

Capstone Project

April 8, 2024

1 Capstone Project

1.1 Neural translation model

1.1.1 Instructions

In this notebook, you will create a neural network that translates from English to German. You will use concepts from throughout this course, including building more flexible model architectures, freezing layers, data processing pipeline and sequence modelling.

This project is peer-assessed. Within this notebook you will find instructions in each section for how to complete the project. Pay close attention to the instructions as the peer review will be carried out according to a grading rubric that checks key parts of the project instructions. Feel free to add extra cells into the notebook as required.

1.1.2 How to submit

When you have completed the Capstone project notebook, you will submit a pdf of the notebook for peer review. First ensure that the notebook has been fully executed from beginning to end, and all of the cell outputs are visible. This is important, as the grading rubric depends on the reviewer being able to view the outputs of your notebook. Save the notebook as a pdf (File -> Download as -> PDF via LaTeX). You should then submit this pdf for review.

1.1.3 Let's get started!

We'll start by running some imports, and loading the dataset. For this project you are free to make further imports throughout the notebook as you wish.

```
In [ ]: import tensorflow as tf
import tensorflow_hub as hub
import unicodedata
import re
```

For the capstone project, you will use a language dataset from <http://www.manythings.org/anki/> to build a neural translation model. This dataset consists of over 200,000 pairs of sentences in English and German. In order to make the training quicker, we will restrict to our dataset to 20,000 pairs. Feel free to change this if you wish - the size of the dataset used is not part of the grading rubric.

Your goal is to develop a neural translation model from English to German, making use of a pre-trained English word embedding module.



Flags overview image

```
In [ ]: # Run this cell to load the dataset
```

```
NUM_EXAMPLES = 20000
data_examples = []
with open('data/deu.txt', 'r', encoding='utf8') as f:
    for line in f.readlines():
        if len(data_examples) < NUM_EXAMPLES:
            data_examples.append(line)
        else:
            break
```

```
In [ ]: # These functions preprocess English and German sentences
```

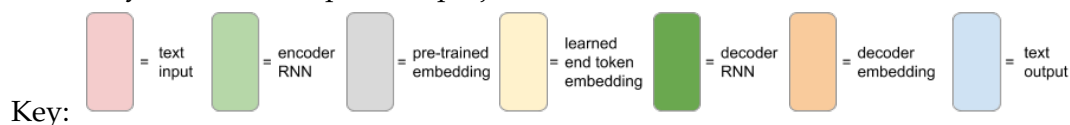
```
def unicode_to_ascii(s):
    return ''.join(c for c in unicodedata.normalize('NFD', s) if unicodedata.category(c) != 'Mn')

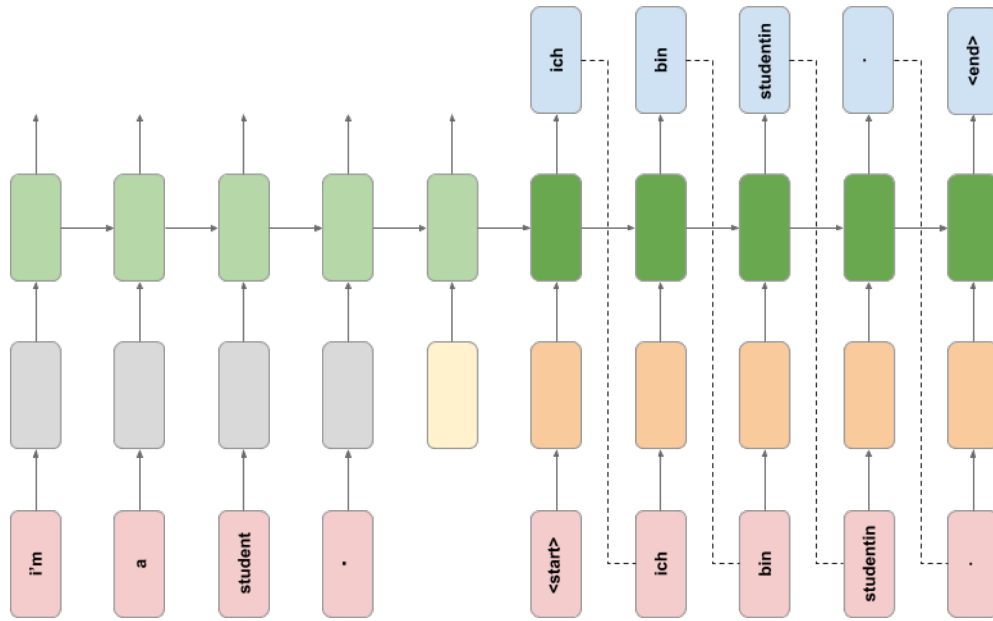
def preprocess_sentence(sentence):
    sentence = sentence.lower().strip()
    sentence = re.sub(r"ü", 'ue', sentence)
    sentence = re.sub(r"ä", 'ae', sentence)
    sentence = re.sub(r"ö", 'oe', sentence)
    sentence = re.sub(r" ", 'ss', sentence)

    sentence = unicode_to_ascii(sentence)
    sentence = re.sub(r"([?.,])", r" \1 ", sentence)
    sentence = re.sub(r"^[^a-z?.,,]+", " ", sentence)
    sentence = re.sub(r"[" "]+" , " ", sentence)

    return sentence.strip()
```

