

Трансформер. Часть 1

Корнет Мария Евгеньевна
старший преподаватель кафедры
Инженерная кибернетика

План

✓ Модели внимания RNN, seq2seq

↑
✓ Self-attention

✓ Трансформеры

Напоминание: RNN

x_t - входной вектор в момент времени

$t = 1, \dots, T$

h_t - вектор скрытого состояния в момент t

y_t - выходной вектор в момент времени t

$$h_t = \sigma(Ux_t + Wh_{t-1} + b_h)$$

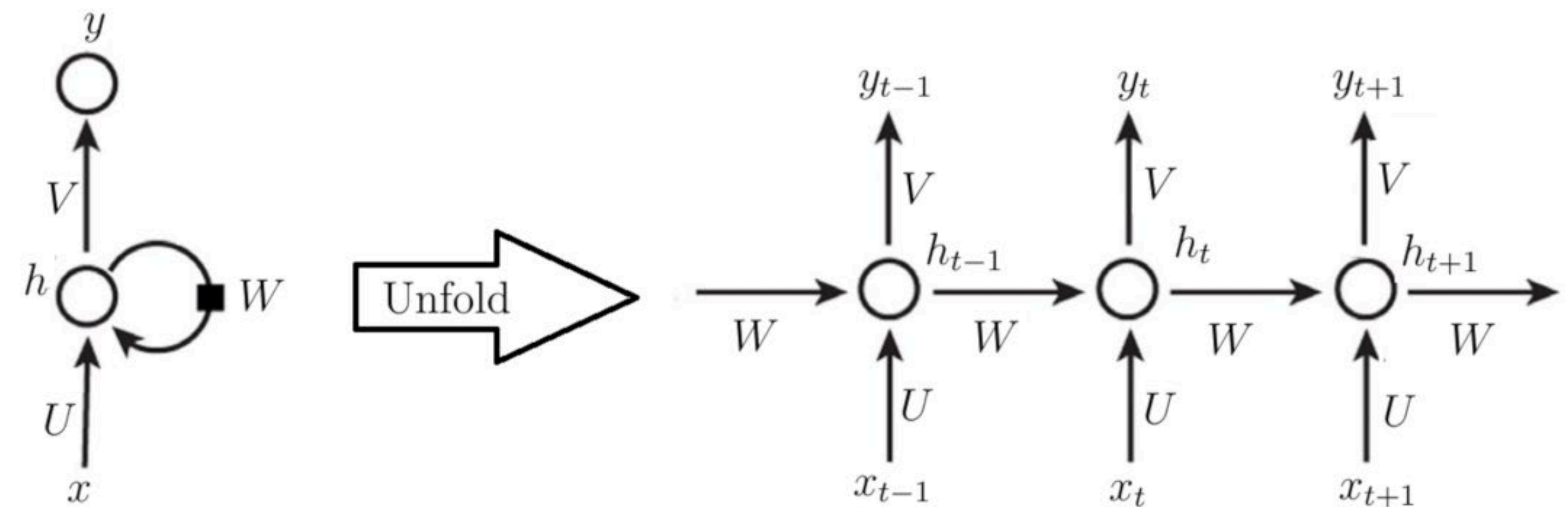
$$y_t = \sigma(Vh_t + b_y)$$

Обучение RNN:

$$\sum_{t=0}^T L(U, V, W) \rightarrow \min_{U, V, W}$$

Недостатки:

- Длина входной и выходной последовательностей должна быть одинаковой, что неестественно для перевода.
- Последовательный однонаправленный просмотр.



