

# **Recurrent networks**

**Recurrent nets unrolling in time**

**Karel Zimmermann**

**Czech Technical University in Prague**

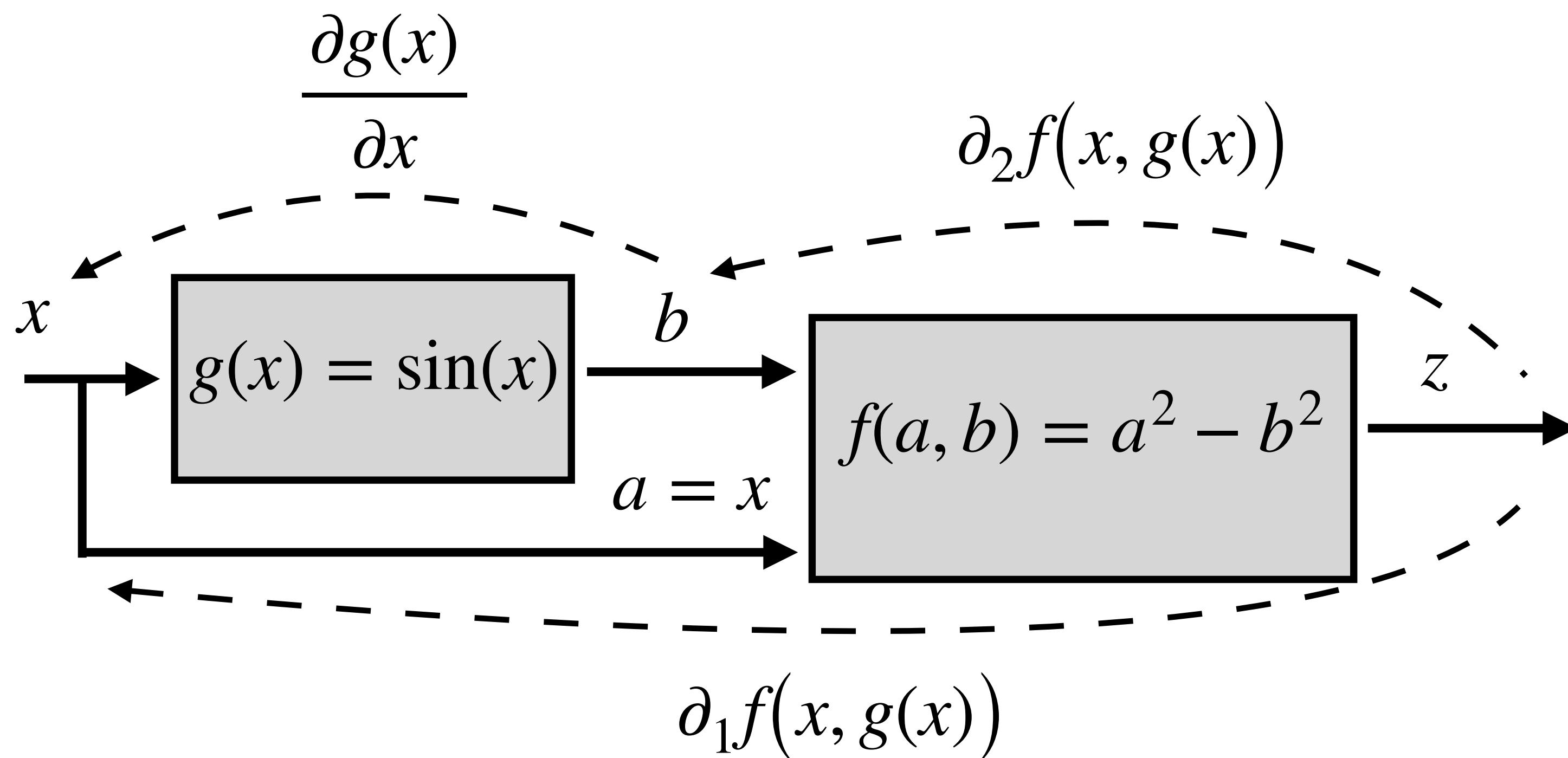
**Faculty of Electrical Engineering, Department of Cybernetics**



Prerequisites: derivative of compound function

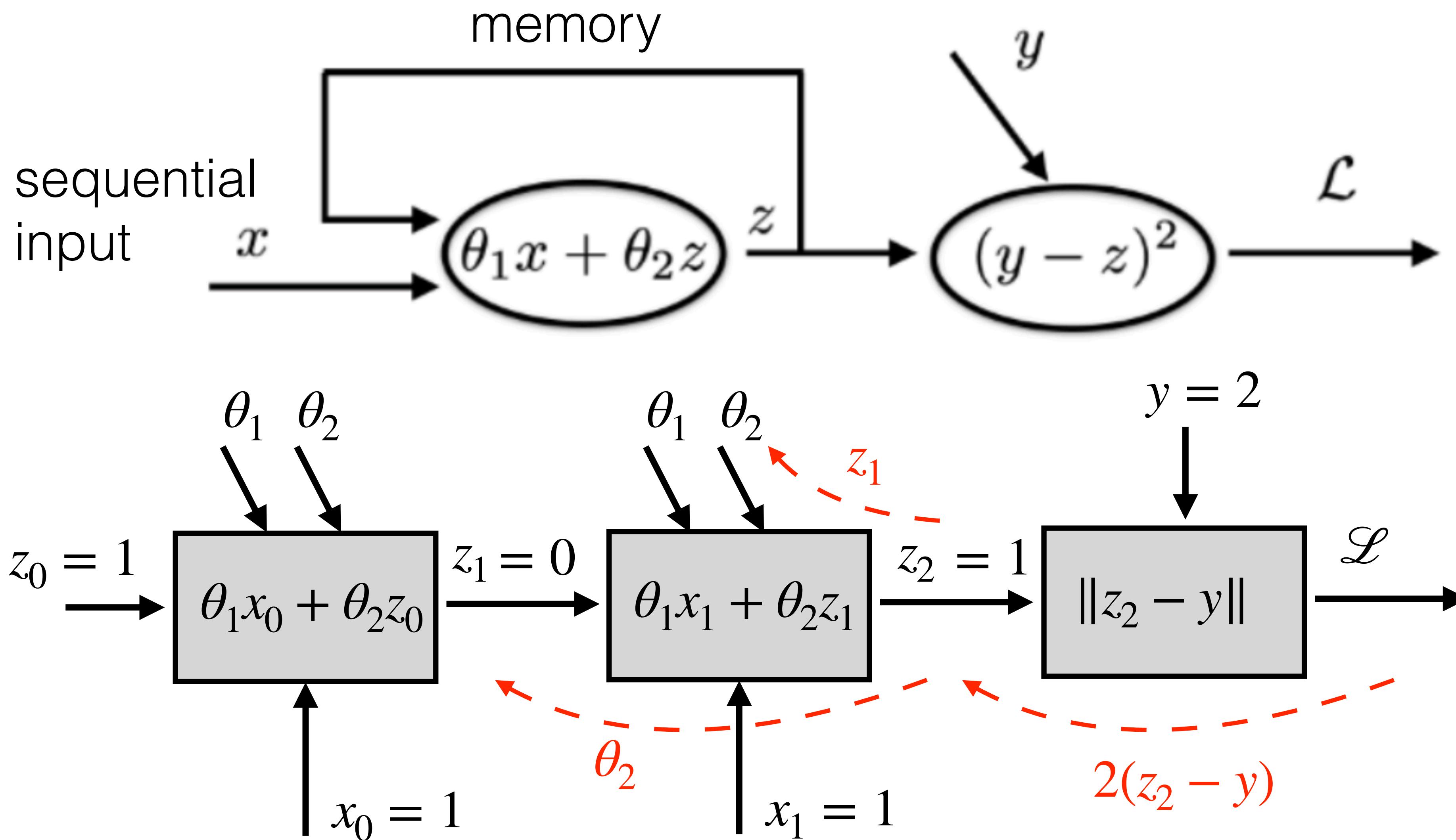
Introduce notation:  $\frac{\partial f(a, b)}{\partial a} = \partial_1 f(a, b)$ ,  $\frac{\partial f(a, b)}{\partial b} = \partial_2 f(a, b)$

$$f(a, b) = a^2 - b^2, \quad g(x) = \sin(x) \quad \Rightarrow \quad f(x, g(x)) = x^2 - \sin^2(x)$$



$$\frac{\partial f(x, g(x))}{\partial x} = \partial_1 f(x, g(x)) + \partial_2 f(x, g(x)) \frac{\partial g(x)}{\partial x} = 2x - 2 \sin(x)\cos(x)$$

# RNN example with backprop



input sequence:

$$x_0 = 1$$

$$x_1 = 1$$

output label:

$$y = 2$$

initial memory:

$$z_0 = 1$$

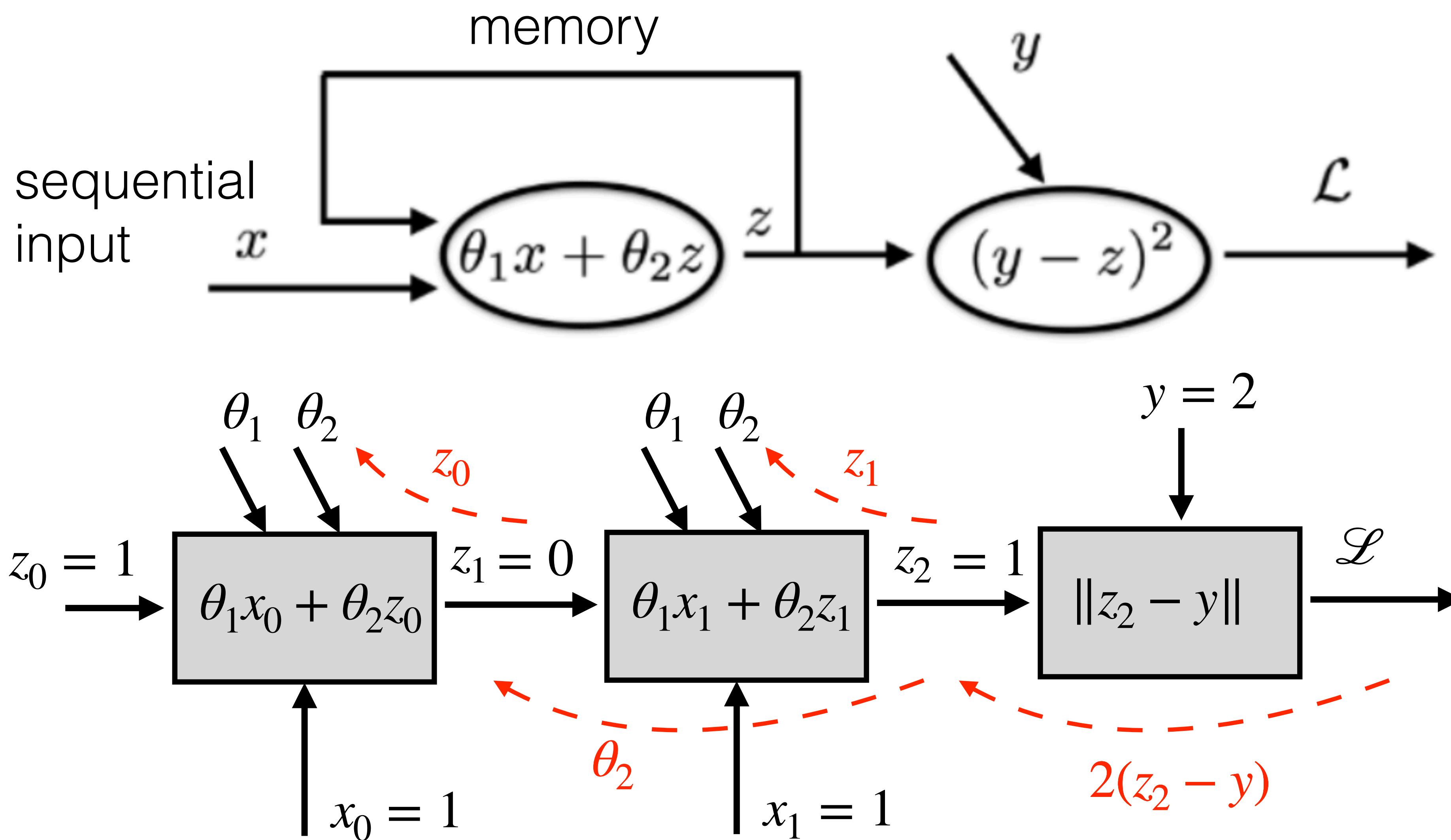
initial parameters:

$$\theta_1 = 1$$

$$\theta_2 = -1$$

$$\frac{\partial \mathcal{L}}{\partial \theta_2} =$$

# RNN example with backprop



$$\frac{\partial \mathcal{L}}{\partial \theta_2} =$$

input sequence:

$$x_0 = 1$$

$$x_1 = 1$$

output label:

$$y = 2$$

initial memory:

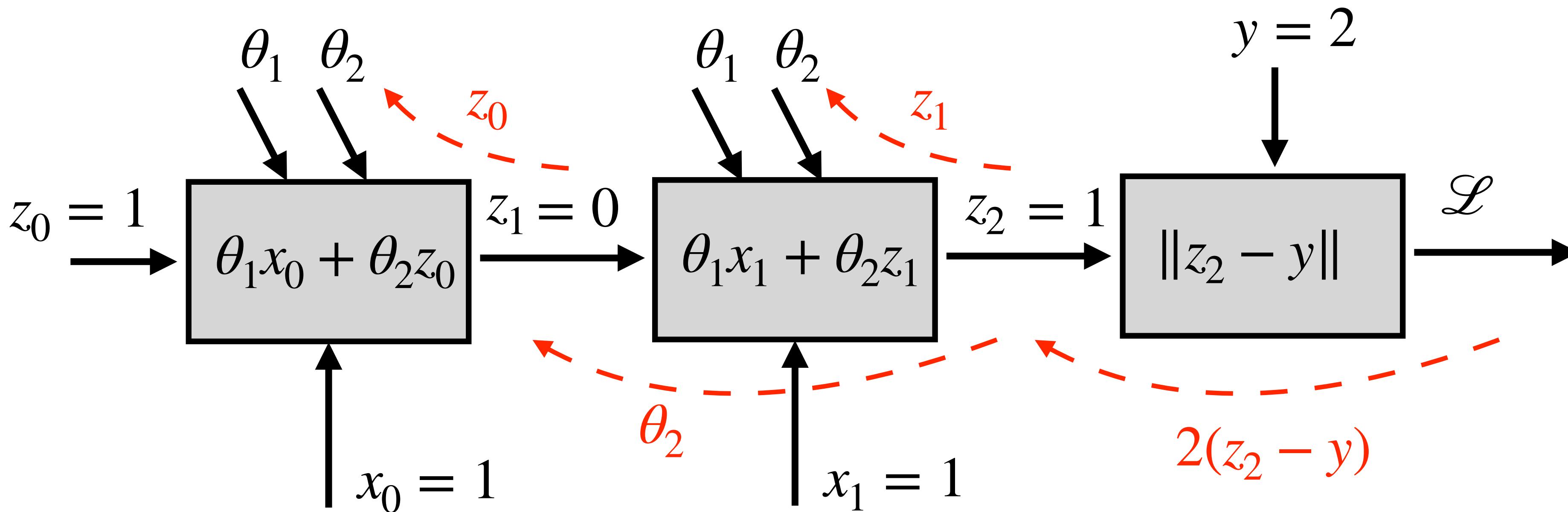
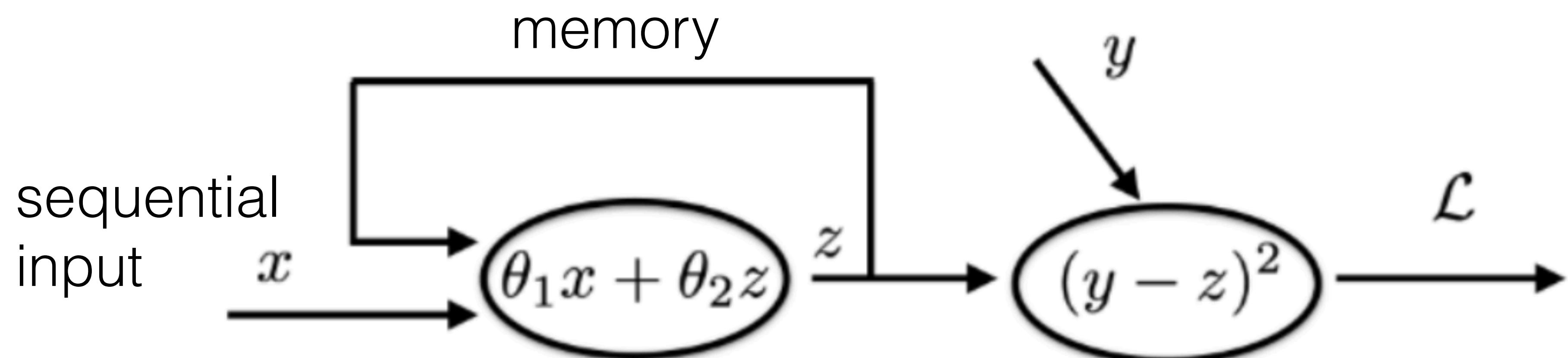
$$z_0 = 1$$

initial parameters:

$$\theta_1 = 1$$

$$\theta_2 = -1$$

# RNN example with backprop



$$\frac{\partial \mathcal{L}}{\partial \theta_2} = 2(z_2 - y) z_1 + 2(z_2 - y) \theta_2 z_0 = 2(1 - 2)0 + 2(1 - 2)(-1)1 = 2$$

input sequence:

$$x_0 = 1$$

$$x_1 = 1$$

output label:

$$y = 2$$

initial memory:

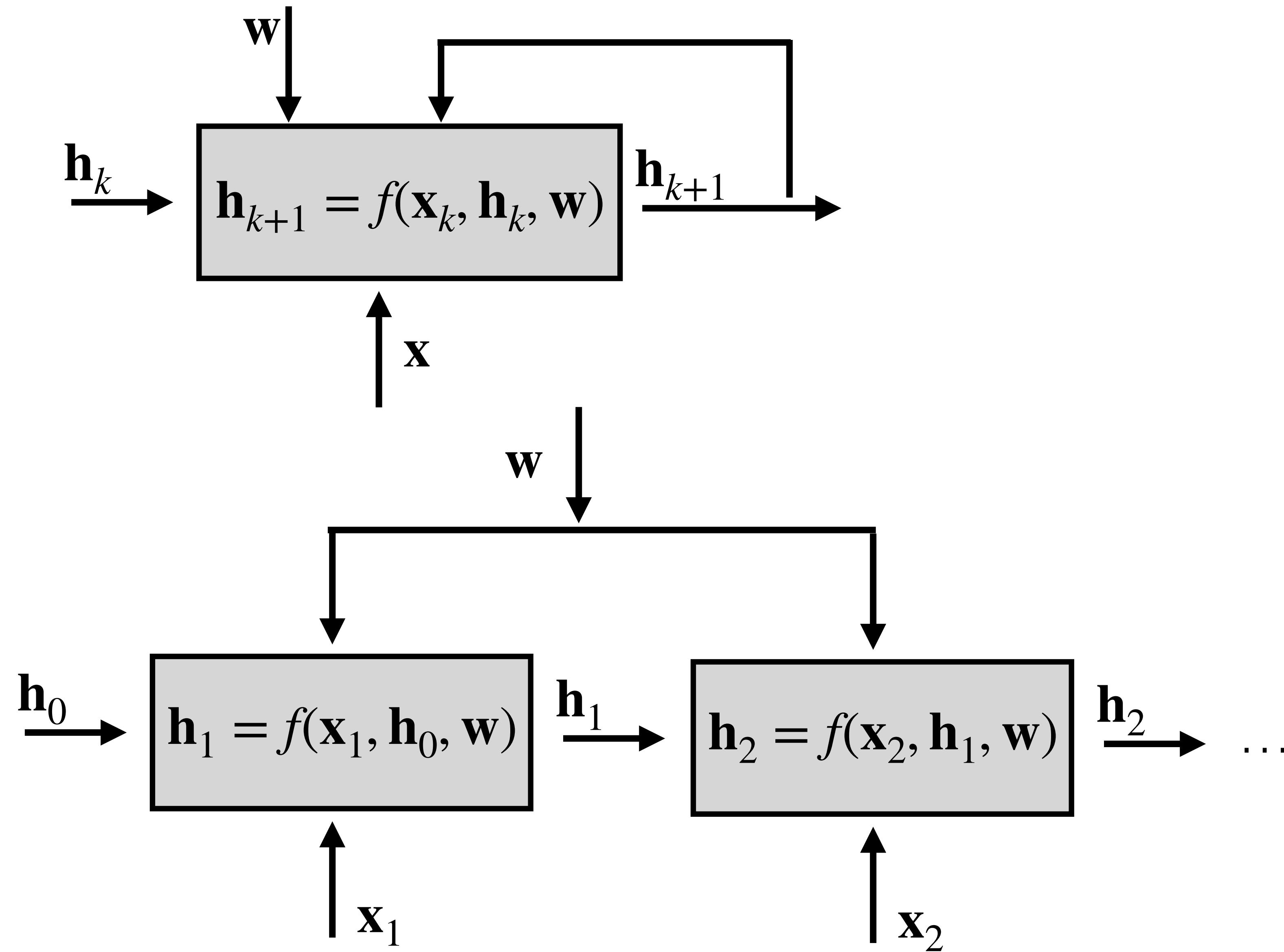
$$z_0 = 1$$

initial parameters:

$$\theta_1 = 1$$

$$\theta_2 = -1$$

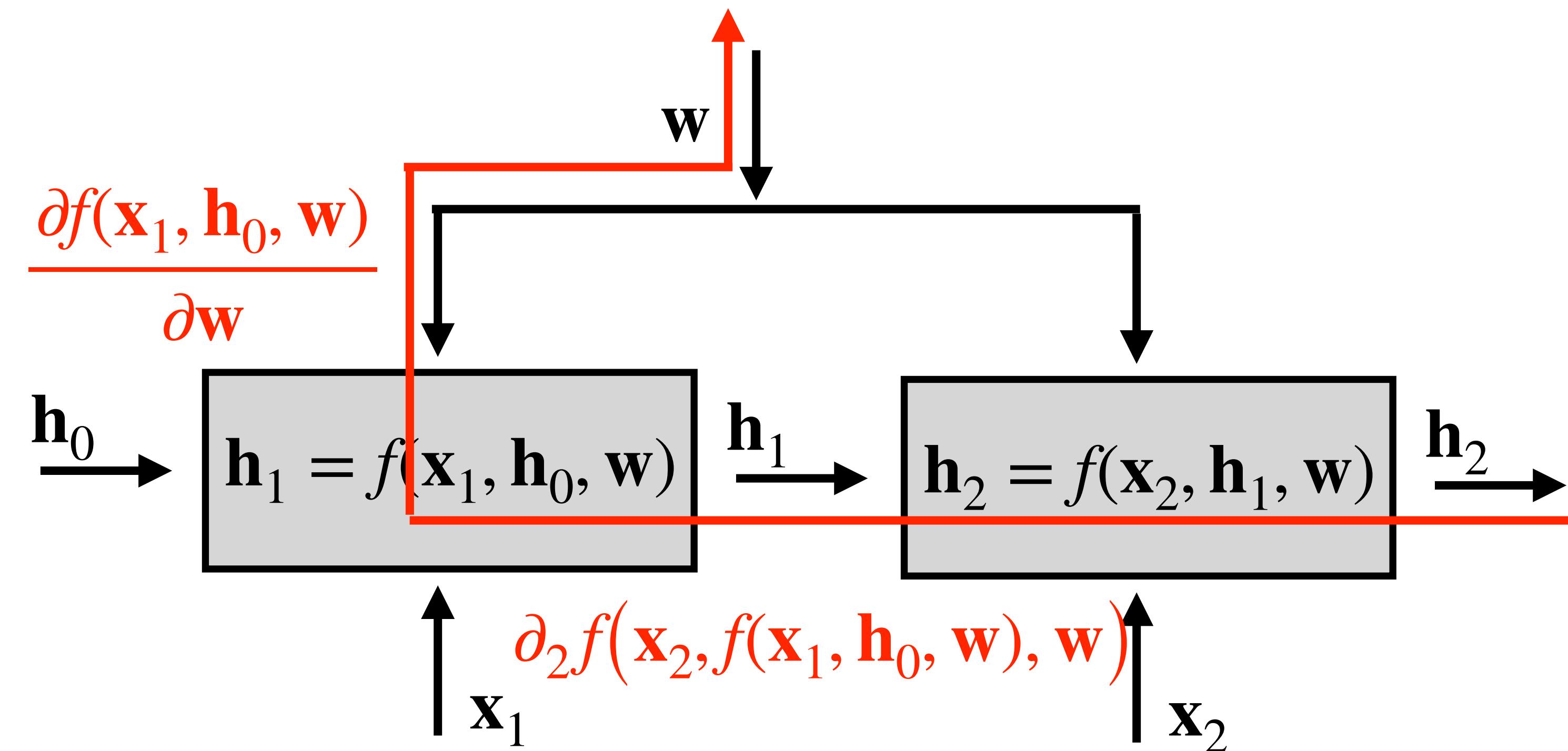
Recurrence resolved by unrolling the computation graph in time



Recurrence resolved by unrolling the computation graph in time

$$\mathbf{h}_1 = f(\mathbf{x}_1, \mathbf{h}_0, \mathbf{w}), \quad \mathbf{h}_2 = f(\mathbf{x}_2, \mathbf{h}_1, \mathbf{w}) \quad \Rightarrow \quad \mathbf{h}_2 = f(\mathbf{x}_2, f(\mathbf{x}_1, \mathbf{h}_0, \mathbf{w}), \mathbf{w})$$

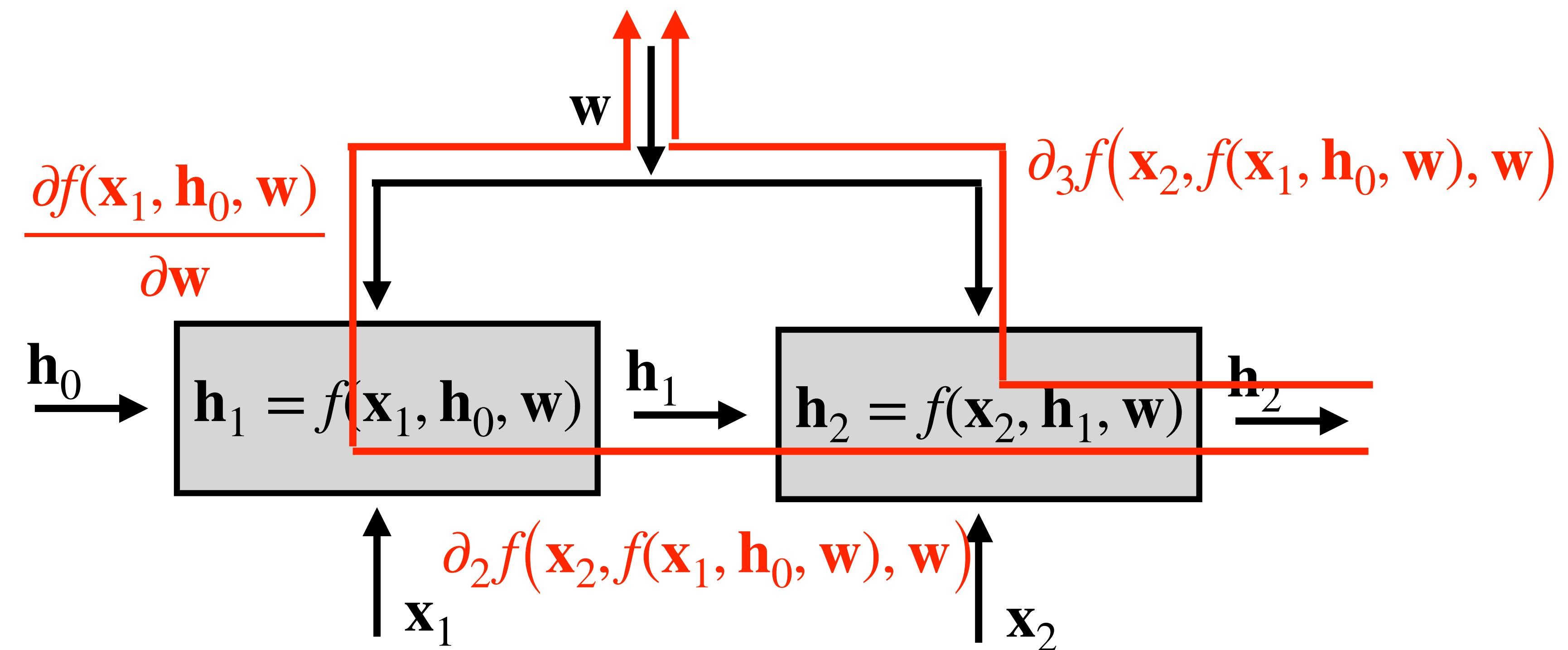
$$\frac{\partial f(\mathbf{x}_2, f(\mathbf{x}_1, \mathbf{h}_0, \mathbf{w}), \mathbf{w})}{\partial \mathbf{w}} = \partial_2 f(\mathbf{x}_2, f(\mathbf{x}_1, \mathbf{h}_0, \mathbf{w}), \mathbf{w}) \frac{\partial f(\mathbf{x}_1, \mathbf{h}_0, \mathbf{w})}{\partial \mathbf{w}}$$



Recurrence resolved by unrolling the computation graph in time

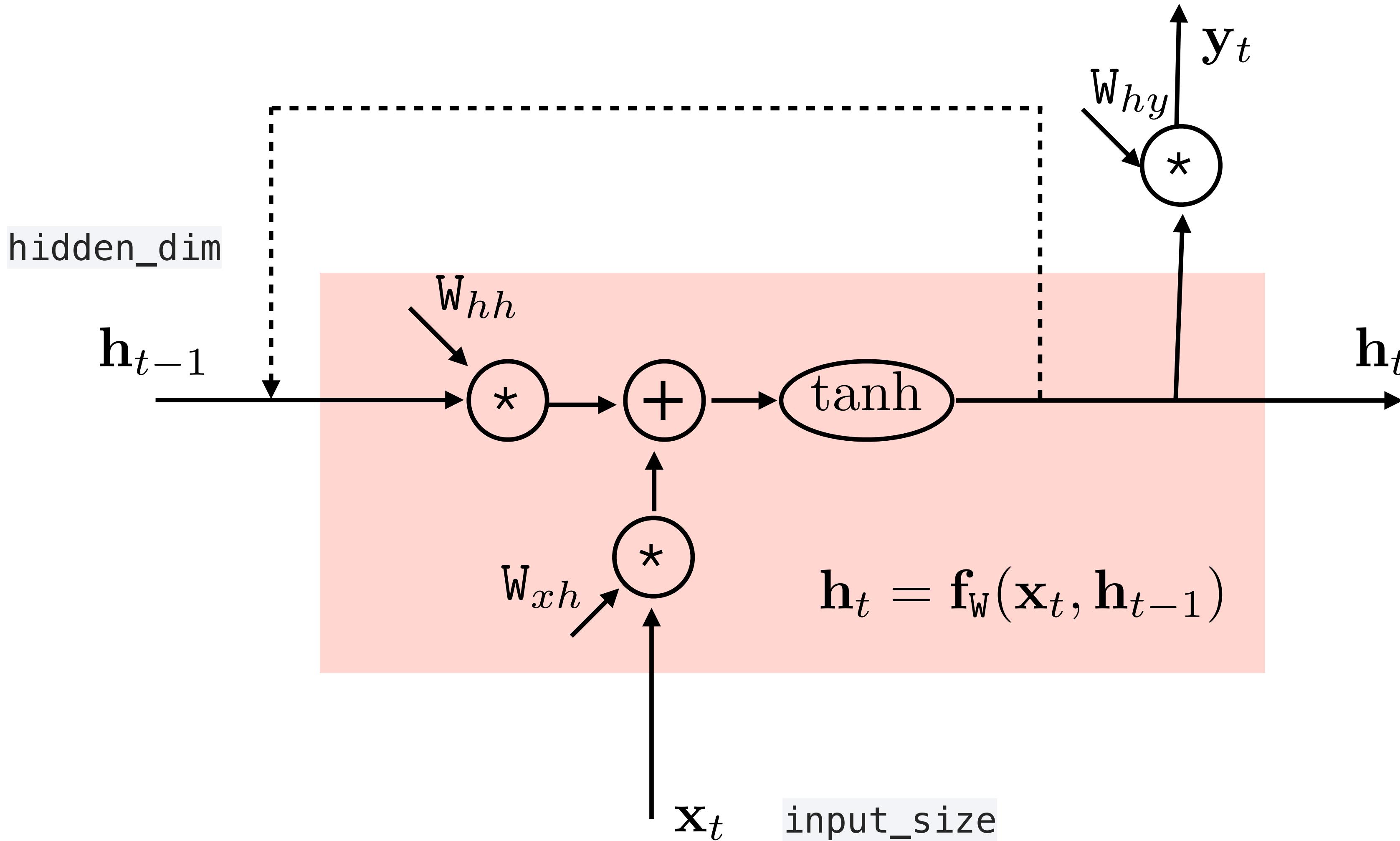
$$\mathbf{h}_1 = f(\mathbf{x}_1, \mathbf{h}_0, \mathbf{w}), \quad \mathbf{h}_2 = f(\mathbf{x}_2, \mathbf{h}_1, \mathbf{w}) \quad \Rightarrow \quad \mathbf{h}_2 = f(\mathbf{x}_2, f(\mathbf{x}_1, \mathbf{h}_0, \mathbf{w}), \mathbf{w})$$

$$\frac{\partial f(\mathbf{x}_2, f(\mathbf{x}_1, \mathbf{h}_0, \mathbf{w}), \mathbf{w})}{\partial \mathbf{w}} = \partial_2 f(\mathbf{x}_2, f(\mathbf{x}_1, \mathbf{h}_0, \mathbf{w}), \mathbf{w}) \frac{\partial f(\mathbf{x}_1, \mathbf{h}_0, \mathbf{w})}{\partial \mathbf{w}} + \partial_3 f(\mathbf{x}_2, f(\mathbf{x}_1, \mathbf{h}_0, \mathbf{w}), \mathbf{w})$$



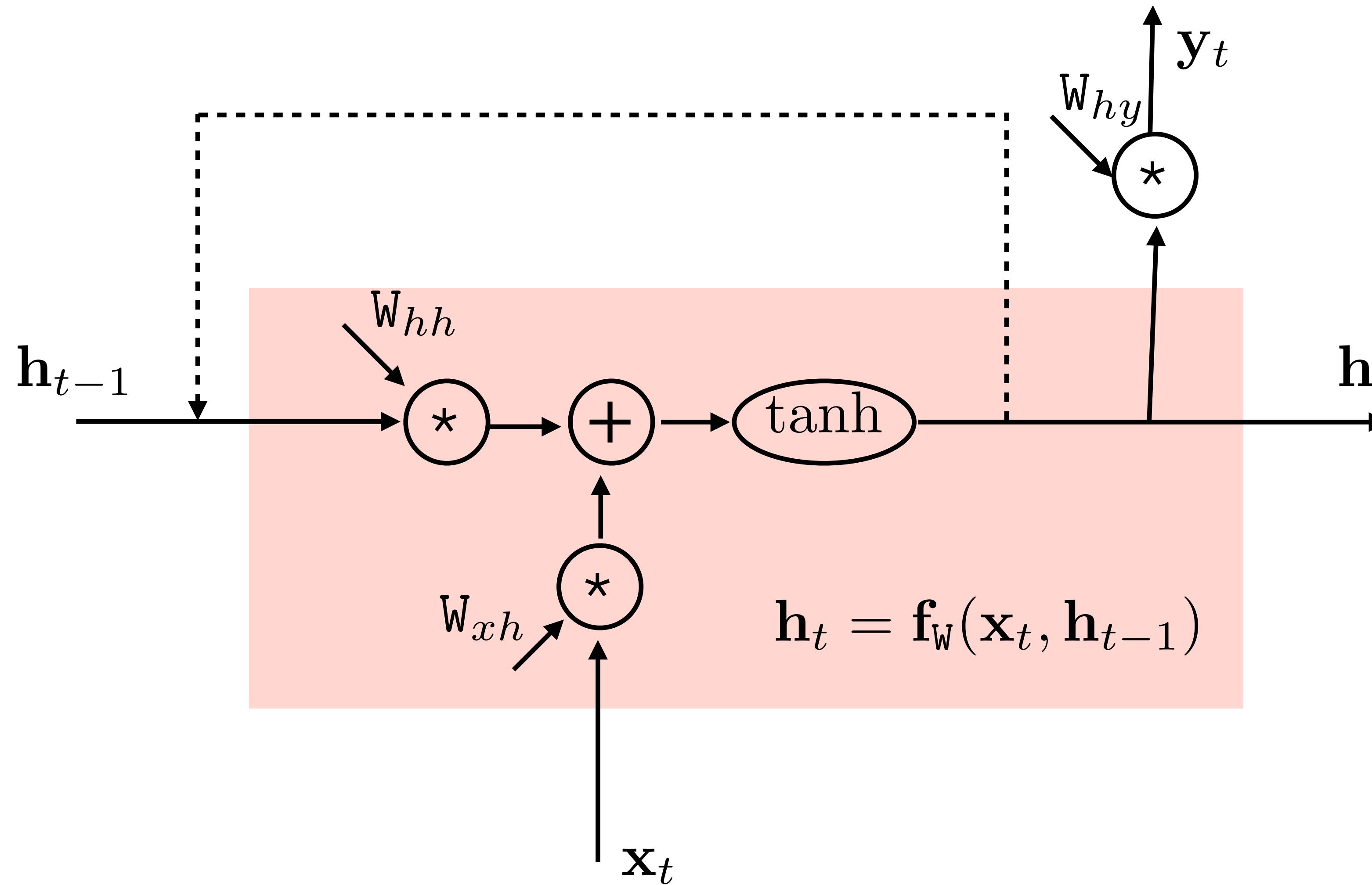
# Simple recurrent block

```
torch.nn.RNN(input_size, hidden_dim, n_layers)
```

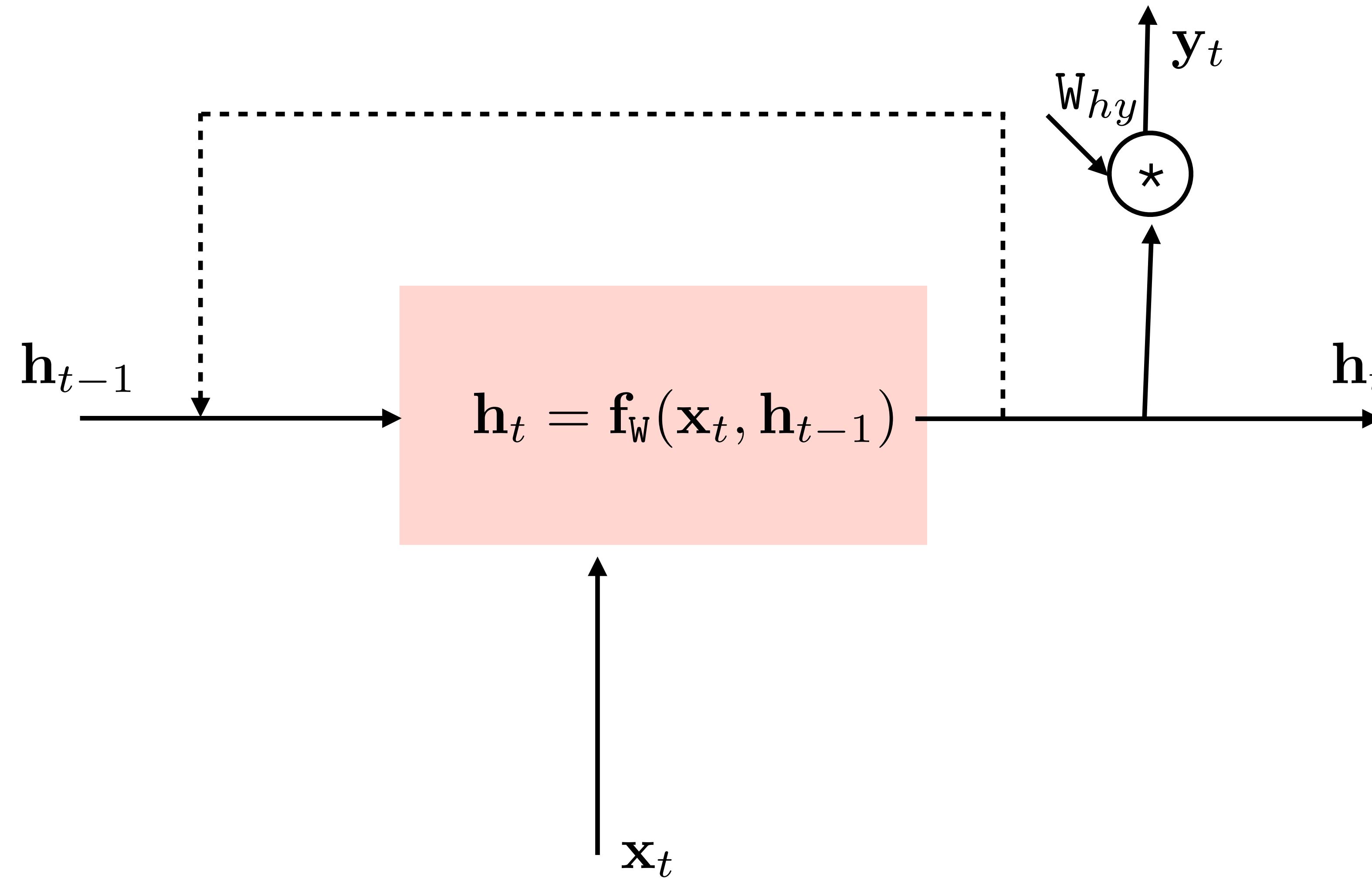


PyTorch: <https://pytorch.org/docs/stable/nn.html>

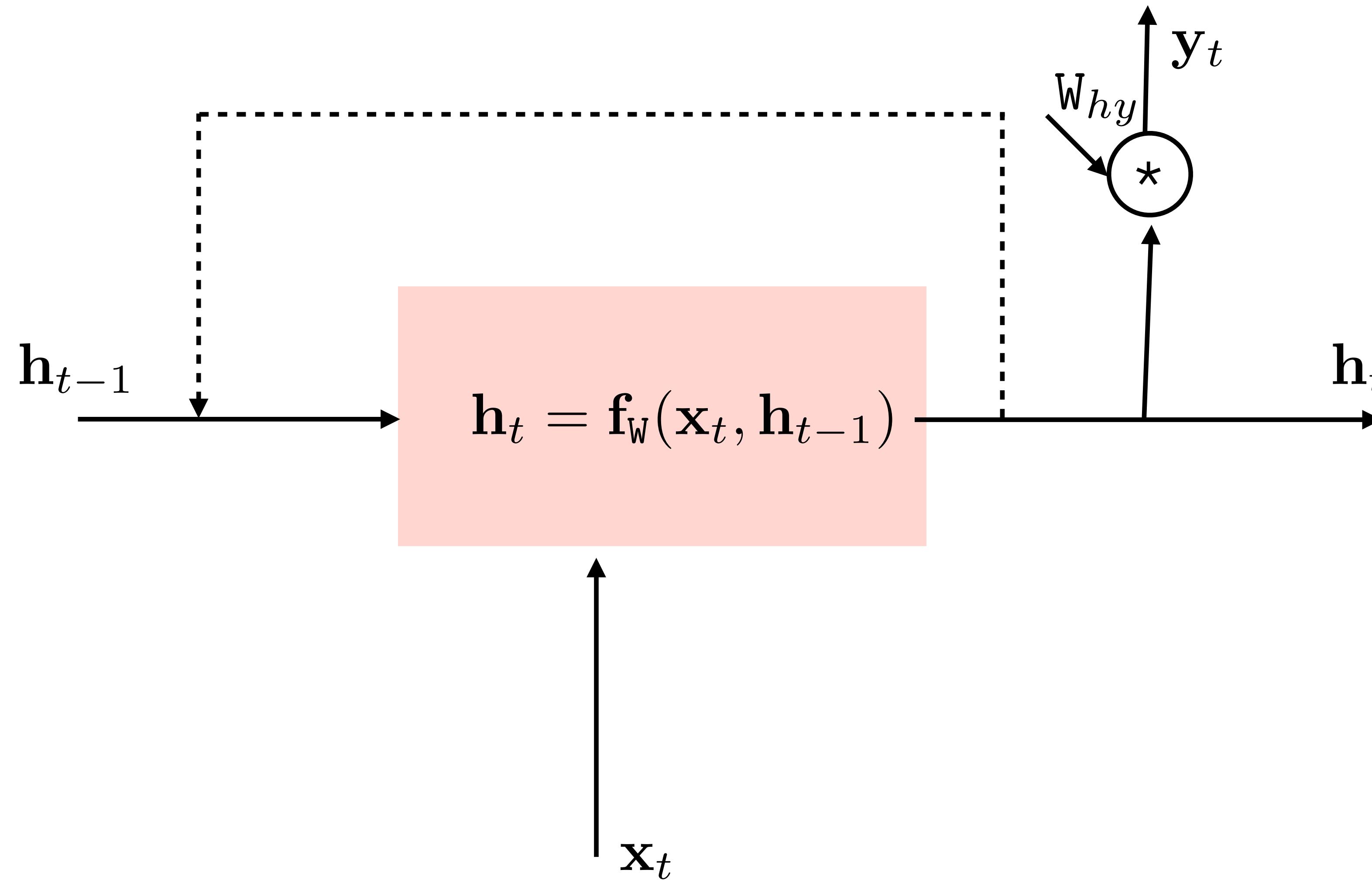
# Simple recurrent block - feed-forward pass



