

Distinguished Engineering

Transformer 4/5

- Transformer, revenge of the fallen

...

BW

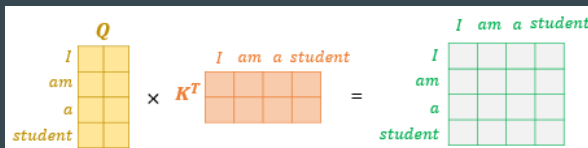
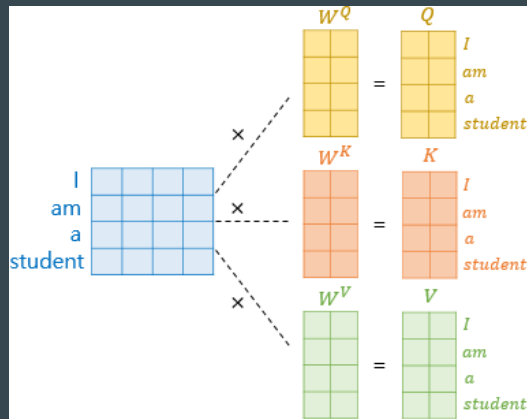
Plan

- Prologue, seq2seq
- Attention, please
- Transformer, a new hope
- Transformer, revenge of the fallen
- Transformer, vision

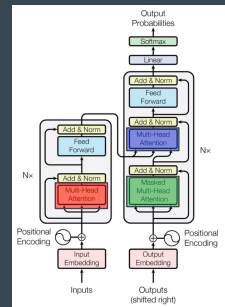
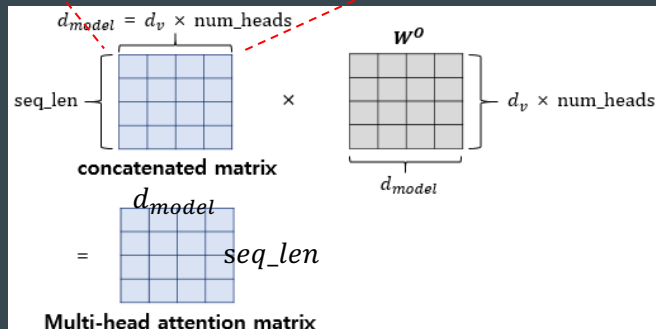
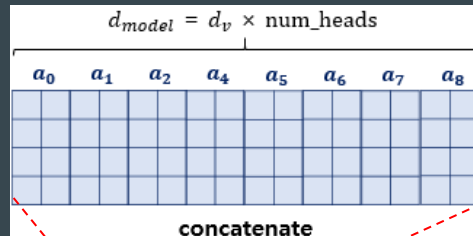
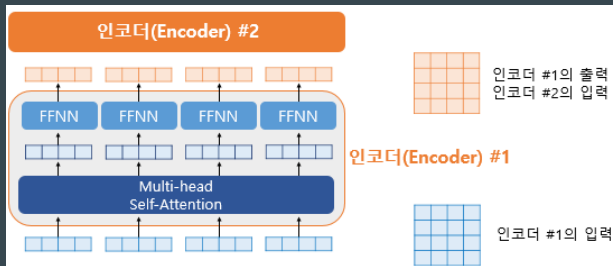
Transformer

- Multi-head Attention

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

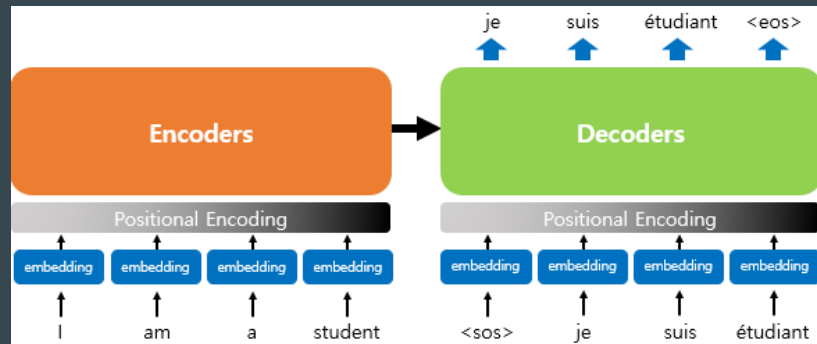


$$\text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right) \times V = \text{Attention Value Matrix } a$$



