



Deep Learning School

# Рекуррентная нейронная сеть

# План занятия

- Рекуррентный слой: идея.  
Рекуррентная нейросеть (RNN);
- Forward pass RNN;
- Обучение RNN (backward pass);
- Функции активации в RNN;
- Bidirectional RNN;
- GRU, LSTM

# План занятия

- **Рекуррентный слой: идея.**  
**Рекуррентная нейросеть (RNN);**
- **Forward pass RNN;**
- Обучение RNN (backward pass);
- Функции активации в RNN;
- Bidirectional RNN;
- GRU, LSTM

# Особенность текста и звука

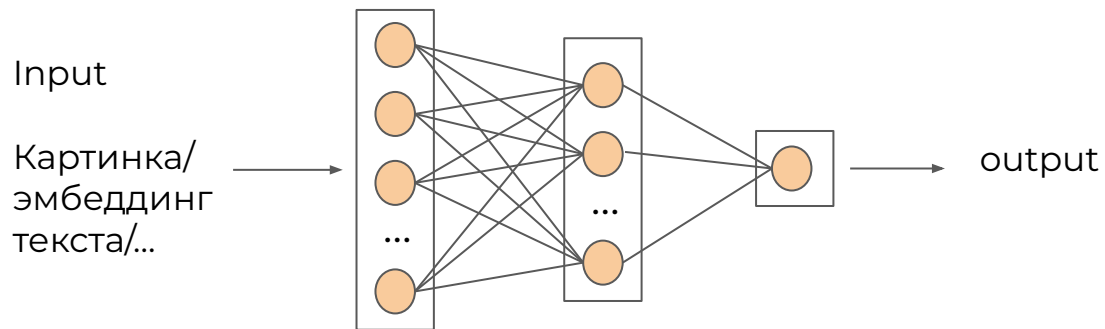
Отличие текста и звука от других типов данных (например, изображений) состоит в наличии временной компоненты.

Мы читаем текст не моментально, а слово за словом, в строго определенном порядке.

Возникает идея придумать идею нейросети, которая учитывала бы эту особенность этих типов данных.

# Рекуррентный слой

Как работает обычная нейросеть:



# Рекуррентный слой

Как работает рекуррентная нейросеть:

Рекуррентная нейросеть  
обрабатывает один токен текста за  
один момент времени.

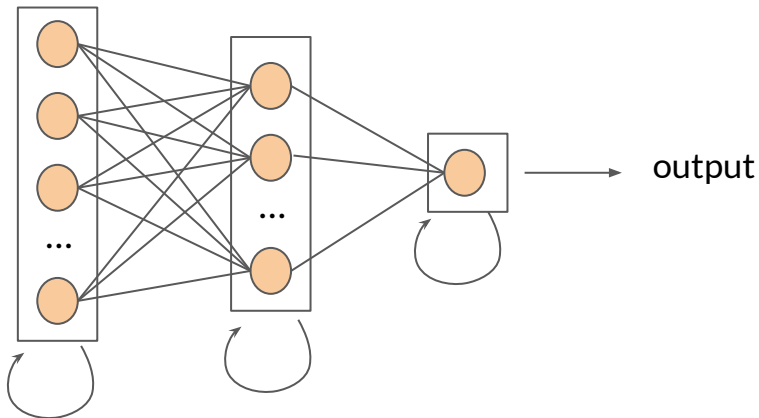
К слою добавляется связь “из себя в  
себя” — “память” слоя

Input

a  
cat  
is  
sitting  
on  
the  
mat

word2vec

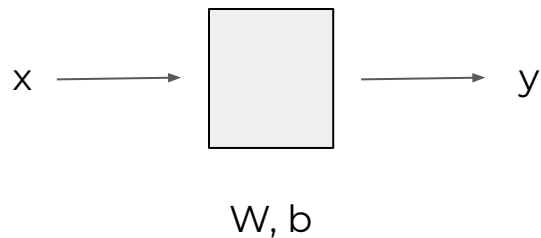
$\begin{pmatrix} 0.45 \\ -1.34 \\ 2.34 \\ \dots \\ -0.45 \end{pmatrix}$



# Рекуррентный слой

Полносвязный слой

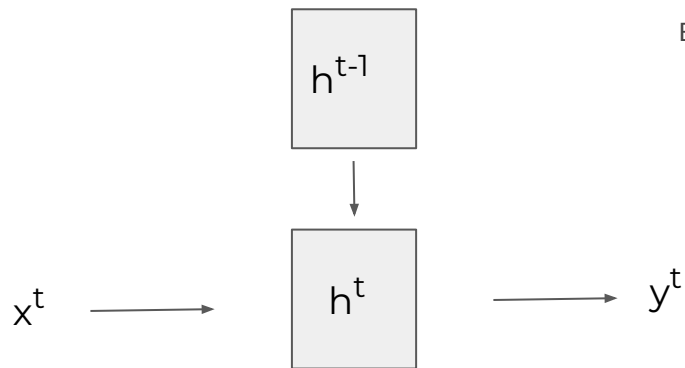
Вычисление выхода слоя:



$$y = \sigma(WX + b)$$

# Рекуррентный слой

Рекуррентный слой



Обновление вектора  
скрытого состояния и  
вычисление выхода слоя:

$$h^t = \sigma(WX^t + Uh^{t-1} + b_h)$$

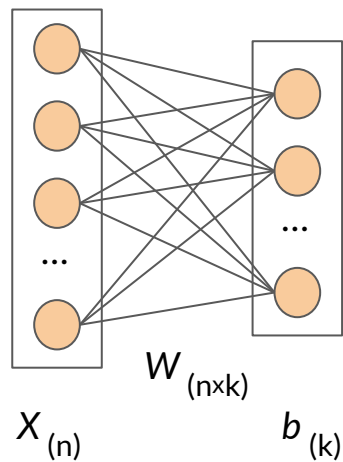
$$y = \sigma(Vh^t + b_y)$$

$W, U, V$   
 $b_h, b_y$



# Рекуррентный слой

Слой полносвязной нейросети



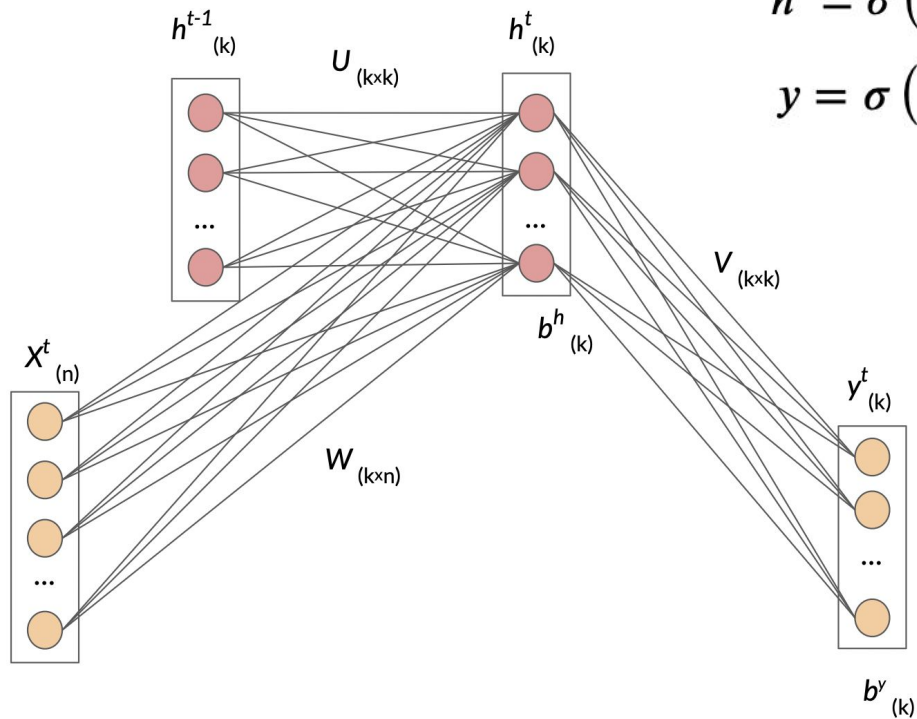
$$y = \sigma(WX + b)$$

## Слой рекуррентной нейросети

Обновление вектора  
скрытого состояния и  
вычисление выхода слоя:

$$h^t = \sigma(WX^t + Uh^{t-1} + b_h)$$

$$y = \sigma(Vh^t + b_y)$$



a cat is sitting  
on the mat

$$h^0_1$$

$$h^0_2$$

...

$$h^0_n$$

$$a \begin{pmatrix} 0.45 \\ -1.34 \\ 2.34 \\ \dots \\ -0.45 \end{pmatrix} \chi^1 \longrightarrow h^1_1$$

время

a cat is sitting  
on the mat

$$h^0_1$$

$$h^0_2$$

...

$$h^0_n$$

$$a \begin{pmatrix} 0.45 \\ -1.34 \\ 2.34 \\ \dots \\ -0.45 \end{pmatrix} \chi^1 \longrightarrow h^1_1 \xrightarrow{y^1_1}$$

время

a cat is sitting  
on the mat

$$h^0_1$$

$$h^0_2$$

...

$$h^0_n$$

$$a \begin{pmatrix} 0.45 \\ -1.34 \\ 2.34 \\ \dots \\ -0.45 \end{pmatrix}$$

$\chi^1$

$$h^1_1$$

$y^1_1$

$$h^1_2$$

время

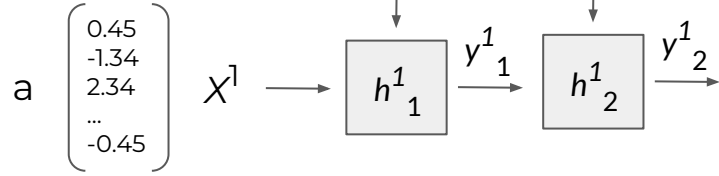
a cat is sitting  
on the mat

$$h^0_1$$

$$h^0_2$$

...

$$h^0_n$$



время



a cat is sitting  
on the mat

$$h^0_1$$

$$h^0_2$$

...

$$h^0_n$$

$$a \begin{pmatrix} 0.45 \\ -1.34 \\ 2.34 \\ \dots \\ -0.45 \end{pmatrix}$$

$\chi^1$

$$h^1_1$$

$y^1_1$

$$h^1_2$$

$y^1_2$

...

$$h^1_n$$

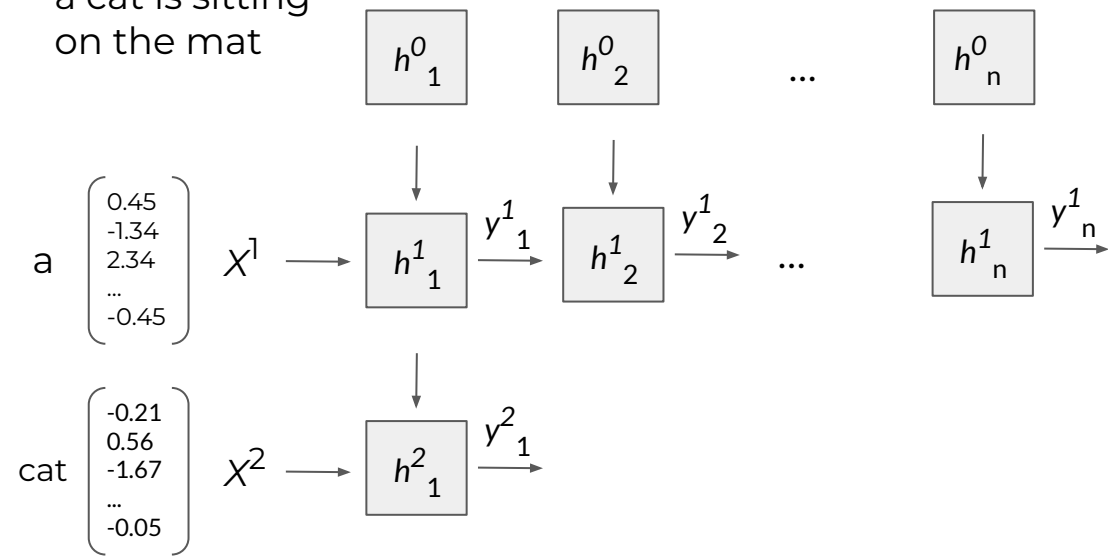
$y^1_n$

$$\text{cat} \begin{pmatrix} -0.21 \\ 0.56 \\ -1.67 \\ \dots \\ -0.05 \end{pmatrix}$$

$\chi^2$

время

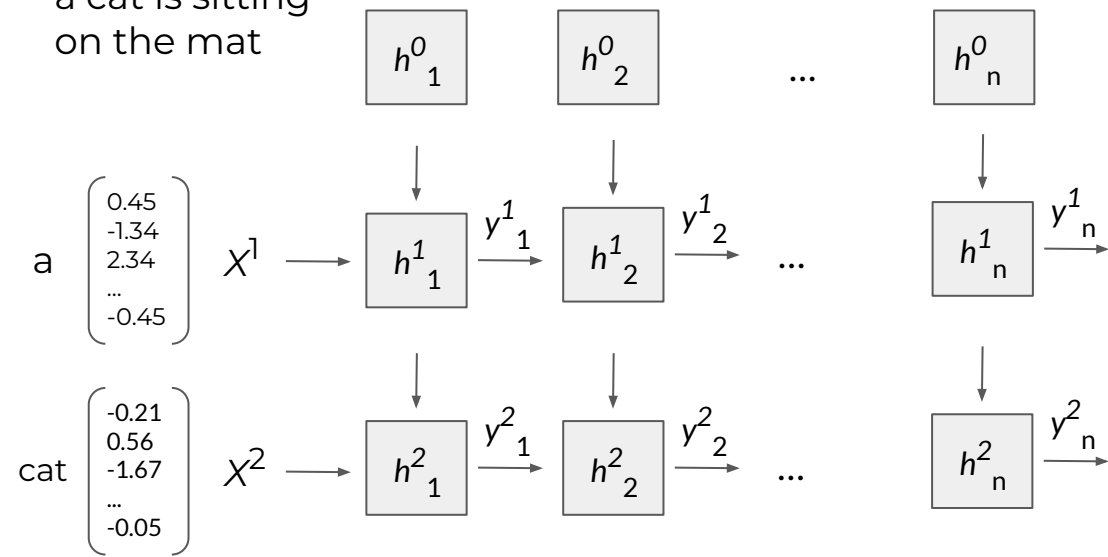
a cat is sitting  
on the mat



время

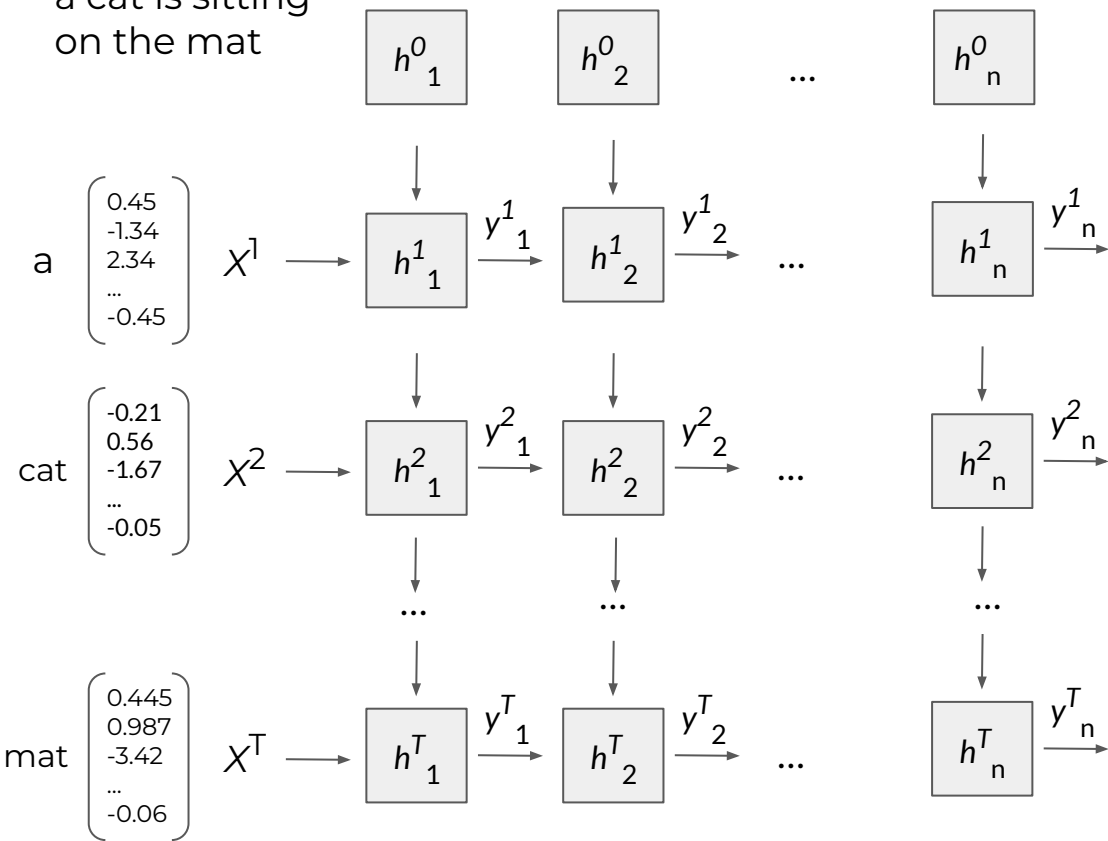


a cat is sitting  
on the mat



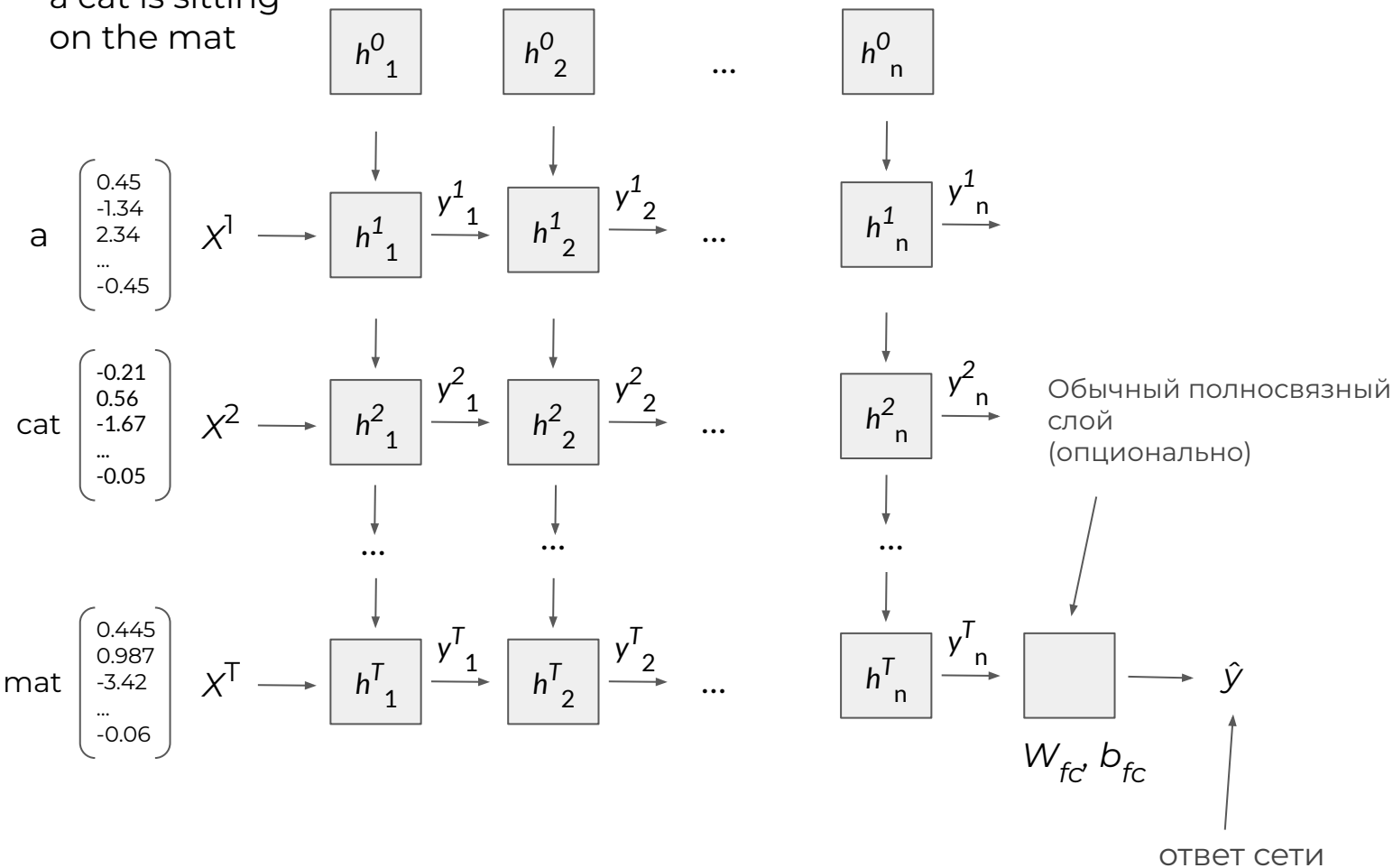
время

a cat is sitting  
on the mat



время

a cat is sitting  
on the mat



время

# Итоги видео

В этом видео мы:

- Познакомились с идеей устройства рекуррентного слоя и рекуррентной нейросети;
- Разобрали forward pass рекуррентной сети.

