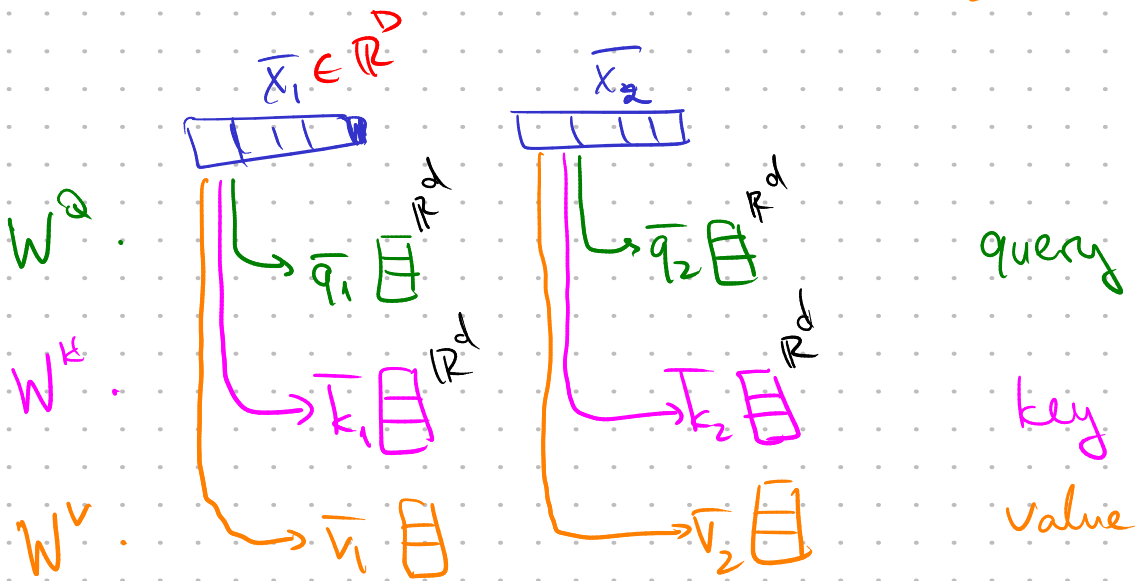
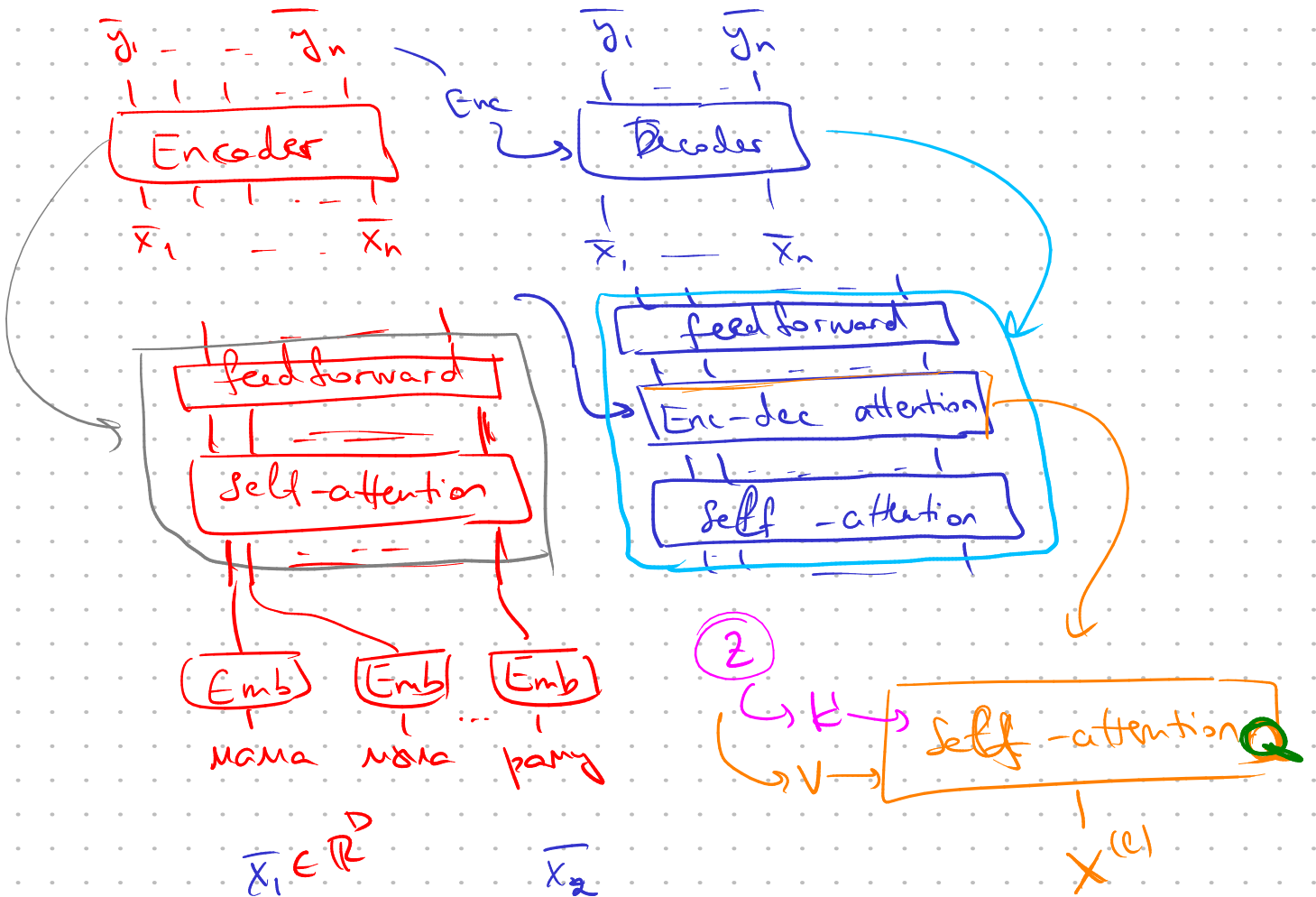
