

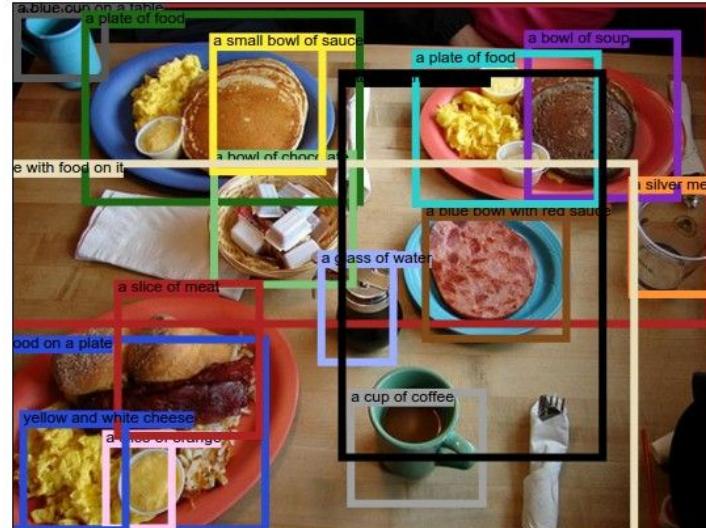
# Цифровая обработка изображения

8. Применение рекурентных сетей в задачах  
анализа изображений

# План занятия

- Рекурентные сети
- Автоматическая аннотация изображения
- Распознавание текста

# Описание изображения



a plate of food. food on a plate. a blue cup on a table. a plate of food. a blue bowl with red sauce. a bowl of soup. a cup of coffee. a bowl of chocolate. a glass of water. a plate of food. a silver metal container. a small bowl of sauce. table with food on it. a slice of orange. a table with food on it. a slice of meat. yellow and white cheese.

# Распознавание текста на изображении



# Распознавание рукописного текста

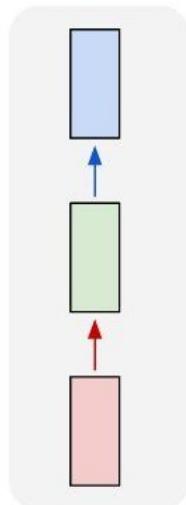
Optical Character Recognition  
is designed to convert your  
handwriting into text.

Optical Character Recognition  
is designed to convert your  
handwriting into text.

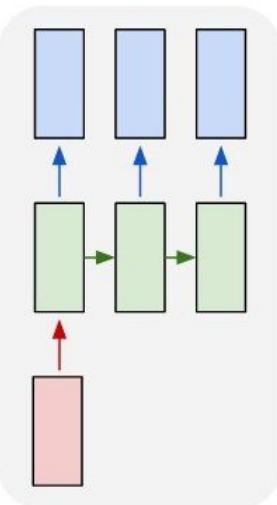
# Рекурентные сети

# Рекурентные сети

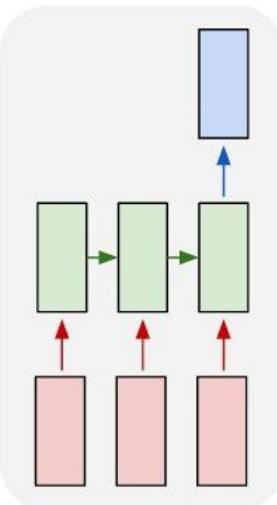
one to one



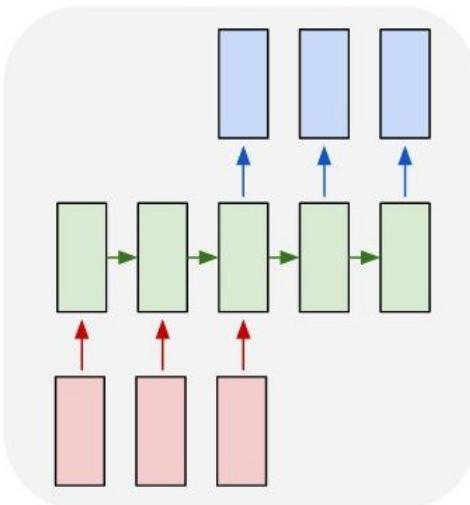
one to many



many to one



many to many



many to many

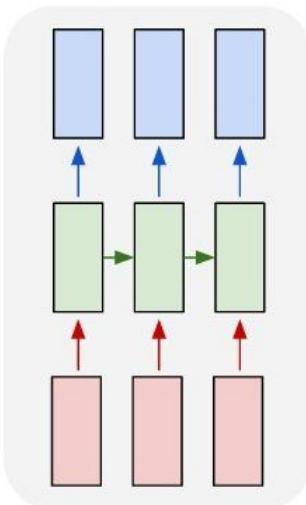


Image in  
Label out

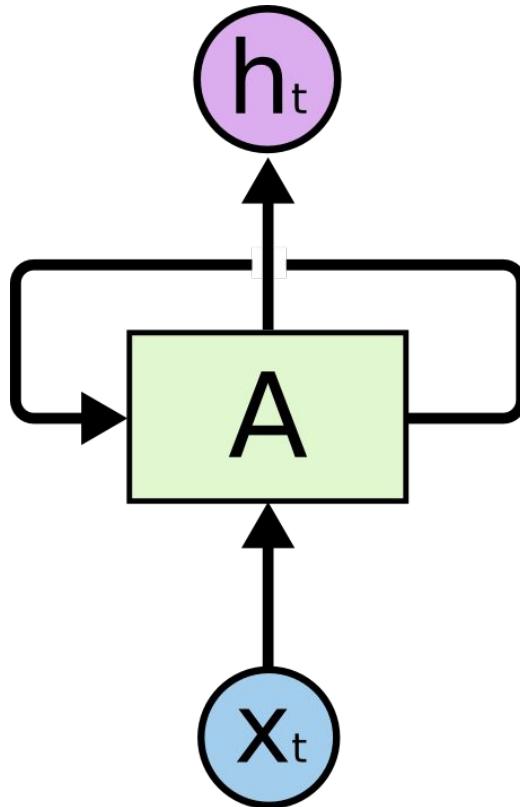
Image in  
Words out

Words in  
Sentiment out

English in  
Portuguese out

Video In  
Labels out

# Рекурентные сети



$$h_t = f_{weights}(h_{t-1}, x_t) \therefore$$
$$h_t = \tanh(W_{hh} \cdot h_{t-1} + W_{xh} \cdot x_t)$$

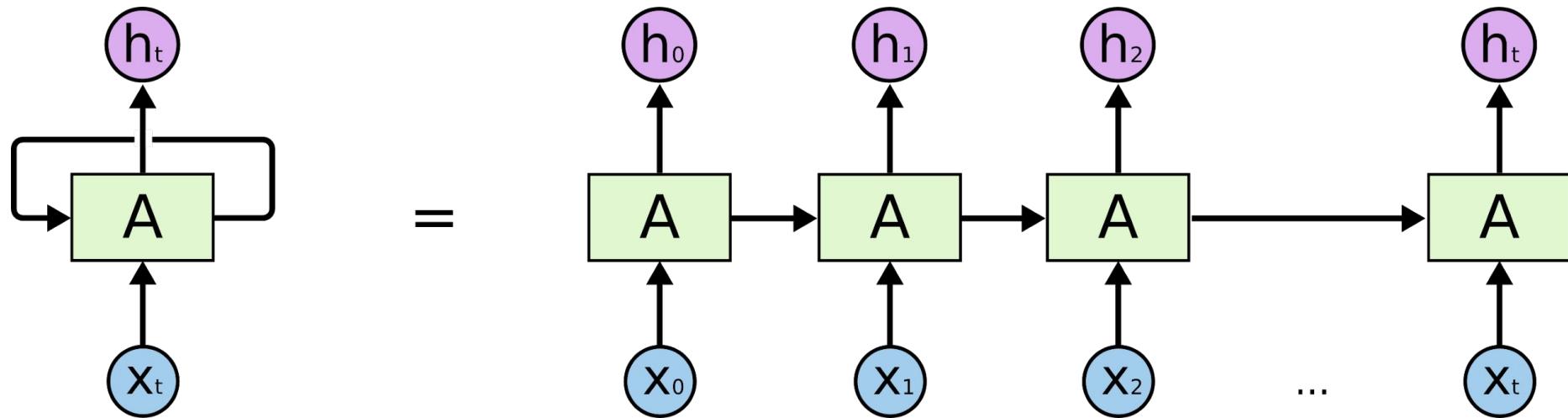
$h_t$  - состояние ячейки в момент времени  $t$

$x_t$  - входной сигнал на шаге  $t$

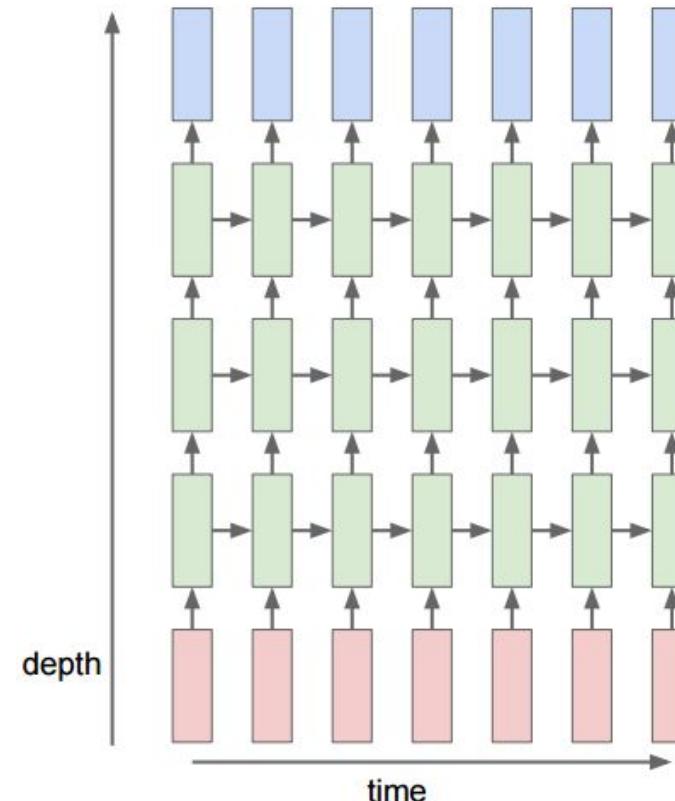
$W_{hh}$  - матрица преобразования состояния

