



# E9: 309 Advanced Deep Learning

## 14-10-2020

**Instructor: Sriram Ganapathy**  
[sriramg@iisc.ac.in](mailto:sriramg@iisc.ac.in)

**Teaching Assistant : Akshara Soman, Prachi Singh, Jaswanth Reddy**  
[aksharas@iisc.ac.in](mailto:aksharas@iisc.ac.in), [prachis@iisc.ac.in](mailto:prachis@iisc.ac.in), [jaswanthk@iisc.ac.in](mailto:jaswanthk@iisc.ac.in)

**Schedule - MW - 430-6pm (Microsoft Teams)**  
<http://leap.ee.iisc.ac.in/sriram/teaching/ADL2020/>



# Housekeeping

- \* Filling the google form in the webpage
  - Contents will be made available to the folks in the creditors mailing lists.
    - ✓ Announcements regarding evaluations and projects will be shared only with creditors as well as video links.
      - ★ Teams channel interaction and TA session for creditors only.
- \* Online registration portal from [academics.iisc.ac.in](https://academics.iisc.ac.in)
  - ✓ Your research/faculty advisor may need to approve also before the deadline (Oct. 20th?)

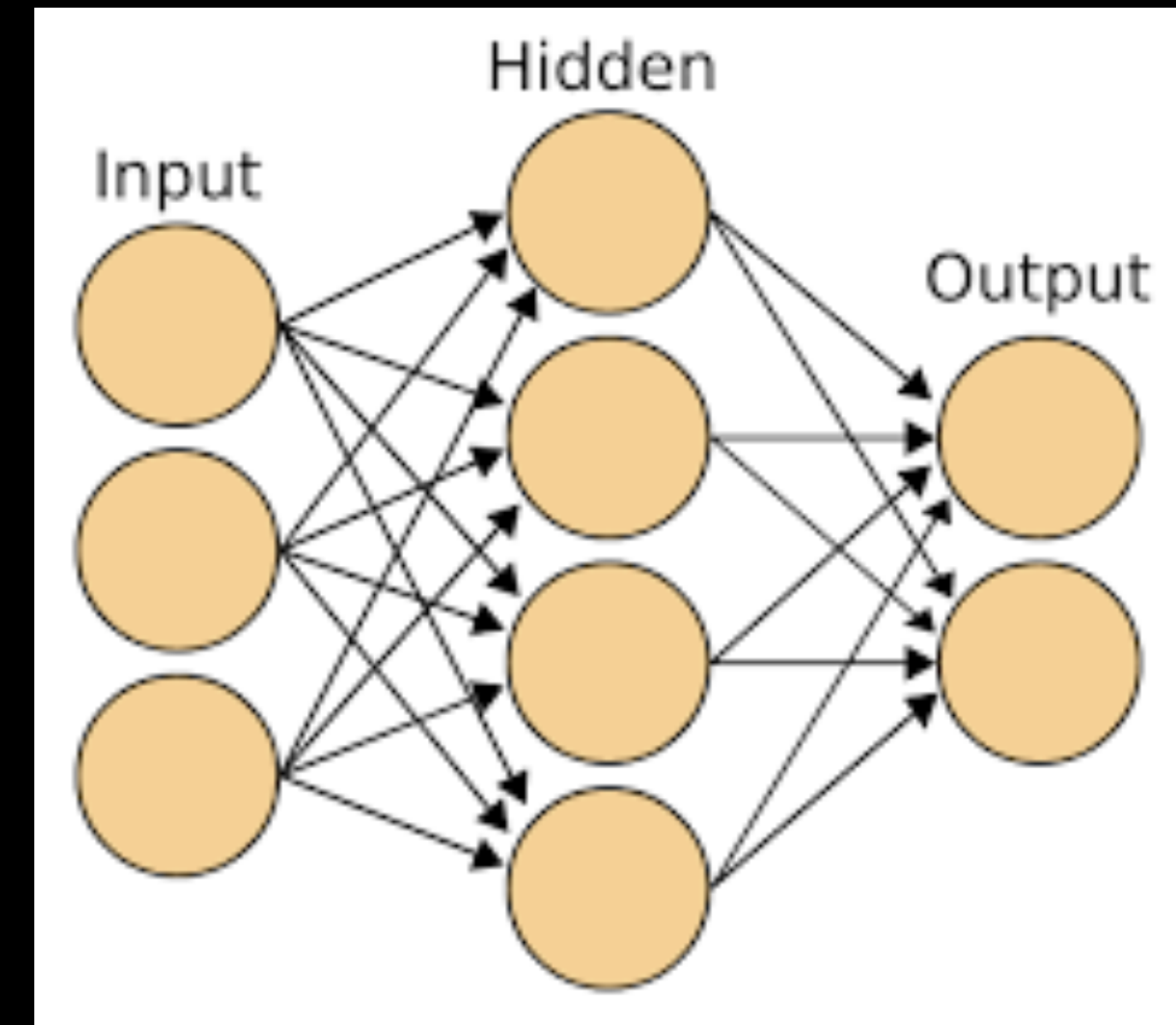


# Recap of previous class



# Some notations

- \*  $\mathbf{x} \in \mathcal{R}^D$  - input data.
- \*  $\mathbf{y} \in \mathcal{R}^C$  - neural network targets.
- \*  $\hat{\mathbf{y}} \in \mathcal{B}^C$  - model outputs.
- \*  $\mathbf{e}, \mathbf{h} \in \mathcal{R}^d$  - hidden model representations or embeddings.
- \*  $\Theta$  - collection of learnable parameters in the model.
- \*  $E(\mathbf{y}, \hat{\mathbf{y}})$  - error function used in the model training.



# Some notations

- \*  $\{\mathbf{x}_1, \dots, \mathbf{x}_N, \mathbf{y}_1, \dots, \mathbf{y}_N\}$  - labeled training data
- \*  $q = \{1 \dots Q\}$  - iteration index.
- \*  $t = \{1 \dots T\}$  - discrete time index.
- \*  $l = \{1 \dots L\}$  - layer index
- \*  $\eta$  - learning rate (hyper-parameter)
- \*  $N_b$  - mini-batch size and  $B$  is the number of mini-batches.



