

Трансформеры, модель BERT: архитектура и применение

ДОКЛАД ПОДГОТОВИЛ
АНДРЕЙ ЛЕБЕДЕВ

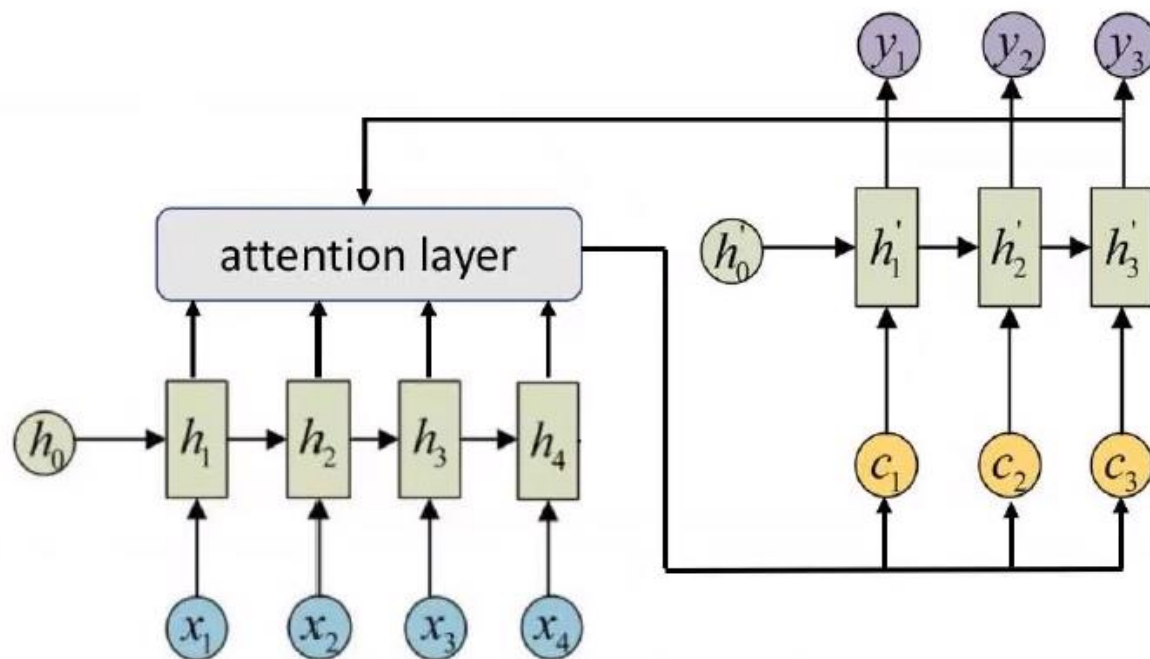


Содержание

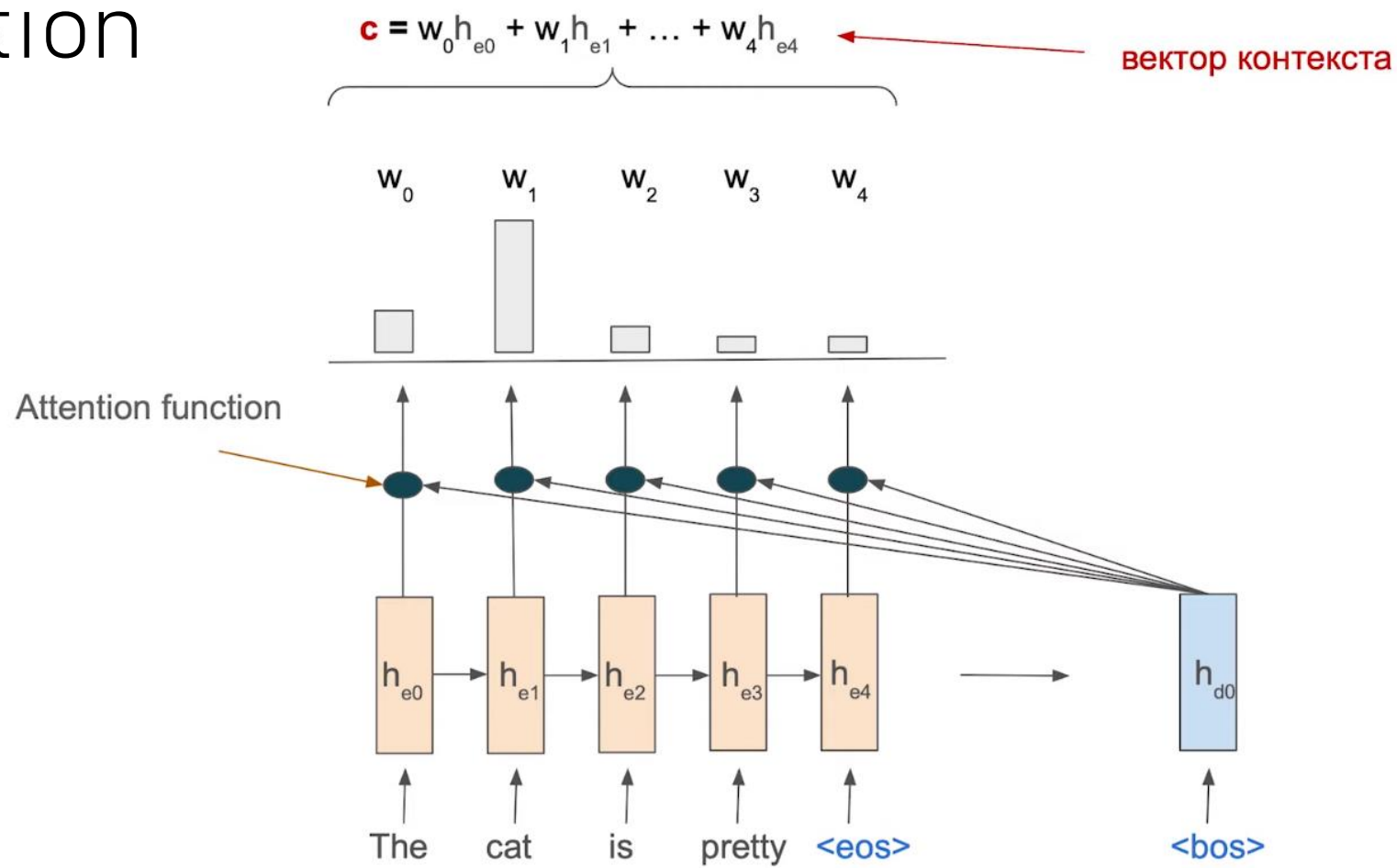
1. Введение и предпосылки к появлению механизма внимания
2. Механизм внимания: attention, self-attention, multi-head attention
3. Архитектура трансформеров
4. Модель BERT: история, устройство и назначение
5. Прикладные задачи, решаемые моделью BERT

