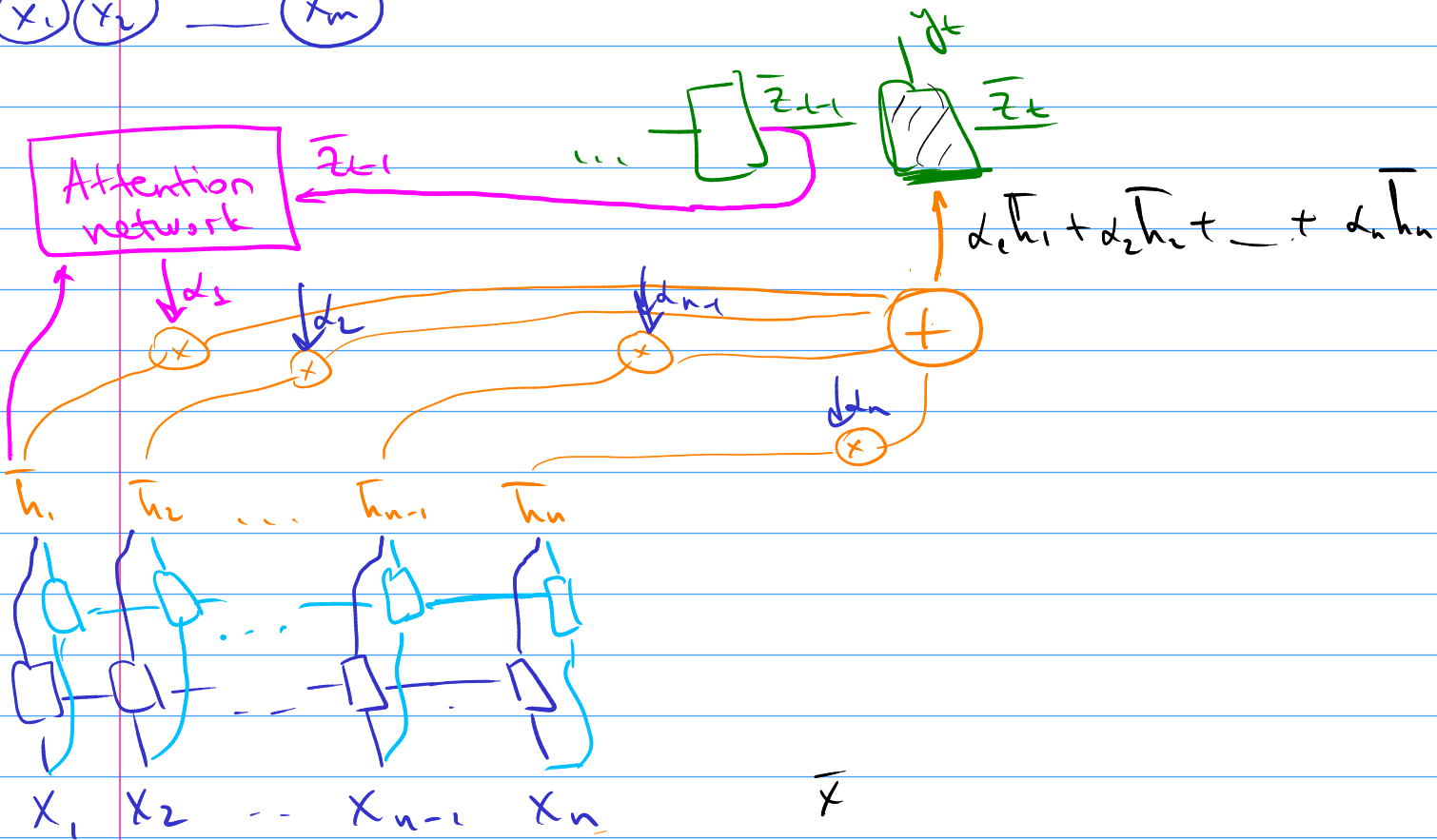
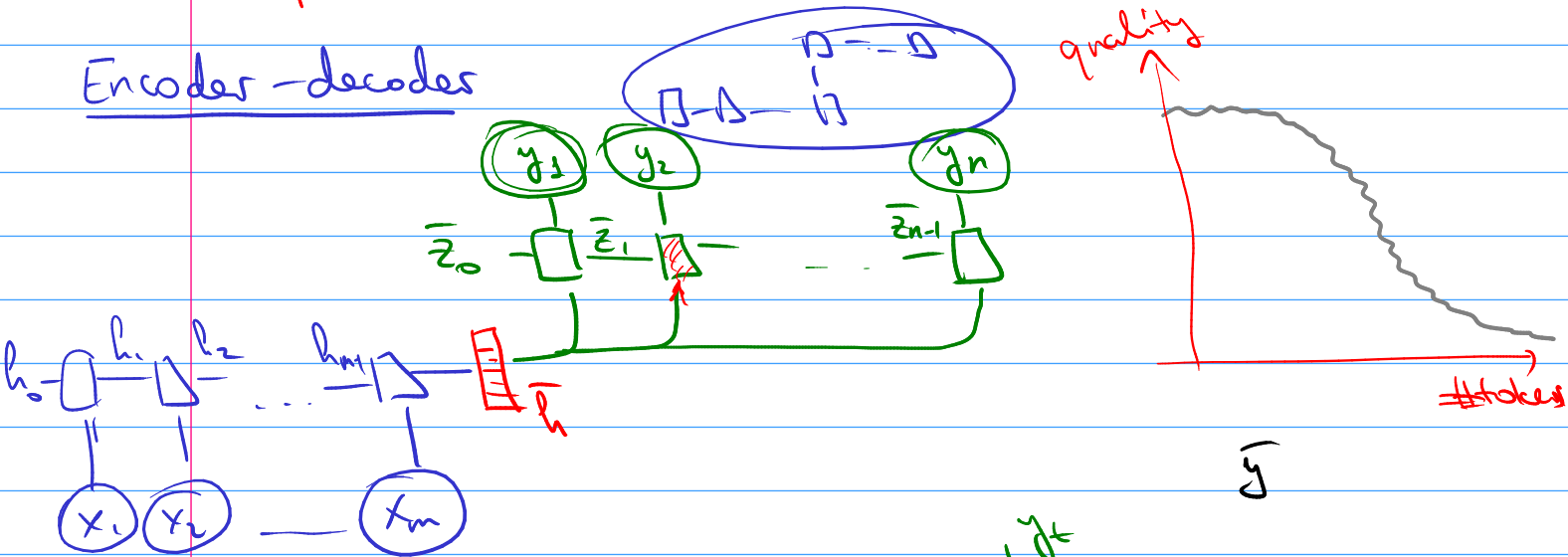
