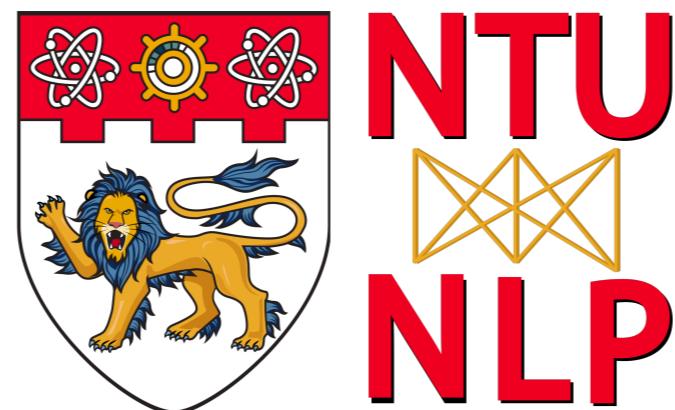


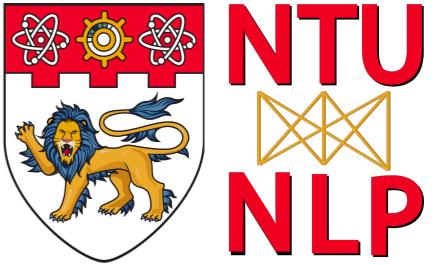
# Deep Learning for Natural Language Processing

Shafiq Joty



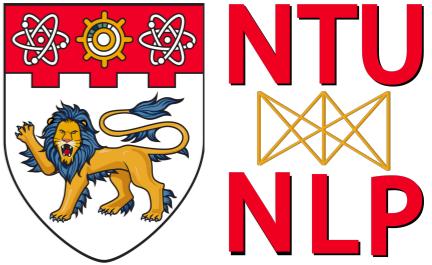
Lecture 9: Transformers

# Recurrent Neural Networks



- Used to be the **default model of choice** for learning variable-length representations.
- Natural fit for sentences (or a sequence of objects)
- LSTMs and GRUs are the dominant recurrent models.

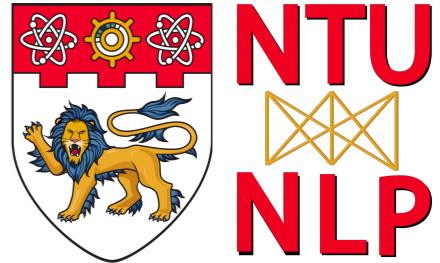
# Recurrent Neural Networks



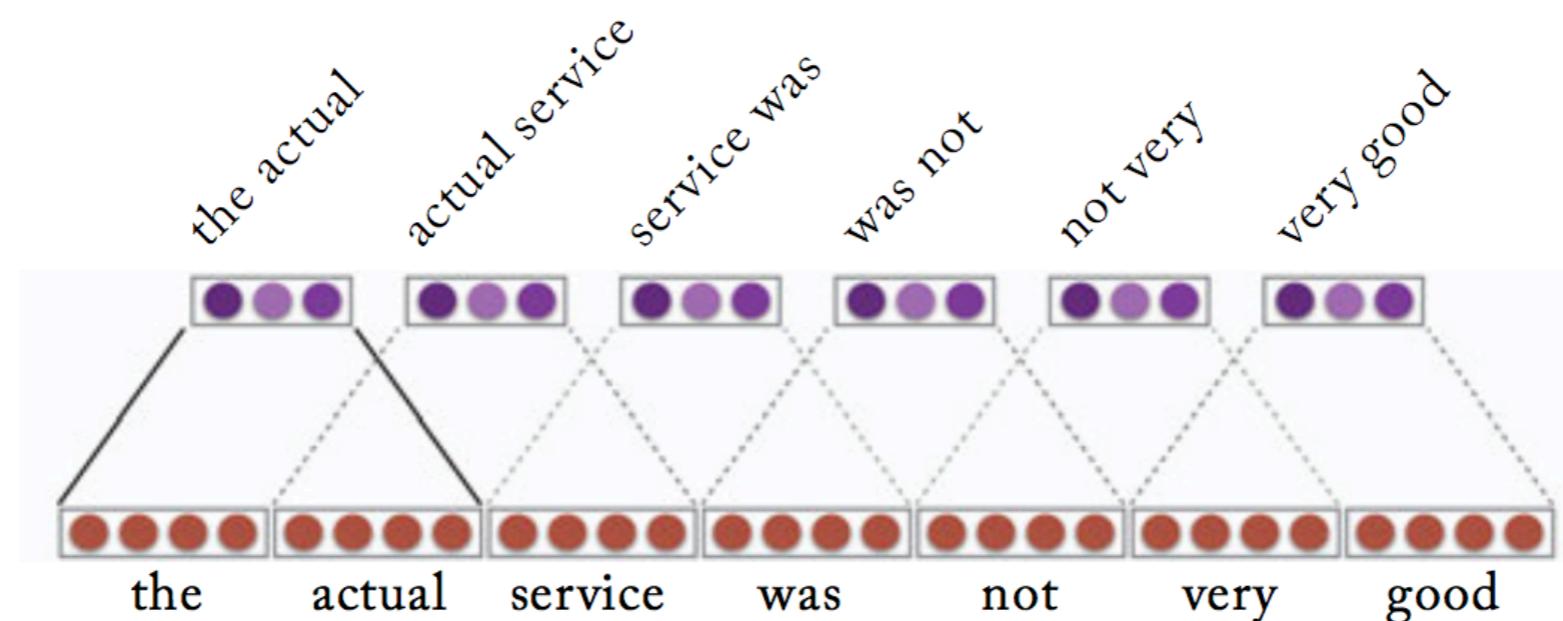
## Issues with RNN Computations

- But sequential computation inhibits parallelization.
- No explicit modeling of long and short range dependencies
- Not many hierarchy of features.
- RNNs are a bit wasteful!

# Convolutional Neural Networks?



- Per-layer Parallelizable
- Exploits local dependencies
- 'Interaction distance' between positions linear or logarithmic
- Modeling long-distance dependencies require many layers

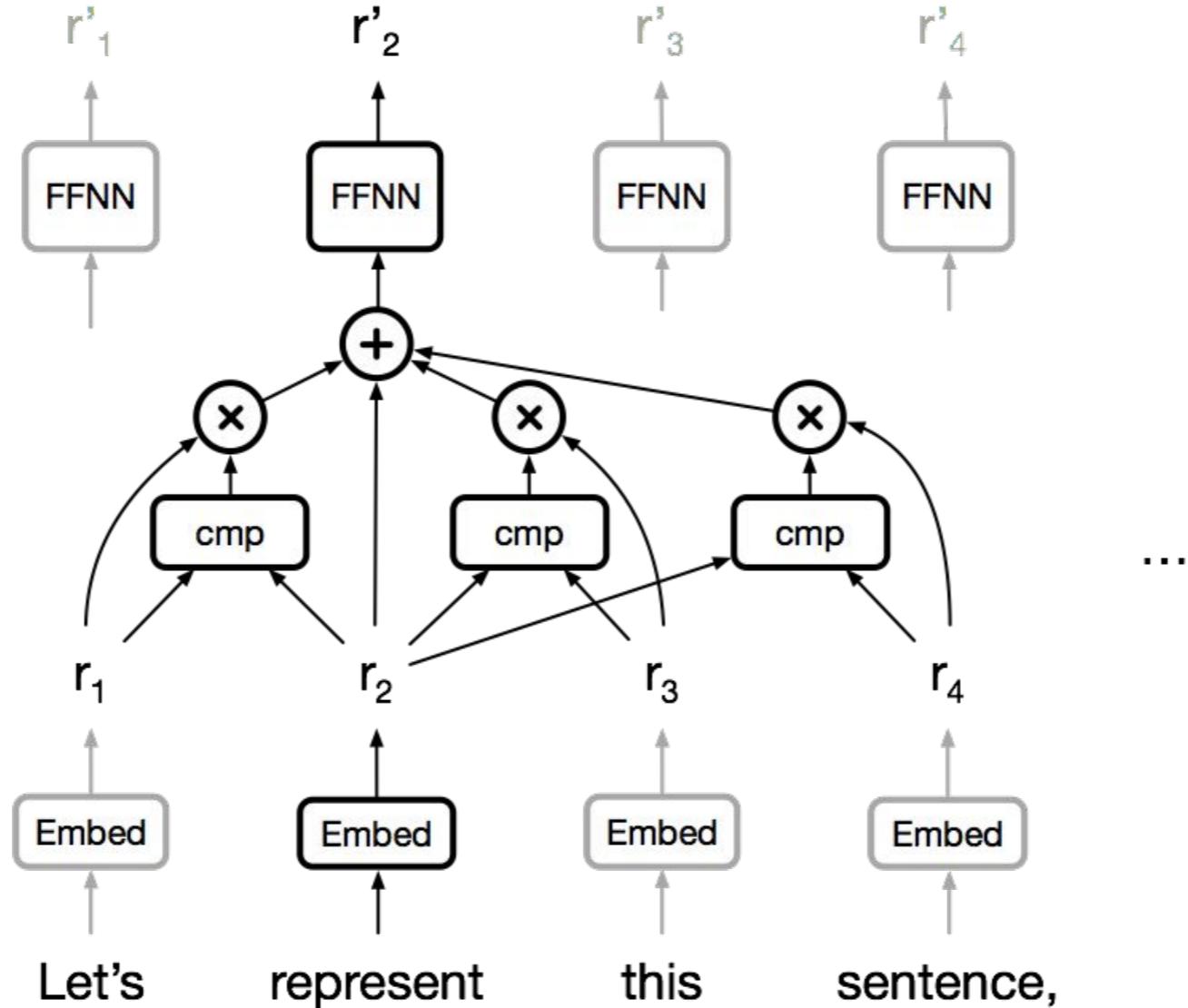


# Attention?

- Attention between encoder and decoder is crucial in NMT

Why not use attention for representations?

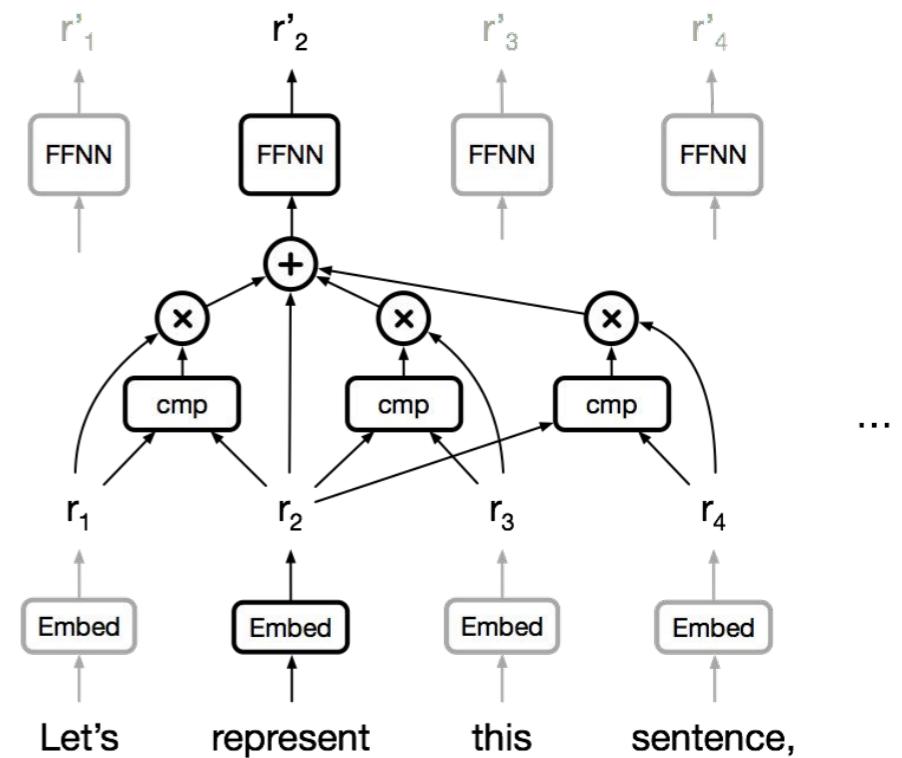
# Self Attentions



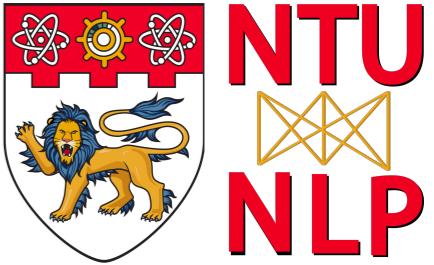
Let's represent this sentence,

# Self Attentions

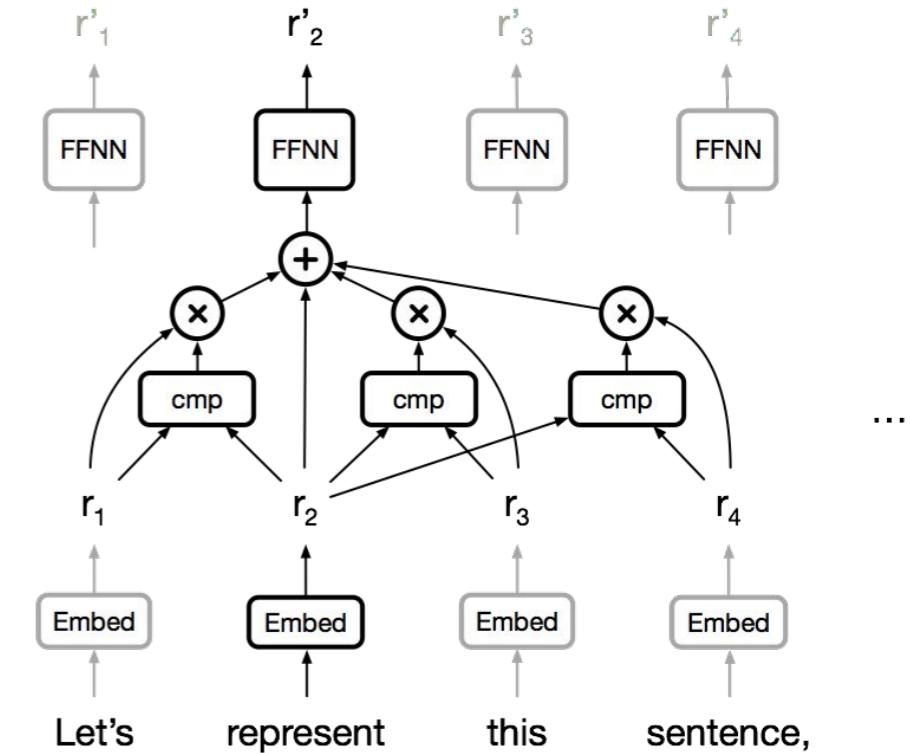
- Constant 'path length' between any two positions.
- (Soft) Gating & multiplicative interactions.
- Easy to parallelize (per layer).
- Can replace sequential computation entirely



