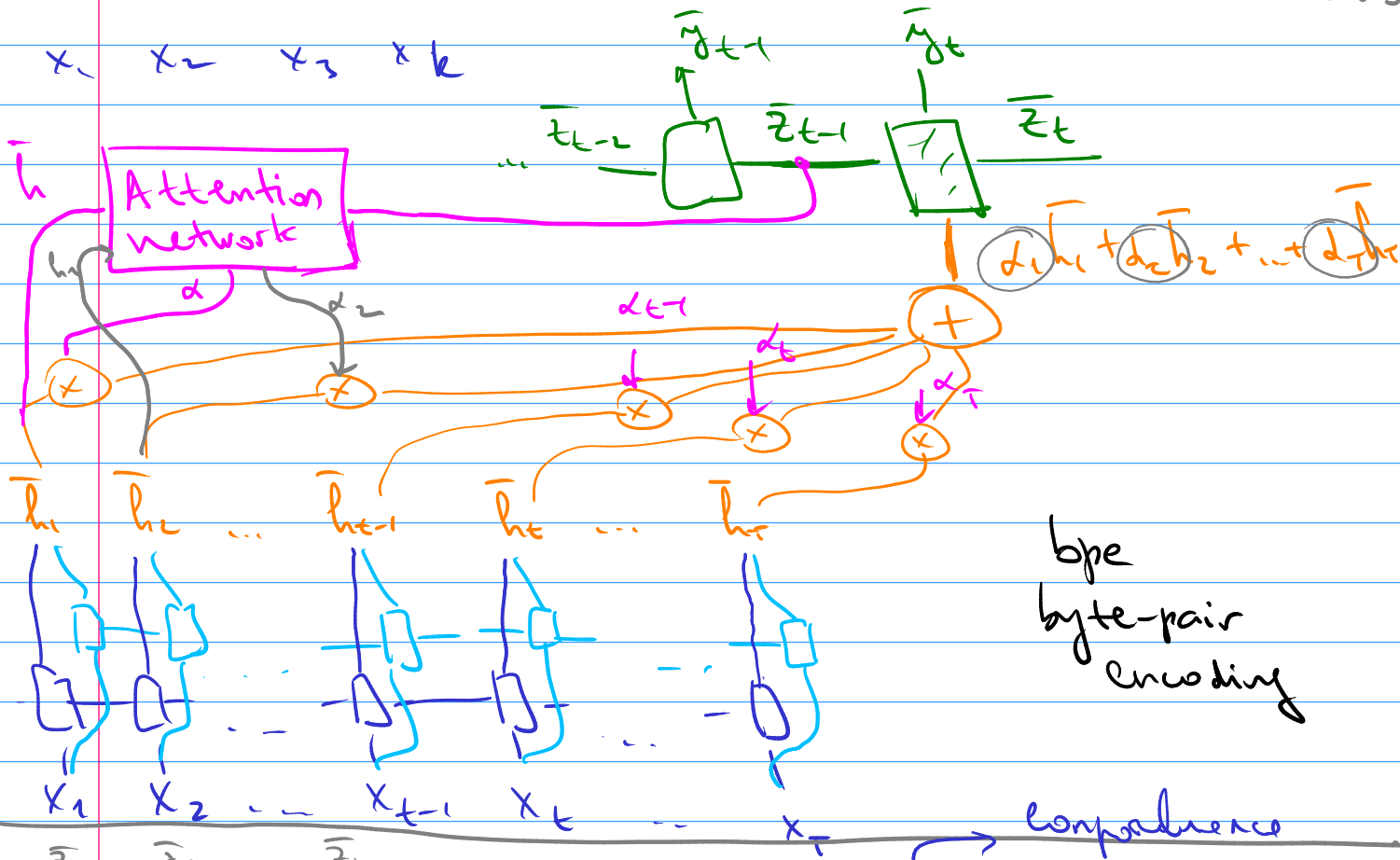
