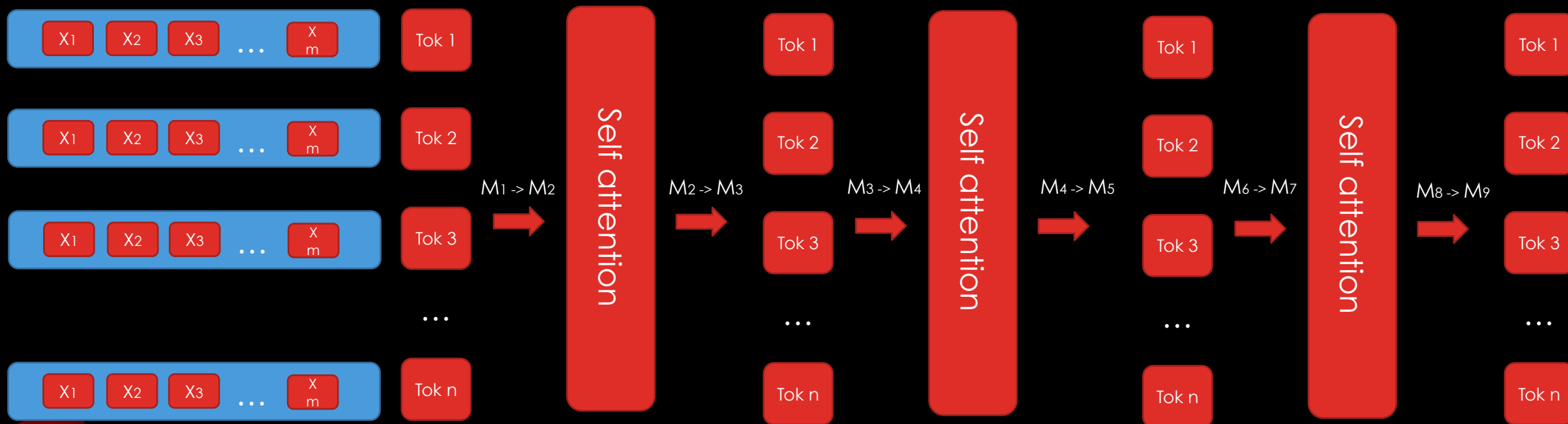


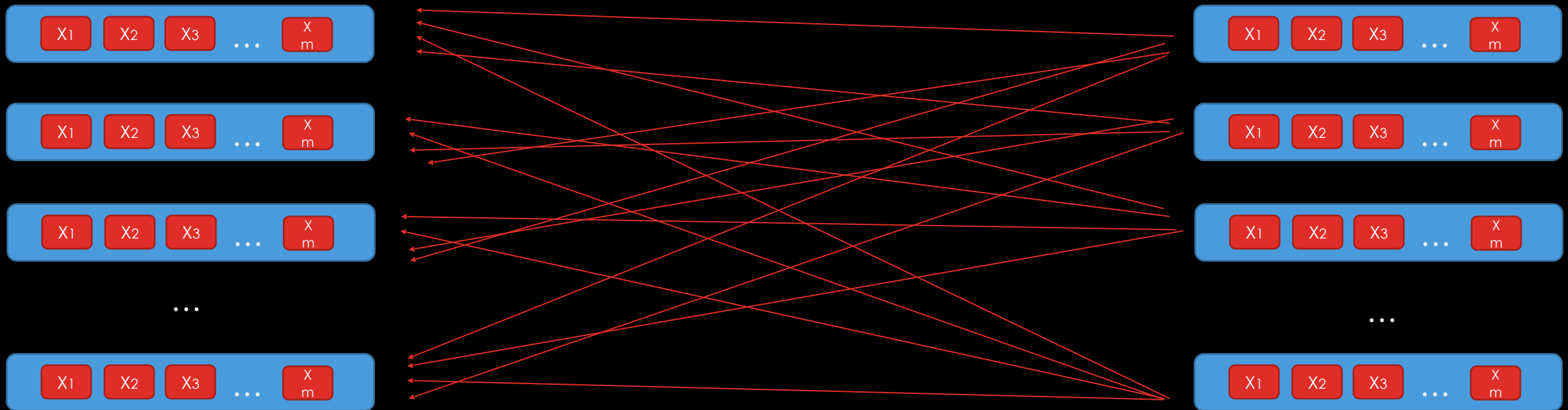
TRANSFORMERS ARE RNNs

(He советует)

