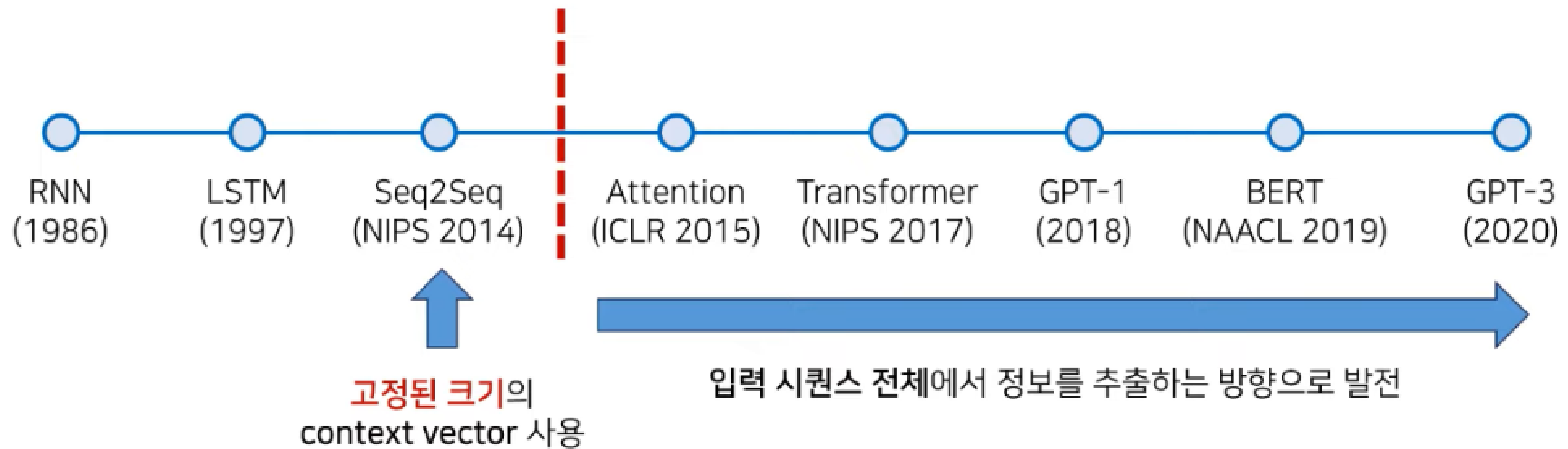


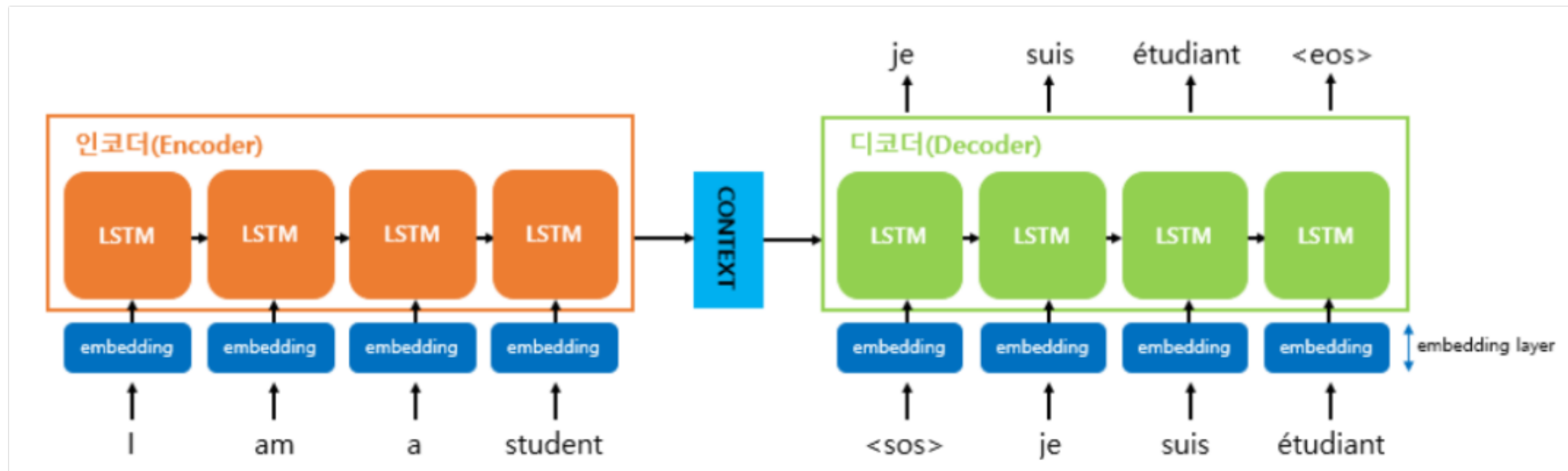
Transformer : Attention Is All You Need

윤주환, 최은서

기계 번역의 발전

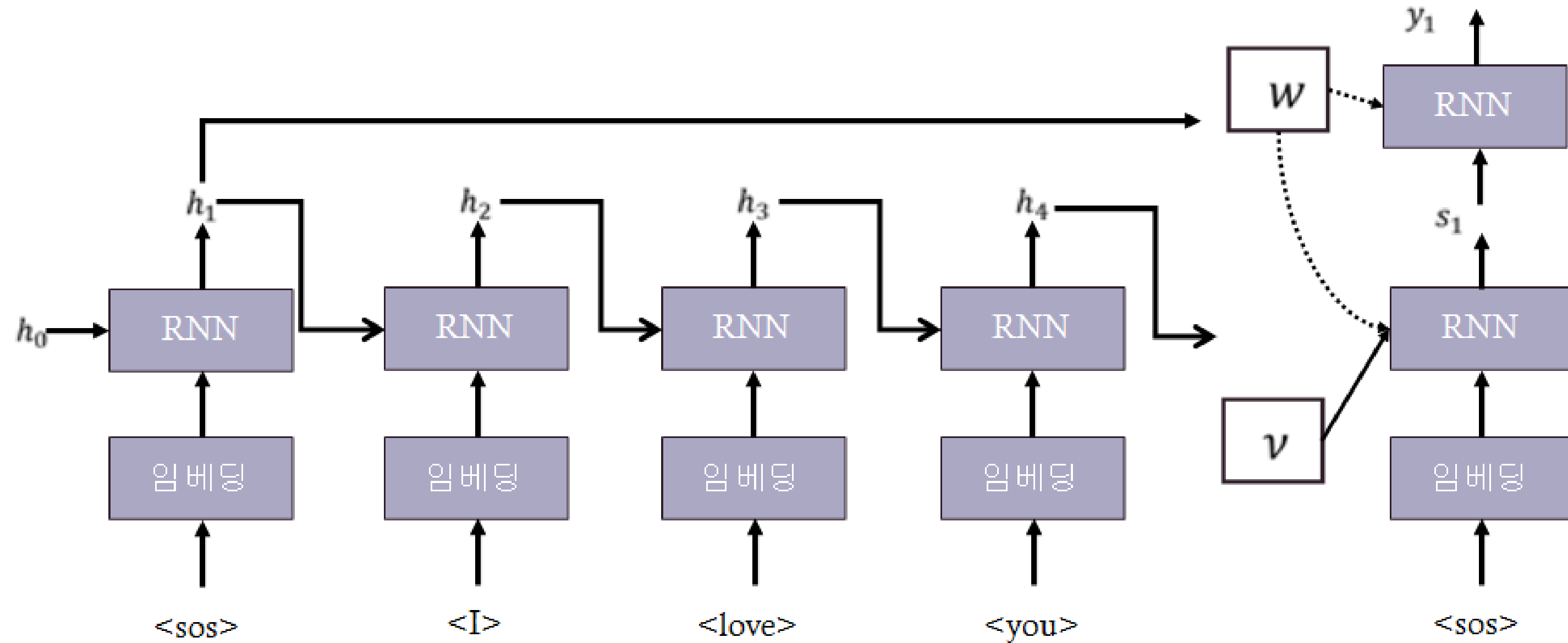


기존 Seq2Seq 모델들의 한계



한계점 : 고정된 크기의 context vector에 소스 문장의 정보를 압축하는
방식으로 학습되기 때문에 병목 현상이 발생하여 성능 하락

Seq2Seq with Attention



Seq2Seq with Attention : Decoder

디코더는 매번 인코더의 모든 출력 중에서 어떤 정보가 중요한지 계산

- i = 현재 디코더가 처리중인 인덱스
- j = 각각의 인코더 출력 인덱스

Energy

$$e_{ij} = \alpha(s_{i-1}, h_j)$$

h_j : Encoder 파트 각각의 state

s_{i-1} : Decoder 이전에 출력했던 단어를 만들기 위해 사용했던 hidden state

Weight

에너지 값에 softmax (0~1 사이의 확률 값)를 취한 값 $\alpha_{ij} = \sum_{k=1}^{T_x} \frac{\exp(e_{ij})}{\exp(e_{ik})}$

The context vector c_i is, then, computed as a weighted sum of these annotations h_i :

$$c_i = \sum_{j=1}^{T_x} \alpha_{ij} h_j. \quad (5)$$

The weight α_{ij} of each annotation h_j is computed by

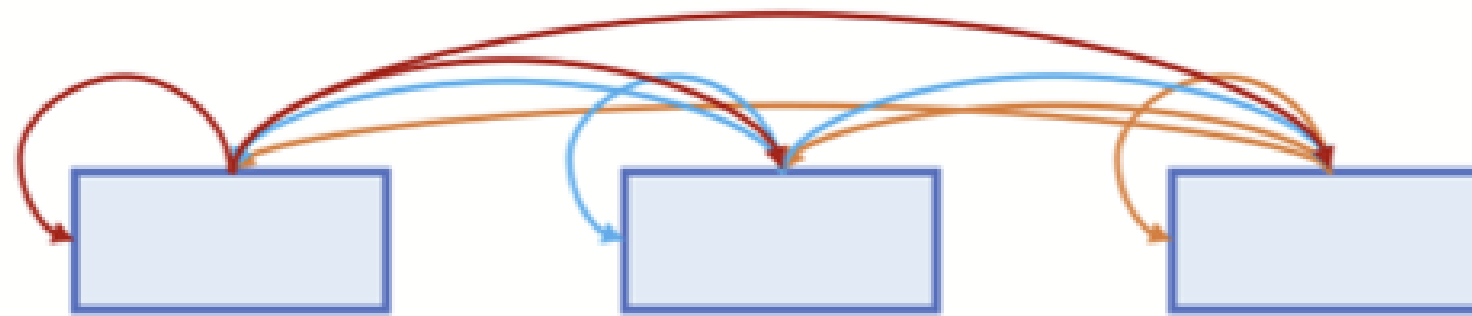
$$\alpha_{ij} = \frac{\exp(e_{ij})}{\sum_{k=1}^{T_x} \exp(e_{ik})},$$

where

$$e_{ij} = a(s_{i-1}, h_j)$$

Attention in this paper

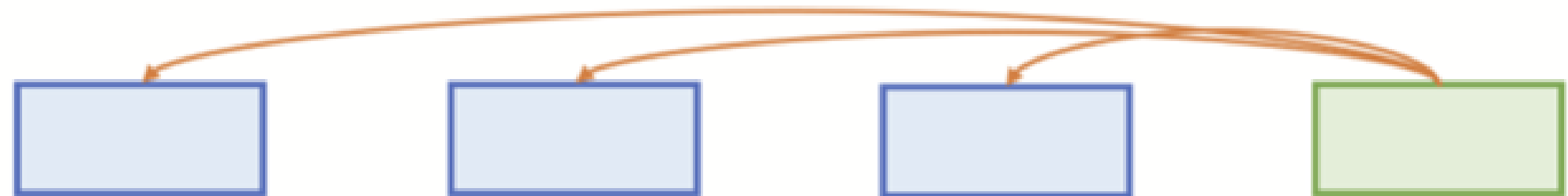
Encoder Self-Attention:



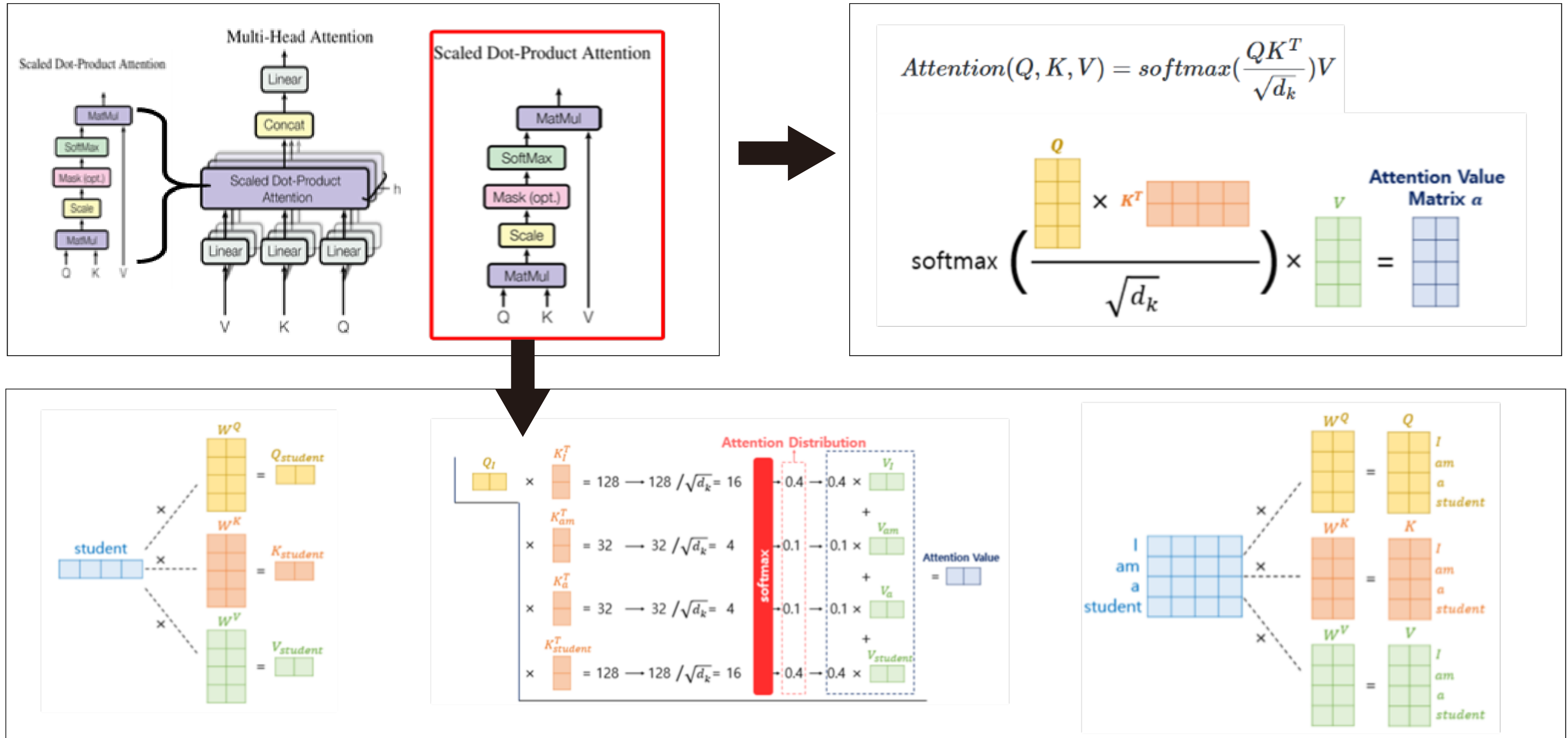
Masked Decoder Self-Attention:



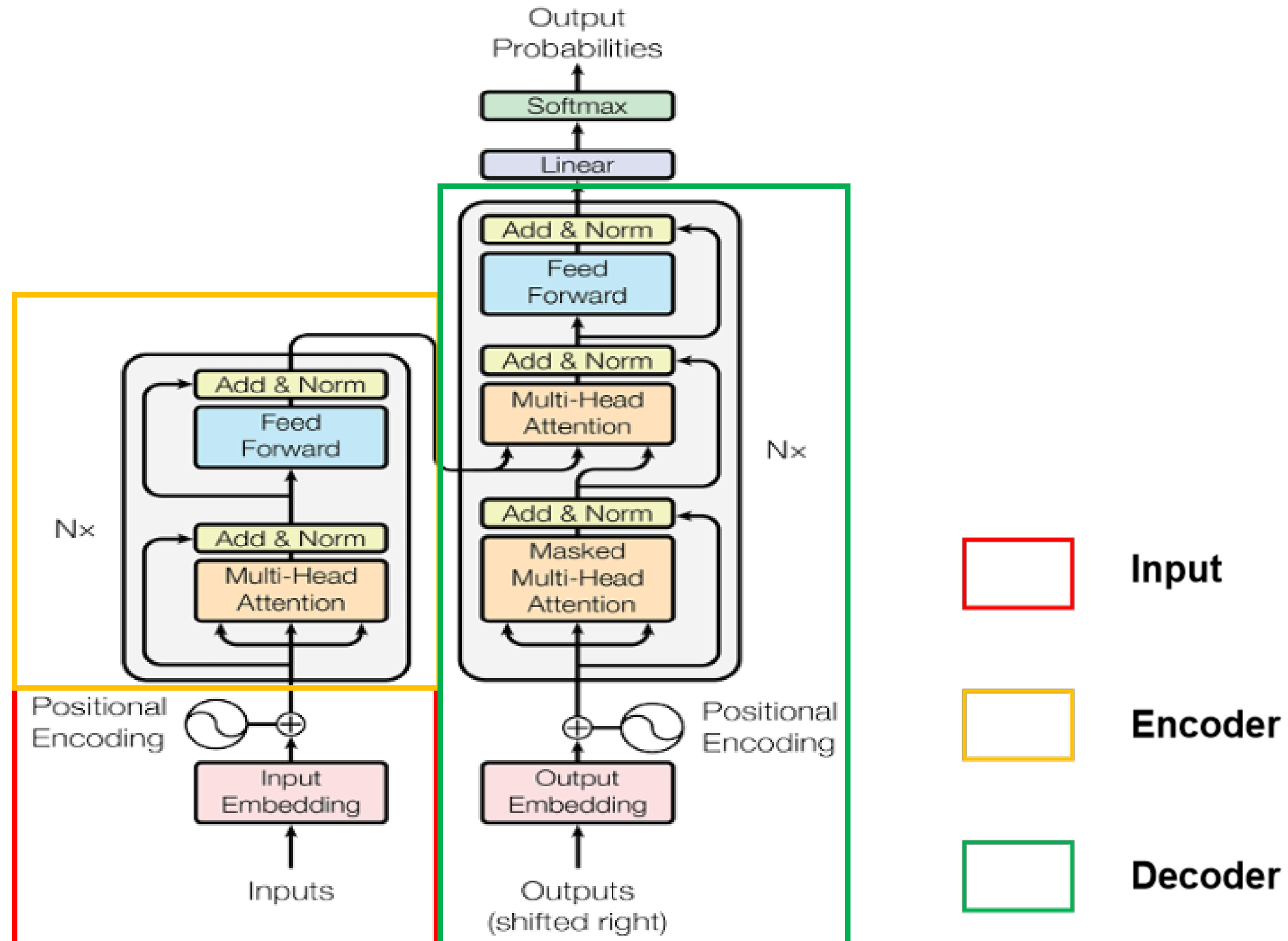
Encoder-Decoder Attention:



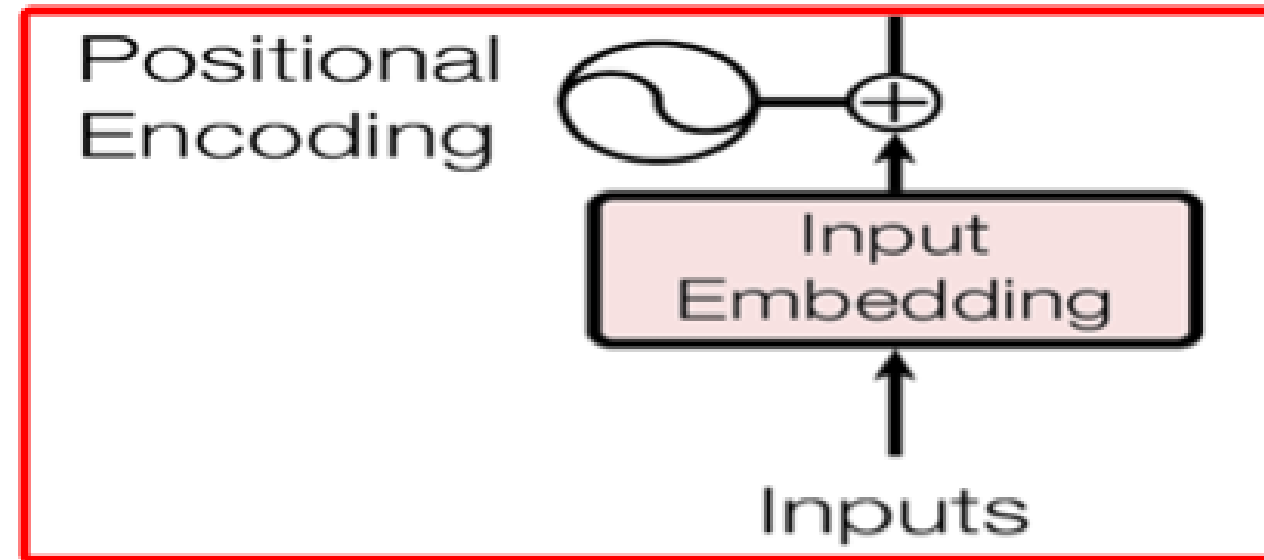
Attention in this paper



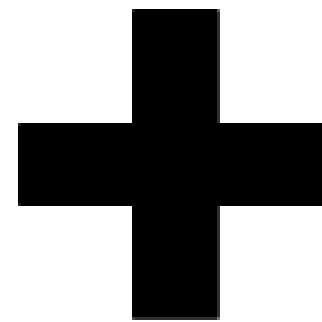
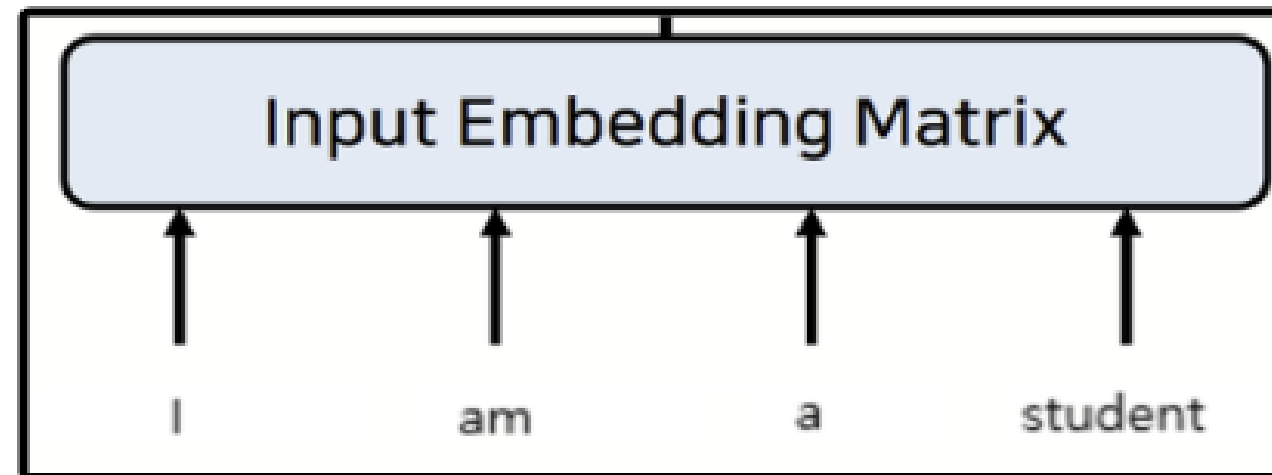
Transformer Structure



