

```
1 from IPython.display import Image
2 import numpy as np
3 import torch
4 import torch.nn as nn
5 import math
6 import copy
7 import torch.nn.functional as F
8 from torch.nn.utils.rnn import pack_padded_sequence,
9 print("Torch version:", torch.__version__)
```

```
Torch version: 1.9.0+cu102
```

Attention Is All You Need (Transformer)

Core Idea of the Paper

Problem

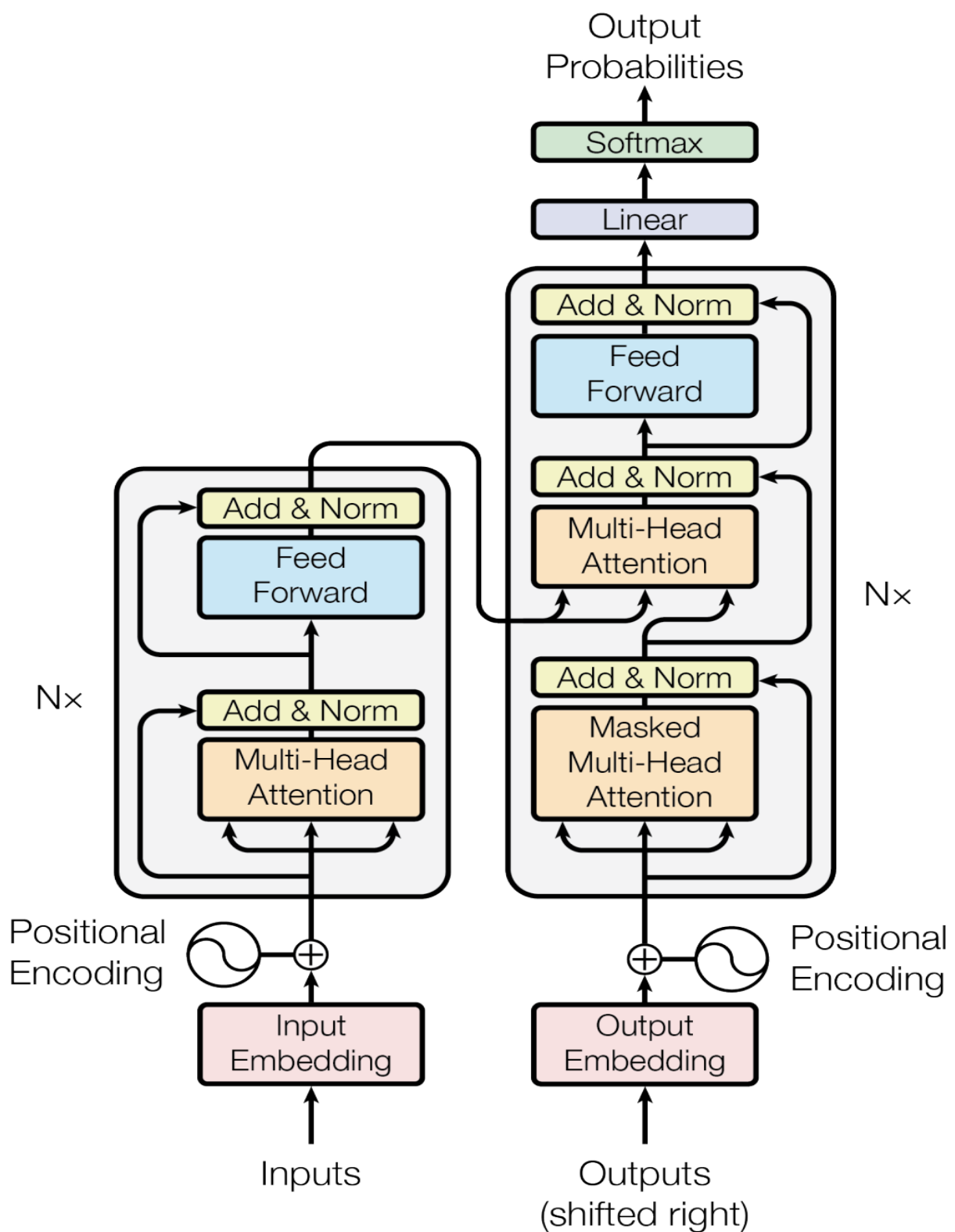
- In sequence-to-sequence problems such as the neural machine translation, the first proposals were based on the use of RNNs in an encoder-decoder architecture.
- The best performing models also connect the encoder and decoder through an attention mechanism.
- These architectures have a great limitation when working with long sequences
- In the encoder, the hidden state in every step is associated with a certain word in the input sentence, usually one of the most recent. Therefore, if the decoder only accesses the last hidden state of the decoder, it will lose relevant information about the first elements of the sequence.