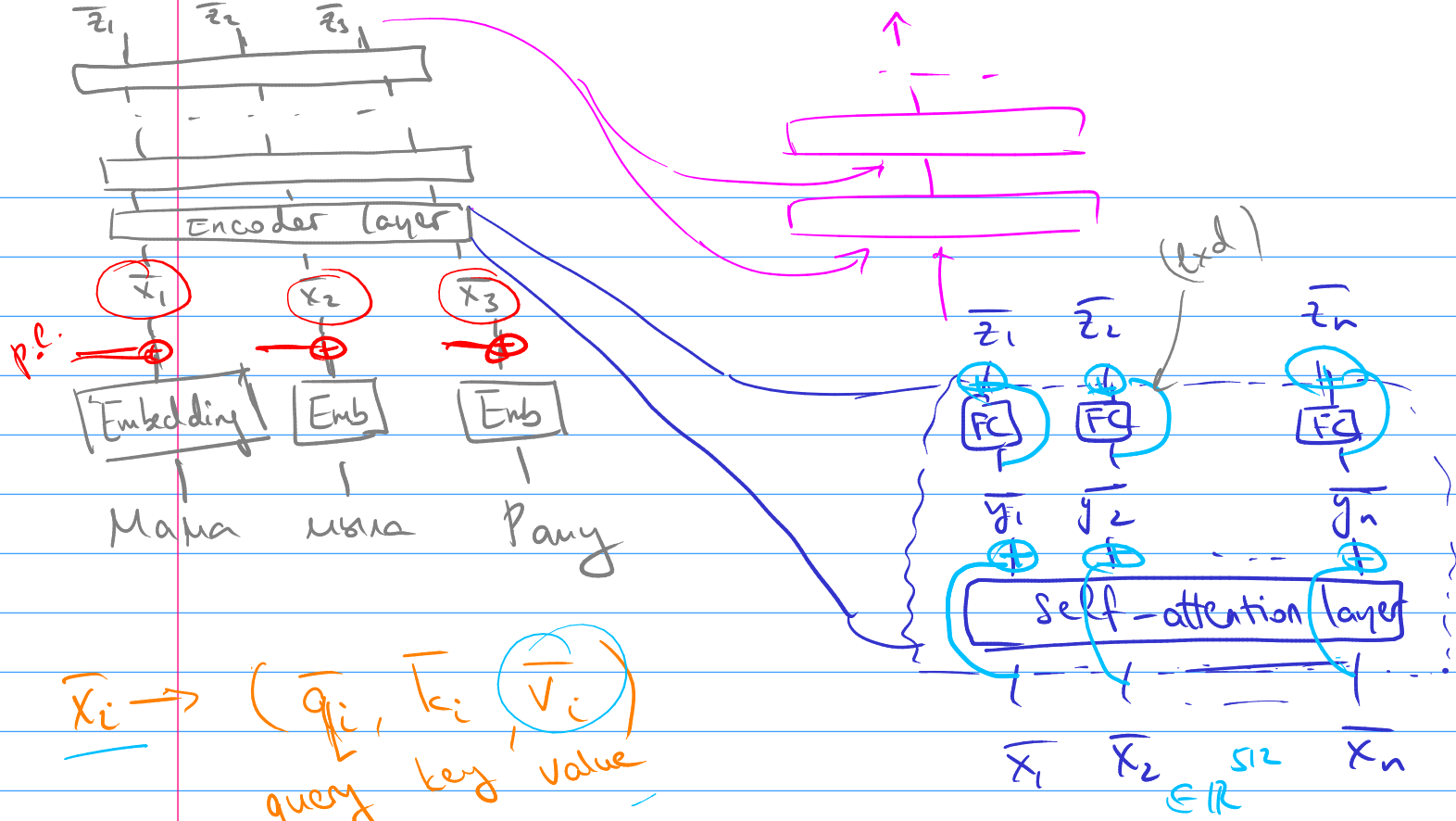


