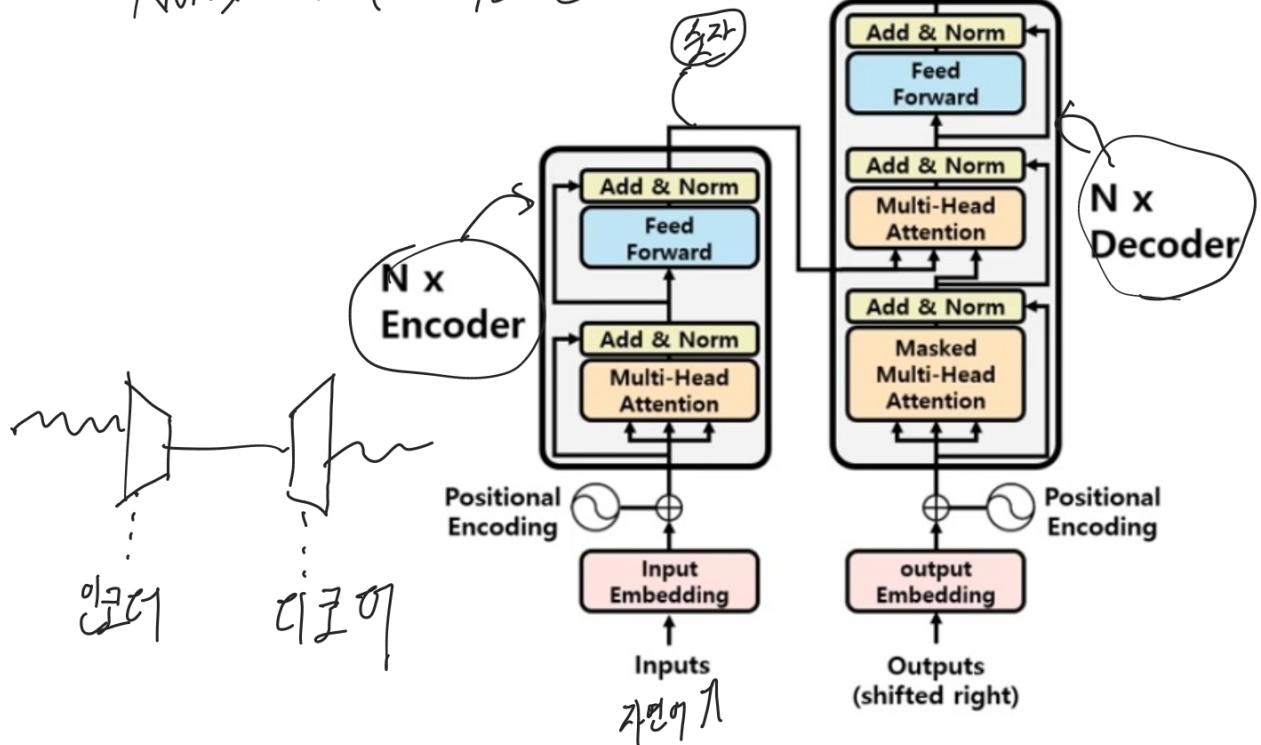


Transformer

↳ 초기에 번역기를 위해 만들어짐

↳ 병렬학습

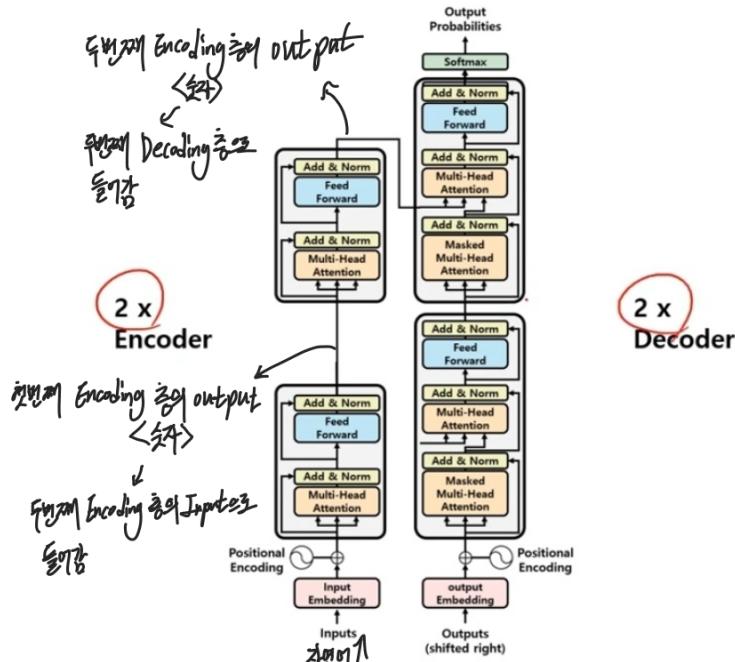
RNN, LSTM → 직렬학습



위에 있는 그림은 $N=1$ 일 때 transformer 모델 그림이고 아래의 그림은

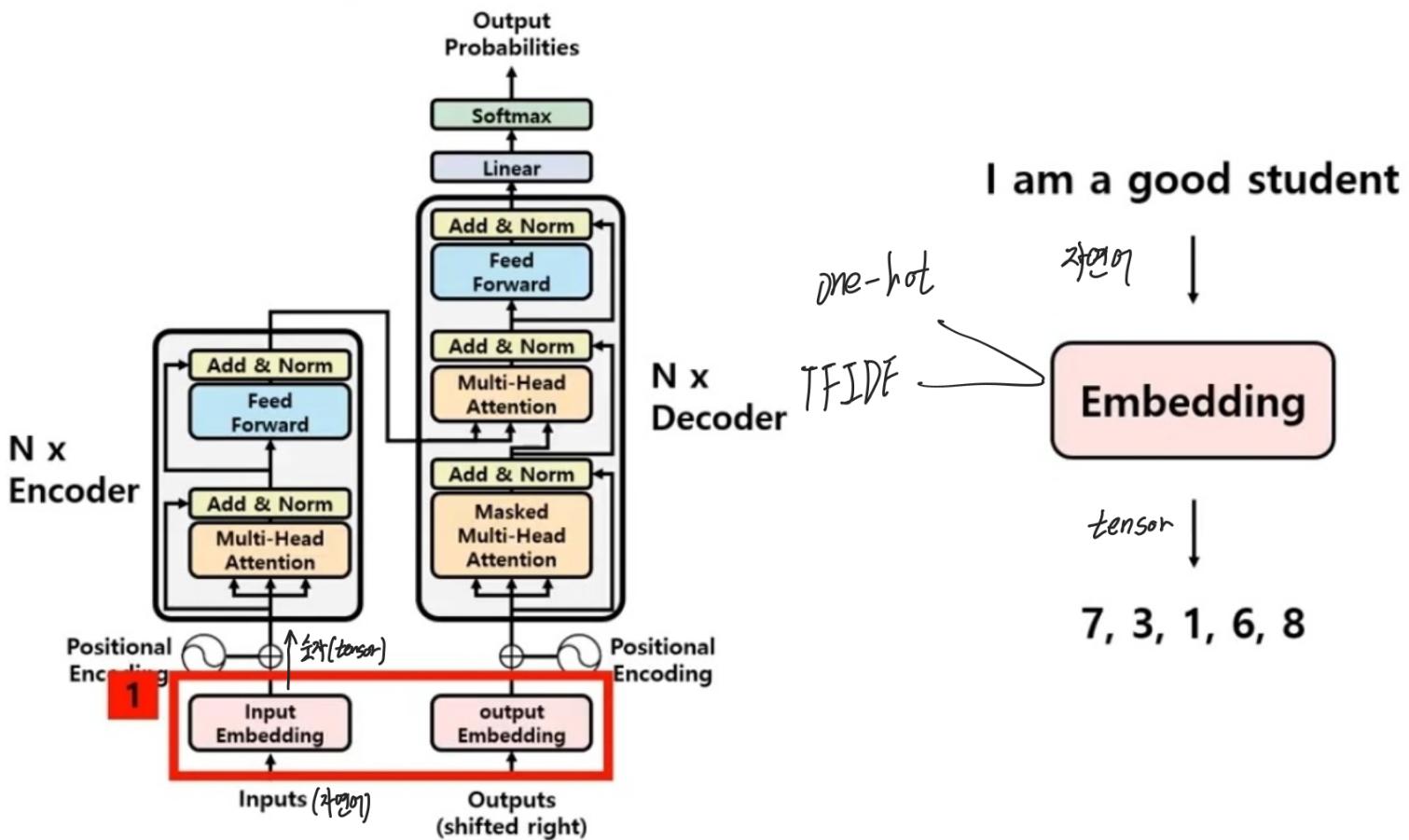
$N=2$ 일 때 transformer 모델의 그림이다.

