# Attention Is All You Need
# (How to implement Transformer)

*2020. 01. 06*

*Juhyuk Lee*

*sehkmg@gmail.com*

KAIST

MLILAB
Machine Learning & Intelligence

# Table of Contents

- **NMT Basics and Encoder-Decoder Model**

- **Abstract View of Transformer Architecture**

- **Dive into the Transformer Architecture**

# Table of Contents

- **NMT Basics and Encoder-Decoder Model**


- Abstract View of Transformer Architecture


- Dive into the Transformer Architecture

# NMT Basics and Encoder-Decoder Model

- Neural Machine Translation
  - Dataset만 이용해서 black box model을 end-to-end로 학습시켜 번역.
  - 기본적인 Deep Learning Framework.

# NMT Basics – Dataset 준비

- Dataset 준비

**Source (Train)**

Je suis etudiant.
Quel mois?

**Target (Train)**

I am a student.
What month?

**Source (Valid)**

Je suis medecin.
Ce mois-ci?

**Target (Valid)**

I am a doctor.
This month?

**Source (Test)**

Je suis enseignant.
Quel jour?

**Target (Test)**

I am a teacher.
What day?

# NMT Basics – Data preprocessing

- Make batches - tokenize

**Source Batch (Train)**

[Je, suis, etudiant, .]
[Quel, mois, ?]

**Target Batch (Train)**

[I, am, a, student, .]
[What, month, ?]

**Source Batch (Valid)**

[Je, suis, medecin, .]
[Ce, mois, -, ci, ?]

**Target Batch (Valid)**

[I, am, a, doctor, .]
[This, month, ?]

**Source Batch (Test)**

[Je, suis, enseignant, .]
[Quel, jour, ?]

# NMT Basics – Data preprocessing

- Make batches – add <sos>, <eos> token

**Source Batch (Train)**　　**Target Batch (Train)**

[Je, suis, etudiant, .]　　[<sos>, I, am, a, student, ., <eos>]
[Quel, mois, ?]　　[<sos>, What, month, ?, <eos>]


**Source Batch (Valid)**　　**Target Batch (Valid)**

[Je, suis, medecin, .]　　[<sos>, I, am, a, doctor, ., <eos>]
[Ce, mois, -, ci, ?]　　[<sos>, This, month, ?, <eos>]


**Source Batch (Test)**

[Je, suis, enseignant, .]
[Quel, jour, ?]

# NMT Basics – Data preprocessing

- Make batches – padding

**Source Batch (Train)**

[Je, suis, etudiant, .]
[Quel, mois, ?, <pad>]

**Target Batch (Train)**

[<sos>, I, am, a, student, ., <eos>]
[<sos>, What, month, ?, <eos>, <pad>, <pad>]

**Source Batch (Valid)**

[Je, suis, medecin, ., <pad>]
[Ce, mois, -, ci, ?]

**Target Batch (Valid)**

[<sos>, I, am, a, doctor, ., <eos>]
[<sos>, This, month, ?, <eos>, <pad>, <pad>]

**Source Batch (Test)**

[Je, suis, enseignant, .]
[Quel, jour, ?, <pad>]

# NMT Basics – Data preprocessing

- Numericalize – make vocabulary from train dataset

**Source (Train)**

Je suis etudiant.
Quel mois?

**Target (Train)**

I am a student.
What month?

**Source Vocab**

0: <sos>
1: <eos>
2: <pad>
3: <unk>
4: Je
5: suis
6: etudiant
7: .
8: Quel
9: mois
10: ?
11: …..

**Target Vocab**

0: <sos>
1: <eos>
2: <pad>
3: <unk>
4: I
5: am
6: a
7: student
8: .
9: What
10: month
11: ?
12: …..

# NMT Basics – Data preprocessing

- Numericalize – handle out-of-vocabulary words

**Source Batch (Train)**

[Je, suis, etudiant, .]
[Quel, mois, ?, <pad>]

**Target Batch (Train)**

[<sos>, I, am, a, student, ., <eos>]
[<sos>, What, month, ?, <eos>, <pad>, <pad>]

**Source Batch (Valid)**

[Je, suis, <unk>, ., <pad>]
[<unk>, mois, <unk>, <unk>, ?]

**Target Batch (Valid)**

[<sos>, I, am, a, <unk>, ., <eos>]
[<sos>, <unk>, month, ?, <eos>, <pad>, <pad>]

**Source Batch (Test)**

[Je, suis, <unk>, .]
[Quel, <unk>, ?, <pad>]

MLILAB
Machine Learning & Intelligence

# NMT Basics – Data preprocessing

- Numericalize – numericalize

**Source Batch (Train)**

[4, 5, 6, 7]
[8, 9, 10, 2]

**Target Batch (Train)**

[0, 4, 5, 6, 7, 8, 1]
[0, 9, 10, 11, 1, 2, 2]

**Source Batch (Valid)**

[4, 5, 3, 7, 2]
[3, 9, 3, 3, 10]

**Target Batch (Valid)**

[0, 4, 5, 6, 3, 8, 1]
[0, 3, 10, 11, 1, 2, 2]
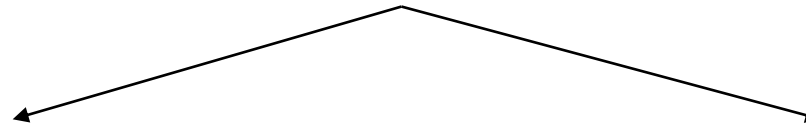
**Source Batch (Test)**

[4, 5, 3, 7]
[8, 3, 10, 2]

# Encoder-Decoder Model - Training

**Source Batch (Train)**

| 4 | 5 | 6 | 7 |
|---|---|---|---|
| 8 | 9 | 10 | 2 |

**Target Batch (Train)**

| 0 | 4 | 5 | 6 | 7 | 8 | 1 |
|---|---|---|---|---|---|---|
| 0 | 9 | 10 | 11 | 1 | 2 | 2 |

**Target Batch (Input)**

| 0 | 4 | 5 | 6 | 7 | 8 | 1 |
|---|---|---|---|---|---|---|
| 0 | 9 | 10 | 11 | 1 | 2 | 2 |

**Target Batch (True Output)**

| 0 | 4 | 5 | 6 | 7 | 8 | 1 |
|---|---|---|---|---|---|---|
| 0 | 9 | 10 | 11 | 1 | 2 | 2 |

