

HDDL - LESSON 4

RECURRENT NEURAL NETWORKS
LSTMs AND GRUs

JOSEBA DALMAU

SEQUENTIAL DATA

The boss said the employee is late.

The employee said the boss is late.



Same words , different meanings

⇒ The order of the words in
the sentence matters !

SEQUENTIAL DATA

- Text
- Audio : speech, music
- Video
- Time series
- Biological data : DNA , RNA

TASKS FOR SEQUENTIAL DATA

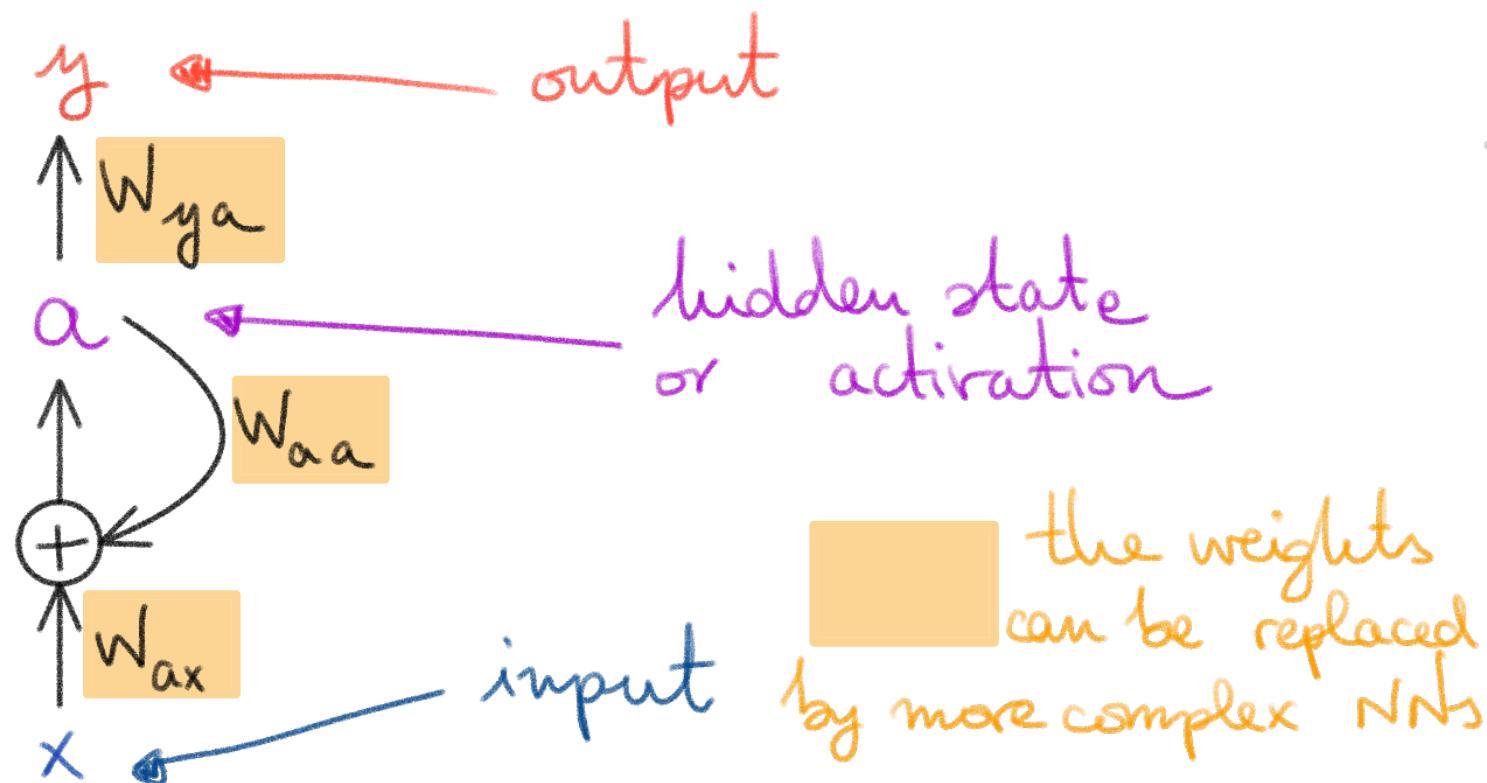
- Prediction → next token prediction in text, audio
→ forecasting for time series
- Classification → sentiment analysis
→ DNA sequence classification
- Generation : text, audio, music, video . . .
- Sequence - to - sequence mapping (seq2seq)
 - machine translation
 - text summarization
 - speech - to - text conversion

NAIVE SOLUTIONS

- Use common feed-forward networks
 - separate parameters for each input
 - have to learn language rules for all positions
 - fixed size input/output
- Use a 1-dimensional convolution
 - parameters are shared across inputs
 - output is a fcn of a small nb of neighbors of the input
 - fixed size input/output

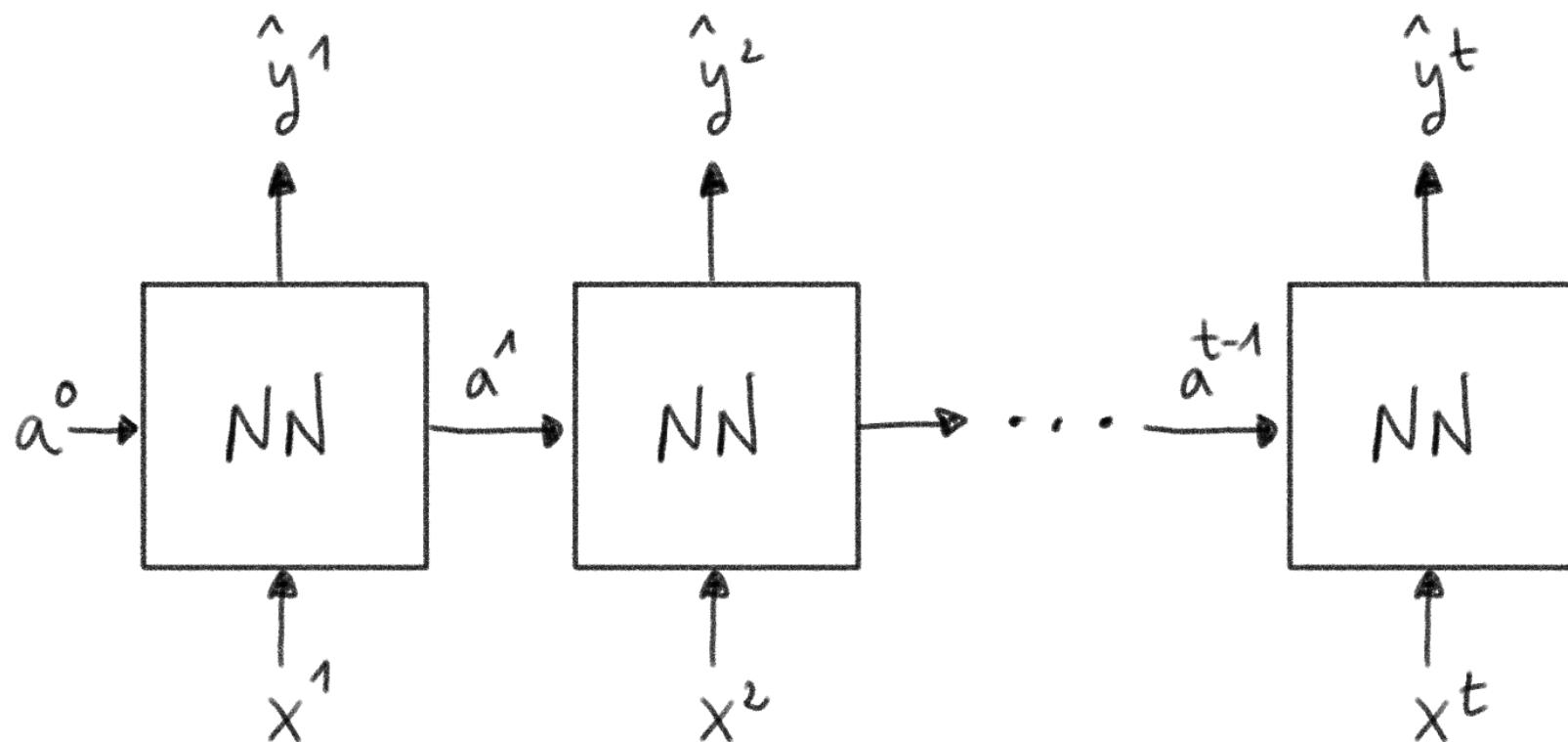
RNNs: BASIC IDEA

Process input sequentially, Keep activations for next input:



RNNs: "UNFOLDING"

Process input sequentially, Keep activations
for next input:



SIMPLEST FORMULAS

$$a^t = g_1 (W_{aa}^{t-1} a^{t-1} + W_{ax} x^t + b_a)$$

activation
functions

weight
matrices

bias
vectors

$$\hat{y}^t = g_2 (W_{ya} a^t + b_y)$$



The coefficients in the RED COLORED
matrices and vectors are the same
for all t !

