



Тверской  
государственный  
технический  
университет

# Интеллектуальные информационные системы

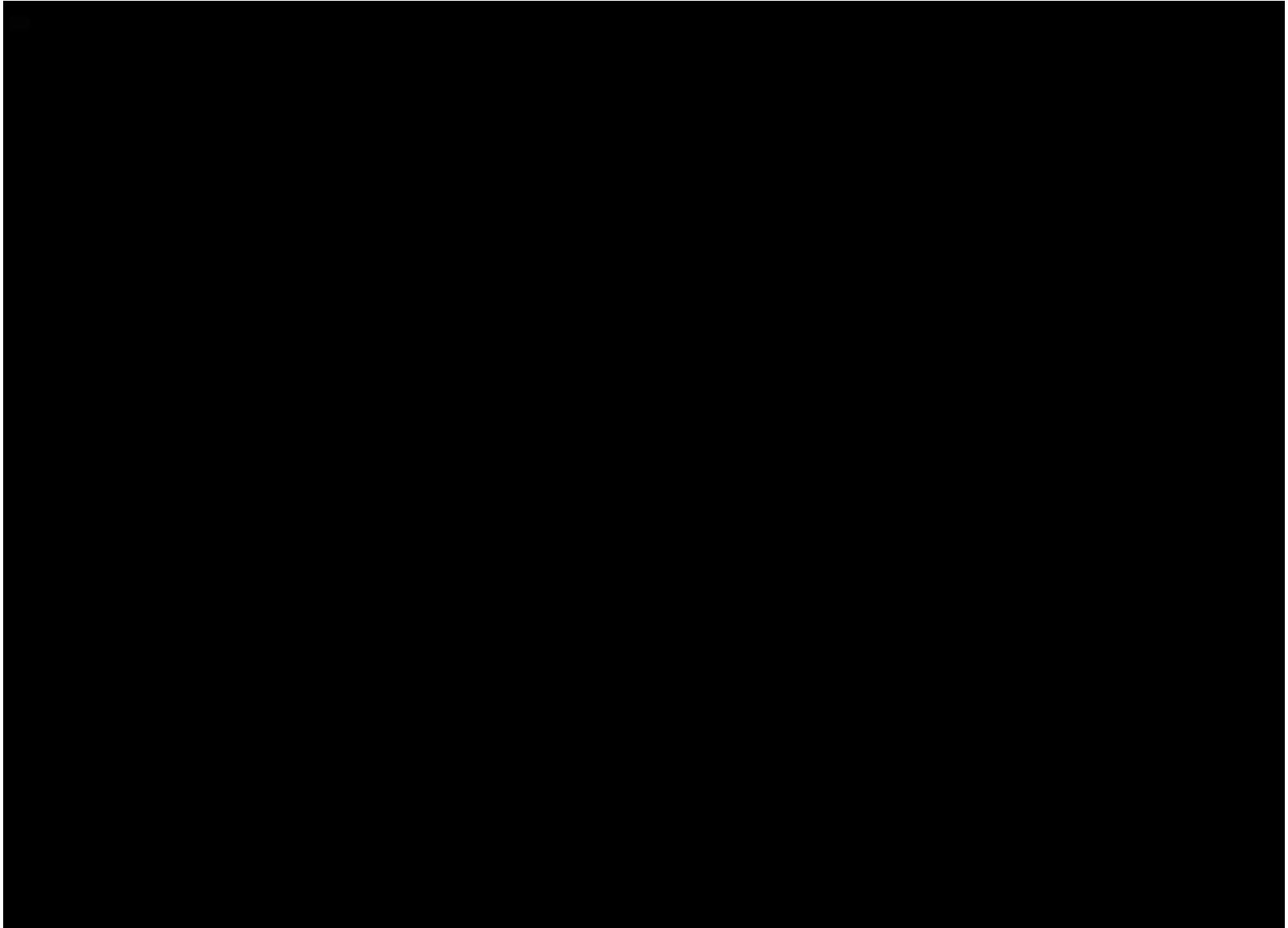
## Сверточные нейронные сети

Материалы курса доступны по ссылке:

<https://github.com/AndreyShpigar/ML-course>

2024 г.

В 1989 году 32-летний Ян ЛеКун представил революционную сверточную нейросеть, способную распознавать цифры в различных стилях написания. Это было логическим продолжением его идеи, предложенной годом ранее, в 1988-м.

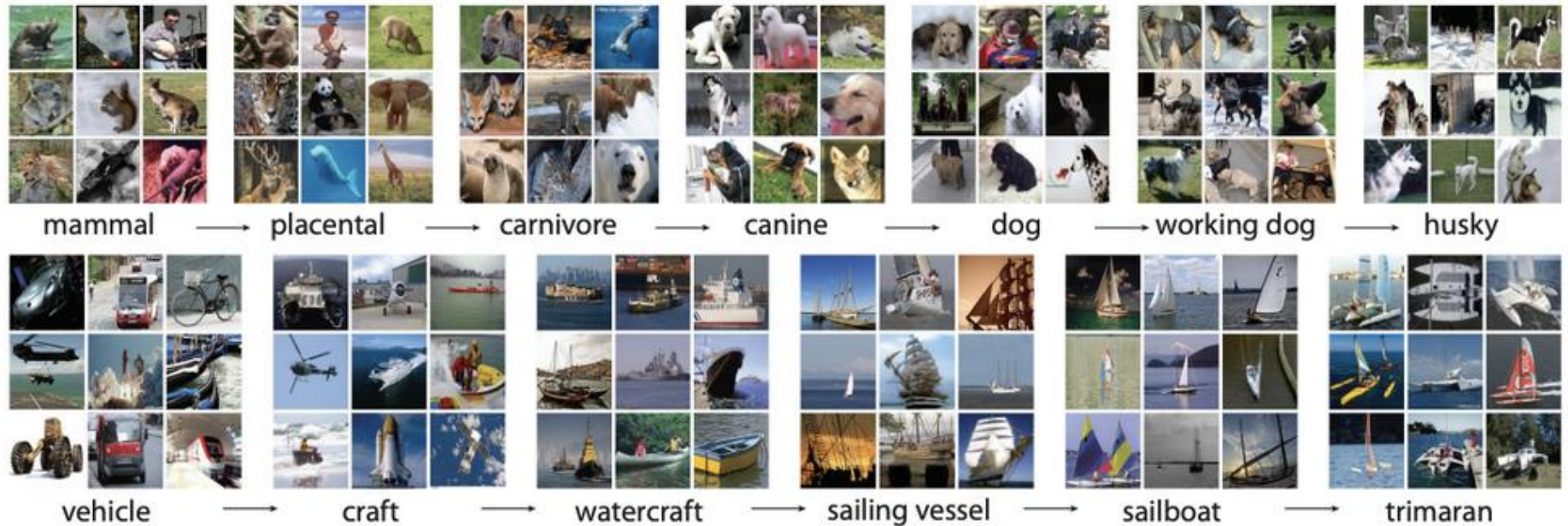


The **ImageNet** project is a large visual database designed for use in visual object recognition software research