TRANSFORMER



SELF ATTENTION



SELF ATTENTION

$$\text{Sim} \left(\begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ \vdots \\ x_m \end{bmatrix}, \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ \vdots \\ x_n \end{bmatrix} \right)$$

$$\text{Sim} (\star, \star) = 10$$

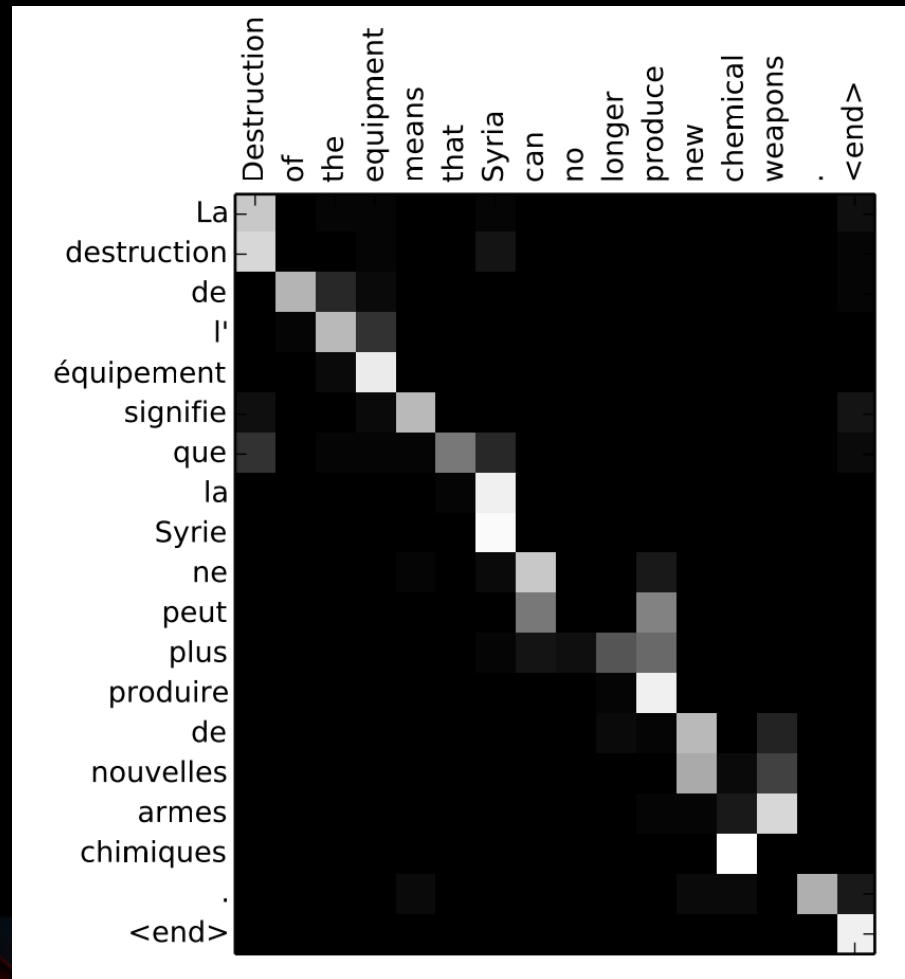
$$\text{Sim} (\star, \text{😊}) = -34$$

