



Jonathan Hui

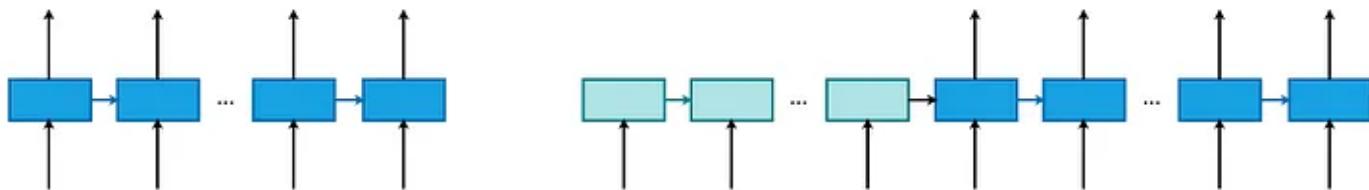
[Follow](#)Feb 1, 2021 · 10 min read · [▶ Listen](#)[Save](#)

...

TensorFlow Sequence to Sequence Model Examples



Sequence-to-sequence models are particularly popular in NLP. This article, as part of the TensorFlow series, will cover examples for the sequence to sequence model.



The examples include:

- Shakespear text generation with GRU,
- Seq2seq language translation with attention,
- Image captioning with visual attention and Inception V3 for feature extraction.

Shakespear text generation with GRU

Given a large file with Shakespear's writings,

```
First Citizen:  
Before we proceed any further, hear me speak.  
  
All:  
Speak, speak.  
  
First Citizen:  
You are all resolved rather to die than to famish?  
  
All:  
Resolved. resolved.  
  
First Citizen:  
First, you know Caius Marcius is chief enemy to the people.
```

we can use its content as a stream of characters to train a model to write like Shakespeare.

For example, during inference, we start with an initial text, say “ROMEO: ”. The model predicts what is the next character once at a time.

```
print(generate_text(model, start_string=u"ROMEO: "))
```

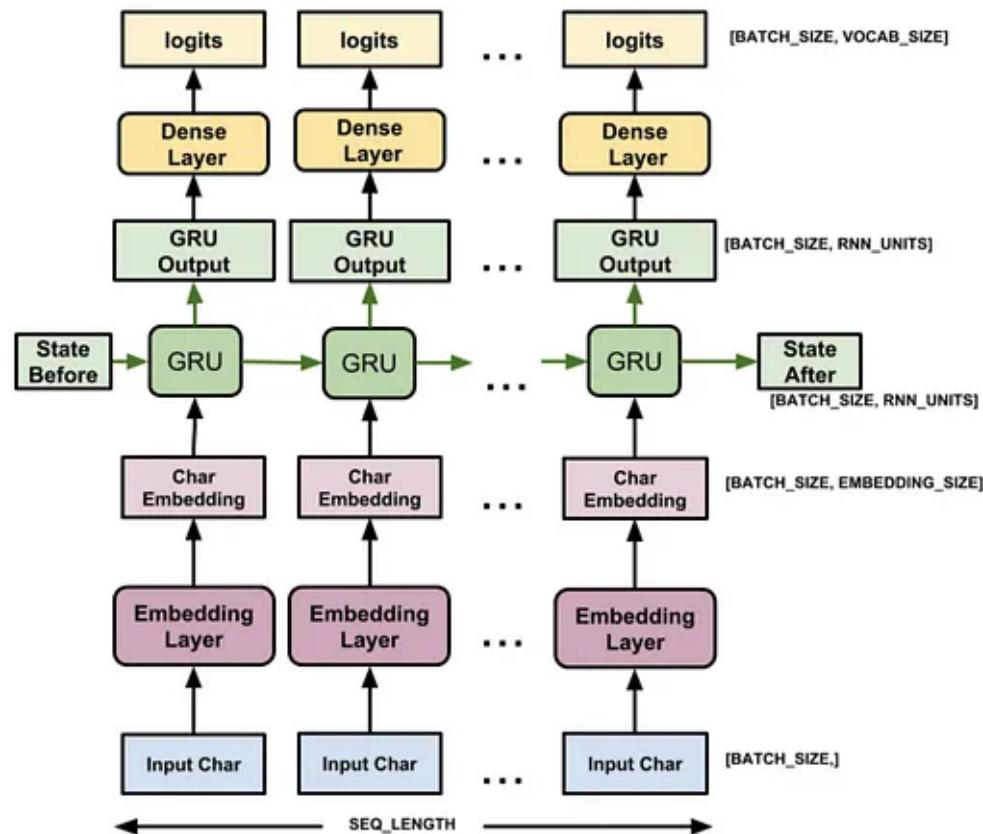
ROMEO: ghast I cut go,
 Know the normander and the wrong:
 To our Morsuis misdress are behiod;
 And after as if no other husion.

VALERIS:
 Your father and of worms?

LADY GREY:
 Your hot can dosta.

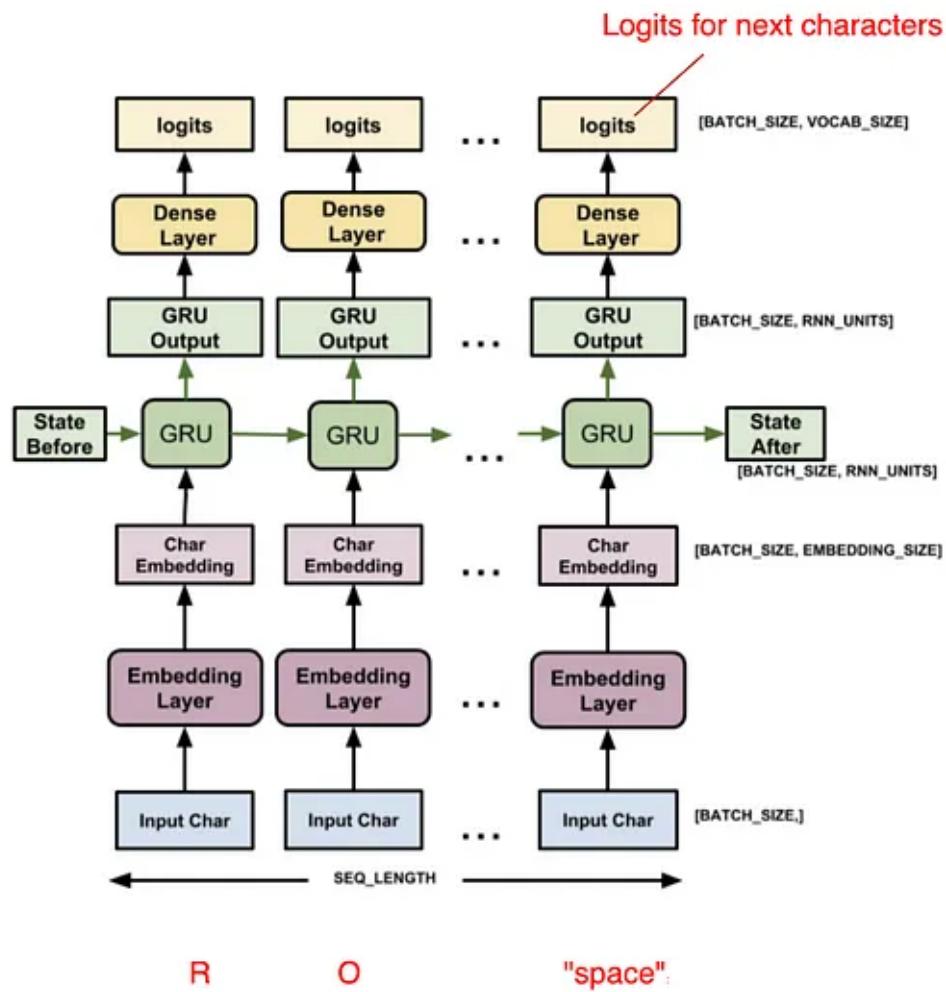
WARWICK:
 Then, atient the bade, truckle aid,

This can be done with the GRU model below with embedding and dense layers. The next character is the character predicted in the last GRU cell.



Source (The round rectangles represent operation and the square rectangles represent Tensors.)

Let trace it from the beginning with the initial text. The logits output of the last cell will be used to determine the next character after “ROMEO:”.



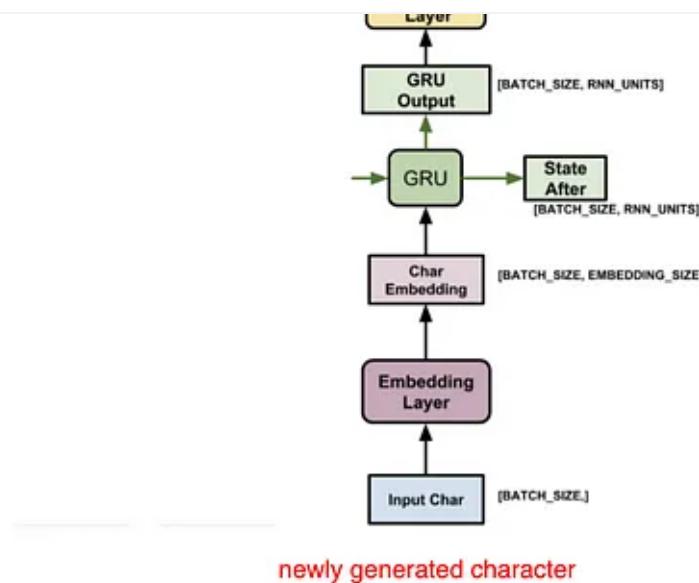
We take the generated character as the input to the model again and generates the next character.

[Open in app](#)[Resume Membership](#)

Search Medium



v



As we keep repeating the process, it generates a script similar to the one below. Even though not grammatically sound, it can be improved with more complex designs. For now, we stick with the simple model.

```

ROMEO: ghast I cut go,
Know the normander and the wrong:
To our Morsuis misdress are behiod;
And after as if no other husion.
  
