

<복제물에 대한 경고>

본 저작물은 **저작권법 제25조 수업목적 저작물 이용 보상금제도**에 의거, **한국복제전송저작권협회와 약정을 체결하고** 적법하게 이용하고 있습니다. 약정범위를 초과하는 사용은 저작권법에 저촉될 수 있으므로

저작물의 재 복제 및 수업 목적 외의 사용을 금지합니다.

2020. 03. 30.

건국대학교(서울)한국복제전송저작권협회

<전송에 대한 경고>

본 사이트에서 수업 자료로 이용되는 저작물은 **저작권법 제25조 수업목적 저작물 이용 보상금제도**에 의거,

한국복제전송저작권협회와 약정을 체결하고 적법하게 이용하고 있습니다.

약정범위를 초과하는 사용은 저작권법에 저촉될 수 있으므로

수업자료의 대중 공개·공유 및 수업 목적 외의 사용을 금지합니다.

2020. 03. 30.

건국대학교(서울)한국복제전송저작권협회

Transformer

Attention Models

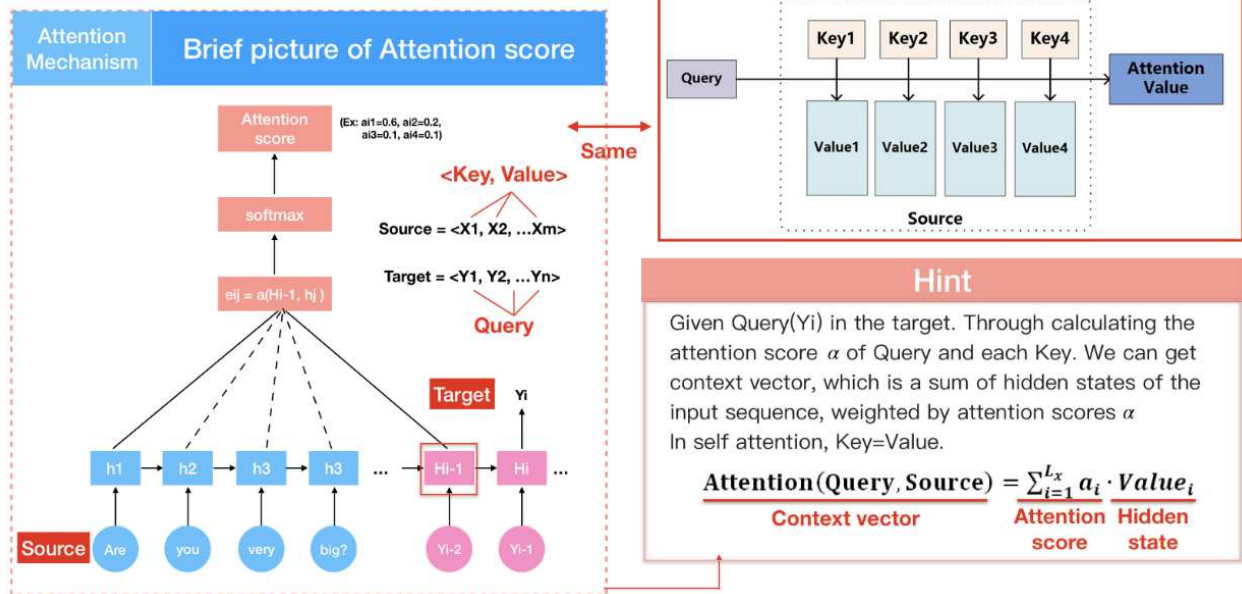


그림 출처: Ta-Chun (Bgg/Gene) Su's blog



Edited by Harksoo Kim

Attention Models in Detail

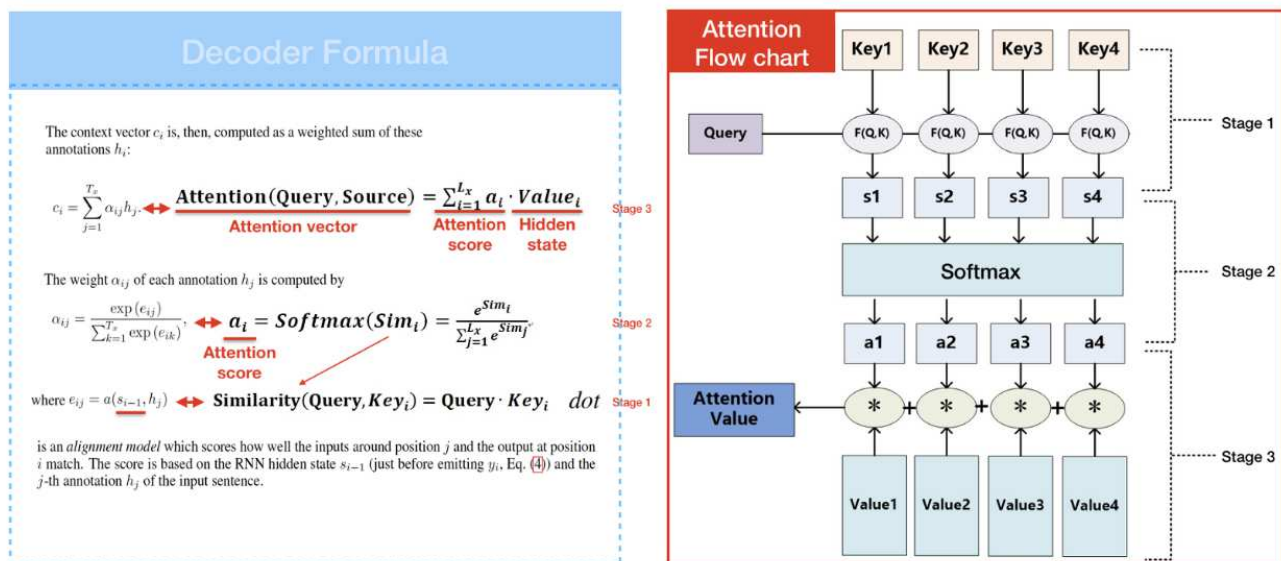


