

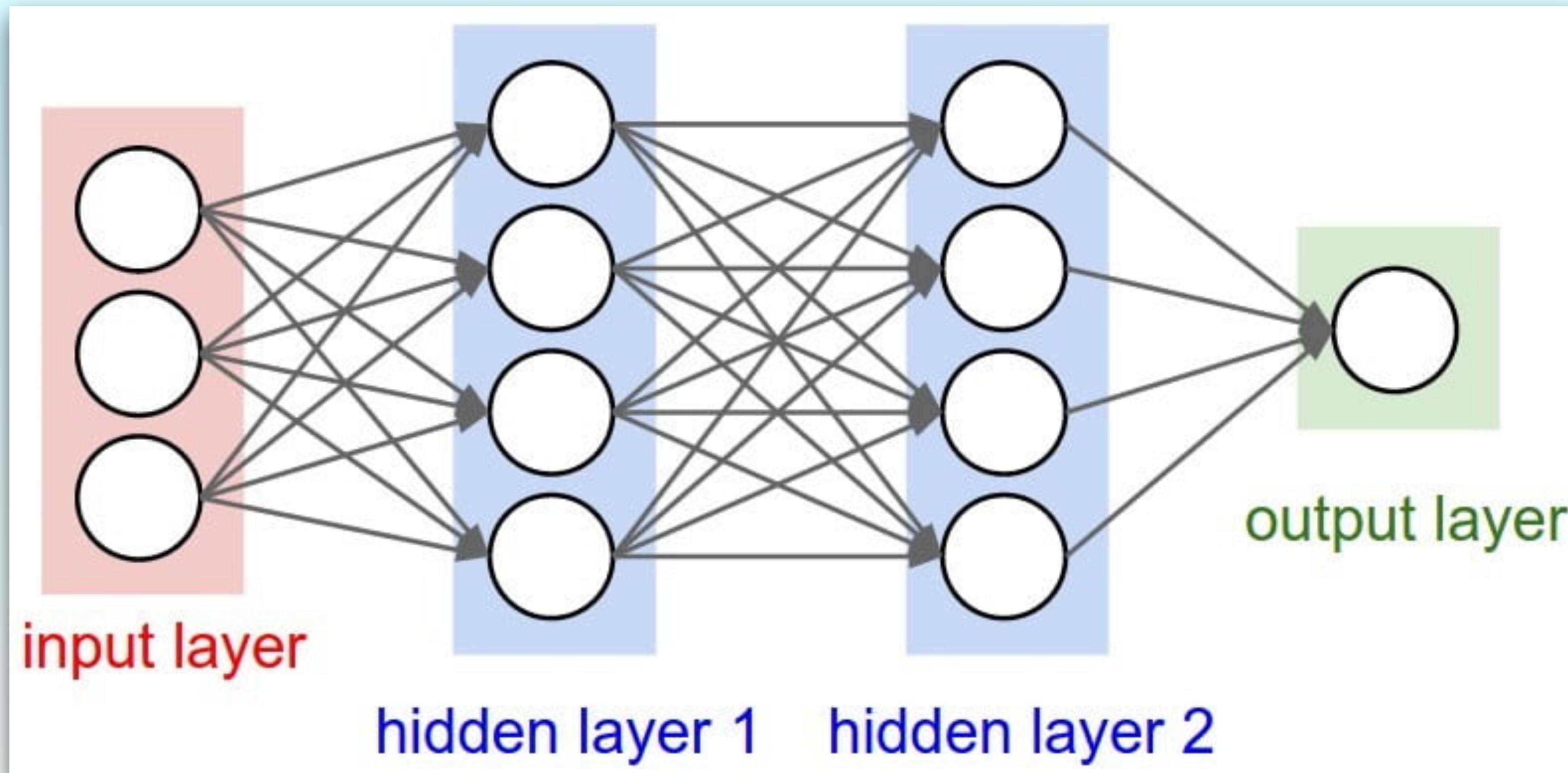
Transformers

LauzHack Deep Learning Summer Bootcamp

Seyed Parsa Neshaei – July 2024

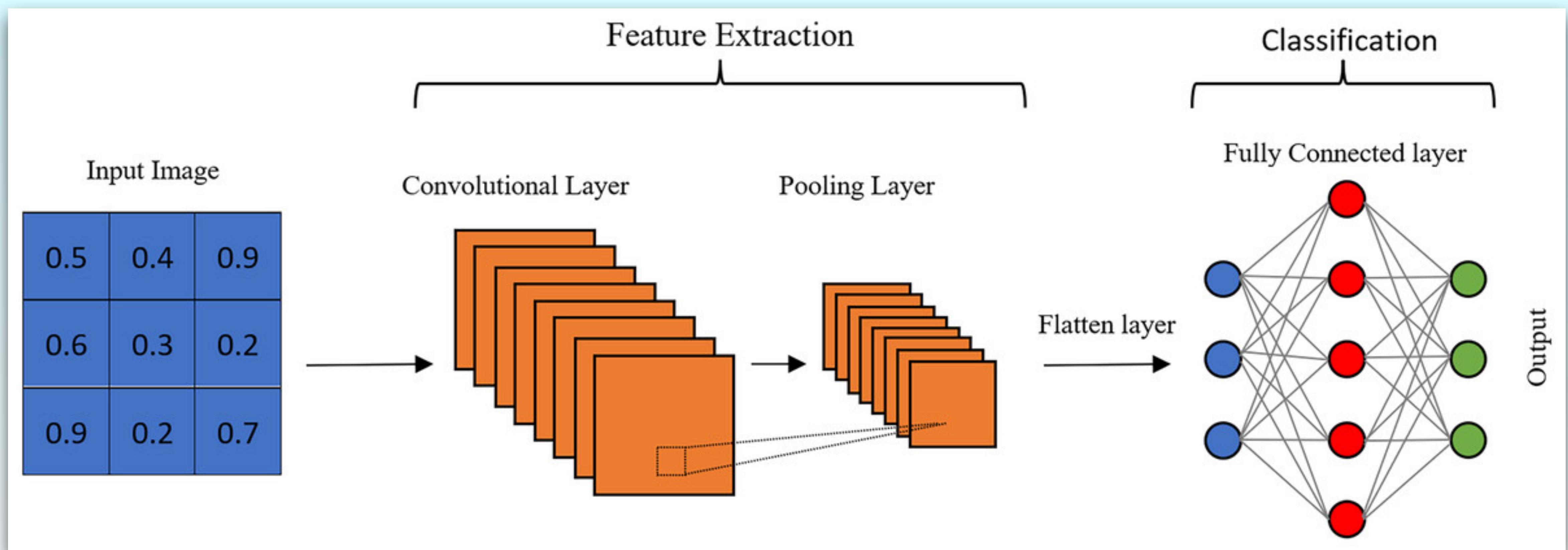
Let's have a recap!

Deep Neural Networks



Let's have a recap!

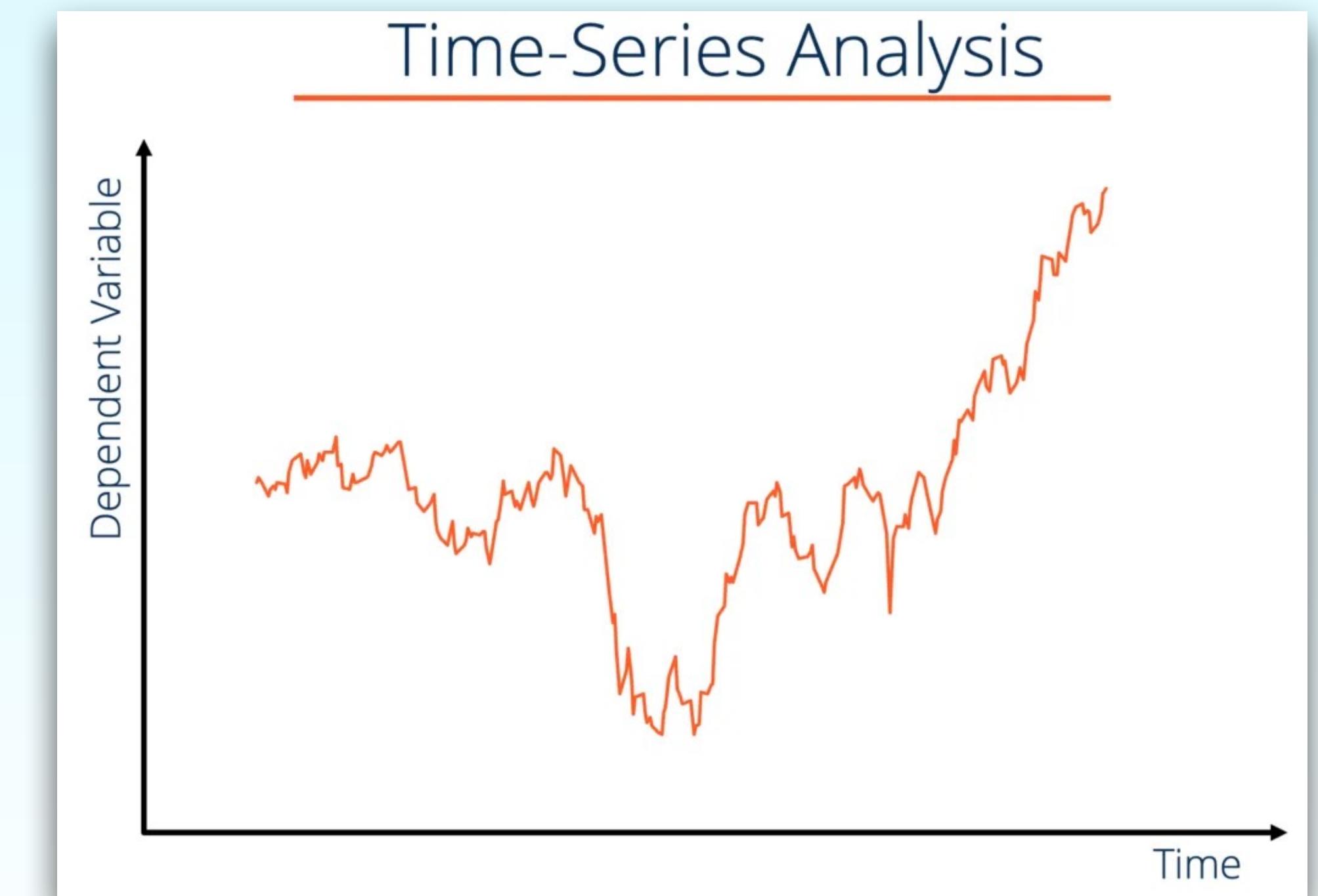
Convolutional Neural Networks



A new task!

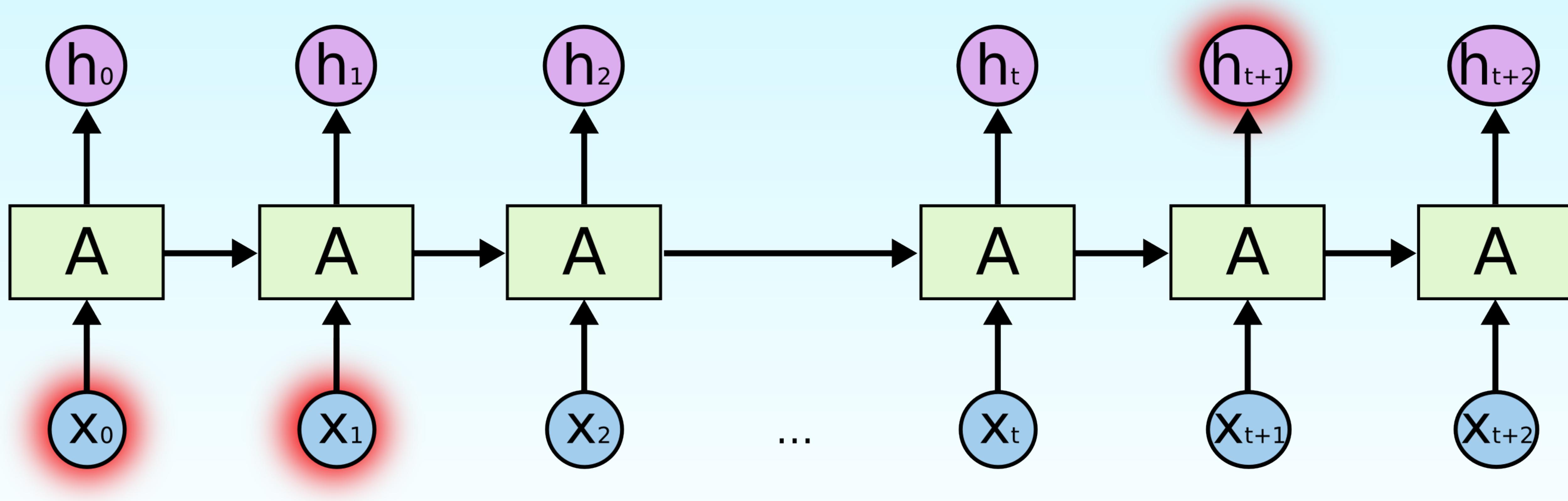
Time series

- Data measured at regular intervals, or step-by-step. Examples:
 - Hourly temperature in a city
 - Heights of ocean tides
 - Heartbeats per minute
 - The stock price of a company
 - Answers of students to exam questions
 - A written text (!): *EPFL I like vs. I like EPFL*
- Tasks: *classification* or *forecasting*
- Tip: if data semantics change by permutation



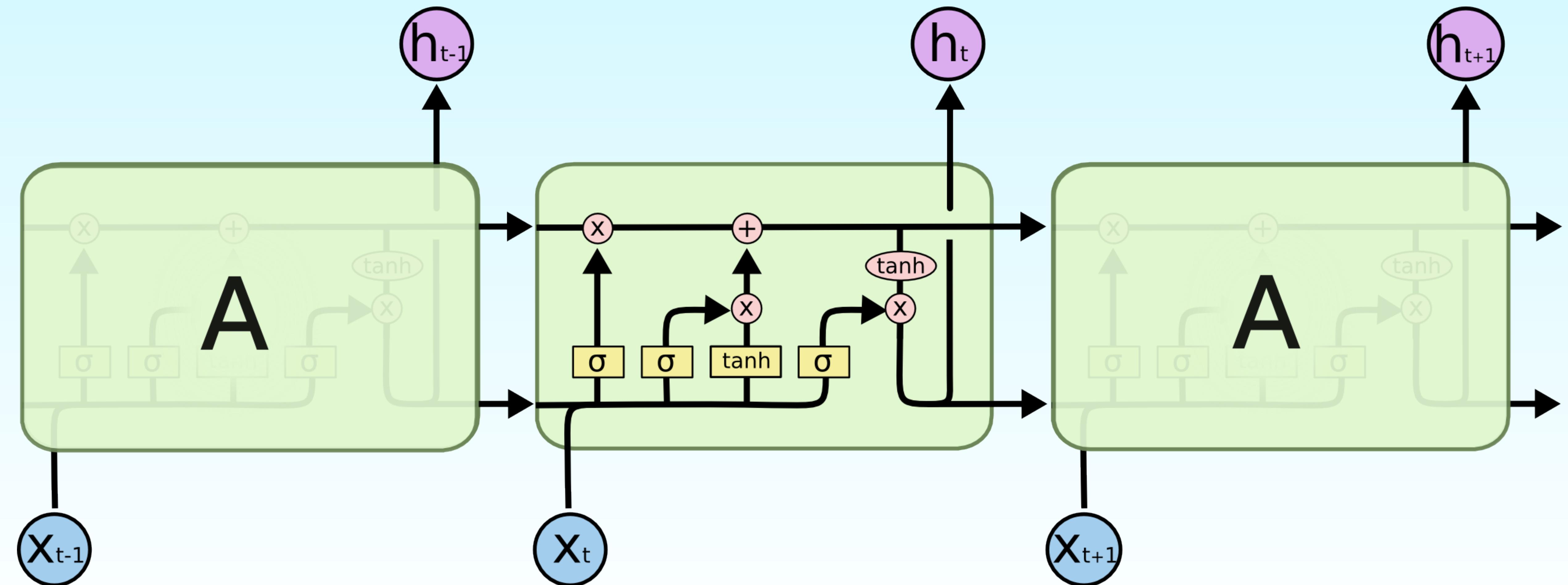
Let's have a recap!

