

Последовательность -> Последовательность

Машинный перевод

Пример:

ITA: Il gatto si e' seduto sul tappetino.



EN: The cat sat on the mat.

Проблемы:

- Выравнивание: последовательности на входе / выходе могут иметь разную длину
- Неопределенность (отображение 1-ко-многим: множество возможных способов перевода)
- Метрика: как автоматически оценивать, означают ли предложения к одно и то же?

Последовательность -> Последовательность

Машинный перевод

Пример:

ITA: Il gatto si e' seduto sul tappetino.



EN: The cat sat on the mat.

Подход:

Имеем один RNN для кодирования исходного предложения, а другой RNN - для предсказания целевого предложения. Целевой RNN учится (мягко) выравнивать предложение с помощью механизма внимания.

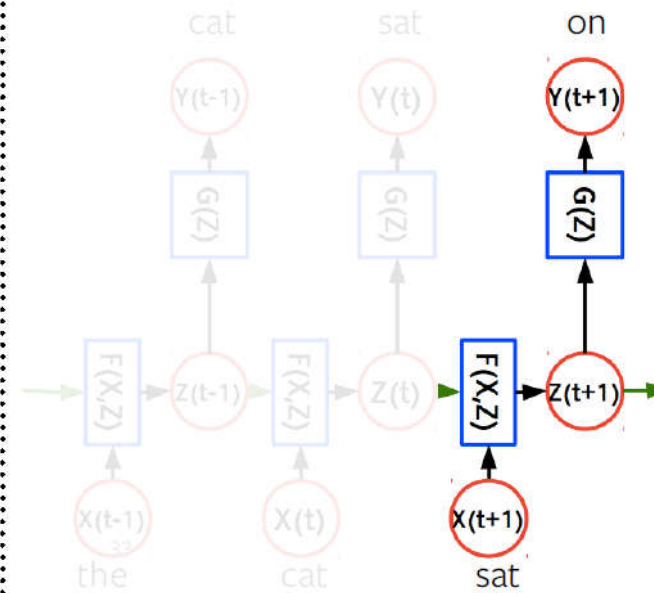
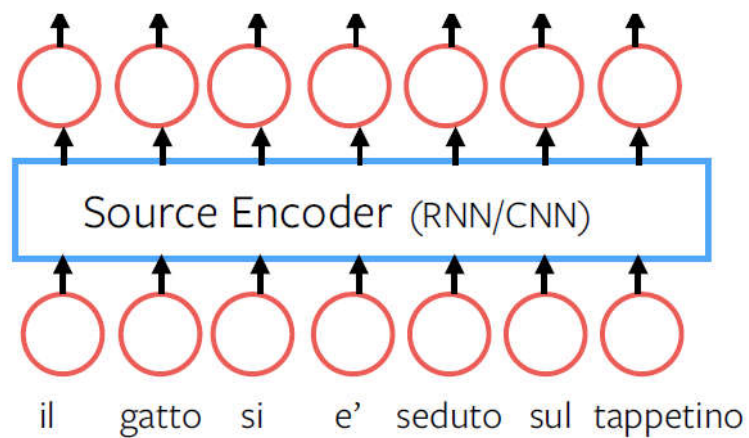
Последовательность -> Последовательность

