₊˙* 2021 SMARCLE Paper Review ´˙*
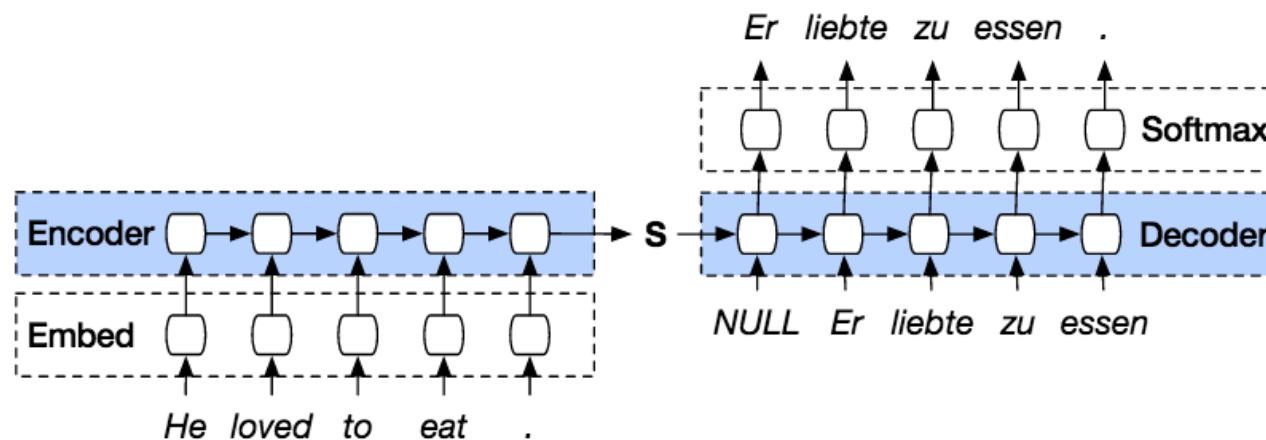
# Transformer - Attention Is All You Need
## (NIPS 2017)

지능기전공학부 무인이동체공학전공
송혜원

SMARCLE

# Seq2Seq의 문제점



계산량 ↑ Bottleneck

RNN (1986)  LSTM (1997)  Seq2Seq (NIPS 2014)  Attention (ICLR 2015)  Transformer (NIPS 2017)  GPT-1 (2018)  BERT (NAACL 2019)  GPT-3 (2020)

고정된 크기의 context vector 사용

입력 시퀀스 전체에서 정보를 추출하는 방향으로 발전

Er liebte zu essen .

Softmax

Encoder ☐→☐→☐→☐→☐ → S → ☐→☐→☐→☐→☐ Decoder

Embed

He loved to eat .
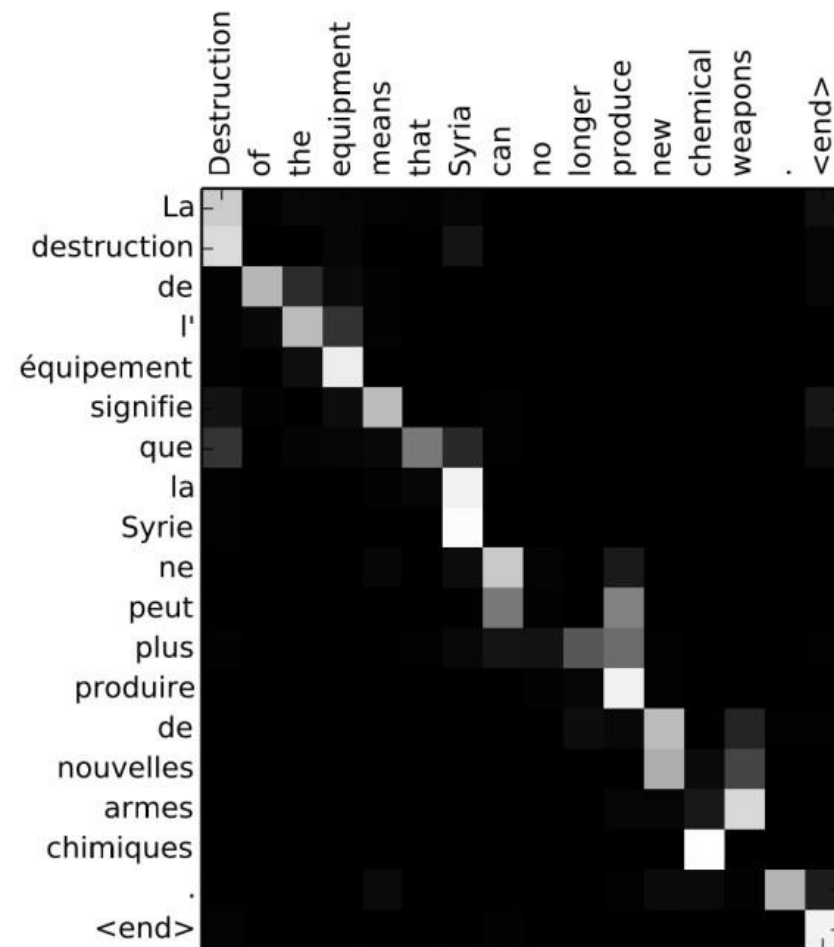
NULL Er liebte zu essen

# Attention

이전엔 입력을
하나의 **context vector**로 압축

⬇

Attention 기법은
문장 자체를 입력으로!

