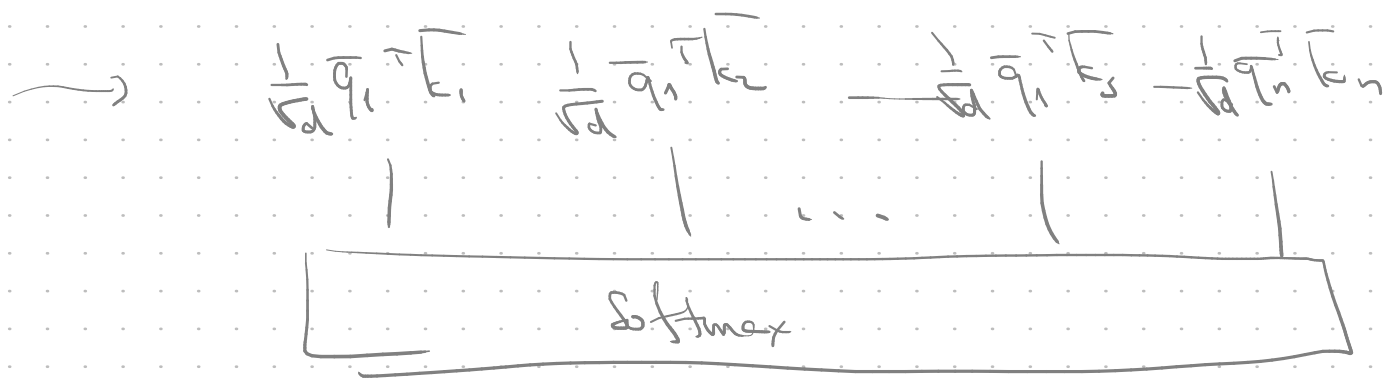
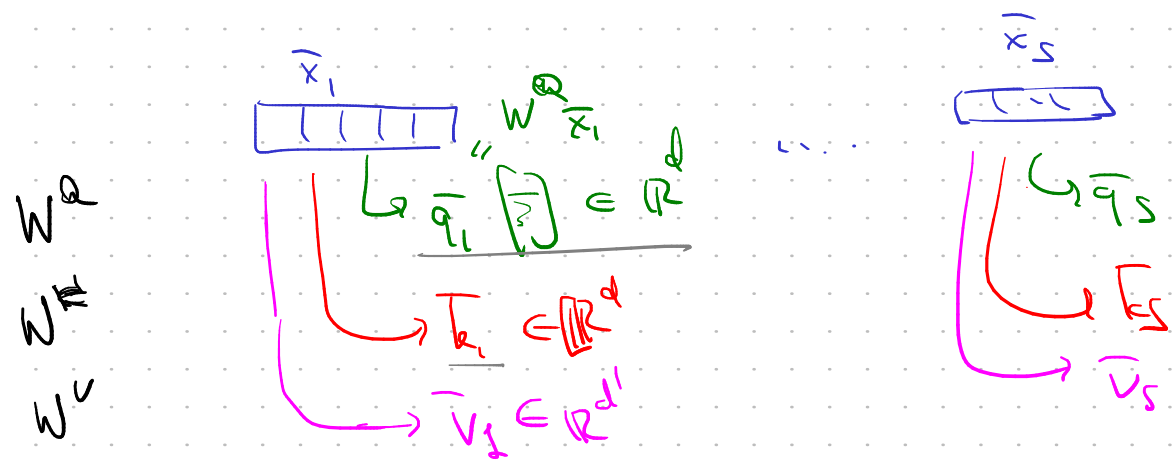
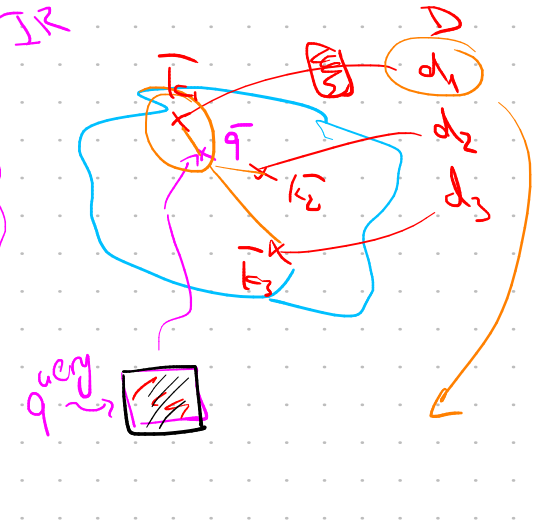
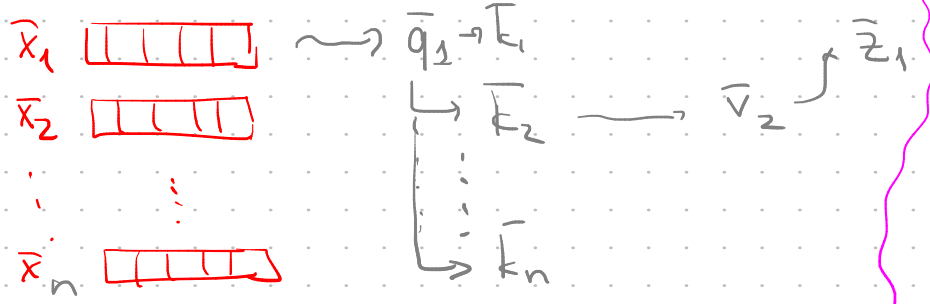