Напоминание: seq2seq

$X = (x_1, \dots, x_n)$ - входная последовательность

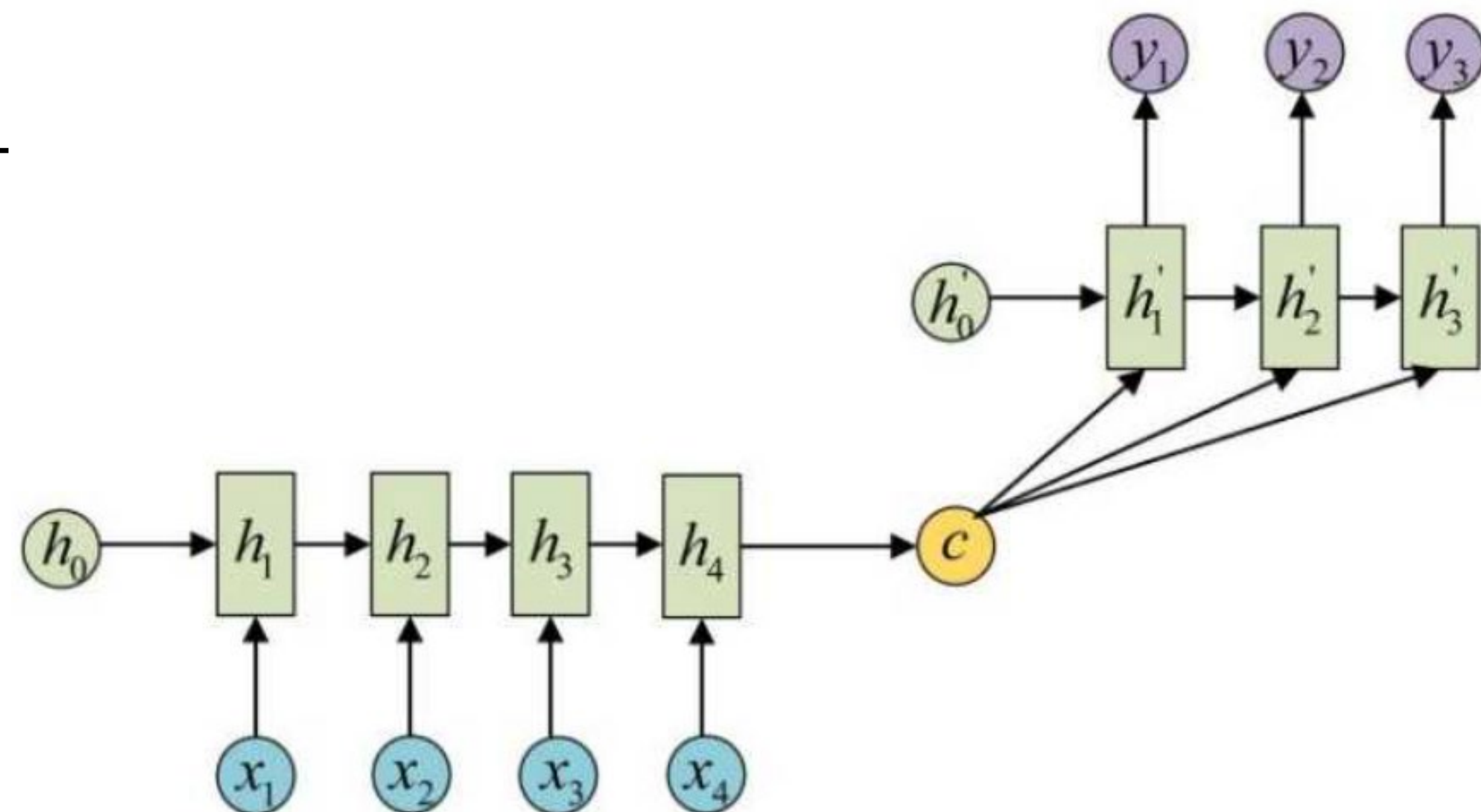
$Y = (y_1, \dots, y_m)$ - выходная последовательность

c — контекстный вектор, используемый декодером (context vector)

$$h_i = f_{in}(x_i, h_{i-1})$$

$$h'_t = f_{out}(h'_{t-1}, y_{t-1}, c)$$

$$y_t = f_y(h'_t, y_{t-1})$$



Недостатки:

- «бутылочное горлышко» вектора c
- затухание, взрывы градиента

RNN с вниманием

$a(h, h')$ - функция сходства состояний
входа h , и выхода h'

c_t - вектор входного контекста для выхода t

α_{t_i} - **веса внимания (attention score)**

важность i - го входа для выхода t

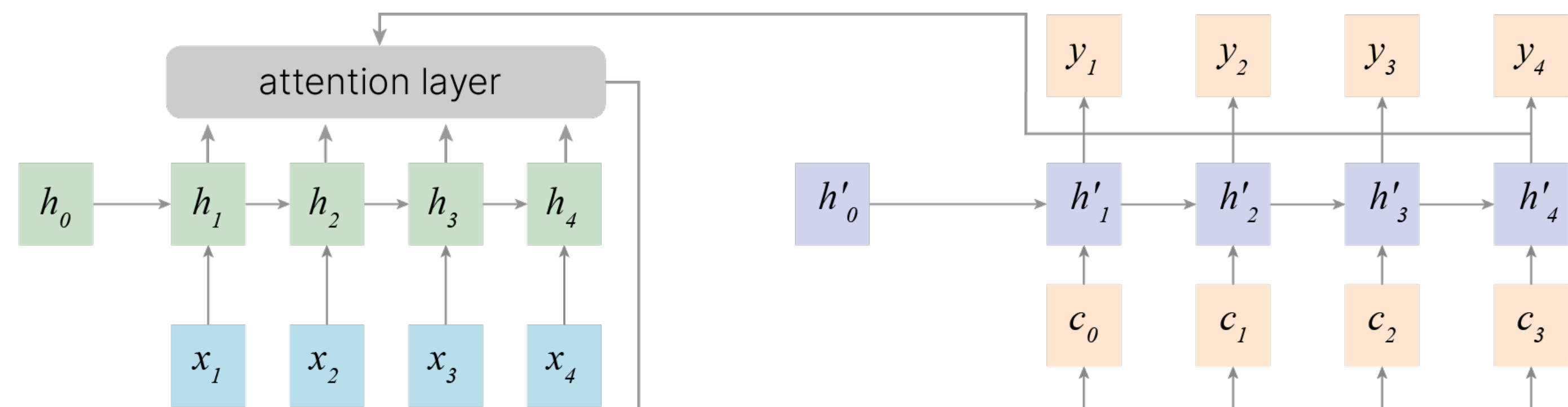
$$h_i = f_{in}(x_i, h_{i-1})$$

$$\alpha_{t_i} = \text{norm } a(h_i, h'_{t-1})$$

$$c_t = \sum_i \alpha_{t_i} h_i$$

$$h'_t = f_{out}(h'_{t-1}, y_{t-1}, c_t)$$

$$y_t = f_{out}(h'_t, y_{t-1}, c_t)$$



Матрица внимания

Матрица внимания — двумерная таблица, в которой каждая ячейка соответствует значению вектора attention score и обновляется во время обучения.

