# Intuitive Understanding of Attention Mechanism in Deep Learning

A TensorFlow Implementation of Neural Machine Translation with Attention

Harshall Lamba  [Follow]

Mar 20 · 11 min read ★

## Caution

This is a slightly advanced tutorial and requires basic understanding of sequence to sequence models using RNNs. Please refer my earlier *blog here* wherein I have explained in detail the concept of Seq2Seq models.

## Table of Contents

## 1. Introduction

Attention is one of the most influential ideas in the Deep Learning community. Even though this mechanism is now used in various problems like image captioning and others,it was initially designed in the context of Neural Machine Translation using Seq2Seq Models. In this blog post I will consider the same problem as the running example to illustrate the concept. We would be using attention to design a system which translates a given English sentence to Marathi, the exact same example I considered in my earlier blog.

So what's wrong with seq2seq models?

The seq2seq models is normally composed of an encoder-decoder architecture, where the encoder processes the input sequence and encodes/compresses/summarizes the information into a context vector (also called as the "thought vector") of a fixed length. This representation is expected to be a good summary of the entire input sequence. The decoder is then initialized with this context vector,

using which it starts generating the transformed output.

*A critical and apparent disadvantage of this fixed-length context vector design is the incapability of the system to remember longer sequences. Often is has forgotten the earlier parts of the sequence once it has processed the entire the sequence. The attention mechanism was born to resolve this problem.*

Let's break this down into finer details. Since I have already explained most of the basic concepts required to understand Attention in my previous *blog*, here I will directly jump into the meat of the issue without any further adieu.
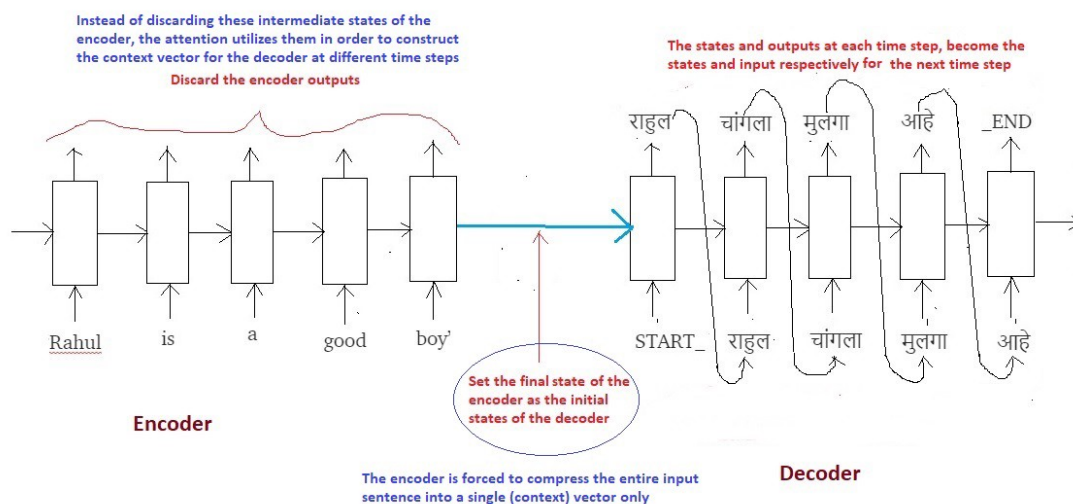
## 2. The central idea behind Attention

For the illustrative purposes, I will borrow the same example that I used to explain Seq2Seq models in my previous *blog*.

Input (English) Sentence: "Rahul is a good boy"

Target (Marathi) Sentence: "राहुल चांगला मुलगा आहे"

The only change will be that instead of an LSTM layer that I used in my previous explanation, here I will use a GRU layer. The reason being that LSTM has two internal states (hidden state and cell state) and GRU has only one internal state (hidden state). This will help simplify the the concept and explanation.

Recall the below diagram in which I summarized the entire process procedure of Seq2Seq modelling.
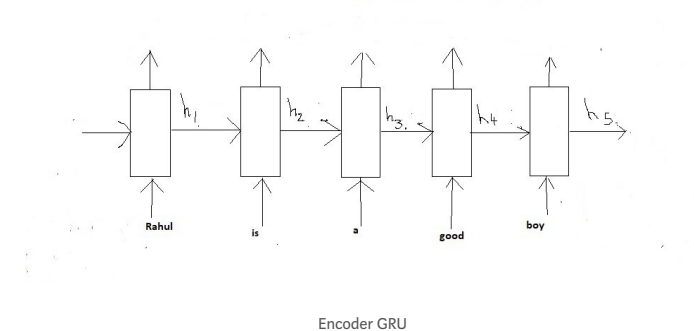


In the traditional Seq2Seq model, we discard all the intermediate states of the encoder and use only its final states (vector) to initialize the decoder. This technique works good for smaller sequences,

however as the length of the sequence increases, a single vector becomes a bottleneck and it gets very difficult to summarize long sequences into a single vector. This observation was made empirically as it was noted that the performance of the system decreases drastically as the size of the sequence increases.

