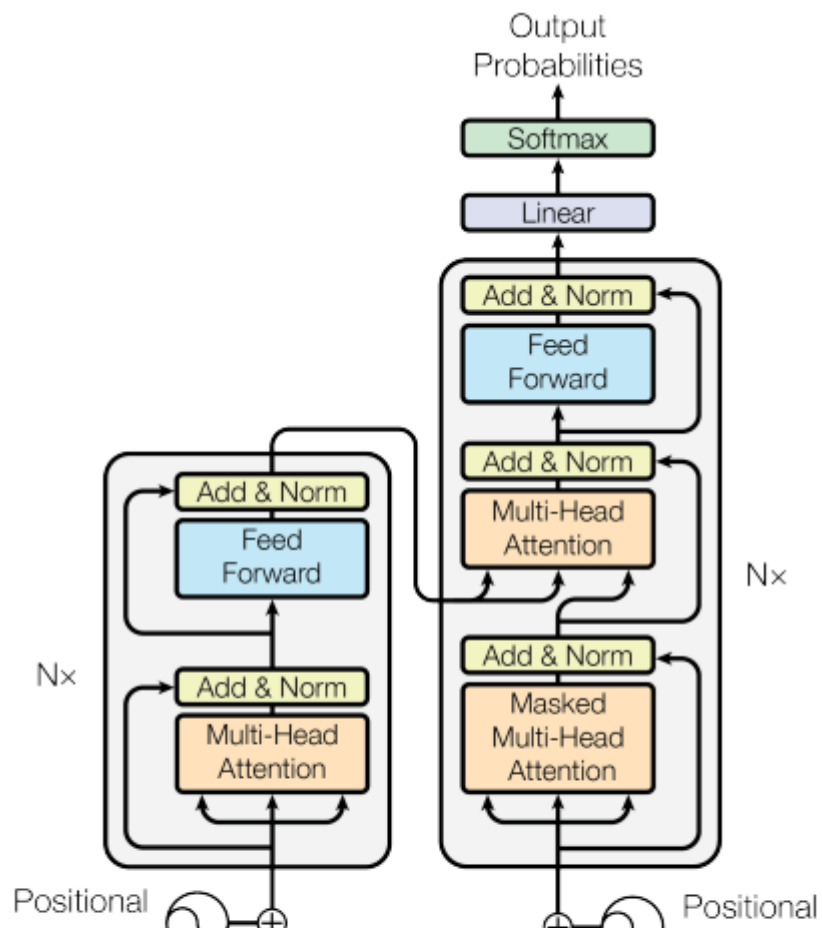


## ▼ 6 - Attention is All You Need

In this notebook we will be implementing a (slightly modified version) of the Transformer model from the [Attention is All You Need](#) paper. All images in this notebook will be taken from the Transformer paper. For more information about the Transformer, [see these three](#) articles.



## ▼ Preparing the Data

As always, let's import all the required modules and set the random seeds for reproducibility.

```

Inputs      Outputs
import torch
import torch.nn as nn
import torch.optim as optim

import torchtext
from torchtext.legacy.datasets import Multi30k
from torchtext.legacy.data import Field, BucketIterator

import matplotlib.pyplot as plt

```



```
Requirement already satisfied: urllib3!=1.25.0,!=1.25.1,<1.26,>=1.21.1 in /usr/local/lib/python3.7/dist-packages (from requests<3.0.0,>=2.13.0->spacy>=
Requirement already satisfied: chardet<4,>=3.0.2 in /usr/local/lib/python3.7/dist-packages (from requests<3.0.0,>=2.13.0->spacy>=2.2.2->de_core_news_sm
Requirement already satisfied: certifi>=2017.4.17 in /usr/local/lib/python3.7/dist-packages (from requests<3.0.0,>=2.13.0->spacy>=2.2.2->de_core_news_s
Requirement already satisfied: idna<3,>=2.5 in /usr/local/lib/python3.7/dist-packages (from requests<3.0.0,>=2.13.0->spacy>=2.2.2->de_core_news_sm==2.2.2)
✓ Download and installation successful
You can now load the model via spacy.load('de_core_news_sm')
```

We'll then create our tokenizers as before.

```
spacy_de = spacy.load('de_core_news_sm')
spacy_en = spacy.load('en_core_web_sm')

def tokenize_de(text):
    """
    Tokenizes German text from a string into a list of strings
    """
    return [tok.text for tok in spacy_de.tokenizer(text)]

def tokenize_en(text):
    """
    Tokenizes English text from a string into a list of strings
    """
    return [tok.text for tok in spacy_en.tokenizer(text)]
```

Our fields are the same as the previous notebook. The model expects data to be fed in with the batch dimension first, so we use `batch_first = True`.

```
SRC = Field(tokenize = tokenize_de,
            init_token = '<sos>',
            eos_token = '<eos>',
            lower = True,
            batch_first = True)

TRG = Field(tokenize = tokenize_en,
```

```
init_token = '<sos>',
eos_token = '<eos>',
lower = True,
batch_first = True)
```

We then load the Multi30k dataset and build the vocabulary.

```
train_data, valid_data, test_data = Multi30k.splits(exts = ('.de', '.en'),
                                                    fields = (SRC, TRG))
```

```
SRC.build_vocab(train_data, min_freq = 2)
TRG.build_vocab(train_data, min_freq = 2)
```

Finally, we define the device and the data iterator.

```
device = torch.device('cuda' if torch.cuda.is_available() else 'cpu')
```

```
BATCH_SIZE = 128
```

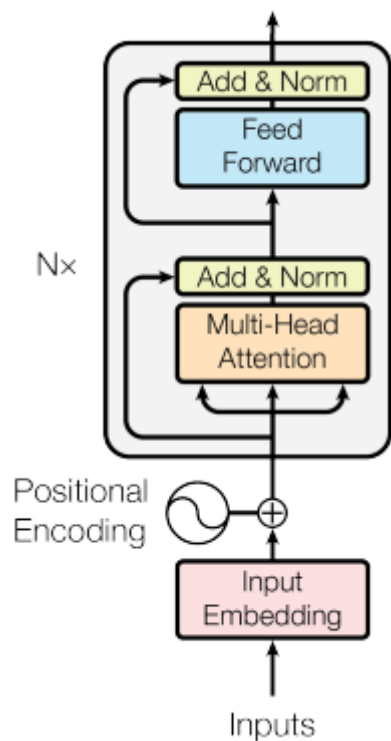
```
train_iterator, valid_iterator, test_iterator = BucketIterator.splits(
    (train_data, valid_data, test_data),
    batch_size = BATCH_SIZE,
    device = device)
```

## ▼ Building the Model

Next, we'll build the model. Like previous notebooks it is made up of an *encoder* and a *decoder*, with the encoder *encoding* the input/source sentence (in German) into *context vector* and the decoder then *decoding* this context vector to output our output/target sentence (in English).

### Encoder

Similar to the ConvSeq2Seq model, the Transformer's encoder does not attempt to compress the entire source sentence,  $X = (x_1, \dots, x_n)$ , into a single context vector,  $z$ . Instead it produces a sequence of context vectors,  $Z = (z_1, \dots, z_n)$ . So, if our input sequence was 5 tokens long we would have  $Z = (z_1, z_2, z_3, z_4, z_5)$ . Why do we call this a sequence of context vectors and not a sequence of hidden states? A hidden state at time  $t$  in an RNN has only seen tokens  $x_t$  and all the tokens before it. However, each context vector here has seen all tokens at all positions within the input sequence.



First, the tokens are passed through a standard embedding layer. Next, as the model has no recurrent it has no idea about the order of the tokens within the sequence. We solve this by using a second embedding layer called a *positional embedding layer*. This is a standard embedding layer where the input is not the token itself but the position of the token within the sequence, starting with the first token, the  $\langle \text{sos} \rangle$  (start of sequence) token, in position 0. The position embedding has a "vocabulary" size of 100, which means our model can accept sentences up to 100 tokens long. This can be increased if we want to handle longer sentences.

The original Transformer implementation from the Attention is All You Need paper does not learn positional embeddings. Instead it uses a fixed static embedding. Modern Transformer architectures, like BERT, use positional embeddings instead, hence we have decided to use them in these tutorials. Check out [this](#) section to read more about the positional embeddings used in the original Transformer model.

Next, the token and positional embeddings are elementwise summed together to get a vector which contains information about the token and also its position within the sequence. However, before they are summed, the token embeddings are multiplied by a scaling factor which is  $\sqrt{d_{model}}$ , where  $d_{model}$  is the hidden dimension size, `hid_dim`. This supposedly reduces variance in the embeddings and the model is difficult to train reliably without this scaling factor. Dropout is then applied to the combined embeddings.

The combined embeddings are then passed through  $N$  *encoder layers* to get  $Z$ , which is then output and can be used by the decoder.