Encoder

Decoder

# Encoder-Decoder Model - Training



| 0 | 4 | 5 | 6 | 7 | 8 | 1 |
| 0 | 9 | 10 | 11 | 1 | ✗ | ✗ |

0: <sos>
1: <eos>
2: <pad>
3: <unk>

Encoder

| 4 | 5 | 6 | 7 |
| 8 | 9 | 10 | 2 |

Decoder

| 0 | 4 | 5 | 6 | 7 | 8 | 1 |
| 0 | 9 | 10 | 11 | 1 | 2 | 2 |

# Encoder-Decoder Model - Validation

# Encoder-Decoder Model – Test (Inference)

<sos>와 <pad>는 제외하고 생성

Decoder

Encoder

| 4 |
| 9 |

0: <sos>
1: <eos>
2: <pad>
3: <unk>

| 4 | 5 | 3 | 7 |
|---|---|---|---|
| 8 | 3 | 10 | 2 |

| 0 |
| 0 |

MLILAB
Machine Learning & Intelligence

# Encoder-Decoder Model – Test (Inference)



<sos>와 <pad>는 제외하고 생성

Decoder

Encoder

| 4 | 5 |
|---|---|
| 9 | 3 |

| 4 | 5 | 3 | 7 |
|---|---|----|---|
| 8 | 3 | 10 | 2 |

| 0 | 4 |
|---|---|
| 0 | 9 |

0: <sos>
1: <eos>
2: <pad>
3: <unk>

# Encoder-Decoder Model – Test (Inference)

<sos>와 <pad>는 제외하고 생성

| 4 | 5 | 6 |
|---|---|---|
| 9 | 3 | 11 |

Decoder

Encoder

0: <sos>
1: <eos>
2: <pad>
3: <unk>

| 4 | 5 | 3 | 7 |
|---|---|---|---|
| 8 | 3 | 10 | 2 |

| 0 | 4 | 5 |
|---|---|---|
| 0 | 9 | 3 |

# Encoder-Decoder Model – Test (Inference)

<sos>와 <pad>는 제외하고 생성

| Encoder | | | |
|---|---|---|---|
| 4 | 5 | 3 | 7 |
| 8 | 3 | 10 | 2 |

| Decoder | | | |
|---|---|---|---|
| 0 | 4 | 5 | 6 |
| 0 | 9 | 3 | 11 |

| | | | |
|---|---|---|---|
| 4 | 5 | 6 | 3 |
| 9 | 3 | 11 | 1 |

0: <sos>
1: <eos>
2: <pad>
3: <unk>

# Encoder-Decoder Model – Test (Inference)

<sos>와 <pad>는 제외하고 생성

| Encoder |
|---|

| 4 | 5 | 3 | 7 |
|---|---|---|---|
| 8 | 3 | 10 | 2 |

| Decoder |
|---|

| 0 | 4 | 5 | 6 | 3 |
|---|---|---|---|---|
| 0 | 9 | 3 | 11 | 1 |

| 4 | 5 | 6 | 3 | 8 |
|---|---|---|---|---|
| 9 | 3 | 11 | 1 | 10 |

0: <sos>
1: <eos>
2: <pad>
3: <unk>

MLILAB
Machine Learning & Intelligence

# Encoder-Decoder Model – Test (Inference)



<sos>와 <pad>는 제외하고 생성

| 4 | 5 | 6 | 3 | 8 | 1 |
|---|---|---|---|---|---|
| 9 | 3 | 11 | 1 | 10 | 11 |

Encoder

Decoder

0: <sos>
1: <eos>
2: <pad>
3: <unk>

| 4 | 5 | 3 | 7 |
|---|---|---|---|
| 8 | 3 | 10 | 2 |

| 0 | 4 | 5 | 6 | 3 | 8 |
|---|---|---|---|---|---|
| 0 | 9 | 3 | 11 | 1 | 10 |

# Encoder-Decoder Model – Test (Inference)

Encoder

Decoder

종료 조건
1. \<eos>가 모든 문장에 나타났을 때.
2. 정해둔 max length에 도달했을 때.

0: \<sos>
1: \<eos>
2: \<pad>
3: \<unk>

| 4 | 5 | 3 | 7 |
|---|---|---|---|
| 8 | 3 | 10 | 2 |

| 0 | 4 | 5 | 6 | 3 | 8 | 1 |
|---|---|---|---|---|---|---|
| 0 | 9 | 3 | 11 | 1 | 10 | 11 |

**Target Vocab**

0: <sos>
1: <eos>
2: <pad>
3: <unk>
4: I
5: am
6: a
7: student
8: .
9: What
10: month
11: ?
12: …..

| 0 | 4 | 5 | 6 | 3 | 8 | 1 |
|---|---|---|---|---|---|---|
| 0 | 9 | 3 | 11 | 1 | 10 | 11 |

↓

| 0 | 4 | 5 | 6 | 3 | 8 | 1 |
|---|---|---|---|---|---|---|
| 0 | 9 | 3 | 11 | 1 | 10 | 11 |

↓

[I, am, a, <unk>, .]
[What, <unk>, ?]

↓

I am a <unk>.
What <unk>?

←  BLEU score  →

**Target (Test)**

I am a teacher.
What day?

# NMT Basics and Encoder-Decoder Model

- **NMT Basics**
  - Dataset 준비
  - Data preprocessing
    - Make batches
      - tokenize
      - add <sos>, <eos> token
      - padding
    - Numericalize
      - make vocabulary from train dataset
      - handle out-of-vocabulary words
      - numericalize
  - Compute BLEU Score

- **Encoder-Decoder Model**
  - Training
  - Validation
  - Test (Inference)

# NMT Basics and Encoder-Decoder Model

- **NMT Basics (Done!)**
  - Dataset 준비: Multi30k English to German Translation Dataset
  - Data preprocessing
    - Make batches
      - tokenize
      - add <sos>, <eos> token
      - padding
    - Numericalize
      - make vocabulary from train dataset
      - handle out-of-vocabulary words
      - numericalize
  - Compute BLEU Score

- **Encoder-Decoder Model**
  - Training
  - Validation
  - Test (Inference)

