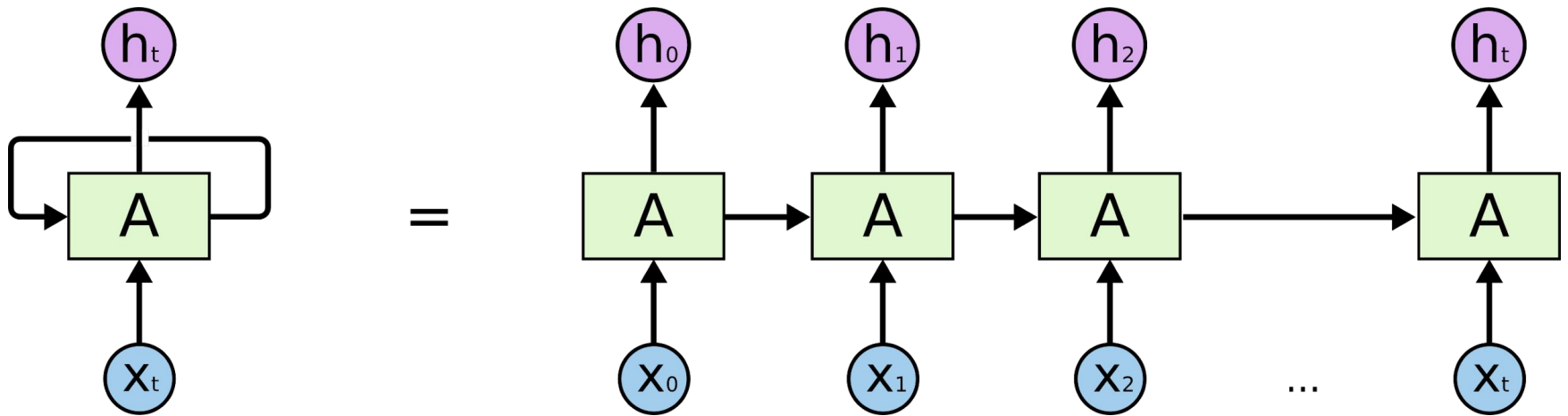


A decorative header at the top of the slide consisting of a series of overlapping triangles in various shades of blue, teal, and orange, creating a colorful, abstract border.

Simple RNN

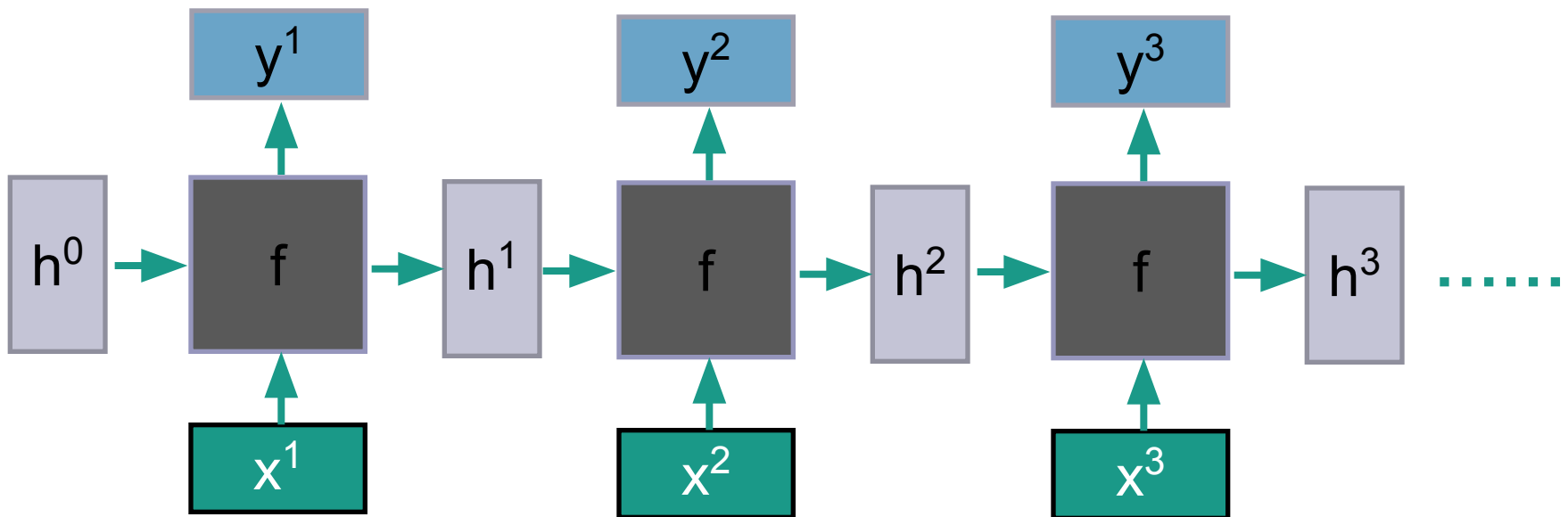
Recurrent Neural Network



How does RNN reduce complexity?

- Given function $f: h', y = f(h, x)$

h and h' are vectors with the same dimension

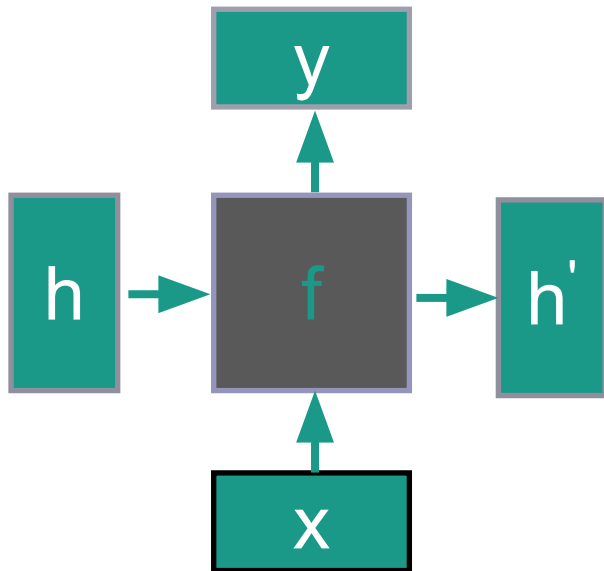


$x^n = \text{vector}$

Only need one function f .

Naïve RNN

- Given function $f: h', y = f(h, x)$



$$h' = \sigma(W^h h + W^i x)$$

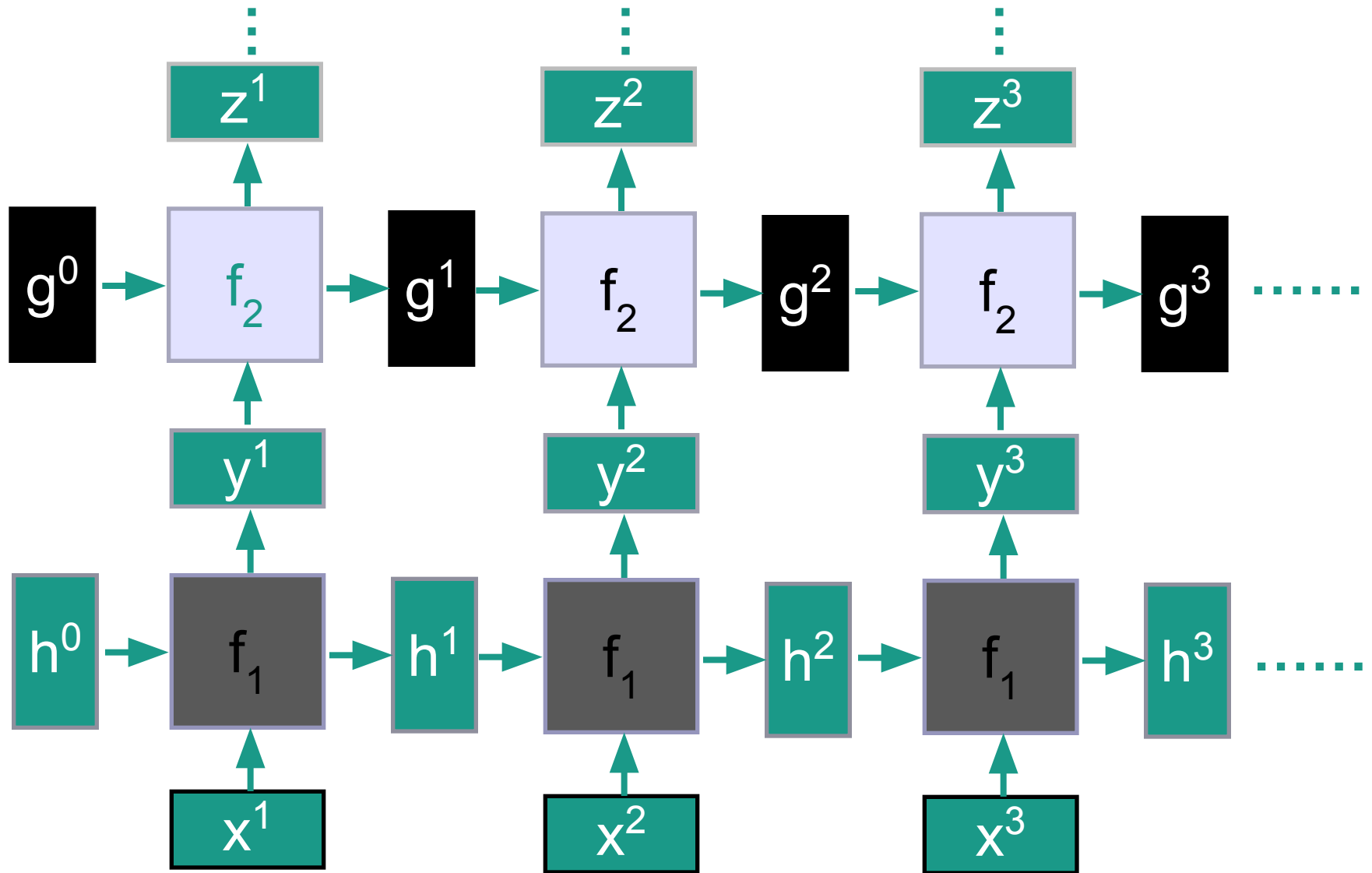
$$y = \sigma(W^o h')$$

Note: y is computed from h'

