

Neural Attention & Transformers



Evangelos Kalogerakis

Attention vs Convolution

Convolution has been extremely popular in 2D/3D vision:

- ability to exploit local dependencies in the input data
- highly parallelizable / efficient to compute on GPUs

Attention vs Convolution

Convolution has been extremely popular in 2D/3D vision:

- ability to exploit local dependencies in the input data
- highly parallelizable / efficient to compute on GPUs

Capturing long-range interactions between pixels, points, voxels... is not trivial with convolution

Attention vs Convolution

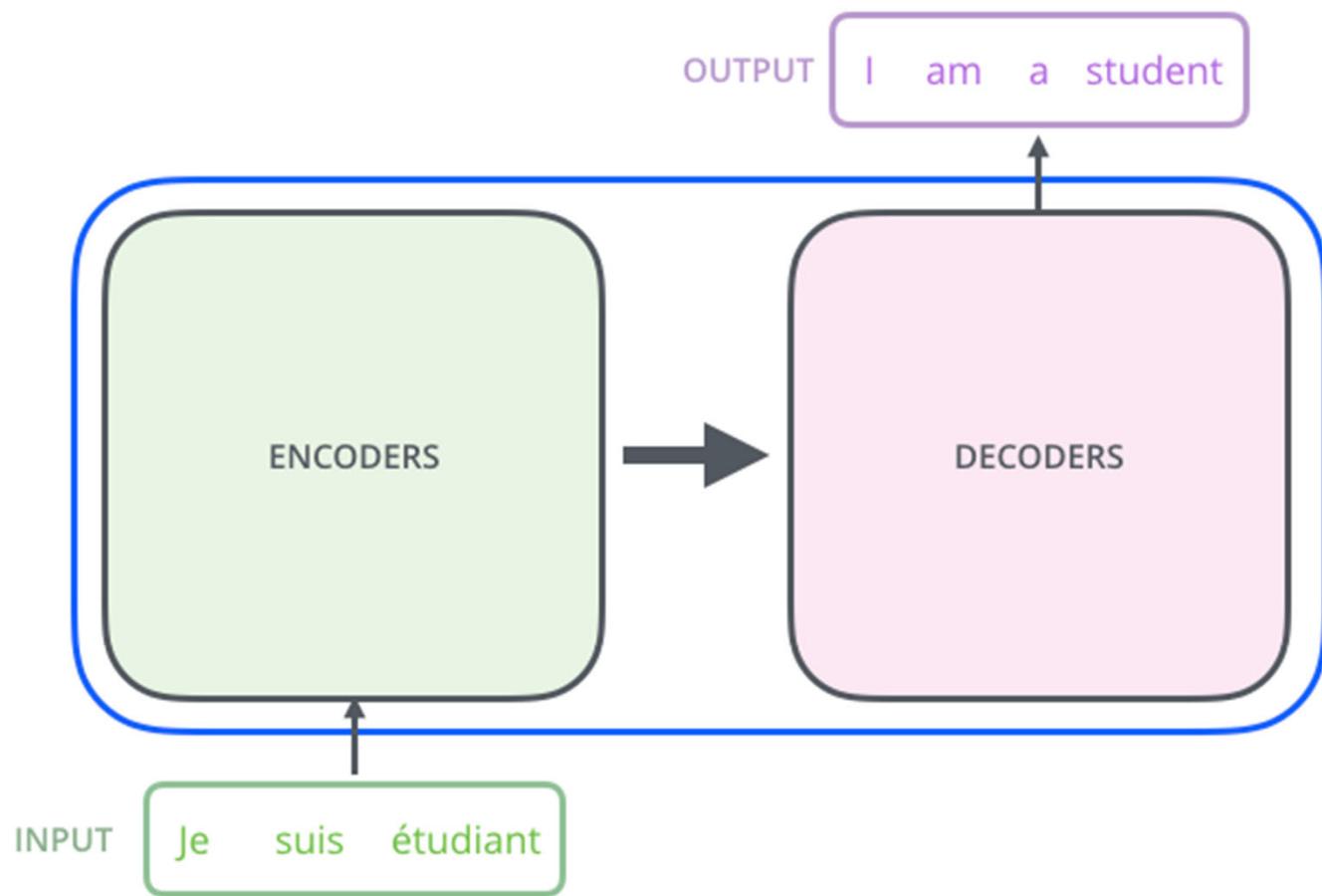
Convolution has been extremely popular in 2D/3D vision:

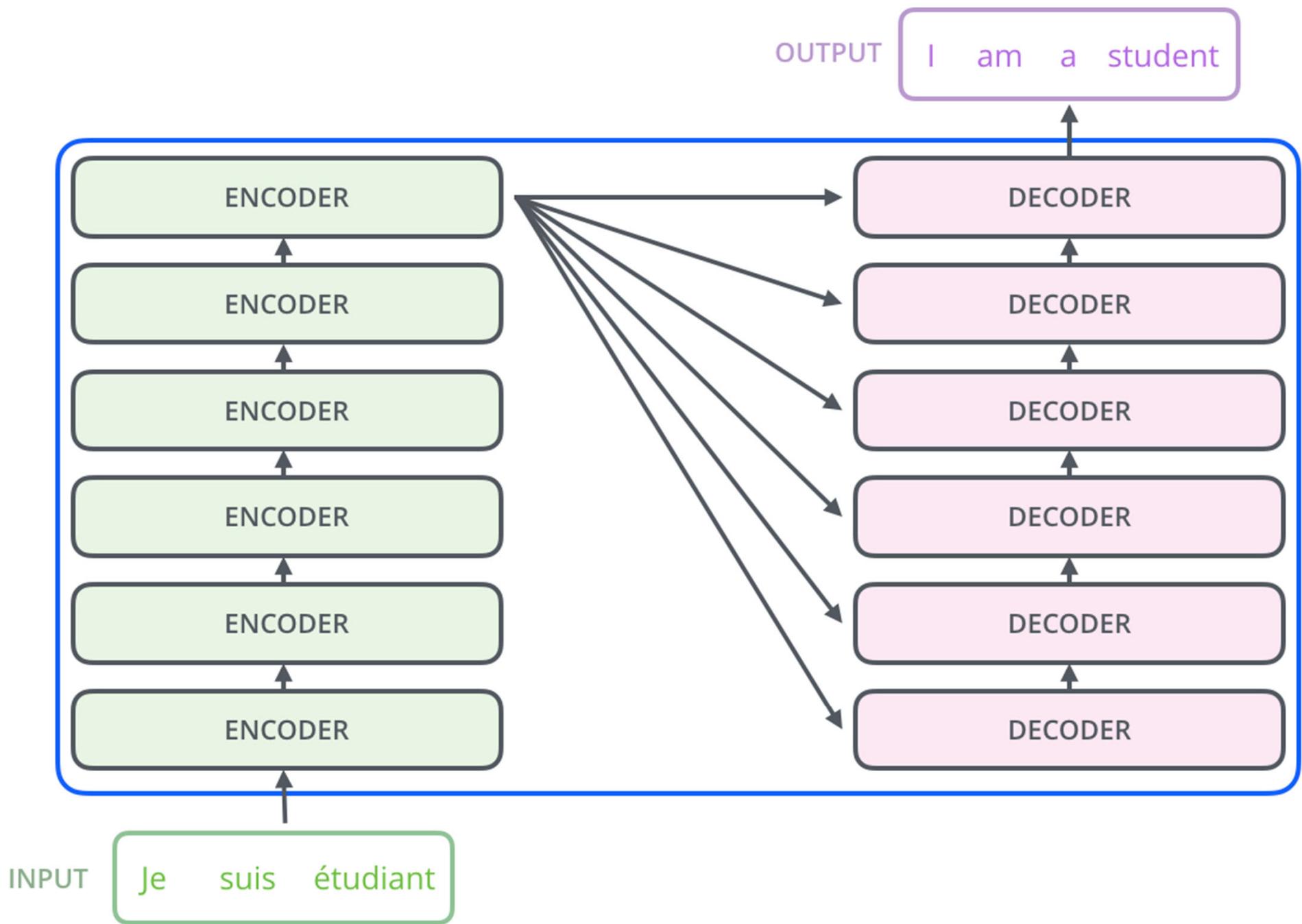
- ability to exploit local dependencies in the input data
- highly parallelizable / efficient to compute on GPUs

Capturing long-range interactions between pixels, points, voxels... is not trivial with convolution

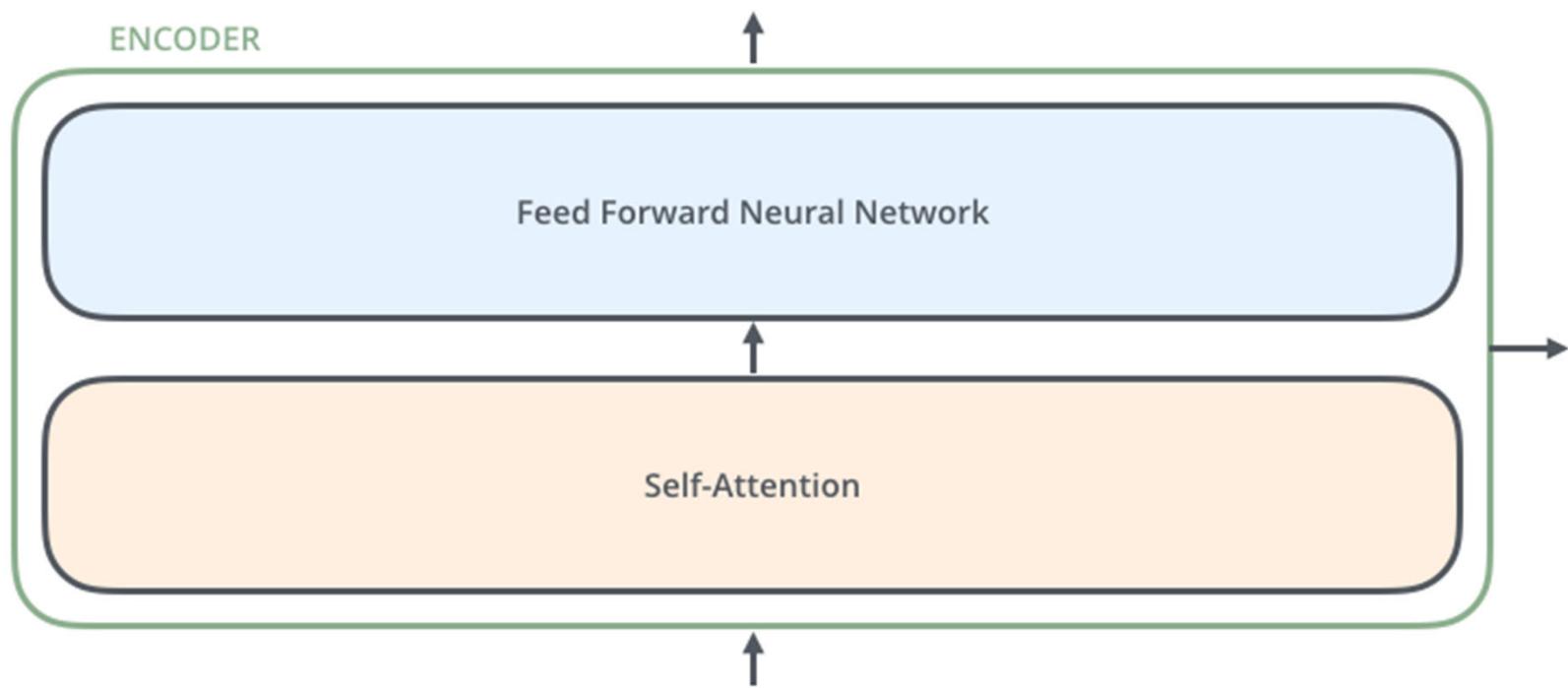
Attention was first introduced in NLP (e.g., machine translation) to capture such interactions!

