

Long-Short Term Memory

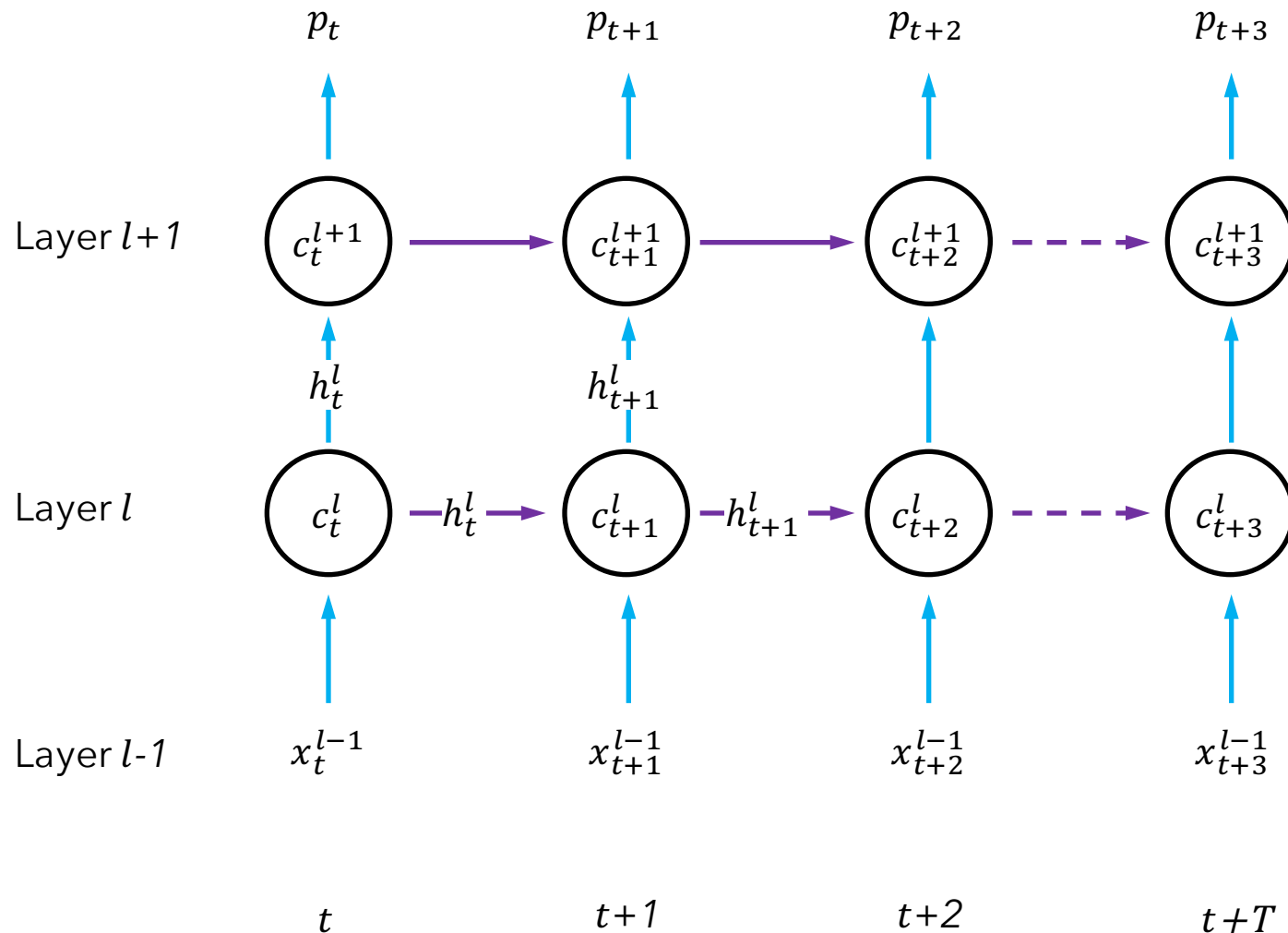
Authors: Hochreiter Sepp, Jürgen Schmidhuber

