

Attention is All You Need [논문리뷰]

12181912 김현수



목차

Table of contents



Attention (Self-Attention, Multi-head Attention)



Encoder



Decoder



Result

Attention



Attention (Self-Attention, Multi-head Attention)



Encoder



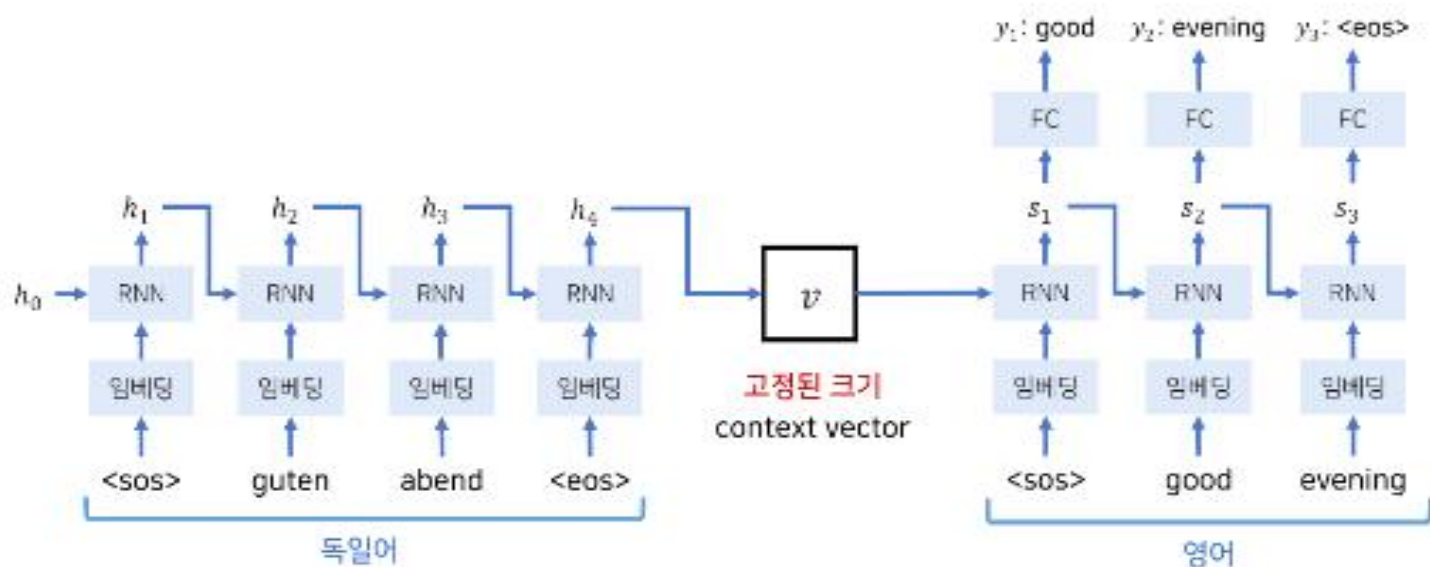
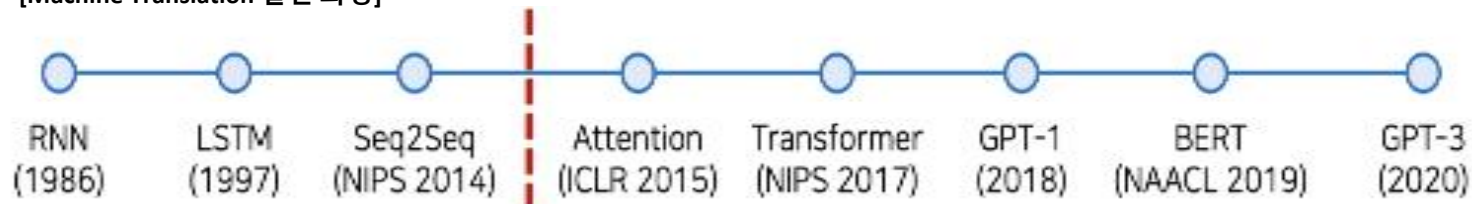
Decoder



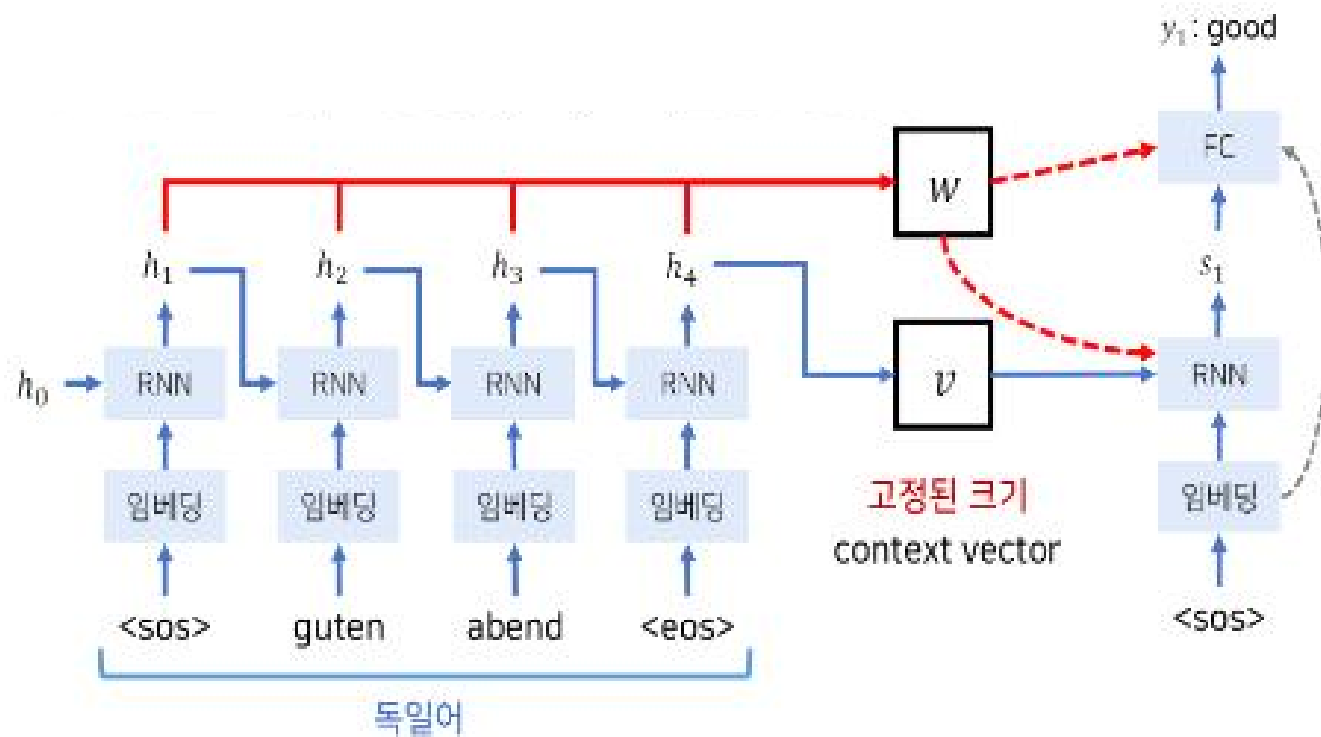
Result

Attention?

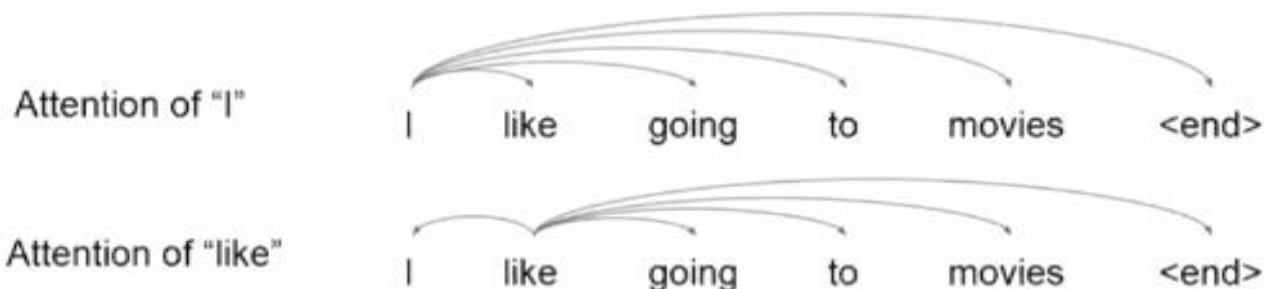
[Machine Translation 발전 과정]



Attention?



