

# Lecture 8

## Self-Attention and Transformers

투빅스 18기  
이화여자대학교 이선민

# Lecture Plan



- 01** From RNN to attention-based NLP models
- 02** The Transformer model
- 03** Great results with Transformers
- 04** Drawbacks and variants of Transformers

## Problem of RNN

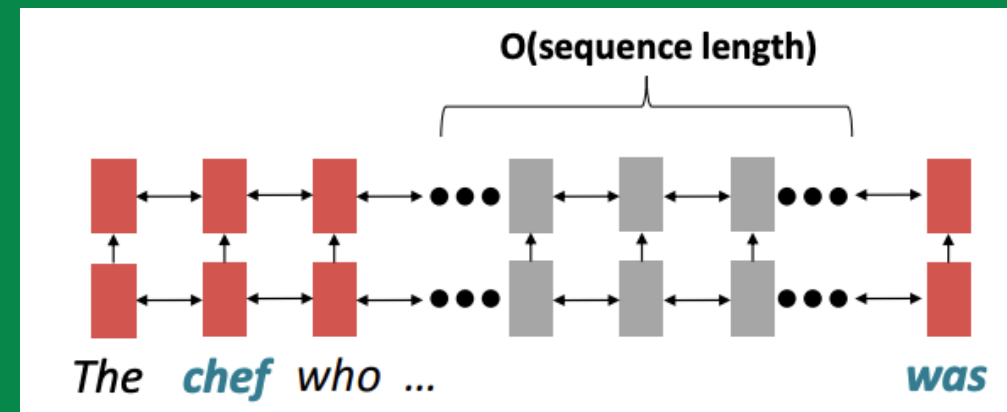
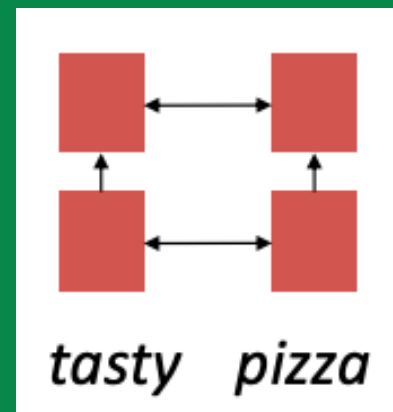
The Transformer

Great results with Transformers

Drawbacks and variants of Transformers

## Linear Interaction Distance

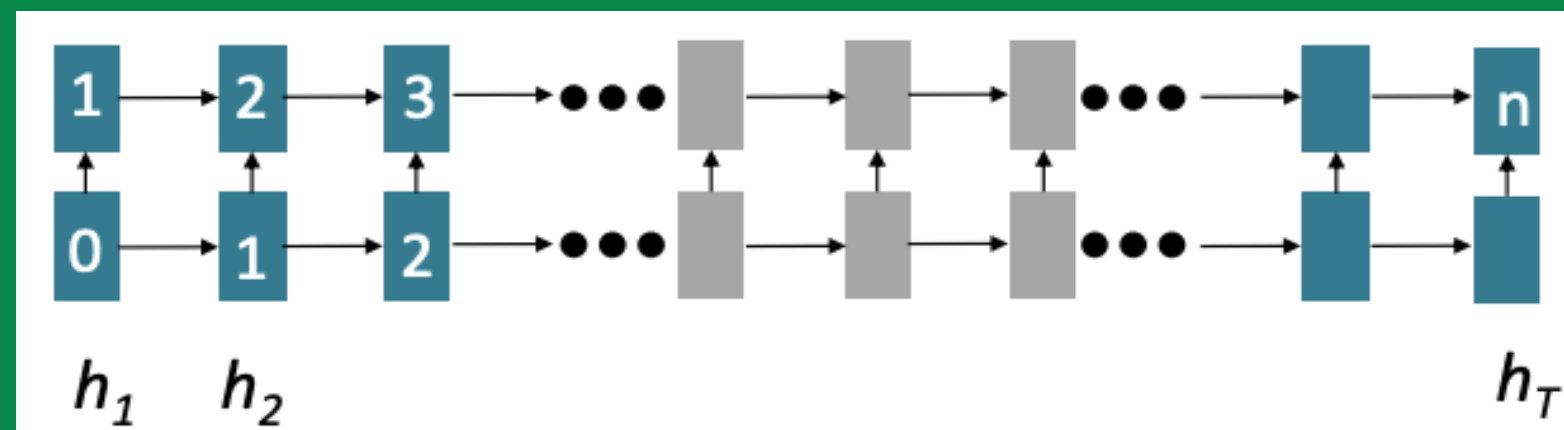
- Previous hidden state -> Future hidden state 방향으로 학습 진행
- 멀리 떨어진 단어 사이의 상호 작용이 어려움



Vanishing Gradient problem !

## Lack of Parallelizability

- Previous hidden state가 연산되어야 Next hidden state 연산 가능
- GPU를 사용한 병렬 연산 불가능



Other Building Blocks

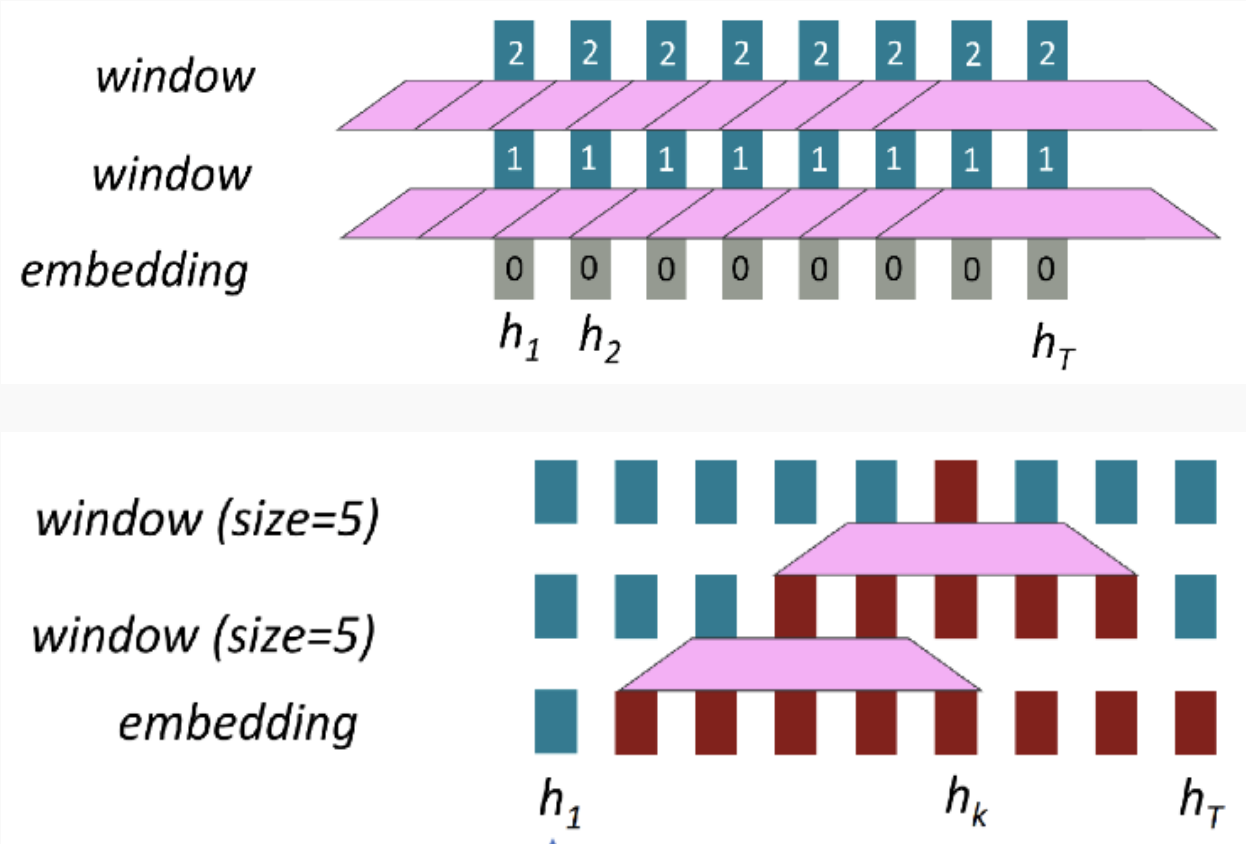
RNN을 대신하여 사용할만한 다른 Building Block들을 살펴보자.