그림 출처: Ta-Chun (Bgg/Gene) Su's blog



Edited by Harksoo Kim

Problems of Attention-Based Models

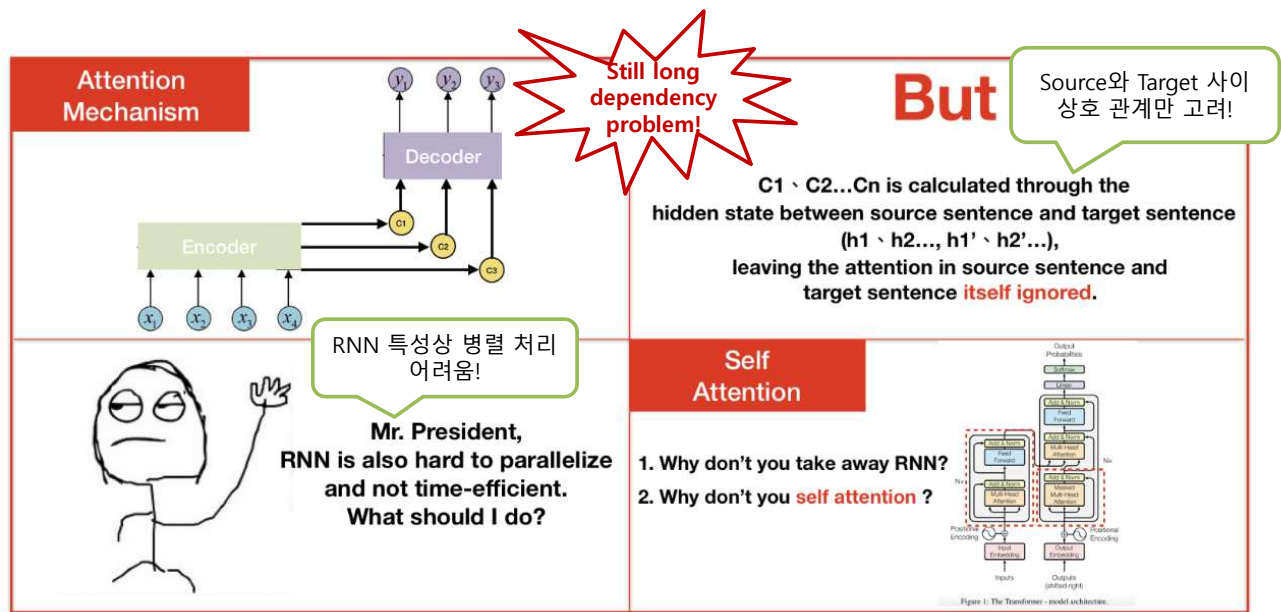


그림 출처: Ta-Chun (Bgg/Gene) Su's blog



Edited by Harksoo Kim

Transformer

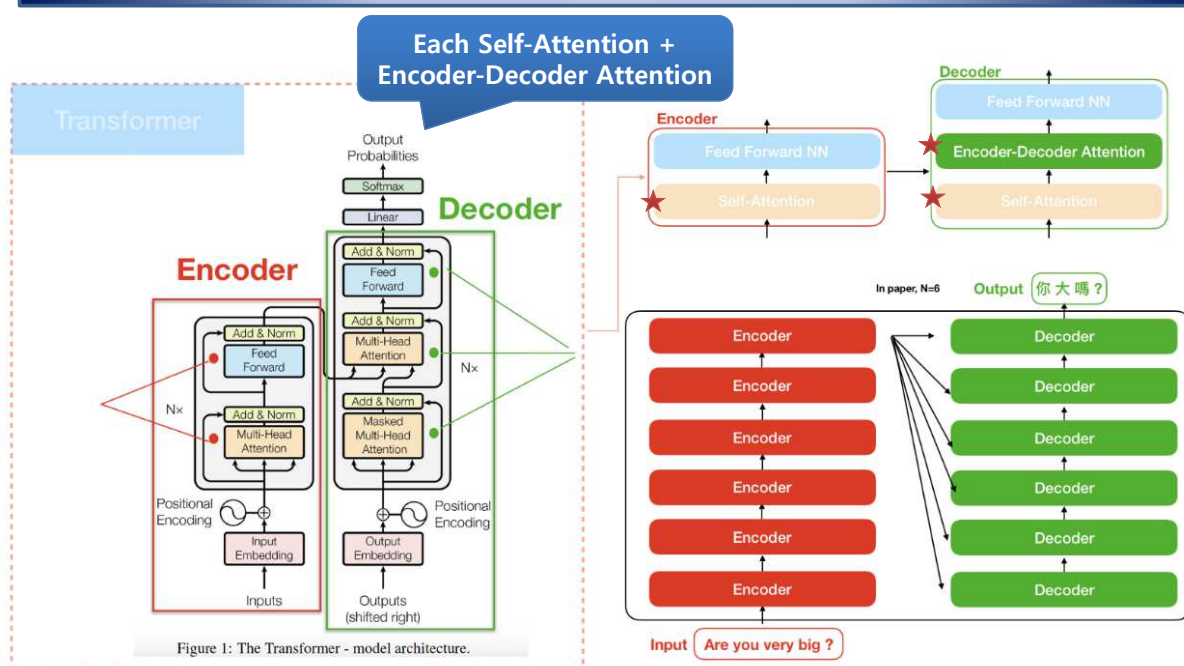
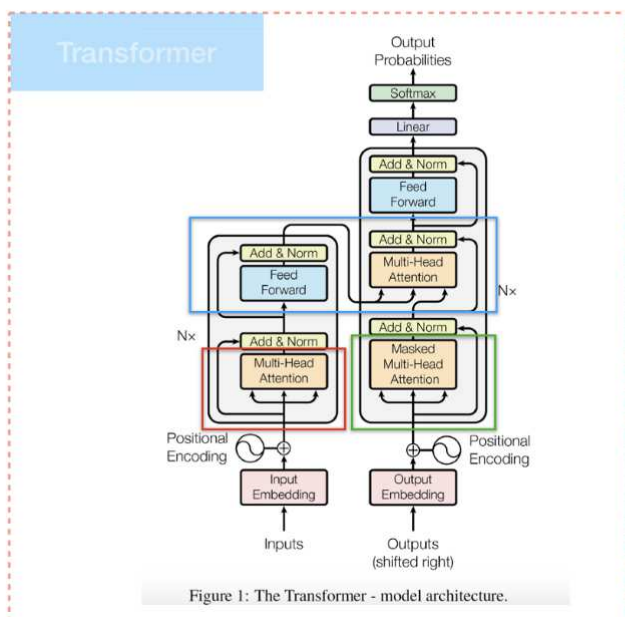