Self-Attention



- $Attention(Q, K, V) = Softmax\left(\frac{QK^T}{\sqrt{d}}\right)V$

0.5	0.1	0.0	0.2	0.2	0.0
0.2	0.6	0.0	0.0	0.1	0.0
...
...
...
...

$$Softmax\left(\frac{QK^T}{\sqrt{d}}\right)$$

I
like
going
to
movies
<end>

V



$0.5*I + 0.1*like + 0.2*to + 0.2* movies$
$0.2*I + 0.6*like + 0.1*movies$
...
...
...
...

$Attention(Q, K, V)$

Self-Attention

Input

Thinking

Machines

Embedding

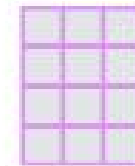
x_1 

x_2 

Queries

q_1 

q_2 

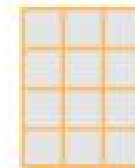


W^Q

Keys

k_1 

k_2 

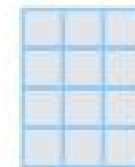


W^K

Values

v_1 

v_2 



W^V

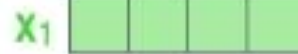
Self-Attention

Input

Thinking

Machines

Embedding



Queries



Keys



Values



Score

$$q_1 \cdot k_1 = 112$$

$$q_1 \cdot k_2 = 96$$

Self-Attention

Input

Embedding

Queries

Keys

Values

Score

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

Softmax

Softmax

X

Value

Sum

Thinking

x_1 

q_1 

k_1 

v_1 

$q_1 \cdot k_1 = 112$

14

0.88

v_1 

z_1 

Machines

x_2 

q_2 

k_2 

v_2 

$q_1 \cdot k_2 = 96$

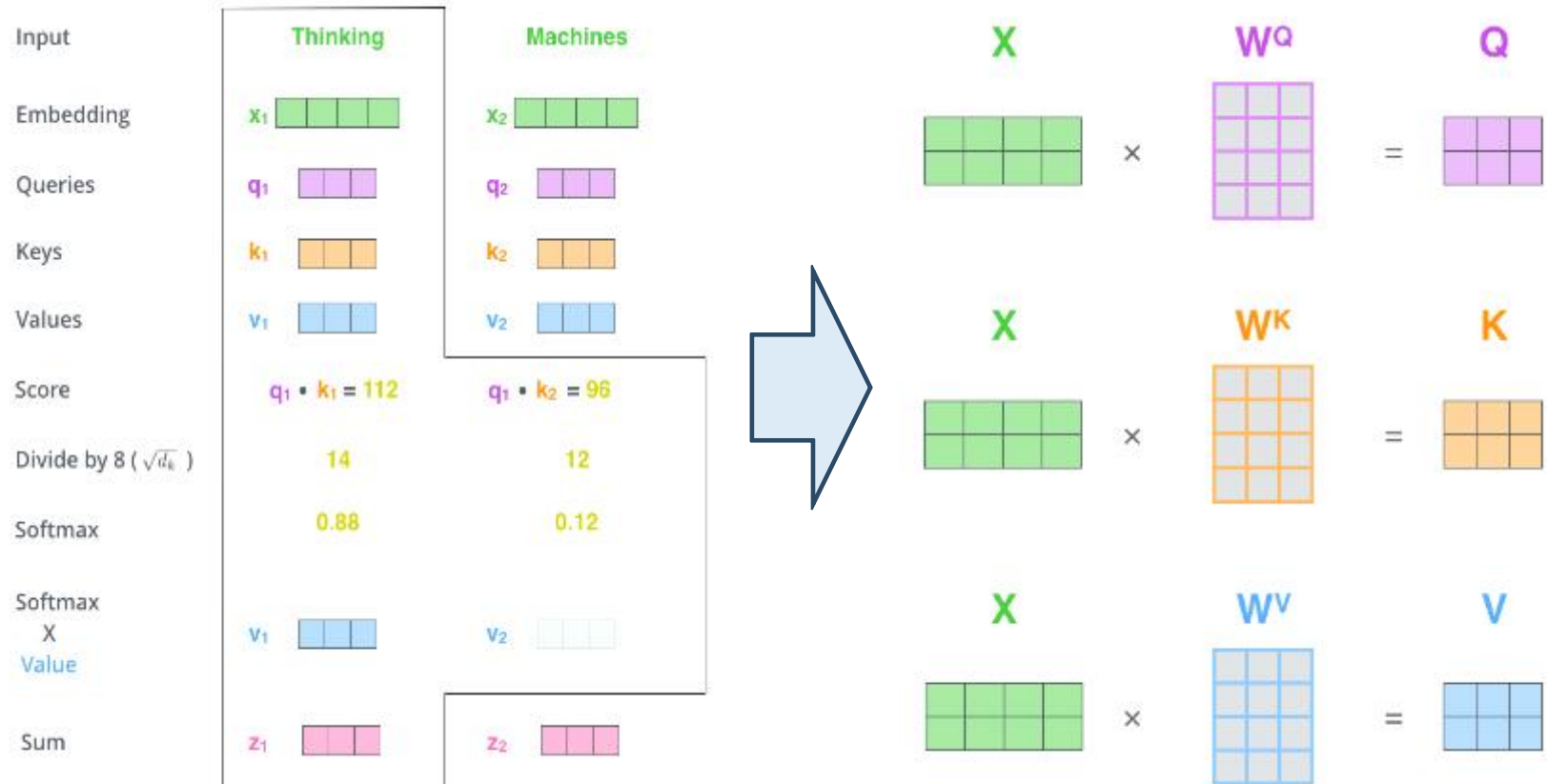
12

