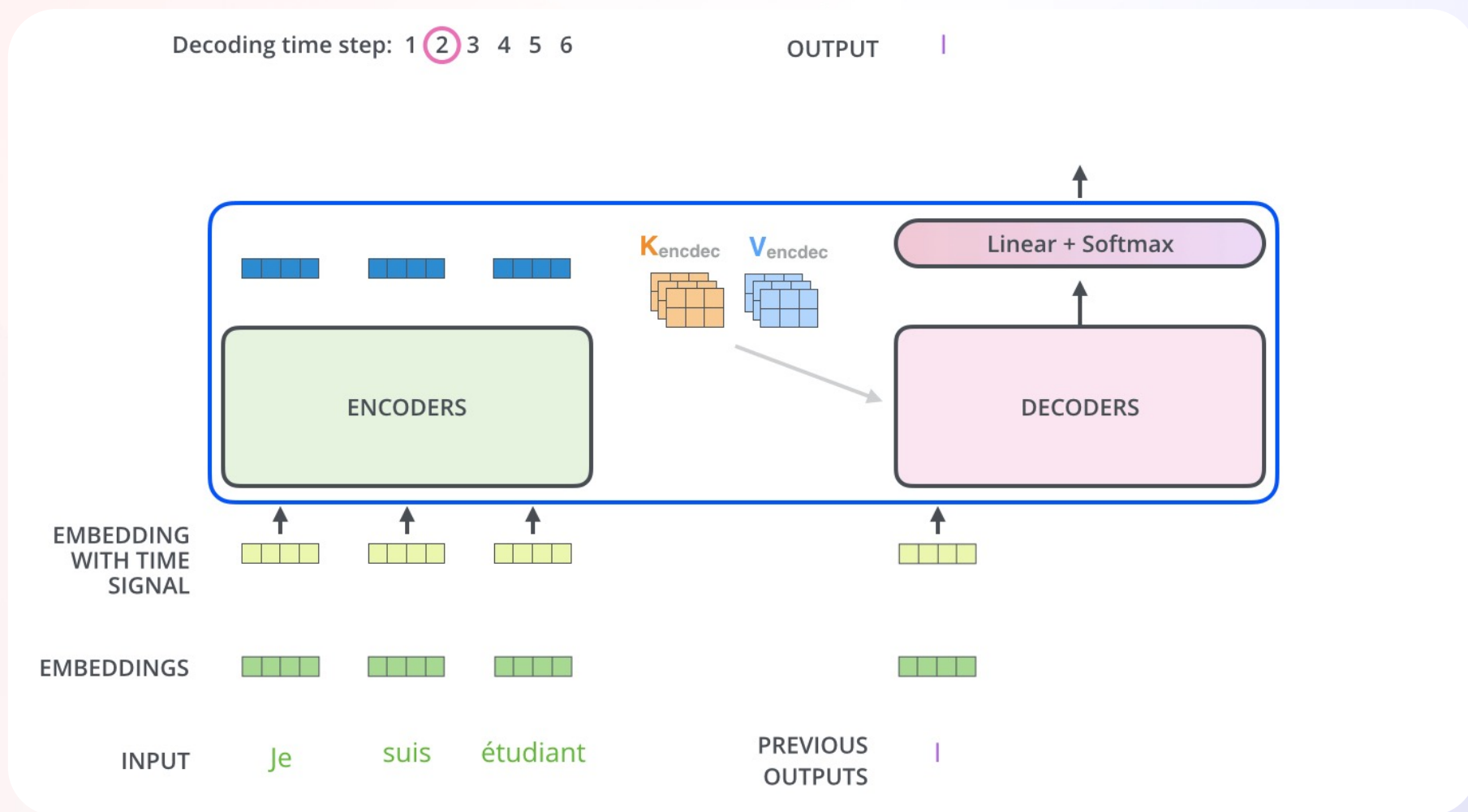
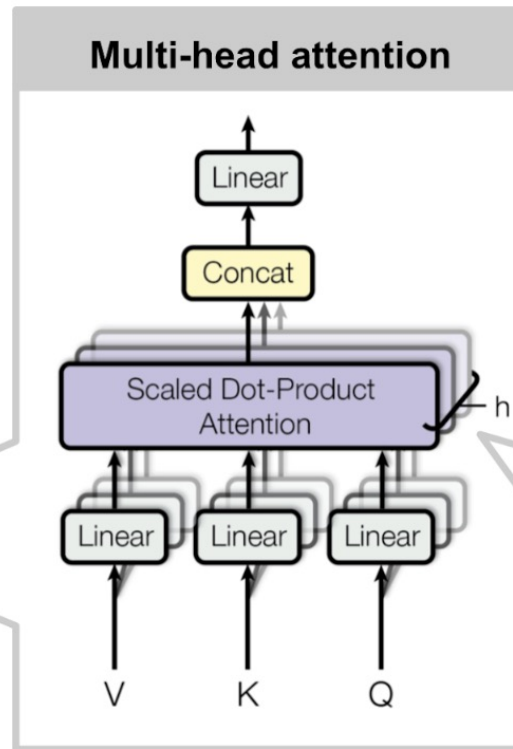
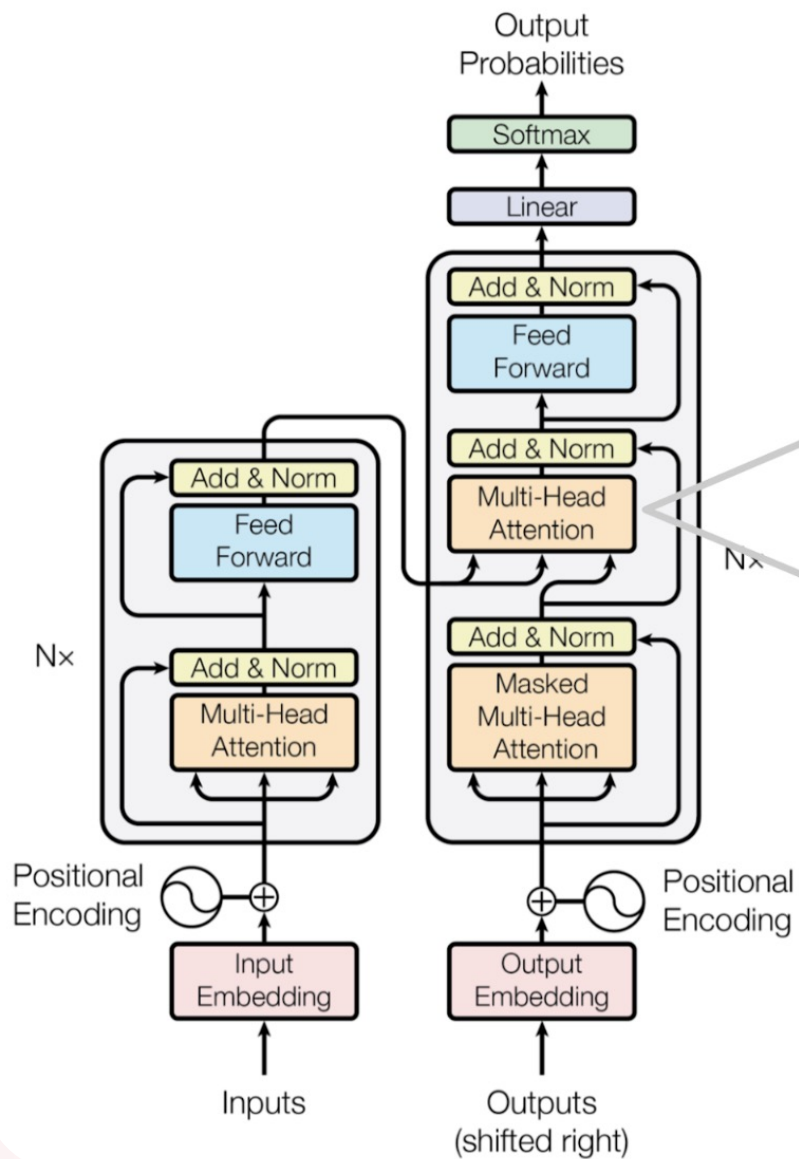


Проблемы трансформеров и их решения

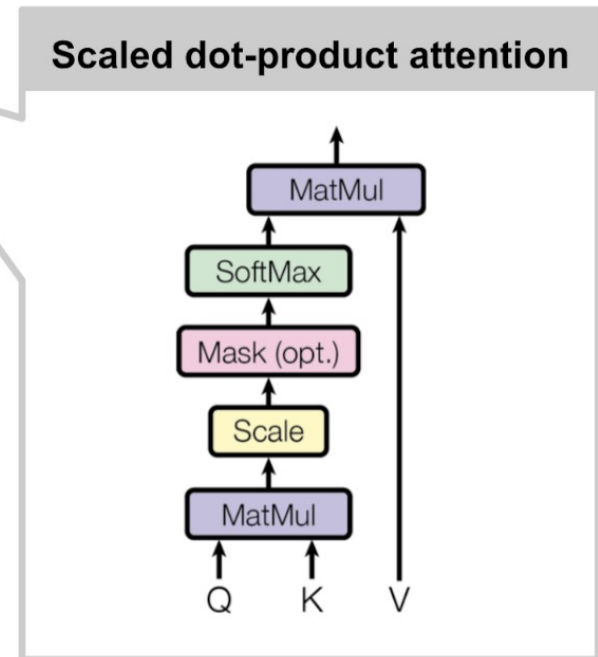
Денис Козлов

Трансформер





Zoom-In!