⬇

어떤 단어가 어떤 단어와
가장 연관성이 있는지 알 수 있게 됨

# Transformer_Input



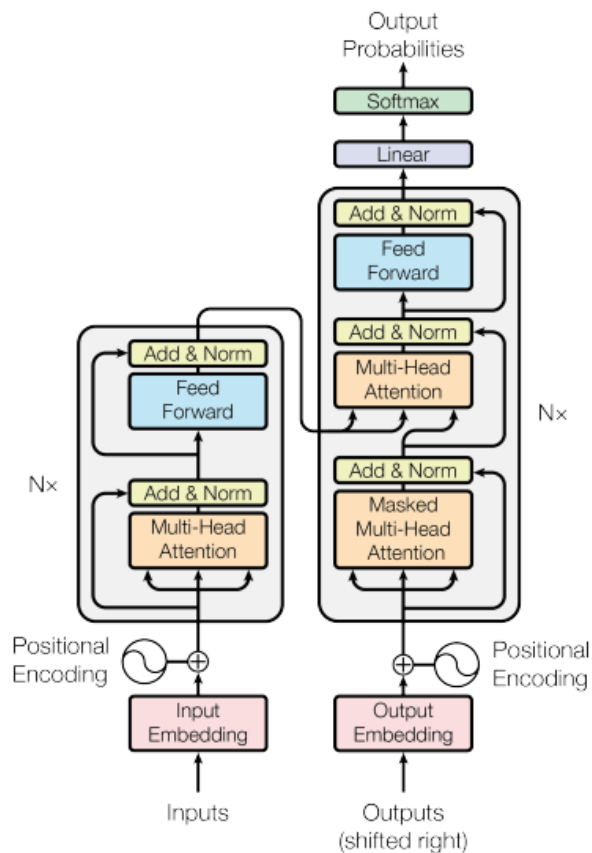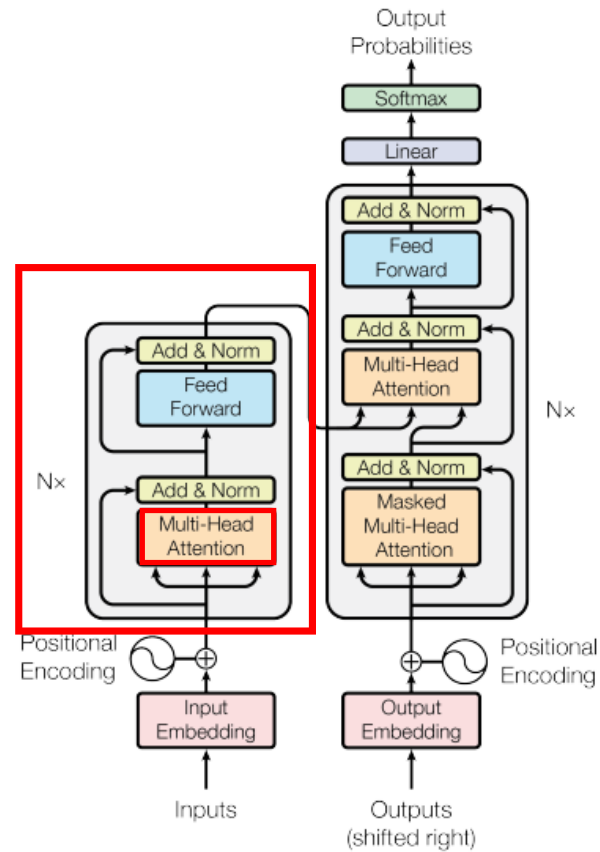Figure 1: The Transformer - model architecture.

Dog loves potato

| Dog | | | | ... | | | |
|---|---|---|---|---|---|---|---|
| loves | | | | ... | | | |
| potato | | | | ... | | | |

512

문장 전체를 통째로 넣다보니 단어의 위치 정보를 알 수 없음

➡️ **Positional Encoding** 통해 위치 정보를 따로 더해줌!

# Transformer _Encoder Self-Attention



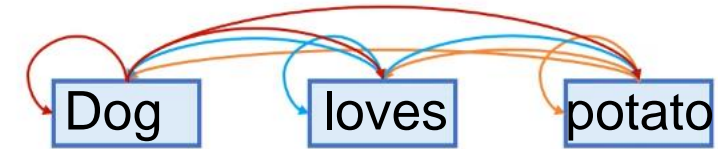Figure 1: The Transformer - model architecture.

Encoder Self-Attention:

Dog    loves    potato
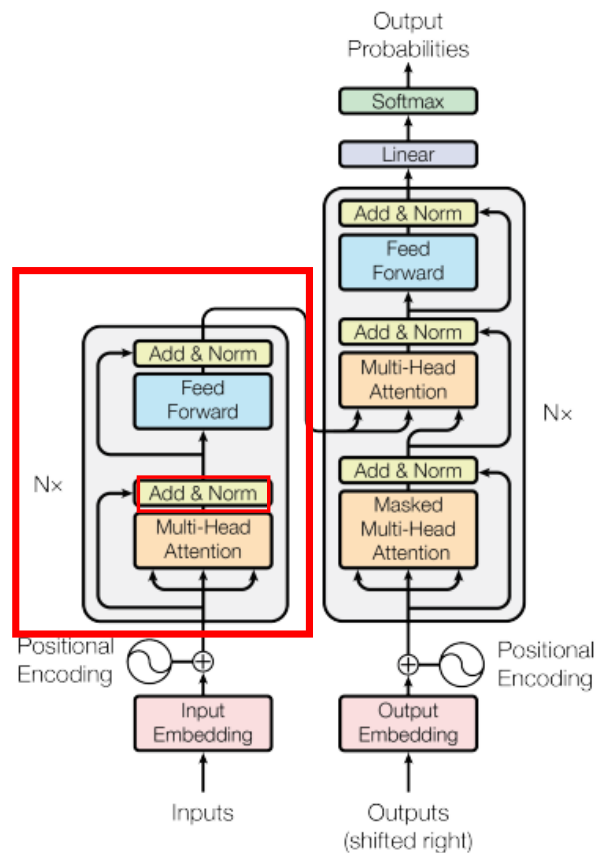
# Transformer_Encoder



Figure 1: The Transformer - model architecture.

Encoder Self-Attention:

Dog    loves    potatoc

## Residual learning

➡ 특정 layer를 건너뜀

Gradient vanishing 완화

Global optima 잘 찾을 수 있음
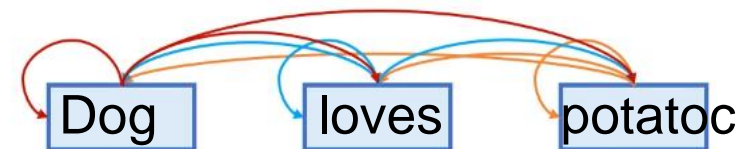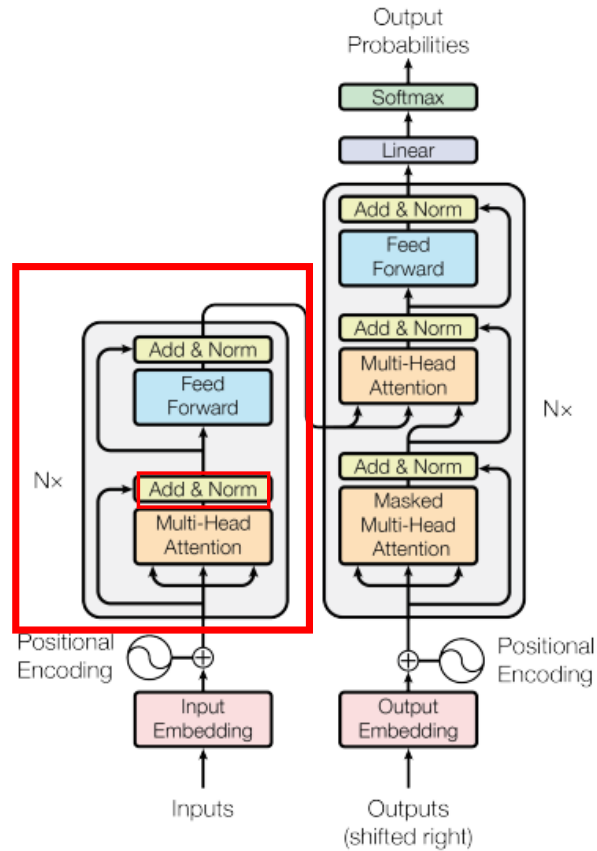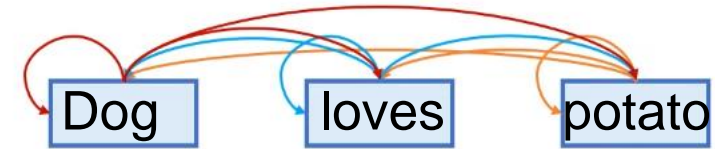
# Transformer _Encoder

Encoder Self-Attention:

Dog → loves → potato

Figure 1: The Transformer - model architecture.