0.12

v_2 

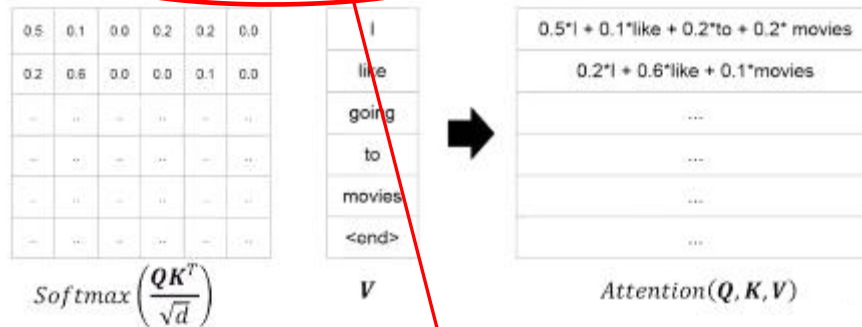
z_2 

Self-Attention



Self-Attention

- $Attention(Q, K, V) = \text{Softmax}\left(\frac{QK^T}{\sqrt{d}}\right)V$

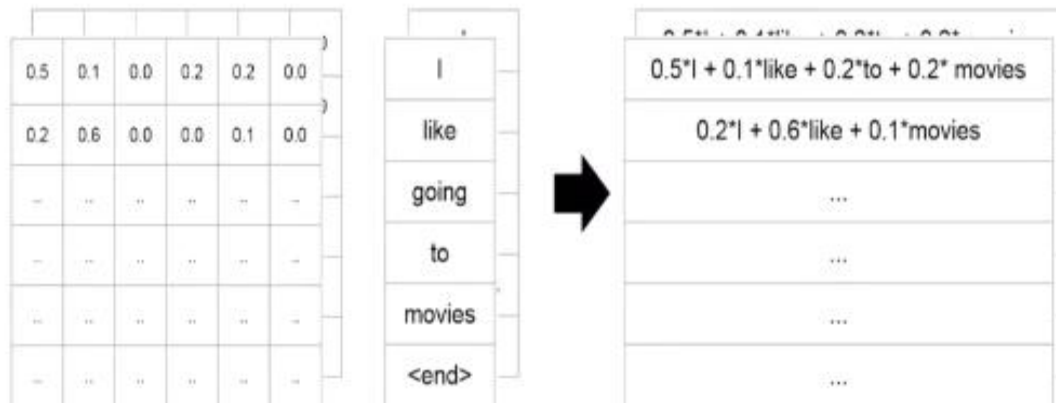
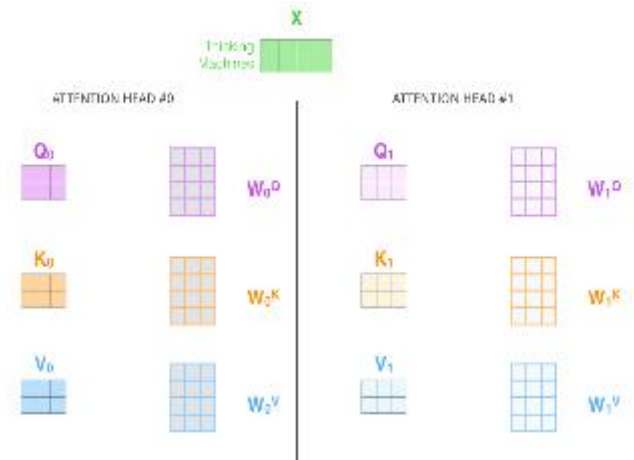
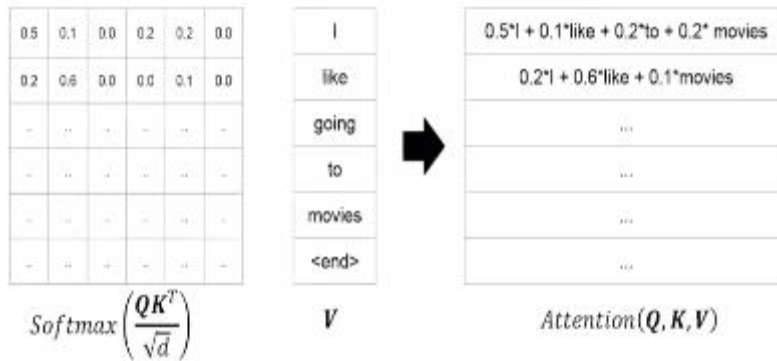


$$\text{softmax}\left(\frac{\begin{matrix} \text{Q} \\ \text{[3x3 grid]} \end{matrix} \times \begin{matrix} \text{K}^T \\ \text{[3x3 grid]} \end{matrix}}{\sqrt{d_k}}\right) \begin{matrix} \text{V} \\ \text{[3x3 grid]} \end{matrix}$$

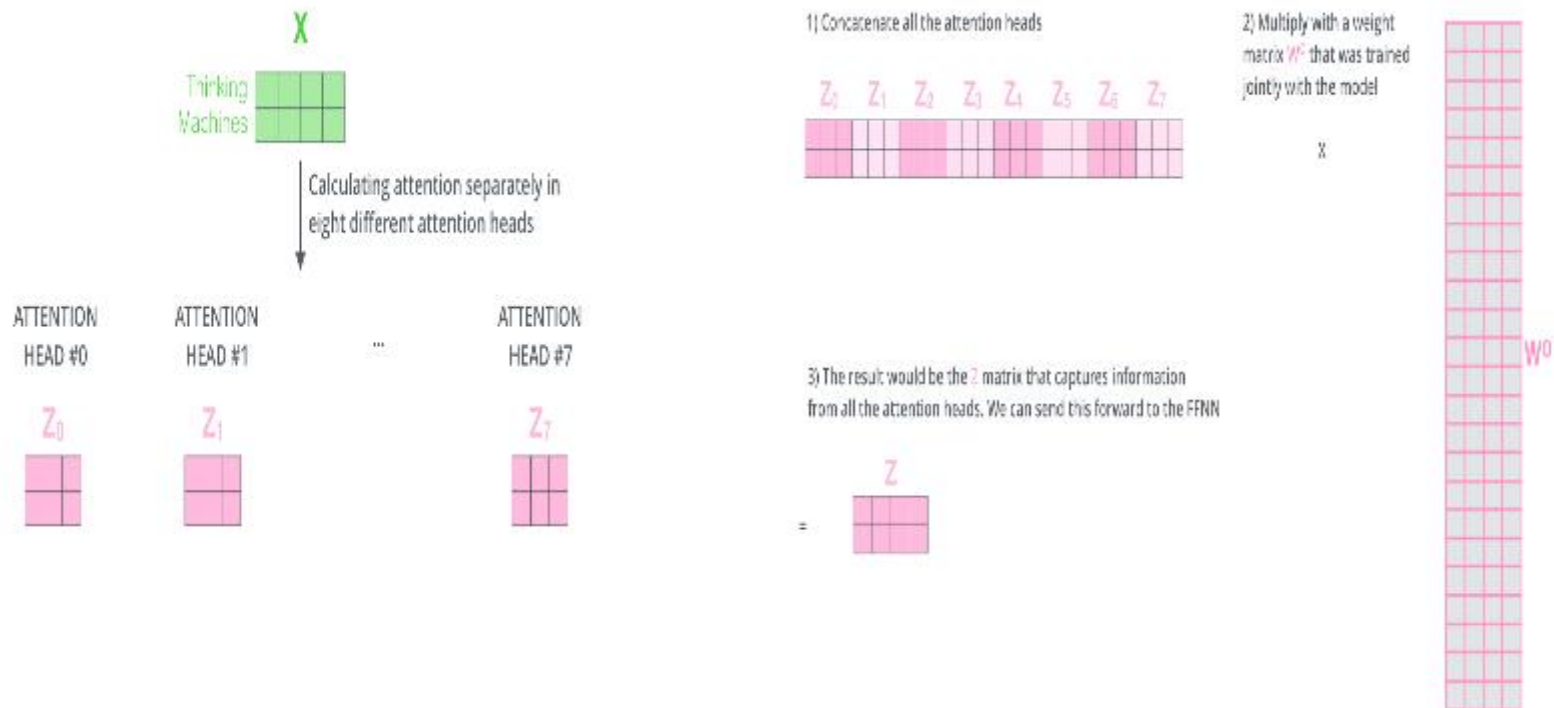
$$= \begin{matrix} \text{Z} \\ \text{[3x3 grid]} \end{matrix}$$

Multi-head Attention

- $Attention(Q, K, V) = Softmax\left(\frac{QK^T}{\sqrt{d}}\right)V$



Multi-head Attention



Multi-head Attention

1) This is our
input sentence*

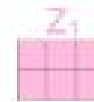
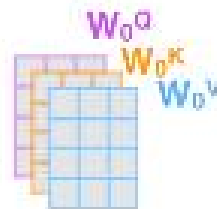
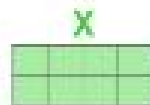
2) We embed
each word*

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

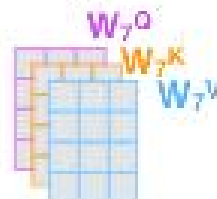
Thinking
Machines



...

...

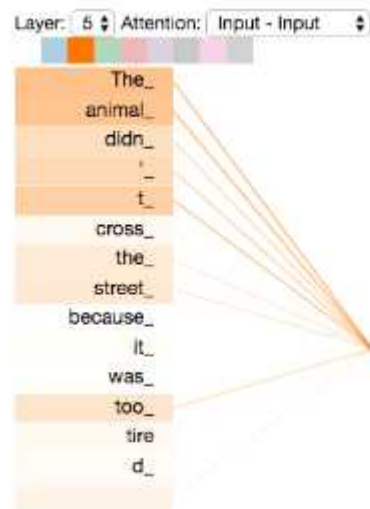
...