The source mask, `src_mask`, is simply the same shape as the source sentence but has a value of 1 when the token in the source sentence is not a `<pad>` token and 0 when it is a `<pad>` token. This is used in the encoder layers to mask the multi-head attention mechanisms, which are used to calculate and apply attention over the source sentence, so the model does not pay attention to `<pad>` tokens, which contain no useful information

```
class Encoder(nn.Module):
    def __init__(self,
                    input_dim,
                    hid_dim,
                    n_layers,
                    n_heads,
                    pf_dim,
                    dropout,
                    device,
                    max_length = 100):
        super().__init__()

        self.device = device

        self.tok_embedding = nn.Embedding(input_dim, hid_dim)
        self.pos_embedding = nn.Embedding(max_length, hid_dim)

        self.layers = nn.ModuleList([EncoderLayer(hid_dim,
                                                    n_heads,
                                                    pf_dim,
                                                    dropout,
                                                    device)
                                      for _ in range(n_layers)])

        self.dropout = nn.Dropout(dropout)
```

```

        self.scale = torch.sqrt(torch.FloatTensor([hid_dim])).to(device)

    def forward(self, src, src_mask):

        #src = [batch size, src len]
        #src_mask = [batch size, 1, 1, src len]

        batch_size = src.shape[0]
        src_len = src.shape[1]

        pos = torch.arange(0, src_len).unsqueeze(0).repeat(batch_size, 1).to(self.device)

        #pos = [batch size, src len]

        src = self.dropout((self.tok_embedding(src) * self.scale) + self.pos_embedding(pos))

        #src = [batch size, src len, hid dim]

        for layer in self.layers:
            src = layer(src, src_mask)

        #src = [batch size, src len, hid dim]

        return src

```

## ▼ Encoder Layer

The encoder layers are where all of the "meat" of the encoder is contained. We first pass the source sentence and its mask into the *multi-head attention layer*, then perform dropout on it, apply a residual connection and pass it through a [Layer Normalization](#) layer. We then pass it through a *position-wise feedforward* layer and then, again, apply dropout, a residual connection and then layer normalization to get the output of this layer which is fed into the next layer. The parameters are not shared between layers.

The multi head attention layer is used by the encoder layer to attend to the source sentence, i.e. it is calculating and applying attention over itself instead of another sequence, hence we call it *self attention*.



[This](#) article goes into more detail about layer normalization, but the gist is that it normalizes the values of the features, i.e. across the hidden dimension, so each feature has a mean of 0 and a standard deviation of 1. This allows neural networks with a larger number of layers, like the Transformer, to be trained easier

```
class EncoderLayer(nn.Module):
    def __init__(self,
                    hid_dim,
                    n_heads,
                    pf_dim,
                    dropout,
                    device):
        super().__init__()

        self.self_attn_layer_norm = nn.LayerNorm(hid_dim)
        self.ff_layer_norm = nn.LayerNorm(hid_dim)
        self.self_attention = MultiHeadAttentionLayer(hid_dim, n_heads, dropout, device)
        self.positionwise_feedforward = PositionwiseFeedforwardLayer(hid_dim,
                                                                        pf_dim,
                                                                        dropout)

        self.dropout = nn.Dropout(dropout)

    def forward(self, src, src_mask):

        #src = [batch size, src len, hid dim]
        #src_mask = [batch size, 1, 1, src len]

        #self attention
        _src, _ = self.self_attention(src, src, src, src_mask)

        #dropout, residual connection and layer norm
        src = self.self_attn_layer_norm(src + self.dropout(_src))

        #src = [batch size, src len, hid dim]

        #positionwise feedforward
        _src = self.positionwise_feedforward(src)

        #dropout, residual and layer norm
        src = self.ff_layer_norm(src + self.dropout(_src))
```

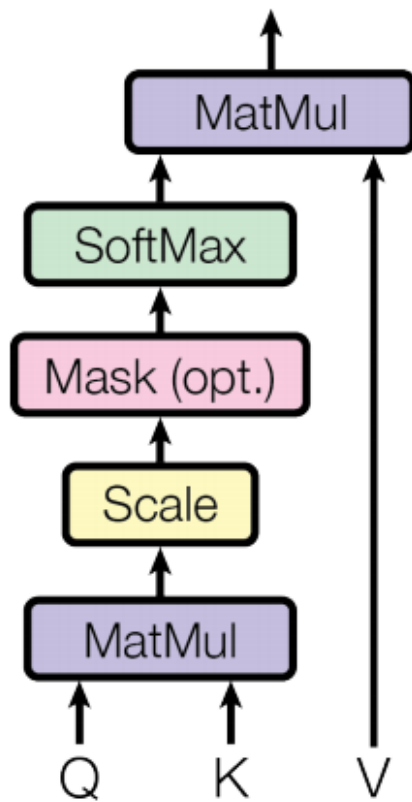
```
#src = [batch size, src len, hid dim]
```

```
return src
```

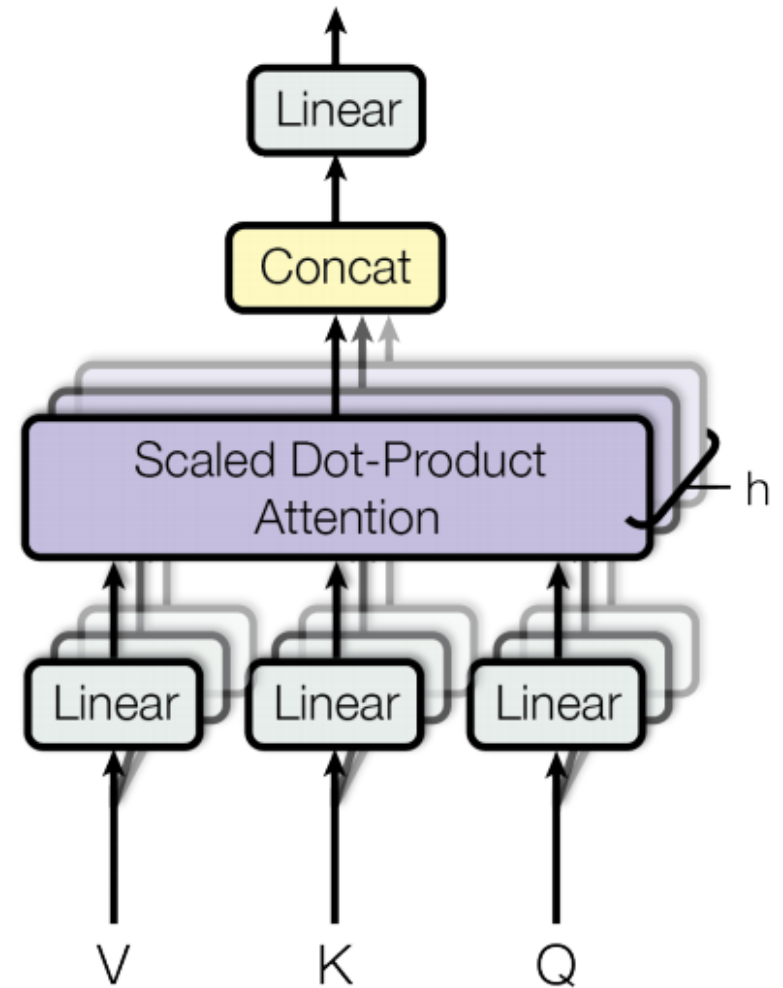
## ▼ Mutli Head Attention Layer

One of the key, novel concepts introduced by the Transformer paper is the *multi-head attention layer*.

### Scaled Dot-Product Attention



### Multi-Head Attention



Attention can be thought of as *queries*, *keys* and *values* - where the query is used with the key to get an attention vector (usually the output of a *softmax* operation and has all values between 0 and 1 which sum to 1) which is then used to get a weighted sum of the values.

The Transformer uses *scaled dot-product attention*, where the query and key are combined by taking the dot product between them, then applying the softmax operation and scaling by  $d_k$  before finally then multiplying by the value.  $d_k$  is the *head dimension*, `head_dim`, which we will shortly explain further.

$$\text{Attention}(Q, K, V) = \text{Softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

This is similar to standard *dot product attention* but is scaled by  $d_k$ , which the paper states is used to stop the results of the dot products growing large, causing gradients to become too small.

However, the scaled dot-product attention isn't simply applied to the queries, keys and values. Instead of doing a single attention application the queries, keys and values have their `hid_dim` split into  $h$  heads and the scaled dot-product attention is calculated over all heads in parallel. This means instead of paying attention to one concept per attention application, we pay attention to  $h$ . We then re-combine the heads into their `hid_dim` shape, thus each `hid_dim` is potentially paying attention to  $h$  different concepts.

$$\text{MultiHead}(Q, K, V) = \text{Concat}(\text{head}_1, \dots, \text{head}_h)W^O$$

$$\text{head}_i = \text{Attention}(QW_i^Q, KW_i^K, VW_i^V)$$

$W^O$  is the linear layer applied at the end of the multi-head attention layer, `fc`.  $W^Q$ ,  $W^K$ ,  $W^V$  are the linear layers `fc_q`, `fc_k` and `fc_v`.

Walking through the module, first we calculate  $QW^Q$ ,  $KW^K$  and  $VW^V$  with the linear layers, `fc_q`, `fc_k` and `fc_v`, to give us `Q`, `K` and `V`. Next, we split the `hid_dim` of the query, key and value into `n_heads` using `.view` and correctly permute them so they can be multiplied together. We then calculate the `energy` (the un-normalized attention) by multiplying `Q` and `K` together and scaling it by the square root of `head_dim`, which is calculated as `hid_dim // n_heads`. We then mask the energy so we do not pay attention over any elements of the sequence we shouldn't, then apply the softmax and dropout. We then apply the attention to the value heads, `V`, before combining the `n_heads` together. Finally, we multiply this  $W^O$ , represented by `fc_o`.

Note that in our implementation the lengths of the keys and values are always the same, thus when matrix multiplying the output of the `softmax`, `attention`, with `V` we will always have valid dimension sizes for matrix multiplication. This multiplication is carried out using `torch.matmul` which, when both tensors are >2-dimensional, does a batched matrix multiplication over the last two dimensions of each tensor.

This will be a **[query len, key len] x [value len, head dim]** batched matrix multiplication over the batch size and each head which provides the **[batch size, n heads, query len, head dim]** result.