Zoom-In!

Нотация

Входная последовательность длины L

Выходная последовательность длины \dot{L}

Эмбединги размера d

Входная последовательность $X \in \mathbb{R}^{L \times d}$

Матрицы весов $W^k, W^q \in \mathbb{R}^{d \times d_k}$; $W^v \in \mathbb{R}^{d \times d_v}$

Как правило $d_k = d_v = d$

Attention

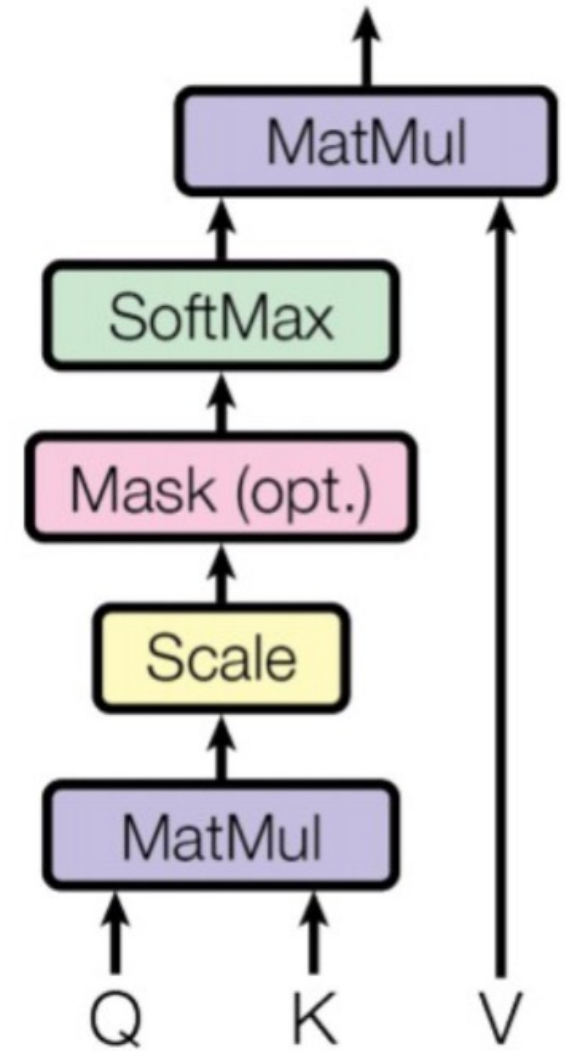
Query $Q = XW^q \in \mathbb{R}^{L \times d_k}$

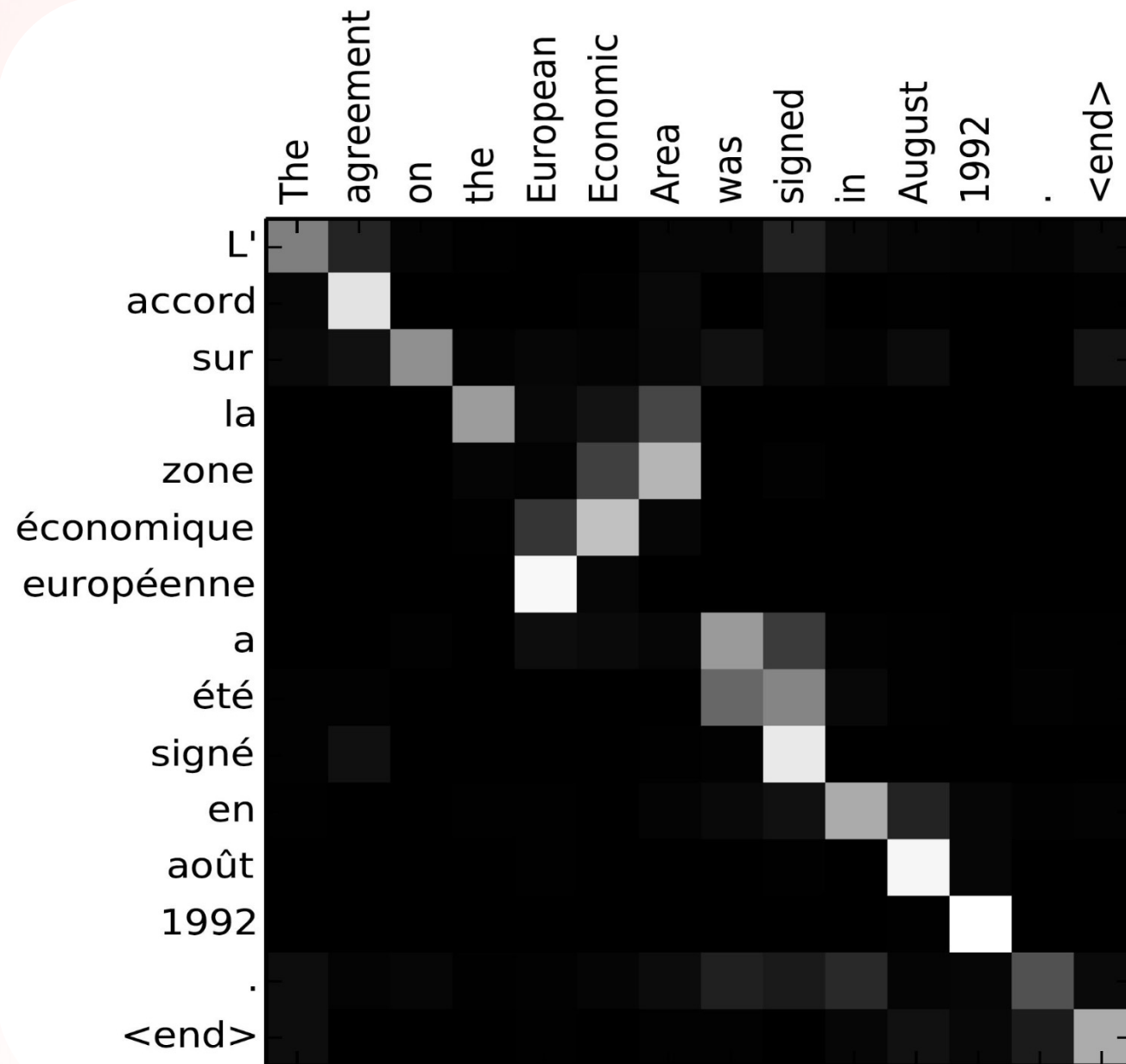
Key $K = XW^k \in \mathbb{R}^{L \times d_k}$

Value $V = XW^v \in \mathbb{R}^{L \times d_v}$

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

$\in \mathbb{R}^{L \times d_v}$





Attention

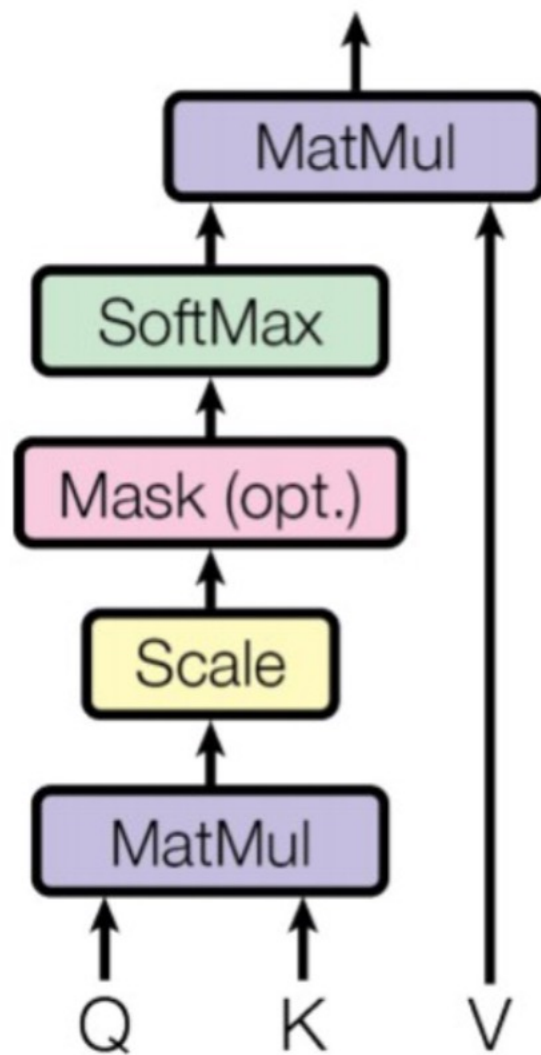
Query $Q = XW^q \in \mathbb{R}^{L \times d_k}$

