Deep Learning for Natural Language Processing

# The LSTM architecture

Marco Kuhlmann

Department of Computer and Information Science

Linköping University

WASP | WALLENBERG AI, AUTONOMOUS SYSTEMS AND SOFTWARE PROGRAM

# Challenges with recurrent neural networks

- In principle, recurrent neural networks are capable of learning long-distance dependencies in input sequences.

- In practice, training recurrent neural networks is challenging due to the large depth of the unrolled networks.

# Vanishing and exploding gradients



$$\delta_k = \frac{\partial E}{\partial z_k} = \frac{\partial E}{\partial y_k}\frac{\partial y_k}{\partial z_k} = \frac{\partial E}{\partial y_k}f'(z_k)$$

$$\delta_j = \frac{\partial E}{\partial z_j} = \frac{\partial y_j}{\partial z_j}\sum_k\frac{\partial E}{\partial z_k}\frac{\partial z_k}{\partial y_j} = f'(z_j)\sum_k\delta_k w_{jk}$$

# Vanishing and exploding gradients

- In backpropagation there is a risk of gradients either vanishing or exploding, depending on the magnitude of the weights.

- This problem is exacerbated in recurrent networks, whose unrolled computation graphs can be very deep.

- Research on recurrent networks has proposed various methods to mitigate this problem.

  weight scaling and clipping, specialised architectures
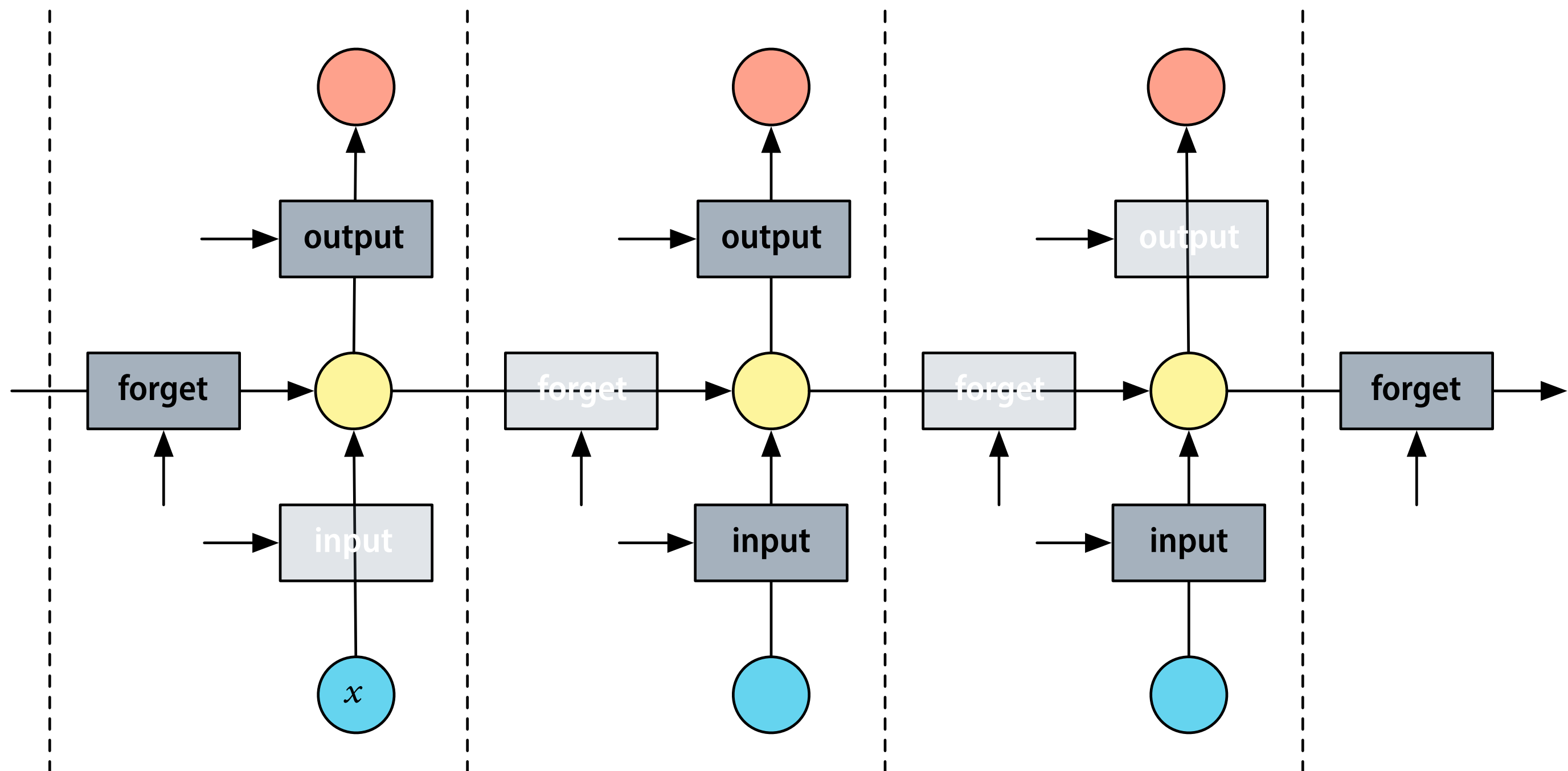
# Long Short-Term Memory (LSTM)