Presented by: Zihang Dai

Outline

- Motivation
- LSTM
- Experiments

Recurrent Neural Network

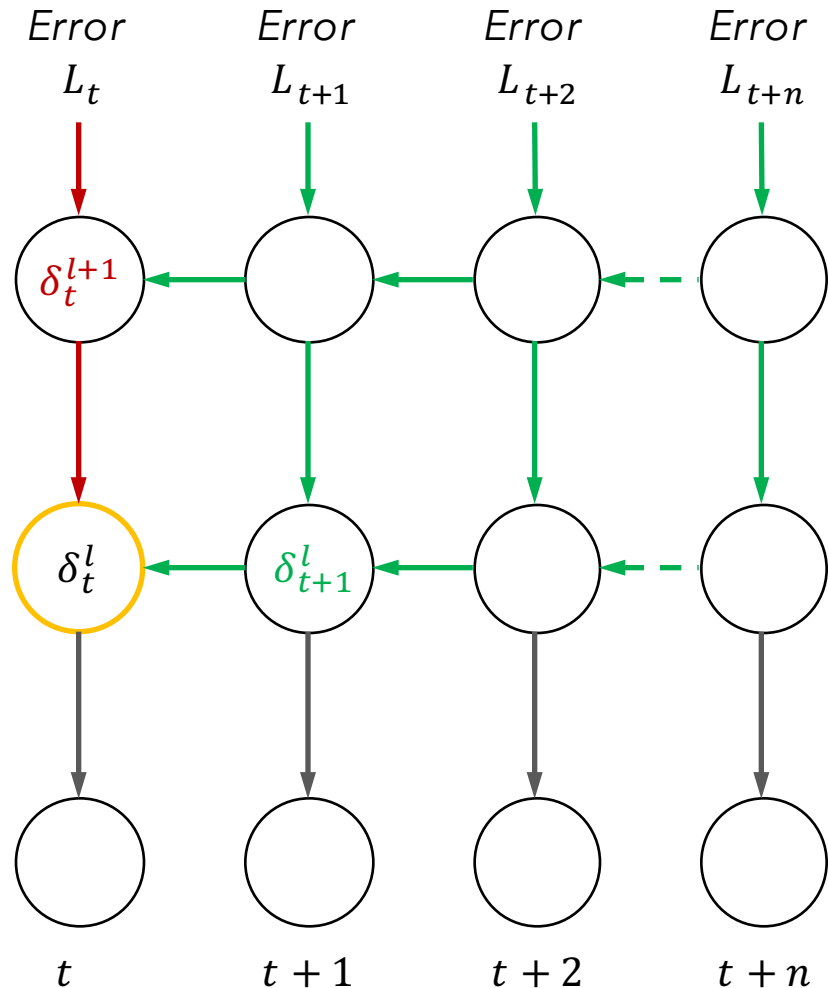


$$c_{t+1}^l = W_v x_{t+1}^l + W_h h_t^l + b$$

$$h_{t+1}^l = \text{sigmoid}(c_{t+1}^l)$$

$$L = \sum_{j=t}^T L_j = \sum_{j=t}^T (y_j - p_j)^2$$

Back-Propagation Through Time (BPTT)



$$\delta_t^l = \frac{\partial}{\partial c_t^l} \sum_{j=t}^T L_j = \underbrace{\frac{\partial L_t}{\partial c_t^l}}_{\text{Spatial term}} + \underbrace{\frac{\partial}{\partial c_t^l} \sum_{j=t+1}^T L_j}_{\text{Temporal term}}$$

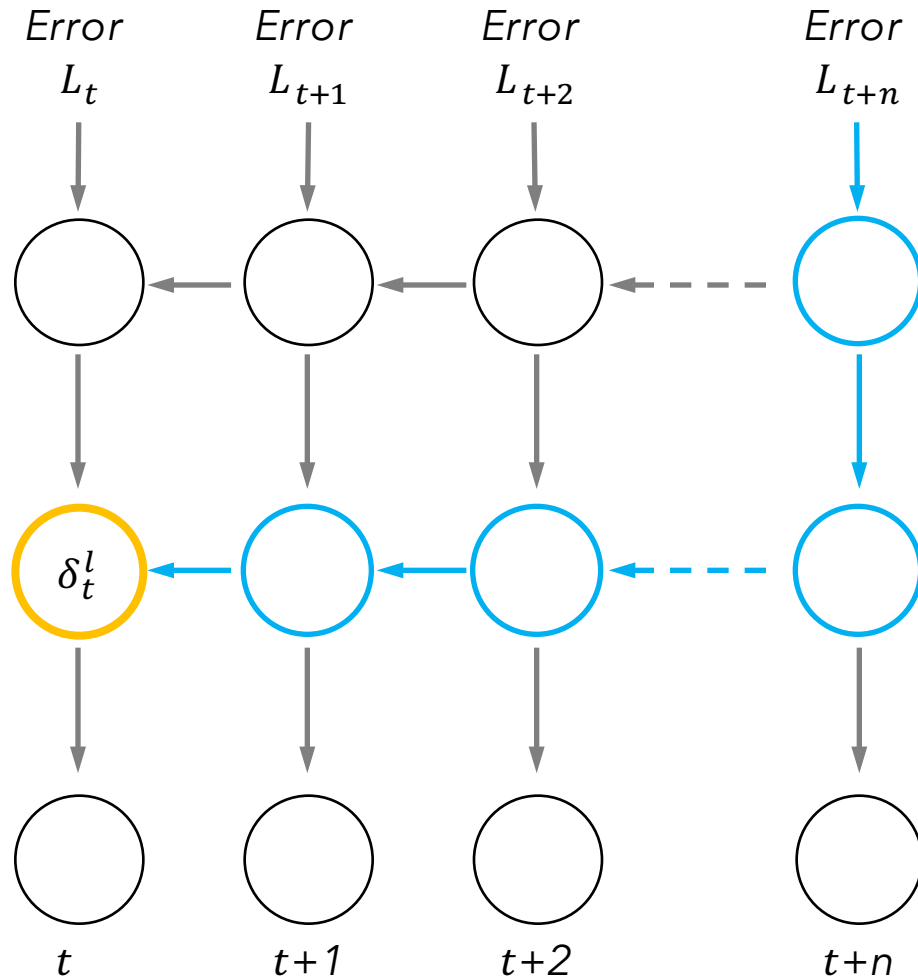
Spatial term:

$$\frac{\partial L_t}{\partial c_t^{l+1}} \cdot \frac{\partial c_t^{l+1}}{\partial c_t^l} = \delta_t^{l+1} \cdot \frac{\partial c_t^{l+1}}{\partial c_t^l}$$

Temporal term:

$$\frac{\partial \sum_{j=t+1}^T L_j}{\partial c_{t+1}^l} \cdot \frac{\partial c_{t+1}^l}{\partial c_t^l} = \delta_{t+1}^l \cdot \frac{\partial c_{t+1}^l}{\partial c_t^l}$$

Vanishing Gradient Problem



$$\frac{\partial L_{t+n}}{\partial c_t^l} = \frac{\partial L_{t+n}}{\partial c_{t+n}^l} \cdot \frac{\partial c_{t+n}^l}{\partial c_{t+n-1}^l} \cdot \dots \cdot \frac{\partial c_{t+1}^l}{\partial c_t^l}$$

$$= \frac{\partial L(t)}{\partial c_{t+n}^l} \cdot \underbrace{\prod_{\tau=t}^{t+n} \frac{\partial c_{\tau+1}^l}{\partial c_{\tau}^l}}_{\text{Sequential Jacobian}}$$

Sequential Jacobian

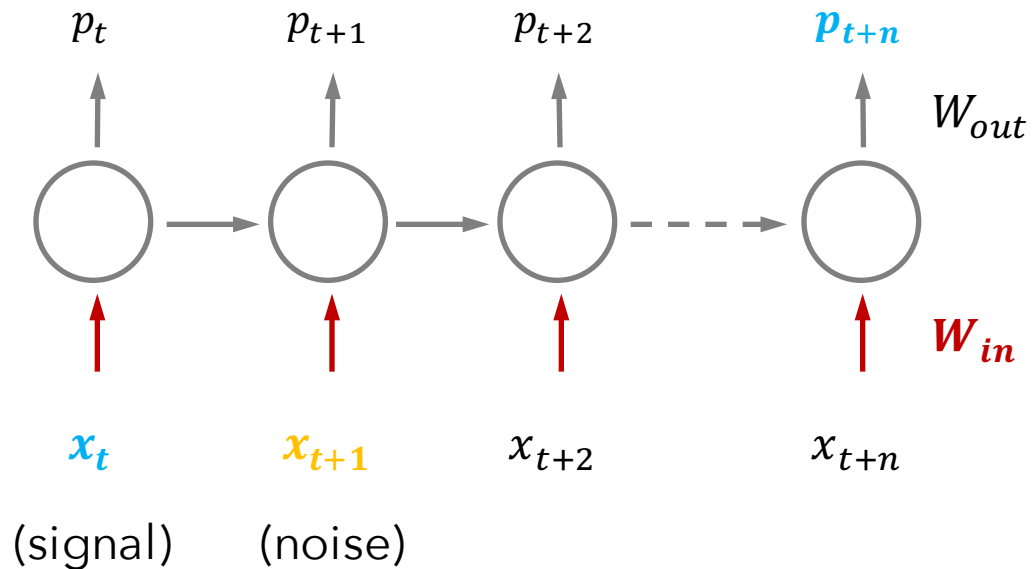
$$\frac{\partial c_{\tau+1}^l}{\partial c_{\tau}^l} = W_h^T \sigma'(c_{\tau}^l) \Rightarrow \left\| \frac{\partial c_{\tau+1}^l}{\partial c_{\tau}^l} \right\| \leq \|W_h\| \underbrace{\|\sigma'(c_{\tau}^l)\|}_{\leq 1/4}$$