$W_{xh}$  - матрица преобразования входного сигнала

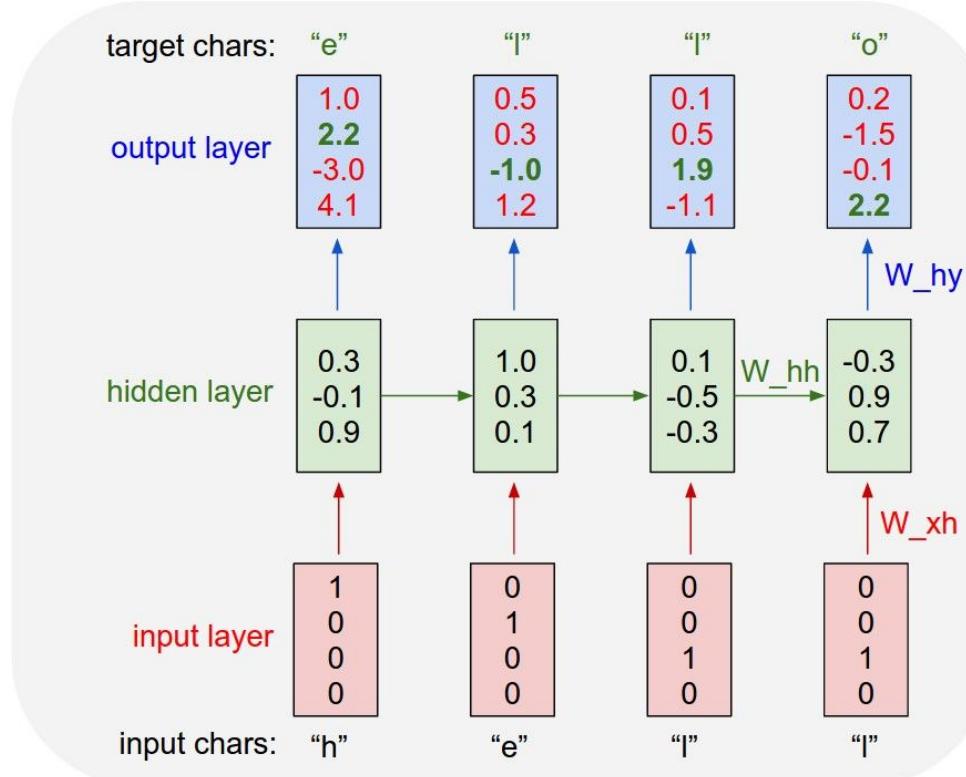
# Рекуррентные сети



# Рекурентные сети



# Генерация текстовой последовательности



# Генерация текстовой последовательности

*Proof.* Omitted.  $\square$

**Lemma 0.1.** Let  $\mathcal{C}$  be a set of the construction.

Let  $\mathcal{C}$  be a gerber covering. Let  $\mathcal{F}$  be a quasi-coherent sheaves of  $\mathcal{O}$ -modules. We have to show that

$$\mathcal{O}_{\mathcal{O}_X} = \mathcal{O}_X(\mathcal{L})$$

*Proof.* This is an algebraic space with the composition of sheaves  $\mathcal{F}$  on  $X_{\text{étale}}$  we have

$$\mathcal{O}_X(\mathcal{F}) = \{\text{morph}_1 \times_{\mathcal{O}_X} (\mathcal{G}, \mathcal{F})\}$$

where  $\mathcal{G}$  defines an isomorphism  $\mathcal{F} \rightarrow \mathcal{F}$  of  $\mathcal{O}$ -modules.  $\square$

**Lemma 0.2.** This is an integer  $\mathcal{Z}$  is injective.

*Proof.* See Spaces, Lemma ??.

This since  $\mathcal{F} \in \mathcal{F}$  and  $x \in \mathcal{G}$  the diagram

$$\begin{array}{ccccc}
 S & \xrightarrow{\quad} & & & \\
 \downarrow & & & & \\
 \xi & \longrightarrow & \mathcal{O}_{X'} & & \\
 & & \uparrow & \searrow & \\
 & & =\alpha' & \longrightarrow & X \\
 & & \uparrow & & \downarrow \\
 & & =\alpha' & \longrightarrow & \text{Spec}(K_\psi) \\
 & & \downarrow & & \text{Mor}_{\text{Sets}} \\
 & & \alpha & & \text{d}(\mathcal{O}_{X_{f/k}}, \mathcal{G})
 \end{array}$$

is a limit. Then  $\mathcal{G}$  is a finite type and assume  $S$  is a flat and  $\mathcal{F}$  and  $\mathcal{G}$  is a finite type  $f_*$ . This is of finite type diagrams, and

- the composition of  $\mathcal{G}$  is a regular sequence,
- $\mathcal{O}_{X'}$  is a sheaf of rings.

$\square$

*Proof.* We have see that  $X = \text{Spec}(R)$  and  $\mathcal{F}$  is a finite type representable by algebraic space. The property  $\mathcal{F}$  is a finite morphism of algebraic stacks. Then the cohomology of  $X$  is an open neighbourhood of  $U$ .  $\square$

*Proof.* This is clear that  $\mathcal{G}$  is a finite presentation, see Lemmas ??.  
A reduced above we conclude that  $U$  is an open covering of  $\mathcal{C}$ . The functor  $\mathcal{F}$  is a “field”

$$\mathcal{O}_{X_{\bar{x}}} \longrightarrow \mathcal{F}_{\bar{x}} \dashv^{-1} (\mathcal{O}_{X_{\text{étale}}}) \longrightarrow \mathcal{O}_{X_{\bar{x}}}^{-1} \mathcal{O}_{X_{\bar{x}}}(\mathcal{O}_{X_{\bar{x}}}^{\bar{\pi}})$$

is an isomorphism of covering of  $\mathcal{O}_{X_{\bar{x}}}$ . If  $\mathcal{F}$  is the unique element of  $\mathcal{F}$  such that  $X$  is an isomorphism.

The property  $\mathcal{F}$  is a disjoint union of Proposition ?? and we can filtered set of presentations of a scheme  $\mathcal{O}_X$ -algebra with  $\mathcal{F}$  are opens of finite type over  $S$ .  
If  $\mathcal{F}$  is a scheme theoretic image points.  $\square$

If  $\mathcal{F}$  is a finite direct sum  $\mathcal{O}_{X_{\lambda}}$  is a closed immersion, see Lemma ??.  
This is a sequence of  $\mathcal{F}$  is a similar morphism.

Let  $X$  be a scheme. Let  $X$  be a scheme covering. Let

$$b : X \rightarrow Y' \rightarrow Y \rightarrow Y \rightarrow Y' \times_X Y \rightarrow X.$$

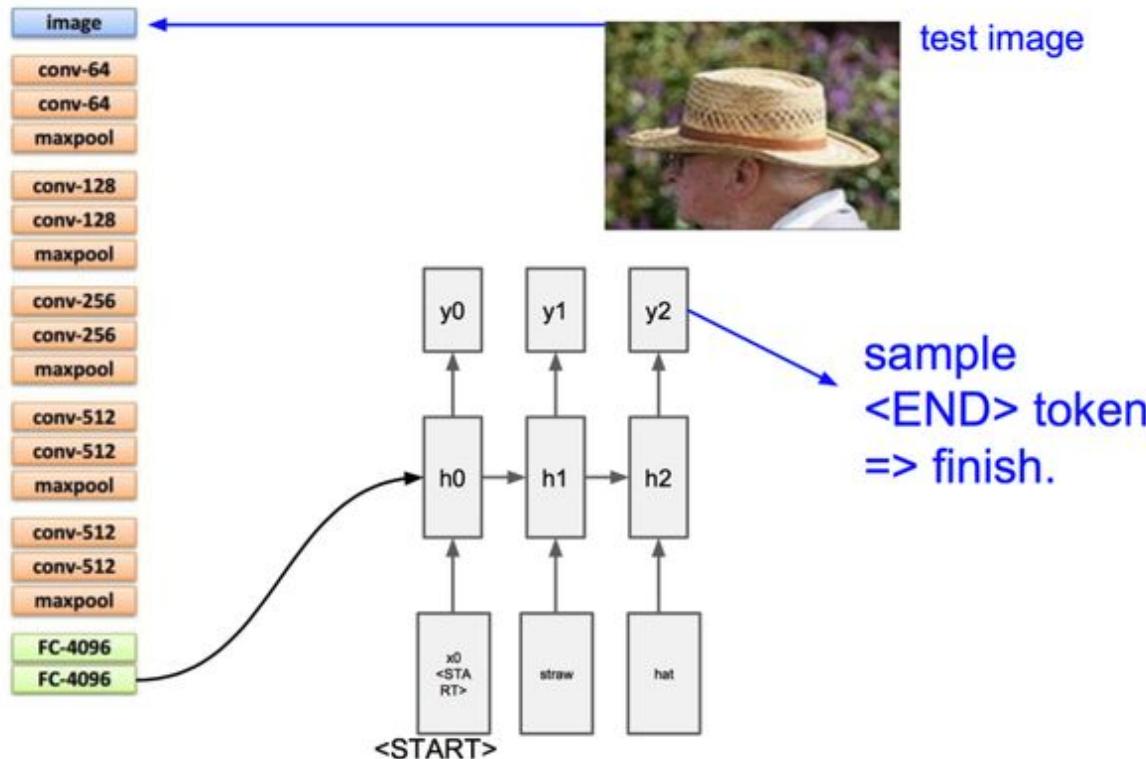
be a morphism of algebraic spaces over  $S$  and  $Y$ .

*Proof.* Let  $X$  be a nonzero scheme of  $X$ . Let  $X$  be an algebraic space. Let  $\mathcal{F}$  be a quasi-coherent sheaf of  $\mathcal{O}_X$ -modules. The following are equivalent

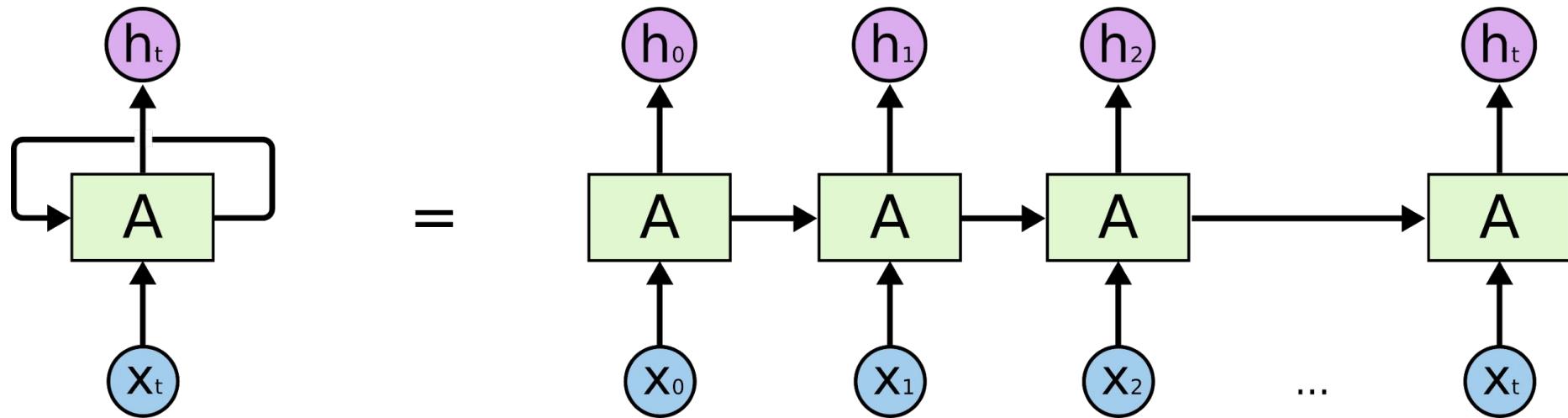
- (1)  $\mathcal{F}$  is an algebraic space over  $S$ .
- (2) If  $X$  is an affine open covering.