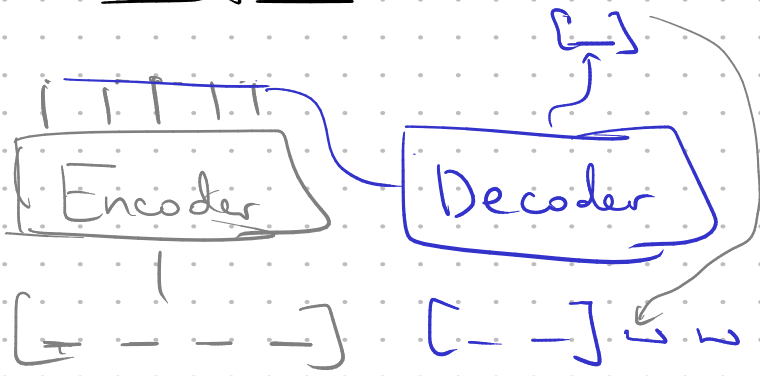
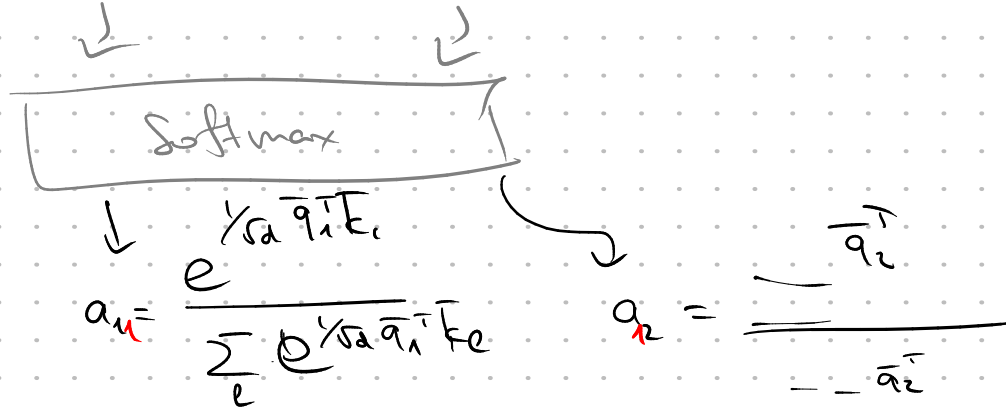