The central idea behind Attention is not to throw away those intermediate encoder states but to utilize all the states in order to construct the context vectors required by the decoder to generate the output sequence.

## 3. Why the name Attention?

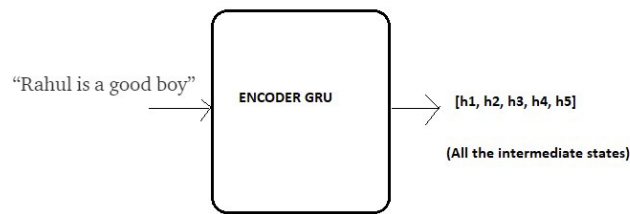Let's name each of the intermediate states of the encoder as below:



Encoder GRU

Notice that since we are using a GRU instead of an LSTM, we only have a single state at each time step and not two states, which thus helps to simplify the illustration. Also note that attention is useful specially in case of longer sequences but for the sake of simplicity we will consider the same above example for illustration.

Recall that these states (h1 to h5) are nothing but vectors of fixed length. To develop some intuition think of these states as vectors which store local information within the sequence. For example;

h1 stores the information present in the start of the sequence (words like 'Rahul' and 'is') while h5 stores the information present in the later part of the sequence (words like 'good' and 'boy').

Lets represent our Encoder GRU with the below simplified diagram:

Compact Representation of Encoder GRU

Now the idea is to utilize all of these local information collectively in order to decide the next sequence while decoding the target sentence.

Imagine you are translating "Rahul is a good boy" to "राहुल चांगला मुलगा आहे". Ask yourself, how do you do it in your mind?

When you predict "राहुल", its obvious that this name is the result of the word "Rahul" present in the input English sentence regardless of the rest of the sentence. We say that while predicting "राहुल", *we pay more attention* to the word "Rahul" in the input sentence.

Similarly while predicting the word "चांगला", we pay more attention to the word "good" in the input sentence.

Similarly while predicting the word "मुलगा", we pay more attention to the word "boy" in the input sentence. And so on..

Hence the name *"ATTENTION".*

As human beings we are quickly able to understand these mappings between different parts of the input sequence and corresponding parts of the output sequence. However its not that straight forward for artificial neural network to automatically detect these mappings.

Thus the Attention mechanism is developed to *"learn"* these mappings through Gradient Descent and Back-propagation.

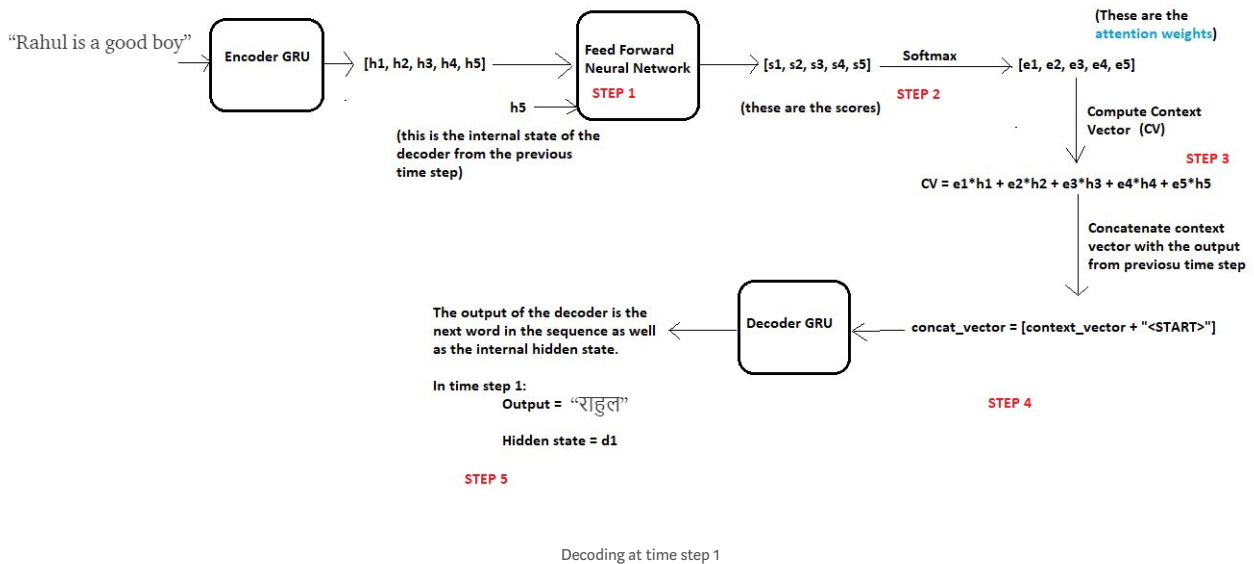## 4. How does Attention work?

Let's get technical and dive into the nitty gritty of Attention mechanism.