PROs AND CONs

Can process input of any length

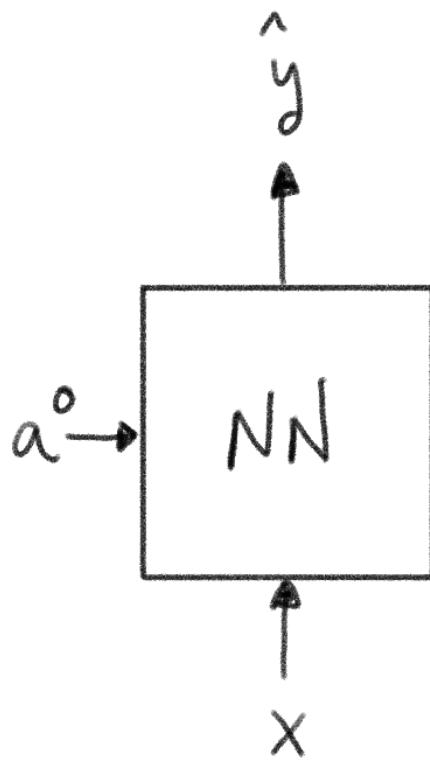
Model size indep. of input length

Slow computation (why?)

Bad with long range dependencies

Cannot consider future input

RNN TYPES: one-to-one

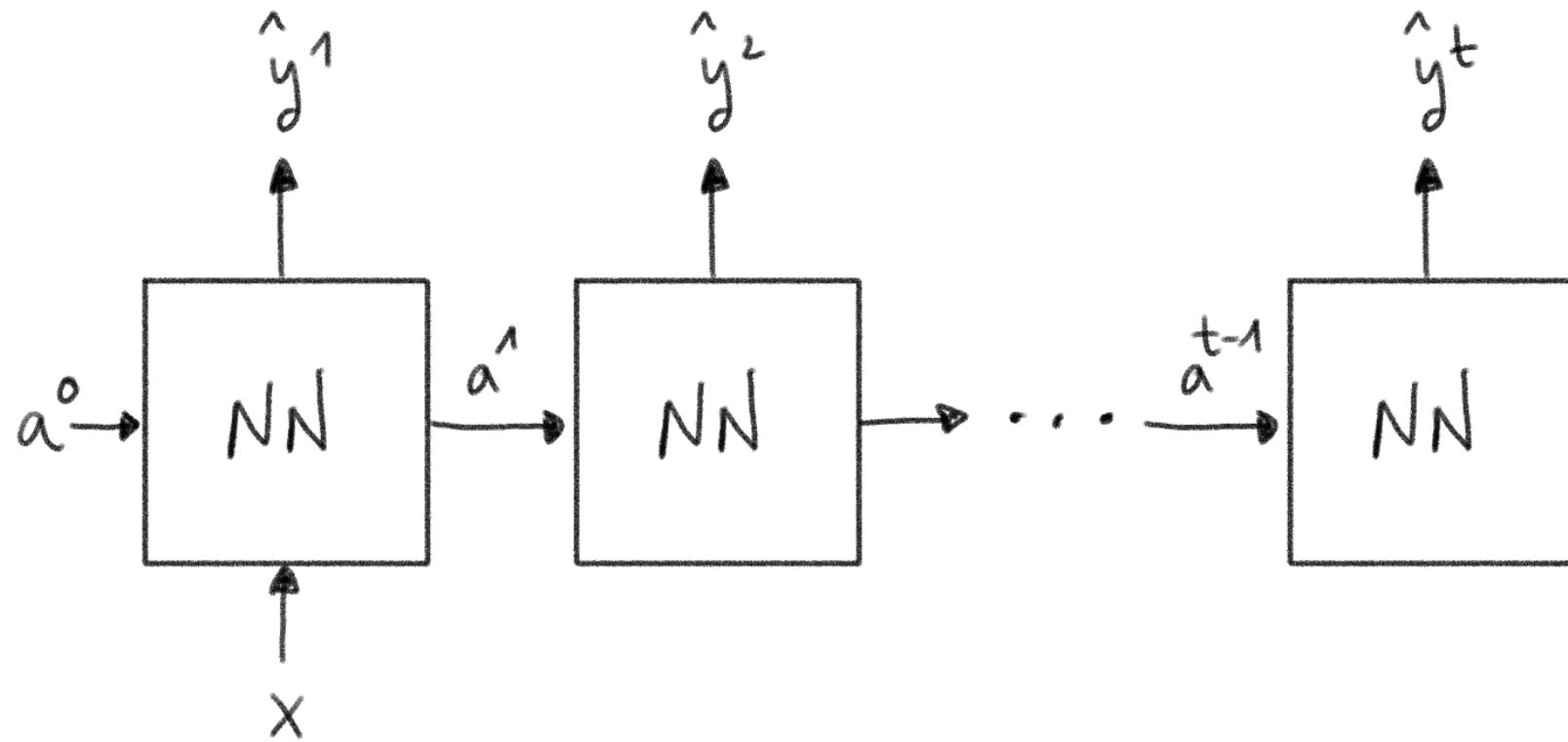


$$T_x = T_y = 1$$

$\uparrow \quad \uparrow$
sequence lengths

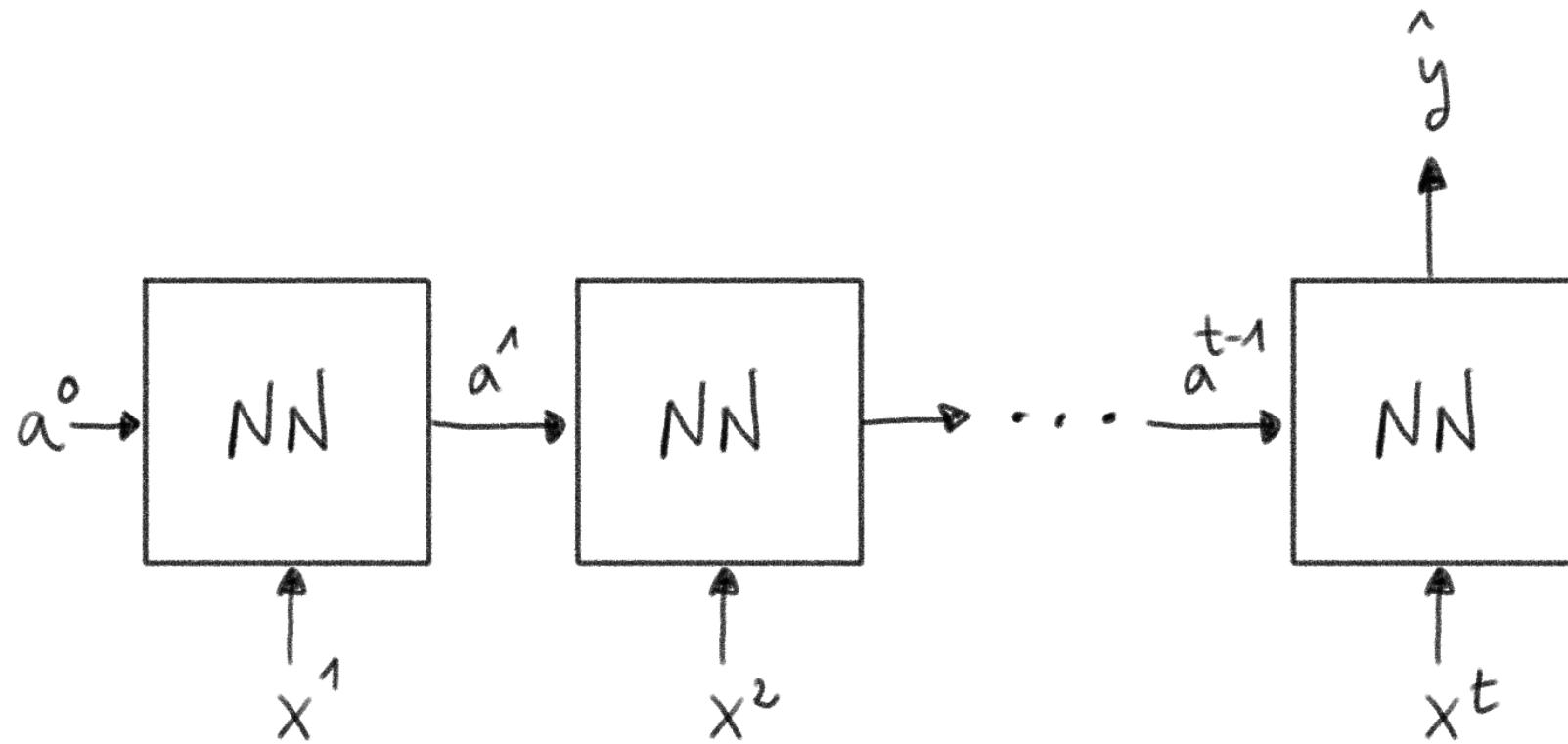
e.g. Fully Connected
or Convolutional NNs

RNN TYPES: one-to-many



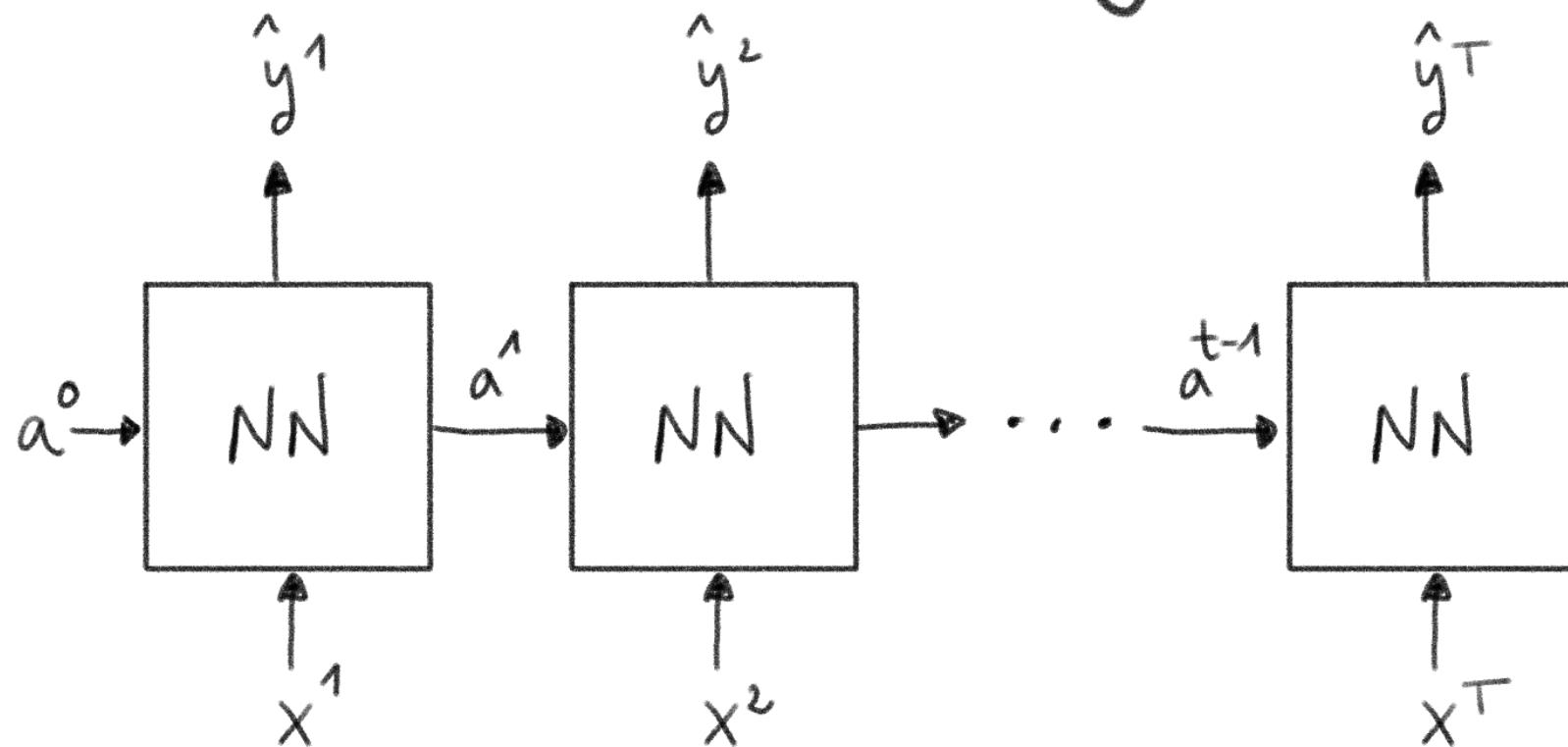