$$\text{Sim}(a, b) = a^T b$$

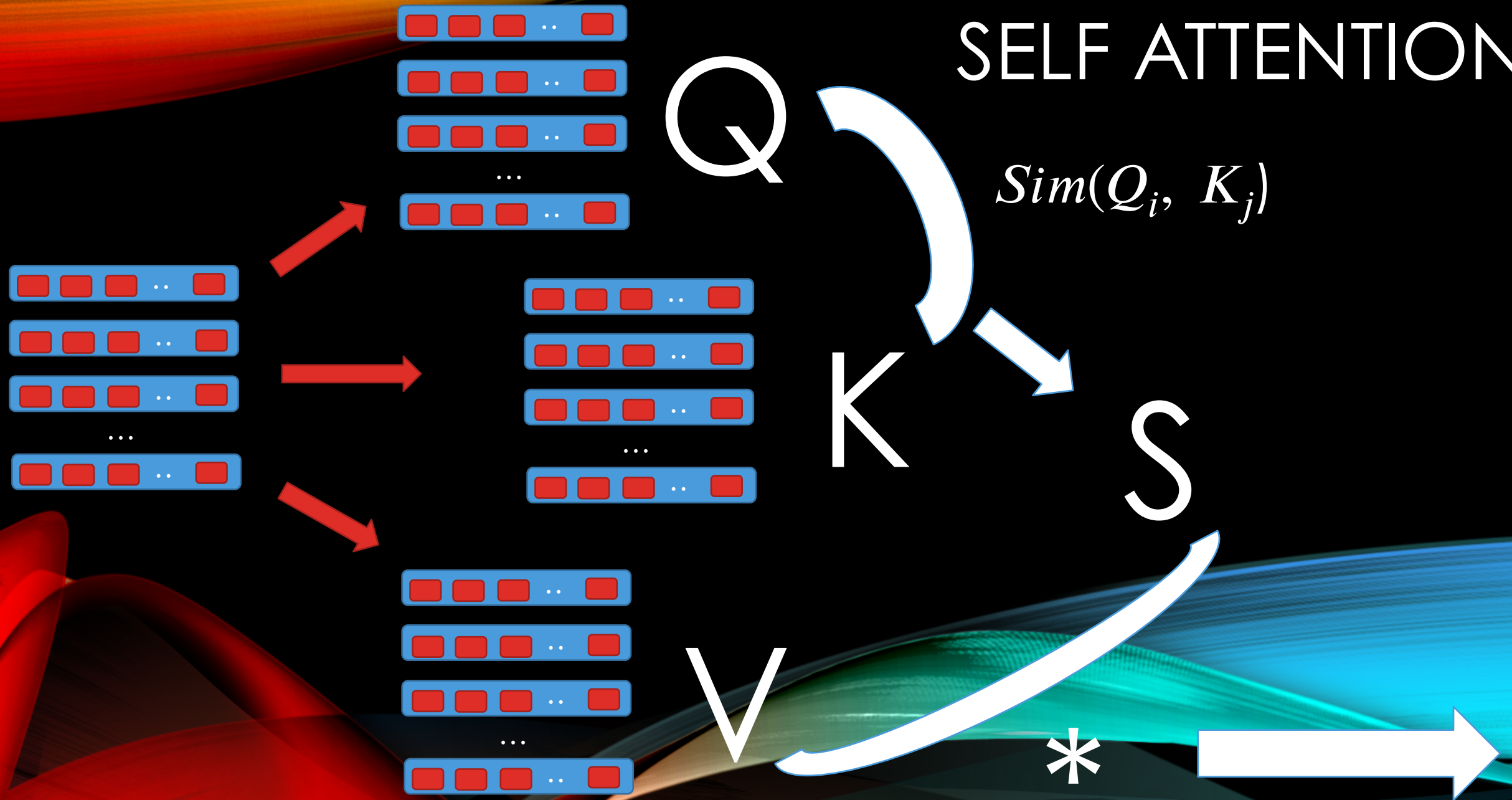
SELF ATTENTION

$$\begin{pmatrix} 10 & 9 & -20 & -34 \\ 10 & -34 & 29 & 0 \\ -29 & -56 & 99 & -3 \\ -22 & 13 & 13 & -33 \end{pmatrix} \xrightarrow{\text{softmax}} \begin{pmatrix} 0.6 & 0.4 & 0 & 0 \\ 0.2 & 0 & 0.8 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0.5 & 0.5 & 0 \end{pmatrix}$$