Машинный перевод

Source

Target

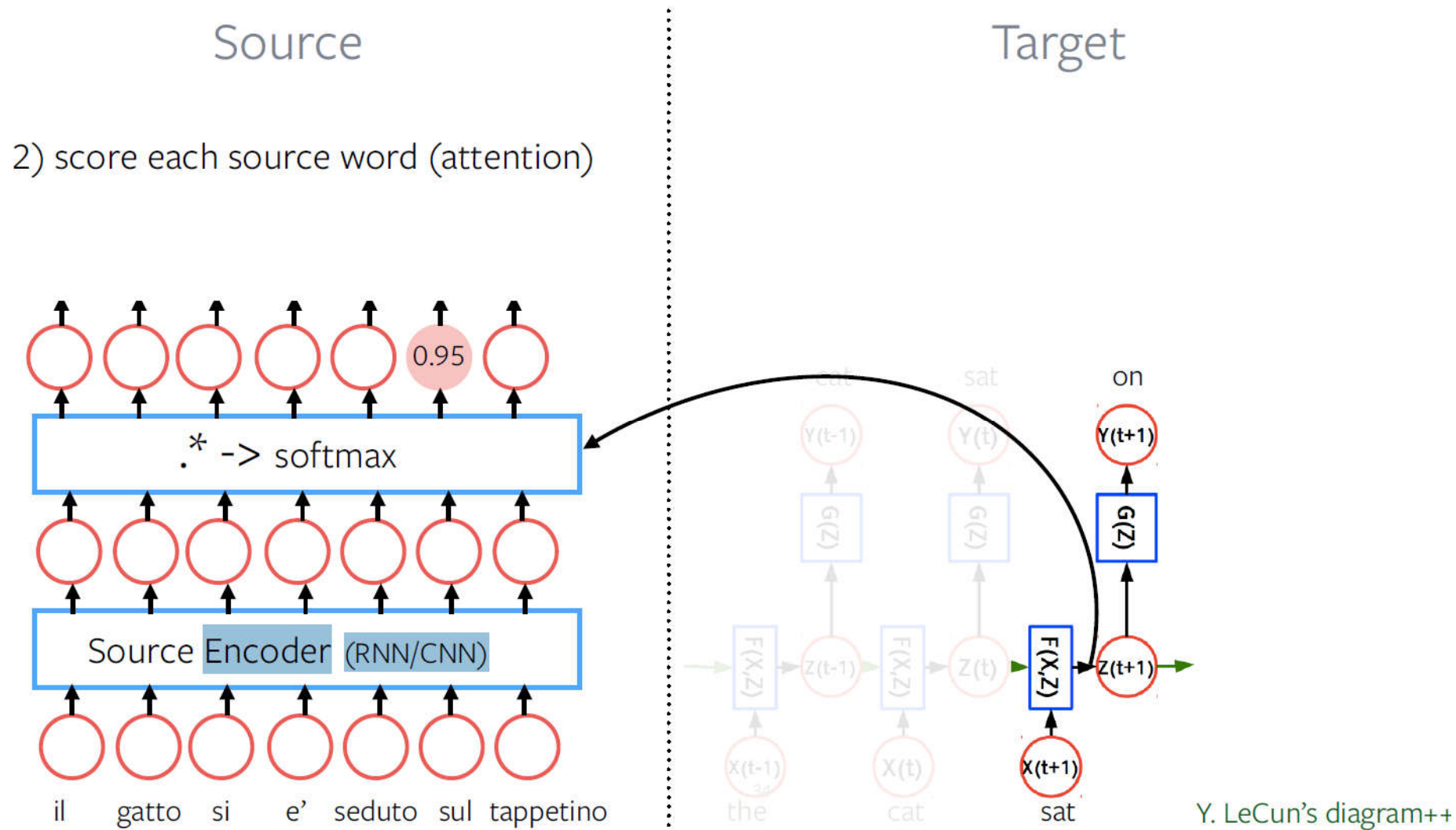
1) Represent source



Y. LeCun's diagram++

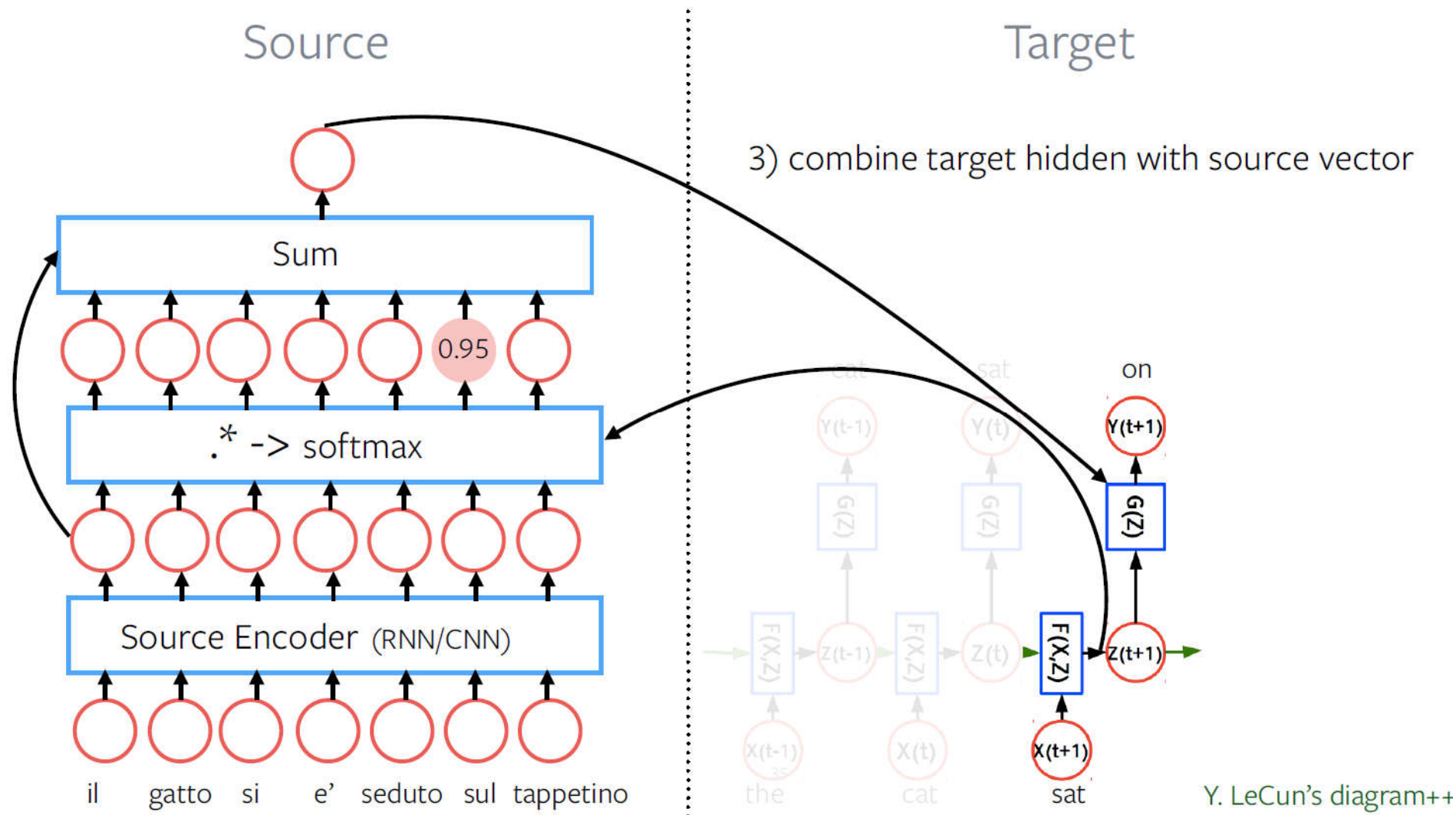
Последовательность -> Последовательность

Машинный перевод



Последовательность -> Последовательность

Машинный перевод



Последовательность -> Последовательность

Машинный перевод

Пример:

ITA: Il gatto si e' seduto sul tappetino.



EN: The cat sat on the mat.

Заметки:

- + Исходное и целевое предложение может иметь любую длину, хорошо работает и при длинных предложениях!
- + Учится неявно выравнивать
- + RNN можно заменить на СНС. *A convolutional encoder model for NMT, Gehring et al. 2016*
- + Генерирует беглые предложения
- Имеет проблемы с редко встречаемыми словами, точным выбором слов
- Обычно обучается как языковая модель (кросс-энтропия), хороша для вычисления оценки, но не для генерации

Последовательность -> Последовательность

Машинный перевод

Заключение:

- + Механизм внимания является довольно общим и может использоваться для:
 - + Работает с входами переменной длины, поскольку «мягко выбирает один»
 - + Неявное выравнивание, которое обучается моделью, по мере необходимости
 - + Для выполнения раундов «объяснений» (например, «переходов» в сетях памяти)
- + Этот же механизм использовался для создания субтитров, суммирования и т. Д.
- Функция потерь на уровне слов (кросс-энтропия для предсказания следующего слова) является субоптимальной для задач генерации текстов

Sequence level training with RNNs, Ranzato et al. ICLR 2016

An actor-critic algorithm for sequence prediction, ICLR 2017

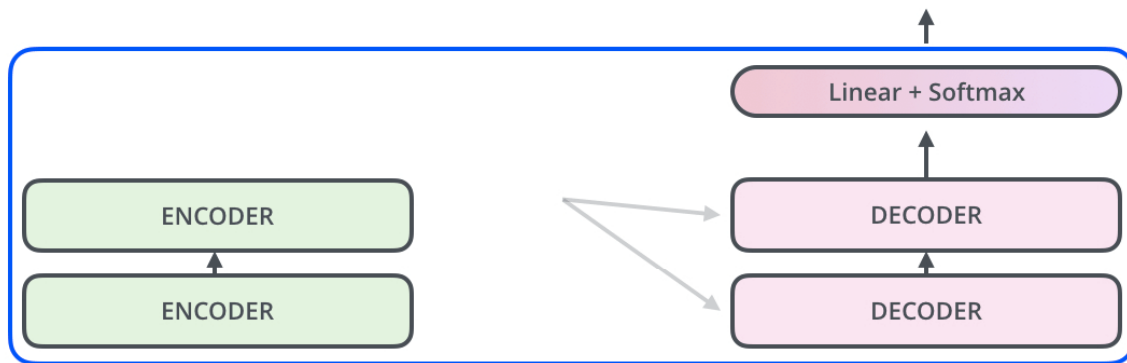
Sequence-to-sequence learning as beam-search optimization, EMNLP 2016