### Decoding at time step 1

Continuing the above example, let's say we now want our decoder to start predicting the first word of the target sequence i.e. "राहुल"

At time step 1, we can break the entire process into **five steps** as

below:



Decoding at time step 1

Before we start decoding, we first need to encode the input sequence into a set of internal states (in our case h1, h2, h3, h4 and h5).

Now the hypothesis is that, the next word in the output sequence is dependent on the current state of the decoder (decoder is also a GRU) as well as on the hidden states of the encoder. Thus at each time step, we consider these two things and follow the below steps:

**Step 1 — Compute a score each encoder state**

Since we are predicting the first word itself, the decoder does not have any current internal state. For this reason, we will consider the last state of the encoder (i.e. h5) as the previous decoder state.

Now using these two components (all the encoder states and the current state of the decoder), we will train a simple feed forward neural network.

Why?

Recall we are trying to predict the first word in the target sequence i.e. "राहुल". As per the idea behind attention, we do not need all the encoder states to predict this word, but we need those encoder states which store information about the word "Rahul" in the input sequence.

As discussed previously these intermediate encoder states store the local information of the input sequence. So it is highly likely that the information of the word "Rahul" will be present in the states, let's say, h1 and h2.

Thus we want our decoder to pay more attention to the states h1 and h2 while paying less attention to the remaining states of the encoder.

For this reason we train a feed forward neural network which will **learn** to identify relevant encoder states by generating a high score for the states for which attention is to be paid while low score for the states which are to be ignored.

Let s1, s2, s3, s4 and s5 be the scores generated for the states h1, h2, h3, h4 and h5 correspondingly. Since we assumed that we need to pay more attention to the states h1 and h2 and ignore h3, h4 and h5 in order to predict "राहुल", we expect the above neural to generate scores such that s1 and s2 are high while s3, s4 and s5 are relatively low.

**Step 2— Compute the attention weights**

Once these scores are generated, we apply a softmax on these scores to produce the attention weights e1, e2, e3 ,e4 and e5 as shown above. The advantage of applying softmax is as below:

a) All the weights lie between 0 and 1, i.e., $0 \leq e1, e2, e3, e4, e5 \leq 1$

b) All the weights sum to 1, i.e., e1+e2+3+e4+e5 = 1

Thus we get a nice probabilistic interpretation of the attention weights.

In our case we would expect values like below: (just for intuition)

e1 = 0.75, e2 = 0.2, e3 = 0.02, e4 = 0.02, e5 = 0.01

This means that while predicting the word "राहुल", the decoder needs to put more attention on the states h1 and h2 (since values of e1 and e2 are high) while ignoring the states h3, h4 and h5 (since the values of e3, e4 and e5 are very small).

**Step 3— Compute the context vector**

Once we have computed the attention weights, we need to compute the context vector (thought vector) which will be used by the decoder in order to predict the next word in the sequence. Calculated as follows:

context_vector = e1 * h1 + e2 * h2 + e3 * h3 + e4 * h4 + e5 * h5

Clearly if the values of e1 and e2 are high and those of e3, e4 and e5 are low then the context vector will contain more information from the states h1 and h2 and relatively less information from the states h3, h4 and h5.

**Step 4— Concatenate context vector with output of previous time step**

Finally the decoder uses the below two input vectors to generate the next word in the sequence

a) The context vector

b) The output word generated from the previous time step.

We simply concatenate these two vectors and feed the merged vector to the decoder. **Note that for the first time step, since there is no output from the previous time step, we use a special <START> token for this purpose**. This concept is already discussed in detail in my previous *blog*.

**Step 5— Decoder Output**

The decoder then generates the next word in the sequence (in this case, it is expected to generate "राहुल") and along with the output, the decoder will also generate an internal hidden state, and lets call it as "d1".

## Decoding at time step 2

Now in order to generate the next word "चांगला", the decoder will repeat the same procedure which can be summarized in the below diagram:

The changes are highlighted in **green circles**



Decoding at time step 2

## Decoding at time step 3

"Rahul is a good boy" → **Encoder GRU** → [h1, h2, h3, h4, h5] → **Feed Forward Neural Network** → [s1, s2, s3, s4, s5] → Softmax → [e1, e2, e3, e4, e5]
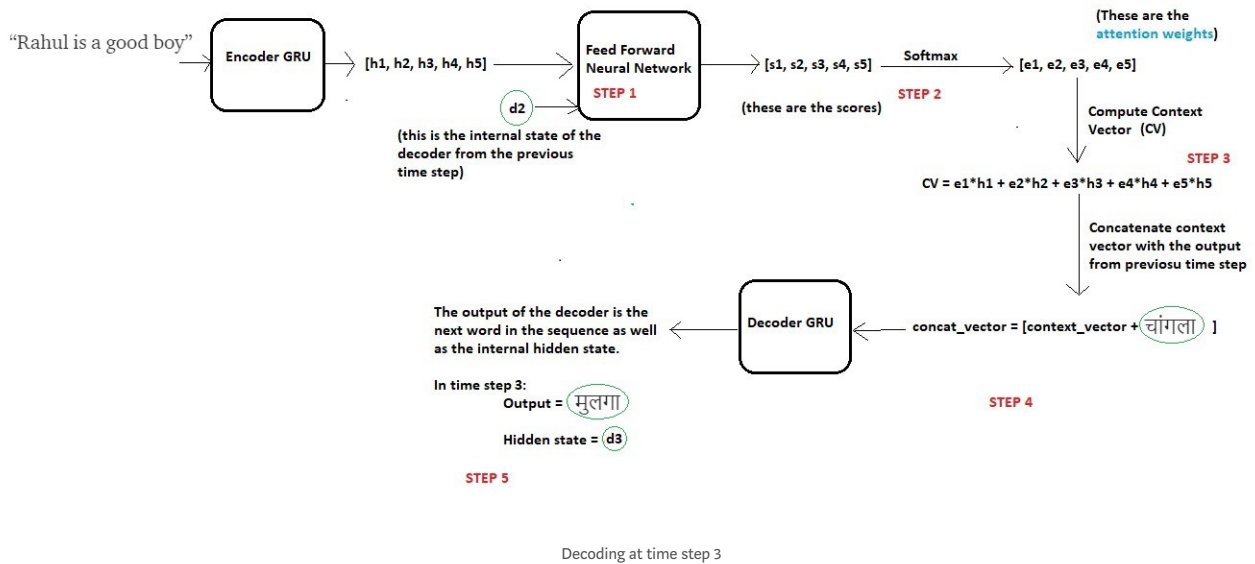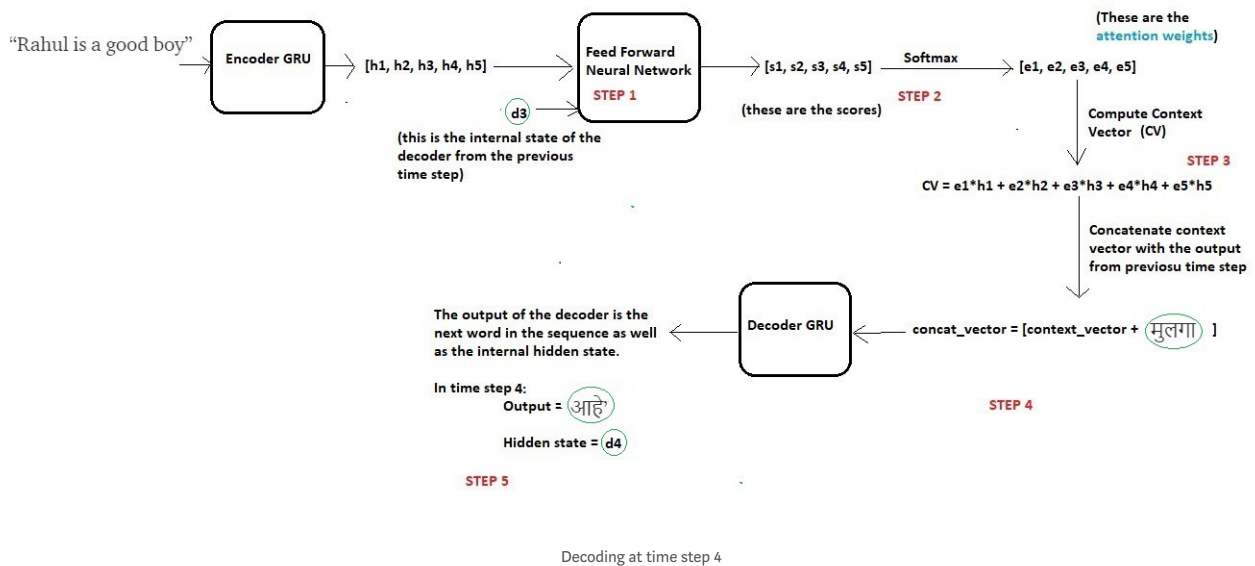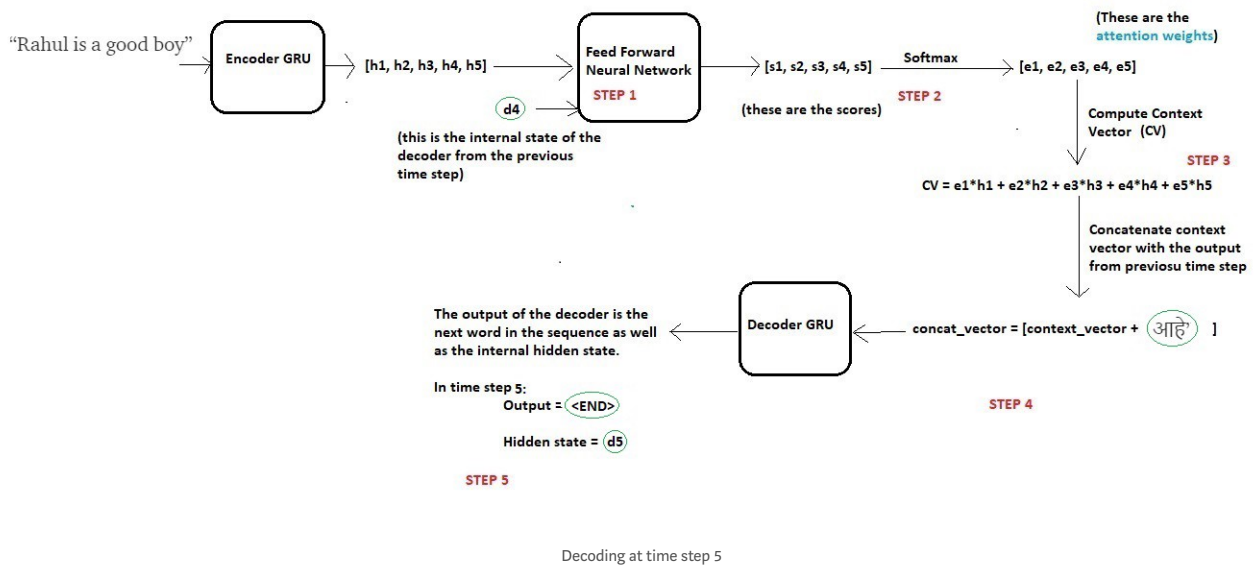
**STEP 1**

