



Master in Computer Vision *Barcelona*

Module M6 : Video Analysis

Lectures 7: Transformers

Lecturer: Ramon Morros

Slides: Javier Ruiz Hidalgo

Outline

- Introduction
- Self-attention
- Multi-head attention
- Deep layers
- Positional encoding
- The Transformer

Motivation for transformers

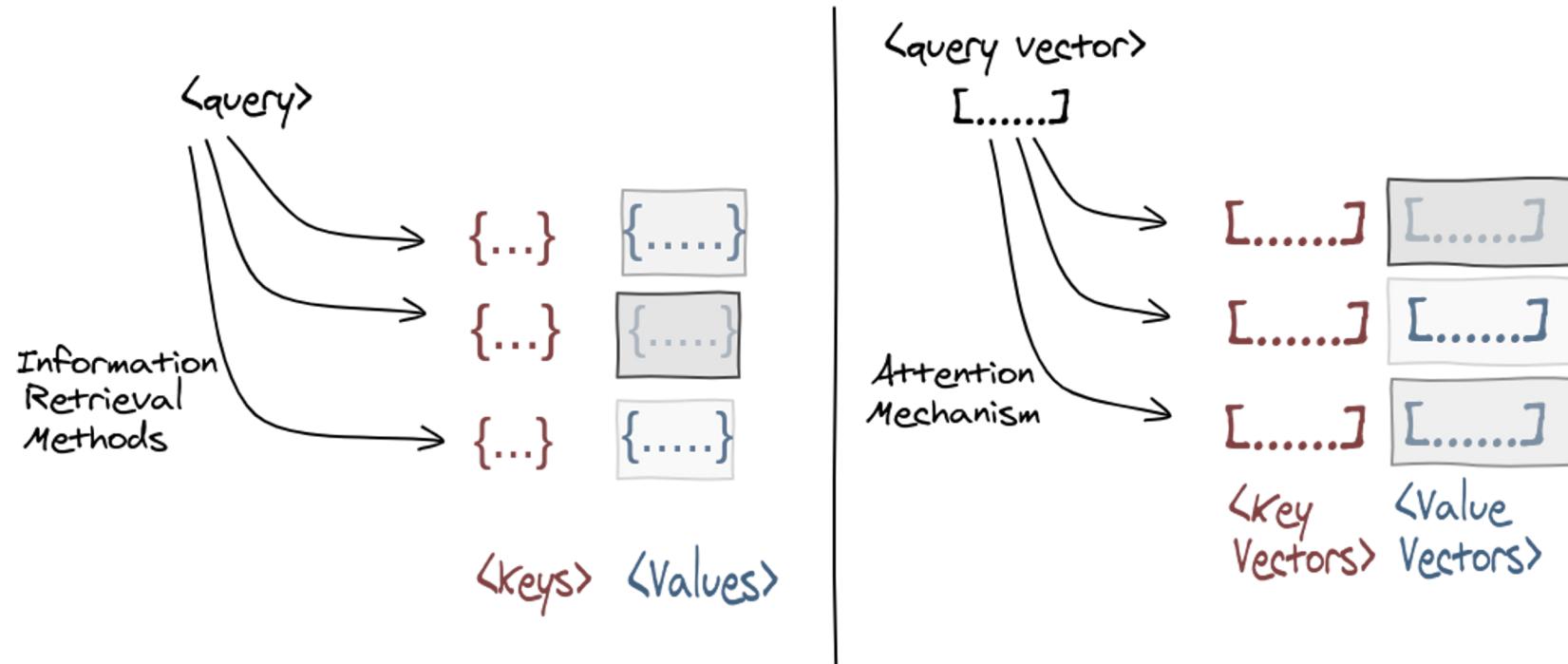
- Despite LSTMs and GRUs, recursive neural networks still need attention mechanism to better deal with dependencies
- But if attention gives us access to any state.. Maybe we don't need the recursion?
- We can go a step forward and ask **is attention all we need?**



- Spoiler alert! Not quite...

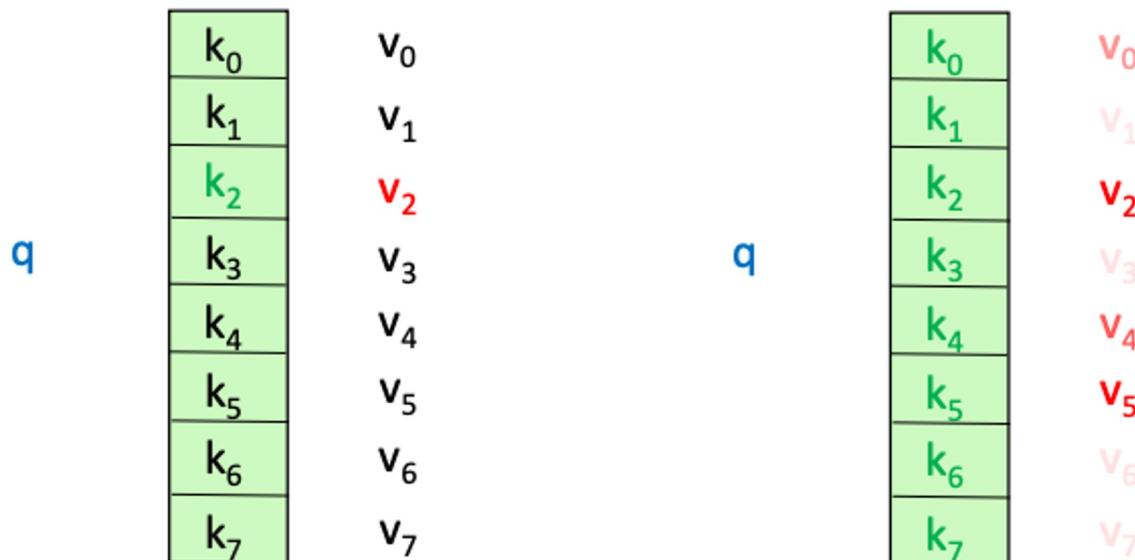
Attention

- Attention is a mechanism to compute a **context vector** for a **query** as a weighted sum of **values**
 - Improves the performance of recurrent neural networks

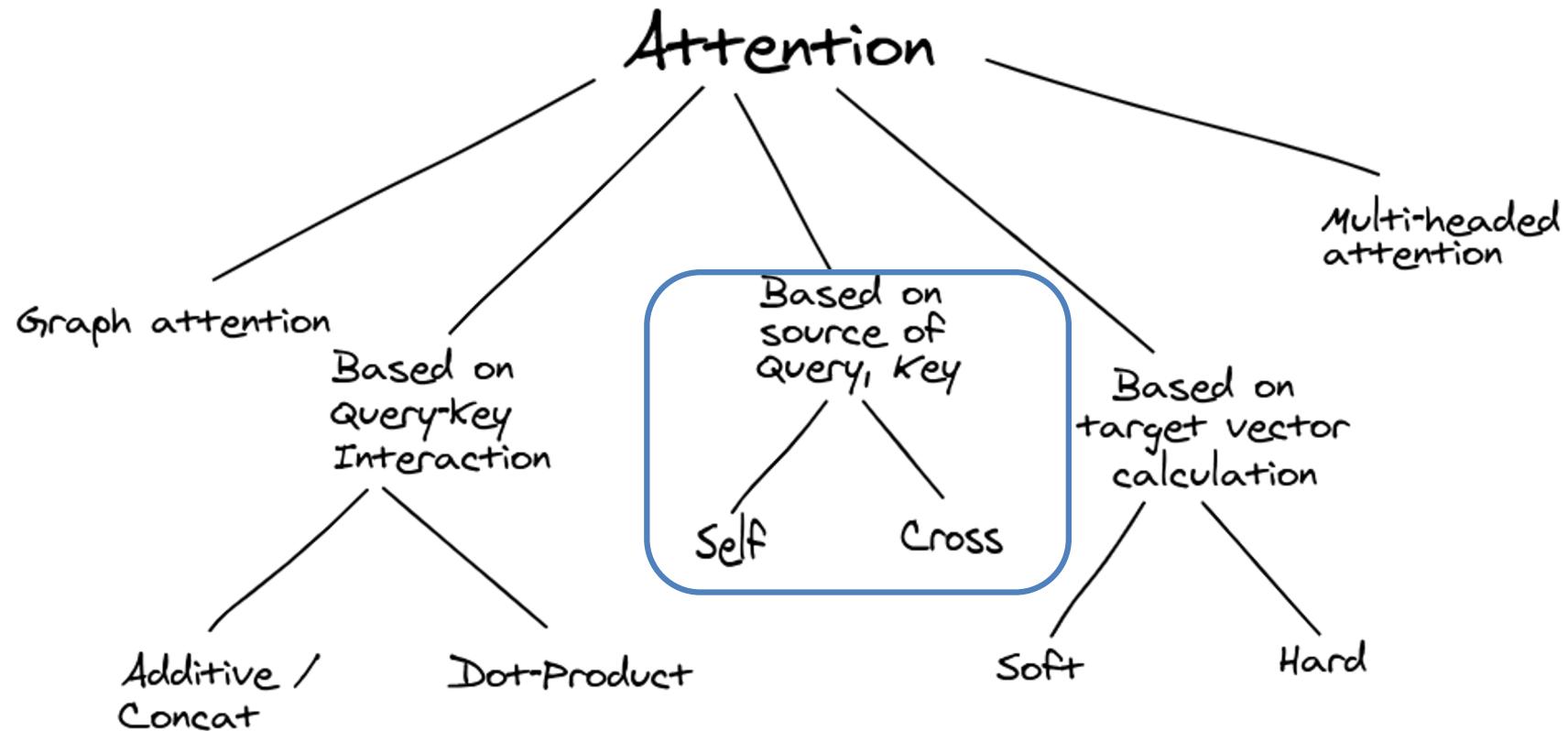


Attention recap (1)

- Let's think of attention as a "fuzzy" or approximate hashtable:
 - To look up a value, we compare a query against keys in a table.
 - In a hashtable (shown on the bottom left):
 - Each query (hash) maps to exactly one key-value pair.
 - In attention (shown on the bottom right):
 - Each query matches each key to varying degrees.
 - We return a sum of values weighted by the query-key match.

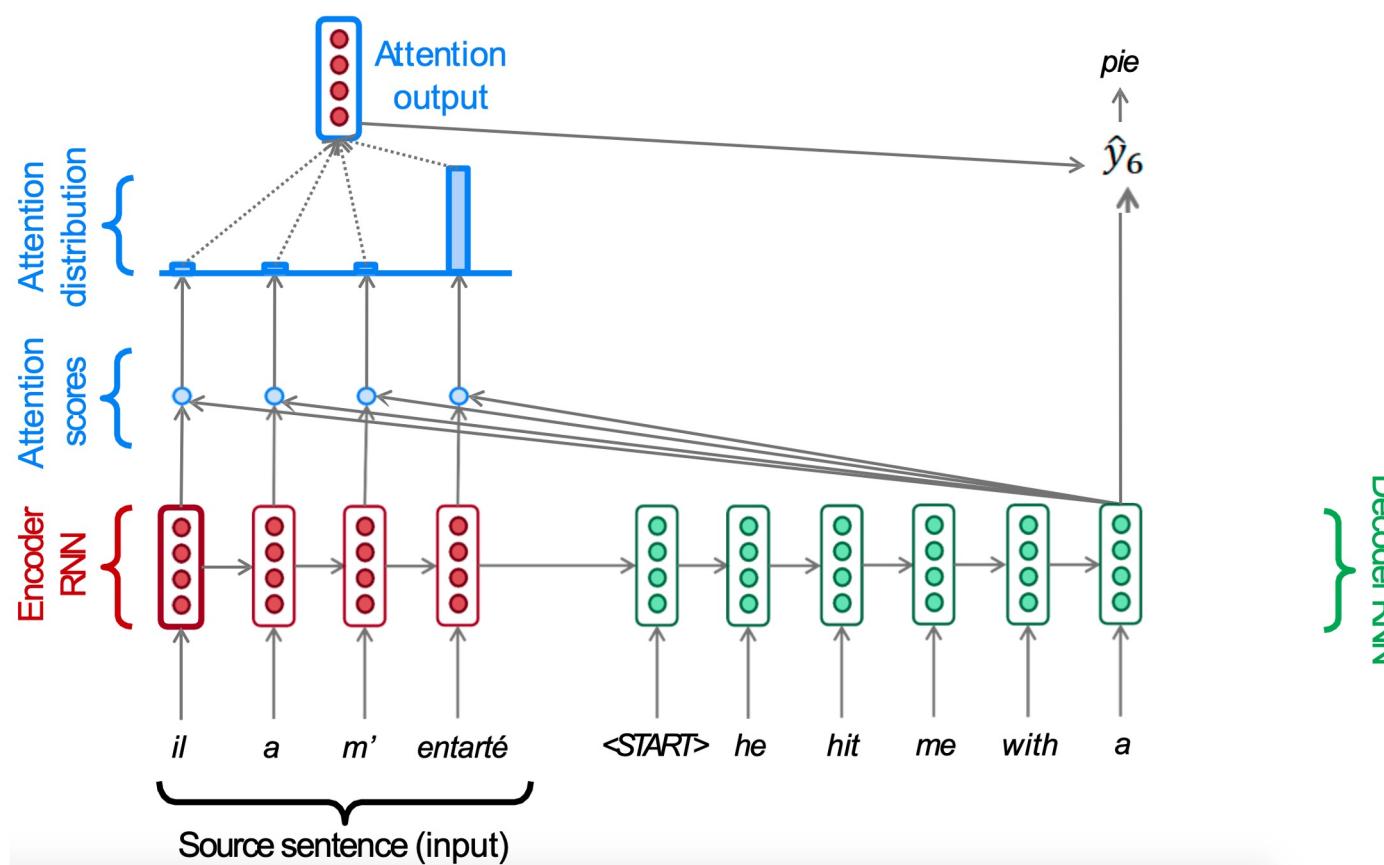


Attention recap (2)



Cross-attention

In this case, **cross-attention** refers to the attention between the encoder and decoder states



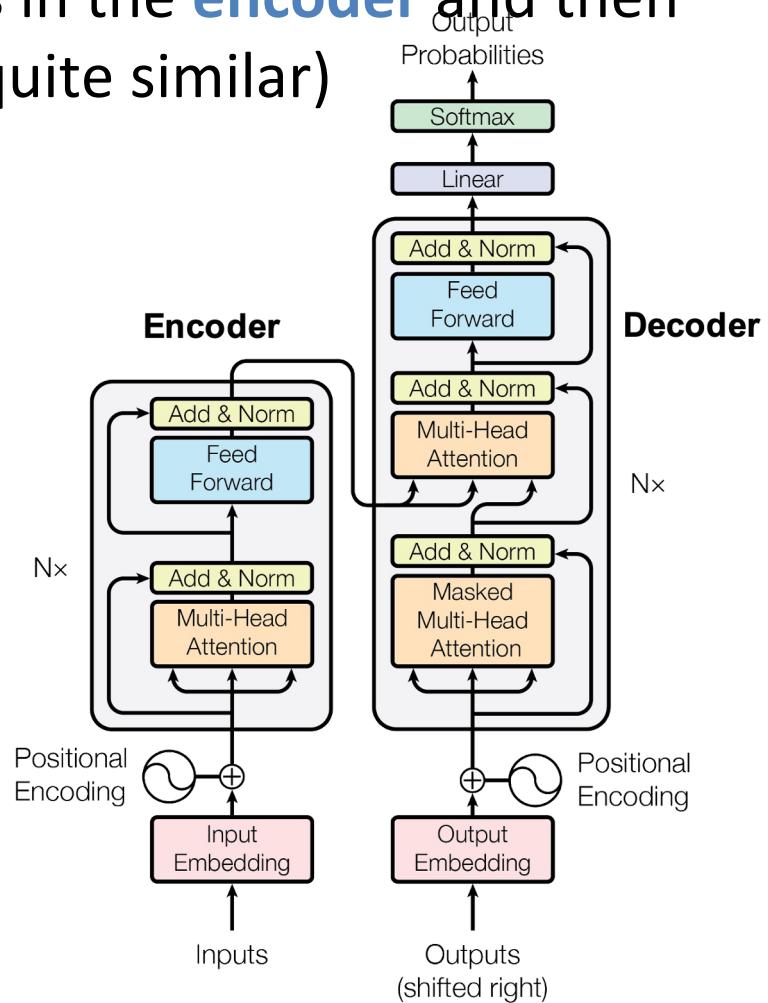
Is attention all you need?