그림 출처: Ta-Chun (Bgg/Gene) Su's blog



Edited by Harksoo Kim

Attentions in Transformer



encoder self attention

1. Multi-head Attention
2. **Q**uery=**K**ey=**V**alue

decoder self attention

1. **M**asked Multi-head Attention
2. **Q**uery=**K**ey=**V**alue

encoder-decoder attention

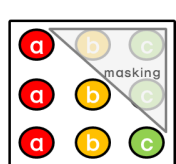
1. Multi-head Attention
2. Encoder Self attention=**K**ey=**V**alue
3. Decoder Self attention=**Q**uery

그림 출처: Ta-Chun (Bgg/Gene) Su's blog



Edited by Harksoo Kim

Scaled Dot-Product Attention



언어 모델
디코딩

Scaled Dot-Product Attention

Macro

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

where $Q \in \mathbb{R}^{n \times d_k}, K \in \mathbb{R}^{m \times d_k}, V \in \mathbb{R}^{m \times d_v}$

$$\text{Sequence } Q : n * d_k \xrightarrow{\text{attention layer}} n * d_v$$

Micro

$$\text{Attention}(q_i, K, V) = \sum_{s=1}^m \frac{1}{Z} \exp\left(\frac{\langle q_i, k_s \rangle}{\sqrt{d_k}}\right) v_s$$

Attention score

Recall Attention score $a_i = \text{Softmax}(\text{Sim}_i) = \frac{e^{\text{Sim}_i}}{\sum_{j=1}^{L_x} e^{\text{Sim}_j}}$

비교대상, 입력, 히든벡터

Hint

Q: Why divide by $\sqrt{d_k}$?

A: $\because q \cdot k = \sum_{i=1}^{d_k} q_i k_i$

Assume q_i, k_i are normal distribution,

$$\because \text{Var}\left(\sum_{i=1}^m X_i\right) = \sum_{i=1}^m \text{Var}(X_i)$$

$q \cdot k = \sum_{i=1}^{d_k} q_i k_i$, has mean 0 and variance d_k .

Sum of the dot products will grow large in magnitude if d_k is large, pushing softmax into regions 0 or 1.

