# Seq 2 Seq



$X$ - input sequence

$$P_\theta(Y|X) = \prod_{i=1}^{n} P_\theta(y_i|X, y_{<i}) \longrightarrow \max_\theta$$
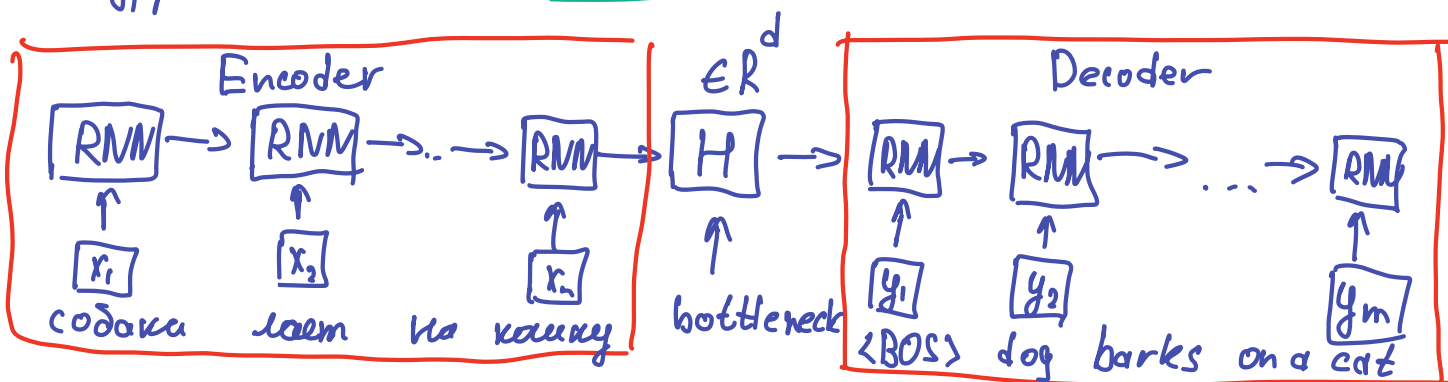
output
sequence

$$\hat{y}_i = \underset{y_i}{\text{argmax}}\; P_\theta(y_i|X, \hat{y}_{<i}) \quad \text{— Жадная генерация}$$

$$P_\theta(\hat{y}|x) = \prod_{i=1}^{n} \max_{y_i} P_\theta(y_i|X, \hat{y}_{<i}) \neq \max_y \prod_{i=1}^{n} P(y_i|X, y_{<i})$$
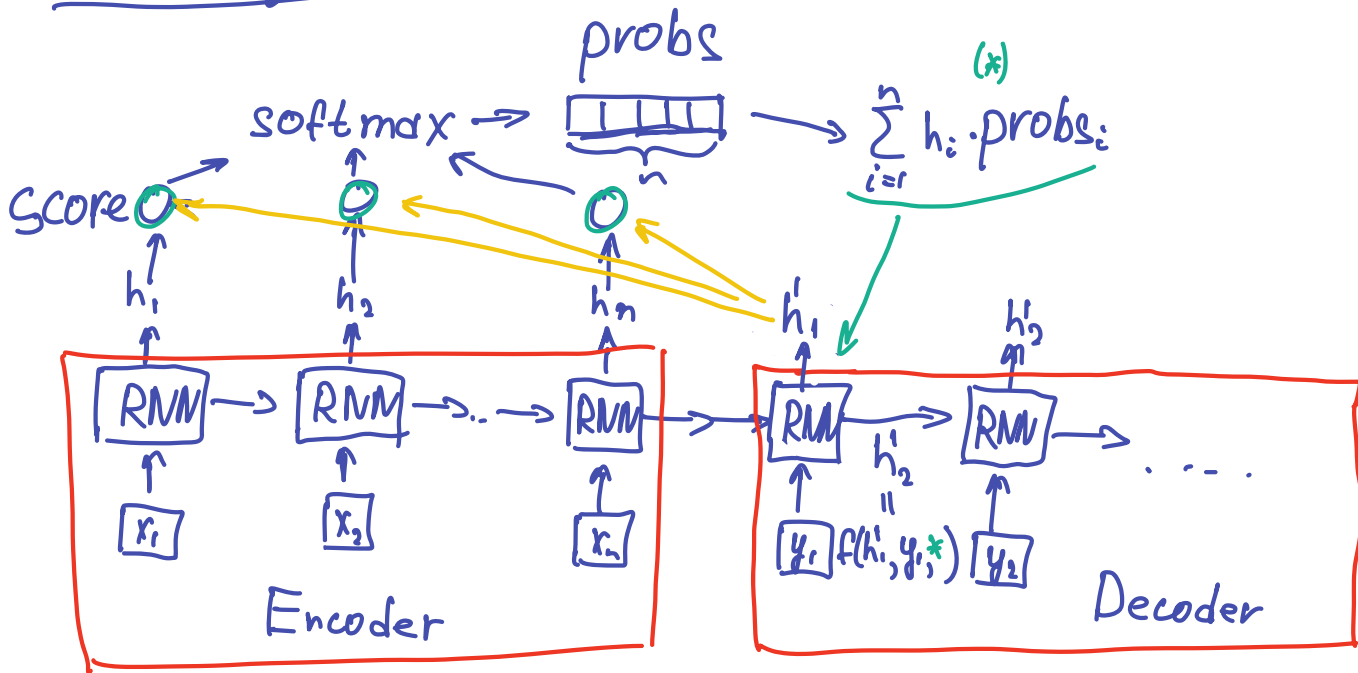
Для генерации лучше **Beam Search**