# Attention is Cheap



FLOPs	
Self-Attention	$O(\text{length}^2 \cdot \text{dim})$
RNN (LSTM)	$O(\text{length} \cdot \text{dim}^2)$
Convolution	$O(\text{length} \cdot \text{dim}^2 \cdot \text{kernel\_width})$

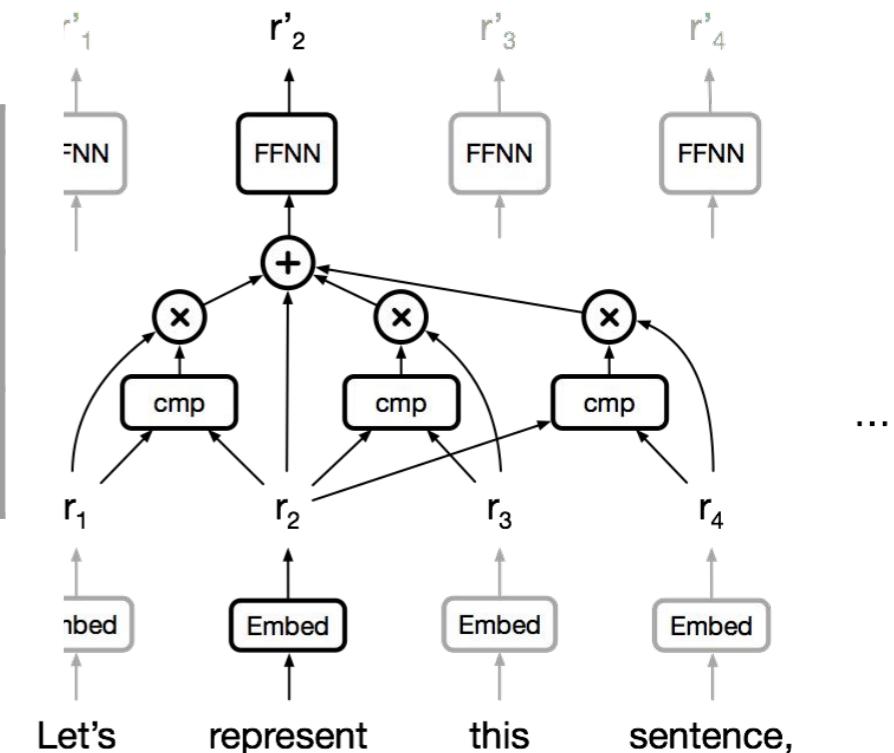


source: Vaswani et al.

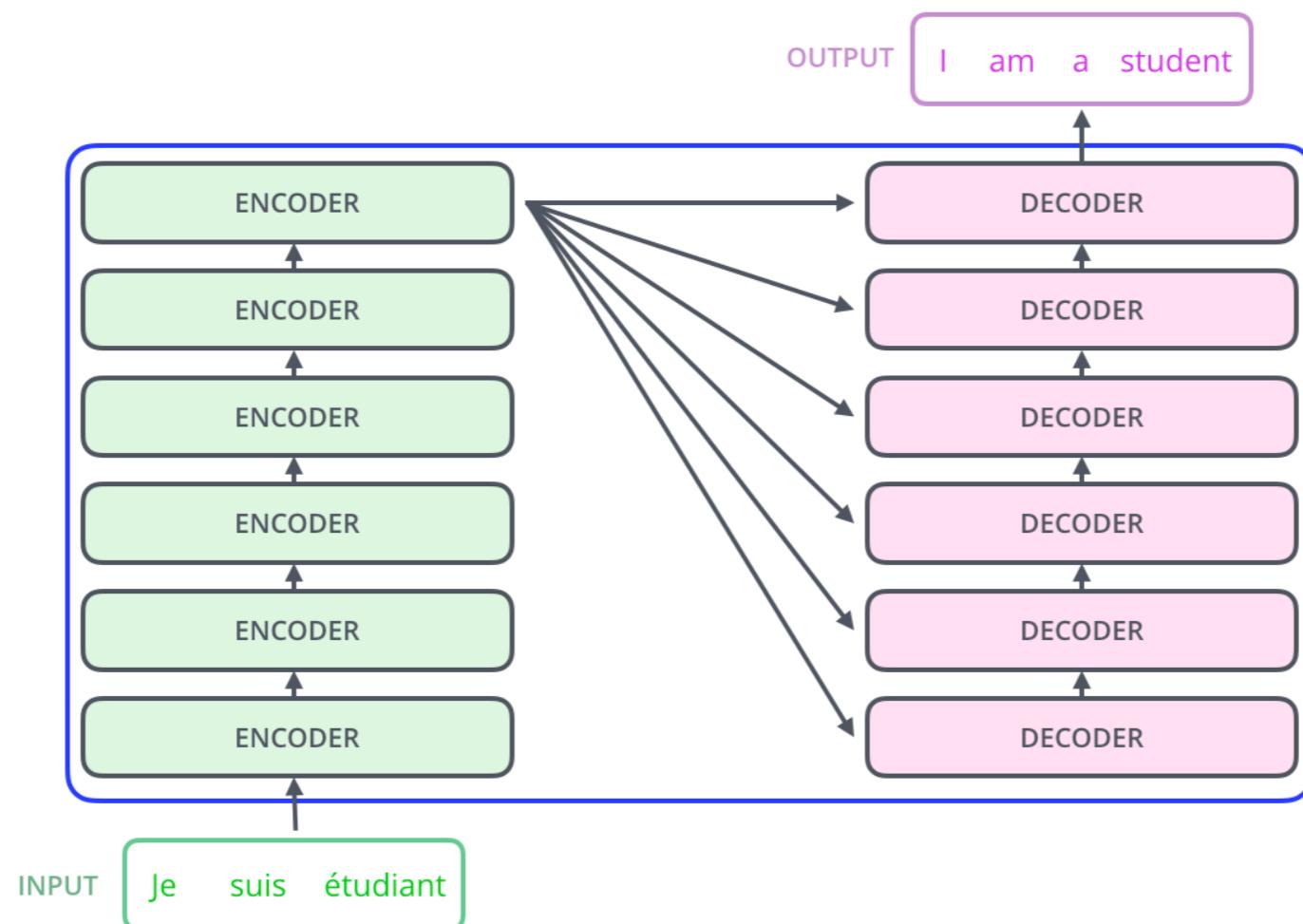
# Attention is Cheap

FLOPs	
Self-Attention	$O(\text{length}^2 \cdot \text{dim})$ = $4 \cdot 10^9$
RNN (LSTM)	$O(\text{length} \cdot \text{dim}^2)$ = $16 \cdot 10^9$
Convolution	$O(\text{length} \cdot \text{dim}^2 \cdot \text{kernel\_width})$ = $6 \cdot 10^9$

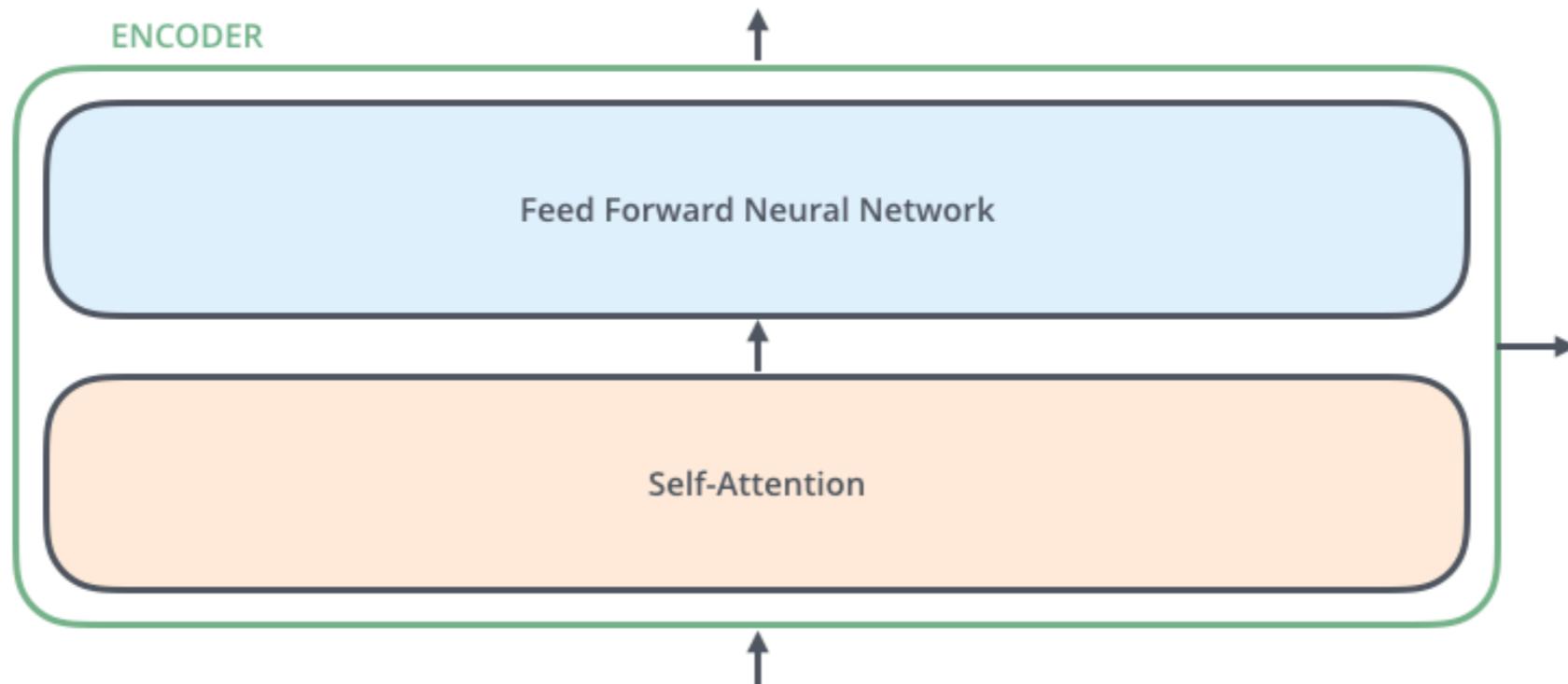
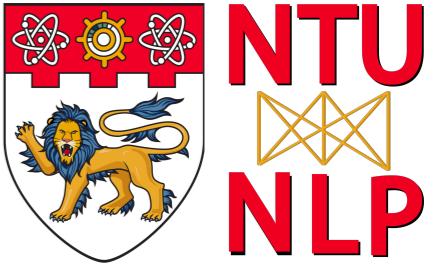
length=1000 dim=1000 kernel\_width=3