https://github.com/sehkmg/NMT_practice

```python
# TODO: train
for epoch in range(args.epochs):
    for src_batch, tgt_batch in train_loader:
        pass

    # TODO: validation
    for src_batch, tgt_batch in valid_loader:
        pass
```

```python
for src_batch, tgt_batch in test_loader:
    # TODO: predict pred_batch from src_b
    pred_batch = tgt_batch

    # every sentences in pred_batch shoul
    # every <pad> token (index: 2) should
    # example of pred_batch:
    # [[0, 5, 6, 7, 1],
    #  [0, 4, 9, 1, 2],
    #  [0, 6, 1, 2, 2]]
```
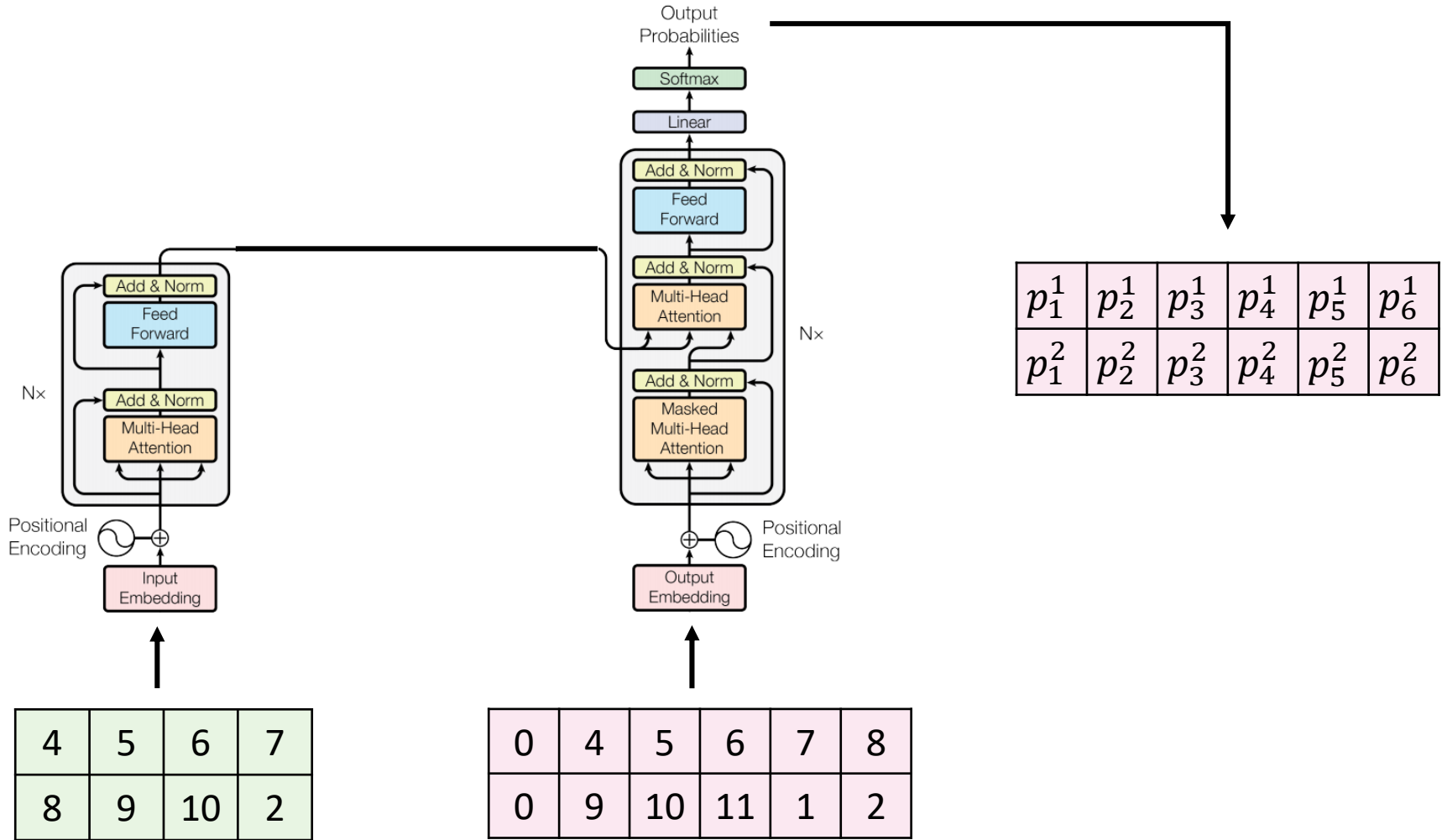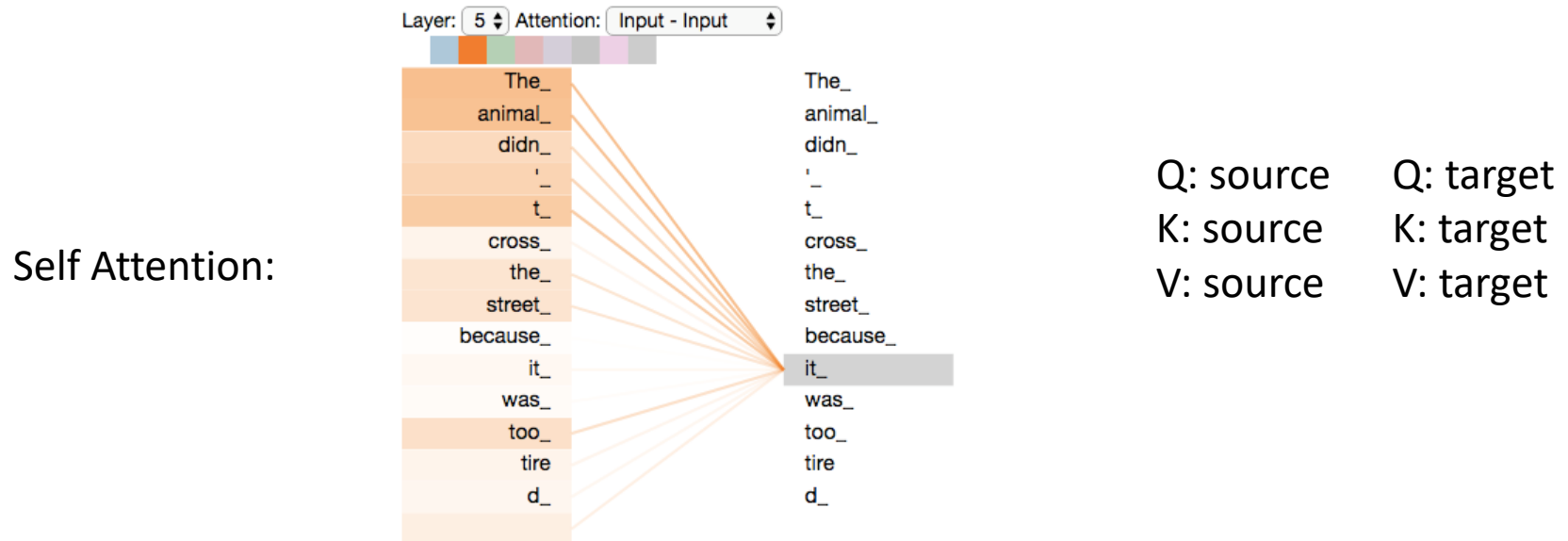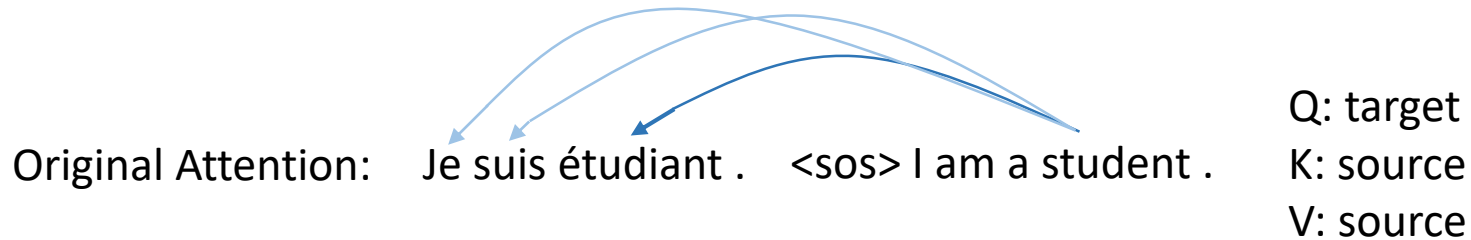
# Table of Contents

# Abstract View of Transformer Architecture

# Self Attention

Original Attention: Je suis étudiant . &lt;sos&gt; I am a student .

Q: target
K: source
V: source

Self Attention:



```
Layer: 5 ▲▼   Attention: Input - Input ▲▼
```

| | |
|---|---|
| The_ | The_ |
| animal_ | animal_ |
| didn_ | didn_ |
| '_ | '_ |
| t_ | t_ |
| cross_ | cross_ |
| the_ | the_ |
| street_ | street_ |
| because_ | because_ |
| it_ | it_ |
| was_ | was_ |
| too_ | too_ |
| tire | tire |
| d_ | d_ |

Q: source    Q: target
K: source    K: target
V: source    V: target

# Multi-Head Attention

$\alpha_1$     $\alpha_2$

$\alpha_3$

Original Attention:    Je suis étudiant .    <sos> I am a student .

Q: target
K: source
V: source

Attention for student: $\alpha_1 \times$ Je + $\alpha_2 \times$ suis + $\alpha_3 \times$ étudiant



Scaled Dot-Product Attention

Multi-Head Attention

# Abstract View of Transformer Architecture



| $p_1^1$ | $p_2^1$ | $p_3^1$ | $p_4^1$ | $p_5^1$ | $p_6^1$ |
|---|---|---|---|---|---|
| $p_1^2$ | $p_2^2$ | $p_3^2$ | $p_4^2$ | $p_5^2$ | $p_6^2$ |

| 4 | 5 | 6 | 7 |
|---|---|---|---|
| 8 | 9 | 10 | 2 |

| 0 | 4 | 5 | 6 | 7 | 8 |
|---|---|---|---|---|---|
| 0 | 9 | 10 | 11 | 1 | 2 |

# Table of Contents

# Attention is a Weighted Sum

- Query: 주인공 문장.

- Key, Value: 주인공 문장이 보는 문장.

