

# Attention, transformer.

## 1. 배경

- 기계번역 작업

\* 개요

• Input: Sequence of words (번역할 문장들)

• Model : 번역모델

• Output : Sequence of words (번역된 문장들)

## \* word Embedding

• 텍스트는 이미지처럼 수치적인 값으로 표현되는 데이터가 아님.

• 단어를 수치적인 자료(벡터)로 표시하기 위한 작업, word embedding이 필요한.

• Word embedding의 방법은 목적이나 모델에 기반하여 매우 다양함

→ 똑같은 단어라도 다른 모델에 따라 임베딩된 벡터가 다름.

• 대용량의 데이터로 잘 학습된 모델이 word embedding 결과를 가져다 사용

word → word embedding

→ dimensionality reduction

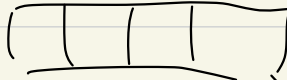
→ visualization of word embeddings in 2D

cat →  $\begin{bmatrix} 0.3 & 0.1 & -0.1 & 0.3 \end{bmatrix}$

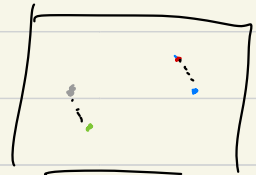
kitten

dog

man  
king  
queen

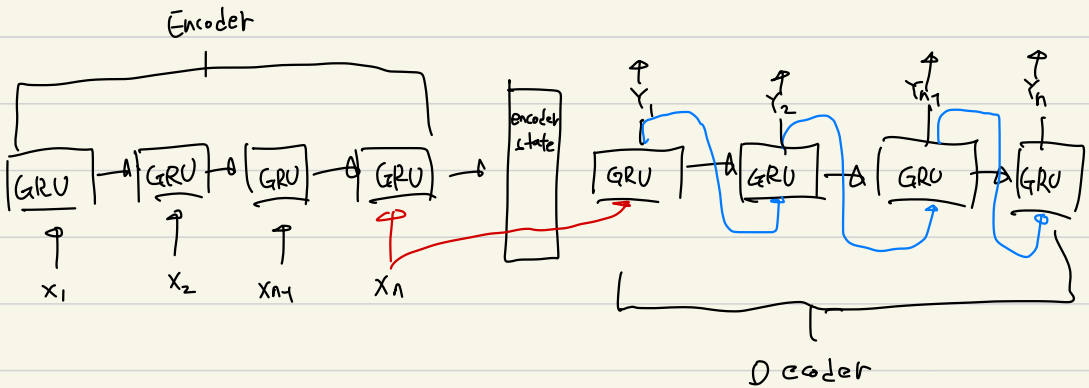


4D → 2D



## \* Seq2Seq

- Seq2Seq (sequence-to-sequence)는 LSTM (또는 GRU) 으로 구성된 인코더-디코더 구성의 신경망 모형
- Input: Sequence of text ( $R^{T \times d_w}$  행렬)
- Output: Sequence of text ( $R^{T' \times d_w}$  행렬)
- 인코더에서 번역할 문장을 차례대로 입력받아 디코더로 넘길 파쳐 (feature,  $h \in R^{k \times d_h}$ )를 학습함
- 디코더에선 인코더에서 넘긴 파쳐 (h)를 입력받아 순차적으로 번역 단어를 생성함