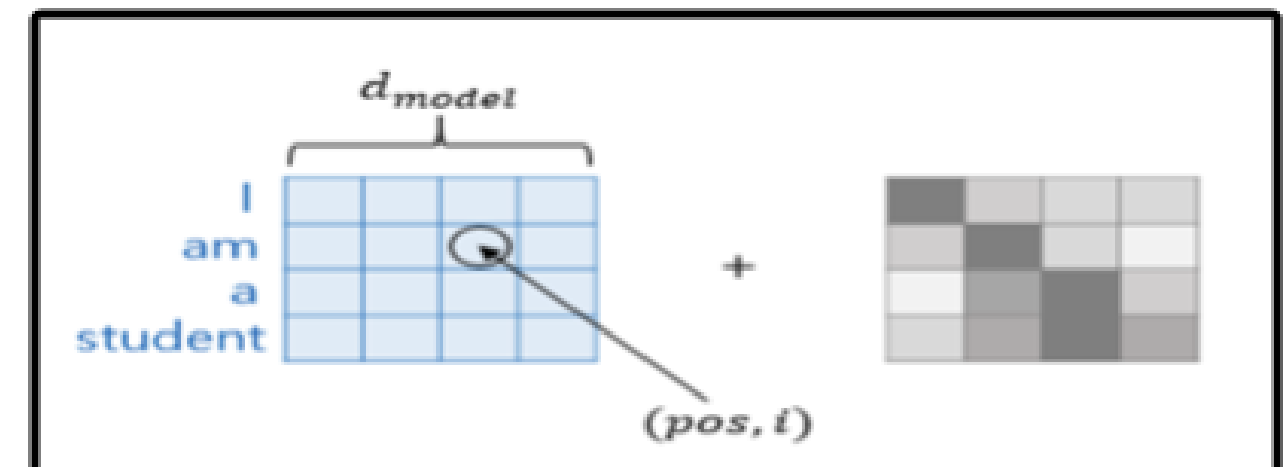
Input into the Encoder



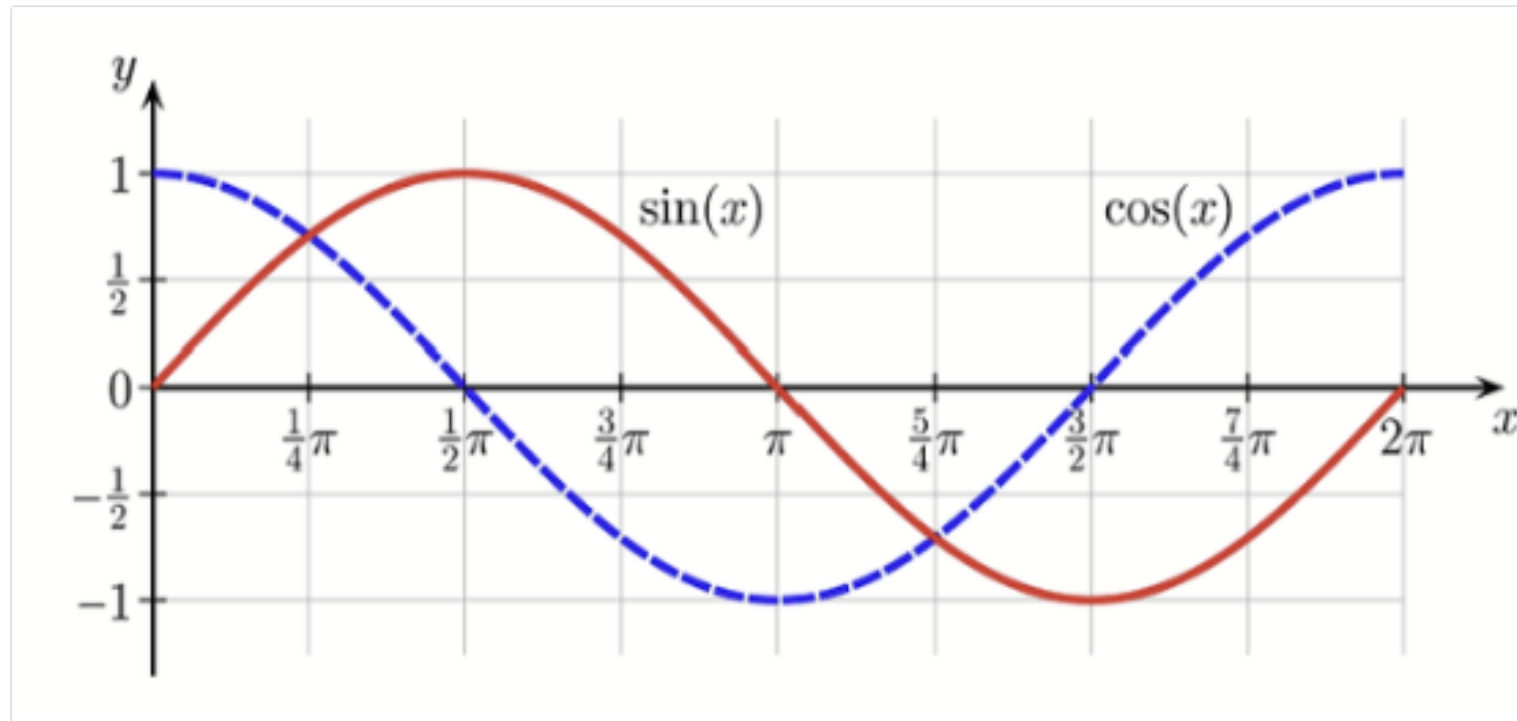
Input Embedding



Positional Encoding

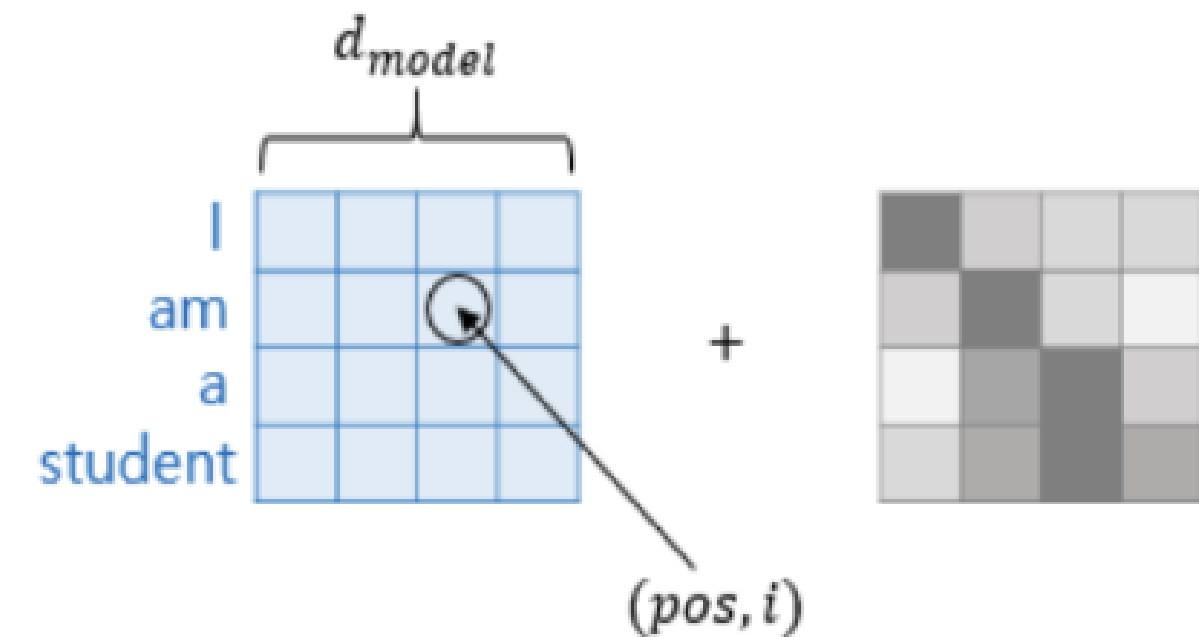
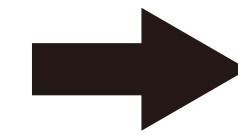


Positional encoding

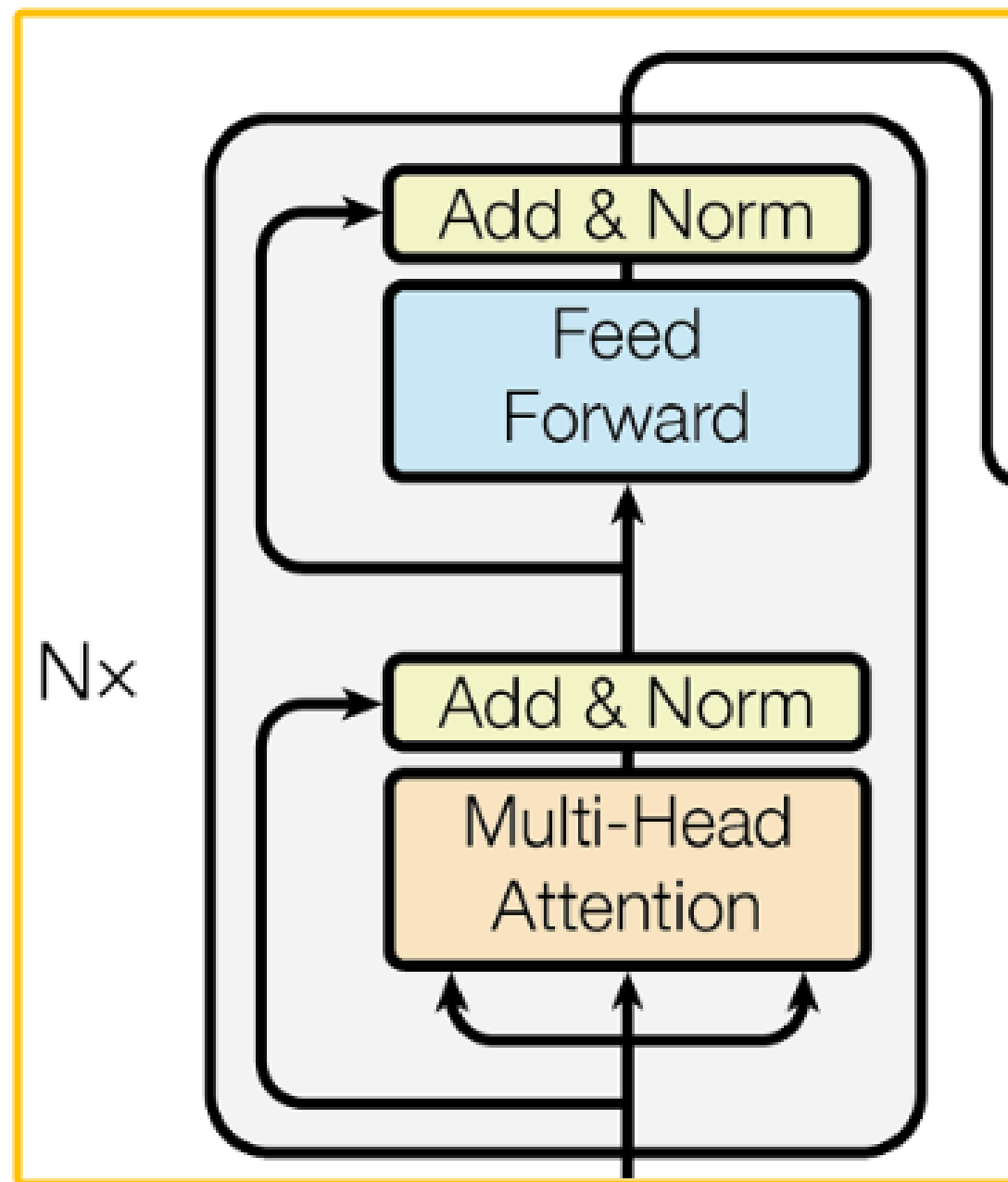


$$PE_{(pos, 2i)} = \sin(pos/10000^{2i/d_{model}})$$

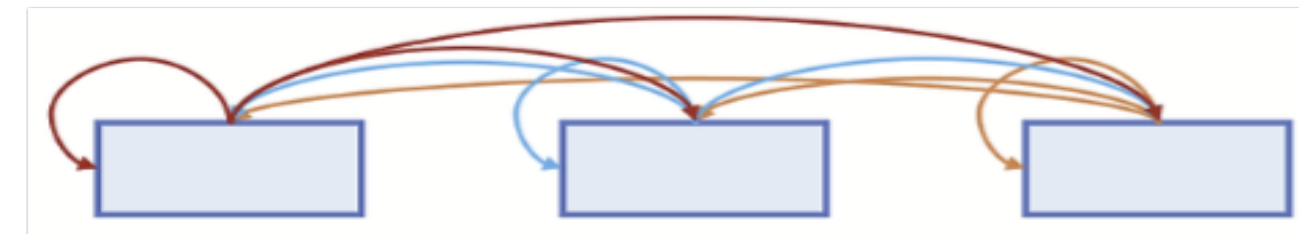
$$PE_{(pos, 2i+1)} = \cos(pos/10000^{2i/d_{model}})$$



