

Transformer (Attention is All You Need)

자연어처리 텍스트마이닝

Transformer

(Attention is All You Need)

Transformer

Transformer 개요 (1)

- Transformer의 가장 큰 특징은 Convolution도, Recurrence도 사용하지 않음
- Since our model contains **no recurrence and no convolution**, in order for the model to make use of the order of the sequence, we must inject some information about the relative or absolute position of the tokens in the sequence. (Vaswani et al., Attention Is All You Need, 2017)



Transformer 개요 (2)

- **Long-term dependency problem**

어떤 정보와 다른 정보 사이의 거리가 멀 때 해당 정보를 이용하지 못하는 것 (RNN의 문제점)

=> Attention mechanism 으로 해결

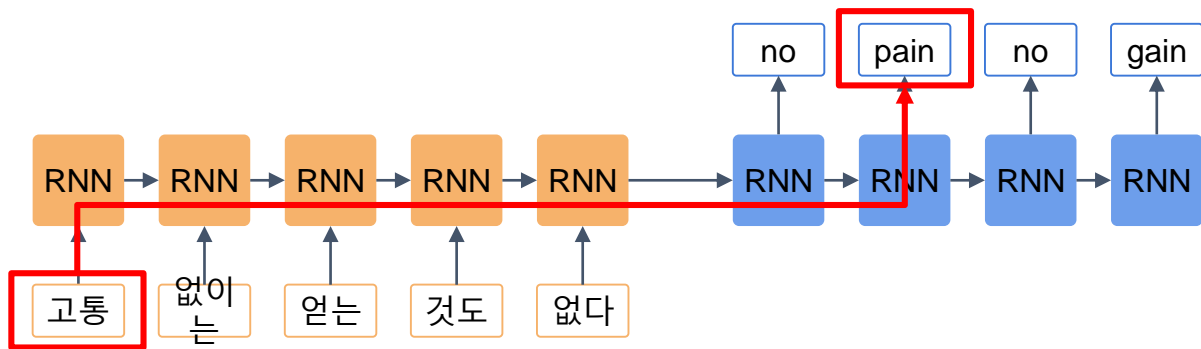
- **Parallelization**

RNN은 이전 hidden state를 사요하으써 순차적으로 계산이 되어야함 (병렬화 불가능)

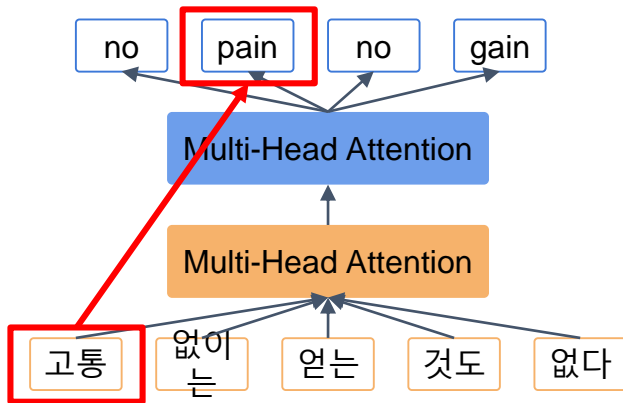


Transformer 개요 (3)

Seq2Seq



트랜스포머

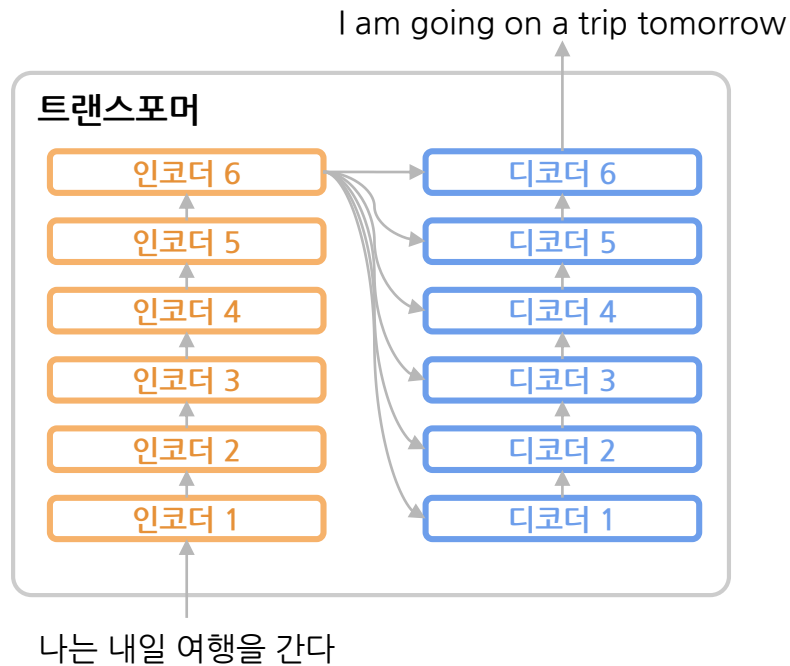
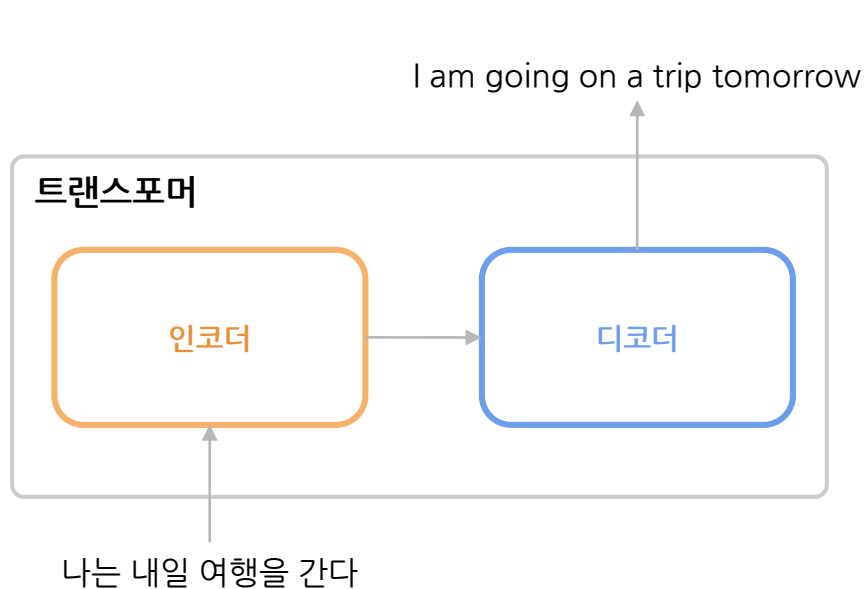


RNN을 사용하지 않고 행렬 병렬연산으로 빠르게 학습이 가능

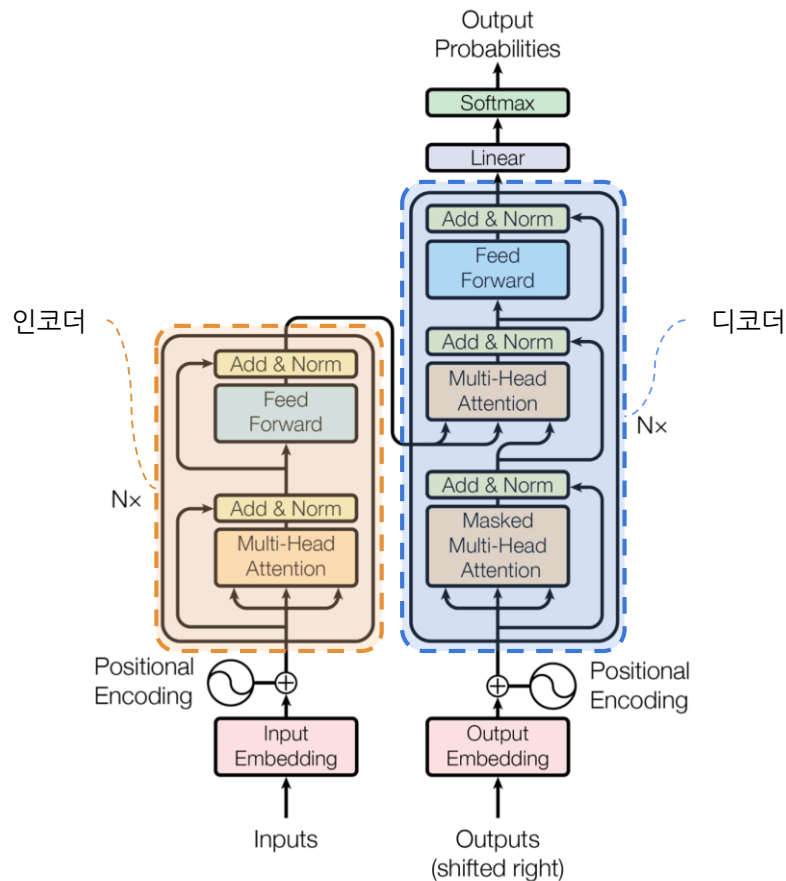
Transformer 개요 (4)



Transformer 개요 (4)



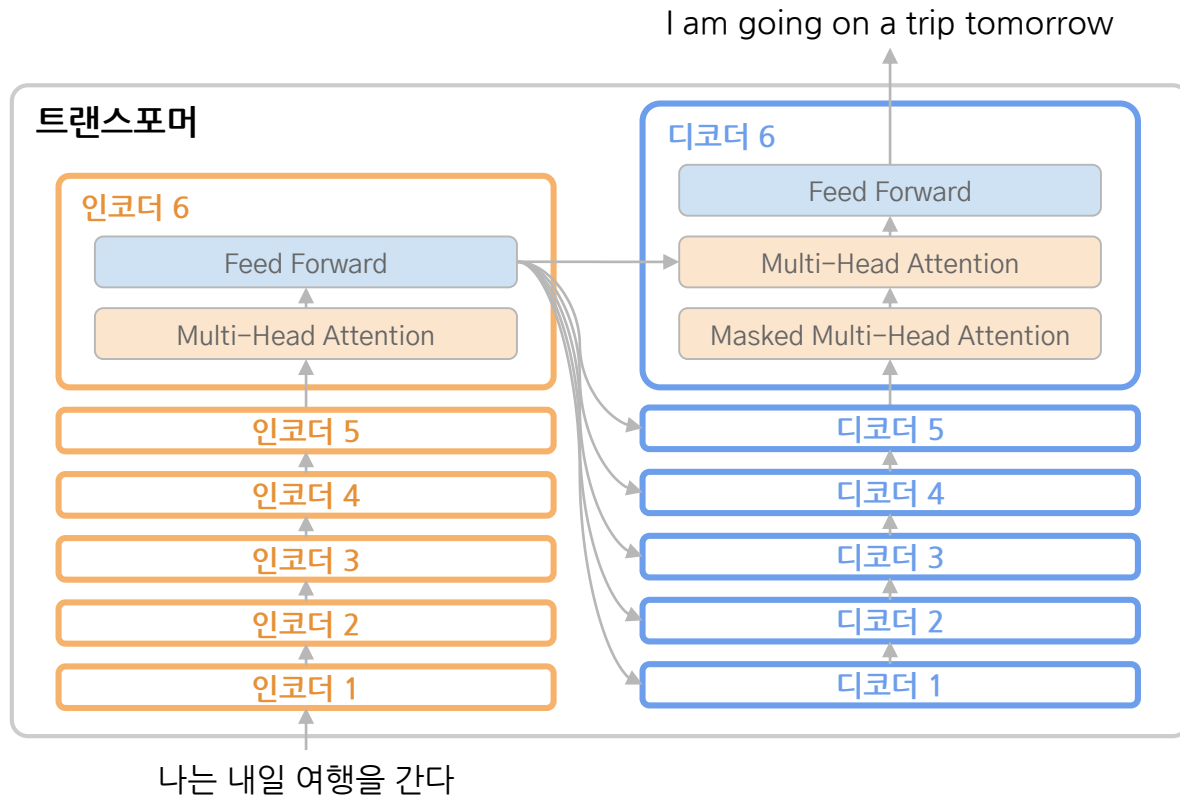
Transformer 개요 (5)



