

Transformer

MokkeMeguru

<2019-08-24 土>

Contents

1 TensorFlow ライブラリのインポート	1
2 データセットのインポート	3
3 Positional Encoding	5
4 Masking	7
5 Scaled dot product attention	8
6 Multi-head attention	12
7 Point wise feed forward network	14
8 Encoder and Decoder	15
8.1 Encoder Layer	15
8.2 Decoder Layer	16
8.3 Encoder	17
8.4 Decoder	19
9 Create the Transformer	20
10 Set hyperparameters	21
11 Optimizer	21
12 Loss and metrics	23
13 Training and checkpointing	24
14 Evaluation	25

1 TensorFlow ライブラリのインポート

```
1 from __future__ import division, absolute_import
2 from __future__ import print_function, unicode_literals
```

```

3  from functools import reduce, partial
4
5  import tensorflow as tf
6  # import tensorflow_hub as hub
7  import tensorflow_datasets as tfds
8  import tensorflow_probability as tfp
9  # from tensorflow_examples.models.pix2pix import pix2pix
10
11 from tensorflow import keras
12 from tensorflow.keras import layers, datasets, models
13 from tensorflow.keras.models import Sequential
14 # import tensorflow_text as text
15
16 import numpy as np
17 import matplotlib.pyplot as plt
18 import matplotlib.ticker as ticker
19
20 import pandas as pd
21 from sklearn.model_selection import train_test_split
22 # import seaborn as sns
23 import os
24 import time
25 # import yaml
26 # import h5py
27 # import pathlib
28 # import random
29 # import IPython.display as display
30 # from IPython.display import clear_output
31 # import PIL.Image as Image
32 # import urllib3
33 import io
34 import tempfile
35 from pprint import pprint
36
37 import unicodedata
38 import re
39 import time
40
41 def print_infos(infolist: list):
42     for info in infolist:
43         print(info)
44
45
46 def pprint_infos(infolist: list):
47     for info in infolist:
48         pprint(info)
49
50
51 print_infos([

```

```

52     '{:25}: {}'.format("tensorflow\'s version", tf.__version__),
53     # '{:25}: {}'.format("tensorflow\'s version", hub.__version__),
54 ] )
55
56 AUTOTUNE = tf.data.experimental.AUTOTUNE
57 # urllib3.disable_warnings(urllib3.exceptions.InsecureRequestWarning)

```

tensorflow's version : 2.0.0-rc0

2 データセットのインポート

```

1  examples, metadata = tfds.load('ted_hrlr_translate/pt_to_en',
2                                with_info=True,
3                                as_supervised=True)
4  train_examples, val_examples = examples['train'], examples['validation']
5
6  for pt_ex, en_ex in train_examples.take(1):
7      print(pt_ex.numpy())
8      print(en_ex.numpy())

```

WARNING: Logging before flag parsing goes to stderr.
W0904 14:16:00.669344 139700101117568 dataset_builder.py:439] Warning: Setting shuffle_files=True because split=TRAIN and shuffle_files=None. This behavior will be deprecated on 2019-08-06, at which point shuffle_files=False will be the default for all splits.

b'os astr33nomos acreditam que cada estrela da gal31xia tem um planeta , e especulam que at39 um quinto deles tem um planeta do tipo da terra que poder31 ter vida , mas ainda n33o vimos nenhum deles .'
b"astronomers now believe that every star in the galaxy has a planet , and they speculate that up to one fifth of them have an earth-like planet that might be able to harbor life , but we have n't seen any of them ."

```

1  tokenizer_en = tfds.features.text.SubwordTextEncoder.build_from_corpus(
2      (en.numpy() for pt, en in train_examples), target_vocab_size=2**13)
3
4  tokenizer_pt = tfds.features.text.SubwordTextEncoder.build_from_corpus(
5      (pt.numpy() for pt, en in train_examples), target_vocab_size=2**13)
6  sample_string = 'Transformer is awesome.'
7  tokenized_string = tokenizer_en.encode(sample_string)
8  original_string = tokenizer_en.decode(tokenized_string)
9  assert original_string == sample_string
10
11 print_infos([
12     'Tokenized string is {}'.format(tokenized_string),
13     'The original string: {}'.format(original_string)
14 ])

```

Tokenized string is [7915, 1248, 7946, 7194, 13, 2799, 7877]
The original string: Transformer is awesome.

```
1 for ts in tokenized_string:
2     print('{} -> {}'.format(ts, tokenizer_en.decode([ts])))
```

7915 -> T
1248 -> ran
7946 -> s
7194 -> former
13 -> is
2799 -> awesome
7877 -> .

```
1 BUFFER_SIZE = 20000
2 BATCH_SIZE = 64
3 def encode(lang1, lang2):
4     # <start> = (tokenizer_pt / tokenizer_en) .vocab_size
5     # <end> = (tokenizer_pt / tokenizer_en) .vocab_size + 1
6     lang1 = [tokenizer_pt.vocab_size] + tokenizer_pt.encode(
7         lang1.numpy()) + [tokenizer_pt.vocab_size + 1]
8     lang2 = [tokenizer_en.vocab_size] + tokenizer_en.encode(
9         lang2.numpy()) + [tokenizer_en.vocab_size + 1]
10    return lang1, lang2
```

```
1 MAX_LENGTH = 40
2
3
4 def filter_max_length(x, y, max_length=MAX_LENGTH):
5     return tf.logical_and(tf.size(x) <= max_length, tf.size(y) <=
6         max_length)
7
8 def tf_encode(pt, en):
9     return tf.py_function(encode, [pt, en], [tf.int64, tf.int64])
```

```
1 train_dataset = train_examples.map(tf_encode).filter(
2     filter_max_length).cache().shuffle(BUFFER_SIZE).padded_batch(
3     BATCH_SIZE, padded_shapes=[[-1], [-1]]).prefetch(AUTOTUNE)
4
```

```

5 val_dataset = val_examples.map(tf_encode).filter(
6     filter_max_length).padded_batch(BATCH_SIZE, padded_shapes=[[-1], [-1]])

```

```

1 pt_batch, en_batch = next(iter(val_dataset))
2 pt_batch, en_batch

```

```

(<tf.Tensor: id=546958, shape=(64, 40), dtype=int64, numpy=
array([[8214, 1259,    5, ...,    0,    0,    0],
      [8214,  299,   13, ...,    0,    0,    0],
      [8214,   59,    8, ...,    0,    0,    0],
      ...,
      [8214,   95,    3, ...,    0,    0,    0],
      [8214, 5157,    1, ...,    0,    0,    0],
      [8214, 4479, 7990, ...,    0,    0,    0]])>,
<tf.Tensor: id=546959, shape=(64, 40), dtype=int64, numpy=
array([[8087,   18,   12, ...,    0,    0,    0],
      [8087,  634,   30, ...,    0,    0,    0],
      [8087,   16,   13, ...,    0,    0,    0],
      ...,
      [8087,   12,   20, ...,    0,    0,    0],
      [8087,   17, 4981, ...,    0,    0,    0],
      [8087,   12, 5453, ...,    0,    0,    0]])>)

```

3 Positional Encoding

ref. https://github.com/tensorflow/examples/blob/master/community/en/position_encoding.ipynb

The formula is here.