# Encoder-Decoder Network

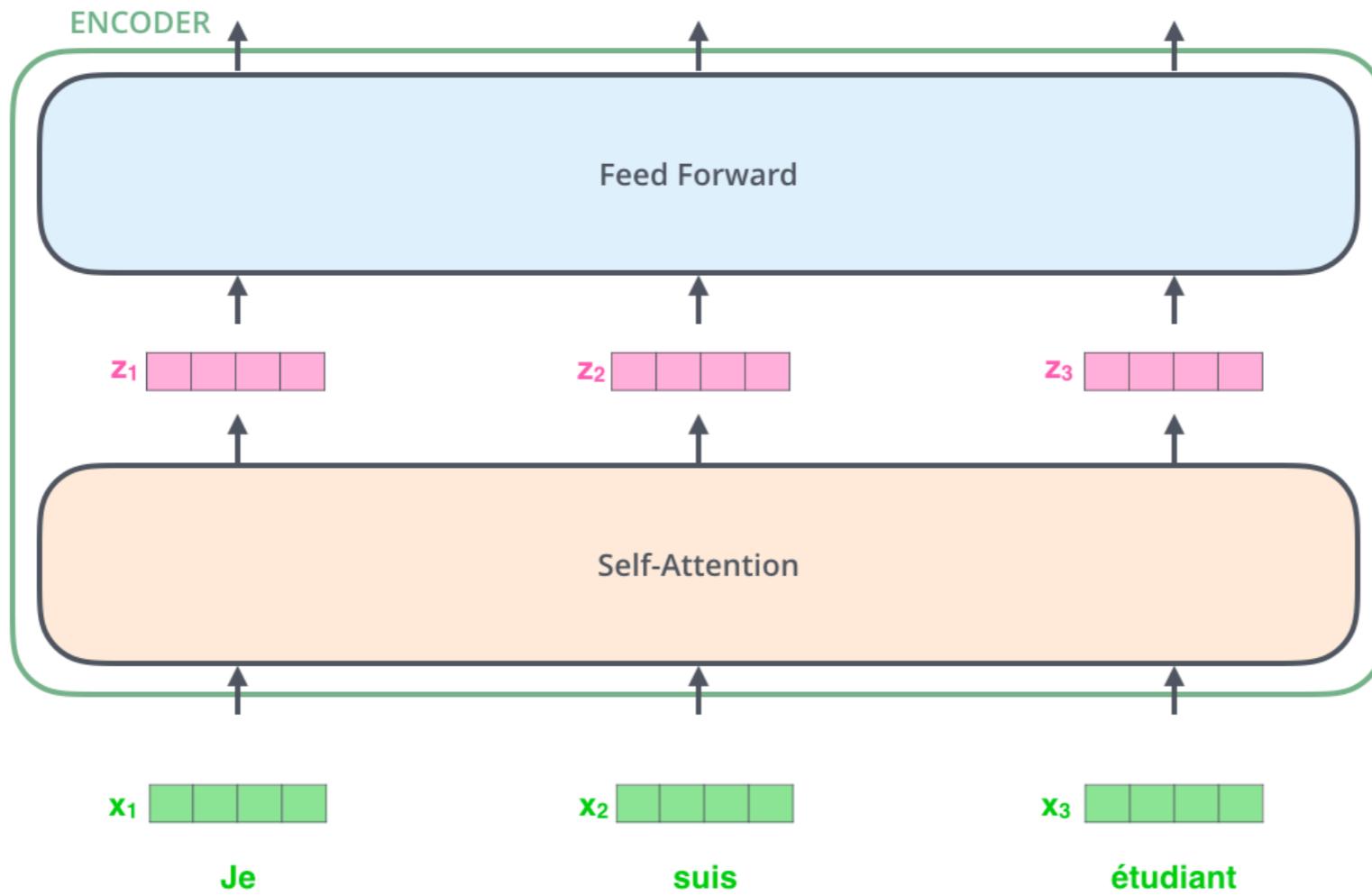


# One Encoder Layer



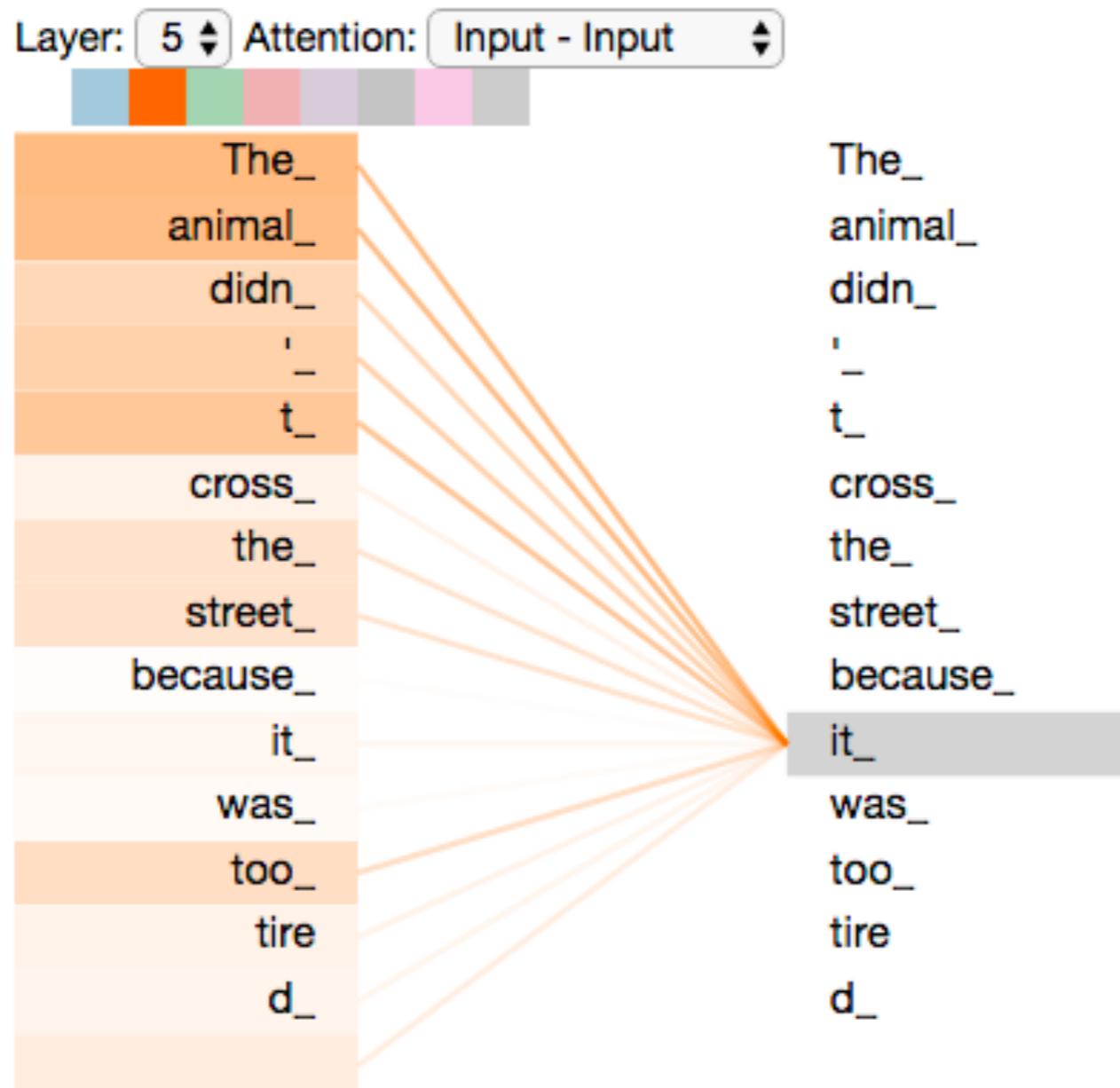
- The encoder's inputs first flow through a self-attention layer
- The outputs of the self-attention layer are fed to a feed-forward neural network. The exact same feed-forward network is independently applied to each position.

# The First Encoder Layer



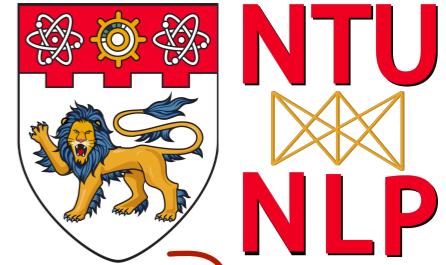
- The dependencies between word vectors are modelled by the self-attention layer.

# One Self-Attention Layer



- The dependencies between word vectors are modelled by the self-attention layer.

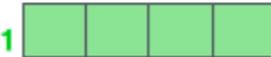
# One Self-Attention Layer



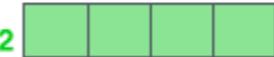
Input

Thinking

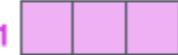
Embedding

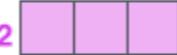
$x_1$  

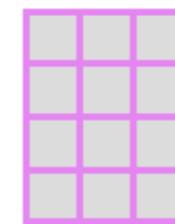
Machines

$x_2$  

Queries

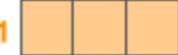
$q_1$  

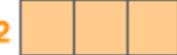
$q_2$  



$WQ$

Keys

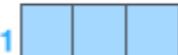
$k_1$  

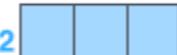
$k_2$  



$WK$

Values

$v_1$  

$v_2$  



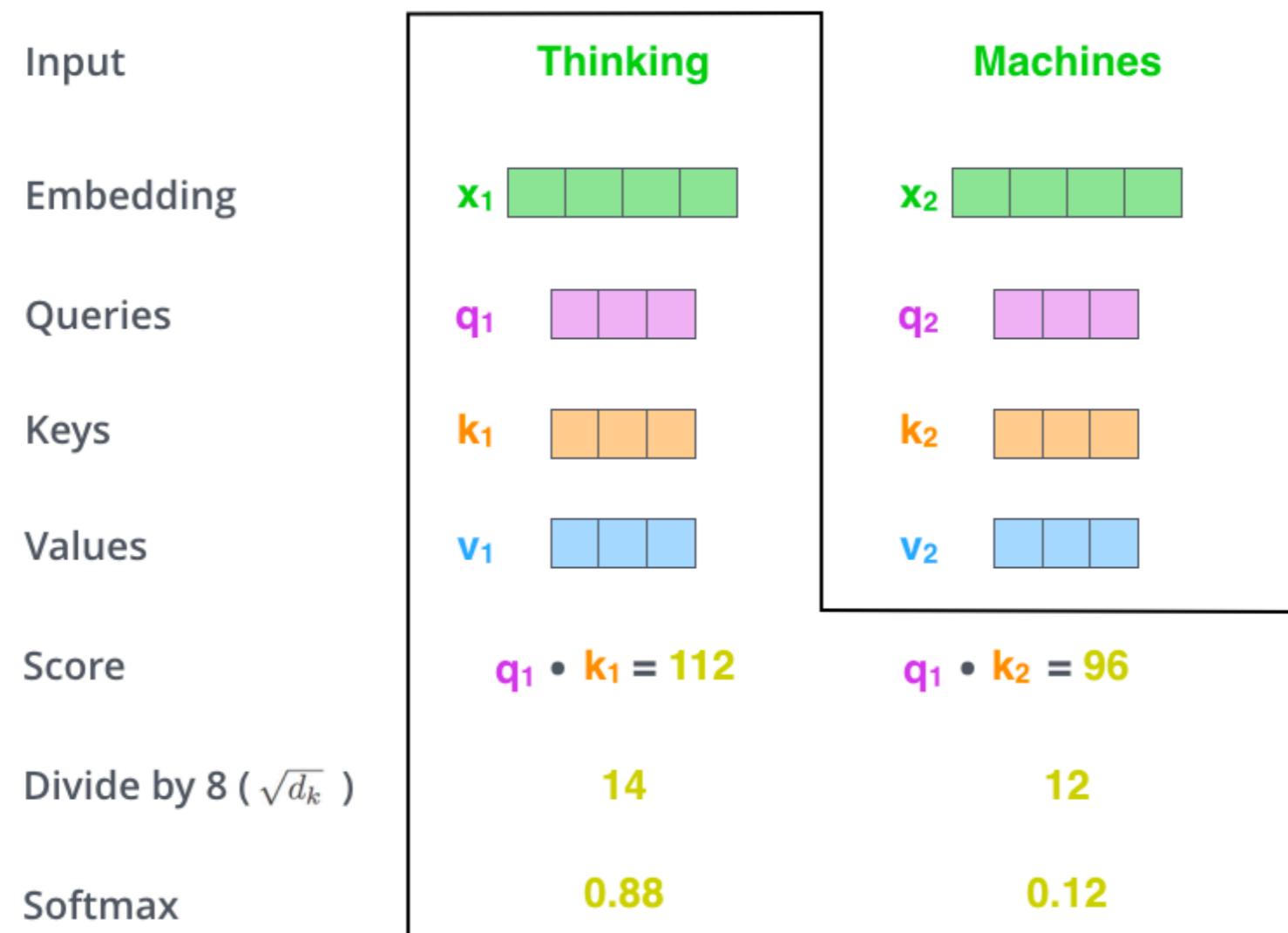
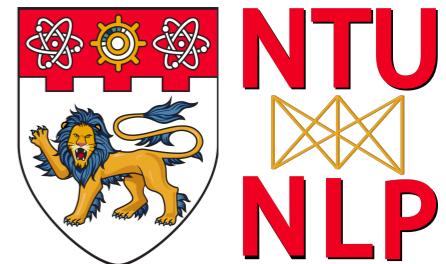
$WV$

$$\text{Soft}\left[\frac{(QW^Q)(K^W K)^T}{\sqrt{d}}(VW V)\right]$$

- The dependencies between word vectors are modelled by the self-attention layer.