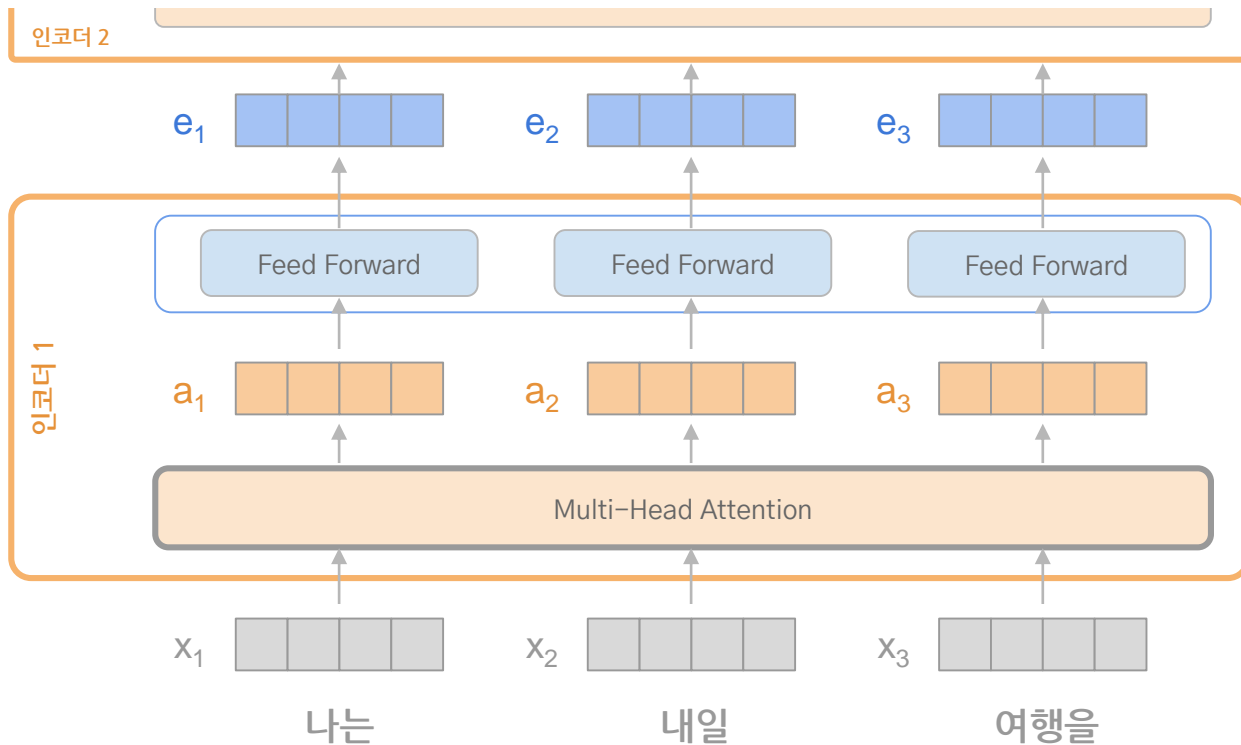
Encoder

Transformer

Transformer 구조



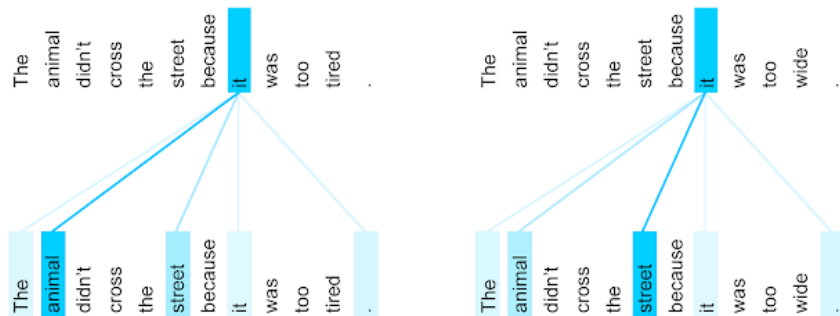
Transformer 인코더 – Multi-Head Attention



Self Attention (1) – Scaled Dot-Product Attention

The animal didn't cross the street because it was too tired.
 ⇒ 동물은 길을 건너지 않았다. 왜냐하면 그것(it)은 너무 피곤하기 때문이다.

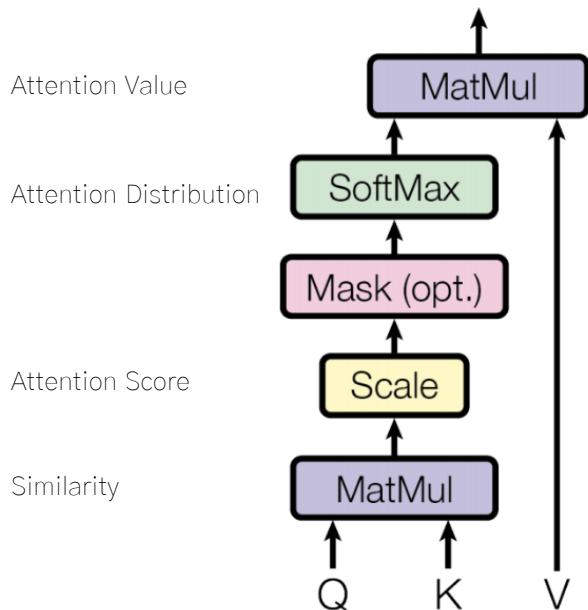
The animal didn't cross the street because it was too wide.
 ⇒ 동물은 길을 건너지 않았다. 왜냐하면 그것(it)이 너무 넓기 때문이다.



<https://ai.googleblog.com/2017/08/transformer-novel-neural-network.htm>

Self Attention (1) – Scaled Dot-Product Attention

Scaled Dot-Product Attention



- 연산 Dependency 가 줄어 빠른 연산 가능
- 병렬화 가능 연산 증가
- long-range의 term들의 dependency도 학습가능
- QK^T : Q(query)와 K(key)의 유사도를 의미
- $\text{sqrt}(d_k)$: K(key)의 차원수로 나누어 scaling

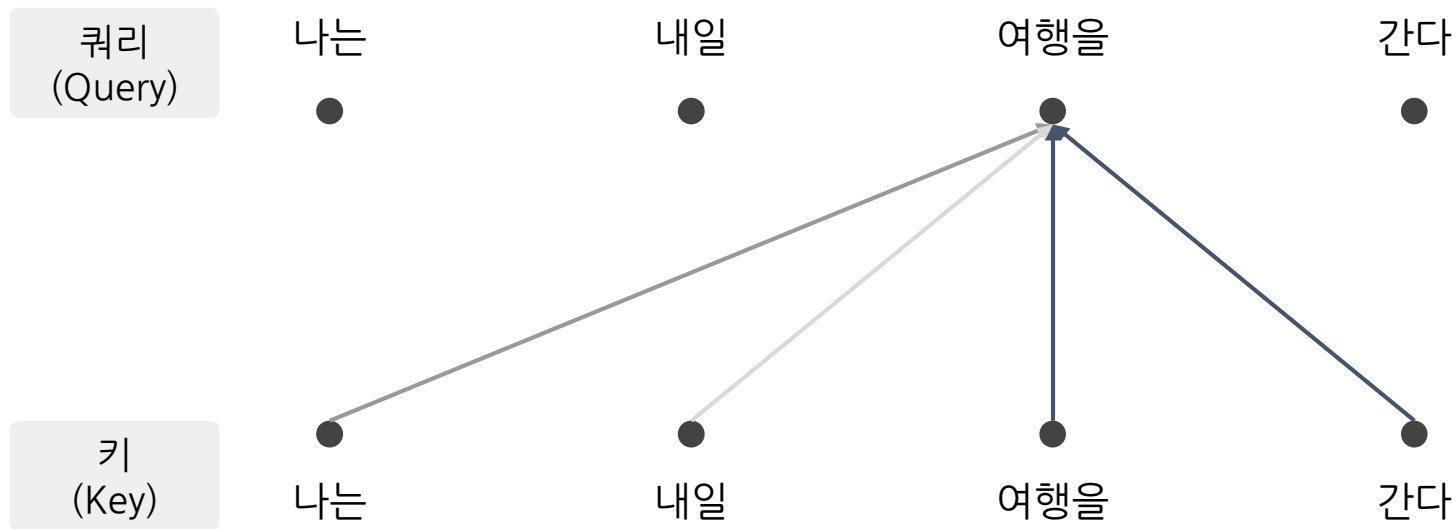
$$Attention(Q, K, V) = \underset{\text{어텐션 분포}}{\text{softmax}_k} \left(\overset{\text{단어간 유사도}}{\frac{QK^T}{\sqrt{d_k}}} \right) \underset{\text{어텐션 스코어}}{V}$$

벡터의 내적과 코사인 유사도

$$A \cdot B = \|A\| \|B\| \cos \theta$$

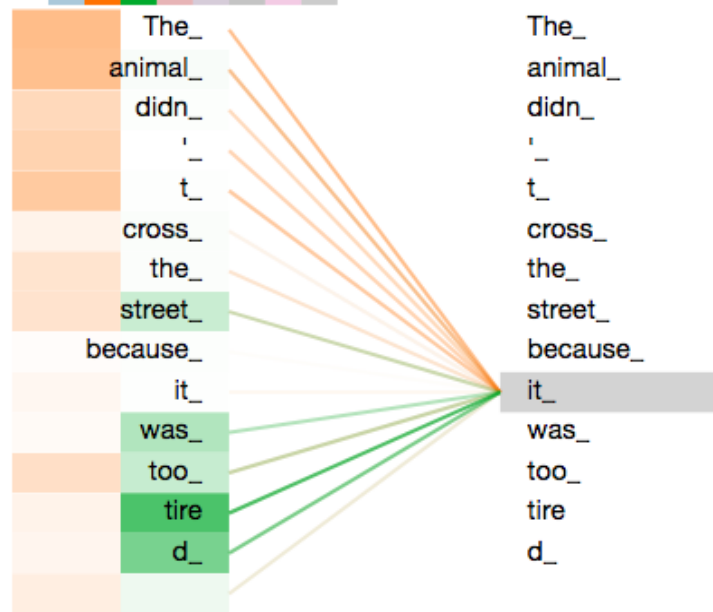
$$\textit{similarity} = \cos(\theta) = \frac{A \cdot B}{\|A\| \cdot \|B\|} = \frac{\sum_{i=1}^n A_i B_i}{\sqrt{\sum_{i=1}^n A_i^2} \sqrt{\sum_{i=1}^n B_i^2}}$$

Self Attention (2) - 예제

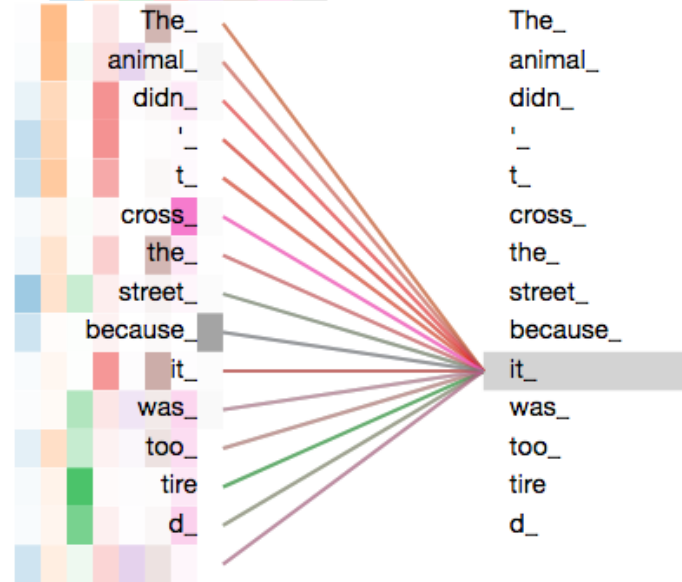


Self Attention (3)