## Residual learning

➡ 특정 layer를 건너뜀

Gradient vanishing 완화

Global optima 잘 찾을 수 있음

인풋의 다양한 특징 attention 위해 n개의 layer 병렬로 쌓음

# Transformer _Masked Decoder Self-Attention



Figure 1: The Transformer - model architecture.
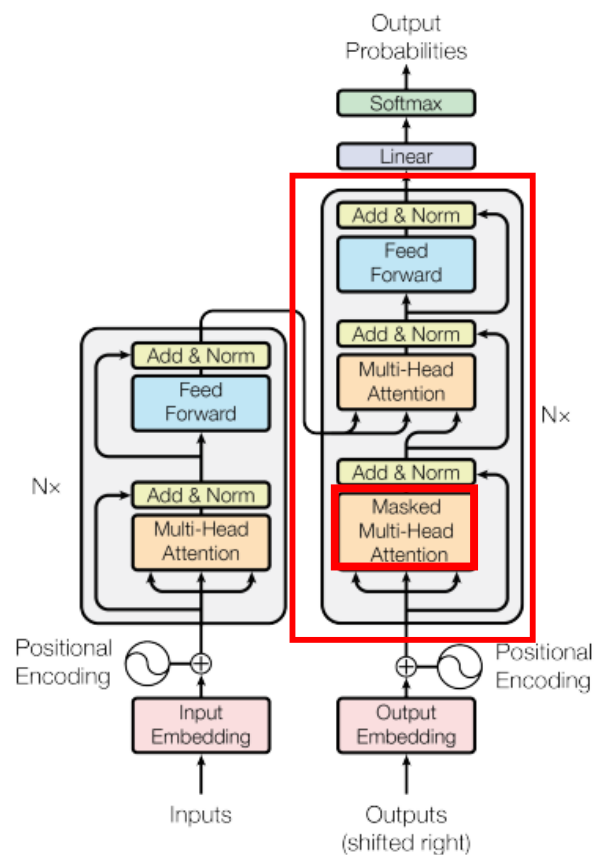
Masked Decoder Self-Attention: 강아지는    감자를    좋아해

# Transformer_Masked Decoder Self-Attention



Figure 1: The Transformer - model architecture.

Masked Decoder Self-Attention: 강아지는   감자를   좋아해

Encoder Self-Attention: Dog   loves   potato

# Transformer_Masked Decoder Self-Attention

Output
Probabilities

Softmax

Linear

Add & Norm

Feed
Forward

Add & Norm

Multi-Head
Attention

N×

Add & Norm

Feed
Forward

N×

Add & Norm

Multi-Head
Attention

Add & Norm

Masked
Multi-Head
Attention

Positional
Encoding

Positional
Encoding

Input
Embedding

Output
Embedding

Inputs

Outputs
(shifted right)

Figure 1: The Transformer - model architecture.

Encoder-Decoder Attention:

| Dog | loves | potato | 강아지 |

· 인코더 파트에서 나온 출력 결과를 디코더마다 적용하여 전적으로 활용

· 시퀀스가 끝날 때까지 반복

# Transformer_Multi-Head Attention?

Q(Query) : 물어보는 주체    K(Key) : 대상     V(Value) : 가중치
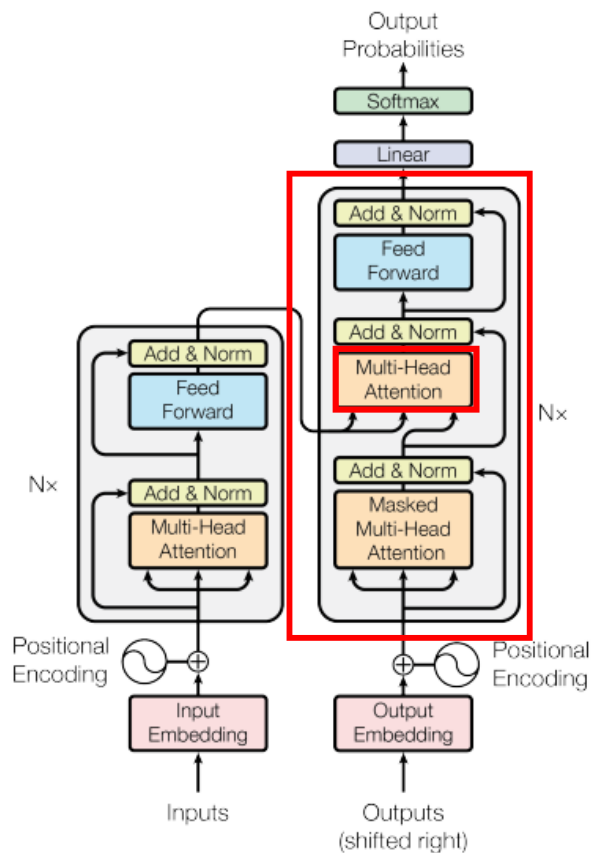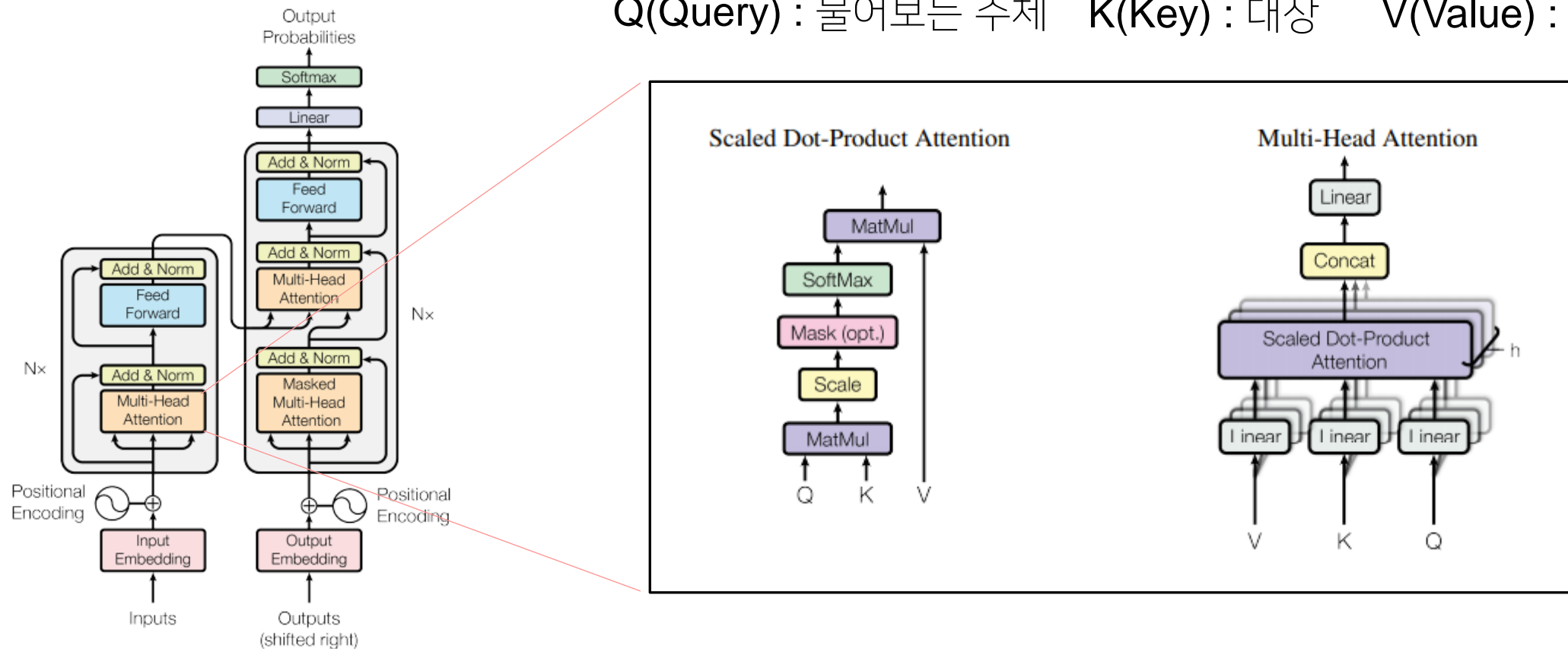


