

Attention and Memory

Karel Zimmermann

<http://cmp.felk.cvut.cz/~zimmerk/>



Vision for Robotics and Autonomous Systems

<https://cyber.felk.cvut.cz/vras/>



Center for Machine Perception

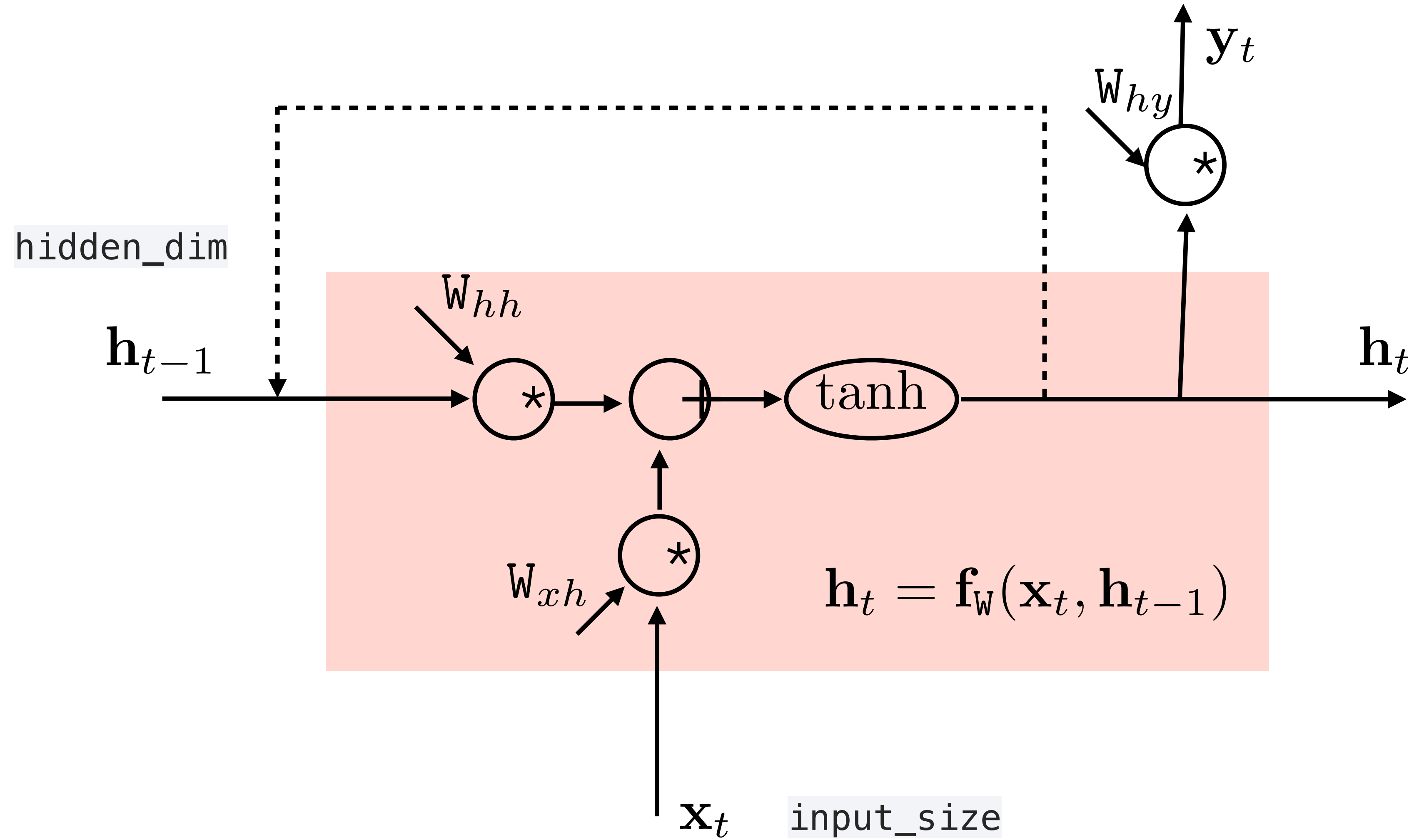
<https://cmp.felk.cvut.cz>



Department for Cybernetics
Faculty of Electrical Engineering
Czech Technical University in Prague

Simple recurrent block

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



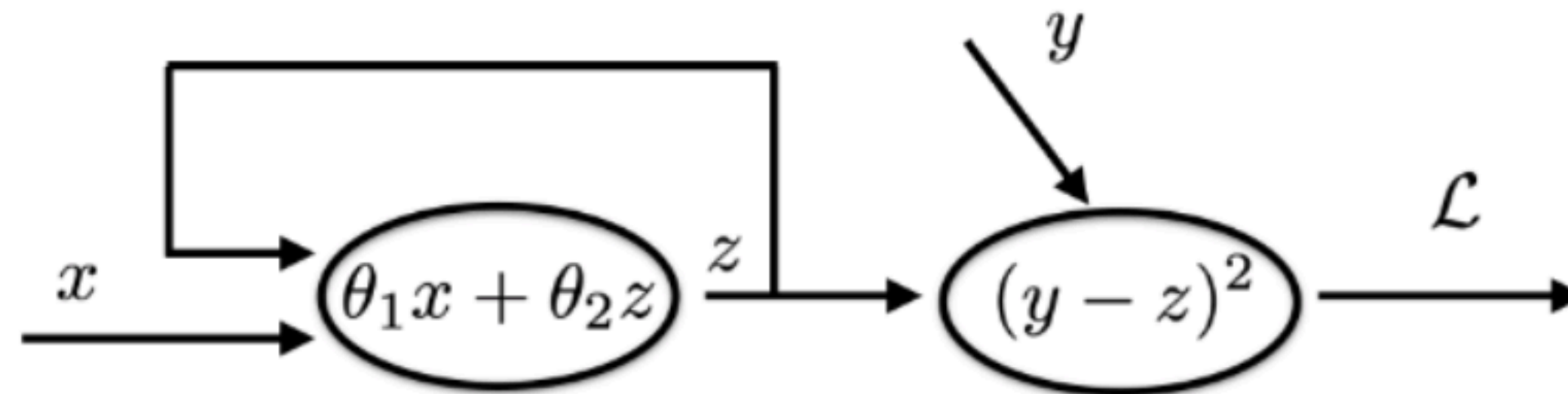
PyTorch: <https://pytorch.org/docs/stable/nn.html>

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:

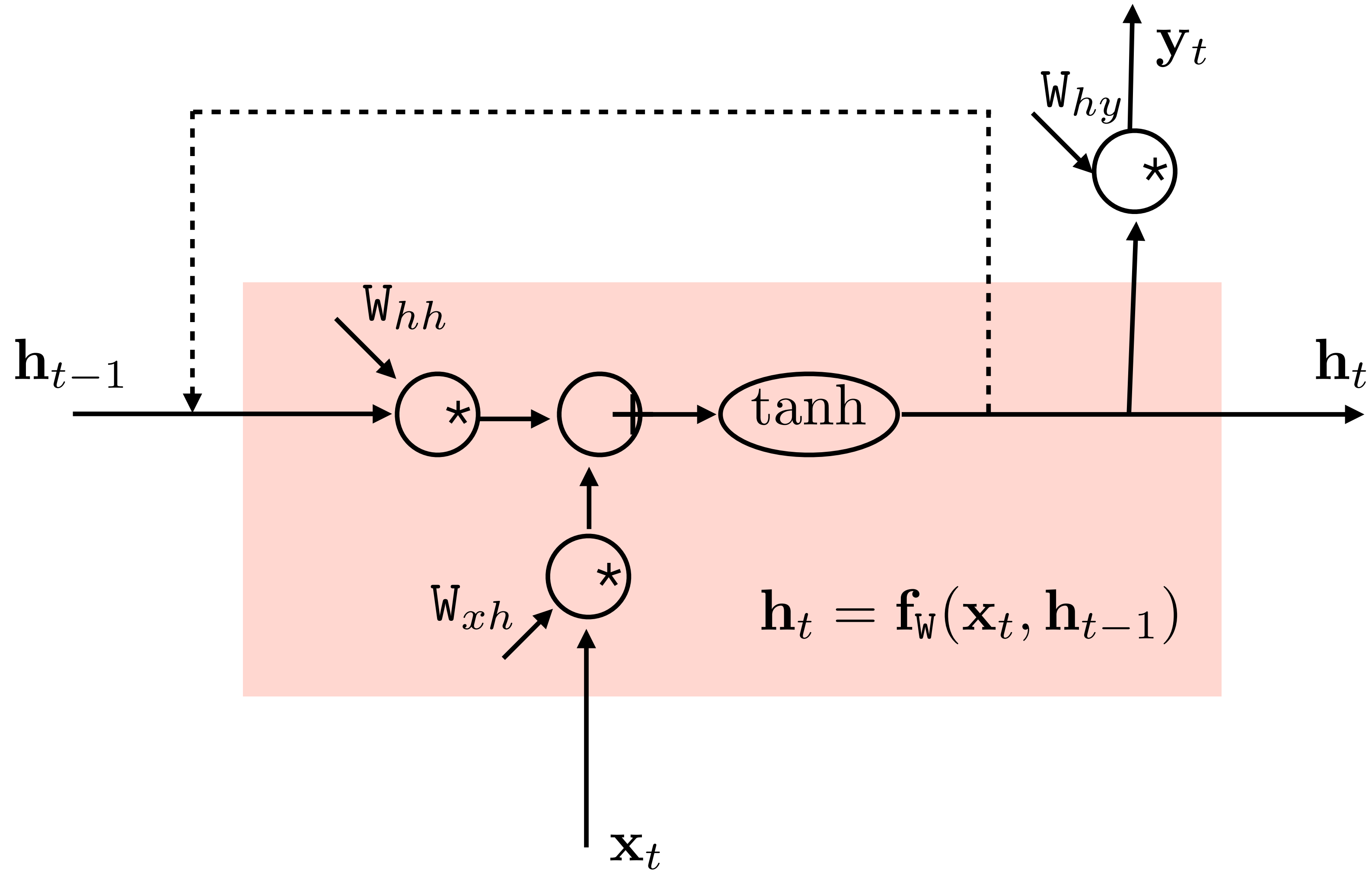
| time=1 | time=2 |
|-----------|-----------|
| $x_1 = 0$ | $x_2 = 1$ |
| $y_1 = 1$ | $y_2 = 3$ |

Estimate gradient of the overall loss (computed over all available outputs y_i for both available times $i = 1, 2$) with respect to θ_1 .

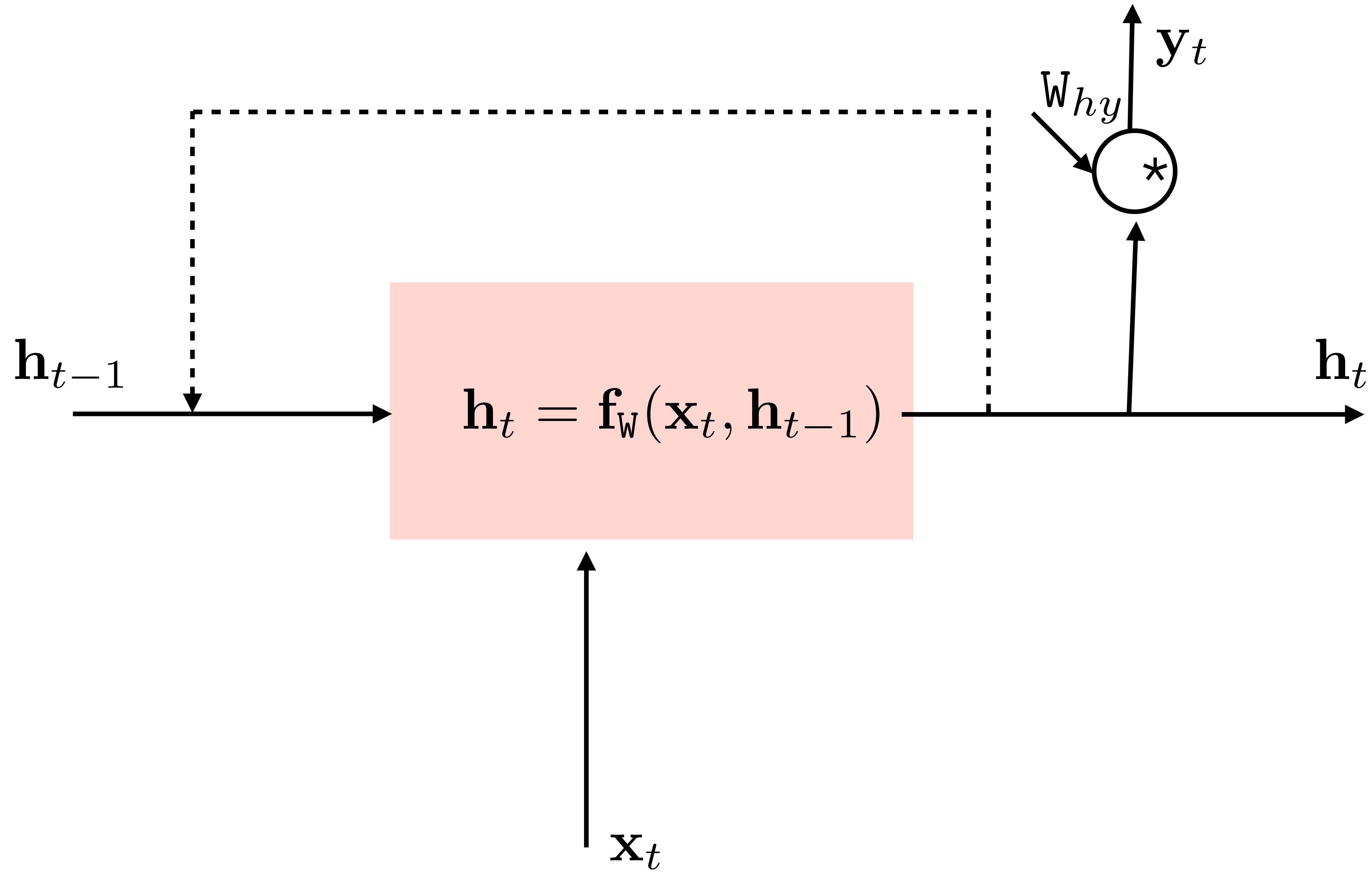


Hint: Unroll the network in time, to obtain a usual feedforward network with two loss nodes. Do the backpropagation as usual.

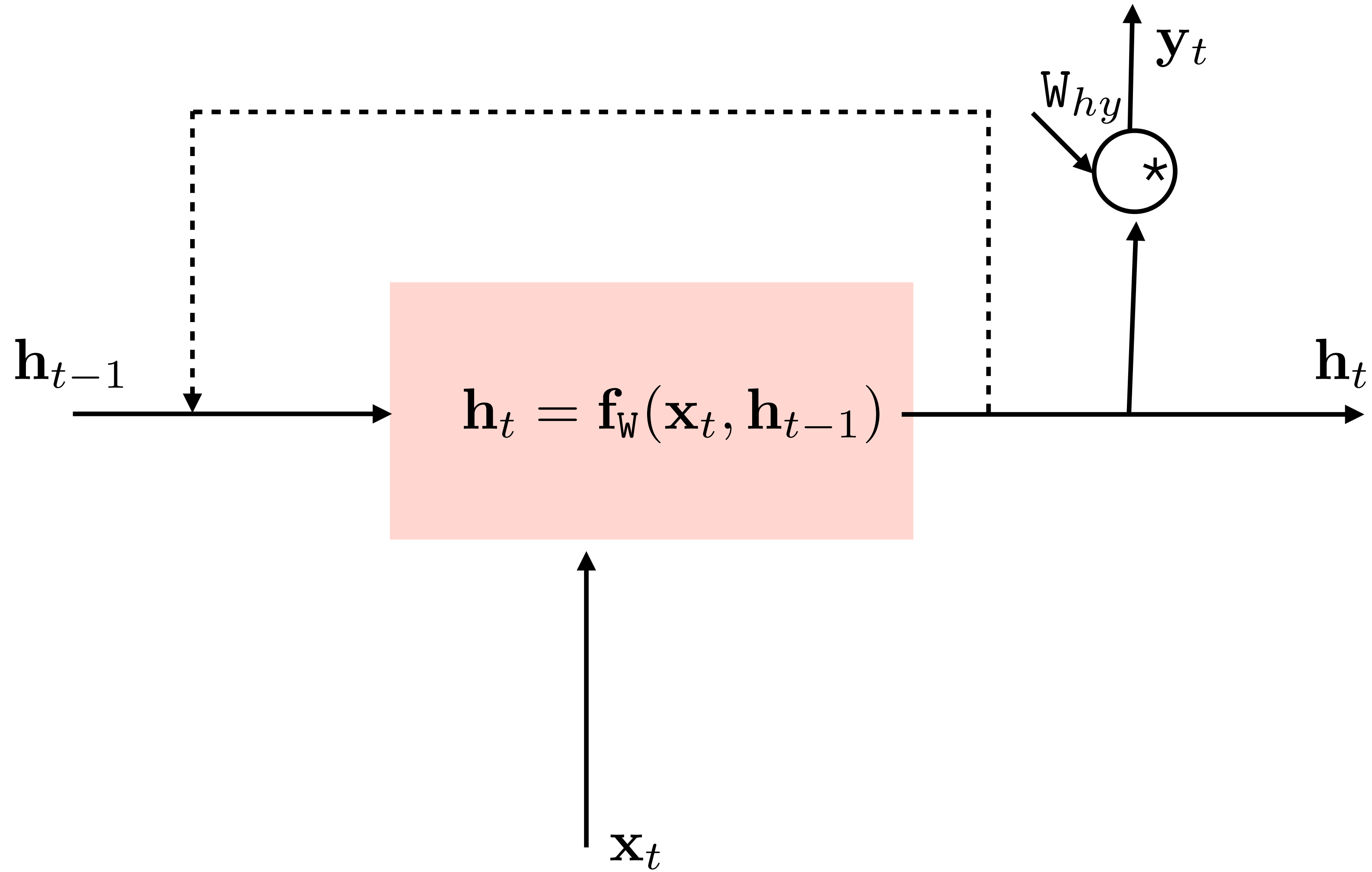
Simple recurrent block - feed-forward pass



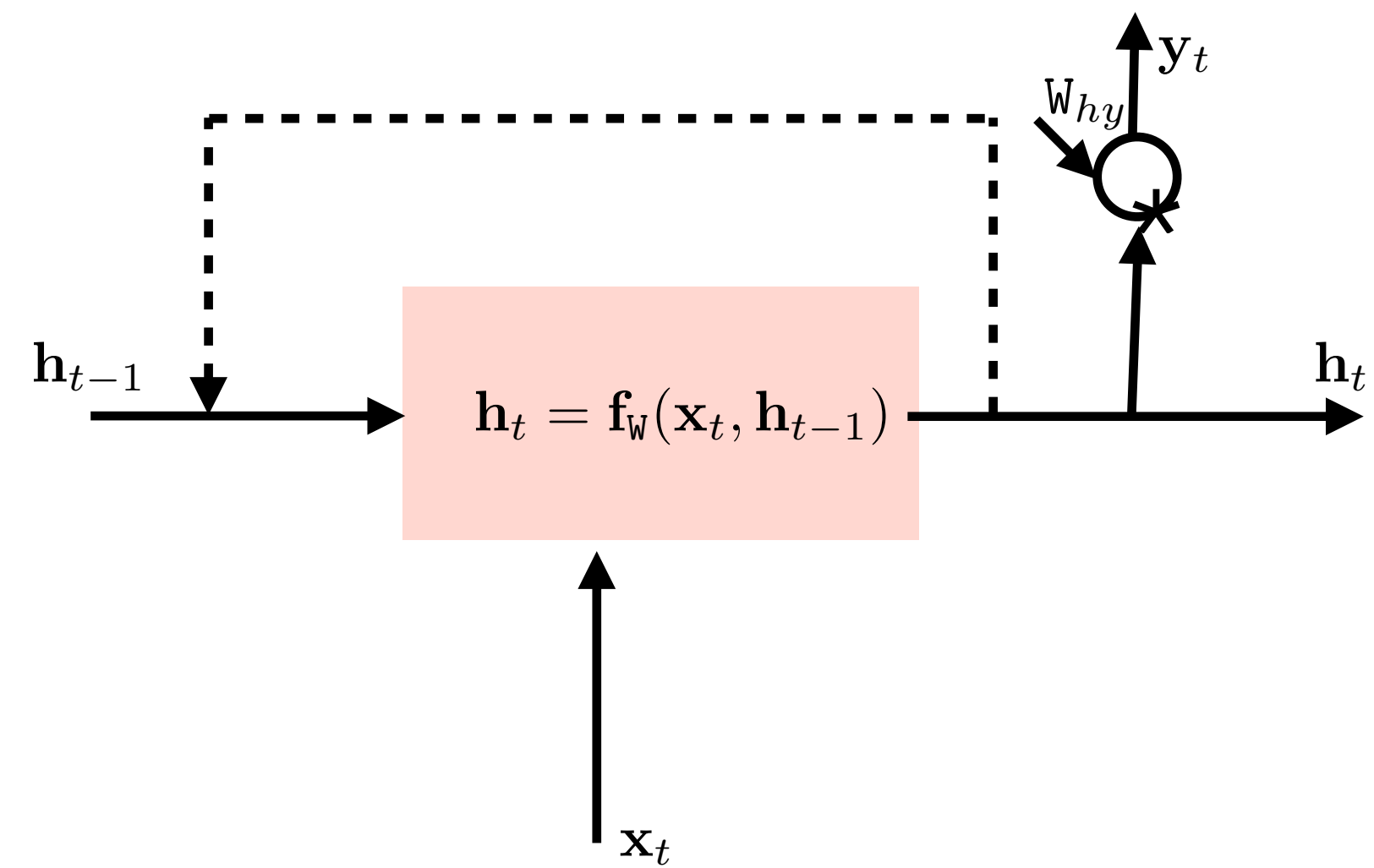
Simple recurrent block - feed-forward pass



Simple recurrent block - feed-forward pass



Simple recurrent block - feed-forward pass

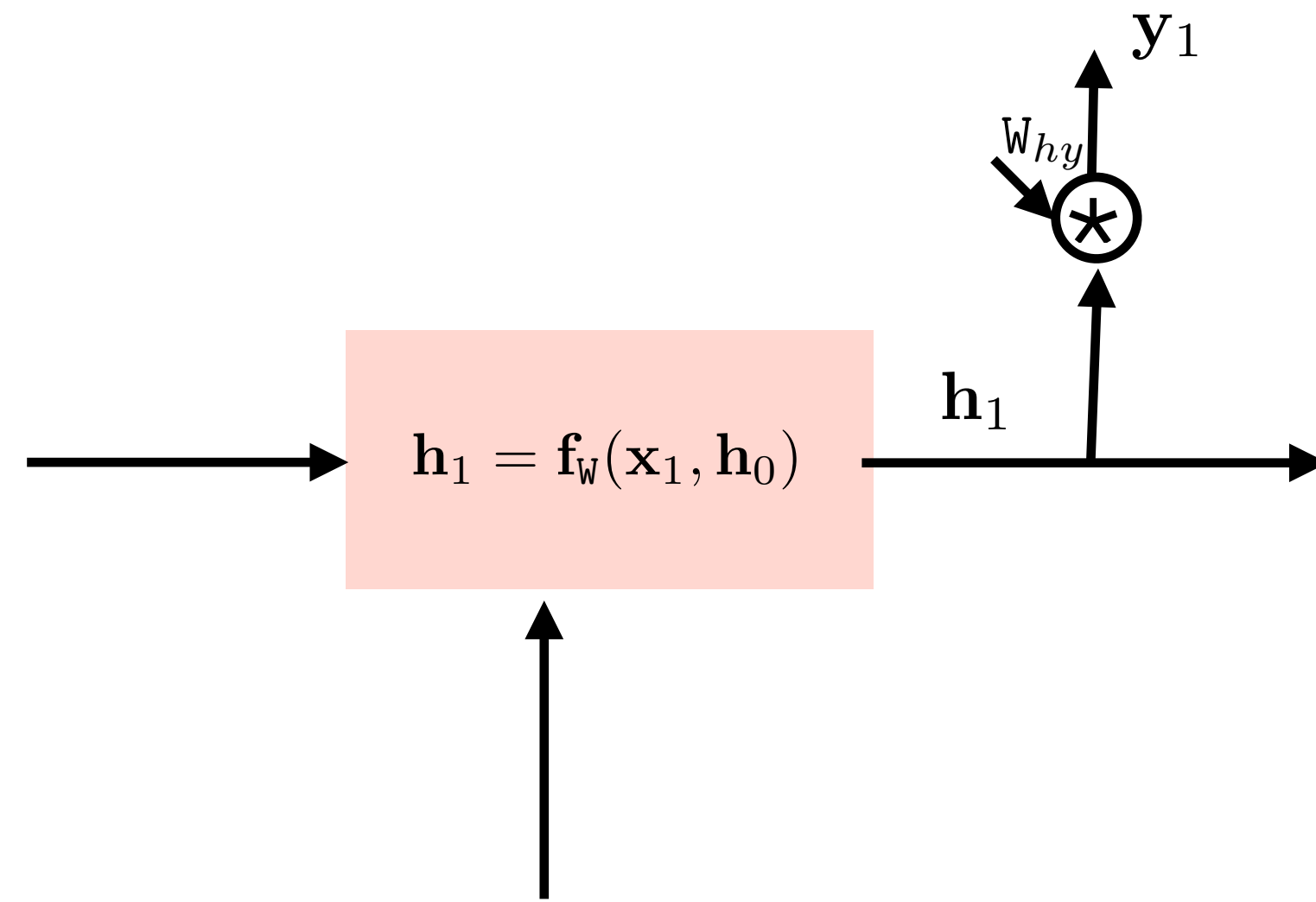


Given a finite input sequence:

we remove the recurrent connection by:

- successive substitution of inputs and
- unrolling the net

Simple recurrent block - feed-forward pass

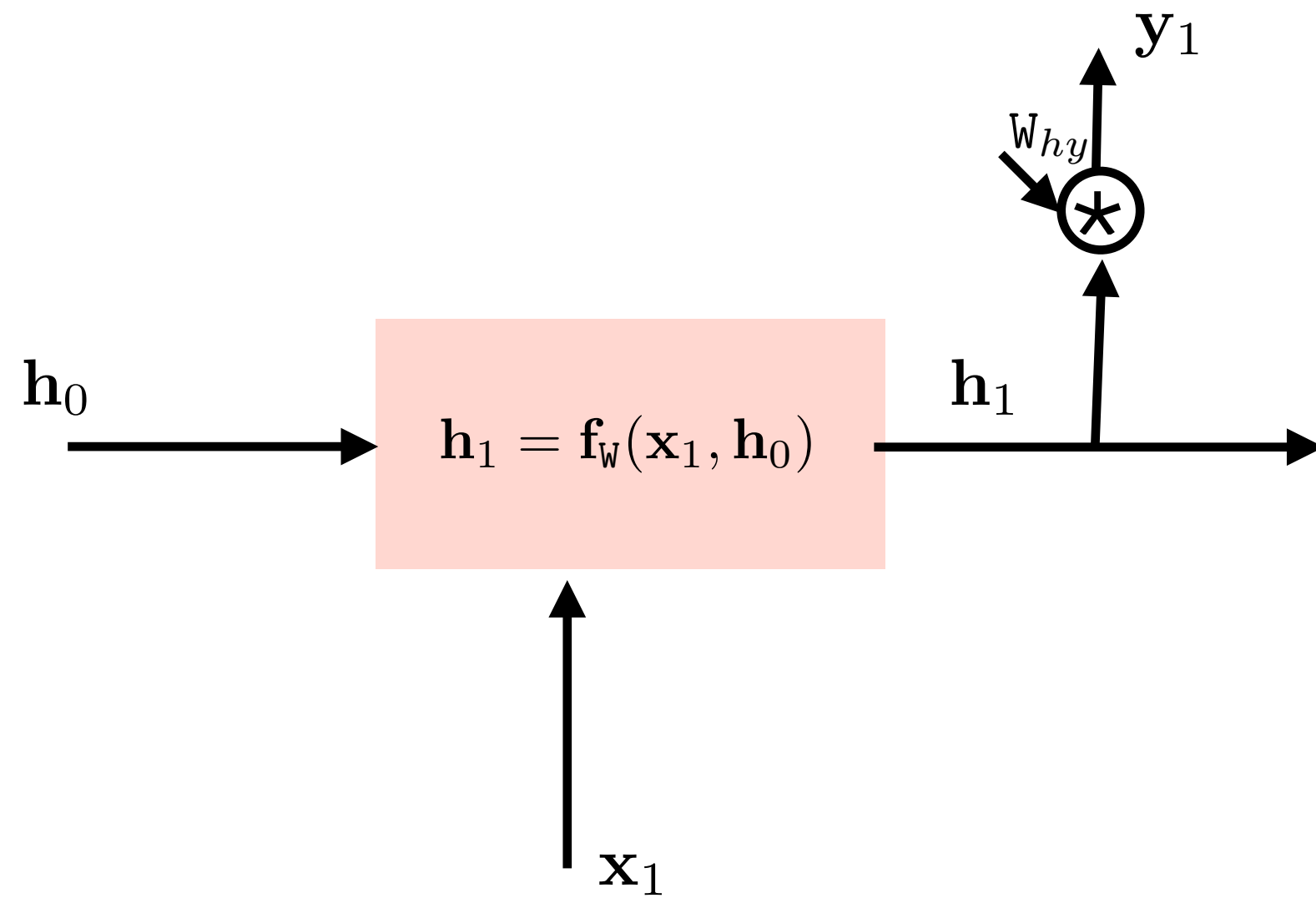


Given a finite input sequence:

we remove the recurrent connection by:

- successive substitution of inputs and
- unrolling the net

Simple recurrent block - feed-forward pass

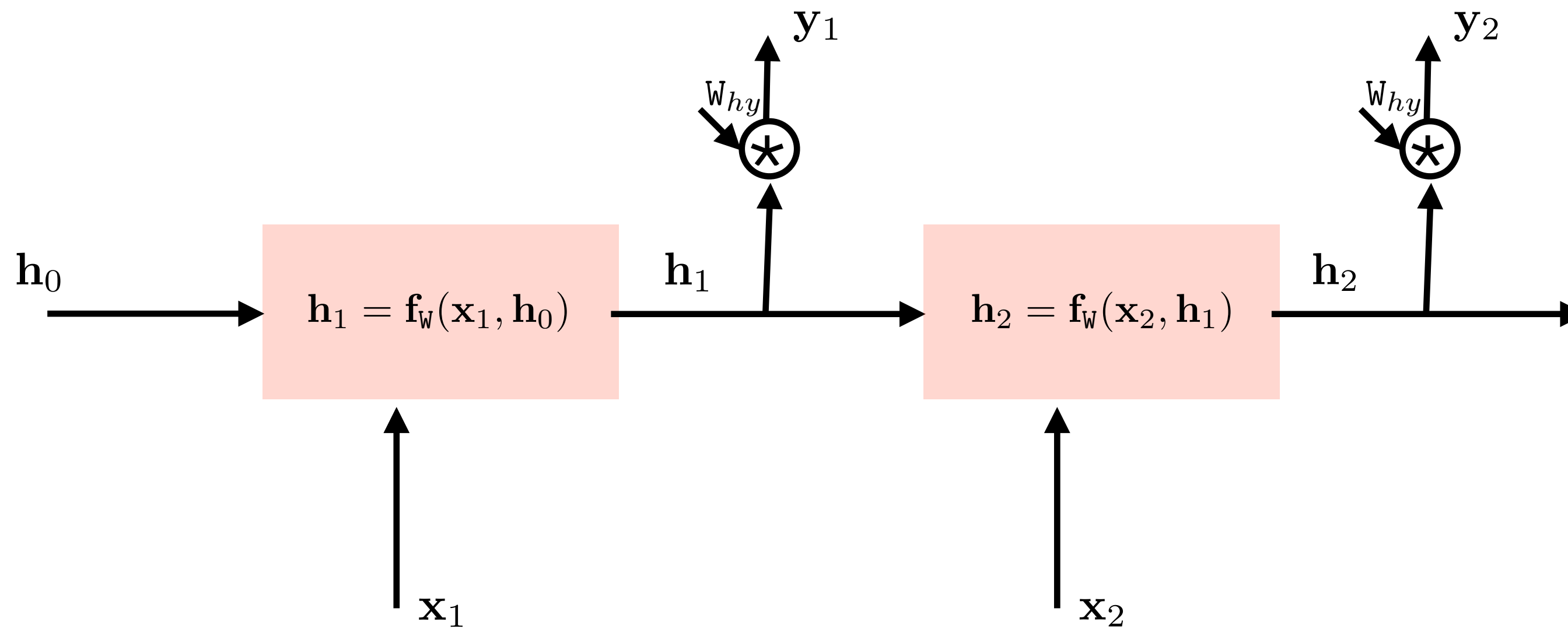


Given a finite input sequence:

we remove the recurrent connection by:

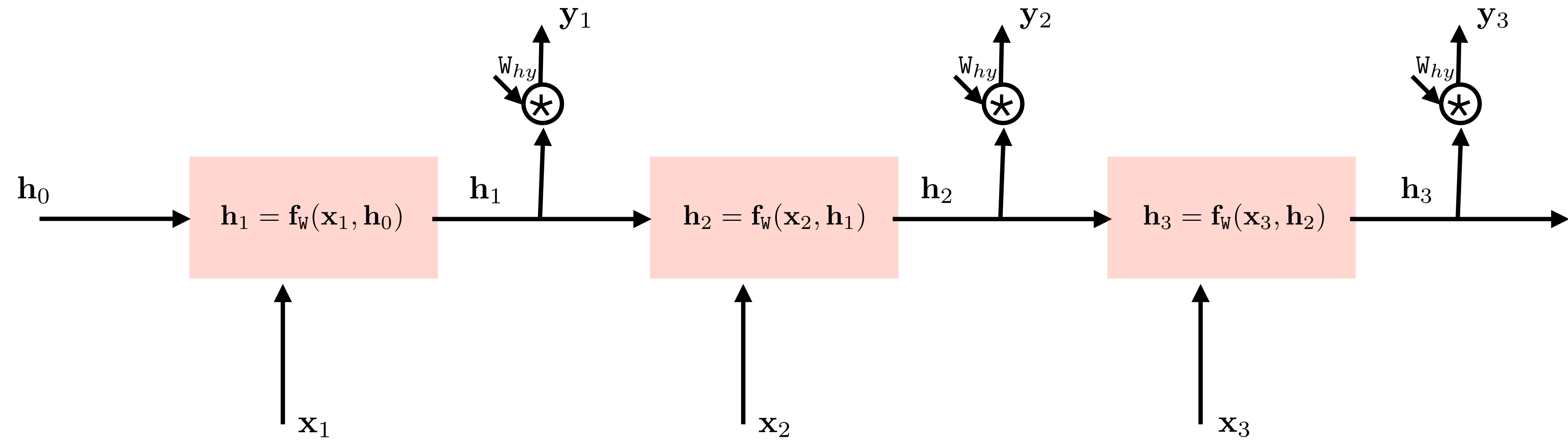
- successive substitution of inputs and
- unrolling the net

Simple recurrent block - feed-forward pass



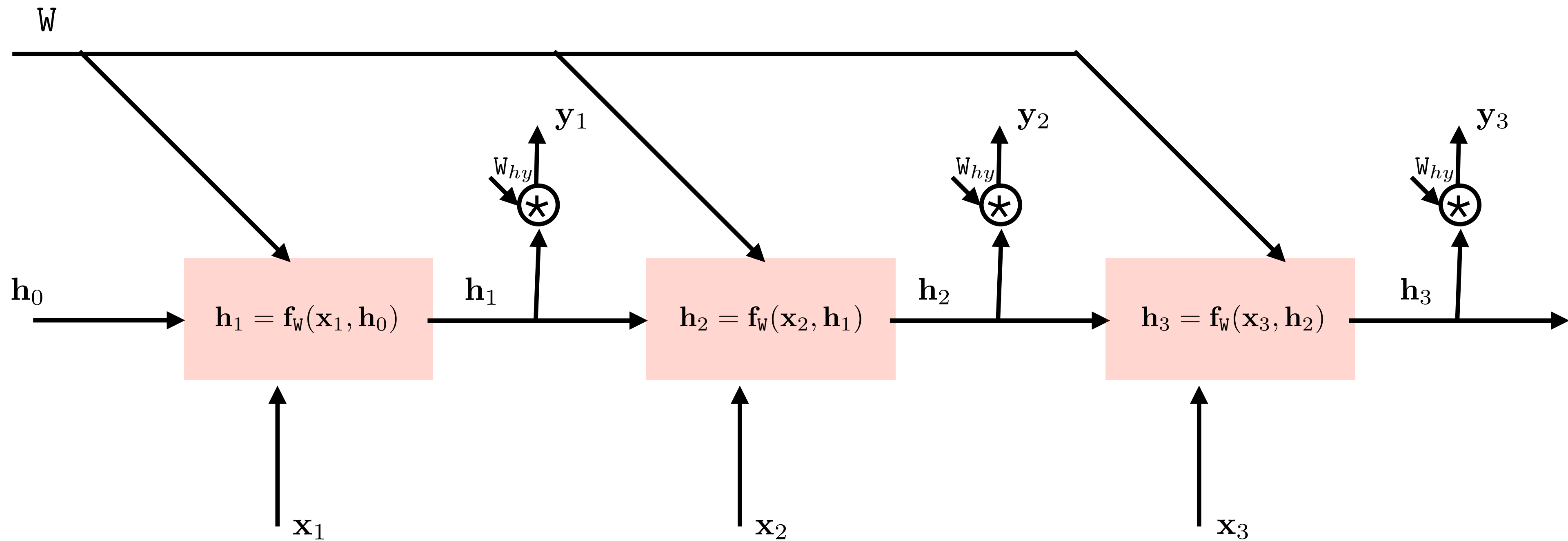
- Given a finite input sequence:
- we remove the recurrent connection by:
- successive substitution of inputs and
 - unrolling the net

Simple recurrent block - feed-forward pass



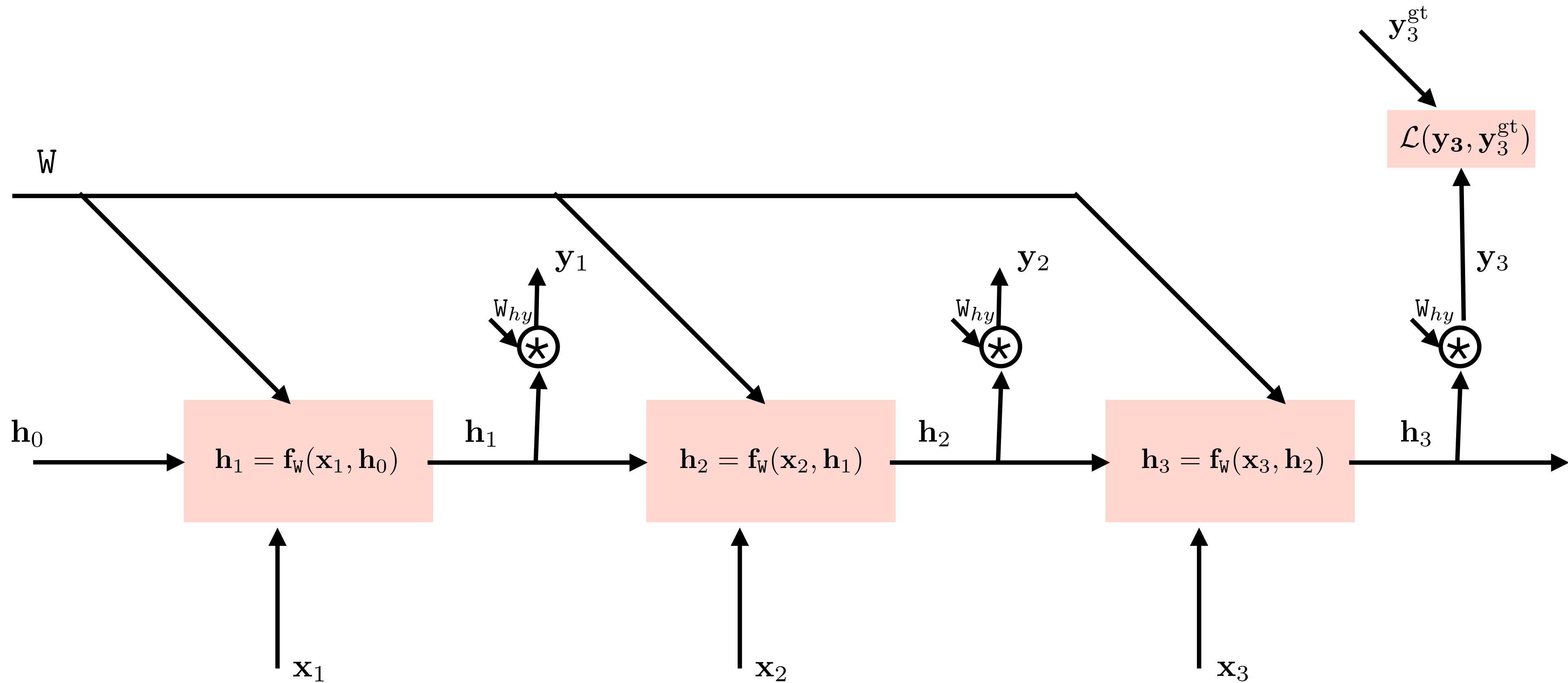
- Unrolled computational graph:
 - it is normal feedforward network

Simple recurrent block - feed-forward pass



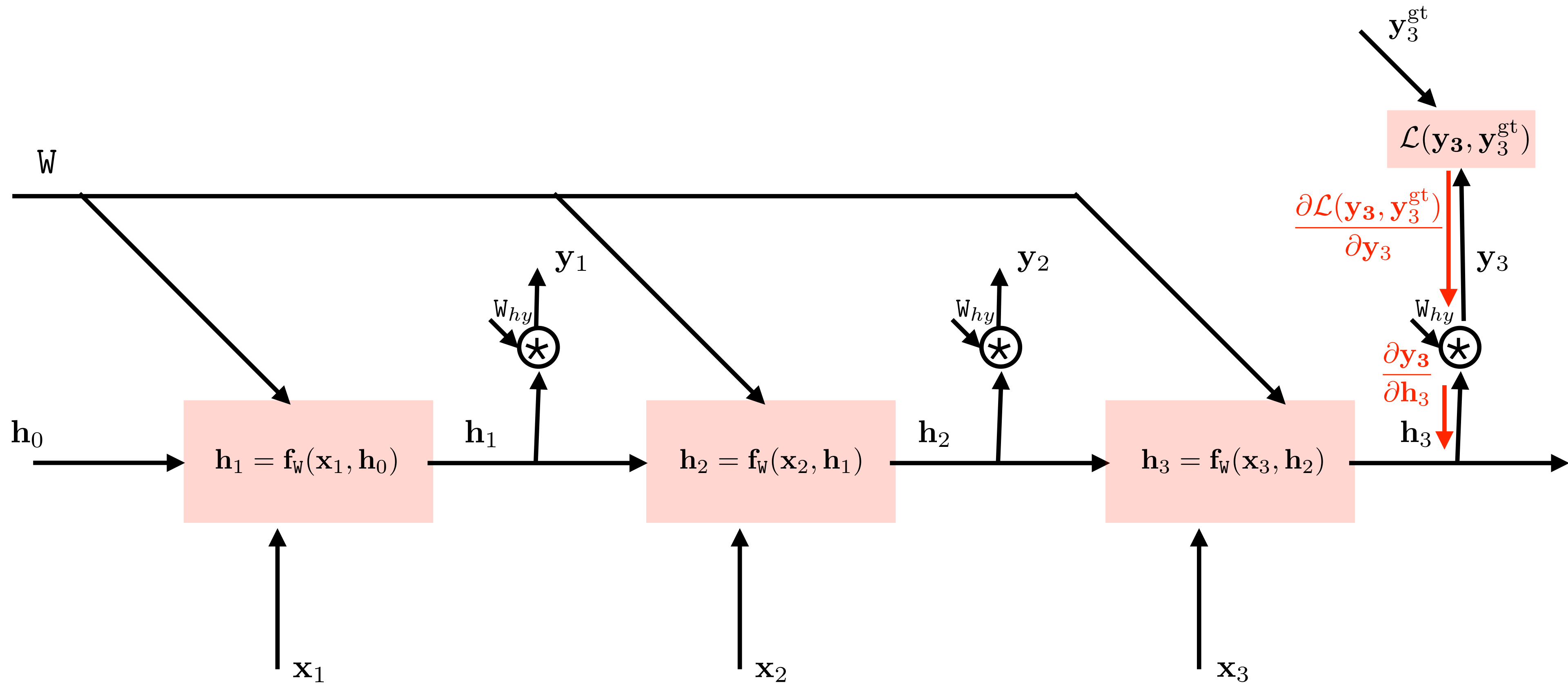
- Unrolled computational graph:
 - it is normal feedforward network
 - it consists of several same blocks with the same weights!

Simple recurrent block - backward pass



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

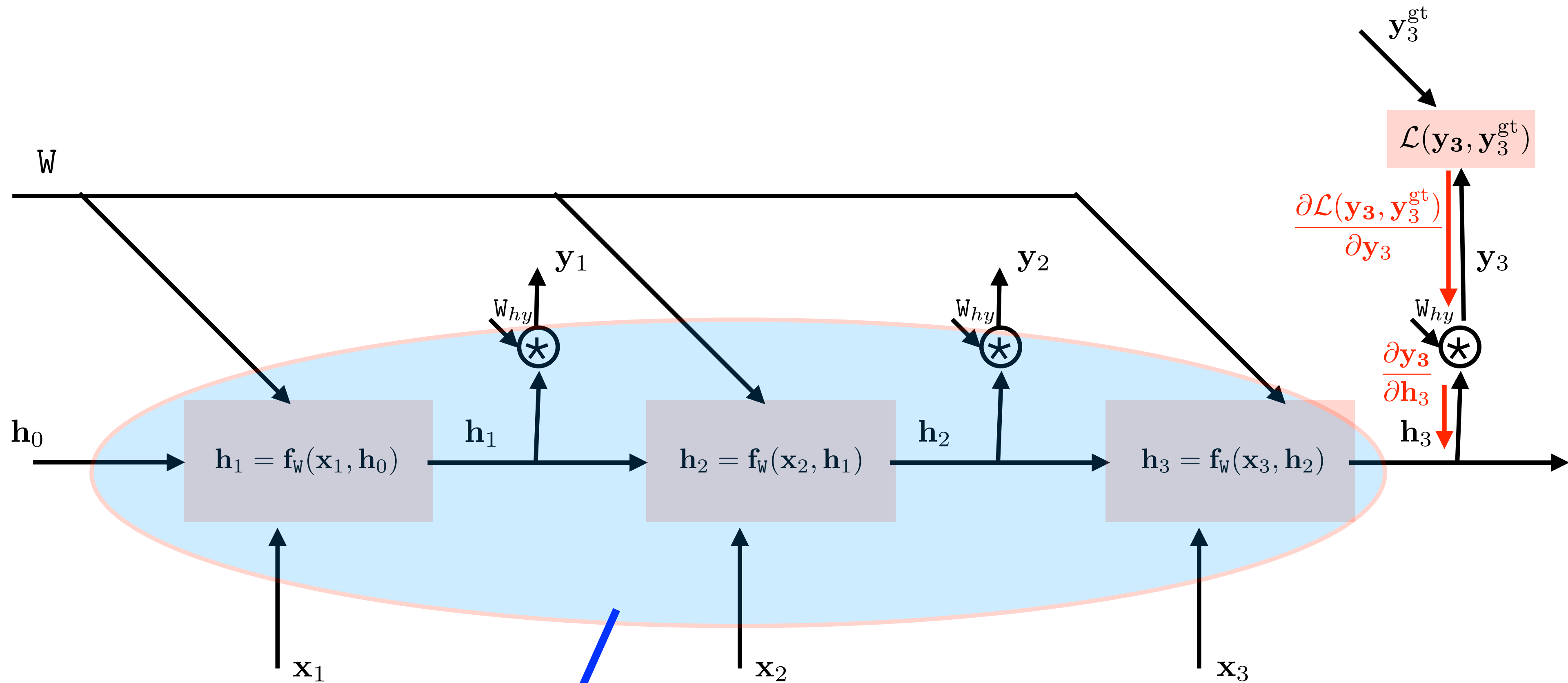
Simple recurrent block - backward pass



- Backprop:

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

Simple recurrent block - backward pass

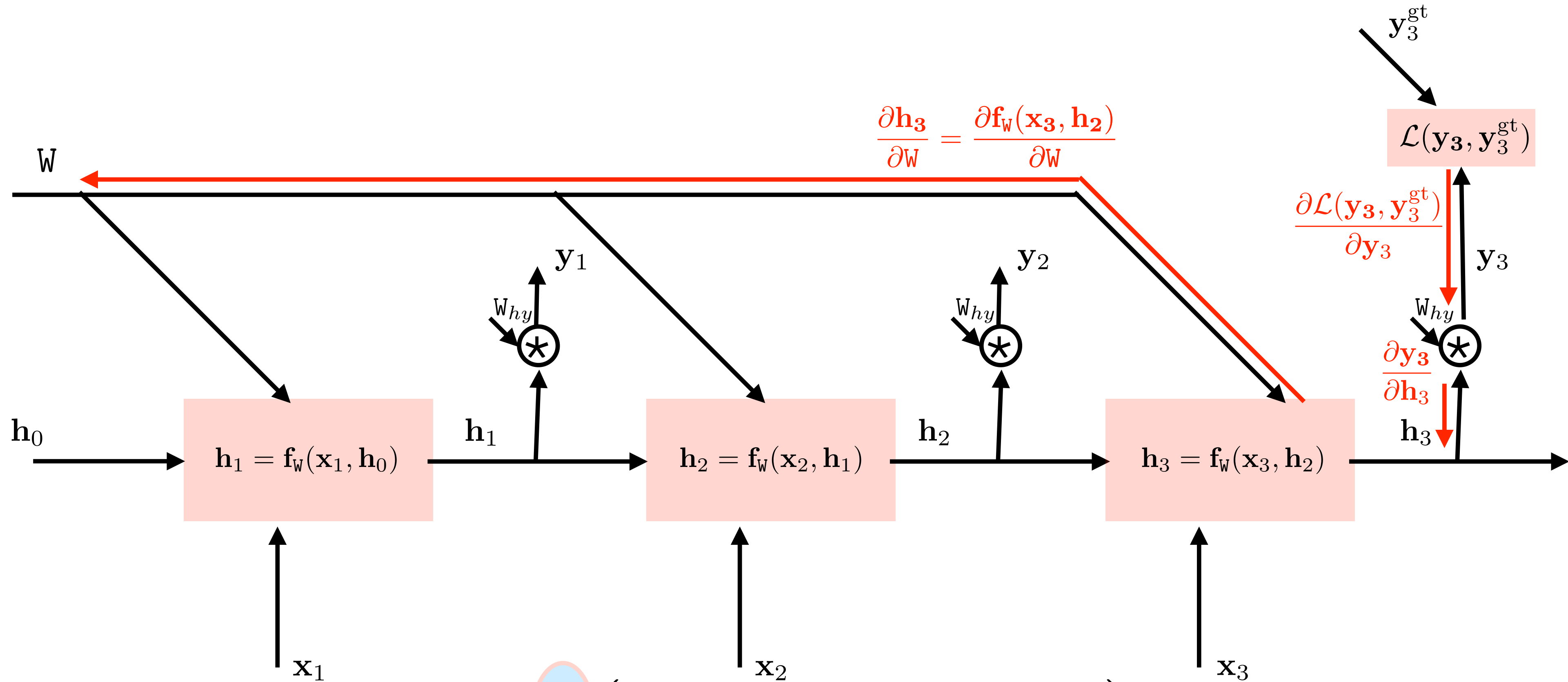


- Backprop:
 - differentiation of multi-dimensional composite function:

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

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

Simple recurrent block - backward pass



$$\frac{\partial \mathbf{h}_3}{\partial W} = \frac{\partial \mathbf{f}_W(\mathbf{x}_3, \mathbf{h}_2)}{\partial W}$$

