

# **Attention**

## **in Deep Learning:**

### **Transformers, visual transformers...**

### **2022**



**Rufin VanRullen**  
Directeur de Recherche,  
Laboratoire Cerveau et Cognition  
(CerCo), Unité mixte CNRS-UPS



# Attention: definitions

A reminder:

- **Attention is how we select things and discard others**
- **Attention can be bottom-up (exogenous, world-dependent), or top-down (endogenous, task-dependent)**

# Models of attention

- **Psychological models**
  - Feature-Integration Theory (Treisman)
  - Guided Search (Wolfe)
- **Neuro-Computational models**
  - Saliency map (Itti/Koch)
  - Selective tuning (Tsotsos)
  - Reentry (Hamker)
- **ML (Deep Learning) models**
  - Transformers in NLP (Vaswani)
  - Transformers in computer vision [←code/notebook](#)

# Attention in Deep Learning

Sources:

Attention? Attention!

The Illustrated Transformer

Self-attention in Computer Vision

(blog: [lilianweng.github.io/lil-log/2018/06/24/attention-attention.html](http://lilianweng.github.io/lil-log/2018/06/24/attention-attention.html))

(blog: [jalammar.github.io/illustrated-transformer/](http://jalammar.github.io/illustrated-transformer/))

(blog: [towardsdatascience.com/self-attention-in-computer-vision-2782727021f6](http://towardsdatascience.com/self-attention-in-computer-vision-2782727021f6))

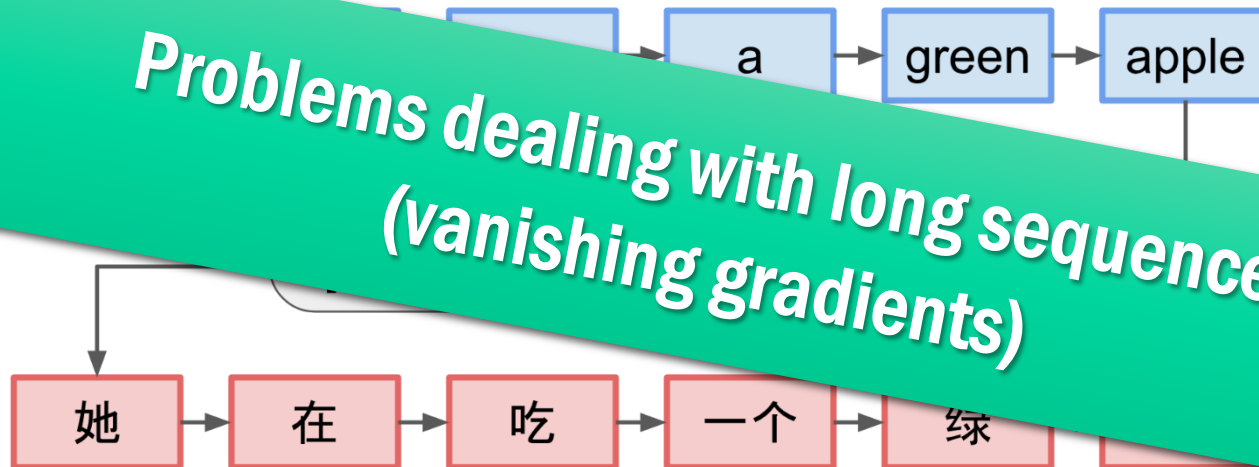
- **Seq2Seq models in Natural Language Processing (NLP)**

- Machine translation
- Question-Answer dialog (chatbots)
- Document summarization
- ...

Encoder  
RNN  
(LSTM/GRU units)

Decoder  
RNN  
(LSTM/GRU units)

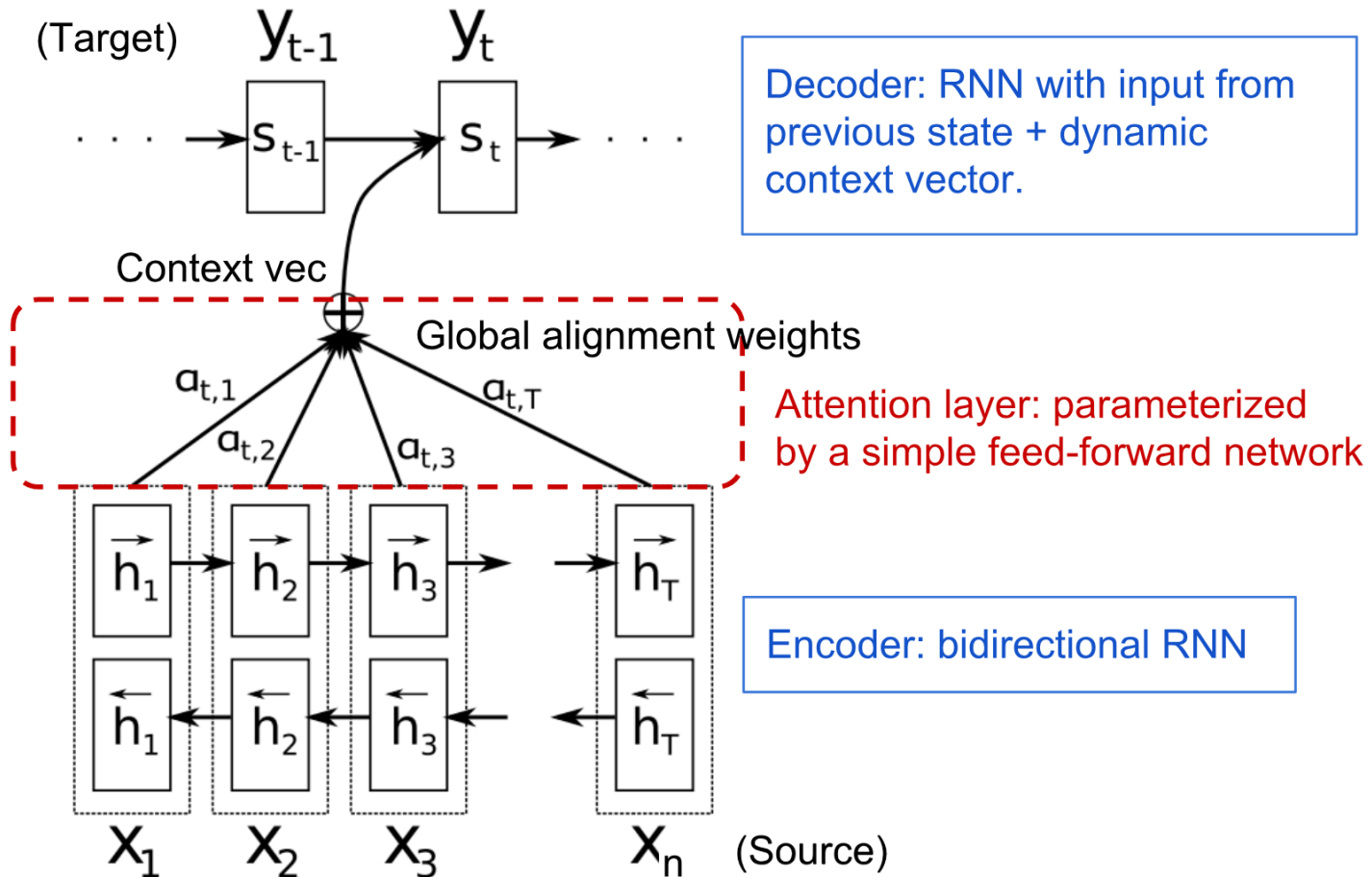
*Problems dealing with long sequences  
(vanishing gradients)*



Sutskever et al (2014)

# Attention in Deep Learning

- Solving the vanishing gradients problem with attention



Bahdanau et al (2015)

# Attention in Deep Learning

- Solving the vanishing gradients problem with attention

Source sequence  $\mathbf{x} = [x_1, x_2, \dots, x_n]$   
 Output Sequence  $\mathbf{y} = [y_1, y_2, \dots, y_m]$

Attention score:  $\text{score}(\mathbf{s}_t, \mathbf{h}_i) = \mathbf{v}_a^\top \tanh(\mathbf{W}_a [\mathbf{s}_t, \mathbf{h}_i])$

where both  $\mathbf{v}_a$  and  $\mathbf{W}_a$  are weight matrices to be learned

Input:  
 Output:  
 Attention:  
 sequence  
 sequence  
 = a form of memory pointer?