Figure 1: The Transformer - model architecture.

# Transformer _Multi-Head Attention?

SMARCLE

Q(Query) : 물어보는 주체    K(Key) : 대상      V(Value) : 가중치



Figure 1: The Transformer - model architecture.

Scaled Dot-Product Attention

Multi-Head Attention

Encoder Self-Attention:

Dog    loves    potato

Transformer - Attention Is All You Need

# Transformer_Multi-Head Attention?

Q(Query) : 물어보는 주체    K(Key) : 대상    V(Value) : 가중치



Figure 1: The Transformer - model architecture.

Scaled Dot-Product Attention

Multi-Head Attention

potato loves

Encoder Self-Attention:

Dog    loves    potato

# Transformer _Multi-Head Attention?

Q(Query) : 물어보는 주체    K(Key) : 대상     V(Value) : 가중치



Figure 1: The Transformer - model architecture.

Scaled Dot-Product Attention

Multi-Head Attention

potato loves

야 loves! 나 너랑 얼마나 관련있어?

Encoder Self-Attention:

Dog    loves    potato

Transformer - Attention Is All You Need

# Transformer _Multi-Head Attention?

Q(Query) : 물어보는 주체    K(Key) : 대상    V(Value) : 가중치



Figure 1: The Transformer - model architecture.

Scaled Dot-Product Attention

Multi-Head Attention

Encoder-Decoder Attention:    Dog    loves    potato    강아지

# Transformer _Multi-Head Attention?

Q(Query) : 물어보는 주체   K(Key) : 대상     V(Value) : 가중치



Figure 1: The Transformer - model architecture.

강아지 loves

야 loves! 나 너랑 얼마나 관련있어?

Encoder-Decoder Attention:

Dog    loves    potato    강아지

# Transformer _Multi-Head Attention?


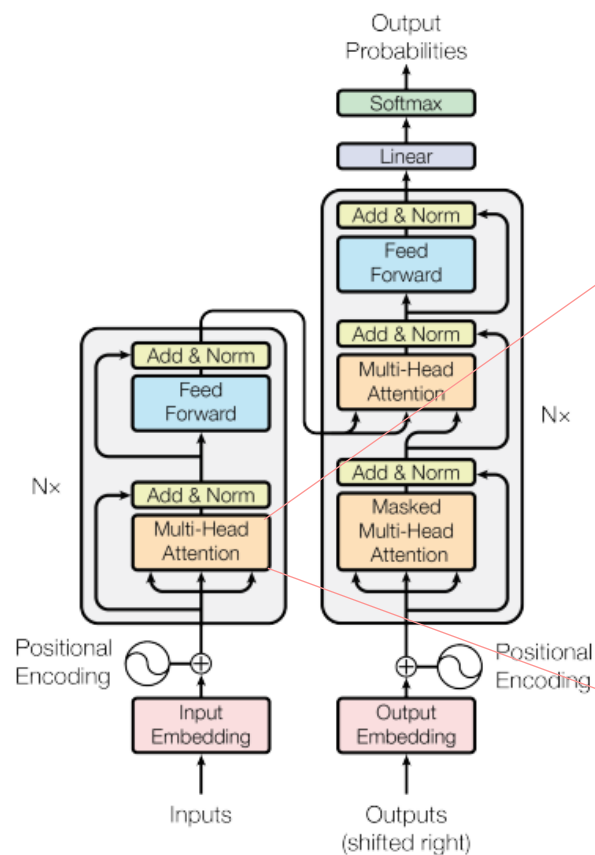
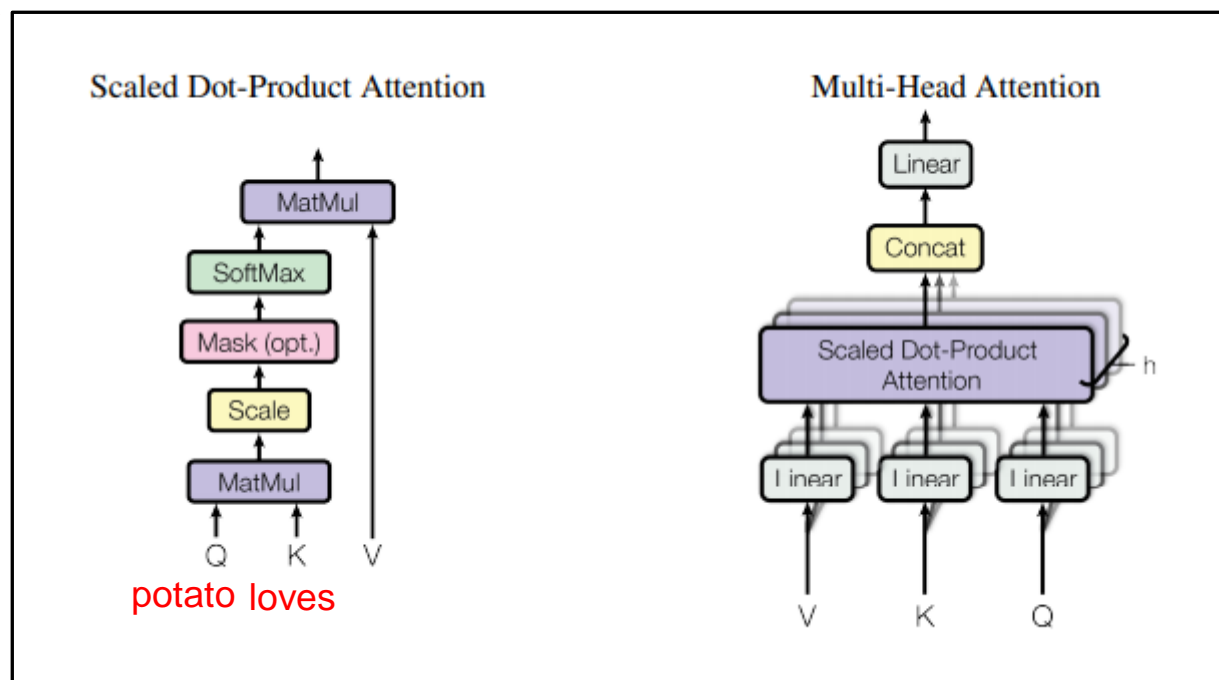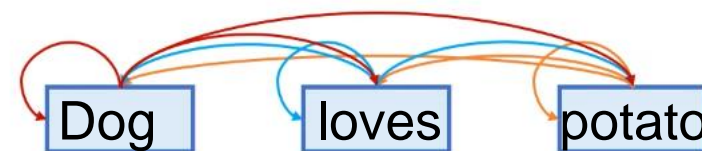Q(Query) : 물어보는 주체    K(Key) : 대상      V(Value) : 가중치