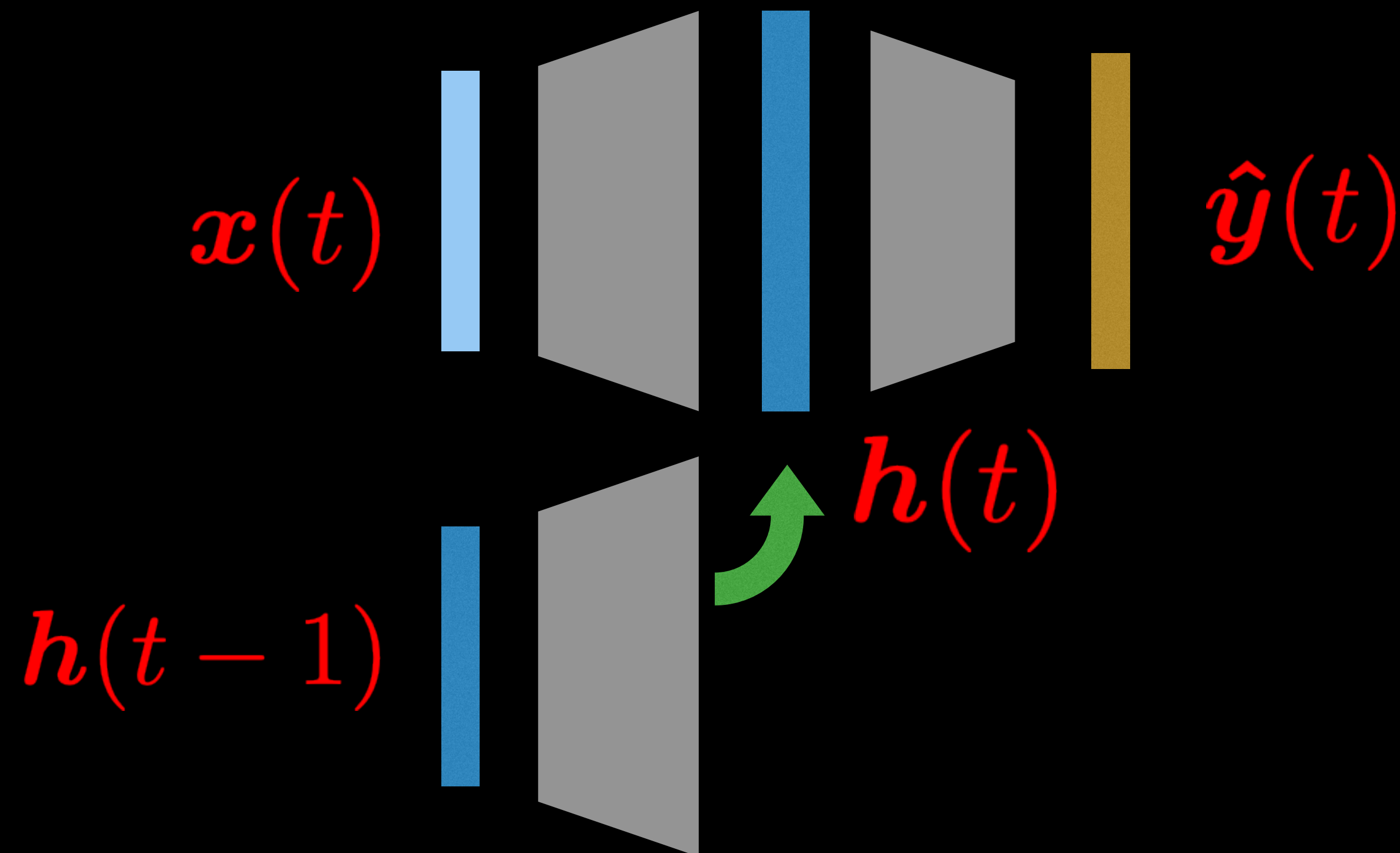
# *Module - I Visual and Time Series Modeling*



# First order recurrence - hidden layer

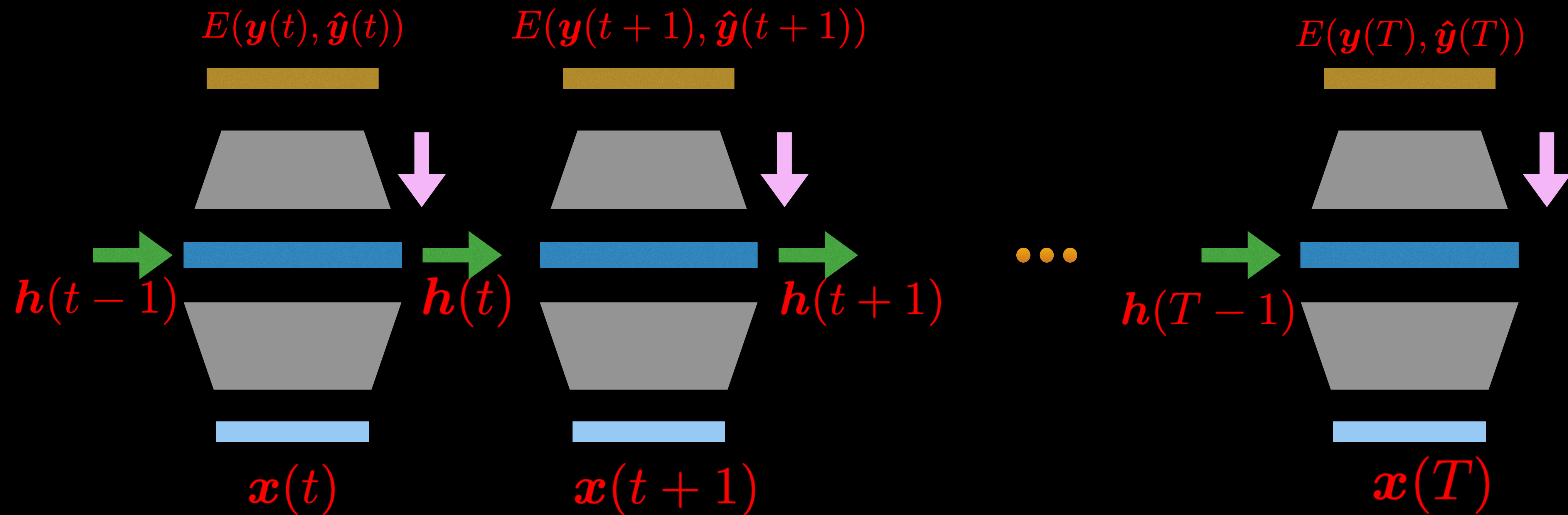
- \* Making the hidden layer a function of the previous outputs from the hidden layer along with the input

$$h(t) = f(h(t-1), x(t))$$



# Error backpropagation

- \* Error functions are computed at every time-instant



- \* Total error 
$$E = \sum_t E(\mathbf{y}(t), \hat{\mathbf{y}}(t))$$



# Error back propagation

- ✓ Output activations

$$\frac{\partial E}{\partial \mathbf{a}^2(t)} = \hat{\mathbf{y}}(t) - \mathbf{y}(t) \quad \text{for } t = 1 \dots T$$

- ✓ Hidden activations at last instant  $T$

$$\frac{\partial E}{\partial \mathbf{h}^1(T)} = (\mathbf{W}^2)^T \frac{\partial E}{\partial \mathbf{a}^2(T)}$$