Появление механизма внимания в RNN

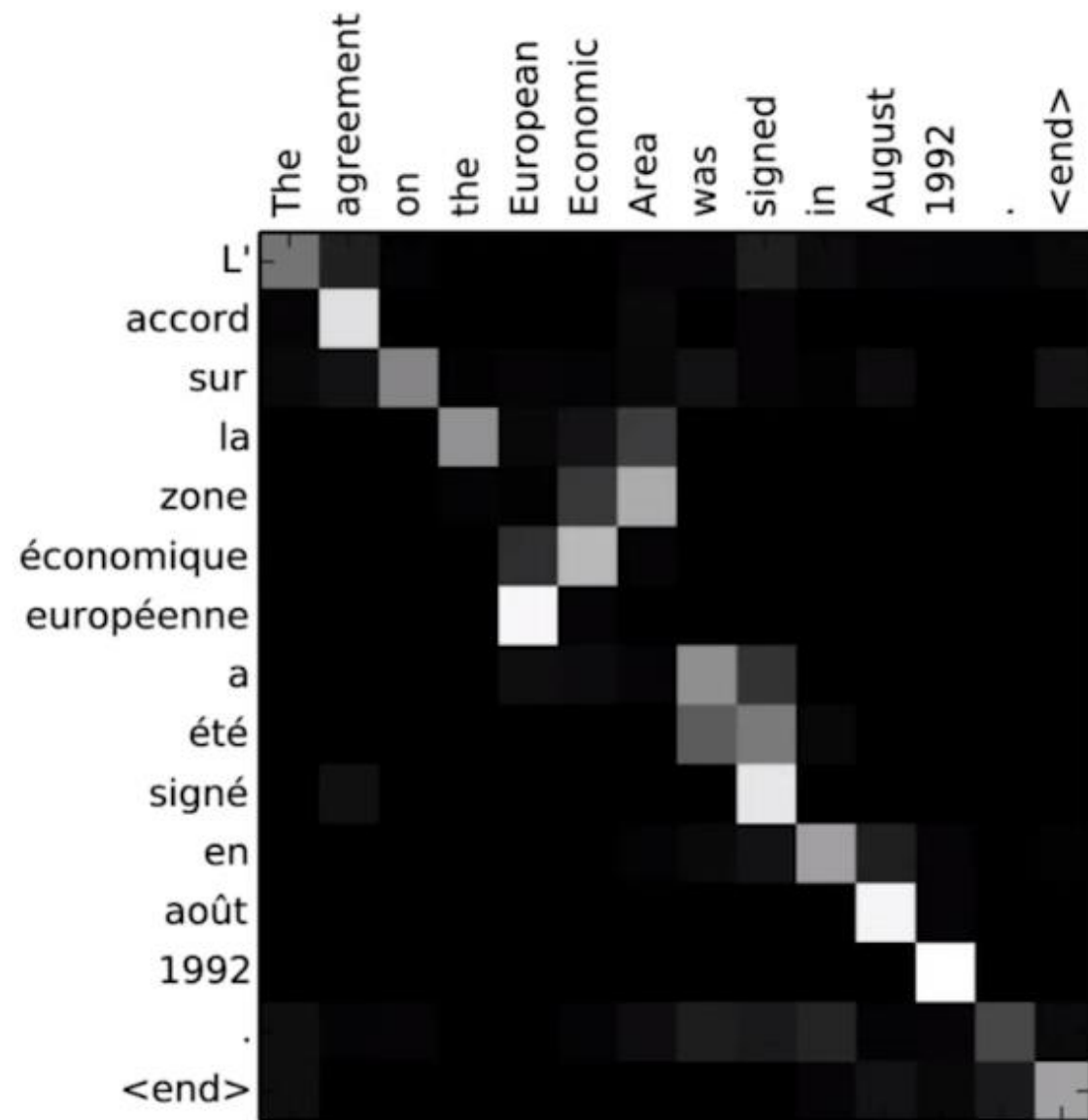
- The animal didn't cross the street because it was too tired.
- The animal didn't cross the street because it was too wide.



Attention

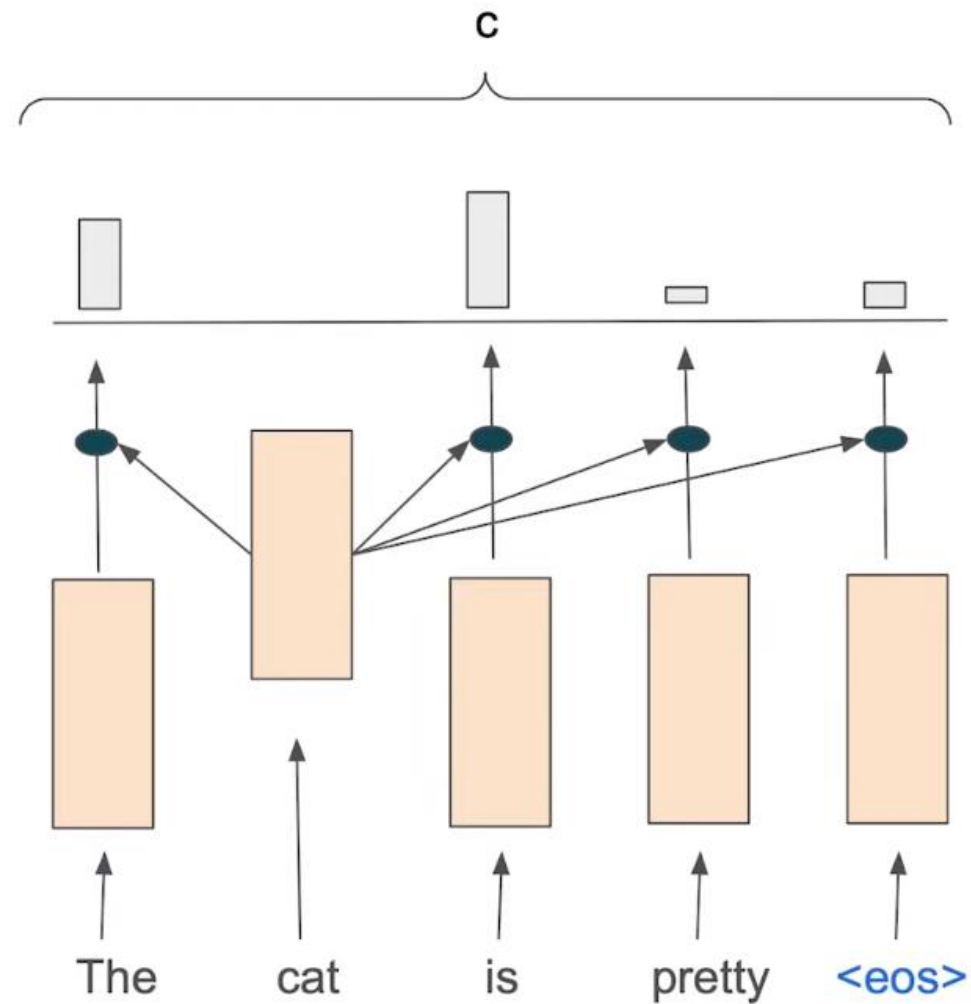


Пример работы attention



Self-attention

- A man in a nice cowboy hat and a parrot on his shoulder has entered the bar.
- Человек в красивой ковбойской шляпе и попугаем на плече ... [зашел? зашла? зашло?]





