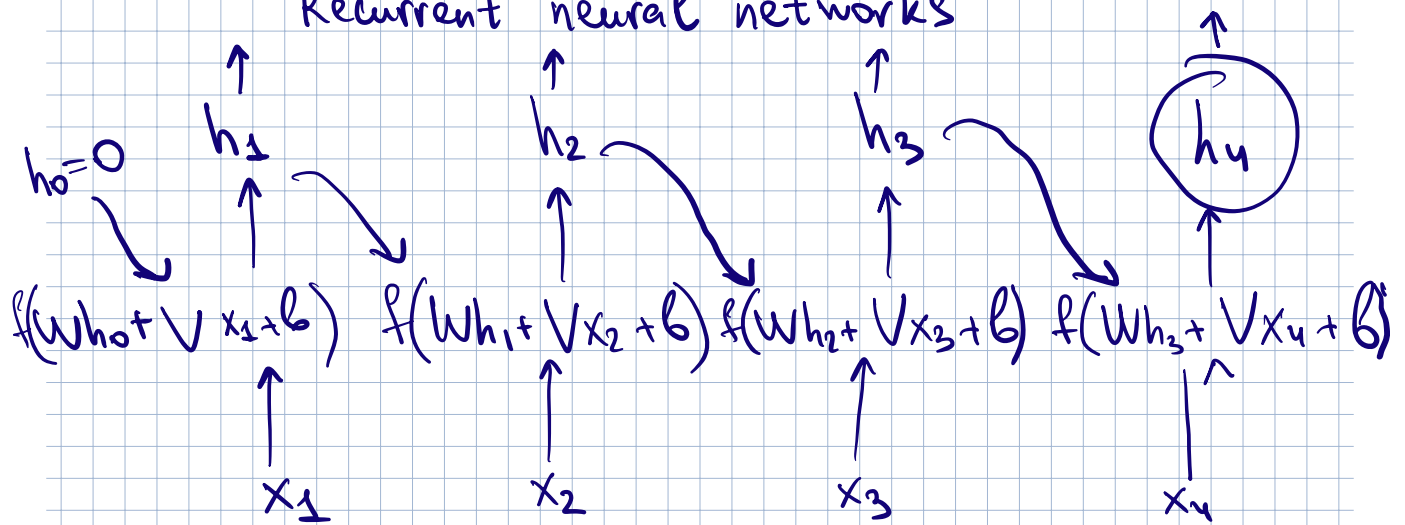


# Recurrent neural networks

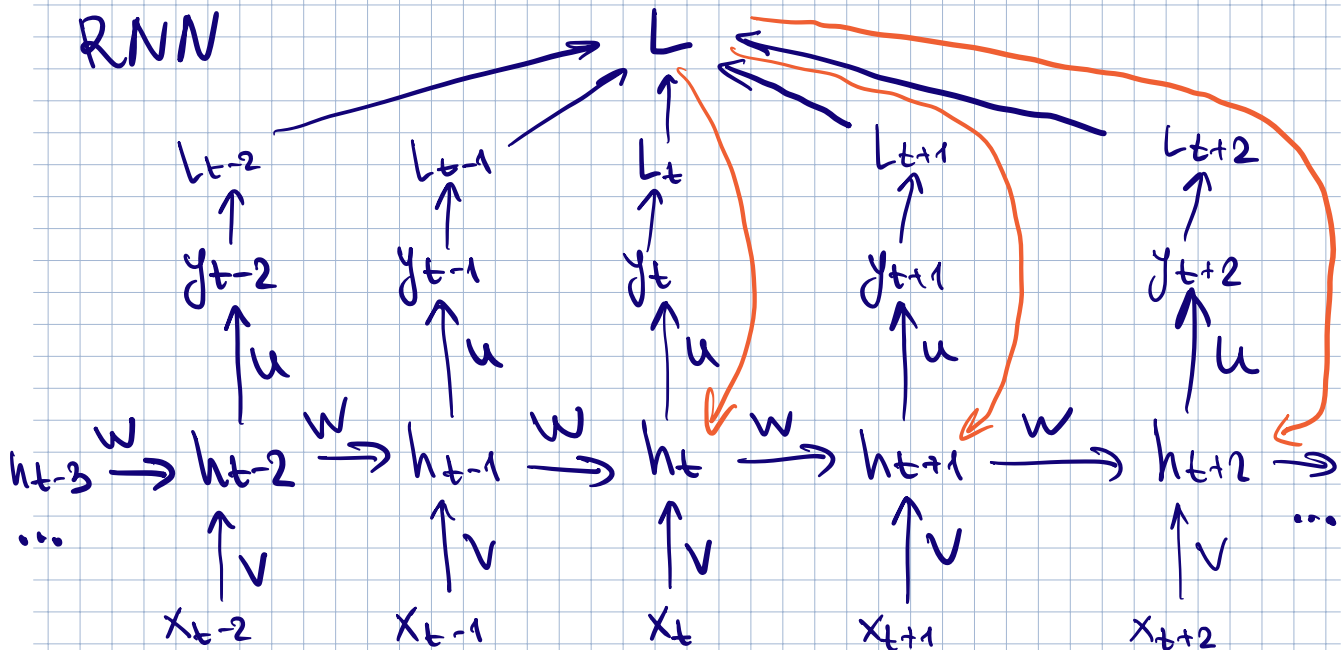


$x_i \in \mathbb{R}^d$ 

$$\begin{pmatrix} V & W \end{pmatrix} \begin{pmatrix} x_t \\ h_{t-1} \end{pmatrix}$$

$$h_t = f(W h_{t-1} + V x_t + b)$$

RNN



$$L = \frac{1}{T} \sum_{t=1}^T L_t$$

$$\frac{dL}{du} = \frac{1}{T} \sum_{t=1}^T \frac{dL_t}{du} = \frac{1}{T} \sum_{t=1}^T \frac{dL_t}{dy_t} \cdot \frac{dy_t}{du}$$

$$\bullet \frac{dL}{dW} = \frac{1}{T} \sum_{t=1}^T \frac{dh_t}{dW} = \frac{1}{T} \sum_{t=1}^T \frac{dh_t}{dy_t} \cdot \frac{dy_t}{dh_t} \cdot \frac{dh_t}{dW}$$

$$h_t = f(Vx_t + W h_{t-1} + b)$$

$$\frac{dh_t}{dW} = \frac{\partial h_t}{\partial W} + \frac{\partial h_t}{\partial h_{t-1}} \cdot \frac{dh_{t-1}}{dW} = \frac{\partial h_t}{\partial W} + \frac{\partial h_t}{\partial h_{t-1}} \cdot$$