$T_x = 1, T_y > 1$ e.g. image captioning

RNN TYPES: many-to-one



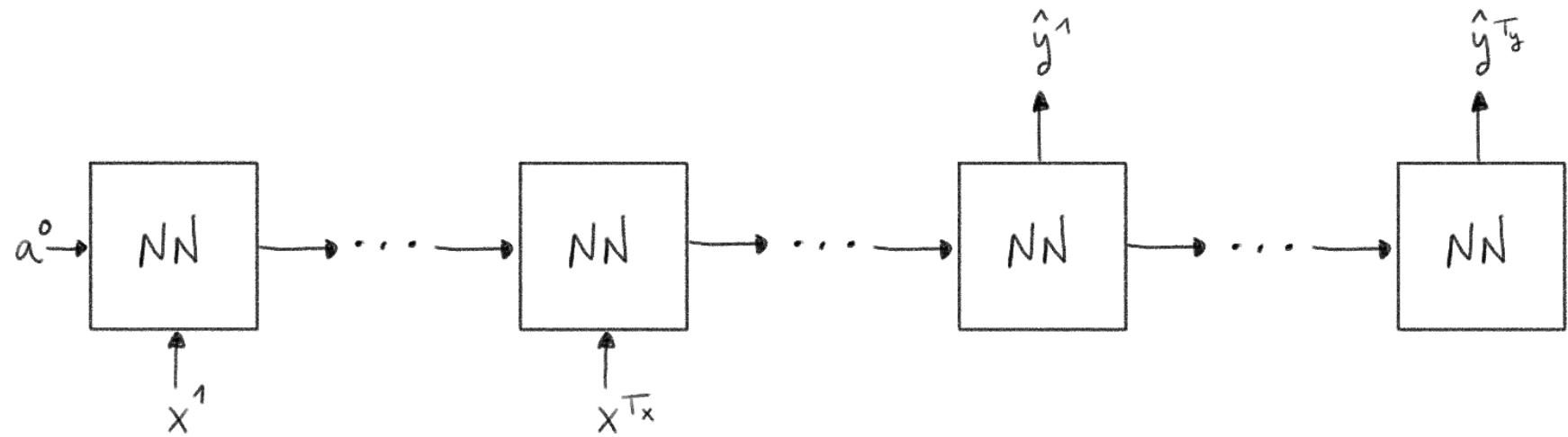
$T_x > 1, T_y = 1$ e.g. sentiment analysis

RNN TYPES: many-to-many
(same length)



$T_x = T_y > 1$ e.g. name entity recognition

RNN TYPES: many-to-many
(different length)



$T_x \neq T_y > 1$ e.g. machine translation

TRAINING

At each time step t :

$$a^t = \tanh(W_{aa}a^{t-1} + W_{ax}x^t + b_a)$$

$$\hat{y}^t = \text{softmax}(W_{ya}a^t + b_y)$$

↑
over what?

