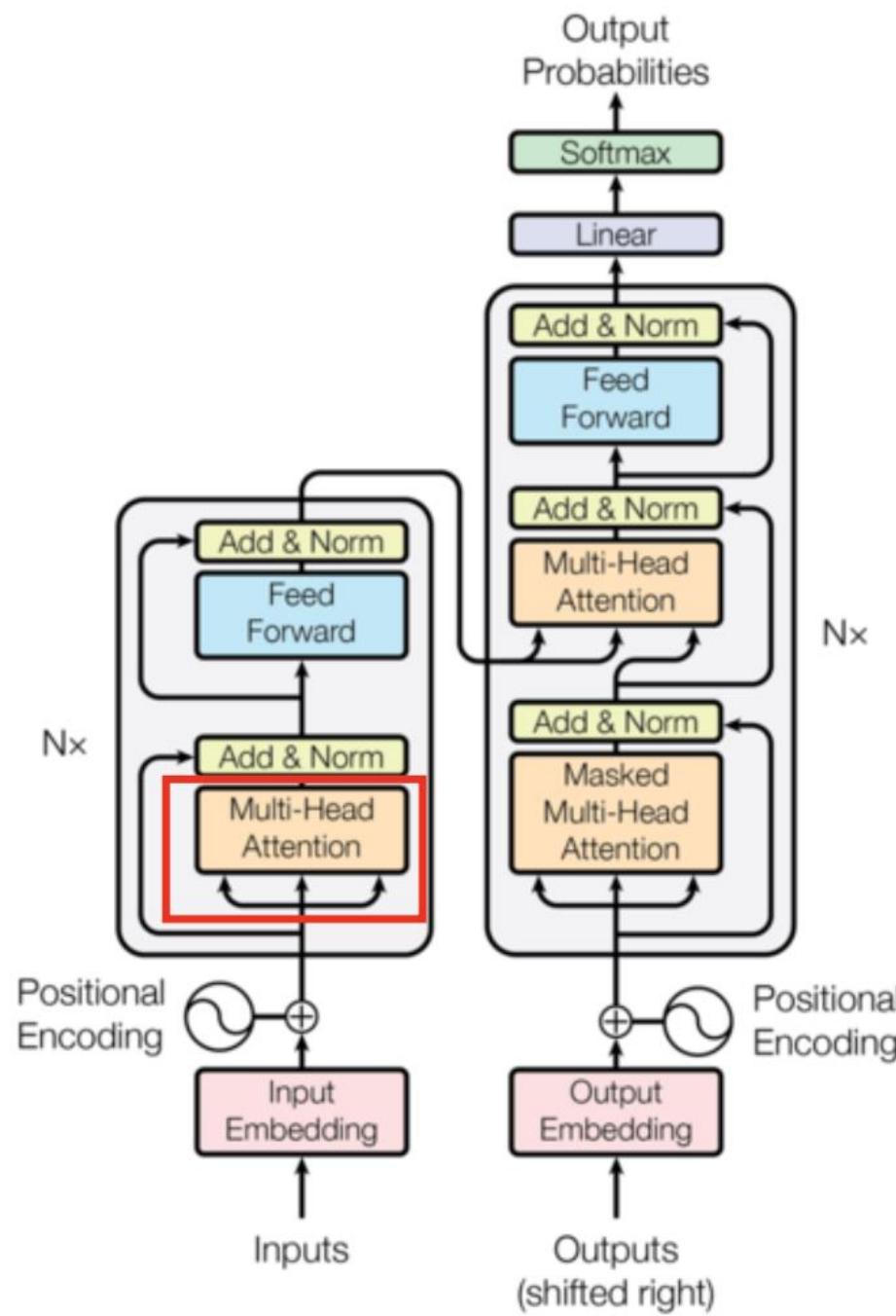


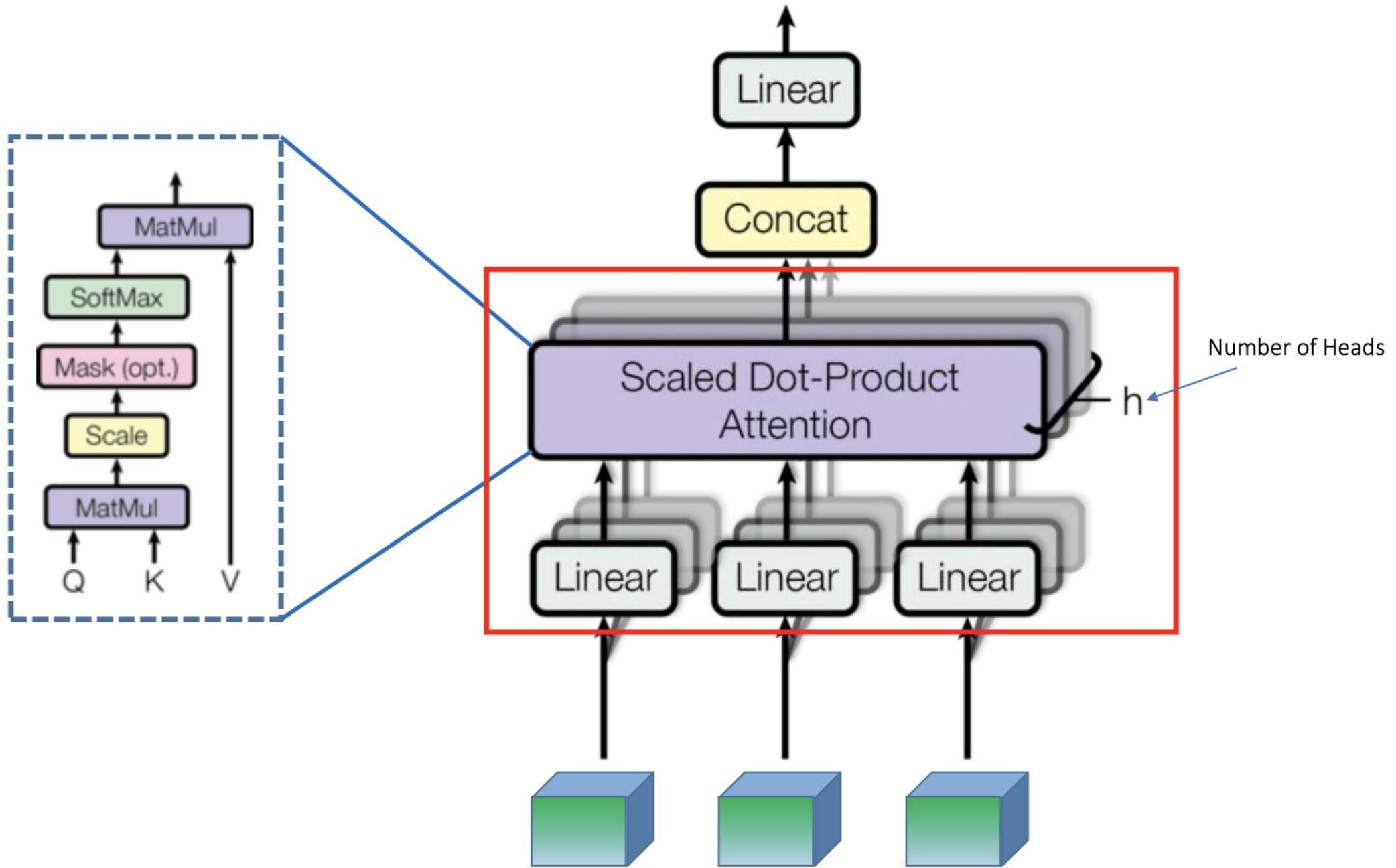
Transformer II

Yuxuan Du

EE 5993 AI Practicum

03/24/2025





Equations for Multi-Head Attention

As we said, we have more than one attention (multi-head attention).

For each head i:

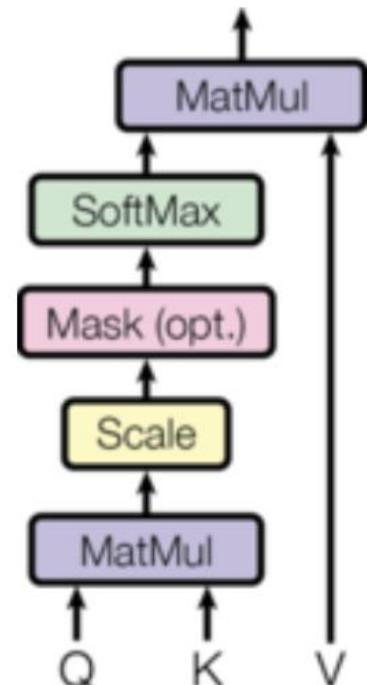
$$\text{head}_i = \text{Attention}(QW_i^Q, KW_i^K, VW_i^V)$$

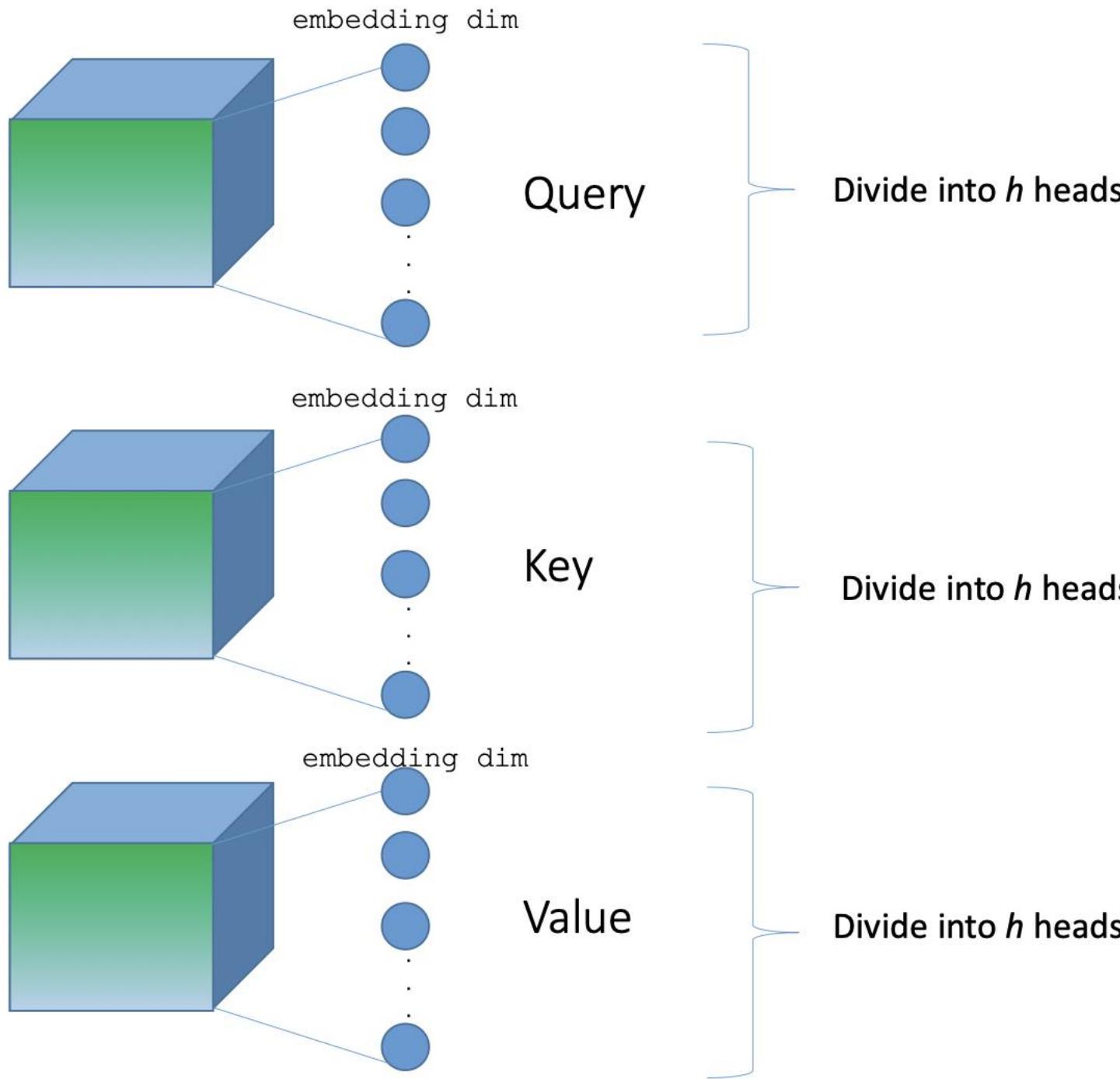
$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