Transformer

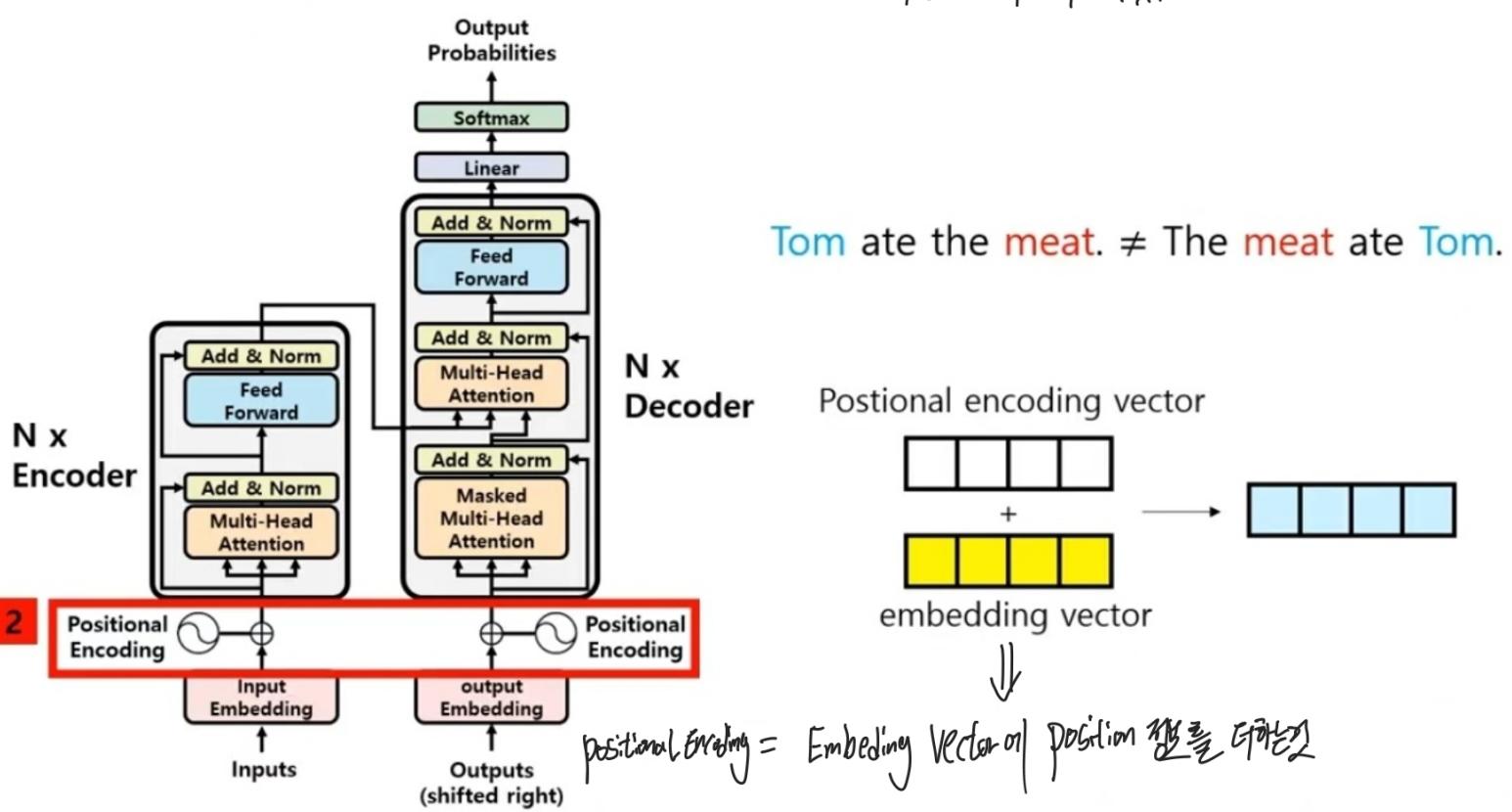


즉, 마지막 끝에 있는 Encoder, Decoder 까지만 연결된다.

Embedding

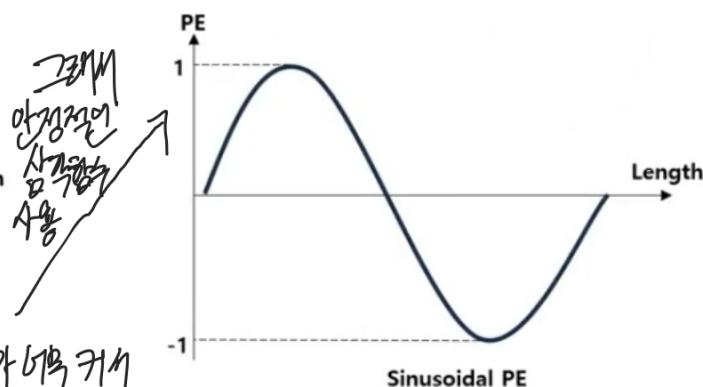
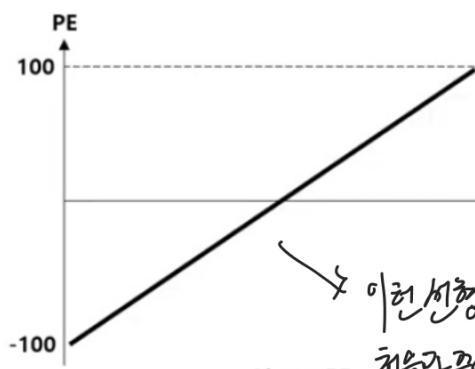


Positional Encoding



어떤 방식으로 position 정보를 줄 것인가?