Трансформер

- Трансформер с механизмом ВНИМАНИЯ
- BERT, GPT

Decoding time step: 1 2 3 4 5 6

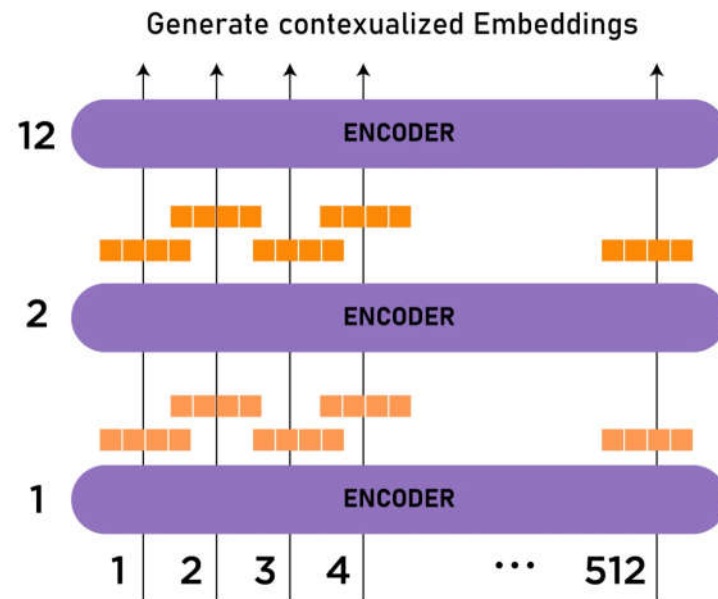
OUTPUT



EMBEDDING WITH TIME SIGNAL

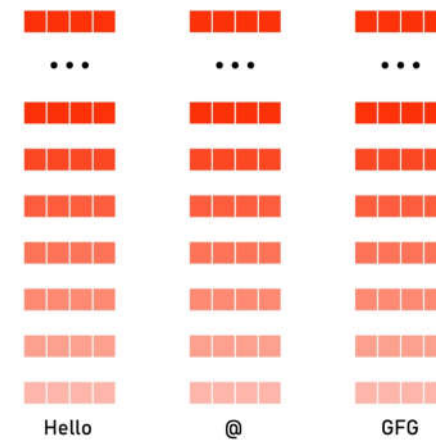
EMBEDDINGS

INPUT Je suis étudiant

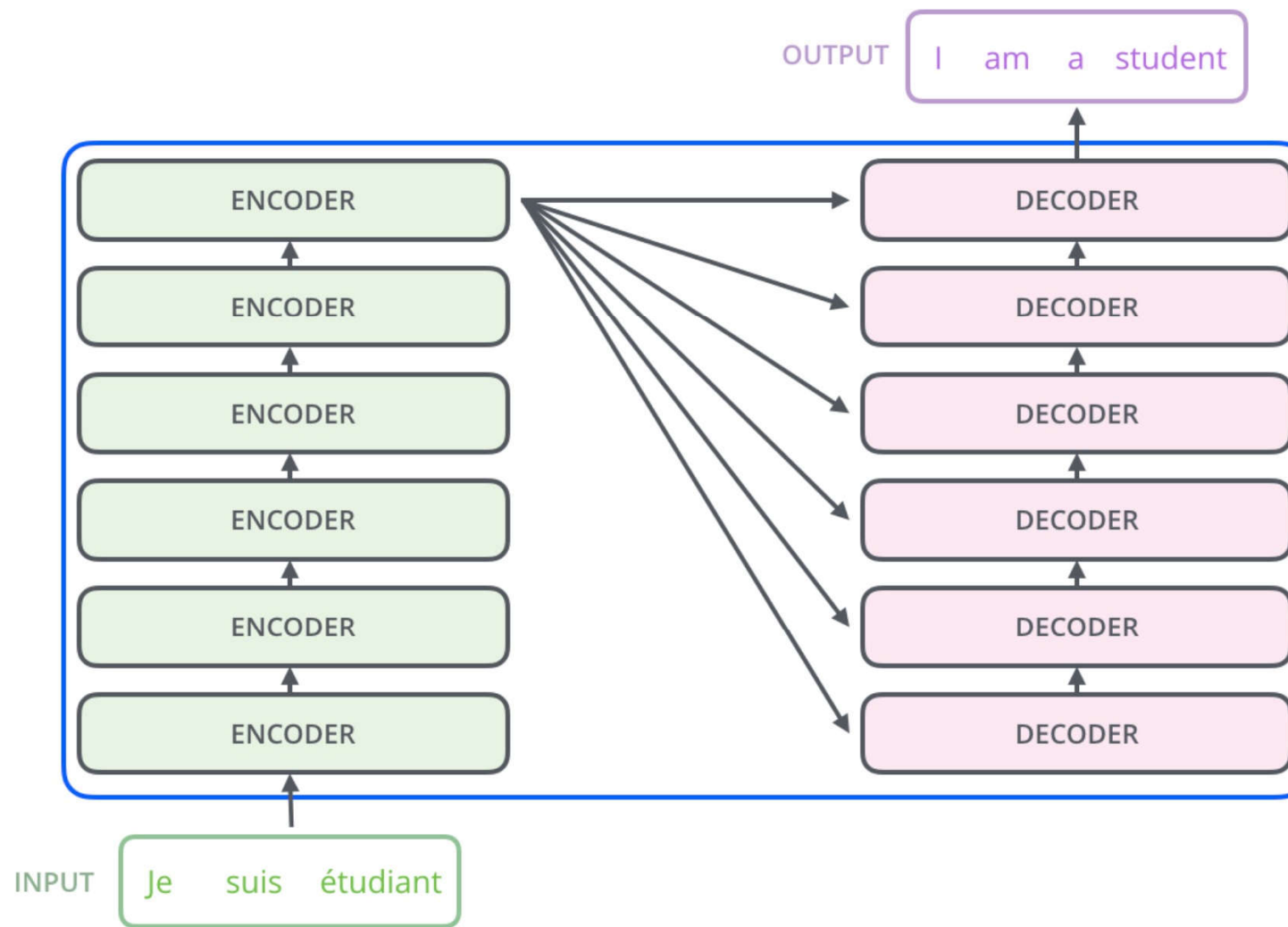


BERT

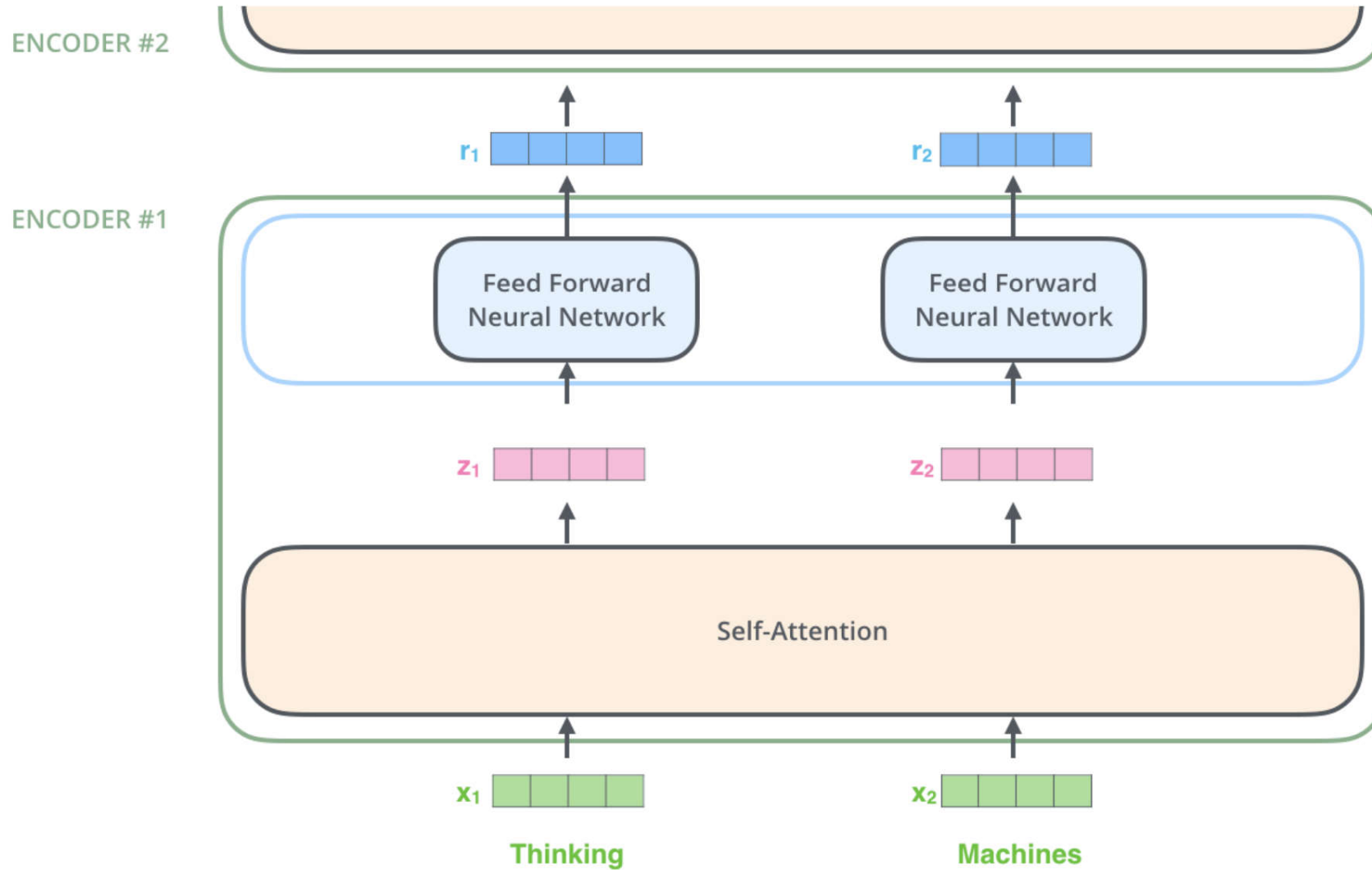
The output of each encoder layer can be used to represent the feature for that token