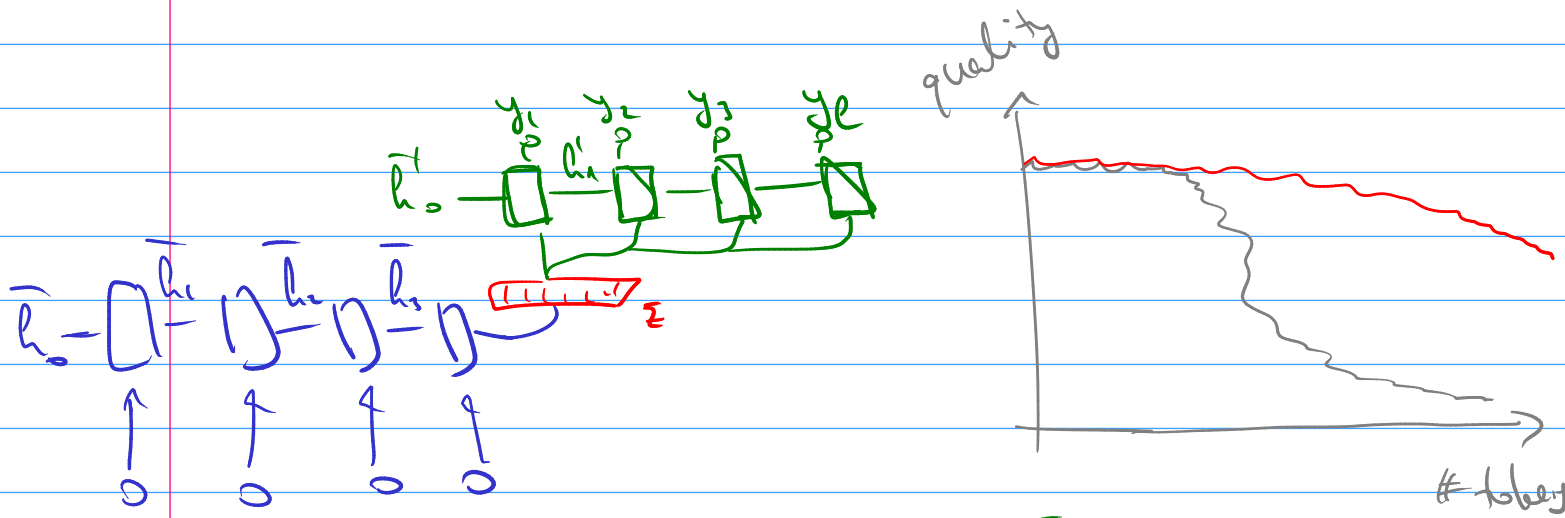
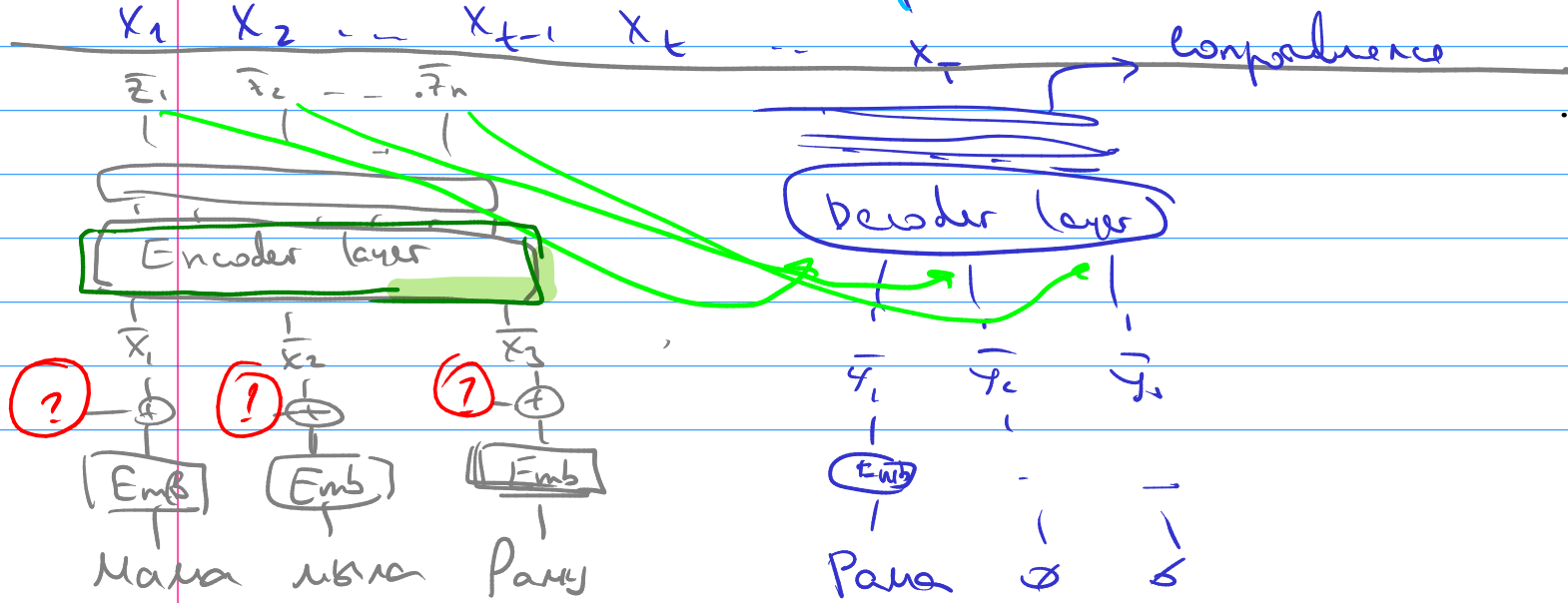