TRAINING: LOSS FUNCTION

$$\mathcal{L}_{\theta}(\hat{y}, y) = \sum_{t=1}^T L_{\theta}(\hat{y}^t, y^t)$$

↑ parameters of the RNN,
i.e., W_{aa} , W_{ax} , W_{ya} , b_a , b_y

Example: Negative log-likelihood

$$L_g(\hat{y}^t, y^t) = - \log \underbrace{P_{\text{model}}(y^t | \{x^1, \dots, x^{t-1}\})}_{= \hat{y}^t(y^t)}$$

TRAINING: BACKPROPAGATION THROUGH TIME (BPTT)

$$\mathcal{L}_{\theta}(\hat{y}, y) = \sum_{t=1}^T L_{\theta}(\hat{y}^t, y^t)$$

$$\nabla_{\theta} \mathcal{L}(\hat{y}, y) = \sum_{t=1}^T \nabla_{\theta} L_{\theta}(\hat{y}^t, y^t)$$



Careful! The dependence of \hat{y}^t on θ is quite complex

→ Runtime of BPTT is $\mathcal{O}(T)$

Memory cost of BPTT is $\mathcal{O}(T)$

EXERCISE : BPTT FOR LINEAR RNN

Assume: $a^t = W_{aa} a^{t-1} + W_{ax} x^t$

$$\hat{y}^t = W_{ya} a^t$$

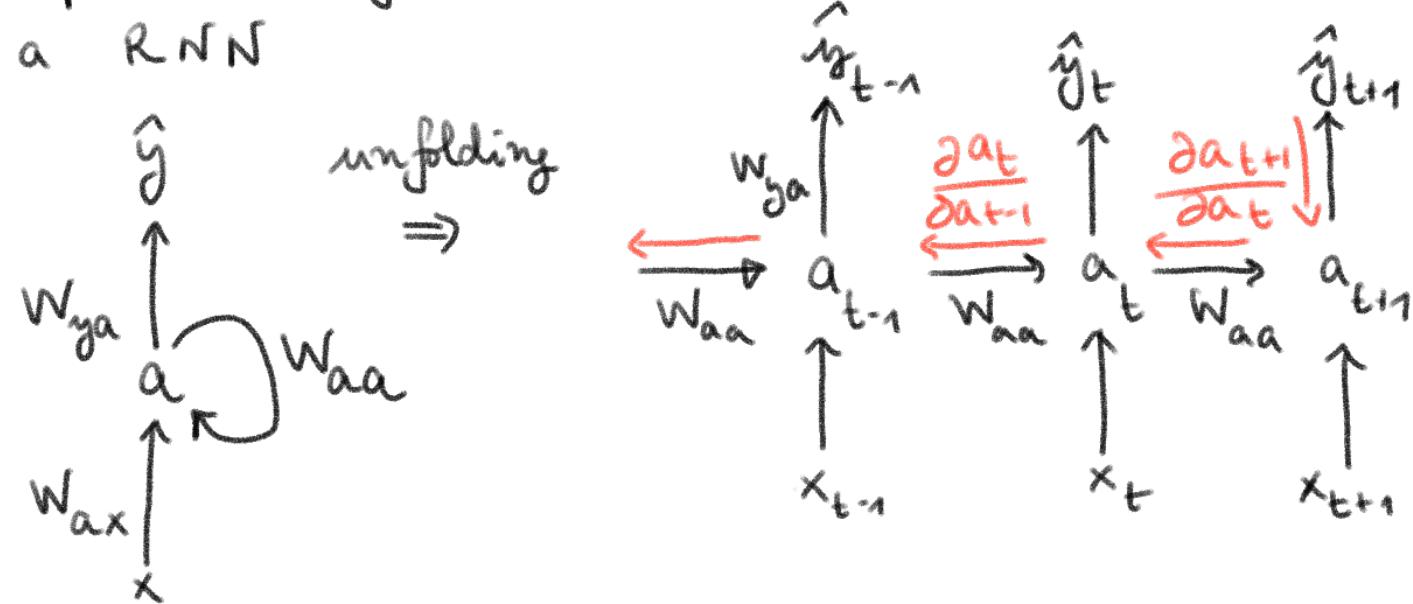
and compute: $\frac{\partial}{\partial a^t} L(\hat{y}^t, y^t)$ for the MSE loss, for $t \leq T$: i.e.