- Query는 Key와 연산하여 weight를 구한다.



$$w_{11} \qquad\qquad w_{12} \qquad\qquad w_{13}$$

# Attention is a Weighted Sum

- $d_k$로 나눠주고 원치 않는 정보를 masking을 통해 지운 후 Softmax를 취한다.

$$\text{Softmax}(\frac{w_{11}}{d_k}, \cancel{\frac{w_{12}}{d_k}}, \frac{w_{13}}{d_k}) = (\alpha_{11}, 0, \alpha_{13})$$

$$\text{Softmax}(\frac{w_{21}}{d_k}, \frac{w_{22}}{d_k}, \cancel{\frac{w_{23}}{d_k}}) = (\alpha_{21}, \alpha_{22}, 0)$$

- Value를 대상으로 Weighted Sum을 한다.

$\alpha_{11} \cdot$ [ ][ ][ ][ ] $+ \; 0 \; \cdot$ [ ][ ][ ][ ] $+ \alpha_{13} \cdot$ [ ][ ][ ][ ]

$\alpha_{21} \cdot$ [ ][ ][ ][ ] $+ \alpha_{22} \cdot$ [ ][ ][ ][ ] $+ \; 0 \; \cdot$ [ ][ ][ ][ ]
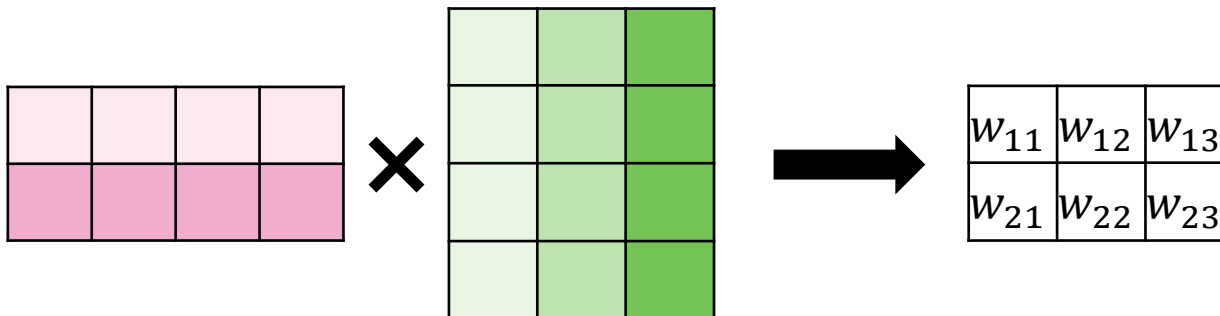
- 최종 output.

[ ][ ][ ][ ]　　　　[ ][ ][ ][ ]

# Attention in Matrix Form

- Query: 주인공 문장.

- Key, Value: 주인공 문장이 보는 문장.

- Query는 Key와 연산하여 weight를 구한다.

# Attention in Matrix Form

- $d_k$로 나눠주고 원치 않는 정보를 masking을 통해 지운 후 Softmax를 취한다.

$$\text{Softmax}\left( \frac{\begin{array}{|c|c|c|} \hline w_{11} & w_{12} & w_{13} \\ \hline w_{21} & w_{22} & w_{23} \\ \hline \end{array}}{d_k} \right) \quad \Longrightarrow \quad \begin{array}{|c|c|c|} \hline \alpha_{11} & 0 & \alpha_{13} \\ \hline \alpha_{21} & \alpha_{22} & 0 \\ \hline \end{array}$$

- Value를 대상으로 Weighted Sum을 한다.

$$\begin{array}{|c|c|c|} \hline \alpha_{11} & 0 & \alpha_{13} \\ \hline \alpha_{21} & \alpha_{22} & 0 \\ \hline \end{array} \quad \times$$

- 최종 output.

# Attention in Implementation

- Query: 주인공 문장.



- Key, Value: 주인공 문장이 보는 문장.

# Multi-Head Attention

- Query

- Key, Value

- Concatenate Outputs

**Out1   Out2   Out3**

- Match the dimension

# Multi-Head Attention

- Multi_Head_Attention(Query, Key, Value, Mask)


Multi_Head_Attention (



$=$



- 위 example에서는 단어 벡터의 차원이 4, Head가 3개, 각 Head의 차원이 3.
- 보통은 단어 벡터의 차원이 512, Head가 8개, 각 Head의 차원이 64.

# Dive into the Transformer Architecture



$$
\begin{array}{|c|c|c|c|c|c|}
\hline
p_1^1 & p_2^1 & p_3^1 & p_4^1 & p_5^1 & p_6^1 \\
\hline
p_1^2 & p_2^2 & p_3^2 & p_4^2 & p_5^2 & p_6^2 \\
\hline
\end{array}
$$

| 4 | 5 | 6 | 7 |
|---|---|---|---|
| 8 | 9 | 10 | 2 |

| 0 | 4 | 5 | 6 | 7 | 8 |
|---|---|---|---|---|---|
| 0 | 9 | 10 | 11 | 1 | 2 |

# Input Embedding



$$\begin{array}{|c|c|c|c|c|c|} \hline p_1^1 & p_2^1 & p_3^1 & p_4^1 & p_5^1 & p_6^1 \\ \hline p_1^2 & p_2^2 & p_3^2 & p_4^2 & p_5^2 & p_6^2 \\ \hline \end{array}$$

| 4 | 5 | 6 | 7 |
|---|---|---|---|
| 8 | 9 | 10 | 2 |

| 0 | 4 | 5 | 6 | 7 | 8 |
|---|---|---|---|---|---|
| 0 | 9 | 10 | 11 | 1 | 2 |

one-hot vector

# Input Embedding

dim=512

Linear map

# Positional Encoding



Output
Probabilities

Softmax

Linear

Add & Norm

Feed
Forward

Add & Norm

Multi-Head
Attention

Add & Norm

Masked
Multi-Head
Attention

Nx

Positional
Encoding

Output
Embedding

Add & Norm

Feed
Forward

Add & Norm

Multi-Head
Attention

Nx

Positional
Encoding

Input
Embedding

| $p_1^1$ | $p_2^1$ | $p_3^1$ | $p_4^1$ | $p_5^1$ | $p_6^1$ |
|---|---|---|---|---|---|
| $p_1^2$ | $p_2^2$ | $p_3^2$ | $p_4^2$ | $p_5^2$ | $p_6^2$ |

| 4 | 5 | 6 | 7 |
|---|---|---|---|
| 8 | 9 | 10 | 2 |

| 0 | 4 | 5 | 6 | 7 | 8 |
|---|---|---|---|---|---|
| 0 | 9 | 10 | 11 | 1 | 2 |

# Positional Encoding

word        embedding                  Positional Encoding

Je   $+ \; \cos(1) \; \sin\left(\frac{1}{10}\right) \cos\left(\frac{1}{10}\right) \sin\left(\frac{1}{10^2}\right) \cos\left(\frac{1}{10^2}\right) \sin\left(\frac{1}{10^3}\right) \cos\left(\frac{1}{10^3}\right) \sin\left(\frac{1}{10^4}\right)$

suis   $+ \; \cos(2) \; \sin\left(\frac{2}{10}\right) \cos\left(\frac{2}{10}\right) \sin\left(\frac{2}{10^2}\right) \cos\left(\frac{2}{10^2}\right) \sin\left(\frac{2}{10^3}\right) \cos\left(\frac{2}{10^3}\right) \sin\left(\frac{2}{10^4}\right)$

étudiant   $+ \; \cos(3) \; \sin\left(\frac{3}{10}\right) \cos\left(\frac{3}{10}\right) \sin\left(\frac{3}{10^2}\right) \cos\left(\frac{3}{10^2}\right) \sin\left(\frac{3}{10^3}\right) \cos\left(\frac{3}{10^3}\right) \sin\left(\frac{3}{10^4}\right)$

# Positional Encoding

| word | embedding | Positional Encoding |
|------|-----------|---------------------|