Архитектура Transformer



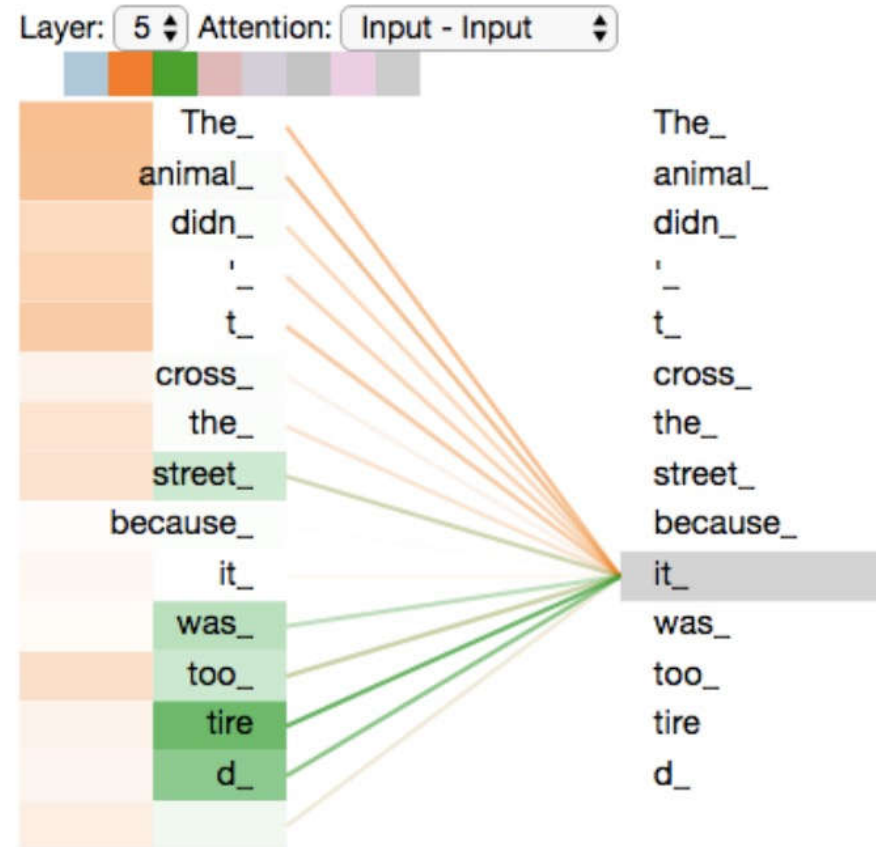
Кодировщик



Внутреннее внимание

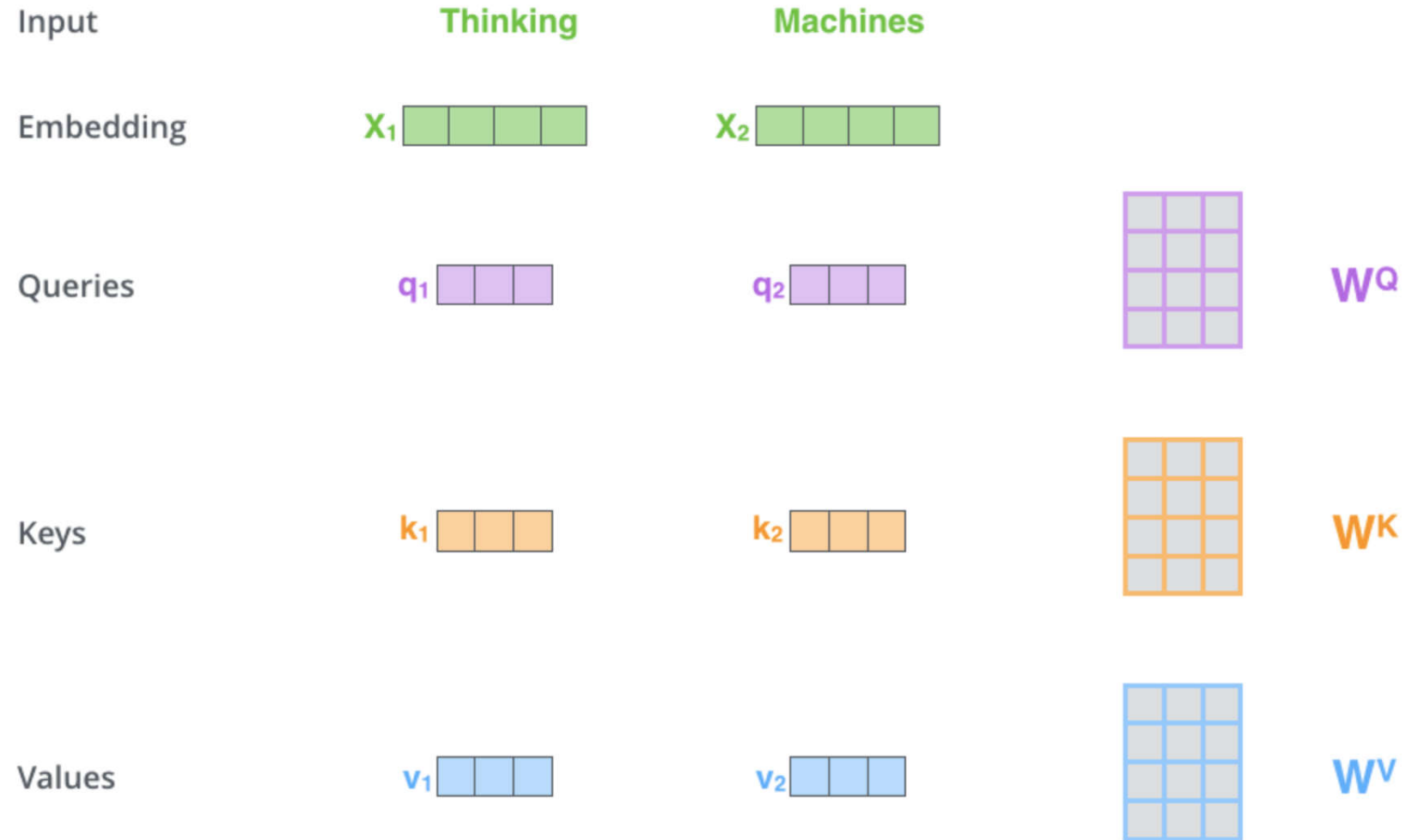
As we encode the word "it", one attention head is focusing most on "the animal", while another is focusing on "tired" -- in a sense.

The model's representation of the word "it" bakes in some of the representation of both "animal" and "tired"