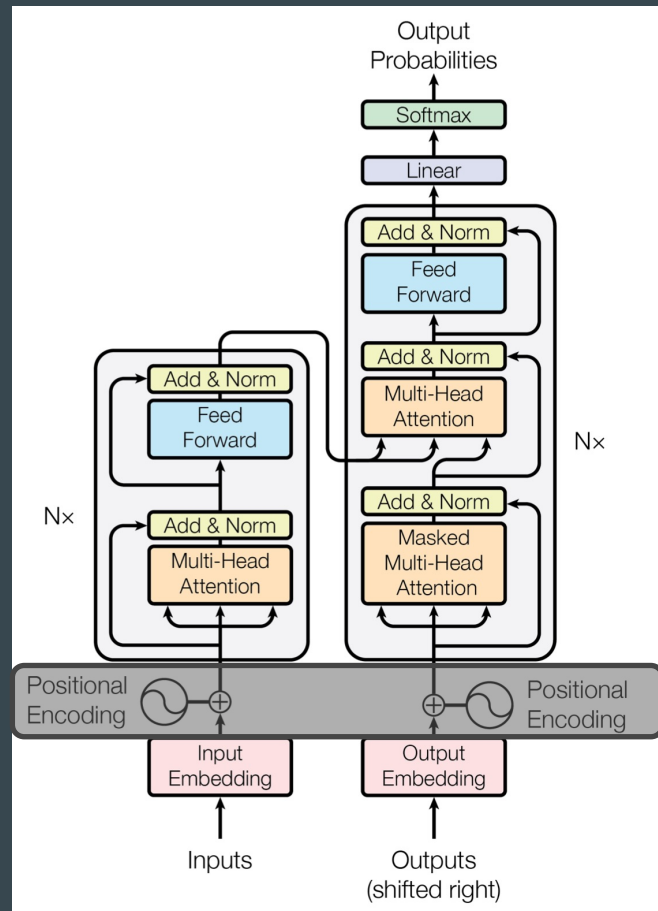
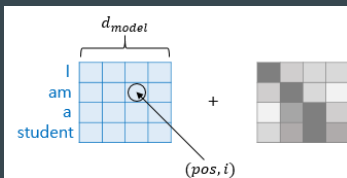
Transformer

