

Applied Text Analytics

(F20/F21AA)

Introduction to Transformers

Heriot Watt University , Dubai Campus

Dr. Radu Mihailescu

Content

- Building blocks of a Transformer
- Encoder-Decoder Architecture
- Attention
 - Self-attention
 - Multi-headed attention
 - Encoder-Decoder attention



Transformers

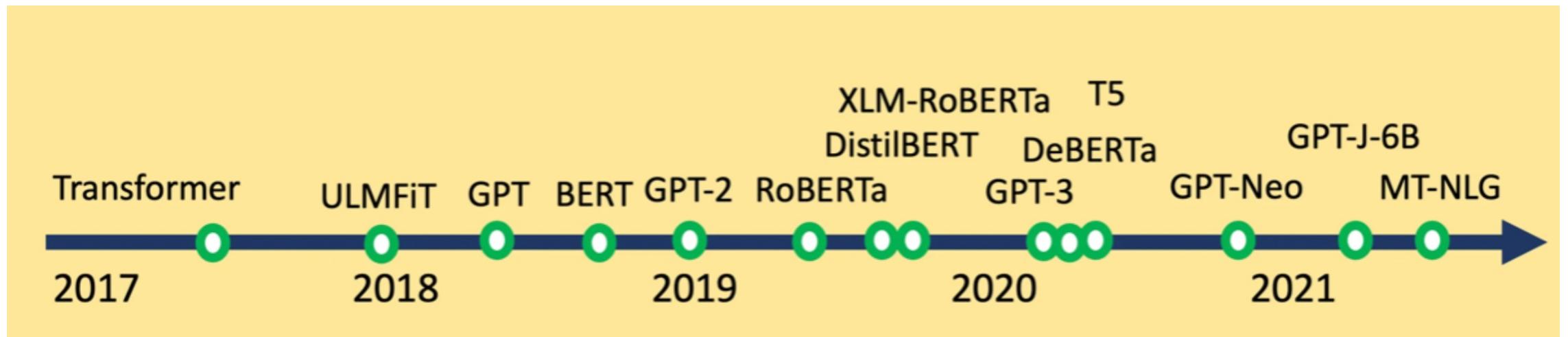
Neural network architecture for sequence modeling

- Based on work A. Vaswani et al., “**Attention Is All You Need**”, (2017)
<https://arxiv.org/pdf/1706.03762.pdf>
- Outperformed recurrent neural networks (RNNs) on machine translation tasks, both in terms of translation quality and training cost.

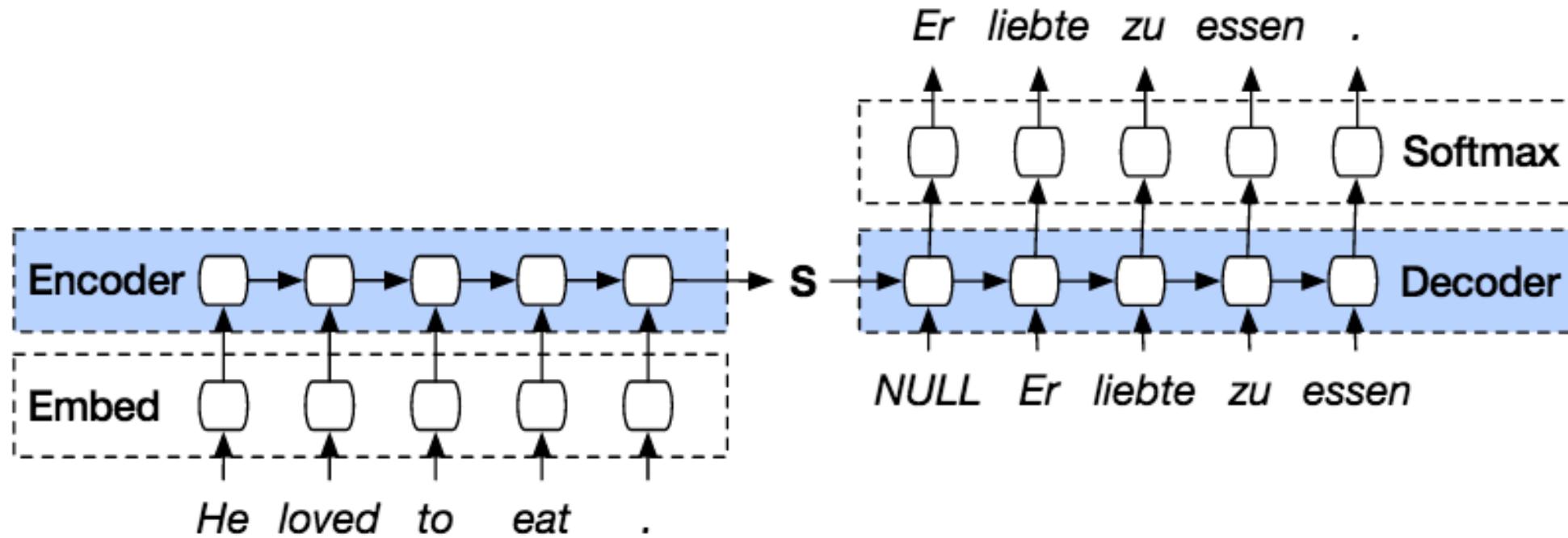
Today's most well-known transformers:

- Generative Pretrained Transformer (GPT)
- Bidirectional Encoder Representations from Transformers (BERT).

Transformers Timeline

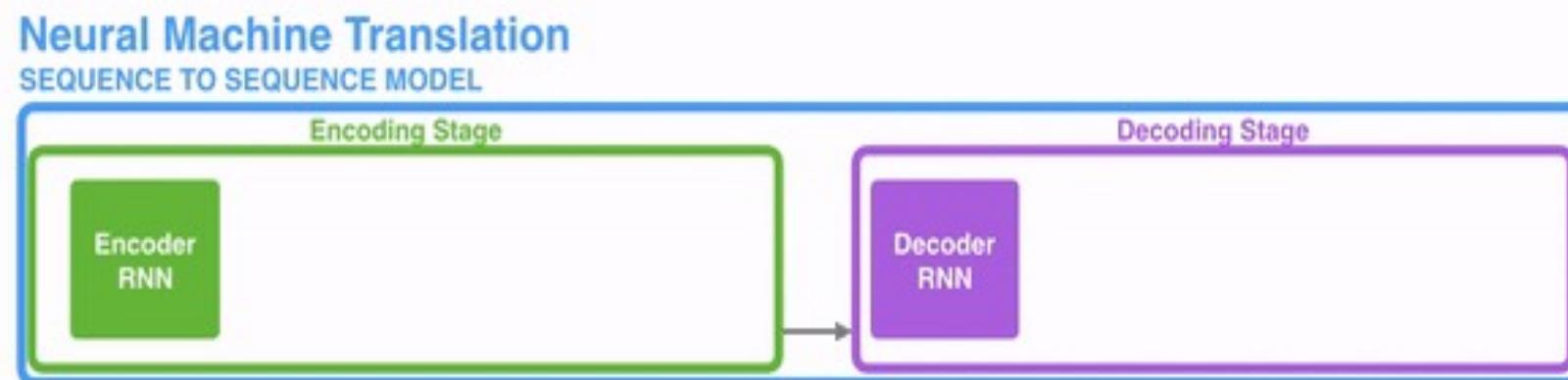


RNN Encoder-Decoder: Language Translation



The idea is to use one LSTM to read the input sequence, *one timestep at a time*, to obtain large fixed-dimensional vector representation (S), and then to use another LSTM to extract the output sequence from that vector

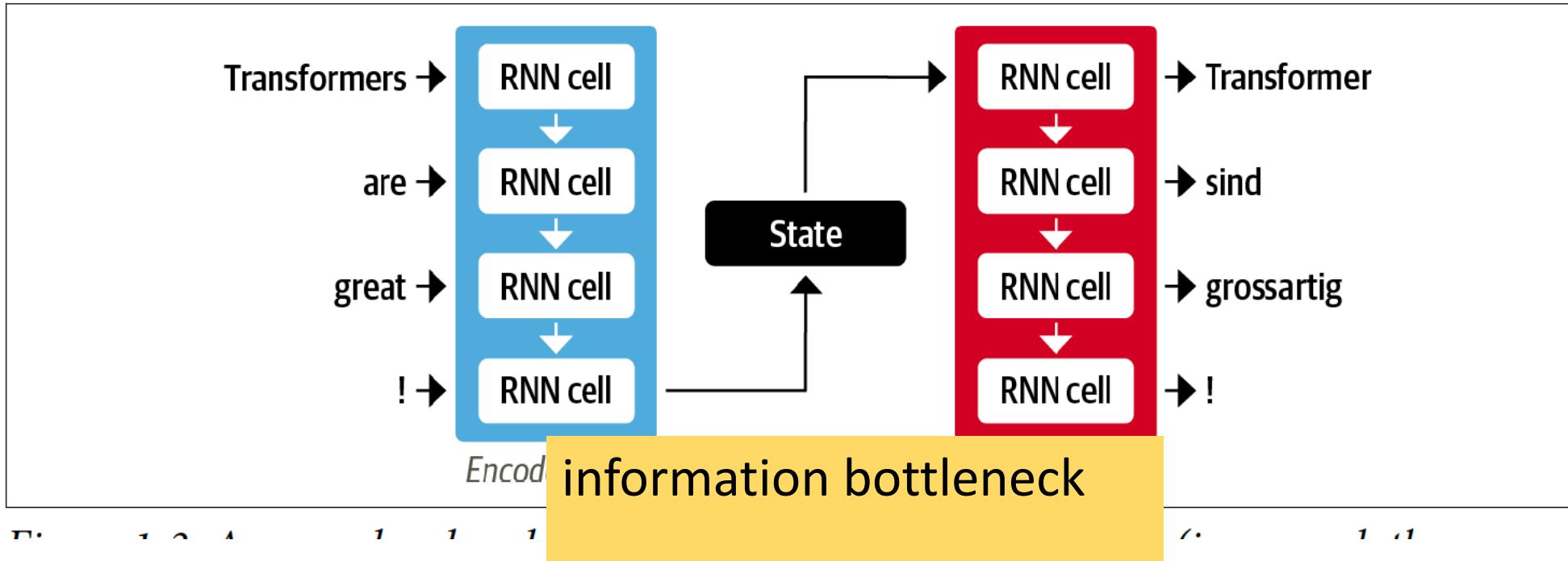
RNN Sequence to Sequence model



Je suis étudiant

<https://jalammar.github.io/visualizing-neural-machine-translation-mechanics-of-seq2seq-models-with-attention/>