Other Building Block

The Transformer

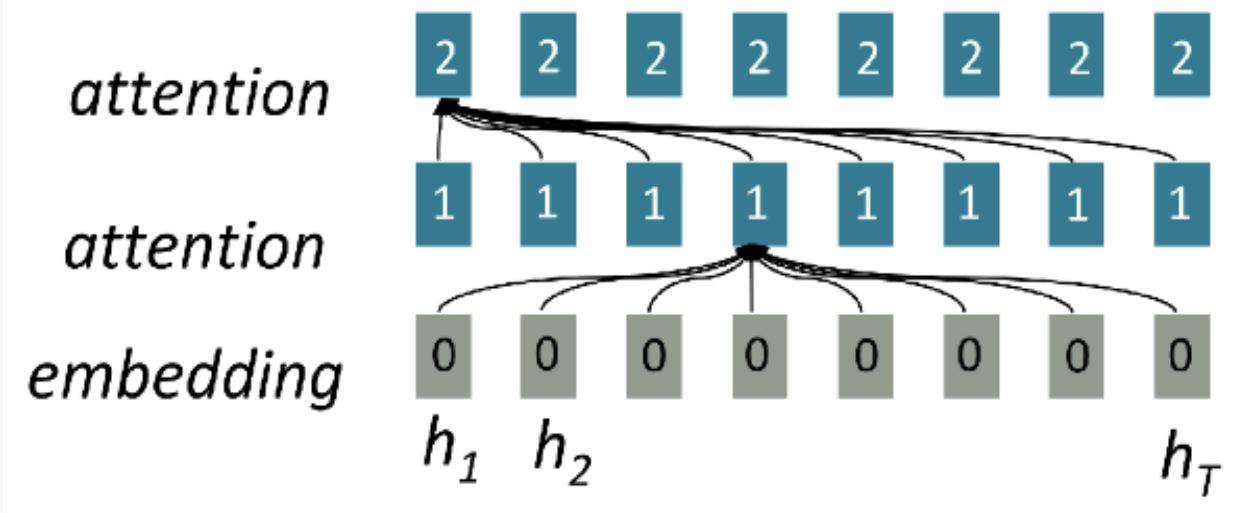
Great results with Transformers

Drawbacks and variants of Transformers



Long Distance Dependency !

Word Window



2 problems Solved !

- Embedding layer : 각각 독립적으로 연산 가능
- Attention layer : Previous layer의 연산만 완료되면, Next layer의 모든 sequence state에 대해 계산 가능

Self-Attention

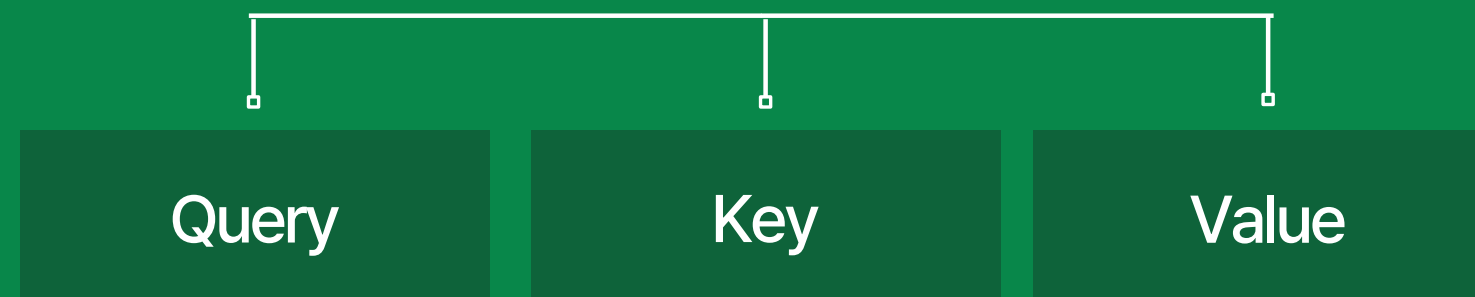
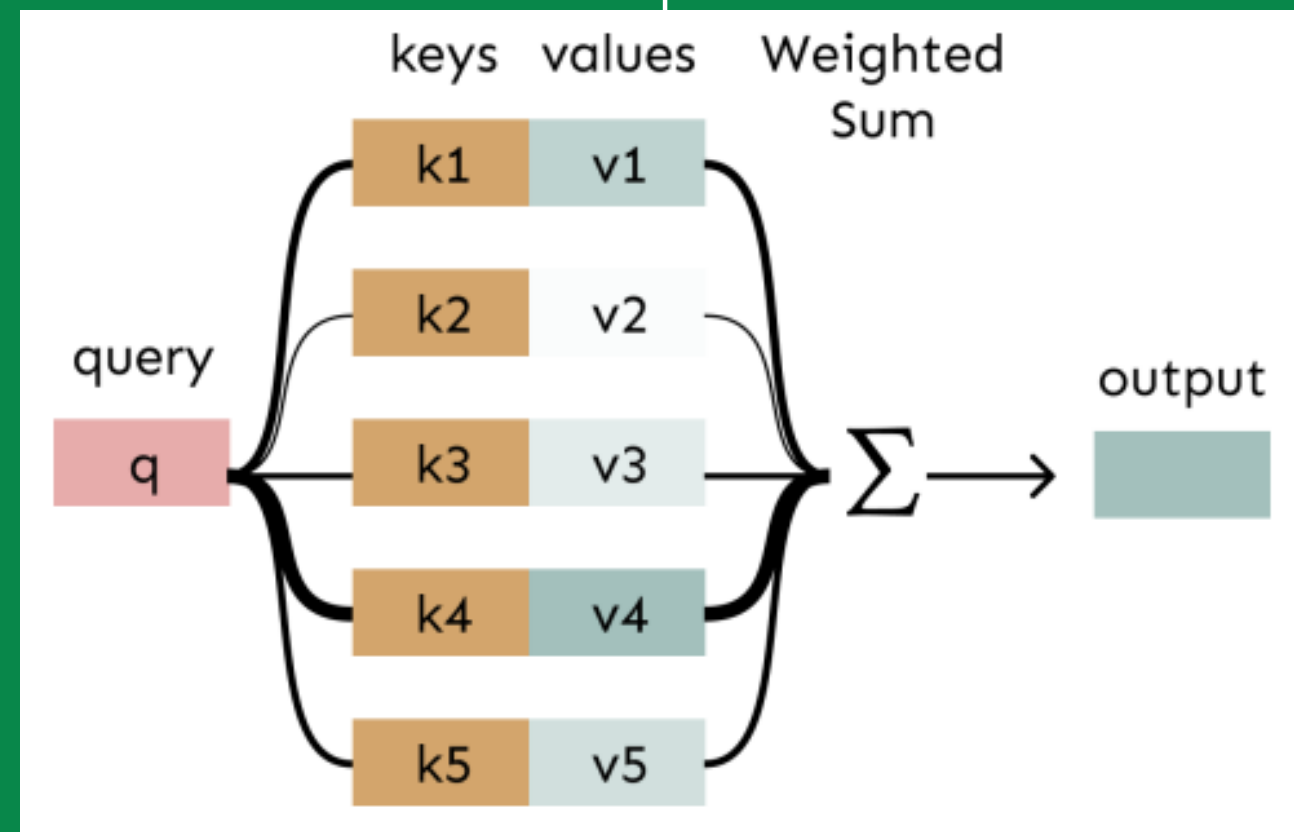
## Self-Attention

The Transformer

Great results with Transformers

Drawbacks and variants of Transformers

# Self-Attention



Query

현재 보고있는 단어의 representation (다른 단어를 평가하는 기준)

Key

Query와 관련된 단어를 찾을 때, label처럼 활용되는 vector

Value

Query와 Key를 통해 탐색하여 실제로 사용할 값

LSTM layer처럼 쌓기만하면 해결될까?

참고:

[https://github.com/yookyungkho/DSBA\\_CS224N\\_2021/blob/main/slides/dsba\\_cs224n2021\\_lec09\\_gunhono.pdf](https://github.com/yookyungkho/DSBA_CS224N_2021/blob/main/slides/dsba_cs224n2021_lec09_gunhono.pdf)

## Self-Attention

The Transformer

Great results with Transformers

Drawbacks and variants of Transformers

## Problem of Self-Attention

3가지 문제 : 순서에 대한 정보 없음, 선형 결합만 존재함, Future sequence data 활용 가능함

### Sequence Order

병렬 처리로 동시에 처리하므로, 순서에 대한 정보 없음

**위치 벡터 사용 !**

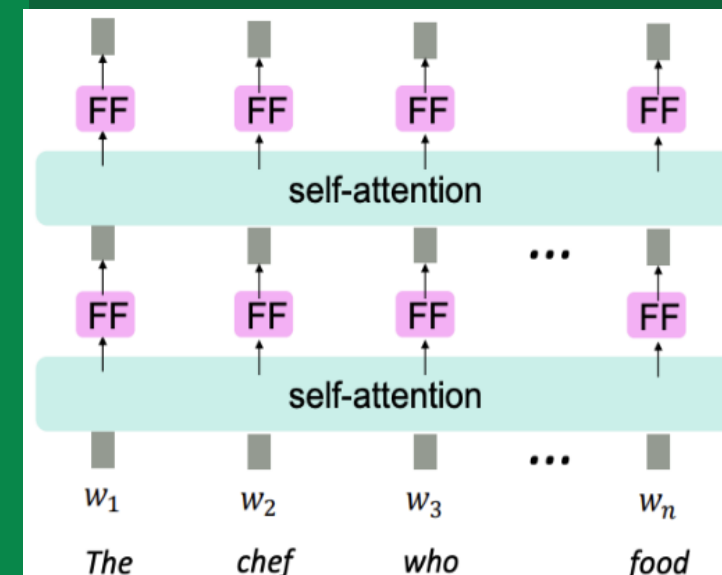
ex. Sinusodial

$$p_i = \begin{pmatrix} \sin(i/10000^{2*1/d}) \\ \cos(i/10000^{2*1/d}) \\ \vdots \\ \sin(i/10000^{2*\frac{d}{2}/d}) \\ \cos(i/10000^{2*\frac{d}{2}/d}) \end{pmatrix}$$