The custom translation model The following is a schematic of the custom translation model architecture you will develop in this project.





Model Schematic

The custom model consists of an encoder RNN and a decoder RNN. The encoder takes words of an English sentence as input, and uses a pre-trained word embedding to embed the words into a 128-dimensional space. To indicate the end of the input sentence, a special end token (in the same 128-dimensional space) is passed in as an input. This token is a TensorFlow Variable that is learned in the training phase (unlike the pre-trained word embedding, which is frozen).

The decoder RNN takes the internal state of the encoder network as its initial state. A start token is passed in as the first input, which is embedded using a learned German word embedding. The decoder RNN then makes a prediction for the next German word, which during inference is then passed in as the following input, and this process is repeated until the special <end> token is emitted from the decoder.

1.2 1. Text preprocessing

- Create separate lists of English and German sentences, and preprocess them using the `preprocess_sentence` function provided for you above.
- Add a special "<start>" and "<end>" token to the beginning and end of every German sentence.
- Use the `Tokenizer` class from the `tf.keras.preprocessing.text` module to tokenize the German sentences, ensuring that no character filters are applied. *Hint: use the `Tokenizer`'s "filter" keyword argument.*
- Print out at least 5 randomly chosen examples of (preprocessed) English and German sentence pairs. For the German sentence, print out the text (with start and end tokens) as well as the tokenized sequence.
- Pad the end of the tokenized German sequences with zeros, and batch the complete set of sequences into one numpy array.

```
In [ ]: import numpy as np
import re
```

```

import unicodedata
from tensorflow.keras.preprocessing.text import Tokenizer
from tensorflow.keras.preprocessing.sequence import pad_sequences
import random

In [ ]: # data_examples
english_sentences = [line.split('\t')[0] for line in data_examples]
german_sentences = ["<start> " + line.split('\t')[1] + " <end>" for line in data_examples]

# preprocess_sentence
english_sentences = [preprocess_sentence(sentence) for sentence in english_sentences]
german_sentences = [preprocess_sentence(sentence) for sentence in german_sentences]

# Tokenizer
tokenizer = Tokenizer(filters='') #
tokenizer.fit_on_texts(german_sentences)
german_sequences = tokenizer.texts_to_sequences(german_sentences)

# 5()
for _ in range(5):
    index = random.randint(0, len(english_sentences) - 1)
    print(f"English: {english_sentences[index]}")
    print(f"German (Text): {german_sentences[index]}")
    print(f"German (Sequence): {german_sequences[index]}\n")

#
german_sequences_padded = pad_sequences(german_sequences, padding='post')

print("Padded German Sequences Shape:", german_sequences_padded.shape)

In [ ]:

In [ ]:

In [ ]:

In [ ]:

```

1.3 2. Prepare the data with tf.data.Dataset objects

Load the embedding layer As part of the dataset preprocessing for this project, you will use a pre-trained English word embedding module from TensorFlow Hub. The URL for the module is <https://tfhub.dev/google/tf2-preview/nnlm-en-dim128-with-normalization/1>. This module has also been made available as a complete saved model in the folder './models/tf2-preview_nnlm-en-dim128_1'.