d2 →

(this is the internal state of the decoder from the previous time step)

(these are the scores)

**STEP 2**

(These are the attention weights)

Compute Context Vector (CV)

**STEP 3**

CV = e1*h1 + e2*h2 + e3*h3 + e4*h4 + e5*h5

Concatenate context vector with the output from previosu time step

The output of the decoder is the next word in the sequence as well as the internal hidden state.

← **Decoder GRU** ← concat_vector = [context_vector + चांगला ]

**STEP 4**

In time step 3:
　Output = मुलगा
　Hidden state = d3

**STEP 5**

Decoding at time step 3

## Decoding at time step 4

"Rahul is a good boy" → **Encoder GRU** → [h1, h2, h3, h4, h5] → **Feed Forward Neural Network** → [s1, s2, s3, s4, s5] → Softmax → [e1, e2, e3, e4, e5]

**STEP 1**

d3 →

(this is the internal state of the decoder from the previous time step)

(these are the scores)

**STEP 2**

(These are the attention weights)

Compute Context Vector (CV)

**STEP 3**

CV = e1*h1 + e2*h2 + e3*h3 + e4*h4 + e5*h5

Concatenate context vector with the output from previosu time step

The output of the decoder is the next word in the sequence as well as the internal hidden state.

← **Decoder GRU** ← concat_vector = [context_vector + मुलगा ]

**STEP 4**

In time step 4:
　Output = आहे
　Hidden state = d4

**STEP 5**

Decoding at time step 4

## Decoding at time step 5

Decoding at time step 5

- Once the decoder outputs the <END> token, we stop the generation process.

  Note that unlike the fixed context vector used for all the decoder time steps in case of the traditional Seq2Seq models, here in case of Attention, we compute a separate context vector for each time step by computing the attention weights every time.

  Thus using this mechanism our model is able to find interesting mappings between different parts of the input sequence and corresponding parts of the output sequence.

  Note that during the training of the network, we use teacher forcing in order to input the actual word rather than the predicted word from the previous time step. This concept also has been explained in my previous *blog*.

## 5. Code Walk through

As in case of any NLP task, after reading the input file, we perform the basic cleaning and preprocessing as follows:

```
1   # Set the file path
2   file_path = './mar.txt' # this might be different in y
3
4   # read the file
5   lines = open(file_path, encoding='UTF-8').read().strip
6
7   # perform basic cleaning
8   exclude = set(string.punctuation) # Set of all special
9   remove_digits = str.maketrans('', '', string.digits) #
10
11  def preprocess_eng_sentence(sent):
12      '''Function to preprocess English sentence'''
13      sent = sent.lower()
14      sent = re.sub("'", '', sent)
15      sent = ''.join(ch for ch in sent if ch not in excl
16      sent = sent.translate(remove_digits)
17      sent = sent.strip()
18      sent = re.sub(" +", " ", sent)
19      sent = '<start> ' + sent + ' <end>'
20
21  def preprocess_mar_sentence(sent):
22      '''Function to preprocess Marathi sentence'''
23      sent = re.sub("'", '', sent)
24      sent = ''.join(ch for ch in sent if ch not in excl
25      sent = re.sub("[२३०८५४७३९४६]", "", sent)
26      sent = sent.strip()
27      sent = re.sub(" +", " ", sent)
```

Create a class to map every word to an index and vice-versa for any given vocabulary:

```
1   # This class creates a word -> index mapping (e.g,. "d
2   # (e.g., 5 -> "dad") for each language,
3   class LanguageIndex():
4       def __init__(self, lang):
5           self.lang = lang
6           self.word2idx = {}
7           self.idx2word = {}
8           self.vocab = set()
9
10          self.create_index()
11
12      def create_index(self):
13          for phrase in self.lang:
14              self.vocab.update(phrase.split(' '))
15
16          self.vocab = sorted(self.vocab)
17
18          self.word2idx['<pad>'] = 0
```

