

# What can('t) we do we ConvNets?

**Recurrent nets, Memory and attention**

**Karel Zimmermann**

**Czech Technical University in Prague**

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



Prerequisites: derivative of compound function

Introduce notation:  $\frac{\partial f(a, b)}{\partial a} = \partial_0 f(a, b)$ ,  $\frac{\partial f(a, b)}{\partial b} = \partial_1 f(a, b)$

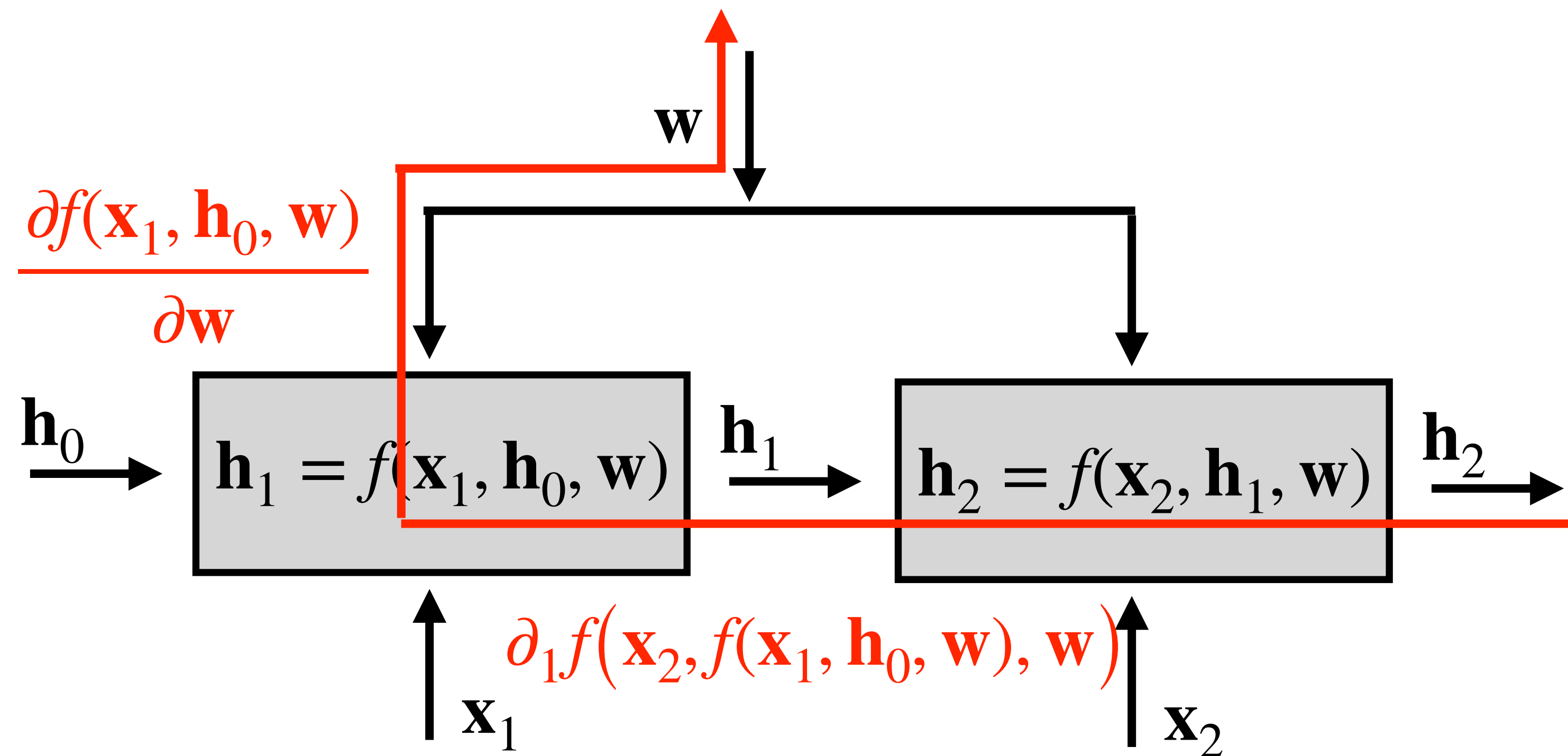
$$f(a, b) = a^2 - b^2, \quad g(x) = \sin(x) \quad \Rightarrow \quad f(x, g(x)) = x^2 - \sin^2(x)$$

$$\frac{\partial f(x, g(x))}{\partial x} = \partial_0 f(x, g(x)) + \partial_1 f(x, g(x)) \frac{\partial g(x)}{\partial x} = 2x - 2 \sin(x) \cos(x)$$

# Prerequisites: derivative of compound function

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

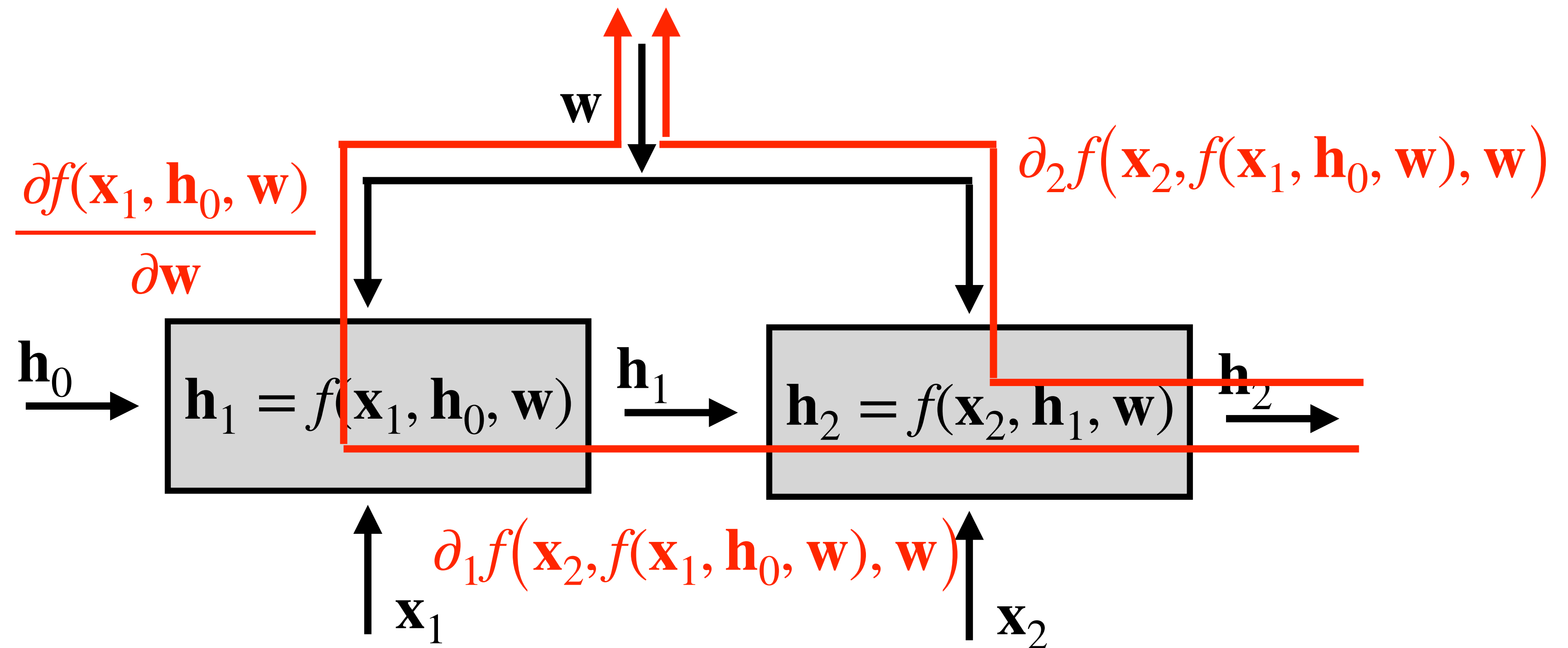
$$\frac{\partial f(\mathbf{x}_2, f(\mathbf{x}_1, \mathbf{h}_0, \mathbf{w}), \mathbf{w})}{\partial \mathbf{w}} = \partial_1 f(\mathbf{x}_2, f(\mathbf{x}_1, \mathbf{h}_0, \mathbf{w}), \mathbf{w}) \frac{\partial f(\mathbf{x}_1, \mathbf{h}_0, \mathbf{w})}{\partial \mathbf{w}}$$



# Prerequisites: derivative of compound function

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

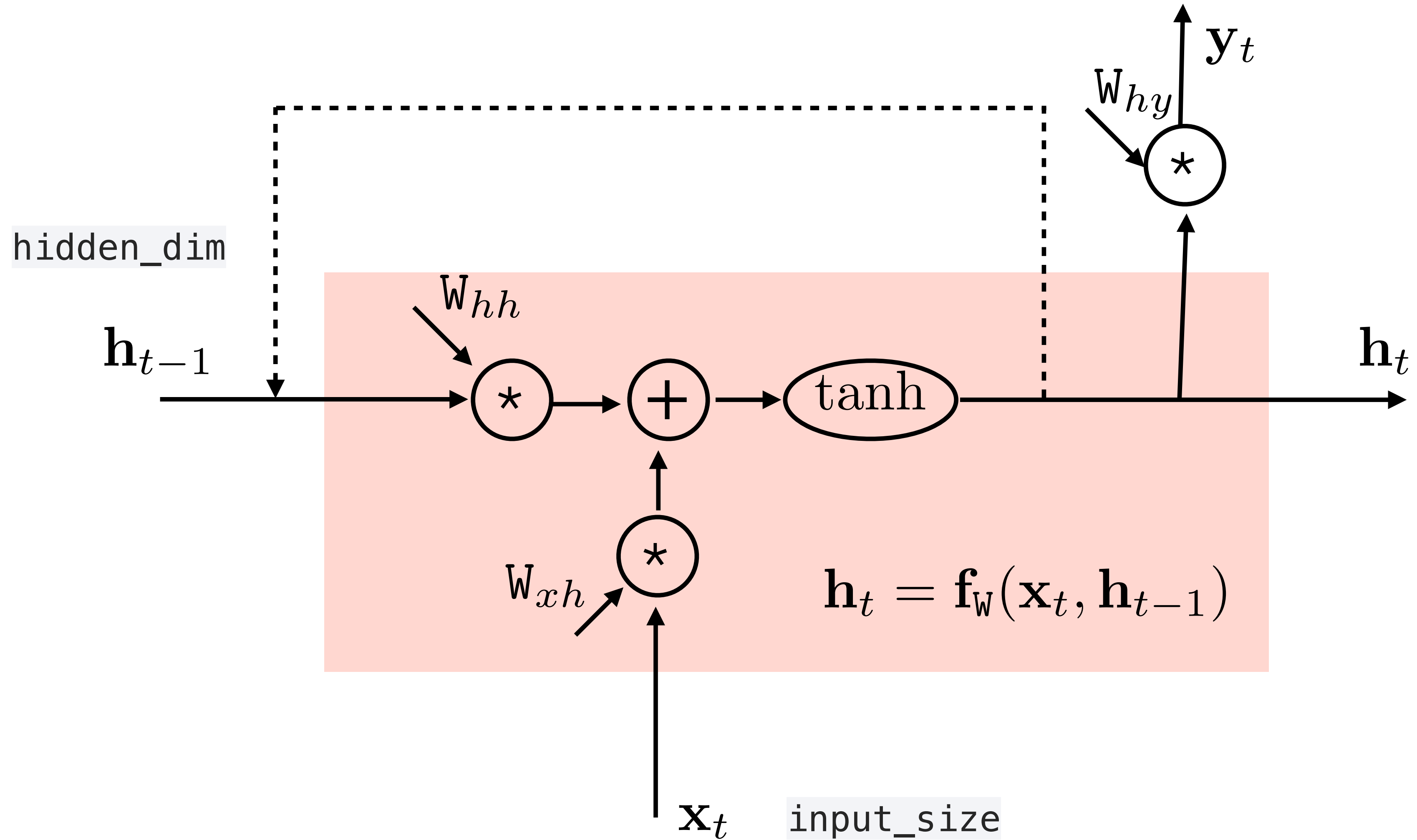
$$\frac{\partial f(\mathbf{x}_2, f(\mathbf{x}_1, \mathbf{h}_0, \mathbf{w}), \mathbf{w})}{\partial \mathbf{w}} = \partial_1 f(\mathbf{x}_2, f(\mathbf{x}_1, \mathbf{h}_0, \mathbf{w}), \mathbf{w}) \frac{\partial f(\mathbf{x}_1, \mathbf{h}_0, \mathbf{w})}{\partial \mathbf{w}} + \partial_2 f(\mathbf{x}_2, f(\mathbf{x}_1, \mathbf{h}_0, \mathbf{w}), \mathbf{w})$$





# Simple recurrent block

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



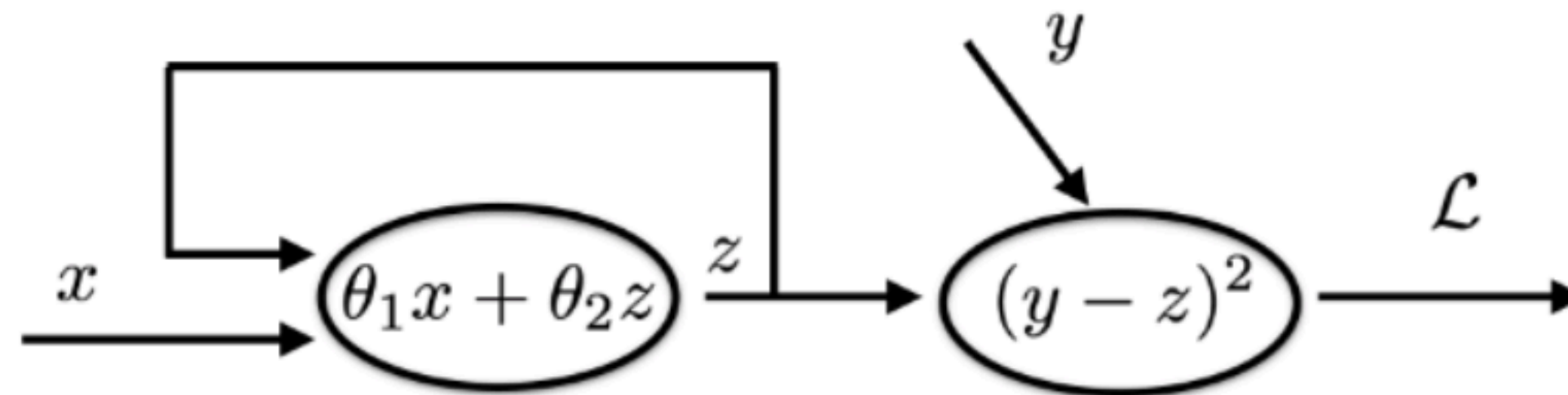
PyTorch: <https://pytorch.org/docs/stable/nn.html>

# RNN example with backprop

Consider linear recurrent neural network with L2 loss depicted on the image below. The network is initialized with parameters  $\theta_1 = 1, \theta_2 = 0, z_0 = 0$ . You are given the following training sequence:

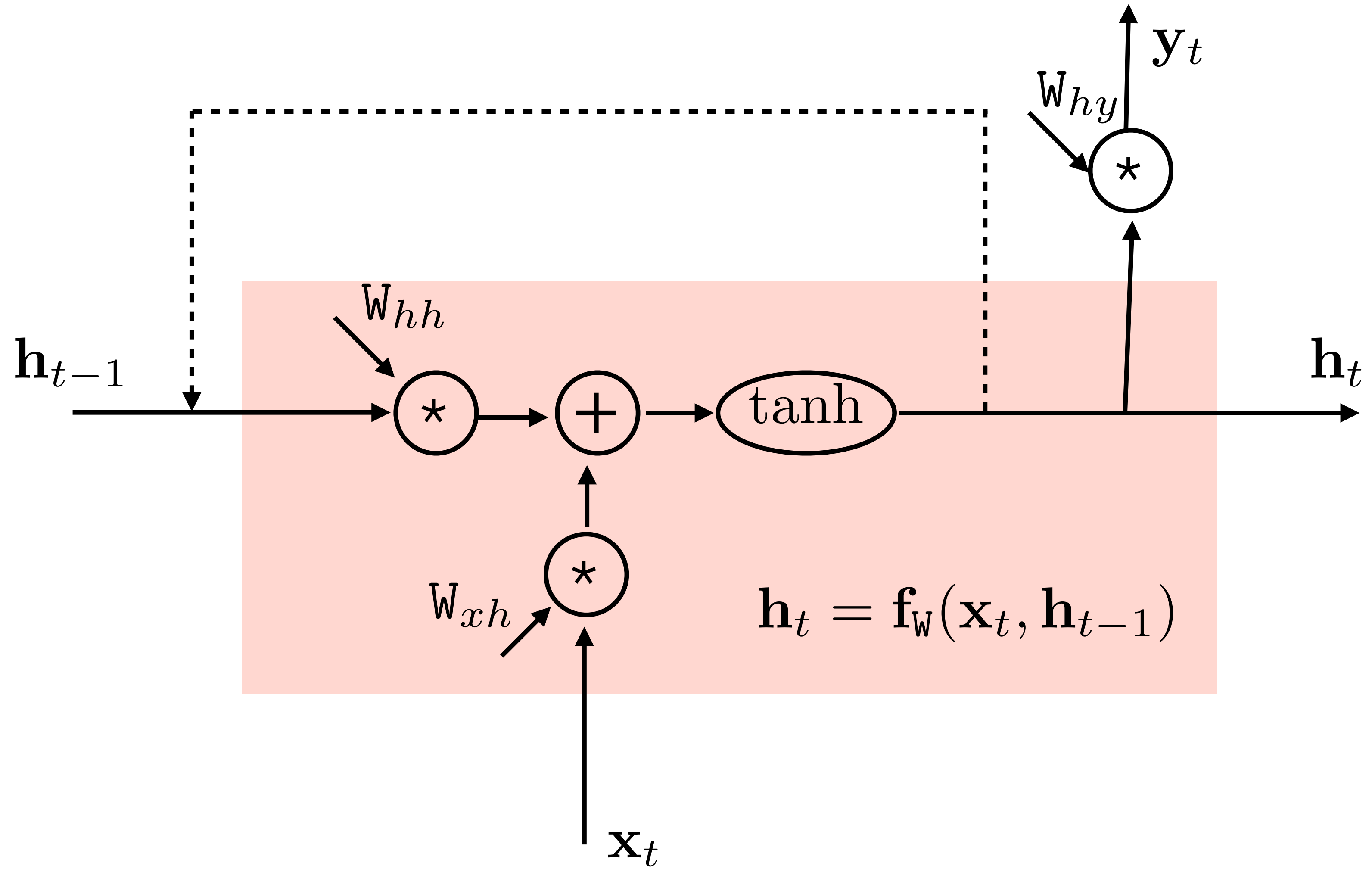
time=1	time=2
$x_1 = 0$	$x_2 = 1$
$y_1 = 1$	$y_2 = 3$

Estimate gradient of the overall loss (computed over all available outputs  $y_i$  for both available times  $i = 1, 2$ ) with respect to  $\theta_1$ .

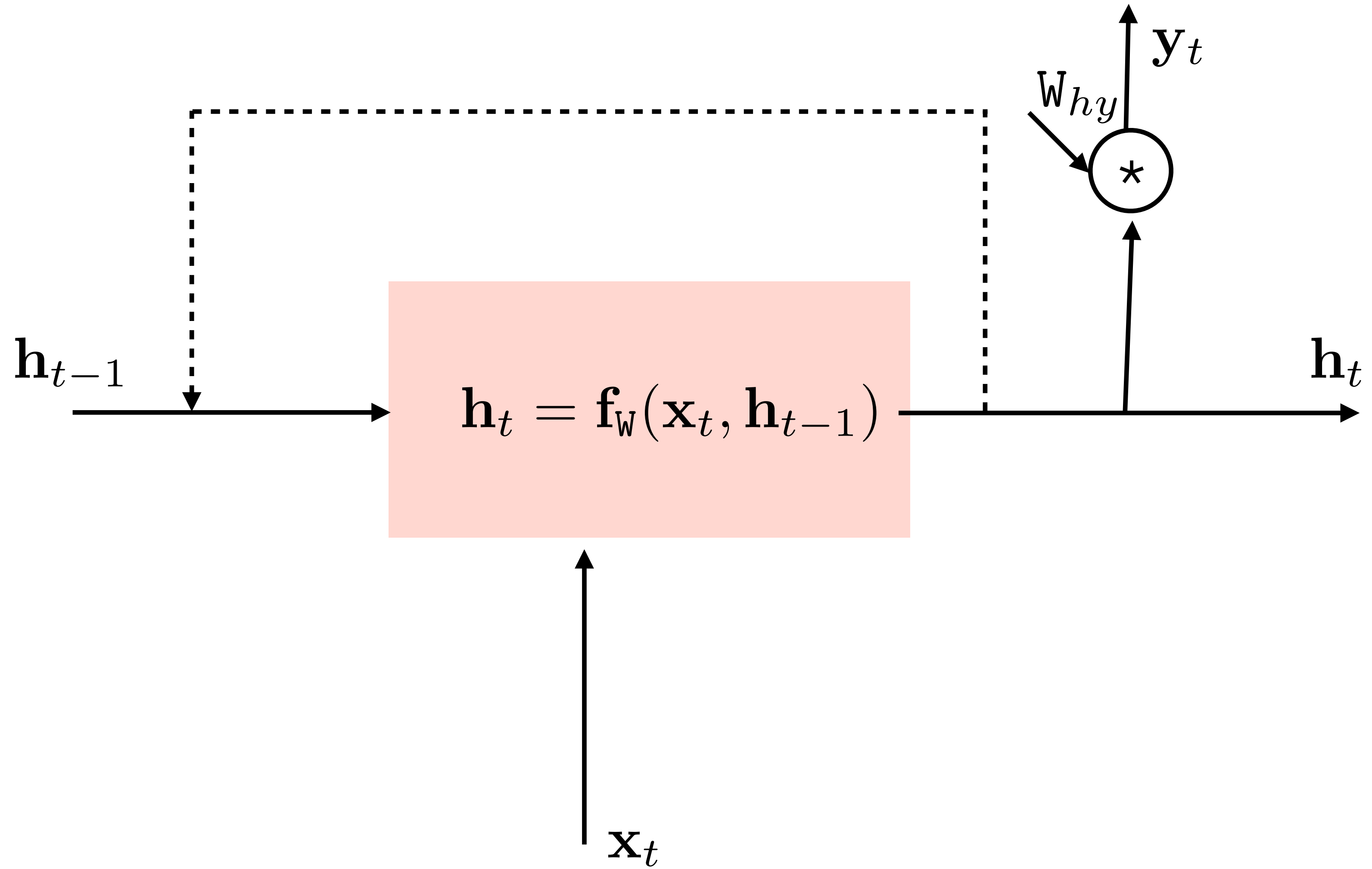


**Hint:** Unroll the network in time, to obtain a usual feedforward network with two loss nodes. Do the backpropagation as usual.

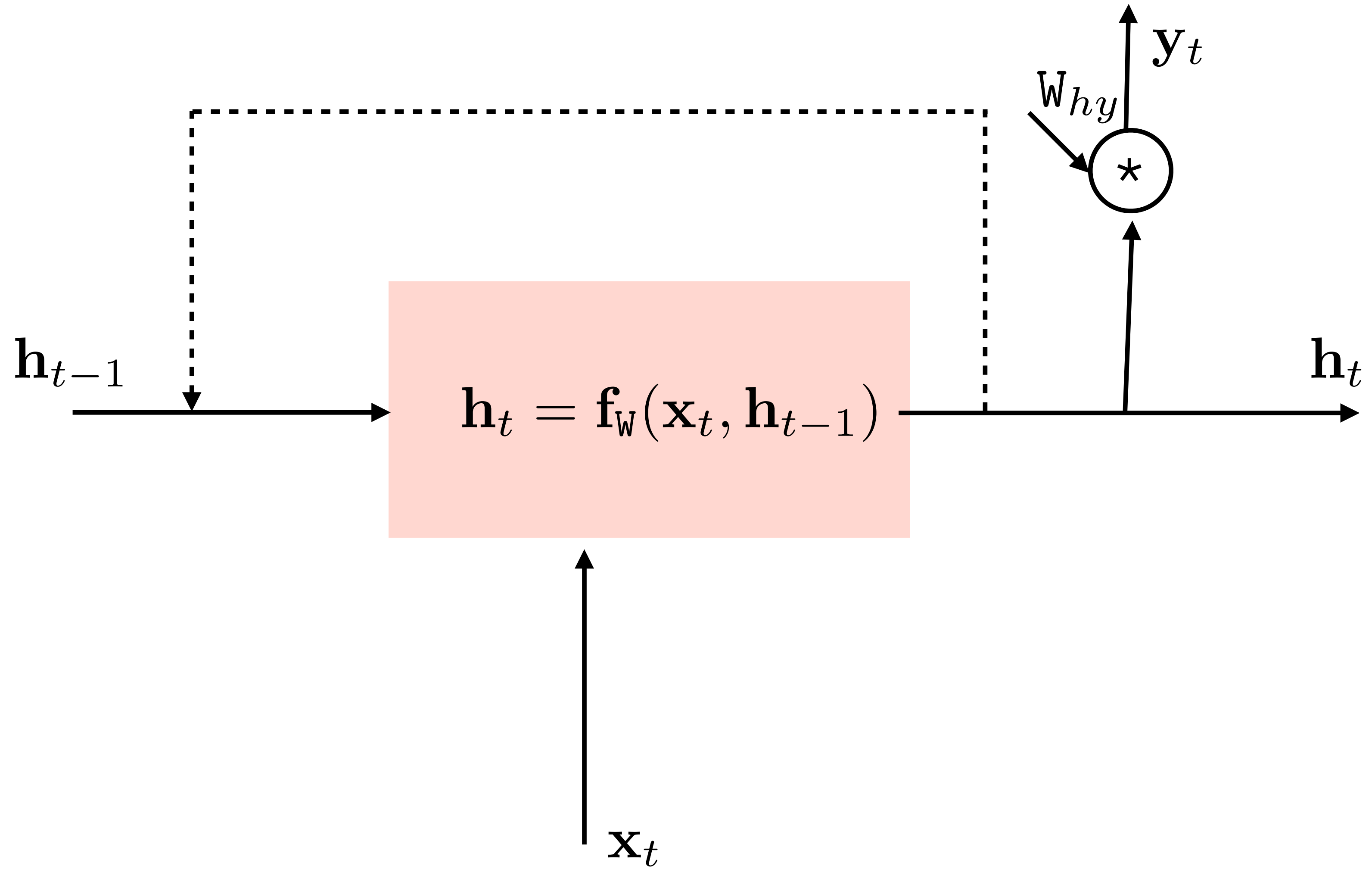
# Simple recurrent block - feed-forward pass



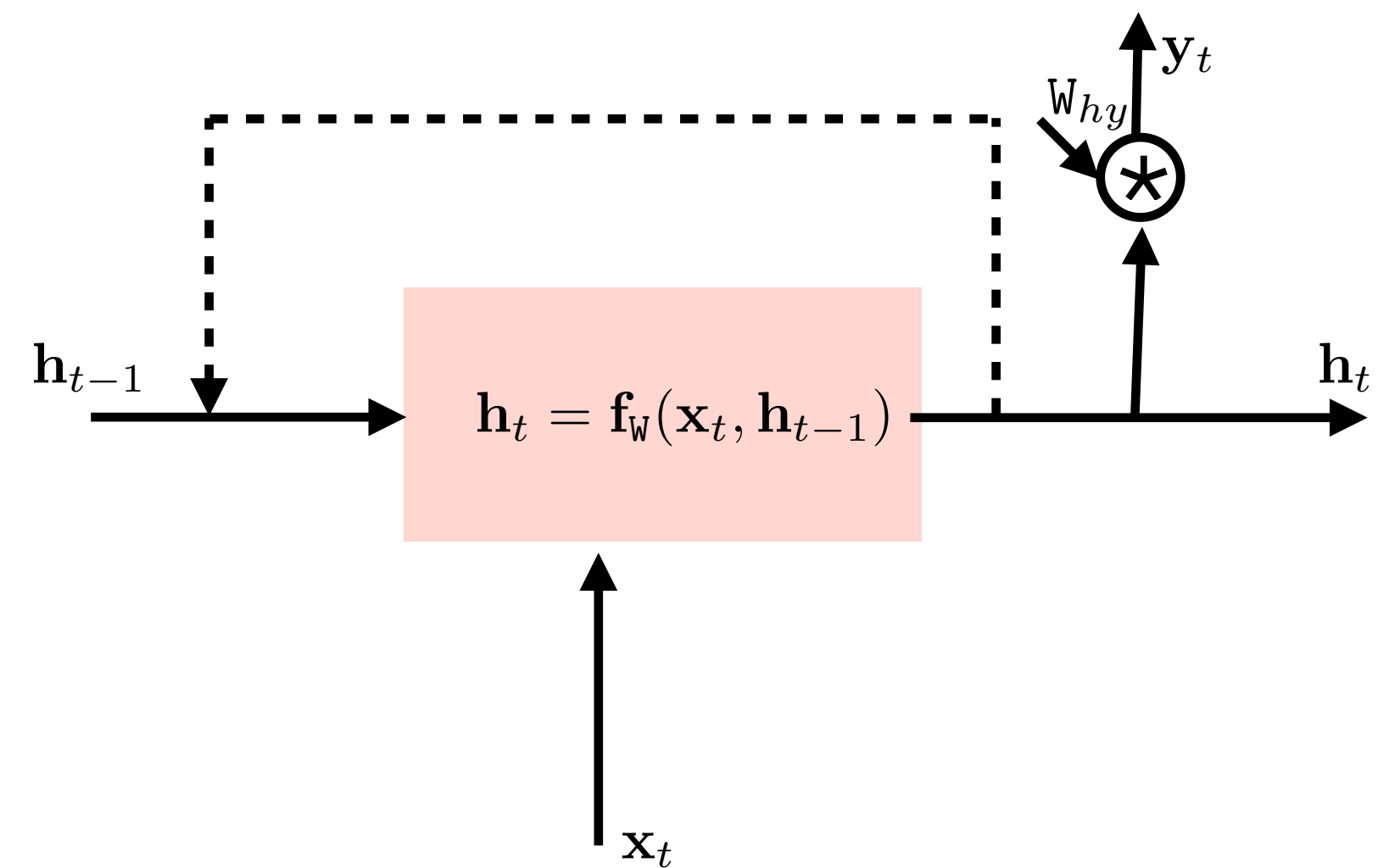
# Simple recurrent block - feed-forward pass