To counteract this effect, we scale the dot products by $\sqrt{d_k}$

그림 출처: Ta-Chun (Bgg/Gene) Su's blog



Edited by Harksoo Kim

Calculation of Attentions

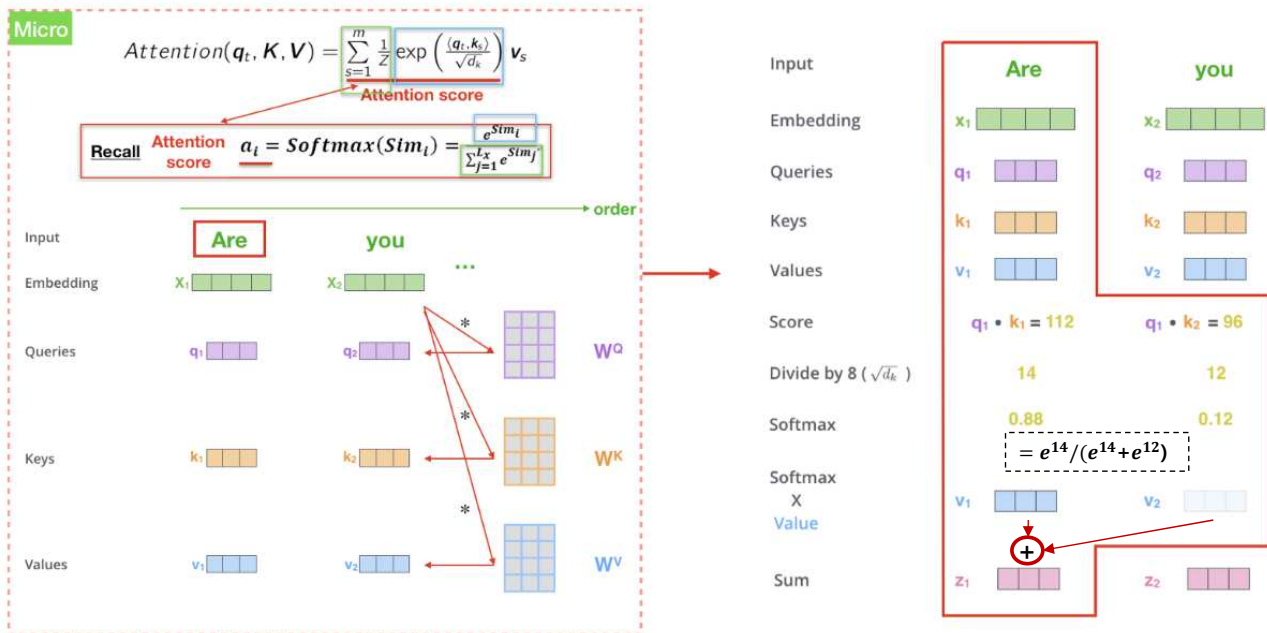


그림 출처: Ta-Chun (Bgg/Gene) Su's blog



Edited by Harksoo Kim

Calculation of Attentions

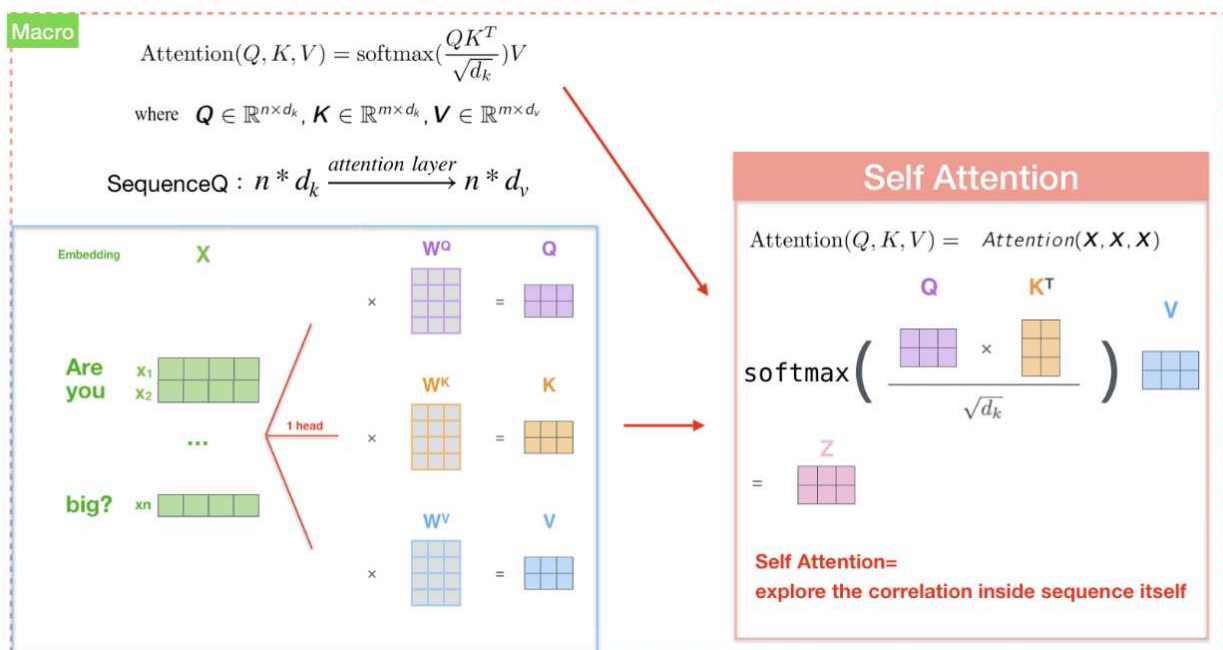


그림 출처: Ta-Chun (Bgg/Gene) Su's blog



Edited by Harksoo Kim

확인 문제

- 다음과 같이 단어 임베딩이 주어졌을 때, self-attention score를 계산하시오. (소수점 이하 두 자리에서 반올림)

- are: [1,1], you: [2,1]

- $\text{root}(2)=1.4$ 로 계산

$$Attention(q_t, K, V) = \sum_{s=1}^m \frac{1}{Z} \exp\left(\frac{(q_t, k_s)}{\sqrt{d_k}}\right) v_s$$