Key $K = XW^k \in \mathbb{R}^{L \times d_k}$

Value $V = XW^v \in \mathbb{R}^{L \times d_v}$

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

$\in \mathbb{R}^{L \times d_v}$



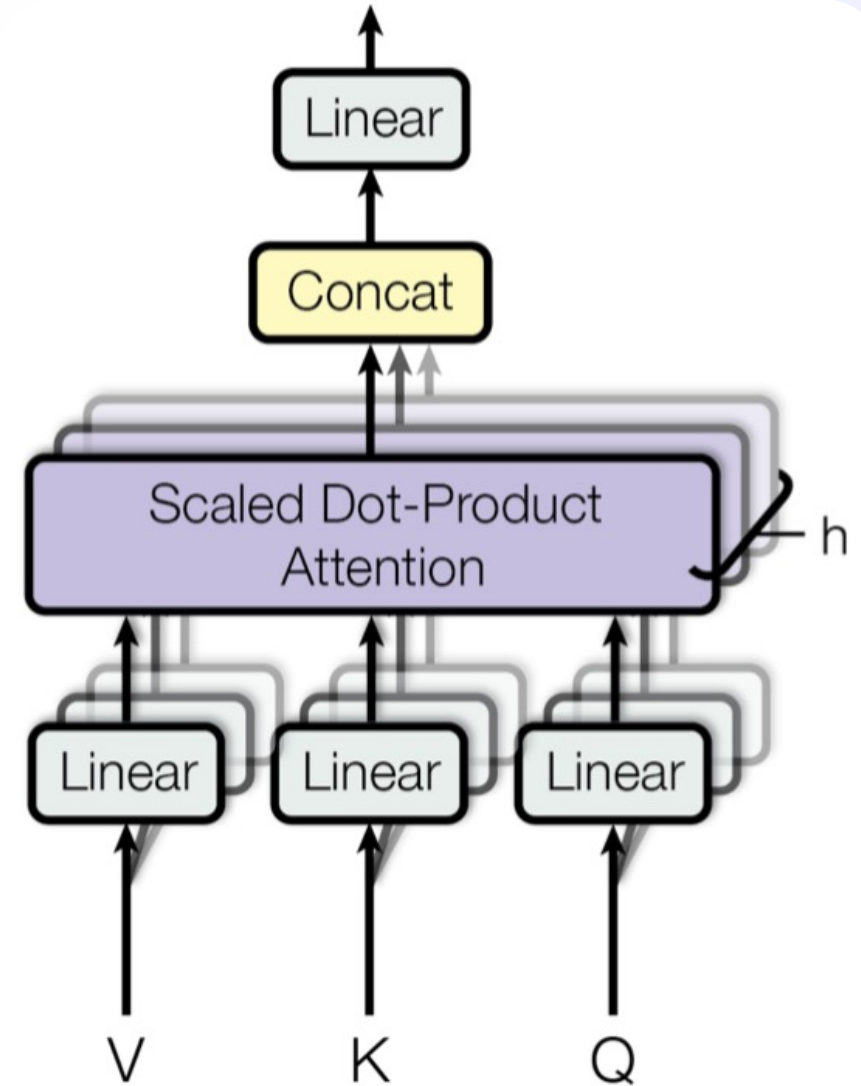
Multi-head attention

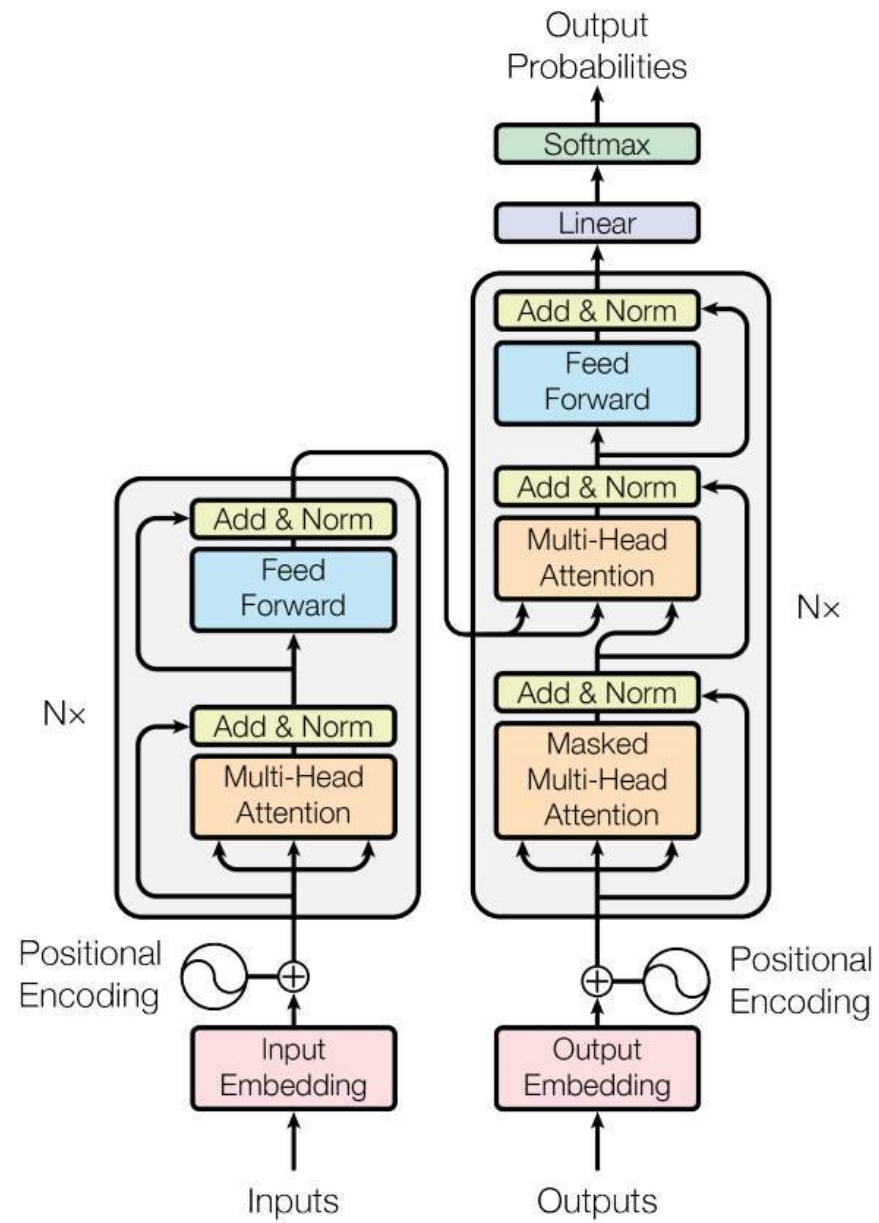
Количество голов h

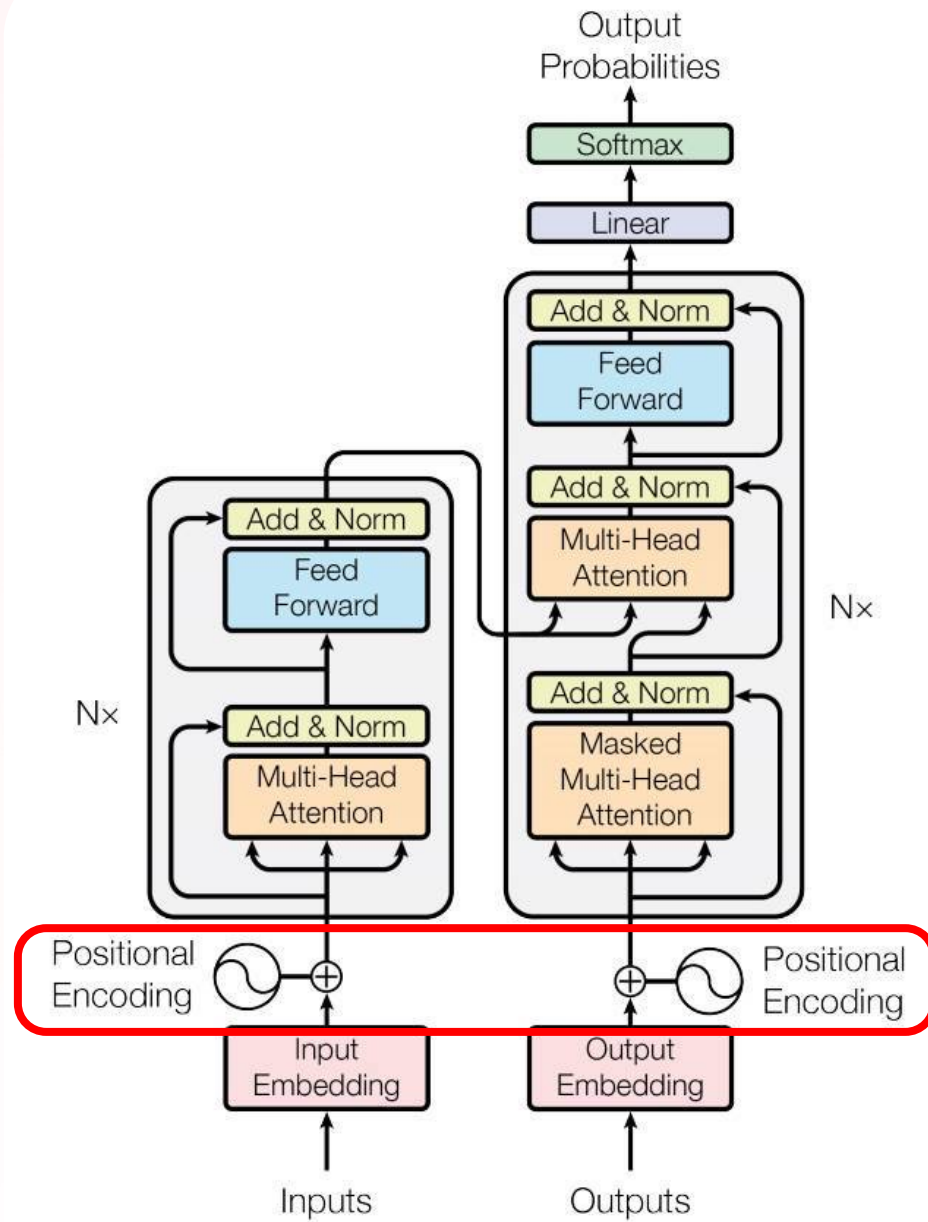
Матрица весов $W^o \in \mathbb{R}^{d_v \times d}$

$$\begin{aligned}\text{MultiHeadAttention}(X) &= \\ &= \text{concat}(\text{head}_1; \dots; \text{head}_h) W^o \\ &\quad \in \mathbb{R}^{L \times d}\end{aligned}$$