We have ignored the bias

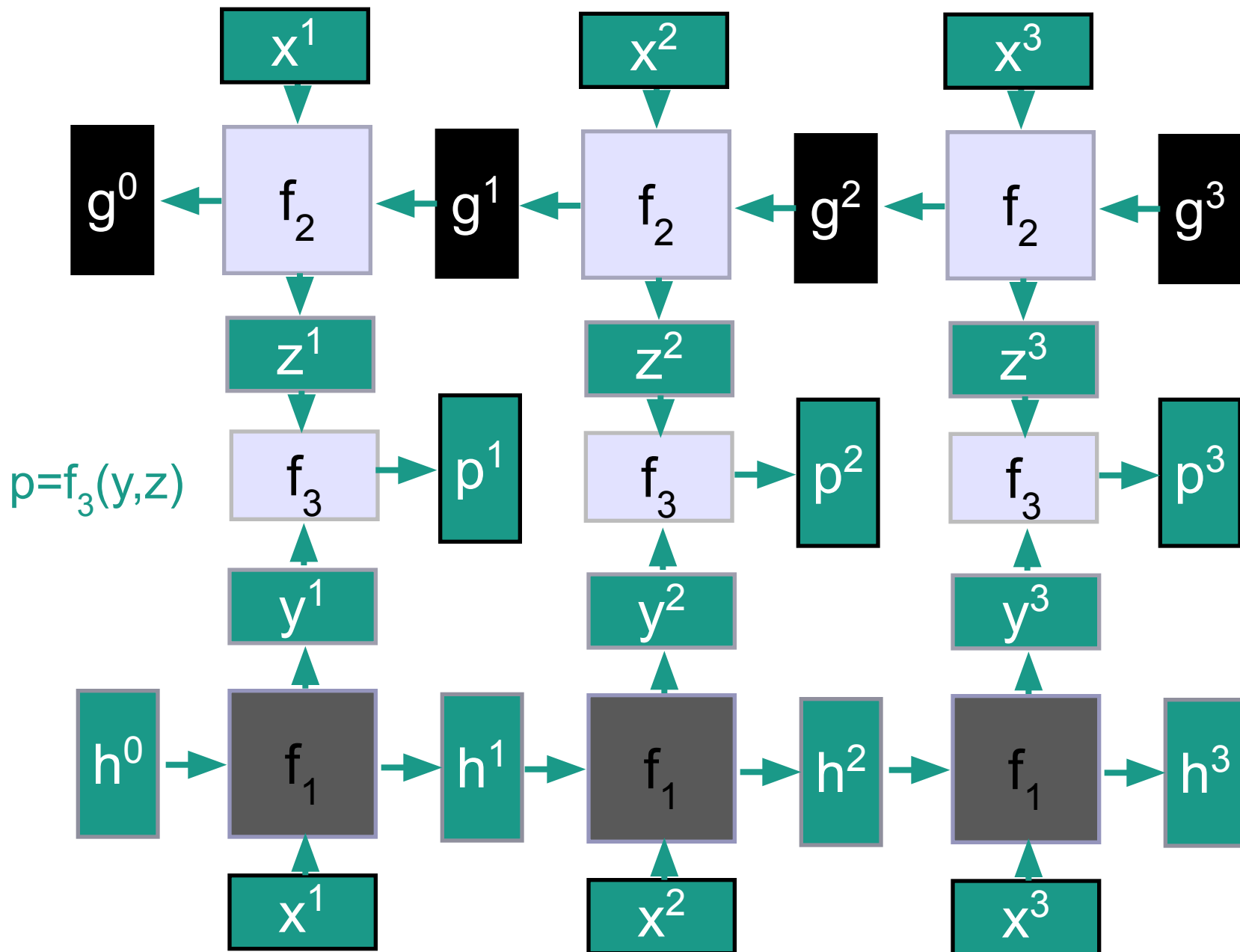
Deep RNN

$$h', y = f_1(h, x), \quad g', z = f_2(g, y) \quad \dots$$



Bidirectional RNN

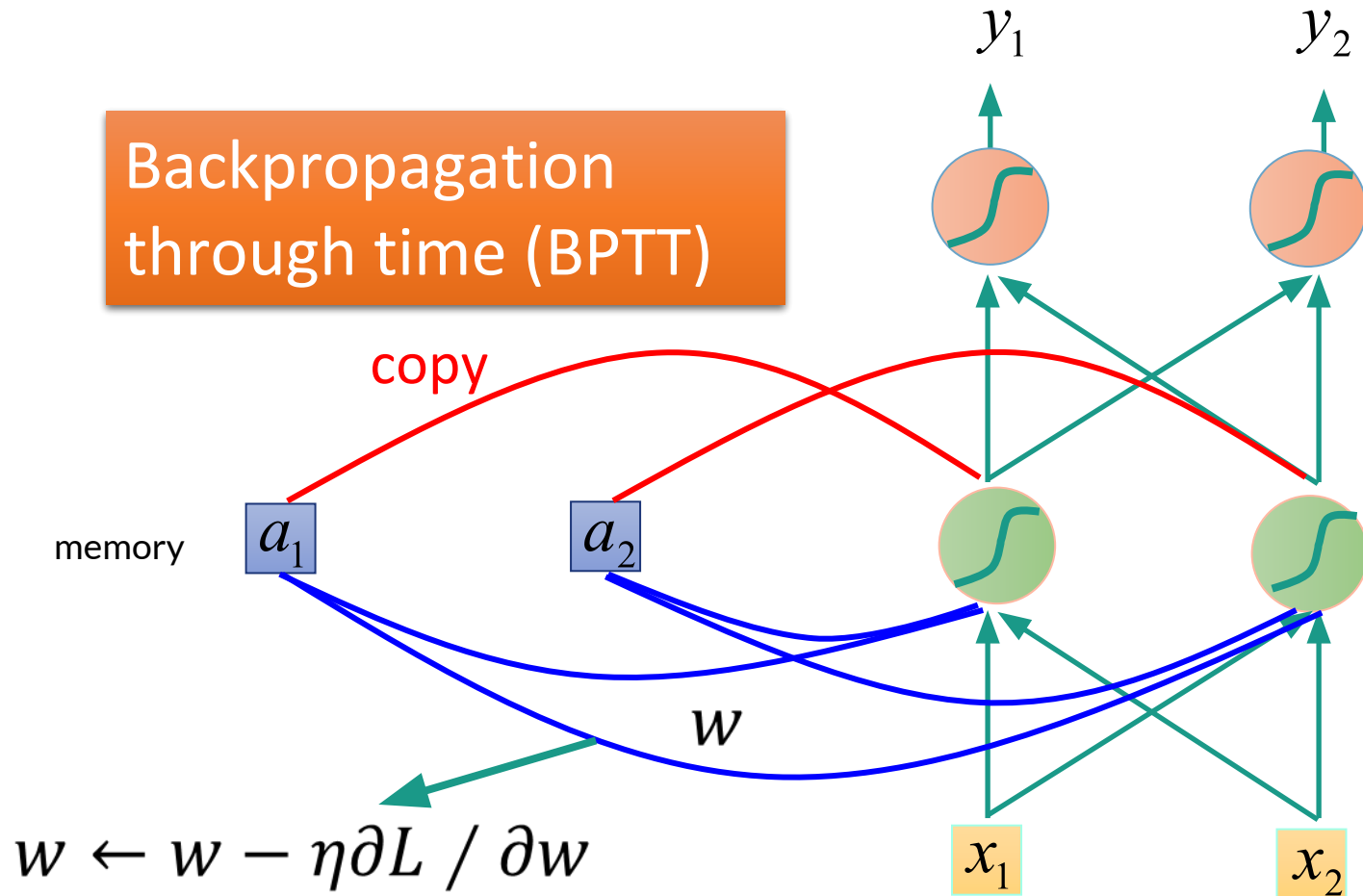
$$y, h = f_1(x, h) \quad z, g = f_2(g, x)$$



Problems with naive RNN

- “memory” too short.
- Gradient vanishing
- Gradient exploding

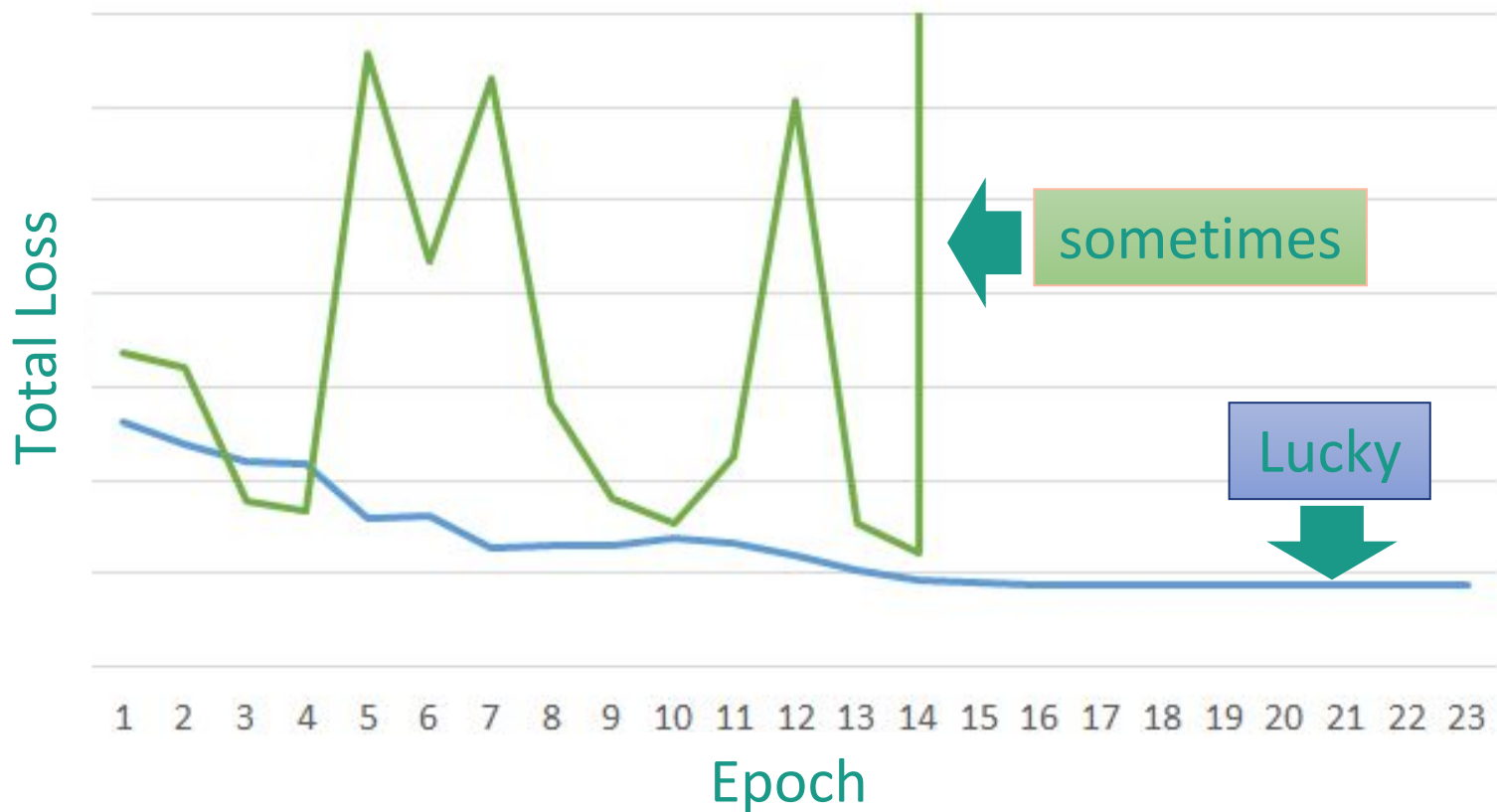
How to Training RNN



RNN-based network

- RNN-based network is not always easy to learn

Real experiments on Language modeling

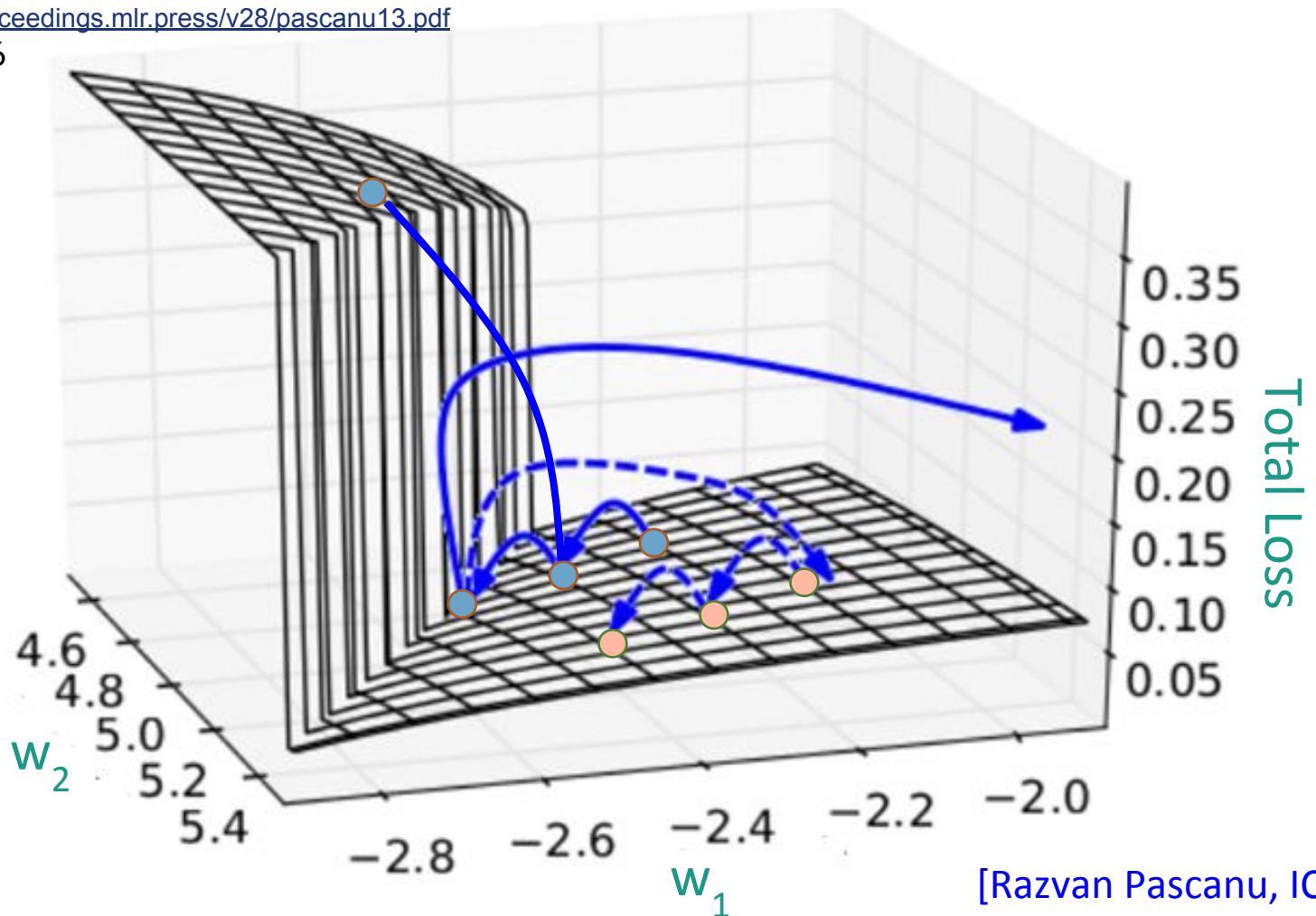


The error surface is rough.

$$w \leftarrow w - \eta \partial L / \partial w$$

<http://proceedings.mlr.press/v28/pascanu13.pdf>

figure.6



[Razvan Pascanu, ICML'13]

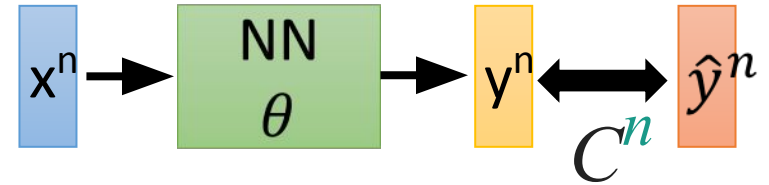
Update Weights

$$w = 1 \quad \longrightarrow \quad y^{1000} = 1$$

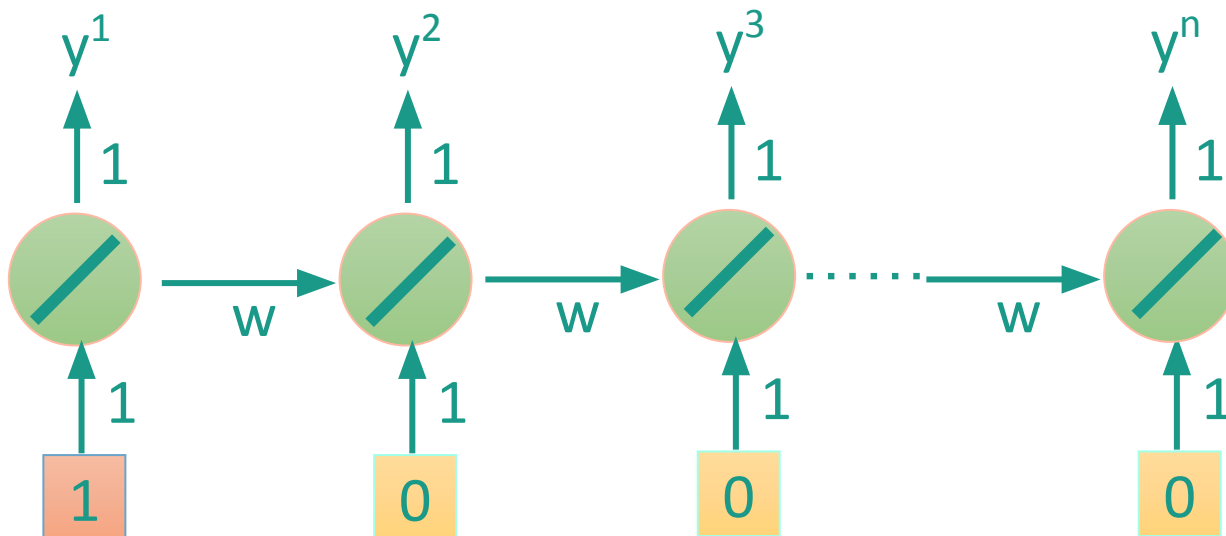
$$w = 1.01 \quad \longrightarrow \quad y^{1000} \approx 20000$$

$$w = 0.99 \quad \longrightarrow \quad y^{1000} \approx 0$$

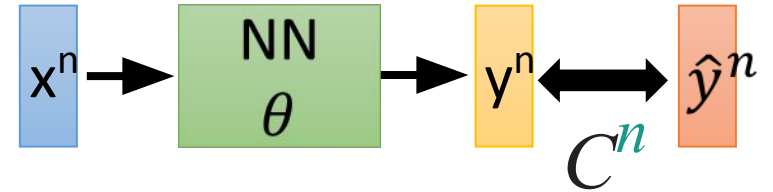
$$w = 0.01 \quad \longrightarrow \quad y^{1000} \approx 0$$