# Simple recurrent block - feed-forward pass

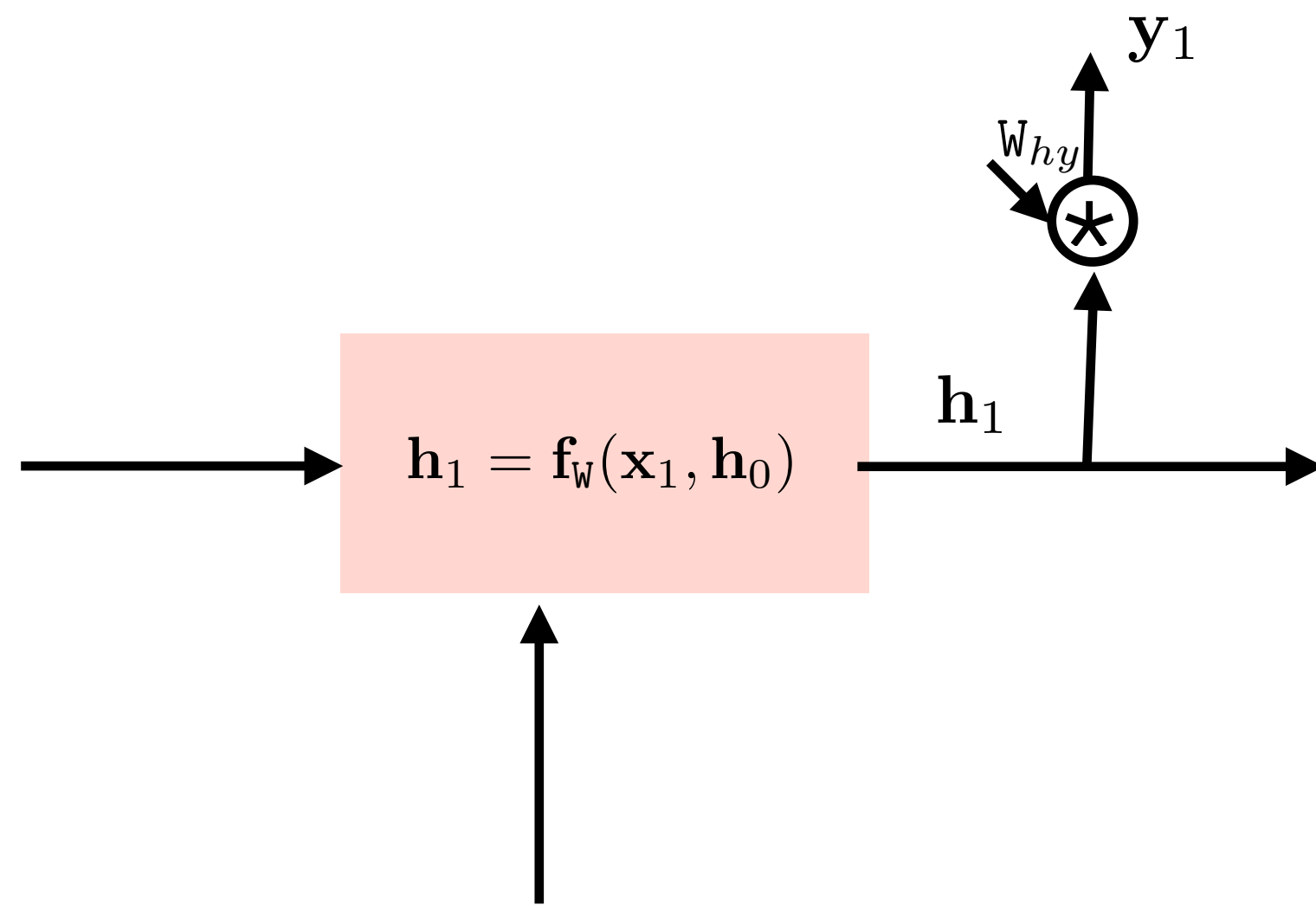


# Simple recurrent block - feed-forward pass



- 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

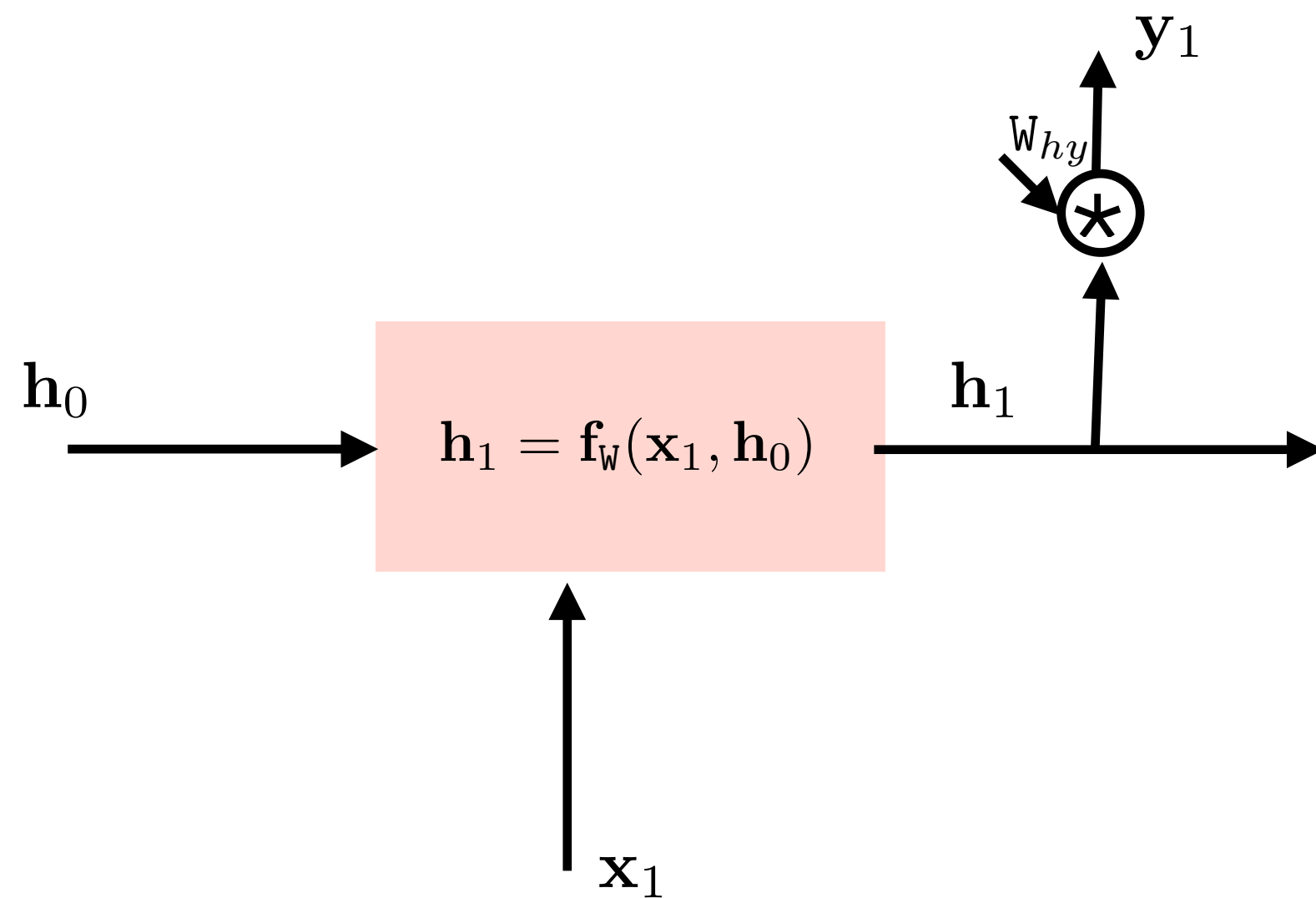
# Simple recurrent block - feed-forward pass



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

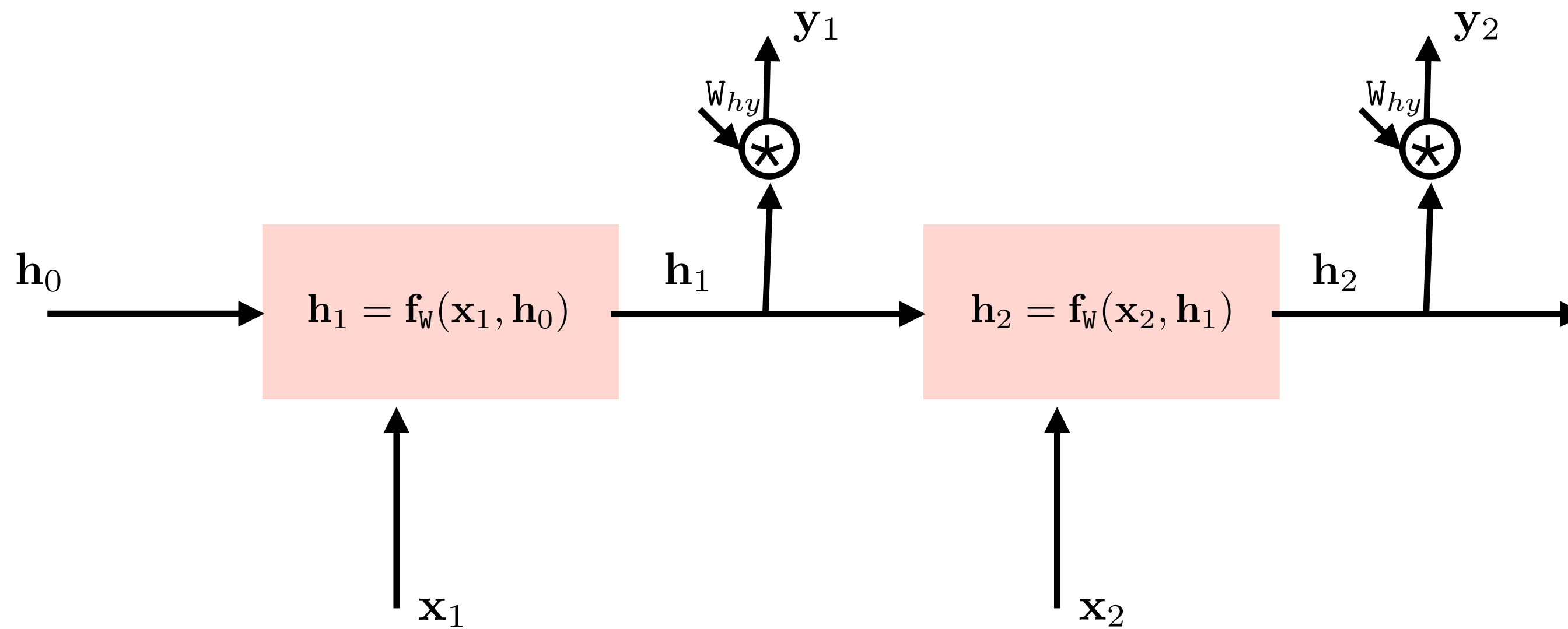


Given a finite input sequence:  $\mathbf{x}_2 \ \mathbf{x}_3$   
we remove the recurrent connection by:

- successive substitution of inputs and
- unrolling the net



# Simple recurrent block - feed-forward pass

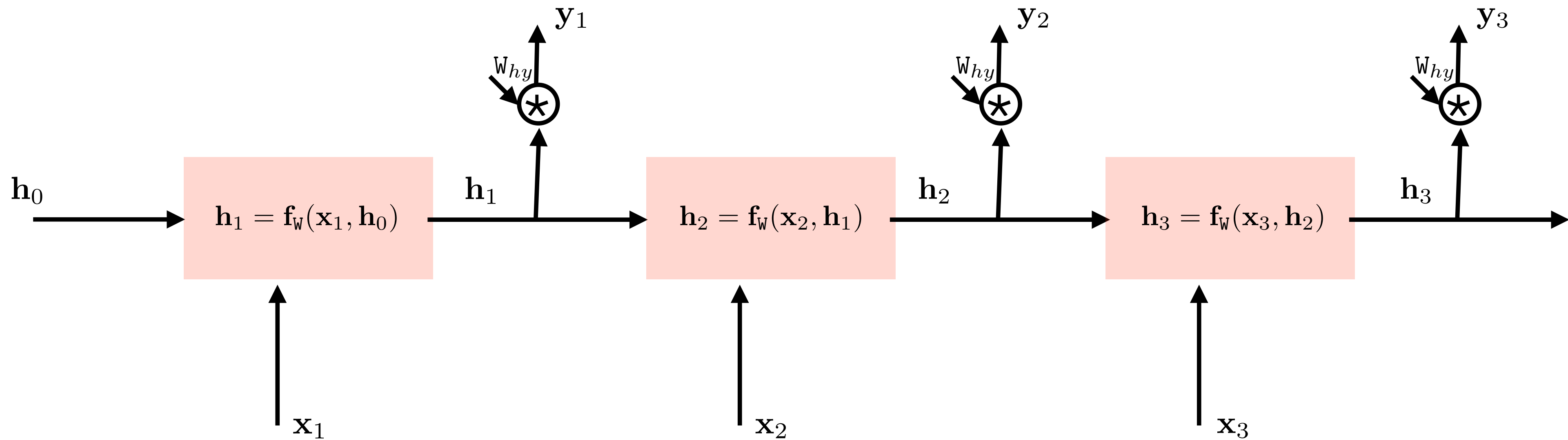


Given a finite input sequence:  $x_3$

we remove the recurrent connection by:

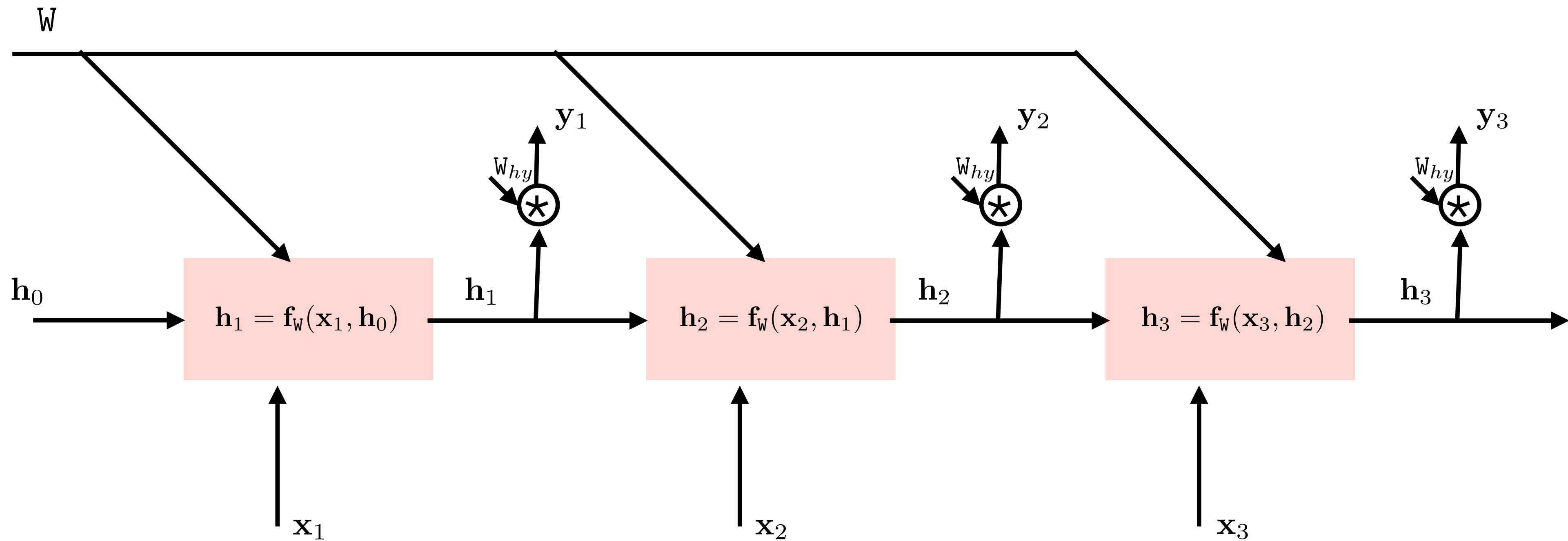
- successive substitution of inputs and
- unrolling the net

# Simple recurrent block - feed-forward pass



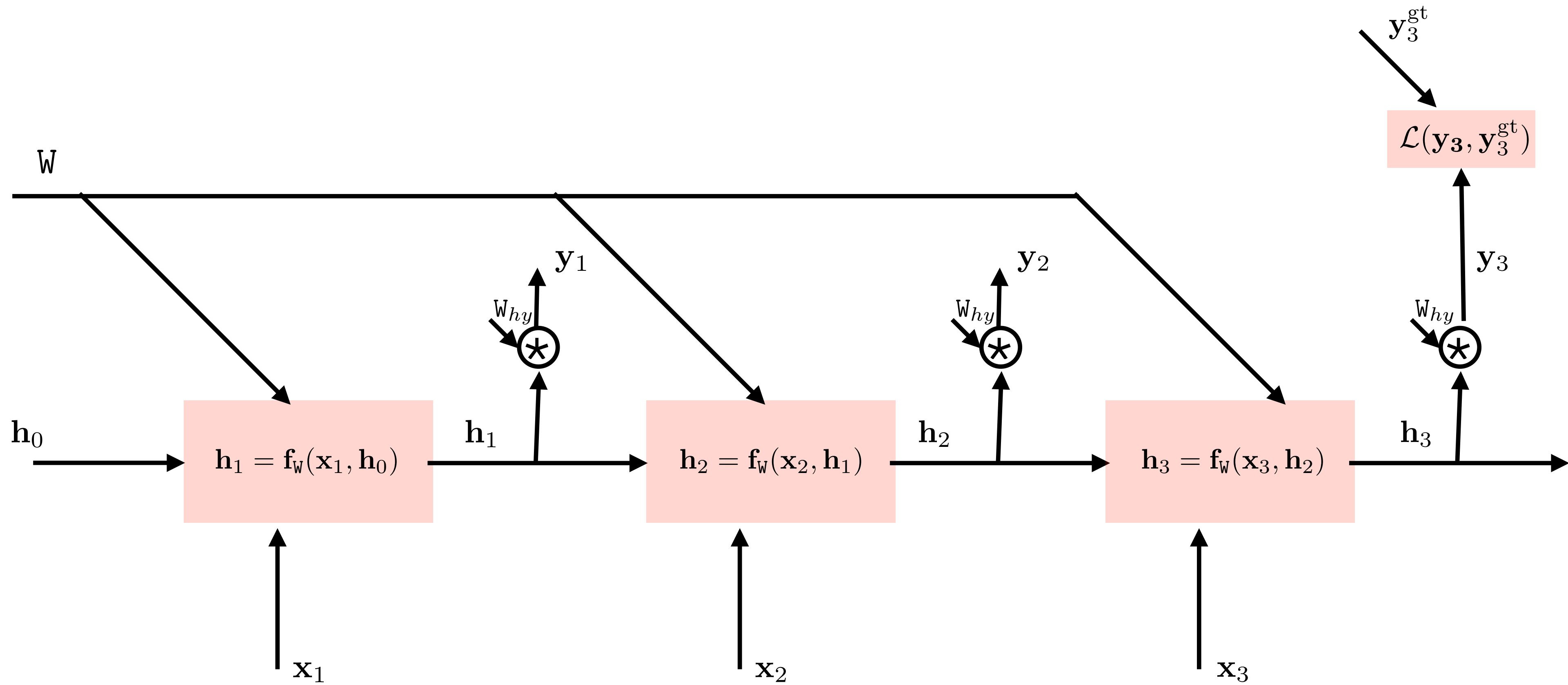
- Unrolled computational graph:
  - it is normal feedforward network

# Simple recurrent block - feed-forward pass



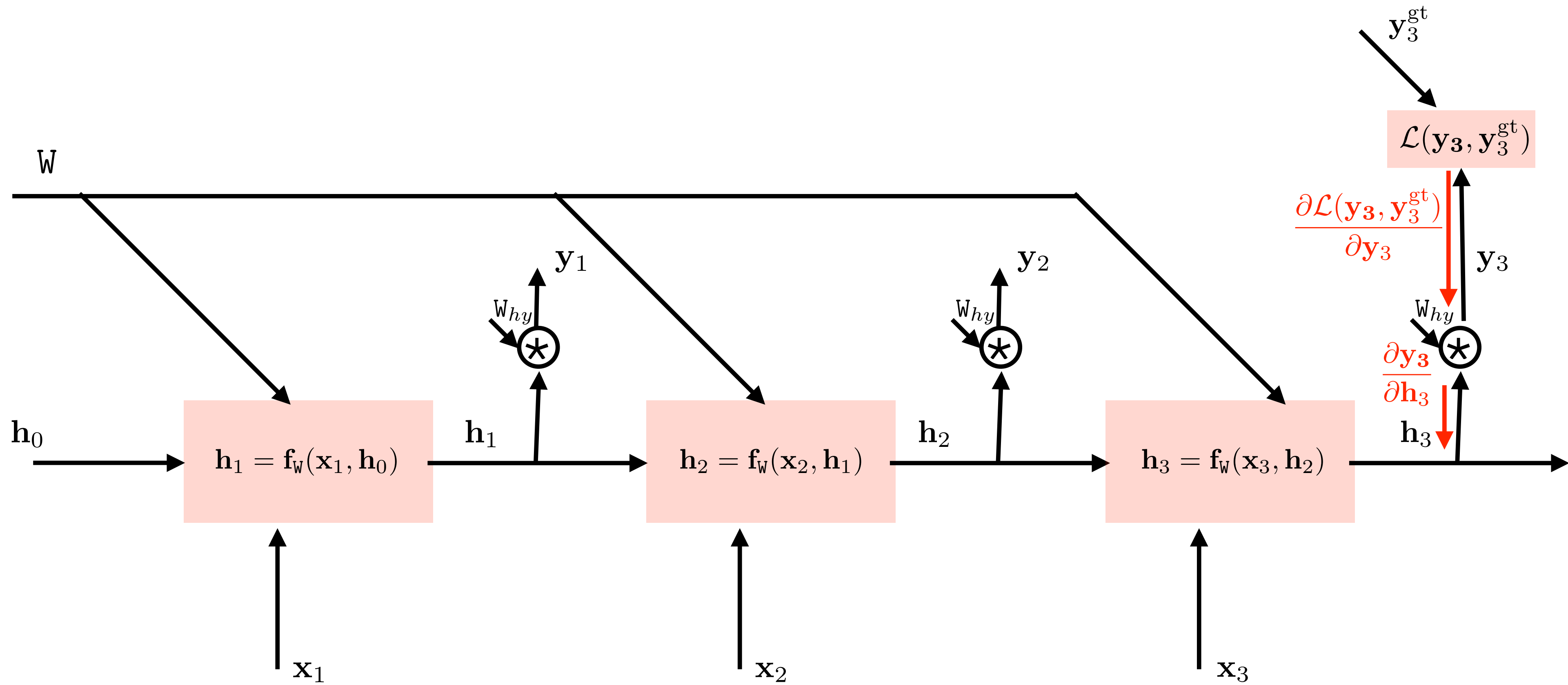
- Unrolled computational graph:
  - it is normal feedforward network
  - it consists of several same blocks with the same weights!

# Simple recurrent block - backward pass



- Loss function:
  - cross-entropy loss on the last output only (for simplicity)

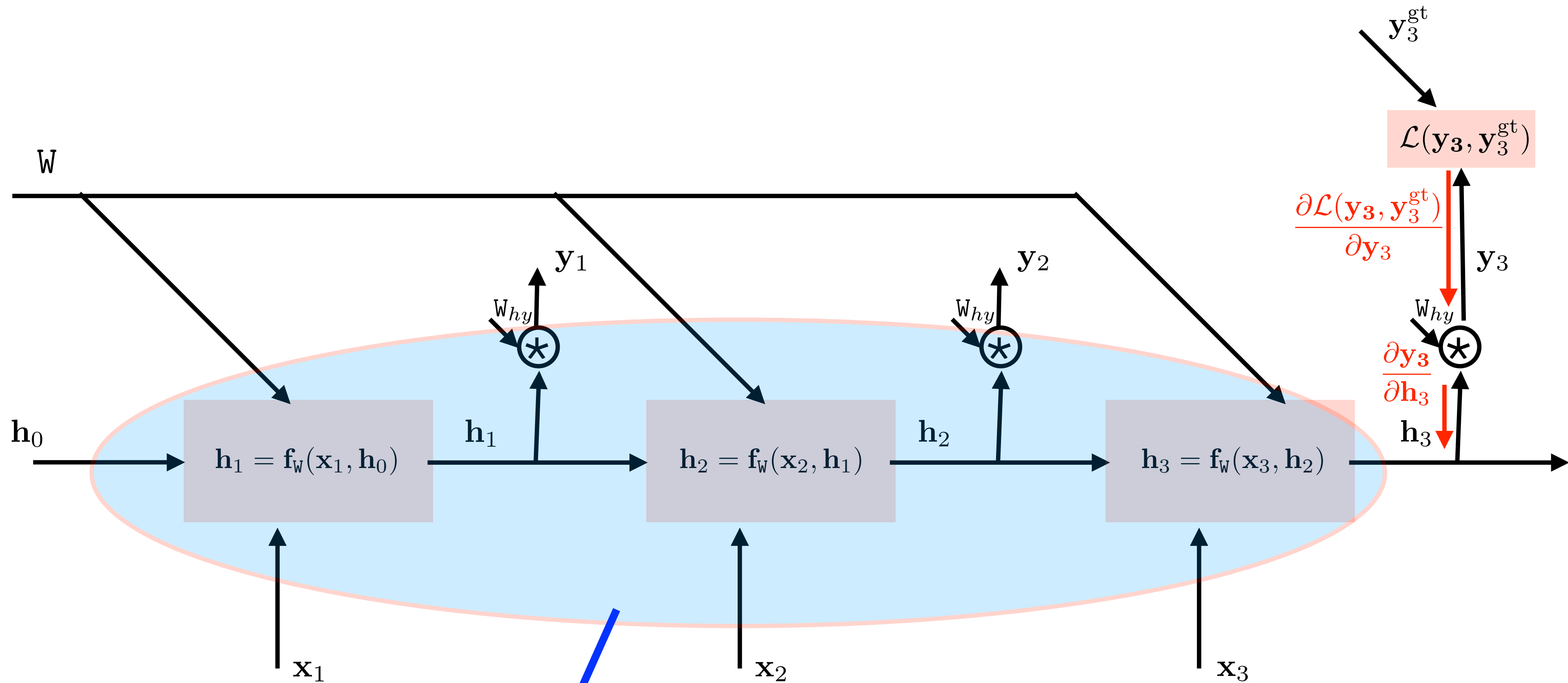
# Simple recurrent block - backward pass



- Backprop:

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

# Simple recurrent block - backward pass

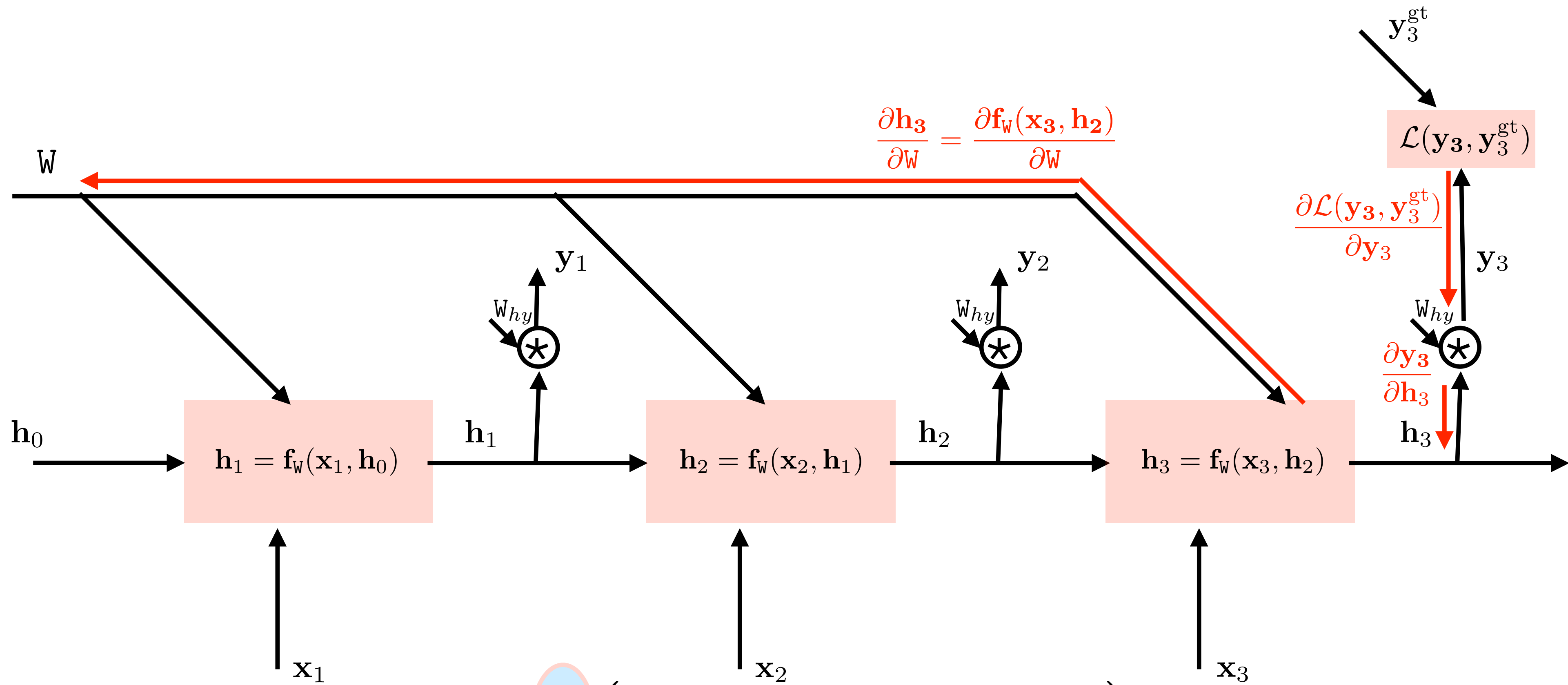


- Backprop:
  - differentiation of multi-dimensional composite function:

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

$$\mathbf{h}_3 = \mathbf{f}_W \left( \mathbf{x}_3, \mathbf{f}_W \left( \mathbf{x}_2, \mathbf{f}_W \left( \mathbf{x}_1, \mathbf{h}_0 \right) \right) \right)$$