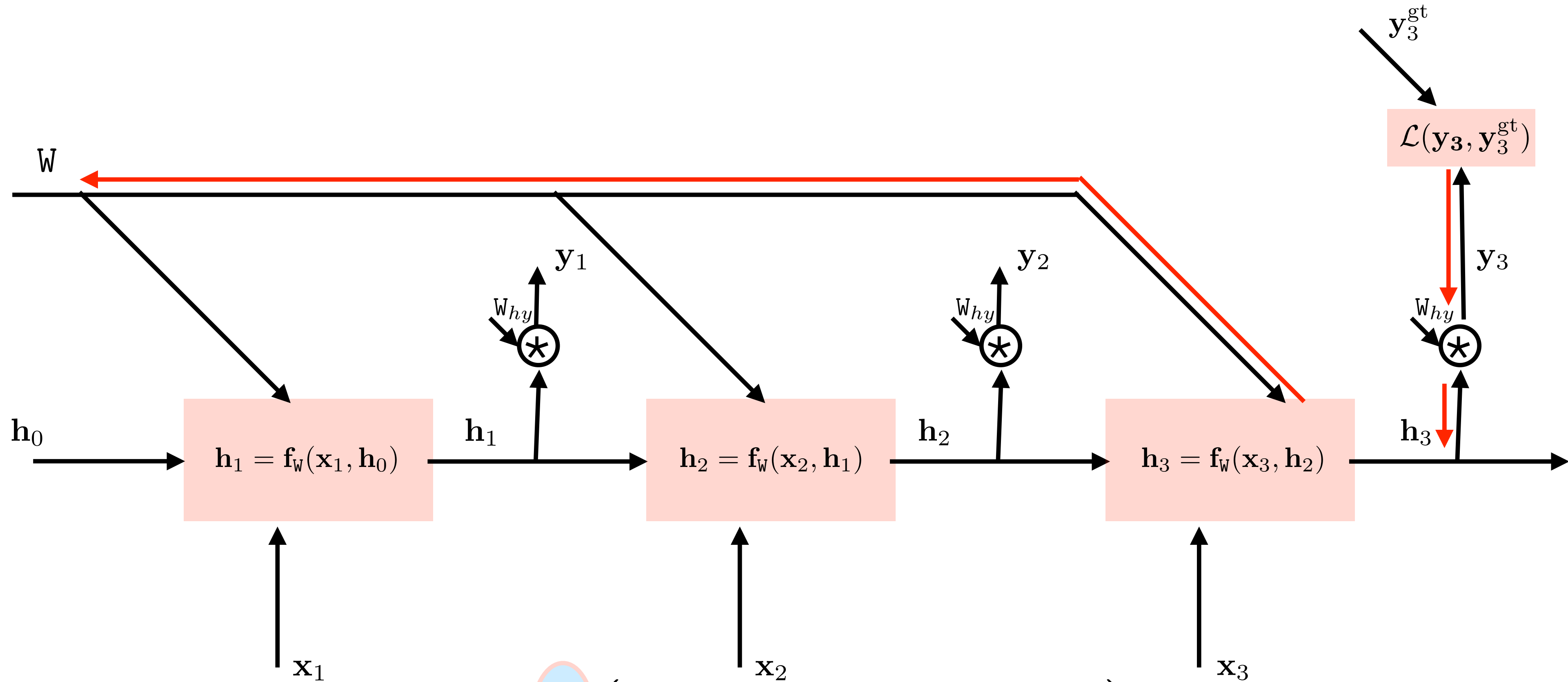
$$\frac{\partial \mathcal{L}(y_3, y_3^{\text{gt}})}{\partial y_3}$$

$$\frac{\partial y_3}{\partial \mathbf{h}_3}$$

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

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

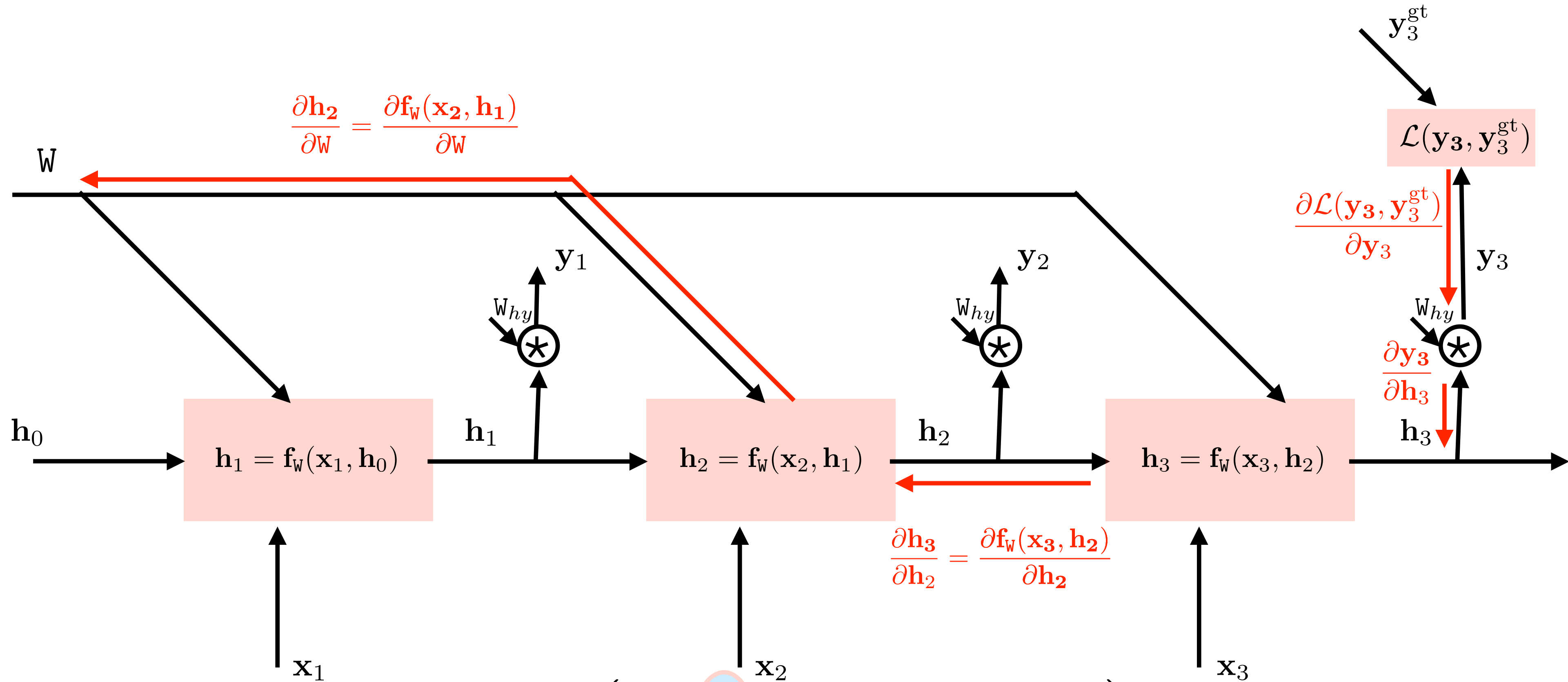
Simple recurrent block - backward pass



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

$$\frac{\partial \mathcal{L}(y_3, y_3^{gt})}{\partial W} = \frac{\partial \mathcal{L}(y_3, y_3^{gt})}{\partial y_3} \frac{\partial y_3}{\partial h_3} \frac{\partial f_W(\mathbf{x}_3, h_2)}{\partial W} +$$

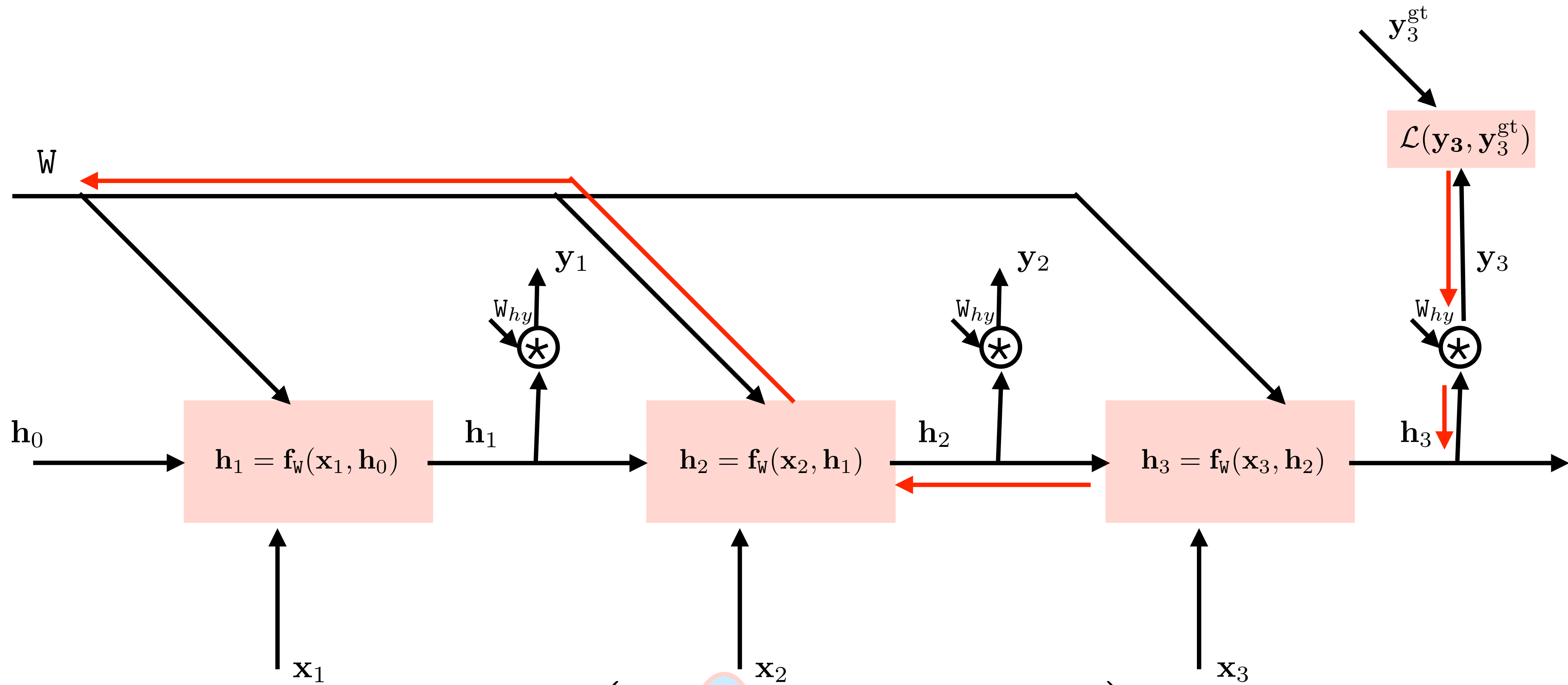
Simple recurrent block - backward pass



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

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

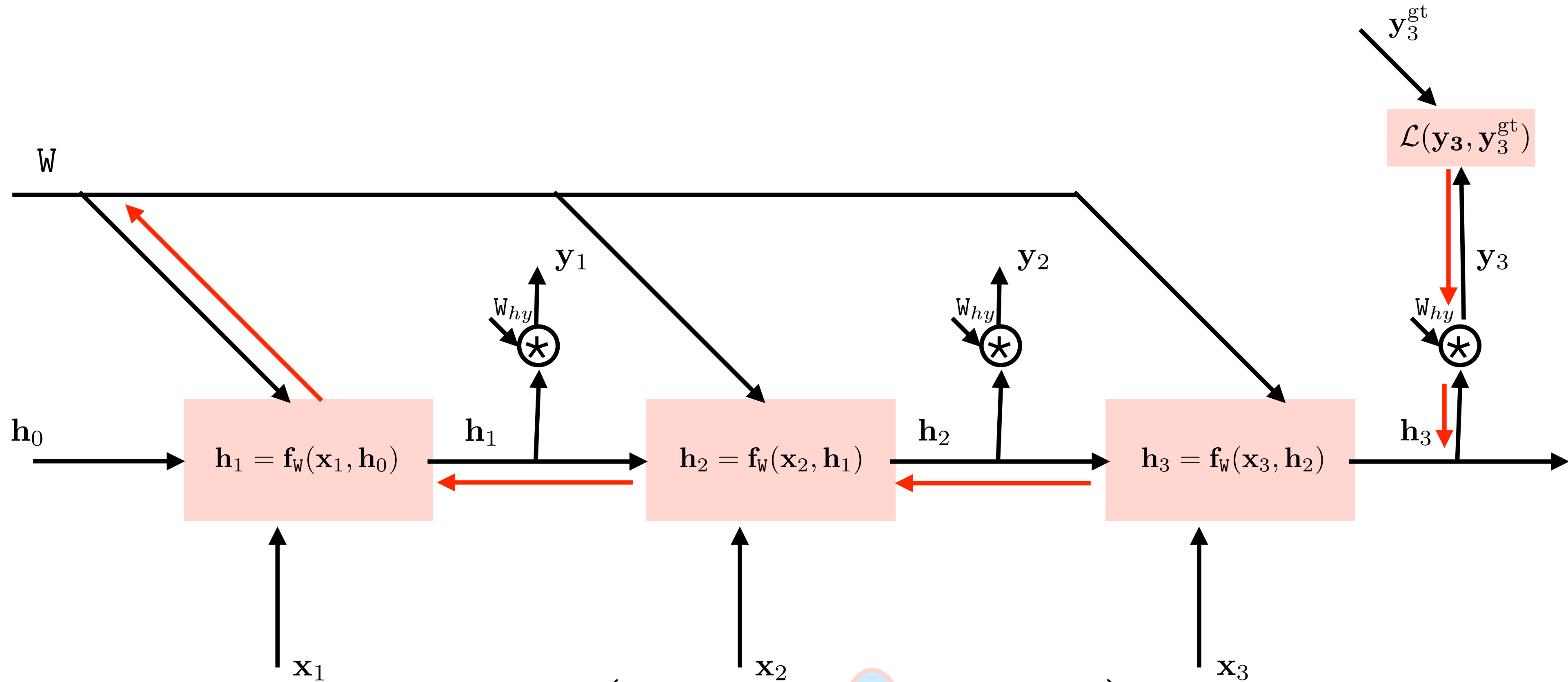
Simple recurrent block - backward pass



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

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

Simple recurrent block - backward pass



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

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

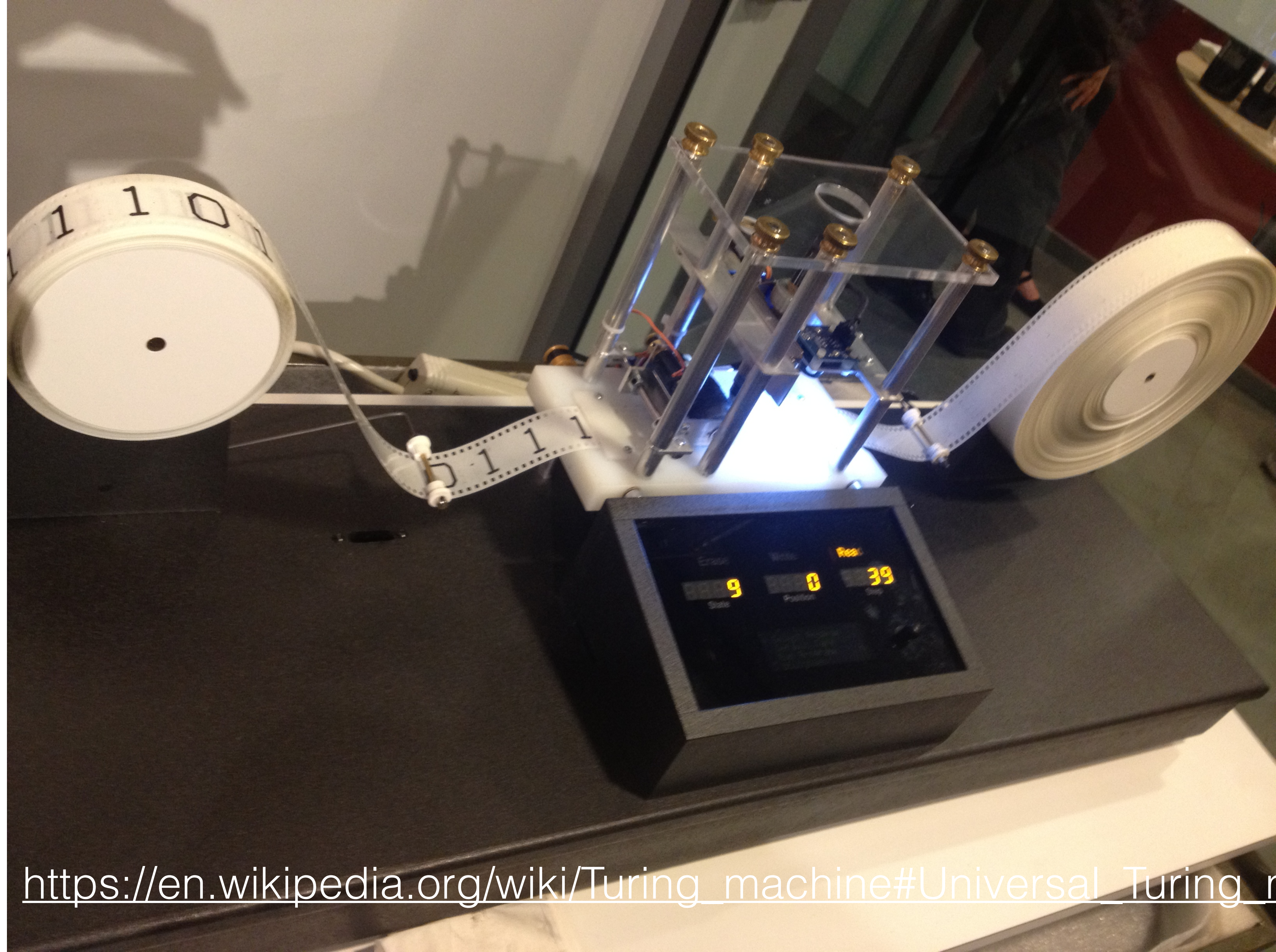
RNN vs feedforward network

$$\mathbf{h}_t = \mathbf{f}_W(\mathbf{x}_t, \mathbf{h}_{t-1})$$

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

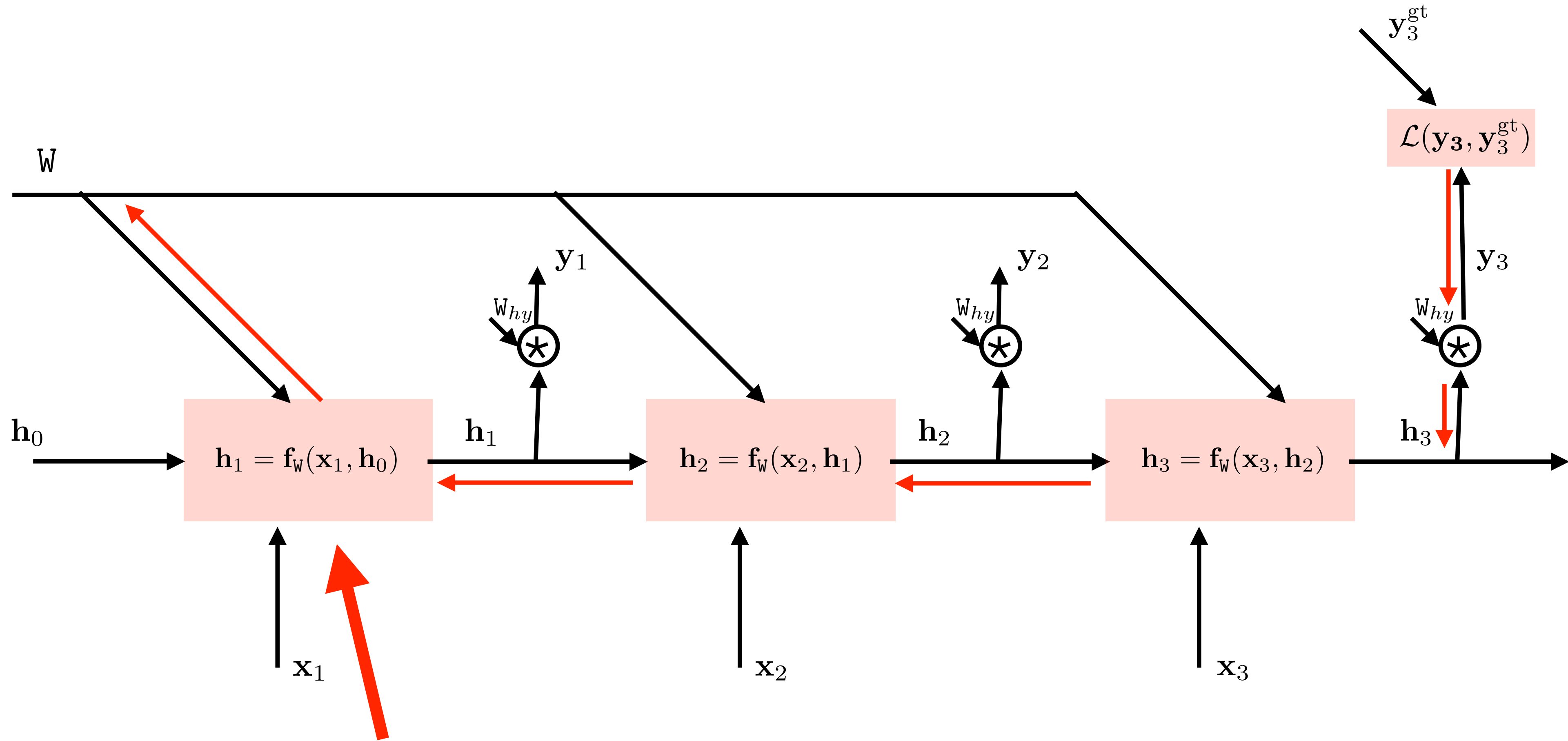
$$\mathbf{h}_t = \mathbf{g}^{(t)}(\mathbf{x}_t, \mathbf{x}_{t-1}, \dots \mathbf{x}_1)$$

- Advantage of RNN wrt stacking the input sequence into a long vector and using a common feedforward network:
 - RNN works for different sequence lengths
 - RNN share weights between different time instances (similarly as convolution on spatial domain).
- RNN is universal (can compute any function computable by Turing machine)



https://en.wikipedia.org/wiki/Turing_machine#Universal_Turing_m

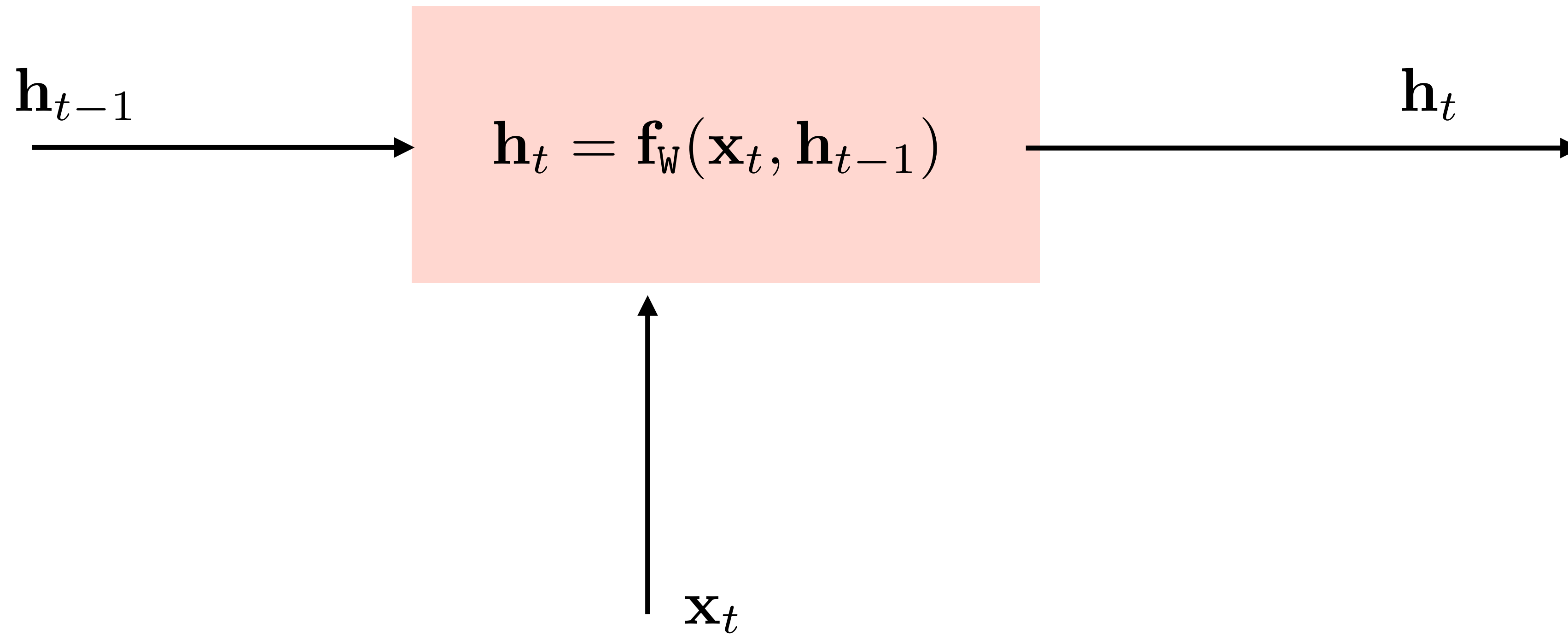
Simple recurrent block - backward pass



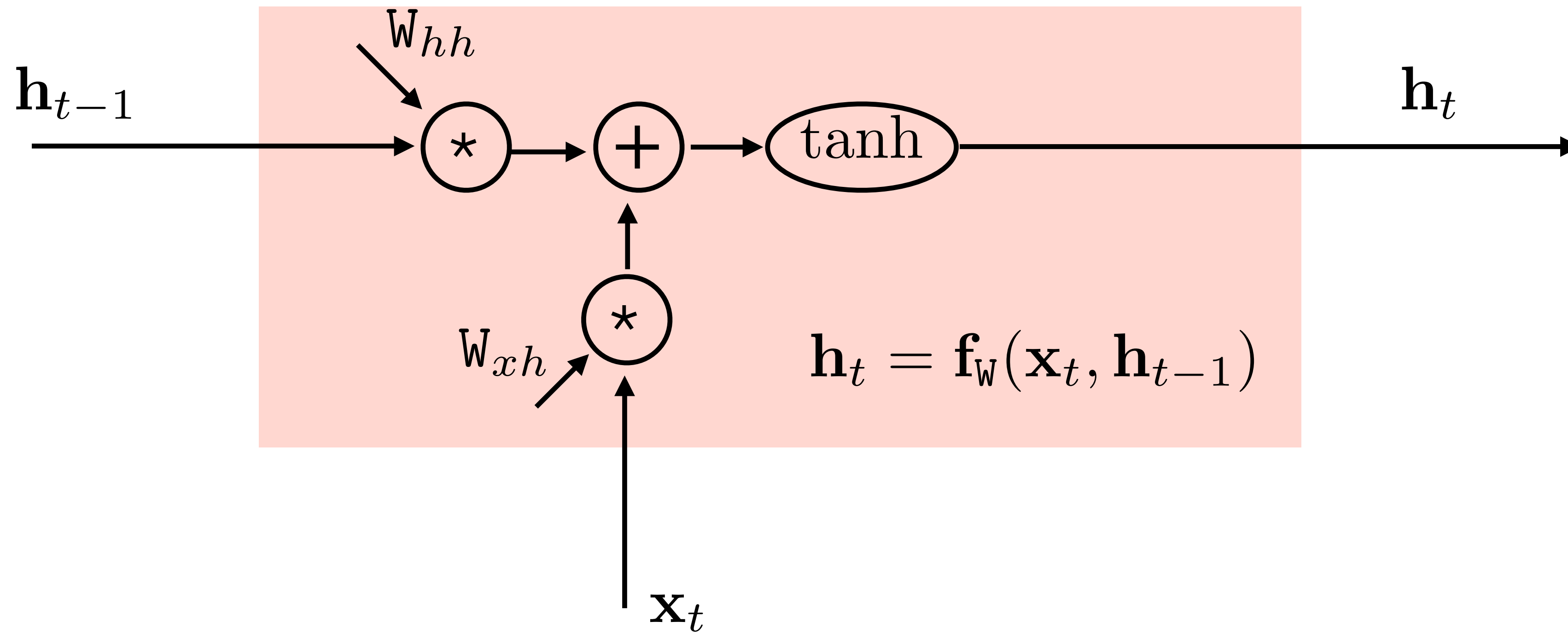
deep blocks often suffer from vanishing gradients
=> better structure needed

LSTM (kind of ResNet for recurrent networks)