- The **transformer** is a very powerful architecture that exploits attention to the max
- We will first review the main ideas in the **encoder** and then the **decoder** (but fortunately it is quite similar)



Courtesy of Paramount Pictures

A. Vaswani, et al. [Attention is all you need](#). NeurIPS, 2017.

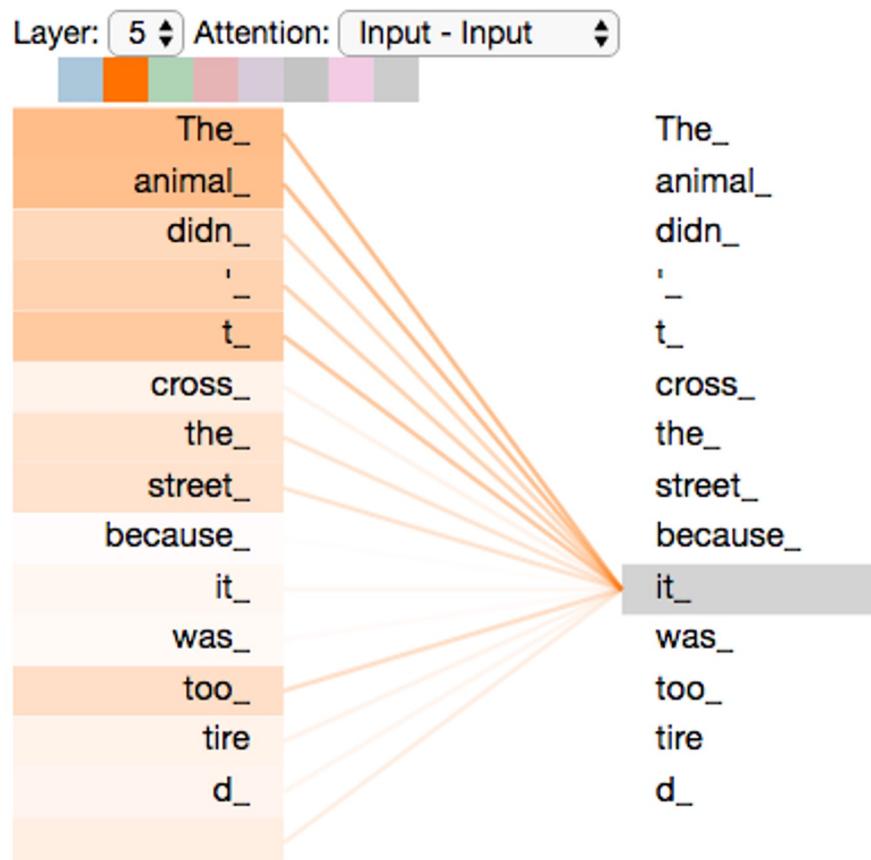


Outline

- Introduction
- **Self-attention**
- Multi-head attention
- Positional encoding
- Deep layers
- The Transformer

Self-attention (1)

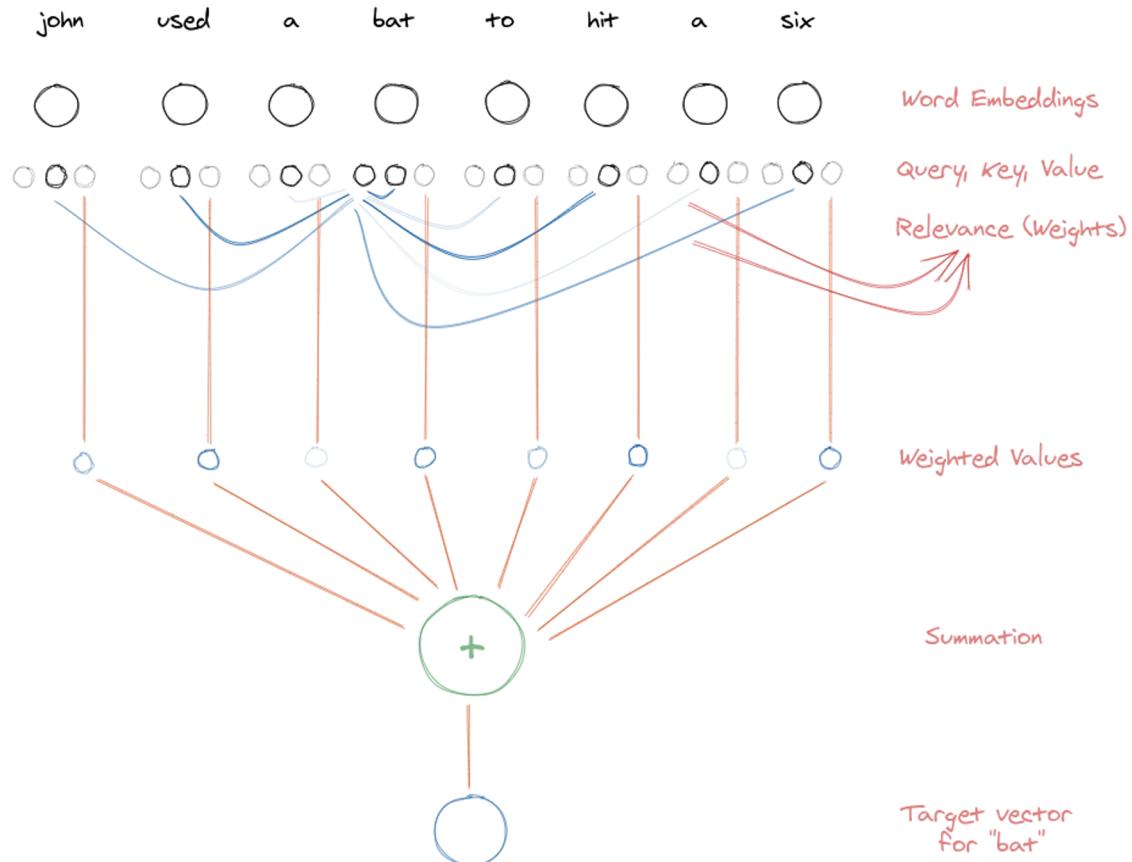
- **Self-attention** refers to attending to other elements from the **same sequence**



Jay Alammar, [The Illustrated Transformer](#), 2018.

Z. Lin, et al. [A structured self-attentive sentence embedding](#), ICLR, 2017.

Self-attention (2)



Queries:

$$q_i = \mathbf{x}_i \mathbf{W}^Q$$

Keys:

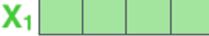
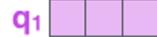
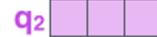
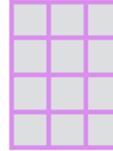
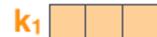
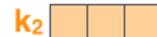
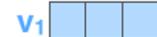
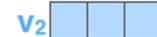
$$k_i = \mathbf{x}_i \mathbf{W}^K$$

Values:

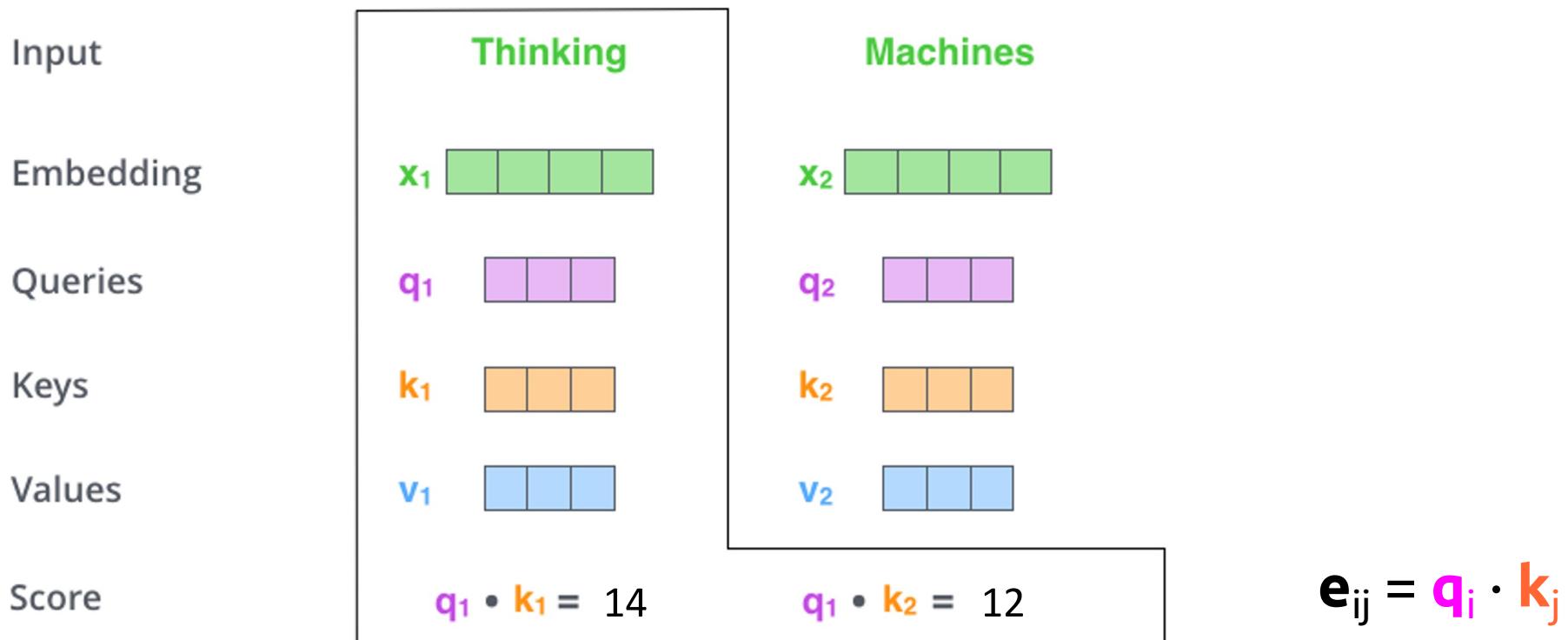
$$v_i = \mathbf{x}_i \mathbf{W}^V$$

\mathbf{W}^Q , \mathbf{W}^K and \mathbf{W}^V are **projection layers shared across all words**

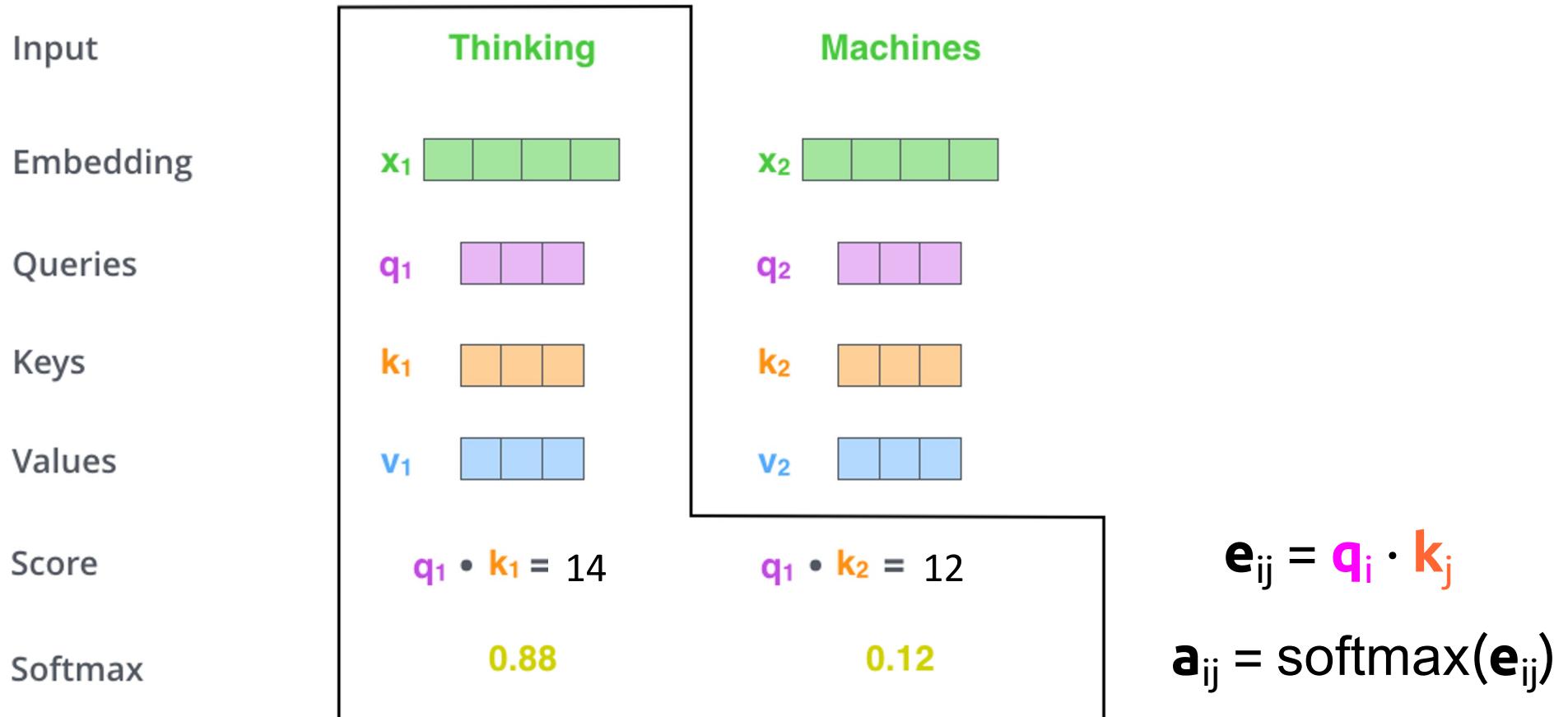
Self-attention (3)

Input	Thinking	Machines		
Embedding	x_1 	x_2 		
Queries	q_1 	q_2 		w^Q $q_i = x_i w^Q$
Keys	k_1 	k_2 		w^K $k_i = x_i w^K$
Values	v_1 	v_2 		w^V $v_i = x_i w^V$

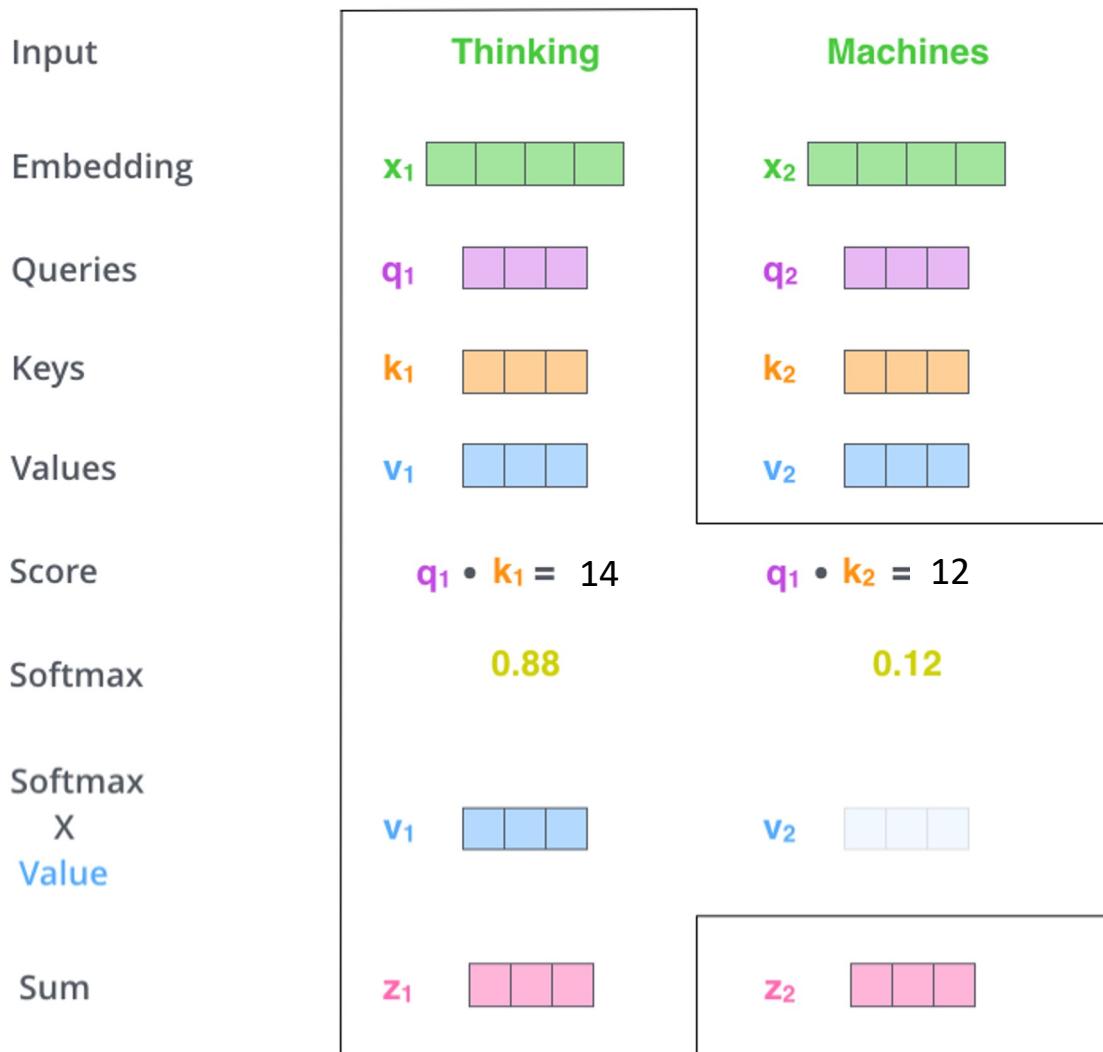
Self-attention (4)



Self-attention (5)



Self-attention (6)



$$\mathbf{e}_{ij} = \mathbf{q}_i \cdot \mathbf{k}_j$$

$$\mathbf{a}_{ij} = \text{softmax}(\mathbf{e}_{ij})$$

$$\mathbf{z}_i = \sum_j \mathbf{a}_{ij} \mathbf{v}_j$$

Self-attention: vectorized steps (1)

1) This is our
input sentence

2) We embed
each word

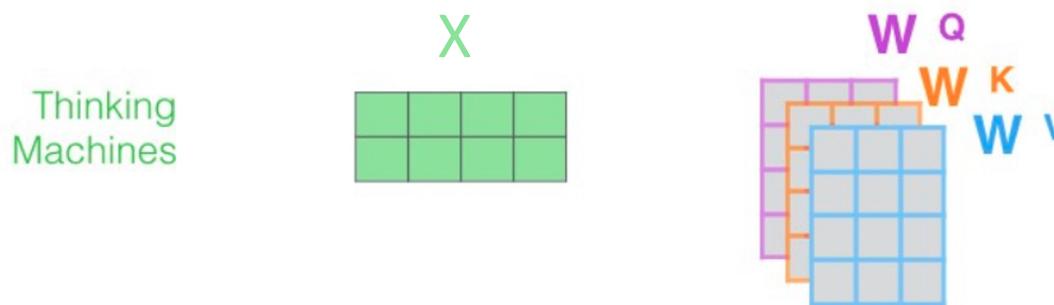


Self-attention: vectorized steps (2)

1) This is our
input sentence

2) We embed
each word

3) We multiply \mathbf{X} with
weight matrices

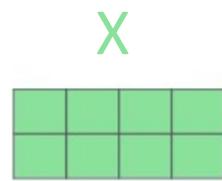


Self-attention: vectorized steps (3)

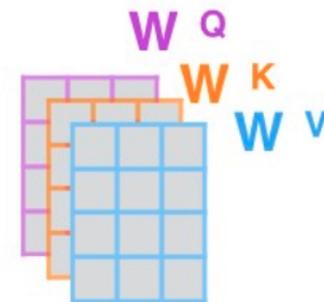
1) This is our input sentence

Thinking
Machines

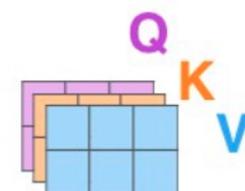
2) We embed each word



3) We multiply X with weight matrices



4) Calculate attention using the resulting $Q/K/V$ matrices

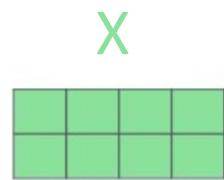


Self-attention: vectorized steps (4)

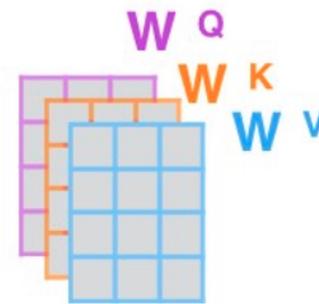
1) This is our input sentence

Thinking
Machines

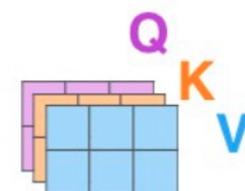
2) We embed each word



3) We multiply X with weight matrices



4) Calculate attention using the resulting $Q/K/V$ matrices



5) Calculate outputs (context)

$$\text{softmax} \left(\begin{matrix} Q \\ \times \\ K^T \end{matrix} \right) V = Z$$

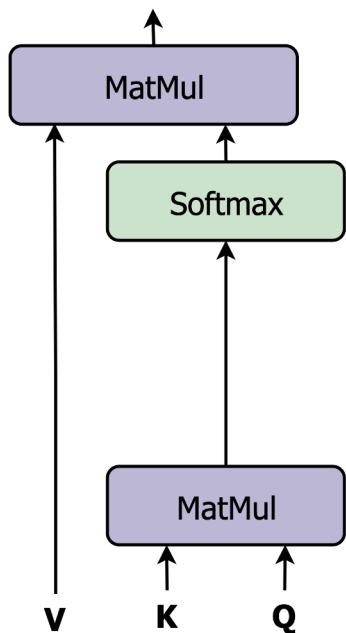
Q K^T V Z

softmax $\left(\begin{matrix} \text{pink} \\ \times \\ \text{orange} \end{matrix} \right) \text{blue} = \text{pink}$

A diagram showing the calculation of context outputs. It starts with three matrices: Q (pink), K^T (orange), and V (blue). These are multiplied together using the softmax function to produce the final output matrix Z (pink).