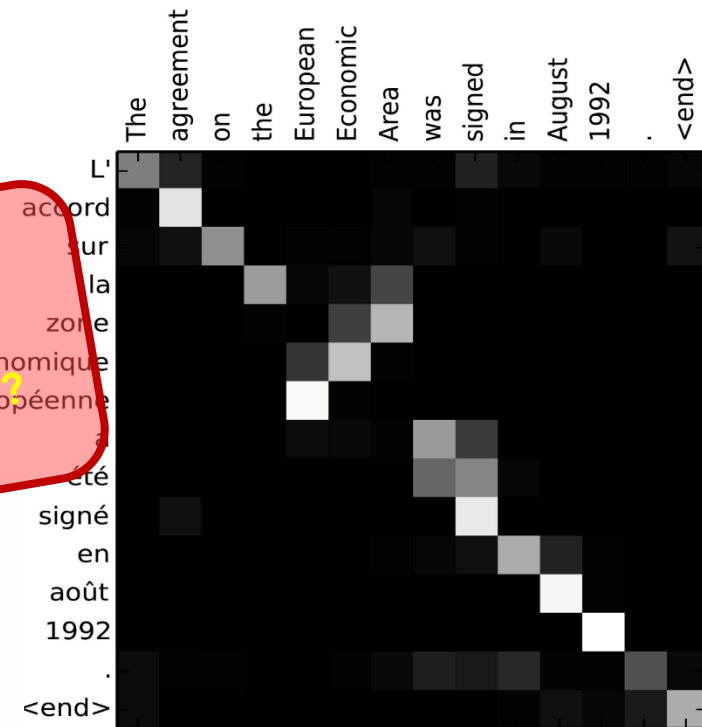
The decoder network has hidden state  $\mathbf{s}_t = f(\mathbf{s}_{t-1}, \mathbf{y}_{t-1}, \mathbf{c}_t)$  for the output word at position  $t$ ,  $t = 1, \dots, m$ , where the context vector  $\mathbf{c}_t$  is a sum of hidden states of the input sequence, weighted by alignment scores:

$$\mathbf{c}_t = \sum_{i=1}^n \alpha_{t,i} \mathbf{h}_i \quad ; \text{Context vector for output } y_t$$

$\alpha_{t,i} = \text{align}(y_t, x_i)$  ; How well two words  $y_t$  and  $x_i$  are aligned.

$$= \frac{\exp(\text{score}(\mathbf{s}_{t-1}, \mathbf{h}_i))}{\sum_{i'=1}^n \exp(\text{score}(\mathbf{s}_{t-1}, \mathbf{h}_{i'}))} \quad ; \text{Softmax of some predefined alignment score..}$$

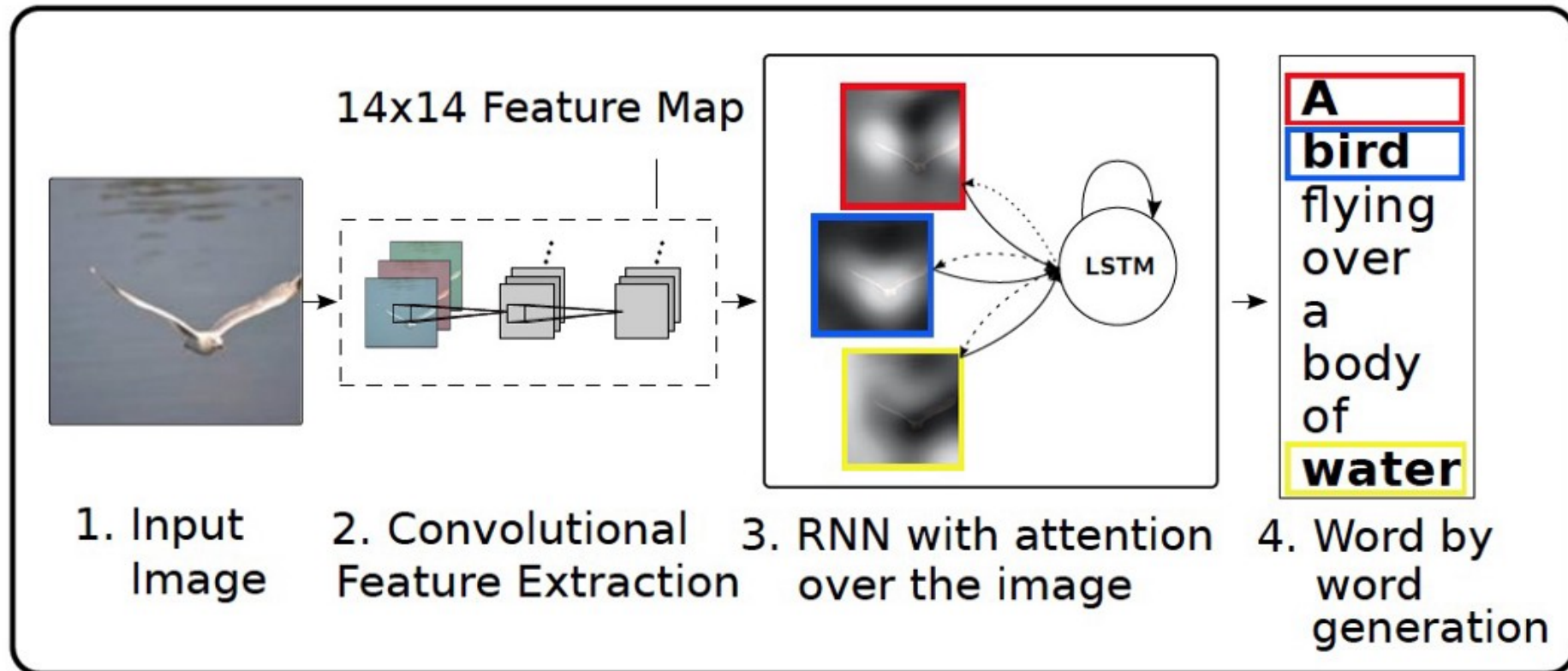
French → English translation



Bahdanau et al (2015)

# Attention in Deep Learning

- Show, attend and tell

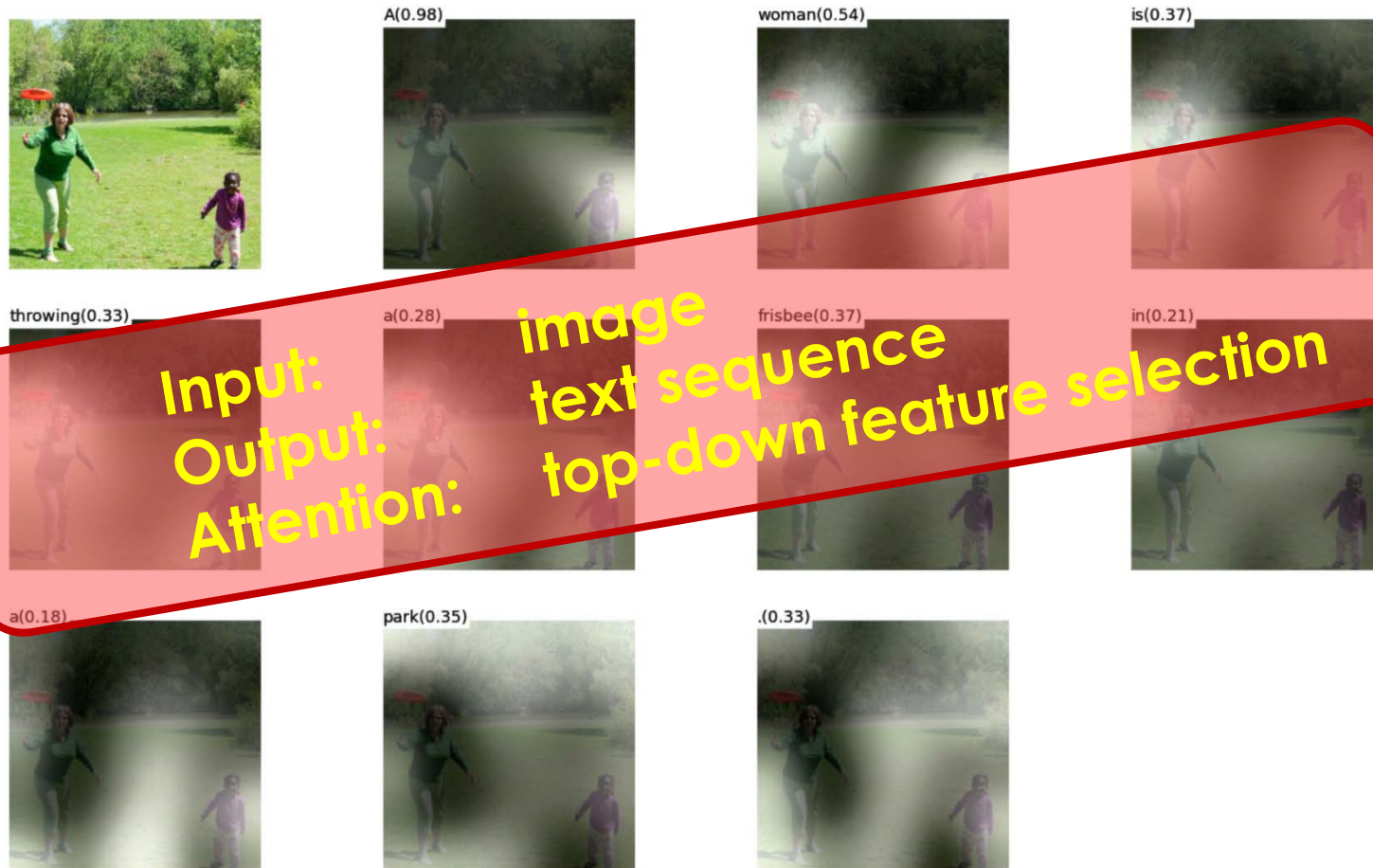


Xu et al (2015)