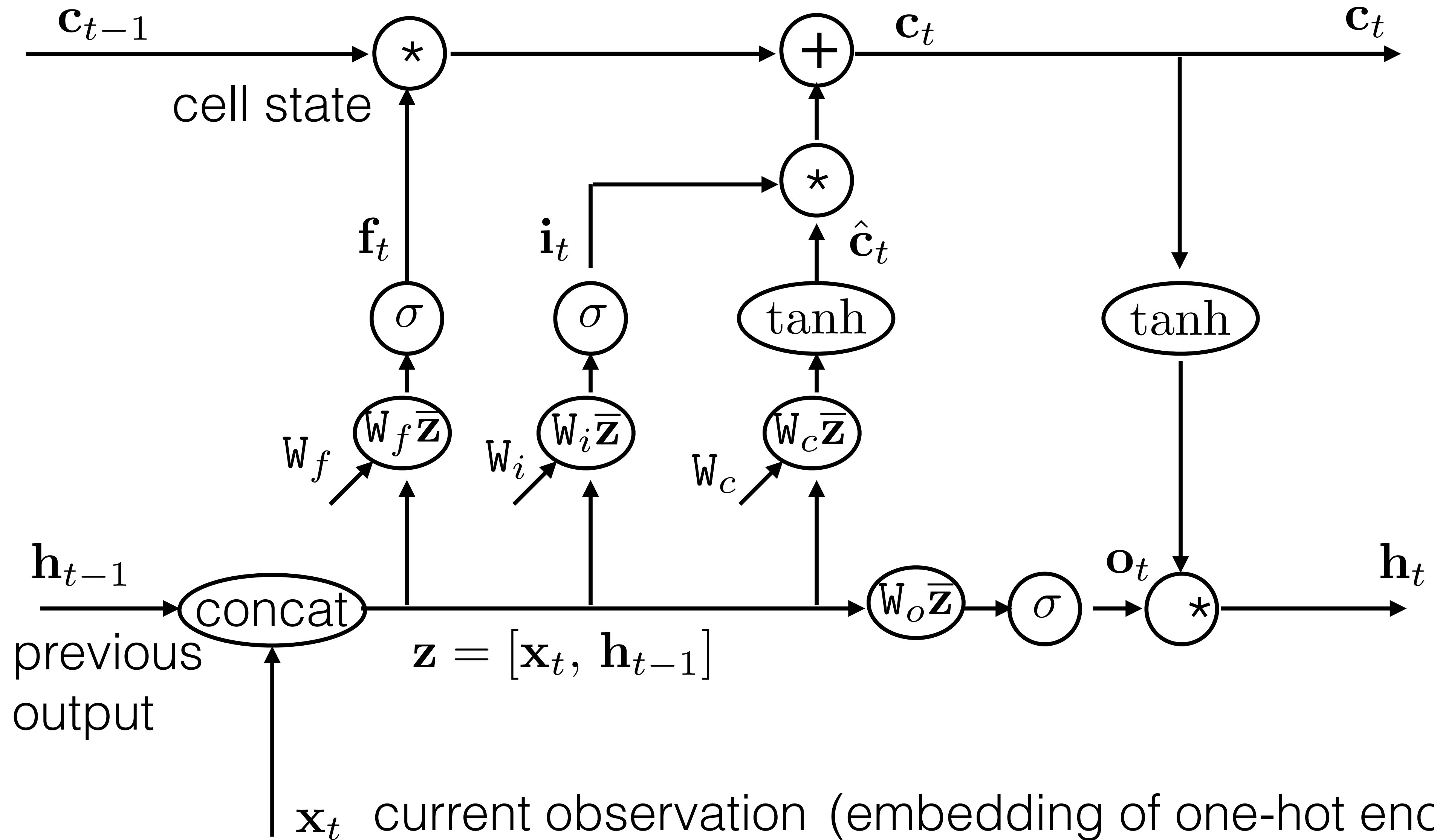
Simple recurrent block



Simple recurrent block

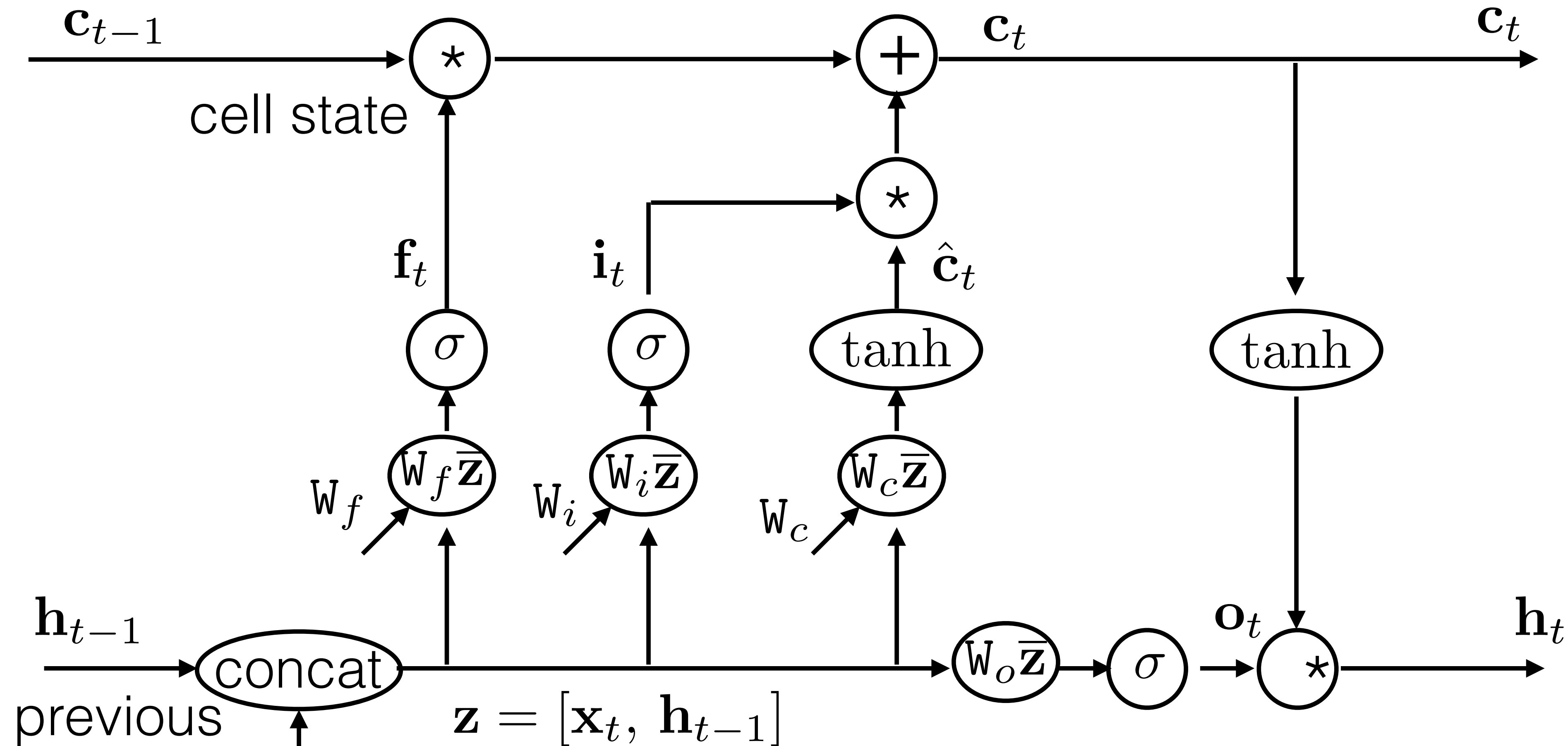


LSTM block



“I live with my parents, ...”

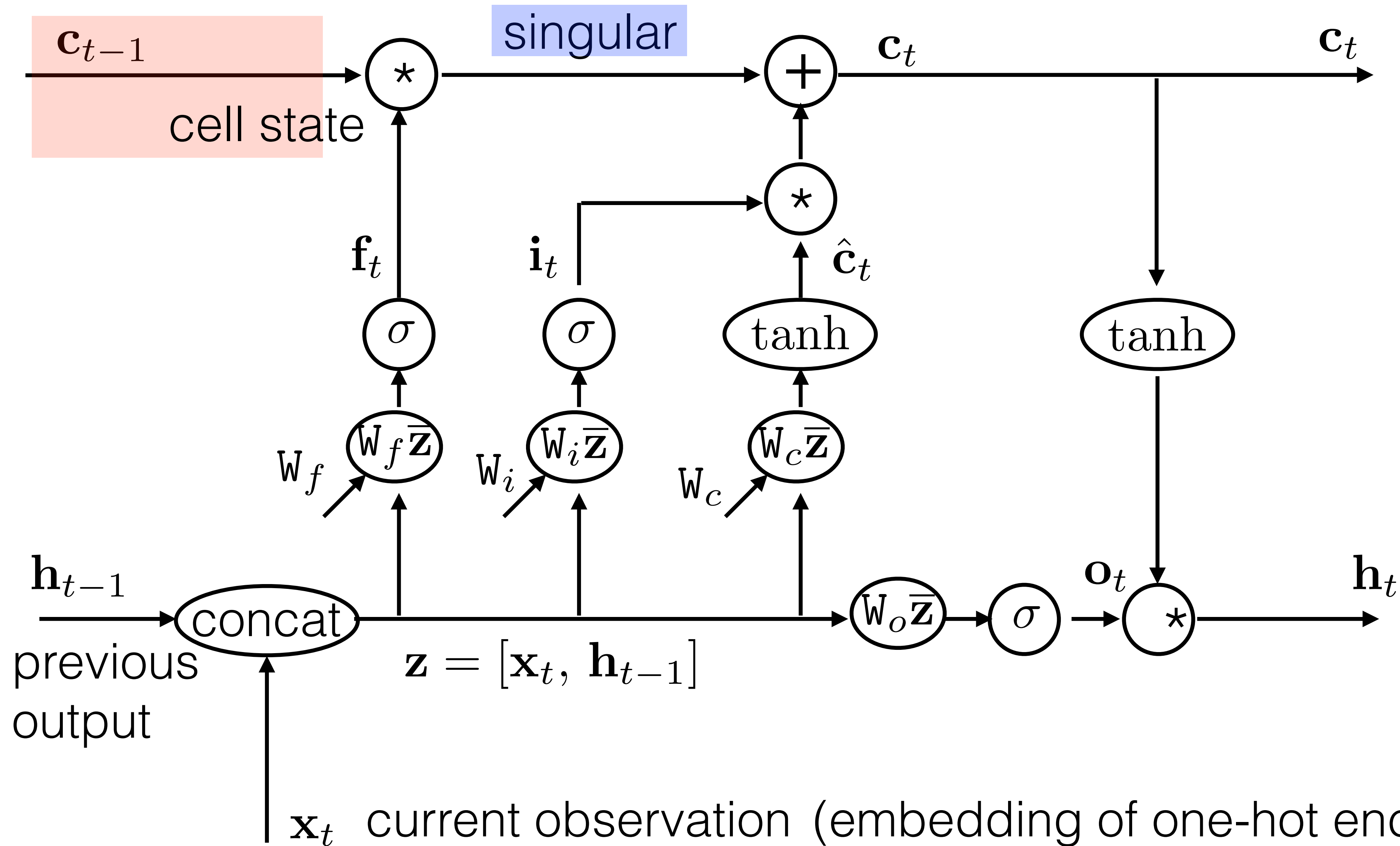
LSTM block



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

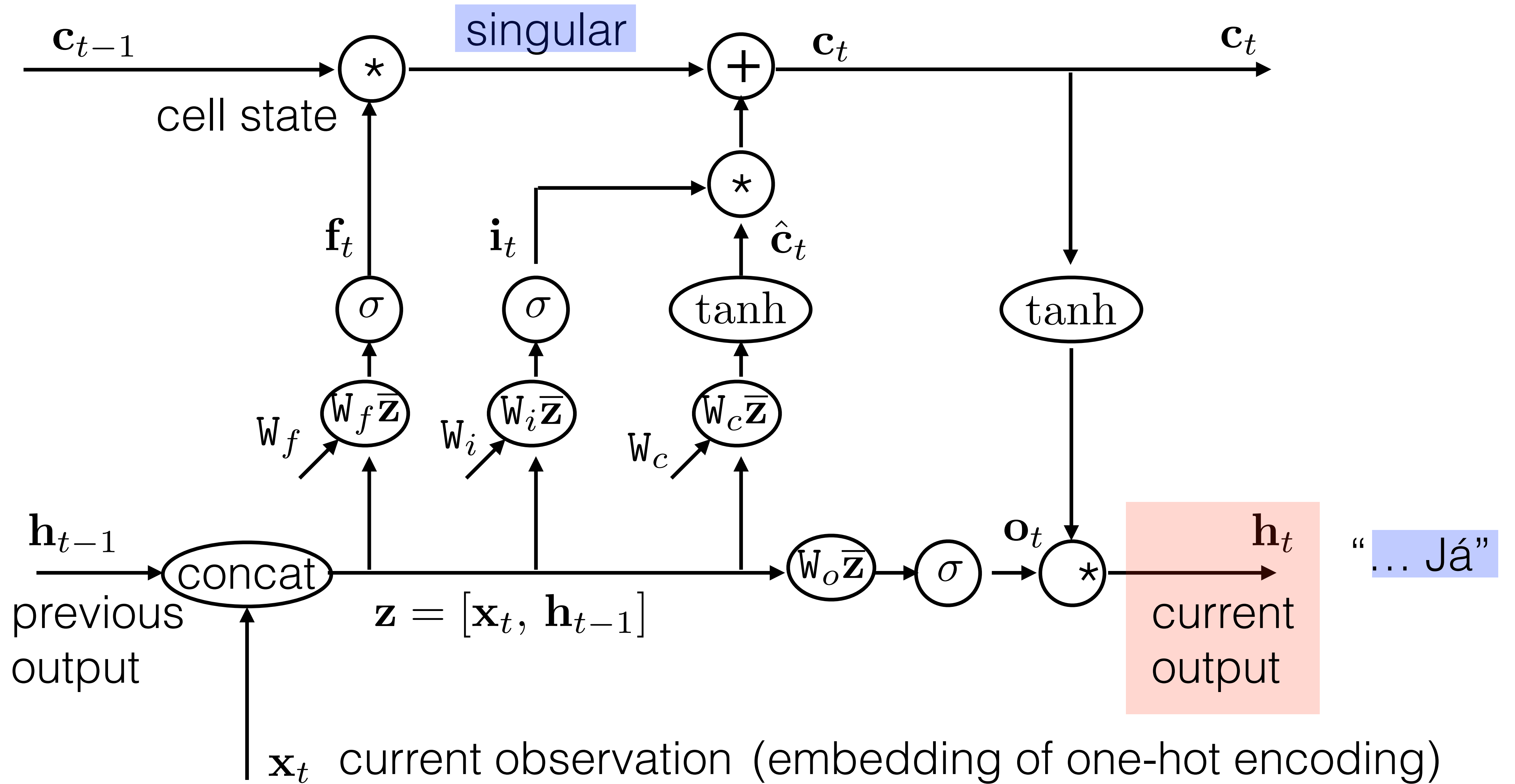
“I live with my parents, ...”

LSTM block

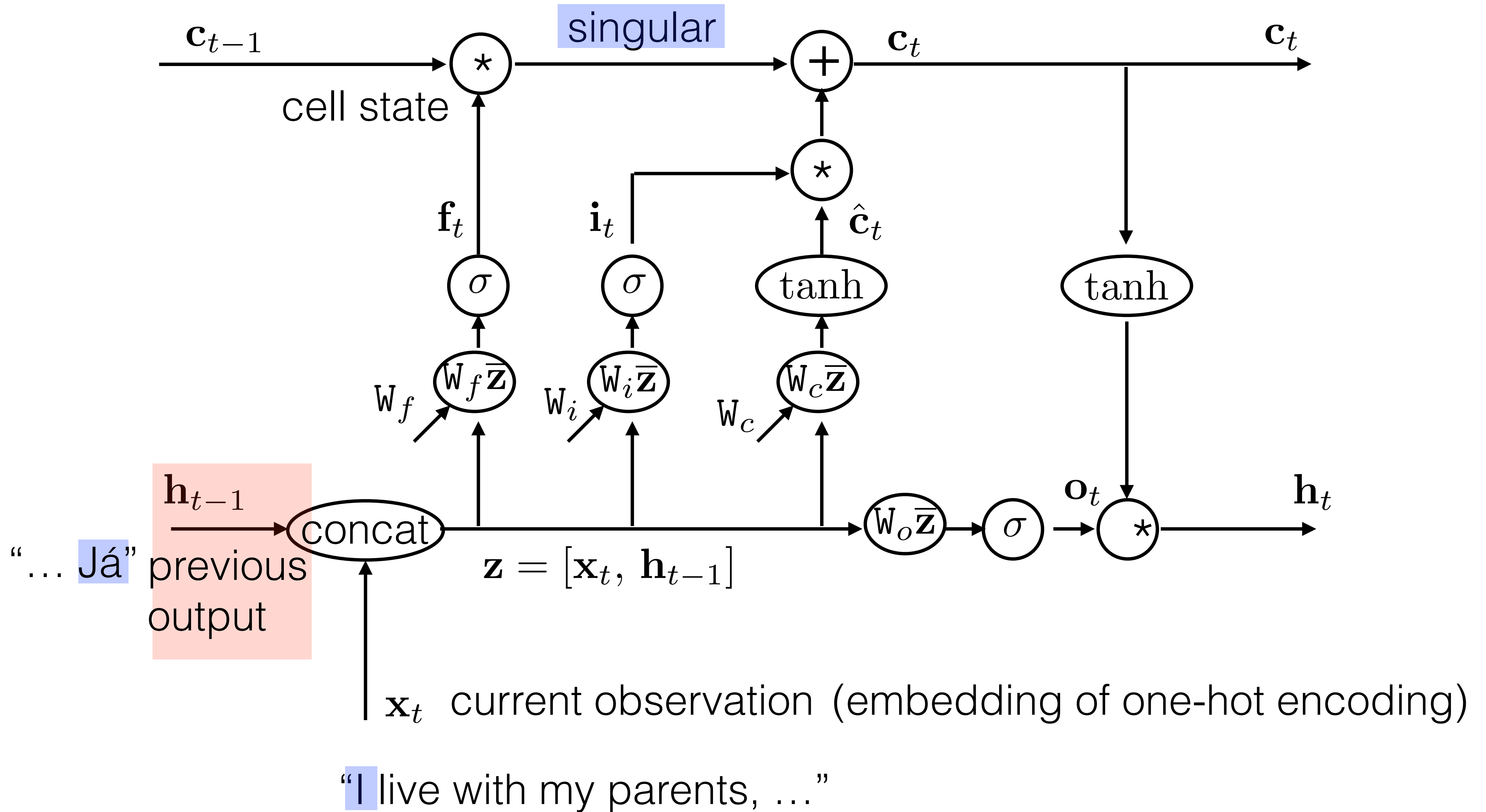


“I live with my parents, ...”

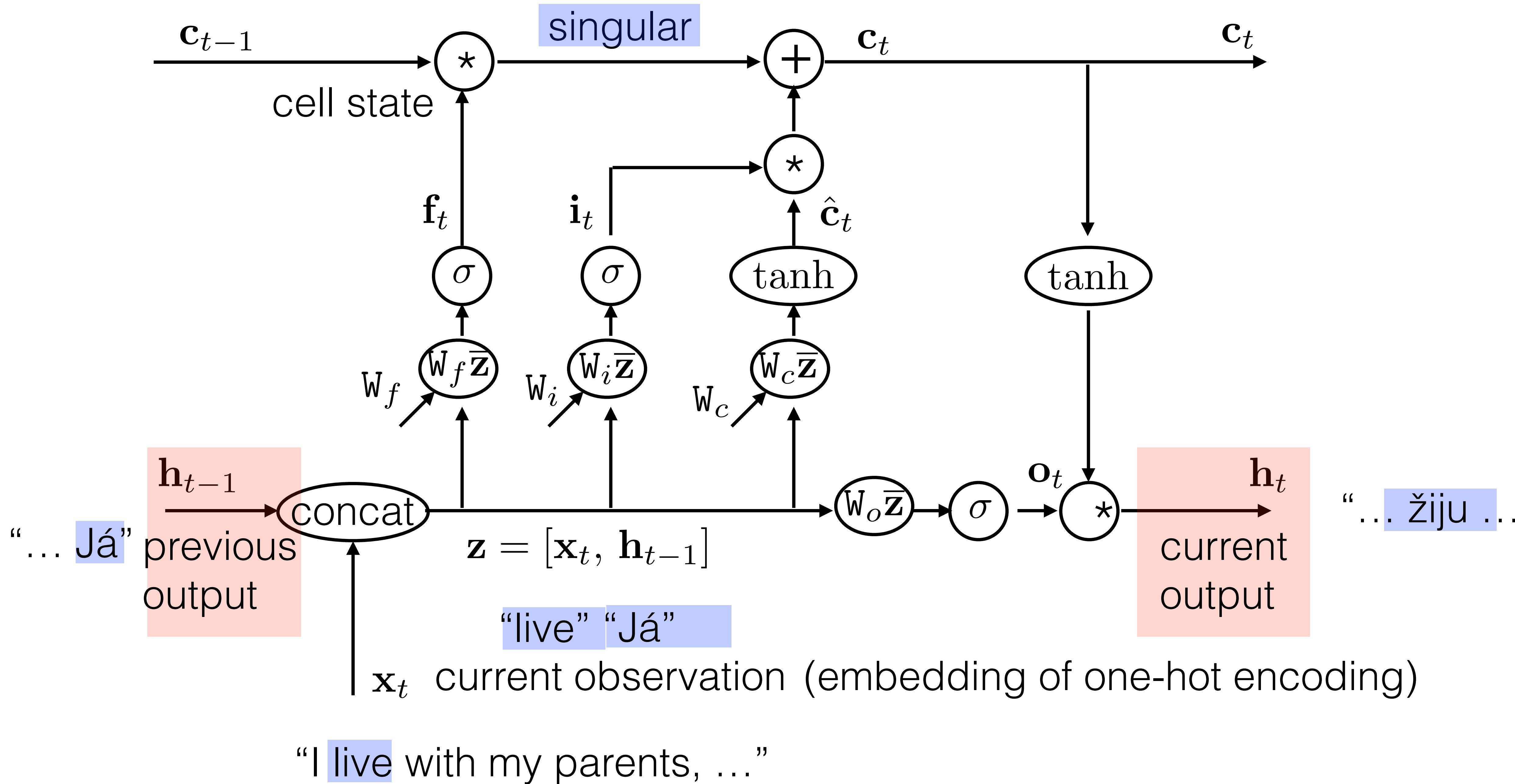
LSTM block



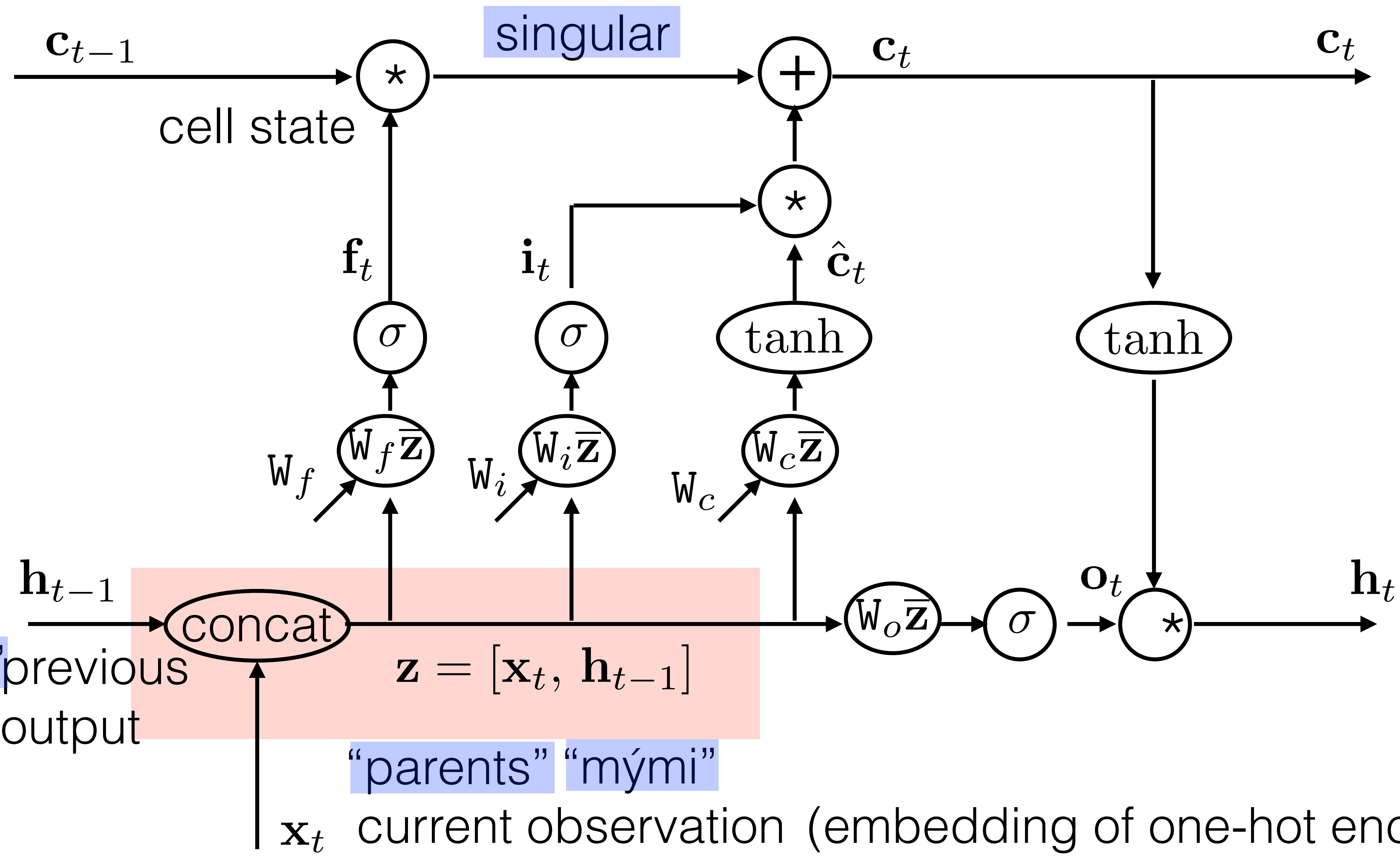
LSTM block



LSTM block



LSTM block



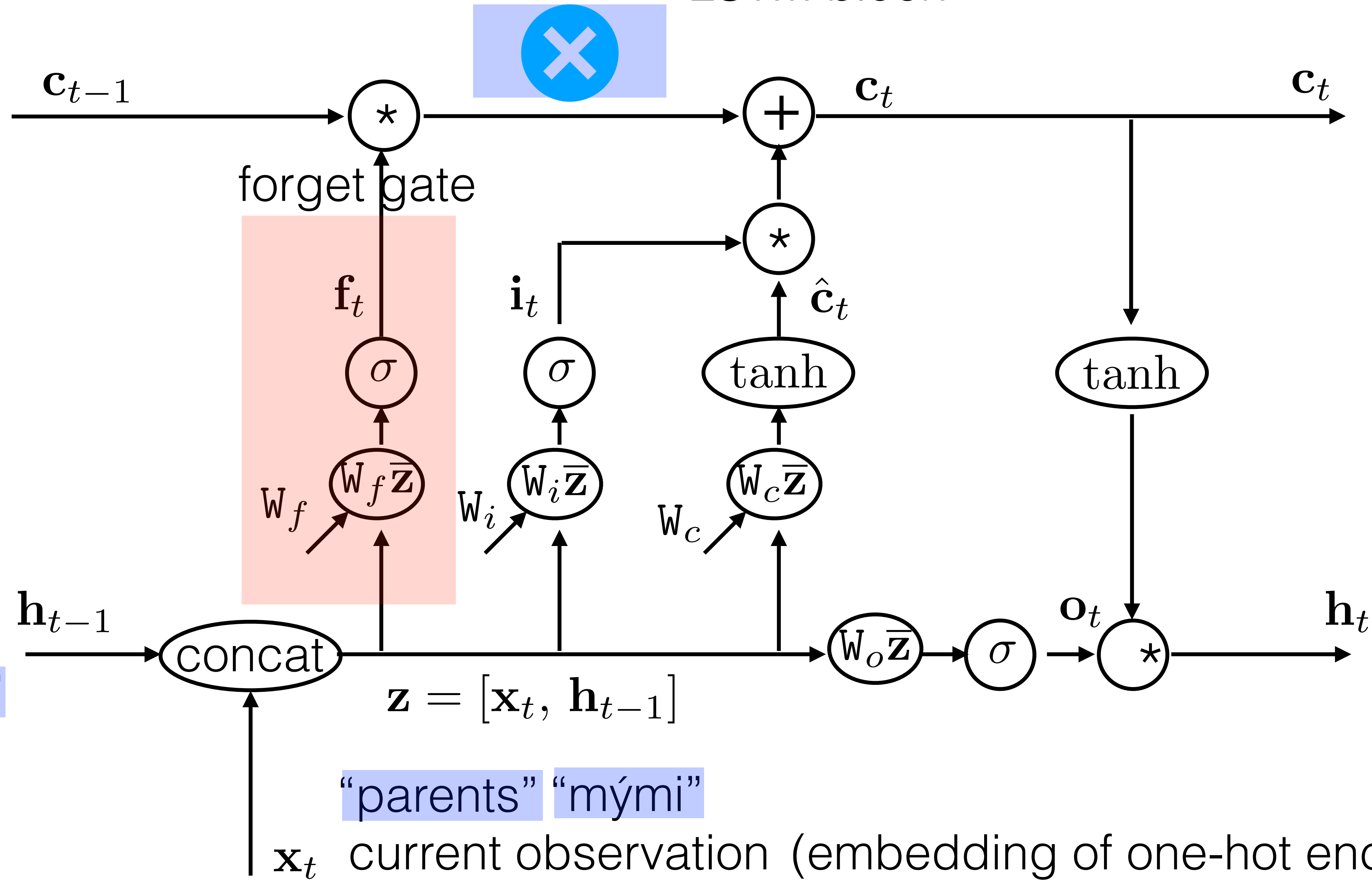
“... **mými**” previous output

“parents” “mými”

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

“I live with my **parents**, ...”