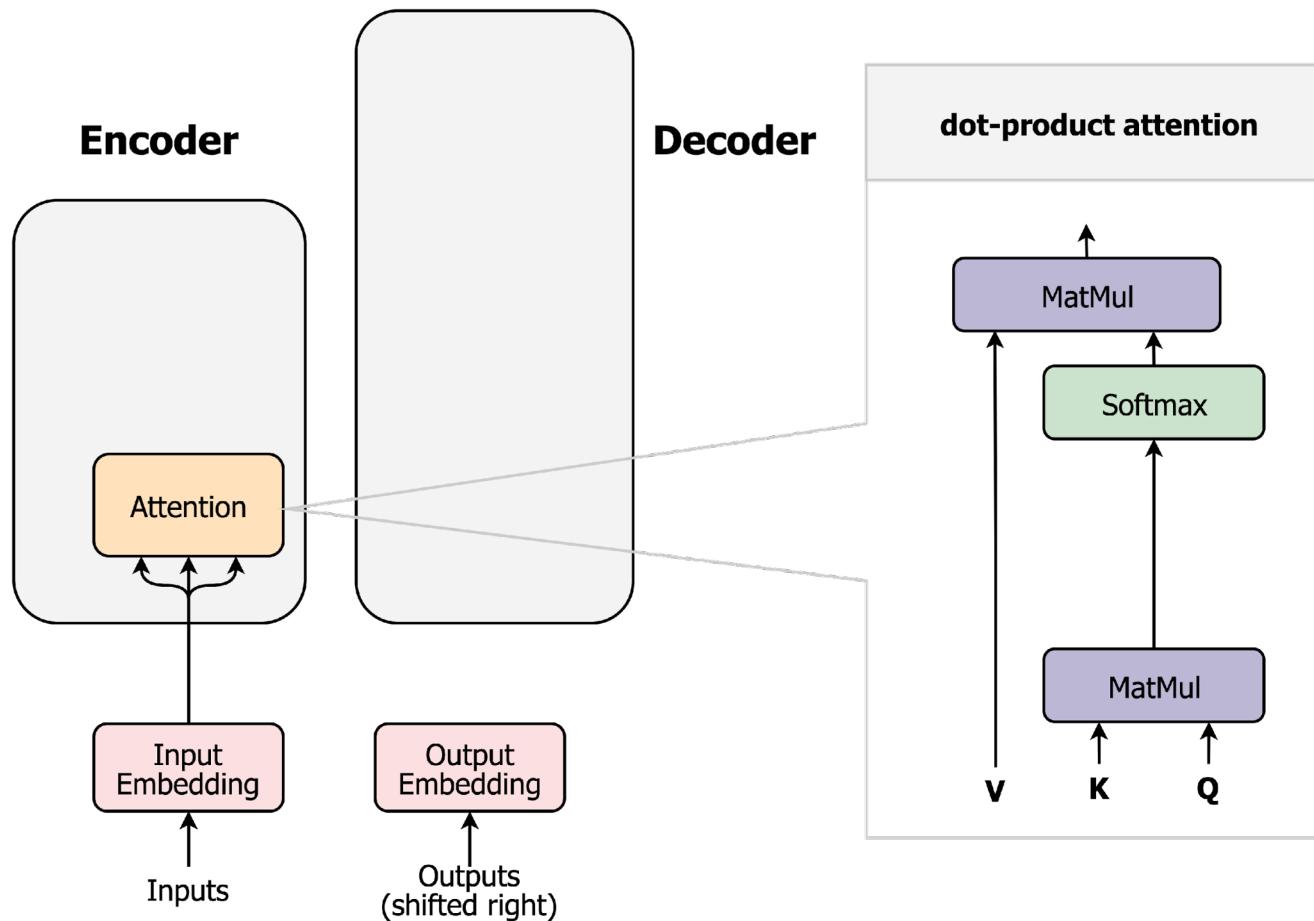
Self-attention in (vectorized) equations

$$\mathbf{Z} = \text{softmax}(\mathbf{Q}\mathbf{K}^T) \quad \mathbf{V} = \text{softmax}(\mathbf{X}\mathbf{W}^Q\mathbf{W}^K\mathbf{X}^T) \quad \mathbf{X}\mathbf{W}^V$$



$$\begin{aligned}\mathbf{W}^Q &\in \mathbb{R}^{d_x \times d_q} \\ \mathbf{W}^K &\in \mathbb{R}^{d_x \times d_k} \quad d_q = d_k \\ \mathbf{W}^V &\in \mathbb{R}^{d_x \times d_v}\end{aligned}$$

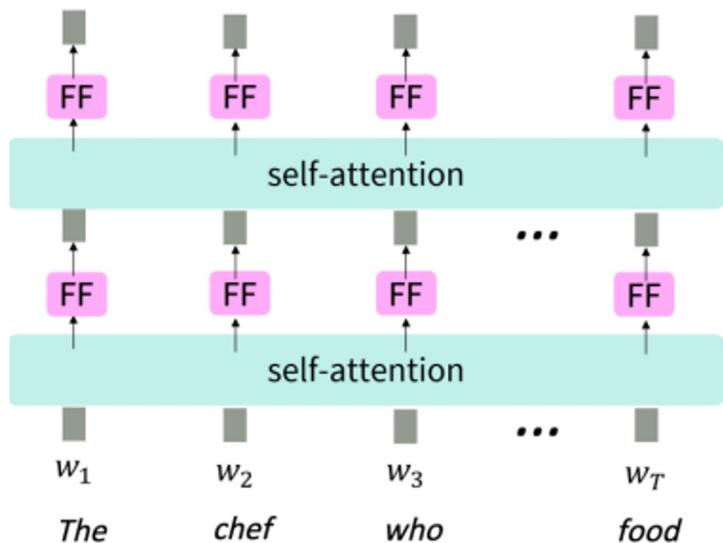
Self-attention in transformer (1)



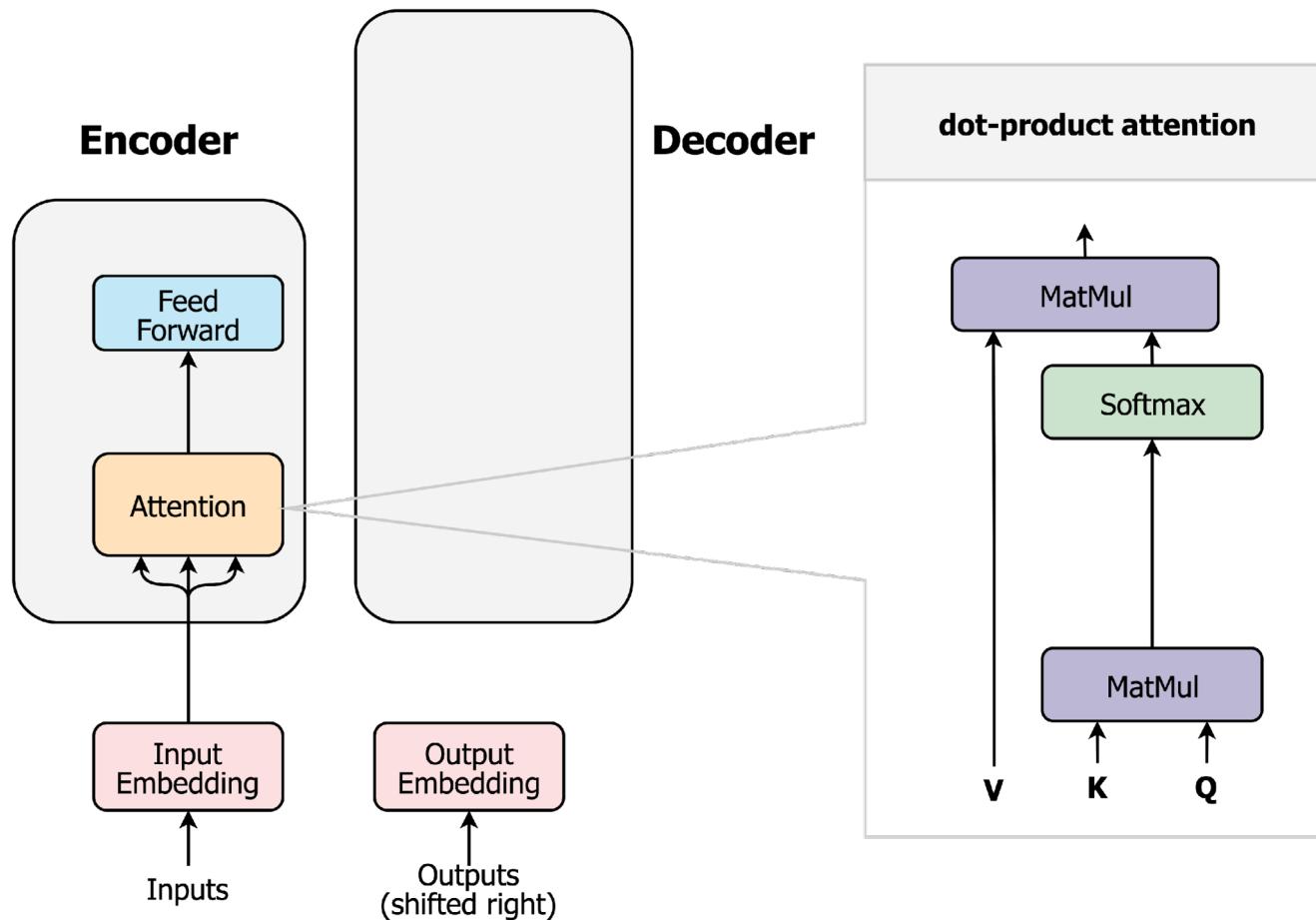
Self-attention in transformer (2)

- **Problem:** Since there are no element-wise non-linearities, self-attention is simply performing a re-averaging of the value vector
- **Solution:** Apply a feed forward layer to each output of attention (each row of the matrix Z), providing non-linear activation and additional expressive power

$$\mathbf{c}_i = \text{MLP}(\mathbf{z}_i) = \mathbf{W}_2 f(\mathbf{W}_1 \mathbf{z}_i + \mathbf{b}_1) + \mathbf{b}_2$$

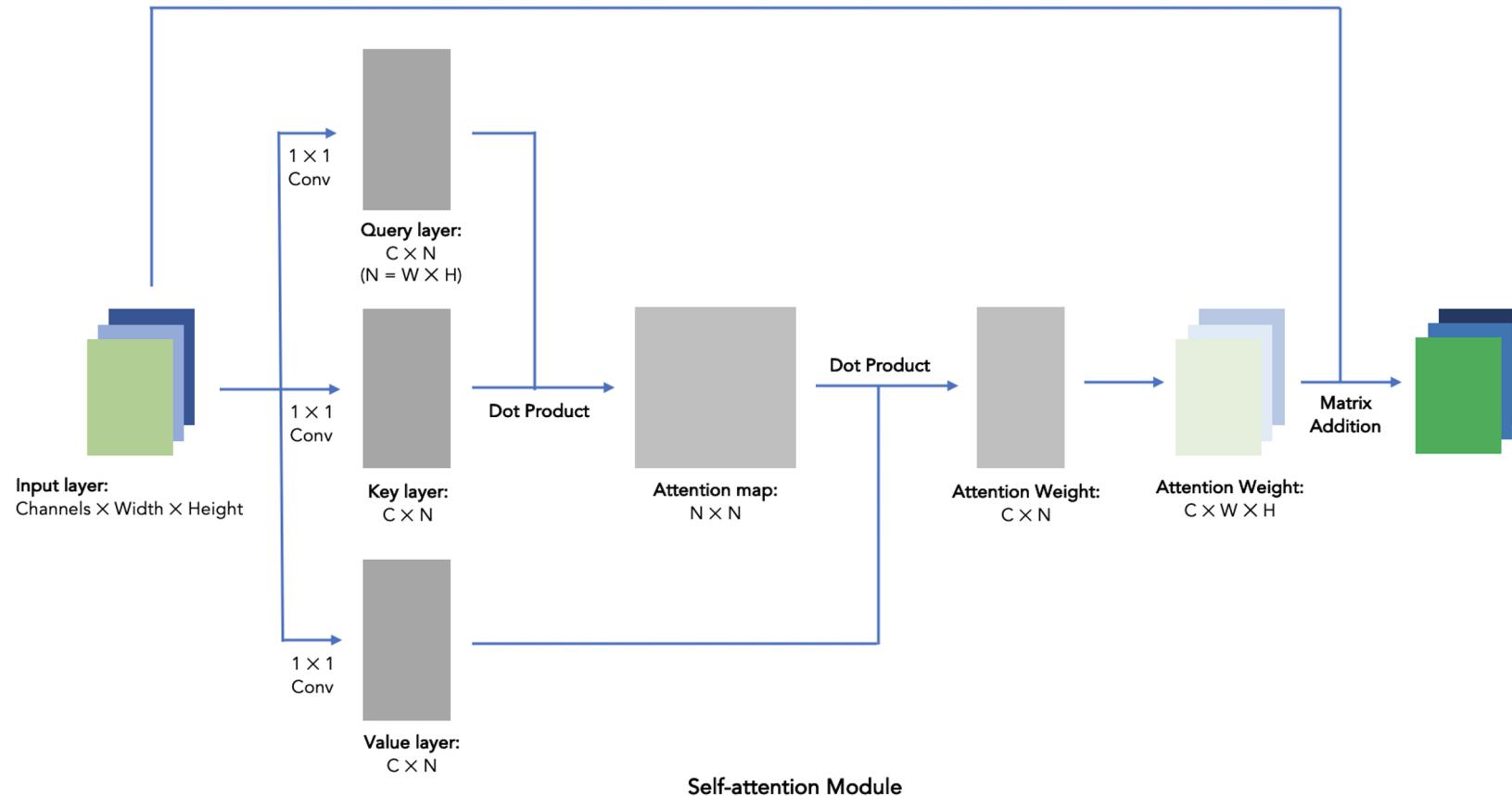


Self-attention in transformer (3)



Study case: Self-attention in image generation (1)

- Details can be generated using cues from all feature locations.
- Can check consistency between features in distant portions of the image



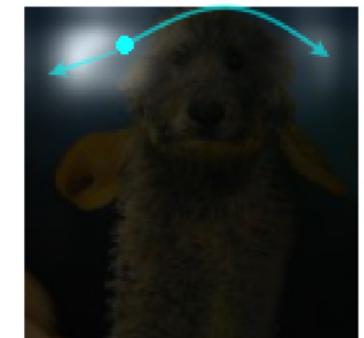
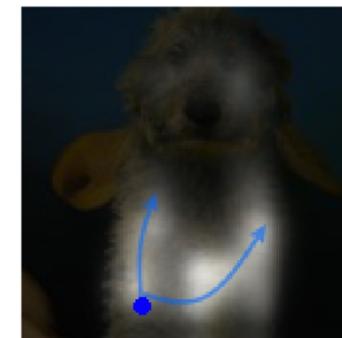
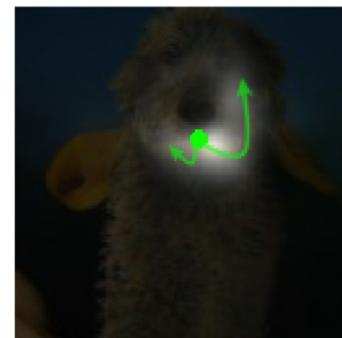
Zhang, Han, Ian Goodfellow, Dimitris Metaxas, and Augustus Odena. ["Self-attention generative adversarial networks."](#) ICML 2019. [\[video\]](#)

Study case: Self-attention in image generation (2)

Query locations



Attention maps for different query locations



Zhang, Han, Ian Goodfellow, Dimitris Metaxas, and Augustus Odena. "[Self-attention generative adversarial networks](#)." ICML 2019. [\[video\]](#)

Outline

- Introduction
- Self-attention
- **Multi-head attention**
- Deep layers
- Positional encoding
- The Transformer

Multi-head attention (1)