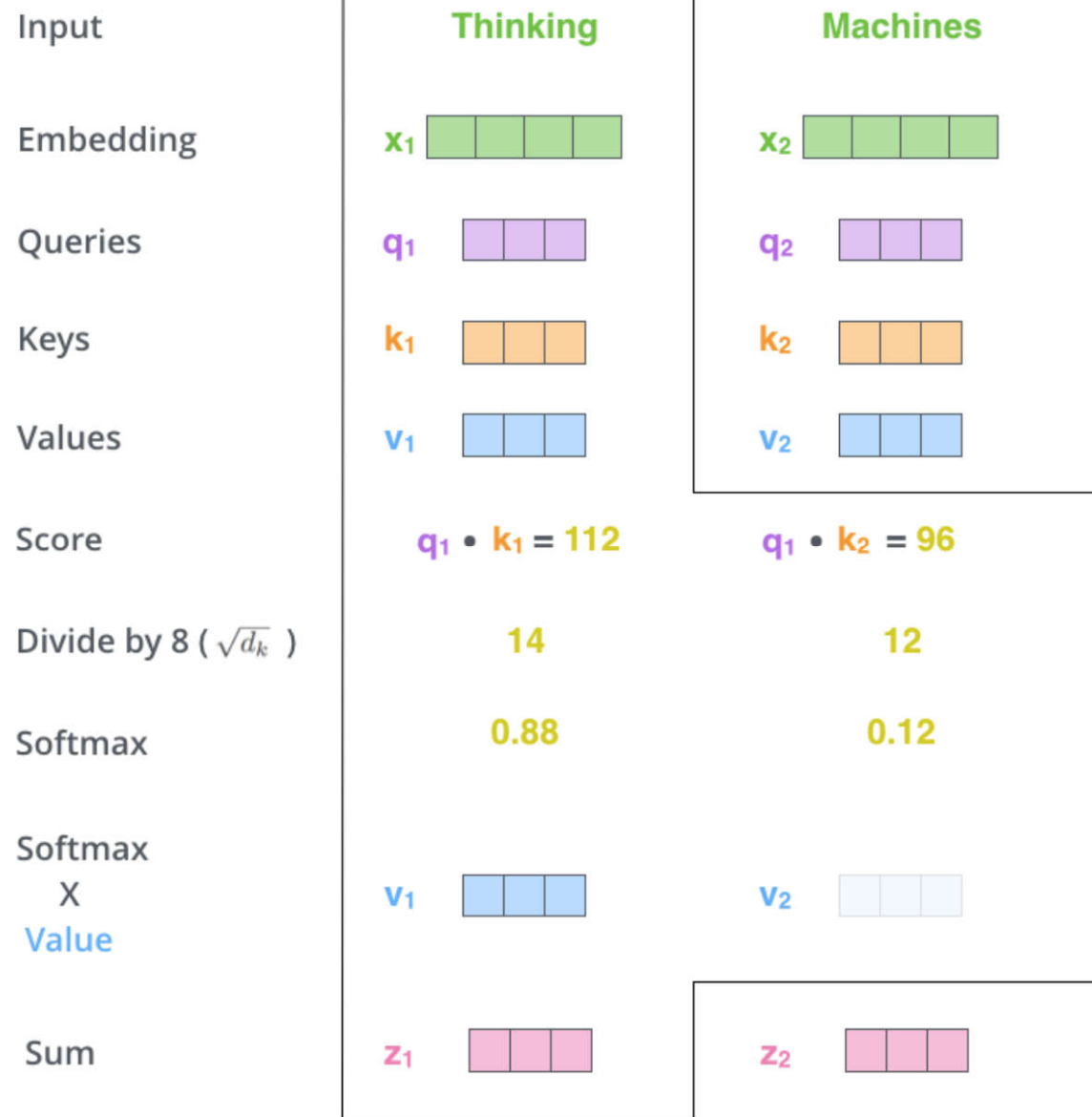


Внутреннее внимание

Multiplying x_1 by the W^Q weight matrix produces q_1 , the "query" vector associated with that word. We end up creating a "query", a "key", and a "value" projection of each word in the input sentence.

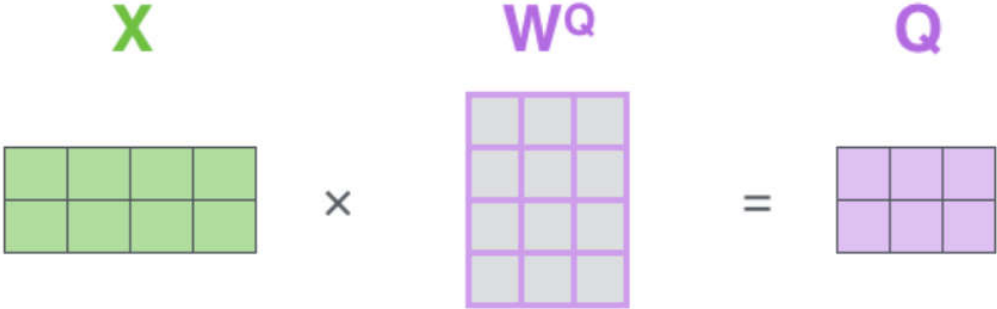


Value



Several words

Every row in the **X** matrix corresponds to a word in the input sentence



Calculation

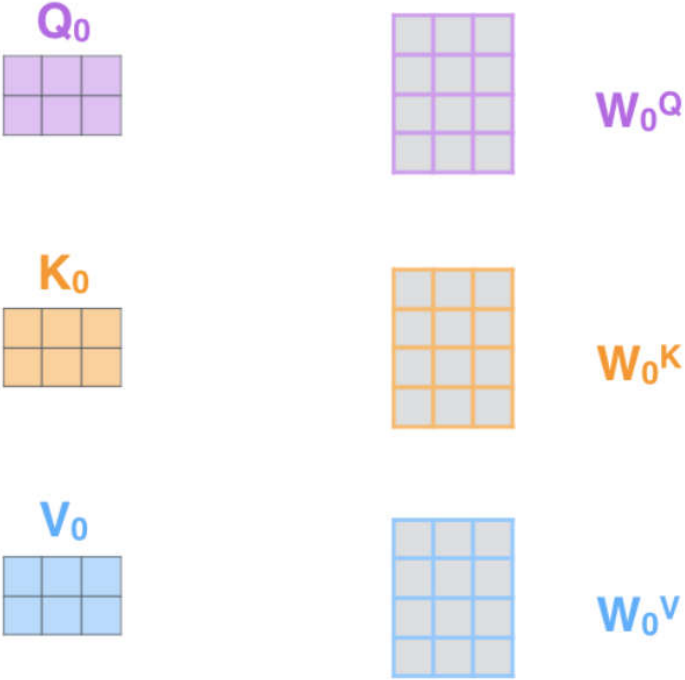
$$\text{softmax} \left(\frac{\begin{matrix} \mathbf{Q} \\ \begin{array}{|c|c|c|} \hline \square & \square & \square \\ \hline \square & \square & \square \\ \hline \end{array} \end{matrix} \times \begin{matrix} \mathbf{K}^T \\ \begin{array}{|c|c|} \hline \square & \square \\ \hline \square & \square \\ \hline \square & \square \\ \hline \end{array} \end{matrix} \right) \begin{matrix} \mathbf{V} \\ \begin{array}{|c|c|c|} \hline \square & \square & \square \\ \hline \square & \square & \square \\ \hline \end{array} \end{matrix}$$
$$= \begin{matrix} \mathbf{Z} \\ \begin{array}{|c|c|c|} \hline \square & \square & \square \\ \hline \square & \square & \square \\ \hline \end{array} \end{matrix}$$

The self-attention calculation in matrix form

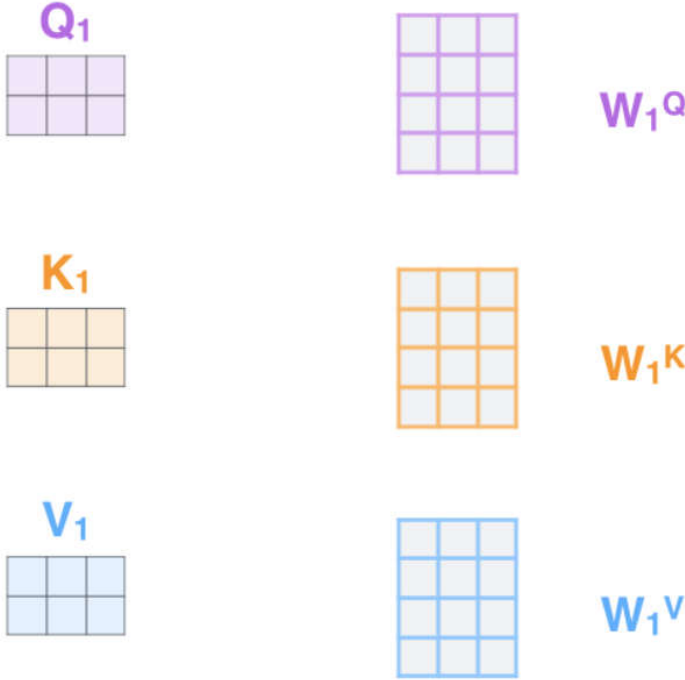
Attention heads



ATTENTION HEAD #0



ATTENTION HEAD #1



With multi-headed attention, we maintain separate Q/K/V weight matrices for each head resulting in different Q/K/V matrices. As we did before, we multiply X by the WQ/WK/WV matrices to produce Q/K/V matrices