# Attention in Deep Learning

- Show, attend and tell

*“A woman is throwing a frisbee in a park.”*

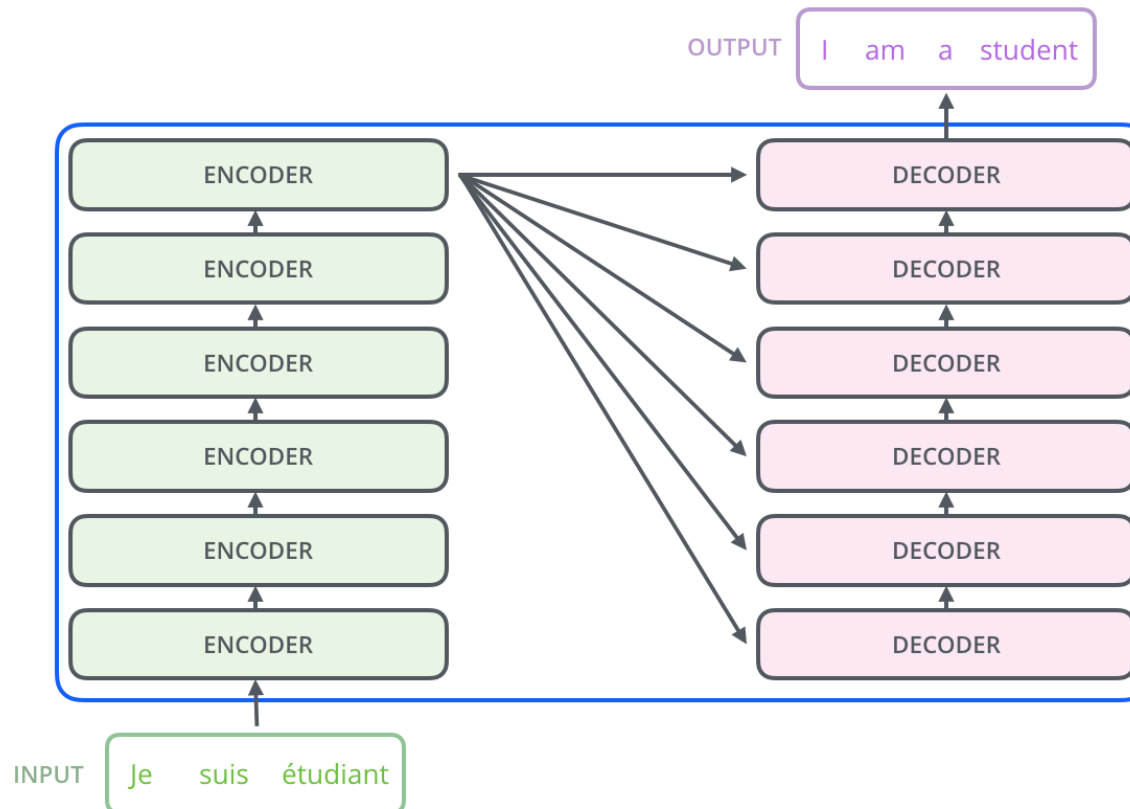


Xu et al (2015)



# Attention in Deep Learning

- **Back to NLP: Attention is all you need (Transformers)!**
  - Main idea: no recurrence, feed-forward architecture, all inputs are provided at once (very large matrix)
  - Attention takes care of long-range dependencies

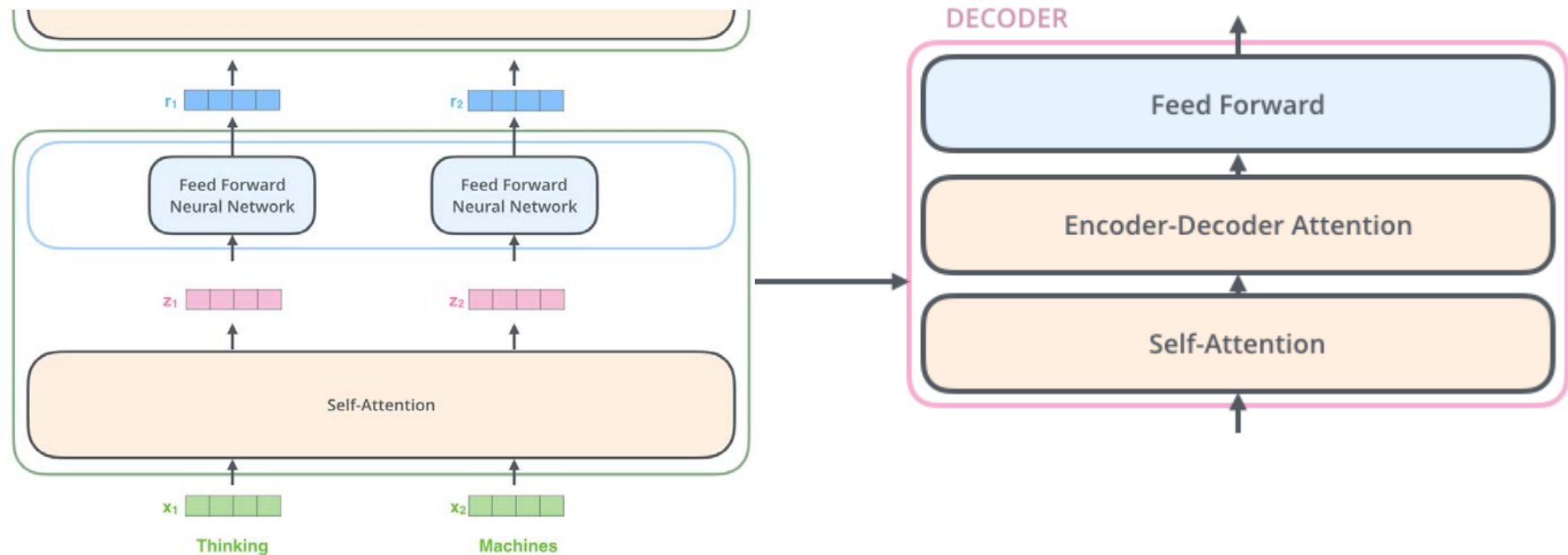


Vaswani et al (2017)

# Attention in Deep Learning

## • Transformer details

- Self-attention layer: every token/position can « attend » to every other in the same layer
- Encoder-decoder-attention: every position in the decoder can attend to every position in the final encoder layer (as in Seq2Seq models)
- The same feed-forward network is applied independently to each position

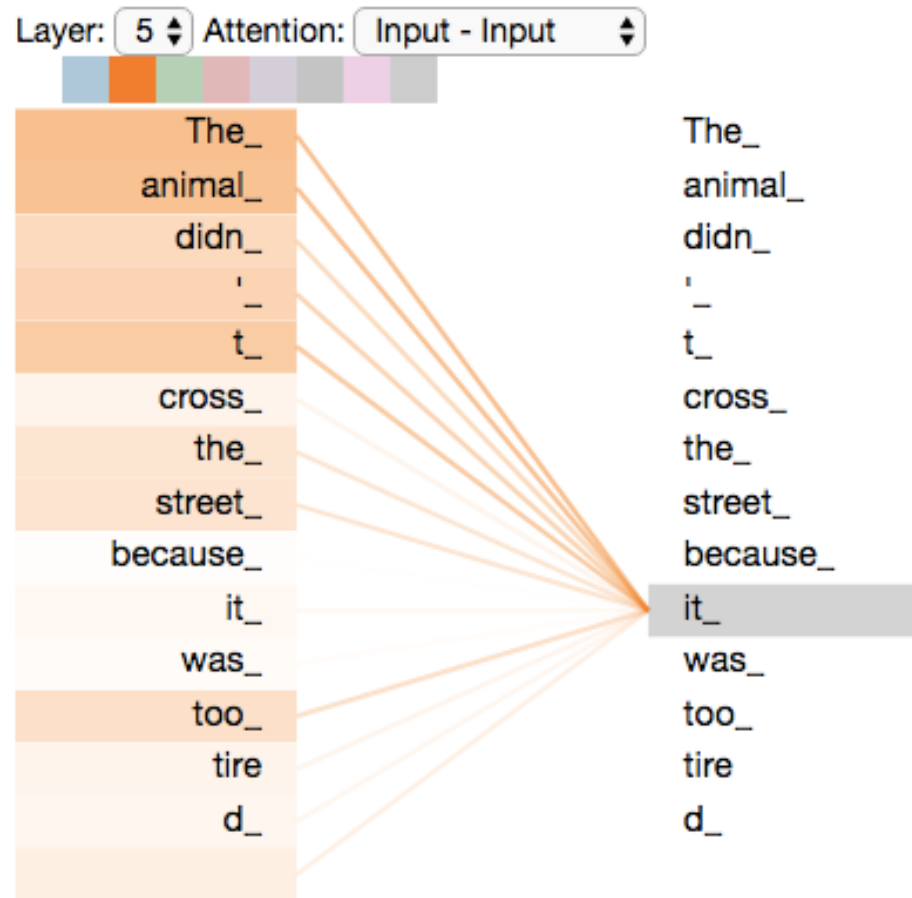


Vaswani et al (2017)

