

MML minor #6

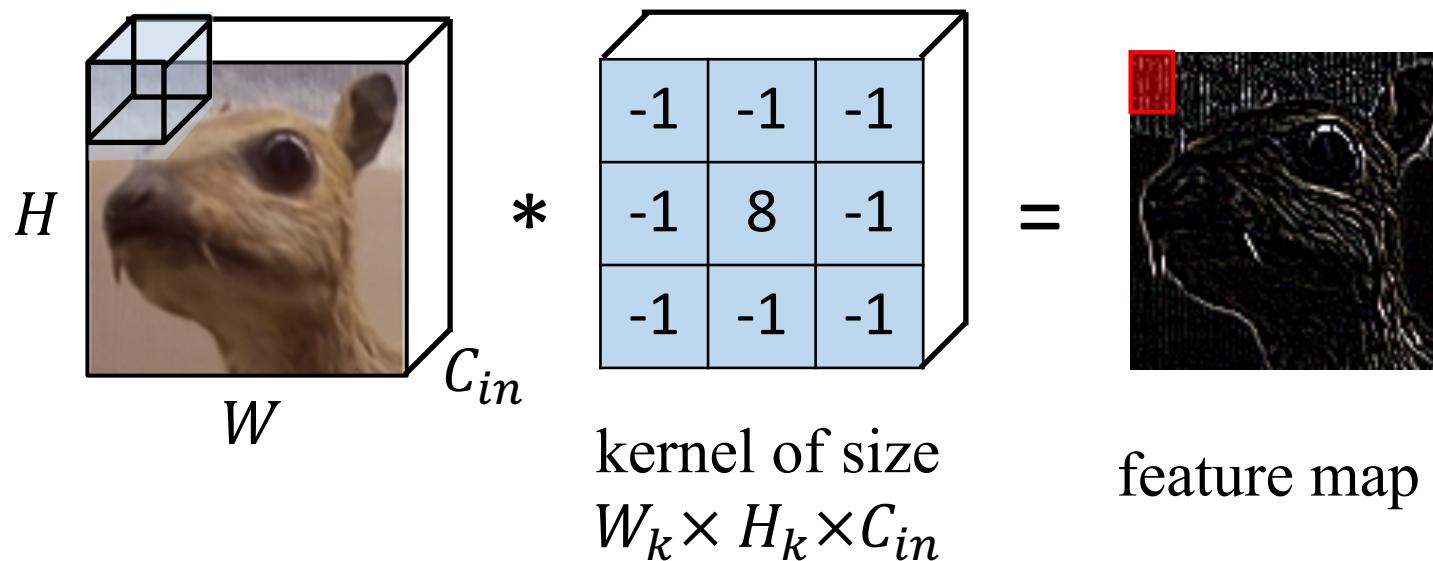
Свёрточные (*продолжение*) и
рекуррентные нейронные сети

A color image input

- Let's say we have a color image as an input, which is $W \times H \times C_{in}$ **tensor** (multidimensional array), where
 - W – is an image width,
 - H – is an image height,
 - C_{in} – is a number of input channels (e.g. 3 **RGB** channels).

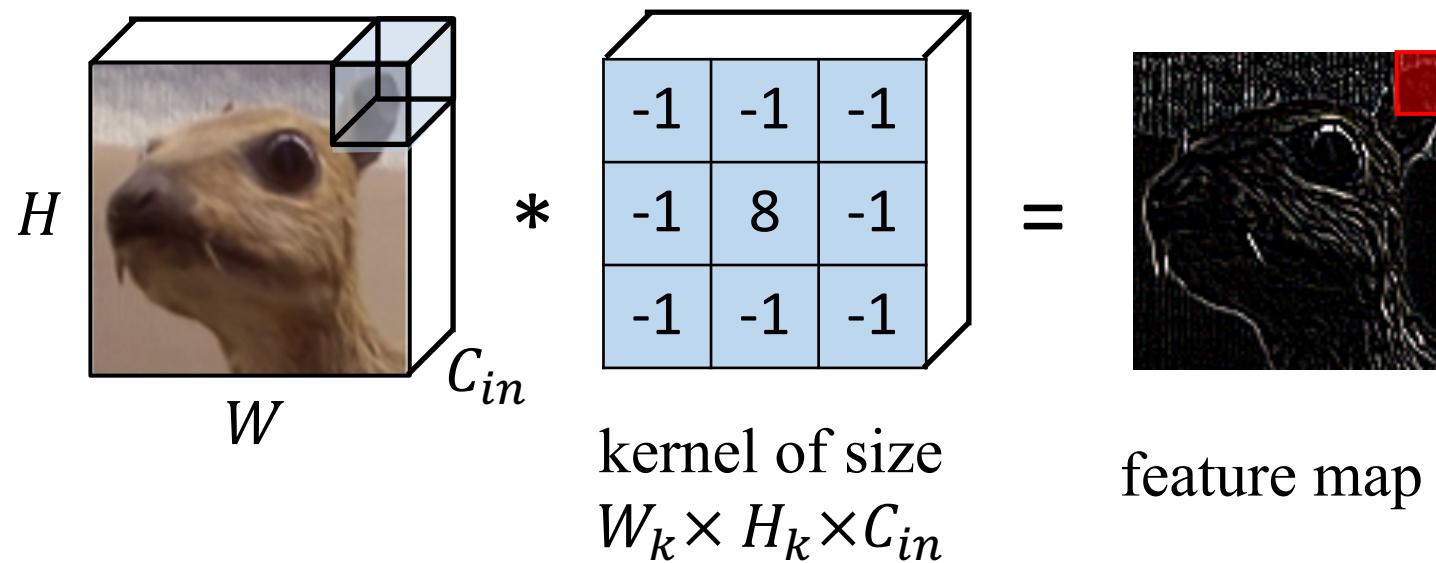
A color image input

- Let's say we have a color image as an input, which is $W \times H \times C_{in}$ **tensor** (multidimensional array), where
 - W – is an image width,
 - H – is an image height,
 - C_{in} – is a number of input channels (e.g. 3 **RGB** channels).



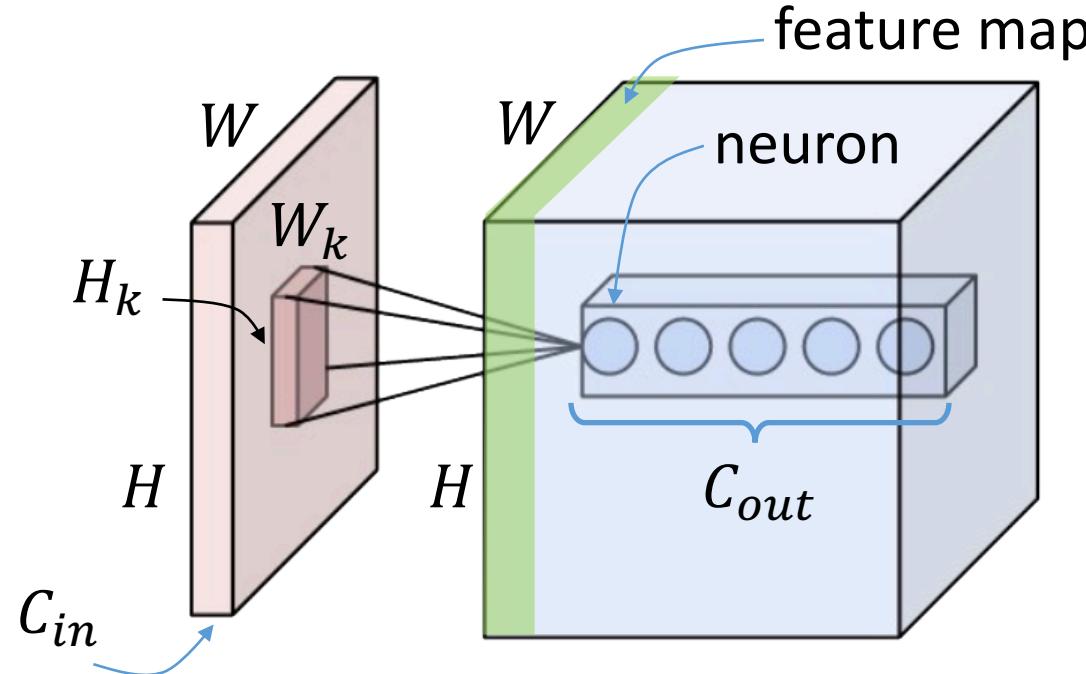
A color image input

- Let's say we have a color image as an input, which is $W \times H \times C_{in}$ **tensor** (multidimensional array), where
 - W – is an image width,
 - H – is an image height,
 - C_{in} – is a number of input channels (e.g. 3 **RGB** channels).



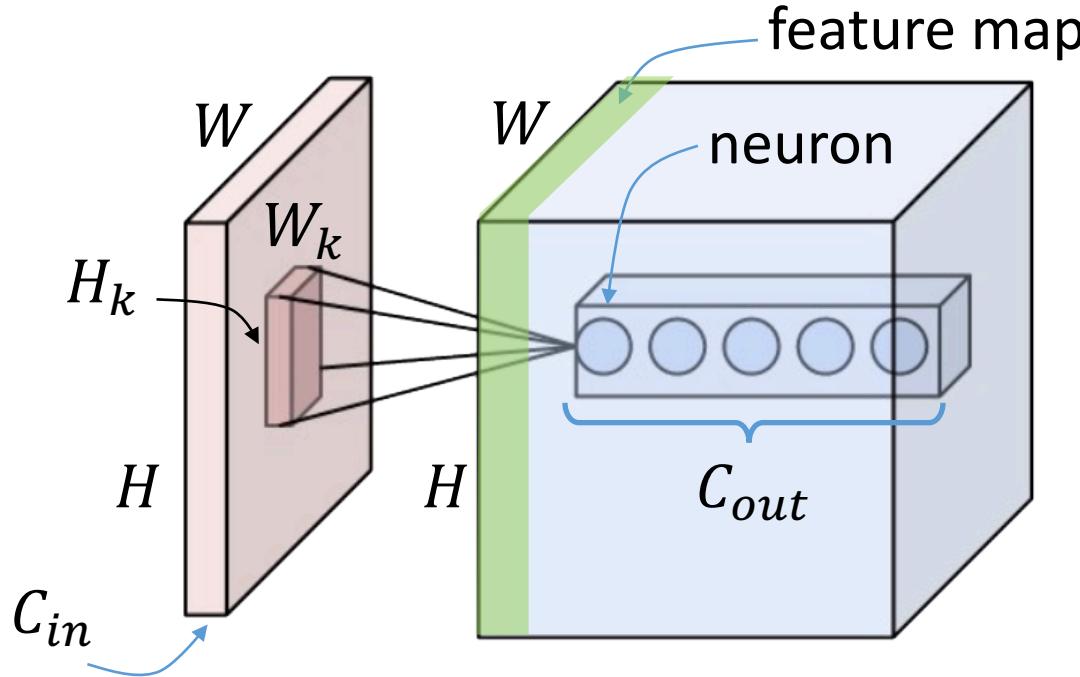
One kernel is not enough!

- We want to train C_{out} kernels of size $W_k \times H_k \times C_{in}$.
- Having a stride of 1 and enough zero padding we can have $W \times H \times C_{out}$ output neurons.



One kernel is not enough!

- We want to train C_{out} kernels of size $W_k \times H_k \times C_{in}$.
- Having a stride of 1 and enough zero padding we can have $W \times H \times C_{out}$ output neurons.



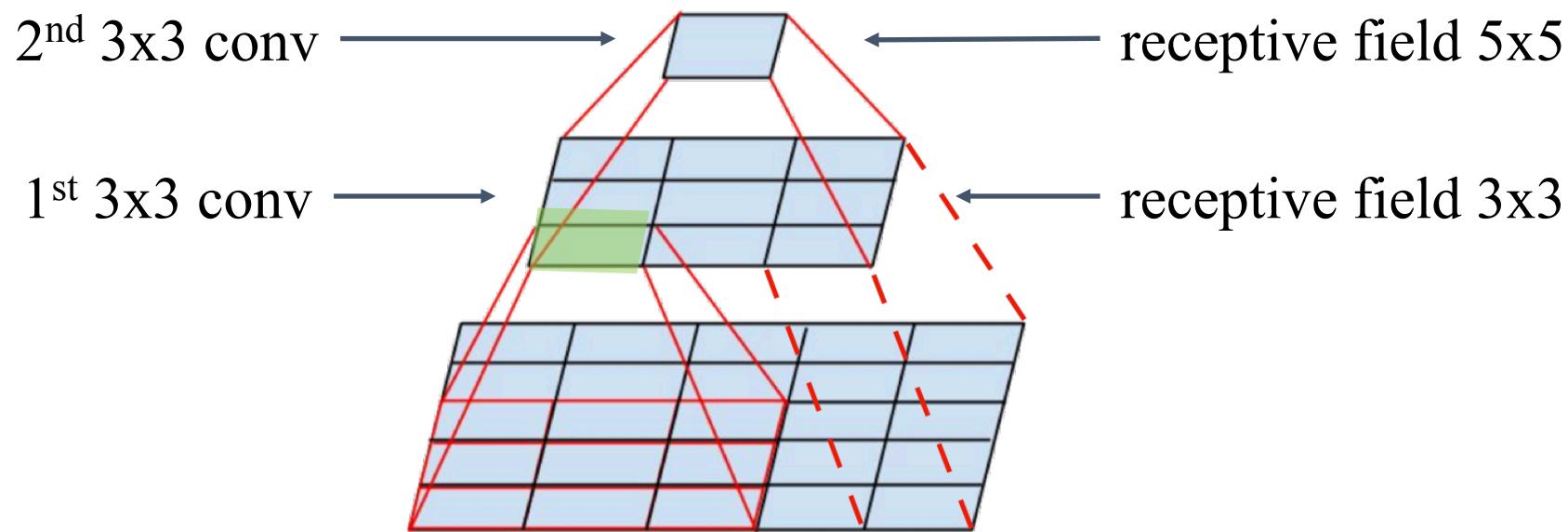
- Using $(W_k * H_k * C_{in} + 1) * C_{out}$ parameters.

One convolutional layer is not enough!

