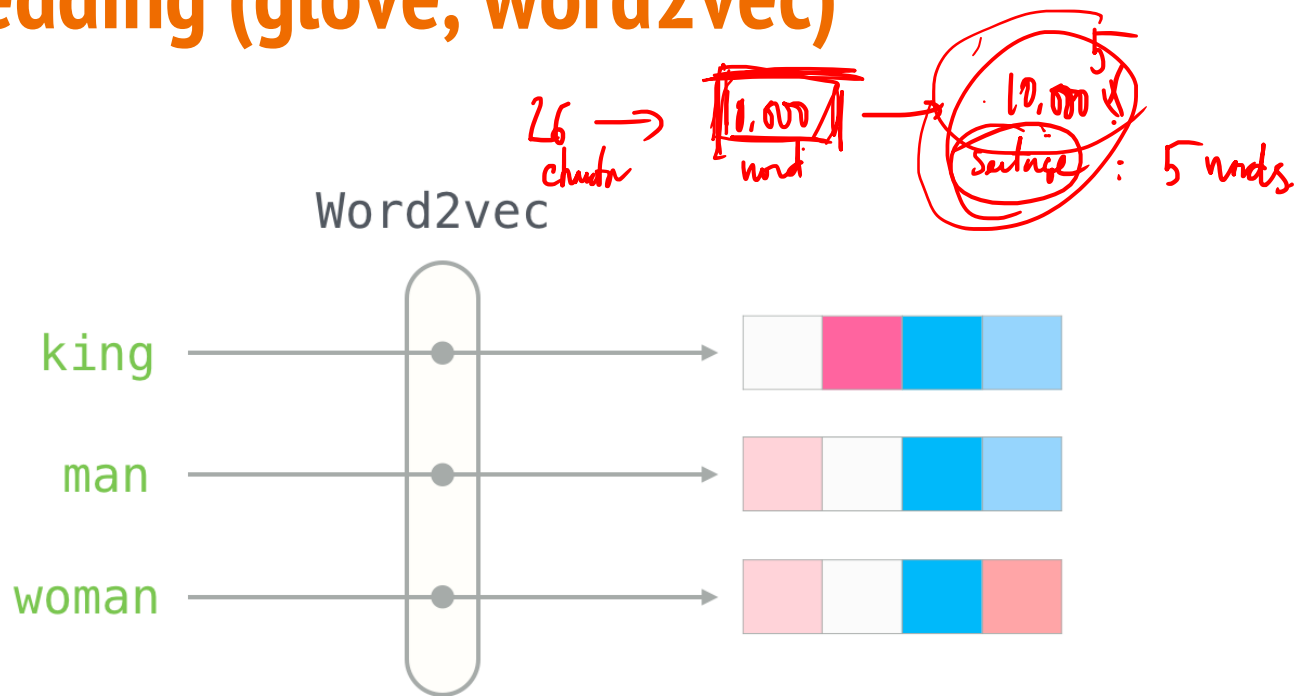
ABPA - ABDA

$\langle \text{END} \rangle$

k nhánh

$\langle \text{END} \rangle$

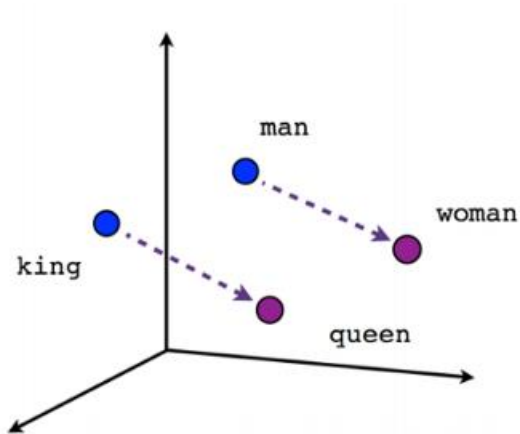
Word embedding (glove, word2vec)



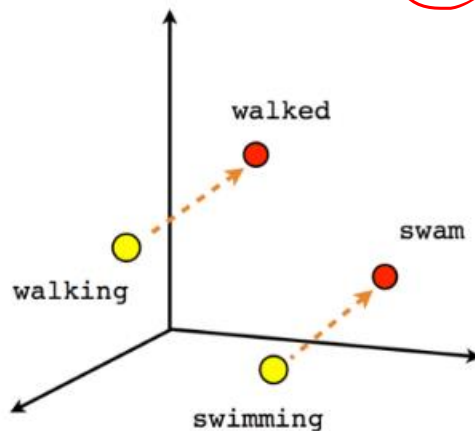
Word semantic

que - kij + man \approx woman

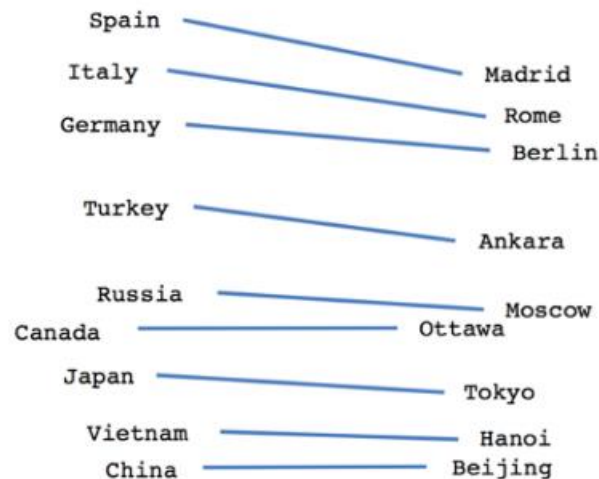
elmo
BERT



Male-Female



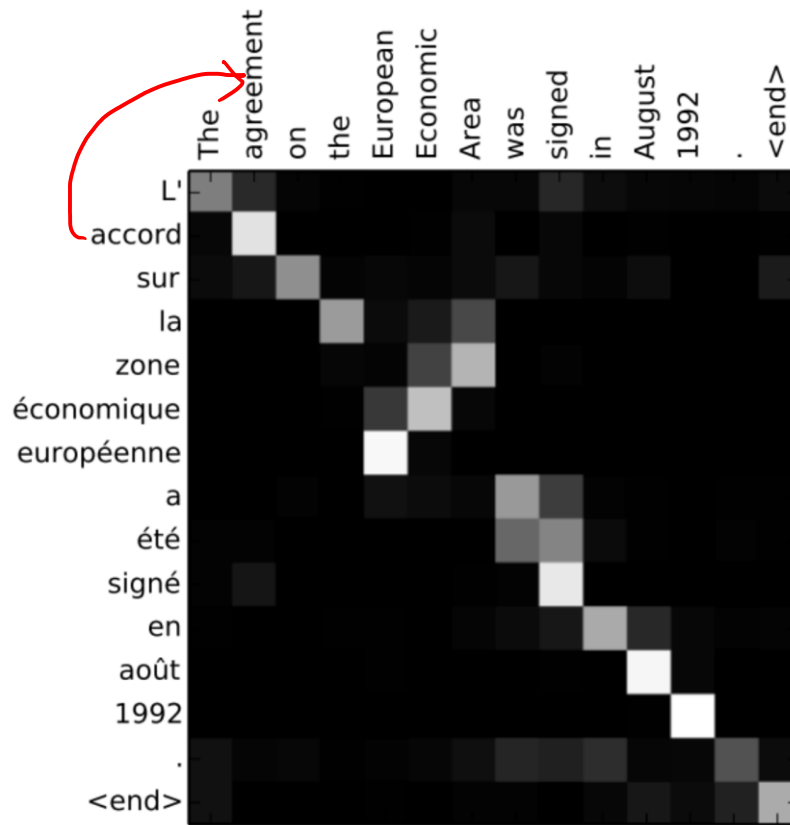
Verb tense



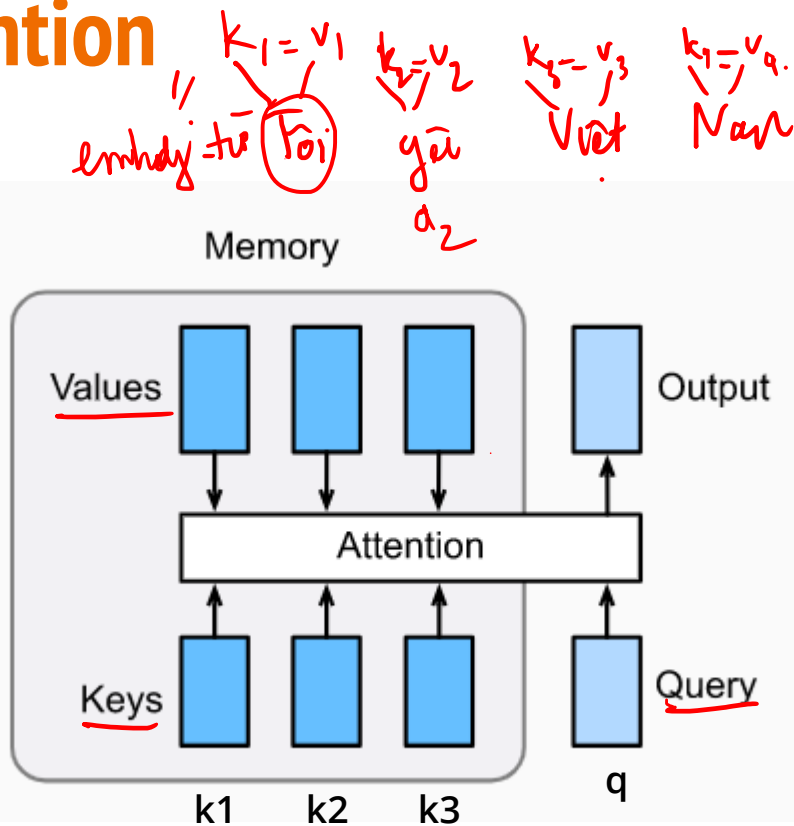
Country-Capital

Attention - motivation

The food is god → pos / neg



Attention



$$\underline{V} = \sum v_i \cdot b_i$$

you love

$$a_i = 2$$

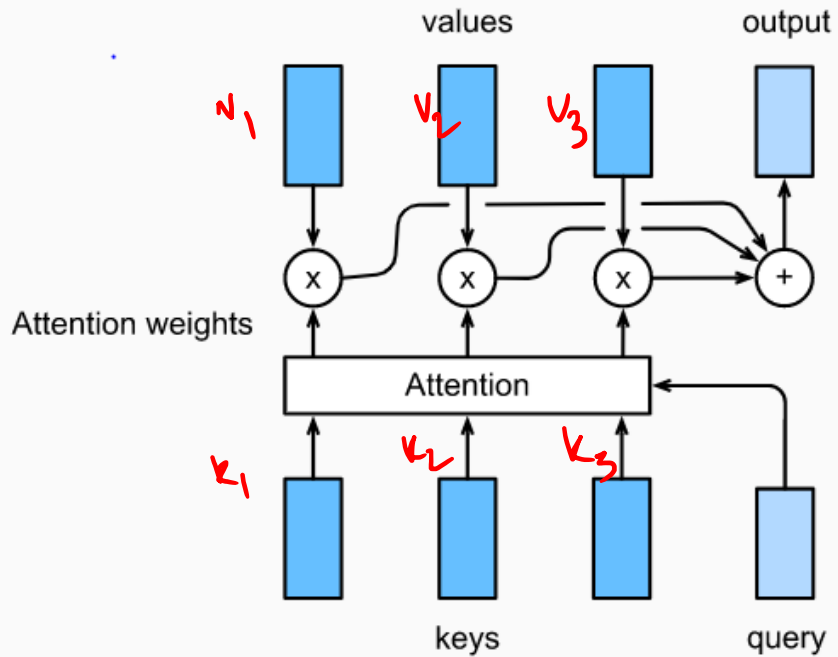
$$\underline{V} = \sum v_i \cdot a_i$$

$\sum a_i > 1$

$\alpha(\mathbf{q}, \mathbf{k}_i)$: độ quan trọng của từ i với từ cần dịch.