Je $\blacksquare\blacksquare\blacksquare\blacksquare\blacksquare\blacksquare\blacksquare\blacksquare$ ➕ $\cos(1) \quad \sin\left(\frac{1}{10}\right) \quad \cos\left(\frac{1}{10}\right) \quad \sin\left(\frac{1}{10^2}\right) \quad \cos\left(\frac{1}{10^2}\right) \quad \sin\left(\frac{1}{10^3}\right) \quad \cos\left(\frac{1}{10^3}\right) \quad \sin\left(\frac{1}{10^4}\right)$

suis $\blacksquare\blacksquare\blacksquare\blacksquare\blacksquare\blacksquare\blacksquare\blacksquare$ ➕ $\cos(2) \quad \sin\left(\frac{2}{10}\right) \quad \cos\left(\frac{2}{10}\right) \quad \sin\left(\frac{2}{10^2}\right) \quad \cos\left(\frac{2}{10^2}\right) \quad \sin\left(\frac{2}{10^3}\right) \quad \cos\left(\frac{2}{10^3}\right) \quad \sin\left(\frac{2}{10^4}\right)$

étudiant $\blacksquare\blacksquare\blacksquare\blacksquare\blacksquare\blacksquare\blacksquare\blacksquare$ ➕ $\cos(3) \quad \sin\left(\frac{3}{10}\right) \quad \cos\left(\frac{3}{10}\right) \quad \sin\left(\frac{3}{10^2}\right) \quad \cos\left(\frac{3}{10^2}\right) \quad \sin\left(\frac{3}{10^3}\right) \quad \cos\left(\frac{3}{10^3}\right) \quad \sin\left(\frac{3}{10^4}\right)$

Je

➕

| $\cos(1)$ | $\sin\left(\frac{1}{10}\right)$ | $\cos\left(\frac{1}{10}\right)$ | $\sin\left(\frac{1}{10^2}\right)$ | $\cos\left(\frac{1}{10^2}\right)$ | $\sin\left(\frac{1}{10^3}\right)$ | $\cos\left(\frac{1}{10^3}\right)$ | $\sin\left(\frac{1}{10^4}\right)$ |
|---|---|---|---|---|---|---|---|

# Positional Encoding

word      embedding            Positional Encoding

Je     ➕ $\cos(1) \; \sin\left(\frac{1}{10}\right) \cos\left(\frac{1}{10}\right) \sin\left(\frac{1}{10^2}\right) \cos\left(\frac{1}{10^2}\right) \sin\left(\frac{1}{10^3}\right) \cos\left(\frac{1}{10^3}\right) \sin\left(\frac{1}{10^4}\right)$

suis     ➕ $\cos(2) \; \sin\left(\frac{2}{10}\right) \cos\left(\frac{2}{10}\right) \sin\left(\frac{2}{10^2}\right) \cos\left(\frac{2}{10^2}\right) \sin\left(\frac{2}{10^3}\right) \cos\left(\frac{2}{10^3}\right) \sin\left(\frac{2}{10^4}\right)$

étudiant     ➕ $\cos(3) \; \boxed{\sin\left(\frac{3}{10}\right)} \cos\left(\frac{3}{10}\right) \sin\left(\frac{3}{10^2}\right) \cos\left(\frac{3}{10^2}\right) \sin\left(\frac{3}{10^3}\right) \cos\left(\frac{3}{10^3}\right) \sin\left(\frac{3}{10^4}\right)$

(3,2)

$$pos = 3$$
$$2i = 2$$
$$d_{model} = 8$$

⬇

$$PE_{pos,2i} = \sin\left(\frac{pos}{10000^{\frac{2i}{d_{model}}}}\right) \quad \Rightarrow \quad PE_{3,2} = \sin\left(\frac{3}{10000^{\frac{2}{8}}}\right) = \sin\left(\frac{3}{10}\right)$$
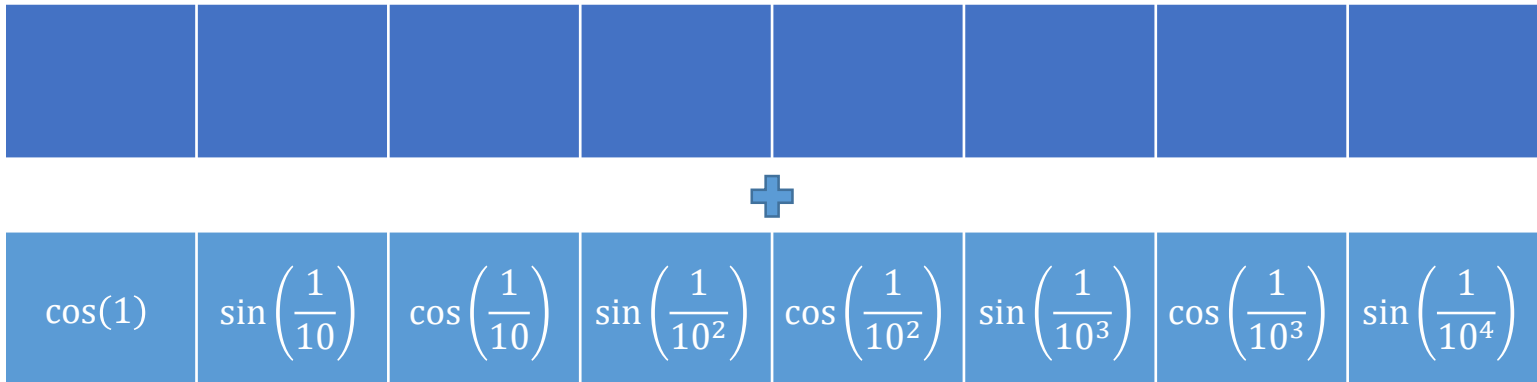
# Positional Encoding

word          embedding                          Positional Encoding

Je
$$\cos(1)\ \sin\left(\frac{1}{10}\right)\cos\left(\frac{1}{10}\right)\sin\left(\frac{1}{10^2}\right)\cos\left(\frac{1}{10^2}\right)\sin\left(\frac{1}{10^3}\right)\cos\left(\frac{1}{10^3}\right)\sin\left(\frac{1}{10^4}\right)$$

$(2,5)$

suis
$$\cos(2)\ \sin\left(\frac{2}{10}\right)\cos\left(\frac{2}{10}\right)\sin\left(\frac{2}{10^2}\right)\boxed{\cos\left(\frac{2}{10^2}\right)}\sin\left(\frac{2}{10^3}\right)\cos\left(\frac{2}{10^3}\right)\sin\left(\frac{2}{10^4}\right)$$