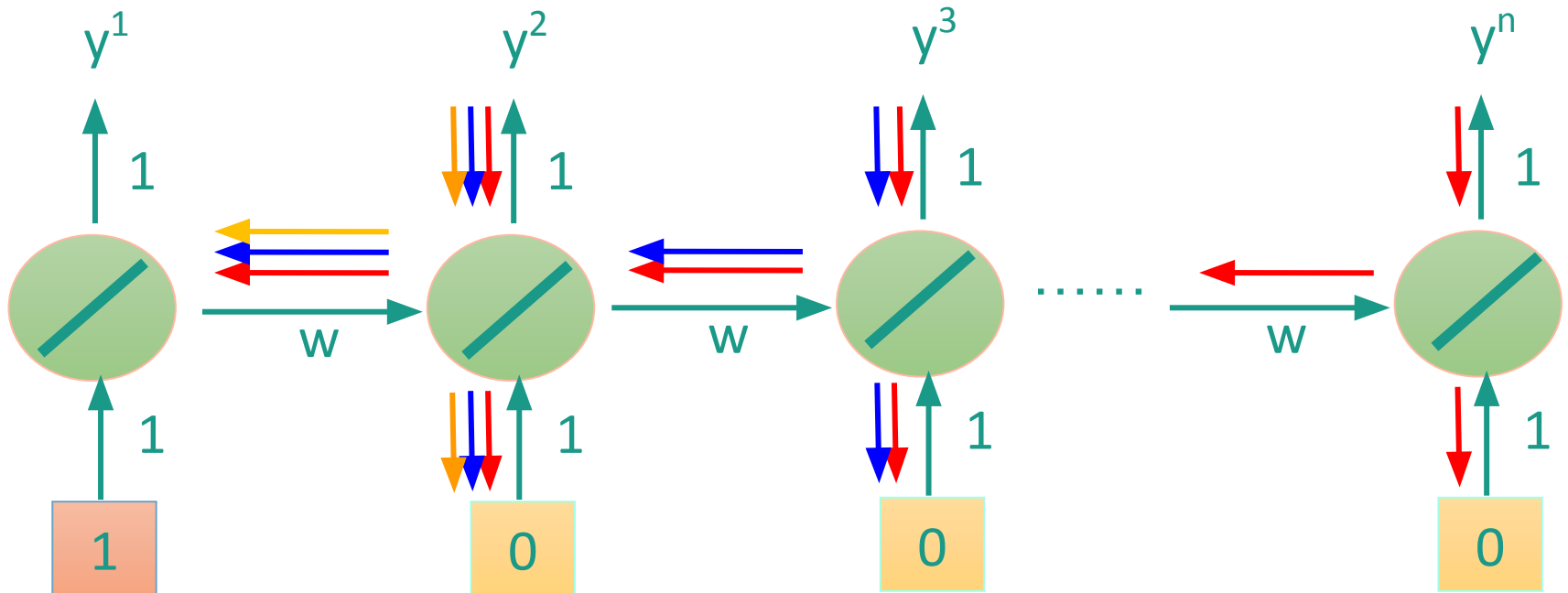
$$y^n = 1 * W^{n-1}$$



Backpropagation



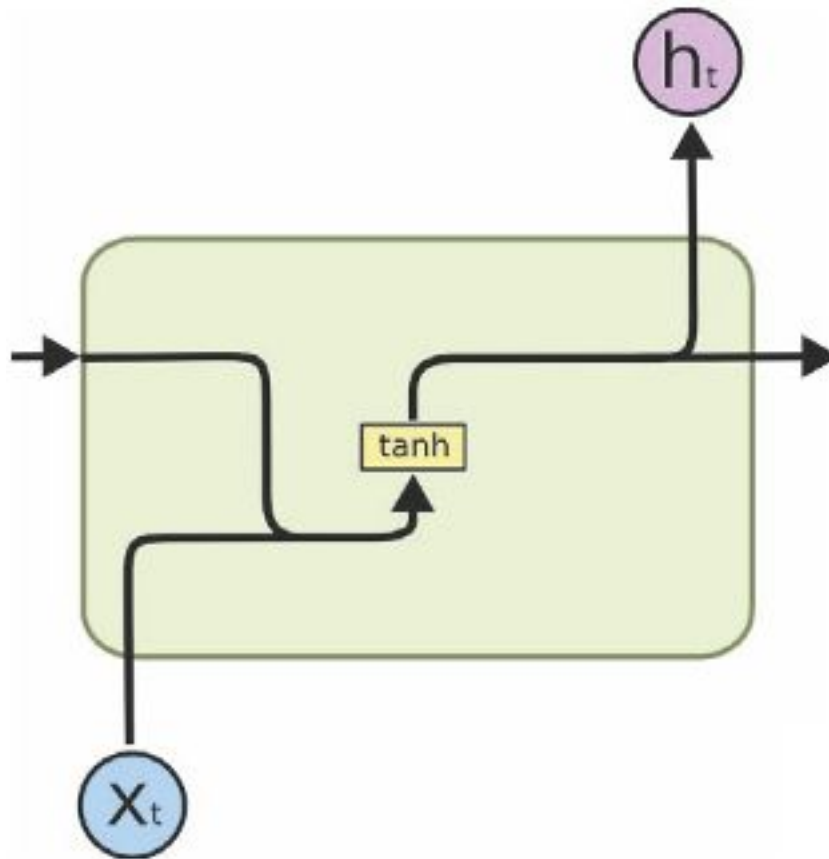
$$\frac{\partial C^n}{\partial W} = \boxed{\frac{\partial C^n}{\partial y^n}} \frac{\partial y^n}{\partial W}$$



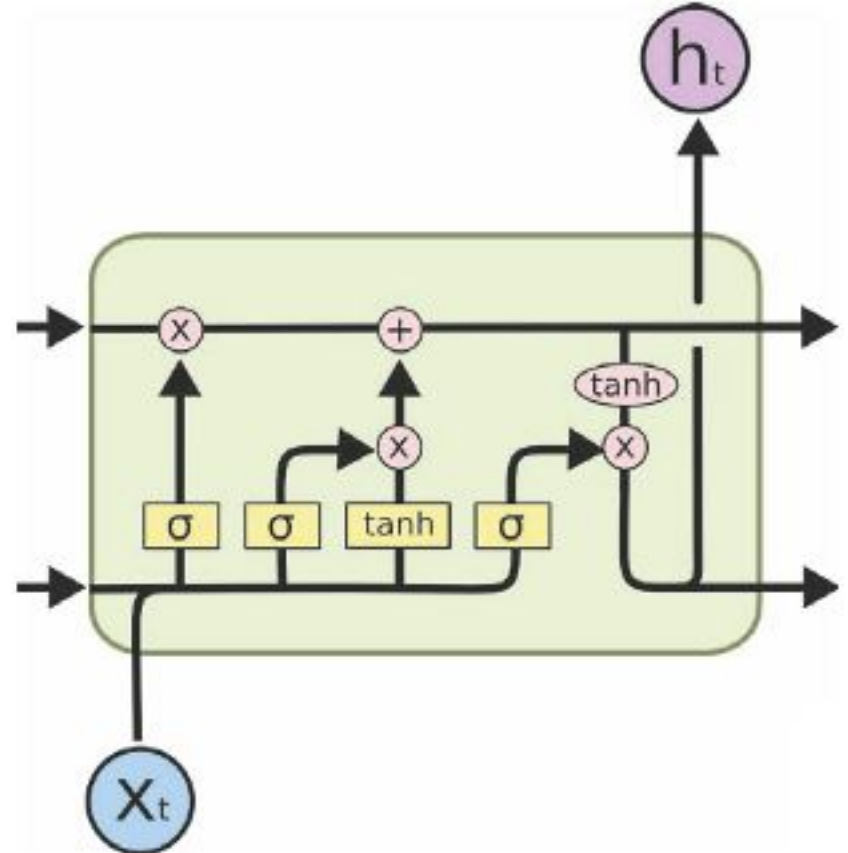
A decorative header consisting of a series of overlapping triangles in various shades of blue, teal, and orange, creating a jagged, mountain-like silhouette.