Recurrent Neural Networks

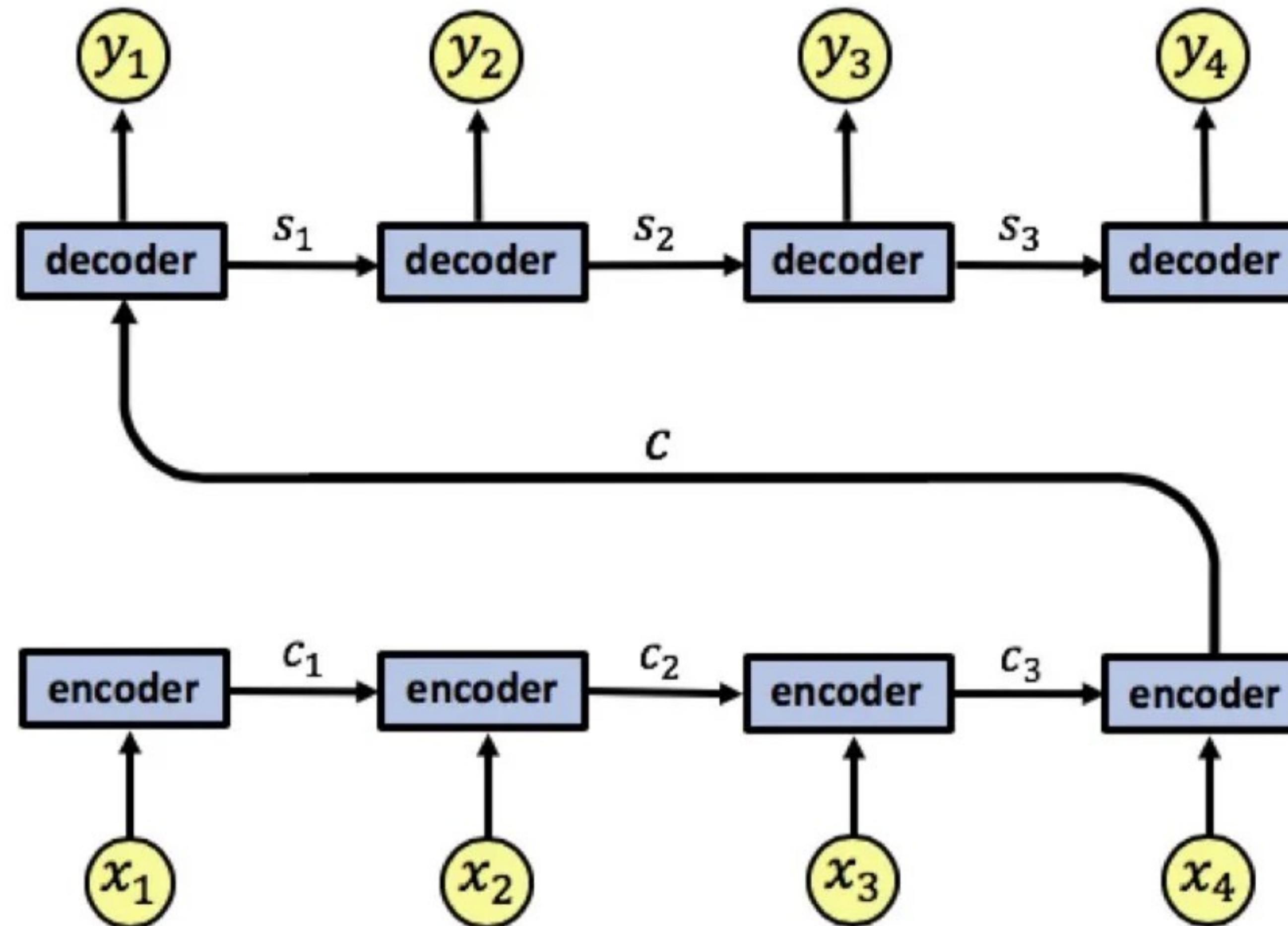


Let's have a recap!

Long-Short Term Memory Networks



Seq2seq Models



A new method today!

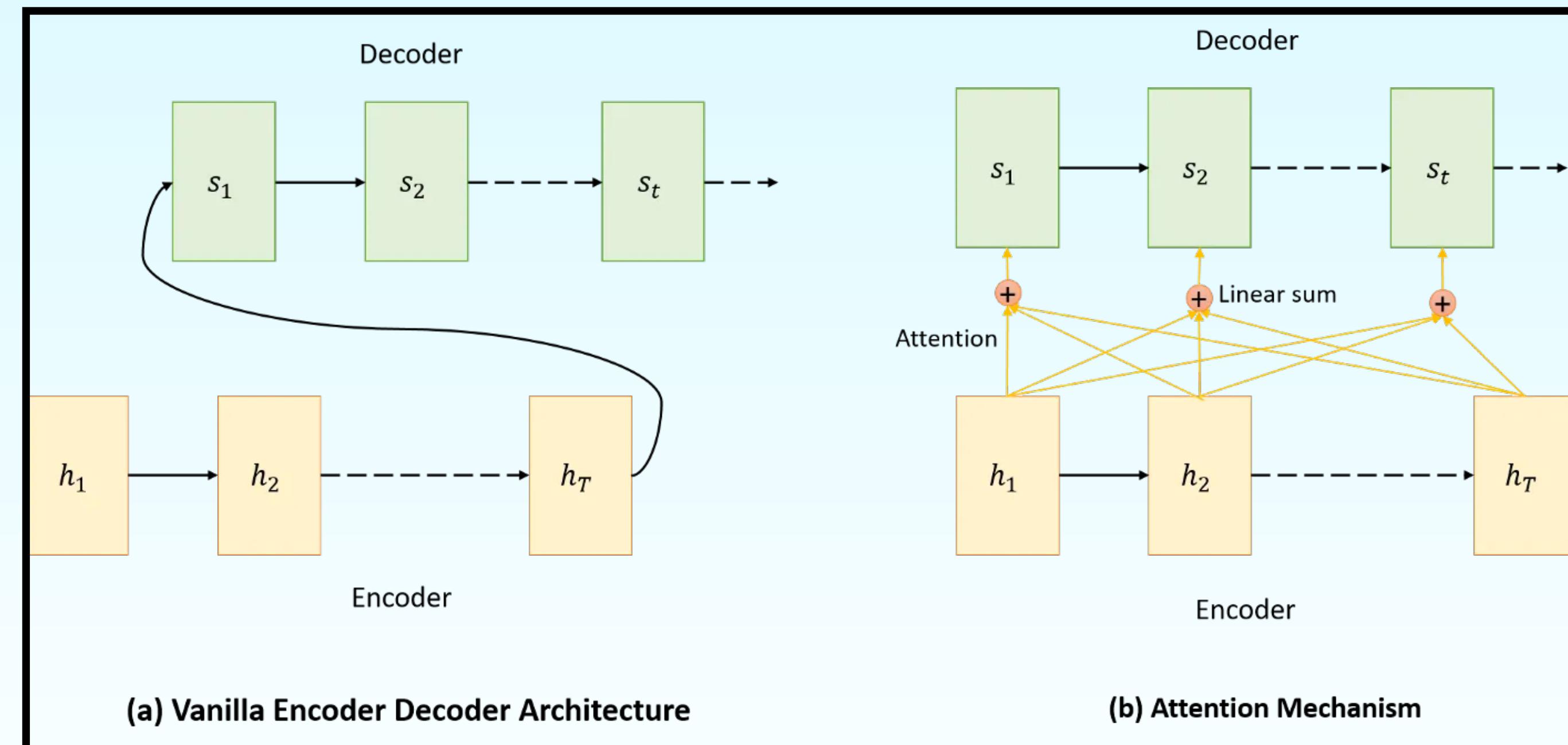
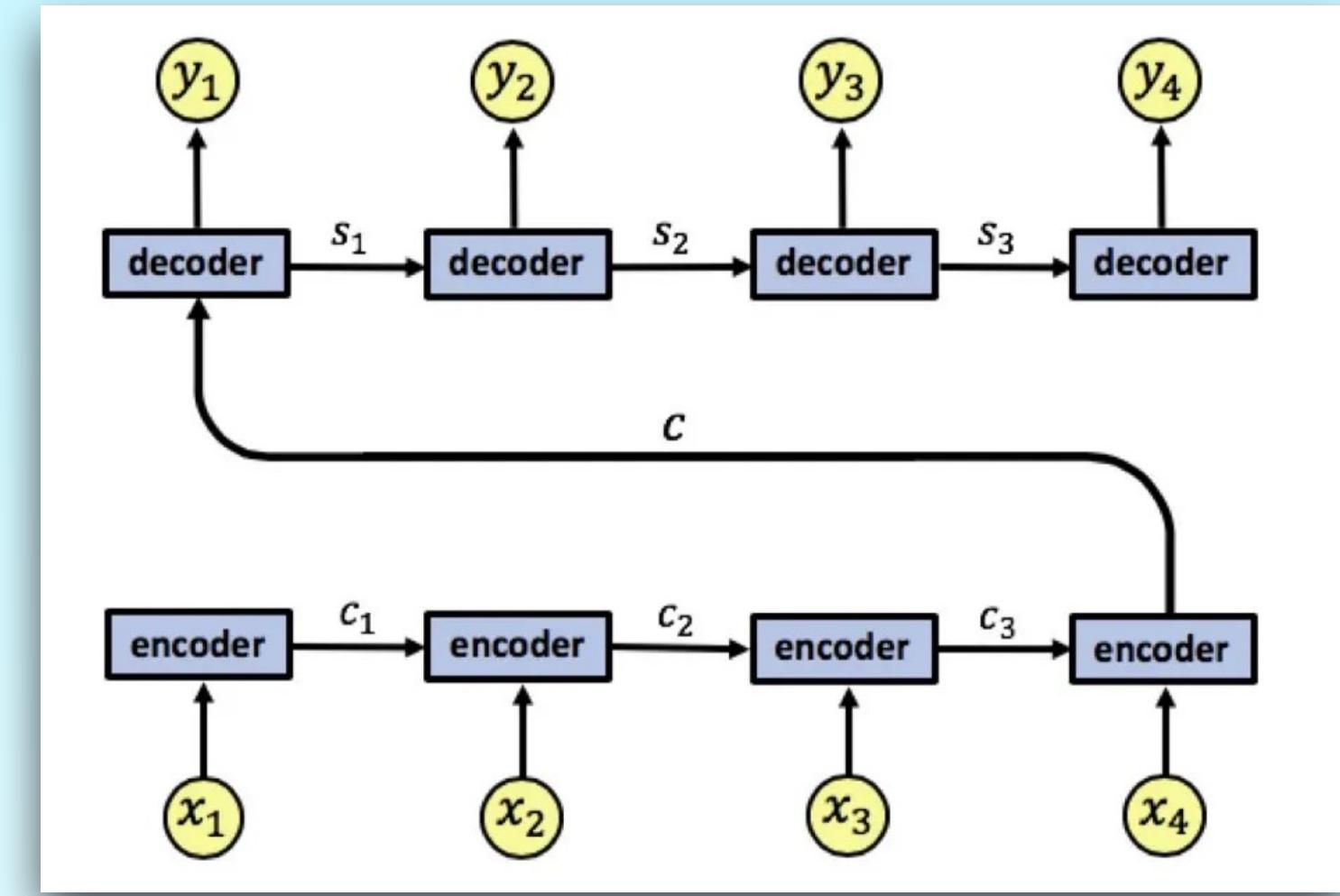
Question: which university/school is this?