Positional Encoding



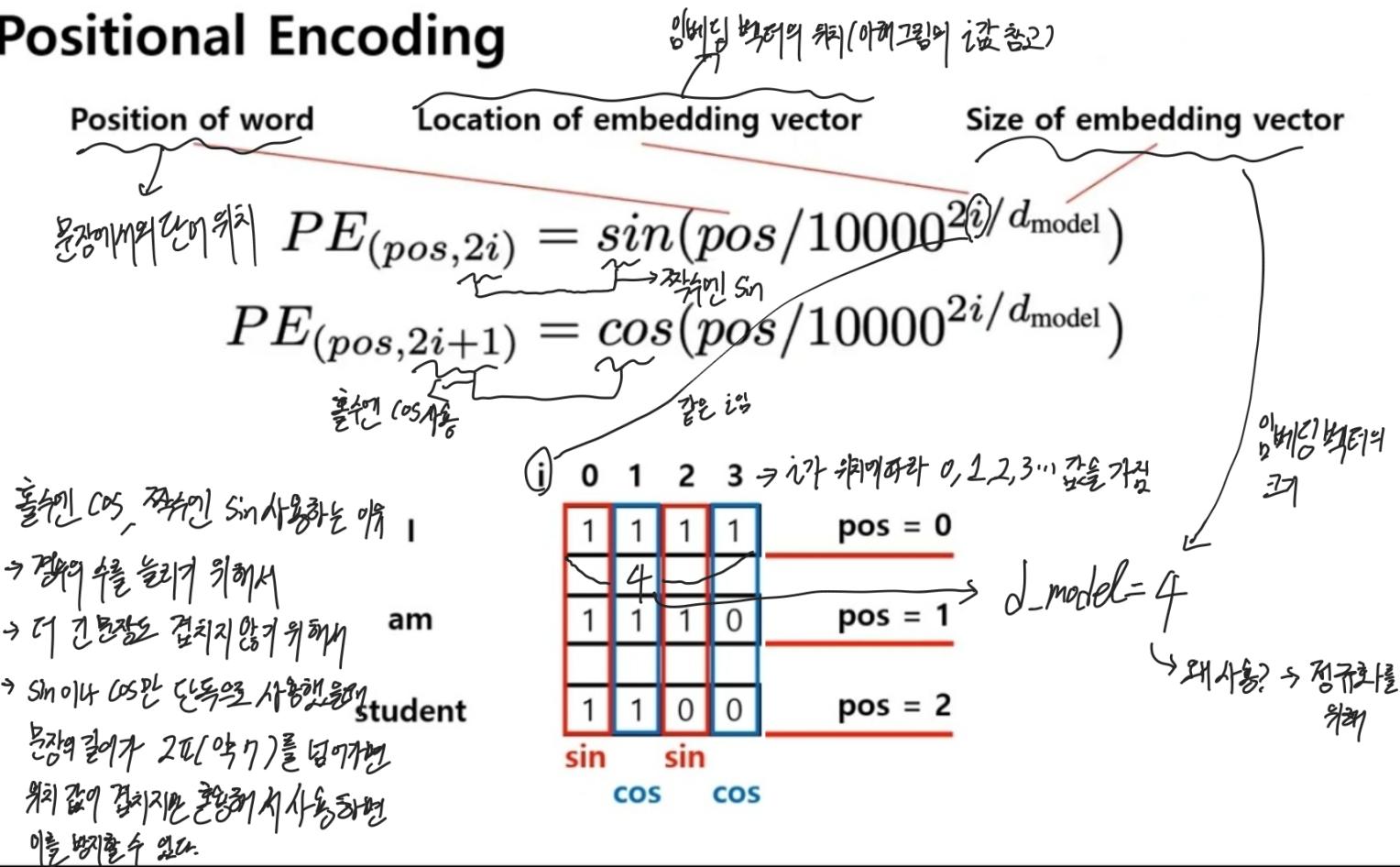
position을 안정적으로 사용할 수 있도록

$$PE_{(pos, 2i)} = \sin(pos/10000^{2i/d_{model}})$$

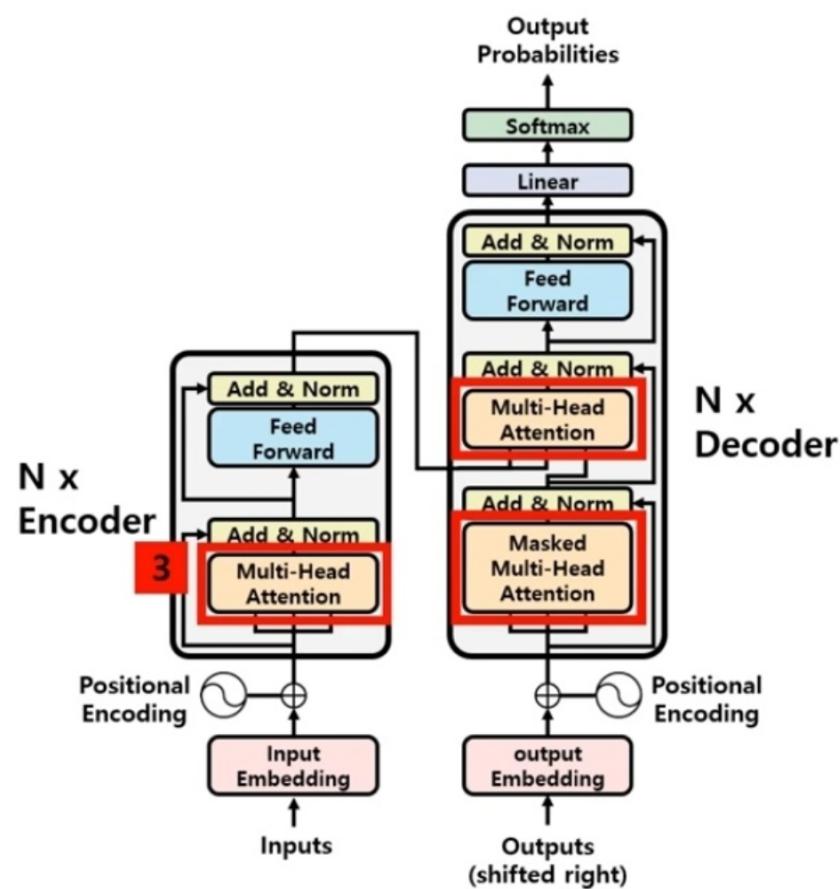
$$PE_{(pos, 2i+1)} = \cos(pos/10000^{2i/d_{model}})$$

- 범위가 $-1 \sim 1$ 로 안정적
- 주기함수로 글자수에 상관없음

Positional Encoding



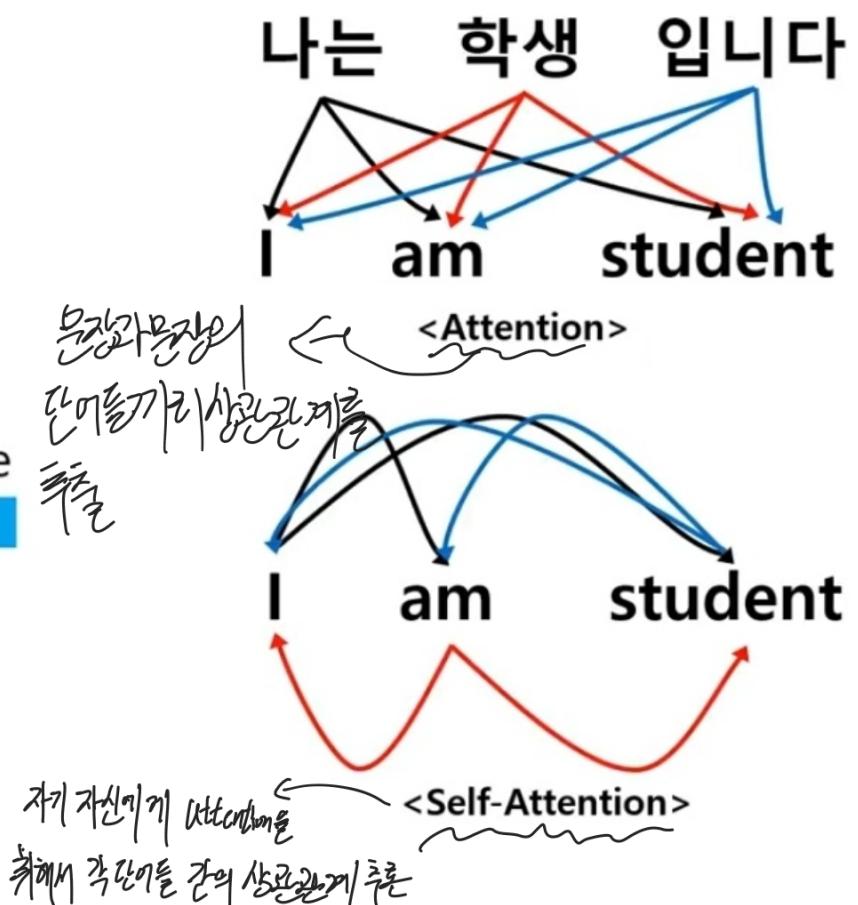
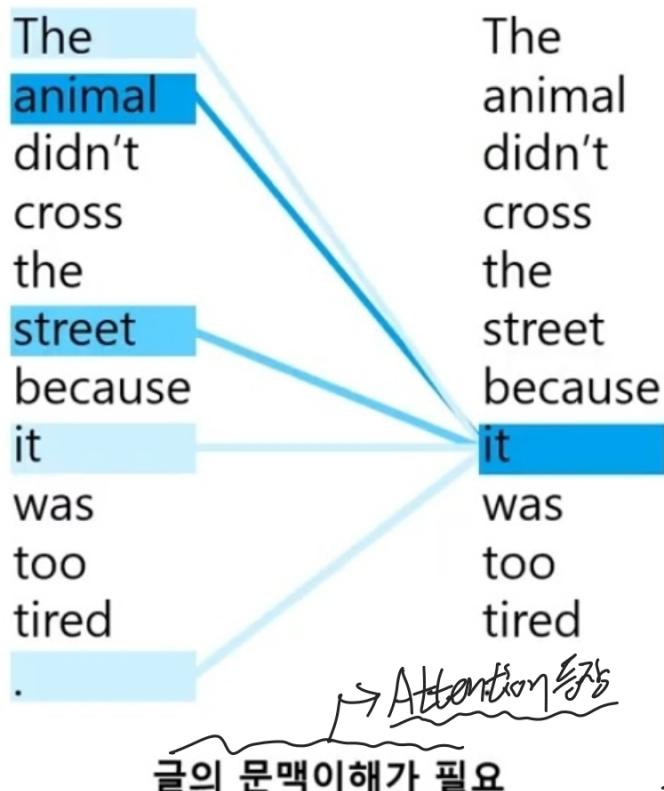
Multi-Head Attention



