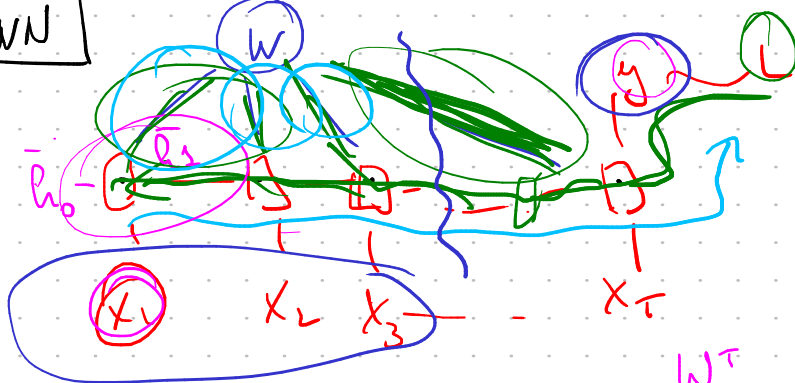
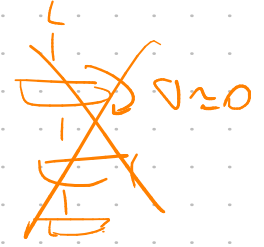


RNN



$$\frac{\partial L}{\partial w_i}$$



$$h_T = W \cdot W \cdot W \cdot W \cdot \dots \cdot W^T h_0$$

- 1) Exploding gradients
- 2) Vanishing gradients

$$\|W\| > 1$$

$$\|W\| < 1$$