лучше, проще, дешевле \bar{x}
 \bar{x} нежелателен, делит на \bar{x}
 \bar{x} нежелателен

$$p(y_1, \dots, y_n | x_1, \dots, x_m) = p(y_1 | \bar{x}) p(y_2 | y_1, \bar{x}) p(y_3 | y_1, y_2, \bar{x}) \dots p(y_n | y_{1:n-1}, \bar{x})$$

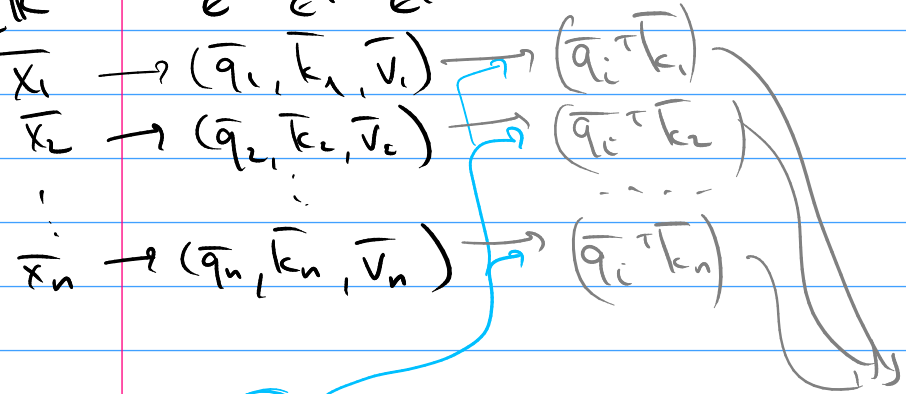
Encoder-decoder



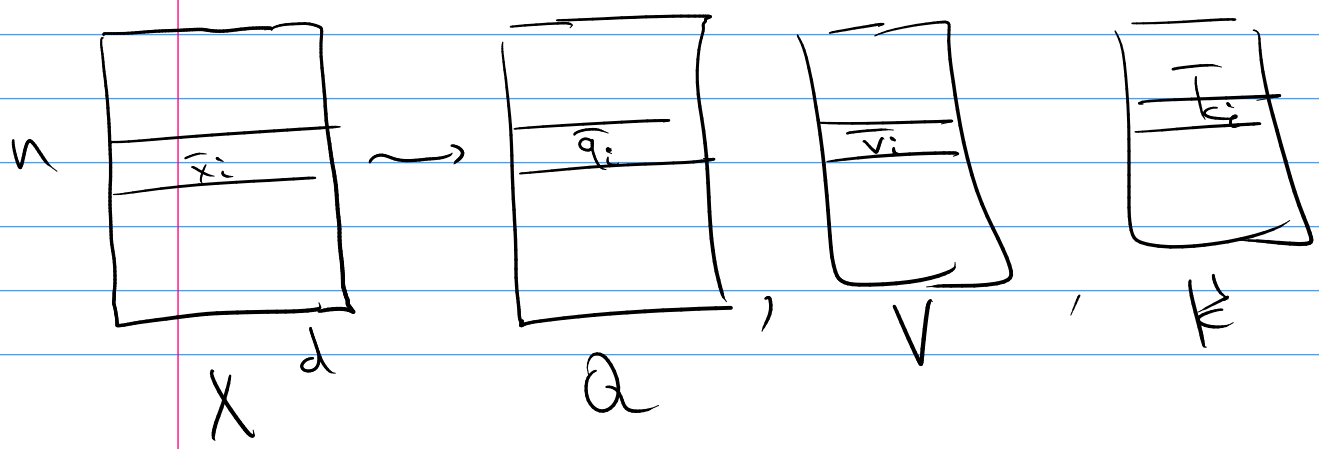