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

```

1 def get_angles(pos, i, d_model):
2     angle_rates = 1 / np.power(1000, (2 * (i // 2)) / np.float32(d_model))
3     return pos * angle_rates
4
5
6 def positional_encoding(position, d_model):
7     angle_rads = get_angles(
8         np.arange(position)[:, np.newaxis],
9         np.arange(d_model)[np.newaxis, :], d_model)
10
11     angle_rads[:, 0::2] = np.sin(angle_rads[:, 0::2])

```

```

12     angle_rads[:, 1::2] = np.cos(angle_rads[:, 0::2])
13     pos_encoding = angle_rads[np.newaxis, ...]
14     return tf.cast(pos_encoding, dtype=tf.float32)

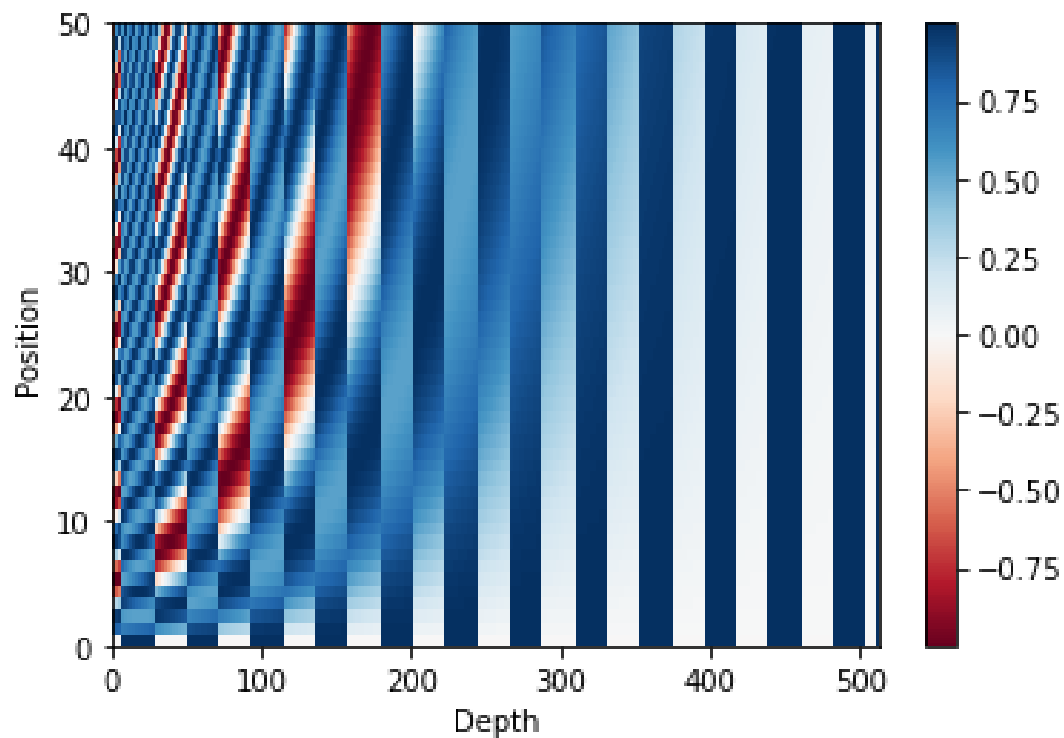
```

```

1 pos_encoding = positional_encoding(50, 512)
2 print(pos_encoding.shape)
3 plt.pcolormesh(pos_encoding[0], cmap='RdBu')
4 plt.xlabel('Depth')
5 plt.xlim((0, 512))
6 plt.ylabel('Position')
7 plt.colorbar()
8 plt.show()

```

(1, 50, 512)



appendix

```

1 print(np.arange(10)[:, np.newaxis])
2 print(np.arange(10)[np.newaxis, :])

```

```

[[0]
 [1]
 [2]
 [3]
 [4]
 [5]
 [6]
 [7]
 [8]
 [9]]
[[0 1 2 3 4 5 6 7 8 9]]

```

4 Masking

for ignoring padding in calculation

```

1 def create_padding_mask(seq):
2     seq = tf.cast(tf.math.equal(seq, 0), tf.float32)
3     # (batch_size, 1, 1, seq_len)
4     return seq[:, tf.newaxis, tf.newaxis, :]
5
6
7 x = tf.constant([[7, 6, 0, 0, 1], [1, 2, 3, 0, 0], [0, 0, 0, 4, 5]])
8 print_infos([x, create_padding_mask(x)])

```

```

tf.Tensor(
[[7 6 0 0 1]
 [1 2 3 0 0]
 [0 0 0 4 5]], shape=(3, 5), dtype=int32)

```

```

tf.Tensor(
[[[0. 0. 1. 1. 0.]]]

```

```

[[[0. 0. 0. 1. 1.]]]

```

```

[[[1. 1. 1. 0. 0.]]], shape=(3, 1, 1, 5), dtype=float32)

```

for ignoring prediction word (in decoding model)

```

1 def create_look_ahead_mask(size):
2     mask = 1 - tf.linalg.band_part(tf.ones((size, size)), -1, 0)
3     return mask
4
5

```

```

6 x = tf.random.uniform((1, 3))
7 print_infos([x, create_look_ahead_mask(x.shape[1])])

```

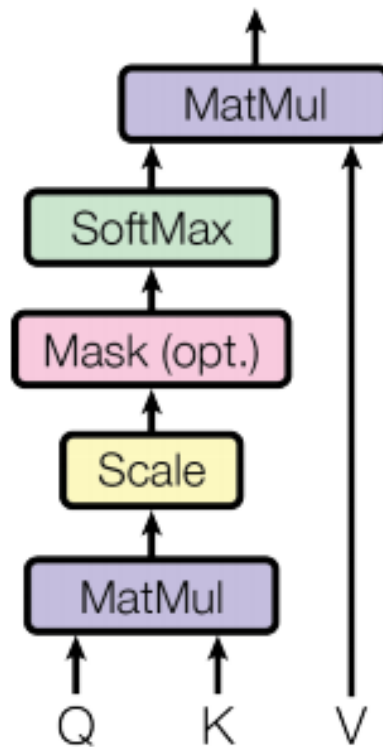
```

tf.Tensor([[0.22582567 0.58736205 0.92260003]], shape=(1, 3), dtype=float32)
tf.Tensor(
[[0. 1. 1.]
 [0. 0. 1.]
 [0. 0. 0.]], shape=(3, 3), dtype=float32)