LSTM

RNN vs LSTM

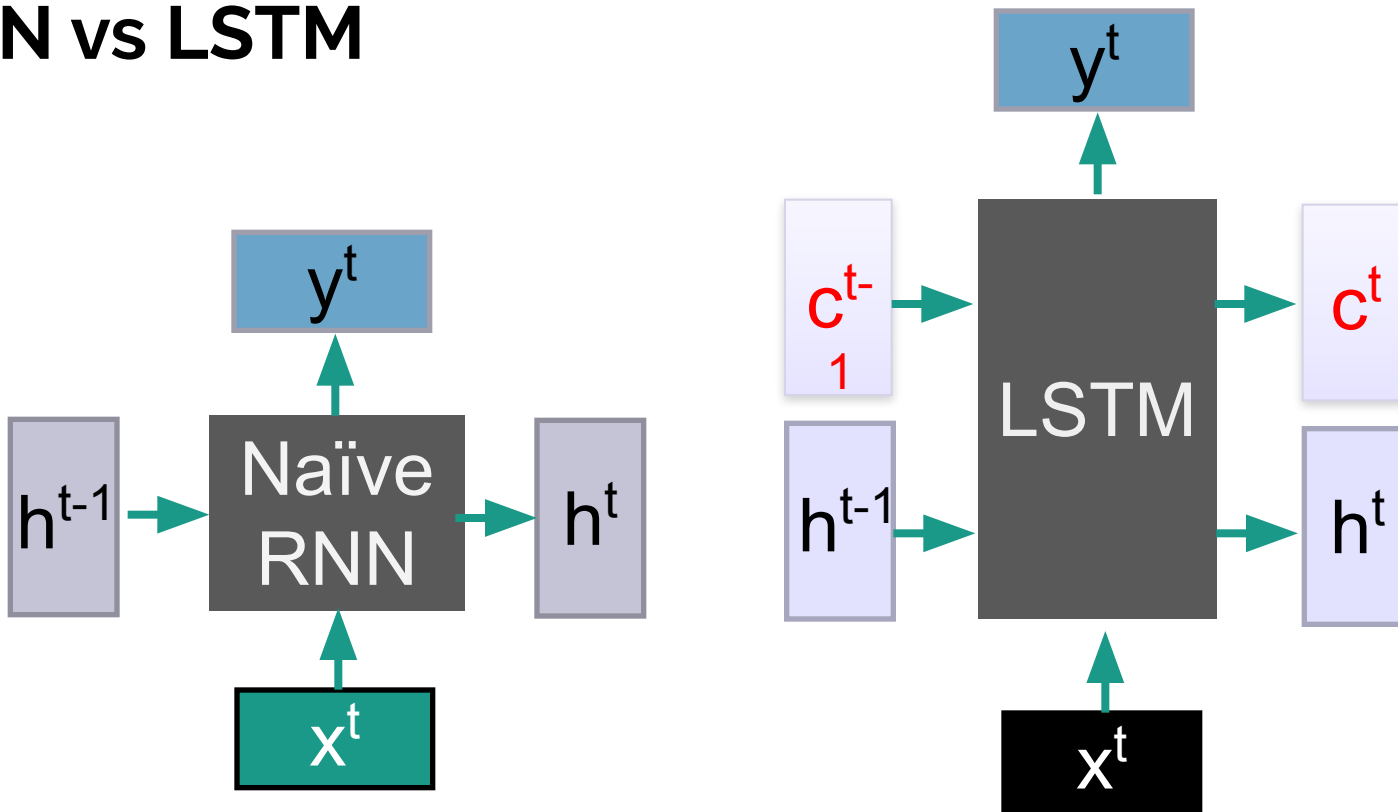


(a) RNN



(b) LSTM

RNN vs LSTM

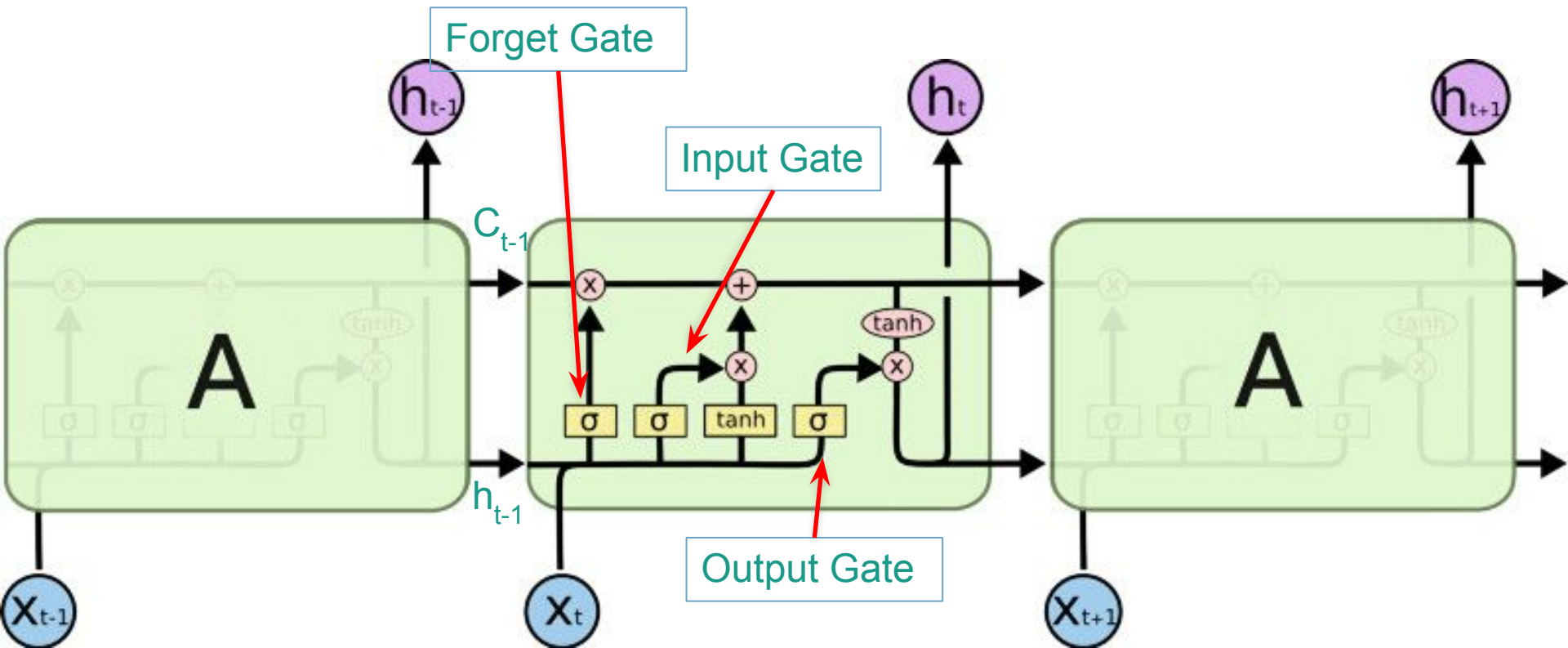


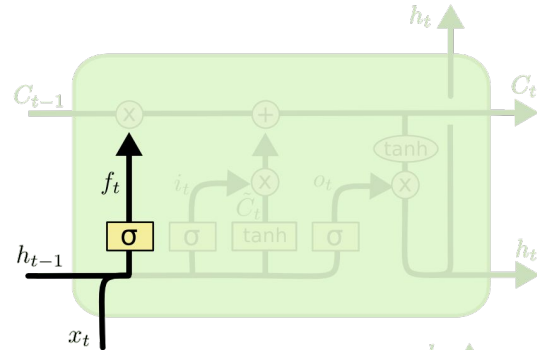
c changes slowly $\longrightarrow c^t$ is c^{t-1} added by something

h changes faster $\longrightarrow h^t$ and h^{t-1} can be very different

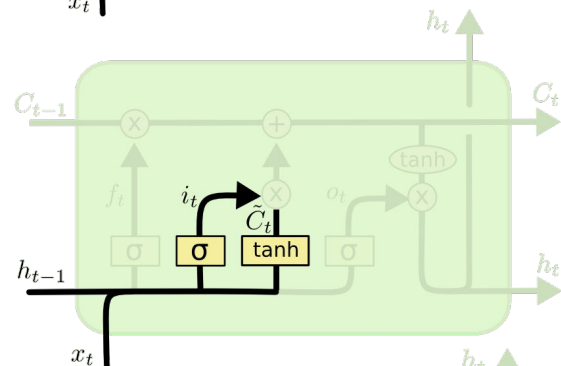
$$h_t = o_t * \tanh(C_t)$$

LSTM





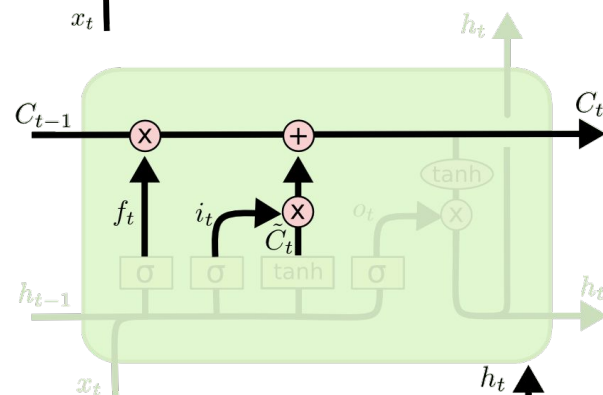
$$f_t = \sigma(W_f \cdot [h_{t-1}, x_t] + b_f)$$



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

$$\tilde{C}_t = \tanh(W_C \cdot [h_{t-1}, x_t] + b_C)$$

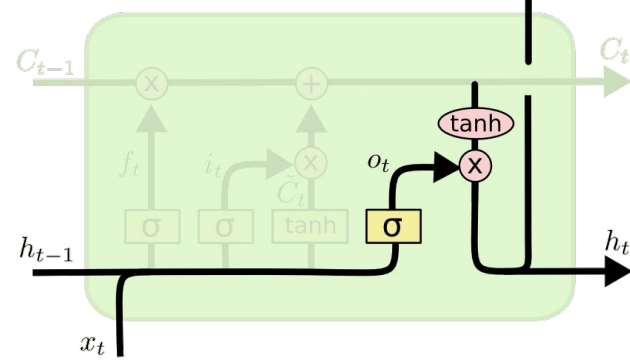
i_t decides what component is to be updated.
 C_t provides change contents



$$C_t = f_t * C_{t-1} + i_t * \tilde{C}_t$$

Updating the cell state

C_t : Cell's memory changed slowly

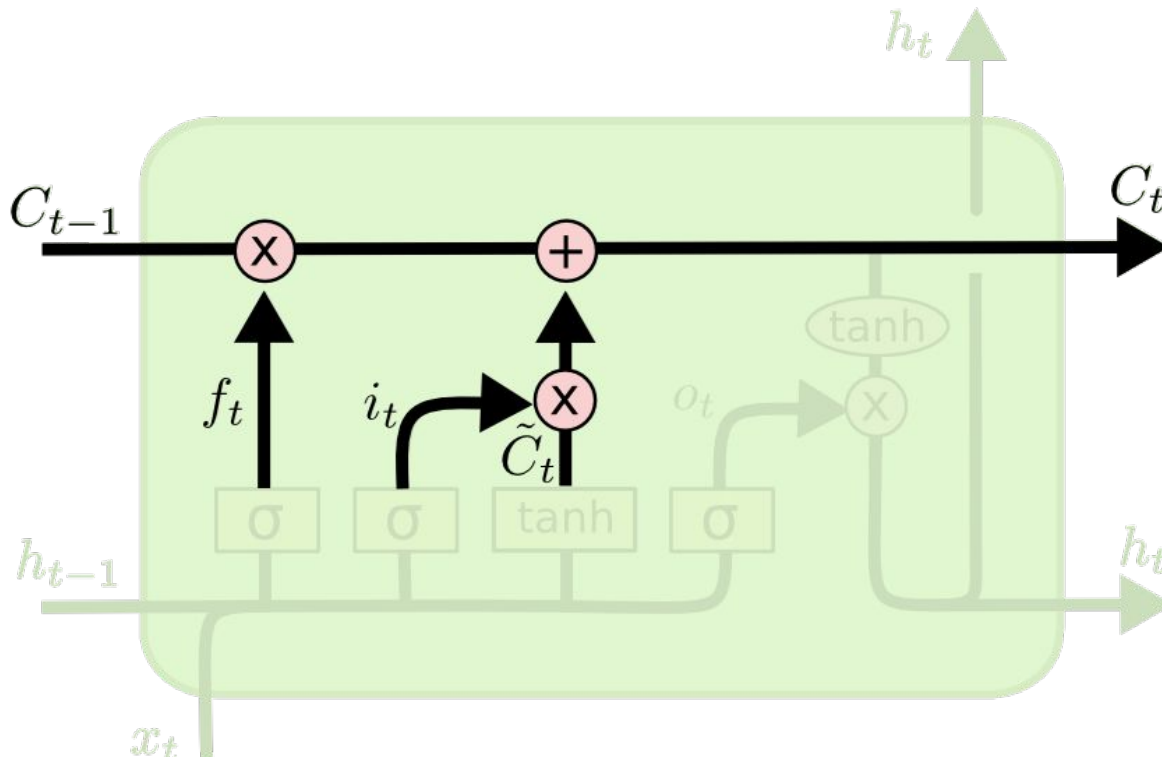


$$o_t = \sigma(W_o [h_{t-1}, x_t] + b_o)$$

$$h_t = o_t * \tanh(C_t)$$

Decide what part of the cell state to output

Cell's memory

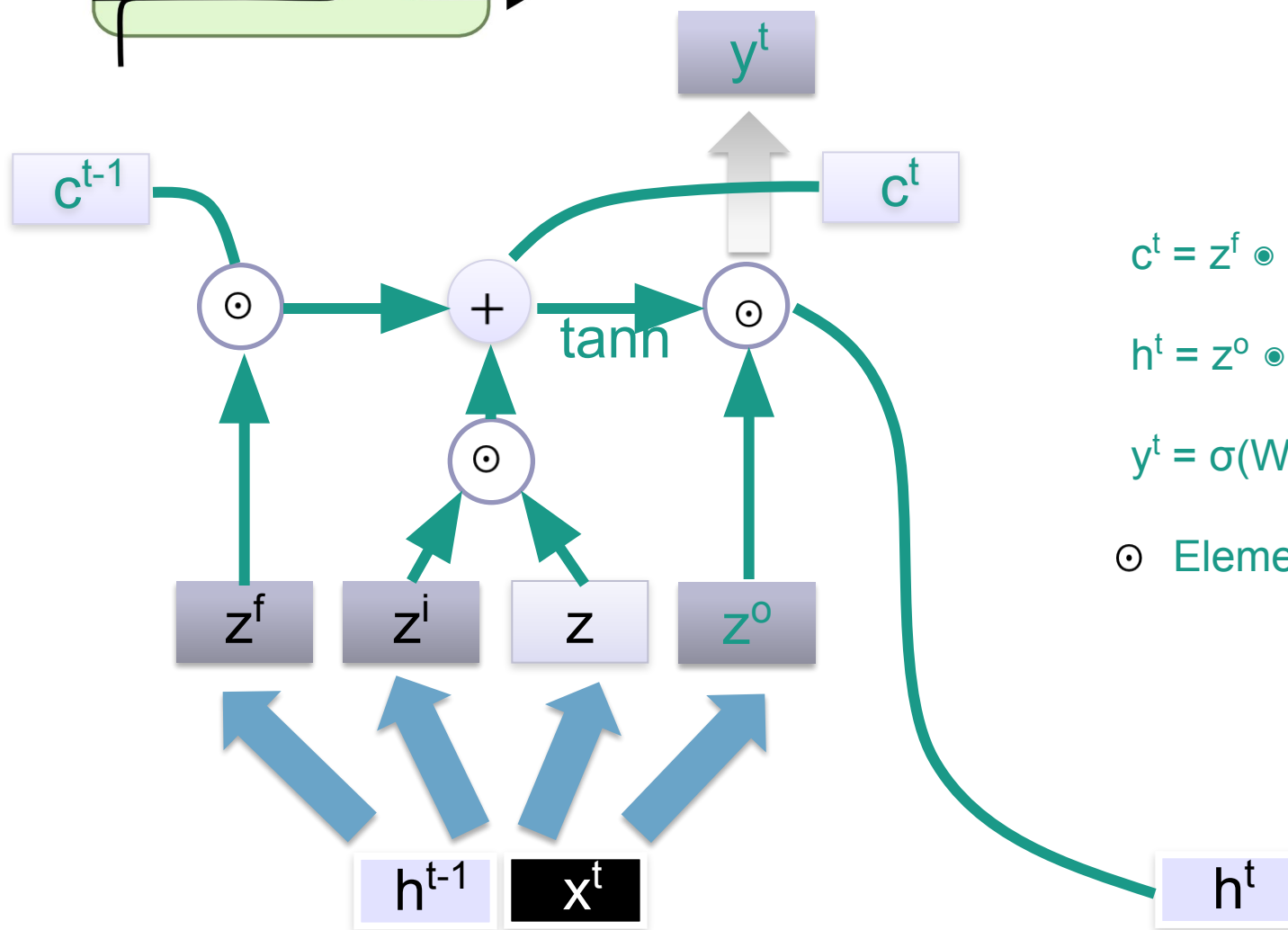
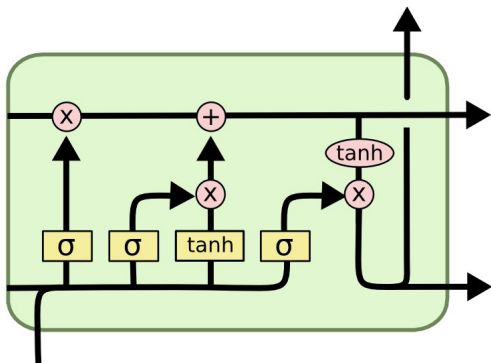


Forget Gate

if $1 > \text{sigmoid} > 0$, remember C^{t-1}
 if $\text{sigmoid} = 0$, forget C^{t-1}

Input Gate

if $1 > \text{sigmoid} > 0$, remember \tilde{C}_t
 if $\text{sigmoid} = 0$, forget \tilde{C}_t