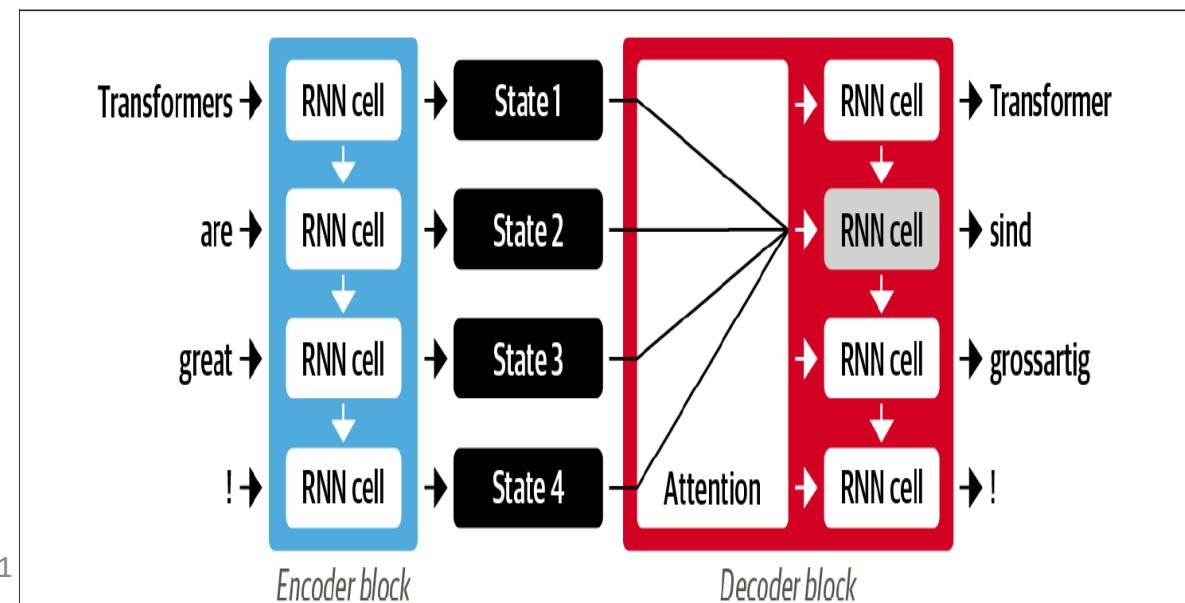
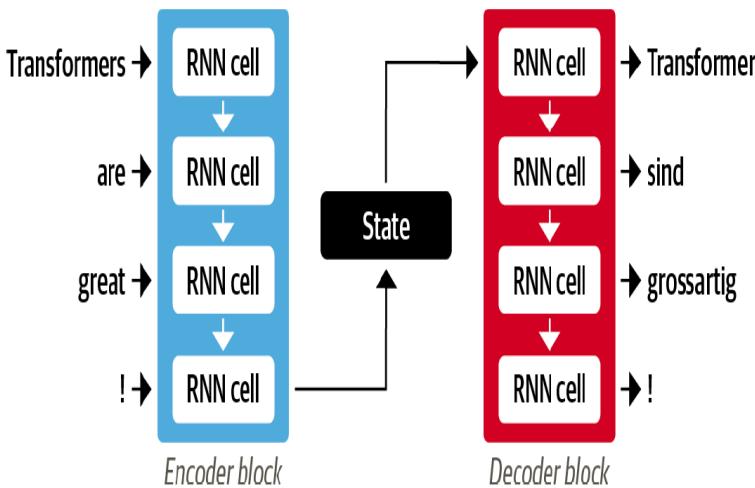
RNN Encoder Decoder



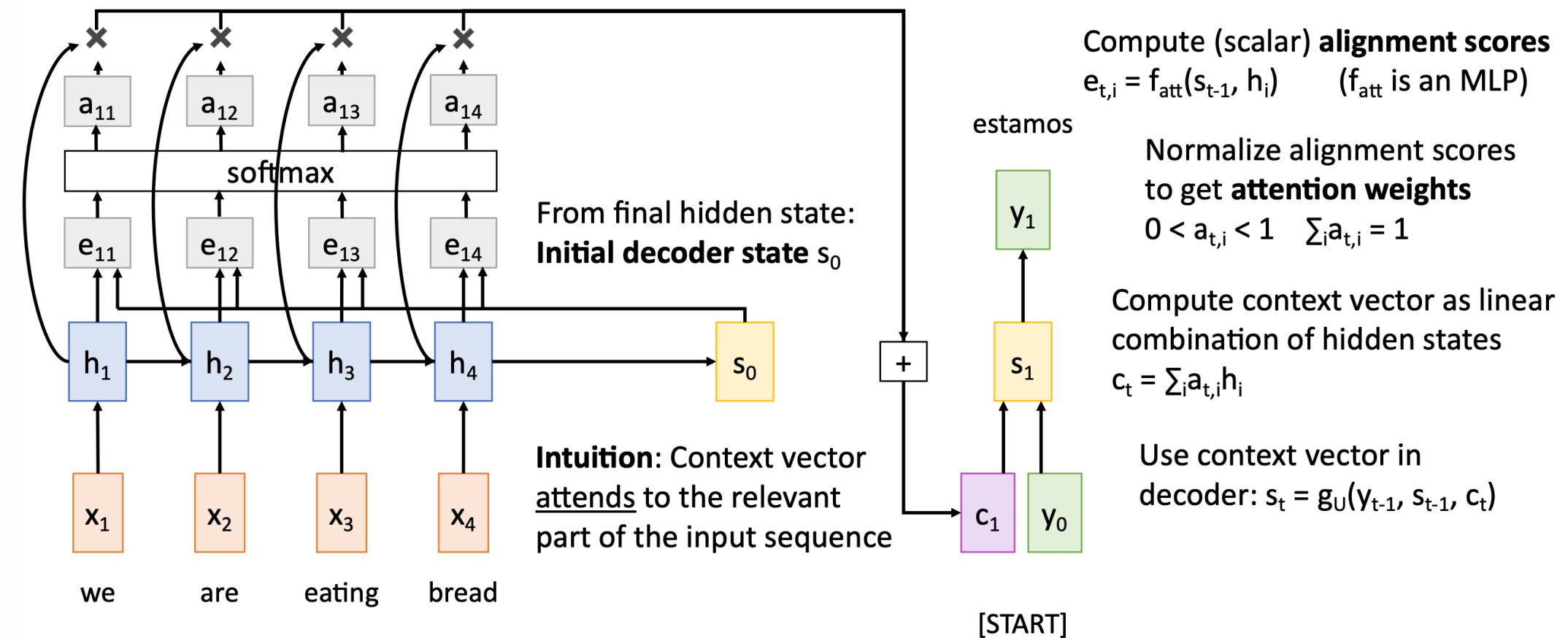
- “state” has to represent the meaning of the whole input sequence
- the decoder has access to only “state” while generating the output.
- challenging for long sequences, so information at the start of the sequence might be lost in the process of compressing everything to a single, fixed representation.

Attention

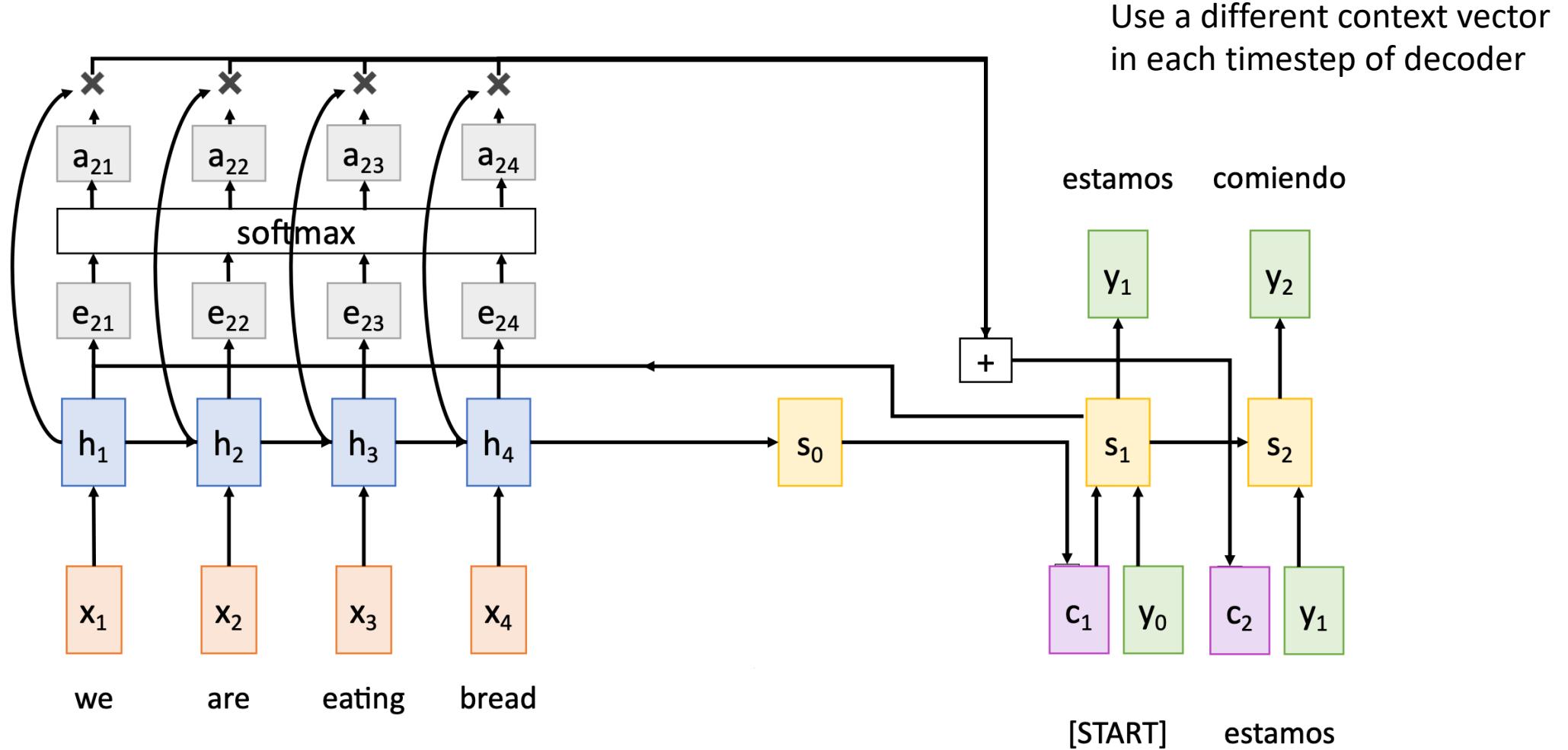
- Attention mechanisms revolutionized neural machine translation (and deep learning in general), allowing a significant improvement in the state of the art, especially for long sentences (e.g., over 30 words).
- instead of producing a single hidden state for the input sequence, the encoder outputs a hidden state at each step that the decoder can access.