- Positional encoding



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

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



Transformer

● Positional encoding

```
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)

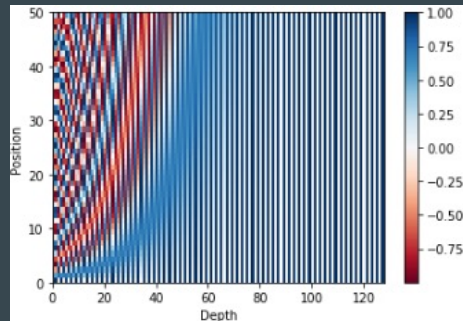
    def get_angles(self, position, i, d_model):
        angles = 1 / tf.pow(10000, (2 * (i // 2)) / tf.cast(d_model, tf.float32))
        return position * angles

    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, ...]

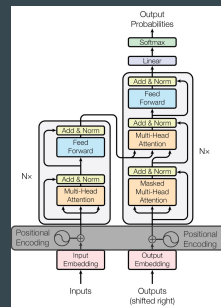
        print(pos_encoding.shape)
        return tf.cast(pos_encoding, tf.float32)

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



Positional Encoding Value

문장의 길이 50, 임베딩 벡터의 차원 128



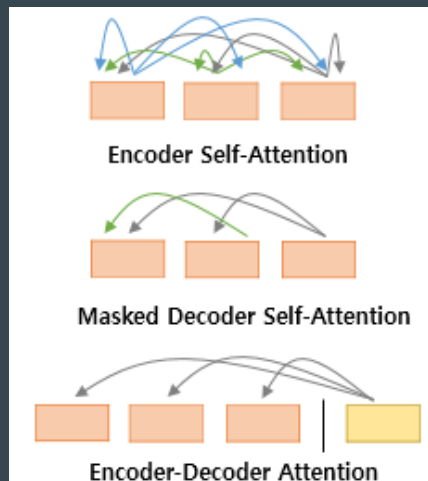
this	is	my	car
0.19	0.01	0.70	0.13
-0.47	0.01	0.34	0.14
-0.77	0.02	0.87	0.16
0.59	0.02	0.11	0.23
		-0.39	0.11
		0.04	0.24
		-0.91	0.07
		0.69	0.06
		0.79	0.04
		-0.25	0.11
		0.44	0.15

↓

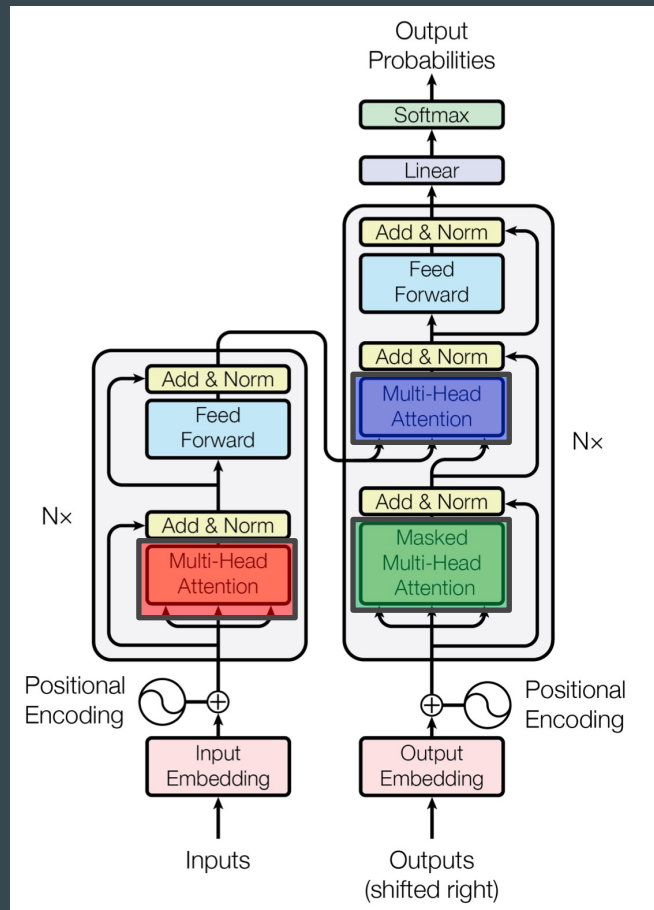