```

5 Scaled dot product attention

Scaled Dot-Product Attention



Attention formula is here (Q is query, K is key, V is value)

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

```

1 def scaled_dot_product_attention(q, k, v, mask):
2     """Calculate the attention weights.
3     q, k, v must have matching leading dimensions.
4     k, v must have matching penultimate dimension, i.e.: seq_len_k =
    ↪ seq_len_v.
5     The mask has different shapes depending on its type(padding or look
    ↪ ahead)
6     but it must be broadcastable for addition.
7
8     q, k, v は次に示される次元である必要があります
9     k, v は 第2次元のサイズを統一されている必要があります (つまり seq_len_k ==
    ↪ seq_len_v)
10    マスクは、そのタイプ (padding / look ahead) に応じてサイズが変わります。
11    しかし (... seq_len_q, seq_len_k) へブロードキャストできるようになっていなければなり
    ↪ ません
12
13    Args:
14        q: query shape == (... , seq_len_q, depth)
15        k: key shape == (... , seq_len_k, depth)
16        v: value shape == (... , seq_len_v, depth_v)
17        mask: Float tensor with shape broadcastable
18              to (... , seq_len_q, seq_len_k). Defaults to None.
19
20    Returns:
21        output, attention_weights
22    """
23
24    matmul_qk = tf.matmul(q, k,
25                           transpose_b=True) # (... , seq_len_q, seq_len_k)
26
27    # scale matmul_qk
28    dk = tf.cast(tf.shape(k)[-1], tf.float32)
29    scaled_attention_logits = matmul_qk / tf.math.sqrt(dk)
30
31    # add the mask to the scaled tensor
32    if mask is not None:
33        scaled_attention_logits += (mask * -1e9)
34
35    # softmax is normalized on the last axis (seq_len_k) so that the scores
36    # add up to 1
37    attention_weights = tf.nn.softmax(scaled_attention_logits, axis=-1)
38    # (... , seq_len_q, seq_len_k)
39
40    output = tf.matmul(attention_weights, v)
41    # (... , seq_len_q, depth_v)

```

```
42     return output, attention_weights
```

```
1  def print_out(q, k, v):
2      temp_out, temp_attn = scaled_dot_product_attention(q, k, v, None)
3      print_infos(['Attention weights are:', temp_attn, 'Output is:',
4                  ↪ temp_out])
5
6  np.set_printoptions(suppress=True)
7  temp_k = tf.constant([[10, 0, 0], [0, 10, 0], [0, 0, 10], [0, 0, 10]],
8                        dtype=tf.float32) # (4, 3)
9  temp_v = tf.constant([[1, 0], [10, 0], [100, 5], [1000, 6]],
10                       dtype=tf.float32) # (4, 2)
11
12  # この query は 2 番目の key と一致するので、2 番目の value が返されます。
13  # => v[k.search(query)]
14  temp_q = tf.constant([[0, 10, 0]], dtype=tf.float32) # (1, 3)
15  print_out(temp_q, temp_k, temp_v)
16  print()
17
18  # この query は 3, 4 番目の key と一致するので、3, 4 番目の value の平均値が返されま
19  ↪ す。
20  temp_q = tf.constant([[0, 0, 10]], dtype=tf.float32) # (1, 3)
21  print_out(temp_q, temp_k, temp_v)
22
23  # この query は 1, 2 番目の key と一致するので、1, 2 番目の value の平均値が返されます。
24  ↪
25  temp_q = tf.constant([[10, 10, 0]], dtype=tf.float32) # (1, 3)
26  print_out(temp_q, temp_k, temp_v)
```

```
Attention weights are:
tf.Tensor([[0. 1. 0. 0.]], shape=(1, 4), dtype=float32)
Output is:
tf.Tensor([[10. 0.]], shape=(1, 2), dtype=float32)
```

```
Attention weights are:
tf.Tensor([[0. 0. 0.5 0.5]], shape=(1, 4), dtype=float32)
Output is:
tf.Tensor([[550. 5.5]], shape=(1, 2), dtype=float32)
Attention weights are:
tf.Tensor([[0.5 0.5 0. 0. ]], shape=(1, 4), dtype=float32)
Output is:
tf.Tensor([[5.5 0. ]], shape=(1, 2), dtype=float32)
```

```
1  # 上の query をすべて行列にまとめて実行すると次のようになります。
2  temp_q = tf.constant([[0, 0, 10], [0, 10, 0], [10, 10, 0]],
```

```
3         dtype=tf.float32) # (3, 3)
4 print_out(temp_q, temp_k, temp_v)
```

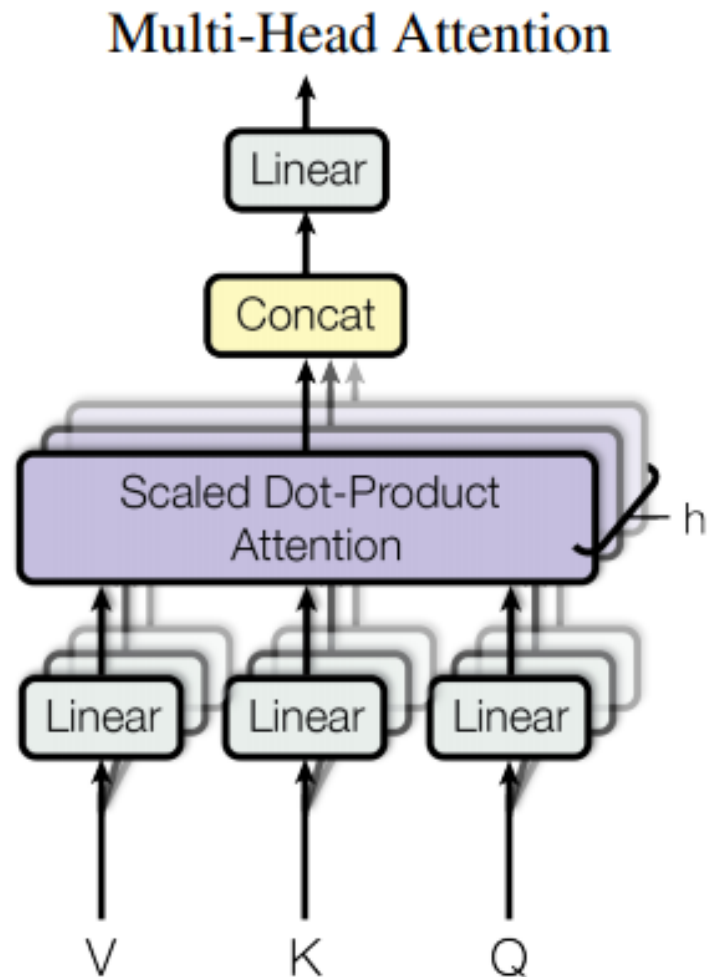
Attention weights are:

```
tf.Tensor(
[[0.  0.  0.5 0.5]
 [0.  1.  0.  0. ]
 [0.5 0.5 0.  0. ]], shape=(3, 4), dtype=float32)
```