### Nonlinearities

단순한 Weighted averages의 선형 결합

**FF Network !**



### Future Sequence

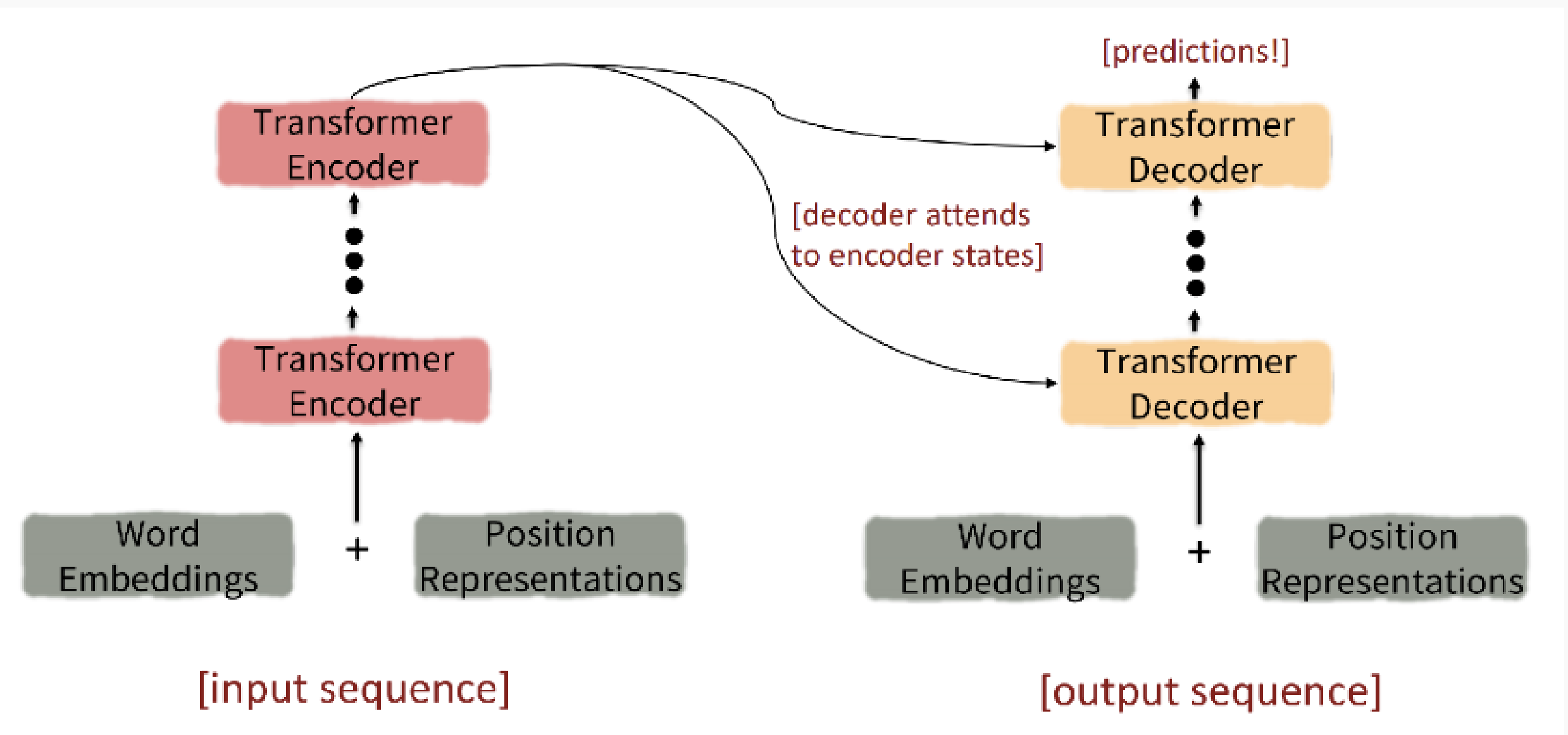
Decoder에서 Future Sequence data를 볼 수 있음

**Masking !**

	[START]	The	chef	who
[START]		$-\infty$	$-\infty$	$-\infty$
The			$-\infty$	$-\infty$
chef				$-\infty$
who				

## Transformer Model

Transformer Encoder & Decoder



Transformer Query, Key, Value

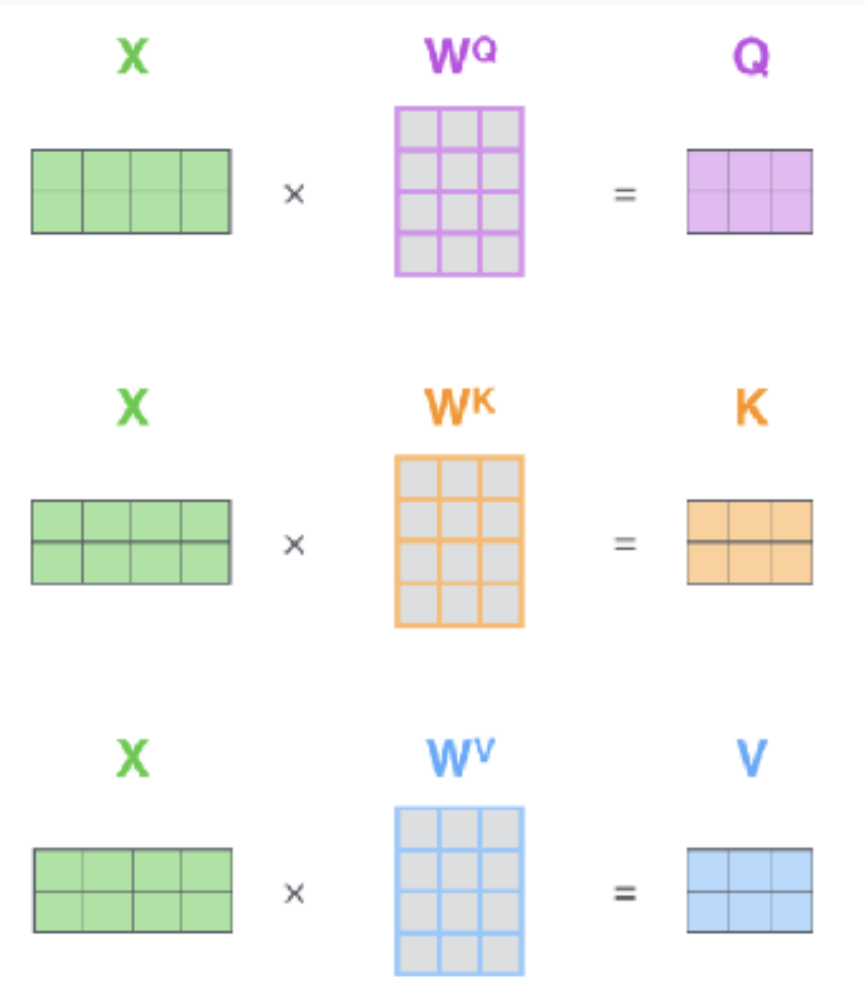
Self-Attention과 달리 학습하는 대상의 Query, Key, Value를 행렬로 나타낸다.

Self-Attention

The Transformer

Great results with Transformers

Drawbacks and variants of Transformers



행렬 Q, K, V  
= 우리가 찾아야하는 미지수

- Input vector가 d dimension을 가질 때, 그를 결합한 행렬을 X라 하자.
- Query, Key, Value를 계산하기 위해 행렬 X와 행렬 K, Q, V를 dot product한다.
- Attention score 계산을 위해 다음의 행렬 연산을 수행한다.  $XQ(XK)^T$
- Input vector가 d dimension을 가질 때, 그를 결합한 행렬을 X라 하자.
- Query, Key, Value를 계산하기 위해 행렬 X와 행렬 K, Q, V를 dot product한다.
- Attention score 계산을 위해 다음의 행렬 연산을 수행한다.  $XQ(XK)^T$

- Attention score 결과를 softmax를 통해 가중치를 얻고 그를 가중합하여 output을 얻는다.

$$X = [x_1; \dots; x_T] \in \mathbb{R}^{T \times d}$$

$$XK \in \mathbb{R}^{T \times d}, XQ \in \mathbb{R}^{T \times d}, XV \in \mathbb{R}^{T \times d}$$

$$XQ \quad K^T X^T = XQK^T X^T \in \mathbb{R}^{T \times T}$$

$$X = [x_1; \dots; x_T] \in \mathbb{R}^{T \times d}$$

$$XK \in \mathbb{R}^{T \times d}, XQ \in \mathbb{R}^{T \times d}, XV \in \mathbb{R}^{T \times d}$$

$$XQ \quad K^T X^T = XQK^T X^T \in \mathbb{R}^{T \times T}$$

$$\text{softmax} \left( XQK^T X^T \right) XV = \text{output} \in \mathbb{R}^{T \times d}_6$$



## Transformer Multi-head Attention

Look at Different things, and construct value vectors differently !