Multi-head attention

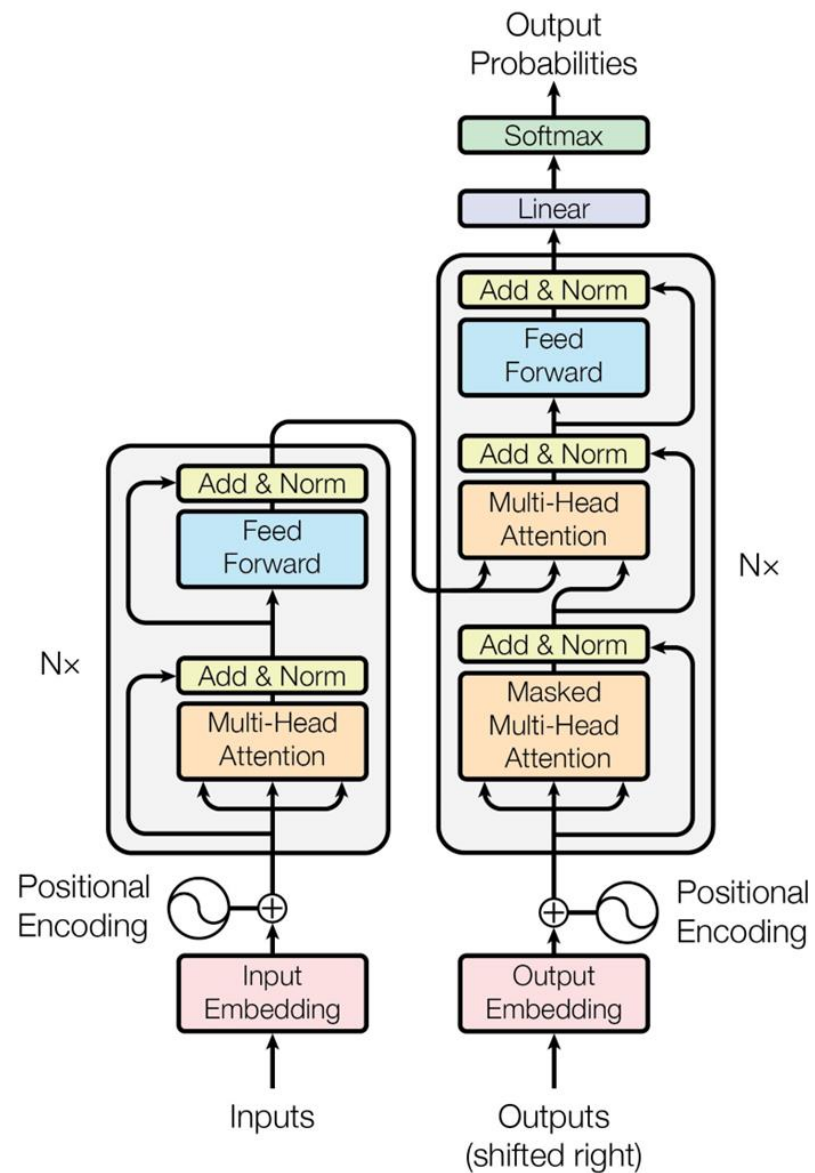
- Если в модели несколько механизмов внимания и само-внимания, такой механизм называется multi-head attention. При этом каждый из них может начать выделять определенную информацию, например, согласование по падежу/лицу и т. д.