```

Let's go back to the code. First, we load the Shakespear file and uses its unique characters (65) to form a vocabulary (this vocabulary contains characters instead of words). Then, we create a mapping between a character and an integer index.

```
1 # Modified from
2 # https://www.tensorflow.org/tutorials/text/text_generation
3
4 import tensorflow as tf
5 import numpy as np
6 import os
7 import time
8
9 path_to_file = tf.keras.utils.get_file('shakespeare.txt',
10                                     'https://storage.googleapis.com/download.t
11
12 # Read, then decode for py2 compat.
13 # text contains 1115394 characters
14 text = open(path_to_file, 'rb').read().decode(encoding='utf-8')
15
16 # The unique characters in the file
17 vocab = sorted(set(text)) # 65 unique characters
18
19 # Creating a mapping from unique characters to indices
20 char2idx = {u: i for i, u in enumerate(vocab)}
21 idx2char = np.array(vocab)
```

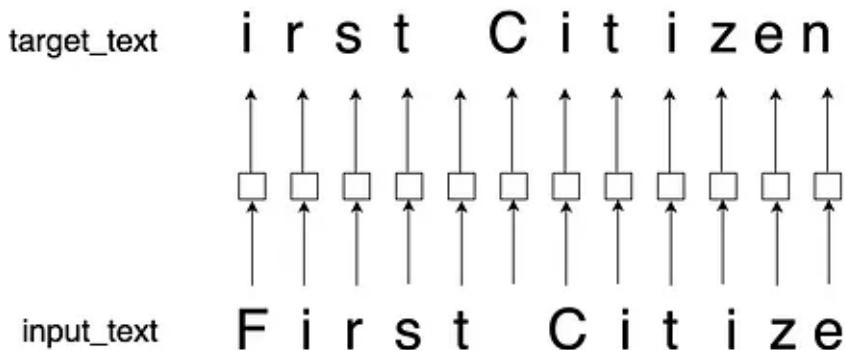
And we convert the characters in the downloaded Shakespear text into a long sequence of integer indexes. We slice and batch it to create the dataset “*sequences*”. Next, we create a *dataset* with an additional field holding the target text in line 39 as the labels.

```

23     text_as_int = np.array([char2idx[c] for c in text])
24
25     # The maximum length sentence you want for a single input in characters
26     seq_length = 100
27     examples_per_epoch = len(text) // (seq_length + 1)
28
29     # Create training examples / targets
30     char_dataset = tf.data.Dataset.from_tensor_slices(text_as_int)
31
32     sequences = char_dataset.batch(seq_length + 1, drop_remainder=True)
33
34
35     # Create the target text by left shift of one character
36     def split_input_target(chunk):
37         input_text = chunk[:-1]
38         target_text = chunk[1:]
39         return input_text, target_text
40
41
42     dataset = sequences.map(split_input_target)

```

Given an input character, the corresponding character in the target text (the character we want to predict) should be the next character of the input text. So these labels are simply formed by shifting the input left by one.



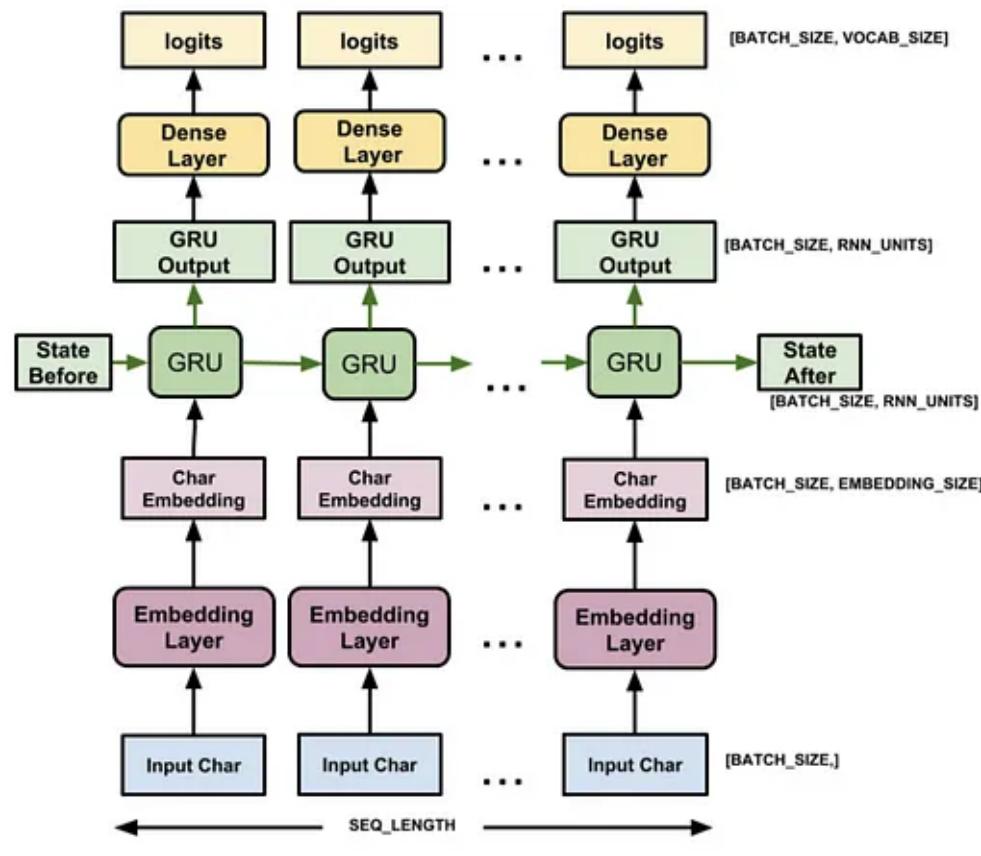
Then, we shuffle and batch the dataset.

```
44     # Create training batches
45     # Batch size
46     BATCH_SIZE = 64
47
48     # Buffer size to shuffle the dataset
49     # (TF data is designed to work with possibly infinite sequences,
50     # so it doesn't attempt to shuffle the entire sequence in memory. Instead,
51     # it maintains a buffer in which it shuffles elements).
52     BUFFER_SIZE = 10000
53
54     dataset = dataset.shuffle(BUFFER_SIZE).batch(BATCH_SIZE, drop_remainder=True)
```

Then we build the model.

```
56     # Build The Model
57     # Length of the vocabulary in chars
58     vocab_size = len(vocab)
59
60     # The embedding dimension
61     embedding_dim = 256
62
63     # Number of RNN units
64     rnn_units = 1024
65
66
67     def build_model(vocab_size, embedding_dim, rnn_units, batch_size):
68         model = tf.keras.Sequential([
69             tf.keras.layers.Embedding(vocab_size, embedding_dim,
70                                     batch_input_shape=[batch_size, None]),
71             tf.keras.layers.GRU(rnn_units,
72                                 return_sequences=True,
73                                 stateful=True,
74                                 recurrent_initializer='glorot_uniform'),
75             tf.keras.layers.Dense(vocab_size)
76         ])
77
78         return model
79
80     model = build_model(
81         vocab_size=len(vocab),
82         embedding_dim=embedding_dim,
83         rnn_units=rnn_units,
84         batch_size=BATCH_SIZE)
```

It contains an embedding layer, a GRU, and a dense layer.



[Source](#)

And here is a summary of the model.

Layer (type)	Output Shape	Param #
<hr/>		
embedding (Embedding)	(64, None, 256)	16640
gru (GRU)	(64, None, 1024)	3938304
dense (Dense)	(64, None, 65)	66625
<hr/>		
Total params: 4,021,569		
Trainable params: 4,021,569		
Non-trainable params: 0		
<hr/>		

Then, we train the model and create checkpoints.

```
87 # Train the model
88 def loss(labels, logits):
89     return tf.keras.losses.sparse_categorical_crossentropy(labels, logits, from_logits=True)
90
91
92 model.compile(optimizer='adam', loss=loss)
93
94 # Directory where the checkpoints will be saved
95 checkpoint_dir = './training_checkpoints'
96 # Name of the checkpoint files
97 checkpoint_prefix = os.path.join(checkpoint_dir, "ckpt_{epoch}")
98
99 checkpoint_callback = tf.keras.callbacks.ModelCheckpoint(
100     filepath=checkpoint_prefix,
101     save_weights_only=True)
102
103 EPOCHS = 10
104
105 history = model.fit(dataset, epochs=EPOCHS, callbacks=[checkpoint_callback])
```

Next, we reload the model from the checkpoint.

```
107 model = build_model(vocab_size, embedding_dim, rnn_units, batch_size=1)
108 model.load_weights(tf.train.latest_checkpoint(checkpoint_dir))
109
110 model.build(tf.TensorShape([1, None]))
```

With the model restored, we start generating text in the Shapespear style with the initial text “ROMEO: ”

```

113     def generate_text(model, start_string):
114         # Evaluation step (generating text using the learned model)
115
116         # Number of characters to generate
117         num_generate = 1000
118
119         # Converting our start string to numbers (vectorizing)
120         input_eval = [char2idx[s] for s in start_string]
121         input_eval = tf.expand_dims(input_eval, 0)
122
123         # Empty string to store our results
124         text_generated = []
125
126         # Low temperature results in more predictable text.
127         # Higher temperature results in more surprising text.
128         # Experiment to find the best setting.
129         temperature = 1.0
130
131         # Here batch size == 1
132         model.reset_states()
133         for i in range(num_generate):
134             predictions = model(input_eval)
135             # remove the batch dimension
136             predictions = tf.squeeze(predictions, 0)
137
138             # using a categorical distribution to predict the character returned by the model
139             predictions = predictions / temperature
140             predicted_id = tf.random.categorical(predictions, num_samples=1)[-1, 0].numpy()
141
142             # Pass the predicted character as the next input to the model
143             # along with the previous hidden state
144             input_eval = tf.expand_dims([predicted_id], 0)
145
146             text_generated.append(idx2char[predicted_id])
147
148         return (start_string + ''.join(text_generated))
149
150
151     print(generate_text(model, start_string=u"ROMEO: "))

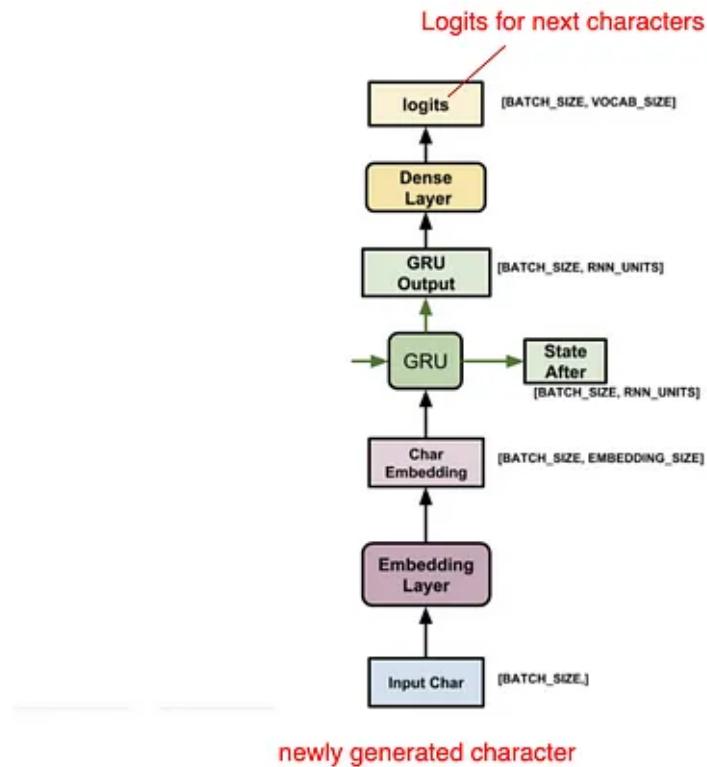
```

We convert this text into the corresponding sequence of integer indexes. Then, we fit it into the model which outputs a Tensor logits with the shape (1, 7, 65). One prediction for each input character in the initial text “ROMEO:” (batch size, sequence length, vocabulary size).