# Attention in Deep Learning

- Transformer details

- Self-attention example

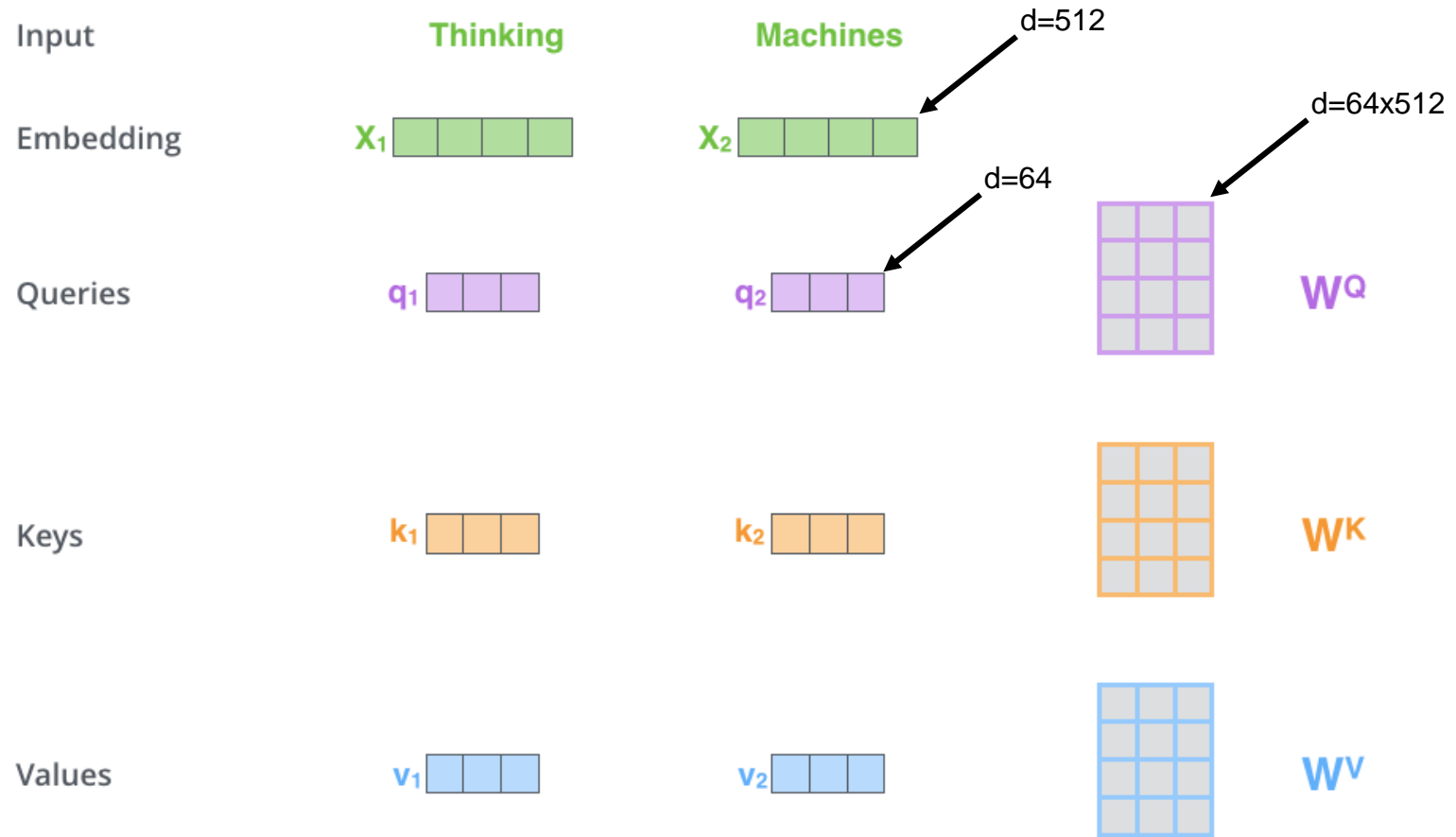


Vaswani et al (2017)

# Attention in Deep Learning

- Transformer details

- Self-attention computation: **Queries, Keys, Values**

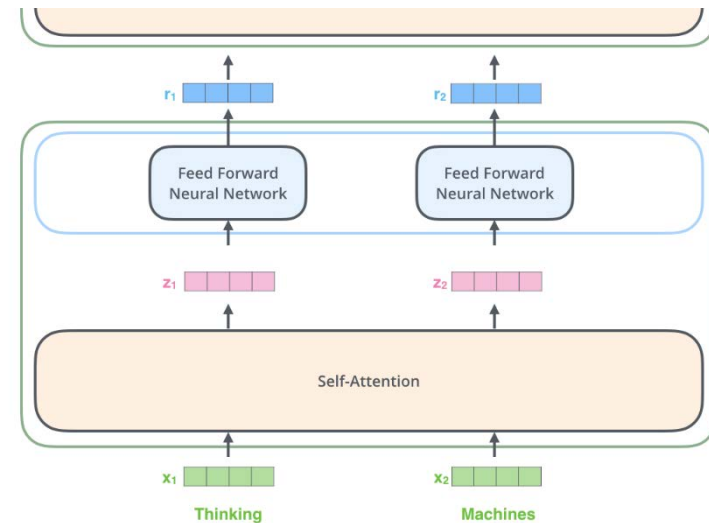
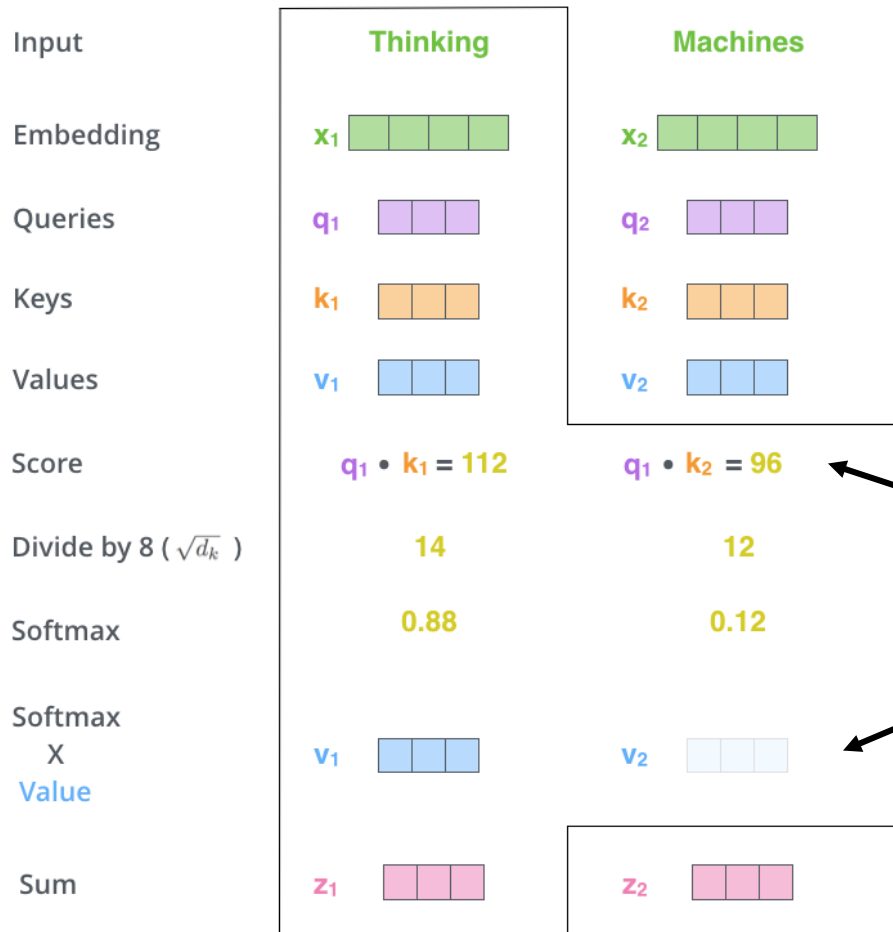


Vaswani et al (2017)

# Attention in Deep Learning

## • Transformer details

- Self-attention computation: **Queries, Keys, Values**



The match between **query** and **key** determines how much of each **value** is included in the final output

Vaswani et al (2017)