Self-Attention

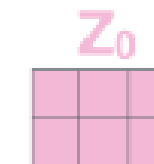
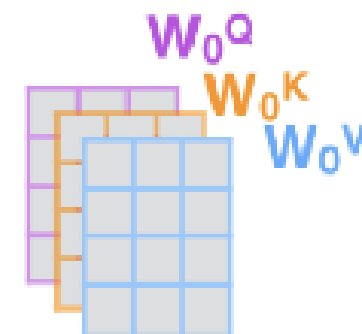
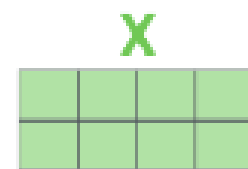
### The Transformer

Great results with  
Transformers

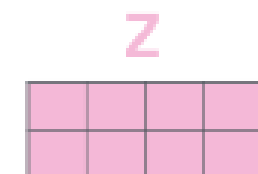
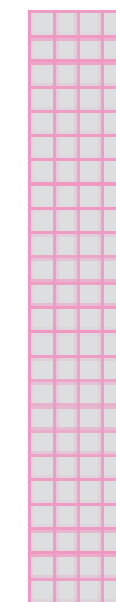
Drawbacks and  
variants of  
Transformers

- 1) This is our input sentence\*
- 2) We embed each word\*
- 3) Split into 8 heads. We multiply  $X$  or  $R$  with weight matrices
- 4) Calculate attention using the resulting  $Q/K/V$  matrices
- 5) Concatenate the resulting  $Z$  matrices, then multiply with weight matrix  $W^O$  to produce the output of the layer

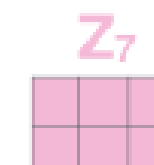
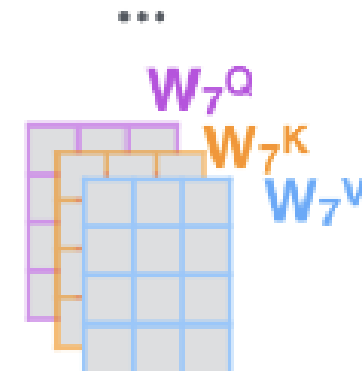
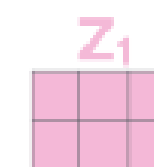
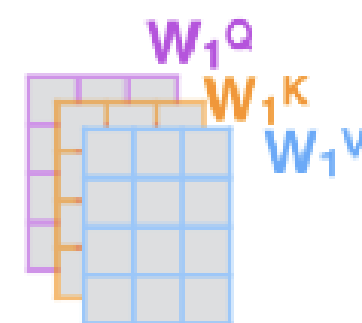
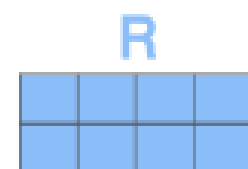
Thinking  
Machines



$W^O$



\* In all encoders other than #0, we don't need embedding. We start directly with the output of the encoder right below this one



## Transformer Multi-head Attention

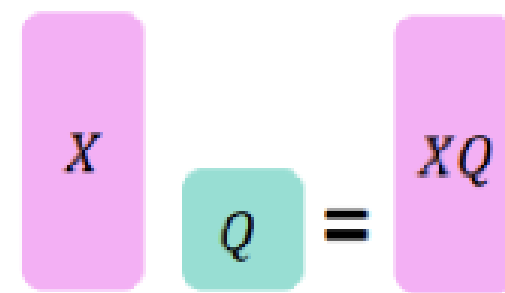
Single-head Attention과 계산 과정이 동일하다.

$$Q_\ell, K_\ell, V_\ell \in \mathbb{R}^{d \times \frac{d}{h}} \quad h \text{ 는 head의 개수, } \ell = 1, 2, \dots, h$$

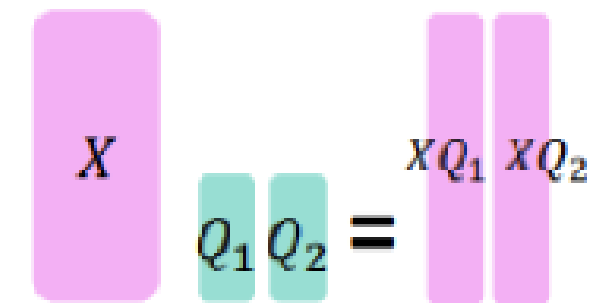
$$\text{output}_\ell = \text{softmax}(XQ_\ell K_\ell^\top X^\top) * XV_\ell, \text{ where } \text{output}_\ell \in \mathbb{R}^{d/h}$$

$$\text{output} = Y[\text{output}_1; \dots; \text{output}_h], \text{ where } Y \in \mathbb{R}^{d \times d}$$

### Single-head attention (just the query matrix)



### Multi-head attention (just two heads here)



독립적으로 동시에 수행되므로, Self-Attention과 계산량이 동일하다.

## Transformer Multi-head Attention

$W^o$  matrix를 곱해주는 이유는 무엇일까?

Self-Attention

### The Transformer

Great results with Transformers

Drawbacks and variants of Transformers

#### Question ?

Multi-head Attention을 할 때, 각 헤드들의 output을 concat 한 후  $W^o$  matrix를 곱해줍니다.

각 layer의 input output 차원을 동일하게 유지하기 위해  $W^o$  matrix를 곱해준다고 이해를 했습니다.

본 논문의 구현상으로는 512차원을 8개의 헤드로 나눠 각각의 헤드가 64차원을 가지고, concat 하면  $d_{model}$ 의 차원과 같아지는데  $W^o$  matrix를 곱해주는 이유는 무엇인가요?

질문을 정리하자면,  $W^o$  matrix의 역할은 무엇인가요?

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

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

#### Correct Answer ?

if  $T = 10, d = 512, h = 8$

input :  $X (R^{T \times d})$   
 $R^{10 \times 512}$

self-attention :  $R^{T \times d}$   
 $R^{10 \times 512}$

multi-head attention :  $R^{T \times \frac{d}{h}}$   
 $R^{10 \times \frac{512}{8}}$   $\xrightarrow{\text{concat}}$   $R^{(\frac{T \times d}{h}) \times h}$   
 $R^{(10 \times \frac{512}{8}) \times 8} = R^{10 \times 512}$

$W^o : R^{(h \times \frac{d}{h}) \times d}$   
목표 차원수 = 64  
 $R^{(8 \times \frac{512}{8}) \times 64} = R^{512 \times 64}$

최종 Output : multi-head attention output  $\times W^o$  =  $R^{10 \times 64}$   
 $R^{10 \times 512} \quad R^{512 \times 64}$

$W^o$  matrix는 목표 차원수로 맞춰주는 역할을 한다.

## Transformer Modeling Improvements

### Residual Connection

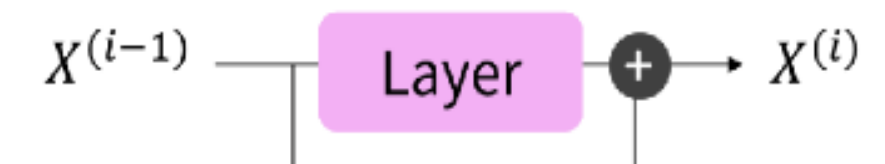
자기 자신을 더해준다.

- 미분 결과값이 너무 작아 gradient propagation이 원활하지 않은 경우, gradient를 보정함
- 기울기 Smoothing 효과 -> local minimum에 빠지지 않도록 함

$$X^{(i)} = \text{Layer}(X^{(i-1)})$$



$$X^{(i)} = X^{(i-1)} + \text{Layer}(X^{(i-1)})$$



Inductive Bias !

이전 state보다 얼마나 달라졌는지 학습한다.

Self-Attention

**The Transformer**

Great results with  
Transformers

Drawbacks and  
variants of  
Transformers

## Transformer Modeling Improvements

Layer Normalization

Gradient를 normalize 한다.