that	is	not	[PAD]
0.19	0.01	0.70	0.13
-0.41	0.01	0.74	0.14
0.12	0.02	0.87	0.16
0.59	0.02	0.11	0.23
		0.13	0.11
		0.04	0.24
		0.18	0.07
		0.74	0.05
		0.76	0.04
		0.13	0.11
		0.00	0.15

Transformer

- Attention_s,

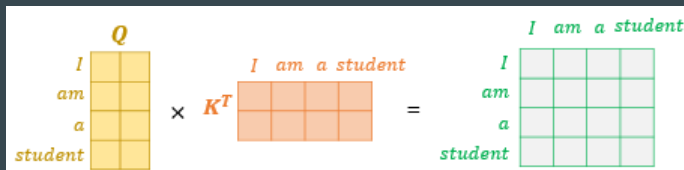
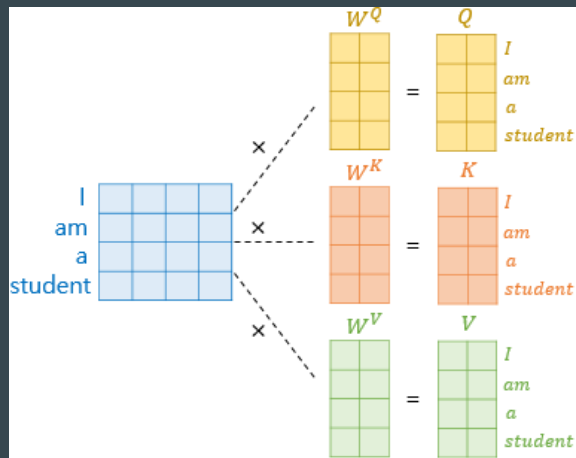
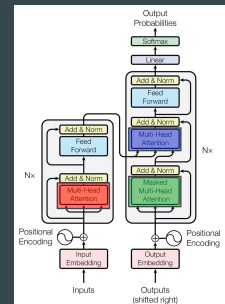


Encoder **Self-Attention** : Query = Key = Value
Decoder **Masked Self-Attention** : Query = Key = Value
Decoder **Encoder-Decoder Attention** : Query : Decoder Vector /
 Key = Value : Encoder Vector



Transformer

- Scaled dot-product Attention → MATRIX!!!!



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

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

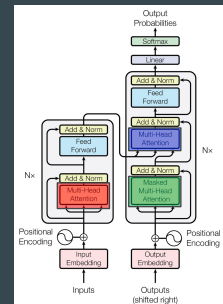
Transformer



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

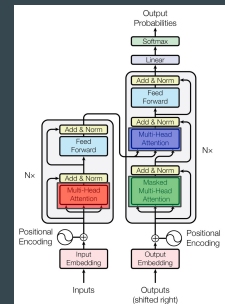
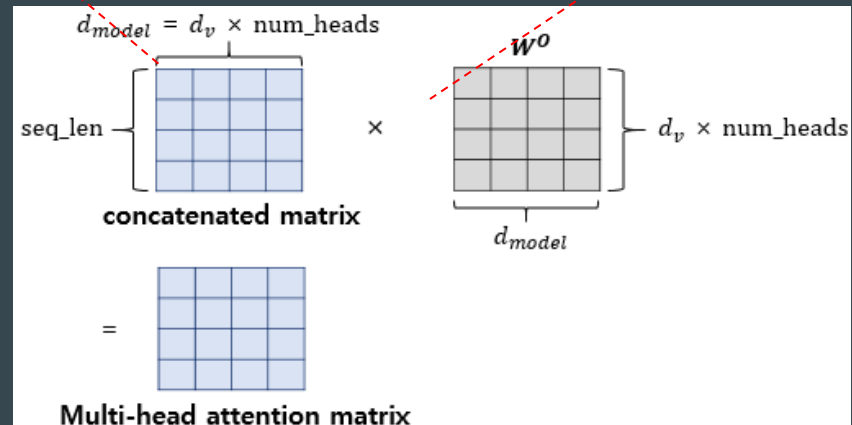
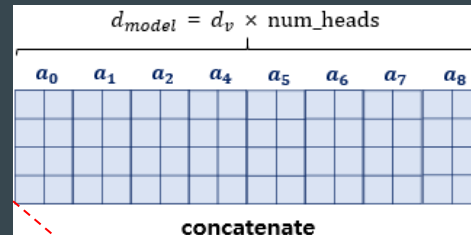
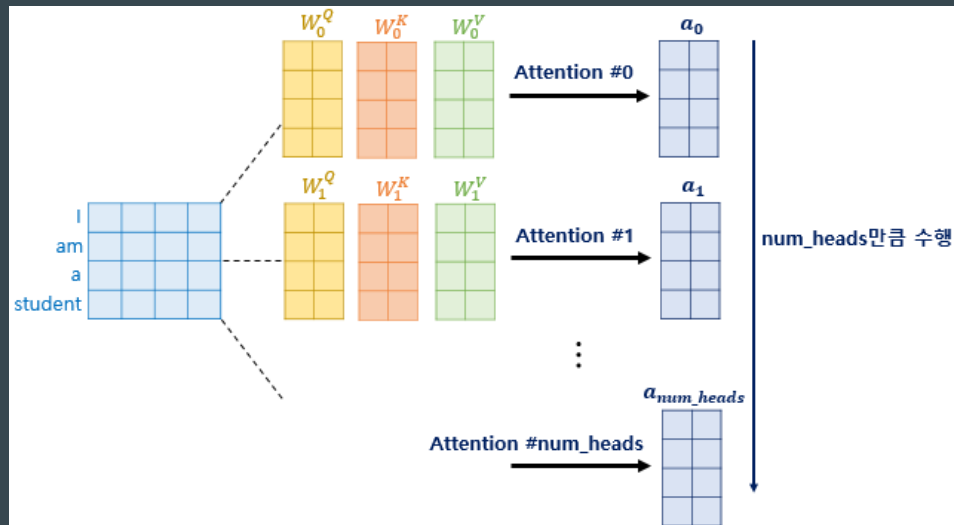
- Encoder Self-Attention