- Seq2Seq 모델의 구조

## 2. Attention

### \* Seq2Seq

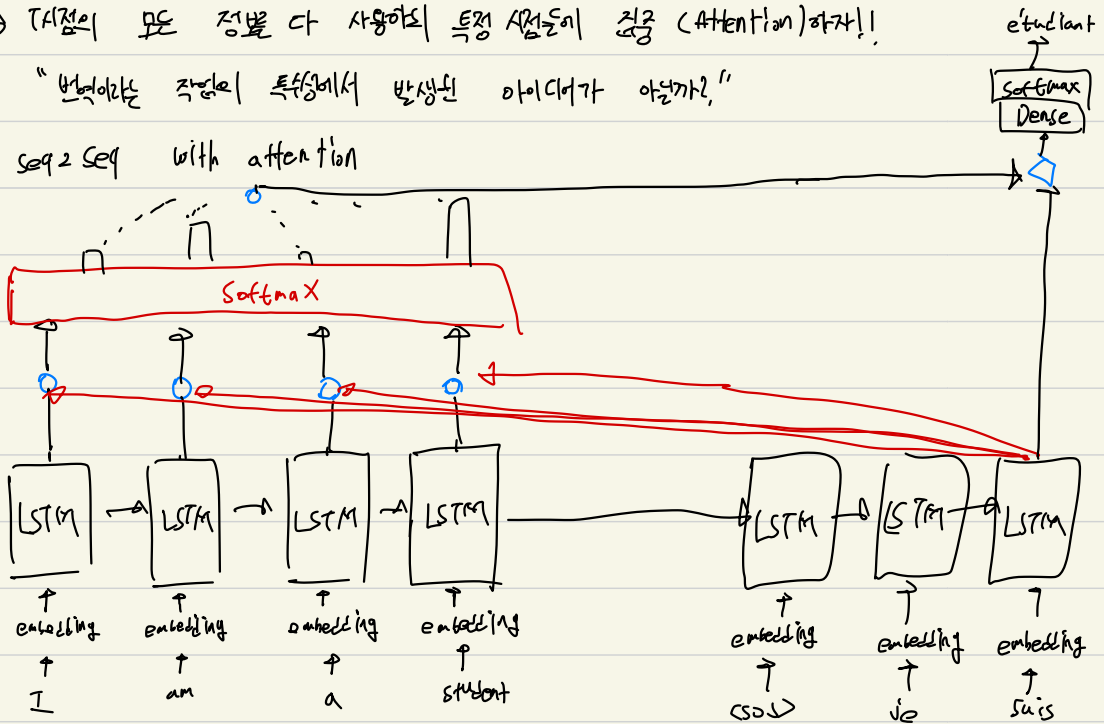
- T개의 단어로 이루어진 input 시퀀스 ( $\in \mathbb{R}^{T \times d_w}$ )
- 하나의 벡터  $h \in \mathbb{R}^{d_h}$ 인 인코더의 결과물
- (T+1)의 압축이 정보손실을 일으켜서 정확도가 떨어짐

### \* 해결책

→ T시점의 모든 정보 다 사용하되 특정 시점들이 중요 (Attention)하자!!

"번역이란 작업에 특정에서 발생한 아이디어가 아닐까?"

### - Seq2Seq with attention



Seq2Seq의 attention 방법 예시

# ① Attention score

: non-normalized attention weight  $\xrightarrow{\text{softmax}}$  attention weight

다음의  $t$ 시점을 기준으로 인코더 시점과 인코더 attention score 계산

- $q_t$ : 디코더  $t$ 시점의 LSTM cell의 hidden state 벡터
  - $k_i, i \in \{1, \dots, T\}$ : 인코더  $i$ 시점의 LSTM cell의 hidden state 벡터
  - $q_t^T k_i$ : 디코더의  $t$ 시점과 인코더의  $i$ 시점의 attention score (스칼라)
  - $S_t = [q_t^T k_1 \dots q_t^T k_T]^T$ : 디코더의  $t$ 시점과 인코더 모든 시점의 attention score 벡터.
  - $w_t = \text{softmax}(S_t) \in \mathbb{R}^T$ : 디코더의  $t$ 시점과 인코더 모든 시점의 attention weight 벡터
- 집중반영

•  $w_t^T \mathbf{1} = 1$

→  $q_t^T k_i$ : 디코더의  $t$ 시점과 인코더의  $i$ 시점의 attention score

- 두 벡터간의 내적을 바탕으로 두 벡터의 크기 관계를 가짐

- 디코더의  $t$ 시점의  $q_t$ 와 인코더의  $i$ 시점의  $k_i$ 가 가중치를 가짐

## ② Attention value

: 인코더의 hidden states와 attention weight의 가중 평균

디코더의 태깅을 기준으로 인코더 시퀀스와의 attention value 계산

•  $K = [k_1, \dots, k_T]^T \in \mathbb{R}^{T \times d_m}$ : 인코더 hidden state들 쌓은 행렬