# Attention in Deep Learning

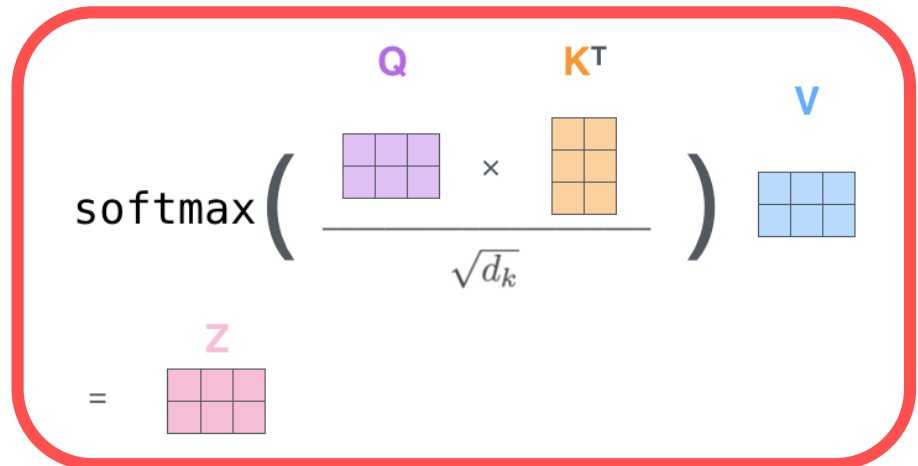
- Transformer details

- Self-attention computation: matrix version

$$X \times W^Q = Q$$


$$X \times W^K = K$$


$$X \times W^V = V$$


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


**Pay attention...** This is what we'll implement in the notebook!

Vaswani et al (2017)

# Attention in Deep Learning

## • Transformer details

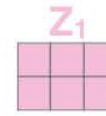
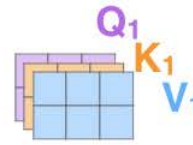
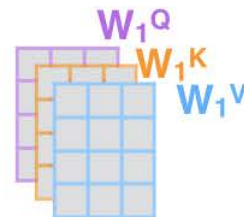
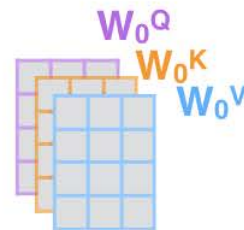
### • Multi-head attention

- 1) This is our input sentence\*
- 2) We embed each word\*
- 3) Split into 8 heads. We multiply  $X$  or  $R$  with weight matrices
- 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



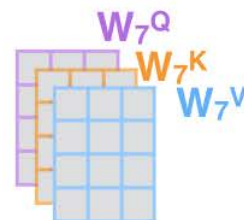
\* In all encoders other than #0, we don't need embedding. We start directly with the output of the encoder right below this one



...

...

...

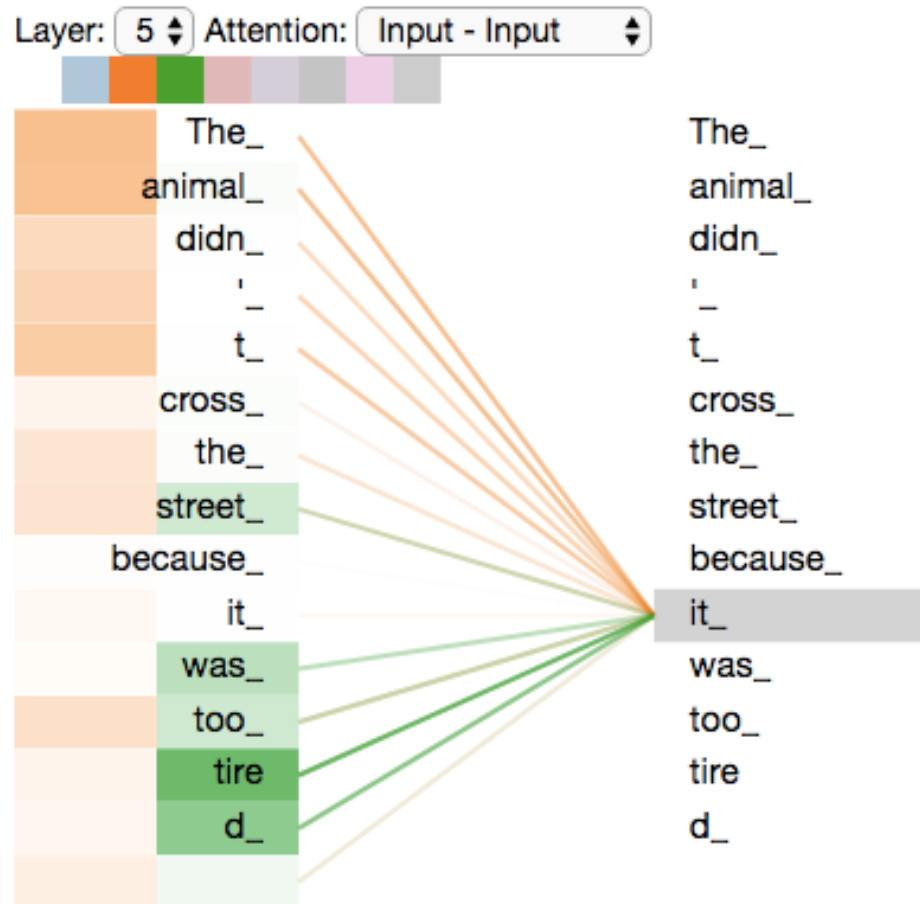


Vaswani et al (2017)

# Attention in Deep Learning

- **Transformer details**

- Multi-head attention



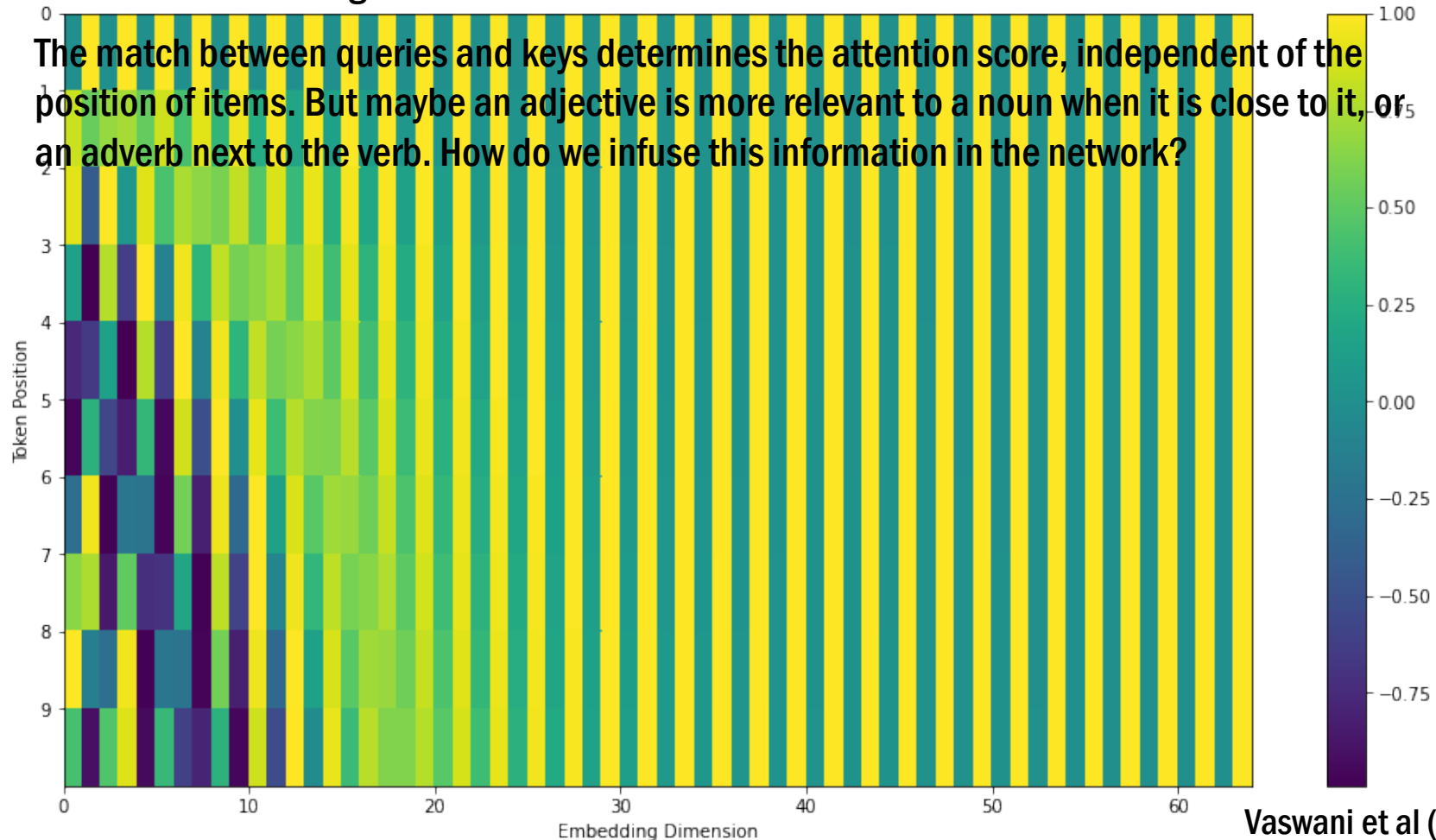
Vaswani et al (2017)