source: <https://jalammar.github.io/>

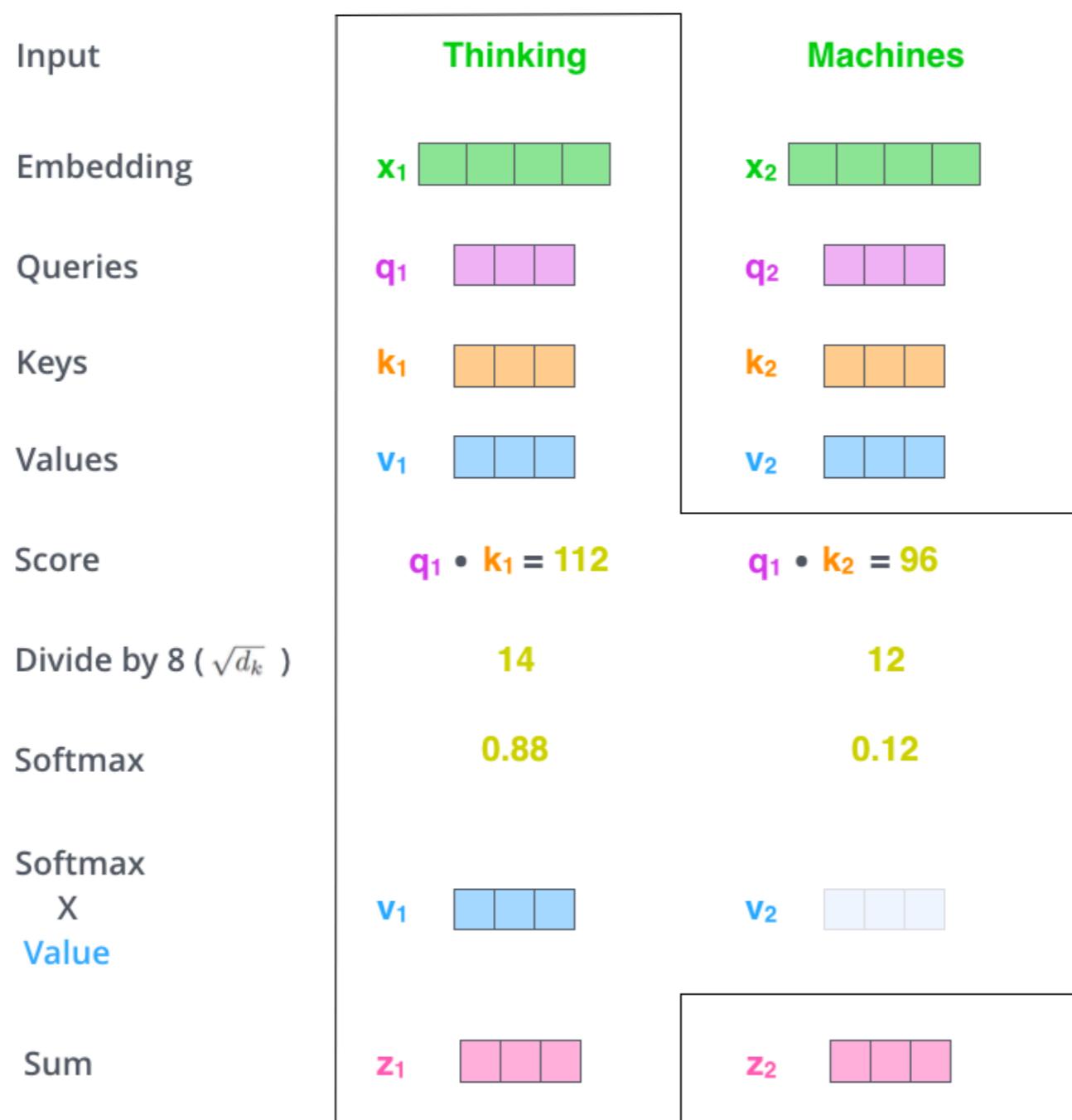
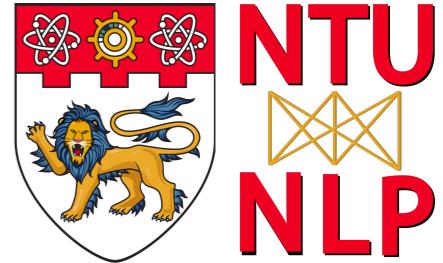
# One Self-Attention Layer



- The dependencies between word vectors are modelled by the self-attention layer.

source: <https://jalammar.github.io/>

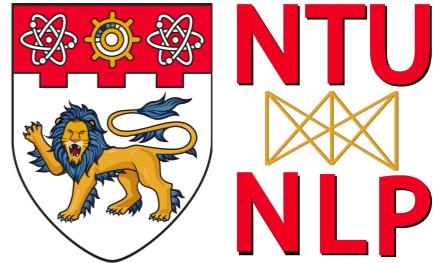
# One Self-Attention Layer



- The dependencies between word vectors are modelled by the self-attention layer.

source: <https://jalammar.github.io/>

# One Self-Attention Layer



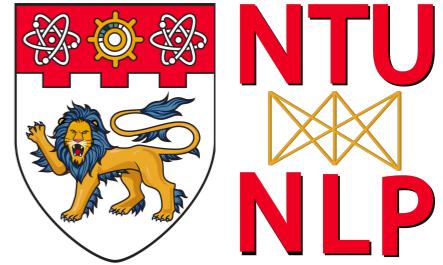
$$\begin{array}{ccc} \text{x/Q'} & & \\ \begin{matrix} \text{---} \\ \text{---} \\ \text{---} \end{matrix} & \times & \begin{matrix} \text{W}^Q \\ \text{---} \\ \text{---} \end{matrix} & = & \begin{matrix} \text{Q} \\ \text{---} \\ \text{---} \end{matrix} \\ \\ \text{x/K'} & & \\ \begin{matrix} \text{---} \\ \text{---} \\ \text{---} \end{matrix} & \times & \begin{matrix} \text{W}^K \\ \text{---} \\ \text{---} \end{matrix} & = & \begin{matrix} \text{K} \\ \text{---} \\ \text{---} \end{matrix} \\ \\ \text{x/V'} & & \\ \begin{matrix} \text{---} \\ \text{---} \\ \text{---} \end{matrix} & \times & \begin{matrix} \text{W}^V \\ \text{---} \\ \text{---} \end{matrix} & = & \begin{matrix} \text{V} \\ \text{---} \\ \text{---} \end{matrix} \end{array}$$

The diagram illustrates the computation of three attention matrices (Q, K, V) from an input matrix (x). Each row of the input matrix is multiplied by a weight matrix (W<sup>Q</sup>, W<sup>K</sup>, or W<sup>V</sup>) to produce the corresponding Q, K, and V matrices. A red line with an arrow points from the label "x/Q'" to the first row of the input matrix.

- The dependencies between word vectors are modelled by the self-attention layer.

source: <https://jalammar.github.io/>

# One Self-Attention Layer



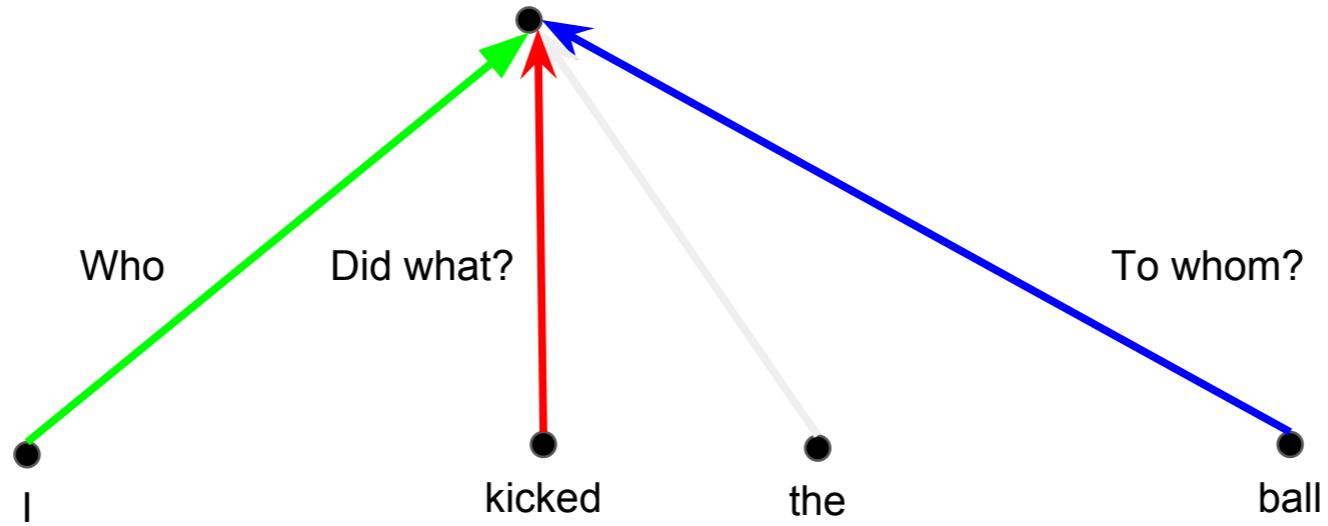
$$\text{softmax} \left( \frac{\begin{matrix} Q \\ \times \\ K^T \end{matrix}}{\sqrt{d_k}} \right) V = Z$$

The diagram illustrates the computation of a self-attention layer. It shows three input matrices:  $Q$  (purple, 3x3),  $K^T$  (orange, 3x3), and  $V$  (blue, 3x3). The  $Q$  matrix is multiplied by the transpose of the  $K$  matrix ( $K^T$ ). The result is then divided by the square root of the dimension  $d_k$ . Finally, this scaled matrix is multiplied by the  $V$  matrix to produce the output  $Z$ .

- The dependencies between word vectors are modelled by the self-attention layer.

# Multi-head Attention

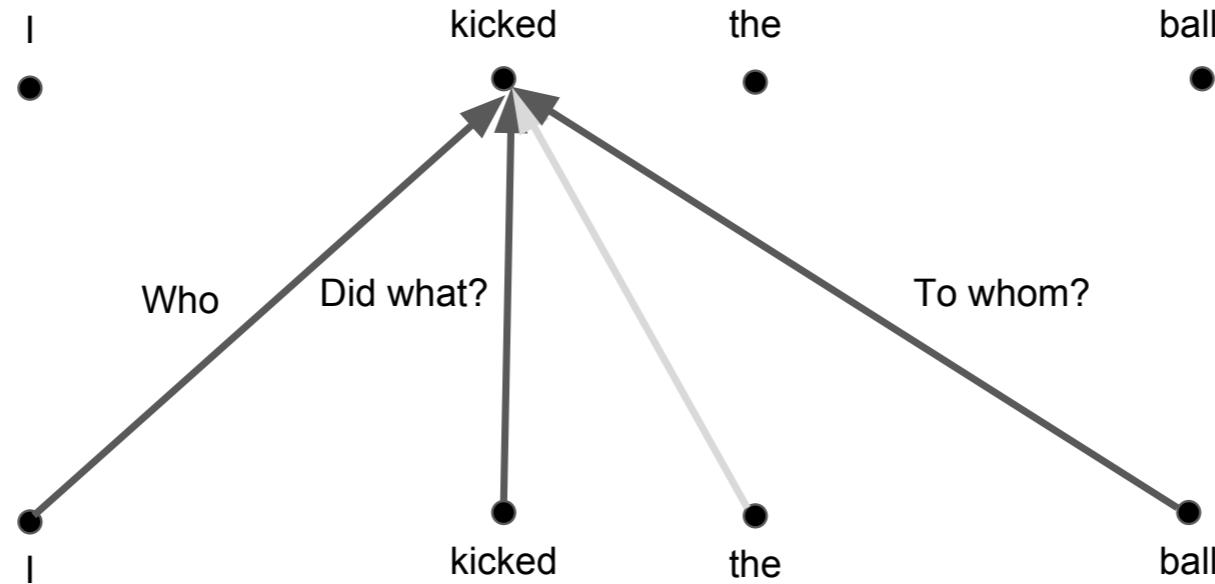
## • Convolution:



- Different kernels extract different features

# Multi-head Attention

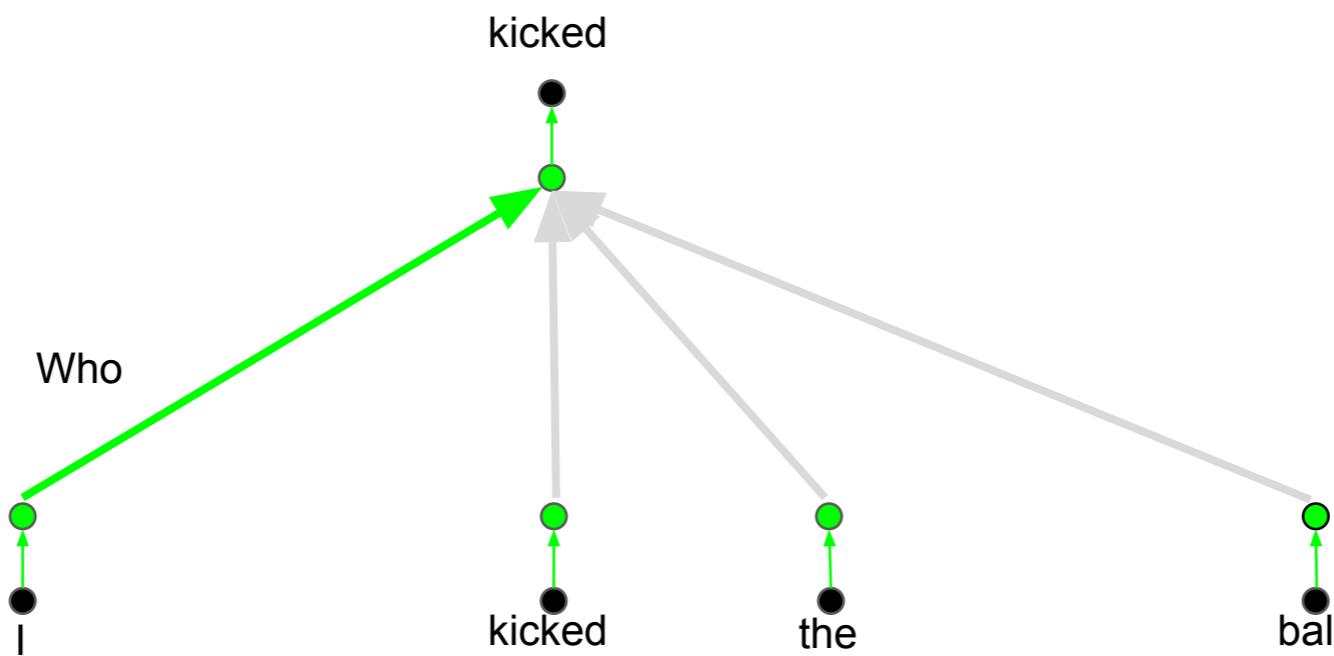
## • Self-attention:



$$\text{Soft}\left[\frac{(QW^Q)(KW^K)^T}{\sqrt{d}}\right]VW$$

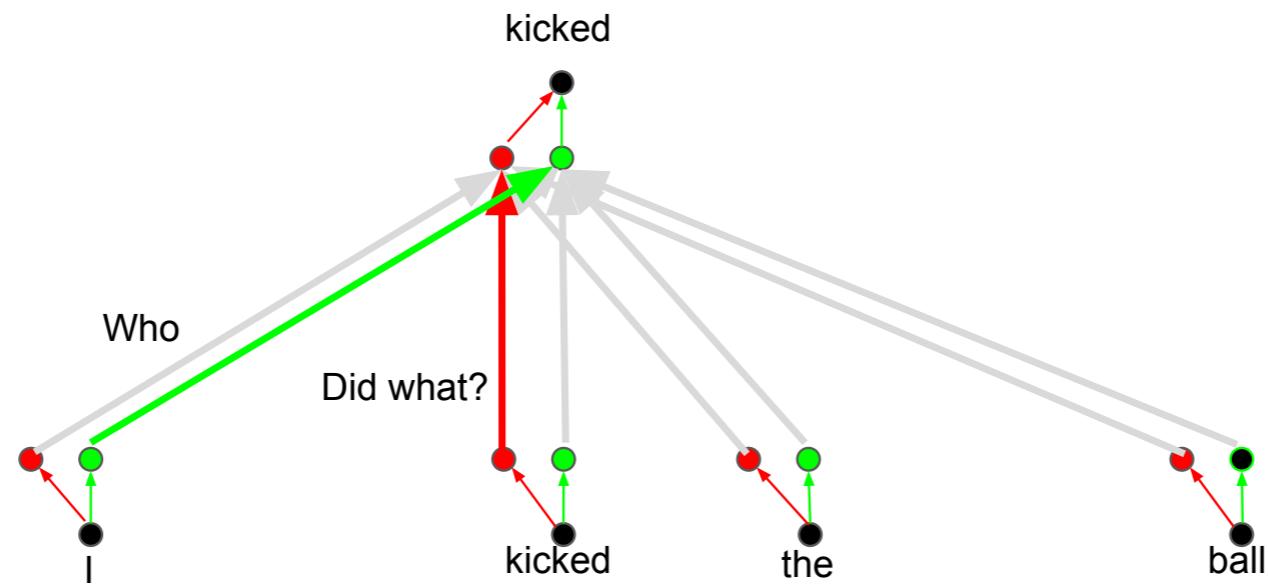
- All these are mixed in a single vector

# Self-attention – Who head



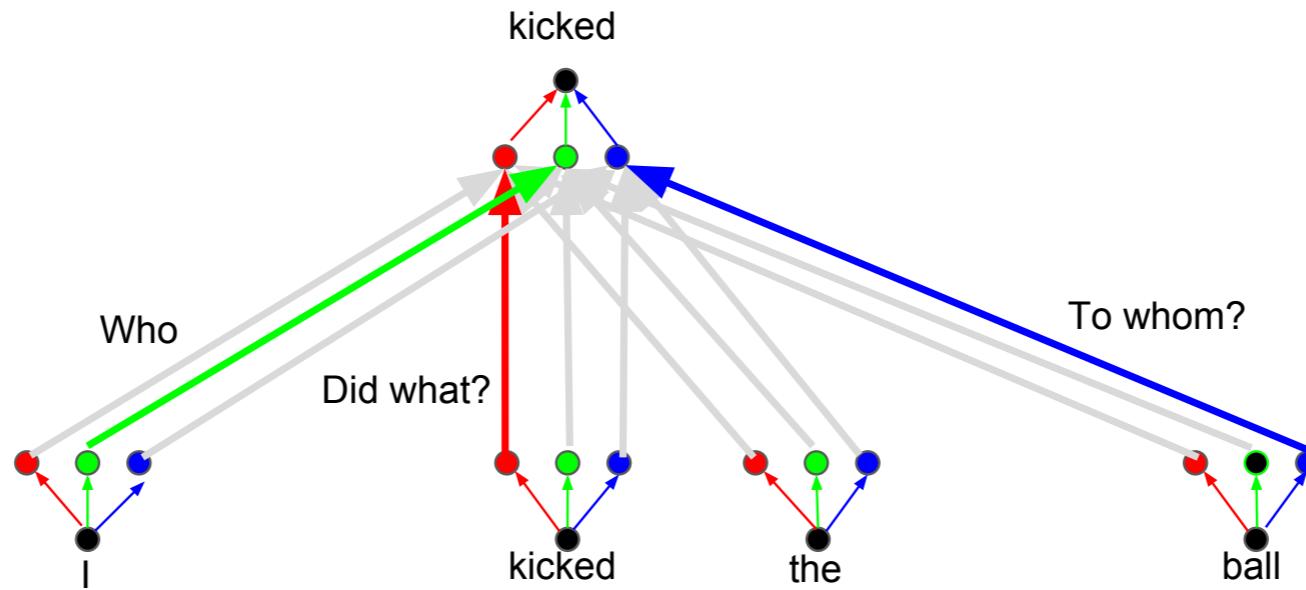
- Who head (kernel) extracts subj information

# Self-attention – Did what head



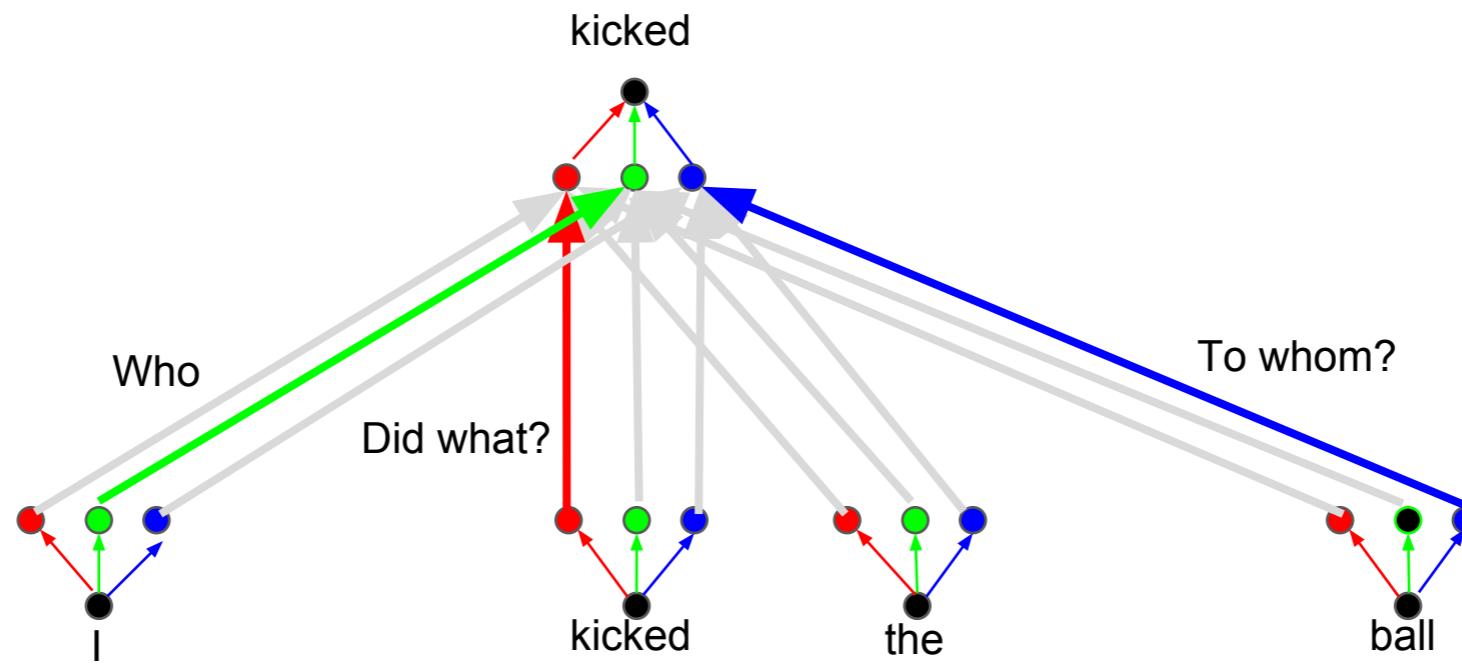
- Did What head (kernel) extracts verb information

# Self-attention – To whom head



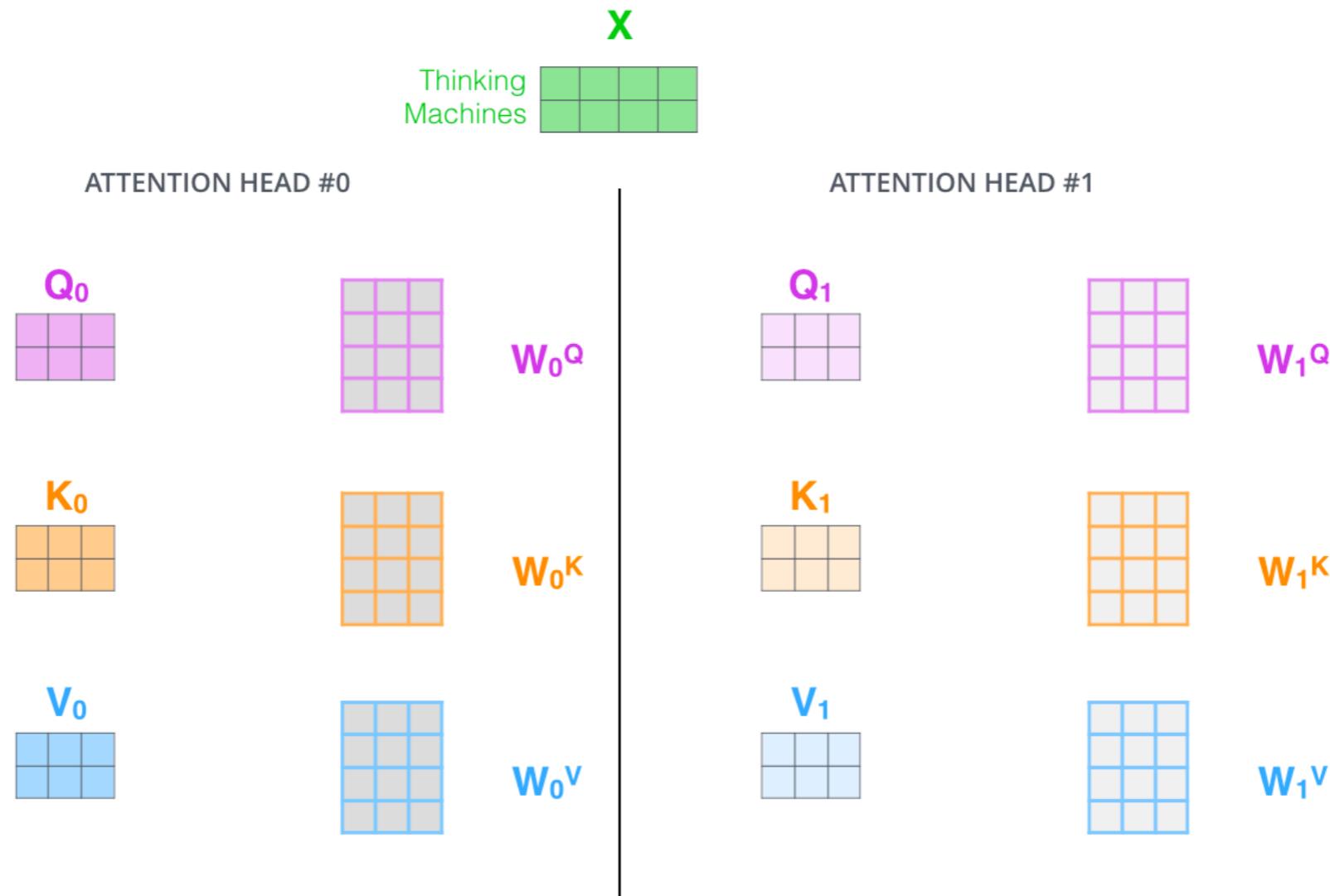
- To Whom head (kernel) extracts object information