- ✓ Hidden activations for previous instances  $t = T-1, \dots, 1$

$$\frac{\partial E}{\partial \mathbf{h}^1(t)} = (\mathbf{W}^2)^T \frac{\partial E}{\partial \mathbf{a}^2(t)} + \frac{\partial \mathbf{h}^1(t+1)}{\partial \mathbf{h}^1(t)} \frac{\partial E}{\partial \mathbf{h}^1(t+1)}$$

$$\begin{aligned} \mathbf{a}^1(t) &= \mathbf{W}^1 \mathbf{x}(t) + \mathbf{U}^1 \mathbf{h}^1(t-1) + b^1 \\ \mathbf{h}^1(t) &= \tanh(\mathbf{a}^1(t)) \\ \mathbf{a}^2(t) &= \mathbf{W}^2 \mathbf{h}^1(t) + b^2 \\ \hat{\mathbf{y}}(t) &= S(\mathbf{a}^2(t)) \end{aligned}$$



# Error back propagation

- ✓ Hidden activations for previous instances  $t = T-1, \dots, 1$

$$\begin{aligned} \mathbf{a}^1(t) &= \mathbf{W}^1 \mathbf{x}(t) + \mathbf{U}^1 \mathbf{h}^1(t-1) + b^1 \\ \mathbf{h}^1(t) &= \tanh(\mathbf{a}^1(t)) \\ \mathbf{a}^2(t) &= \mathbf{W}^2 \mathbf{h}^1(t) + b^2 \\ \hat{\mathbf{y}}(t) &= S(\mathbf{a}^2(t)) \end{aligned}$$

$$\frac{\partial E}{\partial \mathbf{h}^1(t)} = (\mathbf{W}^2)^T \frac{\partial E}{\partial \mathbf{a}^2(t)} + \frac{\partial \mathbf{h}^1(t+1)}{\partial \mathbf{h}^1(t)} \frac{\partial E}{\partial \mathbf{h}^1(t+1)}$$

$$\frac{\partial E}{\partial \mathbf{h}^1(t)} = (\mathbf{W}^2)^T \frac{\partial E}{\partial \mathbf{a}^2(t)} + (\mathbf{U}^1)^T \text{diag}(1 - (\mathbf{h}^1(t+1))^2) \frac{\partial E}{\partial \mathbf{h}^1(t+1)}$$

- ✓ Here, the term  $\text{diag}(1 - \mathbf{h}(t+1)^2)$  comes from the derivative of  $\tanh$
- ✓ and the notation  $\cdot^2$  denotes element wise operation of squaring.





# Error back propagation

- ✓ The derivatives of the output weights.

$$\frac{\partial E}{\partial \mathbf{W}^2} = \sum_t \frac{\partial E}{\partial \mathbf{a}^2(t)} (\mathbf{h}^1(t))^T$$

$$\frac{\partial E}{\partial \mathbf{b}^2} = \sum_t \frac{\partial E}{\partial \mathbf{a}^2(t)}$$

$$\begin{aligned} \mathbf{a}^1(t) &= \mathbf{W}^1 \mathbf{x}(t) + \mathbf{U}^1 \mathbf{h}^1(t-1) + b^1 \\ \mathbf{h}^1(t) &= \tanh(\mathbf{a}^1(t)) \\ \mathbf{a}^2(t) &= \mathbf{W}^2 \mathbf{h}^1(t) + b^2 \\ \hat{\mathbf{y}}(t) &= S(\mathbf{a}^2(t)) \end{aligned}$$

- ✓ Transferring derivatives to the first layer

$$\frac{\partial E}{\partial \mathbf{a}^1(t)} = \text{diag}(1 - \mathbf{h}(t).^2) \frac{\partial E}{\partial \mathbf{h}^1(t)}$$



# Error back propagation

- ✓ The derivatives of the first layer weights.

$$\frac{\partial E}{\partial \mathbf{W}^1} = \sum_t \frac{\partial E}{\partial \mathbf{a}^1(t)} (\mathbf{x}^1(t))^T$$

$$\frac{\partial E}{\partial \mathbf{U}^1} = \sum_t \frac{\partial E}{\partial \mathbf{a}^1(t)} (\mathbf{h}^1(t))^T$$

$$\frac{\partial E}{\partial \mathbf{b}^1} = \sum_t \frac{\partial E}{\partial \mathbf{a}^1(t)}$$

$$\begin{aligned} \mathbf{a}^1(t) &= \mathbf{W}^1 \mathbf{x}(t) + \mathbf{U}^1 \mathbf{h}^1(t-1) + b^1 \\ \mathbf{h}^1(t) &= \tanh(\mathbf{a}^1(t)) \\ \mathbf{a}^2(t) &= \mathbf{W}^2 \mathbf{h}^1(t) + b^2 \\ \hat{\mathbf{y}}(t) &= S(\mathbf{a}^2(t)) \end{aligned}$$





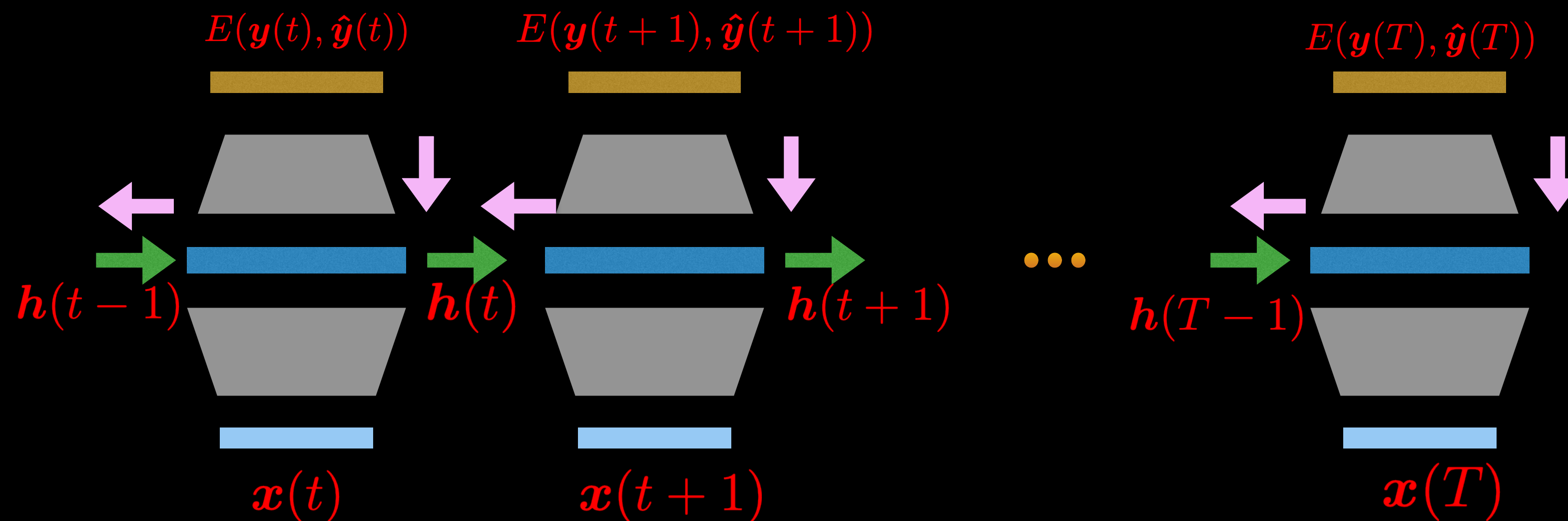
# Error Backpropagation

\* Key equation in the backward direction

$$\frac{\partial E}{\partial \mathbf{h}^1(t)} = (\mathbf{W}^2)^T \frac{\partial E}{\partial \mathbf{a}^2(t)} + (\mathbf{U}^1)^T \text{diag}(1 - (\mathbf{h}^1(t+1))^2) \frac{\partial E}{\partial \mathbf{h}^1(t+1)}$$

\* When the model incorporates a recurrence in the forward direction.