In line 140, we use `tf.random.categorical` to sample a character prediction in each sequence using the logit values. The output will have the shape of (7, 1). We are only

interested in the last GRU cell (what is next after the last character), so we take the last sampled value [-1, 0] as the next generated character.

For the rest of the “for” loop iterations, we just use the newly generated character as input to the model and predict the next character. We repeat the iterations until 1000 characters are generated.



Seq2seq language translation with attention

This example trains a sequence to sequence (seq2seq) model for Spanish to English translation. First, we download the dataset files. The following is part of the Spanish to English translation file.

```

Go.      Ve.
Go.      Vete.
Go.      Vaya.
Go.      Váyase.
Hi.      Hola.
Run!    ;Corre!
Run.     Corred.
Who?    ;Quién?
Fire!   ;Fuego!
Fire!   ;Incendio!
Fire!   ;Disparad!
Help!   ;Ayuda!
Help!   ;Socorro! ;Auxilio!
Help!   ;Auxilio!
Jump!   ;Salta!
Jump.   Salte.
Stop!   ;Parad!
Stop!   ;Para!
Stop!   ;Pare!
Wait!   ;Espera!
Wait.   Esperen.
Go on.  Continúa.
Go on.  Continúe.
Hello!  Hola.
I ran.  Corri.
I ran.  Corria.
I try.  Lo intento.
I won!  ;He ganado!
Oh no!  ;Oh, no!
Relax.  Tomátelo con soda.
Smile.  Sonrie.
Attack! ;Al ataque!
Attack! ;Atacad!
Get up. Levanta.
Go now. Ve ahora mismo.

```

```

1  # Modified from
2  # https://www.tensorflow.org/tutorials/text/nmt_with_attention
3
4  import tensorflow as tf
5
6  import matplotlib.pyplot as plt
7  import matplotlib.ticker as ticker
8  from sklearn.model_selection import train_test_split
9
10 import unicodedata
11 import re
12 import numpy as np
13 import os
14 import io
15 import time
16
17 # Download the file
18 path_to_zip = tf.keras.utils.get_file(
19     'spa-eng.zip', origin='http://storage.googleapis.com/download.tensorflow.org/data/spa-eng.zip',
20     extract=True)
21
22 path_to_file = os.path.dirname(path_to_zip) + "/spa-eng/spa.txt"

```

Next, we create methods to preprocess a sentence.

```
u"May I borrow this book?"      → <start> may i borrow this book ? <end>
u"¿Puedo tomar prestado este libro?" → b'<start> \xc2\xbf puedo tomar prestado este libro ? <end>'
```

```

25     # Converts the unicode file to ascii
26     def unicode_to_ascii(s):
27         return ''.join(c for c in unicodedata.normalize('NFD', s)
28                         if unicodedata.category(c) != 'Mn')
29
30
31     def preprocess_sentence(w):
32         w = unicode_to_ascii(w.lower().strip())
33
34         # creating a space between a word and the punctuation following it
35         # eg: "he is a boy." => "he is a boy ."
36         # Reference:- https://stackoverflow.com/questions/3645931/python-padding-punctuation-with-space
37         w = re.sub(r"([?.!,\xc1])", r" \1 ", w)
38         w = re.sub(r'[" "]+', " ", w)
39
40         # replacing everything with space except (a-z, A-Z, ".", "?", "!", " ", ",")
41         w = re.sub(r"[\xa1-\xd7?.!,\xc1]+", " ", w)
42
43         w = w.strip()
44
45         # adding a start and an end token to the sentence
46         # so that the model know when to start and stop predicting.
47         w = '<start> ' + w + ' <end>'
48
        return w

```

Given a translation file, *create_dataset* generates samples for the target text and the original text and *tokenize* converts the text into a sequence of integers.

```

path_to_file = os.path.dirname(path_to_zip)+"/spa-eng/spa.txt"

en, sp = create_dataset(path_to_file, None)

51  # 1. Remove the accents
52  # 2. Clean the sentences
53  # 3. Return word pairs in the format: [ENGLISH, SPANISH]
54 def create_dataset(path, num_examples):
55     lines = io.open(path, encoding='UTF-8').read().strip().split('\n')

56
57     word_pairs = [[preprocess_sentence(w) for w in l.split('\t')] for l in lines[:num_examples]]

58
59     return zip(*word_pairs)

60
61
62  # Convert text into a sequence of int
63  # Lang contains 30,000 sentences/phrases
64 def tokenize(lang):
65     lang_tokenizer = tf.keras.preprocessing.text.Tokenizer(
66         filters='')
67     lang_tokenizer.fit_on_texts(lang)

68
69     tensor = lang_tokenizer.texts_to_sequences(lang)

70
71     # Shape (30000, 16) int32
72     tensor = tf.keras.preprocessing.sequence.pad_sequences(tensor,
73                                         padding='post')

74
75     return tensor, lang_tokenizer

```

Now, *load_dataset* put them together in loading source (Spanish) and target (English) text into samples of source and target integer sequences. It also returns the tokenizers used.

```

78  # Given a translation file containing source and target text,
79  # return source and target integer sequence and tokenizers
80 def load_dataset(path, num_examples=None):
81     # creating cleaned input, output pairs
82     targ_lang, inp_lang = create_dataset(path, num_examples)

83
84     input_tensor, inp_lang_tokenizer = tokenize(inp_lang)
85     target_tensor, targ_lang_tokenizer = tokenize(targ_lang)

86
87     return input_tensor, target_tensor, inp_lang_tokenizer, targ_lang_tokenizer

```

```

> targ_lang = {tuple: 30000}('<start> go . <end>', '<start> go . <end>', '<start> go . <end>', '<start> go . <end>',
> targ_lang_tokenizer = {Tokenizer} <keras_preprocessing.text.Tokenizer object at 0x7fb99df8cf70>
> target_tensor = {ndarray: (30000, 11)} [[ 1 36 3 ... 0 0 0], [ 1 36 3 ... 0 0 0], [ 1 36 3 ... 0 0 0],

```

Following is a sample of the integer sequences for the source and the target text.

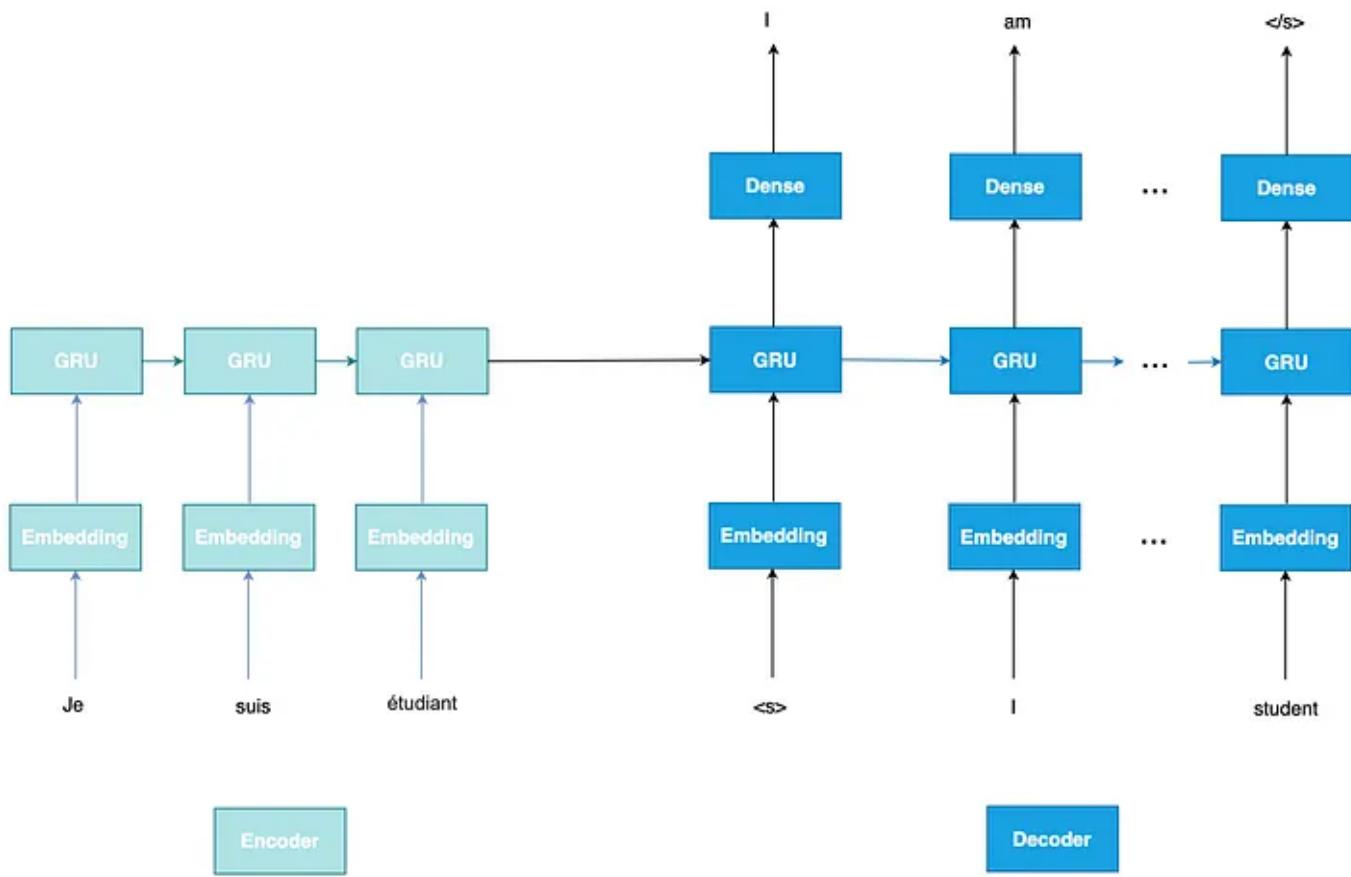
```
Input Language; index to word mapping
1 ----> <start>
6379 ----> dese
395 ----> vuelta
32 ----> ,
22 ----> por
50 ----> favor
3 ----> .
2 ----> <end>
```

```
Target Language; index to word mapping
1 ----> <start>
56 ----> please
205 ----> turn
197 ----> over
3 ----> .
2 ----> <end>
```