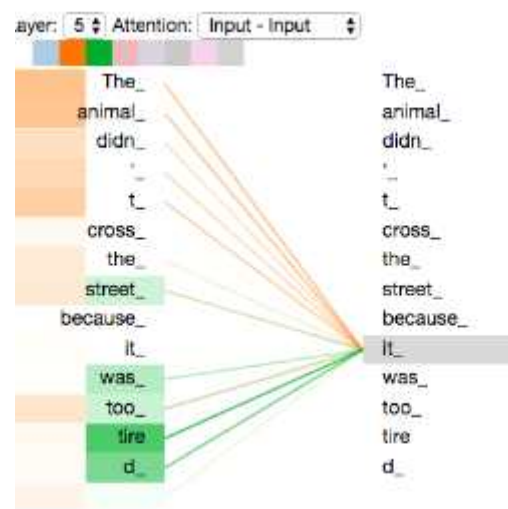
* In all encoders other than #0,
we don't need embedding.
We start directly with the output
of the encoder right below this one



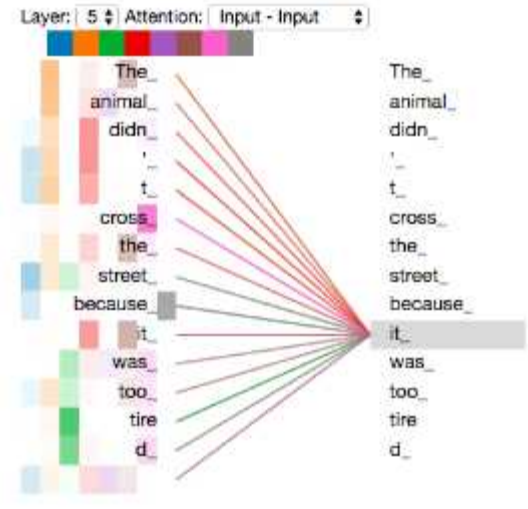
Visualize Attention



1 attention



2 attentions



8 attentions

Encoder



Attention (Self-Attention, Multi-head Attention)



Encoder



Decoder



Result

Encoder

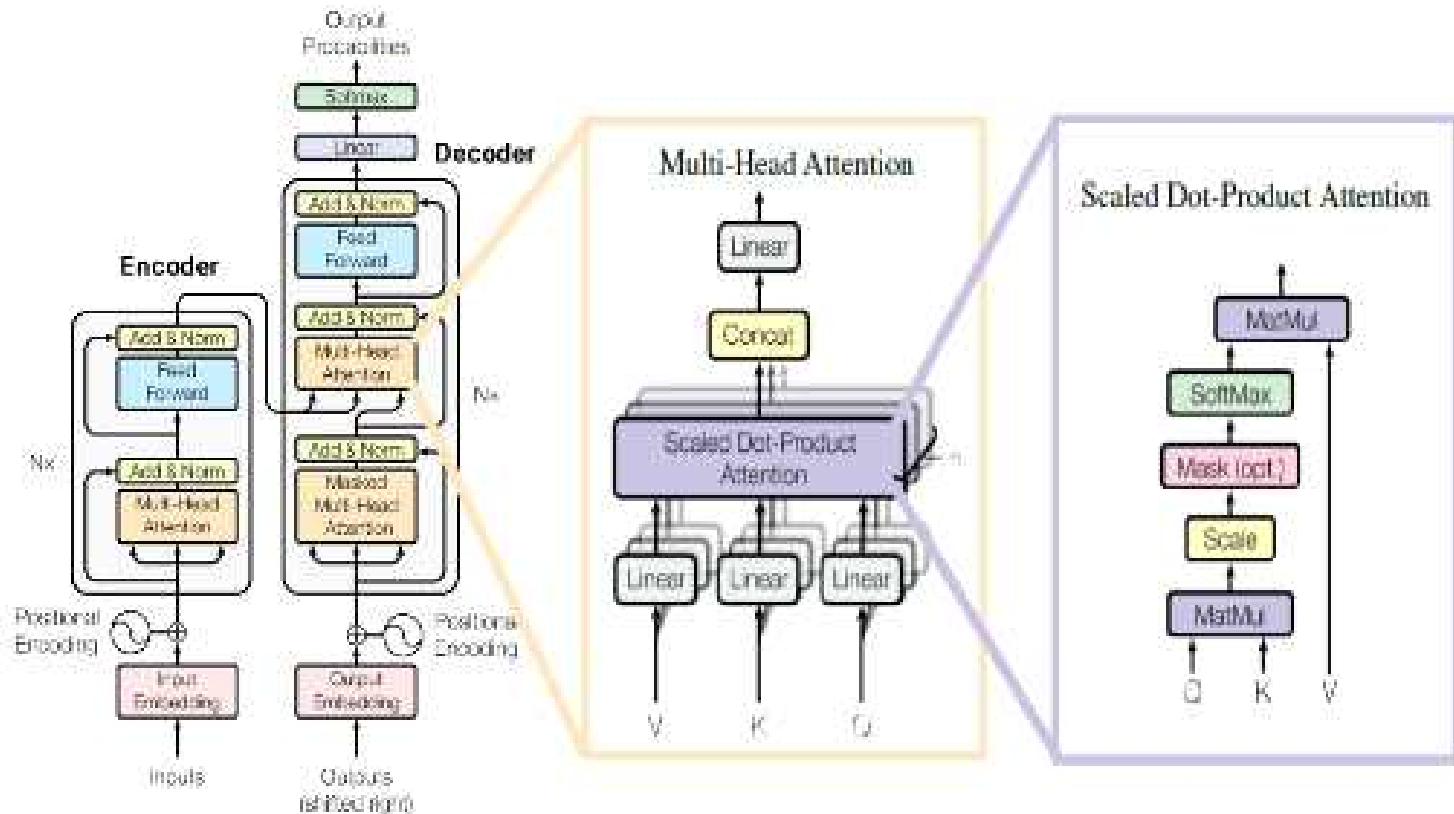
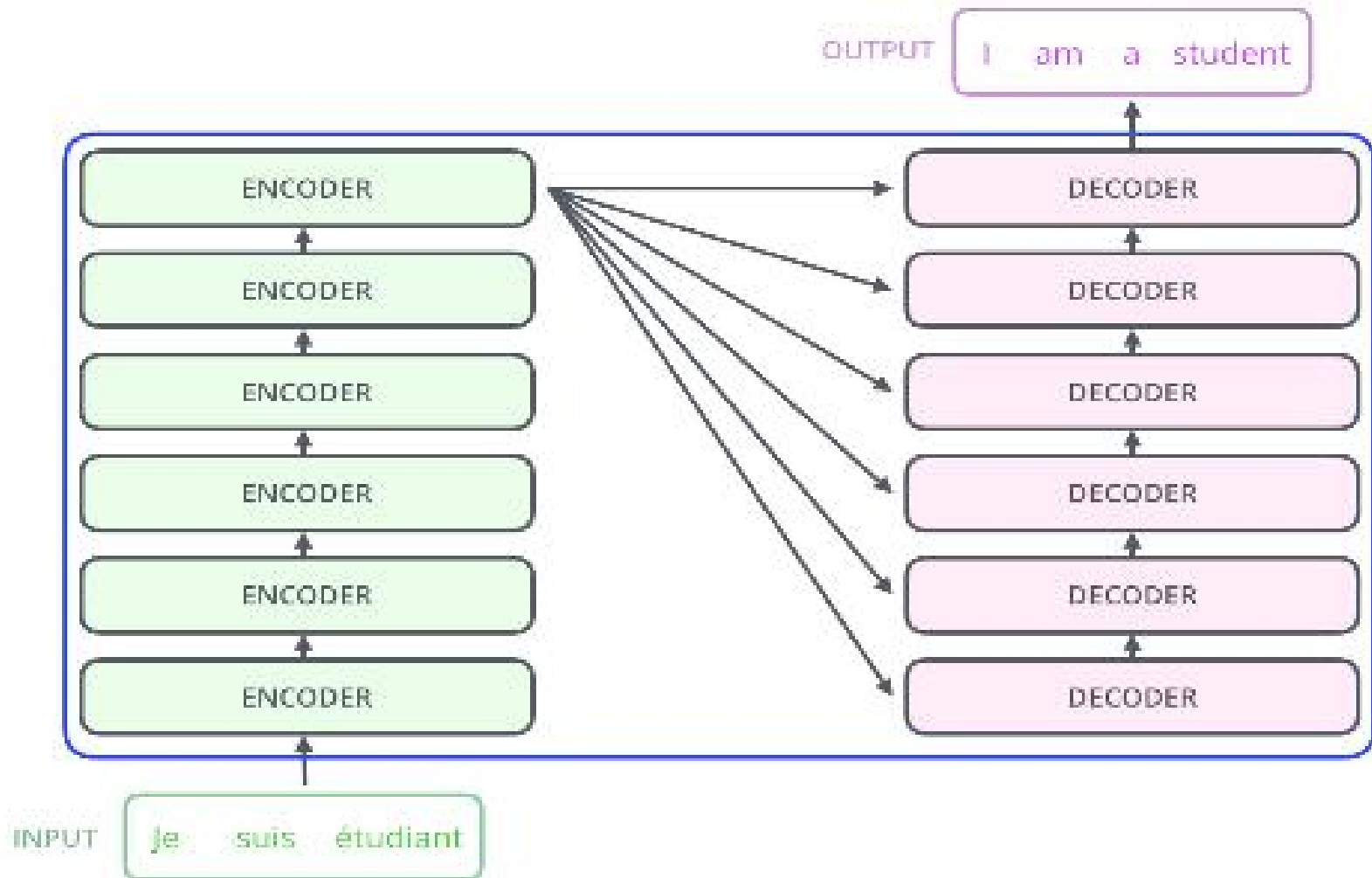
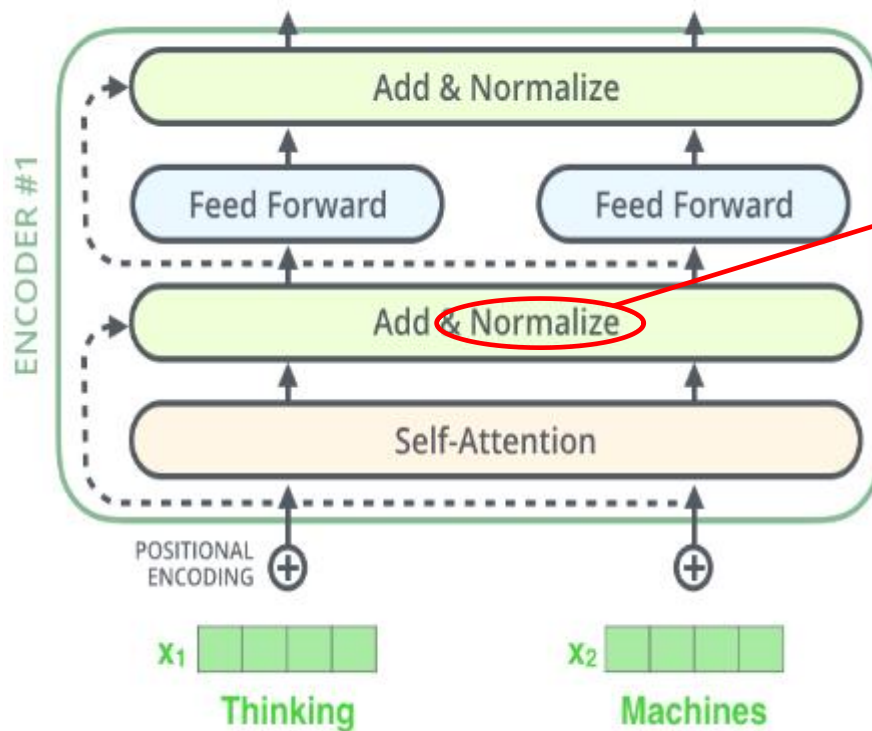


