

김동하 - Transformer & Bert 유튜브 영상보고 transformer 추가공부

Transformer, bert → downstream task 향

• downstream obj?

- upstream & downstream relationship

↳ component이 어떤 성능하는지를 알아가는 방향

Upstream ←

Product

Component

Retail / Marketing

→ downstream

↳ component를 어떤 방식

↓ NLP 향

Upstream : word Embedding 향 어떤 영향되는지

←

Word Analogy

Word Embedding

Translation

→ downstream

: Appel translational 적용시키는

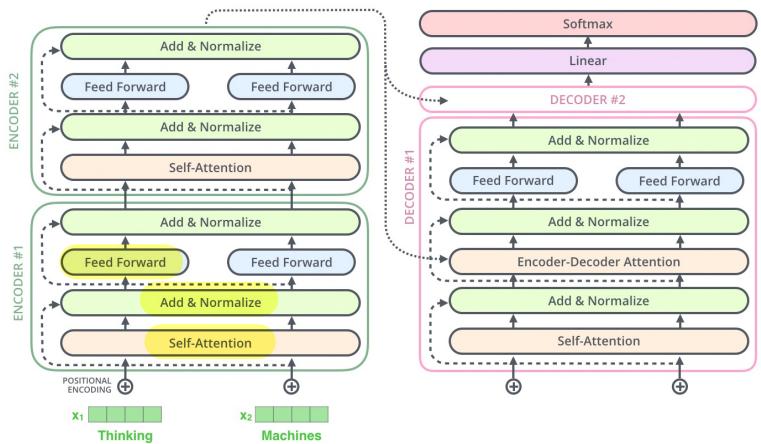
(In CS224N)

[Upstream : Intrinsic (내재)한 향

↓ downstream : Extrinsic (외재)한 향

Transformer

middle
encoder #1, #2



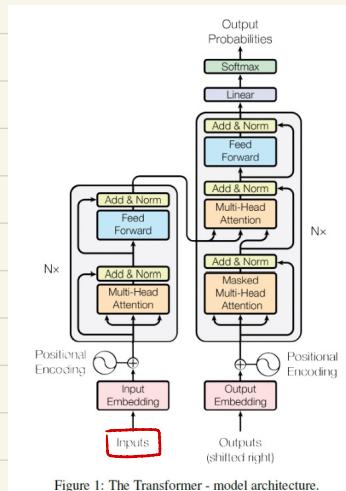
Downstream Task for NLP Translation

인수의 쌍 - Pair Sentence translation은 2 pair sentence!

Attention 헤드 key, Query, Value

K, Q, V는 Encoder-Decoder에서 각각 사용된다.

문장 번역은 디코더입니다.
입력은 단일 텐서가 됩니다.
→ 단일 텐서를 처리하는 디코더를 사용하고
마지막 padding 토큰을 처리하는 디코더를 사용합니다.



Input은 token right : Token Embedding

- Max length = 6 일때 → max sequence length 6 token.

- Sent1 : i am looking for happiness
- Sent2 : I am suhyun
- Sent3 : I am ①*#\$

→ 어떤가 vocabulary를 확장할 수 있는가.

- 인공지, 디코더를 갖기 vocab을 확장

-特殊的 token을 추가 : <PAD>, <UNK>, <EOS>, <BOS>

기본을 확장해주는 풀어쓰는 단어로 사용되는

단어로 사용되는

Figure 1: The Transformer - model architecture.

- Encoder, max_length = 6000H.

Sent1: I am looking for happiness

{ $x_1, x_2, x_3, x_4, x_5, \langle /S \rangle$ }

문장의 끝 token

$\langle /S \rangle$

마지막 토큰에 대한 인덱스

Sent2: I am suhyun

{ $x_1, x_2, x_3, \langle /S \rangle, \langle PAD \rangle, \langle PAD \rangle$ }

Sent3: I am @#\$%

{ $x_1, x_2, \langle UNK \rangle, \langle /S \rangle, \langle PAD \rangle, \langle PAD \rangle$ }

- Decoder, max_length = 6000H.