TRANSFORMER



$$a_{11} \quad a_{12} \quad \dots \quad a_{1n} = \frac{e^{\frac{1}{\sqrt{d}} \bar{q}_1^T \bar{k}_n}}{\sum_l e^{\frac{1}{\sqrt{d}} \bar{q}_1^T \bar{k}_l}}$$

$$\bar{z}_1 = \sum a_{1l} \bar{v}_l = \sum \text{softmax} \left(\frac{1}{\sqrt{d}} \bar{q}_1^T \bar{k}_l \right) \bar{v}_l$$

$$\bar{z}_2 = \sum a_{2l} \bar{v}_l$$

$$\bar{z}_n = \sum a_{nl} \bar{v}_l$$

Multi-head attention

$\bar{x}_1 \quad \bar{x}_2 \quad \dots \quad \bar{x}_n \quad X \quad n \times m$

Head #0

W_0^Q, W_0^K, W_0^V

Head #1

W_1^Q, W_1^K, W_1^V

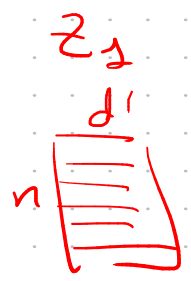
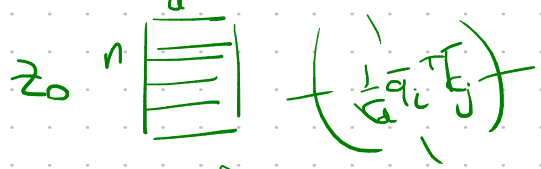
Head #7

W_7^Q, W_7^K, W_7^V

$$\bar{z}_{0i} = \text{softmax}_s \left(\frac{1}{\sqrt{d}} \bar{q}_{0i} \bar{k}_{0s} \right) \bar{v}_{0s}$$

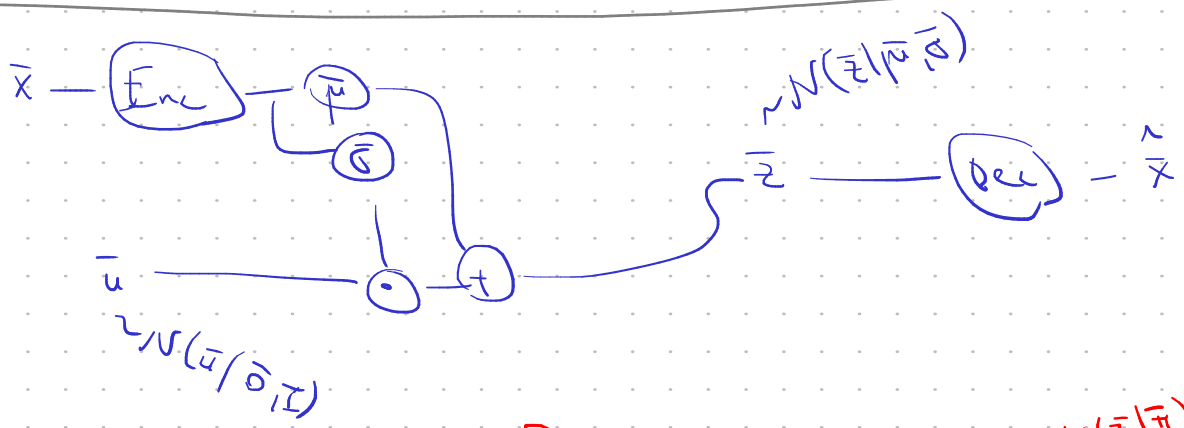
$$\bar{z}_0 = \text{softmax} \left(\frac{1}{\sqrt{d}} Q_0 K_0^T \right) V_0$$

nxd dkn nxd'



$W^0 = Y$

$$Z = \text{Concat}(Z_0 - Z_7) \quad n \times m$$



$$\bar{x} \rightarrow \text{Enc} \rightarrow \bar{\pi} = (\pi_1, \dots, \pi_K)$$

$$\bar{z} \sim \text{Mult}(\bar{z} | \bar{\pi}) \rightarrow \text{Dec} \rightarrow \hat{x}$$

$g_i \sim \text{Gumbel}$

$$z_i = \text{softmax}_s \left(\frac{1}{\tau} (g_i + \log \pi_i) \right)$$

Gumbel-Max trick

$$z \sim \text{Mult}(\bar{\pi}) \Leftrightarrow z = \arg \max_i (g_i + \log \pi_i)$$

$$\bar{\pi} = (\pi_1, \dots, \pi_k)$$

ye $g_i \sim \text{Gumbel}$

$$p(g_i) = e^{-(g_i + e^{-g_i})}, \quad F(g_i) = e^{-e^{-g_i}}$$

Δ Δ Δ !

$$p(z=k) = p(\exists j \quad \underline{g_k + \log \pi_k \geq g_j + \log \pi_j}) =$$

$$= \int_{-\infty}^{\infty} \prod_{j \neq k} p(g_k + \log \pi_k \geq g_j + \log \pi_j) p(g_k) dg_k =$$

$$p(g_j \leq g_k + \log \pi_k - \log \pi_j) = F(g_k + \log \pi_k - \log \pi_j)$$

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

$$= e^{-\sum_{j \neq k} \pi_j} \cdot e^{-g_k - \log \pi_k}$$

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

$$e^{-g_k - \log \pi_k + \log \pi_k} \cdot e^{-g_k - \log \pi_k + \log \pi_k} =$$

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

$$= \pi_k \int e^{-g_k + \log \pi_k} \cdot e^{g_k + \log \pi_k} dg_k = \pi_k$$