Machine Translation

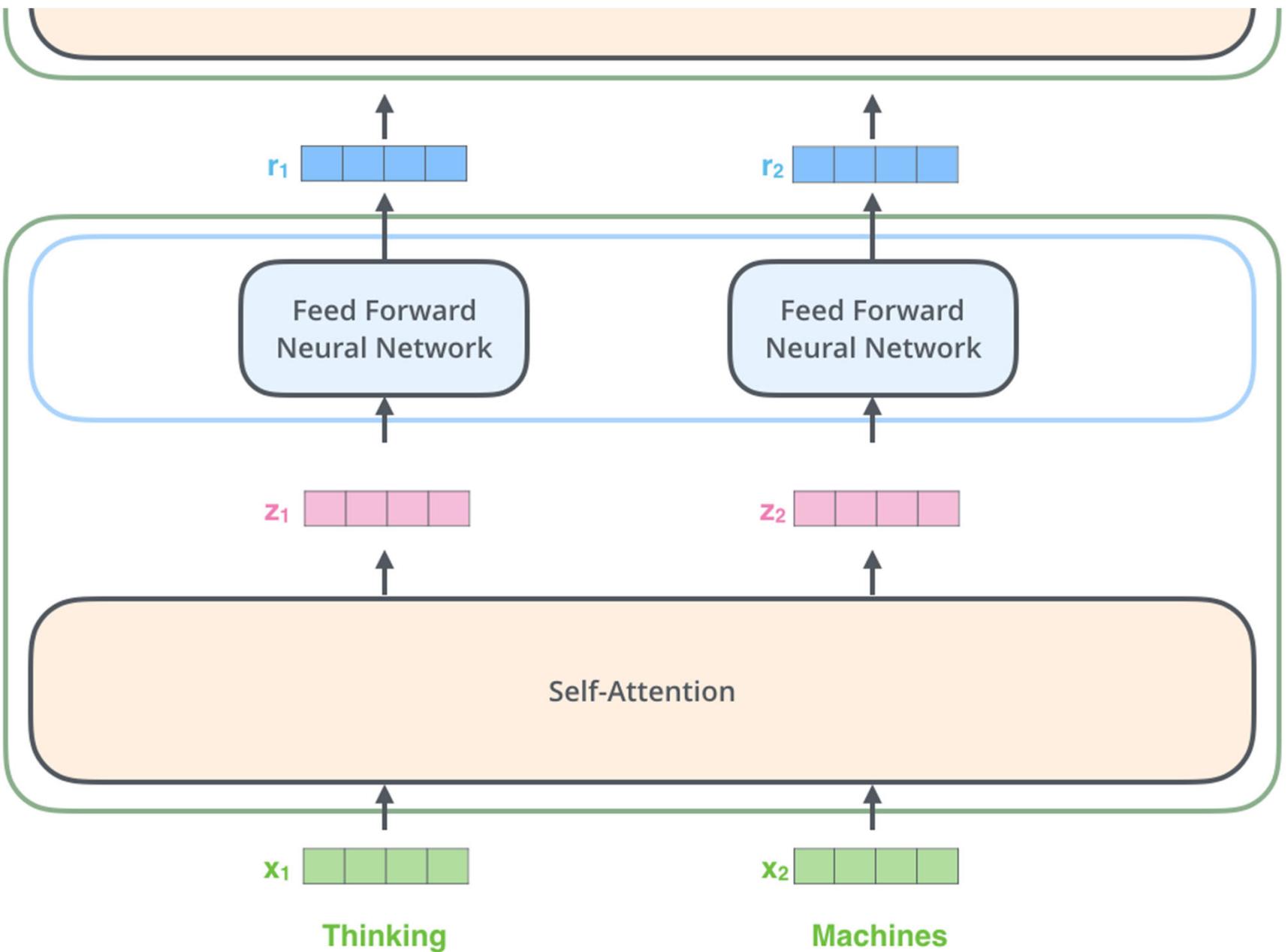




The Encoder



ENCODER #2



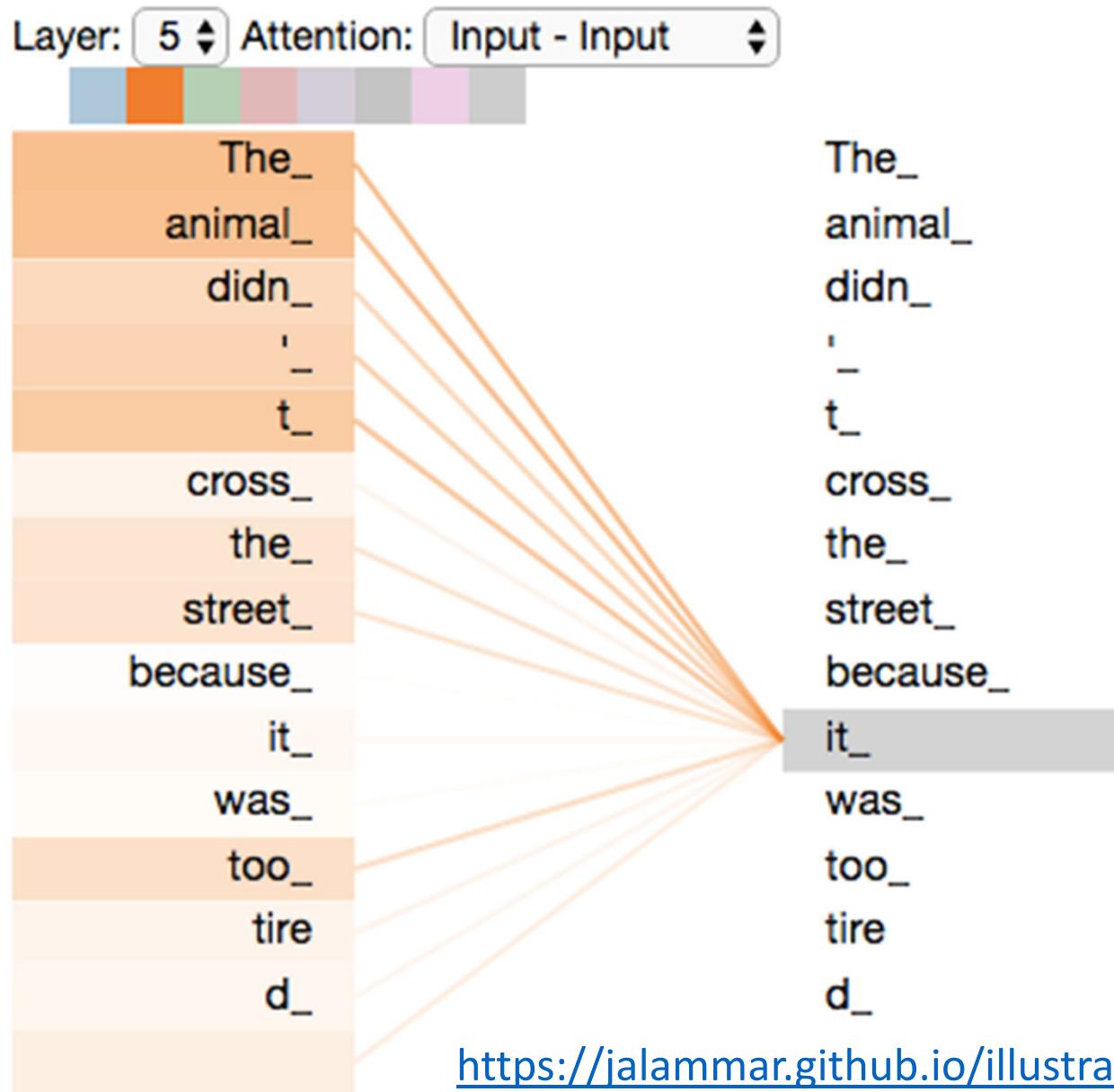
What's the deal with attention?

Say we want to translate the following sentence:

“The animal didn't cross the street because it was too tired”

What does “***it***” in this sentence refer to? Is it referring to the street or to the animal?

Attention: allows it to look at other tokens in the input sequence that can help lead to a better encoding for this token!



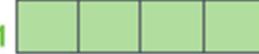
Step A: create three vectors from each of the encoder's input token embeddings:
a Query vector, a Key vector, and a Value vector.

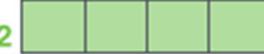
Input

Thinking

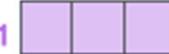
Machines

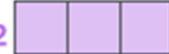
Embedding

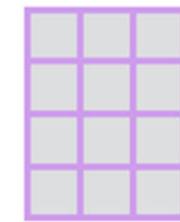
X_1 

X_2 

Queries

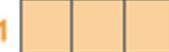
q_1 

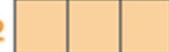
q_2 



W^Q

Keys

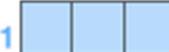
k_1 

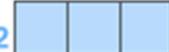
k_2 



W^K

Values

v_1 

v_2 



W^V

How to obtain these **query**, **key**, **value** vectors?
Matrix multiplication with (learned) matrices \mathbf{W}^q , \mathbf{W}^k , \mathbf{W}^v !

$$\mathbf{X} \times \mathbf{W}^q = \mathbf{Q}$$

A diagram illustrating matrix multiplication. On the left, a green 3x4 matrix labeled \mathbf{X} is multiplied by a purple 4x4 matrix labeled \mathbf{W}^q . The result is a purple 3x4 matrix labeled \mathbf{Q} .

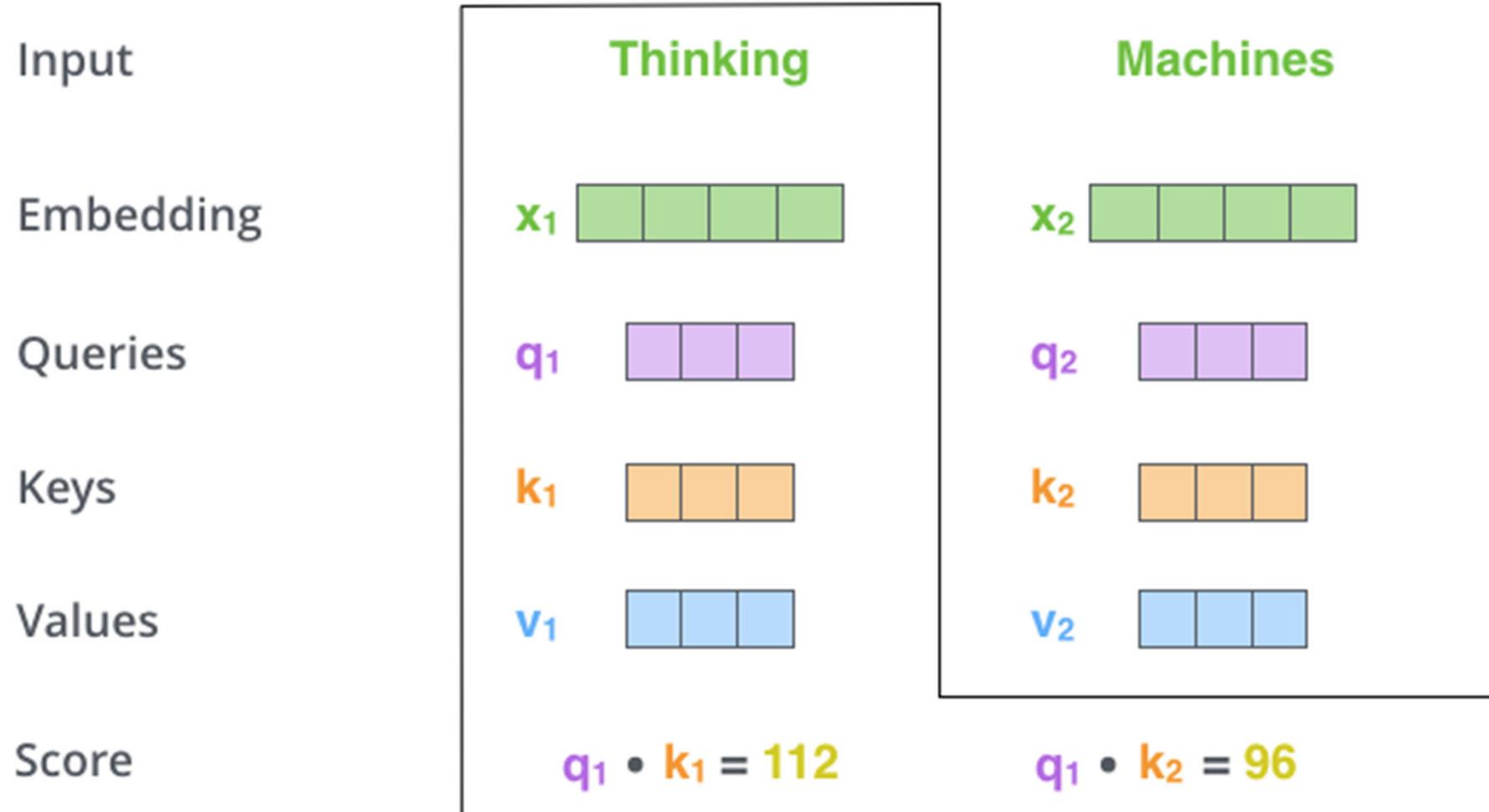
$$\mathbf{X} \times \mathbf{W}^k = \mathbf{K}$$

A diagram illustrating matrix multiplication. On the left, a green 3x4 matrix labeled \mathbf{X} is multiplied by an orange 4x4 matrix labeled \mathbf{W}^k . The result is an orange 3x4 matrix labeled \mathbf{K} .

$$\mathbf{X} \times \mathbf{W}^v = \mathbf{V}$$

A diagram illustrating matrix multiplication. On the left, a green 3x4 matrix labeled \mathbf{X} is multiplied by a blue 4x4 matrix labeled \mathbf{W}^v . The result is a blue 3x4 matrix labeled \mathbf{V} .