```
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
```



Transformer

- Multi-head Attention



Transformer

● 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
```

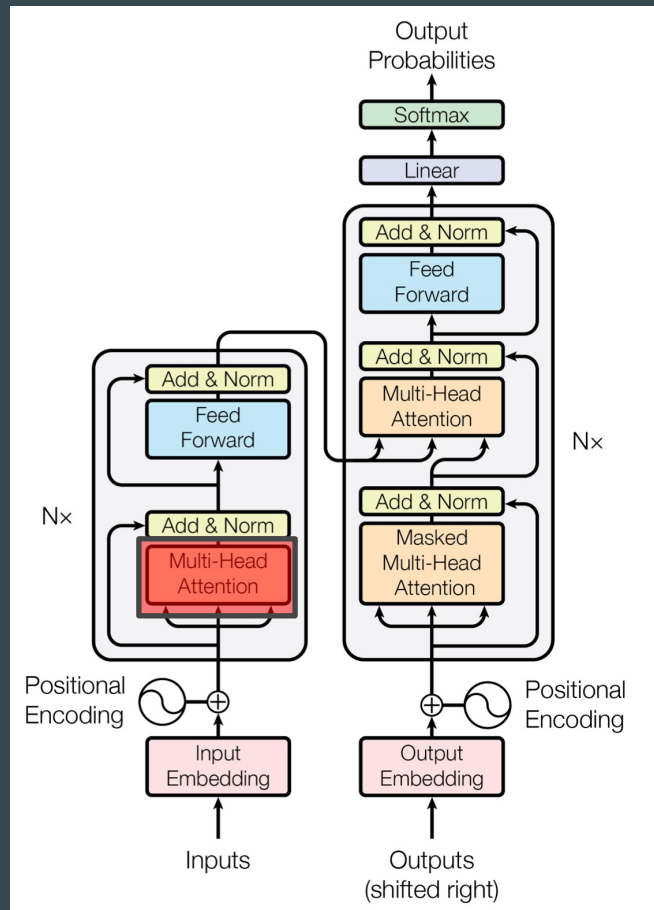
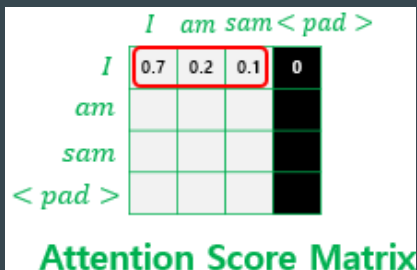
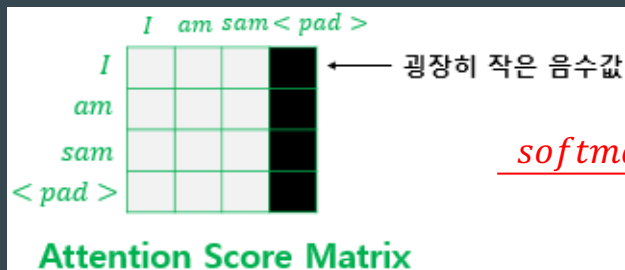
Transformer

- Padding Mask

- Key의 경우에 <PAD> 토큰이 존재한다면 이에 대해서는 유사도를 구하지 않도록

$$\frac{\begin{matrix} Q \\ \begin{matrix} I \\ am \\ sam \\ <pad> \end{matrix} \end{matrix} \times K^T \begin{matrix} \begin{matrix} I & am & sam & <pad> \end{matrix} \\ \begin{matrix} I \\ am \\ sam \\ <pad> \end{matrix} \end{matrix}}{\sqrt{d_k}} = \begin{matrix} \begin{matrix} I & am & sam & <pad> \end{matrix} \\ \begin{matrix} I \\ am \\ sam \\ <pad> \end{matrix} \end{matrix}$$