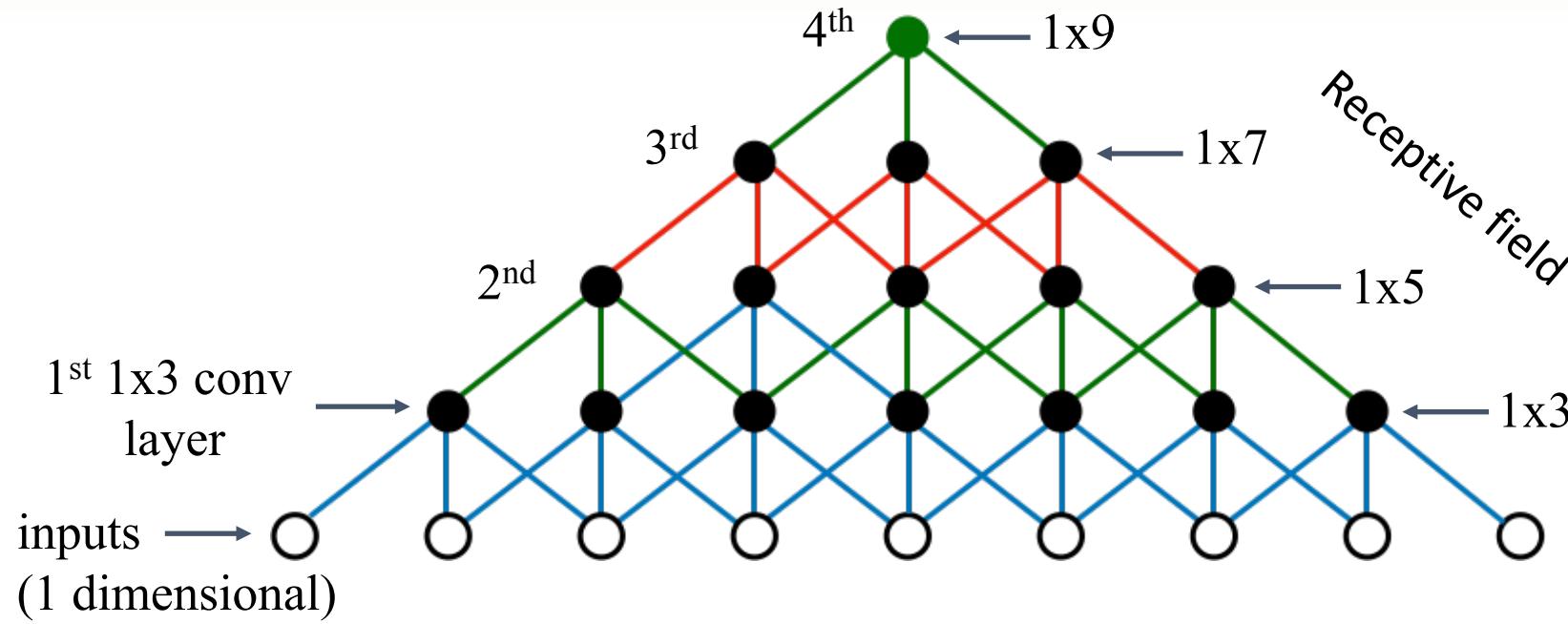
- Let's say neurons of the 1st convolutional layer look at the patches of the image of size 3x3.
- What if an object of interest is bigger than that?
- We need a 2nd convolutional layer on top of the 1st!

One convolutional layer is not enough!

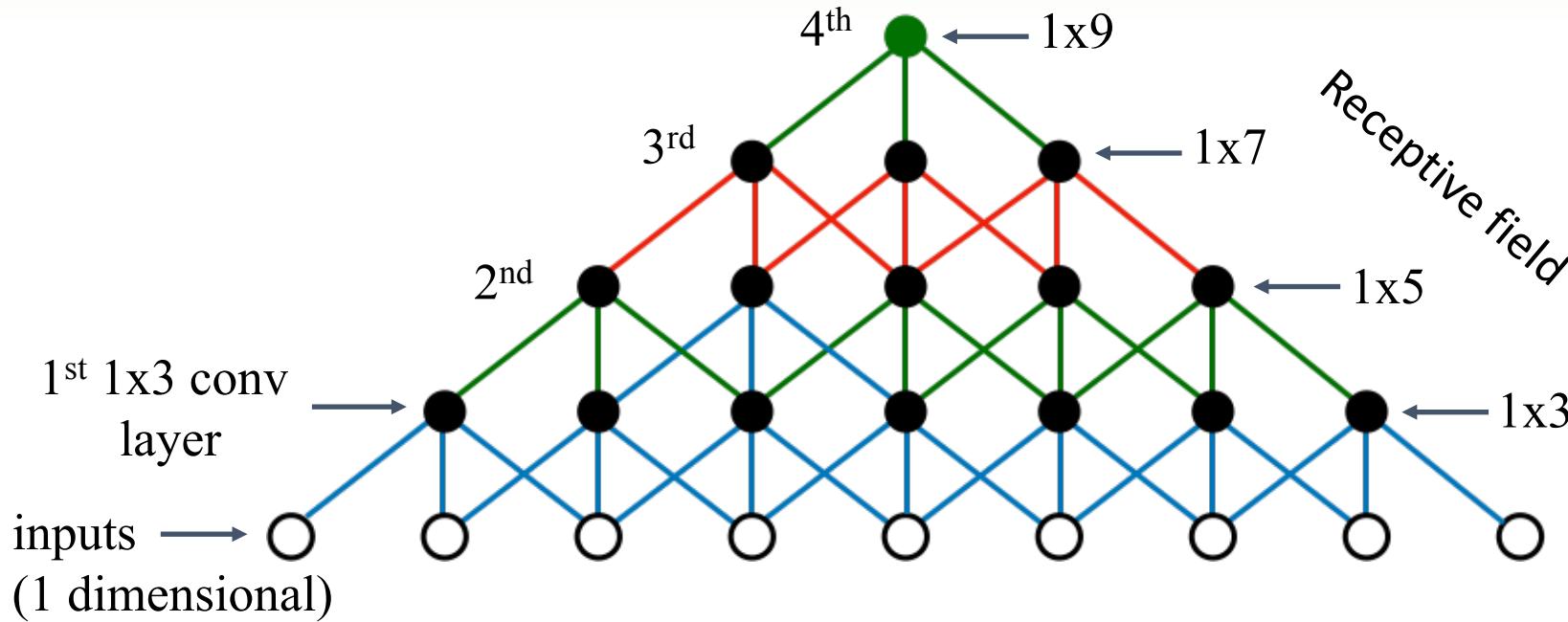
- Let's say neurons of the 1st convolutional layer look at the patches of the image of size 3x3.
- What if an object of interest is bigger than that?
- We need a 2nd convolutional layer on top of the 1st!



Receptive field after N convolutional layers



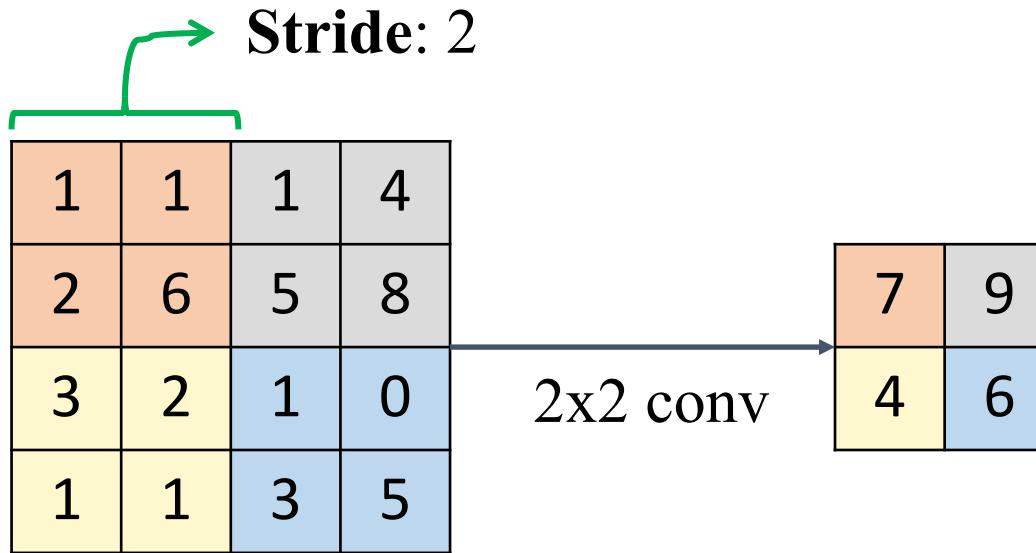
Receptive field after N convolutional layers



- If we stack N convolutional layers with the same kernel size 3x3 the receptive field on N -th layer will be $2N + 1 \times 2N + 1$.
- It looks like we need to stack a lot of convolutional layers! To be able to identify objects as big as the input image **300x300** we will need **150** convolutional layers!

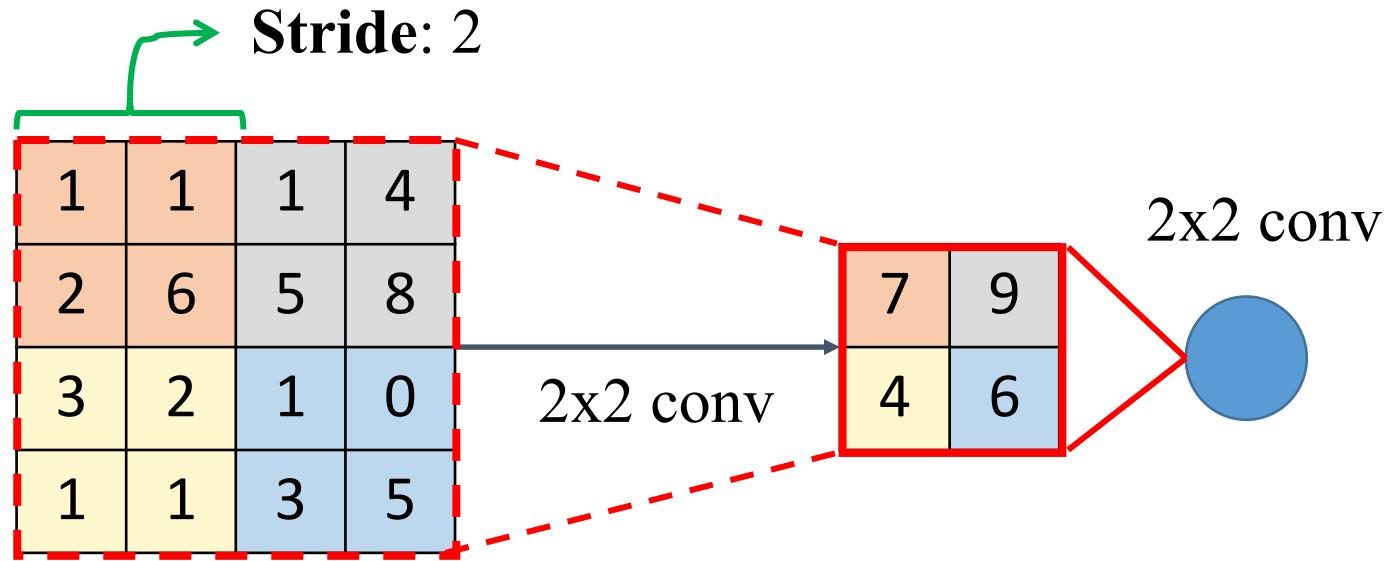
We need to grow receptive field faster!

- We can increase a **stride** in our convolutional layer to reduce the output dimensions!



We need to grow receptive field faster!

- We can increase a **stride** in our convolutional layer to reduce the output dimensions!



Further convolutions will effectively **double** their receptive field!

How do we maintain translation invariance?

0	0	0	0
0	0	0	0
0	0	1	0
0	0	0	1

Input

1	0	0	0
0	1	0	0
0	0	0	0
0	0	0	0

Input

$$\begin{matrix} * & \begin{matrix} 1 & 0 \\ 0 & 1 \end{matrix} & = & \begin{matrix} 0 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 2 \end{matrix} \end{matrix}$$

Kernel Output

Max = 2

Didn't change

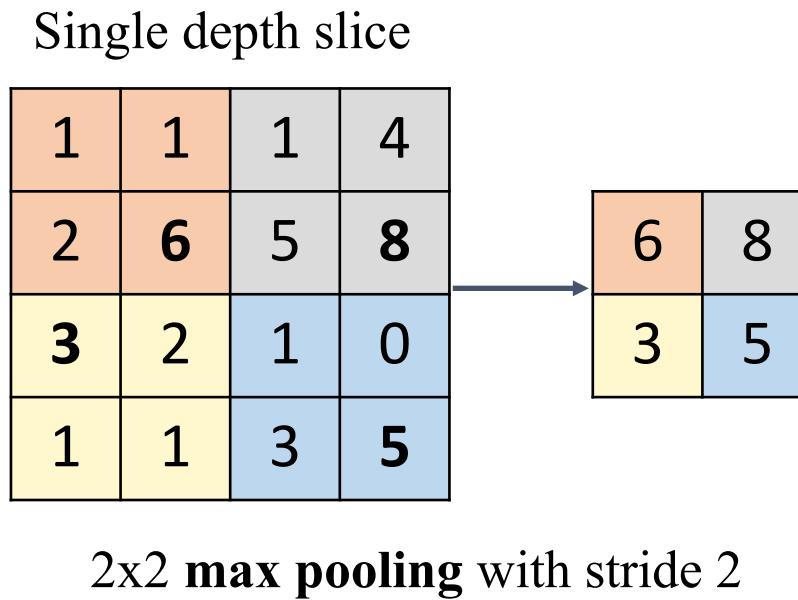
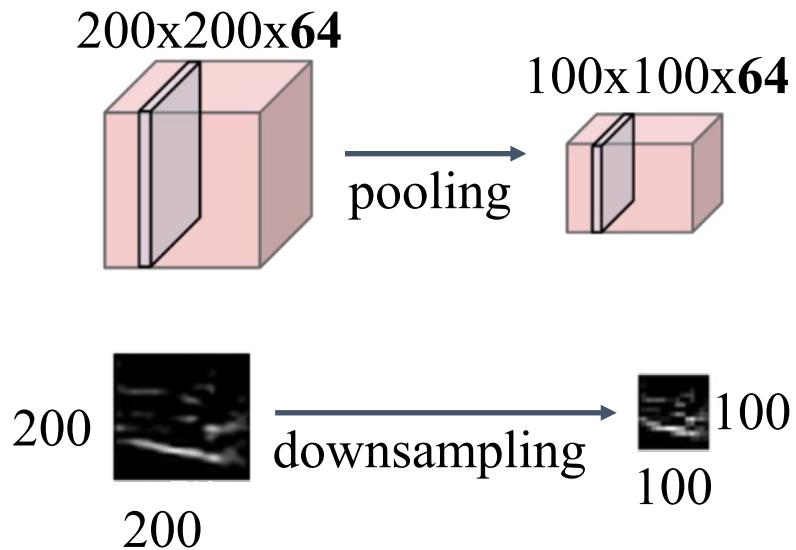
Max = 2

$$\begin{matrix} * & \begin{matrix} 1 & 0 \\ 0 & 1 \end{matrix} & = & \begin{matrix} 2 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \end{matrix} \end{matrix}$$

Kernel Output

Pooling layer will help!