Attention result

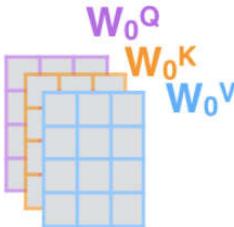
1) This is our input sentence*

Thinking
Machines

2) We embed each word*



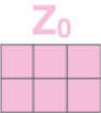
3) Split into 8 heads. We multiply X or R with weight matrices



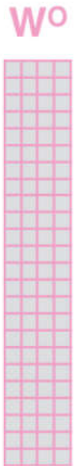
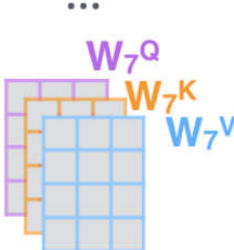
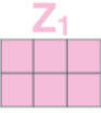
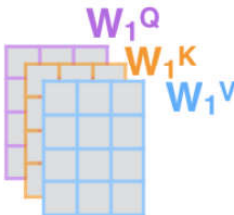
4) Calculate attention using the resulting $Q/K/V$ matrices



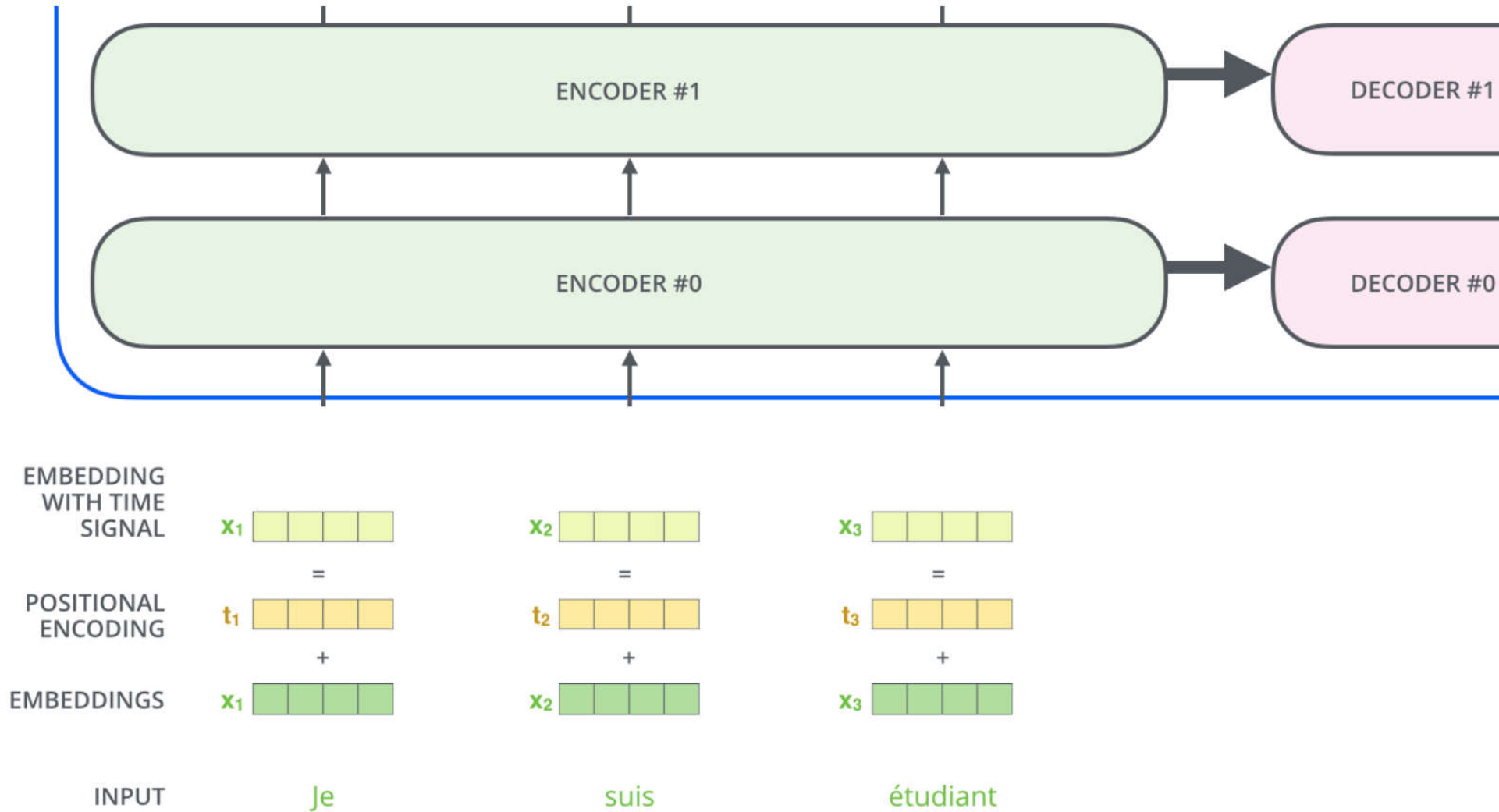
5) Concatenate the resulting Z matrices, then multiply with weight matrix W^O to produce the output of the layer