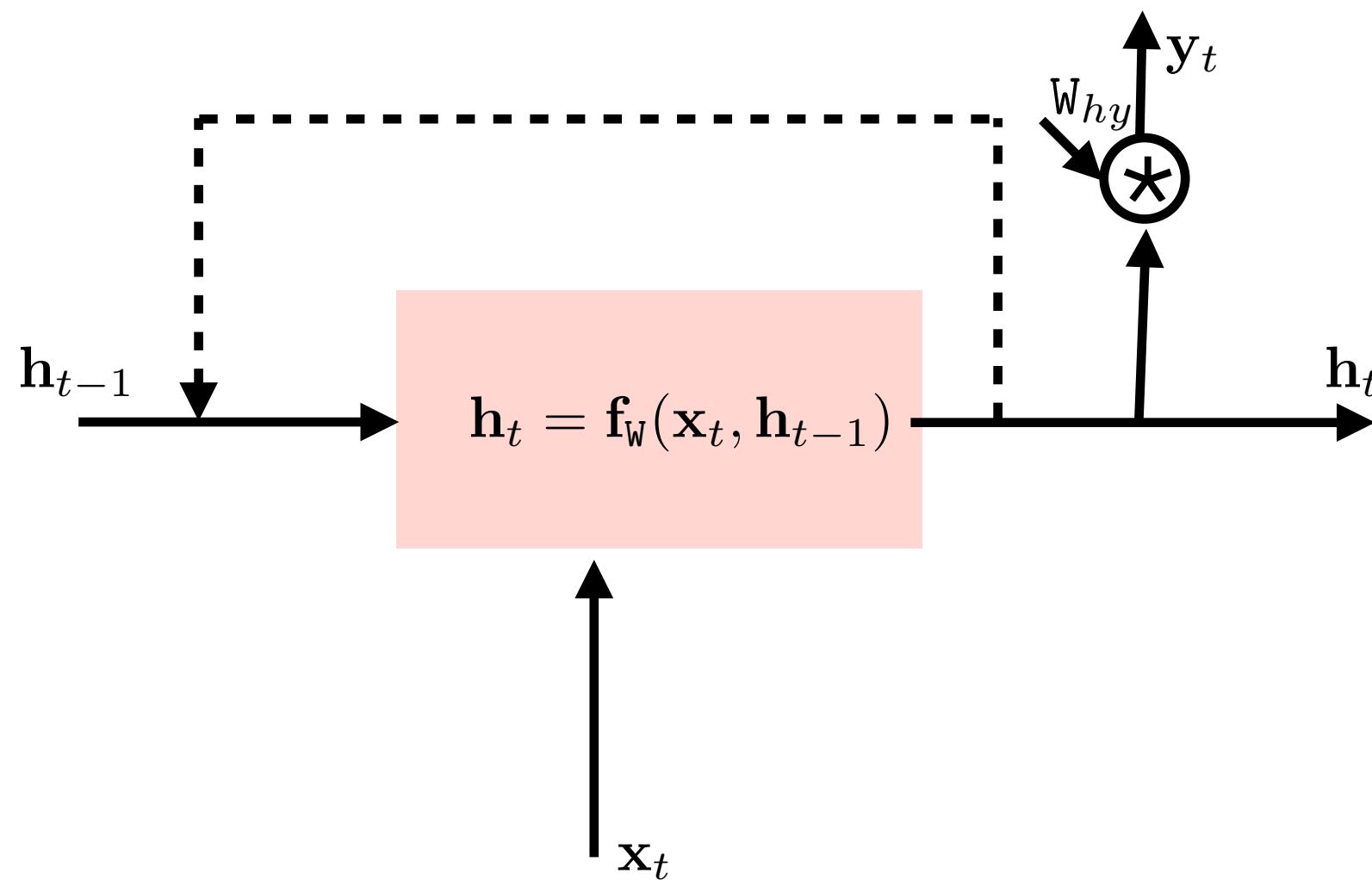
# Simple recurrent block - feed-forward pass



# Simple recurrent block - feed-forward pass



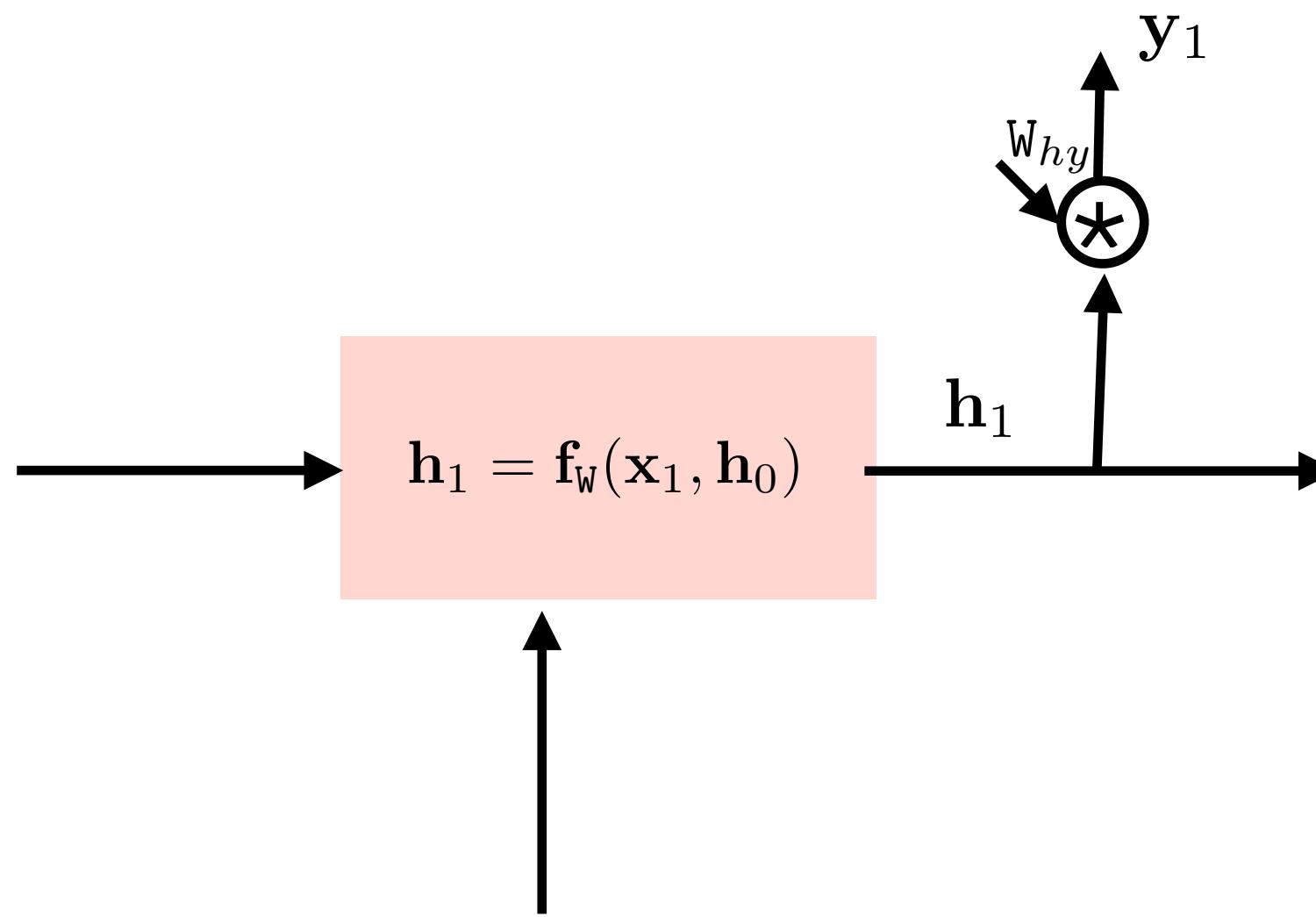
# Simple recurrent block - feed-forward pass



Given a finite input sequence:  $h_0 \ x_1 \ x_2 \ x_3$   
we remove the recurrent connection by:

- successive substitution of inputs and
- unrolling the net

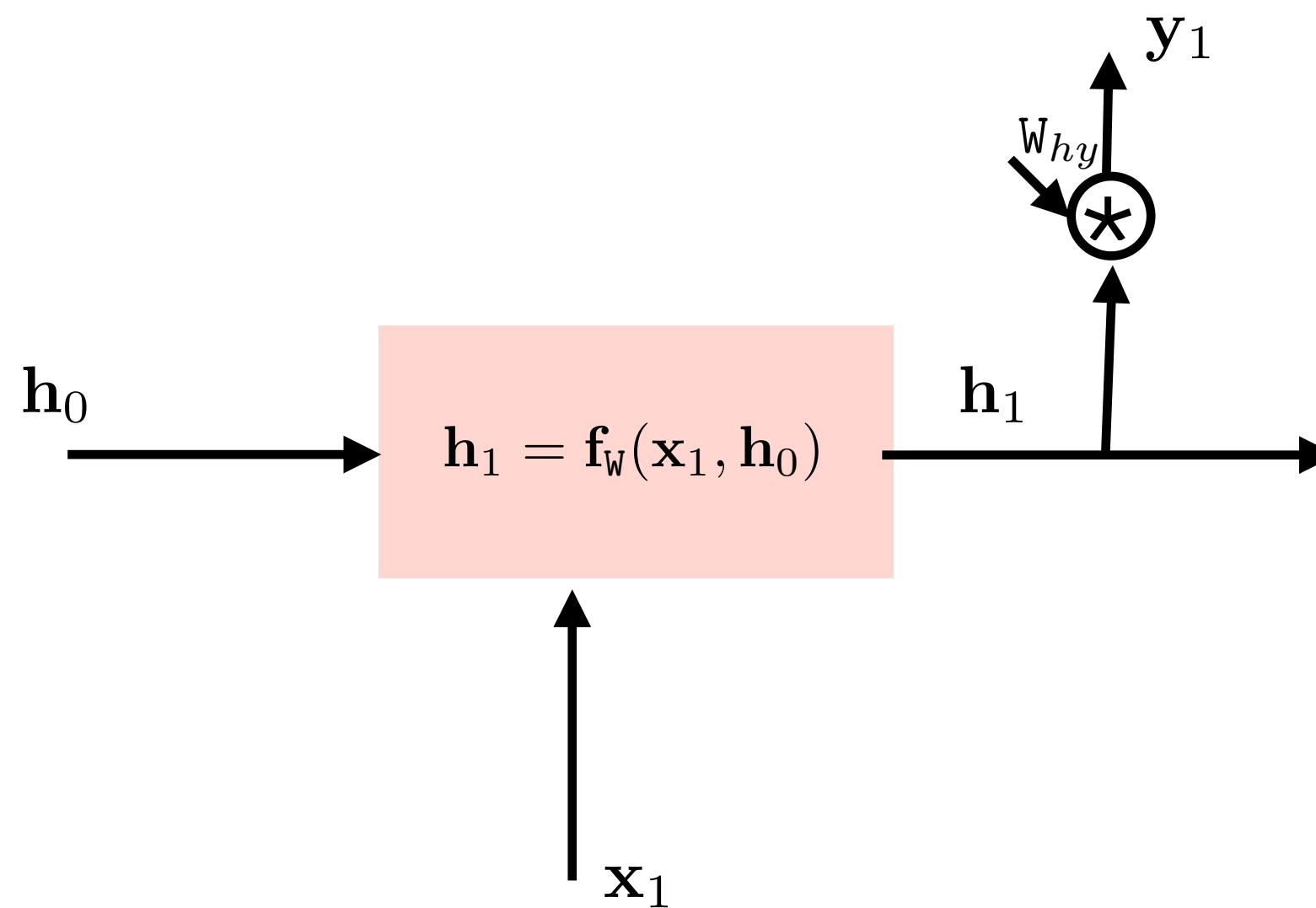
# Simple recurrent block - feed-forward pass



Given a finite input sequence:  $h_0 \ x_1 \ x_2 \ x_3$   
we remove the recurrent connection by:

- successive substitution of inputs and
- unrolling the net

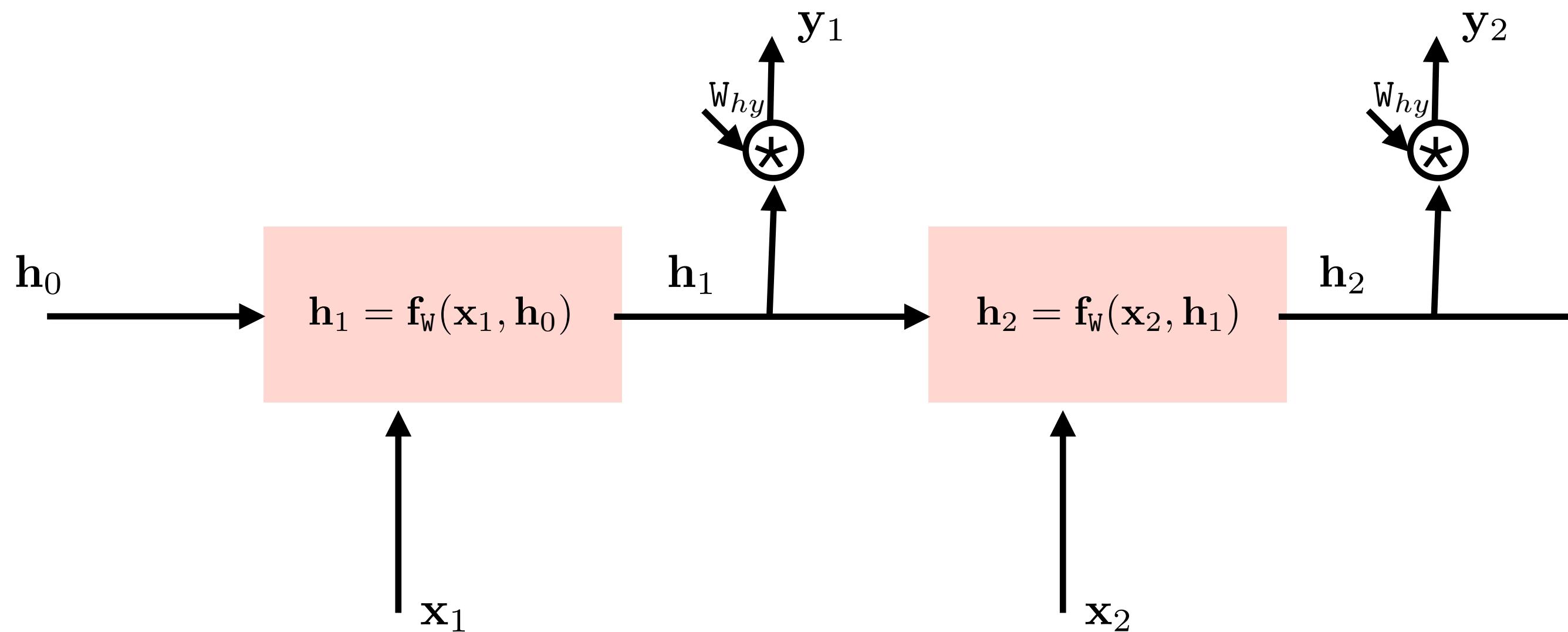
# Simple recurrent block - feed-forward pass



Given a finite input sequence:  $x_2 \ x_3$   
we remove the recurrent connection by:

- successive substitution of inputs and
- unrolling the net

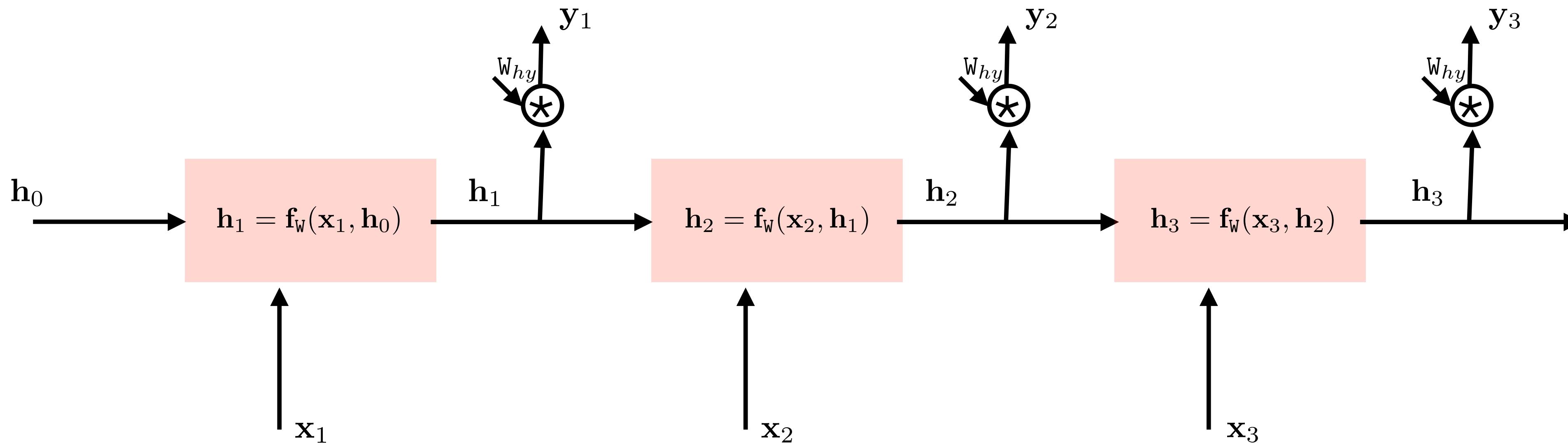
# Simple recurrent block - feed-forward pass



Given a finite input sequence:  
we remove the recurrent connection by:

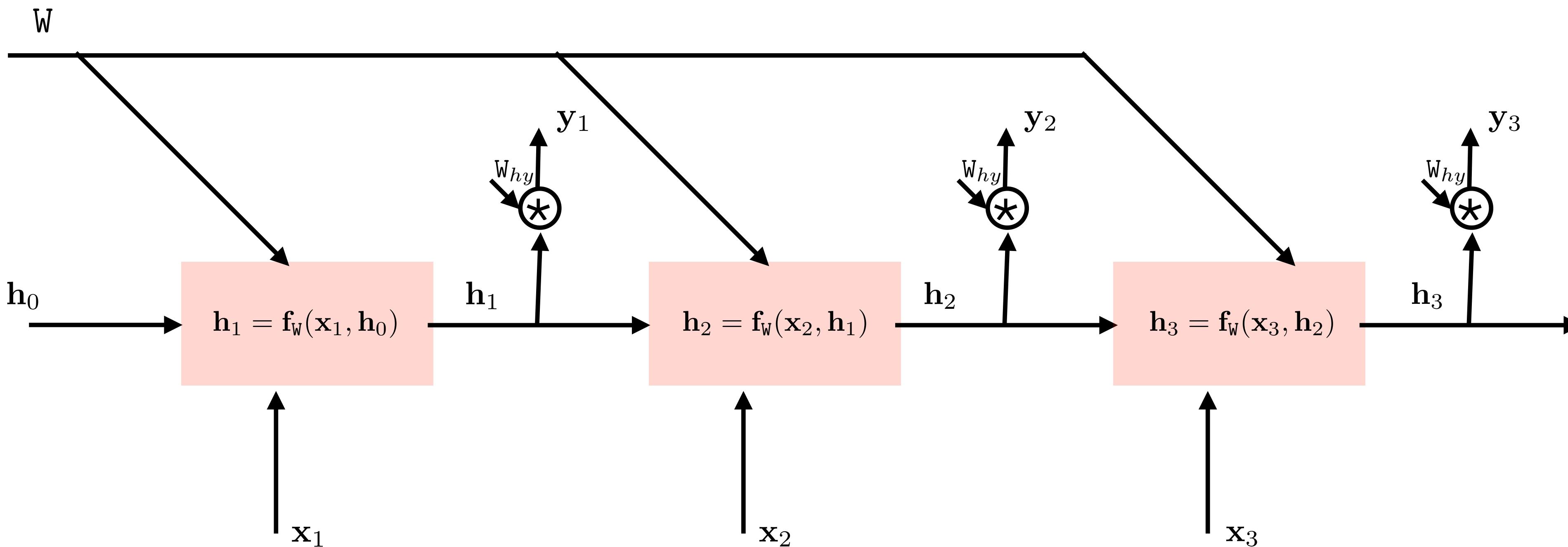
- successive substitution of inputs and
- unrolling the net

# Simple recurrent block - feed-forward pass



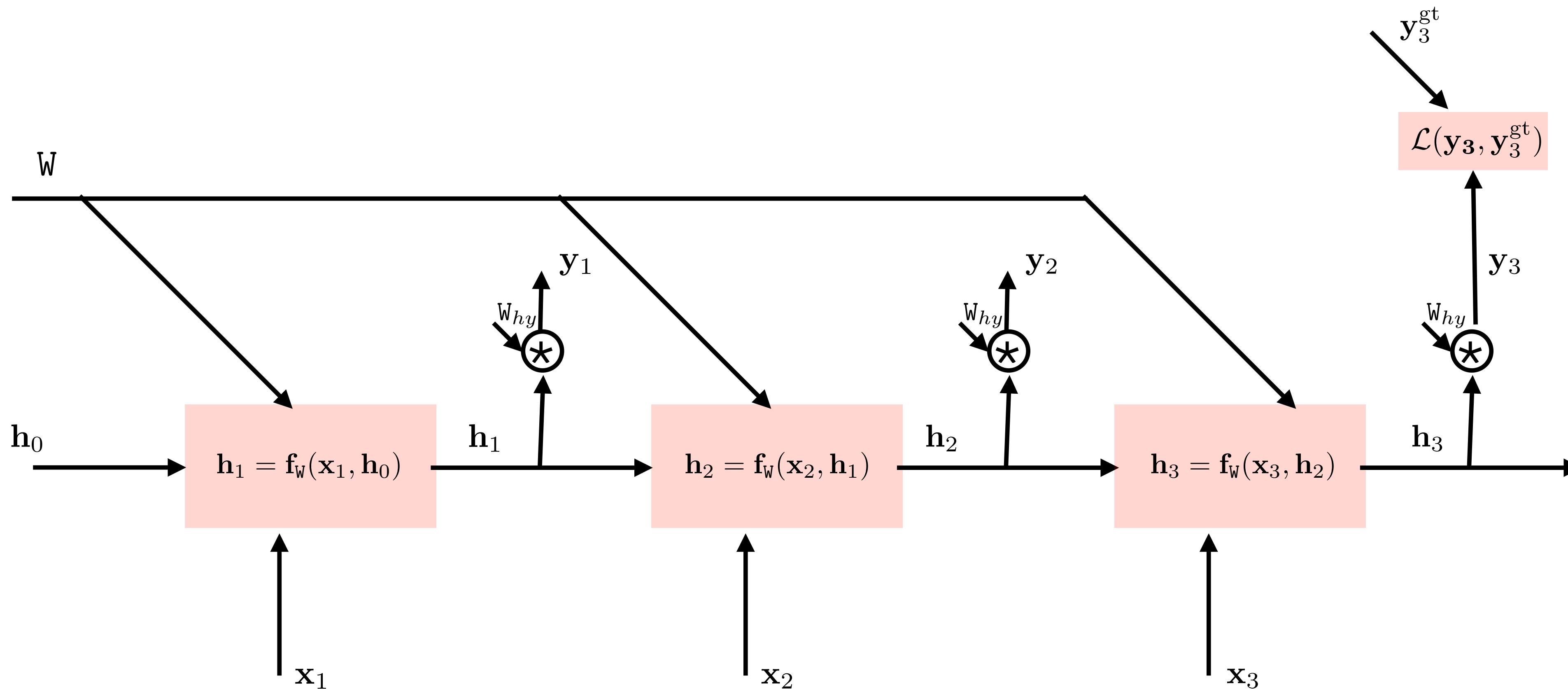
- **Unrolled** computational graph:
  - it is normal feedforward network

# Simple recurrent block - feed-forward pass



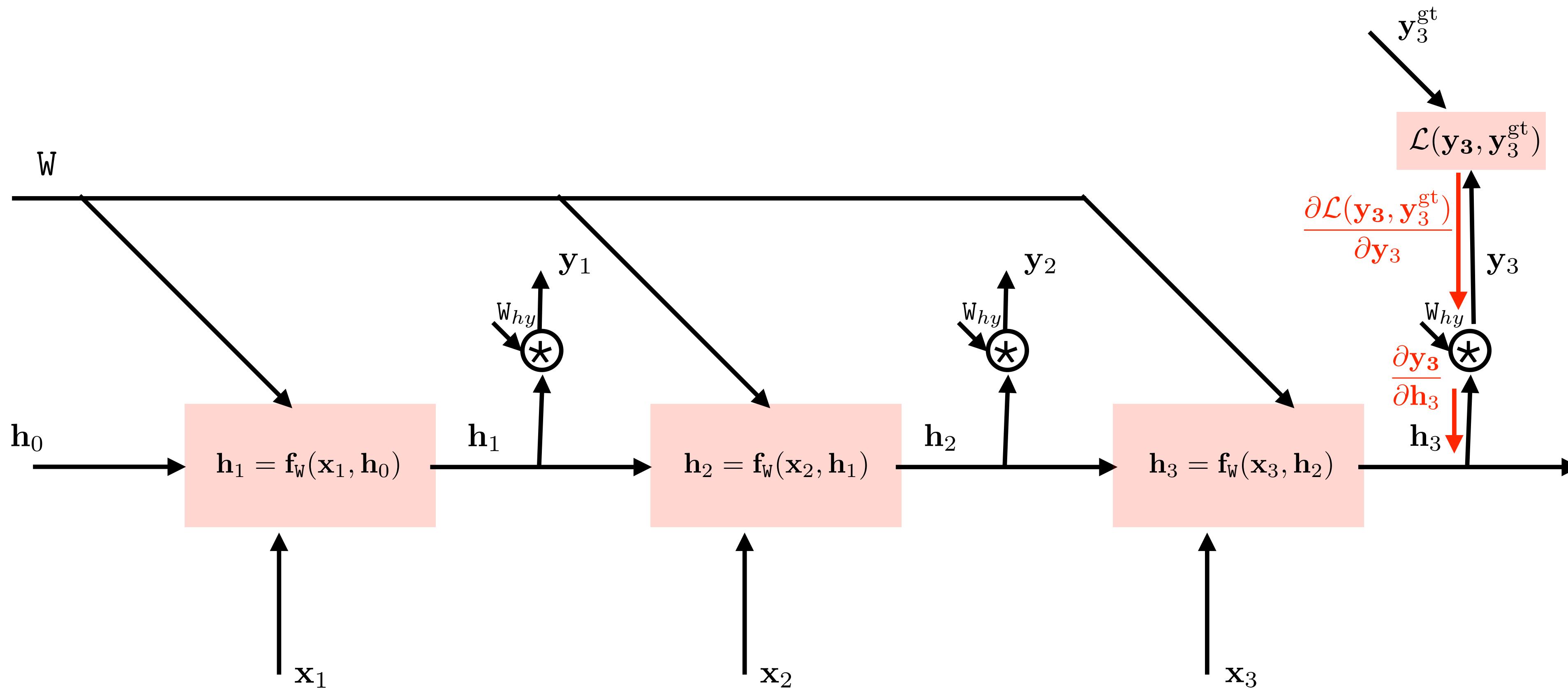
- **Unrolled** computational graph:
  - it is normal feedforward network
  - it consists of several **same blocks** with the **same weights!**

# Simple recurrent block - backward pass



- **Loss function** could be connected to all outputs in all times
- We connect the **loss** to the **final output** only (for the simplicity)

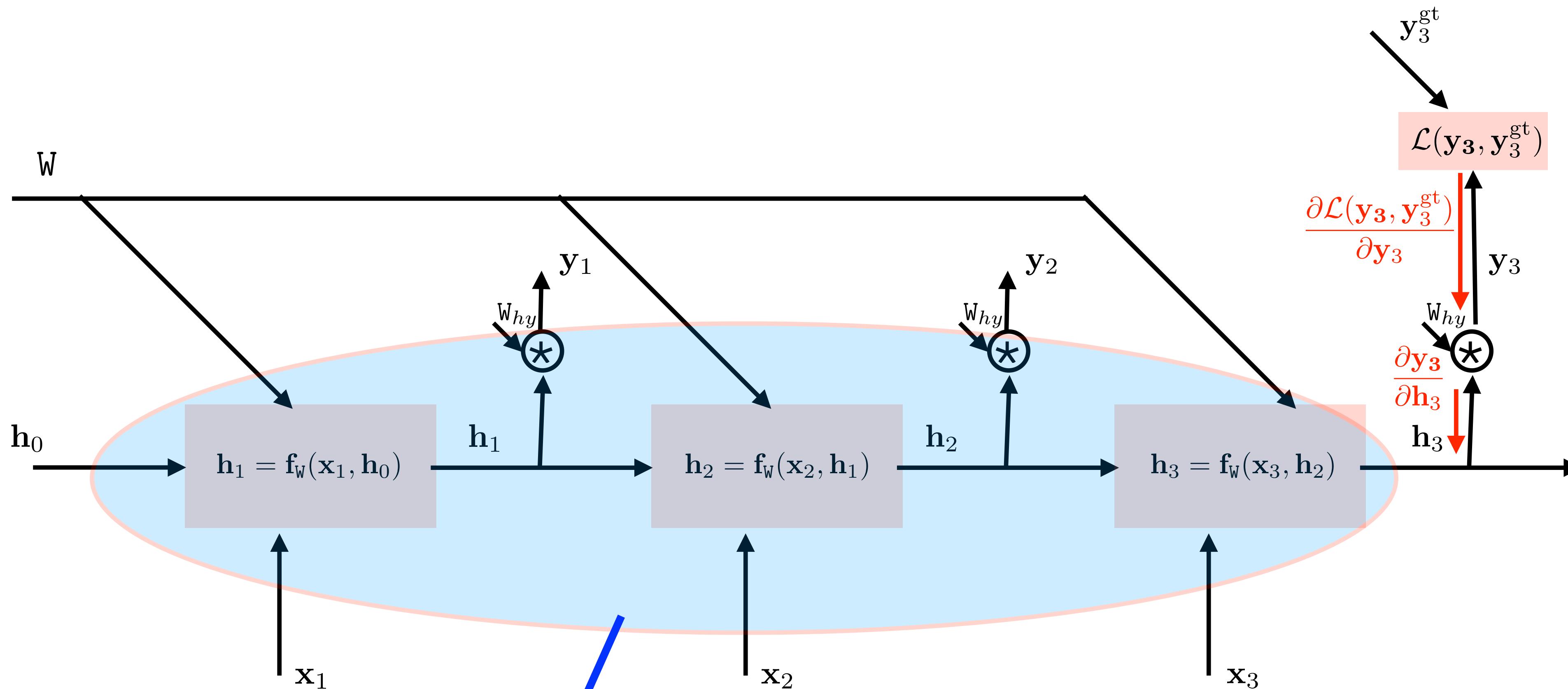
# Simple recurrent block - backward pass



- Backprop:

$$\frac{\partial \mathcal{L}(y_3, y_3^{gt})}{\partial w} = ?$$

# Simple recurrent block - backward pass

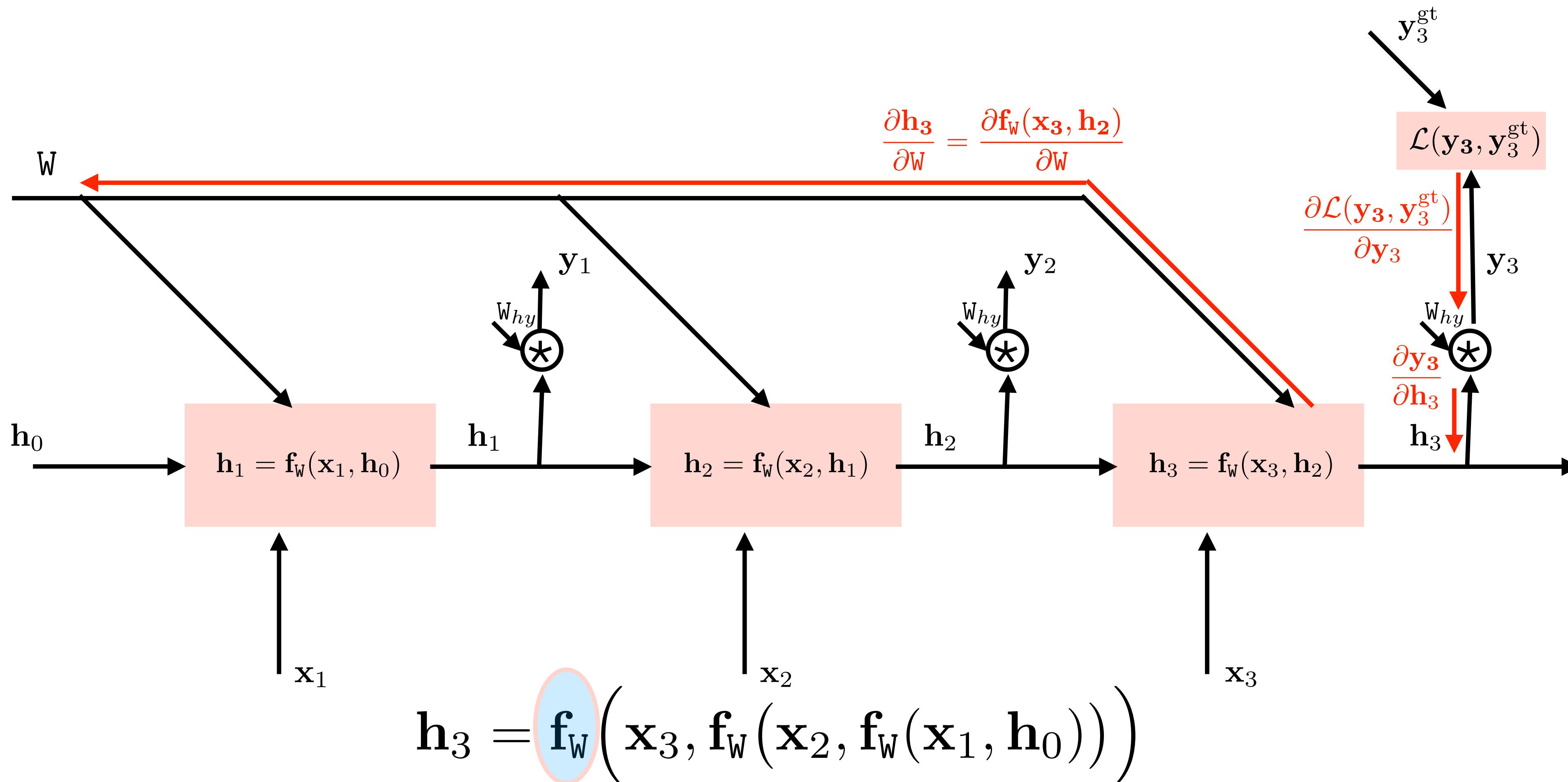


- Backprop:
  - differentiation of multi-dimensional composite function:

$$\frac{\partial \mathcal{L}(y_3, y_3^{gt})}{\partial w} = ?$$

$$h_3 = f_w(x_3, f_w(x_2, f_w(x_1, h_0)))$$

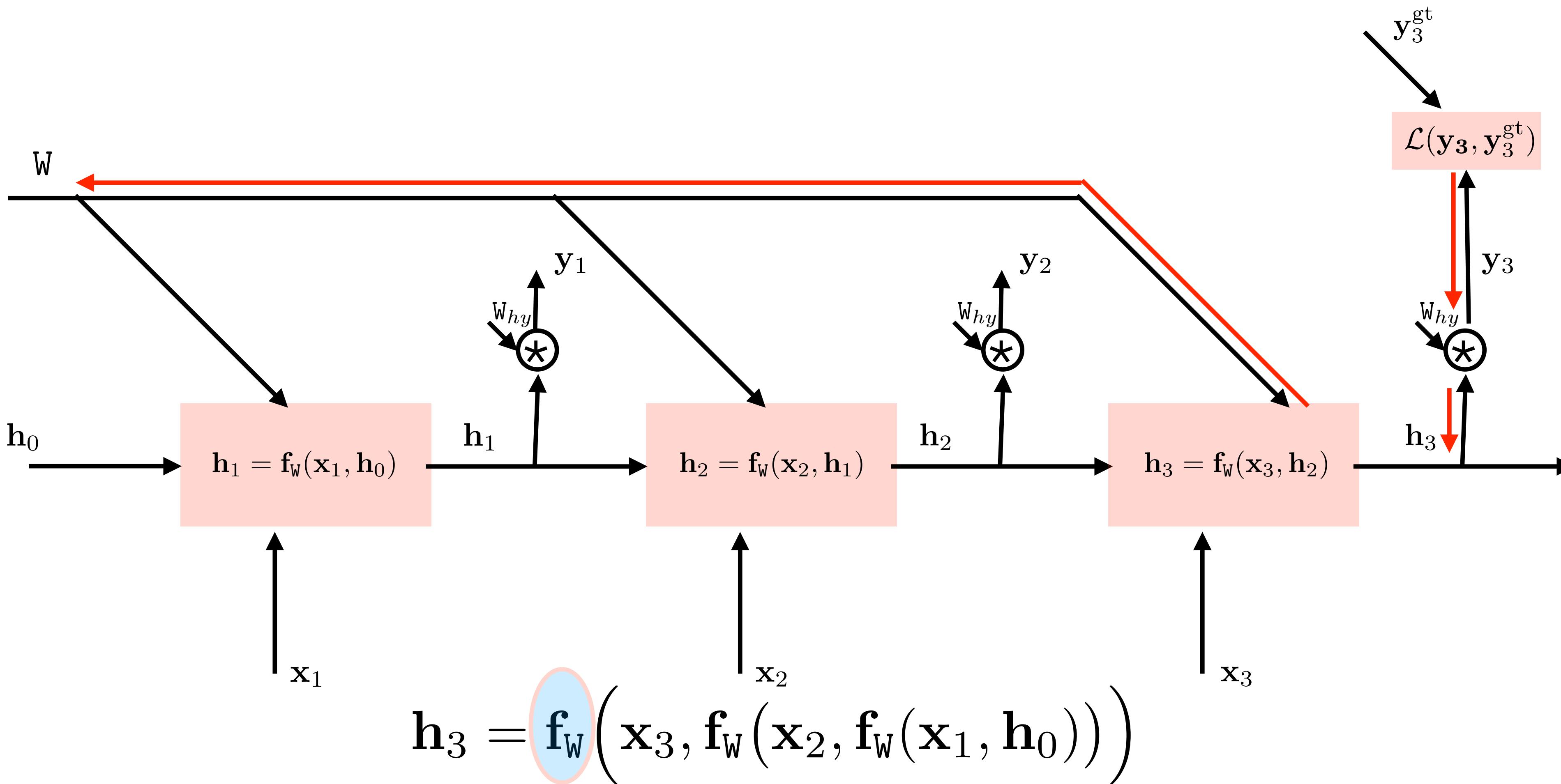
# Simple recurrent block - backward pass



$$\frac{\partial \mathcal{L}(\mathbf{y}_3, \mathbf{y}_3^{gt})}{\partial \mathbf{w}} = ?$$