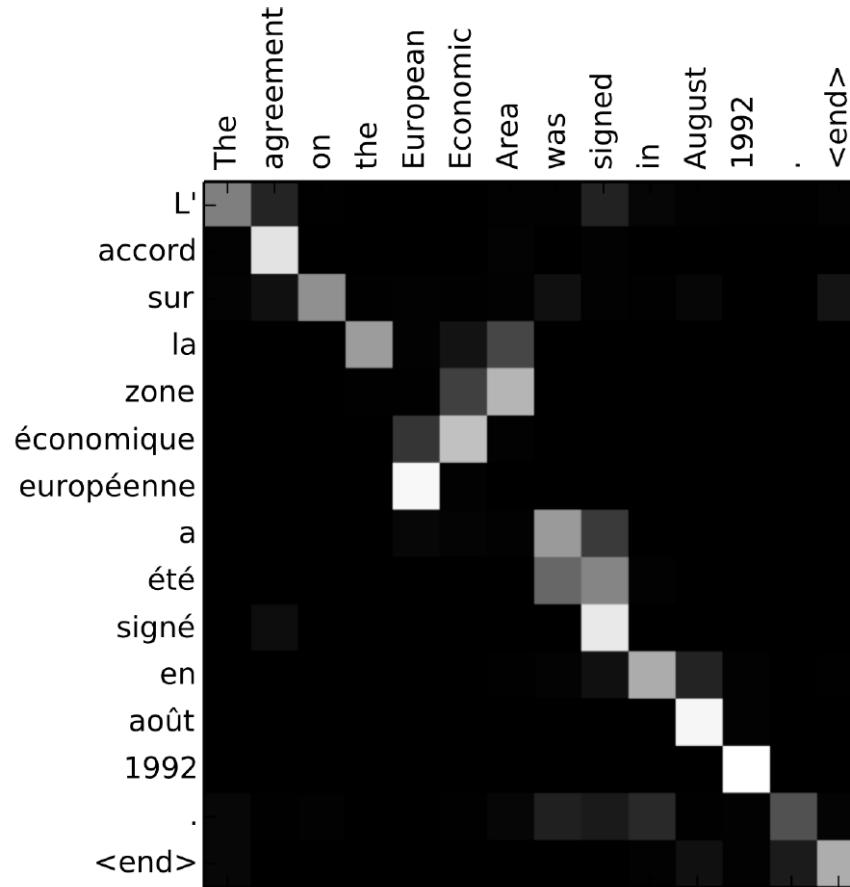
RNN Encoder Decoder (with attention mechanism)



RNN Encoder Decoder (with attention mechanism)



RNN Encoder Decoder (with attention mechanism)



decoder assign a different amount of weight, or “attention,” to each of the encoder states

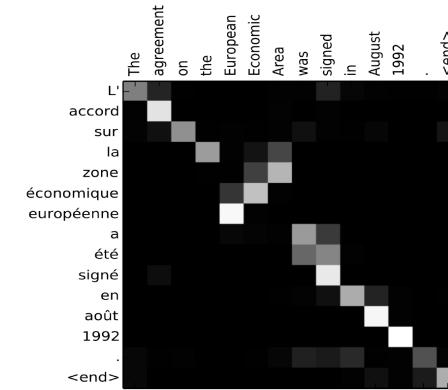
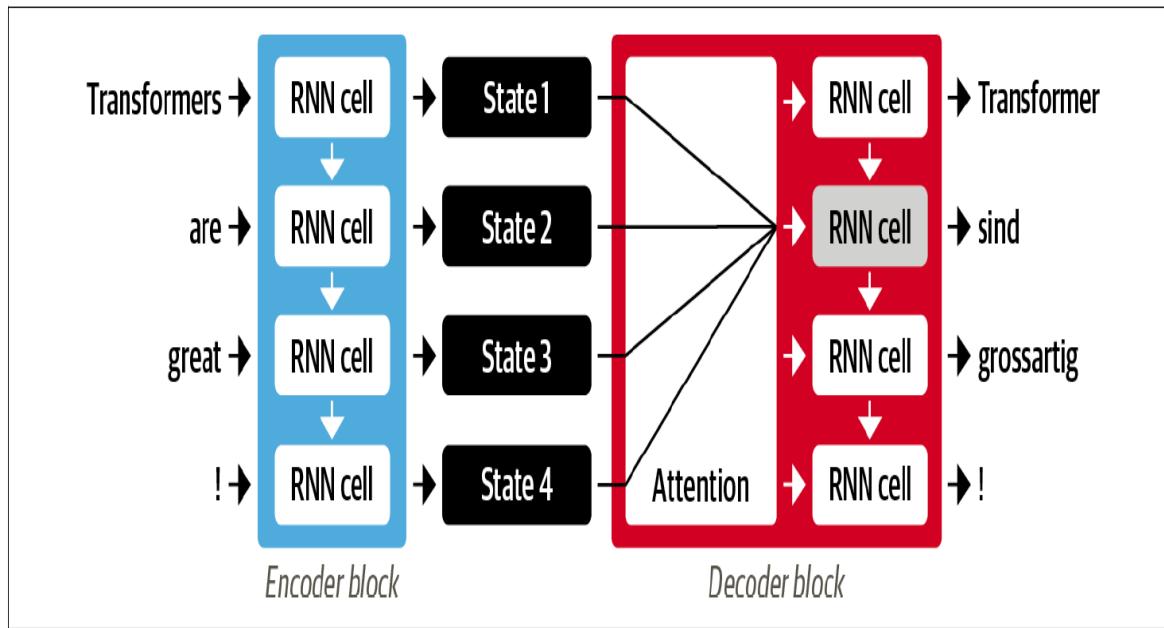
Each pixel denotes a weight.

attention-based models are able to learn nontrivial alignments between the words in a generated translation and those in a source sentence

D. Bahdanau, K. Cho, and Y. Bengio, “[Neural Machine Translation by Jointly Learning to Align and Translate](#)” (2014).

the source sentence (English) and the generated translation (French),

RNN Encoder Decoder (with attention mechanism)



- major shortcoming with using recurrent models for the encoder and decoder: the computations are inherently sequential and cannot be parallelized across the input sequence

Generalizing attention layers

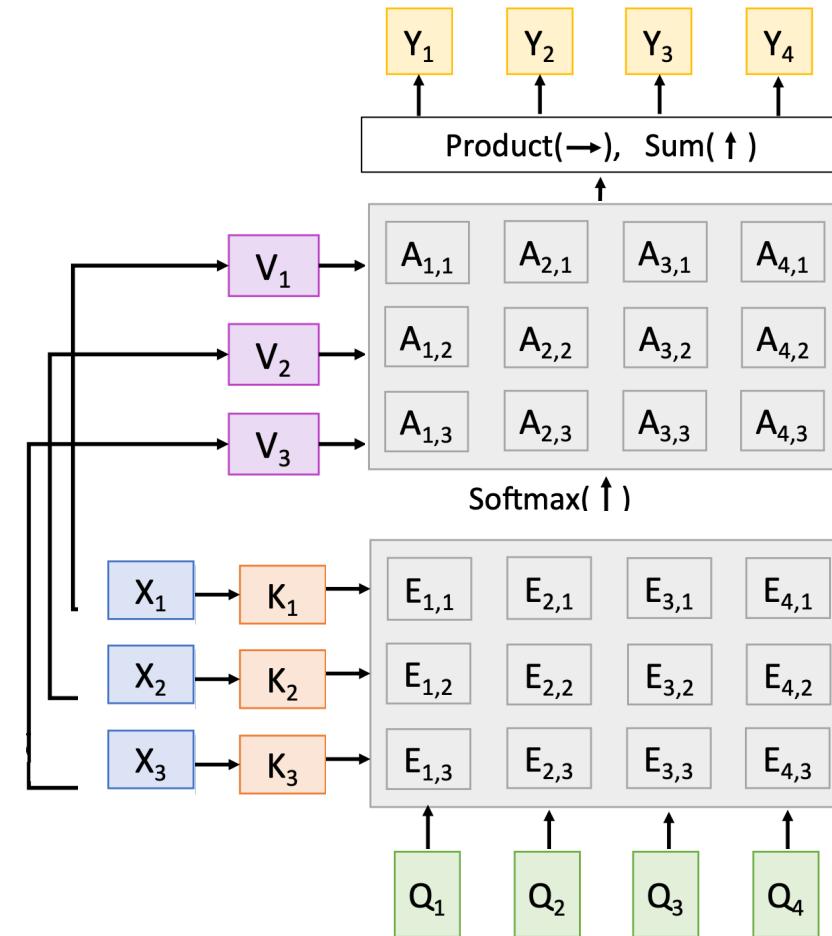
The **query**, **key** and **value** concepts come from **information retrieval syst.**

For example, when you type a query to search for some video on Youtube, the search engine will map your **query** against a set of **keys** (video title, description etc.) associated with candidate videos in the database, then present you the best matched videos (**values**).

Generalizing attention layers

Key matrix: W_K (Shape: $D_X \times D_Q$)
Value matrix: W_V (Shape: $D_X \times D_V$)

Key vectors: $K = XW_K$ (Shape: $N_X \times D_Q$)
Value Vectors: $V = XW_V$ (Shape: $N_X \times D_V$)



Self-attention layers

Every input vector x_i is used in three different ways in the self attention operation:

- Queries
- Key
- Values

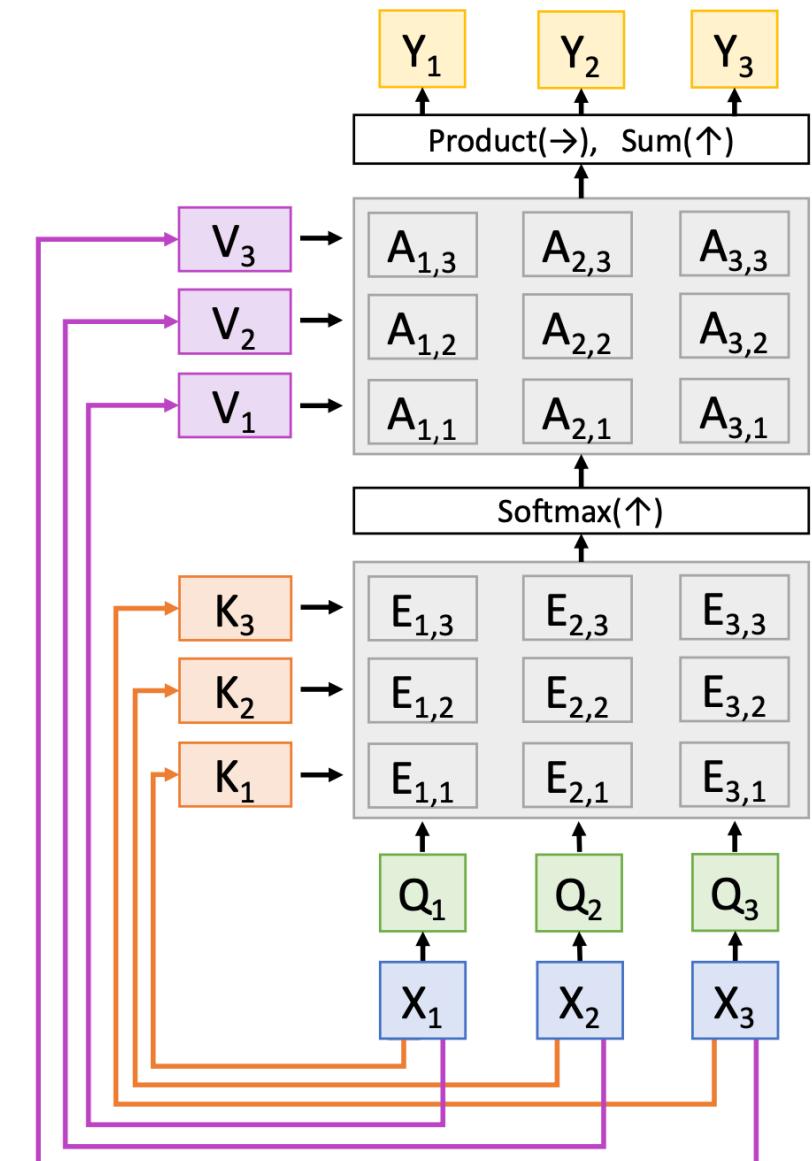
We add three $k \times k$ weight matrices W_q , W_k , W_v and compute three linear transformations of each x_i , for the three different parts of the self attention