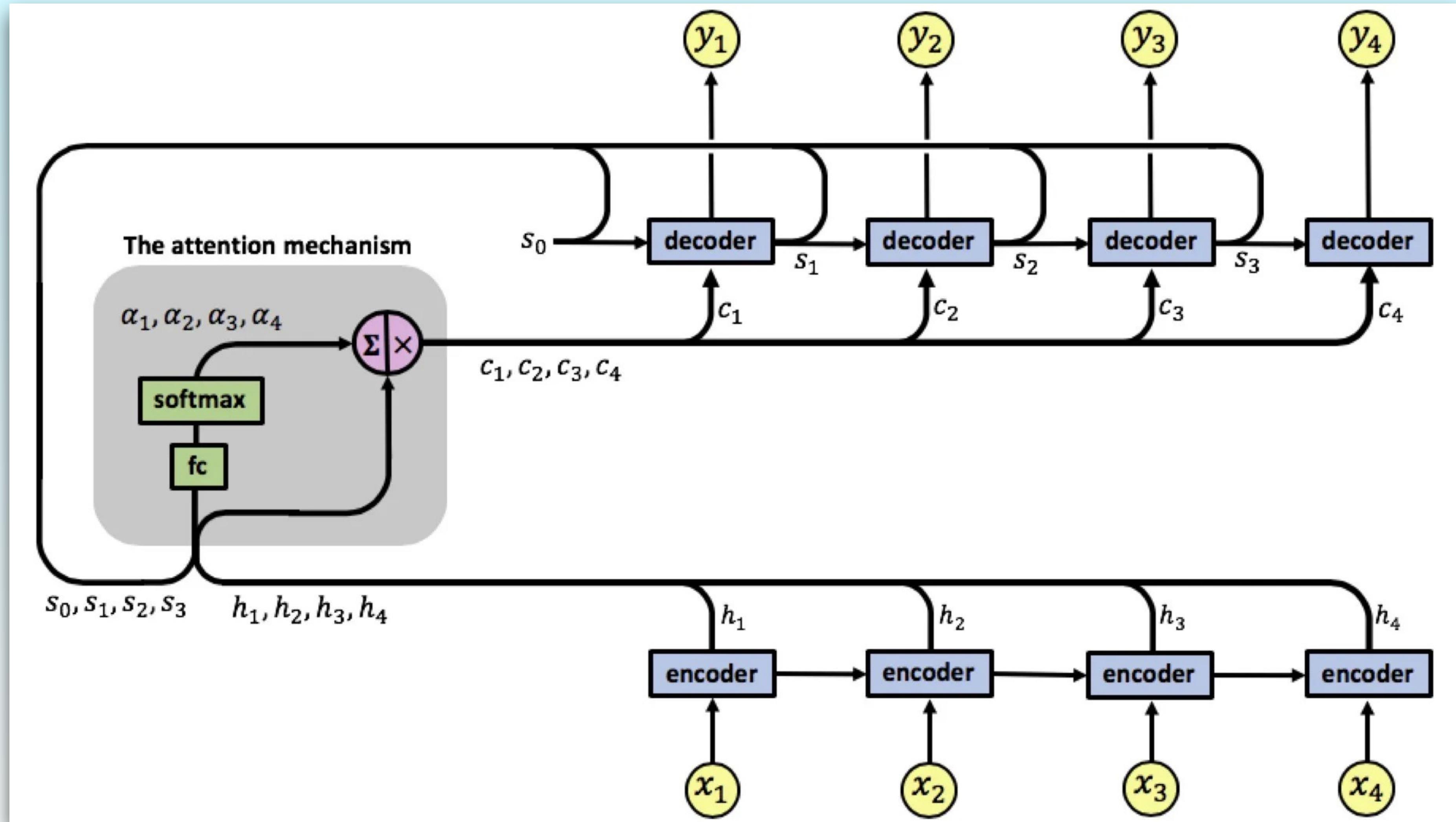
Attention Mechanism

- Until now, the final encoding steps had to represent the whole sequence into one vector

- Information loss + the whole decoding from a single vector
- Idea: add an addition to the architecture, helping it to *attend* to the different inputs in the sequence
- Looks at all the encoder hidden states



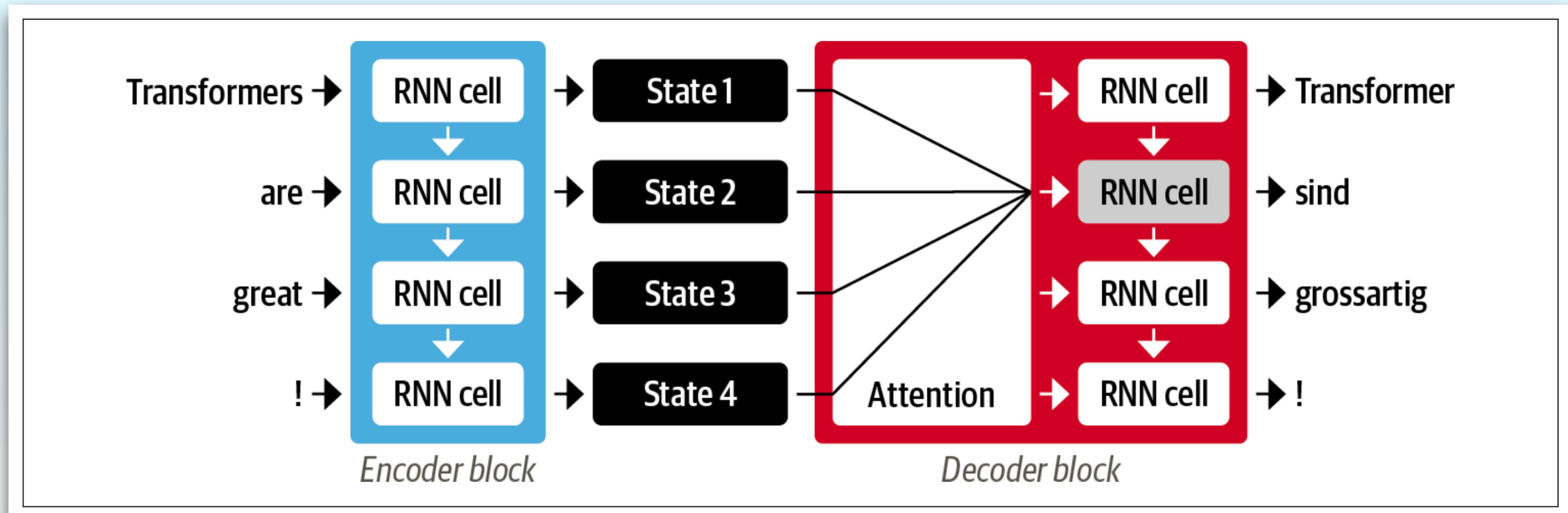
Attention Mechanism



Attention Mechanism

From Tunstall et al.:

“The main idea behind attention is that instead of producing a single hidden state for the input sequence, the encoder outputs a hidden state at each step that the decoder can access.”



Attention Mechanism

Learns non-trivial input alignments

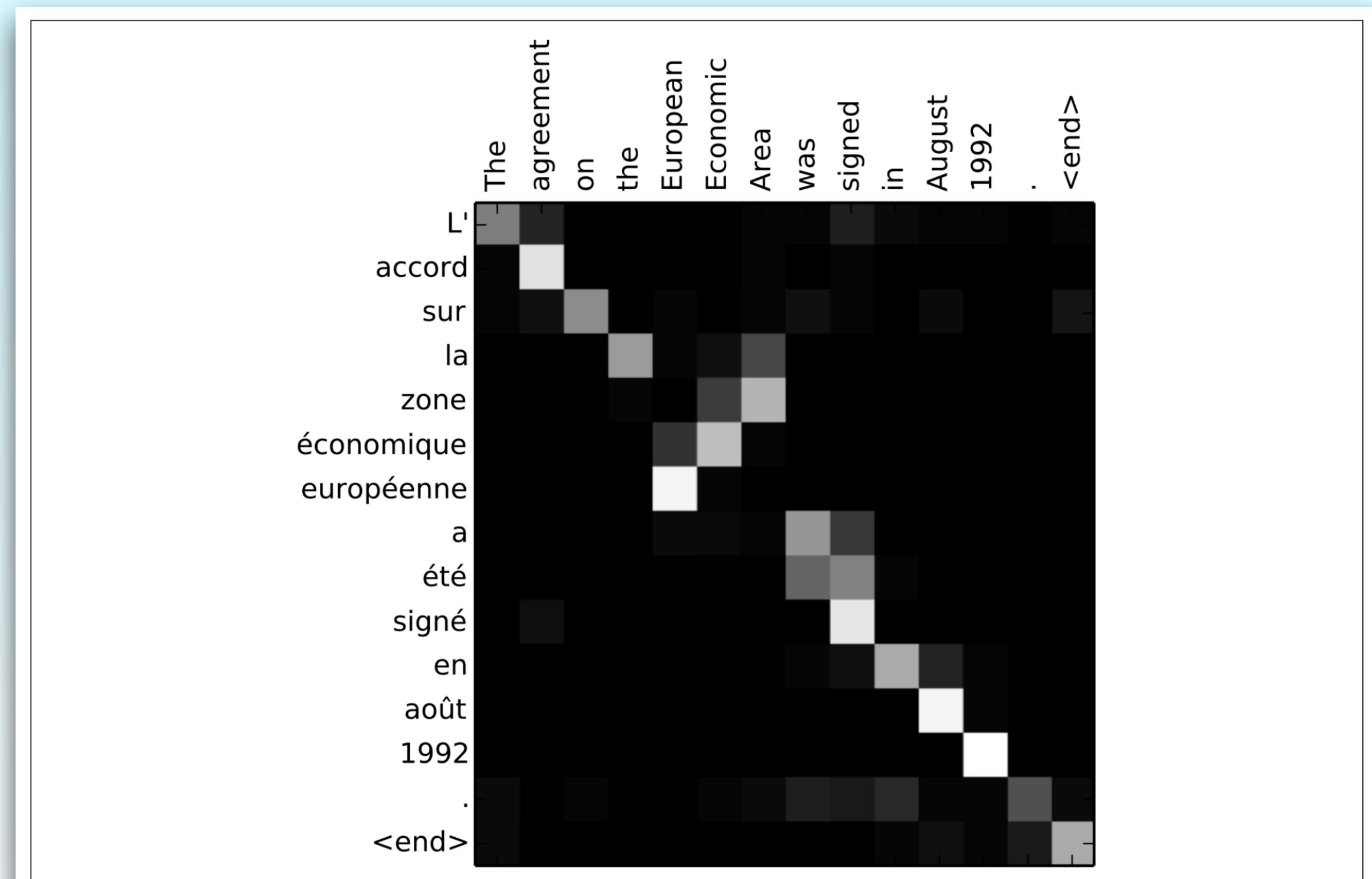
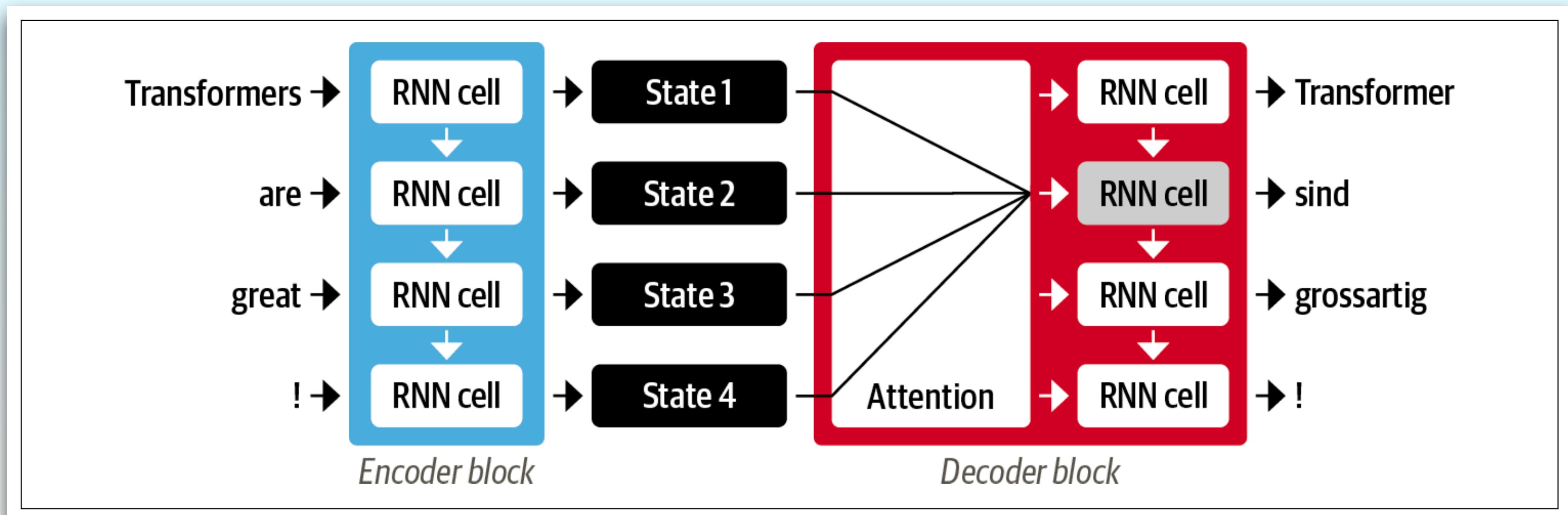


Figure 1-5. RNN encoder-decoder alignment of words in English and the generated translation in French (courtesy of Dzmitry Bahdanau)

Any downside that comes to mind?



Transformer

- Remove the recurrence! (Self-)attention is all you need!
- Encoder and Decoder (can be only one of them; see in NLP week!)
- Input can be text, image, ... (but has small details to take care of!)