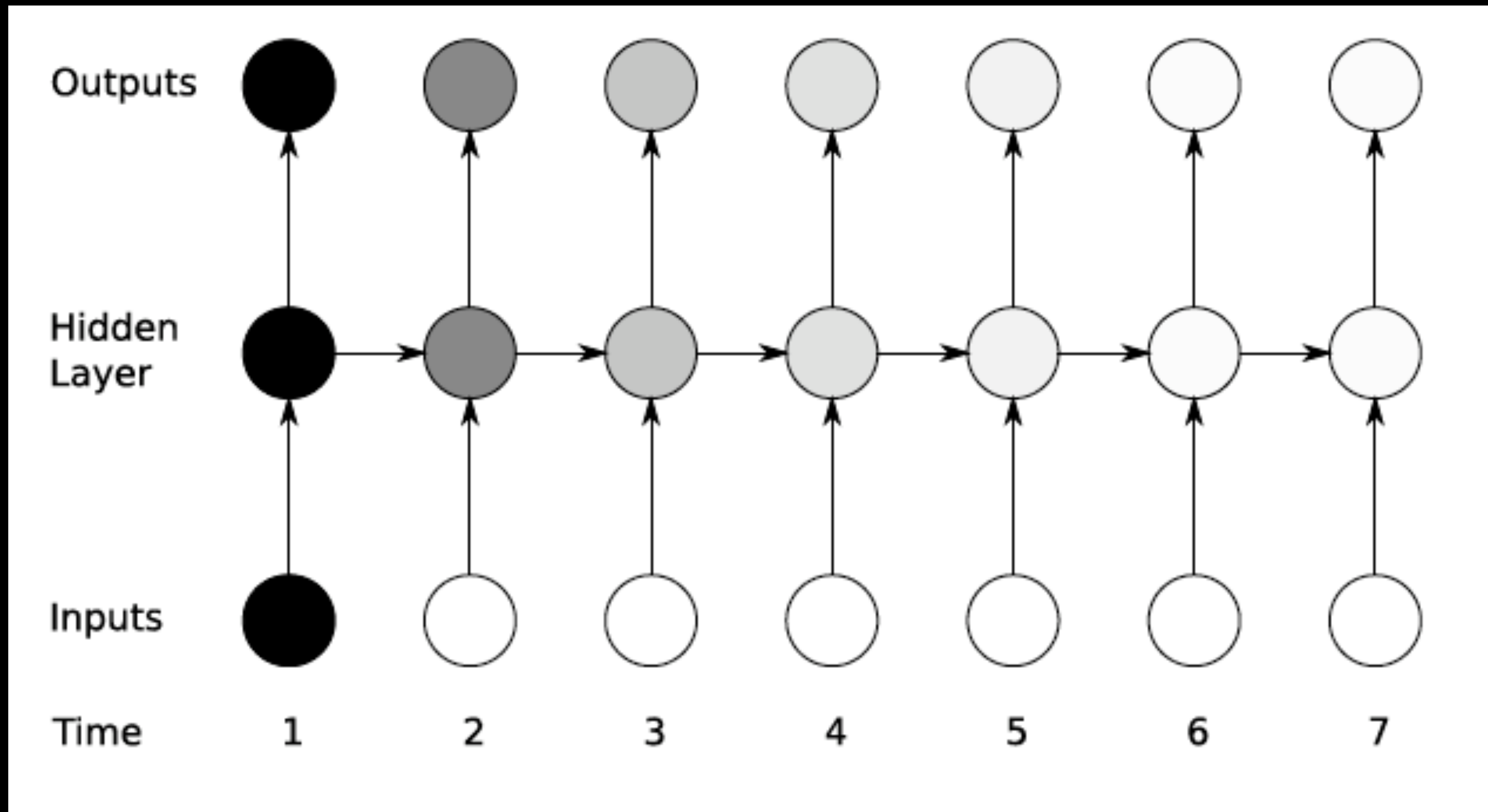
✓ Gradients incorporate a recurrence in the backward direction.





