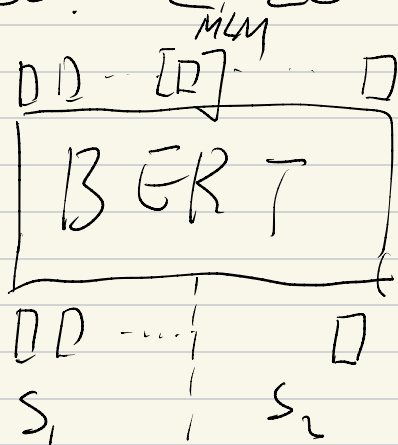
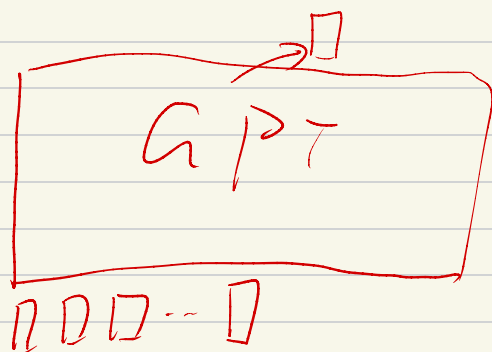


25.02.25



Encoder Only

Easy for transforming
texts into representations



Decoder Only

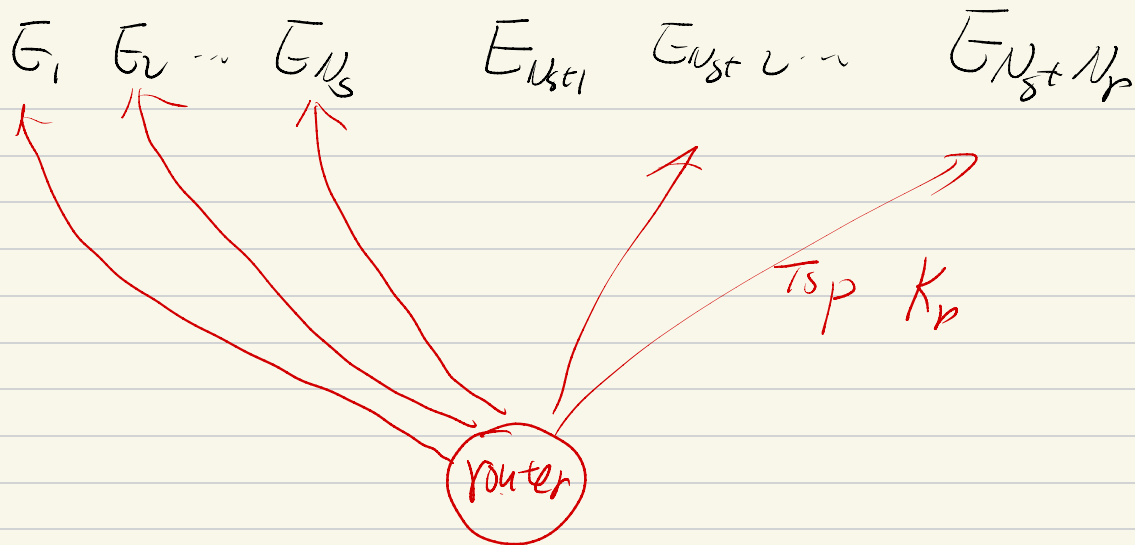
KL Divergence.

$$KL(P|Q) = \mathbb{E}_P \log \left[\frac{P}{Q} \right] = \sum_{i \in V} P_i \log \frac{P_i}{Q_i} \geq 0$$

$P, Q \in \text{Dis}(\{1, 2, \dots, |V|\})$

"=" holds iff
 $P = Q$

$$KL(P|Q) \neq KL(Q|P)$$



$$O_t = h_t + \sum_{i=1}^{N_s} f_i(h_t) + \sum_{i=N_{st}+1}^{N_{st}+N_r} g_i f_i(h_t)$$

↓

MLP for expert j

$$g_j = \frac{g'_{j,t}}{\sum_{j=N_s}^{N_{st}+N_r} g'_{j,t}}$$

$$g'_{j,t} = \begin{cases} s_{j,t} & \text{Top } k \\ 0 & \text{o.w.} \end{cases}$$

$$s_{j,t} = \text{softmax}(h_t^T \cdot e_j) \quad e_j: \text{Trainable}$$

$$g'_{j,t} = \begin{cases} s_{j,t} & \text{if } s_{j,t} + b_j \text{ is Top } k \\ 0 & \text{o.w.} \end{cases}$$