$$b_i = \frac{\exp(a_i)}{\sum_j \exp(a_j)}, \mathbf{b} = [b_1, \dots, b_n]^T$$

Attention

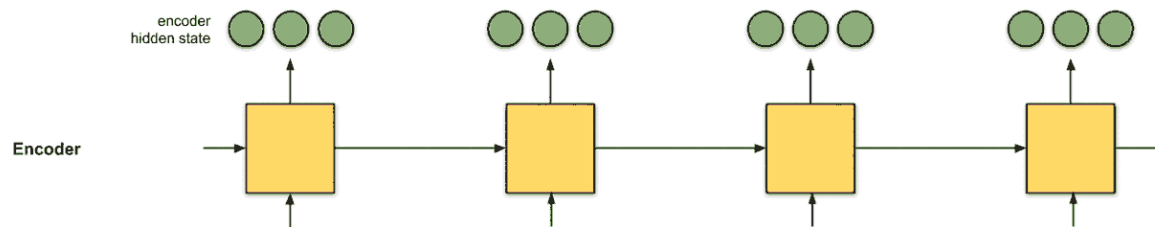
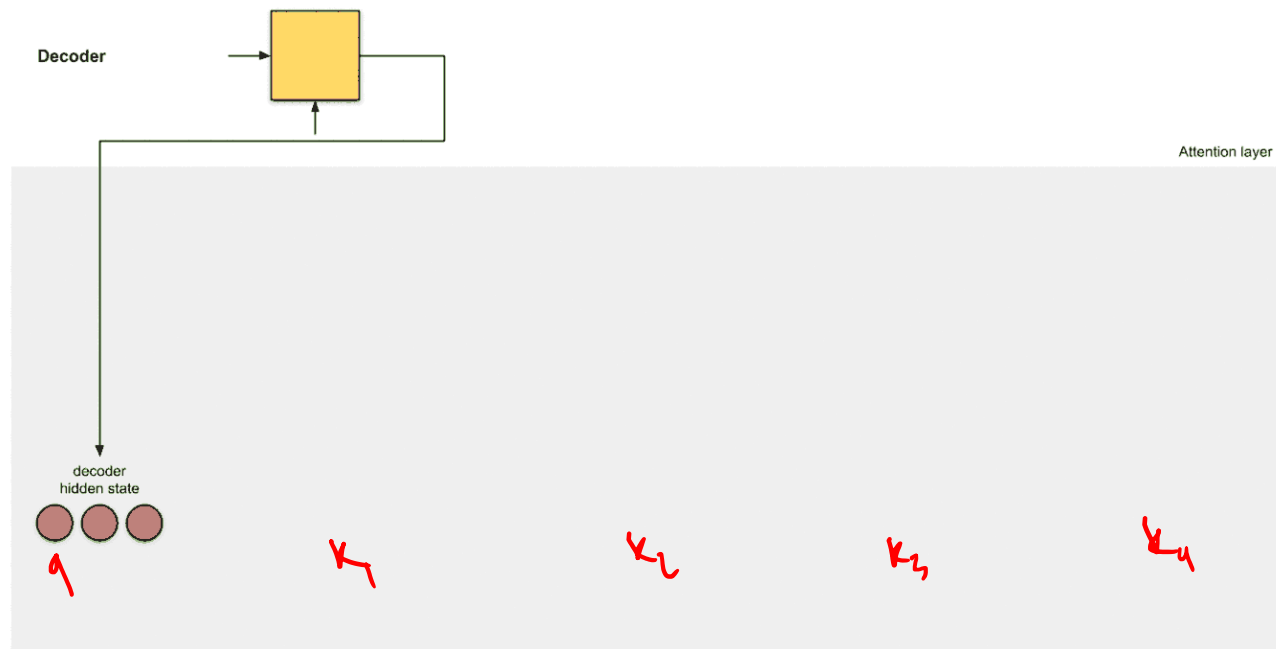


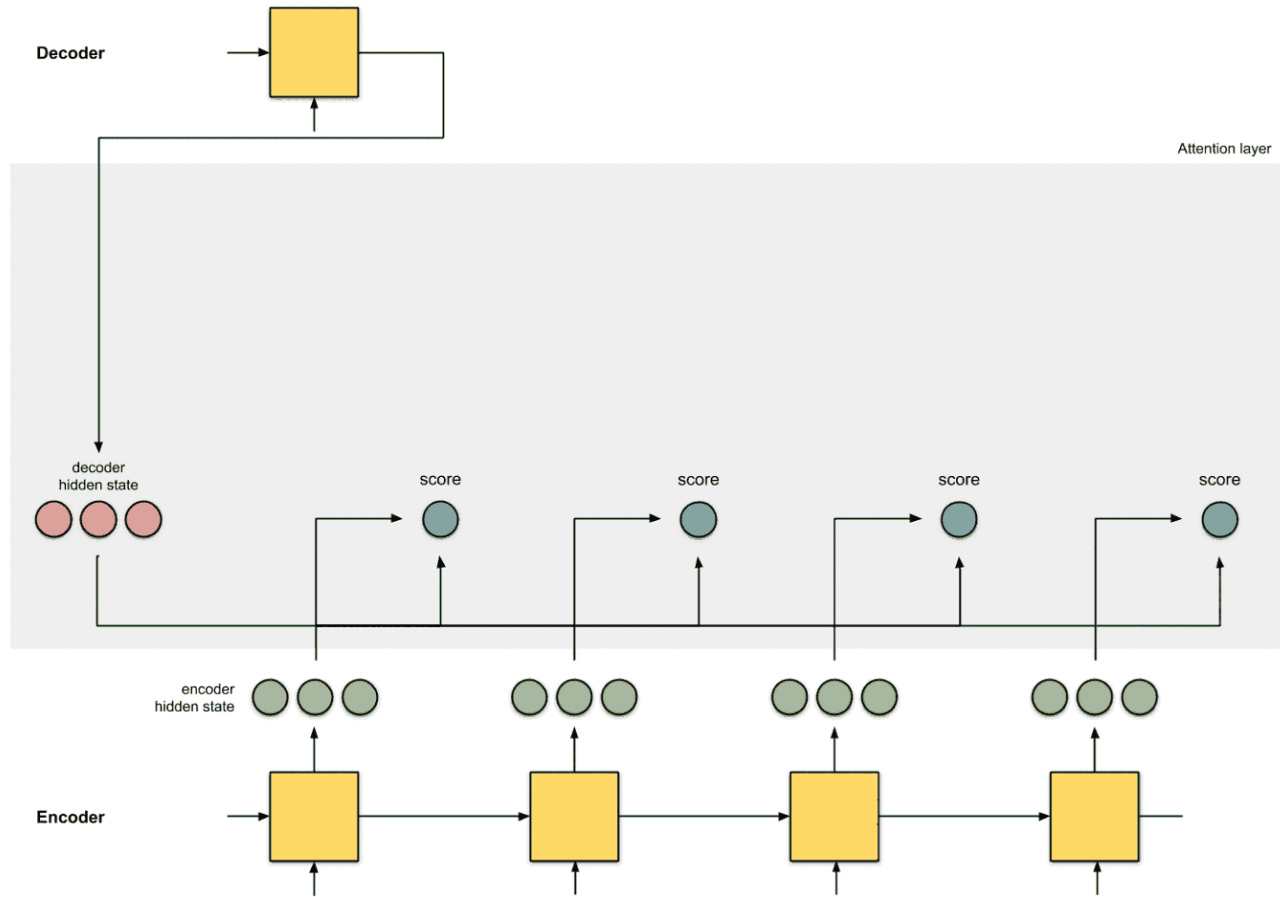
Attention function

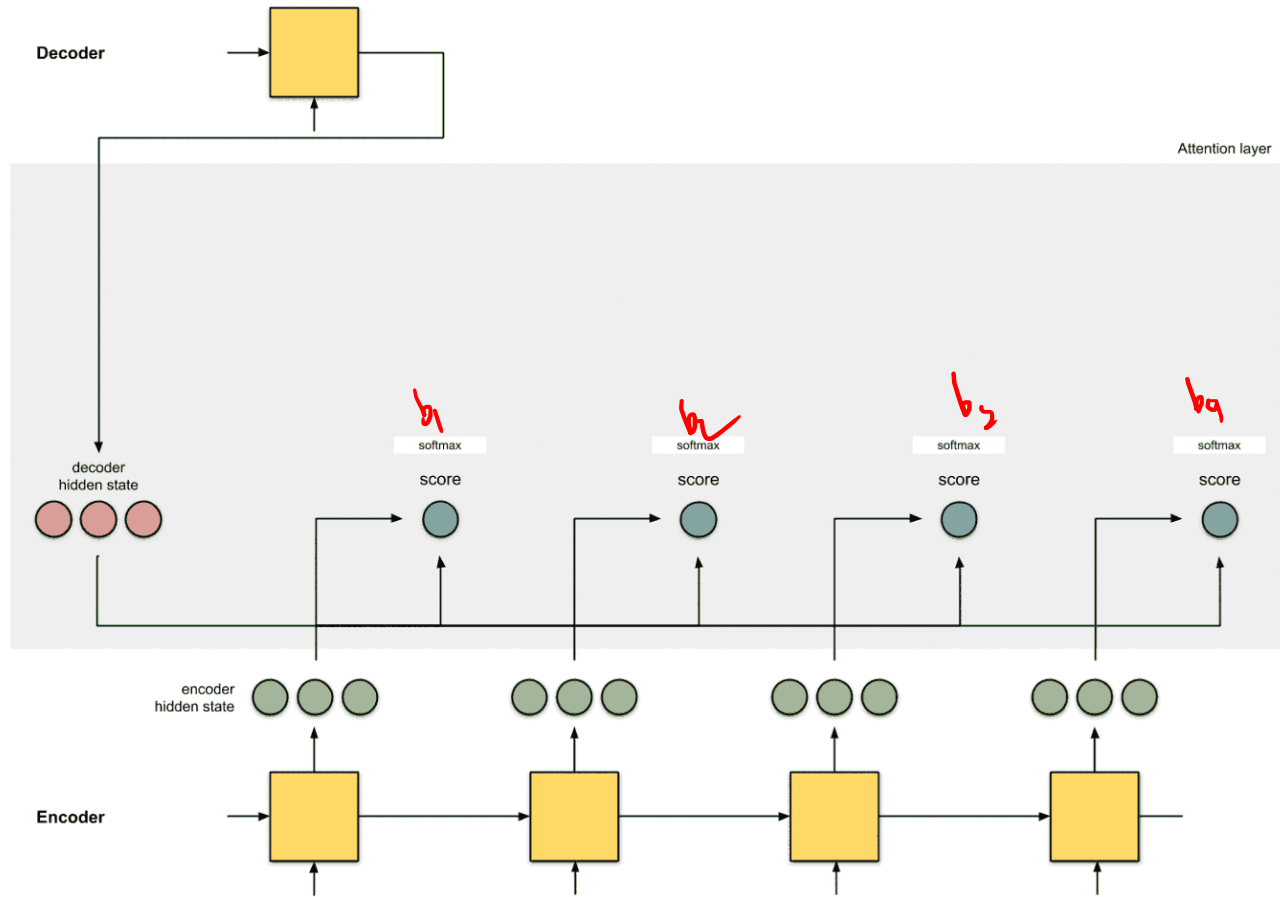
$$a_i = \alpha(q_i, k_i) \quad \alpha: q_i, k_i \rightarrow a_i$$