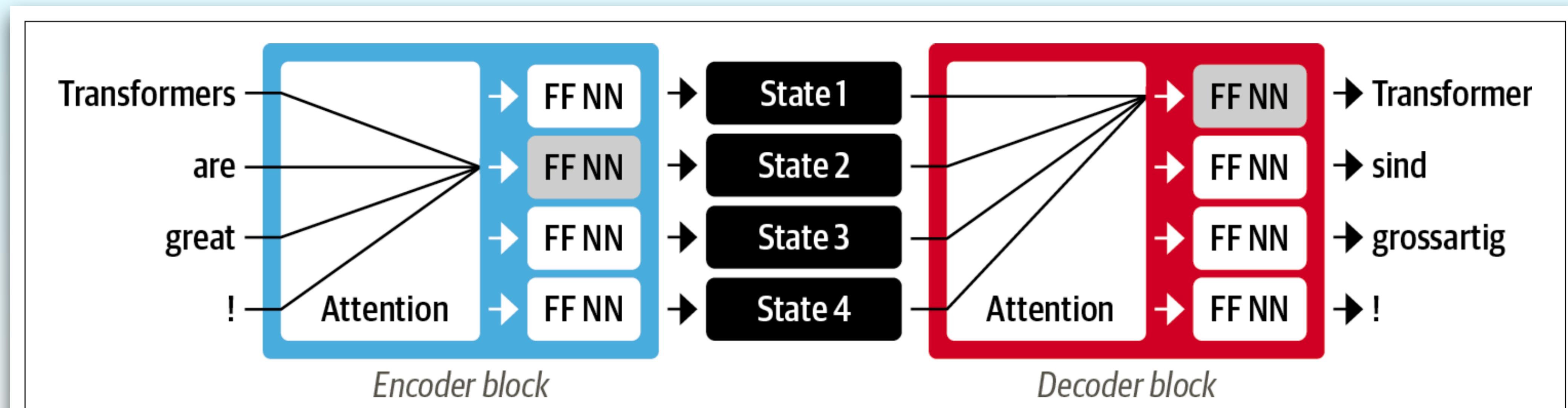
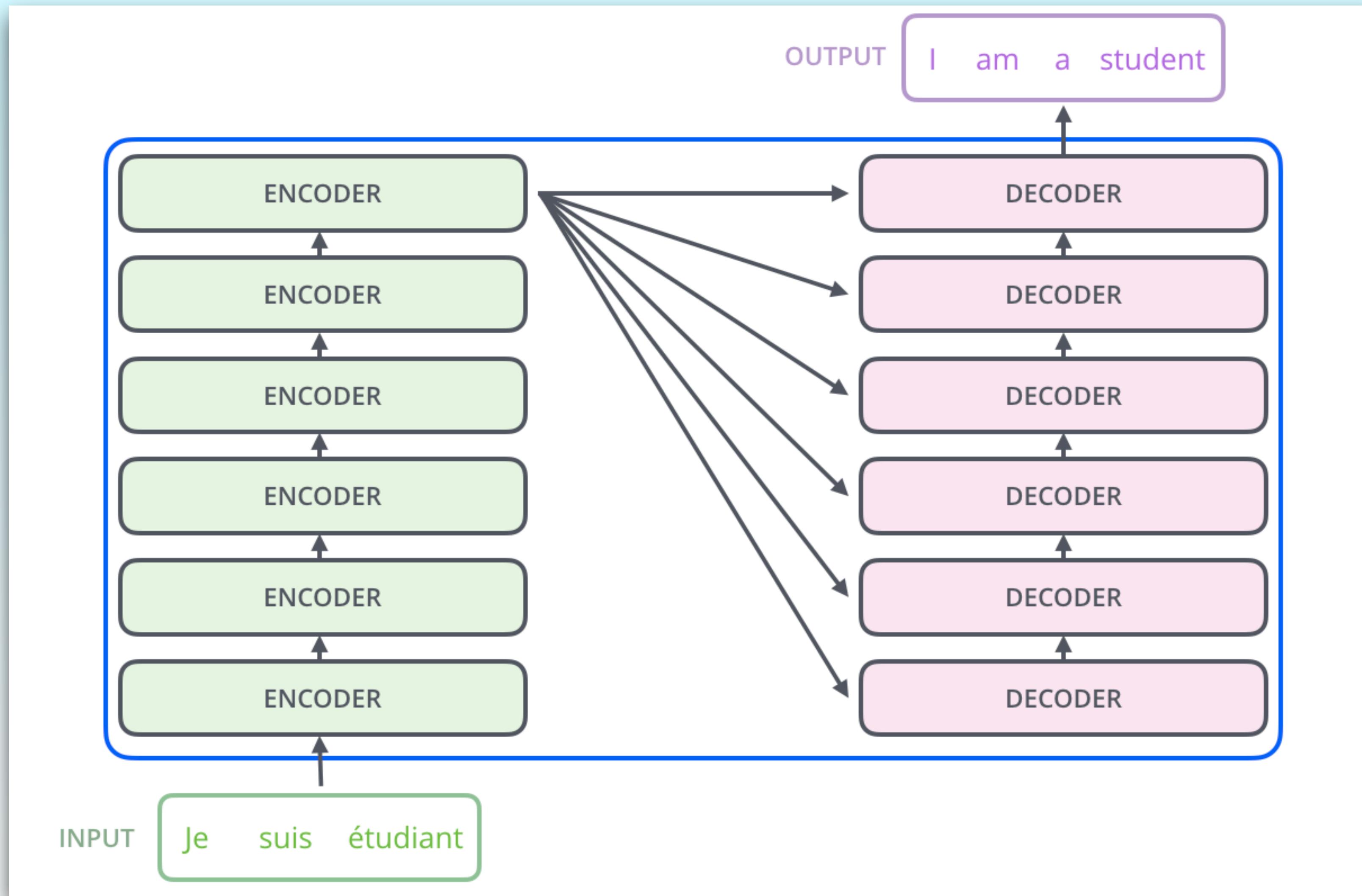


Figure 1-6. Encoder-decoder architecture of the original Transformer

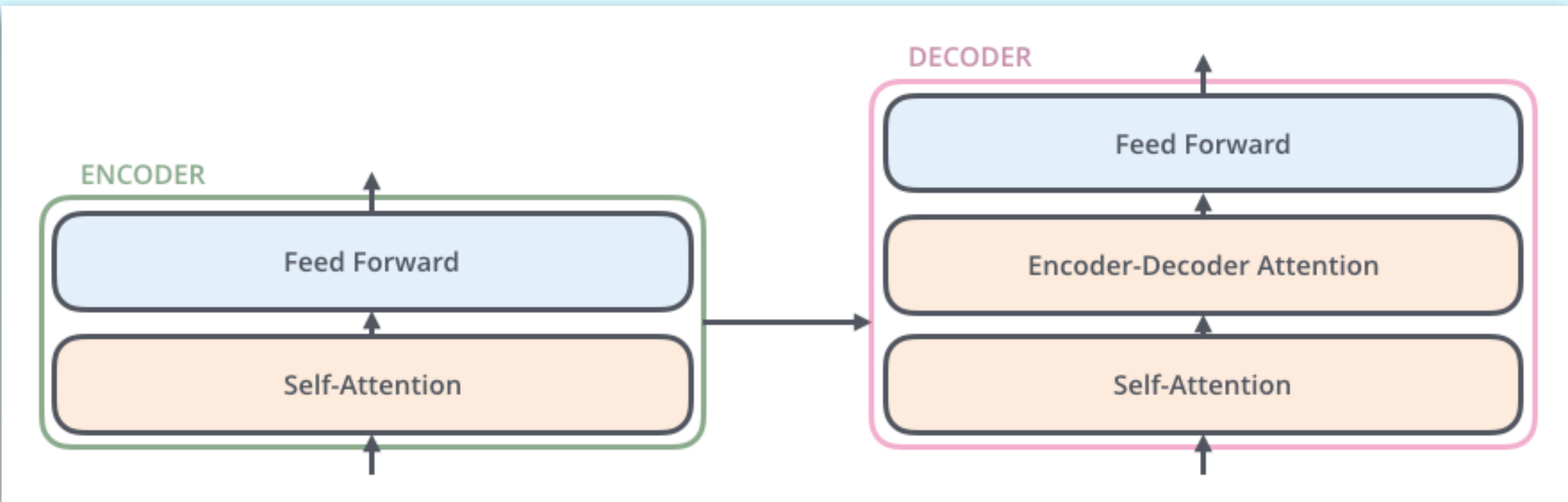
Let's explore how Transformers are designed!

Images and content in this section is from <https://jalammar.github.io/illustrated-transformer/>

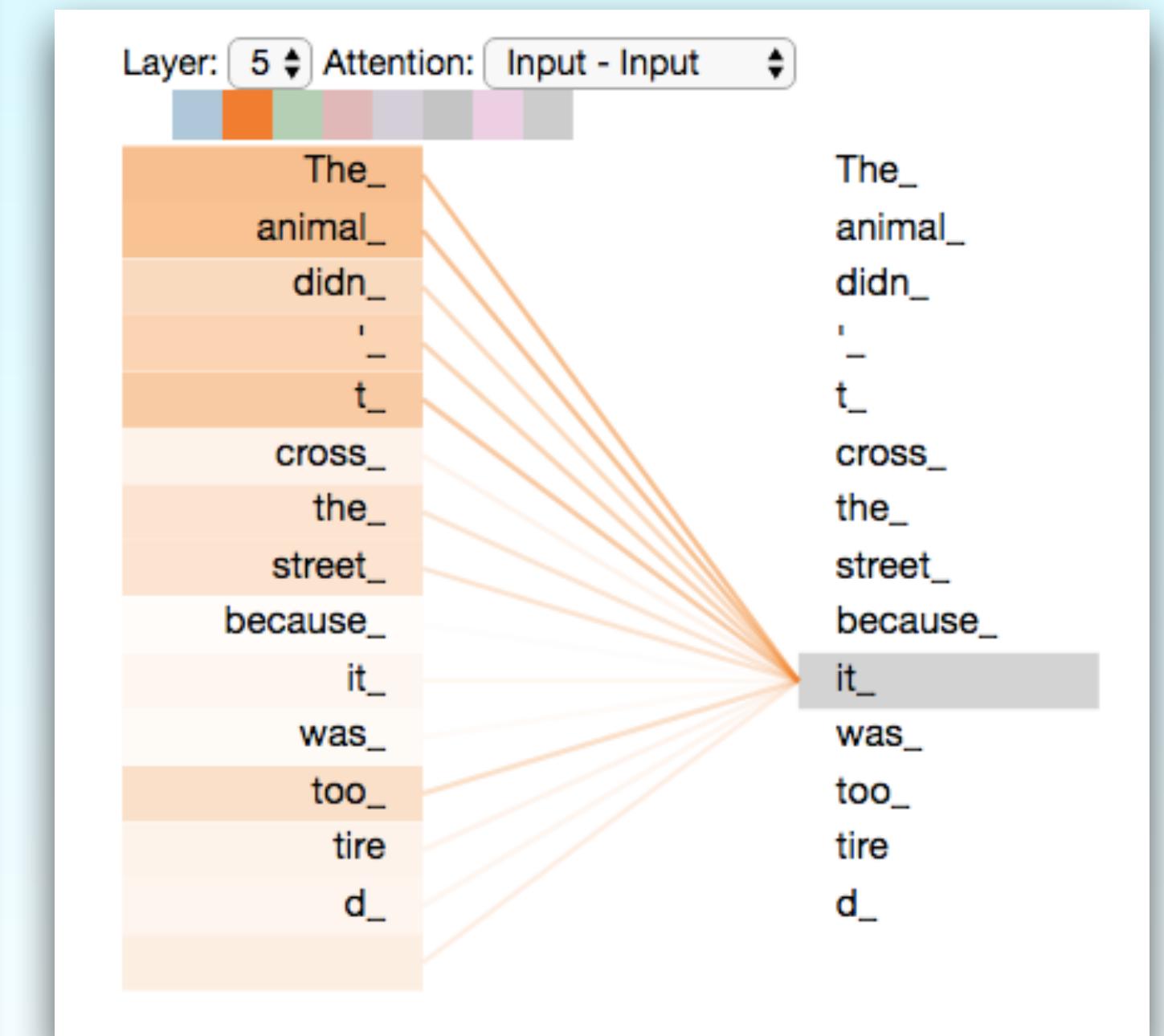
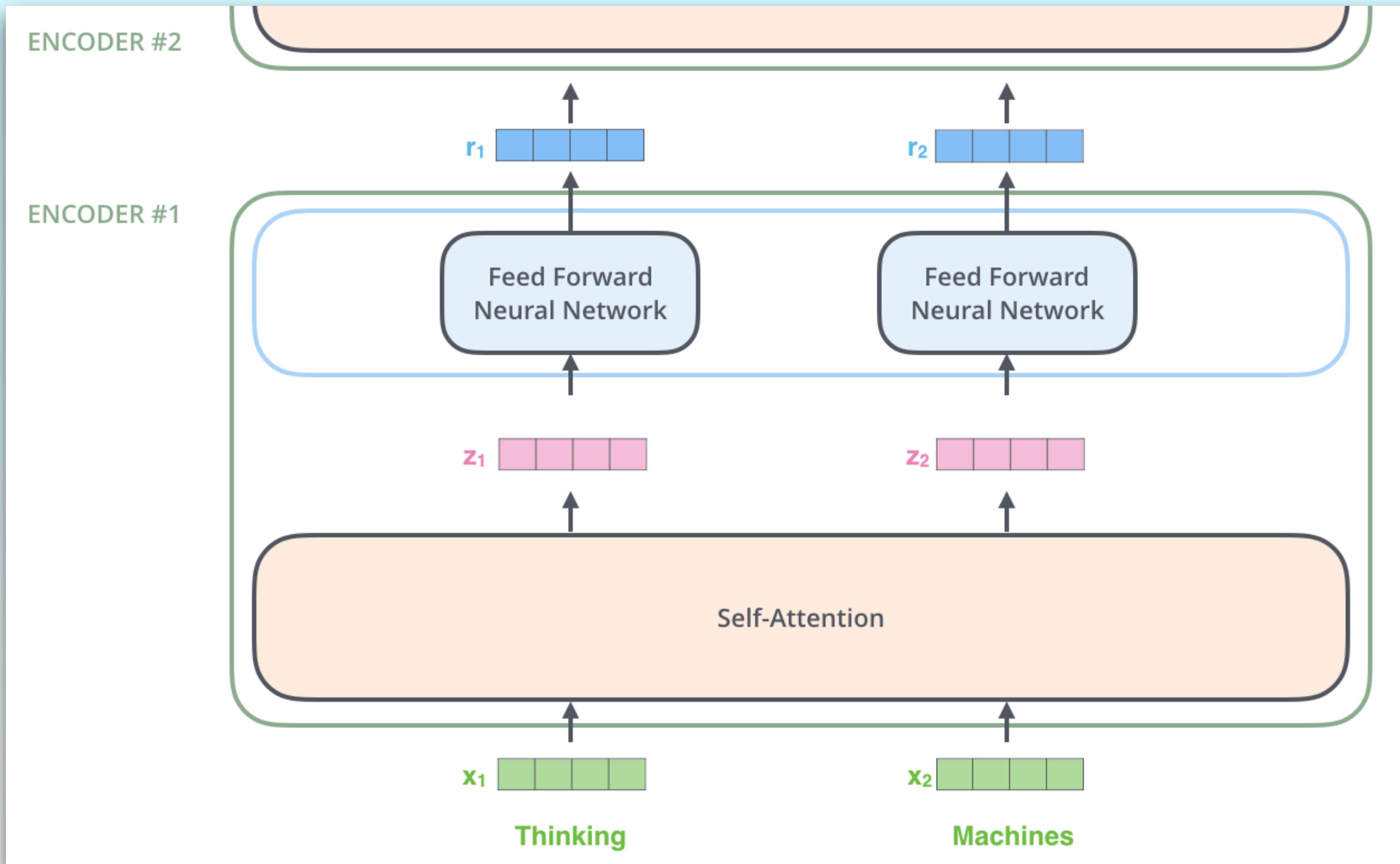
Transformer