# Simple recurrent block - backward pass



$$\frac{\partial h_3}{\partial W} = \frac{\partial f_W(x_3, h_2)}{\partial W}$$

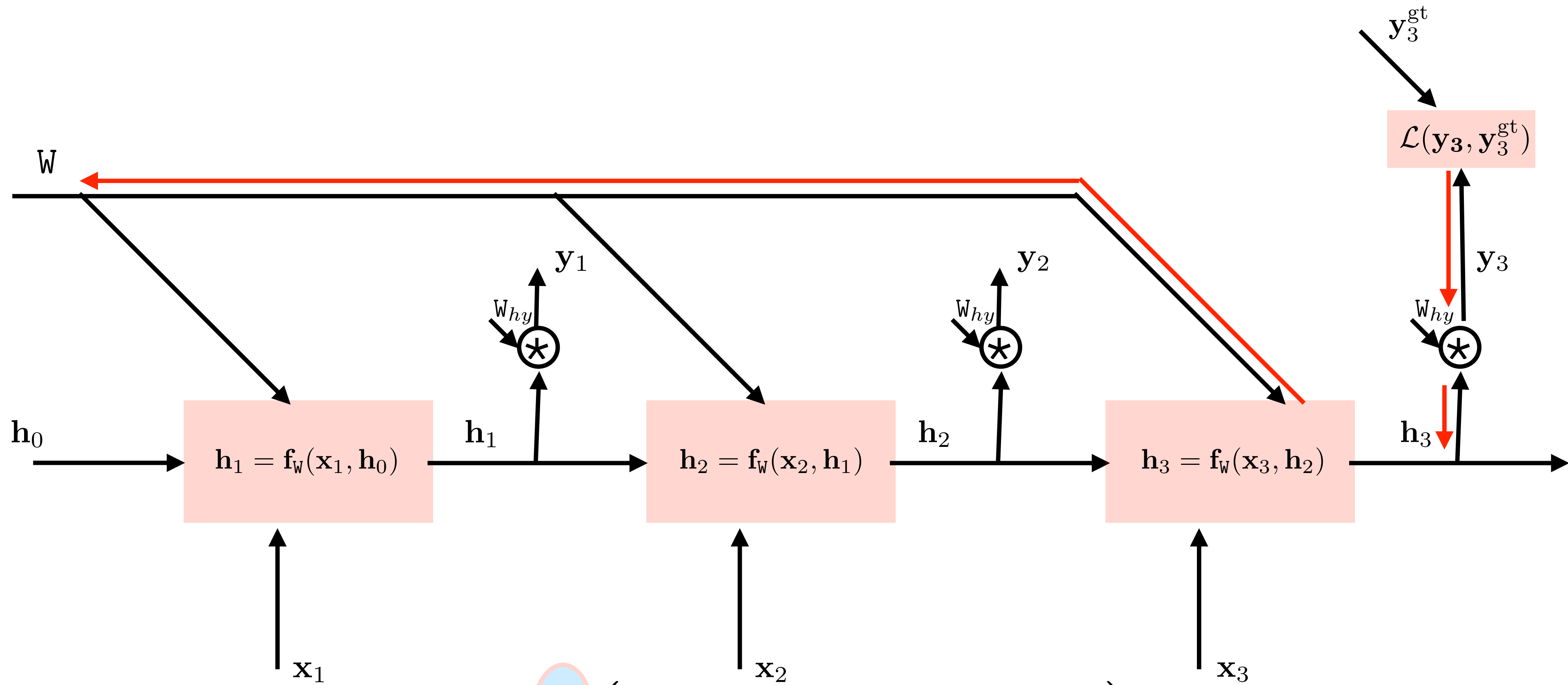
$$\frac{\partial \mathcal{L}(y_3, y_3^{gt})}{\partial y_3}$$

$$\frac{\partial y_3}{\partial h_3}$$

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

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

# Simple recurrent block - backward pass

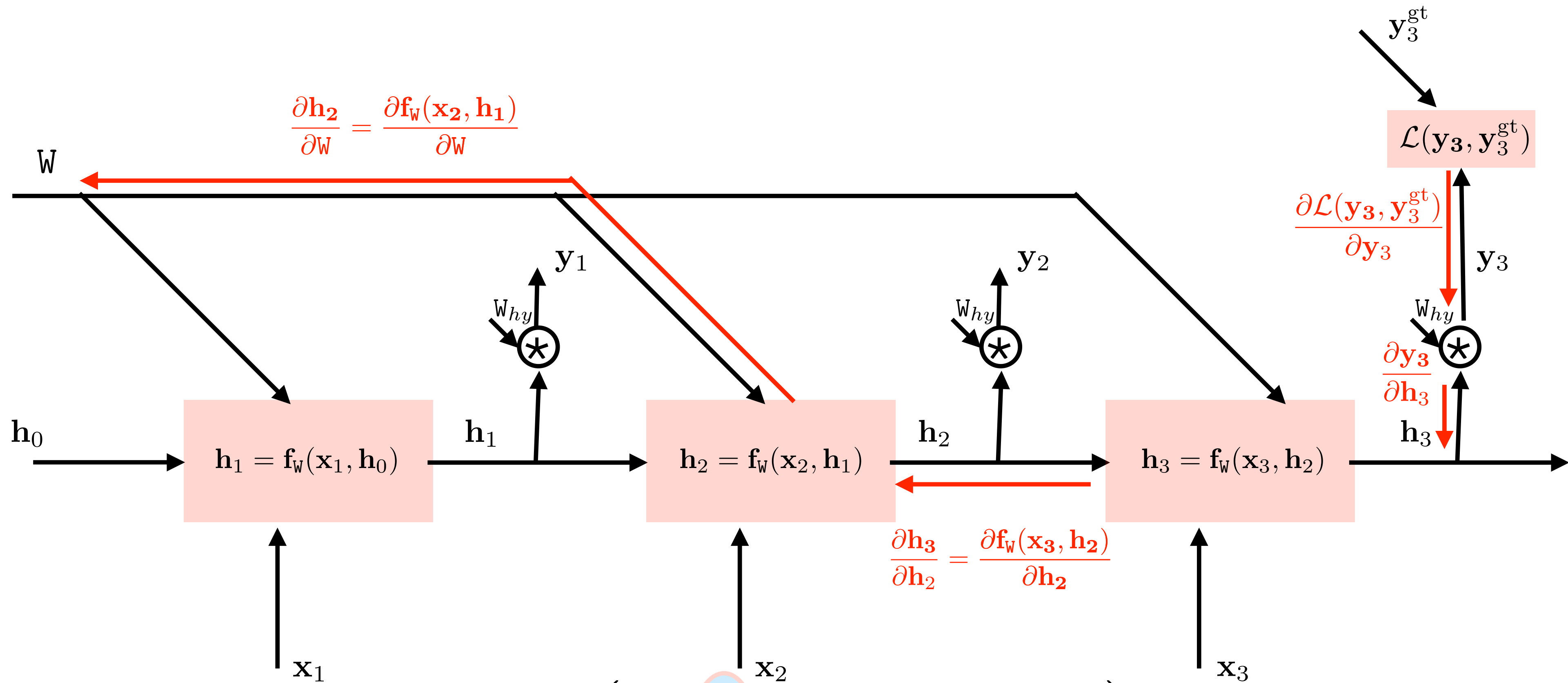


$$\mathbf{h}_3 = \mathbf{f}_W(\mathbf{x}_3, \mathbf{f}_W(\mathbf{x}_2, \mathbf{f}_W(\mathbf{x}_1, \mathbf{h}_0)))$$

$$\frac{\partial \mathcal{L}(y_3, y_3^{gt})}{\partial W} = \frac{\partial \mathcal{L}(y_3, y_3^{gt})}{\partial y_3} \frac{\partial y_3}{\partial h_3} \frac{\partial f_W(\mathbf{x}_3, \mathbf{h}_2)}{\partial W} +$$



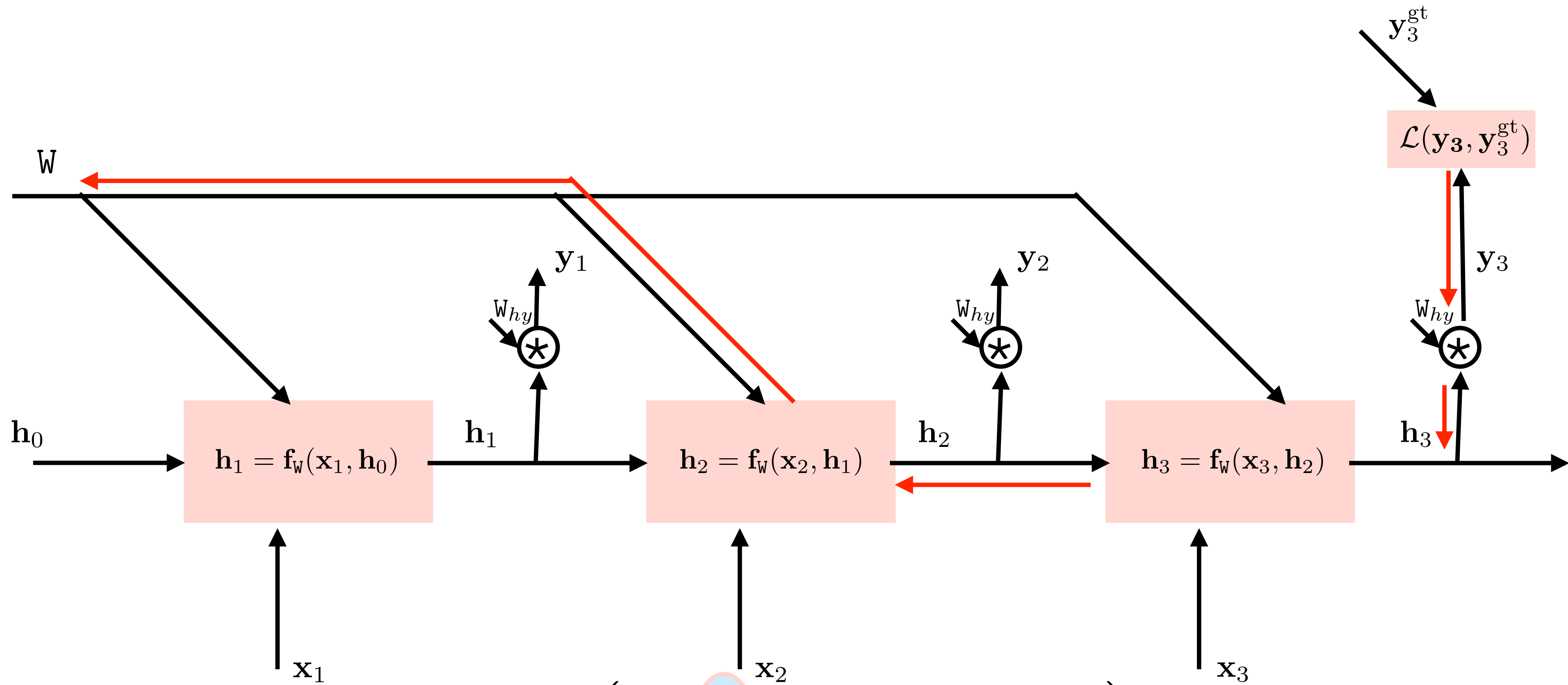
# Simple recurrent block - backward pass



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

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

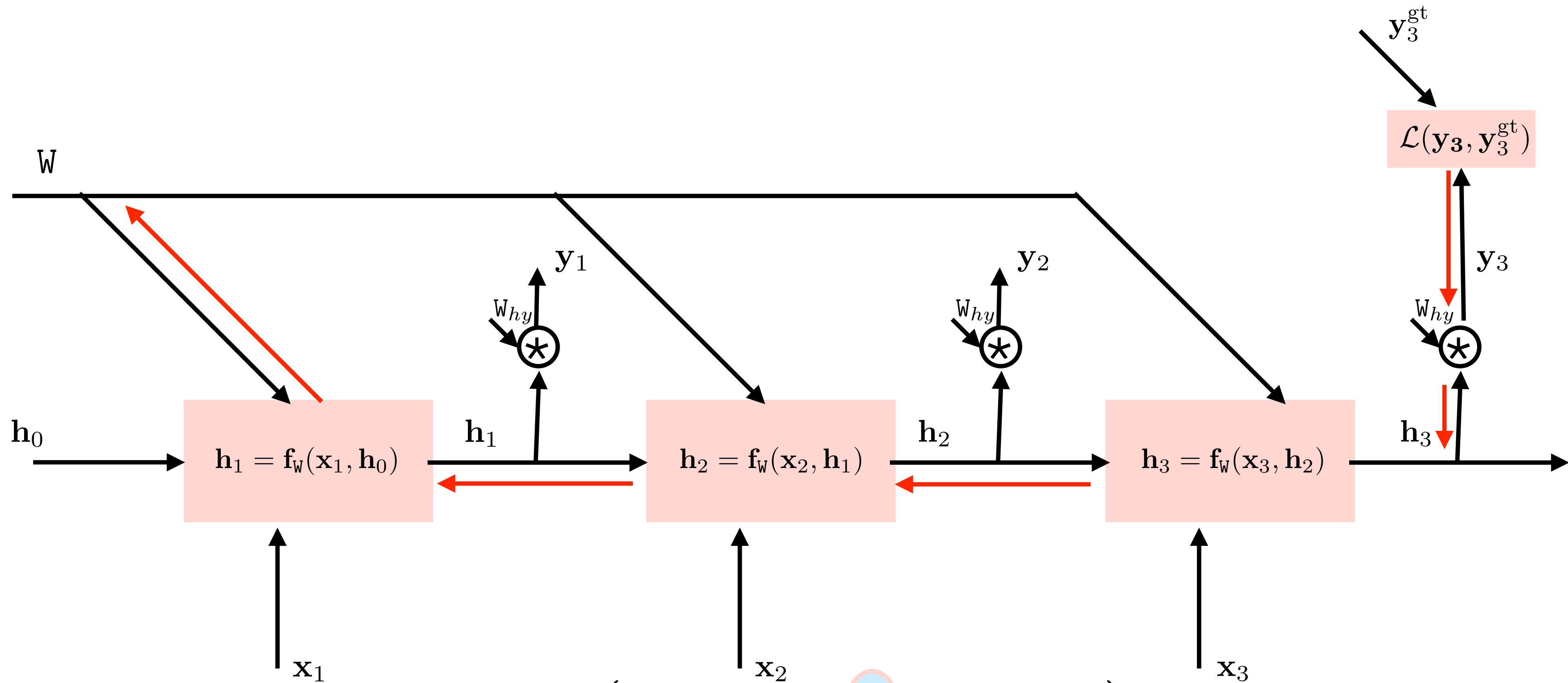
# Simple recurrent block - backward pass



$$\mathbf{h}_3 = \mathbf{f}_W \left( \mathbf{x}_3, \mathbf{f}_W \left( \mathbf{x}_2, \mathbf{f}_W \left( \mathbf{x}_1, \mathbf{h}_0 \right) \right) \right)$$

$$\frac{\partial \mathcal{L}(y_3, y_3^{gt})}{\partial W} = \frac{\partial \mathcal{L}(y_3, y_3^{gt})}{\partial y_3} \frac{\partial y_3}{\partial h_3} \frac{\partial \mathbf{f}_W(\mathbf{x}_3, \mathbf{h}_2)}{\partial W} + \frac{\partial \mathcal{L}(y_3, y_3^{gt})}{\partial y_3} \frac{\partial y_3}{\partial h_3} \frac{\partial \mathbf{f}_W(\mathbf{x}_3, \mathbf{h}_2)}{\partial h_2} \frac{\partial \mathbf{f}_W(\mathbf{x}_2, \mathbf{h}_1)}{\partial W} +$$

# Simple recurrent block - backward pass



$$\mathbf{h}_3 = \mathbf{f}_W \left( \mathbf{x}_3, \mathbf{f}_W \left( \mathbf{x}_2, \mathbf{f}_W \left( \mathbf{x}_1, \mathbf{h}_0 \right) \right) \right)$$

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

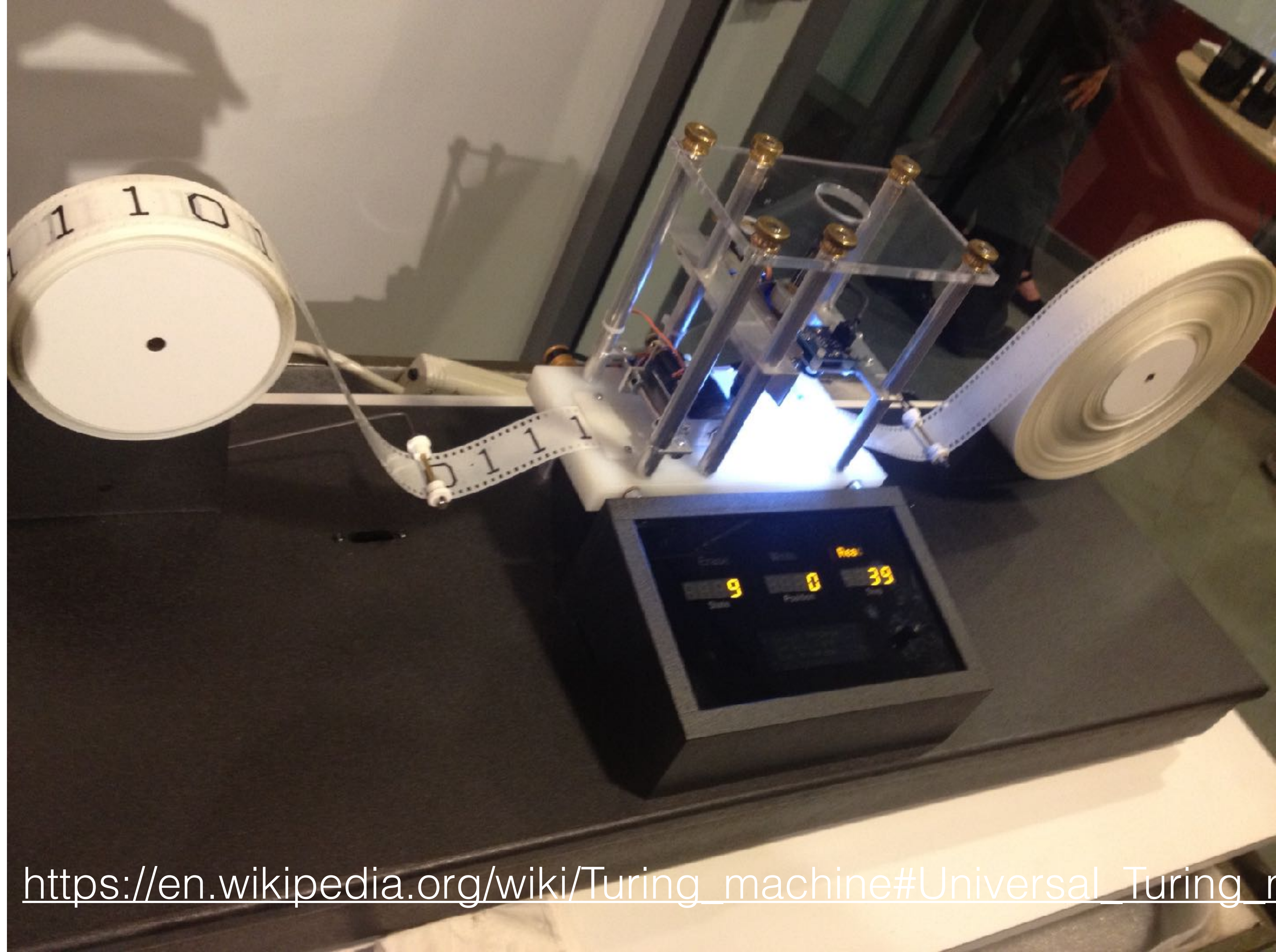
## RNN vs feedforward network

recurrent (RNN):  $\mathbf{h}_t = \mathbf{f}_w(\mathbf{x}_t, \mathbf{f}_w(\mathbf{x}_{t-1}, \dots \mathbf{f}_w(\mathbf{x}_1, \mathbf{h}_0)))$

feedforward:  $\mathbf{h}_t = \mathbf{g}(\mathbf{x}_t, \mathbf{x}_{t-1}, \dots \mathbf{x}_1, \mathbf{w})$  (stacking sequence to long input vector)

- RNN works for different lengths of input sequences
- RNN share weights between different time instances (similarly as convolution on spatial domain).
- Some RNN are spatio-temporal convolutions [Hinton 1988]
- Memory is attention in time [Alex Graves 2020]
- RNN is universal (can compute any function computable by Turing machine)

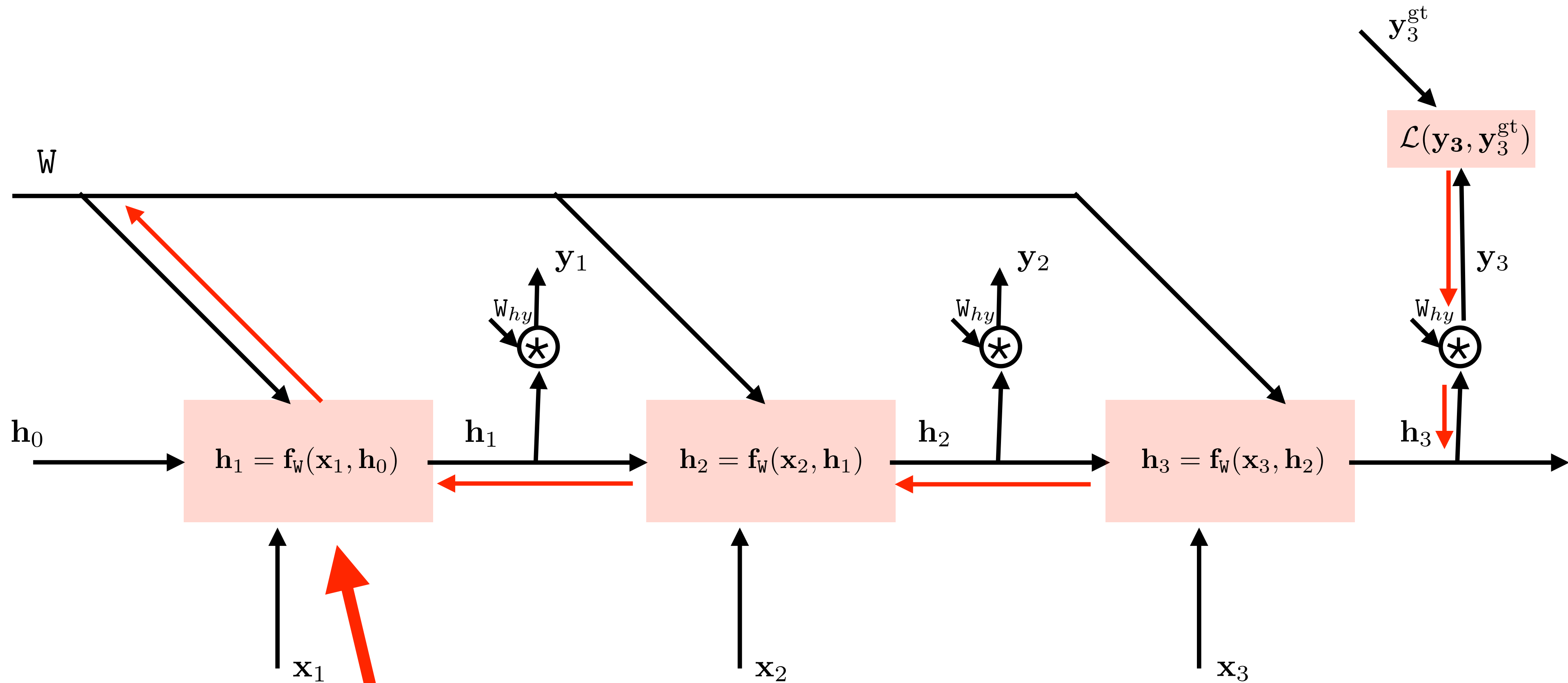




[https://en.wikipedia.org/wiki/Turing\\_machine#Universal\\_Turing\\_m](https://en.wikipedia.org/wiki/Turing_machine#Universal_Turing_m)



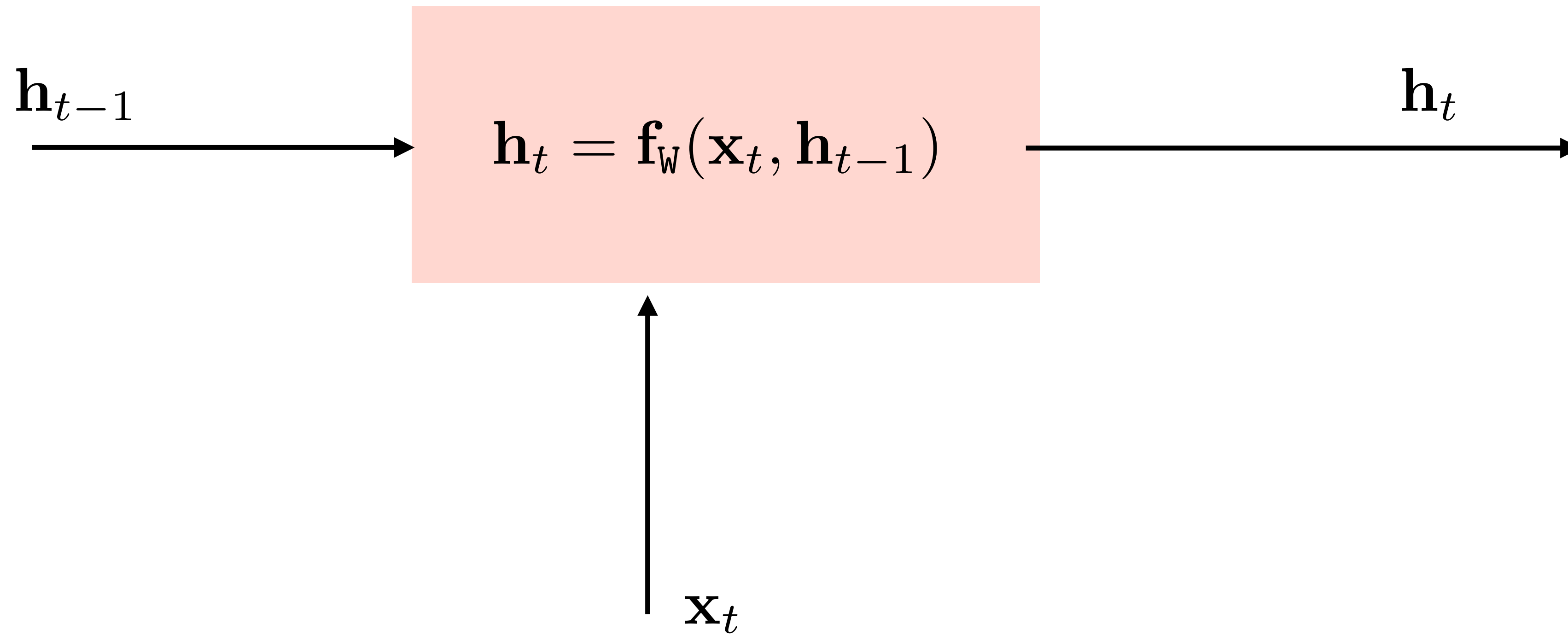
# Simple recurrent block - backward pass



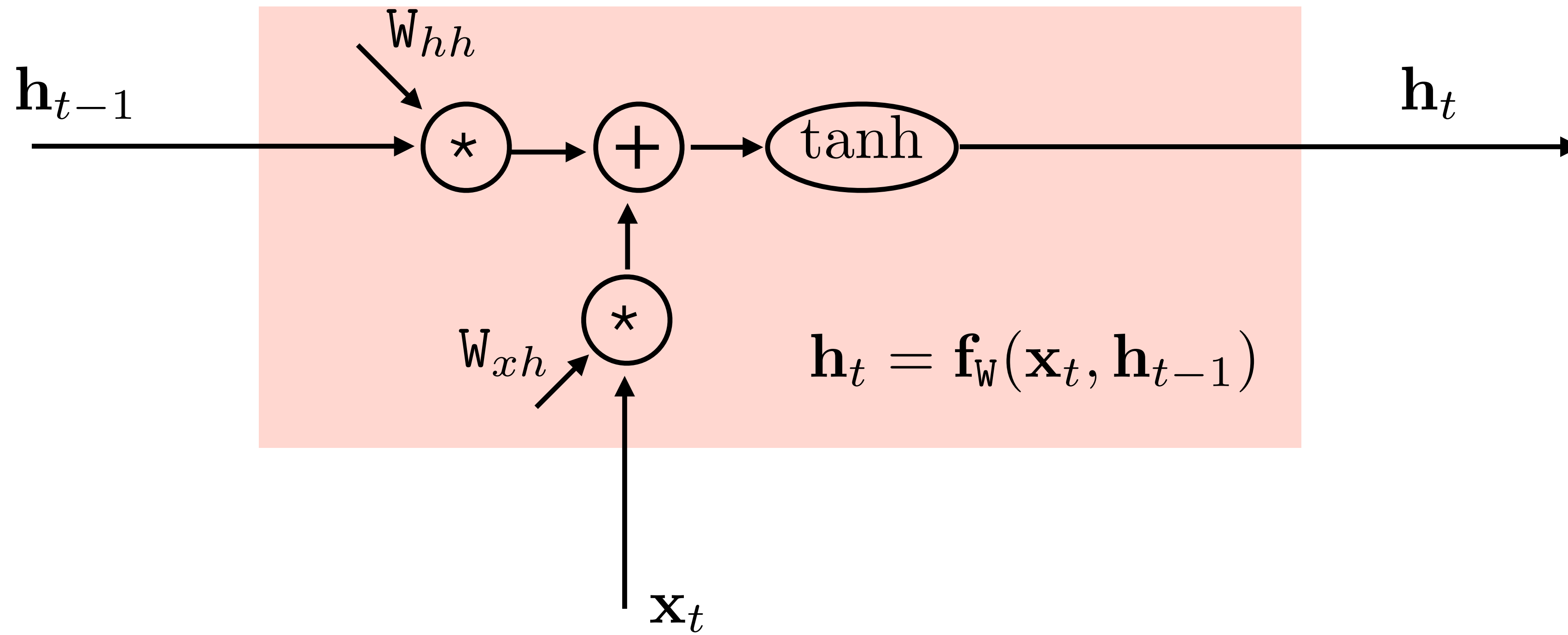
deep blocks often suffer from vanishing gradient  
=> better structure needed

LSTM (kind of ResNet for recurrent networks)