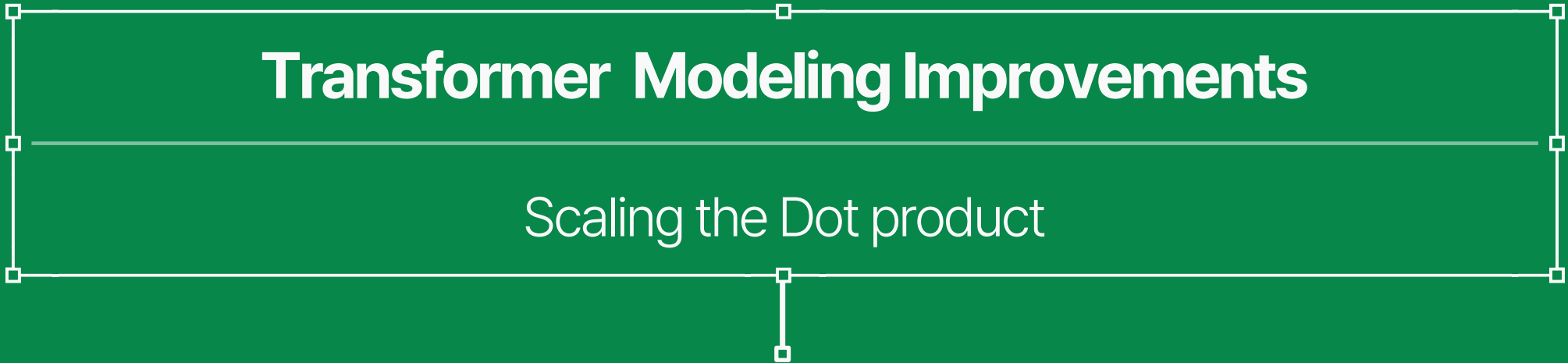
- 한 layer에서 하나의 input sample  $x$ 에 대해 모든 feature에 대한 평균과 분산을 구해 normalization 함

$$\text{output} = \frac{x - \mu}{\sqrt{\sigma + \epsilon}} * \gamma + \beta$$

Normalize by scalar mean and variance

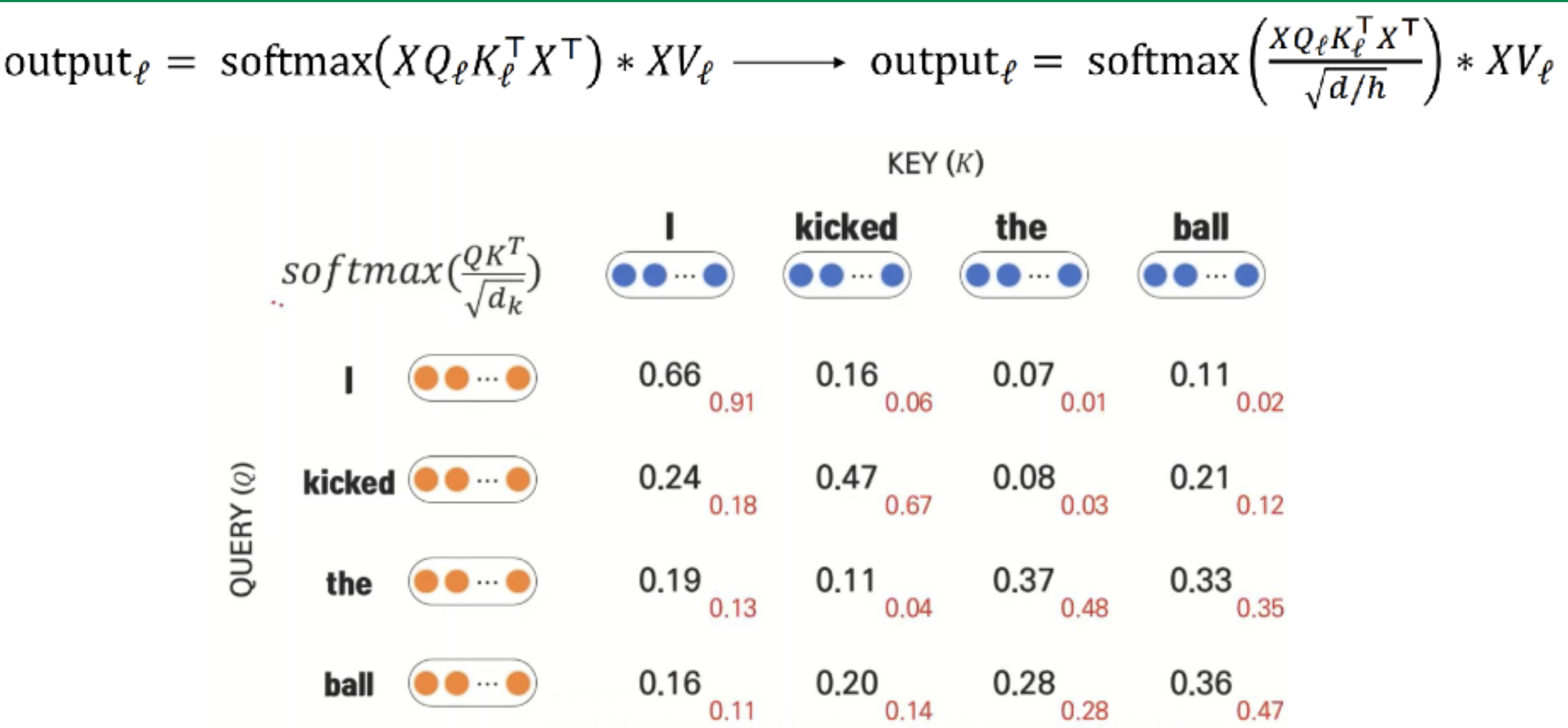
Modulate by learned elementwise gain and bias

- Self-Attention
- The Transformer
- Great results with Transformers
- Drawbacks and variants of Transformers



Dot product의 결과가 너무 커지지 않도록 유지한다.

- 차원 수가 늘어날수록, dot product의 결과가 커짐 = softmax 결과가 치중됨
- 모든 sequence에서 gradient propagating이 잘 되도록 유지해줌
- Attention score을 좀 더 다양한 vector에 분배함



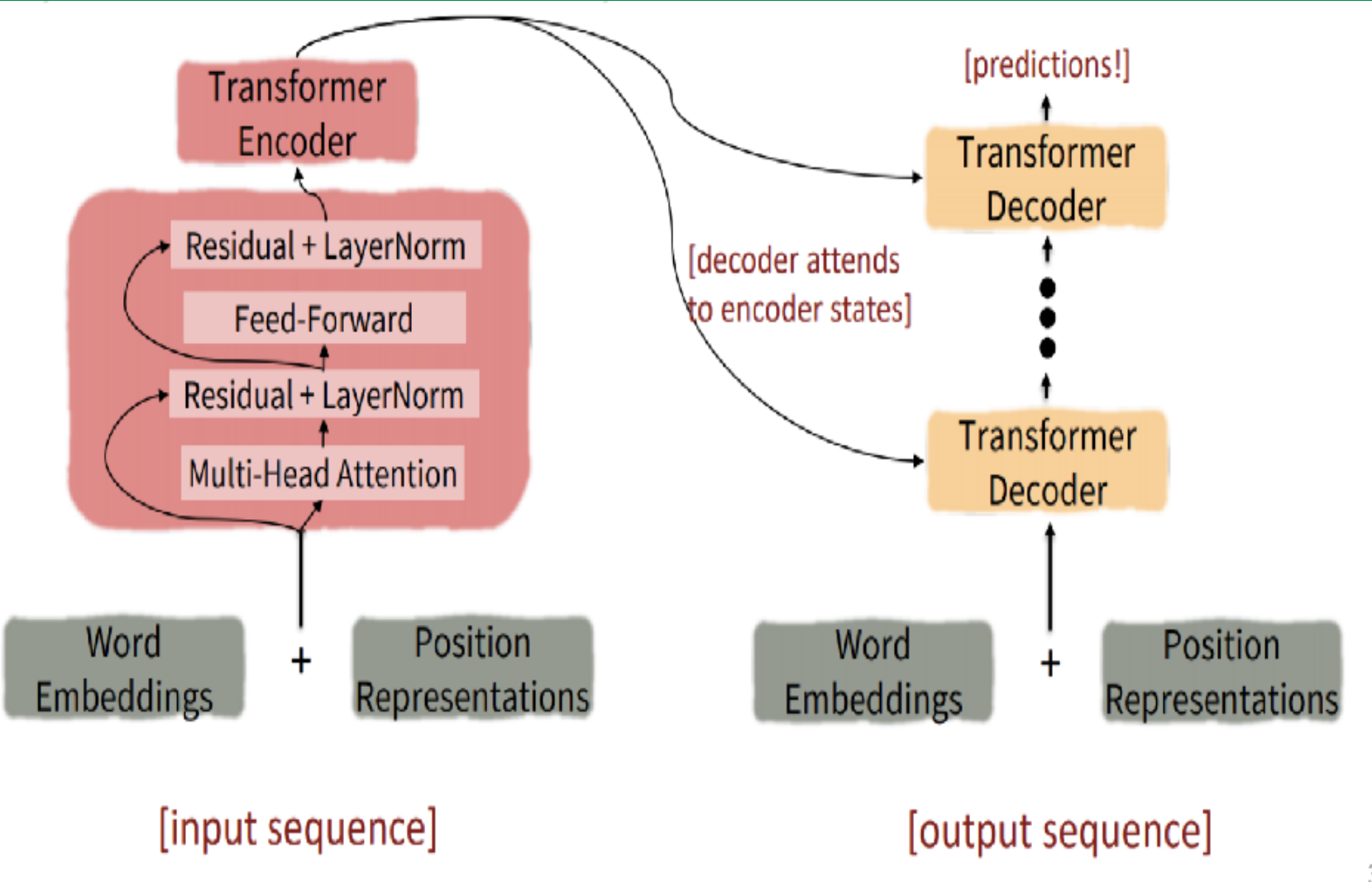
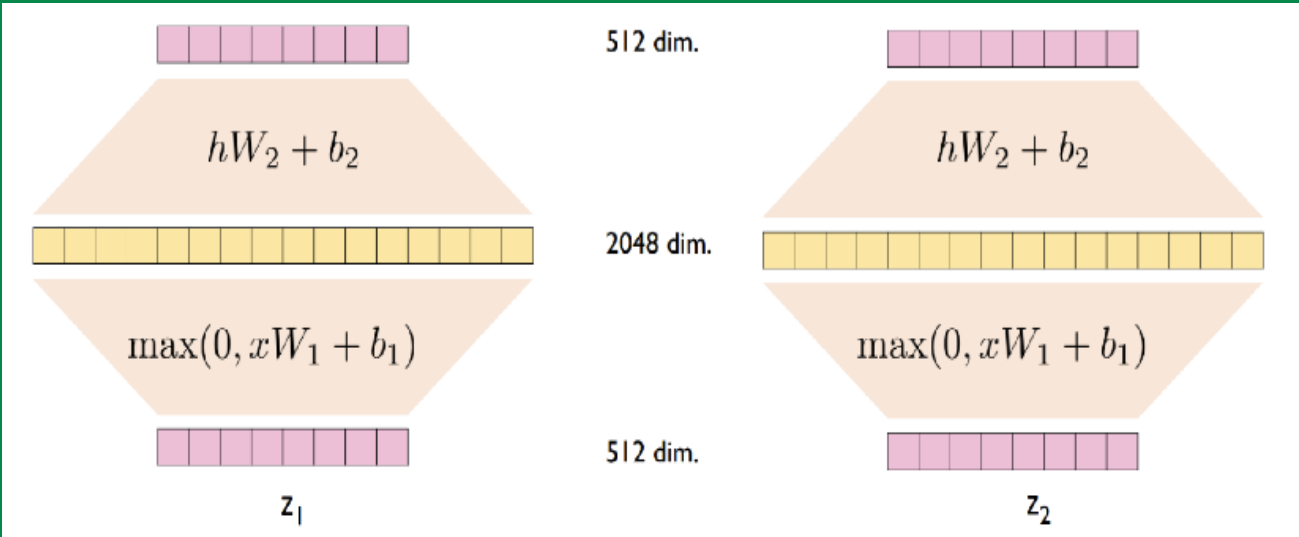
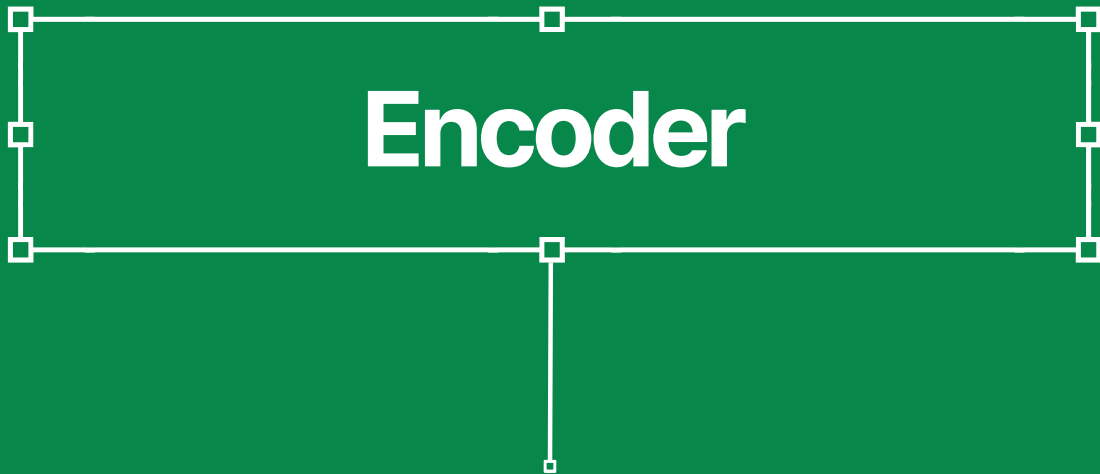