After that, concatenate them back and run them through a linear layer:

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

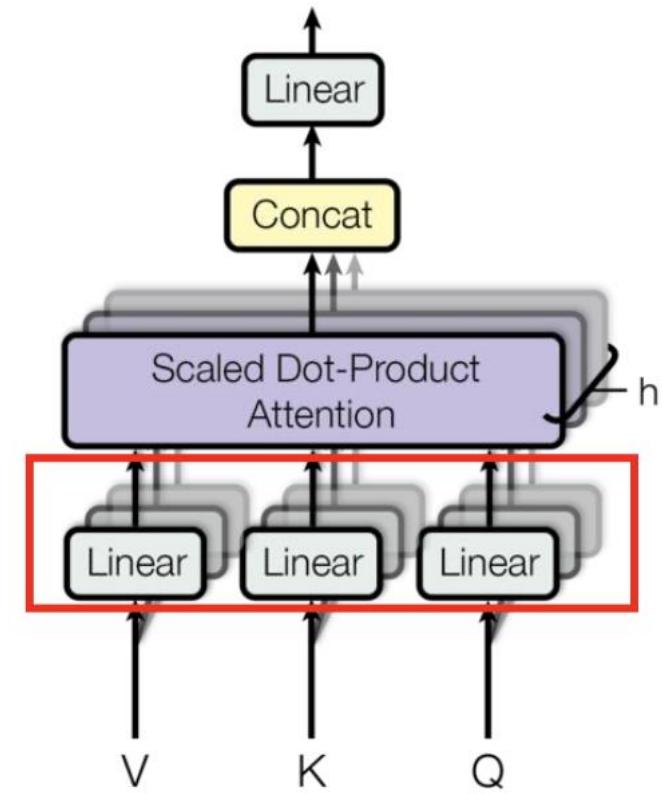
Now Let's go in detail into each one of them.....





The dimensionality of each head d_k is $\text{embedding_dim} / h$

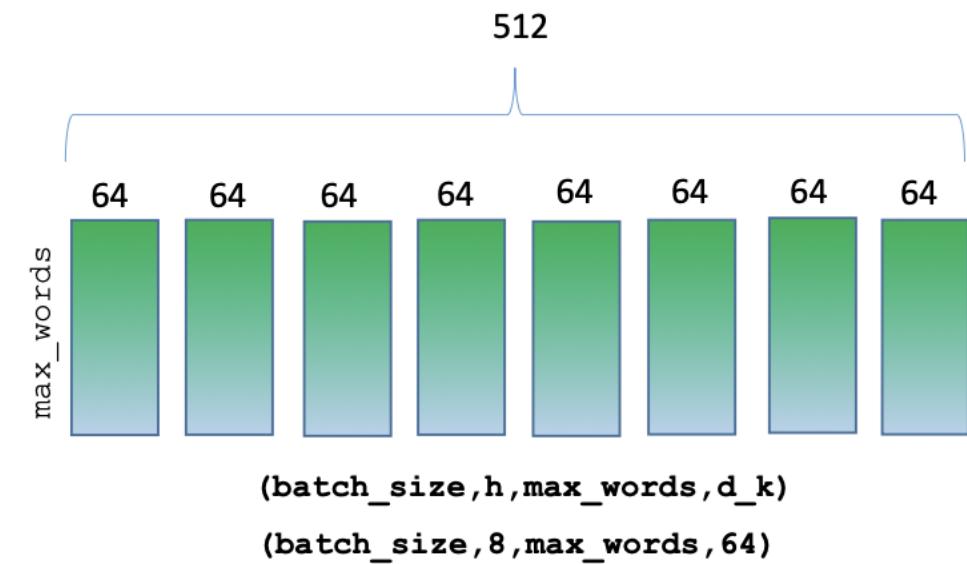
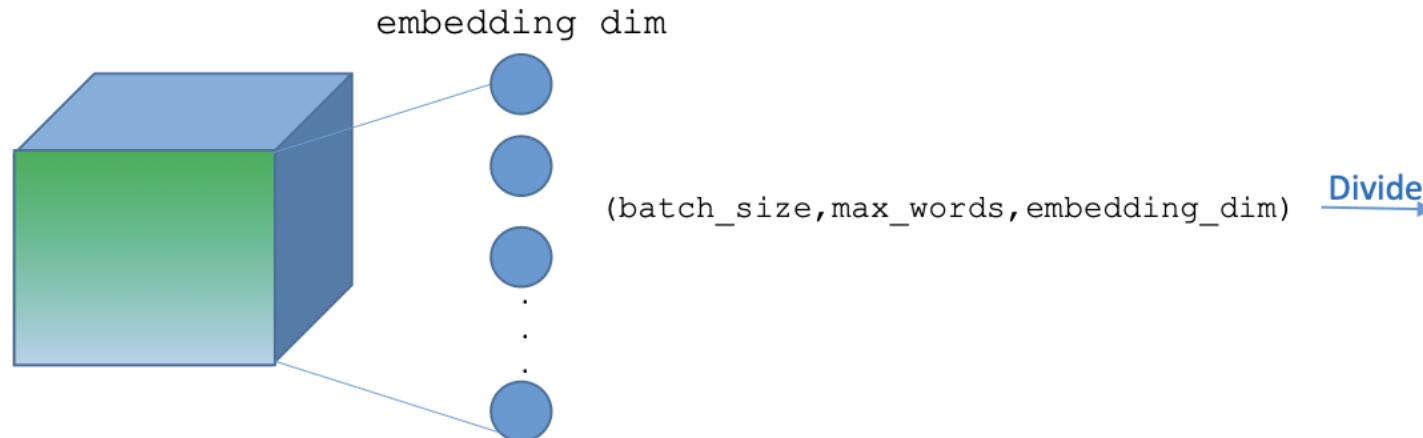
Example: If we use 8 heads, then each head would have a dimension of: $512 / 8 = 64$

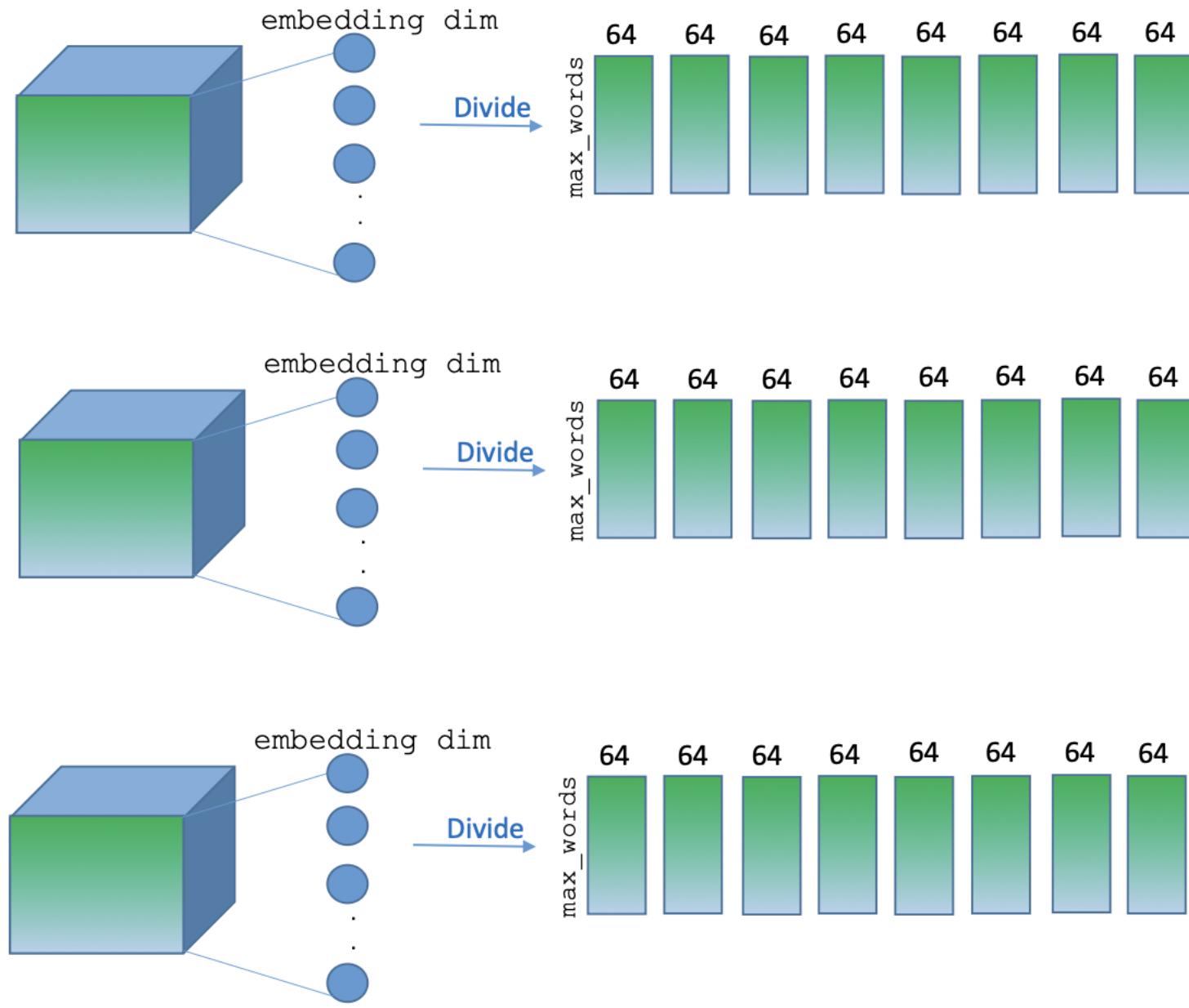


Let's choose:

embedding dimension = 512

Number of heads = 8



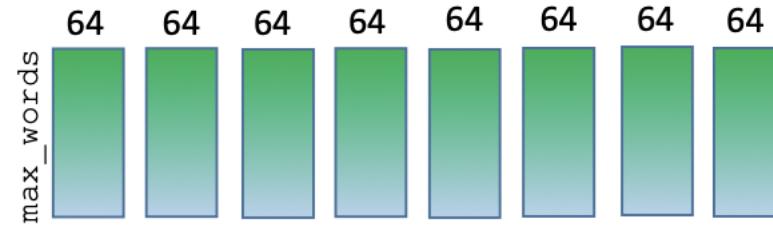


Query

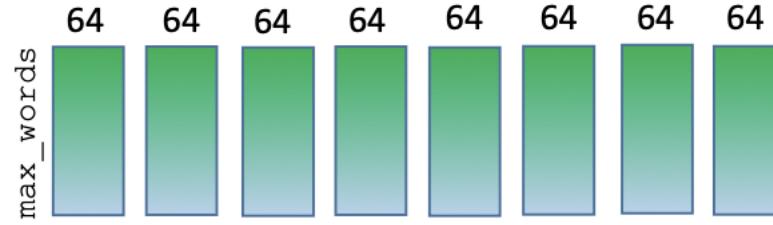
Key

Value

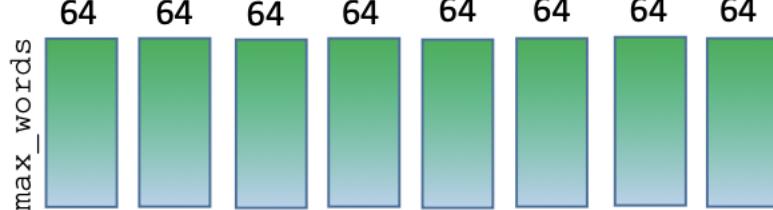
Query



Key



Value



$$(\text{batch_size}, \text{h}, \text{max_len}, d_k) \odot (\text{batch_size}, \text{h}, d_k, \text{max_len})$$

$$(\text{batch_size}, \text{h}, \text{max_len}, \text{max_len})$$

For each

Dot Product and Scale

After dot product, we might get large numbers, so we scale them

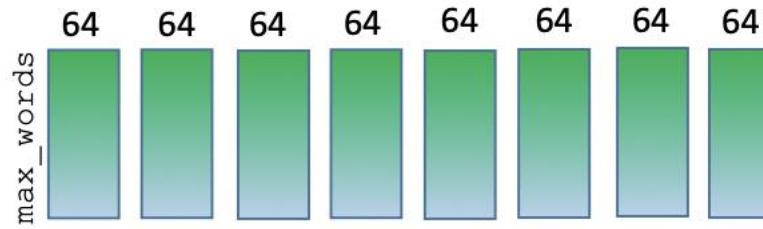
Diagram illustrating the scaling operation. A blue bracket groups the first two terms of the equation. An arrow points from this bracket to a diagram showing a green rectangle of size $\text{max_len} \times \text{max_len}$ being multiplied by a scalar $/ \sqrt{d_k}$. The scalar is represented by a blue rectangle of size $1 \times \text{max_len}$.