$\bar{x}_i \rightarrow (q_i, k_i, v_i)$
 query key value

$a_{ij} \approx (q_i^T k_j)$
 $\in \mathbb{R}^d$ $\in \mathbb{R}^m$ $\in \mathbb{R}^m$ $\in \mathbb{R}$
 "relevance"



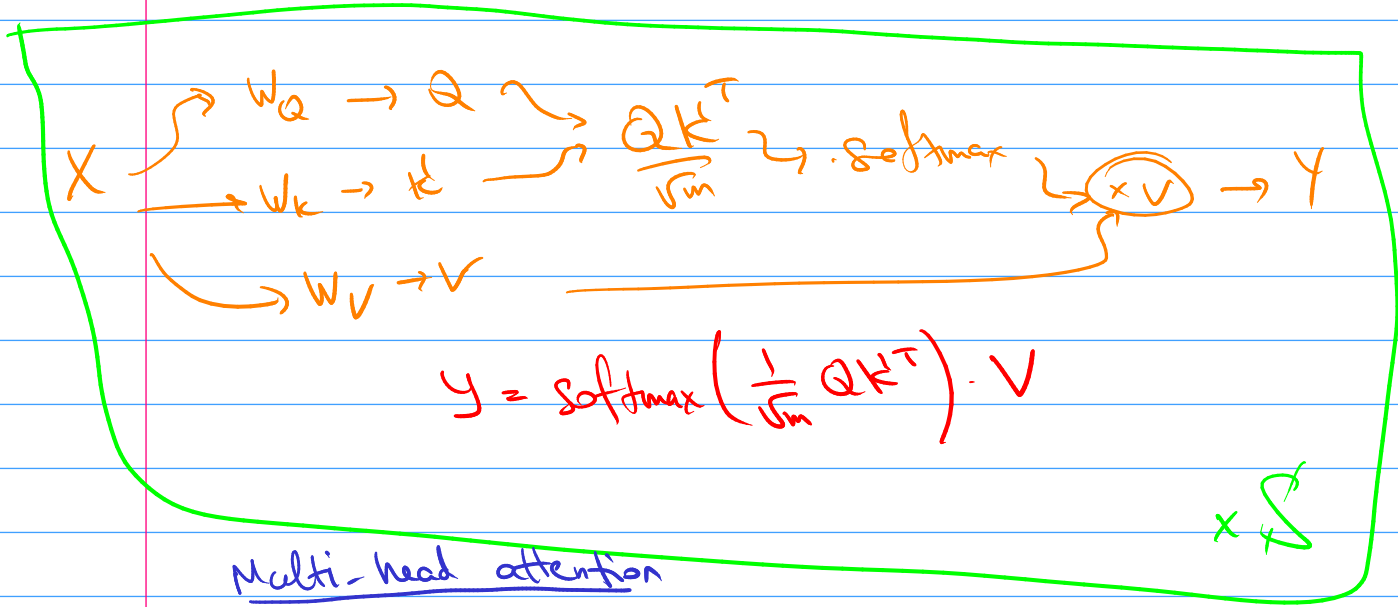
$\bar{x}_i \rightarrow (q_i, k_i, v_i) \rightarrow y_i = \sum_{j=1}^n \text{softmax}\left(\frac{1}{\sqrt{m}} q_i^T k_j\right) v_j$