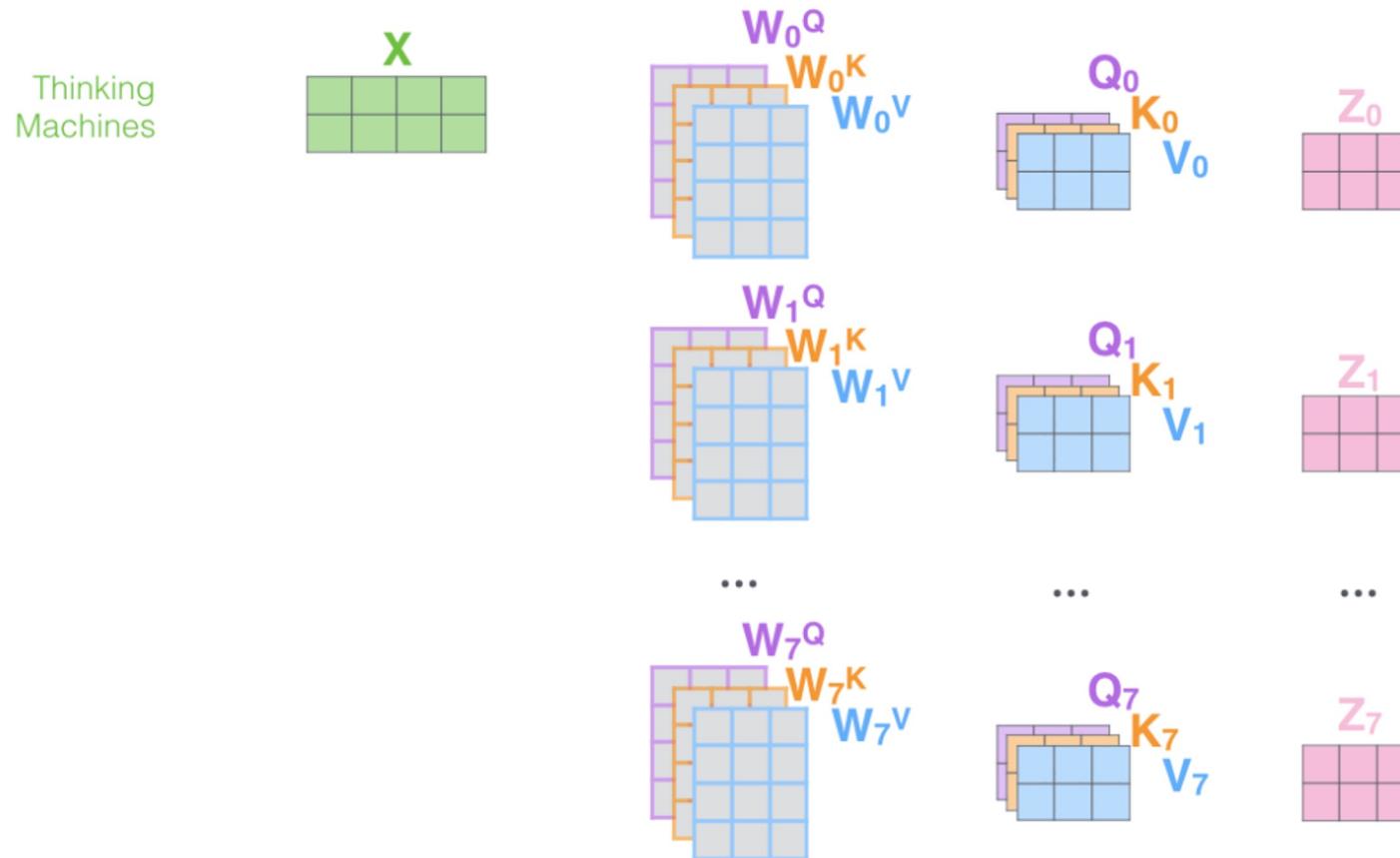
- H heads is better than 1!!
- **Main idea:** perform self-attention multiple times in parallel and combine the results



Wizards of the Coast, Artist: Todd Lockwood

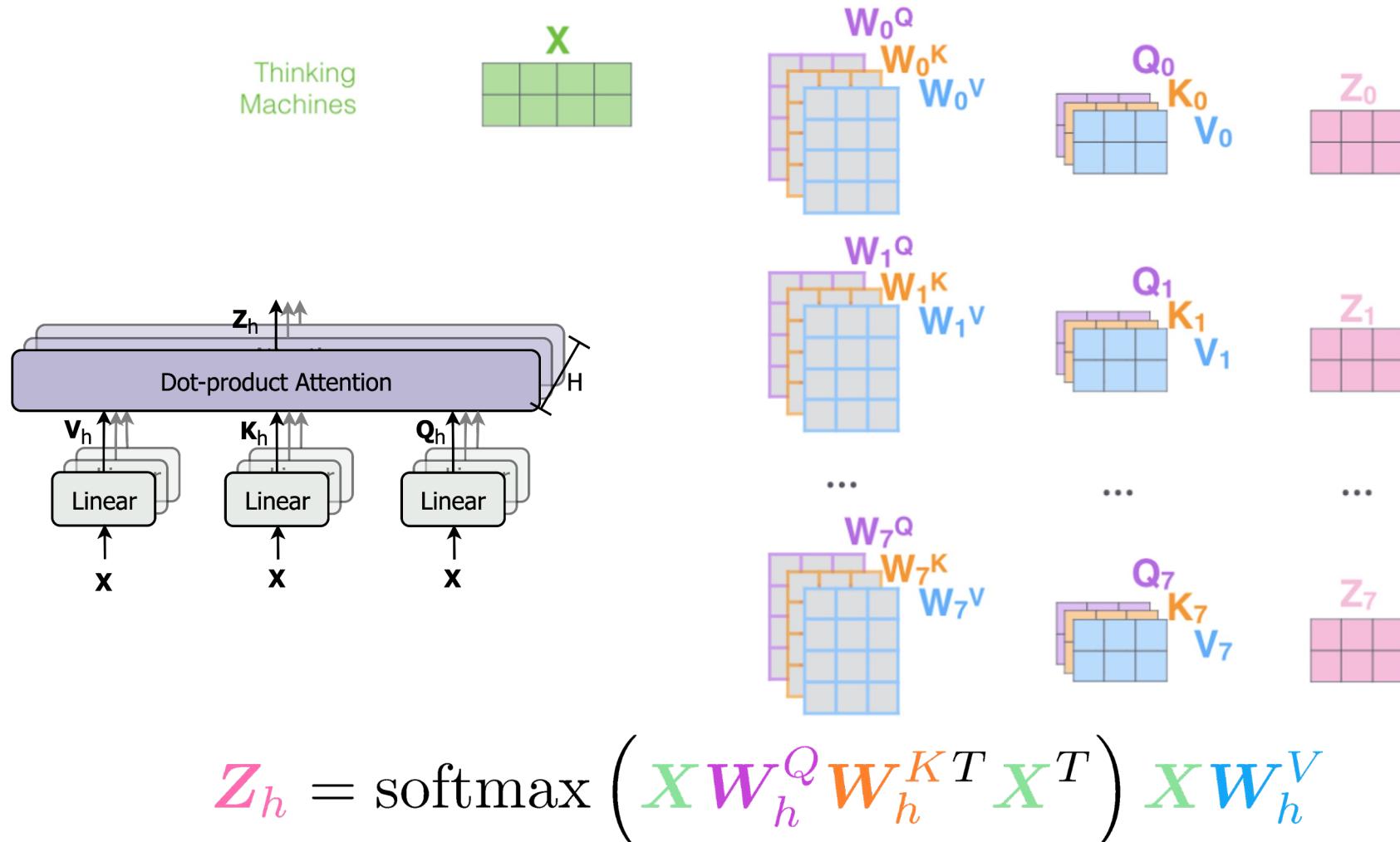
Multi-head attention (2)

Multiple sets of projection matrices to provide **different contextual representations** for the same input token



Multi-head attention (3)

Multiple sets of projection matrices to provide **different contextual representations** for the same input token



Jay Alammar, [The Illustrated Transformer](#), 2018.

Multi-head attention (5)

How to combine all H heads?

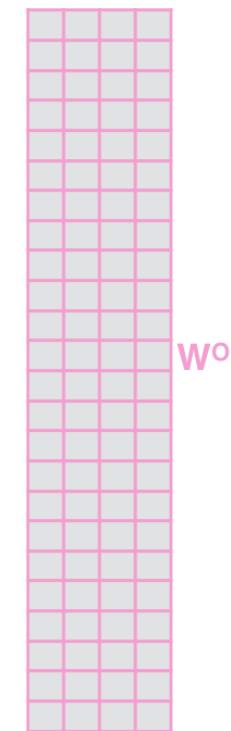
- All Z_h matrices are concatenated and combined in a new Z

1) Concatenate all the attention heads



2) Multiply with a weight matrix W^0 that was trained jointly with the model

X



3) The result would be the Z matrix that captures information from all the attention heads. We can send this forward to the FFNN

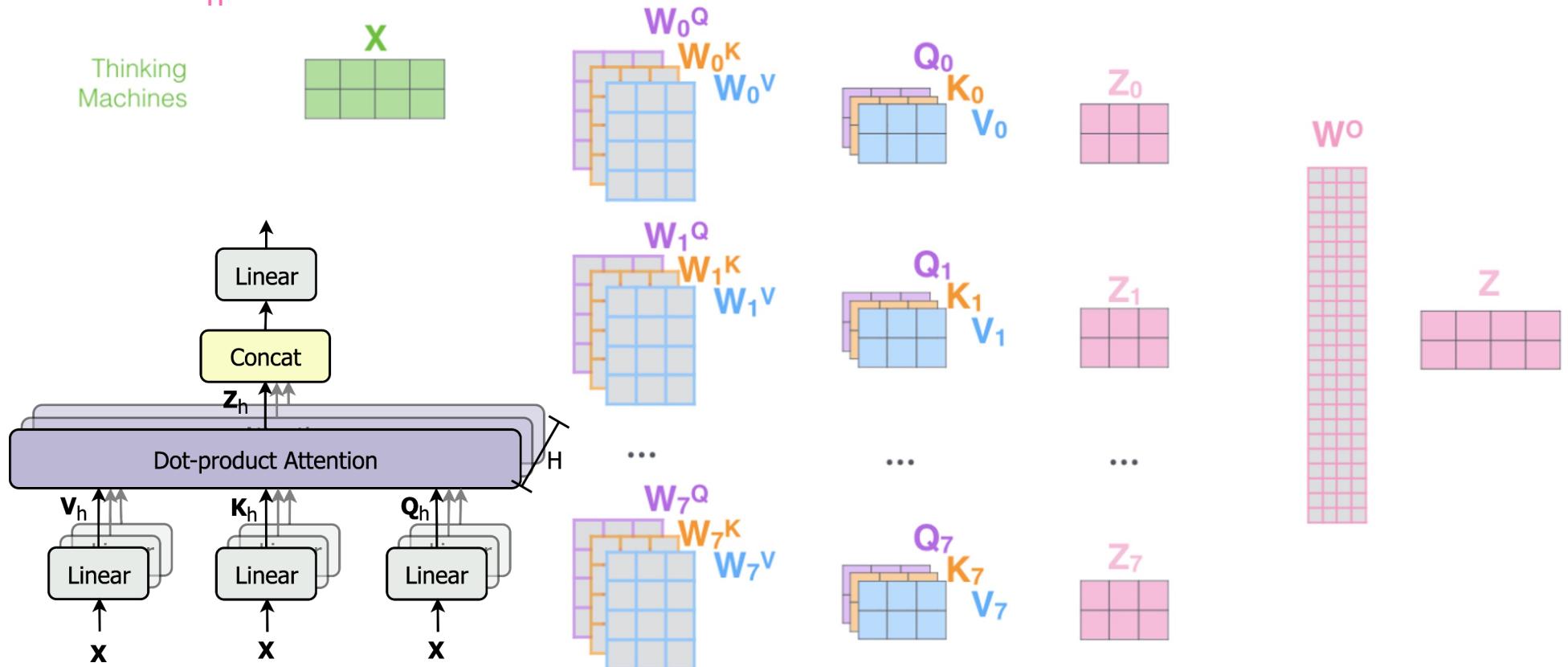
$$= \begin{matrix} Z \\ \hline \end{matrix}$$

$$Z = [Z_0, Z_1, \dots, Z_H] W^0$$

Multi-head attention (5)

How to combine all H heads?

- All Z_h matrices are concatenated and combined in a new Z

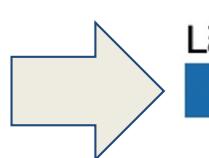


$$Z = [Z_0, Z_1, \dots, Z_H] W^0$$

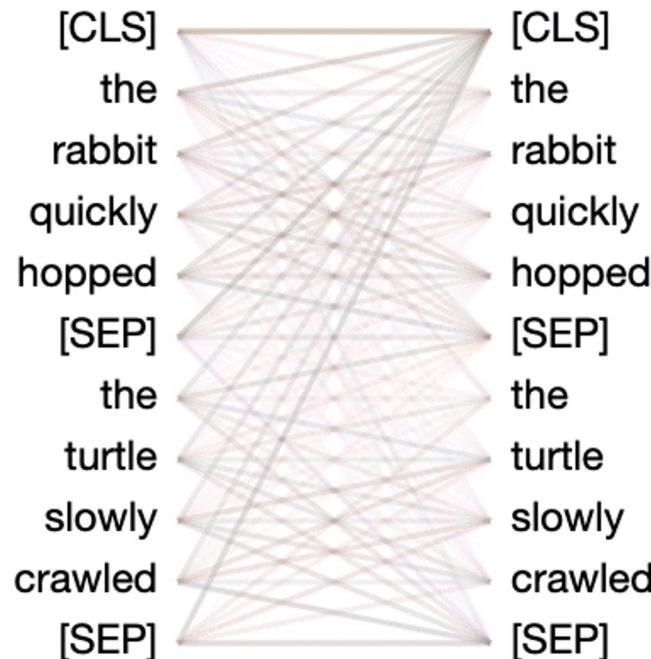
Multi-head attention (6)

Visualization

Each colour corresponds to a head.



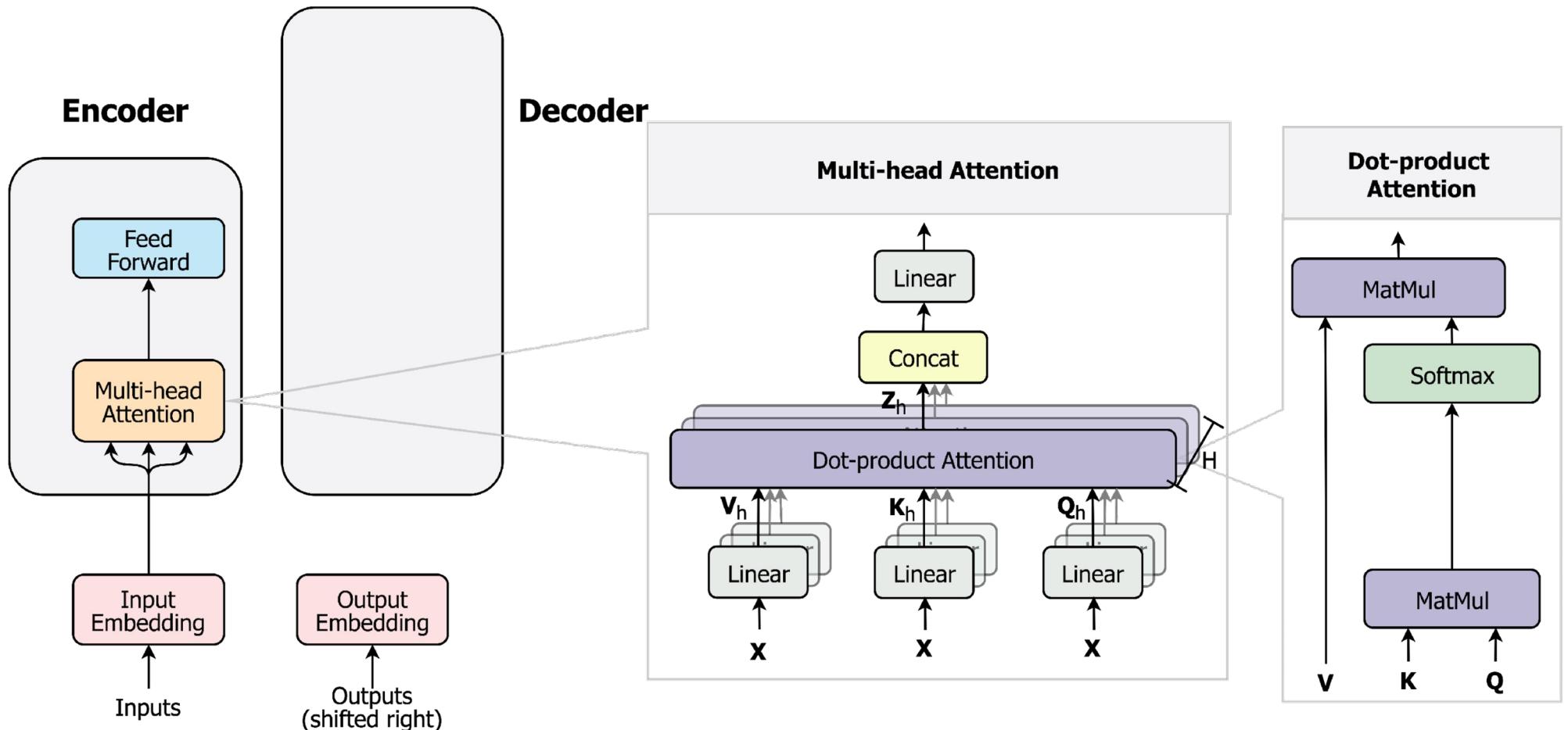
Layer: 0 Attention: All



Blue: First head only.

Multi-color: Multiple heads.

Multi-head attention in transformer



Outline

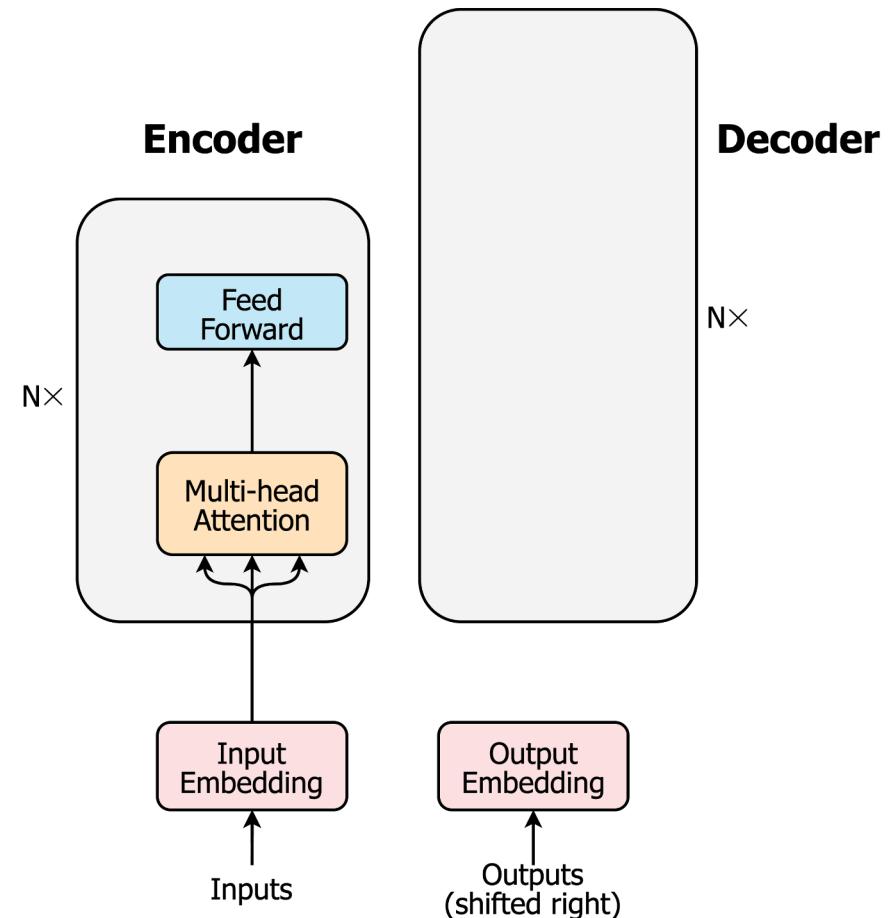
- Introduction
- Self-attention
- Multi-head attention
- **Deep layers**
- Positional encoding
- The Transformer

Let's go deep (1)

Add **N layers** ($N=6$ in the original paper) of multi-head attentions and feed-forward networks.