# Simple recurrent block

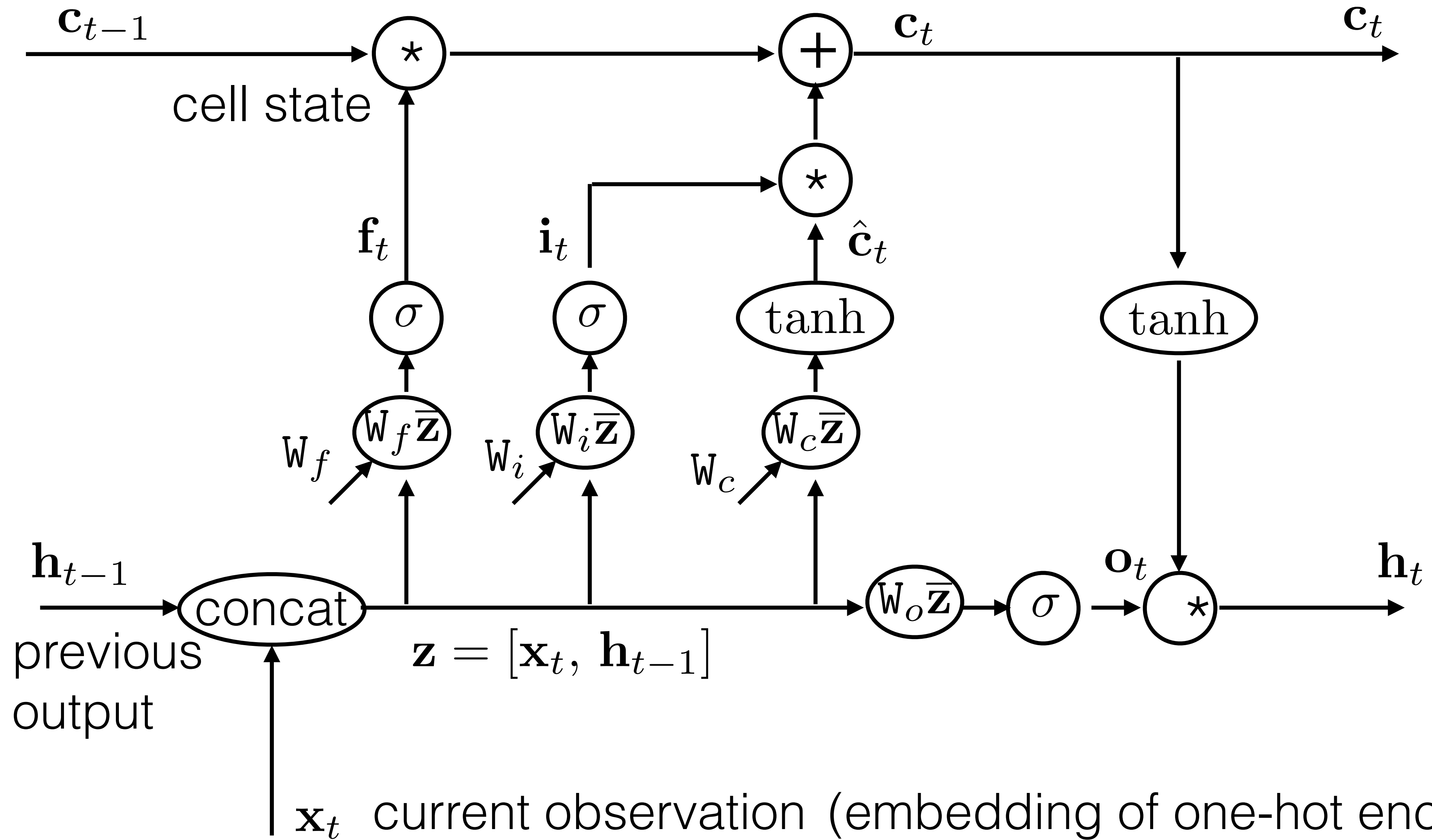


# Simple recurrent block



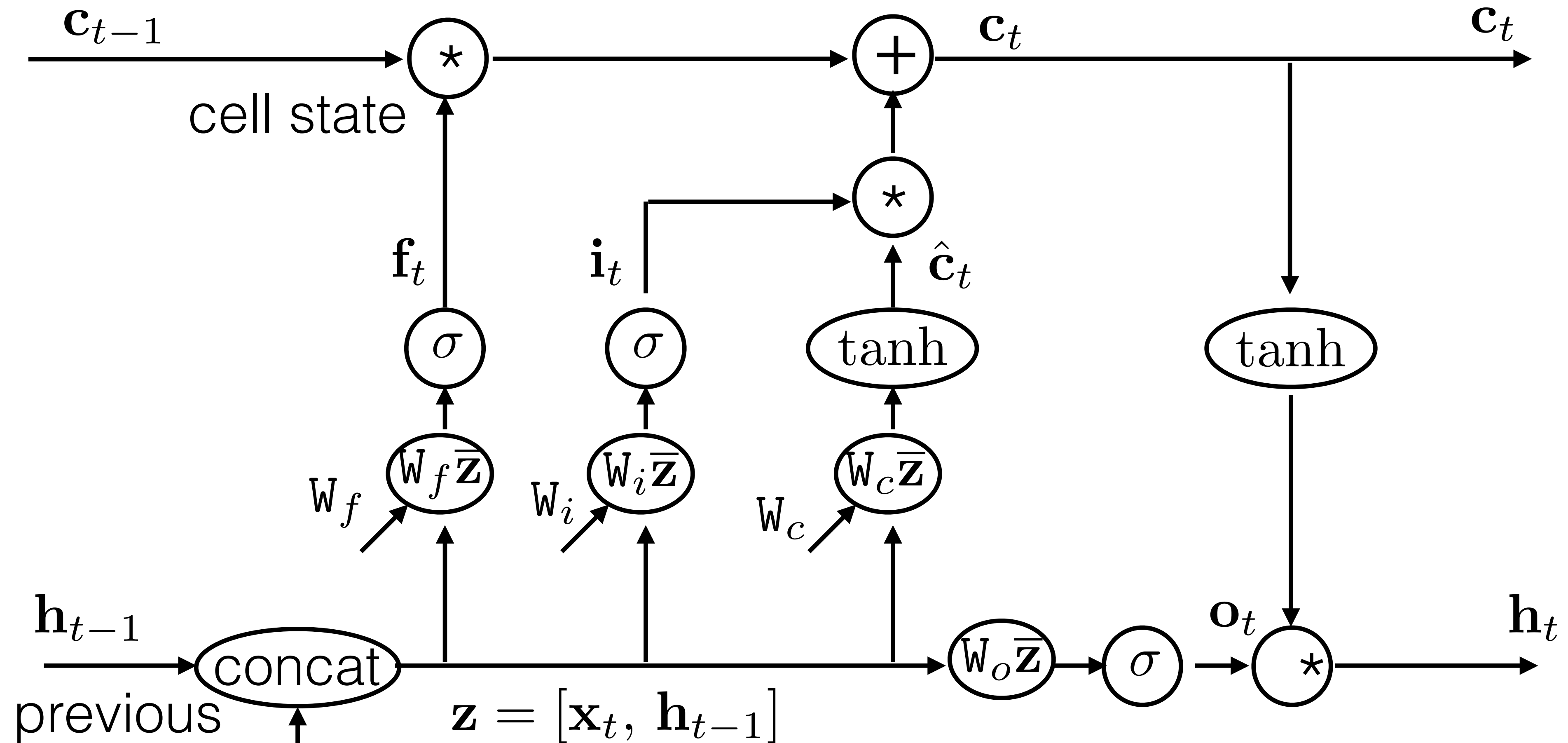


# LSTM block



“I live with my parents, ...”

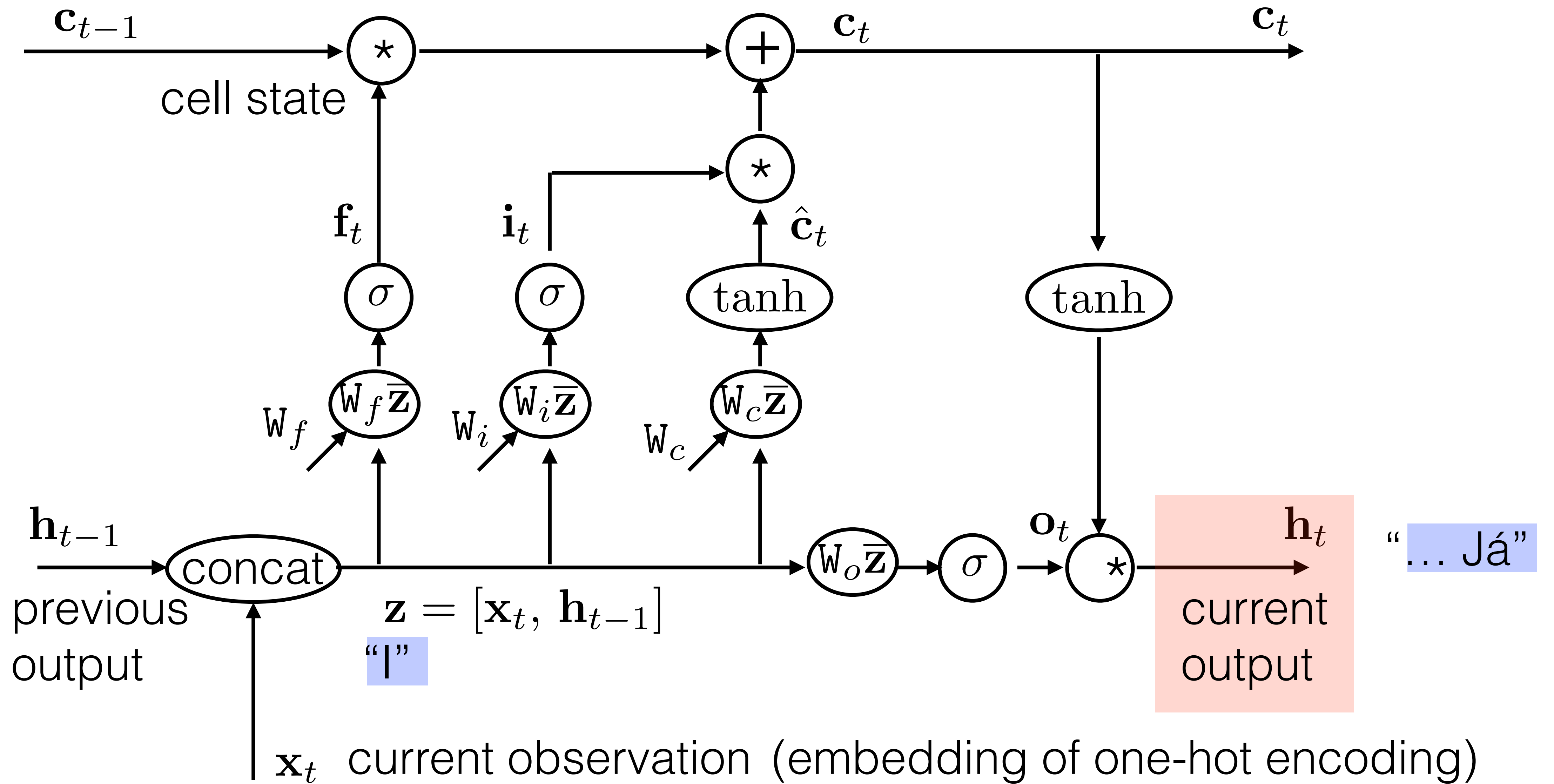
# LSTM block



$\mathbf{x}_t$  current observation (embedding of one-hot encoding)

“I live with my parents, ...”

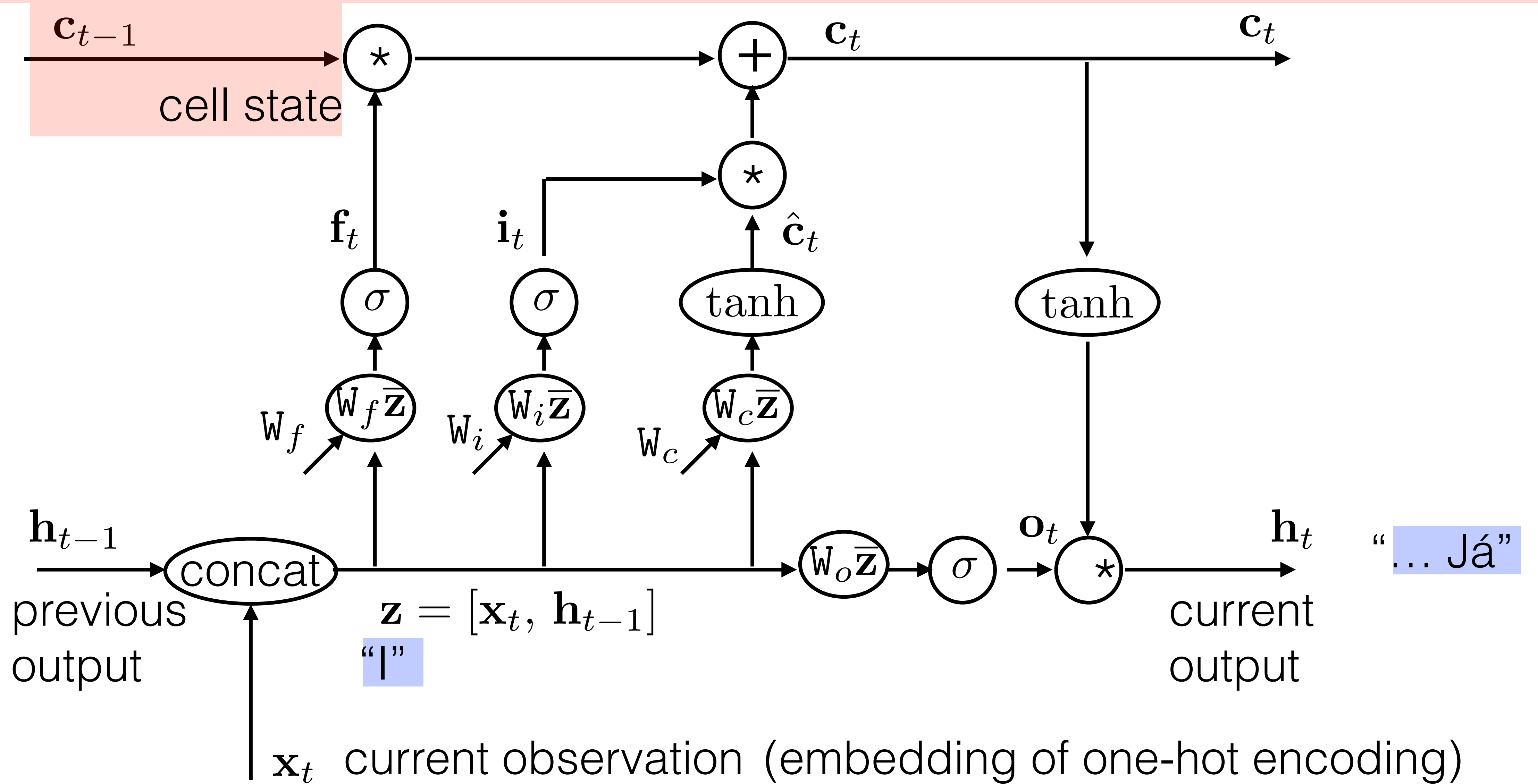
# LSTM block



"I live with my parents, ..."

# LSTM block

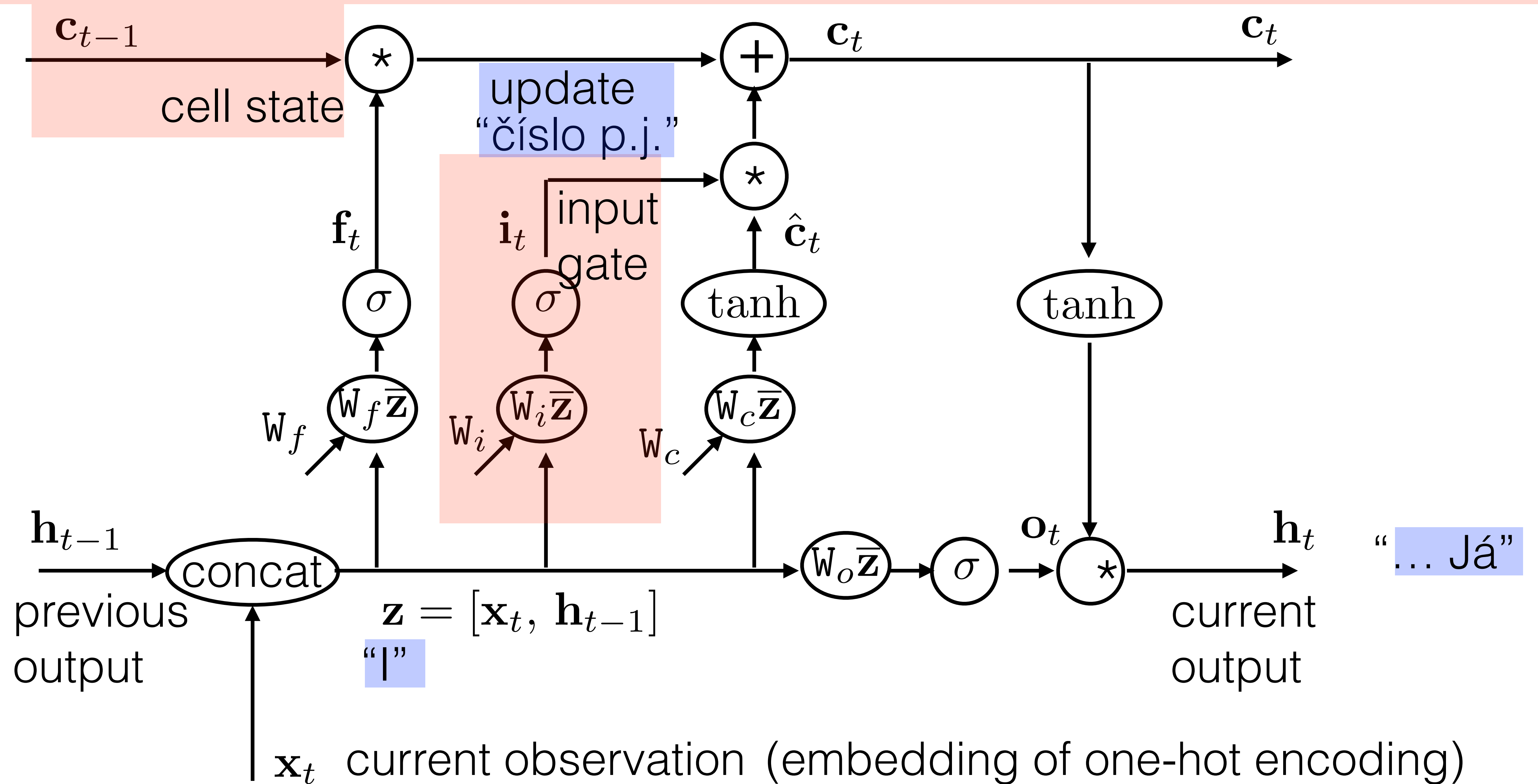
cell state is long term memory= vector [singular/plural, feminine/masculine, case, ...]



"I live with my parents, ..."

# LSTM block

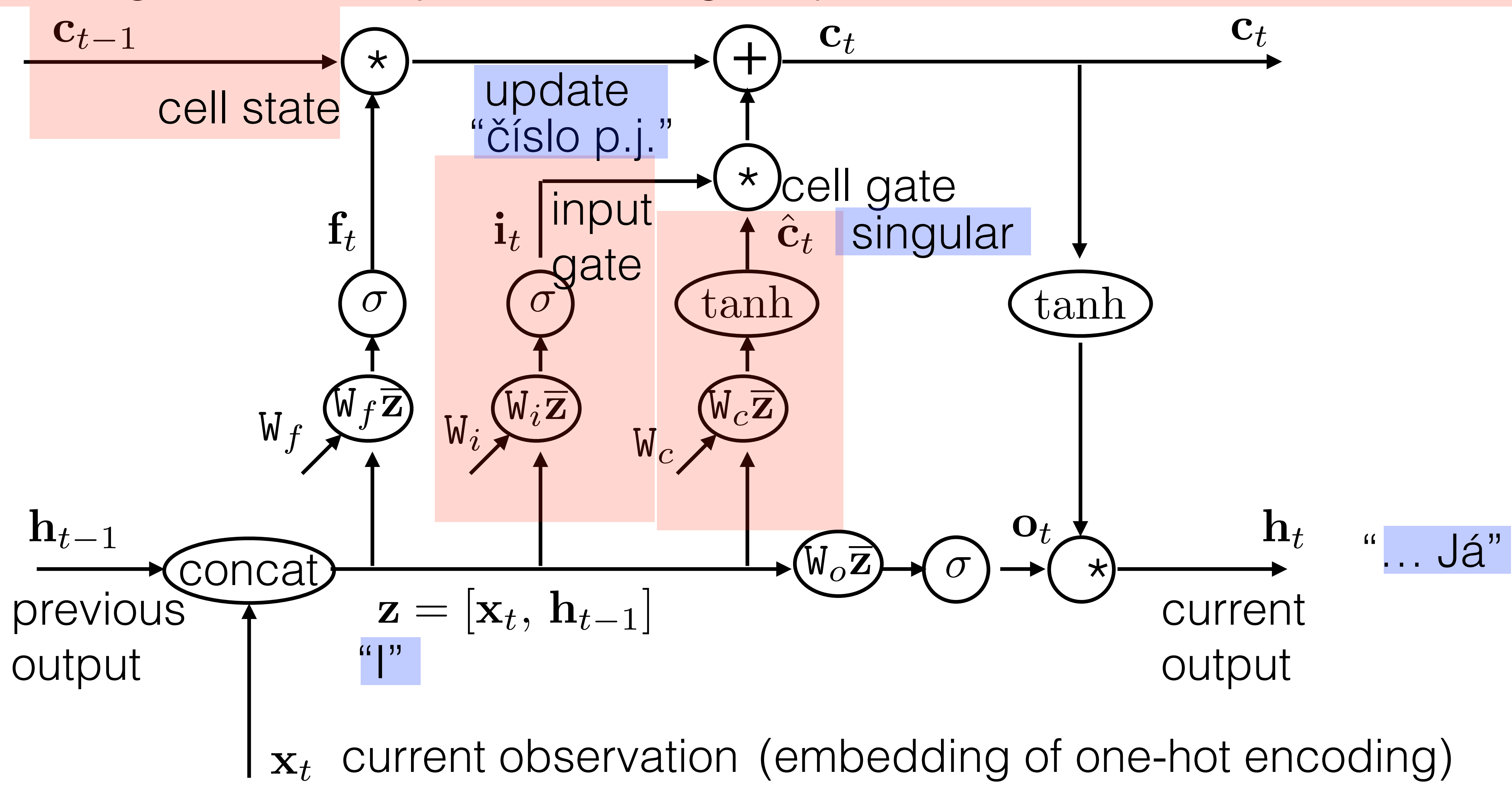
cell state is long term memory= vector [singular/plural, feminine/masculine, case, ...]



"I live with my parents, ..."

# LSTM block

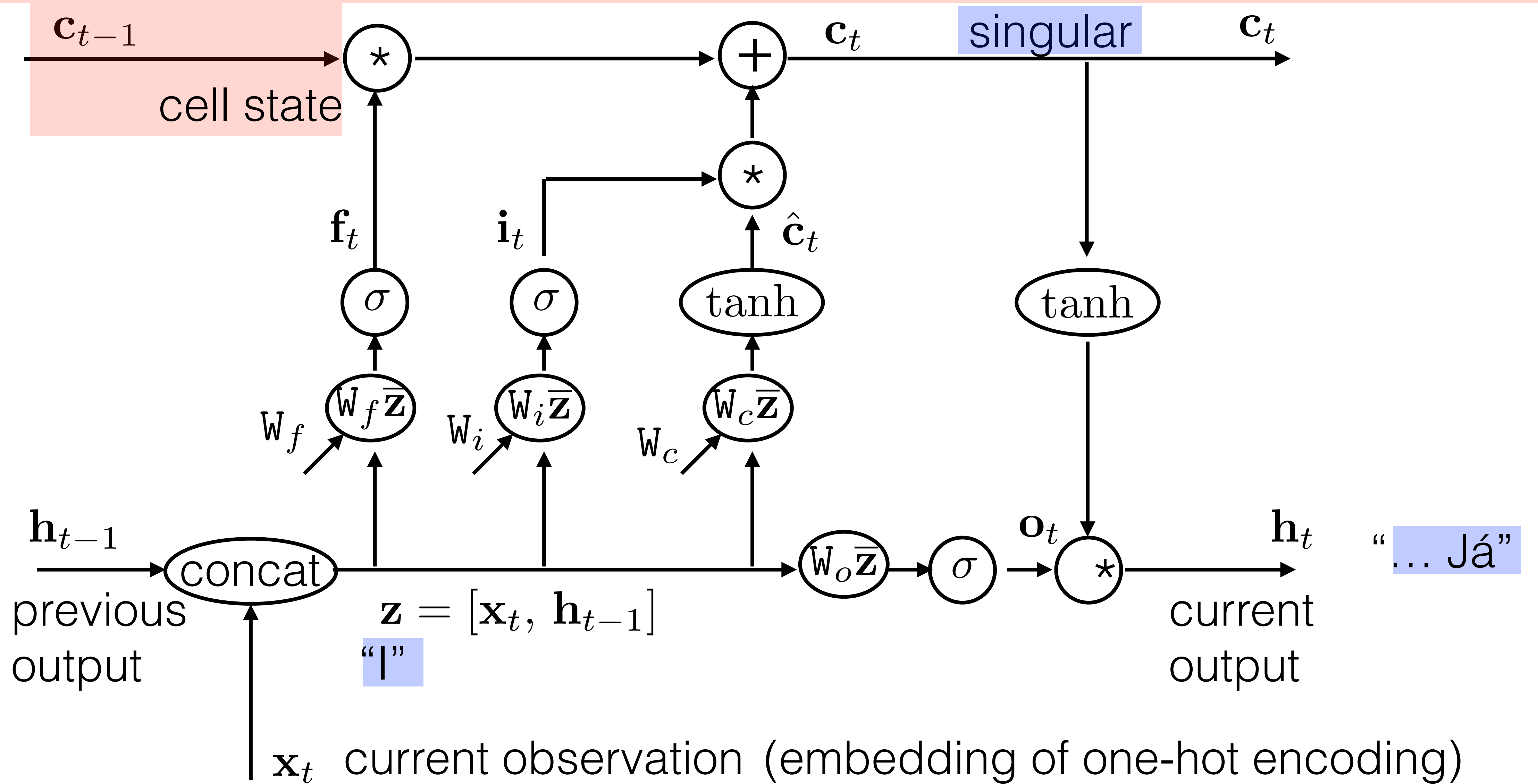
cell state is long term memory= vector [singular/plural, feminine/masculine, case, ...]



"I live with my parents, ..."

# LSTM block

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

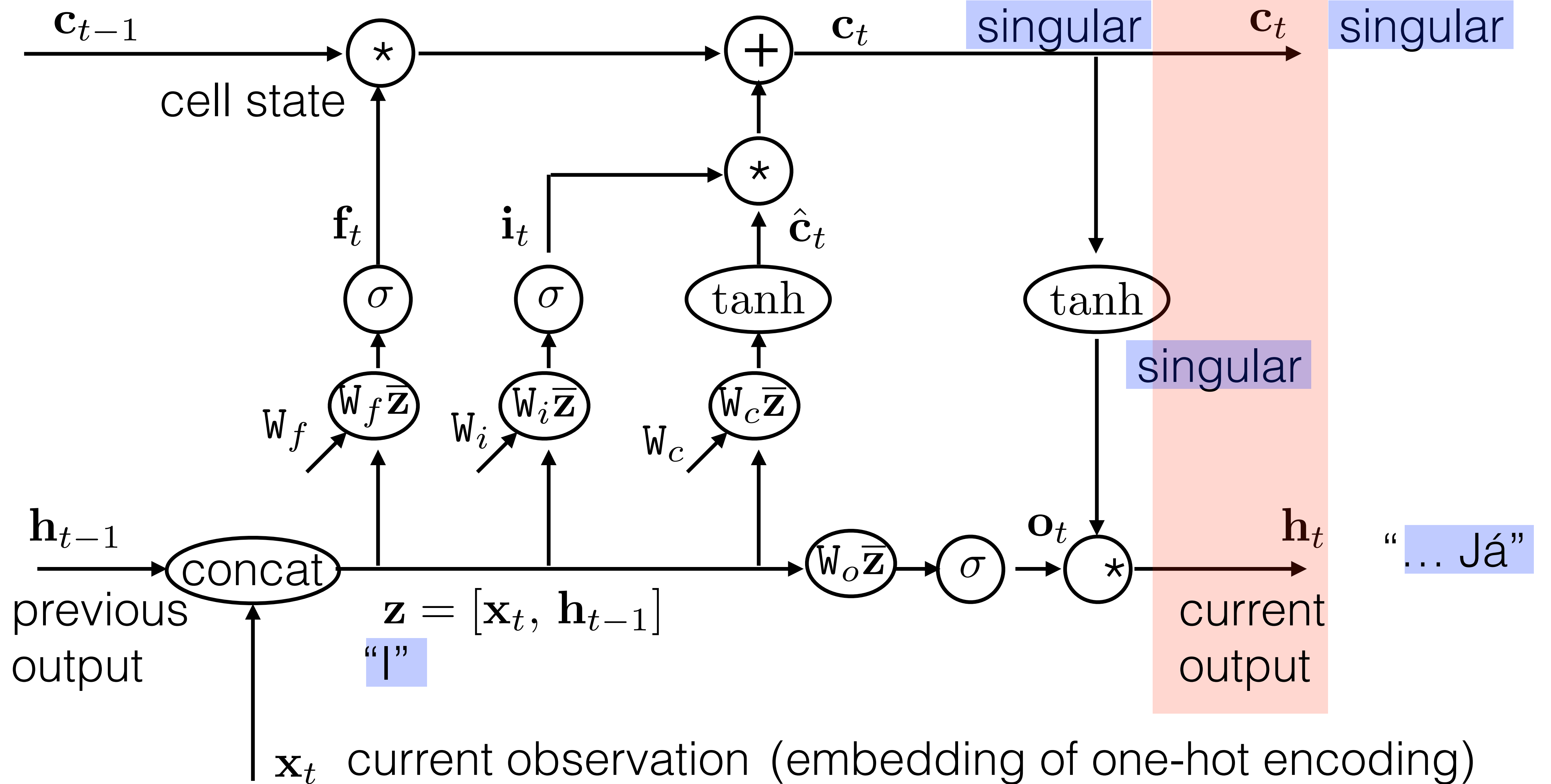


"I live with my parents, ..."