$$\text{head}_i = \text{Attention}(XW_i^q, XW_i^k, XW_i^v)$$







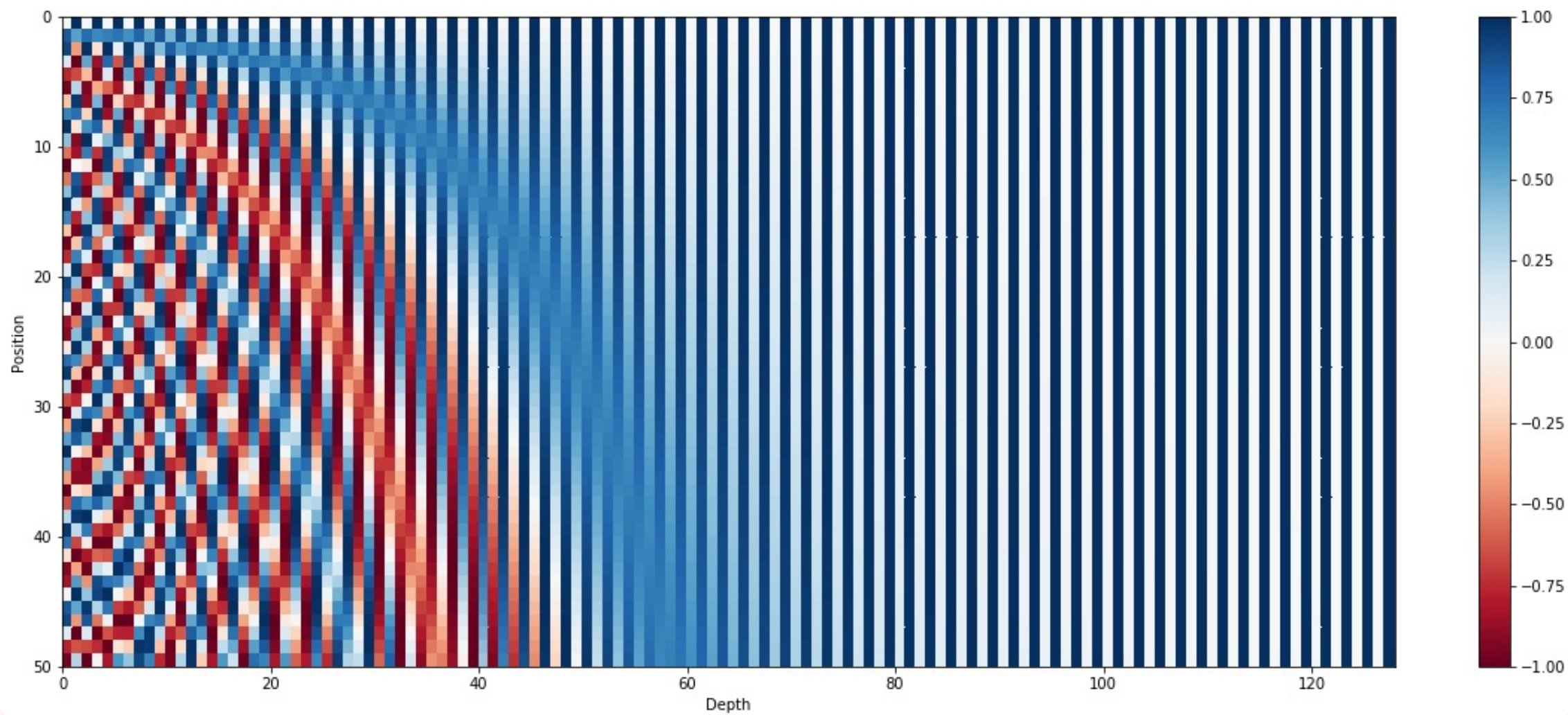
Positional encoding

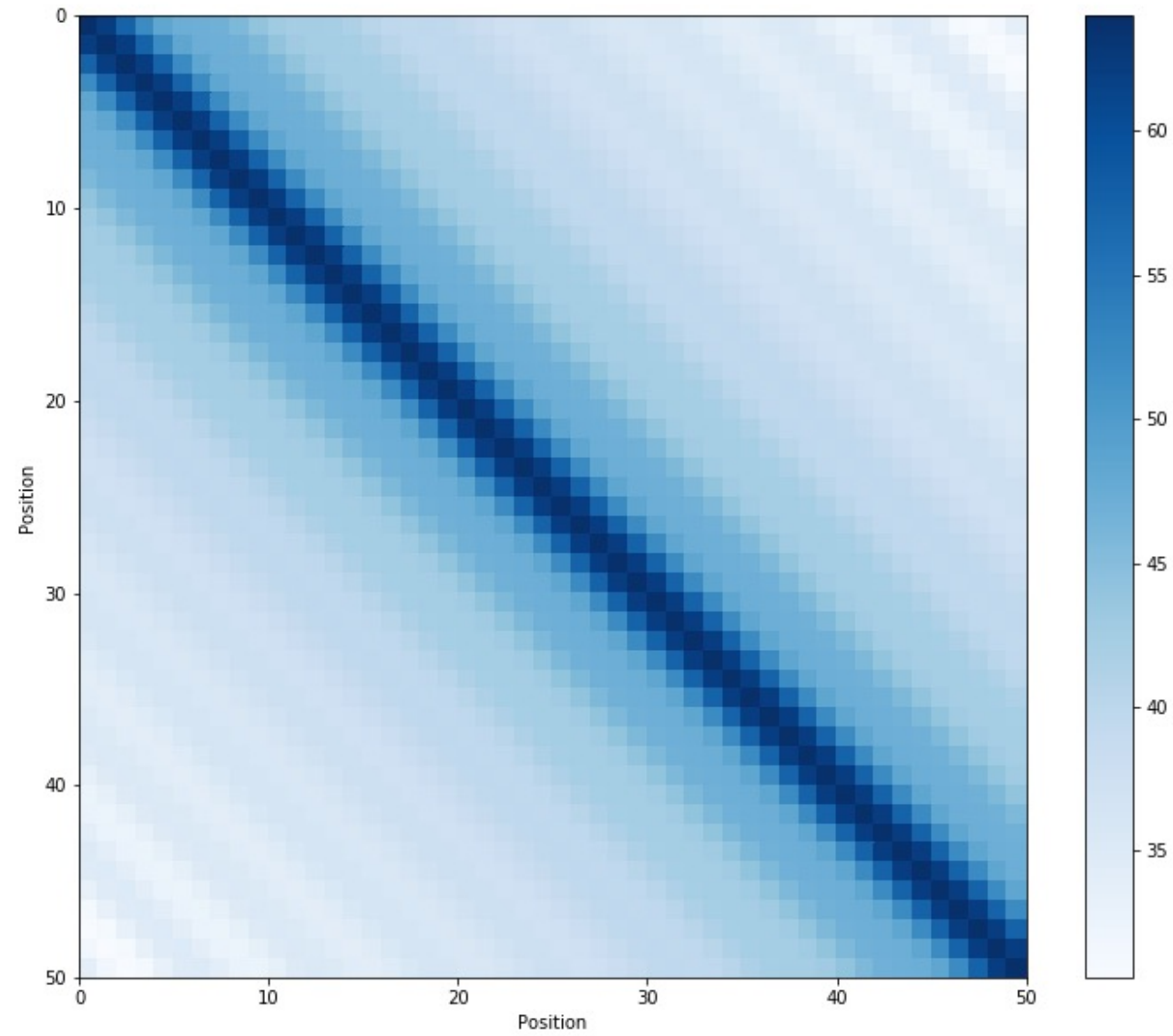
0000
0001
0010
...
1111

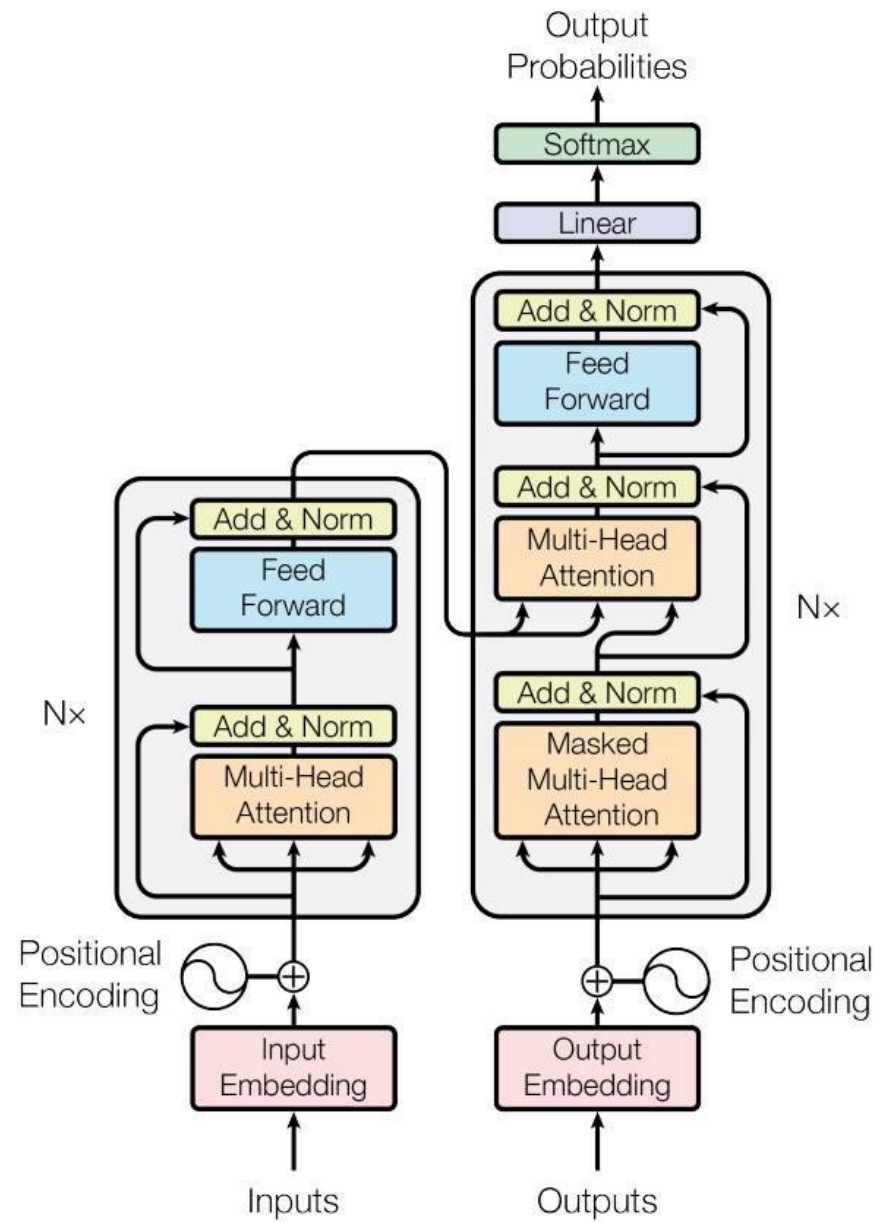
} Но у нас же float!