$$(\text{batch_size}, \text{h}, \text{max_len}, \text{max_len}) = (\text{batch_size}, \text{h}, d_k, \text{max_len}) / \sqrt{d_k}$$

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

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

Query

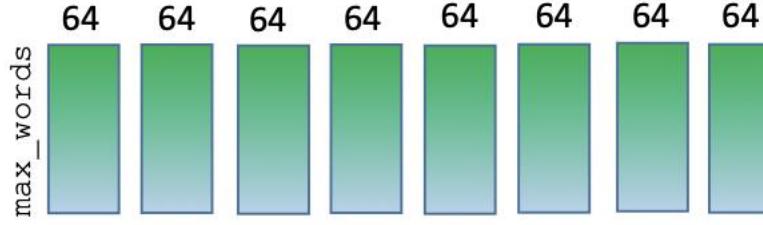


For each

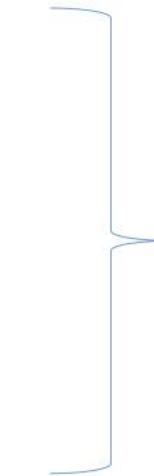
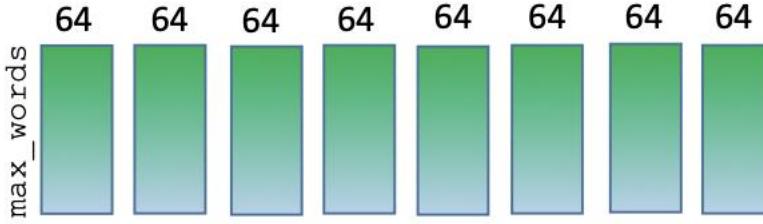


Dot Product and Scale

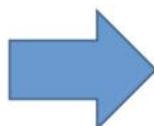
Key



Value

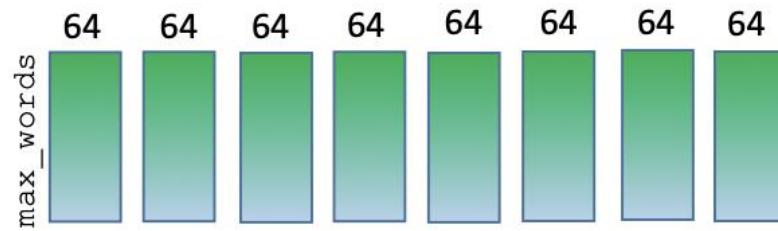


Softmax

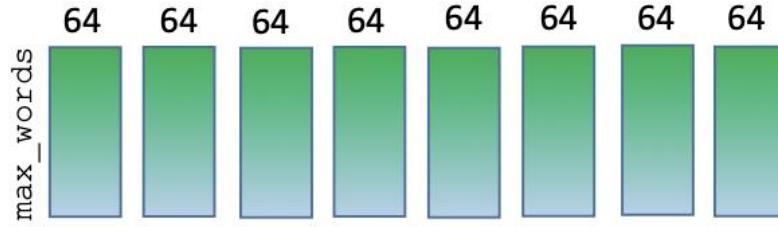


$$\text{Softmax} \rightarrow \text{Matrix} \rightarrow \text{Sparse Matrix}$$
$$\text{Result} = \text{Query} \odot \text{Key}^\top = \text{Query} \cdot \text{Key}^\top / \sqrt{64}$$

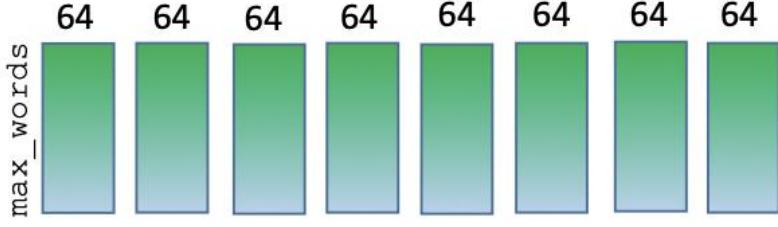
Query



Key



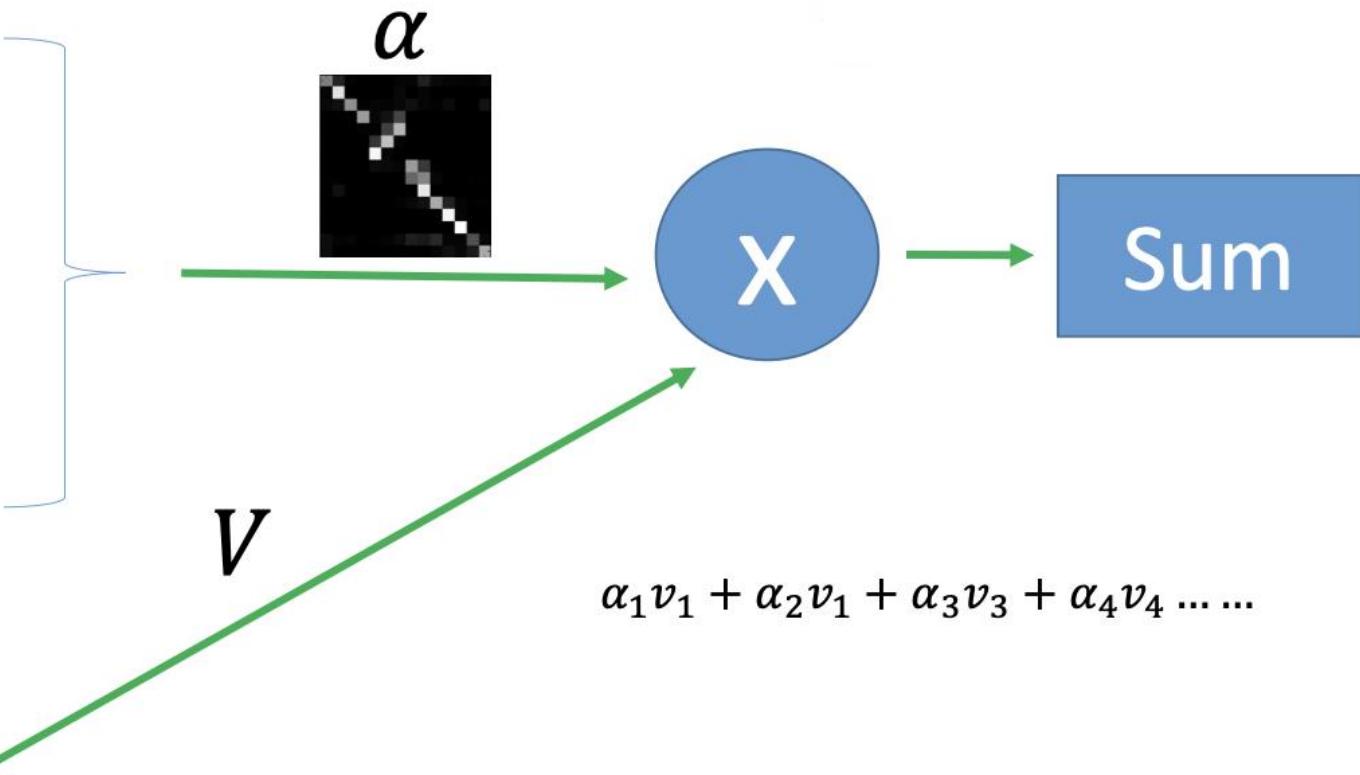
Value



$$(\text{batch_size}, \text{h}, \text{max_len}, \text{max_len}) \odot (\text{batch_size}, \text{h}, \text{max_len}, d_k)$$

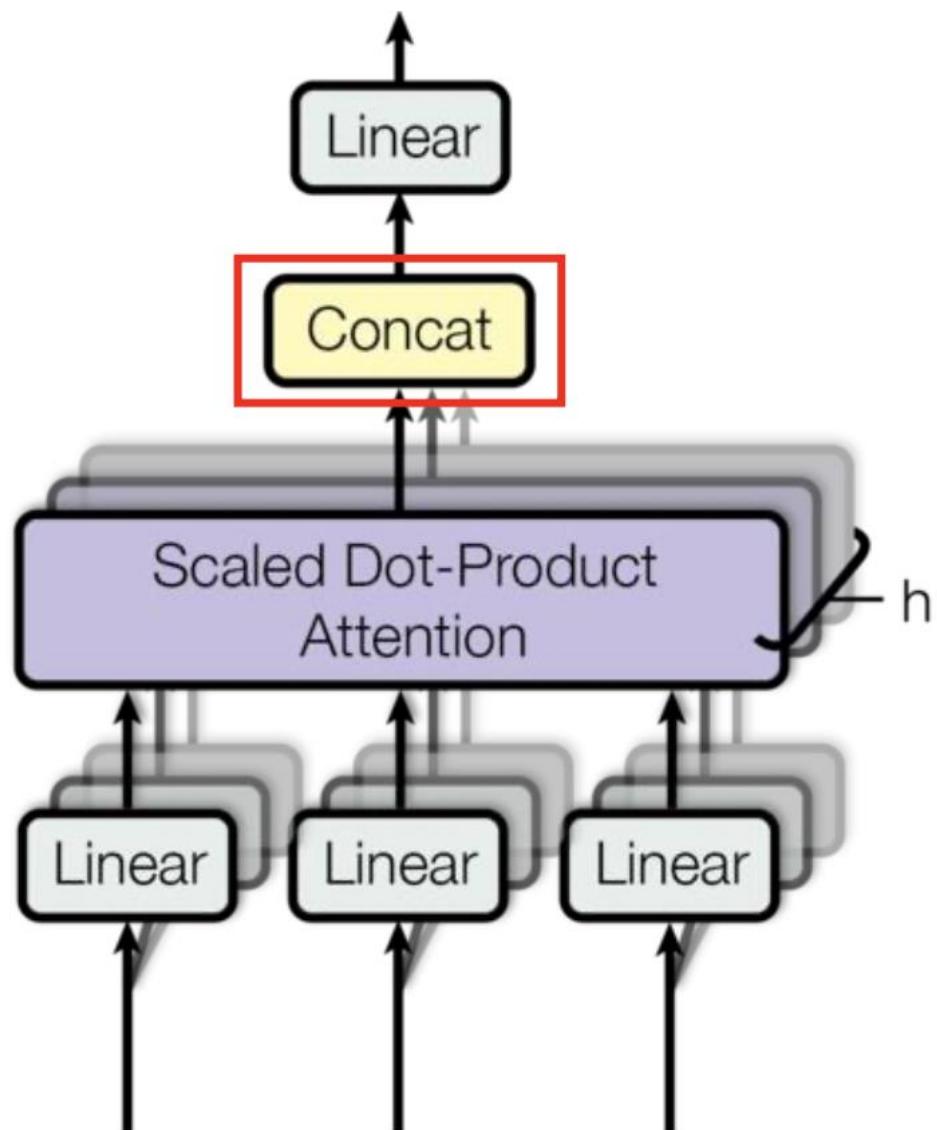
(`batch_size, h, max_len, d_k`)

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

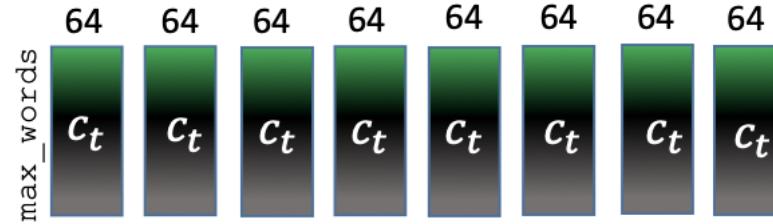


For each head i :