# Self-attention – Multihead



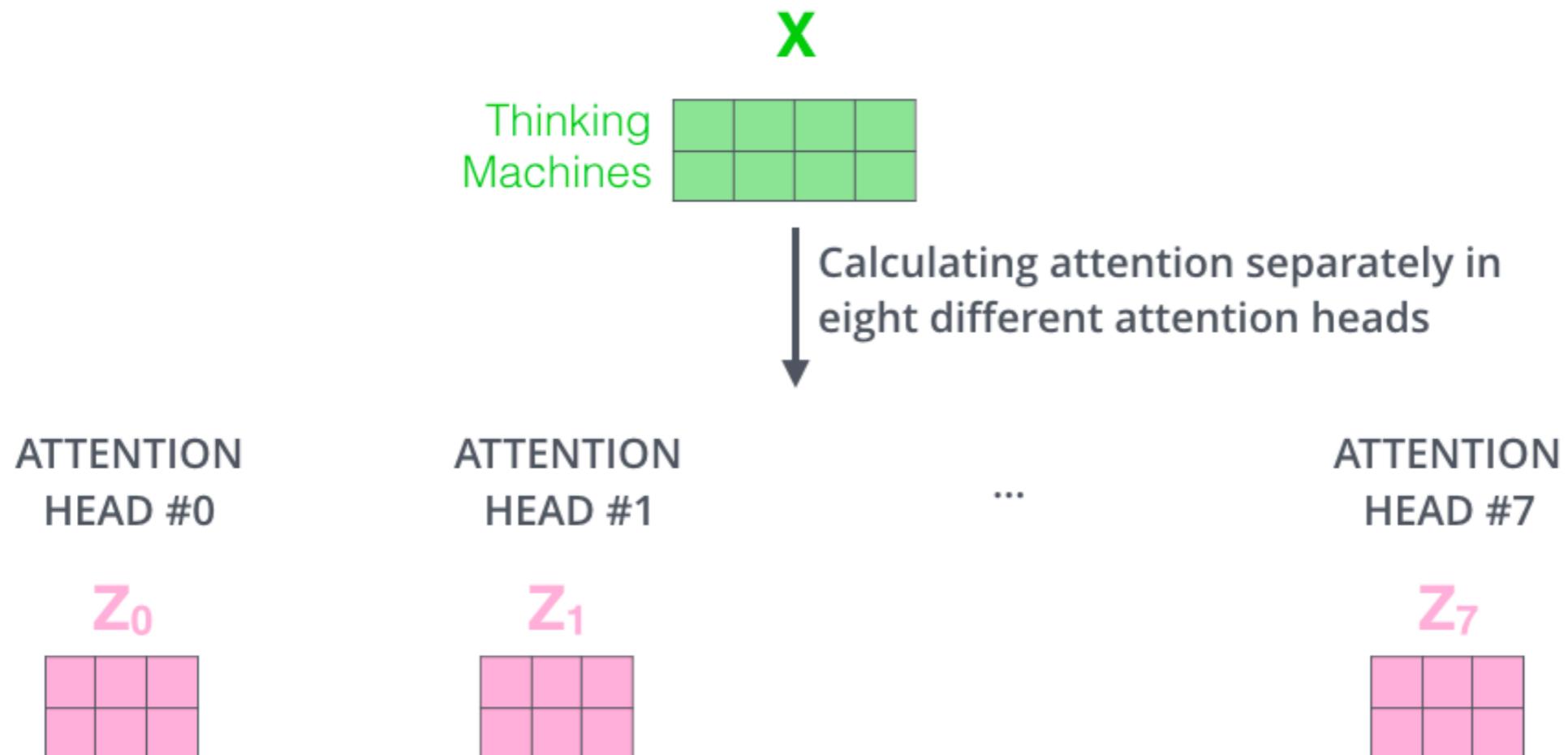
- Multi-head extracts different aspects

# Using Multiple Heads



- Multiple heads expand the model's ability to focus on different positions.
- Each head gives the attention layer multiple representation subspaces

# Using Multiple Heads



- Transformer uses 8 different heads

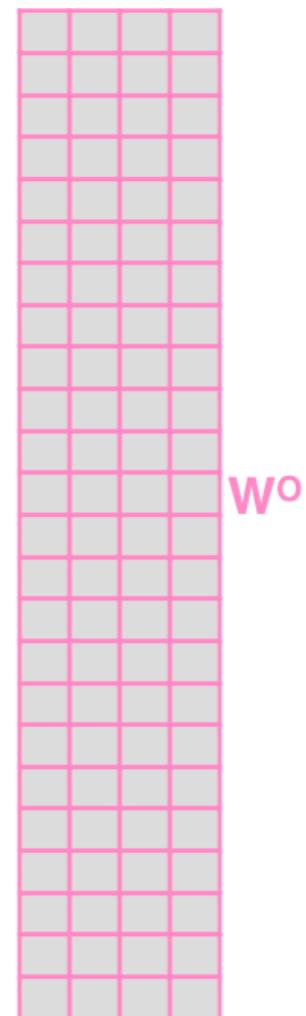
# Using Multiple Heads

1) Concatenate all the attention heads



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

$X$



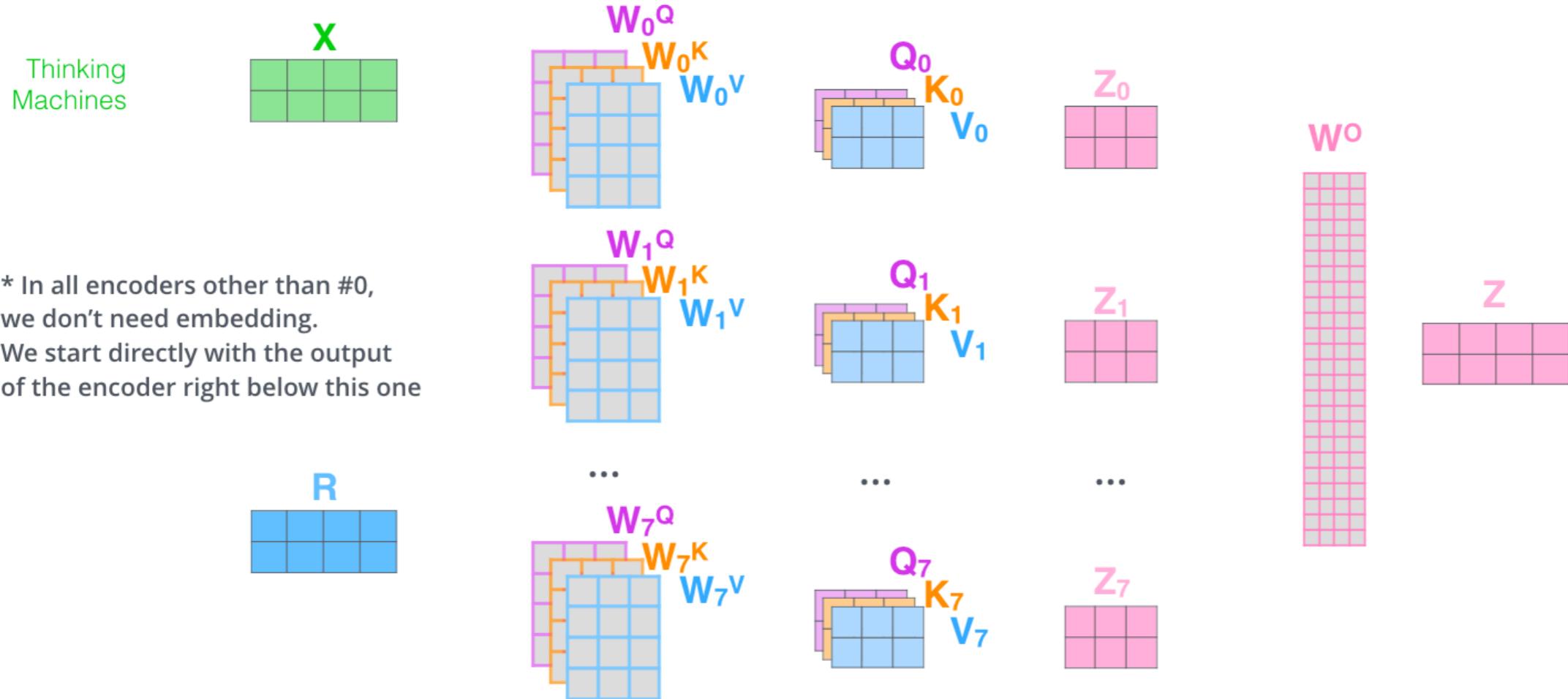
3) The result would be the  $Z$  matrix that captures information from all the attention heads. We can send this forward to the FFNN

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

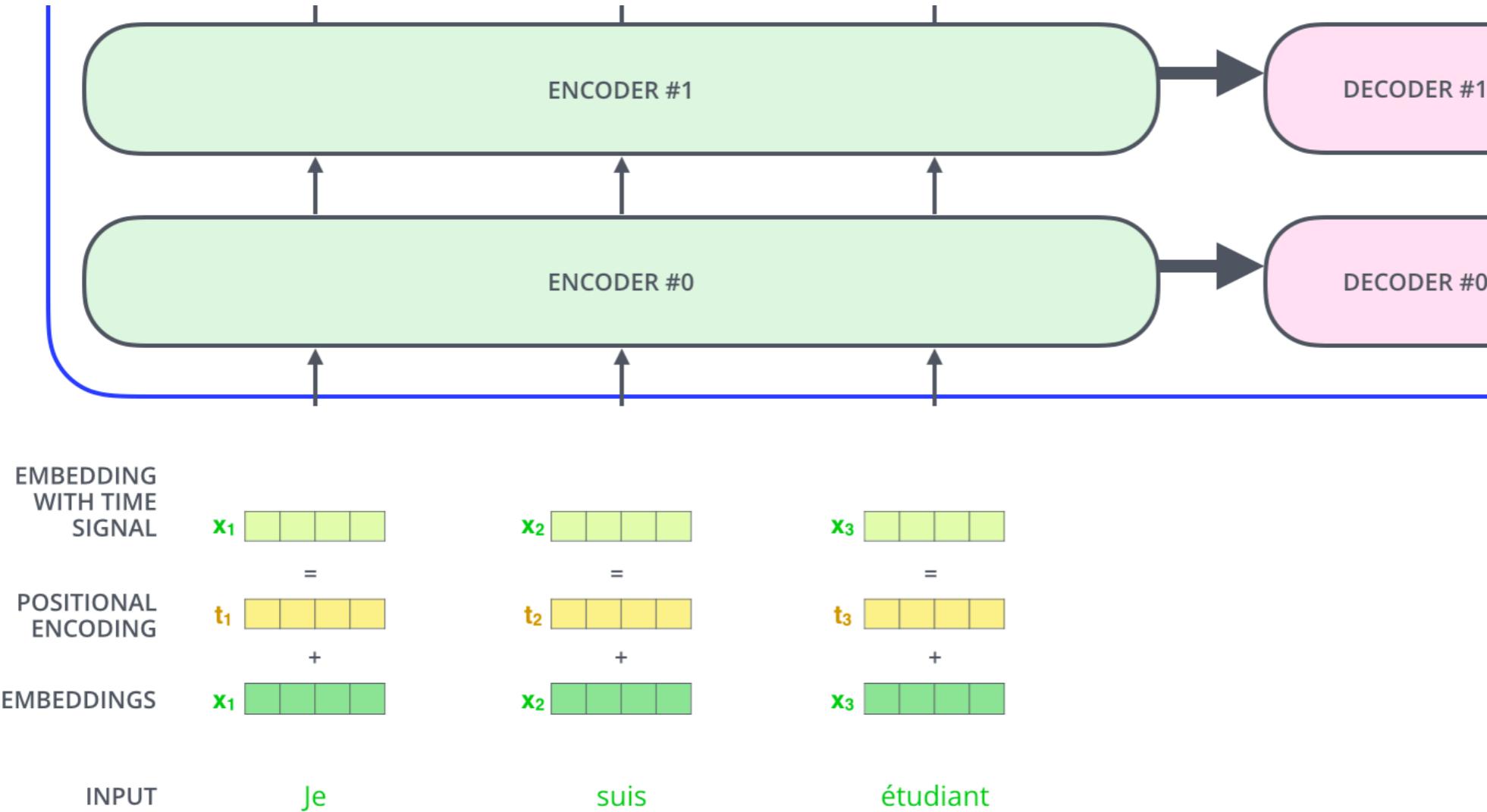
- Concatenate the 8-head's outputs and pass it the FFN

# Using Multiple Heads

- 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



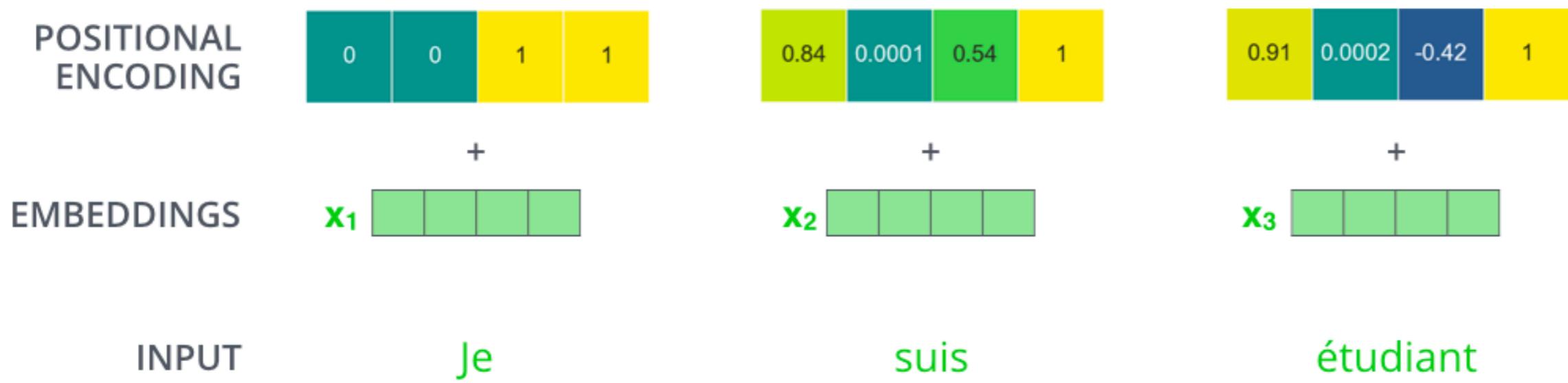
# Injecting Order into Embeddings



To give the model a sense of the order of the words, we add positional encoding vectors -- the values of which follow a specific pattern (sine/cosine).