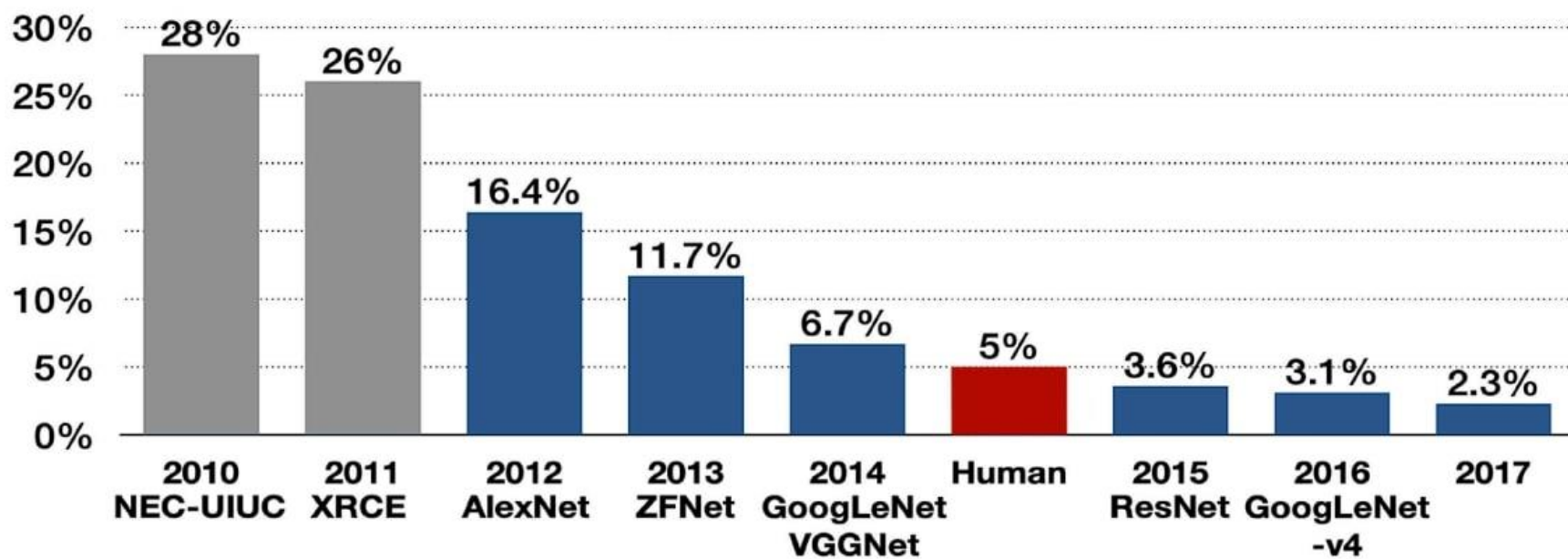


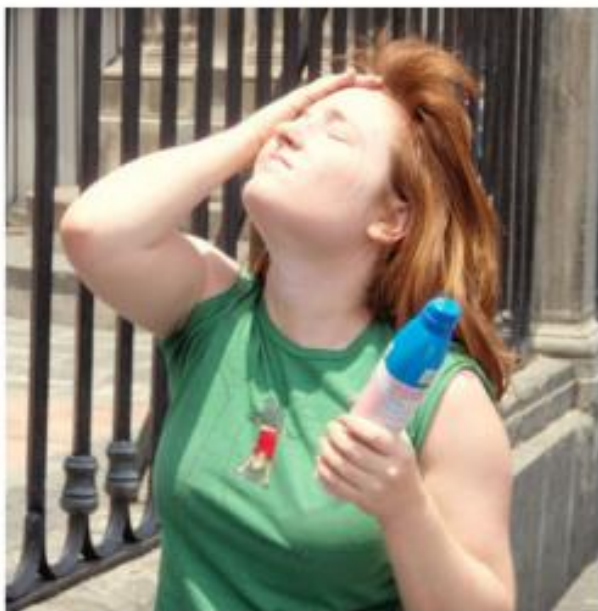


# ImageNet: A Large-Scale Hierarchical Image Database



## Top-5 error





**GT: sunscreen**

1: hair spray

2: ice lolly

3: sunscreen

4: water bottle

5: lotion



**GT: flute**

1: flute

2: oboe

3: panpipe

4: trombone

5: bassoon



**GT: wooden spoon**

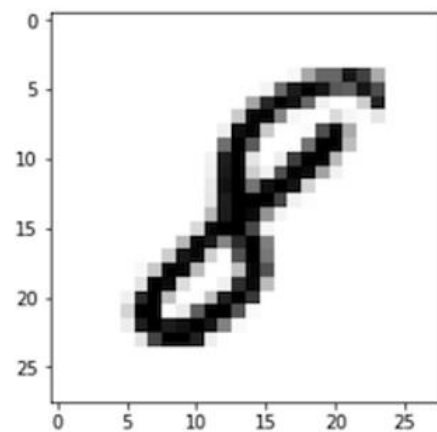
1: wok

2: frying pan

3: spatula

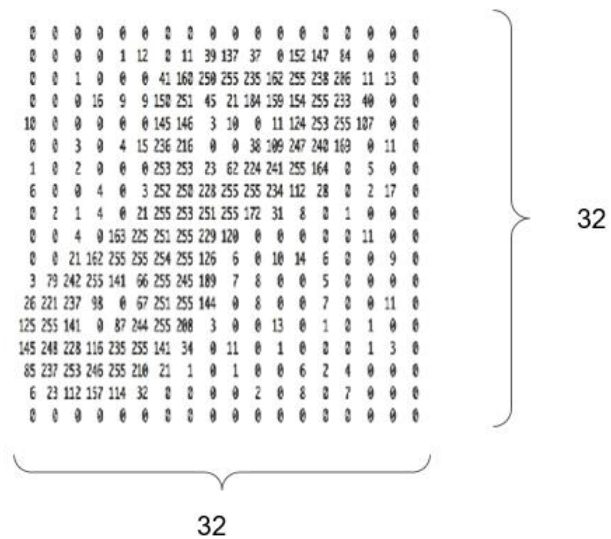
4: wooden spoon

5: hot pot

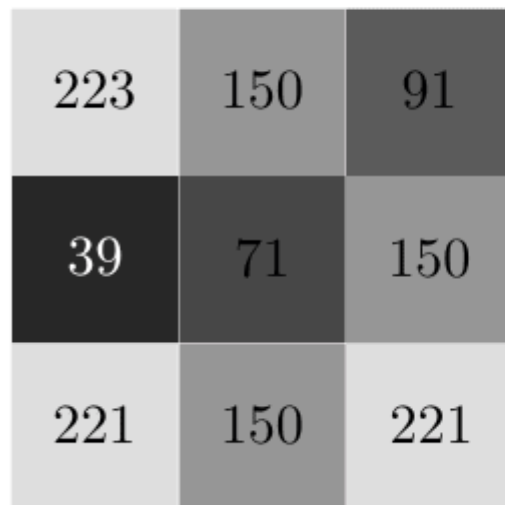