One thing that looks strange at first is that dropout is applied directly to the attention. This means that our attention vector will most probably not sum to 1 and we may pay full attention to a token but the attention over that token is set to 0 by dropout. This is never explained, or even mentioned, in the paper however is used by the [official implementation](#) and every Transformer implementation since, [including BERT](#).

```
class MultiHeadAttentionLayer(nn.Module):
    def __init__(self, hid_dim, n_heads, dropout, device):
        super().__init__()

        assert hid_dim % n_heads == 0

        self.hid_dim = hid_dim
        self.n_heads = n_heads
        self.head_dim = hid_dim // n_heads

        self.fc_q = nn.Linear(hid_dim, hid_dim)
        self.fc_k = nn.Linear(hid_dim, hid_dim)
        self.fc_v = nn.Linear(hid_dim, hid_dim)

        self.fc_o = nn.Linear(hid_dim, hid_dim)

        self.dropout = nn.Dropout(dropout)

        self.scale = torch.sqrt(torch.FloatTensor([self.head_dim])).to(device)

    def forward(self, query, key, value, mask = None):

        batch_size = query.shape[0]

        #query = [batch size, query len, hid dim]
        #key = [batch size, key len, hid dim]
        #value = [batch size, value len, hid dim]

        Q = self.fc_q(query)
```

```

K = self.fc_k(key)
V = self.fc_v(value)

#Q = [batch size, query len, hid dim]
#K = [batch size, key len, hid dim]
#V = [batch size, value len, hid dim]

Q = Q.view(batch_size, -1, self.n_heads, self.head_dim).permute(0, 2, 1, 3)
K = K.view(batch_size, -1, self.n_heads, self.head_dim).permute(0, 2, 1, 3)
V = V.view(batch_size, -1, self.n_heads, self.head_dim).permute(0, 2, 1, 3)

#Q = [batch size, n heads, query len, head dim]
#K = [batch size, n heads, key len, head dim]
#V = [batch size, n heads, value len, head dim]

energy = torch.matmul(Q, K.permute(0, 1, 3, 2)) / self.scale

#energy = [batch size, n heads, query len, key len]

if mask is not None:
    energy = energy.masked_fill(mask == 0, -1e10)

attention = torch.softmax(energy, dim = -1)

#attention = [batch size, n heads, query len, key len]

x = torch.matmul(self.dropout(attention), V)

#x = [batch size, n heads, query len, head dim]

x = x.permute(0, 2, 1, 3).contiguous()

#x = [batch size, query len, n heads, head dim]

x = x.view(batch_size, -1, self.hid_dim)

#x = [batch size, query len, hid dim]

x = self.fc_o(x)

```

```
#x = [batch size, query len, hid dim]

return x, attention
```

## ▼ Position-wise Feedforward Layer

The other main block inside the encoder layer is the *position-wise feedforward layer*. This is relatively simple compared to the multi-head attention layer. The input is transformed from `hid_dim` to `pf_dim`, where `pf_dim` is usually a lot larger than `hid_dim`. The original Transformer used a `hid_dim` of 512 and a `pf_dim` of 2048. The ReLU activation function and dropout are applied before it is transformed back into a `hid_dim` representation.

Why is this used? Unfortunately, it is never explained in the paper.

BERT uses the [GELU](#) activation function, which can be used by simply switching `torch.relu` for `F.gelu`. Why did they use GELU? Again, it is never explained.

```
class PositionwiseFeedforwardLayer(nn.Module):
    def __init__(self, hid_dim, pf_dim, dropout):
        super().__init__()

        self.fc_1 = nn.Linear(hid_dim, pf_dim)
        self.fc_2 = nn.Linear(pf_dim, hid_dim)

        self.dropout = nn.Dropout(dropout)

    def forward(self, x):

        #x = [batch size, seq len, hid dim]

        x = self.dropout(torch.relu(self.fc_1(x)))

        #x = [batch size, seq len, pf dim]

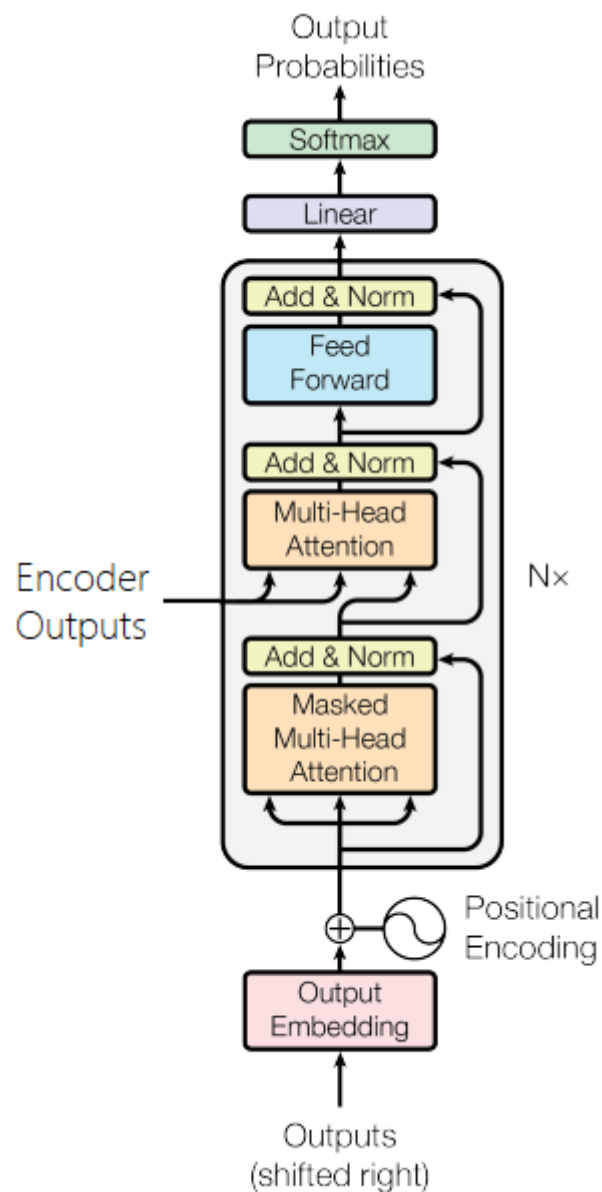
        x = self.fc_2(x)
```

```
#x = [batch size, seq len, hid dim]
```

```
return x
```

## ▼ Decoder

The objective of the decoder is to take the encoded representation of the source sentence,  $Z$ , and convert it into predicted tokens in the target sentence,  $\hat{Y}$ . We then compare  $\hat{Y}$  with the actual tokens in the target sentence,  $Y$ , to calculate our loss, which will be used to calculate the gradients of our parameters and then use our optimizer to update our weights in order to improve our predictions.



The decoder is similar to encoder, however it now has two multi-head attention layers. A *masked multi-head attention* layer over the target sequence, and a multi-head attention layer which uses the decoder representation as the query and the encoder representation as the key and value.

The decoder uses positional embeddings and combines - via an elementwise sum - them with the scaled embedded target tokens, followed by dropout. Again, our positional encodings have a "vocabulary" of 100, which means they can accept sequences up to 100 tokens long. This can





```

                                                                    pf_dim,
                                                                    dropout,
                                                                    device)
for _ in range(n_layers)])

self.fc_out = nn.Linear(hid_dim, output_dim)

self.dropout = nn.Dropout(dropout)

self.scale = torch.sqrt(torch.FloatTensor([hid_dim])).to(device)

def forward(self, trg, enc_src, trg_mask, src_mask):

    #trg = [batch size, trg len]
    #enc_src = [batch size, src len, hid dim]
    #trg_mask = [batch size, 1, trg len, trg len]
    #src_mask = [batch size, 1, 1, src len]

    batch_size = trg.shape[0]
    trg_len = trg.shape[1]

    pos = torch.arange(0, trg_len).unsqueeze(0).repeat(batch_size, 1).to(self.device)

    #pos = [batch size, trg len]

    trg = self.dropout((self.tok_embedding(trg) * self.scale) + self.pos_embedding(pos))

    #trg = [batch size, trg len, hid dim]

    for layer in self.layers:
        trg, attention = layer(trg, enc_src, trg_mask, src_mask)

    #trg = [batch size, trg len, hid dim]
    #attention = [batch size, n heads, trg len, src len]

    output = self.fc_out(trg)

    #output = [batch size, trg len, output dim]

```

```
return output, attention
```

## ▼ Decoder Layer

As mentioned previously, the decoder layer is similar to the encoder layer except that it now has two multi-head attention layers, `self_attention` and `encoder_attention`.

The first performs self-attention, as in the encoder, by using the decoder representation so far as the query, key and value. This is followed by dropout, residual connection and layer normalization. This `self_attention` layer uses the target sequence mask, `trg_mask`, in order to prevent the decoder from "cheating" by paying attention to tokens that are "ahead" of the one it is currently processing as it processes all tokens in the target sentence in parallel.

The second is how we actually feed the encoded source sentence, `enc_src`, into our decoder. In this multi-head attention layer the queries are the decoder representations and the keys and values are the encoder representations. Here, the source mask, `src_mask` is used to prevent the multi-head attention layer from attending to `<pad>` tokens within the source sentence. This is then followed by the dropout, residual connection and layer normalization layers.

Finally, we pass this through the position-wise feedforward layer and yet another sequence of dropout, residual connection and layer normalization.