Attention Score Matrix

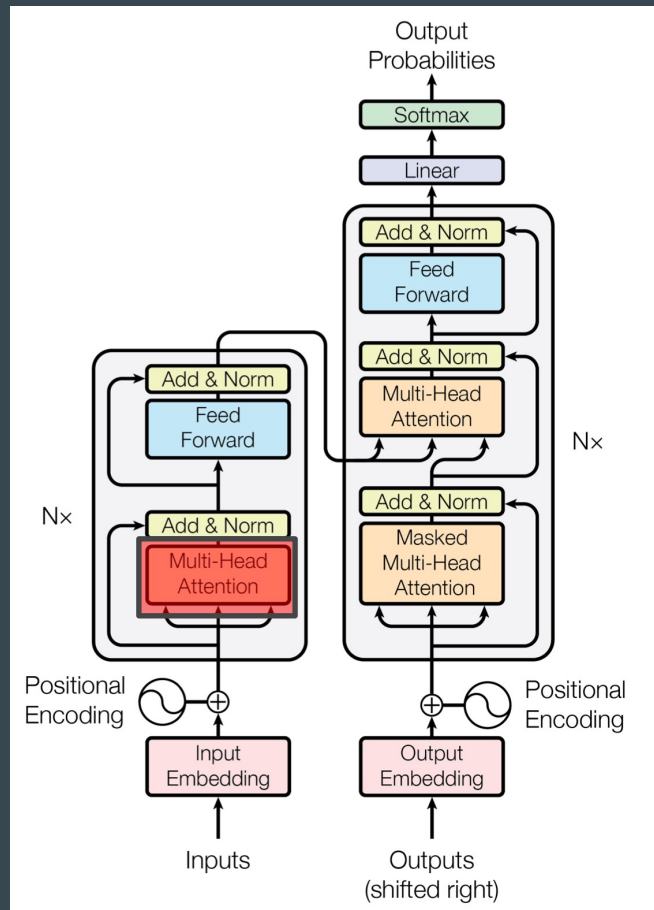
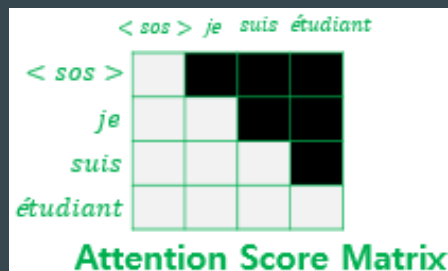


Transformer

- Look-ahead Mask
 - 현재 시점의 예측에서 현재 시점보다 미래에 있는 단어들을 참고하지 못하도록

$$\frac{\begin{matrix} Q \\ \begin{matrix} I \\ am \\ sam \\ <pad> \end{matrix} \end{matrix} \times K^T \begin{matrix} I \ am \ sam \ <pad> \\ \begin{matrix} I \\ am \\ sam \\ <pad> \end{matrix} \end{matrix}}{\sqrt{d_k}} = \begin{matrix} I \ am \ sam \ <pad> \\ \begin{matrix} I \\ am \\ sam \\ <pad> \end{matrix} \end{matrix}$$

Attention Score Matrix



Transformer

- Padding & Look-ahead Mask

인코더/디코더에서 문장의 끝을 Mask하는 함수

```
def create_padding_mask(x):  
    mask = tf.cast(tf.math.equal(x, 0), tf.float32)  
    # (batch_size, 1, 1, key의 문장 길이)  
    return mask[:, tf.newaxis, tf.newaxis, :]
```