# LSTM block

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

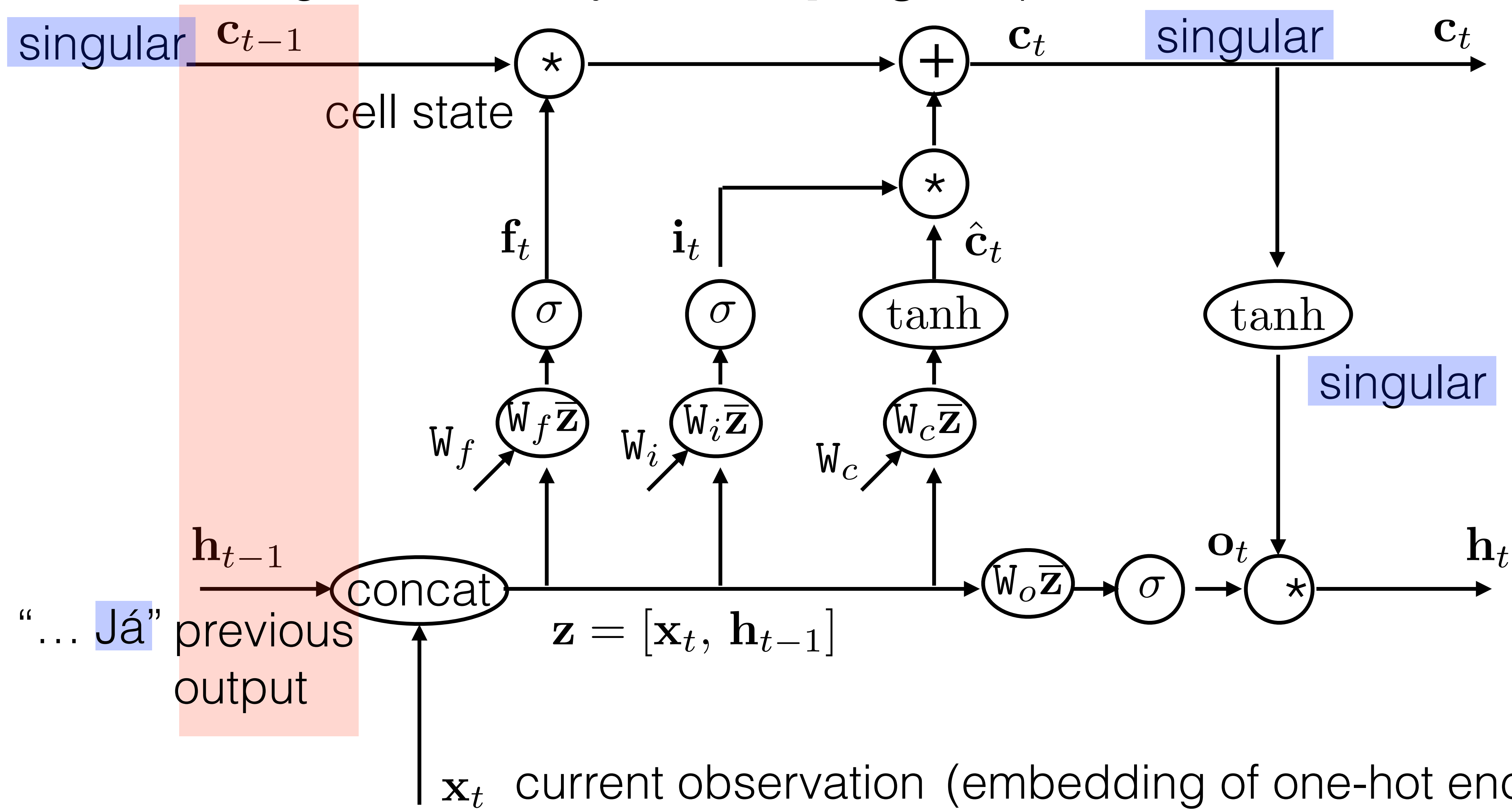


“I live with my parents, ...”



# LSTM block

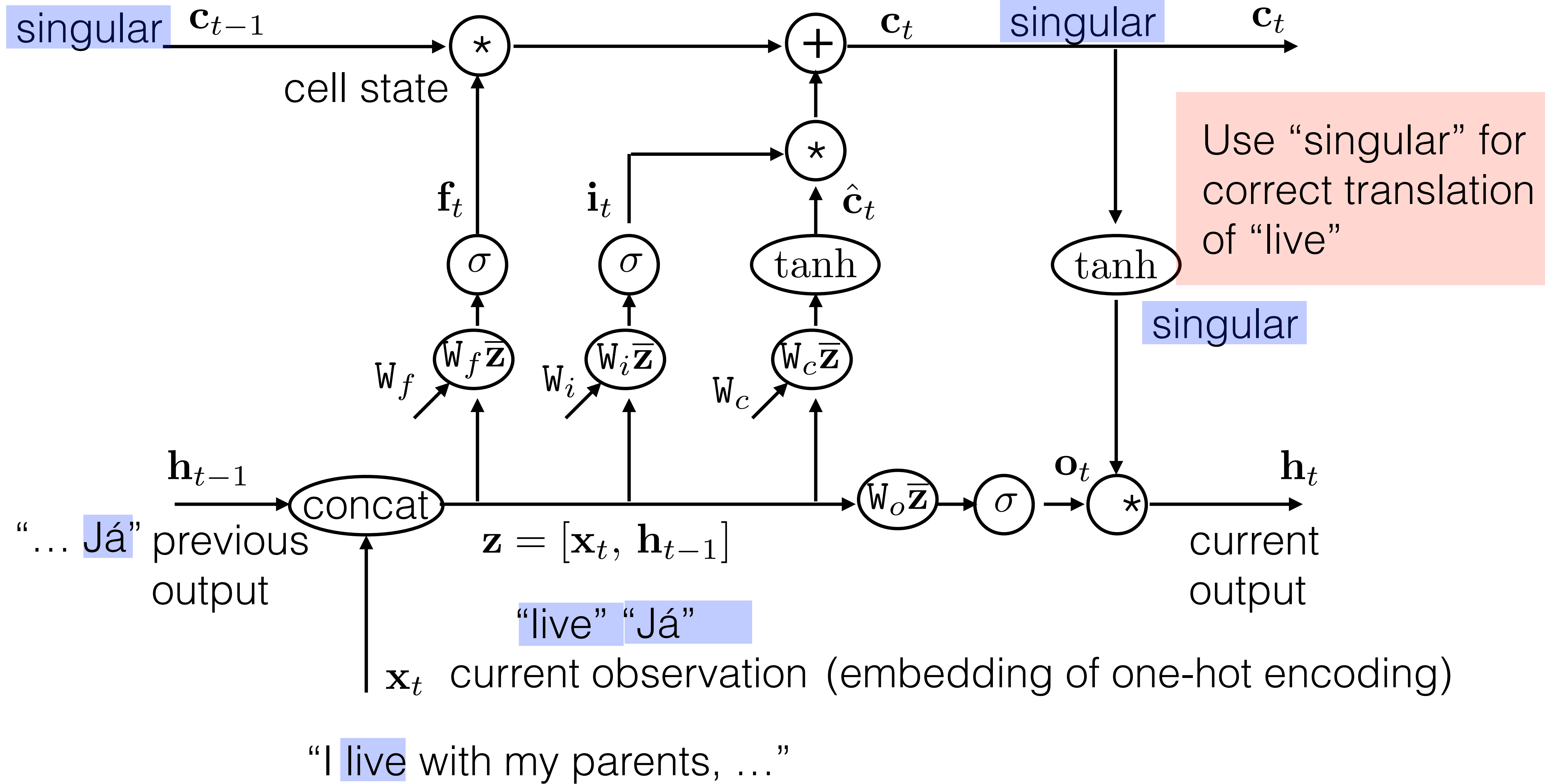
cell state is long term memory= vector [**singular**/plural, feminine/masculine, case, ...]



“I live with my parents, ...”

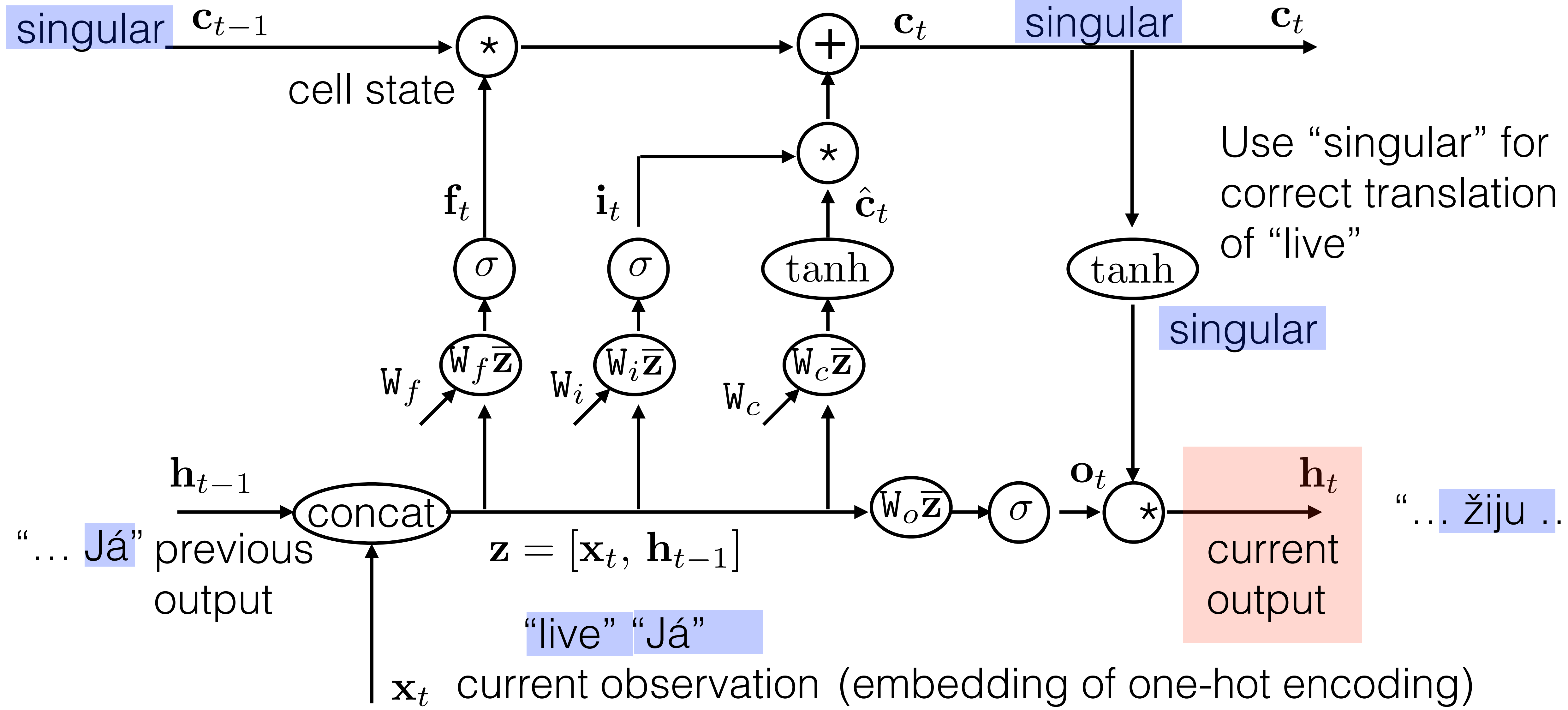
# LSTM block

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



# LSTM block

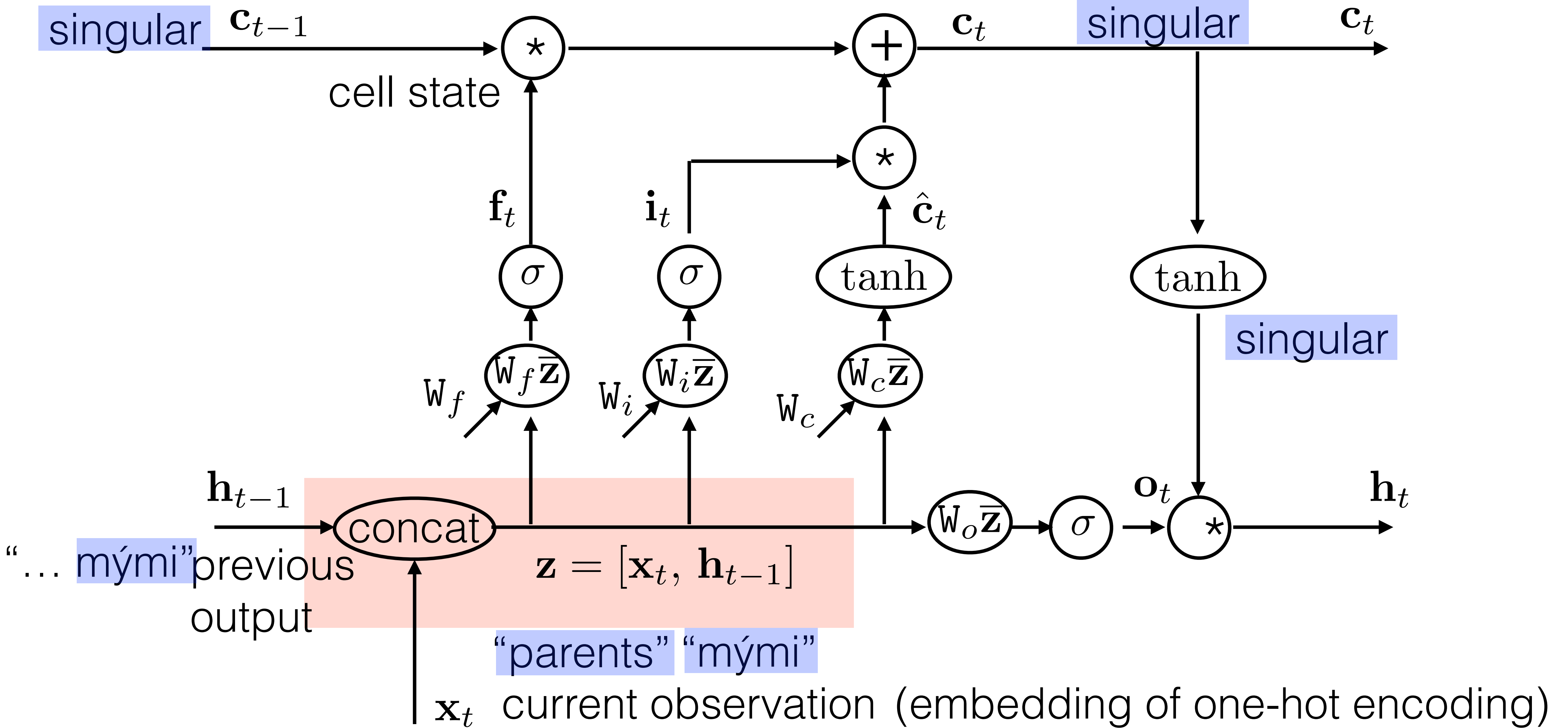
cell state is long term memory= vector [**singular**/plural, feminine/masculine, case, ...]



"I live with my parents, ..."

# LSTM block

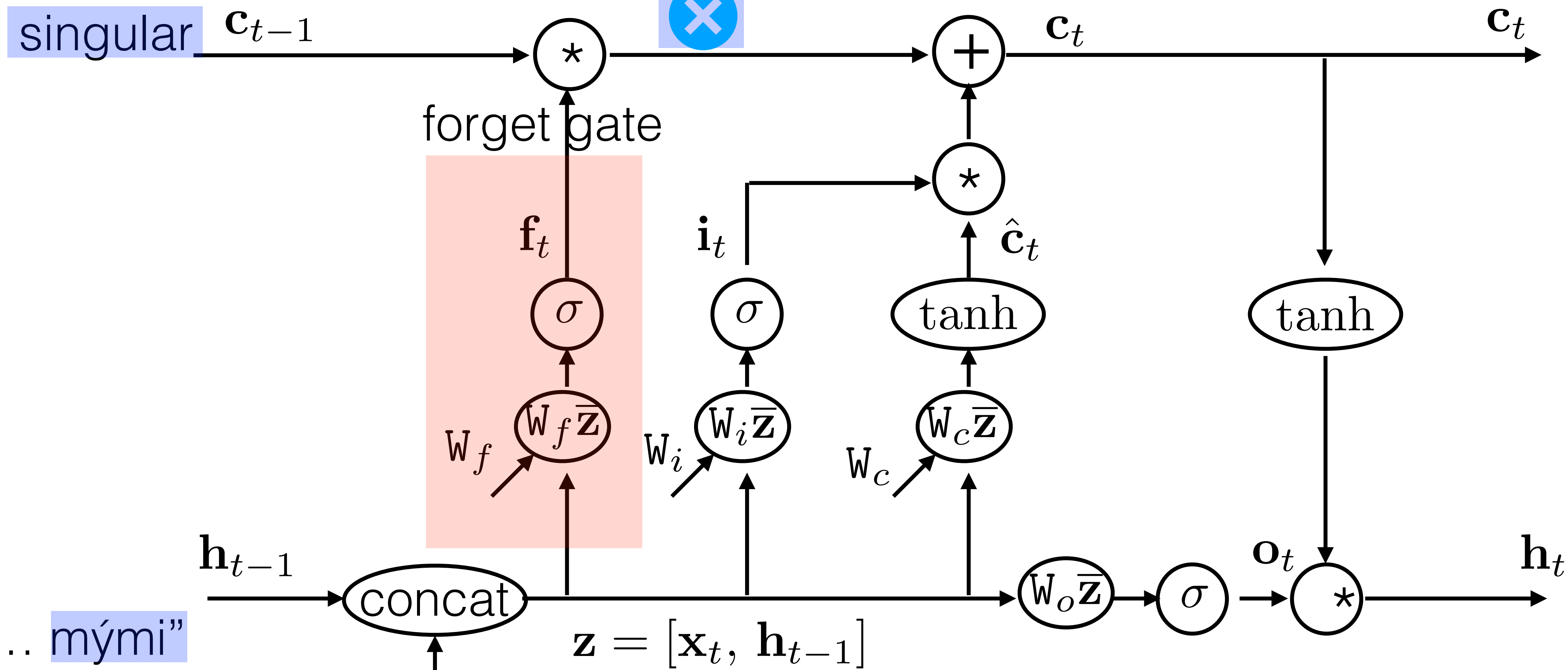
cell state is long term memory= vector [**singular**/plural, feminine/masculine, case, ...]



"I live with my **parents**, ..."

# LSTM block

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



"... mými"

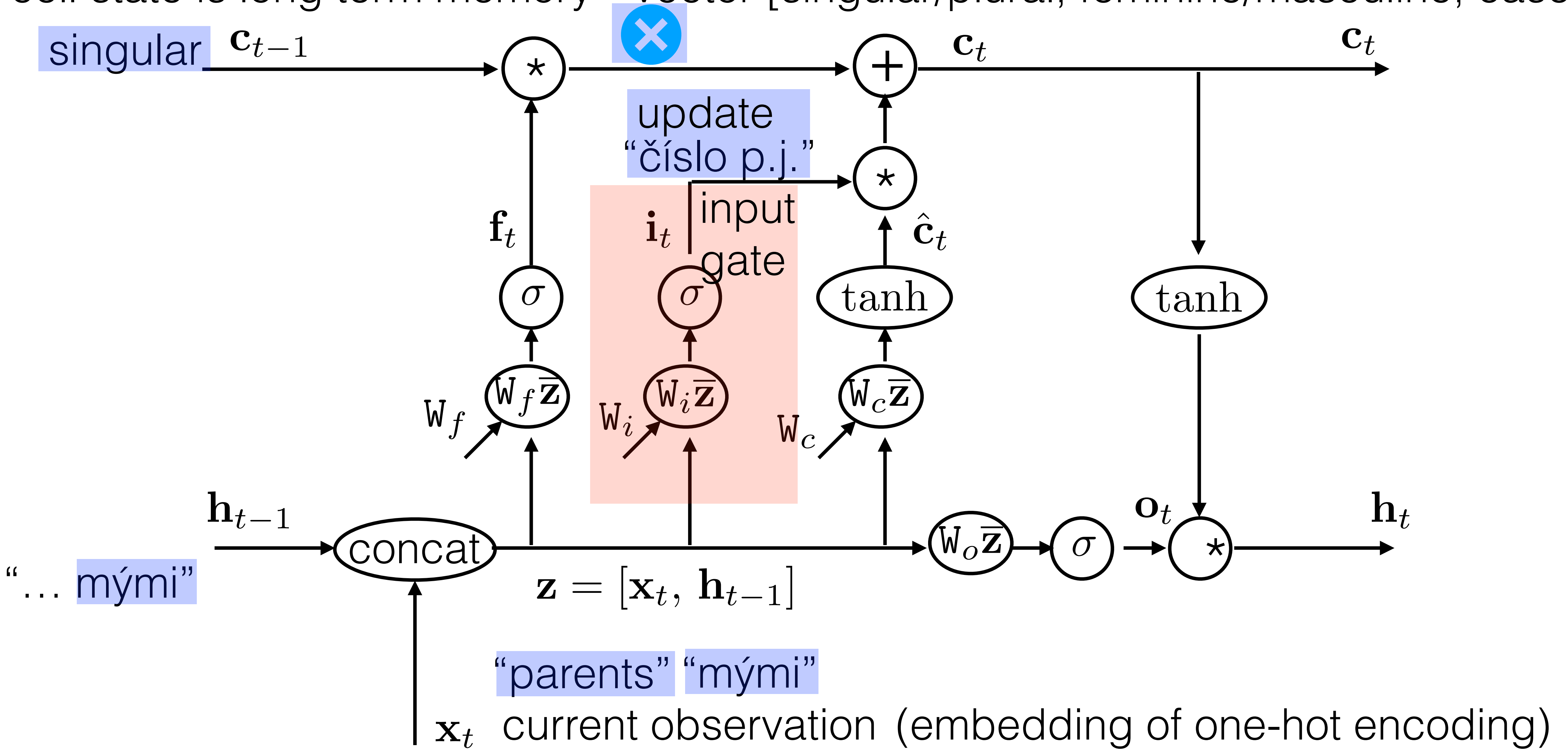
"parents" "mými"

$\mathbf{x}_t$  current observation (embedding of one-hot encoding)

"I live with my parents, ..."

# LSTM block

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

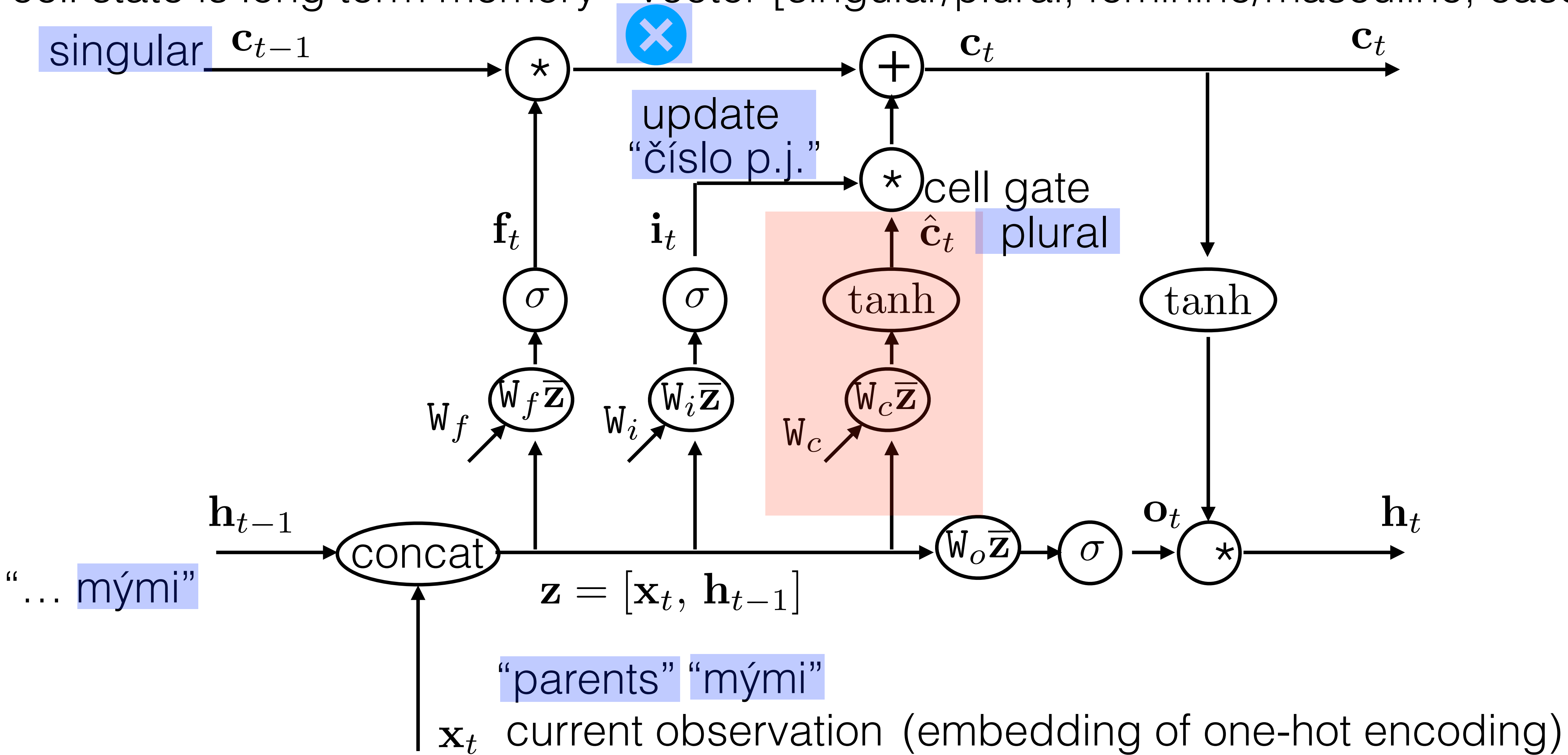


"I live with my parents, ..."



# LSTM block

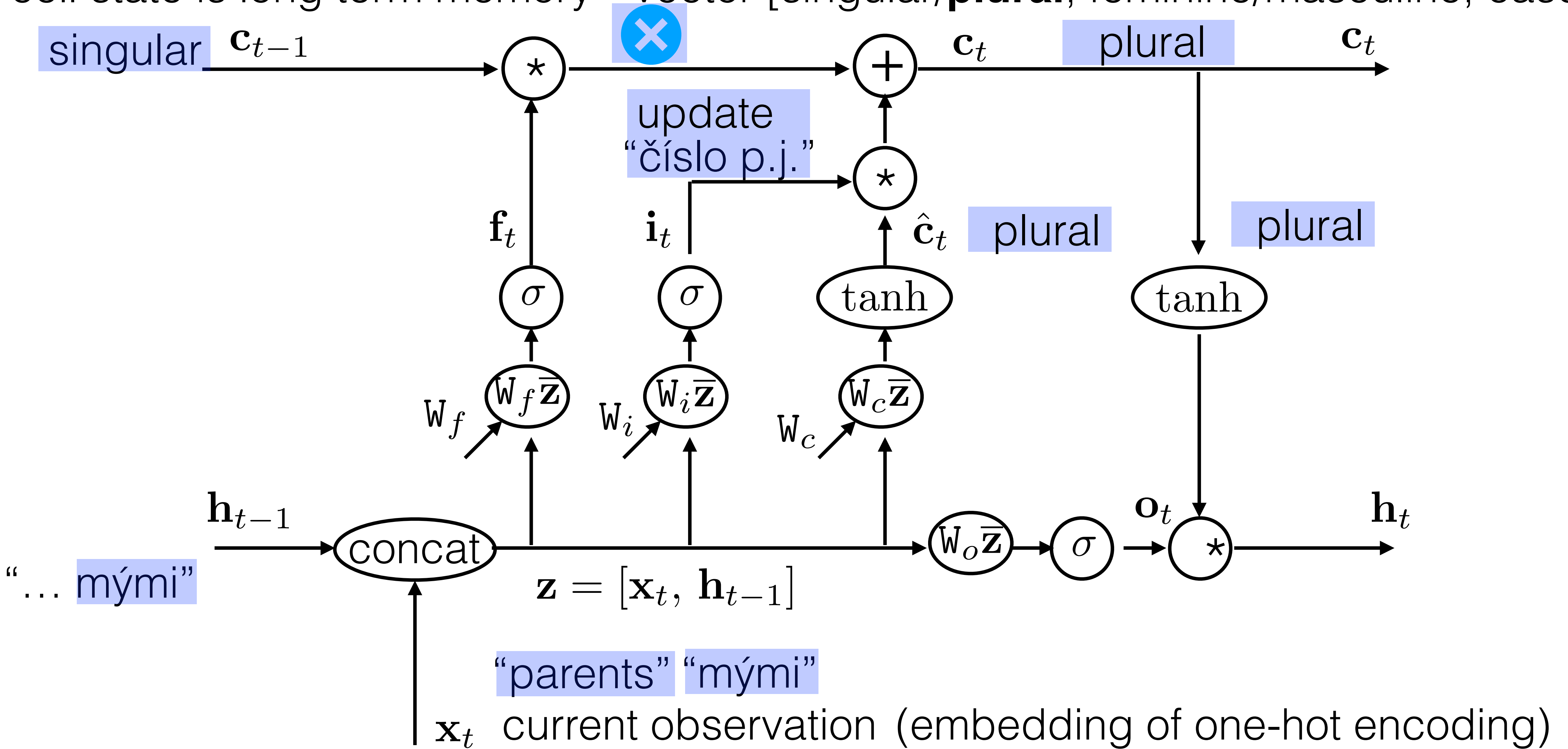
cell state is long term memory= vector [singular/plural, feminine/masculine, case, ...]



"I live with my parents, ..."

# LSTM block

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

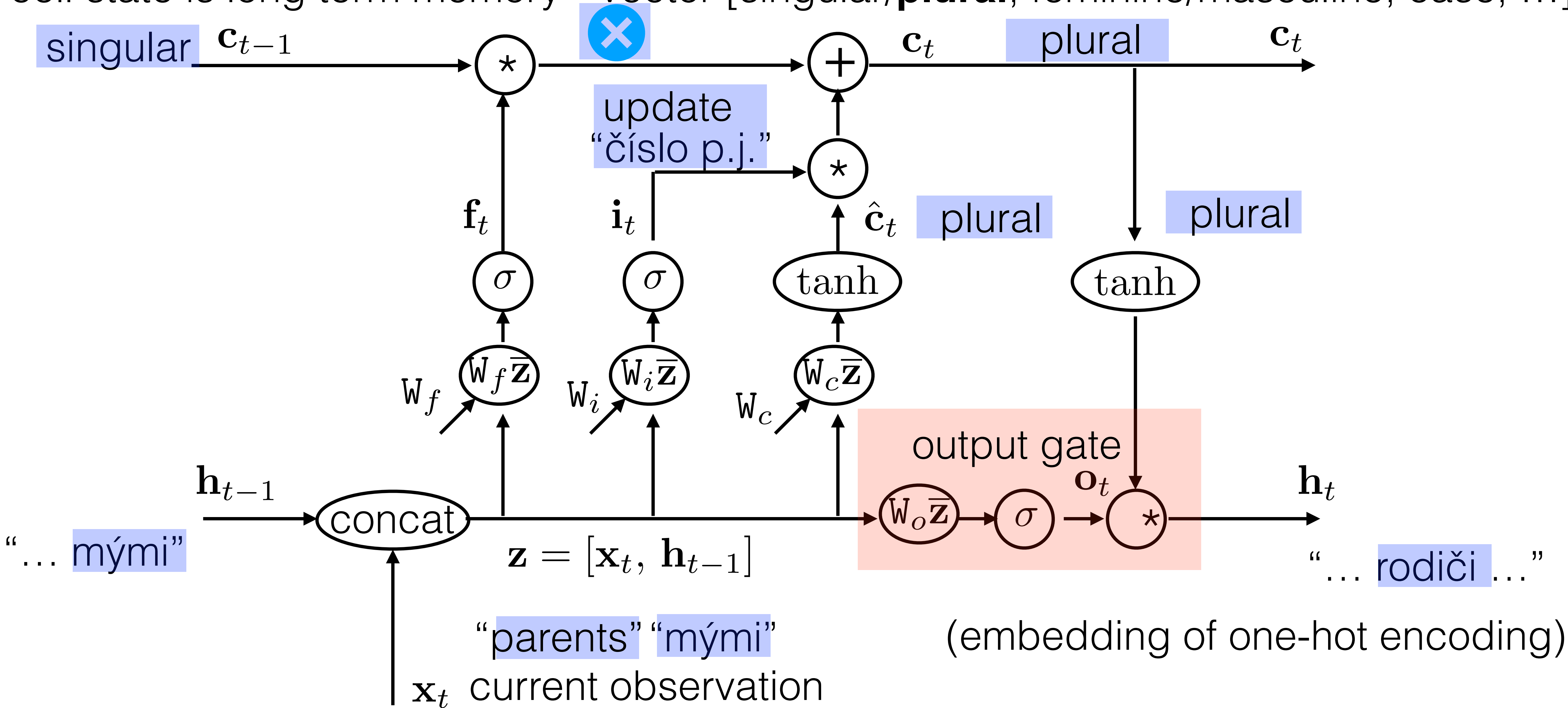


"I live with my **parents**, ..."