étudiant
$$\cos(3)\ \sin\left(\frac{3}{10}\right)\cos\left(\frac{3}{10}\right)\sin\left(\frac{3}{10^2}\right)\cos\left(\frac{3}{10^2}\right)\sin\left(\frac{3}{10^3}\right)\cos\left(\frac{3}{10^3}\right)\sin\left(\frac{3}{10^4}\right)$$
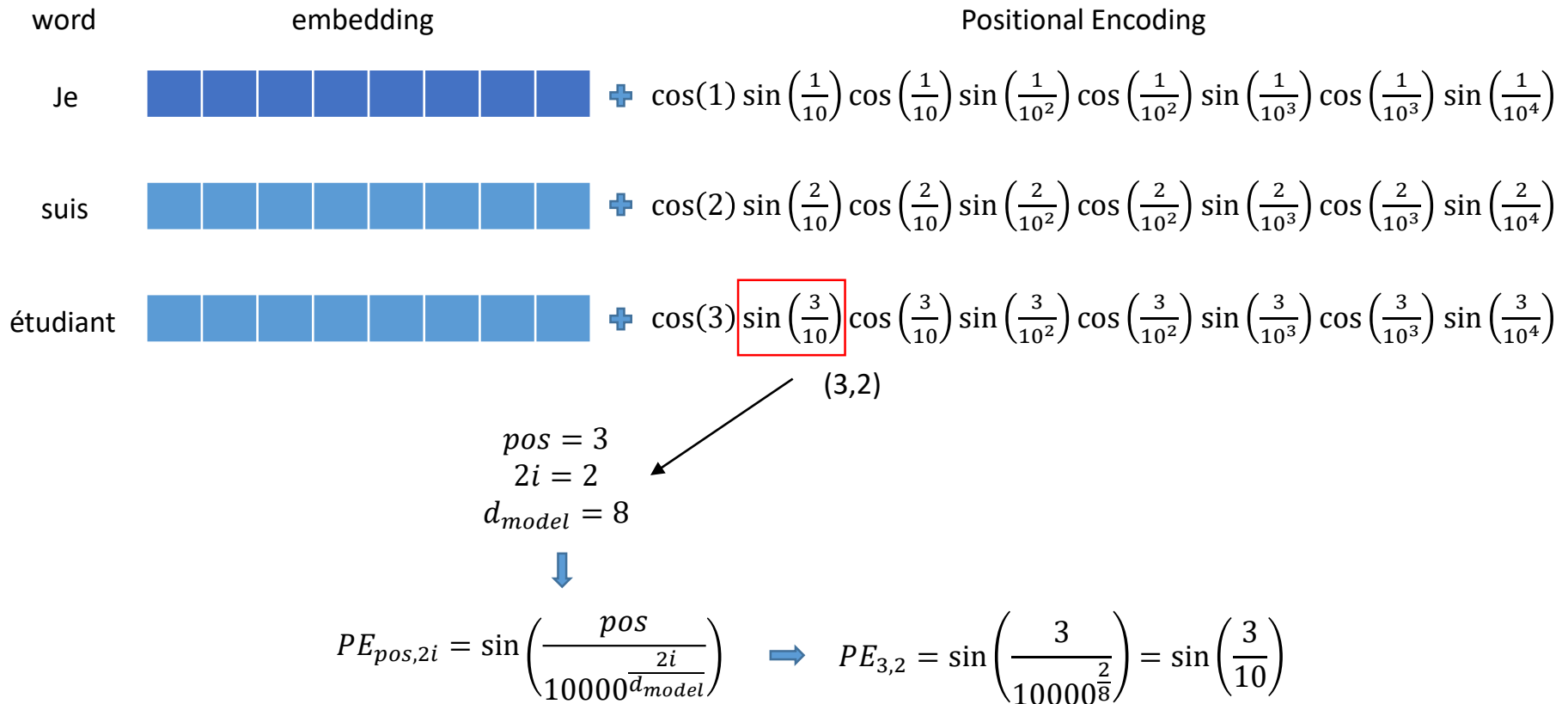
$$pos = 3$$
$$2i + 1 = 5$$
$$d_{model} = 8$$

$$PE_{pos,2i+1} = \cos\left(\frac{pos}{10000^{\frac{2i}{d_{model}}}}\right) \quad \Rightarrow \quad PE_{2,5} = \cos\left(\frac{3}{10000^{\frac{4}{8}}}\right) = \cos\left(\frac{2}{10^2}\right)$$

# Positional Encoding

# Positional Encoding

# Positional Encoding

# Positional Encoding

word        embedding        Positional Encoding

Je $\boxed{\phantom{xxxxxxxx}}$ ➕ $\cos(1)\ \sin\left(\frac{1}{10}\right) \cos\left(\frac{1}{10}\right) \sin\left(\frac{1}{10^2}\right) \cos\left(\frac{1}{10^2}\right) \sin\left(\frac{1}{10^3}\right) \cos\left(\frac{1}{10^3}\right) \sin\left(\frac{1}{10^4}\right)$

suis $\boxed{\phantom{xxxxxxxx}}$ ➕ $\cos(2)\ \sin\left(\frac{2}{10}\right) \cos\left(\frac{2}{10}\right) \sin\left(\frac{2}{10^2}\right) \cos\left(\frac{2}{10^2}\right) \sin\left(\frac{2}{10^3}\right) \cos\left(\frac{2}{10^3}\right) \sin\left(\frac{2}{10^4}\right)$