•  $C_v^T = w_v^T K \in \mathbb{R}^{1 \times d_m}$ : attention value 벡터

→ Attention value는 디코더의 태깅과의 연관성을 기반으로 입력정보를 사용하여 새롭게 생성된 피쳐

→  $C_v$ 와  $K$ 를 곱한 벡터를 가중 곱의 L2 정규화 동원한 과정 수행

Attention은 input과 output 시퀀스 간의 연관성을 고려한 피쳐 학습 방법

(feature learning method)

### 3. Transformer

#### - Transformer 요약

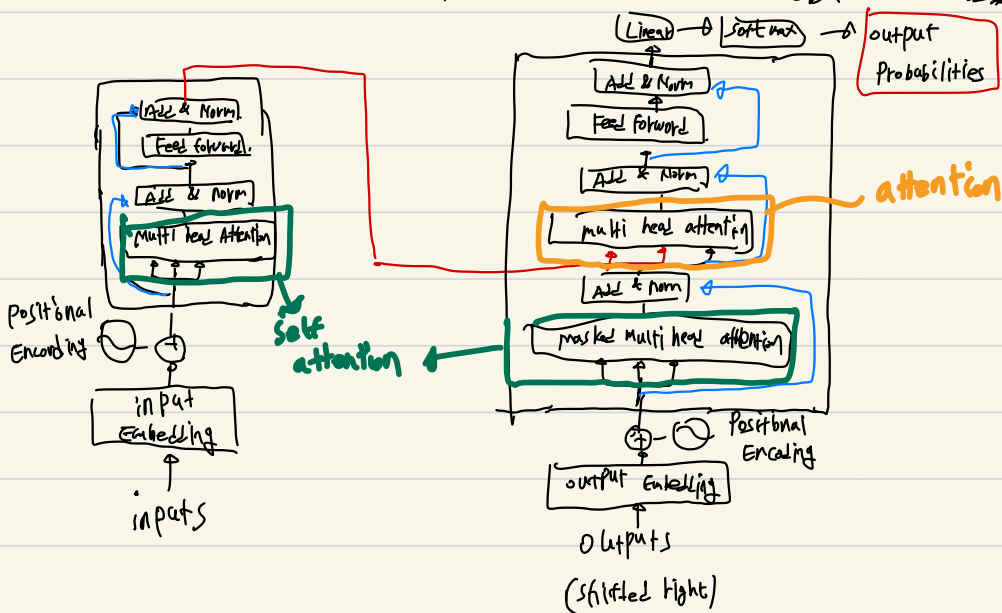
a new simple network architecture, the Transformer, based solely on attention mechanisms, dispensing with recurrence and convolution entirely.

Superior in quality, being more parallelizable.

→ 많은 양의 학습을 위한 많은 GPU 자원들을 해결함과 과부하 방지하기 attention based 함.  
효율성과 더 많은 병렬화가 가능해짐.

#### - 배경

- LSTM과 Convolution 기반의 모델에서 순차적인 계산량은 너무 높고 비효율적임.
- LSTM과 Convolution이 좋은 버리고 병렬화가 가능한 연산구조를 가진 모델 구성
- Self-attention은 병렬화가 가능한 이이들은 자연어처리 분야에서 성공적으로 사용되어 왔음.



# ① positional encoding

Transformer는 기존의 구분이 순차적인 연산을 비평화면서 구조를 과감하게 없앴.

→ positional encoding 이라는 대안을 사용함

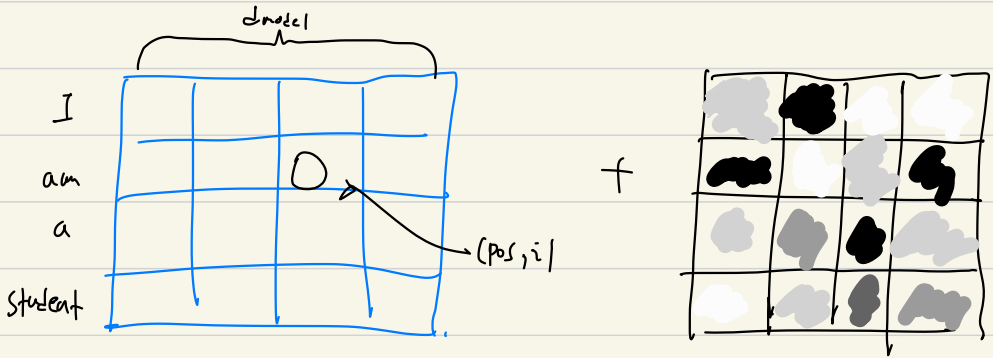
$$PE(pos, 2i) = \sin\left(\frac{pos}{10000^{2i/d_m}}\right)$$