$$q_i = W_q x_i \quad k_i = W_k x_i \quad v_i = W_v x_i$$

$$w'_{ij} = q_i^T k_j$$

$$w_{ij} = \text{softmax}(w'_{ij})$$

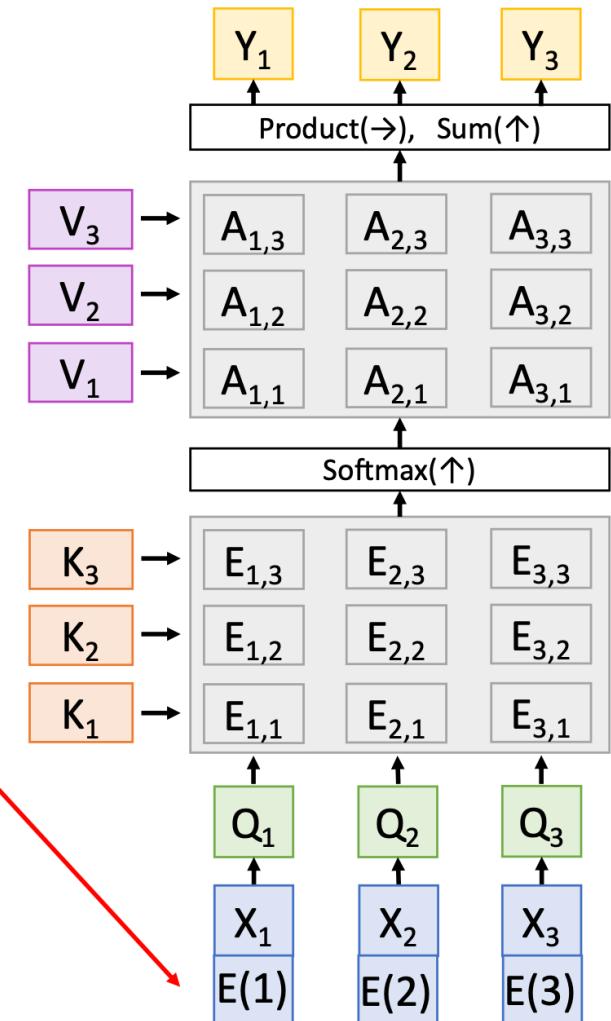
$$y_i = \sum_j w_{ij} v_j .$$



Self-attention layers

Self-attention doesn't "know" the order of the vectors it is processing!

In order to make processing position-aware, concatenate input with **positional encoding**



Positional Embedding

- Used to inject information about token positions into the input to model the sequential nature of text.
- The token embeddings are thus combined with positional embeddings that contain positional information for each token.
- Types of positional embeddings:
 - Absolute positional representations: use static patterns consisting of modulated sine and cosine signals to encode the positions of the tokens.
 - Relative positional representations: take into account the surrounding tokens are most important.

Positional Embedding Example

- <https://machinelearningmastery.com/a-gentle-introduction-to-positional-encoding-in-transformer-models-part-1/>

$$P(k, 2i) = \sin\left(\frac{k}{n^{2i/d}}\right)$$

$$P(k, 2i + 1) = \cos\left(\frac{k}{n^{2i/d}}\right)$$

n: User-defined scalar, originally set to 10,000 by the authors of [Attention Is All You Need](#).

d= Dimension of the output embedding space

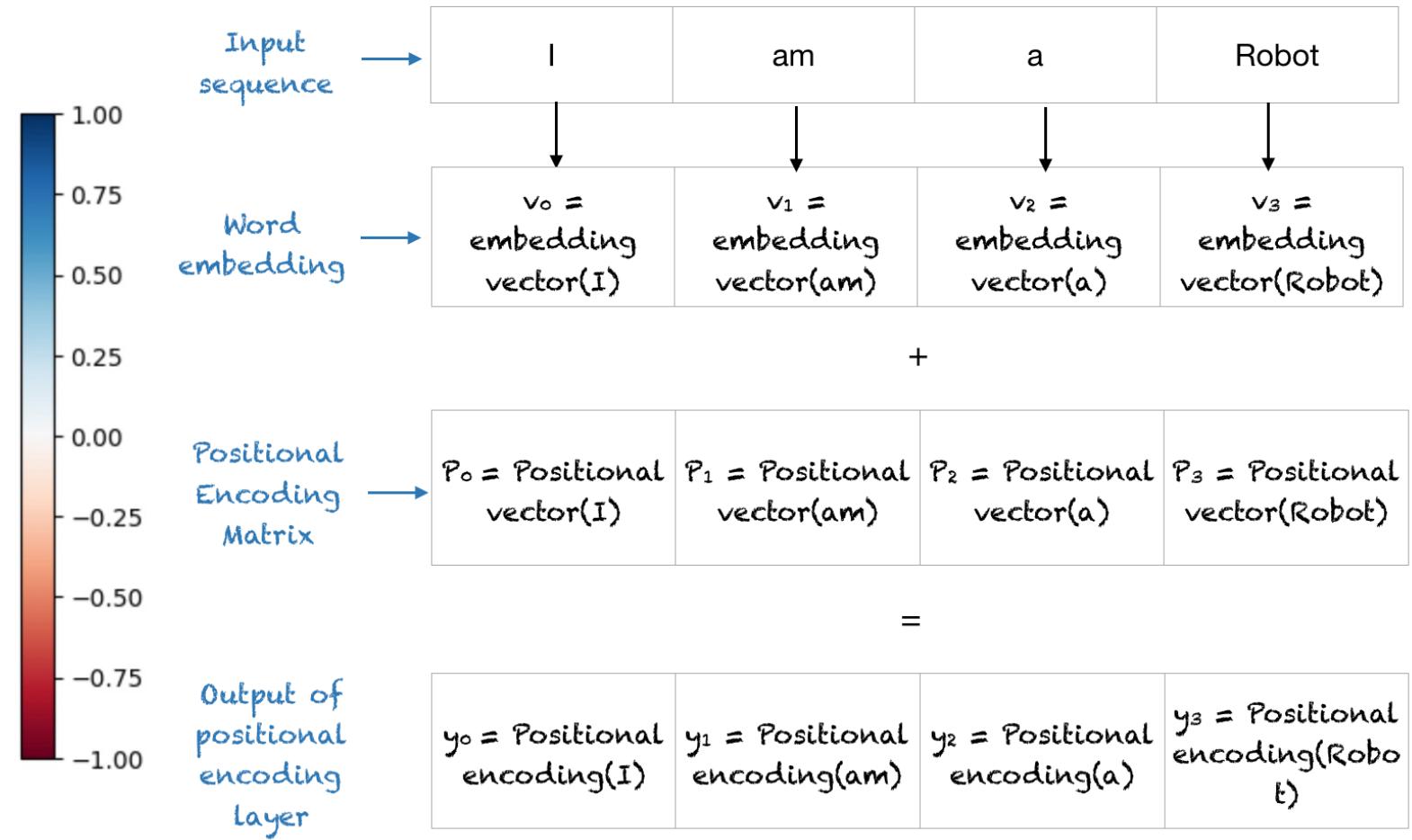
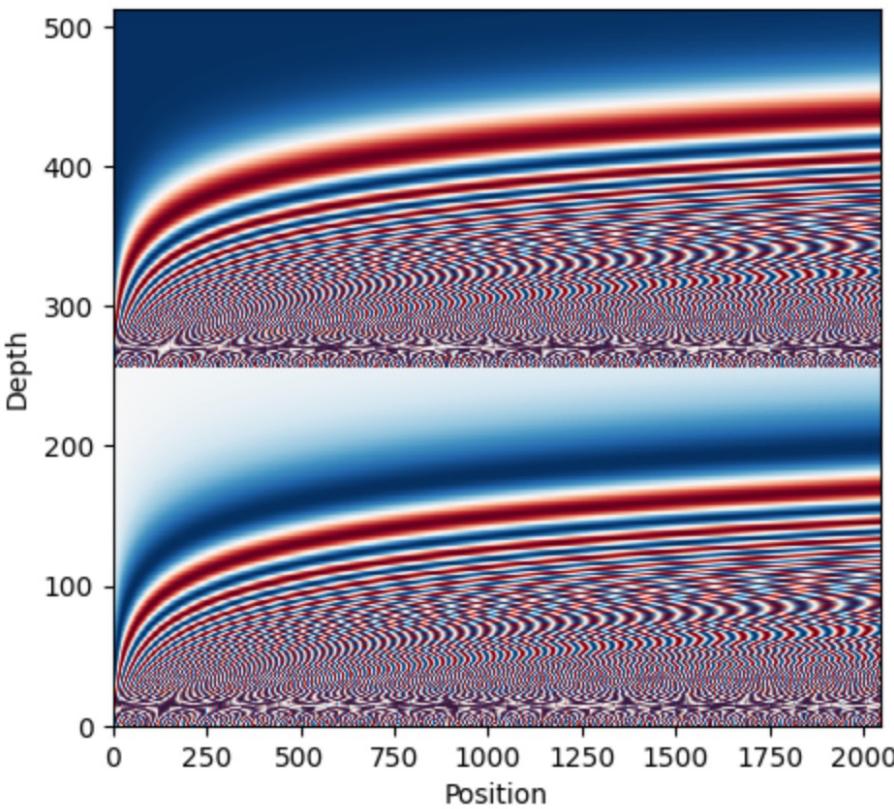
$0 \leq i < d/2$ used for mapping column indices (single value for both sin and cos)

Recent

Sequence	Index of token, k	Positional Encoding Matrix with d=4, n=100			
		i=0	i=0	i=1	i=1
I	0	$P_{00} = \sin(0) = 0$	$P_{01} = \cos(0) = 1$	$P_{02} = \sin(0) = 0$	$P_{03} = \cos(0) = 1$
am	1	$P_{10} = \sin(1/1) = 0.84$	$P_{11} = \cos(1/1) = 0.54$	$P_{12} = \sin(1/10) = 0.10$	$P_{13} = \cos(1/10) = 1.0$
a	2	$P_{20} = \sin(2/1) = 0.91$	$P_{21} = \cos(2/1) = -0.42$	$P_{22} = \sin(2/10) = 0.20$	$P_{23} = \cos(2/10) = 0.98$
Robot	3	$P_{30} = \sin(3/1) = 0.14$	$P_{31} = \cos(3/1) = -0.99$	$P_{32} = \sin(3/10) = 0.30$	$P_{33} = \cos(3/10) = 0.96$

Positional Encoding Matrix for the sequence 'I am a robot'