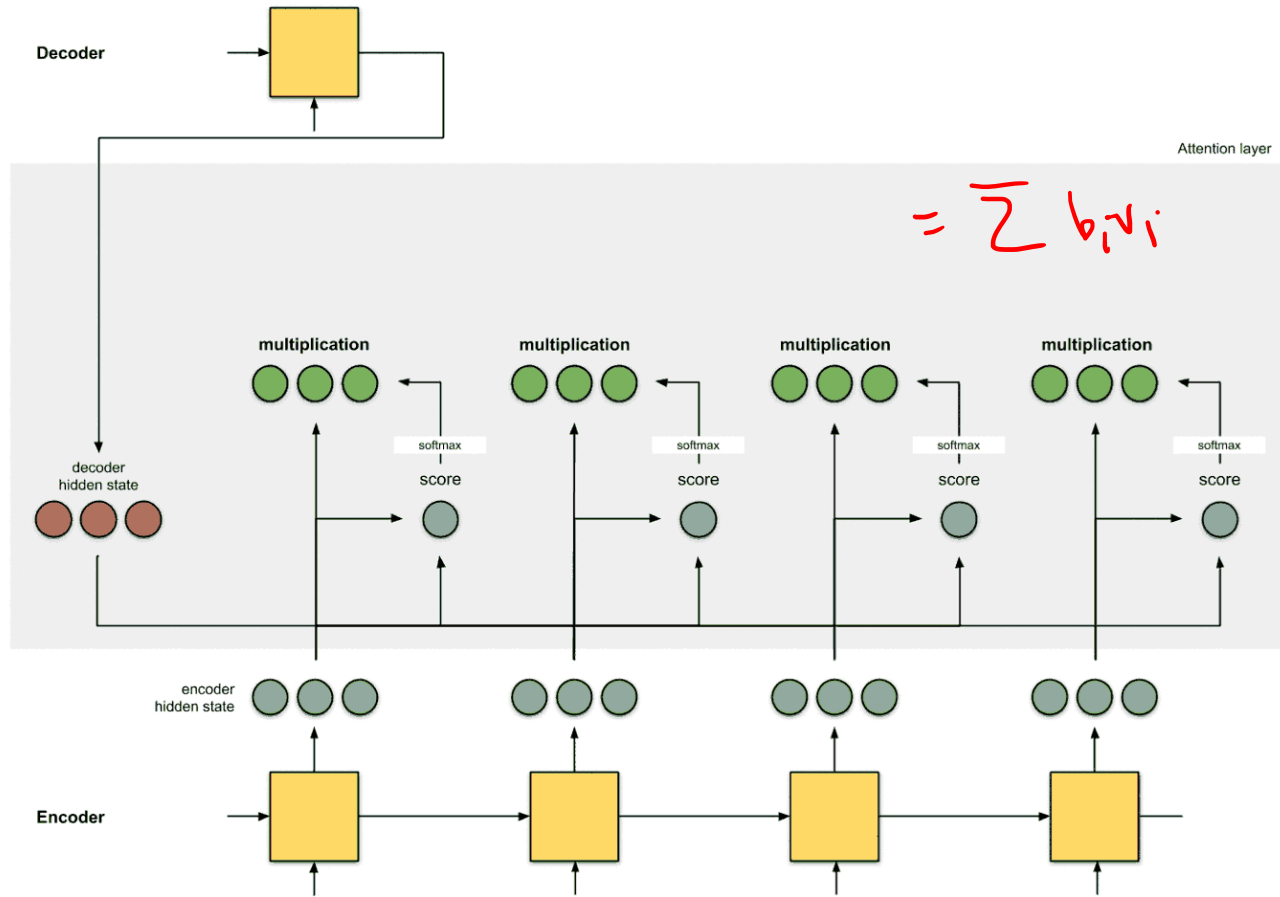
Name	Alignment score function
Content-base attention	$\text{score}(\mathbf{s}_t, \mathbf{h}_i) = \text{cosine}[\mathbf{s}_t, \mathbf{h}_i]$
Additive(*)	$\text{score}(\mathbf{s}_t, \mathbf{h}_i) = \mathbf{v}_a^\top \tanh(\mathbf{W}_a[\mathbf{s}_t; \mathbf{h}_i])$
Location-Base	$\alpha_{t,i} = \text{softmax}(\mathbf{W}_a \mathbf{s}_t)$ Note: This simplifies the softmax alignment to only depend on the target position.
General	$\text{score}(\mathbf{s}_t, \mathbf{h}_i) = \mathbf{s}_t^\top \mathbf{W}_a \mathbf{h}_i$ where \mathbf{W}_a is a trainable weight matrix in the attention layer.
Dot-Product	$\text{score}(\mathbf{s}_t, \mathbf{h}_i) = \mathbf{s}_t^\top \mathbf{h}_i$
Scaled Dot-Product(^)	$\text{score}(\mathbf{s}_t, \mathbf{h}_i) = \frac{\mathbf{s}_t^\top \mathbf{h}_i}{\sqrt{n}}$ Note: very similar to the dot-product attention except for a scaling factor; where n is the dimension of the source hidden state.

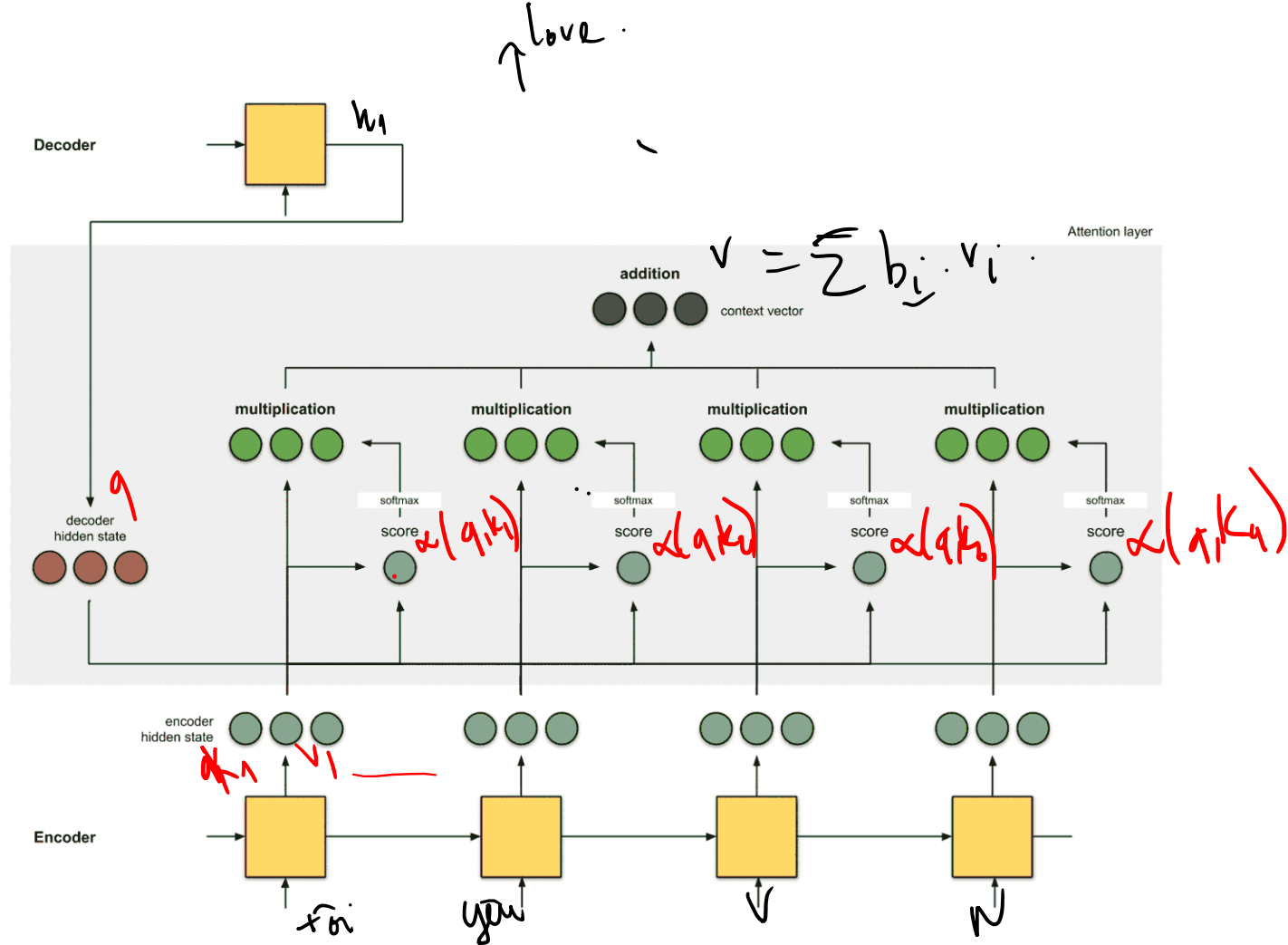
Encoder





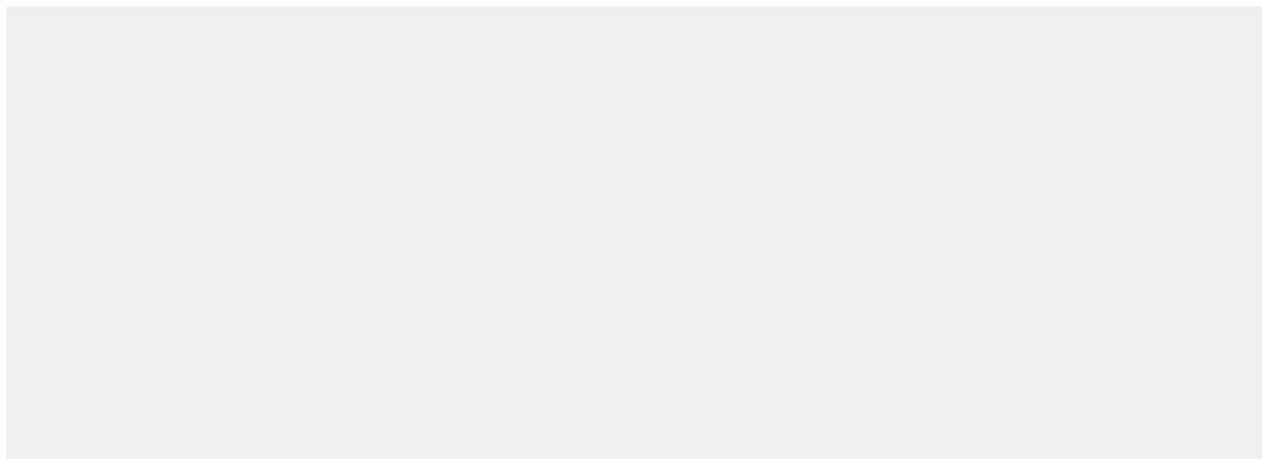






Self attention : hoc high-level feature củ câu.

Self-attention



input #1

1	0	1	0
---	---	---	---

input #2

0	2	0	2
---	---	---	---

input #3

1	1	1	1
---	---	---	---

Self-attention

LSA . 24 .

K, Q, V
 128
 $K, 256$

\downarrow

128

256 .