Attention

Multi-head Attention을 잘가능합니다.
Attention을 잘가능합니다.

Attention



Attention

고객아이디	고객이름	나이	등급	직업	적립금
apple	김현준	20	gold	학생	1000
banana	정소화	25	vip	간호사	2500
carrot	원유선	28	gold	교사	4500
orange	정지영	22	silver	학생	0



Query



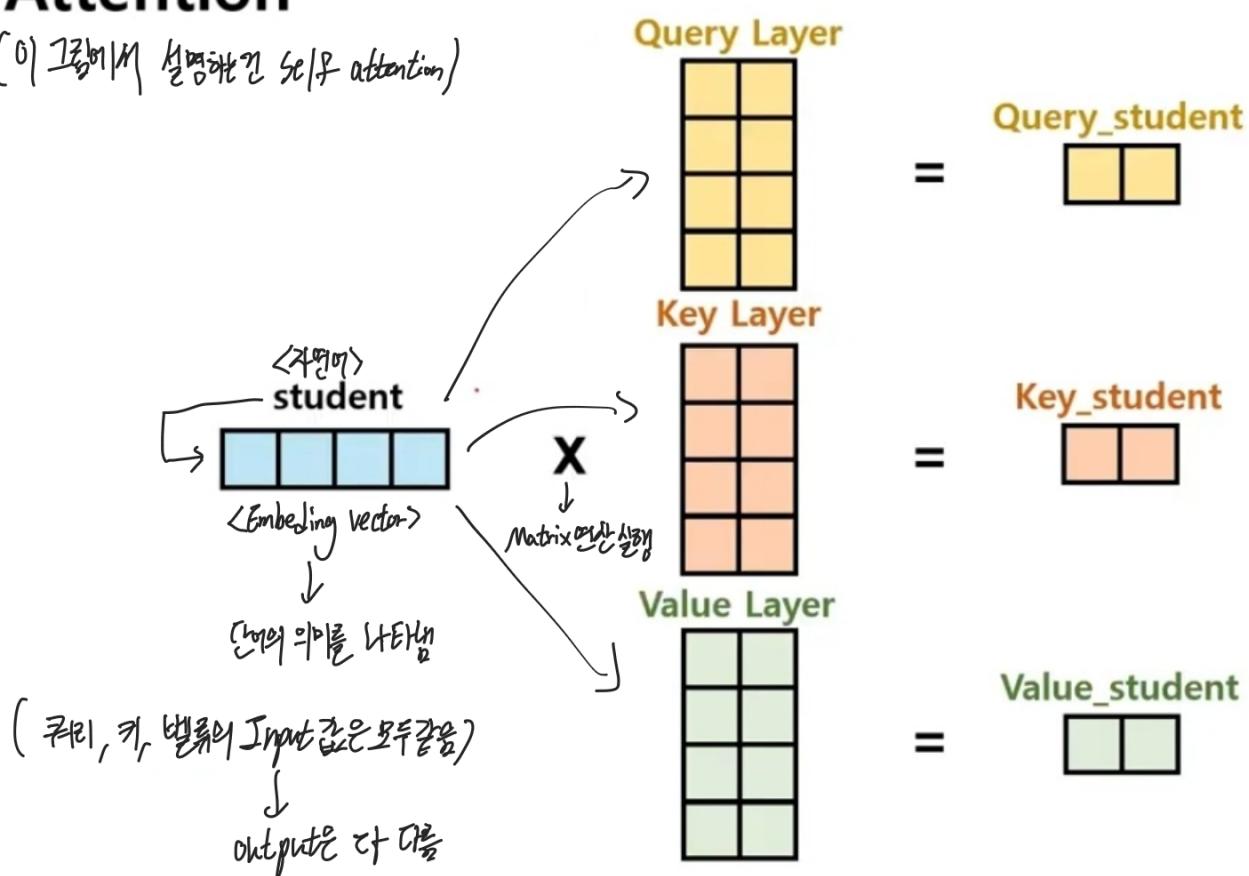
Key



Value

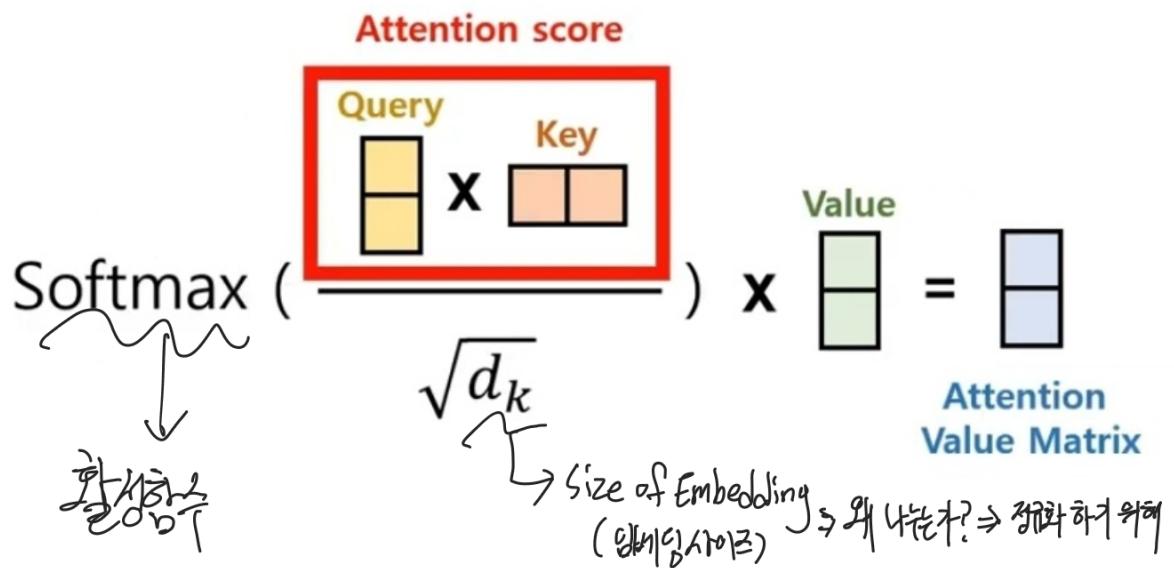
Attention

(0) 그림에서 설명해준 Self attention)