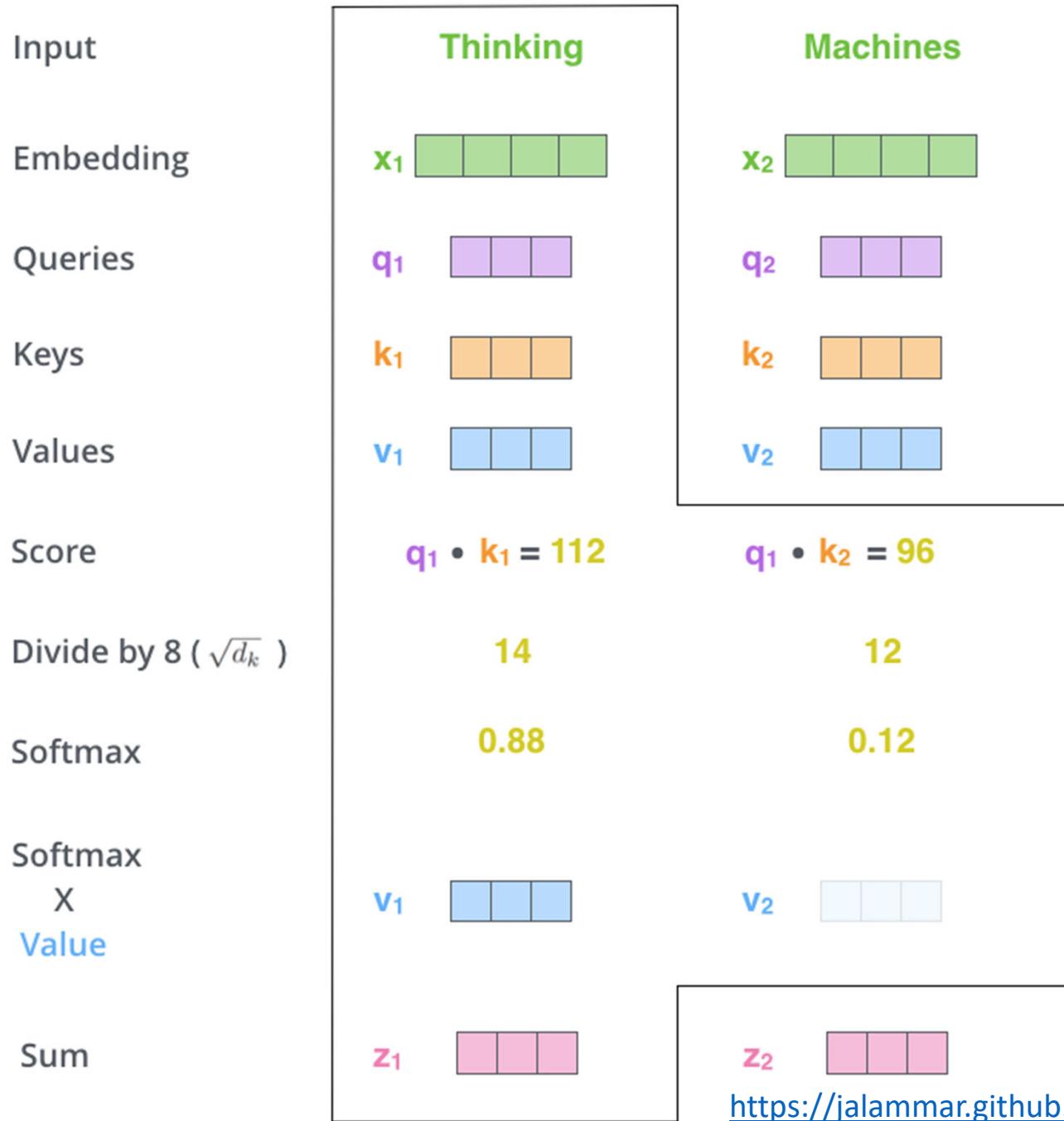
Step B: calculate a score determining how much focus to place on other parts of the input sentence as we encode a token at a certain **position**.



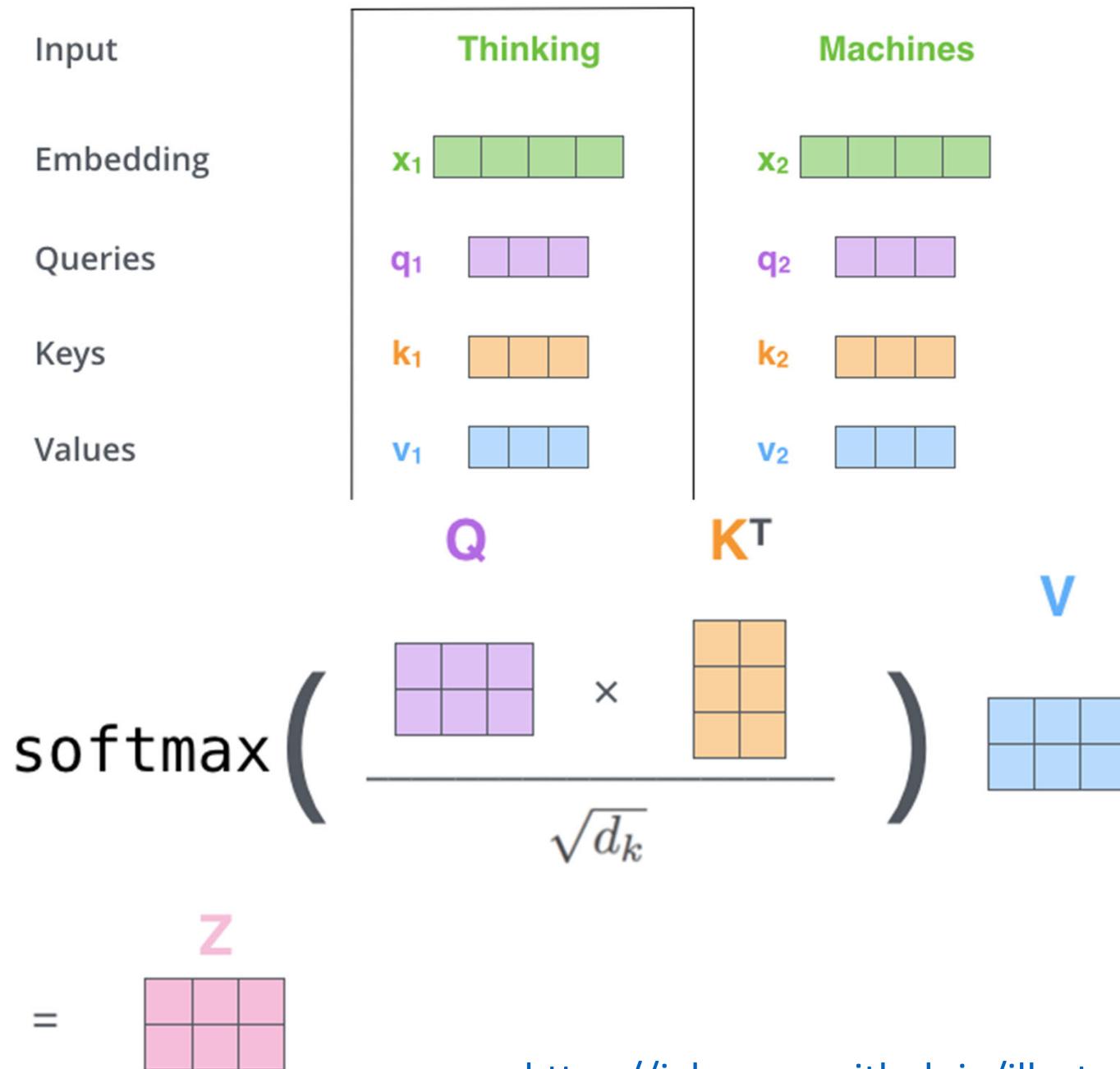
Step C: divide the scores with the sqrt of the dimension of the vectors (64) -- helps with stable gradients & pass them through softmax so they're positive & sum up to 1.

Input	Thinking		Machines	
Embedding	x_1		x_2	
Queries	q_1		q_2	
Keys	k_1		k_2	
Values	v_1		v_2	
Score		$q_1 \cdot k_1 = 112$		$q_1 \cdot k_2 = 96$
Divide by 8 ($\sqrt{d_k}$)		14		12
Softmax		0.88		0.12

Step D: Multiply each value vector by the softmax score (keep intact the values of the tokens we want to focus on & drown-out irrelevant ones). Sum up the weighted value vectors.

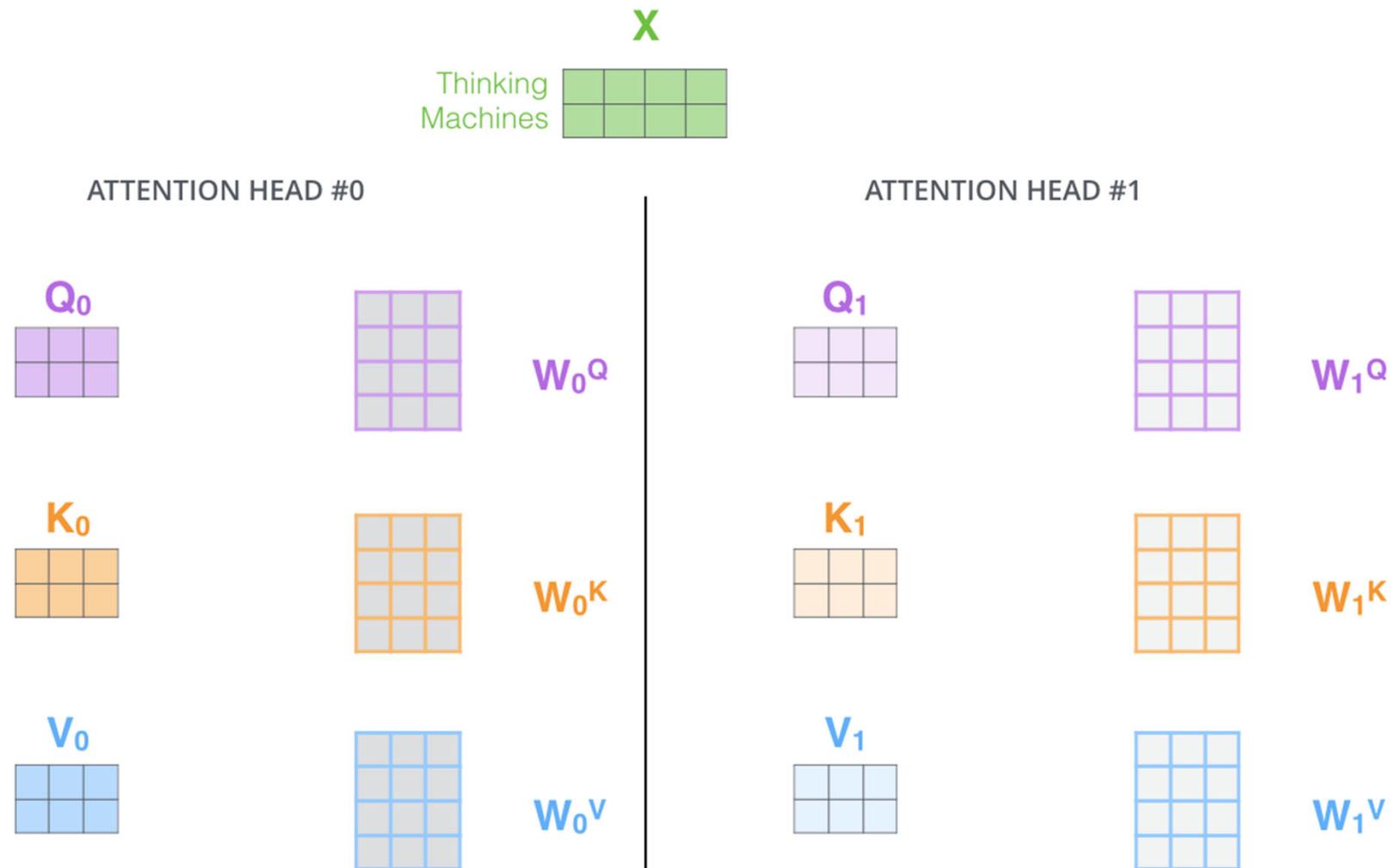


Step D: Multiply each value vector by the softmax score (keep intact the values of the tokens we want to focus on & drown-out irrelevant ones). Sum up the weighted value vectors.