LSTM block



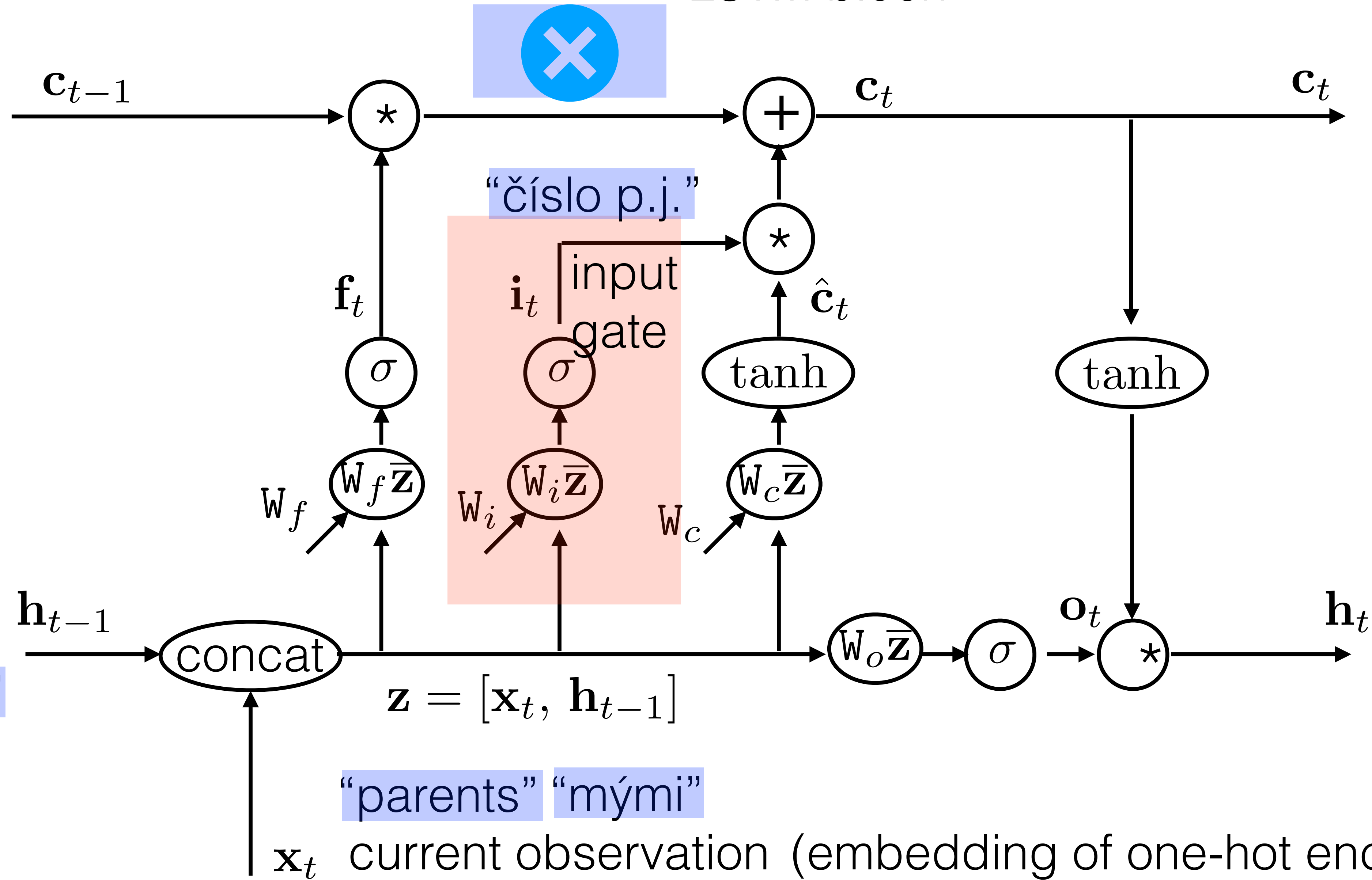
“... **mými**”

“**parents**” “**mými**”

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

“I live with my **parents**, ...”

LSTM block

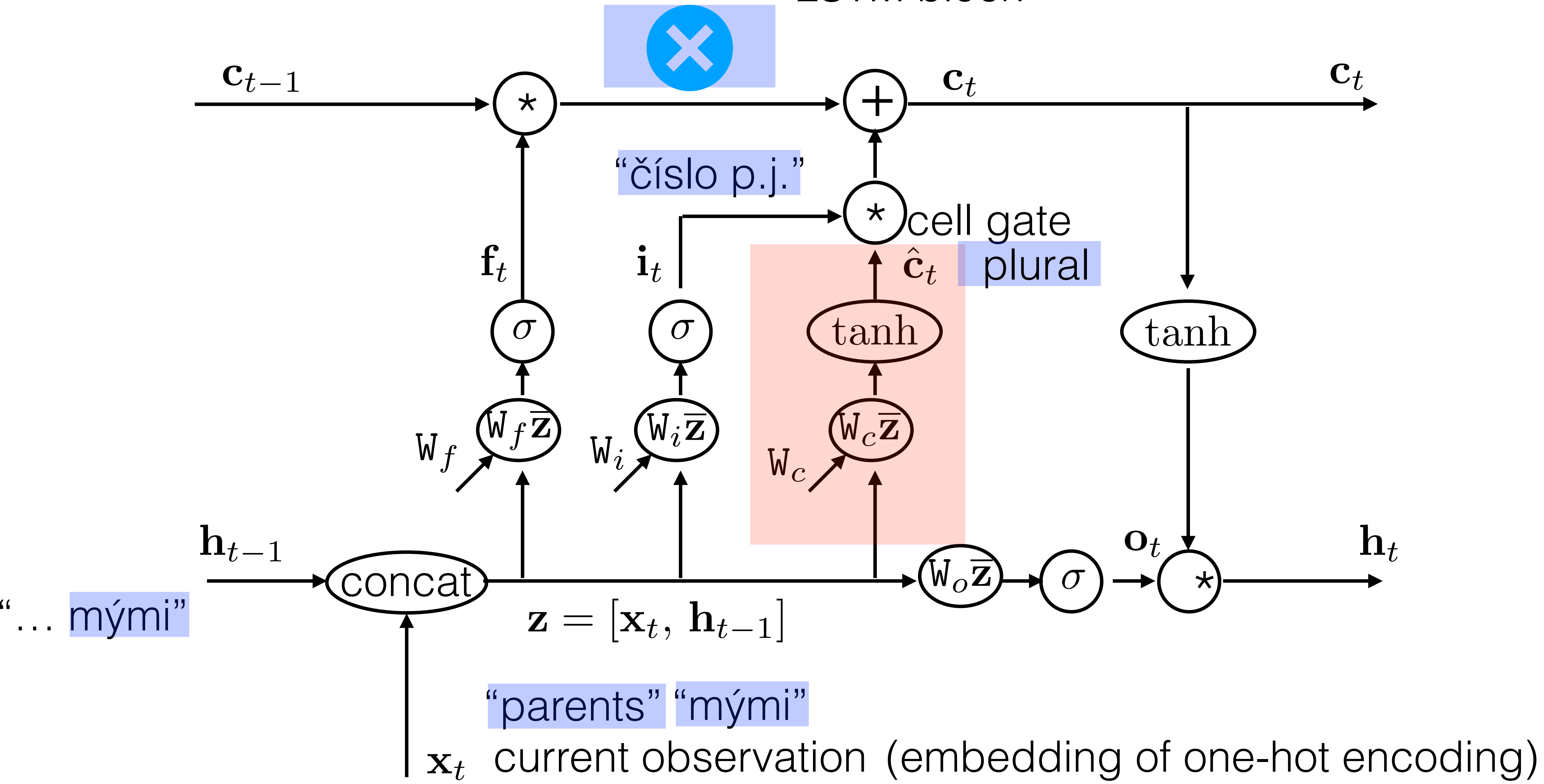


“parents” “mými”

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

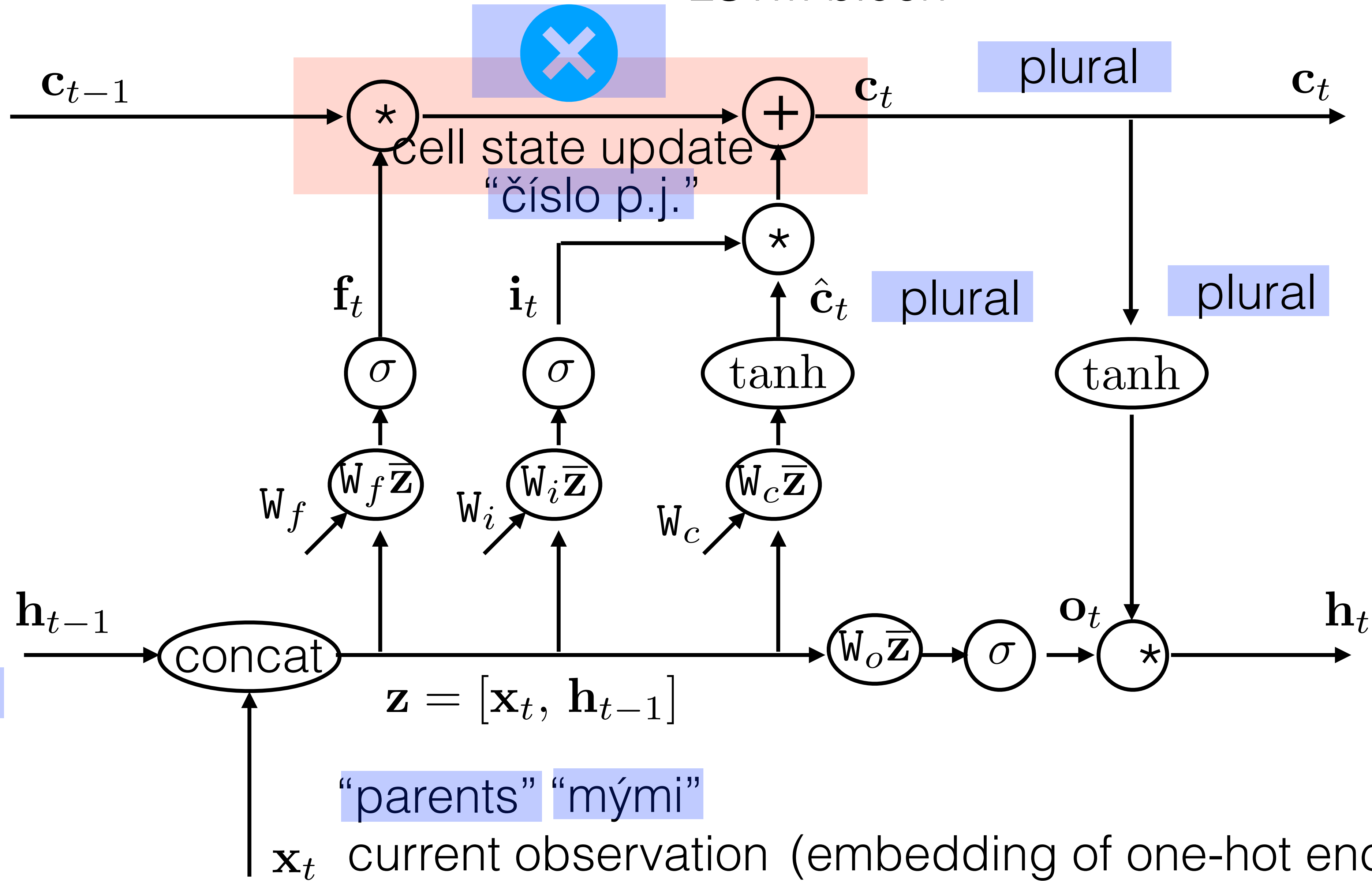
“I live with my parents, ...”

LSTM block



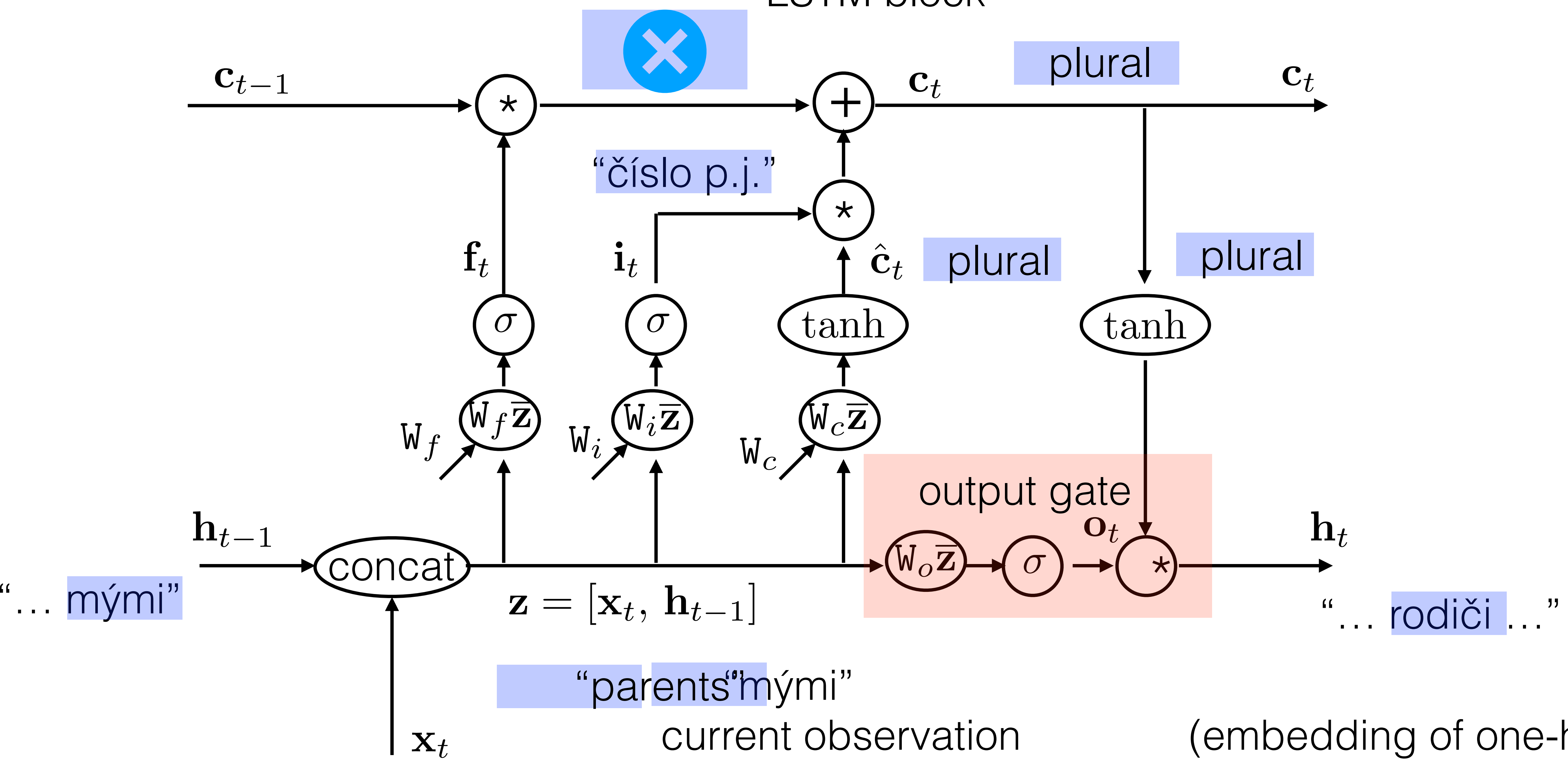
"I live with my parents, ..."

LSTM block



"I live with my parents, ..."

LSTM block



"I live with my parents, ..."

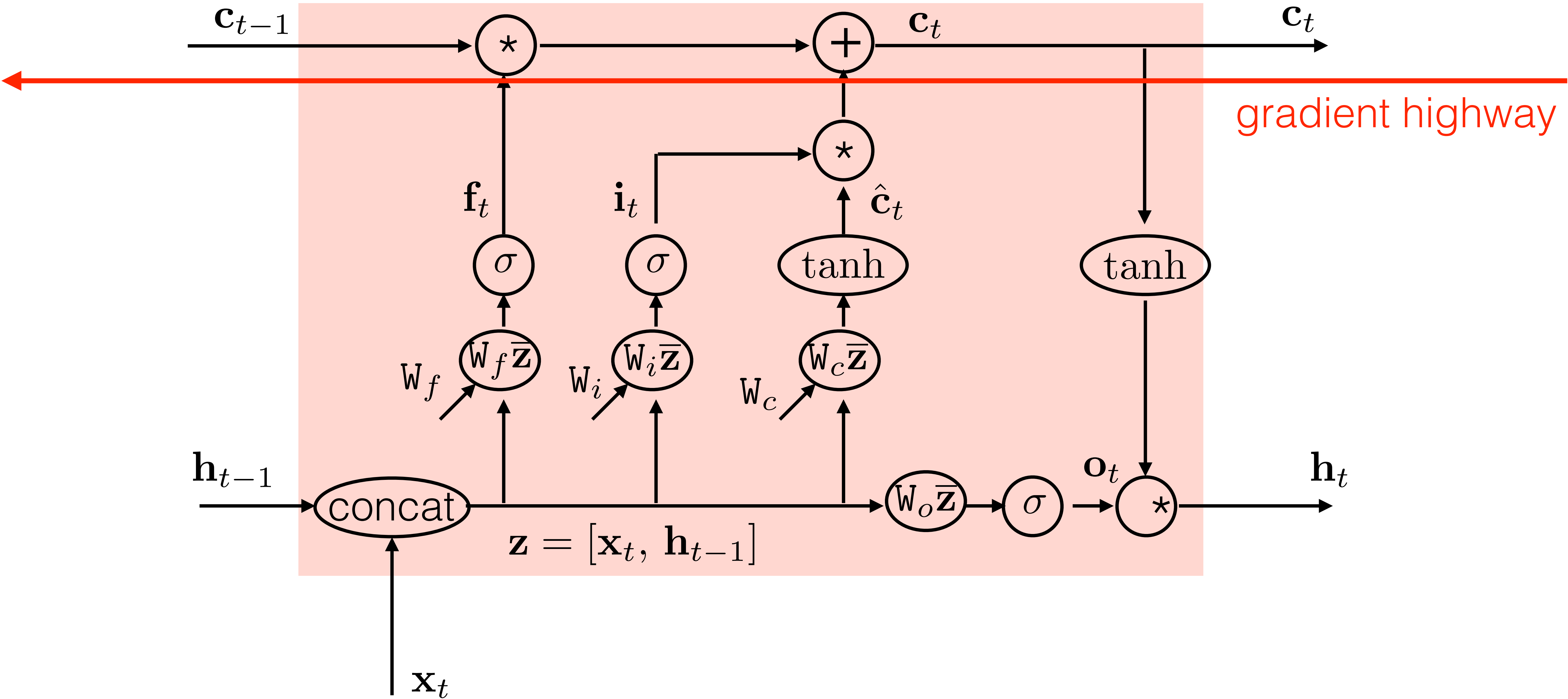
"... mými"

"parents" mými

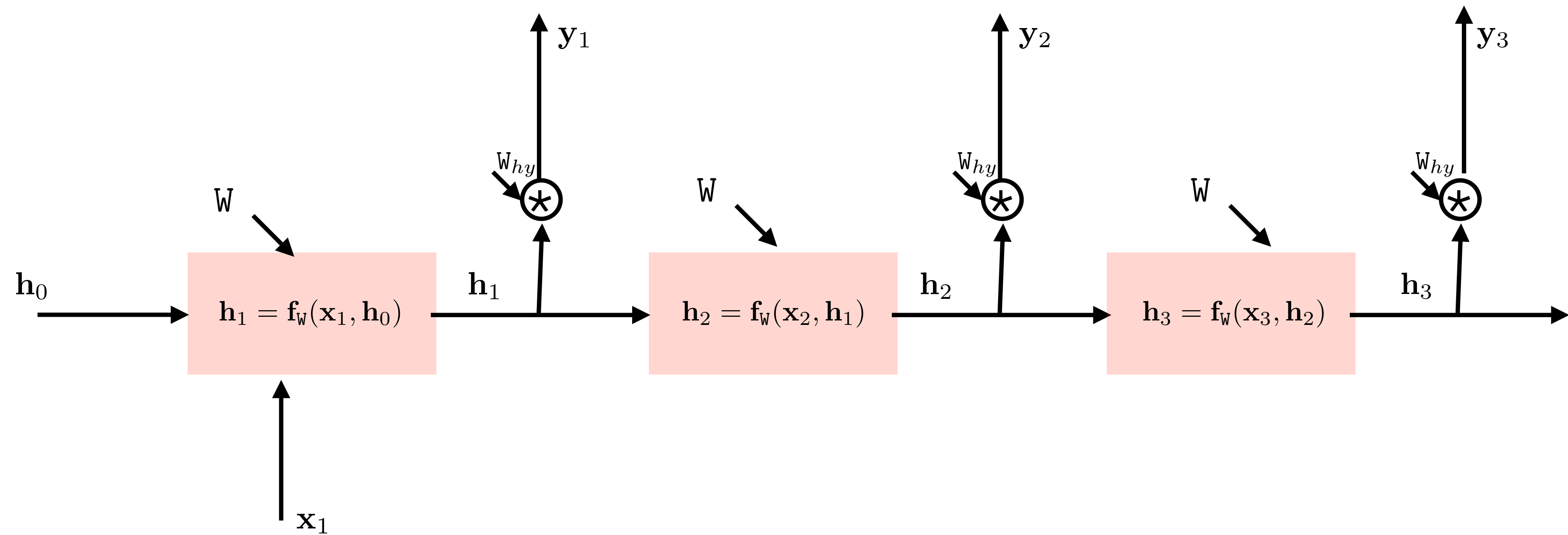
"... rodiči ..."