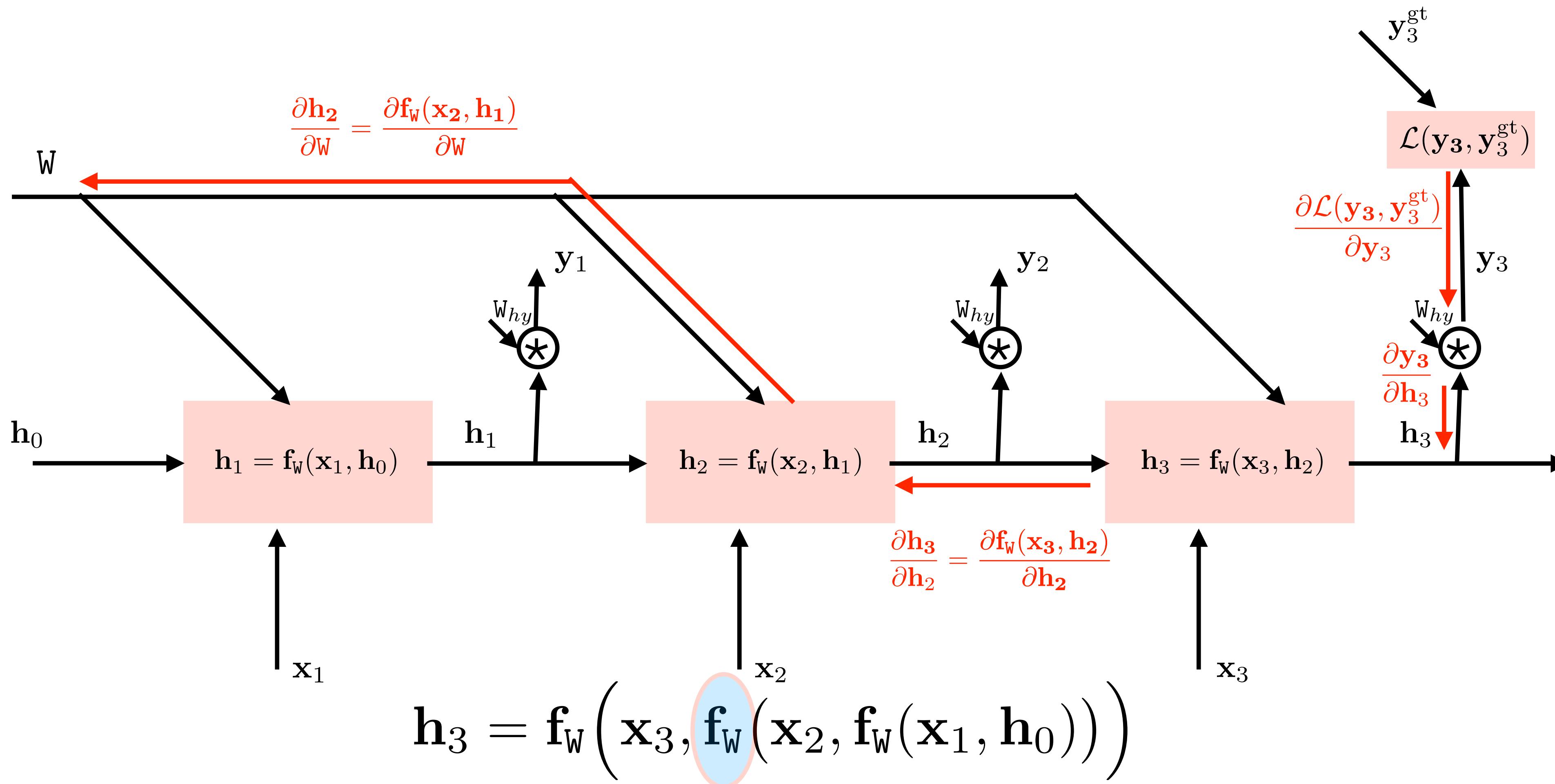
$$h_3 = f_w\left(x_3, f_w\left(x_2, f_w\left(x_1, h_0\right)\right)\right)$$

# Simple recurrent block - backward pass



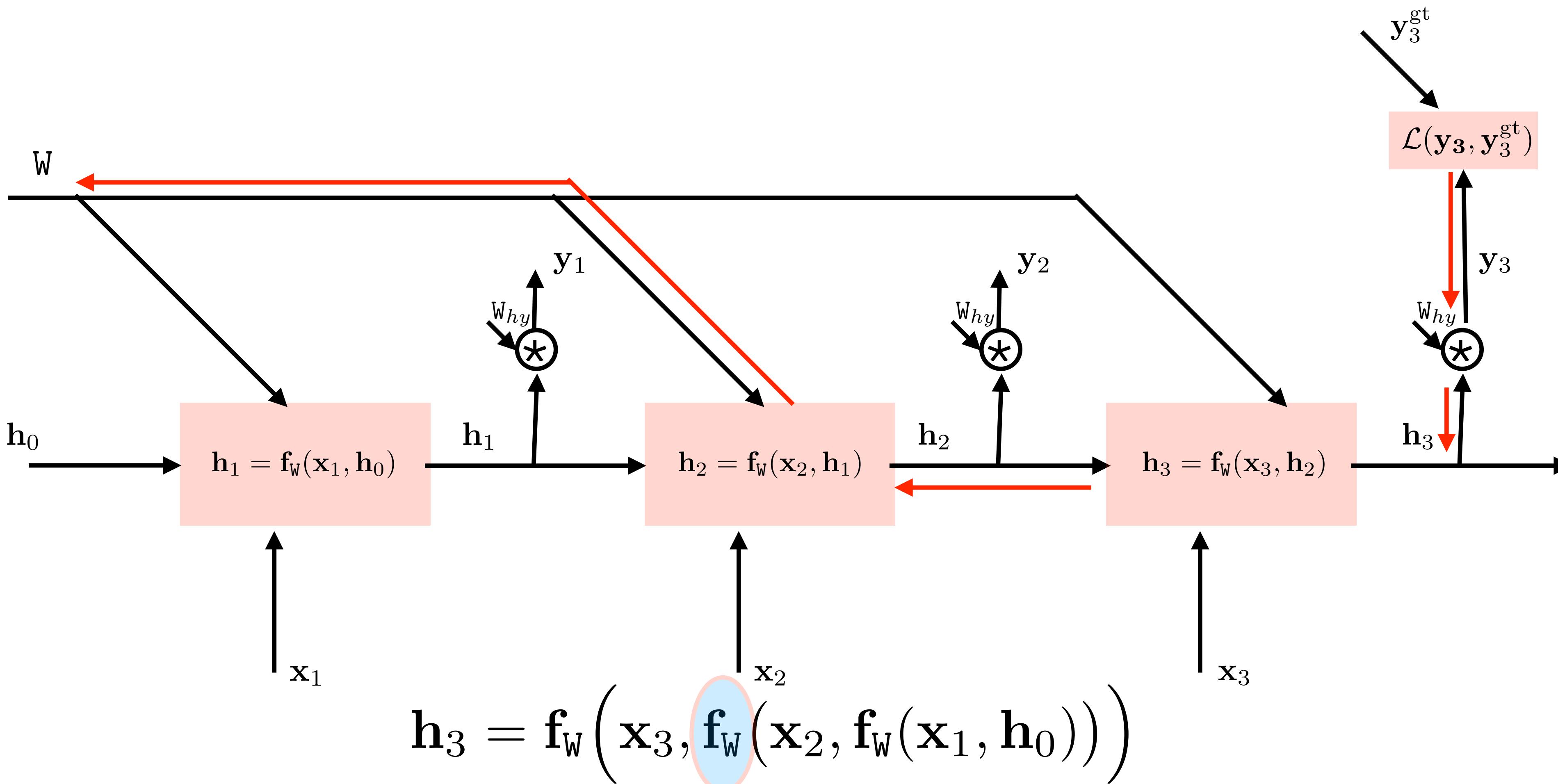
$$\frac{\partial \mathcal{L}(y_3, y_3^{gt})}{\partial W} = \frac{\partial \mathcal{L}(y_3, y_3^{gt})}{\partial y_3} \frac{\partial y_3}{\partial h_3} \frac{\partial f_w(x_3, h_2)}{\partial W} +$$

# Simple recurrent block - backward pass



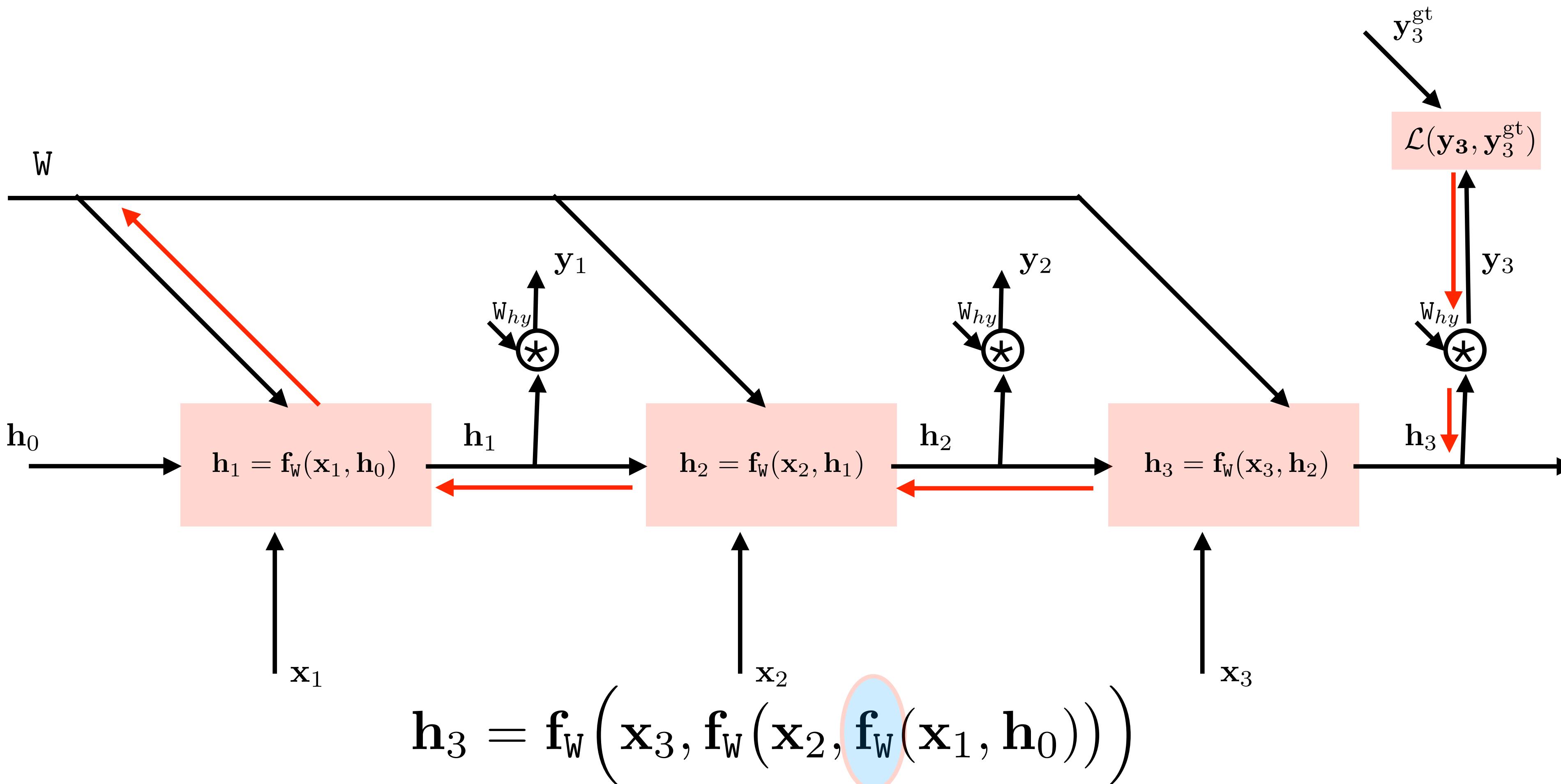
$$\frac{\partial \mathcal{L}(y_3, y_3^{gt})}{\partial W} = \frac{\partial \mathcal{L}(y_3, y_3^{gt})}{\partial y_3} \frac{\partial y_3}{\partial h_3} \frac{\partial f_w(x_3, h_2)}{\partial W} +$$

# Simple recurrent block - backward pass



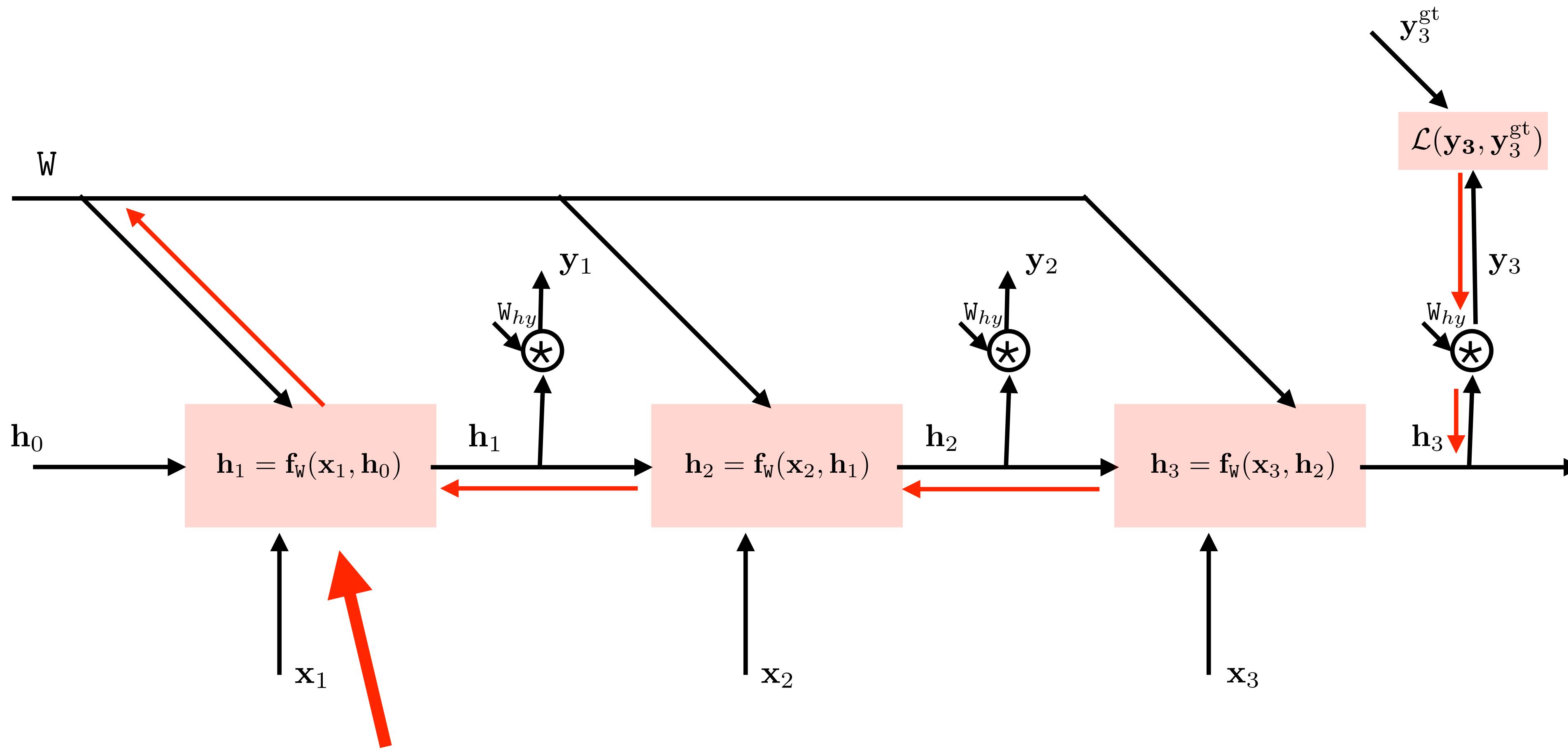
$$\frac{\partial \mathcal{L}(y_3, y_3^{\text{gt}})}{\partial W} = \frac{\partial \mathcal{L}(y_3, y_3^{\text{gt}})}{\partial y_3} \frac{\partial y_3}{\partial h_3} \frac{\partial f_w(x_3, h_2)}{\partial W} + \frac{\partial \mathcal{L}(y_3, y_3^{\text{gt}})}{\partial y_3} \frac{\partial y_3}{\partial h_3} \frac{\partial f_w(x_3, h_2)}{\partial h_2} \frac{\partial f_w(x_2, h_1)}{\partial W} +$$

# Simple recurrent block - backward pass



$$\frac{\partial \mathcal{L}(y_3, y_3^{gt})}{\partial w} = \frac{\partial \mathcal{L}(y_3, y_3^{gt})}{\partial y_3} \frac{\partial y_3}{\partial h_3} \frac{\partial f_w(x_3, h_2)}{\partial w} + \frac{\partial \mathcal{L}(y_3, y_3^{gt})}{\partial y_3} \frac{\partial y_3}{\partial h_3} \frac{\partial f_w(x_3, h_2)}{\partial h_2} \frac{\partial f_w(x_2, h_1)}{\partial w} + \dots$$

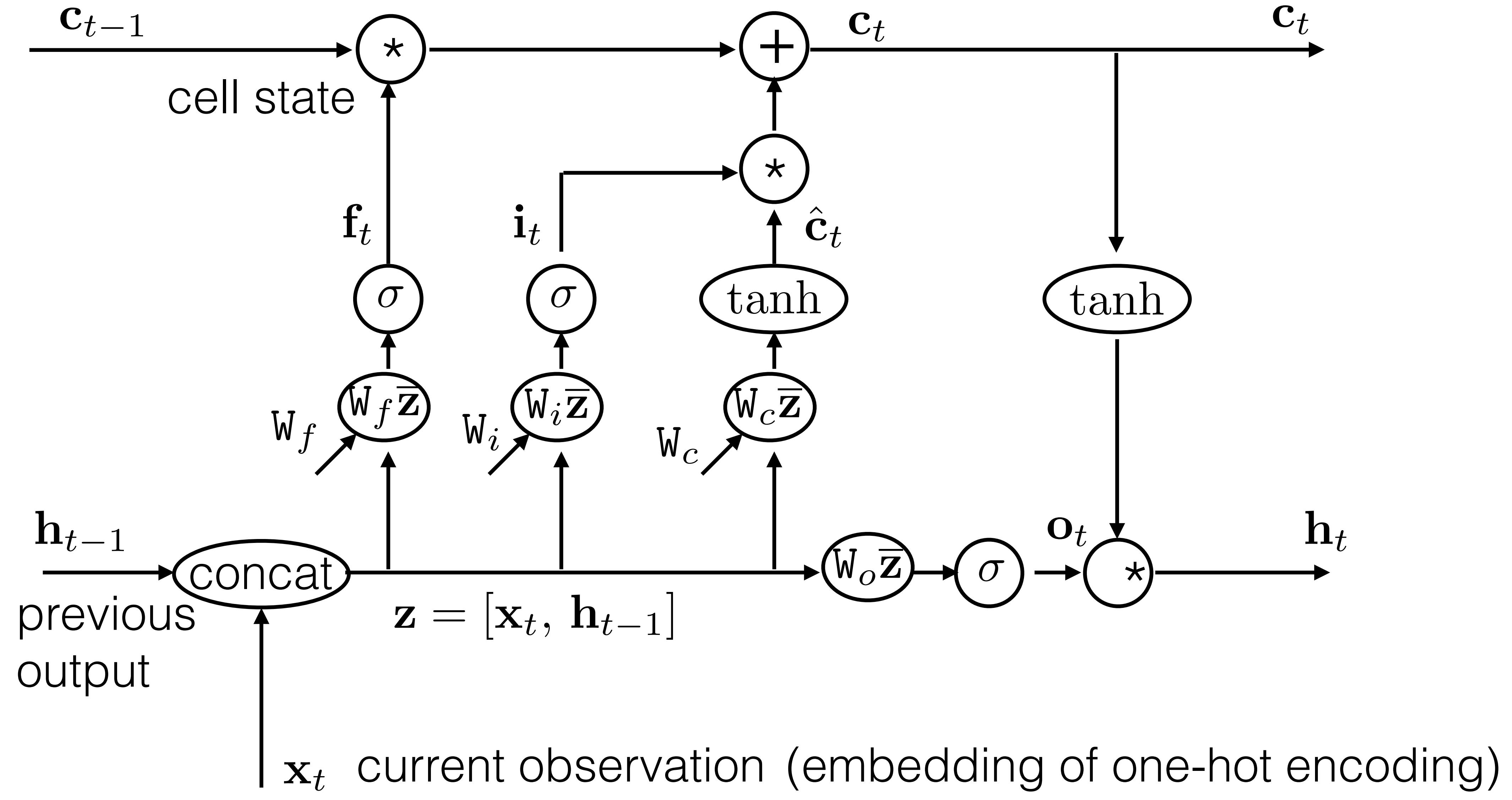
# Simple recurrent block - backward pass



deep blocks often suffer from vanishing gradient  
=> better structure needed

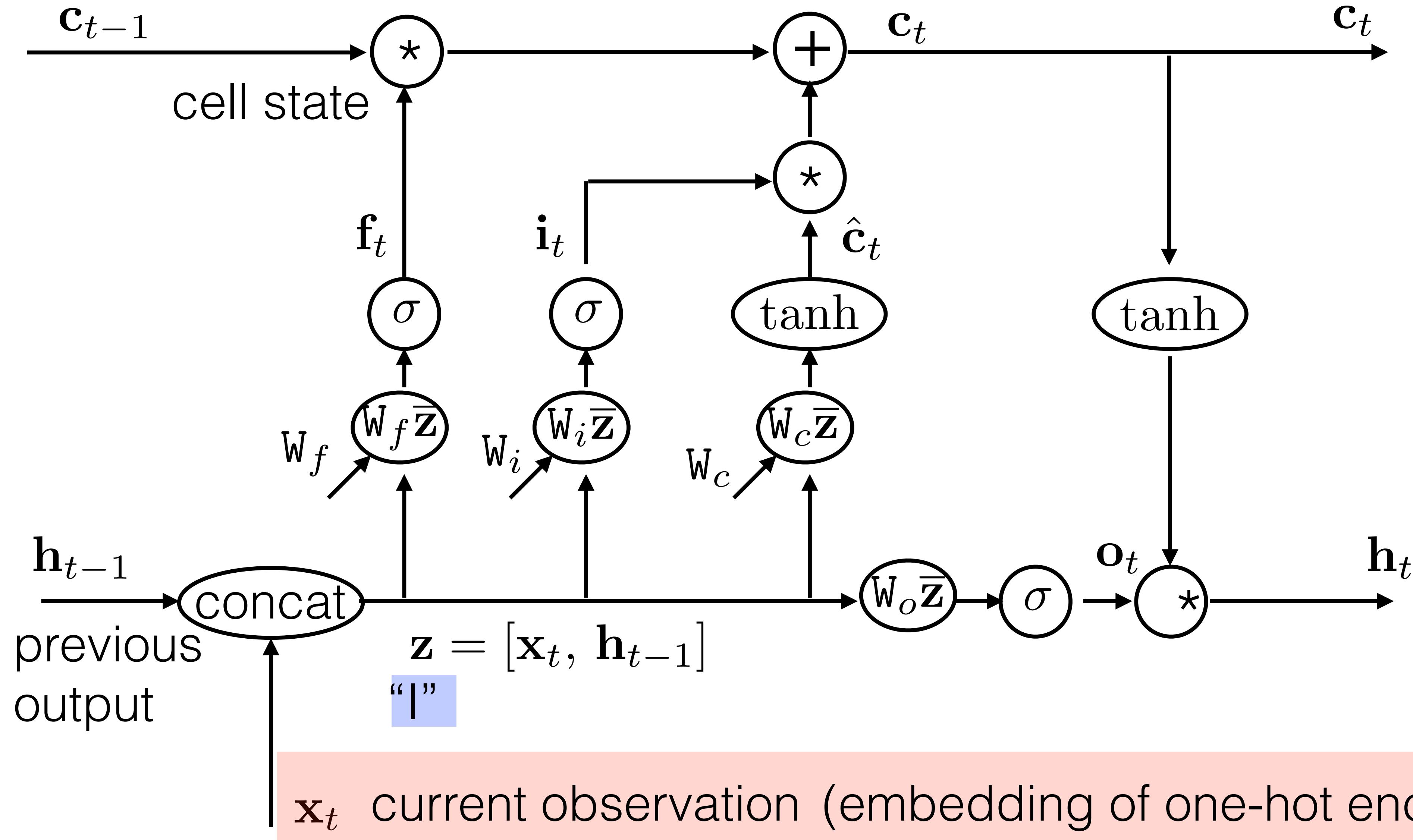
LSTM (kind of ResNet for recurrent networks)