Матрица показывает на какие слова входной последовательности x_i модель обращает внимание, генерируя слова выходной последовательности y_i

	- Каждый	- охотник	- желает	- знать	- где	- сидит	- фазан
Каждый -	0.16	0.36	0.11	0.51	0.02	0.12	0.22
охотник -	0.56	0.36	0.11	0.52	0.11	0.19	0.23
желает -	0.06	0.22	0.59	0.61	0.18	0.02	0.50
знать -	0.81	0.31	0.51	0.40	0.27	0.12	0.22
где -	0.14	0.46	0.11	0.29	0.28	0.07	0.52
сидит -	0.34	0.76	0.27	0.56	0.78	0.09	0.42
фазан -	0.76	0.26	0.40	0.06	0.58	0.70	0.12

Функции сходства

1. Скалярное произведение:

$$a(h, h') = h^T h'$$

2. Экспоненциальное скалярное произведение:

$$a(h, h') = \exp(h^T h')$$

3. Билинейное внимание:

$$a(h, h') = h^T W h'$$

W — обучаемая матрица, которая позволяет моделировать сложные взаимодействия

4. Аддитивное внимание:

$$a(h, h') = w^T \tanh(Uh + Vh')$$

w, U, V — обучаемые параметры

Query-Key-Value attention

Функция внимания вычисляет взвешенную сумму значений, при этом веса определяются мерой соответствия между запросом и ключами:

$$a(h_i, h'_{t-1}) = \frac{(\mathbf{W}_k h_i)^T (\mathbf{W}_q h'_{t-1})}{\sqrt{d}},$$

d — размерность векторов (для стабилизации градиентов)

$$\alpha_{t_i} = \text{softmax } a(h_i, h'_{t-1})$$

$$c_t = \sum_i \alpha_{t_i} \mathbf{W}_v h_i$$

$\mathbf{W}_k, \mathbf{W}_q, \mathbf{W}_v$ - матрицы **key, query, value**

Три вектора Query, Key, Value:

Query (Запрос) – текущее состояние декодера. Можно представить как вопрос, который токен задаёт сети: «Что мне нужно знать, чтобы понять контекст?».

Key (Ключ) - это то, что содержит данный токен, его содержимое.

Value (Значение) - хранит информацию, которая может быть передана другим токенам.

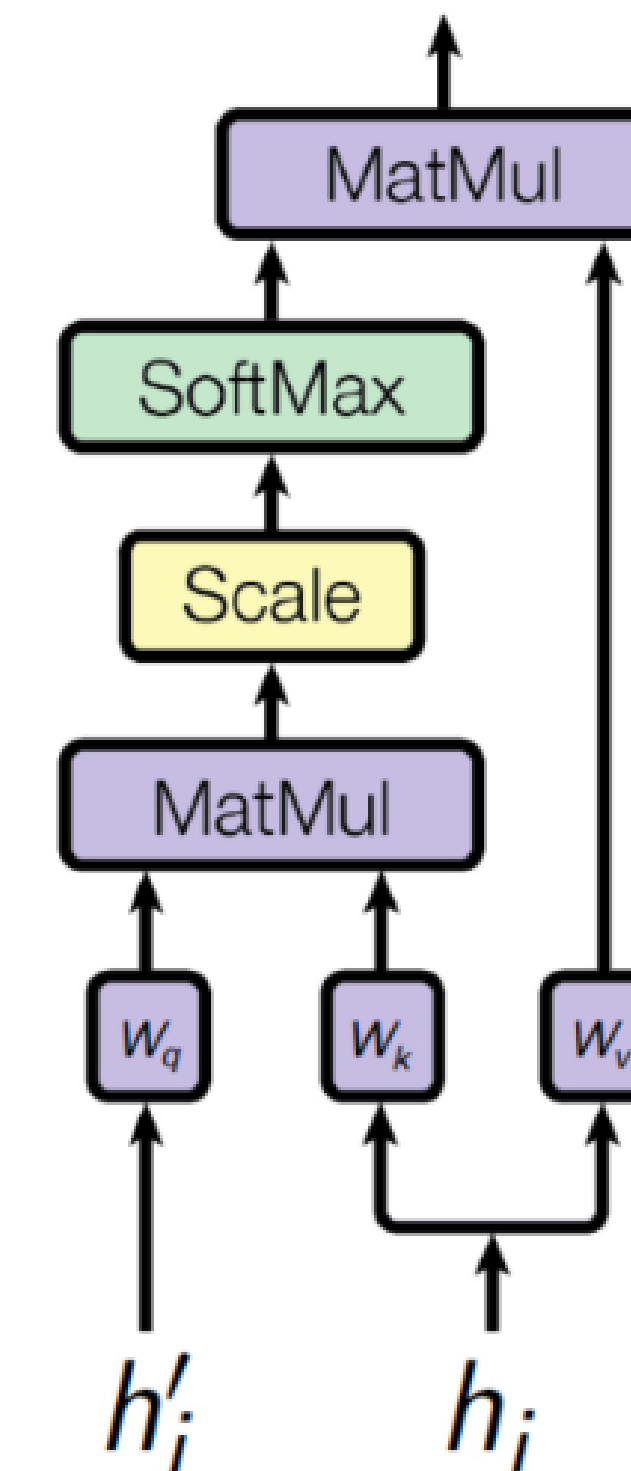
Key и **Value** — это выходы декодера.