The decoder layer isn't introducing any new concepts, just using the same set of layers as the encoder in a slightly different way.

```
class DecoderLayer(nn.Module):
    def __init__(self,
                    hid_dim,
                    n_heads,
                    pf_dim,
                    dropout,
                    device):
        super().__init__()

        self.self_attn_layer_norm = nn.LayerNorm(hid_dim)
        self.enc_attn_layer_norm = nn.LayerNorm(hid_dim)
        self.ff_layer_norm = nn.LayerNorm(hid_dim)
        self.self_attention = MultiHeadAttentionLayer(hid_dim, n_heads, dropout, device)
```

```

self.encoder_attention = MultiHeadAttentionLayer(hid_dim, n_heads, dropout, device)
self.positionwise_feedforward = PositionwiseFeedforwardLayer(hid_dim,

                                                                    pf_dim,
                                                                    dropout)

self.dropout = nn.Dropout(dropout)

def forward(self, trg, enc_src, trg_mask, src_mask):

    #trg = [batch size, trg len, hid dim]
    #enc_src = [batch size, src len, hid dim]
    #trg_mask = [batch size, 1, trg len, trg len]
    #src_mask = [batch size, 1, 1, src len]

    #self attention
    _trg, _ = self.self_attention(trg, trg, trg, trg_mask)

    #dropout, residual connection and layer norm
    trg = self.self_attn_layer_norm(trg + self.dropout(_trg))

    #trg = [batch size, trg len, hid dim]

    #encoder attention
    _trg, attention = self.encoder_attention(trg, enc_src, enc_src, src_mask)

    #dropout, residual connection and layer norm
    trg = self.enc_attn_layer_norm(trg + self.dropout(_trg))

    #trg = [batch size, trg len, hid dim]

    #positionwise feedforward
    _trg = self.positionwise_feedforward(trg)

    #dropout, residual and layer norm
    trg = self.ff_layer_norm(trg + self.dropout(_trg))

    #trg = [batch size, trg len, hid dim]
    #attention = [batch size, n heads, trg len, src len]

    return trg, attention

```

## Seq2Seq

Finally, we have the `Seq2Seq` module which encapsulates the encoder and decoder, as well as handling the creation of the masks.

The source mask is created by checking where the source sequence is not equal to a `<pad>` token. It is 1 where the token is not a `<pad>` token and 0 when it is. It is then unsqueezed so it can be correctly broadcast when applying the mask to the `energy`, which of shape ***[batch size, n heads, seq len, seq len]***.

The target mask is slightly more complicated. First, we create a mask for the `<pad>` tokens, as we did for the source mask. Next, we create a "subsequent" mask, `trg_sub_mask`, using `torch.tril`. This creates a diagonal matrix where the elements above the diagonal will be zero and the elements below the diagonal will be set to whatever the input tensor is. In this case, the input tensor will be a tensor filled with ones. So this means our `trg_sub_mask` will look something like this (for a target with 5 tokens):

1	0	0	0	0
1	1	0	0	0
1	1	1	0	0
1	1	1	1	0
1	1	1	1	1

This shows what each target token (row) is allowed to look at (column). The first target token has a mask of ***[1, 0, 0, 0, 0]*** which means it can only look at the first target token. The second target token has a mask of ***[1, 1, 0, 0, 0]*** which it means it can look at both the first and second target tokens.

The "subsequent" mask is then logically anded with the padding mask, this combines the two masks ensuring both the subsequent tokens and the padding tokens cannot be attended to. For example if the last two tokens were `<pad>` tokens the mask would look like:

1	0	0	0	0
1	1	0	0	0
1	1	1	0	0
1	1	1	0	0
1	1	1	0	0

After the masks are created, they are used with the encoder and decoder along with the source and target sentences to get our predicted target sentence, `output`, along with the decoder's attention over the source sequence.

```
class Seq2Seq(nn.Module):
    def __init__(self,
                  encoder,
                  decoder,
                  src_pad_idx,
                  trg_pad_idx,
                  device):
        super().__init__()

        self.encoder = encoder
        self.decoder = decoder
        self.src_pad_idx = src_pad_idx
        self.trg_pad_idx = trg_pad_idx
        self.device = device

    def make_src_mask(self, src):

        #src = [batch size, src len]

        src_mask = (src != self.src_pad_idx).unsqueeze(1).unsqueeze(2)

        #src_mask = [batch size, 1, 1, src len]

        return src_mask

    def make_trg_mask(self, trg):

        #trg = [batch size, trg len]

        trg_pad_mask = (trg == self.trg_pad_idx).unsqueeze(1).unsqueeze(2)

        #trg_pad_mask = [batch size, 1, 1, trg len]
```

```

trg_len = trg.shape[1]

trg_sub_mask = torch.tril(torch.ones((trg_len, trg_len), device = self.device)).bool()

#trg_sub_mask = [trg_len, trg_len]

trg_mask = trg_pad_mask & trg_sub_mask

#trg_mask = [batch_size, 1, trg_len, trg_len]

return trg_mask

def forward(self, src, trg):

    #src = [batch_size, src_len]
    #trg = [batch_size, trg_len]

    src_mask = self.make_src_mask(src)
    trg_mask = self.make_trg_mask(trg)

    #src_mask = [batch_size, 1, 1, src_len]
    #trg_mask = [batch_size, 1, trg_len, trg_len]

    enc_src = self.encoder(src, src_mask)

    #enc_src = [batch_size, src_len, hid_dim]

    output, attention = self.decoder(trg, enc_src, trg_mask, src_mask)

    #output = [batch_size, trg_len, output_dim]
    #attention = [batch_size, n_heads, trg_len, src_len]

    return output, attention

```

## ▼ Training the Seq2Seq Model

We can now define our encoder and decoders. This model is significantly smaller than Transformers used in research today, but is able to be

```
INPUT_DIM  = len(SRC.vocab)
OUTPUT_DIM  = len(TRG.vocab)
HID_DIM    = 256
ENC_LAYERS  = 3
DEC_LAYERS  = 3
ENC_HEADS   = 8
DEC_HEADS   = 8
ENC_PF_DIM  = 512
DEC_PF_DIM  = 512
ENC_DROPOUT = 0.1
DEC_DROPOUT = 0.1

enc = Encoder(INPUT_DIM,
               HID_DIM,
               ENC_LAYERS,
               ENC_HEADS,
               ENC_PF_DIM,
               ENC_DROPOUT,
               device)

dec = Decoder(OUTPUT_DIM,
               HID_DIM,
               DEC_LAYERS,
               DEC_HEADS,
               DEC_PF_DIM,
               DEC_DROPOUT,
               device)
```

Then, use them to define our whole sequence-to-sequence encapsulating model.

```
SRC_PAD_IDX = SRC.vocab.stoi[SRC.pad_token]
TRG_PAD_IDX = TRG.vocab.stoi[TRG.pad_token]

model = Seq2Seq(enc, dec, SRC_PAD_IDX, TRG_PAD_IDX, device).to(device)
```



We can check the number of parameters, noticing it is significantly less than the 37M for the convolutional sequence-to-sequence model.

```
def count_parameters(model):
    return sum(p.numel() for p in model.parameters() if p.requires_grad)

print(f'The model has {count_parameters(model):,} trainable parameters')

The model has 9,038,853 trainable parameters
```

The paper does not mention which weight initialization scheme was used, however Xavier uniform seems to be common amongst Transformer models, so we use it here.

```
def initialize_weights(m):
    if hasattr(m, 'weight') and m.weight.dim() > 1:
        nn.init.xavier_uniform_(m.weight.data)

model.apply(initialize_weights);
```

The optimizer used in the original Transformer paper uses Adam with a learning rate that has a "warm-up" and then a "cool-down" period. BERT and other Transformer models use Adam with a fixed learning rate, so we will implement that. Check [this](#) link for more details about the original Transformer's learning rate schedule.

Note that the learning rate needs to be lower than the default used by Adam or else learning is unstable.

```
LEARNING_RATE = 0.0005
```

```
optimizer = torch.optim.Adam(model.parameters(), lr = LEARNING_RATE)
```

Next, we define our loss function, making sure to ignore losses calculated over `<pad>` tokens.

```
criterion = nn.CrossEntropyLoss(ignore_index = TRG_PAD_IDX)
```

Then, we'll define our training loop. This is the exact same as the one used in the previous tutorial.

As we want our model to predict the `<eos>` token but not have it be an input into our model we simply slice the `<eos>` token off the end of the sequence. Thus:

$$\begin{aligned}\text{trg} &= [\text{sos}, x_1, x_2, x_3, \text{eos}] \\ \text{trg}[:-1] &= [\text{sos}, x_1, x_2, x_3]\end{aligned}$$

$x_i$  denotes actual target sequence element. We then feed this into the model to get a predicted sequence that should hopefully predict the `<eos>` token:

$$\text{output} = [y_1, y_2, y_3, \text{eos}]$$

$y_i$  denotes predicted target sequence element. We then calculate our loss using the original `trg` tensor with the `<sos>` token sliced off the front, leaving the `<eos>` token:

$$\begin{aligned}\text{output} &= [y_1, y_2, y_3, \text{eos}] \\ \text{trg}[1:] &= [x_1, x_2, x_3, \text{eos}]\end{aligned}$$

We then calculate our losses and update our parameters as is standard.

```
def train(model, iterator, optimizer, criterion, clip):

    model.train()

    epoch_loss = 0

    for i, batch in enumerate(iterator):

        src = batch.src
        trg = batch.trg

        optimizer.zero_grad()

        output, _ = model(src, trg[:, :-1])

        #output = [batch size, trg len - 1, output dim]
```

```

#trg = [batch size, trg len]

output_dim = output.shape[-1]

output = output.contiguous().view(-1, output_dim)
trg = trg[:,1:].contiguous().view(-1)

#output = [batch size * trg len - 1, output dim]
#trg = [batch size * trg len - 1]

loss = criterion(output, trg)

loss.backward()

torch.nn.utils.clip_grad_norm_(model.parameters(), clip)

optimizer.step()

epoch_loss += loss.item()

return epoch_loss / len(iterator)

```

The evaluation loop is the same as the training loop, just without the gradient calculations and parameter updates.

```

def evaluate(model, iterator, criterion):

    model.eval()

    epoch_loss = 0

    with torch.no_grad():

        for i, batch in enumerate(iterator):

            src = batch.src
            trg = batch.trg