input #1

1	0	1	0
---	---	---	---

input #2

0	2	0	2
---	---	---	---

input #3

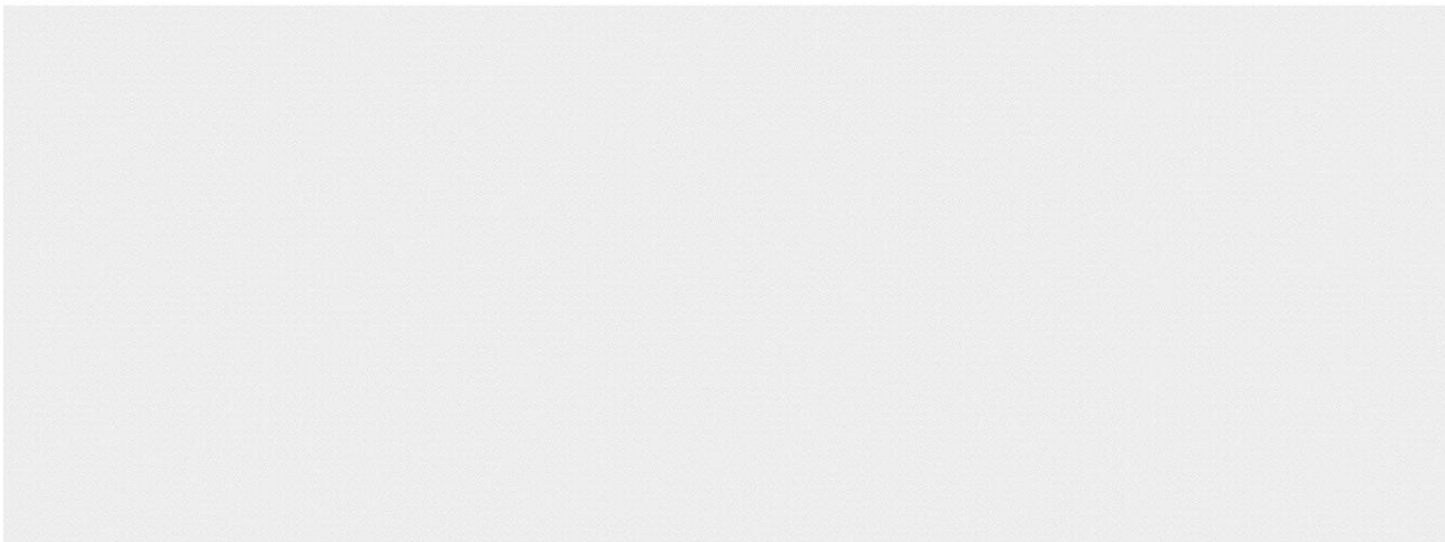
1	1	1	1
---	---	---	---

I

love

you

Self-attention



input #1

1	0	1	0
---	---	---	---

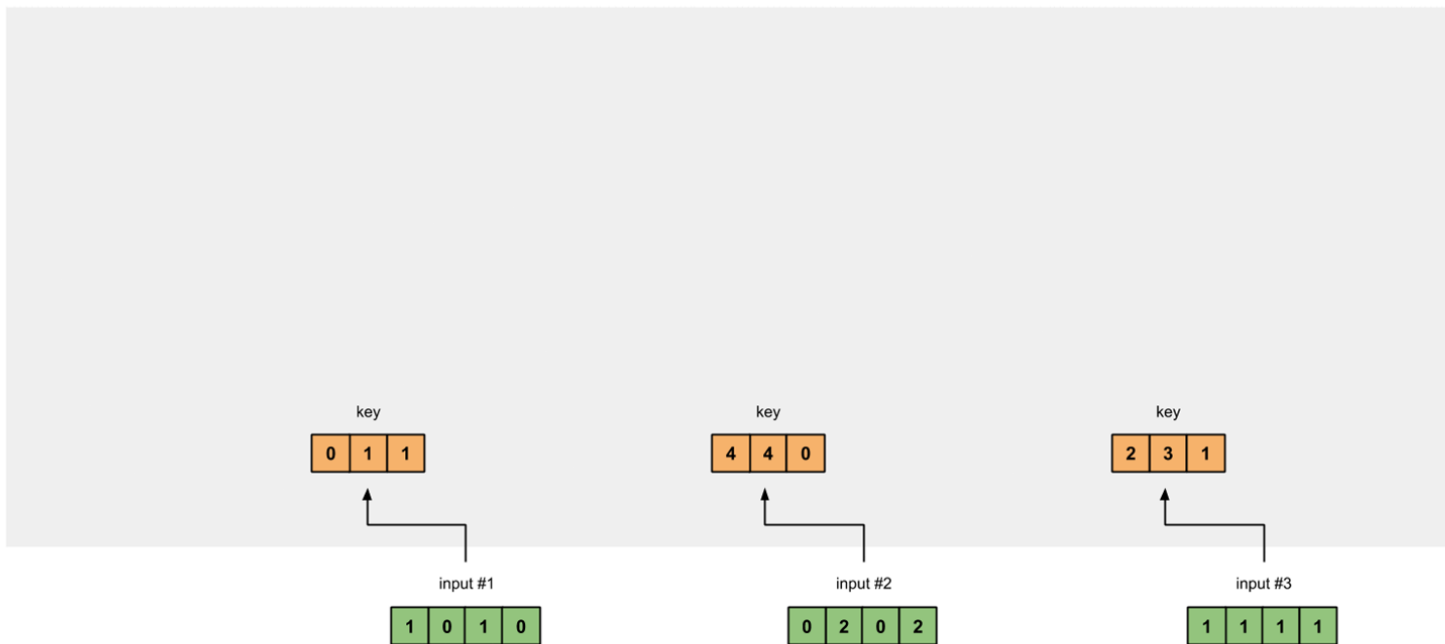
input #2

0	2	0	2
---	---	---	---

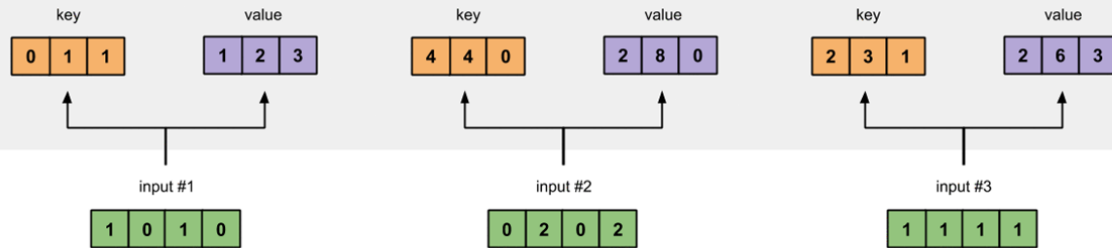
input #3

1	1	1	1
---	---	---	---

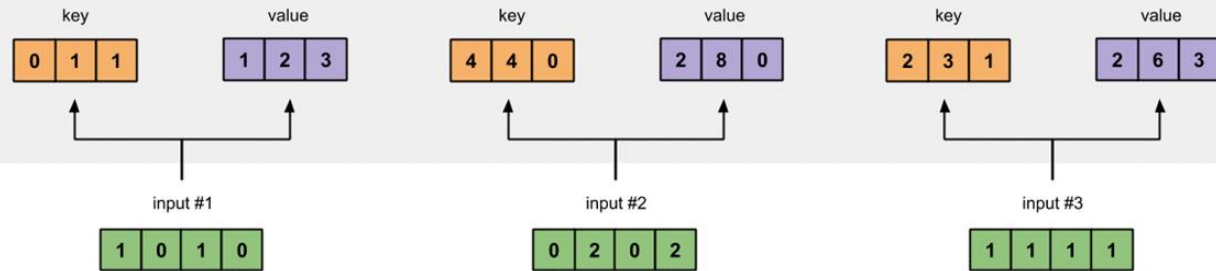
Self-attention



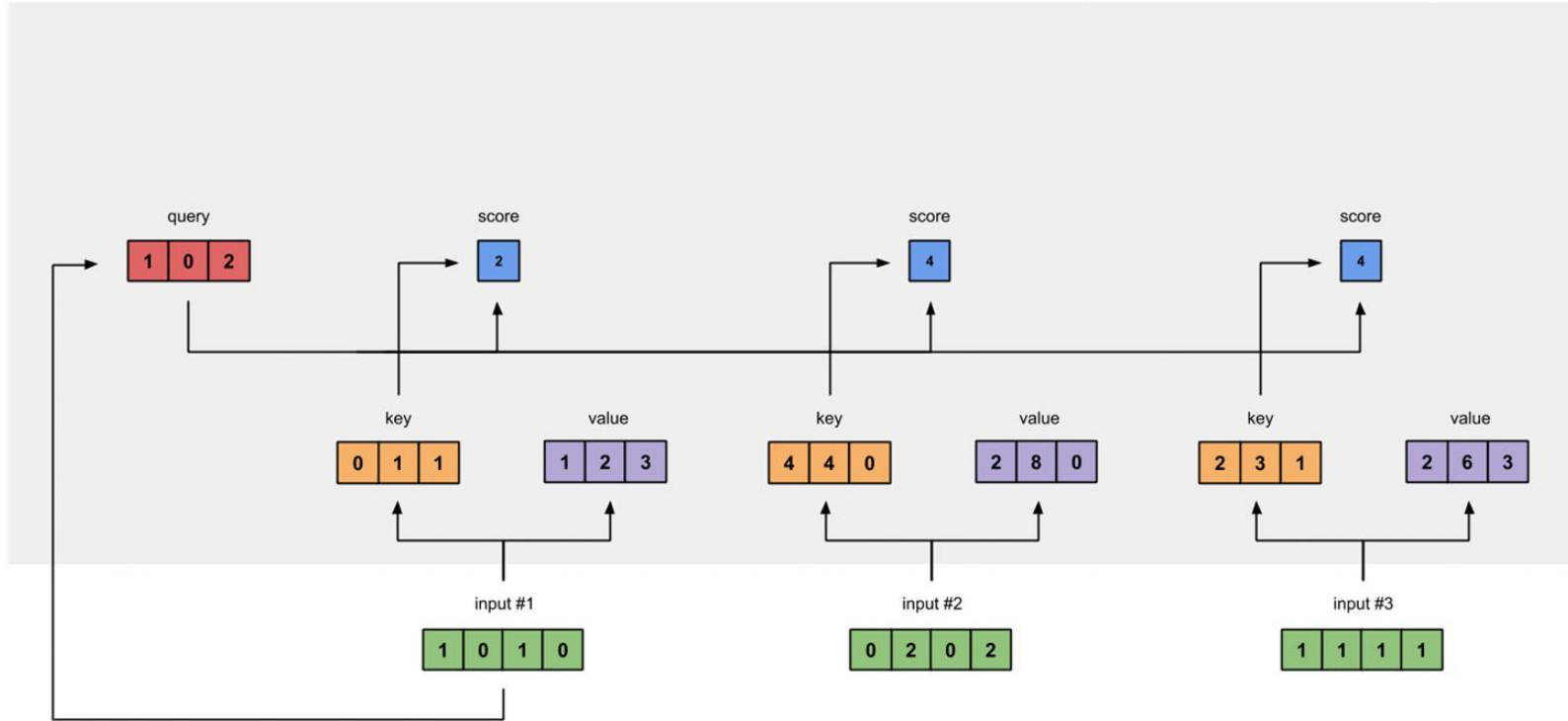
Self-attention



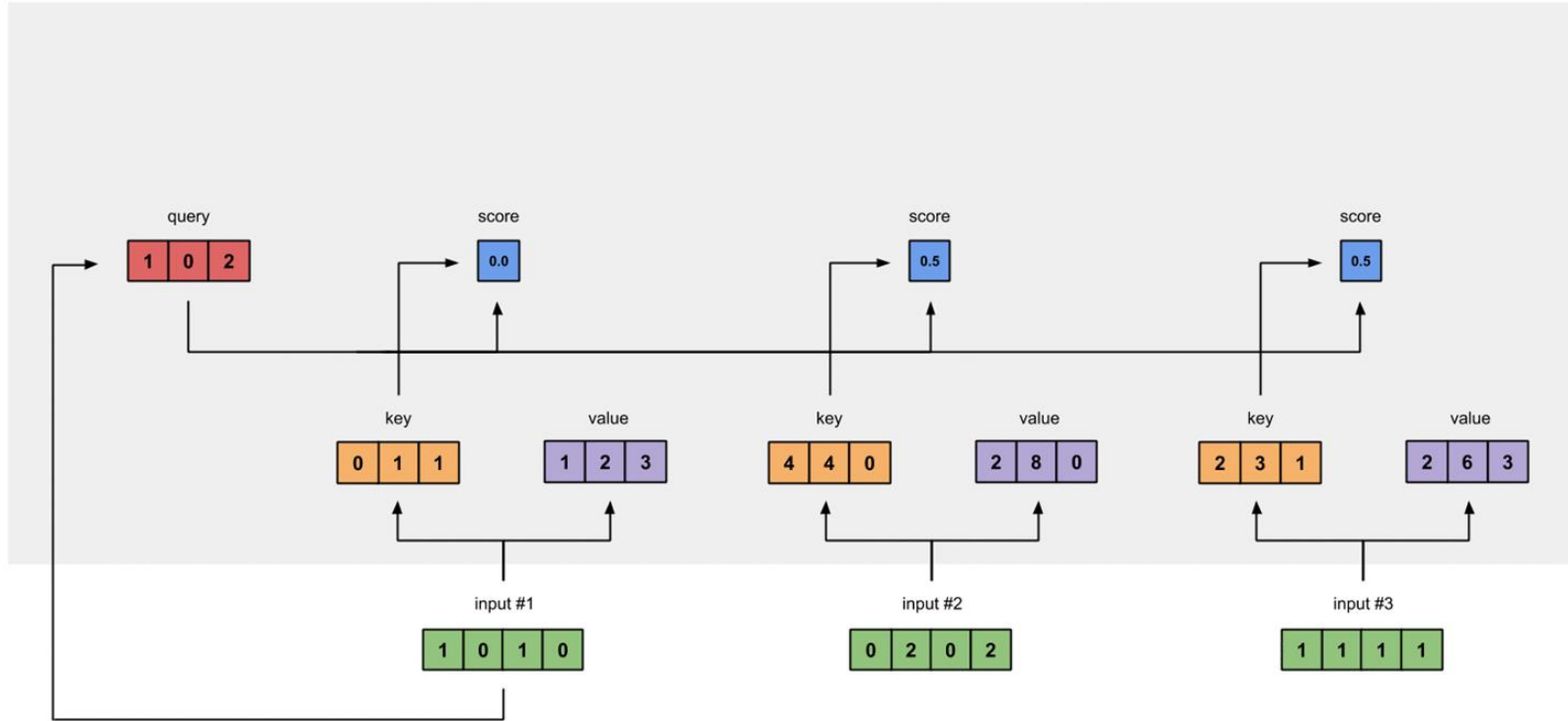
Self-attention



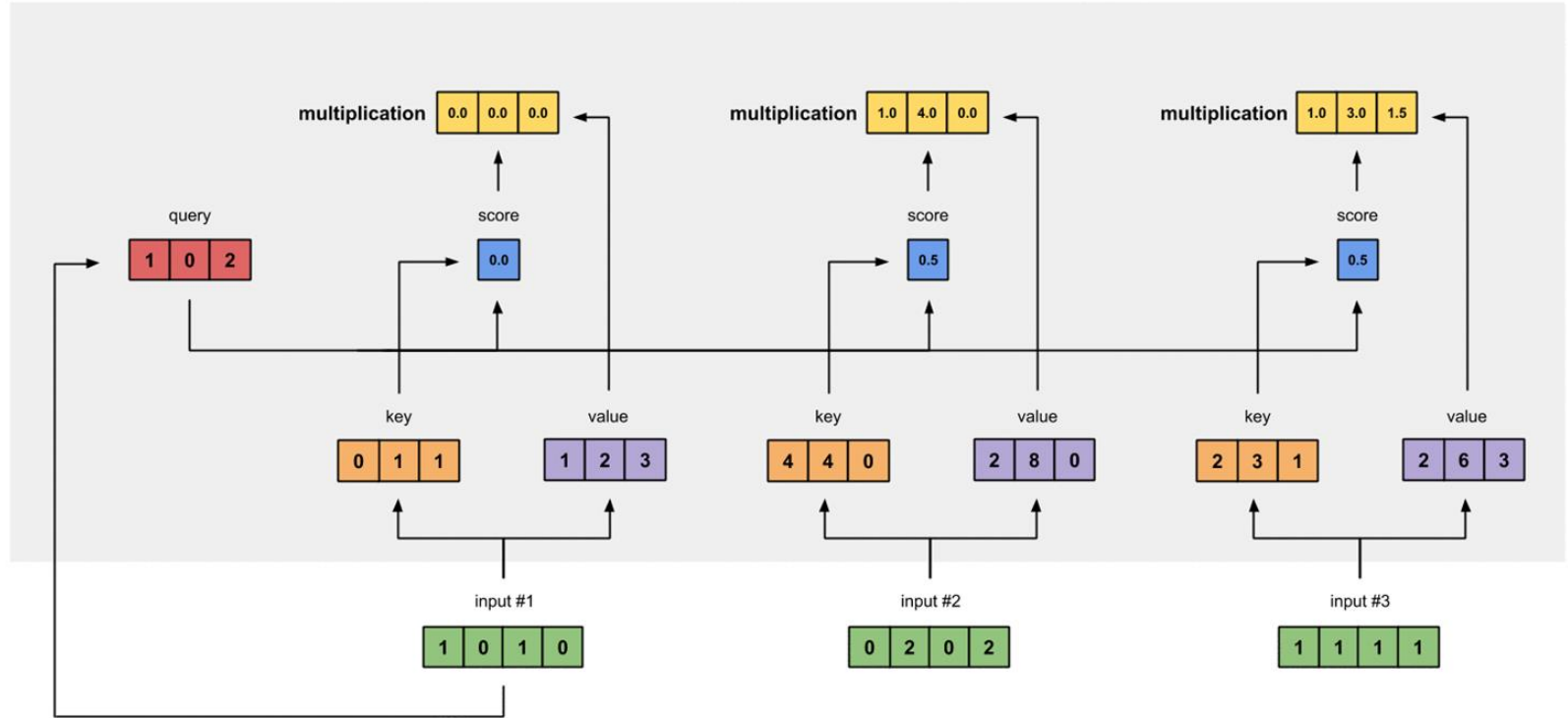
Self-attention



Self-attention



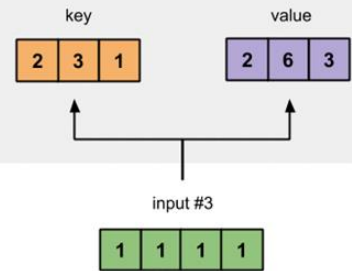
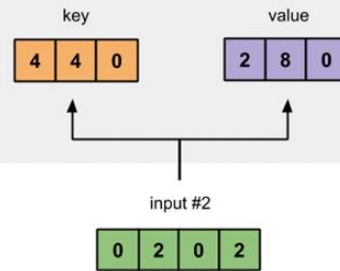
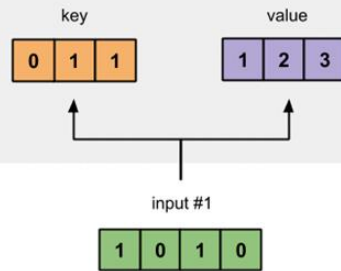
Self-attention