LSTM block

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

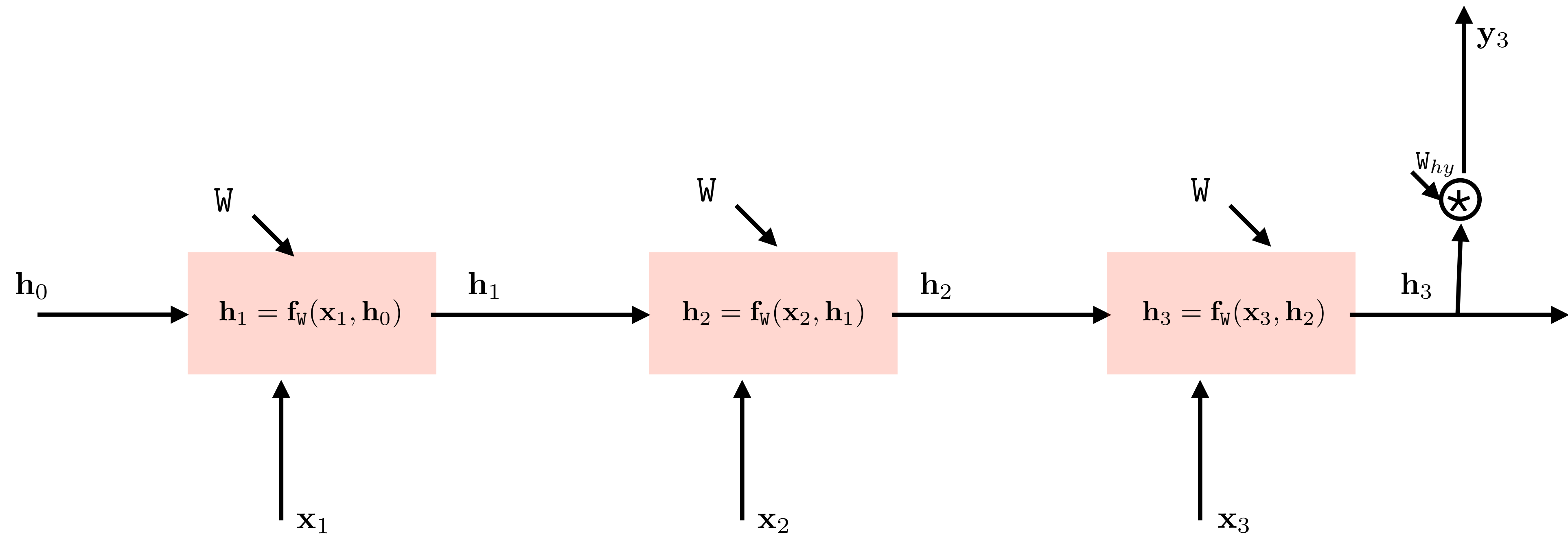


RNN architectures: one-to-many



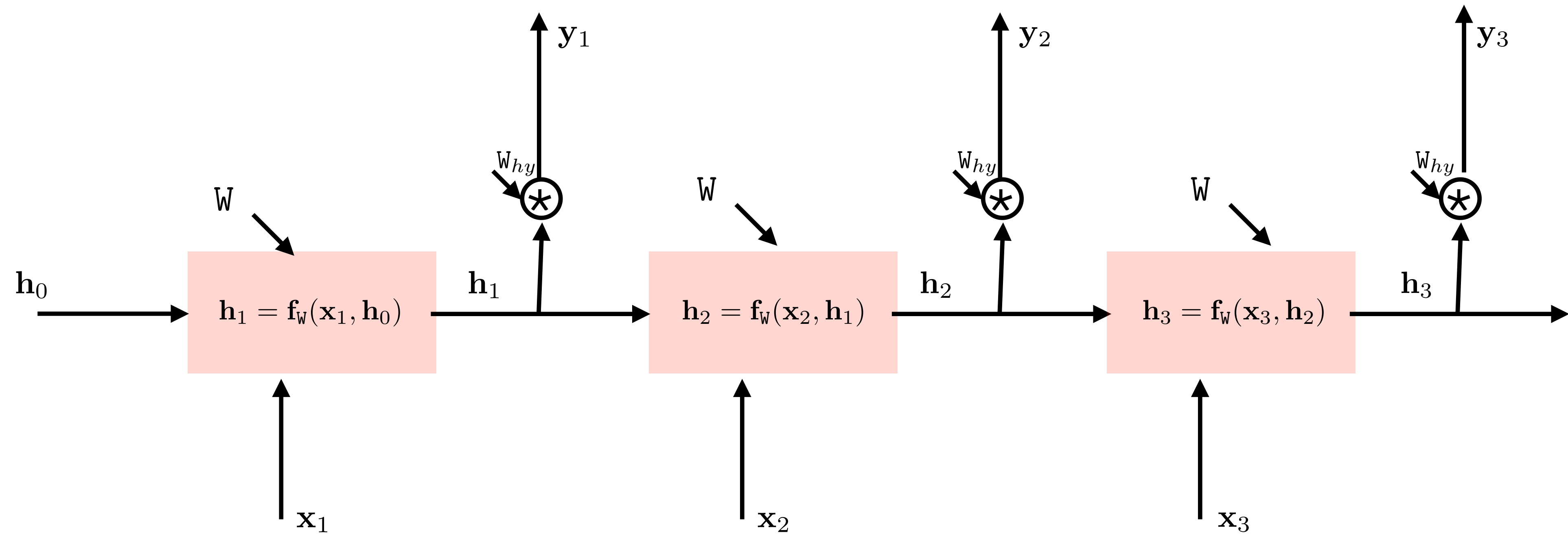
- Music generation
- Image captioning

RNN architectures: many-to-one



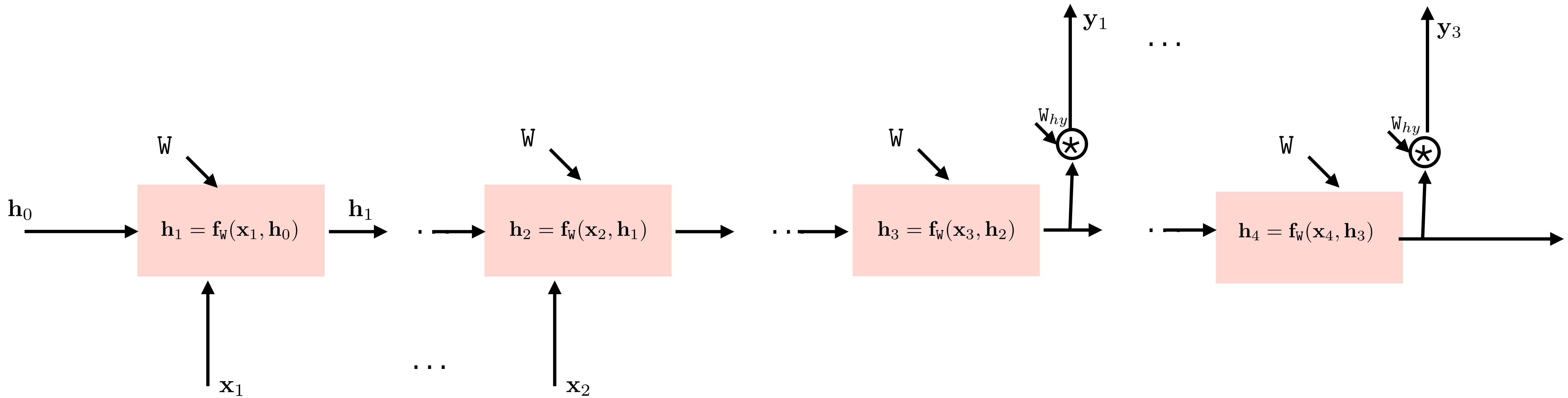
- Sentiment classification
- Action recognition

RNN architectures: many-to-many



- Named-entity recognition
- Speech recognition

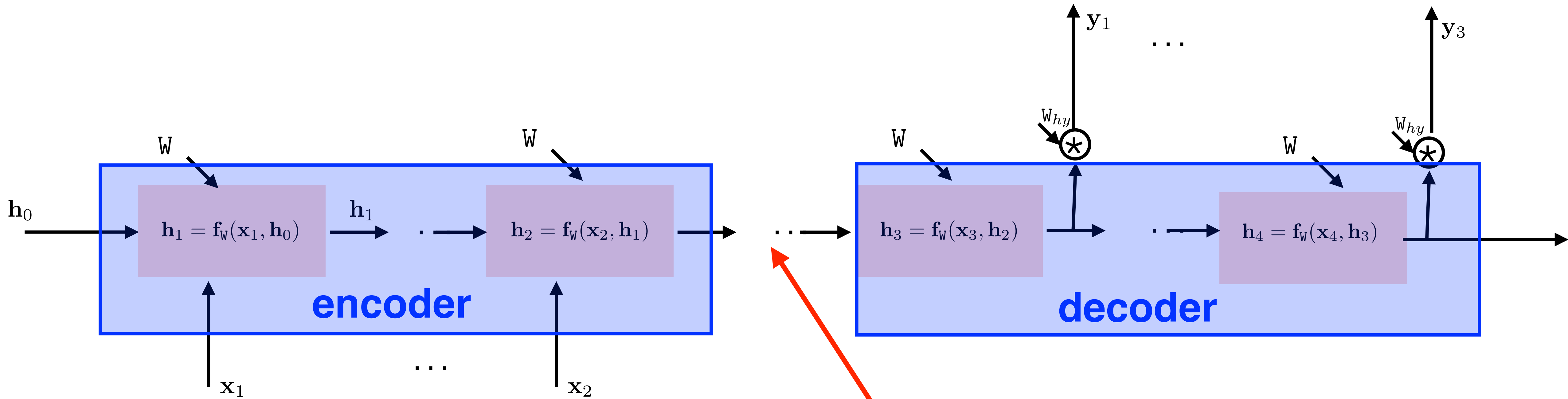
RNN architectures: many-to-many



- Machine translation
- Question answering

RNN architectures: many-to-many

output: variable-size sequence



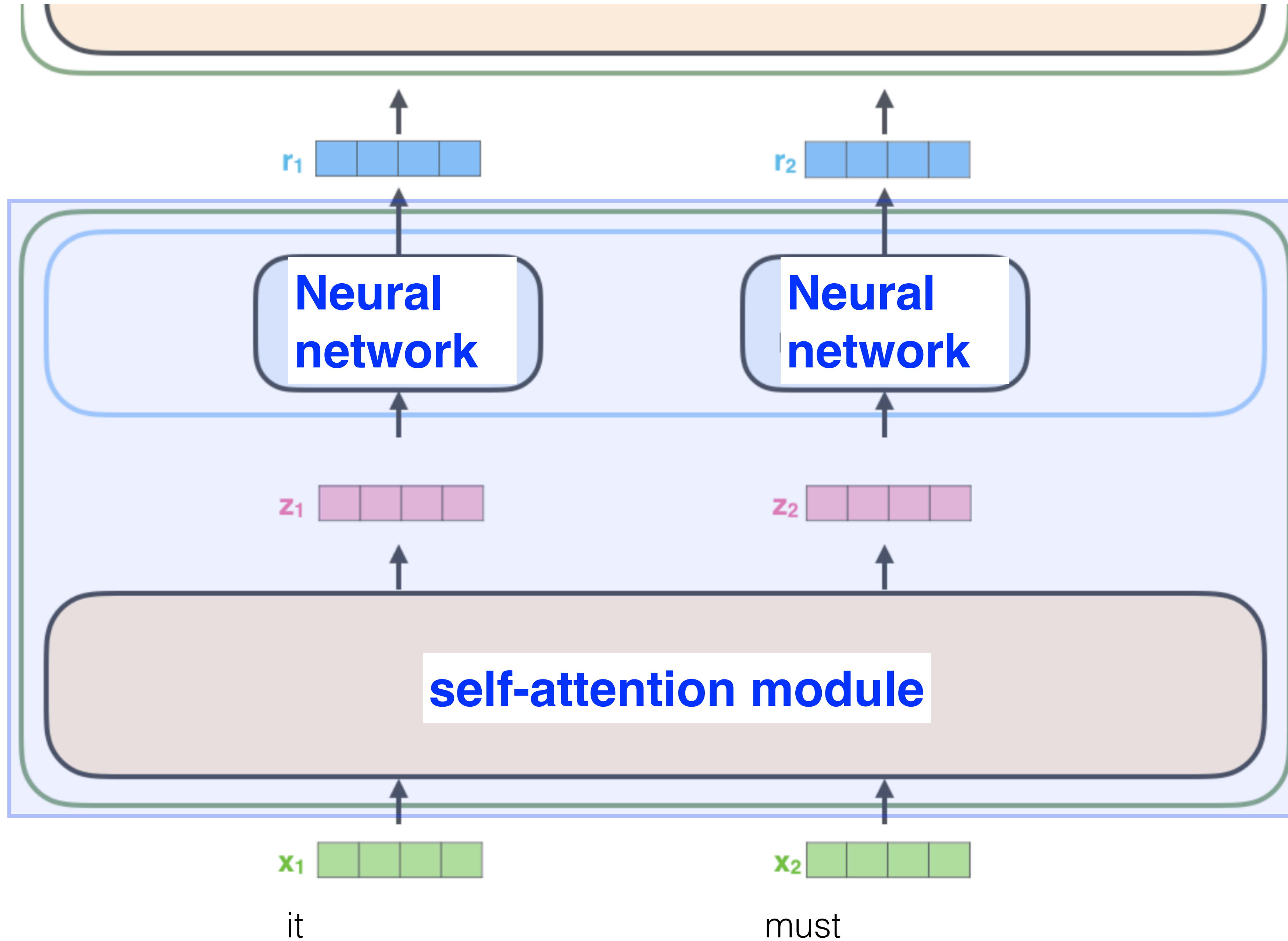
input: variable-size sequence

context: fixed-size semantic summary of input sequence

<https://papers.nips.cc/paper/7181-attention-is-all-you-need.pdf>
[Waswani, NIPS 2018]

encoder #2

encoder #1



self-attention module

Input

it

Embedding



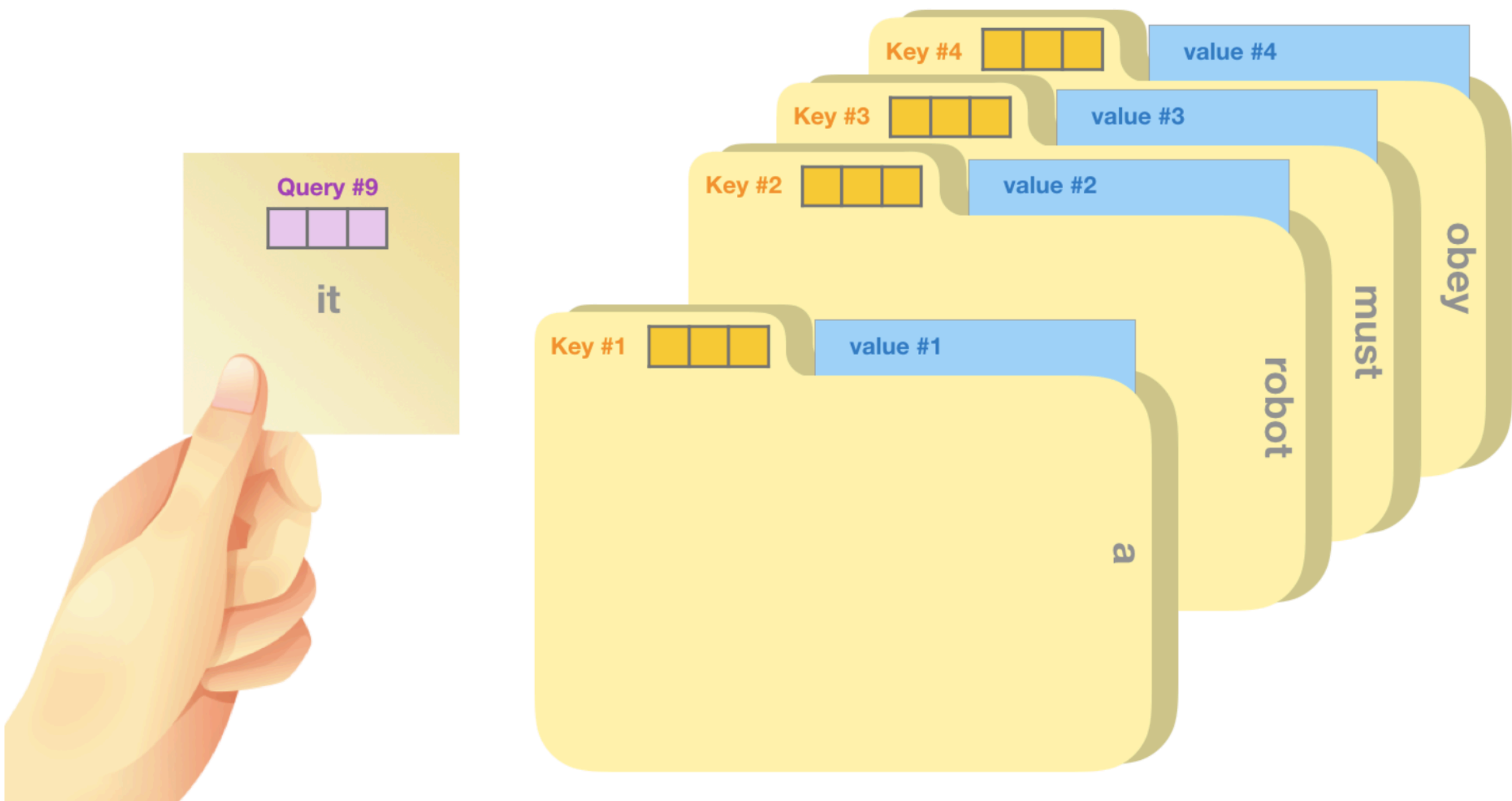
Queries



Keys



Values



self-attention module

Input

it

robot

must

Embedding



Queries

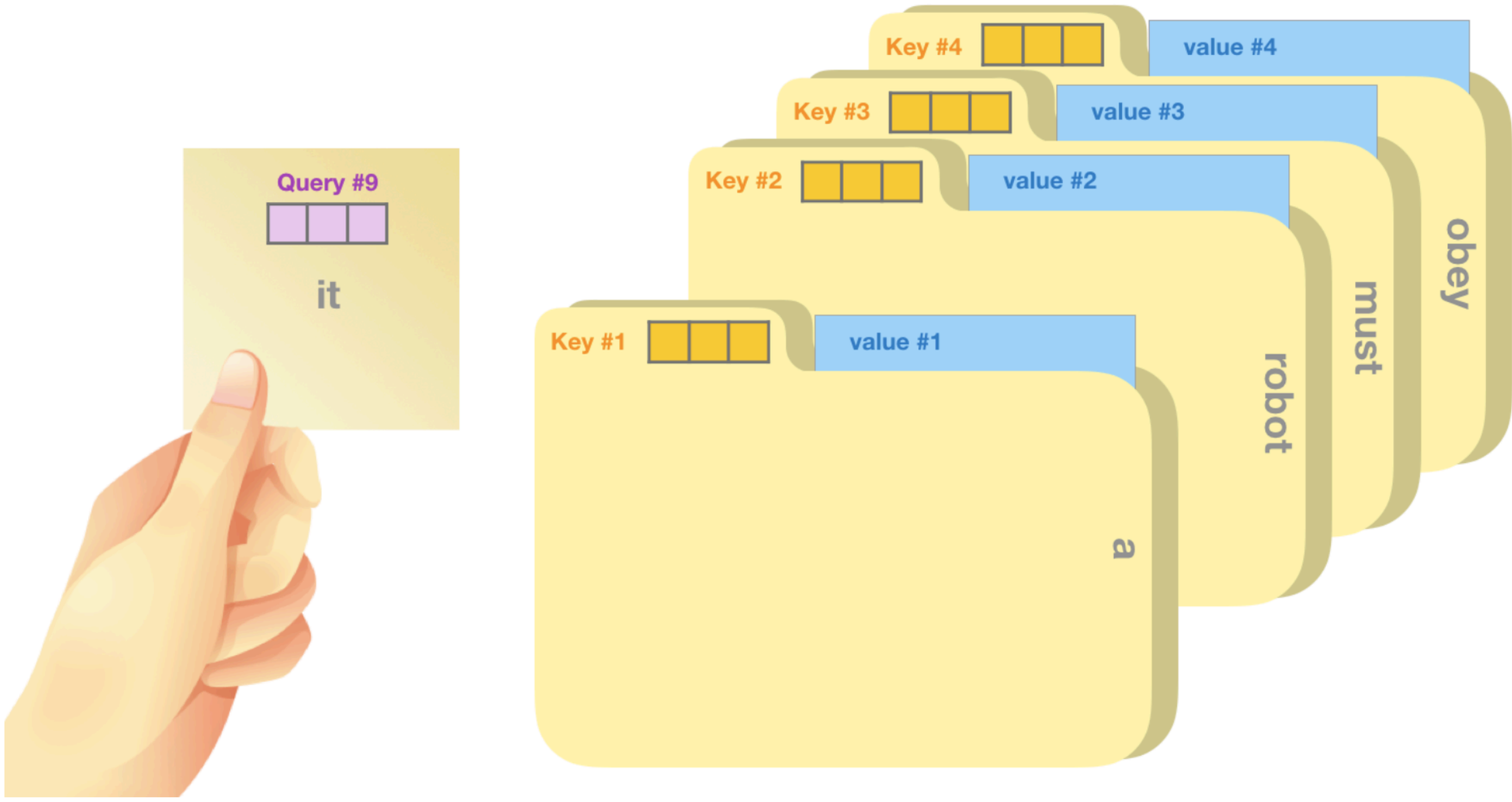


....

Keys



Values



self-attention module

Input

it

robot

Embedding



Queries



Keys



Values



Score

$q_1 \cdot k_1 = 72$

$q_1 \cdot k_2 = 120$

Divide by 8 ($\sqrt{d_k}$)

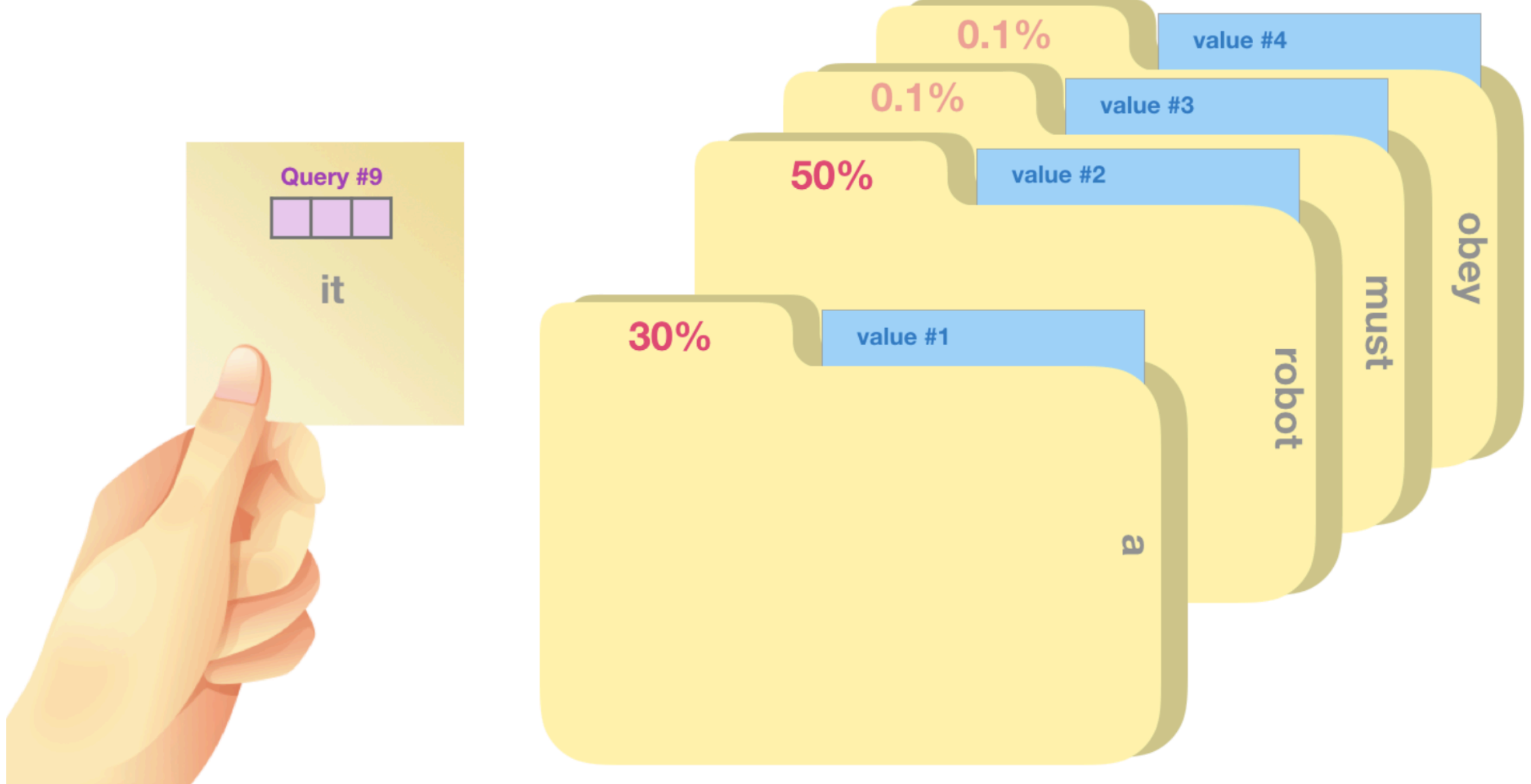
9

15

Softmax

0.3

0.5



self-attention module

Input

it

robot

Embedding



Queries



Keys



Values



Score

$q_1 \cdot k_1 = 72$

$q_1 \cdot k_2 = 120$

Divide by 8 ($\sqrt{d_k}$)

9

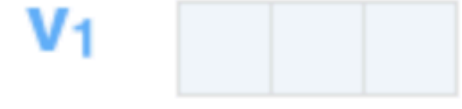
15

Softmax

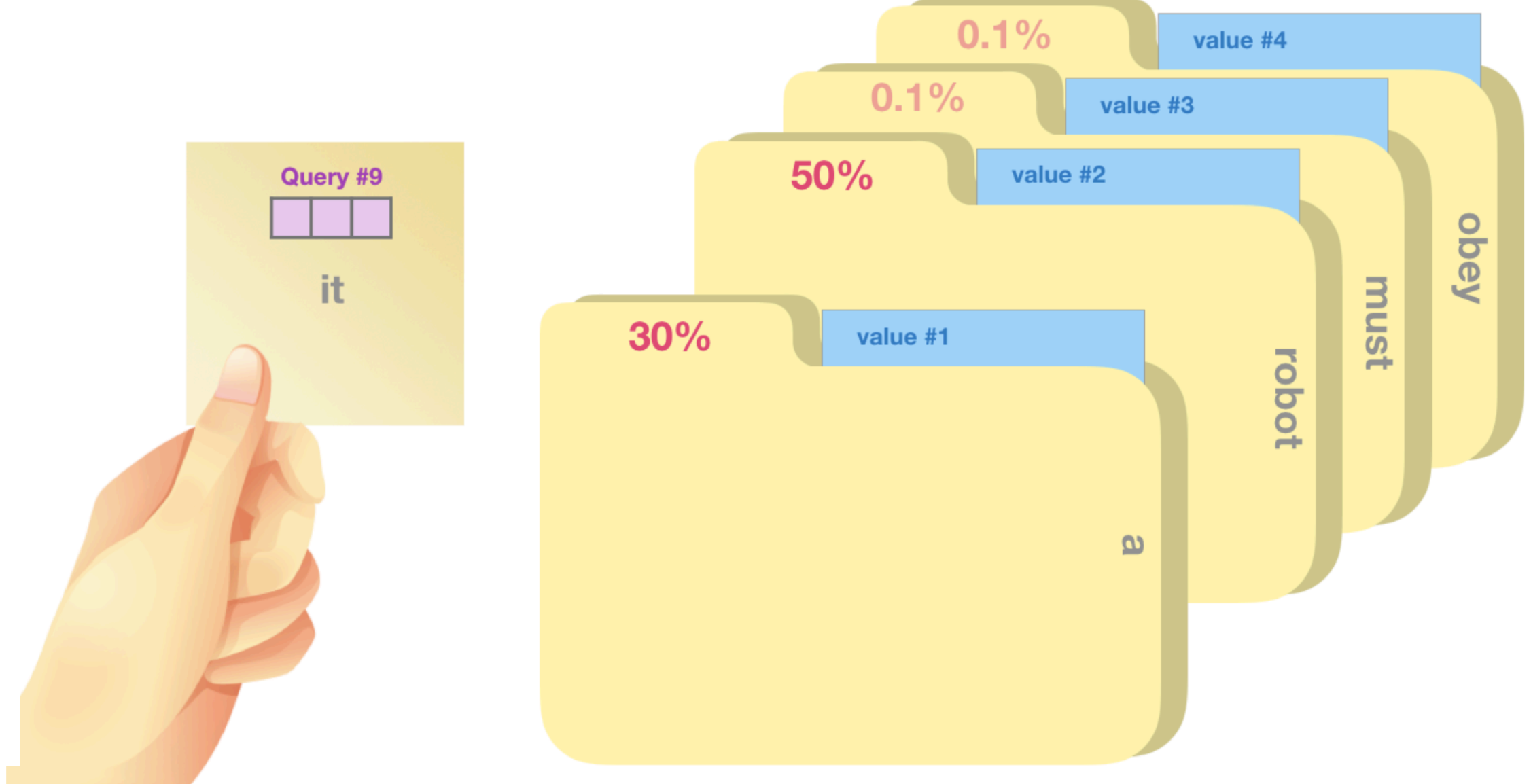
0.3

0.5

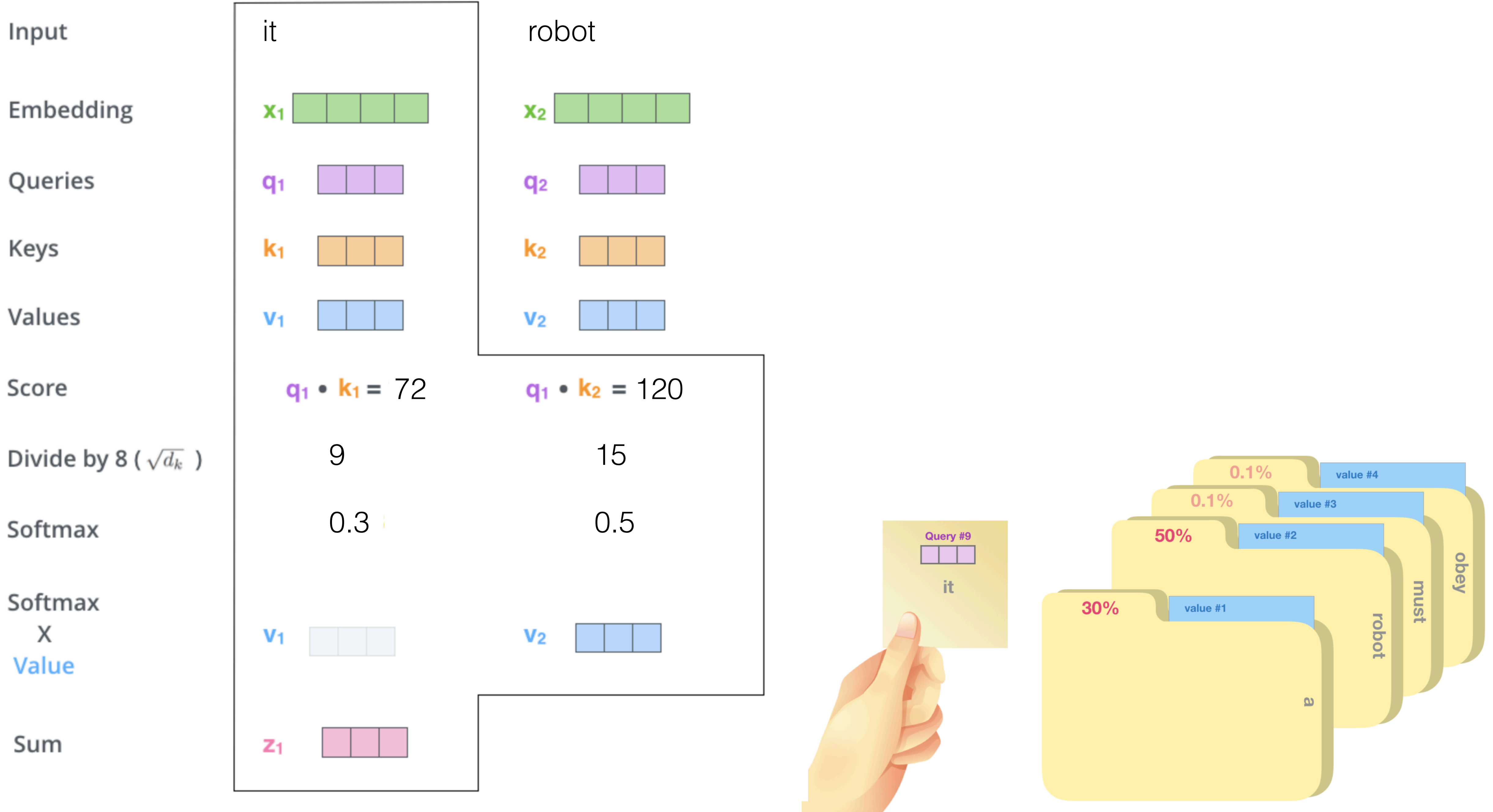
Softmax
X
Value



Sum



self-attention module



self-attention module

Input

it

robot

Embedding



Queries



Keys



Values



Score

$q_1 \cdot k_1 = 72$

$q_1 \cdot k_2 = 120$

Divide by 8 ($\sqrt{d_k}$)