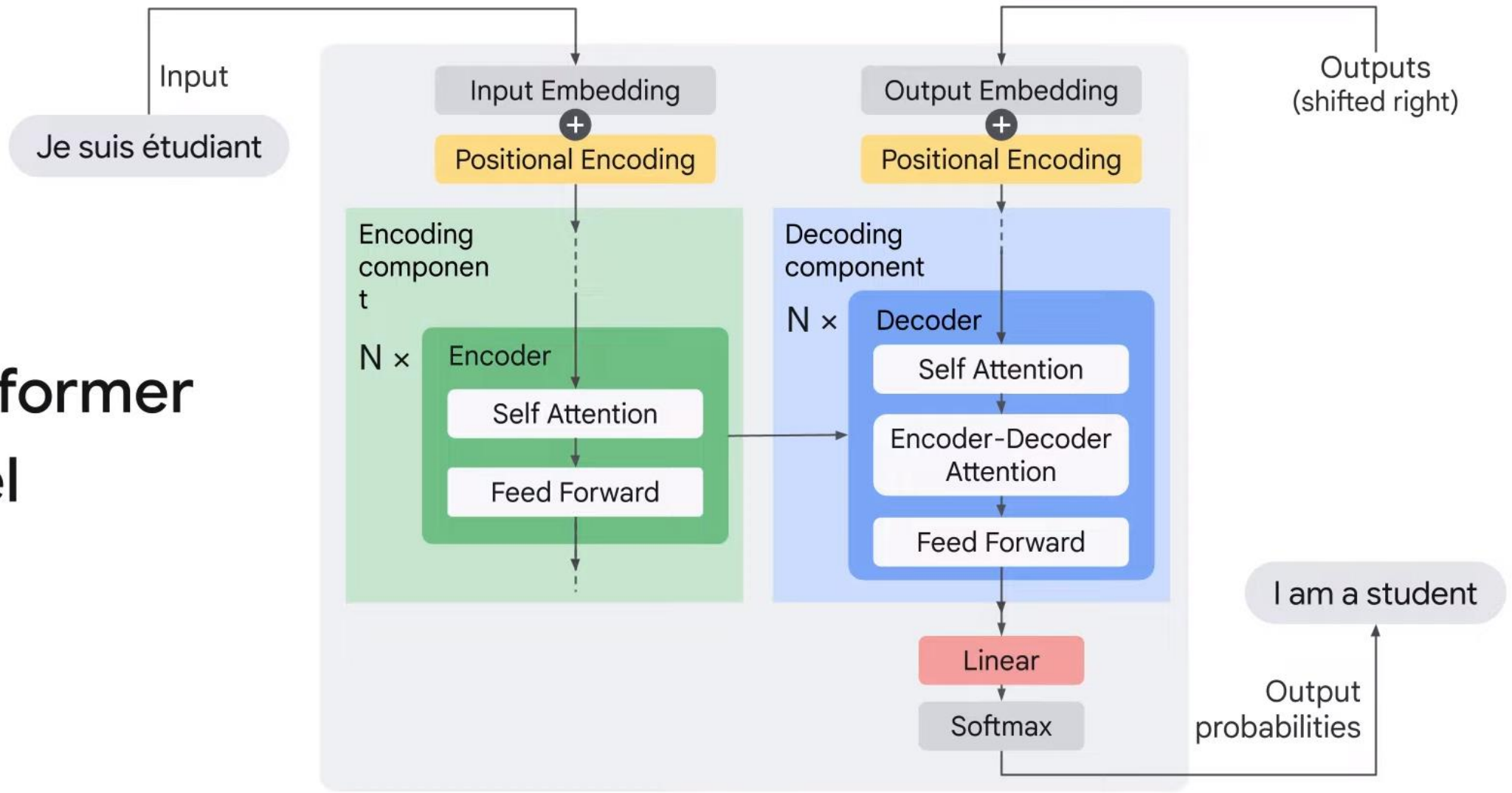
Трансформеры

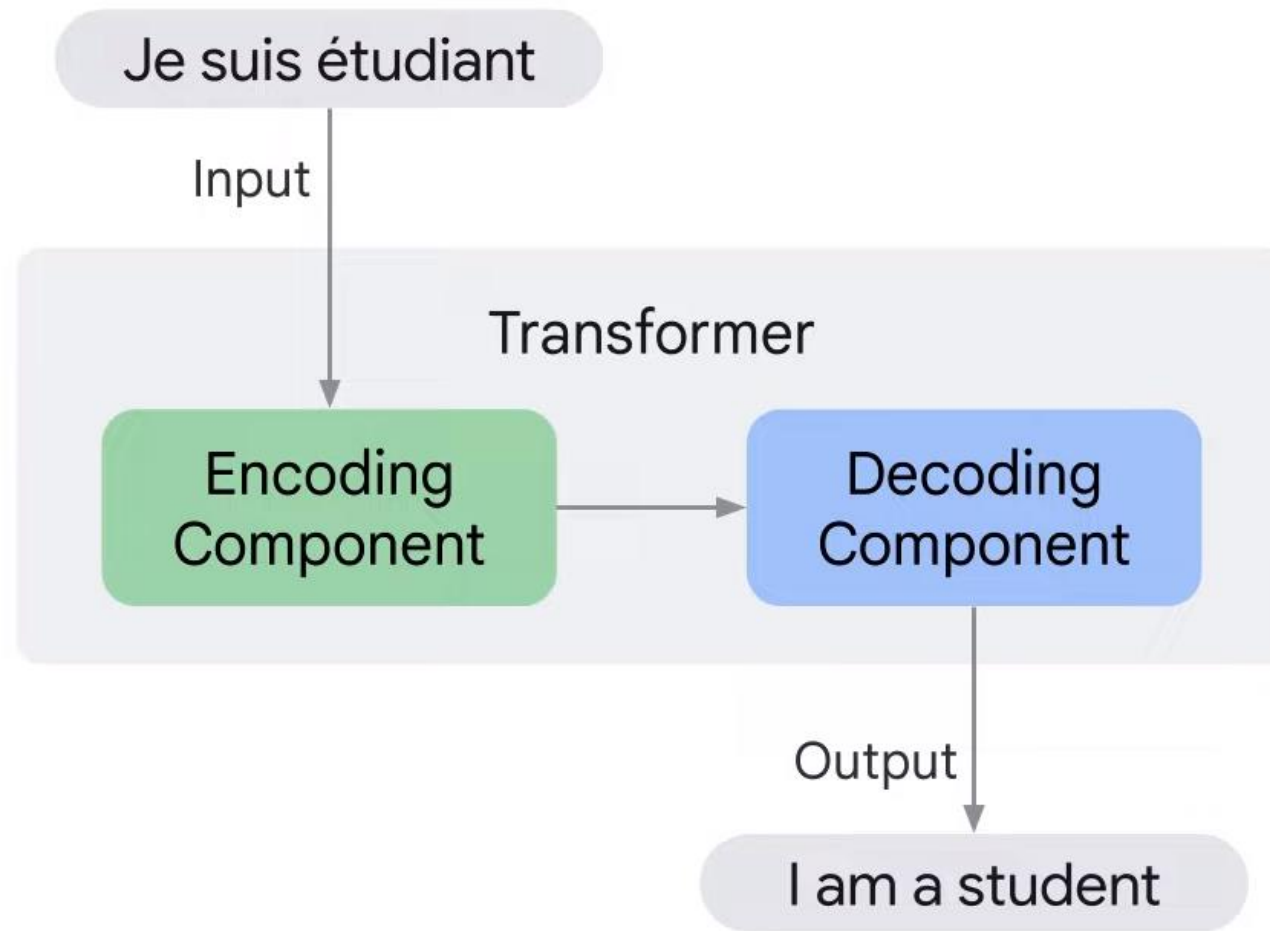
Трансформер

Трансформер - нейросетевая архитектура на основе внимания и полносвязных слоев.