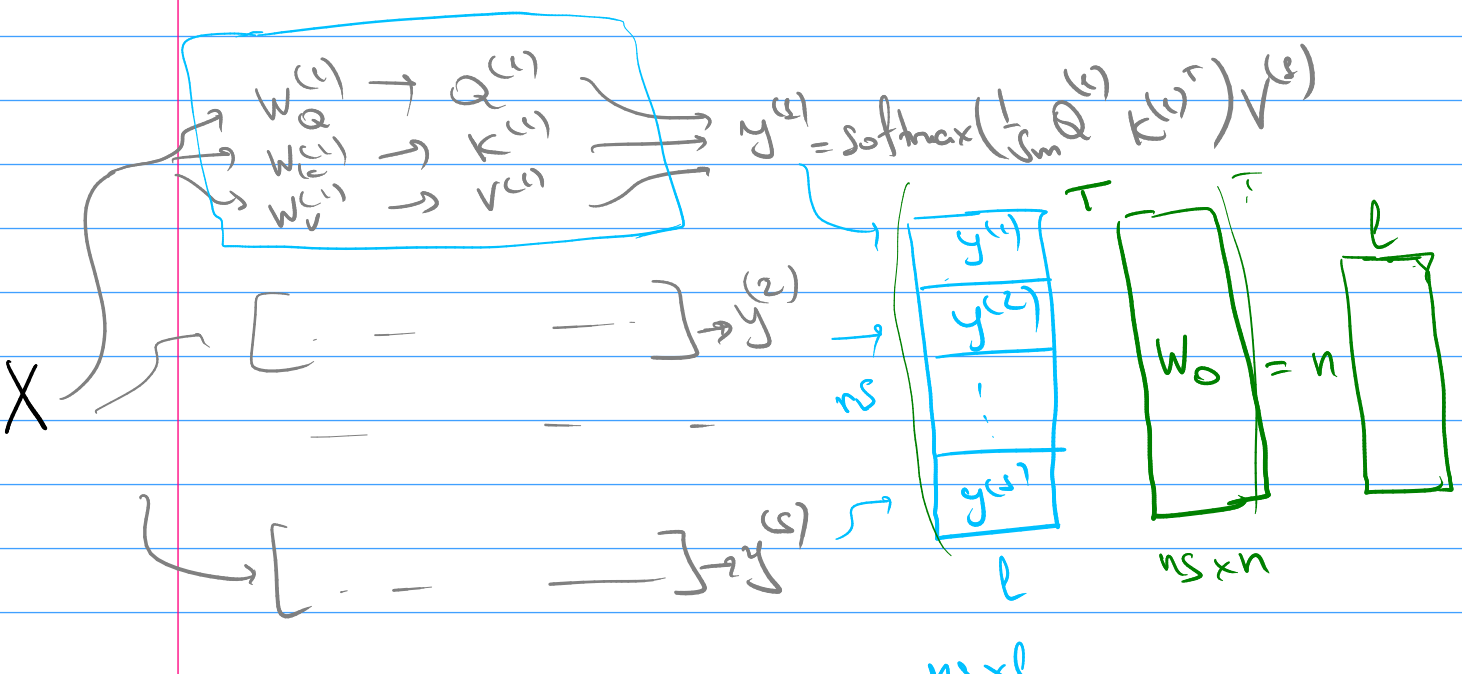
$$W_Q: \quad \bar{q}_i = W_Q \bar{x}_i, \quad Q = W_Q X$$