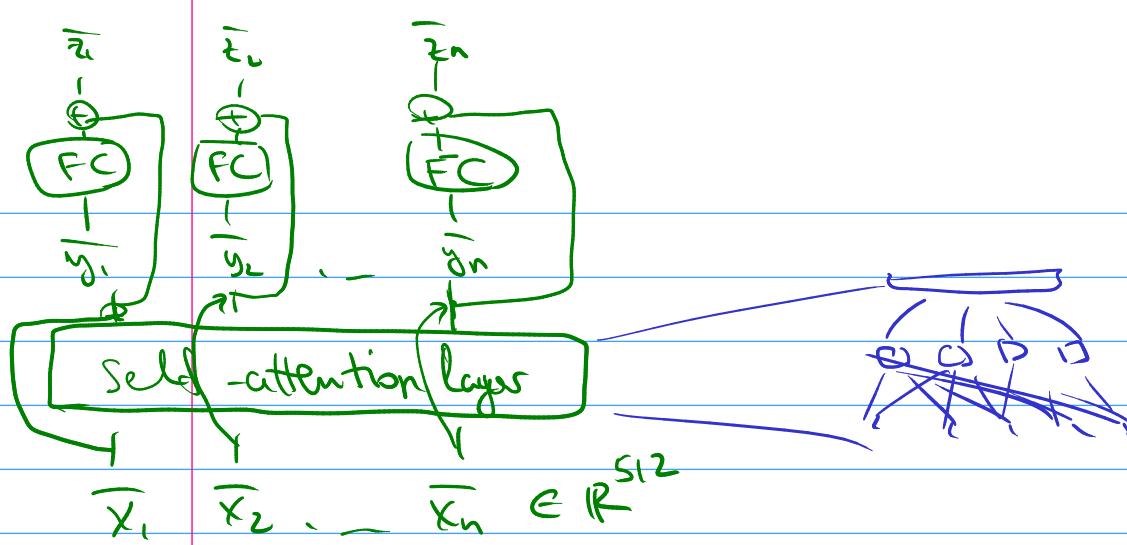