Self-attention

output #1
addition

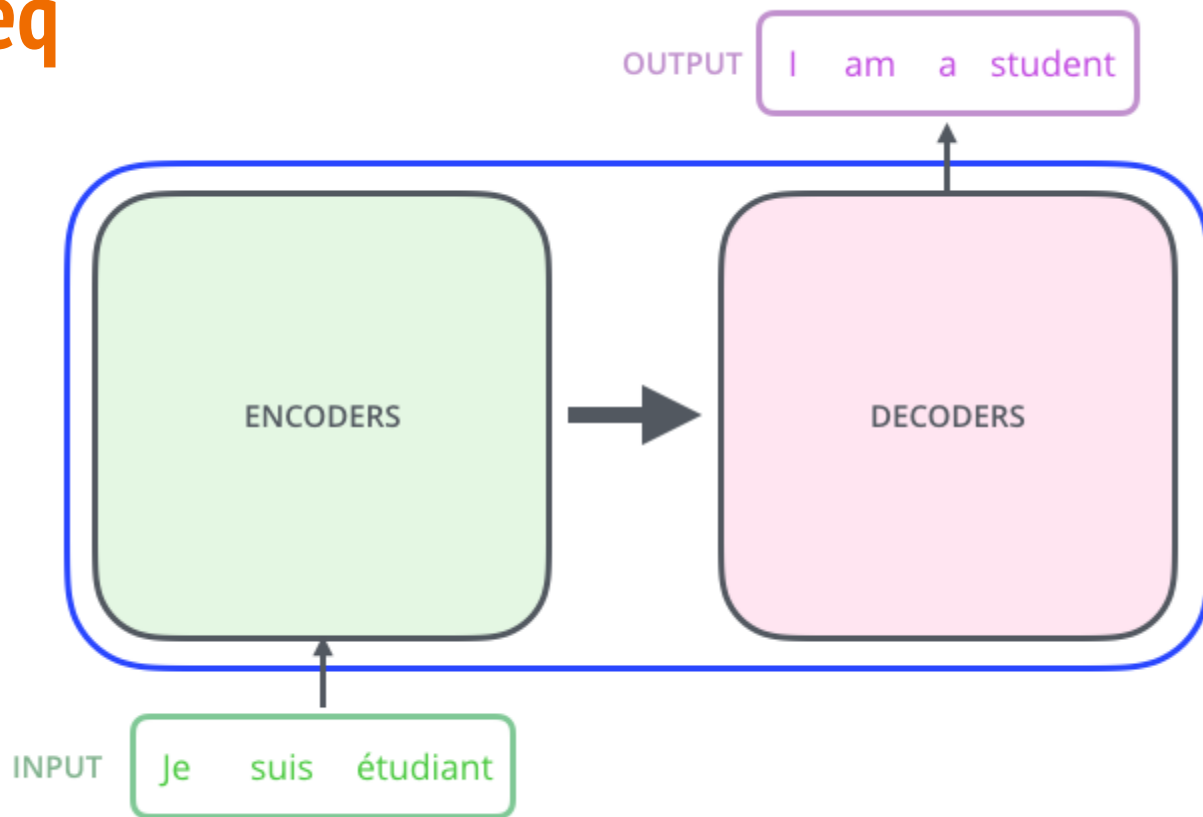
2.0	7.0	1.5
-----	-----	-----



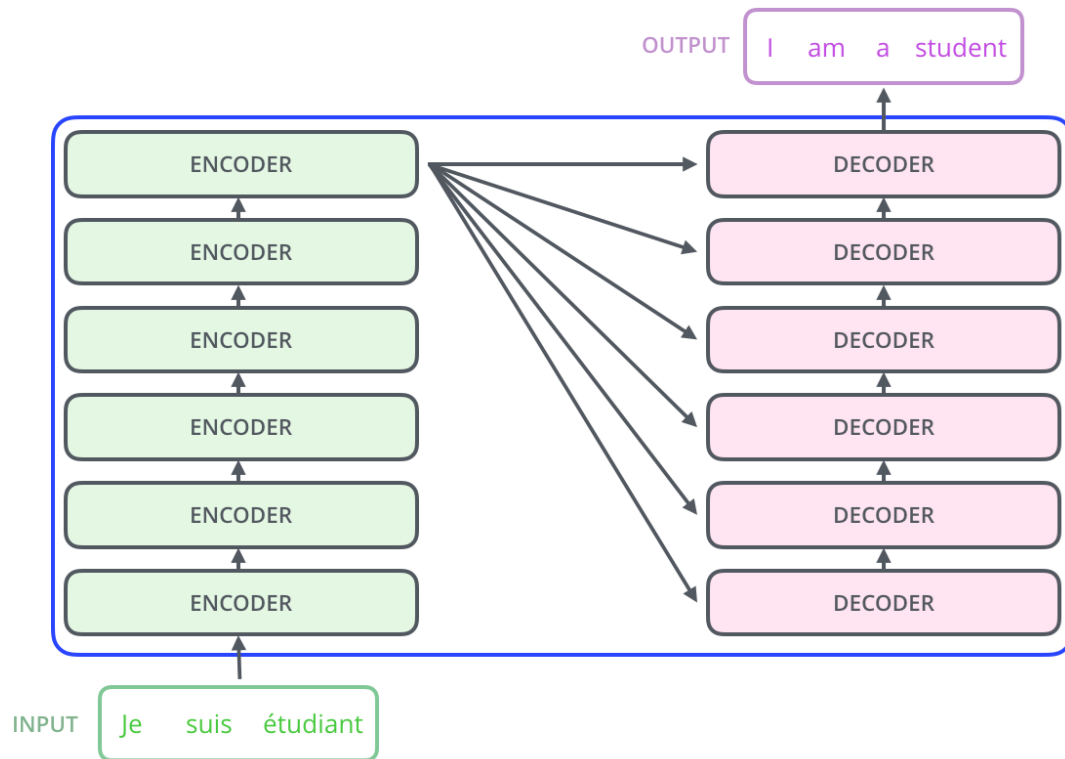
Transformer

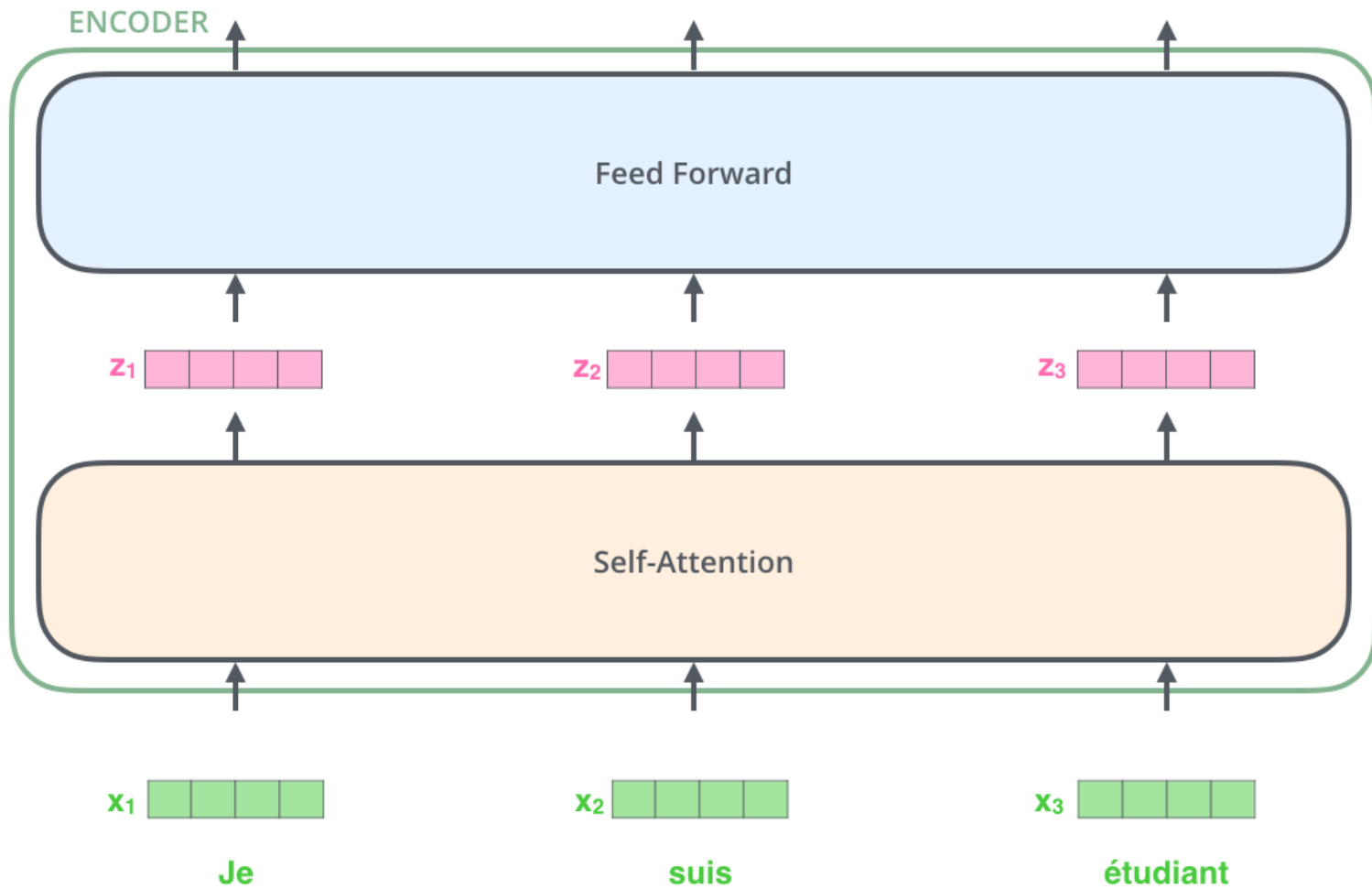


Seq2seq



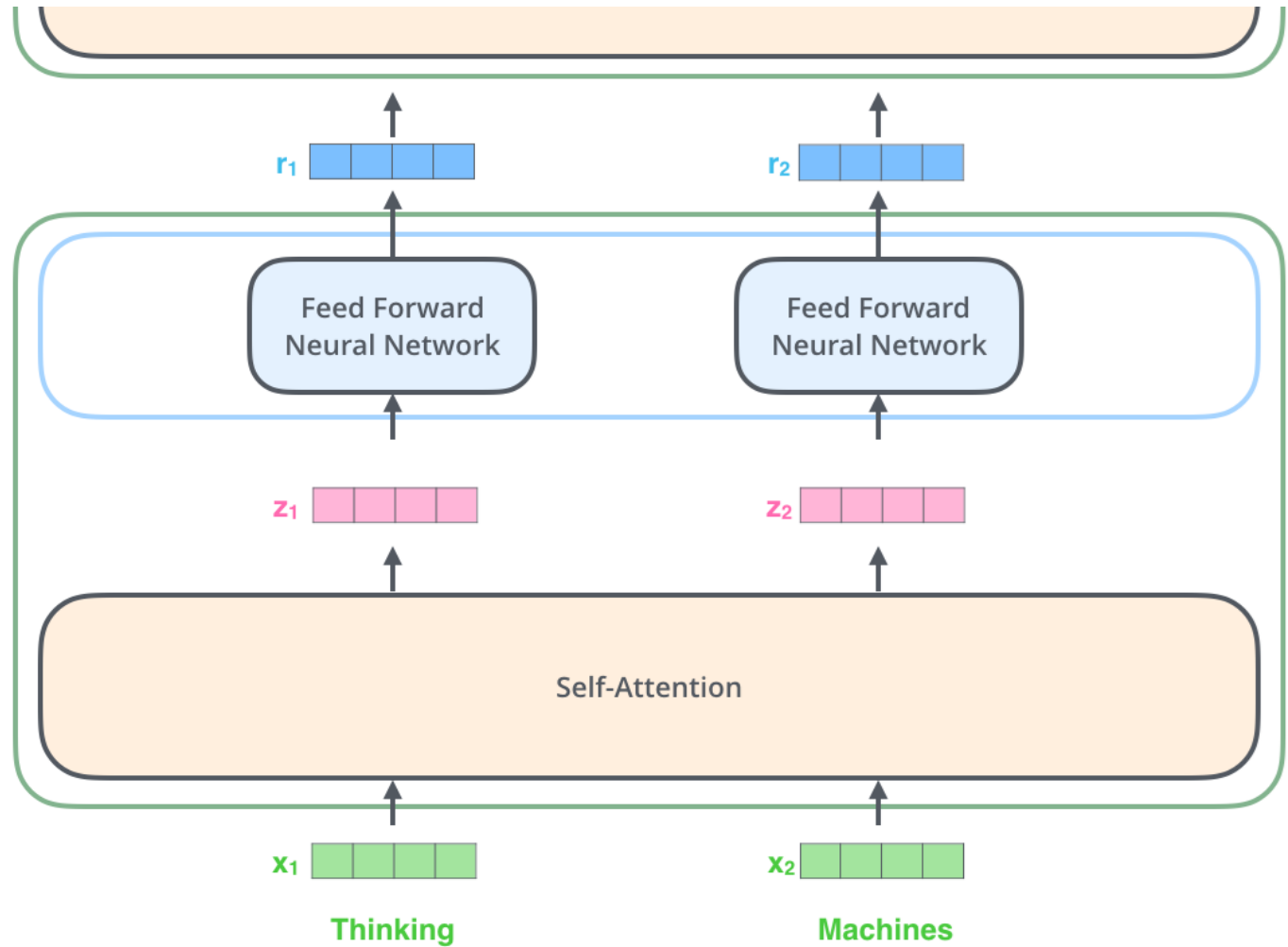
Transformer

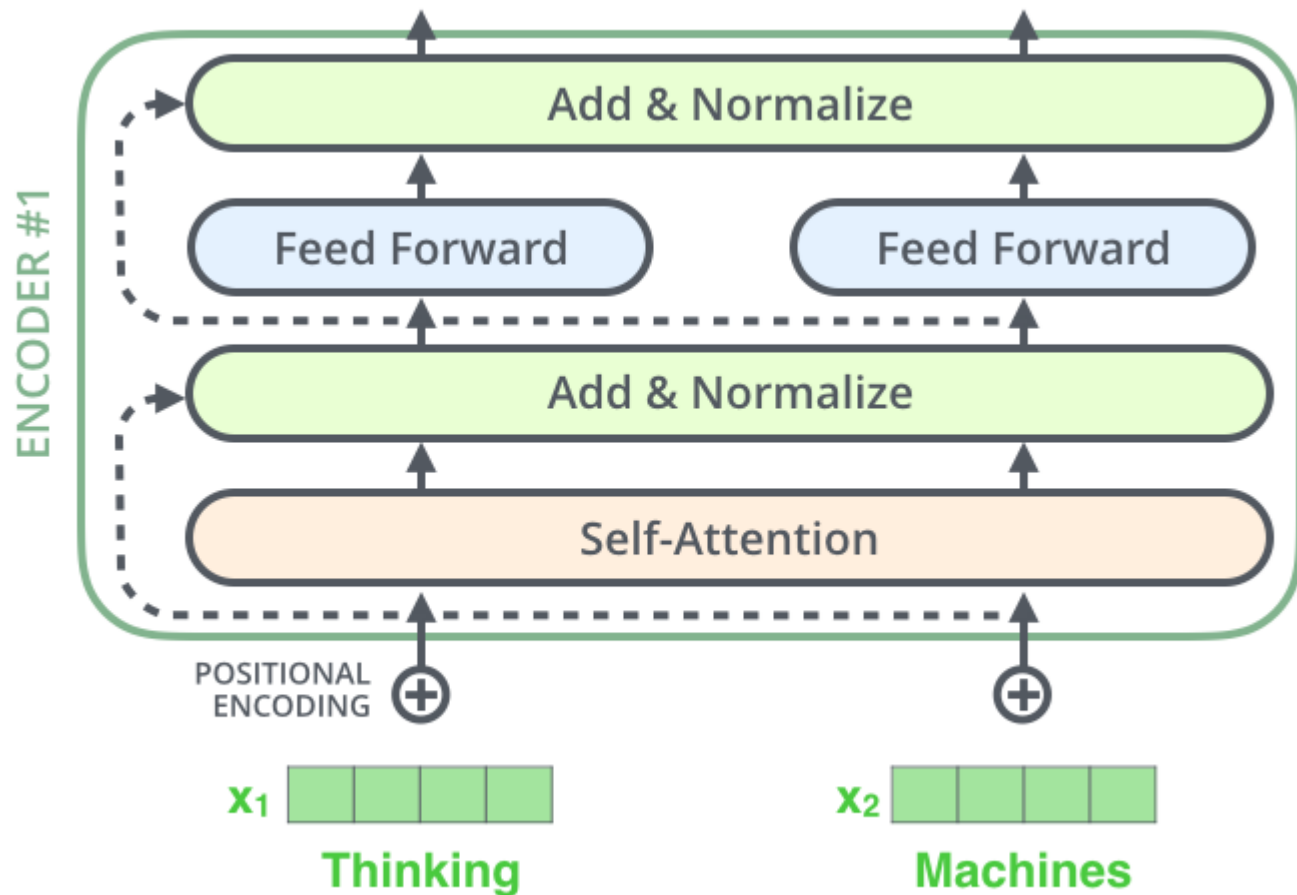


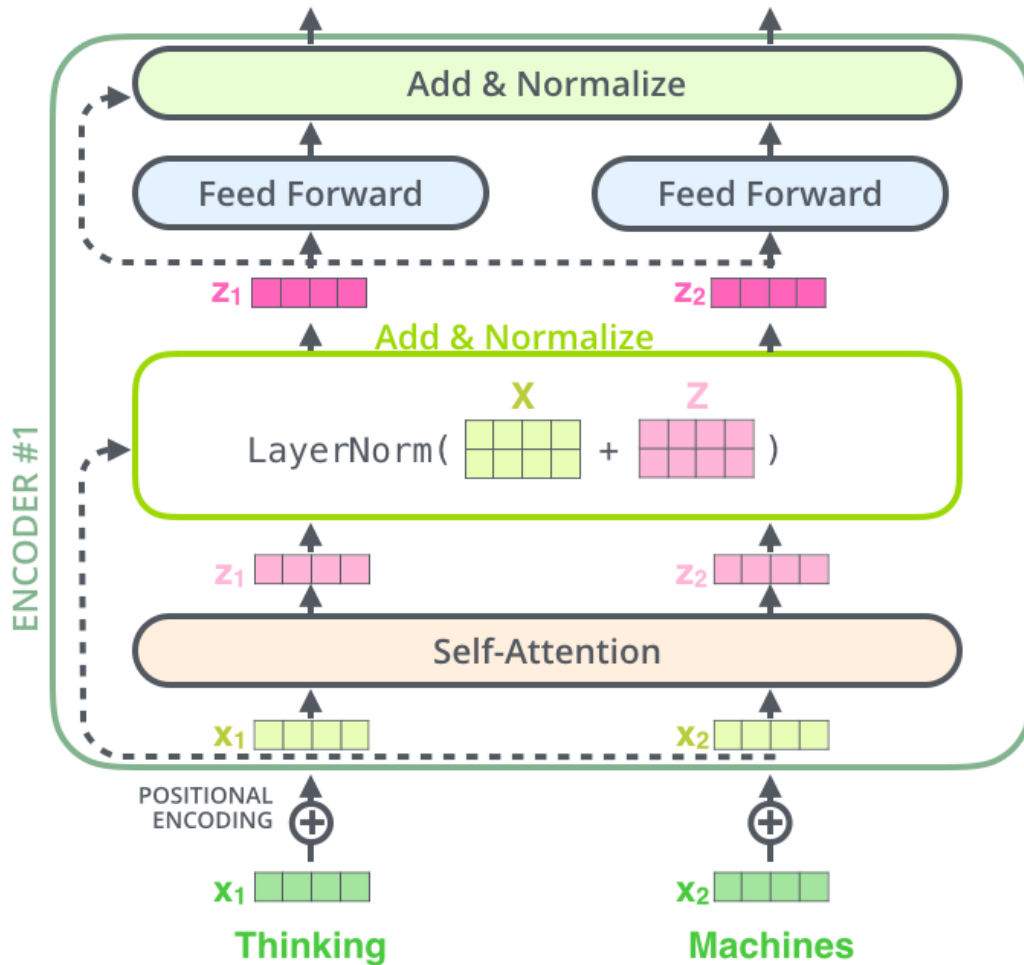