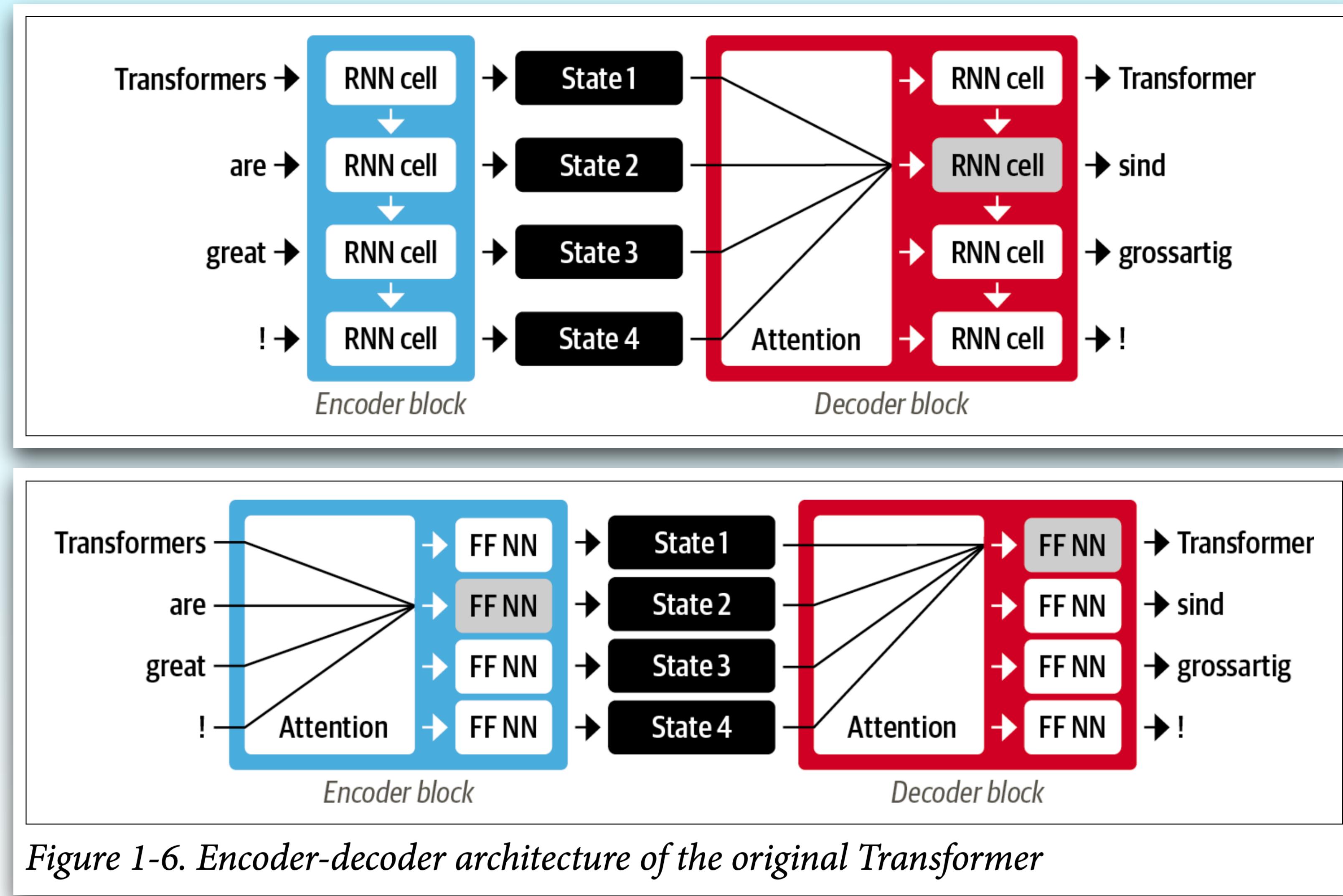
Transformer



Transformer Encoder



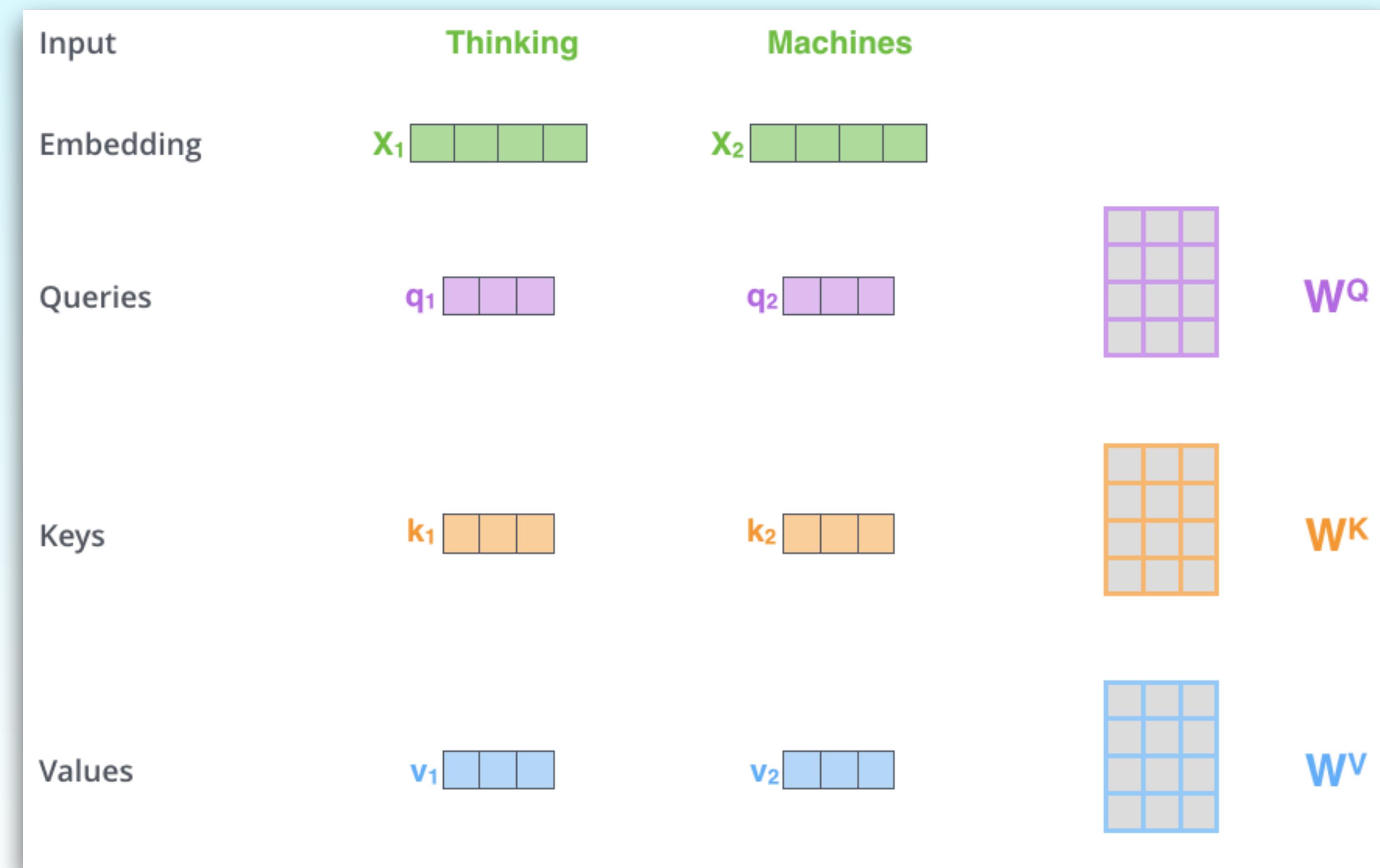
Self-attention – find the differences!



Now let's get a bit more technical

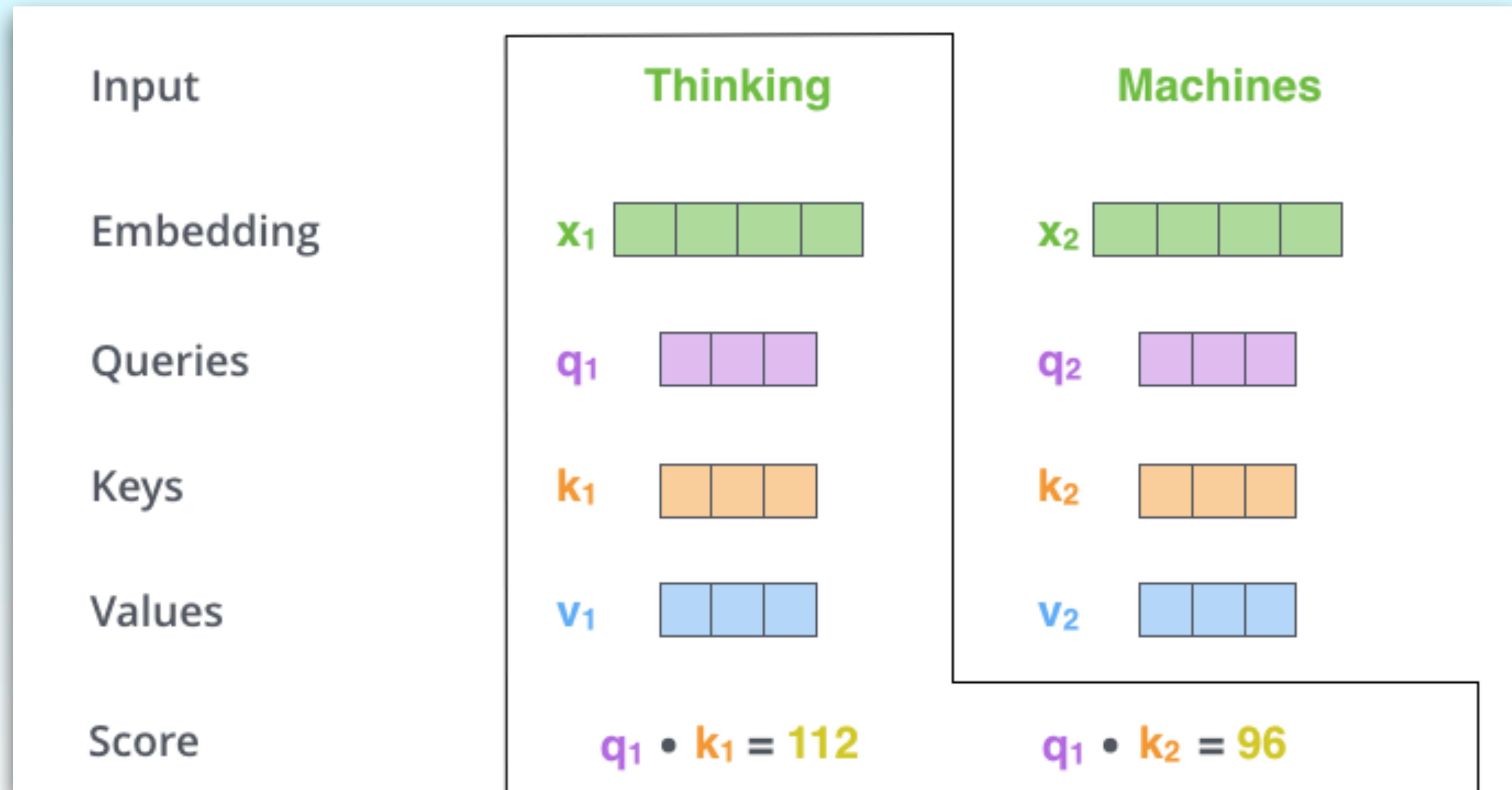
Self-attention Process

1) Create three vectors from each input of the encoder (can be smaller than input data!)



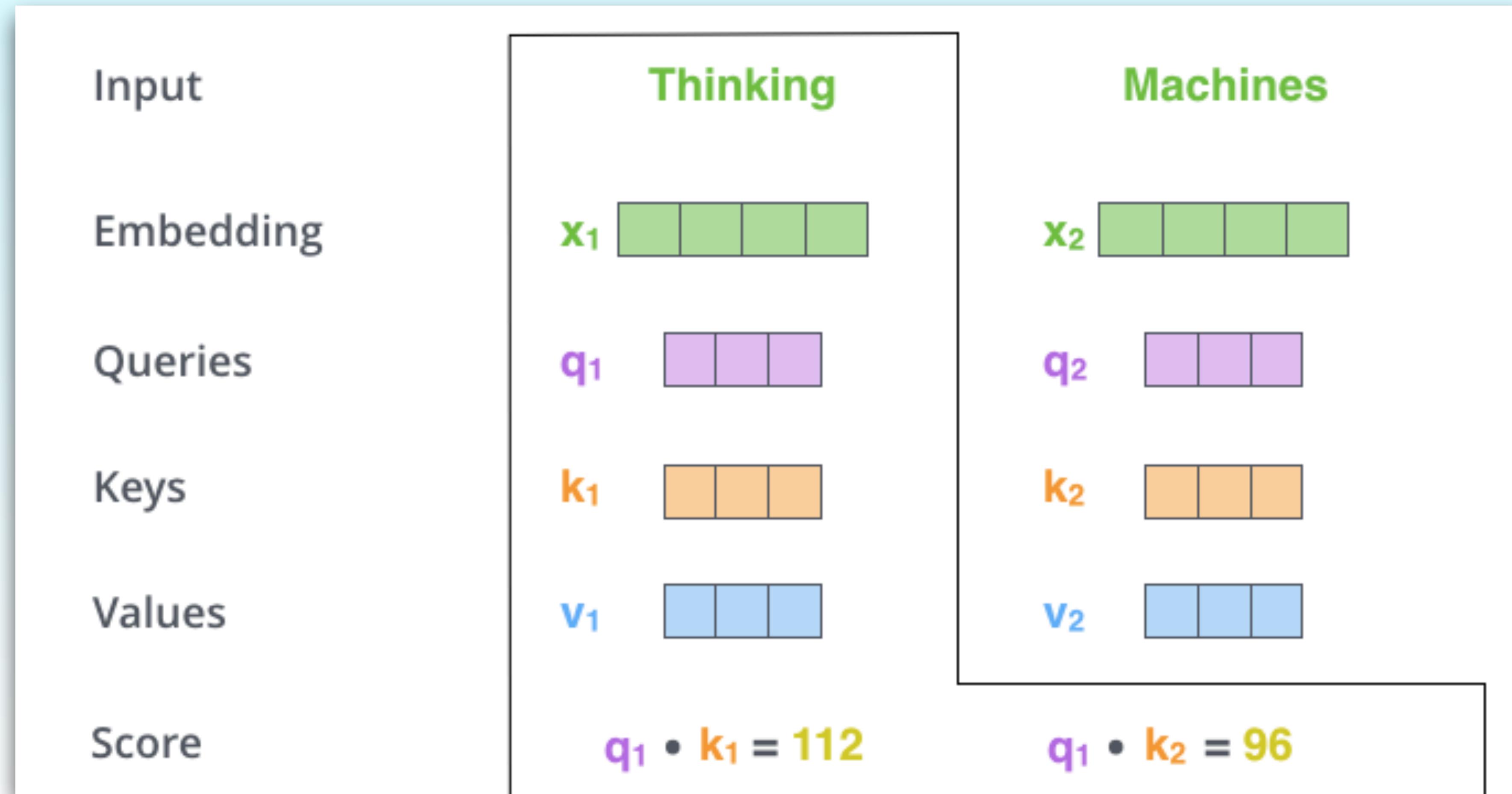
Self-attention Process

2) For each input, multiply its query with all the keys



Self-attention Process

3) Divide scores by sqrt of dim of key vectors (for stable gradients) and put into softmax