В следующем видео мы узнаем, как RNN обучается.

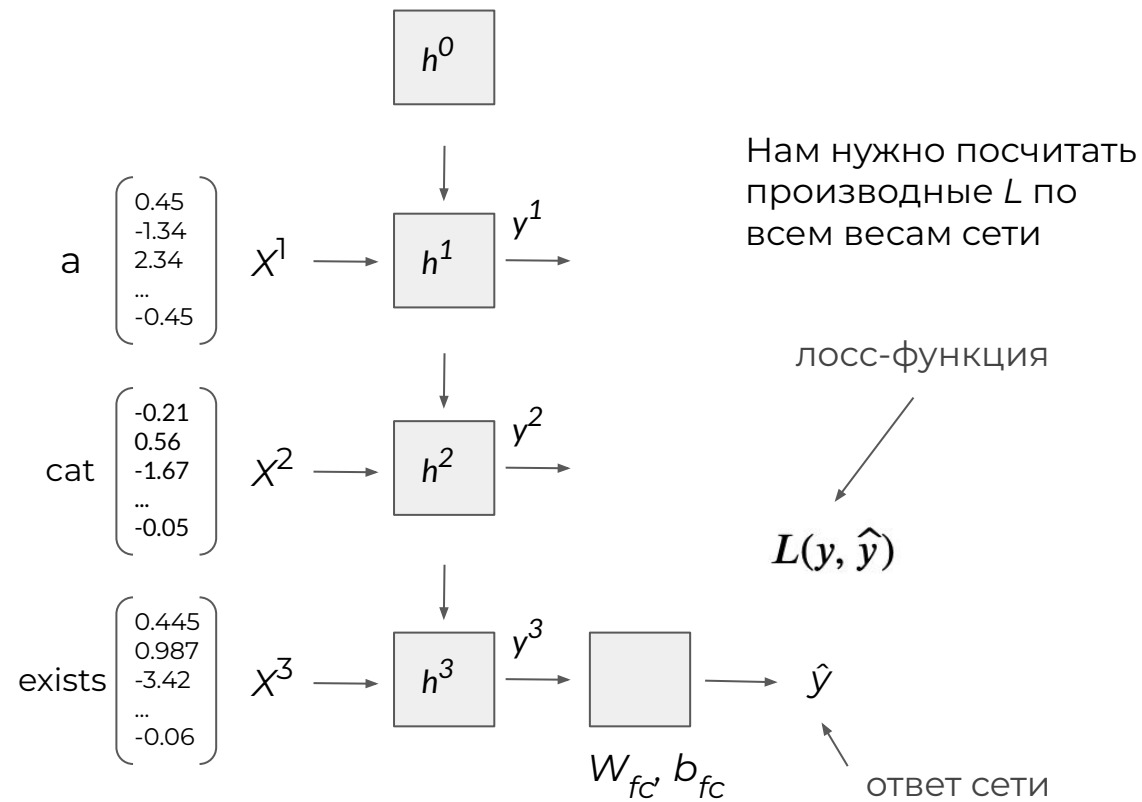


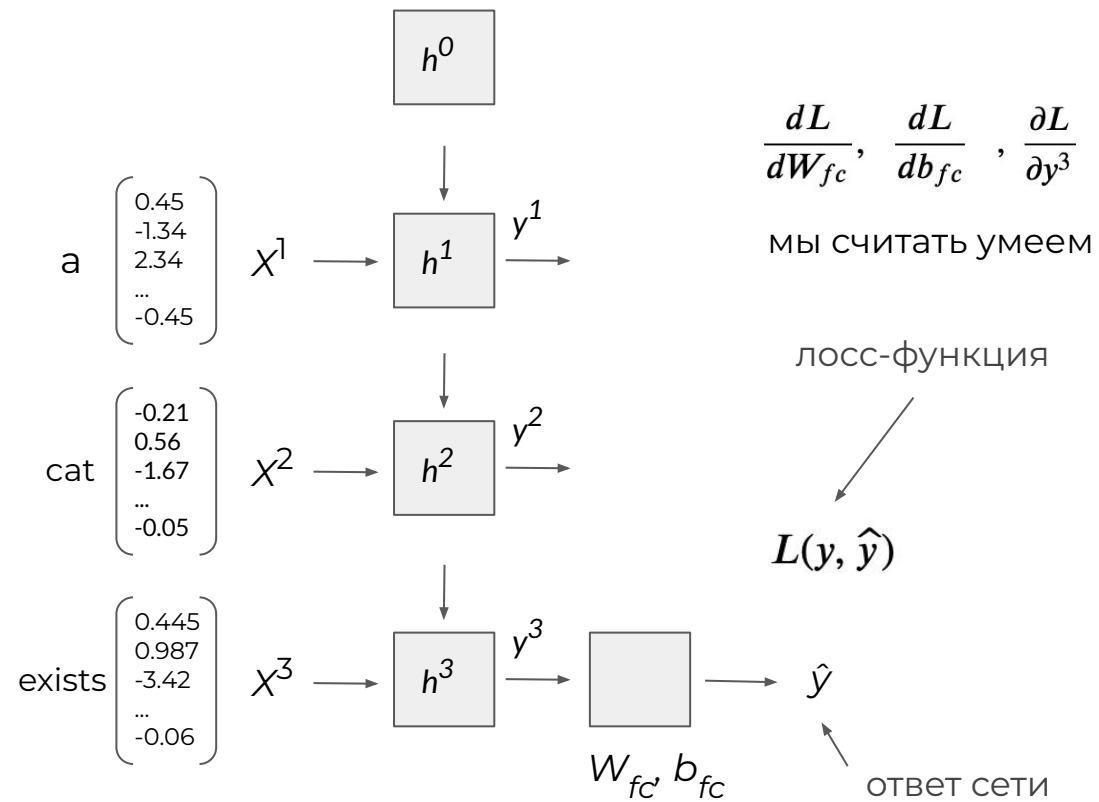
Deep Learning School

# Обучение RNN

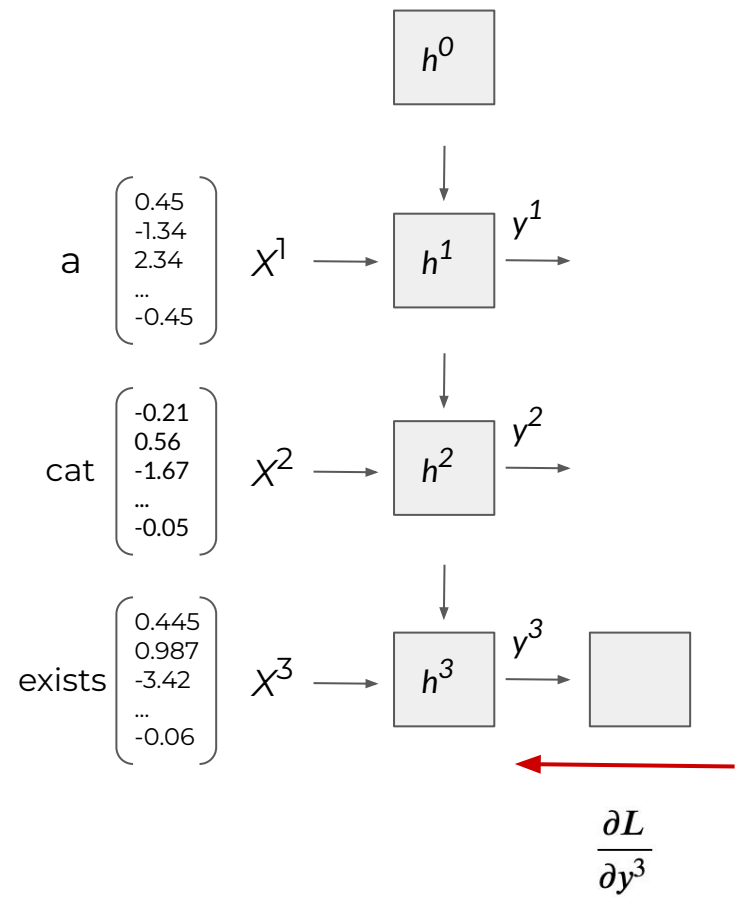
# План занятия

- Рекуррентный слой: идея.  
Рекуррентная нейросеть (RNN);
- Forward pass RNN;
- **Обучение RNN (backward pass);**
- Функции активации в RNN;
- Bidirectional RNN;
- GRU, LSTM







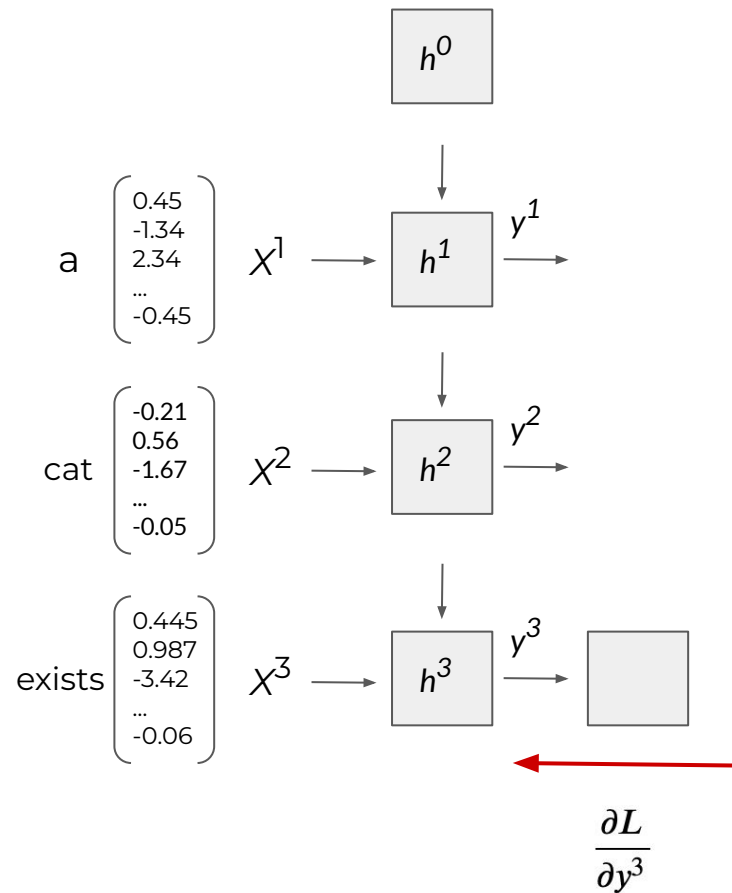


Обучаемые параметры слоя:

$$W, U, V, b_h, b_y$$

$$h^t = \sigma(WX^t + Uh^{t-1} + b_h)$$

$$y^t = \sigma(Vh^t + b_y)$$

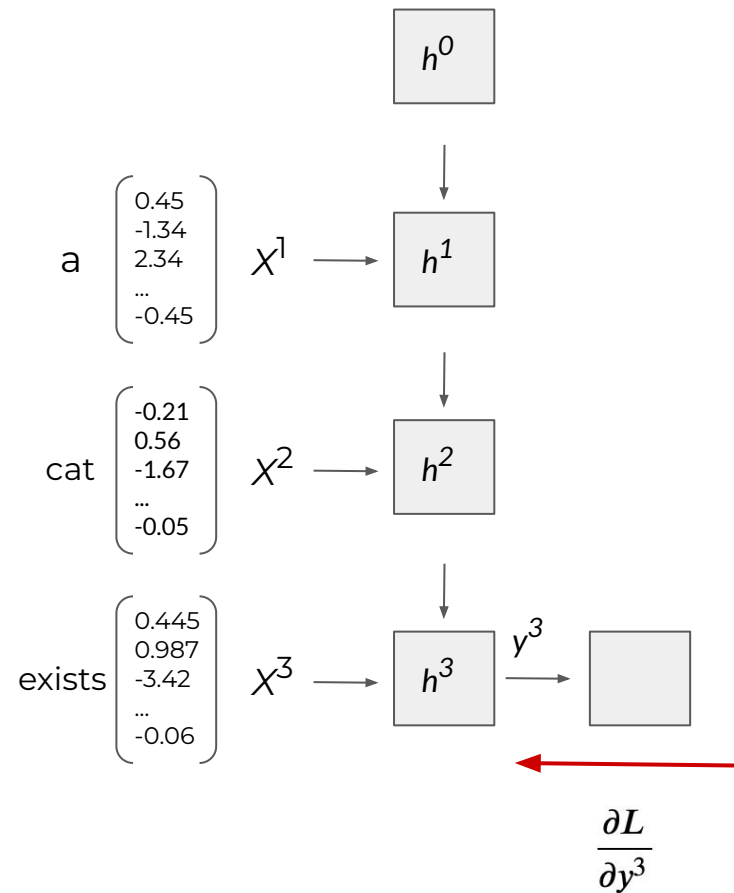


Обучаемые параметры слоя:

$$W, U, V, b_h, b_y$$

$$h^t = \sigma(WX^t + Uh^{t-1} + b_h)$$

$$y^t = \sigma(Vh^t + b_y)$$



Обучаемые параметры слоя:

$$W, U, V, b_h, b_y$$

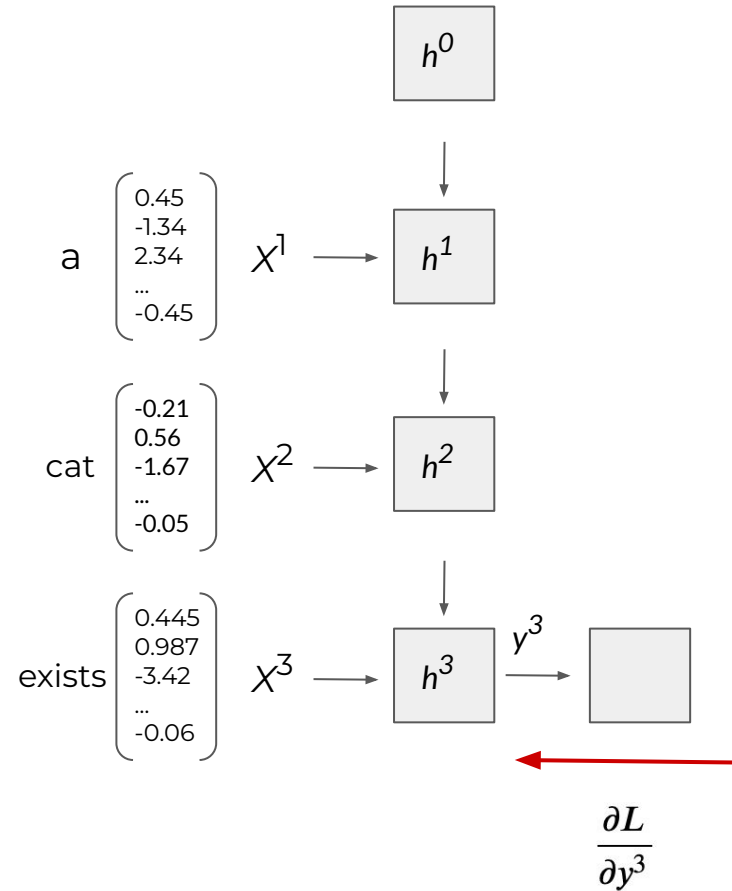
Для  $t = 1, 2$ :

$$h^t = \sigma(WX^t + Uh^{t-1} + b_h)$$

Для  $t = 3$ :

$$h^3 = \sigma(WX^3 + Uh^2 + b_h)$$

$$y^3 = \sigma(Vh^3 + b_y)$$



Обучаемые параметры слоя:

$$W, U, V, b_h, b_y$$

Для  $t = 1, 2$ :