```

```

output, _ = model(src, trg[:, :-1])

#output = [batch size, trg len - 1, output dim]
#trg = [batch size, trg len]

output_dim = output.shape[-1]

output = output.contiguous().view(-1, output_dim)
trg = trg[:, 1:].contiguous().view(-1)

#output = [batch size * trg len - 1, output dim]
#trg = [batch size * trg len - 1]

loss = criterion(output, trg)

epoch_loss += loss.item()

return epoch_loss / len(iterator)

```

We then define a small function that we can use to tell us how long an epoch takes.

```

def epoch_time(start_time, end_time):
    elapsed_time = end_time - start_time
    elapsed_mins = int(elapsed_time / 60)
    elapsed_secs = int(elapsed_time - (elapsed_mins * 60))
    return elapsed_mins, elapsed_secs

```

Finally, we train our actual model. This model is almost 3x faster than the convolutional sequence-to-sequence model and also achieves a lower validation perplexity!

```

N_EPOCHS = 10
CLIP = 1

```

```

best_valid_loss = float('inf')

```

```

for epoch in range(N_EPOCHS):

    start_time = time.time()

    train_loss = train(model, train_iterator, optimizer, criterion, CLIP)
    valid_loss = evaluate(model, valid_iterator, criterion)

    end_time = time.time()

    epoch_mins, epoch_secs = epoch_time(start_time, end_time)

    if valid_loss < best_valid_loss:
        best_valid_loss = valid_loss
        torch.save(model.state_dict(), 'tut6-model.pt')

    print(f'Epoch: {epoch+1:02} | Time: {epoch_mins}m {epoch_secs}s')
    print(f'\tTrain Loss: {train_loss:.3f} | Train PPL: {math.exp(train_loss):7.3f}')
    print(f'\t Val. Loss: {valid_loss:.3f} | Val. PPL: {math.exp(valid_loss):7.3f}')

```

```

Epoch: 01 | Time: 12m 29s
    Train Loss: 4.220 | Train PPL: 68.002
    Val. Loss: 3.009 | Val. PPL: 20.273
Epoch: 02 | Time: 12m 24s
    Train Loss: 2.806 | Train PPL: 16.548
    Val. Loss: 2.294 | Val. PPL: 9.914
Epoch: 03 | Time: 12m 20s
    Train Loss: 2.235 | Train PPL: 9.346
    Val. Loss: 1.983 | Val. PPL: 7.265
Epoch: 04 | Time: 12m 20s
    Train Loss: 1.887 | Train PPL: 6.603
    Val. Loss: 1.816 | Val. PPL: 6.149
Epoch: 05 | Time: 12m 38s
    Train Loss: 1.645 | Train PPL: 5.179
    Val. Loss: 1.713 | Val. PPL: 5.545
Epoch: 06 | Time: 12m 49s
    Train Loss: 1.458 | Train PPL: 4.298
    Val. Loss: 1.650 | Val. PPL: 5.208
Epoch: 07 | Time: 12m 51s
    Train Loss: 1.306 | Train PPL: 3.692
    Val. Loss: 1.624 | Val. PPL: 5.076
Epoch: 08 | Time: 12m 42s

```

```

Train Loss: 1.180 | Train PPL: 3.255
Val. Loss: 1.622 | Val. PPL: 5.063
Epoch: 09 | Time: 12m 47s
Train Loss: 1.068 | Train PPL: 2.911
Val. Loss: 1.639 | Val. PPL: 5.148
Epoch: 10 | Time: 12m 43s
Train Loss: 0.975 | Train PPL: 2.651
Val. Loss: 1.637 | Val. PPL: 5.140

```

We load our "best" parameters and manage to achieve a better test perplexity than all previous models.

```

model.load_state_dict(torch.load('tut6-model.pt'))

test_loss = evaluate(model, test_iterator, criterion)

print(f'| Test Loss: {test_loss:.3f} | Test PPL: {math.exp(test_loss):7.3f} |')

| Test Loss: 1.682 | Test PPL: 5.377 |

```

## ▼ Inference

Now we can can translations from our model with the `translate_sentence` function below.

The steps taken are:

- tokenize the source sentence if it has not been tokenized (is a string)
- append the `<sos>` and `<eos>` tokens
- numericalize the source sentence
- convert it to a tensor and add a batch dimension
- create the source sentence mask
- feed the source sentence and mask into the encoder
- create a list to hold the output sentence, initialized with an `<sos>` token
- while we have not hit a maximum length
  - convert the current output sentence prediction into a tensor with a batch dimension

- create a target sentence mask
- place the current output, encoder output and both masks into the decoder
- get next output token prediction from decoder along with attention
- add prediction to current output sentence prediction
- break if the prediction was an `<eos>` token
- convert the output sentence from indexes to tokens
- return the output sentence (with the `<sos>` token removed) and the attention from the last layer

```
def translate_sentence(sentence, src_field, trg_field, model, device, max_len = 50):
```

```
    model.eval()
```

```
    if isinstance(sentence, str):
```

```
        nlp = spacy.load('de_core_news_sm')
```

```
        tokens = [token.text.lower() for token in nlp(sentence)]
```

```
    else:
```

```
        tokens = [token.lower() for token in sentence]
```

```
    tokens = [src_field.init_token] + tokens + [src_field.eos_token]
```

```
    src_indexes = [src_field.vocab.stoi[token] for token in tokens]
```

```
    src_tensor = torch.LongTensor(src_indexes).unsqueeze(0).to(device)
```

```
    src_mask = model.make_src_mask(src_tensor)
```

```
    with torch.no_grad():
```

```
        enc_src = model.encoder(src_tensor, src_mask)
```

```
    trg_indexes = [trg_field.vocab.stoi[trg_field.init_token]]
```

```
    for i in range(max_len):
```

```
        trg_tensor = torch.LongTensor(trg_indexes).unsqueeze(0).to(device)
```

```
        trg_mask = model.make_trg_mask(trg_tensor)
```

```

with torch.no_grad():
    output, attention = model.decoder(trg_tensor, enc_src, trg_mask, src_mask)

    pred_token = output.argmax(2)[:,-1].item()

    trg_indexes.append(pred_token)

    if pred_token == trg_field.vocab.stoi[trg_field.eos_token]:
        break

trg_tokens = [trg_field.vocab.itos[i] for i in trg_indexes]

return trg_tokens[1:], attention

```

We'll now define a function that displays the attention over the source sentence for each step of the decoding. As this model has 8 heads our model we can view the attention for each of the heads.

```

def display_attention(sentence, translation, attention, n_heads = 8, n_rows = 4, n_cols = 2):

    assert n_rows * n_cols == n_heads

    fig = plt.figure(figsize=(15,25))

    for i in range(n_heads):

        ax = fig.add_subplot(n_rows, n_cols, i+1)

        _attention = attention.squeeze(0)[i].cpu().detach().numpy()

        cax = ax.matshow(_attention, cmap='bone')

        ax.tick_params(labelsize=12)
        ax.set_xticklabels(['']+['<sos>']+ [t.lower() for t in sentence]+['<eos>'],
                           rotation=45)
        ax.set_yticklabels(['']+translation)

```



```

ax.xaxis.set_major_locator(ticker.MultipleLocator(1))
ax.yaxis.set_major_locator(ticker.MultipleLocator(1))

plt.show()
plt.close()

```

First, we'll get an example from the training set.

```

example_idx = 8

src = vars(train_data.examples[example_idx])['src']
trg = vars(train_data.examples[example_idx])['trg']

print(f'src = {src}')
print(f'trg = {trg}')

src = ['eine', 'frau', 'mit', 'einer', 'großen', 'geldbörse', 'geht', 'an', 'einem', 'tor', 'vorbei', '.']
trg = ['a', 'woman', 'with', 'a', 'large', 'purse', 'is', 'walking', 'by', 'a', 'gate', '.']

```

Our translation looks pretty good, although our model changes *is walking by* to *walks by*. The meaning is still the same.

```

translation, attention = translate_sentence(src, SRC, TRG, model, device)

print(f'predicted trg = {translation}')

predicted_trg = ['a', 'woman', 'with', 'a', 'large', 'purse', 'walks', 'past', 'a', 'gate', '.', '<eos>']