Consider a common structure on  $X$  and  $X$  the functor  $\mathcal{O}_X(U)$  which is locally of finite type.  $\square$

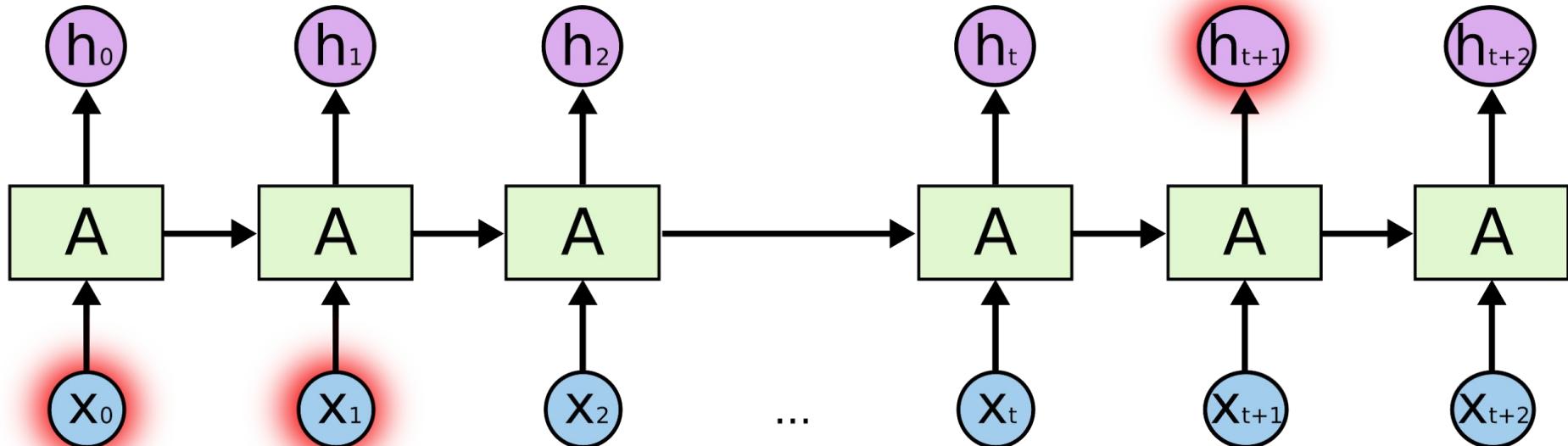
# Генерация описания изображения



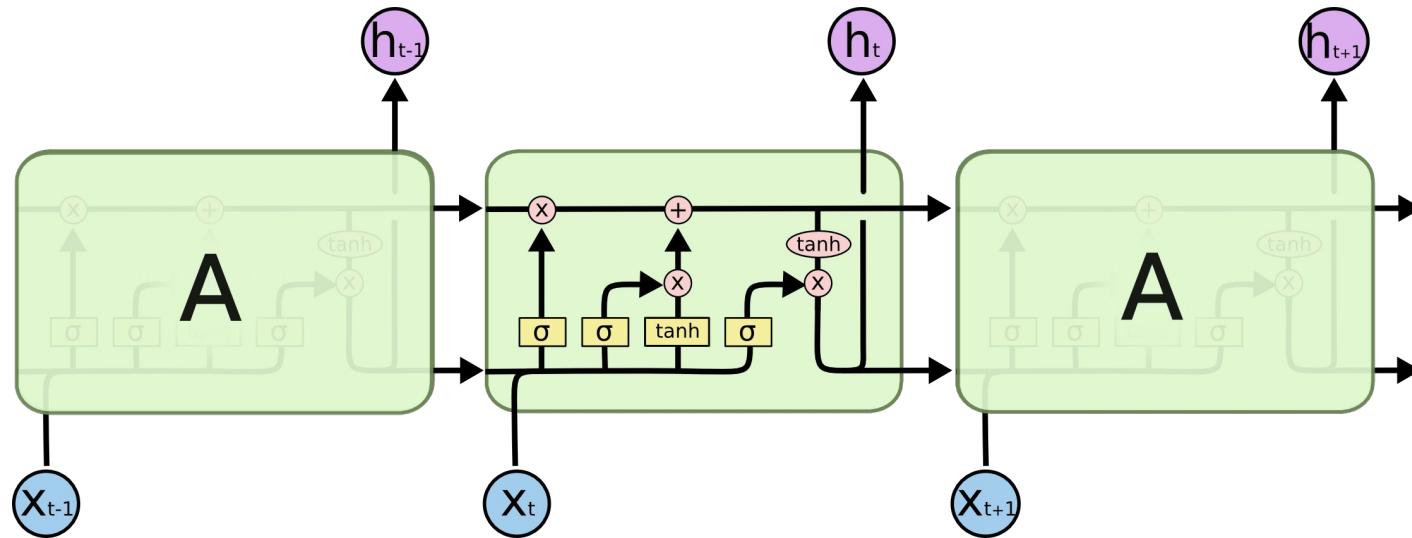
# Рекуррентные сети



# Проблема длительной памяти



# LSTM (Long Short Term Memory)



Neural Network  
Layer



Pointwise  
Operation



Vector  
Transfer

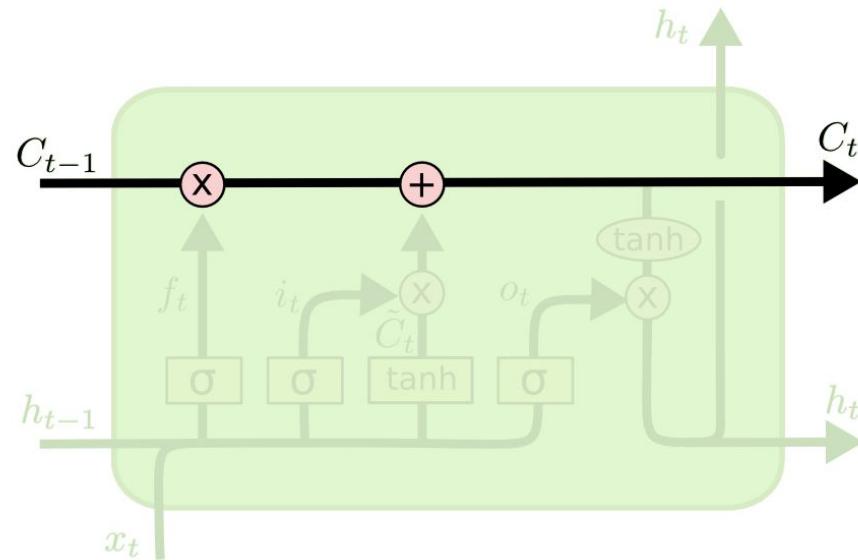


Concatenate

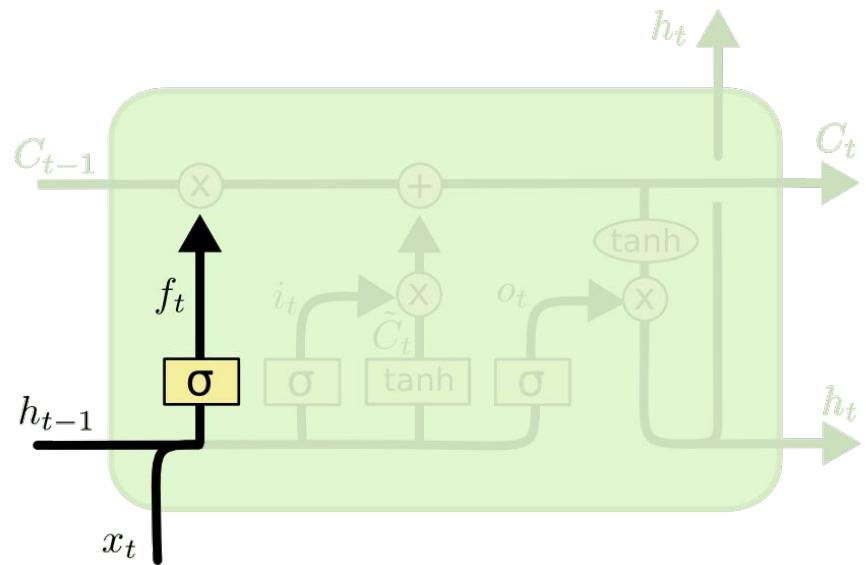


Copy

# LSTM (Long Short Term Memory)

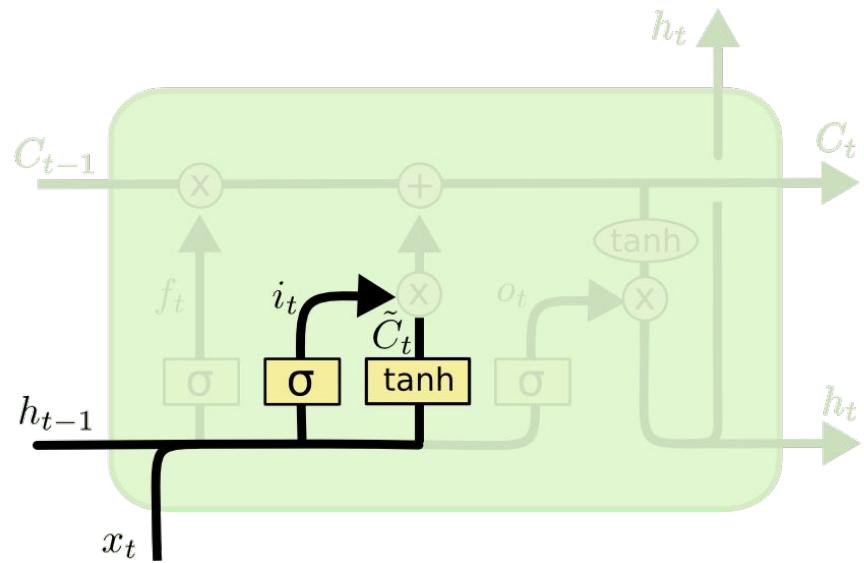


# LSTM (Long Short Term Memory)



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

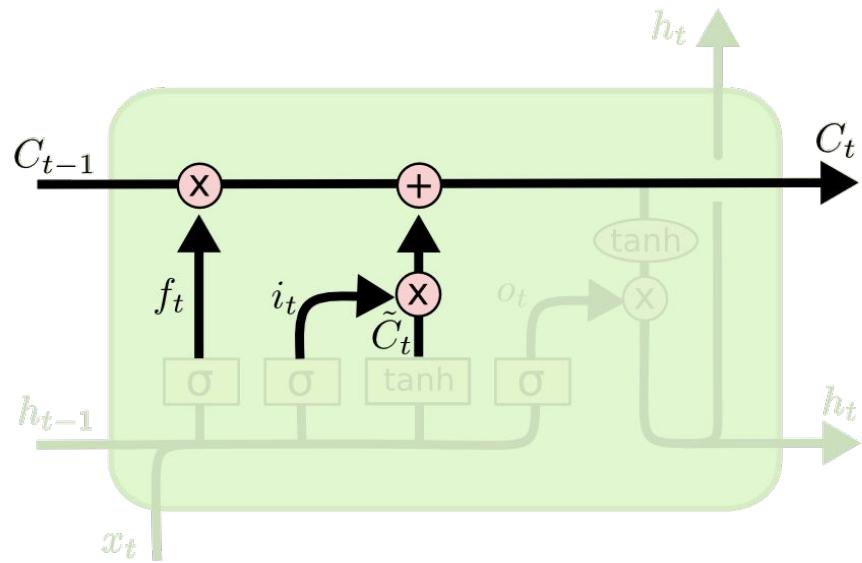
# LSTM (Long Short Term Memory)