Scale dot product attention

На вход подаются **запросы** и **ключи** (размерностью d_k) и **значения** (размерностью d_v) .

Далее считается *скалярное произведение* (*Matmul*) запросов со всеми ключами.

Scale: каждый результат делится на \sqrt{d} , потом применяется *Softmax* чтобы получить веса значений.



Модель внимания

2017 Google «Attantion is all you need»

Искомый вектор контекста релевантного запросу:

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

Q, K, V – вектора запроса (query), ключа (key), значения (value).

self-attention – частный случай, когда входная и выходная последовательности совпадают.

$$c = \text{Attention}(q, k, v) = \sum_i v_i \text{Softmax}_i a(k_i, q)$$

$a(k_i, q)$ – функция сходства, оценка релевантности ключа k_i запросу q

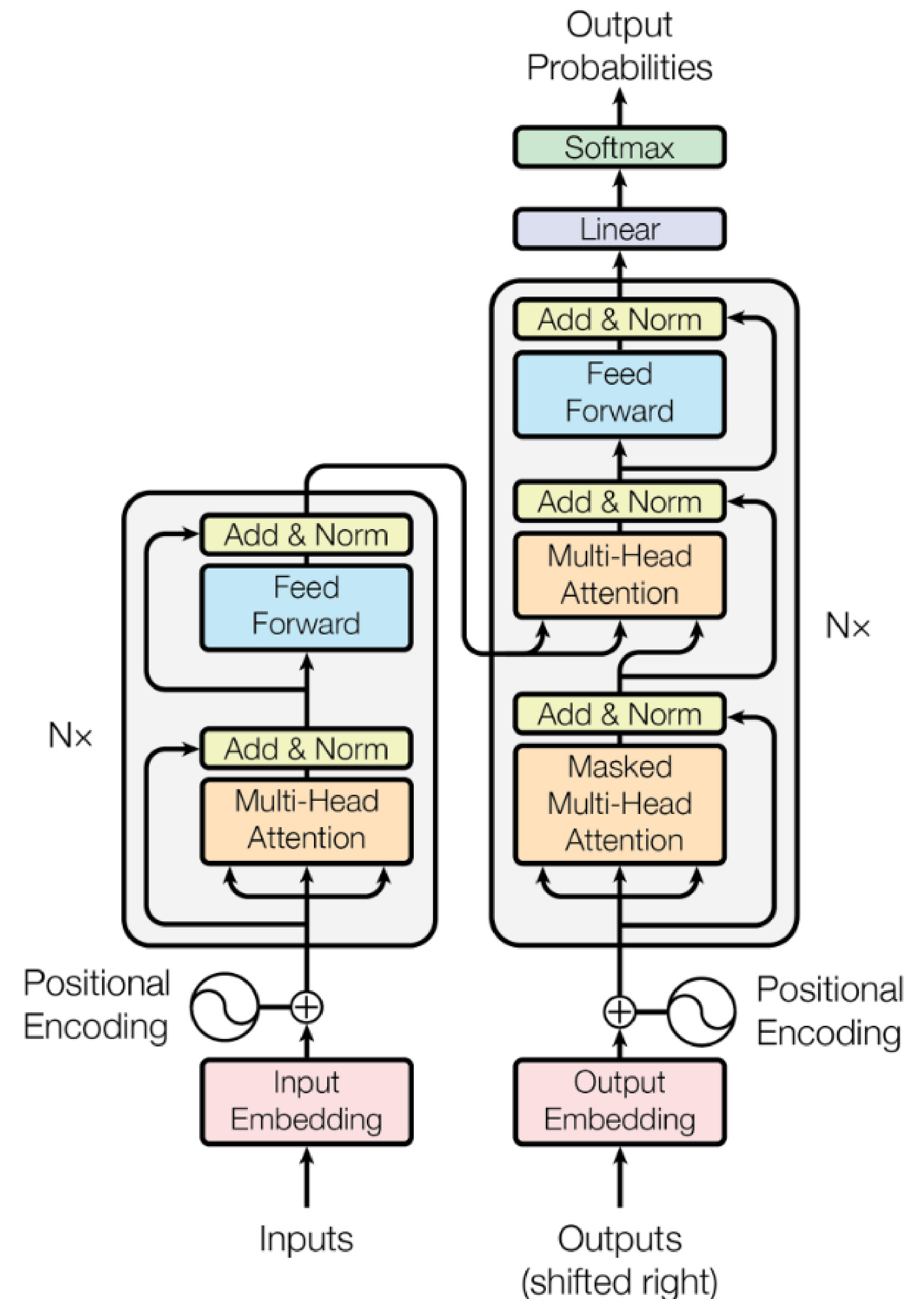
Архитектура Transformer

Энкодер:

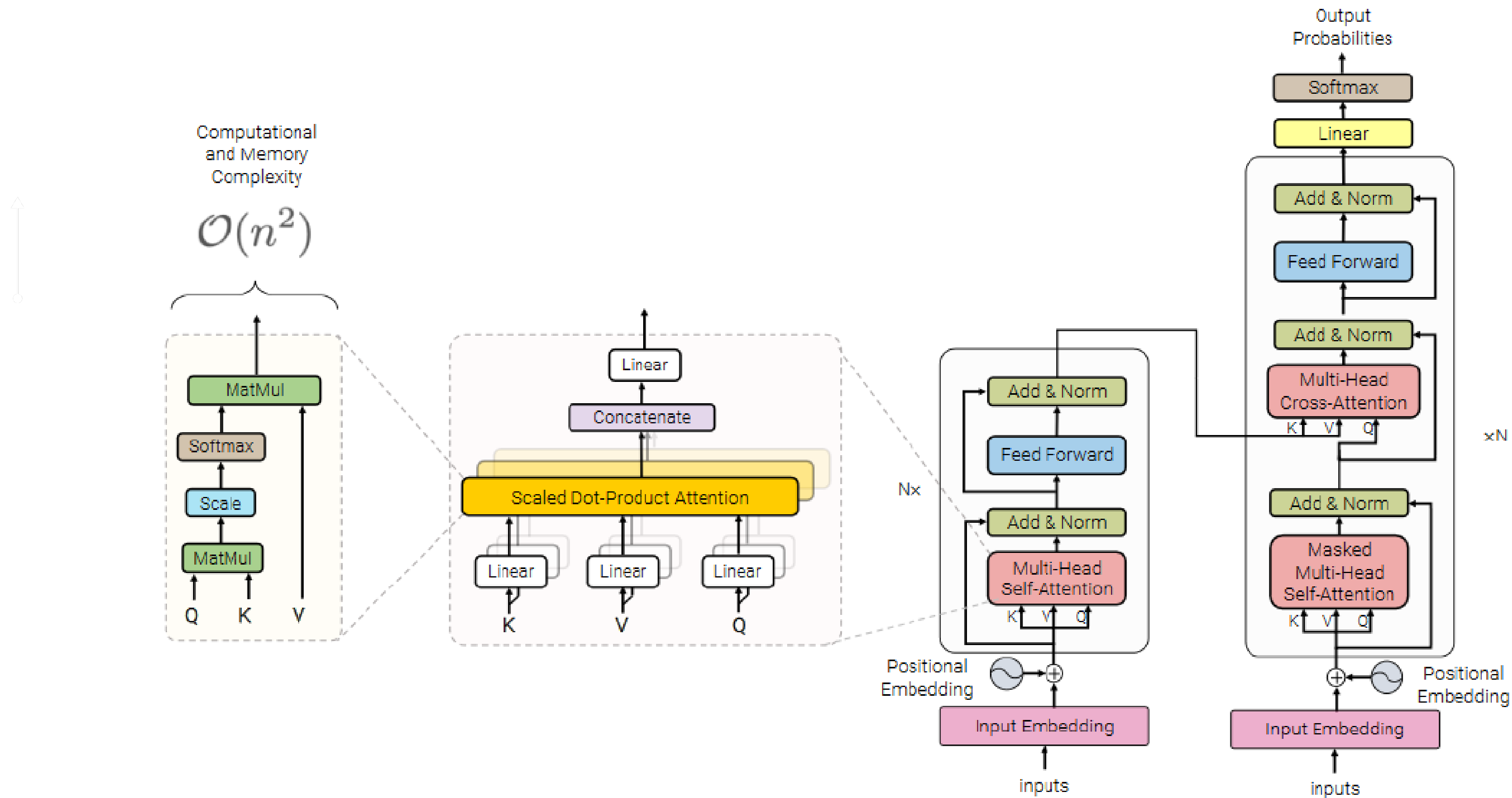
1. Векторное представление входной последовательности
2. Positional encoding
3. Multi-head self-attention слой:
 - Полносвязный слой (feed-forward network);
 - Нормализацию и остаточные связи (LayerNorm + residual connections).