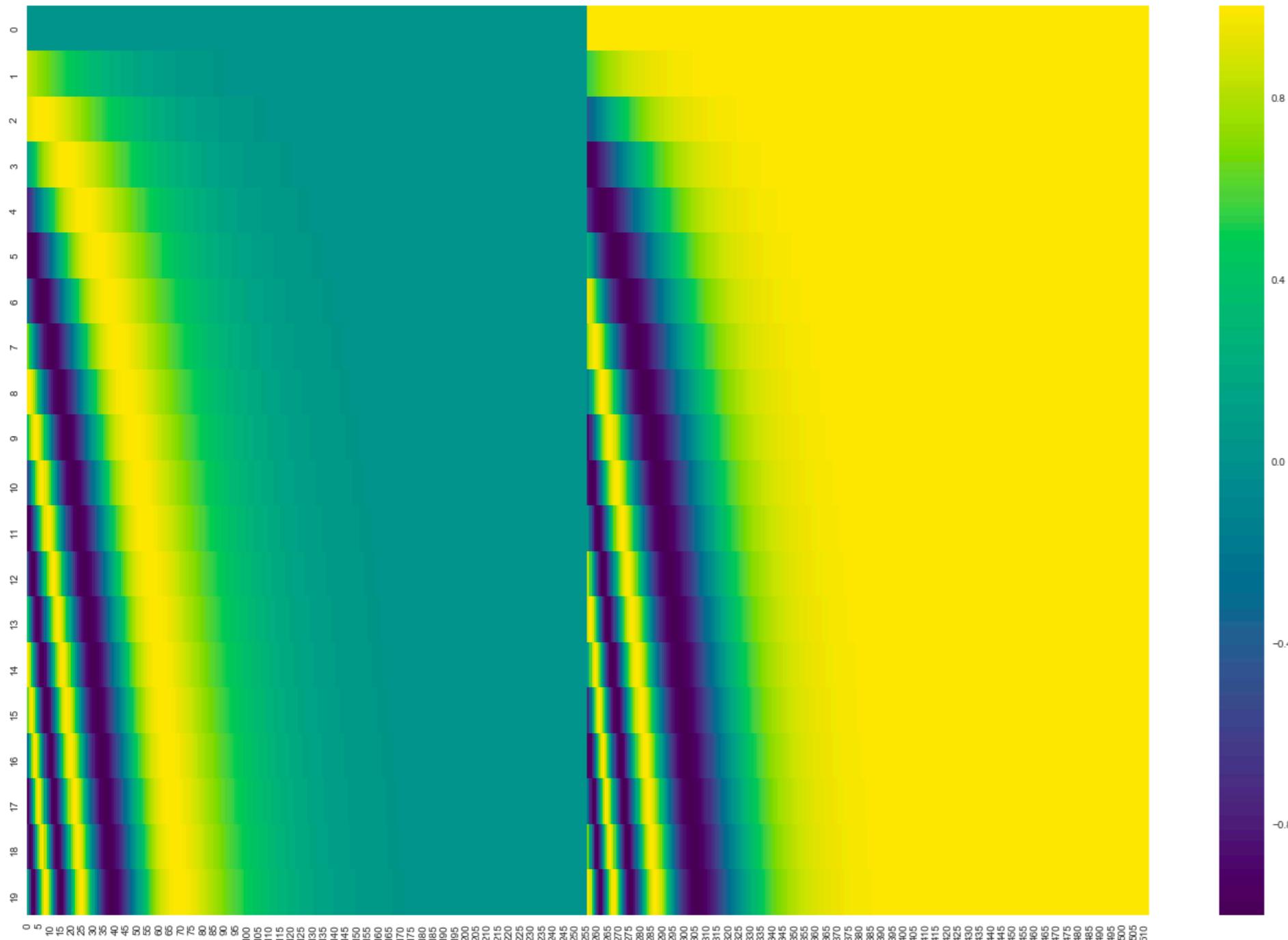
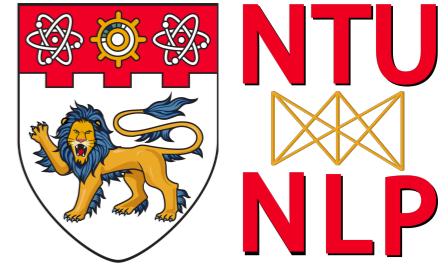
# Injecting Order into Embeddings

- Actual word representations are byte-pair encodings
  - As seen in last lecture



Example of positional encoding with a toy embedding size of 4

# Injecting Order into Embeddings

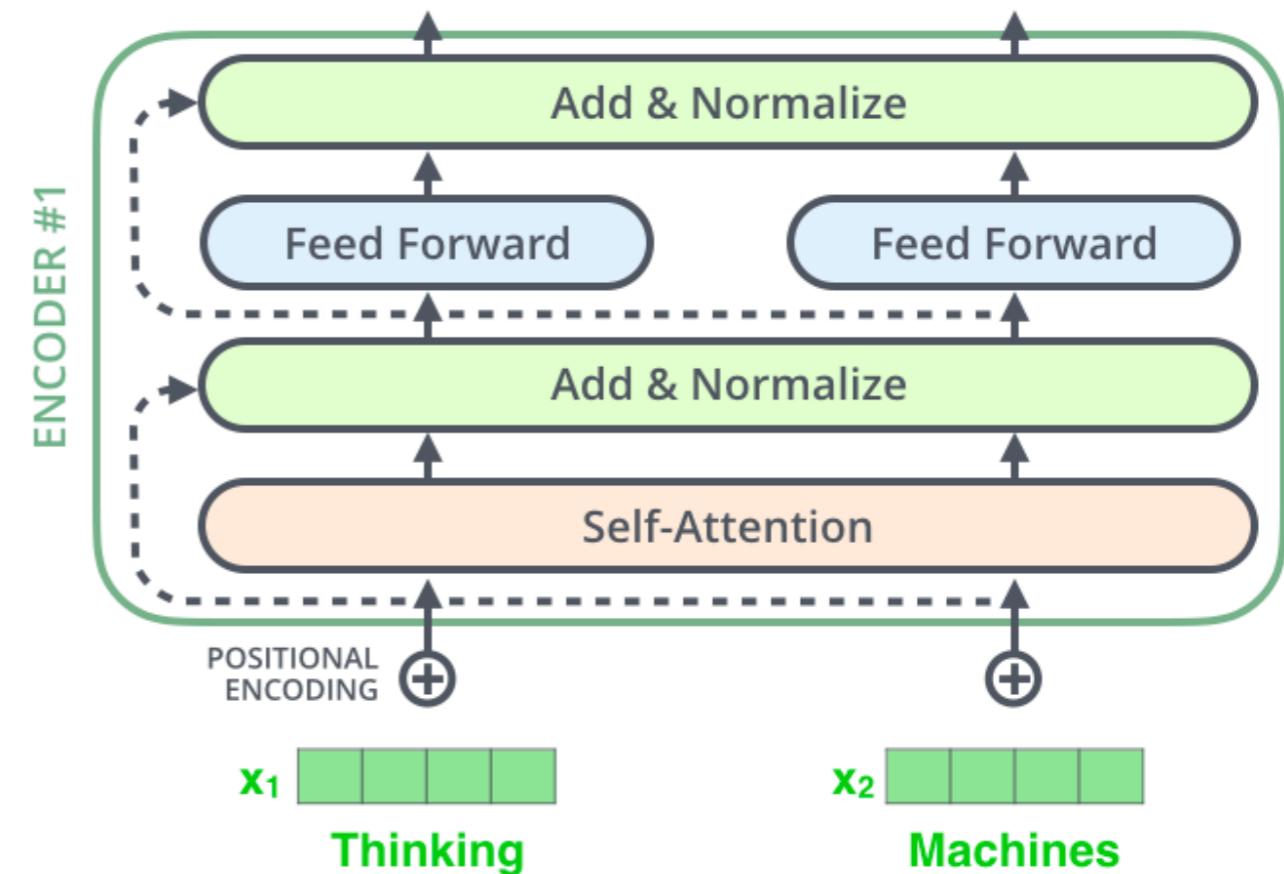


real example of positional encoding for 20 words (rows) with an embedding size of 512 (columns)

source: <https://jalammar.github.io/>

# The Residuals

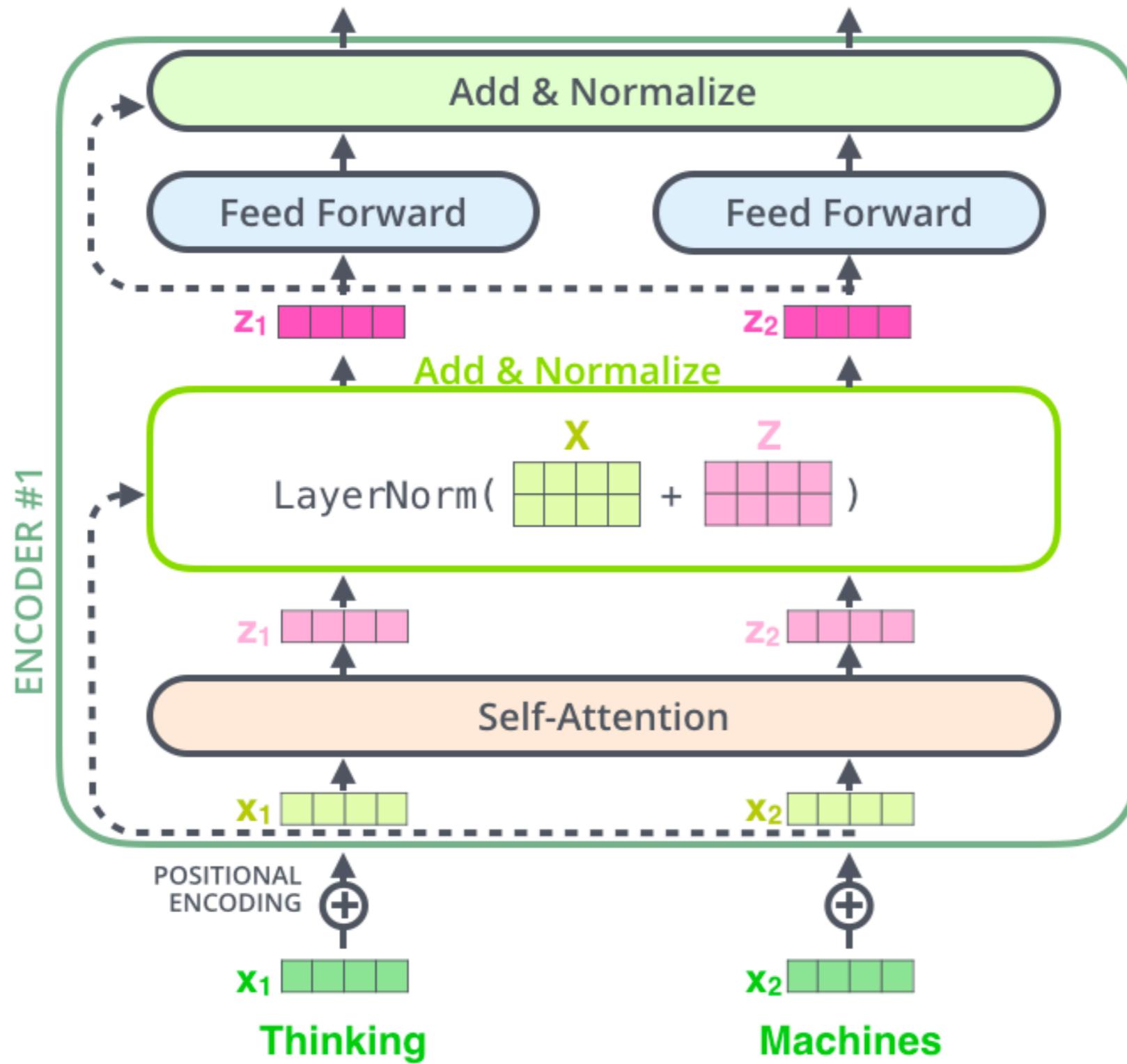
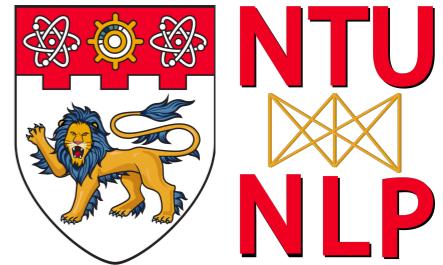
- Each sub-layer (self-attention, ffnn) in each encoder has a residual connection around it, and is followed by a layer-normalization step
- Layernorm changes input to have mean 0 and variance 1, per layer and per training point (and adds two more parameters)