\_\_\_\_\_

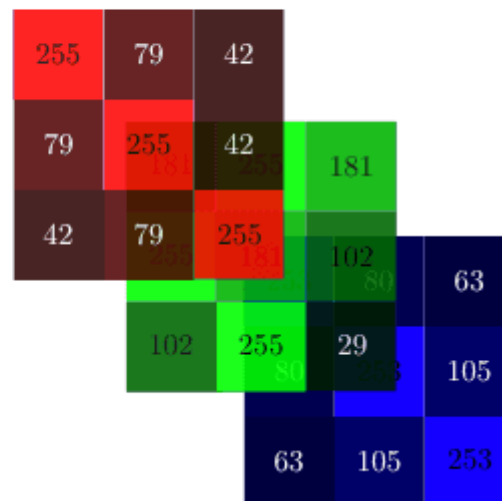
\_\_\_\_\_



## Grayscale



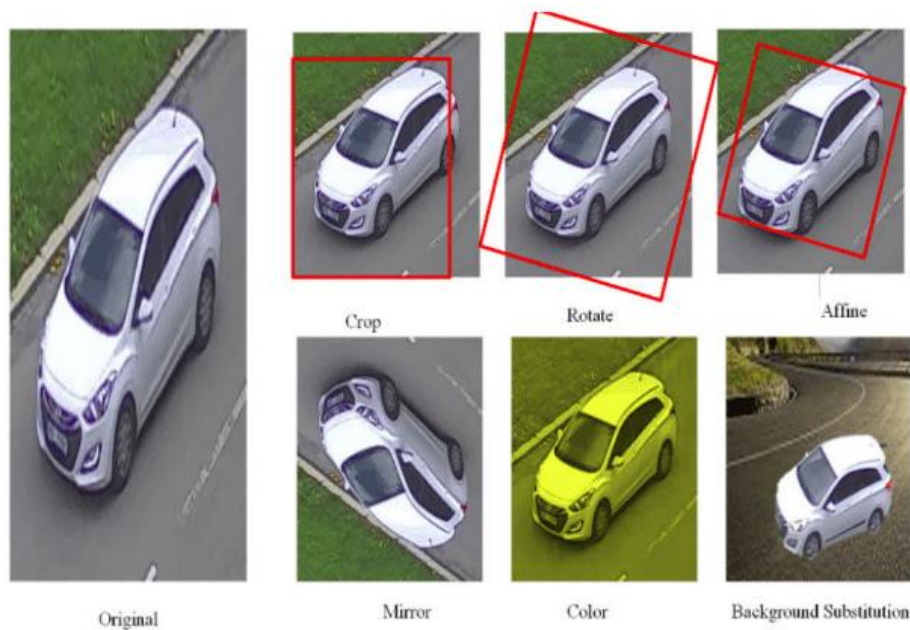
RGB





## Классификация с помощью полносвязной НС с кросс-энтропийной функцией потерь:

- Никак не учитывается структура данных - хочется получить инвариантность к различным преобразованиям исходного изображения

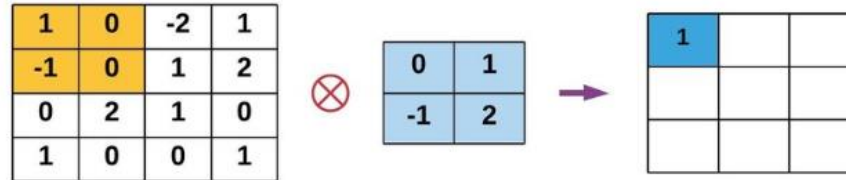


- Большое количество параметров (нейронов в первом слое сети). Например, RGB изображение размером 128x128 это вектор из 49152 элементов