# LSTM block



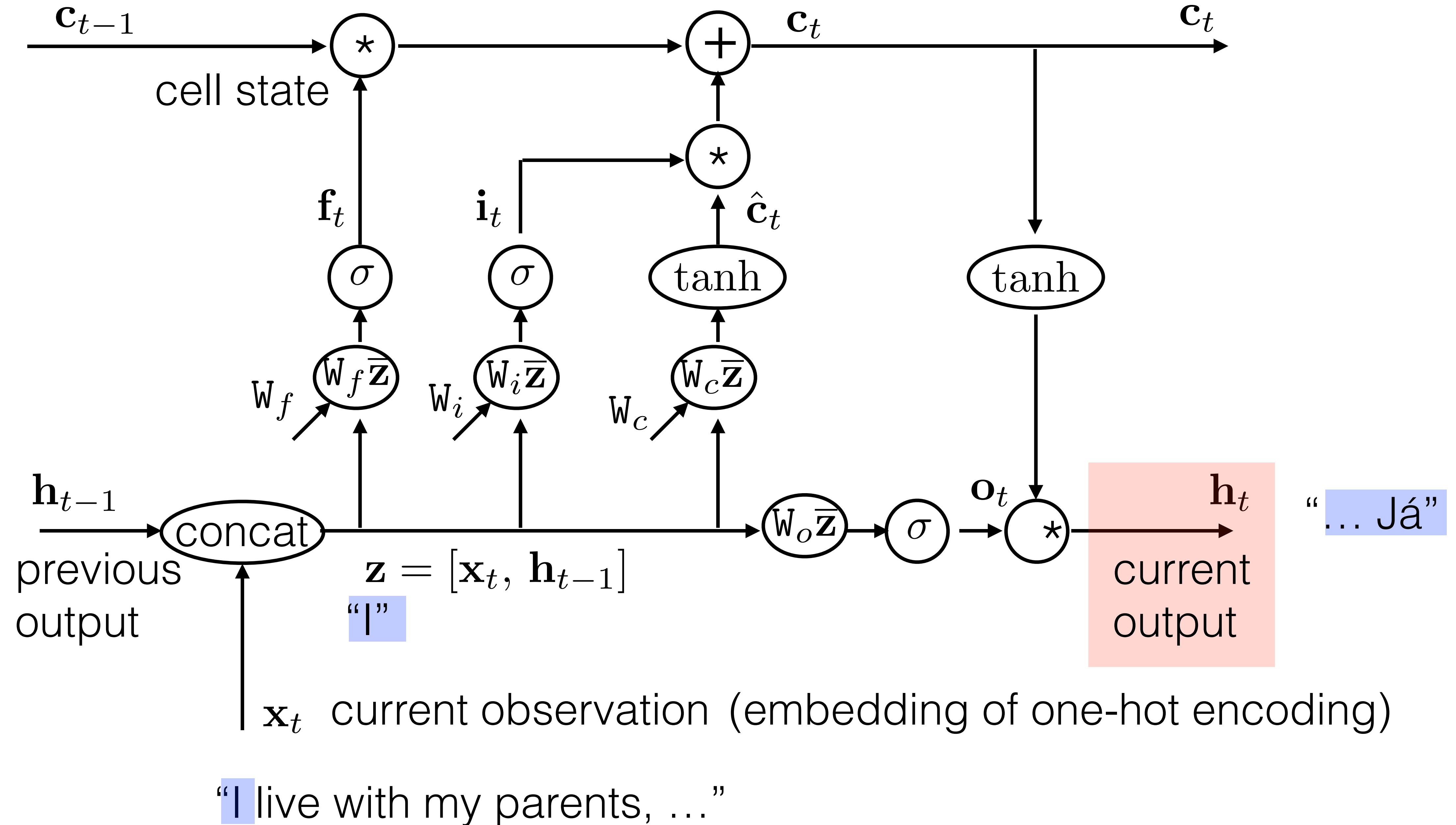
“I live with my parents, ...”

# LSTM block



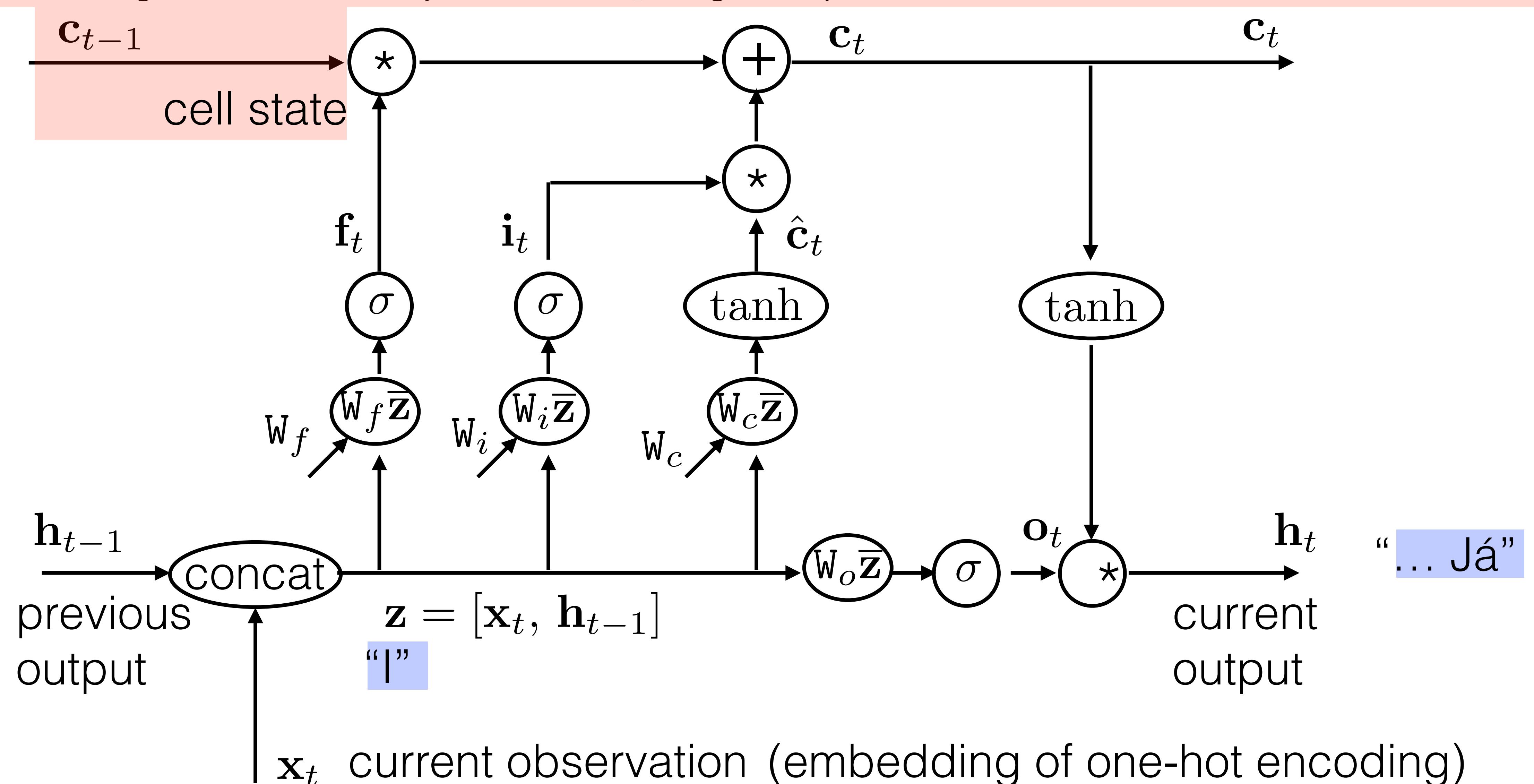
“I live with my parents, ...”

# LSTM block



## LSTM block

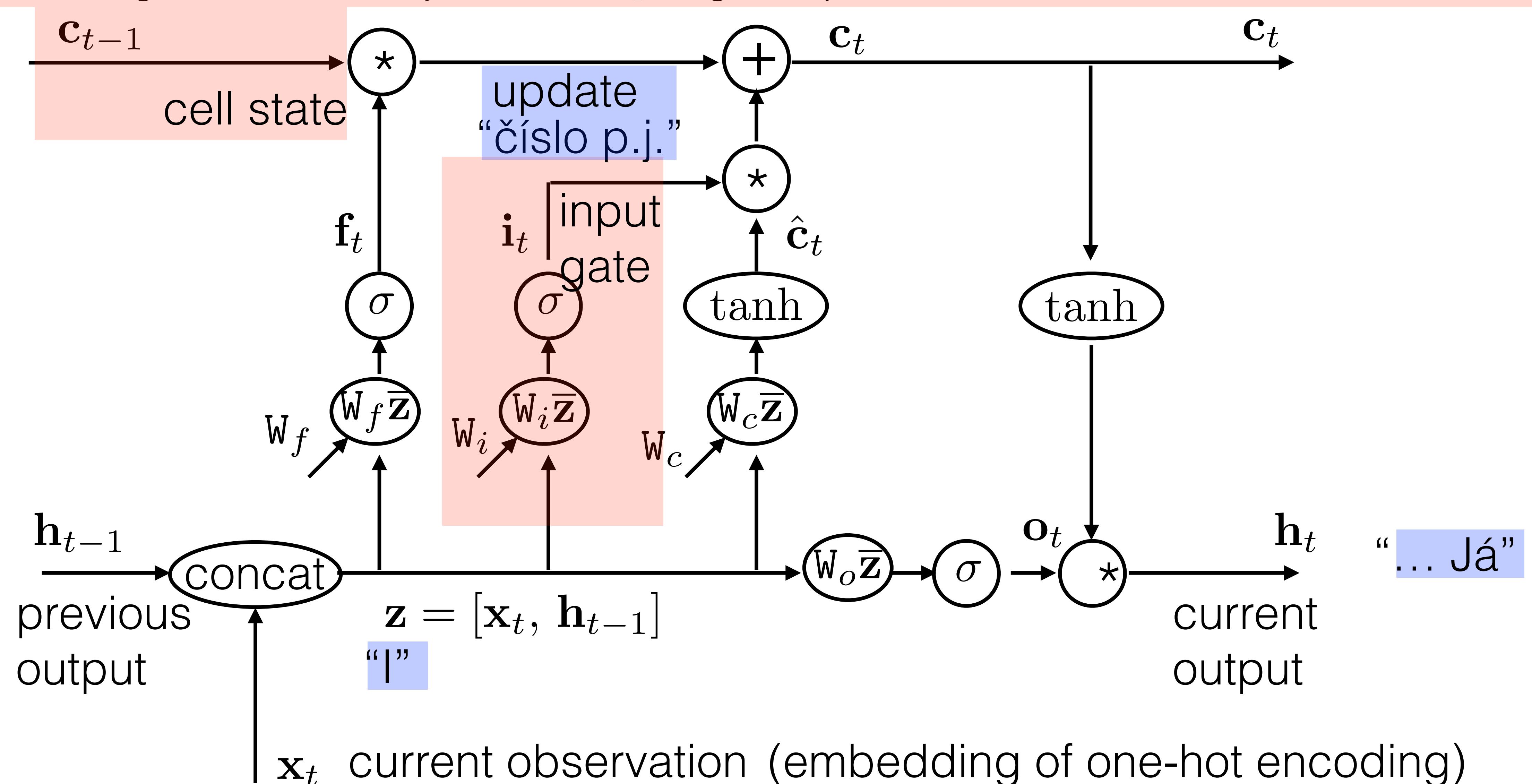
cell state is long term memory= vector [singular/plural, feminine/masculine, case, ...]



“I live with my parents, ...”

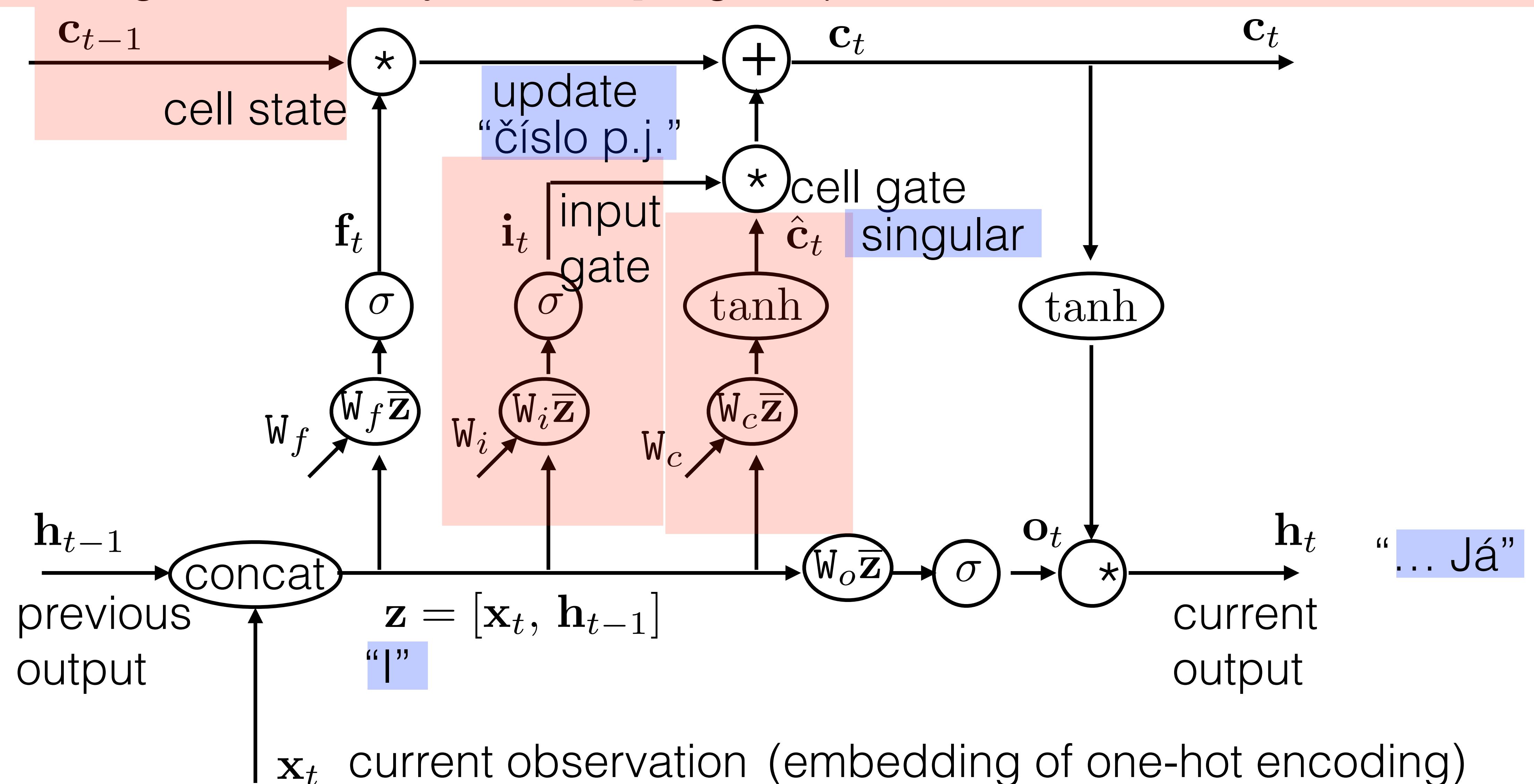
# LSTM block

cell state is long term memory= vector [singular/plural, feminine/masculine, case, ...]



# LSTM block

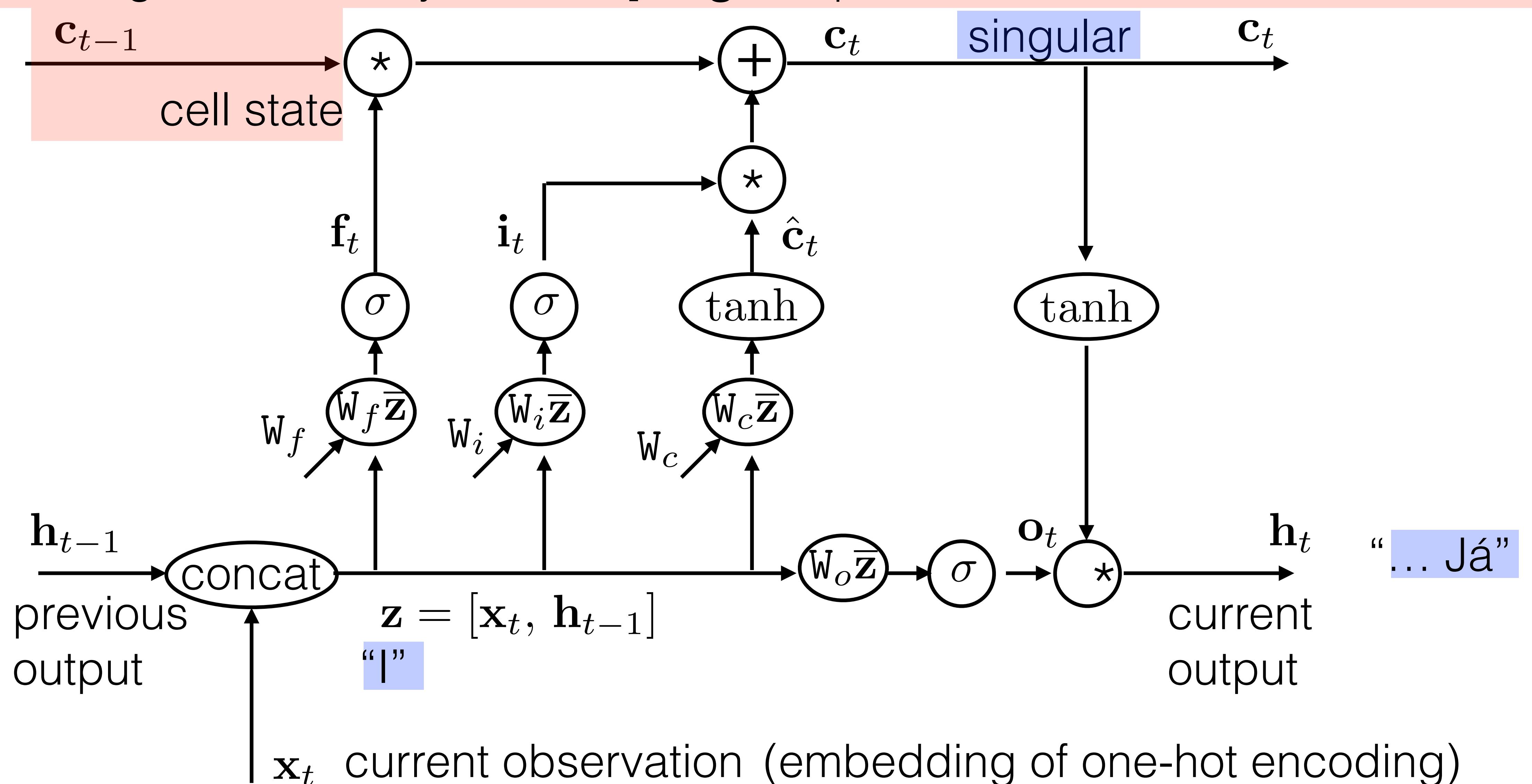
cell state is long term memory= vector [singular/plural, feminine/masculine, case, ...]