$$\bullet \left( \frac{\partial h_{t-1}}{\partial W} + \frac{\partial h_{t-1}}{\partial h_{t-2}} \cdot \frac{dh_{t-2}}{dW} \right) =$$

$$= \frac{\partial h_t}{\partial W} + \frac{\partial h_t}{\partial h_{t-1}} \cdot \frac{\partial h_{t-1}}{\partial W} + \frac{\partial h_t}{\partial h_{t-1}} \cdot \frac{\partial h_{t-1}}{\partial h_{t-2}} \cdot \frac{dh_{t-2}}{dW}$$

$$= \frac{\partial h_t}{\partial W} + \sum_{j=1}^{t-1} \left( \prod_{k=t-j}^{t-1} \frac{\partial h_{k+1}}{\partial h_k} \right) \cdot \frac{\partial h_{t-j}}{\partial W} =$$

$$= \sum_{j=0}^{t-1} \left( \prod_{k=t-j}^{t-1} \frac{\partial h_{k+1}}{\partial h_k} \right) \cdot \frac{\partial h_{t-j}}{\partial W}$$

$$\frac{dL}{dW} = \frac{1}{T} \sum_{t=1}^T \frac{dL_t}{dy_t} \cdot \frac{dy_t}{dh_t} \left[ \sum_{j=0}^{t-1} \left( \prod_{k=t-j}^{t-1} \frac{\partial h_{k+1}}{\partial h_k} \right) \frac{\partial h_{t-j}}{\partial W} \right]$$