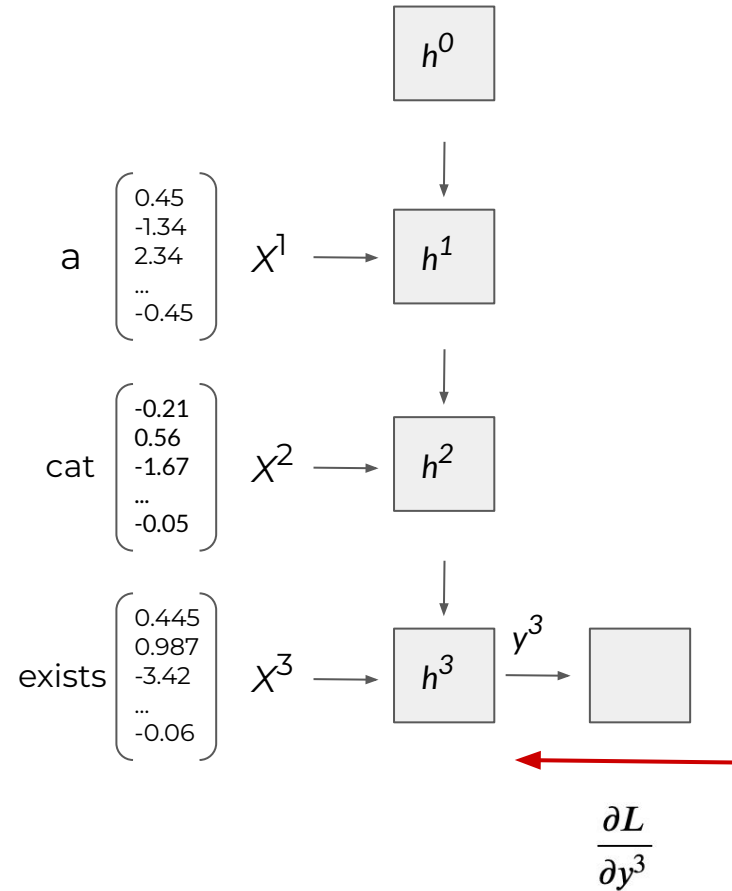
$$h^t = \sigma(WX^t + Uh^{t-1} + b_h)$$

Для  $t = 3$ :

$$h^3 = \sigma(WX^3 + Uh^2 + b_h)$$

$$y^3 = \sigma(Vh^3 + b_y)$$

$$\frac{dL}{dV} = \frac{\partial L}{\partial y^3} \frac{\partial y^3}{\partial V}, \quad \frac{dL}{db_y} = \frac{\partial L}{\partial y^3} \frac{\partial y^3}{\partial b_y}$$



Обучаемые параметры слоя:

$$W, U, V, b_h, b_y$$

Для  $t = 1, 2$ :

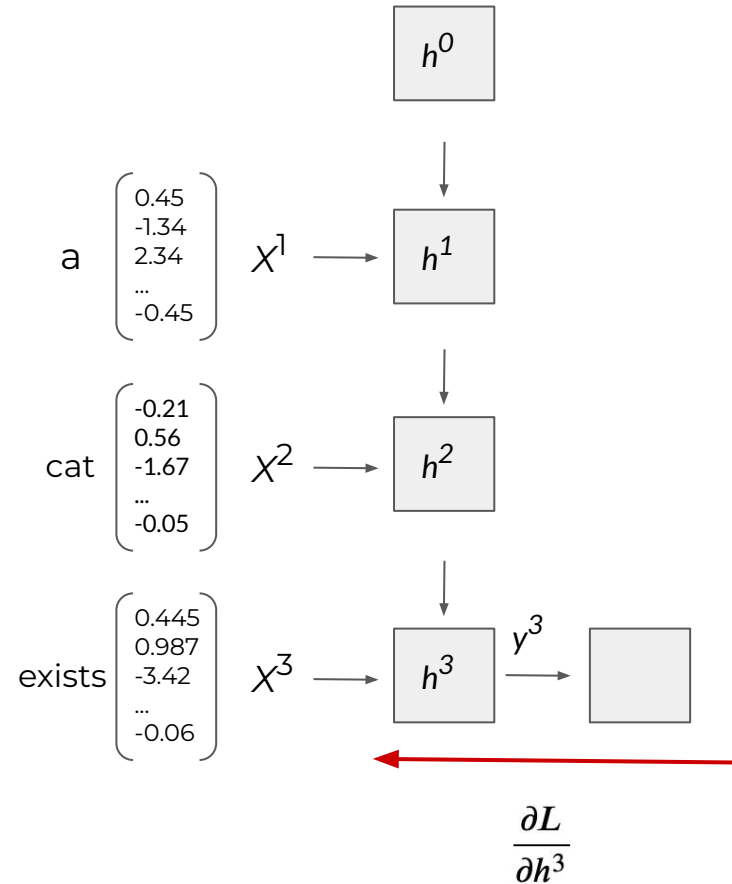
$$h^t = \sigma(WX^t + Uh^{t-1} + b_h)$$

Для  $t = 3$ :

$$h^3 = \sigma(WX^3 + Uh^2 + b_h)$$

$$y^3 = \sigma(Vh^3 + b_y)$$

$$\frac{dL}{dW} = \frac{\partial L}{\partial y^3} \frac{dy^3}{dW} = \frac{\partial L}{\partial y^3} \frac{\partial y^3}{\partial h^3} \frac{dh^3}{dW} = \frac{\partial L}{\partial h^3} \frac{dh^3}{dW}$$

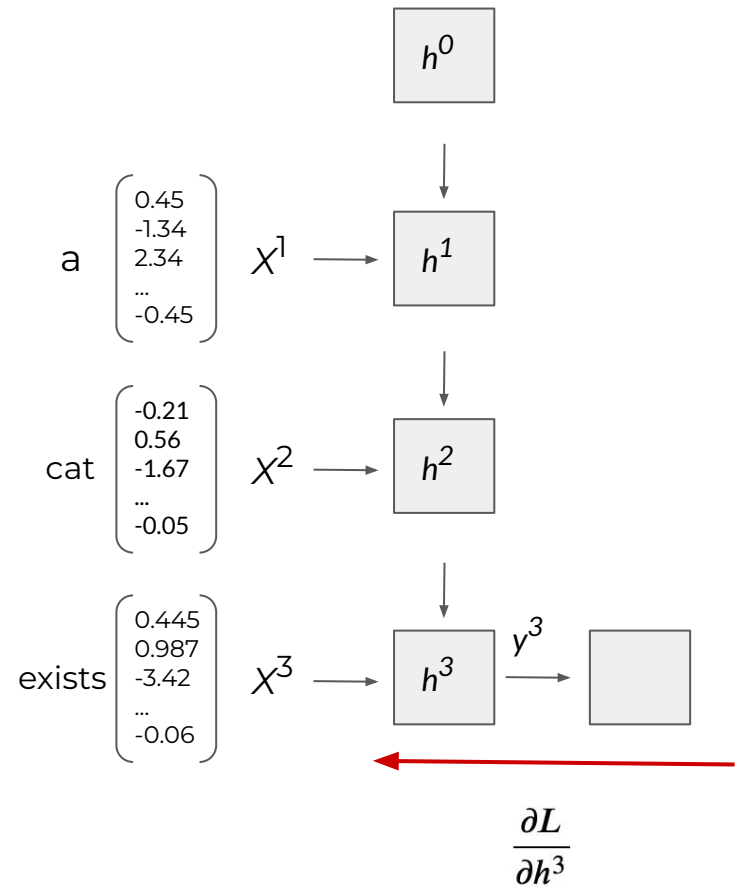


$$h^t = \sigma(WX^t + Uh^{t-1} + b_h)$$

$$\frac{dL}{dW} = \frac{\partial L}{\partial h^3} \boxed{\frac{dh^3}{dW}}$$

||

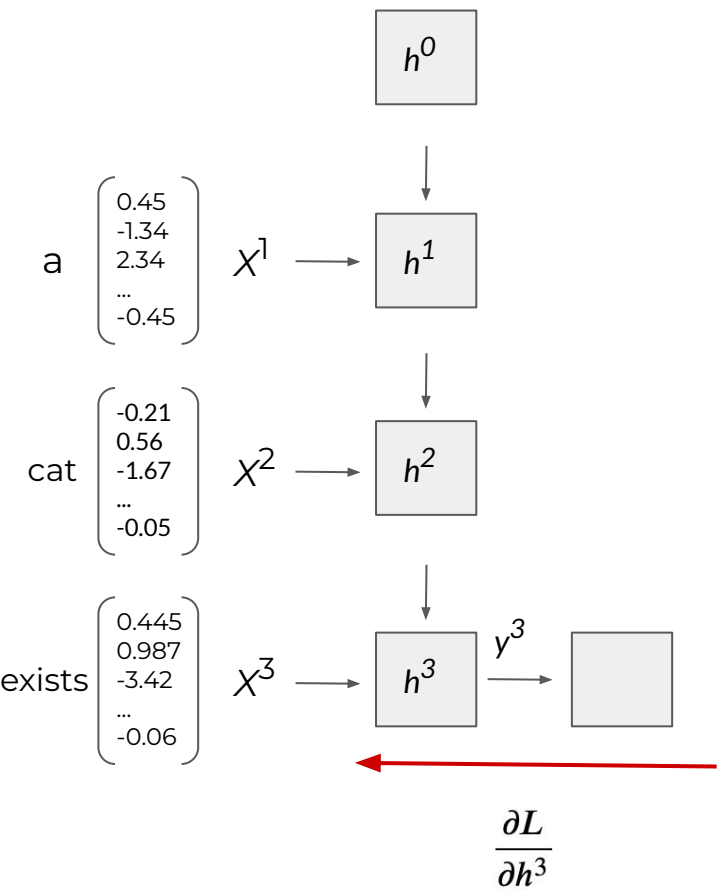
$$\frac{\partial h^3}{\partial W} + \frac{\partial h^3}{\partial h^2} \frac{dh^2}{dW}$$



$$h^t = \sigma(WX^t + Uh^{t-1} + b_h)$$

$$\frac{dL}{dW} = \frac{\partial L}{\partial h^3} \frac{dh^3}{dW} =$$

$$= \frac{\partial L}{\partial h^3} \frac{\partial h^3}{\partial W} + \frac{\partial L}{\partial h^3} \frac{\partial h^3}{\partial h^2} \frac{dh^2}{dW}$$

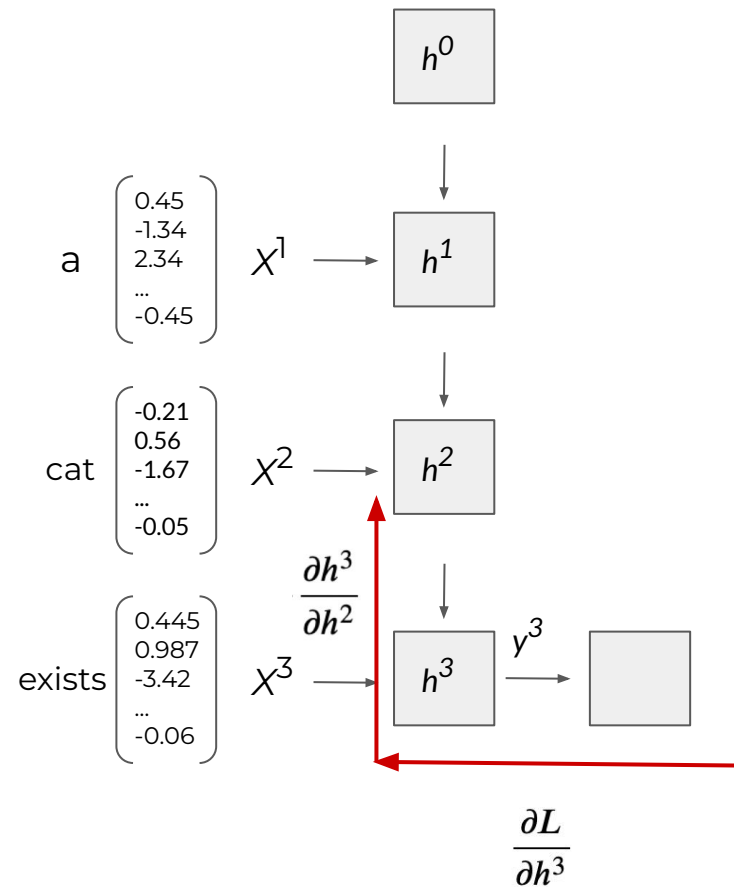




$$h^t = \sigma(WX^t + Uh^{t-1} + b_h)$$

$$\frac{dL}{dW} = \frac{\partial L}{\partial h^3} \frac{dh^3}{dW} =$$

$$= \frac{\partial L}{\partial h^3} \frac{\partial h^3}{\partial W} + \frac{\partial L}{\partial h^3} \boxed{\frac{\partial h^3}{\partial h^2}} \frac{dh^2}{dW}$$



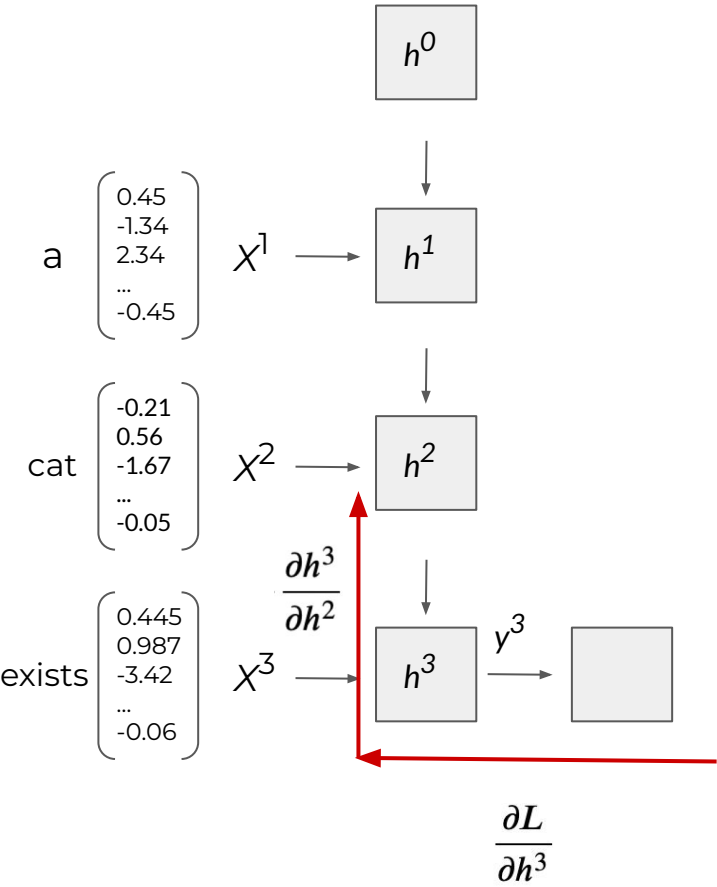
$$h^t = \sigma(WX^t + Uh^{t-1} + b_h)$$

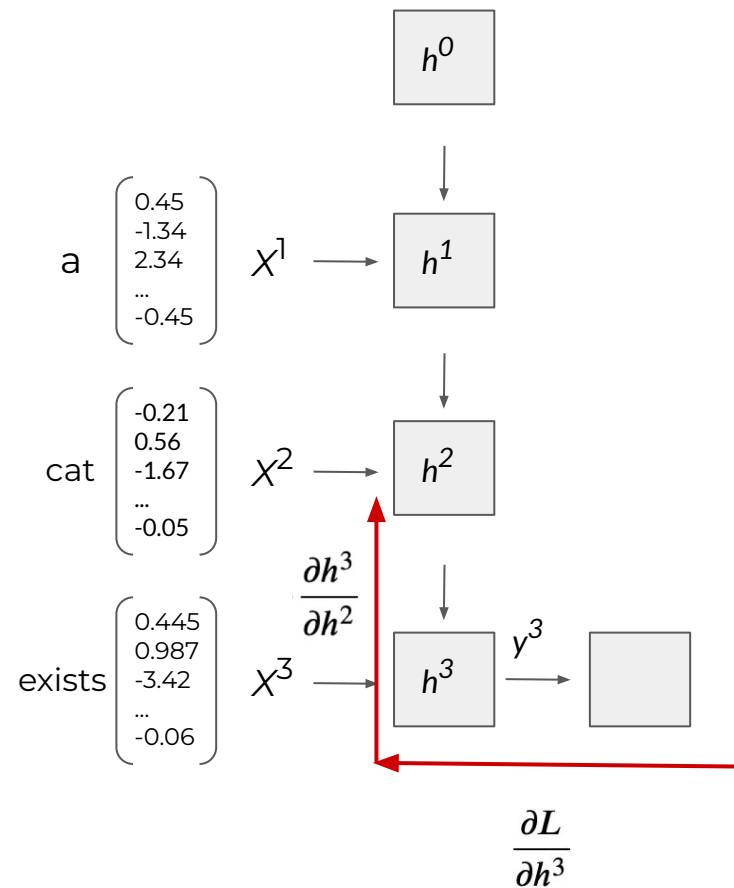
$$\frac{dL}{dW} = \frac{\partial L}{\partial h^3} \frac{dh^3}{dW} =$$

$$= \frac{\partial L}{\partial h^3} \frac{\partial h^3}{\partial W} + \frac{\partial L}{\partial h^3} \frac{\partial h^3}{\partial h^2} \boxed{\frac{dh^2}{dW}}$$

||

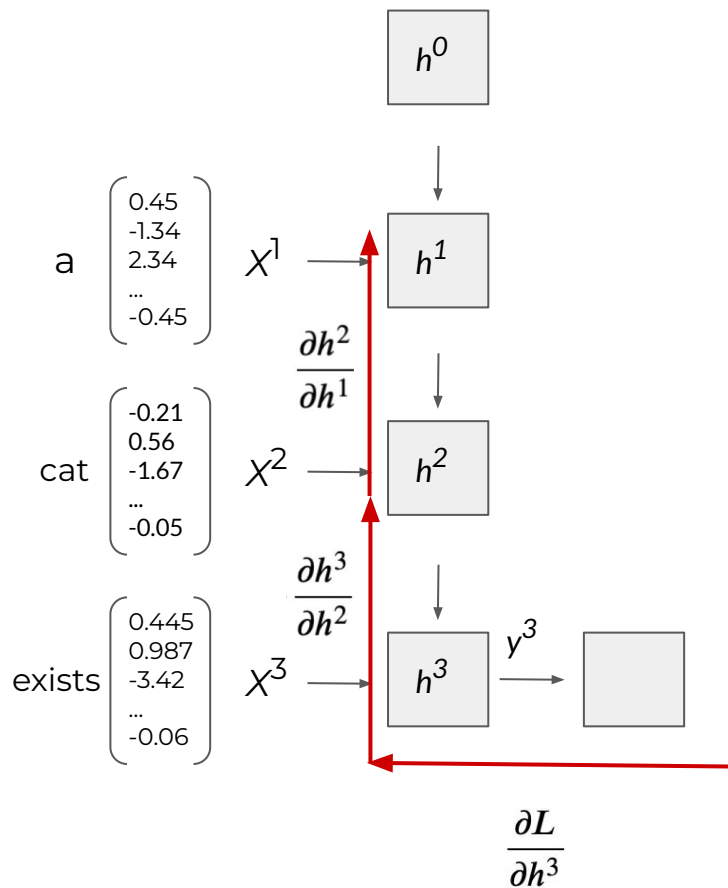
$$\frac{\partial h^2}{\partial W} + \frac{\partial h^2}{\partial h^1} \frac{dh^1}{dW}$$