디코더의 첫번째 서브층(sublayer)에서 미래 토큰을 Mask하는 함수

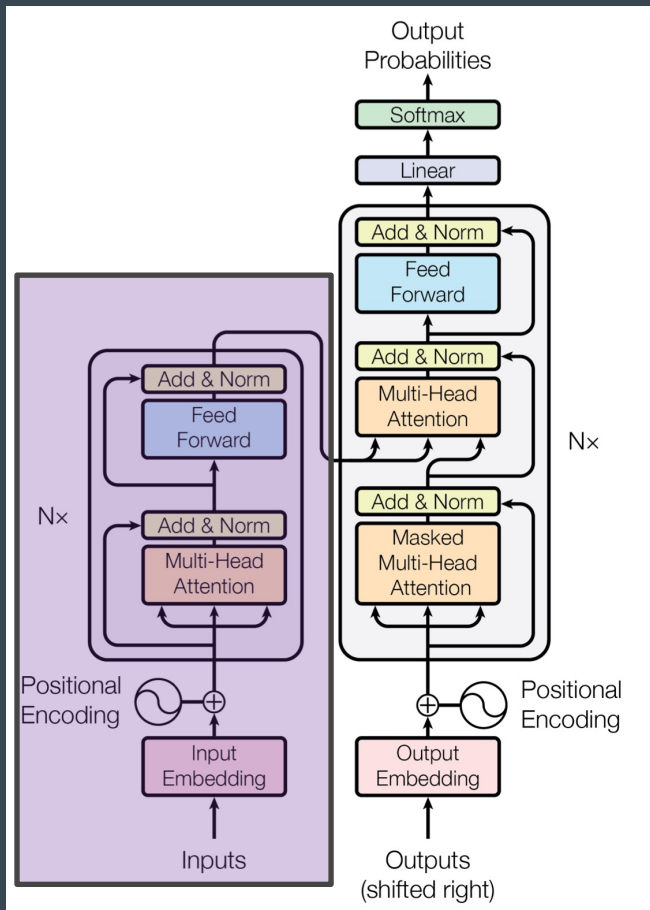
```
def create_look_ahead_mask(x):  
    seq_len = tf.shape(x)[1]  
    look_ahead_mask = 1 - tf.linalg.band_part(tf.ones((seq_len, seq_len)), -1, 0)  
    padding_mask = create_padding_mask(x) # 패딩 마스크도 포함  
    return tf.maximum(look_ahead_mask, padding_mask)
```

```
IPython: Users/bw  
In [3]: def create_padding_mask(x):  
...:     mask = tf.cast(tf.math.equal(x, 0), tf.float32)  
...:     # (batch_size, 1, 1, key의 문장 길이)  
...:     return mask[:, tf.newaxis, tf.newaxis, :]  
...:  
  
In [4]: print(create_padding_mask(tf.constant([[1, 21, 777, 0, 0]])))  
2022-12-12 00:07:16.486789: I tensorflow/core/platform/cpu_feature_guard.cc:193] This TensorFlow binary is optimized with oneAPI Deep Neural Network Library (on  
eDNN) to use the following CPU instructions in performance-critical operations:  
AVX2 FMA  
To enable them in other operations, rebuild TensorFlow with the appropriate comp  
iler flags.  
tf.Tensor([[[[0. 0. 0. 1. 1.]]]], shape=(1, 1, 1, 5), dtype=float32)  
  
In [5]:
```

```
IPython: Users/bw  
In [5]: def create_look_ahead_mask(x):  
...:     seq_len = tf.shape(x)[1]  
...:     look_ahead_mask = 1 - tf.linalg.band_part(tf.ones((seq_len, seq_len)), -1, 0)  
...:     padding_mask = create_padding_mask(x) # 패딩 마스크도 포함  
...:     return tf.maximum(look_ahead_mask, padding_mask)  
...:  
  
In [6]: print(create_look_ahead_mask(tf.constant([[1, 2, 0, 4, 5]])))  
tf.Tensor(  
[[[0. 1. 1. 1. 1.]  
  [0. 0. 1. 1. 1.]  
  [0. 0. 1. 1. 1.]  
  [0. 0. 1. 0. 1.]  
  [0. 0. 1. 0. 0.] ]]], shape=(1, 1, 5, 5), dtype=float32)  
  
In [7]:  
  
In [7]: print(create_look_ahead_mask(tf.constant([[1, 2, 0, 4, 5, 0]])))  
tf.Tensor(  
[[[0. 1. 1. 1. 1. 1.]  
  [0. 0. 1. 1. 1. 1.]  
  [0. 0. 1. 1. 1. 1.]  
  [0. 0. 1. 0. 1. 1.]  
  [0. 0. 1. 0. 0. 1.]  
  [0. 0. 1. 0. 0. 1.] ]]], shape=(1, 1, 6, 6), dtype=float32)  
  
In [8]:
```