лучше неограниченно
 or price не 60
 or words



bpe
 byte-pair
 encoding





\mathbb{R}^d

$\bar{q}_i \in \mathbb{R}^m$ — query

$\bar{k}_i \in \mathbb{R}^m$ — key

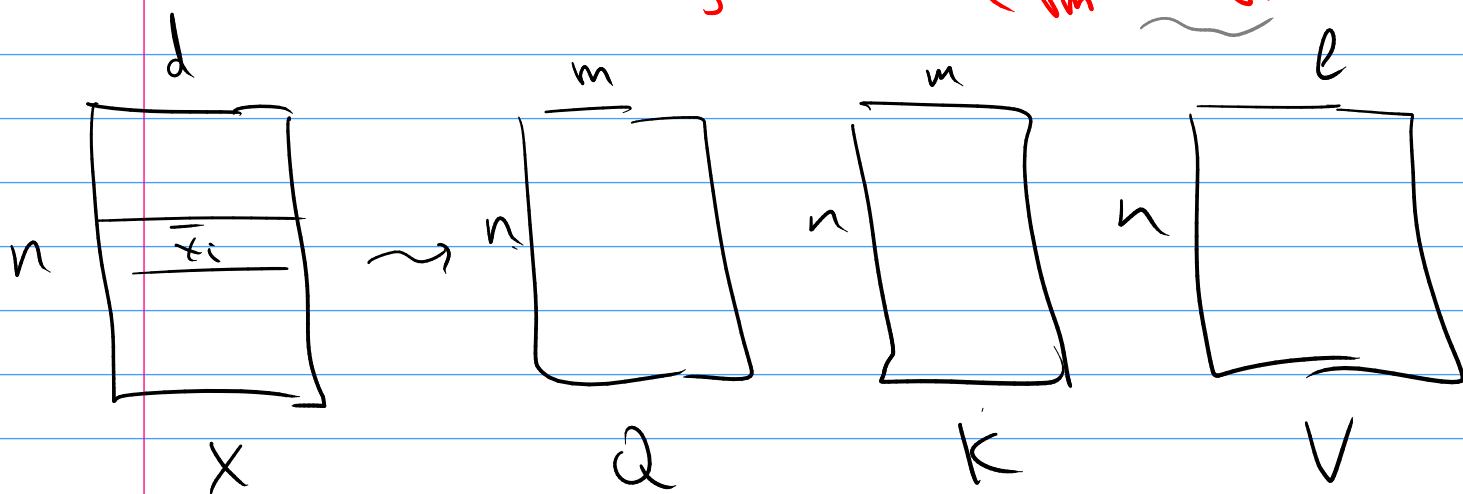
$\bar{v}_i \in \mathbb{R}^l$ — value

$$\begin{aligned} \bar{x}_1 &\rightarrow (\bar{q}_1, \bar{k}_1, \bar{v}_1) \rightarrow (\bar{q}_1^T \bar{k}_1) \\ \bar{x}_2 &\rightarrow (\bar{q}_2, \bar{k}_2, \bar{v}_2) \rightarrow (\bar{q}_2^T \bar{k}_2) \\ &\vdots \\ \bar{x}_n &\rightarrow (\bar{q}_n, \bar{k}_n, \bar{v}_n) \rightarrow (\bar{q}_n^T \bar{k}_n) \end{aligned}$$

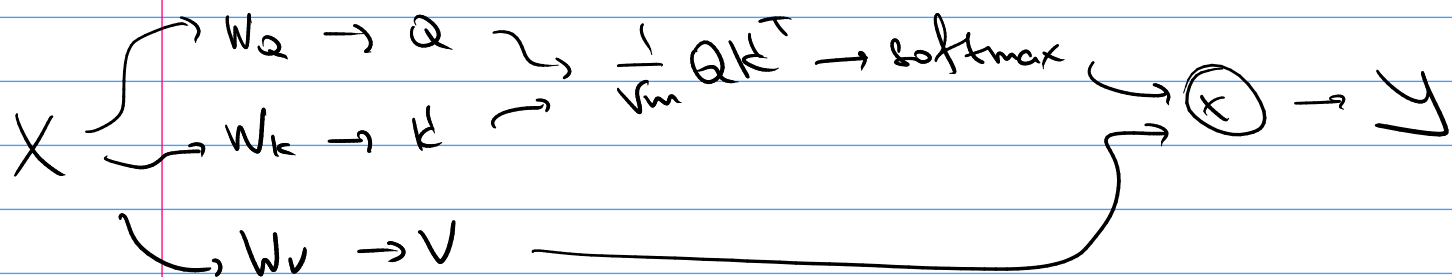
$$a_{ij} \approx \bar{q}_i^T \bar{k}_j \quad \text{„насколько } \bar{x}_i \text{ похож на } \bar{x}_j \text{“}$$

$$\bar{x}_i \rightarrow (\bar{q}_i, \bar{k}_i, \bar{v}_i) \rightarrow \bar{y}_i = \sum_{j=1}^n a_{ij} \bar{v}_j,$$