- The **Long Short-Term Memory (LSTM)** architecture was specifically designed to adress the vanishing gradients problem.

- Metaphor: The hidden state of the neural network can be considered as a short-term memory.

- The LSTM architecture tries to make this short-term memory last as long as possible by preventing vanishing gradients.

# Memory cell and gating mechanism

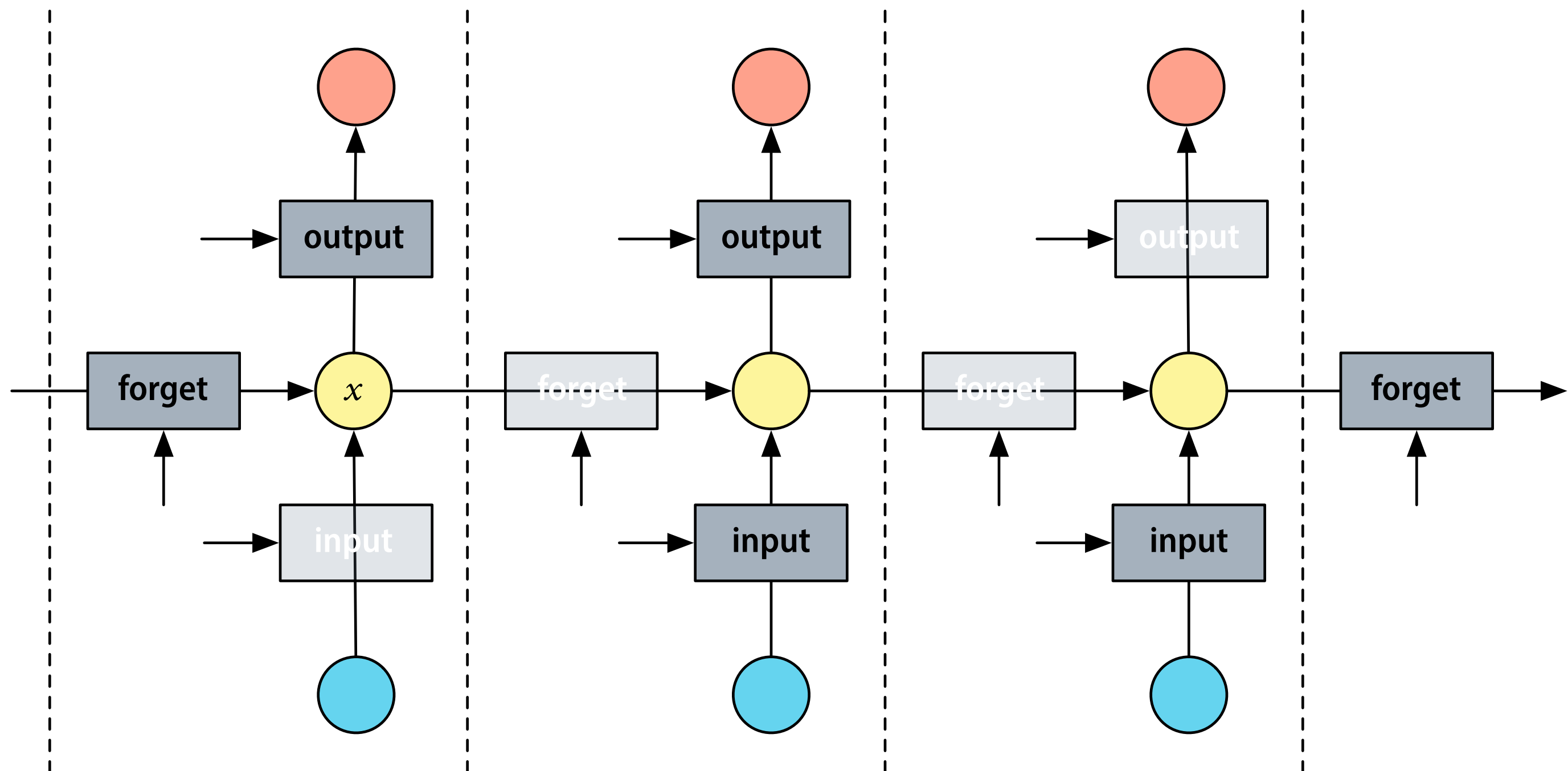The crucial innovation in an LSTM is the design of its memory cell.