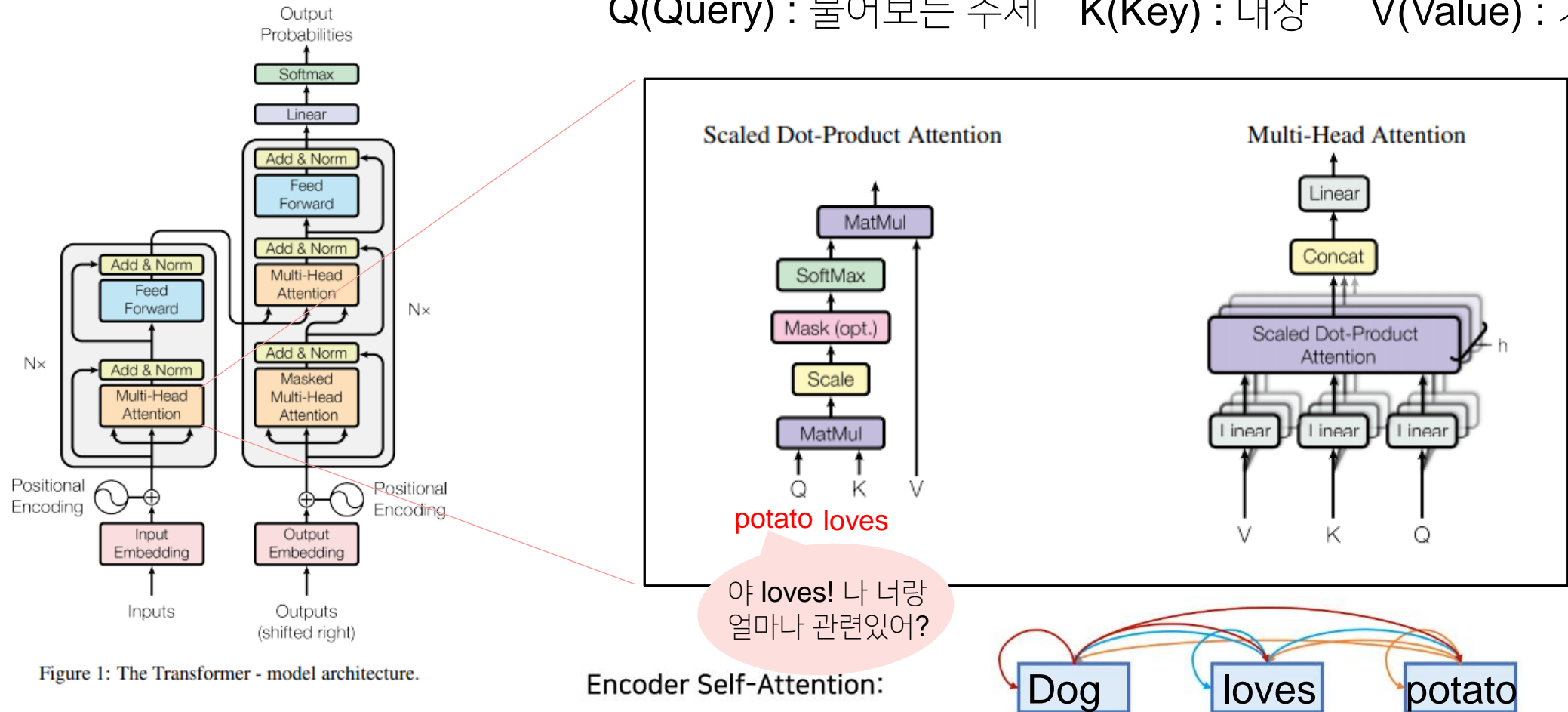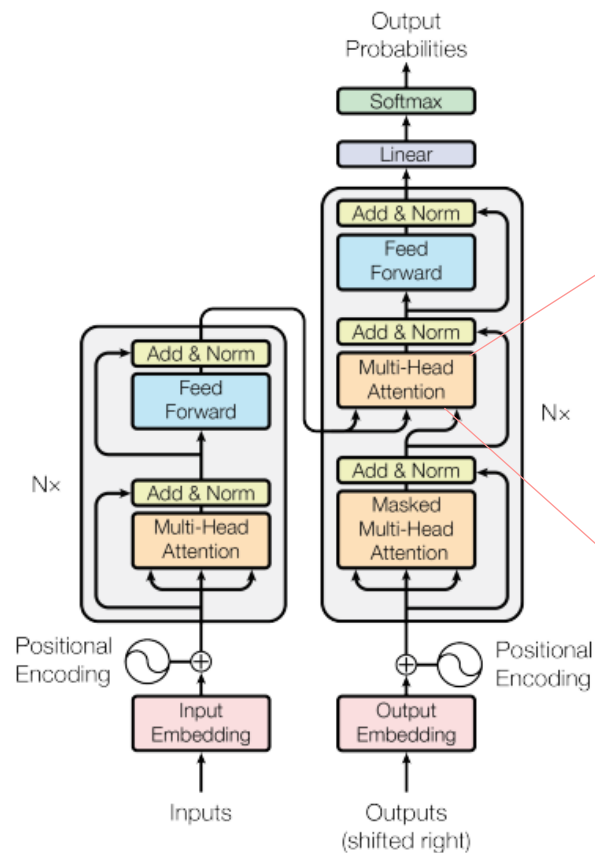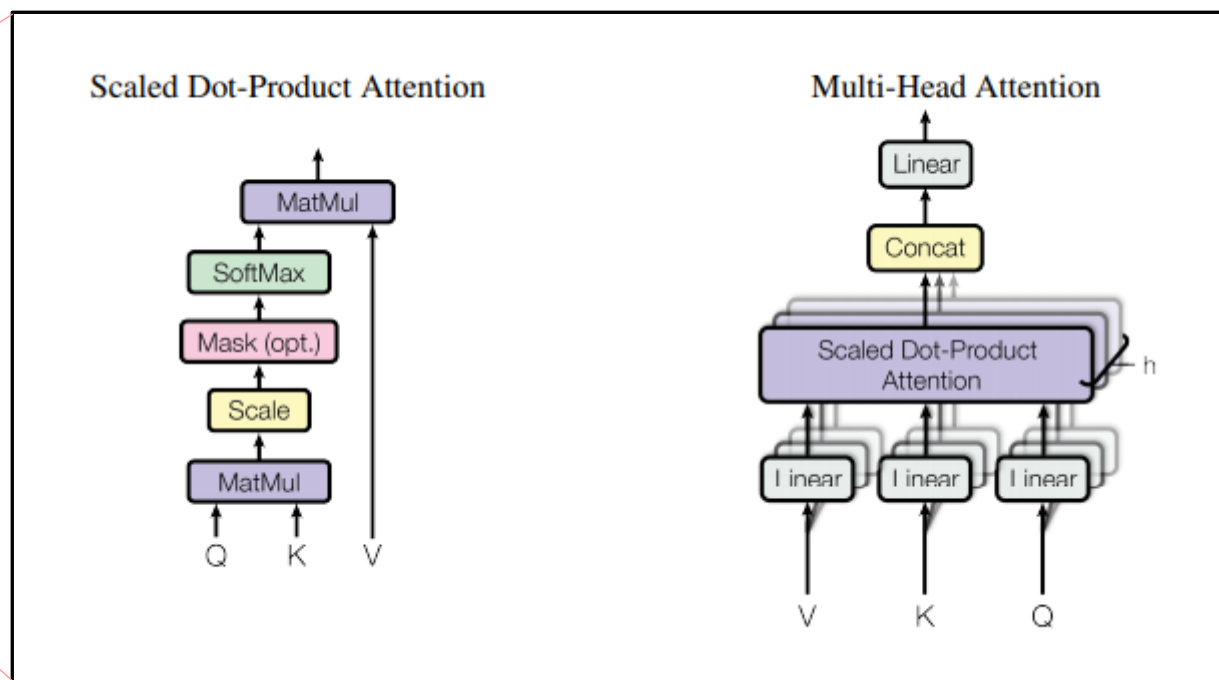