- **Exponential** decayed error message
- **Long-term dependency** cannot be learned

Weight Conflict Problem

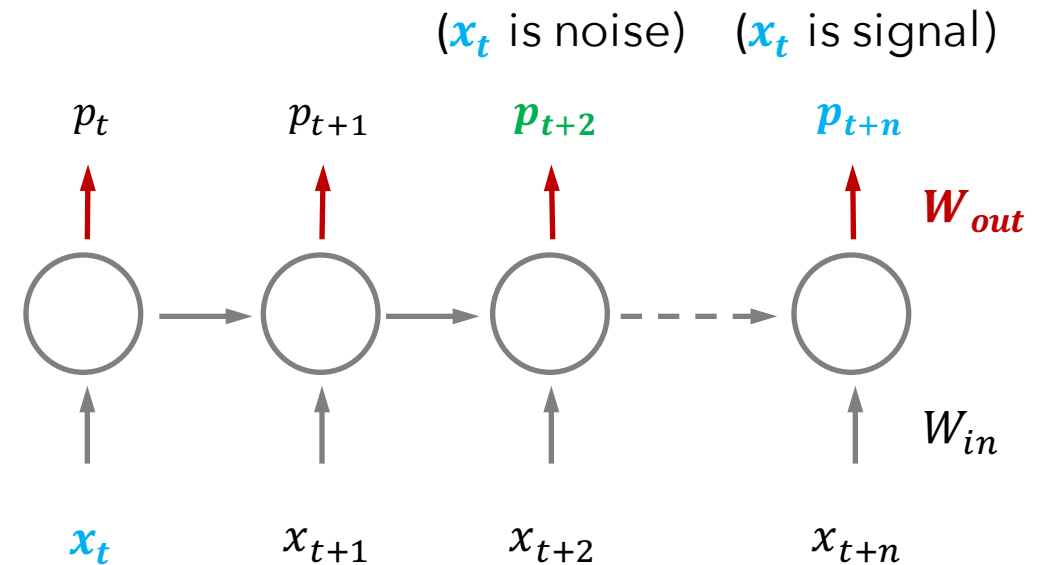
Two **conflict** roles of W_{in}

- **Absorb** useful signal x_t
- **Reject** harmful noise x_{t+1}



Two **conflict** roles of W_{out}

- **Reject** useless memory of x_t for p_{t+2}
- **Retrieve** useful memory of x_t for p_{t+n}



Treating Vanishing Gradient: Constant Error Carrousel (CEC)

Error signal doesn't vanish $\xrightarrow{\text{i.e.}}$ $\frac{\partial c_{\tau+1}^l}{\partial c_{\tau}^l} = W_h^T \sigma'(c_{\tau}^l) \approx \mathbf{I}$ \rightarrow **CEC**

$\xrightarrow{\text{e.g.}}$ Let $W = \mathbf{I}, f(c_{\tau}^l) = c_{\tau}^l$
Then, $W_h^T \sigma'(c_{\tau}^l) = \mathbf{I}$

Problem with this idea

- No **non-linearity** (network won't be powerful)

Treating Weight Conflict: Gating Function

Core Idea



Learn

1. what to store in the memory
2. what to retrieve from the memory

Gated Input

$$in_t = f_t^{in} \otimes x_t$$

- f_t^{in} is the input gating function
- $[f_t^{in}]_i \in [0, 1]$ (each element within $[0, 1]$)