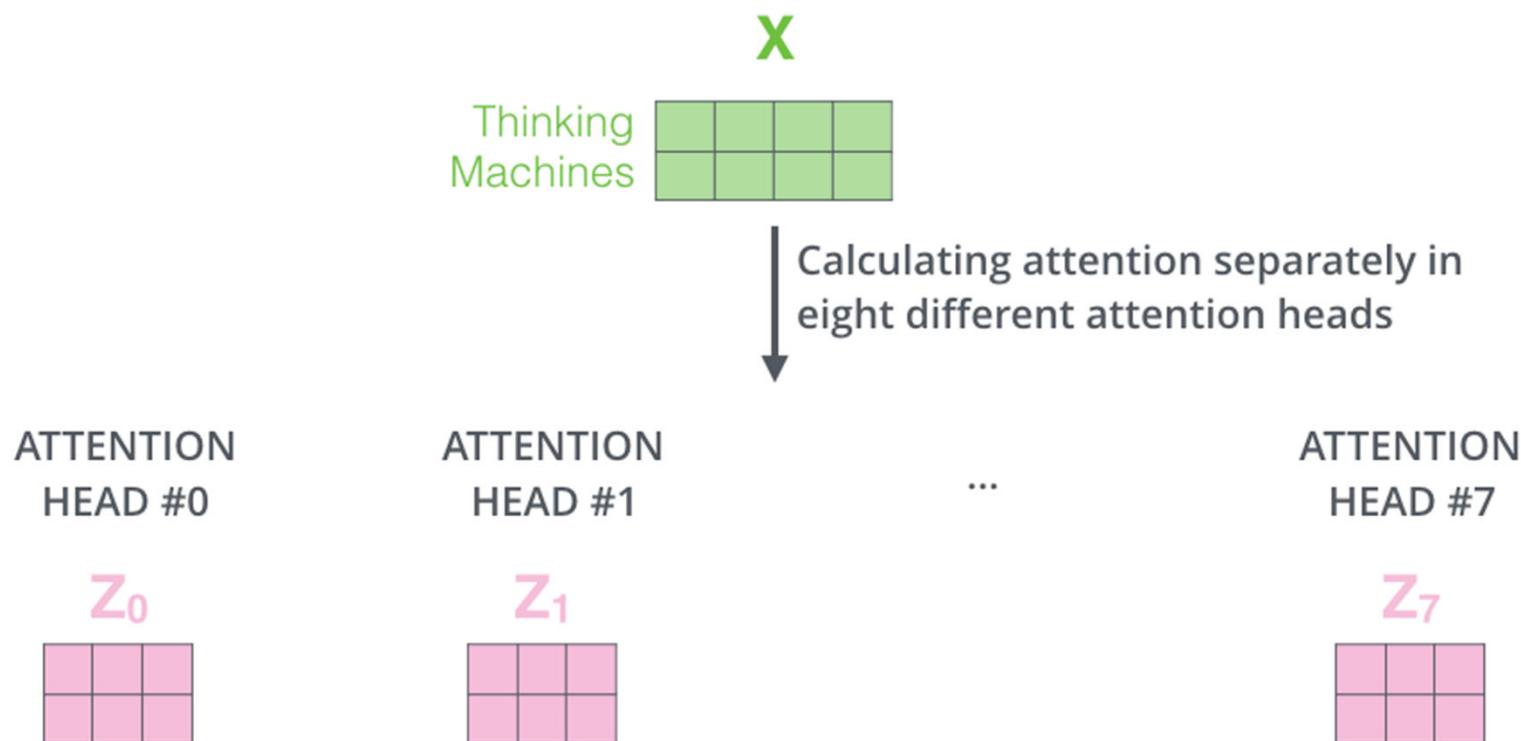


The beast with many “heads”!

Use multiple attention transformations to encode a token in multiple “representation” contexts e.g., in the sentence: ***“The animal didn't cross the street because it was too tired”*** encode the context of “it” wrt “animal” & also encode its context wrt “was too tired”



Processing with multiple attention heads...



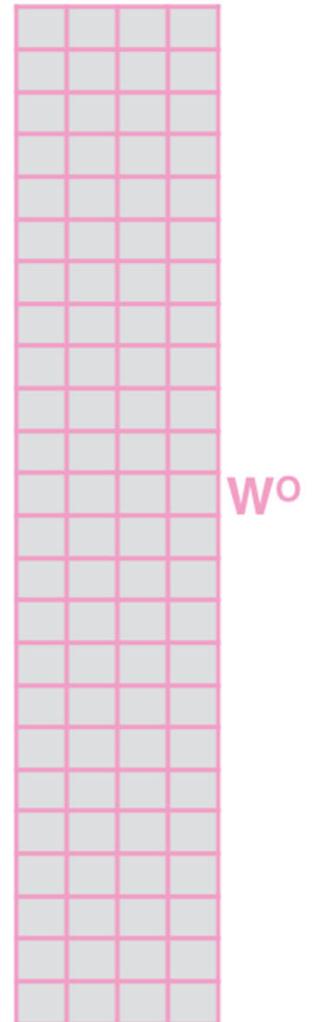
Processing with multiple attention heads...

1) Concatenate all the attention heads



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

\times



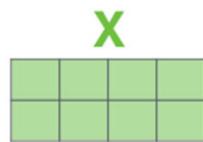
3) The result would be the Z matrix that captures information from all the attention heads. We can send this forward to the FFNN

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

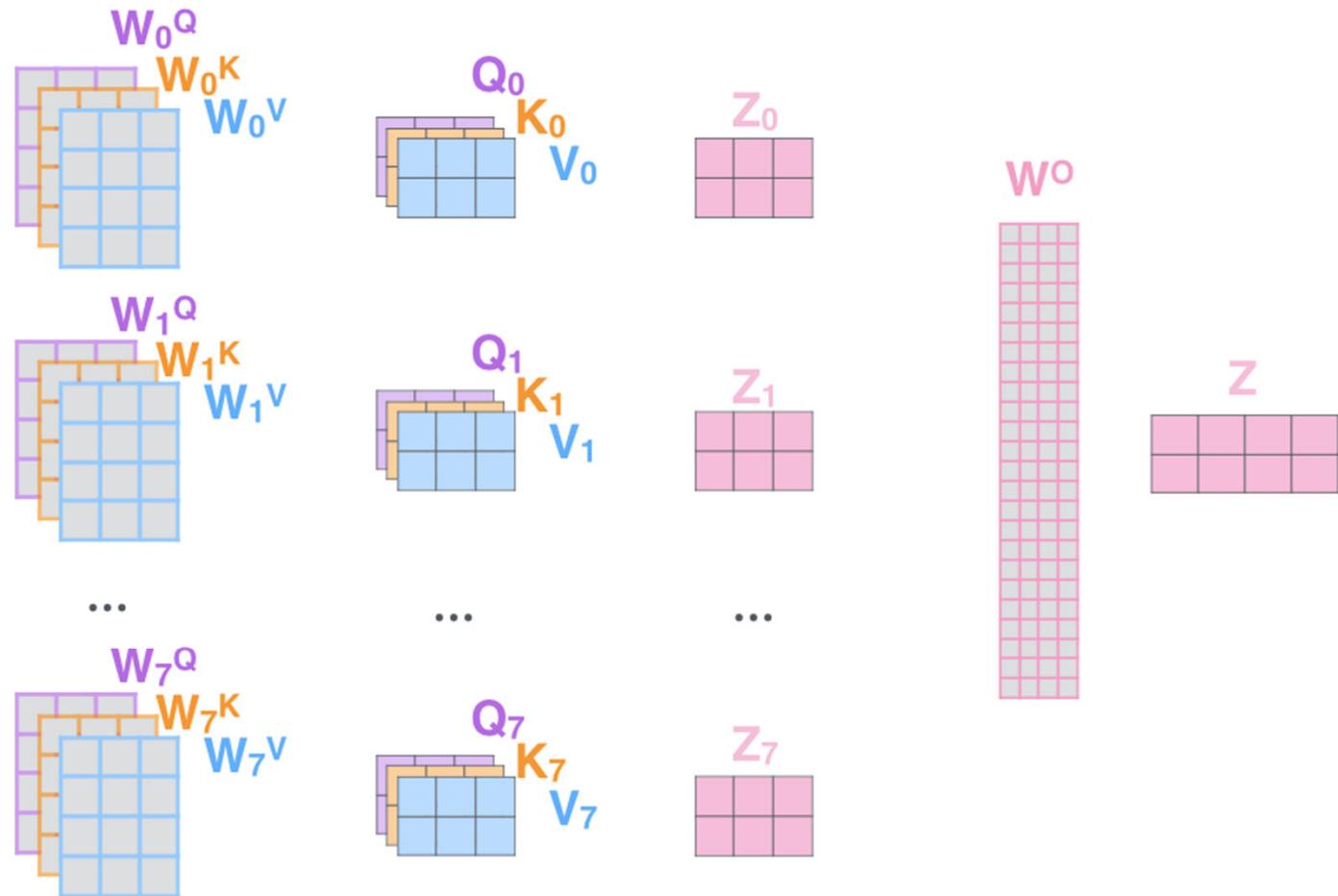
Summary...

- 1) This is our input sentence* X
- 2) We embed each word* R
- 3) Split into 8 heads. We multiply X or R with weight matrices W_0^Q, W_0^K, W_0^V
- 4) Calculate attention using the resulting $Q/K/V$ matrices
- 5) Concatenate the resulting Z matrices, then multiply with weight matrix W^O to produce the output of the layer

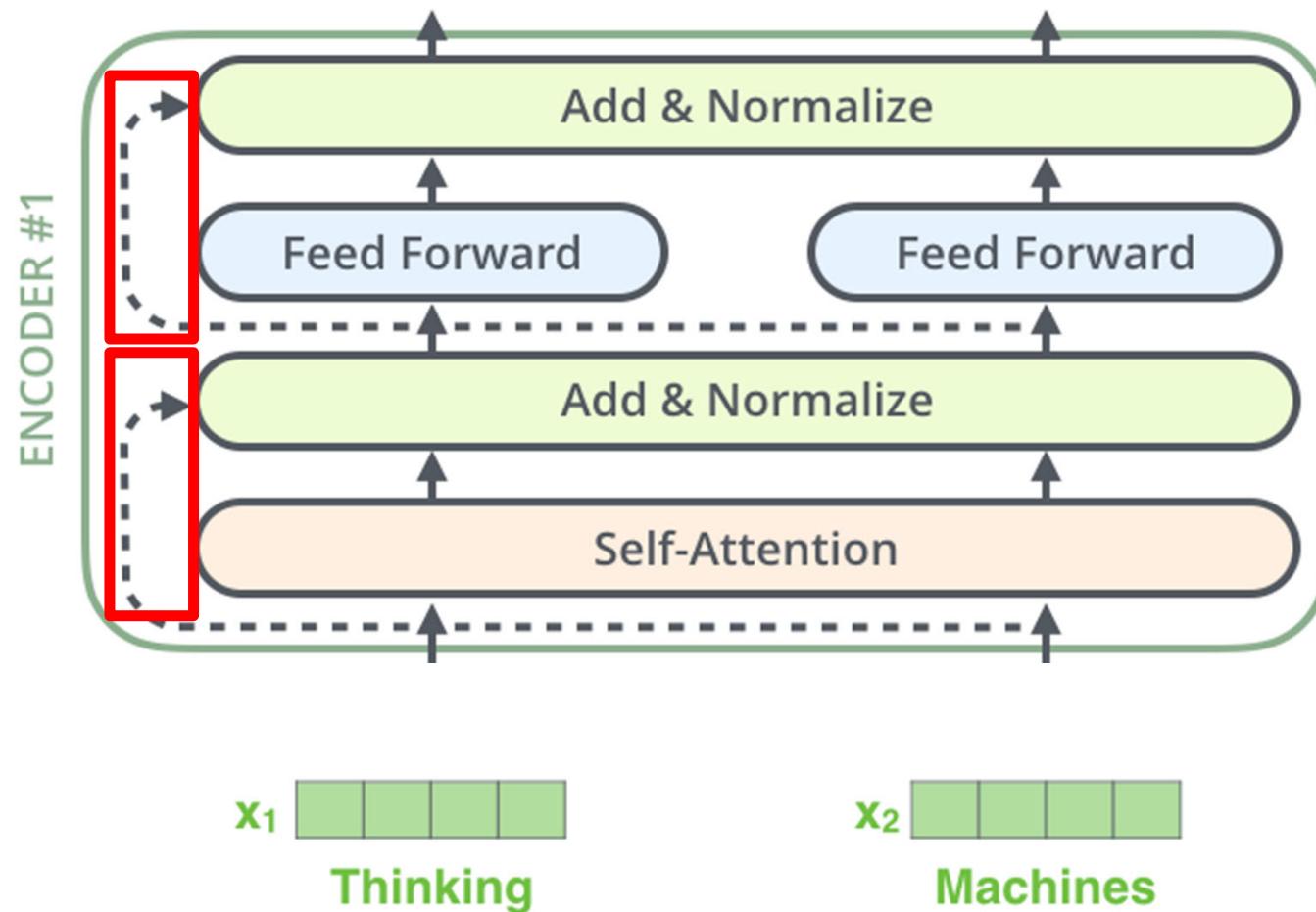
Thinking
Machines



* In all encoders other than #0, we don't need embedding. We start directly with the output of the encoder right below this one