Декодер:

- + masked self-attention
- + attention к выходу декодера
- + линейный слой и softmax

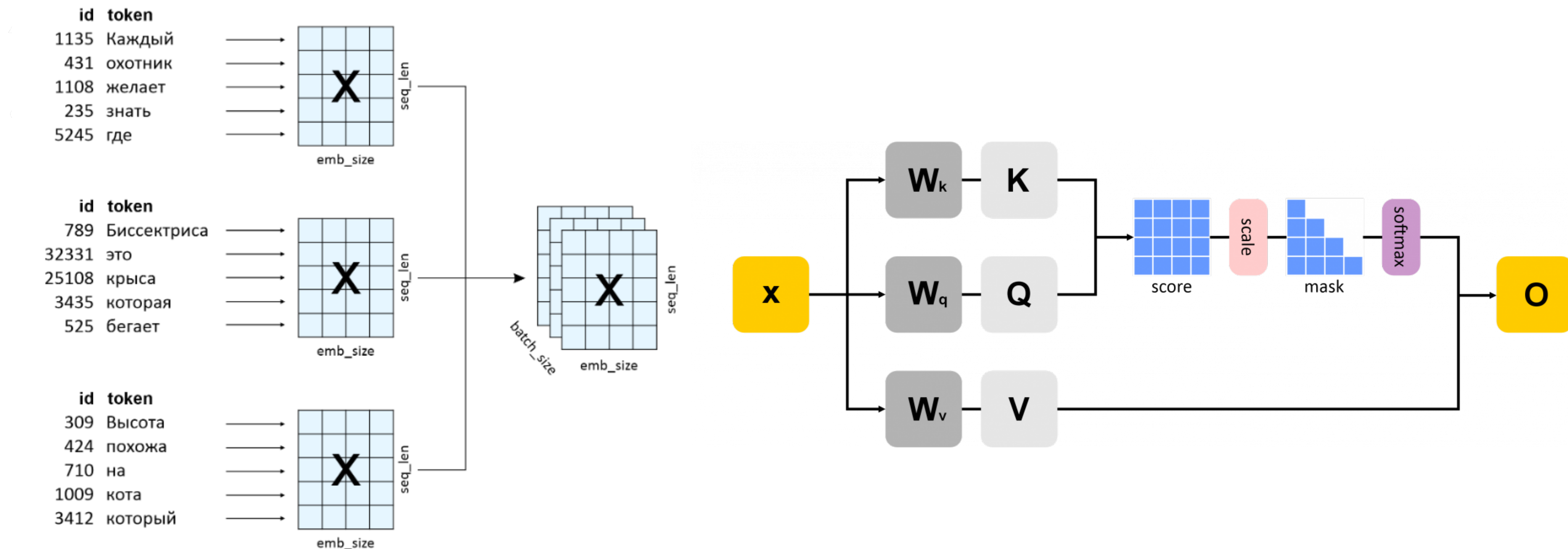