- This layer works like a convolutional layer but doesn't have kernel, instead it calculates **maximum** or **average** of input patch values.



Backpropagation for max pooling layer

Strictly speaking: maximum is not a differentiable function!

Backpropagation for max pooling layer

Strictly speaking: maximum is not a differentiable function!

6	8
3	5

Maximum = 8

7	8
3	5

Maximum = 8

There is no gradient with respect to non maximum patch neurons,
since changing them slightly
does not affect the output.

Backpropagation for max pooling layer

Strictly speaking: maximum is not a differentiable function!

6	8
3	5

Maximum = 8

7	8
3	5

Maximum = 8

There is no gradient with respect to non maximum patch neurons,
since changing them slightly
does not affect the output.

6	8
3	5

Maximum = 8

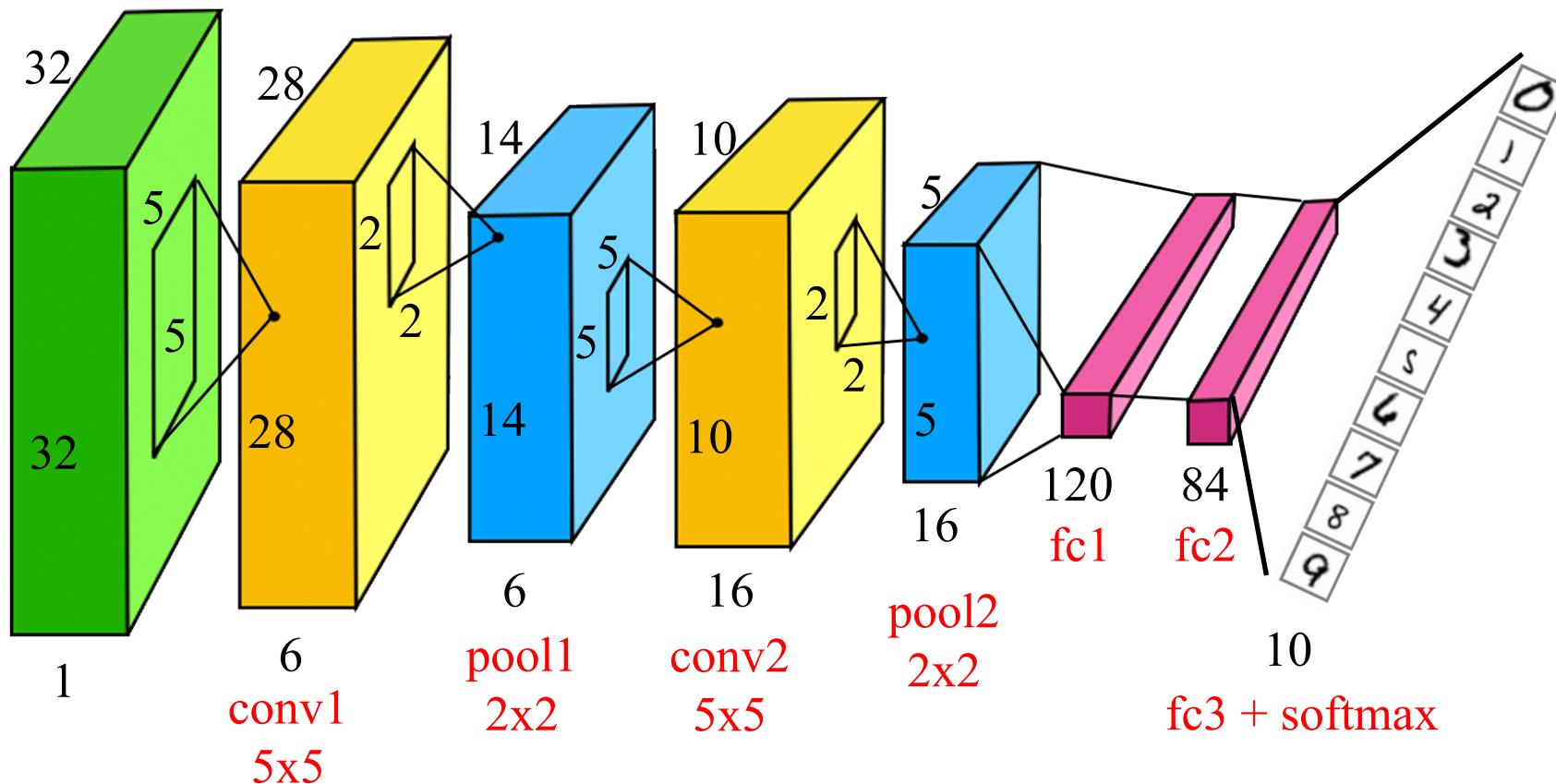
7	9
3	5

Maximum = 9

For the maximum patch neuron we have a gradient of 1.

Putting it all together into a simple CNN

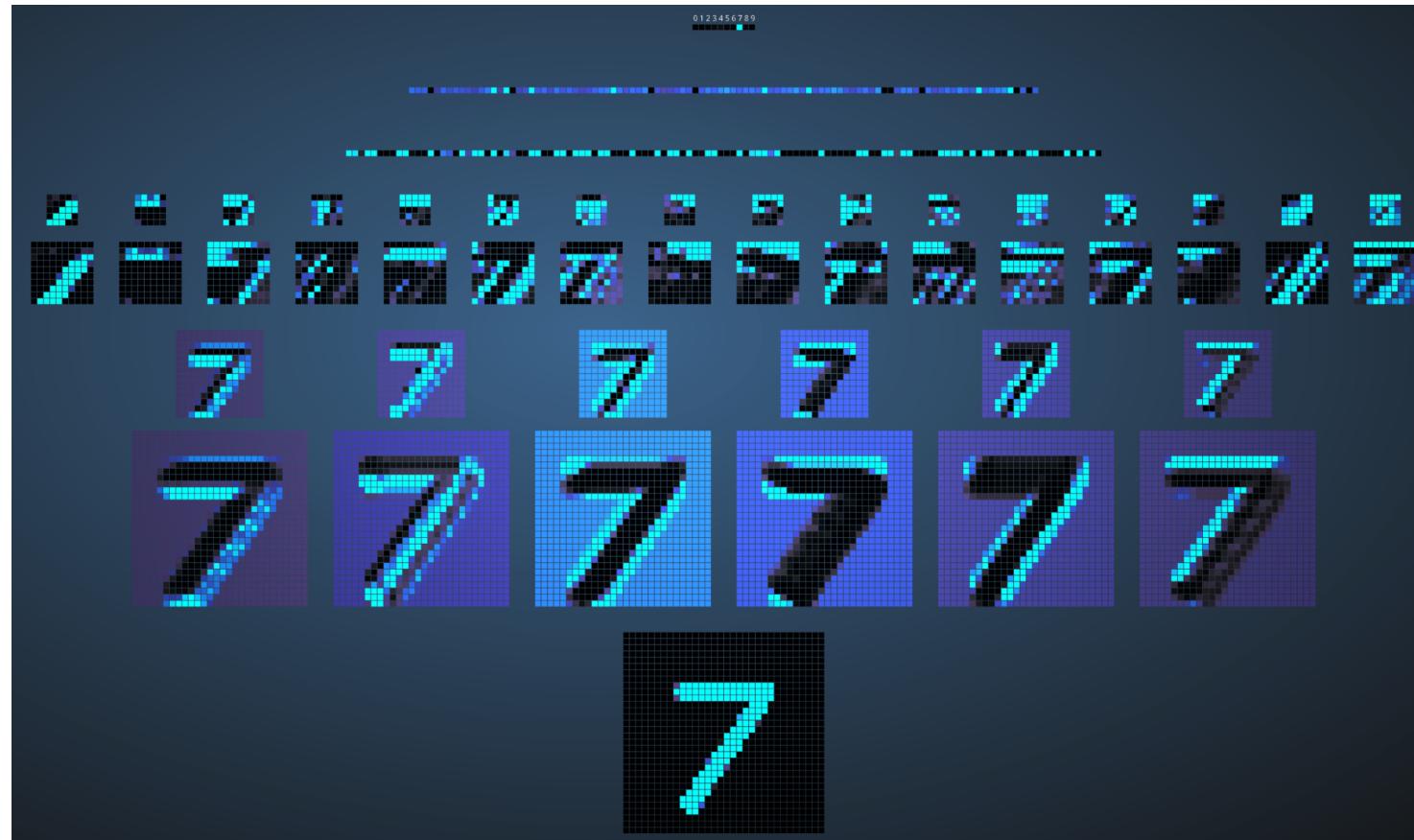
- LeNet-5 architecture (1998) for handwritten digits recognition on MNIST dataset:



<http://yann.lecun.com/exdb/publis/pdf/lecun-98.pdf>

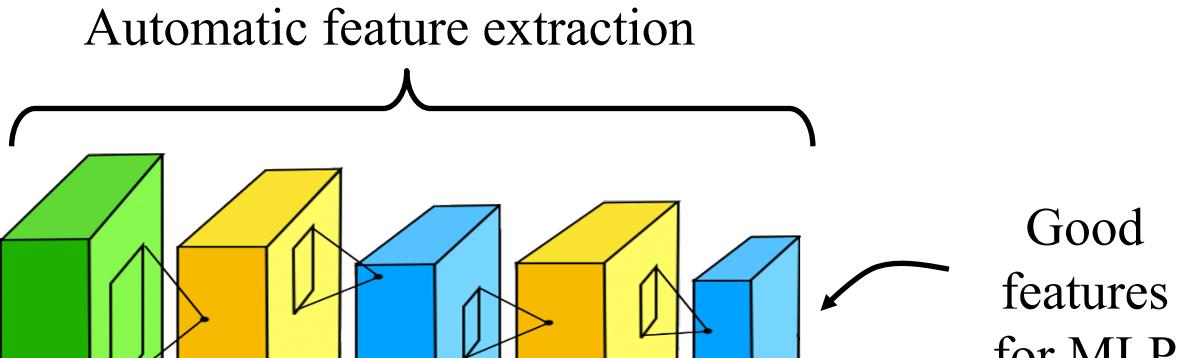
Визуализация обученной сети для MNIST

- <http://scs.ryerson.ca/~aharley/vis/conv/flat.html>



Learning deep representations

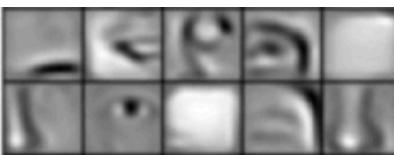
- Neurons of deep convolutional layers learn complex representations that can be used as features for classification with MLP.



Inputs that provide highest activations:



conv1



conv2



conv3

Современные применения: перенос стиля

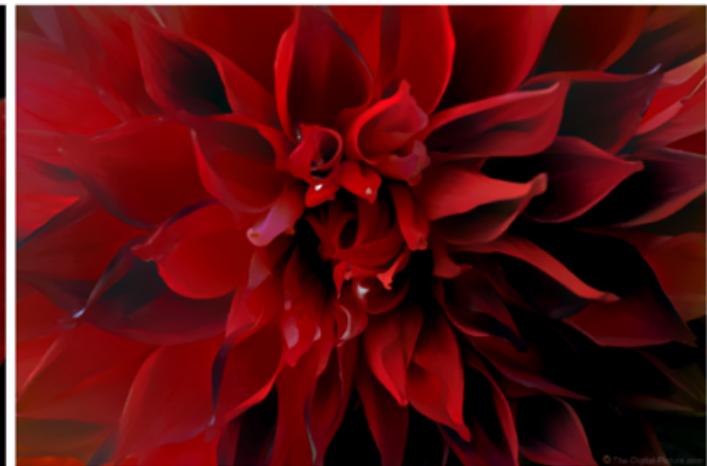
Картина



Стиль



Картина со стилем



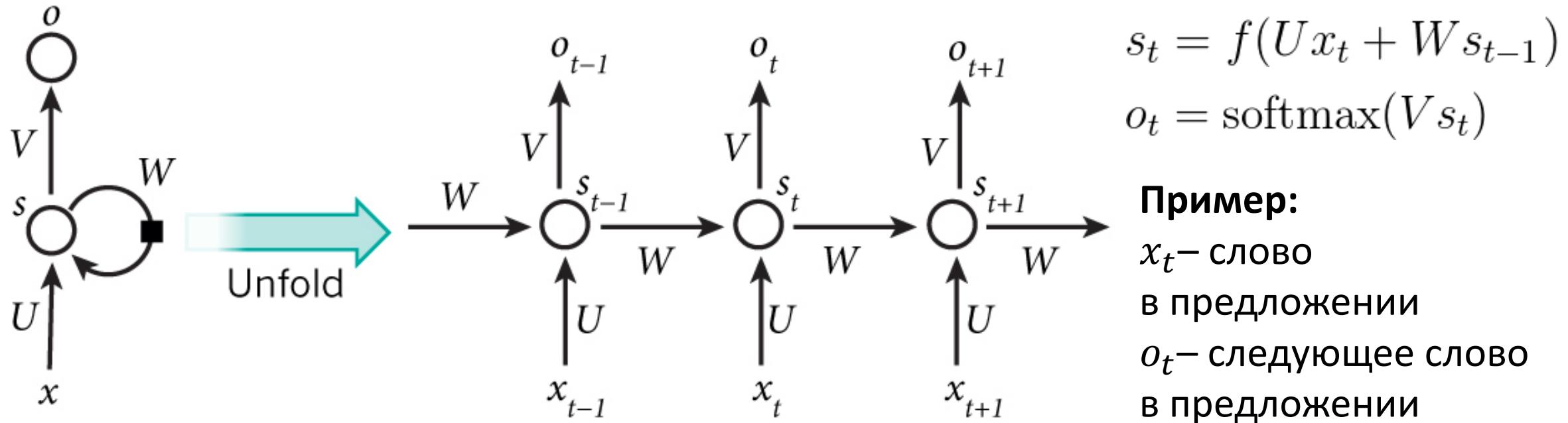
Еще примеры применений

- <https://quickdraw.withgoogle.com>
- <https://affinelayer.com/pixsrv/>

Recurrent Neural Networks (RNN)

- Работают с последовательностями
 - Слов в предложении
 - Букв в предложении
 - Отсчетов в аудио сигнале (амплитуда, частота)
 - Пикселей изображения
 - ...