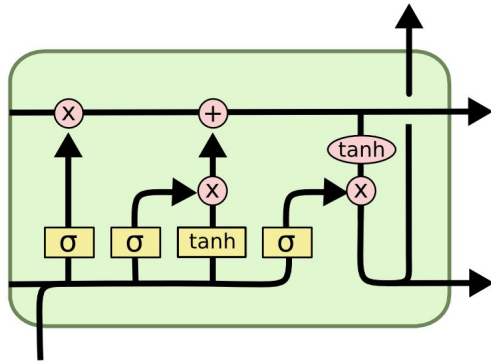


$$c^t = z^f \odot c^{t-1} + z^i \odot z$$

$$h^t = z^o \odot \tanh(c^t)$$

$$y^t = \sigma(W' h^t)$$

\odot Element-wise multiply



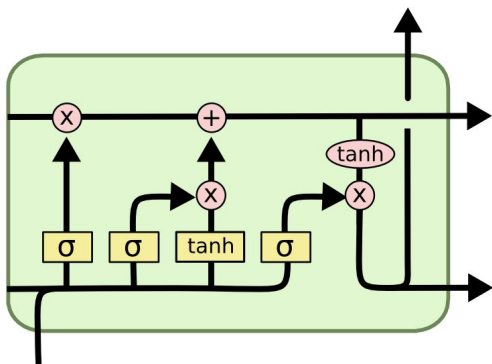
$$z = \left(\text{Transform} \begin{pmatrix} x^t \\ h^{t-1} \end{pmatrix} \right)$$

Transform = matrix

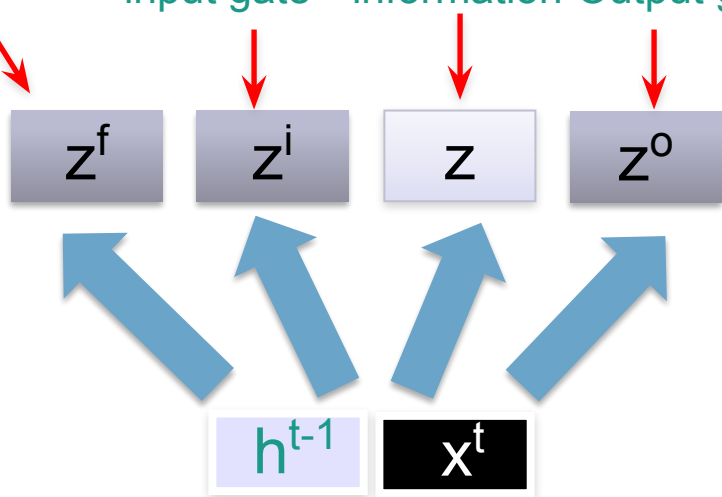
z^0 z^f z^i obtained by the same way



← Vector



Controls forget gate Controls input gate Updating information Controls Output gate



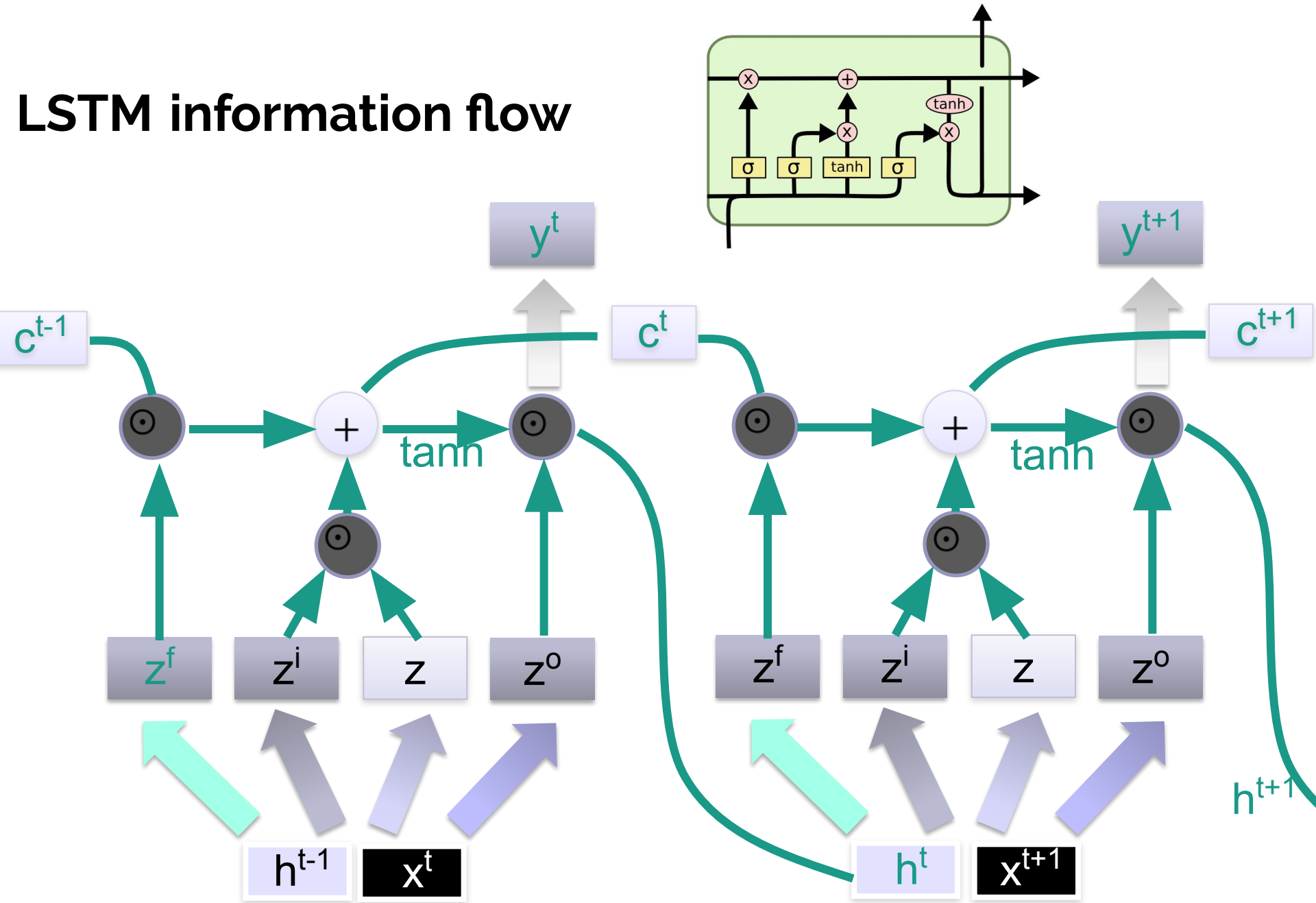
$$z = \tanh(W \begin{bmatrix} x^t \\ h^{t-1} \end{bmatrix})$$

$$z^i = \sigma(W^i \begin{bmatrix} x^t \\ h^{t-1} \end{bmatrix})$$

$$z^f = \sigma(W^f \begin{bmatrix} x^t \\ h^{t-1} \end{bmatrix})$$

$$z^o = \sigma(W^o \begin{bmatrix} x^t \\ h^{t-1} \end{bmatrix})$$

LSTM information flow



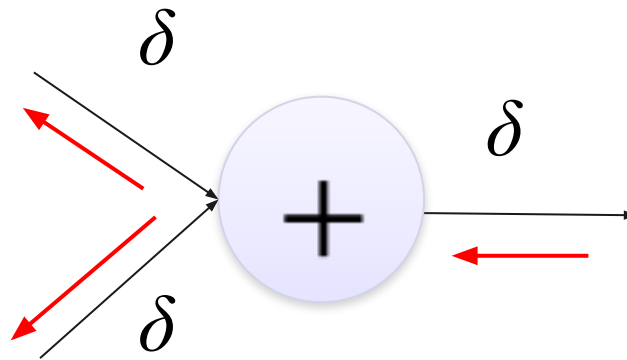
Constant Error Carousels

Gradient Vanishing / Exploding in RNN

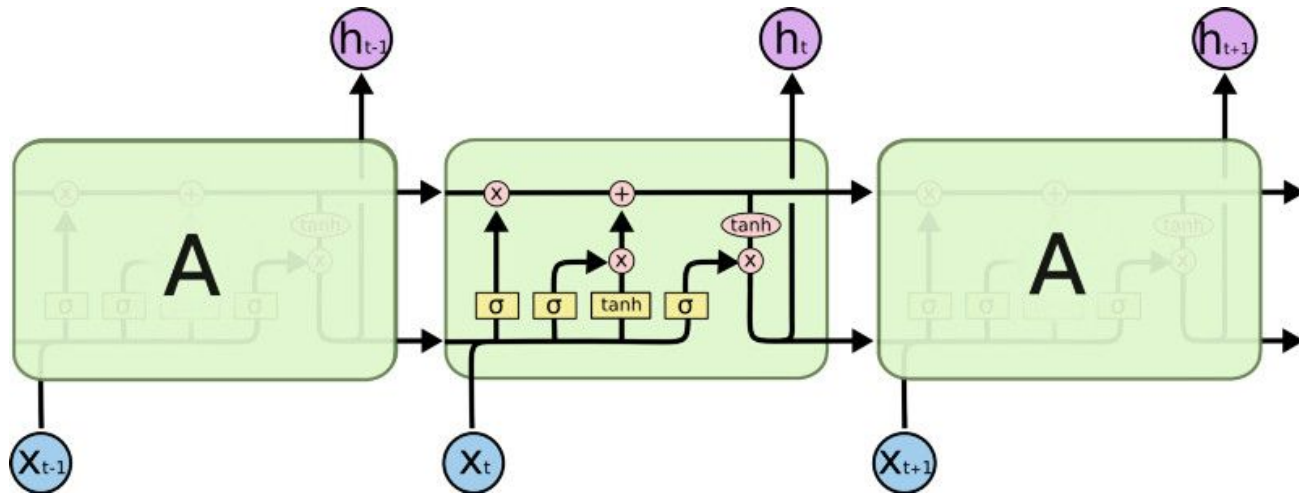
the repeated application of the recurrent weight matrix

Gradients Vanish in LSTM

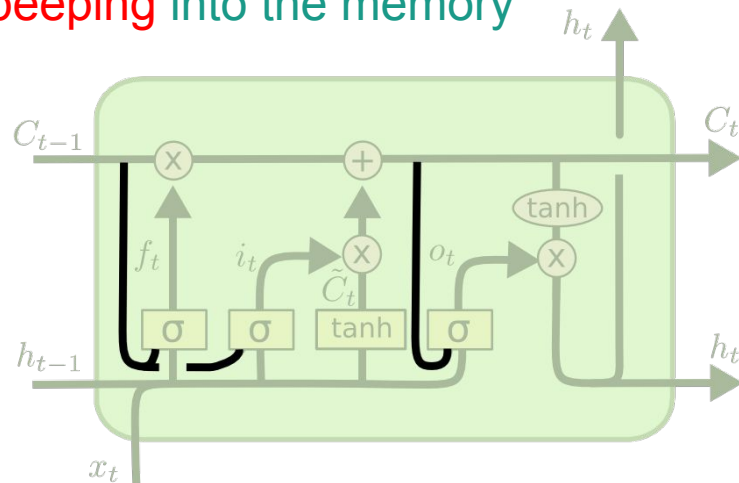
Only some of them, gradient information local in time will still be present.



Peephole LSTM



Allows “peeping into the memory”



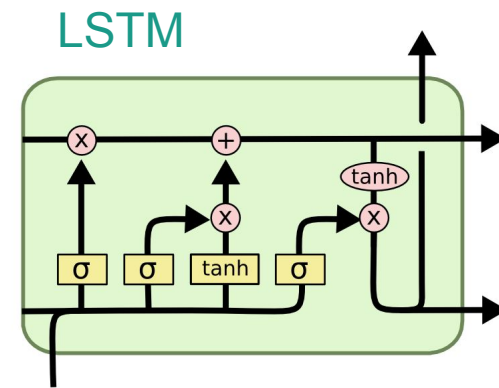
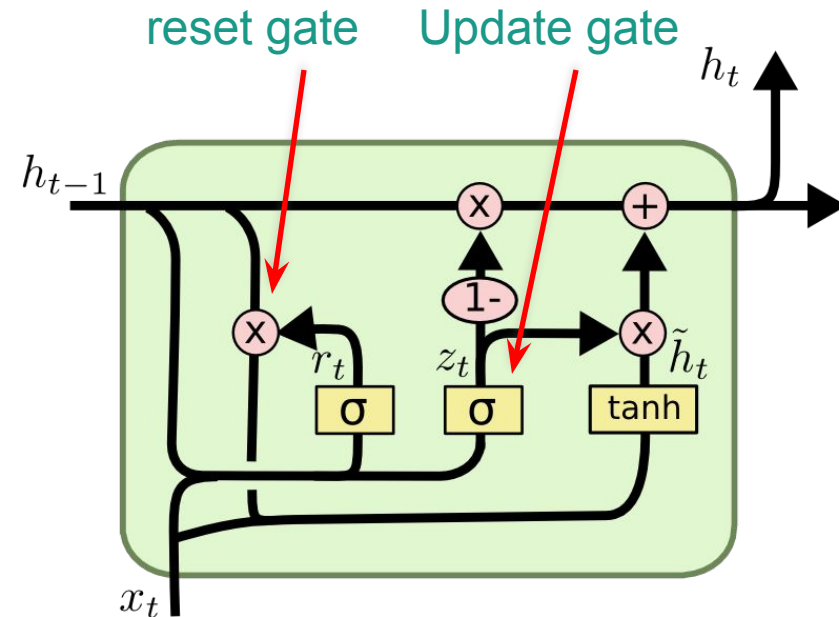
$$f_t = \sigma(W_f \cdot [C_{t-1}, h_{t-1}, x_t] + b_f)$$

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

$$o_t = \sigma(W_o \cdot [C_t, h_{t-1}, x_t] + b_o)$$

GRU – gated recurrent unit

(more compression)



