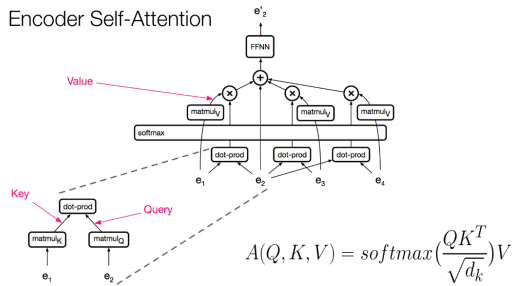


Attention

Encoder Self-Attention



```
def scaled_dot_product_attention(q, k, v, mask):
    """Calculate the attention weights.
    q, k, v must have matching leading dimensions.
    k, v must have matching penultimate dimension, i.e.: seq_len_k = seq_len_v.
    The mask has different shapes depending on its type(padding or look ahead)
    but it must be broadcastable for addition.

    Args:
        q: query shape == (..., seq_len_q, depth)
        k: key shape == (..., seq_len_k, depth)
        v: value shape == (..., seq_len_v, depth_v)
        mask: Float tensor with shape broadcastable
            to (..., seq_len_q, seq_len_k). Defaults to None.

    Returns:
        output, attention_weights
    """

    matmul_qk = tf.matmul(q, k, transpose_b=True) # (..., seq_len_q, seq_len_k)

    # scale matmul_qk
    dk = tf.cast(tf.shape(k)[-1], tf.float32)
    scaled_attention_logits = matmul_qk / tf.math.sqrt(dk)

    # add the mask to the scaled tensor.
    if mask is not None:
        scaled_attention_logits += (mask * -1e9)

    # softmax is normalized on the last axis (seq_len_k) so that the scores
    # add up to 1.
    attention_weights = tf.nn.softmax(scaled_attention_logits, axis=-1) # (..., seq_len_q, seq_len_k)

    output = tf.matmul(attention_weights, v) # (..., seq_len_q, depth_v)

    return output, attention_weights
```

```
temp_k = tf.constant([[10,0,0], [0,10,0], [0,0,10], [0,0,10]], dtype=tf.float32) # (4, 3)
temp_v = tf.constant([[1,0], [10,0], [100,5], [1000,6]], dtype=tf.float32) # (4, 2)

# This 'query' aligns with the second 'key',
# so the second 'value' is returned.
temp_q = tf.constant([[0, 10, 0]], dtype=tf.float32) # (1, 3)
print_out(temp_q, temp_k, temp_v)
```

```
def print_out(q, k, v):
    temp_out, temp_attn = scaled_dot_product_attention(
        q, k, v, None)
    print ('Attention weights are:')
    print (temp_attn)
    print ('Output is:')
    print (temp_out)
```

ONE query → ONE RESPONSE

```
Attention weights are:
tf.Tensor([[0. 0. 0.5 0.5]], shape=(1, 4), dtype=float32)
Output is:
tf.Tensor([[550. 5.5]], shape=(1, 2), dtype=float32)
```

Response - is weighted representation of values (words) according to query

Let's say :

w1 = a cat
w2 = playing with
w3 = a ball
w4 = and a strip

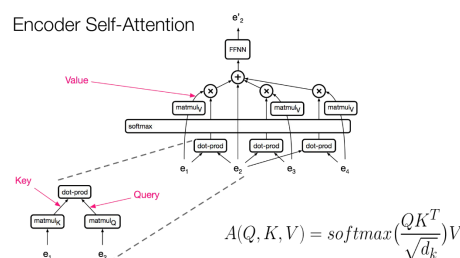
Query -> what cat is playing with?

Response -> $0.5 * \text{ball}$ and $0.5 * \text{strip}$
 $0.5 * [100, 5] + 0.5 * [1000, 6] = [550, 5.5]$

```
temp_q = tf.constant([[0, 0, 10], [0, 10, 0], [10, 10, 0]], dtype=tf.float32) # (3, 3)
print_out(temp_q, temp_k, temp_v)
```

```
Attention weights are:
tf.Tensor(
[[0. 0. 0.5 0.5]
 [0. 1. 0. 0. ]
 [0.5 0.5 0. 0. ]], shape=(3, 4), dtype=float32)
Output is:
tf.Tensor(
[[550. 5.5]
 [ 10. 0. ]
 [ 5.5 0. ]], shape=(3, 2), dtype=float32)
```

Encoder Self-Attention



For each word we get a Key and return a Response based on Query

In first layers we can create a transformations where each word is represented through itself
 it means $\text{Len(Query)} = \text{Len(Keys)}$

However, number of Queries to one sentence should not be one to one.