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



Positions

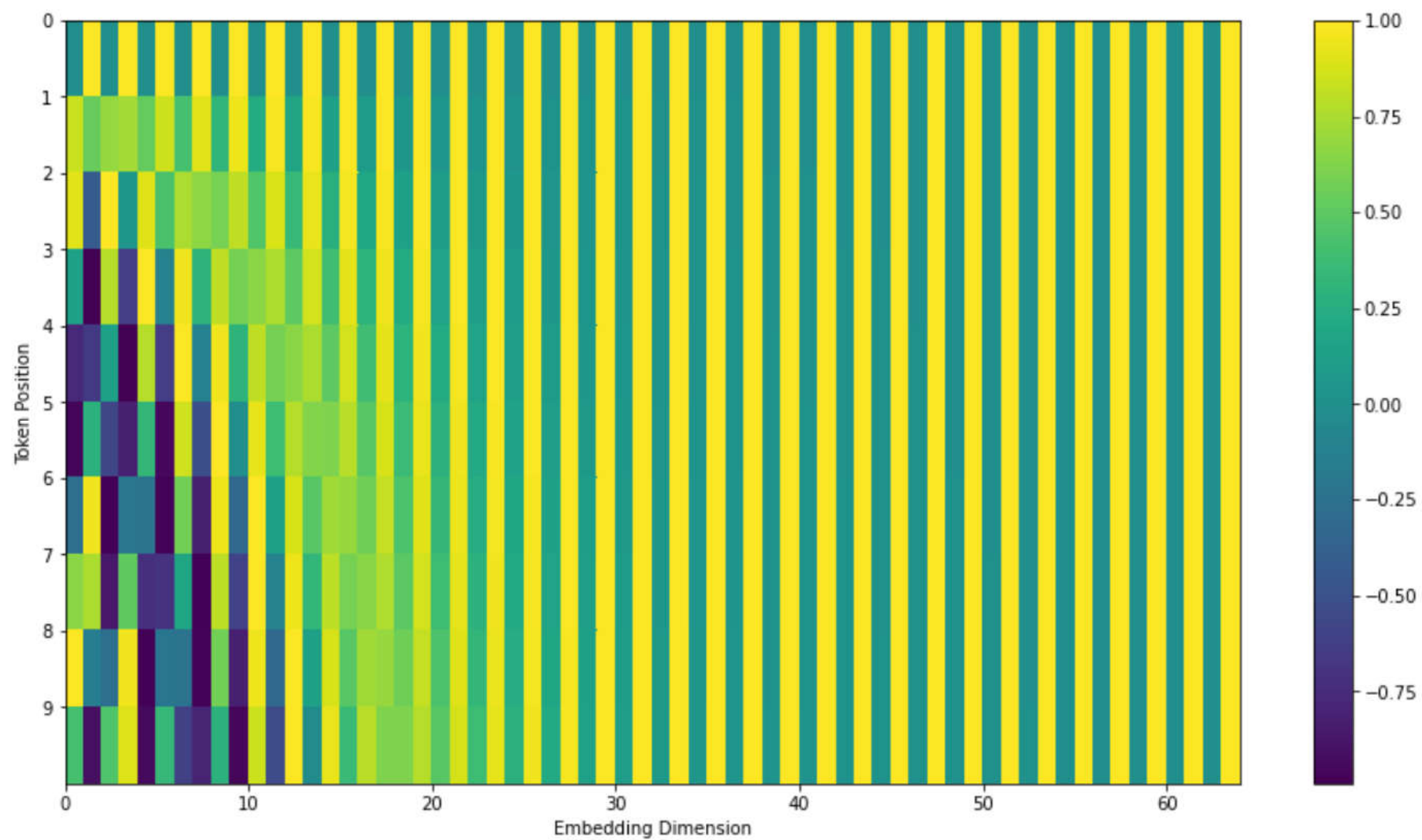


Positional encoding

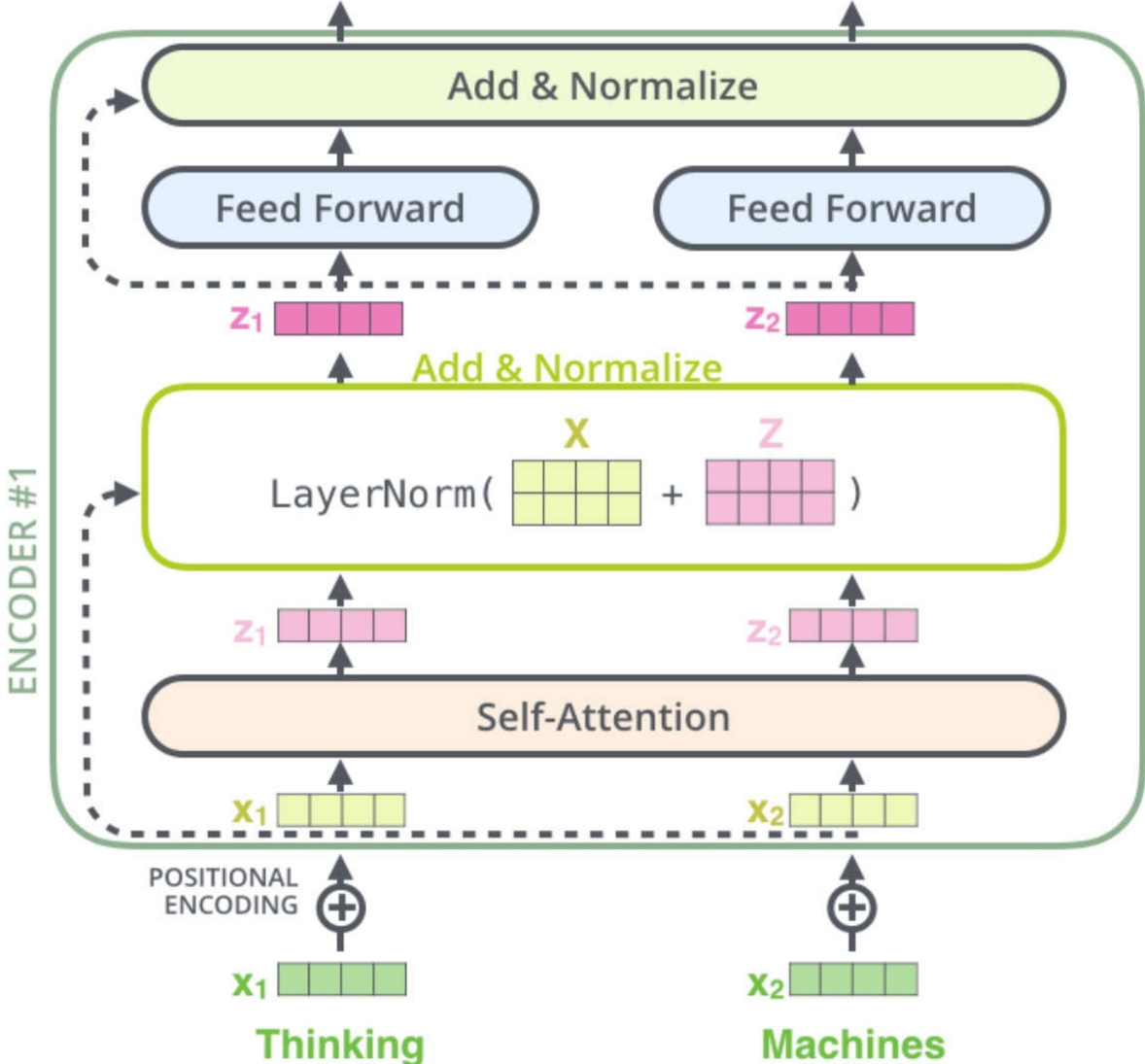
A real example of positional encoding with a toy embedding size of 4