Layer: 5 Attention: Input - Input

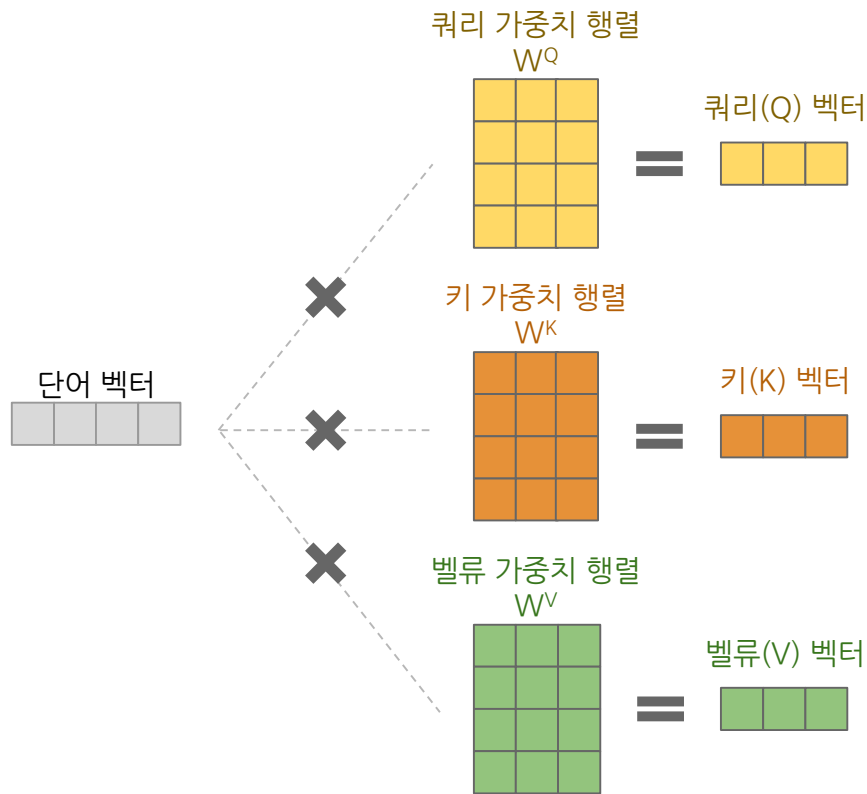


Layer: 5 Attention: Input - Input

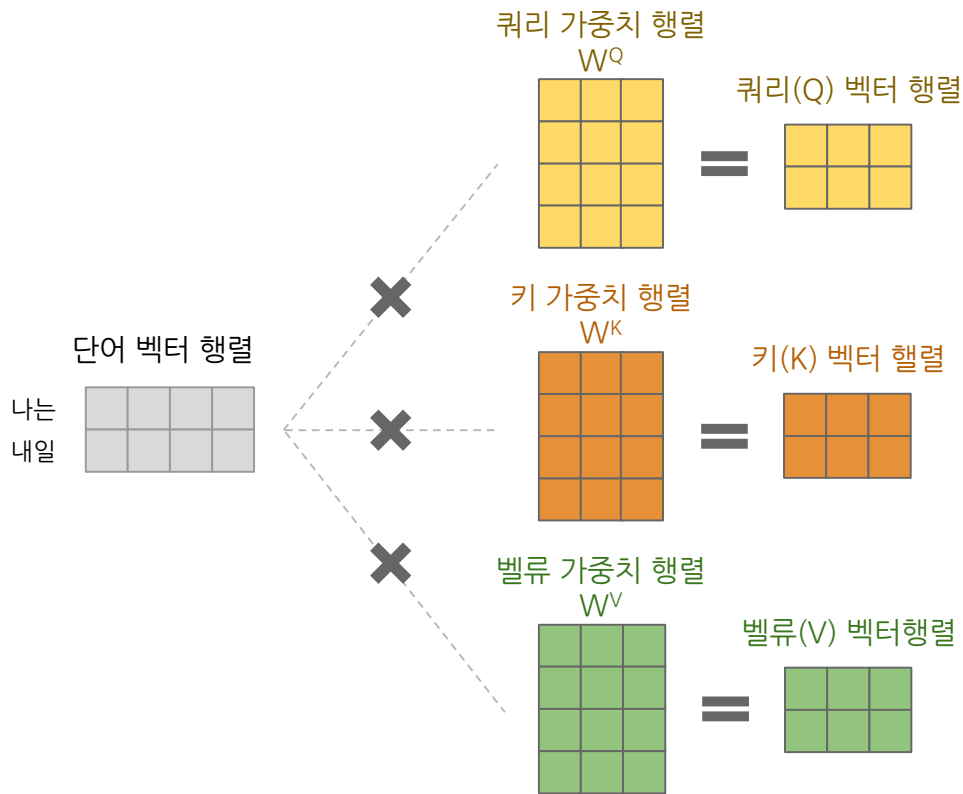


https://colab.research.google.com/github/tensorflow/tensor2tensor/blob/master/tensor2tensor/notebooks/hello_t2t.ipynb#scrollTo=QJKU36QAfqQC

Scaled Dot-Product Attention (1) – Query, Key, Value

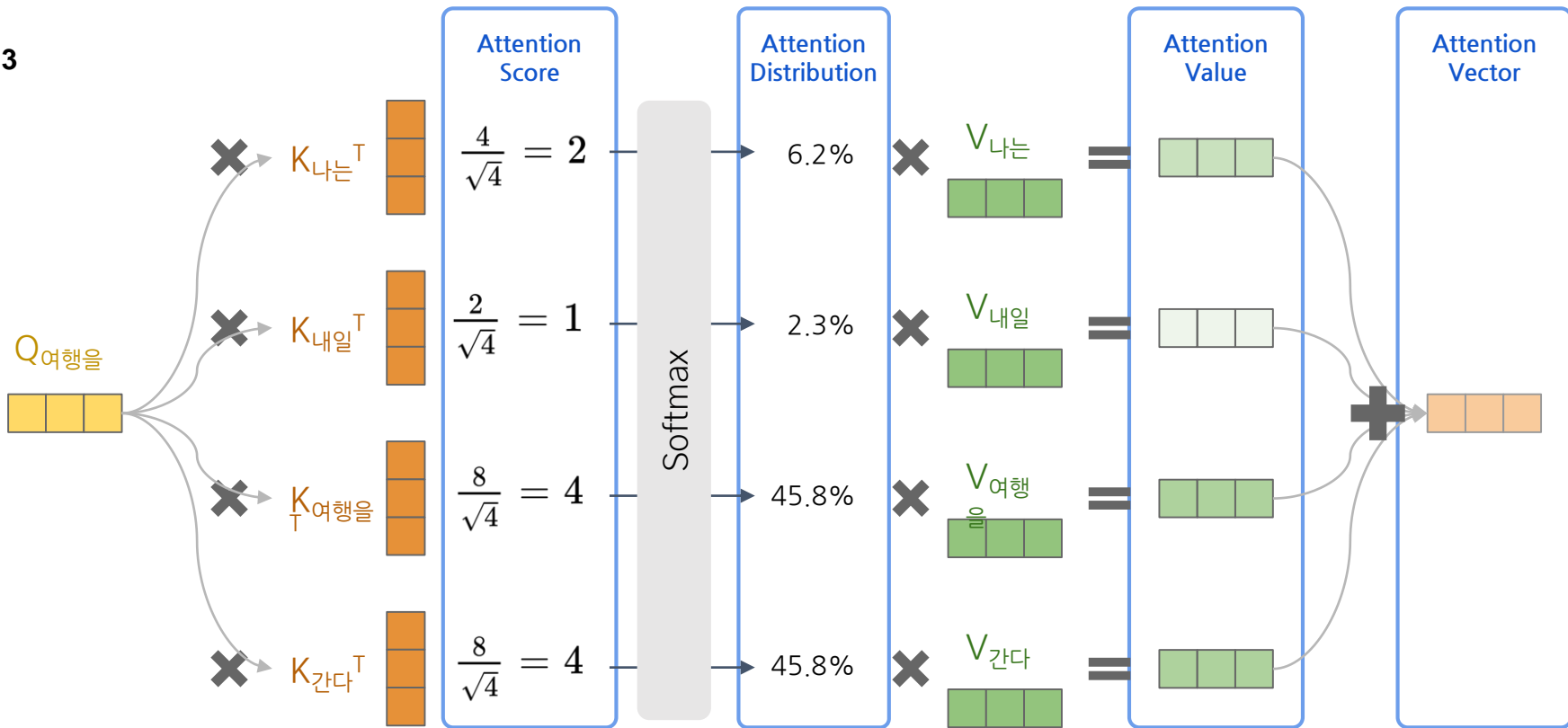


Scaled Dot-Product Attention (2) – Query, Key, Value



Scaled Dot-Product Attention (3) – Query, Key, Value

3



Scaled Dot-Product Attention (4) – Query, Key, Value

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

The diagram illustrates the Scaled Dot-Product Attention mechanism, showing the calculation of the Attention Matrix A from the Query (Q), Key (K), and Value (V) matrices.

Attention Distribution: The process involves calculating the Attention Score, which is the dot product of the Query matrix (Q) and the transposed Key matrix (K^T), scaled by the square root of the dimensionality of the key ($\sqrt{d_k}$).

Attention Score Calculation: The Query matrix (Q) is a 4x4 matrix (rows: 나는, 내일, 여행을, 간다) and the Key matrix (K^T) is a 4x4 matrix (columns: 나는, 내일, 여행을, 간다). The dot product is calculated as:

$$Q \times K^T = \begin{bmatrix} \text{나는} \\ \text{내일} \\ \text{여행을} \\ \text{간다} \end{bmatrix} \times \begin{bmatrix} \text{나는} & \text{내일} & \text{여행을} & \text{간다} \end{bmatrix}$$

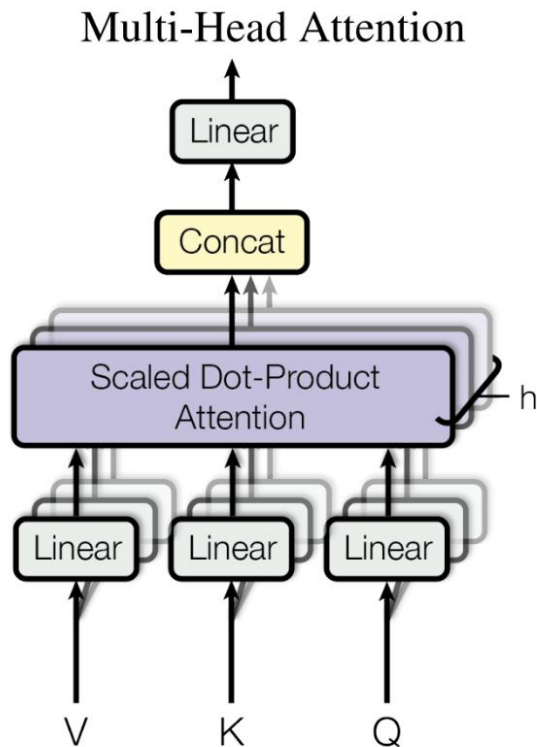
The result is a 4x4 matrix, which is then divided by $\sqrt{d_k}$ to produce the Attention Score matrix.