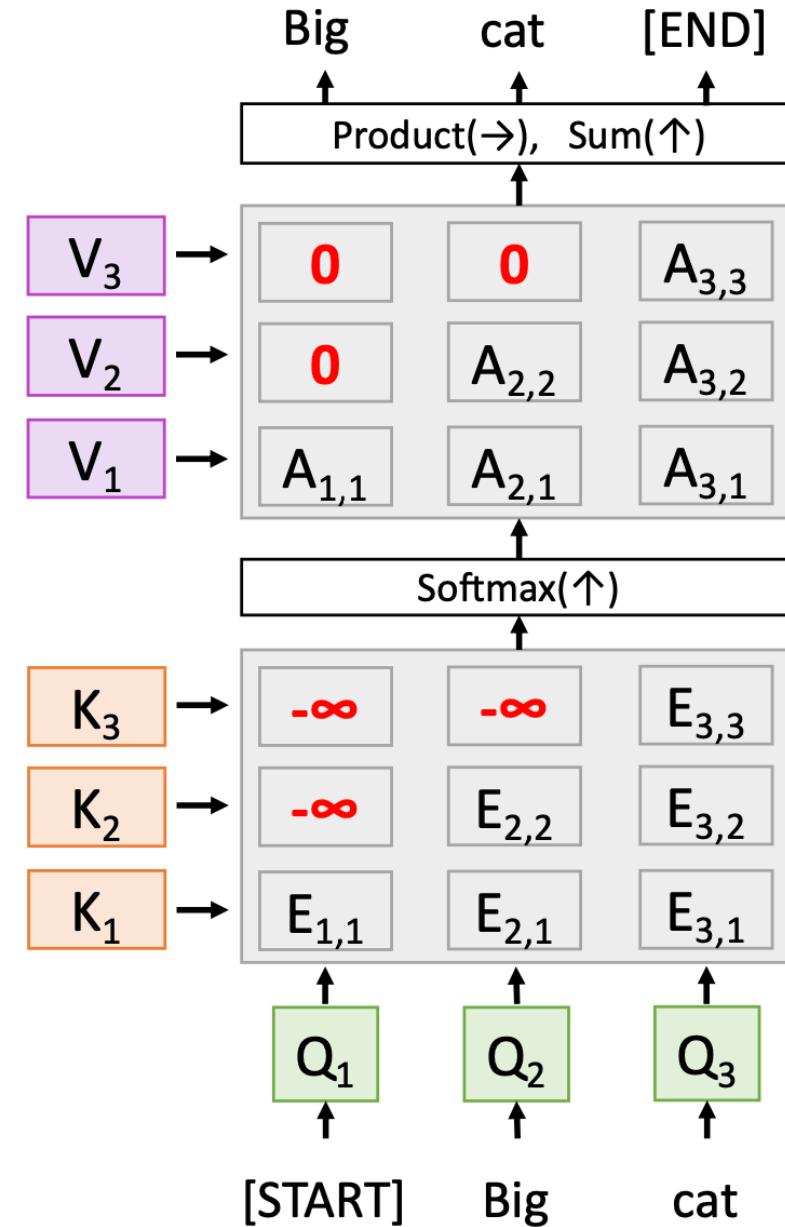
Positional Embedding Example



Masked Self-attention

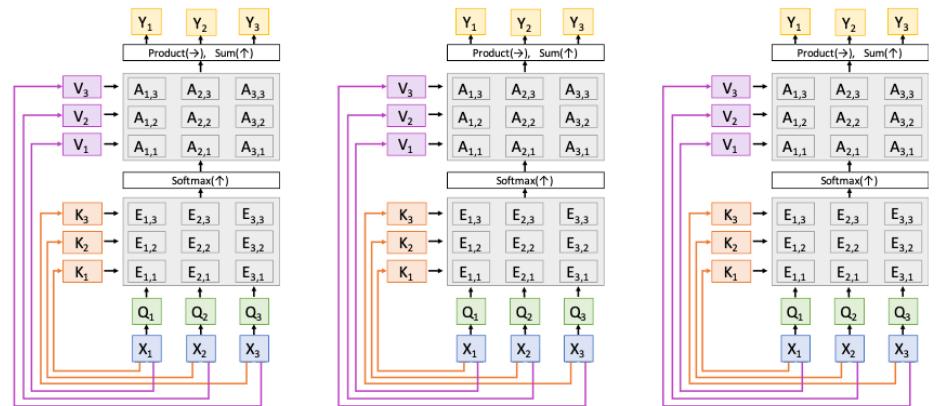
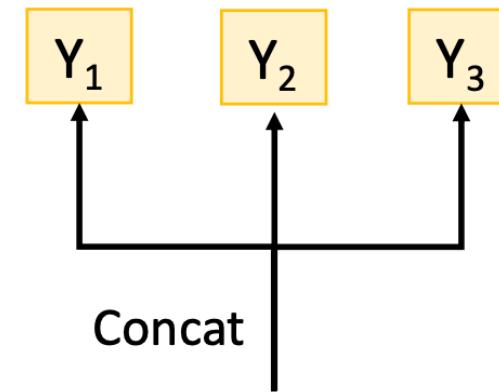
Don't let vectors “look ahead” in the sequence

Used for language modelling
(predict next word)



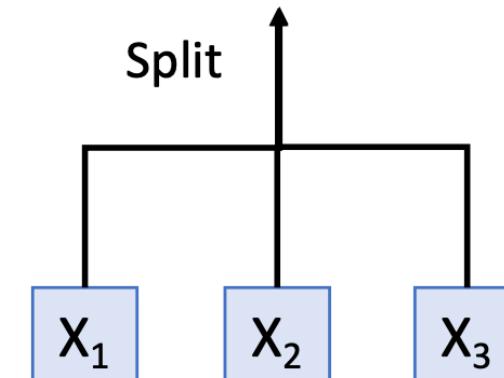
Multi-head Self-attention

Use H independent “Attention Heads” in parallel (input split into equal chunks)



Hyperparameters:

- Query dimension D_Q
- Number of heads H



Self-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

Transformer Neural Networks

Transformers (original architecture)

- No recurrence (sequential processing)
- Rely on Multi-head self-attention
- Allow attention to operate on all the states in the same layer of the neural network

The Transformer Architecture

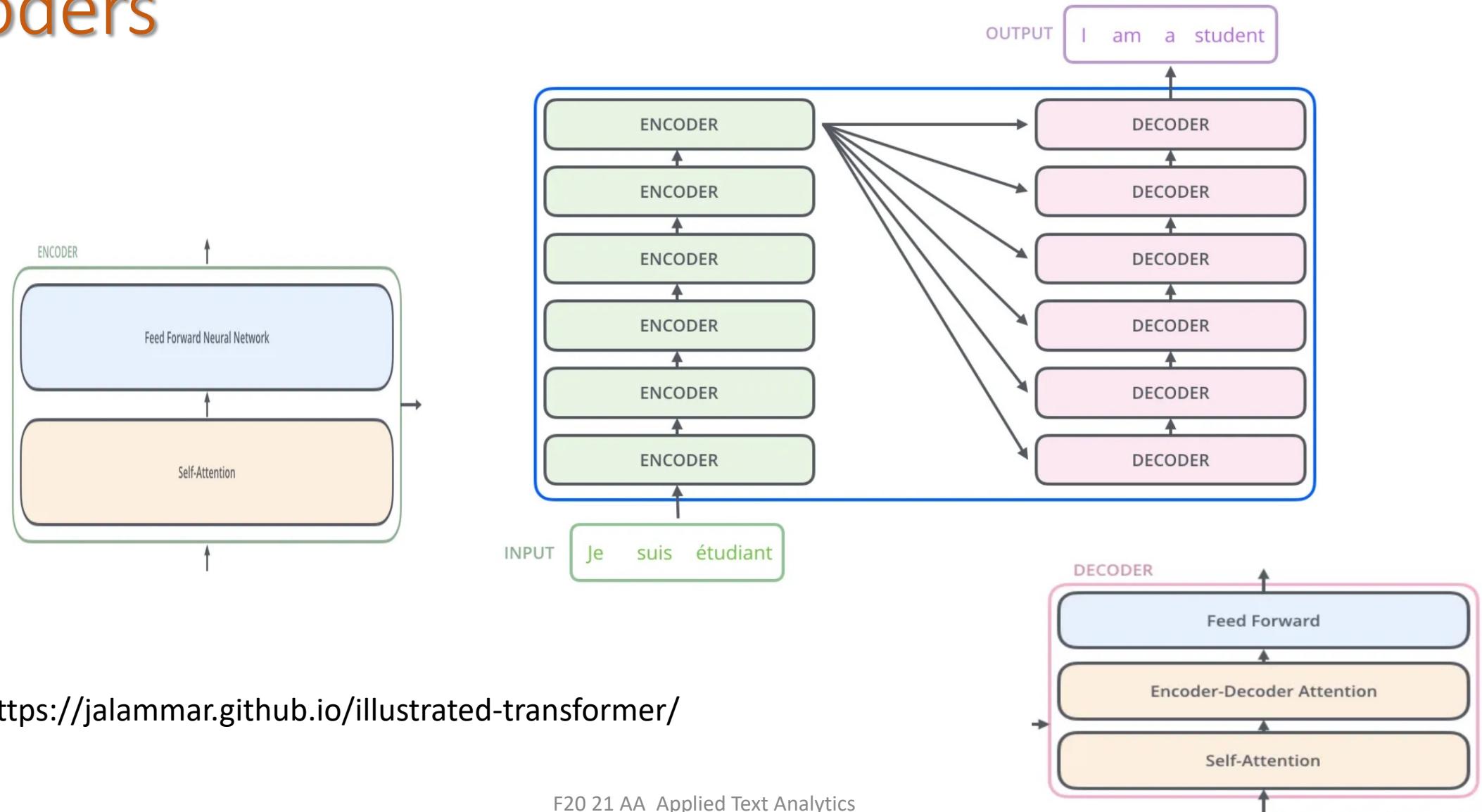
Encoder

- Converts an input sequence of tokens into a sequence of embedding vectors, often called the hidden state or context