- Information is written into the cell if its INPUT gate is open.

- Information stays in the cell as long as its FORGET gate is closed.

- Information is read from the cell if its READ gate is open.
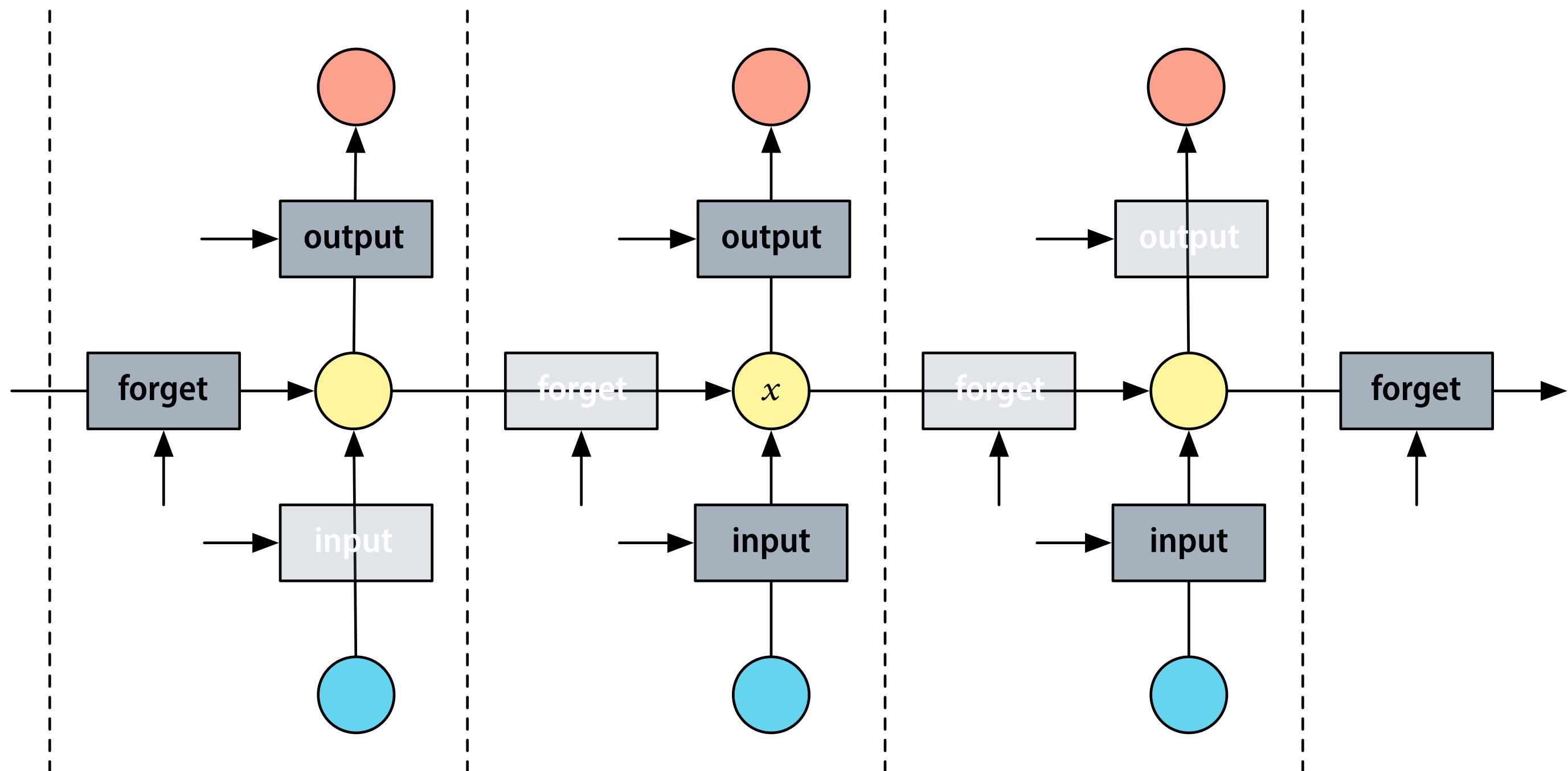
# Information flow in an LSTM