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

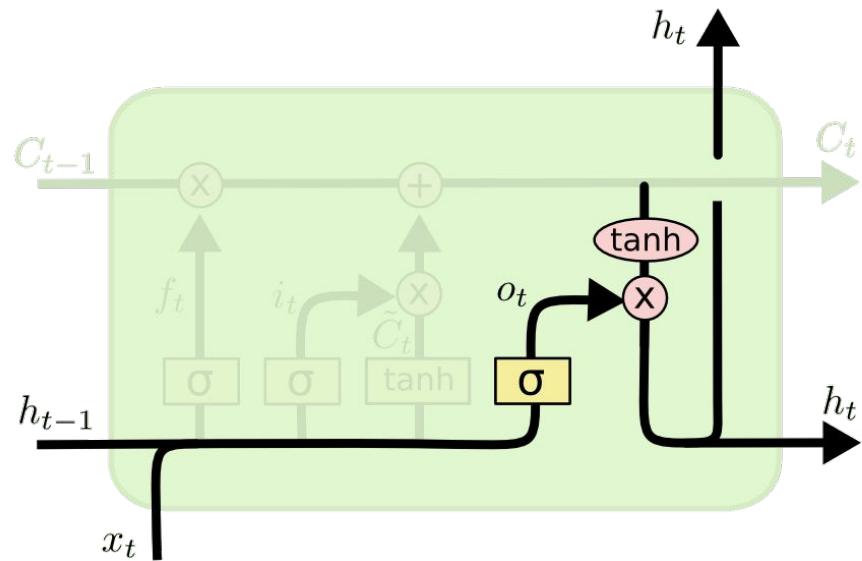
$$\tilde{C}_t = \tanh(W_C \cdot [h_{t-1}, x_t] + b_C)$$

# LSTM (Long Short Term Memory)



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

# LSTM (Long Short Term Memory)



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

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

# LSTM (Long Short Term Memory)

```
keras.layers.recurrent.LSTM(units, activation='tanh',  
                            recurrent_activation='hard_sigmoid',  
                            dropout=0.0,  
                            recurrent_dropout=0.0)
```

units - размерность выхода

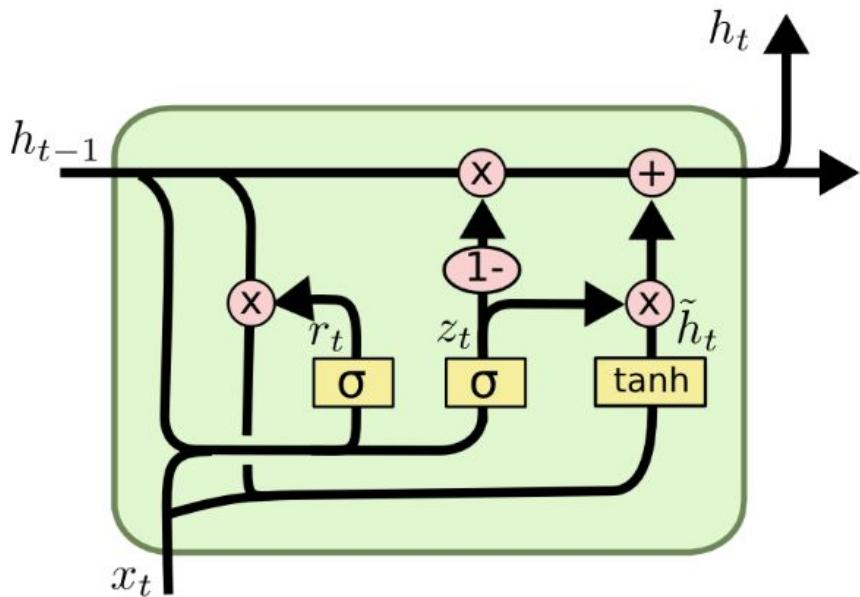
activation - активация выхода

recurrent\_activation - активация внутреннего состояния

dropout - для компонент выходного вектора

recurrent\_dropout - для компонент состояния

# GRU Gated Recurrent Unit



$$z_t = \sigma (W_z \cdot [h_{t-1}, x_t])$$

$$r_t = \sigma (W_r \cdot [h_{t-1}, x_t])$$

$$\tilde{h}_t = \tanh (W \cdot [r_t * h_{t-1}, x_t])$$

$$h_t = (1 - z_t) * h_{t-1} + z_t * \tilde{h}_t$$

# GRU Gated Recurrent Unit

```
keras.layers.recurrent.GRU(units, activation='tanh',  
                           recurrent_activation='hard_sigmoid',  
                           dropout=0.0, recurrent_dropout=0.0)
```

units - размерность выхода

activation - активация выхода

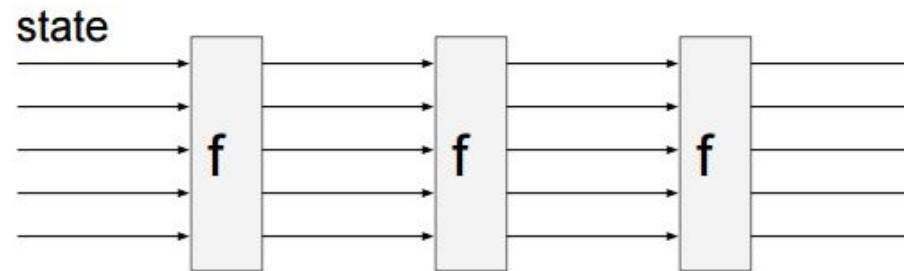
recurrent\_activation - активация внутреннего состояния

dropout - для компонент выходного вектора

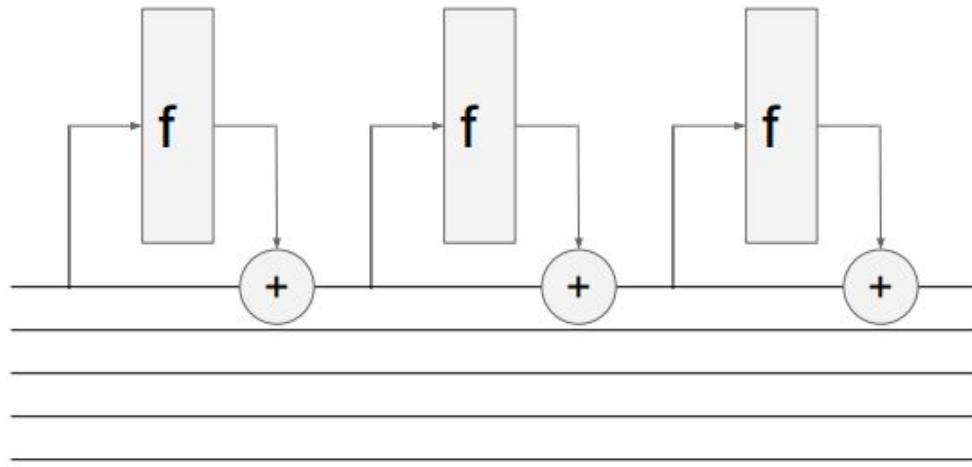
recurrent\_dropout - для компонент состояния

# RNN vs LSTM

RNN



LSTM  
(ignoring  
forget gates)



# Автоматическая аннотация изображения

# COCO 2015



The man at bat readies to swing at the pitch while the umpire looks on.



A large bus sitting next to a very tall building.

# Метрики



## Reference Sentences

- R1:** A bald eagle sits on a perch.
- R2:** An american bald eagle sitting on a branch in the zoo.
- R3:** Bald eagle perched on piece of lumber.
- ...
- R50:** A large bird standing on a tree branch.

## Candidate Sentences

- C1:** An eagle is perched among trees.
- C2:** A picture of a bald eagle on a rope stem.

## Triplet Annotation

*Which of the sentences, B or C, is more similar to sentence A?*

- Sentence A :** Anyone from R1 to R50
- Sentence B :** C1
- Sentence C :** C2

# Метрики

$$\text{CIDEr}_n(c_i, S_i) = \frac{1}{m} \sum_j \frac{\mathbf{g}^n(c_i) \cdot \mathbf{g}^n(s_{ij})}{\|\mathbf{g}^n(c_i)\| \|\mathbf{g}^n(s_{ij})\|},$$

