

```

from keras.layers import Layer

class FeedForward(Layer):
    def __init__(self, d_ff, d_model, **kwargs):
        super(FeedForward, self).__init__(**kwargs)
        self.fully_connected1 = Dense(d_ff) # First fully connected layer
        self.fully_connected2 = Dense(d_model) # Second fully connected layer
        self.activation = ReLU() # ReLU activation layer

    def call(self, x):
        # The input is passed into the two fully-connected layers, with a ReLU in between
        x_fc1 = self.fully_connected1(x)

        return self.fully_connected2(self.activation(x_fc1))

class AddNormalization(Layer):
    def __init__(self, **kwargs):
        super(AddNormalization, self).__init__(**kwargs)
        self.layer_norm = LayerNormalization() # Layer normalization layer

    def call(self, x, sublayer_x):
        # The sublayer input and output need to be of the same shape to be summed
        add = x + sublayer_x

        # Apply layer normalization to the sum
        return self.layer_norm(add)

class EncoderLayer(Layer):
    def __init__(self, h, d_k, d_v, d_model, d_ff, rate, **kwargs):
        super(EncoderLayer, self).__init__(**kwargs)
        self.multihead_attention = MultiHeadAttention(h, d_k, d_v, d_model)
        self.dropout1 = Dropout(rate)
        self.add_norm1 = AddNormalization()
        self.feed_forward = FeedForward(d_ff, d_model)
        self.dropout2 = Dropout(rate)
        self.add_norm2 = AddNormalization()

    def call(self, x, padding_mask, training):
        # Multi-head attention layer
        multihead_output = self.multihead_attention(x, x, x, padding_mask)
        # Expected output shape = (batch_size, sequence_length, d_model)

        # Add in a dropout layer
        multihead_output = self.dropout1(multihead_output, training=training)

        # Followed by an Add & Norm layer
        addnorm_output = self.add_norm1(x, multihead_output)
        # Expected output shape = (batch_size, sequence_length, d_model)

        # Followed by a fully connected layer
        feedforward_output = self.feed_forward(addnorm_output)
        # Expected output shape = (batch_size, sequence_length, d_model)

        # Add in another dropout layer
        feedforward_output = self.dropout2(feedforward_output, training=training)

        # Followed by another Add & Norm layer
        return self.add_norm2(addnorm_output, feedforward_output)

class Encoder(Layer):
    def __init__(self, vocab_size, sequence_length, h, d_k, d_v, d_model, d_ff, n, rate, **kwargs):
        super(Encoder, self).__init__(**kwargs)
        self.pos_encoding = PositionEmbeddingFixedWeights(sequence_length, vocab_size, d_model)
        self.dropout = Dropout(rate)
        self.encoder_layer = [EncoderLayer(h, d_k, d_v, d_model, d_ff, rate) for _ in range(n)]
        ...

    def call(self, input_sentence, padding_mask, training):
        # Generate the positional encoding
        pos_encoding_output = self.pos_encoding(input_sentence)
        # Expected output shape = (batch_size, sequence_length, d_model)

        # Add in a dropout layer
        x = self.dropout(pos_encoding_output, training=training)

```

```

        # Pass on the positional encoded values to each encoder layer
        for i, layer in enumerate(self.encoder_layer):
            x = layer(x, padding_mask, training)

    return x

from tensorflow import math, matmul, reshape, shape, transpose, cast, float32
from tensorflow.keras.layers import Dense, Layer
from keras.backend import softmax

# Implementing the Scaled-Dot Product Attention
class DotProductAttention(Layer):
    def __init__(self, **kwargs):
        super(DotProductAttention, self).__init__(**kwargs)

    def call(self, queries, keys, values, d_k, mask=None):
        # Scoring the queries against the keys after transposing the latter, and scaling
        scores = matmul(queries, keys, transpose_b=True) / math.sqrt(cast(d_k, float32))

        # Apply mask to the attention scores
        if mask is not None:
            scores += -1e9 * mask

        # Computing the weights by a softmax operation
        weights = softmax(scores)

        # Computing the attention by a weighted sum of the value vectors
        return matmul(weights, values)

# Implementing the Multi-Head Attention
class MultiHeadAttention(Layer):
    def __init__(self, h, d_k, d_v, d_model, **kwargs):
        super(MultiHeadAttention, self).__init__(**kwargs)
        self.attention = DotProductAttention() # Scaled dot product attention
        self.heads = h # Number of attention heads to use
        self.d_k = d_k # Dimensionality of the linearly projected queries and keys
        self.d_v = d_v # Dimensionality of the linearly projected values
        self.d_model = d_model # Dimensionality of the model
        self.W_q = Dense(d_k) # Learned projection matrix for the queries
        self.W_k = Dense(d_k) # Learned projection matrix for the keys
        self.W_v = Dense(d_v) # Learned projection matrix for the values
        self.W_o = Dense(d_model) # Learned projection matrix for the multi-head output

    def reshape_tensor(self, x, heads, flag):
        if flag:
            # Tensor shape after reshaping and transposing: (batch_size, heads, seq_length, -1)
            x = reshape(x, shape=(shape(x)[0], shape(x)[1], heads, -1))
            x = transpose(x, perm=(0, 2, 1, 3))
        else:
            # Reverting the reshaping and transposing operations: (batch_size, seq_length, d_k)
            x = transpose(x, perm=(0, 2, 1, 3))
            x = reshape(x, shape=(shape(x)[0], shape(x)[1], self.d_k))
        return x

    def call(self, queries, keys, values, mask=None):
        # Rearrange the queries to be able to compute all heads in parallel
        q_resaped = self.reshape_tensor(self.W_q(queries), self.heads, True)
        # Resulting tensor shape: (batch_size, heads, input_seq_length, -1)

        # Rearrange the keys to be able to compute all heads in parallel
        k_resaped = self.reshape_tensor(self.W_k(keys), self.heads, True)
        # Resulting tensor shape: (batch_size, heads, input_seq_length, -1)

        # Rearrange the values to be able to compute all heads in parallel
        v_resaped = self.reshape_tensor(self.W_v(values), self.heads, True)
        # Resulting tensor shape: (batch_size, heads, input_seq_length, -1)

        # Compute the multi-head attention output using the reshaped queries, keys and values
        o_resaped = self.attention(q_resaped, k_resaped, v_resaped, self.d_k, mask)
        # Resulting tensor shape: (batch_size, heads, input_seq_length, -1)

        # Rearrange back the output into concatenated form
        output = self.reshape_tensor(o_resaped, self.heads, False)
        # Resulting tensor shape: (batch_size, input_seq_length, d_v)

        # Apply one final linear projection to the output to generate the multi-head attention
        # Resulting tensor shape: (batch_size, input_seq_length, d_model)
        return self.W_o(output)