Transformer : Encoder



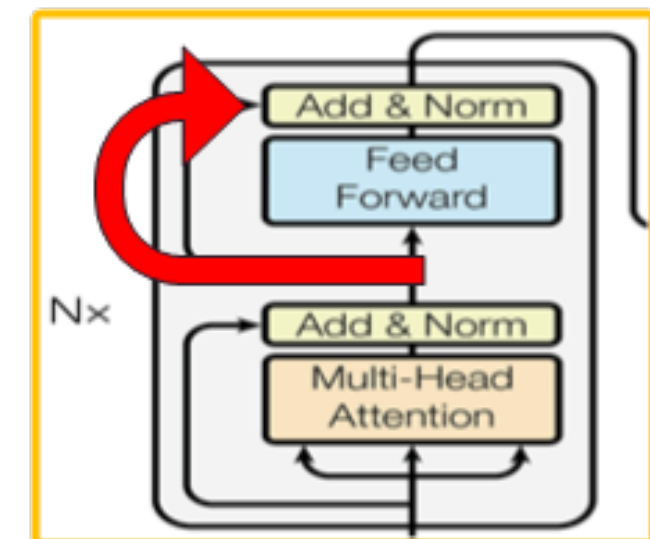
Encoder Self attention



Position-wise Feed-Forward Networks

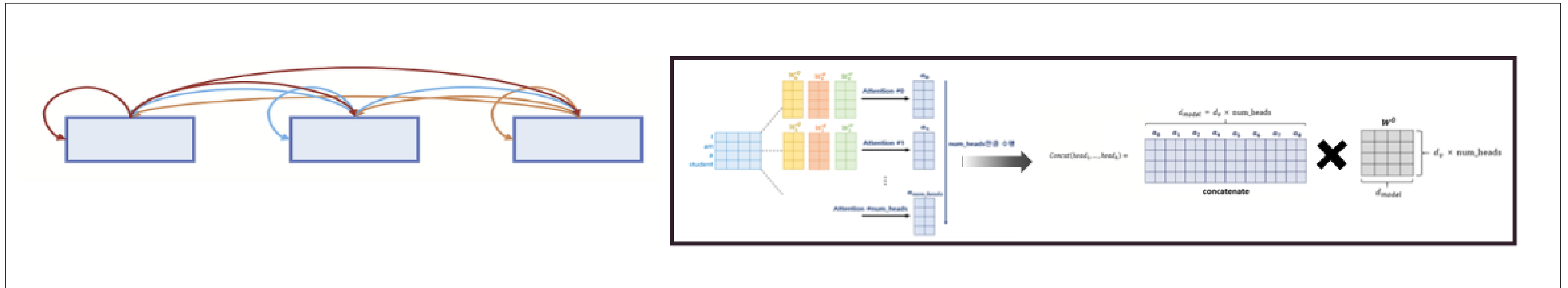


Residual connection

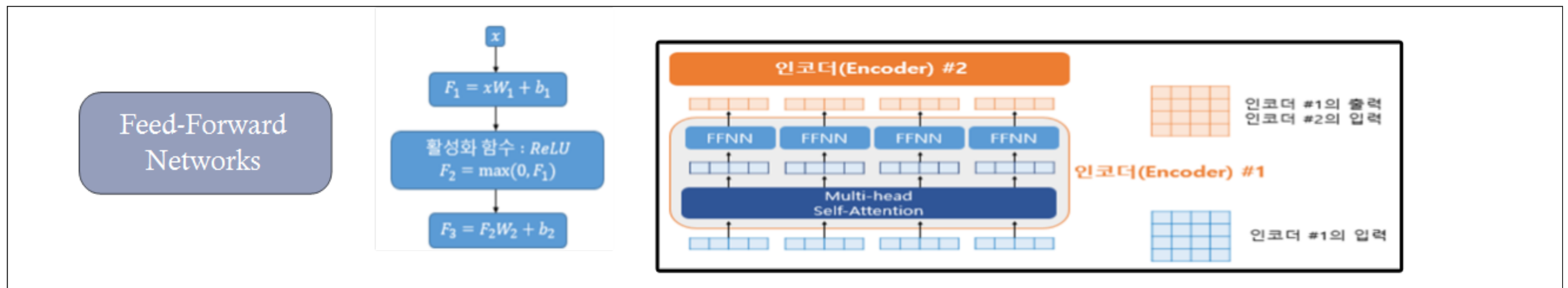


Transformer : Encoder

Encoder Self attention

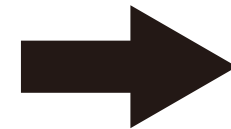
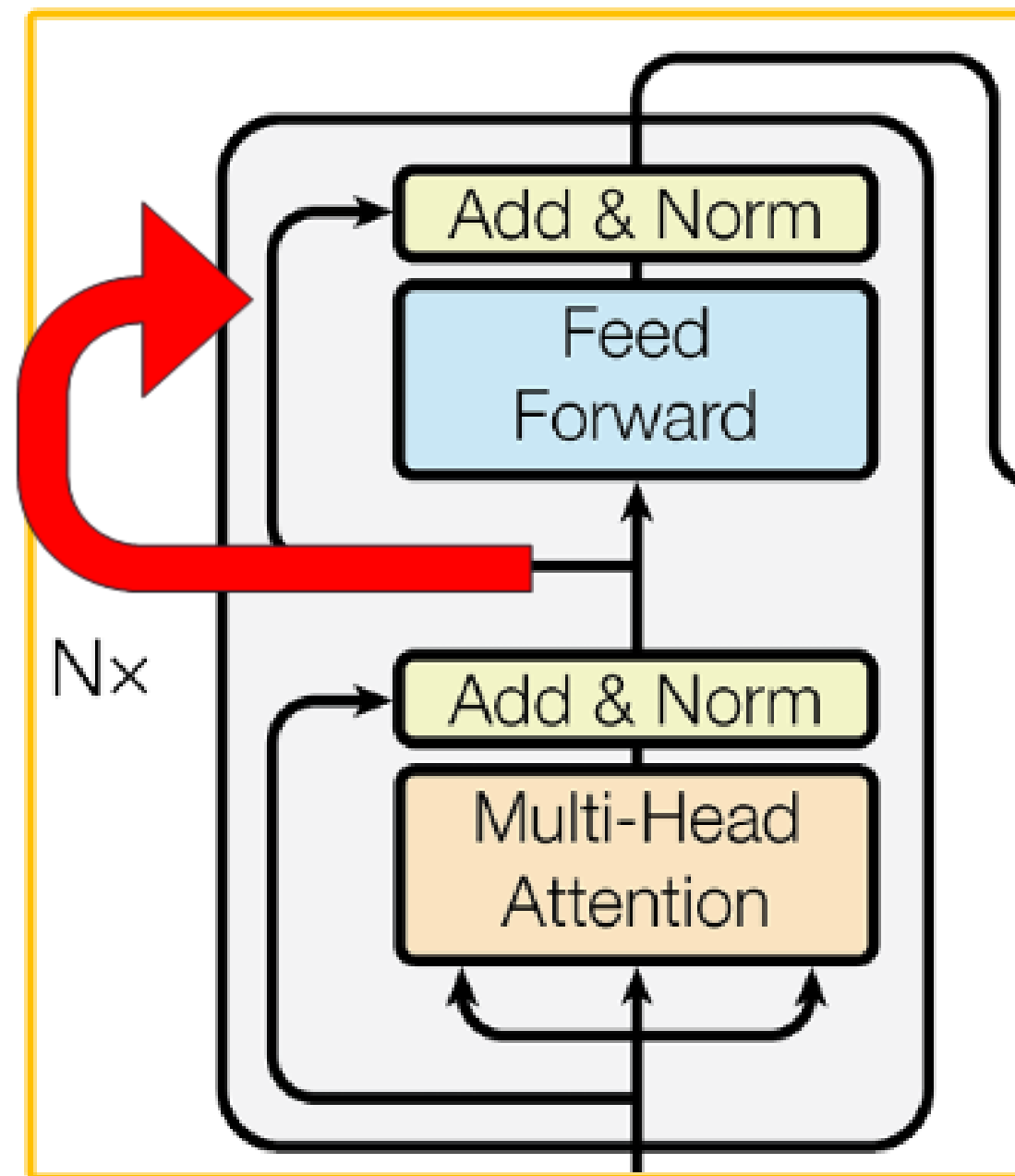


Position-wise Feed-Forward Networks

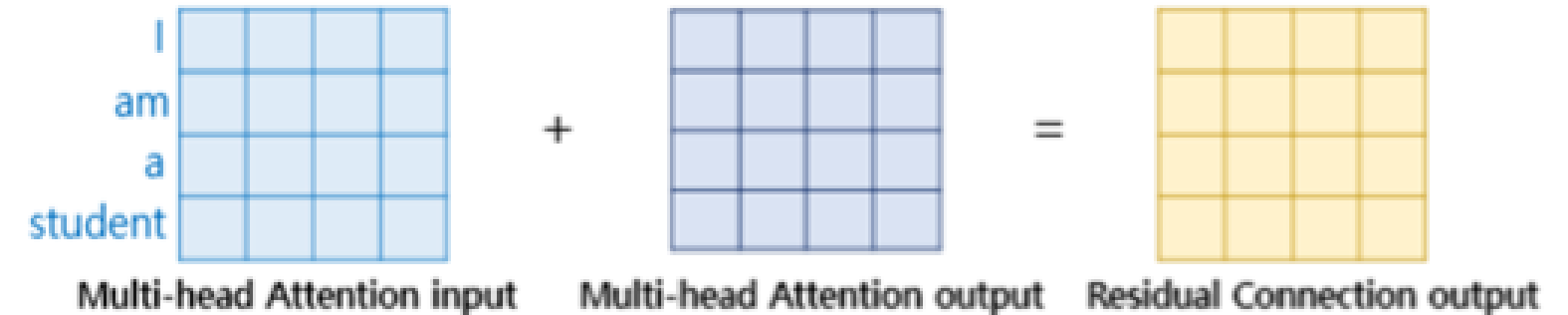


Transformer : Encoder

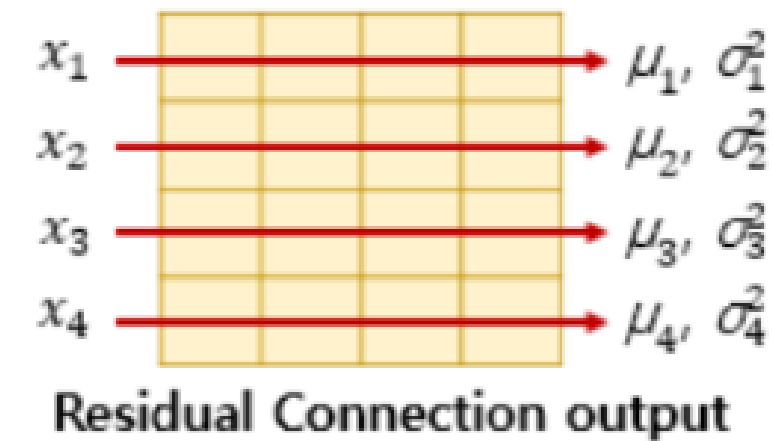
Residual connection



$$H(x) = x + \text{Multi-head Attention}(x)$$

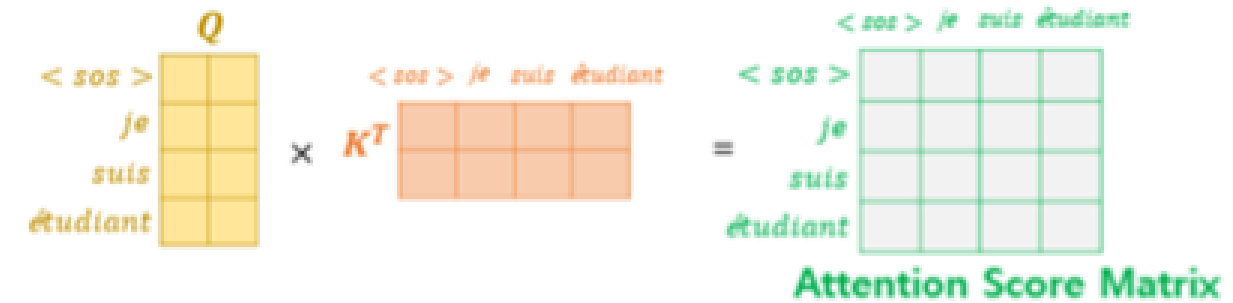
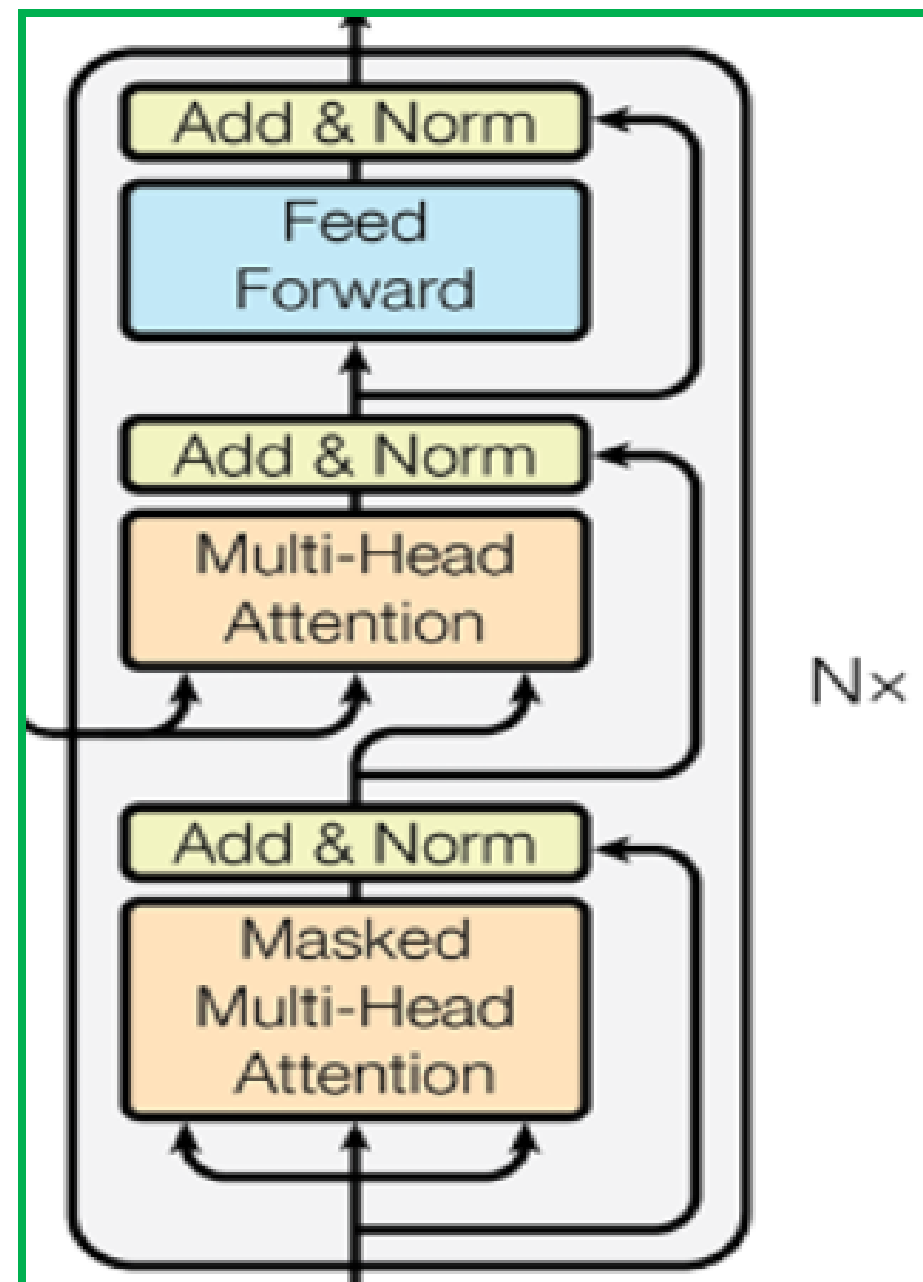


$$LN = \text{LayerNorm}(x + \text{Sublayer}(x))$$



Transformer : Decoder

Self-Attention



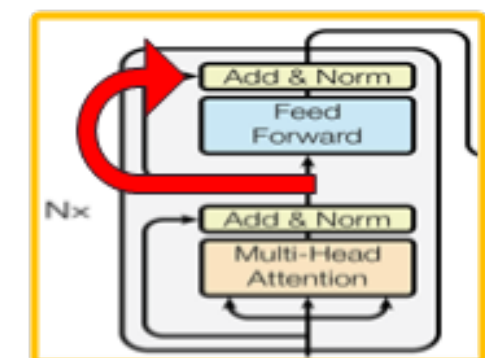
Encoder-Decoder attention & Look-Ahead Mask