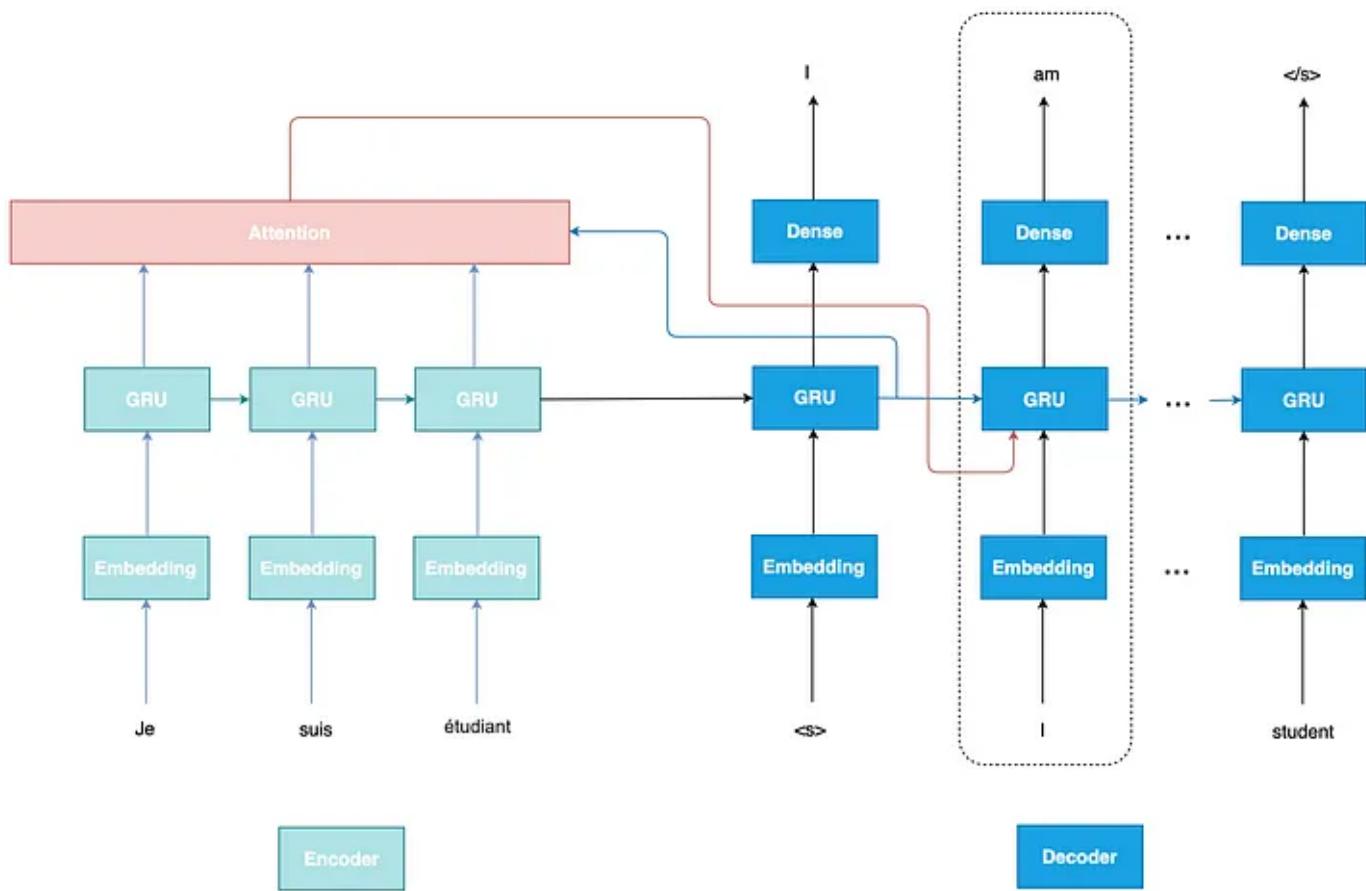
Here is the boilerplate code for creating a smaller dataset for training and validation.

```
90      # Smaller sample size for faster training
91      # Try experimenting with the size of that dataset
92      num_examples = 30000
93      input_tensor, target_tensor, inp_lang, targ_lang = load_dataset(
94                                  path_to_file, num_examples)
95
96      # Calculate max_length of the target tensors
97      max_length_targ, max_length_inp = target_tensor.shape[1], input_tensor.shape[1]
98
99      # Creating training and validation sets using an 80-20 split
100     # 24000 samples for training 6000 for validation
101     input_tensor_train, input_tensor_val, target_tensor_train, target_tensor_val = \
102         train_test_split(input_tensor, target_tensor, test_size=0.2)
103
104     BUFFER_SIZE = len(input_tensor_train)
105     BATCH_SIZE = 64
106     steps_per_epoch = len(input_tensor_train) // BATCH_SIZE
107     embedding_dim = 256
108     units = 1024
109     vocab_inp_size = len(inp_lang.word_index) + 1
110     vocab_tar_size = len(targ_lang.word_index) + 1
111
112     dataset = tf.data.Dataset.from_tensor_slices(
113         (input_tensor_train, target_tensor_train)).shuffle(BUFFER_SIZE)
114     dataset = dataset.batch(BATCH_SIZE, drop_remainder=True)
```

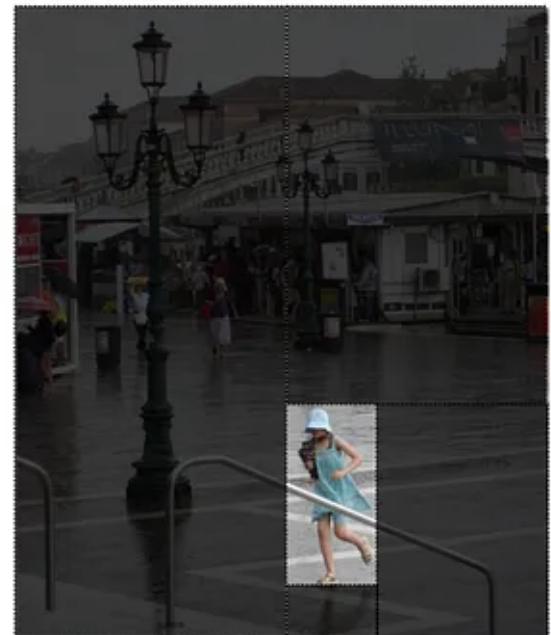
The diagram below contains an encoder to extract the **input context** for “Je suis étudiant” and a decoder to translate this context into “I am a student”.



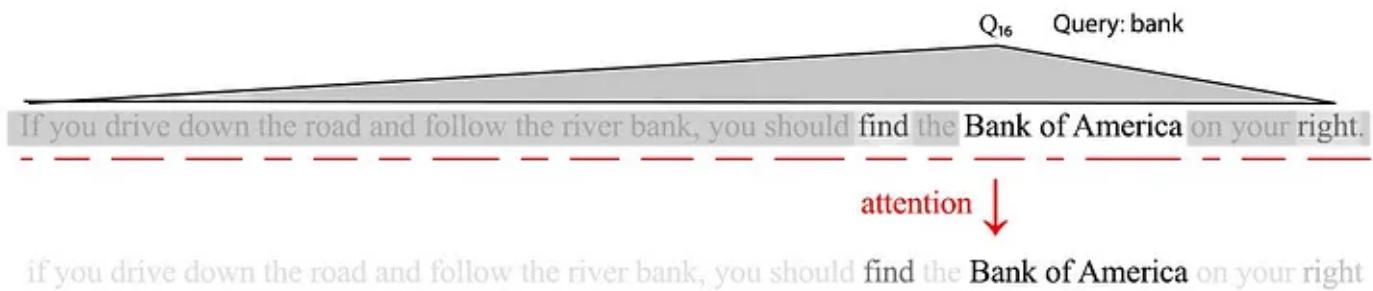
But our example has an additional attention module. Without attention, the decoder predicts the next word based on the predicted word and the GRU hidden states in the last step.



With attention, we also use the input context. But not the complete input context, we use the last GRU hidden states to select the part of the input context that we should pay attention to. In computer vision, we use attention to mask out information that is not important at the current stage.



Here is the conceptual visualization of applying attention in NLP. If we use the 16th word (Bank) as the query, it will pay more attention to the phrase “Bank of America” rather than the “river bank”.



Encoder

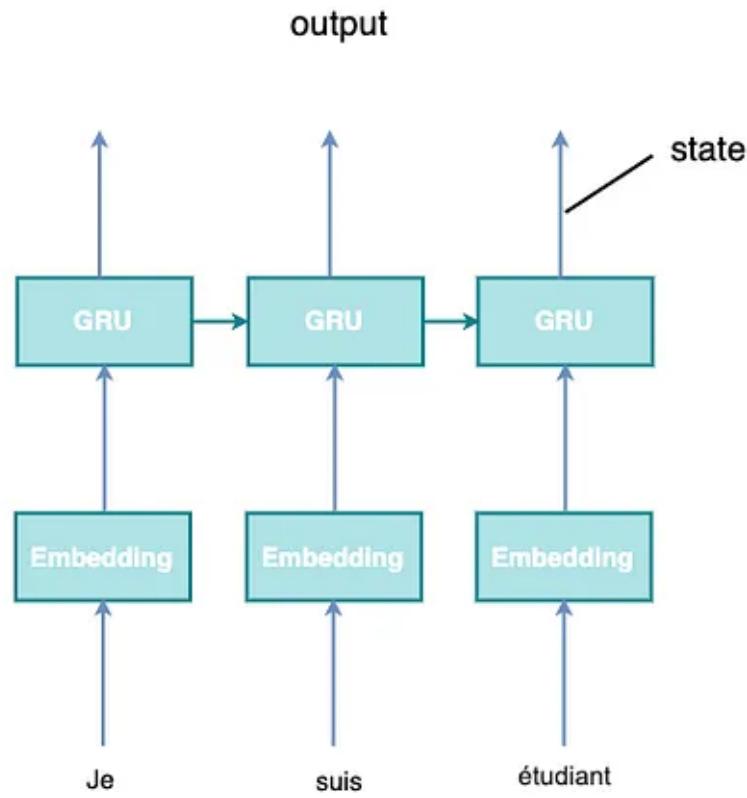
Here is the encoder. It contains an embedding layer and a GRU.

```

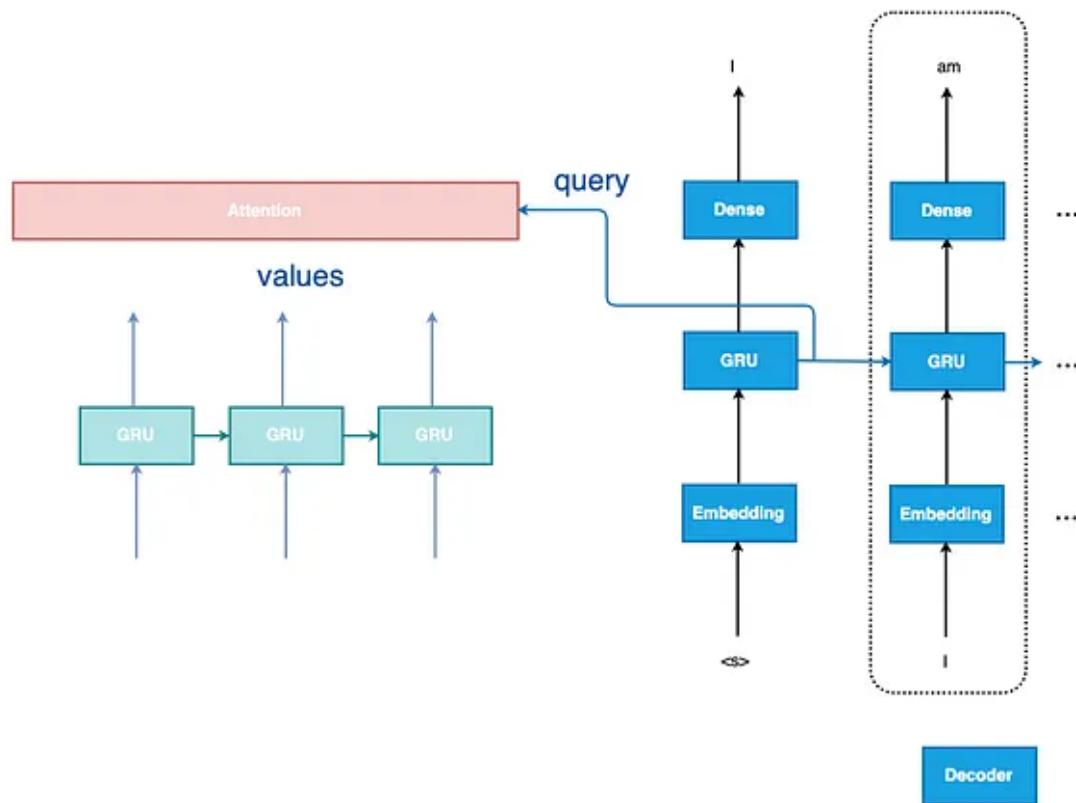
117     class Encoder(tf.keras.Model):
118         def __init__(self, vocab_size, embedding_dim, enc_units, batch_sz):
119             super(Encoder, self).__init__()
120             self.batch_sz = batch_sz
121             self.enc_units = enc_units
122             self.embedding = tf.keras.layers.Embedding(vocab_size, embedding_dim)
123             self.gru = tf.keras.layers.GRU(self.enc_units,
124                                         return_sequences=True,
125                                         return_state=True,
126                                         recurrent_initializer='glorot_uniform')
127
128         def call(self, x, hidden):
129             x = self.embedding(x)
130             output, state = self.gru(x, initial_state=hidden)
131             return output, state
132
133         def initialize_hidden_state(self):
134             return tf.zeros((self.batch_sz, self.enc_units))
135
136
137     encoder = Encoder(vocab_inp_size, embedding_dim, units, BATCH_SIZE)