We use the tf.data input pipeline to create the dataset and then load it later in mini batches. To read more about the input pipeline in TensorFlow, go through the official documentations *here* and *here*.

```
1   def load_dataset(pairs, num_examples):
2       # pairs => already created cleaned input, output p
3
4       # index language using the class defined above
5       inp_lang = LanguageIndex(en for en, ma in pairs)
6       targ_lang = LanguageIndex(ma for en, ma in pairs)
7
8       # Vectorize the input and target languages
9
10      # English sentences
11      input_tensor = [[inp_lang.word2idx[s] for s in en.
12
13      # Marathi sentences
14      target_tensor = [[targ_lang.word2idx[s] for s in m
15
16      # Calculate max_length of input and output tensor
17      # Here, we'll set those to the longest sentence in
18      max_length_inp, max_length_tar = max_length(input_
19
20      # Padding the input and output tensor to the maxim
21      input_tensor = tf.keras.preprocessing.sequence.pad
22
23
24
25      target_tensor = tf.keras.preprocessing.sequence.pa
26
27
28
29      return input_tensor, target_tensor, inp_lang, targ
30
31  # Create the tensors
```

Now using the model sub-classing API of TensorFlow, we define the model as follows. To read more about model sub classing, read the official documentation *here*.

**Note**: Please read the comments in the below section of the code to get better understanding using the concepts we discussed above. Most of the important lines of the code point to the corresponding section of the explanation given above.

```python
def gru(units):
  # If you have a GPU, we recommend using CuDNNGRU(pro
  # the code automatically does that.
    if tf.test.is_gpu_available():
        return tf.keras.layers.CuDNNGRU(units,
                                        return_sequenc
                                        return_state=T
                                        recurrent_init
    else:
        return tf.keras.layers.GRU(units,
                                   return_sequences=Tr
                                   return_state=True,
                                   recurrent_activatio
                                   recurrent_initializ

class Encoder(tf.keras.Model):
    def __init__(self, vocab_size, embedding_dim, enc_
        super(Encoder, self).__init__()
        self.batch_sz = batch_sz
        self.enc_units = enc_units
        self.embedding = tf.keras.layers.Embedding(voc
        self.gru = gru(self.enc_units)

    def call(self, x, hidden):
        x = self.embedding(x)
        output, state = self.gru(x, initial_state = hi
        return output, state

    def initialize_hidden_state(self):
        return tf.zeros((self.batch_sz, self.enc_units

class Decoder(tf.keras.Model):
    def __init__(self, vocab_size, embedding_dim, dec_
        super(Decoder, self).__init__()
        self.batch_sz = batch_sz
        self.dec_units = dec_units
        self.embedding = tf.keras.layers.Embedding(voc
        self.gru = gru(self.dec_units)
        self.fc = tf.keras.layers.Dense(vocab_size)

        # used for attention
        self.W1 = tf.keras.layers.Dense(self.dec_units
        self.W2 = tf.keras.layers.Dense(self.dec_units
        self.V = tf.keras.layers.Dense(1)

    def call(self, x, hidden, enc_output):
        # enc_output shape == (batch_size, max_length,

        # hidden shape == (batch_size, hidden size)
        # hidden_with_time_axis shape == (batch_size,
        # we are doing this to perform addition to cal
        hidden_with_time_axis = tf.expand_dims(hidden,

        # score shape == (batch_size, max_length, 1)
        # we get 1 at the last axis because we are app
        # this is the step 1 described in the blog to
```

Define Optimizer, Loss Function and Checkpoints

```
1   optimizer = tf.train.AdamOptimizer()
2
3   def loss_function(real, pred):
4       mask = 1 - np.equal(real, 0)
5       loss_ = tf.nn.sparse_softmax_cross_entropy_with_lo
6       return tf.reduce_mean(loss_)
7
8   checkpoint_dir = './training_checkpoints'
9   checkpoint_prefix = os.path.join(checkpoint_dir, "ckp
```