Self-Attention

The Transformer

Great results with Transformers

Drawbacks and variants of Transformers





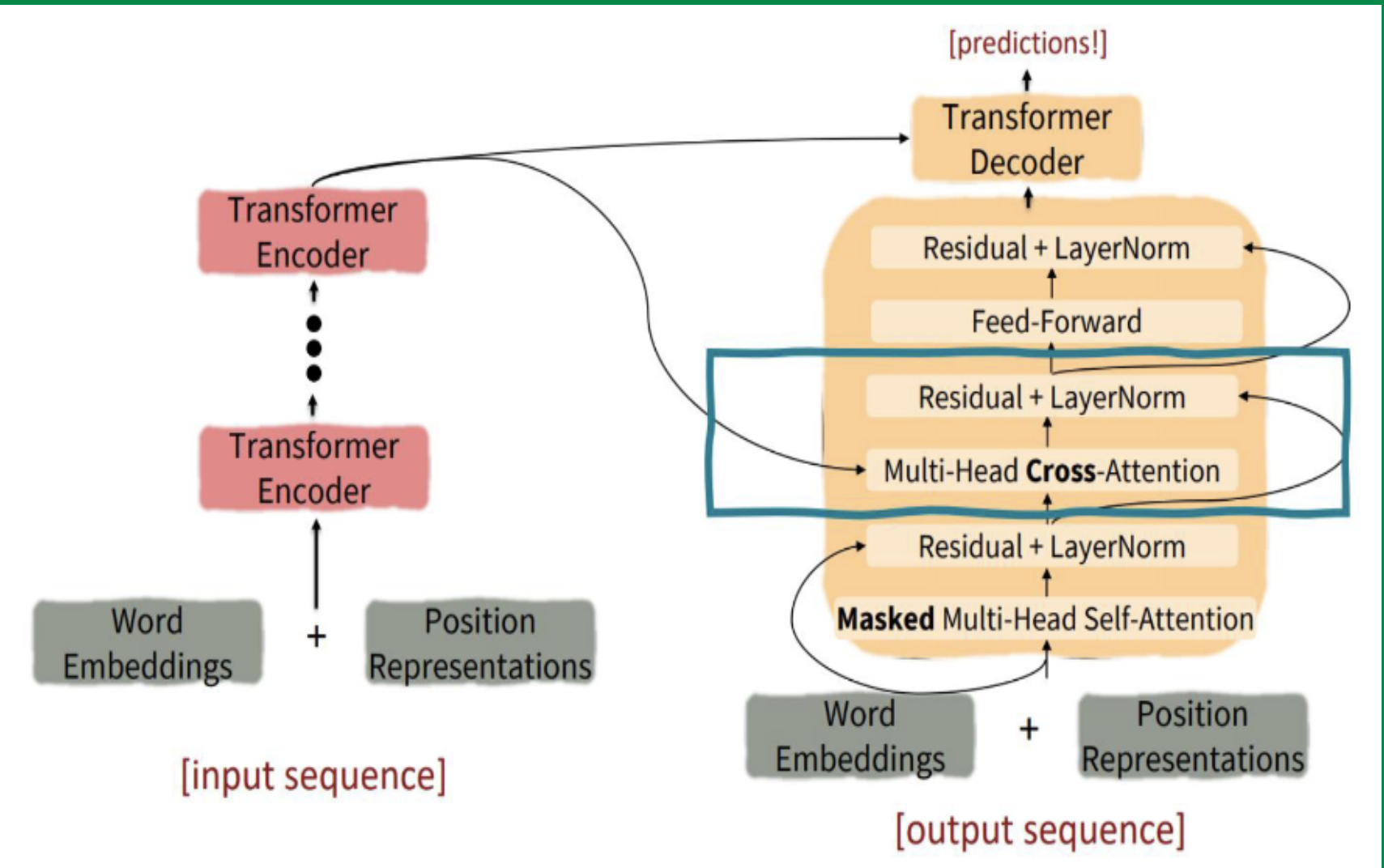
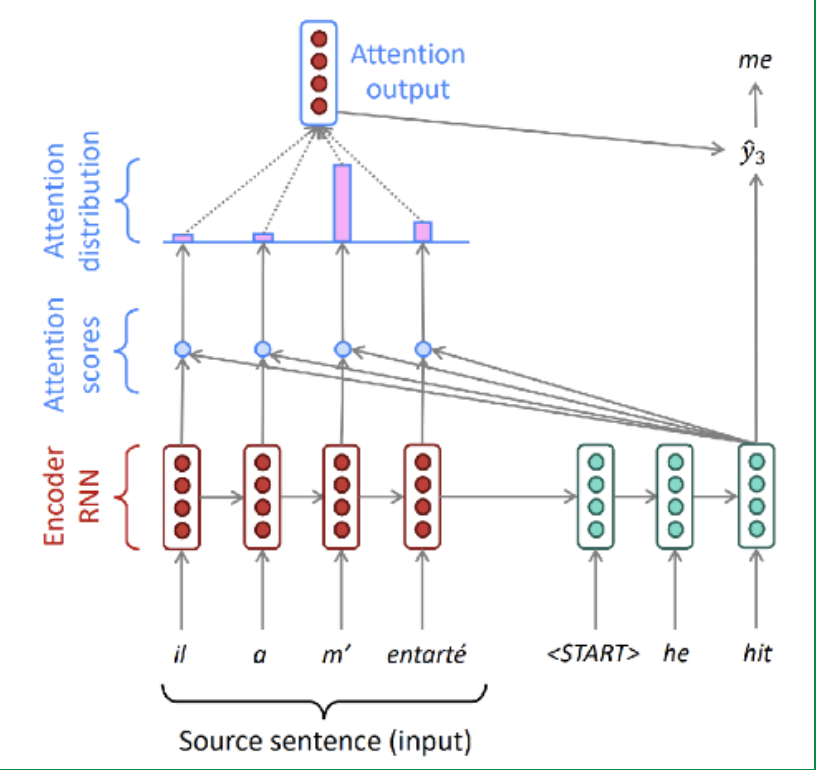
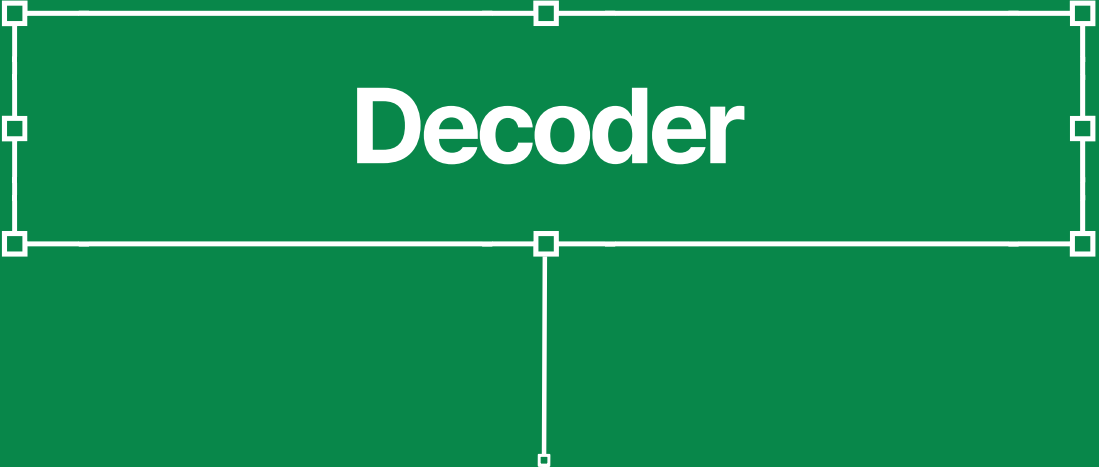
# ToBigs