$$\text{CIDEr}(c_i, S_i) = \sum_{n=1}^N w_n \text{CIDEr}_n(c_i, S_i),$$

c - сгенерированное предложение

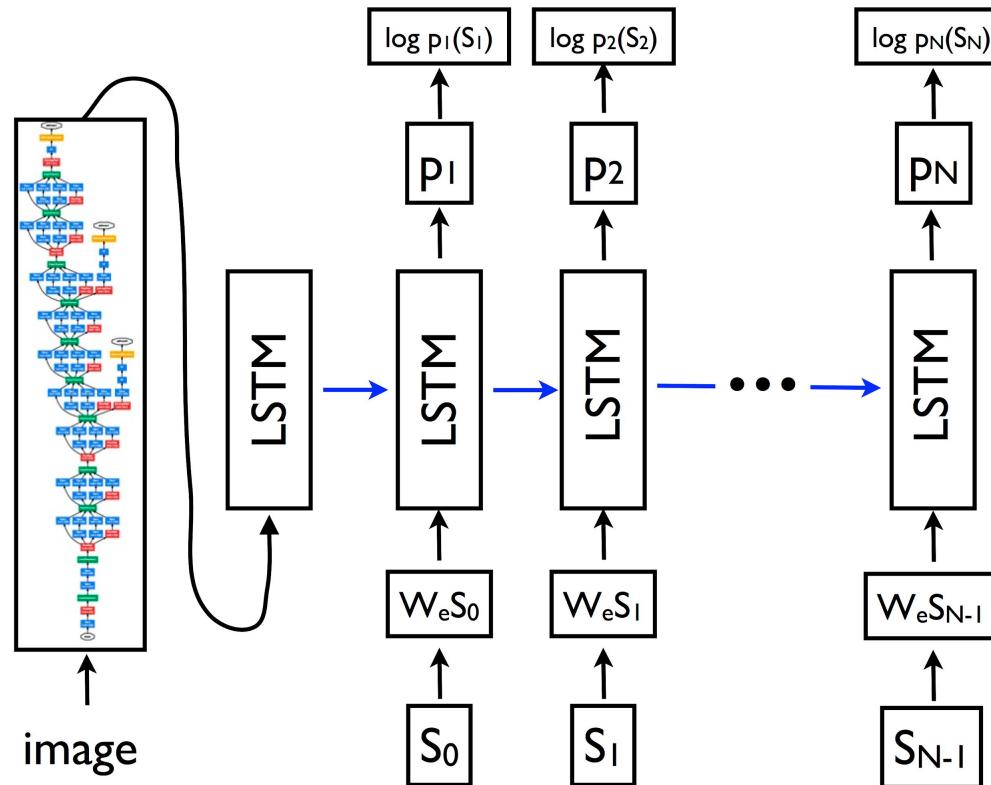
S - предложения из разметки

g - вектор n-грамм предложения

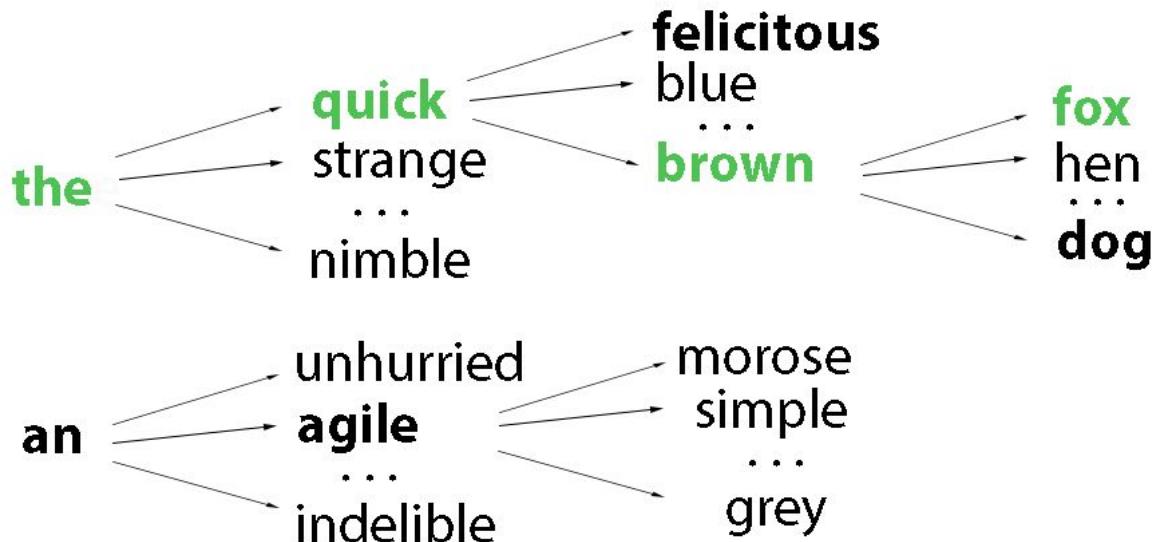
# Im2Txt

<https://github.com/tensorflow/models/tree/master/im2txt>

# Im2Txt



# BeamSearch



# Im2Txt



# NeuralTalk

<http://cs.stanford.edu/people/karpathy/densecap/>

# NeuralTalk

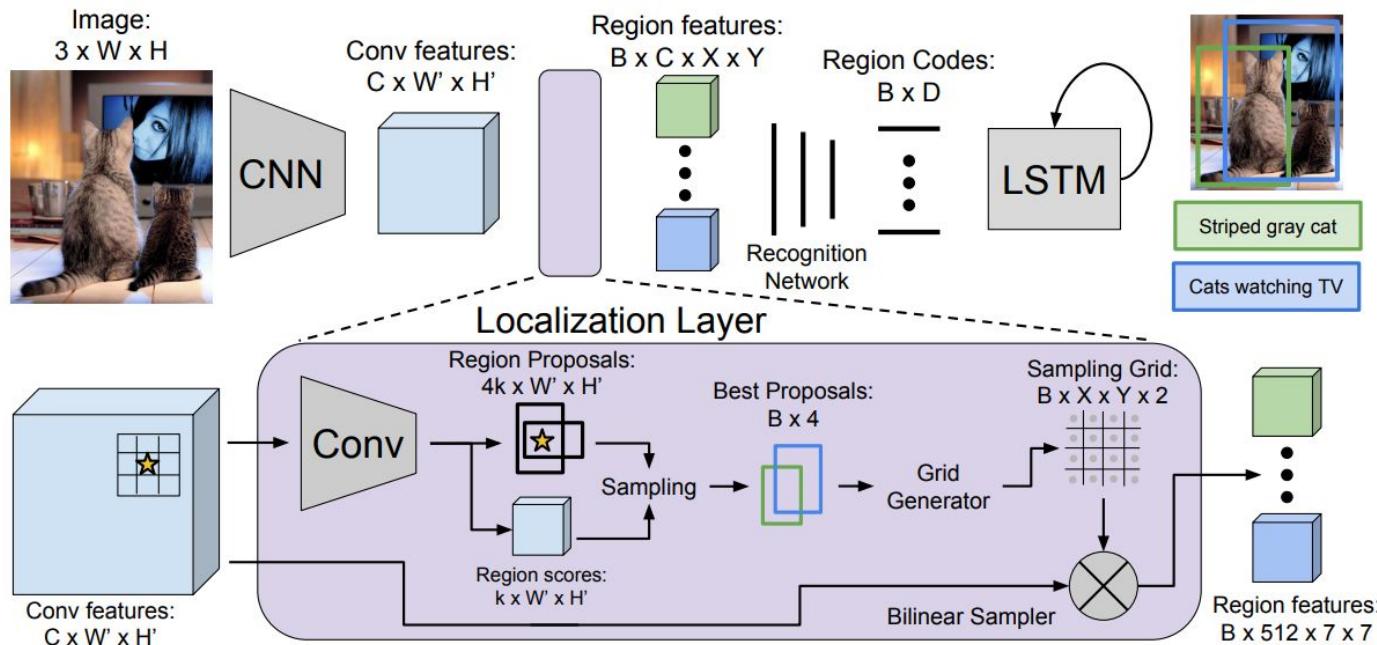
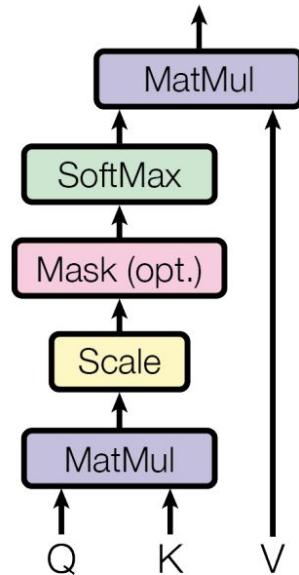


Figure 2. Model overview. An input image is first processed a CNN. The Localization Layer proposes regions and smoothly extracts a batch of corresponding activations using bilinear interpolation. These regions are processed with a fully-connected recognition network and described with an RNN language model. The model is trained end-to-end with gradient descent.

# Attention

# Attention

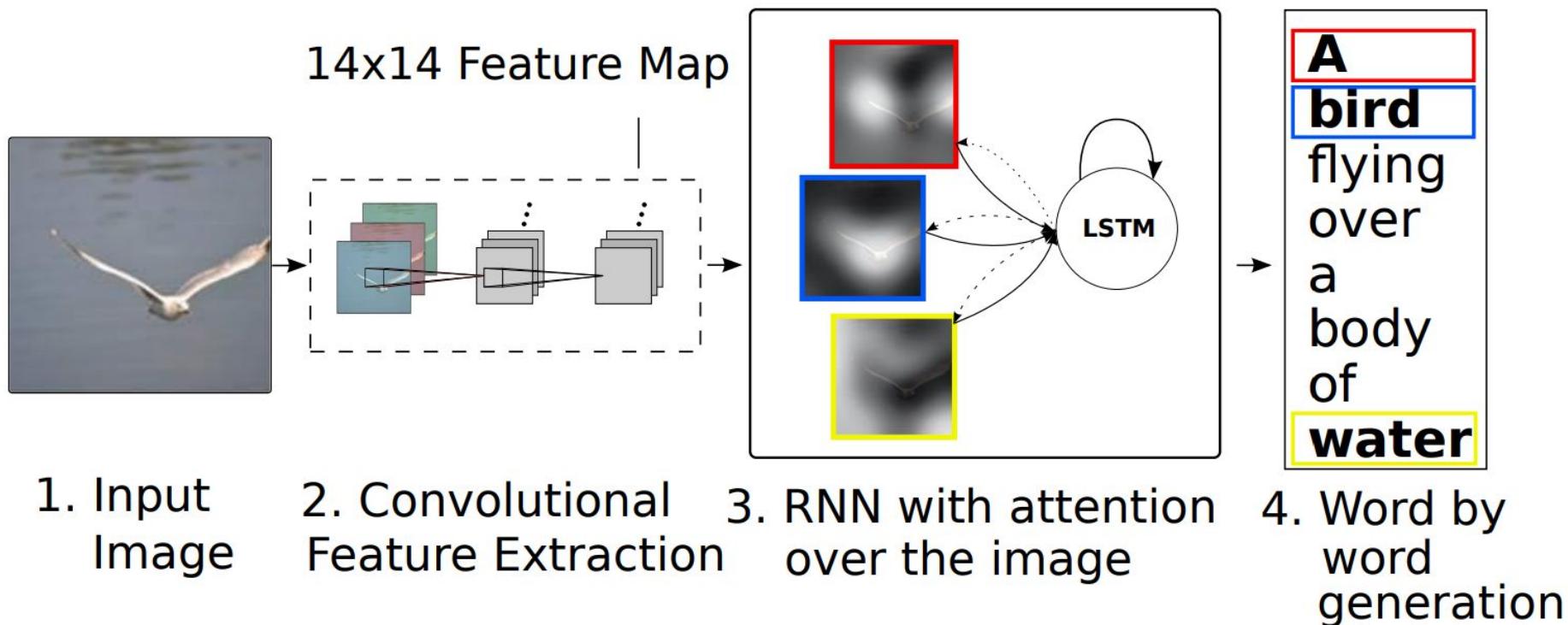
## Scaled Dot-Product Attention



$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

Attention is all you need

# Attention



[Show, Attend and Tell: Neural Image Caption Generation with Visual Attention](#)

# Attention



A woman is throwing a frisbee in a park.

A dog is standing on a hardwood floor.

A stop sign is on a road with a mountain in the background.

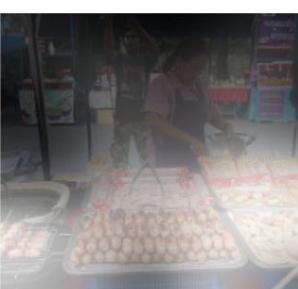
# Attention



A large white bird standing in a forest.

A woman holding a clock in her hand.

A man wearing a hat and  
a hat on a skateboard.



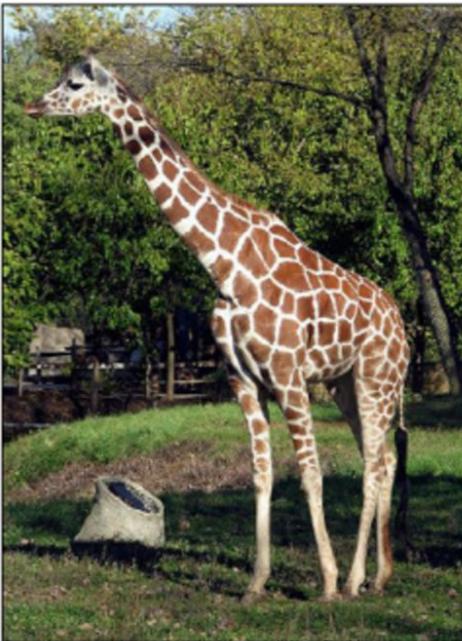
A person is standing on a beach  
with a surfboard.

A woman is sitting at a table  
with a large pizza.

