

```
1 from google.colab import drive
2 drive.mount('/content/drive')
```

➞ Go to this URL in a browser: [https://accounts.google.com/o/oauth2/auth?client\\_id=947318989803-6bn6qk8qdgf4n4g3pfee649](https://accounts.google.com/o/oauth2/auth?client_id=947318989803-6bn6qk8qdgf4n4g3pfee649)

Enter your authorization code:

.....

Mounted at /content/drive

```
1 # for google colab only
2 %tensorflow_version 2.x
3
4 import tensorflow as tf
5 from tensorflow.keras.preprocessing.text import Tokenizer
6 from tensorflow.keras.preprocessing import sequence
7
8 import matplotlib.pyplot as plt
9 import matplotlib.ticker as ticker
10 from sklearn.model_selection import train_test_split
11
12 import unicodedata
13 import re
14 import numpy as np
15 import os
16
17 ### Reference article ###
18 ### https://www.tensorflow.org/tutorials/text/nmt\_with\_attention
```

➞ TensorFlow 2.x selected.

```
1 # this path is for google drive mounted in google colab. If you want to use file from local directory
2 # I request you to modify this path
3
4 file_path = "/content/drive/My Drive/Stevens/Fall 2019/NLP/fra.txt"
```

```
1 # Converts the unicode file to ascii
```

```
2 def unicode_to_ascii(s):
3     return ''.join(c for c in unicodedata.normalize('NFD', s)
4         if unicodedata.category(c) != 'Mn')
5
6 # data preprocessing
7 def data_preprocessing(w):
8     w = unicode_to_ascii(w.lower().strip())
9
10    # adding extra space before punctuation
11    w = re.sub(r"([?!.,:;])", r" \1 ", w)
12    w = re.sub(r'[" "]', " ", w)
13
14    # removing unnecessary characters
15    w = re.sub(r"[^a-zA-Z?!.,:;]+", " ", w)
16
17    w = w.rstrip().strip()
18
19    # adding a start and end token
20    w = '<start> ' + w + ' <end>'
21    return w
22
23 # dataset creation
24 def getPreprocessedLanguagePairs(path, sample_size):
25     with open(path, encoding='UTF-8') as f:
26         lines = f.read().strip().split('\n')
27         lang_pairs = [[data_preprocessing(w) for w in l.split('\t')] for l in lines[:sample_size]]
28
29     return zip(*lang_pairs)
30
31 def tokenize_corpus(corpus):
32     tokenizer = Tokenizer(filters='')
33     tokenizer.fit_on_texts(corpus)
34     tokenized_text = tokenizer.texts_to_sequences(corpus)
35     tokenized_text = sequence.pad_sequences(tokenized_text, padding='post')
36
37     return tokenized_text, tokenizer
38
39 def getTokenizedPairs(path, sample_size=None):
40     eng set, fr set, = getPreprocessedLanguagePairs(path, sample size)
```

```

41 eng_tokenized_text, eng_set_tokenizer = tokenize_corpus(eng_set)
42 fr_tokenized_text, fr_set_tokenizer = tokenize_corpus(fr_set)
43
44 return eng_tokenized_text, fr_tokenized_text, eng_set_tokenizer, fr_set_tokenizer

1 # Google colab is crashing when using whole dataset, so taking small sample out of it
2 sample_size = 50000
3 eng_tokenized_text, fr_tokenized_text, eng_set_tokenizer, fr_set_tokenizer = getTokenizedPairs(file_path,
4                                                                                               sample_size)
5
6 max_length_fr = max(len(text) for text in fr_tokenized_text)
7 max_length_eng = max(len(text) for text in eng_tokenized_text)

1 # Splitting set into train and test
2 X_train, X_test, y_train, y_test = train_test_split(eng_tokenized_text, fr_tokenized_text, test_size=0.2, random_st
3
4 # Show length
5 print(len(X_train), len(y_train), len(X_test), len(y_test))

☞ 40000 40000 10000 10000

1 # Hyperparameters
2 BUFFER_SIZE = len(X_train)
3 BATCH_SIZE = 64
4 steps_per_epoch = len(X_train)//BATCH_SIZE
5 embedding_dim = 100
6 units = 512
7 vocab_en_size = len(eng_set_tokenizer.word_index)+1
8 vocab_fr_size = len(fr_set_tokenizer.word_index)+1
9
10 dataset = tf.data.Dataset.from_tensor_slices((X_train, y_train)).shuffle(BUFFER_SIZE)
11 dataset = dataset.batch(BATCH_SIZE, drop_remainder=True)
12
13 input_batch, target_batch = next(iter(dataset))

1 # Encoder RNN Class
2 class Encoder(tf.keras.Model):