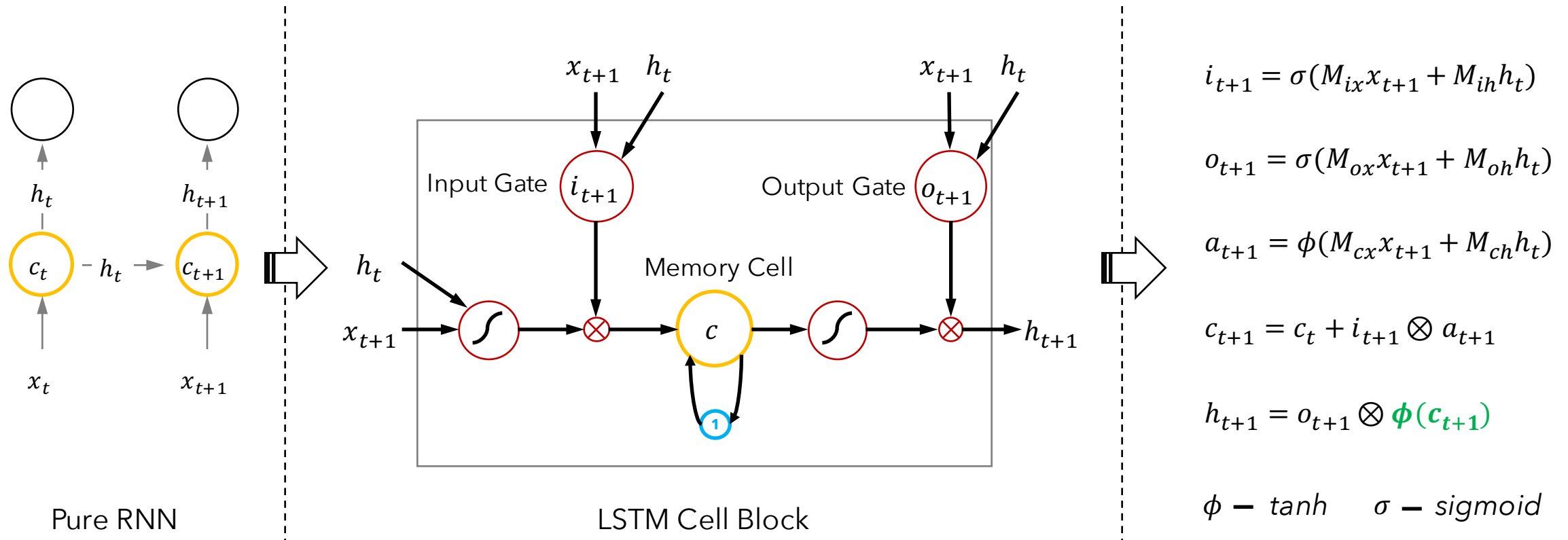
Gated Output

$$out_t = f_t^{out} \otimes c_t$$

- f_t^{out} is the output gating function
- $[f_t^{out}]_i \in [0, 1]$ (each element within $[0, 1]$)

($\otimes \rightarrow$ element-wise multiplication)

CEC + Gates \rightarrow Long-Short Term Memory (LSTM)



- 1 Constant Error Carrousel (CEC)
- ⊗ Element-wise multiplication
- S Non-linearity
- \$i_{t+1}\$ Input gate
- \$o_{t+1}\$ Output gate