Attention

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



Attention

Self-Attention Map

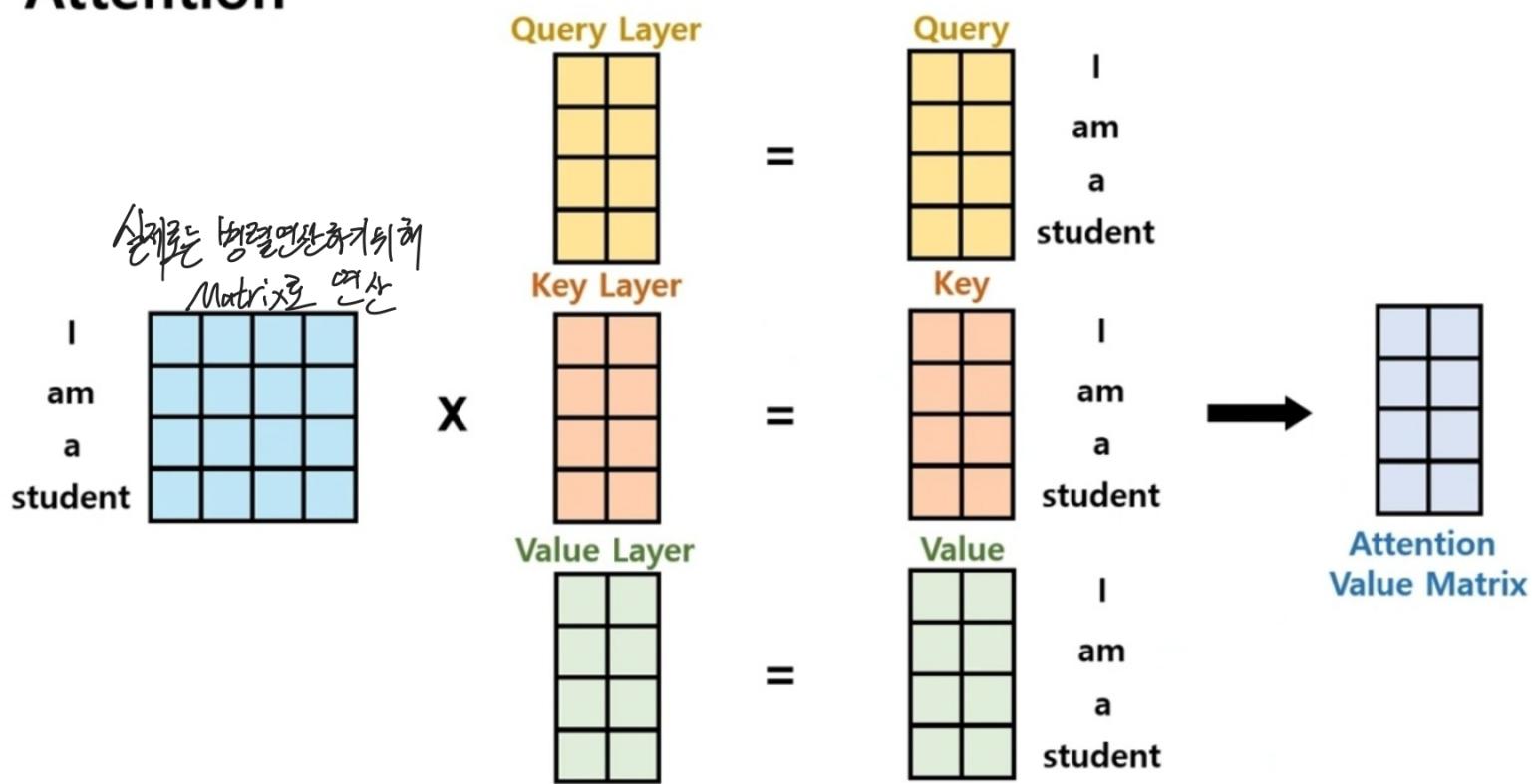
	\	am	a	student
I	1.0	0.5	0.3	1.0
am	0.5	1.0	0.4	0.2
a	0.3	0.4	1.0	0.1
student	1.0	0.2	0.1	1.0

I = Student라는 정보를 알려주는 유의미한 정보

[나에게] 그가 자신의 모든 무의미한 정보

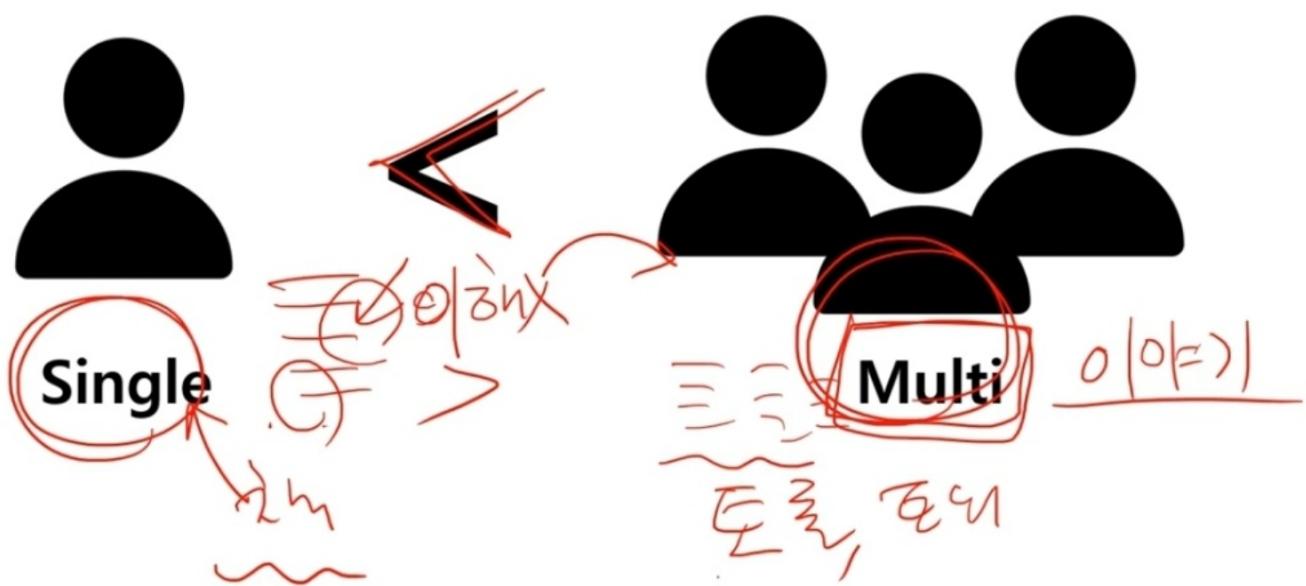
$$\frac{QK^T}{\sqrt{d_k}} = \frac{3}{\sqrt{4}} = 1.5$$