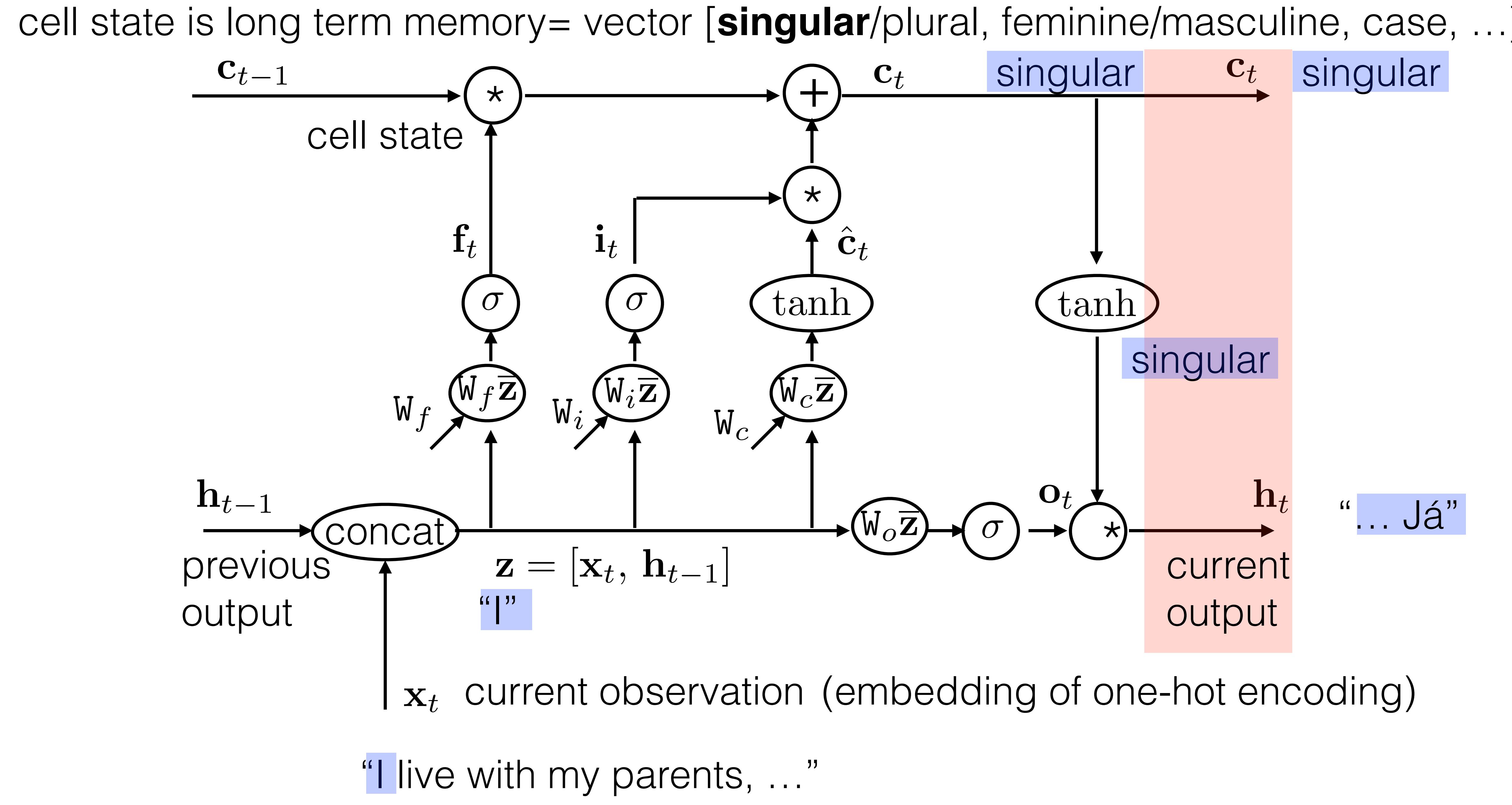
"I live with my parents, ..."

## LSTM block

cell state is long term memory= vector [singular/plural, feminine/masculine, case, ...]

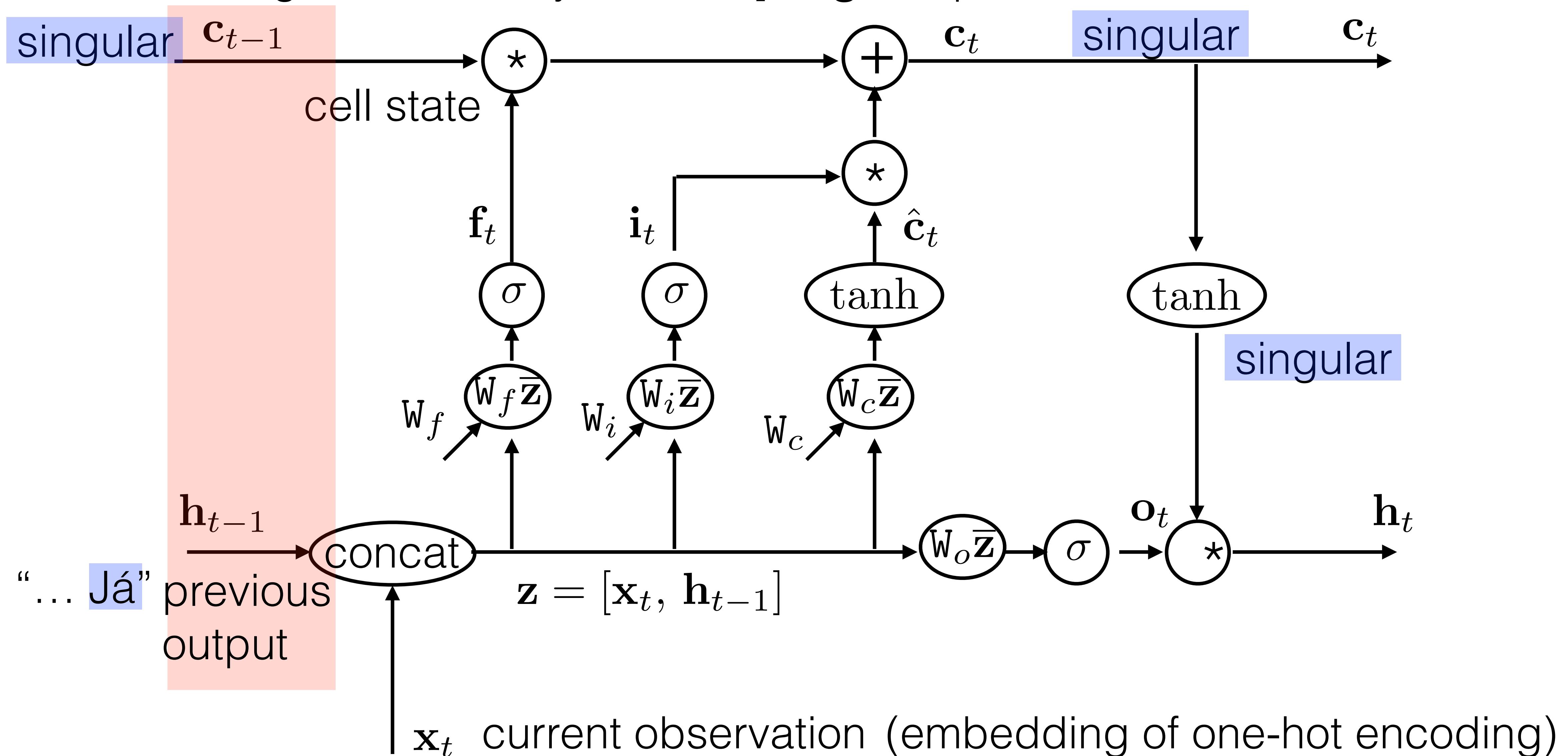


## LSTM block



## LSTM block

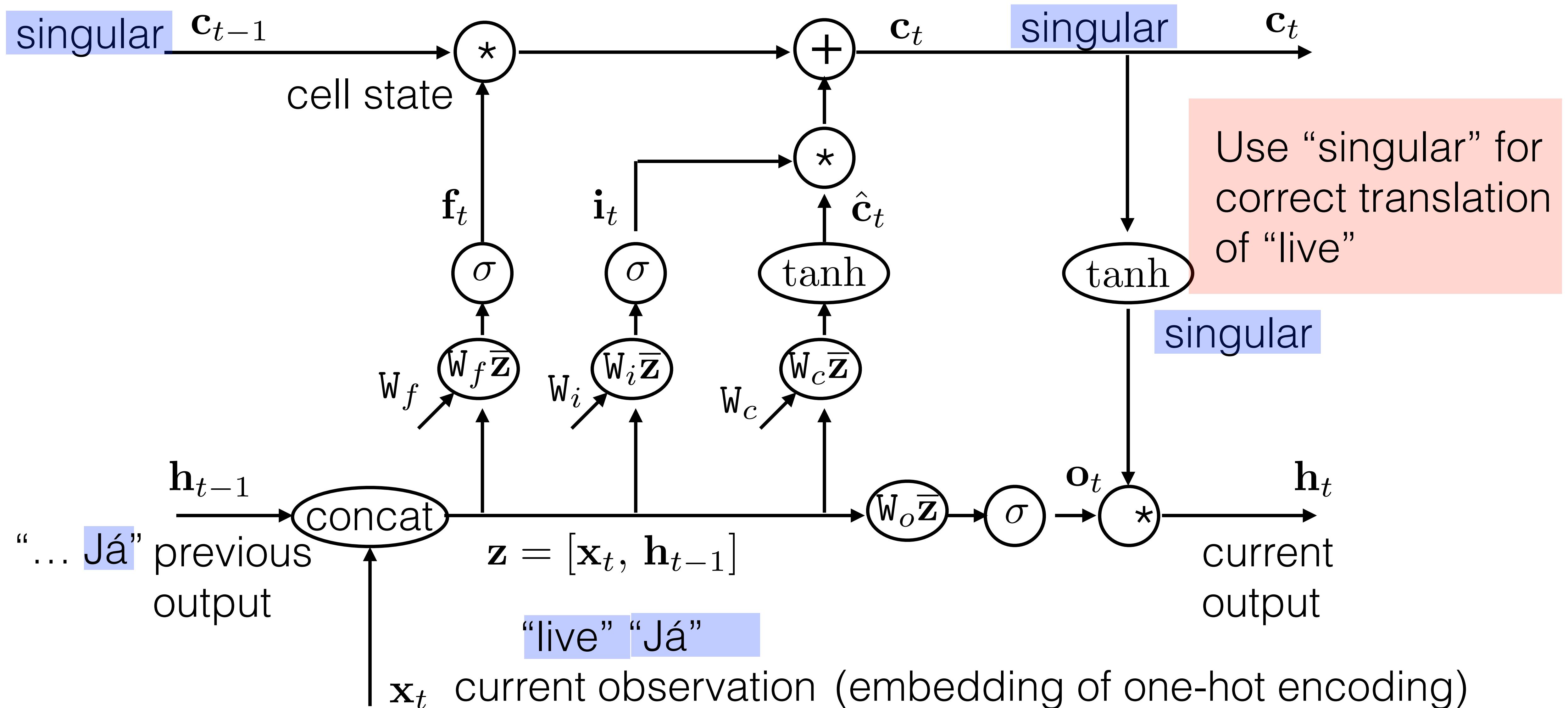
cell state is long term memory= vector [singular/plural, feminine/masculine, case, ...]



"I live with my parents, ..."

## LSTM block

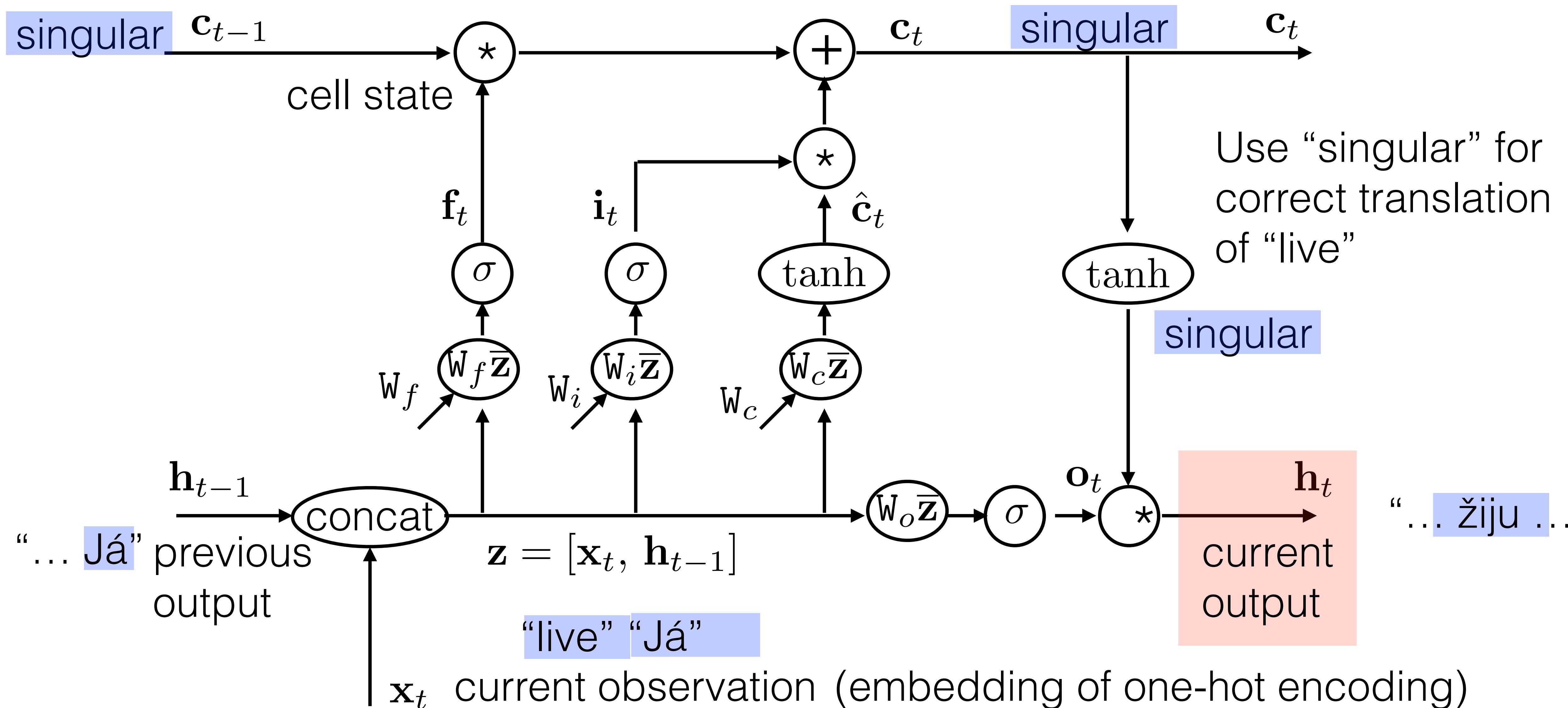
cell state is long term memory= vector [singular/plural, feminine/masculine, case, ...]



“I live with my parents, ...”

## LSTM block

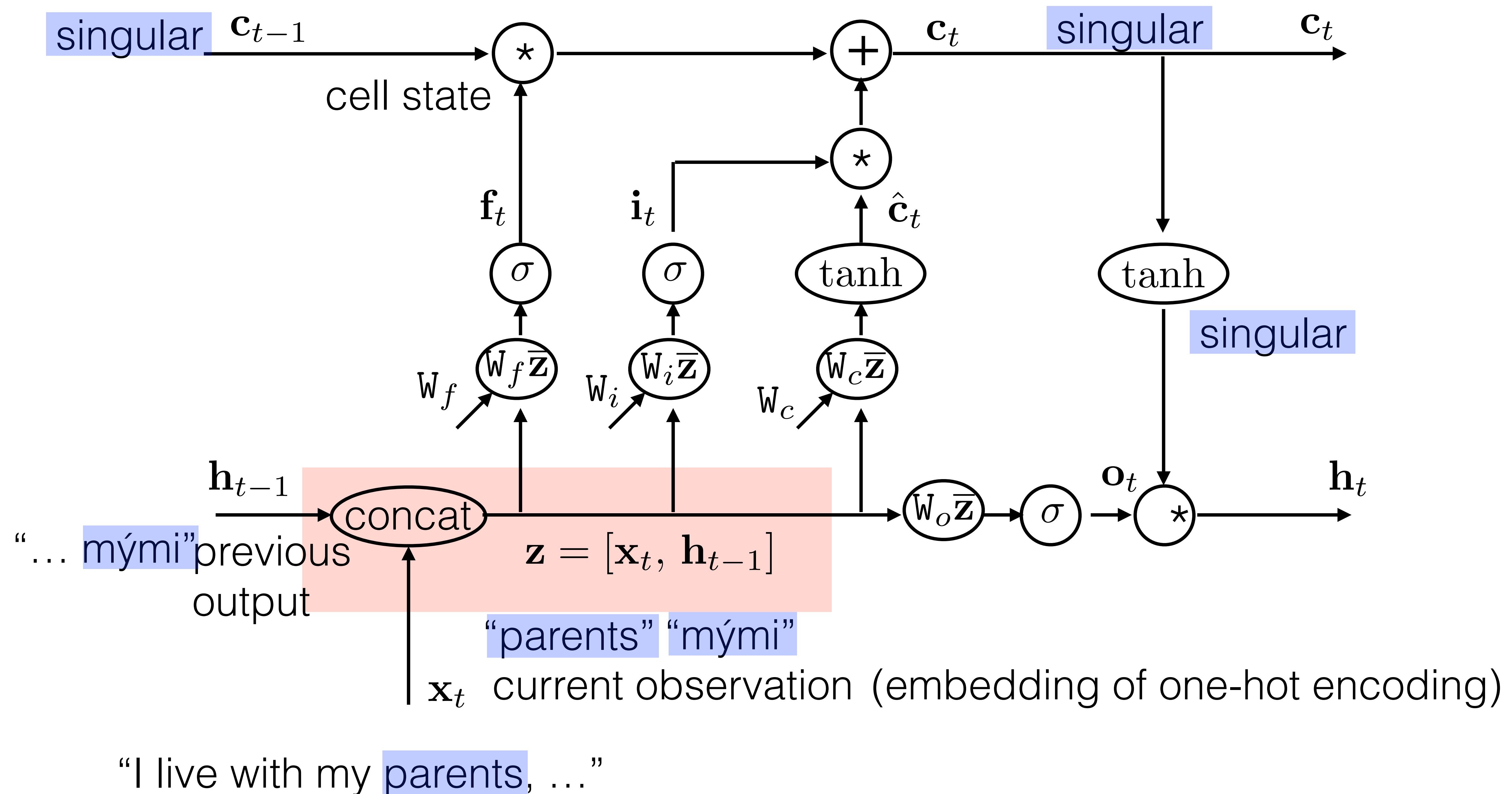
cell state is long term memory= vector [singular/plural, feminine/masculine, case, ...]



“I live with my parents, ...”