Solution

- This paper propose a new simple network architecture, the Transformer, based solely on attention mechanisms to draw global dependencies between input and output.
- Instead of paying attention to the last state of the encoder as is usually done with RNNs, in each step of the decoder we look at all the states of the encoder, being able to access information about all the elements of the input sequence.
- The total computational complexity per layer
- The amount of computation that can be parallelized, as measured by the minimum number of sequential operations required.
- The path length between long-range dependencies in the network. Learning long-range dependencies is a key challenge in many sequence transduction tasks.
- List item

▼ ***Model Architecture***

The Transformer - model architecture.



▼ *Input Embedding and Softmax*

With this layer, we convert the input tokens and output tokens to vectors of dimension d_{model} using a learned embedding. We also use the usual learned linear transformation and softmax function to convert the decoder output to predicted next-token probabilities. In the embedding layers, we multiply those weights by $\sqrt{d_{\text{model}}}$.

```
1 class Embeddings(nn.Module):
2     def __init__(self, d_model, vocab):
```

```

3     super(Embeddings, self).__init__()
4     self.lut = nn.Embedding(vocab, d_model)
5     self.d_model = d_model
6
7     def forward(self, x):
8         return self.lut(x) * math.sqrt(self.d_model)

```

▼ *Positional Encoding*

In order for the model to make use of the order of the sequence, we must inject some information about the relative or absolute position of the tokens in the senquence. In this work, we use sine and cosine functions of different frequencies:

$$PE_{(pos,2i)} = \sin(pos/10000^{2i/d_{\text{model}}})$$

$$PE_{(pos,2i+1)} = \cos(pos/10000^{2i/d_{\text{model}}})$$

where pos is the position and i is the dimension. That is, each dimension of the positional encoding corresponds to a sinusoid. The wavelengths form a geometric progression from 2π to $10000 \cdot 2\pi$. We chose this function because we hypothesized **it would allow the model to easily learn to attend by relative positions, since for any fixed offset k , PE_{pos+k} can be represented as a linear function of PE_{pos} .**

```

1 class PositionalEncoding(nn.Module):
2     "Implement the PE function."
3     def __init__(self, d_model, dropout, max_len=5000):
4         super(PositionalEncoding, self).__init__()
5         self.dropout = nn.Dropout(p=dropout)
6
7         PE = torch.zeros(max_len, d_model)
8         position = torch.arange(0, max_len).unsqueeze(1)
9         div_term = torch.exp(torch.arange(0, d_model, 2).float() * (-math.log(2) / (d_model / 2)))
10        PE[:, 0::2] = torch.sin(position * div_term)
11        PE[:, 1::2] = torch.cos(position * div_term)
12        PE = PE.unsqueeze(0)
13        self.register_buffer('PE', PE)

```

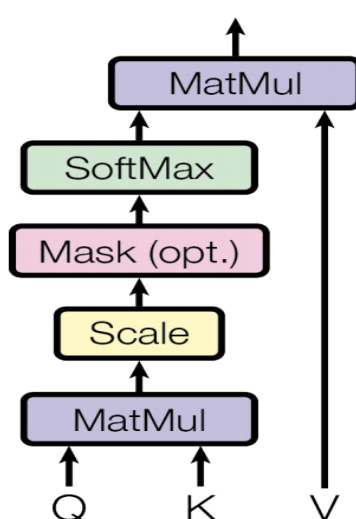
```

13         self.register_buffer('PE', torch.zeros(1, 1))
14
15     def forward(self, x):
16         x = x + self.PE[:, :x.size(1)]
17         return self.dropout(x)

```

▼ *Scaled Dot-Product Attention*

Scaled Dot-Product Attention



An attention function can be described as mapping a query and a set of key-value pairs to an output, where the query, keys, values, and output are all vectors. The output is computed as a weighted sum of the values, where the weight assigned to each value is computed by a compatibility function of the query with the corresponding key.