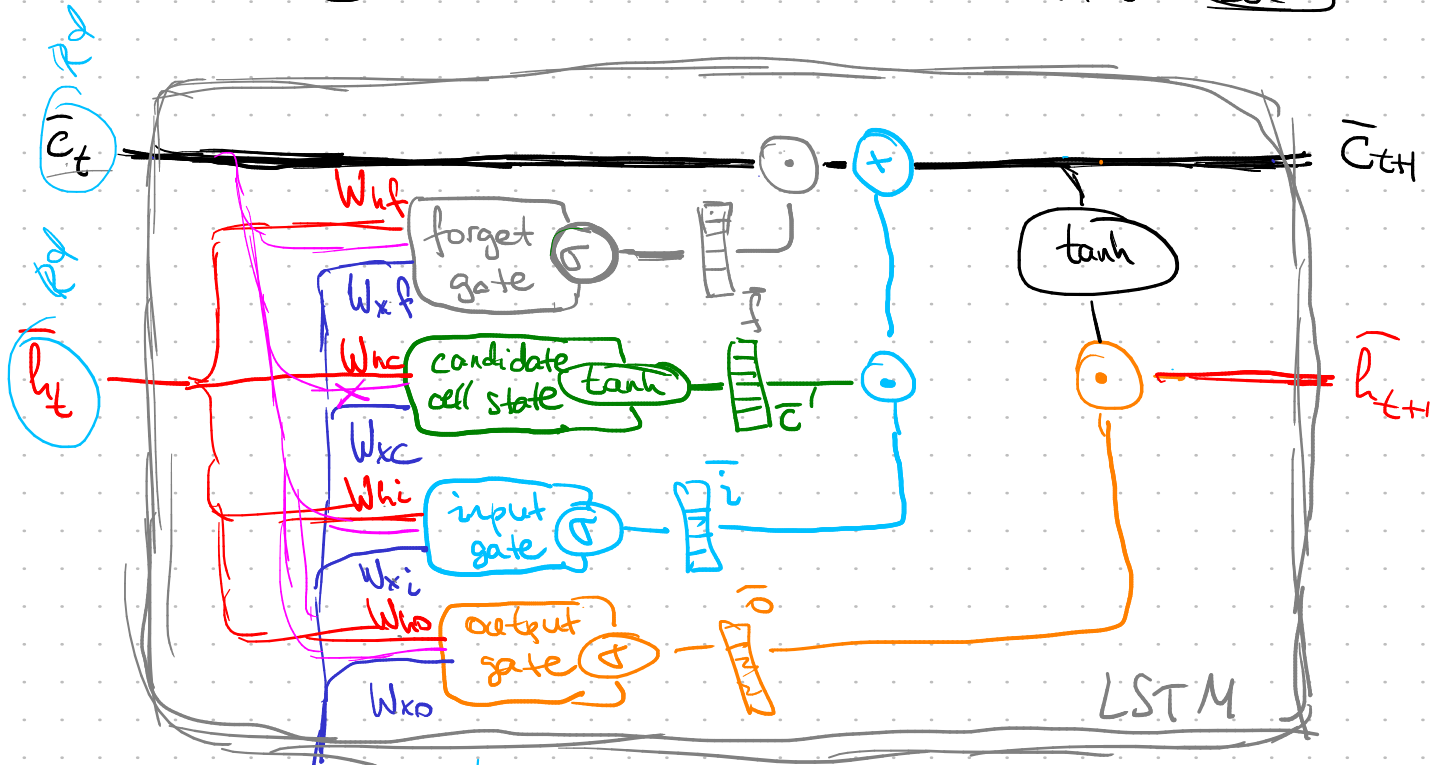
$$\|x\| < \|Wx\|$$

$$\|Wx\| \approx \alpha \|x\|$$

$\alpha < 1$

LSTM - long short-term memory

Hochreiter / Schmidhuber
1995 2000



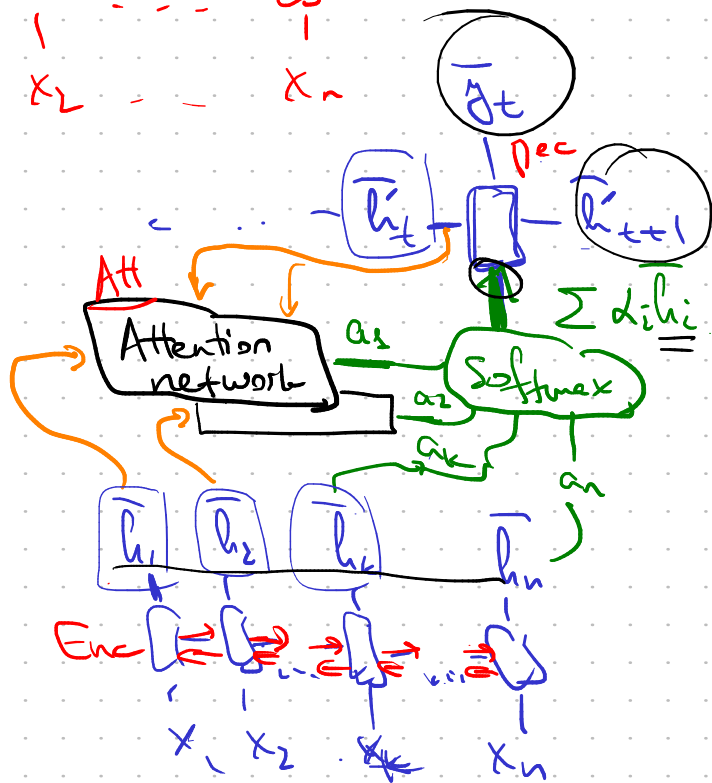
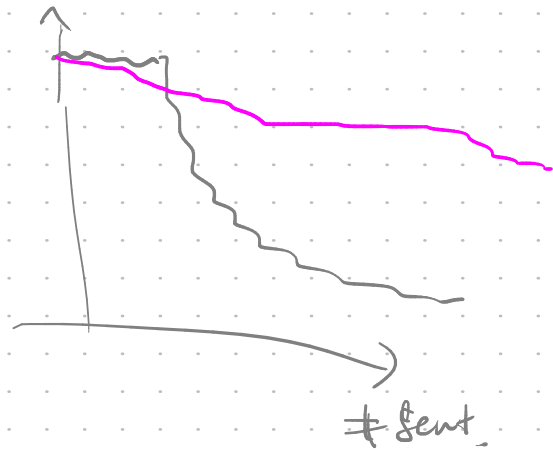
$$W_{xx} : d \times d$$

