

훈련의 목표 : 입력 문장 + 출력 문장을 받아서  
출력 문장의 각 위치에서 다음 토큰을 예측

입력 문장

: BOS + I + like + apples + EOS  
+ (PAD)

출력 문장

: BOS + 나는 + 사과를 + 좋아한다 + EOS

model에 input으로 들어감

정답으로 쓸 Label

# Encoder

하나의 batch별로 처리됨  $B$ : 배치 크기

초기 입력:  $(B, S)$   $S$ : 배치 내에서 문장의  
token 최대 길이  
index가 저장  $\hookrightarrow$  짧은 애들은  
PAD로 채움

1. 임베딩

$(B, S, d\text{-model})$

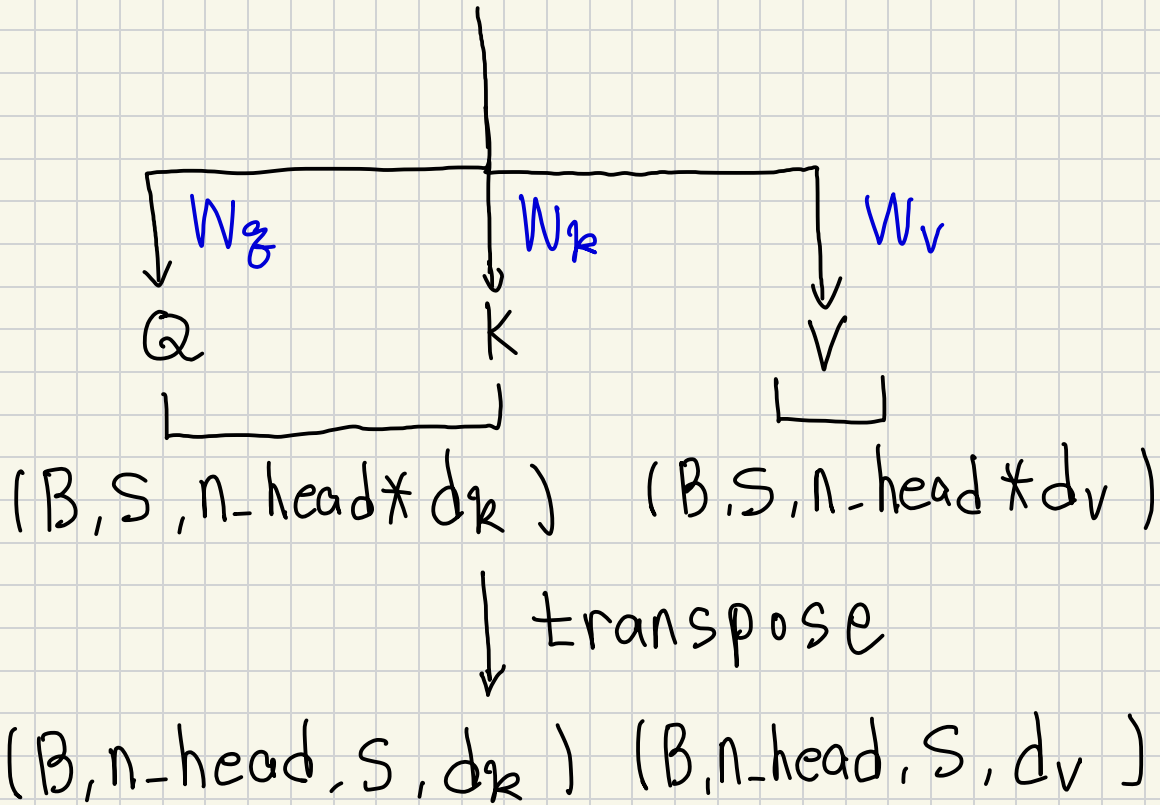
각 토큰이  $d\text{-model}$  길이의  
embedding vector로 변환