Using Eager Execution, we train the network for 10 epochs. To read more about Eager Execution, refer the official documentation *here*.

```
1   EPOCHS = 10
2
3   for epoch in range(EPOCHS):
4       start = time.time()
5
6       hidden = encoder.initialize_hidden_state()
7       total_loss = 0
8
9       for (batch, (inp, targ)) in enumerate(dataset):
10          loss = 0
11
12          with tf.GradientTape() as tape:
13              enc_output, enc_hidden = encoder(inp, hidd
14
15              dec_hidden = enc_hidden
16
17              dec_input = tf.expand_dims([targ_lang.word
18
19              # Teacher forcing - feeding the target as
20              for t in range(1, targ.shape[1]):
21                  # passing enc_output to the decoder
22                  predictions, dec_hidden, _ = decoder(d
23
24                  loss += loss_function(targ[:, t], pred
25
26                  # using teacher forcing
27                  dec_input = tf.expand_dims(targ[:, t],
28
29          batch_loss = (loss / int(targ.shape[1]))
30
31          total_loss += batch_loss
```
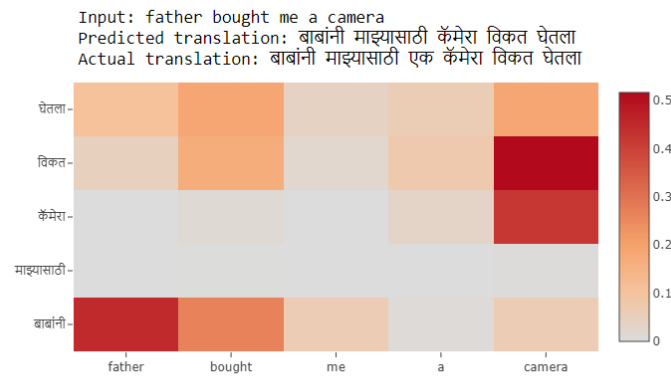
Inference setup and testing:

```python
 1   def evaluate(inputs, encoder, decoder, inp_lang, targ_
 2
 3       attention_plot = np.zeros((max_length_targ, max_le
 4       sentence = ''
 5       for i in inputs[0]:
 6           if i == 0:
 7               break
 8           sentence = sentence + inp_lang.idx2word[i] + '
 9       sentence = sentence[:-1]
10
11       inputs = tf.convert_to_tensor(inputs)
12
13       result = ''
14
15       hidden = [tf.zeros((1, units))]
16       enc_out, enc_hidden = encoder(inputs, hidden)
17
18       dec_hidden = enc_hidden
19       dec_input = tf.expand_dims([targ_lang.word2idx['<s
20
21       # start decoding
22       for t in range(max_length_targ): # limit the lengt
23           predictions, dec_hidden, attention_weights = d
24
25           # storing the attention weights to plot later
26           attention_weights = tf.reshape(attention_weigh
27           attention_plot[t] = attention_weights.numpy()
28
29           predicted_id = tf.argmax(predictions[0]).numpy
30
31           result += targ_lang.idx2word[predicted_id] + '
32
33           # stop decoding if '<end>' is predicted
34           if targ_lang.idx2word[predicted_id] == '<end>'
35               return result, sentence, attention_plot
36
37           # the predicted ID is fed back into the model
38           dec_input = tf.expand_dims([predicted_id], 0)
39
40       return result, sentence, attention_plot
41
42
43   def predict_random_val_sentence():
44       actual_sent = ''
```
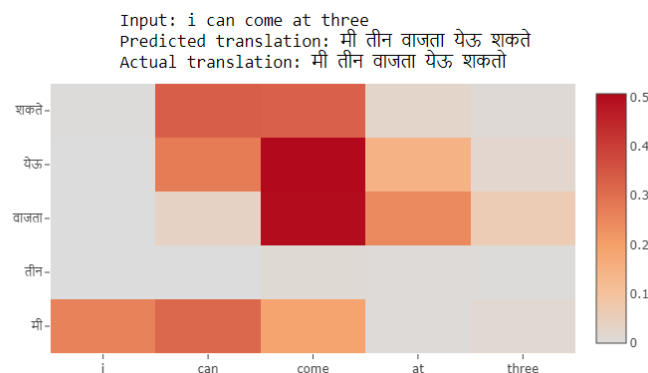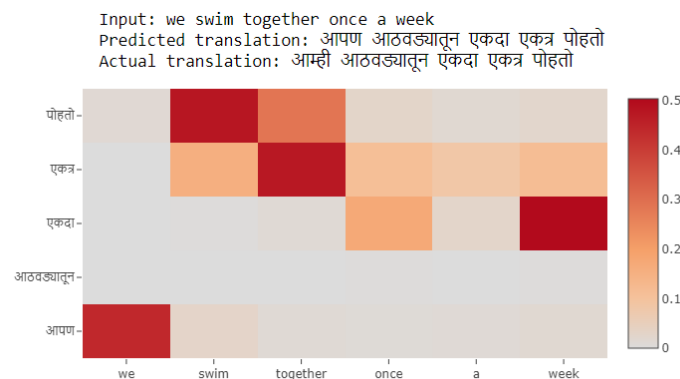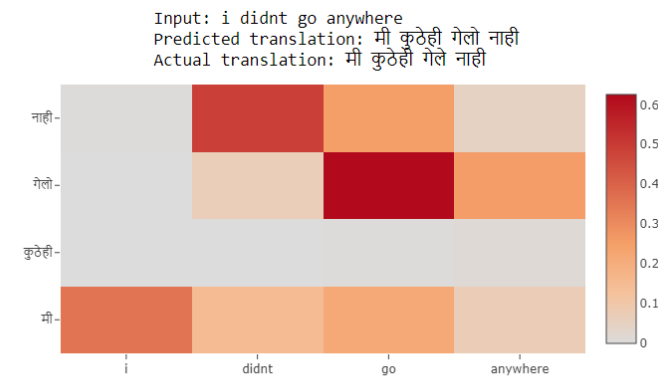
## 6. Visualizing the Results

