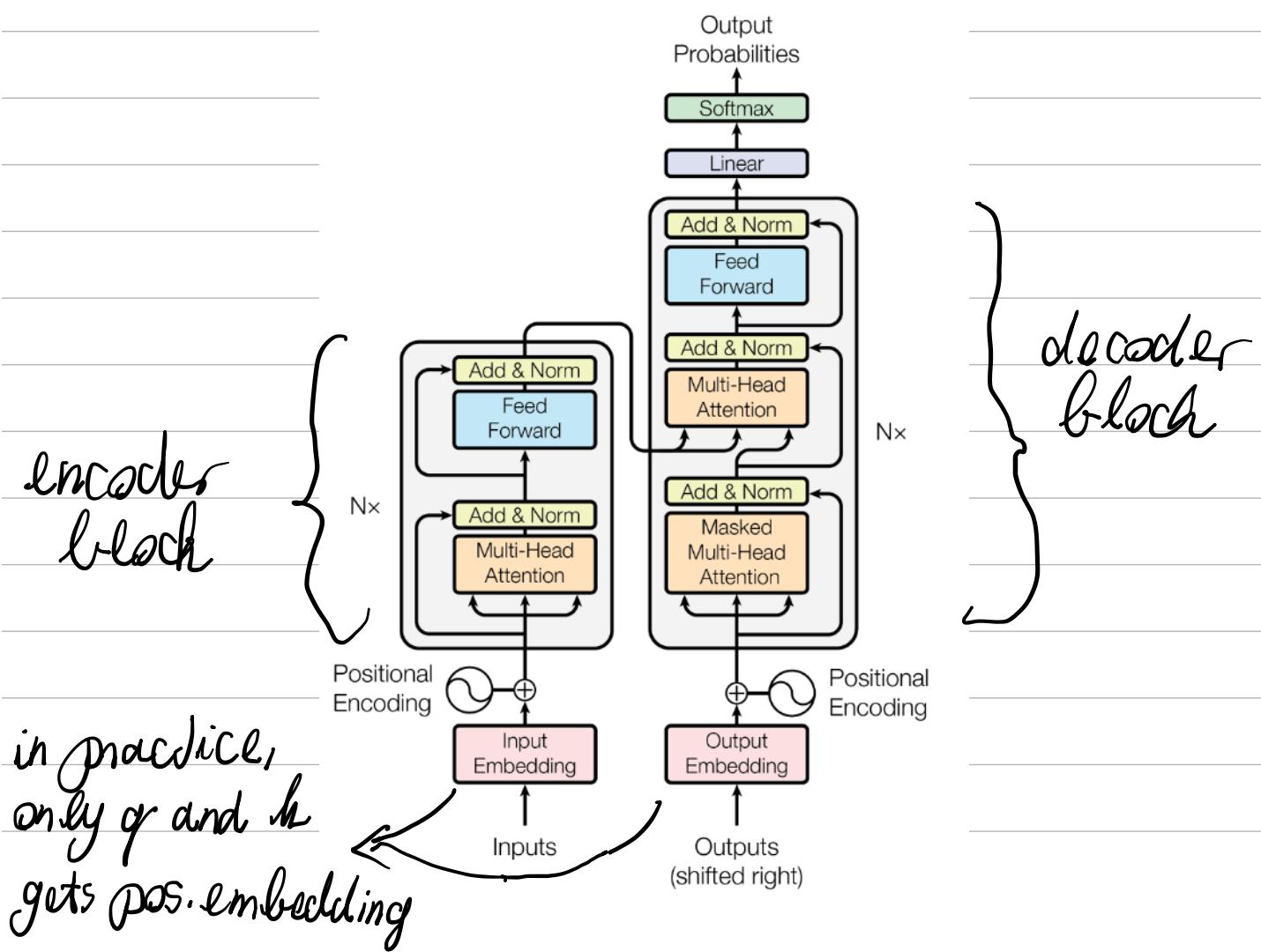


# Transformer - Attention is all you need

## Model architecture

- ↳ Encoder - decoder
- ↳ encoder, maps input  $(x_1, \dots x_n)$  to  
 $Z = (z_1 \dots z_n)$
- ↳ decoder maps  $Z$  to  $(y_1, \dots y_n)$
- ↳ auto-regressive, meaning it consumes previously generated symbols as additional input.