Sent1: 우는 딸의 흐느낌

문장의 토큰을

로그인하여 풀어보기

{ $\langle S \rangle, x_1, x_2, x_3, x_4, \langle /S \rangle$ }

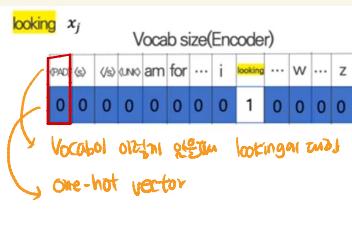
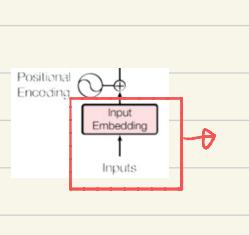
Sent2: 나를 사랑해

{ $\langle S \rangle, x_1, x_2, \langle /S \rangle, \langle PAD \rangle, \langle PAD \rangle$ }

Sent3: 우는 @#\$%

{ $\langle S \rangle, x_1, \langle UNK \rangle, \langle /S \rangle, \langle PAD \rangle, \langle PAD \rangle$ }

> I am looking for happiness { $x_1, x_2, x_3, x_4, x_5, \langle /S \rangle$ }



Embed size

$\langle PAD \rangle$	0.0	0.0	0.0	0.0	0.0
:	0.0	0.3	0.7	0.0	
looking	0.2	0.1	0.7	0.2	
:	0.9	0.1	0.7	0.9	
z	0.5	0.7	0.1	0.5	
	0.2	0.9	0.4	0.2	
	0.2	0.0	0.7	0.2	

0.2 0.1 0.7 0.2 w_j looking
↳ input embeddings은 절대
But 원래에는 주파수로 scaling

□: $\langle PAD \rangle$ 같은 경우는 같은 단어에 대해서도 같은 토큰으로

구별하지 않고!

don'ts vector
비슷한 토큰

- Token embedding Scaling

$$[0.2 \ 0.1 \ 0.7 \ 0.2] \times \frac{1}{\text{Embed size}}$$

: 흥미 gradient를 잘.. 하기 위함

- positional encoding

: 임베딩이 시각적 주제를 추가하기 위해 사용

방법 2번이 ① Xarker initialize

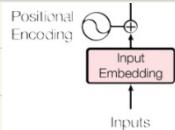
② Sinusoid initialize

① Xavier initialize

Position vocab					
3rd position	0	1	2	3	4
"looking"	0	0	1	0	0

Embed size	0.1	0.4	0.7	0.1
0.0	0.0	0.3	0.7	0.0
0.1	0.1	0.1	0.9	0.1
...
0.2	0.0	0.7	0.2	0.0

3th position embedding
0.1 0.1 0.9 0.1



: positional targ one-hot vector \in weight matrixer 係り受け position embedding 係り受け

: **Input position embedding**은 **Input Embedding**의 **Output Embedding**입니다.