Transformer

- Encoder



Transformer

● Encoder

```
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)
```

```
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)

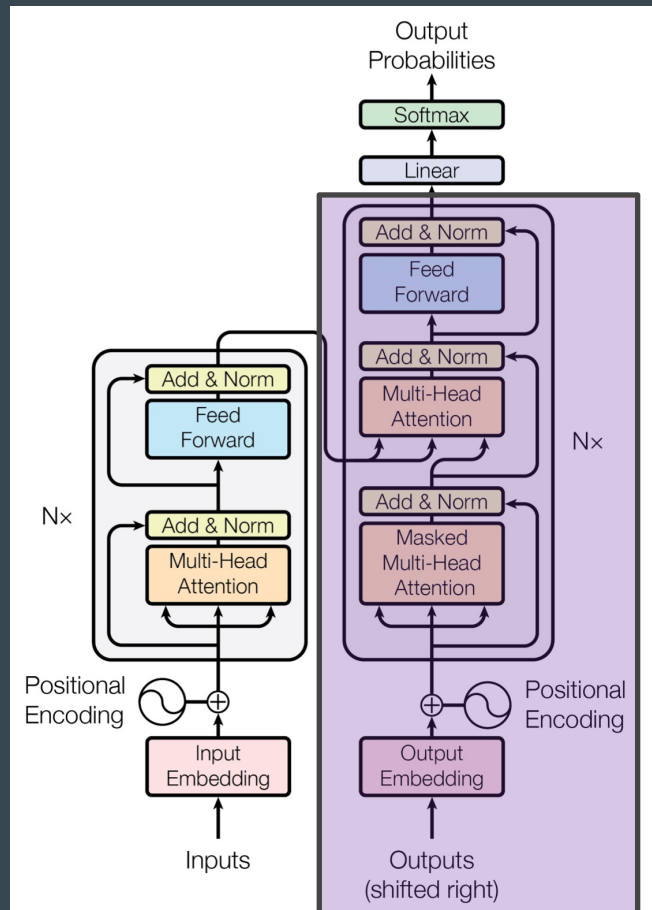
    # 포지션 와이즈 피드 포워드 신경망 (두번째 서브층)
    outputs = tf.keras.layers.Dense(units=dff, activation='relu')(attention)
    outputs = tf.keras.layers.Dense(units=d_model)(outputs)

    # 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).
```

Transformer

- Decoder



Transformer

● Decoder

```
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)
```

```
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")

    # 록어헤드 마스크(첫번째 서브층)
    look_ahead_mask = tf.keras.Input(
        shape=(1, None, None), name="look_ahead_mask")

    # 패딩 마스크(두번째 서브층)
    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)

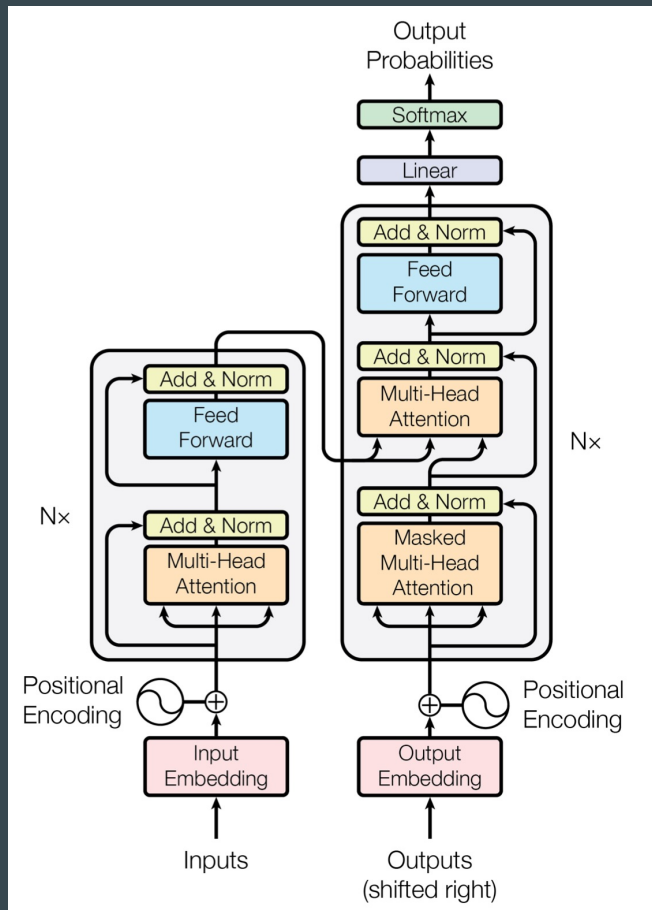
    # 포지션 와이즈 피드 포워드 신경망 (세번째 서브층)
    outputs = tf.keras.layers.Dense(units=dff, activation='relu')(attention2)
    outputs = tf.keras.layers.Dense(units=d_model)(outputs)

    # 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)
```

Transformer

- All together



Transformer

```
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 룩어헤드 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)
```

Everything can be found in

<https://wikidocs.net/book/2155>

<https://github.com/ukairia777/tensorflow-nlp-tutorial>