## Encoder

- ↳ They use 6 blocks
  - ↳ Each block has 2 sub-layers
    - ↳ multi-head self-attention
    - ↳ position-wise fully connected feed forward network
  - ↳ residual connections, as in the image above
  - ↳ Uses layer norm
  - ↳  $d_{model} = 512$

## Decoder

- ↳ same way, using 6 blocks
  - ↳ blocks are similar, but has an additional sub-layer
    - ↳ multi-head attention over the output of the encoder stack

↳ also, The self-attention is modified (masked) to prevent from looking at future positions

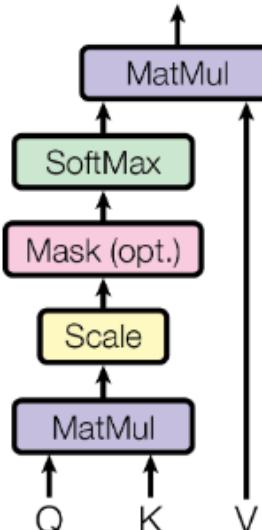
## Attention

↳ Scaled dot product attention

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

Scaled Dot-Product Attention

It's the same as above, just visualized



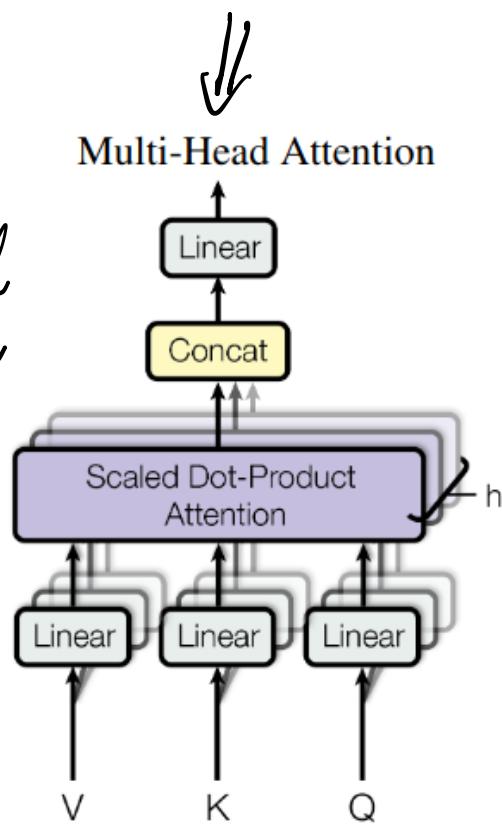
Scaling

↓  
to counter softmax creating small gradients when the values are large

↳ Multi-head attention

↳ meaning, we linearly project  $Q, K, V$   $N$  times and do this parallel, then concat

allows the model to jointly attend to information from different representation subspaces at different positions.



## Position-wise Feed Forward Networks

↳ Each encoder, decoder block contains a fully-connected FFN

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

ReLU

↳ Inner dim is  $\rightarrow d_{ff} = 2048$ , while  $d_{model} = 512$

## Embedding and Softmax

↳ Learned embeddings convert input and output tokens to dim  $d_{model}$

↳ Learned linear transform +  
softmax is used to convert decoder  
output to next-token prob.

↳ Shared weight matrix between  
the 2 embedding layers + pre-softmax  
linear transform

↳  $W_e \rightarrow$  used to embed input  
tokens

$W_e' \rightarrow$  used to embed output  
tokens (during decoding)

$W_{out} \rightarrow$  decoder output  
projection to logits

$$W_e = W_e' = W$$
$$W_{out} = W^T$$

## Positional encoding

↳ Order of sequence is important, but  
the architecture doesn't have any  
information on it



This is why we need positional emb.

↳ same sized vector as  $d_{model}$ , and is added to  $K, Q$

↳ In the original implementation

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

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

where,

$pos$  = position in the sequence

$i$  = index inside the pos-emb. vector

$\text{range}(i) \Rightarrow 0 \rightarrow \frac{d_{model}}{2} - 1 \rightarrow 0.$  starting index  
1/2 because of sin+cos

This way, the wavelengths range from  $2\pi$  to  $10000 \cdot 2\pi$



The Positional encoding at  $pos+h$  can be expressed as a linear func. of the pos. encoding at position  $pos$  when  $h$  is a fixed offset



- 1 Linearity makes it easy for the model to learn the relationship between others.