```

```

class PositionEmbeddingFixedWeights(Layer):
    def __init__(self, sequence_length, vocab_size, output_dim, **kwargs):
        super(PositionEmbeddingFixedWeights, self).__init__(**kwargs)
        word_embedding_matrix = self.get_position_encoding(vocab_size, output_dim)
        position_embedding_matrix = self.get_position_encoding(sequence_length, output_dim)
        self.word_embedding_layer = Embedding(
            input_dim=vocab_size, output_dim=output_dim,
            weights=[word_embedding_matrix],
            trainable=False
        )
        self.position_embedding_layer = Embedding(
            input_dim=sequence_length, output_dim=output_dim,
            weights=[position_embedding_matrix],
            trainable=False
        )

    def get_position_encoding(self, seq_len, d, n=10000):
        P = np.zeros((seq_len, d))
        for k in range(seq_len):
            for i in np.arange(int(d/2)):
                denominator = np.power(n, 2*i/d)
                P[k, 2*i] = np.sin(k/denominator)
                P[k, 2*i+1] = np.cos(k/denominator)
        return P

    def call(self, inputs):
        position_indices = tf.range(tf.shape(inputs)[-1])
        embedded_words = self.word_embedding_layer(inputs)
        embedded_indices = self.position_embedding_layer(position_indices)
        return embedded_words + embedded_indices

from tensorflow.keras.layers import LayerNormalization, Layer, Dense, ReLU, Dropout
import tensorflow as tf
from tensorflow.keras.layers import Layer, Embedding, Dropout, Dense, LayerNormalization, LSTM
import numpy as np

# from multihead_attention import MultiHeadAttention
# from positional_encoding import PositionEmbeddingFixedWeights

# Implementing the Add & Norm Layer
class AddNormalization(Layer):
    def __init__(self, **kwargs):
        super(AddNormalization, self).__init__(**kwargs)
        self.layer_norm = LayerNormalization() # Layer normalization layer

    def call(self, x, sublayer_x):
        # The sublayer input and output need to be of the same shape to be summed
        add = x + sublayer_x

        # Apply layer normalization to the sum
        return self.layer_norm(add)

# Implementing the Feed-Forward Layer
class FeedForward(Layer):
    def __init__(self, d_ff, d_model, **kwargs):
        super(FeedForward, self).__init__(**kwargs)
        self.fully_connected1 = Dense(d_ff) # First fully connected layer
        self.fully_connected2 = Dense(d_model) # Second fully connected layer
        self.activation = ReLU() # ReLU activation layer

    def call(self, x):
        # The input is passed into the two fully-connected layers, with a ReLU in between
        x_fc1 = self.fully_connected1(x)

        return self.fully_connected2(self.activation(x_fc1))

# Implementing the Encoder Layer
class EncoderLayer(Layer):
    def __init__(self, d_model, num_heads, d_k, d_v, d_ff, rate):
        super(EncoderLayer, self).__init__()
        self.mha = tf.keras.layers.MultiHeadAttention(key_dim=d_k, value_dim=d_v, num_heads=num_heads, dropout=rate)
        self.ffn = tf.keras.Sequential([Dense(d_ff, activation='relu'), Dense(d_model)])
        self.layer_norm1 = LayerNormalization(epsilon=1e-6)
        self.layer_norm2 = LayerNormalization(epsilon=1e-6)
        self.dropout1 = Dropout(rate)
        self.dropout2 = Dropout(rate)

    def call(self, x, padding_mask, training):
        attn_output = self.mha(x, x, x, attention_mask=padding_mask)