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

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

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

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

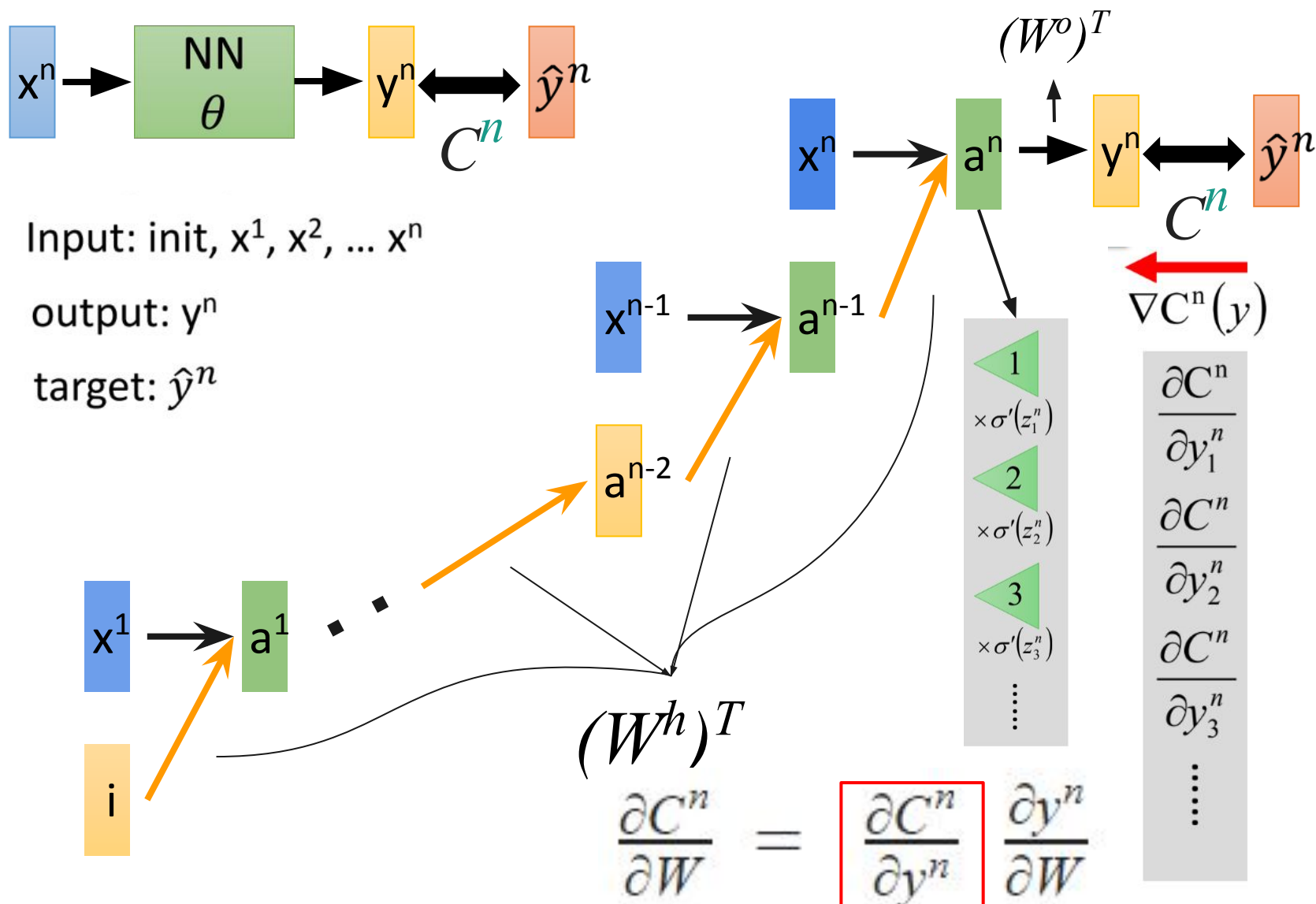
It combines the **forget** and **input** into a single **update gate**.
It also merges the cell state and hidden state. This is simpler than LSTM. There are many other variants too.

$X, *$: element-wise multiply

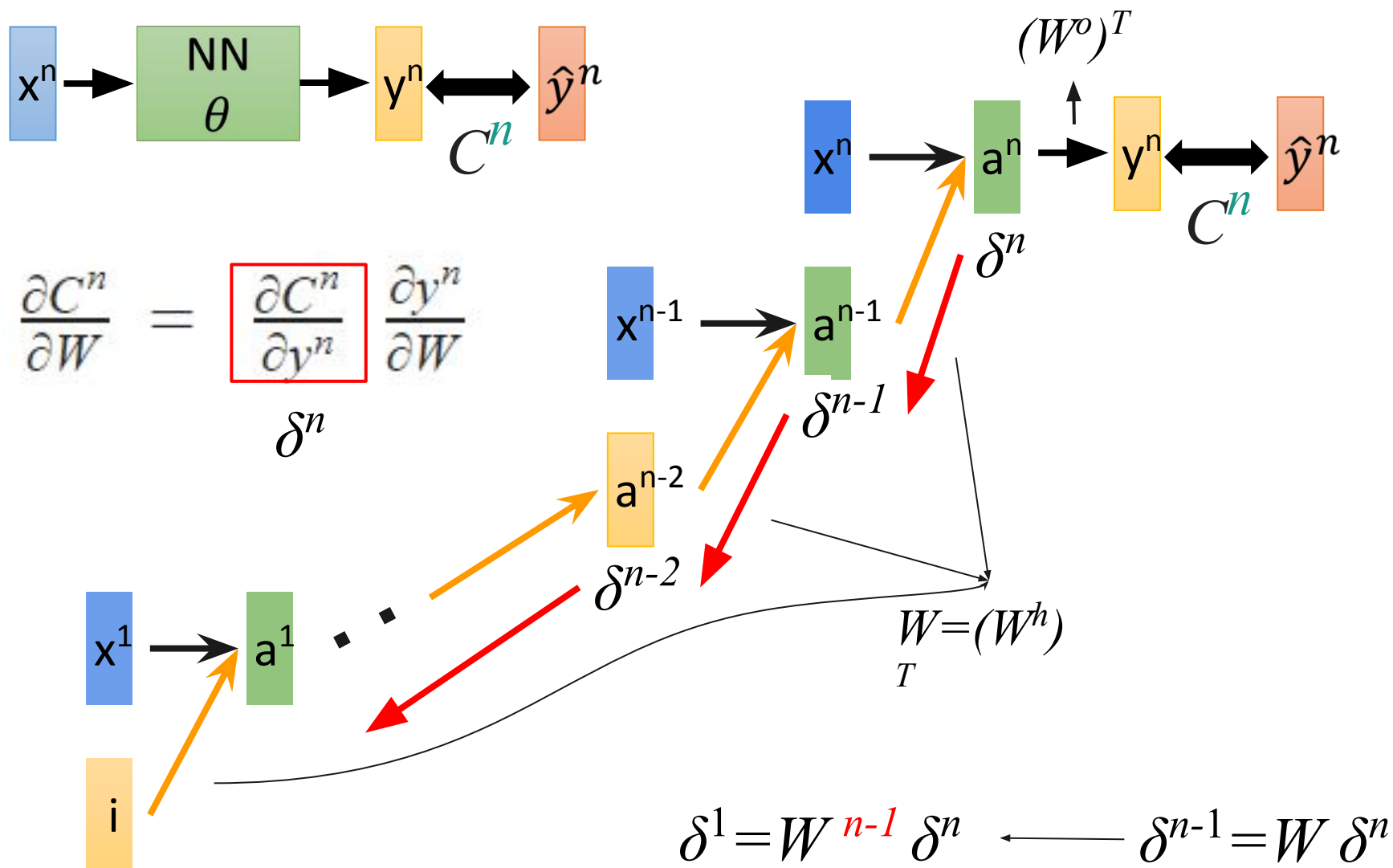
A decorative header consisting of a series of overlapping triangles in various shades of blue, teal, and orange, creating a jagged, mountain-like silhouette.

BPTT

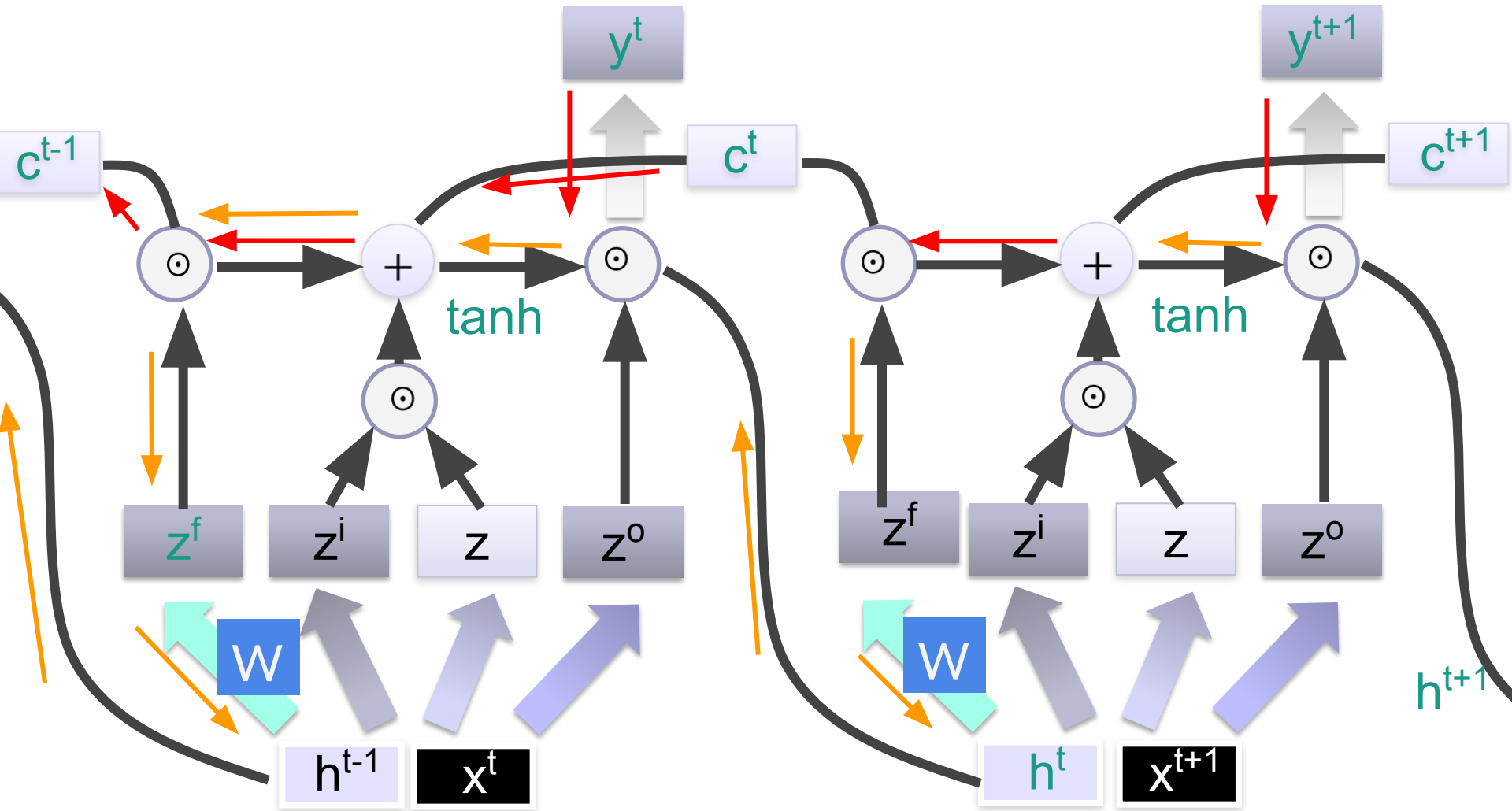
Backpropagation Through Time



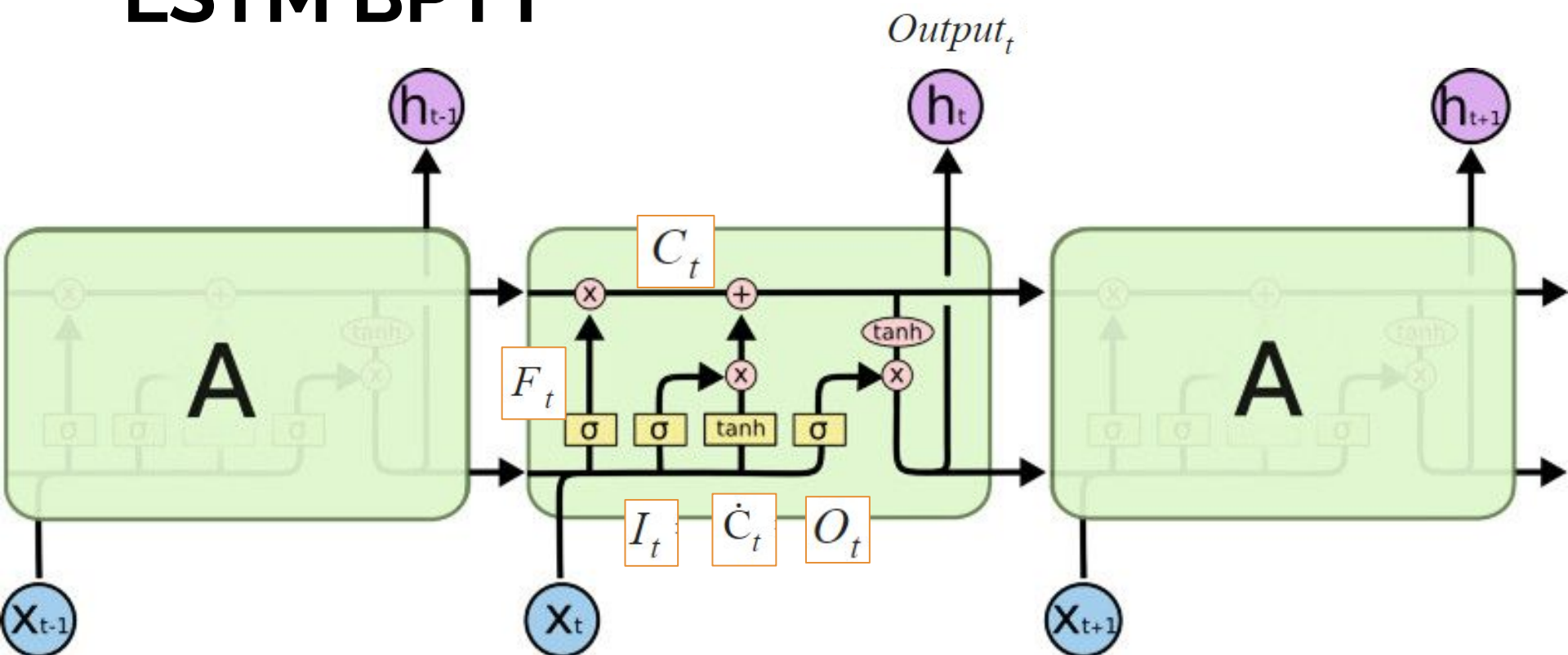
Gradient Exploding / Vanishing



LSTM BPTT flow



LSTM BPTT



$$F_t = \sigma(W_f X_t + W_{recf} h_{t-1} + b_f) \quad \dot{C}_t = \tanh(W_{\dot{C}} X_t + W_{rec\dot{C}} h_{t-1} + b_{\dot{C}})$$