$$\vec{p}_t^{(i)} = f(t)^{(i)} := \begin{cases} \sin(\omega_k \cdot t), & \text{if } i = 2k \\ \cos(\omega_k \cdot t), & \text{if } i = 2k + 1 \end{cases}$$

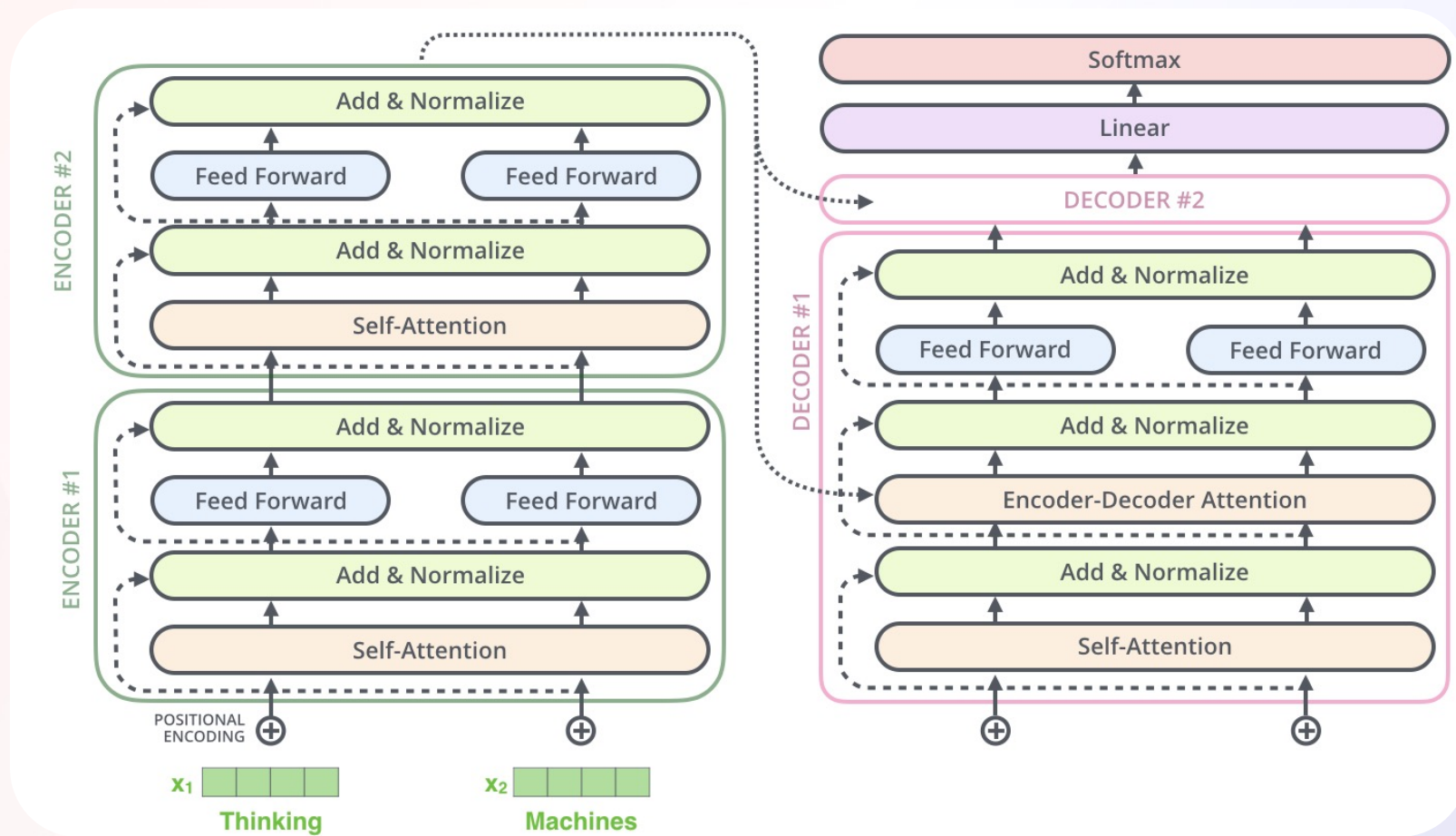
$$\omega_k = \frac{1}{10000^{2k/d}}$$







Много уровней



Все ли хорошо?



Все ли хорошо?