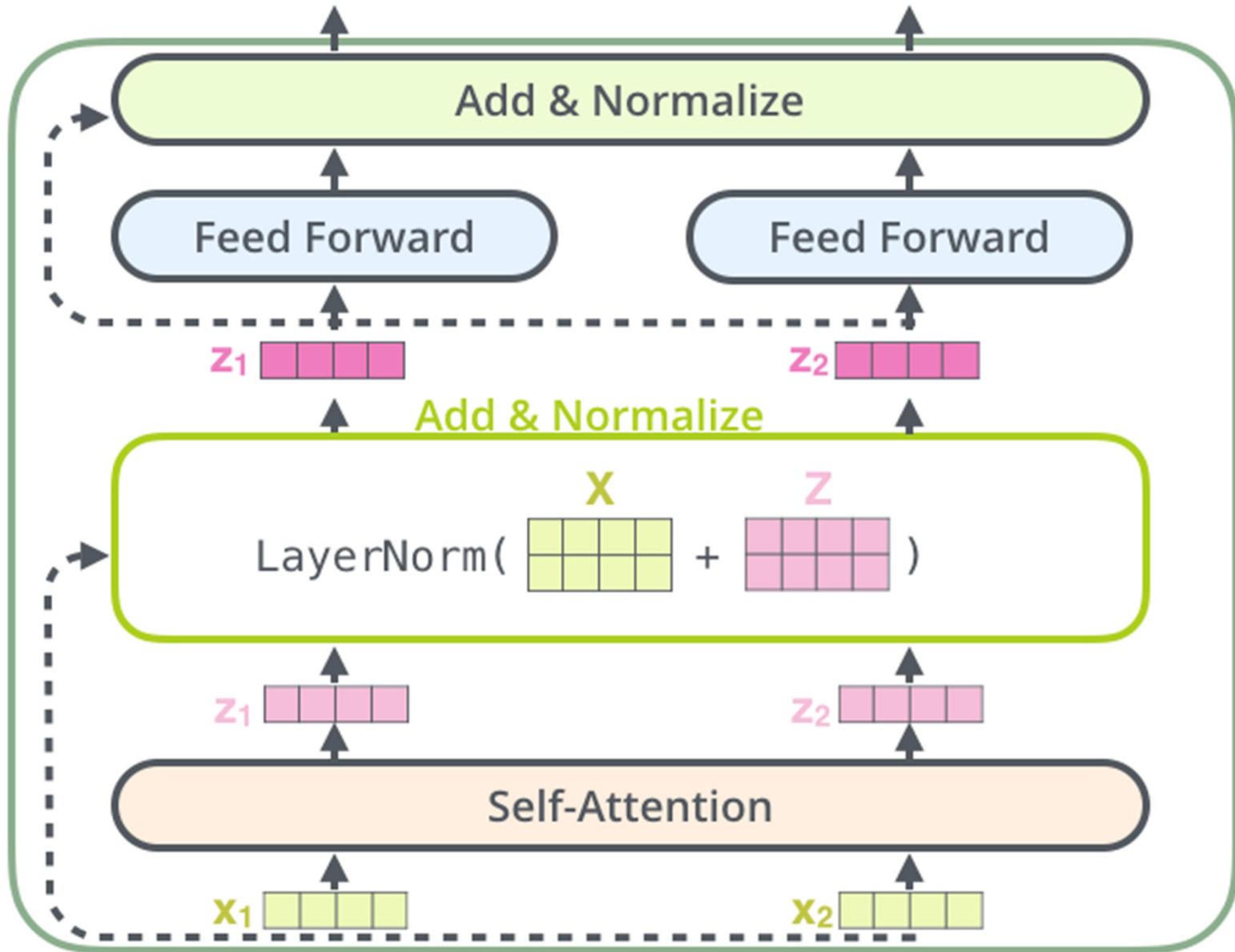


More details: residual connections



More details: layer normalization

ENCODER #1



Thinking

Machines

Another detail: so far the model has no idea of what the **position** of the token is in a sentence!

The animal didn't cross the street because it was too tired

0 1 2 3 4 5 6 7 8 9 10 ?

Thinking Machines

0 1 ?

Another detail: so far the model has no idea of what the **position** of the token is in a sentence!

The animal didn't cross the street because it was too tired

0 1 2 3 4 5 6 7 8 9 10 ?

Thinking Machines

0 1 ?

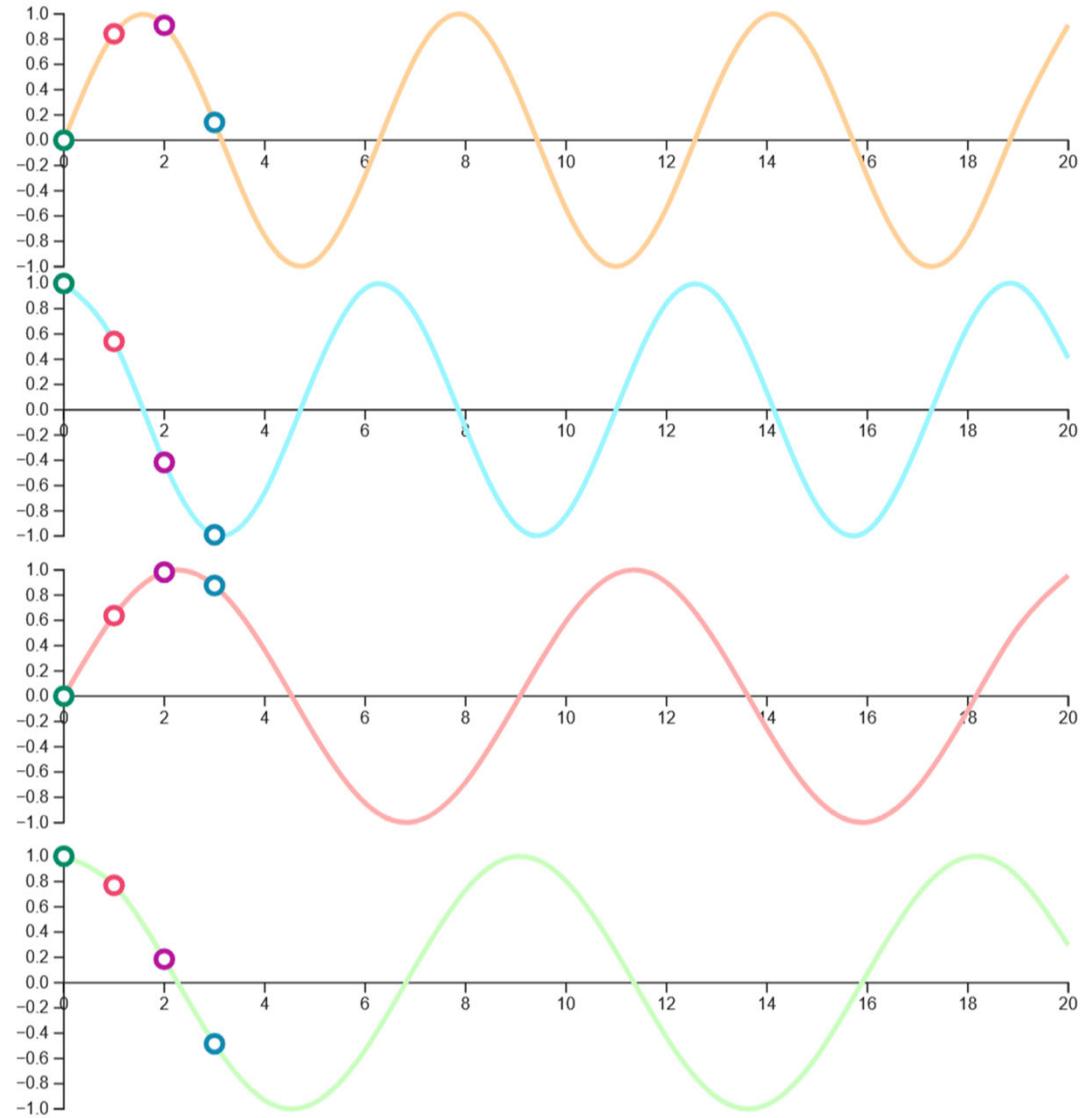
Problem: if the largest sentence is 50 tokens long during training, how can the model generalize to sentences longer than 50 during testing?

Another detail: so far the model has no idea of what the **position** of the token is in a sentence!

$$PE_{(pos,2i)} = \sin\left(\frac{pos}{10000^{2i/d_{\text{model}}}}\right)$$

$$PE_{(pos,2i+1)} = \cos\left(\frac{pos}{10000^{2i/d_{\text{model}}}}\right)$$

Positional encoding visualization



	p_0	p_1	p_2	p_3	$i=0$
	0.000	0.841	0.909	0.141	
	1.000	0.540	-0.416	-0.990	$i=1$
	0.000	0.638	0.983	0.875	
	1.000	0.770	0.186	-0.484	$i=3$

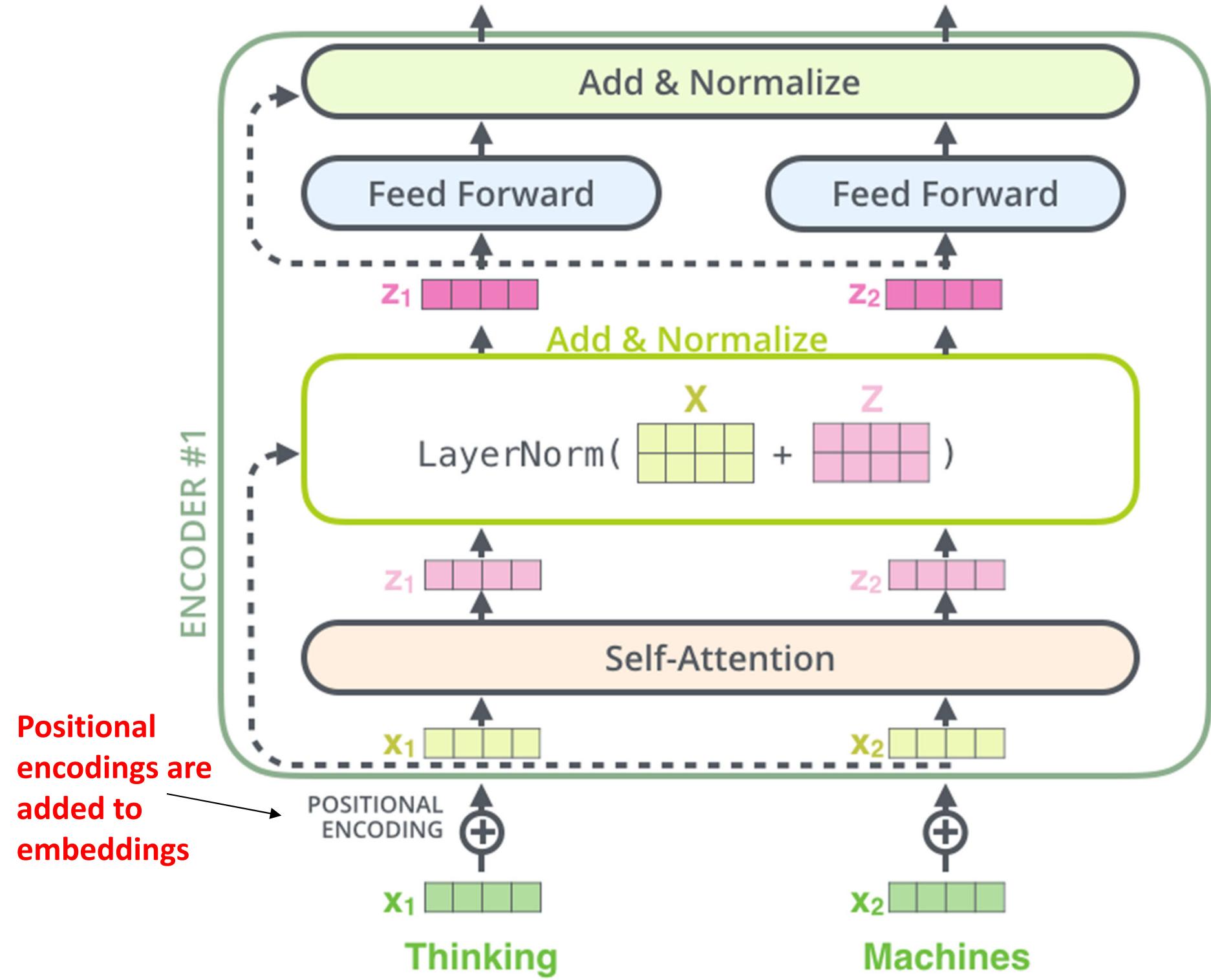
Positional Encoding

$$PE_{(pos,2i)} = \sin\left(\frac{pos}{10000^{2i/d_{model}}}\right)$$

$$PE_{(pos,2i+1)} = \cos\left(\frac{pos}{10000^{2i/d_{model}}}\right)$$

Settings: $d = 50$

The value of each positional encoding depends on the *position* (*pos*) and *dimension* (*d*). We calculate result for every *index* (*i*) to get the whole vector.