✓ „Proof“

$$\omega_i = \frac{1}{10000 \cdot d_{\text{mache}}}$$

$$\sin(A+B) = \sin A \cos B + \cos A \sin B$$

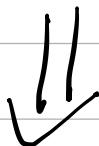
$$\cos(A+B) = \cos A \cos B - \sin A \sin B$$

$$PE(pos, 2i) = \sin(\omega_i \cdot pos)$$

$$PE(pos, 2i+1) = \cos(\omega_i \cdot pos)$$

$$PE(pos+h, 2i) = \sin(\omega_i (pos+h)) = \sin(\omega_i pos + \omega_i h)$$

$$PE(pos+h, 2i+1) = \cos(\omega_i (pos+h)) = \cos(\omega_i pos + \omega_i h)$$



$$PE(pos+h, 2i) = \sin(\omega_i pos + \omega_i h) =$$

$$= \sin(\omega_i pos) \cos(\omega_i h) + \cos(\omega_i pos) \sin(\omega_i h)$$

$$= PE(pos, 2i) \cos(\omega_i h) + PE(pos, 2i+1) \sin(\omega_i h)$$

$$PE(pos+h, 2i+1) = \cos(\omega_i pos + \omega_i h) =$$

$$= \cos(\omega_i pos) \cos(\omega_i h) - \sin(\omega_i pos) \sin(\omega_i h)$$

$$= PE(pos, 2i+1) \cos(\omega_i h) - PE(pos, 2i) \sin(\omega_i h)$$

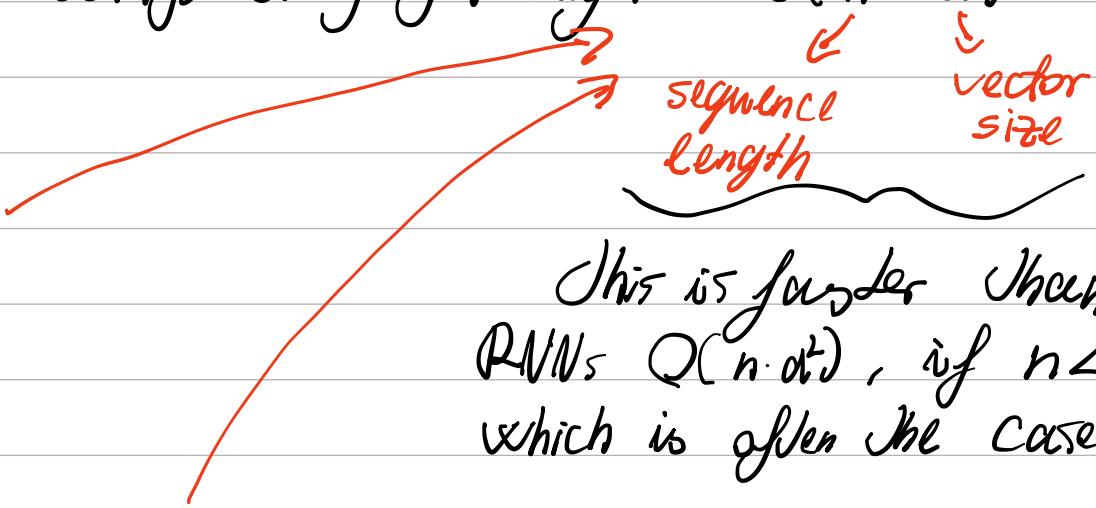
As we can see,  $pos+h$  is indeed a lin. comb of  $pos$

## Why self-attention?

↳ constant number of sequentially executed operations, regardless of sequence length  
↳ because all positions run in parallel

↳ complexity per layer  $\rightarrow O(n^2 \cdot d)$

$n^2$ , because we check everything for everything (all vectors)



This is faster than RNNs  $O(n \cdot d^2)$ , if  $n \ll d$ , which is often the case

↳ Proposed problem

↳ dealing with long sequences ( $n^2$ )

Explains  $n \cdot n$  vs  $r \cdot r$

Restricted self-attention: could only see  $r$  positions around itself.

reduced complexity  $\rightarrow O(n \cdot r \cdot d)$