Attention Matrix: The Attention Score matrix is passed through a softmax operation (indicated by softmax_k) to produce the Attention Matrix (A), which is a 4x4 matrix (rows: 나는, 내일, 여행을, 간다).

Value Matrix: The Value matrix (V) is a 4x4 matrix (rows: 나는, 내일, 여행을, 간다).

The final result is the Attention Matrix (A) multiplied by the Value matrix (V), resulting in the Attention Matrix (A).

Multi-Head Attention (1)



$$MultiHead(Q, K, V) = Concat(head_1, \dots, head_h)W^O$$

where $head_i = Attention(QW_i^Q, KW_i^K, VW_i^V)$

$$QW_i^Q = [d_Q \times d_{model}] \times [d_{model} \times d_k] = [d_Q \times d_k]$$

$$KW_i^K = [d_K \times d_{model}] \times [d_{model} \times d_k] = [d_K \times d_k]$$

$$VW_i^V = [d_V \times d_{model}] \times [d_{model} \times d_v] = [d_V \times d_v]$$

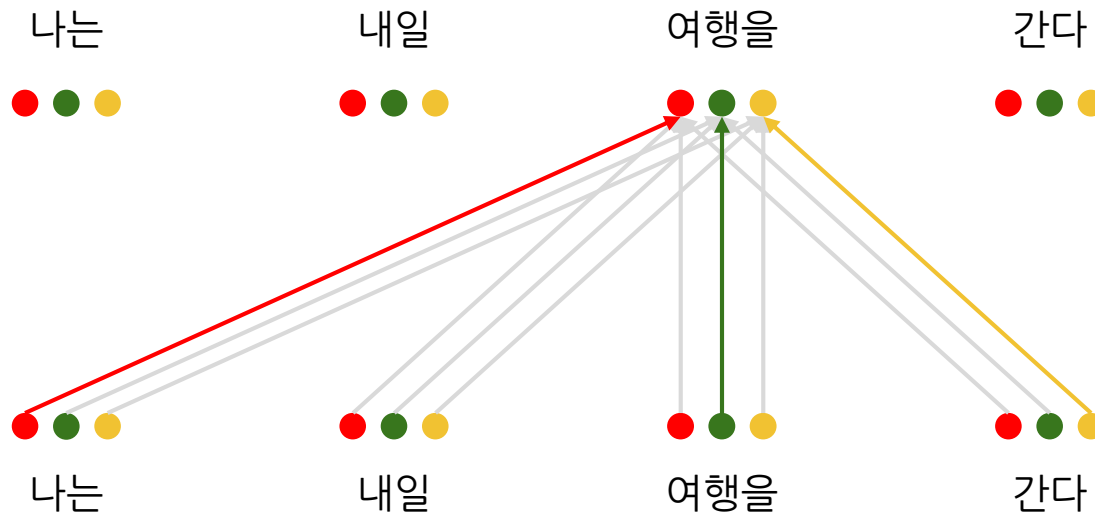


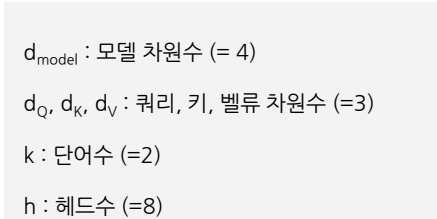
$$Attention(QW_i^Q, KW_i^K, VW_i^V) = [d_V \times d_v]$$



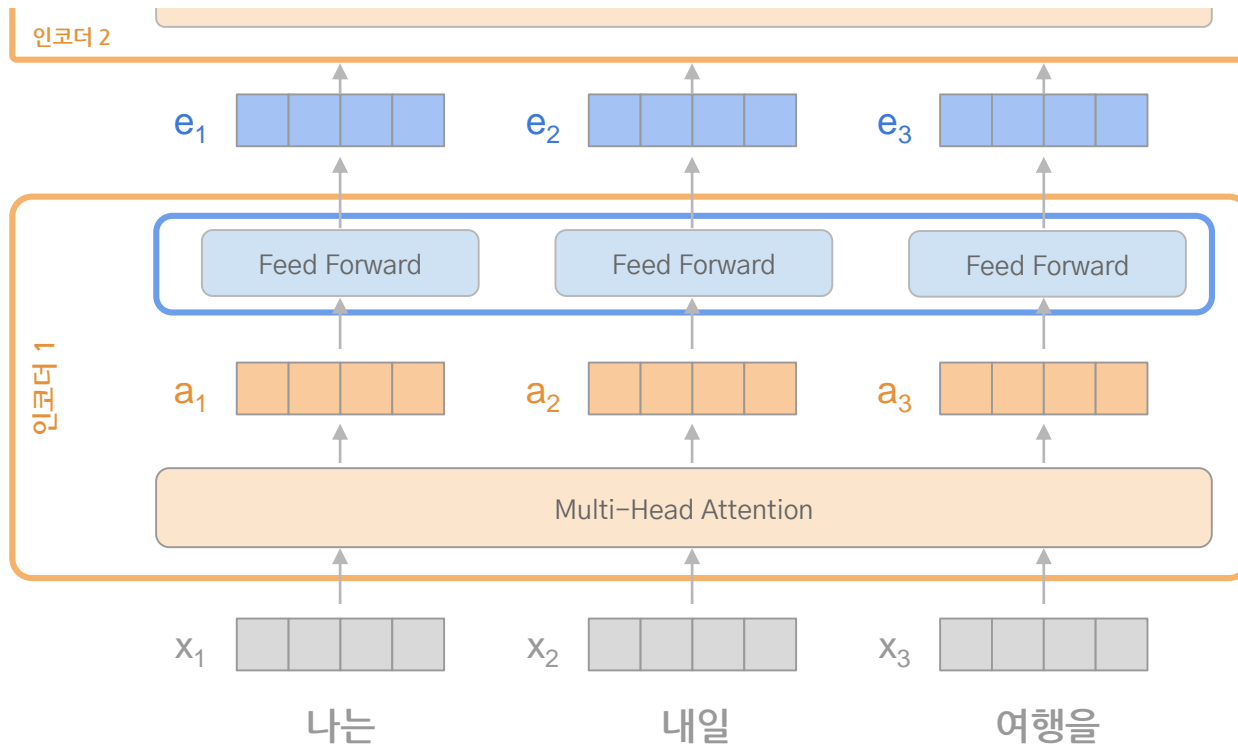
$$Concat(QW_i^Q, KW_i^K, VW_i^V)W^O = [d_V \times h d_v] \times [h d_v \times d_{model}] = [d_V \times d_{model}]$$

Multi-Head Attention (2) – 예제

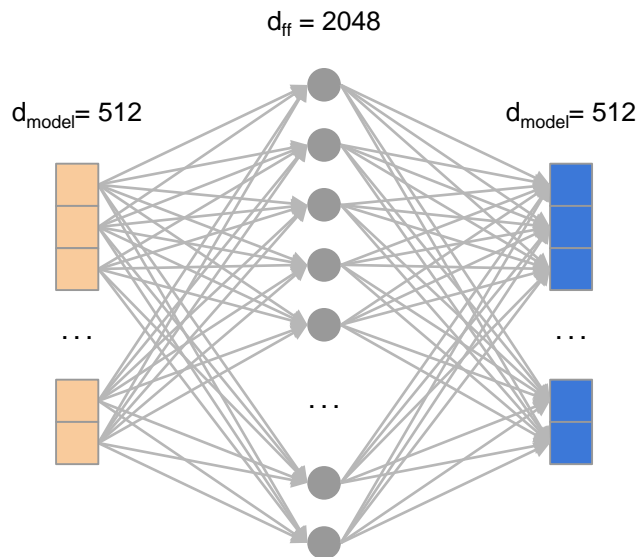




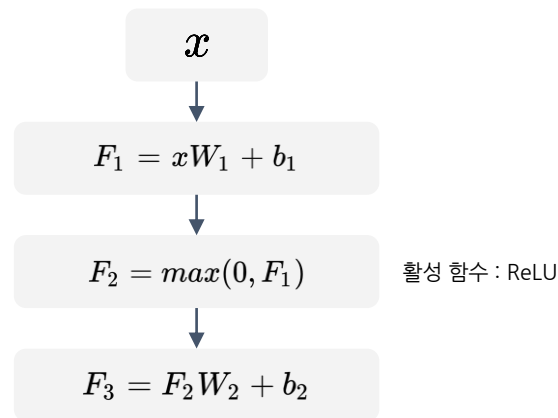
Transformer 인코더 – Feed Forward



Position-wise Feed-Forward Networks



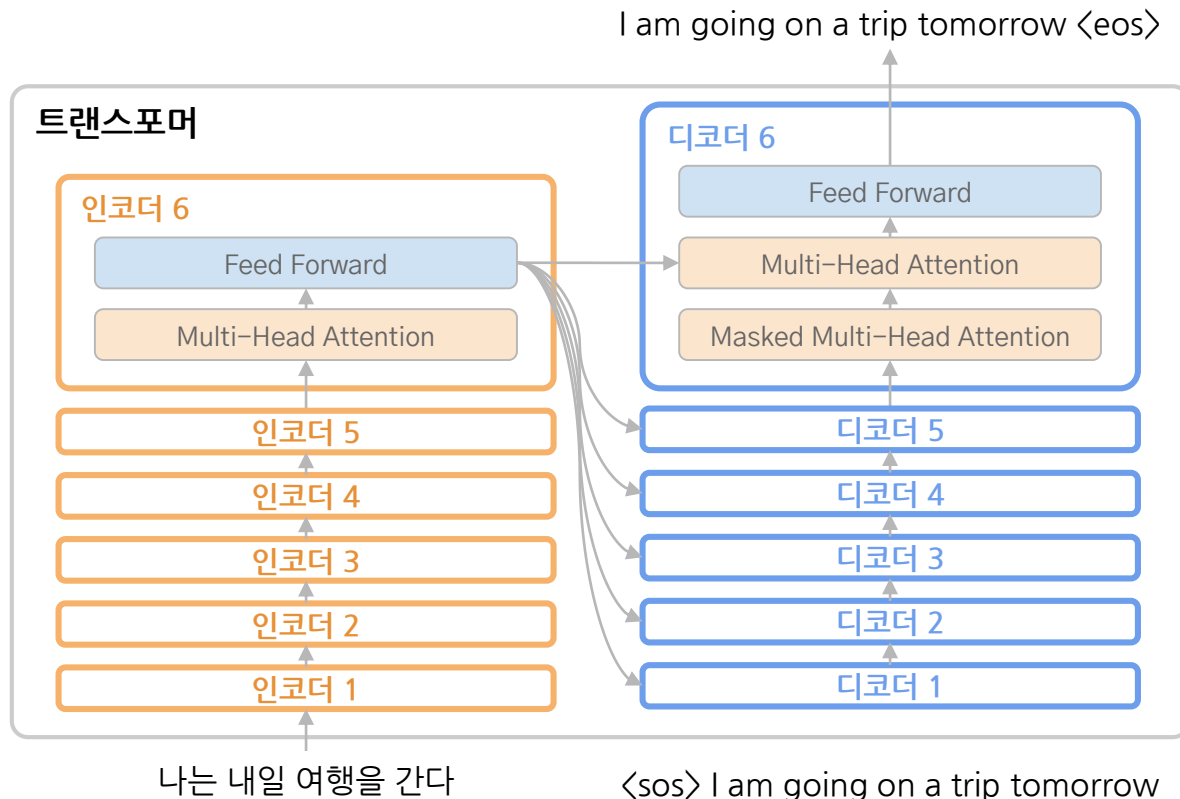
$$FFN(x) = \max(0, xW_1 + b_1)W_2 + b_2$$