# Information flow in an LSTM

# Information flow in an LSTM

# Information flow in an LSTM

# Information flow in an LSTM

# Gating mechanism

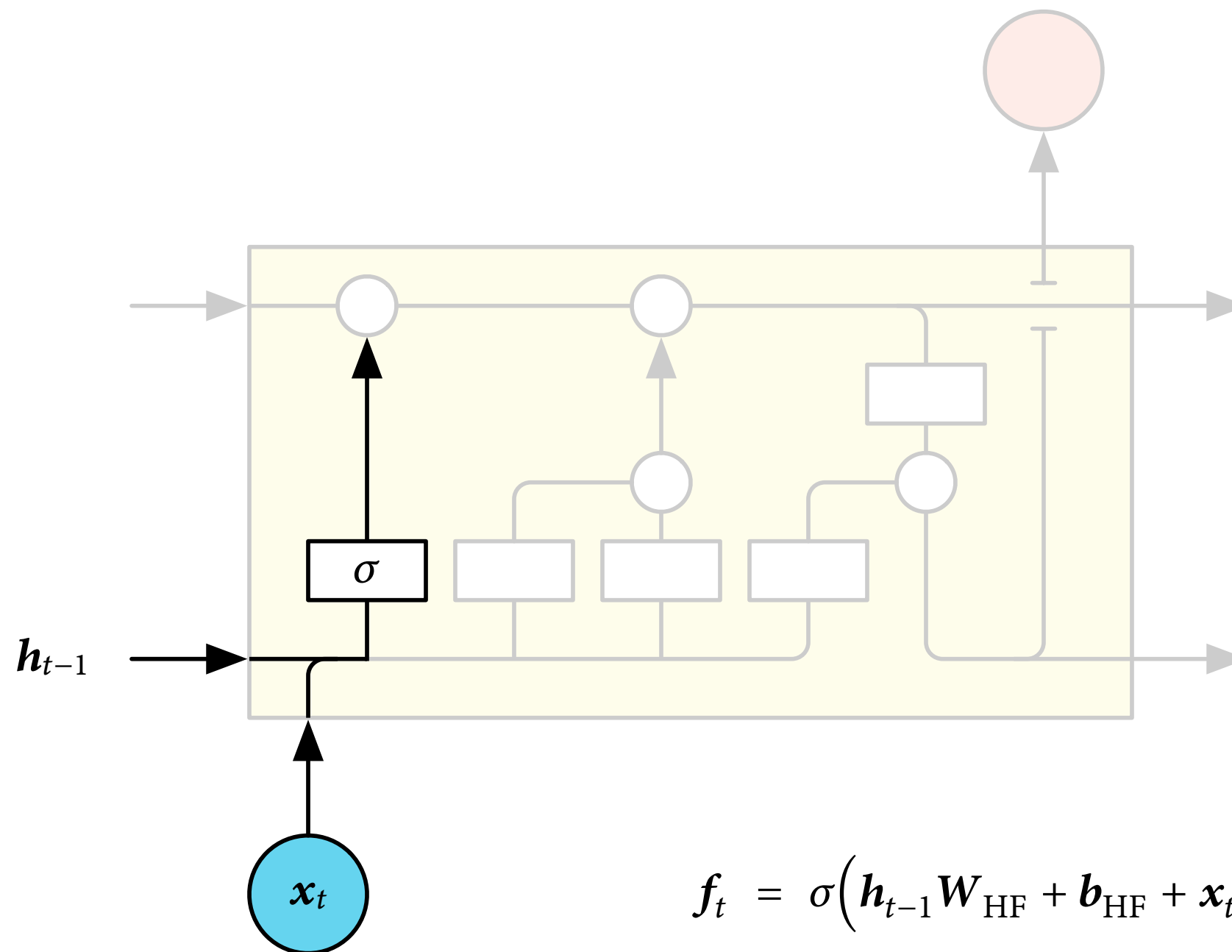$$\begin{bmatrix} 1 \\ 2 \\ 3 \\ 4 \end{bmatrix} \odot \begin{bmatrix} 0 \\ 1 \\ 1 \\ 0 \end{bmatrix} + \begin{bmatrix} 5 \\ 6 \\ 7 \\ 8 \end{bmatrix} \odot \begin{bmatrix} 1 \\ 0 \\ 0 \\ 1 \end{bmatrix} = \begin{bmatrix} 5 \\ 2 \\ 3 \\ 8 \end{bmatrix}$$