```

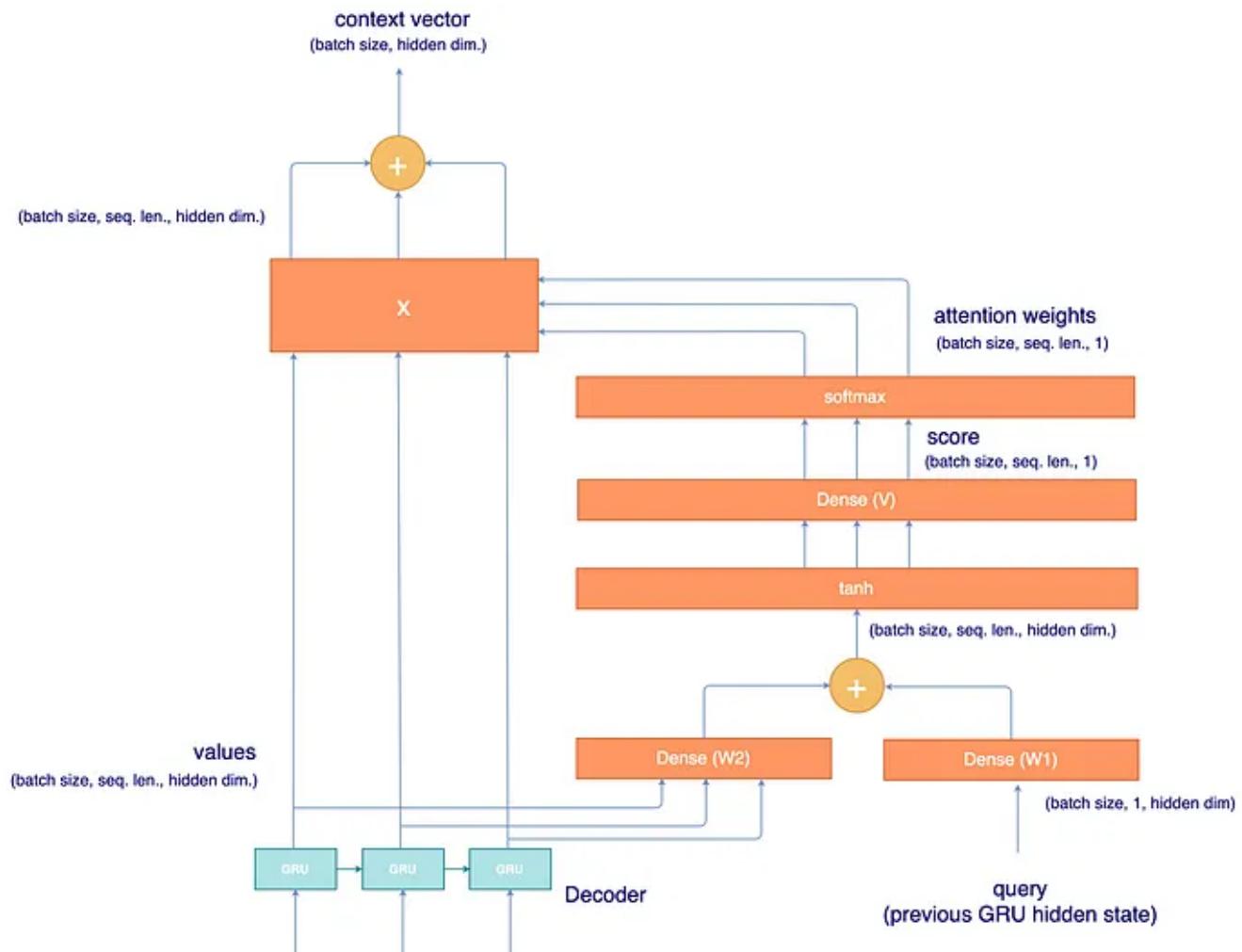
Encoder produces “*output*”, the hidden states of the complete sequence, and “*state*”, the last hidden state.



Let's elaborate further on how attention works. The query is the previous hidden GRU state of the decoder and values are the output of the encoder. The attention keeps the input context that is relevant to the query.



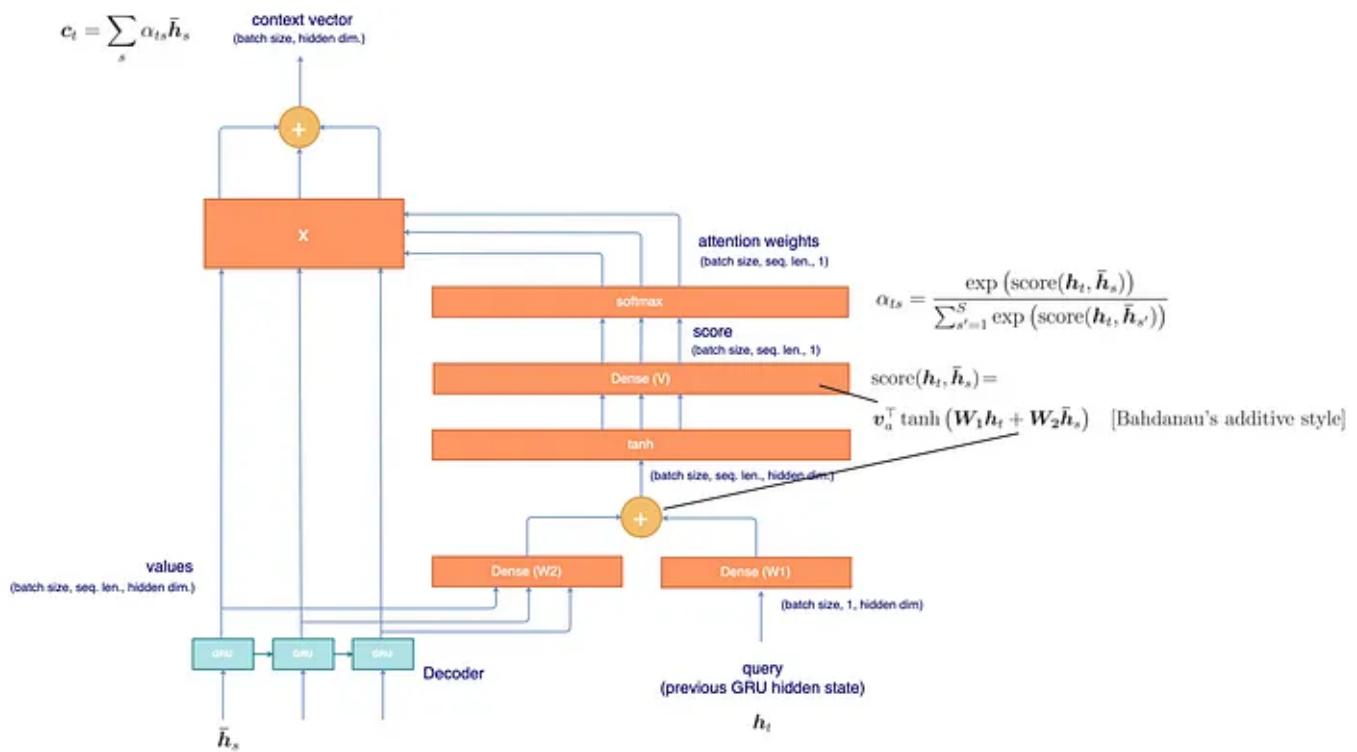
This is the detailed model diagram for the attention module.