$$L(\hat{y}^t, y^t) = \frac{1}{2} \| \hat{y}^t - y^t \|^2$$

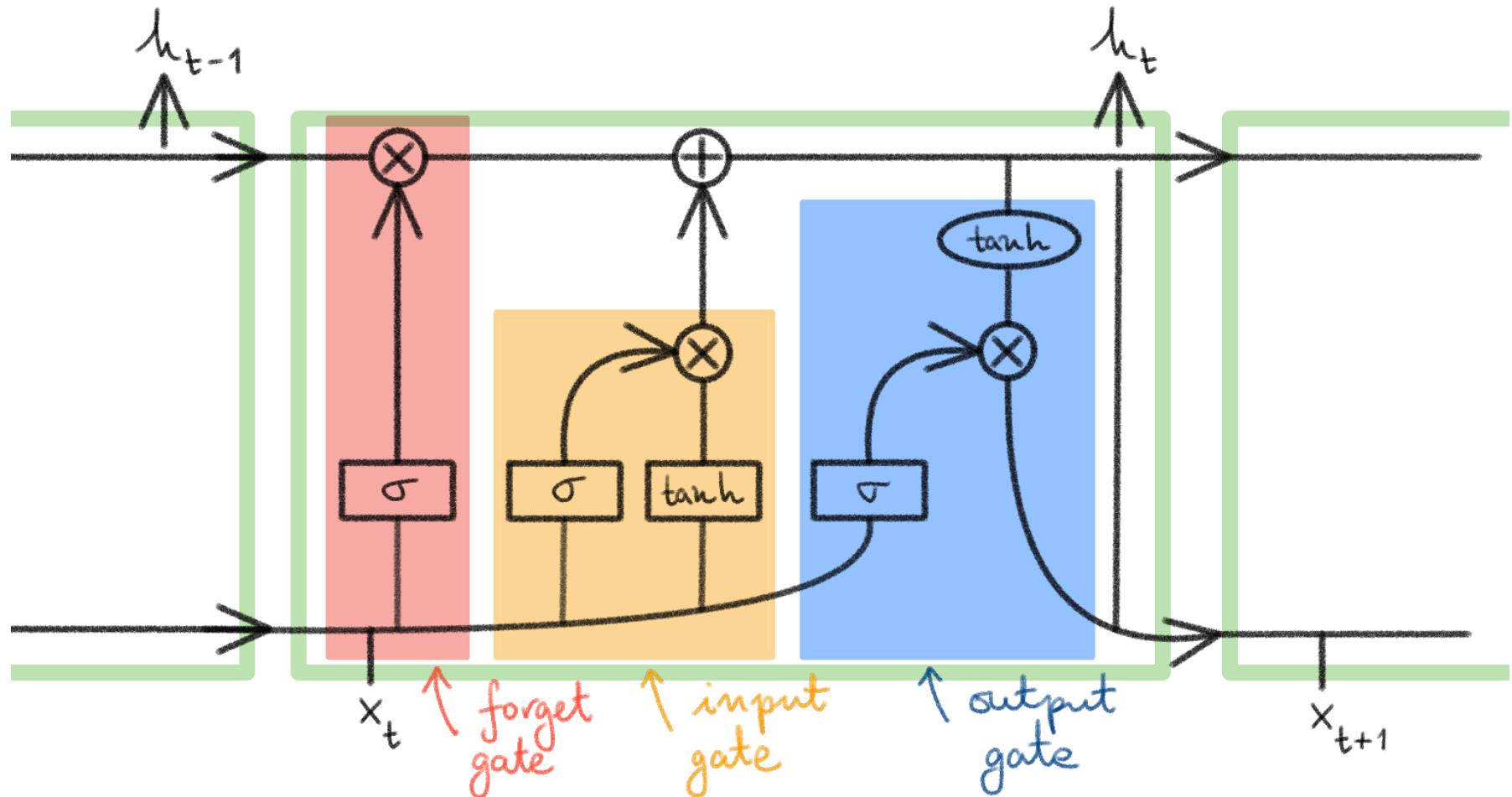
PROBLEMS WITH BPPT

- Exploding gradient \Rightarrow very bad for learning
- Vanishing gradient \Rightarrow very short memory

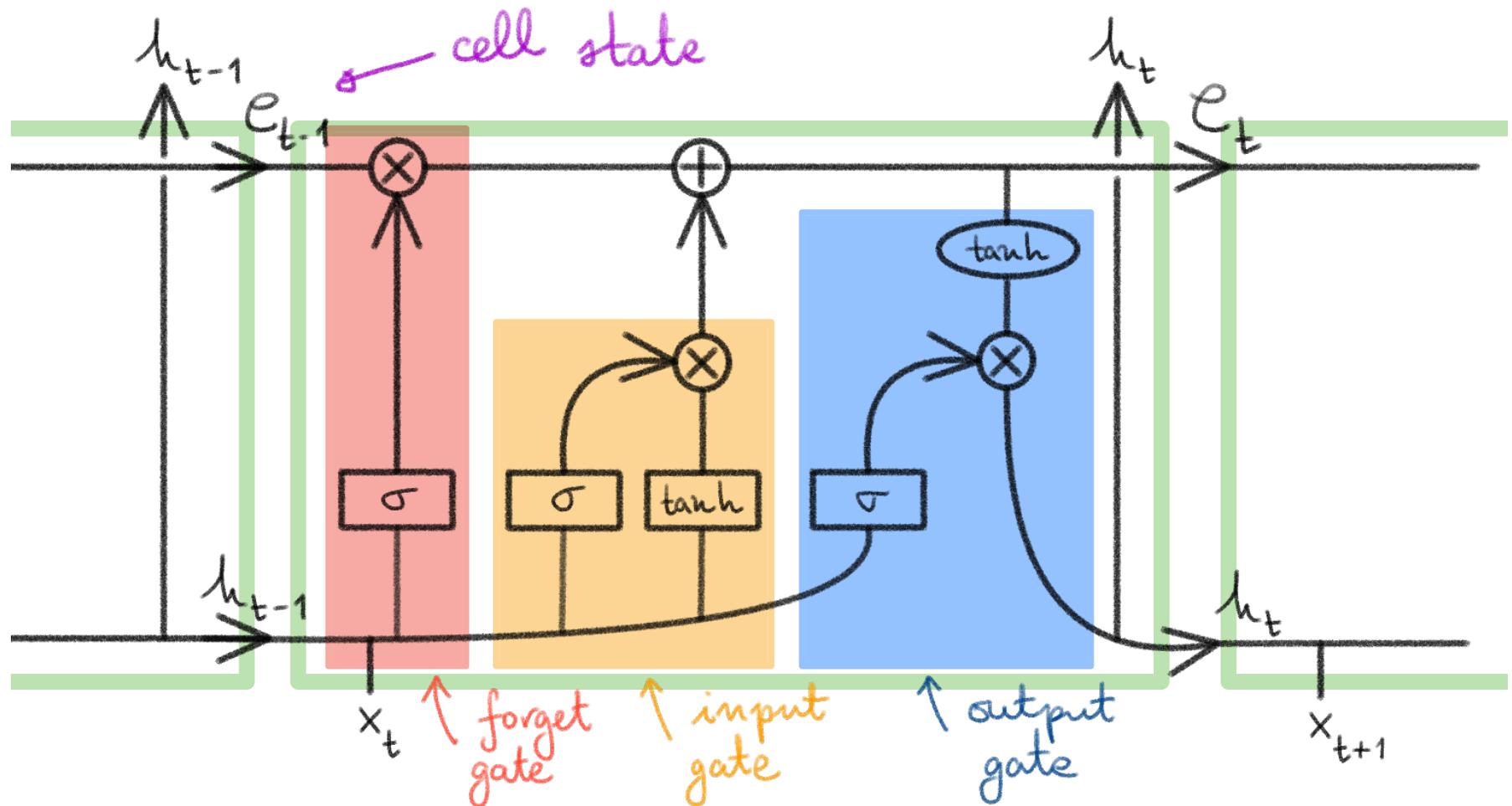
Graph view of
a RNN



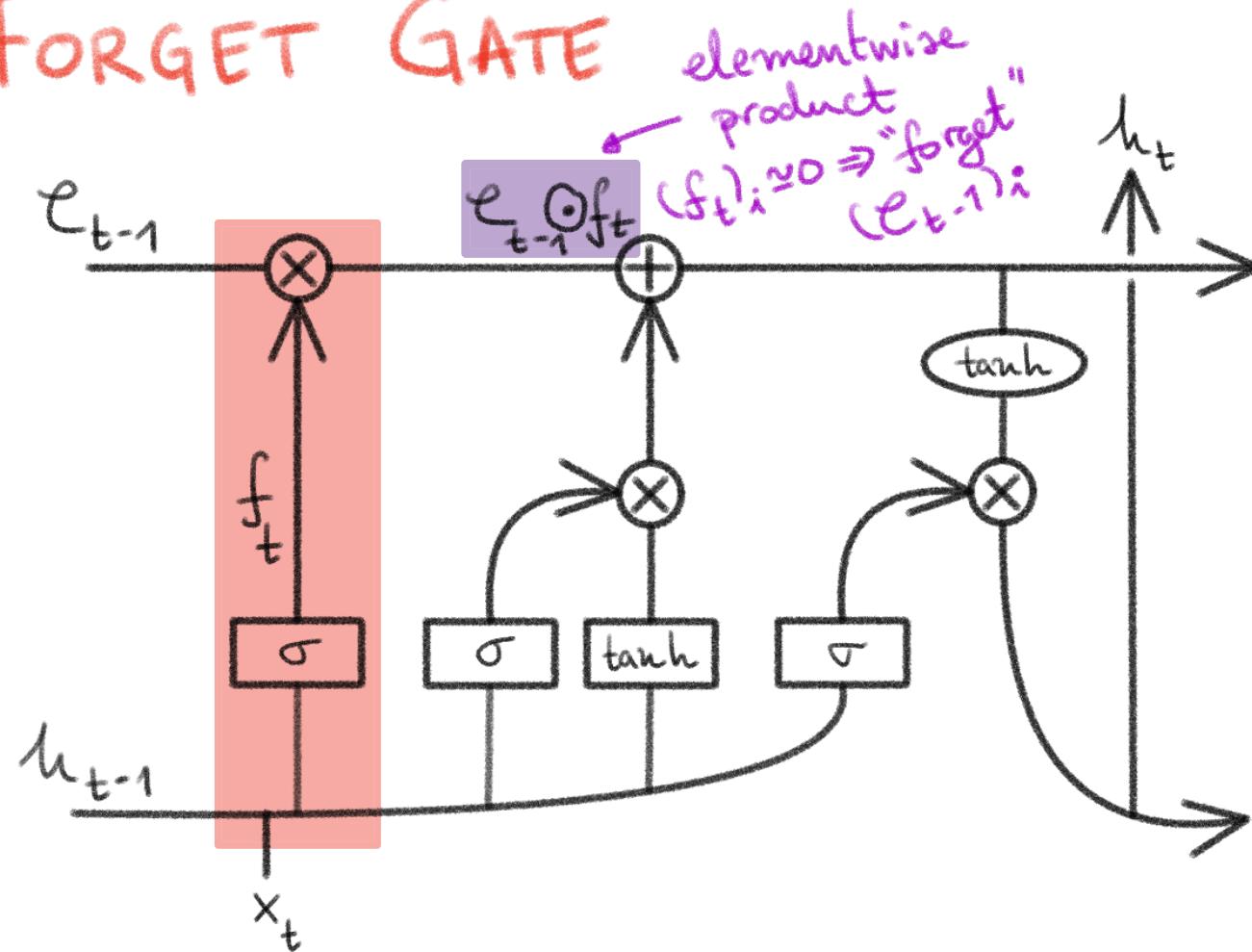
LONG SHORT TERM MEMORY (LSTM)