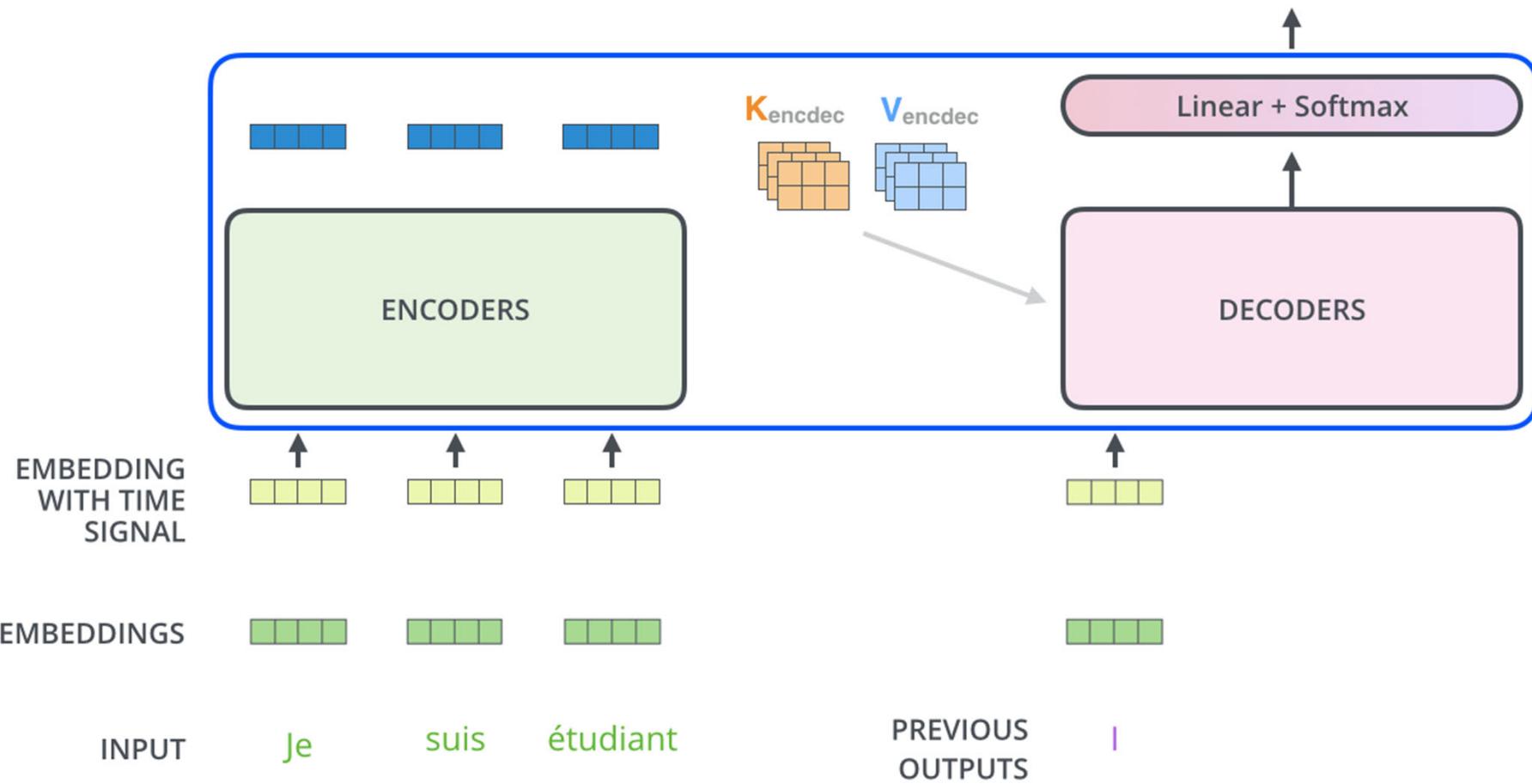


The decoder

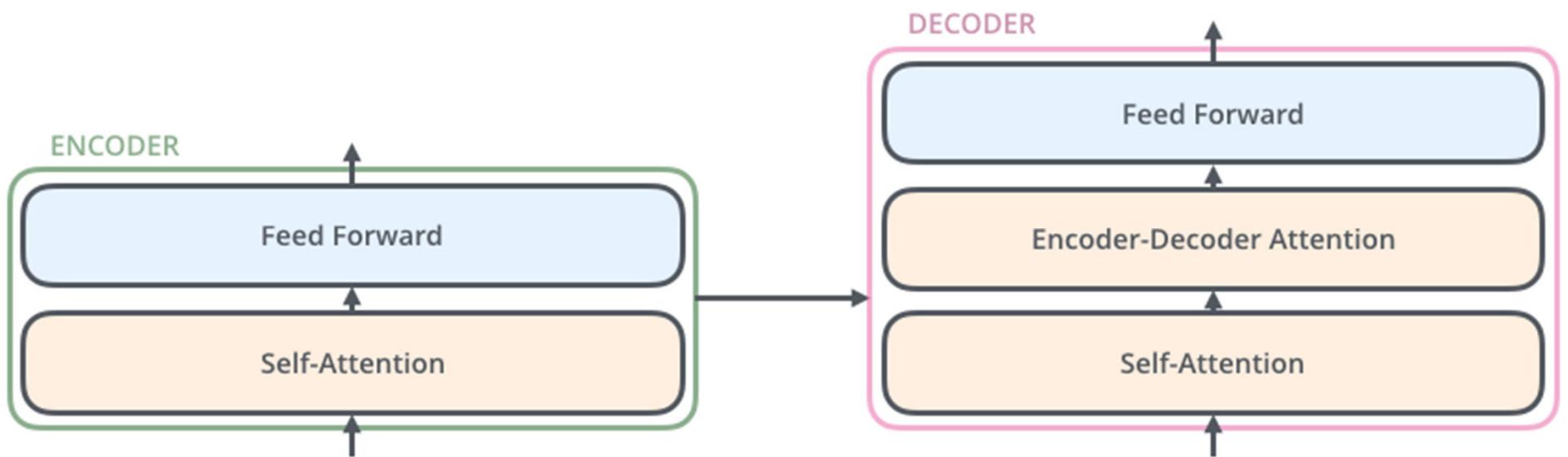
Decoding time step: 1 2 3 4 5 6

OUTPUT

|



The decoder

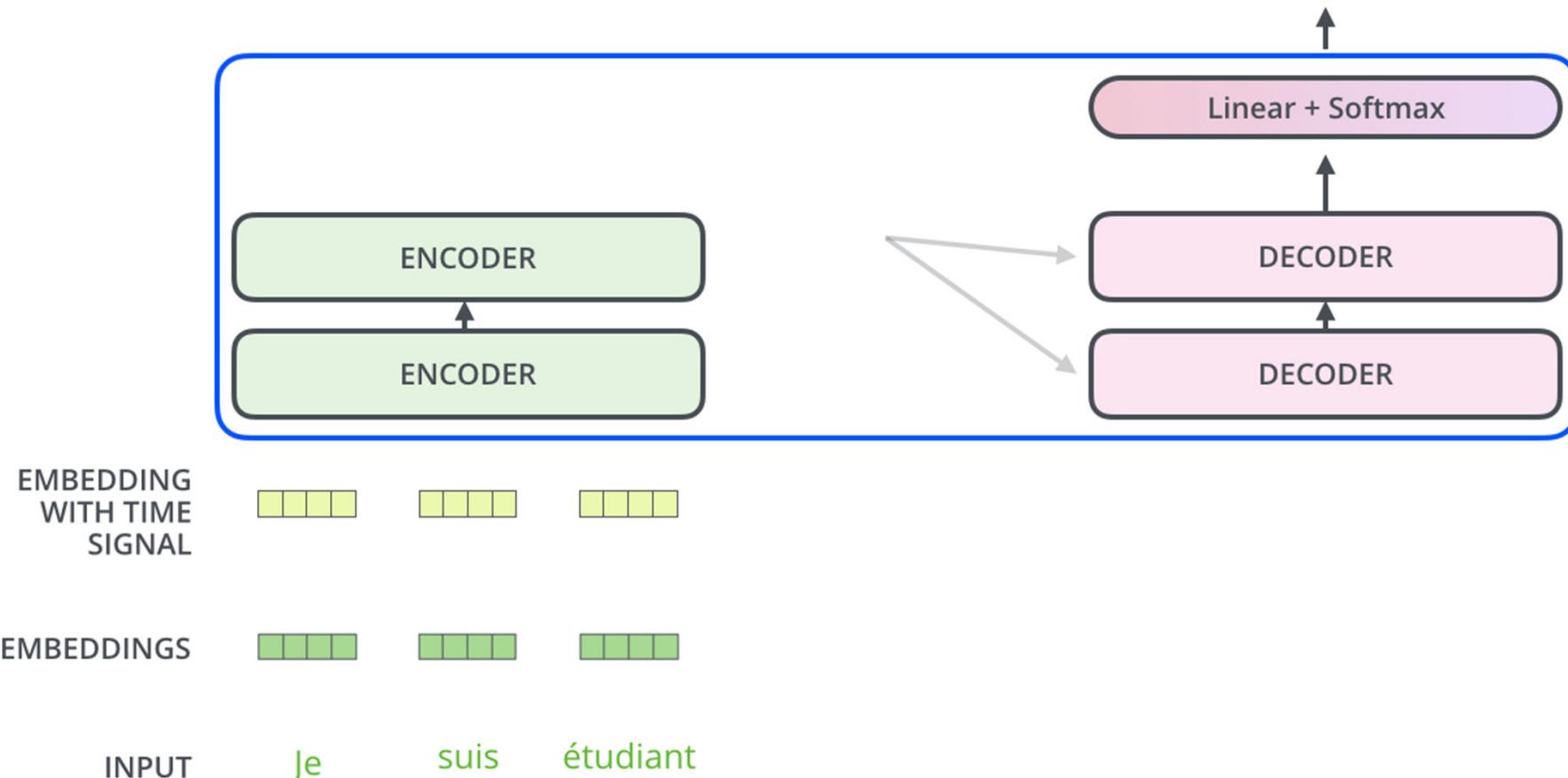


the self-attention layer is similar to the encoder's self attention, yet attends only to earlier positions in the output sequence

The output of the last encoder block provides an attention vector \mathbf{K} and \mathbf{V} for each **input token**. These are used by the decoder blocks in the “**encoder-decoder attention**” layers helping the decoder focus on appropriate places in the input sentence!

Decoding time step: 1 2 3 4 5 6

OUTPUT



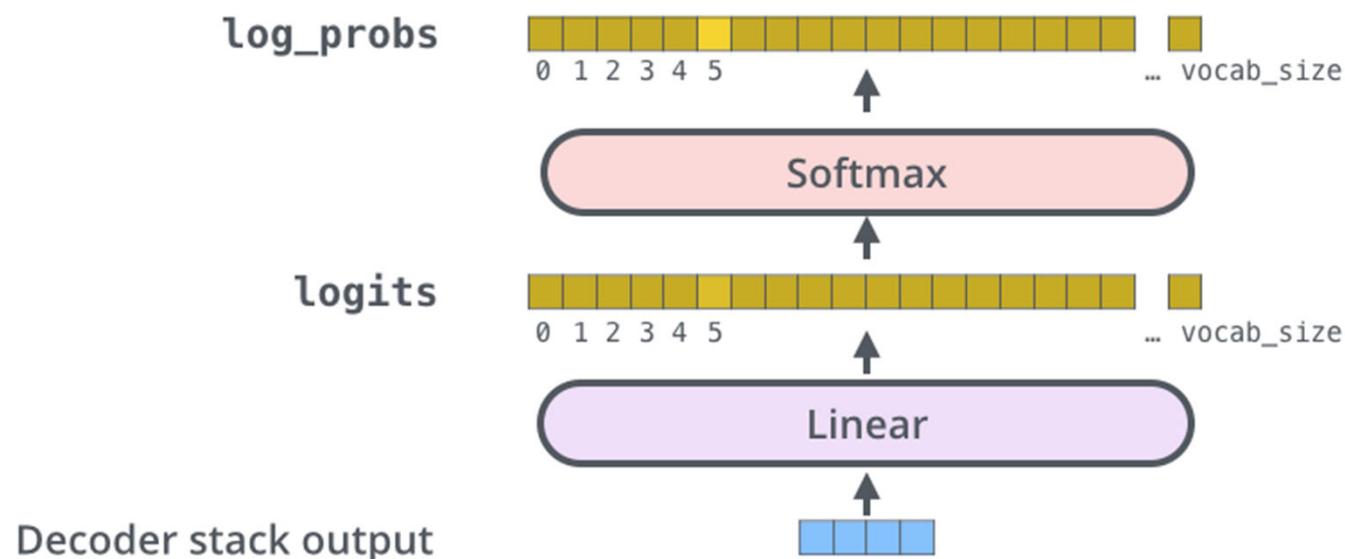
The final layer of the decoder

Which word in our vocabulary
is associated with this index?

am

Get the index of the cell
with the highest value
(argmax)

5



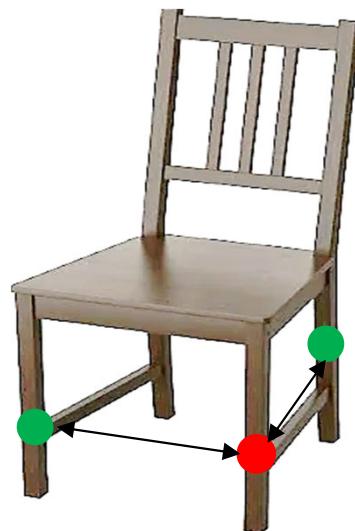
Attention for visual data (images/shapes)

Help encoder look at other pixels/patches in image while encoding a pixel/patch due to various relations that may exist in objects or scenes in general