Архитектура Transformer

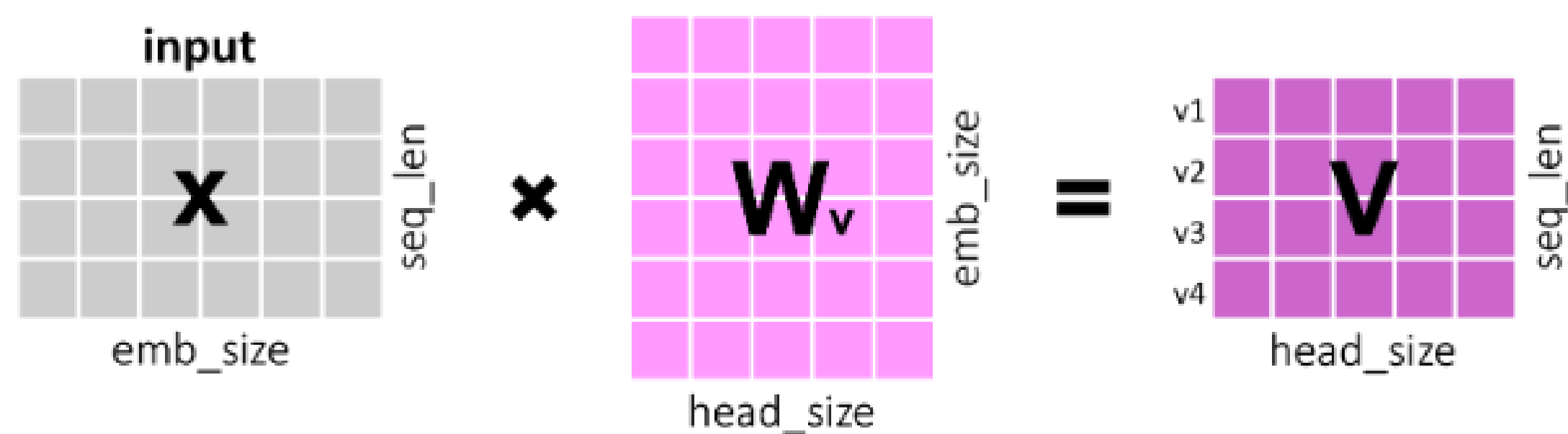
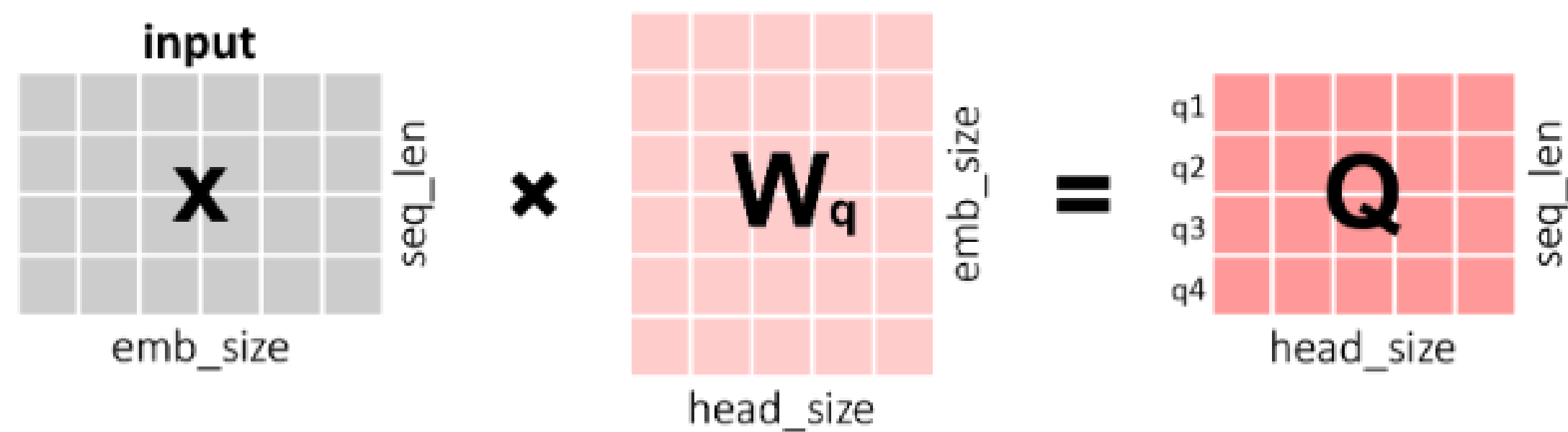
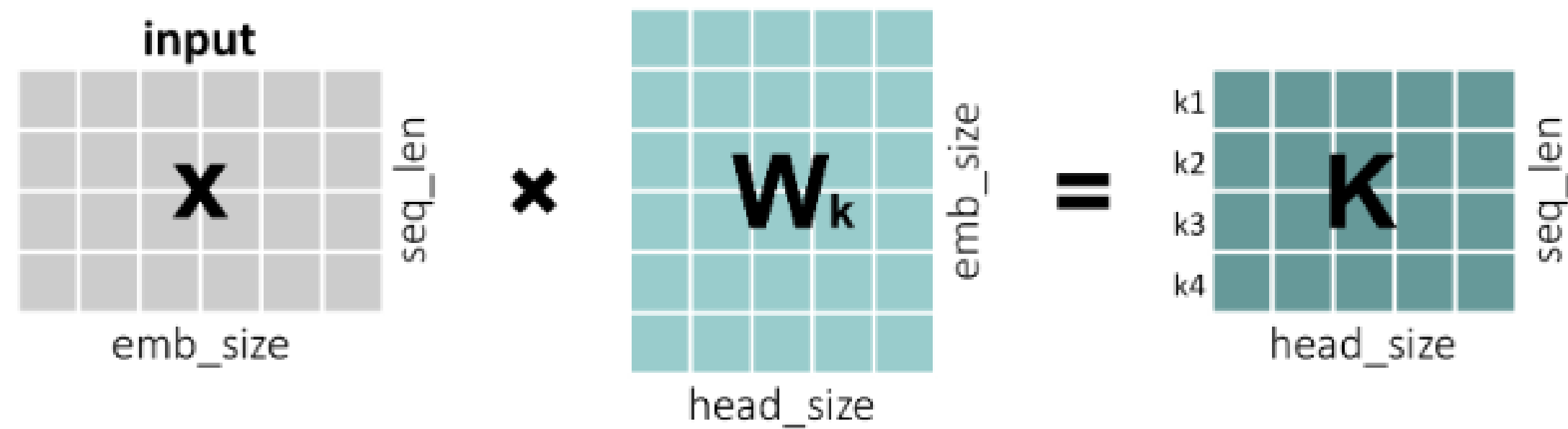