9

15

Softmax

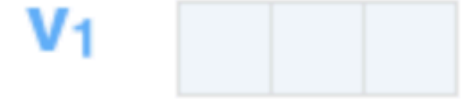
0.3

0.5

Softmax

X

Value



Sum

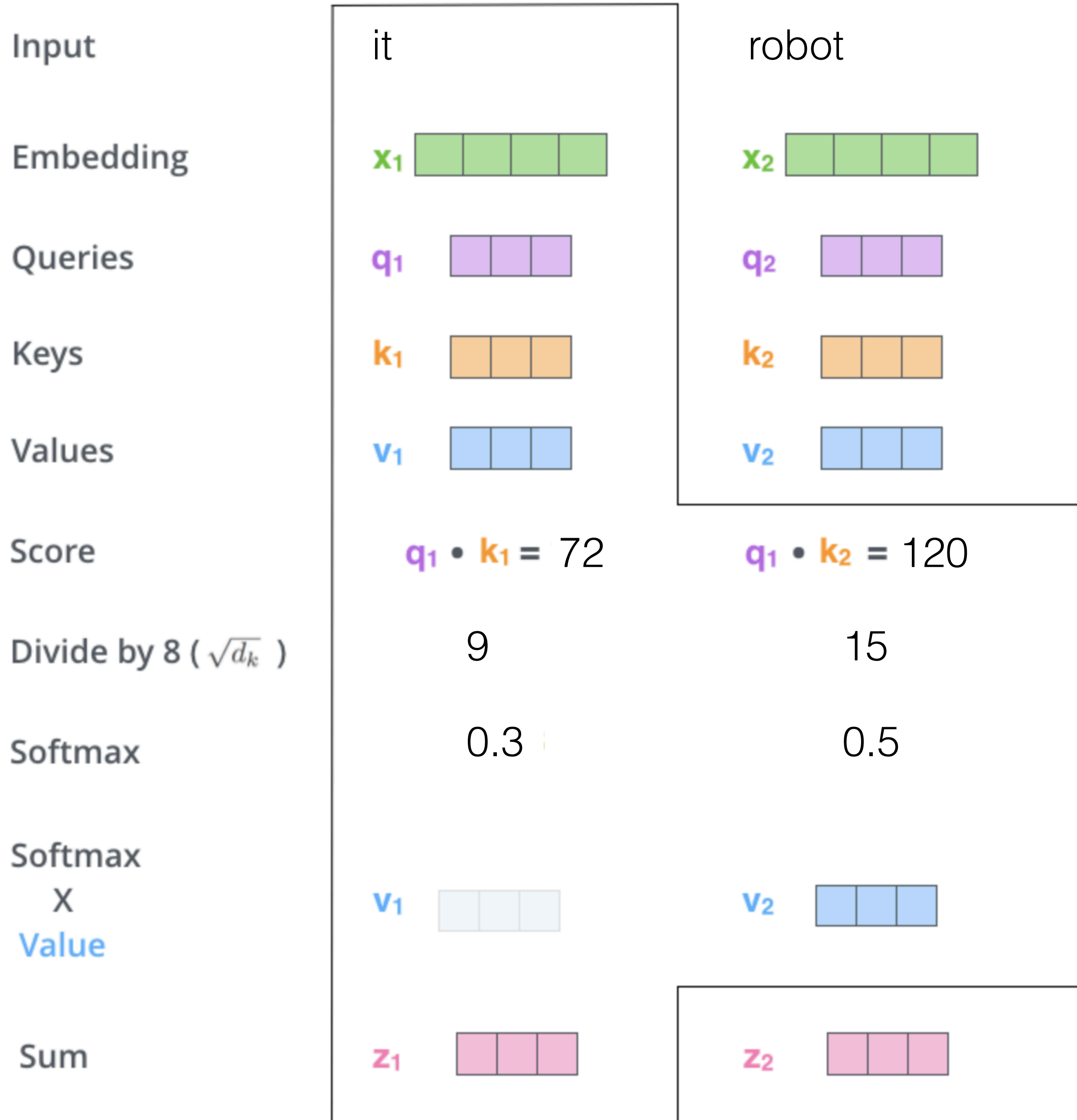


$$\text{softmax} \left(\frac{Q \times K^T}{\sqrt{d_k}} \right) V$$

= Z

output "z" is weighted

self-attention module



$$\text{softmax} \left(\frac{Q \times K^T}{\sqrt{d_k}} \right) V$$

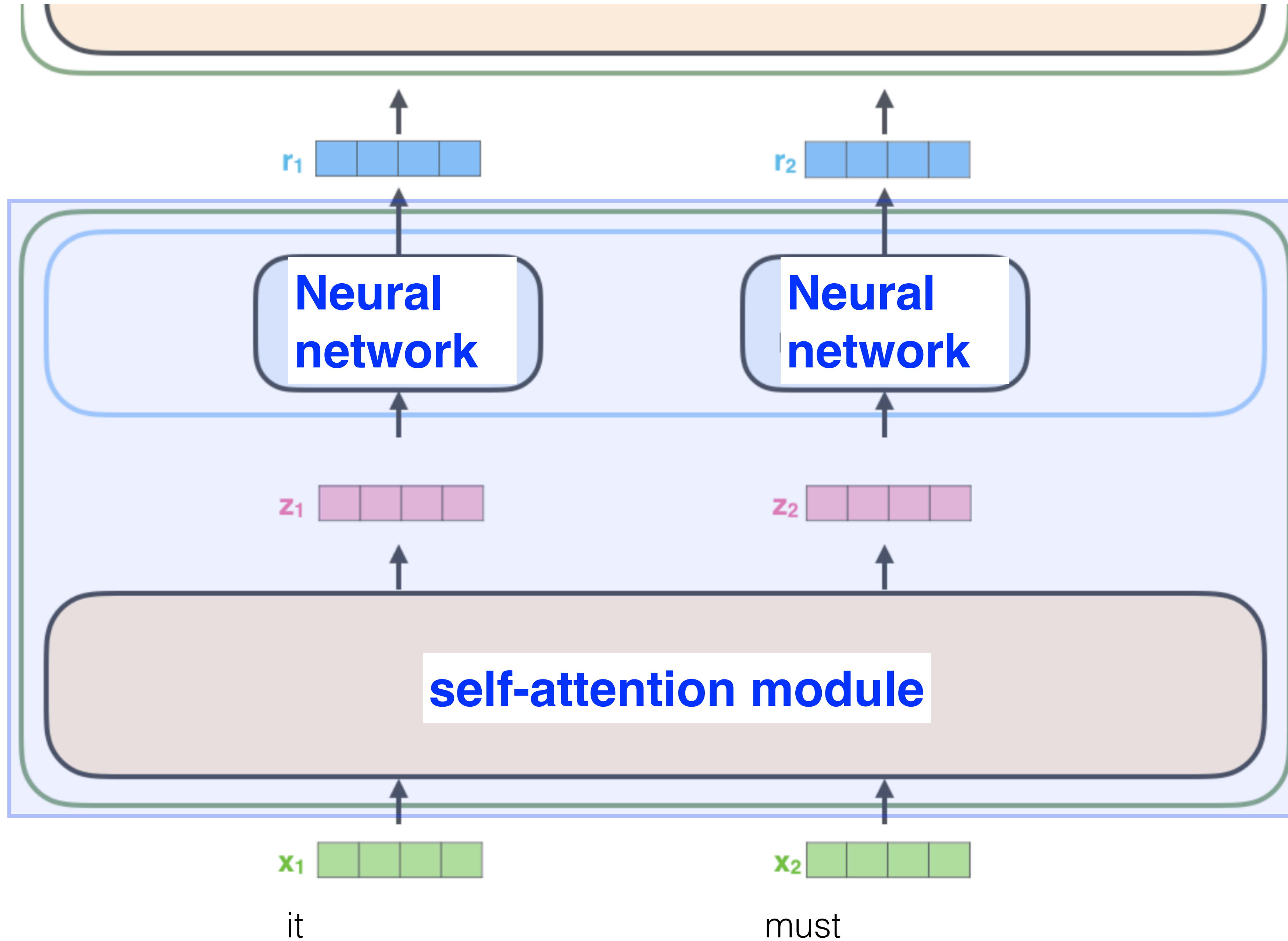
$$= Z$$

Proceed similarly for every single word.

<https://papers.nips.cc/paper/7181-attention-is-all-you-need.pdf>
[Waswani, NIPS 2018]

encoder #2

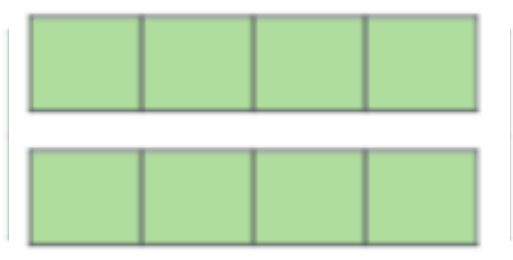
encoder #1



encoder

it
must

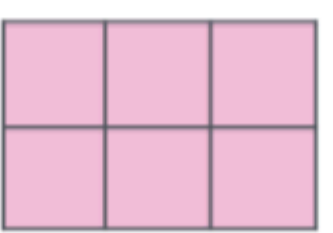
X



Calculating attention separately in
eight different attention heads

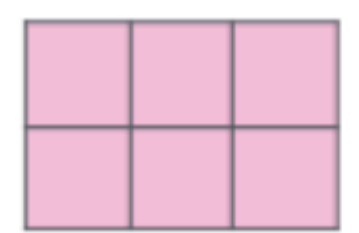
ATTENTION
HEAD #0

Z₀



ATTENTION
HEAD #1

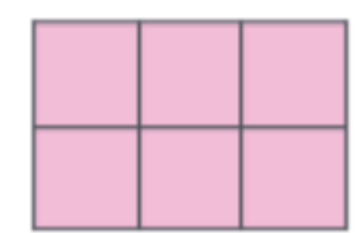
Z₁



...

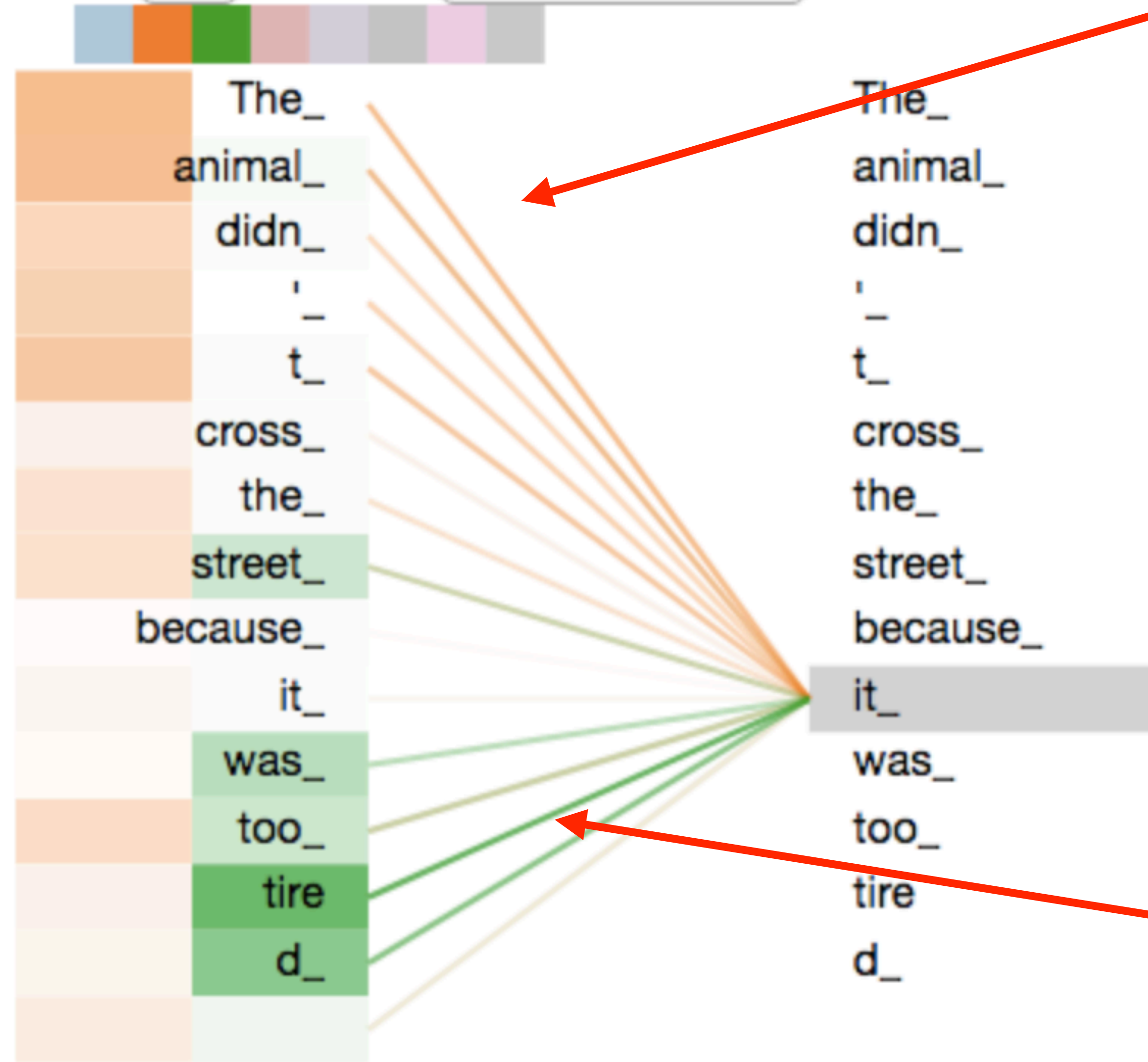
ATTENTION
HEAD #7

Z₇



Tensor2tensor:

Layer: 5 Attention: Input - Input



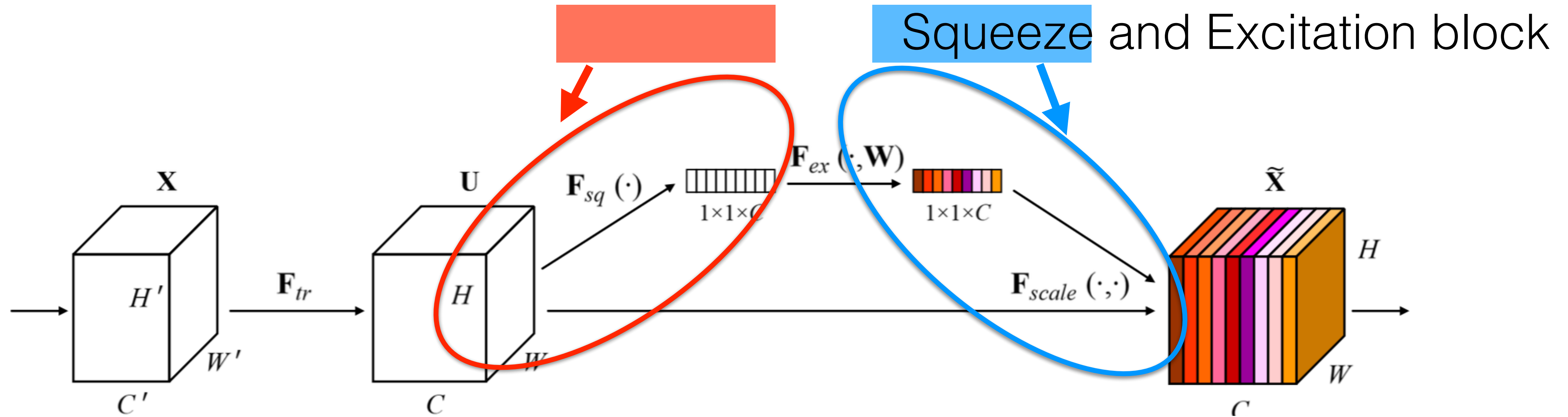
attention weights of **orange** attention head

attention weights of **green** attention head

https://colab.research.google.com/github/tensorflow/tensor2tensor/blob/master/tensor2tensor/notebooks/hello_t2t.ipynb

Attention in images

Channel attention



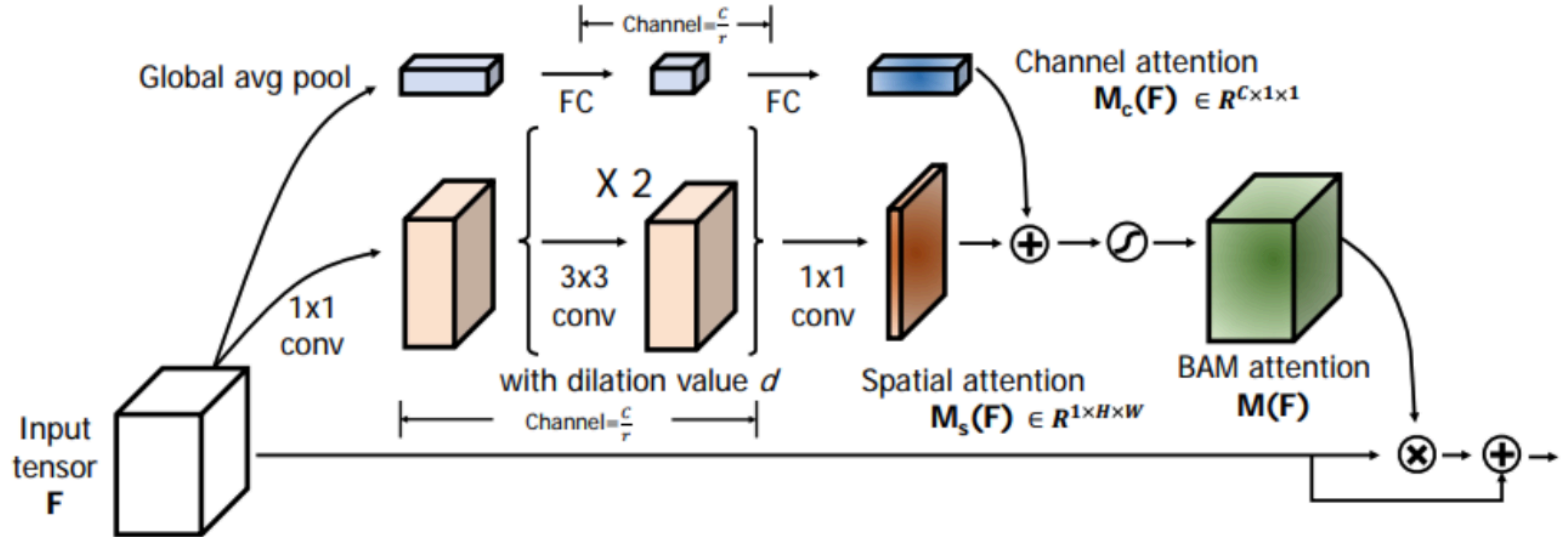
- Enhancement of ResNet, InceptionNet and DenseNet architectures by SE blocks consistently decrease error on ImageNet, COCO, ...

Squeeze and Excitation Networks [Hu et al, CVPR oral, 2017]

<https://arxiv.org/pdf/1709.01507.pdf>

Attention in images

Channel+spatial attention









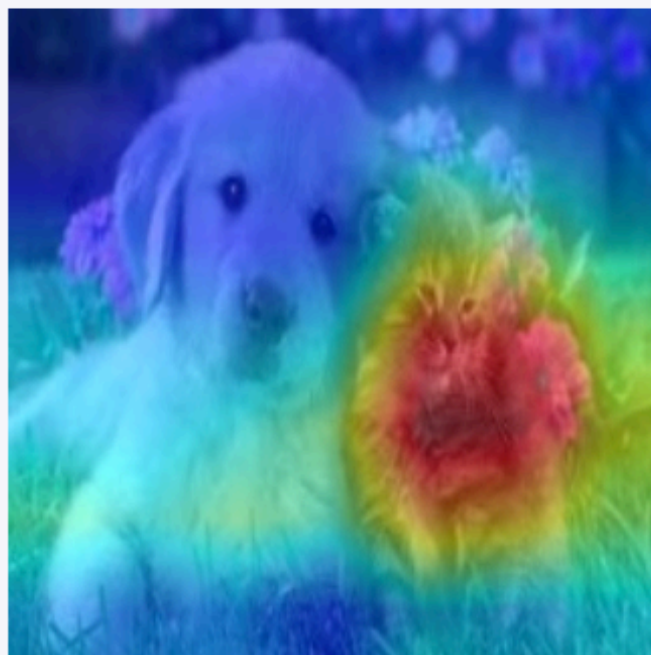

$$\mathbf{M}(\mathbf{F}) = \sigma(\mathbf{M}_c(\mathbf{F}) + \mathbf{M}_s(\mathbf{F})),$$

Attention modules [Woo et al, ECCV, 2018]
<https://arxiv.org/pdf/1807.06521v2.pdf>

GradCAM [Selvaraju et al, ICCV, 2018]

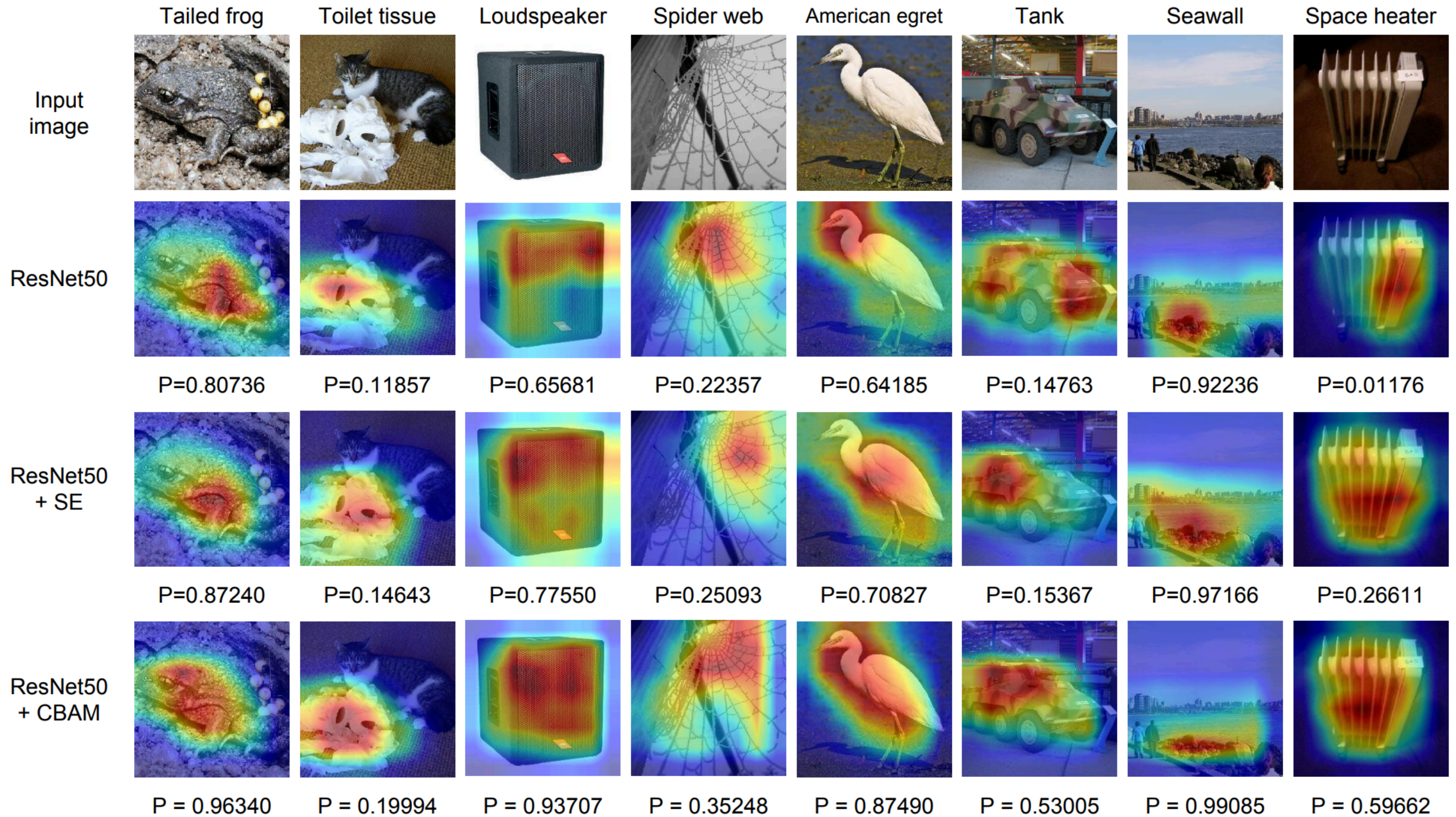
<https://github.com/jacobgil/pytorch-grad-cam>



| Category | Image | GradCAM | AblationCAM | ScoreCAM |
|----------|--|---|---|---|
| Dog |  |  |  |  |
| Cat |  |  |  |  |