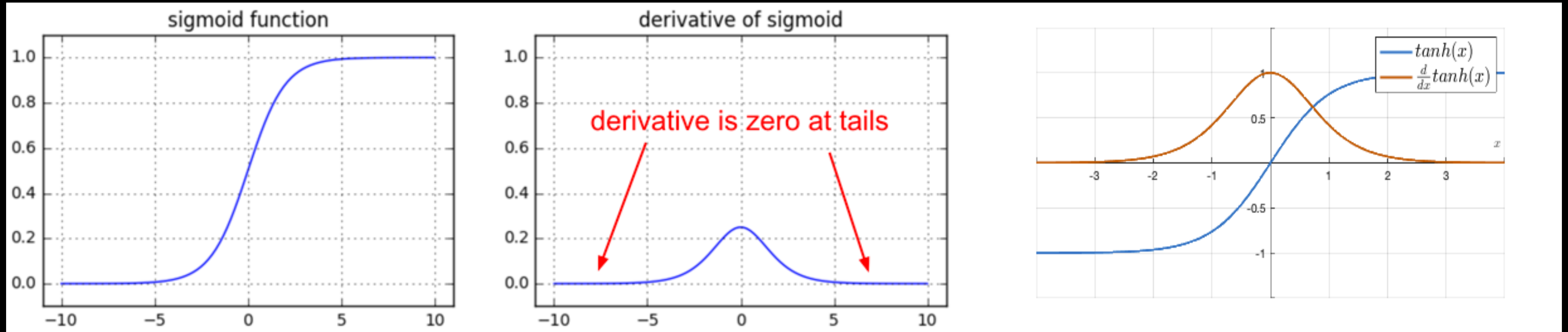


# Long-term dependency issues



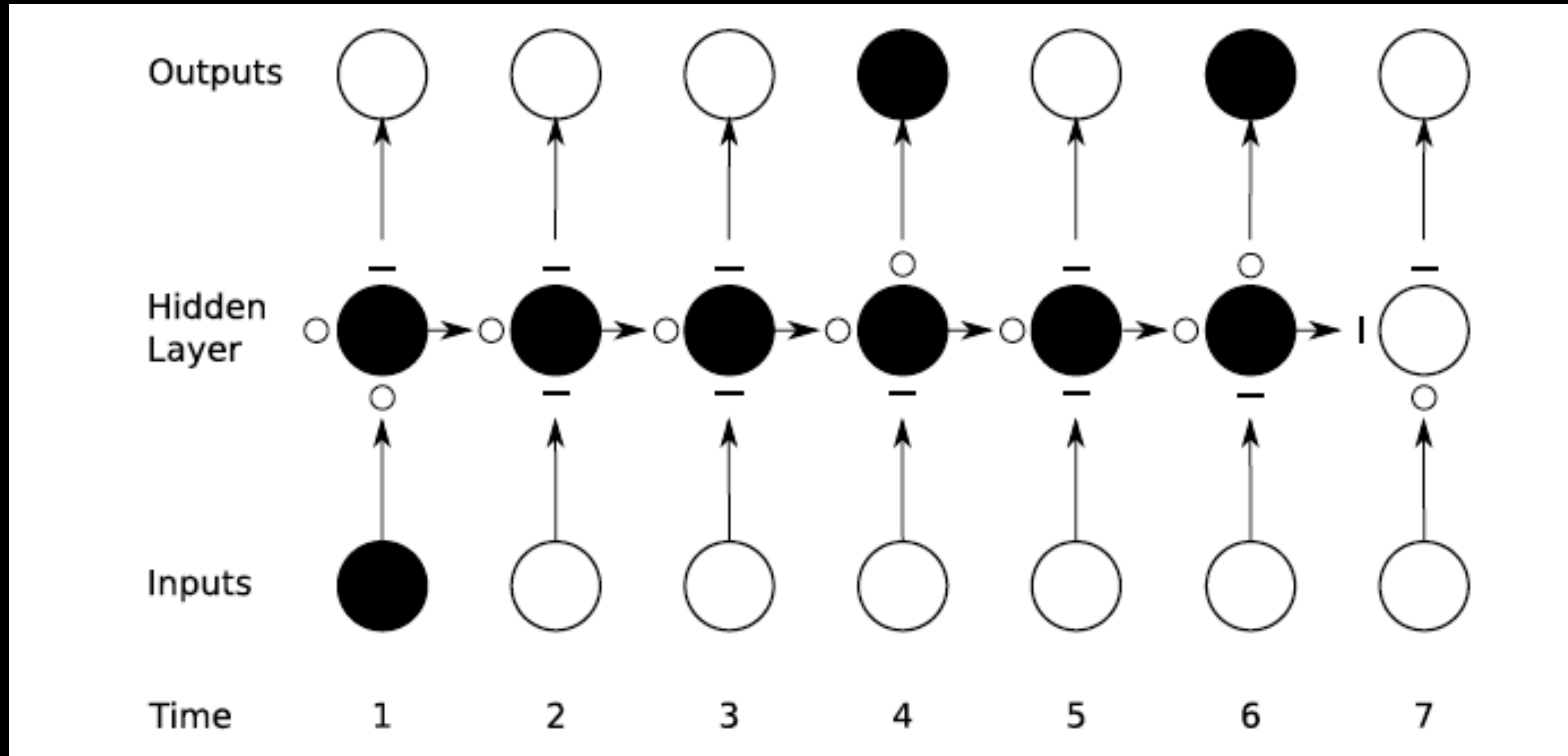
# Long-term dependency issues

- ✓ Gradients tend to vanish or explode



- ✓ Initial frames may not have impact in the final predictions.

# Long short term memory (LSTM) idea



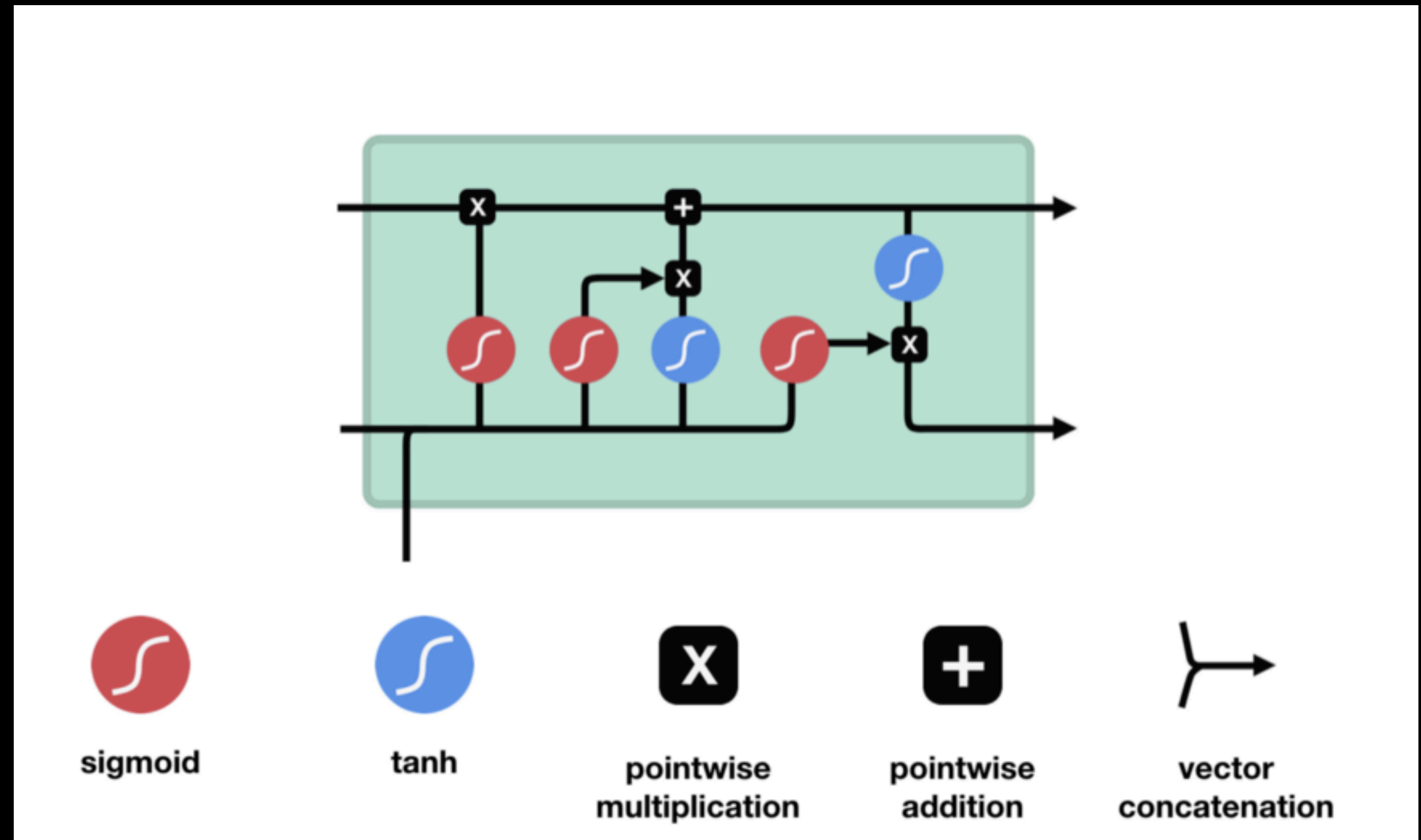


# Modeling questions

- \* How can we make adaptable gates with neural networks
  - How can we make gates dependent on the data itself.
    - ✓ Gates can be implemented as neural layers with sigmoidal outputs ?
      - ★ Sigmoids can approximate 0-1 functions
  - ◉ Modulate the gate output with inputs, hidden layer outputs or outputs