```

```
3 def __init__(self, vocab_size, embedding_dim, encoder_units, batch_size):
4     super(Encoder, self).__init__()
5     self.batch_size = batch_size
6     self.encoder_units = encoder_units
7     self.embedding = tf.keras.layers.Embedding(vocab_size, embedding_dim)
8     self.gru = tf.keras.layers.GRU(self.encoder_units,
9                                     return_sequences=True,
10                                    return_state=True)
11
12 def call(self, x, hidden):
13     x = self.embedding(x)
14     output, state = self.gru(x, initial_state = hidden)
15     return output, state
16
17 def initialize_hidden_state(self):
18     return tf.zeros((self.batch_size, self.encoder_units))
19
20 # Attention Class
21 class AttentionLayer(tf.keras.layers.Layer):
22     def __init__(self, units):
23         super(AttentionLayer, self).__init__()
24         self.W1 = tf.keras.layers.Dense(units)
25         self.W2 = tf.keras.layers.Dense(units)
26         self.V = tf.keras.layers.Dense(1)
27
28     def call(self, query, values):
29         hidden_with_time_axis = tf.expand_dims(query, 1)
30
31         score = self.V(tf.nn.tanh(
32             self.W1(values) + self.W2(hidden_with_time_axis)))
33
34         # attention_weights.shape = (batch_size, max_length, 1)
35         attention_weights = tf.nn.softmax(score, axis=1)
36
37         # context_vector.shape = (batch_size, hidden_size)
38         context_vector = attention_weights * values
39         context_vector = tf.reduce_sum(context_vector, axis=1)
40
41         return context_vector, attention_weights
```

```

42
43 # Decoder RNN class
44 class Decoder(tf.keras.Model):
45     def __init__(self, vocab_size, embedding_dim, decoder_units, batch_size):
46         super(Decoder, self).__init__()
47         self.batch_size = batch_size
48         self.decoder_units = decoder_units
49         self.embedding = tf.keras.layers.Embedding(vocab_size, embedding_dim)
50         self.gru = tf.keras.layers.GRU(self.decoder_units,
51                                         return_sequences=True,
52                                         return_state=True)
53         self.fc = tf.keras.layers.Dense(vocab_size)
54
55         # attention
56         self.attention = AttentionLayer(self.decoder_units)
57
58     def call(self, x, hidden, enc_output):
59         context_vector, attention_weights = self.attention(hidden, enc_output)
60
61         x = self.embedding(x)
62         x = tf.concat([tf.expand_dims(context_vector, 1), x], axis=-1)
63
64         # passing the concatenated vector to the LSTM
65         output, state = self.gru(x)
66         output = tf.reshape(output, (-1, output.shape[2]))
67         x = self.fc(output)
68
69         return x, state, attention_weights

```

```

1 encoder = Encoder(vocab_en_size, embedding_dim, units, BATCH_SIZE)
2
3 # sample input
4 encoder_hidden = encoder.initialize_hidden_state()
5 encoder_output, encoder_hidden = encoder(input_batch, encoder_hidden)
6
7 attention_layer = AttentionLayer(10)
8 attention_result, attention_weights = attention_layer(encoder_hidden, encoder_output)
9
10 # ... (rest of the code) ...