Рекуррентные сети <u>не позволяют</u> извлечь <u>всю</u> информацию

Обучение: Cross-Entropy

Механизм внимания. (2014)



probs

softmax →

$(*)$

$\sum_{i=1}^{n} h_i \cdot probs_i$

Score

$h_1$   $h_2$   $h_n$   $h'_1$   $h'_2$

RNN → RNN → ... → RNN → RNN → RNN → ....

$x_1$   $x_2$   $x_n$   $y_1$ $f(h'_i, y_i, *)$ $y_2$

Encoder     Decoder

$h'_2 \parallel$

Score functions:

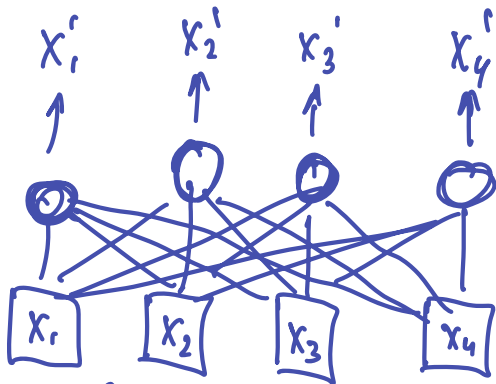$\in \mathbb{R}^d$   $\in \mathbb{R}^d$

1) Dot-product: $s(h, h') = h^T h'$

2) Bilinear: $s(h, h') = \underline{h^T W_s h'}$

$\mathbb{R}^d$   $\mathbb{R}^{d'}$

3) MLP: $s(h, h') = \vec{w_s}^T \tanh \left( W_s \begin{bmatrix} h^T \\ h'^T \end{bmatrix} \right)$

## Self-attention

Encoder

$$q_i = W_Q x_i$$
$$k_i = W_K x_i$$
$$v_i = W_V x_i$$

$x_1'$   $x_2'$   $x_3'$   $x_4'$

$x_1$   $x_2$   $x_3$   $x_4$

На двери замка висит замок

pos (0    1    2    3    4)
embed.

$$score_{ij} = \frac{q_i^T k_j}{\sqrt{d}}$$

$$attention_i = softmax(score_i)$$

$$x_i' = \sum_{j=1}^{n} attention_{ij} \underset{\uparrow \atop scalar}{v_j}$$

## Multi-head Attention

Идея: разобьём $q, k, v$ на $n_{heads}$ векторов. Каждая голова будет извлекать свою информацию. Потом склеим все выходы вместе.

$q^1$   $k^3$   $v^1$
$q^2$   $k^2$   $v^2$
$q^3$   $k^1$   $v^3$

$q$    $k$    $v$

$$attention_i^s = softmax\left(\frac{q_i^{s\,T} k_j^s}{\sqrt{d}}\right)$$

$$x_i^s = \sum_{j=1}^{n} attention_i^s v_j^s$$

$$\begin{bmatrix} x_i^{1'} \\ x_i^{2'} \\ x_i^{3'} \end{bmatrix} = x_i'$$

# Masked Self-Attention (Decoder)

$$\text{attention}_{ij} = \text{softmax}\left(\frac{q_i^T k_j}{\sqrt{d}} + \text{mask}_{ij}\right)$$

$$\text{mask}_{ij} = \begin{cases} 0, & i \leq j \\ -\infty, & i > j \end{cases}$$

$$\Rightarrow \forall i > j : \text{attention}_{ij} = 0$$

Зануляем внимание так, чтобы декодер при обучении не мог смотреть вперёд