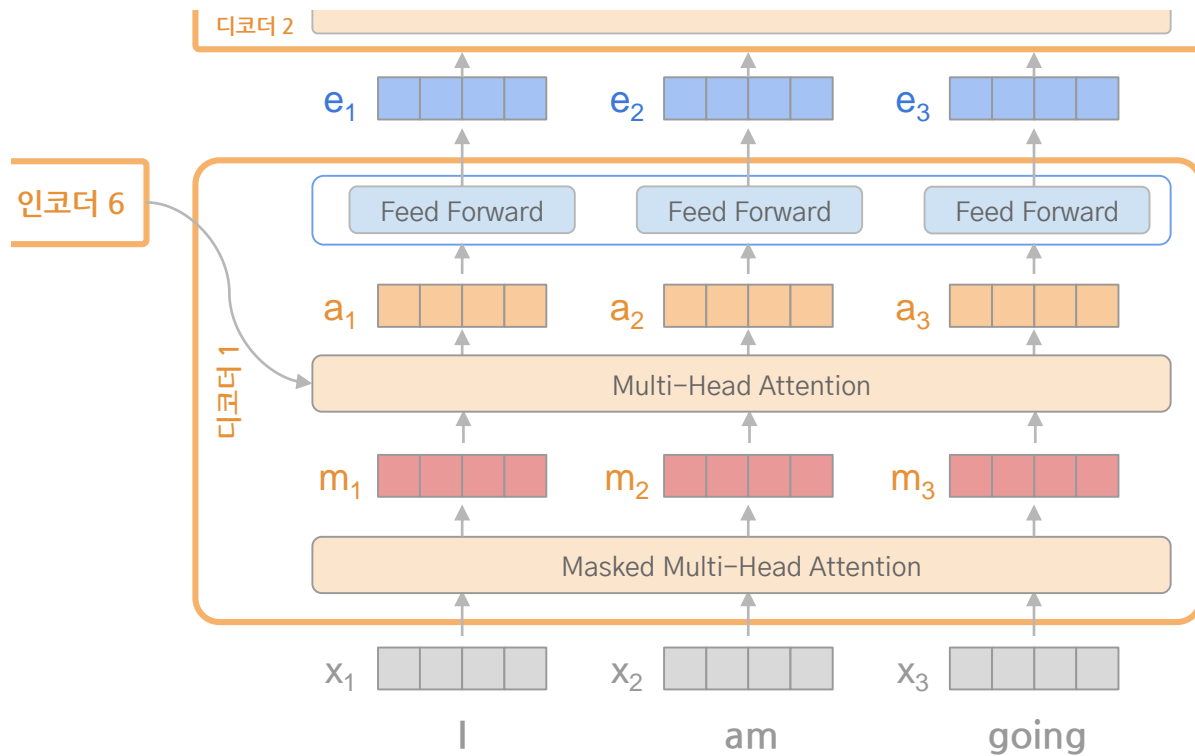
Decoder

Transformer

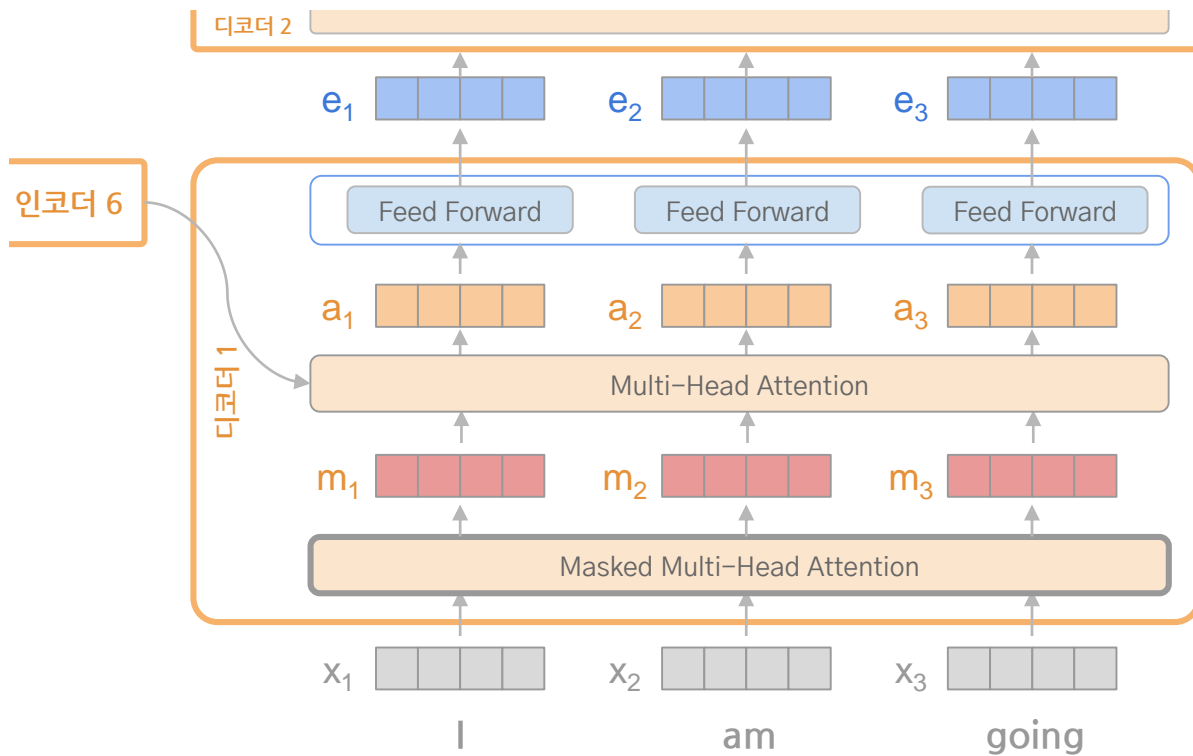
Transformer 구조 - 학습



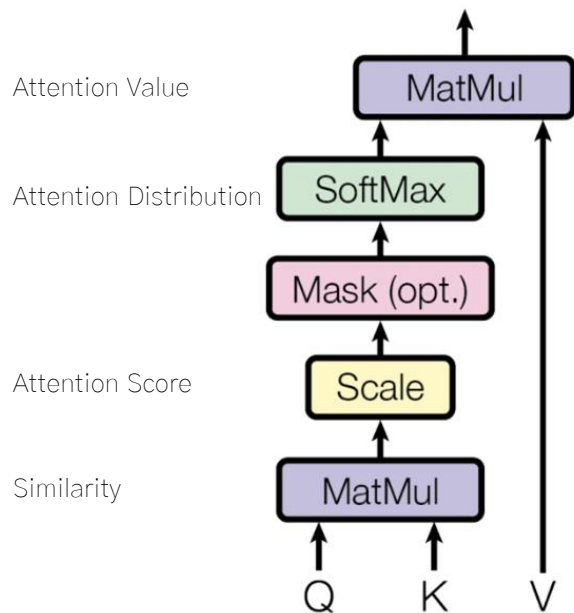
Transformer 디코더



Masked Multi-Head Attention



Scaled Dot Product Attention (Masking)



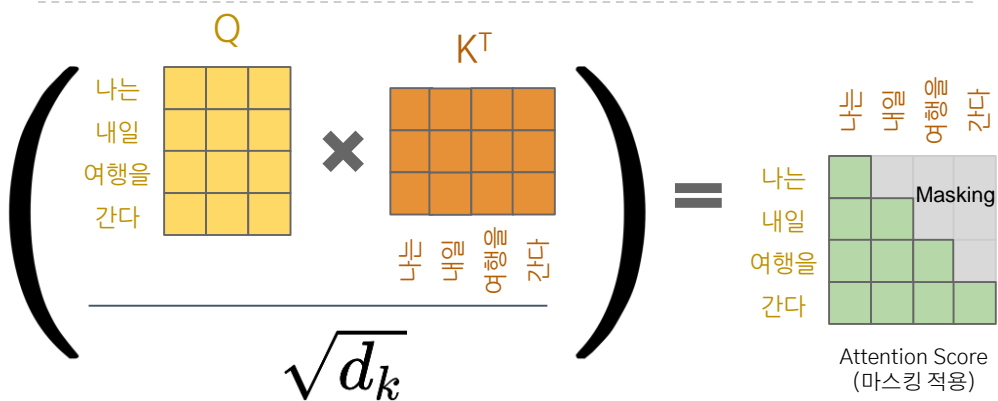
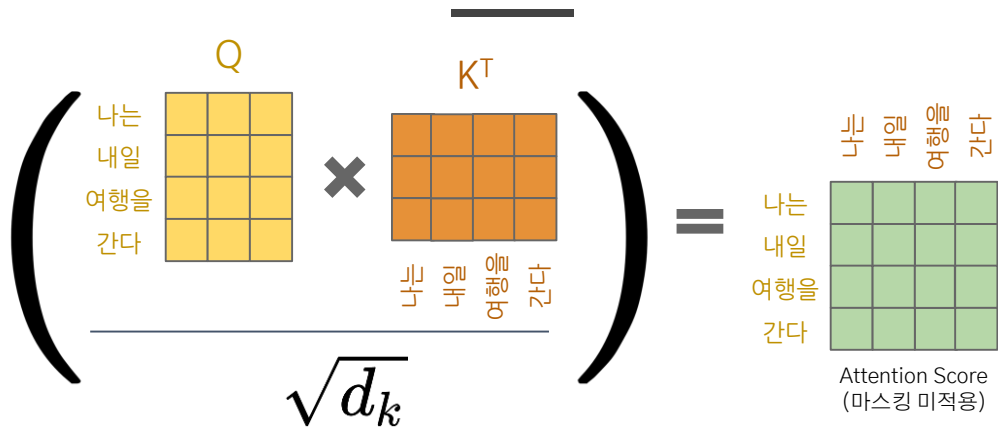
$$Attention(Q, K, V) = softmax_k \left(\frac{QK^T}{\sqrt{d_k}} \right) V$$

어텐션 분포 어텐션 스코어 단어간 유사도

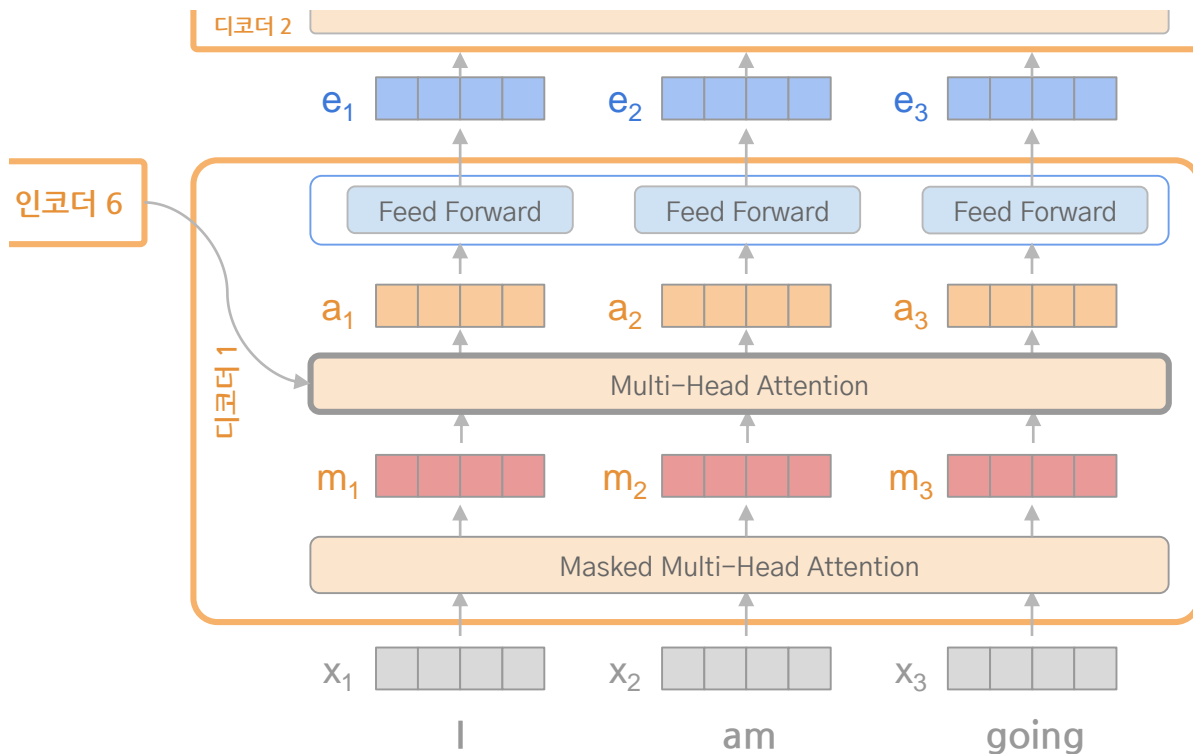
The diagram shows the matrix multiplication of Q (4x4, yellow) and K^T (4x4, orange) to produce a 4x4 matrix. The result is then divided by $\sqrt{d_k}$. The resulting matrix is shown with a **Masking** operation, where the diagonal elements are highlighted in green and the off-diagonal elements are grayed out. The labels for the matrices are: Q (나는, 내일, 여행을, 간다) and K^T (나는, 내일, 여행을, 간다). The resulting matrix is labeled with the same words and a **Masking** label.