```

```

10 decoder = Decoder(vocab_tr_size, embedding_dim, units, BATCH_SIZE)
11
12 sample_decoder_output, _, _ = decoder(tf.random.uniform((64, 1)),
13                                     encoder_hidden, encoder_output)

1 # Loss
2 optimizer = tf.keras.optimizers.Adam()
3 loss_object = tf.keras.losses.SparseCategoricalCrossentropy(from_logits=True, reduction='none')
4
5 def loss_function(real, pred):
6     mask = tf.math.logical_not(tf.math.equal(real, 0))
7     loss_ = loss_object(real, pred)
8
9     mask = tf.cast(mask, dtype=loss_.dtype)
10    loss_ *= mask
11
12    return tf.reduce_mean(loss_)

1 # Model training
2 @tf.function
3 def training(en, fr, enc_hidden):
4     loss = 0
5
6     with tf.GradientTape() as tape:
7         enc_output, enc_hidden = encoder(en, enc_hidden)
8         dec_hidden = enc_hidden
9         dec_input = tf.expand_dims([fr_set_tokenizer.word_index['<start>']] * BATCH_SIZE, 1)
10
11        # taking target as input to next
12        for t in range(1, fr.shape[1]):
13            # passing enc_output to the decoder
14            predictions, dec_hidden, _ = decoder(dec_input, dec_hidden, enc_output)
15            loss += loss_function(fr[:, t], predictions)
16            dec_input = tf.expand_dims(fr[:, t], 1)
17
18        # calculate batch loss
19        batch_loss = (loss / int(fr.shape[1]))
20
21

```

```

21 variables = encoder.trainable_variables + decoder.trainable_variables
22 gradients = tape.gradient(loss, variables)
23 optimizer.apply_gradients(zip(gradients, variables))
24
25 return batch_loss

1 EPOCHS = 201
2
3 for epoch in range(1, EPOCHS):
4     # print("Epoch #%d"%epoch)
5     enc_hidden = encoder.initialize_hidden_state()
6     total_loss = 0
7
8     for (batch, (en, fr)) in enumerate(dataset.take(steps_per_epoch)):
9         batch_loss = training(en, fr, enc_hidden)
10        total_loss += batch_loss
11
12    # printing loss after every 20 epochs
13    if(epoch%20==0):
14        print('Epoch {} Loss {:.4f}'.format(epoch,
15                                              total_loss / steps_per_epoch))

↳ Epoch 20 Loss 0.0957
Epoch 40 Loss 0.0663
Epoch 60 Loss 0.0590
Epoch 80 Loss 0.0551
Epoch 100 Loss 0.0534
Epoch 120 Loss 0.0511
Epoch 140 Loss 0.0491
Epoch 160 Loss 0.0471
Epoch 180 Loss 0.0458
Epoch 200 Loss 0.0456

1 def translate(sentence):
2     inputs = [eng_set_tokenizer.word_index[i] for i in sentence.split(' ')]
3     inputs = tf.keras.preprocessing.sequence.pad_sequences([inputs],
4                                                             maxlen=max_length_eng,
5                                                             padding='post')
6     inputs = tf.convert_to_tensor(inputs)
7

```

```

/
8     result = ''
9
10    hidden = [tf.zeros((1, units))]
11    enc_out, enc_hidden = encoder(inputs, hidden)
12
13    dec_hidden = enc_hidden
14    dec_input = tf.expand_dims([fr_set_tokenizer.word_index['<start>']], 0)
15
16    for t in range(max_length_fr):
17        predictions, dec_hidden, attention_weights = decoder(dec_input,
18                                                            dec_hidden,
19                                                            enc_out)
20
21        attention_weights = tf.reshape(attention_weights, (-1, ))
22
23        predicted_id = tf.argmax(predictions[0]).numpy()
24
25        result += fr_set_tokenizer.index_word[predicted_id] + ' '
26
27        if fr_set_tokenizer.index_word[predicted_id] == '<end>':
28            return result
29
30        # the predicted ID is fed back into the model
31        dec_input = tf.expand_dims([predicted_id], 0)
32
33    return result


1 # train-test set split in raw sentence format
2 full_en_, full_fr_, _ = getPreprocessedLanguagePairs(file_path, sample_size)
3
4 X_train_en, X_test_en, y_train_fr, y_test_fr = train_test_split(list(full_en_), list(full_fr_), test_size=0.2, rand
5
6 len(y_test_fr), len(X_test_en)

☐➔ (10000, 10000)