Figure 1: The Transformer - model architecture.

Encoder



Layer Normalization



Batch Normalization

batch

1	3	6
2	2	2
0	1	5
4	6	1
5	2	3
1	0	1

mean

3
2
3
4
3
1

std

3
0
3
3
2
1

Same for all training examples

Layer Normalization

batch

1	3	6
2	2	2
0	1	5
4	6	1
5	2	3
1	0	1

mean

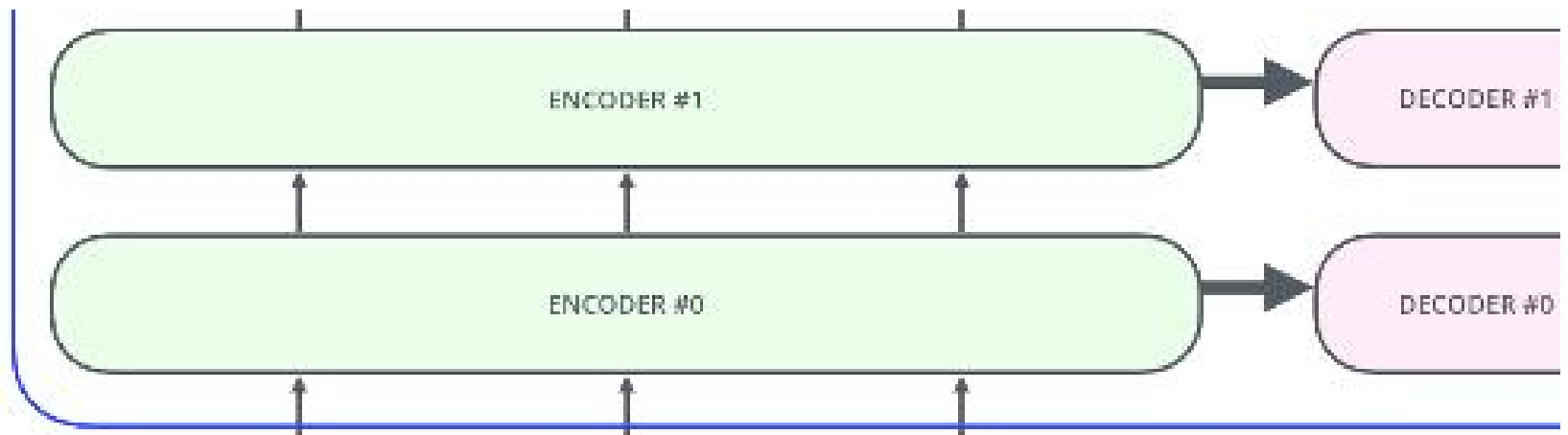
2
3
3

std

2
2
2

Same for all feature dimensions

Positional Encoding



EMBEDDING
WITH TIME
SIGNAL

x_1

x_2

x_3

=

=

=

POSITIONAL
ENCODING

t_1

t_2

t_3

+

+

+

EMBEDDINGS

x_1

x_2

x_3

INPUT

Je

suis

étudiant

Positional Encoding

<How to Positional Encoding>

1. 각 time-step(문장에서 단어의 위치)에 대해 고유한 인코딩을 출력해야 합니다.
2. 일정 time-step 사이의 거리는 길이가 다른 문장끼리 일정해야 한다.
3. test data에서 train data에서보다 더 긴 시퀀스가 들어왔을때도 처리할 수 있어야 한다. (더 긴 문장도 일반화가 가능해야함)
→ 상한(upper bound)이 필요하다.
4. 정확하게 위치를 결정할 수 있어야 한다.(확률값 x)