This embedding takes a batch of text tokens in a 1-D tensor of strings as input. It then embeds the separate tokens into a 128-dimensional space.

The code to load and test the embedding layer is provided for you below.

NB: this model can also be used as a sentence embedding module. The module will process each token by removing punctuation and splitting on spaces. It then averages the word embeddings over a sentence to give a single embedding vector. However, we will use it only as a word embedding module, and will pass each word in the input sentence as a separate token.

```
In [6]: embedding_layer = hub.KerasLayer("./models/tf2-preview_nnlm-en-dim128_1",  
                                         output_shape=[128], input_shape=[], dtype=tf.string)
```

```
In [7]: # Test the layer
```

```
embedding_layer(tf.constant(["these", "aren't", "the", "droids", "you're", "looking"],
```

```
Out [7]: TensorShape([7, 128])
```

You should now prepare the training and validation Datasets.

- Create a random training and validation set split of the data, reserving e.g. 20% of the data for validation (NB: each English dataset example is a single sentence string, and each German dataset example is a sequence of padded integer tokens).
- Load the training and validation sets into a `tf.data.Dataset` object, passing in a tuple of English and German data for both training and validation sets.
- Create a function to map over the datasets that splits each English sentence at spaces. Apply this function to both Dataset objects using the map method. *Hint: look at the `tf.strings.split` function.*
- Create a function to map over the datasets that embeds each sequence of English words using the loaded embedding layer/model. Apply this function to both Dataset objects using the map method.
- Create a function to filter out dataset examples where the English sentence is more than 13 (embedded) tokens in length. Apply this function to both Dataset objects using the filter method.
- Create a function to map over the datasets that pads each English sequence of embeddings with some distinct padding value before the sequence, so that each sequence is length 13. Apply this function to both Dataset objects using the map method. *Hint: look at the `tf.pad` function. You can extract a Tensor shape using `tf.shape`; you might also find the `tf.math.maximum` function useful.*
- Batch both training and validation Datasets with a batch size of 16.
- Print the `element_spec` property for the training and validation Datasets.
- Using the Dataset `.take(1)` method, print the shape of the English data example from the training Dataset.