LONG SHORT TERM MEMORY (LSTM)

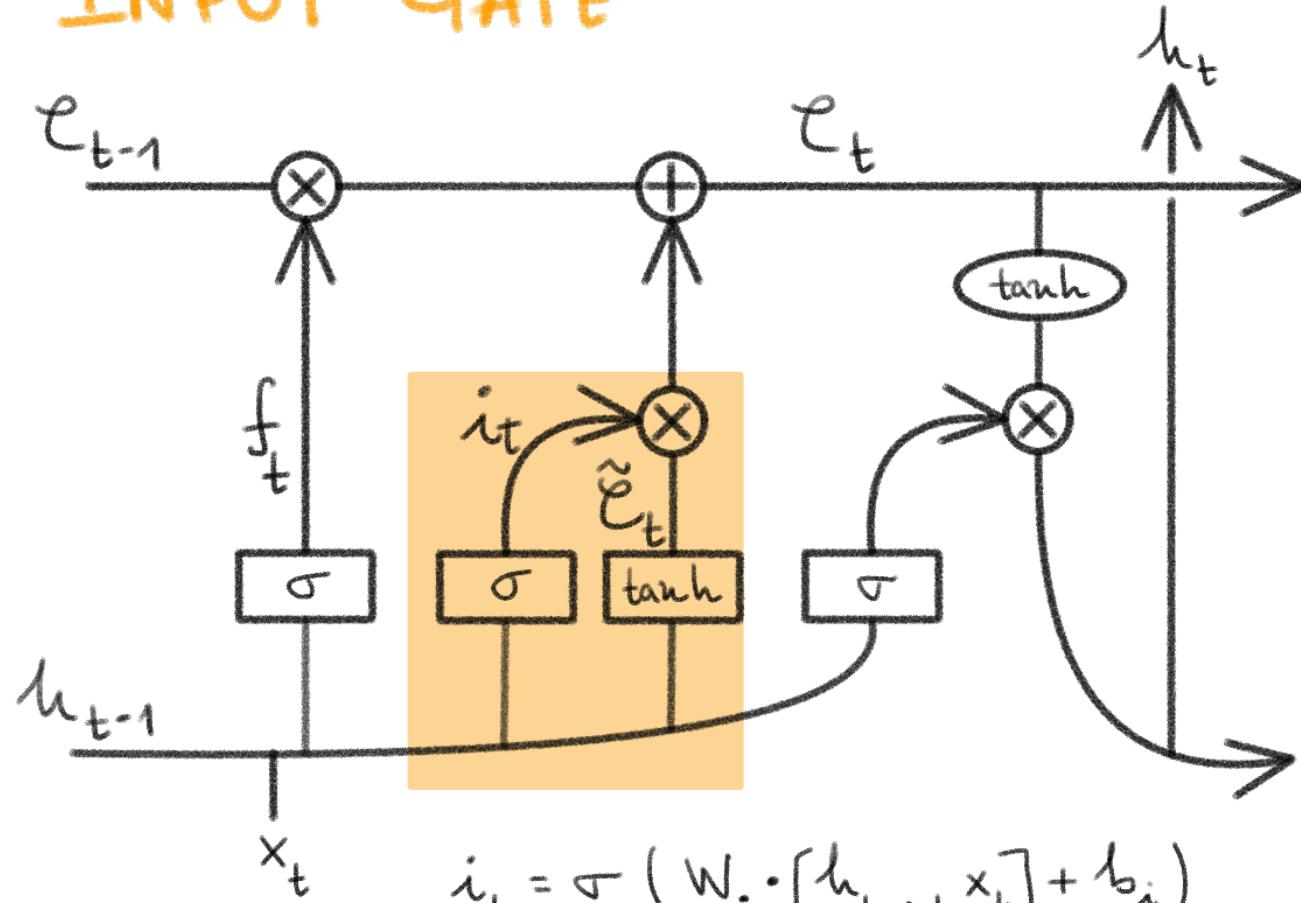


FORGET GATE



$$f_t = \sigma (W_f \cdot [h_{t-1}, x_t] + b_f)$$

INPUT GATE

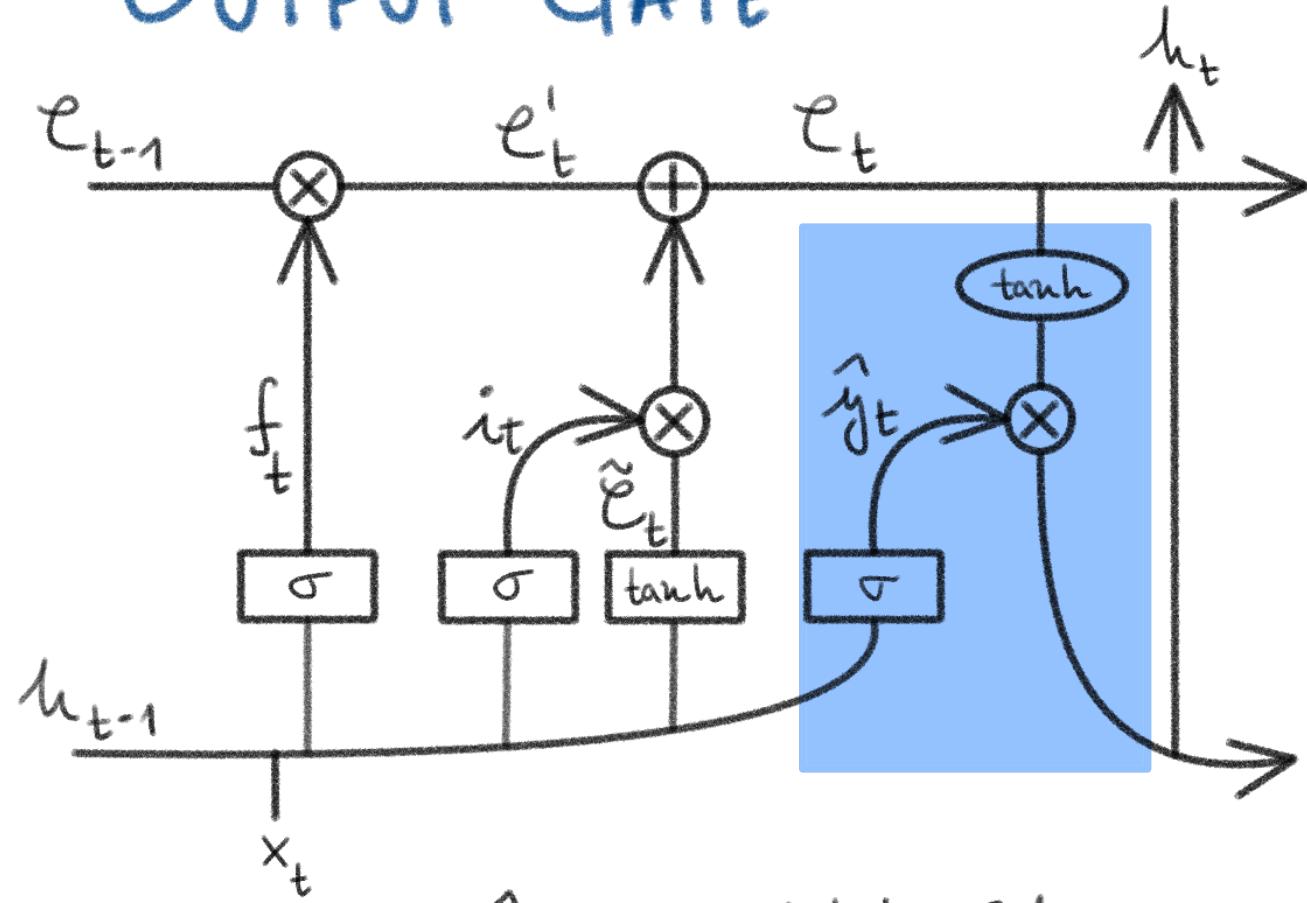


$$i_t = \sigma(W_i \cdot [h_{t-1}, x_t] + b_i)$$

$$\tilde{c}_t = \tanh(W_c \cdot [h_{t-1}, x_t] + b_c)$$