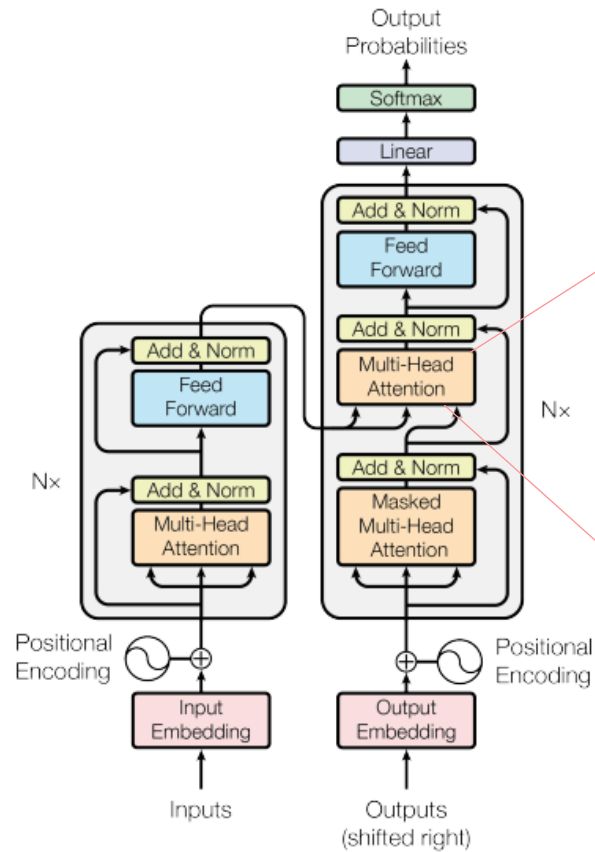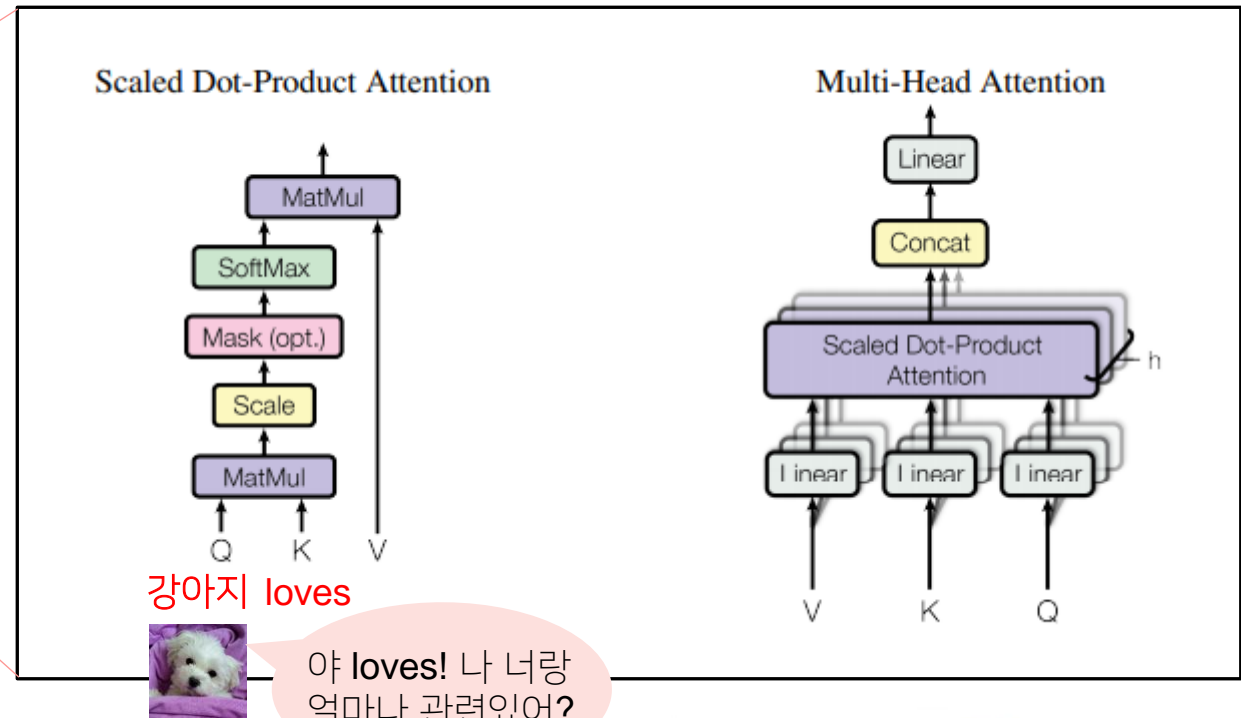