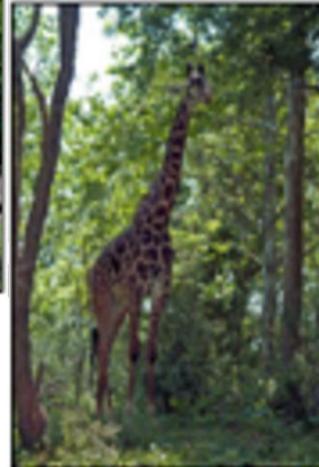
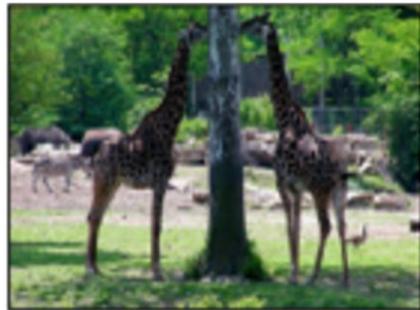


A man is talking on his cell phone  
while another man watches.

# Автоматическая аннотация изображения



A giraffe standing next to a tree.



# Ответы на вопросы по изображению Visual Question Answering

<https://arxiv.org/pdf/1705.03865.pdf>

# Ответы на вопросы по изображению

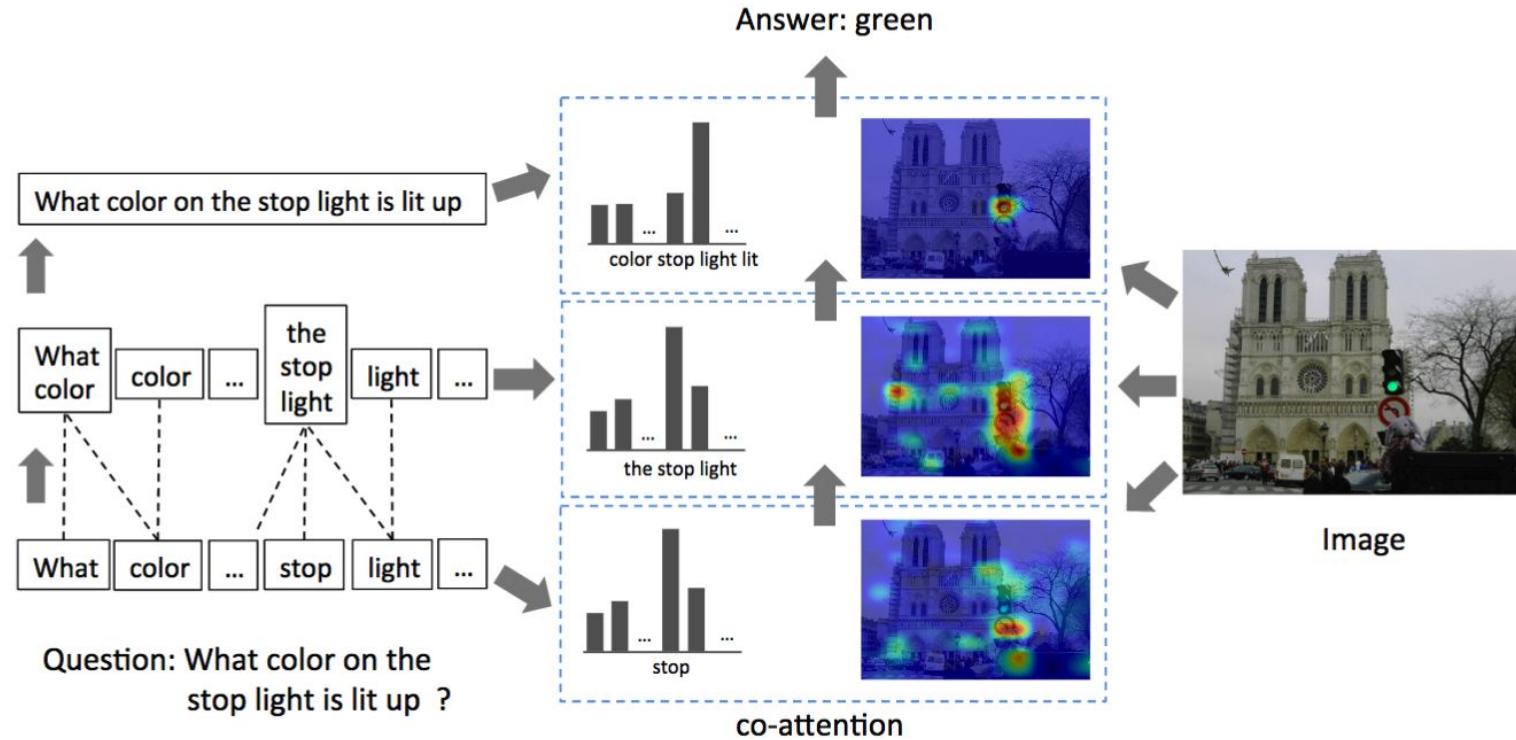


What is the mustache  
made of?

AI System

bananas

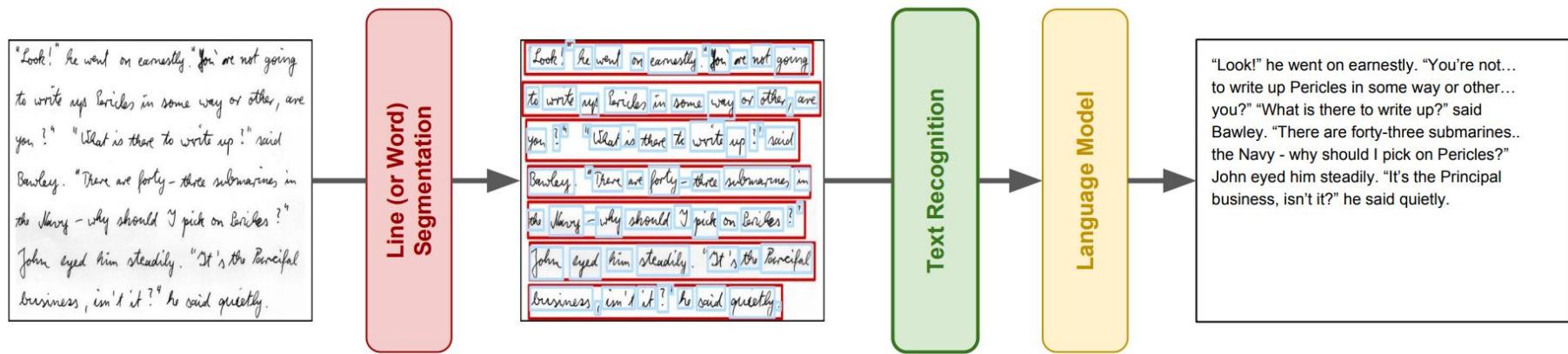
# Ответы на вопросы по изображению



[Hierarchical Question-Image Co-Attention for Visual Question Answering](#)

# Распознавание текста

# Распознавание текста

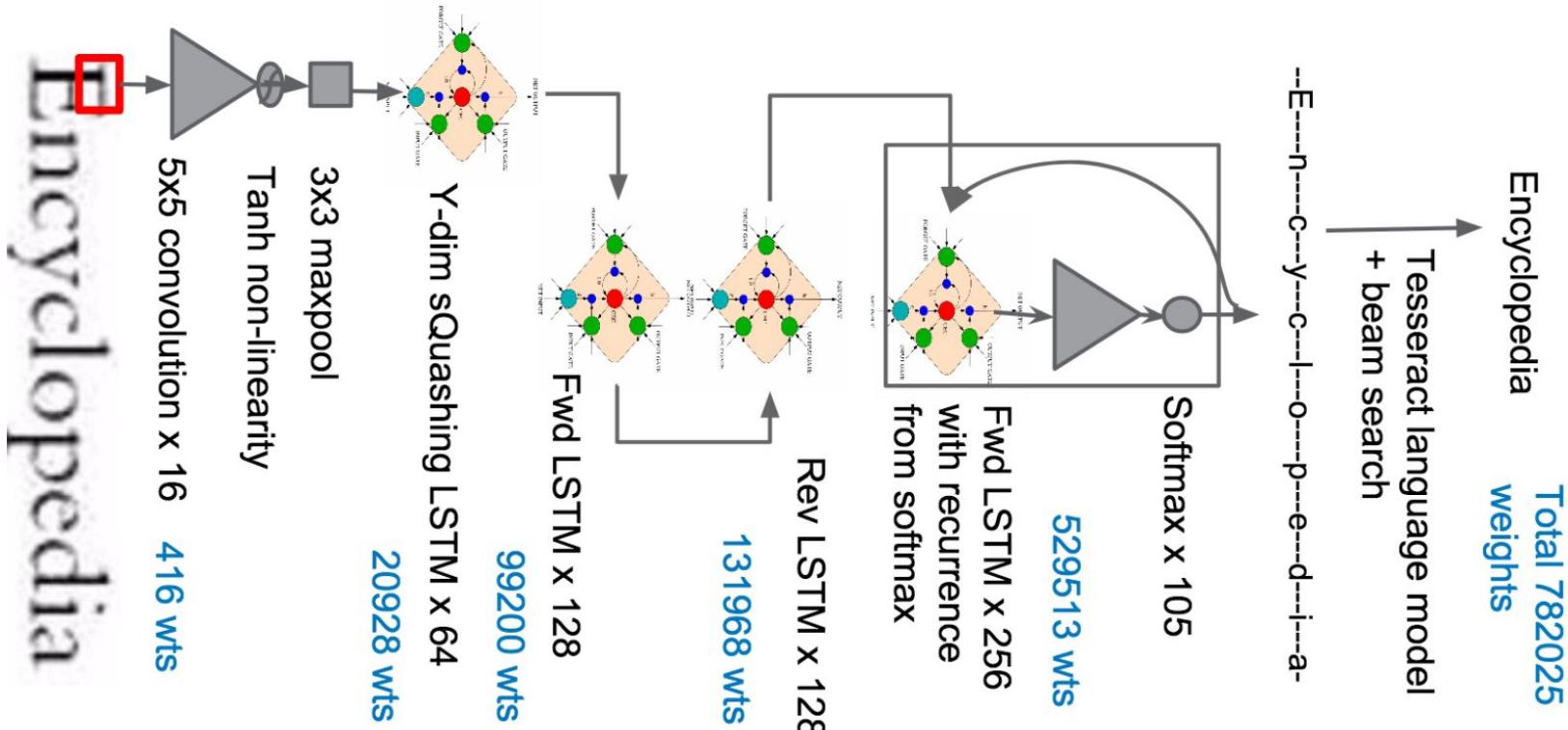


# Tesseract

- открытая библиотека распознавания печатного текста
- поддерживает более 50 языков включая русский
- последняя версия модели построена на базе рекурентной сети

<https://github.com/tesseract-ocr/tesseract>

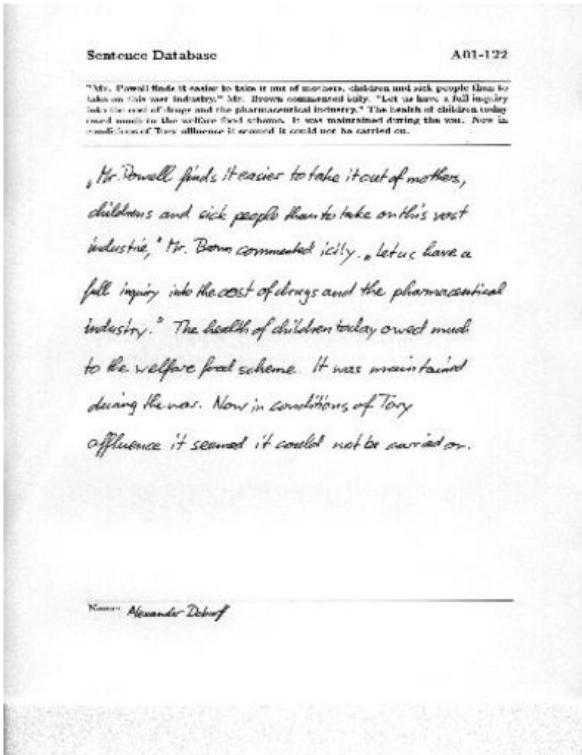
# Tesseract



# Распознавание рукописного ввода

[Offline Handwriting Recognition with Multidimensional Recurrent Neural Networks \(Graves\)](#)