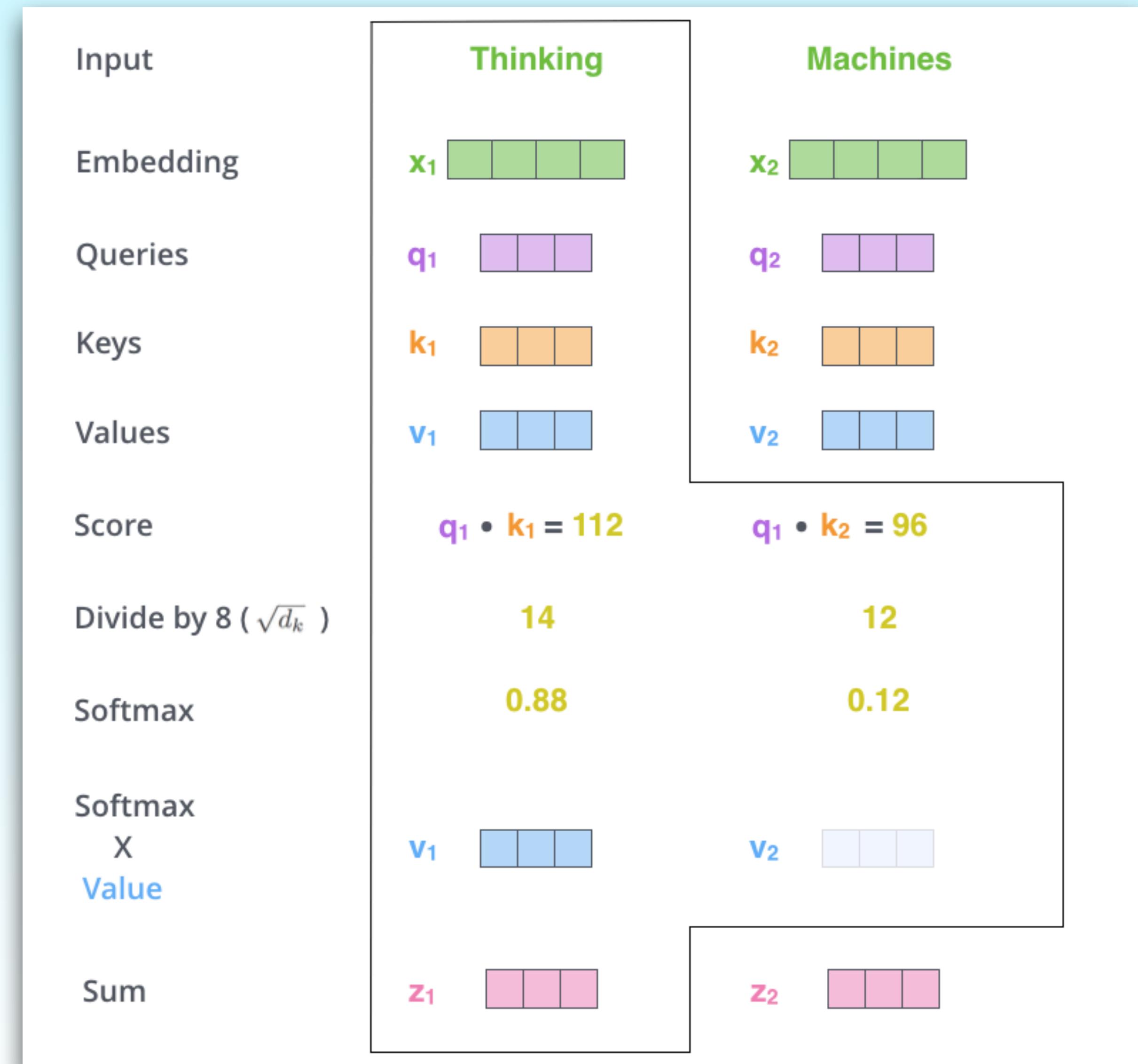


$$s(x_i) = \frac{e^{x_i}}{\sum_{j=1}^n e^{x_j}}$$

Self-attention Process

4) Multiply *value* vectors by softmax
(i.e., removing irrelevants), then sum up

Done!



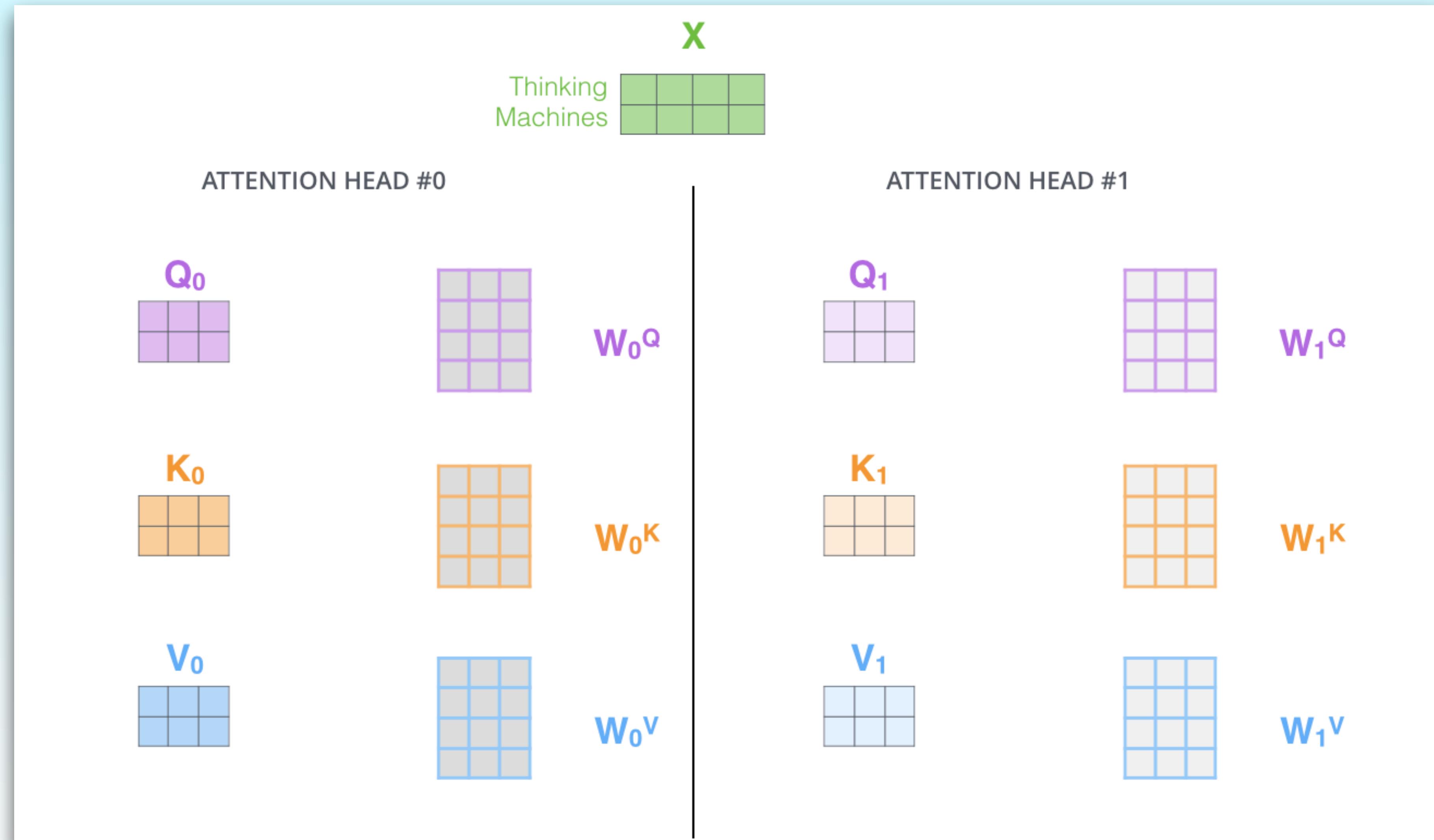
Self-attention In Practice: Matrices

All embeddings into X!

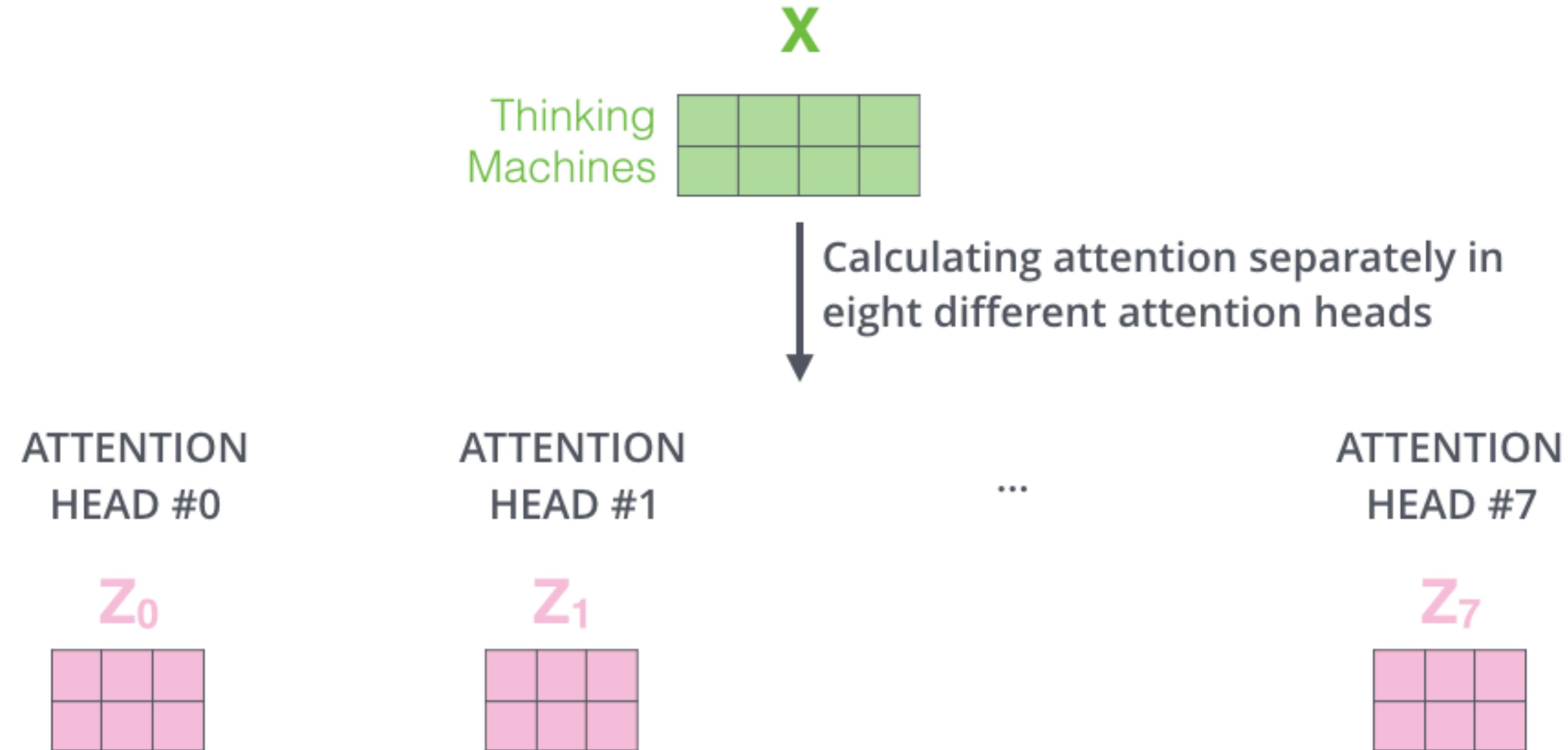
$$\begin{array}{c} \mathbf{X} \quad \mathbf{W}^Q \quad \mathbf{Q} \\ \begin{matrix} \text{---} \\ \text{---} \\ \text{---} \end{matrix} \\ \begin{matrix} | & | & | & | \\ | & | & | & | \\ | & | & | & | \end{matrix} \times \begin{matrix} | & | & | & | \\ | & | & | & | \\ | & | & | & | \end{matrix} = \begin{matrix} | & | & | & | \\ | & | & | & | \end{matrix} \end{array}$$
$$\begin{array}{c} \mathbf{X} \quad \mathbf{W}^K \quad \mathbf{K} \\ \begin{matrix} \text{---} \\ \text{---} \\ \text{---} \end{matrix} \\ \begin{matrix} | & | & | & | \\ | & | & | & | \\ | & | & | & | \end{matrix} \times \begin{matrix} | & | & | & | \\ | & | & | & | \\ | & | & | & | \end{matrix} = \begin{matrix} | & | & | & | \\ | & | & | & | \end{matrix} \end{array}$$
$$\begin{array}{c} \mathbf{X} \quad \mathbf{W}^V \quad \mathbf{V} \\ \begin{matrix} \text{---} \\ \text{---} \\ \text{---} \end{matrix} \\ \begin{matrix} | & | & | & | \\ | & | & | & | \\ | & | & | & | \end{matrix} \times \begin{matrix} | & | & | & | \\ | & | & | & | \\ | & | & | & | \end{matrix} = \begin{matrix} | & | & | & | \\ | & | & | & | \end{matrix} \end{array}$$

$$\begin{aligned} & \text{softmax} \left(\frac{\mathbf{Q} \times \mathbf{K}^T}{\sqrt{d_k}} \right) \mathbf{V} \\ &= \mathbf{z} \end{aligned}$$