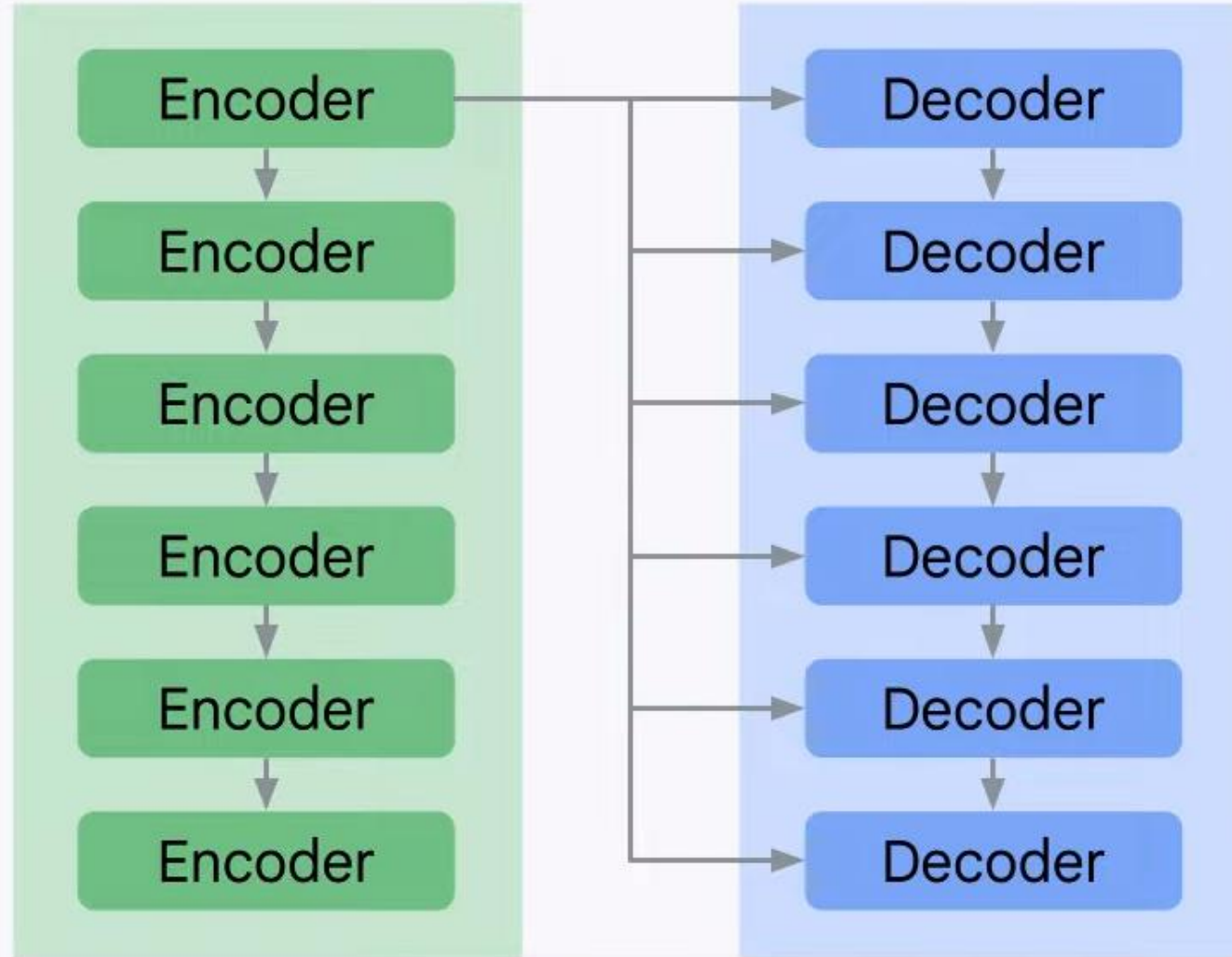


Transformer model



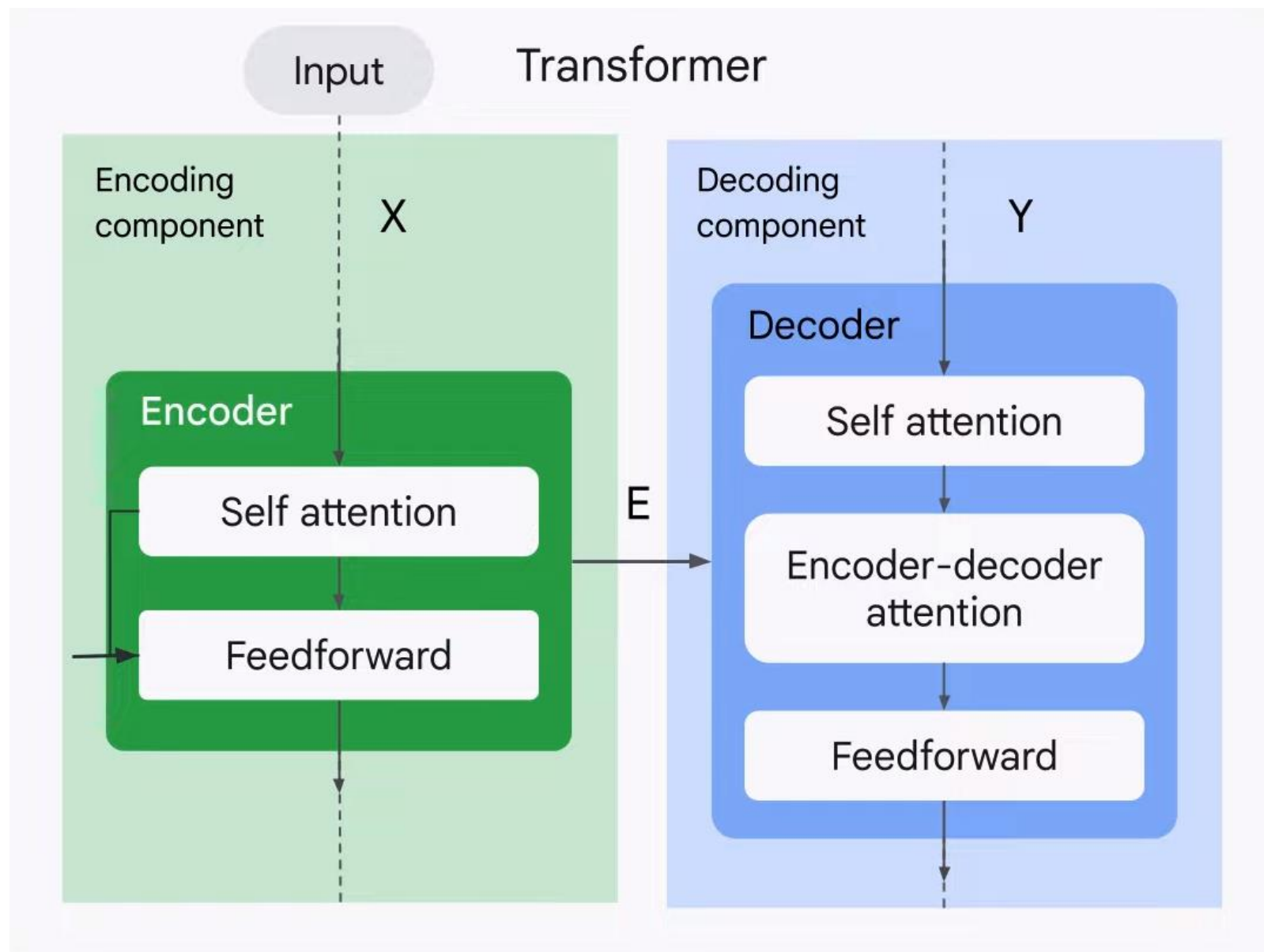


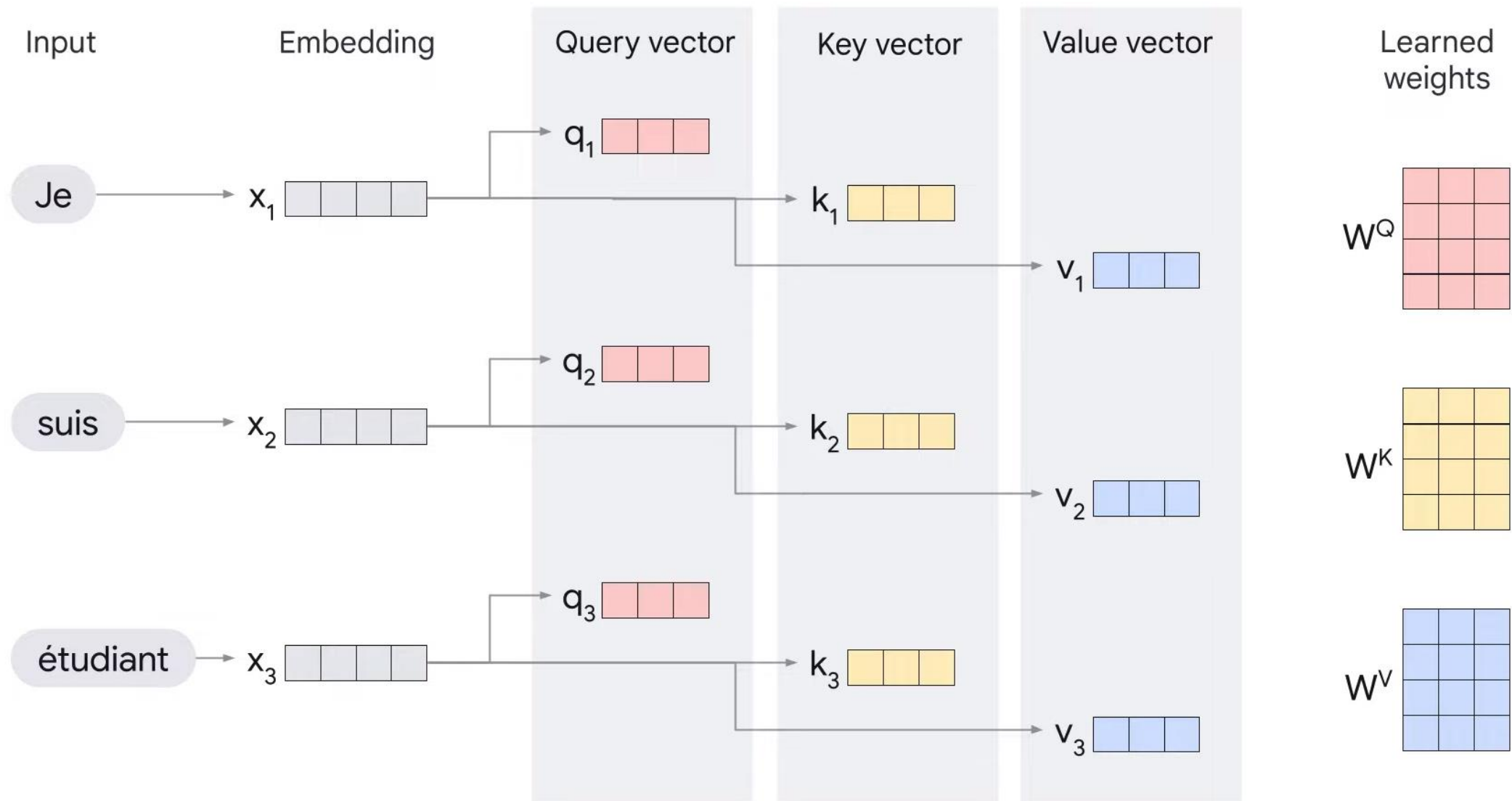
Transformer



Feedforward – полносвязная сеть, позволяющая сделать обучаемое преобразование выходных данных со слоя внимания и подготовить их к следующему этапу кодирования:

$$h_i''' = W_2 \text{ReLU}(W_1 h_i'' + b_1) + b_2$$





Input

Je

suis

étudiant

Embedding

x

x_1				
x_2				
x_3				

Query vector

Q

q_1			
q_2			
q_3			

Key vector

K

k_1			
k_2			
k_3			

Value vector

V

v_1			
v_2			
v_3			

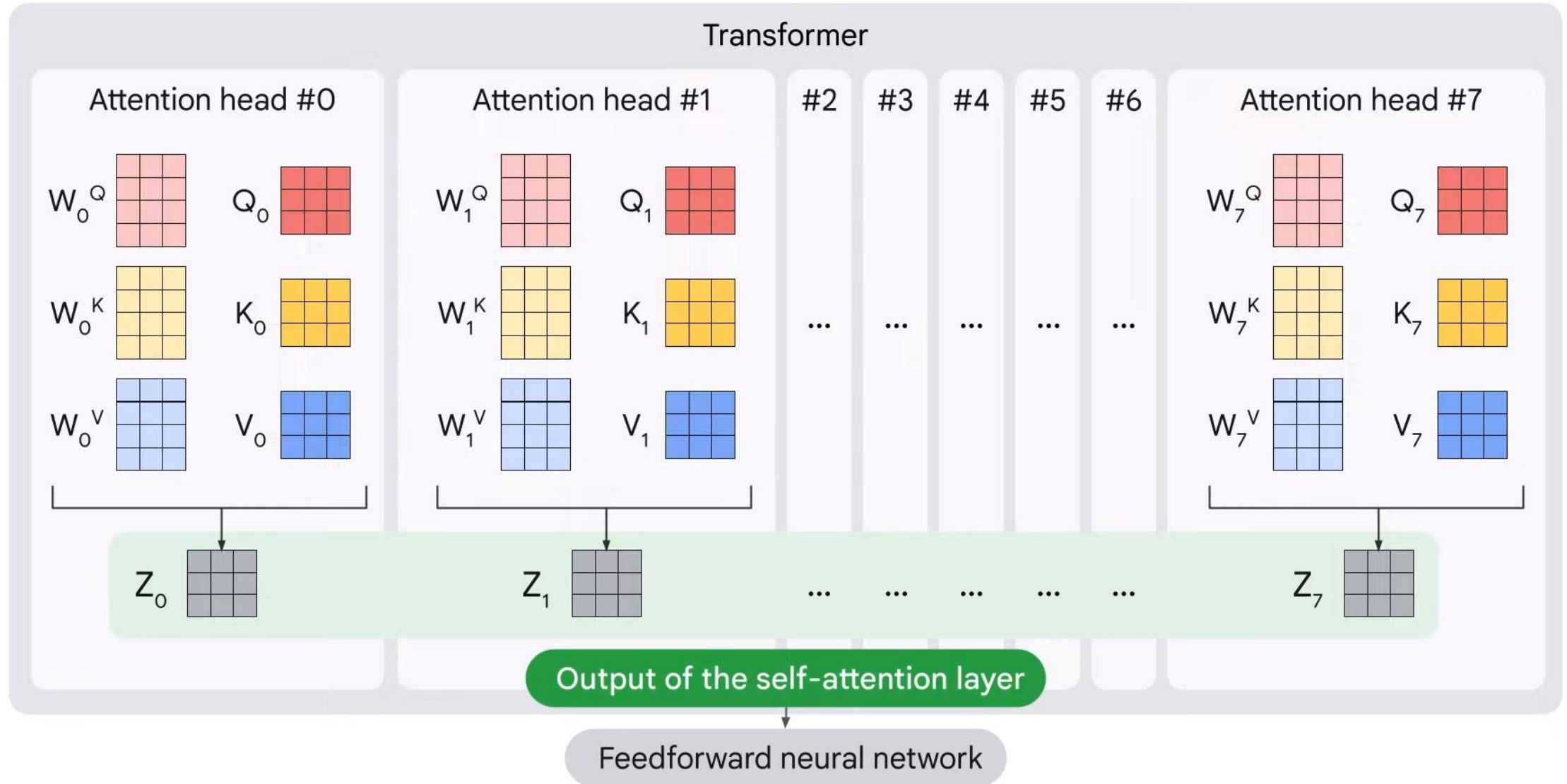
Learned weights

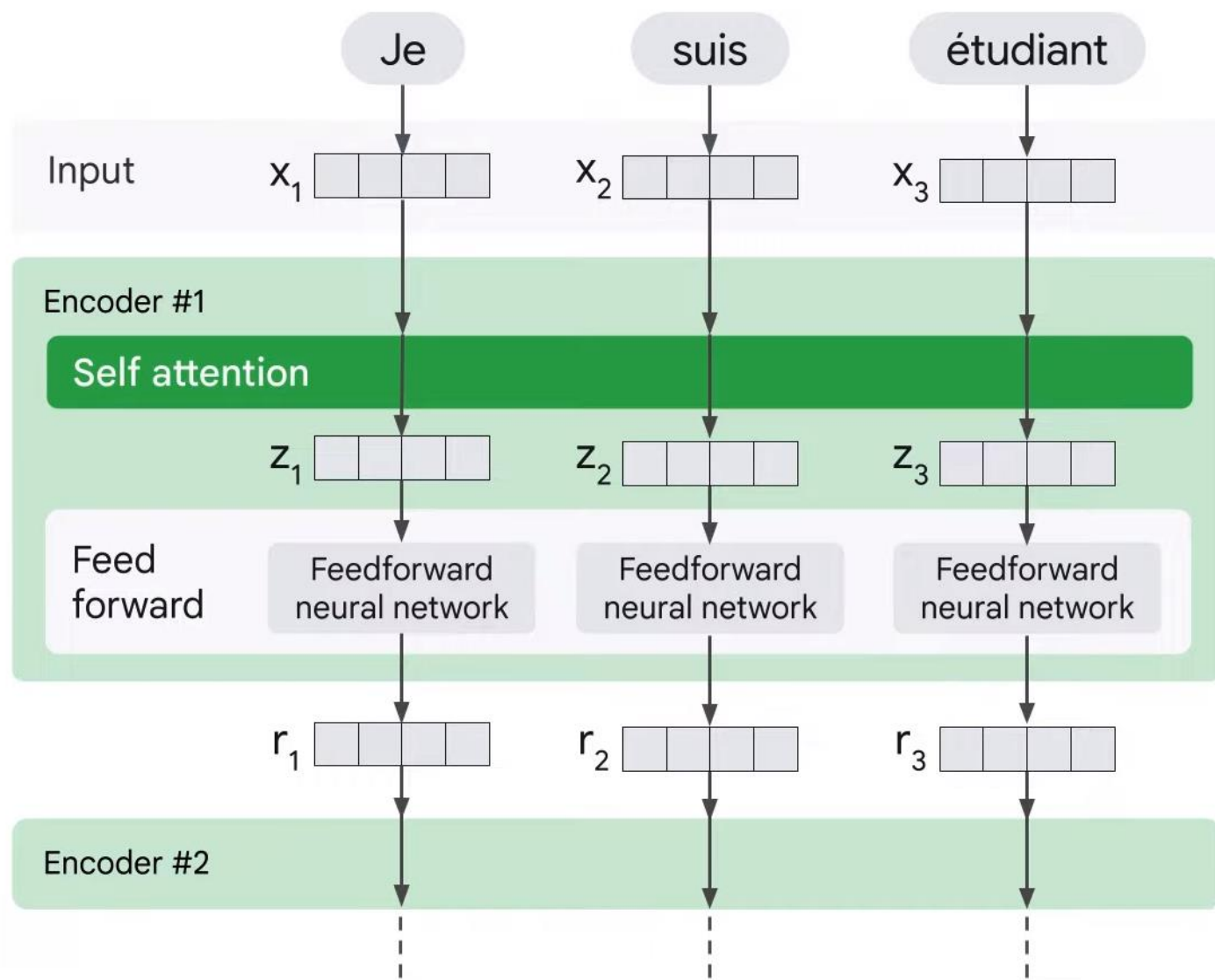
W^Q

W^K

W^V

$$\text{softmax} \left[\frac{\begin{matrix} Q & K^T \\ \begin{matrix} \text{3x3 red grid} \end{matrix} & \begin{matrix} \text{3x3 yellow grid} \end{matrix} \\ \times \end{matrix} \right] \frac{1}{\sqrt{d_k}} = \begin{matrix} V \\ \begin{matrix} \text{3x3 blue grid} \end{matrix} \end{matrix} = \begin{matrix} Z \\ \begin{matrix} \text{3x3 gray grid} \end{matrix} \end{matrix}$$





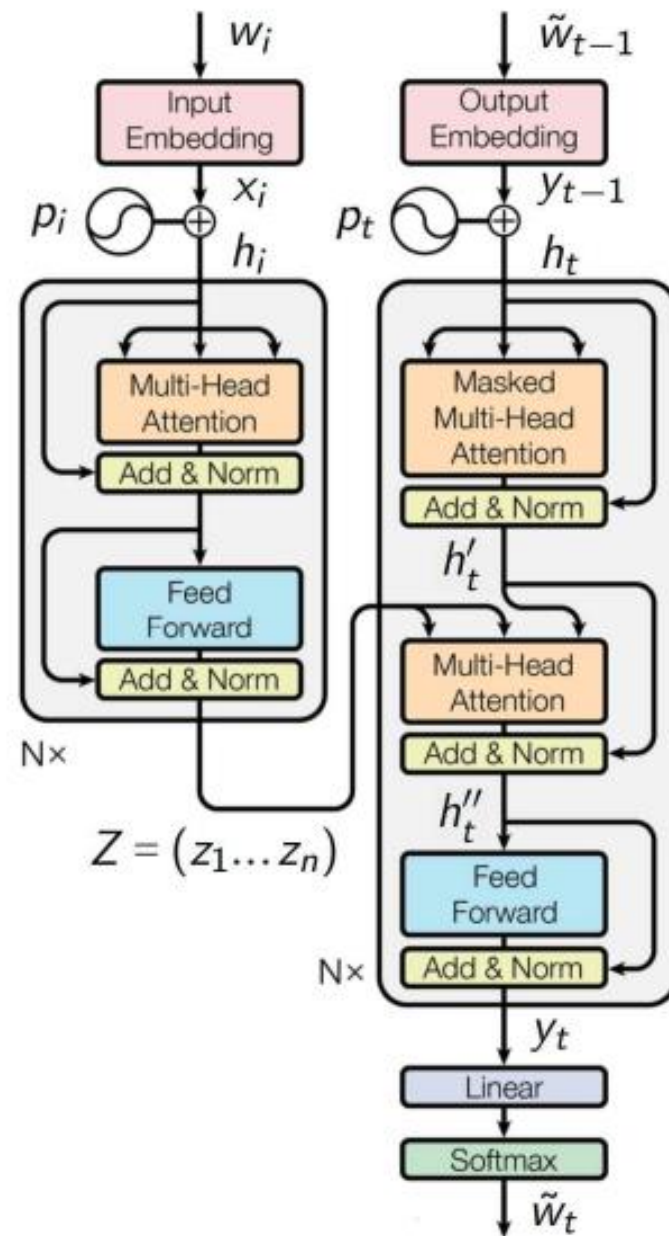


Процесс получения эмбедингов z на выходе из энкодера

- Ввод предложения на естественном языке
- Векторизация каждого слова
- Выполнение multi-head attention и умножение векторизованных слов на соответствующие матрицы весов
- Вычисление функции внимания, используя полученные матрицы QKV
- Конкатенация матриц для получения выходной матрицы, которая имеет тот же размер, что и конечная матрица

Diagram illustrating the multiplication of a matrix Z (shaded gray) by a matrix W^0 (shaded green) to produce a matrix Z (shaded gray). The matrix Z is partitioned into blocks $Z_0, Z_1, Z_2, \dots, Z_7$. The matrix W^0 is partitioned into blocks $W^0_0, W^0_1, W^0_2, \dots, W^0_7$. The result Z is partitioned into blocks $Z_0, Z_1, Z_2, \dots, Z_7$.