In practice, we compute the attention function on a set of queries simultaneously, packed together into a matrix Q . The keys and values are also packed together into matrices K and V . We compute the matrix of outputs as:

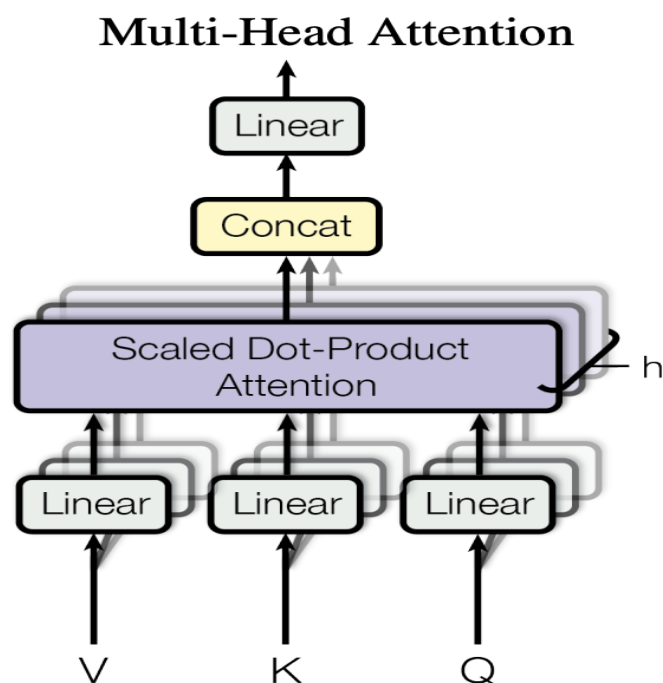
$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

```

1 class ScaledDotProductAttention(nn.Module):
2     ''' Scaled Dot-Product Attention '''
3
4     def __init__(self, d_k, dropout=0.1):
5         super().__init__()
6         self.d_k = d_k
7         self.dropout = nn.Dropout(dropout)
8
9     def forward(self, q, k, v, mask=None):
10
11         attn = torch.matmul(q / math.sqrt(self.d_k), k.transpose(-2, -1))
12
13         if mask is not None:
14             attn = attn.masked_fill(mask == 0, -1e9)
15
16         attn = self.dropout(F.softmax(attn, dim=-1))
17         output = torch.matmul(attn, v)
18
19         return output, attn

```

▼ Multi-Head Attention



Multi-head attention allows the model to jointly attend to information from different representation subspaces at different positions. With a single attention head, averaging inhibits this.

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

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

Where the projections are parameter matrices

$$W_i^Q \in \mathbb{R}^{d_{\text{model}} \times d_k}, W_i^K \in \mathbb{R}^{d_{\text{model}} \times d_k}, W_i^V \in \mathbb{R}^{d_{\text{model}} \times d_v} \text{ and } W^O \in \mathbb{R}^{hd_v \times d_{\text{model}}}.$$

```

1 class MultiHeadAttention(nn.Module):
2     ''' Multi-Head Attention module '''
3
4     def __init__(self, n_head, d_model, d_k, d_v, dropout):
5         super().__init__()
6
7         self.n_head = n_head
8         self.d_k = d_k
9         self.d_v = d_v
10
11         self.w_q = nn.Linear(d_model, n_head * d_k)
12         self.w_k = nn.Linear(d_model, n_head * d_k)
13         self.w_v = nn.Linear(d_model, n_head * d_v)
14         self.FC = nn.Linear(n_head * d_v, d_model)
15
16         self.attention = ScaledDotProductAttention(d_k)
17
18         self.dropout = nn.Dropout(dropout)
19         self.layerNorm = nn.LayerNorm(d_model, eps=1e-6)
20
21
22     def forward(self, Q, K, V, mask=None):
23
24         residual = Q
25
26         # Pass through the pre-attention projection: b x lq x n x dv
27         # Separate different heads: b x lq x n x dv
28         Q = self.w_q(Q).view(Q.size(0), Q.size(1), self.n_head, self.d_k)
29         K = self.w_k(K).view(K.size(0), K.size(1), self.n_head, self.d_k)