ATTENTION



SELF ATTENTION



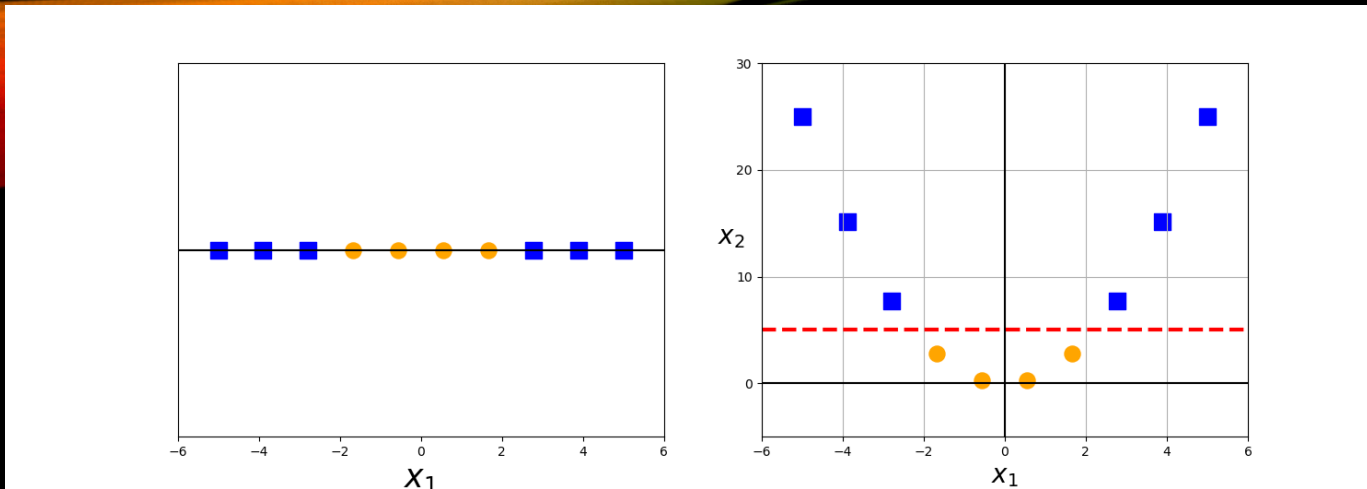
SELF ATTENTION

$$\text{Softmax}(Q^T K)V [O(n^2 d)]$$

SELF ATTENTION

$$\frac{Sim(Q, K)}{\sum_i Sim(Q, K)_i} V \quad Sim(a, b) = e^{a^T b}$$

KERNELS



$$\varphi(x) = (x, x^2)$$

$$K(x, y) = \varphi(x)^T \varphi(y) = (x, x^2)^T (y, y^2) = xy + x^2 y^2$$

TRANSFORMERS ARE RNNs

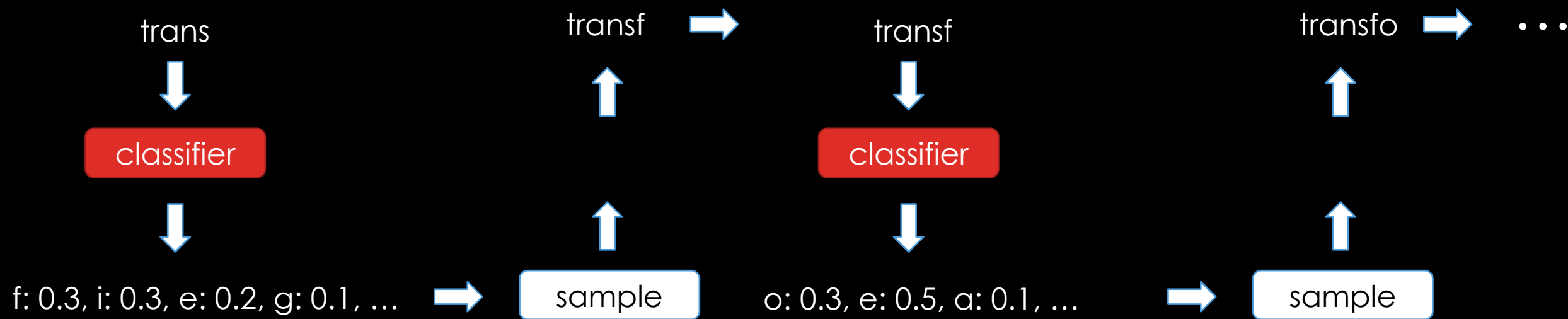
$$\frac{Sim(Q, K)}{\sum_i Sim(Q, K)_i} V \quad \frac{\varphi(Q)^T \varphi(K)}{\sum_i \varphi(Q)_i^T \varphi(K)_i} V$$

TRANSFORMERS ARE RNNs

$$\begin{aligned} X_k &= \frac{\varphi(Q)^T \varphi(K)}{\sum_i \varphi(Q)_i^T \varphi(K)_i} V_k = \frac{\sum_j \varphi(Q_k)^T \varphi(K_j) V_j}{\sum_j \varphi(Q_k)^T \varphi(K_j)} \\ &= \frac{\varphi(Q_k)^T \sum_j \varphi(K_j) V_j^T}{\varphi(Q_k)^T \sum_j \varphi(K_j)} [O(n^2 d) \rightarrow O(nd)] \end{aligned}$$