• Training sample

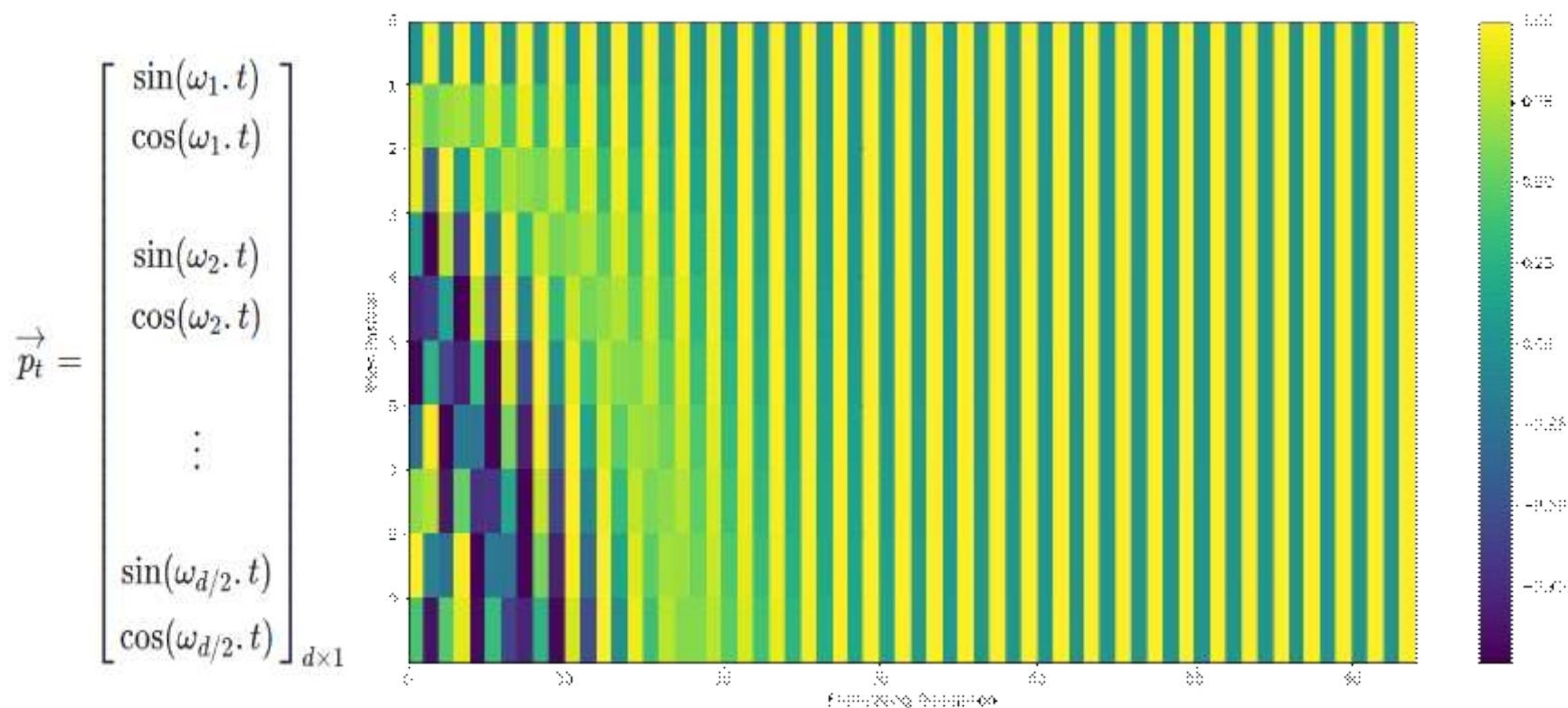
Sentence	Dark	horses	are	faster	than	white	horses
Pos 1	1	2	3	4	5	6	7
Pos 2	0.14	0.28	0.42	0.56	0.70	0.84	1.00

• Test sample

Sentence	Dark	horses	might	be	faster	than	white	horses
Pos 1	1	2	3	4	5	6	7	8
Pos 2	0.12	0.25	0.37	0.50	0.62	0.75	0.87	1.00

Positional Encoding

$$\vec{p}_t^{(i)} = f(t)^{(i)} := \begin{cases} \sin(\omega_k \cdot t), & \text{if } i = 2k \\ \cos(\omega_k \cdot t), & \text{if } i = 2k + 1 \end{cases} \quad \omega_k = \frac{1}{10000^{2k/d}}$$



Positional Encoding

$$\vec{p}_t = \begin{bmatrix} \sin(\omega_1 \cdot t) \\ \cos(\omega_1 \cdot t) \\ \\ \sin(\omega_2 \cdot t) \\ \cos(\omega_2 \cdot t) \\ \\ \vdots \\ \\ \sin(\omega_{d/2} \cdot t) \\ \cos(\omega_{d/2} \cdot t) \end{bmatrix}_{d \times 1}$$

```
import math
import matplotlib.pyplot as plt

n = 4 # 단어(word)의 개수
dim = 8 # 임베딩(embedding) 차원

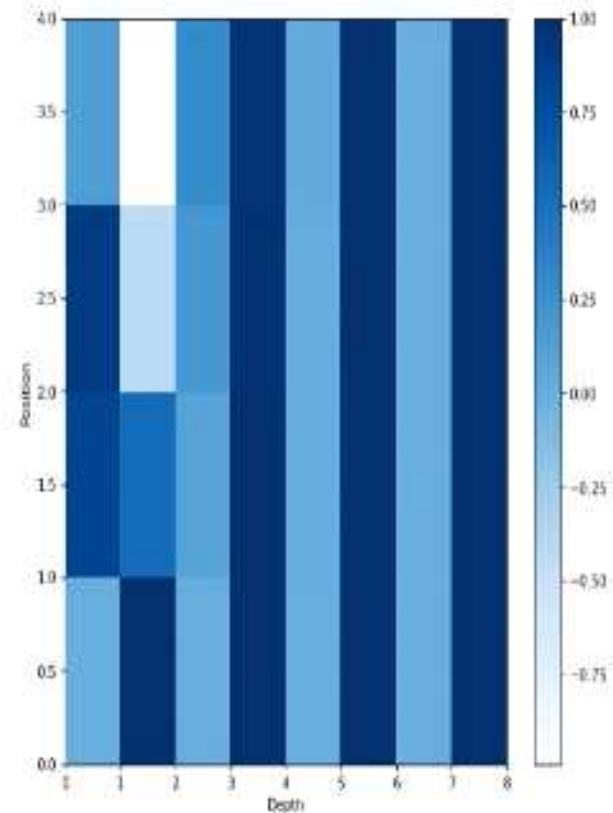
def get_angles(pos, i, dim):
    angles = 1 / math.pow(10000, (2 * (i // 2)) / dim)
    return pos * angles

def get_positional_encoding(pos, i, dim):
    if i % 2 == 0: # 짝수인 경우 사인 함수
        return math.sin(get_angles(pos, i, dim))
    # 홀수인 경우 코사인 함수
    return math.cos(get_angles(pos, i, dim))

result = [[0] * dim for _ in range(n)]

for i in range(n):
    for j in range(dim):
        result[i][j] = get_positional_encoding(i, j, dim)
```

출력 결과: plt.pcolormesh(result, cmap='Blues')



Decoder



Attention (Self-Attention, Multi-head Attention)