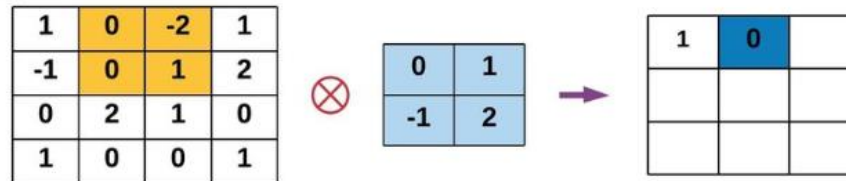


# Операция свертки - поэлементное умножение и суммирование

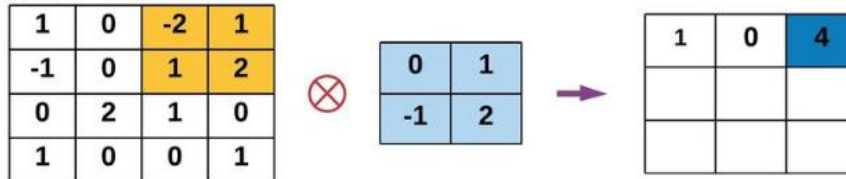
Step-1



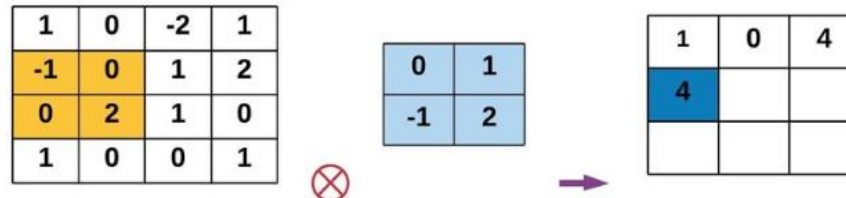
Step-2



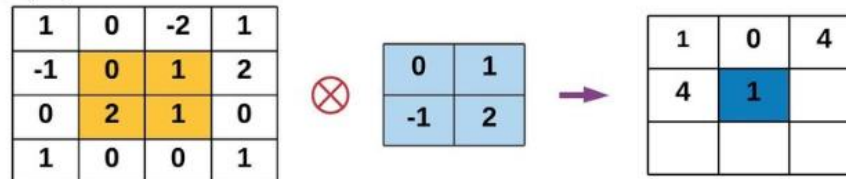
Step-3



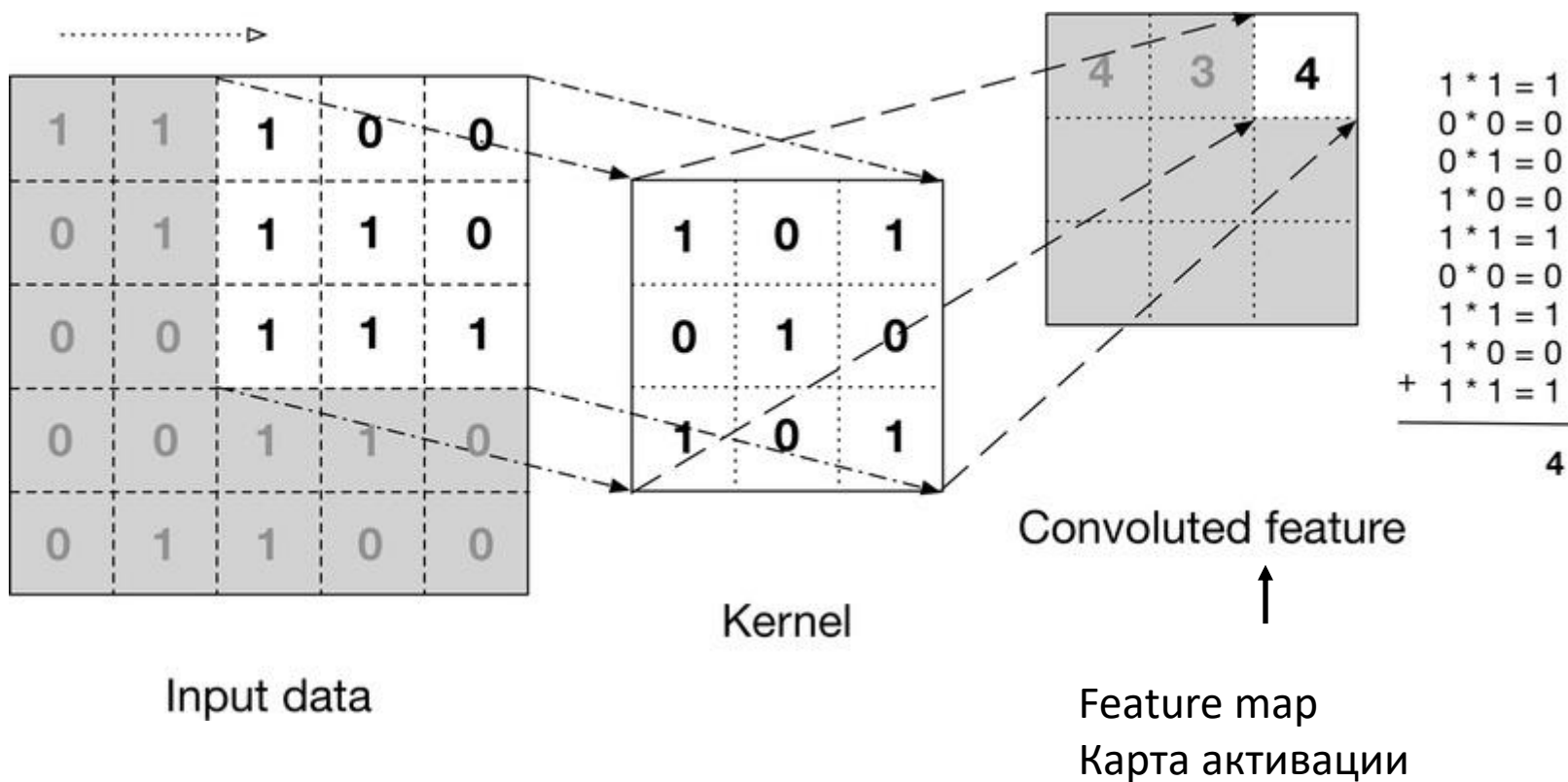
Step-4



Step-5



С чем сворачиваем? С ядром свертки (фильтром)



# Strided Convolution



Input Image

1	9	8	4	4	5	7
4	8	6	7	9	1	7
4	0	5	9	3	8	4
7	3	5	9	0	5	4
7	4	1	1	8	1	2
7	6	6	9	8	7	6
3	6	3	5	4	2	7

Kernel

0	-1	0
-1	5	-1
0	-1	0

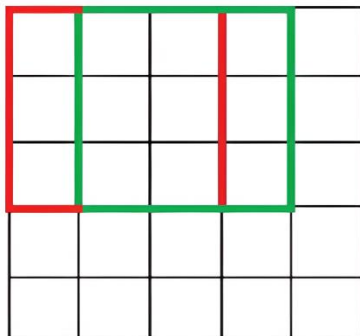
\*

=

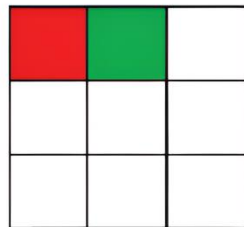
Feature Map


- Stride – величина смещения фильтра

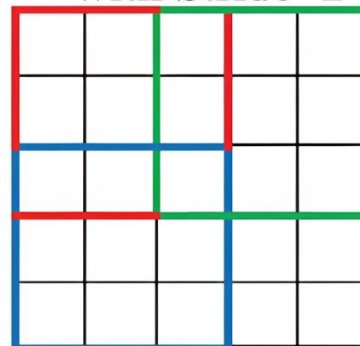
Convolution  
with Stride=1



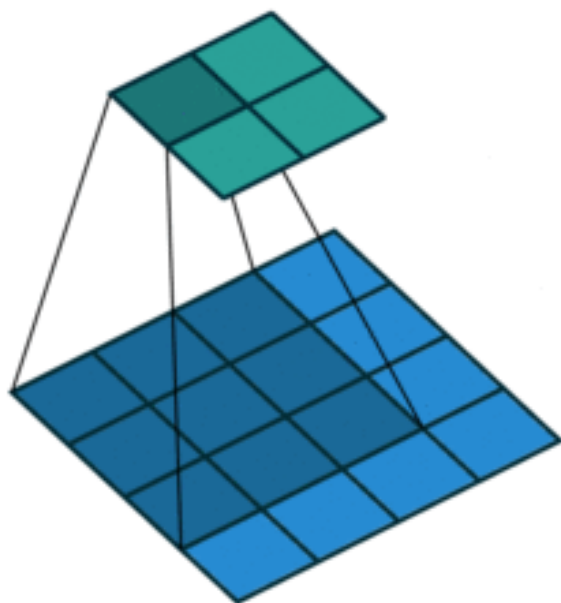
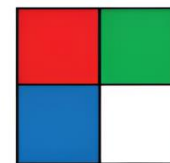
Output



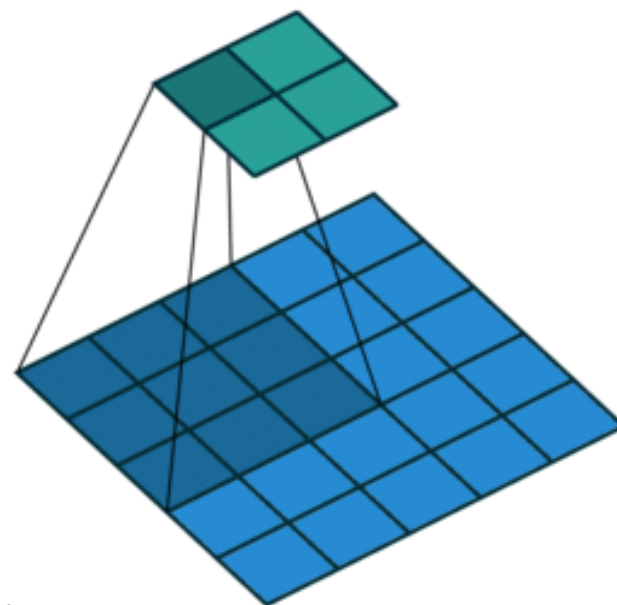
Convolution  
with Stride=2



Output



Stride = 1



Stride = 2



- Padding – дополнение исходного изображения пикселями (рамка вокруг изображения)