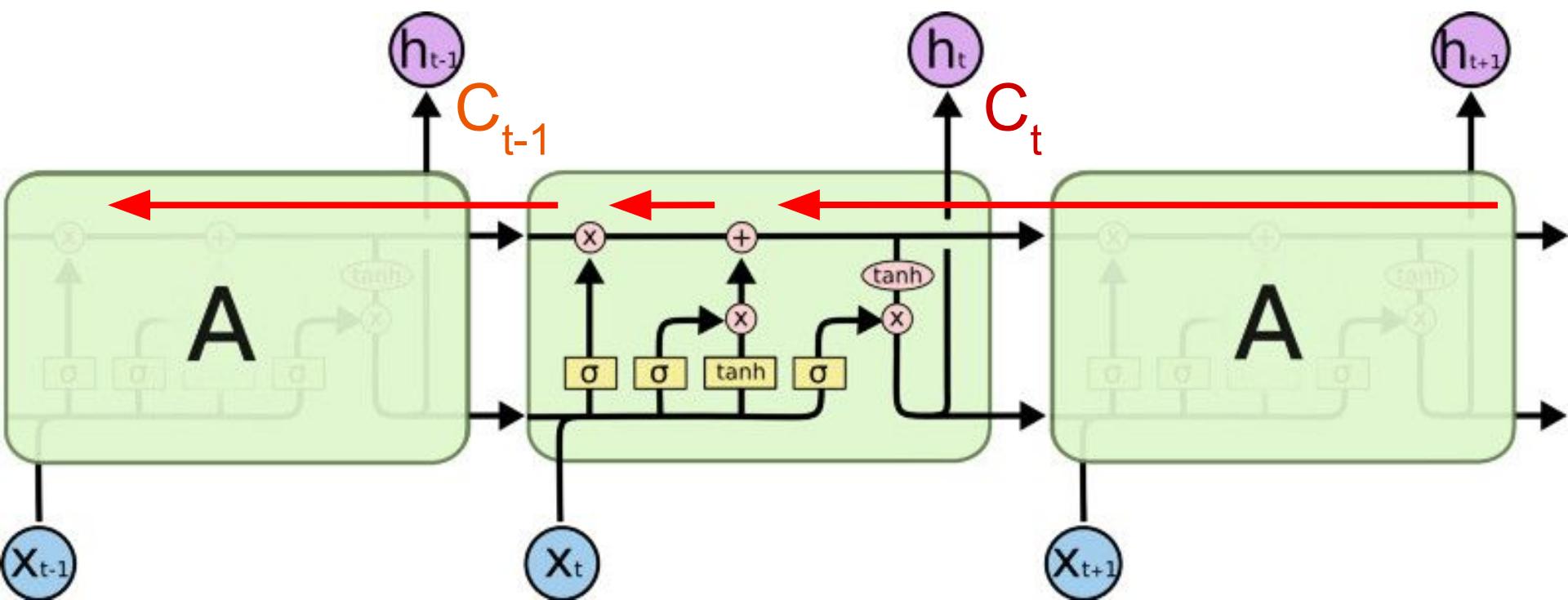
$$I_t = \sigma(W_i X_t + W_{reci} h_{t-1} + b_i) \quad C_t = C_{t-1} F_t + I_t \dot{C}_t$$

$$O_t = \sigma(W_o X_t + W_{reco} h_{t-1} + b_o)$$

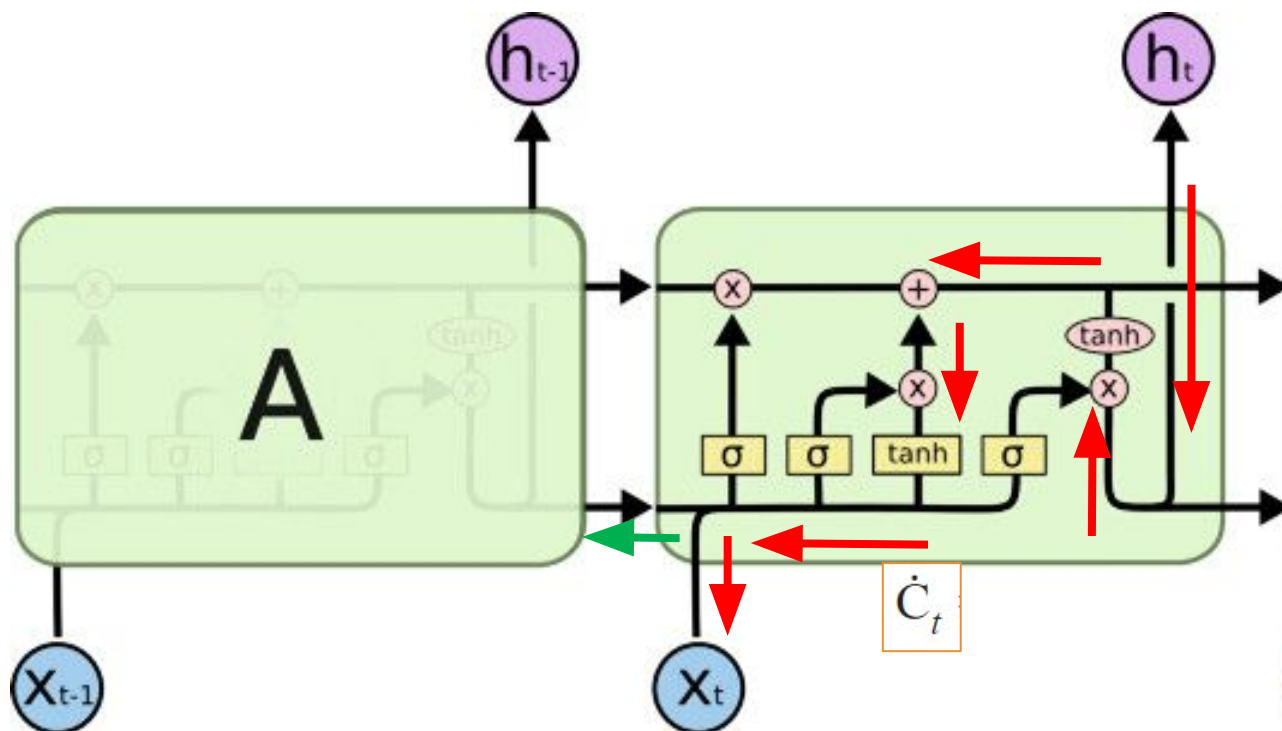
$$Output_t = O_t \tanh(C_t)$$

$$Cost_t = \frac{1}{t} (Output_t - \hat{y}_t)^2$$

Cell memory_(t=2)



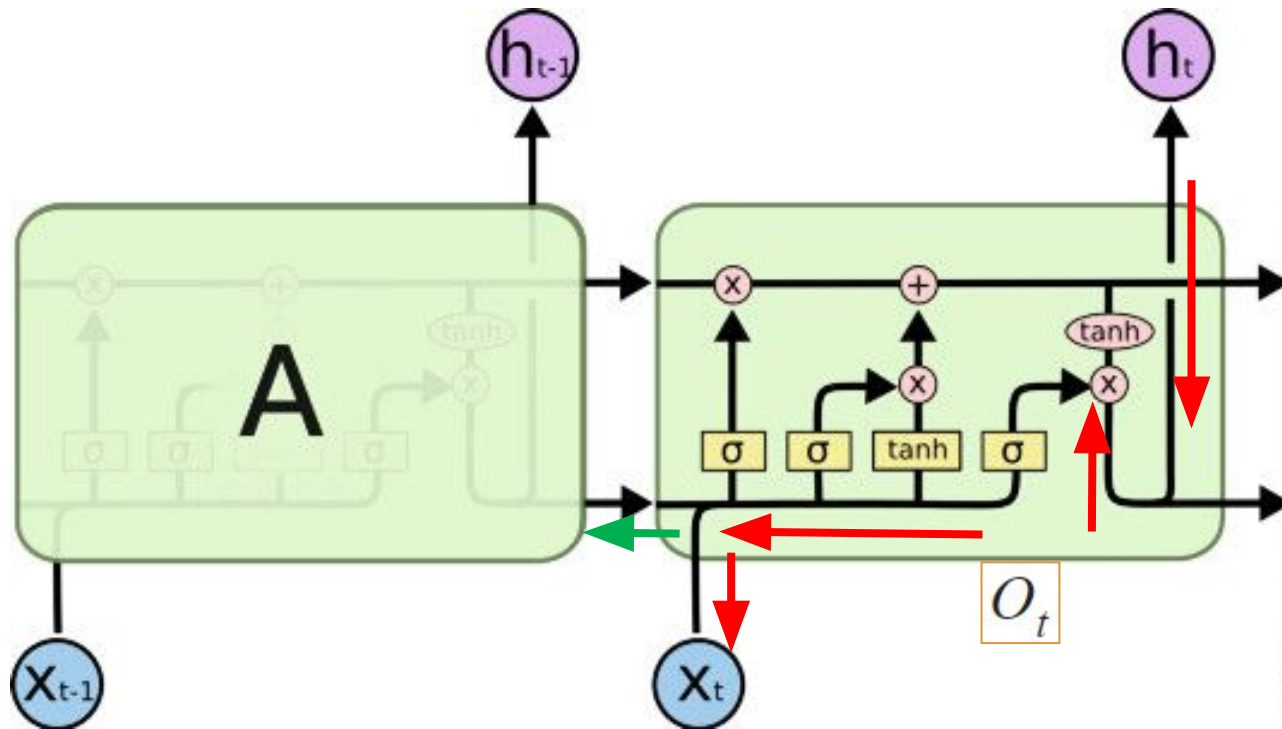
Memory Backpropagation (t = 2)



$$\begin{aligned} \frac{dCost_2}{dW_{\dot{C}}} &= \frac{dCost_2}{dOutput_2} \frac{dOutput_2}{dC_2(out)} \frac{dC_2(out)}{dC_2(in)} \frac{dC_2(in)}{d\dot{C}_2(out)} \frac{d\dot{C}_2(out)}{d\dot{C}_2(in)} \frac{d\dot{C}_2(in)}{dW_{\dot{C}}} \\ &= (Output_2 - \hat{y}_2) O_2 (1 - \tanh(C_2)^2) I_2 (1 - \tanh()^2) X_2 \end{aligned}$$

$$\begin{aligned} \frac{dCost_2}{dW_{rec\dot{C}}} &= \frac{dCost_2}{dOutput_2} \frac{dOutput_2}{dC_2(out)} \frac{dC_2(out)}{dC_2(in)} \frac{dC_2(in)}{d\dot{C}_2(out)} \frac{d\dot{C}_2(out)}{d\dot{C}_2(in)} \frac{d\dot{C}_2(in)}{dW_{rec\dot{C}}} \\ &= (Output_2 - \hat{y}_2) O_2 (1 - \tanh(C_2)^2) I_2 (1 - \tanh()^2) Output_1 \end{aligned}$$