- Using the Dataset `.take(1)` method, print the German data example Tensor from the validation Dataset.

```
In [8]: from sklearn.model_selection import train_test_split  
        english_train, english_val, german_train, german_val = train_test_split(  
            english_sentences, german_sequences_padded, test_size=0.2, random_state=42)  
        train_dataset = tf.data.Dataset.from_tensor_slices((english_train, german_train))  
        val_dataset = tf.data.Dataset.from_tensor_slices((english_val, german_val))
```

```
In [9]: def split_english(sentence, german):  
        return tf.strings.split(sentence, ' '), german
```

```

train_dataset = train_dataset.map(split_english)
val_dataset = val_dataset.map(split_english)

In [10]: def embed_english(english, german):
          return embedding_layer(english), german

          train_dataset = train_dataset.map(embed_english)
          val_dataset = val_dataset.map(embed_english)

In [11]: def filter_long_sentences(english, german):
          return tf.shape(english)[0] <= 13

          train_dataset = train_dataset.filter(filter_long_sentences)
          val_dataset = val_dataset.filter(filter_long_sentences)

In [12]: def pad_english(english, german):
          paddings = [[13 - tf.shape(english)[0], 0], [0, 0]]
          return tf.pad(english, paddings, "CONSTANT"), german

          train_dataset = train_dataset.map(pad_english)
          val_dataset = val_dataset.map(pad_english)

In [13]: train_dataset = train_dataset.batch(16)
          val_dataset = val_dataset.batch(16)

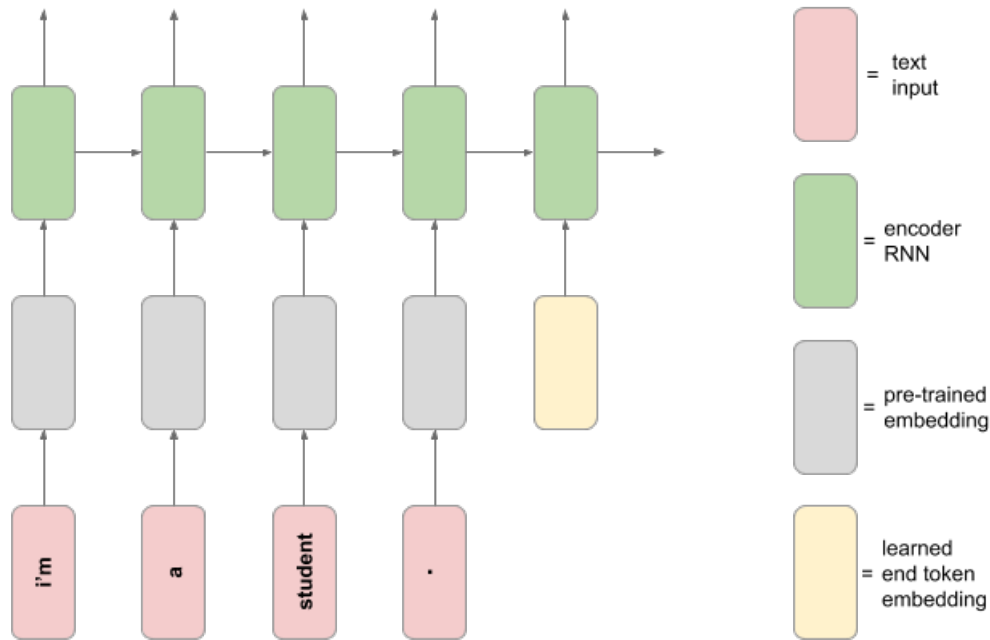
          # 8. element_spec
          print("Training dataset element_spec:", train_dataset.element_spec)
          print("Validation dataset element_spec:", val_dataset.element_spec)

Training dataset element_spec: (TensorSpec(shape=(None, None, 128), dtype=tf.float32, name=None),
Validation dataset element_spec: (TensorSpec(shape=(None, None, 128), dtype=tf.float32, name=None),

In [ ]: for english, german in train_dataset.take(1):
          print("Shape of the English data example from training dataset:", english.shape)
          for english, german in val_dataset.take(1):
              print("German data example Tensor from validation dataset:", german)

Shape of the English data example from training dataset: (16, 13, 128)
German data example Tensor from validation dataset: tf.Tensor(
[[ 1  4 18 20 115 178 3 2 0 0 0 0 0 0]
 [ 1  4 61 315 3 2 0 0 0 0 0 0 0 0]
 [ 1  5 6 831 3 2 0 0 0 0 0 0 0 0]
 [ 1 42 457 9 2 0 0 0 0 0 0 0 0 0]
 [ 1 78 4750 3 2 0 0 0 0 0 0 0 0 0]
 [ 1 5 414 20 75 36 3 2 0 0 0 0 0 0]
 [ 1 6 47 312 5 7 2 0 0 0 0 0 0 0]
 [ 1 120 21 76 90 9 2 0 0 0 0 0 0 0]