Как устроена простая RNN

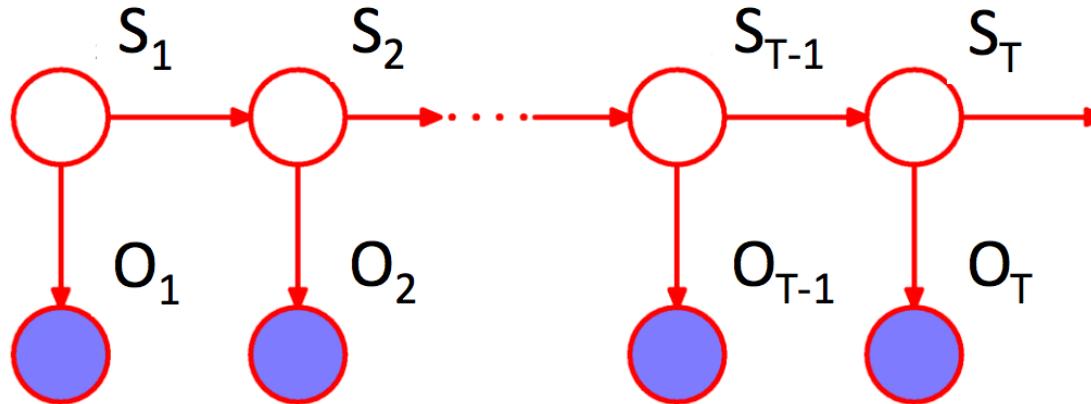


- Работает **одинаково** для каждого элемента последовательности x_t , но вычисления зависят от предыдущих элементов x_t
- Можно сказать, что у RNN есть память (**скрытое состояние** s_t), в которой хранится информация о предыдущих элементах последовательности

А как же Hidden Markov Model (HMM)?

$O_t \in \{y_1, y_2, \dots, y_K\}$

$S_t \in \{1, \dots, I\}$



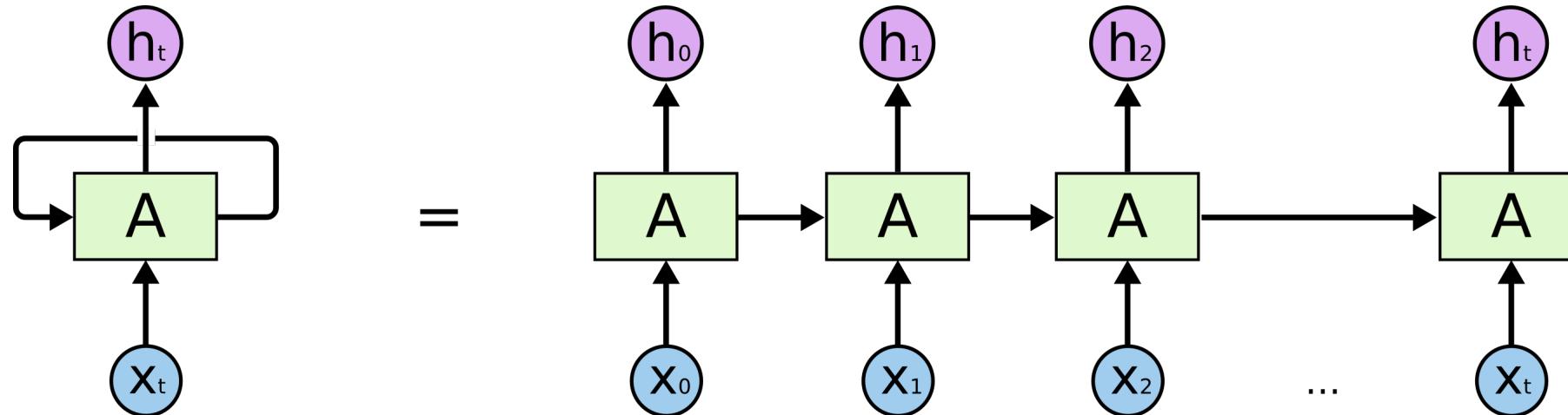
Пример:
 S_t – POS tag
 O_t – слово в
предложении

$$p(S_1, \dots, S_T, O_1, \dots, O_T) = \prod_{t=1}^T p(O_t | S_t) \prod_{t=1}^T p(S_t | S_{t-1})$$

1st order Markov assumption on hidden states $\{S_t\}$ $t = 1, \dots, T$
(can be extended to higher order).

Backpropagation through time (BPTT)

- Если развернуть сеть по времени t , то получим обычную feed-forward сеть с shared параметрами
- Применяем backpropagation, считаем градиенты для каждого параметра
- Далее (так же, как в случае сверточных CNN) суммируем градиенты по shared параметрам и делаем шаг SGD



Посчитаем производную для 2 шагов

$$s_1 = f(Ux_1 + Ws_0)$$
$$a_2 \quad \downarrow$$
$$s_2 = f(Ux_2 + Ws_1) = f(Ux_2 + Wf(Ux_1 + Ws_0))$$

$$\frac{\partial s_2}{\partial W} = \frac{\partial f}{\partial a_2} \left(W \frac{\partial s_1}{\partial W} + s_1 \right) = \boxed{\frac{\partial f}{\partial a_2}} W \frac{\partial s_1}{\partial W} + \boxed{\frac{\partial f}{\partial a_2}} s_1$$

$$\frac{\partial s_1}{\partial W} = \frac{\partial s_1}{\partial W_*}$$

$$\frac{\partial s_2}{\partial s_1} \quad \frac{\partial s_2}{\partial W_*}$$

$$\frac{\partial s_2}{\partial W} = \frac{\partial s_2}{\partial s_1} \frac{\partial s_1}{\partial W_*} + \frac{\partial s_2}{\partial W_*}$$

Обозначение
в предположении,
что s_1 не зависит от w

Посчитаем производную для 3 шагов

$$s_3 = f(Ux_3 + Ws_2)$$

$$\frac{\partial s_2}{\partial W} = \frac{\partial s_2}{\partial s_1} \frac{\partial s_1}{\partial W_*} + \frac{\partial s_2}{\partial W_*}$$

$$\frac{\partial s_3}{\partial s_2}$$

$$\frac{\partial s_3}{\partial W_*}$$

$$\frac{\partial s_3}{\partial W} = \frac{\partial f}{\partial a_3} \left(W \frac{\partial s_2}{\partial W} + s_2 \right) = \boxed{\frac{\partial f}{\partial a_3}} W \left(\frac{\partial s_2}{\partial s_1} \frac{\partial s_1}{\partial W_*} + \frac{\partial s_2}{\partial W_*} \right) + \boxed{\frac{\partial f}{\partial a_3}} s_2$$

$$\frac{\partial s_3}{\partial W} = \frac{\partial s_3}{\partial s_2} \frac{\partial s_2}{\partial s_1} \frac{\partial s_1}{\partial W_*} + \frac{\partial s_3}{\partial s_2} \frac{\partial s_2}{\partial W_*} + \frac{\partial s_3}{\partial W_*}$$

Обозначение
в предположении,
что s_2 не зависит от w

Индукцией легко показать, что...

$$\frac{\partial s_2}{\partial W} = \frac{\partial s_2}{\partial s_1} \frac{\partial s_1}{\partial W_*} + \frac{\partial s_2}{\partial W_*}$$

$$\frac{\partial s_3}{\partial W} = \frac{\partial s_3}{\partial s_2} \frac{\partial s_2}{\partial s_1} \frac{\partial s_1}{\partial W_*} + \frac{\partial s_3}{\partial s_2} \frac{\partial s_2}{\partial W_*} + \frac{\partial s_3}{\partial W_*}$$

$$\frac{\partial s_k}{\partial W} = \sum_{i=1}^k \left(\prod_{j=i+1}^k \frac{\partial s_j}{\partial s_{j-1}} \right) \frac{\partial s_i}{\partial W_*}$$