Why LSTM solves the problem

$$i_{t+1} = \sigma(M_{ix}x_{t+1} + M_{ih}h_t)$$

$$o_{t+1} = \sigma(M_{ox}x_{t+1} + M_{oh}h_t)$$

$$a_{t+1} = \phi(M_{cx}x_{t+1} + M_{ch}h_t)$$

$$c_{t+1} = c_t + i_{t+1} \otimes a_{t+1}$$

$$h_{t+1} = o_{t+1} \otimes \phi(c_{t+1})$$



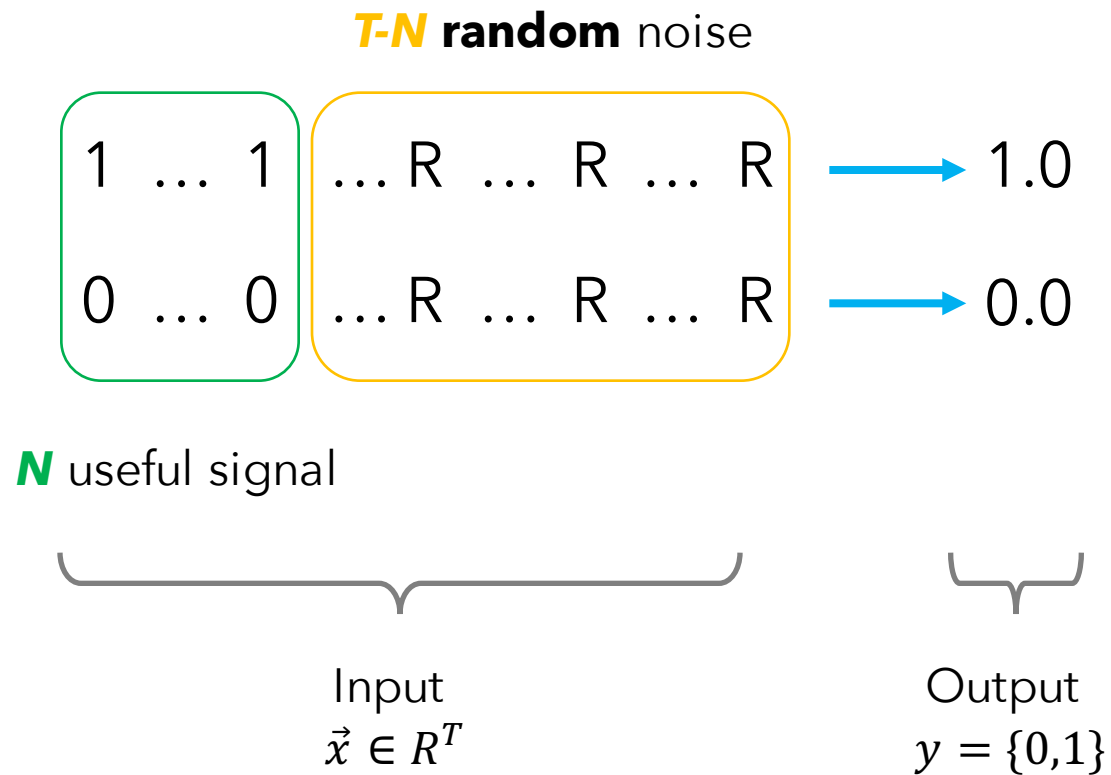
$$\frac{\partial c_{\tau+1}^l}{\partial c_{\tau}^l} = I + \frac{\partial c_{\tau+1}^l}{\partial a_{\tau+1}^l} \cdot \frac{\partial a_{\tau+1}^l}{\partial h_{\tau}^l} \cdot \frac{\partial h_{\tau}^l}{\partial c_{\tau}^l} = \overbrace{I}^{\text{CEC}} + \underbrace{i_{t+1} \dots f_t \dots}_{\text{Gated Error}}$$



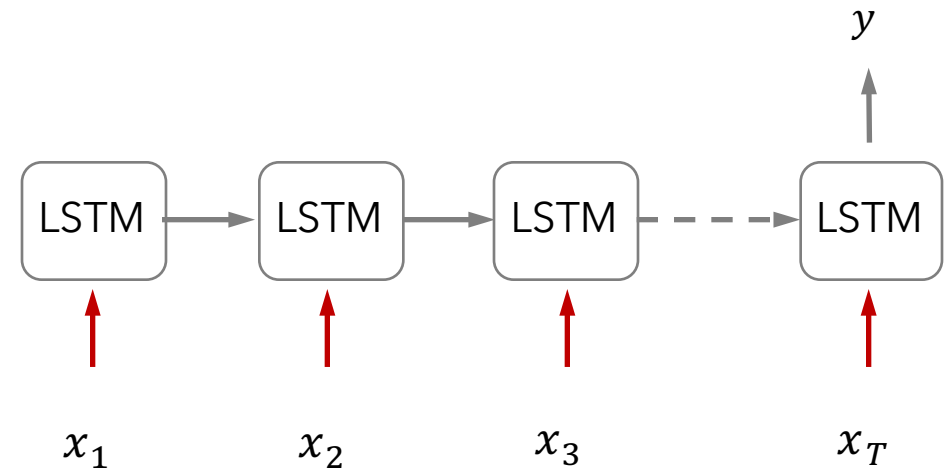
$$\frac{\partial L_{t+n}}{\partial c_t^l} = \frac{\partial L_{t+n}}{\partial c_{t+n}^l} \prod_{\tau=t}^{t+n} \frac{\partial c_{\tau+1}^l}{\partial c_{\tau}^l} = \frac{\partial L_{t+n}}{\partial c_{t+n}^l} \prod_{j=1}^n (I + i_{t+1} \dots f_t \dots)$$

Experiment 3: two-sequence problem

Problem



Model



Experiment 4 & 5: adding/multiplication problem

Second dimension used as a marker

Adding Problem:

$$\dots \begin{bmatrix} X_1 \\ -1 \end{bmatrix} \dots \begin{bmatrix} R \\ 1 \end{bmatrix} \dots \begin{bmatrix} X_2 \\ -1 \end{bmatrix} \dots \begin{bmatrix} R \\ 1 \end{bmatrix} \dots \longrightarrow 0.5 + \frac{X_1 + X_2}{4.0}$$

Multiplication Problem:

$$\dots \begin{bmatrix} X_1 \\ -1 \end{bmatrix} \dots \begin{bmatrix} R \\ 1 \end{bmatrix} \dots \begin{bmatrix} X_2 \\ -1 \end{bmatrix} \dots \begin{bmatrix} R \\ 1 \end{bmatrix} \dots \longrightarrow X_1 \times X_2$$