```

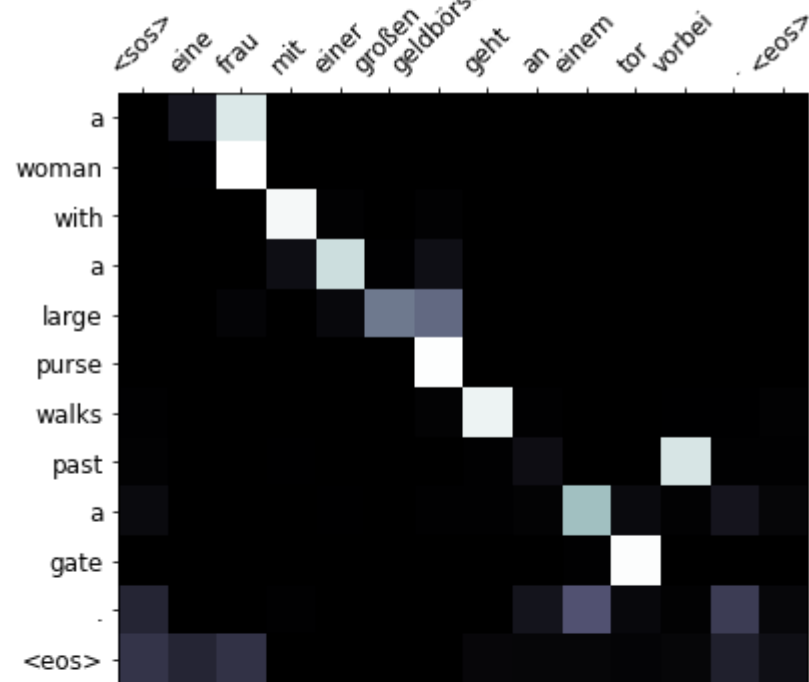
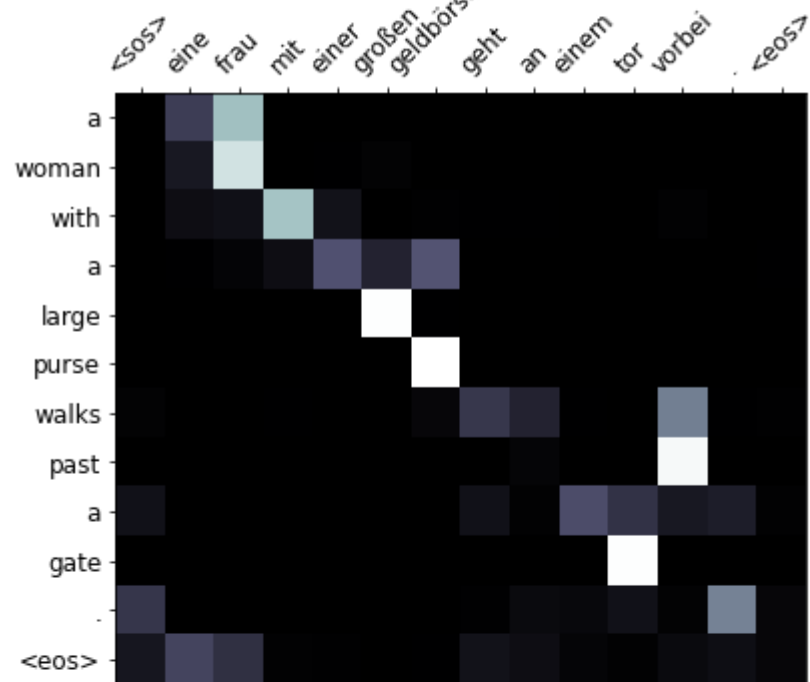
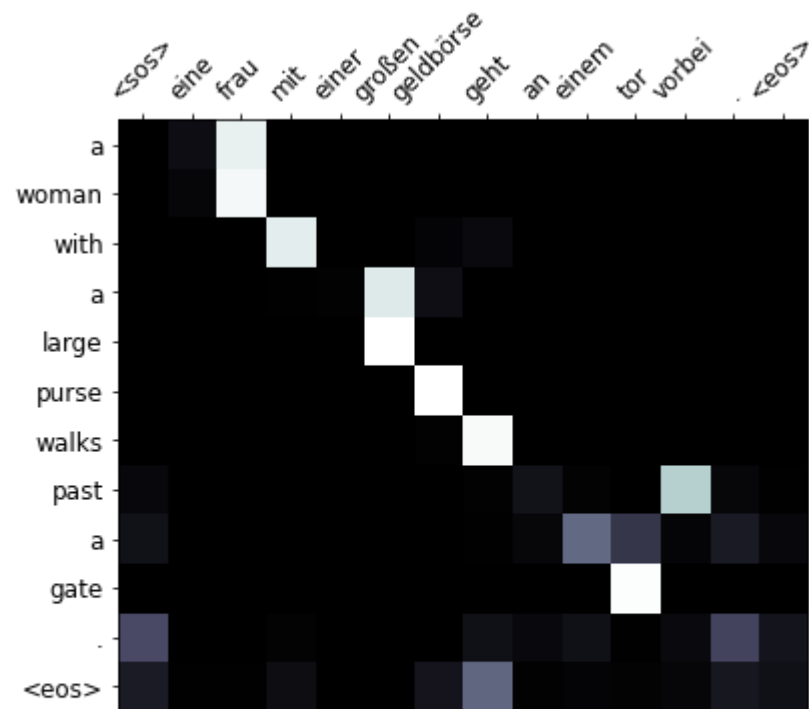
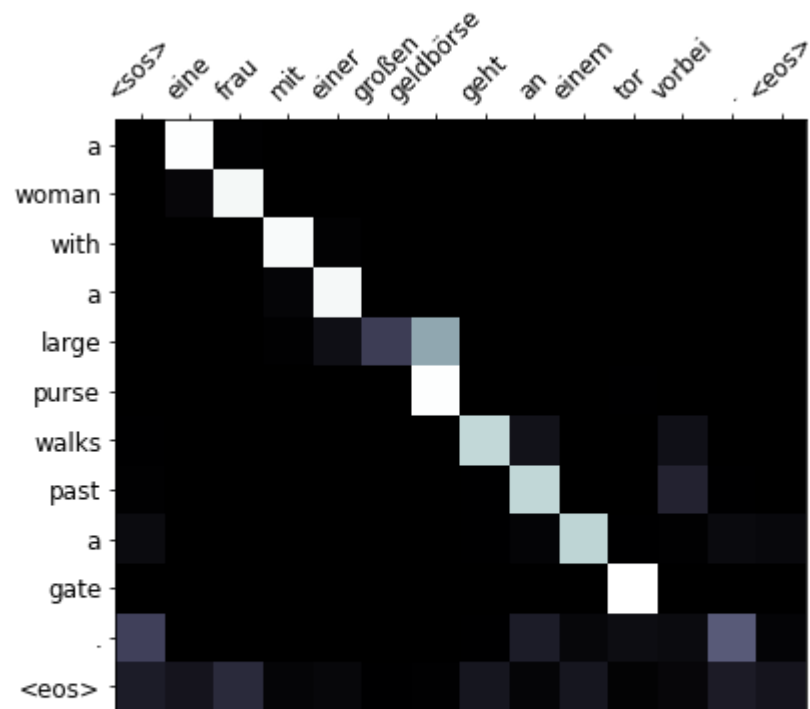
We can see the attention from each head below. Each is certainly different, but it's difficult (perhaps impossible) to reason about what head has actually learned to pay attention to. Some heads pay full attention to "eine" when translating "a", some don't at all, and some do a little. They all seem to follow the similar "downward staircase" pattern and the attention when outputting the last two tokens is equally spread over the final two tokens in the input sentence.

```

display_attention(src, translation, attention)

```





Next, let's get an example the model has not been trained on from the validation set.



```
example_idx = 6
```

```
src = vars(valid_data.examples[example_idx])['src']
```

```
trg = vars(valid_data.examples[example_idx])['trg']
```

```
print(f'src = {src}')
```

```
print(f'trg = {trg}')
```

```
src = ['ein', 'brauner', 'hund', 'rennt', 'dem', 'schwarzen', 'hund', 'hinterher', '.']
```

```
trg = ['a', 'brown', 'dog', 'is', 'running', 'after', 'the', 'black', 'dog', '.']
```



The model translates it by switching *is running* to just *runs*, but it is an acceptable swap.



```
translation, attention = translate_sentence(src, SRC, TRG, model, device)
```

```
print(f'predicted trg = {translation}')
```

```
predicted_trg = ['a', 'brown', 'dog', 'runs', 'after', 'the', 'black', 'dog', '.', '<eos>']
```