Encoder



Decoder



Result

Decoder

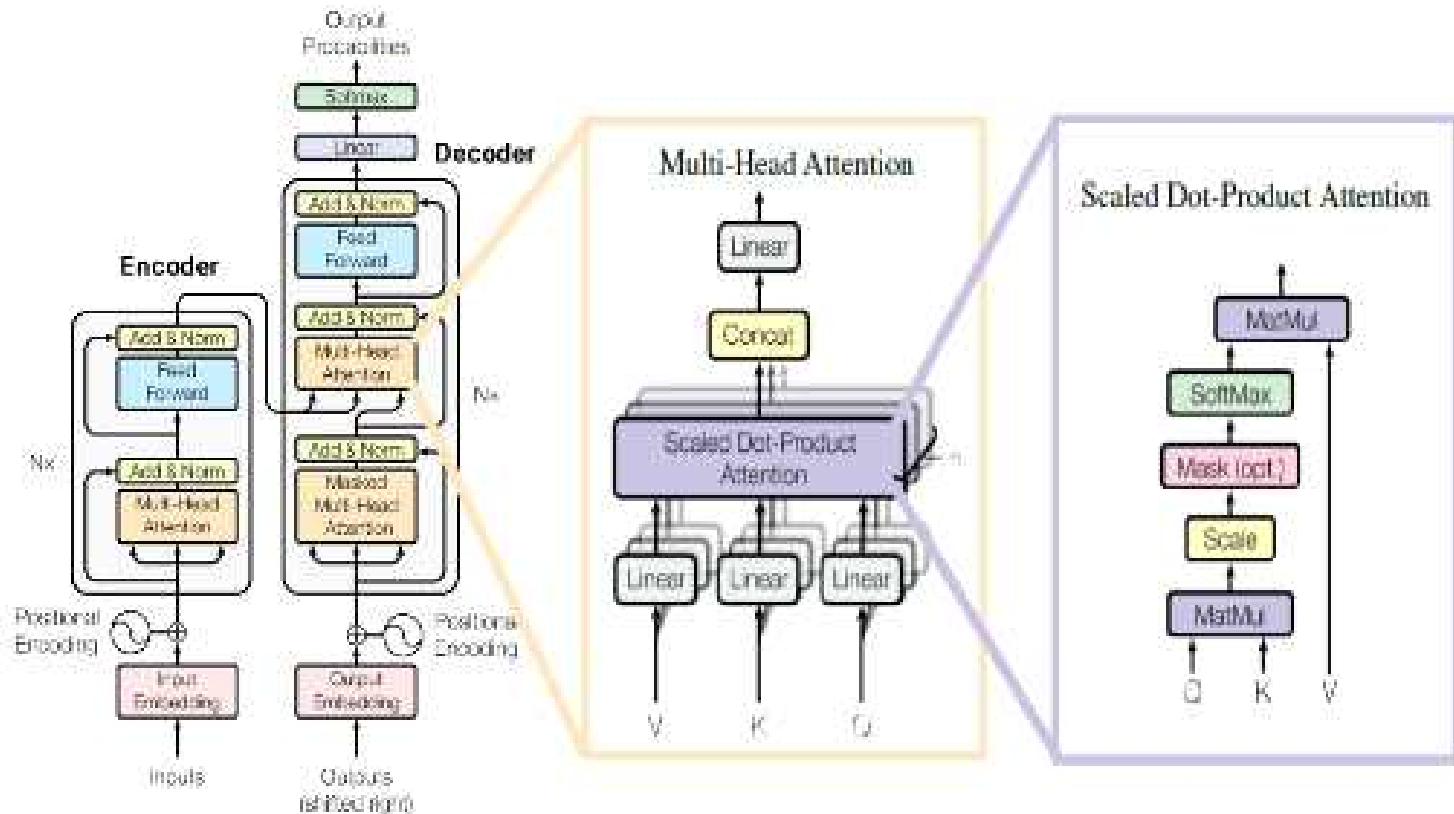
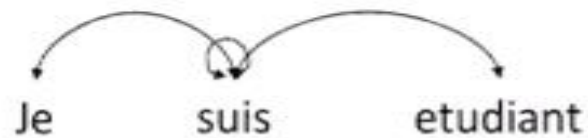


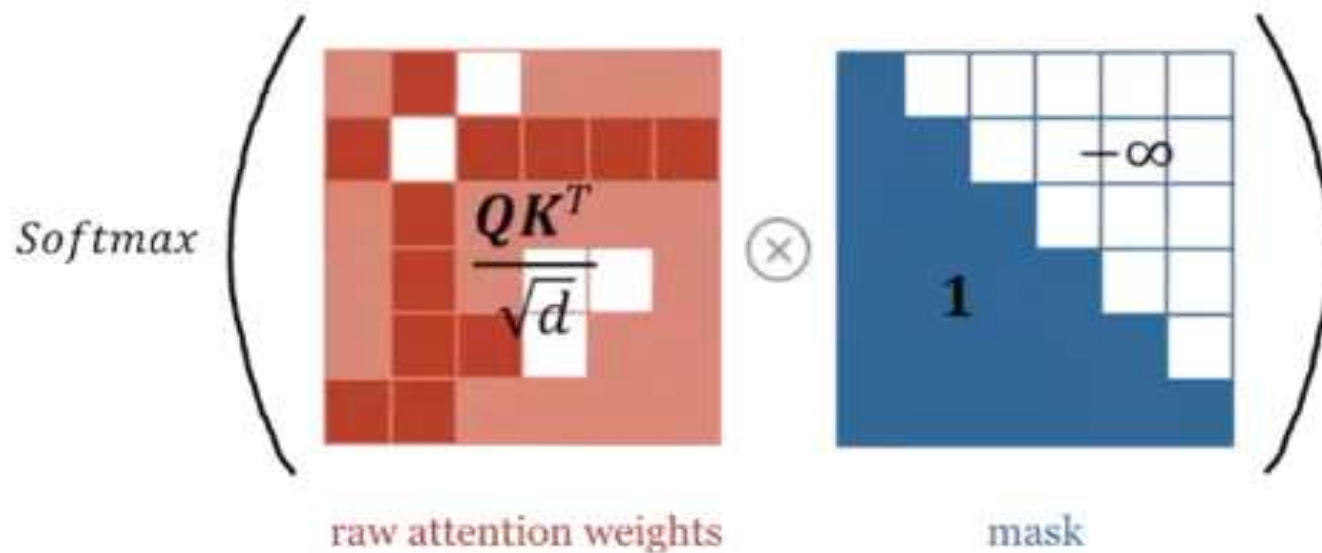
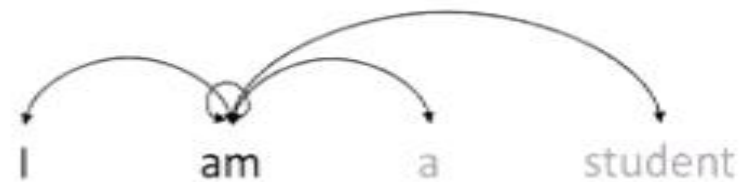
Figure 1: The Transformer - model architecture.

Masked Multi-Head Attention

Encoder attention from "suis"



Decoder attention from "am"