Self-Attention

## The Transformer

Great results with Transformers

Drawbacks and variants of Transformers



Input

Word Embedding + Position Representation

Multi head Attention

Masked Multi head Attention

Residual + Layer Nomalization

Transformer의 성능을 높여줌

Multi head Attention

Multi head Cross Attention

Residual + Layer Nomalization

Transformer의 성능을 높여줌

Transformer Encoder

여러 개의 Encoder를 거친 후, Decoder로 들어감



## Encoder & Decoder

Self-Attention

The Transformer

Great results with Transformers

Drawbacks and variants of Transformers

- $h$ : encoder의 output vector  
 $z$ : decoder의 input vector
- Decoder에서 현재 처리하는 단어의 **Query**를 가져오고, Encoder의 **Key**로 탐색한 결과를 Encoder의 **Value**로 가중 합한다.
- Attention score 계산을 위해 다음의 행렬 연산을 수행한다.  $ZQ(HK)^T$
- Attention score 결과를 softmax를 통해 가중치를 얻고 그를 가중합하여 output을 얻는다.

$$\text{Let } H = [h_1; \dots; h_T] \in \mathbb{R}^{T \times d}$$

$$\text{Let } Z = [z_1; \dots; z_T] \in \mathbb{R}^{T \times d}$$

$$ZQ \quad K^T H^T = ZQK^T H^T \in \mathbb{R}^{T \times T}$$

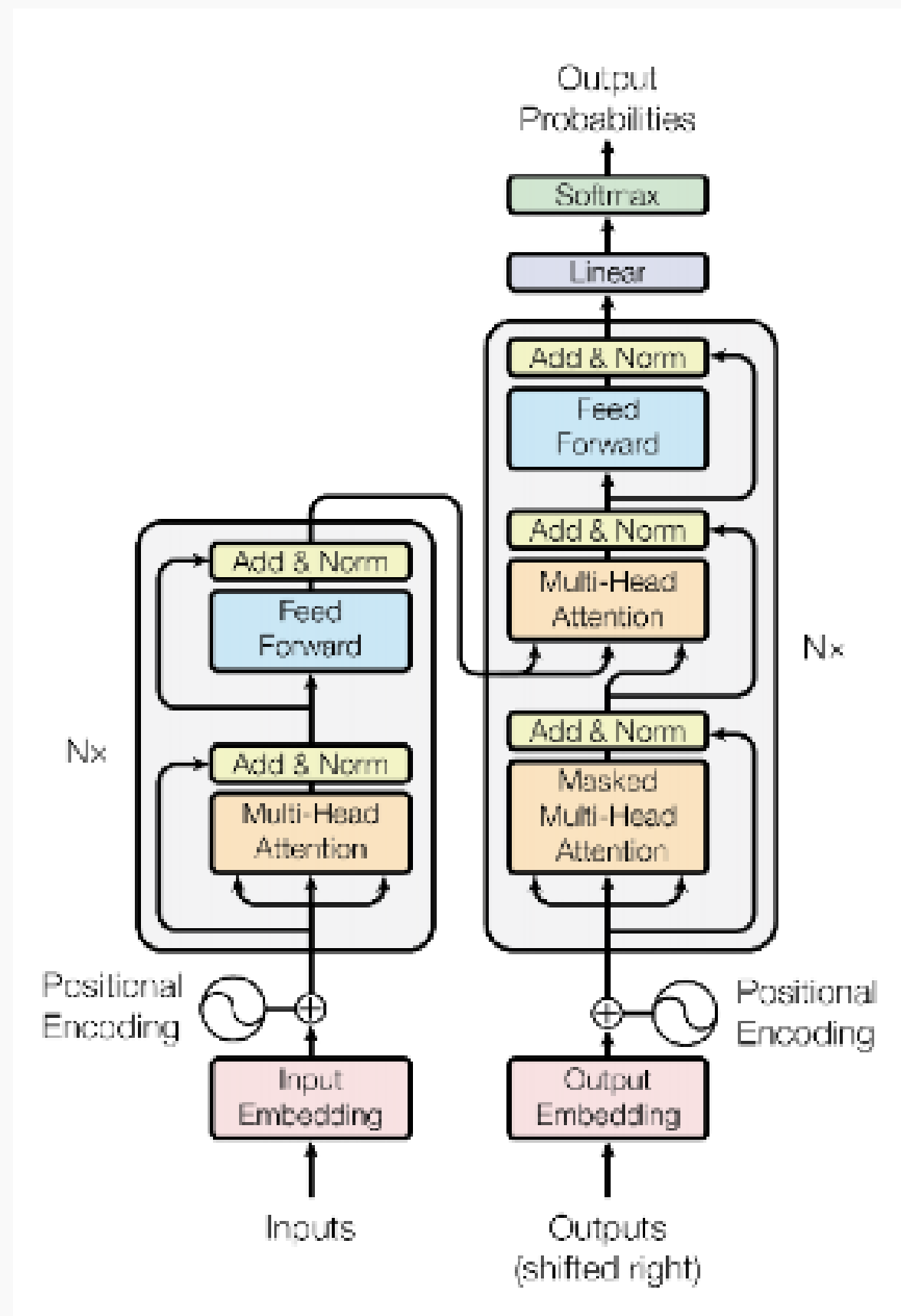
$$\text{softmax} \left( ZQK^T H^T \right) HV = \text{output} \in \mathbb{R}^{T \times d}$$

Self-Attention

## The Transformer

Great results with  
Transformers

Drawbacks and  
variants of  
Transformers



## Transformer Summary

### seq2seq

Self-Attention을 이용하여 문장 내 다른 단어들로부터  
힌트를 받아 현재 단어 Encoding

### Positional Encoding

Positional representation을 추가하여  
새로운 vector 생성

### Encoder-Decoder

- 동일한 개수 사용
- Multi head Attention 사용
- FF network, Residual connection, Layer normalization 사용

Self-Attention

The Transformer

Great results with Transformers

Drawbacks and variants of Transformers

Perplexity

언어 모델을 평가하기 위한 평가 지표

## Transformer Ability

다양한 분야에서 좋은 성능을 보이고 있다.