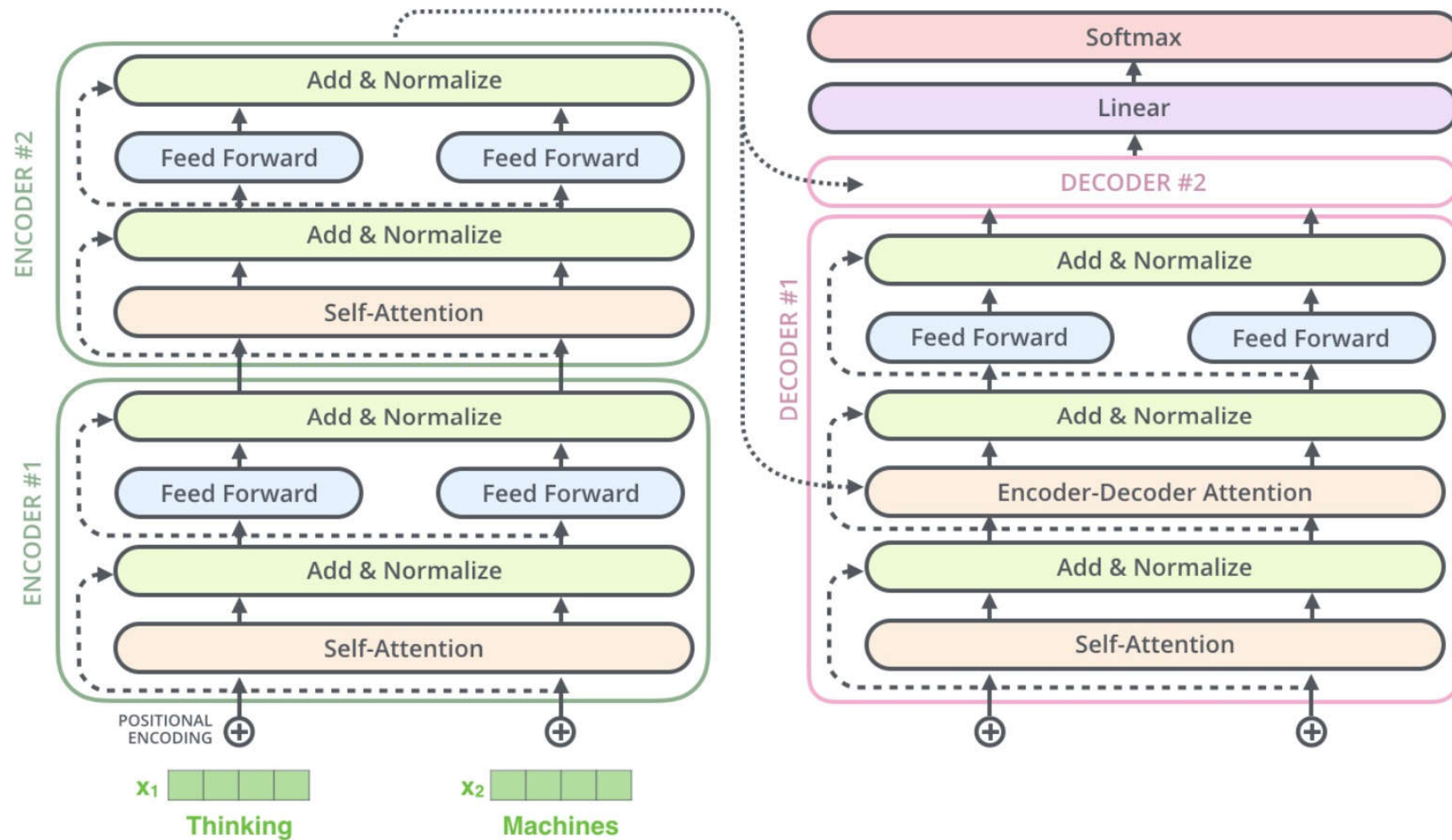
Positional encoding



Add and Normalize



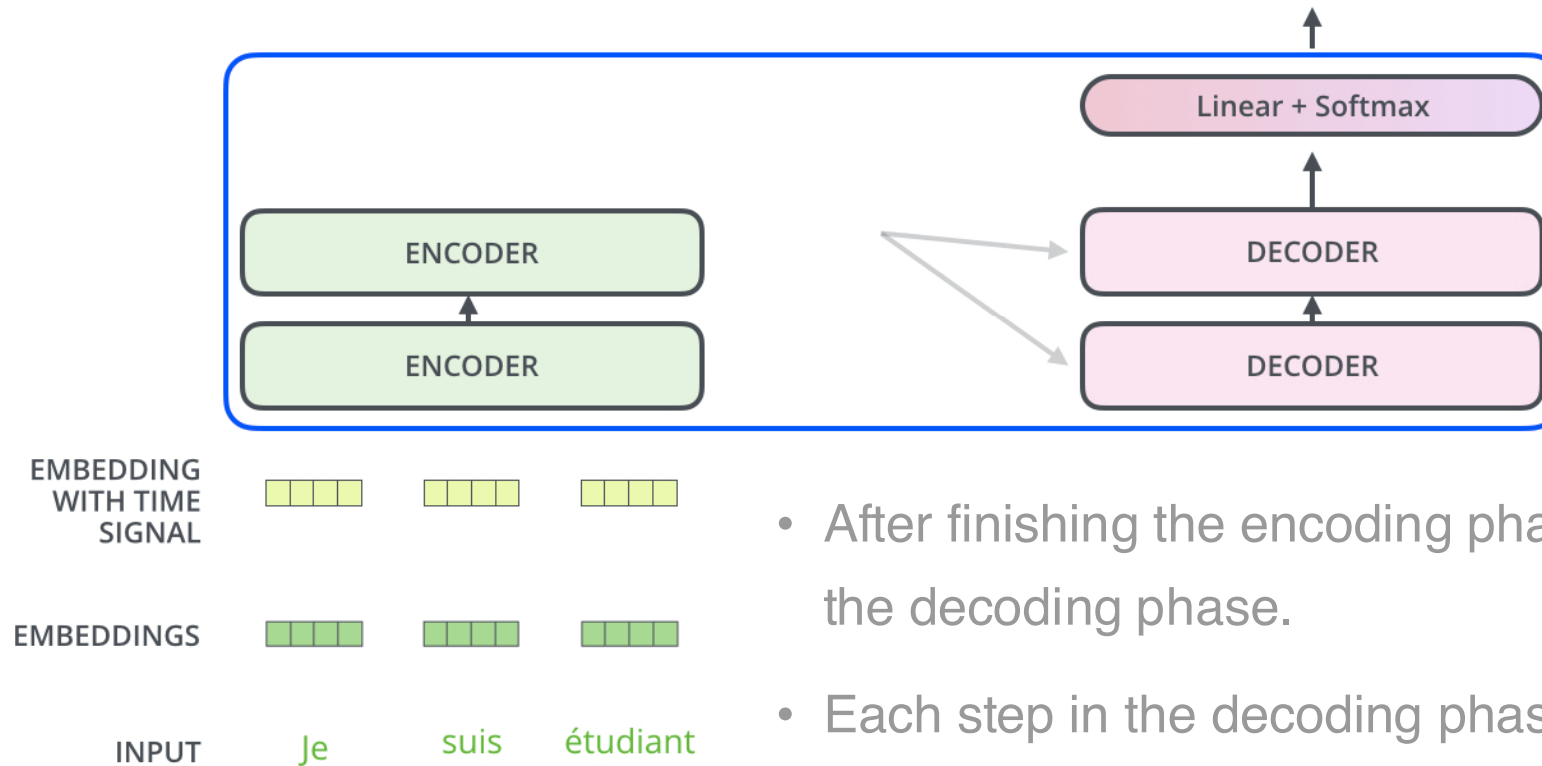
Encoder-decoder connection



Decoder

Decoding time step: ① 2 3 4 5 6

OUTPUT



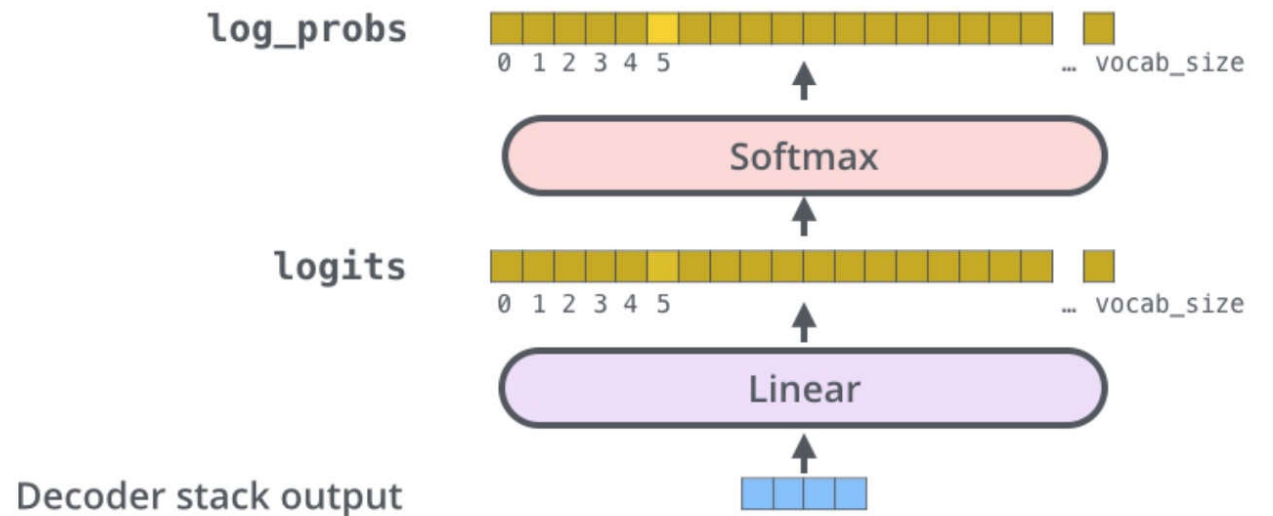
- After finishing the encoding phase, we begin the decoding phase.
- Each step in the decoding phase outputs an element from the output sequence (the English translation sentence in this case)

ИТОГОВОЕ СЛОВО

- Let's assume that our model knows 10,000 unique English words (our model's "output vocabulary") that it's learned from its training dataset.
- This would make the logits vector 10,000 cells wide – each cell corresponding to the score of a unique word.
- That is how we interpret the output of the model followed by the Linear layer.

Which word in our vocabulary is associated with this index?

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



Model output