'1 T u2v1 T u3v3 T u4v4'



Attention
Result

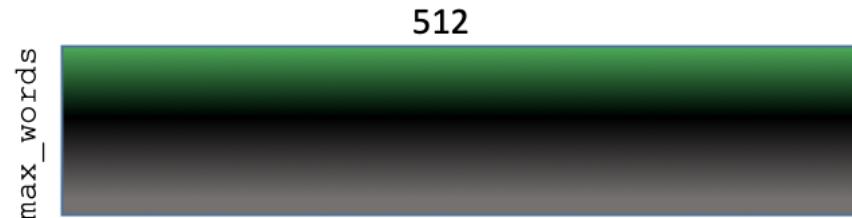


(batch_size, h, max_len, d_k)

So now we have 8 different context vectors

Now concatenate them back again to the original shape:

Final
Attention
Result

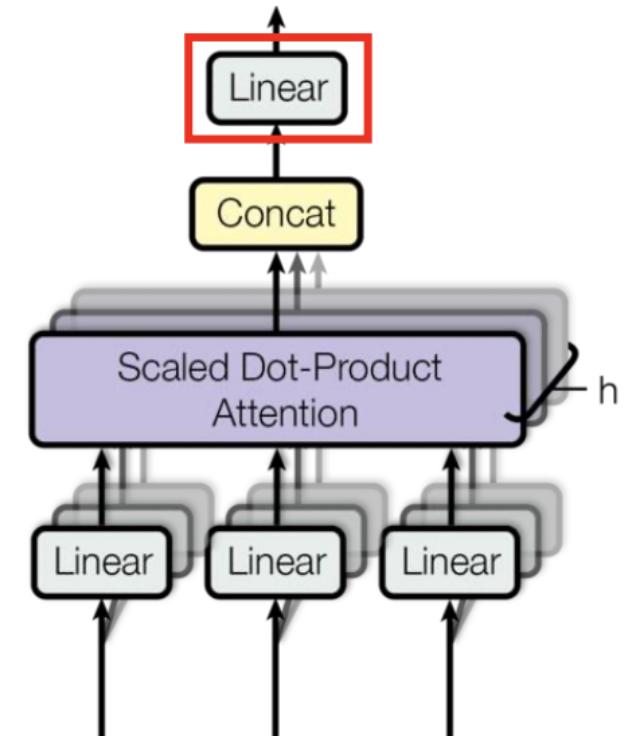
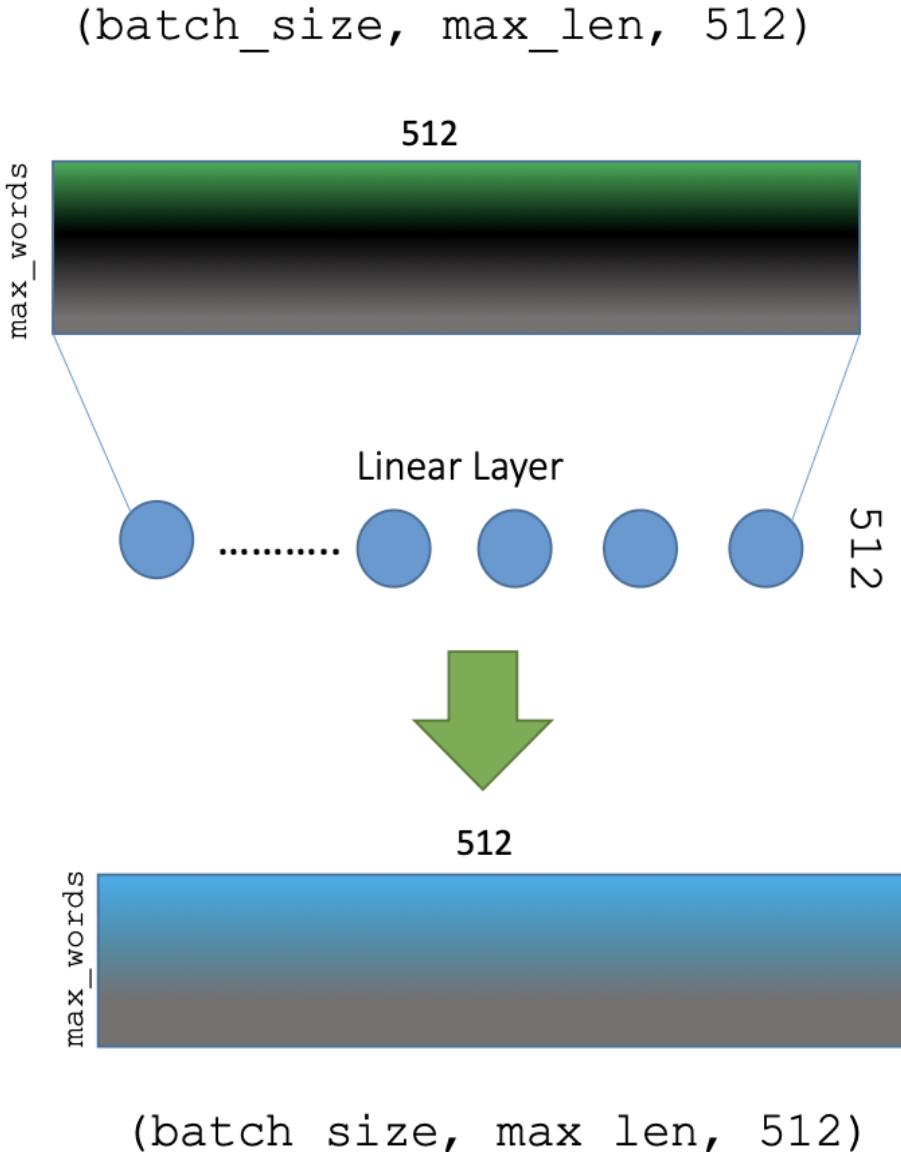
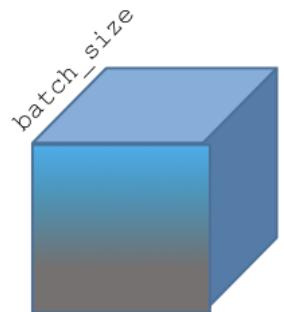


(batch_size, max_len, h * d_k)

(batch_size, max_len, 8 * 64)

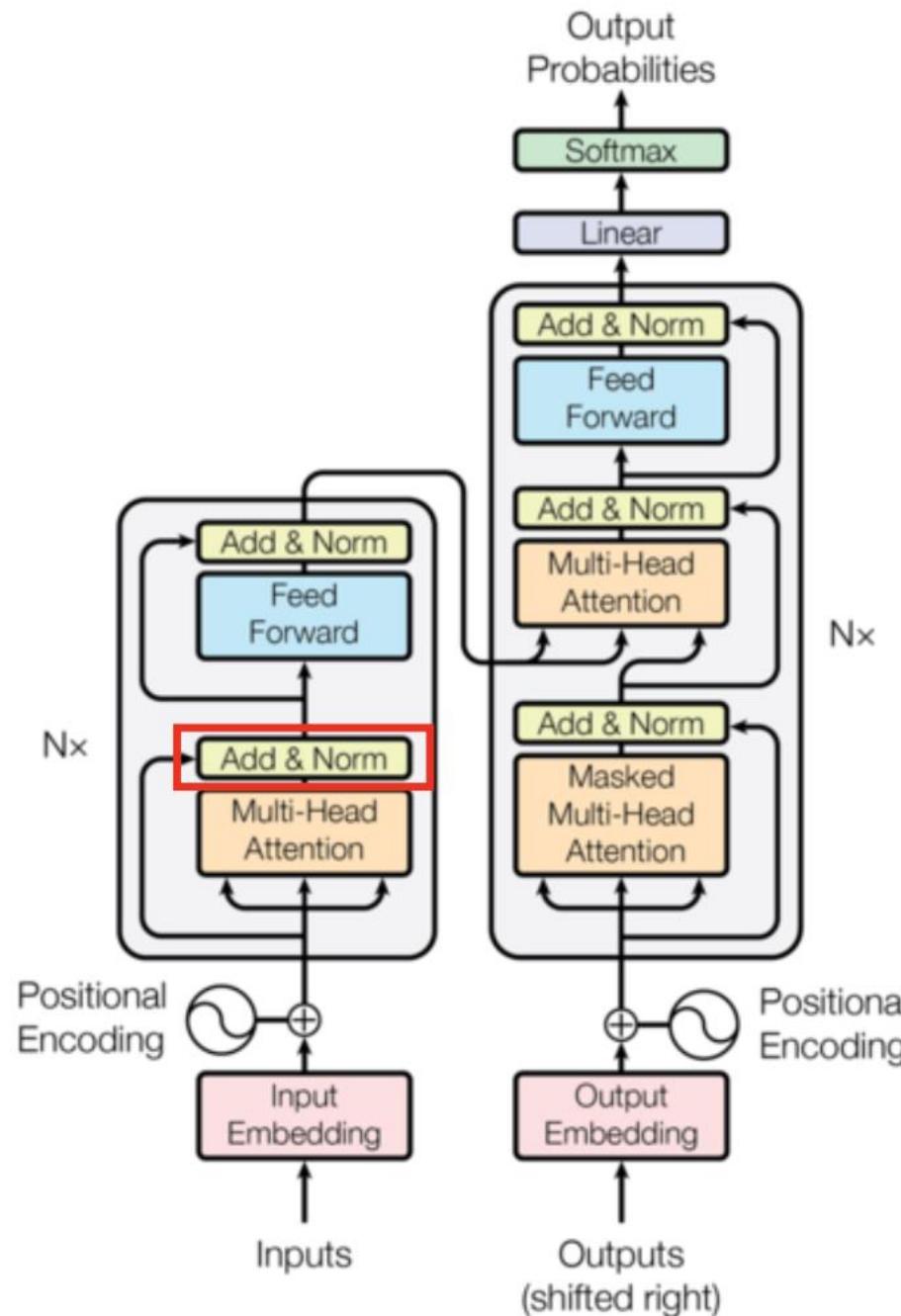
(batch_size, max_len, 512)

Final
Attention
Result



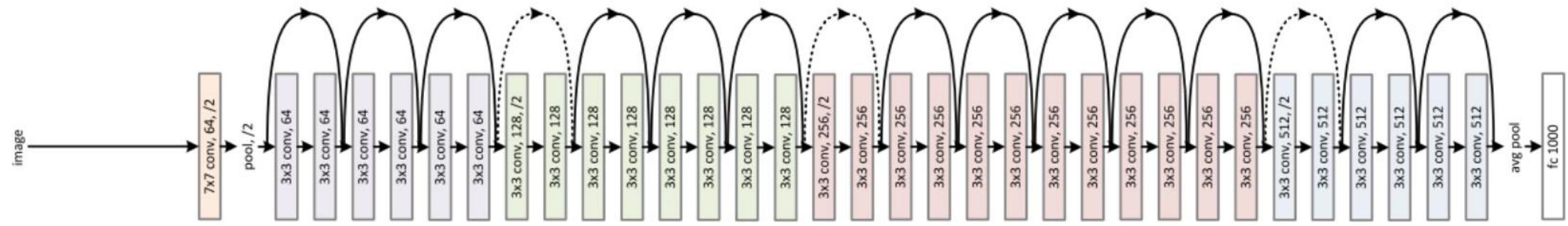
Why do we need a Linear Layer?

We computed 8 different attention results. After we concatenate them, we want them to interact with each other. So we run the concatenated result through a single feed forward (linear) layer.