```



Encoder schematic

```
[ 1  4 766  11 113  3  2  0  0  0  0  0  0  0]
[ 1  8 308 129  3  2  0  0  0  0  0  0  0  0]
[ 1  4 580  3  2  0  0  0  0  0  0  0  0  0]
[ 1 64  22 100 45  3  2  0  0  0  0  0  0  0]
[ 1  5  16  19 955  3  2  0  0  0  0  0  0  0]
[ 1  5 136 186 5141  3  2  0  0  0  0  0  0  0]
[ 1  5 165 271  20 54  3  2  0  0  0  0  0  0]
[ 1 21  6  56 232  3  2  0  0  0  0  0  0  0]], shape=(16, 14), dtype=
```

1.4 3. Create the custom layer

You will now create a custom layer to add the learned end token embedding to the encoder model:

You should now build the custom layer. * Using layer subclassing, create a custom layer that takes a batch of English data examples from one of the Datasets, and adds a learned embedded 'end' token to the end of each sequence. * This layer should create a TensorFlow Variable (that will be learned during training) that is 128-dimensional (the size of the embedding space). *Hint: you may find it helpful in the call method to use the `tf.tile` function to replicate the end token embedding across every element in the batch.* * Using the Dataset `.take(1)` method, extract a batch of English data examples from the training Dataset and print the shape. Test the custom layer by calling the layer on the English data batch Tensor and print the resulting Tensor shape (the layer should increase the sequence length by one).

```
In [ ]: class AddEndTokenLayer(tf.keras.layers.Layer):
        def __init__(self, **kwargs):
            super(AddEndTokenLayer, self).__init__(**kwargs)
            # "end"
```

```

        self.end_token_embedding = tf.Variable(initial_value=tf.random.uniform([128]),

def call(self, inputs):
    #
    batch_size = tf.shape(inputs)[0]
    # tf.tile"end"
    end_tokens = tf.tile(tf.reshape(self.end_token_embedding, shape=(1, 1, 128)),
    # "end"
    return tf.concat([inputs, end_tokens], axis=1)

In [ ]: for english, _ in train_dataset.take(1):
    print("Original English data shape:", english.shape)

    #
    add_end_token_layer = AddEndTokenLayer()
    english_with_end_token = add_end_token_layer(english)

    print("English data shape after adding end token:", english_with_end_token.shape)

Original English data shape: (16, 13, 128)
English data shape after adding end token: (16, 14, 128)

```

```
In [ ]:
```

1.5 4. Build the encoder network

The encoder network follows the schematic diagram above. You should now build the RNN encoder model. * Using the functional API, build the encoder network according to the following spec: * The model will take a batch of sequences of embedded English words as input, as given by the Dataset objects. * The next layer in the encoder will be the custom layer you created previously, to add a learned end token embedding to the end of the English sequence. * This is followed by a Masking layer, with the `mask_value` set to the distinct padding value you used when you padded the English sequences with the Dataset preprocessing above. * The final layer is an LSTM layer with 512 units, which also returns the hidden and cell states. * The encoder is a multi-output model. There should be two output Tensors of this model: the hidden state and cell states of the LSTM layer. The output of the LSTM layer is unused. * Using the Dataset `.take(1)` method, extract a batch of English data examples from the training Dataset and test the encoder model by calling it on the English data Tensor, and print the shape of the resulting Tensor outputs. * Print the model summary for the encoder network.

```

In [ ]: import tensorflow as tf
        from tensorflow.keras.models import Model
        from tensorflow.keras.layers import Input, LSTM, Masking
        from tensorflow.keras.preprocessing.sequence import pad_sequences

        inputs = Input(shape=(None, 128)) # 128

        # "end"

```



```

x = AddEndTokenLayer()(inputs)

# Masking
# Omask_value
x = Masking(mask_value=0.0)(x)

# LSTM
_, hidden_state, cell_state = LSTM(512, return_state=True)(x)

#
encoder = Model(inputs=inputs, outputs=[hidden_state, cell_state])

In [ ]: for english_batch, _ in train_dataset.take(1):
        hidden, cell = encoder(english_batch)
        print("Hidden state shape:", hidden.shape)
        print("Cell state shape:", cell.shape)
        encoder.summary()

```

Hidden state shape: (16, 512)

