

# 엔코더 디코더

정형화된 neural sequence transduction models는 encoder-decoder 구조.

encoder:  $(x_1, \dots, x_n)$   $\xrightarrow{\text{mapping}}$   $(z_1, \dots, z_n)$  auto-regressive: consuming previously generated symbols as additional input.

decoder:  $(z_1, \dots, z_n)$   $\xrightarrow{\text{output}}$   $(y_1, \dots, y_m)$

→ Transformer는 이걸 구현 by using (stacked-self attention + pointwise fully connected layers.) ← 인코더, 디코더 구조.

1) encoder & Decoder stacks

encoder: 6개의 layers stacked.  
 each has 2 sub layers.  $\left\{ \begin{array}{l} \text{multi-head self attention : 공유 히드는 <sup>연산</sup> 동일} \\ \text{pointwise fully connected.} \end{array} \right.$

1차원 input  
 → 6개의 필로링 레이어 층.  
 각 레이어 내에서 공유 히드 공유  
 각각 다른 방식으로  
 입력 정보 학습 수,  
 필터링은 하나의 벡터로 인코딩  
 (디코딩에서)

residual connection & layer normalization.

(LayerNorm(at Sublayer(x)))

$d_{model}$

각 레이어 포워드 인코딩 출력  $d_{model} = 512$  ; residual connection 동일하게 하면.

decoder: 6개의 layers

has 3 sub layers  $\left\{ \begin{array}{l} \text{multi-head self attention} \leftarrow \text{reversed : 이후 포지션이 학습되어 있으므로, 즉, 미리 데이터를 못보게} \\ \text{각 포지션이 (이전에서 생성한 단어들)에 mask} \\ \text{각 레이어의 ~~이전~~ 다른 multi head attention.} \end{array} \right.$

각 레이어 LayerNorm(x → Sublayer(x))

- ① 멀티헤드 어텐션수행  
 → 여러 서브레이를 통해 정보 병합
- ② 8개를 병합  
 → concatenate + 가중치 부여
- ③ FFN 통과시켜서 최종 출력 벡터 생성
- ④ softmax 적용해서 단어 출력 변환  
 → 가장 높은 출력 단어 선택

## 2) Attention

$Q = (K, V) \Rightarrow \text{output}$   
mapping

output: Value의 가중합.

이때 가정치;  $Q$ - $K$ 간의 유사도 측정을 반영할 수 있는 방법  
compatibility function

### • Scaled Dot product Attention

$Q, K \in \text{dim} = d_k$

$V \in \text{dim} = d_v$

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

$Q$ - $K$  중 ~~가장~~ 높은 점수를 선택하여 곱하기

$d_k$  작으면 additional attention 정보가 더보다 줄 (  $d_k$  스케일링  $\propto$  )

$d_k$  커지면 ~~배경 크기 커져서~~ Softmax의 가중치를 너무 작게 만든다.  $\downarrow \text{softmax} = y(1-y)$

input과 관련된 배정량도 커진다.  $\hookrightarrow$  backprop 이터의 배정량도  $\downarrow$  or 학습 불능  
gradient vanishing

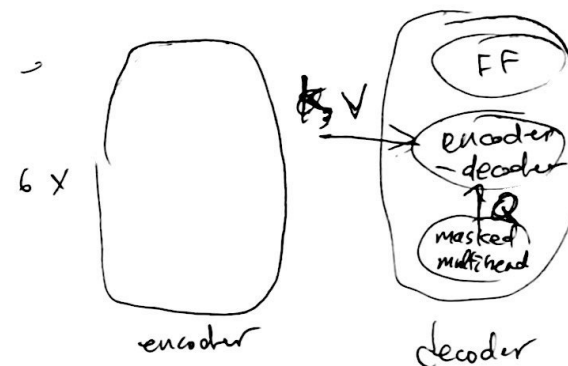
$\Rightarrow \frac{1}{\sqrt{d_k}}$ 로 스케일링하여 해결!

### • Multi-head Attention

$\text{Multi-head}(Q, K, V) = \text{Concat}(\text{head}_1, \dots, \text{head}_h)W^O$

where  $\text{head}_i = \text{Attention}(QW_i^Q, KW_i^K, VW_i^V)$

$$d_k, d_v = \frac{d_{\text{model}}}{h} = 64$$



### • Fully Connected $d_{\text{FF}} = 2048$ (inner layer)

$$\text{FFN}(x) = \max(0, xW_1 + b_1)W_2 + b_2$$

ReLU, input, output dim 역시 2048.

positional encoding  $\leftarrow$  임베딩에 추가될  
즉,  $d_{\text{model}} = 512$

각 서브레이어에  $\wedge$  이후 dropout: residual dropout  
(임베딩과 서브레이어 간의 연결을 통해 학습을 돕는다)  
기본 인코더는 0.1 비율

· 배치치: 여러 가능한 쿼리 시퀀스 중 선택(비트 유망 후보 선택)

· 빔크기: 한번에 고려할 수 있는 시퀀스 길이 개수 설정

· 길이 제한: 생성 시퀀스 길이에 따라 제한/차단

가장 긴 시퀀스 선택되지 않게 한다.

### Result

( $k, V$ )에 따라  $h$ 의 variation 주며 model variance 감소.

·  $d_k \downarrow \rightarrow$  모델 성능  $\downarrow$  : dot product보다 지능적인 코어생성 필요하진도?  
V에 대한 가중치 계산이 필요

· 모델 클러스터 성능  $\uparrow$

· dropout  $\leftarrow$  과적합 방지(효과적)

· sinusoidal positional Encoding  $\approx$  learned positional Encoding: base model과 잘 맞지 않음 X

### Conclusion

· recurrent 제거한 attention 기반 Transformer.

· Sota model 보다 효과적 (성능, 비용, 시간면에서)

· 제약: 픽셀은 이미 인코딩된 Transformer 활용.

Further Tasks  $\rightarrow$  large input (이미지/비디오) 처리하려면 local, restricted ~~한~~ 어텐션 메커니즘 필요.