ВРТТ на примере

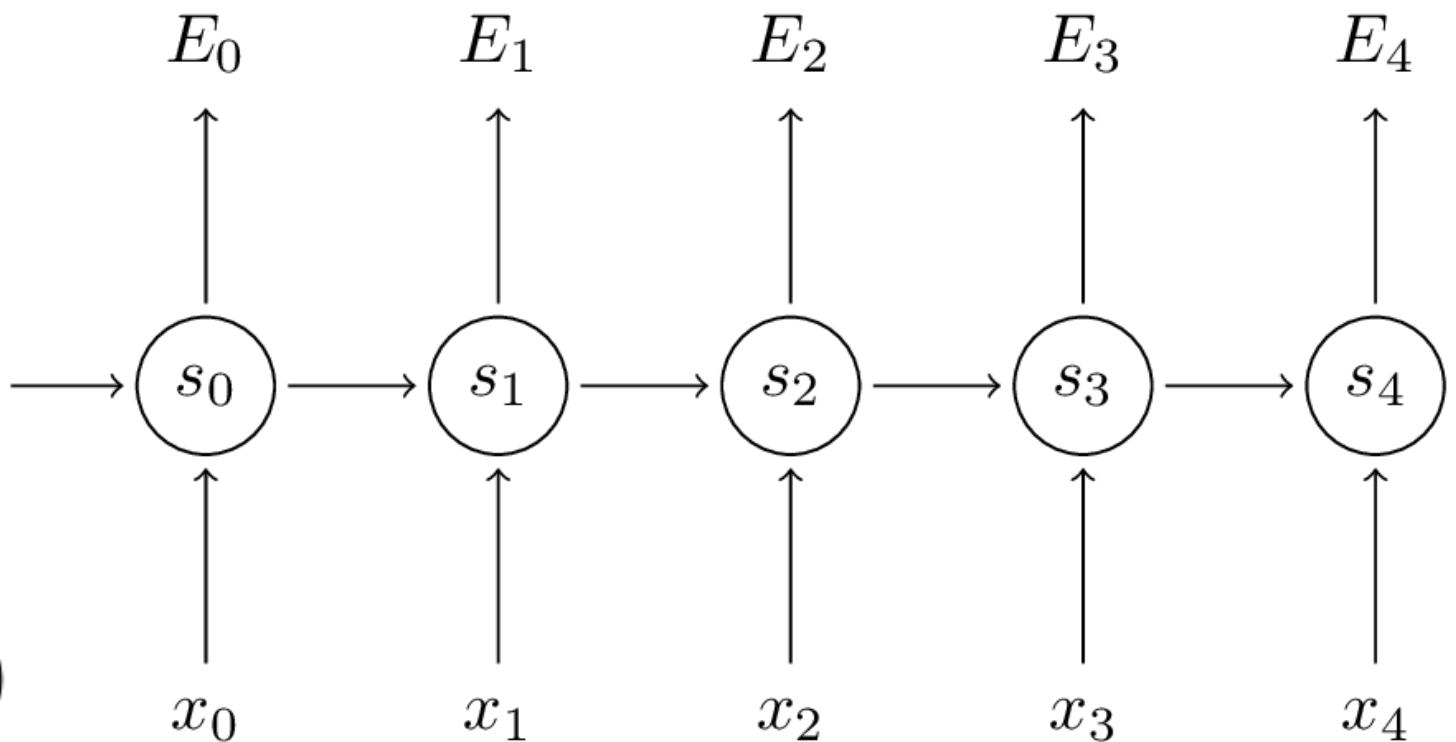
$$s_t = \tanh(Ux_t + Ws_{t-1})$$

$$\hat{y}_t = \text{softmax}(Vs_t)$$

$$E_t(y_t, \hat{y}_t) = -y_t \log \hat{y}_t$$

$$E(y, \hat{y}) = \sum_t E_t(y_t, \hat{y}_t)$$

$$= - \sum_t y_t \log \hat{y}_t$$

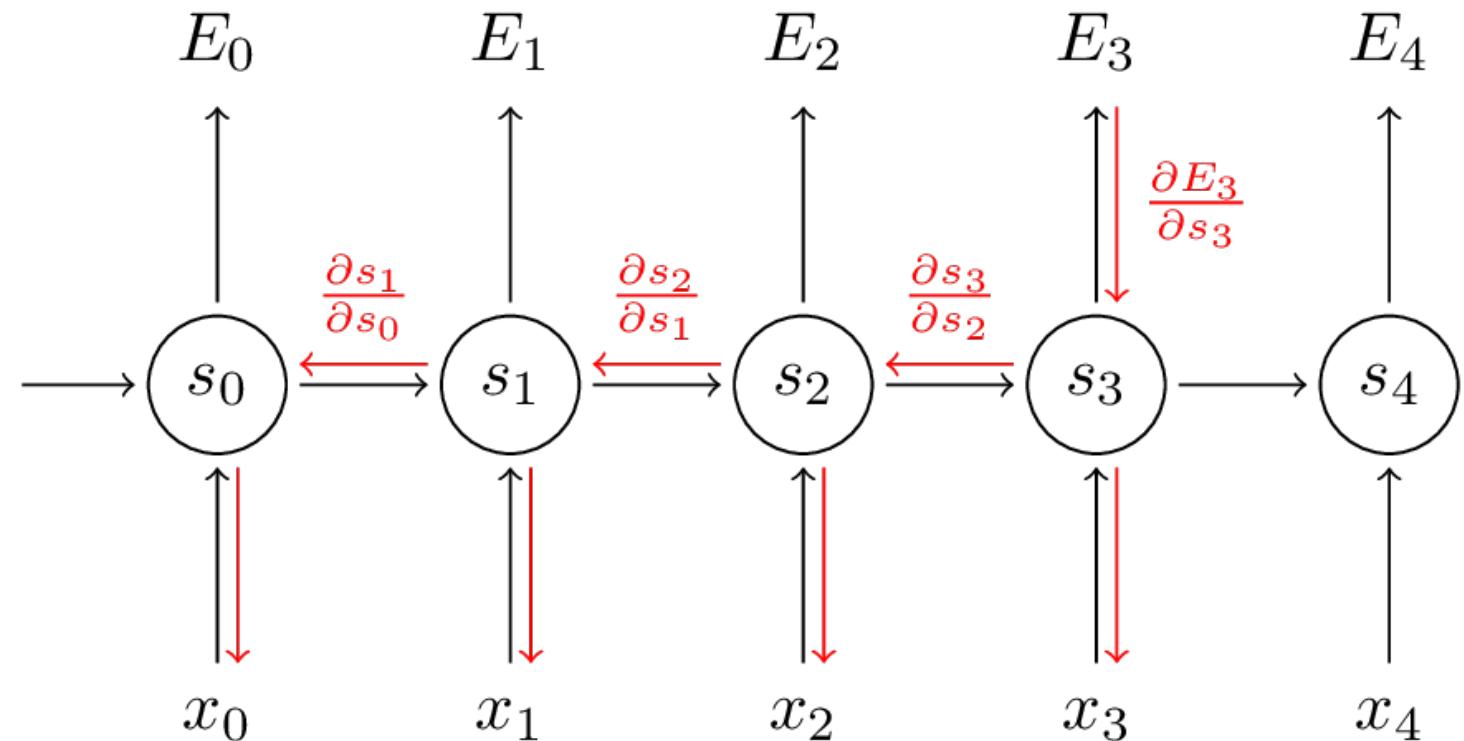


BPTT на примере

$$s_t = \tanh(Ux_t + Ws_{t-1})$$

$$\hat{y}_t = \text{softmax}(Vs_t)$$

$$\begin{aligned} \frac{\partial E_3}{\partial W} &= \frac{\partial E_3}{\partial \hat{y}_3} \frac{\partial \hat{y}_3}{\partial s_3} \frac{\partial s_3}{\partial W} = \\ &= \frac{\partial E_3}{\partial \hat{y}_3} \frac{\partial \hat{y}_3}{\partial s_3} \sum_{i=1}^3 \left(\prod_{j=i+1}^3 \frac{\partial s_j}{\partial s_{j-1}} \right) \frac{\partial s_i}{\partial W_*} \end{aligned}$$



Проблема с затухающими градиентами

$$\frac{\partial E_3}{\partial W} = \frac{\partial E_3}{\partial \hat{y}_3} \frac{\partial \hat{y}_3}{\partial s_3} \sum_{i=1}^3 \left(\prod_{j=i+1}^3 \frac{\partial s_j}{\partial s_{j-1}} \right) \frac{\partial s_i}{\partial W_*}$$

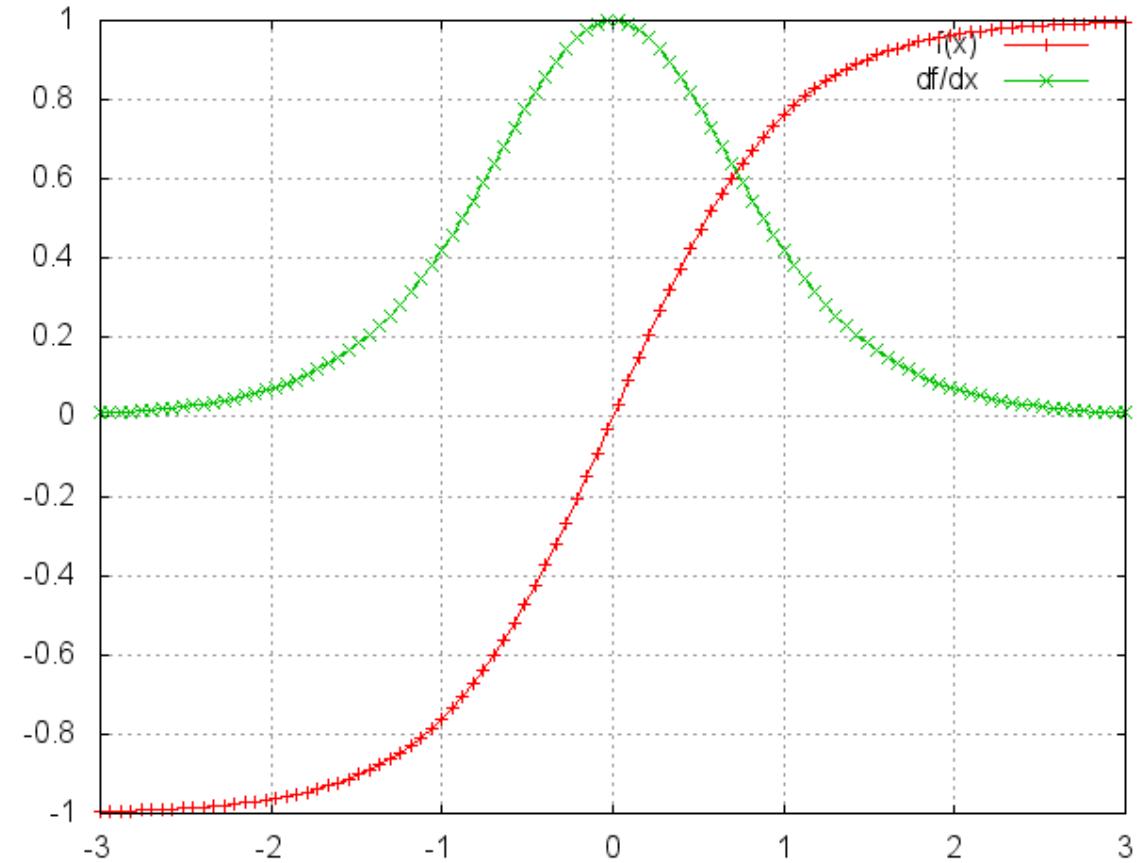
$$s_t = \tanh(Ux_t + Ws_{t-1})$$

$$\frac{\partial s_j}{\partial s_{j-1}} = \frac{\partial f}{\partial a_j} W$$

$$f(x) = \tanh(x)$$

$$\prod_{j=i+1}^3 \frac{\partial s_j}{\partial s_{j-1}}$$

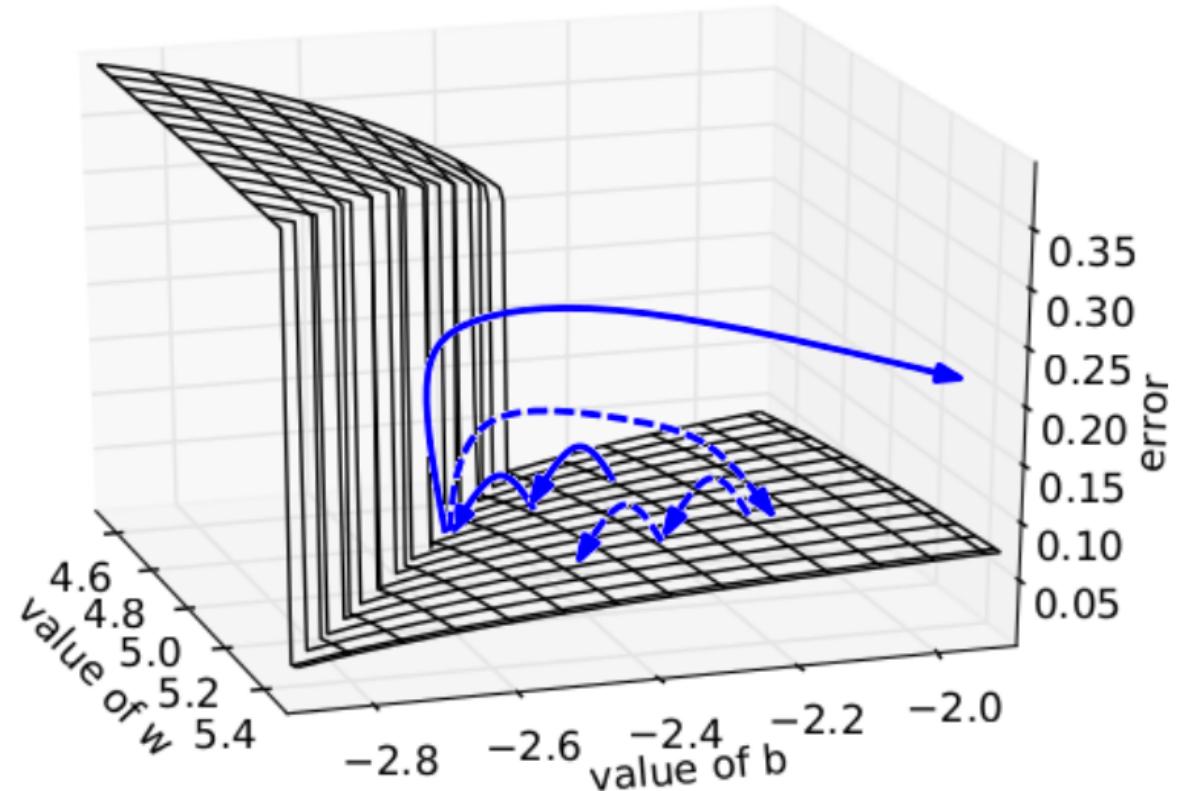
Быстро убывает
или взрывается



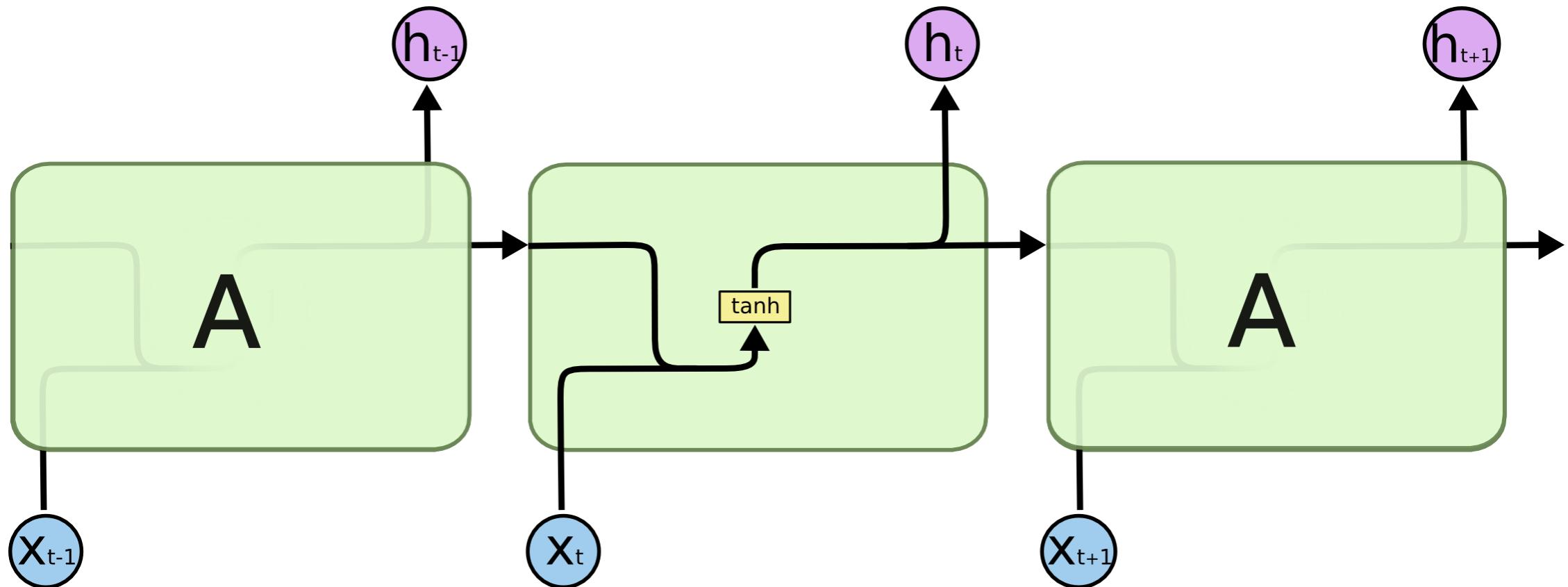
Gradient clipping

- Градиенты могут взорваться – ограничим их норму

```
 $\hat{g} \leftarrow \frac{\partial \varepsilon}{\partial \theta}$ 
if  $\|\hat{g}\| \geq threshold$  then
     $\hat{g} \leftarrow \frac{threshold}{\|\hat{g}\|} \hat{g}$ 
end if
```