Decoder

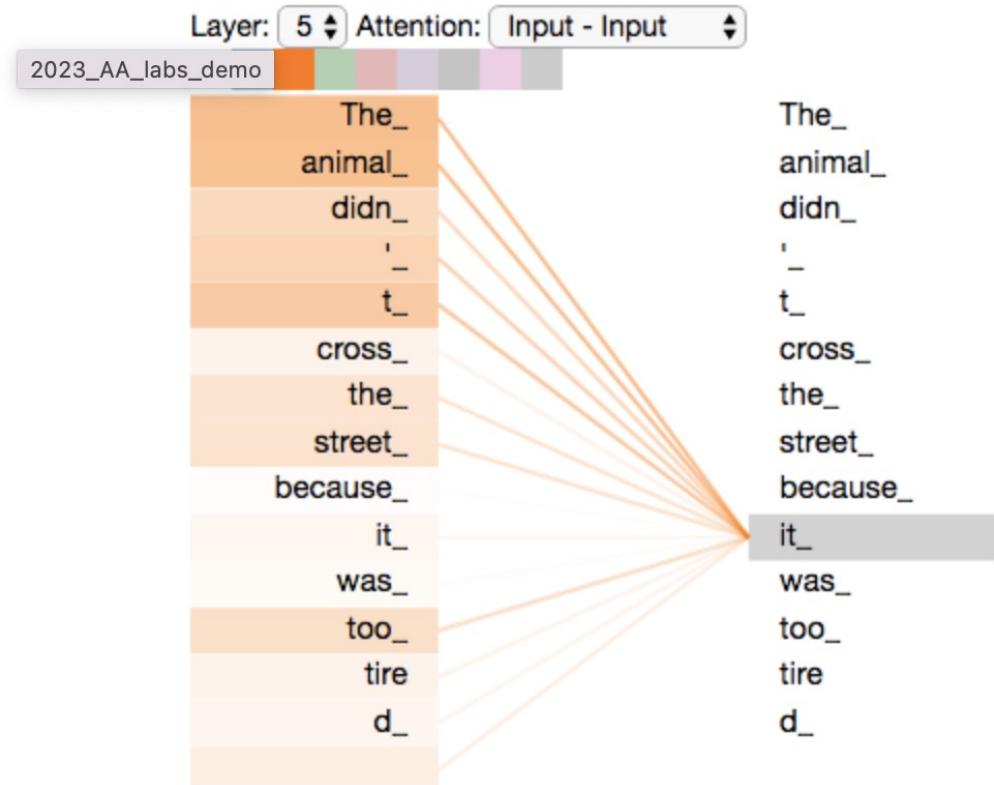
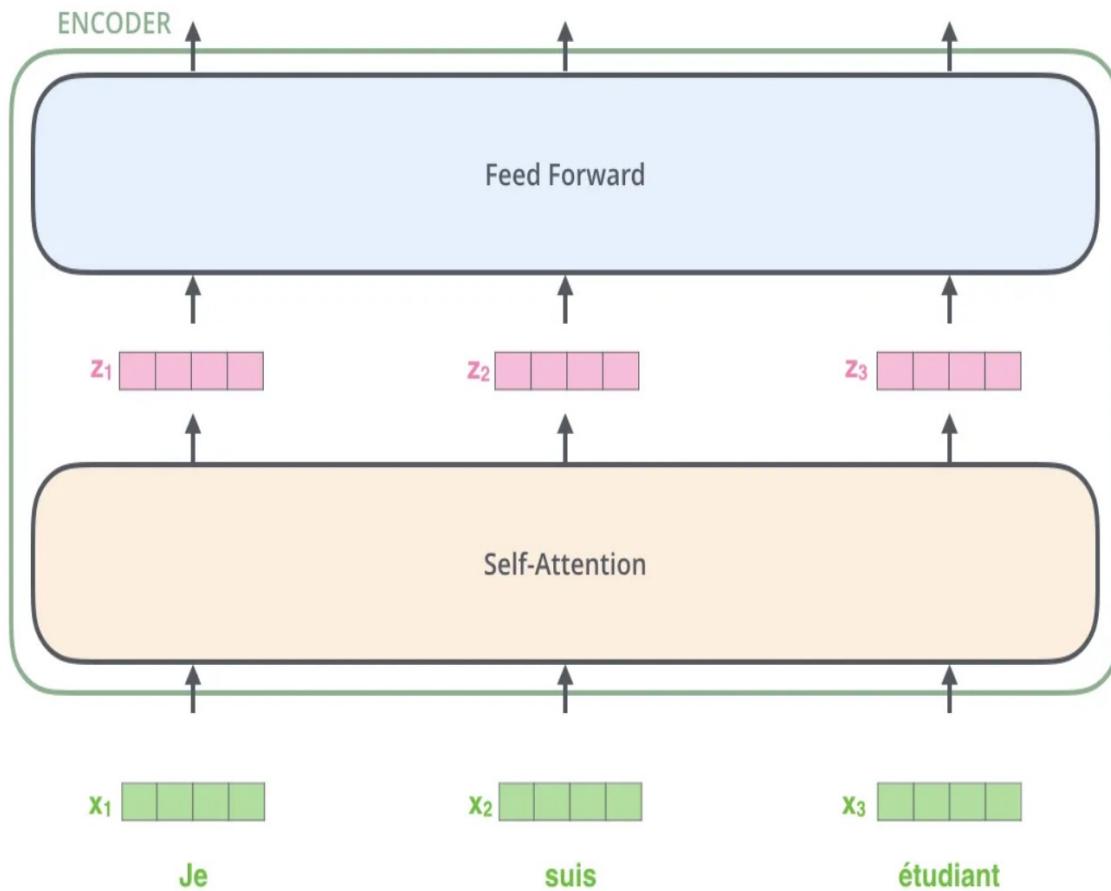
- Uses the encoder's hidden state to *iteratively* generate an output sequence of tokens, one token at a time
- What are Transformer Encoders & Decoders composed of?

Transformer consists of six encoders and six decoders



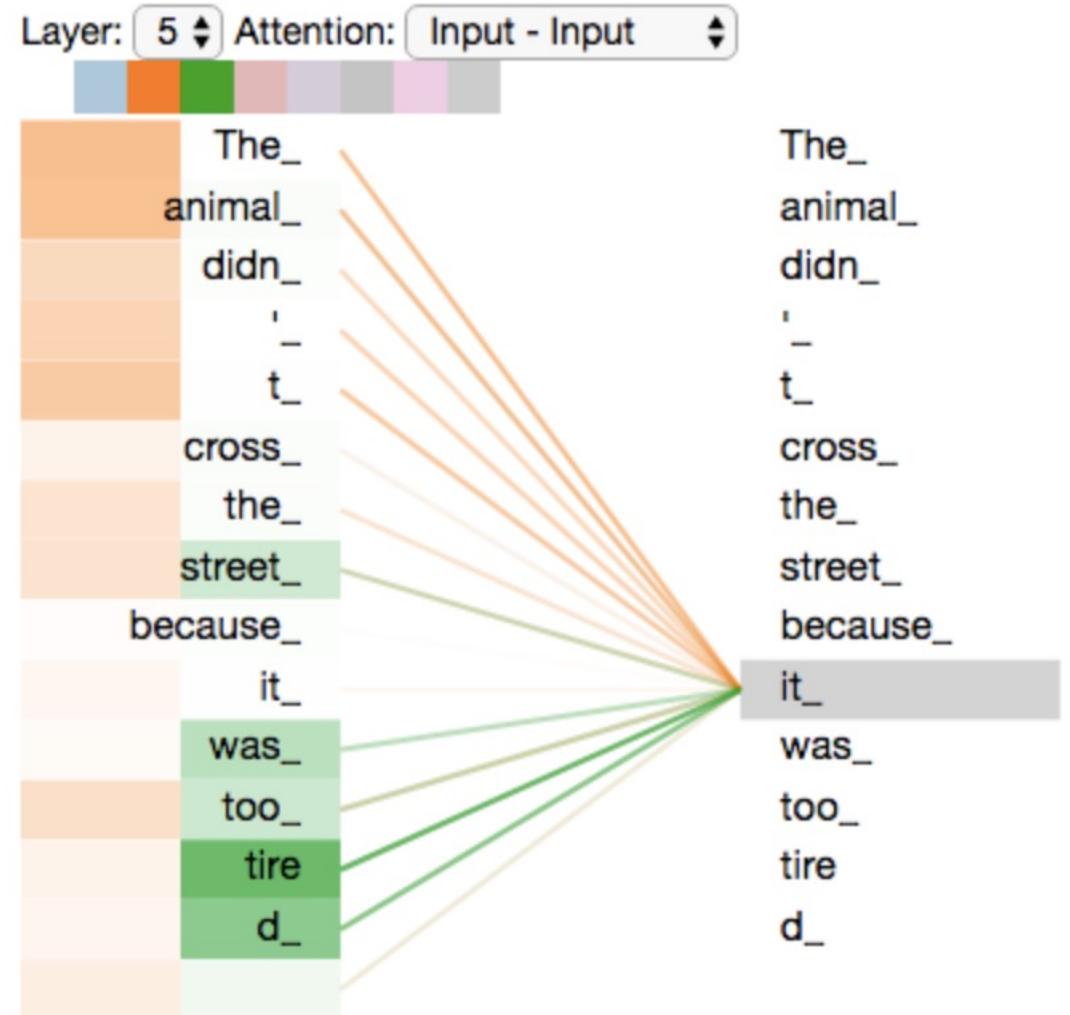
<https://jalammar.github.io/illustrated-transformer/>

Encoder self-attention



Multi-headed attention

- It gives the attention layer multiple “representation subspaces”
- Therefore multiple sets of Query/Key/Value weight matrices
- The Transformer uses eight attention heads=eight sets for each encoder/decoder



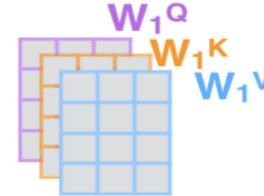
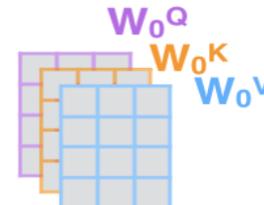
Encoder self-attention (putting it all together)

1) This is our input sentence* 2) We embed each word*

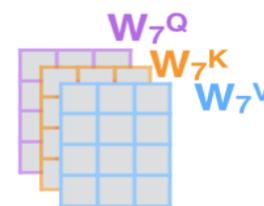
Thinking
Machines



3) Split into 8 heads.
We multiply **X** or **R** with weight matrices



...