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

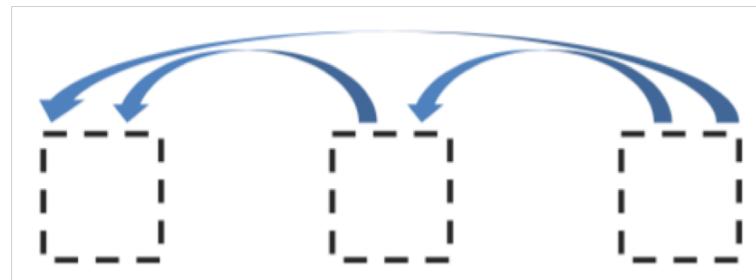
$$h_i = f\left(\frac{g_i}{\sigma_i} (a_i - \mu_i) + b_i\right)$$

# The Residuals

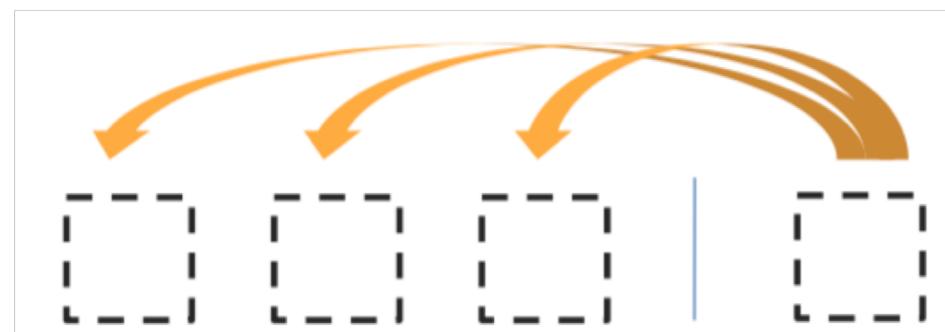


# Decoder Layers

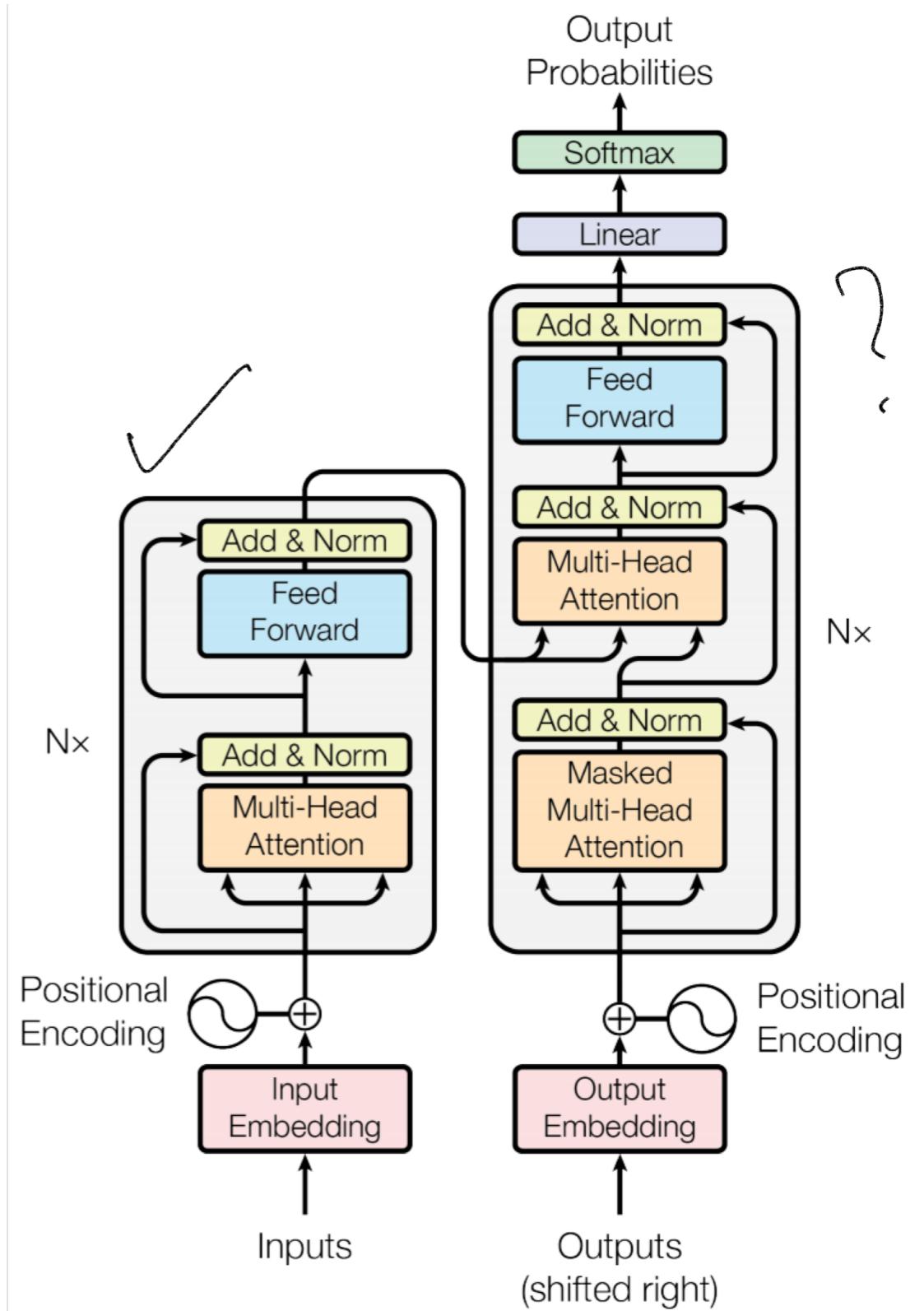
- 2 sublayer changes in decoder
- Masked decoder self-attention on previously generated outputs:



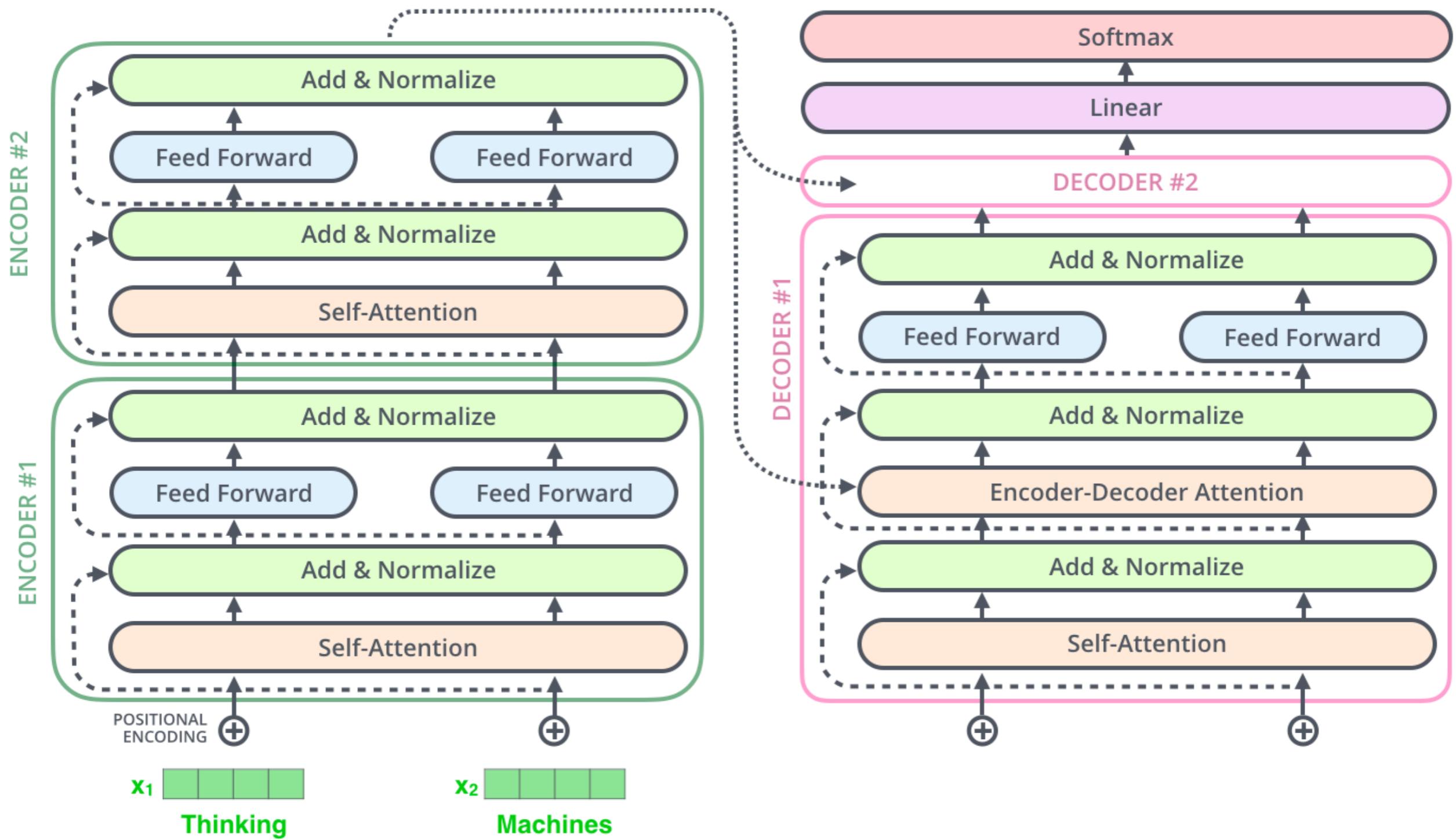
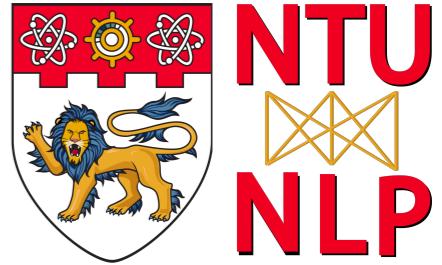
- Encoder-Decoder Attention, where queries come from previous decoder layer and keys and values come from output of encoder



- Blocks repeated 6 times also



# Whole Encoder-Decoder Model



# Transformer (Minor) Details

- Details (in paper and/or later lectures):
- Byte-pair encodings
- Checkpoint averaging
- ADAM optimizer with learning rate changes
- Dropout during training at every layer just before adding residual
- Label smoothing
- Auto-regressive decoding with beam search and length penalties
- Use of transformers is spreading but they are hard to optimize and unlike LSTMs don't usually just work out of the box and they don't play well yet with other building blocks on tasks.

# Transformer Resources

Learning about transformers on your own?

- Key recommended resource:
  - <http://nlp.seas.harvard.edu/2018/04/03/attention.html>
- The Annotated Transformer by SashaRush  
Jupyter Notebook using PyTorch that explains everything!

# Results (Translation)

Model	BLEU		Training Cost (FLOPs)	
	EN-DE	EN-FR	EN-DE	EN-FR
ByteNet [18]	23.75			
Deep-Att + PosUnk [39]		39.2		$1.0 \cdot 10^{20}$
GNMT + RL [38]	24.6	39.92	$2.3 \cdot 10^{19}$	$1.4 \cdot 10^{20}$
ConvS2S [9]	25.16	40.46	$9.6 \cdot 10^{18}$	$1.5 \cdot 10^{20}$
MoE [32]	26.03	40.56	$2.0 \cdot 10^{19}$	$1.2 \cdot 10^{20}$
Deep-Att + PosUnk Ensemble [39]		40.4		$8.0 \cdot 10^{20}$
GNMT + RL Ensemble [38]	26.30	41.16	$1.8 \cdot 10^{20}$	$1.1 \cdot 10^{21}$
ConvS2S Ensemble [9]	26.36	<b>41.29</b>	$7.7 \cdot 10^{19}$	$1.2 \cdot 10^{21}$
Transformer (base model)	27.3	38.1		$3.3 \cdot 10^{18}$
Transformer (big)	<b>28.4</b>	<b>41.8</b>		$2.3 \cdot 10^{19}$

# Results (Parsing)

<b>Parser</b>	<b>Training</b>	<b>WSJ 23 F1</b>
Vinyals & Kaiser el al. (2014) [37]	WSJ only, discriminative	88.3
Petrov et al. (2006) [29]	WSJ only, discriminative	90.4
Zhu et al. (2013) [40]	WSJ only, discriminative	90.4
Dyer et al. (2016) [8]	WSJ only, discriminative	91.7
Transformer (4 layers)	WSJ only, discriminative	91.3
Zhu et al. (2013) [40]	semi-supervised	91.3
Huang & Harper (2009) [14]	semi-supervised	91.3
McClosky et al. (2006) [26]	semi-supervised	92.1
Vinyals & Kaiser el al. (2014) [37]	semi-supervised	92.1
Transformer (4 layers)	semi-supervised	92.7
Luong et al. (2015) [23]	multi-task	93.0
Dyer et al. (2016) [8]	generative	93.3

# Summary