2. Positional Encoding

$(B, S, d\text{-model}) + (1, S, d\text{-model})$

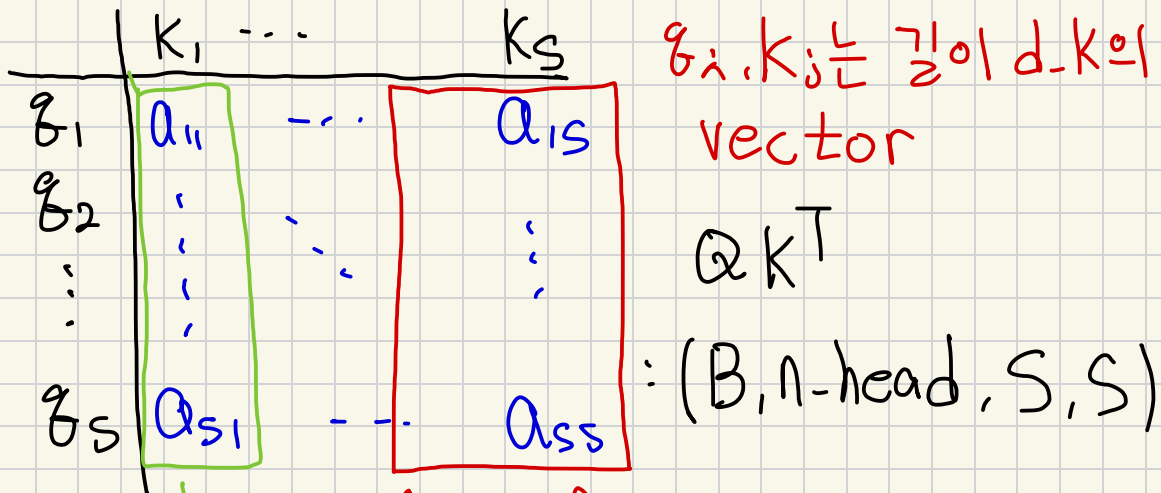
문장에서 token의 위치 pos에 따라  
각 token의 embedding vector에  
positional vector를 더함

3. Multi-head attention (Enc self-attn)  
start (B, S, d\_model)



$$\text{attn} = \text{softmax} \left( \frac{QK^T}{\sqrt{d_k}} \right) V$$

앞의  $(B, n\text{-head})$ 는 유지된다고 생각



이 방향은 Softmax

$\leftarrow -\text{inf}$   $\rightarrow (B, 1, 1, S)$

mask를 통해 PAD와의 attention은 영향 안 주도록

$$\text{softmax} \left( \frac{QK^T}{\sqrt{d_k}} \right) \cdot V$$

$$(B, n\text{-head}, S, S) \quad (B, n\text{-head}, S, d_v)$$

$$\Downarrow$$

$$(B, n\text{-head}, S, d_v)$$

$$z_1 = a_{11}V_1 + a_{12}V_2 + \dots \quad (\text{PAD 전까지만})$$

$$z_2 = a_{21}V_1 + a_{22}V_2 + \dots \quad ( \quad " \quad )$$

$(B, n\text{-head}, S, d_v)$

↓ transpose

$(B, S, n\text{-head} \times d_v)$

↓  $W$

$(B, S, d\text{-model})$  end

$\Rightarrow$  각 token의 새로운 embedding vector

마지막으로 이 결과를 이전의 embedding

과 더해주면 끝 (start + end)