$$PE(pos, 2i+1) = \cos\left(\frac{pos}{10000^{2i/d_m}}\right)$$

여기서 pos는 단어의 지점,  $2i$ 는 단어 임베딩 벡터에서의 인덱스

positional encoding 값으로 만든 행렬을 임베딩된 행렬에 더함.

→ 같은 단어라도 위치에 따라 다른 임베딩 벡터를 가짐.



positional encoding과 임베딩 행렬의 합.

## ② Self-attention

서로 다른 input과 output 간의 연관성이 아닌 동일한 문장 안에서 단어들 간의 연관성을 통해 서로를 파악 생성

- Self-attention을 하기 위해서는 하나의 입력 행렬( $X$ )에 각각 다른 FF를 거쳐 query( $Q$ ), key( $K$ ), value( $V$ )를 생성
- $Q = XW_Q$
- $K = XW_K$
- $V = XW_V$
- $Q$ 와  $K$ 를 가지고 attention scores를 계산하고 attention weights를 계산
- $V$ 의 attention weight만큼의 가중 평균으로 새롭게 파악 생성



③ Seq2Seq의 attention과 Transformer의 Self-attention의 차이점.

1. Seq2Seq에서는 디코더 타겟에 attention 연산을 위해서는 타 시점까지의 attention 연산이 선행되어야 함.

→ 디코더의 모든 시점에서 동시에 진행될 수 없어서 타겟에 hidden state (벡터)에서만 계산함.

2. Transformer는 순차적인 연산이 없어서 모든 시점에 대해 동시에 attention을 진행될 수 있음.

→ 타겟에 벡터만이 아닌 모든 시점  $(1, \dots, T)$ 의 행렬을 가지고 계산

3. Query

· Attention : 디코더의 타겟에서 hidden state (벡터)

· Self-Attention : 모든 시점의 모든 행렬  $X$ 와  $w_q$ 의 곱 (행렬)

4. Key

· Attention : 인코더의 모든 시점의 hidden state를 모든 행렬 (행렬)

· Self-Attention : 모든 시점의 모든 행렬  $X$ 와  $w_k$ 의 곱 (행렬)

5. Value

· Attention : 인코더의 모든 시점의 hidden state를 모든 행렬 (행렬)

· Self-Attention : 모든 시점의 모든 행렬  $X$ 와  $w_v$ 의 곱 (행렬)

## ④ Self-attention 연산 과정

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

attention weight

여기서  $d_k$ 은 행렬  $Q$ 의 연산값이 너무 커지는 것을 방지하기 위한 조정 값

- $QK^T / \sqrt{d_k} \in \mathbb{R}^{T \times T} = [s_1, \dots, s_T]^T$  : 모든 시점 간의 attention score를 담은 행렬
- $W = \text{Softmax}(QK^T / \sqrt{d_k}) = [w_1, \dots, w_T]^T \in \mathbb{R}^{T \times T}$
- $\text{Attention}(Q, K, V)$ 의 첫 번째 행:

$$C_1 = \sum_{t=1}^T w_{1,t} V_t \quad \sum_{t=1}^T w_{1,t} = 1$$

Self-attention의 의미

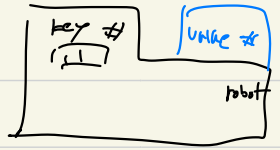
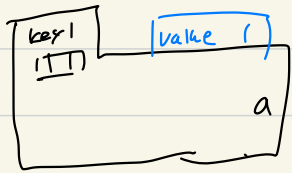
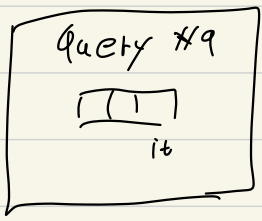
: 내배 단어들 간의 연관성(문맥)을 고려한 피쳐 생성

ex) Tom likes a dog. He also likes a cat.