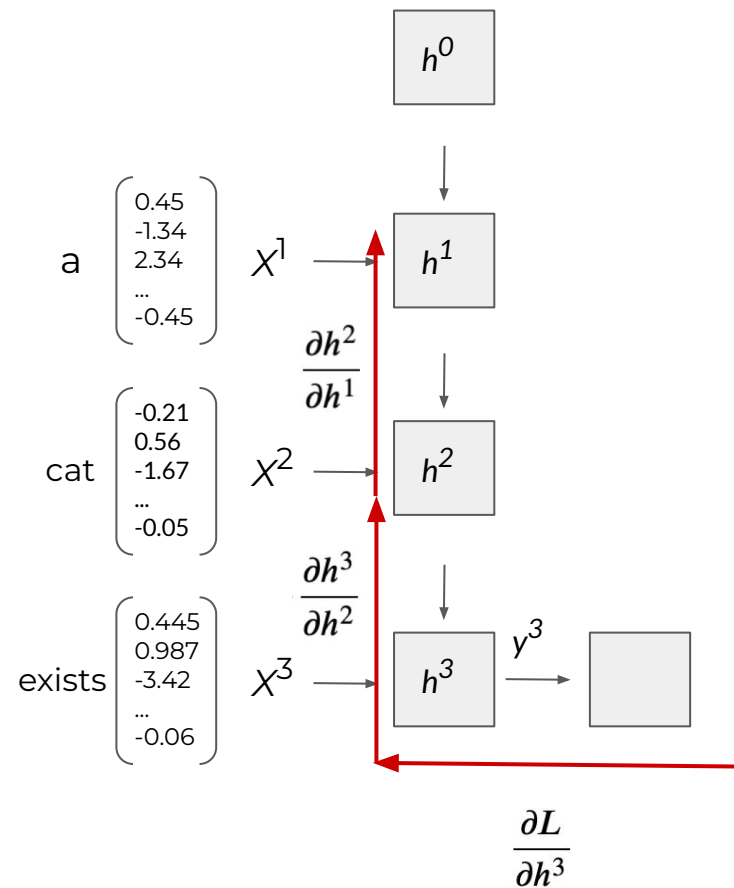
$$h^t = \sigma(WX^t + Uh^{t-1} + b_h)$$

$$\begin{aligned} \frac{dL}{dW} &= \frac{\partial L}{\partial h^3} \frac{dh^3}{dW} = \\ &= \frac{\partial L}{\partial h^3} \frac{\partial h^3}{\partial W} + \frac{\partial L}{\partial h^3} \frac{\partial h^3}{\partial h^2} \frac{dh^2}{dW} = \\ &= \frac{\partial L}{\partial h^3} \frac{\partial h^3}{\partial W} + \frac{\partial L}{\partial h^3} \frac{\partial h^3}{\partial h^2} \frac{\partial h^2}{\partial W} + \\ &+ \frac{\partial L}{\partial h^3} \frac{\partial h^3}{\partial h^2} \frac{\partial h^2}{\partial h^1} \frac{dh^1}{dW} \end{aligned}$$



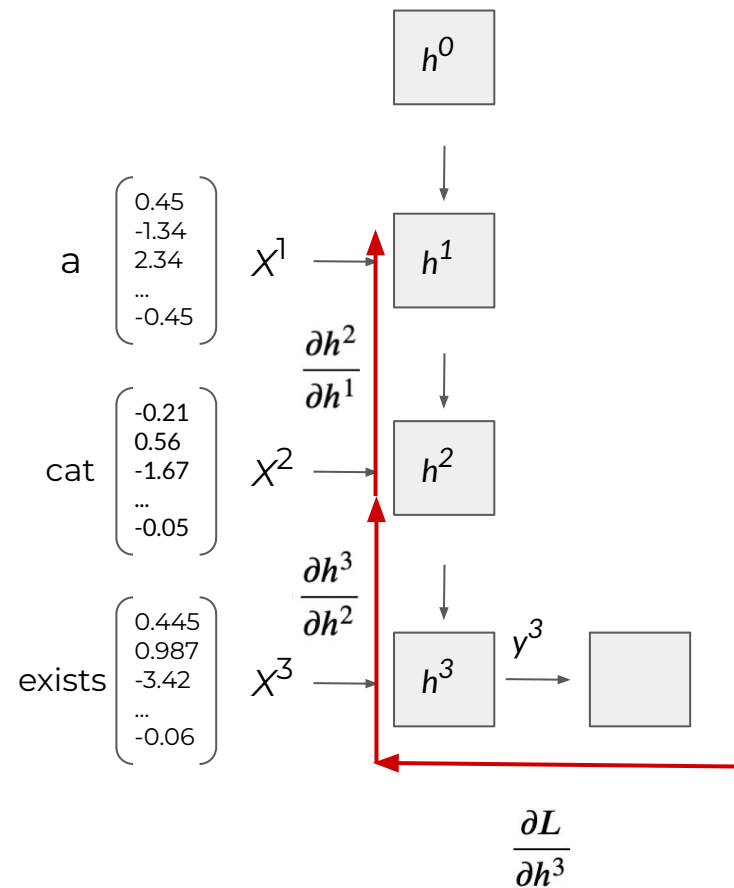
$$h^t = \sigma(WX^t + Uh^{t-1} + b_h)$$

$$\begin{aligned} \frac{dL}{dW} &= \frac{\partial L}{\partial h^3} \frac{dh^3}{dW} = \\ &= \frac{\partial L}{\partial h^3} \frac{\partial h^3}{\partial W} + \frac{\partial L}{\partial h^3} \frac{\partial h^3}{\partial h^2} \frac{dh^2}{dW} = \\ &= \frac{\partial L}{\partial h^3} \frac{\partial h^3}{\partial W} + \frac{\partial L}{\partial h^3} \frac{\partial h^3}{\partial h^2} \frac{\partial h^2}{\partial W} + \\ &+ \frac{\partial L}{\partial h^3} \frac{\partial h^3}{\partial h^2} \boxed{\frac{\partial h^2}{\partial h^1}} \frac{dh^1}{dW} \end{aligned}$$



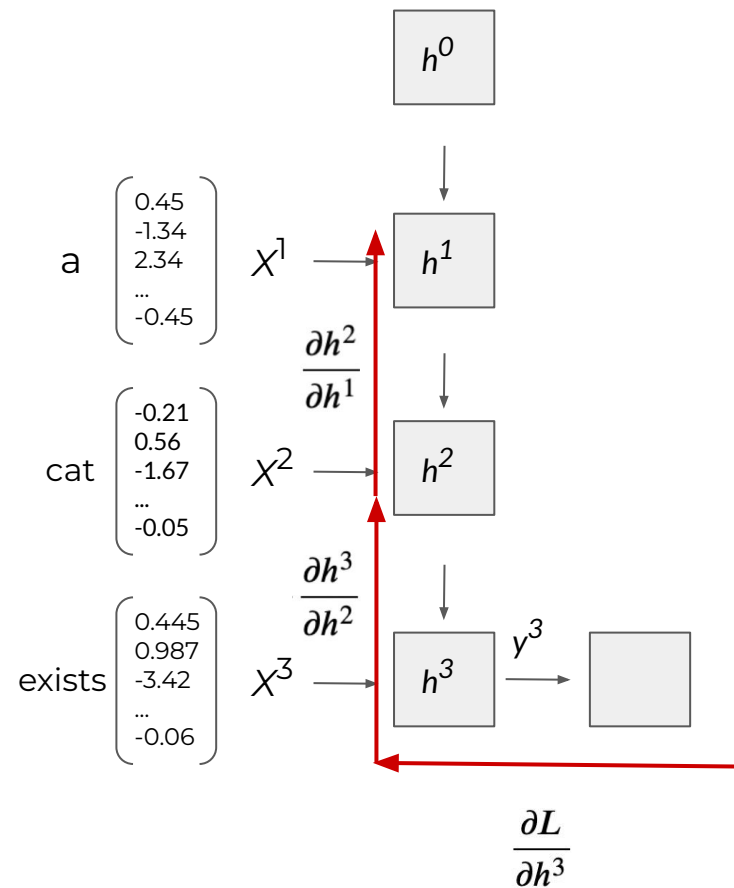
$$h^t = \sigma(WX^t + Uh^{t-1} + b_h)$$

$$\begin{aligned} \frac{dL}{dW} &= \frac{\partial L}{\partial h^3} \frac{dh^3}{dW} = \\ &= \frac{\partial L}{\partial h^3} \frac{\partial h^3}{\partial W} + \frac{\partial L}{\partial h^3} \frac{\partial h^3}{\partial h^2} \frac{dh^2}{dW} = \\ &= \frac{\partial L}{\partial h^3} \frac{\partial h^3}{\partial W} + \frac{\partial L}{\partial h^3} \frac{\partial h^3}{\partial h^2} \frac{\partial h^2}{\partial W} + \\ &+ \frac{\partial L}{\partial h^3} \frac{\partial h^3}{\partial h^2} \frac{\partial h^2}{\partial h^1} \boxed{\frac{dh^1}{dW}} \end{aligned}$$



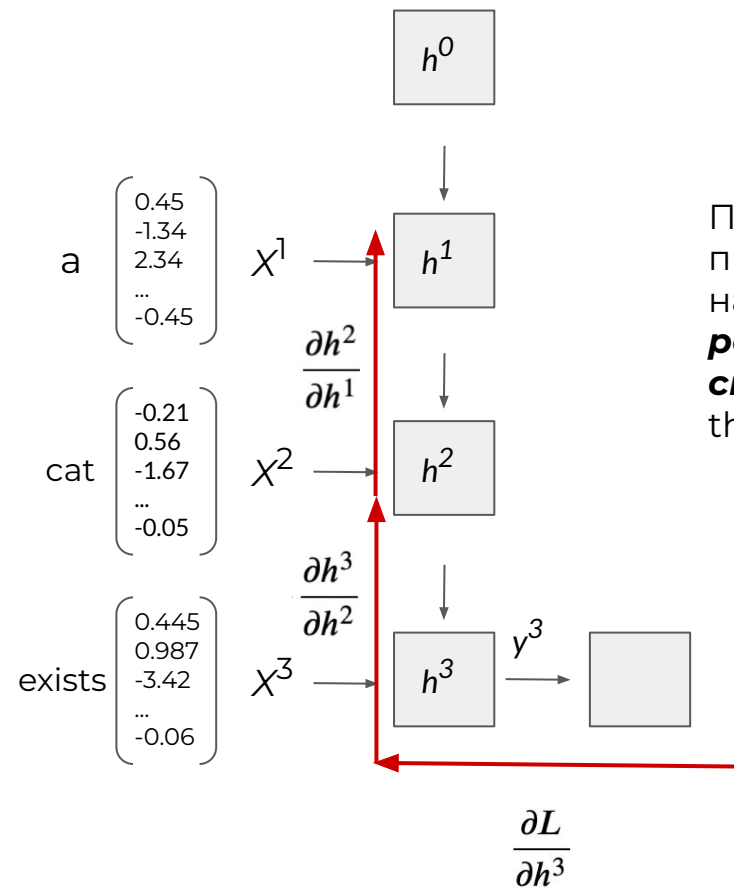
$$h^t = \sigma(WX^t + Uh^{t-1} + b_h)$$

$$\begin{aligned} \frac{dL}{dW} &= \frac{\partial L}{\partial h^3} \frac{dh^3}{dW} = \\ &= \frac{\partial L}{\partial h^3} \frac{\partial h^3}{\partial W} + \frac{\partial L}{\partial h^3} \frac{\partial h^3}{\partial h^2} \frac{dh^2}{dW} = \\ &= \frac{\partial L}{\partial h^3} \frac{\partial h^3}{\partial W} + \frac{\partial L}{\partial h^3} \frac{\partial h^3}{\partial h^2} \frac{\partial h^2}{\partial W} + \\ &+ \frac{\partial L}{\partial h^3} \frac{\partial h^3}{\partial h^2} \frac{\partial h^2}{\partial h^1} \boxed{\frac{dh^1}{dW}} \\ &\parallel \\ &\frac{\partial h^1}{\partial W} \end{aligned}$$



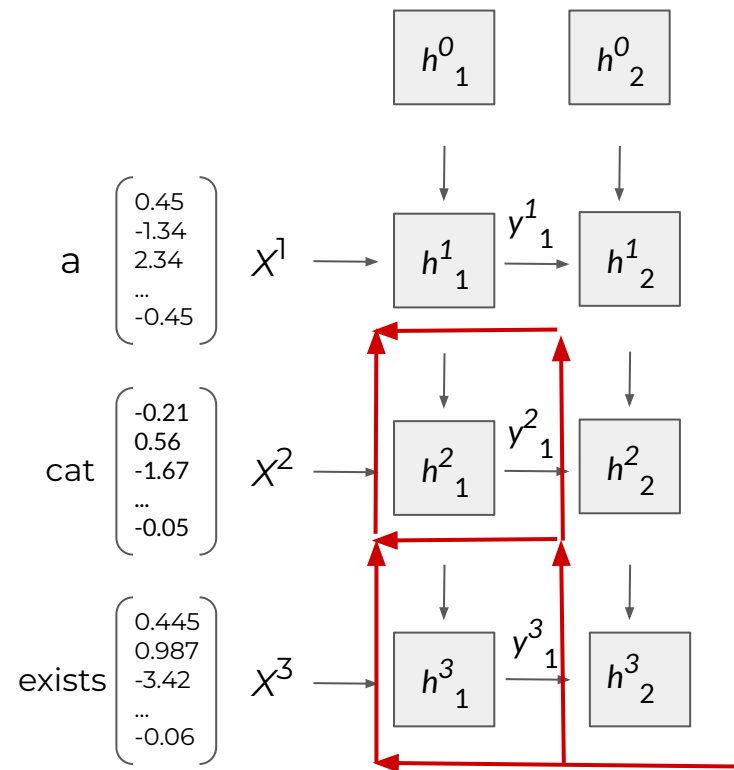
$$h^t = \sigma(WX^t + Uh^{t-1} + b_h)$$

$$\begin{aligned} \frac{dL}{dW} &= \frac{\partial L}{\partial h^3} \frac{dh^3}{dW} = \\ &= \frac{\partial L}{\partial h^3} \frac{\partial h^3}{\partial W} + \frac{\partial L}{\partial h^3} \frac{\partial h^3}{\partial h^2} \frac{dh^2}{dW} = \\ &= \frac{\partial L}{\partial h^3} \frac{\partial h^3}{\partial W} + \frac{\partial L}{\partial h^3} \frac{\partial h^3}{\partial h^2} \frac{\partial h^2}{\partial W} + \\ &+ \frac{\partial L}{\partial h^3} \frac{\partial h^3}{\partial h^2} \frac{\partial h^2}{\partial h^1} \frac{\partial h^1}{\partial W} \end{aligned}$$



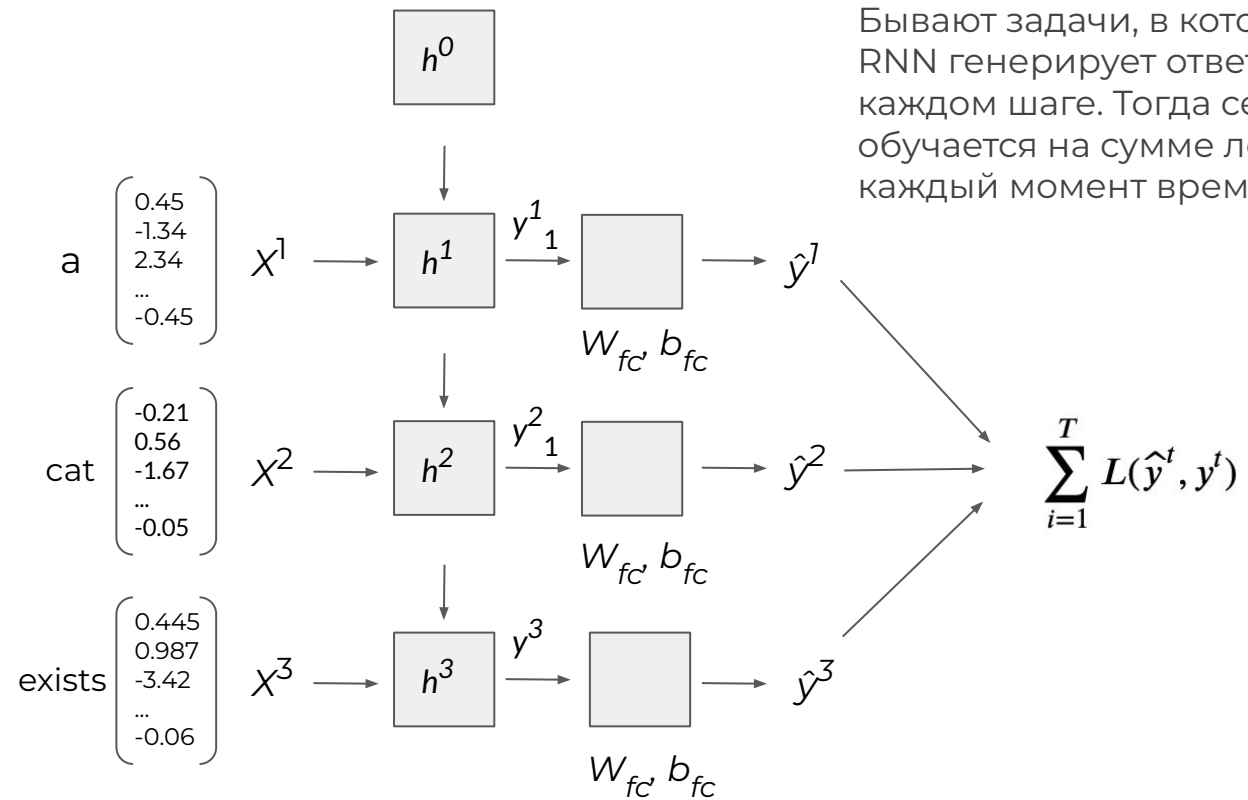
Процесс вычисления производных весов RNN называется **обратное распространение ошибки сквозь время** (backpropagation through time)





При добавлении слоев в RNN подсчет градиентов усложняется

Бывают задачи, в которых RNN генерирует ответ на каждом шаге. Тогда сеть обучается на сумме потерь в каждый момент времени



# Итоги видео

В этом видео мы обсудили работу алгоритма обновления весов нейросети: обратное распространение ошибки сквозь время.

В следующем видео мы обсудим несколько нюансов RNN.