$$a_{ij} = \text{softmax}\left(\frac{1}{\sqrt{m}} \bar{q}_i^T \bar{k}_j\right)$$

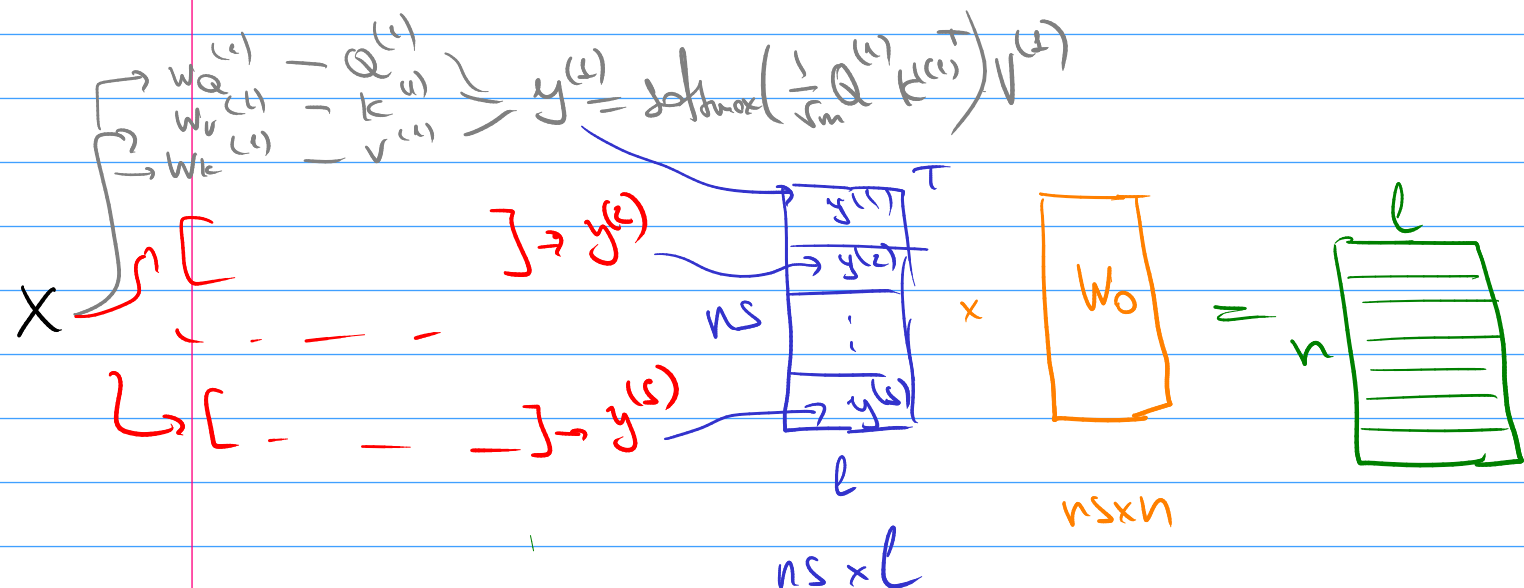


$$\begin{aligned} \underline{W_Q}: \quad \bar{q}_i &= W_Q \bar{x}_i & Q &= W_Q X \\ \underline{W_K}: \quad \bar{k}_i &= W_K \bar{x}_i & K &= W_K X \\ \underline{W_V}: \quad \bar{v}_i &= W_V \bar{x}_i & V &= W_V X \end{aligned}$$



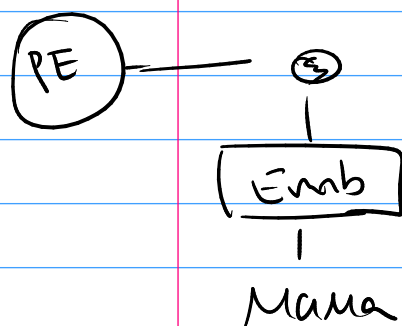
Multi-head attention

$$(W_Q^{(1)}, W_K^{(1)}, W_V^{(1)}) \dots (W_Q^{(s)}, W_K^{(s)}, W_V^{(s)})$$



$$\begin{aligned} W_O & \quad ns \times n \\ W_K, W_Q & \quad [d \times m] \\ W_V & \quad [d \times l] \end{aligned}$$

$$\mathcal{O}(n^2 + dm + dl) \text{ heads}$$



$$\begin{aligned} PE(\text{pos}, 2i) &= \sin\left(\left(\frac{\text{pos}}{10000}\right)^{2i/d}\right) \\ PE(\text{pos}, 2i+1) &= \cos\left(\left(\frac{\text{pos}}{10000}\right)^{2i/d}\right) \end{aligned}$$

Decoder