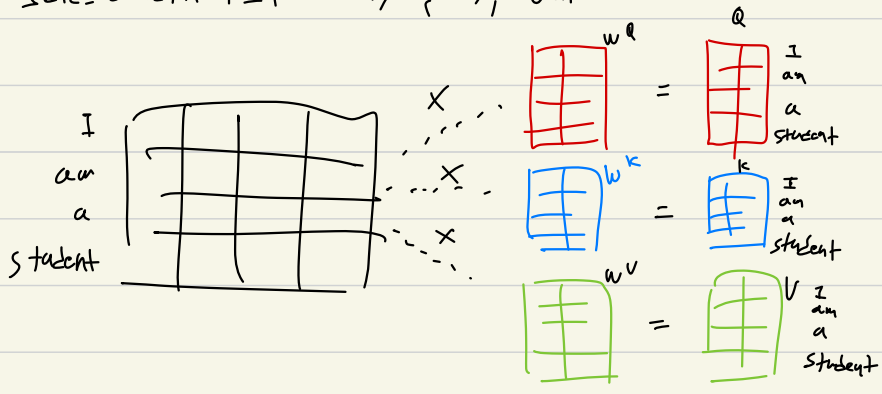
"He" 단어의 문맥을 의미는 주변의 단어들을 통해 생성

- 7월 2일 오전 10시

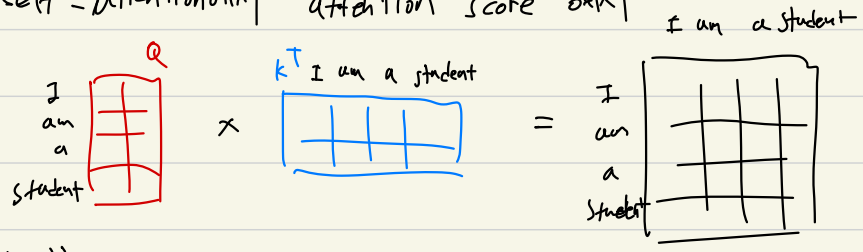
① Database or key | query, key, value.



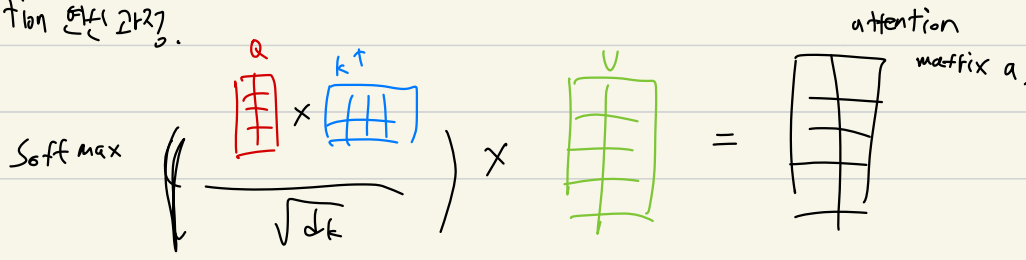
② Self-attention or query, key, value.



③ Self-Attention or key | attention score or key

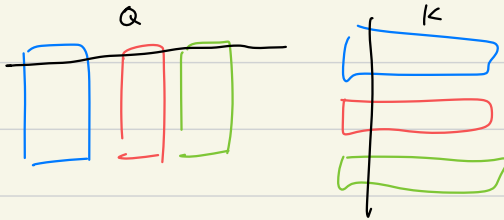


④ attention or key.



## 5. Multihead attention : 어텐션 병렬화

- Attention 연산은 행렬로 임베딩 좌표의 방향을 나눠서 연산을 수행할 수 있음.
- 나눠서 연산을 수행한 다음 합쳐서 다음 작업 수행



우리가 관망하는 시계열 예측을 위한 Transformer 모델 A MO Transformer (Arxiv)