Repeat this structure N times:

- First layer's input are word embeddings
- Further layers' input is the context vector output of the previous layer

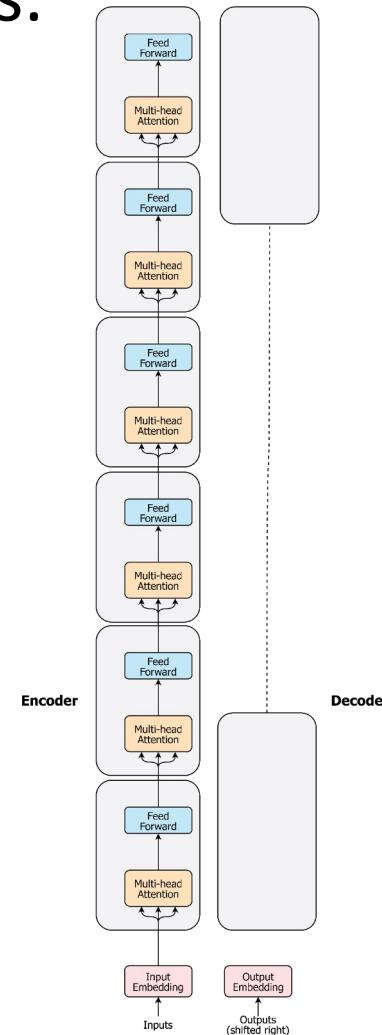


Let's go deep (2)

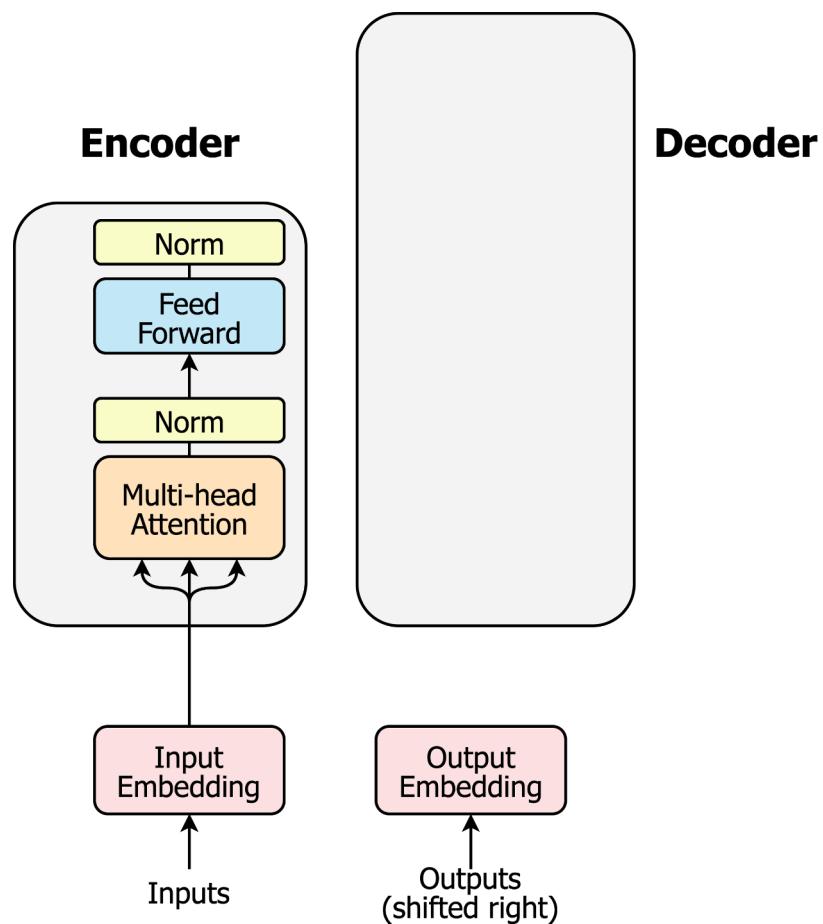
Add N layers (N=6 in the original paper) of multi-head attentions and feed-forward networks.

Repeat this structure N times:

- First layer's input are word embeddings
- Further layers' input is the context vector output of the previous layer



Layer normalization



- **Problem:** Difficult to train the parameters of a given layer because its input from the layer beneath keeps shifting
- **Solution:** Reduce uninformative variation by normalizing to zero mean and standard deviation of one within each layer [[Ba et al., 2016](#)]

$$\text{Mean: } \mu^l = \frac{1}{H} \sum_{i=1}^H a_i^l \quad \text{Standard Deviation: } \sigma^l = \sqrt{\frac{1}{H} \sum_{i=1}^H (a_i^l - \mu^l)^2}$$

$$x^{\ell'} = \frac{x^\ell - \mu^\ell}{\sigma^\ell + \epsilon}$$

Scaled dot-product attention (1)

- After Layer **Normalization**, the mean and variance of vector elements is 0 and 1, respectively
- However, the dot product still tends to take on extreme values, as its variance scales with dimensionality d_k

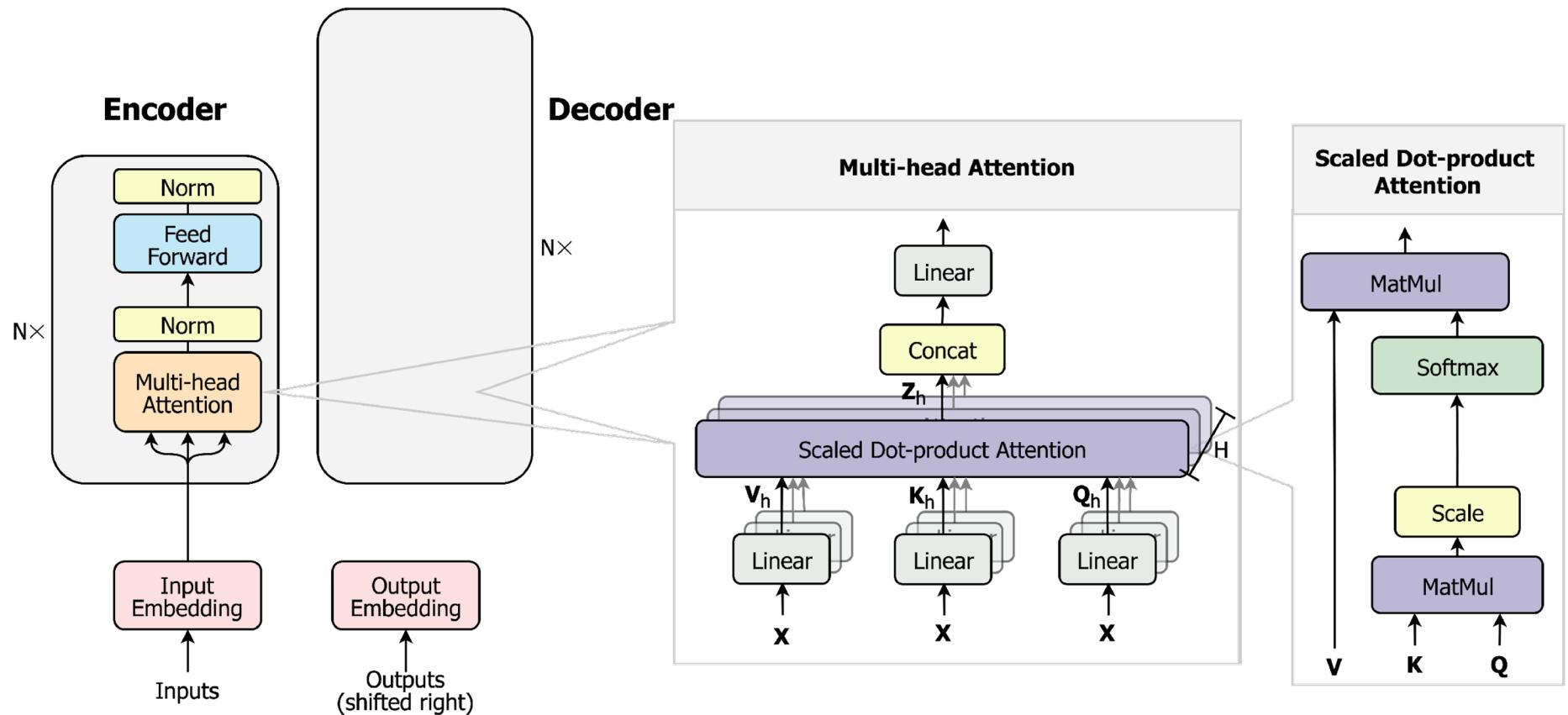
Quick Statistics Review:

- Mean of sum = sum of means = $d_k * 0 = 0$
- Variance of sum = sum of variances = $d_k * 1 = d_k$
- To set the variance to 1, simply divide by $\sqrt{d_k}$!

5) Updated self-attention equation (context)

$$\text{softmax} \left(\frac{\begin{array}{c} \text{Q} \quad \text{K}^T \\ \begin{matrix} \text{pink} & \times & \text{orange} \end{matrix} \\ \hline \sqrt{d_k} \end{array}}{\begin{array}{c} \text{V} \\ \text{blue} \end{array}} \right) = \begin{array}{c} \text{Z} \\ \text{pink} \end{array}$$

Scaled dot-product attention (2)



$$\mathbf{Z} = \text{MHSAs}(\mathbf{X}) = [\mathbf{Z}_0, \dots, \mathbf{Z}_H] \mathbf{W}^0$$

$$\mathbf{Z}_h = \text{softmax} \left(\frac{\mathbf{Q}_h \mathbf{K}_h^T}{\sqrt{d_k}} \right) \mathbf{V}_h = \text{softmax} \left(\frac{\mathbf{X} \mathbf{W}_h^Q \mathbf{W}_h^{K^T} \mathbf{X}^T}{\sqrt{d_k}} \right) \mathbf{X} \mathbf{W}_h^V$$

Residual connections (1)

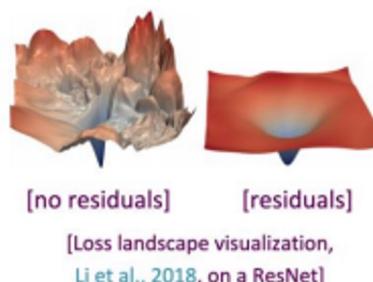
Residual connections are a simple but powerful technique [[He et al., 2016](#)]

- Deep networks are surprisingly bad at learning the identity function
- Directly passing "raw" embeddings to the next layer can actually be very helpful

$$x_l = F(x_{l-1}) + x_{l-1}$$

- This prevents the network from "forgetting" or distorting important information as it is processed by many layers.

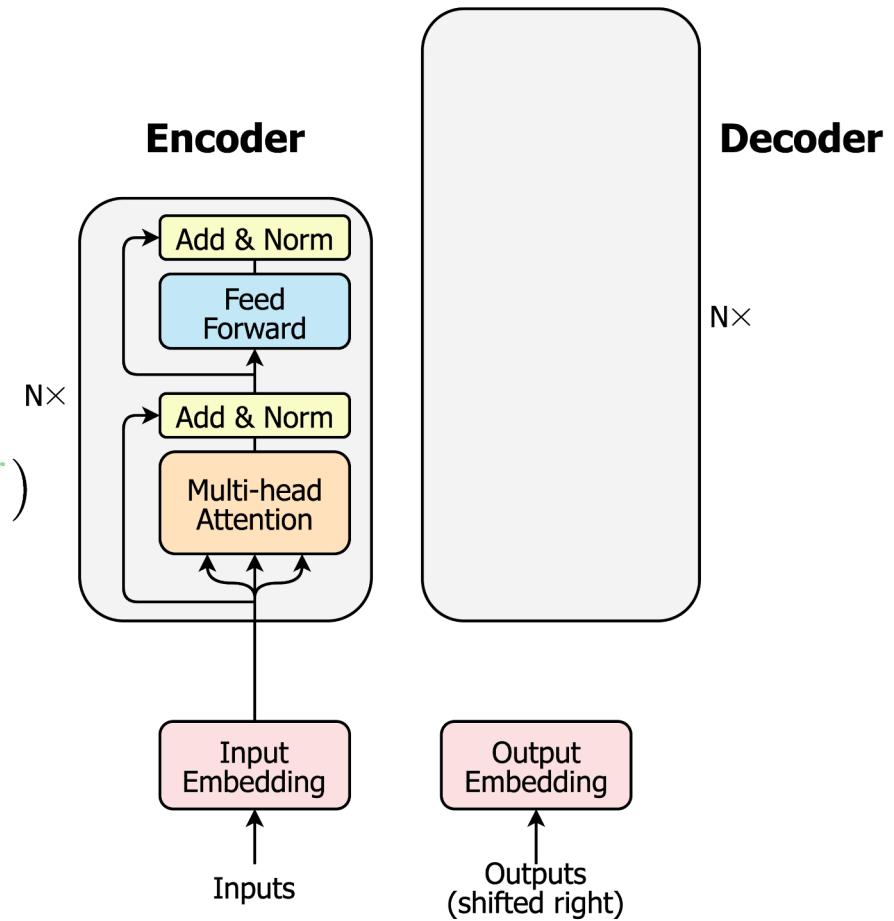
Residual connections are also thought to smooth the loss landscape and make training easier!



Residual connections (2)

$$C = \text{MLP}(\mathbf{Z}) + \mathbf{Z}$$

$$\mathbf{Z} = \text{LayerNorm}(\text{MHSAs}(\mathbf{X}) + \mathbf{X})$$

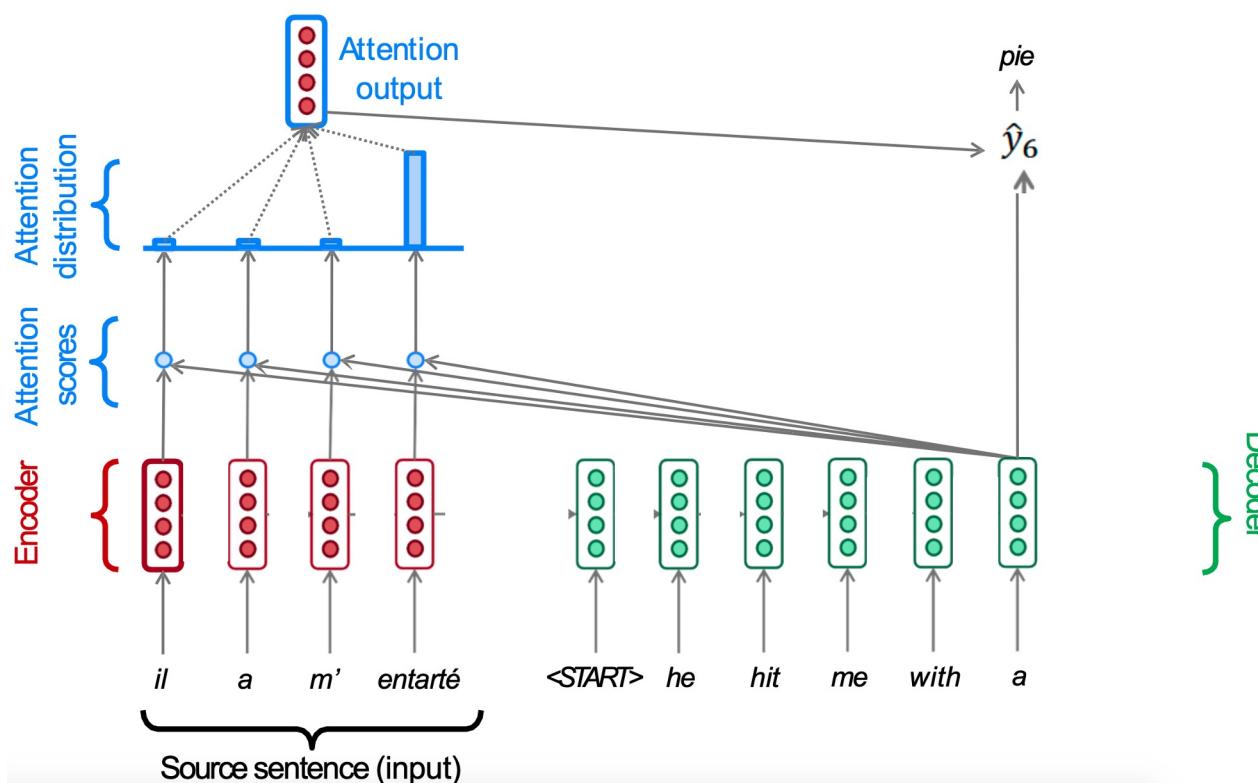


Outline

- Introduction
- Self-attention
- Multi-head attention
- Deep layers
- **Positional encoding**
- The Transformer

Positional encoding (1)

- Given that the attention mechanism allows accessing all input (and output) tokens, we no longer need a memory through recurrent layers



Positional encoding (2)

- Consider the following two sequences:

elephants hate mice
 x_1 x_2 x_3

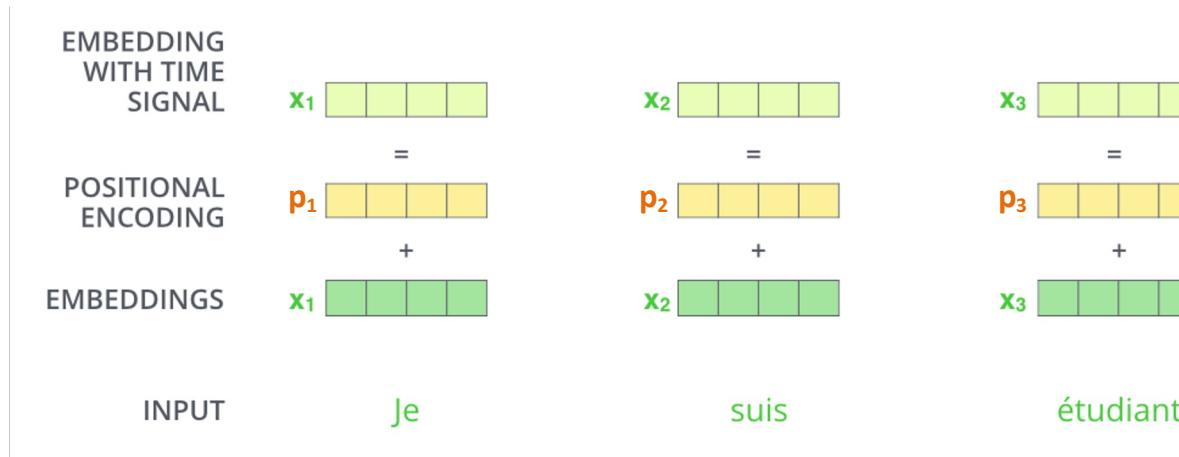
Mice hate elephants
 x_1 x_2 x_3

Are self-attention scores different in both cases?