“values” in our example contains a sequence of 3 integer vectors representing the input “Je suis étudiant”. Each vector has “hidden dim.” components. We train dense layers (W_1 and W_2) to suppress the vectors that are not relevant to the query. Its output passes through a *tanh* function and further trained to score each vector with a dense layer V . Then, the scores pass through a softmax function. Each vector in the sequence is multiplied with the corresponding softmax result. The model will be trained to zero out irrelevant vectors in the input sequence according to the current query. Finally, we sum up all vectors. The attention output, the context vector, will have the shape (batch size, hidden dimension). Here is the model code for the attention module.

```
140     class BahdanauAttention(tf.keras.layers.Layer):
141         def __init__(self, units):
142             super(BahdanauAttention, self).__init__()
143             self.W1 = tf.keras.layers.Dense(units)
144             self.W2 = tf.keras.layers.Dense(units)
145             self.V = tf.keras.layers.Dense(1)
146
147         def call(self, query, values):
148             # query hidden state shape == (batch_size, hidden size)
149             # query_with_time_axis shape == (batch_size, 1, hidden size)
150             # values shape == (batch_size, max_len, hidden size)
151             # we are doing this to broadcast addition along the time axis to calculate the score
152             query_with_time_axis = tf.expand_dims(query, 1)
153
154             # score shape == (batch_size, max_length, 1)
155             # we get 1 at the last axis because we are applying score to self.V
156             # the shape of the tensor before applying self.V is (batch_size, max_length, units)
157             score = self.V(tf.nn.tanh(
158                 self.W1(query_with_time_axis) + self.W2(values)))
159
160             # attention_weights shape == (batch_size, max_length, 1)
161             attention_weights = tf.nn.softmax(score, axis=1)
162
163             # context_vector shape after sum == (batch_size, hidden_size)
164             context_vector = attention_weights * values
165             context_vector = tf.reduce_sum(context_vector, axis=1)
166
167             return context_vector, attention_weights
168
169
170     attention_layer = BahdanauAttention(10)
```

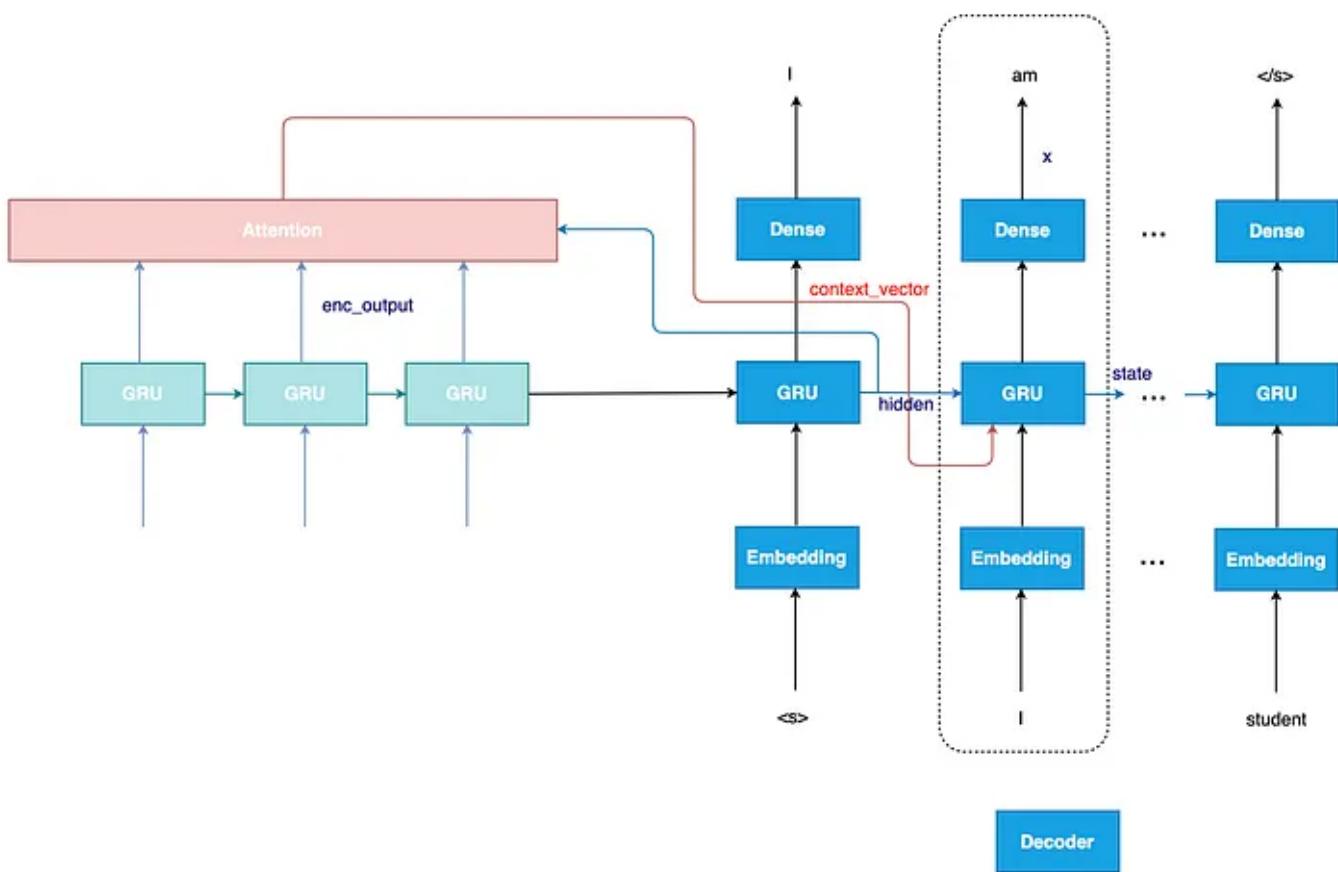
And the equations that we apply.



Finally, this is the decoder.

```
173     class Decoder(tf.keras.Model):
174         def __init__(self, vocab_size, embedding_dim, dec_units, batch_sz):
175             super(Decoder, self).__init__()
176             self.batch_sz = batch_sz
177             self.dec_units = dec_units
178             self.embedding = tf.keras.layers.Embedding(vocab_size, embedding_dim)
179             self.gru = tf.keras.layers.GRU(self.dec_units,
180                                         return_sequences=True,
181                                         return_state=True,
182                                         recurrent_initializer='glorot_uniform')
183             self.fc = tf.keras.layers.Dense(vocab_size)
184
185             # used for attention
186             self.attention = BahdanauAttention(self.dec_units)
187
188         def call(self, x, hidden, enc_output):
189             # enc_output shape == (batch_size, max_length, hidden_size)
190             context_vector, attention_weights = self.attention(hidden, enc_output)
191
192             # x shape after passing through embedding == (batch_size, 1, embedding_dim)
193             x = self.embedding(x)
194
195             # x shape after concatenation == (batch_size, 1, embedding_dim + hidden_size)
196             x = tf.concat([tf.expand_dims(context_vector, 1), x], axis=-1)
197
198             # passing the concatenated vector to the GRU
199             output, state = self.gru(x)
200
201             # output shape == (batch_size * 1, hidden_size)
202             output = tf.reshape(output, (-1, output.shape[2]))
203
204             # output shape == (batch_size, vocab)
205             x = self.fc(output)
206
207             return x, state, attention_weights
208
209
210     decoder = Decoder(vocab_tar_size, embedding_dim, units, BATCH_SIZE)
```

At the second time step in the decoder, here are the notation used in the code above.



Here are the Adam optimizer and the loss function defined.

```

212     optimizer = tf.keras.optimizers.Adam()
213     loss_object = tf.keras.losses.SparseCategoricalCrossentropy(
214         from_logits=True, reduction='none')
215
216
217     def loss_function(real, pred):
218         mask = tf.math.logical_not(tf.math.equal(real, 0))
219         loss_ = loss_object(real, pred)
220
221         mask = tf.cast(mask, dtype=loss_.dtype)
222         loss_ *= mask
223
224     return tf.reduce_mean(loss_)

```

Here is the training step:

```
227 @tf.function
228 def train_step(inp, targ, enc_hidden):
229     loss = 0
230
231     with tf.GradientTape() as tape:
232         # enc_output: the hidden stats for the whole sequence
233         # shape: (batch size, seq. len, hidden dim.)
234         # enc_hidden: the last sequence in enc_output
235         # shape: (batch size, hidden dim.)
236         enc_output, enc_hidden = encoder(inp, enc_hidden)
237
238         # Use the last encoder hidden state to initialize decoder
239         dec_hidden = enc_hidden
240
241         # Feed the start token as the first input sequence to the decoder
242         dec_input = tf.expand_dims([targ_lang.word_index['<start>']] * BATCH_SIZE, 1)
243
244         # Teacher forcing - feeding the target as the next input
245         for t in range(1, targ.shape[1]):
246             # passing enc_output to the decoder
247             predictions, dec_hidden, _ = decoder(dec_input, dec_hidden, enc_output)
248
249             loss += loss_function(targ[:, t], predictions)
250
251             # using teacher forcing
252             # use the target label as the next input
253             dec_input = tf.expand_dims(targ[:, t], 1)
254
255         batch_loss = (loss / int(targ.shape[1]))
256
257         variables = encoder.trainable_variables + decoder.trainable_variables
258
259         gradients = tape.gradient(loss, variables)
260
261         optimizer.apply_gradients(zip(gradients, variables))
262
263     return batch_loss
```

And the training itself:

```
266 EPOCHS = 10
267
268 for epoch in range(EPOCHS):
269     start = time.time()
270
271     enc_hidden = encoder.initialize_hidden_state()
272     total_loss = 0
273
274     for (batch, (inp, targ)) in enumerate(dataset.take(steps_per_epoch)):
275         batch_loss = train_step(inp, targ, enc_hidden)
276         total_loss += batch_loss
277
278         if batch % 100 == 0:
279             print('Epoch {} Batch {} Loss {:.4f}'.format(epoch + 1,
280                                              batch,
281                                              batch_loss.numpy()))
282
283     print('Epoch {} Loss {:.4f}'.format(epoch + 1,
284                                         total_loss / steps_per_epoch))
285     print('Time taken for 1 epoch {} sec\n'.format(time.time() - start))
```

After the model is trained, we can evaluate the model. The evaluation step is similar to the training, except it uses the previous prediction as input, rather than the target word. When we predict the <end> token, the output is complete for that sample.

```

288     def evaluate(sentence):
289         attention_plot = np.zeros((max_length_targ, max_length_inp))
290
291         sentence = preprocess_sentence(sentence)
292
293         inputs = [inp_lang.word_index[i] for i in sentence.split(' ')]
294         inputs = tf.keras.preprocessing.sequence.pad_sequences([inputs],
295                                         maxlen=max_length_inp,
296                                         padding='post')
297         inputs = tf.convert_to_tensor(inputs)
298
299         result = ''
300
301         hidden = [tf.zeros((1, units))]
302         enc_out, enc_hidden = encoder(inputs, hidden)
303
304         dec_hidden = enc_hidden
305         dec_input = tf.expand_dims([targ_lang.word_index['<start>']], 0)
306
307         for t in range(max_length_targ):
308             predictions, dec_hidden, attention_weights = decoder(dec_input,
309                                         dec_hidden,
310                                         enc_out)
311
312             # storing the attention weights to plot later on
313             attention_weights = tf.reshape(attention_weights, (-1,))
314             attention_plot[t] = attention_weights.numpy()
315
316             predicted_id = tf.argmax(predictions[0]).numpy()
317
318             result += targ_lang.index_word[predicted_id] + ' '
319
320             if targ_lang.index_word[predicted_id] == '<end>':
321                 return result, sentence, attention_plot
322
323             # the predicted ID is fed back into the model
324             dec_input = tf.expand_dims([predicted_id], 0)
325
326     return result, sentence, attention_plot

```

Finally, we plot the weights applied in the attention module when we translate “hace mucho frio aqui”.

```
329     # function for plotting the attention weights
330     def plot_attention(attention, sentence, predicted_sentence):
331         fig = plt.figure(figsize=(10, 10))
332         ax = fig.add_subplot(1, 1, 1)
333         ax.matshow(attention, cmap='viridis')
334
335         fontdict = {'fontsize': 14}
336
337         ax.set_xticklabels([''] + sentence, fontdict=fontdict, rotation=90)
338         ax.set_yticklabels([''] + predicted_sentence, fontdict=fontdict)
339
340         ax.xaxis.set_major_locator(ticker.MultipleLocator(1))
341         ax.yaxis.set_major_locator(ticker.MultipleLocator(1))
342
343         plt.show()
344
345
346     def translate(sentence):
347         result, sentence, attention_plot = evaluate(sentence)
348
349         print('Input: %s' % (sentence))
350         print('Predicted translation: {}'.format(result))
351
352         attention_plot = attention_plot[:len(result.split(' ')), :len(sentence.split(' '))]
353         plot_attention(attention_plot, sentence.split(' '), result.split(' '))
354         plt.show()
355
356
357     translate(u'hace mucho frio aqui.')
```

The plot indicates the context where the model focuses on at different stages. At the time when it should predict the word “cold”, the input context at that stage focuses on the word “frio” — cold in Spanish.