$$c_t = c_{t-1} \odot f_t + i_t \odot \tilde{c}_t$$

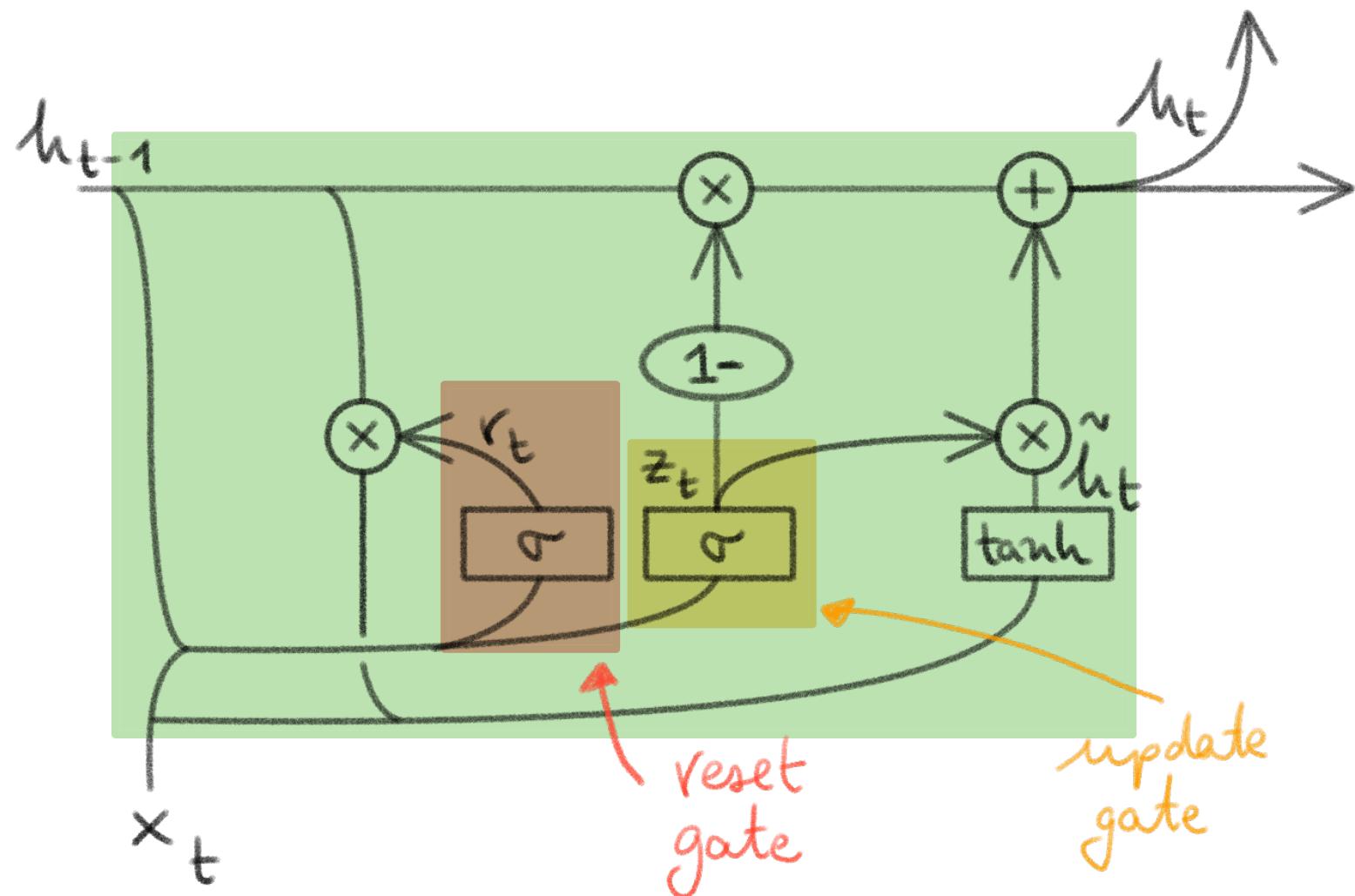
OUTPUT GATE



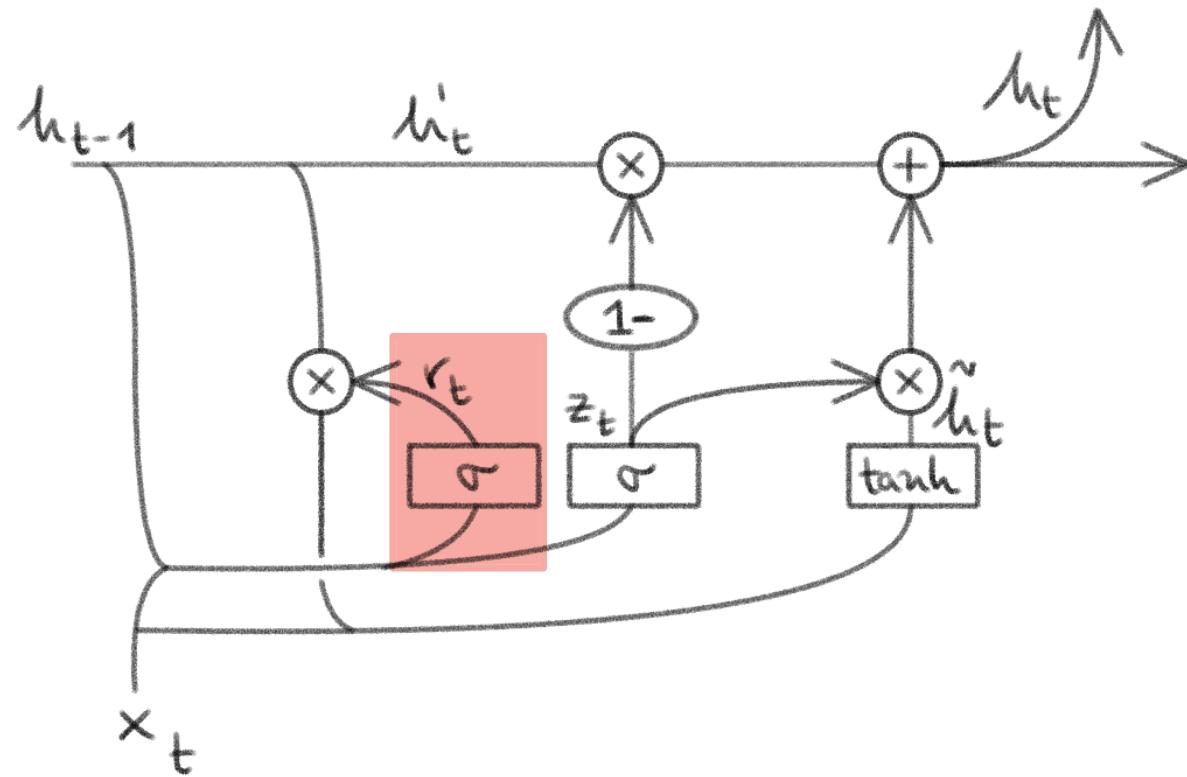
$$\hat{y}_t = \sigma (W_y \odot [h_{t-1}, x_t] + b_y)$$

$$h_t = \hat{y}_t \odot \tanh (c_t)$$

GATED RECURRENT UNIT (GRU)