Cell state shape: (16, 512)

Model: "model"

Layer (type)	Output Shape	Param #
input_1 (InputLayer)	[(None, None, 128)]	0
add_end_token_layer_1 (AddEn	(None, None, 128)	128
masking (Masking)	(None, None, 128)	0
lstm (LSTM)	[(None, 512), (None, 512)	1312768

Total params: 1,312,896

Trainable params: 1,312,896

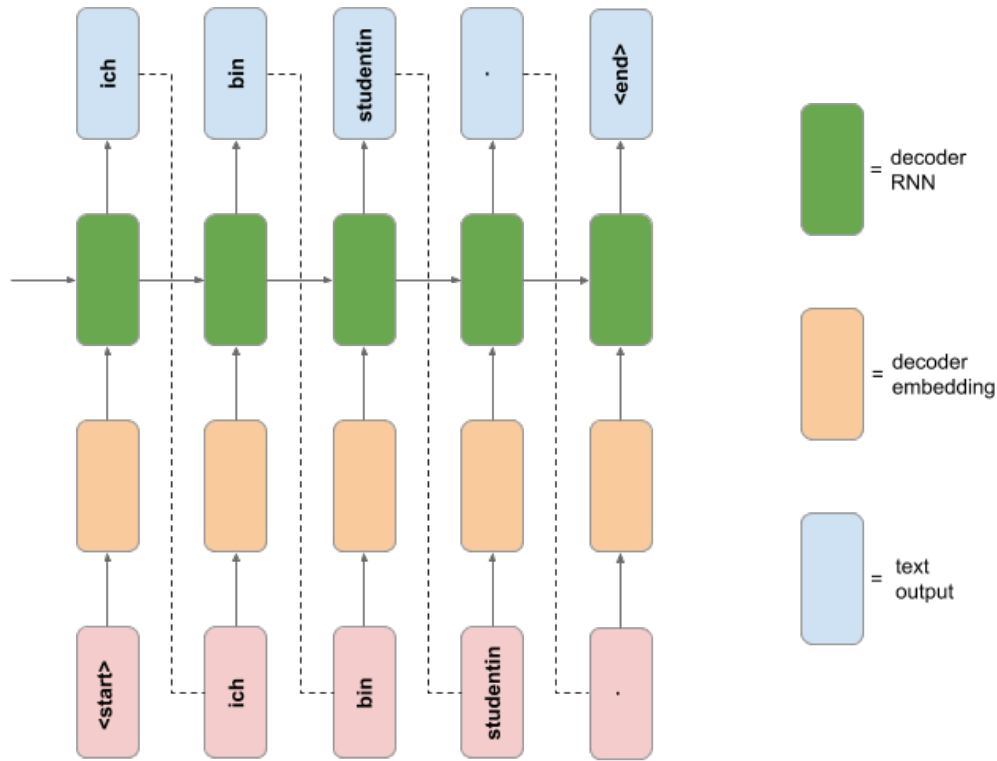
Non-trainable params: 0

```
In [ ]:
```

1.6 5. Build the decoder network

The decoder network follows the schematic diagram below.

You should now build the RNN decoder model. * Using Model subclassing, build the decoder network according to the following spec: * The initializer should create the following layers: * An Embedding layer with vocabulary size set to the number of unique German tokens, embedding dimension 128, and set to mask zero values in the input. * An LSTM layer with 512 units, that returns its hidden and cell states, and also returns sequences. * A Dense layer with number of units equal to the number of unique German tokens, and no activation function. * The call