Output is:

```
tf.Tensor(
[[550.    5.5]
 [ 10.    0. ]
 [  5.5   0. ]], shape=(3, 2), dtype=float32)
```

6 Multi-head attention



Multi head Attention はの 4 つのパートから構成されます。

- Linear layer と 複数の head への分割
- Scaled dot-product attention
- heads の集約
- final linear layer

```

1 class MultiHeadAttention(layers.Layer):
2     def __init__(self, d_model, num_heads):
3         super(MultiHeadAttention, self).__init__()
4         self.num_heads = num_heads
5         self.d_model = d_model
6
7         assert d_model % self.num_heads == 0
8
9         self.depth = d_model // self.num_heads
10
11        self.wq = layers.Dense(d_model)
12        self.wk = layers.Dense(d_model)
13        self.wv = layers.Dense(d_model)
14
15        self.dense = layers.Dense(d_model)
16
17        def split_heads(self, x, batch_size):
18            """Split the last dimension into (num_heads, depth).
19            Transpose the result such that the shape is (batch_size, num_heads,
20            ↪ seq_len, depth)"""
21
22            最後の次元である d_model を (num_heads, depth) へ分割します。
23            また出力時には
24            (batch_size, seq_len, num_heads, depth) -> (batch_size, num_heads,
25            ↪ seq_len, depth) します。
26
27            つまり
28            (batch_size, seq_len, d_model)
29            -> (batch_size, seq_len, num_heads, depth) (次元分割)
30            -> (batch_size, num_heads, seq_len, depth) (軸入れ替え)
31            """
32            x = tf.reshape(x, (batch_size, -1, self.num_heads, self.depth))
33            return tf.transpose(x, perm=[0, 2, 1, 3])
34
35        def call(self, v, k, q, mask):
36            batch_size = tf.shape(q)[0]
37
38            q = self.wq(q) # (... , seq_len, d_model)
39            k = self.wk(k) # (... , seq_len, d_model)
40            v = self.wv(v) # (... , seq_len, d_model)
41
42            q = self.split_heads(
43                q, batch_size) # (batch_size, num_heads, seq_len_q, depth)
44            k = self.split_heads(
45                k, batch_size) # (batch_size, num_heads, seq_len_k, depth)
46            v = self.split_heads(
47                v, batch_size) # (batch_size, num_heads, seq_len_v, depth)

```

```

47     # scaled_attention.shape == (batch_size, num_heads, seq_len_q,
    ↪ depth)
48     # attention_weights.shape == (batch_size, num_heads, seq_len_q,
    ↪ seq_len_k)
49     scaled_attention, attention_weights = scaled_dot_product_attention(
50         q, k, v, mask)
51
52     # (batch_size, seq_len_q, num_heads, depth)
53     scaled_attention = tf.transpose(scaled_attention, perm=[0, 2, 1,
    ↪ 3])
54
55     concat_attention = tf.reshape(
56         scaled_attention,
57         (batch_size, -1, self.d_model)) # (batch_size, seq_len_q,
    ↪ d_model)
58
59     output = self.dense(
60         concat_attention) # (batch_size, seq_len_q, d_model)
61
62     return output, attention_weights

```

```

1 temp_mha = MultiHeadAttention(d_model=512, num_heads=8)
2 y = tf.random.uniform([1, 60, 512]) # (batch_size, encoder_sequence,
    ↪ d_model)
3 out, attn = temp_mha(y, k=y, q=y, mask=None)
4 out.shape, attn.shape

```

(TensorShape([1, 60, 512]), TensorShape([1, 8, 60, 60]))

7 Point wise feed forward network

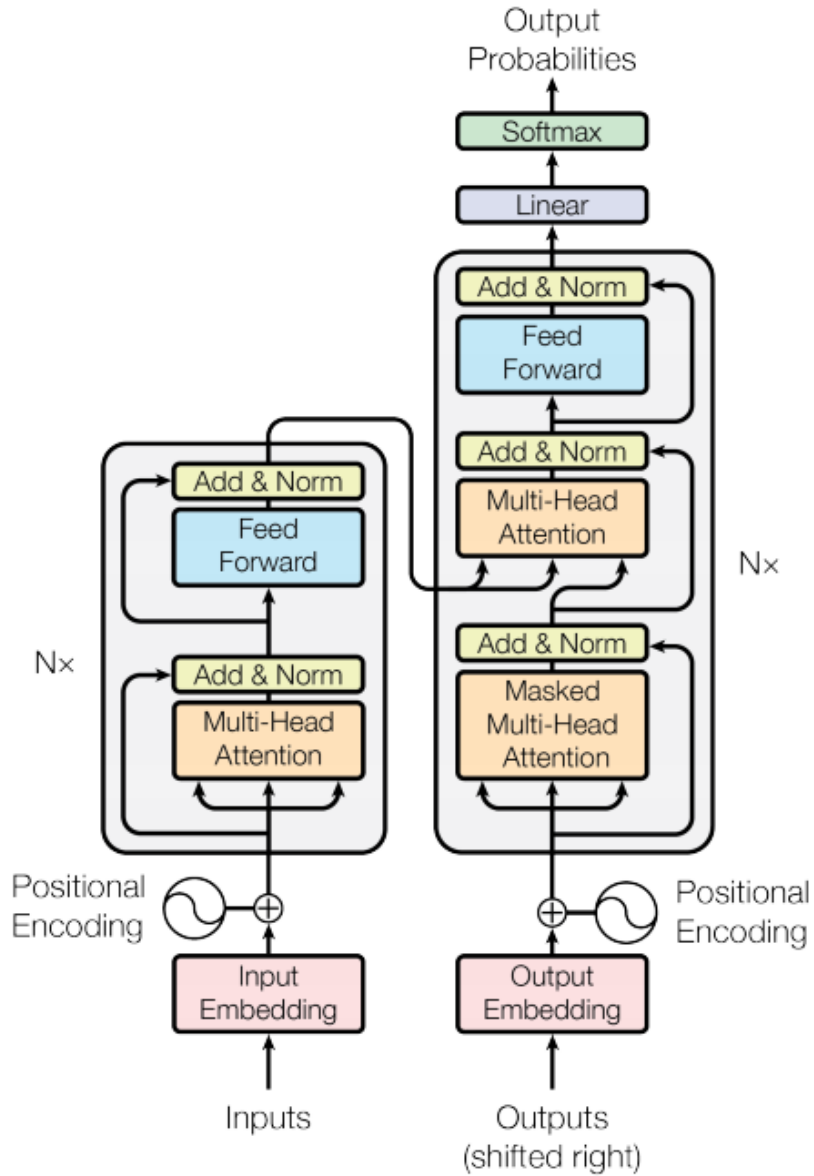
```

1 def point_wise_feed_forward_network(d_model, dff):
2     return Sequential(
3         [layers.Dense(dff, activation='relu'),
4          layers.Dense(d_model)])
5
6
7 sample_ffn = point_wise_feed_forward_network(512, 2048)
8 sample_ffn(tf.random.uniform([64, 50, 512])).shape