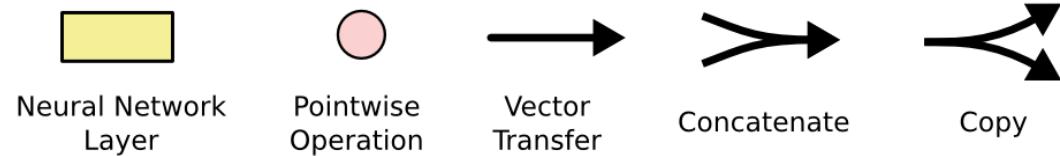
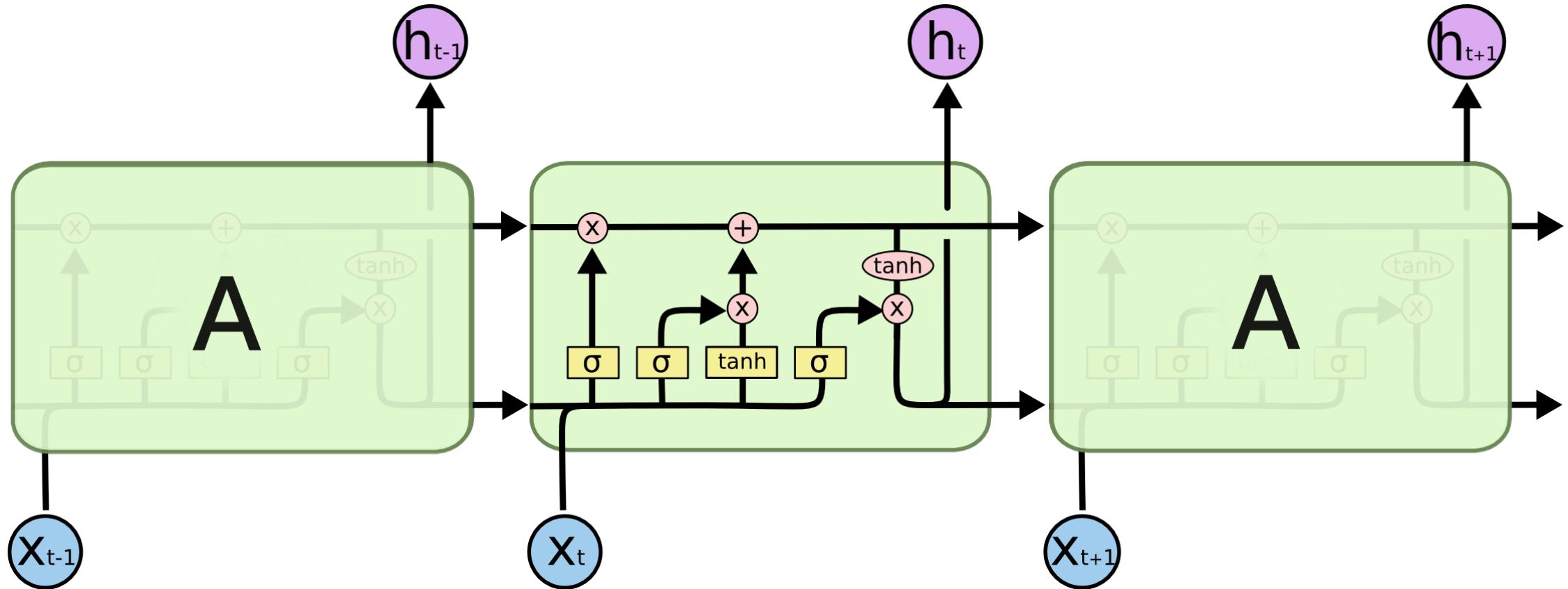


Простая RNN плохо работает



$$h_t = \tanh(Ux_t + Wh_{t-1})$$

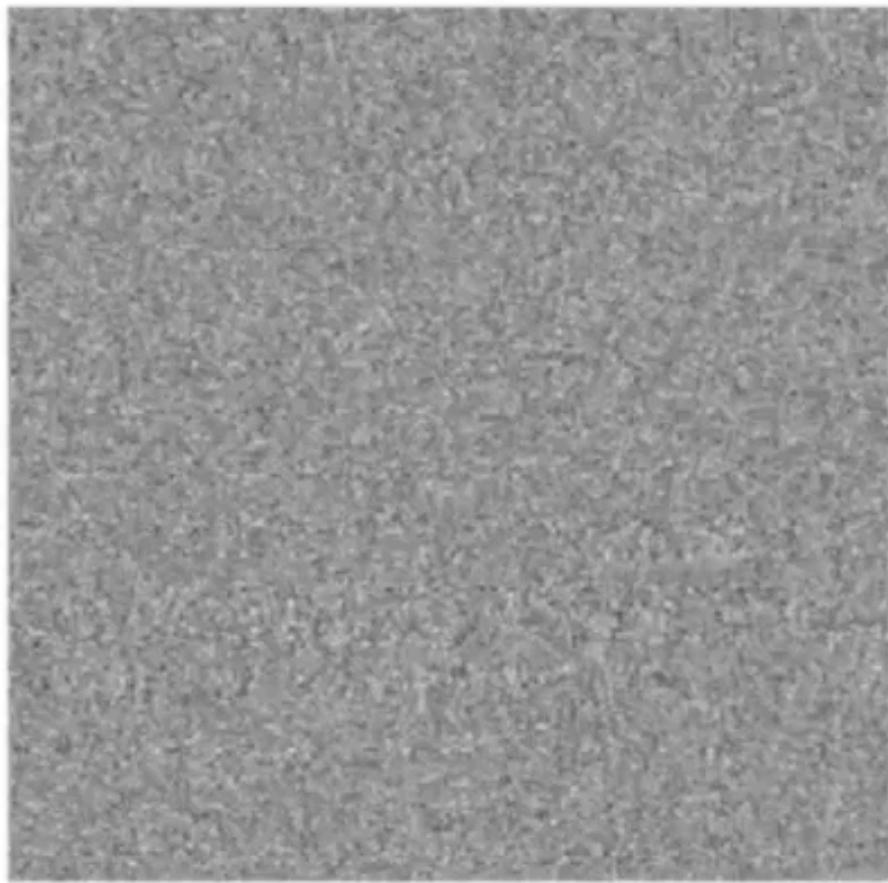
Спасет LSTM



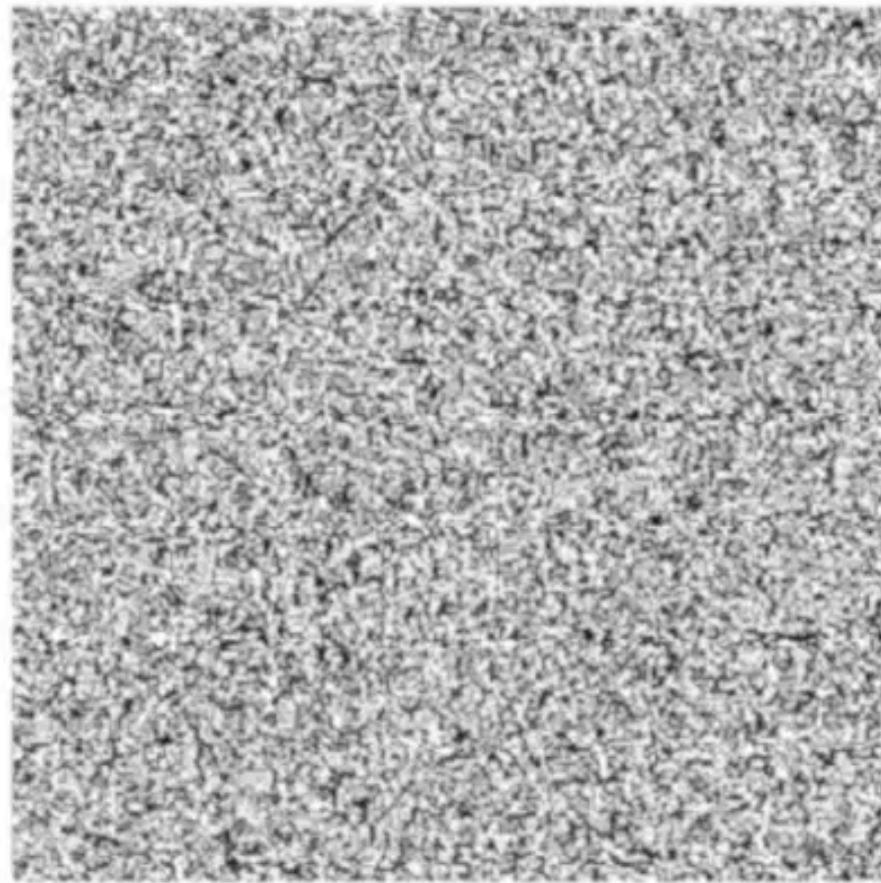
LSTM решает проблему затухания градиентов

Распространение ошибки на 128 шаге

127

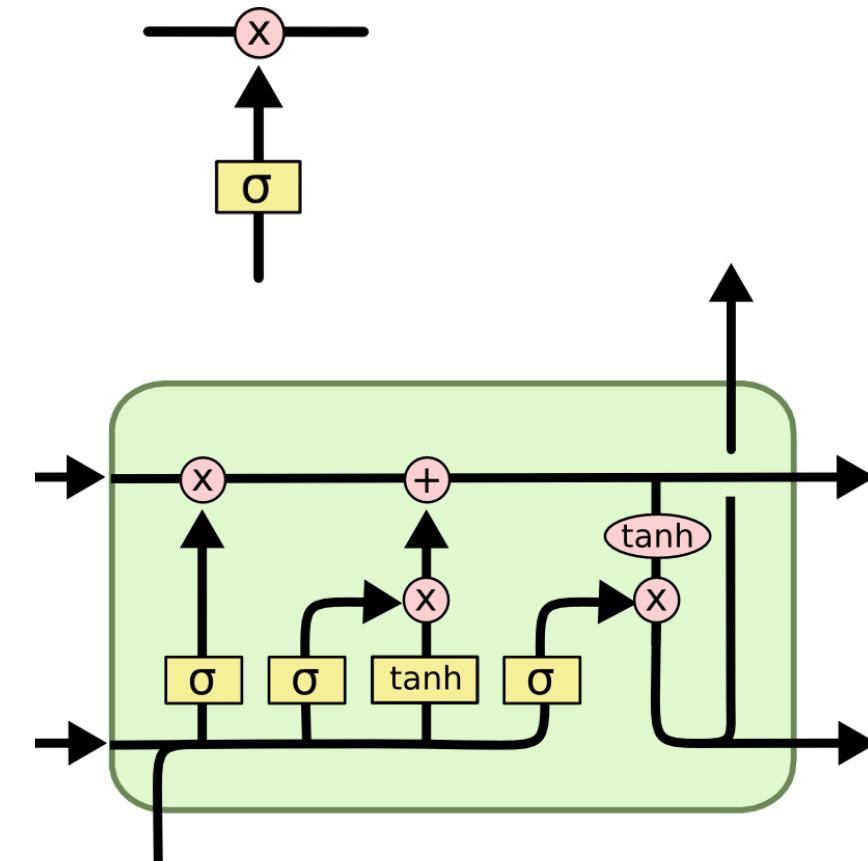
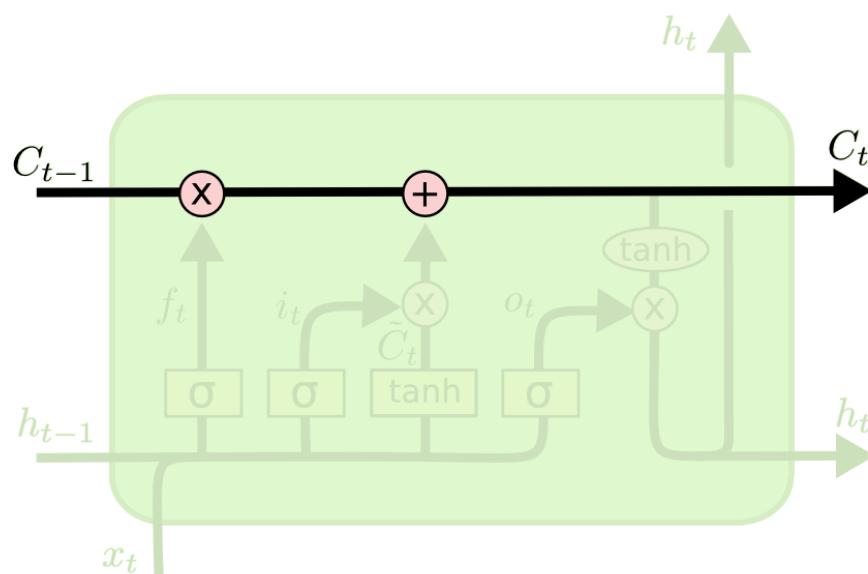


127

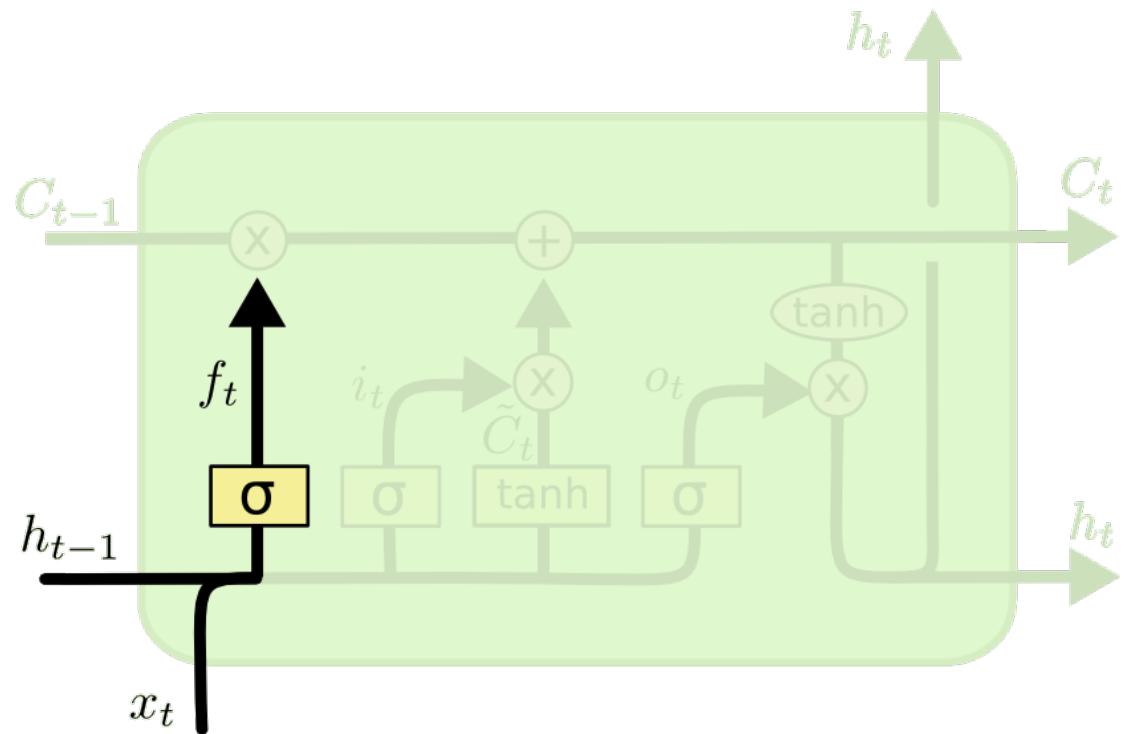


LSTM cell state

- C_t проходит через все ячейки, LSTM может забывать или добавлять информацию в C_t
- Gates учат маски для забывания (**forget**), добавления (**input**) и вывода (**output**) вектора состояний C_t