étudiant $\boxed{\phantom{xxxxxxxx}}$ ➕ $\cos(3)\ \sin\left(\frac{3}{10}\right) \cos\left(\frac{3}{10}\right) \sin\left(\frac{3}{10^2}\right) \cos\left(\frac{3}{10^2}\right) \sin\left(\frac{3}{10^3}\right) \cos\left(\frac{3}{10^3}\right) \sin\left(\frac{3}{10^4}\right)$
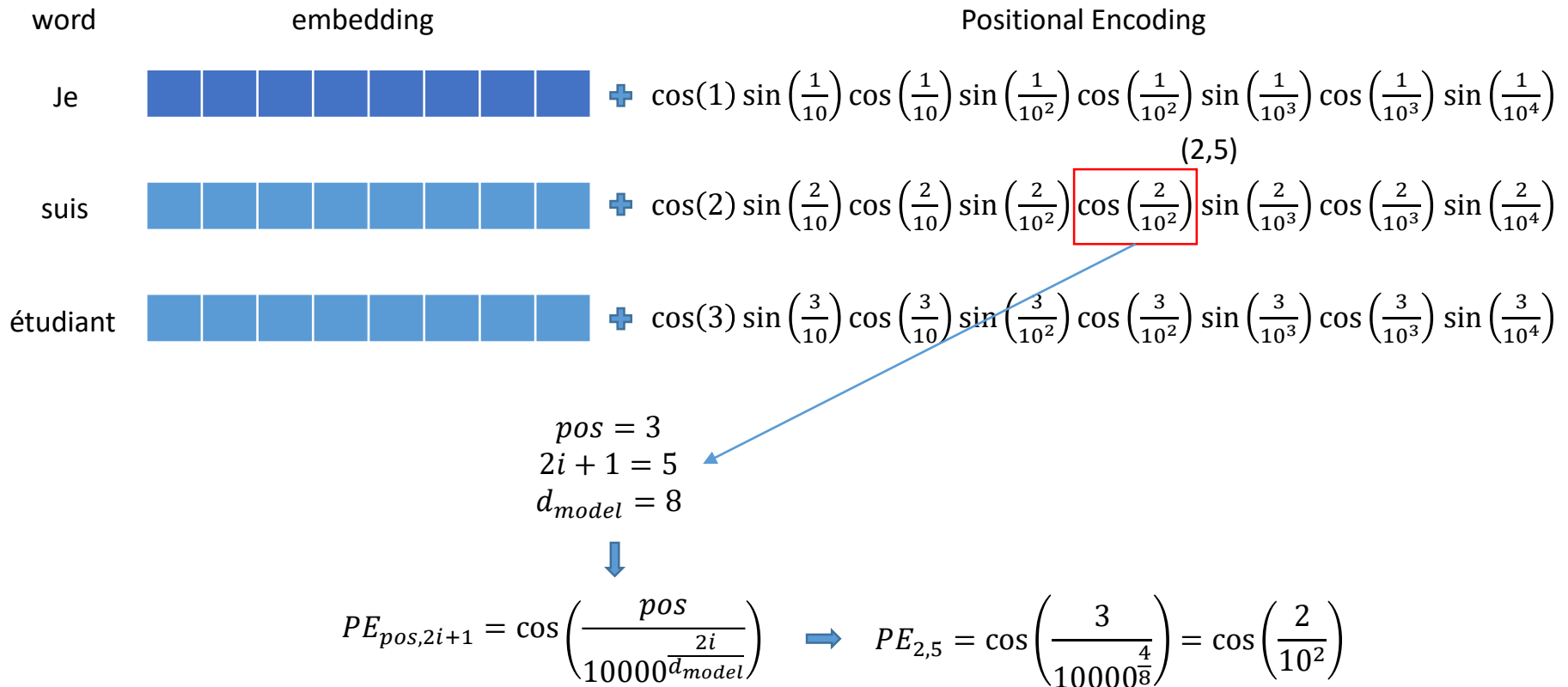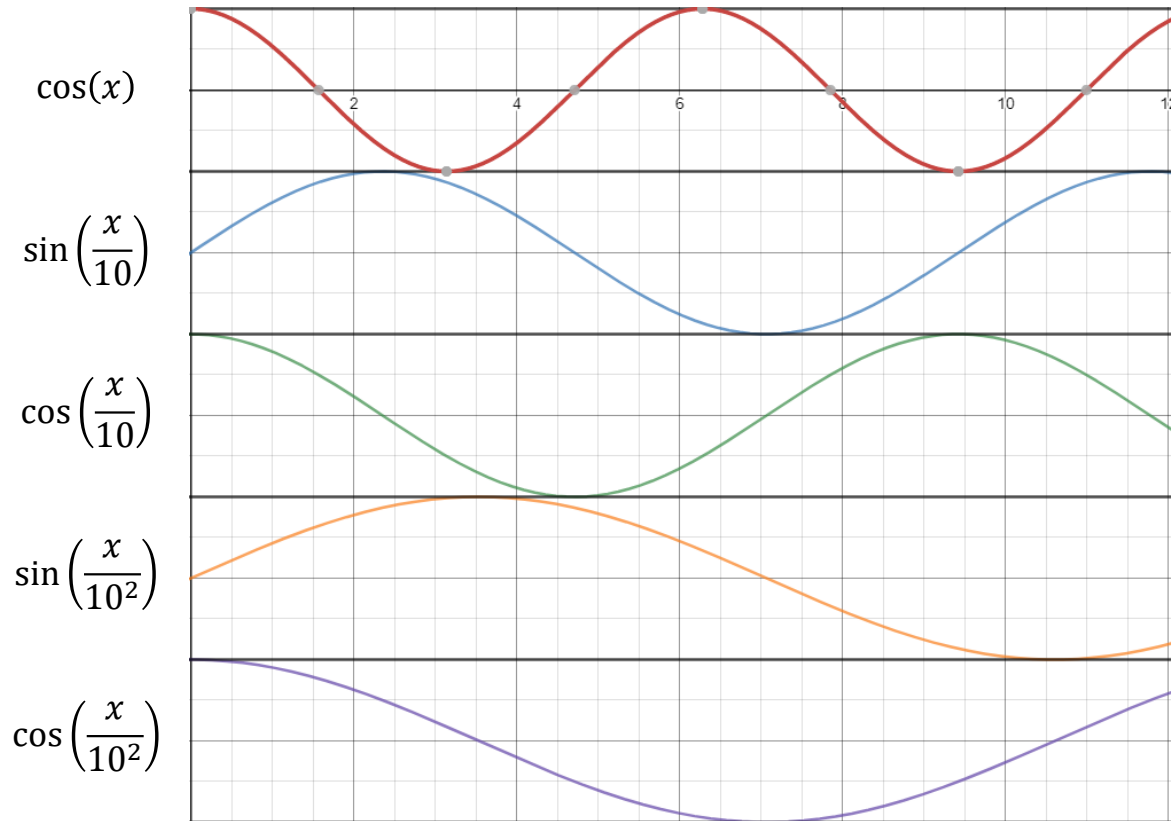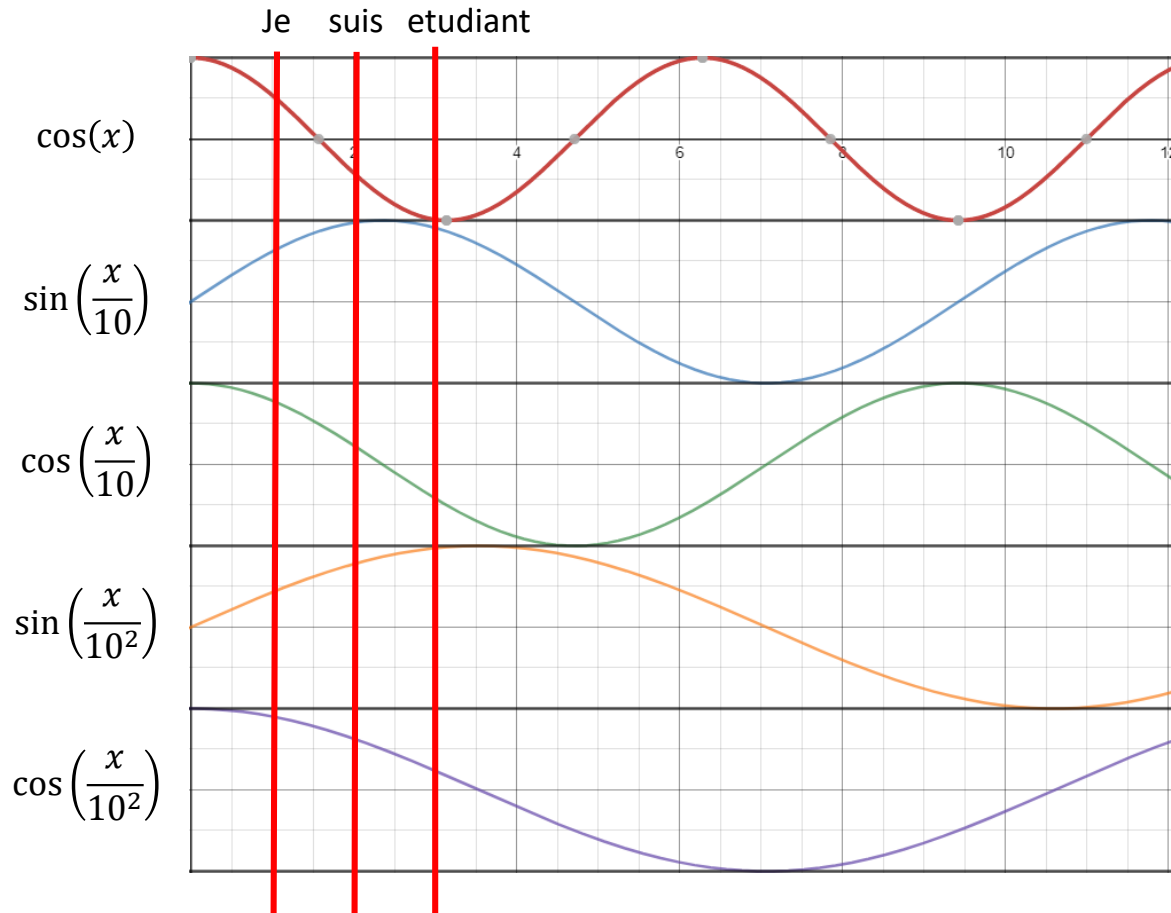
➕

| $\cos(1)$ | $\sin\left(\frac{1}{10}\right)$ | $\cos\left(\frac{1}{10}\right)$ | $\sin\left(\frac{1}{10^2}\right)$ | $\cos\left(\frac{1}{10^2}\right)$ | $\sin\left(\frac{1}{10^3}\right)$ | $\cos\left(\frac{1}{10^3}\right)$ | $\sin\left(\frac{1}{10^4}\right)$ |