```

TensorShape([64, 50, 512])

8 Encoder and Decoder



8.1 Encoder Layer

```
1 class EncoderLayer(layers.Layer):  
2     def __init__(self, d_model, num_heads, dff, rate=0.1):  
3         super(EncoderLayer, self).__init__()
```

```

4         self.mha = MultiHeadAttention(d_model, num_heads)
5         self.layer_norm1 = layers.LayerNormalization(epsilon=1e-6)
6
7         self.ffn = point_wise_feed_forward_network(d_model, dff)
8         self.layer_norm2 = layers.LayerNormalization(epsilon=1e-6)
9
10        self.dropout1 = layers.Dropout(rate)
11        self.dropout2 = layers.Dropout(rate)
12
13    def call(self, x, training, padding_mask):
14        # (... , input_seq_len, d_model)
15        attn_output, _ = self.mha(x, x, x, padding_mask)
16        attn_output = self.dropout1(attn_output, training=training)
17        # (... , input_seq_len, d_model)
18        out1 = self.layer_norm1(x + attn_output)
19
20        # (... , input_seq_len, d_model)
21        ffn_output = self.ffn(out1)
22        ffn_output = self.dropout2(ffn_output, training=training)
23        # (... , input_seq_len, d_model)
24        out2 = self.layer_norm2(out1 + ffn_output)
25
26    return out2

```

```

1 sample_encoder_layer = EncoderLayer(d_model=512, num_heads=8, dff=2048)
2 sample_encoder_layer_output = sample_encoder_layer(
3     tf.random.uniform((64, 42, 512)), False, None)
4 sample_encoder_layer_output.shape

```

TensorShape([64, 42, 512])

8.2 Decoder Layer

```

1 class DecoderLayer(layers.Layer):
2     def __init__(self, d_model, num_heads, dff, rate=0.1):
3         super(DecoderLayer, self).__init__()
4
5         self.mha1 = MultiHeadAttention(d_model, num_heads)
6         self.layernorm1 = layers.LayerNormalization(epsilon=1e-6)
7
8         self.mha2 = MultiHeadAttention(d_model, num_heads)
9         self.layernorm2 = layers.LayerNormalization(epsilon=1e-6)
10
11        self.ffn = point_wise_feed_forward_network(d_model, dff)
12        self.layernorm3 = layers.LayerNormalization(epsilon=1e-6)
13

```



```

14         self.dropout1 = layers.Dropout(rate)
15         self.dropout2 = layers.Dropout(rate)
16         self.dropout3 = layers.Dropout(rate)
17
18     def call(self, x, enc_output, training, look_ahead_mask, padding_mask):
19         # enc_output.shape == (... , input_seq_len, d_model)
20
21         #(... , target_seq_len, d_model)
22         attn1, attn_weights_block1 = self.mha1(x, x, x,
23         ↪ look_ahead_mask)
24         attn1 = self.dropout1(attn1, training=training)
25         out1 = self.layernorm1(attn1 + x)
26
27         #(... , target_seq_len, d_model)
28         # これは 'self attention algorithm'
29         #       v = enc_output, k=enc_output, q=x
30         #       => enc_output でできた v, k に query out1 でアクセスする
31         attn2, attn_weights_block2 = self.mha2(enc_output, enc_output,
32         ↪ out1,
33         padding_mask)
34         attn2 = self.dropout1(attn2, training=training)
35         out2 = self.layernorm1(attn2 + out1)
36
37         # (batch_size, target_seq_len, d_model)
38         ffn_output = self.ffn(out2)
39         ffn_output = self.dropout3(
40         ↪ ffn_output,
41         training=training) # (batch_size, target_seq_len, d_model)
42         out3 = self.layernorm3(ffn_output + out2)
43
44         return out3, attn_weights_block1, attn_weights_block2

```

```

1 sample_decoder_layer = DecoderLayer(512, 8, 2048)
2 sample_decoder_layer_output, _, _ = sample_decoder_layer(
3     tf.random.uniform((64, 50, 512)), sample_encoder_layer_output, False,
4     ↪ None,
5     None)
6 sample_decoder_layer_output.shape

```

TensorShape([64, 50, 512])

8.3 Encoder

```

1 class Encoder(layers.Layer):
2     def __init__(self,

```

```

3         num_layers,
4         d_model,
5         num_heads,
6         dff,
7         input_vocab_size,
8         rate=0.1):
9     super(Encoder, self).__init__()
10    self.d_model = d_model
11    self.num_layers = num_layers
12    self.embedding = layers.Embedding(input_vocab_size, self.d_model)
13
14    # assumption: input_vocab_size > seq_len
15    self.pos_encoding = positional_encoding(input_vocab_size,
16    ↪ self.d_model)
17    self.enc_layers = [
18        EncoderLayer(d_model, num_heads, dff, rate)
19        for _ in range(num_layers)
20    ]
21    self.dropout = layers.Dropout(rate)
22
23    def call(self, x, training, mask):
24        # x.shape == (... , seq_len)
25        seq_len = tf.shape(x)[1]
26
27        # (batch_size, input_seq_len, d_model)
28        x = self.embedding(x)
29        x *= tf.math.sqrt(tf.cast(self.d_model, tf.float32))
30        x += self.pos_encoding[:, :seq_len, :]
31
32        self.dropout(x, training=training)
33
34        for i in range(self.num_layers):
35            x = self.enc_layers[i](x, training, mask)
36
37        # (batch_size, input_seq_len, d_model)
38        return x