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



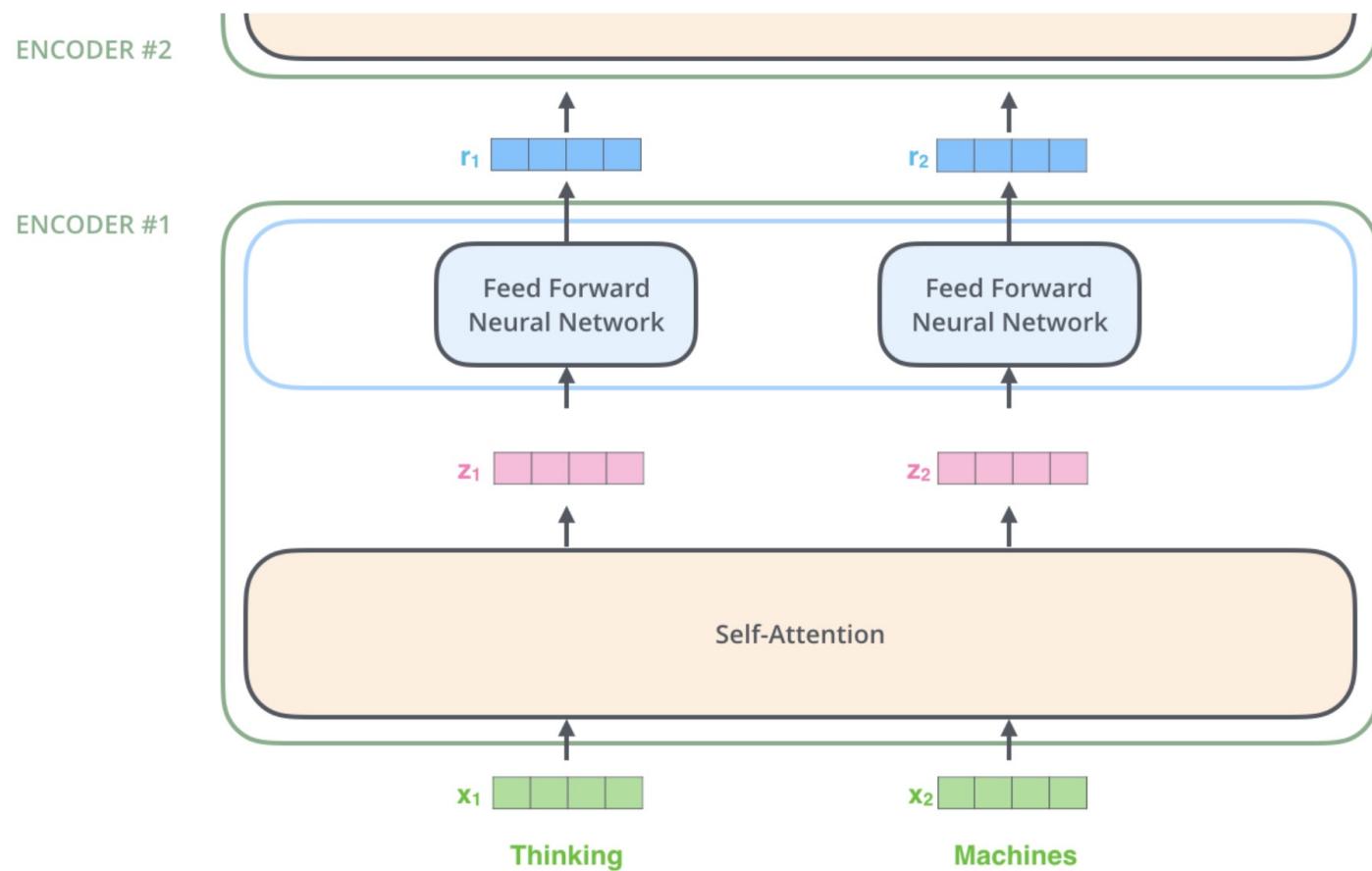
5) Concatenate the resulting **Z** matrices, then multiply with weight matrix **W^O** to produce the output of the layer



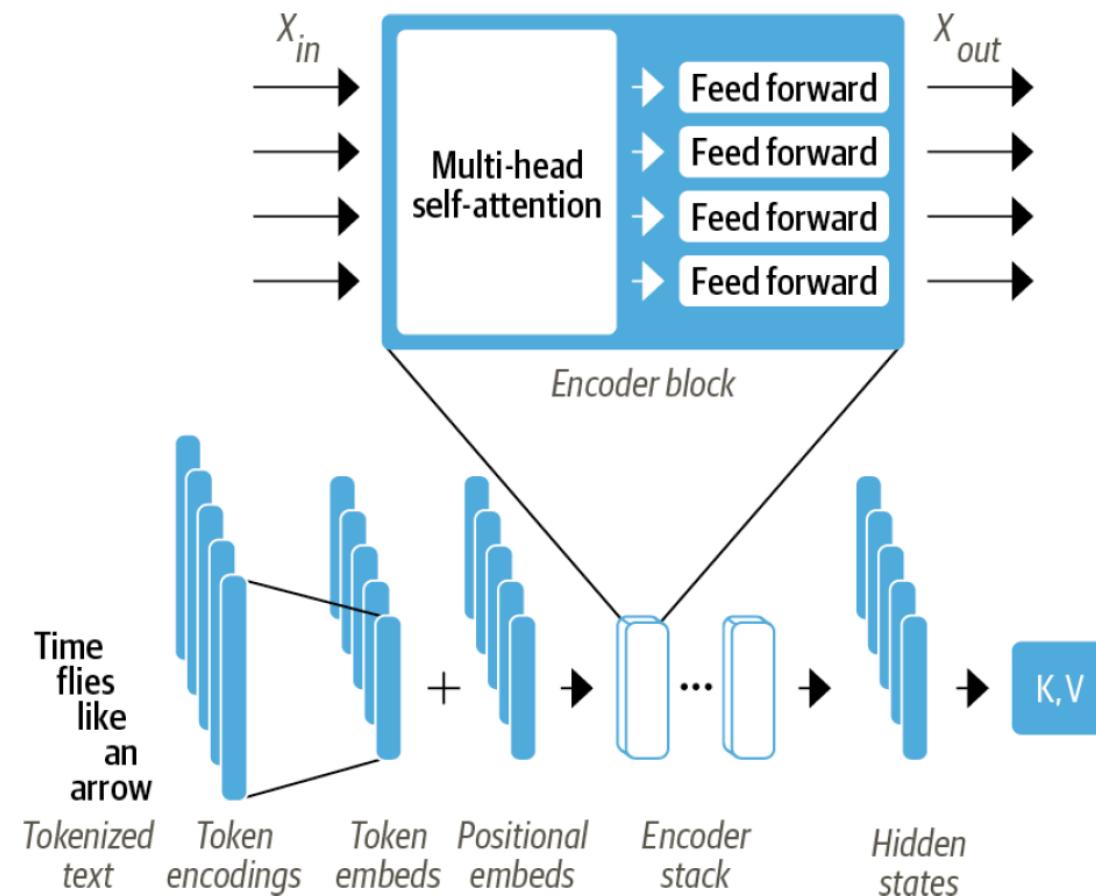
* In all encoders other than #0, we don't need embedding.
We start directly with the output of the encoder right below this one



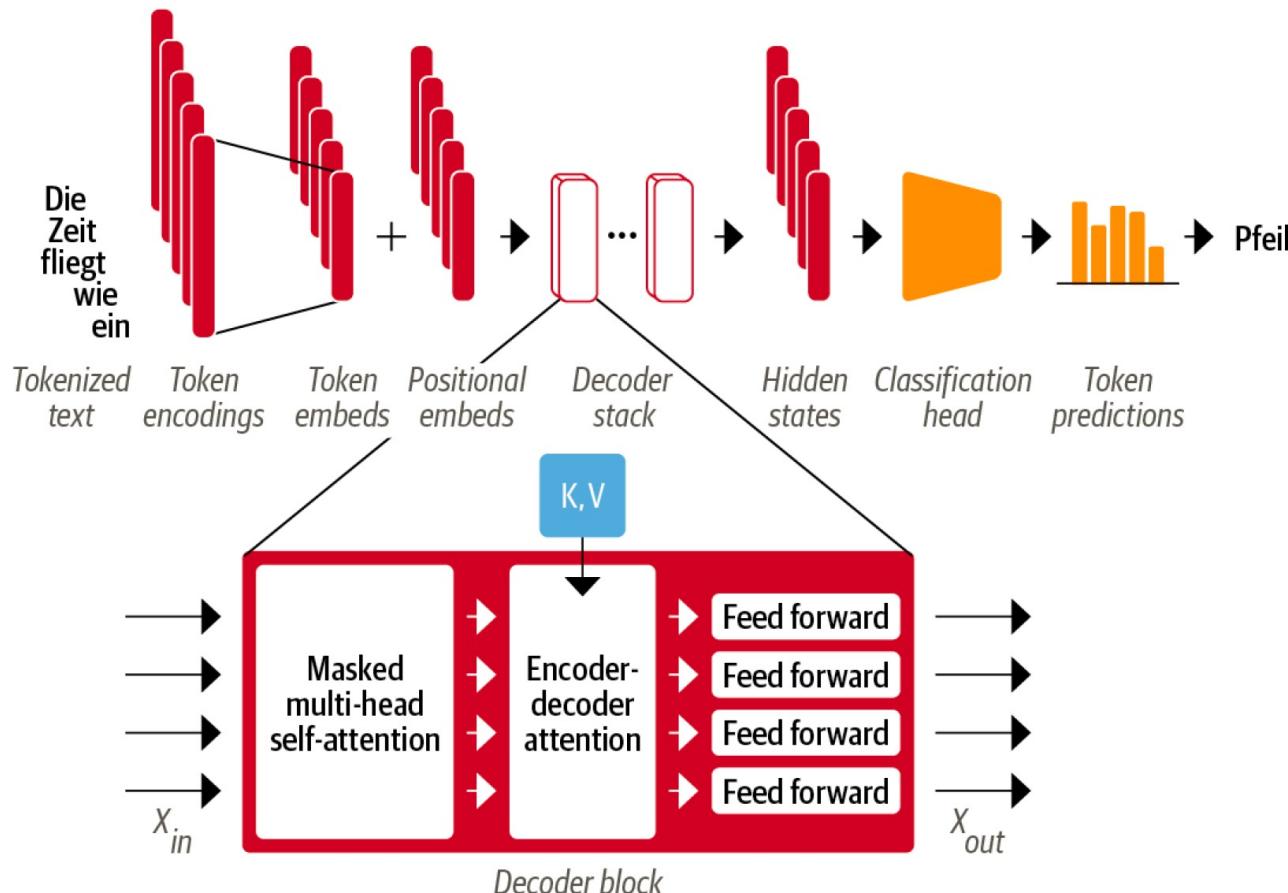
Encoder self-attention (putting it all together)