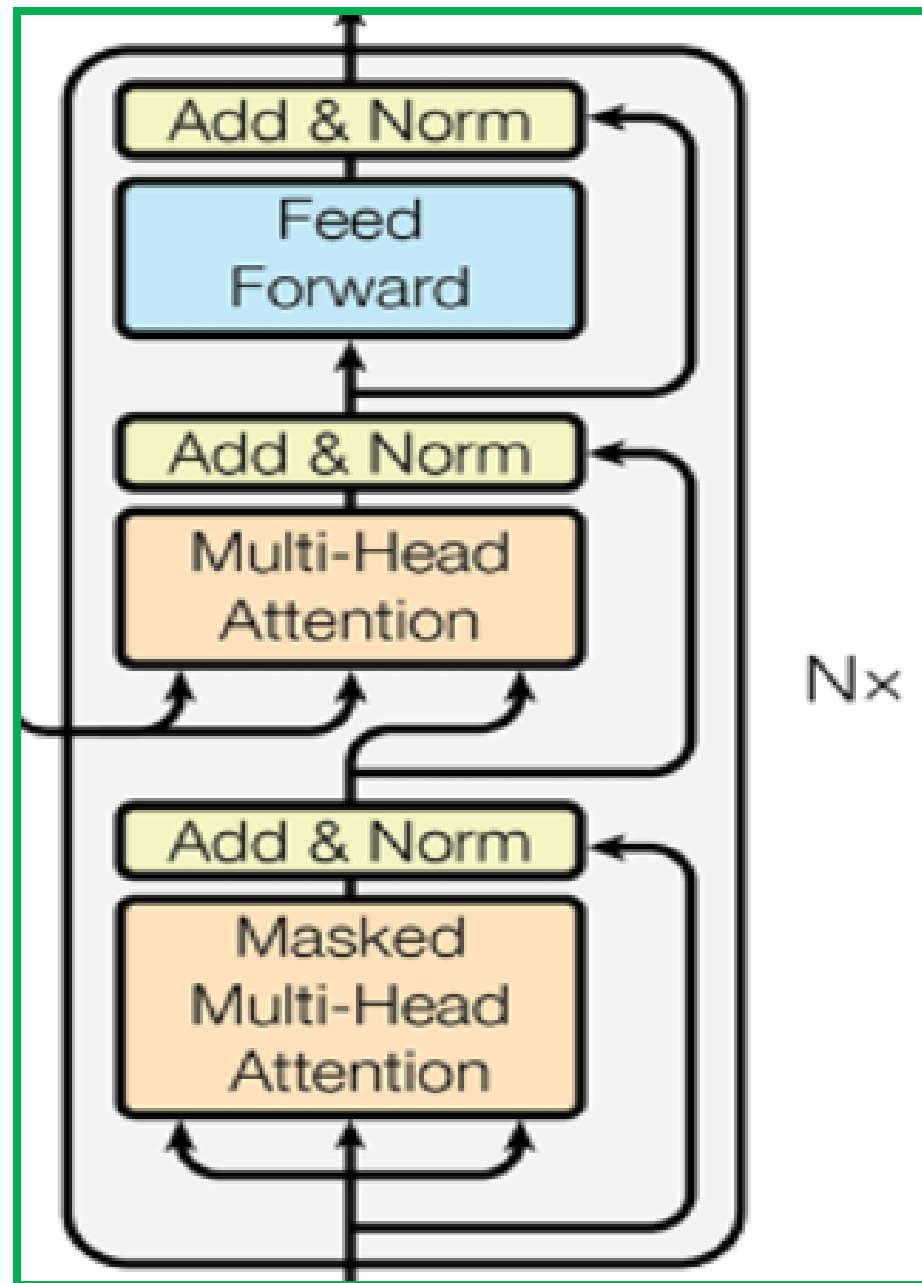
Position-wise Feed-Forward Networks

Feed-Forward Networks

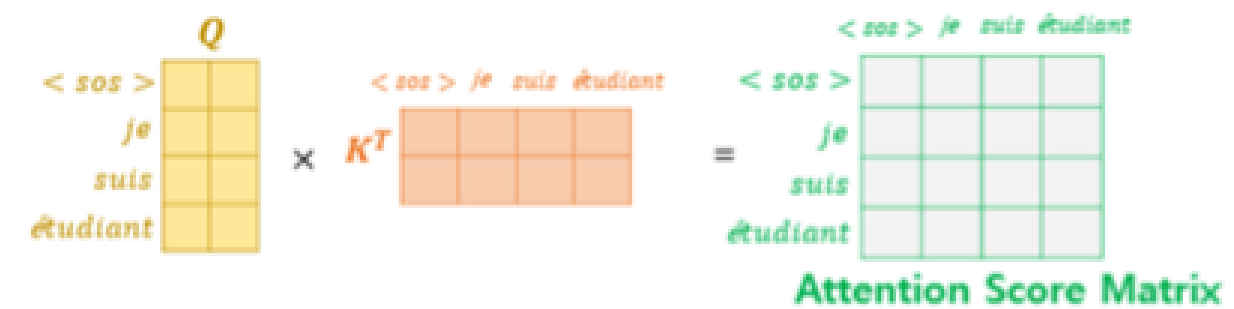
Residual connection



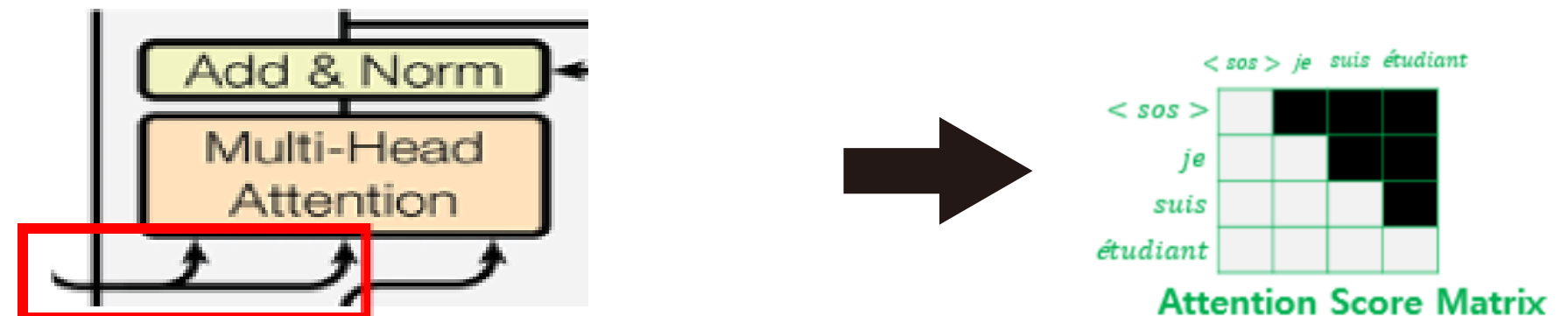
Transformer : Decoder



Self-Attention



Encoder-Decoder attention & Look-head Mask



인코더의 첫번째 서브층 : Query = Key = Value

디코더의 첫번째 서브층 : Query = Key = Value

디코더의 두번째 서브층 : Query : 디코더 행렬 / Key = value : 인코더 행렬

Results

Table 1: Maximum path lengths, per-layer complexity and minimum number of sequential operations for different layer types. n is the sequence length, d is the representation dimension, k is the kernel size of convolutions and r the size of the neighborhood in restricted self-attention.

Layer Type	Complexity per Layer	Sequential Operations	Maximum Path Length
Self-Attention	$O(n^2 \cdot d)$	$O(1)$	$O(1)$
Recurrent	$O(n \cdot d^2)$	$O(n)$	$O(n)$
Convolutional	$O(k \cdot n \cdot d^2)$	$O(1)$	$O(\log_k(n))$
Self-Attention (restricted)	$O(r \cdot n \cdot d)$	$O(1)$	$O(n/r)$

Table 2: The Transformer achieves better BLEU scores than previous state-of-the-art models on the English-to-German and English-to-French newstest2014 tests at a fraction of the training cost.

Model	BLEU		Training Cost (FLOPs)	
	EN-DE	EN-FR	EN-DE	EN-FR
ByteNet [18]	23.75			
Deep-Att + PosUnk [39]		39.2		$1.0 \cdot 10^{20}$
GNMT + RL [38]	24.6	39.92	$2.3 \cdot 10^{19}$	$1.4 \cdot 10^{20}$
ConvS2S [9]	25.16	40.46	$9.6 \cdot 10^{18}$	$1.5 \cdot 10^{20}$
MoE [32]	26.03	40.56	$2.0 \cdot 10^{19}$	$1.2 \cdot 10^{20}$
Deep-Att + PosUnk Ensemble [39]		40.4		$8.0 \cdot 10^{20}$
GNMT + RL Ensemble [38]	26.30	41.16	$1.8 \cdot 10^{20}$	$1.1 \cdot 10^{21}$
ConvS2S Ensemble [9]	26.36	41.29	$7.7 \cdot 10^{19}$	$1.2 \cdot 10^{21}$
Transformer (base model)	27.3	38.1	$3.3 \cdot 10^{18}$	
Transformer (big)	28.4	41.8	$2.3 \cdot 10^{19}$	

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}	ϵ_{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
(D)			4096							4.75	26.2	90
							0.0			5.77	24.6	
							0.2			4.95	25.5	
								0.0		4.67	25.3	
(E)								0.2		5.47	25.7	
		positional embedding instead of sinusoids								4.92	25.7	
big	6	1024	4096	16			0.3		300K	4.33	26.4	213

A : 결과를 통해 적당한 Head 개수를 가질 때 성능이 가장 좋음

B : 결과를 통해 Attention Key Size를 줄이면 모델의 성능이 크게 하락함

C, D : 결과를 통해 사이즈가 큰 모델의 성능이 좋고, Dropout은 overfitting을 완화해줌

E : 결과를 통해 Learning Positional Encoding을 사용해도 결과가 비슷함

Code : Transformer의 주요 파라미터

$d_{model} = 128$

트랜스포머의 인코더와 디코더에서의 정해진 입력과 출력의 크기를 의미

$num_layers = 4$

하나의 인코더와 디코더를 층으로 생각하였을 때, 트랜스포머 모델에서 인코더와 디코더가 총 몇 층으로 구성되었는지를 의미

$num_heads = 4$

어텐션을 사용할 때, 한 번 하는 것 보다 여러 개로 분할해서 병렬로 어텐션을 수행하고 결과값을 다시 하나로 합치는 방식 → 이 때, 병렬의 개수를 의미

$d_{ff} = 512$

트랜스포머 내부에는 피드 포워드 신경망이 존재하며 해당 신경망의 은닉층의 크기를 의미

Code : Positional Encoding

```
# 최종 버전
class PositionalEncoding(tf.keras.layers.Layer):
    def __init__(self, position, d_model):
        super(PositionalEncoding, self).__init__()
        self.pos_encoding = self.positional_encoding(position, d_model)

    def get_angles(self, position, i, d_model):
        angles = 1 / tf.pow(10000, (2 * (i // 2)) / tf.cast(d_model, tf.float32))
        return position * angles

    def positional_encoding(self, position, d_model):
        angle_rads = self.get_angles(
            position=tf.range(position, dtype=tf.float32)[:], tf.newaxis],
            i=tf.range(d_model, dtype=tf.float32)[tf.newaxis, :],
            d_model=d_model)

        # 배열의 짝수 인덱스(2i)에는 사인 함수 적용
        sines = tf.math.sin(angle_rads[:, 0::2])

        # 배열의 홀수 인덱스(2i+1)에는 코사인 함수 적용
        cosines = tf.math.cos(angle_rads[:, 1::2])

        angle_rads = np.zeros(angle_rads.shape)
        angle_rads[:, 0::2] = sines
        angle_rads[:, 1::2] = cosines
        pos_encoding = tf.constant(angle_rads)
        pos_encoding = pos_encoding[tf.newaxis, ...]

        print(pos_encoding.shape)
        return tf.cast(pos_encoding, tf.float32)

    def call(self, inputs):
        return inputs + self.pos_encoding[:, :tf.shape(inputs)[1], :]
```