```

```

30     v = self.w_v(V).view(V.size(0), V.size(1), self.nhead)
31
32     # Transpose for attention dot product: b x n x ld
33     Q, K, V = Q.transpose(1, 2), K.transpose(1, 2), V
34
35     if mask is not None:
36         mask = mask.unsqueeze(1)    # For head axis broadcasting
37
38     Q, attn = self.attention(Q, K, V, mask=mask)
39
40     # Transpose to move the head dimension back: b x n x ld
41     # Combine the last two dimensions to concatenate
42     Q = Q.transpose(1, 2).contiguous().view(Q.size(0), -1)
43     Q = self.FC(Q)
44     Q = self.dropout(Q)
45
46     return Q, attn

```

▼ ***Position-wise Feed-Forward Networks***

In addition to attention sub-layers, each of the layers in our encoder and decoder contains a fully connected feed-forward network, which is applied to each position separately and identically. This consists of two linear transformations with a ReLU activation in between.

$$\text{FFN}(x) = \max(0, xW_1 + b_1)W_2 + b_2$$

```

1 class PositionwiseFeedForward(nn.Module):
2     ''' A two-feed-forward-layer module '''
3
4     def __init__(self, d_model, d_hid, dropout=0.1):
5         super().__init__()
6         self.W_1 = nn.Linear(d_model, d_hid)
7         self.W_2 = nn.Linear(d_hid, d_model)
8         self.dropout = nn.Dropout(dropout)
9
10    def forward(self, x):
11        x = self.W_1(x)
12        x = F.relu(x)

```



```

12         x = self.W_2(x)
13         x = self.dropout(x)
14
15     return x
16

```

▼ *Built our Encoder*

Since the core encoder contains N=6 encoder sublayer, we will define an EncoderLayer class for each sublayer.

```

1 class EncoderLayer(nn.Module):
2     ''' Compose with two layers '''
3
4     def __init__(self, d_model, d_hid, n_head, d_k, d_v):
5         super(EncoderLayer, self).__init__()
6         self.MHA = MultiHeadAttention(n_head, d_model, d_k, d_v)
7         self.layerNorm = nn.LayerNorm(d_model, eps=1e-6)
8         self.FFN = PositionwiseFeedForward(d_model, d_hid)
9
10    def forward(self, encoder_input, mask=None):
11        res1 = encoder_input
12        output1, MHA = self.MHA(encoder_input, encoder_input, mask)
13        output1 = res1 + self.layerNorm(output1)
14        res2 = output1
15        output2 = self.FFN(output1)
16        output2 = res2 + self.layerNorm(output2)
17        encoder_output = output2
18        return encoder_output, MHA

```

```

1 class Encoder(nn.Module):
2     "Core encoder is a stack of N layers"
3     def __init__(self, layer, N):
4         super(Encoder, self).__init__()
5         self.layers = nn.ModuleList([copy.deepcopy(layer) for _ in range(N)])
6         self.layerNorm = nn.LayerNorm(d_model, eps=1e-6)
7
8     def forward(self, x, mask):
9         "Pass the input (and mask) through each layer in the stack"
10        for layer in self.layers:
11            x, MHA = layer(x, mask)

```

```

11         x = layer(x, mask)
12     return self.layerNorm(x)