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

Query $Q = XW^q \in R^{L \times d_k}$

Key $K = XW^k \in R^{L \times d_k}$

Value $V = XW^v \in R^{L \times d_v}$

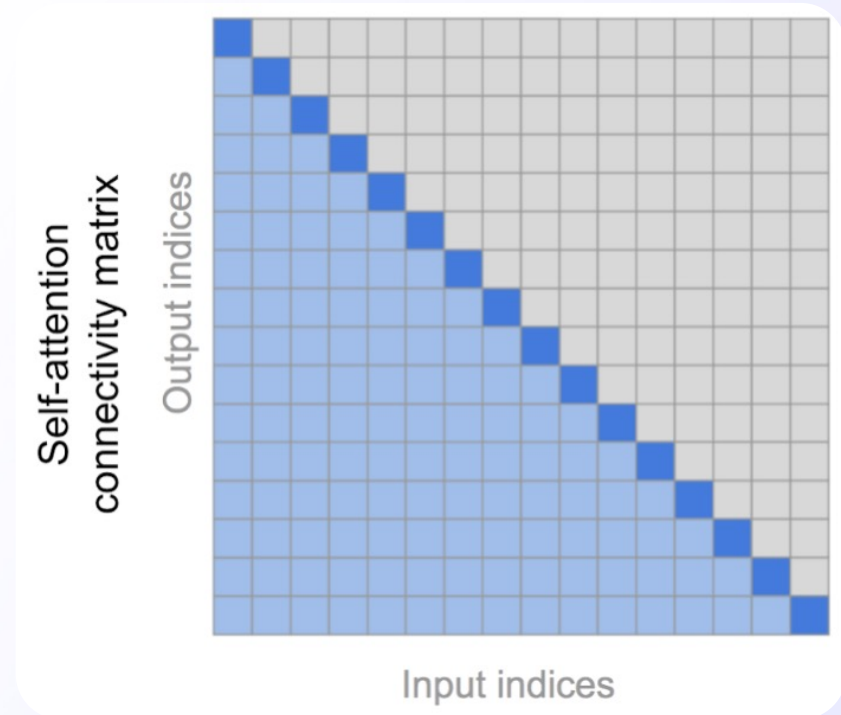
Требуется $O(L^2)$ памяти!

Для простоты примем $n = L$

Тяжелый attention

Decoder-only Transformer (GPT)

Для предсказания смотрим
на все предыдущие элементы
последовательности



Тяжелый attention

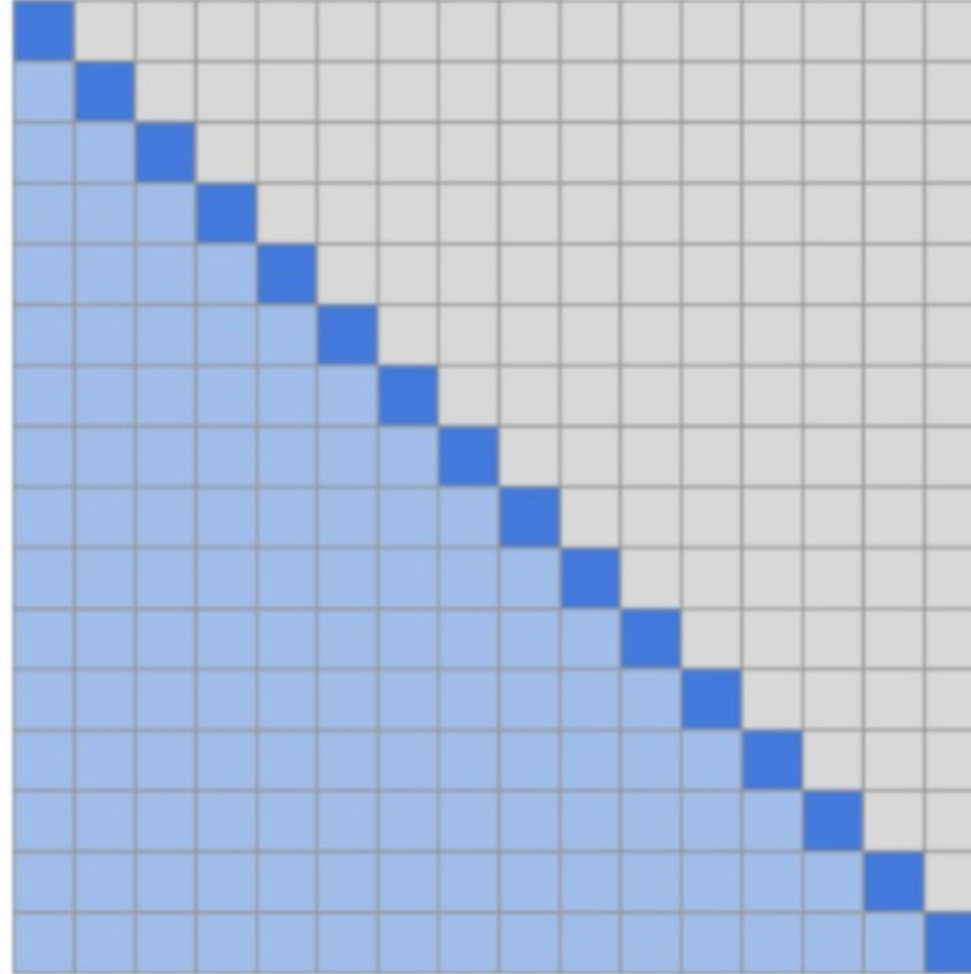
$$\text{Attend}(X) = \text{softmax} \left(\frac{(XW^q)(XW^k)^T}{\sqrt{d_k}} * \text{mask} \right) XW^v$$

$$\begin{aligned} \text{Attend}(\mathbf{X}, \mathcal{S}) &= \left(a(\mathbf{x}_i, S_i) \right)_{i \in \{1, \dots, L\}} \\ \mathcal{S} &= \{S_1, \dots, S_n\} \\ \text{where } a(\mathbf{x}_i, S_i) &= \text{softmax} \left(\frac{(\mathbf{x}_i \mathbf{W}^q)(\mathbf{x}_j \mathbf{W}^k)^T_{j \in S_i}}{\sqrt{d_k}} \right) (\mathbf{x}_j \mathbf{W}^v)_{j \in S_i} \end{aligned}$$