Positional encoding (3)

- Since self-attention does not build in order information, we need to encode the order in the keys, queries and values
- Consider representing each **sequence index** as a **vector** (positional encoding): $p_i \in \mathbb{R}^{d_x}$ for $i \in \{1, 2, \dots, T\}$
- **Add** the positional encoding to the word embedding (concatenating also a possibility)

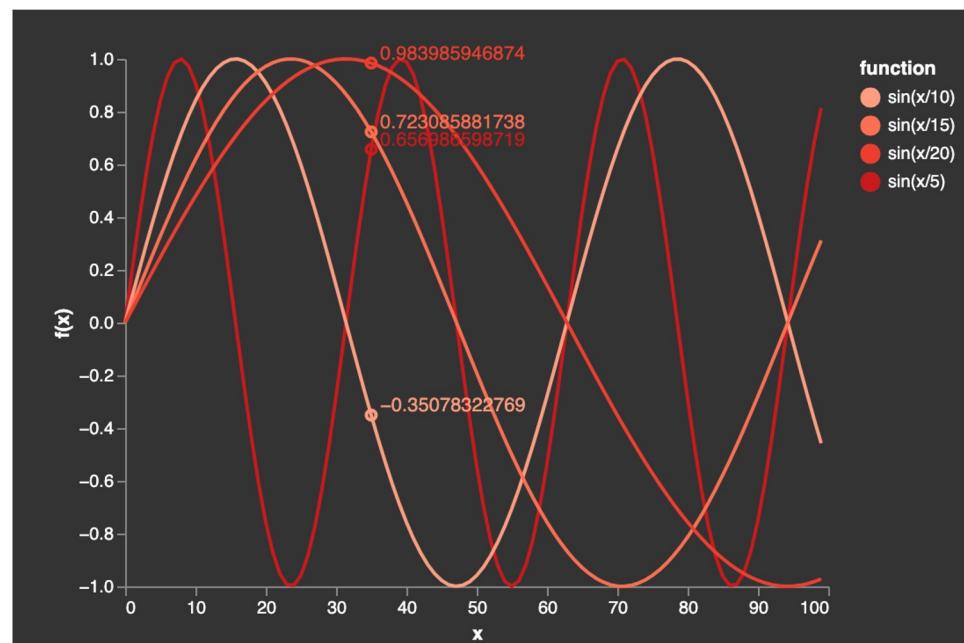
$$x_i = \tilde{x}_i + p_i$$



Positional encoding (4)

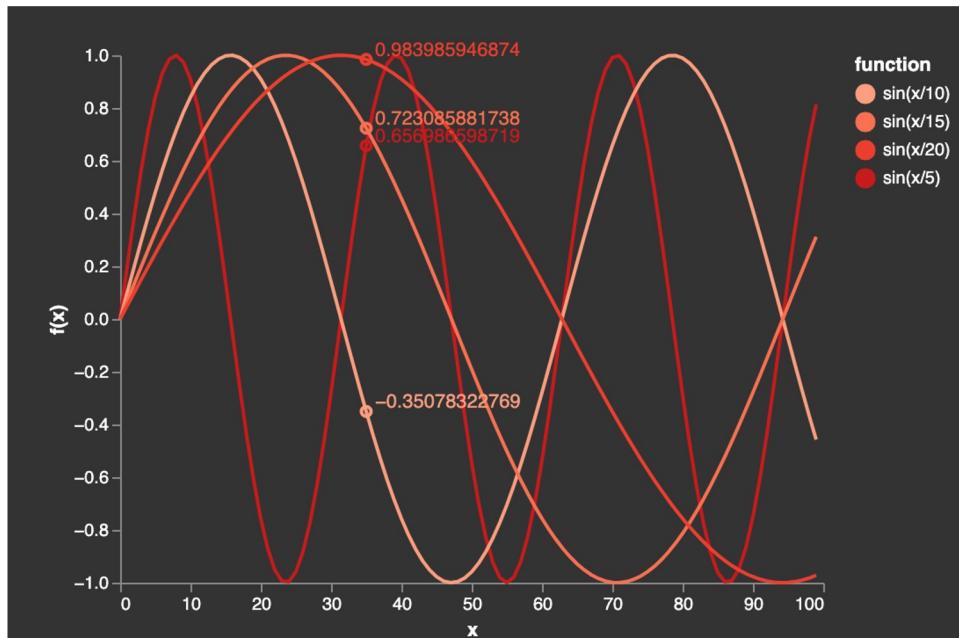
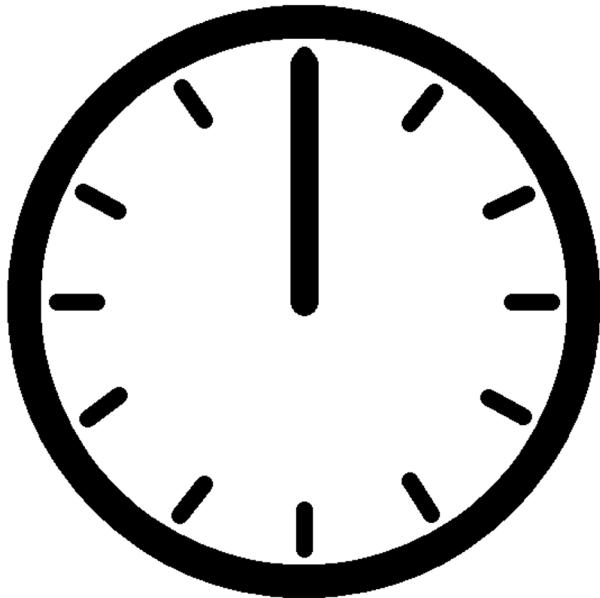
- **Sinusoidal position representations:** concatenate sinusoidal functions of varying periods:

$$p_i = \begin{pmatrix} \sin\left(\frac{i}{10000^2 \frac{1}{d_x}}\right) \\ \cos\left(\frac{i}{10000^2 \frac{1}{d_x}}\right) \\ \vdots \\ \sin\left(\frac{i}{10000^2 \frac{d_x/2}{d_x}}\right) \\ \cos\left(\frac{i}{10000^2 \frac{d_x/2}{d_x}}\right) \end{pmatrix}$$



Positional encoding (5)

- **Sinusoidal position representations:** concatenate sinusoidal functions of varying periods:

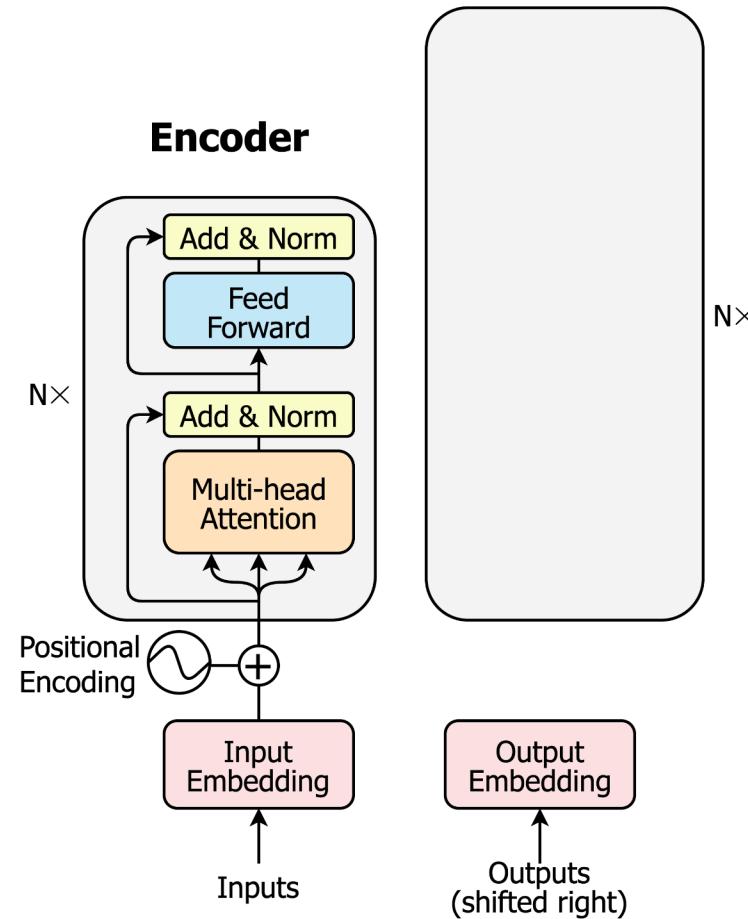


- **Advantages:**
 - Periodicity indicates that maybe “absolute position” isn’t as important
 - Maybe can extrapolate to longer sequences as periods restart
- **Disadvantages:**
 - Not learnable (also, the extrapolation doesn’t really work)

Positional encoding (6)

- **Learned absolute position representations:** Let all \mathbf{p}_i be learnable parameters
 - Learn a matrix $\mathbf{P} \in \mathbb{R}^{T \times d_x}$ and let each \mathbf{p}_i be a row of that matrix
- **Advantages:**
 - Flexibility: each position gets to be learned to fit the data
- **Disadvantages:**
 - Definitely can't extrapolate to indices outside $1, \dots, T$
 - Most systems use this

Positional encoding in transformer

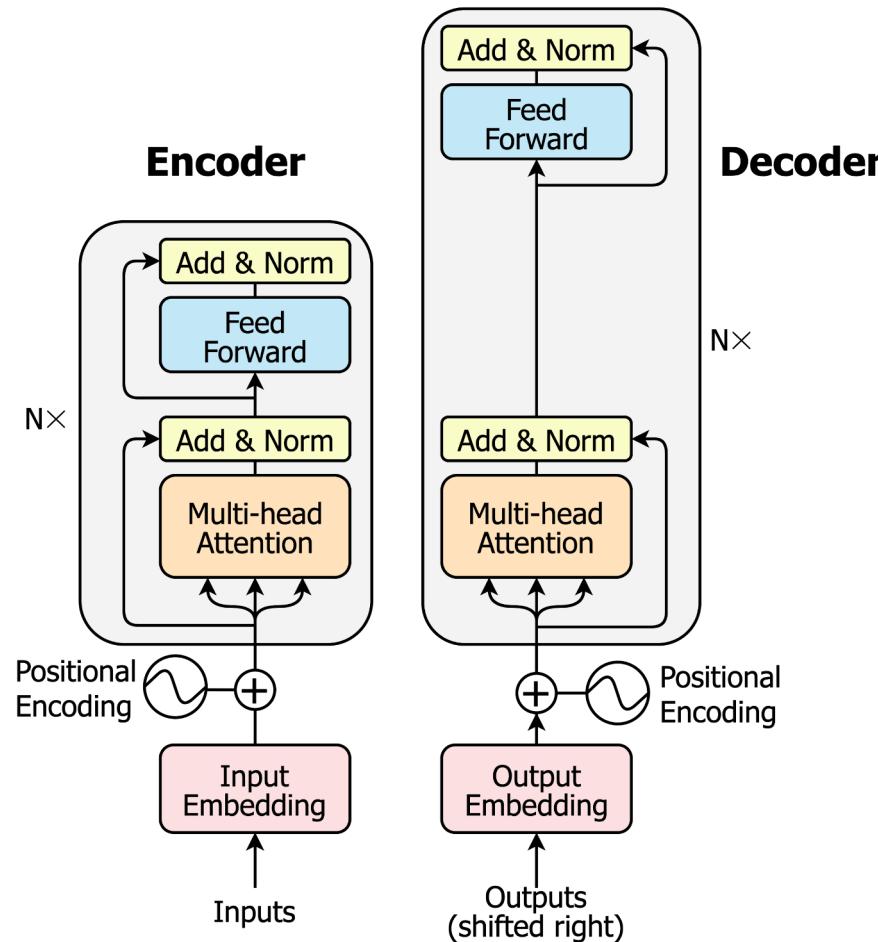


Outline

- Introduction
- Self-attention
- Multi-head attention
- Deep layers
- Positional encoding
- **The Transformer**

Transformer: decoder

- We can also exploit self-attention in the decoder



Self-attention in the decoder

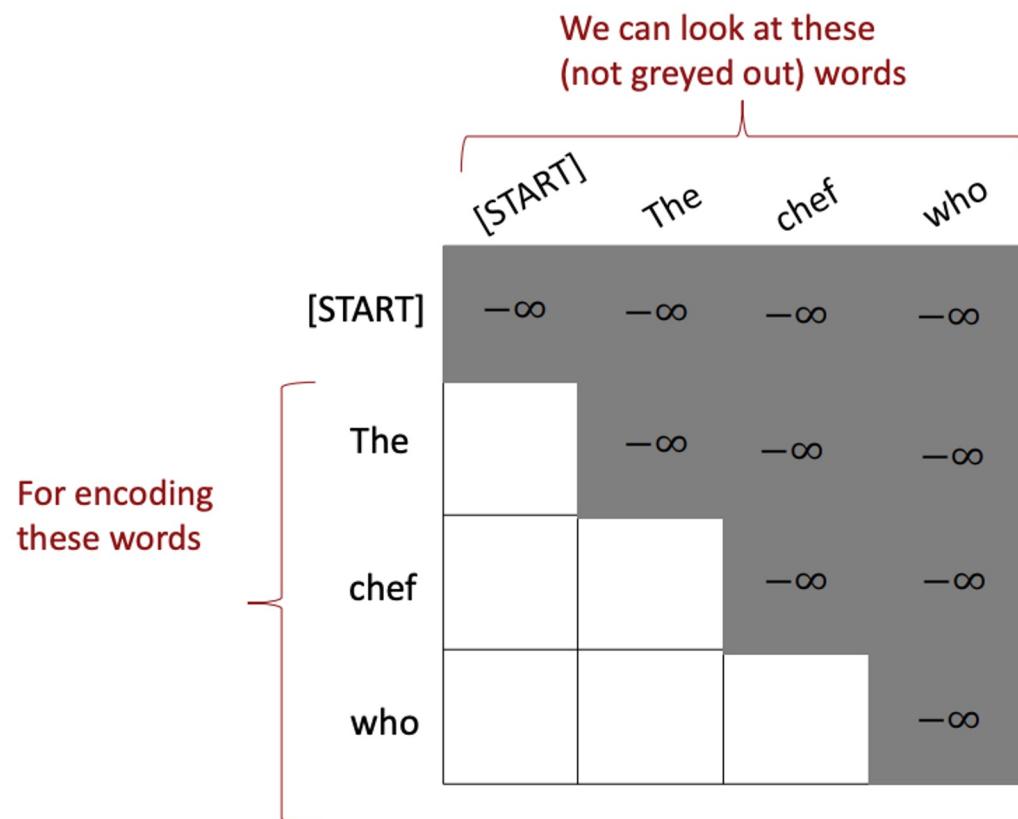
However!