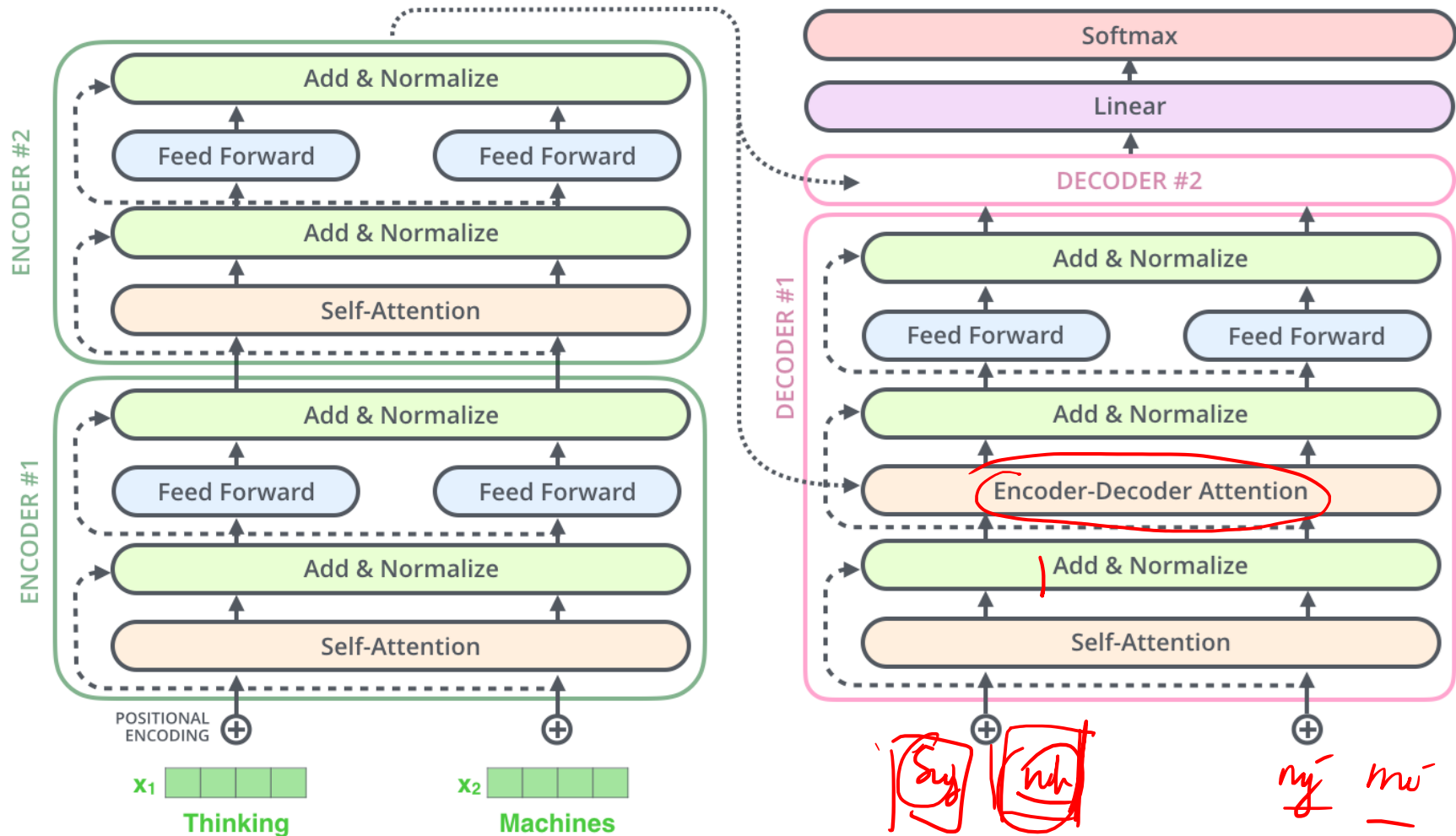
ENCODER #2

ENCODER #1



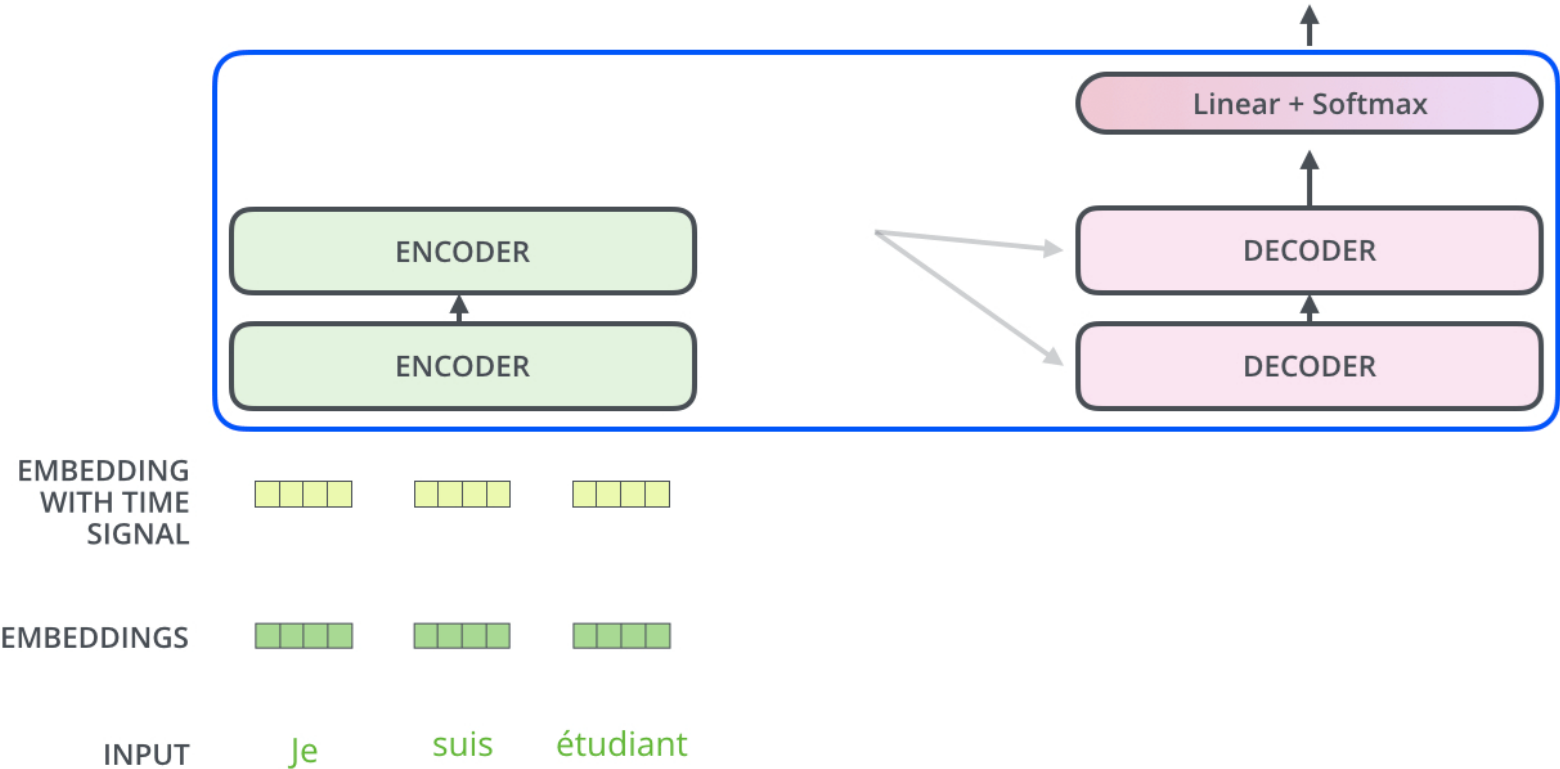






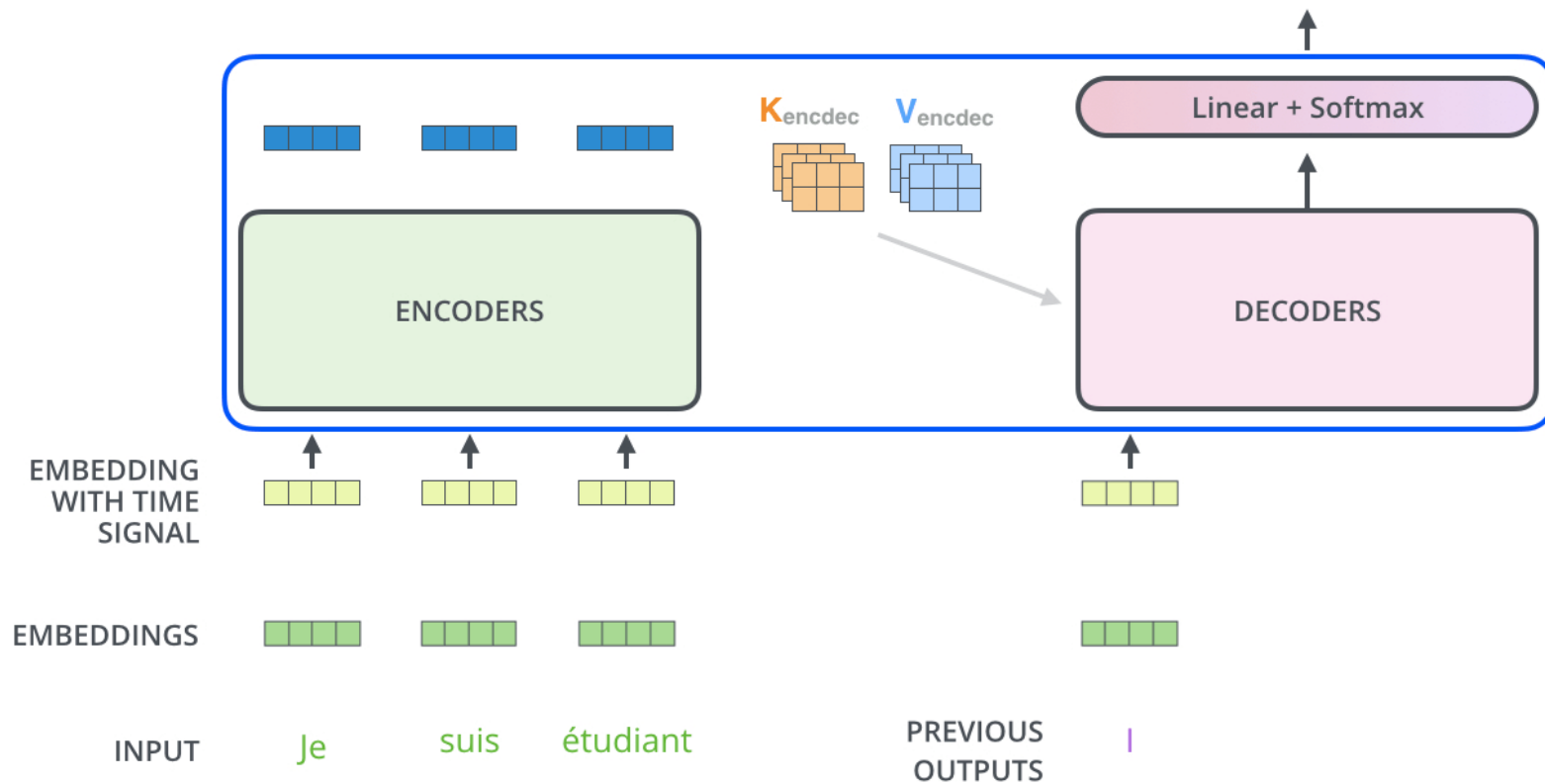
Decoding time step: 1 2 3 4 5 6

OUTPUT

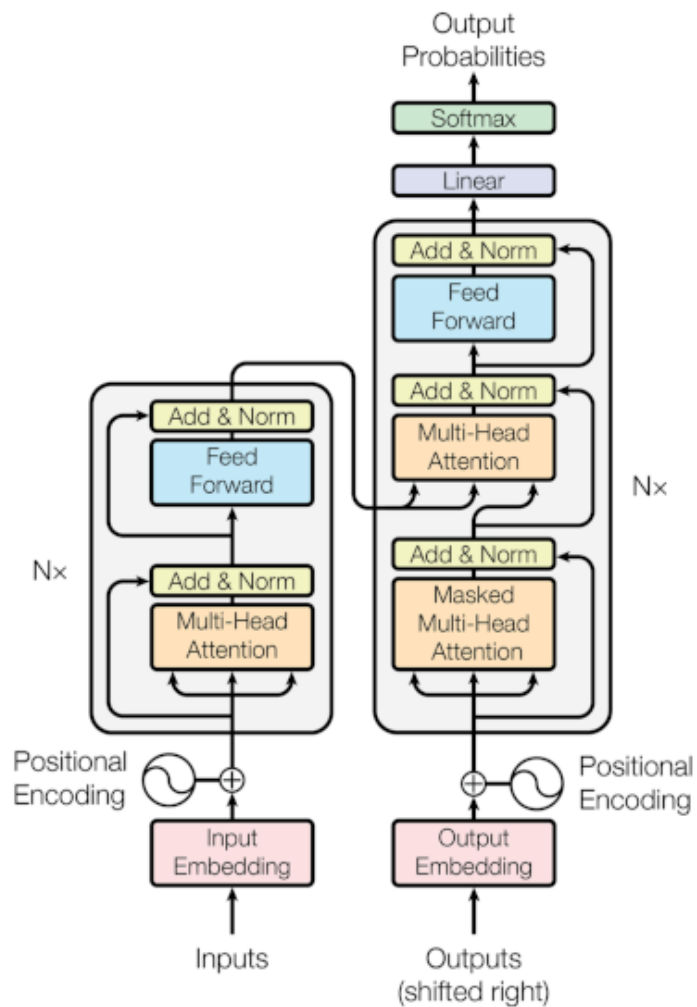


Decoding time step: 1 2 3 4 5 6

OUTPUT |



Transformer

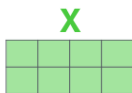


Multi-head attention

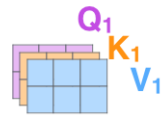
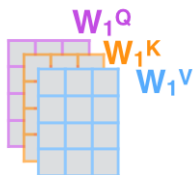
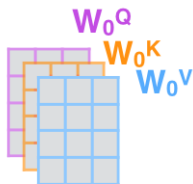
1) This is our input sentence*

Thinking
Machines

2) We embed each word*



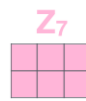
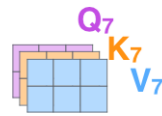
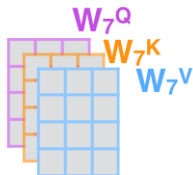
3) Split into 8 heads.
We multiply X or R with weight matrices



...