0	0	0	0	0	0	0	0
0	1	2	3	4	5	6	0
0	7	8	9	10	11	12	0
0	1	2	3	4	5	6	0
0	7	8	9	10	11	12	0
0	1	2	3	4	5	6	0
0	7	8	9	10	11	12	0
0	0	0	0	0	0	0	0

11	10	9	10	11	12	11	10
7	6	5	6	7	8	7	6
3	2	1	2	3	4	3	2
7	6	5	6	7	8	7	6
11	10	9	10	11	12	11	10
15	14	13	14	15	16	15	14
11	10	9	10	11	12	11	10
7	6	5	6	7	8	7	6

Zero padding with a one-pixel thick boundary

1	1	1	2	3	4	4	4
1	1	1	2	3	4	4	4
1	1	1	2	3	4	4	4
5	5	5	6	7	8	8	8
1	1	1	2	3	4	4	4
5	5	5	6	7	8	8	8
5	5	5	8	9	8	8	8
5	5	5	8	9	8	8	8



zero padding

Mirror padding with a two-pixel thick boundary



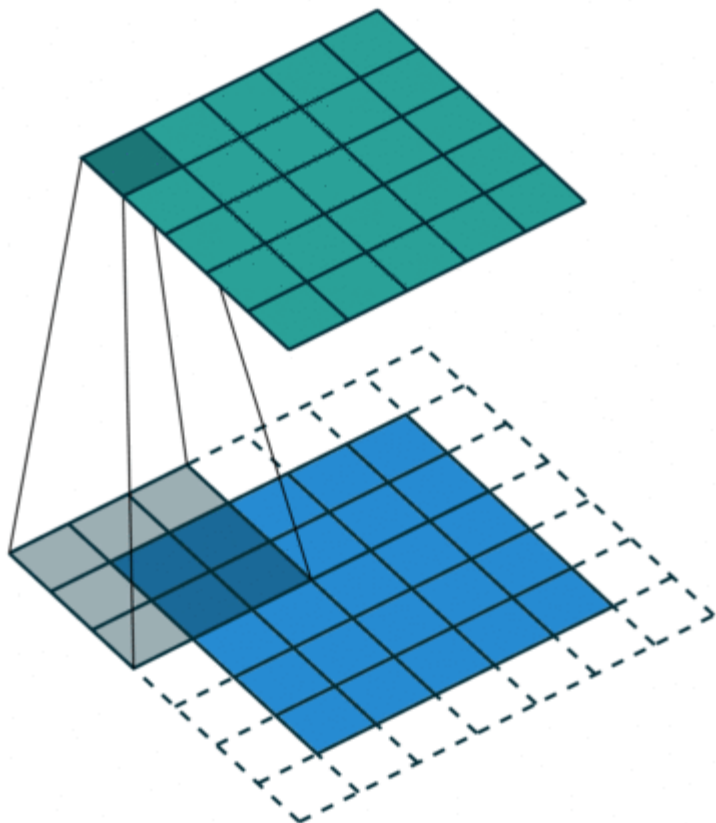
mirror padding



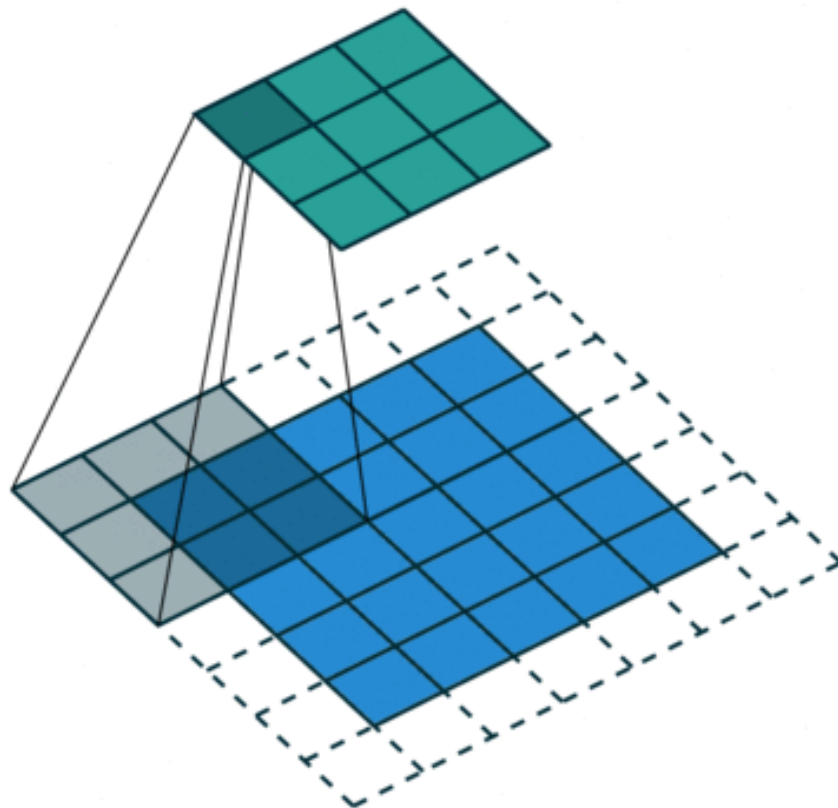
circular padding

Replicate padding with a two-pixel thick boundary

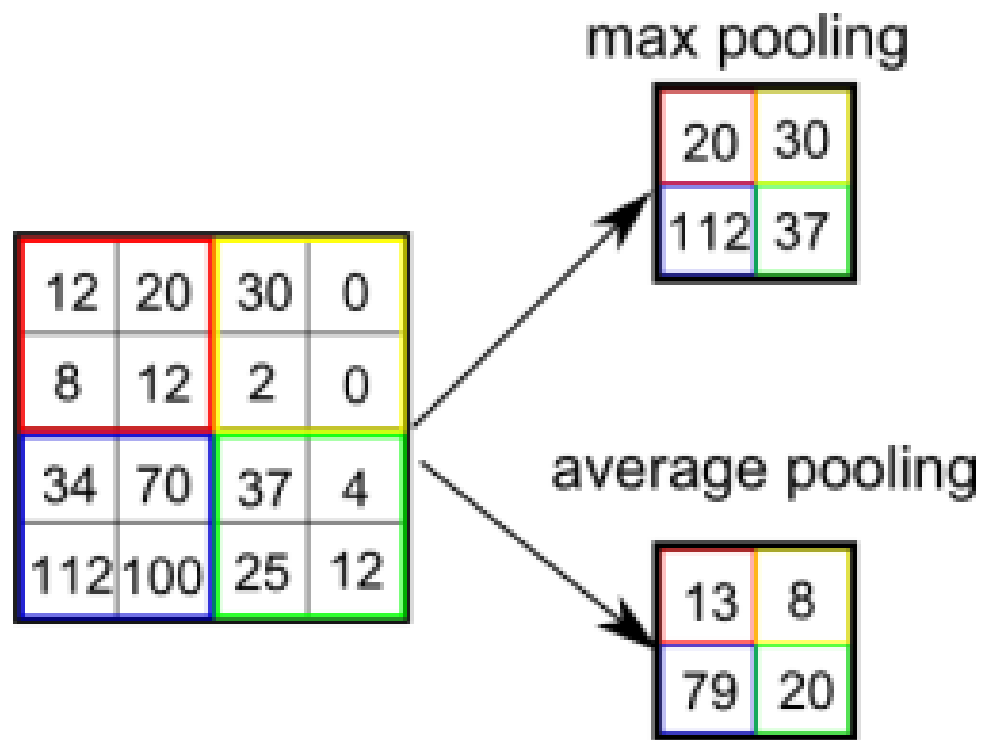
Padding = 1  
Stride = 1



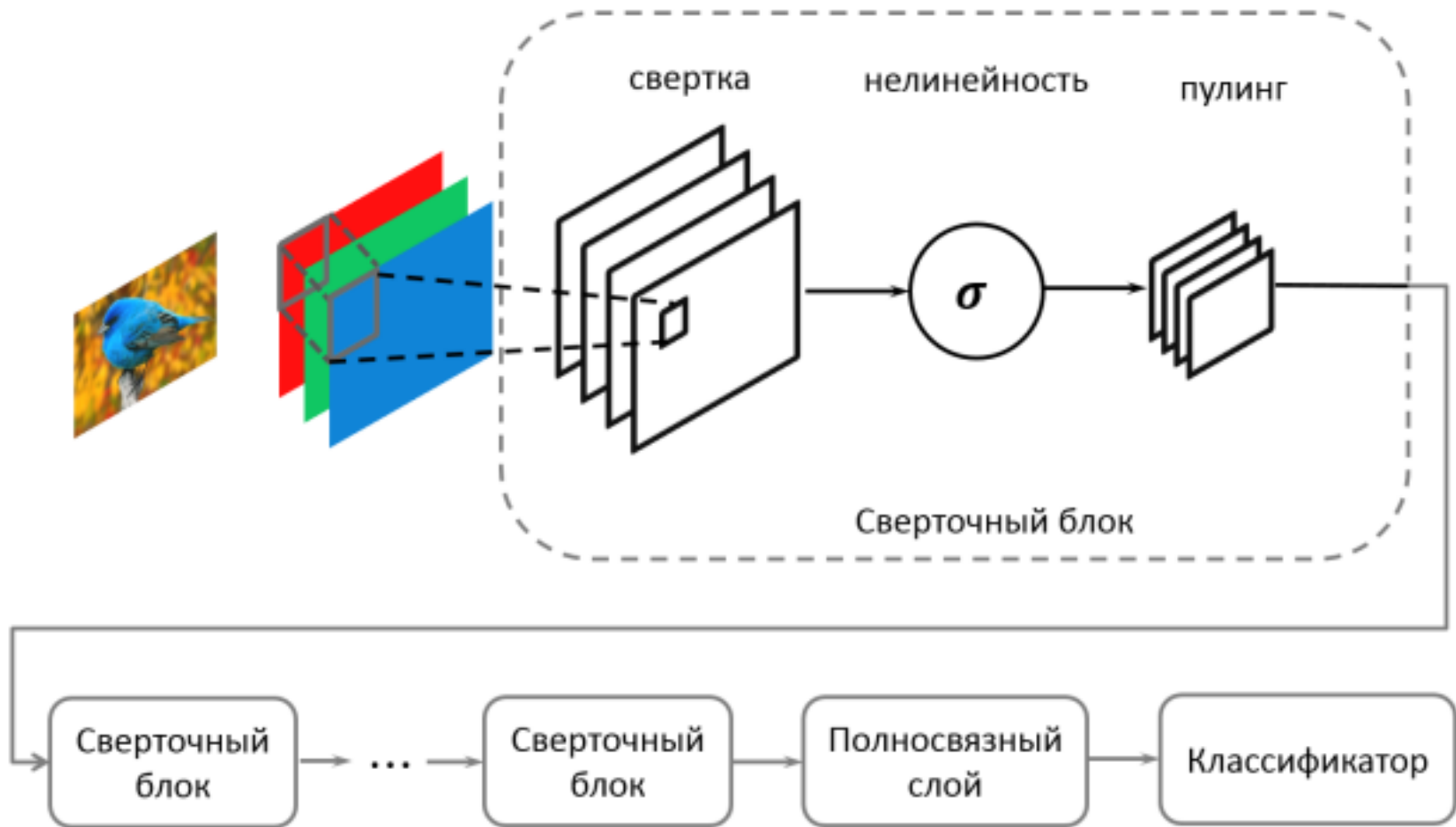
Padding = 1  
Stride = 2



- Pooling – техника уменьшения размерности (downsampling) карт активации. Если некая область содержит ярко выраженные свойства, то мы можем отказаться от поиска других свойств в этой области