Residual Networks

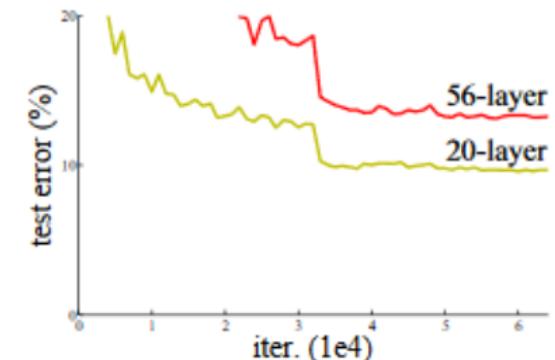
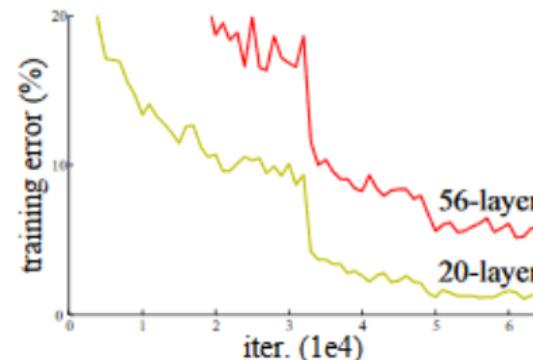
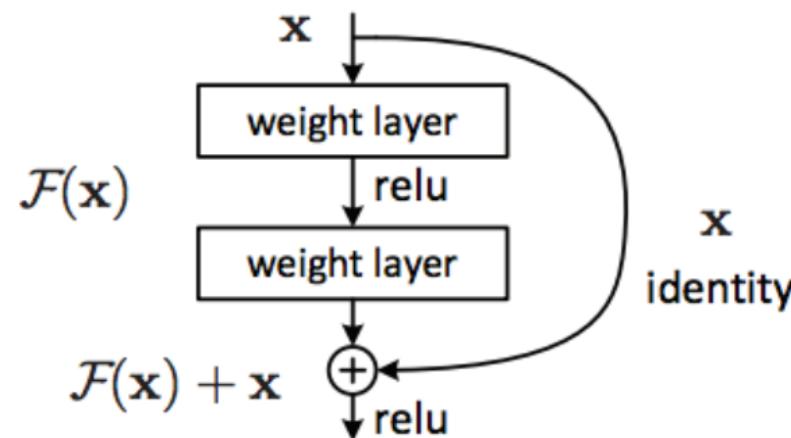
34-layer residual



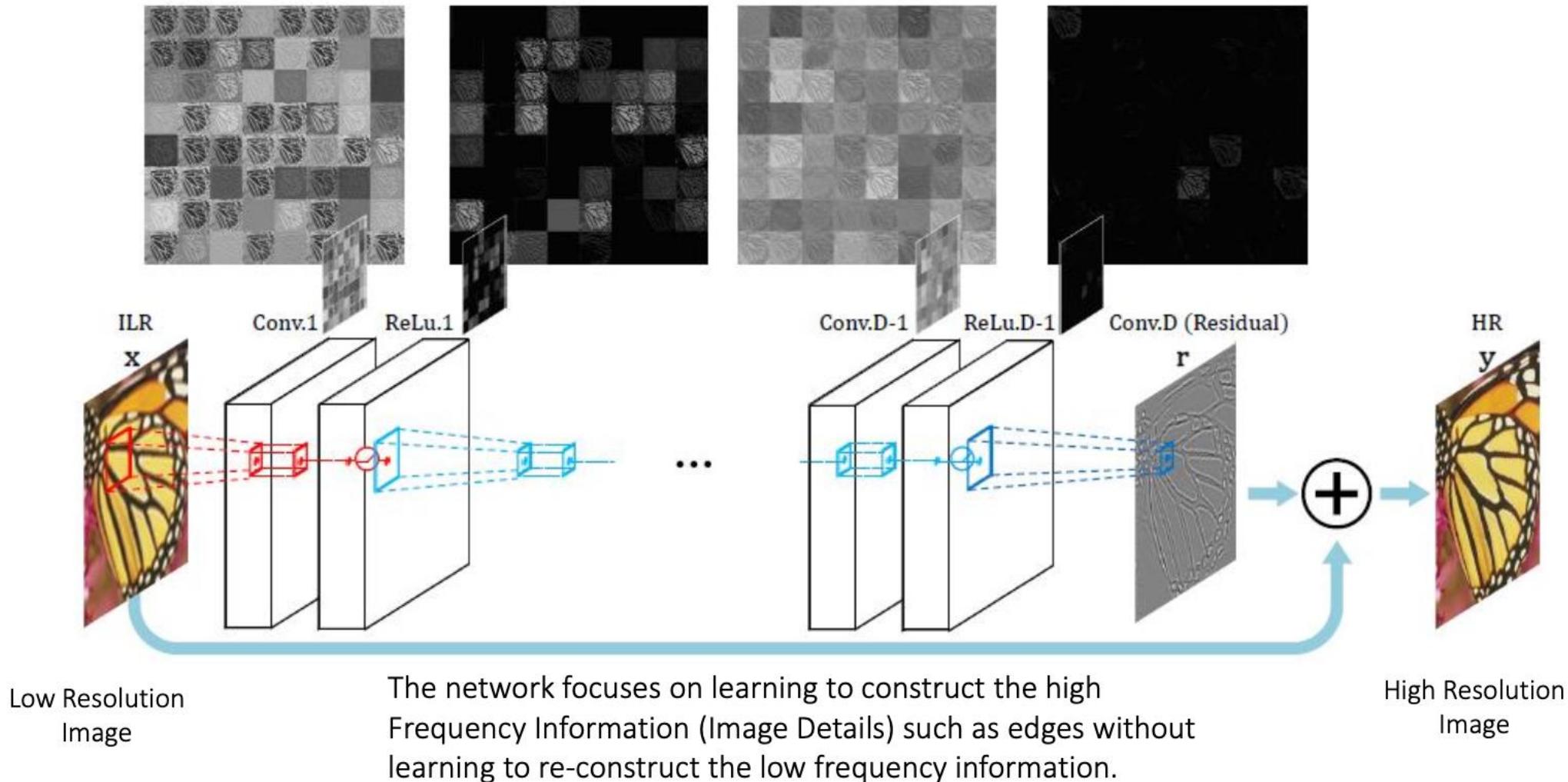
The idea behind Residual Learning

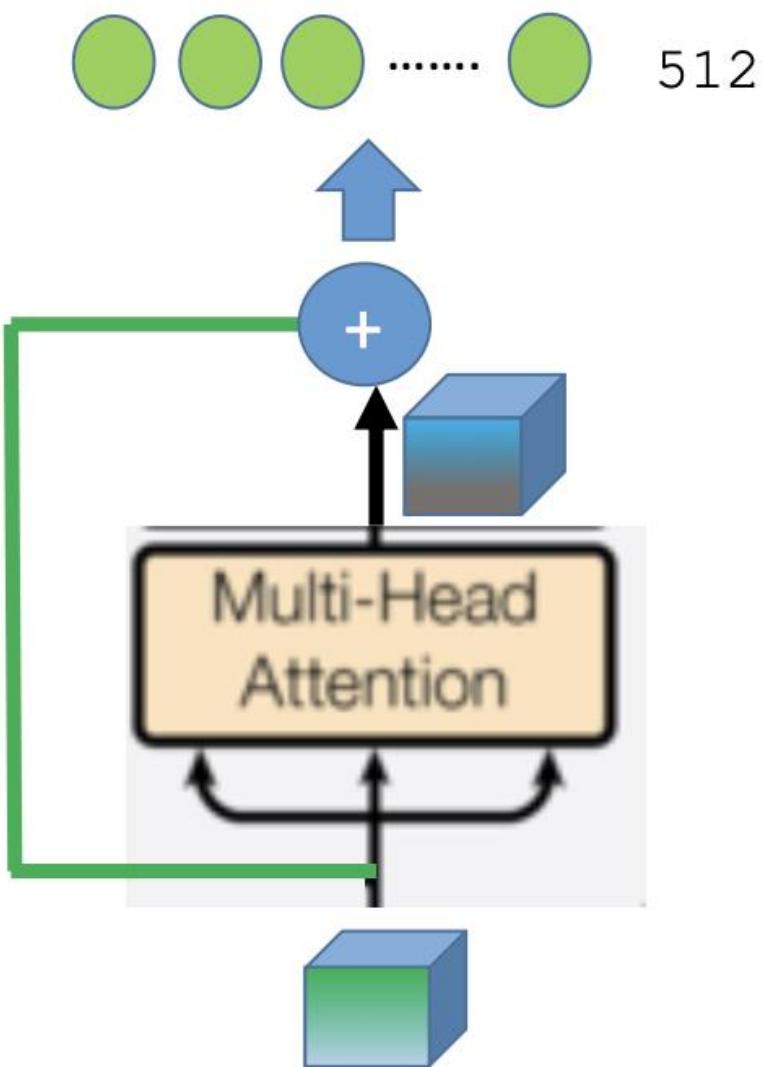
We are learning what's left of (residual), without learning a new representation. You are learning the "remaining" only. If the block doesn't learn anything, then your $F(x)$ would be 0, and that it what makes the training go much faster, since learning a completely new representation is omitted. Therefore, the model can default to using the identity function if the layer is not beneficial.

→ Either learn something useful, or don't learn anything!



A nice way to understand the concept of Residual Learning is through an application





Layer Normalization

Normalize the outputs of a layer

At the output of a feedforward layer:

$$\mu_B \leftarrow \frac{1}{m} \sum_{i=1}^m x_i$$

Mean over all neurons in the layer

$$\sigma_B^2 \leftarrow \frac{1}{m} \sum_{i=1}^m (x_i - \mu_B)^2$$

Variance over all neurons in the layer

For each neuron in the output layer:

$$\hat{x}_i \leftarrow \frac{x_i - \mu_B}{\sqrt{\sigma_B^2 + \epsilon}}$$

Normalize: subtract mean and divide by standard deviation

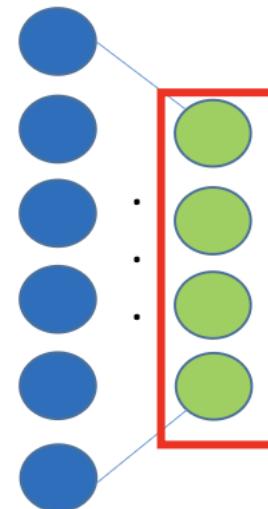
$$y_i \leftarrow \gamma \hat{x}_i + \beta$$

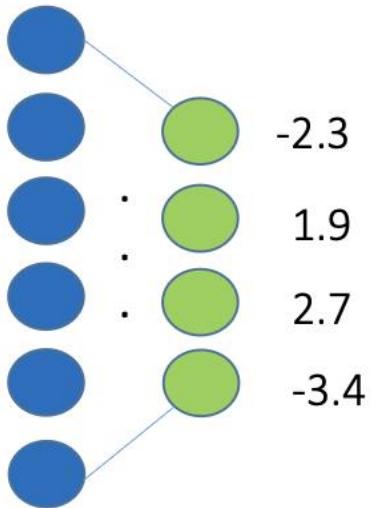
Scale by a learnable value and shift by a learnable value

Learnable Parameters

Why do we need Layer Normalization?

In order to prevent the values of the outputs from becoming bigger. We have performed a lot of operations which may cause the values of the layer output to become bigger. So we use Layer Norm to normalize them back again.