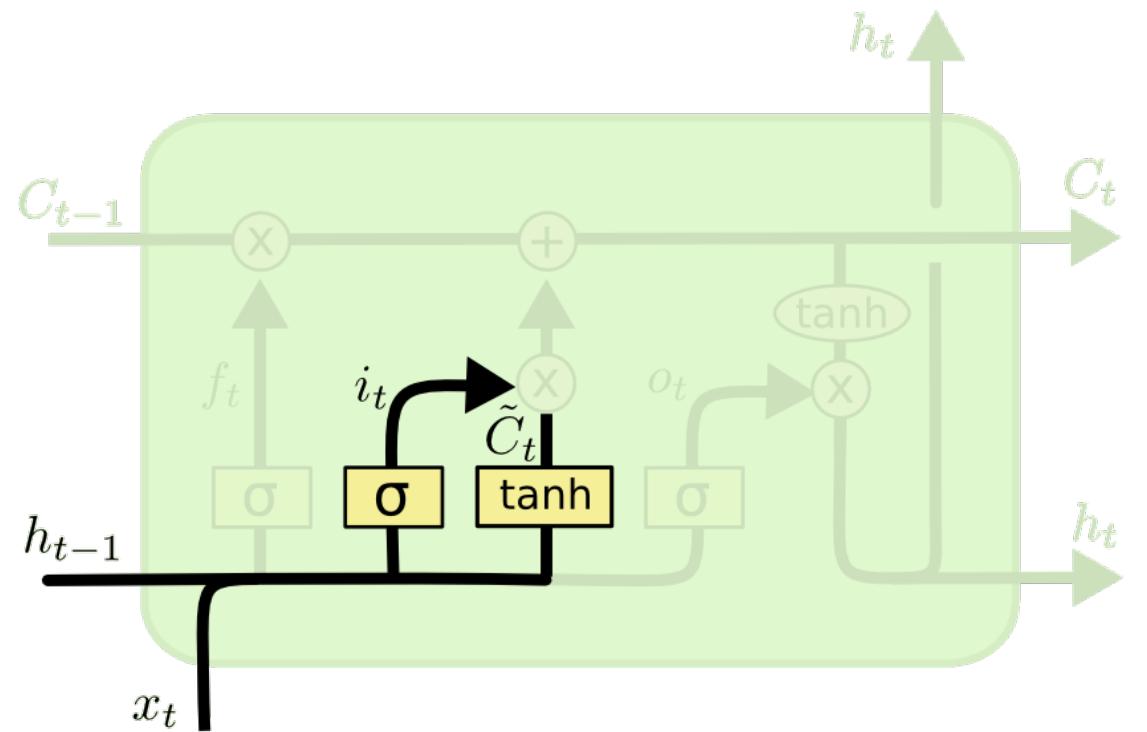


LSTM forget gate



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

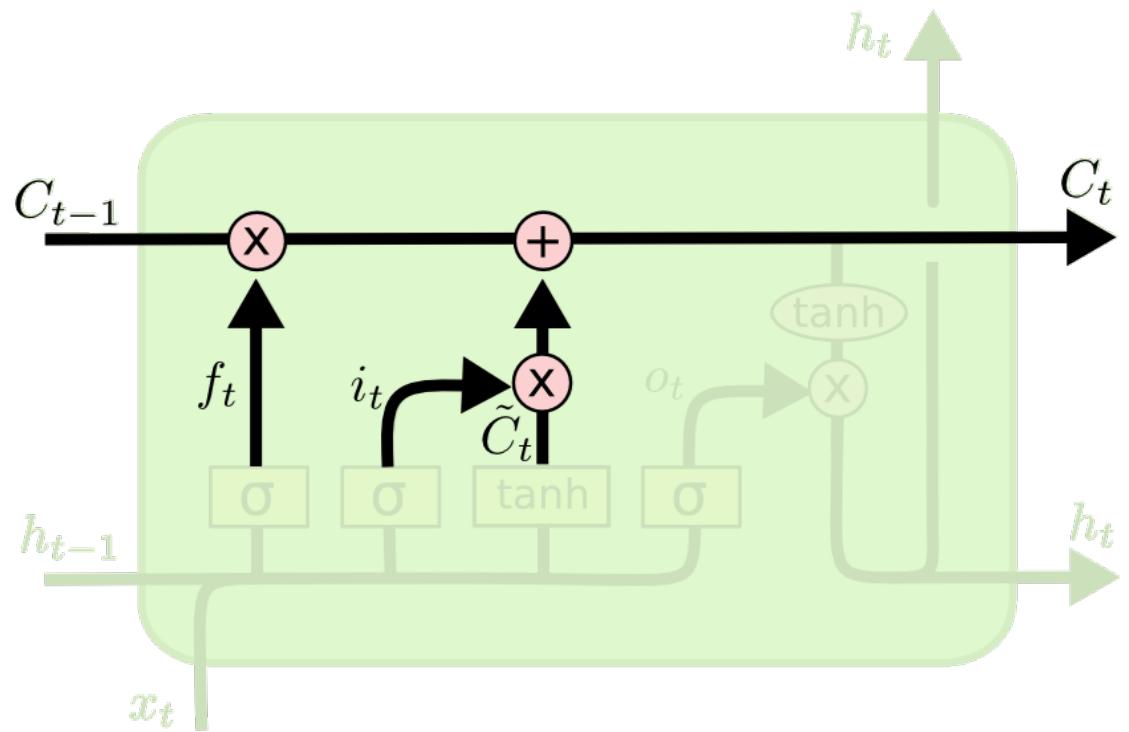
LSTM input gate и вклад ячейки в состояние



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

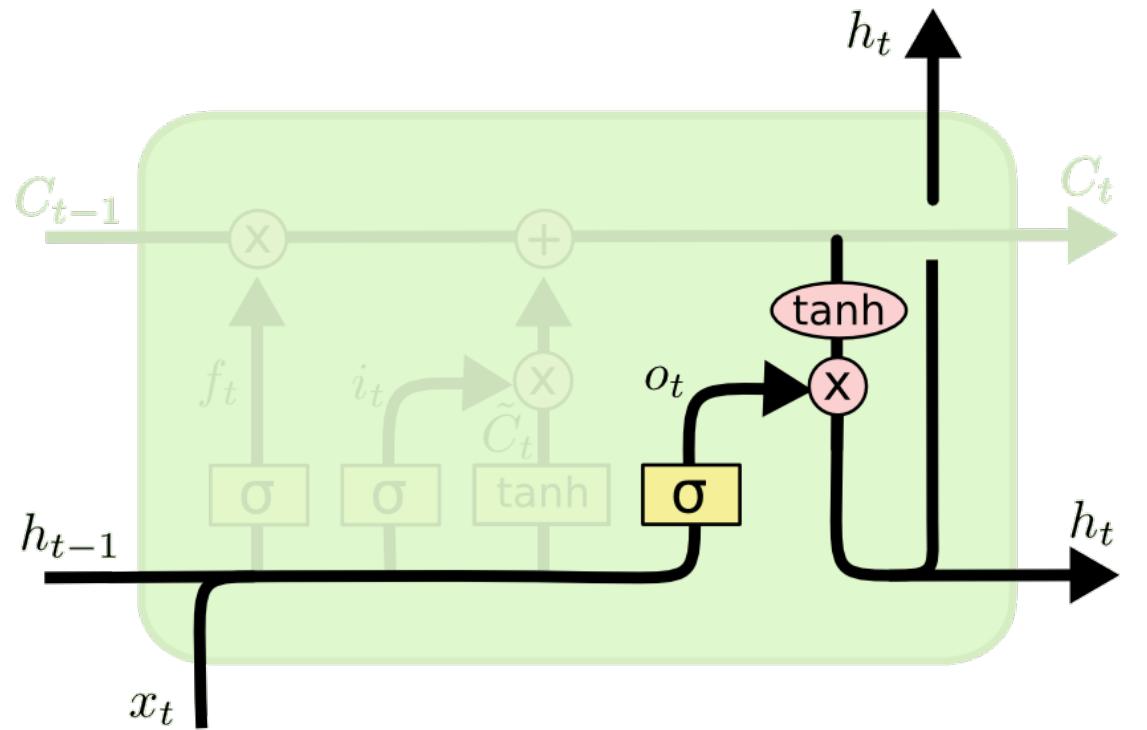
$$\tilde{C}_t = \tanh(W_C \cdot [h_{t-1}, x_t] + b_C)$$

Обновление состояния



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

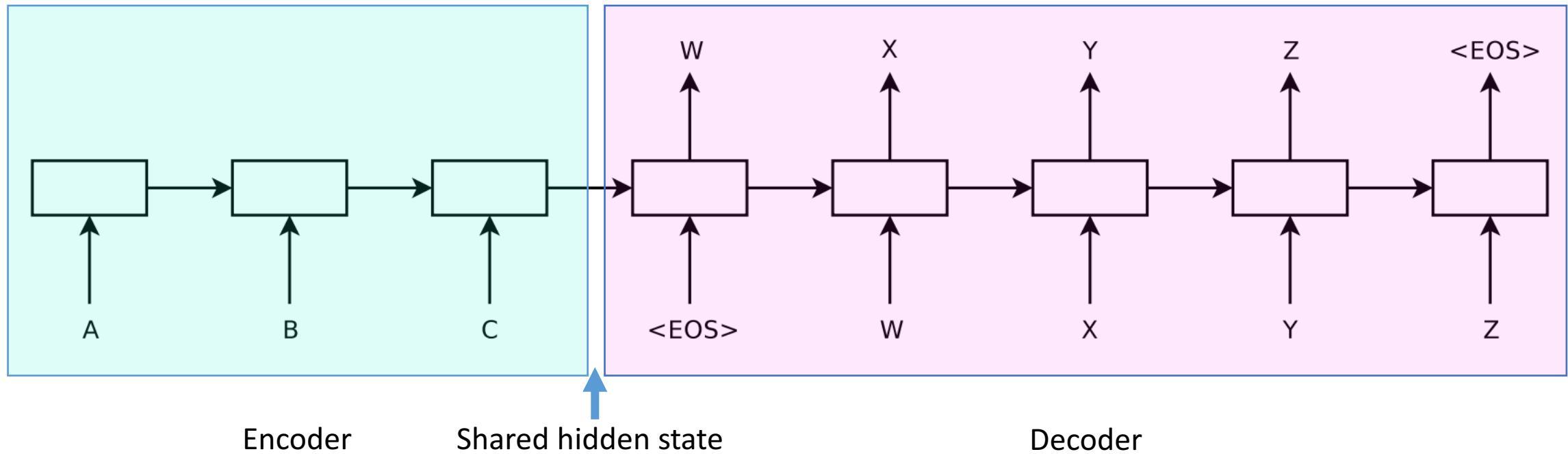
Выдача скрытого состояния наружу



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

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

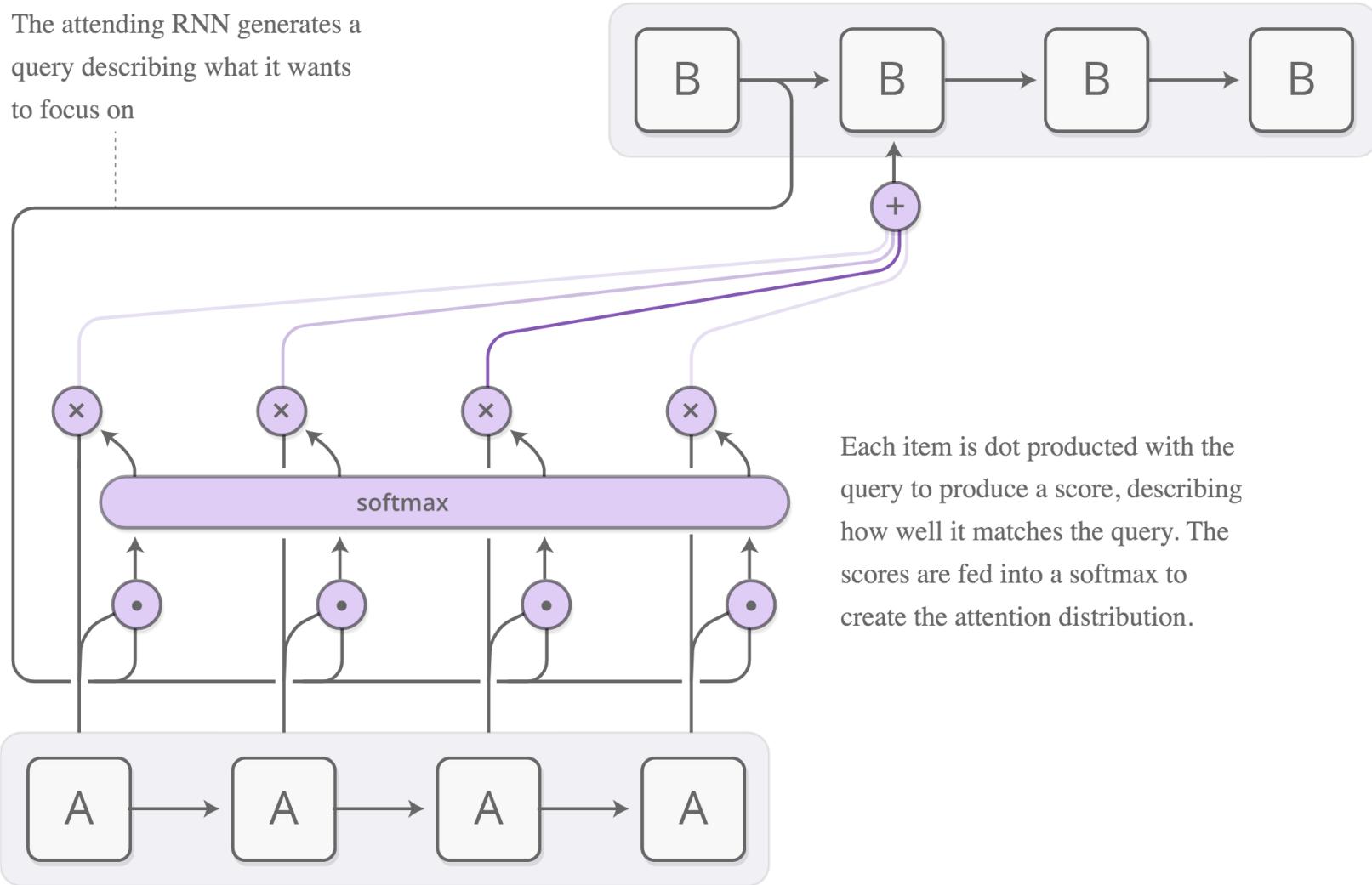
Encoder-decoder (seq2seq)