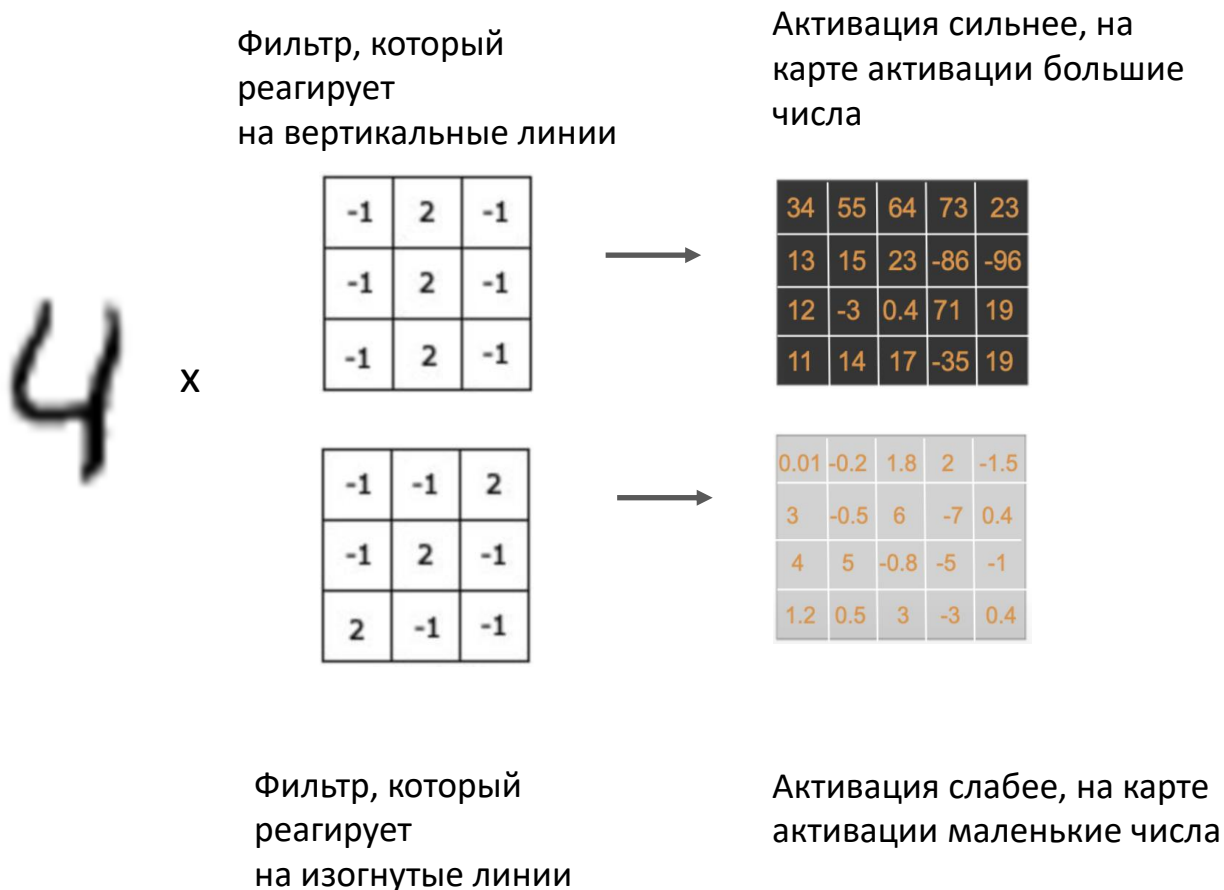


# Базовая схема сверточной сети

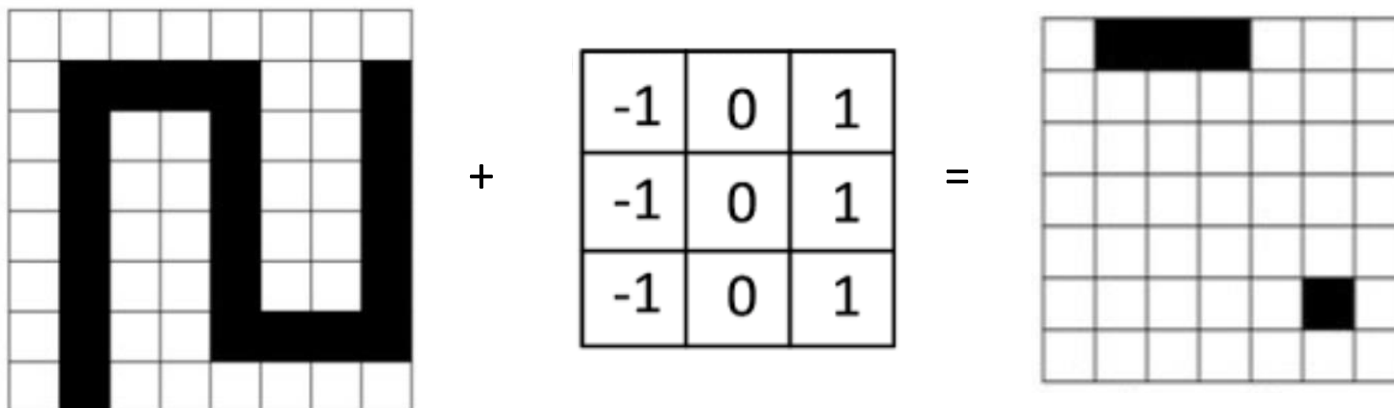




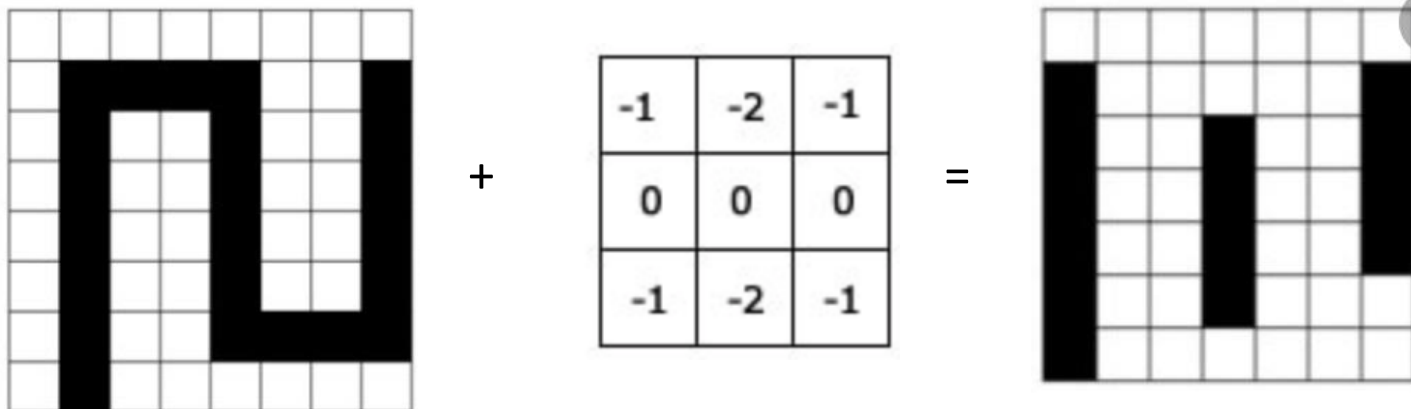
Фильтры “реагируют” на паттерны на изображении. Если паттерн присутствует на изображении, то карта активации после соотв. фильтра будет содержать большие числа