Пока что $S_i = \{j | j \leq i\}$

Self-attention connectivity matrix

Output indices



Input indices

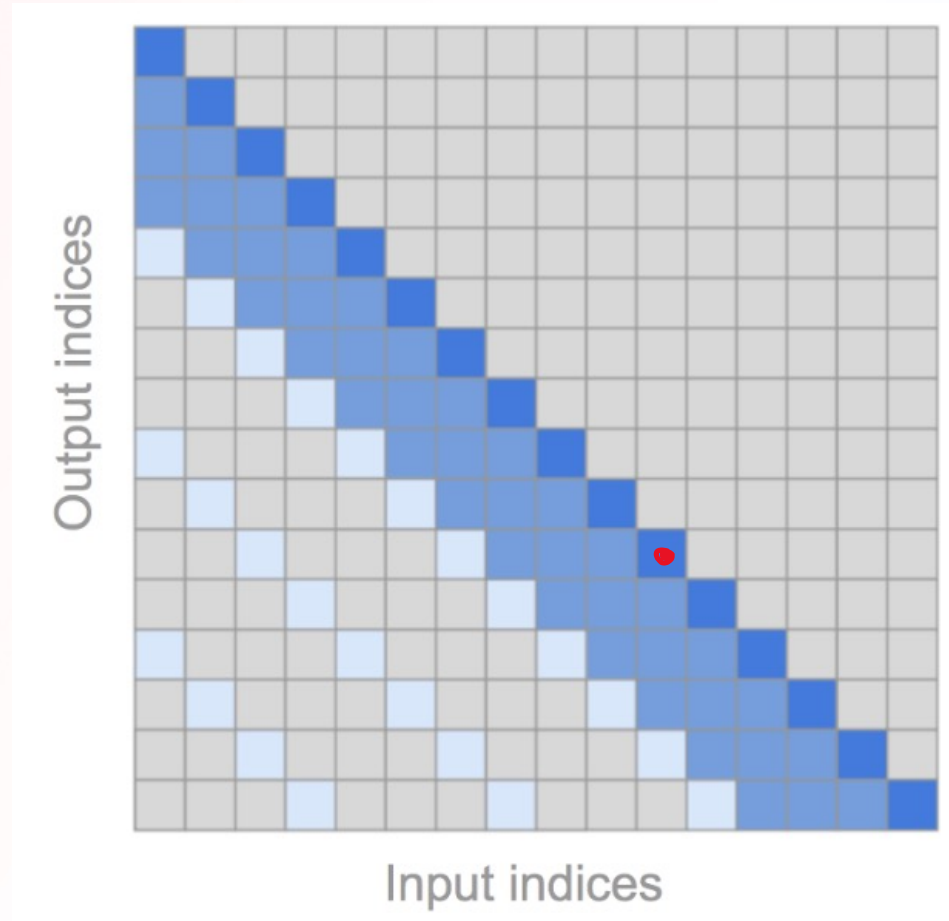
Почти Factorized self-attention

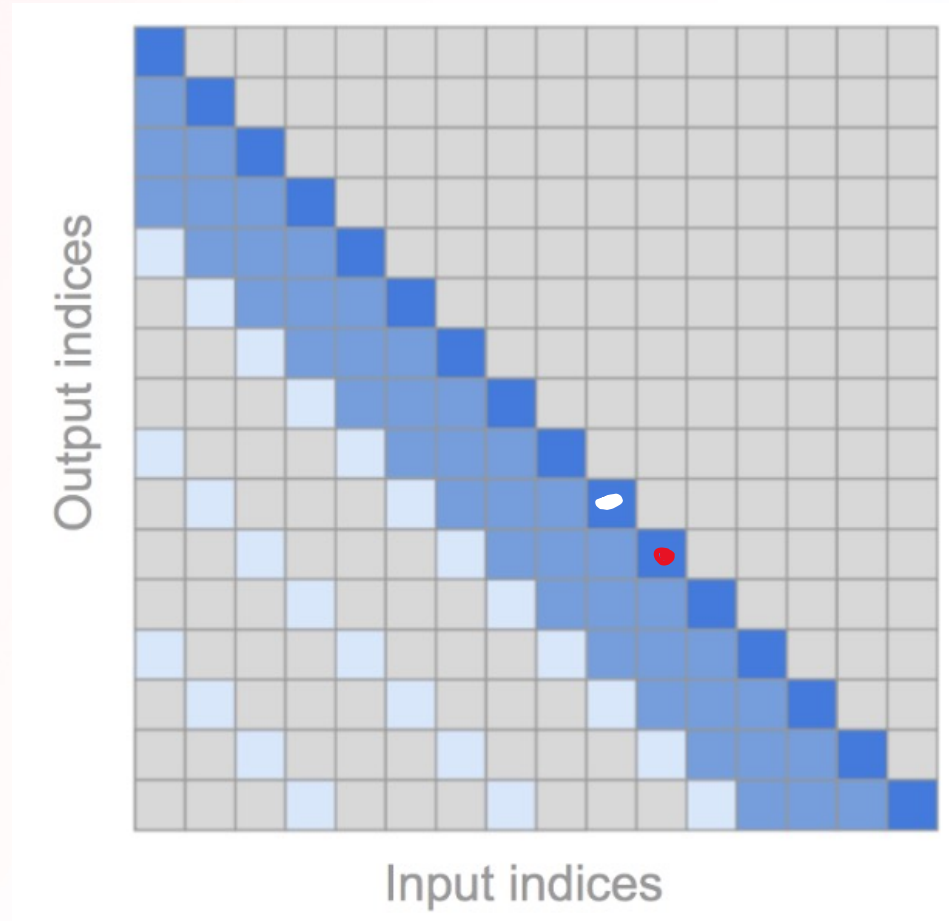
$S_i \subseteq \{j | j < i\}$, и верно:

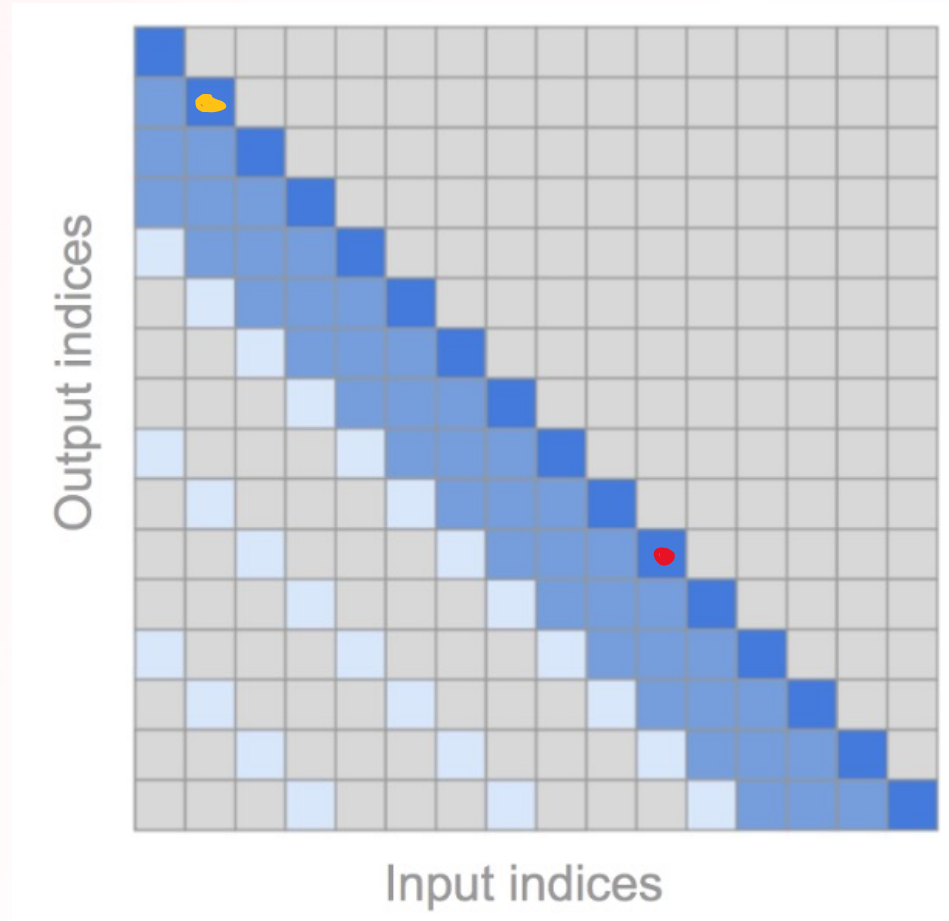
для каждой пары $j \leq i$ существует «путь» от i до j

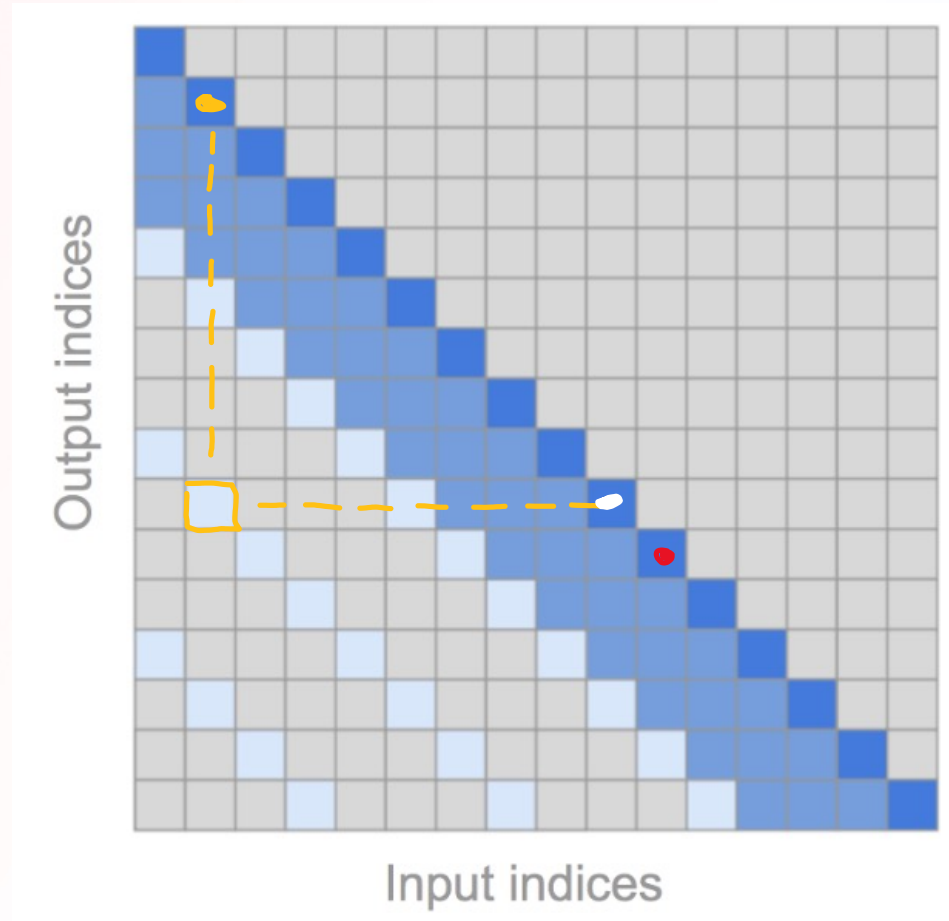
К примеру:

- $j \in S_i$ — *напрямую*
- $j \in S_{t_1}, t_1 \in S_{t_2}, \dots, t_k \in S_i$ — *есть маршрут*









Почти Factorized self-attention

$S_i \subseteq \{j | j < i\}$, и верно:

для каждой пары $j \leq i$ существует «путь» от i до j

К примеру:

- $j \in S_i$ — *напрямую*
- $j \in S_{t_1}, t_1 \in S_{t_2}, \dots, t_k \in S_i$ — *есть маршрут*

Factorized self-attention

$S_i \subseteq \{j | j < i\}$ делится на p непересекающихся множеств A_i^1, \dots, A_i^p и верно, что:

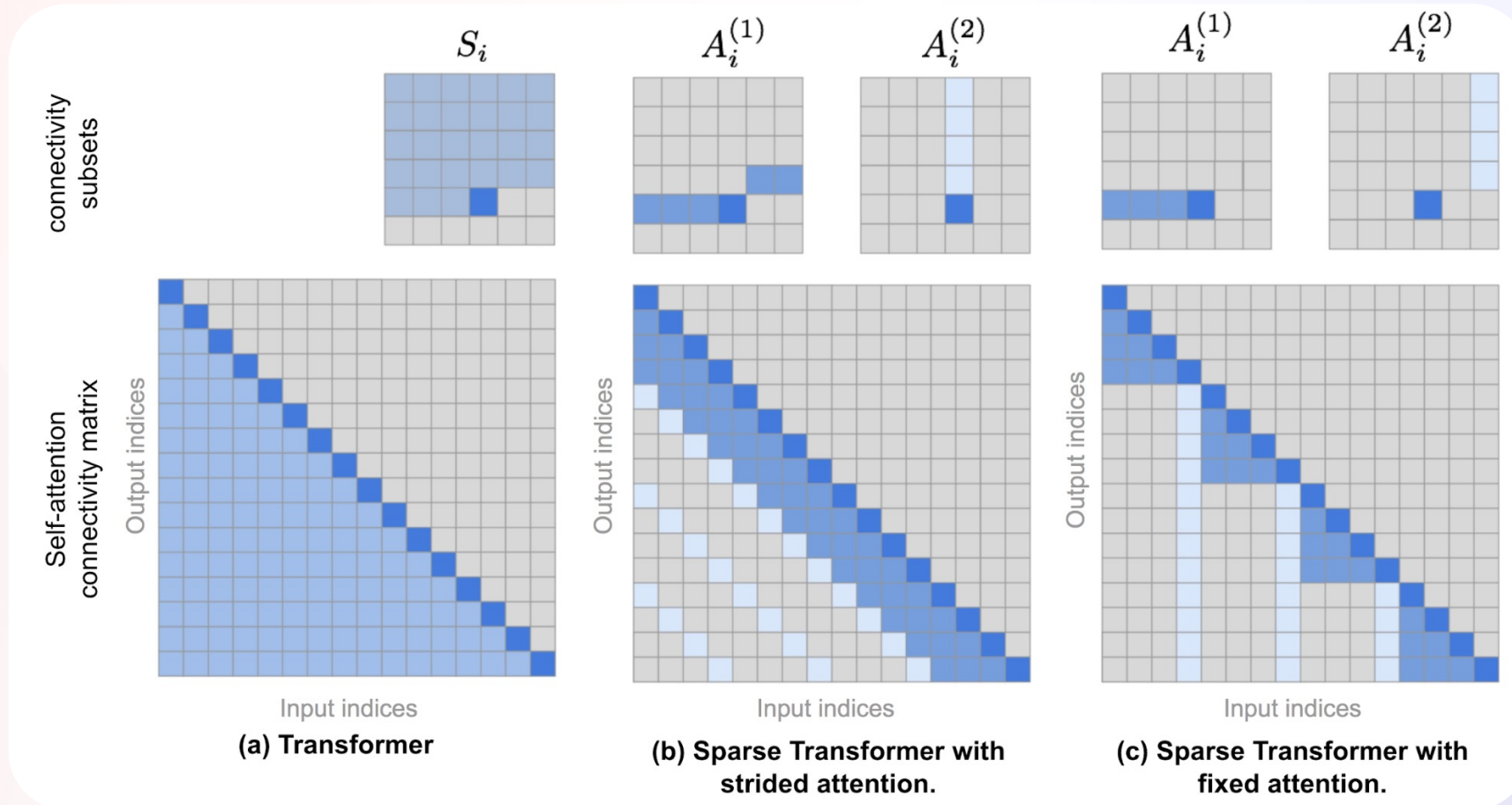
для каждой пары $j \leq i$ существует «путь» от i до j

К примеру:

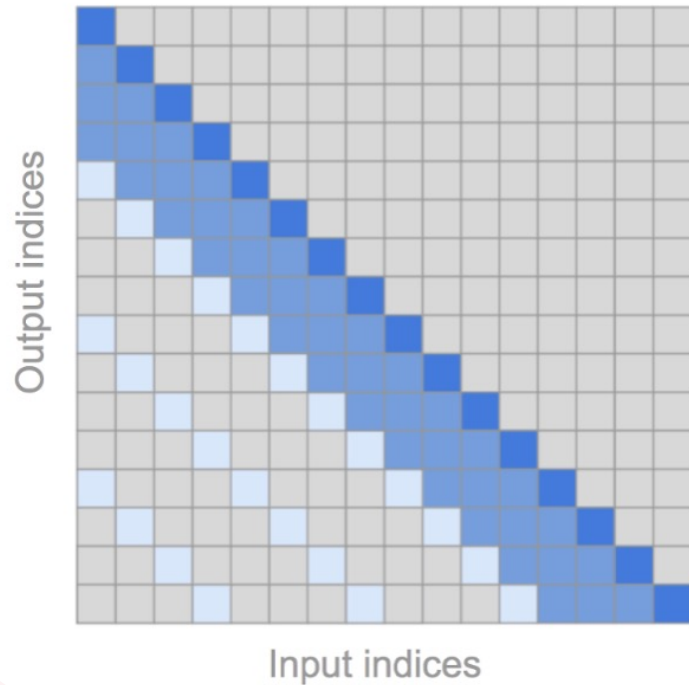
- $j \in S_i$ — напрямую
- $j \in A_{t_1}^1, t_1 \in A_{t_2}^2, \dots, t_k \in A_i$
есть маршрут $(j, t_1, t_2, \dots, t_k, i)$ не длиннее $p + 1$

Посмотрим $p = 2$

Factorized self-attention



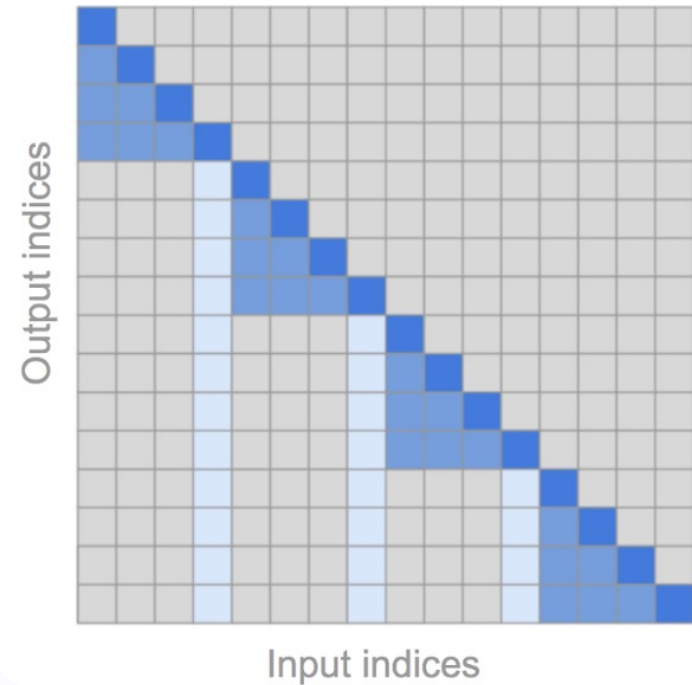
Factorized self-attention



$$A_i^{(1)} = \{t, t+1, \dots, i\}, \text{ where } t = \max(0, i - \ell)$$

$$A_i^{(2)} = \{j : (i - j) \bmod \ell = 0\}$$

$$\ell \sim \sqrt{n}$$



$$A_i^{(1)} = \{j : \lfloor \frac{j}{\ell} \rfloor = \lfloor \frac{i}{\ell} \rfloor\}$$

$$A_i^{(2)} = \{j : j \bmod \ell \in \{\ell - c, \dots, \ell - 1\}\}$$

Factorized self-attention

- $\text{Attention}(X) = \text{Attend}(X, A^{n \% p})W^o$
- $\text{Attention}(X) = \text{Attend}(X, \bigcup_{m=1}^p A^m)W^o$

Результаты