# Attention in Deep Learning

- Transformer details

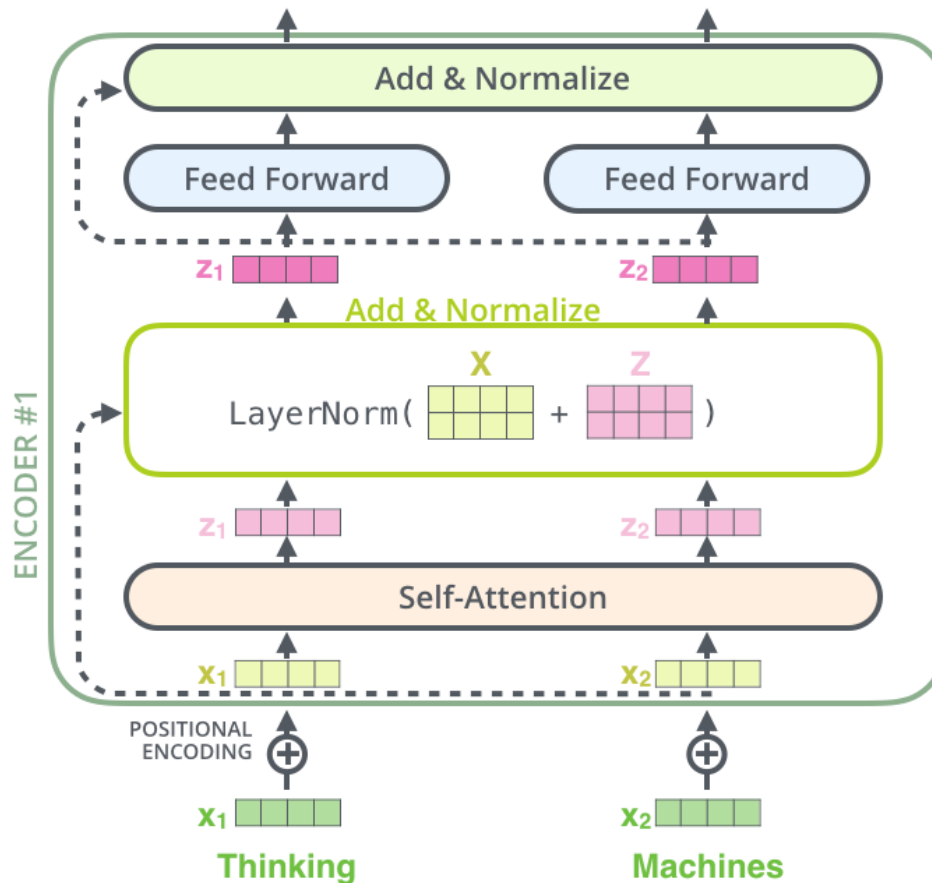
- Positional encoding



# Attention in Deep Learning

- **Transformer details**

- Residual connections

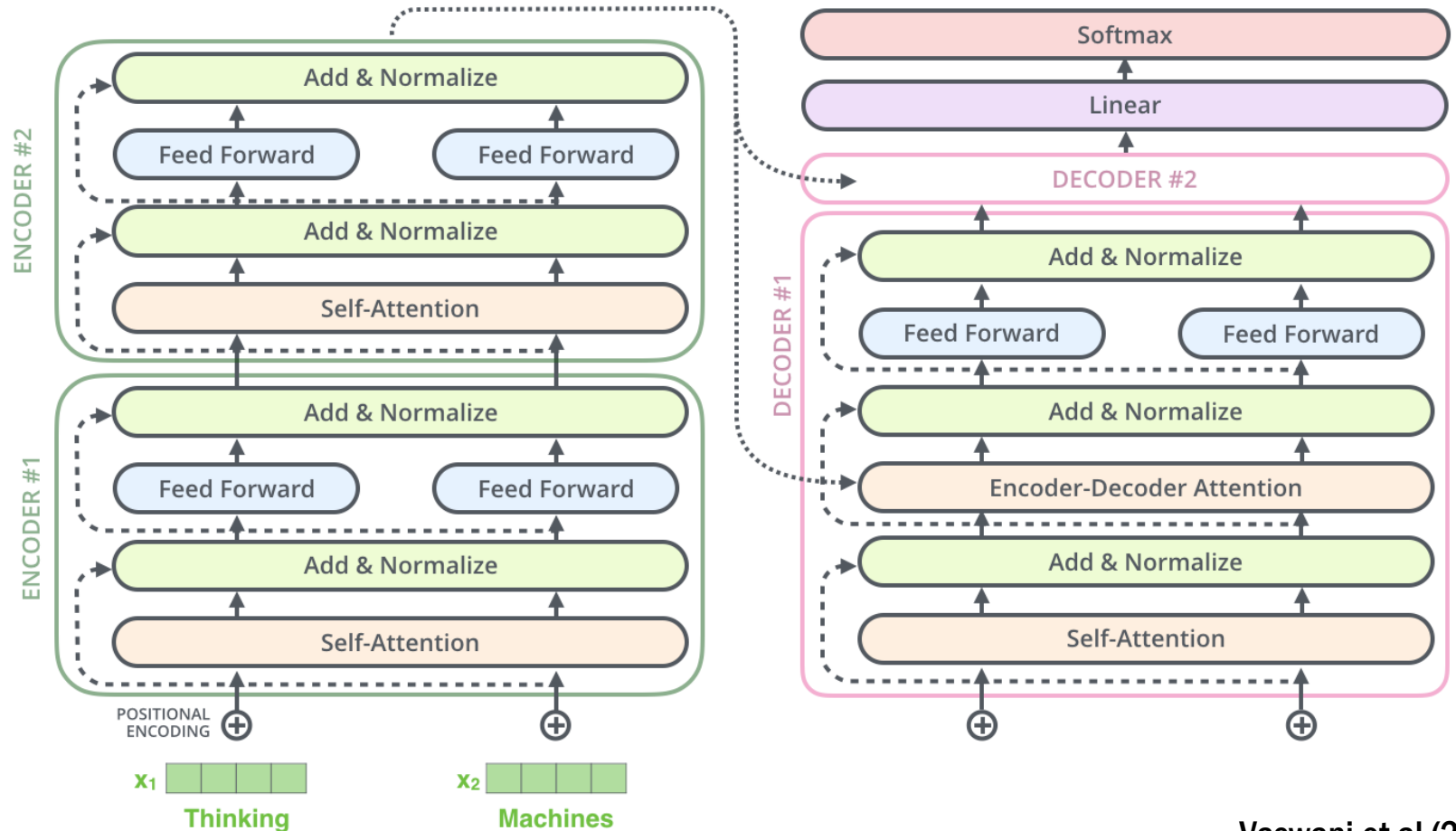


Vaswani et al (2017)