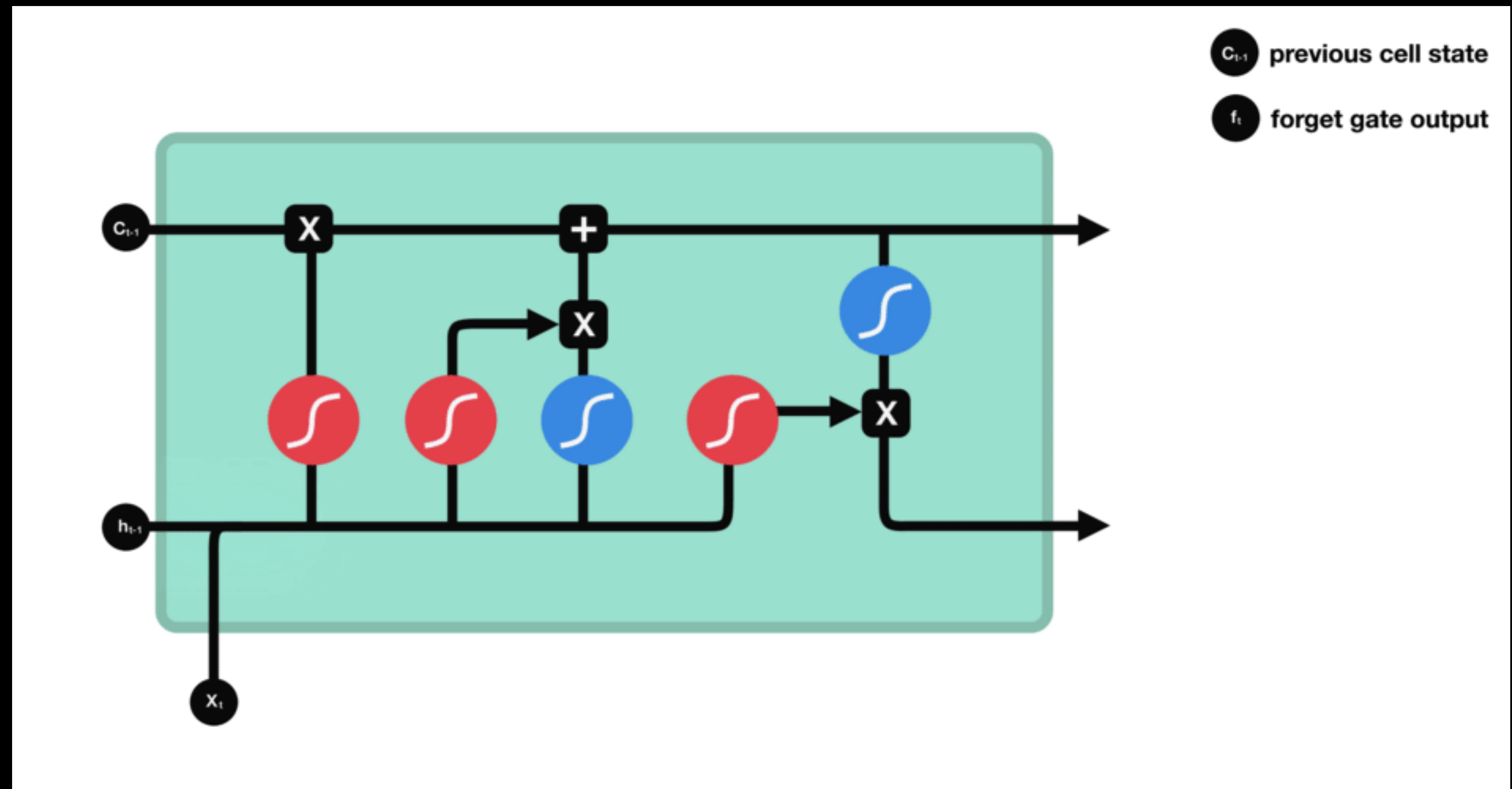
# Long-short term memory - idea



# Long short term memory - idea

\* Forget gate

$$a_{\phi}^t = \sum_{i=1}^I w_{i\phi} x_i^t + \sum_{h=1}^H w_{h\phi} b_h^{t-1}$$
$$b_{\phi}^t = f(a_{\phi}^t)$$