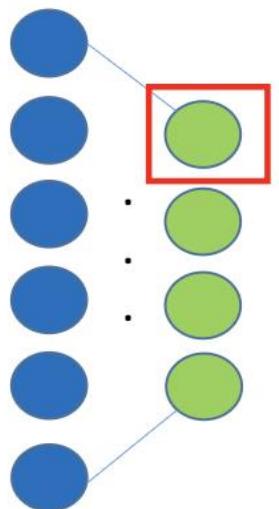


$$\mu_{\mathcal{B}} \leftarrow \frac{1}{m} \sum_{i=1}^m x_i$$

$$\frac{-2.3 + 1.9 + 2.7 - 3.4}{4} = -0.275$$

$$\sigma_{\mathcal{B}}^2 \leftarrow \frac{1}{m} \sum_{i=1}^m (x_i - \mu_{\mathcal{B}})^2$$

$$\frac{(-2.3 + 0.275)^2 + (1.9 + 0.275)^2 + (2.7 + 0.275)^2 + (-3.4 + 0.275)^2}{4} = 6.861$$

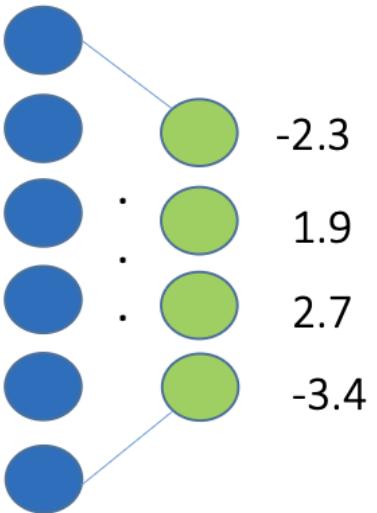


$$\hat{x}_i \leftarrow \frac{x_i - \mu_{\mathcal{B}}}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}}$$

$$\frac{-2.3 + 0.275}{\sqrt{6.861 + 10^{-9}}} = -0.77$$

$$y_i \leftarrow \gamma \hat{x}_i + \beta$$

$$-0.77\theta_1 + \theta_2$$

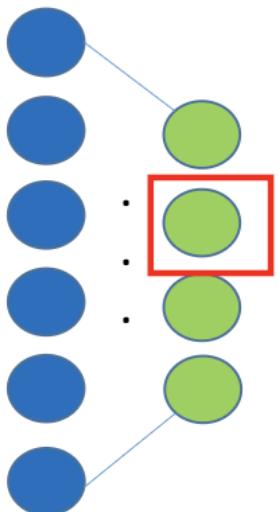


$$\mu_{\mathcal{B}} \leftarrow \frac{1}{m} \sum_{i=1}^m x_i$$

$$\frac{-2.3 + 1.9 + 2.7 - 3.4}{4} = -0.275$$

$$\sigma_{\mathcal{B}}^2 \leftarrow \frac{1}{m} \sum_{i=1}^m (x_i - \mu_{\mathcal{B}})^2$$

$$\frac{(-2.3 + 0.275)^2 + (1.9 + 0.275)^2 + (2.7 + 0.275)^2 + (-3.4 + 0.275)^2}{4} = 6.861$$

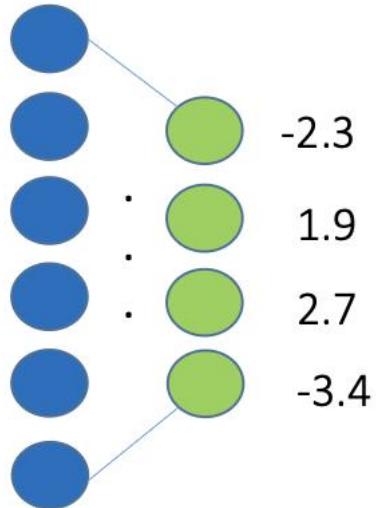


$$\hat{x}_i \leftarrow \frac{x_i - \mu_{\mathcal{B}}}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}}$$

$$\frac{1.9 + 0.275}{\sqrt{6.861 + 10^{-9}}} = 0.83$$

$$y_i \leftarrow \gamma \hat{x}_i + \beta$$

$$0.83\theta_1 + \theta_2$$

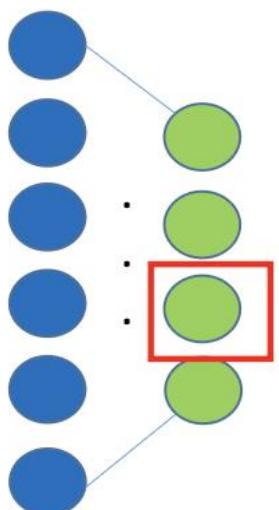


$$\mu_{\mathcal{B}} \leftarrow \frac{1}{m} \sum_{i=1}^m x_i$$

$$\frac{-2.3 + 1.9 + 2.7 - 3.4}{4} = -0.275$$

$$\sigma_{\mathcal{B}}^2 \leftarrow \frac{1}{m} \sum_{i=1}^m (x_i - \mu_{\mathcal{B}})^2$$

$$\frac{(-2.3 + 0.275)^2 + (1.9 + 0.275)^2 + (2.7 + 0.275)^2 + (-3.4 + 0.275)^2}{4} = 6.861$$

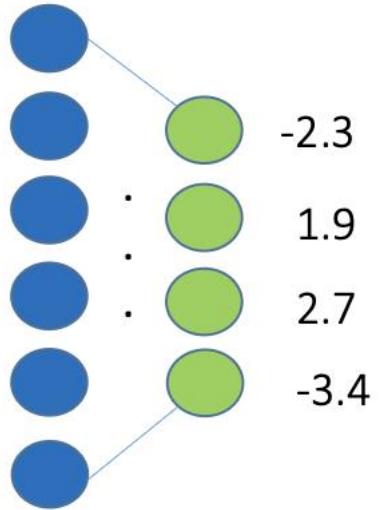


$$\hat{x}_i \leftarrow \frac{x_i - \mu_{\mathcal{B}}}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}}$$

$$\frac{2.7 + 0.275}{\sqrt{6.861 + 10^{-9}}} = 1.13$$

$$y_i \leftarrow \gamma \hat{x}_i + \beta$$

$$1.13\theta_1 + \theta_2$$

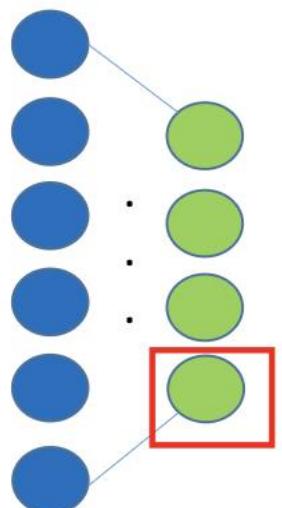


$$\mu_{\mathcal{B}} \leftarrow \frac{1}{m} \sum_{i=1}^m x_i$$

$$\frac{-2.3 + 1.9 + 2.7 - 3.4}{4} = -0.275$$

$$\sigma_{\mathcal{B}}^2 \leftarrow \frac{1}{m} \sum_{i=1}^m (x_i - \mu_{\mathcal{B}})^2$$

$$\frac{(-2.3 + 0.275)^2 + (1.9 + 0.275)^2 + (2.7 + 0.275)^2 + (-3.4 + 0.275)^2}{4} = 6.861$$

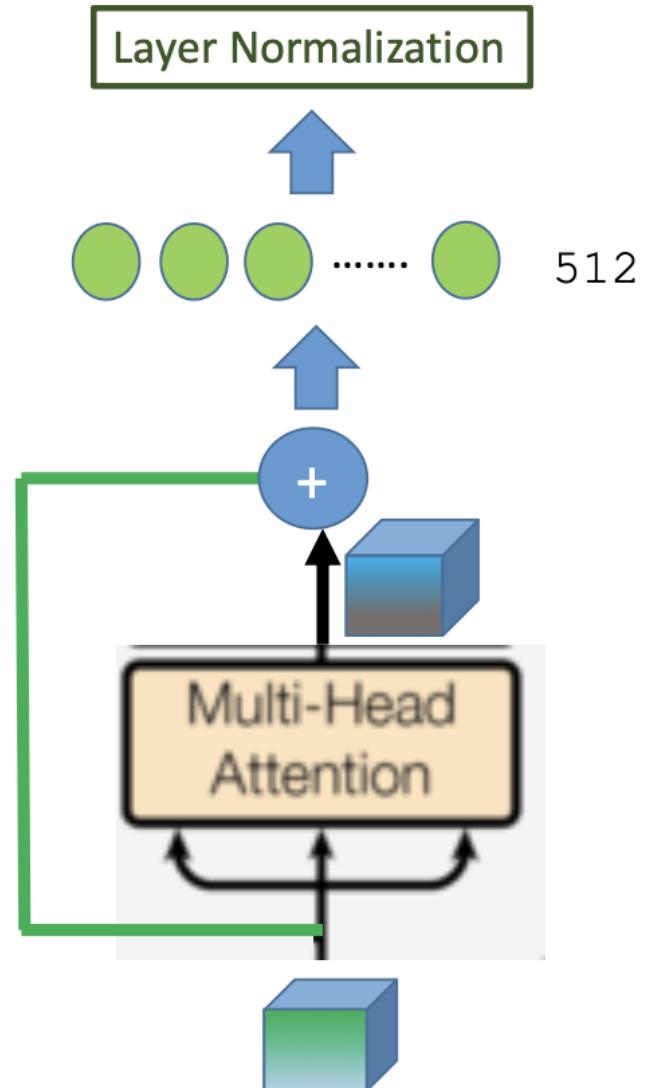


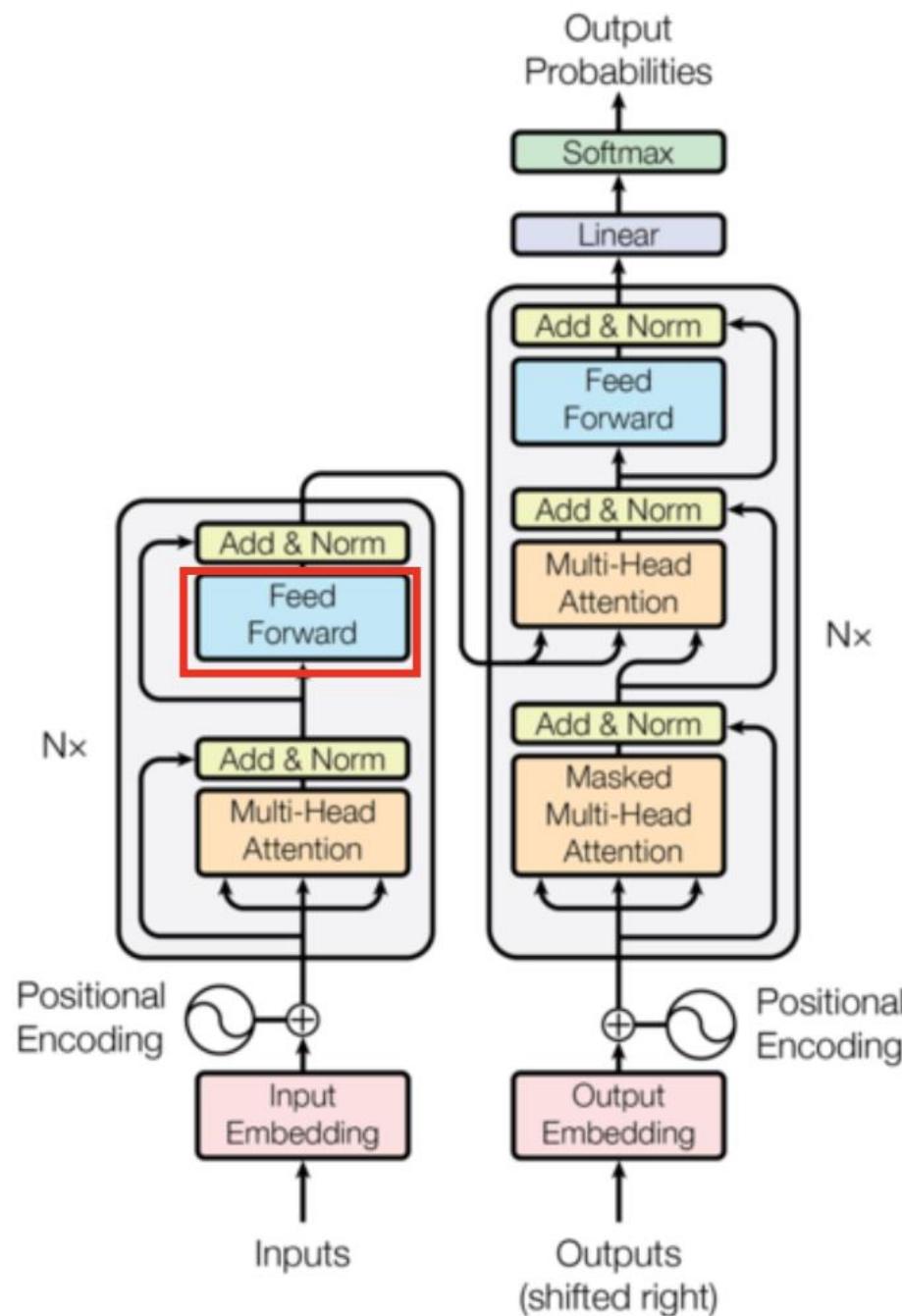
$$\hat{x}_i \leftarrow \frac{x_i - \mu_{\mathcal{B}}}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}}$$