$$\varphi(x) = \text{elu}(x) + 1$$

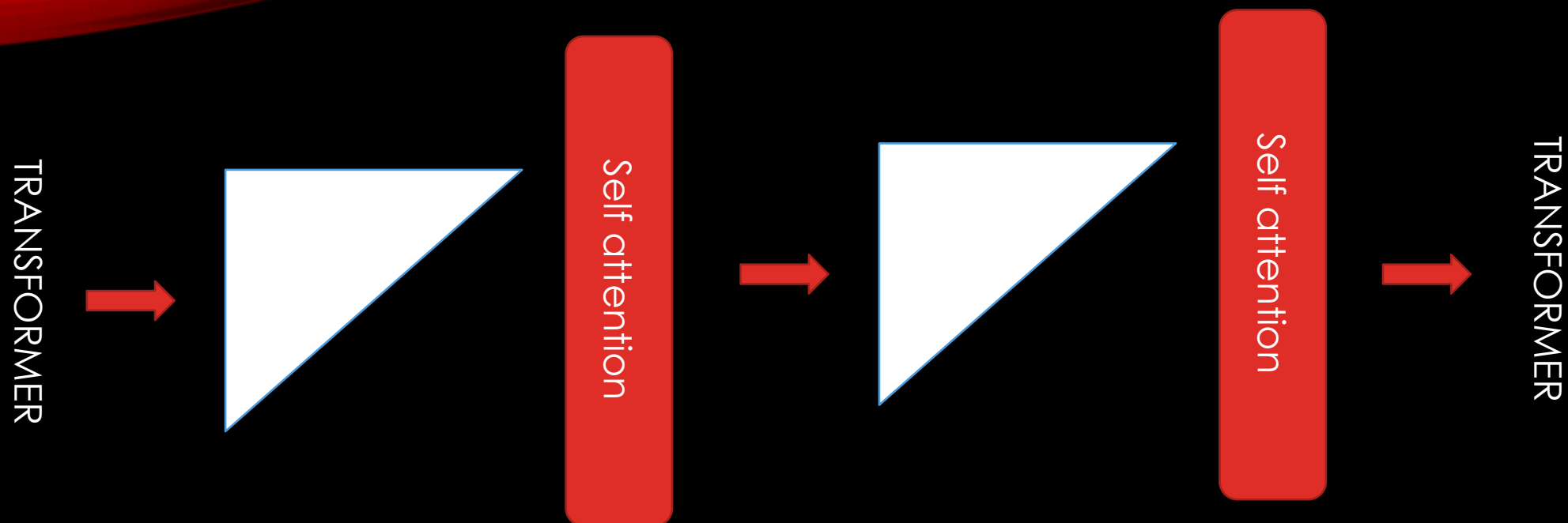
AUTOREGRESSIVE MODELS



$O(n^2d)$ - наивная имплементация kernel трансформером

$O(n^3d)$ - наивная имплементация softmax трансформером

TEACHER FORCING + CAUSAL MASKING

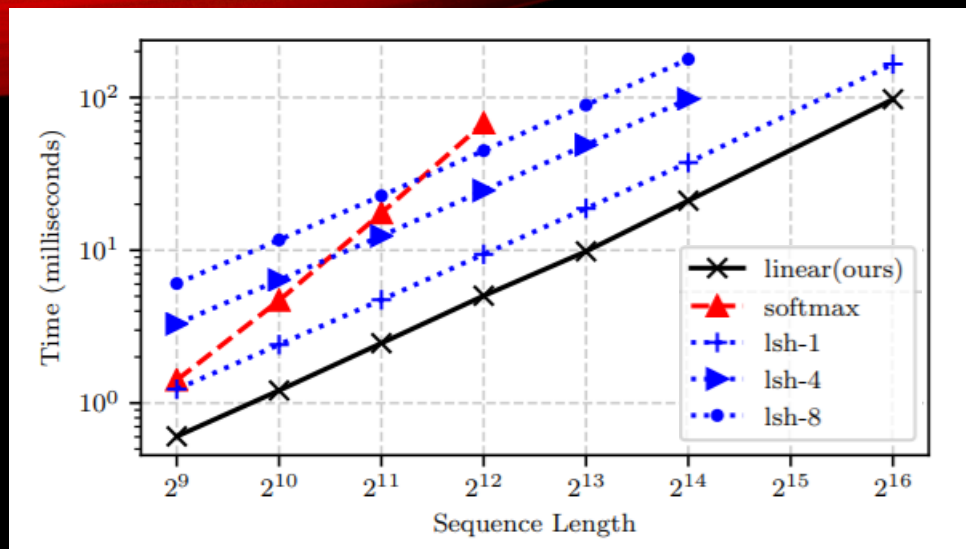


$O(nd)$ - наивная имплементация kernel трансформером

$O(n^2d)$ - наивная имплементация softmax трансформером

$O(nd)$ - RNN

TRANSFORMERS ARE RNNs



Speech recognition

Method	Validation PER	Time/epoch (s)
Bi-LSTM	10.94	1047
Softmax	5.12	2711
LSH-4	9.33	2250
Linear (ours)	8.08	824