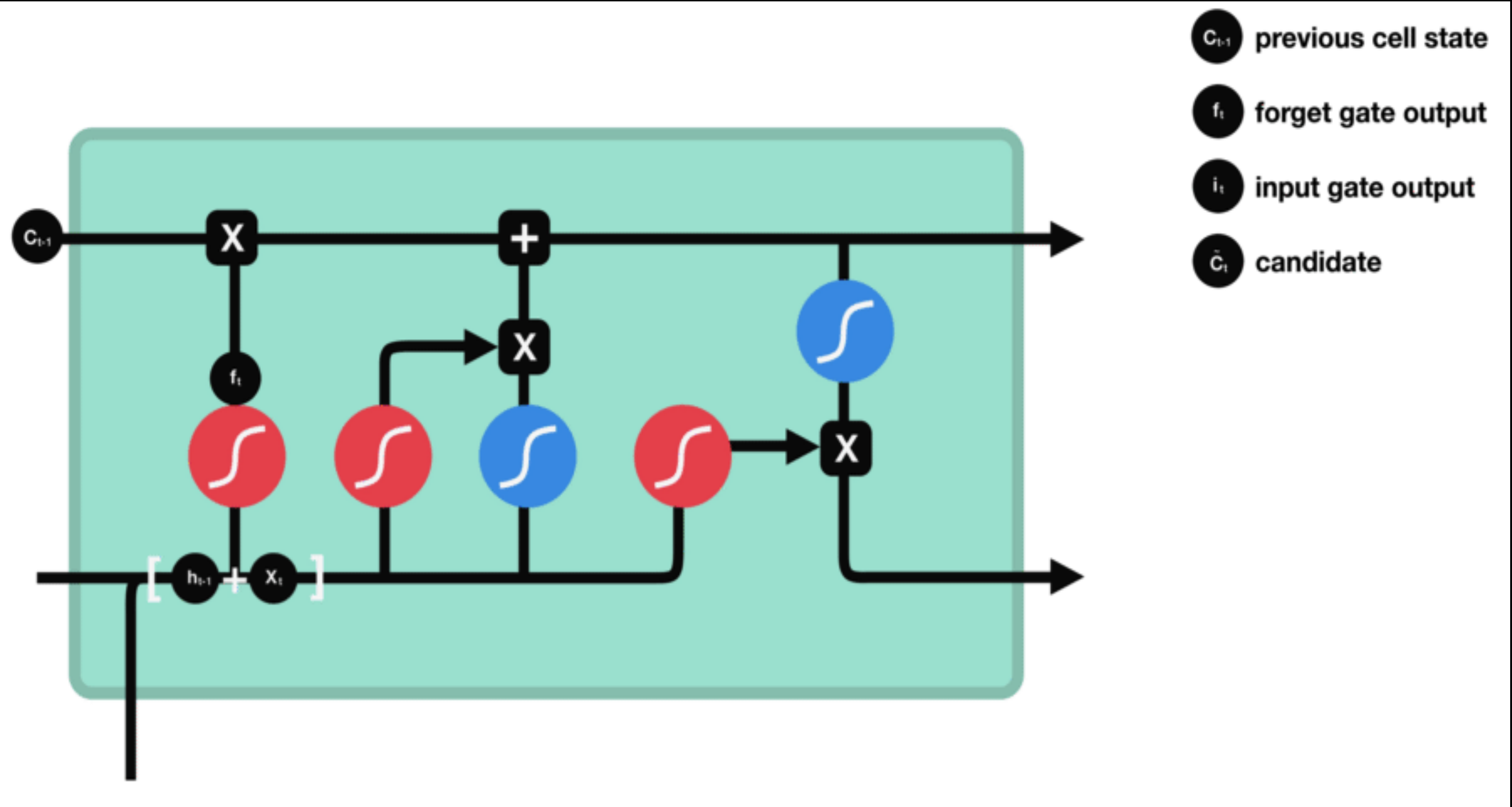




# Long short term memory - idea

## \* Input gate

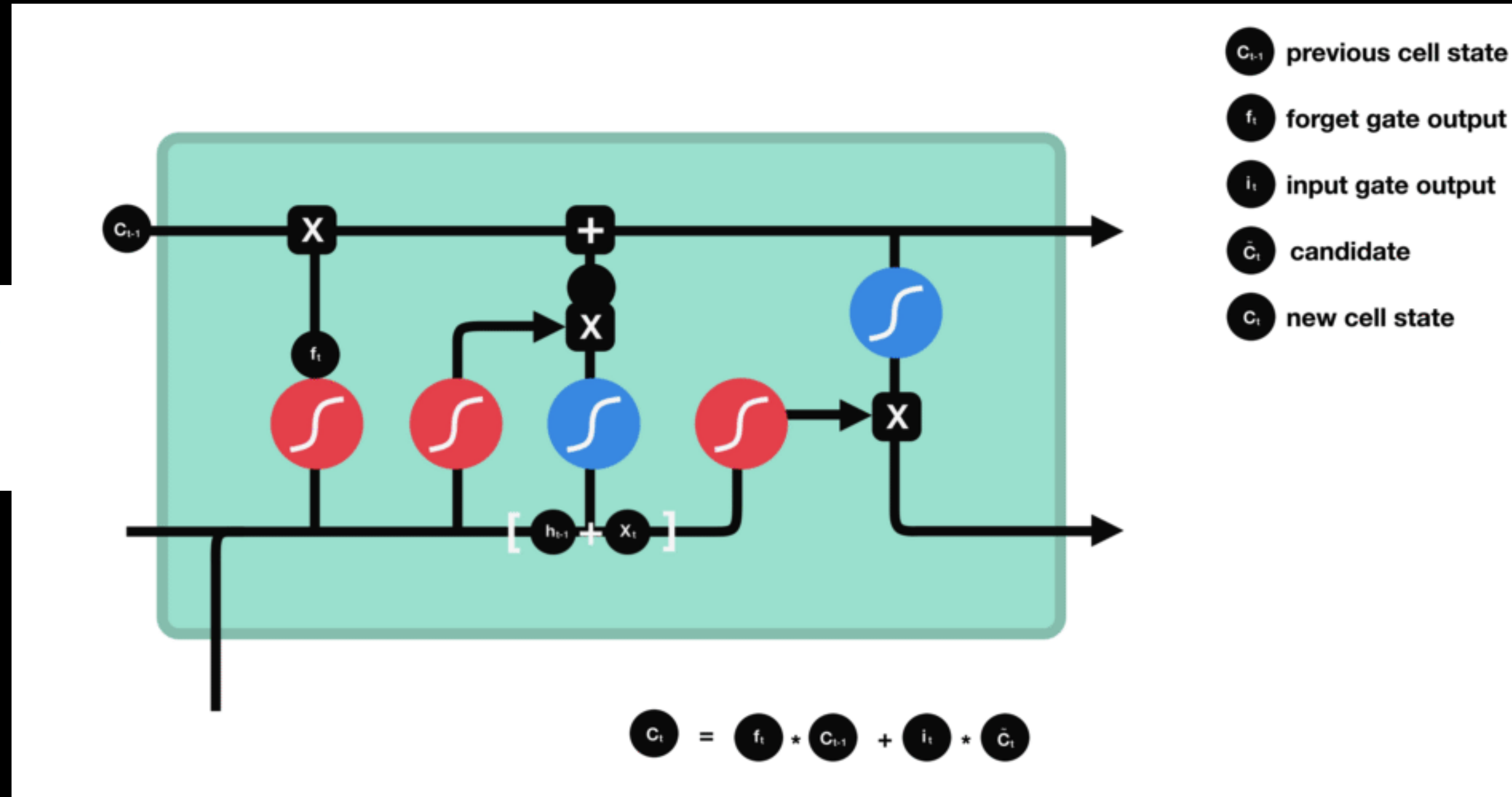
$$a_i^t = \sum_{i=1}^I w_{ii} x_i^t + \sum_{h=1}^H w_{hi} b_h^{t-1}$$
$$b_i^t = f(a_i^t)$$