② Sinusoid initializer

$$PE(pos_{2i}) = \sin(pos / 10^{4-2i/d_{model}})$$

$$PE(pos, \omega_{i+1}) = \cos(pos / 10^{42i/\text{dimodel}})$$

→ d_{model} : embedding size

i : embedding feature 영역

pos : 단어의 인덱스 → 총 6개의 index

```

import numpy as np
import matplotlib.pyplot as plt

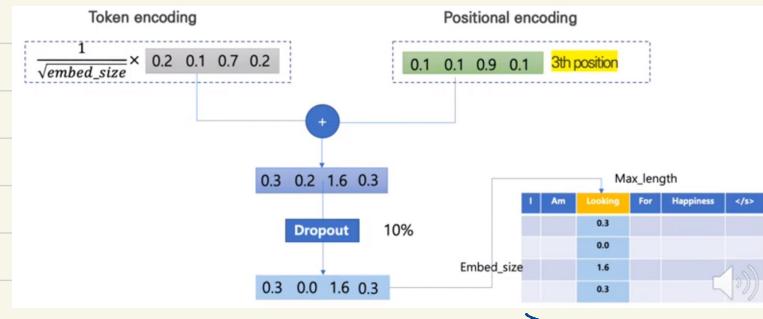
max_length = 512
embed_size = 1024
position_enc = np.array([
    [pos / np.power(10000, 2*i/embed_size) for i in range(embed_size)]
    if pos != 0 else np.zeros(embed_size) for pos in range(max_length)])

position_enc[:, 0:2] = np.sin(position_enc[:, 0:2]) # dim 2i
position_enc[:, 1:2] = np.cos(position_enc[:, 1:2]) # dim 2i+1

img = plt.imshow(position_enc, cmap='hot', interpolation='nearest', aspect='auto')
cmap = plt.get_cmap('hot')
plt.colorbar(img, cmap=cmap)
plt.show()