- At training time, the decoder is feed the complete output sequence
- How do we keep the decoder from cheating (not being able to look-ahead)?
- Solution: Masked attention
 - Hide information about future tokens from the model

The little black dog [masked] [masked]
[masked]

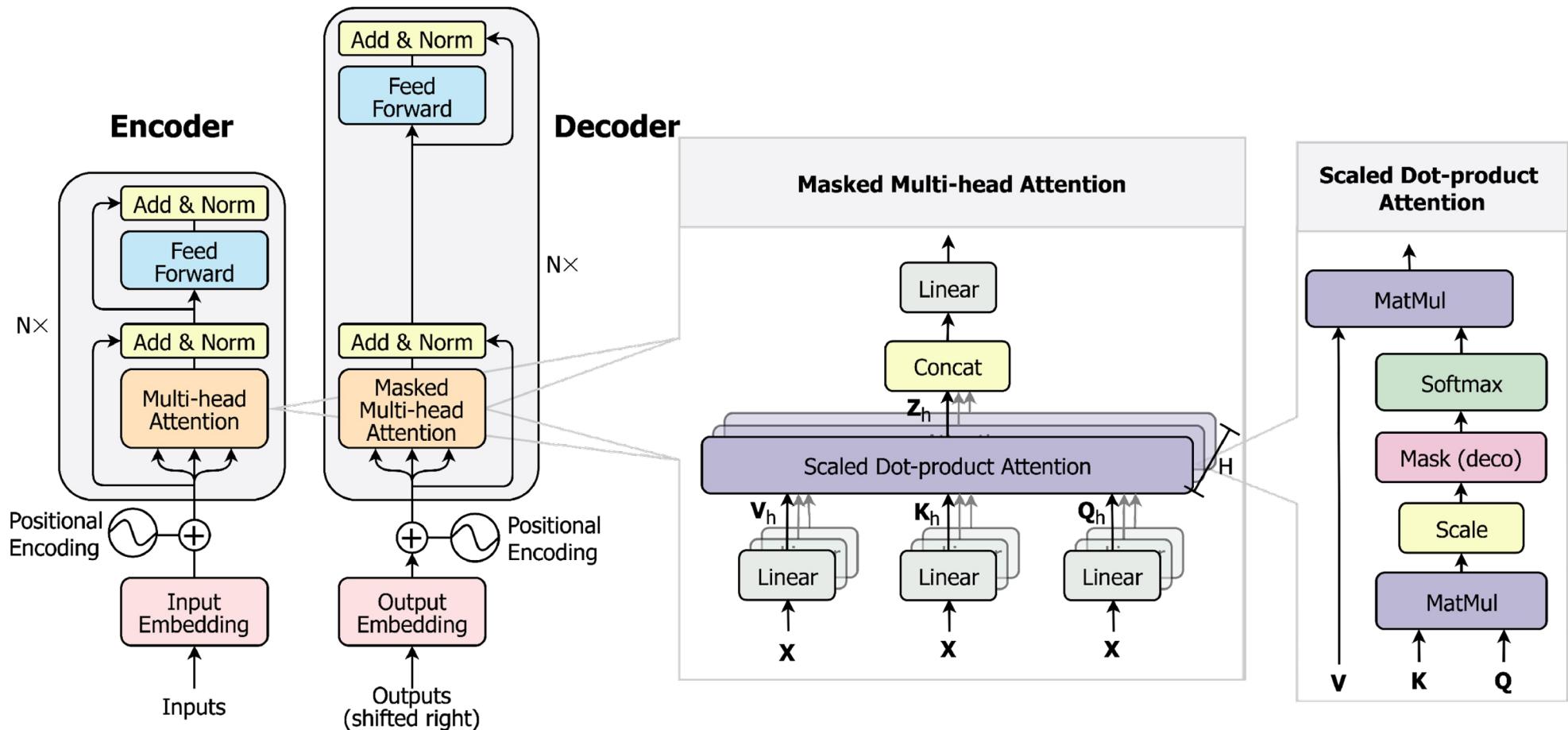
Masked self-attention (1)

- To enable parallelization, mask attention to future tokens by setting attention scores to $-\infty$



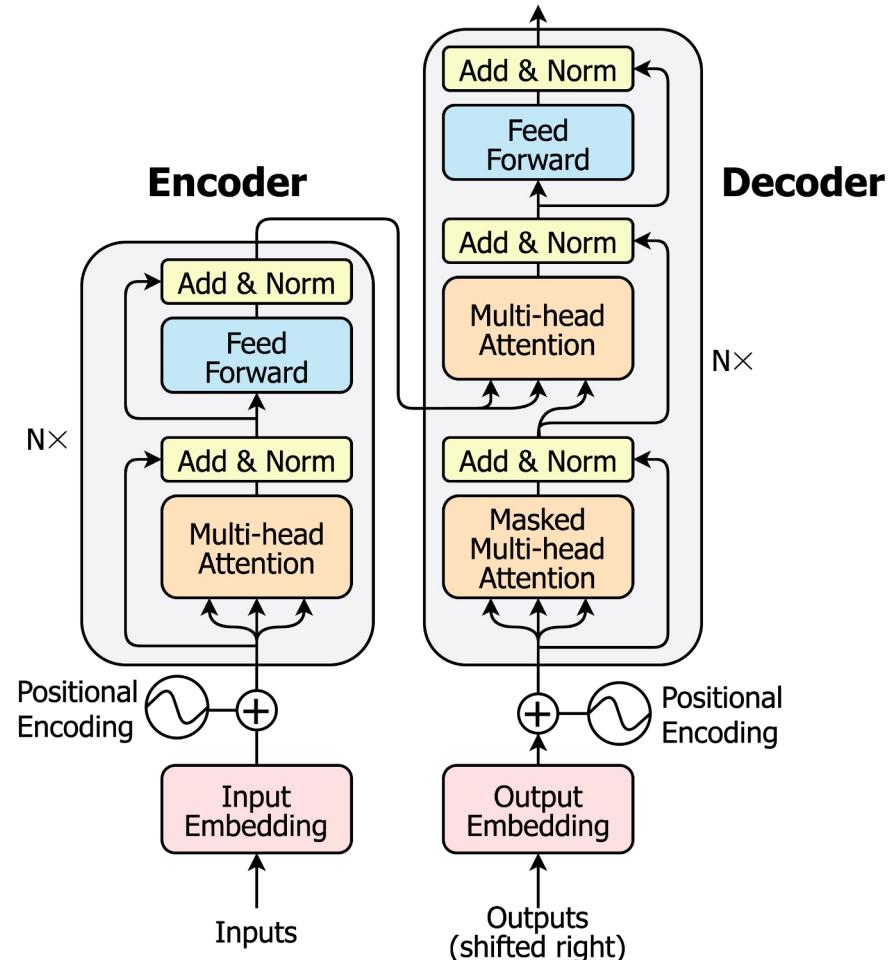
$$e_{ij} = \begin{cases} q_i \cdot k_j & j < i \\ -\infty & j \geq i \end{cases}$$

Masked self-attention (2)



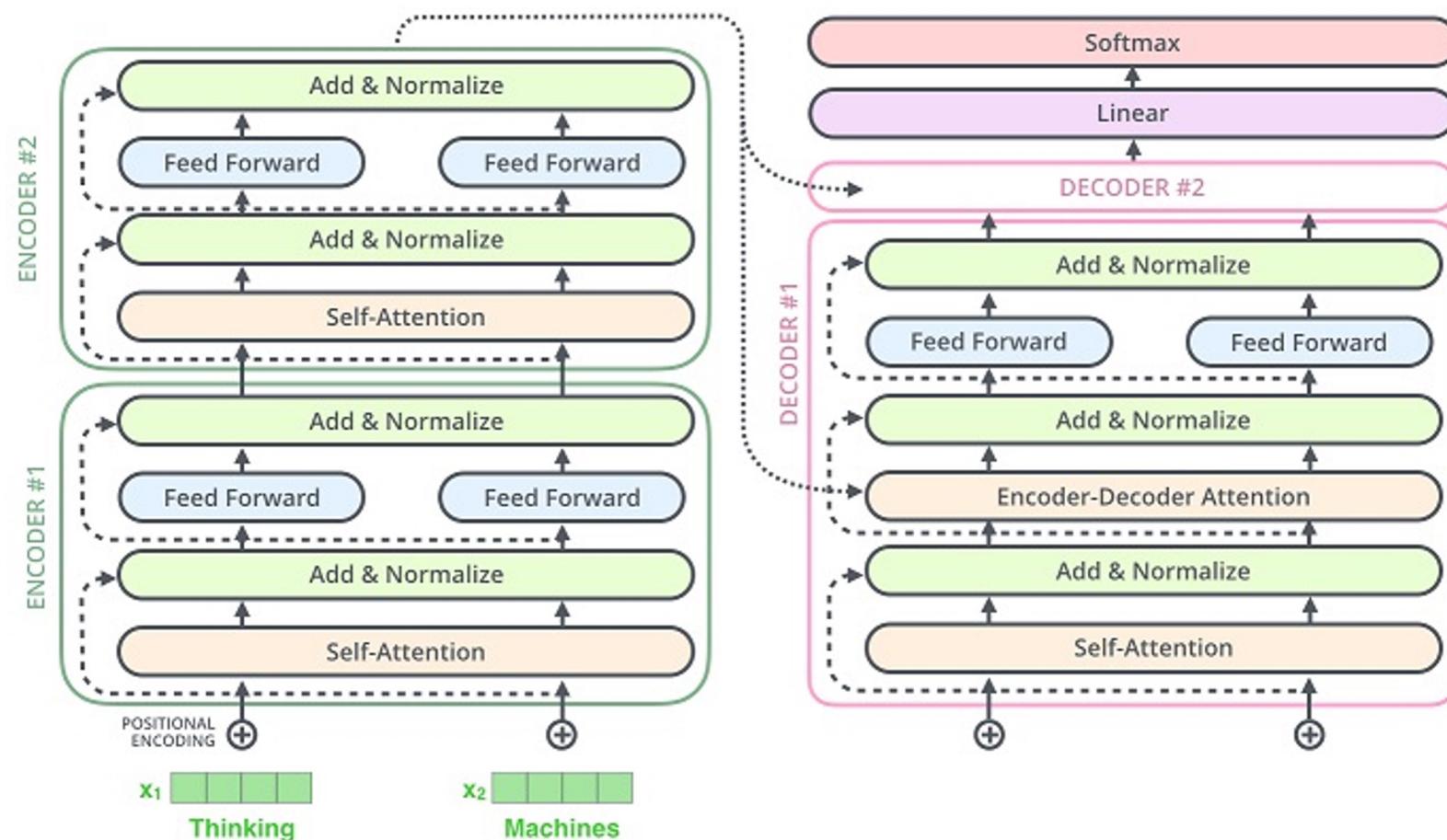
Encoder-Decoder attention (1)

- Use cross-attention in the decoder, where:
 - Keys and values are obtained from the encoder (like a memory)
 - Queries are obtained from the decoder

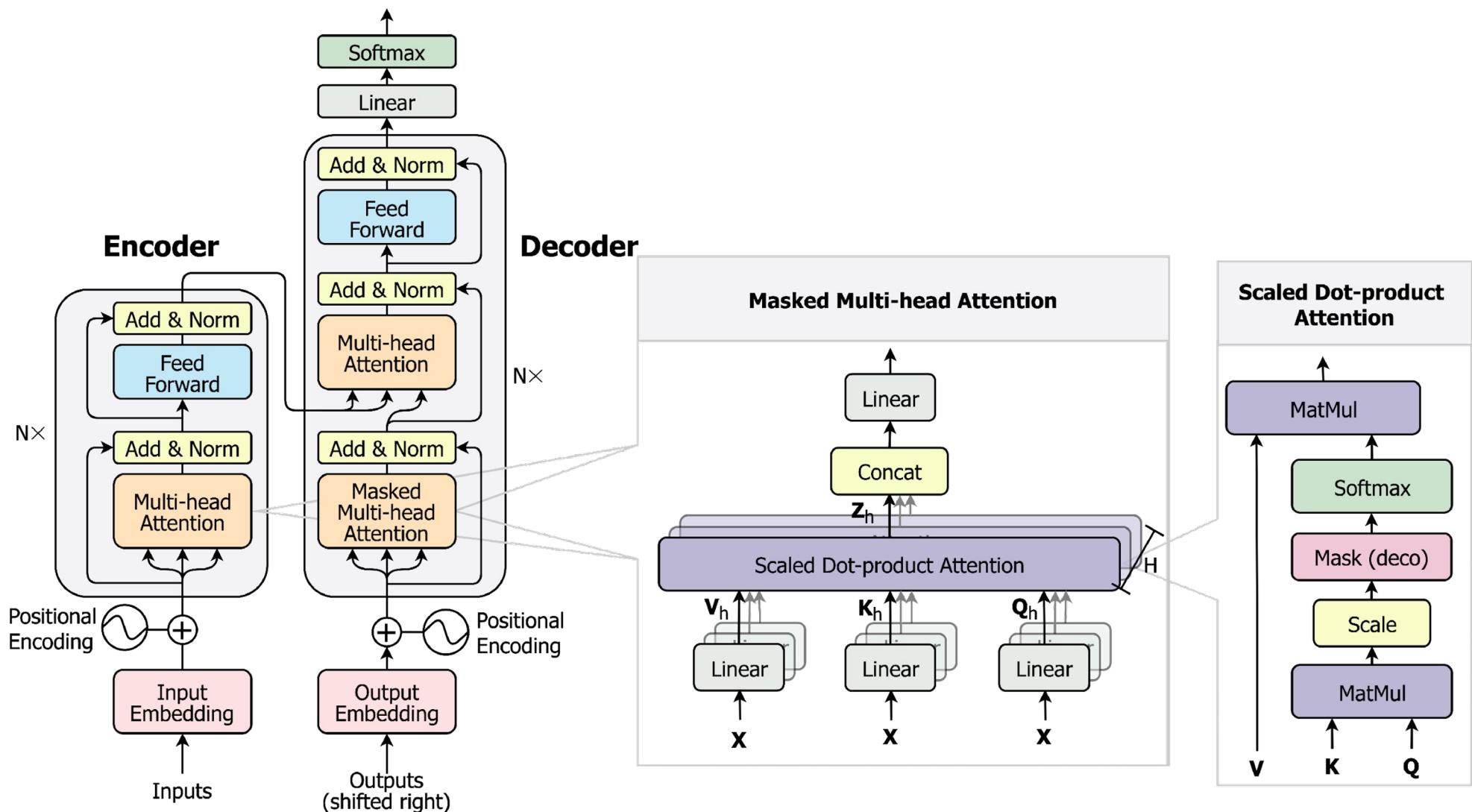


Encoder-Decoder attention (2)

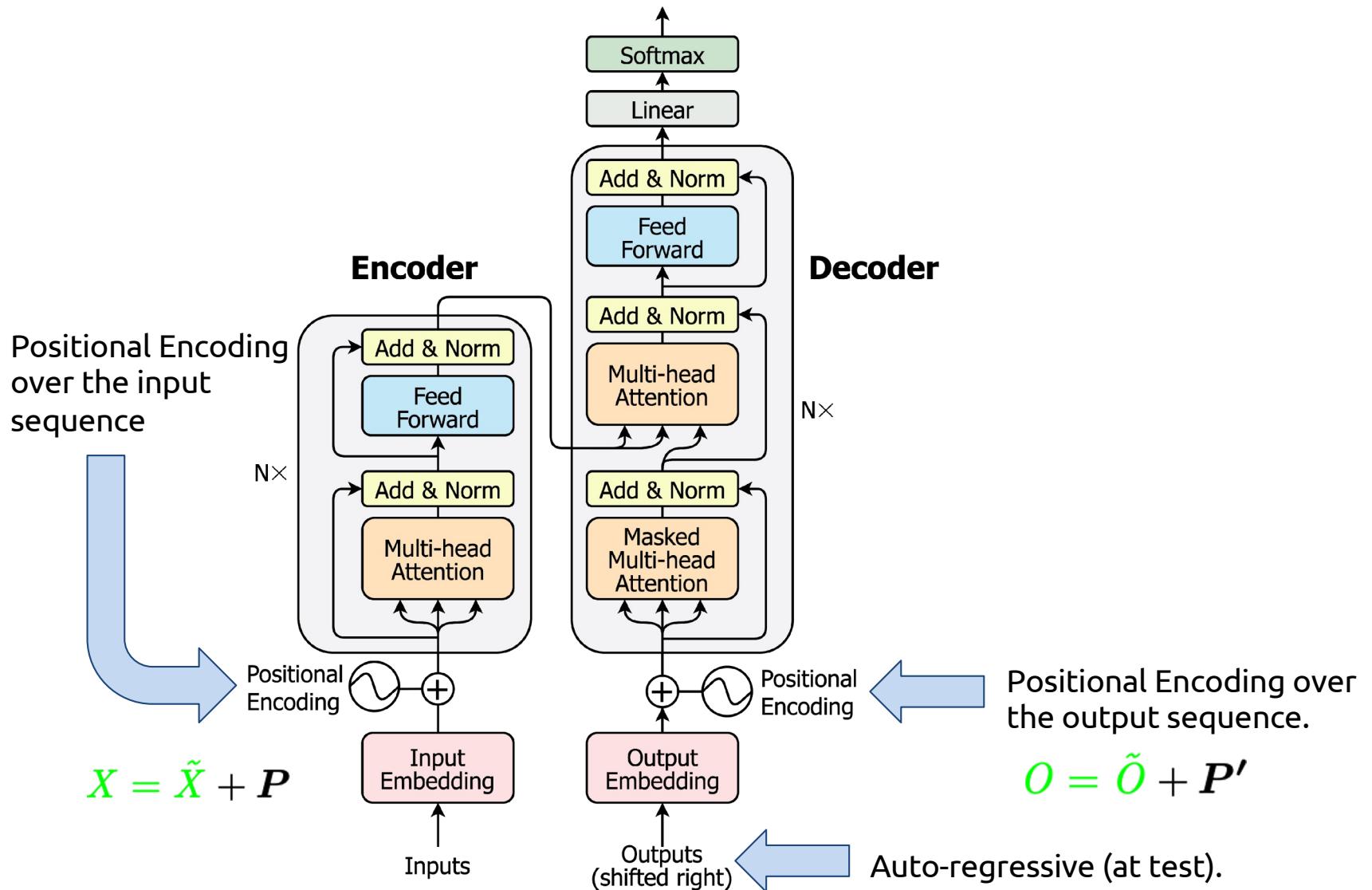
- The last output of the encoder is fed into all decoder layers



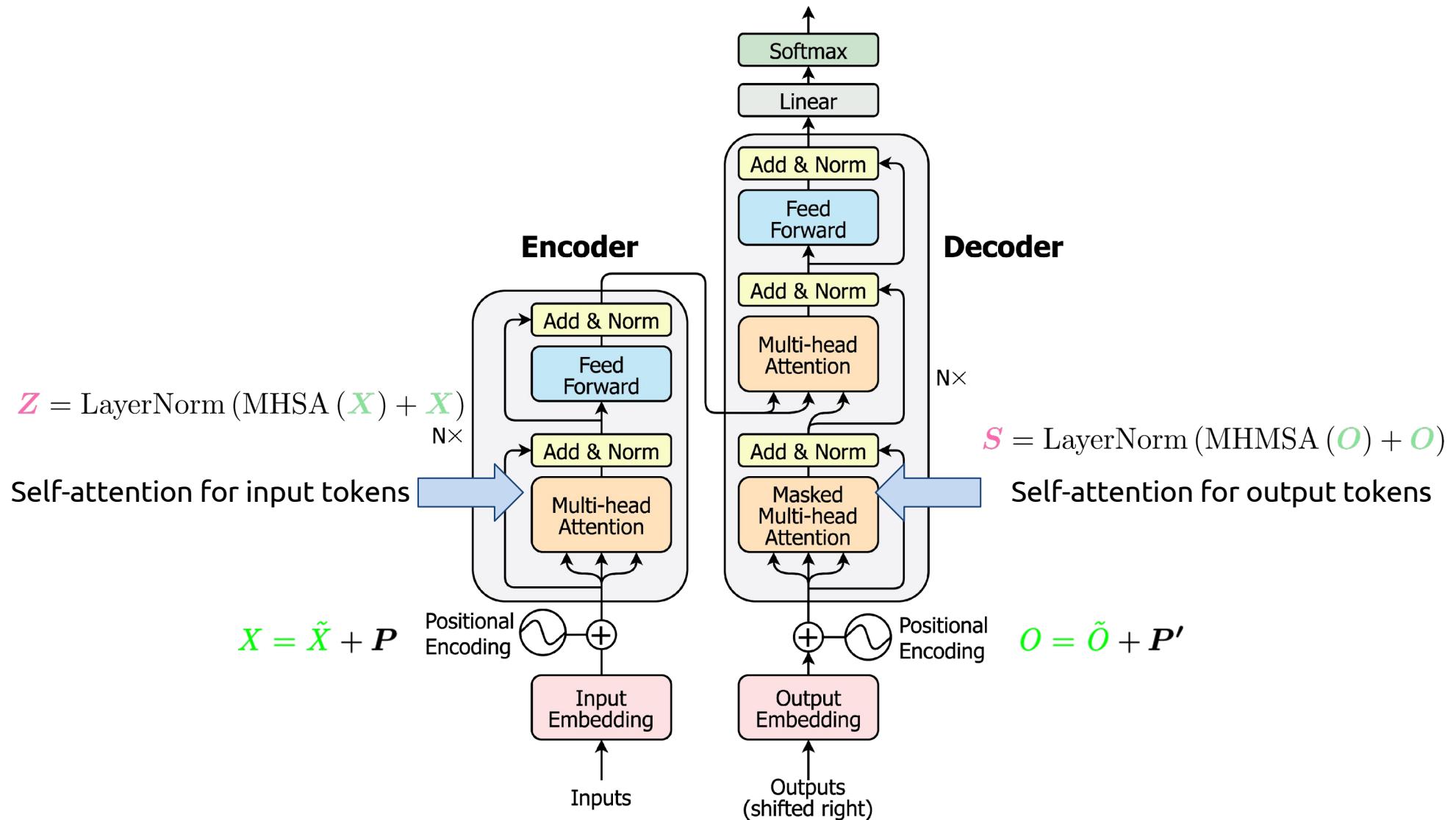
Complete transformer block (1)