↑ 이 과정이 Residual  
정규화는 생략

## 4. FFN (Feed-Forward Network)

$(B, S, d_{\text{model}})$

↓  $W_1$

$(B, S, d_{\text{hid}})$

↓ Relu

$(B, S, d_{\text{hid}})$

↓  $W_2$

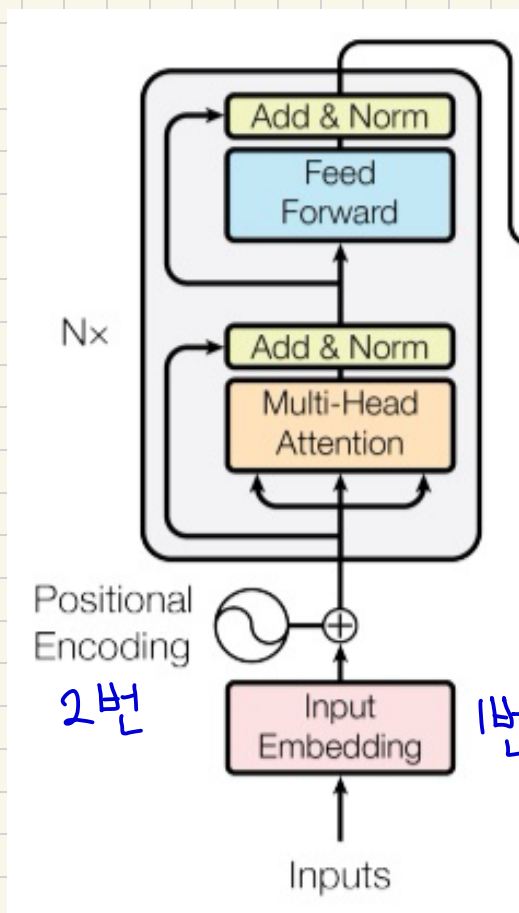
$(B, S, d_{\text{model}})$

FFN의 의미: attention에서 주지 못하는  
비선형 학습관계를 제공

여기에 정규화 + Residual

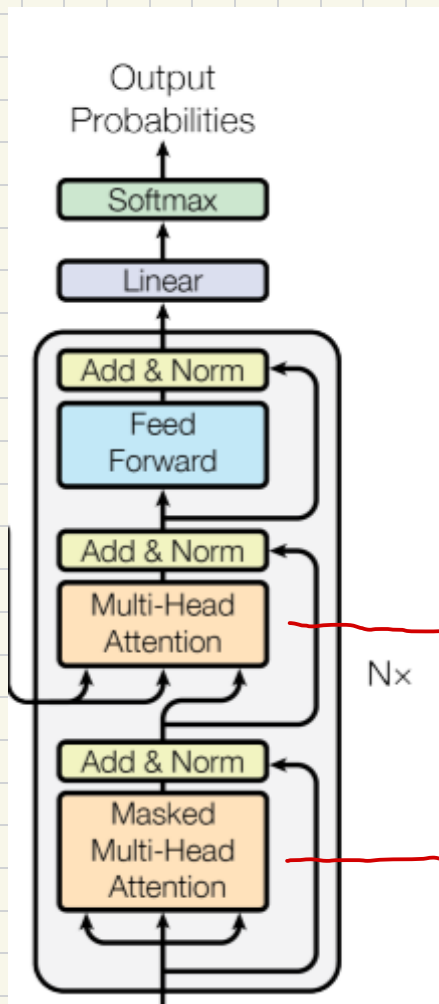


최종 output:  $(B, S, d\text{-model})$



# Decoder

Embedding과 Positional Encoding은  
동일하니까 제외



Enc-Dec attention

Dec self-attention

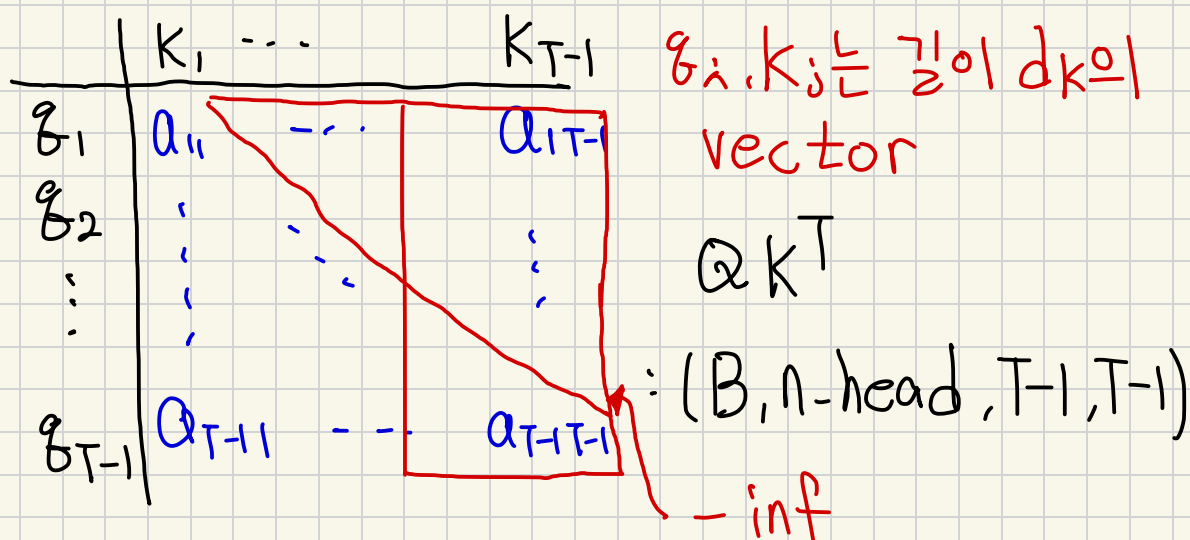
두 attention을 중심으로 보자

# 1. Dec-self attention

$$Q: (B, n\text{-head}, T-1, d_k)$$

$$K: (B, n\text{-head}, T-1, d_k)$$

$$V: (B, n\text{-head}, T-1, d_v)$$



PAD mask +  $q_i$  입장에서 미래 시점의 토큰인  $k_j$ 는 참조 불가능하도록 함 ( $j > i$ )