Decoding 시 다음 단어의 Attention을 고려하지 못하도록 Masking

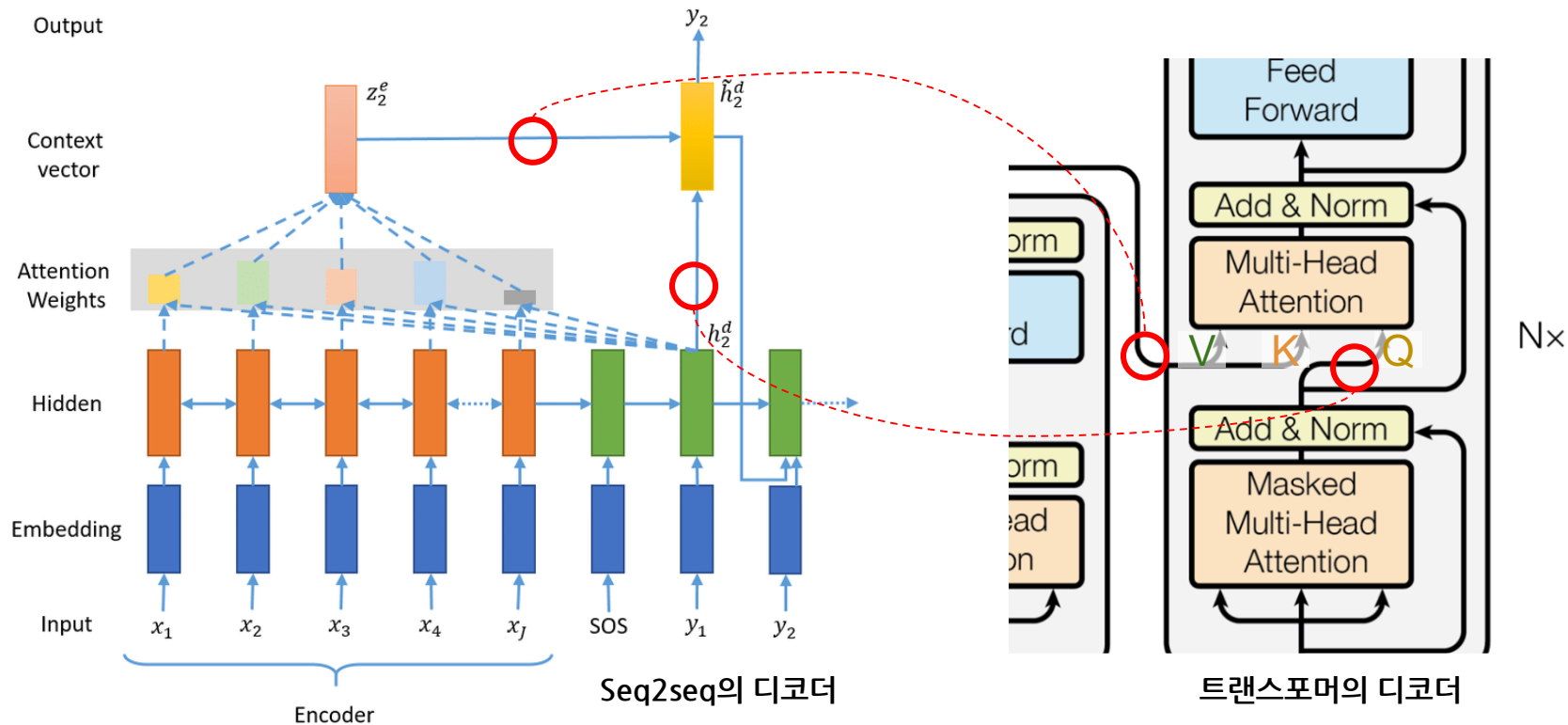
Scaled Dot Product Attention (Masking)



Encoder-Decoder Multi-Head Attention



Transformer 개요 (6)



기존 Seq2Seq with Attention과 구조는 다르지만 유사한 적용

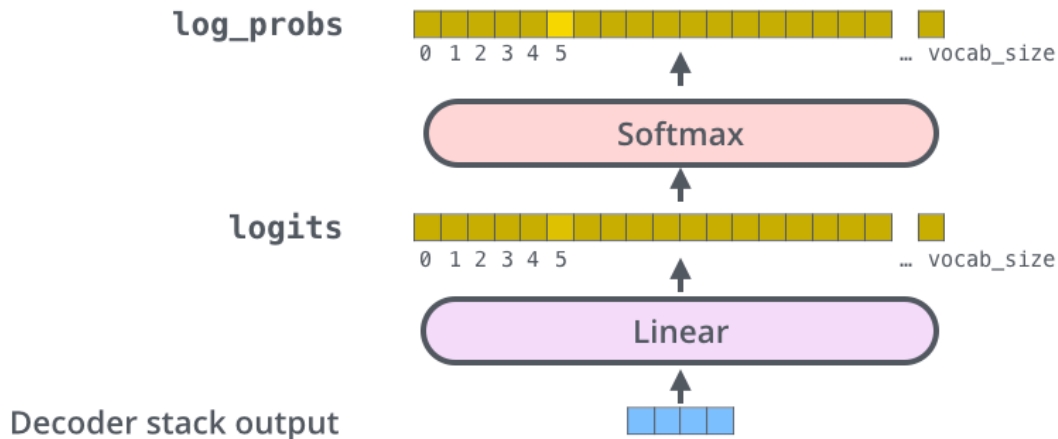
Transformer Review (7)

Which word in our vocabulary
is associated with this index?

am

Get the index of the cell
with the highest value
(argmax)

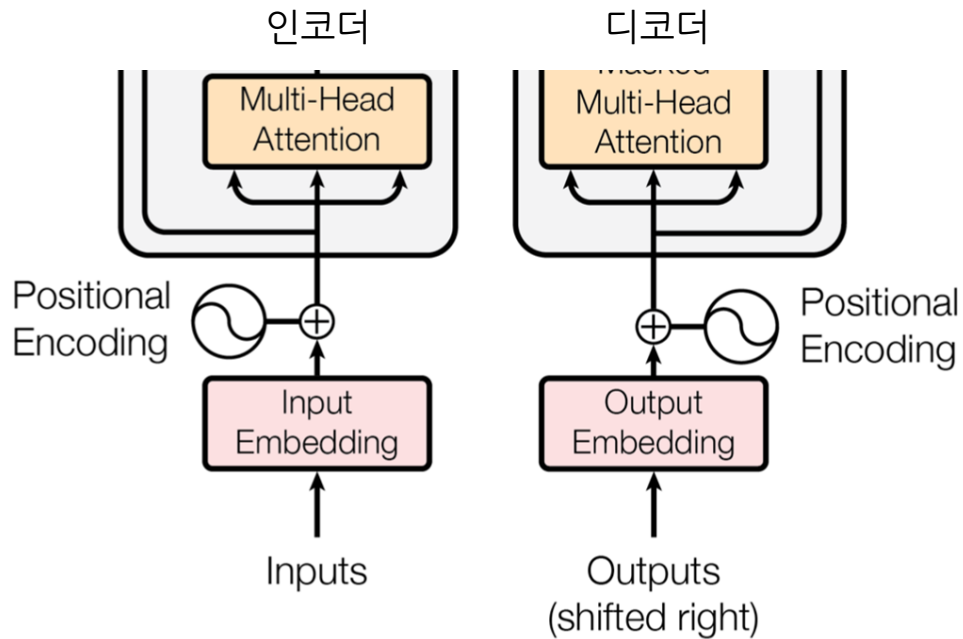
5



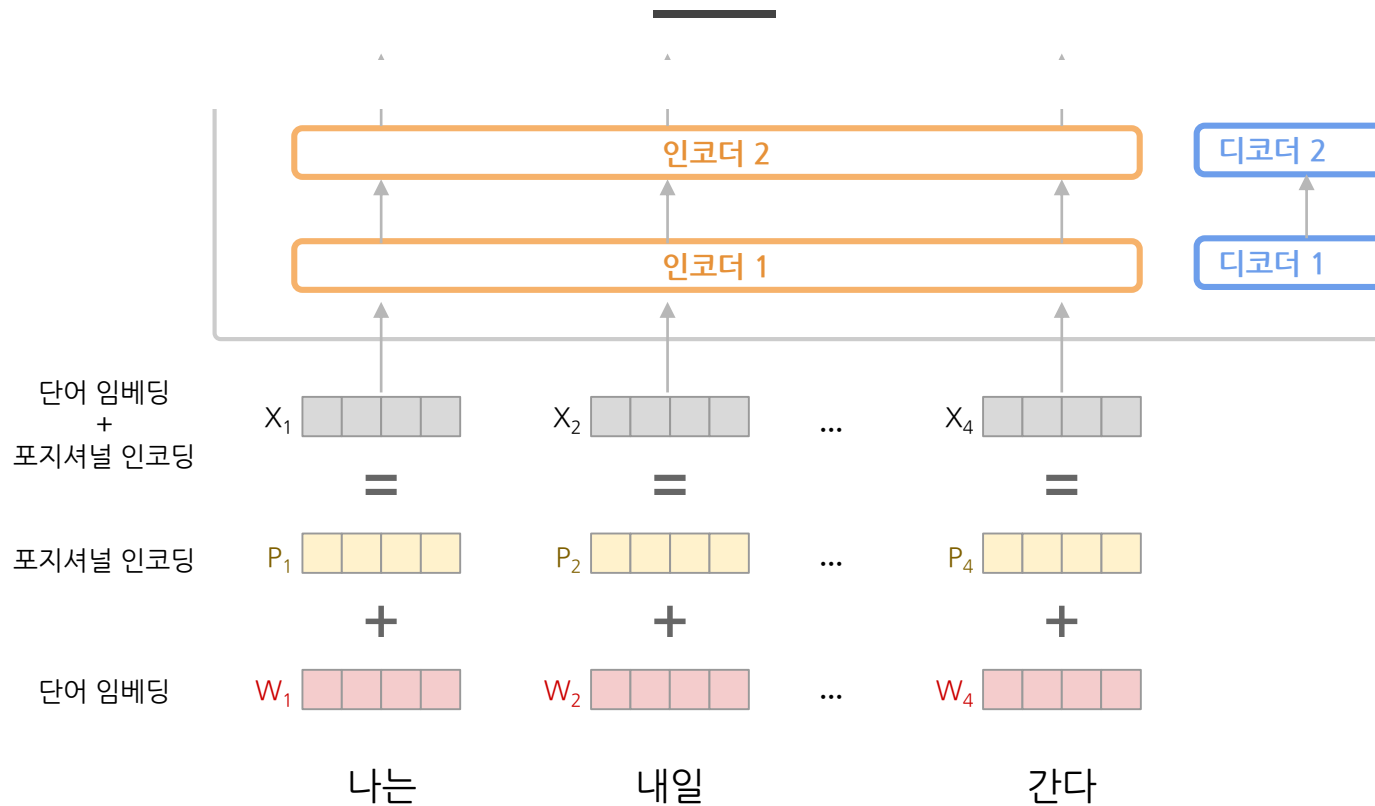
Input Embedding

Transformer

Transformer 개요 (5)