```

• Dropout

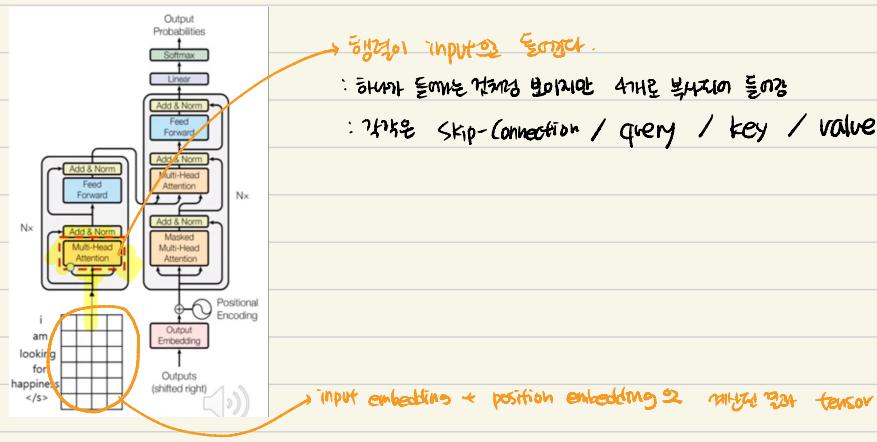


; feature의 개수를 줄이는 것. → overfitting을 방지하기 위해

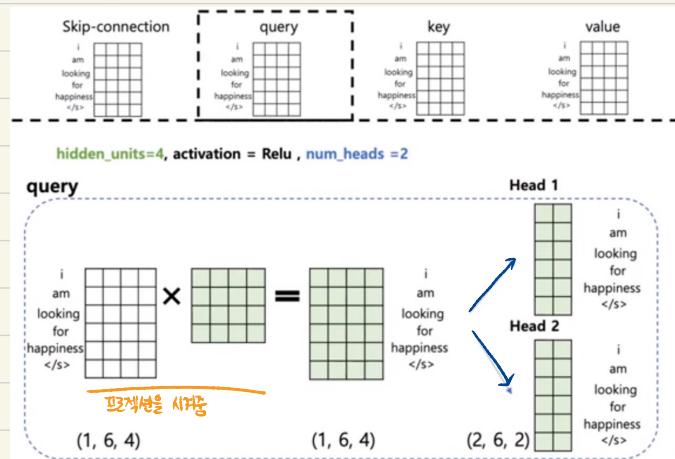
↳ layer가 많아질 때 백버러는 대표적인 문제점

26 clones examined

• Multi-head Attention



① Query

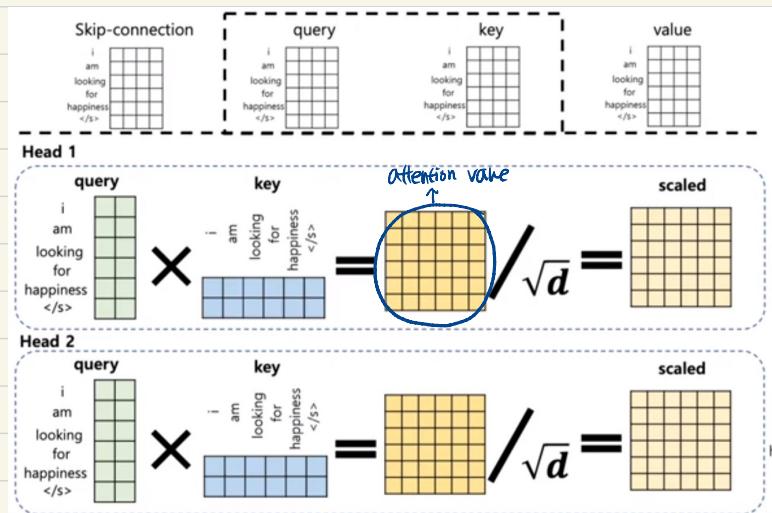


각각의 이유 = 한 번의 훈련으로 나눈 이유
: attention weight를 좀 더 다양하게
할 수 있는 목적으로 추가함

- ② key
- ③ value

\rightarrow Query와 같은 맥락!

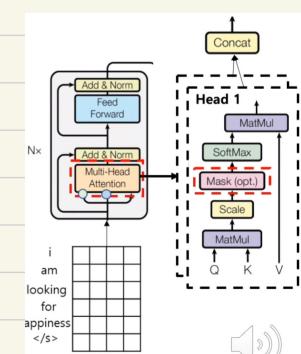
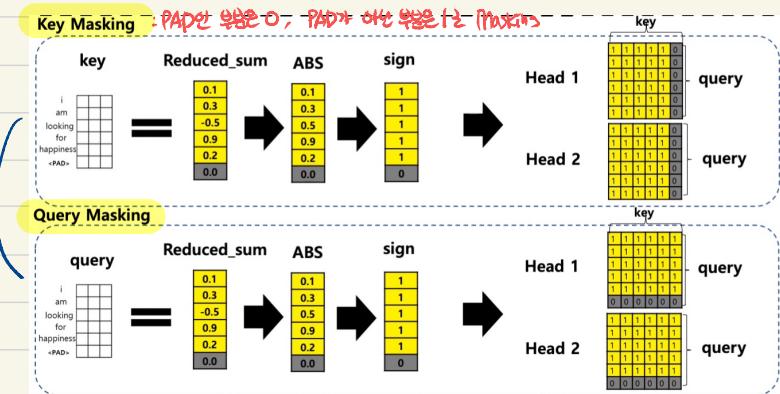
Head 1 2H \rightarrow 2Head Q, K, V



· Query or key^T 을 연산을 통해 Attention value / weight 계산
But, hidden size가 커지면 gradient가 vanishing
그리고 문제 발생 기울기값이 사라지는 문제
그러므로 $/ \sqrt{d}$ 를 적용
↳ scaling 하는 것

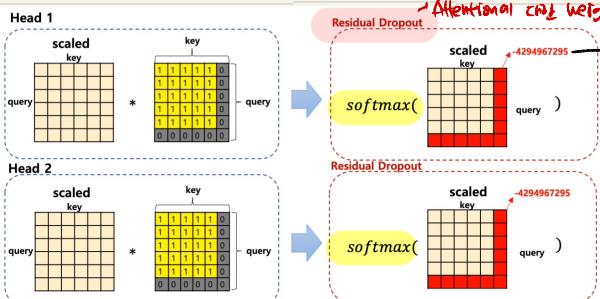
· 단어 Sequence의 <PAD> 가 있는 경우.