### Machine translation

기존 SOTA 모델보다 훨씬 높은 성능을 보임

Model	BLEU		Training Cost (FLOPs)	
	EN-DE	EN-FR	EN-DE	EN-FR
ByteNet [18]	23.72			
Deep-Attn + PostNet [24]		39.2		$1.0 \cdot 10^{10}$
GNMT + RL [24]	34.8	39.92	$1.5 \cdot 10^{10}$	$1.4 \cdot 10^{10}$
CoreNLP [9]	25.16	40.46	$9.6 \cdot 10^{10}$	$1.5 \cdot 10^{10}$
MoE [32]	26.03	40.56	$2.0 \cdot 10^{10}$	$1.2 \cdot 10^{10}$
Deep-Attn + PostNet Ensemble [24]		40.4		$8.0 \cdot 10^{10}$
GNMT + RL Ensemble [38]	26.30	41.05	$1.8 \cdot 10^{10}$	$1.1 \cdot 10^{11}$
CoreNLP Ensemble [5]	26.36	41.29	$7.7 \cdot 10^{10}$	$1.3 \cdot 10^{11}$

### Nonlinearities

기존보다 훨씬 낮은 Perplexity = 높은 성능

Model	Test perplexity	ROUGE-L
seq2seq-attention, $L = 500$	5.04852	12.7
Transformer-ED, $L = 500$	2.46645	34.2
Transformer-D, $L = 4000$	2.22216	33.6
Transformer-BMCA, w/ MoE-layer, $L = 11000$	2.05159	36.2
Transformer-BMCA, MoE-128, $L = 11000$	1.92871	37.9
Transformer-BMCA, MoE-256, $L = 7500$	1.90325	38.8

### Future Sequence

현재 1위  
Microsoft  
alexander v-team

Transformer 사용!

Rank	Team	Model	Score
1	Microsoft Alexander V-team	Transformer	38.8

GLUE Benchmark

<https://gluebenchmark.com/leaderboard>

Self-Attention

The Transformer

Great results with  
Transformers

**Drawbacks and  
variants of  
Transformers**

## Drawbacks of Transformer

몇 가지 한계점이 존재한다.

### Quadratic Computation

Sequence length가  
증가함에 따라  
계산량이 2차식으로 증가

- 짧은 문장 = not that big a deal !
- 긴 문장 (Document) =  $T^2$ 이 기하급수적으로 증가하여 computing에 부담됨

$$O(T^2 d)$$

### Position Representation

절대적 위치만 나타내는  
Sinosoidal 사용

- Relative linear position attention : 상대적인 위치
- Dependency syntax-based position : 구조적인 정보 고려
- Rotary Position Embedding : 매 layer마다 위치 정보 고려

[Rotary Position Embedding](#)

<https://arxiv.org/abs/2104.09864>

Self-Attention

The Transformer

Great results with Transformers

**Drawbacks and variants of Transformers**

**Inference time**  
딥러닝 모델이 커지면서,  
빠른 Inference time이  
매우 중요해짐

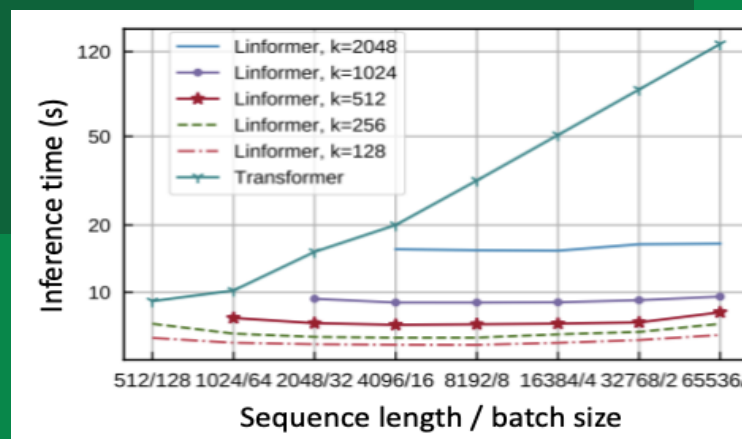
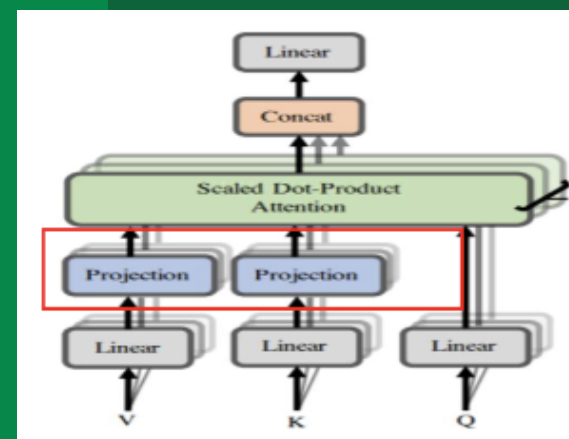
## Variants of Transformer

변형을 통해 한계점을 해결할 수 있다.

### Linformer

Sequence length의 차원을  
낮춰 계산량을 줄임

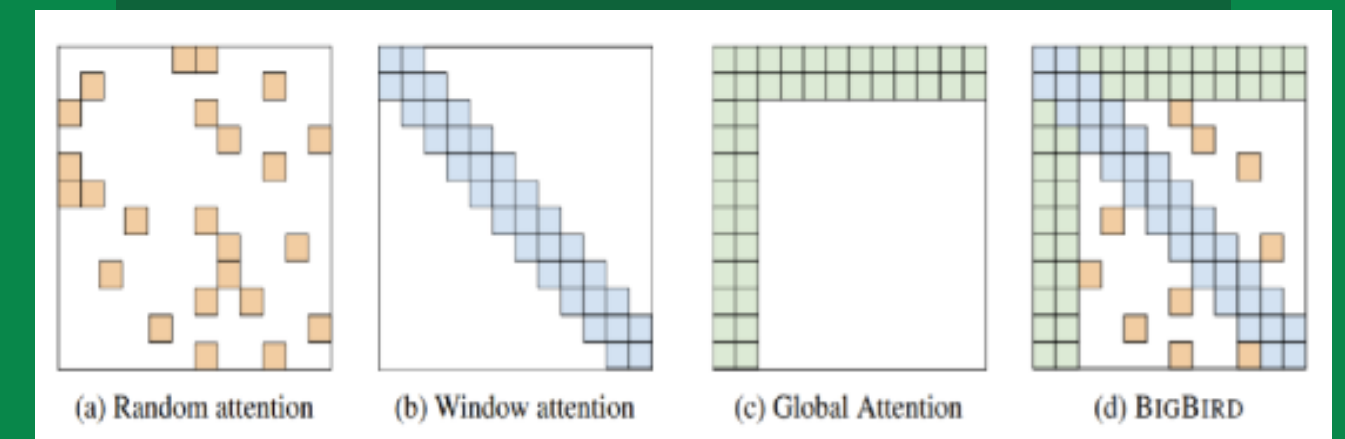
- Projection을 통해 Key, Value의 Sequence length dimension을 낮춤
- Inference time을 보면 높음



### BigBird

최적의 조합만 계산함

- 모든 pair의 Attention을 계산하지 않음
- 최적의 조합만 계산하여, 계산량 줄임



Self-Attention

The Transformer

Great results with  
Transformers

**Drawbacks and  
variants of  
Transformers**

## Variants of Transformer

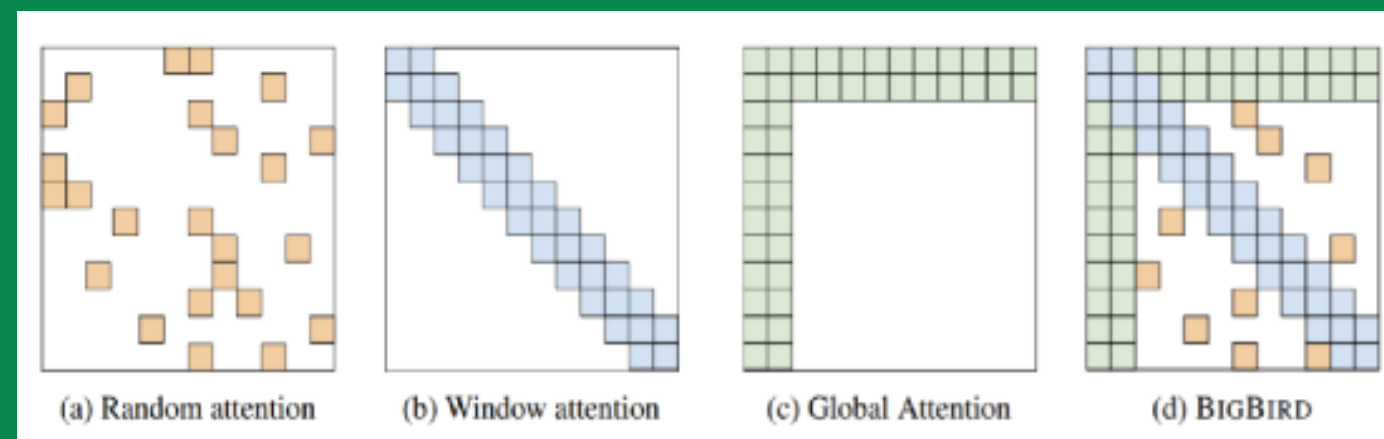
변형을 통해 한계점을 해결할 수 있다.

Question ?

BigBird

어텐션에 대한 그림으로 4종류가 있는데,  
이것들이 의미하는 바가 무엇인지  
궁금합니다.

- (a) **Random attention** : Query와  $r$ 개의 random keys 간의 attention
- (b) **Window attention** : Query 양옆  $w$ 개의 keys와의 attention
- (c) **Global attention** : Query와  $g$ 개의 global token들과의 attention
- (d) **BigBird** : 위 3개를 모두 합친 방식 (가장 성능이 좋게 나옴)



감사합니다.