$\boldsymbol{h}_{t-1}$ $\quad$ $\boldsymbol{g}$ $\quad$ $\boldsymbol{x}_t$ $\quad$ $1-\boldsymbol{g}$ $\quad$ $\boldsymbol{h}_t$

The gating masks $\boldsymbol{g}$ are learned values between 0 and 1.

# A look inside an LSTM cell

# Forget gate



$$f_t \;=\; \sigma\Big(\boldsymbol{h}_{t-1}\boldsymbol{W}_{\mathrm{HF}} + \boldsymbol{b}_{\mathrm{HF}} + \boldsymbol{x}_t\boldsymbol{W}_{\mathrm{XF}} + \boldsymbol{b}_{\mathrm{XF}}\Big)$$

# Input gate



$$i_t = \sigma\Big(h_{t-1}W_{\text{HI}} + b_{\text{HI}} + x_t W_{\text{XI}} + b_{\text{XI}}\Big)$$

# Update candidate



$$\tilde{\boldsymbol{c}}_t \;=\; \tanh\!\Big(\boldsymbol{h}_{t-1}\boldsymbol{W}_{\text{HC}} + \boldsymbol{b}_{\text{HC}} + \boldsymbol{x}_t\boldsymbol{W}_{\text{XC}} + \boldsymbol{b}_{\text{XC}}\Big)$$

# Memory cell update



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

# Output gate



$$\boldsymbol{o}_t = \sigma\Big(\boldsymbol{h}_{t-1}\boldsymbol{W}_{\mathrm{HO}} + \boldsymbol{b}_{\mathrm{HO}} + \boldsymbol{x}_t\boldsymbol{W}_{\mathrm{XO}} + \boldsymbol{b}_{\mathrm{XO}}\Big)$$

# Output



$$h_t = o_t \odot \tanh(c_t)$$

# A look inside an LSTM cell

# Gated Recurrent Unit (GRU)