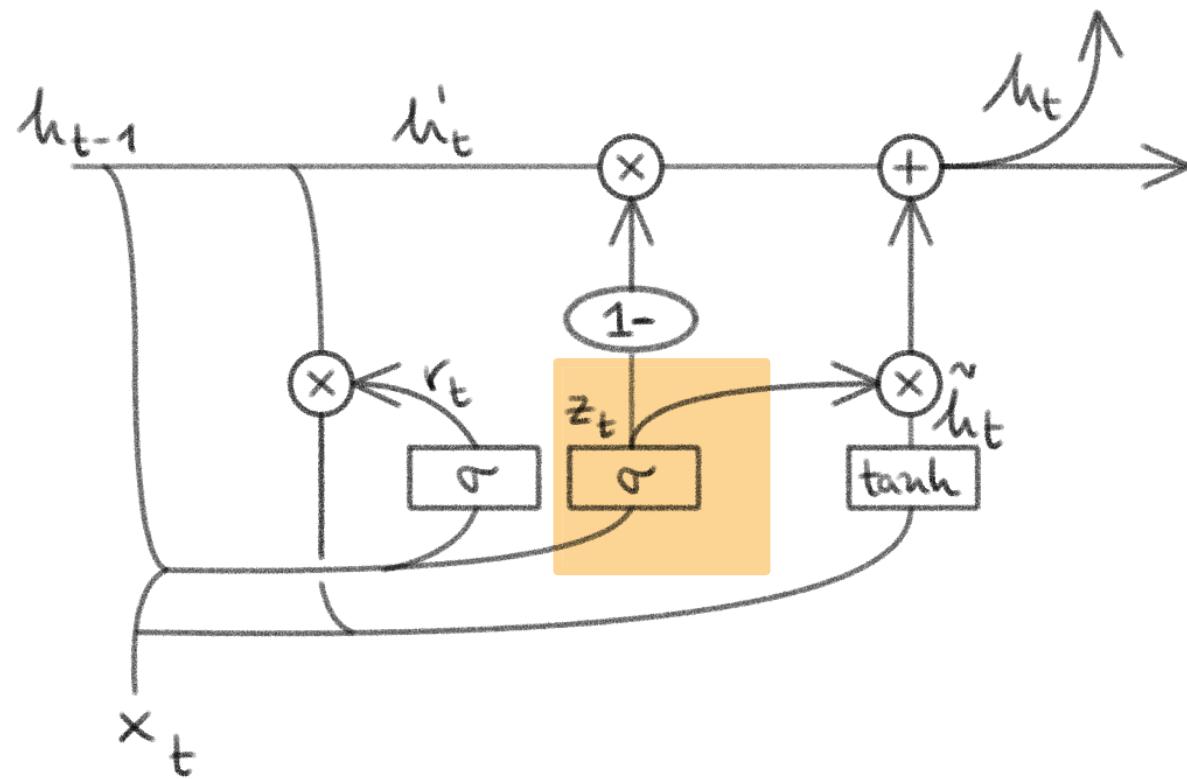


RESET GATE



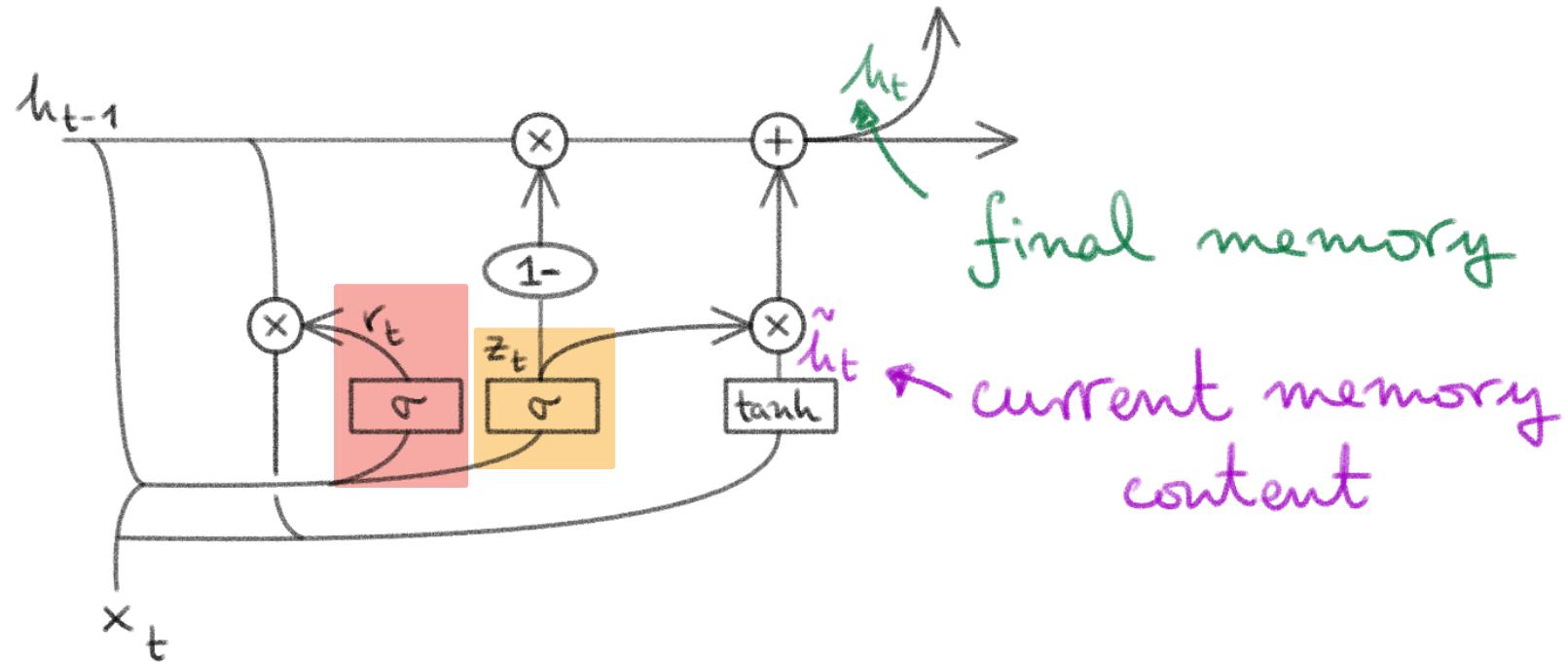
$$r_t = \sigma (W_r \cdot [h_{t-1}, x_t] + b_r)$$

UPDATE GATE



$$z_t = \sigma(W_z \cdot [h_{t-1}, x_t] + b_z)$$

GATED RECURRENT UNIT (GRU)



$$\tilde{h}_t = \tanh(W \cdot [r_t \odot h_{t-1}, x_t] + b_h)$$

$$h_t = (1 - z_t) \odot h_{t-1} + z_t \odot \tilde{h}_t$$

Kaloot

Time!

ARE RNNs / LSTMs STILL USED?

Yes, in the following contexts:

- Time-series forecasting
(e.g. finance, trading)
- On smaller dataset
- When hardware/latency constraints

CONCLUSION

- Draw the (unfolded) graph of a 1-to-many RNN
- Write down the formulas for the 0-th, t-th and T-th step in a 1-to-many RNN
- Write down a question about something you did not fully understand or you would like to know more about.

FURTHER READING

"Introduction to RNNs" in Jeremy Jordan's blog

"The Unreasonable Effectiveness of RNNs"
in Andrej Karpathy's blog