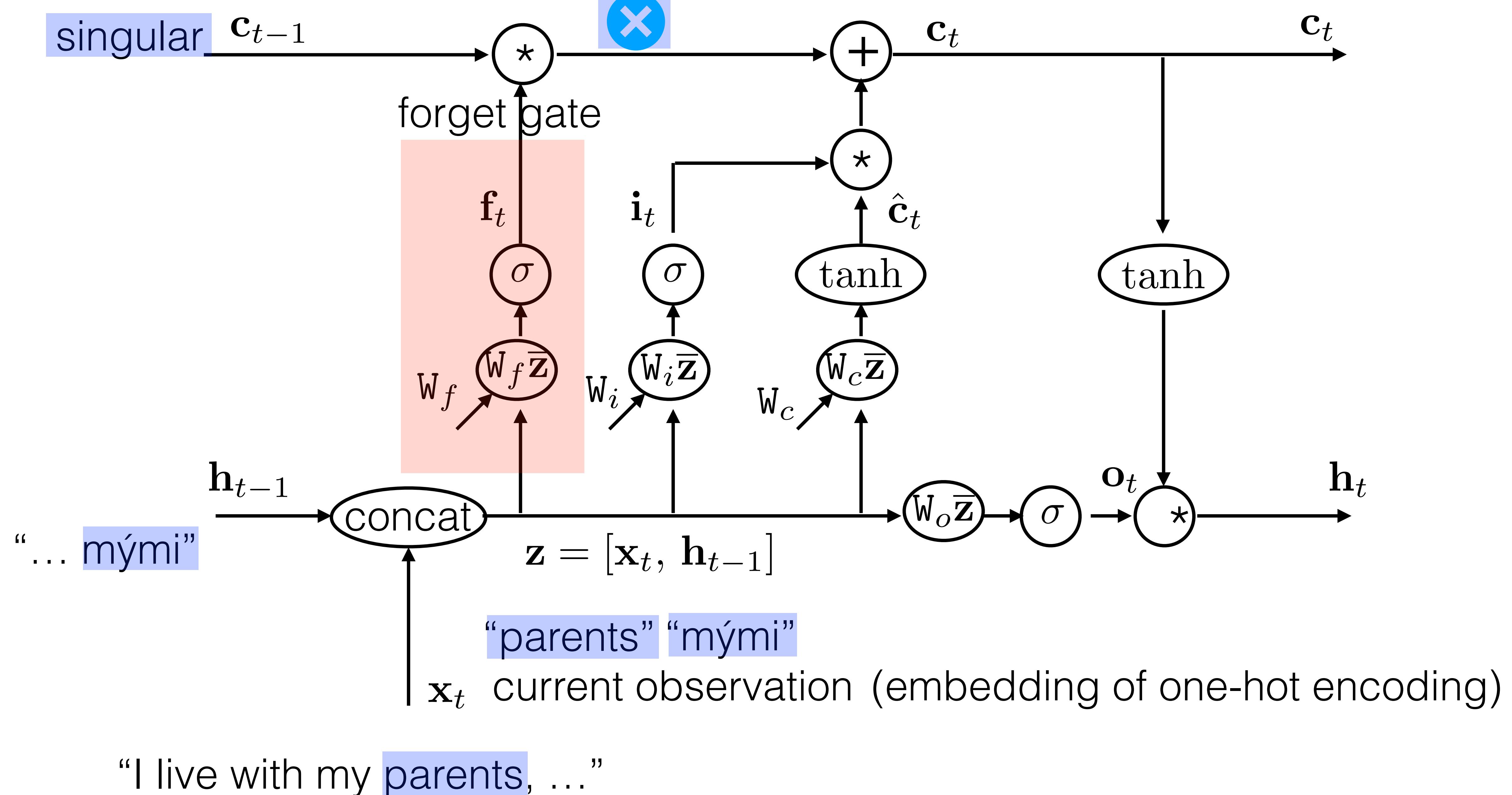
## LSTM block

cell state is long term memory= vector [singular/plural, feminine/masculine, case, ...]



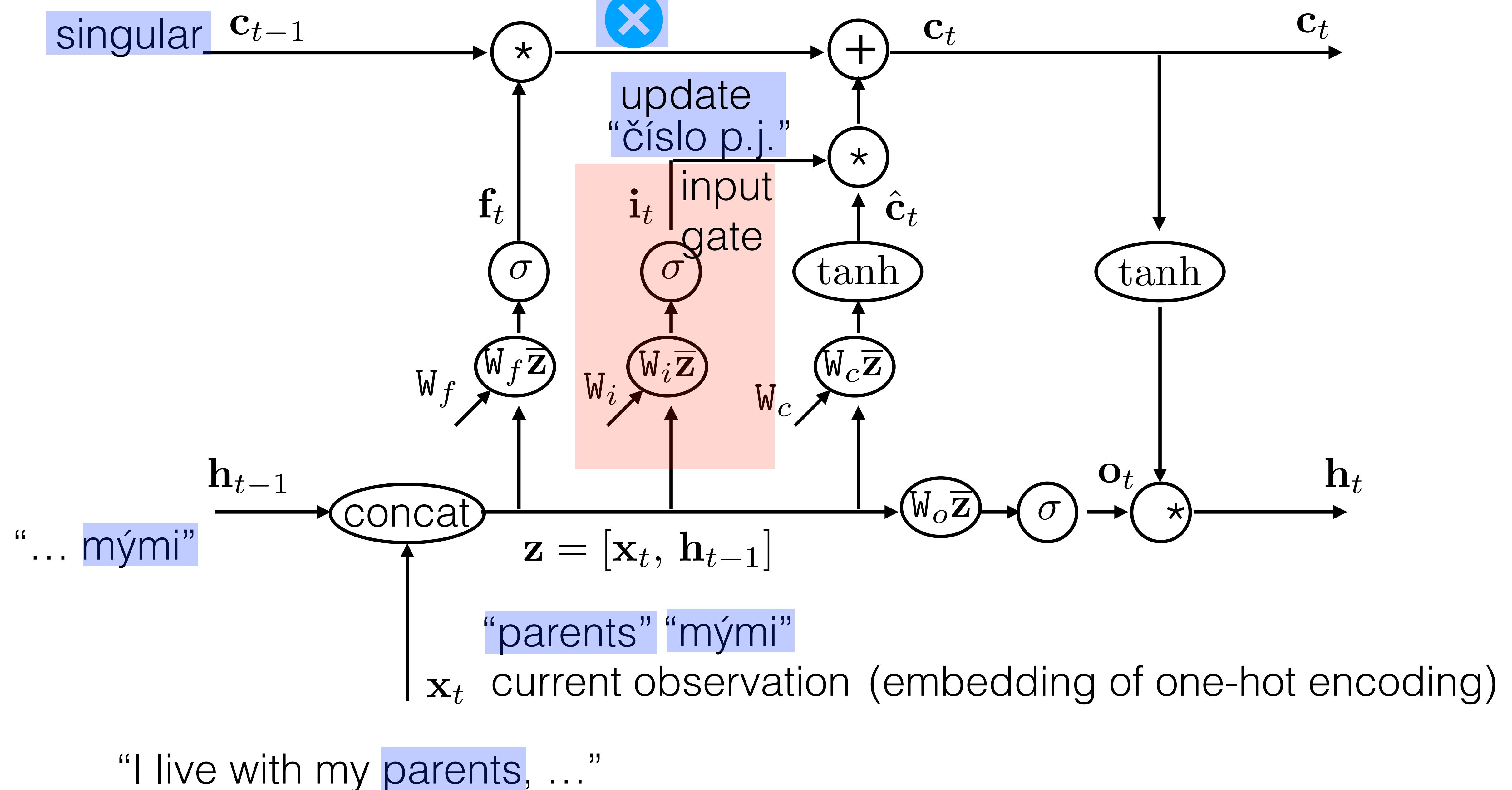
## LSTM block

cell state is long term memory= vector [singular/plural, feminine/masculine, case, ...]



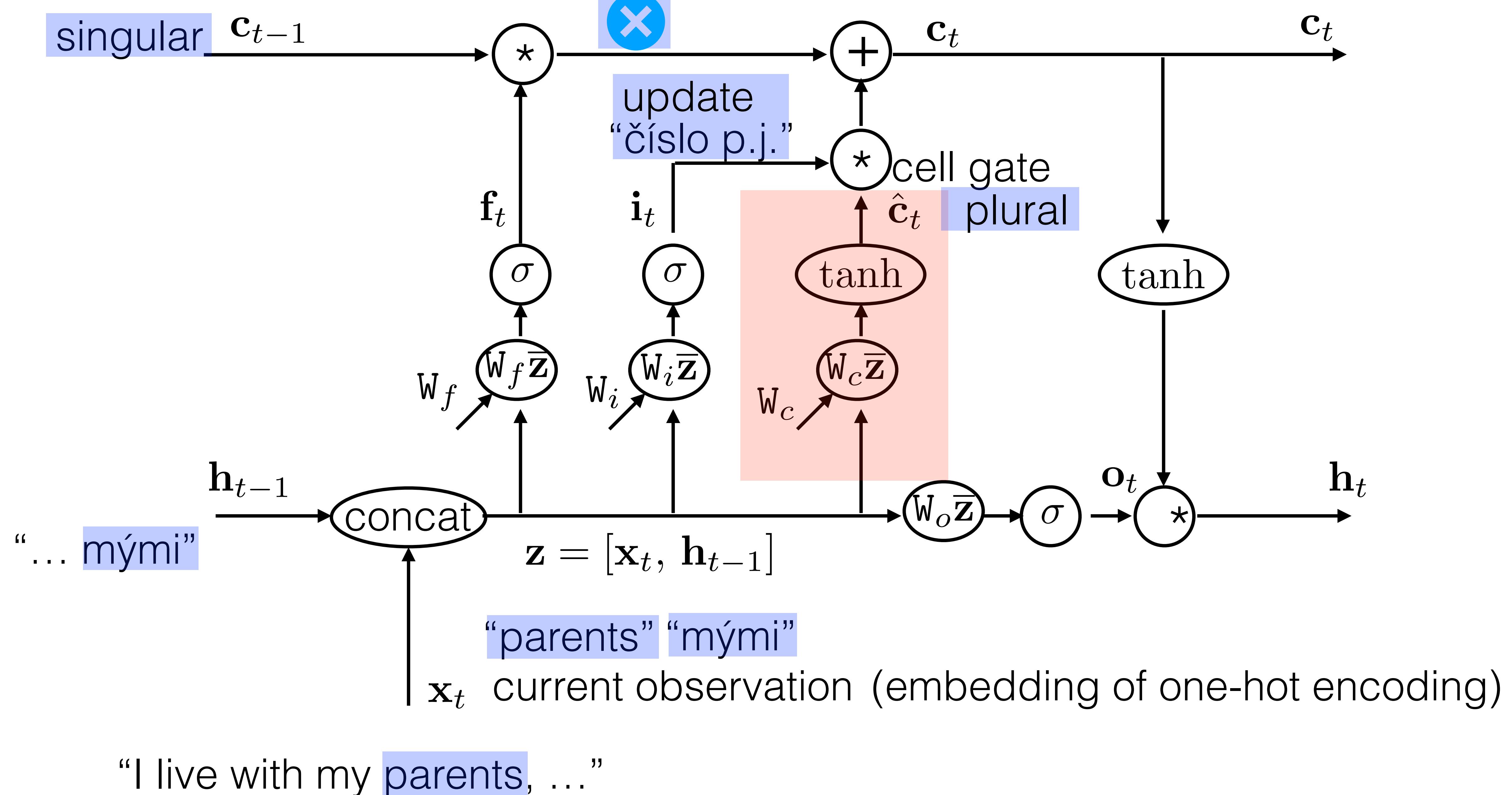
## LSTM block

cell state is long term memory= vector [singular/plural, feminine/masculine, case, ...]



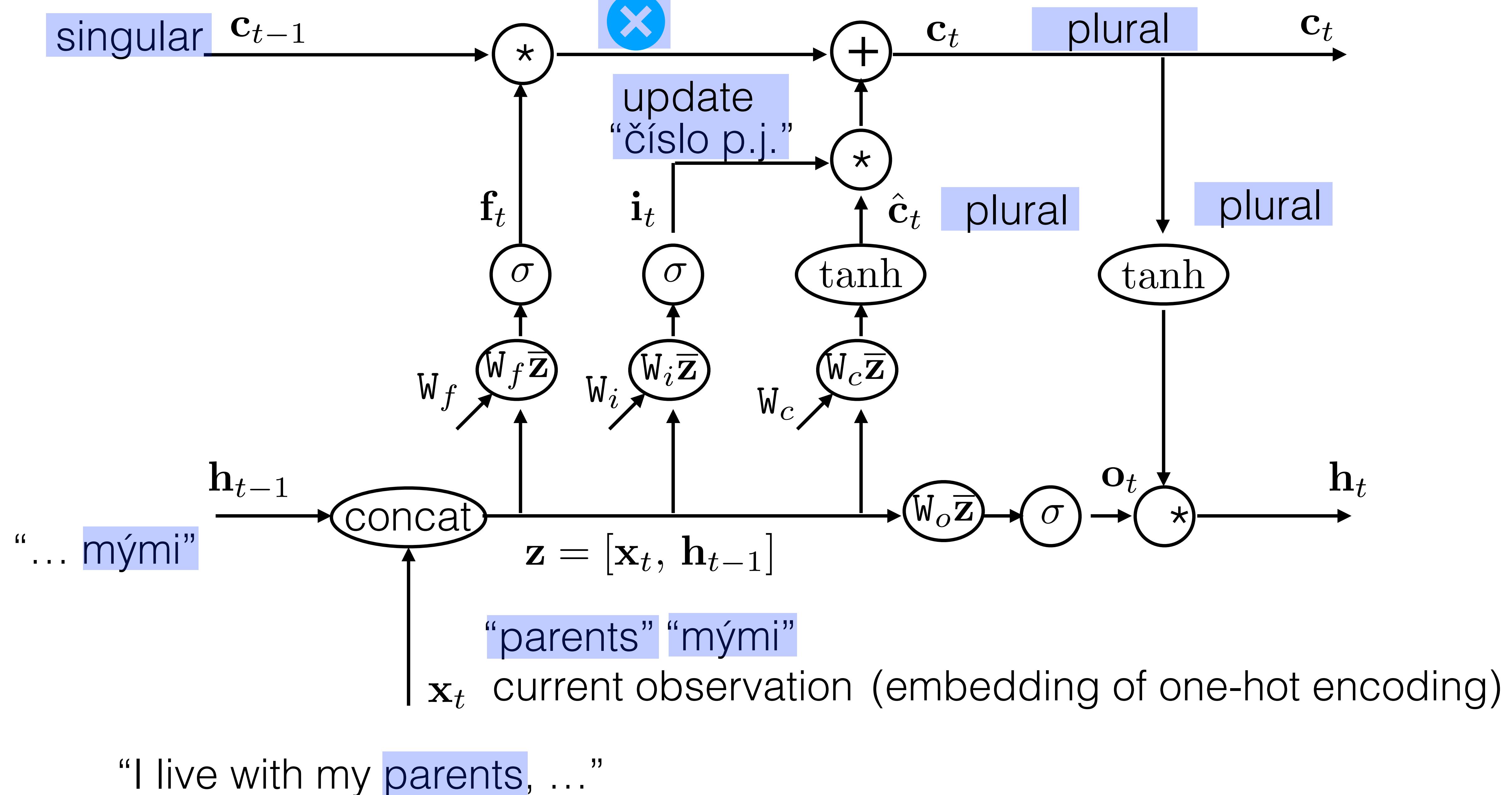
## LSTM block

cell state is long term memory= vector [singular/plural, feminine/masculine, case, ...]



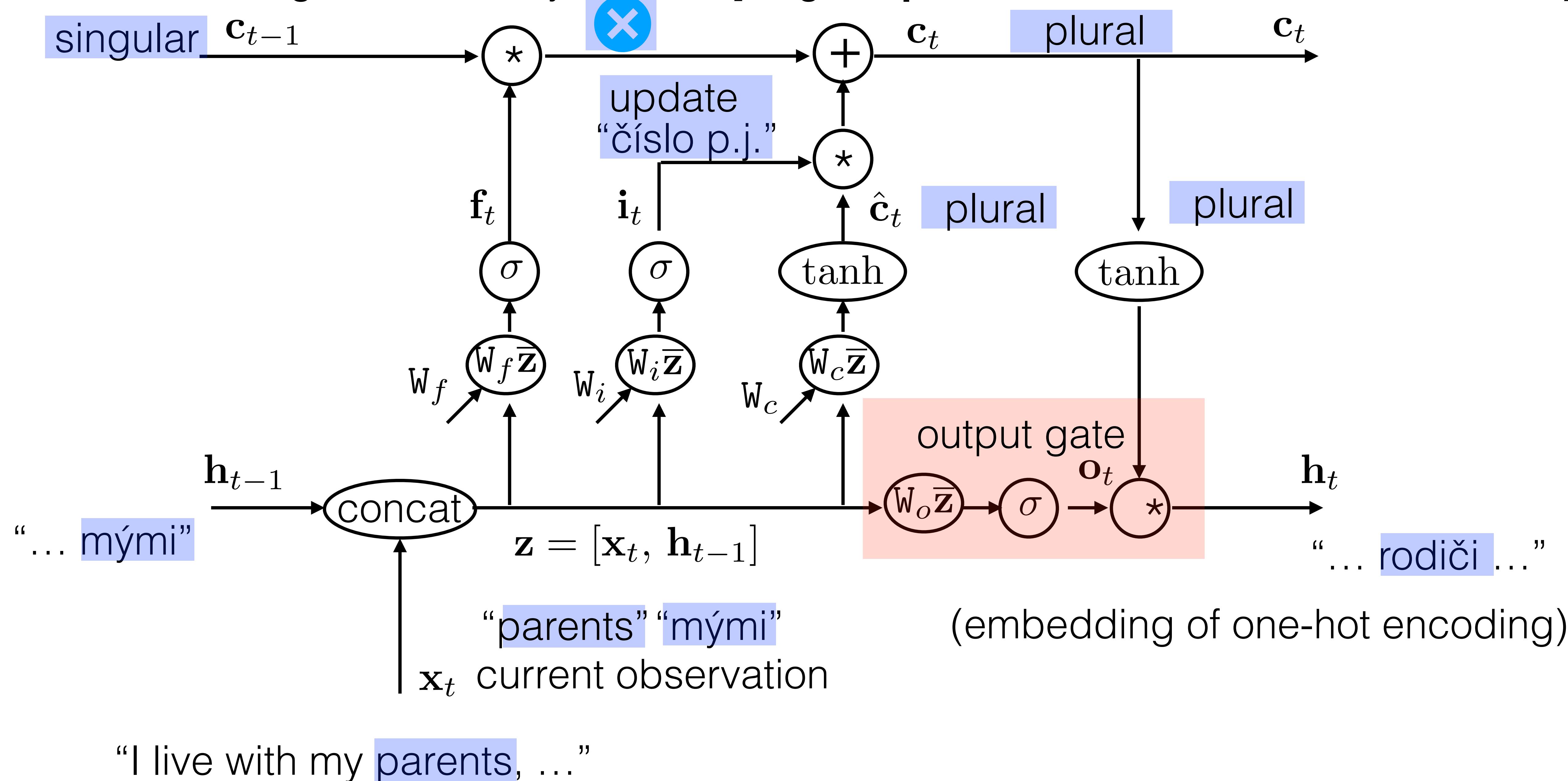
## LSTM block

cell state is long term memory= vector [singular/**plural**, feminine/masculine, case, ...]



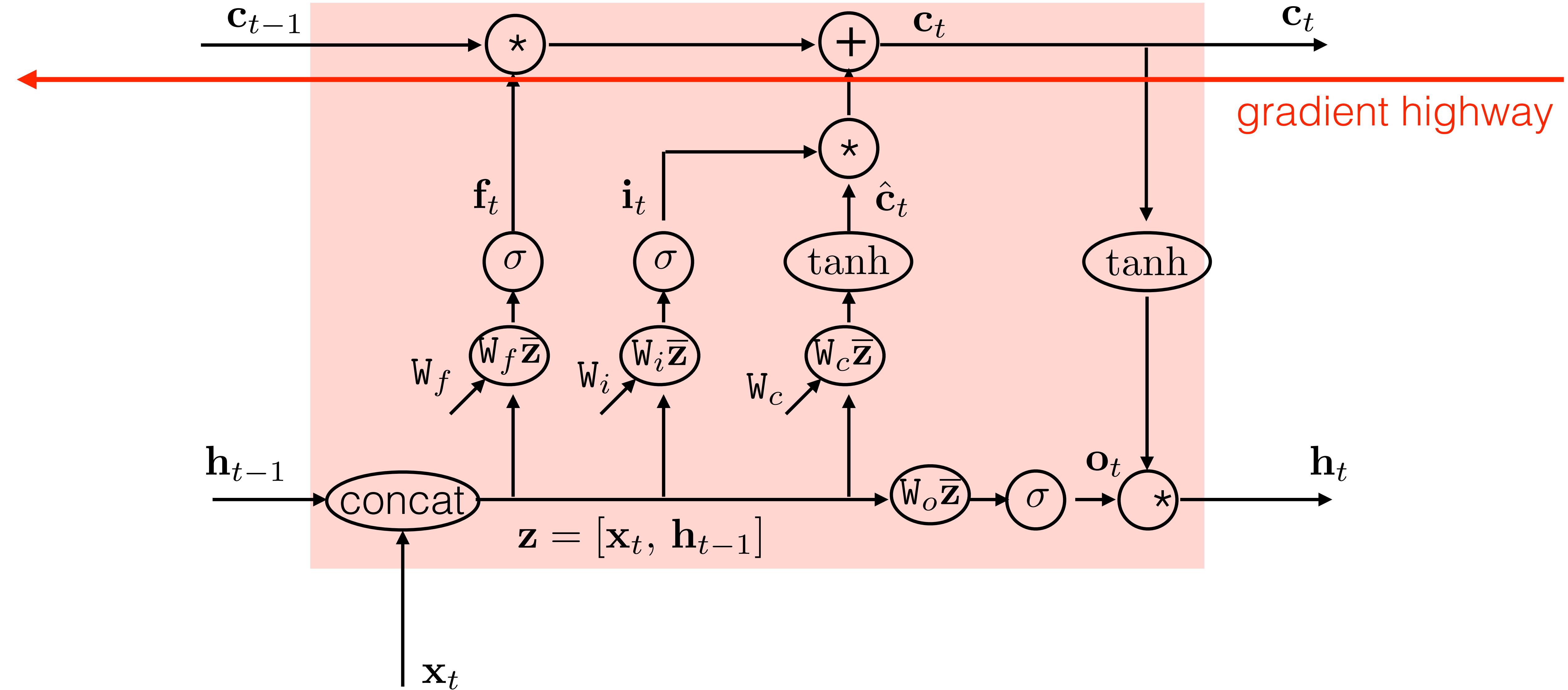
## LSTM block

cell state is long term memory= vector [singular/**plural**, feminine/masculine, case, ...]



## LSTM block

```
torch.nn.LSTM(input_size, hidden_dim, n_layers)
```



## Summary

- recurrence is **resolved by unrolling** the computational graph in time.
- memory is attention through time [Alex Graves 2020]