Transformer



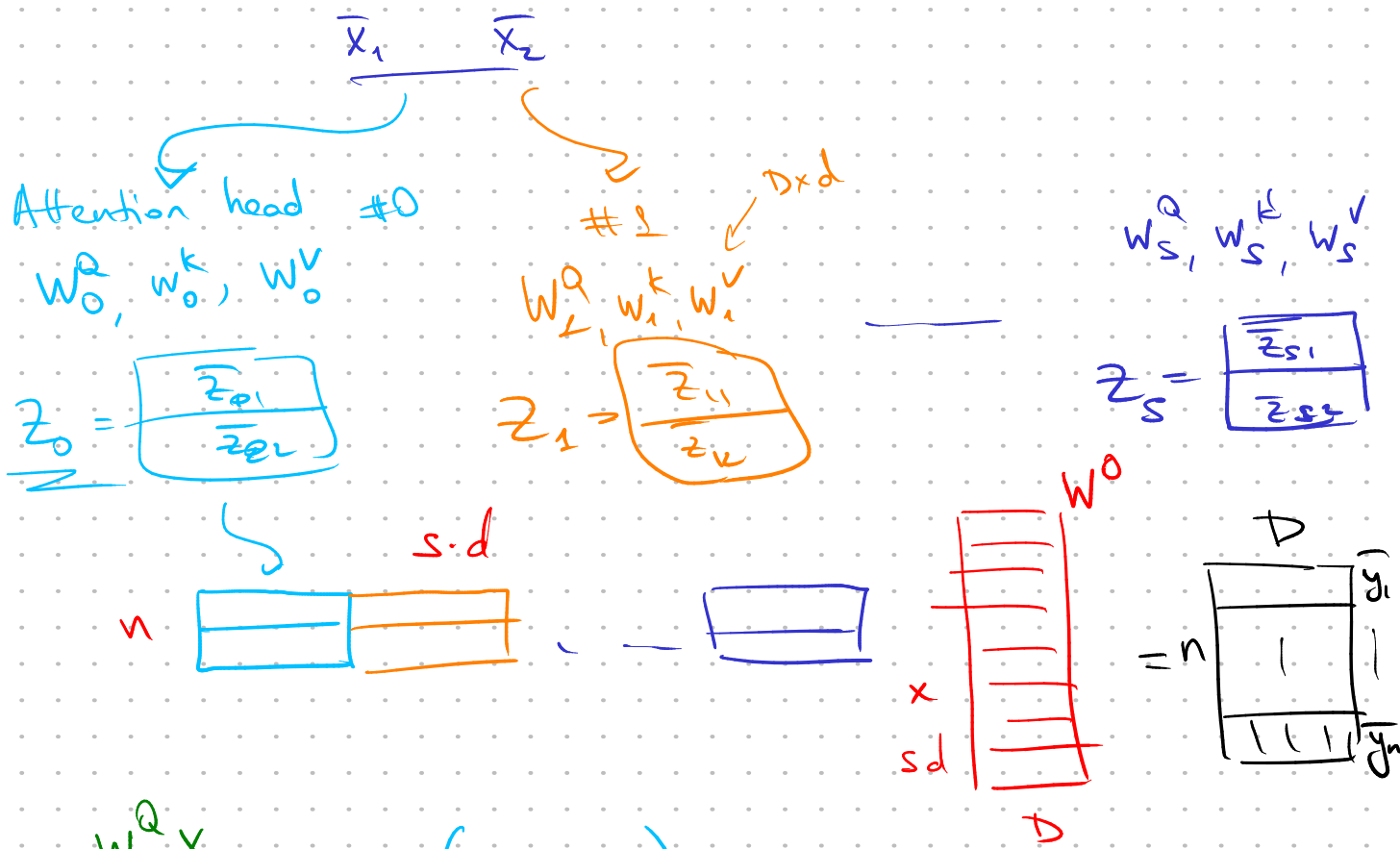
$$\frac{1}{\sqrt{d}} q_1^T k_2$$

$$\frac{1}{\sqrt{d}} q_2^T k_1$$



$$\bar{z}_1 = \sum a_{1e} \bar{v}_e = \sum \text{softmax}(q_1^T K_e) \cdot \bar{v}_e$$

$$\bar{z}_2 = \sum a_{2e} \bar{v}_e$$



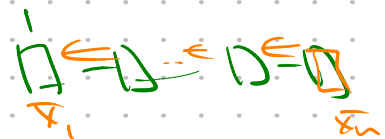
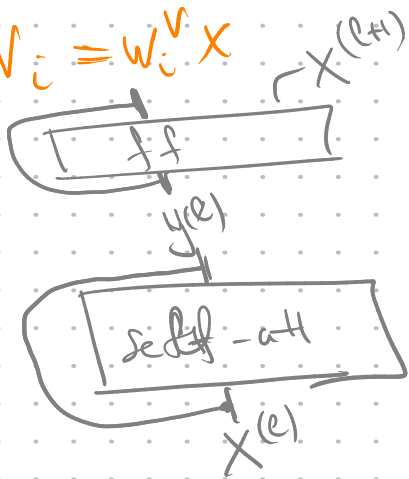
$$Q_i = W_i^Q X$$