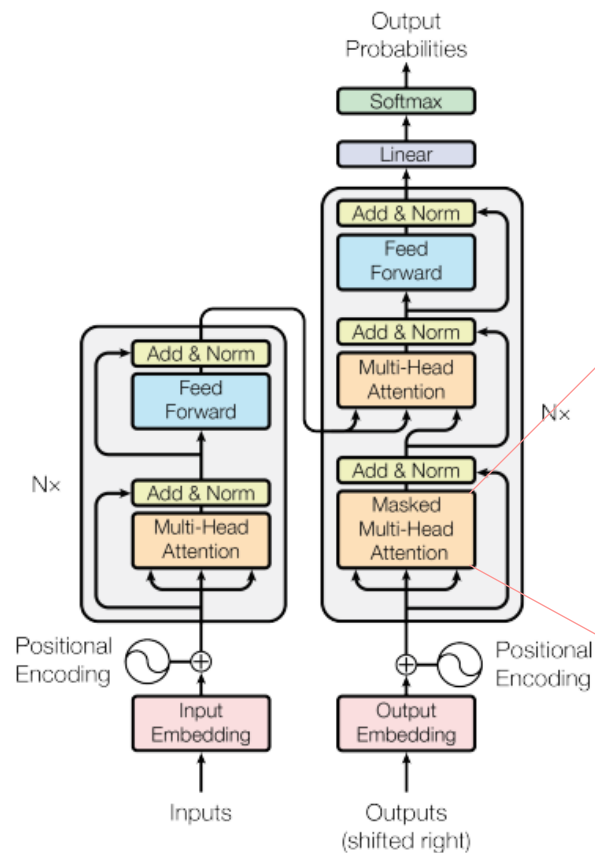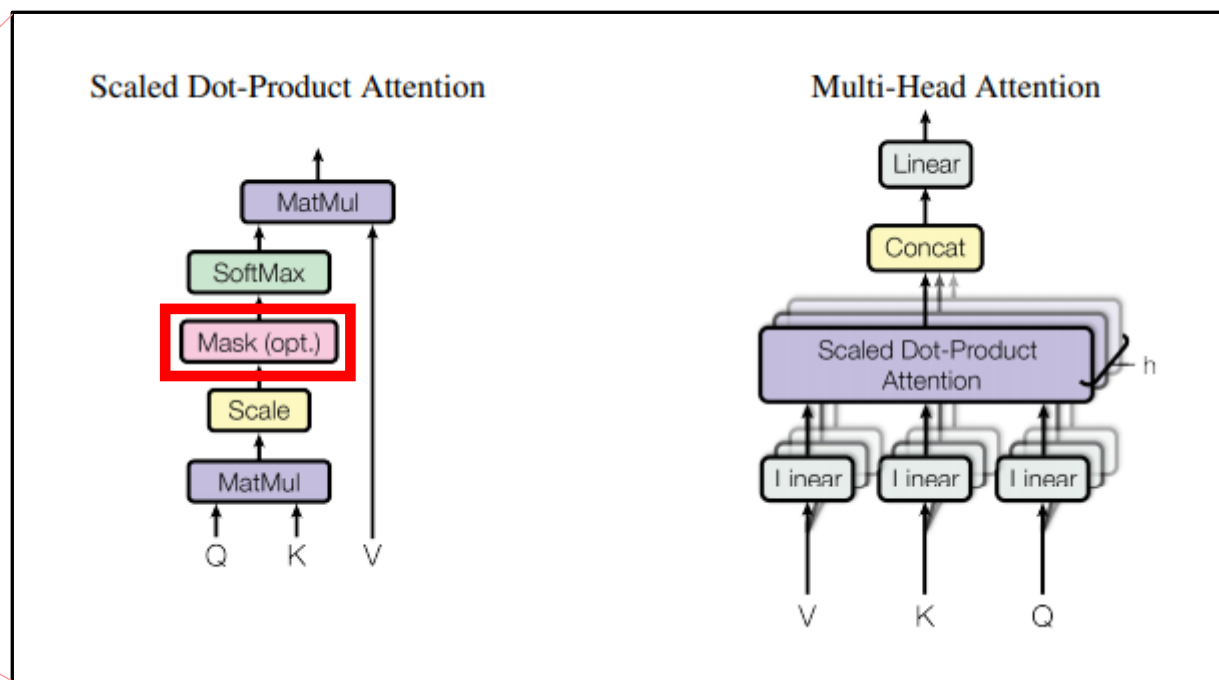Figure 1: The Transformer - model architecture.