$$mask: (B, 1, T-1, T-1)$$

나머지는 Enc self-attention과 동일함

Output은 역시  $(B, T-1, d\text{-model})$

## 2. Enc-Dec attention

Dec embedding  
( $B, T-1, d_{\text{model}}$ )

$\downarrow W_q, \text{transpose}$   
 $Q$

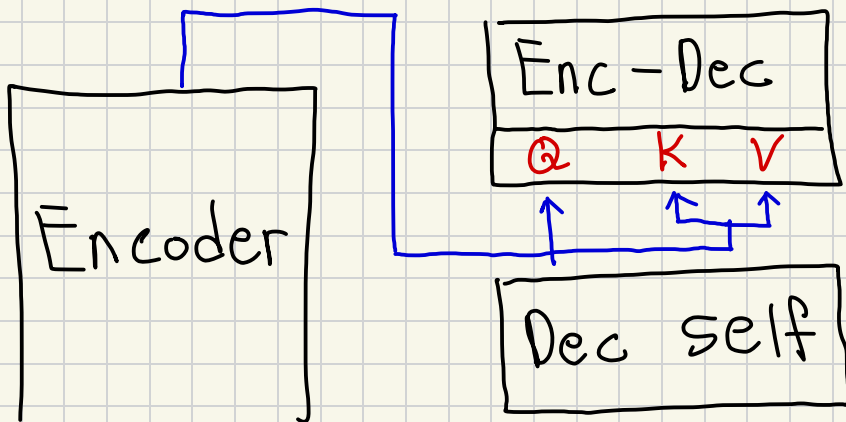
Enc output  
( $B, S, d_{\text{model}}$ )

$\downarrow W_k, \text{transpose} \rightarrow K$   
 $\downarrow W_v, \text{transpose} \rightarrow V$

$Q: (B, n\text{-head}, T-1, d_k)$

$K: (B, n\text{-head}, S, d_k)$

$V: (B, n\text{-head}, S, d_v)$



	$k_1$	...	$k_S$
$q_1$	$a_{11}$	...	$a_{1S}$
$q_2$	$\vdots$	$\ddots$	$\vdots$
$\vdots$	$\vdots$	$\ddots$	$\vdots$
$q_{T-1}$	$a_{T-1,1}$	...	$a_{T-1,S}$

$q_i, k_j$ 는 길이  $d_k$ 의 vector

$$QK^T$$

$:(B, n\text{-head}, T-1, S)$

$\uparrow -inf$

마찬가지로 PAD인 Key는 attention 계산에 반영하지 않음

(Query쪽을 안 막는 이유는 어짜피 PAD의 embedding은 쓰이지 않기 때문)

mask:  $(B, 1, T-1, S)$

역시나 Output은  $(B, T-1, d\text{-model})$

(Decoder쪽 token의 embedding vector를 업데이트하는 과정이므로)

# Residual을 쓰는 이유 (노피셜)

우리는 Enc-Dec attention에서

$g_1 = a_{11}v_1 + a_{22}v_2 + \dots$  처럼 decoder쪽 token의 embedding을 encoder output으로 만들어진  $v_i$ 들의 선형결합으로 바꾸고 있다

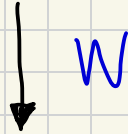
$\Rightarrow$  원래  $g_1$ 이 지니던 decoder token 간의 관계를 나타내던 embedding은 사라지는 느낌 (가중치에만 반영되므로)

$\Rightarrow$  Residual로 이를 살려준다

## 최종 Loss 계산

Decoder에서 최종 Output:  $(B, T-1, d\_model)$

$(B, T-1, d\_model)$



$(B, T-1, \text{Vocab\_size})$



Cross-Entropy로 label과의 Loss를 계산