

연필 아카이브

정규화 neural sequence transduction model, encoder-decoder 73.

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

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

→ Transformers는  $\frac{1}{\sqrt{d}}$  이동을 by using (stacked-self attention + pointwise fully connected layers.) ← 인코더, 디코더 구조.

encoder: 6 layers stacked.  
 each has 2 sub layers.  
 if sub layers { multi-head self attention : 804  $\times$  11  $\leq$  2  $\times$  802  $\times$  802  
 position wise fully connected.

Input

→ 624의 **원근법** 세이어 층사.

각 세로로 나에서 8744 8123  
77 다른 882

은행 정보 학습

팔티 켄데는 하나의 벡터로  $\vec{v}$  (디리레기서)

residual connector  $\oplus$ , layer normalization.

(LayerNorm (dt Sublayer(x)))

$d_{\text{model}}$

$\frac{2 \times 10^6}{2} = 10^6$  이고  $\frac{2 \times 10^6}{2} = 10^6$  output dim = 5/2 ∴ residual connection 용도(보통 1개).

decoder : 64 layers

layers has 3 sub layers { multi-head self attention  $\leftarrow$  reversed : 이쪽 프리젠테이션이 중요하지 않음, 즉, 미래 데이터를 못본채 각 프리젠테이션이 (이전까지 생성한 단어들은 전부) masking.  
 2nd sub layer 이쪽 프리젠테이션은 multi head attention. ① 멀티헤드는 어휘선 수임

7/ 21/017 LayerNorm ( $x \rightarrow \text{SublayerNorm}(x)$ )

- ① 멀티헤드 어텐션 수렴  
→ 여러 서로다른 입력 정보들
- ② 8개를 하나로 합친다  
→ concatenate + 가중치 부여
- ③ FFN 통과시켜서 최종 출력 벡터 생성
- ④ Softmax 적용해서 단어 확률로 변환  
→ 가장 높은 확률 단어 선택

## 2) Attention

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

output: Values의 차등

이때 가정치:  $Q$ - $K$ 간의 유사도 측정을 반영할 수 있는 것  
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 것 선택  
→ 약화시켜 표현

$d_k$  작으면 additional attention 생기기보다  $d_k$  스케일링  $\propto$   $\frac{1}{\sqrt{d_k}}$  입력에도

$d_k$  커지면 배정 크기 커져서 Softmax의 기울기를 너무 작게 만든다.  $\frac{d \text{ Softmax}}{d x} = y(1-y)$

input과 관련된 배정량도 커졌다.  $\hookrightarrow$  backprop에서의 배정량도  $\downarrow$  or 학습 불능  
gradient vanishing

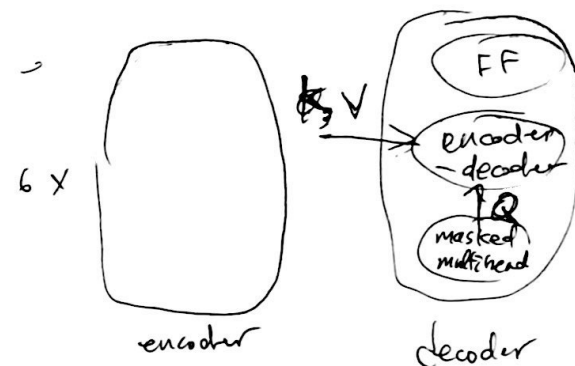
$\Rightarrow \frac{1}{\sqrt{d_k}}$  스케일링시켜야 배정!

### • 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 $d_{FF} = 2048$ (inner layer)

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

ReLU, input dim 역시 2048.

positional encoding  $\leftarrow$  임베딩에 더해짐  
 $\frac{2\pi}{d_m}, d_m = 512$

각 서브레이어에  $\frac{1}{4}$  dropout: residual dropout  
(임베딩과 서브레이어 각각  $\frac{1}{4}$ 씩)  
기분 안정화시키는 0.1 비율

· 배치지: 여러 가능한 후보들 중에서 최상위 유망 후보 선정

· 빔크기: 한번에 고려할 수 있는 시퀀스 길이 제한

· 길이 제한: 생성 시퀀스 길이에 따라 제한을 받음

가장 긴 시퀀스 선택되어 다음으로.

### Result

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

·  $d_k \downarrow \rightarrow$  모델 성능  $\downarrow$  : dot product보다 지능적인 코어레이션 필요함도?  
V에 대한 가중치 계산이 필요함

· 모델 클러스터 성능  $\uparrow$

· dropout  $\leftarrow$  과적합 방지 (효과적)

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

### 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) 처리를 위한 local, restricted ~~한~~ 어텐션 메커니즘 필요.