초기화 (__init__)

position : 인코딩이 지원할 최대 시퀀스 길이

d_model : 임베딩 벡터의 차원 및 레이어의 예상 출력 차원

get_angles

위치 인코딩에 사용되는 위치 계산

positional_encoding

짝수 인덱스에는 사인 함수를 적용, 홀수에는 코사인 함수를 적용
마지막으로, pos_encoding을 텐서로 변환하여 반환

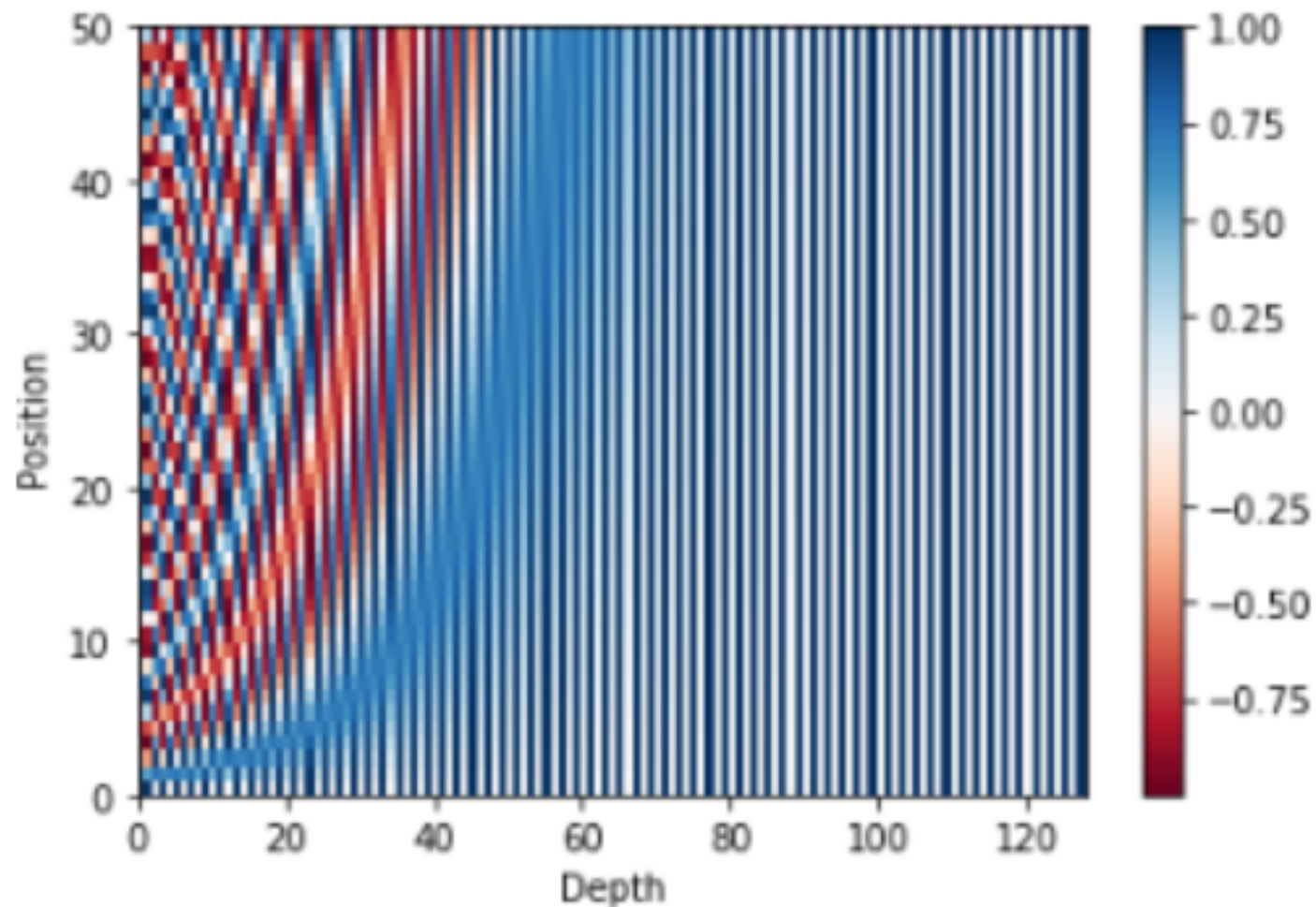
call

위치 인코딩은 입력의 길이까지만 추가,
이는 tf.shape(inputs)[1]로 표시

Code : Positional Encoding

```
sample_pos_encoding = PositionalEncoding(50, 128)

plt.pcolormesh(sample_pos_encoding.pos_encoding.numpy()[0], cmap='RdBu')
plt.xlabel('Depth')
plt.xlim((0, 128))
plt.ylabel('Position')
plt.colorbar()
plt.show()
```



50 × 128의 크기를 가지는 포지셔널 인코딩 행렬을 시각화하여 어떤 형태를 가지는지 확인.

이는 입력 문장의 단어가 50개이면서, 각 단어가 128차원의 임베딩 벡터를 가질 때 사용할 수 있는 행렬

Code : Scaled_Dot_Product_Attention

```
def scaled_dot_product_attention(query, key, value, mask):
    # query 크기 : (batch_size, num_heads, query의 문장 길이, d_model/num_heads)
    # key 크기 : (batch_size, num_heads, key의 문장 길이, d_model/num_heads)
    # value 크기 : (batch_size, num_heads, value의 문장 길이, d_model/num_heads)
    # padding_mask : (batch_size, 1, 1, key의 문장 길이)

    # Q와 K의 곱. 어텐션 스코어 행렬.
    matmul_qk = tf.matmul(query, key, transpose_b=True)

    # 스케일링
    # dk의 루트값으로 나눠준다.
    depth = tf.cast(tf.shape(key)[-1], tf.float32)
    logits = matmul_qk / tf.math.sqrt(depth)

    # 마스킹. 어텐션 스코어 행렬의 마스킹 할 위치에 매우 작은 음수값을 넣는다.
    # 매우 작은 값이므로 소프트맥스 함수를 지나면 행렬의 해당 위치의 값은 0이 된다.
    if mask is not None:
        logits += (mask * -1e9)

    # 소프트맥스 함수는 마지막 차원인 key의 문장 길이 방향으로 수행된다.
    # attention weight : (batch_size, num_heads, query의 문장 길이, key의 문장 길이)
    attention_weights = tf.nn.softmax(logits, axis=-1)

    # output : (batch_size, num_heads, query의 문장 길이, d_model/num_heads)
    output = tf.matmul(attention_weights, value)

    return output, attention_weights
```

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

tf.shape(key)[-1]은 key 텐서의 마지막 차원의 크기를 반환
tf.cast → tf.shape(key)[-1]에서 가져온 마지막 차원의 크기를
tf.float32로 변환

softmax

- mask : 뒤에 자세히 설명
- attention_weights : softmax를 사용하여 확률값을 매김

Code : Multi-Head Attention

```
class MultiHeadAttention(tf.keras.layers.Layer):

    def __init__(self, d_model, num_heads, name="multi_head_attention"):
        super(MultiHeadAttention, self).__init__(name=name)
        self.num_heads = num_heads
        self.d_model = d_model

        assert d_model % self.num_heads == 0

        # d_model을 num_heads로 나눈 값.
        # 논문 기준 : 64
        self.depth = d_model // self.num_heads

        # WQ, WK, WV에 해당하는 밀집층 정의
        self.query_dense = tf.keras.layers.Dense(units=d_model)
        self.key_dense = tf.keras.layers.Dense(units=d_model)
        self.value_dense = tf.keras.layers.Dense(units=d_model)

        # WO에 해당하는 밀집층 정의
        self.dense = tf.keras.layers.Dense(units=d_model)

        # num_heads 개수만큼 q, k, v를 split하는 함수
    def split_heads(self, inputs, batch_size):
        inputs = tf.reshape(
            inputs, shape=(batch_size, -1, self.num_heads, self.depth))
        return tf.transpose(inputs, perm=[0, 2, 1, 3])

    def call(self, inputs):
        query, key, value, mask = inputs['query'], inputs['key'], inputs[
            'value'], inputs['mask']
        batch_size = tf.shape(query)[0]
```

```
        # 1. WQ, WK, WV에 해당하는 밀집층 지나기
        # q : (batch_size, query의 문장 길이, d_model)
        # k : (batch_size, key의 문장 길이, d_model)
        # v : (batch_size, value의 문장 길이, d_model)
        # 참고) 인코더(k, v)-디코더(q) 어텐션에서는 query 길이와 key, value의 길이는 다를 수 있다.
        query = self.query_dense(query)
        key = self.key_dense(key)
        value = self.value_dense(value)

        # 2. 헤드 나누기
        # q : (batch_size, num_heads, query의 문장 길이, d_model/num_heads)
        # k : (batch_size, num_heads, key의 문장 길이, d_model/num_heads)
        # v : (batch_size, num_heads, value의 문장 길이, d_model/num_heads)
        query = self.split_heads(query, batch_size)
        key = self.split_heads(key, batch_size)
        value = self.split_heads(value, batch_size)

        # 3. 스케일드 닷 프로덕트 어텐션. 앞서 구현한 함수 사용.
        # (batch_size, num_heads, query의 문장 길이, d_model/num_heads)
        scaled_attention, _ = scaled_dot_product_attention(query, key, value, mask)
        # (batch_size, query의 문장 길이, num_heads, d_model/num_heads)
        scaled_attention = tf.transpose(scaled_attention, perm=[0, 2, 1, 3])

        # 4. 헤드 연결(concatenate)하기
        # (batch_size, query의 문장 길이, d_model)
        concat_attention = tf.reshape(scaled_attention,
                                       (batch_size, -1, self.d_model))

        # 5. WO에 해당하는 밀집층 지나기
        # (batch_size, query의 문장 길이, d_model)
        outputs = self.dense(concat_attention)

        return outputs
```

Code : 패딩 마스크(Padding Mask)

```
# 마스크. 어텐션 스코어 행렬의 마스크 할 위치에 매우 작은 음수값을 넣는다.  
# 매우 작은 값이므로 소프트맥스 함수를 지나면 행렬의 해당 위치의 값은 0이 된다.  
if mask is not None:  
    logits += (mask * -1e9)
```

mask라는 값을 인자로 받아서, 이 mask값에다가
-1e9라는 아주 작은 음수값을 곱한 후
어텐션 스코어 행렬에 더해주는 연산



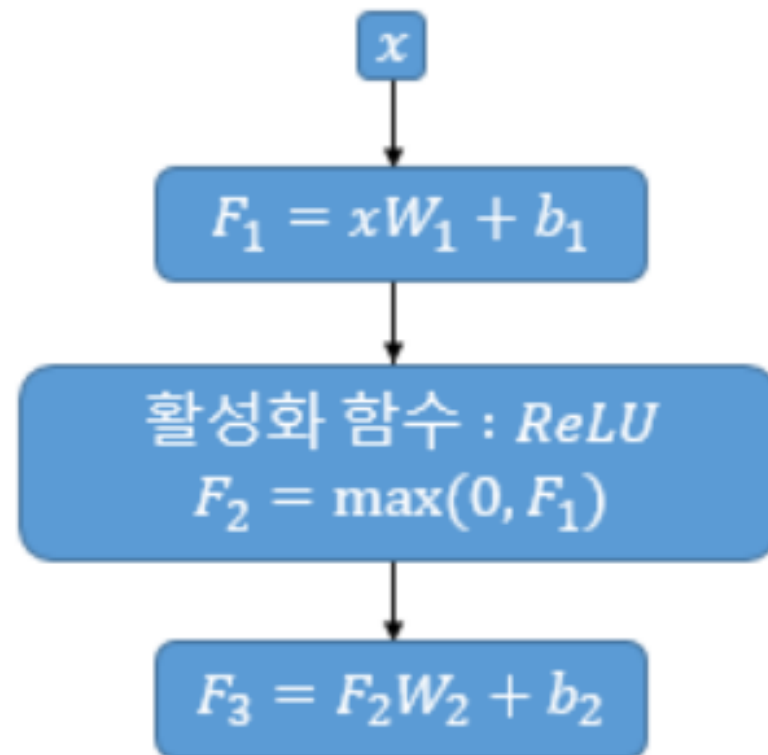
이는 입력 문장에 <PAD> 토큰이 있을 경우
어텐션에서 사실상 제외하기 위한 연산

PAD?