Реализация self-attention

Вход: тензор x размерностью $batch_size \times seq_len \times emb_size$

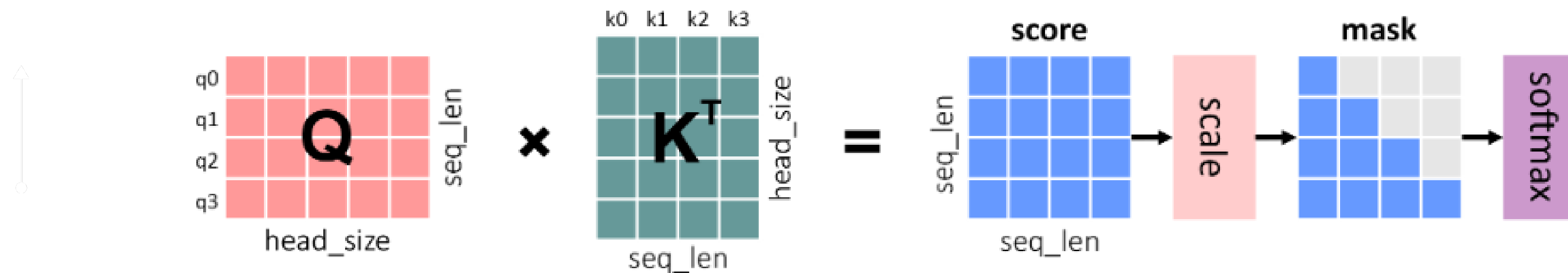


Реализация self-attention

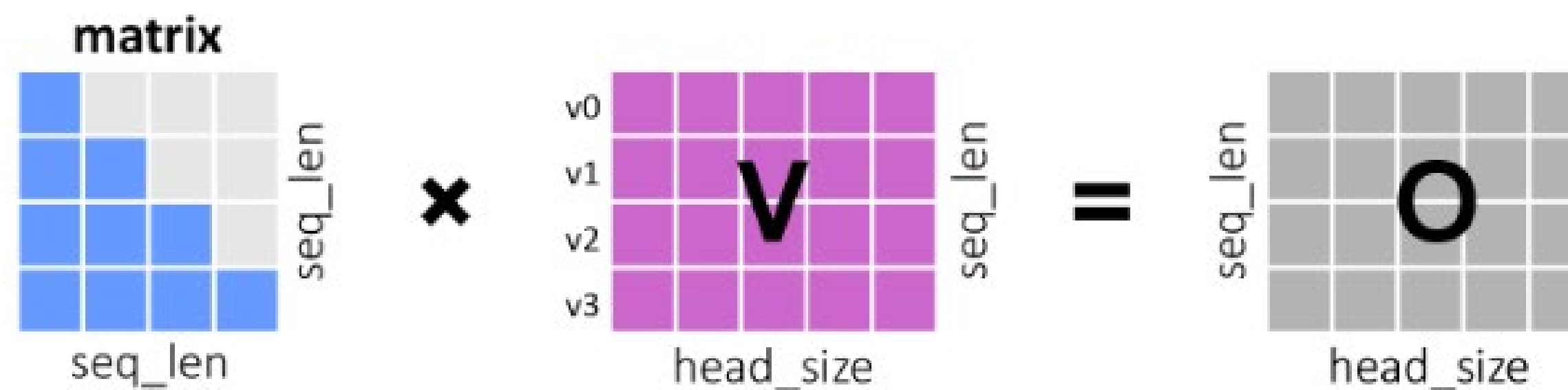


Реализация self-attention

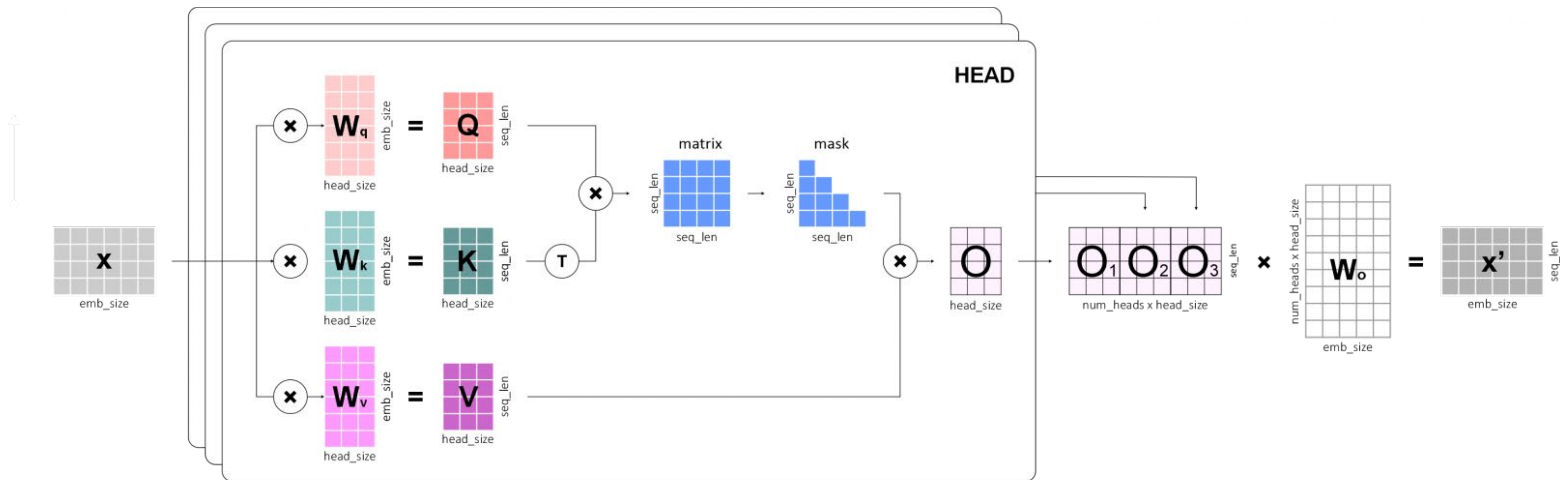
Матрица внимания:



Выходной тензор:



Multi-head attention





**Спасибо
за внимание!**