```

```

        attn_output = self.dropout1(attn_output, training=training)
        out1 = self.layer_norm1(x + attn_output)

        ffn_output = self.ffn(out1)
        ffn_output = self.dropout2(ffn_output, training=training)
        out2 = self.layer_norm2(out1 + ffn_output)

        return out2

class Encoder(Layer):
    def __init__(self, vocab_size, sequence_length, d_model, num_heads, d_k, d_v, d_ff, n, rate):
        super(Encoder, self).__init__()
        self.embedding = Embedding(vocab_size, d_model)
        self.dropout = Dropout(rate)
        self.lstm_layer = LSTM(d_model, return_sequences=True) # Set return_sequences=True for LSTM output
        self.encoder_layer = [EncoderLayer(d_model, num_heads, d_k, d_v, d_ff, rate) for _ in range(n)]

    def call(self, input_sentence, padding_mask, training):
        # Embed the input sentence
        x = self.embedding(input_sentence)
        # Expected output shape = (batch_size, sequence_length, d_model)

        # Add in a dropout layer
        x = self.dropout(x, training=training)

        # Pass through the LSTM layer
        x = self.lstm_layer(x)

        # Pass on the LSTM encoded values to each encoder layer
        for i, layer in enumerate(self.encoder_layer):
            x = layer(x, padding_mask, training)

        return x

from numpy import random
import numpy as np
from keras.layers import Embedding
import tensorflow as tf

enc_vocab_size = 20 # Vocabulary size for the encoder
input_seq_length = 5 # Maximum length of the input sequence
h = 8 # Number of self-attention heads
d_k = 64 # Dimensionality of the linearly projected queries and keys
d_v = 64 # Dimensionality of the linearly projected values
d_ff = 2048 # Dimensionality of the inner fully connected layer
d_model = 512 # Dimensionality of the model sub-layers' outputs
n = 6 # Number of layers in the encoder stack

batch_size = 64 # Batch size from the training process
dropout_rate = 0.1 # Frequency of dropping the input units in the dropout layers

input_seq = random.random((batch_size, input_seq_length))

encoder = Encoder(enc_vocab_size, input_seq_length, h, d_k, d_v, d_model, d_ff, n, dropout_rate)
print(encoder(input_seq, None, True))

[-1.5030246 ]]

[[-0.34877568  1.9845133  0.5497115  ...  0.8198611 -0.95096403
 -1.2634088 ]
 [-0.21195862  1.9619902  0.5660541  ...  0.6644101 -0.83642197
 -1.501505 ]
 [-0.26460993  1.8980888  0.5059202  ...  0.7049494 -0.6726516
 -1.7042773 ]
 [-0.2648147  1.9425071  0.53728694 ...  0.53116846 -0.6231487
 -1.7524109 ]
 [-0.40019843  1.9137564  0.6893785  ...  0.31915146 -0.52713406
 -1.7854747 ]]

[[-0.14925282  2.1031408  0.6289158  ...  0.54821575 -0.9769976
 -1.1765658 ]

```

```

[[-0.07305208  1.8525075  0.531602  ...  0.94143164 -1.0262085
  -1.2804625 ]
 [-0.24272148  1.9749895  0.42825136 ...  0.7151026  -0.79410243
  -1.5513829 ]
 [ 0.04938693  1.9081374  0.4691741  ...  0.7746823  -0.87582046
  -1.5132338 ]
 [-0.15292965  1.897566  0.41157877 ...  0.57161134 -0.6417308
  -1.8375701 ]
 [-0.31671885  1.9746364  0.47018653 ...  0.5163497  -0.60056853
  -1.7420424 ]]

[[-0.09563909  1.8791151  0.6772185  ...  0.88283867 -0.8389759
  -1.3676262 ]
 [-0.14158905  2.0193555  0.52083296 ...  0.5879242  -0.84512633
  -1.4669443 ]
 [-0.4182232  1.9178753  0.5560412  ...  0.67514837 -0.643989
  -1.6846435 ]
 [-0.3341174  1.888315  0.47670078 ...  0.55996275 -0.66979635
  -1.7794807 ]
 [-0.41337007  1.9312357  0.43267813 ...  0.3185696  -0.40980154
  -1.8136287 ]]

[[ 0.29617316  2.1797907 -1.0161338  ... -0.2990103  -0.35619932
  -1.2462692 ]
 [ 0.12402245  2.3038228 -0.13626449 ...  0.26252392 -0.71775293
  -1.3289509 ]
 [ 0.10294637  2.1530223  0.08216639 ...  0.5819097  -0.71855336
  -1.4973611 ]
 [ 0.04791594  2.0326815  0.2432545  ...  0.4788594  -0.780803
  -1.6656587 ]
 [-0.1387699  2.0461552  0.25876433 ...  0.48415136 -0.59909636
  -1.7274647 ]]], shape=(64, 5, 8), dtype=float32)

```