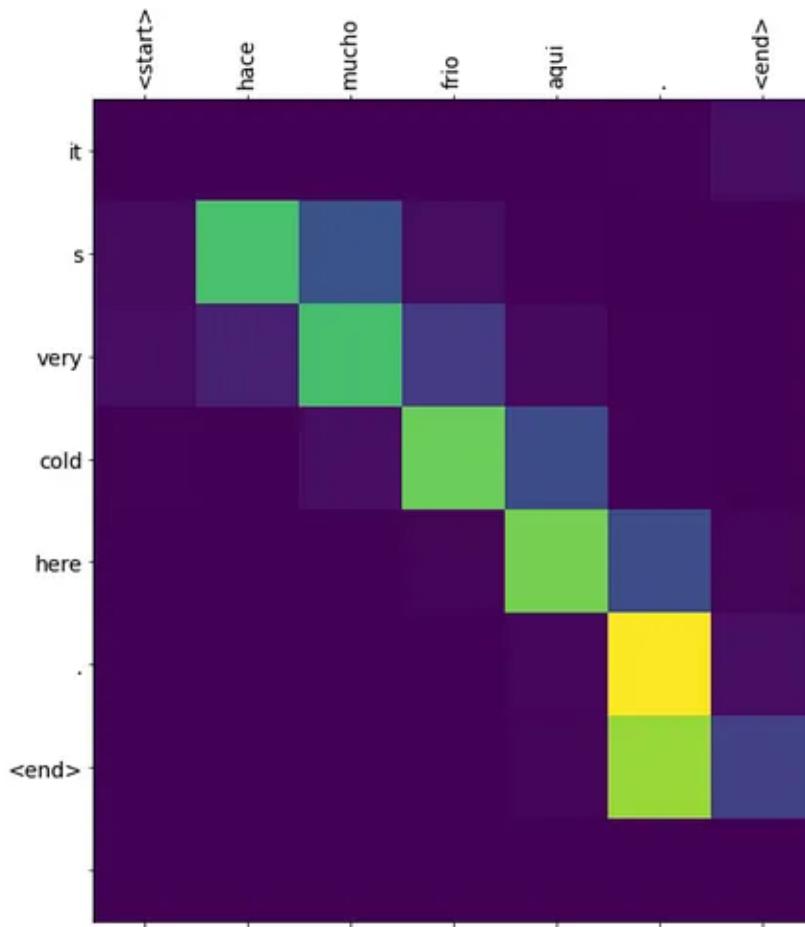
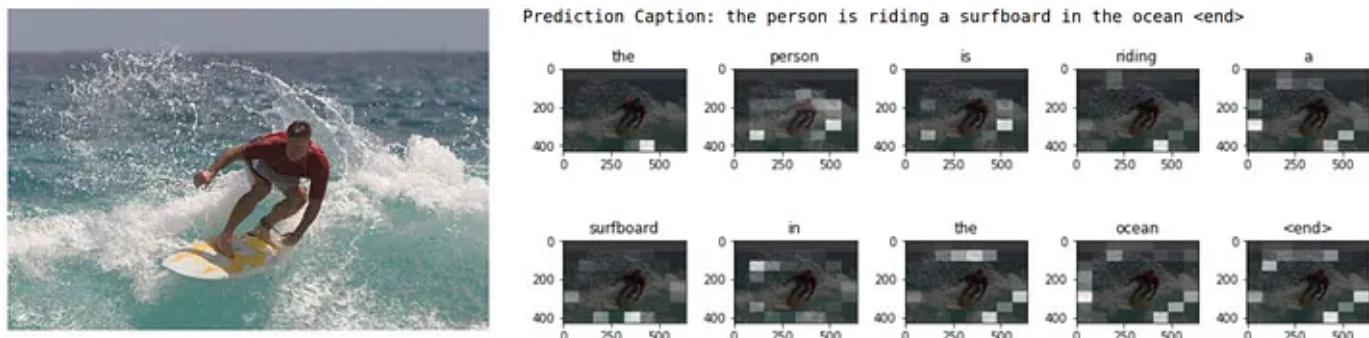


Image captioning with visual attention and Inception V3 for feature extraction

In this example, we generate image caption directly from an image.



Again, we will apply the attention mechanism to pay attention to smaller areas of images in generating the next word. The attention and the decoder module will be similar to the translation example. The major difference is on preprocessing the data. We will use a Inception v3 model to extract image features and save them to files. The

input to the encoder will be these extracted image features. We will go through the code quickly since most parts are similar to the previous translation example.

First, we download the annotation files containing the captions and about 82,000 MSCOCO images (13GB).

```
\\"sky."}, {"image_id": 71691, "id": 137454, "caption": "A woman sitting on a bench outside of a building."}, {"image_id": 182349, "id": 137465, "caption": "A traffic light in front of a fence next to a forest."}, {"image_id": 85657, "id": 137478, "caption": "A herd of \\"
```

```

1  # Modified from
2  # https://www.tensorflow.org/tutorials/text/image_captioning
3  import tensorflow as tf
4
5  # You'll generate plots of attention in order to see which parts of an image
6  # our model focuses on during captioning
7  import matplotlib.pyplot as plt
8
9  import collections
10 import random
11 import numpy as np
12 import os
13 import time
14 import json
15 from PIL import Image
16
17 # Download caption annotation files
18 annotation_folder = '/annotations/'
19 if not os.path.exists(os.path.abspath('..') + annotation_folder):
20     annotation_zip = tf.keras.utils.get_file('captions.zip',
21                                              cache_subdir=os.path.abspath('..'),
22                                              origin='http://images.cocodataset.org/annotations/ann',
23                                              extract=True)
24     annotation_file = os.path.dirname(annotation_zip) + '/annotations/captions_train2014.json'
25     os.remove(annotation_zip)
26 else:
27     annotation_file = os.path.abspath('..') + '/annotations/captions_train2014.json'
28
29 # Download image files
30 image_folder = '/train2014/'
31 if not os.path.exists(os.path.abspath('..') + image_folder):
32     image_zip = tf.keras.utils.get_file('train2014.zip',
33                                         cache_subdir=os.path.abspath('..'),
34                                         origin='http://images.cocodataset.org/zips/train2014.zip',
35                                         extract=True)
36     PATH = os.path.dirname(image_zip) + image_folder
37     os.remove(image_zip)
38 else:
39     PATH = os.path.abspath('..') + image_folder
40

```

But for faster training, we take 6,000 images only. And prepare a list of about 30,011 captions and a list of 30,011 corresponding file names. Each image can have multiple captions.

```
41 # Optional: limit the size of the training set
42 with open(annotation_file, 'r') as f:
43     annotations = json.load(f)
44
45 # Group all captions together having the same image ID.
46 image_path_to_caption = collections.defaultdict(list)
47 for val in annotations['annotations']:
48     caption = f"<start> {val['caption']} <end>"
49     image_path = PATH + 'COCO_train2014_' + '%012d.jpg' % (val['image_id'])
50     image_path_to_caption[image_path].append(caption)
51
52 image_paths = list(image_path_to_caption.keys())
53 random.shuffle(image_paths)
54
55 # Select the first 6000 image_paths from the shuffled set.
56 # Approximately each image id has 5 captions associated with it, so that will
57 # lead to 30,000 examples.
58 train_image_paths = image_paths[:6000]
59
60 # Train_captions hold a list of captions.
61 # img_name_vector holds a list of the corresponding file.
62 # An image can have multiple captions.
63 train_captions = []
64 img_name_vector = []
65
66 for image_path in train_image_paths:
67     caption_list = image_path_to_caption[image_path]
68     train_captions.extend(caption_list)
69     img_name_vector.extend([image_path] * len(caption_list))
```

We will resize the image and use a pre-built Inception v3 to extract features. Then the features are stored into a file.

```
72     def load_image(image_path):
73         img = tf.io.read_file(image_path)
74         img = tf.image.decode_jpeg(img, channels=3)
75         img = tf.image.resize(img, (299, 299))
76         img = tf.keras.applications.inception_v3.preprocess_input(img)
77         return img, image_path
78
79
80     image_model = tf.keras.applications.InceptionV3(include_top=False,
81                                                 weights='imagenet')
82     new_input = image_model.input
83     hidden_layer = image_model.layers[-1].output
84
85     image_features_extract_model = tf.keras.Model(new_input, hidden_layer)
86
87     # Get unique images
88     encode_train = sorted(set(img_name_vector))
89
90     # Feel free to change batch_size according to your system configuration
91     image_dataset = tf.data.Dataset.from_tensor_slices(encode_train)
92     image_dataset = image_dataset.map(
93         load_image, num_parallel_calls=tf.data.experimental.AUTOTUNE).batch(16)
94
95     for img, path in image_dataset:
96         batch_features = image_features_extract_model(img)
97         # flatten: (16, 8, 8, 2048) -> (16, 64, 2048)
98         batch_features = tf.reshape(batch_features,
99                                     (batch_features.shape[0], -1, batch_features.shape[3]))
100
101         # Save the features into a file
102         for bf, p in zip(batch_features, path):
103             path_of_feature = p.numpy().decode("utf-8")
104             np.save(path_of_feature, bf.numpy())
```

Next, we preprocess and tokenize the captions. We will prepare a vocabulary for the top 5,000 words and create the mapping between a word and the word integer index. We also pad all sequences to the same length as the longest one.

```
106 # Preprocess and tokenize the captions #
107
108
109 # Find the maximum length of any caption in our dataset
110 def calc_max_length(tensor):
111     return max(len(t) for t in tensor)
112
113
114 # Choose the top 5000 words from the vocabulary
115 top_k = 5000
116 tokenizer = tf.keras.preprocessing.text.Tokenizer(num_words=top_k,
117                                                 oov_token=<unk>,
118                                                 filters='!"#$%&()/*-./,:=?@[\]^`{|}~')
119 tokenizer.fit_on_texts(trainCaptions)
120 train_seqs = tokenizer.texts_to_sequences(trainCaptions)
121
122 tokenizer.word_index['<pad>'] = 0
123 tokenizer.index_word[0] = '<pad>'
124
125 # Create the tokenized vectors
126 train_seqs = tokenizer.texts_to_sequences(trainCaptions)
127
128 # Pad each vector to the max_length of the captions
129 # If you do not provide a max_length value, pad_sequences calculates it automatically
130 cap_vector = tf.keras.preprocessing.sequence.pad_sequences(train_seqs, padding='post')
131
132 # Calculates the max_length, which is used to store the attention weights
133 max_length = calc_max_length(train_seqs)
```