BPTT (Backpropagation Through Time)

$$\bullet \left\| \frac{\partial h_{k+1}}{\partial h_k} \right\| > 1 \text{ - exploding gradients divergence, NaNs}$$

Gradient clipping

$$g \leftarrow \min \left( 1, \frac{\text{thr}}{\|grad\|} \right) \cdot grad$$

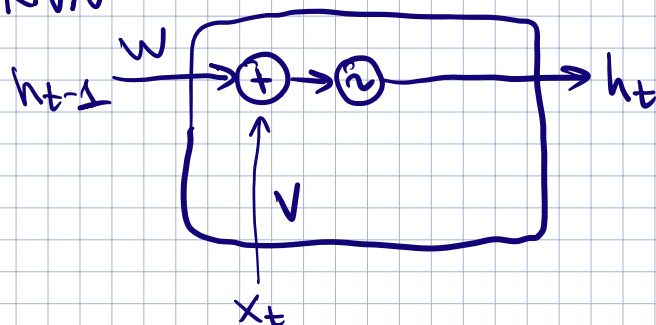
$$g \leftarrow \text{clip}(\text{grad}, -c, c)$$

$$\bullet \left\| \frac{\partial h_{k+1}}{\partial h_k} \right\| < 1 \text{ - vanishing gradients}$$

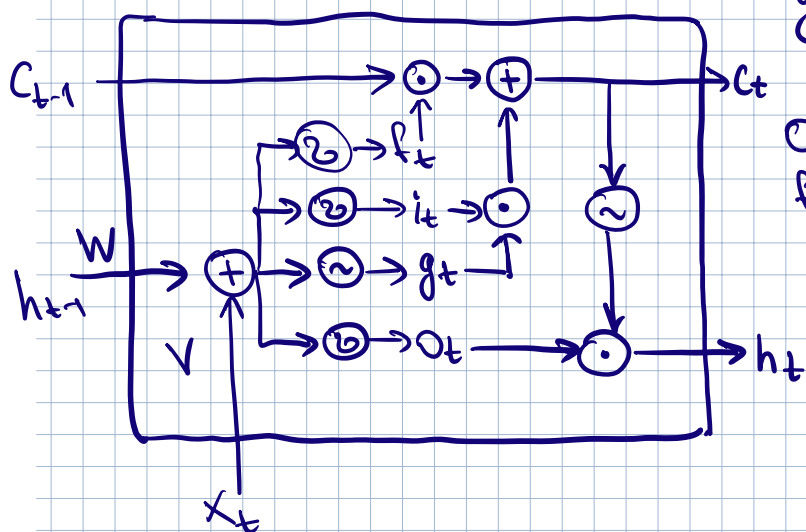
$$h_{k+1} = f(\underbrace{Wh_k + Vx_{k+1} + b}_{z_k})$$

$$\left\| \frac{\partial h_{k+1}}{\partial h_k} \right\| = \left\| \text{diag}(f'(z_k)) \right\| \cdot \|W\| \rightarrow \text{init with orthogonal matrix}$$

RNN

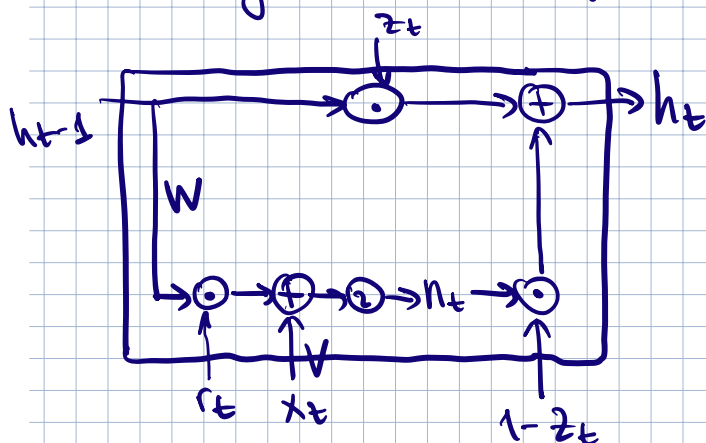


LSTM (long short-term memory)