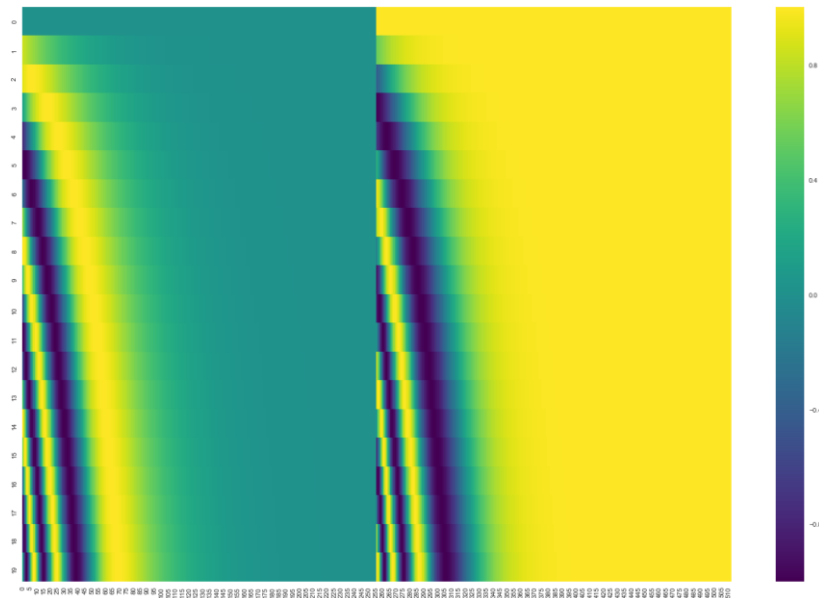


Transformer 구조 - 학습



위치(Position)에 대한 절대적 위치를 표현하는 것이 아니라 대변할 수 있는 encoding

Positional Encoding (3)

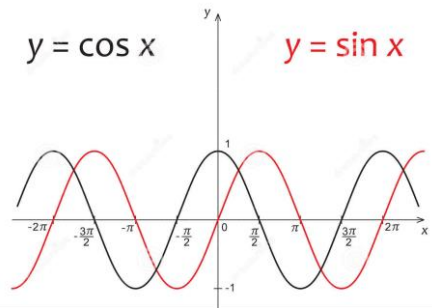


$$PE_{(pos, 2i)} = \sin(pos/10000^{2i/d_{model}})$$

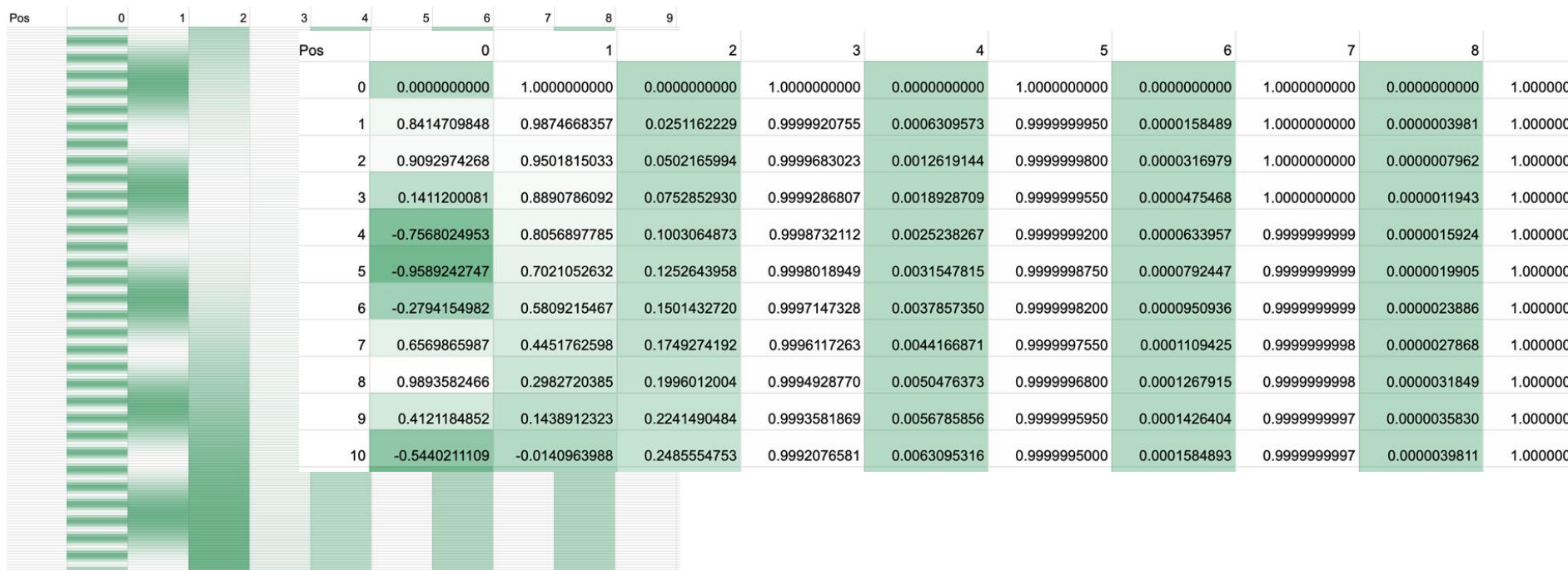
$$PE_{(pos, 2i+1)} = \cos(pos/10000^{2i/d_{model}})$$

$$\left[\sin\left(\frac{pos}{10000^0}\right), \cos\left(\frac{pos}{10000^0}\right), \sin\left(\frac{pos}{10000^{2/4}}\right), \cos\left(\frac{pos}{10000^{2/4}}\right) \right]$$

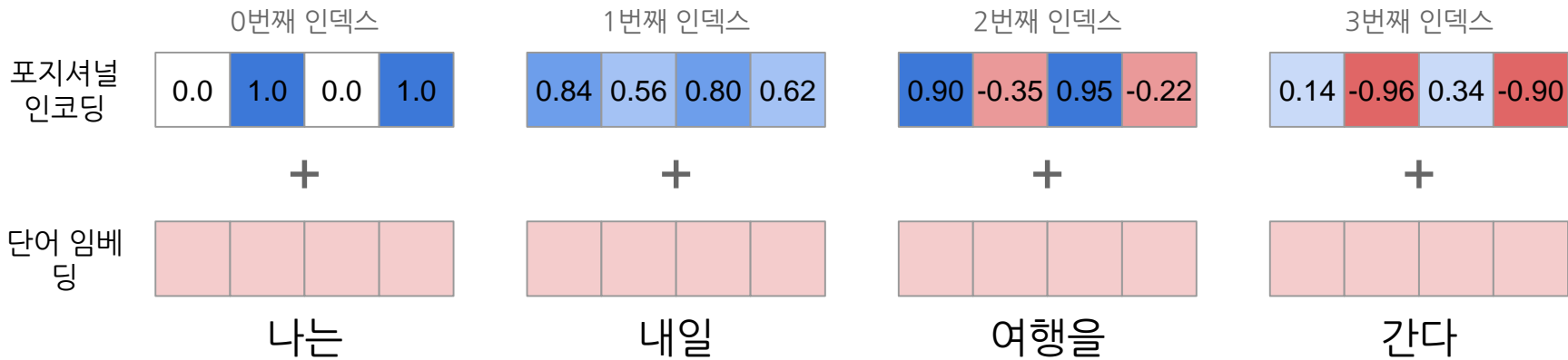
$$\left[\sin(pos), \cos(pos), \sin\left(\frac{pos}{100}\right), \cos\left(\frac{pos}{100}\right) \right]$$



Positional Encoding (4)



위치(Position)에 대한 Unique한 벡터가 생성됨. 동일 위치는 같은 벡터가 생성됨



잔차연결 & 정규화

Transformer

잔차연결 & 정규화

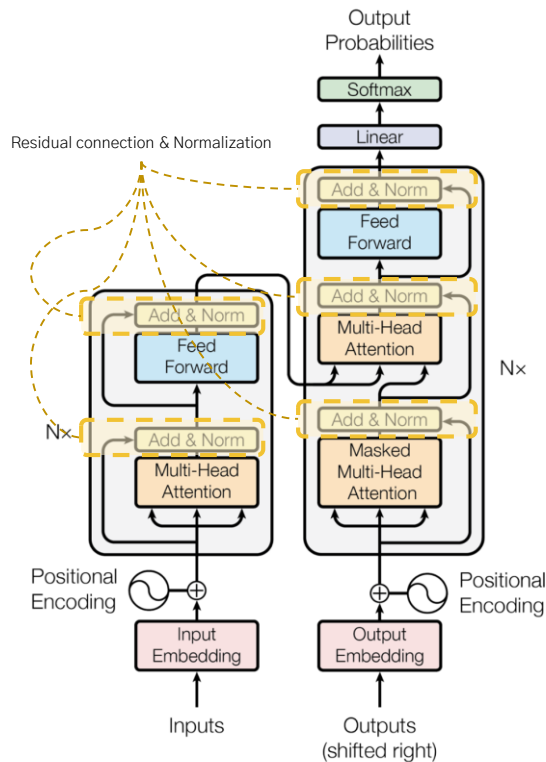
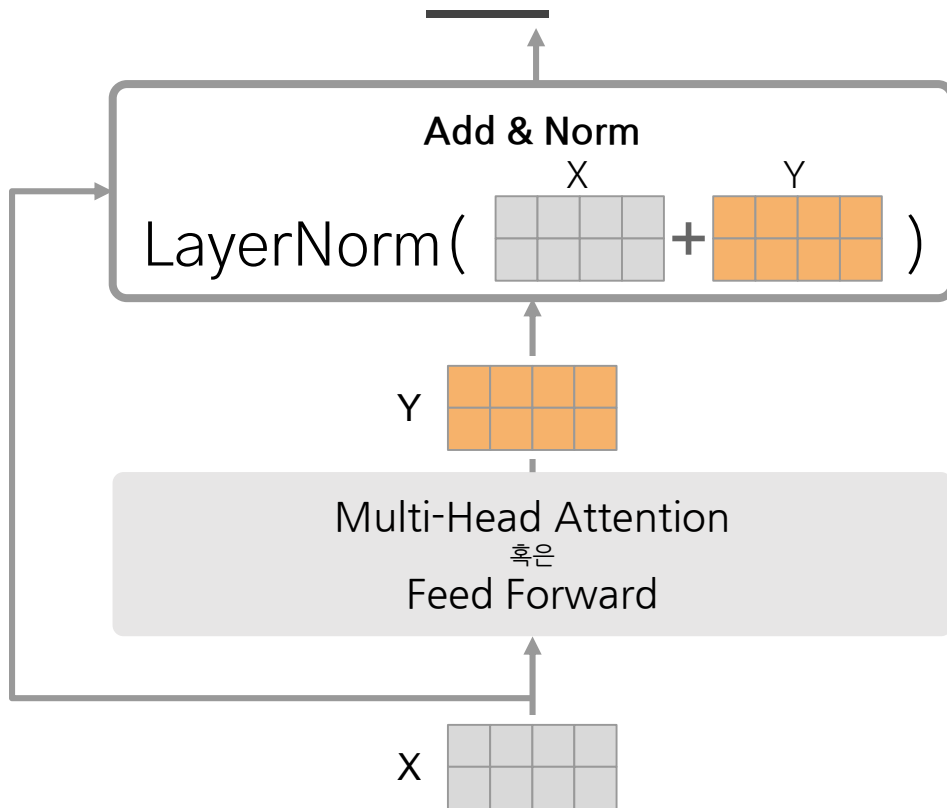