# LSTM block

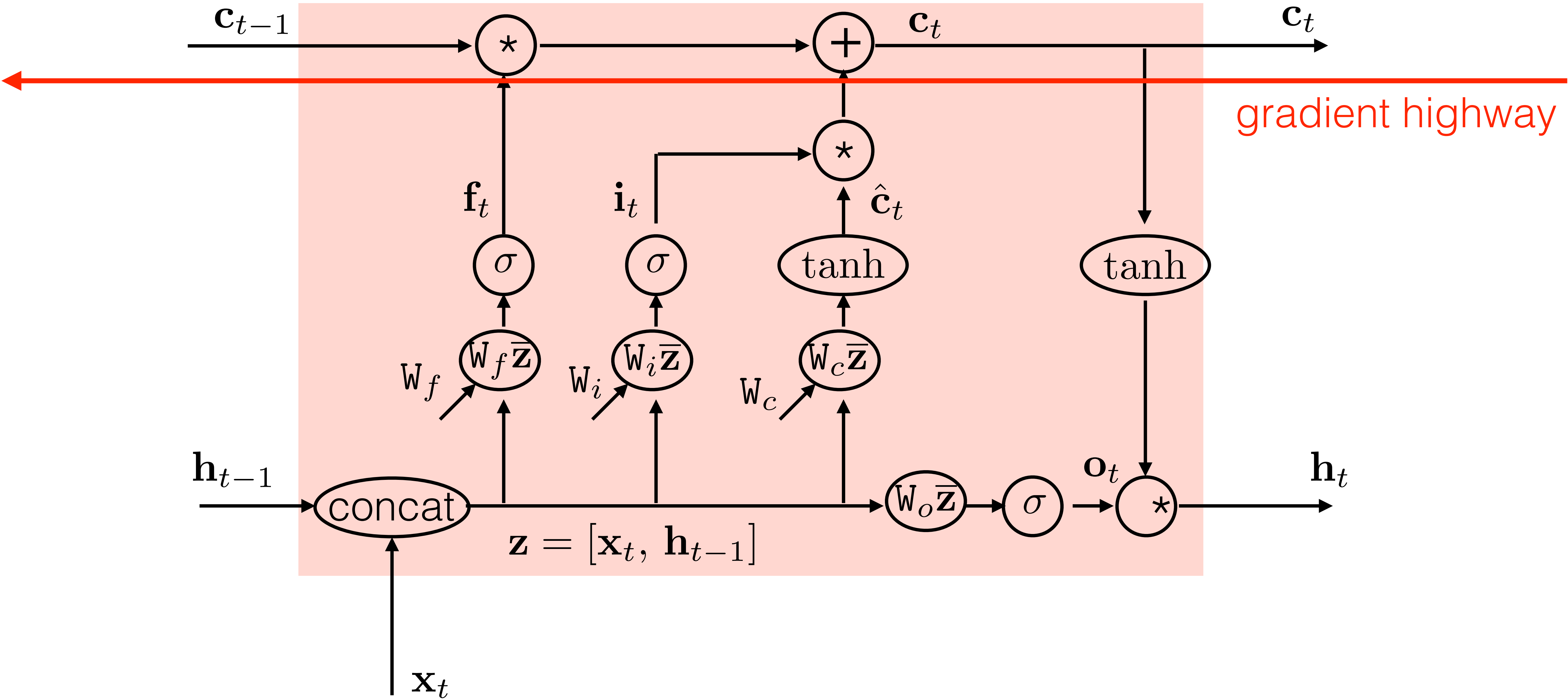
cell state is long term memory= vector [singular/**plural**, feminine/masculine, case, ...]



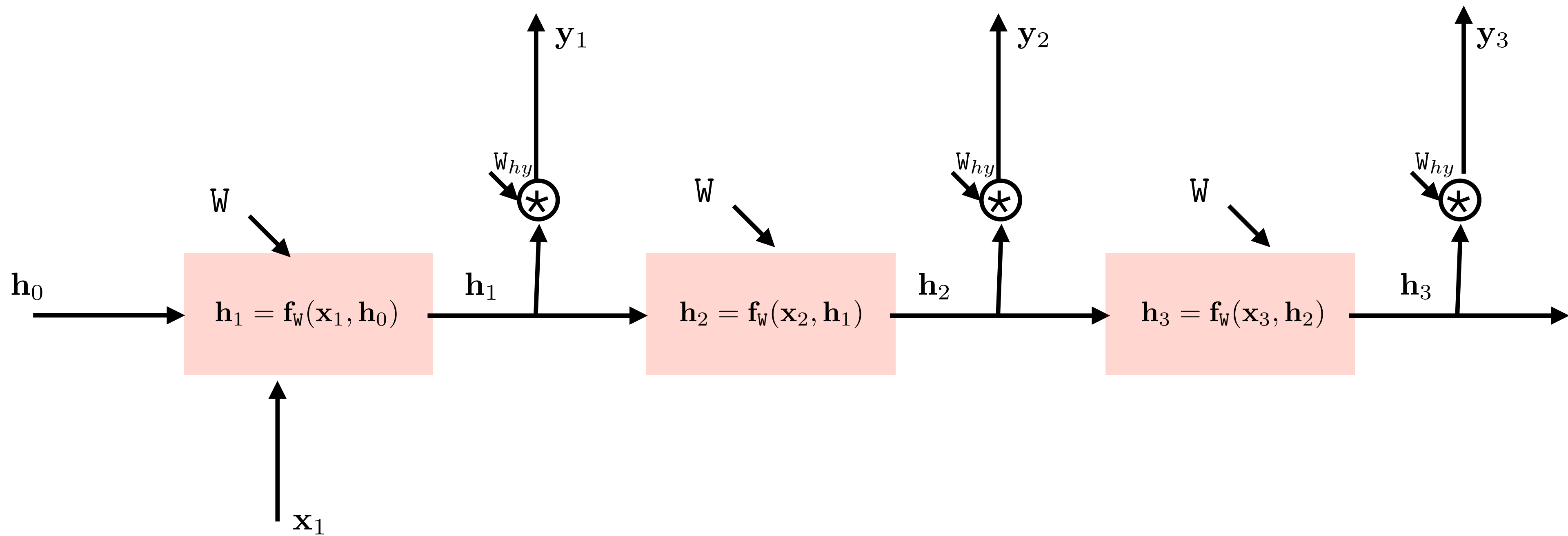
"I live with my **parents**, ..."

# LSTM block

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

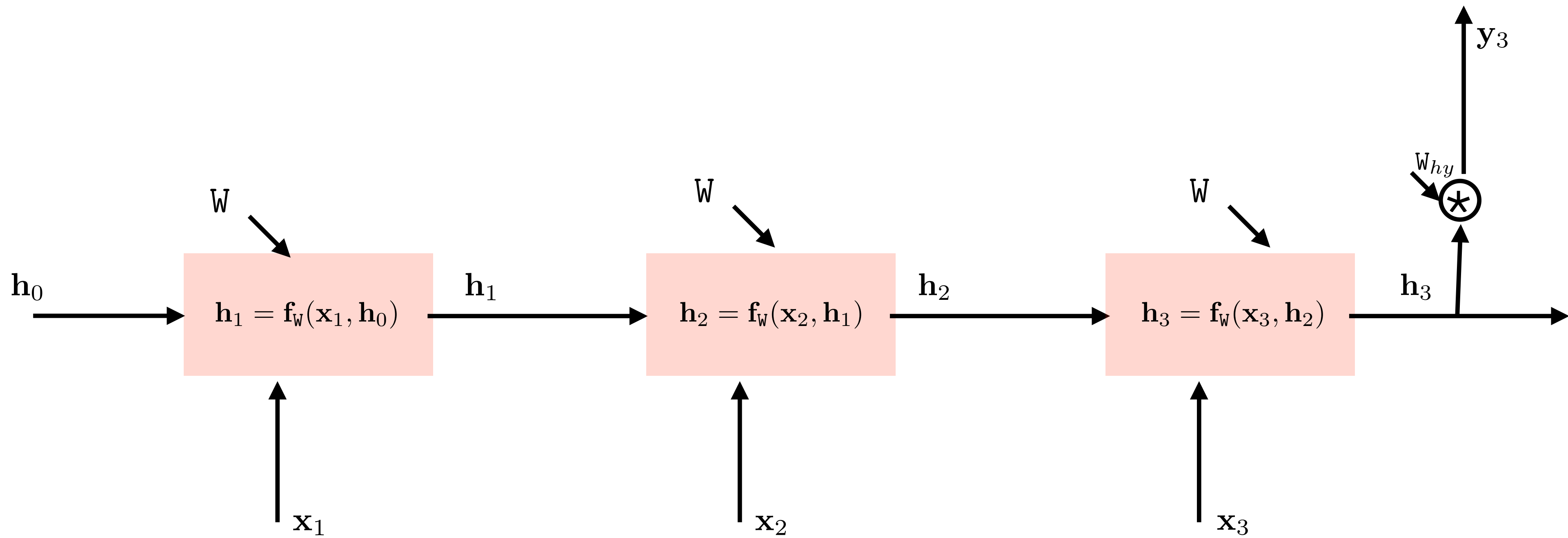


# RNN architectures: one-to-many



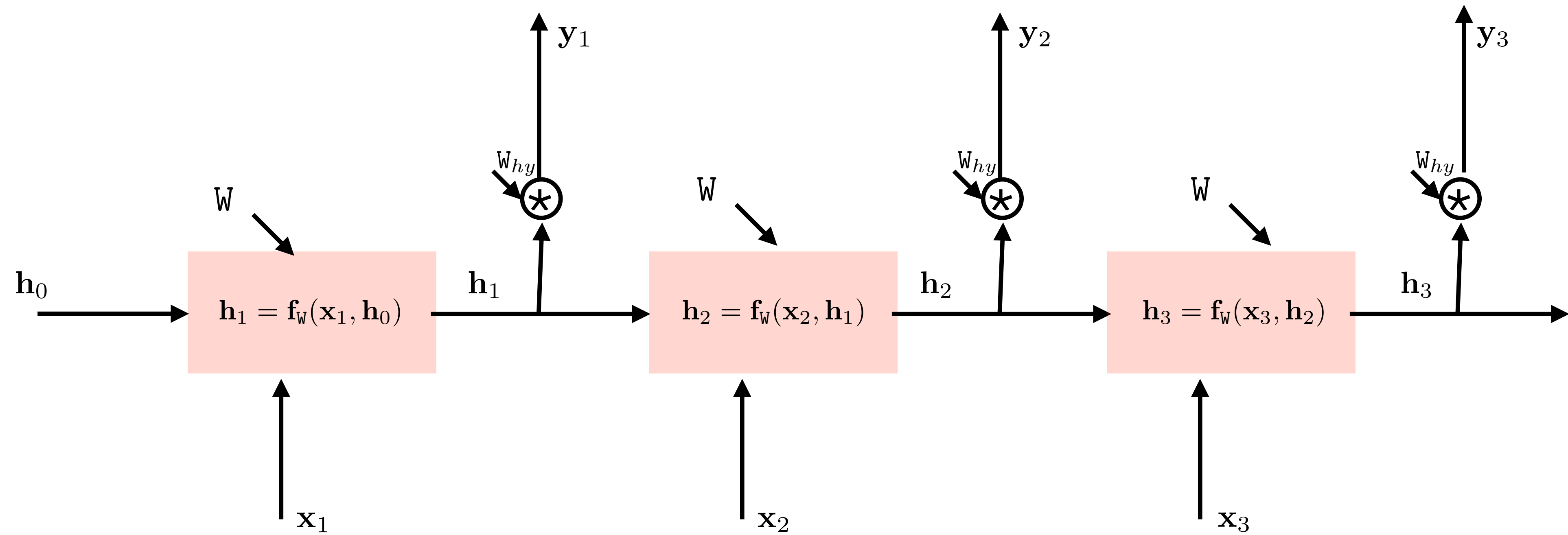
- Music generation
- Image captioning

# RNN architectures: one-to-many



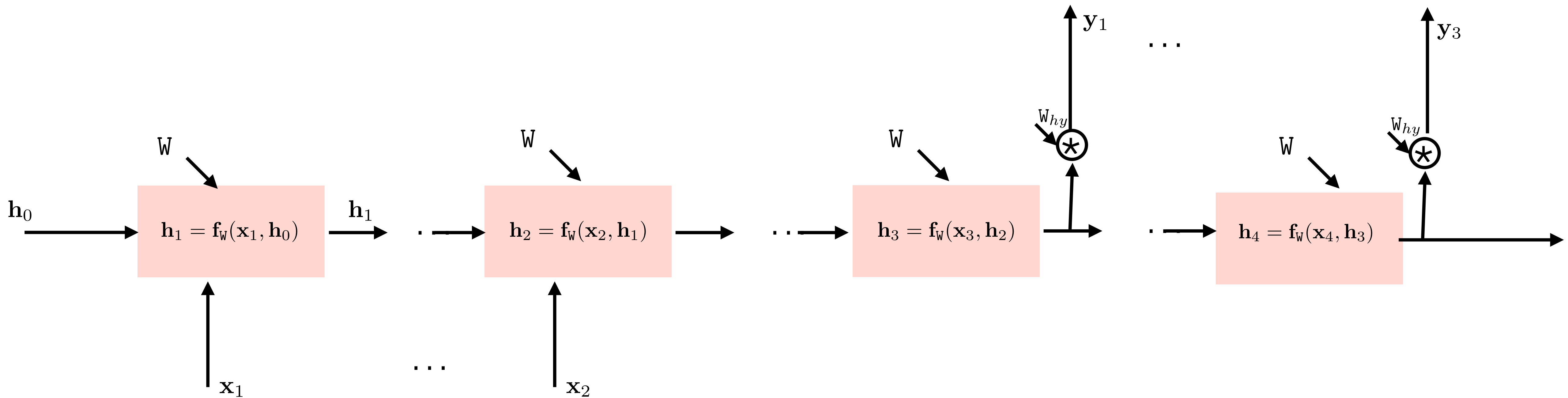
- Sentiment classification
- Action recognition

# RNN architectures: one-to-many



- Coloring/enhancing video sequences
- Named-entity recognition
- Speech recognition

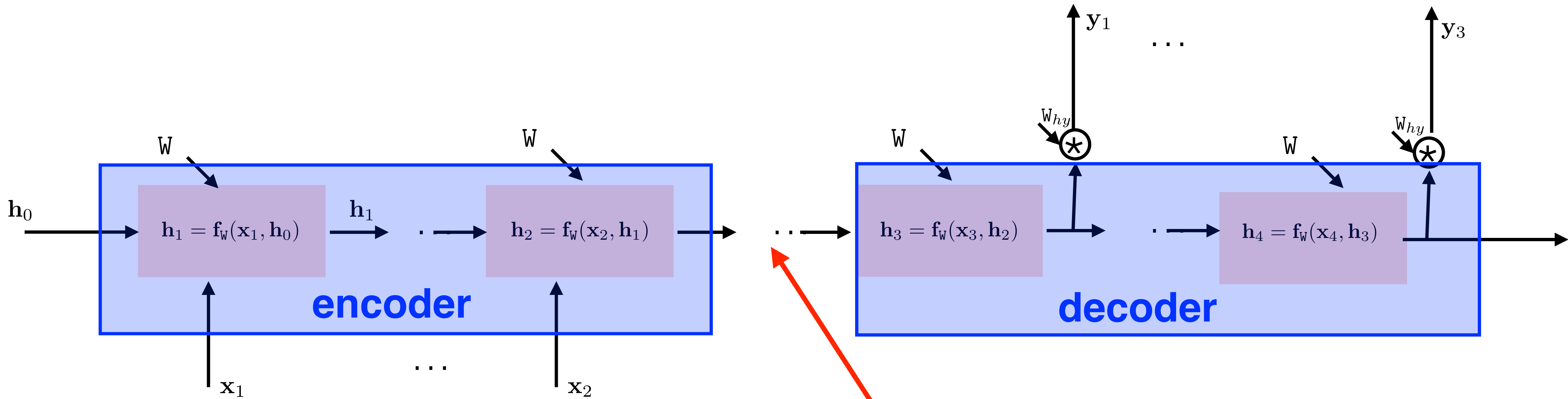
# RNN architectures: one-to-many



- Machine translation
- Question answering

# RNN architectures: many-to-many

**output:** variable-size sequence

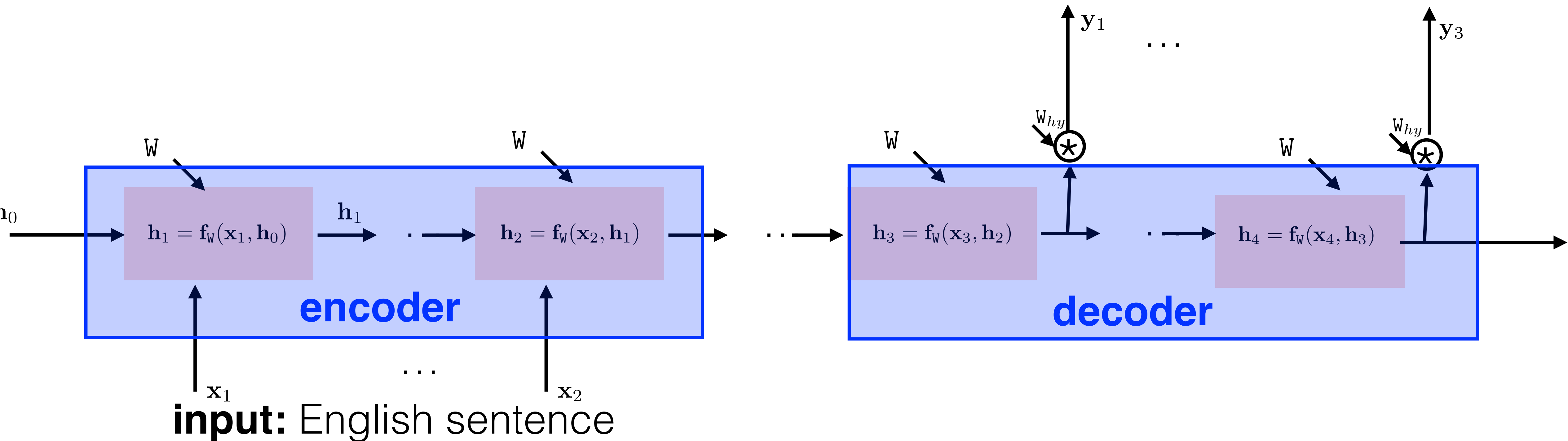


**input:** variable-size sequence

**context:** fixed-size semantic summary of input sequence

# RNN architectures: Machine translation

**output:** Czech sentence

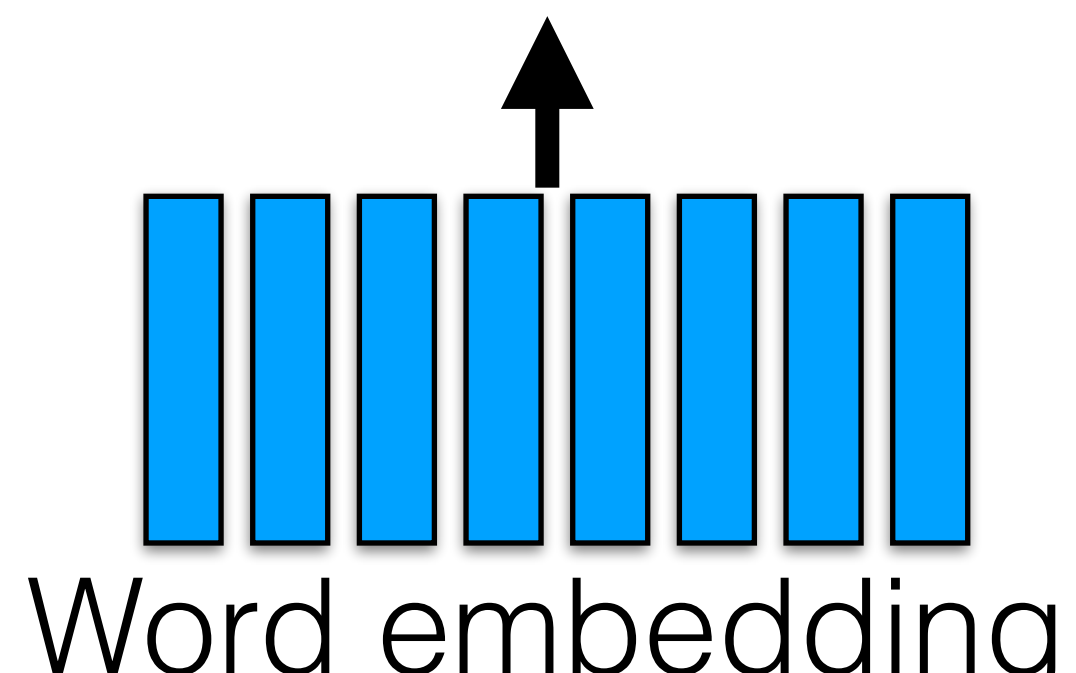


RNN can theoretically **remember everything** important but in practice it suffers from **catastrophic forgetting** (important relations could be very far).

Let's process the whole sentence at once through **transformer!**

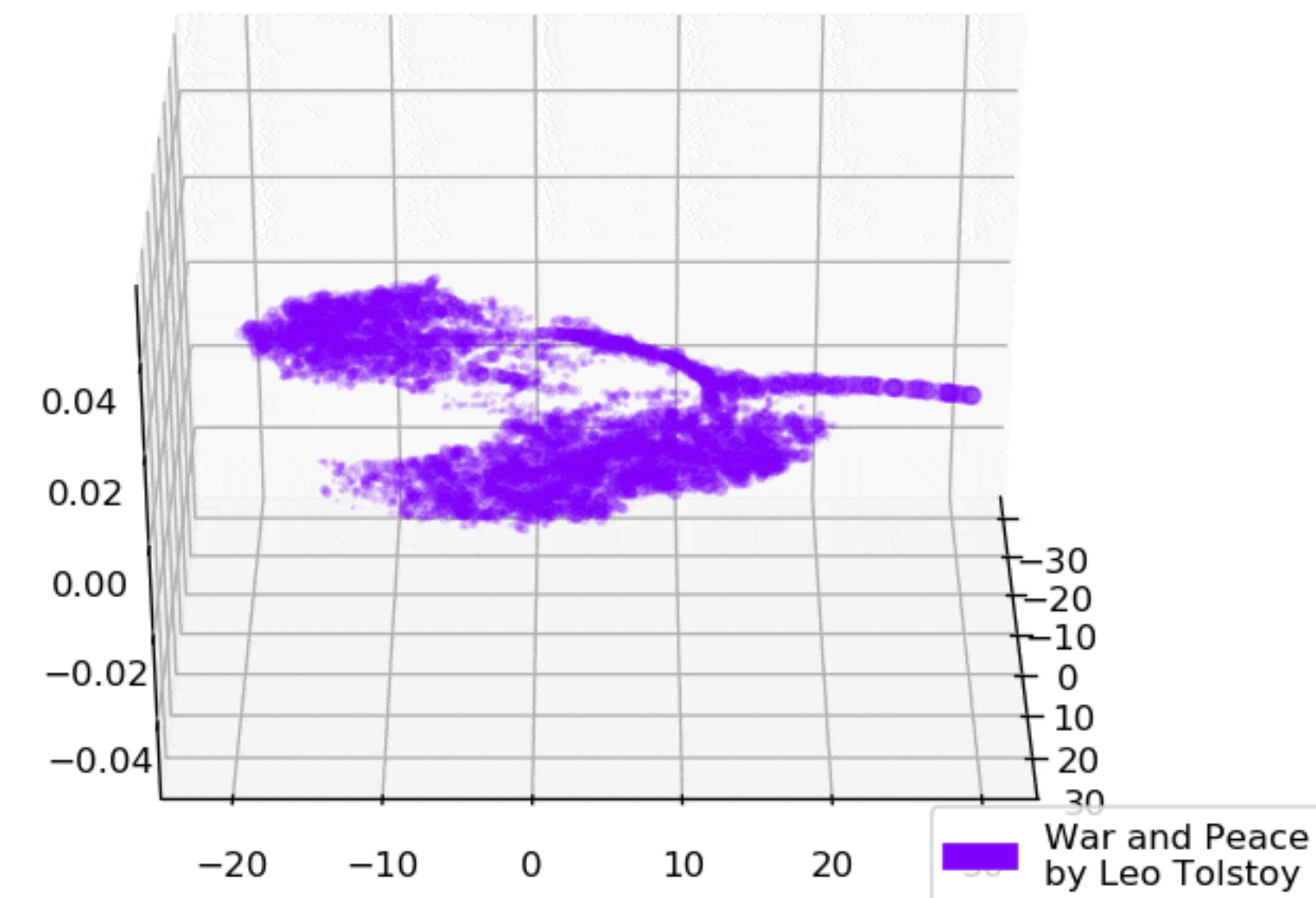


encoder

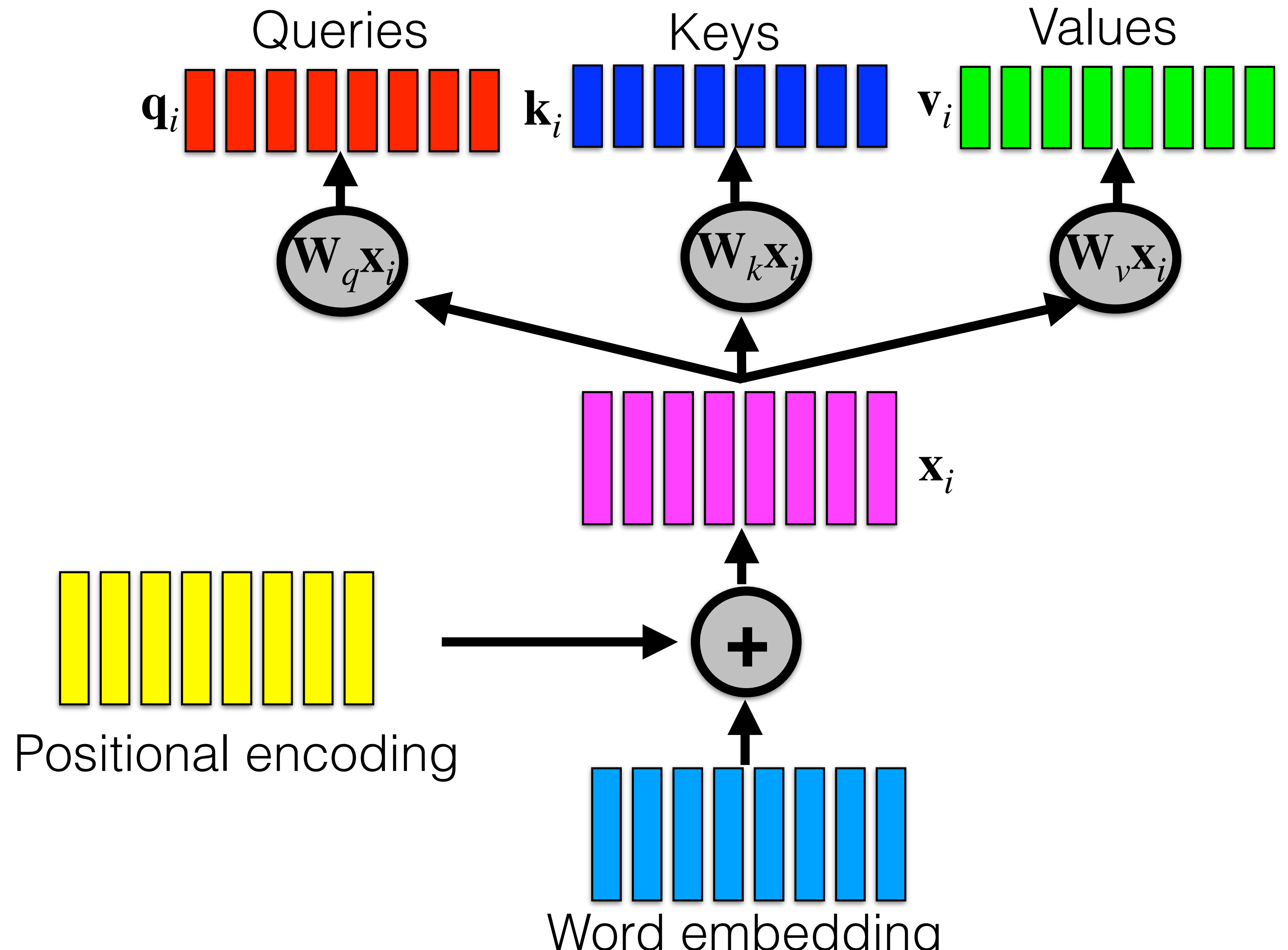


[Mikolov NIPS 2013]

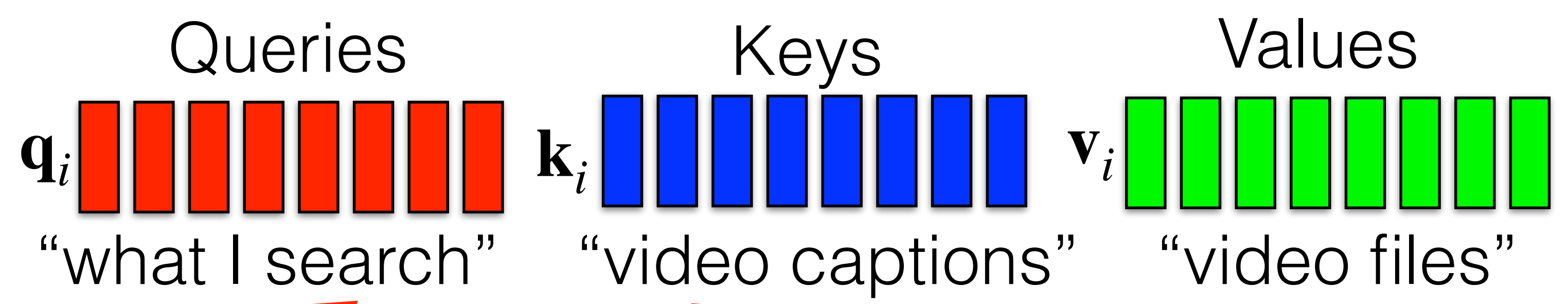
Visualizing Word Embeddings using t-SNE



encoder



# encoder



Scalar product measures similarity between vectors.