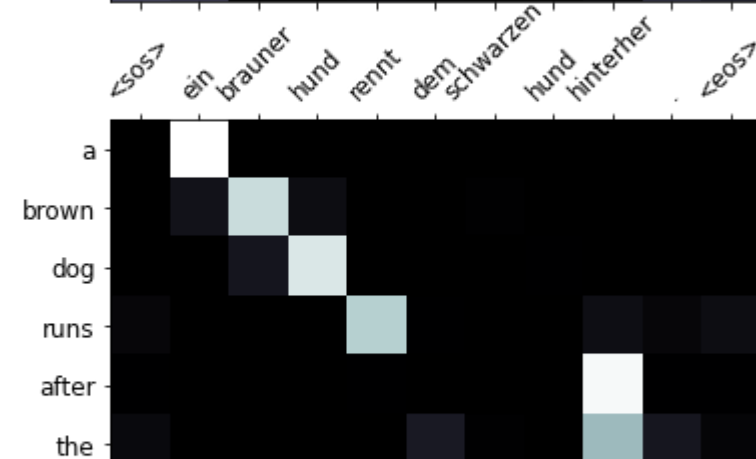
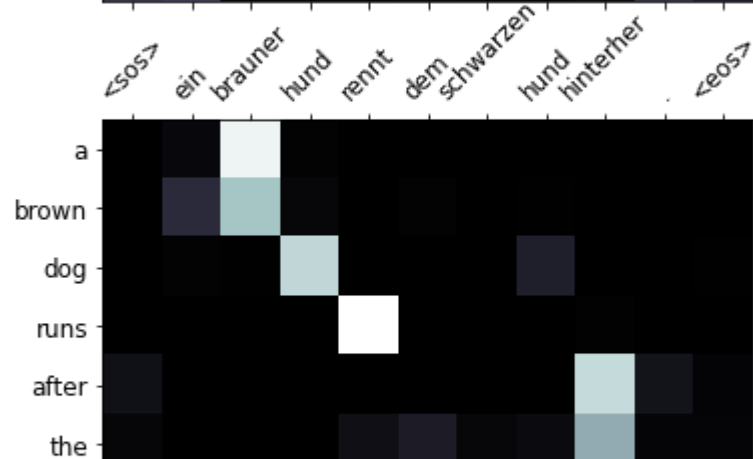
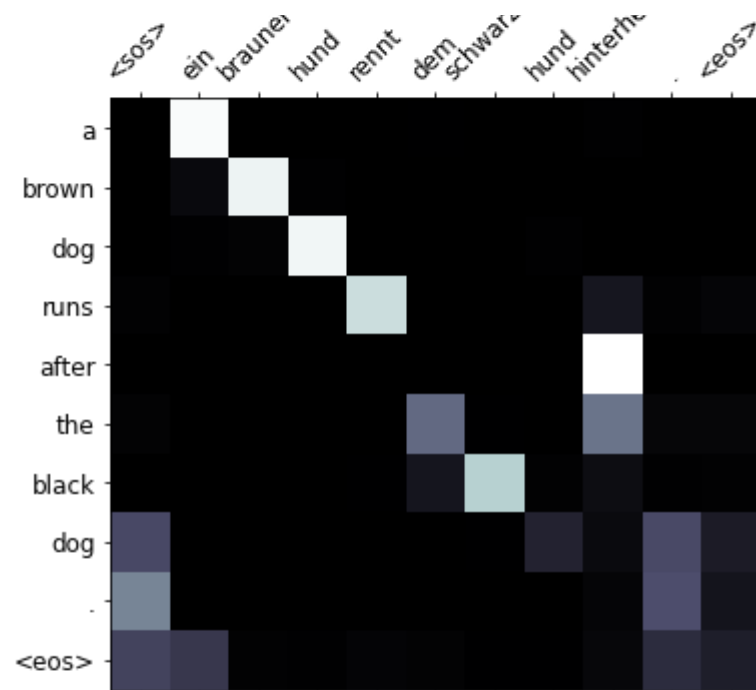
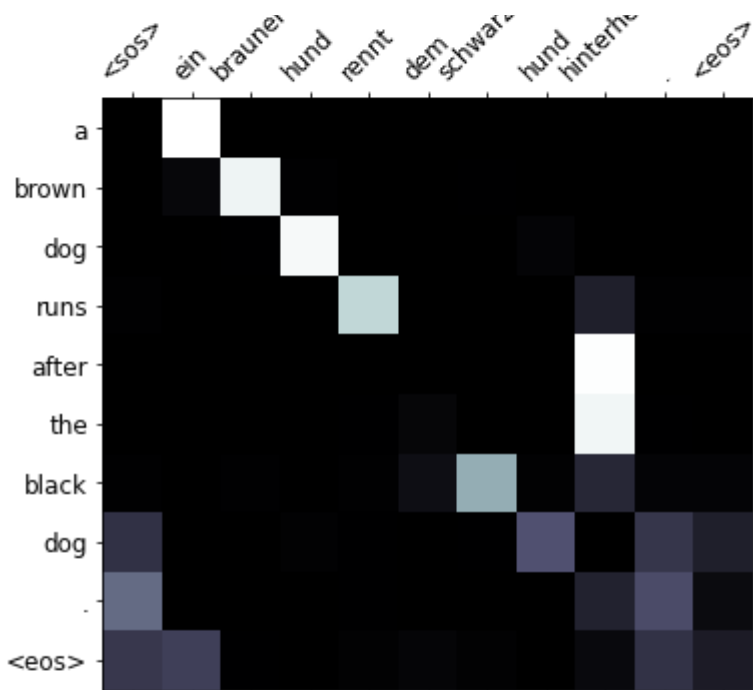


Again, some heads pay full attention to "ein" whilst some pay no attention to it. Again, most of the heads seem to spread their attention over both the period and `<eos>` tokens in the source sentence when outputting the period and `<eos>` sentence in the predicted target sentence, though some seem to pay attention to tokens from near the start of the sentence.



```
display_attention(src, translation, attention)
```





Finally, we'll look at an example from the test data.



example\_idx = 90

```
src = vars(test_data.examples[example_idx])['src']
trg = vars(test_data.examples[example_idx])['trg']
```

```
print(f'src  =  {src}')
print(f'trg  =  {trg}')
```

```
src = ['eine', 'schwarz', 'gekleidete', 'frau', 'mit', 'rosa', 'haaren', 'spricht', 'mit', 'einem', 'mann', '.']
trg = ['a', 'woman', 'with', 'pink', 'hair', 'dressed', 'in', 'black', 'talks', 'to', 'a', 'man', '.']
```

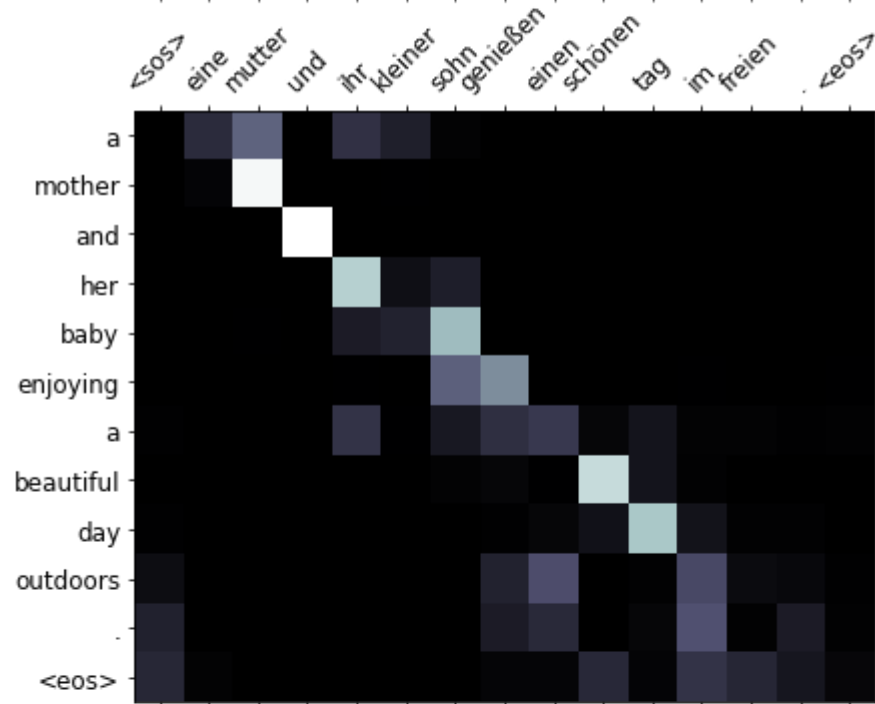
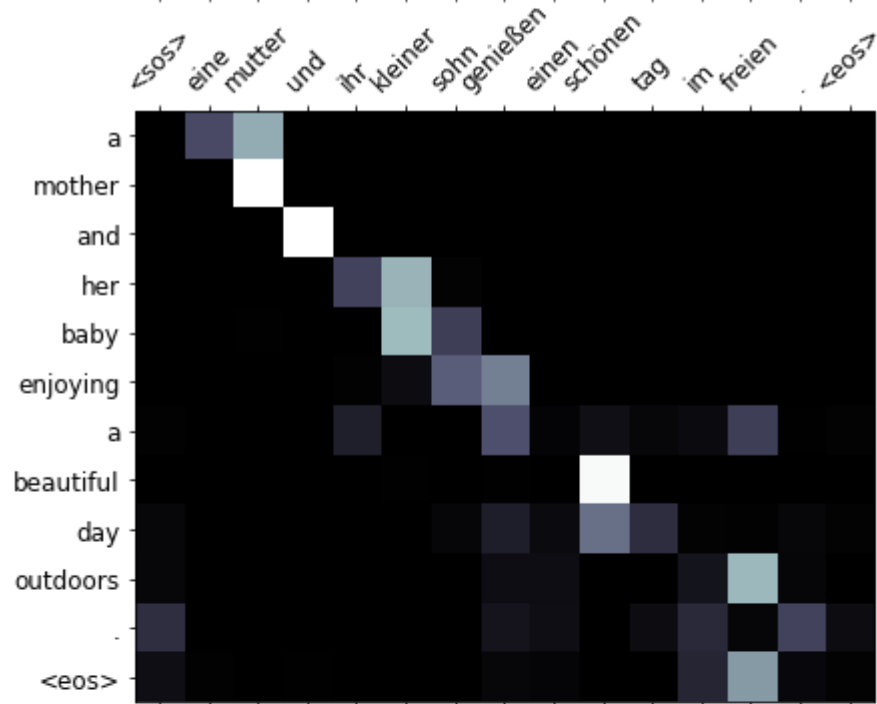
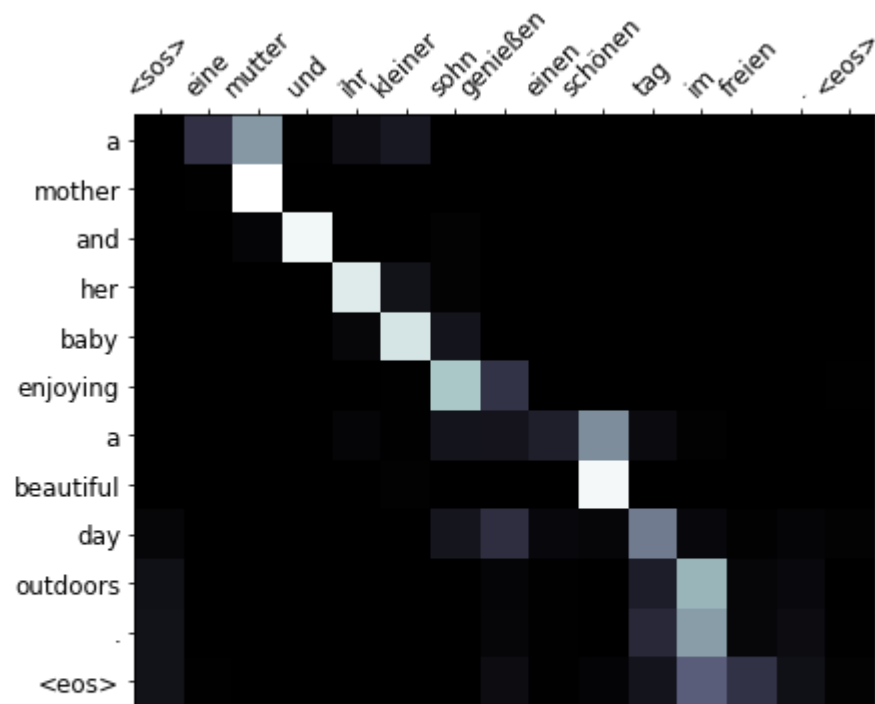
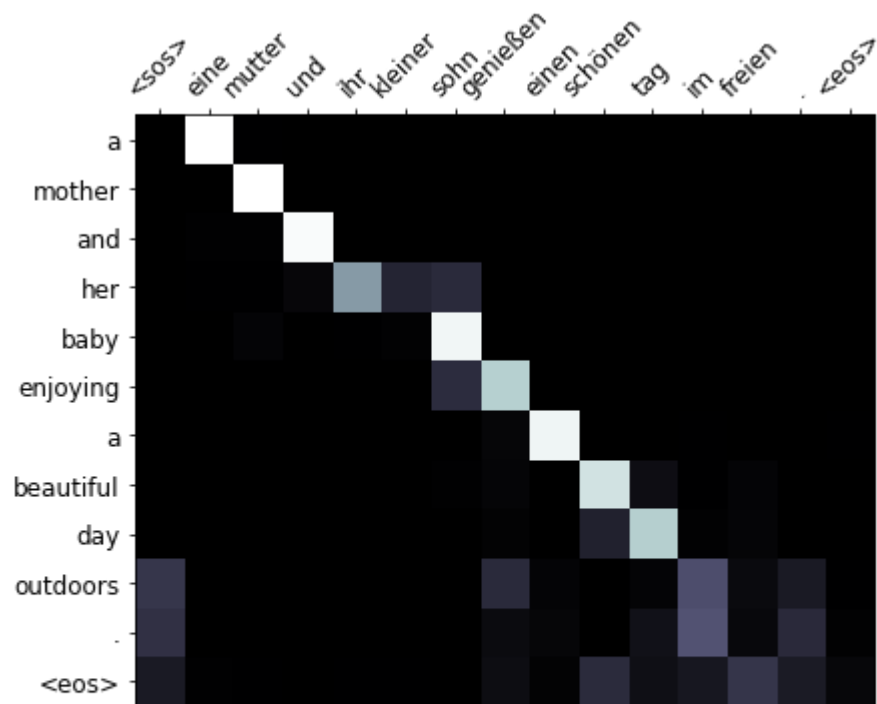
## A perfect translation!

```
translation, attention = translate_sentence(src, SRC, TRG, model, device)
```