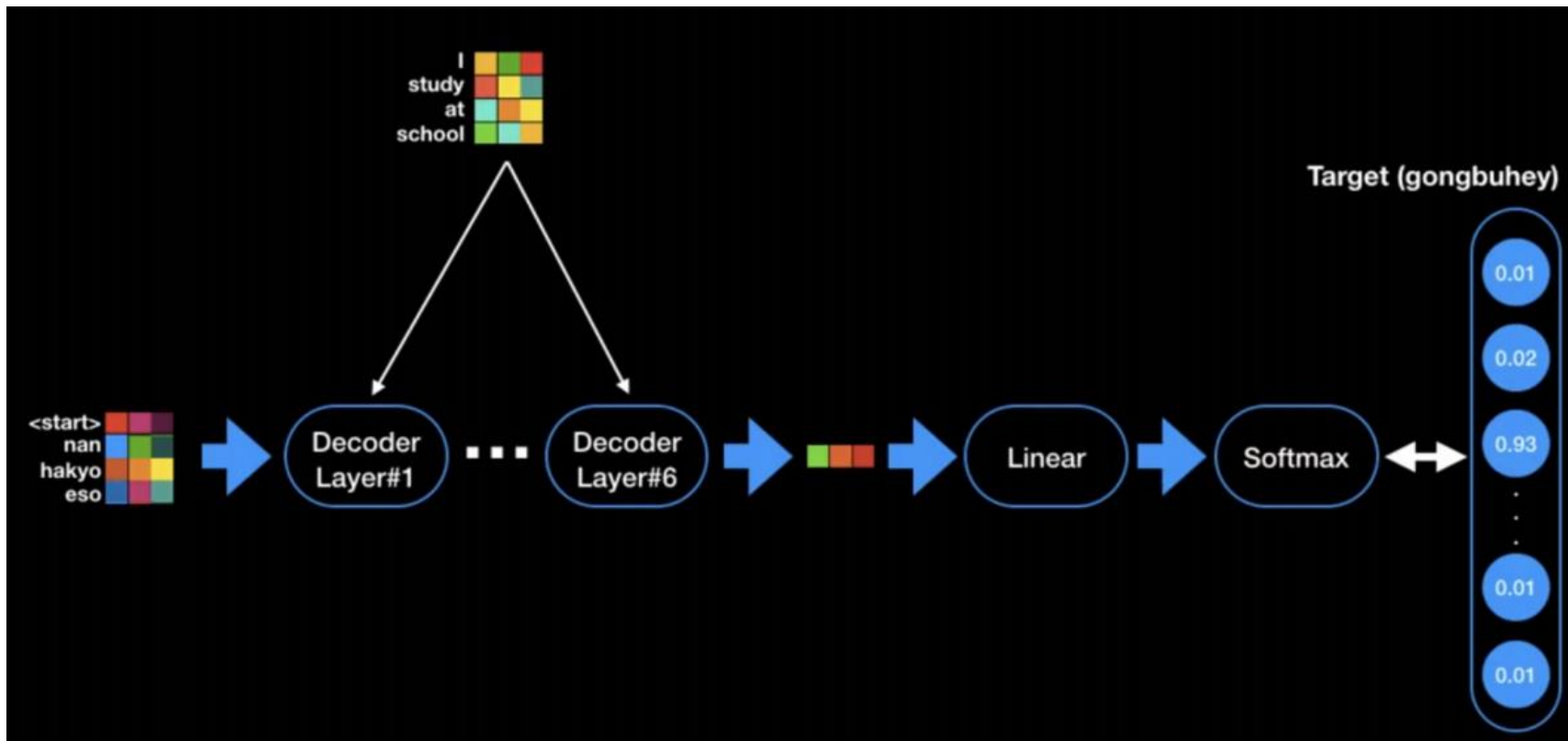
Figure 1: The Transformer - model architecture.

잔차연결 & 정규화



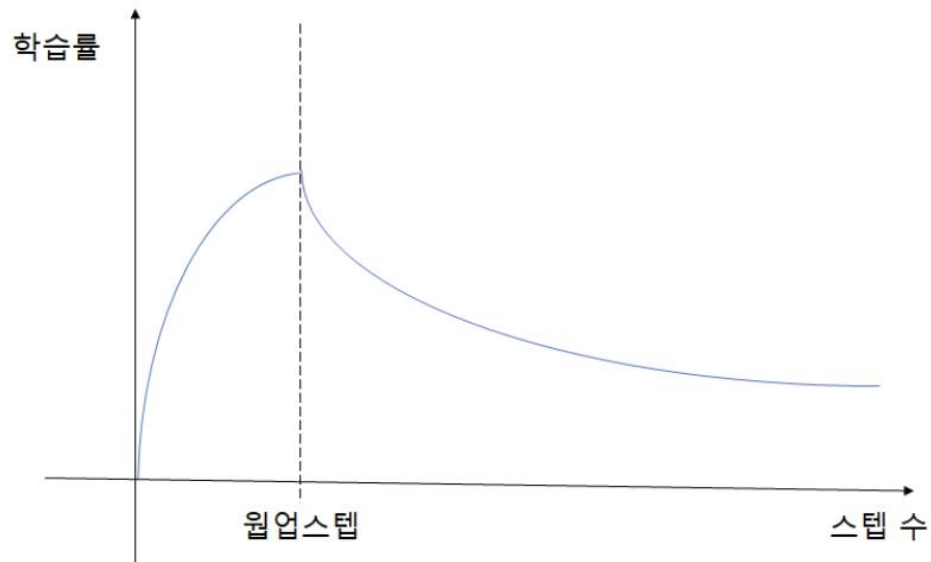
입력 행렬이 서브레이어를 통과했을때 정보가 유실되는 것을 막고자 잔차연결

라벨 스무딩



최적화

$$lrate = d_{model}^{-0.5} \cdot \min(stepnum^{-0.5}, stepnum \cdot warmupsteps^{-1.5})$$



트랜스포머 리뷰

BERT

Transformer Review (1)

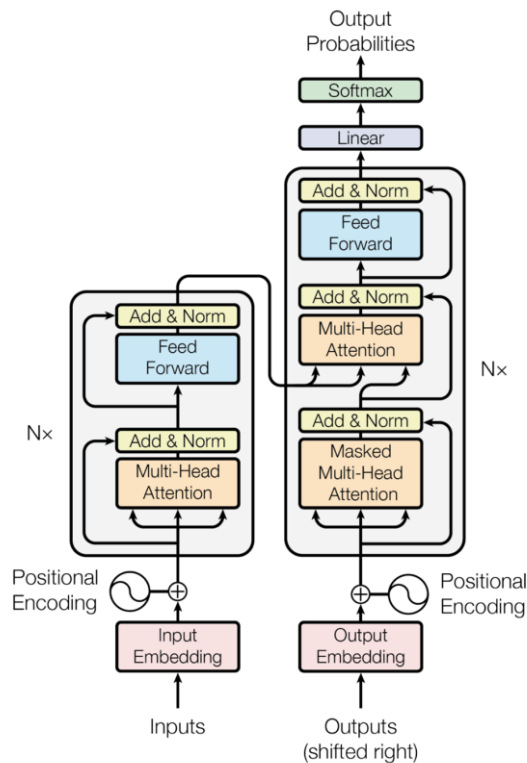
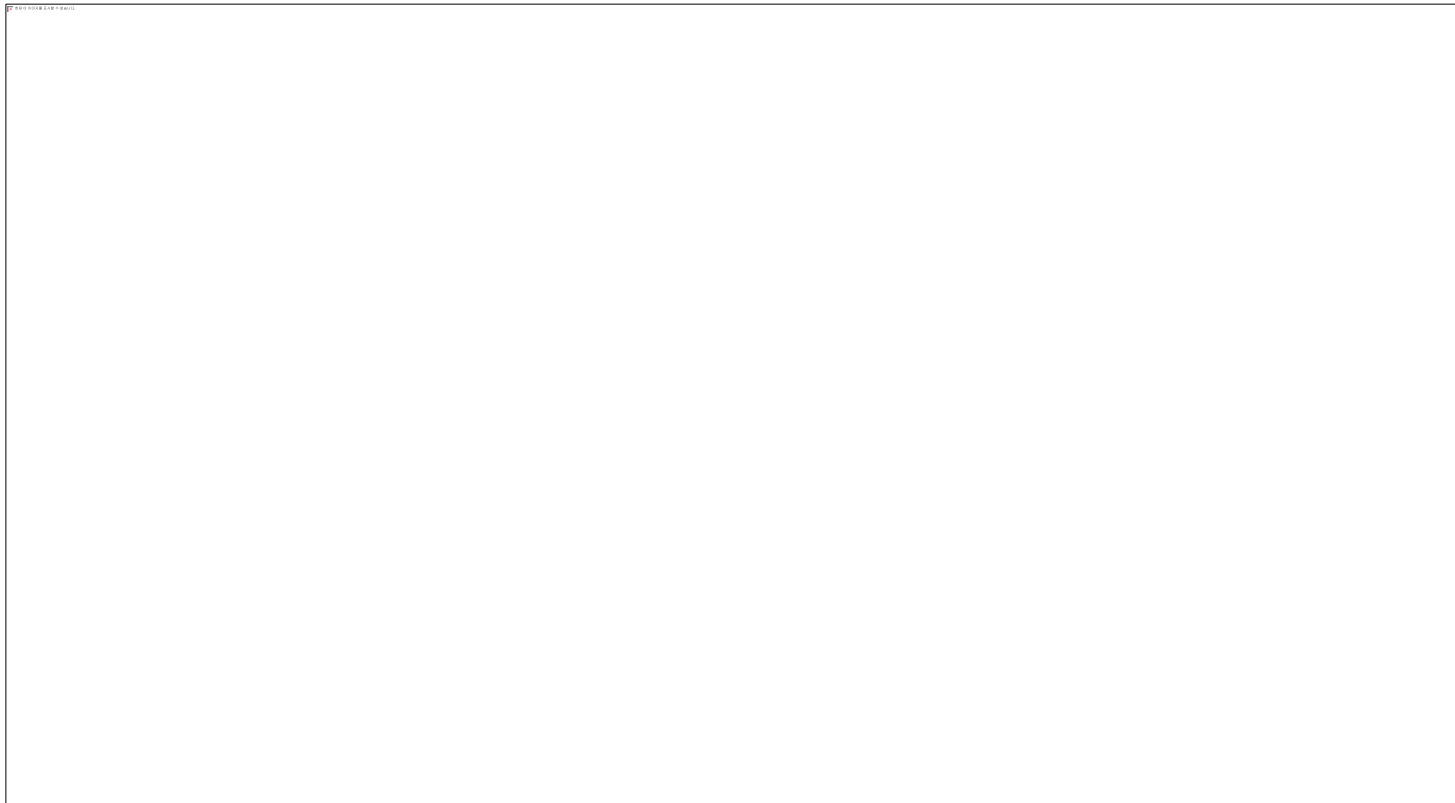


Figure 1: The Transformer - model architecture.

Transformer Review (2)



Transformer Review (3)