$$W_{xx'} : d \times d'$$

gate: $\frac{A}{B} = \frac{w \times d}{R \times d'}$

$$f(x, y) = f(Ax + By + c)$$

Encoders



$$d_i > 0, \sum d_i = 1$$

$$h_t = \text{Enc}(x_t, \bar{s}_t, \bar{s}_t')$$

$$h_k = \text{Enc}(x_k, \bar{s}_k, \bar{s}_k')$$

$$a_k = \text{Att}(h_t, h_k)$$

$$\bar{a} = \text{Att}(H, h_t)$$

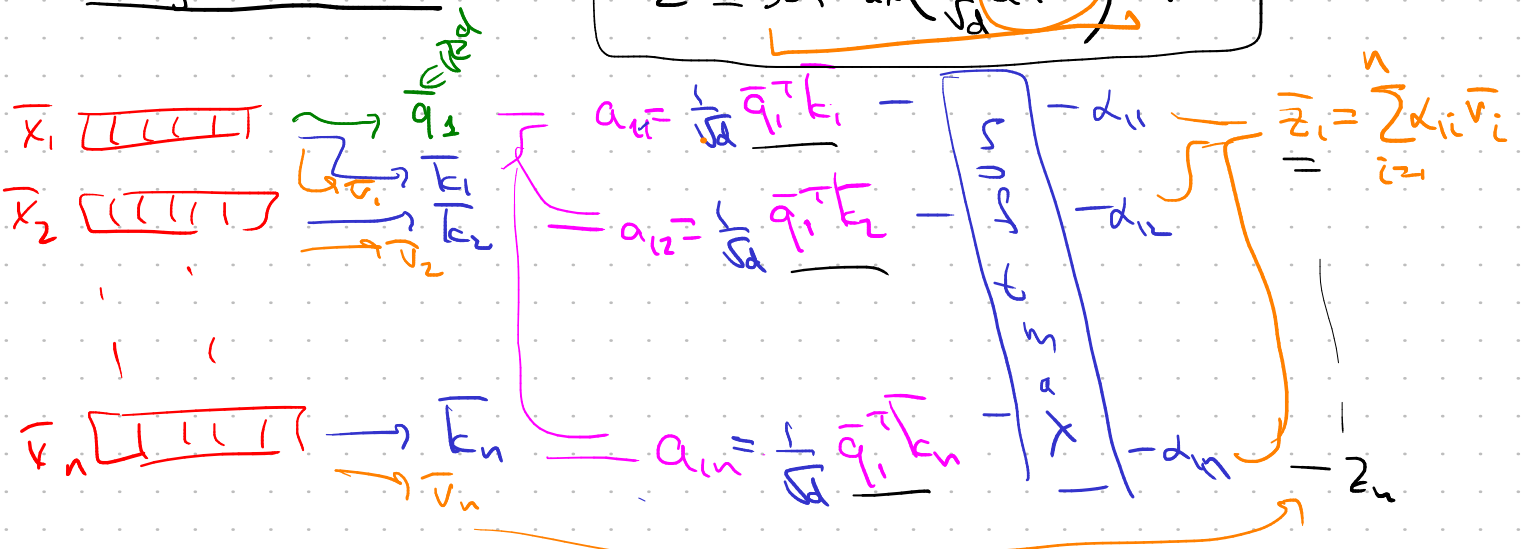
$$\bar{d} = \text{softmax}(\bar{a})$$

$$\bar{y}_t = \text{Dec}(h_t, H \cdot \bar{d})$$

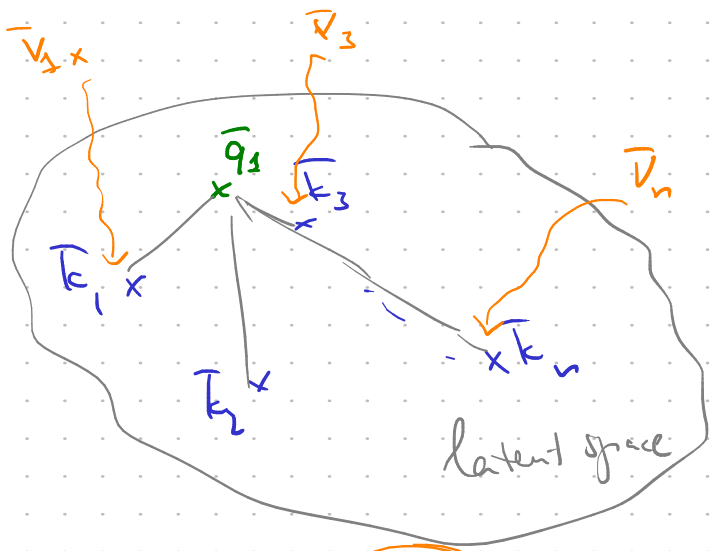