Attention

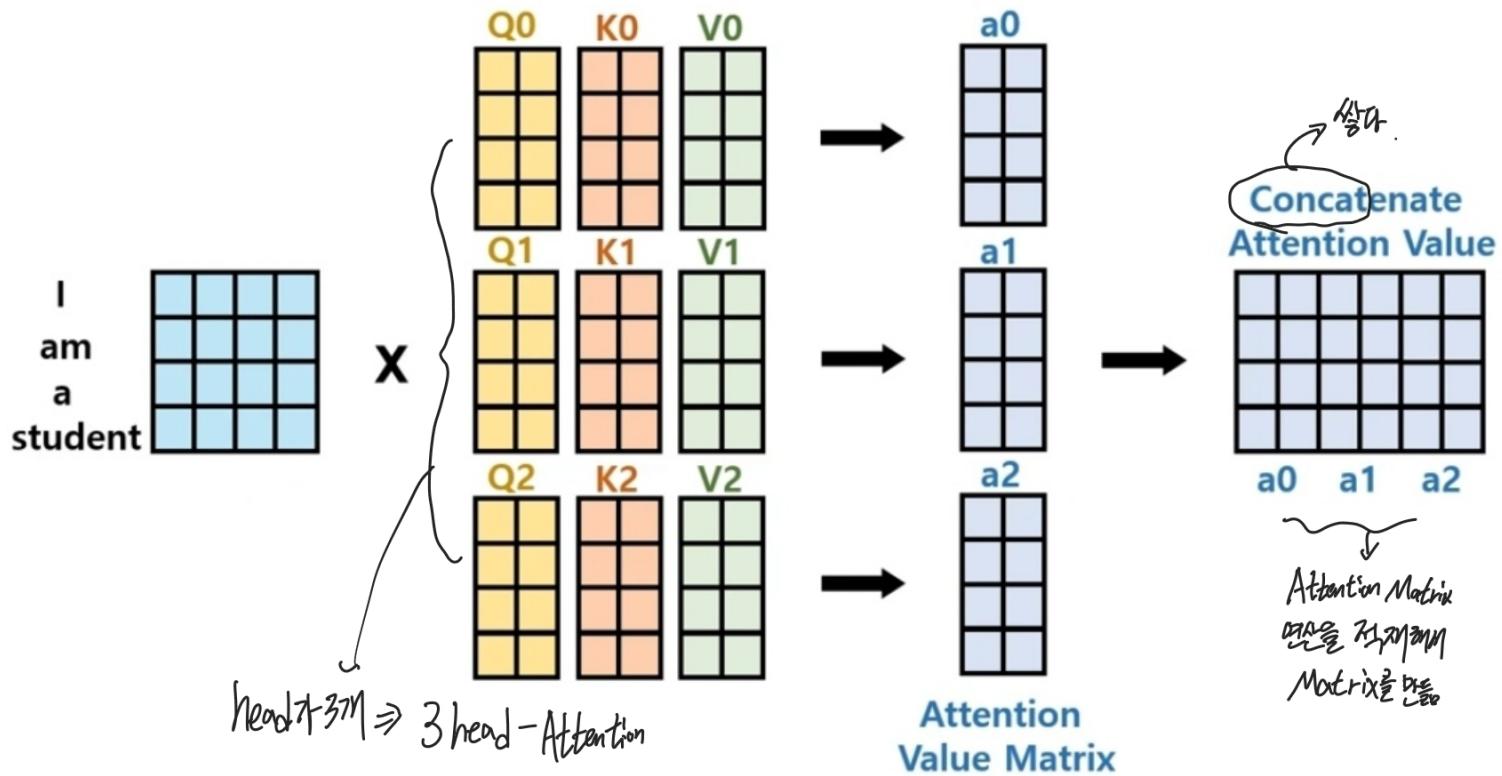


Multi-head Attention

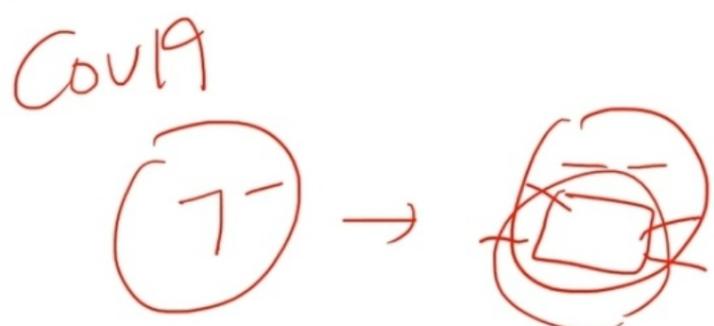
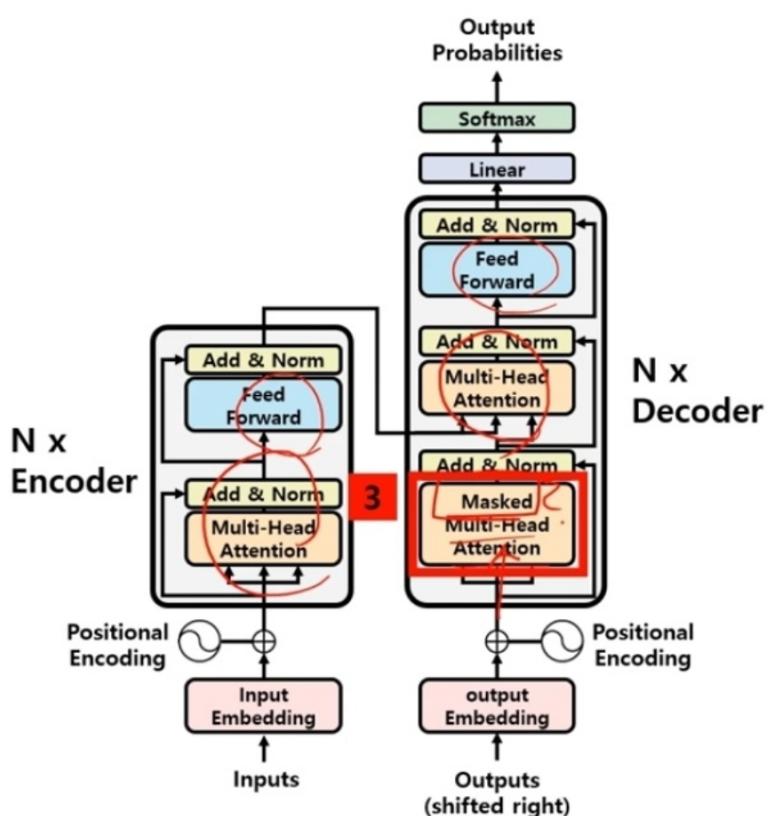
기능: 주제 찾기
영화, 뉴스↑.