Scaled dot-product

	are	you
are	?	
you		

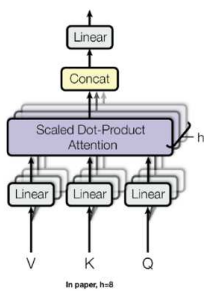
Self-attention score

	are	you
are	?	
you		

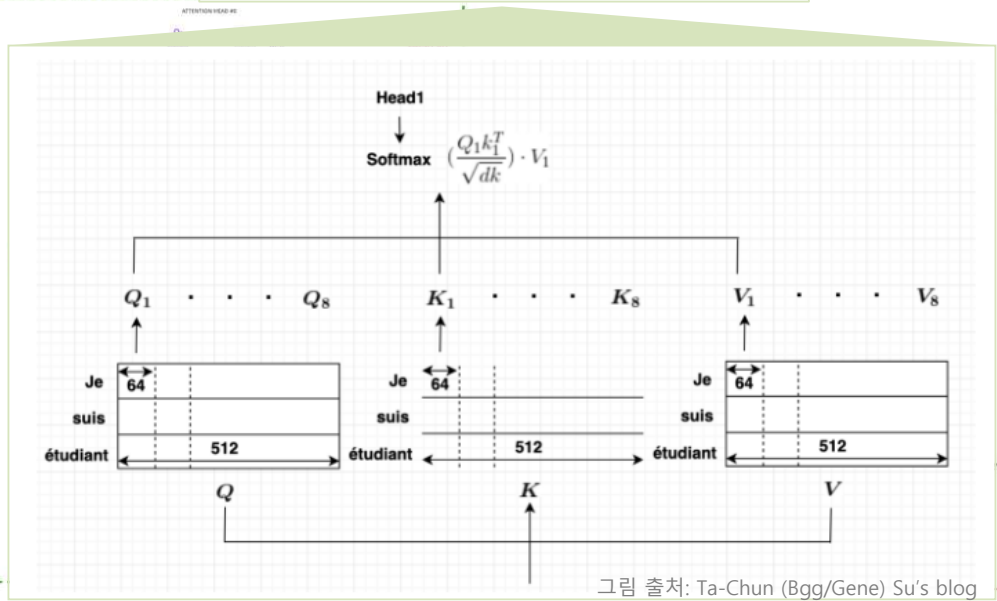


Edited by Harksoo Kim

Multi-Head Attention



1. 작은 범위에서 softmax → 뚜렷한 특징만 살아남는 효과
2. 다양한 관점에서 문장을 바라보는 효과



Edited by Harksoo Kim

Layer Norm. & Residual Conn.

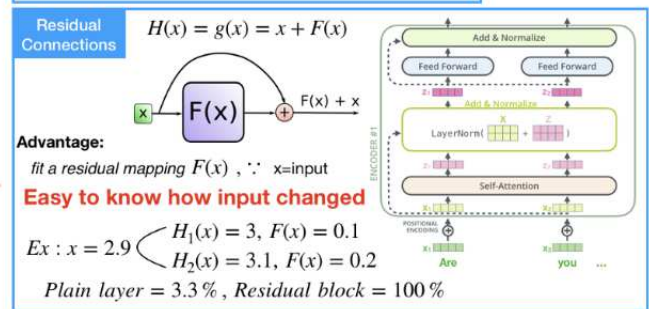
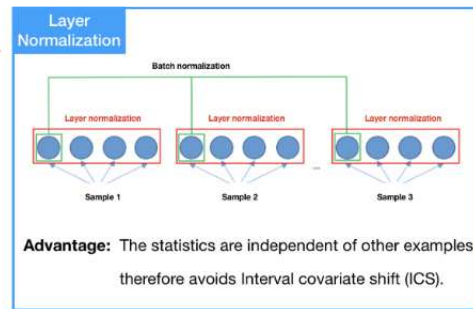
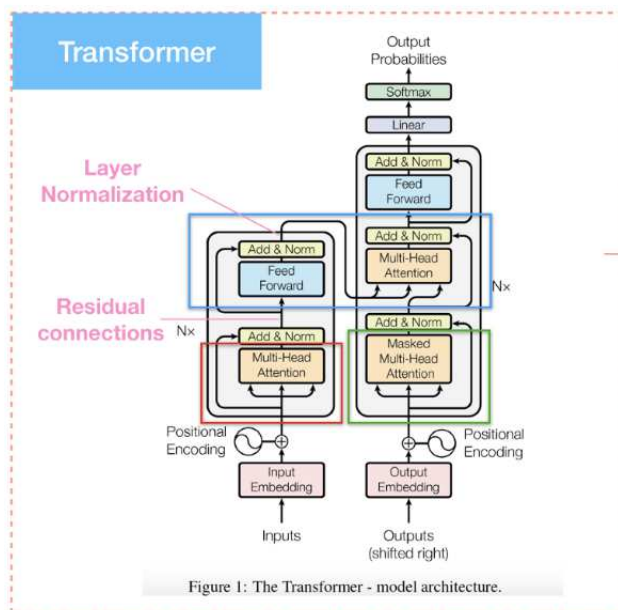