quadratic complexity

self-attention

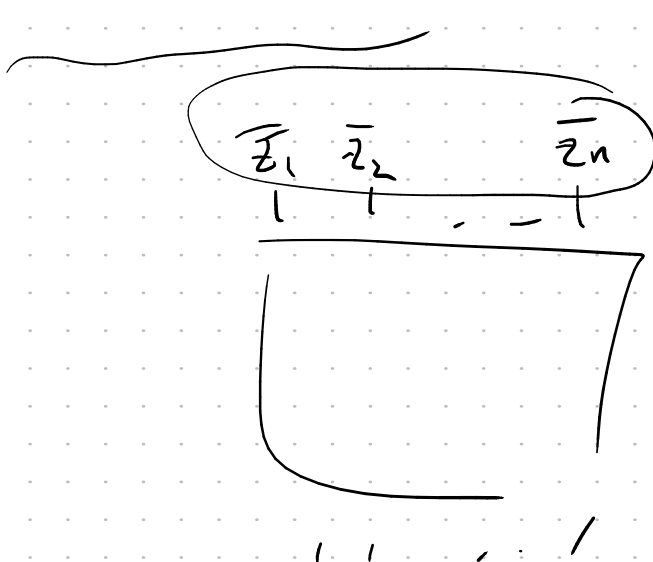
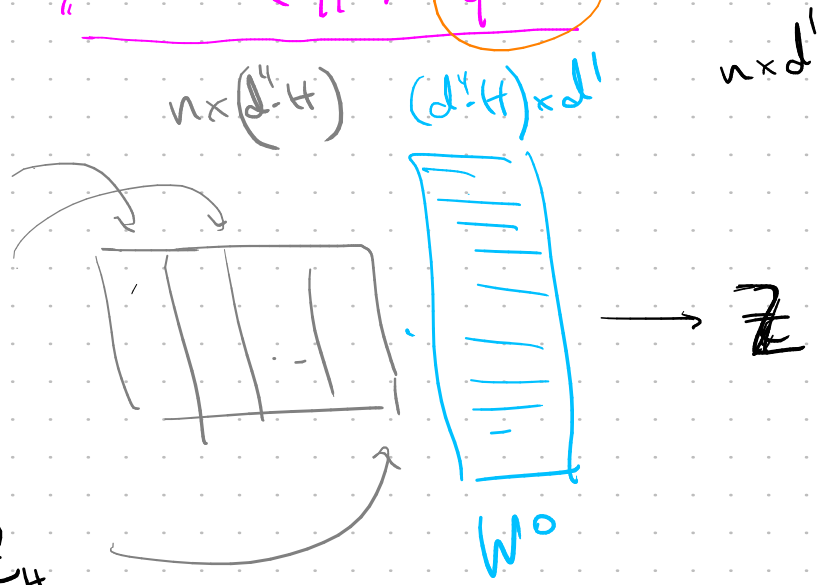
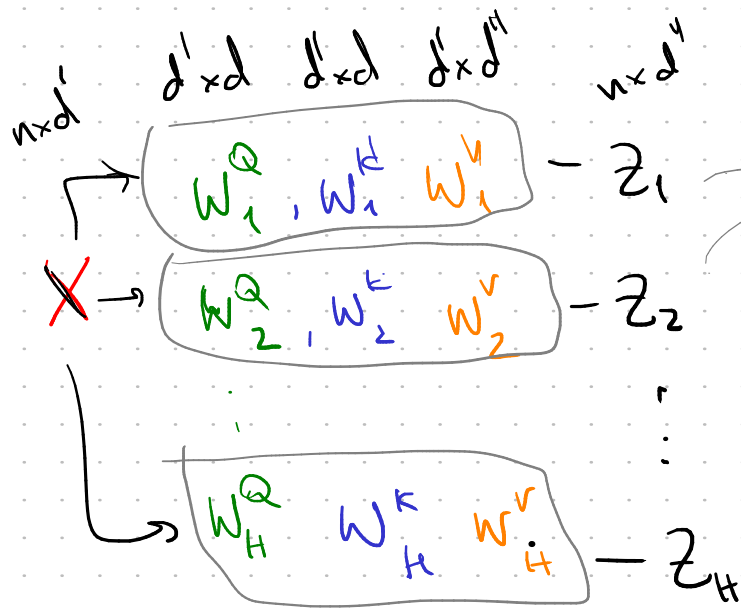
$$z = \text{softmax}\left(\frac{1}{\sqrt{d}} Q K^T\right) \cdot V$$



$\bar{x} = [1 \ 1 \ 1 \ 1 \ 1]$
 $\bar{q} = W^Q \bar{x}$
 $\bar{k} = W^K \bar{x}$
 $\bar{v} = W^V \bar{x}$



"relevance" $(\bar{q}_i, \bar{k}) = \bar{q}_i^T \bar{k}$



$k = W^K \cdot z$
 $v = W^V \cdot z$
 $Q = W^Q \cdot y$
 y_1, y_2, \dots, y_k

ROP E