$$W_K: \quad K = W_K X$$

$$W_V: \quad V = W_V X$$

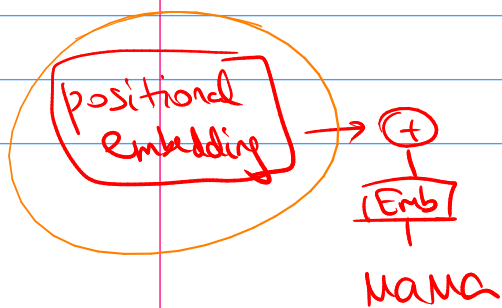


$$(W_Q^{(1)}, W_K^{(1)}, W_V^{(1)}) \quad \dots \quad (W_Q^{(s)}, W_K^{(s)}, W_V^{(s)}) \quad \dots$$



$$\left. \begin{array}{l} W_0 - ns \times n \\ W_Q, W_K - d \times m \times s \\ W_V - d \times l \times s \end{array} \right\} \rightarrow S(n^2 + dm + dl) \text{ params}$$

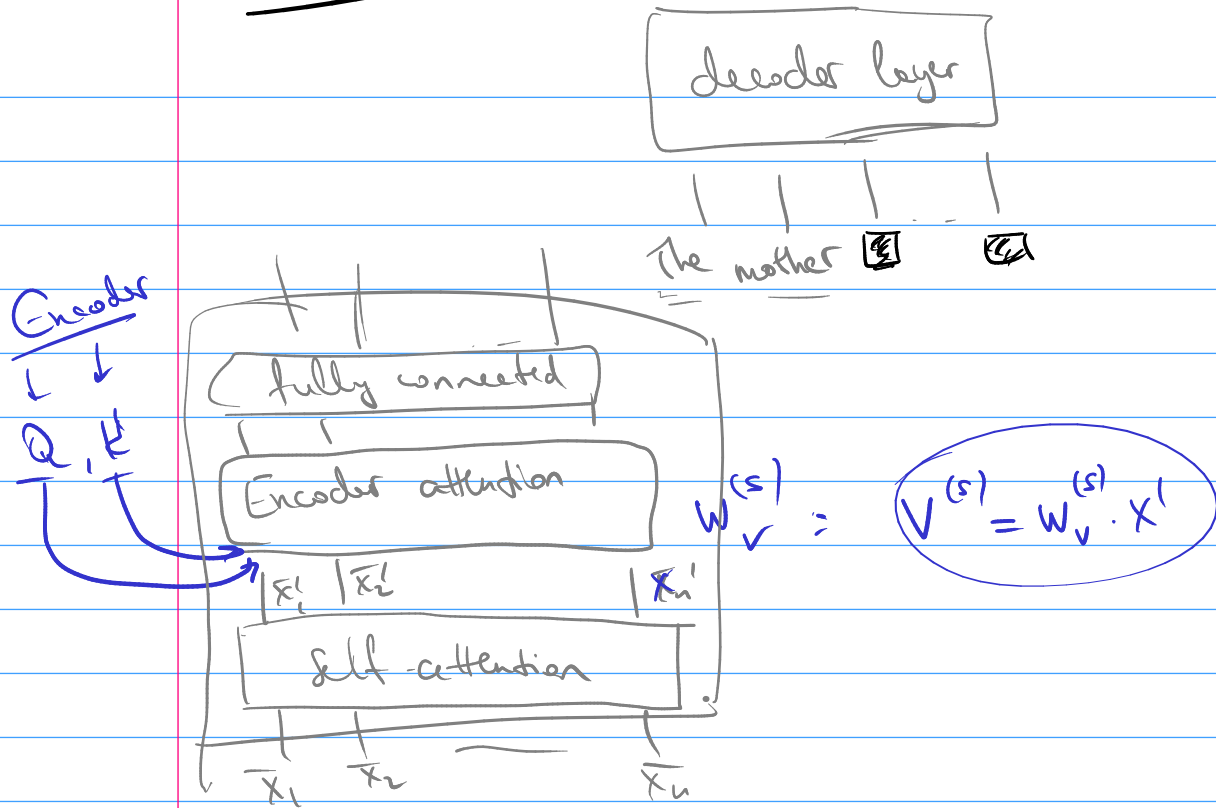
$\begin{matrix} \uparrow & \uparrow & \uparrow & \uparrow \\ 8 & 512 & 64 & 64 \end{matrix}$



$$PE(\text{pos}, 2i) = \sin\left(\left(\frac{\text{pos}}{10000}\right)^{2i/d}\right)$$

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

Decoder



Transformer Enc + Dec

GPT - Dec
language models

BERT Enc