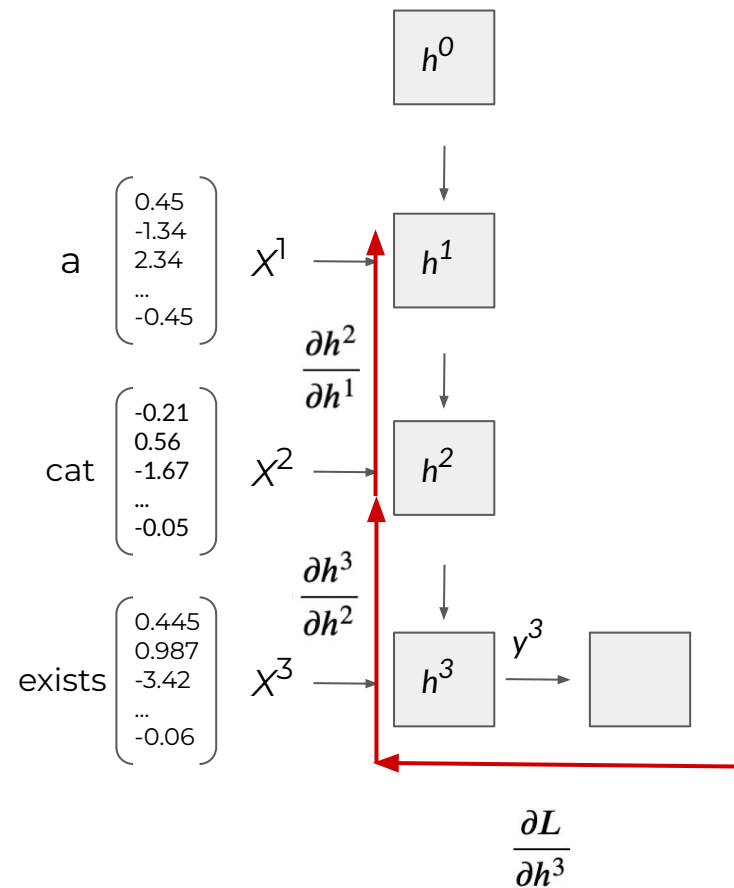


Deep Learning School

# Функции активации в RNN. Bidirectional RNN

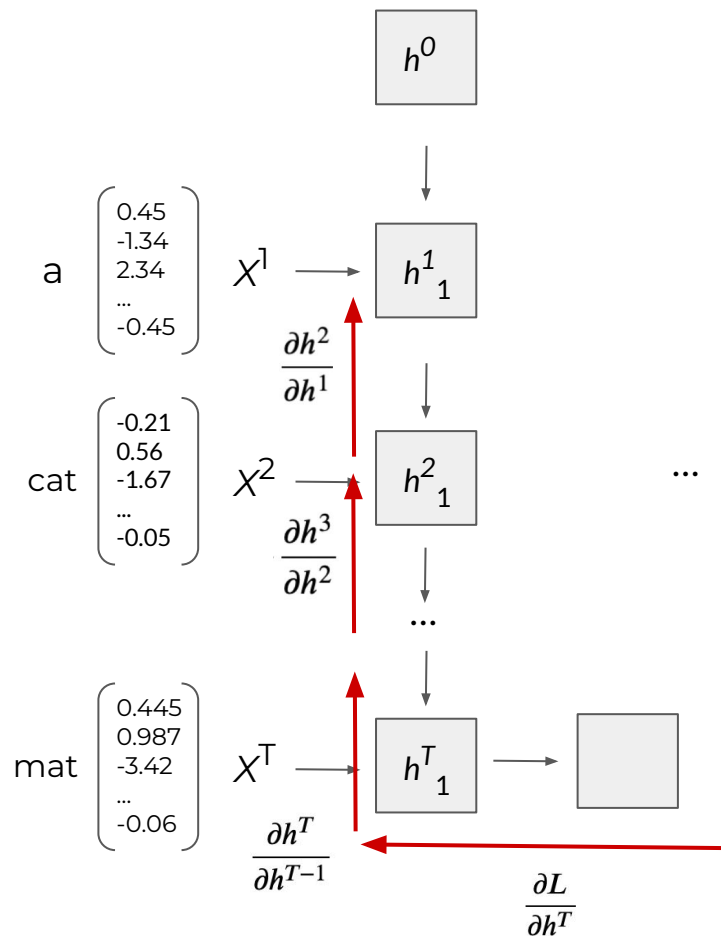
# План занятия

- Рекуррентный слой: идея.  
Рекуррентная нейросеть (RNN);
- Forward pass RNN;
- Обучение RNN (backward pass);
- **Функции активации в RNN;**
- **Bidirectional RNN;**
- GRU, LSTM



$$h^t = \sigma(WX^t + Uh^{t-1} + b_h)$$

$$\begin{aligned} \frac{dL}{dW} &= \frac{\partial L}{\partial h^3} \frac{dh^3}{dW} = \\ &= \frac{\partial L}{\partial h^3} \frac{\partial h^3}{\partial W} + \frac{\partial L}{\partial h^3} \frac{\partial h^3}{\partial h^2} \frac{dh^2}{dW} = \\ &= \frac{\partial L}{\partial h^3} \frac{\partial h^3}{\partial W} + \frac{\partial L}{\partial h^3} \frac{\partial h^3}{\partial h^2} \frac{\partial h^2}{\partial W} + \\ &+ \frac{\partial L}{\partial h^3} \frac{\partial h^3}{\partial h^2} \frac{\partial h^2}{\partial h^1} \frac{\partial h^1}{\partial W} \end{aligned}$$

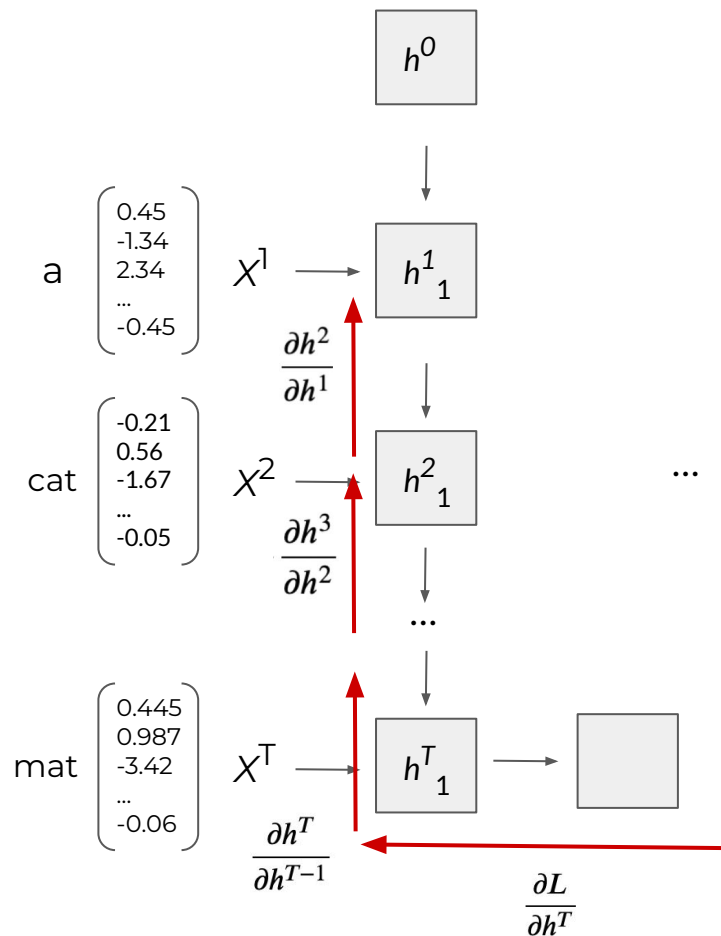


$$h^t = \sigma(WX^t + Uh^{t-1} + b_h)$$

$$\frac{dL}{dW} = \frac{\partial L}{\partial h^T} \frac{dh^T}{dW} =$$

$$= \frac{\partial L}{\partial h^T} \frac{\partial h^T}{\partial W} + \frac{\partial L}{\partial h^T} \frac{\partial h^T}{\partial h^{T-1}} \frac{\partial h^{T-1}}{\partial W} + \dots$$

$$\dots + \underbrace{\frac{\partial L}{\partial h^T} \frac{\partial h^T}{\partial h^{T-1}} \frac{\partial h^{T-1}}{\partial h^{T-2}} \dots \frac{dh^2}{dh^1} \frac{dh^1}{dW}}_{T-1}$$



$$h^t = \sigma(WX^t + Uh^{t-1} + b_h)$$

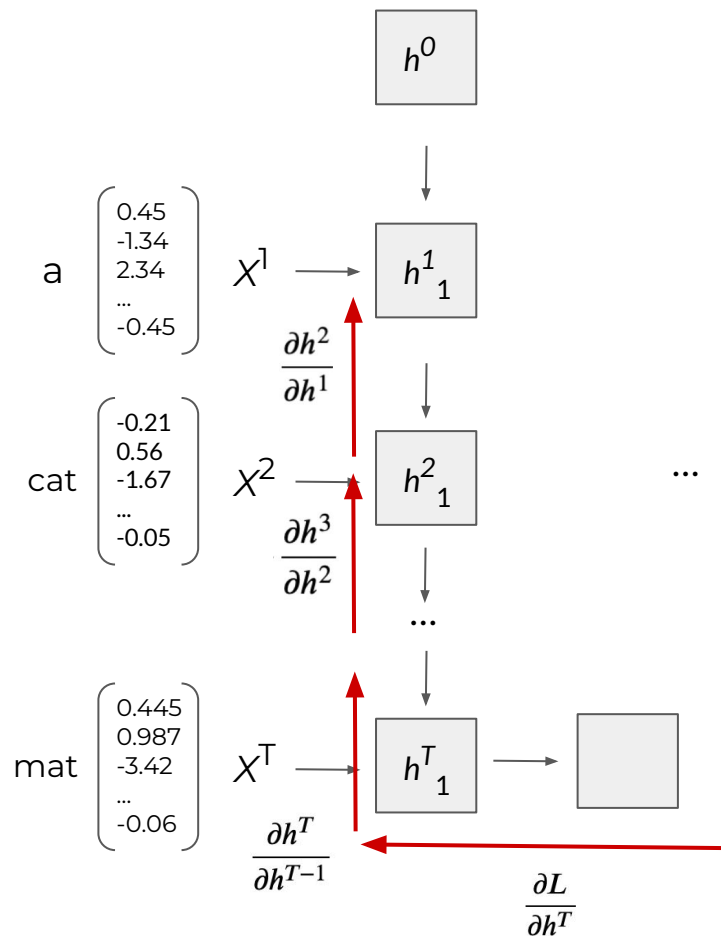
$$\frac{dL}{dW} = \frac{\partial L}{\partial h^T} \frac{dh^T}{dW} =$$

$$= \frac{\partial L}{\partial h^T} \frac{\partial h^T}{\partial W} + \frac{\partial L}{\partial h^T} \frac{\partial h^T}{\partial h^{T-1}} \frac{\partial h^{T-1}}{\partial W} + \dots$$

$$\dots + \underbrace{\frac{\partial L}{\partial h^T} \frac{\partial h^T}{\partial h^{T-1}} \frac{\partial h^{T-1}}{\partial h^{T-2}} \dots \frac{dh^2}{dh^1} \frac{dh^1}{dW}}_{T-1}$$

Для длинных входных последовательностей характерна проблема **затухания градиентов**





$$h^t = \sigma(WX^t + Uh^{t-1} + b_h)$$

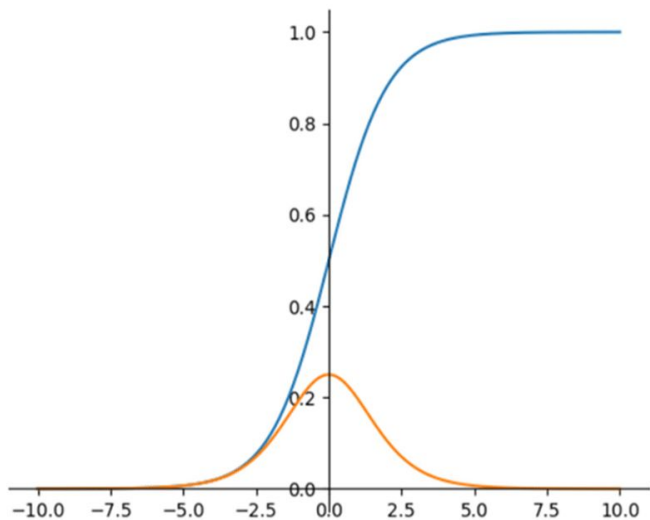
$$\frac{dL}{dW} = \frac{\partial L}{\partial h^T} \frac{dh^T}{dW} =$$

$$= \frac{\partial L}{\partial h^T} \frac{\partial h^T}{\partial W} + \frac{\partial L}{\partial h^T} \frac{\partial h^T}{\partial h^{T-1}} \frac{\partial h^{T-1}}{\partial W} + \dots$$

$$\dots + \frac{\partial L}{\partial h^T} \frac{\partial h^T}{\partial h^{T-1}} \frac{\partial h^{T-1}}{\partial h^{T-2}} \dots \frac{dh^2}{dh^1} \frac{dh^1}{dW}$$

В этих множителях  
содержится производная  
функции активации  $\sigma$

$$\sigma(x) = \frac{1}{1 + e^{-x}}$$



Сигмоидная функция  
активации и ее производная

$$h^t = \sigma(WX^t + Uh^{t-1} + b_h)$$

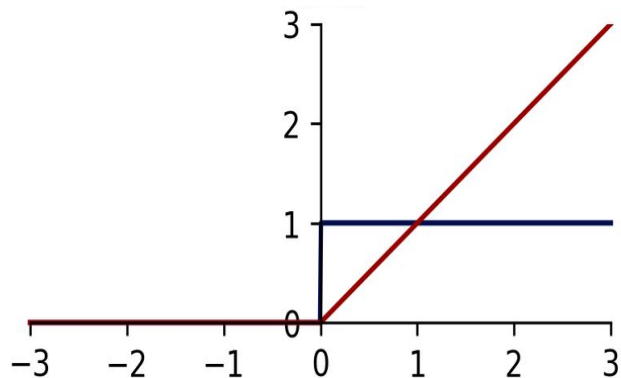
$$\frac{dL}{dW} = \frac{\partial L}{\partial h^T} \frac{dh^T}{dW} =$$

$$= \frac{\partial L}{\partial h^T} \frac{\partial h^T}{\partial W} + \frac{\partial L}{\partial h^T} \frac{\partial h^T}{\partial h^{T-1}} \frac{\partial h^{T-1}}{\partial W} + \dots$$

$$+ \underbrace{\frac{\partial L}{\partial h^T} \frac{\partial h^T}{\partial h^{T-1}} \frac{\partial h^{T-1}}{\partial h^{T-2}} \dots \frac{dh^2}{dh^1} \frac{dh^1}{dW}}_{\text{В этих множителях содержится производная функции активации } \sigma}$$

В этих множителях  
содержится производная  
функции активации  $\sigma$

$$f'(x) = \begin{cases} 0 & \text{for } x < 0 \\ 1 & \text{for } x \geq 0 \end{cases}$$



Функция активации ReLU  
и ее производная

$$h^t = \sigma(WX^t + Uh^{t-1} + b_h)$$

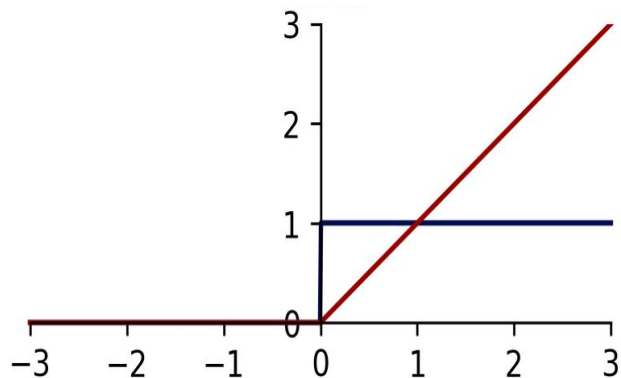
$$\frac{dL}{dW} = \frac{\partial L}{\partial h^T} \frac{dh^T}{dW} =$$

$$= \frac{\partial L}{\partial h^T} \frac{\partial h^T}{\partial W} + \frac{\partial L}{\partial h^T} \frac{\partial h^T}{\partial h^{T-1}} \frac{\partial h^{T-1}}{\partial W} + \dots$$

$$\dots + \underbrace{\frac{\partial L}{\partial h^T} \frac{\partial h^T}{\partial h^{T-1}} \frac{\partial h^{T-1}}{\partial h^{T-2}} \dots \frac{dh^2}{dh^1} \frac{dh^1}{dW}}_{U^{T-1}}$$

Значение элементов  
этой матрицы может  
быть очень большим

$$f'(x) = \begin{cases} 0 & \text{for } x < 0 \\ 1 & \text{for } x \geq 0 \end{cases}$$



Функция активации ReLU  
и ее производная

$$h^t = \sigma(WX^t + Uh^{t-1} + b_h)$$

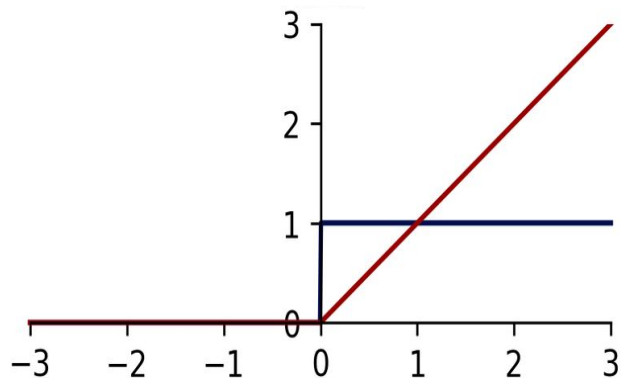
$$\frac{dL}{dW} = \frac{\partial L}{\partial h^T} \frac{dh^T}{dW} =$$

$$= \frac{\partial L}{\partial h^T} \frac{\partial h^T}{\partial W} + \frac{\partial L}{\partial h^T} \frac{\partial h^T}{\partial h^{T-1}} \frac{\partial h^{T-1}}{\partial W} + \dots$$

$$\dots + \underbrace{\frac{\partial L}{\partial h^T} \frac{\partial h^T}{\partial h^{T-1}} \frac{\partial h^{T-1}}{\partial h^{T-2}} \dots \frac{dh^2}{dh^1} \frac{dh^1}{dW}}_{U^{T-1}}$$

Для длинных входных последовательностей с функцией активации ReLU характерна проблема **взрыва градиентов**

$$f'(x) = \begin{cases} 0 & \text{for } x < 0 \\ 1 & \text{for } x \geq 0 \end{cases}$$



Функция активации ReLU  
и ее производная

$$h^t = \sigma(WX^t + Uh^{t-1} + b_h)$$

$$\frac{dL}{dW} = \frac{\partial L}{\partial h^T} \frac{dh^T}{dW} =$$

$$= \frac{\partial L}{\partial h^T} \frac{\partial h^T}{\partial W} + \frac{\partial L}{\partial h^T} \frac{\partial h^T}{\partial h^{T-1}} \frac{\partial h^{T-1}}{\partial W} + \dots$$

$$\dots + \frac{\partial L}{\partial h^T} \frac{\partial h^T}{\partial h^{T-1}} \frac{\partial h^{T-1}}{\partial h^{T-2}} \dots \frac{dh^2}{dh^1} \frac{dh^1}{dW}$$

Еще один минус RELU — он не ограничивает распределение выхода слоя.

Распределение вектора  $h^t$  может меняться в течение времени, что ухудшает работу нейросети.

$$h^t = \sigma(WX^t + Uh^{t-1} + b_h)$$

Как бороться?