```

```

1 sample_encoder = Encoder(num_layers=2,
2                           d_model=512,
3                           num_heads=8,
4                           dff=2048,
5                           input_vocab_size=8500)
6 sample_encoder_output = sample_encoder(tf.random.uniform((64, 62)),
7                                       training=False,
8                                       mask=None)
9
10 # (batch_size, input_seq_len, d_model)
11 print(sample_encoder_output.shape)

```

(64, 62, 512)

8.4 Decoder

```
1 class Decoder(layers.Layer):
2     def __init__(self,
3                 num_layers,
4                 d_model,
5                 num_heads,
6                 dff,
7                 target_vocab_size,
8                 rate=0.1):
9         super(Decoder, self).__init__()
10        self.d_model = d_model
11        self.num_layers = num_layers
12        self.embedding = layers.Embedding(target_vocab_size, d_model)
13
14        # assumption: target_vocab_size > seq_len
15        self.pos_encoding = positional_encoding(target_vocab_size, d_model)
16        self.dec_layers = [
17            DecoderLayer(d_model, num_heads, dff, rate)
18            for _ in range(num_layers)
19        ]
20        self.dropout = layers.Dropout(rate)
21
22    def call(self, x, enc_output, training, look_ahead_mask, padding_mask):
23        seq_len = tf.shape(x)[1]
24        attention_weights = {}
25        # (batch_size, target_seq_len, d_model)
26        x = self.embedding(x)
27        x *= tf.math.sqrt(tf.cast(self.d_model, tf.float32))
28        x += self.pos_encoding[:, :seq_len, :]
29
30        x = self.dropout(x, training=training)
31        for i in range(self.num_layers):
32            x, block1, block2 = self.dec_layers[i](x, enc_output, training,
33                                                  look_ahead_mask,
34                                                  padding_mask)
35
36            attention_weights['decoder_layer{}_block1'.format(i + 1)] =
37                ↪ block1
38            attention_weights['decoder_layer{}_block2'.format(i + 1)] =
39                ↪ block2
40
41        # x.shape == (batch_size, target_seq_len, d_model)
42        return x, attention_weights
```

```

1 sample_decoder = Decoder(num_layers=2,
2                           d_model=512,
3                           num_heads=8,
4                           dff=2048,
5                           target_vocab_size=8000)
6
7 output, attn = sample_decoder(tf.random.uniform((64, 26)),
8                               enc_output=sample_encoder_output,
9                               training=False,
10                              look_ahead_mask=None,
11                              padding_mask=None)
12
13 output.shape, attn['decoder_layer2_block2'].shape

```

```

(TensorShape([64, 26, 512]), TensorShape([64, 8, 26, 62]))

```

9 Create the Transformer

```

1 class Transformer(keras.Model):
2     def __init__(self,
3                   num_layers,
4                   d_model,
5                   num_heads,
6                   dff,
7                   input_vocab_size,
8                   target_vocab_size,
9                   rate=0.1):
10         super(Transformer, self).__init__()
11         self.encoder = Encoder(num_layers, d_model, num_heads, dff,
12                                input_vocab_size, rate)
13
14         self.decoder = Decoder(num_layers, d_model, num_heads, dff,
15                                target_vocab_size, rate)
16         self.final_layer = layers.Dense(target_vocab_size)
17
18     def call(self, inp, tar, training, enc_padding_mask, look_ahead_mask,
19             dec_padding_mask):
20         # inp.shape == (batch_size, tar_seq_len, target_vocab_size)
21
22         # (batch_size, tar_seq_len, target_vocab_size)
23         enc_output = self.encoder(inp, training, enc_padding_mask)
24
25         # dec_output.shape == (batch_size, tar_seq_len, d_model)
26         dec_output, attention_weights = self.decoder(tar, enc_output,
27                                                     ↪ training,

```

```

look_ahead_mask,

```

```

28                                     dec_padding_mask)
29     # (batch_size, tar_seq_len, target_vocab_size)
30     final_output = self.final_layer(dec_output)
31     return final_output, attention_weights

```

```

1  sample_transformer = Transformer(num_layers=2,
2                                  d_model=512,
3                                  num_heads=8,
4                                  dff=2048,
5                                  input_vocab_size=8500,
6                                  target_vocab_size=8000)
7
8  temp_input = tf.random.uniform((64, 62))
9  temp_target = tf.random.uniform((64, 26))
10
11 fn_out, _ = sample_transformer(temp_input,
12                                temp_target,
13                                training=False,
14                                enc_padding_mask=None,
15                                look_ahead_mask=None,
16                                dec_padding_mask=None)
17
18 # (batch_size, tar_seq_len, target_vocab_size)
19 fn_out.shape

```

```
TensorShape([64, 26, 8000])
```

10 Set hyperparameters

```

1  num_layers = 4
2  d_model = 128
3  dff = 512
4  num_heads = 8
5
6  input_vocab_size = tokenizer_pt.vocab_size + 2
7  target_vocab_size = tokenizer_en.vocab_size + 2
8  dropout_rate = 0.1

```

11 Optimizer

custom Adam optimizer ref. <https://arxiv.org/abs/1706.03762> (Attention is All You Need)
This formula is here.

$$lrate = d_{model}^{-0.5} * \min(step_num^{-0.5}, step_num * warmup_steps^{-1.5})$$

```

1 class CustomSchedule(keras.optimizers.schedules.LearningRateSchedule):
2     def __init__(self, d_model, warmup_steps=4000):
3         super(CustomSchedule, self).__init__()
4         self.d_model = d_model
5         self.d_model = tf.cast(self.d_model, tf.float32)
6         self.warmup_steps = warmup_steps
7
8     def __call__(self, step):
9         arg1 = step**-0.5
10        arg2 = step * (self.warmup_steps**-1.5)
11        return (self.d_model**-0.5) * tf.math.minimum(arg1, arg2)