그림 출처: Ta-Chun (Bgg/Gene) Su's blog

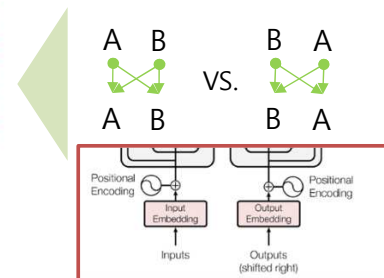


Edited by Harksoo Kim

Position Encoding

Problem

The multi-head attention network **cannot** naturally **make use of** the position of the words in the input sequence.
 The output of the multi-head attention network would be **the same** for the same sentences in different order.



Sol

Positional Encoding

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

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

if shape of enc, dec = $[T, d_{model}]$ pos: position of the word
 \rightarrow then $pos \in [0, T), i \in [0, d_{model})$ i: Element i in d_{model}

Advantage:

$PE[pos+k]$ can be represented as a linear function of $PE[pos]$, so the relative position between different embeddings can be easily inferred

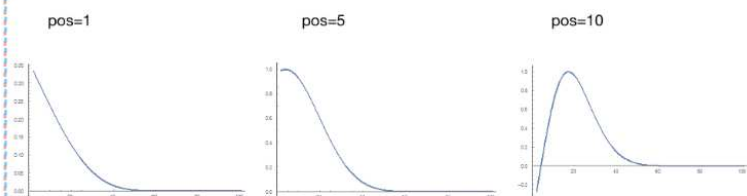
$$\begin{aligned} \sin(\alpha + \beta) &= \sin\alpha\cos\beta + \cos\alpha\sin\beta \\ \cos(\alpha + \beta) &= \cos\alpha\cos\beta - \sin\alpha\sin\beta \end{aligned} \quad \because \text{Trigonometric Periodicity}$$

\therefore **Learned the relation between relative and absolute position**

Then Word embedding + positional encoding

(sum, concat...)

if $d_{model} = 512$:



\rightarrow **Three plots are the same by doing linear transformation**

그림 출처: Ta-Chun (Bgg/Gene) Su's blog



Edited by Harksoo Kim

Linear & Softmax

Linear + Softmax

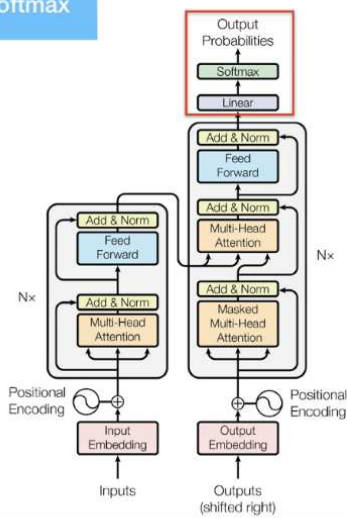


Figure 1: The Transformer - model architecture.

Which word in our vocabulary is associated with this index?
Get the index of the cell with the highest value (argmax)

big

5

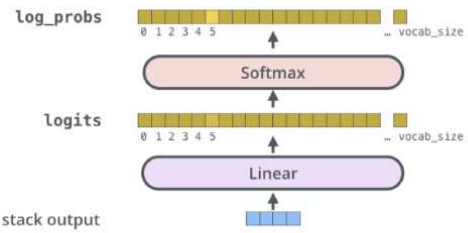


그림 출처: Ta-Chun (Bgg/Gene) Su's blog



Edited by Harksoo Kim

질의응답

Q&A

Homepage: <http://nlp.konkuk.ac.kr>
E-mail: nlpdrkim@konkuk.ac.kr



Edited by Harksoo Kim