Input: father bought me a camera
Predicted translation: बाबांनी माझ्यासाठी कॅमेरा विकत घेतला
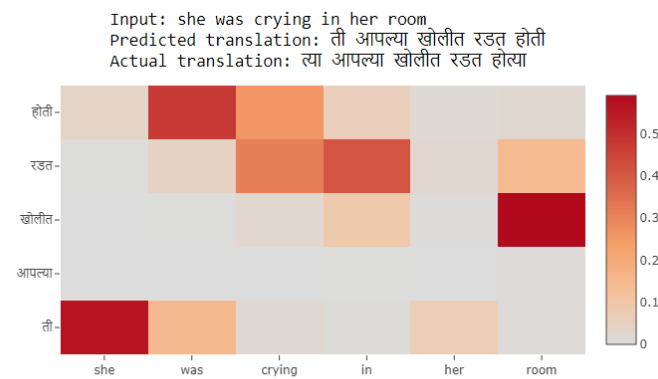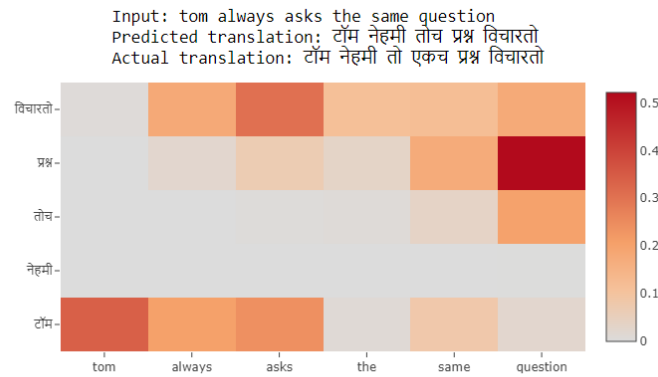Actual translation: बाबांनी माझ्यासाठी एक कॅमेरा विकत घेतला



If you are new to heat maps, this is how you can interpret the above plot:

Notice that the cell at the intersection of "father" and "बाबांनी" is pretty dark This means when the decoder predicts the word "बाबांनी", it is paying more attention to the input word "father" (which is what we wanted).

Similarly while predicting the word "कॅमेरा", the decoder pays a lot of attention to the input word "camera". And so on.

Input: we swim together once a week
Predicted translation: आपण आठवड्यातून एकदा एकत्र पोहतो
Actual translation: आम्ही आठवड्यातून एकदा एकत्र पोहतो



Input: i can come at three
Predicted translation: मी तीन वाजता येऊ शकते
Actual translation: मी तीन वाजता येऊ शकतो

Input: tom always asks the same question
Predicted translation: टॉम नेहमी तोच प्रश्न विचारतो
Actual translation: टॉम नेहमी तो एकच प्रश्न विचारतो

Input: she was crying in her room
Predicted translation: ती आपल्या खोलीत रडत होती
Actual translation: त्या आपल्या खोलीत रडत होत्या

Input: i didnt go anywhere
Predicted translation: मी कुठेही गेलो नाही
Actual translation: मी कुठेही गेले नाही

## Conclusion

The first thing to be noted is that the translation results are much better than the results from my previous *blog*. Secondly the model is able to find the correct local mappings between the input and the output sequences which do match with our intuition.

Given more data and with more hyper parameter tuning, the results and mappings will definitely improve by a good margin.

Using LSTM layers in place of GRU and adding Bidirectional wrapper

on the encoder will also help in improved performance.

Deep Learning models are generally considered as black boxes, meaning that they do not have the ability to explain their outputs. However, Attention is one of the successful methods that helps to make our model interpretable and explain why it does what it does.

The only disadvantage of the Attention mechanism is that it is a very time consuming and hard to parallelize system. To solve this problem, Google Brain came up with the "Transformer Model" which uses only Attention and gets rid of all the Convolutional and Recurrent Layers, thus making it highly parallelizable and compute efficient.

## 7. References

https://arxiv.org/abs/1409.0473 (Original Paper)

https://github.com/tensorflow/tensorflow/blob/master /tensorflow/contrib/eager/python/examples /nmt_with_attention/nmt_with_attention.ipynb (TensorFlow Implementation available on their official website as a tutorial)

https://www.coursera.org/lecture/nlp-sequence-models /attention-model-lSwVa (Andrew Ng's Explanation on Attention)

https://jalammar.github.io/visualizing-neural-machine-translation-mechanics-of-seq2seq-models-with-attention/

https://www.tensorflow.org/xla/broadcasting (Broadcasting in TensorFlow)

Dataset: http://www.manythings.org/anki/ (mar-eng.zip)

PS: For complete implementation, refer my GitHub repository *here*.