Key Masking : PAD인 부분은 0, PAD가 아닌 부분은 1을 Masking

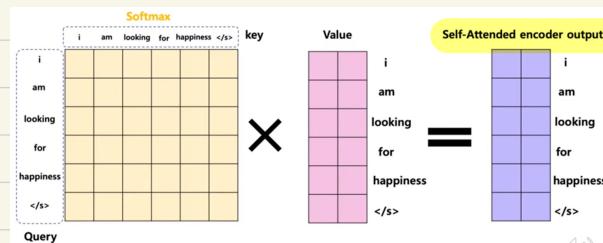
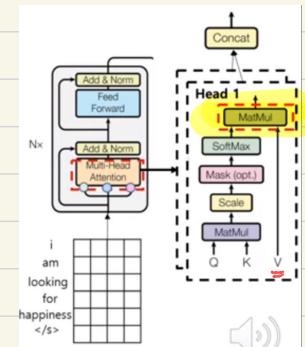
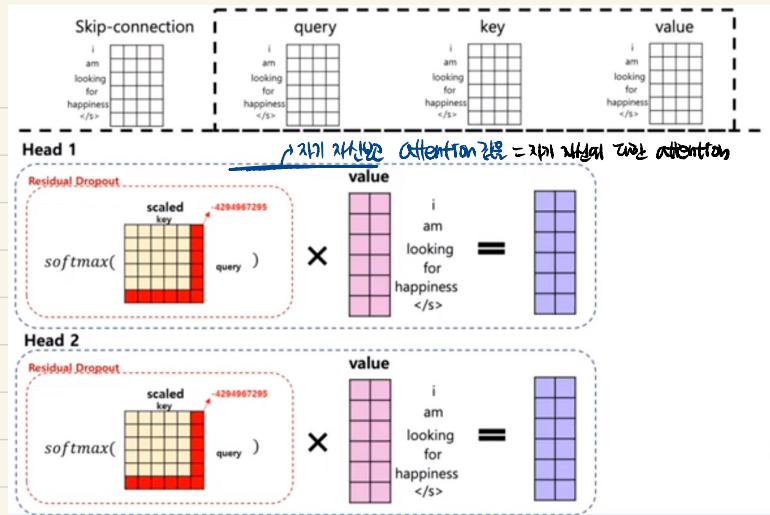


→ 이 둘은 유사하게 PAD(<PAD>)가 아닌 부분에만 1인 key, Query Masking을 만들수 0 \rightarrow 그 Scaled attention weight on 328

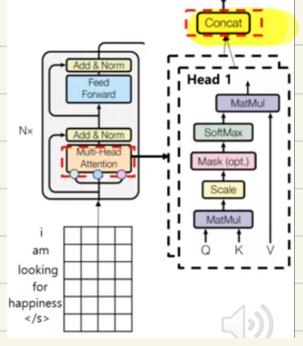
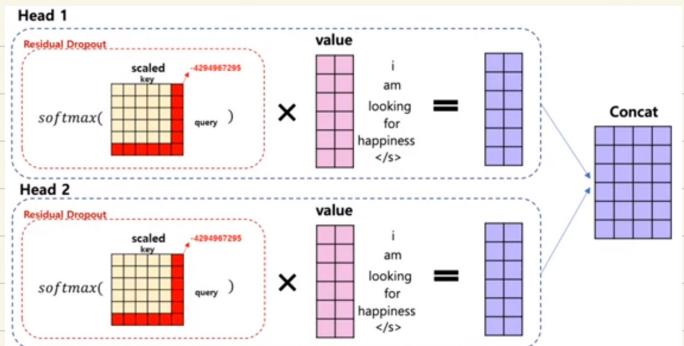
- Attentional out weight를 더해줌