Capture visual context

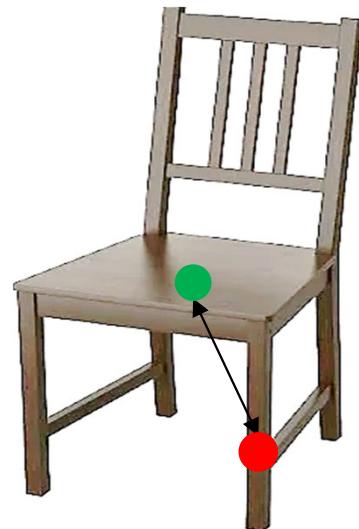
Help encoder look at other pixels/patches in image while encoding a pixel/patch due to various relations that may exist in objects or scenes in general



e.g., symmetries

Capture visual context

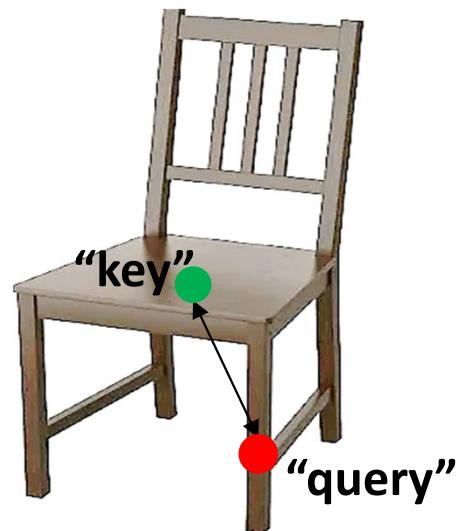
Help encoder look at other pixels/patches in image while encoding a pixel/patch due to various relations that may exist in objects or scenes in general



... or any pixel/patch relations may help recognize objects

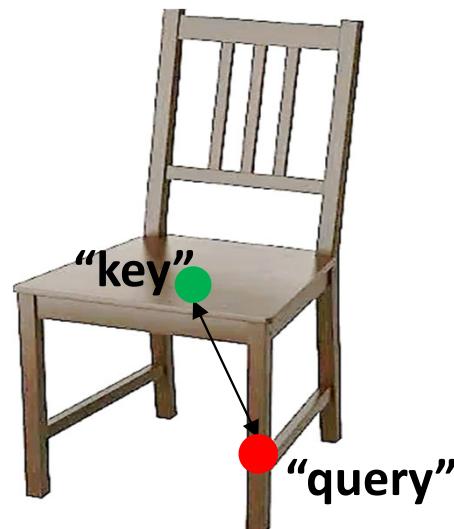
Attention transformations

How much a pixel/patch (query) is related to others (keys)?



Attention transformations

How much a pixel/patch (query) is related to others (keys)?



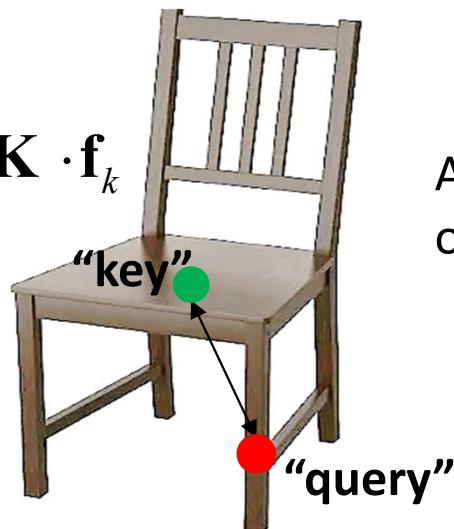
$$q(\mathbf{f}_q) = \mathbf{Q} \cdot \mathbf{f}_q$$

A linear transformation on the input feature vector of the query pixel/patch
(e.g., MLP on input point/patch raw features)

Attention transformations

How much a pixel/patch (query) is related to others (keys)?

$$k(\mathbf{f}_k) = \mathbf{K} \cdot \mathbf{f}_k$$



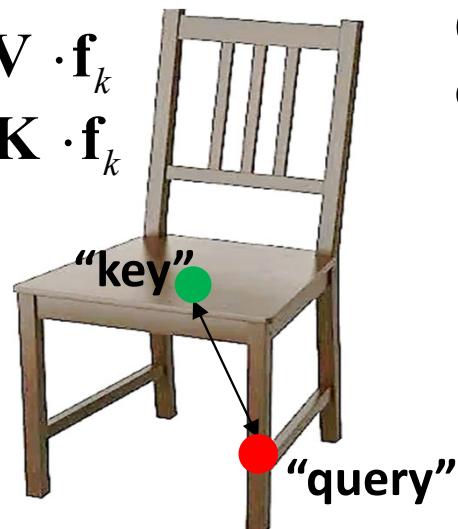
Another transformation on the input feature vector
of the key point

$$q(\mathbf{f}_q) = \mathbf{Q} \cdot \mathbf{f}_q$$

Attention transformations

How much a pixel/patch (query) is related to others (keys)?

$$v(\mathbf{f}_k) = \mathbf{V} \cdot \mathbf{f}_k$$
$$k(\mathbf{f}_k) = \mathbf{K} \cdot \mathbf{f}_k$$



One more transformation on feature vector
of the key point i.e., we have a “key-value” pair

$$q(\mathbf{f}_q) = \mathbf{Q} \cdot \mathbf{f}_q$$

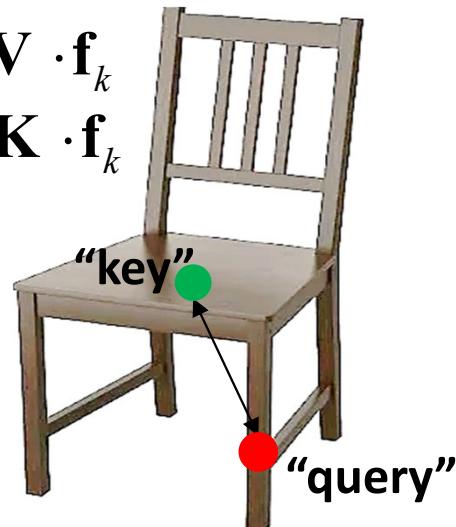
Attention transformations

How much a pixel/patch (query) is related to others (keys)?

Attention “score”: $a(\mathbf{f}_q, \mathbf{f}_k) = k(\mathbf{f}_k) \cdot q(\mathbf{f}_q)$

(dot product between key-query vectors)

$$v(\mathbf{f}_k) = \mathbf{V} \cdot \mathbf{f}_k$$
$$k(\mathbf{f}_k) = \mathbf{K} \cdot \mathbf{f}_k$$



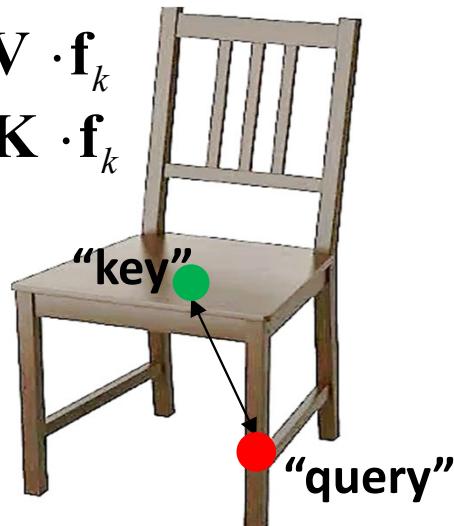
$$q(\mathbf{f}_q) = \mathbf{Q} \cdot \mathbf{f}_q$$

Attention score

How much a pixel/patch (query) is related to others (keys)?

Attention “score”: $\hat{a}(\mathbf{f}_q, \mathbf{f}_k) = \frac{\exp\{ a(\mathbf{f}_q, \mathbf{f}_k) \}}{\sum_{k'} \exp\{ a(\mathbf{f}_q, \mathbf{f}_{k'}) \}}$
(i.e., use softmax)

$$v(\mathbf{f}_k) = \mathbf{V} \cdot \mathbf{f}_k$$
$$k(\mathbf{f}_k) = \mathbf{K} \cdot \mathbf{f}_k$$

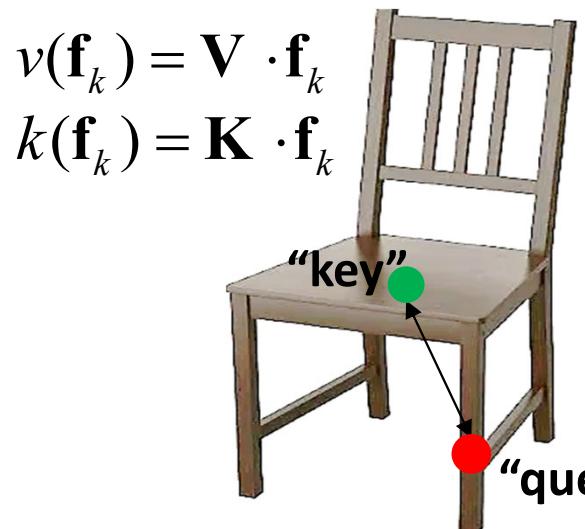


$$q(\mathbf{f}_q) = \mathbf{Q} \cdot \mathbf{f}_q$$

Attention score

How much a pixel/patch (query) is related to others (keys)?

Attention “score”: $\hat{a}(\mathbf{f}_q, \mathbf{f}_k) = \frac{\exp\{ a(\mathbf{f}_q, \mathbf{f}_k) \}}{\sum_{k'} \exp\{ a(\mathbf{f}_q, \mathbf{f}_{k'}) \}}$
(i.e., use softmax)



Problem: As the number of dimensions D of the input feature vector gets large, the variance of the dot product increased ...the input to softmax gets high values... the softmax is too peaked ...
=> tiny gradients

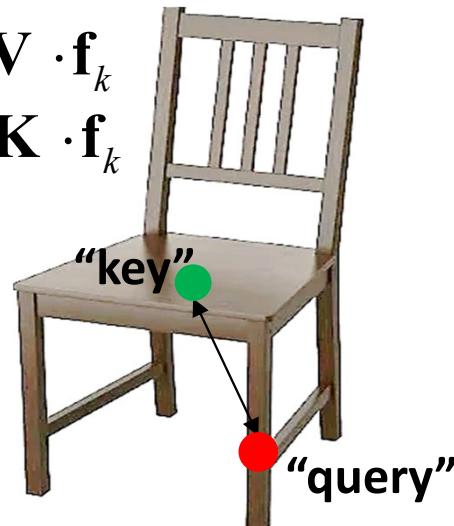
$$q(\mathbf{f}_q) = \mathbf{Q} \cdot \mathbf{f}_q$$

Scaled dot-product attention

How much a pixel/patch (query) is related to others (keys)?

Attention “score”: $\hat{a}(\mathbf{f}_q, \mathbf{f}_k) = \frac{\exp\{ a(\mathbf{f}_q, \mathbf{f}_k) / \sqrt{D} \}}{\sum_{k'} \exp\{ a(\mathbf{f}_q, \mathbf{f}_{k'}) / \sqrt{D} \}}$
(i.e., use softmax)

$$v(\mathbf{f}_k) = \mathbf{V} \cdot \mathbf{f}_k$$
$$k(\mathbf{f}_k) = \mathbf{K} \cdot \mathbf{f}_k$$



$$q(\mathbf{f}_q) = \mathbf{Q} \cdot \mathbf{f}_q$$

<https://ai.stackexchange.com/questions/41861/why-use-a-square-root-in-the-scaled-dot-product>

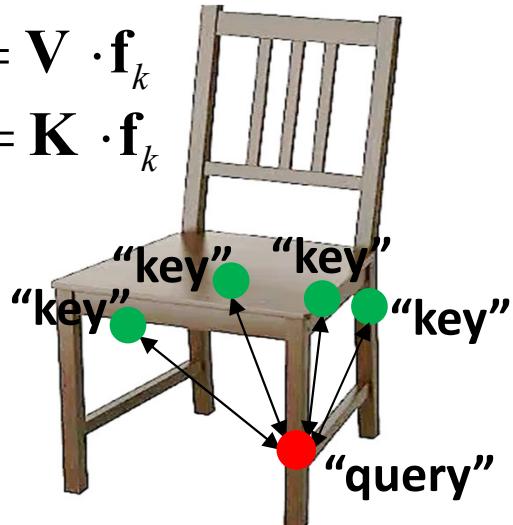
Vaswani, Attention Is All You Need, 2017

New feature representations

Compute a new encoding for query pixel/patch as a **weighted sum of values of key points**:

$$\text{New features : } \mathbf{f}_q' = \sum_k \hat{a}(\mathbf{f}_q, \mathbf{f}_k) v(\mathbf{f}_k)$$

$$v(\mathbf{f}_k) = \mathbf{V} \cdot \mathbf{f}_k$$
$$k(\mathbf{f}_k) = \mathbf{K} \cdot \mathbf{f}_k$$



... these can be added back to the input feature vector as a residual and further processed by an MLP (fully connected layer)

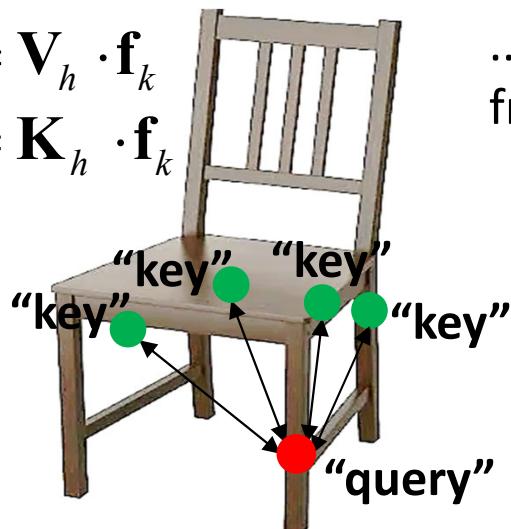
$$q(\mathbf{f}_q) = \mathbf{Q} \cdot \mathbf{f}_q$$

Multi-head Attention

Learn different query, key, value transformations for each attention “head”.

$$\text{Multi-head attention : } \mathbf{f}_{q,h}' = \sum_k \hat{a}_h(\mathbf{f}_q, \mathbf{f}_k) v_h(\mathbf{f}_k)$$

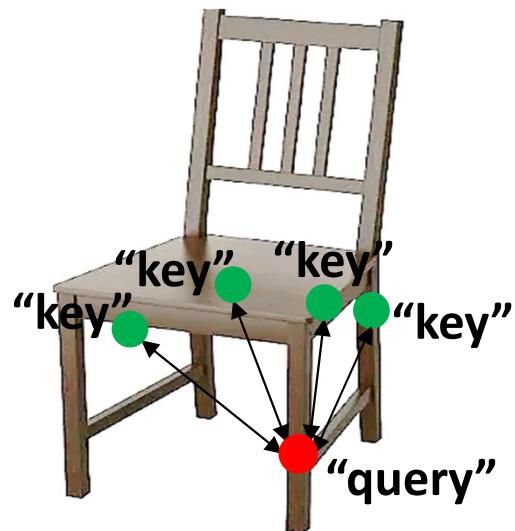
$v(\mathbf{f}_k) = \mathbf{V}_h \cdot \mathbf{f}_k$
 $k(\mathbf{f}_k) = \mathbf{K}_h \cdot \mathbf{f}_k$
... concatenate the resulting features
from all heads, and process them with the MLP



$$q(\mathbf{f}_q) = \mathbf{Q}_h \cdot \mathbf{f}_q$$

Quadratic complexity!