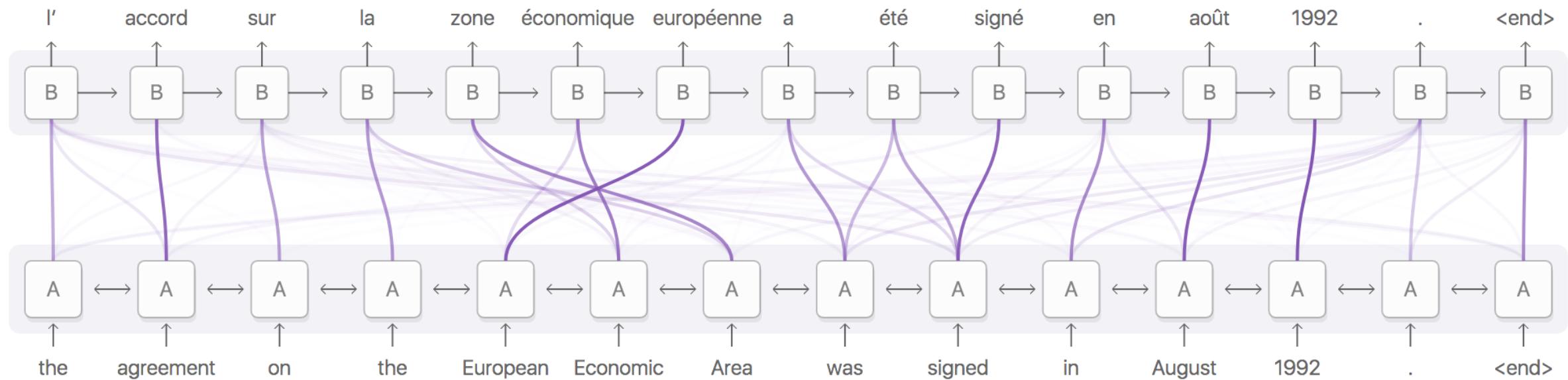
- Слабое место – связь между encoder и decoder только через последний скрытый слой encoder

Поможет attention

The attending RNN generates a query describing what it wants to focus on

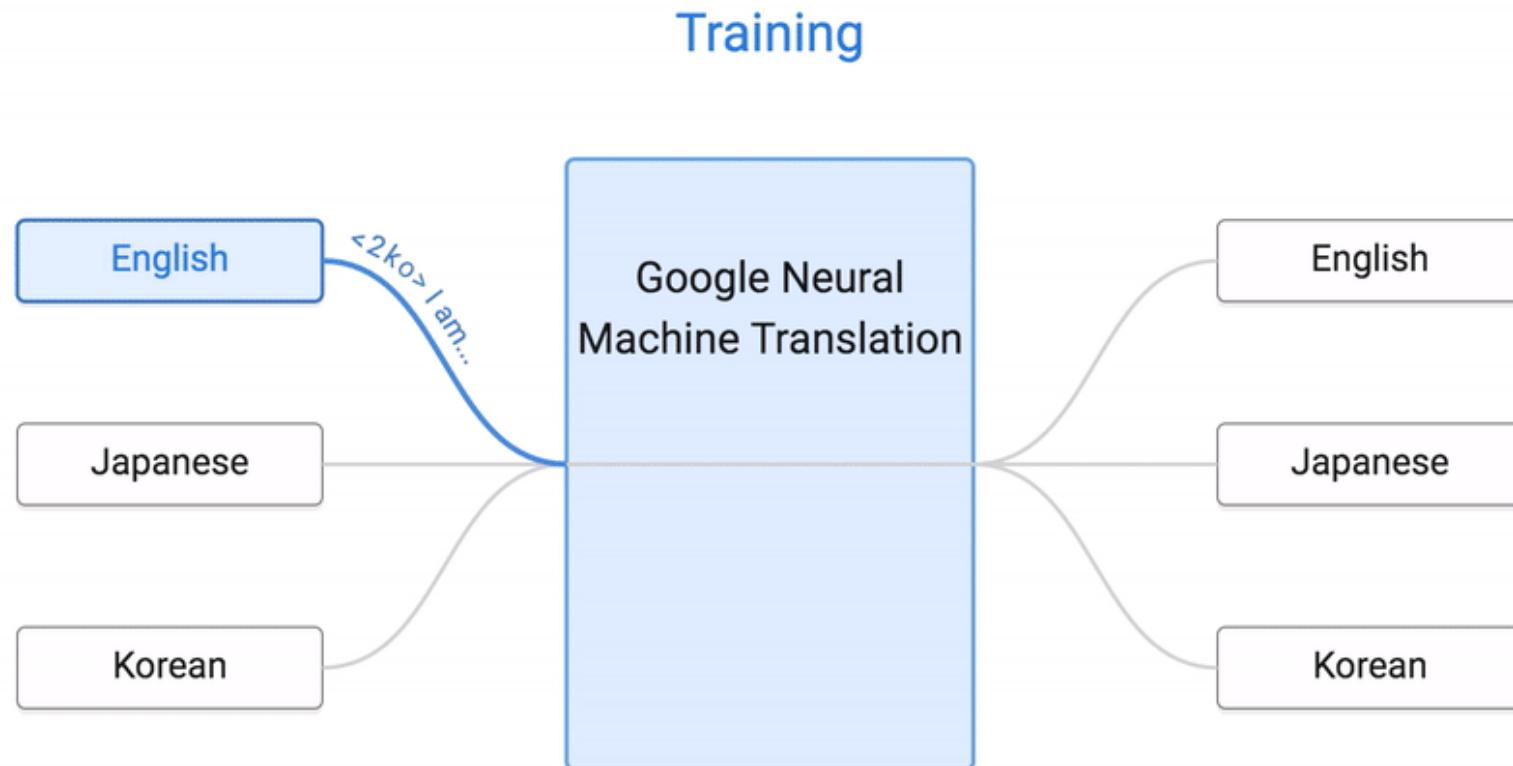


Машинный перевод (seq2seq + attention)



Multi-language перевод от Google

- Обучение encoder и decoder для каждого языка на разных парах языков одновременно



Распознавание голоса (seq2seq + attention)

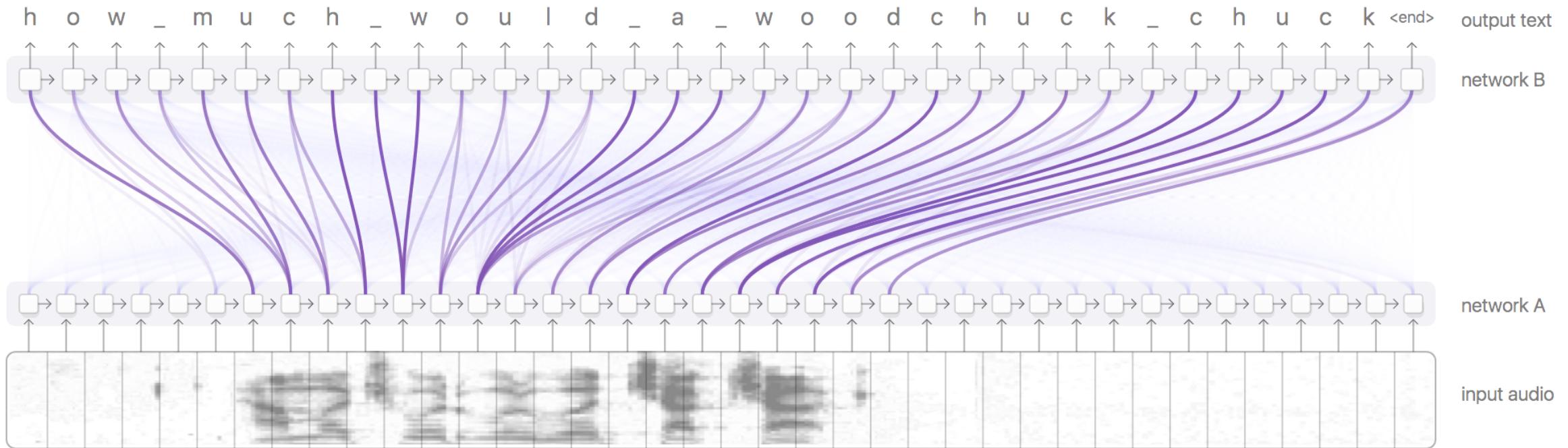
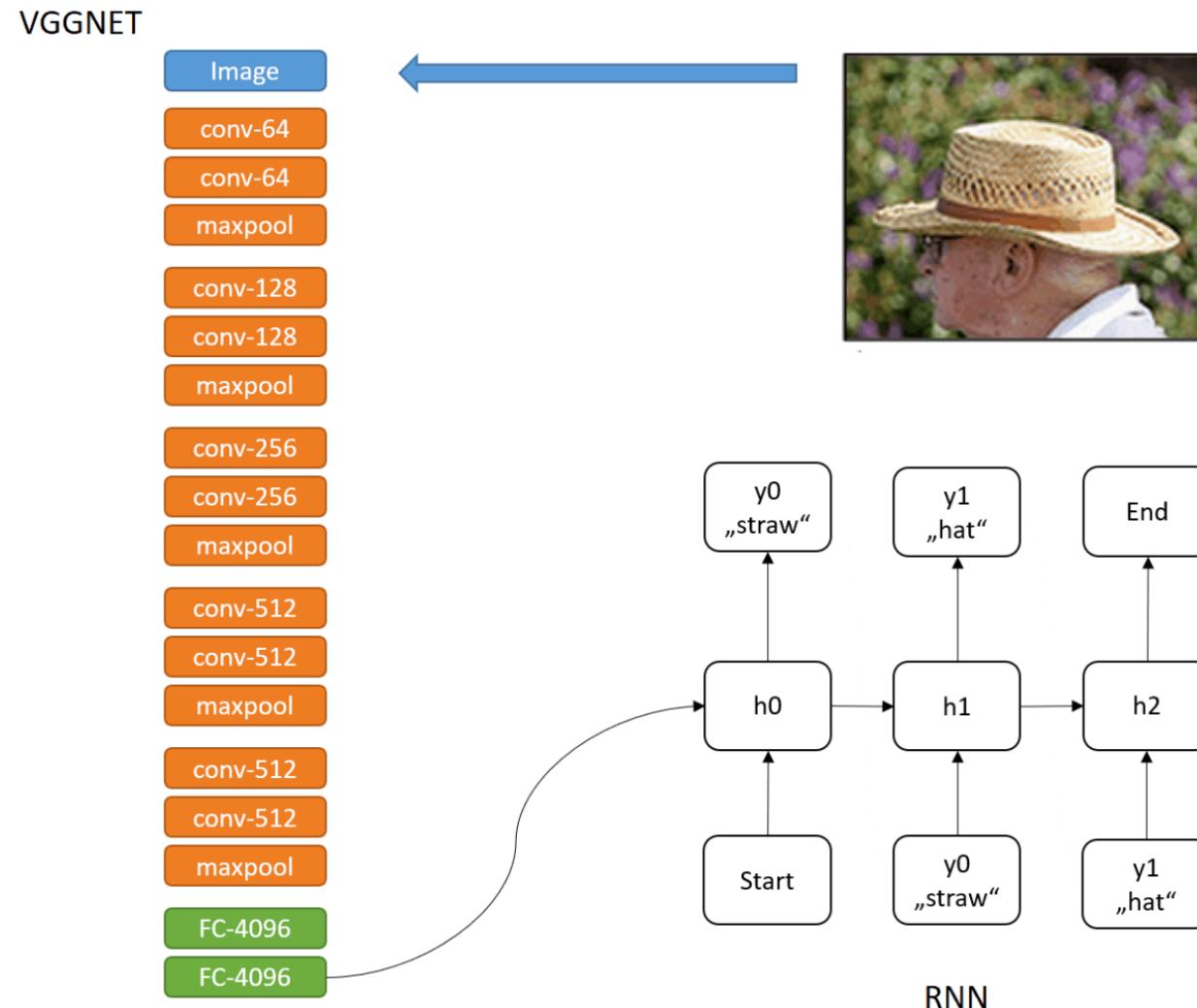


Image captioning (CNN encoder – RNN decoder)



Пример: распознавание рукописного текста

Выход

Скрытый слой

Вход



Ссылки

- <http://colah.github.io/posts/2015-08-Understanding-LSTMs/>
- <http://www.wildml.com/2015/09/recurrent-neural-networks-tutorial-part-1-introduction-to-rnns/>
- http://www.fit.vutbr.cz/research/groups/speech/publi/2010/mikolov_interspeech2010_IS100722.pdf
- <http://distill.pub/2016/augmented-rnns/>
- Multilingual Neural Machine Translation <https://arxiv.org/abs/1611.04558>
- <https://deepmind.com/blog/wavenet-generative-model-raw-audio/>
- <https://www.tensorflow.org/tutorials/recurrent>
- <https://www.tensorflow.org/tutorials/seq2seq>