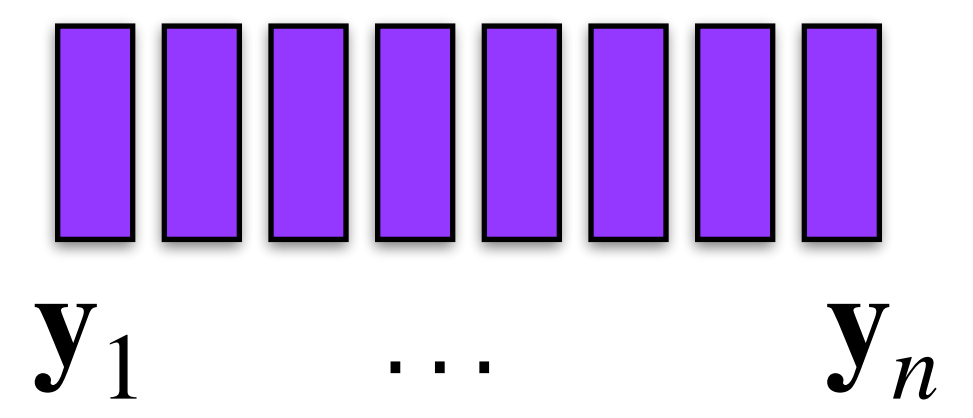
Get attention weights

$$s \left( \begin{matrix} \text{red bar} \\ \mathbf{q}_i^\top \end{matrix} \times \begin{matrix} \text{blue bars} \\ \mathbf{k}_1 \quad \dots \quad \mathbf{k}_n \end{matrix} \right) = \begin{matrix} \text{orange bars} \\ \mathbf{a}_1 \quad \dots \quad \mathbf{a}_n \end{matrix}$$

Attention-weighted sum of values

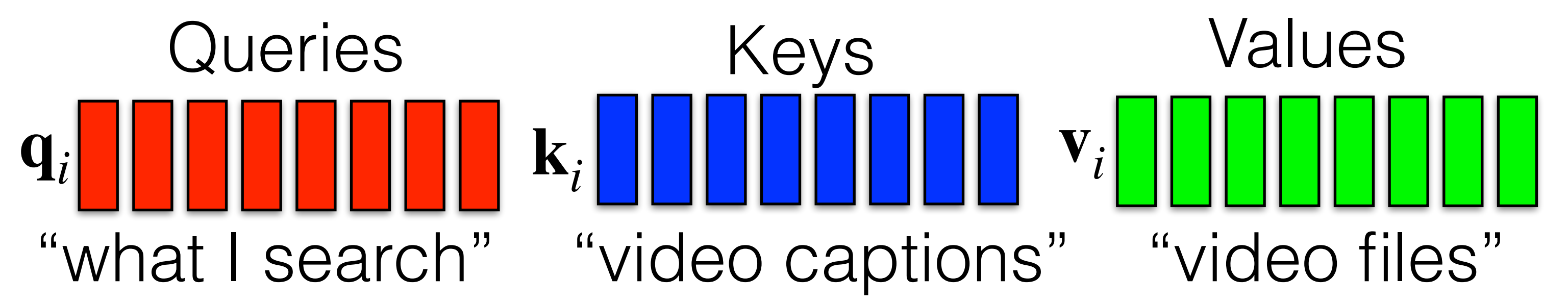
$$\begin{matrix} \text{orange bars} \\ \mathbf{a}_1 \quad \dots \quad \mathbf{a}_n \end{matrix} \begin{matrix} \text{green bars} \\ \mathbf{v}_1 \\ \vdots \\ \mathbf{v}_n \end{matrix} = \text{orange} \cdot \text{green} + \dots = \text{purple bar} \mathbf{y}_1$$

Outputs:

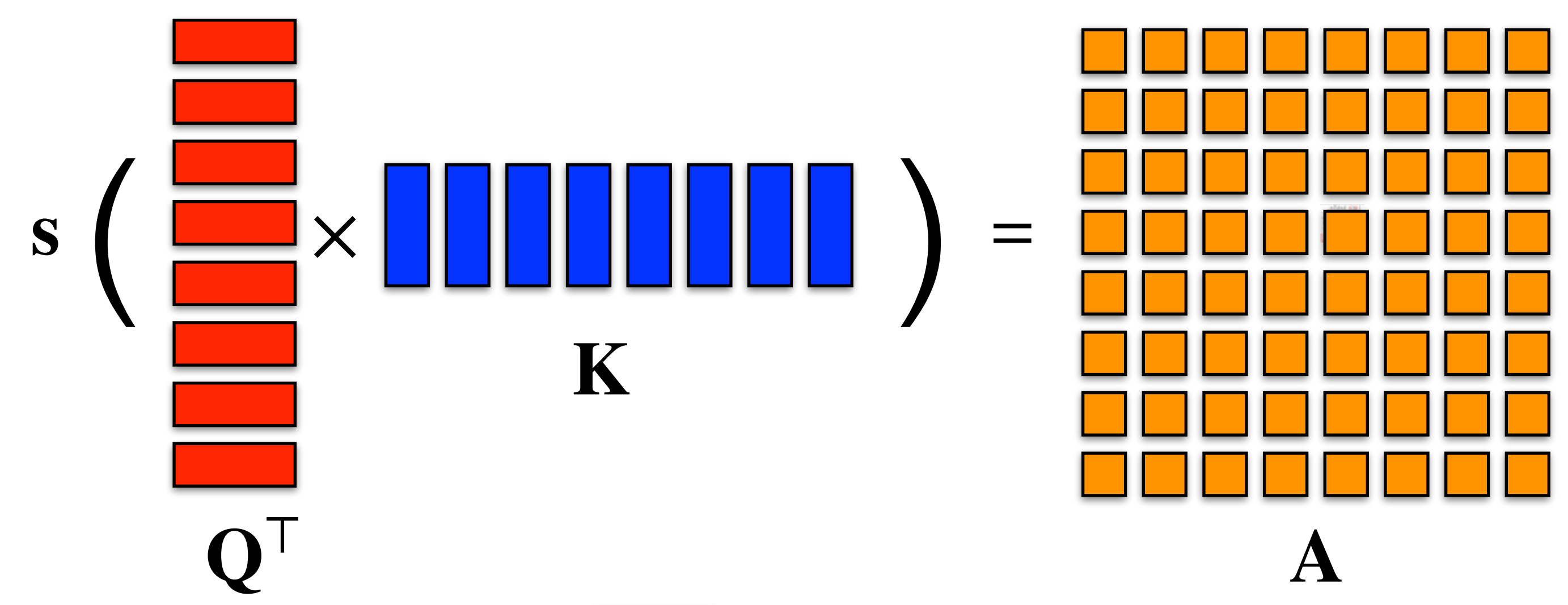


Avoid for-loop by smart matrix multiplication

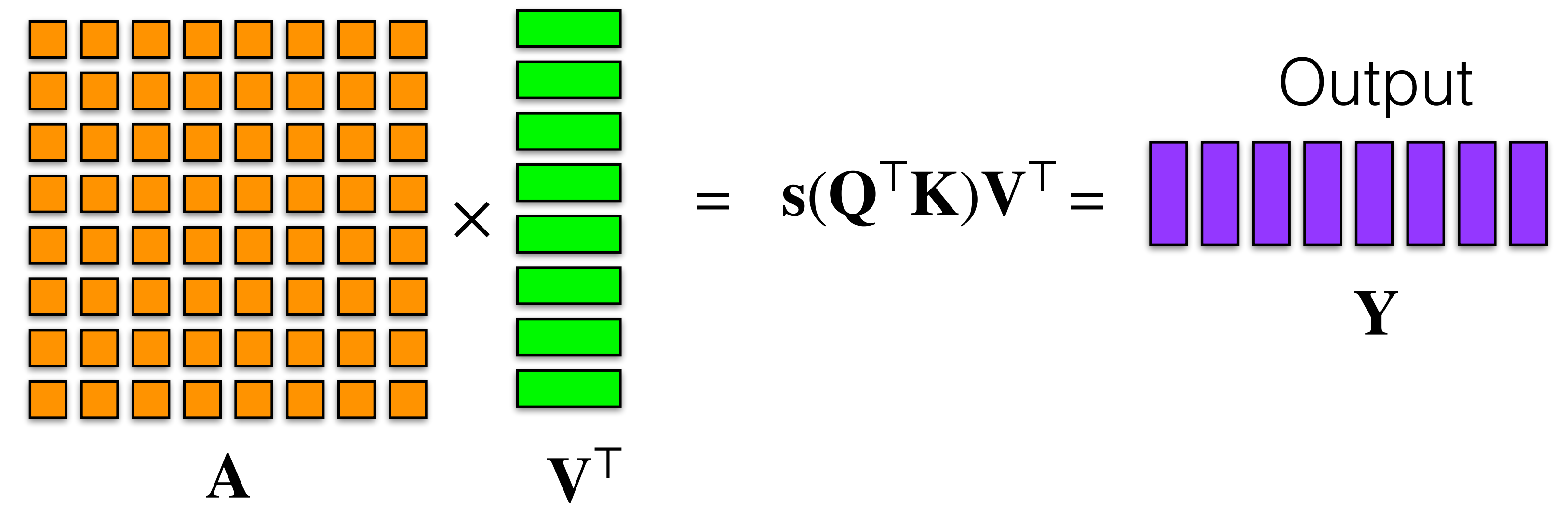
# encoder



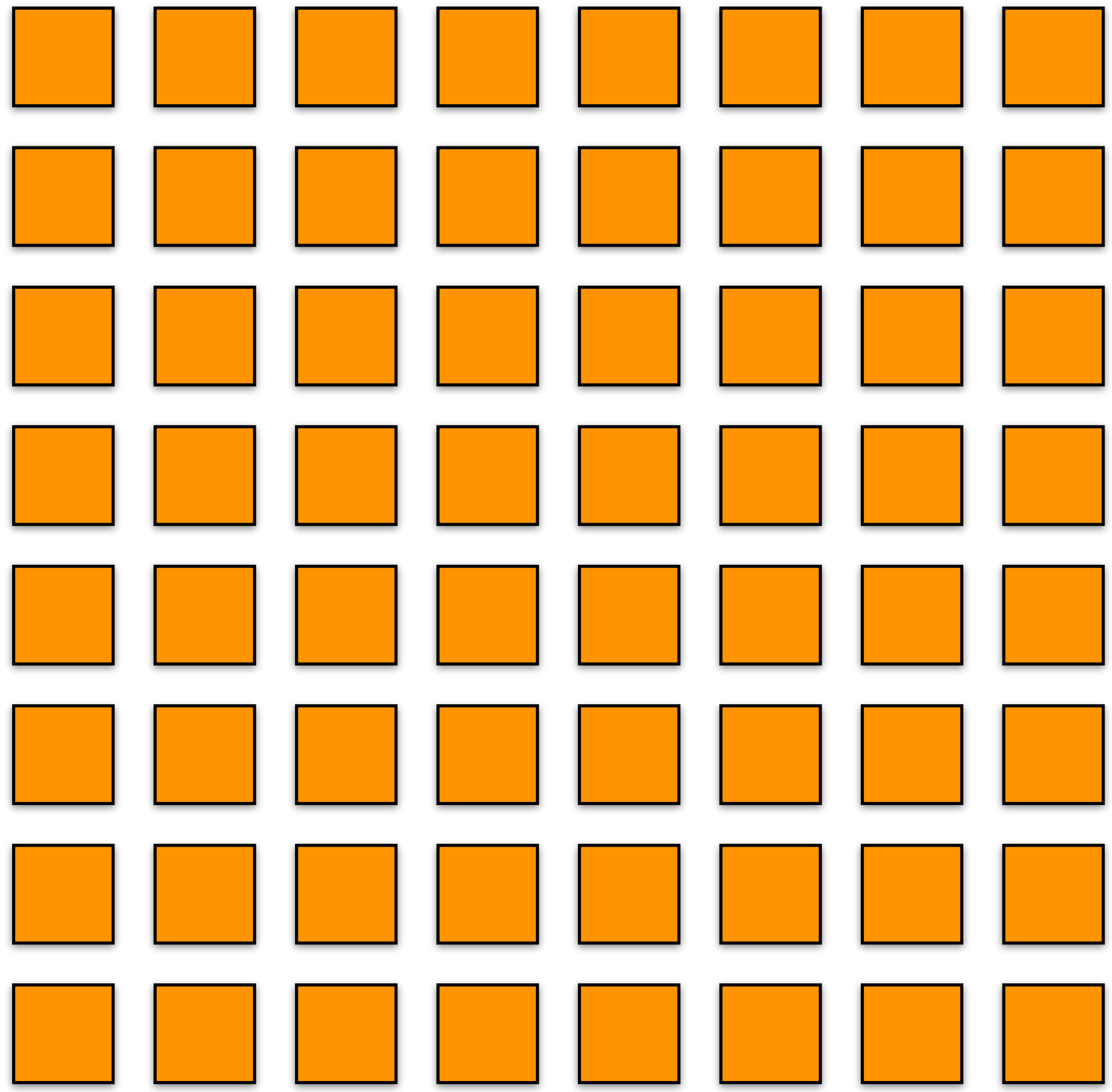
Get attention weights



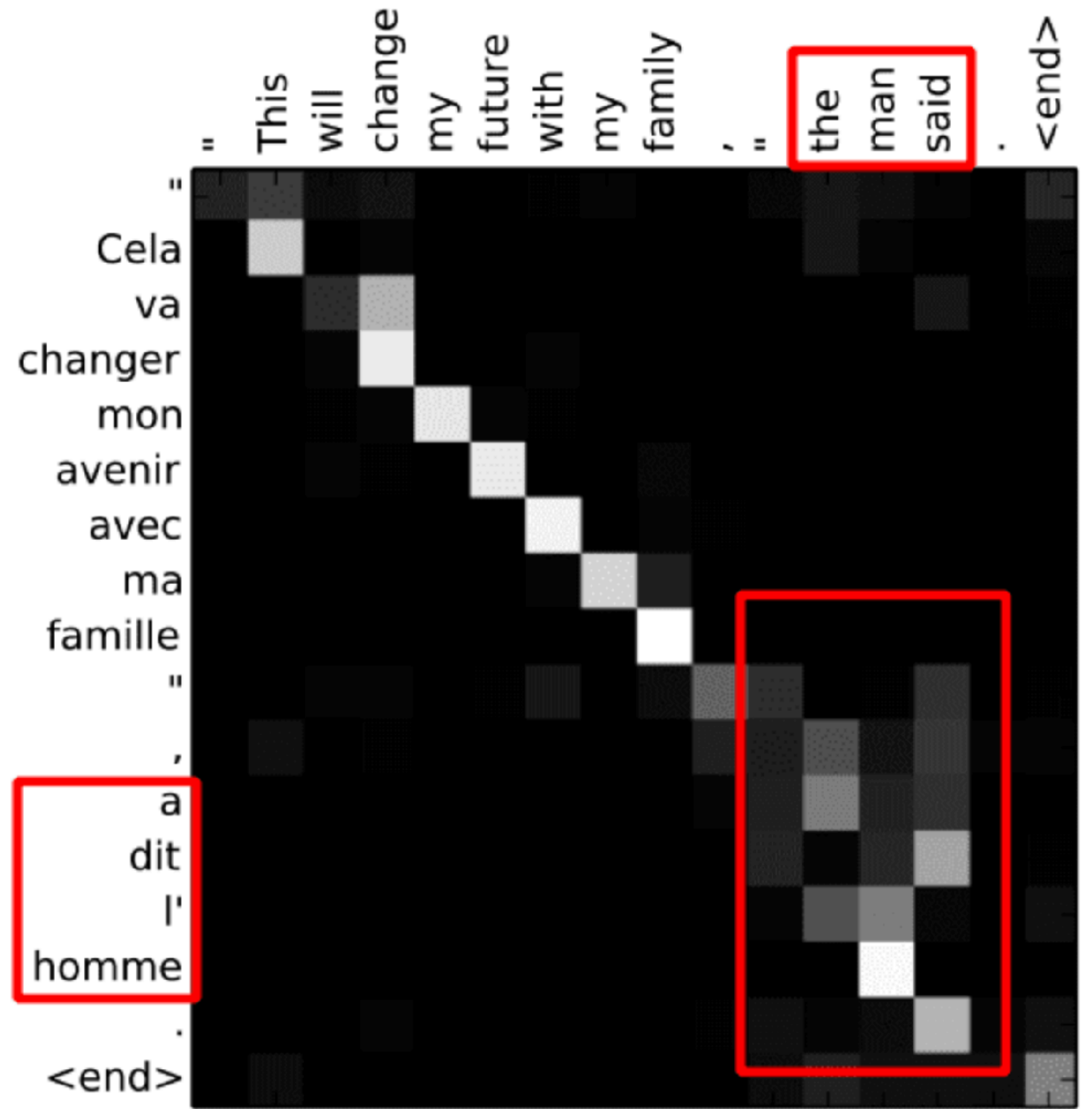
Attention-weighted sum of values



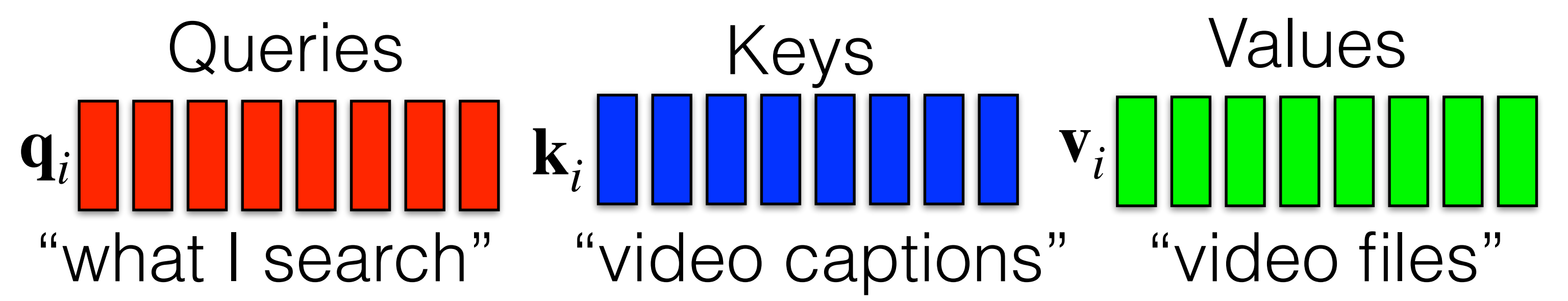
# encoder



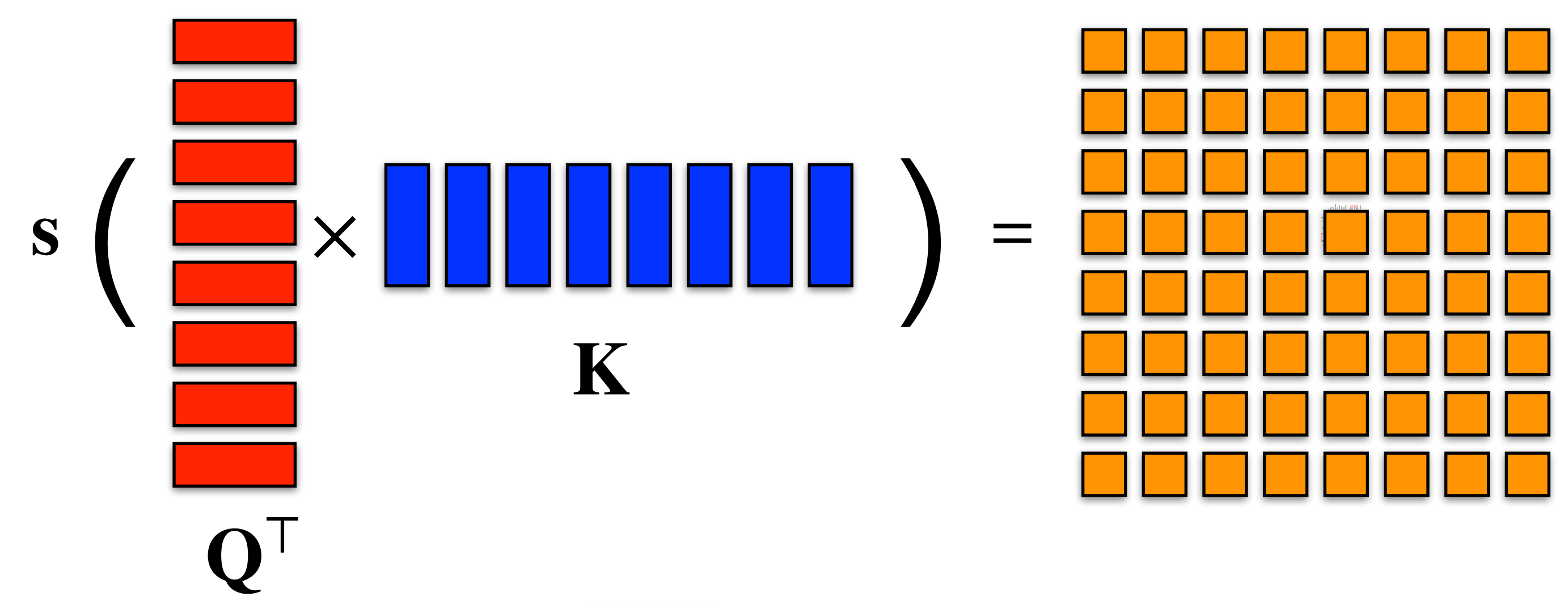
=



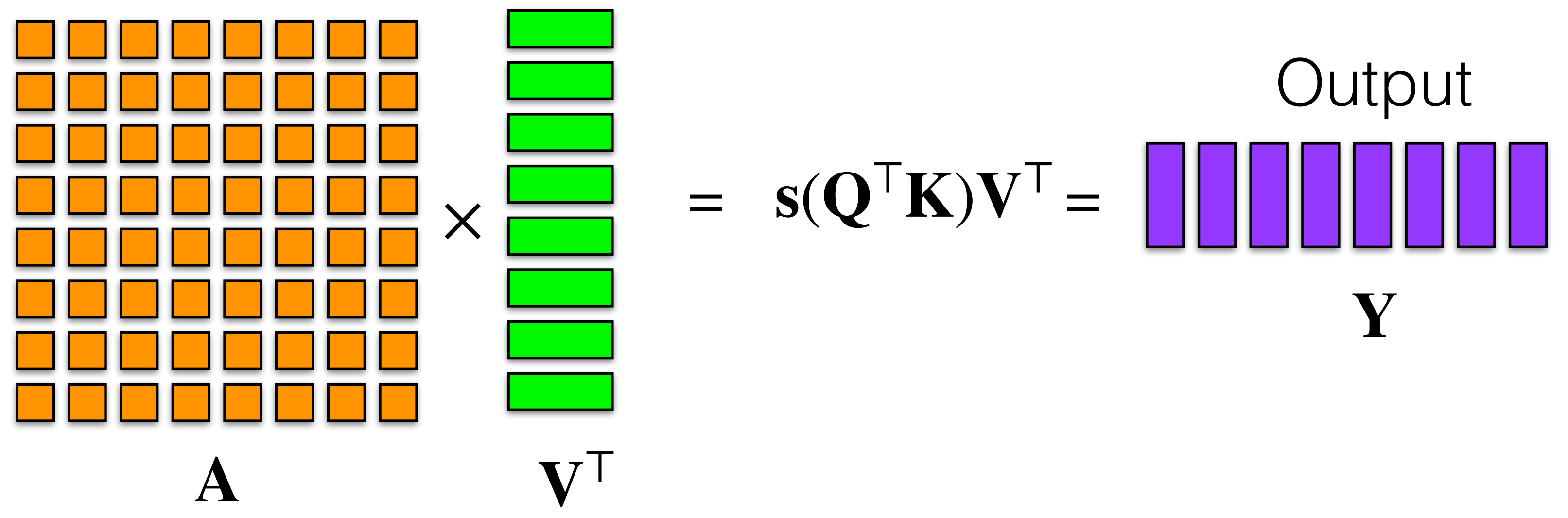
# encoder



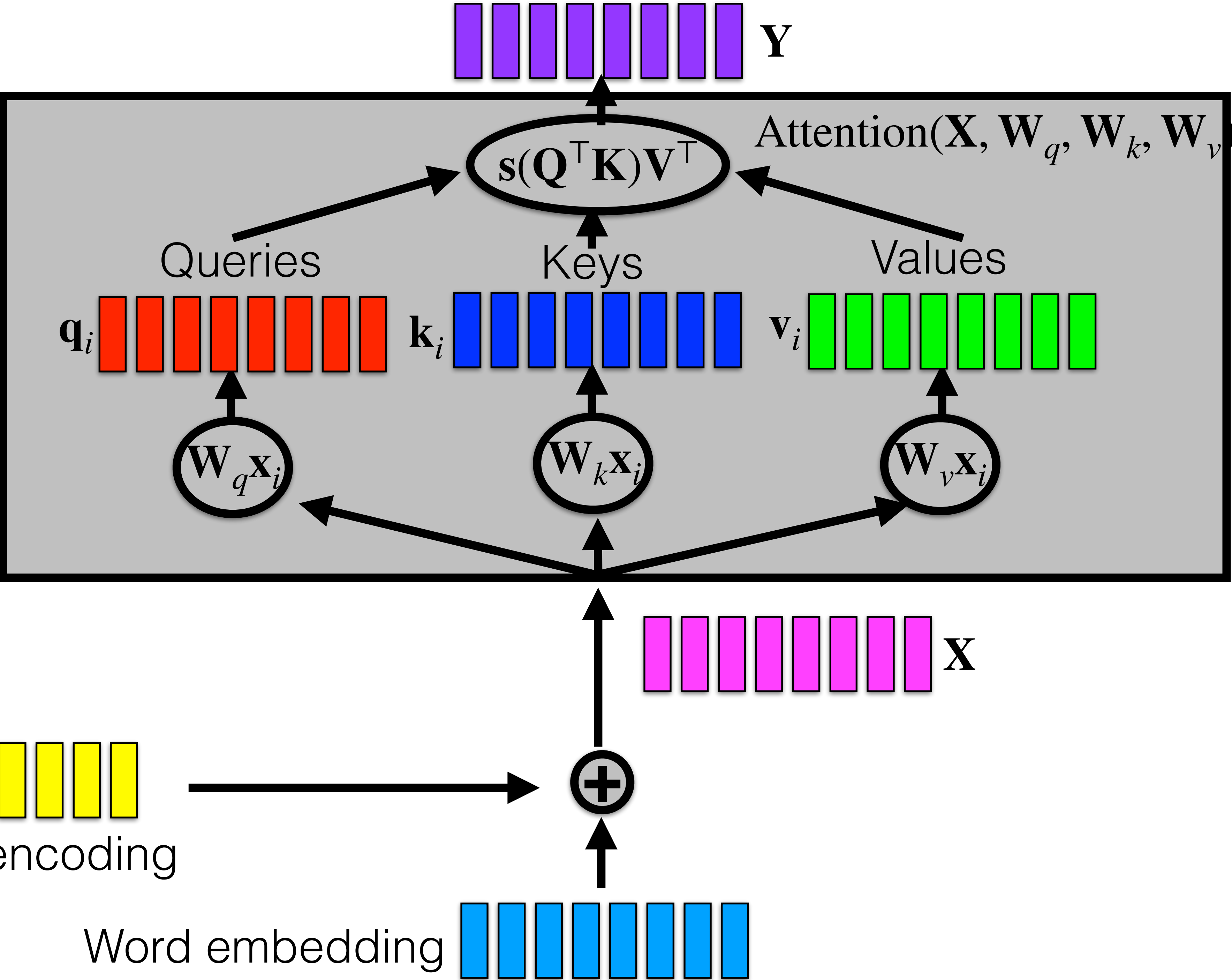
Get attention weights

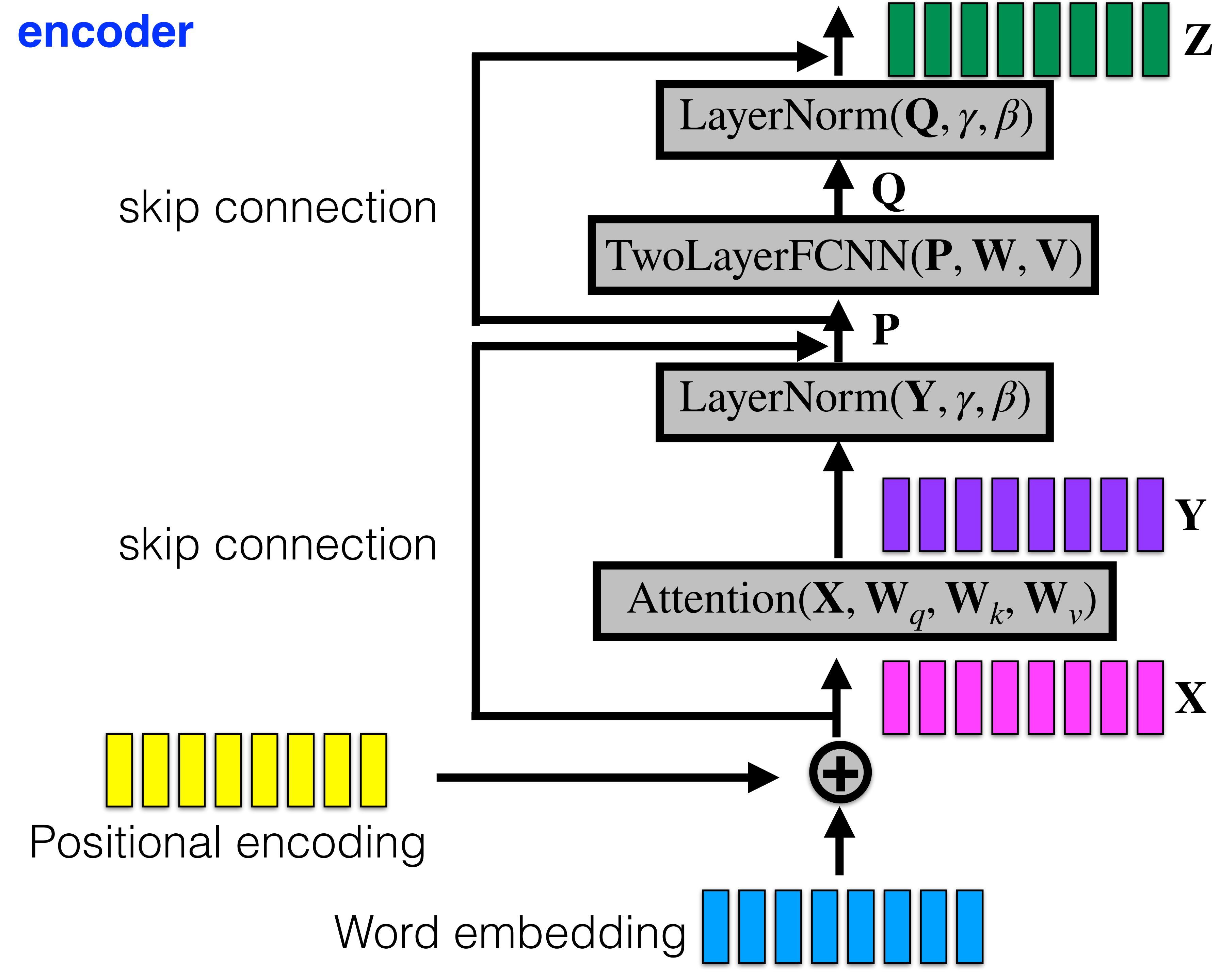


Attention-weighted sum of values

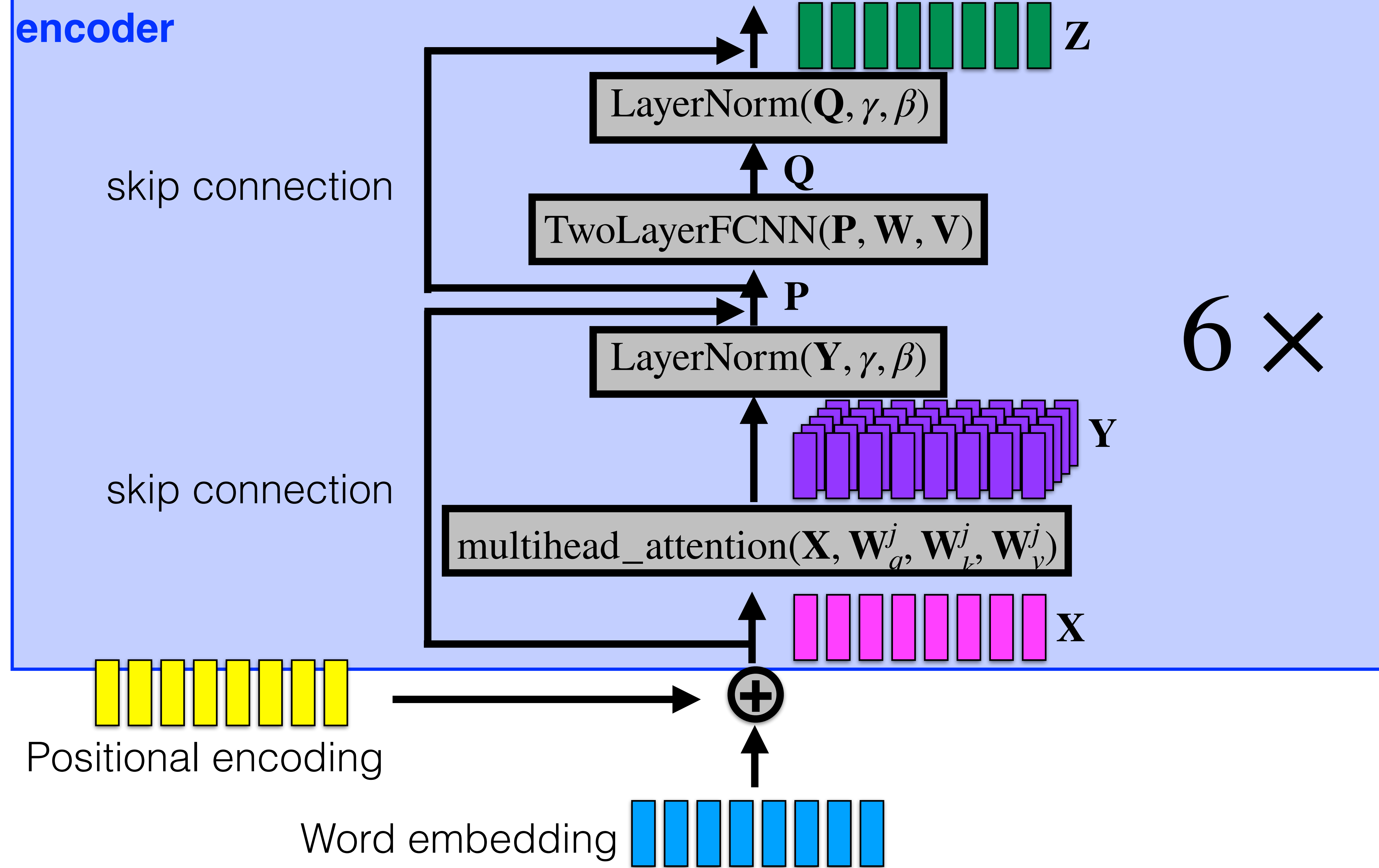


encoder



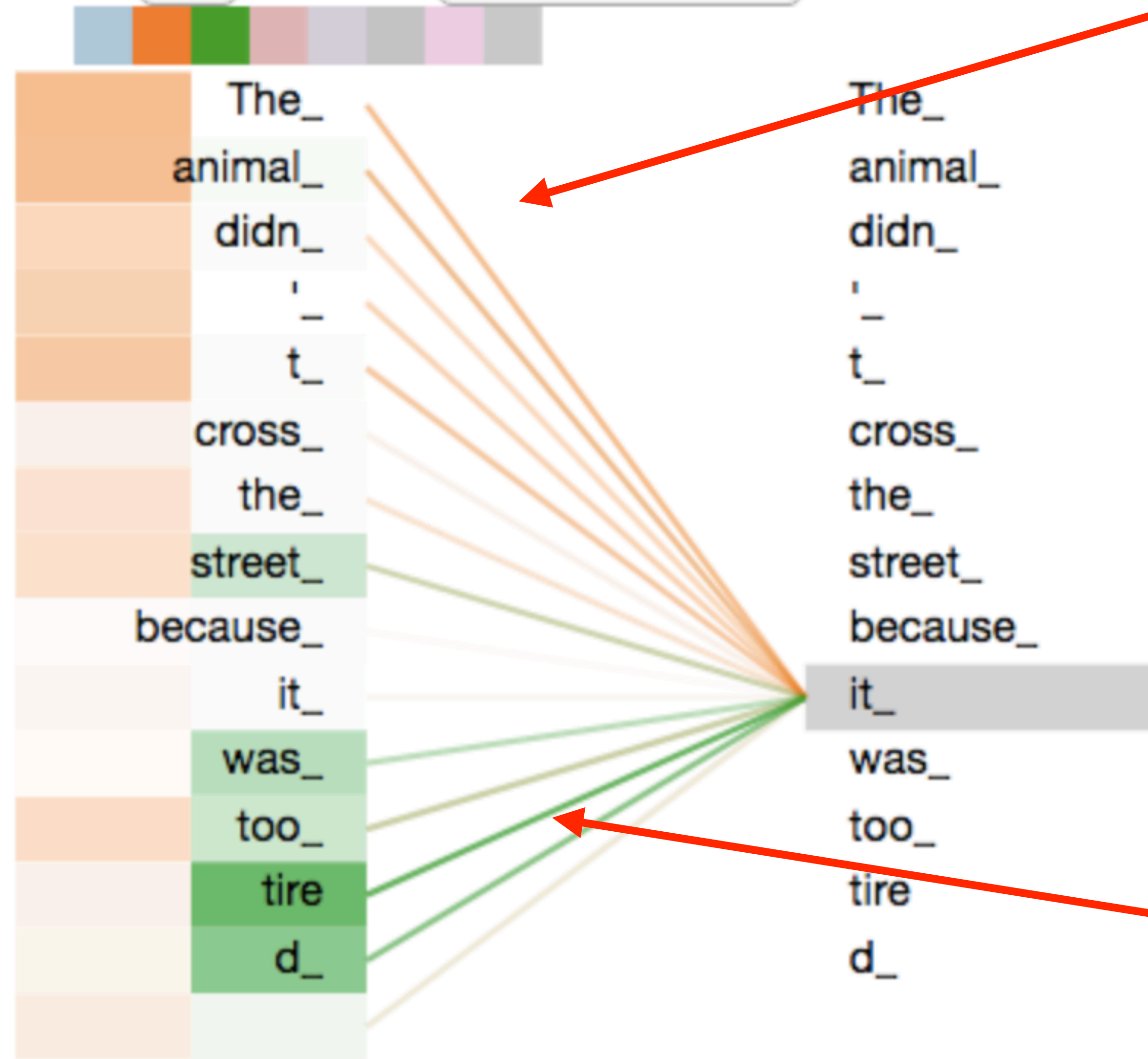






# BertViz

Layer: 5 Attention: Input - Input



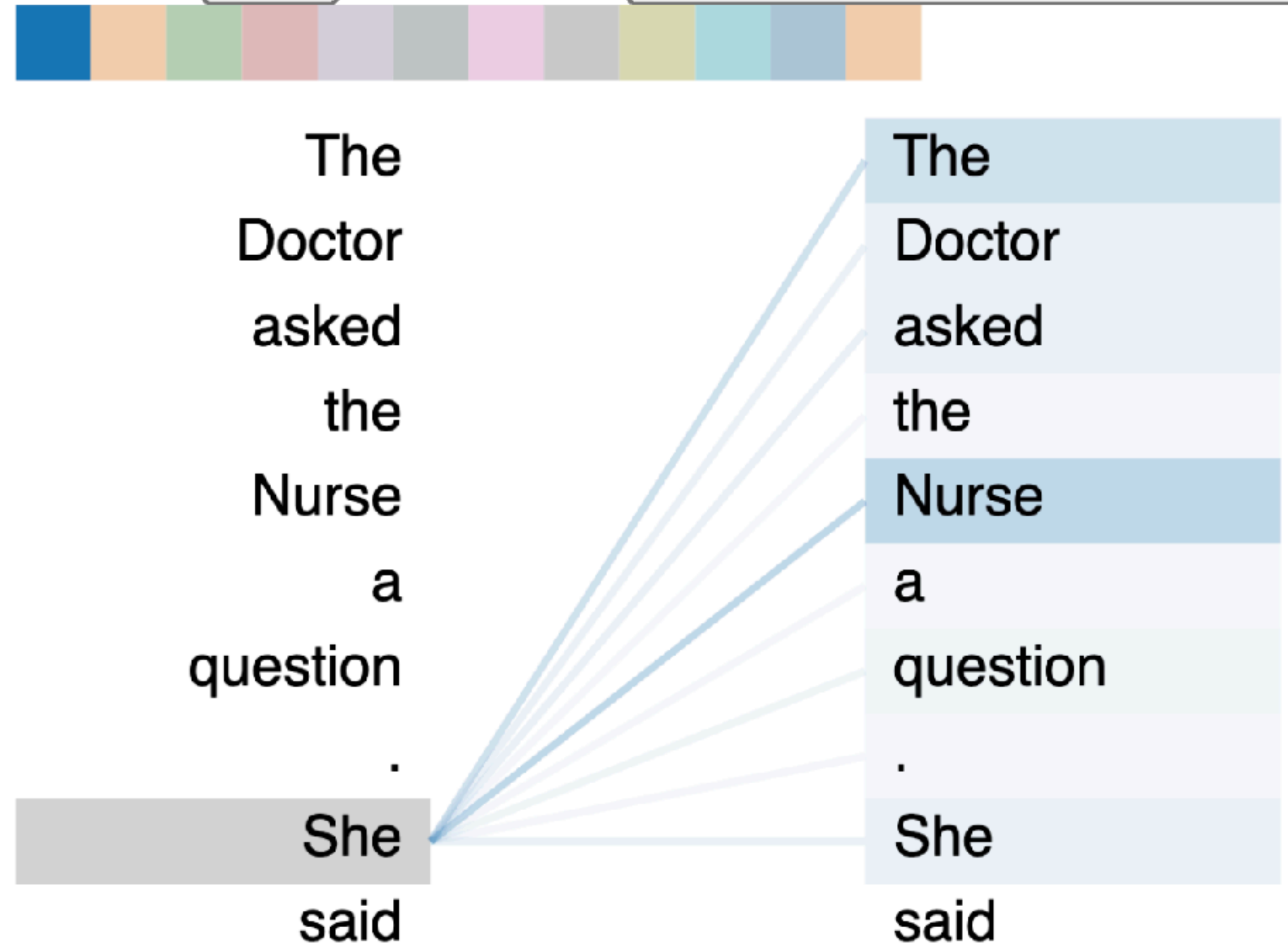
attention weights of **orange** attention head

“it=animal” vs “it=street”??

attention weights of **green** attention head

[https://colab.research.google.com/github/tensorflow/tensor2tensor/blob/master/tensor2tensor/notebooks/hello\\_t2t.ipynb](https://colab.research.google.com/github/tensorflow/tensor2tensor/blob/master/tensor2tensor/notebooks/hello_t2t.ipynb)

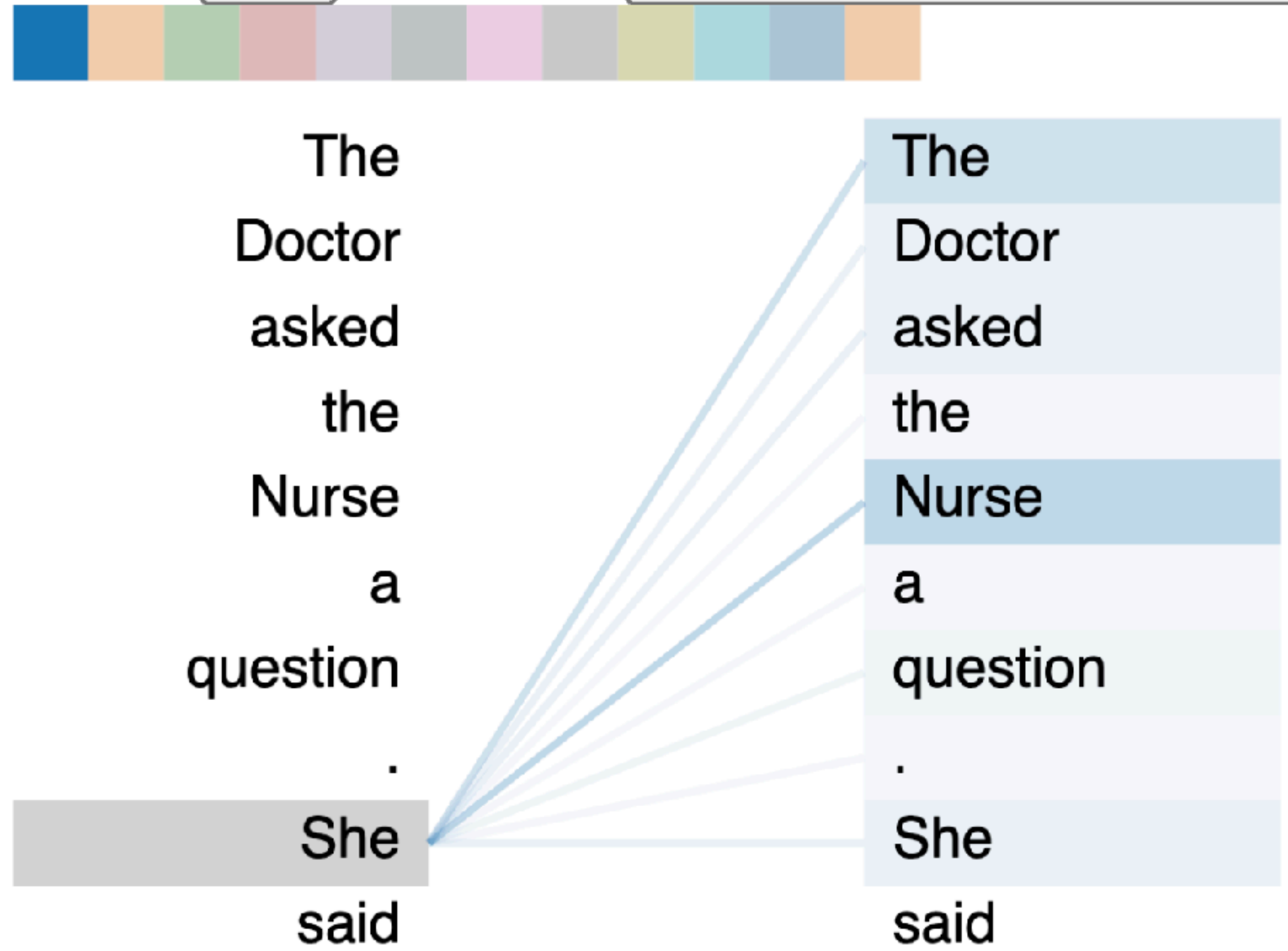
Layer: 0 Attention: All



Model assumes "she=nurse"

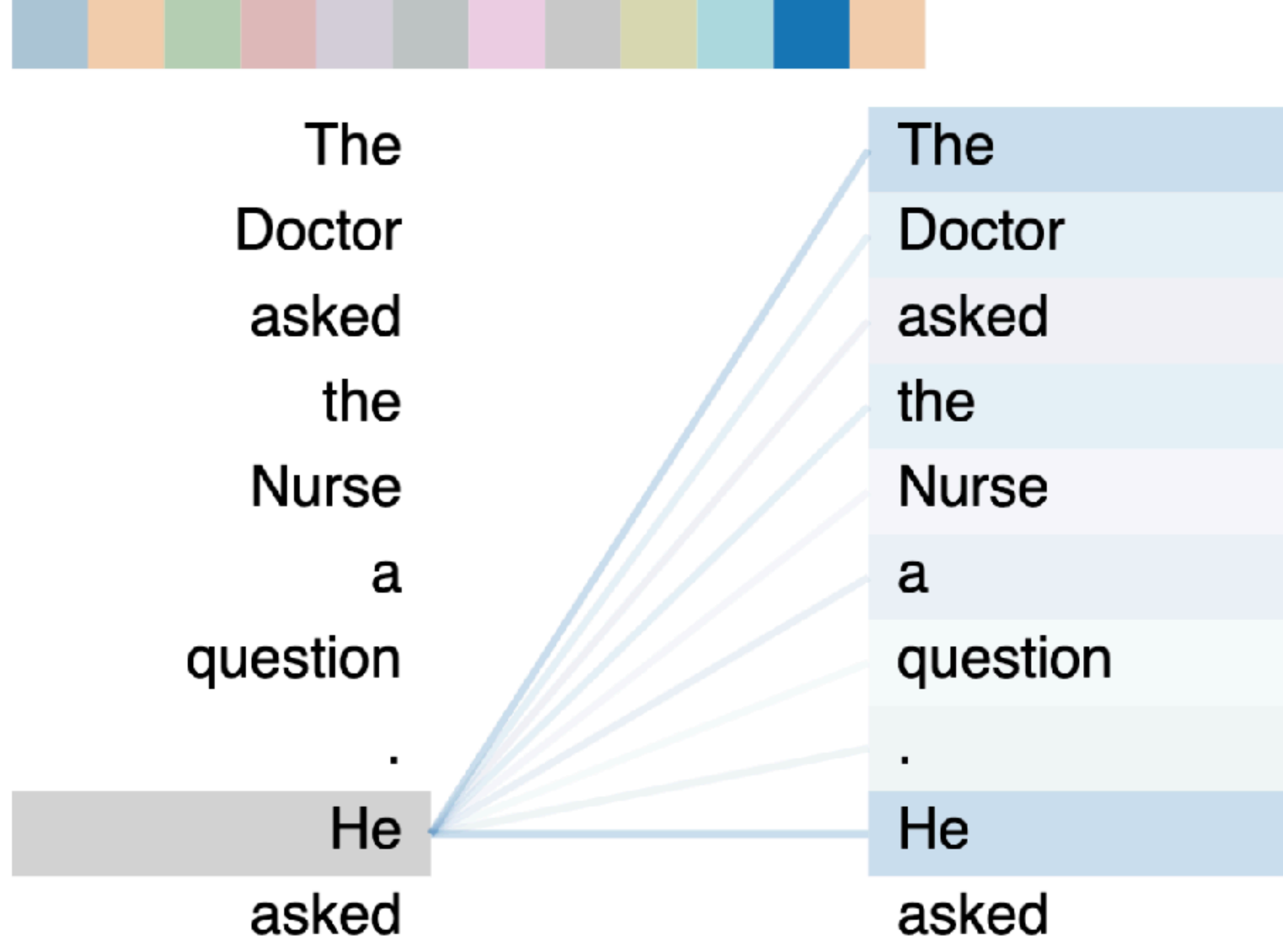
# BertViz (GPT2 model)

Layer: 0 Attention: All



Model assumes "she=nurse"

Layer: 0 Attention: All



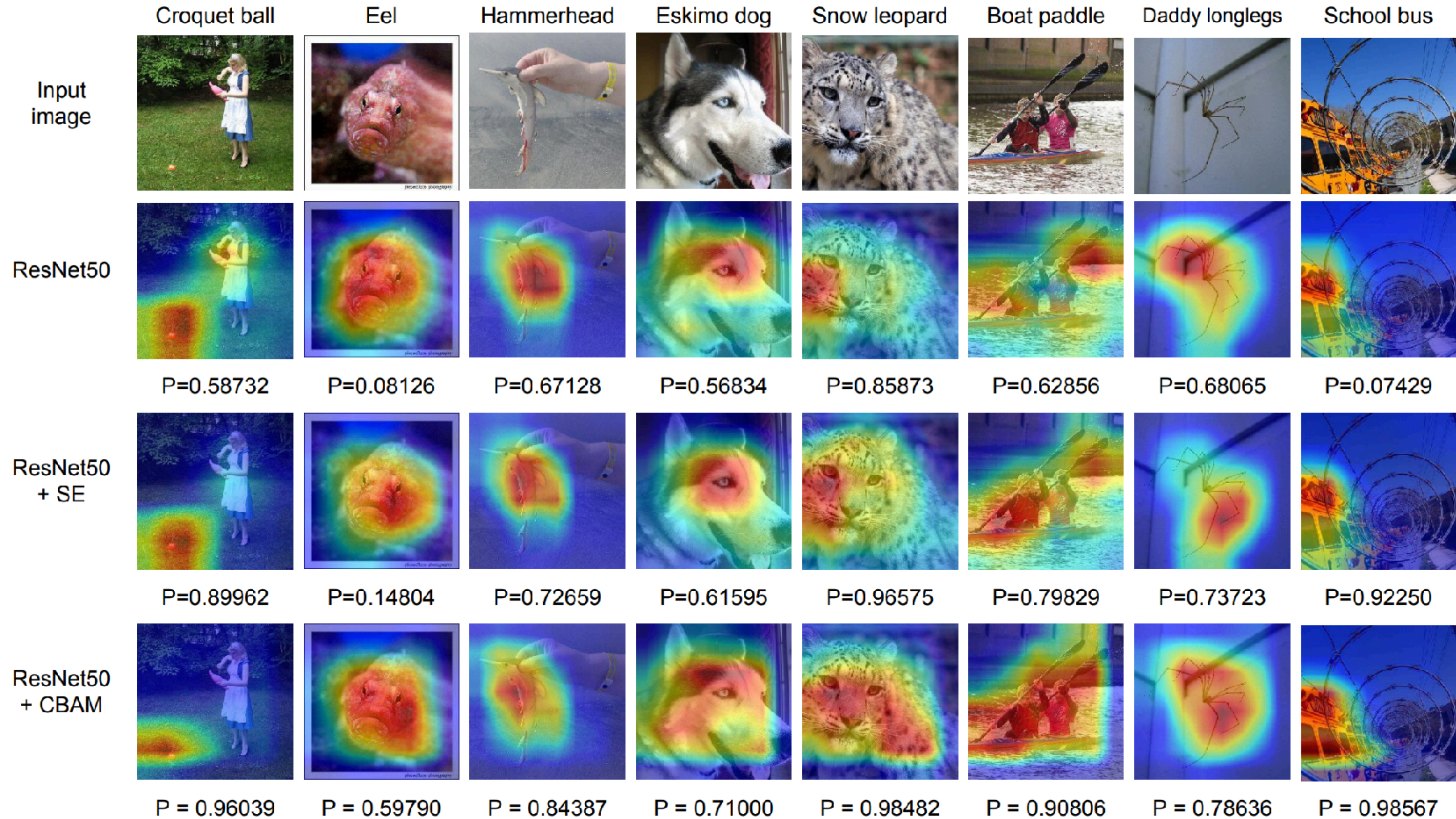
Model assumes "he=doctor"

# Transformers and Attention in images

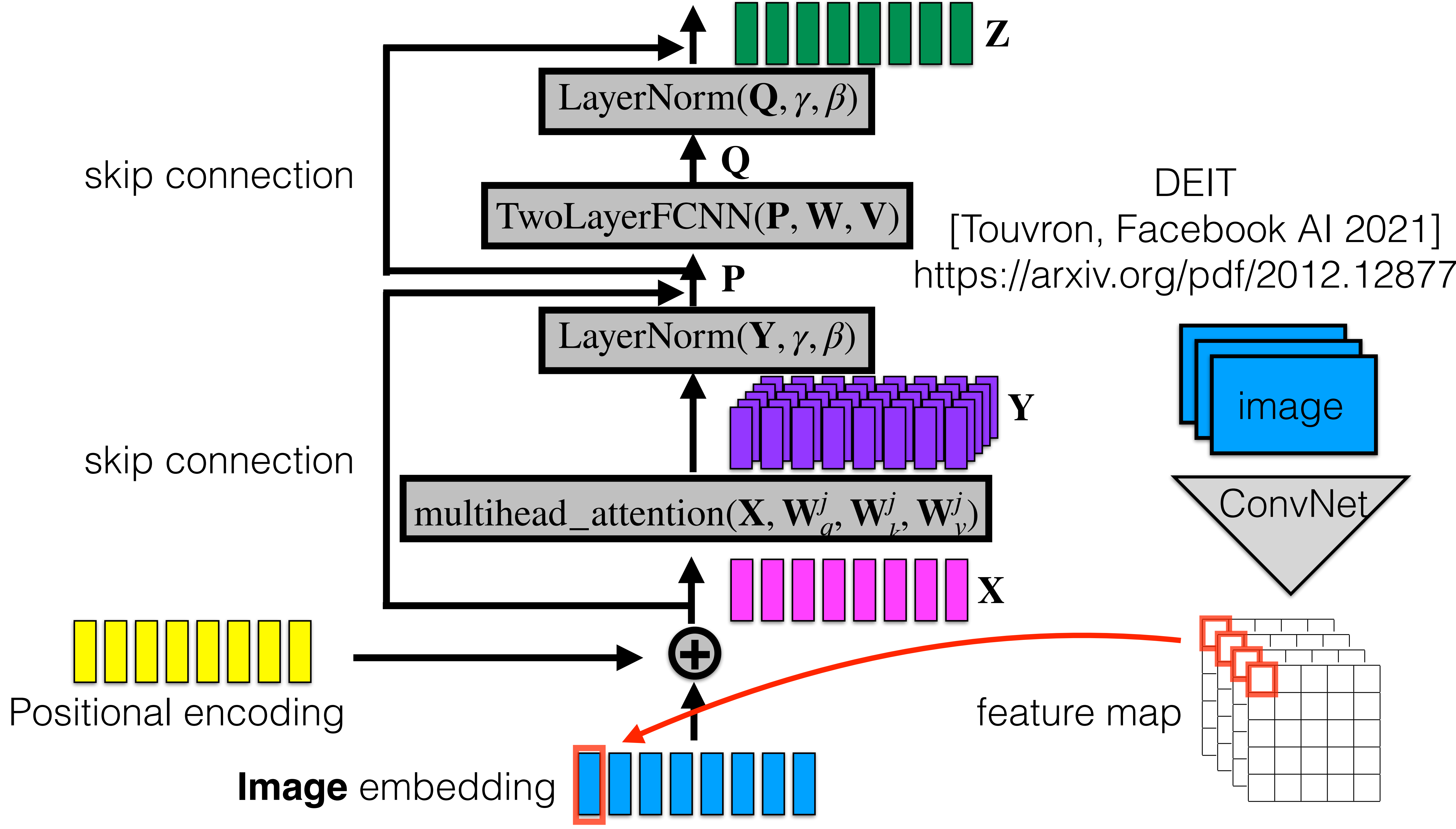


# Attention modules [Woo et al, ECCV, 2018]

<https://arxiv.org/pdf/1807.06521v2.pdf>

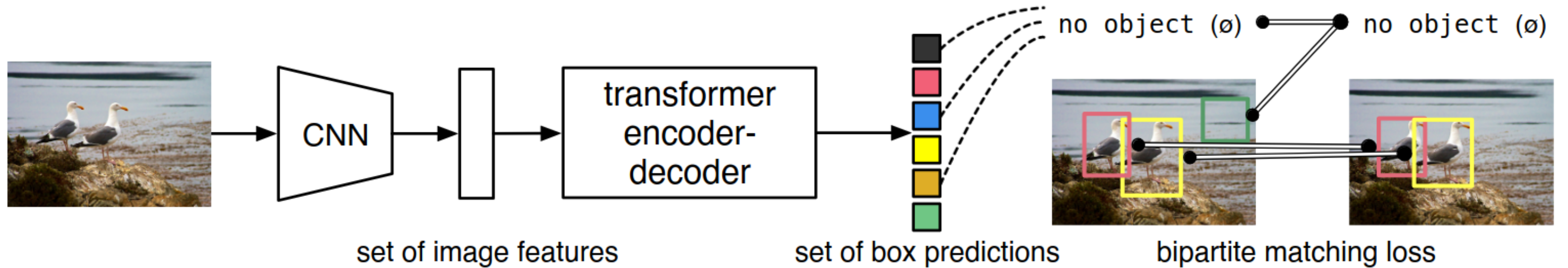






# DETR: transformers for object detection

<https://arxiv.org/abs/2005.12872>



- No anchors
- Output set of N bbs (ordering does not matter)
- Matching loss matches each predicted bb with ground truth bb, or enforces  $\emptyset$



# Attention in RL

Value function  $V(x)$



Advantage function  $A(x,u)$



## Summary

- self-attention overfits (requires large dataset) => combining with hard explicit attention may work better
- memory is attention through time [Alex Graves 2020]
- pyTorch library: <https://github.com/The-AI-Summer/self-attention-cv>  
`model = MultiHeadSelfAttention(dim=64)`

# Future?

- Most of predictions were wrong
  - 1954 IBM predicted that natural language processing will be solved in 3 years
  - 1965 Herbert Simon: machines will replace humans in all manual works
  - 1970 Marvin Minsky: machines will have general AI comparable with humans
  - 2014 Ray Kurzweil: the same for for 2029, now talks about 2045
- Rodney Brooks prediction score card:  
<https://rodneybrooks.com/predictions-scorecard-2021-january-01/>
- False generalization
  - AI is better in solving particular instances (image processing, stabilization)
  - Rather carefully isolated successes than exponentially growing general AI