<PAD>는 자연어 처리에서 패딩(padding) 토큰을
나타

패딩 토큰을 무시하고 실제 단어에만 주의를 기울이
도록 모델을 학습시키는 것이 일반적

Code : Position-wise FFNN



$$FFNN(x) = MAX(0, xW_1 + b_1)W_2 + b_2$$

각 포지션에 대한 정보를 고려하여 피드 포워드
신경망이 적용됨

비선형성을 도입하여 입력에 대한 복잡한 비선형 변환을 수행.
이는 모델이 더욱 복잡한 패턴을 학습할 수 있도록 도움을 줌

```
outputs = tf.keras.layers.Dense(units=d_ff, activation='relu')(attention)
outputs = tf.keras.layers.Dense(units=d_model)(outputs)
```

다음의 코드는 인코더와 디코더 내부에서 사용할 예정

Code : Encoder 구현

```
def encoder_layer(dff, d_model, num_heads, dropout, name="encoder_layer"):
    inputs = tf.keras.Input(shape=(None, d_model), name="inputs")

    # 인코더는 패딩 마스크 사용
    padding_mask = tf.keras.Input(shape=(1, 1, None), name="padding_mask")

    # 멀티-헤드 어텐션 (첫번째 서브층 / 셀프 어텐션)
    attention = MultiHeadAttention(
        d_model, num_heads, name="attention")({
            'query': inputs, 'key': inputs, 'value': inputs, # Q = K = V
            'mask': padding_mask # 패딩 마스크 사용
        })

    # 드롭아웃 + 잔차 연결과 층 정규화
    attention = tf.keras.layers.Dropout(rate=dropout)(attention)
    attention = tf.keras.layers.LayerNormalization(
        epsilon=1e-6)(inputs + attention)

    # 포지션 와이즈 피드 포워드 신경망 (두번째 서브층)
    outputs = tf.keras.layers.Dense(units=dff, activation='relu')(attention)
    outputs = tf.keras.layers.Dense(units=d_model)(outputs)

    # 드롭아웃 + 잔차 연결과 층 정규화
    outputs = tf.keras.layers.Dropout(rate=dropout)(outputs)
    outputs = tf.keras.layers.LayerNormalization(
        epsilon=1e-6)(attention + outputs)

    return tf.keras.Model(
        inputs=[inputs, padding_mask], outputs=outputs, name=name)
```

인코더의 입력으로 들어가는 문장에는 패딩이 있을 수 있으므로,
어텐션 시 패딩 토큰을 제외하도록 패딩 마스크를 사용.

인코더는 총 두 개의 서브층으로 이루어짐

→ 멀티 헤드 어텐션 & 피드 포워드 신경망

각 서브층 이후에는 드롭 아웃, 잔차 연결과 층 정규화 수행

Code : Encoder 쌓기

```
def encoder(vocab_size, num_layers, dff,
            d_model, num_heads, dropout,
            name="encoder"):
    inputs = tf.keras.Input(shape=(None,), name="inputs")

    # 인코더는 패딩 마스크 사용
    padding_mask = tf.keras.Input(shape=(1, 1, None), name="padding_mask")

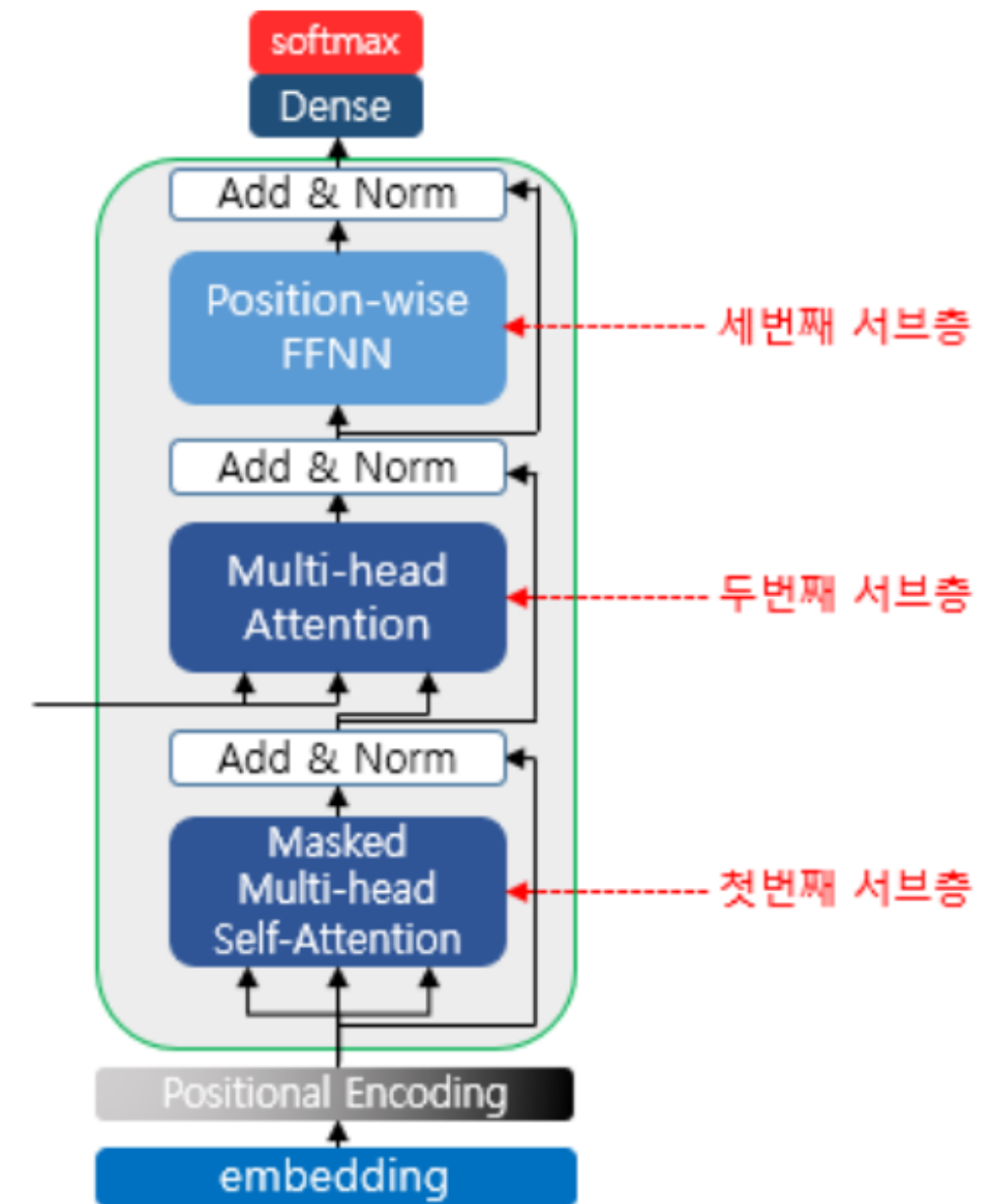
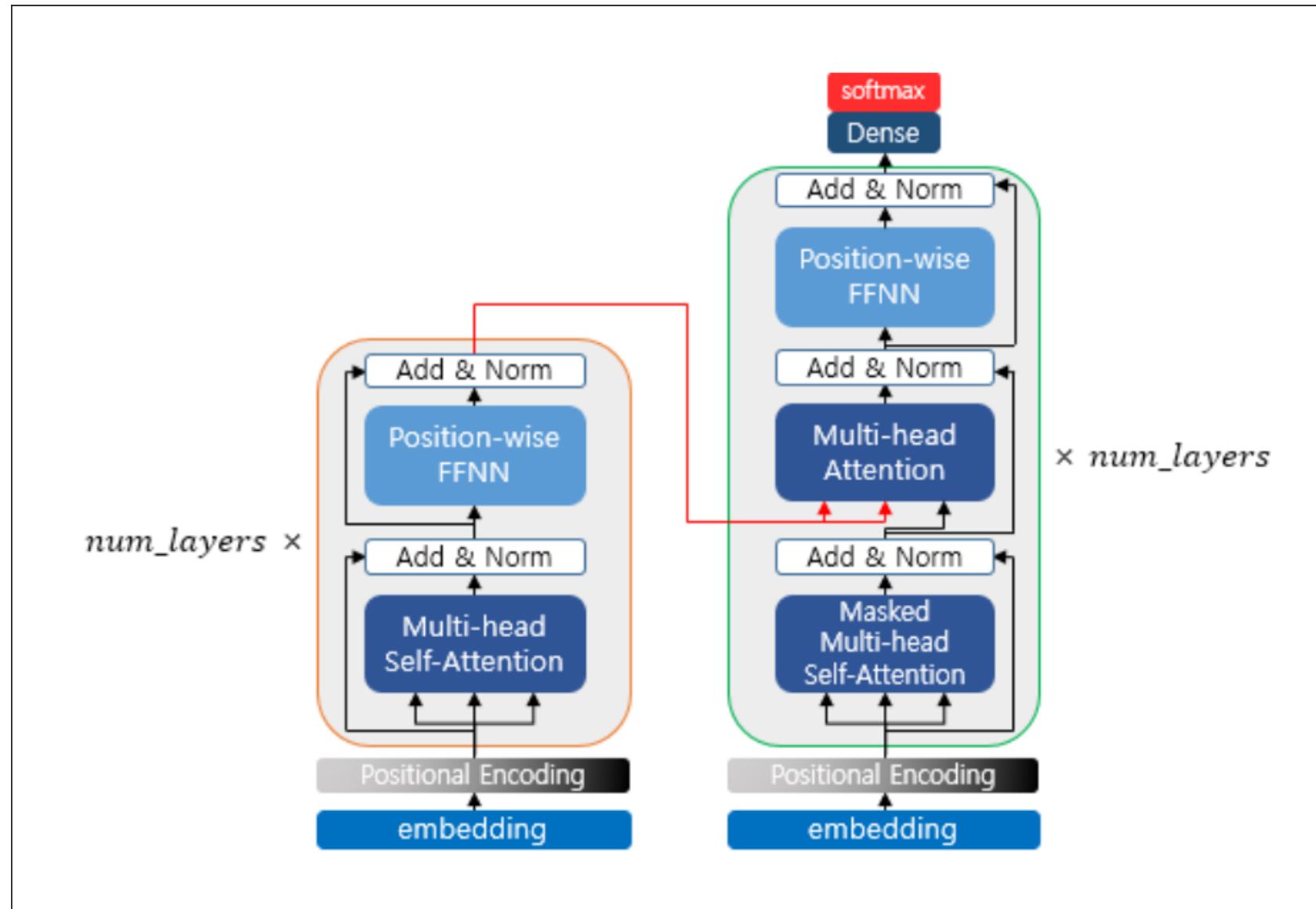
    # 포지셔널 인코딩 + 드롭아웃
    embeddings = tf.keras.layers.Embedding(vocab_size, d_model)(inputs)
    embeddings *= tf.math.sqrt(tf.cast(d_model, tf.float32))
    embeddings = PositionalEncoding(vocab_size, d_model)(embeddings)
    outputs = tf.keras.layers.Dropout(rate=dropout)(embeddings)

    # 인코더를 num_layers개 쌓기
    for i in range(num_layers):
        outputs = encoder_layer(dff=dff, d_model=d_model, num_heads=num_heads,
                                dropout=dropout, name="encoder_layer_{}".format(i),
                                )([outputs, padding_mask])

    return tf.keras.Model(
        inputs=[inputs, padding_mask], outputs=outputs, name=name)
```

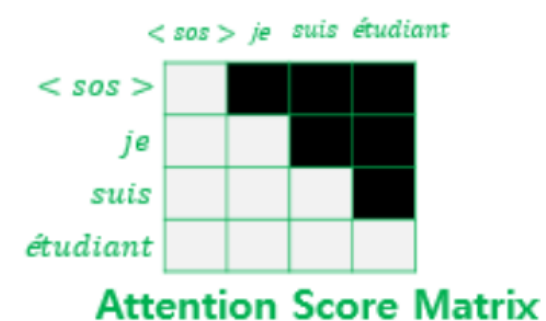
인코더 층을 num_layers개만큼 쌓는 코드로,
마지막 인코더 층에서 얻는 (seq_len, d_model) 크기의 행렬을
디코더로 보내주므로써 트랜스포머 인코더의 인코딩 연산이
끝나게 됨.

Code : From Encoder To Decoder



look-ahead mask

미리보기에 대한 마스크



Code : Decoder 구현

```
def decoder_layer(dff, d_model, num_heads, dropout, name="decoder_layer"):
    inputs = tf.keras.Input(shape=(None, d_model), name="inputs")
    enc_outputs = tf.keras.Input(shape=(None, d_model), name="encoder_outputs")

    # 디코더는 록어헤드 마스크(첫번째 서브층)와 패딩 마스크(두번째 서브층) 둘 다 사용.
    look_ahead_mask = tf.keras.Input(
        shape=(1, None, None), name="look_ahead_mask")
    padding_mask = tf.keras.Input(shape=(1, 1, None), name='padding_mask')

    # 멀티-헤드 어텐션 (첫번째 서브층 / 마스크드 셀프 어텐션)
    attention1 = MultiHeadAttention(
        d_model, num_heads, name="attention_1")(inputs={
            'query': inputs, 'key': inputs, 'value': inputs, # Q = K = V
            'mask': look_ahead_mask # 록어헤드 마스크
        })

    # 잔차 연결과 층 정규화
    attention1 = tf.keras.layers.LayerNormalization(
        epsilon=1e-6)(attention1 + inputs)

    # 멀티-헤드 어텐션 (두번째 서브층 / 디코더-인코더 어텐션)
    attention2 = MultiHeadAttention(
        d_model, num_heads, name="attention_2")(inputs={
            'query': attention1, 'key': enc_outputs, 'value': enc_outputs, # Q != K = V
            'mask': padding_mask # 패딩 마스크
        })
```

```
# 드롭아웃 + 잔차 연결과 층 정규화
attention2 = tf.keras.layers.Dropout(rate=dropout)(attention2)
attention2 = tf.keras.layers.LayerNormalization(
    epsilon=1e-6)(attention2 + attention1)

# 포지션 와이즈 피드 포워드 신경망 (세번째 서브층)
outputs = tf.keras.layers.Dense(units=dff, activation='relu')(attention2)
outputs = tf.keras.layers.Dense(units=d_model)(outputs)

# 드롭아웃 + 잔차 연결과 층 정규화
outputs = tf.keras.layers.Dropout(rate=dropout)(outputs)
outputs = tf.keras.layers.LayerNormalization(
    epsilon=1e-6)(outputs + attention2)

return tf.keras.Model(
    inputs=[inputs, enc_outputs, look_ahead_mask, padding_mask],
    outputs=outputs,
    name=name)
```