# Attention in Deep Learning

- Transformer details

- Encoder + Decoder

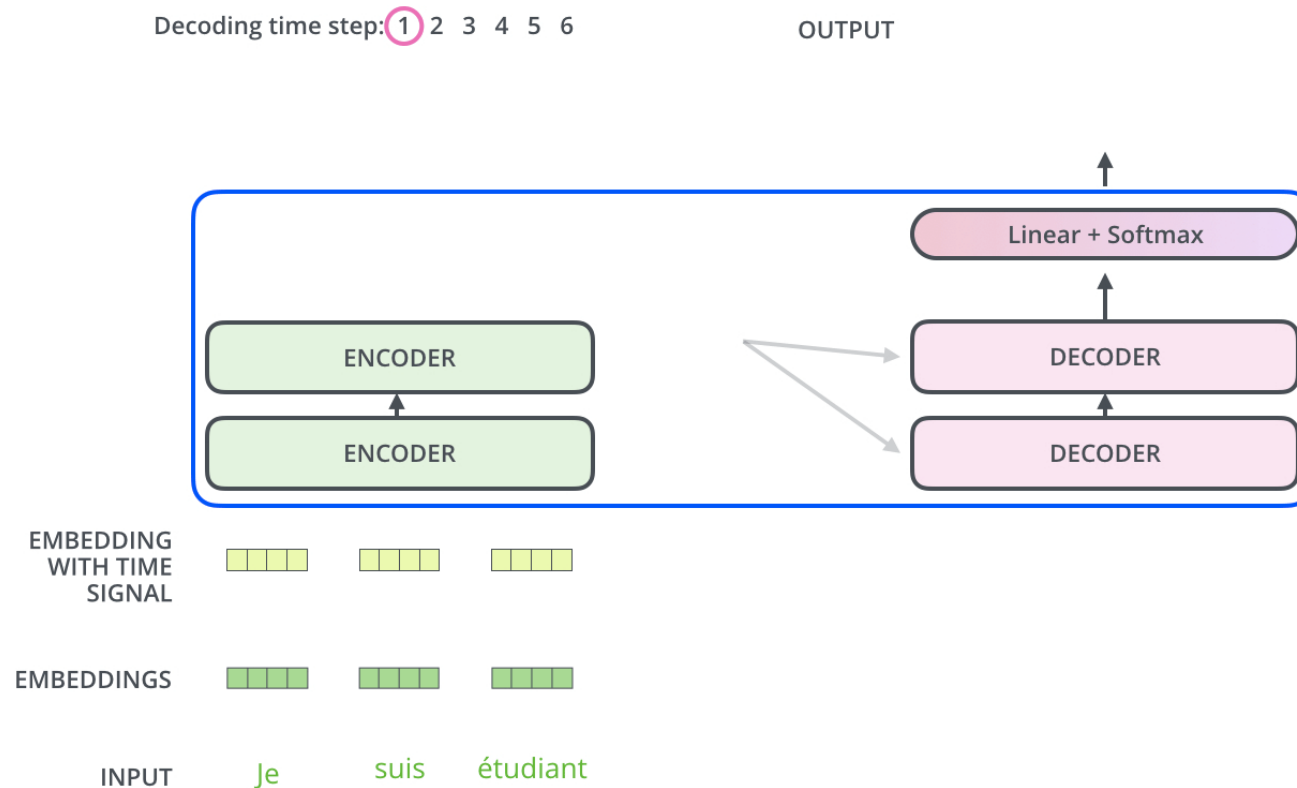


Vaswani et al (2017)

# Attention in Deep Learning

- Transformer details

- The transformer at work... (encoding)

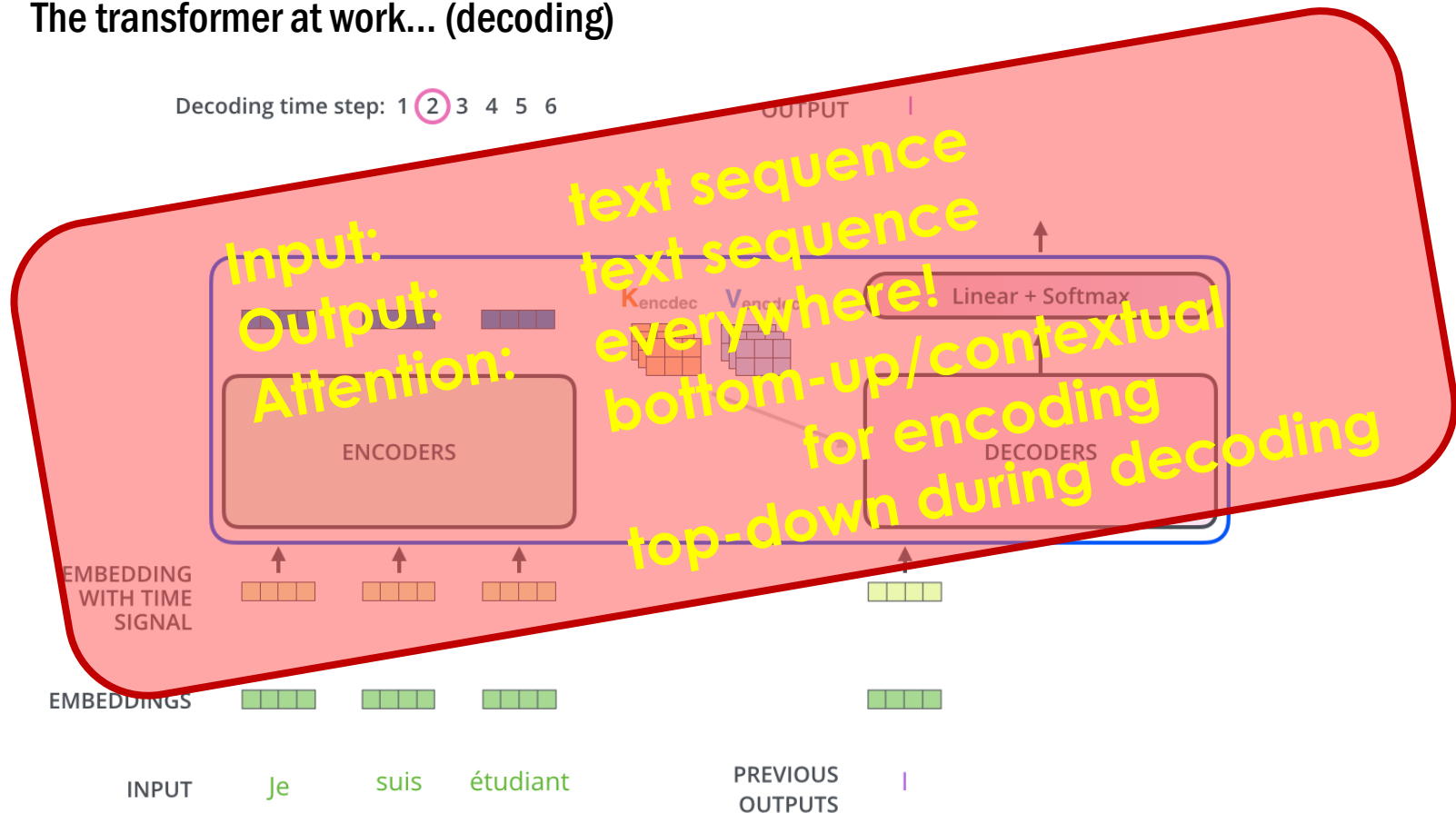


Vaswani et al (2017)

# Attention in Deep Learning

- Transformer details

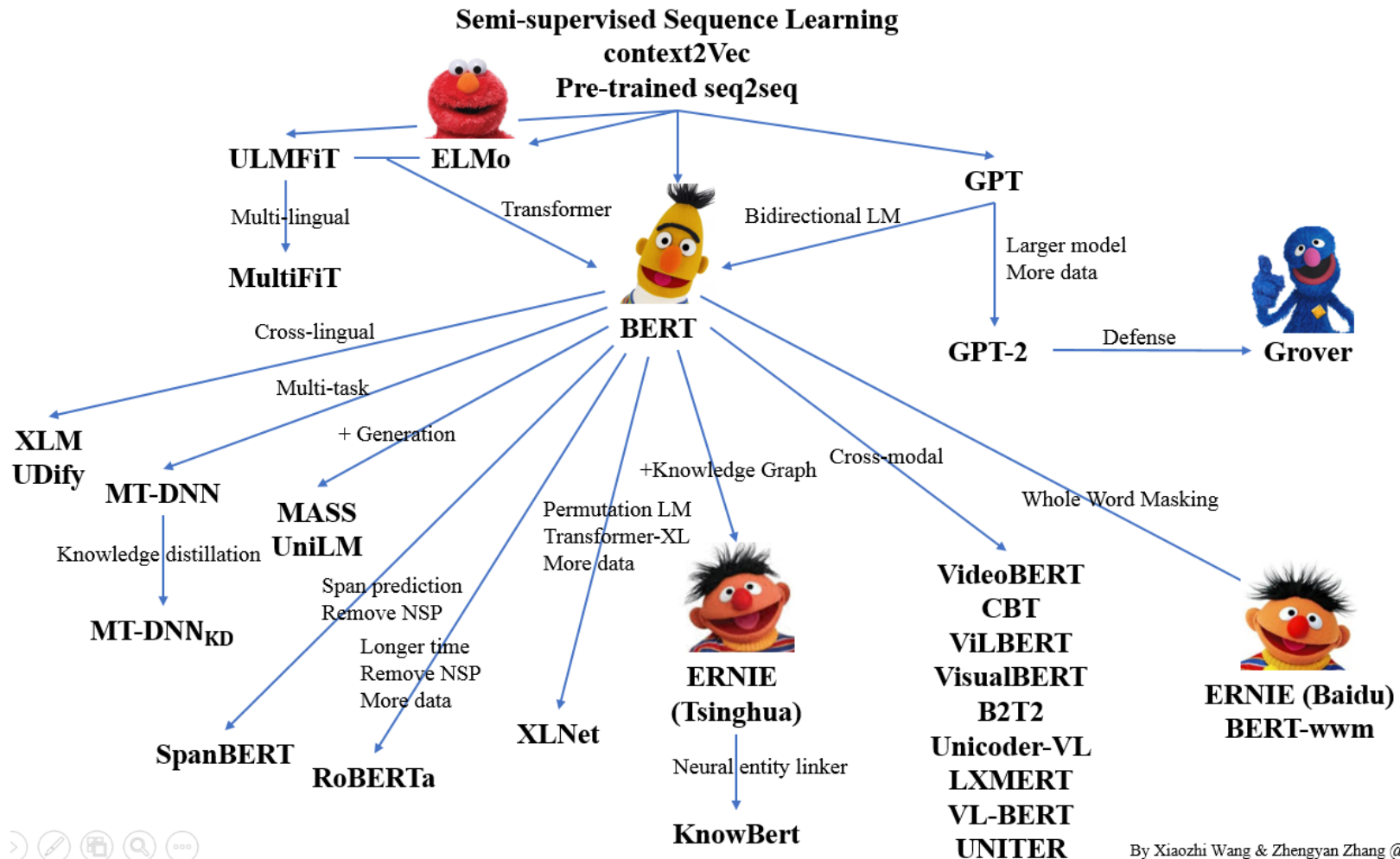
- The transformer at work... (decoding)



Vaswani et al (2017)

# Attention in Deep Learning

- Family tree of NLP Transformers:



By Xiaozhi Wang & Zhengyan Zhang @THUNLP

# Attention in Deep Learning

- **Transformers (attention) in computer vision**

- Many recent papers have thought to augment or replace convolution operations in deep networks with a self-attention mechanism
- One example: Ramachandran et al (2019), Stand-alone self-attention in vision models, NeurIPS.