# Multi-Head Attention (Self Attention)



$$
\begin{array}{|c|c|c|c|c|c|}
\hline
p_1^1 & p_2^1 & p_3^1 & p_4^1 & p_5^1 & p_6^1 \\
\hline
p_1^2 & p_2^2 & p_3^2 & p_4^2 & p_5^2 & p_6^2 \\
\hline
\end{array}
$$

| 4 | 5 | 6 | 7 |
|---|---|---|---|
| 8 | 9 | 10 | 2 |

| 0 | 4 | 5 | 6 | 7 | 8 |
|---|---|---|---|---|---|
| 0 | 9 | 10 | 11 | 1 | 2 |

# Multi-Head Attention (Self Attention)

Multi_Head_Attention (

| | | | |
|---|---|---|---|
| | | | |
| | | | |
| | | | |

| | | | |
|---|---|---|---|
| | | | |
| | | | |
| | | | |

| | | | |
|---|---|---|---|
| | | | |
| | | | |
| | | | |

| 0 | 0 | 0 | 0 |
|---|---|---|---|
| 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 0 |

)

Multi_Head_Attention (

| | | | |
|---|---|---|---|
| | | | |
| | | | |
| | | | |

| | | | |
|---|---|---|---|
| | | | |
| | | | |
| | | | |

| | | | |
|---|---|---|---|
| | | | |
| | | | |
| | | | |

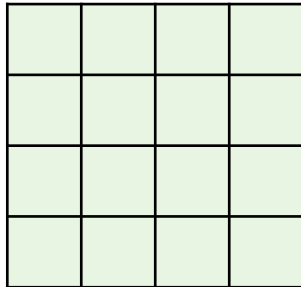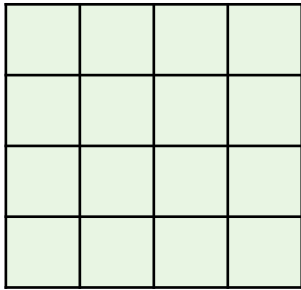| 0 | 0 | 0 | 1 |
|---|---|---|---|
| 0 | 0 | 0 | 1 |
| 0 | 0 | 0 | 1 |
| 0 | 0 | 0 | 1 |

)

# Position-wise Feed-Forward Networks

# Position-wise Feed-Forward Networks

Position-wise Feed-Forward Networks

ReLU + DropOut

512    2048    512

# Output Embedding

# Positional Encoding



| $p_1^1$ | $p_2^1$ | $p_3^1$ | $p_4^1$ | $p_5^1$ | $p_6^1$ |
|---------|---------|---------|---------|---------|---------|
| $p_1^2$ | $p_2^2$ | $p_3^2$ | $p_4^2$ | $p_5^2$ | $p_6^2$ |

| 4 | 5 | 6 | 7 |
|---|---|---|---|
| 8 | 9 | 10 | 2 |

| 0 | 4 | 5 | 6 | 7 | 8 |
|---|---|---|---|---|---|
| 0 | 9 | 10 | 11 | 1 | 2 |

# Multi-Head Attention (Masked Self Attention)



| $p_1^1$ | $p_2^1$ | $p_3^1$ | $p_4^1$ | $p_5^1$ | $p_6^1$ |
|---------|---------|---------|---------|---------|---------|
| $p_1^2$ | $p_2^2$ | $p_3^2$ | $p_4^2$ | $p_5^2$ | $p_6^2$ |

| 4 | 5 | 6 | 7 |
|---|---|---|---|
| 8 | 9 | 10 | 2 |

| 0 | 4 | 5 | 6 | 7 | 8 |
|---|---|---|---|---|---|
| 0 | 9 | 10 | 11 | 1 | 2 |

# Multi-Head Attention (Masked Self Attention)

Multi_Head_Attention (

| | | | | | |
|---|---|---|---|---|---|
| 0 | 1 | 1 | 1 | 1 | 1 |
| 0 | 0 | 1 | 1 | 1 | 1 |
| 0 | 0 | 0 | 1 | 1 | 1 |
| 0 | 0 | 0 | 0 | 1 | 1 |
| 0 | 0 | 0 | 0 | 0 | 1 |
| 0 | 0 | 0 | 0 | 0 | 0 |

)

Multi_Head_Attention (

| | | | | | |
|---|---|---|---|---|---|
| 0 | 1 | 1 | 1 | 1 | 1 |
| 0 | 0 | 1 | 1 | 1 | 1 |
| 0 | 0 | 0 | 1 | 1 | 1 |
| 0 | 0 | 0 | 0 | 1 | 1 |
| 0 | 0 | 0 | 0 | 0 | 1 |
| 0 | 0 | 0 | 0 | 0 | 1 |

)

# Multi-Head Attention (Original Attention)



| $p_1^1$ | $p_2^1$ | $p_3^1$ | $p_4^1$ | $p_5^1$ | $p_6^1$ |
|---------|---------|---------|---------|---------|---------|
| $p_1^2$ | $p_2^2$ | $p_3^2$ | $p_4^2$ | $p_5^2$ | $p_6^2$ |

| 4 | 5 | 6 | 7 |
|---|---|---|---|
| 8 | 9 | 10 | 2 |

| 0 | 4 | 5 | 6 | 7 | 8 |
|---|---|---|---|---|---|
| 0 | 9 | 10 | 11 | 1 | 2 |

# Multi-Head Attention (Original Attention)

Multi_Head_Attention (

| 0 | 0 | 0 | 0 |
|---|---|---|---|
| 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 0 |

)

Multi_Head_Attention (

| 0 | 0 | 0 | 1 |
|---|---|---|---|
| 0 | 0 | 0 | 1 |
| 0 | 0 | 0 | 1 |
| 0 | 0 | 0 | 1 |
| 0 | 0 | 0 | 1 |
| 0 | 0 | 0 | 1 |

)

# Position-wise Feed-Forward Networks



| $p_1^1$ | $p_2^1$ | $p_3^1$ | $p_4^1$ | $p_5^1$ | $p_6^1$ |
|---|---|---|---|---|---|
| $p_1^2$ | $p_2^2$ | $p_3^2$ | $p_4^2$ | $p_5^2$ | $p_6^2$ |

| 4 | 5 | 6 | 7 |
|---|---|---|---|
| 8 | 9 | 10 | 2 |

| 0 | 4 | 5 | 6 | 7 | 8 |
|---|---|---|---|---|---|
| 0 | 9 | 10 | 11 | 1 | 2 |

# Prediction



$$
\begin{array}{|c|c|c|c|c|c|}
\hline
p_1^1 & p_2^1 & p_3^1 & p_4^1 & p_5^1 & p_6^1 \\
\hline
p_1^2 & p_2^2 & p_3^2 & p_4^2 & p_5^2 & p_6^2 \\
\hline
\end{array}
$$

| 4 | 5 | 6 | 7 |
|---|---|---|---|
| 8 | 9 | 10 | 2 |

| 0 | 4 | 5 | 6 | 7 | 8 |
|---|---|---|---|---|---|
| 0 | 9 | 10 | 11 | 1 | 2 |

# References

1) Attention Is All You Need (NeurIPS'17)

    - https://papers.nips.cc/paper/7181-attention-is-all-you-need.pdf

2) Neural Machine Translation by Jointly Learning to Align and Translate (ICLR'15)

    - https://arxiv.org/pdf/1409.0473.pdf

3) Transformers From Scratch

    - http://www.peterbloem.nl/blog/transformers

4) The Illustrated Transformer

    - https://jalammar.github.io/illustrated-transformer/

5) The Annotated Transformer

    - https://nlp.seas.harvard.edu/2018/04/03/attention.html

6) A Brief Overview of Attention Mechanism

    - https://medium.com/syncedreview/a-brief-overview-of-attention-mechanism-13c578ba9129

# Thank you !

## Any Questions ?