Multi-head Attention



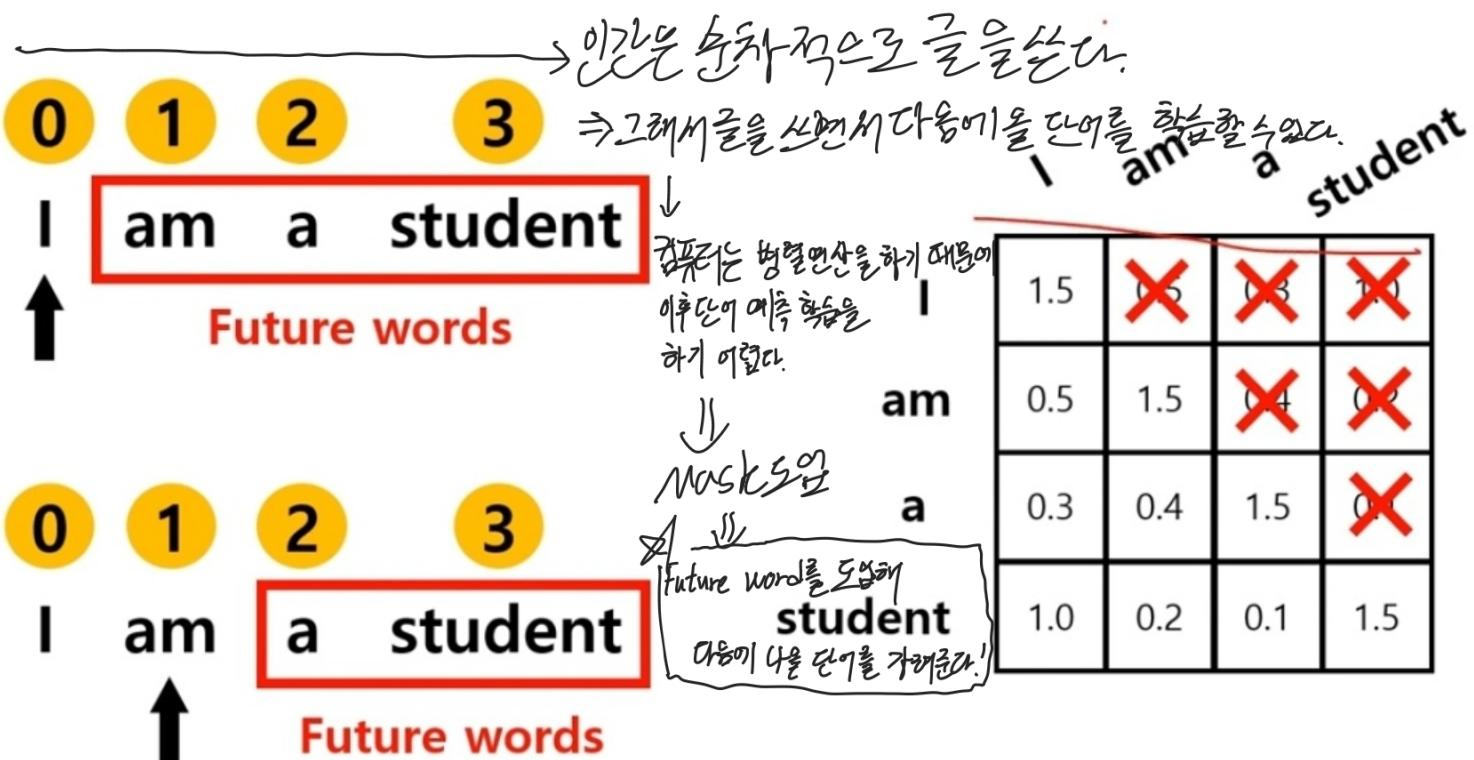
Masked Multi-Head Attention



Masked ???

9/01 → 9/2/02

Masked Multi-Head Attention



Masked Multi-Head Attention

↳ 어떤 방식으로 Future word를 가져줄까?

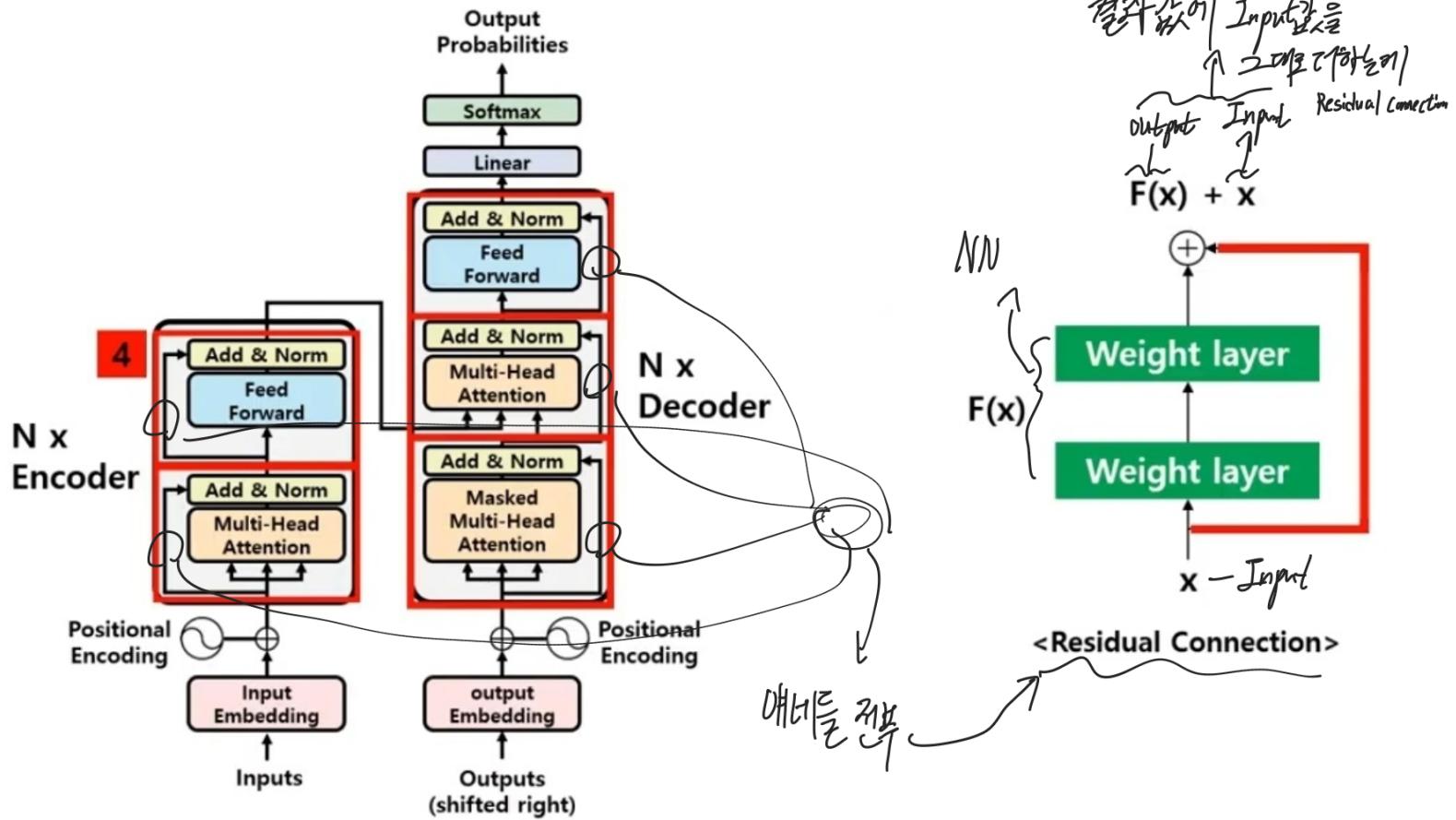
	\	am	a	student
I	10	-\infty	-\infty	-\infty
am	0	10	-\infty	-\infty
a	-5	-0.1	10	-\infty
student	10	-7	-8	10

Softmax

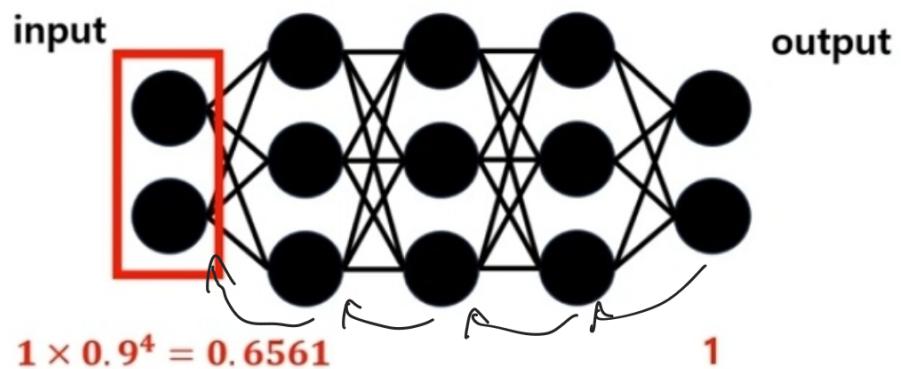
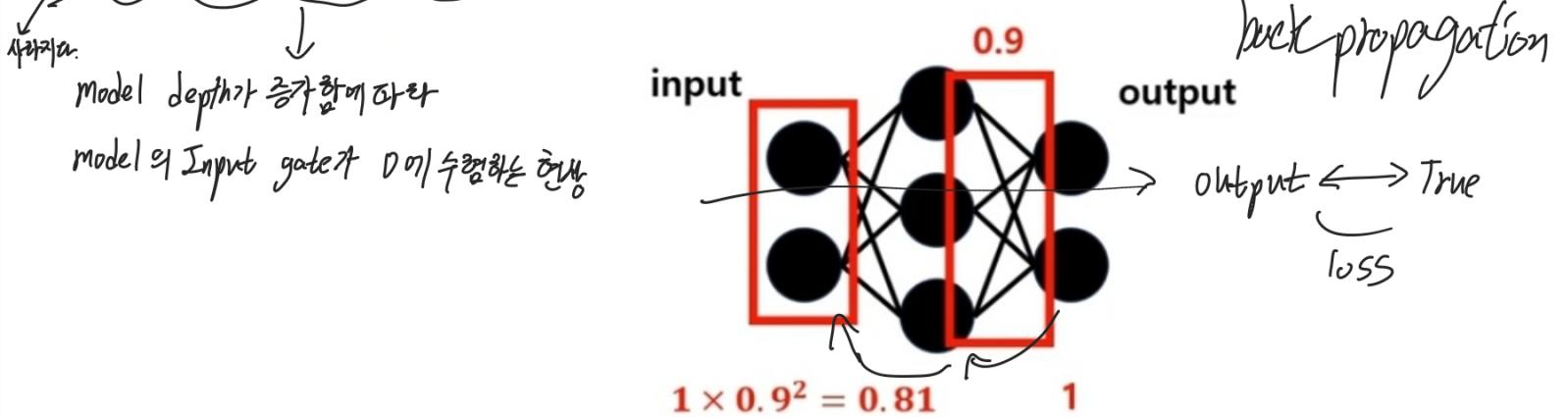
	\	am	a	student
I	1.0	0.0	0.0	0.0
am	0.5	1.0	0.0	0.0
a	0.3	0.4	1.0	0.0
student	1.0	0.2	0.1	1.0