Can be multi-headed



Can be multi-headed



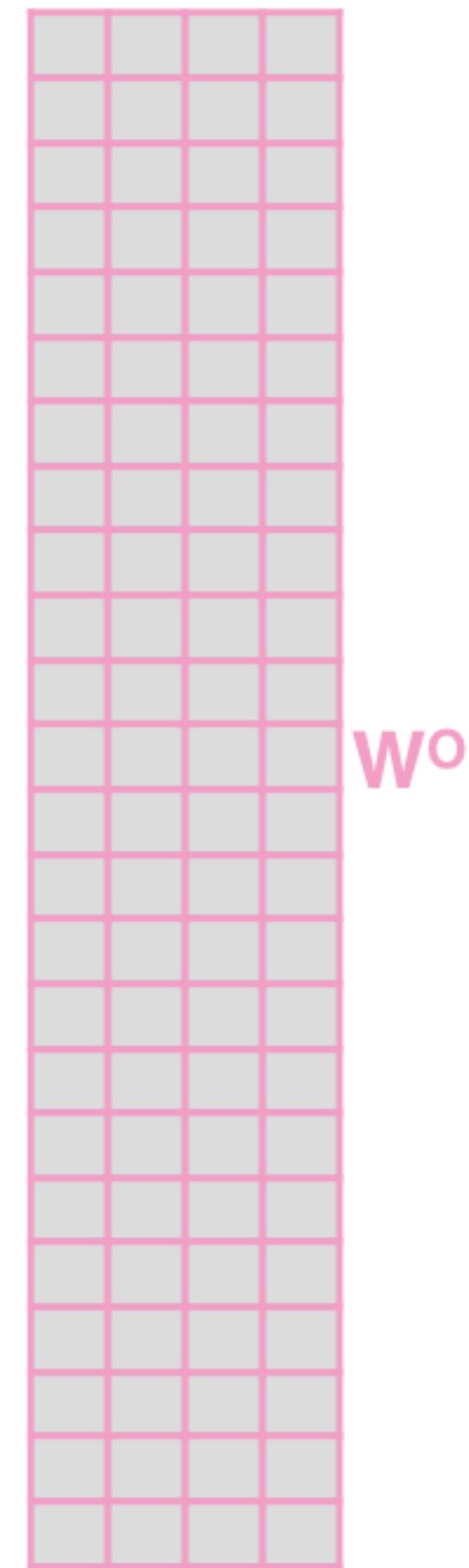
Can be multi-headed

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

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

What is one problem with this approach?

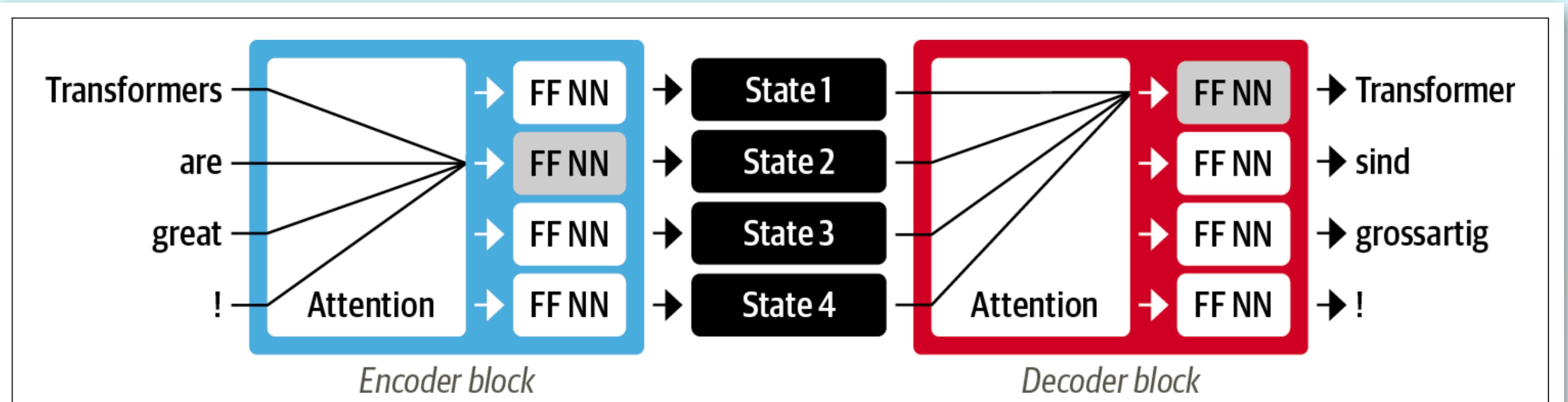
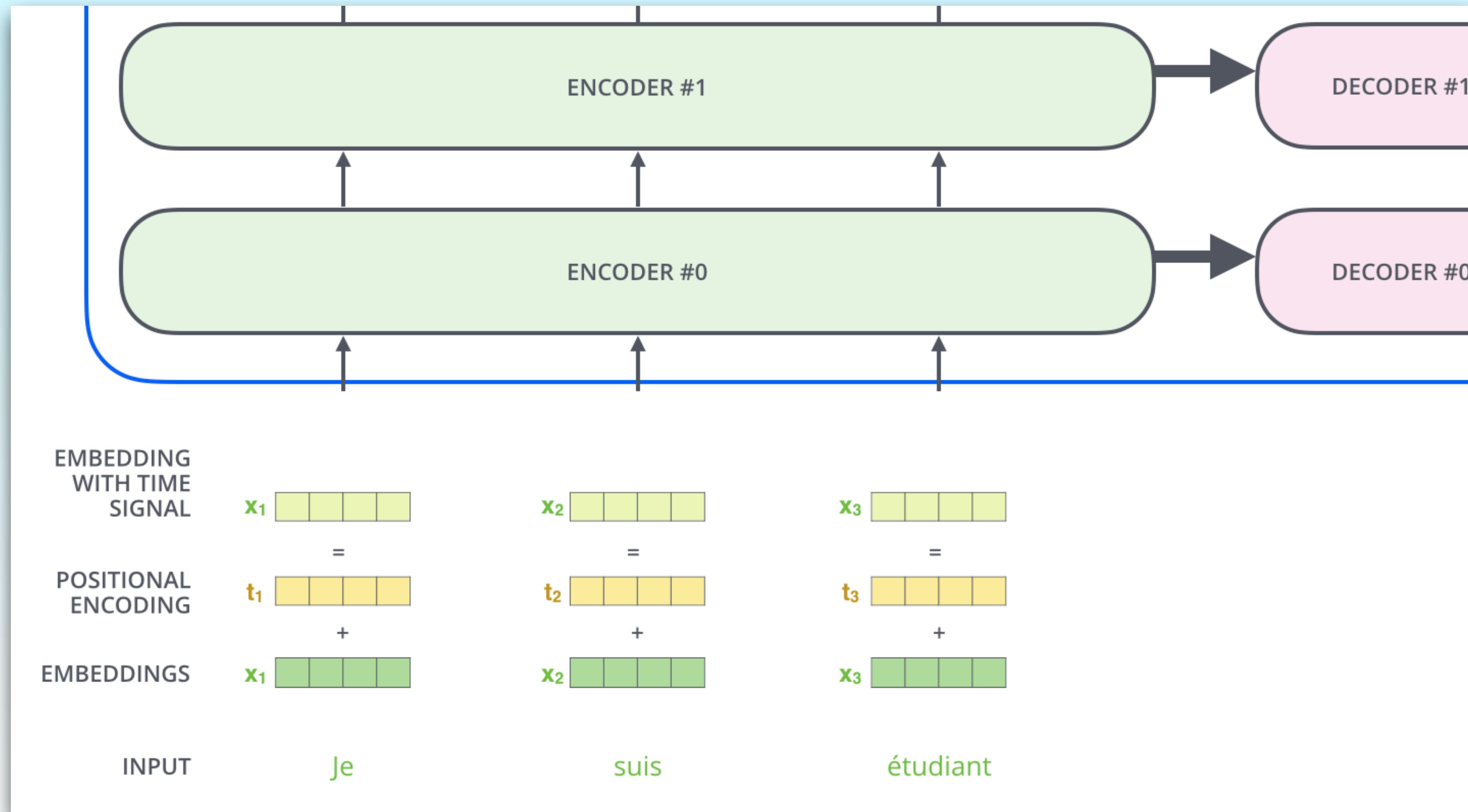
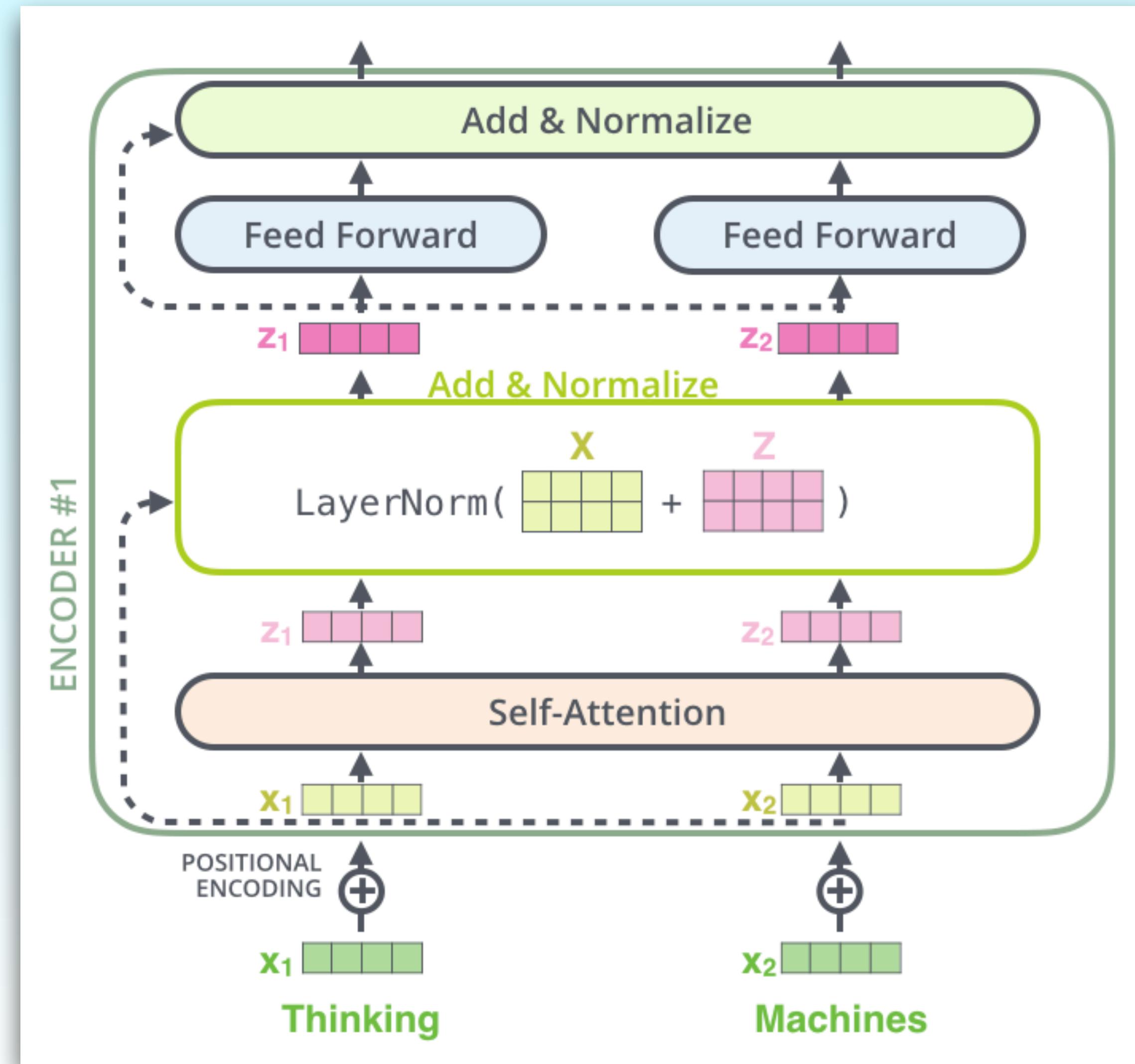