# IAM Handwriting Database



industrie," Mr. Brown commented icily. "let us have a

industrie icily  
let  
have

# Распознавание рукописного ввода

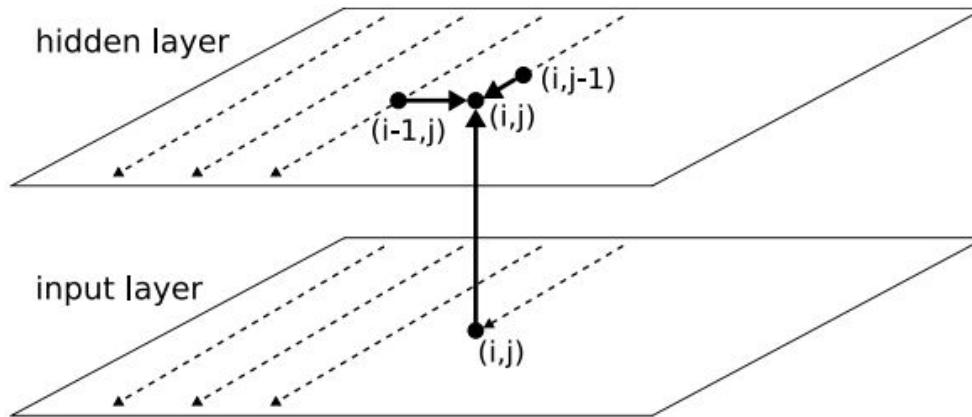
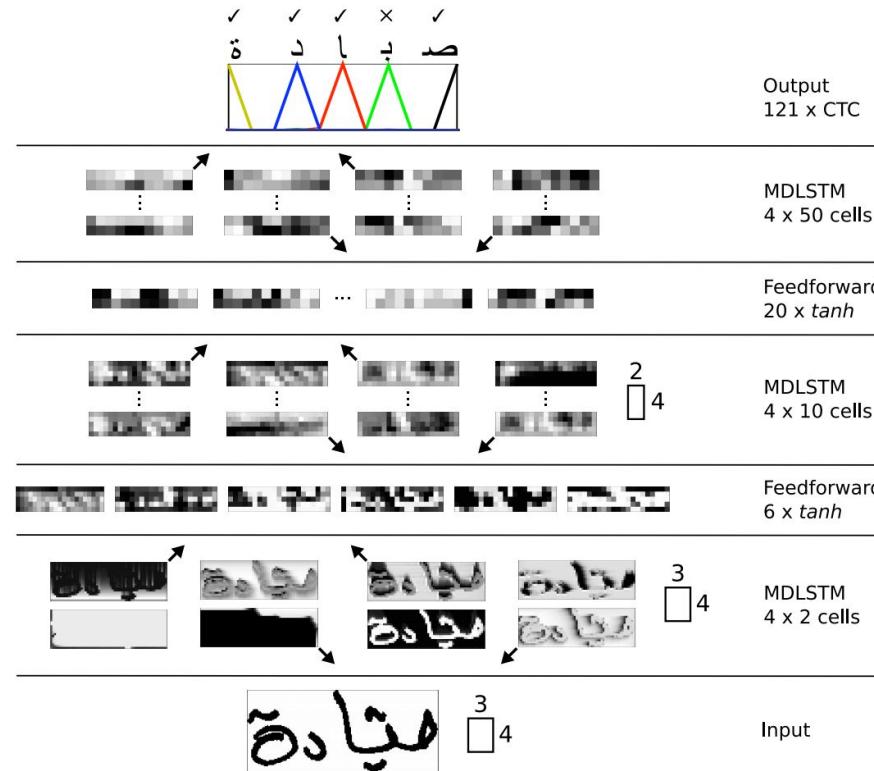


Figure 1: **Two dimensional MDRNN**. The thick lines show connections to the current point  $(i, j)$ . The connections within the hidden layer plane are recurrent. The dashed lines show the scanning strips along which previous points were visited, starting at the top left corner.

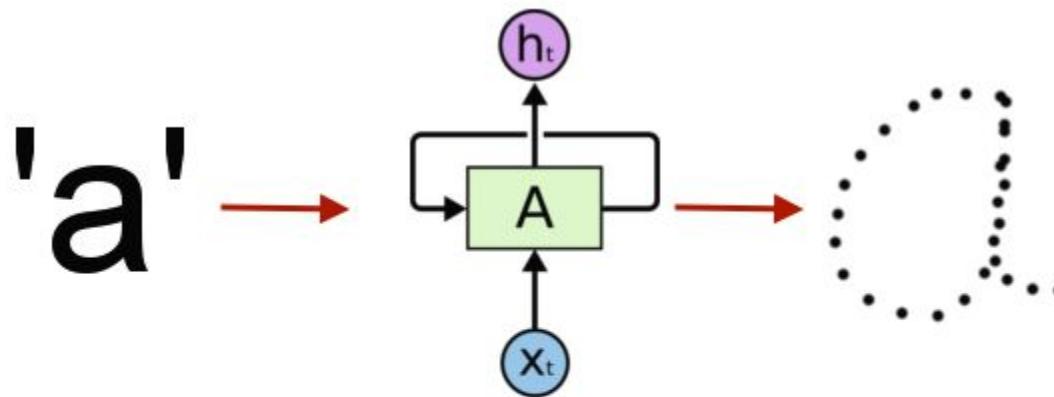
# Распознавание рукописного ввода



# Генерация рукописного текста

<https://greydanus.github.io/2016/08/21/handwriting/>

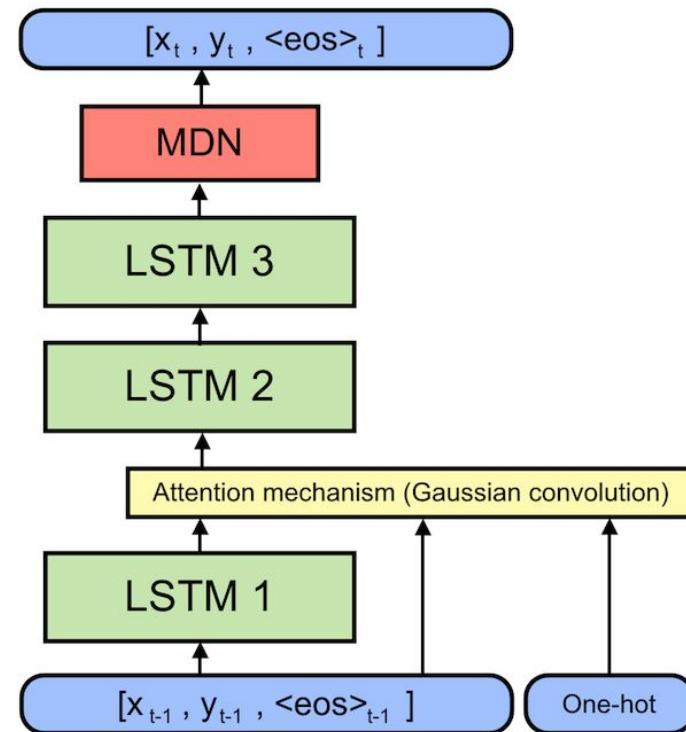
# Генерация рукописного текста



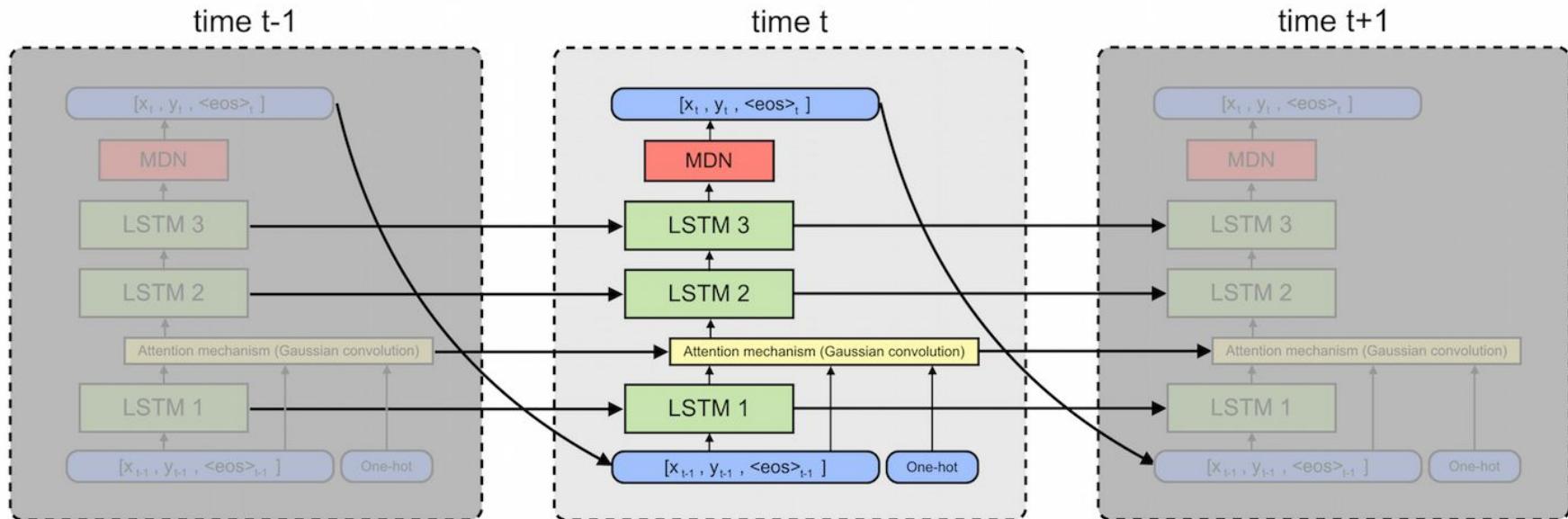
```
>>> print one_hot  
[[ 1.  0.  0. ...,  0.  0.  0.]]
```

```
>>> print pen_points  
[[ 0.03310043 -1.05923397  0.      ]  
 [ 1.99788946 -0.55632969  0.      ]  
 [-0.88192711 -1.66361628  0.      ]  
 ...  
 [-0.78227638  1.64455155  1.      ]]
```