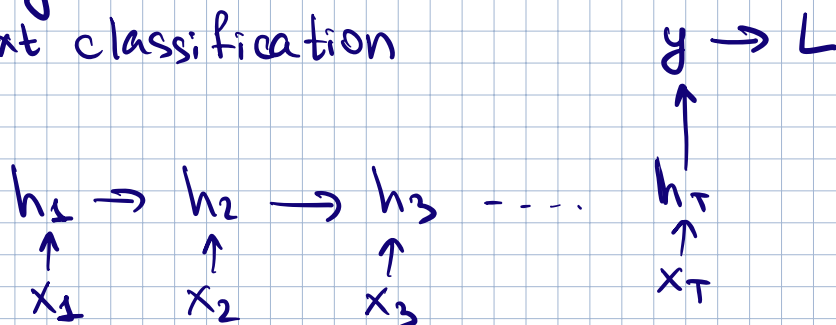
$$\begin{aligned} g_t &= f(W_g h_{t-1} + V_g x_t + b_g) \\ i_t &= \sigma(W_i h_{t-1} + V_i x_t + b_i) \\ o_t &= \sigma(W_o h_{t-1} + V_o x_t + b_o) \\ f_t &= \sigma(W_f h_{t-1} + V_f x_t + b_f) \\ C_t &= C_{t-1} \odot f_t + i_t \odot g_t \\ h_t &= f(C_t) \odot o_t \end{aligned}$$

## GRU (gated recurrent unit)



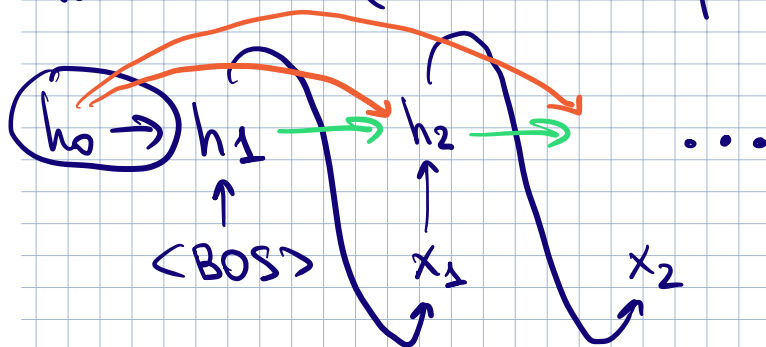
$$\begin{aligned} r_t &= \sigma(V_r x_t + W_r h_{t-1} + b_r) \\ z_t &= \sigma(V_z x_t + W_z h_{t-1} + b_z) \\ n_t &= f(V_n x_t + r_t \odot (W_n h_{t-1} + b_n)) \\ h_t &= (1 - z_t) \odot n_t + z_t \odot h_{t-1} \end{aligned}$$

- Many-to-one
- Text classification



- One-to-many
- Conditional generation (image captioning, genre music generation)

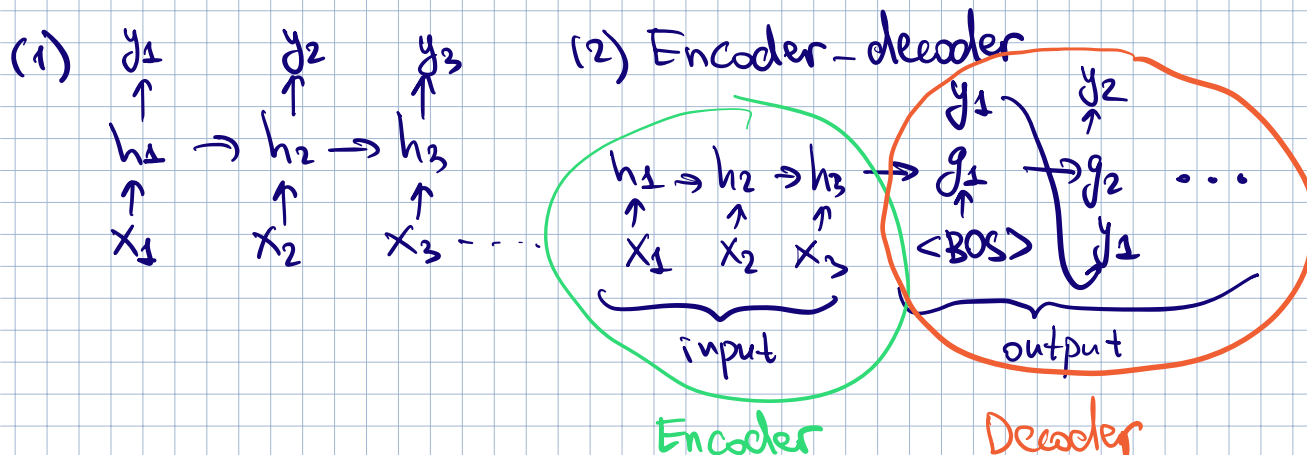
$h_0 = \text{Linear}(\text{Encoder (input)})$