MNIST autoregressive

Method	Bits/dim	Images/sec
Softmax	0.621	0.45 (1×)
LSH-1	0.745	0.68 (1.5×)
LSH-4	0.676	0.27 (0.6×)
Linear (ours)	0.644	142.8 (317×)

CIFAR-10 autoregressive

Method	Bits/dim	Images/sec
Softmax	3.47	0.004 (1×)
LSH-1	3.39	0.015 (3.75×)
LSH-4	3.51	0.005 (1.25×)
Linear (ours)	3.40	17.85 (4,462×)

TRANSFORMERS ARE RNNs

Unconditional samples



Image completion



(a)

(b)

(c)

Unconditional samples

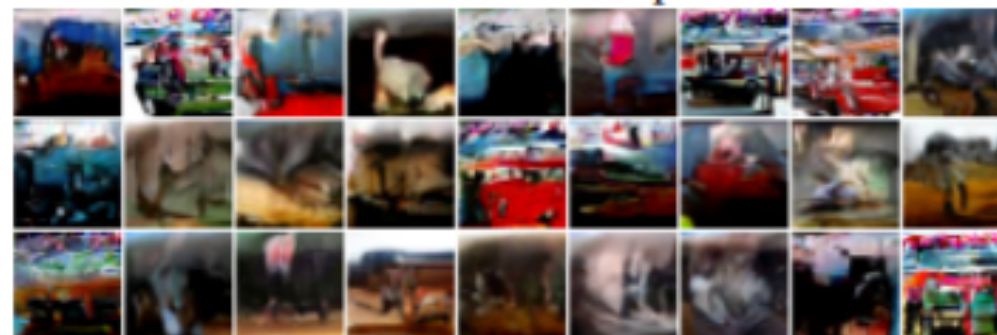
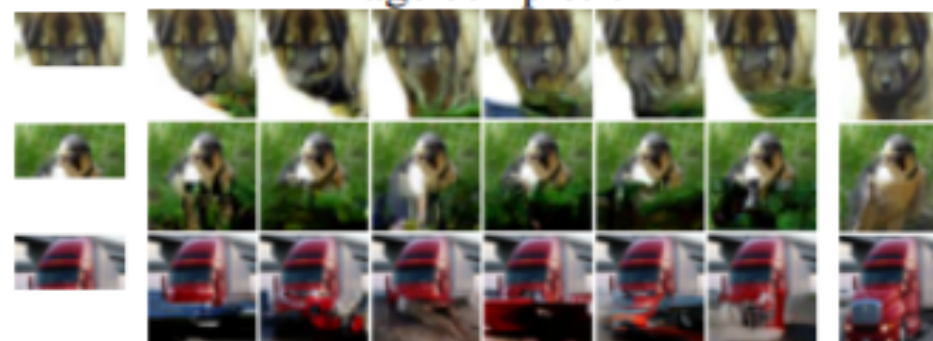


Image completion



(a)

(b)

(c)

TRANSFORMERS ARE RNNS

$$s_0 = 0,$$

$$z_0 = 0,$$

$$s_i = s_{i-1} + \phi(x_i W_K) (x_i W_V)^T,$$

$$z_i = z_{i-1} + \phi(x_i W_K),$$

$$y_i = f_l \left(\frac{\phi(x_i W_Q)^T s_i}{\phi(x_i W_Q)^T z_i} + x_i \right).$$

QUESTIONS

- 1) какова асимптотическая сложность обычного трансформера в зависимости от длины последовательности токенов n ? какова асимптотическая сложность rnn в такой-же ситуации? какова асимптотическая сложность трансформера из статьи?
- 2) выписать формулу attention стандартного softmax трансформера
- 3) выписать формулу attention произвольного трансформера через функцию похожести sim .
- 4) выписать формулу attention произвольного трансформера, через функцию добавляющую новые признаки соответствующие ядру трансформера (ядро наоборот) ϕ .

SOURCES

- 1) <https://arxiv.org/abs/2006.16236>
- 2) <https://arxiv.org/abs/1706.03762>
- 3) <https://www.youtube.com/watch?v=hAooAOFRsYc>