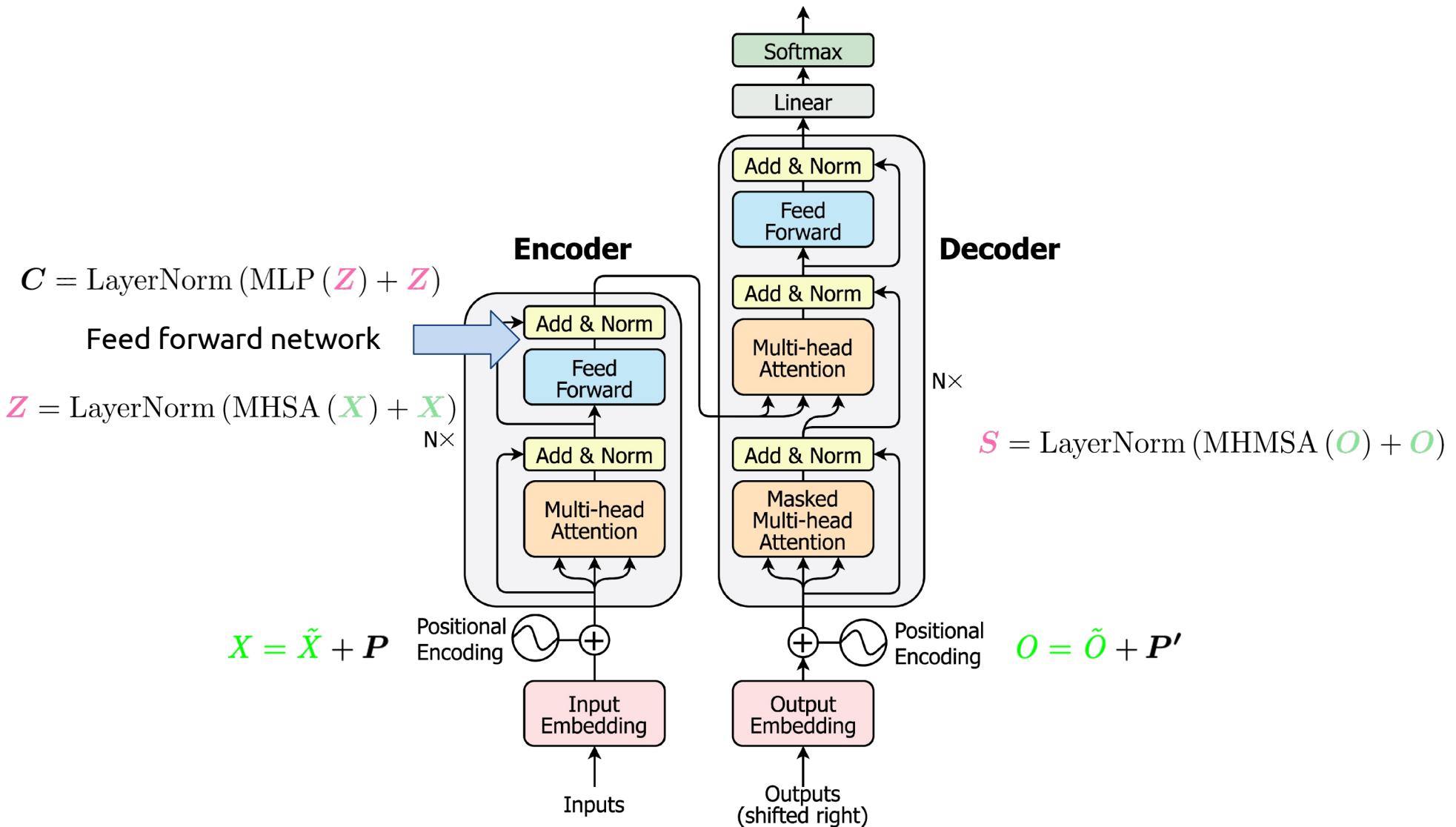
Complete transformer block (2)



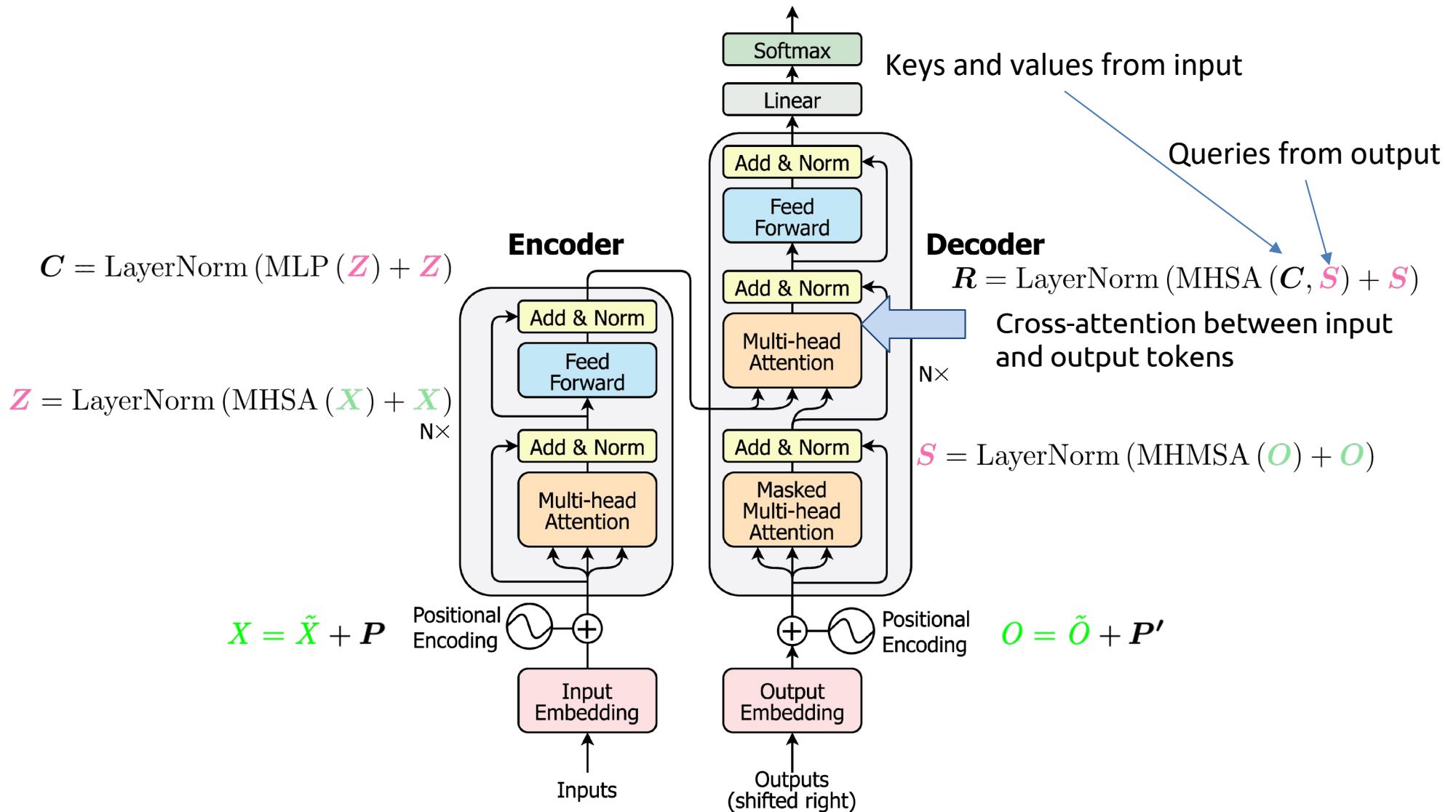
Complete transformer block (3)



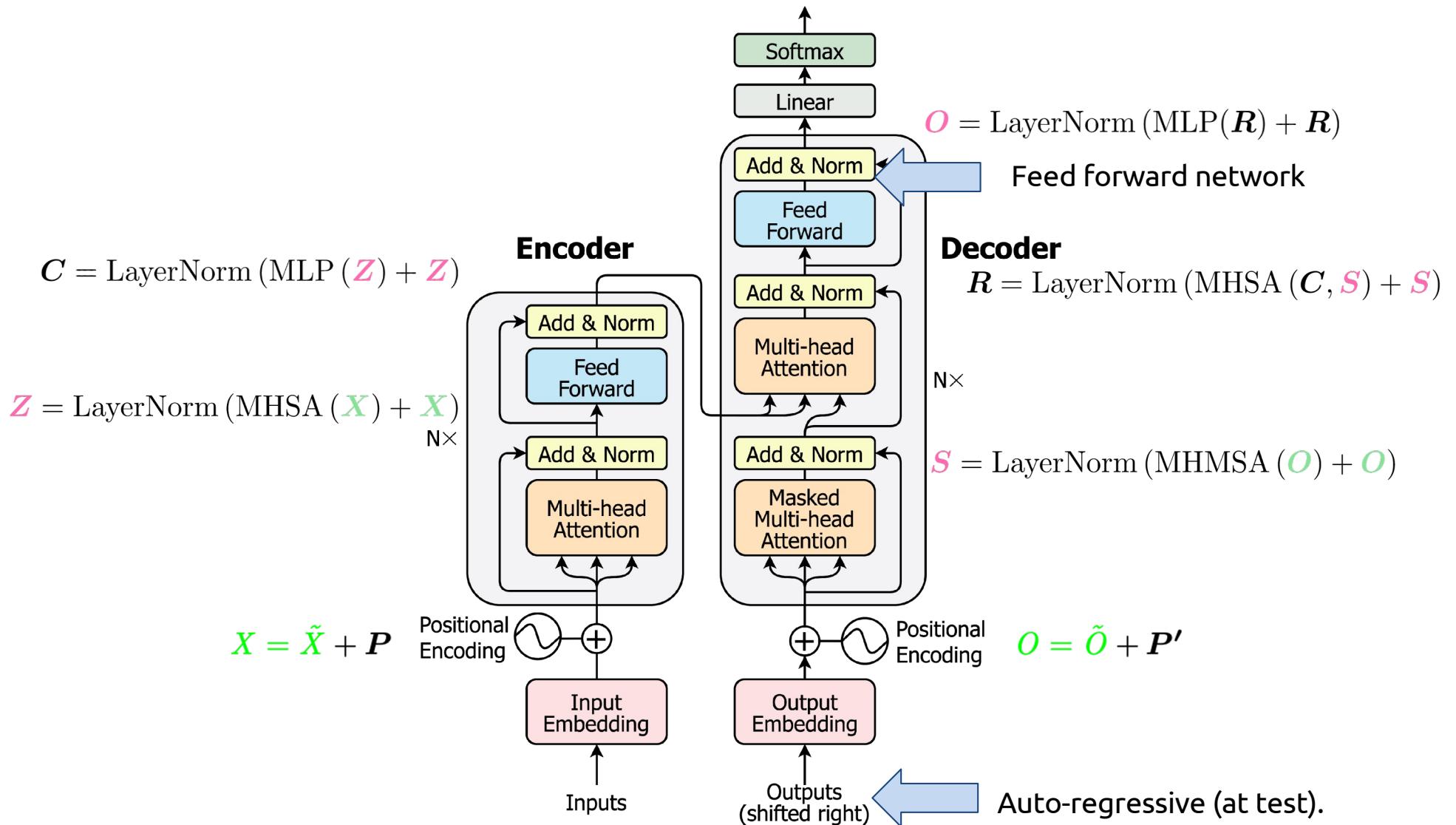
Complete transformer block (4)



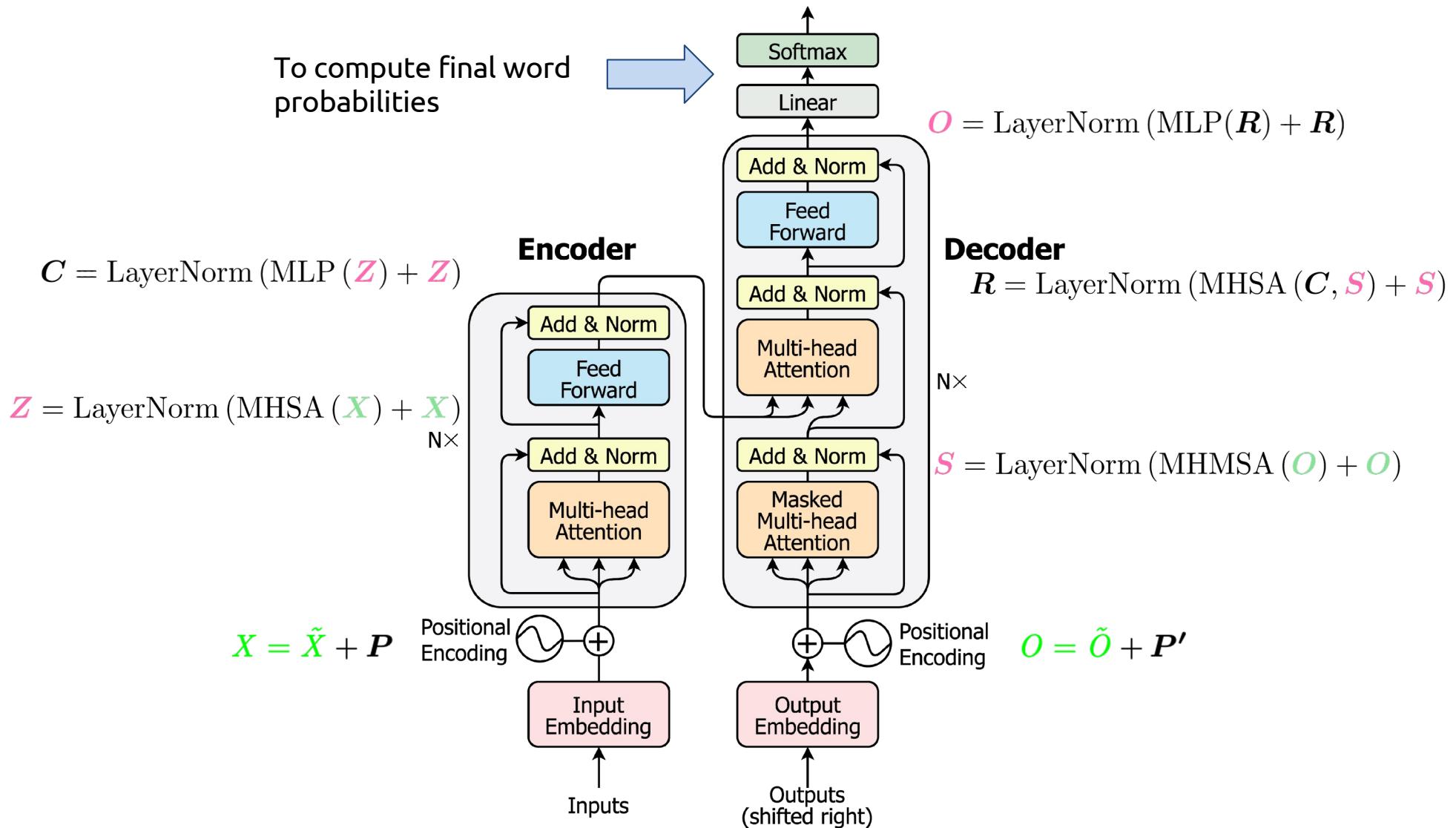
Complete transformer block (5)



Complete transformer block (6)



Complete transformer block (7)



Summary

- Transformers are a **powerful tool** that exploits attention:
 - Each "word" interacts with each other, so maximum interaction distance is **O(1)**
- Their architecture has been ported to **many domains**:
 - language, audio, vision
- **Interpretability** is easier with attention
- They are **hard to optimize** and do not play well with other blocks
- Computing all pairs of interactions means computation **grows quadratically** with sequence length
- More hyperparameters to optimize (num_heads, q/k/v embedding size, etc.)
- The use of **position** representation is still **weak**

Want to know more?

- Tutorials
 - Sebastian Ruder, [Deep Learning for NLP Best Practices # Attention](#), 2017
 - Chris Olah, et al. [Attention and Augmented Recurrent Neural Networks](#) distill.pub 2016
 - Lilian Weg, [The Transformer Family](#). Lil'Log 2020
- Scientific publications
 - **#Perceiver** A. Jaegle, et al. [Perceiver: General perception with iterative attention](#), arXiv:2103.03206, 2021
 - **#Longformer** Beltagy, et al. [Longformer: The long-document transformer.](#)" arXiv preprint arXiv:2004.05150, 2020
 - A. Katharopoulos, et al. [Transformers are RNS: Fast autoregressive transformers with linear attention](#), ICML, 2020.
 - M. Siddhant et al. [Multiplicative Interactions and Where to Find Them](#)". ICLR, 2020
 - Self-attention in language
 - J. Cheng, et al. [Long short-term memory-networks for machine reading](#). arXiv preprint arXiv:1601.06733, 2016
 - Self-attention in images
 - N. Parmar, et al. [Image transformer](#). ICML, 2018
 - X. Wang, et al. [Non-local neural networks](#), CVPR, 2018

Questions?



UNIVERSITAT POLITÈCNICA DE CATALUNYA
BARCELONATECH

Department of Signal Theory
and Communications

Image Processing Group