# Long-short term memory - idea

## \* Cell state

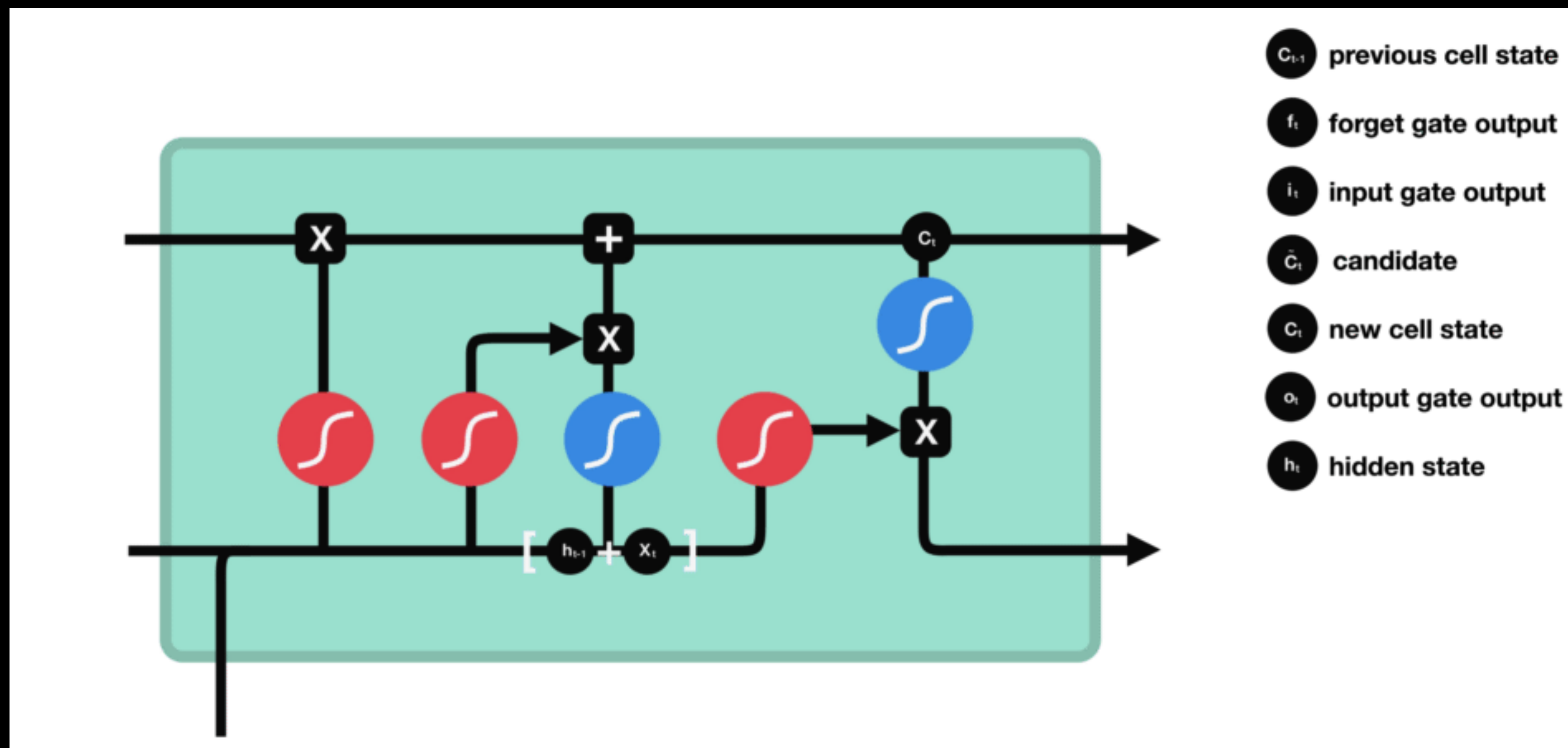
$$a_c^t = \sum_{i=1}^I w_{ic} x_i^t + \sum_{h=1}^H w_{hc} b_h^{t-1}$$
$$s_c^t = b_{\phi}^t s_c^{t-1} + b_i^t g(a_c^t)$$



# Long-short term memory - idea

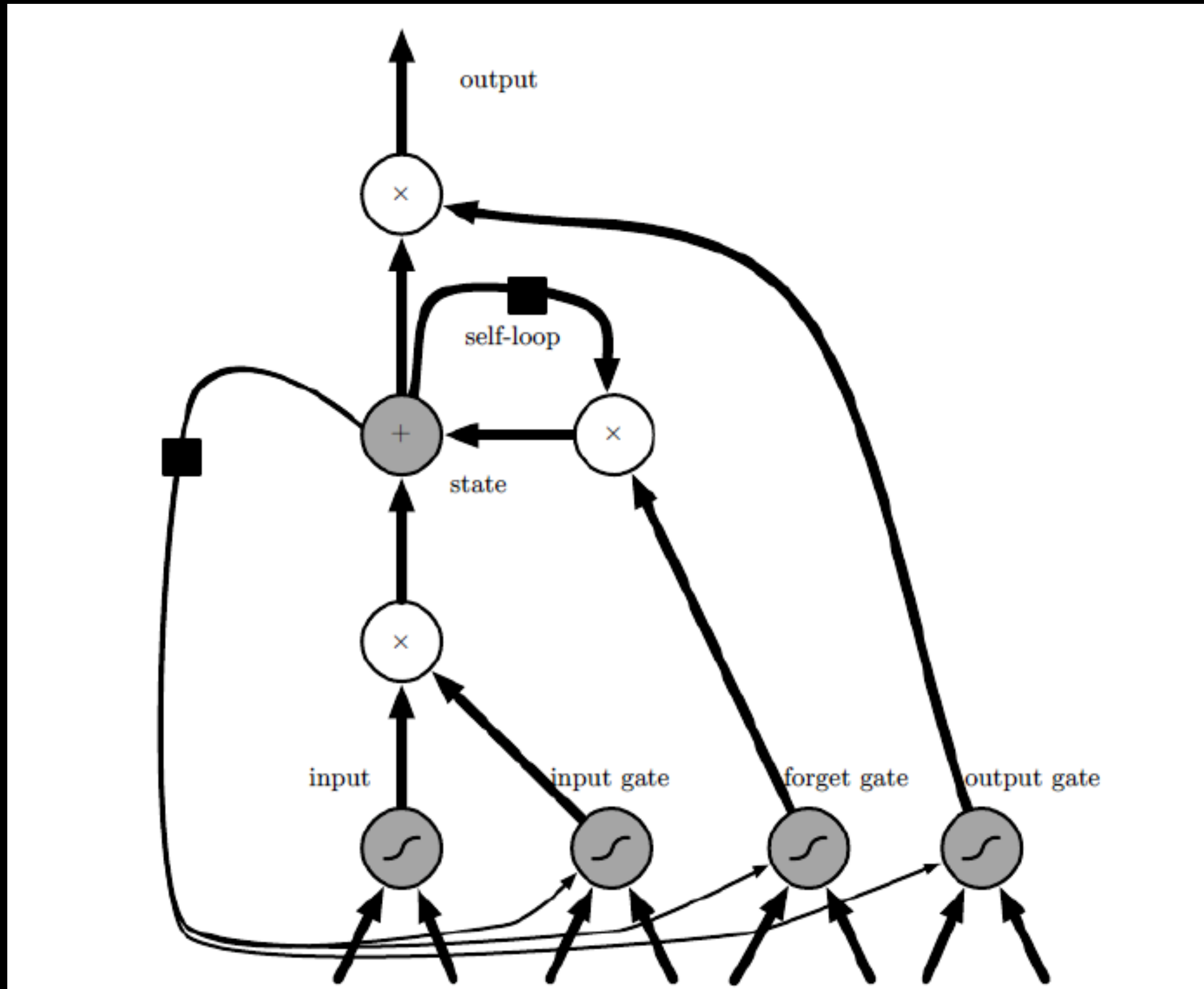
## \* Output gate

$$a_{\omega}^t = \sum_{i=1}^I w_{i\omega} x_i^t + \sum_{h=1}^H w_{h\omega} b_h^{t-1}$$
$$b_{\omega}^t = f(a_{\omega}^t)$$





# Long short-term memory - idea



# Long-short term memory and GRUs