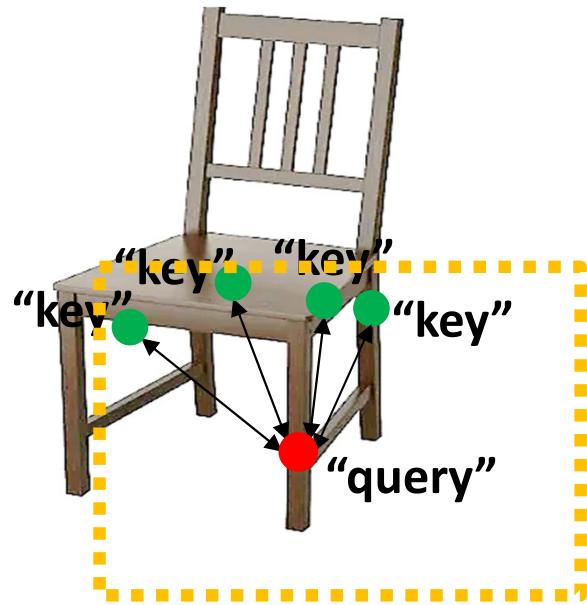
Comparing each query (N pixels) with every key (N pixels)
yields **quadratic complexity**!



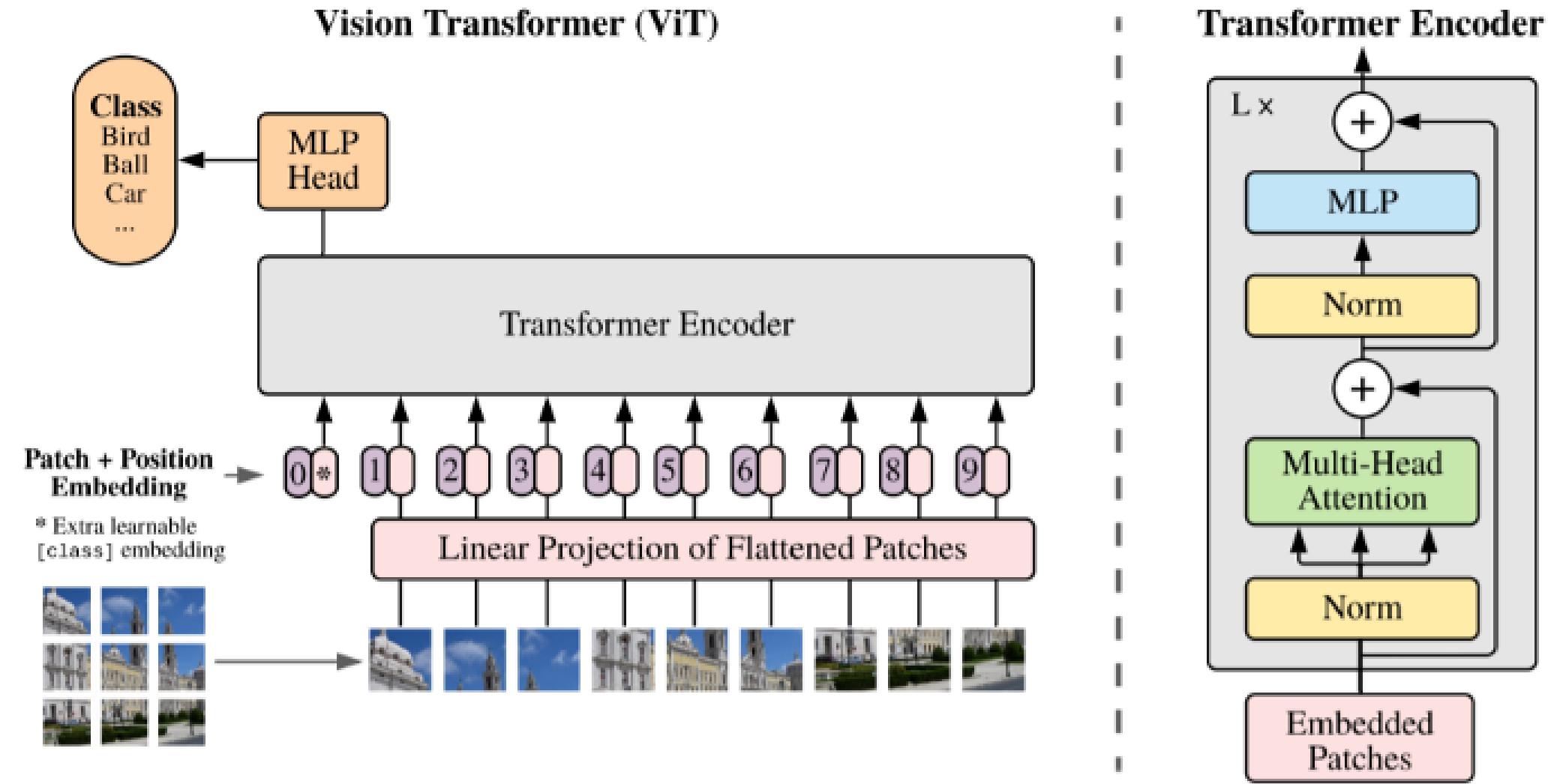
Quadratic complexity!

Solutions:

- (a) Compute attention between patches instead of pixels
(common in 2D vision transformers)
- (b) Limit attention within regions e.g., 3D Euclidean balls
(common in point transformers)



Vision Transformers



"An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale", 2021

Transformers as Autoregressive Models

Explicitly models data distribution by assuming that
our data consists of individual elements

$$X = \{x_1, x_2, x_3, x_4 \dots\}$$

e.g., an image consists of a (flattened) series of pixels, or a mesh consists of a series of triangles, or a sentence is made out of tokens ...

Transformers as Autoregressive Models

Explicitly models data distribution by assuming that
our data consists of individual elements

$$\mathbf{X} = \{\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3, \mathbf{x}_4 \dots\}$$

e.g., an image consists of a (flattened) series of pixels, or a mesh consists of a series of triangles, or a sentence is made out of tokens ...

Data distribution is modeled as:

$$P(\mathbf{X}) = P(\mathbf{x}_1) \cdot P(\mathbf{x}_2 | \mathbf{x}_1) \cdot P(\mathbf{x}_3 | \mathbf{x}_1, \mathbf{x}_2) \cdot P(\mathbf{x}_4 | \mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3) \dots$$

Transformers as Autoregressive Models

Explicitly models data distribution by assuming that
our data consists of individual elements

$$\mathbf{X} = \{\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3, \mathbf{x}_4 \dots\}$$

e.g., an image consists of a (flattened) series of pixels, or a mesh consists of a series of triangles, or a sentence is made out of tokens ...

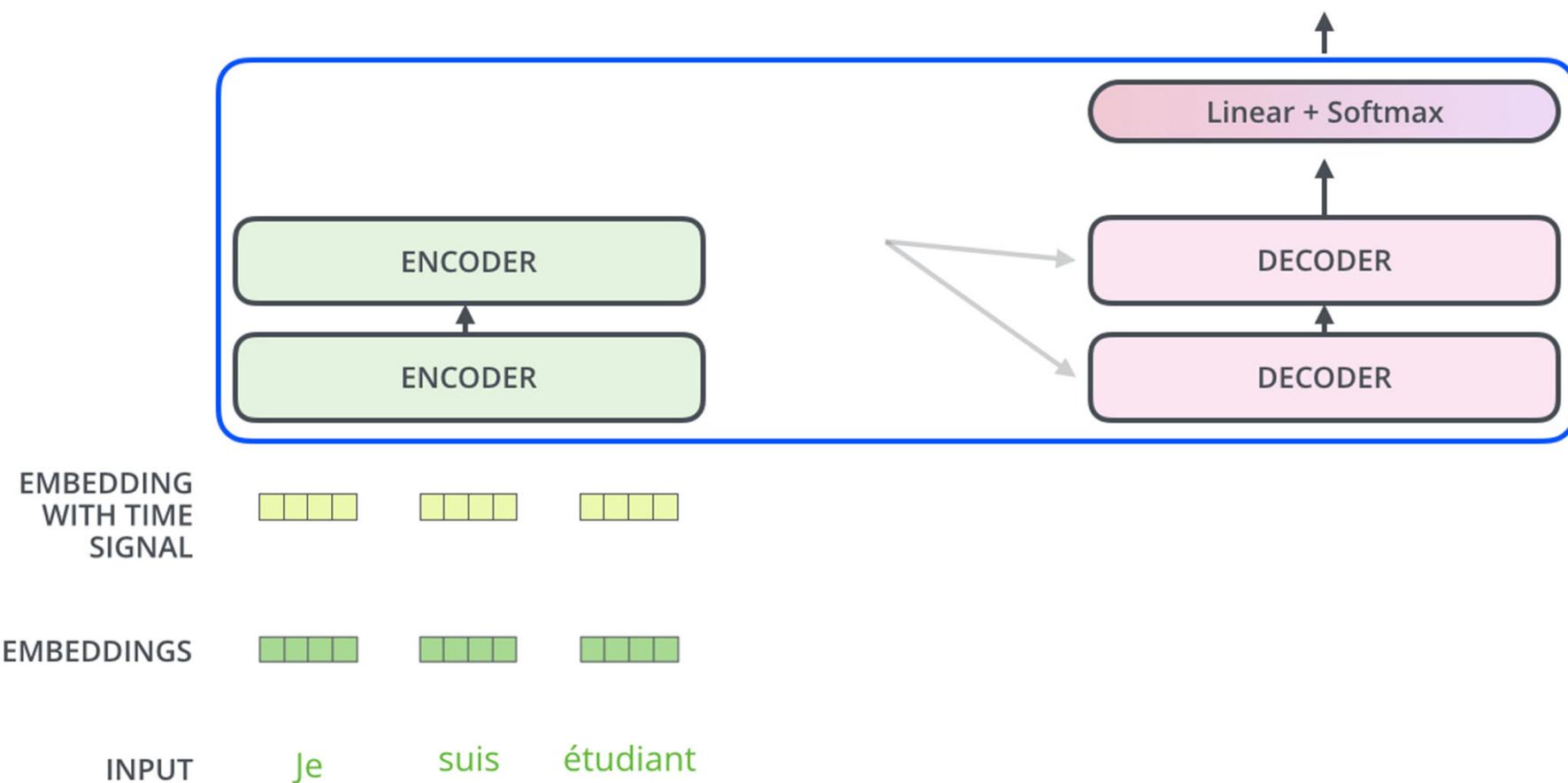
Data distribution is modeled as:

$$P(\mathbf{X}) = \prod_{t=0}^T P(\mathbf{x}_{t+1} | \mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3, \dots, \mathbf{x}_t)$$

Transformers as Autoregressive Models

Decoding time step: 1 2 3 4 5 6

OUTPUT



Generative models (to be discussed)

- **Autoregressive models [done]**
- Variational Autoencoders [TODO this/next week]
- Diffusion Models [TODO this/next week]
- Generative Adversarial Networks [lower priority, after Easter]

Announcements

- **Lecture to replenish lost hours**

Fri Apr 11, 3-5pm (will also discuss midterm as well)
Room: [145Π58]

- Assignment 1 to be released next week
- List of papers to choose for presentation