```

```

1 learning_rate = CustomSchedule(d_model)
2 optimizer = keras.optimizers.Adam(learning_rate,
3                                     beta_1=0.9,
4                                     beta_2=0.98,
5                                     epsilon=1e-9)

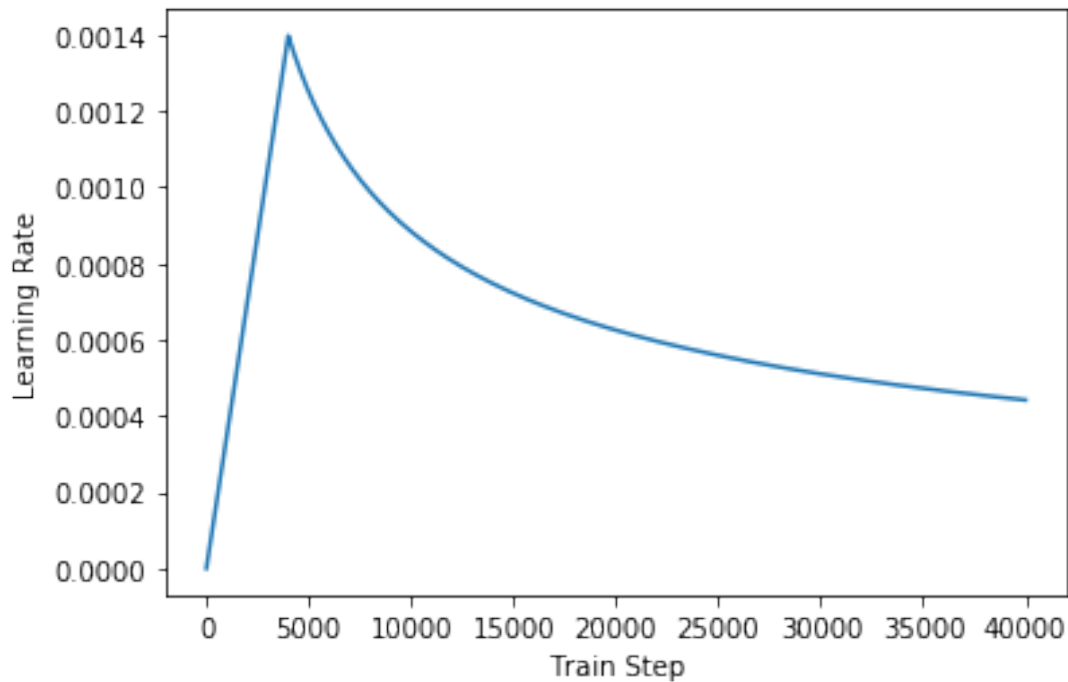
```

```

1 temp_learning_rate_schedule = CustomSchedule(d_model)
2 plt.plot(temp_learning_rate_schedule(tf.range(40000, dtype=tf.float32)))
3 plt.ylabel('Learning Rate')
4 plt.xlabel('Train Step')

```

Text(0.5, 0, 'Train Step')



12 Loss and metrics

```
1 loss_object = keras.losses.SparseCategoricalCrossentropy(from_logits=True,  
2                                                           reduction='none')  
3  
4  
5 def loss_function(real, pred):  
6     # 0 == padding  
7     mask = tf.math.logical_not(tf.math.equal(real, 0))  
8     loss_ = loss_object(real, pred)  
9  
10    mask = tf.cast(mask, dtype=loss_.dtype)  
11    loss_ *= mask  
12    return tf.reduce_mean(loss_)  
13  
14  
15 train_loss = keras.metrics.Mean(name='train_loss')  
16 train_acc = keras.metrics.SparseCategoricalAccuracy(name='train_accuracy')
```

13 Training and checkpointing

```
1 transformer = Transformer(num_layers, d_model, num_heads, dff,
2                           input_vocab_size, target_vocab_size,
3                           ↪ dropout_rate)
```

```
1 def create_masks(inp, tar):
2     # Encoder padding mask
3     enc_padding_mask = create_padding_mask(inp)
4
5     # Used in the 2nd attention block in the decoder.
6     # This padding mask is used to mask the encoder outputs.
7     dec_padding_mask = create_padding_mask(inp)
8
9     # Used in the 1st attention block in the decoder.
10    # It is used to pad and mask future tokens in the input received
11    # by the decoder.
12    look_ahead_mask = create_look_ahead_mask(tf.shape(tar)[1])
13    dec_target_padding_mask = create_padding_mask(tar)
14    combined_mask = tf.maximum(dec_target_padding_mask, look_ahead_mask)
15
16    return enc_padding_mask, combined_mask, dec_padding_mask
```

```
1 checkpoint_path = "./checkpoints/train"
2
3 ckpt = tf.train.Checkpoint(transformer=transformer, optimizer=optimizer)
4 ckpt_manager = tf.train.CheckpointManager(ckpt, checkpoint_path,
5     ↪ max_to_keep=5)
6
7 # if a checkpoint exists, restore the latest checkpoint.
8 if ckpt_manager.latest_checkpoint:
9     ckpt.restore(ckpt_manager.latest_checkpoint)
10    print('Latest checkpoint restored!!')
```

```
1 EPOCHS = 20
2 train_step_signature = [
3     tf.TensorSpec(shape=(None, None), dtype=tf.int64),
4     tf.TensorSpec(shape=(None, None), dtype=tf.int64),
5 ]
6
7 @tf.function(input_signature=train_step_signature)
8 def train_step(inp, tar):
9     tar_inp = tar[:, :-1]
10    tar_real = tar[:, 1:]
```



```

11
12     enc_padding_mask, combined_mask, dec_padding_mask = create_masks(
13         inp, tar_inp)
14
15     with tf.GradientTape() as tape:
16         predictions, _ = transformer(inp, tar_inp, True, enc_padding_mask,
17                                     combined_mask, dec_padding_mask)
18         loss = loss_function(tar_real, predictions)
19
20     gradients = tape.gradient(loss, transformer.trainable_variables)
21     optimizer.apply_gradients(zip(gradients,
22                                   ↪ transformer.trainable_variables))
23
24     train_loss(loss)
25     train_acc(tar_real, predictions)

```

```

1  # training
2  for epoch in range(EPOCHS):
3      start = time.time()
4      train_loss.reset_states()
5      train_acc.reset_states()
6      # inp -> portuguese, tar -> english
7      for (batch, (inp, tar)) in enumerate(train_dataset):
8          train_step(inp, tar)
9          if batch % 50 == 0:
10              print('Epoch {} Batch {} Loss {:.4f} Acc {:.4f}'.format(
11                  epoch + 1, batch, train_loss.result(), train_acc.result()))
12
13  if (epoch + 1) % 5 == 0:
14      ckpt_save_path = ckpt_manager.save()
15      print('Saving checkpoint for epoch {} at {}'.format(
16          epoch + 1, ckpt_save_path))
17      print('Epoch {} Loss {:.4f} Accuracy {:.4f}'.format(
18          epoch + 1, train_loss.result(), train_acc.result()))
19      print('Time taken for 1 epoch: {} secs\n'.format(time.time() -
20                                                         start))