Figure 1-6. Encoder-decoder architecture of the original Transformer

Positional Encoding

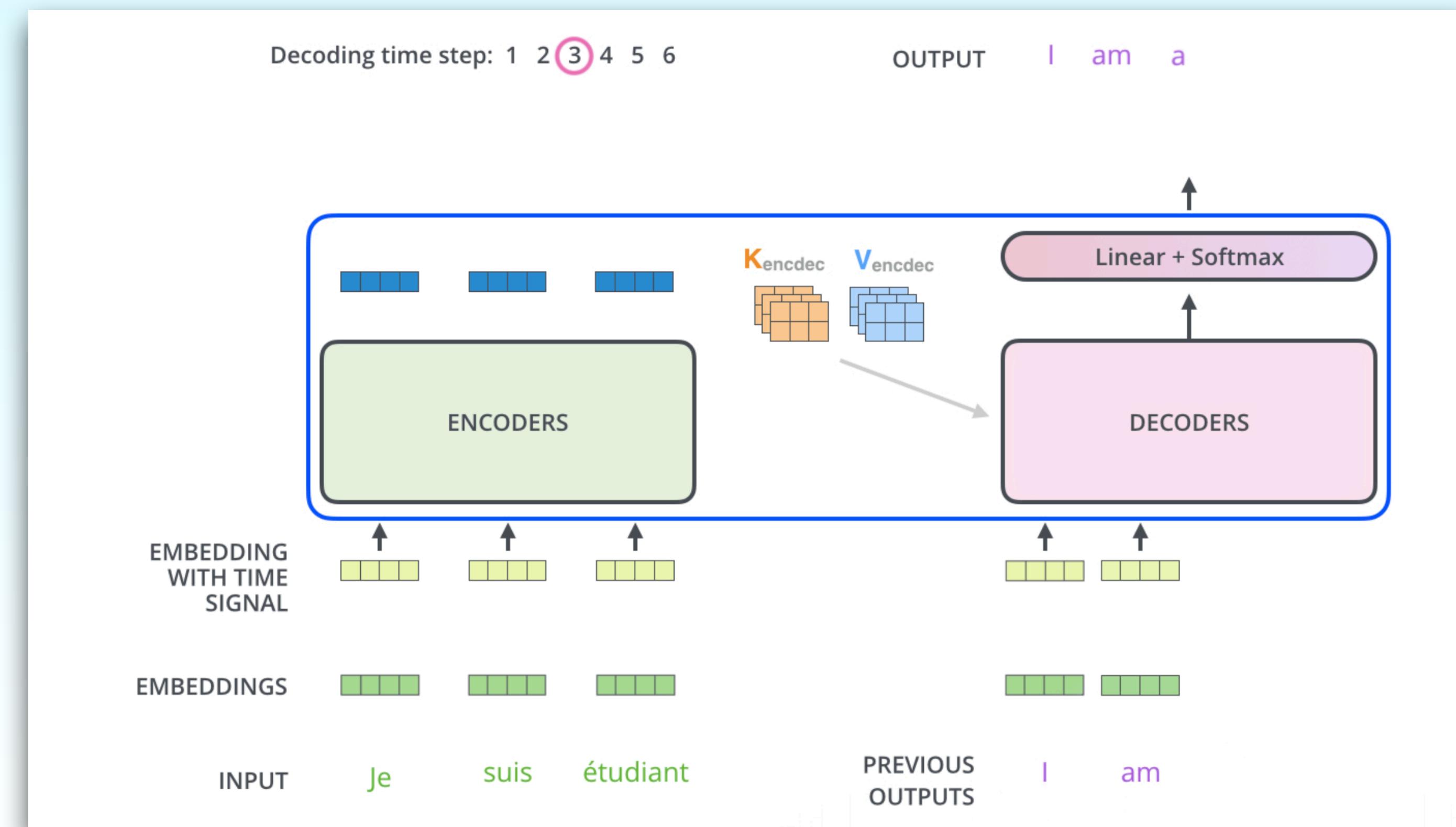


Full Encoder Together!

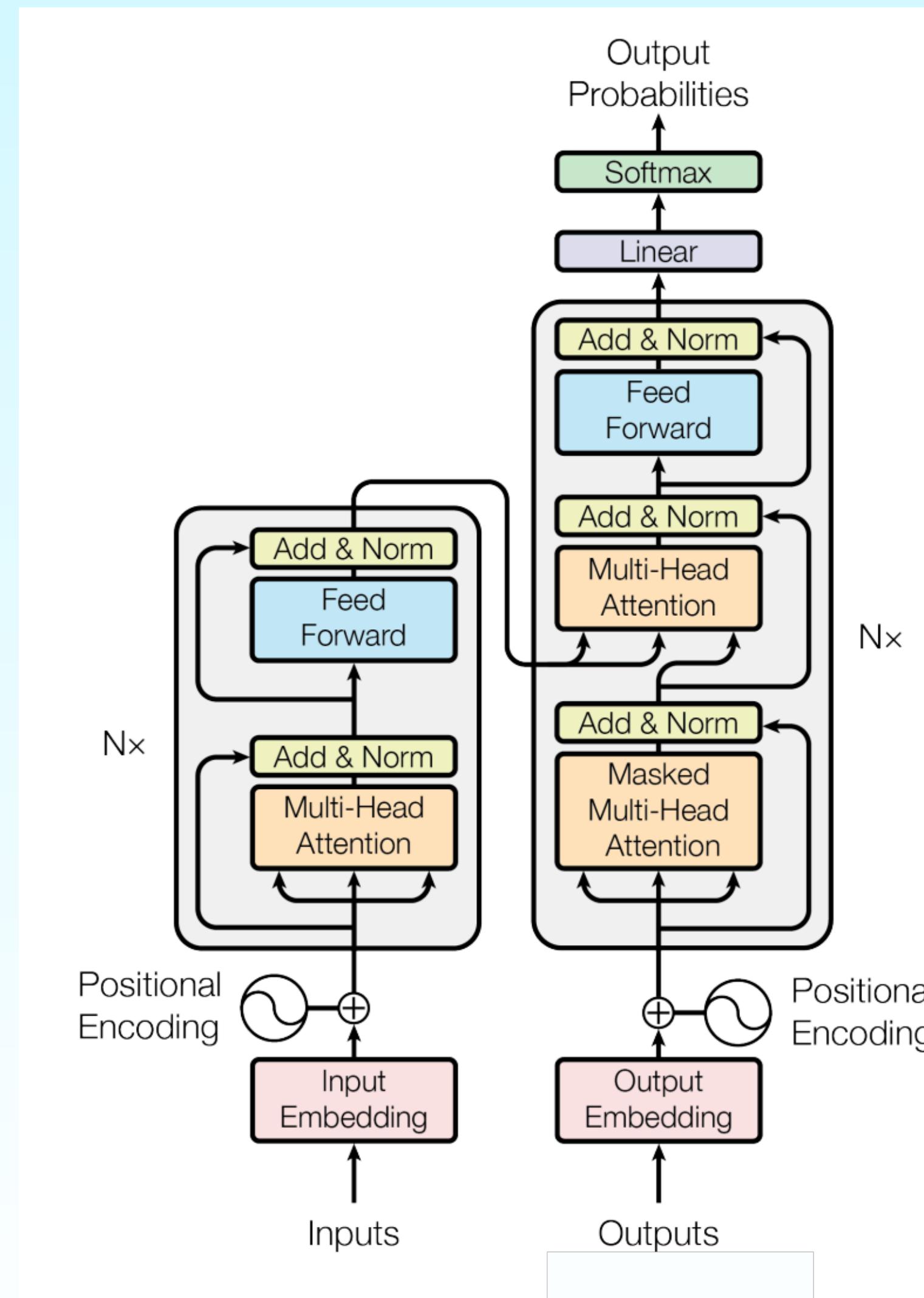


Transformer Decoder

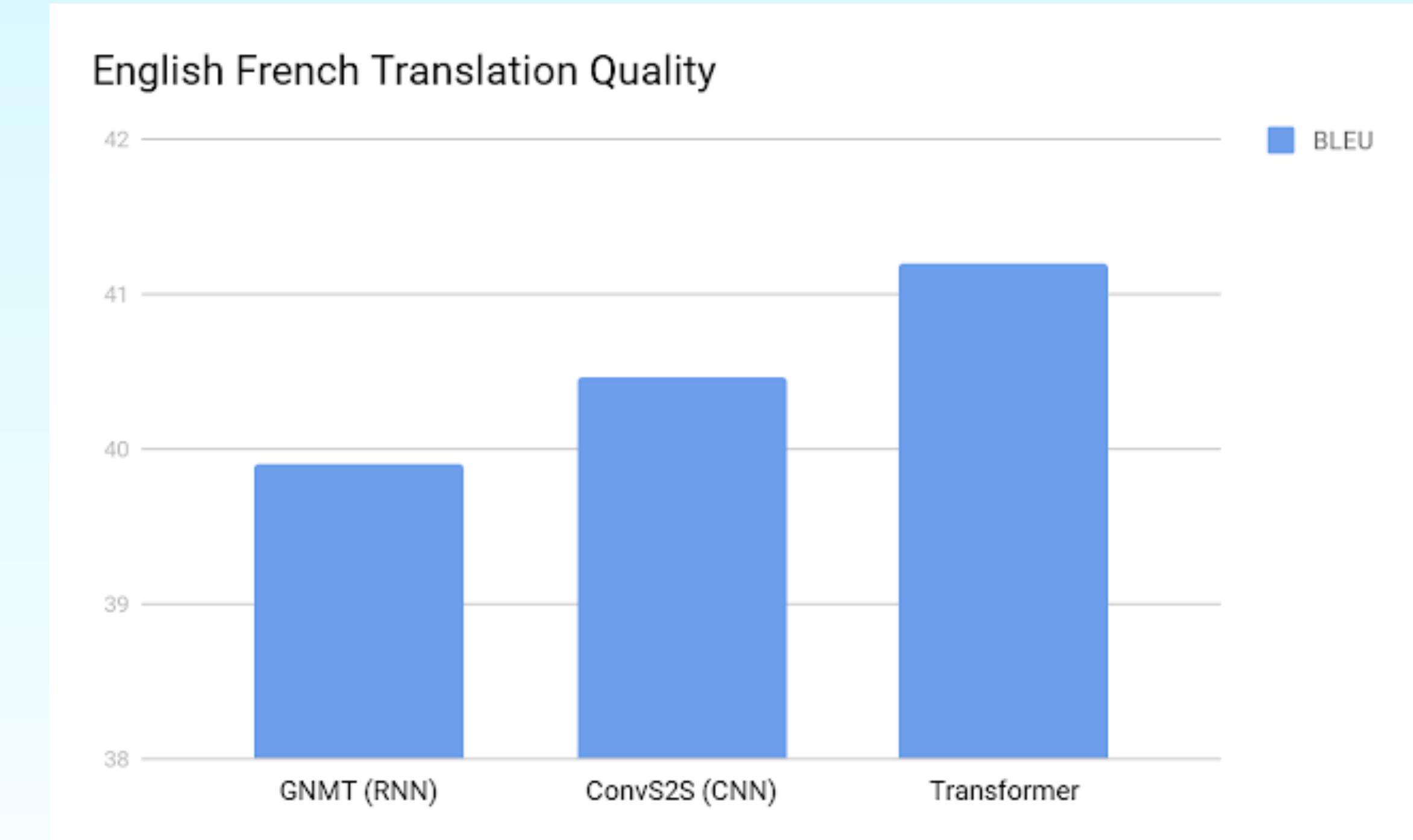
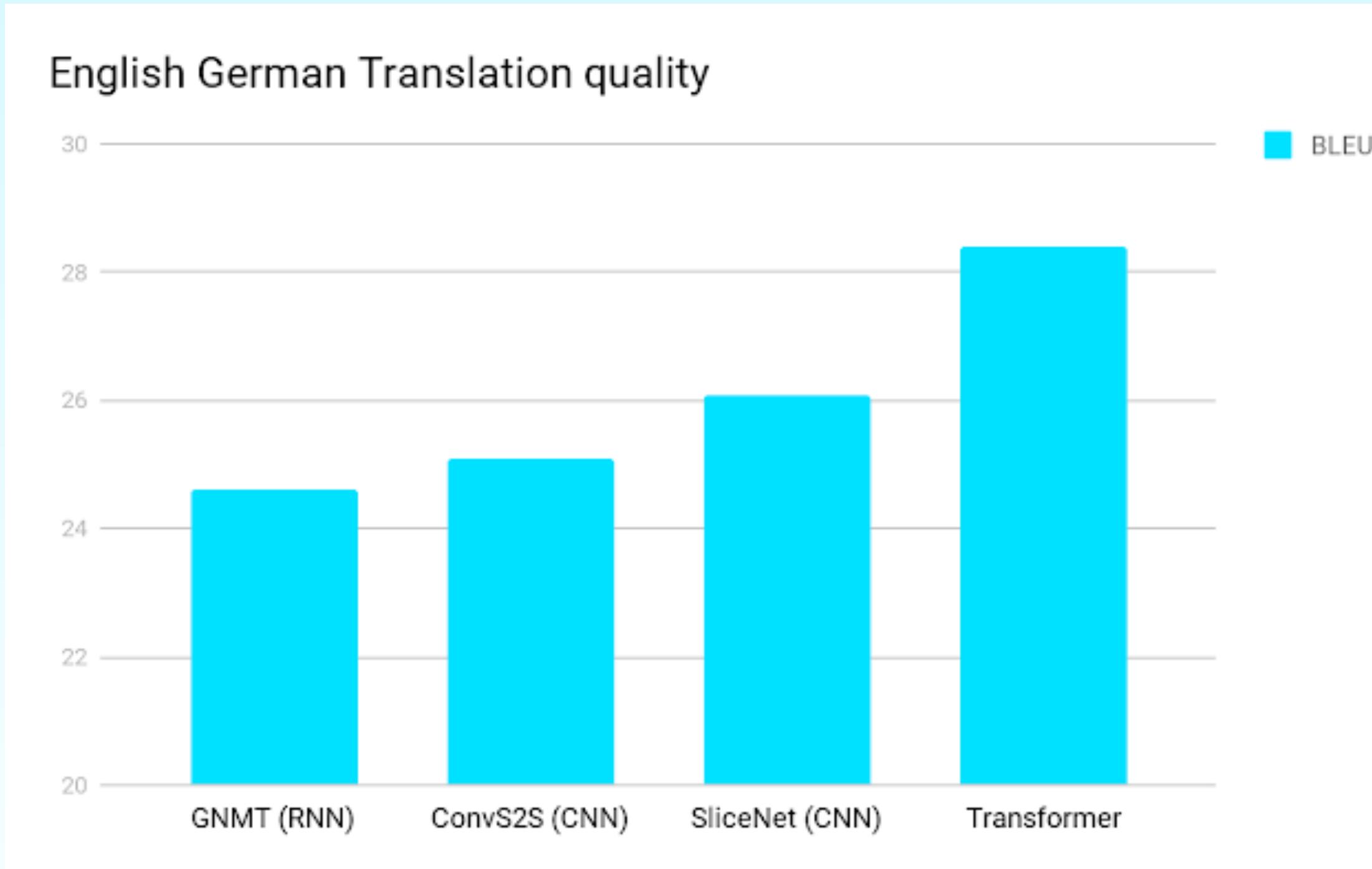
The output of the last decoder turns into K and V – *masked self-attention + encoder-decoder attention (keys and values from encoder outputs)* -> each cell one word, then softmax



Transformer in one view!



What did they achieve?



Let's see a code together!