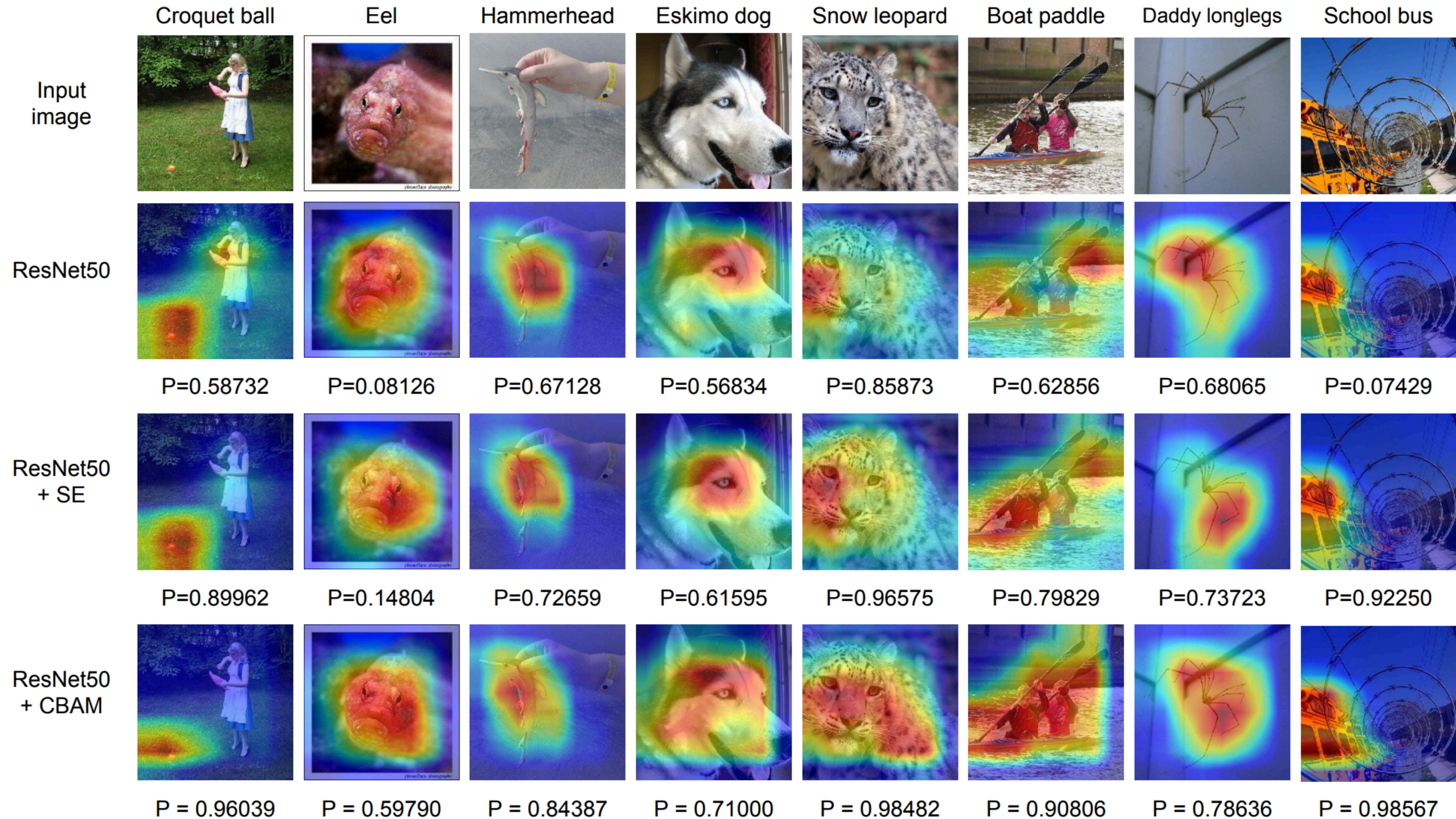
Attention modules [Woo et al, ECCV, 2018]

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

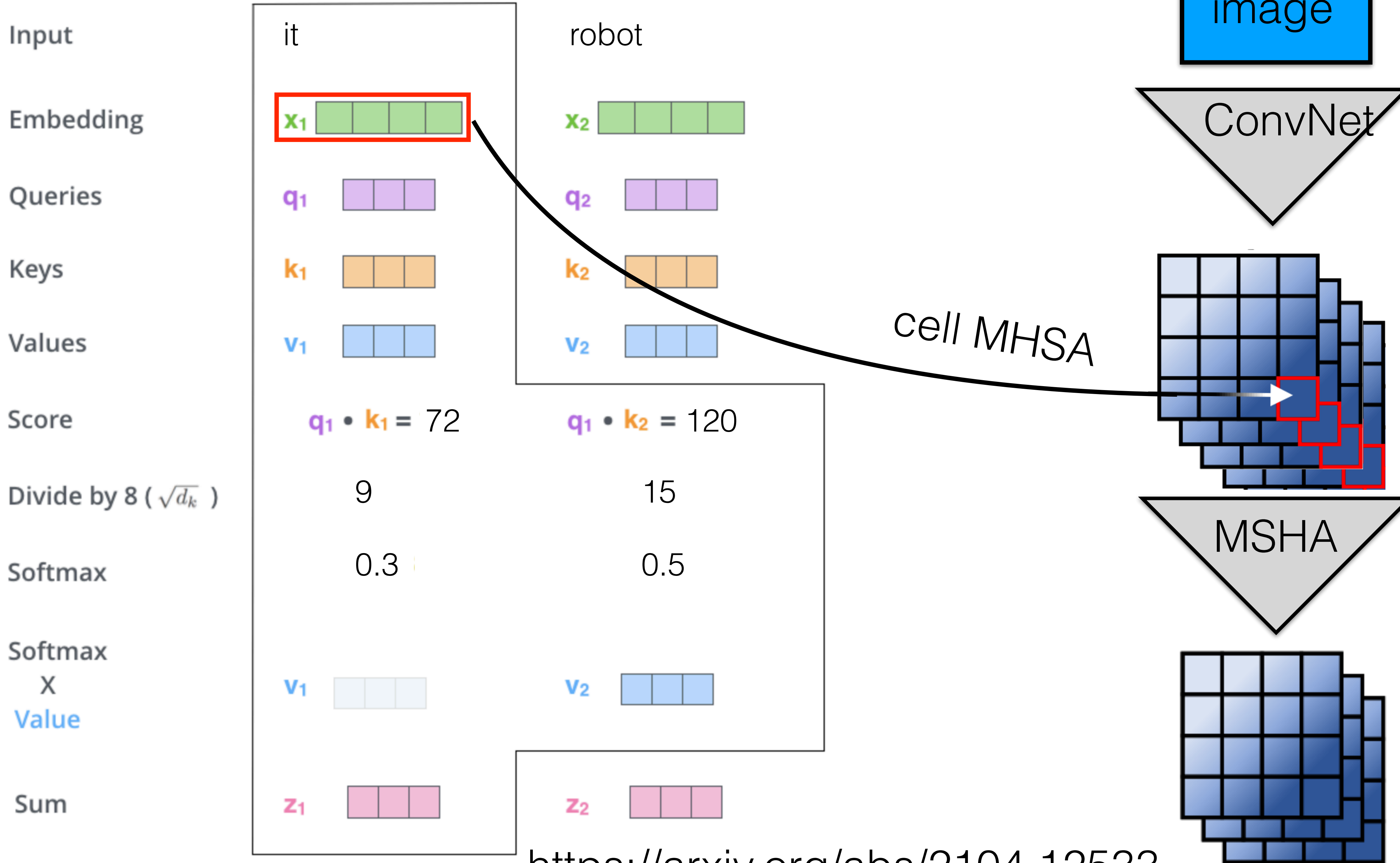


Attention modules [Woo et al, ECCV, 2018]

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

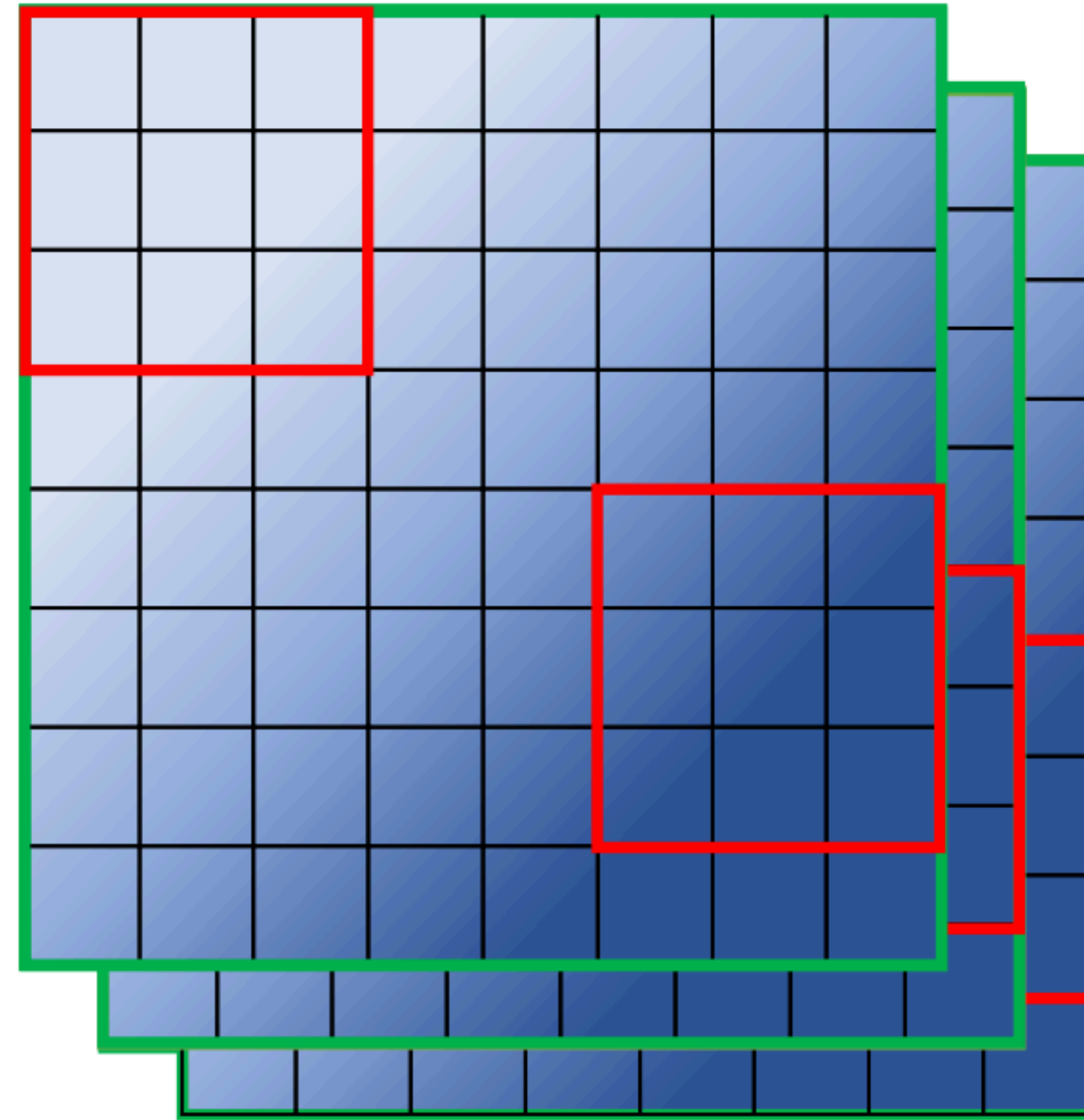


Visformer



<https://arxiv.org/abs/2104.12533>

Visformer



3x3 Convolution

Global Self-attention

Visformer

| Network | | ResNet-50 |
|----------------|---------------|-----------|
| FLOPs (G) | | 4.1 |
| Parameters (M) | | 25.6 |
| Full data | base setting | 77.43 |
| | elite setting | 78.73 |

Visformer

| Network | | ResNet-50 |
|----------------|---------------|-----------|
| FLOPs (G) | | 4.1 |
| Parameters (M) | | 25.6 |
| Full data | base setting | 77.43 |
| | elite setting | 78.73 |

ResNet has no significant benefit from introducing 20M augmentations due to its limited expressing power

Visformer

| Network | | ResNet-50 | Visformer-S |
|----------------|---------------|-----------|-------------|
| FLOPs (G) | | 4.1 | 4.9 |
| Parameters (M) | | 25.6 | 40.2 |
| Full data | base setting | 77.43 | 77.20 |
| | elite setting | 78.73 | 82.19 |

Visformer (ResNet with MHSA) exhibit improvement in elite setting

Visformer

| Network | | ResNet-50 | Visformer-S | DeiT-S |
|----------------|---------------|-----------|-------------|--------|
| FLOPs (G) | | 4.1 | 4.9 | 4.6 |
| Parameters (M) | | 25.6 | 40.2 | 21.8 |
| Full data | base setting | 77.43 | 77.20 | 63.12 |
| | elite setting | 78.73 | 82.19 | 80.07 |

Direct usage of “language transformers” without convolutions is more prone to overfitting

Attention in RL



0000
1 200

Q-value function



0000
1 200

Advantage function

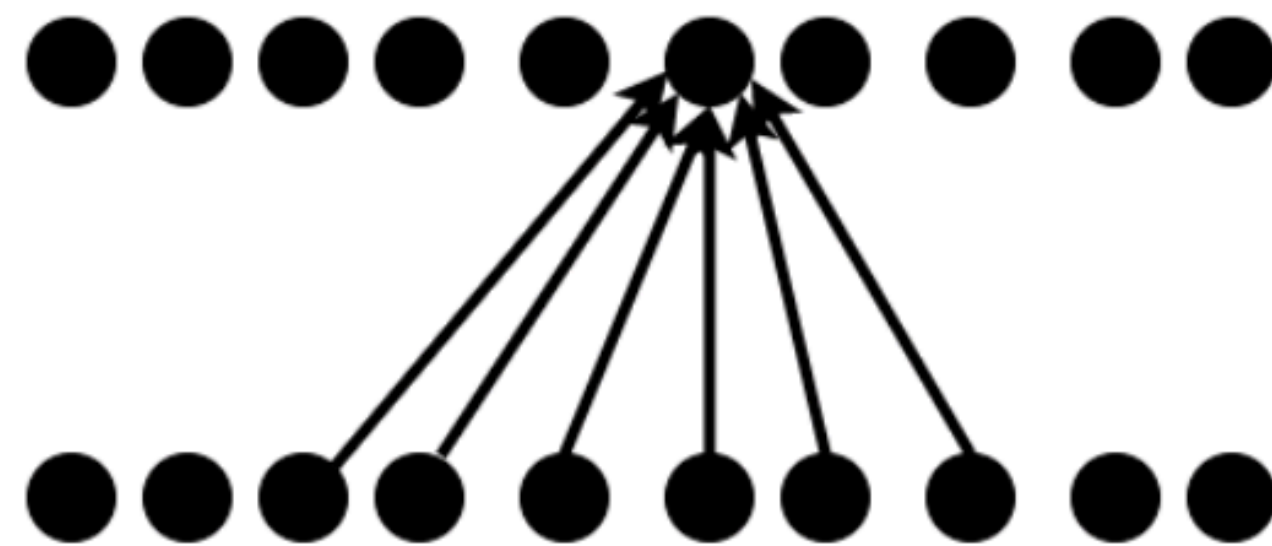


Summary

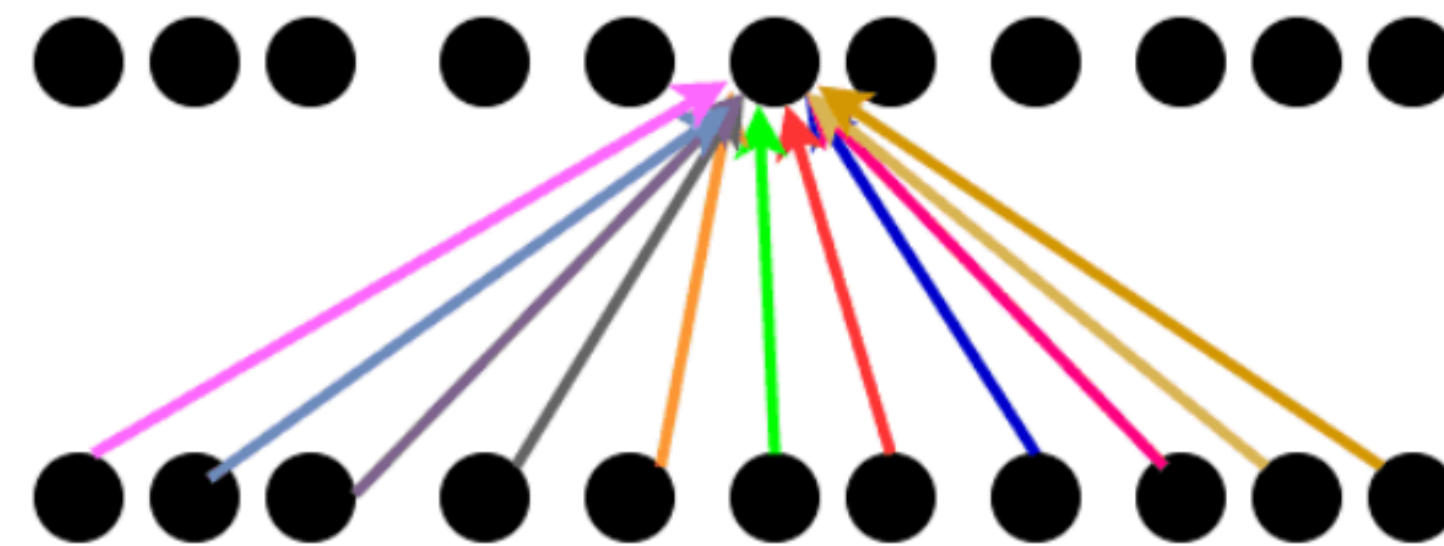
- self-attention overfits (requires large dataset)
- memory is attention through time [Alex Graves 2020]
- pyTorch library: <https://github.com/The-AI-Summer/self-attention-cv>

```
model = MultiHeadSelfAttention(dim=64)
```

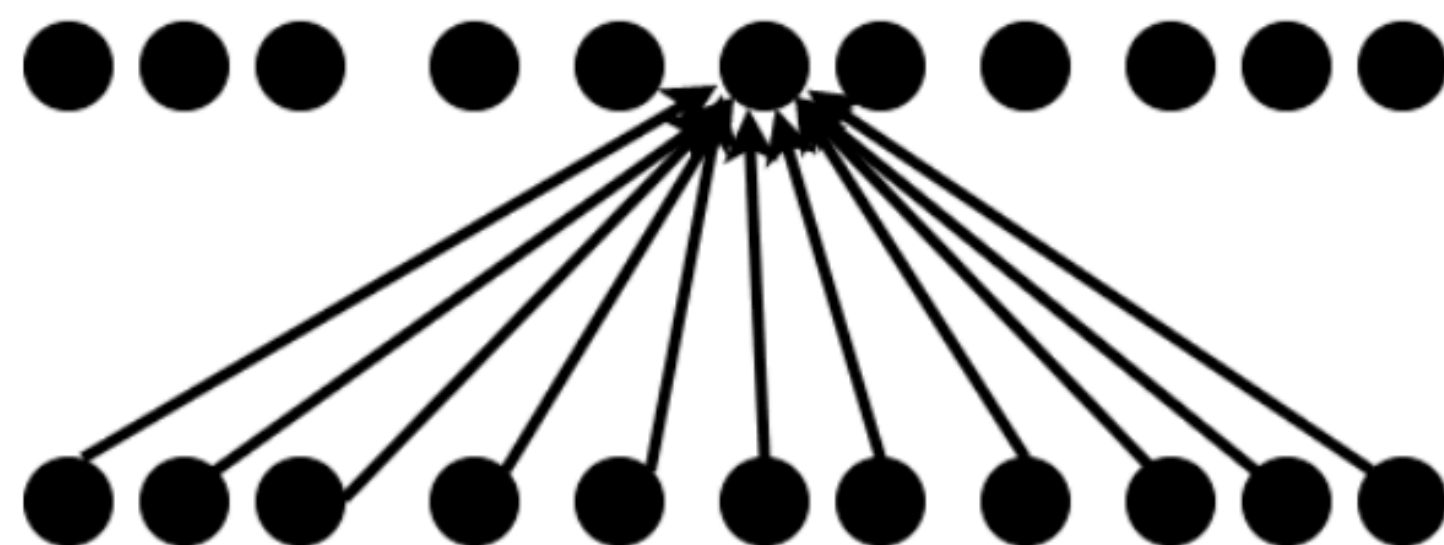
Convolution



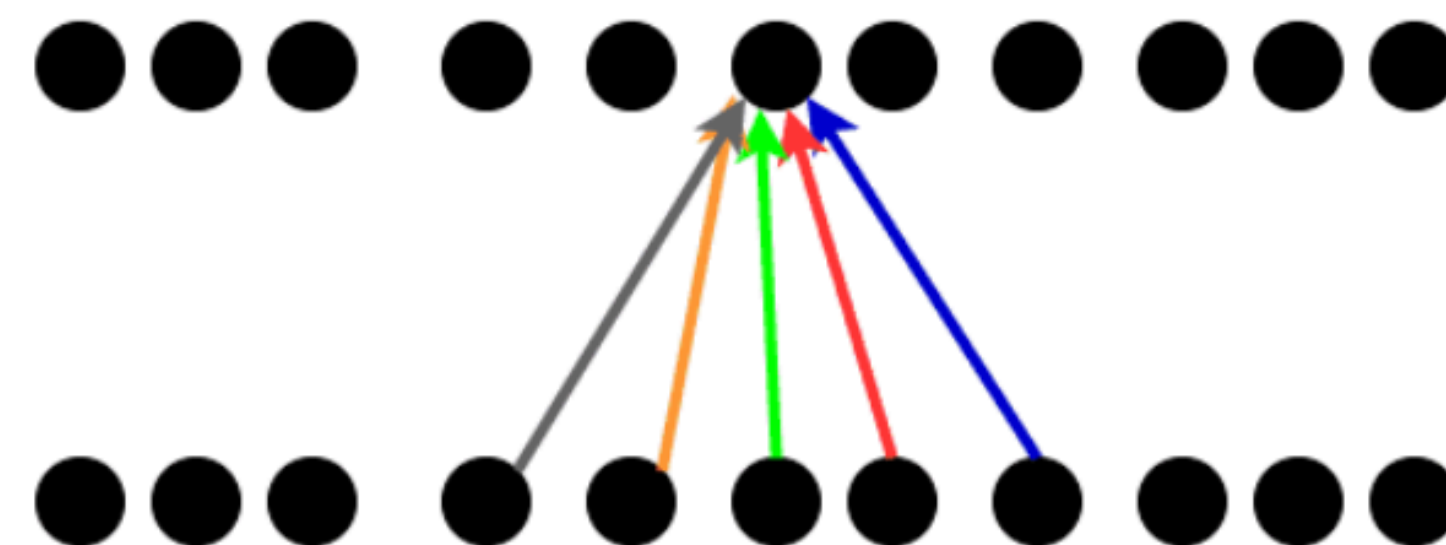
Global attention



Fully Connected layer



Local attention

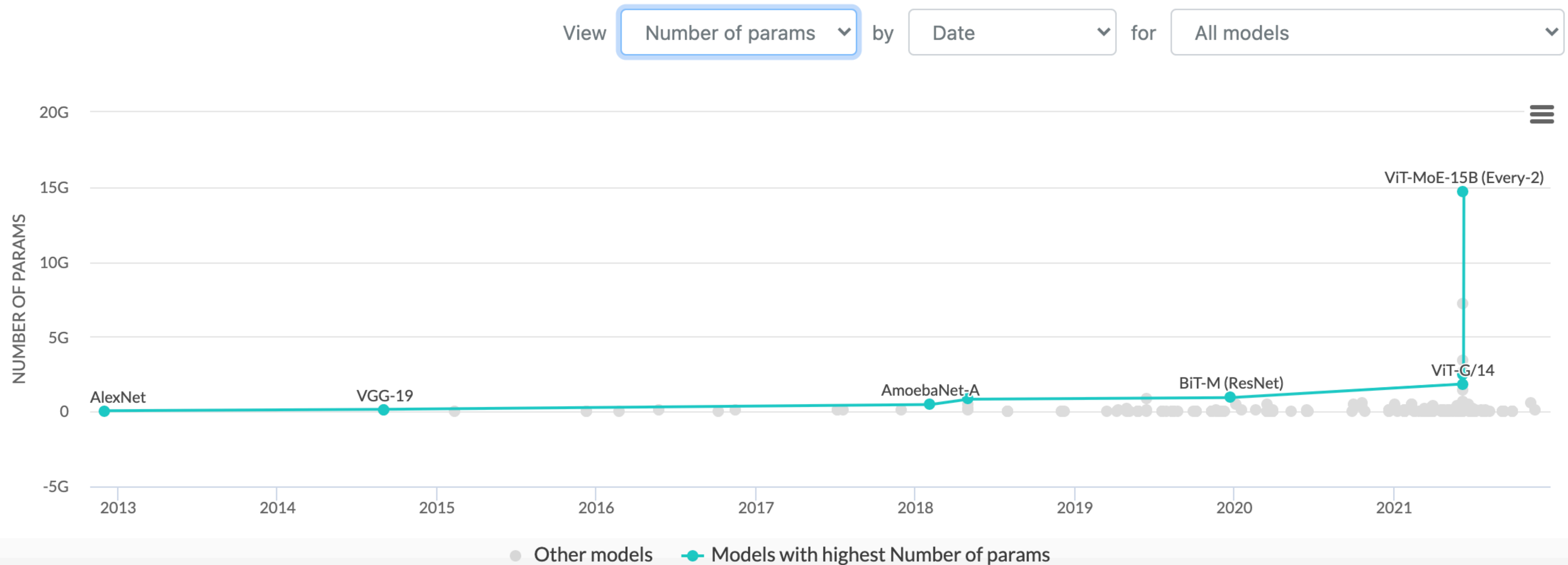


Future of robotics?

- 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 unique successes than exponentially growing start general AI

Future of robotics?

- Moore law:
 - number of tranzistors in integrated circuit will double every year
- Observation:
 - size of models (learning time => number of GPUs+memory) grows explonentially



Future of robotics?

- Moore law:
 - number of tranzistors in integrated circuit will double every year
- Observation:
 - size of models (learning time => number of GPUs+memory) grows explonentially
- What we get from it?
 - Improvement in accuracy of AI models slows down