$$\frac{-3.4 + 0.275}{\sqrt{6.861 + 10^{-9}}} = -1.19$$

$$y_i \leftarrow \gamma \hat{x}_i + \beta$$

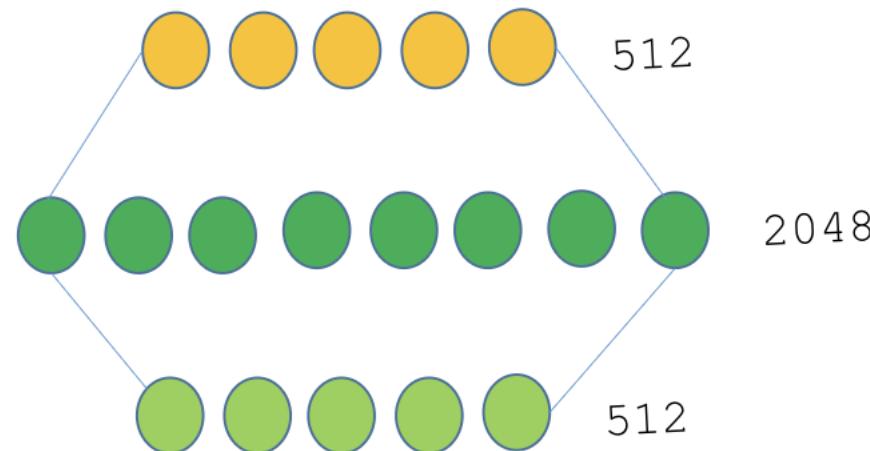
$$-1.19\theta_1 + \theta_2$$

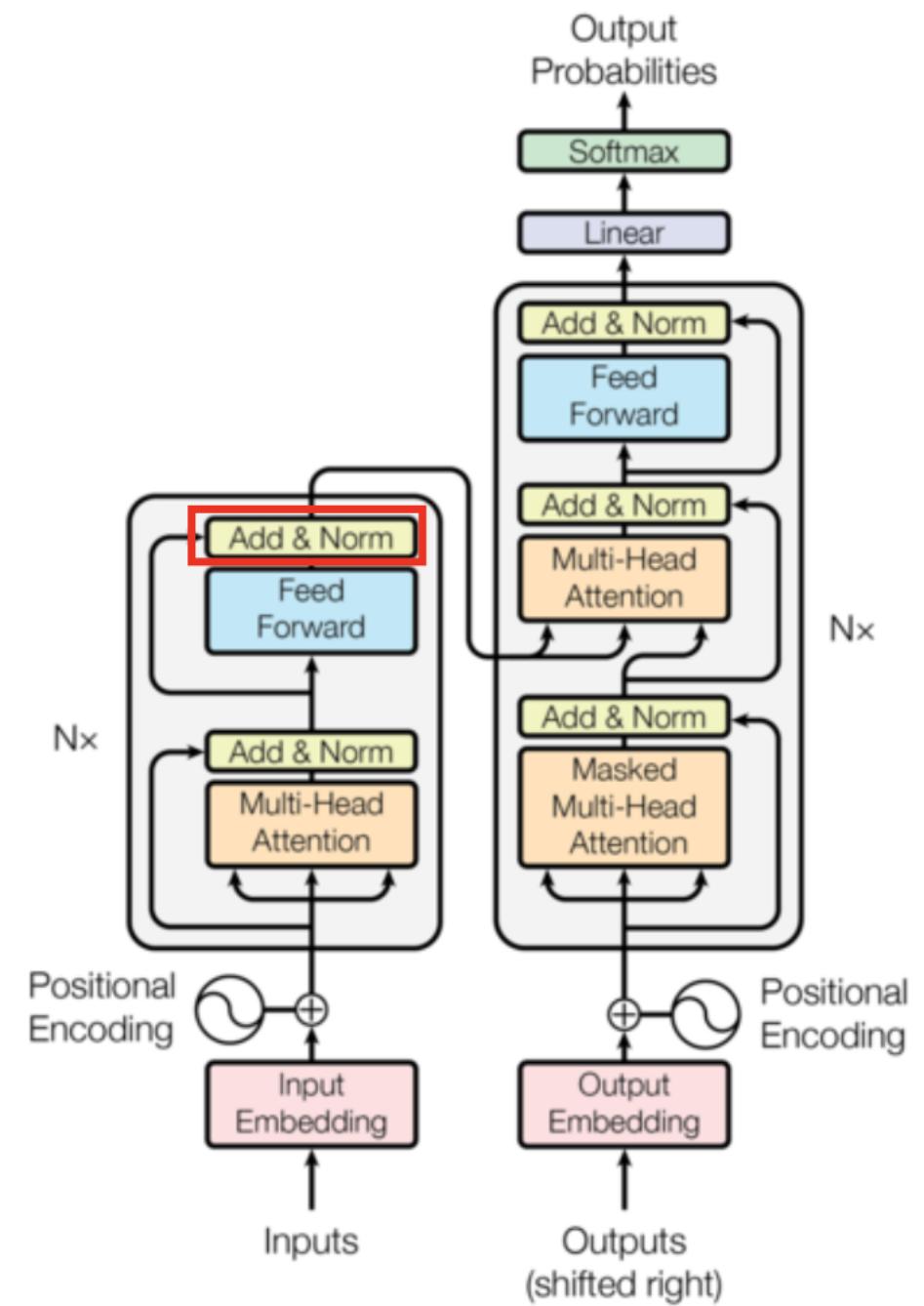
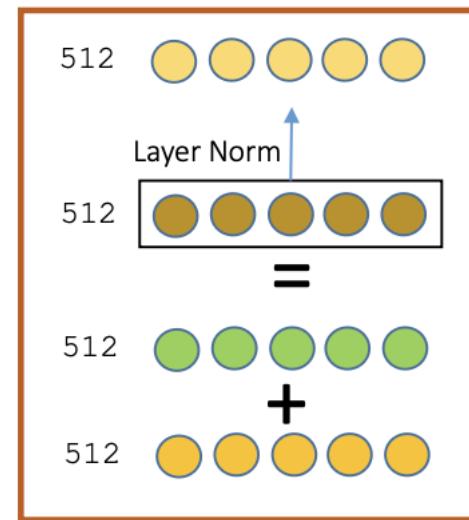
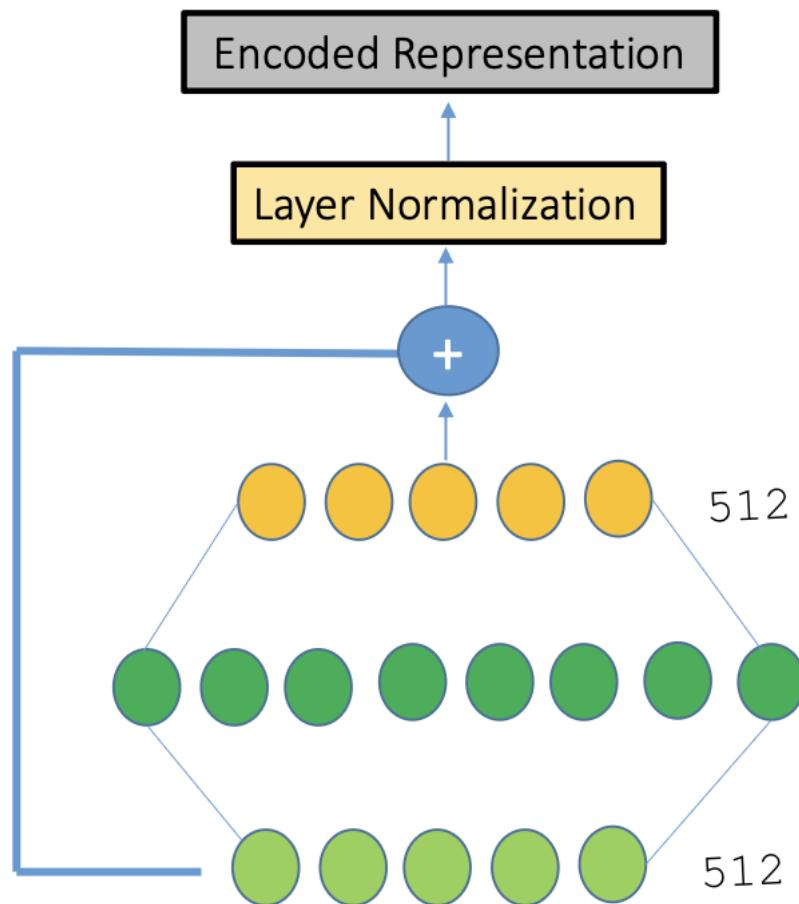




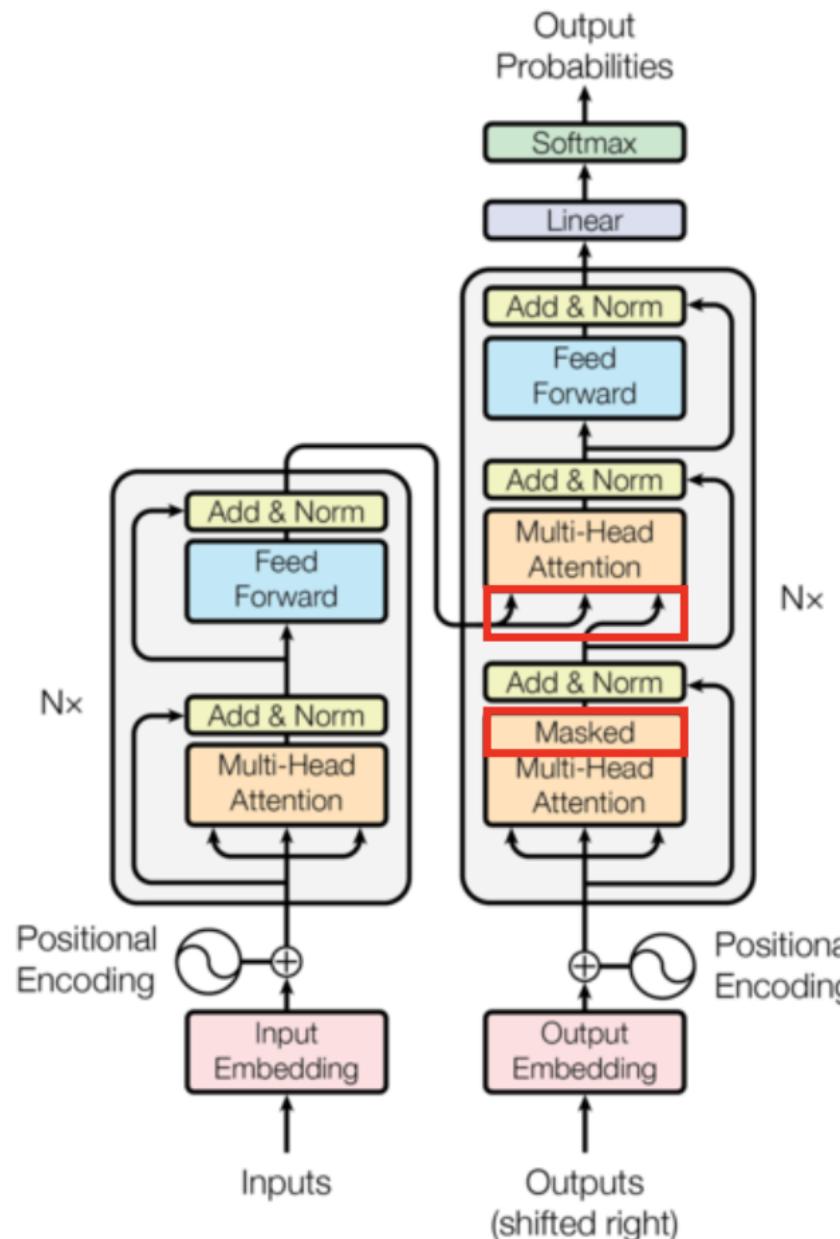
Feed Forward

- 2 linear (feed forward) layers to learn more encoding information.
- The dimensions set in the paper are $512 \rightarrow 2048$ for the first layer, and $2048 \rightarrow 512$ for the second layer.