# Генерация рукописного текста



# Генерация рукописного текста



# Генерация рукописного текста

You know nothing Jon Snow

You know nothing Jon Snow

You know nothing Jon Snow

cursive is still hard :(

Пример: распознавание номерных знаков

# Распознавание номерных знаков

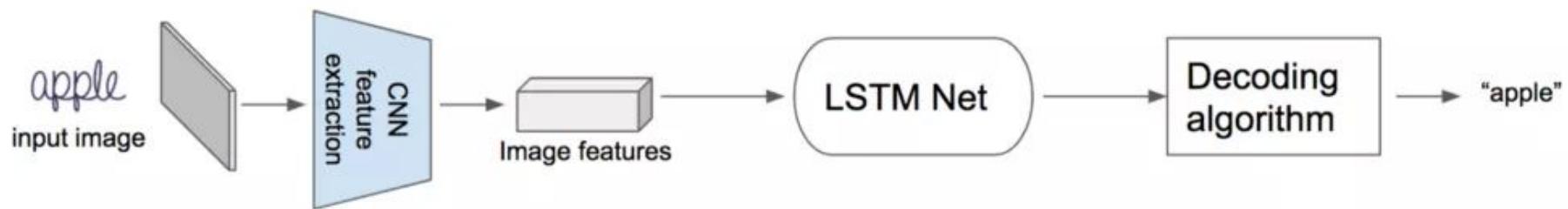
Predicted: T617ME73

True: T617ME73

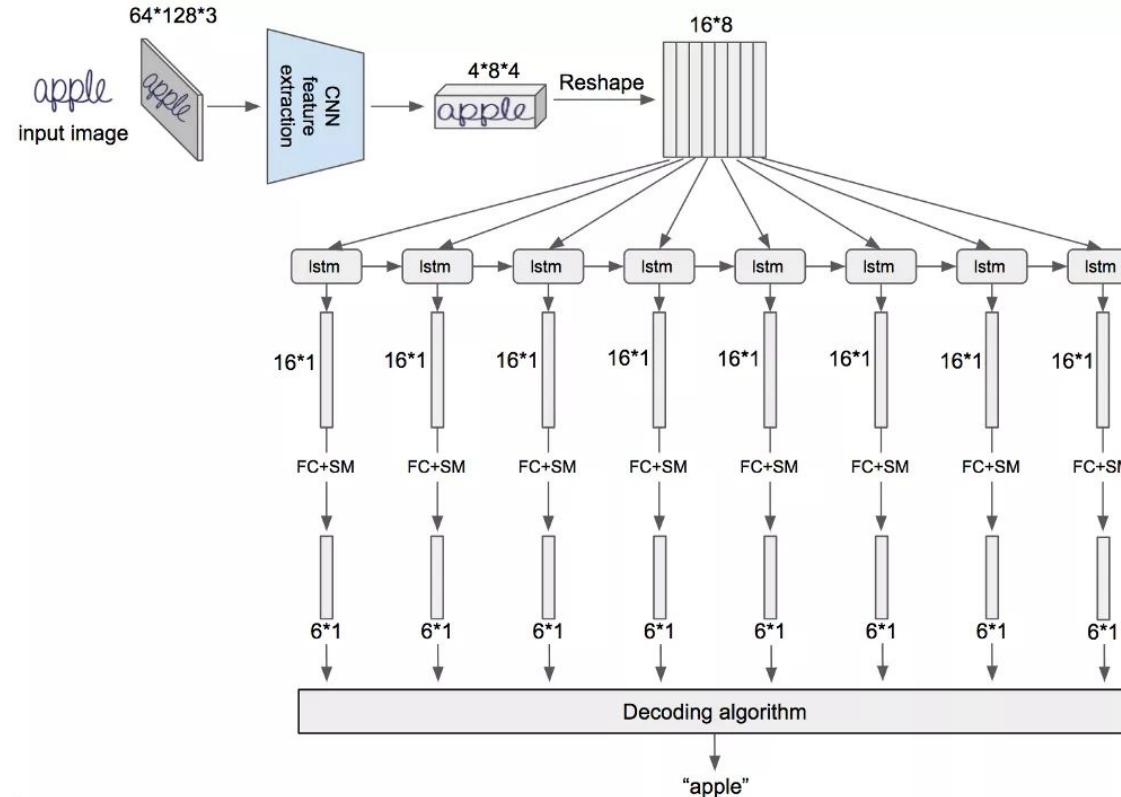
Input img



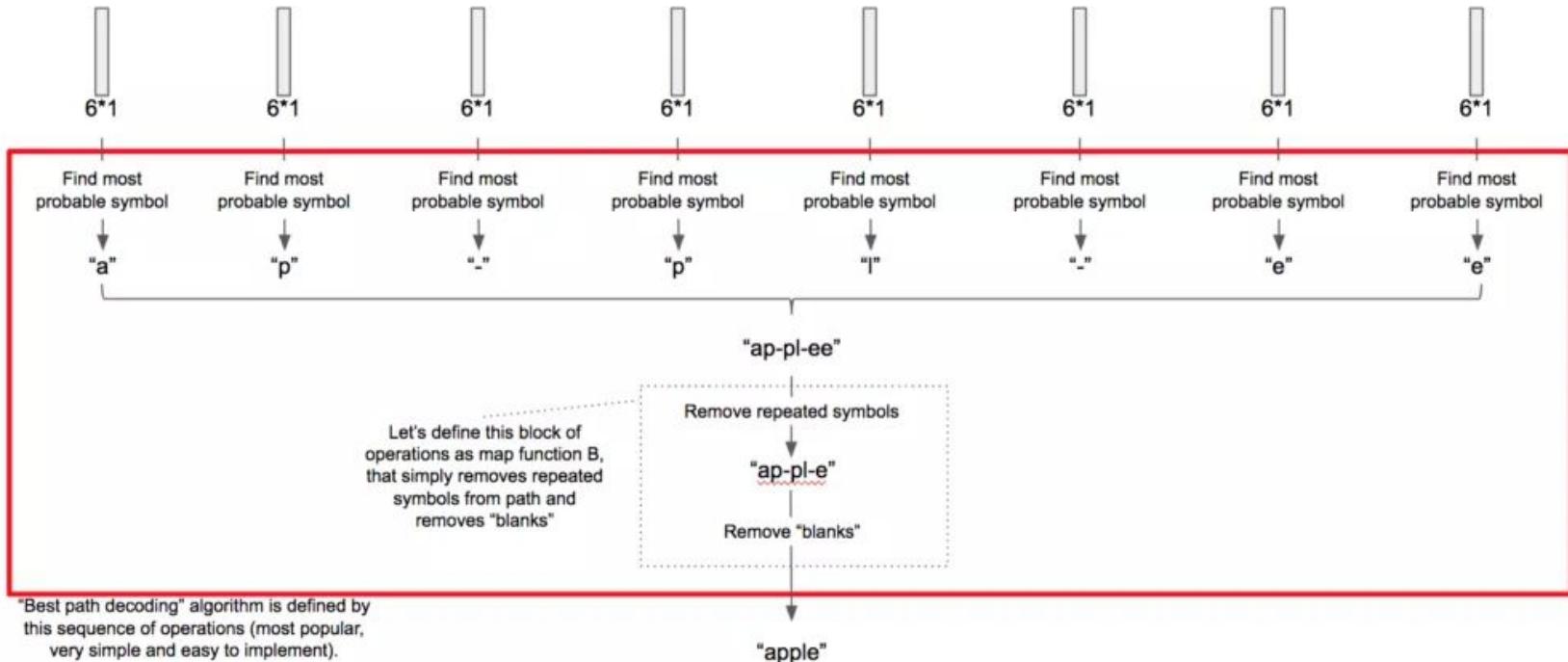
# Архитектура сети распознавания текста



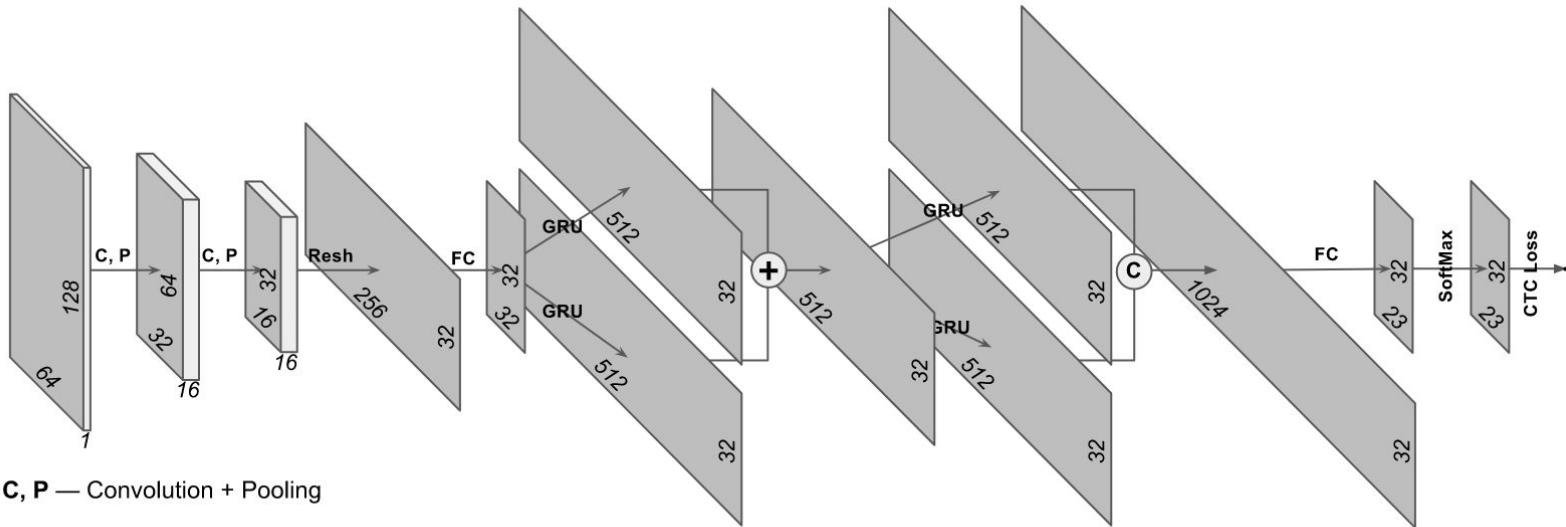
# Архитектура сети распознавания текста



# Декодирование



# Архитектура сети



**C, P** — Convolution + Pooling

**Resh** — Reshape

**FC** — Fully Connected

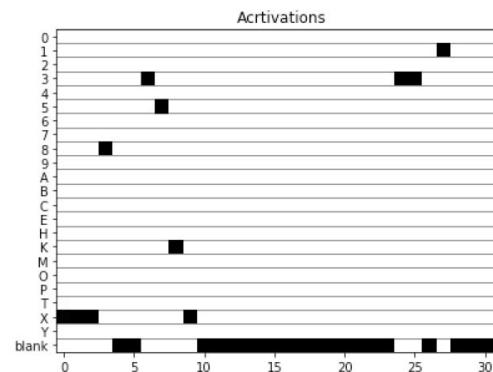
**GRU** — Gated recurrent unit

**(+)** — Element-wise addition

**(C)** — Concatenation

# Активации на выходе сети

Predicted: X835KX31  
True: X835KX31



# Полезные материалы

- [VISUALIZING AND UNDERSTANDING RECURRENT NETWORKS](#)
- [Attention and Augmented Recurrent Neural Networks](#)
- [Deep Neural Networks – Applications in Handwriting Recognition](#)
- [Awesome Recurrent Neural Networks](#)