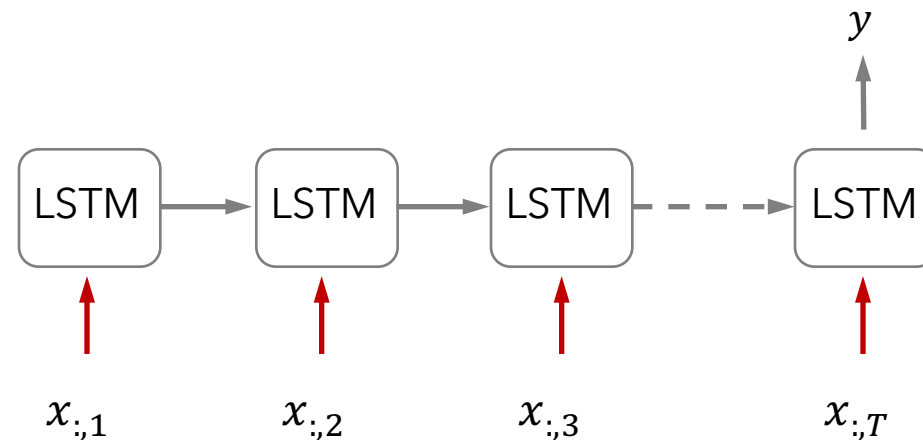


Input: $X \in R^{2 \times T}$



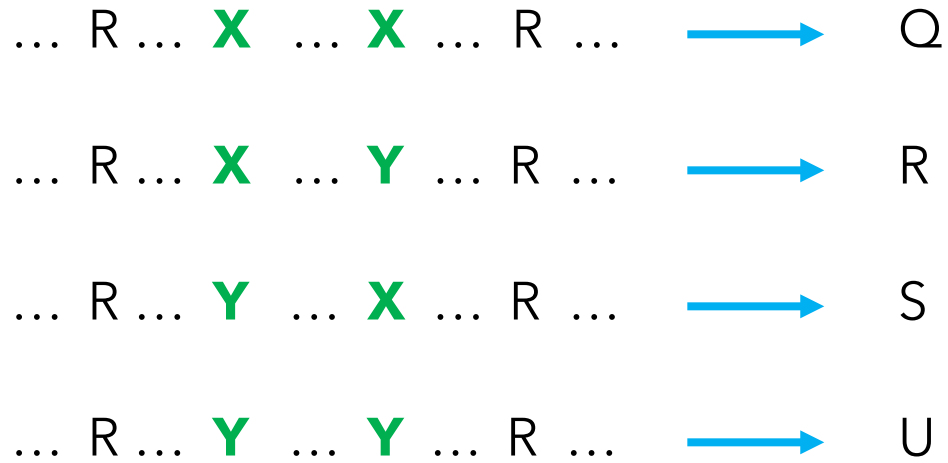
Output: $y \in R$

Model



Experiment 6: temporal order problem

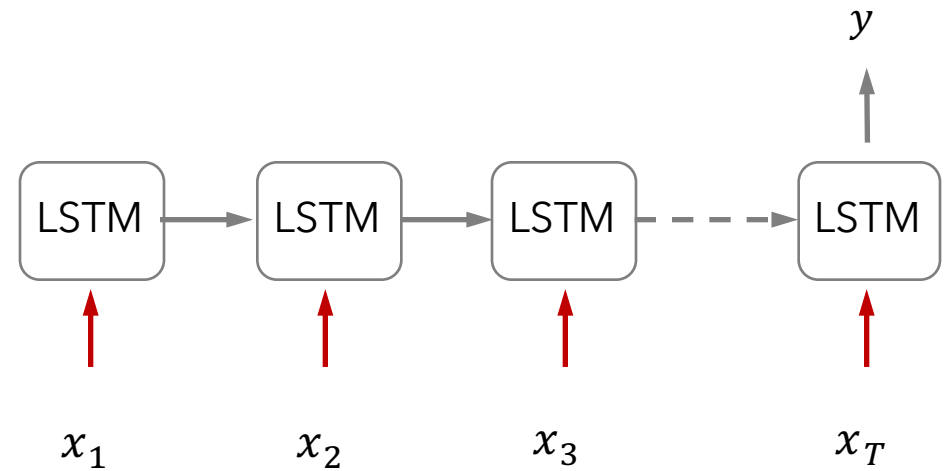
Problem



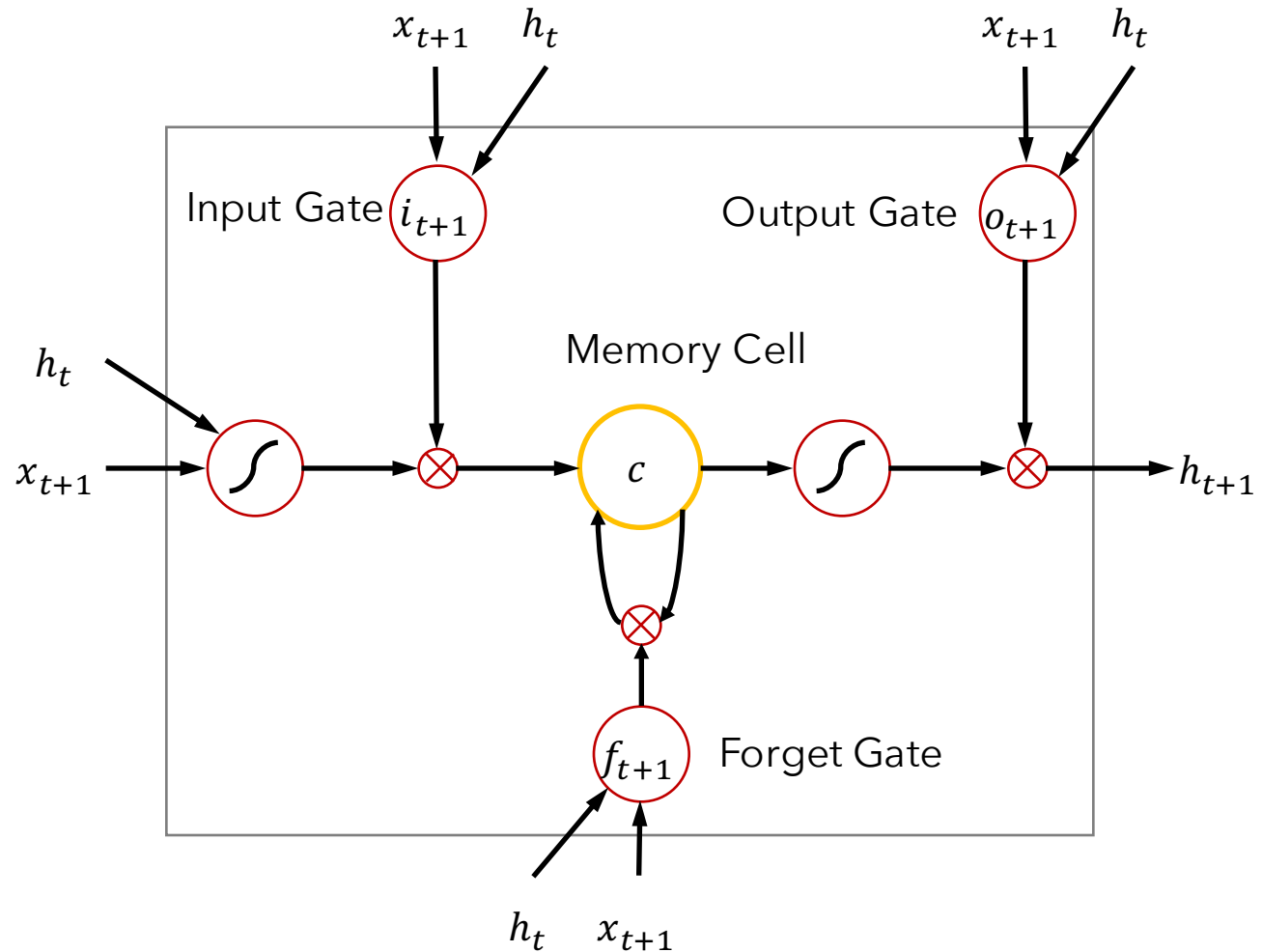
Input
 $\vec{x} \in R^T$

Output
 $y = \{Q, R, S, U\}$

Model



Introducing Forget Gate



$$i_{t+1} = \sigma(M_{ix}x_{t+1} + M_{ih}h_t)$$

$$f_{t+1} = \sigma(M_{fx}x_{t+1} + M_{fh}h_t)$$

$$o_{t+1} = \sigma(M_{ox}x_{t+1} + M_{oh}h_t)$$

$$a_{t+1} = \phi(M_{cx}x_{t+1} + M_{ch}h_t)$$

$$c_{t+1} = f_{t+1} \otimes c_t + i_{t+1} \otimes a_{t+1}$$

$$h_{t+1} = o_{t+1} \otimes \phi(c_{t+1})$$

Thanks & Questions