Против взрыва градиентов:

- Gradient clipping;

Против затухания градиентов:

- Модели нейрона GRU,  
LSTM

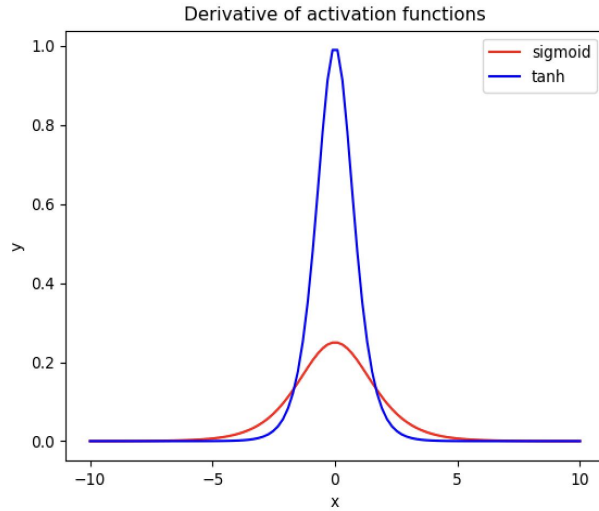
$$\frac{dL}{dW} = \frac{\partial L}{\partial h^T} \frac{dh^T}{dW} =$$

$$= \frac{\partial L}{\partial h^T} \frac{\partial h^T}{\partial W} + \frac{\partial L}{\partial h^T} \frac{\partial h^T}{\partial h^{T-1}} \frac{\partial h^{T-1}}{\partial W} + \dots$$

$$\dots + \frac{\partial L}{\partial h^T} \frac{\partial h^T}{\partial h^{T-1}} \frac{\partial h^{T-1}}{\partial h^{T-2}} \dots \frac{dh^2}{dh^1} \frac{dh^1}{dW}$$

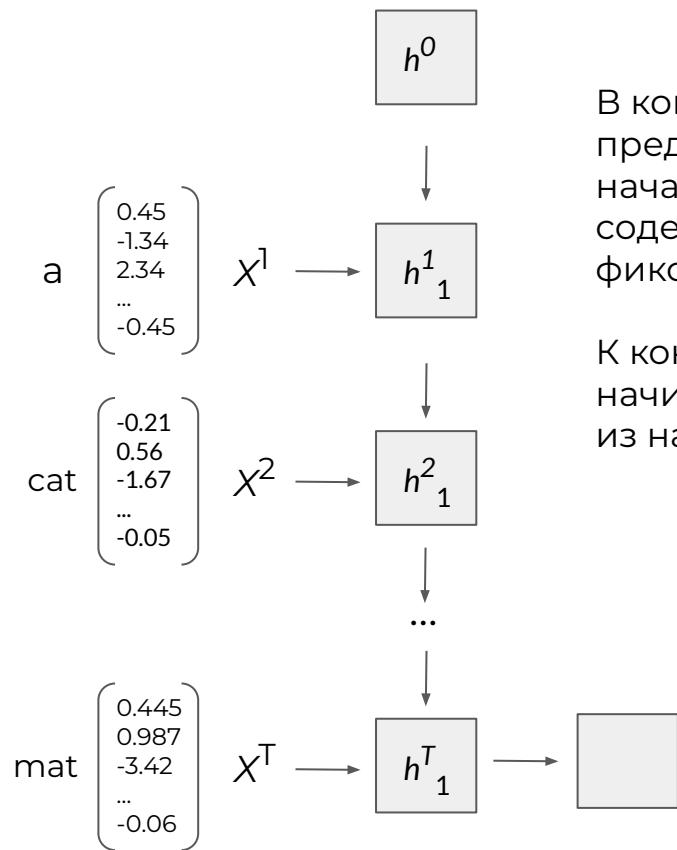
В целом: использовать  
функцию активации Tanh

$$h^t = \sigma(WX^t + Uh^{t-1} + b_h)$$



Производные  
сигмоиды и tanh

$$\begin{aligned} \frac{dL}{dW} &= \frac{\partial L}{\partial h^T} \frac{dh^T}{dW} = \\ &= \frac{\partial L}{\partial h^T} \frac{\partial h^T}{\partial W} + \frac{\partial L}{\partial h^T} \frac{\partial h^T}{\partial h^{T-1}} \frac{\partial h^{T-1}}{\partial W} + \dots \\ &\dots + \frac{\partial L}{\partial h^T} \frac{\partial h^T}{\partial h^{T-1}} \frac{\partial h^{T-1}}{\partial h^{T-2}} \dots \frac{dh^2}{dh^1} \frac{dh^1}{dW} \end{aligned}$$

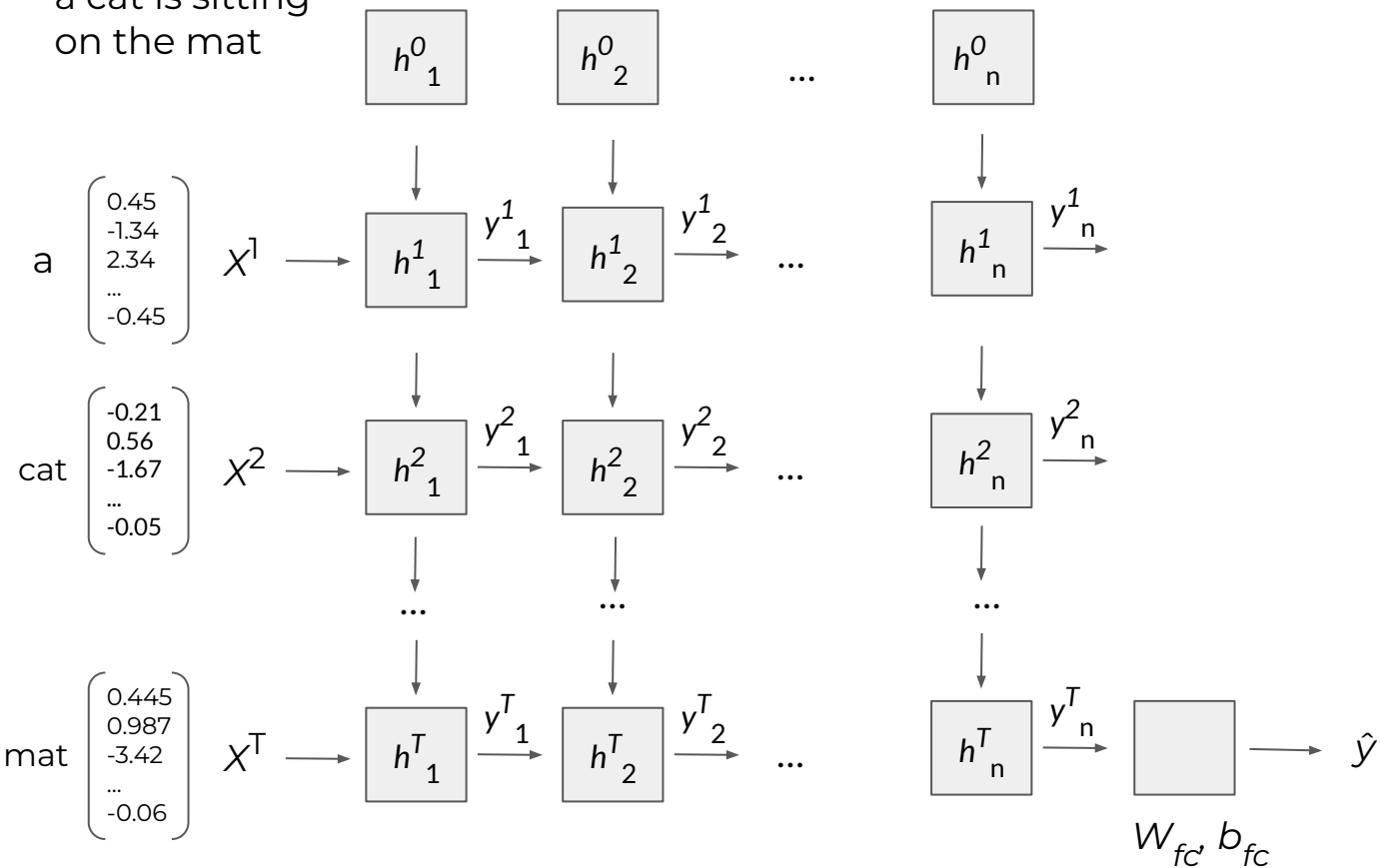


В конце обработки длинного предложения вся информация о начале предложения должна содержаться в одном векторе  $h$  фиксированного размера.

К концу длинного предложения  $h$  начинает “забывать” информацию из начала.

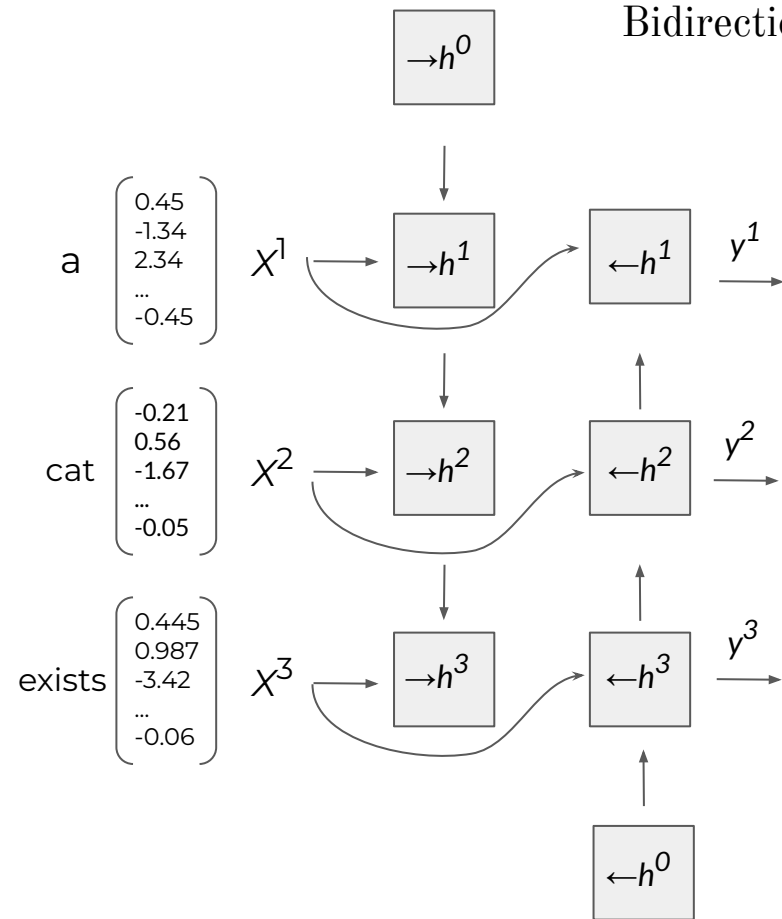


a cat is sitting  
on the mat



время

# Bidirectional RNN

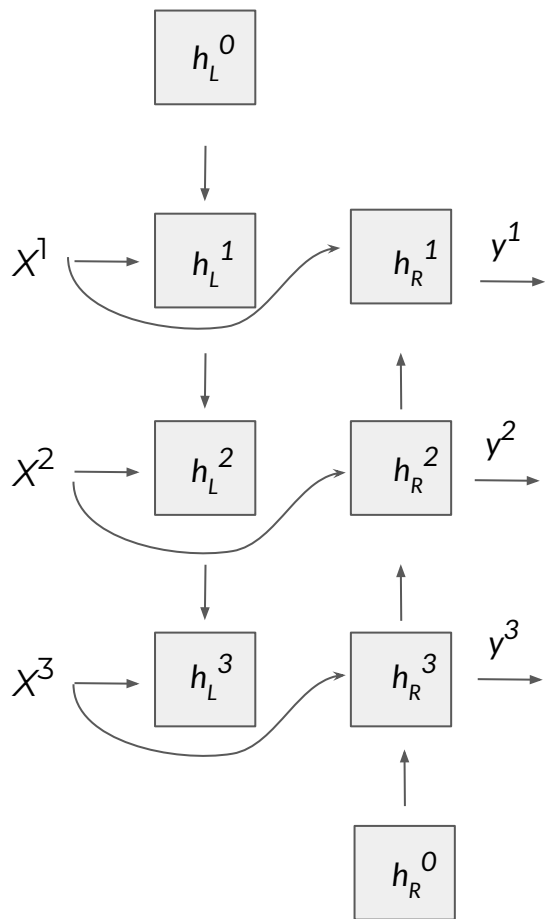


$$\vec{h}^t = \sigma(W^r X^t + U^r h^{t-1} + b_h^r)$$

$$\overleftarrow{h}^t = \sigma(W^l X^t + U^l h^{t-1} + b_h^l)$$

$$h^t = \text{concatenate}([\vec{h}^t, \overleftarrow{h}^t])$$

$$y^t = \sigma(V h^t + b_y)$$



$$h_L^t = \sigma(W_L X^t + U_L h_L^{t-1} + b_L)$$

$$h_R^t = \sigma(W_R X^t + U_R h_R^{t-1} + b_R)$$

$$h^t = [h_L^t, h_R^t]$$

$$y^t = \sigma(V h^t + b_y)$$

# Итоги видео

В этом видео мы обсудили некоторые нюансы RNN:

- Проблемы затухания и взрыва градиентов;
- Выбор функции активации;
- “Забывание” сети.

А также идею борьбы с проблемой забывания:  
bidirectional RNN.

В следующем видео мы рассмотрим идеи устройства  
GRU и LSTM вариантов слоев рекуррентной сети.



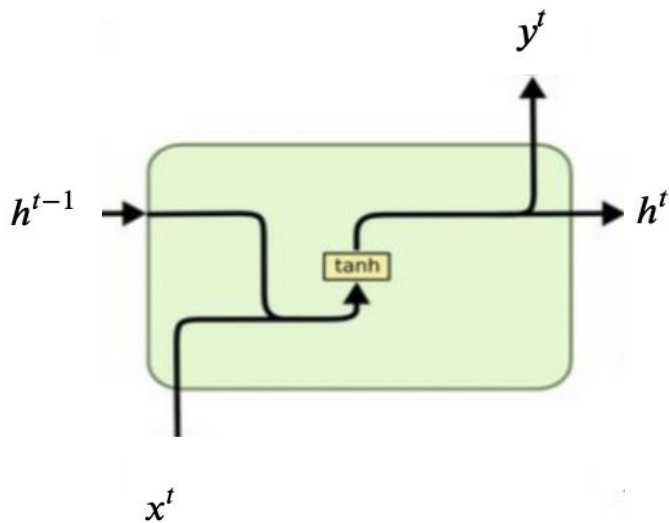
Deep Learning School

# LSTM, GRU

# План занятия

- Рекуррентный слой: идея.  
Рекуррентная нейросеть (RNN);
- Forward pass RNN;
- Обучение RNN (backward pass);
- Функции активации в RNN;
- Bidirectional RNN;
- **GRU, LSTM**

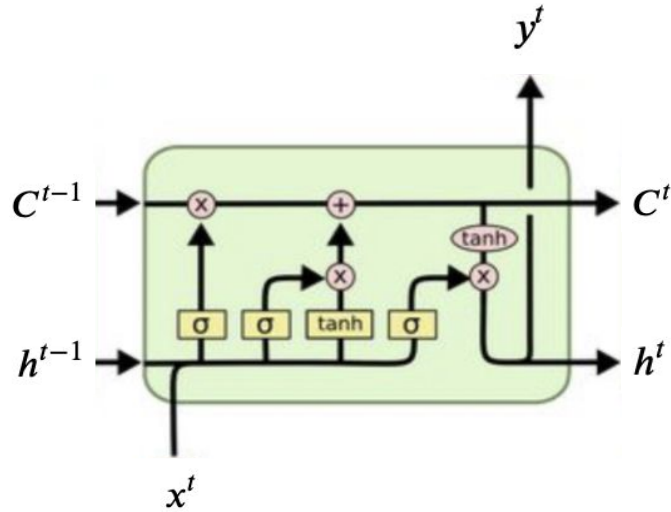
# Vanilla RNN



$$h^t = \tanh(WX^t + Uh^{t-1} + b_h)$$

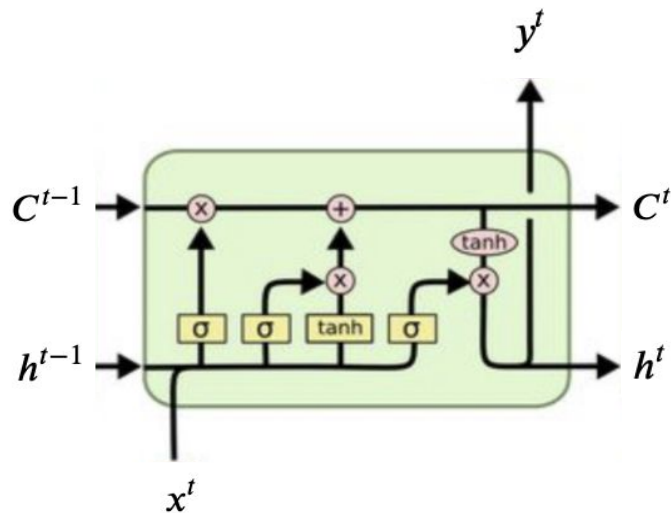
$$y^t = \sigma(W_y h^t + b_y)$$

# LSTM (Long Short Term Memory)





# LSTM (Long Short Term Memory)

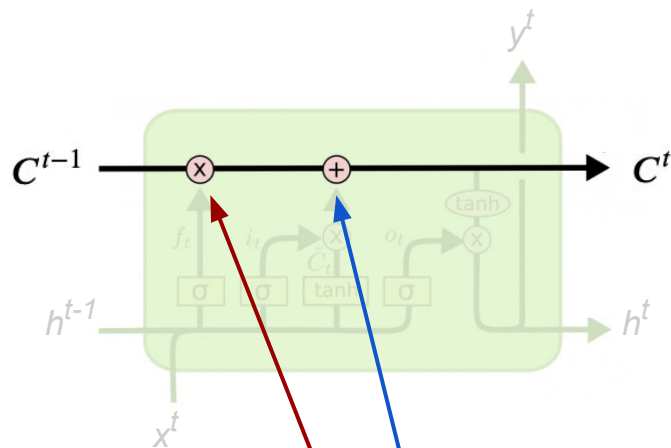


$C^t$  (cell) — “(долгосрочная) память”

$h^t$  — “краткосрочная память” или  
“текущее состояние слоя”

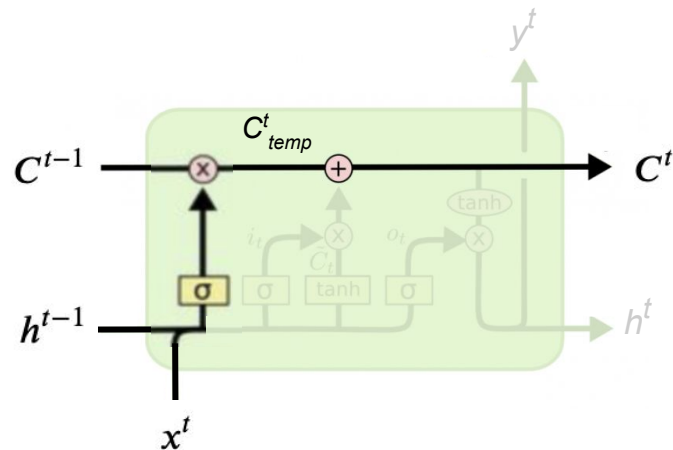
Оба вектора имеют тот же  
размер, что и  $x^t$

# LSTM



LSTM может удалять информацию из памяти  $C_t$  и записывать в нее новую информацию