Code : Decoder 쌓기

```
def decoder(vocab_size, num_layers, dff,
            d_model, num_heads, dropout,
            name='decoder'):
    inputs = tf.keras.Input(shape=(None,), name='inputs')
    enc_outputs = tf.keras.Input(shape=(None, d_model), name='encoder_outputs')

    # 디코더는 록어헤드 마스크(첫번째 서브층)와 패딩 마스크(두번째 서브층) 둘 다 사용.
    look_ahead_mask = tf.keras.Input(
        shape=(1, None, None), name='look_ahead_mask')
    padding_mask = tf.keras.Input(shape=(1, 1, None), name='padding_mask')

    # 포지셔널 인코딩 + 드롭아웃
    embeddings = tf.keras.layers.Embedding(vocab_size, d_model)(inputs)
    embeddings *= tf.math.sqrt(tf.cast(d_model, tf.float32))
    embeddings = PositionalEncoding(vocab_size, d_model)(embeddings)
    outputs = tf.keras.layers.Dropout(rate=dropout)(embeddings)

    # 디코더를 num_layers개 쌓기
    for i in range(num_layers):
        outputs = decoder_layer(dff=dff, d_model=d_model, num_heads=num_heads,
                                dropout=dropout, name='decoder_layer_{}'.format(i),
                                )(inputs=[outputs, enc_outputs, look_ahead_mask, padding_mask])

    return tf.keras.Model(
        inputs=[inputs, enc_outputs, look_ahead_mask, padding_mask],
        outputs=outputs,
        name=name)
```

포지셔널 인코딩 후 디코더 층을 num_layers의
개수만큼 쌓는 코드

Code : Transformer 전체 Code

Hyperparameter in code

```
def transformer(vocab_size, num_layers, dff,
               d_model, num_heads, dropout,
               name="transformer"):

    # 인코더의 입력
    inputs = tf.keras.Input(shape=(None,), name='inputs')

    # 디코더의 입력
    dec_inputs = tf.keras.Input(shape=(None,), name='dec_inputs')

    # 인코더의 패딩 마스크
    enc_padding_mask = tf.keras.layers.Lambda(
        create_padding_mask, output_shape=(1, 1, None),
        name='enc_padding_mask')(inputs)

    # 디코더의 록어헤드 마스크(첫번째 서브층)
    look_ahead_mask = tf.keras.layers.Lambda(
        create_look_ahead_mask, output_shape=(1, None, None),
        name='look_ahead_mask')(dec_inputs)

    # 디코더의 패딩 마스크(두번째 서브층)
    dec_padding_mask = tf.keras.layers.Lambda(
        create_padding_mask, output_shape=(1, 1, None),
        name='dec_padding_mask')(inputs)

    # 인코더의 출력은 enc_outputs, 디코더로 전달된다.
    enc_outputs = encoder(vocab_size=vocab_size, num_layers=num_layers, dff=dff,
                          d_model=d_model, num_heads=num_heads, dropout=dropout,
                          )(inputs=[inputs, enc_padding_mask]) # 인코더의 입력은 입력 문장과 패딩 마스크

    # 디코더의 출력은 dec_outputs, 출력층으로 전달된다.
    dec_outputs = decoder(vocab_size=vocab_size, num_layers=num_layers, dff=dff,
                          d_model=d_model, num_heads=num_heads, dropout=dropout,
                          )(inputs=[dec_inputs, enc_outputs, look_ahead_mask, dec_padding_mask])

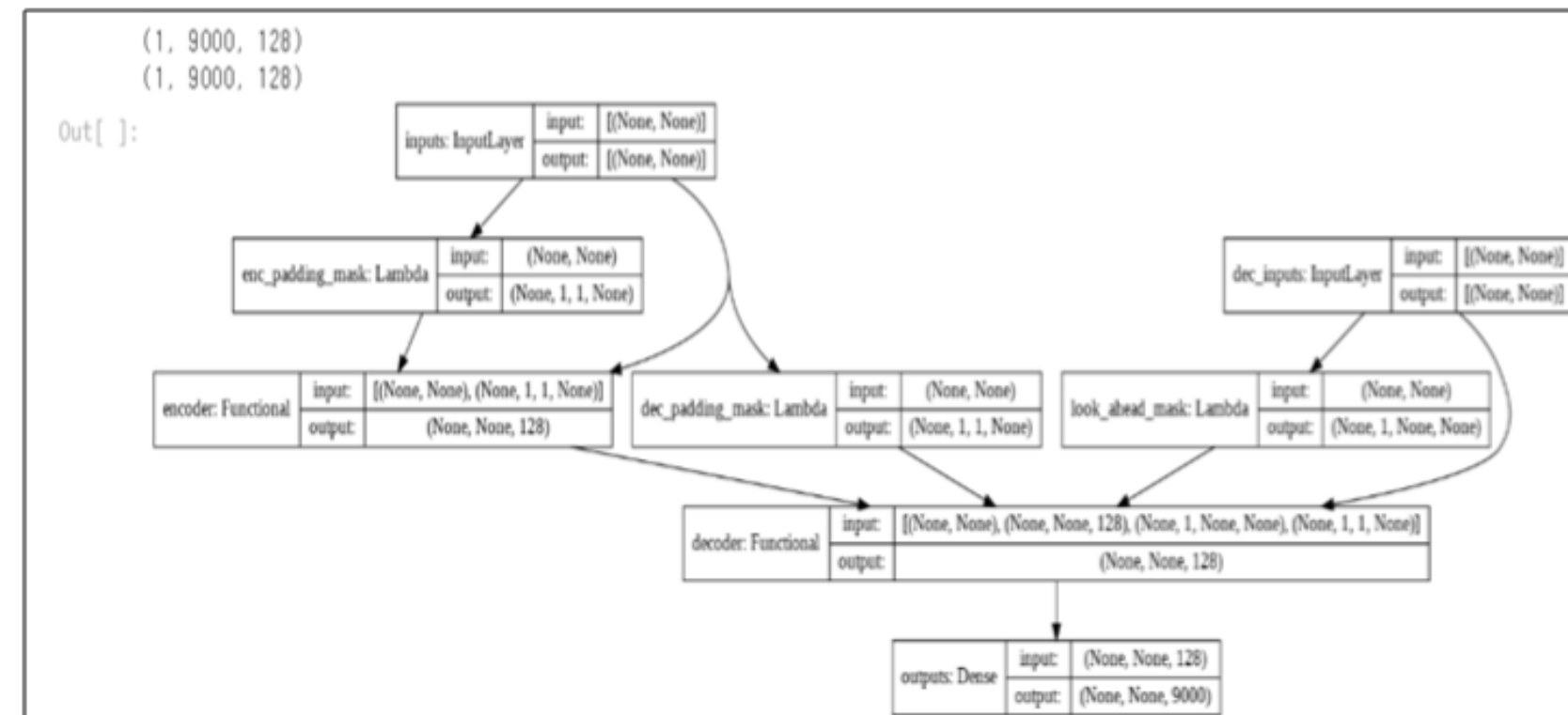
    # 다음 단어 예측을 위한 출력층
    outputs = tf.keras.layers.Dense(units=vocab_size, name="outputs")(dec_outputs)

    return tf.keras.Model(inputs=[inputs, dec_inputs], outputs=outputs, name=name)
```

```
small_transformer = transformer(
    vocab_size = 9000,
    num_layers = 4,
    dff = 512,
    d_model = 128,
    num_heads = 4,
    dropout = 0.3,
    name="small_transformer")

tf.keras.utils.plot_model(
    small_transformer, to_file='small_transformer.png', show_shapes=True)
```

Output



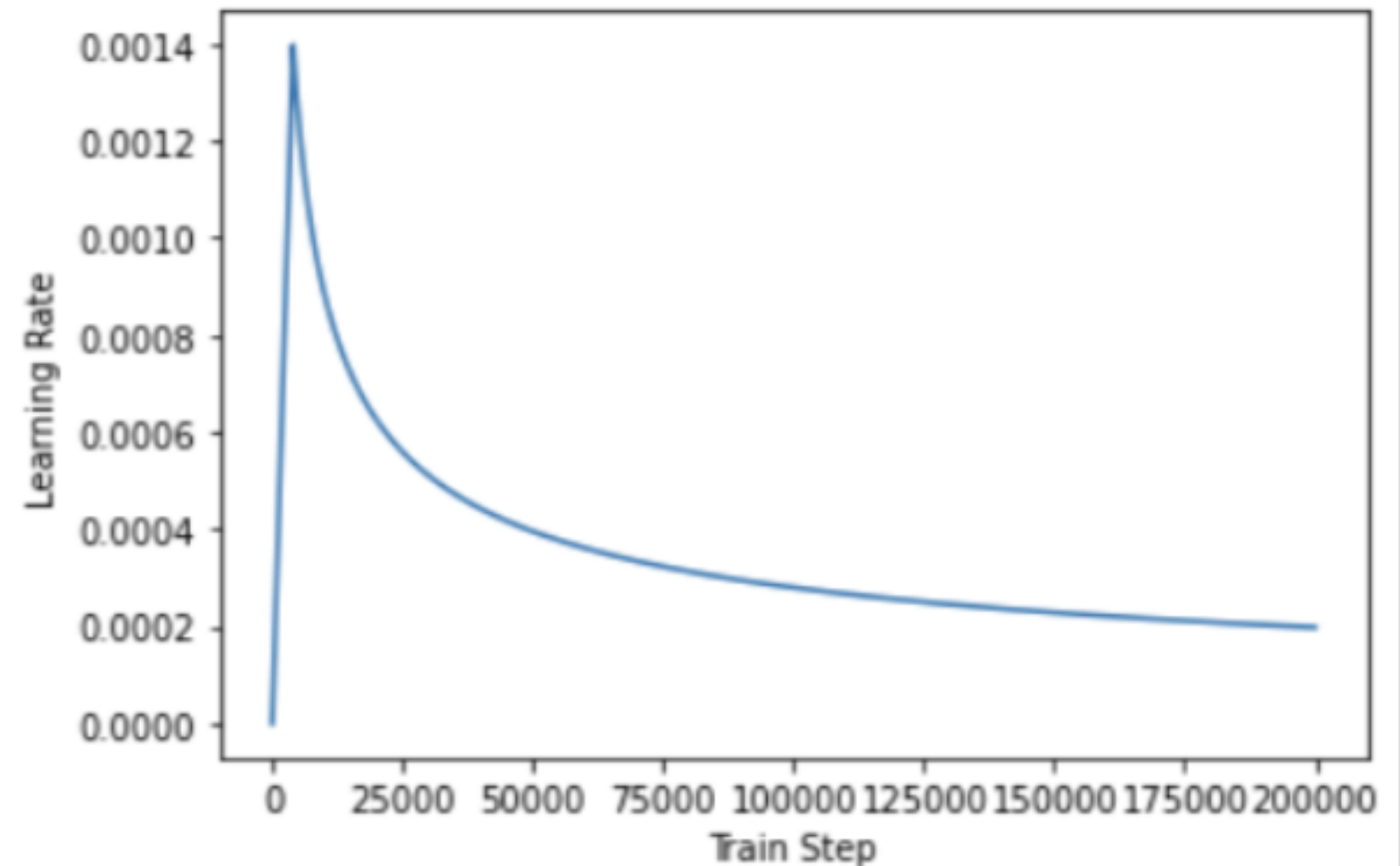
Code : Loss function definition and Learning rate scheduler

```
def loss_function(y_true, y_pred):  
    y_true = tf.reshape(y_true, shape=(-1, MAX_LENGTH - 1))  
  
    loss = tf.keras.losses.SparseCategoricalCrossentropy(  
        from_logits=True, reduction='none')(y_true, y_pred)  
  
    mask = tf.cast(tf.not_equal(y_true, 0), tf.float32)  
    loss = tf.multiply(loss, mask)  
  
    return tf.reduce_mean(loss)
```

```
class CustomSchedule(tf.keras.optimizers.schedules.LearningRateSchedule):  
  
    def __init__(self, d_model, warmup_steps=4000):  
        super(CustomSchedule, self).__init__()  
        self.d_model = d_model  
        self.d_model = tf.cast(self.d_model, tf.float32)  
        self.warmup_steps = warmup_steps  
  
    def __call__(self, step):  
        arg1 = tf.math.rsqrt(step)  
        arg2 = step * (self.warmup_steps** -1.5)  
  
        return tf.math.rsqrt(self.d_model) * tf.math.minimum(arg1, arg2)
```

```
sample_learning_rate = CustomSchedule(d_model=128)  
  
plt.plot(sample_learning_rate(tf.range(200000, dtype=tf.float32)))  
plt.ylabel("Learning Rate")  
plt.xlabel("Train Step")
```

Out[]: Text(0.5, 0, 'Train Step')



Reference

https://github.com/ukairia777/tensorflow-transformer/blob/c55a78a26ad970a3298ea388f06c407e6a69d955/Transformer_Korean_Chatbot.ipynb

<https://wikidocs.net/31379>

<https://www.youtube.com/watch?v=AA621UofTUA&t=2623s>