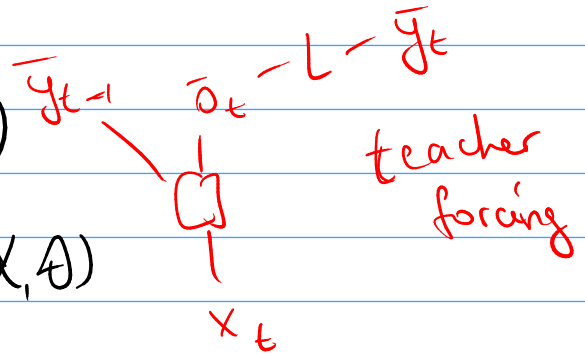


$$p(\bar{y}_1, \dots, \bar{y}_T | \bar{x}_1, \dots, \bar{x}_T, \bar{\theta}) =$$

Jordan networks

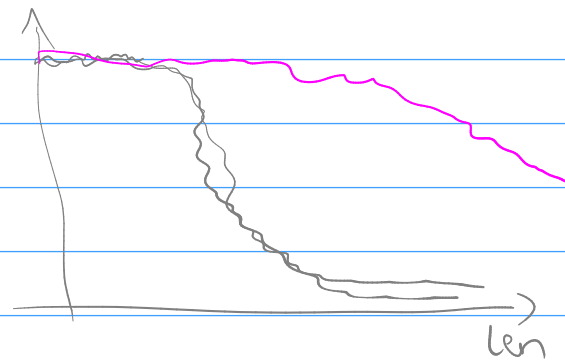
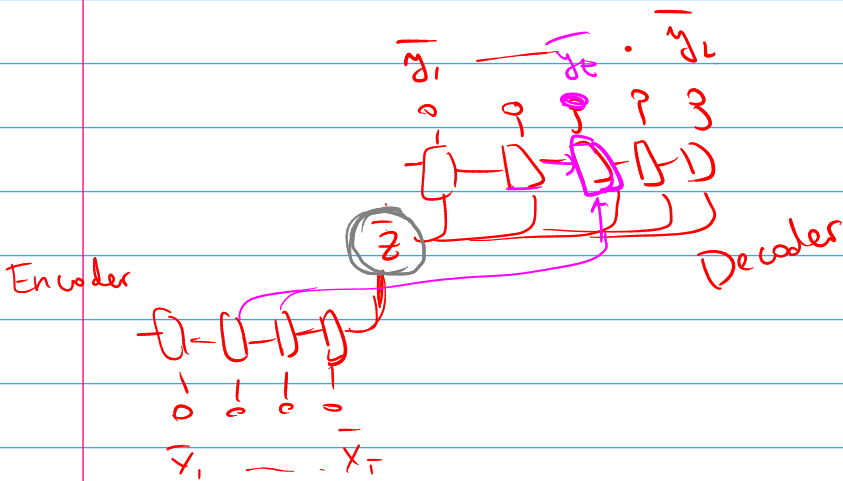
$$= p(\bar{y}_1 | \bar{x}_1, \bar{x}_1, \bar{\theta}) p(\bar{y}_2 | \bar{y}_1, \bar{x}_2, \bar{\theta}) p(\bar{y}_3 | \bar{y}_1, \bar{y}_2, \bar{x}_3, \bar{\theta}) \dots p(\bar{y}_T | \bar{y}_1, \dots, \bar{y}_{T-1}, \bar{x}_T, \bar{\theta})$$

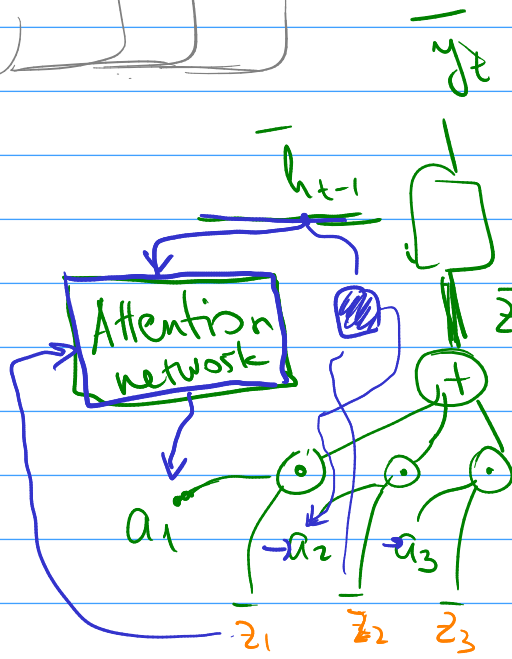
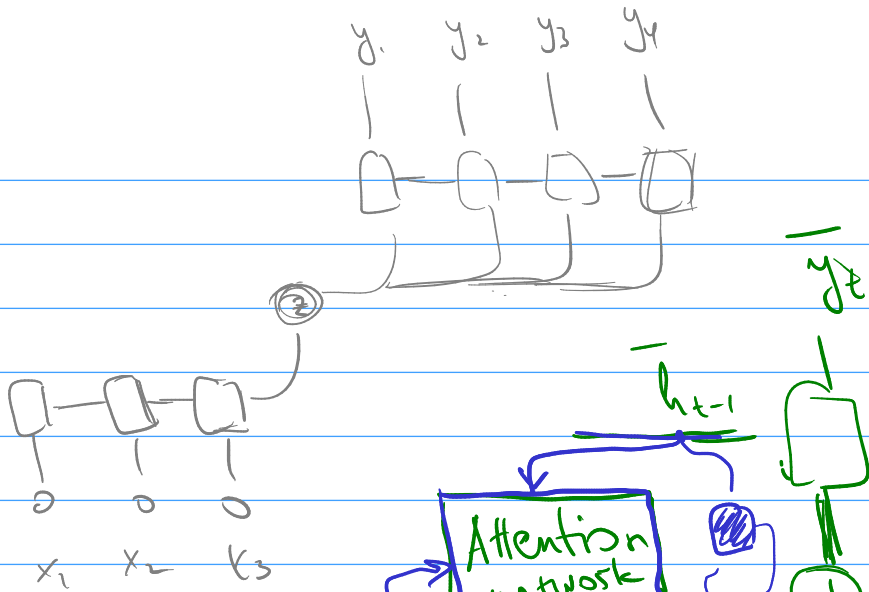


$$p(\bar{y}_1, \bar{y}_T | X, \theta) = p(\bar{y}_1 | \bar{x}_1, \bar{h}_0, \bar{\theta}) \cdot p(\bar{h}_1 | \bar{x}_1, \bar{h}_0, \bar{\theta}) \cdot$$

$$\dots p(\bar{y}_2 | \bar{x}_2, \bar{h}_2, \bar{\theta}) p(\bar{h}_2 | \bar{x}_2, \bar{h}_1, \bar{\theta}) \dots$$

$$p(\bar{y}_1, \bar{y}_T | X, \theta) = p(\bar{y}_1 | \bar{x}_1, \bar{\theta}) p(\bar{y}_2 | \bar{x}_1, \bar{x}_2, \bar{y}_1, \bar{\theta}) \dots$$





$$\bar{z} = \sum_i a_i \bar{z}_i$$

$$\hat{a}_i = f_a(\bar{z}_i, \bar{h}_{t-1})$$

$$\bar{a}_i = \text{softmax}(\hat{a}_i)$$