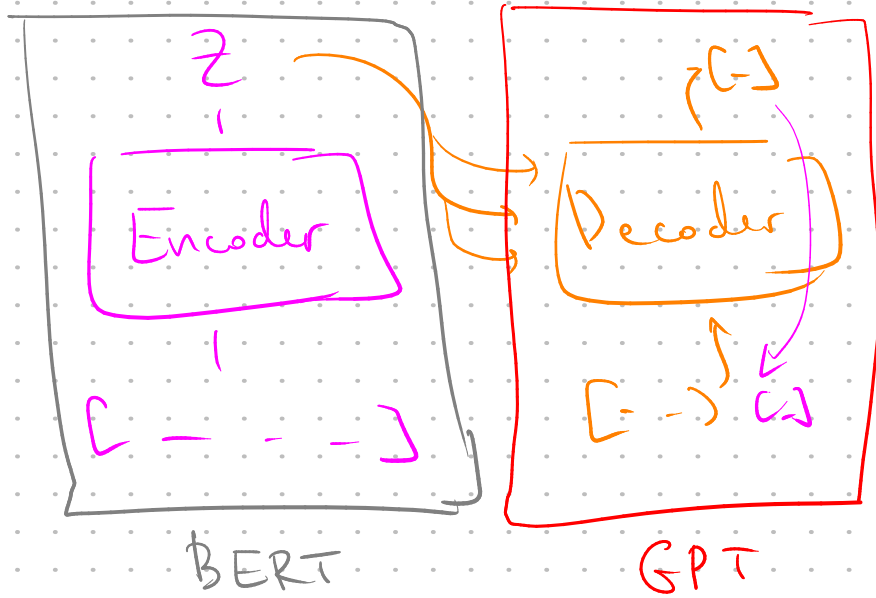
$$K_i = W_i^K X$$

$$V_i = W_i^V X$$

$$\text{softmax}\left(\frac{1}{\sqrt{d}} Q_i K_i^T\right) \cdot V_i = z_i$$

$$\text{concat}(z_1, \dots, z_s) \cdot W^0 = y$$





$$\bar{u} \sim \mathcal{N}(\bar{\mu}, \bar{\Sigma}) \quad \rightarrow \quad \odot \sigma_x \oplus \mu_x \rightarrow \bar{z} \sim \mathcal{N}(\bar{\mu}_x, \bar{\Sigma}_x)$$

Gumbel-Max trick

$$\bar{z} \sim \text{Mult}(\bar{\pi})$$

Then $z \sim \text{Mult}(\bar{\pi}) \Leftrightarrow z = \underset{i}{\text{argmax}} (g_i + \log \pi_i)$

$$p(g_i) = e^{-(g_i + e^{-g_i})} \quad \text{if } g_i \sim \text{Gumbel}$$

$$F(g_i) = e^{-e^{-g_i}}$$

Δ. ∞ $p(z=k) = p(g_k + \log \pi_k \geq g_j + \log \pi_j \quad \forall j) =$

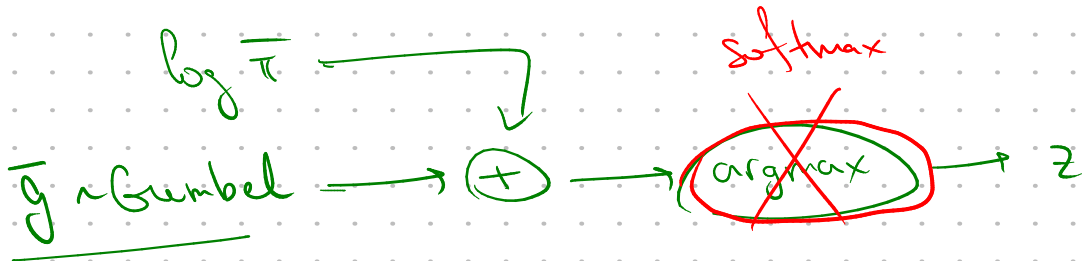
$$= \int \prod_{j \neq k} p(g_k + \log \pi_k \geq g_j + \log \pi_j | g_k) \cdot p(g_k) dg_k =$$

" $F(g_k + \log \pi_k - \log \pi_j)$

$$= \int \prod_{j \neq k} e^{-e^{-g_k - \log \pi_k + \log \pi_j}} \cdot e^{-g_k + e^{-g_k}} dg_k$$

$$= \int e^{-\sum_{j \neq k} \pi_j \cdot e^{-g_k - \log \pi_k}} \cdot \pi_k e^{-(g_k + \log \pi_k + \pi_k e^{-g_k - \log \pi_k})} dg_k$$

$$= \pi_k \int_{-\infty}^{\infty} e^{-g_k - \log \pi_k} = e^{-g_k - \log \pi_k} (\pi_k + \sum_{j \neq k} \pi_j) dg_k = \pi_k$$



Gumbel-Softmax

$$z = \underset{i}{\text{softmax}} \left(\frac{1}{\tau} (g_i + \log \pi_i) \right)$$