...

...



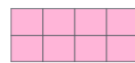
4) Calculate attention using the resulting $Q/K/V$ matrices

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

W^O



Z

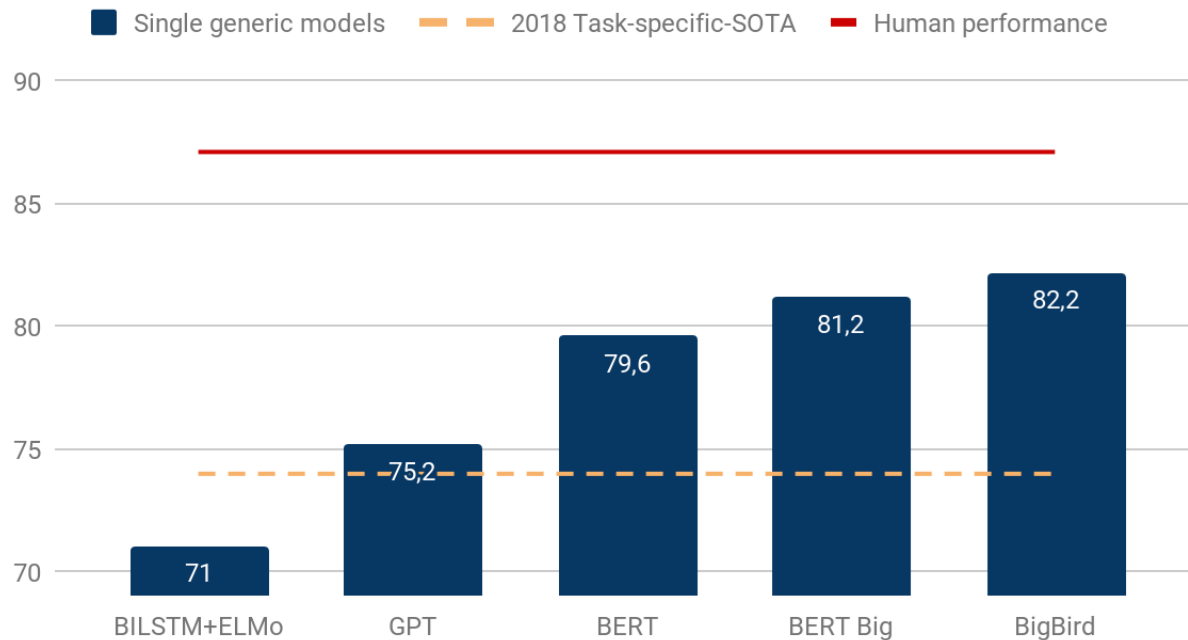


* In all encoders other than #0, we don't need embedding. We start directly with the output of the encoder right below this one

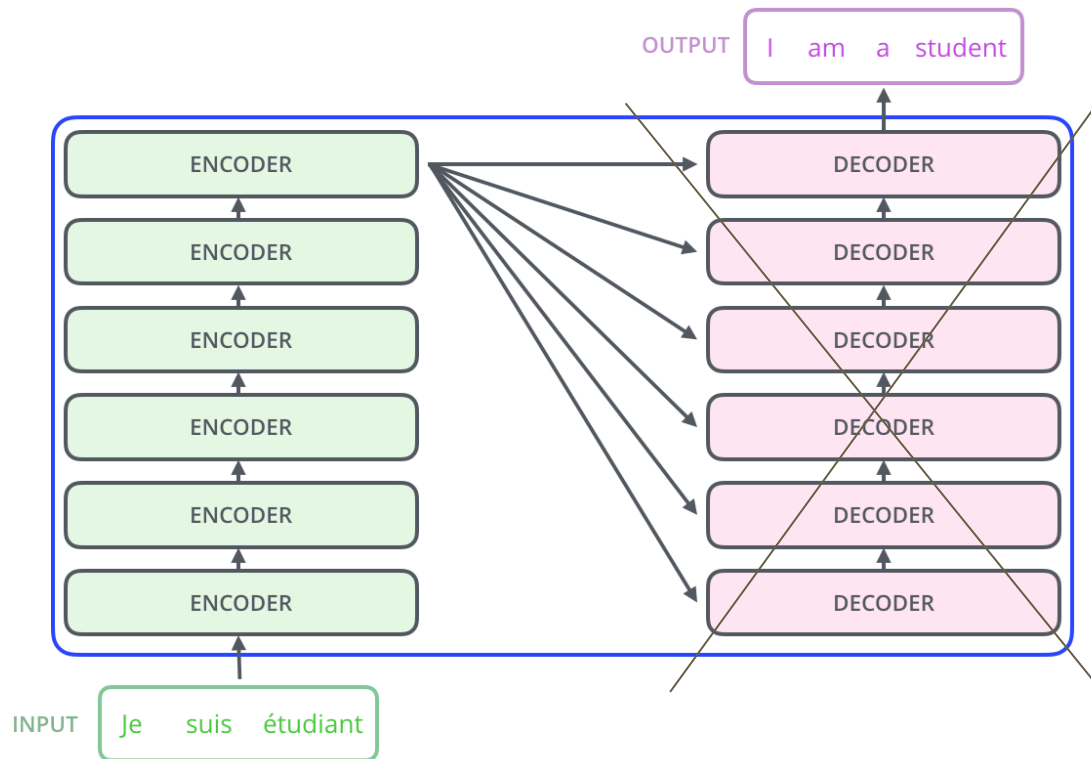


BERT

GLUE scores evolution over 2018-2019

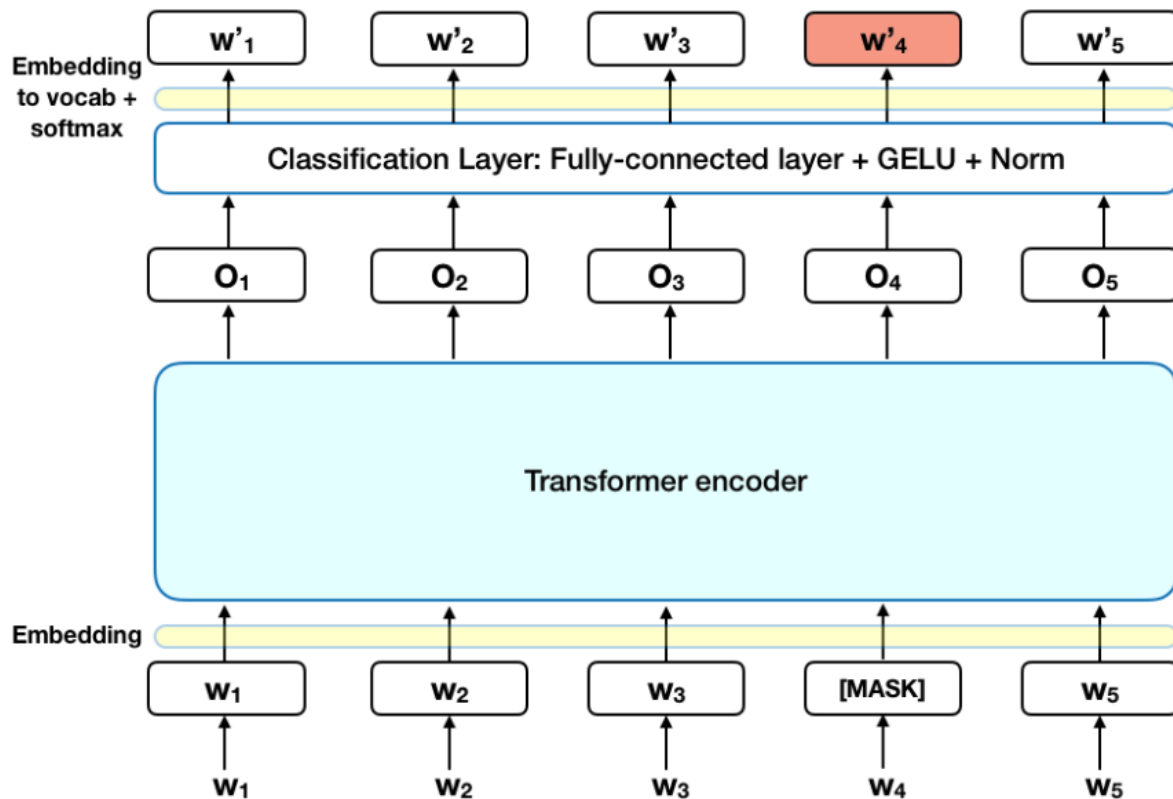


BERT model

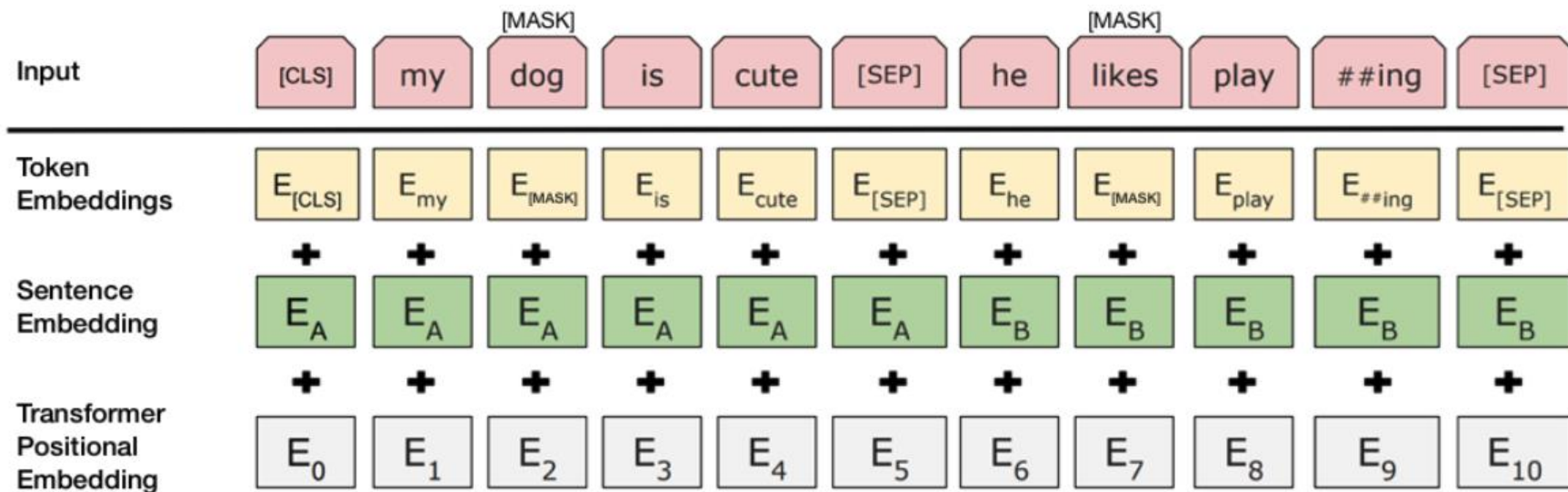


- Mask Language Modeling (MLM)
- Next Sentence Prediction (NSP)

Masked LM (MLM)



Next sentence prediction



Fine-tuning BERT

- Classification tasks such as sentiment analysis.
- In Question Answering tasks (e.g. SQuAD v1.1).
- In Named Entity Recognition (NER).

Q&A