Decoder schematic

method should include the usual `inputs` argument, as well as the additional keyword arguments `hidden_state` and `cell_state`. The default value for these keyword arguments should be `None`.
 * The call method should pass the inputs through the Embedding layer, and then through the LSTM layer. If the `hidden_state` and `cell_state` arguments are provided, these should be used for the initial state of the LSTM layer. *Hint: use the `initial_state` keyword argument when calling the LSTM layer on its input.*
 * The call method should pass the LSTM output sequence through the Dense layer, and return the resulting Tensor, along with the hidden and cell states of the LSTM layer.
 * Using the `Dataset .take(1)` method, extract a batch of English and German data examples from the training Dataset. Test the decoder model by first calling the encoder model on the English data Tensor to get the hidden and cell states, and then call the decoder model on the German data Tensor and hidden and cell states, and print the shape of the resulting decoder Tensor outputs. *
 Print the model summary for the decoder network.

```
In [ ]: from tensorflow.keras.layers import Embedding, LSTM, Dense
        from tensorflow.keras import Model
```

```
class Decoder(Model):
    def __init__(self, vocab_size, **kwargs):
        super(Decoder, self).__init__(**kwargs)
        self.embedding = Embedding(input_dim=vocab_size, output_dim=128, mask_zero=True)
        self.lstm = LSTM(512, return_sequences=True, return_state=True)
        self.dense = Dense(vocab_size)

    def call(self, inputs, hidden_state=None, cell_state=None):
```

```

x = self.embedding(inputs)
if hidden_state is not None and cell_state is not None:
    x, hidden_state, cell_state = self.lstm(x, initial_state=[hidden_state, ce
else:
    x, hidden_state, cell_state = self.lstm(x)
x = self.dense(x)
return x, hidden_state, cell_state

```

```
In [ ]: vocab_size = len(tokenizer.word_index) + 1
```

```
decoder = Decoder(vocab_size=vocab_size)
```

```

#
for english_batch, german_batch in train_dataset.take(1):
    hidden, cell = encoder(english_batch)
    decoder_output, _, _ = decoder(german_batch, hidden, cell)
    print("Decoder output shape:", decoder_output.shape)

```

```
decoder.summary()
```

```
Decoder output shape: (16, 14, 5744)
```

```
Model: "decoder"
```

Layer (type)	Output Shape	Param #
embedding (Embedding)	multiple	735232
lstm_1 (LSTM)	multiple	1312768
dense (Dense)	multiple	2946672
Total params: 4,994,672		
Trainable params: 4,994,672		
Non-trainable params: 0		

```
In [ ]:
```

1.7 6. Make a custom training loop

You should now write a custom training loop to train your custom neural translation model.

- * Define a function that takes a Tensor batch of German data (as extracted from the training Dataset), and returns a tuple containing German inputs and outputs for the decoder model (refer to schematic diagram above).
- * Define a function that computes the forward and backward pass for your translation model. This function should take an English input, German input and German output as arguments, and should do the following:
 - * Pass the English input into the encoder, to get the hidden and cell states of the encoder LSTM.
 - * These hidden and cell states are then passed

into the decoder, along with the German inputs, which returns a sequence of outputs (the hidden and cell state outputs of the decoder LSTM are unused in this function). * The loss should then be computed between the decoder outputs and the German output function argument. * The function returns the loss and gradients with respect to the encoder and decoder's trainable variables. * Decorate the function with `@tf.function` * Define and run a custom training loop for a number of epochs (for you to choose) that does the following: * Iterates through the training dataset, and creates decoder inputs and outputs from the German sequences. * Updates the parameters of the translation model using the gradients of the function above and an optimizer object. * Every epoch, compute the validation loss on a number of batches from the validation and save the epoch training and validation losses. * Plot the learning curves for loss vs epoch for both training and validation sets.