Next, we will split the samples between training and validation.

```
135 img_to_cap_vector = collections.defaultdict(list)
136 for img, cap in zip(img_name_vector, cap_vector):
137     img_to_cap_vector[img].append(cap)
138
139 # Create training and validation sets using an 80-20 split randomly.
140 img_keys = list(img_to_cap_vector.keys())
141 random.shuffle(img_keys)
142
143 slice_index = int(len(img_keys) * 0.8)
144 img_name_train_keys, img_name_val_keys = img_keys[:slice_index], img_keys[slice_index:]
145
146 img_name_train = []
147 cap_train = []
148 for imgt in img_name_train_keys:
149     capt_len = len(img_to_cap_vector[imgt])
150     img_name_train.extend([imgt] * capt_len)
151     cap_train.extend(img_to_cap_vector[imgt])
152
153 img_name_val = []
154 cap_val = []
155 for imgv in img_name_val_keys:
156     capv_len = len(img_to_cap_vector[imgv])
157     img_name_val.extend([imgv] * capv_len)
158     cap_val.extend(img_to_cap_vector[imgv])
```

And we will create the datasets.

```
161 # Feel free to change these parameters according to your system's configuration
162
163 BATCH_SIZE = 64
164 BUFFER_SIZE = 1000
165 embedding_dim = 256
166 units = 512
167 vocab_size = top_k + 1
168 num_steps = len(img_name_train) // BATCH_SIZE
169 # Shape of the vector extracted from InceptionV3 is (64, 2048)
170 # These two variables represent that vector shape
171 features_shape = 2048
172 attention_features_shape = 64
173
174
175 # Load the numpy files
176 def map_func(img_name, cap):
177     img_tensor = np.load(img_name.decode('utf-8') + '.npy')
178     return img_tensor, cap
179
180
181 dataset = tf.data.Dataset.from_tensor_slices((img_name_train, cap_train))
182
183 # Use map to load the numpy files in parallel
184 dataset = dataset.map(lambda item1, item2: tf.numpy_function(
185     map_func, [item1, item2], [tf.float32, tf.int32]),
186     num_parallel_calls=tf.data.experimental.AUTOTUNE)
187
188 # Shuffle and batch
189 dataset = dataset.shuffle(BUFFER_SIZE).batch(BATCH_SIZE)
190 dataset = dataset.prefetch(buffer_size=tf.data.experimental.AUTOTUNE)
191
```

The attention model will be the same as what we already discussed.

```
193     class BahdanauAttention(tf.keras.Model):
194         def __init__(self, units):
195             super(BahdanauAttention, self).__init__()
196             self.W1 = tf.keras.layers.Dense(units)
197             self.W2 = tf.keras.layers.Dense(units)
198             self.V = tf.keras.layers.Dense(1)
199
200         def call(self, features, hidden):
201             # features(CNN_encoder output) shape == (batch_size, 64, embedding_dim)
202
203             # hidden shape == (batch_size, hidden_size)
204             # hidden_with_time_axis shape == (batch_size, 1, hidden_size)
205             hidden_with_time_axis = tf.expand_dims(hidden, 1)
206
207             # attention_hidden_layer shape == (batch_size, 64, units)
208             attention_hidden_layer = (tf.nn.tanh(self.W1(features) +
209                                         self.W2(hidden_with_time_axis)))
210
211             # score shape == (batch_size, 64, 1)
212             # This gives you an unnormalized score for each image feature.
213             score = self.V(attention_hidden_layer)
214
215             # attention_weights shape == (batch_size, 64, 1)
216             attention_weights = tf.nn.softmax(score, axis=1)
217
218             # context_vector shape after sum == (batch_size, hidden_size)
219             context_vector = attention_weights * features
220             context_vector = tf.reduce_sum(context_vector, axis=1)
221
222             return context_vector, attention_weights
```

Because our inputs are extracted features already, the encoder is just a simple dense layer followed by ReLU.

```
225     class CNN_Encoder(tf.keras.Model):
226         # Since you have already extracted the features and dumped it using pickle
227         # This encoder passes those features through a Fully connected layer
228         def __init__(self, embedding_dim):
229             super(CNN_Encoder, self).__init__()
230             # shape after fc == (batch_size, 64, embedding_dim)
231             self.fc = tf.keras.layers.Dense(embedding_dim)
232
233     @tf.function
234     def call(self, x):
235         x = self.fc(x)
236         x = tf.nn.relu(x)
237         return x
```

The decoder is similar to the language translation example, with one more dense layer and a couple ReLU layers. It uses an attention layer also.

```
239     class RNN_Decoder(tf.keras.Model):
240         def __init__(self, embedding_dim, units, vocab_size):
241             super(RNN_Decoder, self).__init__()
242             self.units = units
243
244             self.embedding = tf.keras.layers.Embedding(vocab_size, embedding_dim)
245             self.gru = tf.keras.layers.GRU(self.units,
246                                         return_sequences=True,
247                                         return_state=True,
248                                         recurrent_initializer='glorot_uniform')
249             self.fc1 = tf.keras.layers.Dense(self.units)
250             self.fc2 = tf.keras.layers.Dense(vocab_size)
251
252             self.attention = BahdanauAttention(self.units)
253
254     def call(self, x, features, hidden):
255         # defining attention as a separate model
256         context_vector, attention_weights = self.attention(features, hidden)
257
258         # x shape after passing through embedding == (batch_size, 1, embedding_dim)
259         x = self.embedding(x)
260
261         # x shape after concatenation == (batch_size, 1, embedding_dim + hidden_size)
262         x = tf.concat([tf.expand_dims(context_vector, 1), x], axis=-1)
263
264         # passing the concatenated vector to the GRU
265         output, state = self.gru(x)
266
267         # shape == (batch_size, max_length, hidden_size)
268         x = self.fc1(output)
269
270         # x shape == (batch_size * max_length, hidden_size)
271         x = tf.reshape(x, (-1, x.shape[2]))
272
273         # output shape == (batch_size * max_length, vocab)
274         x = self.fc2(x)
275
276         return x, state, attention_weights
277
278     def reset_state(self, batch_size):
279         return tf.zeros((batch_size, self.units))
```

We define the model, an Adam optimizer, the loss function and the CheckpointManager.

```
282 encoder = CNN_Encoder(embedding_dim)
283 decoder = RNN_Decoder(embedding_dim, units, vocab_size)
284
285 optimizer = tf.keras.optimizers.Adam()
286 loss_object = tf.keras.losses.SparseCategoricalCrossentropy(
287     from_logits=True, reduction='none')
288
289
290 def loss_function(real, pred):
291     mask = tf.math.logical_not(tf.math.equal(real, 0))
292     loss_ = loss_object(real, pred)
293
294     mask = tf.cast(mask, dtype=loss_.dtype)
295     loss_ *= mask
296
297     return tf.reduce_mean(loss_)
298
299
300 checkpoint_path = "./checkpoints/train"
301 ckpt = tf.train.Checkpoint(encoder=encoder,
302                             decoder=decoder,
303                             optimizer=optimizer)
304 ckpt_manager = tf.train.CheckpointManager(ckpt, checkpoint_path, max_to_keep=5)
305
306 start_epoch = 0
307 if ckpt_manager.latest_checkpoint:
308     start_epoch = int(ckpt_manager.latest_checkpoint.split('-')[-1])
309     # restoring the latest checkpoint in checkpoint_path
310     ckpt.restore(ckpt_manager.latest_checkpoint)
311
312 # adding this in a separate cell because if you run the training cell
313 # many times, the loss_plot array will be reset
314 loss_plot = []
```

In a training step, we encode the extracted image features. Then we use the decoder with the attention to predict one word at a time.

```
317     @tf.function
318     def train_step(img_tensor, target):
319         loss = 0
320
321         # initializing the hidden state for each batch
322         # because the captions are not related from image to image
323         hidden = decoder.reset_state(batch_size=target.shape[0])
324
325         dec_input = tf.expand_dims([tokenizer.word_index['<start>']] * target.shape[0], 1)
326
327         with tf.GradientTape() as tape:
328             features = encoder(img_tensor)
329
330             for i in range(1, target.shape[1]):
331                 # passing the features through the decoder
332                 predictions, hidden, _ = decoder(dec_input, features, hidden)
333
334                 loss += loss_function(target[:, i], predictions)
335
336                 # using teacher forcing
337                 dec_input = tf.expand_dims(target[:, i], 1)
338
339             total_loss = (loss / int(target.shape[1]))
340
341             trainable_variables = encoder.trainable_variables + decoder.trainable_variables
342
343             gradients = tape.gradient(loss, trainable_variables)
344
345             optimizer.apply_gradients(zip(gradients, trainable_variables))
346
347         return loss, total_loss
```

Here is the training loop,

```
350 EPOCHS = 20
351
352 for epoch in range(start_epoch, EPOCHS):
353     start = time.time()
354     total_loss = 0
355
356     for (batch, (img_tensor, target)) in enumerate(dataset):
357         batch_loss, t_loss = train_step(img_tensor, target)
358         total_loss += t_loss
359
360         if batch % 100 == 0:
361             print('Epoch {} Batch {} Loss {:.4f}'.format(
362                 epoch + 1, batch, batch_loss.numpy() / int(target.shape[1])))
363             # storing the epoch end loss value to plot later
364             loss_plot.append(total_loss / num_steps)
365
366         if epoch % 5 == 0:
367             ckpt_manager.save()
368
369         print('Epoch {} Loss {:.6f}'.format(epoch + 1,
370                                         total_loss / num_steps))
371     print('Time taken for 1 epoch {} sec\n'.format(time.time() - start))
```

and the evaluation.

```
374     def evaluate(image):
375         attention_plot = np.zeros((max_length, attention_features_shape))
376
377         hidden = decoder.reset_state(batch_size=1)
378
379         temp_input = tf.expand_dims(load_image(image)[0], 0)
380         img_tensor_val = image_features_extract_model(temp_input)
381         img_tensor_val = tf.reshape(img_tensor_val, (img_tensor_val.shape[0], -1, img_tensor_val.shape[3]))
382
383         features = encoder(img_tensor_val)
384
385         dec_input = tf.expand_dims([tokenizer.word_index['<start>']], 0)
386         result = []
387
388         for i in range(max_length):
389             predictions, hidden, attention_weights = decoder(dec_input, features, hidden)
390
391             attention_plot[i] = tf.reshape(attention_weights, (-1,)).numpy()
392
393             predicted_id = tf.random.categorical(predictions, 1)[0][0].numpy()
394             result.append(tokenizer.index_word[predicted_id])
395
396             if tokenizer.index_word[predicted_id] == '<end>':
397                 return result, attention_plot
398
399             dec_input = tf.expand_dims([predicted_id], 0)
400
401         attention_plot = attention_plot[:len(result), :]
402         return result, attention_plot
```

Credits and References

All the source code is originated or modified from the TensorFlow [tutorial](#).

Machine Learning

Artificial Intelligence

Deep Learning

Software Development

Data Science