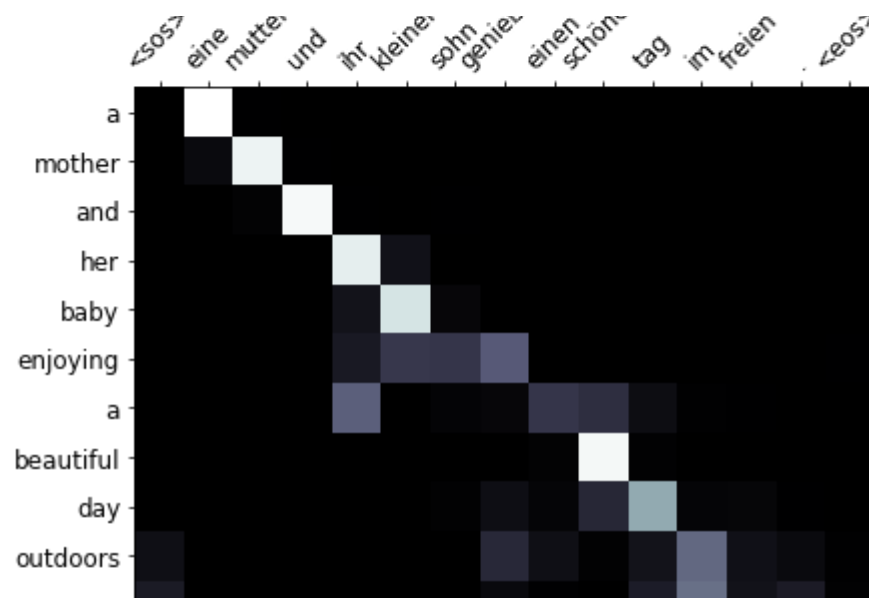
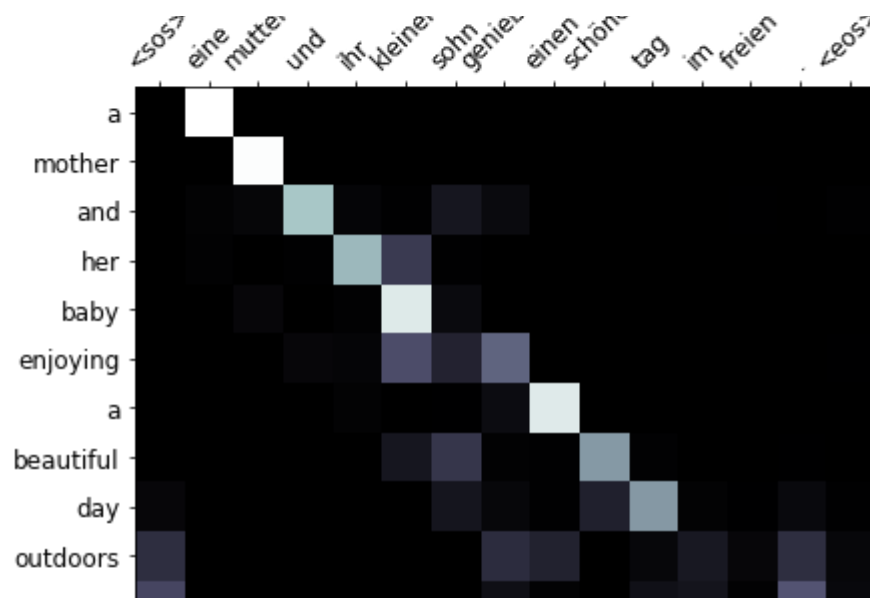
```
print(f'predicted trg = {translation}')
```

```
predicted_trg = ['people', 'walking', 'on', 'a', 'sidewalk', 'next', 'to', 'stores', '.', '<eos>']
```

```
display_attention(src, translation, attention)
```







## ▼ BLEU

Finally we calculate the BLEU score for the Transformer.



```
from torchtext.data.metrics import bleu_score
```

```
def calculate_bleu(data, src_field, trg_field, model, device, max_len = 50):
```

```
    trgs = []
```

```
    pred_trgs = []
```

```
    for datum in data:
```

```
        src = vars(datum)['src']
```

```
        trg = vars(datum)['trg']
```

```
        pred_trg, _ = translate_sentence(src, src_field, trg_field, model, device, max_len)
```

```
        #cut off <eos> token
```

```
        pred_trg = pred_trg[:-1]
```

```

        pred_trgs.append(pred_trg)
        trgs.append([trg])

    return bleu_score(pred_trgs, trgs)

```

We get a BLEU score of 36.52, which beats the ~34 of the convolutional sequence-to-sequence model and ~28 of the attention based RNN model. All this whilst having the least amount of parameters and the fastest training time!

```

bleu_score = calculate_bleu(test_data, SRC, TRG, model, device)

print(f'BLEU score = {bleu_score*100:.2f}')

BLEU score = 35.53

```

Congratulations for finishing these tutorials! I hope you've found them useful.

If you find any mistakes or want to ask any questions about any of the code or explanations used, feel free to submit a GitHub issue and I will try to correct it ASAP.

## ▼ Appendix

The `calculate_bleu` function above is unoptimized. Below is a significantly faster, vectorized version of it that should be used if needed. Credit for the implementation goes to [@azadyasar](#).

```

def translate_sentence_vectorized(src_tensor, src_field, trg_field, model, device, max_len=50):
    assert isinstance(src_tensor, torch.Tensor)

    model.eval()
    src_mask = model.make_src_mask(src_tensor)

    with torch.no_grad():
        enc_src = model.encoder(src_tensor, src_mask)
    # enc_src = [batch_sz, src_len, hid_dim]

```

```

trg_indexes = [[trg_field.vocab.stoi[trg_field.init_token]] for _ in range(len(src_tensor))]
# Even though some examples might have been completed by producing a <eos> token
# we still need to feed them through the model because other are not yet finished
# and all examples act as a batch. Once every single sentence prediction encounters
# <eos> token, then we can stop predicting.
translations_done = [0] * len(src_tensor)
for i in range(max_len):
    trg_tensor = torch.LongTensor(trg_indexes).to(device)
    trg_mask = model.make_trg_mask(trg_tensor)
    with torch.no_grad():
        output, attention = model.decoder(trg_tensor, enc_src, trg_mask, src_mask)
    pred_tokens = output.argmax(2)[:,-1]
    for i, pred_token_i in enumerate(pred_tokens):
        trg_indexes[i].append(pred_token_i)
        if pred_token_i == trg_field.vocab.stoi[trg_field.eos_token]:
            translations_done[i] = 1
    if all(translations_done):
        break

# Iterate through each predicted example one by one;
# Cut-off the portion including the after the <eos> token
pred_sentences = []
for trg_sentence in trg_indexes:
    pred_sentence = []
    for i in range(1, len(trg_sentence)):
        if trg_sentence[i] == trg_field.vocab.stoi[trg_field.eos_token]:
            break
        pred_sentence.append(trg_field.vocab.itos[trg_sentence[i]])
    pred_sentences.append(pred_sentence)

return pred_sentences, attention

```

```

from torchtext.data.metrics import bleu_score

```

```

def calculate_bleu_alt(iterator, src_field, trg_field, model, device, max_len = 50):
    trgs = []
    pred_trgs = []
    with torch.no_grad():

```