Decoder: Same as encoder with 2 slight differences



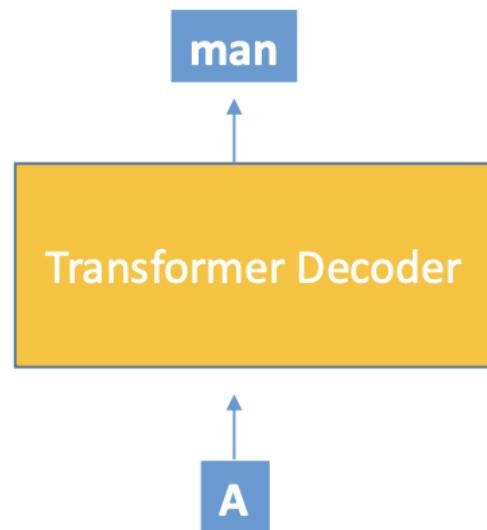
Masked Multi-Head Attention: *Subsequent Mask*

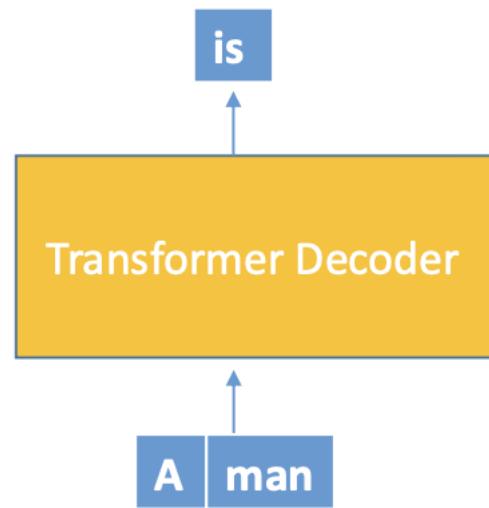
- For self-attention, we don't want our decoder to attend to future word. Otherwise, the model will cheat and learn to look at future words. At testing time, we don't have future words! We are predicting one word at a time (autoregression) (running the decoder for a number of timesteps, just like an LSTM at testing time). So this will be incompatible during testing (inference). Therefore, the decoder is only allowed to attend to earlier positions. During testing time, it can only attend to what has been generated so far. So we need to resemble the testing time scenario during training as well.

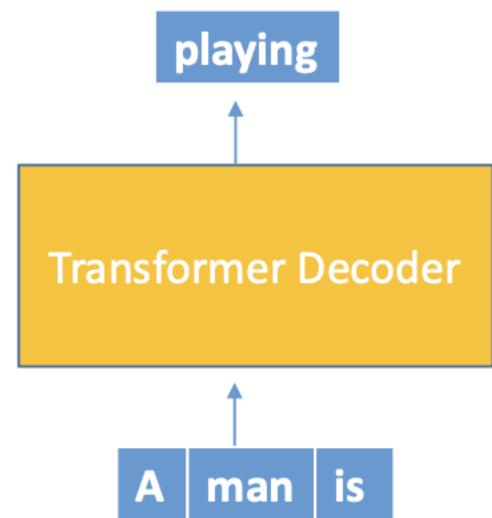
5 words
(5, 5)

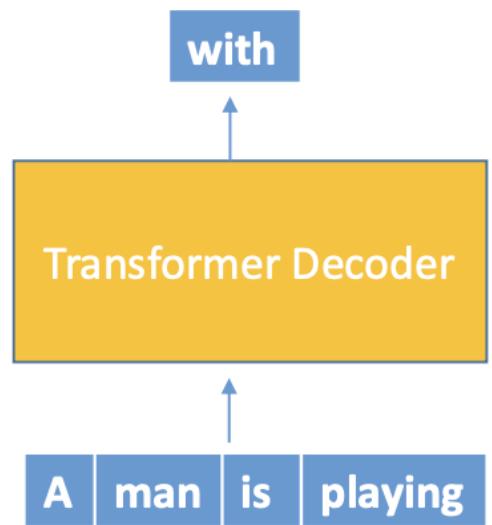
Word 1 `[[[1, 0, 0, 0, 0],`
Word 2 `[1, 1, 0, 0, 0],`
Word 3 `[1, 1, 1, 0, 0],`
Word 4 `[1, 1, 1, 1, 0],`
Word 5 `[1, 1, 1, 1, 1]]]`.

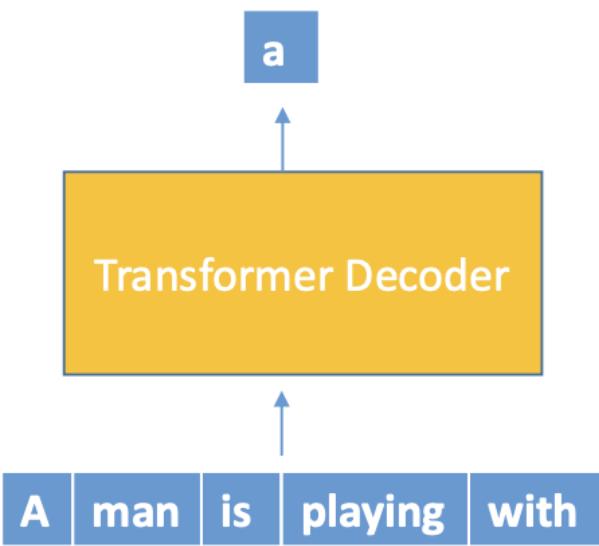
At testing time

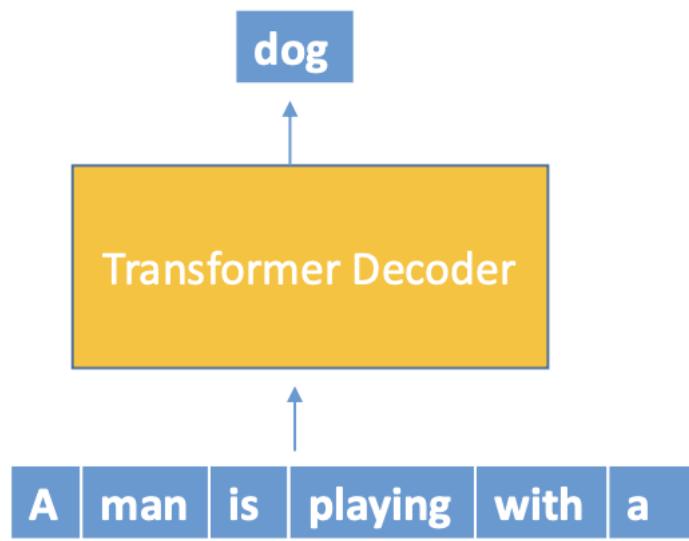


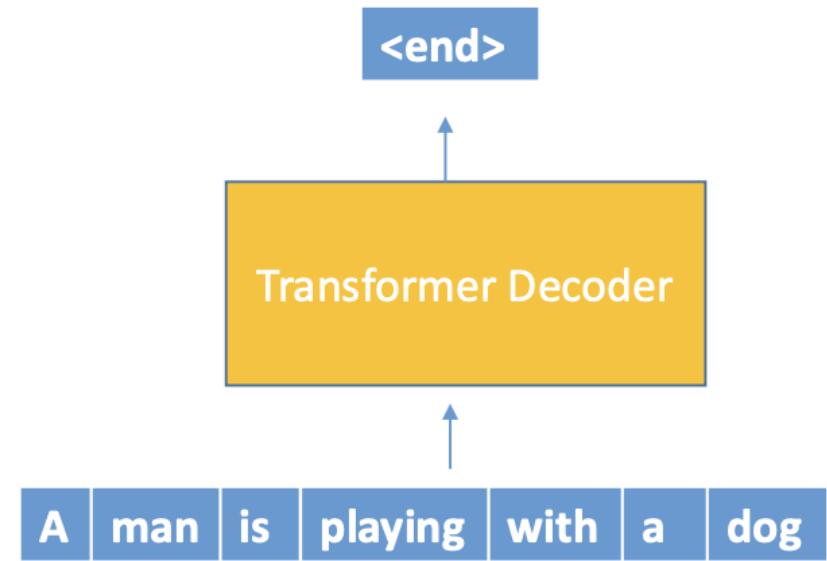










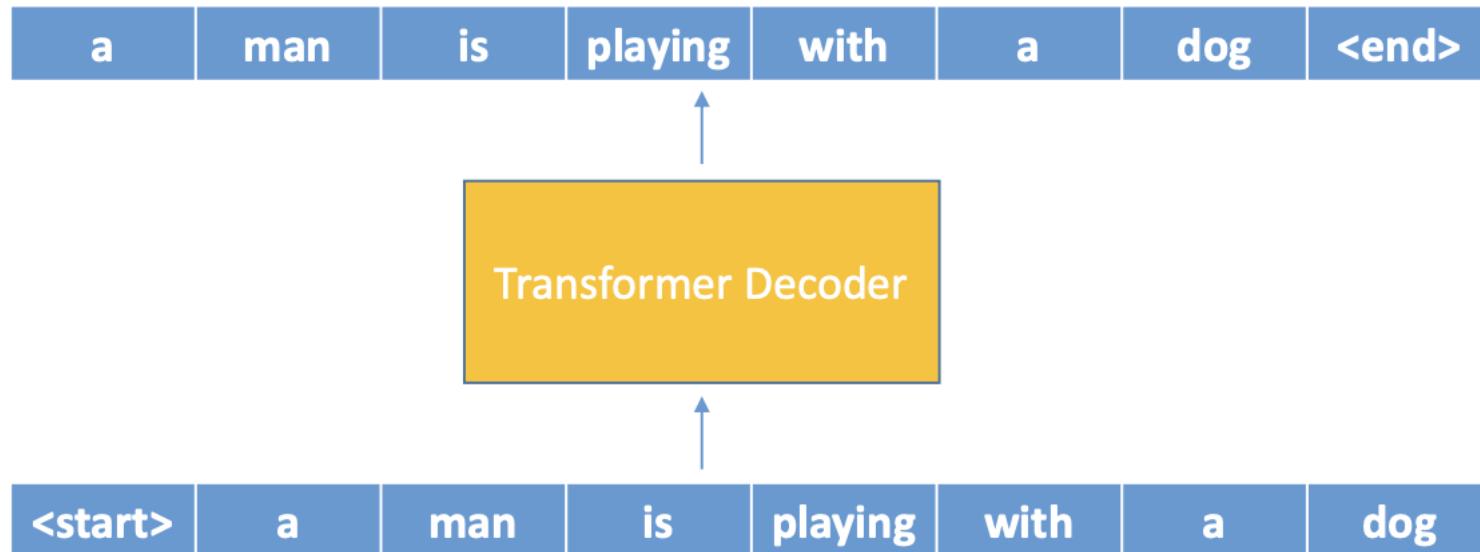


At training time

Sentence: <start> A man is playing with a dog <end>

Input: <start> A man is playing with a dog

Label: A man is playing with a dog <end>

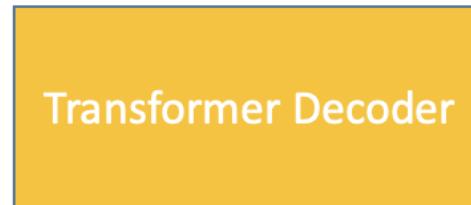


At training time

Sentence: <start> A man is playing with a dog <end>

Input: <start> A man is playing with a dog

Label: A man is playing with a dog <end>



<start>	0	0	0	0	0	0	0
<start>	A	0	0	0	0	0	0
<start>	A	man	0	0	0	0	0
<start>	A	man	is	0	0	0	0
<start>	A	man	is	playing	0	0	0
<start>	A	man	is	playing	with	0	0
<start>	A	man	is	playing	with	a	0
<start>	A	man	is	playing	with	a	dog

Decoder: Queries, Keys and Values of Multi-Head Attention

The Key and Values are the encoded representation we got from the encoder.

The Query is the result we got from self-attention.

This is similar to LSTM attention, except that in LSTM attention, the query comes from decoder LSTM output rather than the decoder self-attention output in a transformer.

The Keys and Values come from the encoder (they are the encoded representation), exactly like in LSTM attention.