Transformer는 (-100000)정도 넓는다.

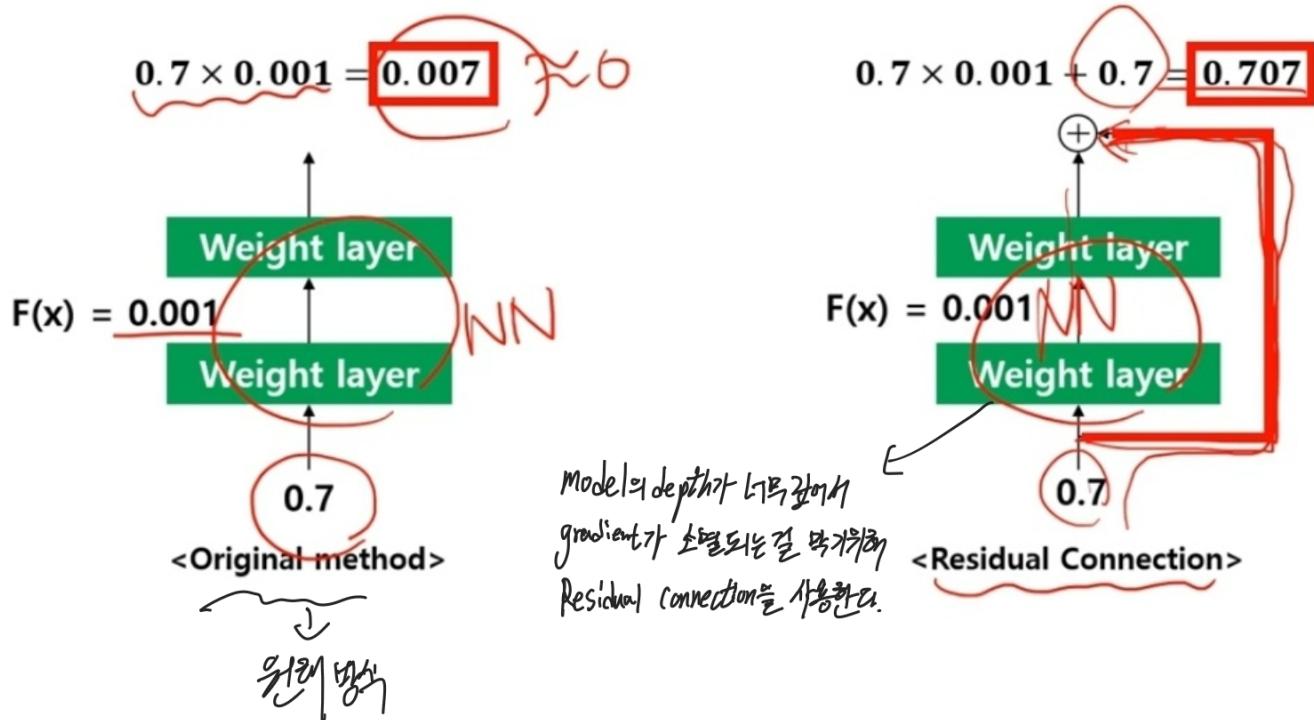
Residual connection



Vanishing gradient \Rightarrow Residual connection은 1이 유지됩니다.



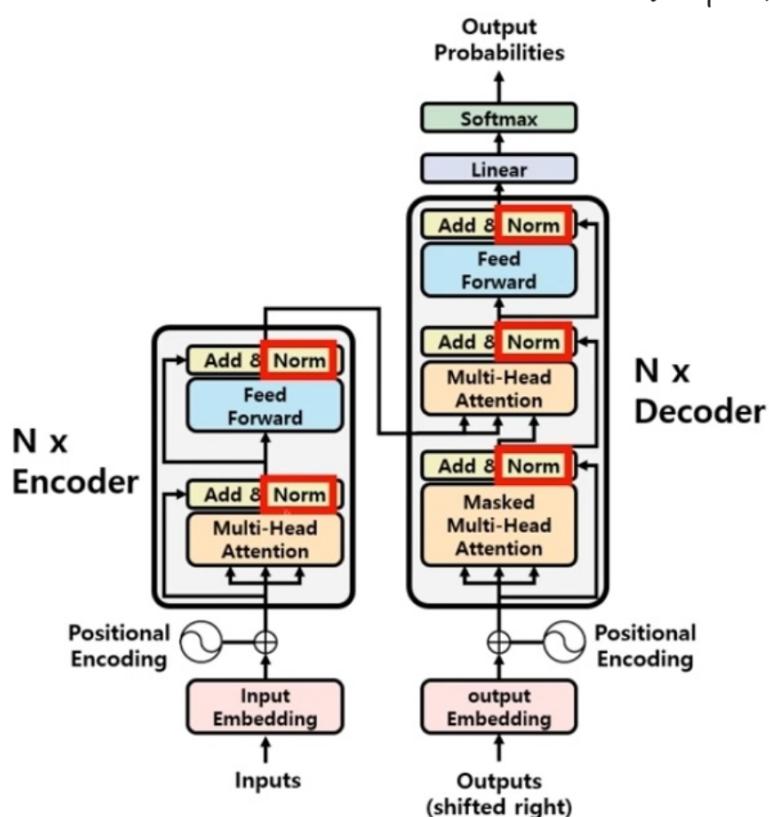
Residual connection



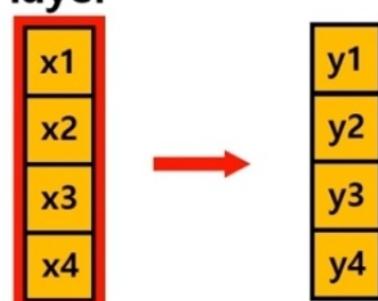
Layer normalization

AI에서 가장 중요한 것 \Rightarrow 정규화 \Rightarrow 왜? \Rightarrow Data들의 상대적인 차이를 인식해서 학습하기 위해

Layer normalization



layer



$$y = \frac{x - E[x]}{\sqrt{V[x] + \epsilon}} * \gamma + \beta$$

Note: Handwritten Korean text below the equation says: "단위 제작을 위해" and "**Layer normalization**".