

## LSTM - NETWORKS

→ Solve the vanishing gradient issue of RNN.

→ Thus better suited for long-term dependencies.

→ Weight matrices:

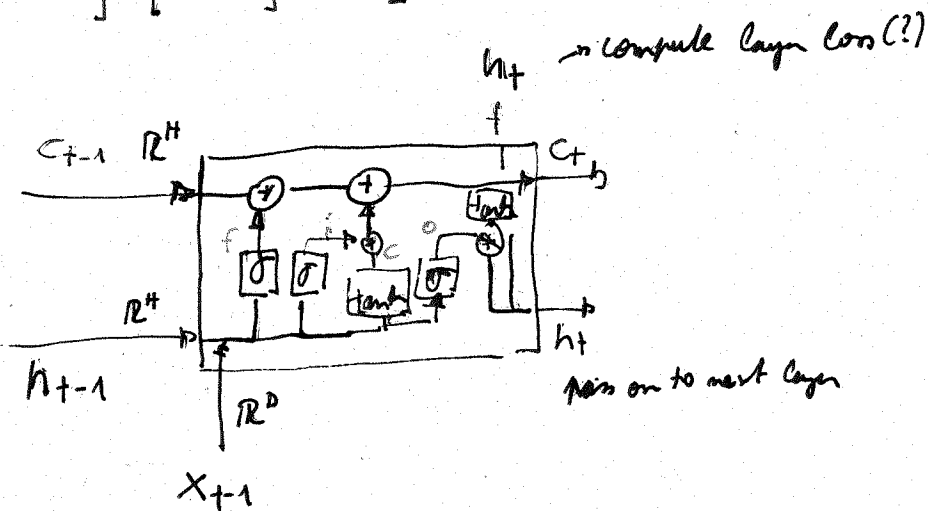
- $W_f$  - forget gate
- $W_i$  - input gate
- $W_c$  - cell state
- $W_o$  - output gate

This results in a  
input-to-hidden matrix  
and a hidden-to-hidden  
matrix.

$$W_x \in \mathbb{R}^{4H \times D}$$

$$W_h \in \mathbb{R}^{4H \times H}$$

$$\begin{matrix} (4H, D) & (N, D) & (N, D) & (D, 4H) & (N, 4H) \\ [W_x] & \cdot [X_t] & \rightarrow [X_t] \cdot [W_x^T] = [A] \end{matrix}$$



## LSTM Backwards

I get some loss at every layer and <sup>gradient</sup> from a previous layer (through the upstream cell state loss and hidden loss). I need to compute the <sup>gradient</sup> loss with regards ~~to~~ to

$C_{t-1}$ ,  $h_{t-1}$ ,  $W_x$ ,  $W_h$ .  $W_h$  is ~~used~~ used  $t$ -times and thus gradients add up.  $W_x$  is used only the first time thus it is computed only at the end.