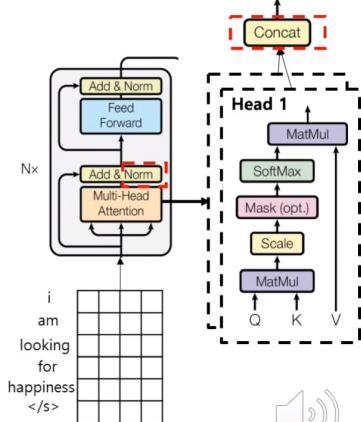
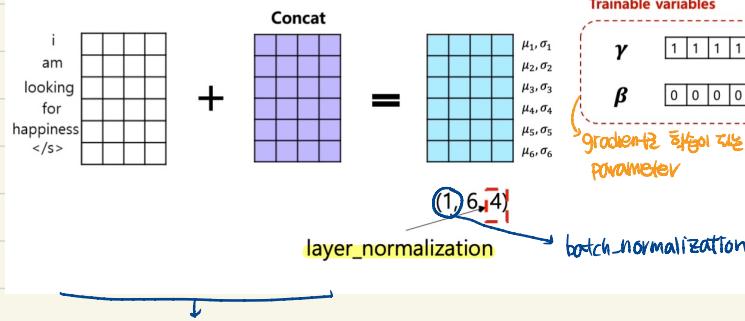
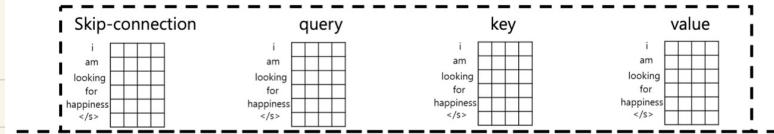
- 이 부분이 0이지만 강장화를 통한 값을 넣어줌
softmax를 계산해 \rightarrow 이 부분의 값이 0이라
마침내 나온다 하기 때문에



Self-Attention's output \oplus Concat!



Add & Norm



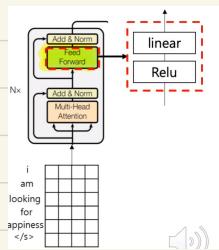
Concat한 결과를 원래의 input과 대비 : Add

Feed Forward

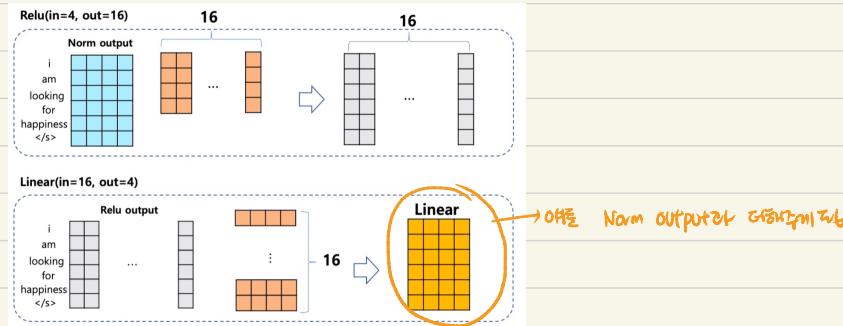
: Add & Norm 및 residual connection 합친 버전

여기서 ① Relu Function or ② linear

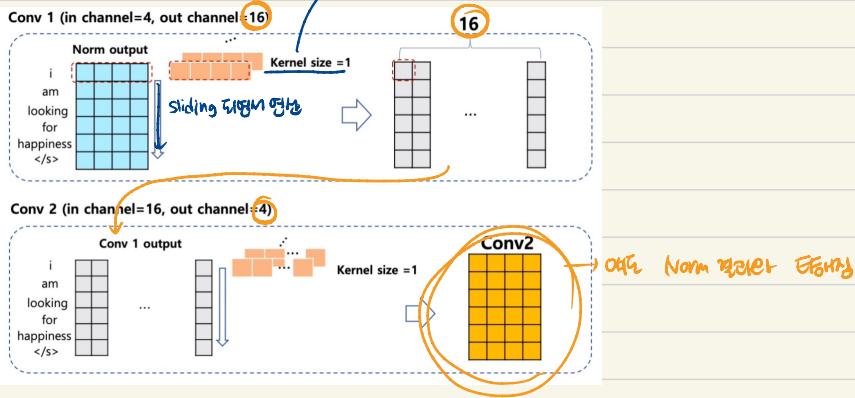
② convolution



① Relu Function → linear



② Convolution



• Aff & Norm

Trainable variables

$$\gamma = [1 \ 1 \ 1 \ 1]$$

$$\beta = [0 \ 0 \ 0 \ 0]$$

$$ln[0,:] = \gamma \frac{x_{[0,:]} - \mu_1}{\sigma_1} + \beta$$

layer_normalization