Резюмирование



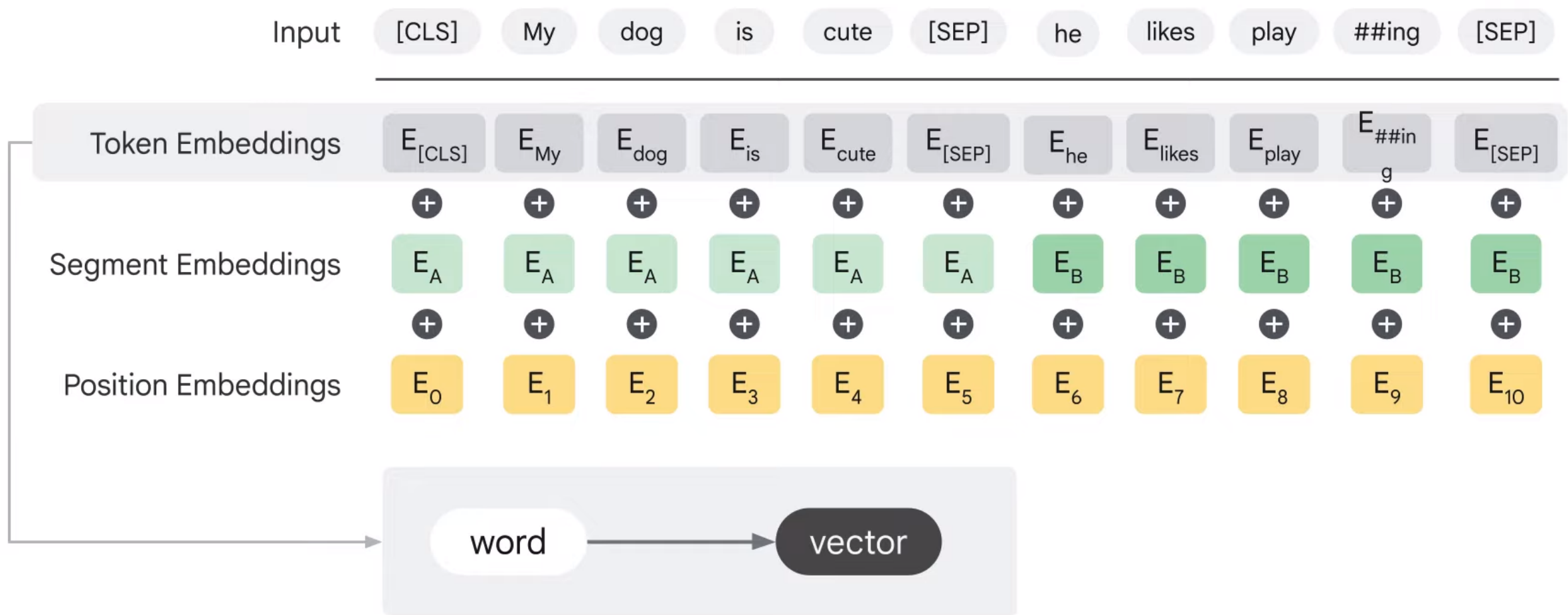


BERT (Bidirectional Encoder Representations from Transformers)



Общие положения

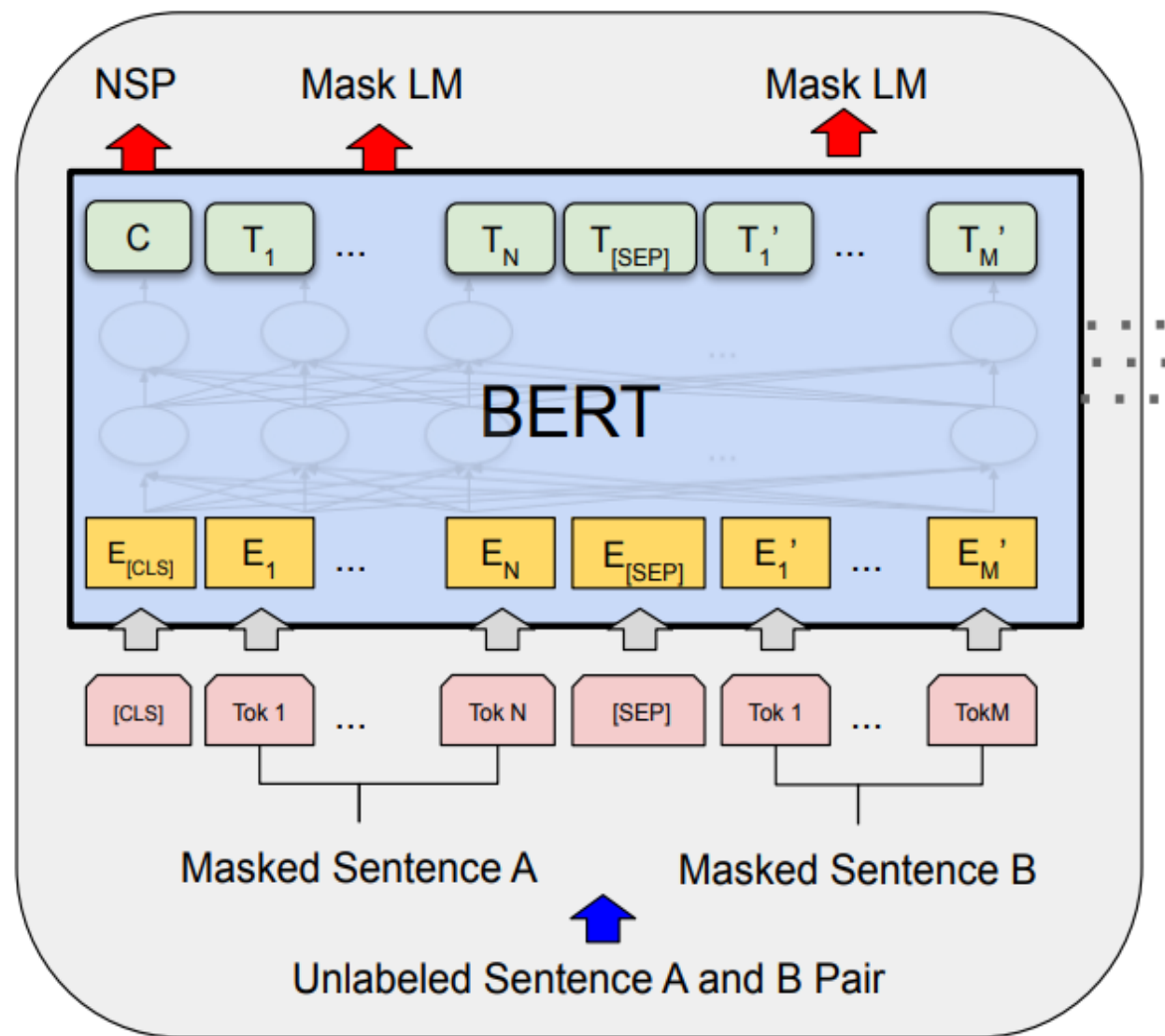
- Компания Google разработала языковую модель BERT в 2018 году. Она состоит из простого набора блоков-трансформеров, который был предварительно обучен на большом корпусе текстов общего характера.
- Отличается двунаправленностью внимания: это значит, что при обработке входной последовательности все токены могут использовать информацию друг о друге.





Предобучение

- BERT обучается одновременно на двух задачах — предсказание следующего предложения и генерация пропущенного токена.
- На вход BERT подаются токенизированные пары предложений, в которых некоторые токены скрыты. Таким образом сеть обучается глубокому двунаправленному представлению языка, учится понимать контекст предложения.
- Задача же предсказания следующего предложения есть задача бинарной классификации — является ли второе предложение продолжением первого. Благодаря ей сеть можно обучить различать наличие связи между предложениями в тексте.





Какие задачи решает BERT?

1. Поисковый алгоритм Google
2. Классификация текстов
3. Ответы на вопросы
4. Генерация текста
5. Суммаризация текста



Спасибо за внимание!