Transformer Encoder



Transformer Decoder

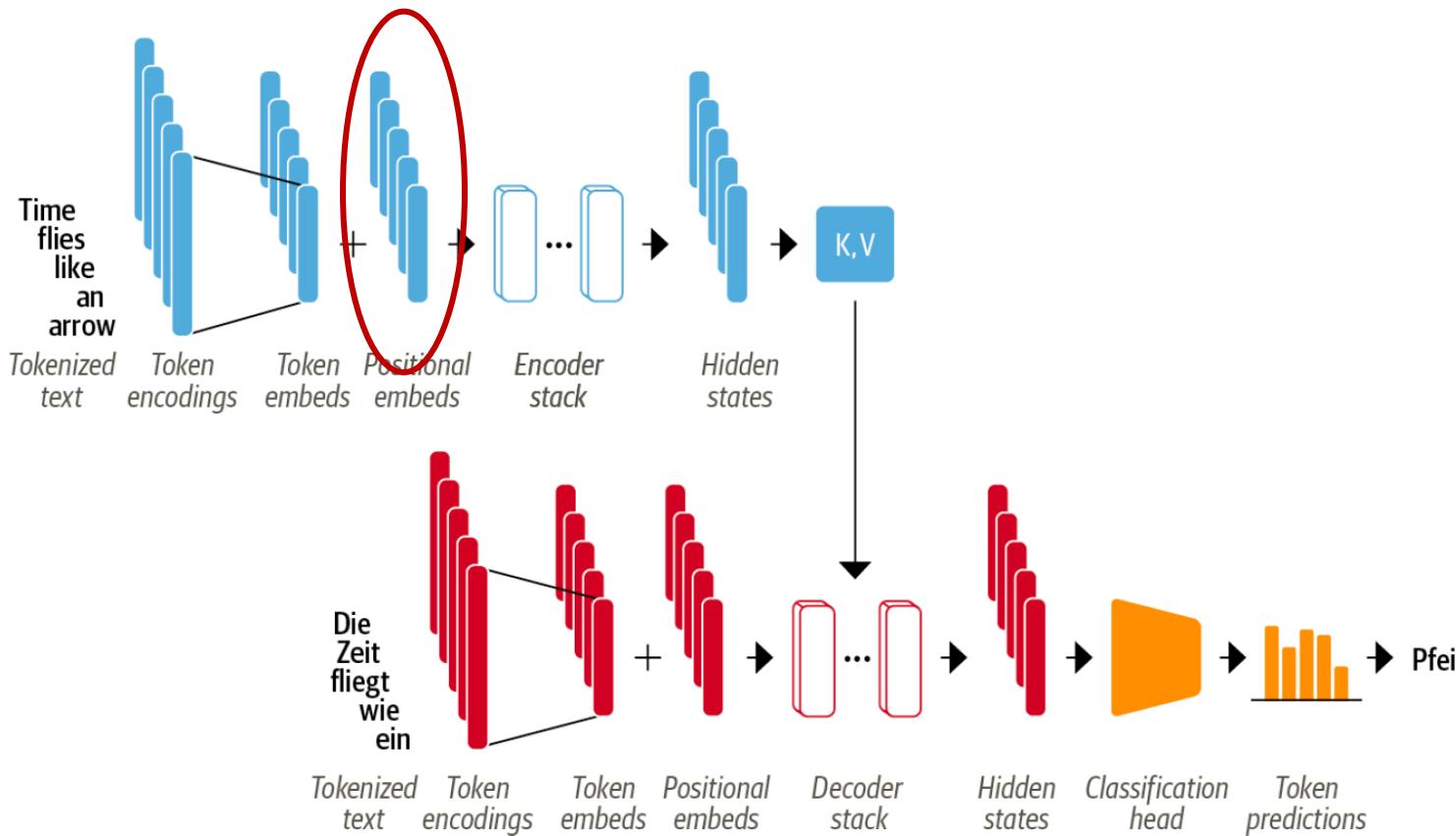


Encoder-decoder attention layer

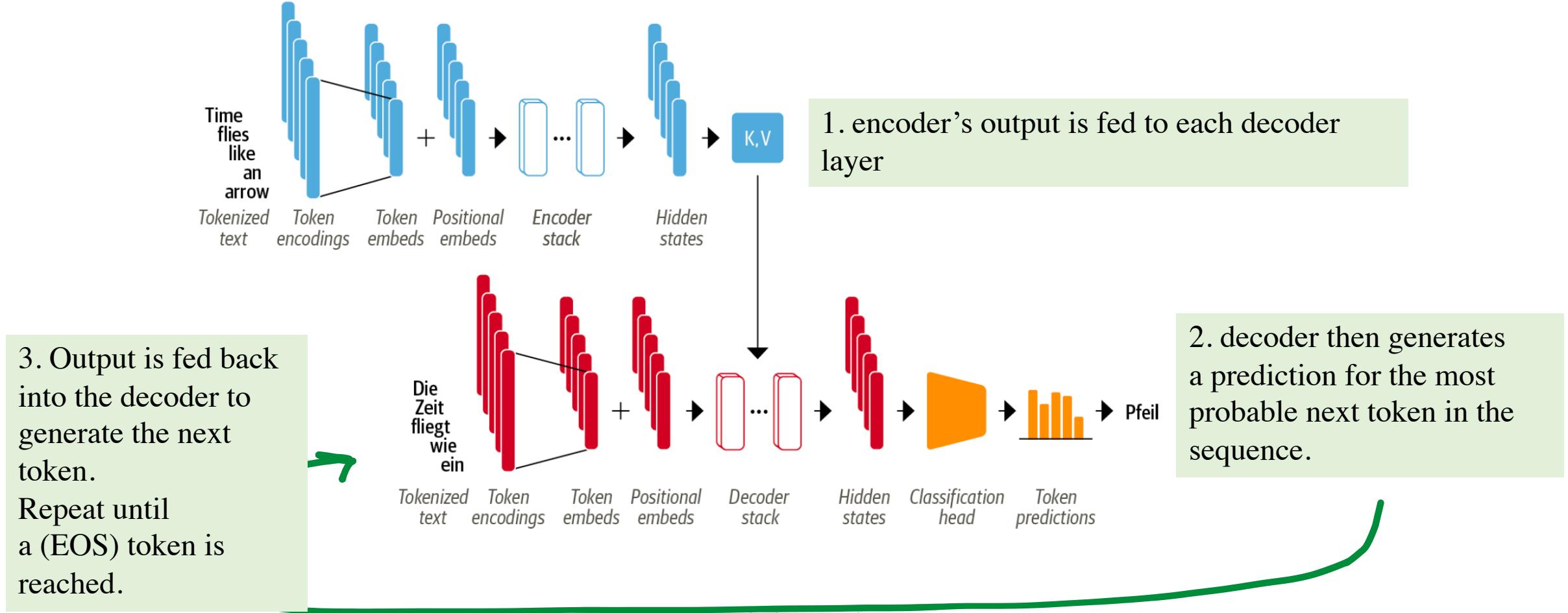
Performs **multi-head attention** over the output **key** and **value** vectors of the encoder stack, with the intermediate representations of the decoder acting as the queries.

This way the encoder-decoder attention layer learns how to relate tokens from two different sequences, such as two different languages. The decoder has access to the encoder keys and values in each block.

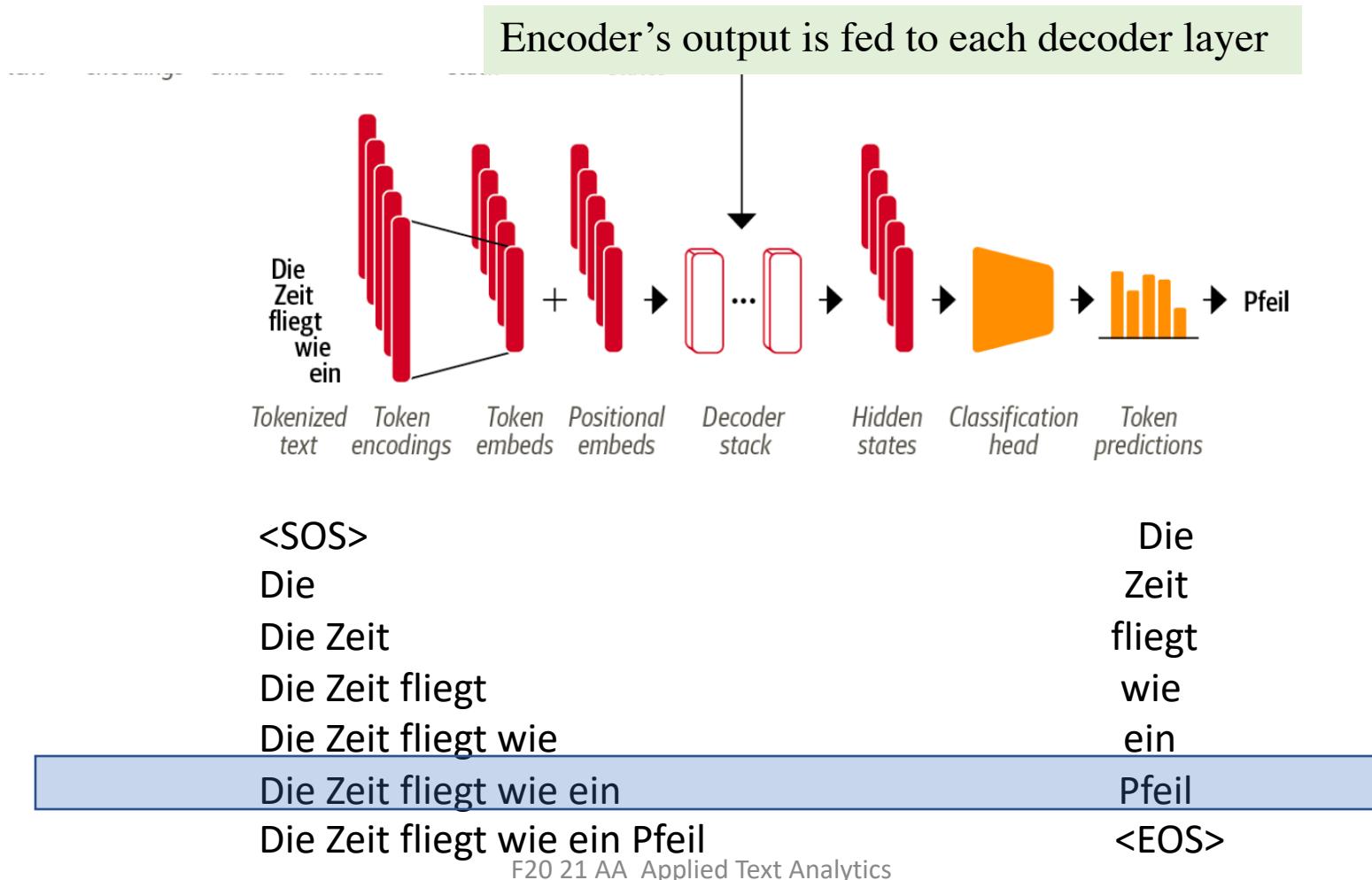
The Transformer Architecture



The Transformer Architecture



The Transformer Architecture- Decoder side



The Transformer Categories

Originally designed for sequence-to-sequence tasks like machine translation, but both the encoder and decoder blocks were soon adapted as standalone models

Encoder-only

- convert an input sequence of text into a rich numerical representation
- Used for text classification or named entity recognition.
- BERT and its variants, like RoBERTa and DistilBERT

Decoder-only

- The family of GPT models belong to this class.
- models will autocomplete the sequence by iteratively predicting the most probable next word
- **Further Encoders in addition to Transformers:** BART and T5 models
 - suitable for machine translation and summarization tasks

Your best way for a gentle start !

Gentle introduction to Transformers Jay Allamar ([channel](#))

<https://jalammar.github.io/illustrated-transformer/> (**Great READ**)

- 8 video series of Transformers & NLP
<https://www.youtube.com/watch?v=jo1NZ3vCS90&list=PLaJCKi8Nk1hwaMUYxJMjM3jTB2o58A6WY>
- ***Beautifully Illustrated: NLP Models from RNN to Transformer, oct 2022***
- <https://towardsdatascience.com/beautifully-illustrated-nlp-models-from-rnn-to-transformer-80d69faf2109>
- ***All you need to know about ‘Attention’ and ‘Transformers’ – In-depth Understanding – Part 1 (Feb 2022)***<https://towardsdatascience.com/all-you-need-to-know-about-attention-and-transformers-in-depth-understanding-part-1-552f0b41d021>

Practical Part

- Recommended Tensorflow tutorial & case studies
- Neural machine translation with attention (translate text with seq2seq models)
[https://www.tensorflow.org/text/tutorials/nmt with attention](https://www.tensorflow.org/text/tutorials/nmt_with_attention)
- Neural machine translation with a Transformer and Keras (translate text with transformer models)
<https://www.tensorflow.org/text/tutorials/transformer>

Thank you!