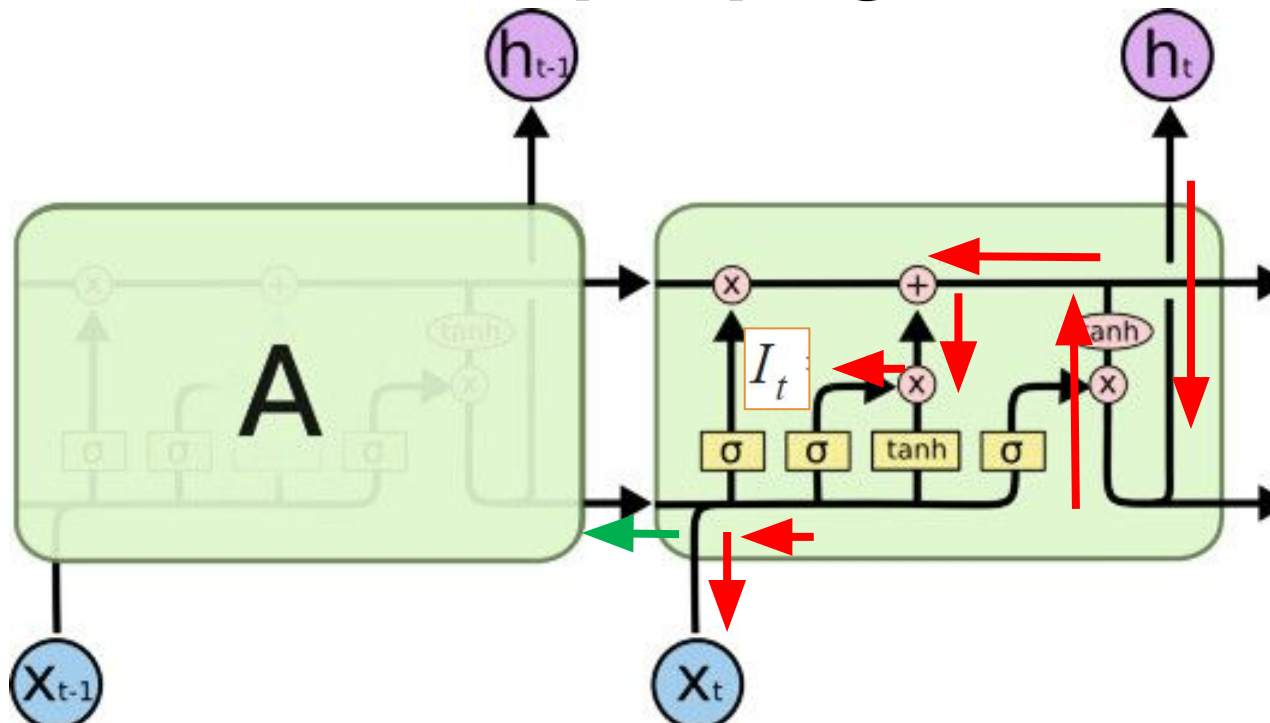
Output Gate Backpropagation ($t = 2$)



$$\begin{aligned} \frac{dCost_2}{dW_o} &= \frac{dCost_2}{dOutput_2} \frac{dOutput_2}{dO_2(out)} \frac{dO_2(out)}{dO_2(in)} \frac{dO_2(in)}{dW_o} \\ &= (Output_2 - \hat{y}_2) \tanh(C_2) (\sigma()) (1 - \sigma())) X_2 \end{aligned}$$

$$\begin{aligned} \frac{dCost_2}{dW_{reco}} &= \frac{dCost_2}{dOutput_2} \frac{dOutput_2}{dO_2(out)} \frac{dO_2(out)}{dO_2(in)} \frac{dO_2(in)}{dW_{reco}} \\ &= (Output_2 - \hat{y}_2) \tanh(C_2) (\sigma()) (1 - \sigma())) Output_1 \end{aligned}$$

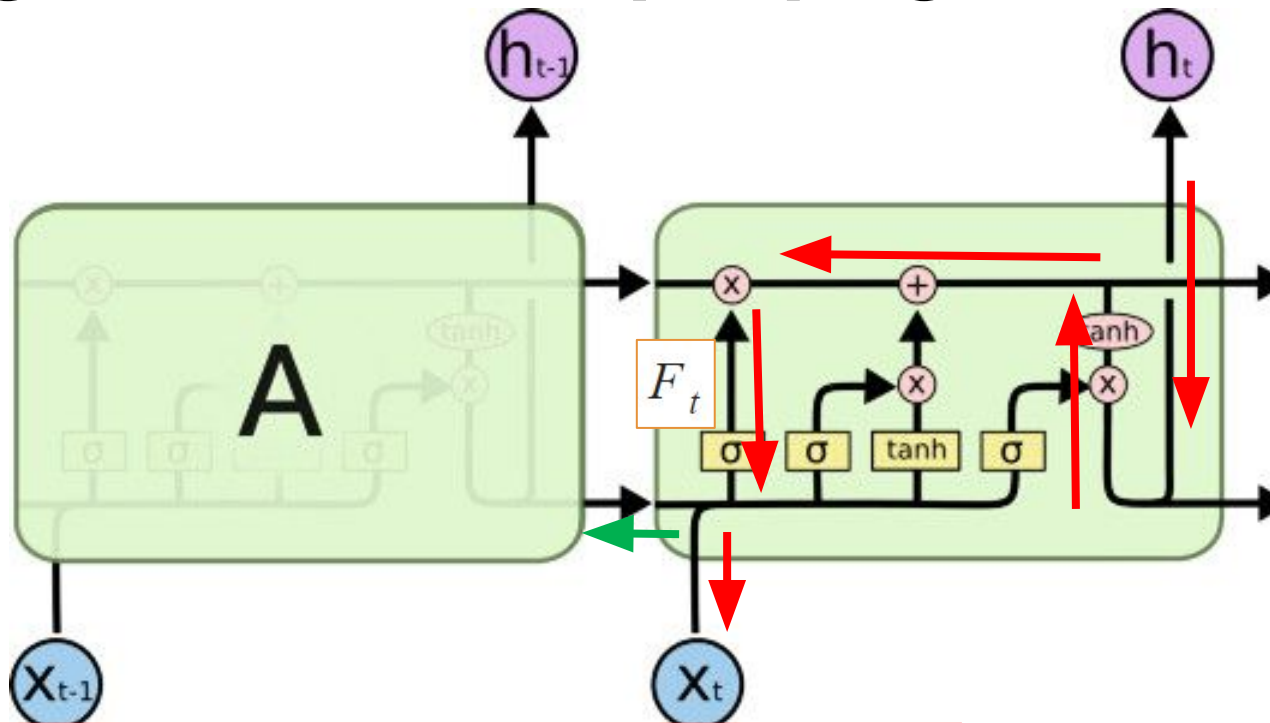
Input Gate Backpropagation (t = 2)



$$\begin{aligned} \frac{dCost_2}{dW_I} &= \frac{dCost_2}{dOutput_2} \frac{dOutput_2}{dC_2(out)} \frac{dC_2(out)}{dC_2(in)} \frac{dC_2(in)}{dI_2(out)} \frac{dI_2(out)}{dI_2(in)} \frac{dI_2(in)}{dW_I} \\ &= (Output_2 - \hat{y}_2) O_2 (1 - \tanh(C_2)^2) \dot{C}_2(\sigma()) (1 - \sigma()) X_2 \end{aligned}$$

$$\begin{aligned} \frac{dCost_2}{dW_{recI}} &= \frac{dCost_2}{dOutput_2} \frac{dOutput_2}{dC_2(out)} \frac{dC_2(out)}{dC_2(in)} \frac{dC_2(in)}{dI_2(out)} \frac{dI_2(out)}{dI_2(in)} \frac{dI_2(in)}{dW_{recI}} \\ &= (Output_2 - \hat{y}_2) O_2 (1 - \tanh(C_2)^2) \dot{C}_2(\sigma()) (1 - \sigma()) Output_1 \end{aligned}$$

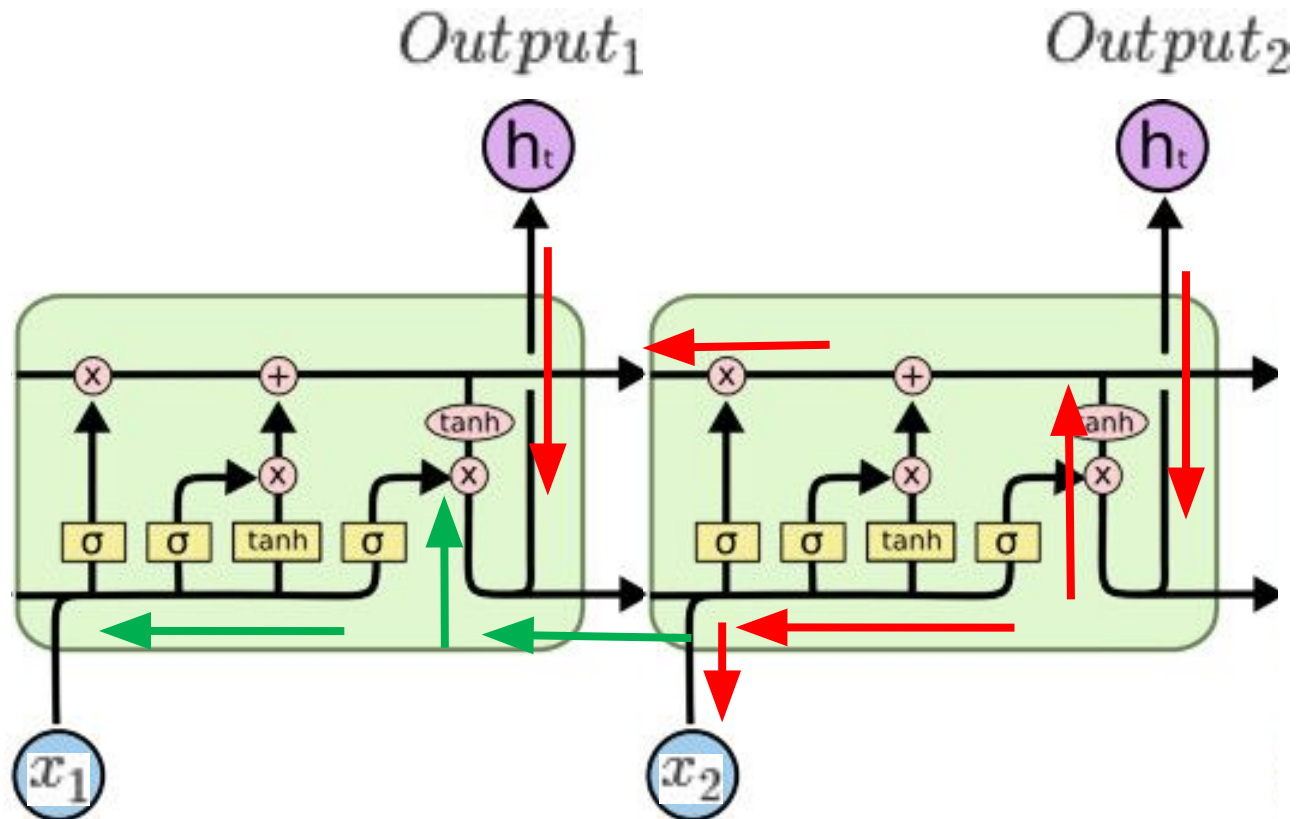
Forget Gate Backpropagation (t = 2)



$$\begin{aligned} \frac{dCost_2}{dW_F} &= \frac{dCost_2}{dOutput_2} \frac{dOutput_2}{dC_2(out)} \frac{dC_2(out)}{dC_2(in)} \frac{dC_2(in)}{dF_2(out)} \frac{dF(out)}{dF_2(in)} \frac{dF_2(in)}{dW_F} \\ &= (Output_2 - \hat{y}_2) O_2 (1 - \tanh(C_2)^2) C_1(\sigma()) (1 - \sigma()) X_2 \end{aligned}$$

$$\begin{aligned} \frac{dCost_2}{dW_{recF}} &= \frac{dCost_2}{dOutput_2} \frac{dOutput_2}{dC_2(out)} \frac{dC_2(out)}{dC_2(in)} \frac{dC_2(in)}{dF_2(out)} \frac{dF(out)}{dF_2(in)} \frac{dF_2(in)}{dW_{recF}} \\ &= (Output_2 - \hat{y}_2) O_2 (1 - \tanh(C_2)^2) C_1(\sigma()) (1 - \sigma()) Output_1 \end{aligned}$$

Backpropagation to Cell1 ($t = 1$)



Backpropagation to Cell1 Output Gate_(t=1)

$$\frac{dCost_{total}}{dW_o} = \frac{dCost_1}{dW_o} + \boxed{\frac{dCost_2}{dW_o}}$$

$$= \frac{dCost_2}{dOutput_1} \frac{dOutput_1}{dO_1(out)} \frac{dO_1(out)}{dO_1(in)} \frac{dO_1(in)}{dW_o} +$$

$$\frac{dCost_2}{dOutput_2} \frac{dOutput_2}{dO_2(out)} \frac{dO_2(out)}{dO_2(in)} \frac{dO_2(in)}{dOutput_1} \frac{dOutput_1}{dO_1(out)} \frac{dO_1(in)}{dW_o} +$$

$$\frac{dCost_2}{dOutput_2} \frac{dOutput_2}{dC_2(out)} \frac{dC_2(out)}{dC_2(in)} \frac{dC_2(in)}{dI_2(out)} \frac{dI_2(out)}{dI_2(in)} \frac{dI_2(in)}{dOutput_1} \frac{dOutput_1}{dO_1(out)} \frac{dO_1(in)}{dW_o} +$$

$$\frac{dCost_2}{dOutput_2} \frac{dOutput_2}{dC_2(out)} \frac{dC_2(out)}{dC_2(in)} \frac{dC_2(in)}{d\dot{C}_2(out)} \frac{d\dot{C}_2(out)}{d\dot{C}_2(in)} \frac{d\dot{C}_2(in)}{dOutput_1} \frac{dOutput_1}{dO_1(out)} \frac{dO_1(in)}{dW_o} +$$

$$\frac{dCost_2}{dOutput_2} \frac{dOutput_2}{dC_2(out)} \frac{dC_2(out)}{dC_2(in)} \frac{dC_2(in)}{dF_2(out)} \frac{dF_2(out)}{dF_2(in)} \frac{dF_2(in)}{dOutput_1} \frac{dOutput_1}{dO_1(out)} \frac{dO_1(in)}{dW_o}$$

Simplify Backpropagation_(t = 1)

$$\frac{dCost_{total}}{dW_o} = \frac{dCost_1}{dW_o} + \frac{dCost_2}{dW_o}$$

$$= (Output_2 - \hat{y}_1) \tanh(C_1) (\sigma() (1 - \sigma())) X_1 +$$

$$(Output_2 - \hat{y}_2) \tanh(C_2) (\sigma() (1 - \sigma())) W_{recO} \tanh(C_1) (\sigma() (1 - \sigma())) X_1 +$$

$$(Output_2 - \hat{y}_2) O_2 (1 - \tanh(C_2)^2) I_2 (1 - \tanh()^2) W_{rec\dot{C}} \tanh(C_1) (\sigma() (1 - \sigma())) X_1 +$$

$$(Output_2 - \hat{y}_2) O_2 (1 - \tanh(C_2)^2) \dot{C}_2 (\sigma() (1 - \sigma())) W_{recI} \tanh(C_1) (\sigma() (1 - \sigma())) X_1 +$$

$$(Output_2 - \hat{y}_2) O_2 (1 - \tanh(C_2)^2) C_1 (\sigma() (1 - \sigma())) W_{recF} \tanh(C_1) (\sigma() (1 - \sigma())) X_1$$

Thank you