# LSTM



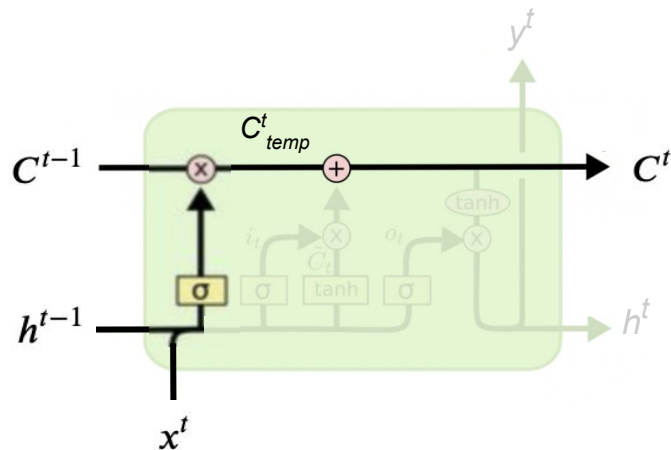
“Ворота забывания” (“forget gate”):

$$f^t = \sigma(W_f \cdot [h^{t-1}, x^t] + b_f)$$

$$C_{temp}^t = f^t * C^{t-1}$$

Берем текущее состояние краткосрочной памяти ( $h^{t-1}$ ) и новую информацию, пришедшую на вход ( $x^t$ ). На их основе понимаем, какую информацию из долгосрочной памяти ( $C^{t-1}$ ) уже можно выкинуть

# LSTM



“Ворота забывания” (“forget gate”):

$$f^t = \sigma(W_f \cdot [h^{t-1}, x^t] + b_f)$$

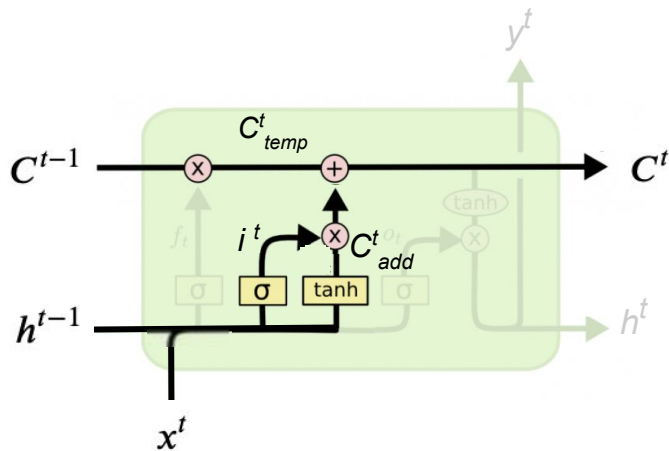
сигмоида

$$C_{temp}^t = f^t * C^{t-1}$$

Вектор значений от 0 до 1.  
Размер тот же, что у  $C^{t-1}$

Каждый элемент вектора  $C^{t-1}$  умножается на значение от 0 до 1, т.е. часть информации из всех элементов исчезает

# LSTM



“Ворота входа” (“input gate”):

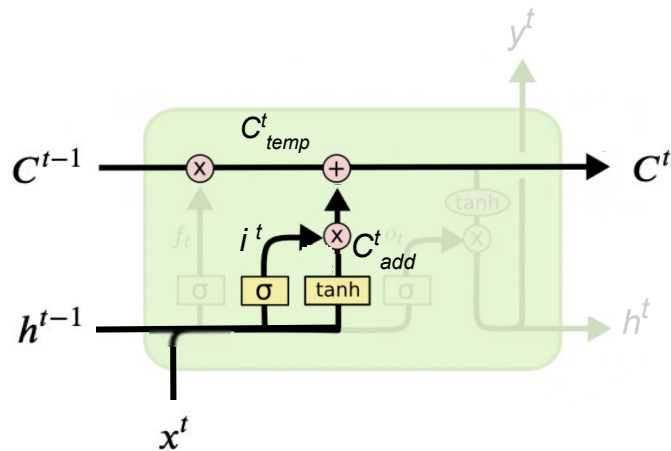
$$i^t = \sigma(W_i \cdot [h^{t-1}, x^t] + b_i)$$

$$C_{add}^t = \tanh(W_C \cdot [h^{t-1}, x^t] + b_C)$$

$$C^t = C_{temp}^t + i^t * C_{add}^t$$

Берем текущее состояние краткосрочной памяти ( $h^{t-1}$ ) и новую информацию, пришедшую на вход ( $x^t$ ). Понимаем, какую информацию из них нам надо сохранить в долгосрочную память ( $C^{t-1}$ )

# LSTM



“Ворота входа” (“input gate”):

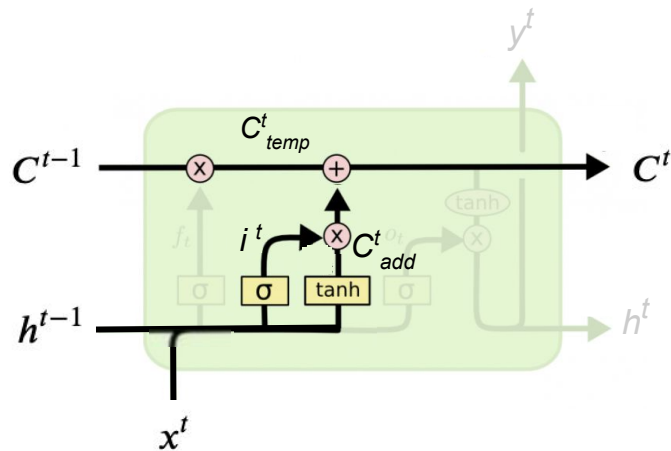
$$i^t = \sigma(W_i \cdot [h^{t-1}, x^t] + b_i)$$

$$C_{add}^t = \tanh(W_C \cdot [h^{t-1}, x^t] + b_C)$$

$$C^t = C_{temp}^t + i^t * C_{add}^t$$

На основе  $h^{t-1}$  и  $x^t$  понимаем, какую новую информацию хотим добавить в  $C^t$

# LSTM



“Ворота входа” (“input gate”):

$$i^t = \sigma(W_i \cdot [h^{t-1}, x^t] + b_i)$$

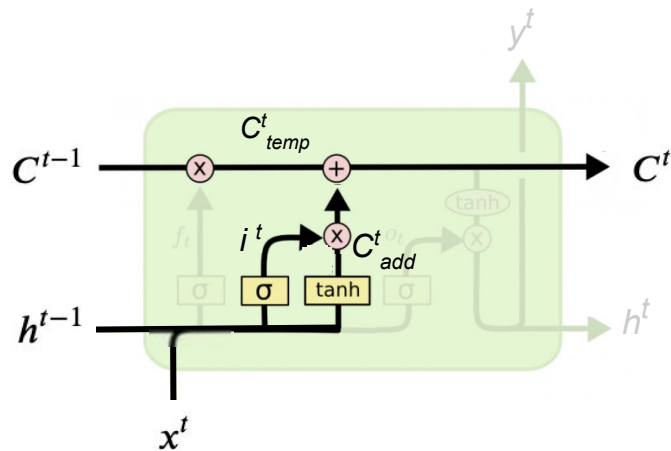
сигмоида

Вектор значений от 0 до 1.

$$C_{add}^t = \tanh(W_C \cdot [h^{t-1}, x^t] + b_C)$$

$$C^t = C_{temp}^t + i^t * C_{add}^t$$

# LSTM



“Ворота входа” (“input gate”):

$$i^t = \sigma(W_i \cdot [h^{t-1}, x^t] + b_i)$$

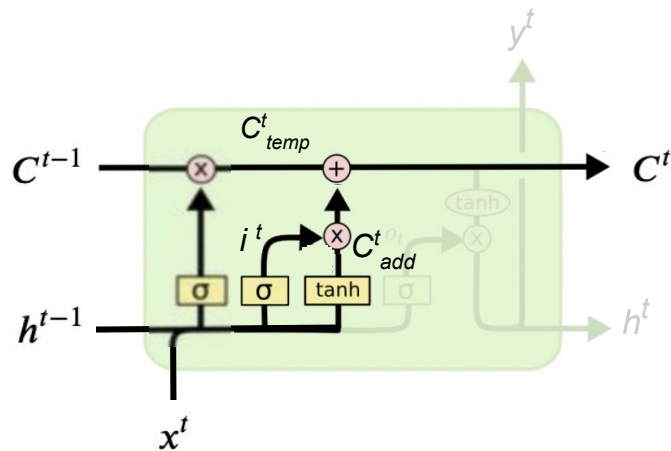
$$C_{add}^t = \tanh(W_C \cdot [h^{t-1}, x^t] + b_C)$$

$$C^t = C_{temp}^t + i^t * C_{add}^t$$

Добавляем в  $C^t$  новую информацию. Каждый элемент вектора умножается на число от 0 до 1: так регулируется, сколько именно этой информации поступит в вектор  $C^t$



# LSTM



Полное обновление вектора  $C^t$ :

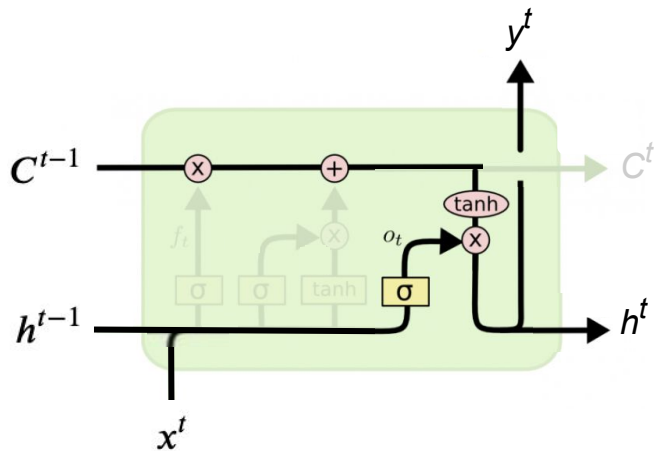
$$f^t = \sigma(W_f \cdot [h^{t-1}, x^t] + b_f)$$

$$i^t = \sigma(W_i \cdot [h^{t-1}, x^t] + b_i)$$

$$C_{add}^t = \tanh(W_C \cdot [h^{t-1}, x^t] + b_C)$$

$$C^t = f^t * C^{t-1} + i^t * C_{add}^t$$

# LSTM



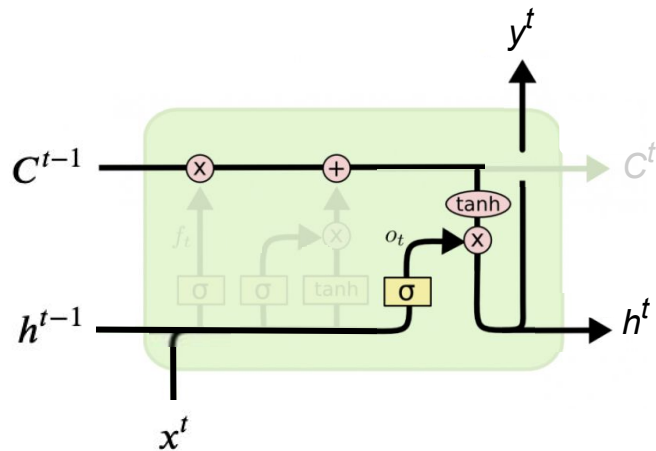
Обновление вектора  $h^t$  и  
получение выхода ячейки  $y^t$  :

$$o^t = \sigma(W_o \cdot [h^{t-1}, x^t] + b_o)$$

$$h^t = o^t * \tanh(C^t)$$

$$y^t = \sigma(W_y h^t + b_y)$$

# LSTM



Обновление вектора  $h^t$  и  
получение выхода ячейки  $y^t$  :

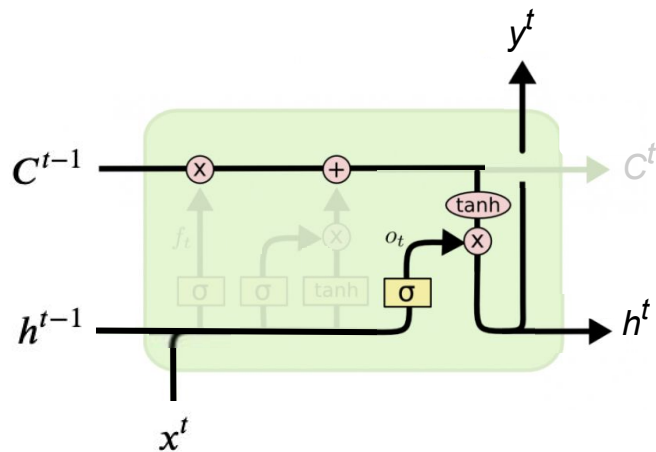
$$o^t = \sigma(W_o \cdot [h^{t-1}, x^t] + b_o)$$

$$h^t = o^t * \tanh(C^t)$$

$$y^t = \sigma(W_y h^t + b_y)$$

Переносим часть информации из  
долгосрочной памяти (  $C^t$  ) в  
краткосрочную (  $h^t$  ) и получаем ответ

# LSTM



Обновление вектора  $h^t$  и  
получение выхода ячейки  $y^t$  :

$$o^t = \sigma(W_o \cdot [h^{t-1}, x^t] + b_o)$$

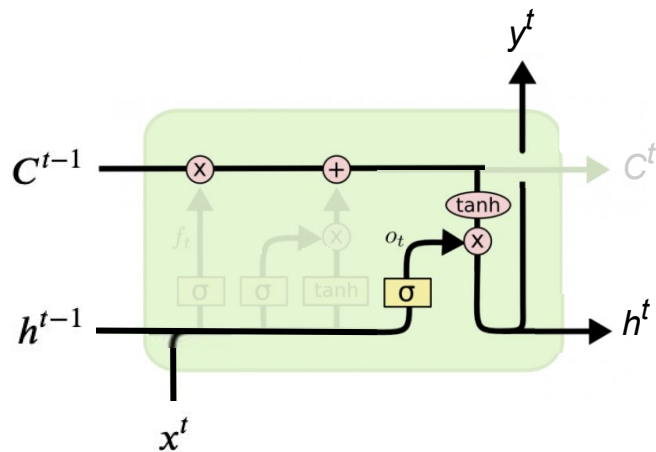
сигмоида

$$h^t = o^t * \tanh(C^t)$$

Вектор из 0 и 1

$$y^t = \sigma(W_y h^t + b_y)$$

# LSTM



Обновление вектора  $h^t$  и  
получение выхода ячейки  $y^t$  :

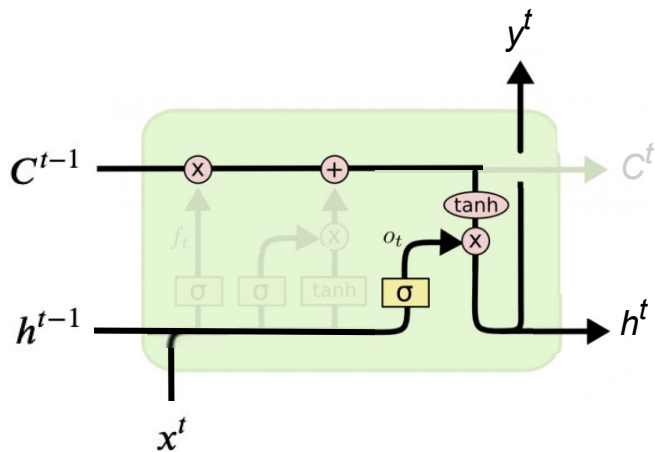
$$o^t = \sigma(W_o \cdot [h^{t-1}, x^t] + b_o)$$

$$h^t = o^t * \tanh(C^t)$$

$$y^t = \sigma(W_y h^t + b_y)$$

Часть информации из  $C^t$   
переносится в  
краткосрочную память

# LSTM



Обновление вектора  $h^t$  и  
получение выхода ячейки  $y^t$  :

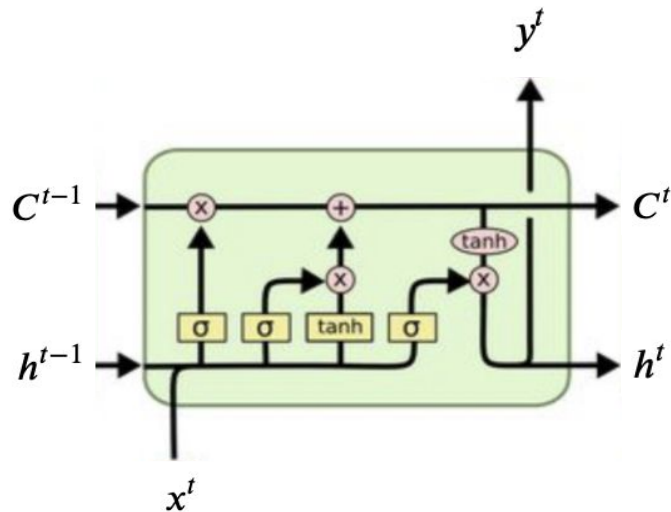
$$o^t = \sigma(W_o \cdot [h^{t-1}, x^t] + b_o)$$

$$h^t = o^t * \tanh(C^t)$$

$$y^t = \sigma(W_y h^t + b_y)$$

Считаем выход ячейки

# LSTM



$$f^t = \sigma(W_f \cdot [h^{t-1}, x^t] + b_f)$$

$$i^t = \sigma(W_i \cdot [h^{t-1}, x^t] + b_i)$$

$$C_{add}^t = \tanh(W_C \cdot [h^{t-1}, x^t] + b_C)$$

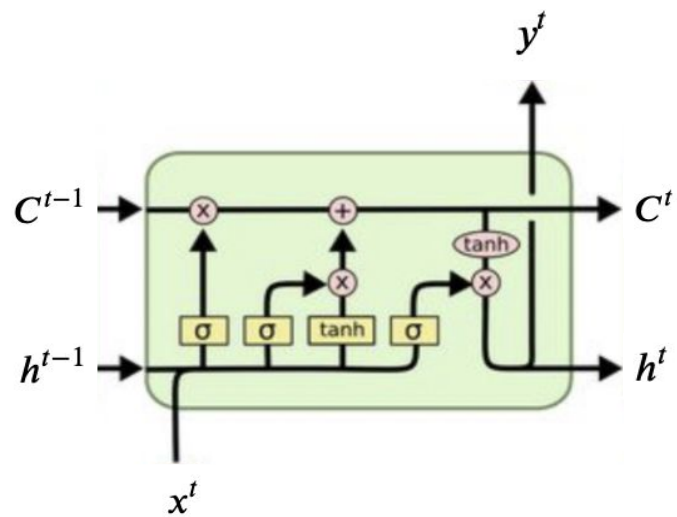
$$C^t = f^t * C^{t-1} + i^t * C_{add}^t$$