Masked Decoder Self-Attention:  강아지는    감자를    좋아해

# Transformer _Multi-Head Attention?

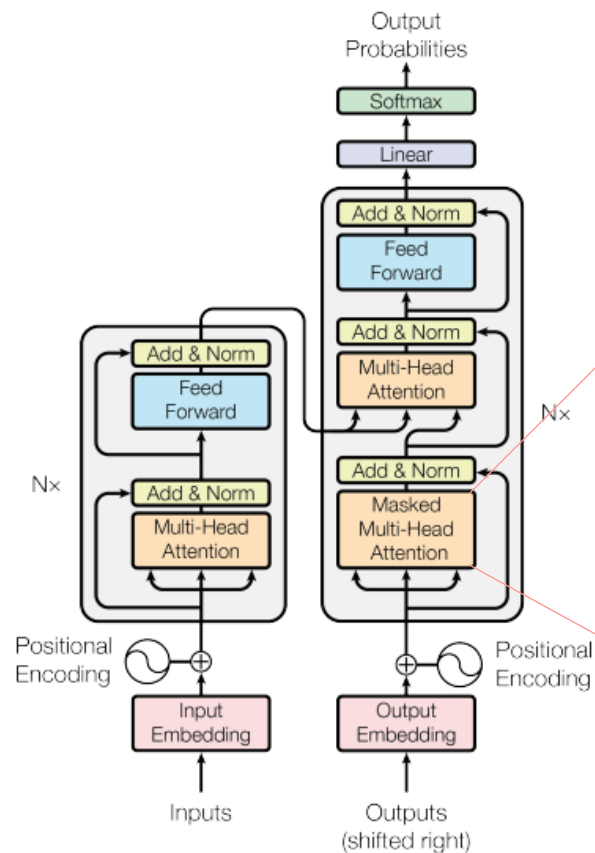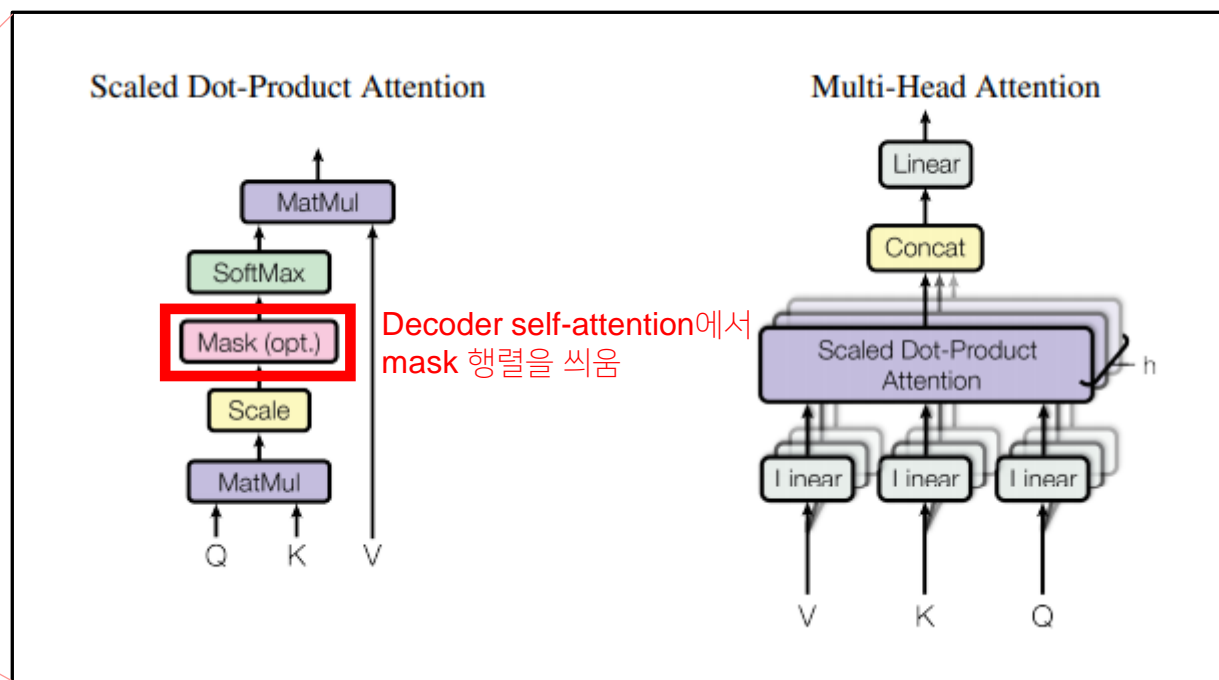Q(Query) : 물어보는 주체   K(Key) : 대상   V(Value) : 가중치



Figure 1: The Transformer - model architecture.

Decoder self-attention에서 mask 행렬을 씌움

Masked Decoder Self-Attention: 강아지는   감자를   좋아해

# Transformer _Positional Encoding



Output
Probabilities

Softmax

Linear

Add & Norm

Feed
Forward

Add & Norm

Multi-Head
Attention

Nx

Add & Norm

Masked
Multi-Head
Attention

Add & Norm

Feed
Forward

Nx

Add & Norm

Multi-Head
Attention

Positional
Encoding

Input
Embedding

Inputs

Positional
Encoding

Output
Embedding

Outputs
(shifted right)

Figure 1: The Transformer - model architecture.

In this work, we use sine and cosine functions of different frequencies:

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

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

where $pos$ is the position and $i$ is the dimension. That is, each dimension of the positional encoding corresponds to a sinusoid. The wavelengths form a geometric progression from $2\pi$ to $10000 \cdot 2\pi$. We chose this function because we hypothesized it would allow the model to easily learn to attend by relative positions, since for any fixed offset $k$, $PE_{pos+k}$ can be represented as a linear function of $PE_{pos}$.

We also experimented with using learned positional embeddings [9] instead, and found that the two versions produced nearly identical results (see Table 3 row (E)). We chose the sinusoidal version because it may allow the model to extrapolate to sequence lengths longer than the ones encountered during training.

- 본문에서는 주기함수를 이용해서 단어의 위치를 정해줌
- 이후 모델들은 학습이 가능한 임베딩 레이어를 넣어주는 추세

# Transformer _Conclusion

Table 3: Variations on the Transformer architecture. Unlisted values are identical to those of the base model. All metrics are on the English-to-German translation development set, newstest2013. Listed perplexities are per-wordpiece, according to our byte-pair encoding, and should not be compared to per-word perplexities.

| | $N$ | $d_{model}$ | $d_{ff}$ | $h$ | $d_k$ | $d_v$ | $P_{drop}$ | $\epsilon_{ls}$ | train steps | PPL (dev) | BLEU (dev) | params $\times 10^6$ |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| base | 6 | 512 | 2048 | 8 | 64 | 64 | 0.1 | 0.1 | 100K | 4.92 | 25.8 | 65 |
| (A) | | | | 1 | 512 | 512 | | | | 5.29 | 24.9 | |
| | | | | 4 | 128 | 128 | | | | 5.00 | 25.5 | |
| | | | | 16 | 32 | 32 | | | | 4.91 | 25.8 | |
| | | | | 32 | 16 | 16 | | | | 5.01 | 25.4 | |
| (B) | | | | | 16 | | | | | 5.16 | 25.1 | 58 |
| | | | | | 32 | | | | | 5.01 | 25.4 | 60 |
| (C) | 2 | | | | | | | | | 6.11 | 23.7 | 36 |
| | 4 | | | | | | | | | 5.19 | 25.3 | 50 |
| | 8 | | | | | | | | | 4.88 | 25.5 | 80 |
| | | 256 | | | 32 | 32 | | | | 5.75 | 24.5 | 28 |
| | | 1024 | | | 128 | 128 | | | | 4.66 | 26.0 | 168 |
| | | | 1024 | | | | | | | 5.12 | 25.4 | 53 |
| | | | 4096 | | | | | | | 4.75 | 26.2 | 90 |
| (D) | | | | | | | 0.0 | | | 5.77 | 24.6 | |
| | | | | | | | 0.2 | | | 4.95 | 25.5 | |
| | | | | | | | | 0.0 | | 4.67 | 25.3 | |
| | | | | | | | | 0.2 | | 5.47 | 25.7 | |
| (E) | | positional embedding instead of sinusoids | | | | | | | | 4.92 | 25.7 | |
| big | 6 | 1024 | 4096 | 16 | | | 0.3 | | 300K | 4.33 | 26.4 | 213 |

기본

key dimension (B)

model 크기 (C)

dropout (D)

p.e 성능 (E)

↓ 숫자 bad

↑ 숫자 good

good

base랑 차이 거의 x

# Transformer_TBC..



고정된 크기의 context vector 사용

입력 시퀀스 전체에서 정보를 추출하는 방향으로 발전

# 감사합니다 :-3

SMARCLE