Masked Multi-Head Attention

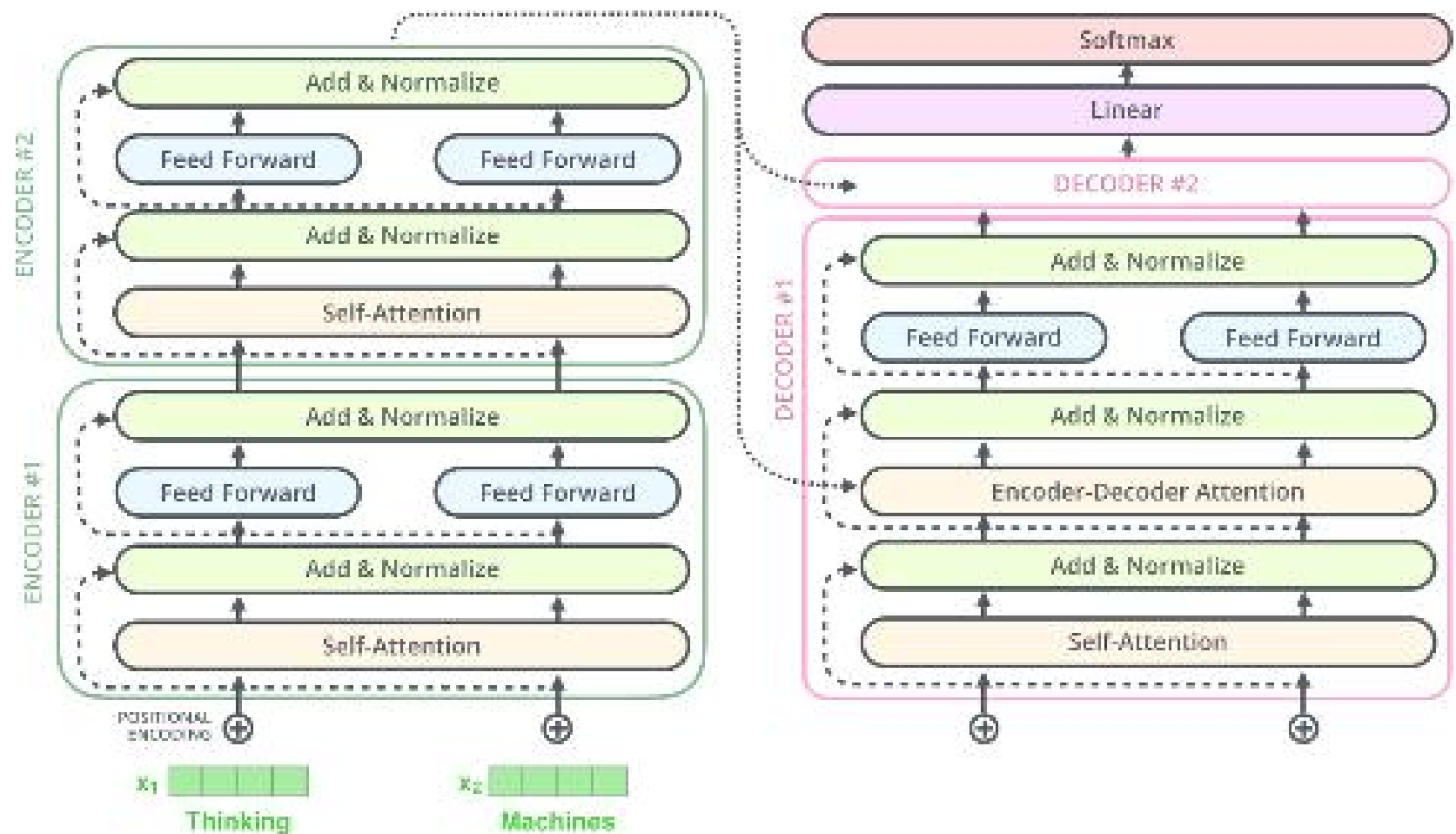
1.0	0	0	0	0	0
0.3	0.7	0	0	0	0
0.2	0.3	0.5	0	0	0
0.1	0.2	0.1	0.6	0	0
...	0
...

I
like
going
to
movies
.

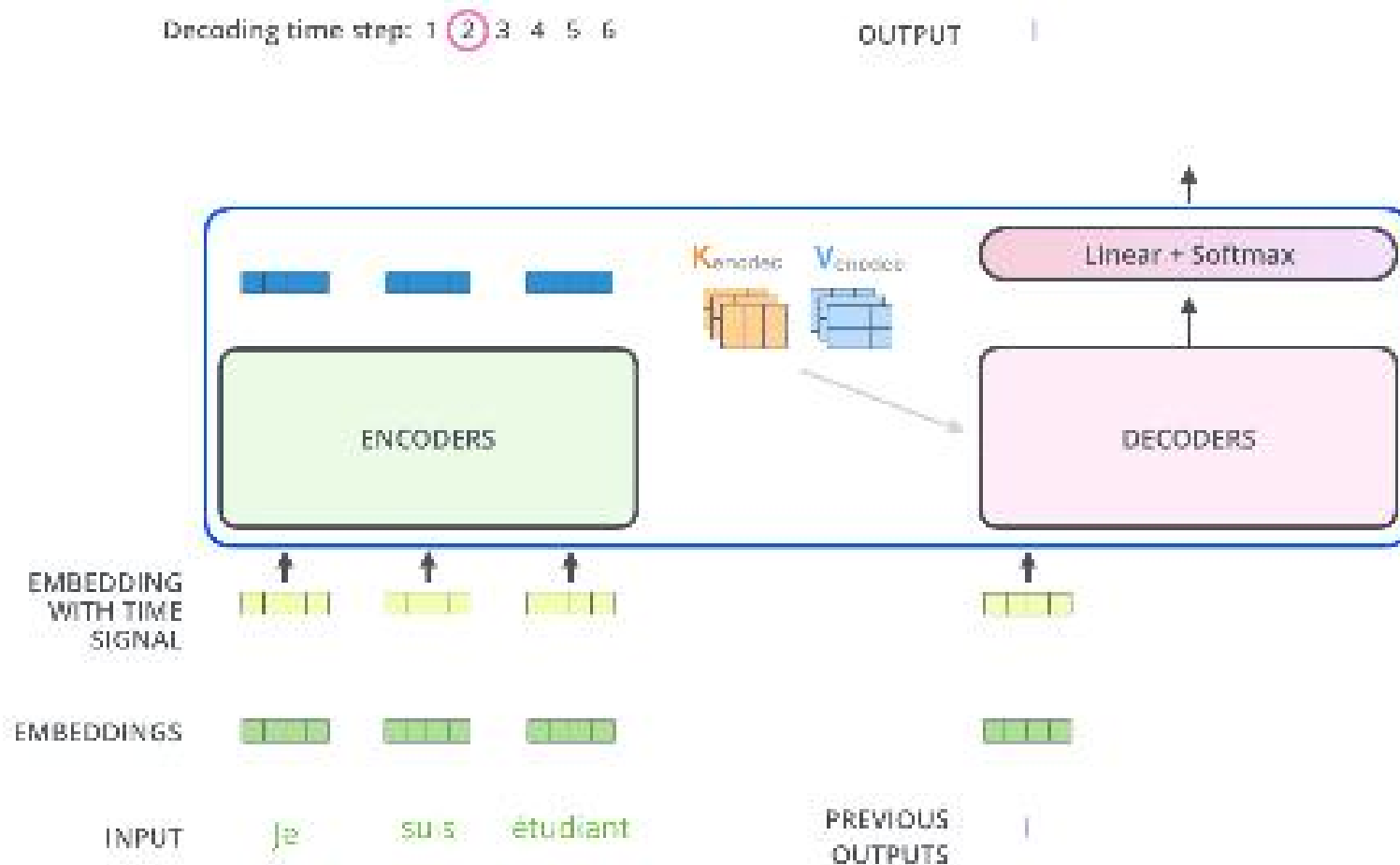


$1.0 * I$
$0.3 * I + 0.7 * \text{like}$
$0.2 * I + 0.3 * \text{like} + 0.5 * \text{going}$
$0.1 * I + 0.2 * \text{like} + 0.1 * \text{going} + 0.6 * \text{to}$
...
...

Decoder



Decoder

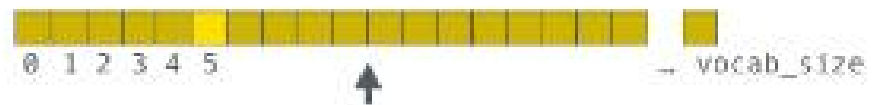


Decoder

Which word in our vocabulary
is associated with this index?

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

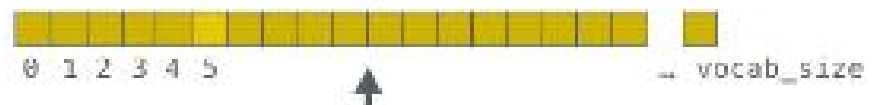
log_probs



am

5

logits



Softmax

Linear

Decoder stack output



Reference

- <https://arxiv.org/pdf/1706.03762.pdf>
- https://kazemnejad.com/blog/transformer_architecture_positional_encoding/#what-is-positional-encoding-and-why-do-we-need-it-in-the-first-place
- <https://www.youtube.com/watch?v=AA621UofTUA>
- <https://www.youtube.com/watch?v=h8avp8yDKV4&list=PLLENHvsRRLjDHllrXj0B8sz5-4xVbisBL&index=24>



Q & A

Thank you.