```

```

Epoch 20 Batch 700 Loss 0.4377 Acc 0.3641
Saving checkpoint for epoch 20 at ./checkpoints/train/ckpt-4
Epoch 20 Loss 0.4379 Accuracy 0.3641
Time taken for 1 epoch: 229.03856348991394 secs

```

14 Evaluation

```

1  def evaluate(inp_sentence):
2      start_token = [tokenizer_pt.vocab_size]

```

```

3     end_token = [tokenizer_pt.vocab_size + 1]
4
5     # inp sentence is portuguese, hence adding the start and end token
6     inp_sentence = start_token + tokenizer_pt.encode(inp_sentence) +
7         ↪ end_token
8     encoder_input = tf.expand_dims(inp_sentence, 0)
9
10    # as the target is english, the first word to the transformer should be
11    ↪ the
12    # english start token.
13    decoder_input = [tokenizer_en.vocab_size]
14    output = tf.expand_dims(decoder_input, 0)
15
16    for i in range(MAX_LENGTH):
17        enc_padding_mask, combined_mask, dec_padding_mask = create_masks(
18            encoder_input, output)
19
20        # predictions.shape == (batch_size, seq_len, vocab_size)
21        predictions, attention_weights = transformer(encoder_input, output,
22                                                    False,
23                                                    ↪ enc_padding_mask,
24                                                    combined_mask,
25                                                    dec_padding_mask)
26
27        # select the last word from the seq_len dimension
28        predictions = predictions[:, -1:, :] # (batch_size, 1, vocab_size)
29
30        predicted_id = tf.cast(tf.argmax(predictions, axis=-1), tf.int32)
31
32        # return the result if the predicted_id is equal to the end token
33        if predicted_id == tokenizer_en.vocab_size + 1:
34            return tf.squeeze(output, axis=0), attention_weights
35
36        # concatenate the predicted_id to the output which is given to the
37        ↪ decoder
38        # as its input.
39        output = tf.concat([output, predicted_id], axis=-1)
40
41    return tf.squeeze(output, axis=0), attention_weights

```

```

1 def plot_attention_weights(attention, sentence, result, layer):
2     fig = plt.figure(figsize=(16, 8))
3
4     sentence = tokenizer_pt.encode(sentence)
5
6     attention = tf.squeeze(attention[layer], axis=0)
7
8     for head in range(attention.shape[0]):

```

```

9     ax = fig.add_subplot(2, 4, head + 1)
10
11     # plot the attention weights
12     ax.matshow(attention[head][: -1, :], cmap='viridis')
13
14     fontdict = {'fontsize': 10}
15
16     ax.set_xticks(range(len(sentence) + 2))
17     ax.set_yticks(range(len(result)))
18
19     ax.set_ylim(len(result) - 1.5, -0.5)
20
21     ax.set_xticklabels(['<start>'] +
22                        [tokenizer_pt.decode([i])
23                         for i in sentence] + ['<end>'],
24                        fontdict=fontdict,
25                        rotation=90)
26
27     ax.set_yticklabels([
28         tokenizer_en.decode([i])
29         for i in result if i < tokenizer_en.vocab_size
30     ],
31                        fontdict=fontdict)
32
33     ax.set_xlabel('Head {}'.format(head + 1))
34
35 plt.tight_layout()
36 plt.show()

```

```

1 def translate(sentence, plot=''):
2     result, attention_weights = evaluate(sentence)
3
4     predicted_sentence = tokenizer_en.decode(
5         [i for i in result if i < tokenizer_en.vocab_size])
6
7     print('Input: {}'.format(sentence))
8     print('Predicted translation: {}'.format(predicted_sentence))
9
10    if plot:
11        plot_attention_weights(attention_weights, sentence, result, plot)

```

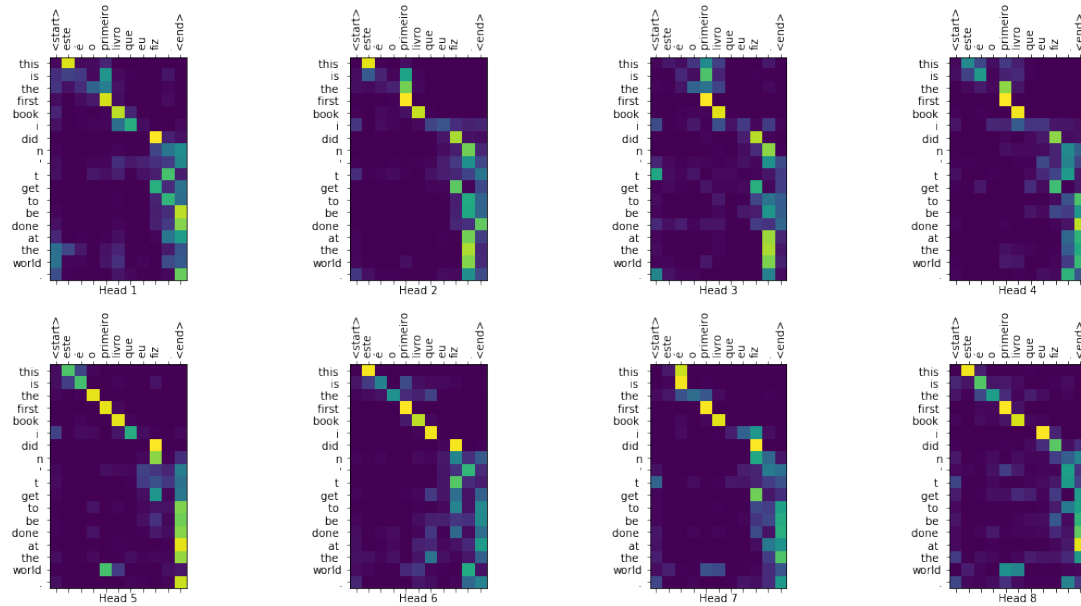
decoder layer's multi-head attention (num_heads == 8)

```

1 translate("este ^c3^a9 o primeiro livro que eu fiz.",
2           ↪ plot='decoder_layer4_block2')
3 print("Real translation: this is the first book i've ever done.")

```

Input: este é o primeiro livro que eu fiz.
 Predicted translation: this is the first book i did n't get to be done at the world .



Real translation: this is the first book i've ever done.