```

▼ Built our Decoder

Since the core decoder contains N=6 decoder sublayer, we will define an DecoderLayer class for each sublayer.

```

1 class DecoderLayer(nn.Module):
2     ''' Compose with three layers '''
3
4     def __init__(self, d_model, d_hid, n_head, d_k, d_v):
5         super(DecoderLayer, self).__init__()
6         self.MMHA = MultiHeadAttention(n_head, d_model, d_k, d_v)
7         self.MHA = MultiHeadAttention(n_head, d_model, d_k, d_v)
8         self.layerNorm = nn.LayerNorm(d_model, eps=1e-6)
9         self.FFN = PositionwiseFeedForward(d_model, d_hid)
10
11     def forward(self, decoder_input, encoder_output, src_mask):
12         res1 = decoder_input
13         output1, MMHA = self.MMHA(decoder_input, decoder_input, src_mask)
14         output1 = res1 + self.layerNorm(output1)
15         res2 = output1
16         output2, MHA = self.MHA(output1, encoder_output, None)
17         output2 = res2 + self.layerNorm(output2)
18         res3 = output2
19         output3 = self.FFN(output2)
20         decoder_output = res3 + output3
21
22     return decoder_output, MMHA, MHA

```

```

1 class Decoder(nn.Module):
2     "Generic N layer decoder with masking."
3     def __init__(self, layer, N):
4         super(Decoder, self).__init__()
5         self.layers = nn.ModuleList([copy.deepcopy(layer) for _ in range(N)])
6
7     def forward(self, x, memory, src_mask, tgt_mask):
8         for layer in self.layers:
9             x = layer(x, memory, src_mask, tgt_mask)
10
11     return x

```

```
10         return x
```

▼ Transformer

```
1 def get_pad_mask(seq, pad_idx):
2     return (seq != pad_idx).unsqueeze(-2)
3
4
5 def get_subsequent_mask(seq):
6     ''' For masking out the subsequent info. '''
7     sz_b, len_s = seq.size()
8     subsequent_mask = (1 - torch.triu(
9         torch.ones((1, len_s, len_s), device=seq.device),
10        return subsequent_mask
```

```
1 class Transformer(nn.Module):
2     ''' A sequence to sequence model with attention mechanism '''
3
4     def __init__(self, n_src_vocab, n_trg_vocab, src_pad_idx,
5                   d_word_vec=512, d_model=512, d_inner=2048,
6                   n_layers=6, n_head=8, d_k=64, d_v=64, dropout=0.1):
7
8
9         super().__init__()
10
11         self.src_pad_idx, self.trg_pad_idx = src_pad_idx, trg_pad_idx
12         self.d_model = d_model
13
14         self.encoder = Encoder(
15             n_src_vocab=n_src_vocab, n_position=n_position,
16             d_word_vec=d_word_vec, d_model=d_model, d_inner=d_inner,
17             n_layers=n_layers, n_head=n_head, d_k=d_k, d_v=d_v,
18             pad_idx=src_pad_idx, dropout=dropout, scalar=1.0)
19
20         self.decoder = Decoder(
21             n_trg_vocab=n_trg_vocab, n_position=n_position,
22             d_word_vec=d_word_vec, d_model=d_model, d_inner=d_inner,
23             n_layers=n_layers, n_head=n_head, d_k=d_k, d_v=d_v,
24             pad_idx=trg_pad_idx, dropout=dropout, scalar=1.0)
25
```

```
26         self.trg_word_prj = nn.Linear(d_model, n_trg_vocab)
27
28     def forward(self, src_seq, trg_seq):
29
30         src_mask = get_pad_mask(src_seq, self.src_pad_idx)
31         trg_mask = get_pad_mask(trg_seq, self.trg_pad_idx)
32
33         enc_output, *_ = self.encoder(src_seq, src_mask)
34         dec_output, *_ = self.decoder(trg_seq, trg_mask, enc_output)
35         seq_logit = self.trg_word_prj(dec_output)
36         if self.scale_prj:
37             seq_logit *= self.d_model ** -0.5
38
39         return seq_logit.view(-1, seq_logit.size(2))
```

Train model

> Todo

Evaluation

> Todo