Hint: This model is computationally demanding to train. The quality of the model or length of training is not a factor in the grading rubric. However, to obtain a better model we recommend using the GPU accelerator hardware on Colab.

```
In [ ]: def prepare_german_data(german_batch):
    #
    german_input = german_batch[:, :-1]
    #
    german_output = german_batch[:, 1:]
    return german_input, german_output

In [ ]: @tf.function
def train_step(encoder, decoder, english_input, german_input, german_output, loss_function):
    with tf.GradientTape() as tape:
        #
        encoder_hidden, encoder_cell = encoder(english_input)
        #
        decoder_output, _, _ = decoder(german_input, hidden_state=encoder_hidden, cell=encoder_cell)

        loss = loss_function(german_output, decoder_output)

        trainable_vars = encoder.trainable_variables + decoder.trainable_variables
        gradients = tape.gradient(loss, trainable_vars)

        optimizer.apply_gradients(zip(gradients, trainable_vars))

    return loss

In [ ]: epochs = 10
optimizer = tf.keras.optimizers.Adam()
loss_function = tf.keras.losses.SparseCategoricalCrossentropy(from_logits=True)

train_loss_results = []
val_loss_results = []

for epoch in range(epochs):
```

```

epoch_loss_avg = tf.keras.metrics.Mean()
val_loss_avg = tf.keras.metrics.Mean()

for english_batch, german_batch in train_dataset:
    german_input, german_output = prepare_german_data(german_batch)
    loss = train_step(encoder, decoder, english_batch, german_input, german_output)
    epoch_loss_avg.update_state(loss)

for english_batch, german_batch in val_dataset:
    german_input, german_output = prepare_german_data(german_batch)
    encoder_hidden, encoder_cell = encoder(english_batch)
    decoder_output, _, _ = decoder(german_input, hidden_state=encoder_hidden, cell=encoder_cell)
    loss = loss_function(german_output, decoder_output)
    val_loss_avg.update_state(loss)

train_loss_results.append(epoch_loss_avg.result())
val_loss_results.append(val_loss_avg.result())

print(f"Epoch {epoch+1}: Training Loss: {epoch_loss_avg.result()}, Validation Loss: {val_loss_avg.result()}")

```

```
In [ ]: import matplotlib.pyplot as plt
```

```

plt.figure(figsize=(8, 6))
plt.plot(train_loss_results, label='Training Loss')
plt.plot(val_loss_results, label='Validation Loss')
plt.xlabel('Epoch')
plt.ylabel('Loss')
plt.legend()
plt.title('Loss vs. Epoch')
plt.show()

```

```
In [ ]:
```

```
In [ ]:
```

1.8 7. Use the model to translate

Now it's time to put your model into practice! You should run your translation for five randomly sampled English sentences from the dataset. For each sentence, the process is as follows: * Preprocess and embed the English sentence according to the model requirements. * Pass the embedded sentence through the encoder to get the encoder hidden and cell states. * Starting with the special "<start>" token, use this token and the final encoder hidden and cell states to get the one-step prediction from the decoder, as well as the decoder's updated hidden and cell states. * Create a loop to get the next step prediction and updated hidden and cell states from the decoder, using the most recent hidden and cell states. Terminate the loop when the "<end>" token is emitted, or when the sentence has reached a maximum length. * Decode the output token sequence into German text and print the English text and the model's German translation.

```
In [ ]: import numpy as np
```

```
def translate_sentence(sentence):
    # 1:
    sentence = preprocess_sentence(sentence)
    embedded_sentence = embedding_layer(tf.constant([sentence])) # embedding_layer
    # 2:
    encoder_hidden, encoder_cell = encoder(embedded_sentence)

    # "<start>"
    decoder_input = tf.expand_dims([tokenizer.word_index['<start>']], 0)
    result = []

    for _ in range(100):
        # 34:
        predictions, decoder_hidden, decoder_cell = decoder(decoder_input,
                                                              hidden_state=encoder_hidden,
                                                              cell_state=encoder_cell)

        predicted_id = tf.argmax(predictions[0]).numpy()
        result.append(tokenizer.index_word[predicted_id])

        # "<end>"
        if tokenizer.index_word[predicted_id] == '<end>':
            return ' '.join(result[:-1]) # "<end>"

        # ID
        decoder_input = tf.expand_dims([predicted_id], 0)

    return ' '.join(result)
```

```
In [ ]: random_sentences = np.random.choice(english_sentences, 5)
        for sentence in random_sentences:
            print(f'English: {sentence}')
            print(f'German Translation: {translate_sentence(sentence)}\n')
```

```
In [ ]:
```