```



## ▼ Translation Demo

```

1 for english_sent in X_test_en[:10]:
2     french_translation = translate(english_sent)
3     print('English Sentence: %s' % (english_sent))
4     print('Predicted French translation: {}'.format(french_translation))

☐ English Sentence: <start> add a little milk . <end>
Predicted French translation: ajoute un peu de lait . <end>
English Sentence: <start> i think it s funny . <end>
Predicted French translation: je pense que c est drôle . <end>
English Sentence: <start> i liked the movie . <end>
Predicted French translation: j ai aimé le film . <end>
English Sentence: <start> don t judge me . <end>
Predicted French translation: ne me jugez pas . <end>
English Sentence: <start> please come this way . <end>
Predicted French translation: par ici , je vous prie . <end>
English Sentence: <start> can you do it faster ? <end>
Predicted French translation: arrives tu a le faire plus vite ? <end>
English Sentence: <start> the airplane is ready . <end>
Predicted French translation: le vent est prêt . <end>
English Sentence: <start> no one does that . <end>
Predicted French translation: personne ne le sait . <end>
English Sentence: <start> this is tasteless . <end>
Predicted French translation: c est absurde . <end>
English Sentence: <start> may i join you ? <end>
Predicted French translation: puis je me joindre a vous ? <end>

```

## ▼ BLEU Calculation for test set

```

1 from nltk.translate.bleu_score import corpus_bleu
2 from nltk.translate import bleu_score
3 import nltk
4 nltk.download('punkt')
5
6 def remove_unnecessary_tokens(sentence):
7
8     # remove <start> and <end> tokens
9     w = re.sub(r'<start>', "", sentence)

```

```

10     w = re.sub(r'<end>', "", w)
11
12     # removing unnecessary characters
13     w = re.sub(r"^[a-zA-Z]", " ", w)
14
15     return w.strip()

```

```

[ ] [nltk_data] Downloading package punkt to /root/nltk_data...
[nltk_data]   Unzipping tokenizers/punkt.zip.

```

```

1 references = []
2 hypotheses = []
3
4 for i in range(len(X_test_en)):
5     french_translation = translate(X_test_en[i])
6
7     #remove <start> and <end> tokens and other unnecessary characters
8     fr_ref = nltk.word_tokenize(remove_unnecessary_tokens(y_test_fr[i]))
9     fr_trans = nltk.word_tokenize(remove_unnecessary_tokens(french_translation))
10
11     references.append([fr_ref])
12     hypotheses.append(fr_trans)

```

```
1 references[:10]
```

```

[ ] [[['ajoutez', 'un', 'peu', 'de', 'lait']],
    [['je', 'pense', 'que', 'c', 'est', 'drole']],
    [['j', 'ai', 'aime', 'le', 'film']],
    [['ne', 'me', 'juge', 'pas']],
    [['par', 'ici', 's', 'il', 'vous', 'plait']],
    [['arrivez', 'vous', 'a', 'le', 'faire', 'plus', 'vite']],
    [['l', 'avion', 'est', 'pret']],
    [['personne', 'ne', 'fait', 'cela']],
    [['ca', 'n', 'a', 'aucun', 'gout']],
    [['puis', 'je', 'me', 'joindre', 'a', 'vous']]]

```

```
1 hypotheses[:10]
```

```
[ ]
```

```
[['ajoute', 'un', 'peu', 'de', 'lait'],  
 ['je', 'pense', 'que', 'c', 'est', 'drole'],  
 ['j', 'ai', 'aime', 'le', 'film'],  
 ['ne', 'me', 'jugez', 'pas'],  
 ['par', 'ici', 'je', 'vous', 'prie'],  
 ['arrives', 'tu', 'a', 'le', 'faire', 'plus', 'vite'],  
 ['le', 'vent', 'est', 'pret'],  
 ['personne', 'ne', 'le', 'sait'],  
 ['c', 'est', 'absurde'],  
 ['puis', 'je', 'me', 'joindre', 'a', 'vous']]
```

```
1 corpus_bleu(references, hypotheses)
```

```
↳ 0.3106569217107635
```

```
1
```