Single-headed attention for computing the pixel output  $y_{ij} \in \mathbb{R}^{d_{out}}$  is then computed as follows (see Figure 3):

$$y_{ij} = \sum_{a,b \in \mathcal{N}_k(i,j)} \text{softmax}_{ab} (q_{ij}^\top k_{ab}) v_{ab} \quad (2)$$

where the *queries*  $q_{ij} = W_Q x_{ij}$ , *keys*  $k_{ab} = W_K x_{ab}$ , and *values*  $v_{ab} = W_V x_{ab}$  are linear transformations of the pixel in position  $ij$  and the neighborhood pixels.  $\text{softmax}_{ab}$  denotes a softmax applied to all logits computed in the neighborhood of  $ij$ .  $W_Q, W_K, W_V \in \mathbb{R}^{d_{out} \times d_{in}}$  are all learned transforms.

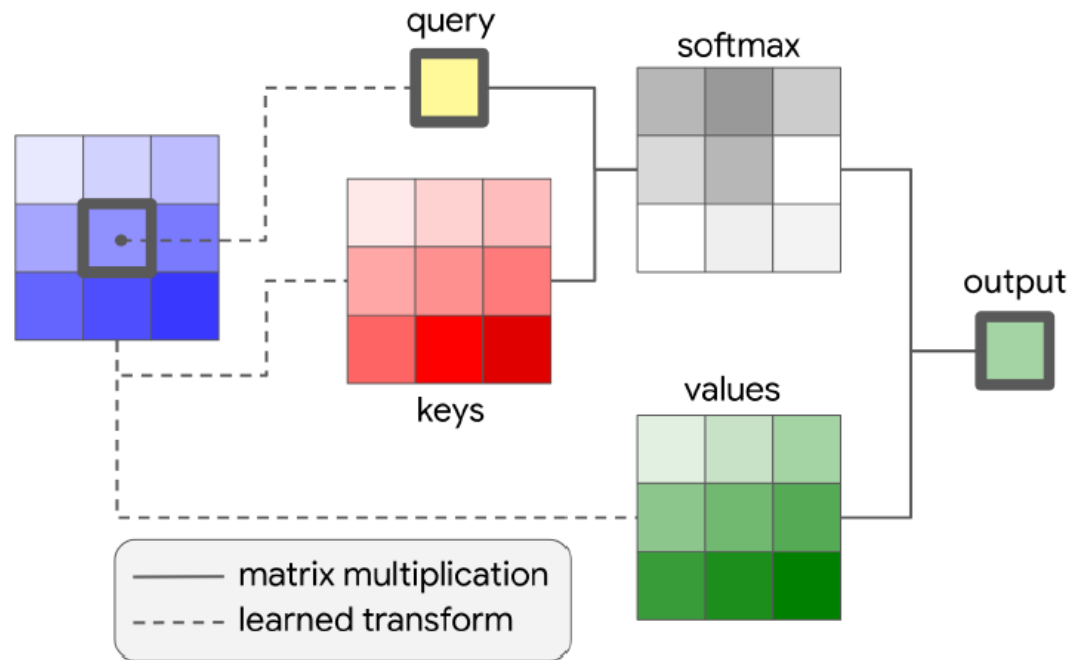


**Pay attention... This is what we'll implement in the notebook!**

# Attention in Deep Learning

- **Transformers (attention) in computer vision**

- Many recent papers have thought to augment or replace convolution operations in deep networks with a self-attention mechanism
- One example: Ramachandran et al (2019), Stand-alone self-attention in vision models, NeurIPS.





# Attention in Deep Learning

- **Transformers (attention) in computer vision**

- Many recent papers have thought to augment or replace convolution operations in deep networks with a self-attention mechanism
- One example: Ramachandran et al (2019), Stand-alone self-attention in vision models, NeurIPS.
- Attention is permutation-invariant → position encoding (as in NLP transformers)

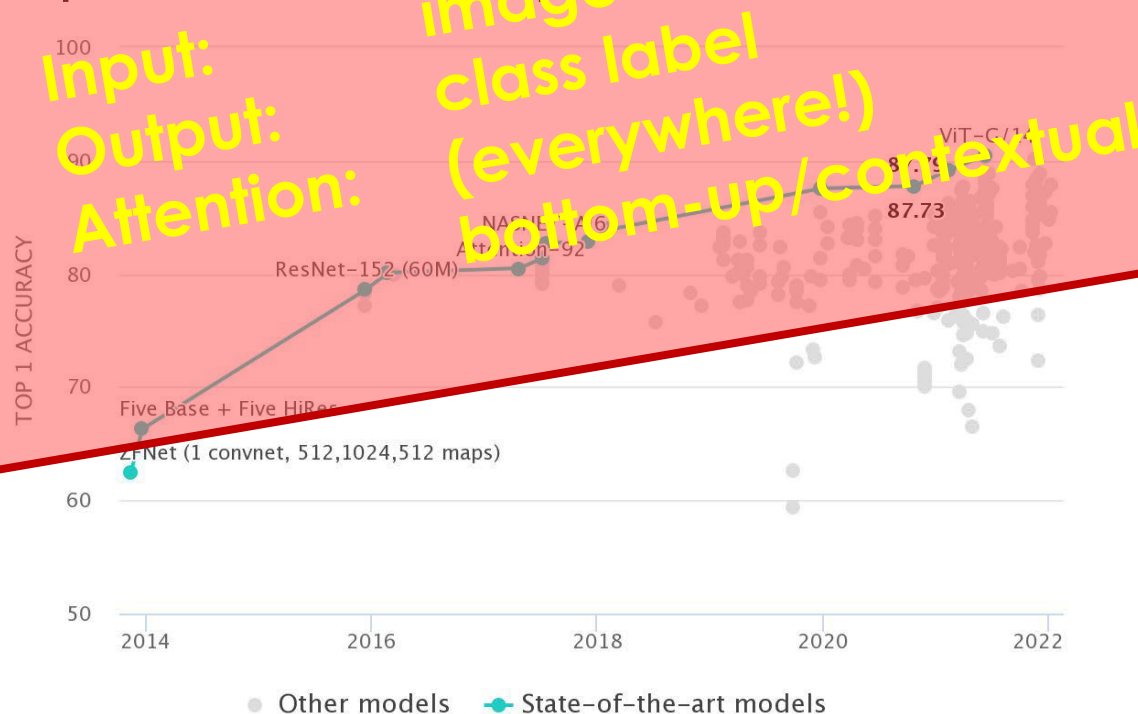
$$y_{ij} = \sum_{a,b \in \mathcal{N}_k(i,j)} \text{softmax}_{ab} \left( q_{ij}^\top k_{ab} + q_{ij}^\top r_{a-i,b-j} \right) v_{ab} \quad (3)$$

Thus, the logit measuring the similarity between the query and an element in  $\mathcal{N}_k(i, j)$  is modulated both by the content of the element and the relative distance of the element from the query. Note that by infusing relative position information, self-attention also enjoys translation equivariance, similar to convolutions.

# Attention in Deep Learning

- **Transformers (attention) in computer vision**

- Many recent papers have thought to augment or replace convolution operations in deep networks with a self-attention mechanism
- One example: Ramachandran et al (2019), Stand-alone self-attention in vision models, NeurIPS.
- Accuracy on ImageNet (ResNet50 backbone, varying channel numbers):
- Comparable performance with  $\frac{1}{2}$  the parameters!



# Attention in Deep Learning

- **Transformers (attention) in computer vision**

- **Colab Notebook:**

[https://github.com/rufinv/DL-Attention-Class/blob/main/VanRullen\\_AttentionClass\\_CIFAR10\\_2022.ipynb](https://github.com/rufinv/DL-Attention-Class/blob/main/VanRullen_AttentionClass_CIFAR10_2022.ipynb)

**Implement, train and test a  
Deep Convolutional Neural Network (CNN)  
for image classification (CIFAR)  
with/without attention**

**go to link now...**