$$o^t = \sigma(W_o \cdot [h^{t-1}, x^t] + b_o)$$

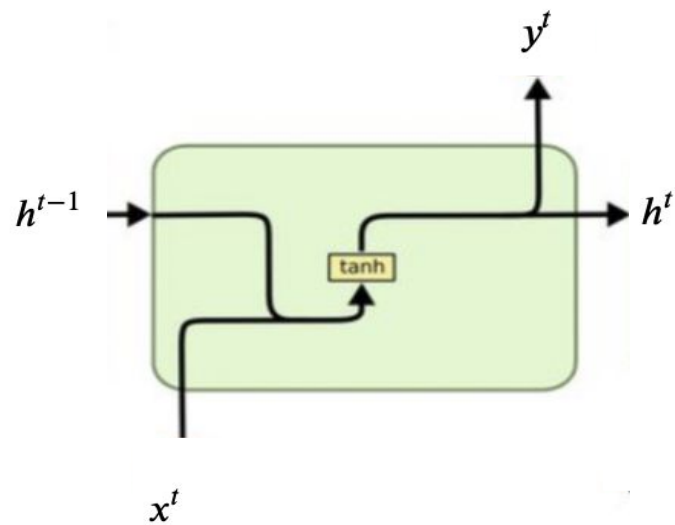
$$h^t = o^t * \tanh(C^t)$$

$$y^t = \sigma(W_y h^t + b_y)$$

# LSTM



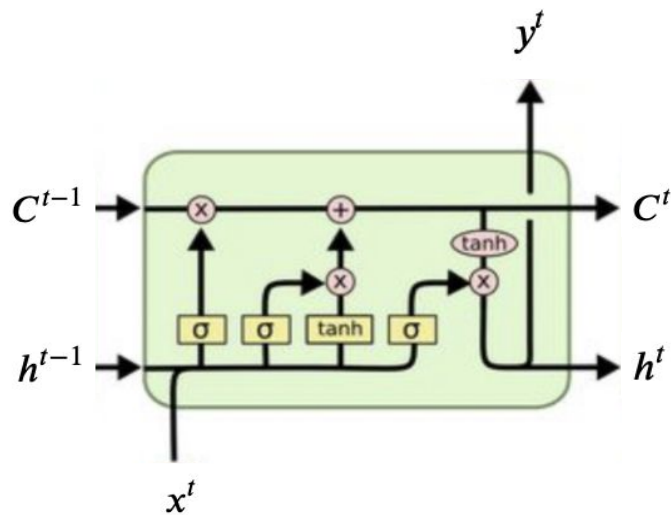
# Vanilla RNN



$$h^t = \tanh(WX^t + Uh^{t-1} + b_h)$$



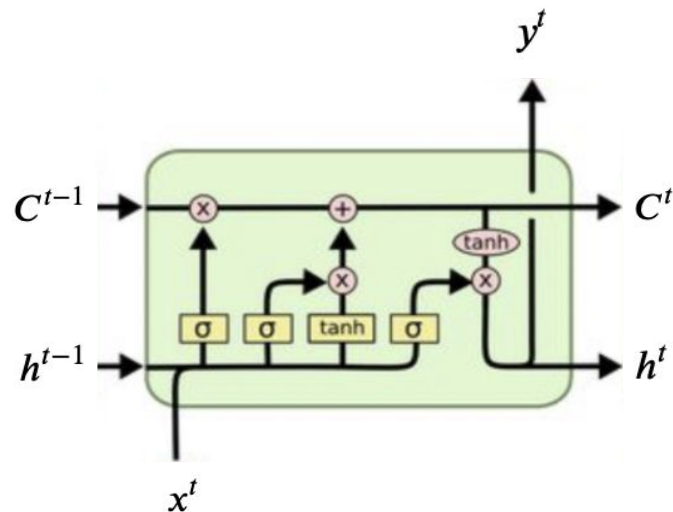
## LSTM



LSTM реализует механизм, похожий на skip connection в ResNet.

Это помогает в борьбе с затуханием градиентов.

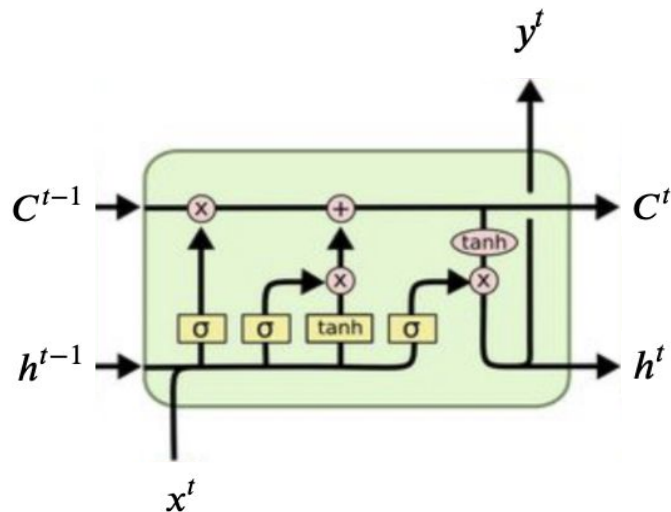
# LSTM



Пример: пусть мы решаем задачу классификации текста по грамматической правильности (правильный/неправильный текст с точки зрения грамматики)

*“Человек в шляпе с красивым попугаем на плече, которого вчера видели на пересечении Косого переулка и улицы Роз, зашел в бар”*

# LSTM

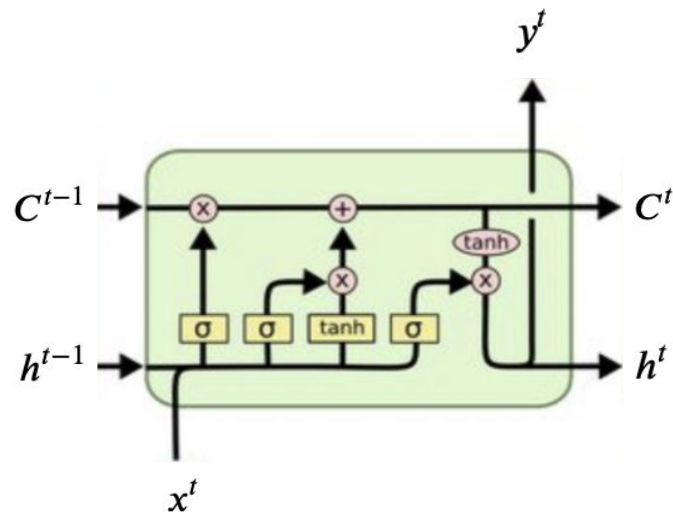


Пример: пусть мы решаем задачу классификации текста по грамматической правильности (правильный/неправильный текст с точки зрения грамматики)

*“Человек в иллье с красивым попугаем на плече, которого вчера видели на пересечении Косого переулка и улицы Роз, зашел в бар”*

Когда мы стоим на слове “человек”, мы должны **добавить** в память  $C^t$  информацию об этом слове. Например, что оно мужского рода.

# LSTM

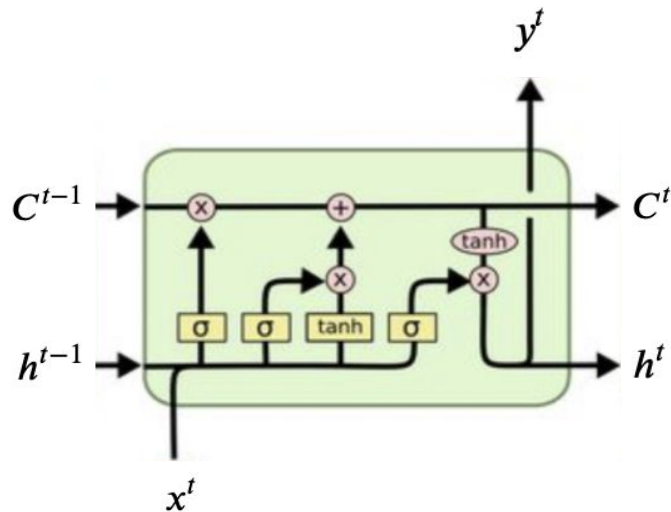


Пример: пусть мы решаем задачу классификации текста по грамматической правильности (правильный/неправильный текст с точки зрения грамматики)

*“Человек в шляпе с красивым попугаем на плече, которого вчера видели на пересечении Косого переулка и улицы Роз, зашел в бар”*

Когда мы стоим на слове “красивым”, мы также должны **добавить** в память  $C^t$  информацию об этом слове. Что оно мужского рода.

# LSTM

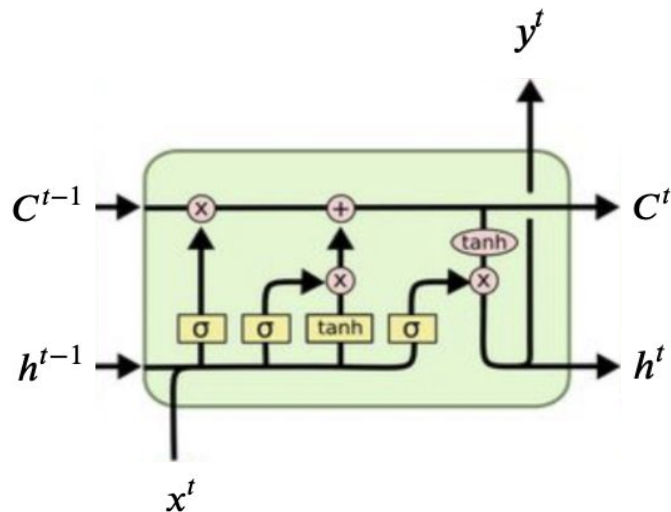


Пример: пусть мы решаем задачу классификации текста по грамматической правильности (правильный/неправильный текст с точки зрения грамматики)

*“Человек в иллье с красивым **попугаем** на плече, которого вчера видели на пересечении Косого переулка и улицы Роз, зашел в бар”*

Когда мы стоим на слове “попугаем”, то при формировании вектора  $h^t$  из  $C^t$  мы хотим вытащить информацию о форме слова “красивым”, но не хотим добавлять информацию о форме слова “человек”

# LSTM

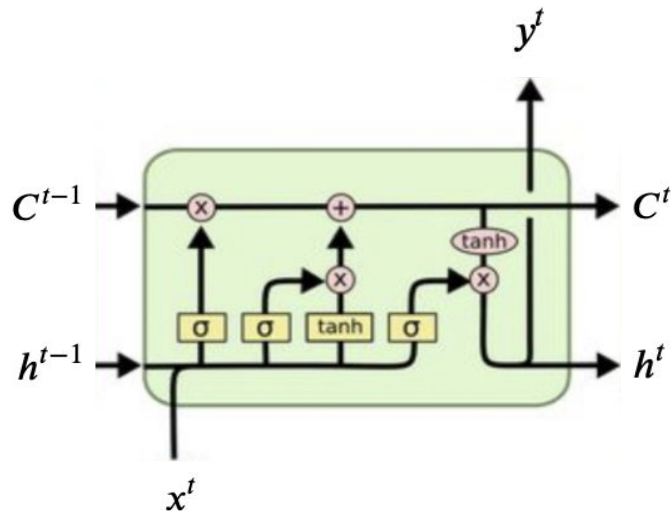


Пример: пусть мы решаем задачу классификации текста по грамматической правильности (правильный/неправильный текст с точки зрения грамматики)

*“Человек в шляпе с красивым попугаем **на** плече, которого вчера видели на пересечении Косого переулка и улицы Роз, зашел в бар”*

Когда мы стоим на слове “на”, то мы можем удалить из  $C^t$  информацию о словах “красивым попугаем”, и добавить информацию о слове “на”

# LSTM

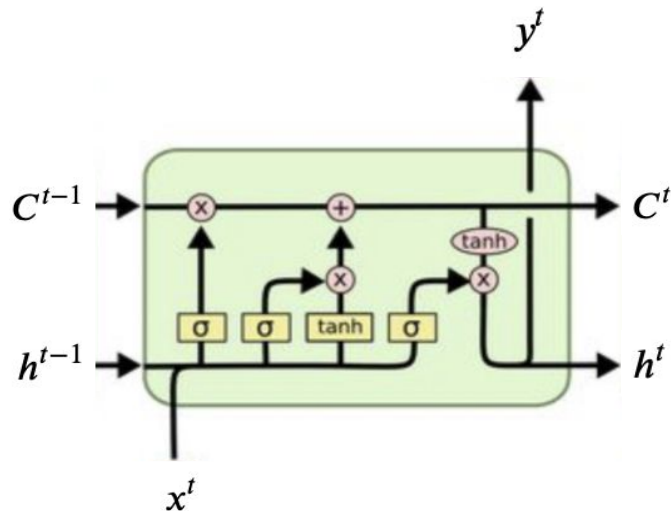


Пример: пусть мы решаем задачу классификации текста по грамматической правильности (правильный/неправильный текст с точки зрения грамматики)

*“Человек в иллье с красивым попугаем на  
плече, которого вчера видели на пересечении  
Косого переулка и улицы Роз, зашел в бар”*

Когда мы стоим на слове “плече”, то при формировании вектора  $h^t$  из  $C^t$  мы хотим вытащить информацию о форме слова “на”, но не хотим добавлять информацию о форме слова “человек”

# LSTM



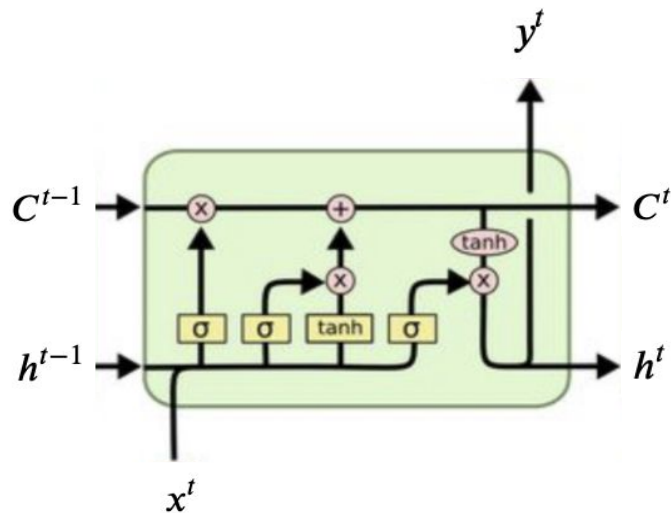
Пример: пусть мы решаем задачу классификации текста по грамматической правильности (правильный/неправильный текст с точки зрения грамматики)

*“Человек в иллье с красивым попугаем на  
плече, которого вчера видели на пересечении  
Косого переулка и улицы Роз, зашел в бар”*

Когда мы стоим на слове “плече”, то при формировании вектора  $h^t$  из  $C^t$  мы хотим вытащить информацию о форме слова “на”, но не хотим добавлять информацию о форме слова “человек”



# LSTM

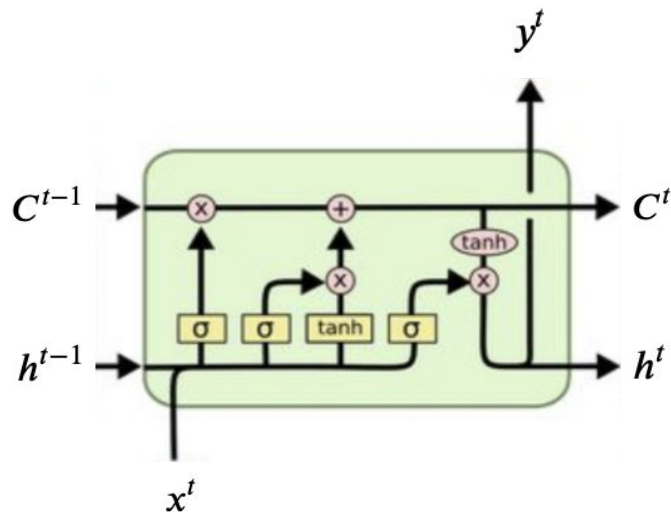


Пример: пусть мы решаем задачу классификации текста по грамматической правильности (правильный/неправильный текст с точки зрения грамматики)

*“Человек в шляпе с красивым попугаем на плече, которого вчера видели на пересечении Косого переулка и улицы Роз, зашел в бар”*

Когда мы стоим на запятой, мы должны **добавить** в память информацию о том, что начался причастный оборот

# LSTM

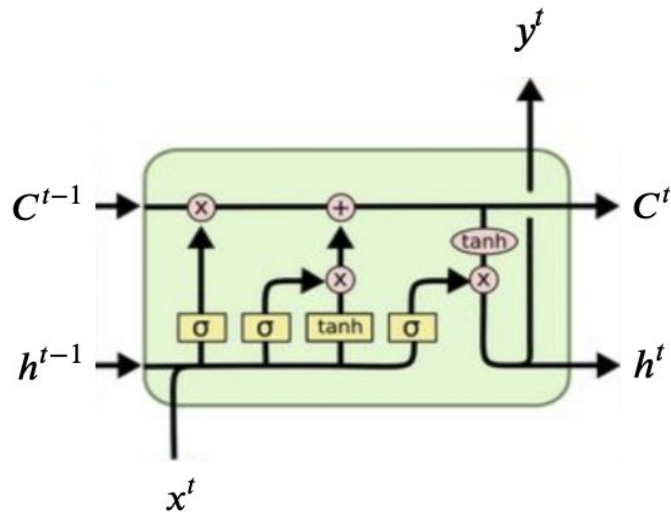


Пример: пусть мы решаем задачу классификации текста по грамматической правильности (правильный/неправильный текст с точки зрения грамматики)

*“Человек в иллье с красивым попугаем на плече, **которого** вчера видели на пересечении Косого переулка и улицы Роз, зашел в бар”*

Когда мы стоим на слове “которого”, то при формировании вектора  $h^t$  из  $C^t$  мы хотим вытащить информацию о форме слове “человек” и о том, что сейчас идет причастный оборот

# LSTM



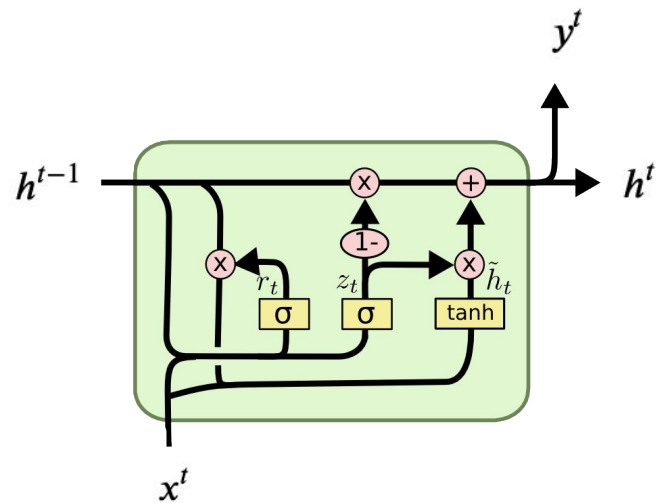
Пример: пусть мы решаем задачу классификации текста по грамматической правильности (правильный/неправильный текст с точки зрения грамматики)

*“Человек в шляпе с красивым попугаем на плече, которого вчера видели на пересечении Косого переулка и улицы Роз, зашел в бар. После ...”*

Когда мы стоим на слове “после”, то мы можем удалить из  $C^t$  информацию о предыдущем предложении, и добавить информацию о слове “после”

# GRU

“Облегченный вариант” LSTM



$$z^t = \sigma(W_z \cdot [h^{t-1}, x^t] + b_z)$$

$$r^t = \sigma(W_r \cdot [h^{t-1}, x^t] + b_r)$$

$$h_{add}^t = \tanh(W \cdot [r^t * h^{t-1}, x^t])$$

$$h^t = (1 - z^t) * h^{t-1} + z^t * h_{add}^t$$

# Итоги видео

В этом видео мы познакомились с идеей устройства новых рекуррентных слоев: LSTM и GRU.