- Many-to-many (seq-to-seq)

- Part of speech tagging (1)

- Machine translation (2)



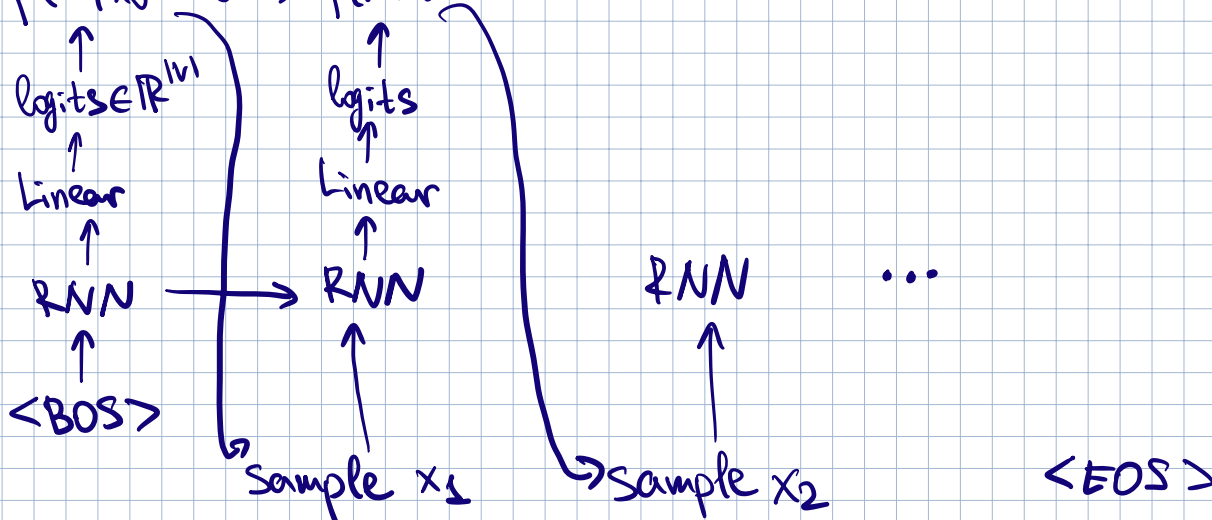
## Language model

Autoregressive generation

$$x = (x_0, x_1, x_2, \dots, x_T)$$

$$p(x) = p(x_0) p(x_1 | x_0) \cdot p(x_2 | x_0, x_1) \cdot \dots \cdot p(x_T | x_0, \dots, x_{T-1})$$

$$p(x | x_0 = \langle \text{BOS} \rangle) \quad p(x | x_0 = \langle \text{BOS} \rangle, x_1 = x_1)$$



$\Theta$  - model parameters

$$P_{\theta}(x_t | x_0, x_1, \dots, x_{t-1})$$

$X = (x_0, \dots, x_T)$  - real sequence from dataset

$$P_{\theta}(x) = \underbrace{p(x_0)}_{=1} \cdot P_{\theta}(x_1 | x_0) \cdot \dots \cdot P(x_T | x_0, \dots, x_{T-1})$$

$$= \sum_{t=1}^T \log P_{\theta}(x_t | x_0, \dots, x_{t-1}) \rightarrow \min_{\theta}$$