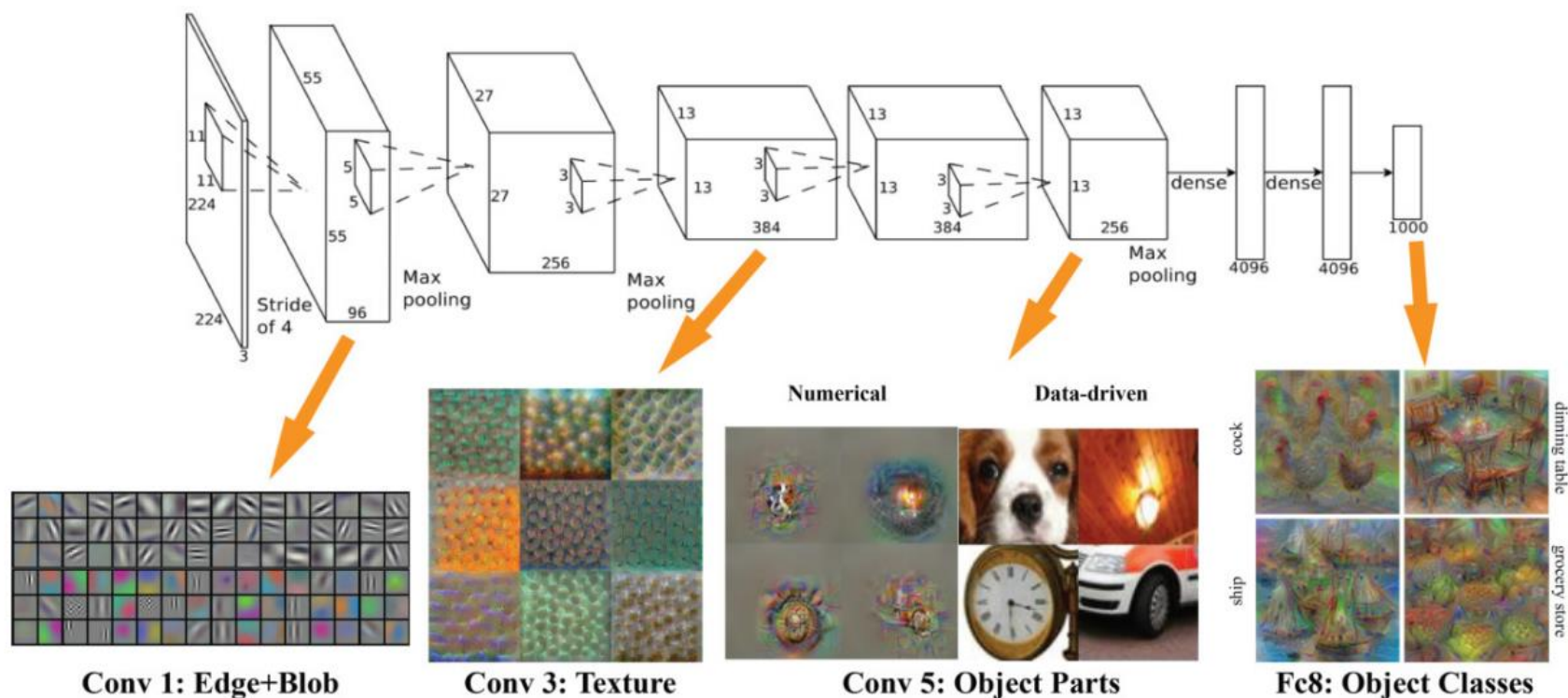
Как мог бы выглядеть фильтр, реагирующий на горизонтальные линии



Как мог бы выглядеть фильтр, реагирующий на вертикальные линии



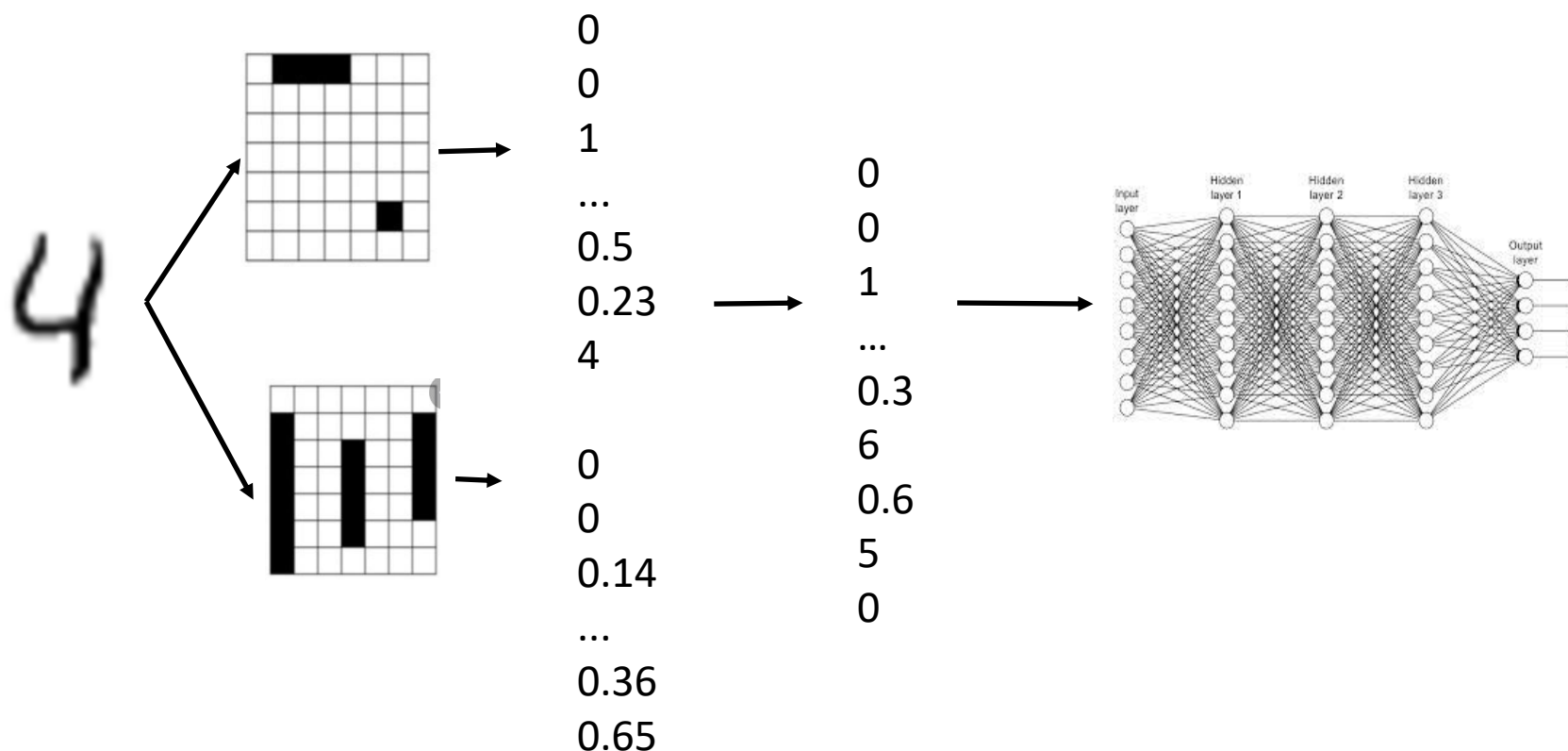
- Сверточная сеть обучается извлечению признаков. Чем выше слой, тем более крупные и сложные элементы изображений он способен распознавать



Krizhevsky A., Sutskever I., Hinton G. ImageNet classification with deep convolutional neural networks. 2012.

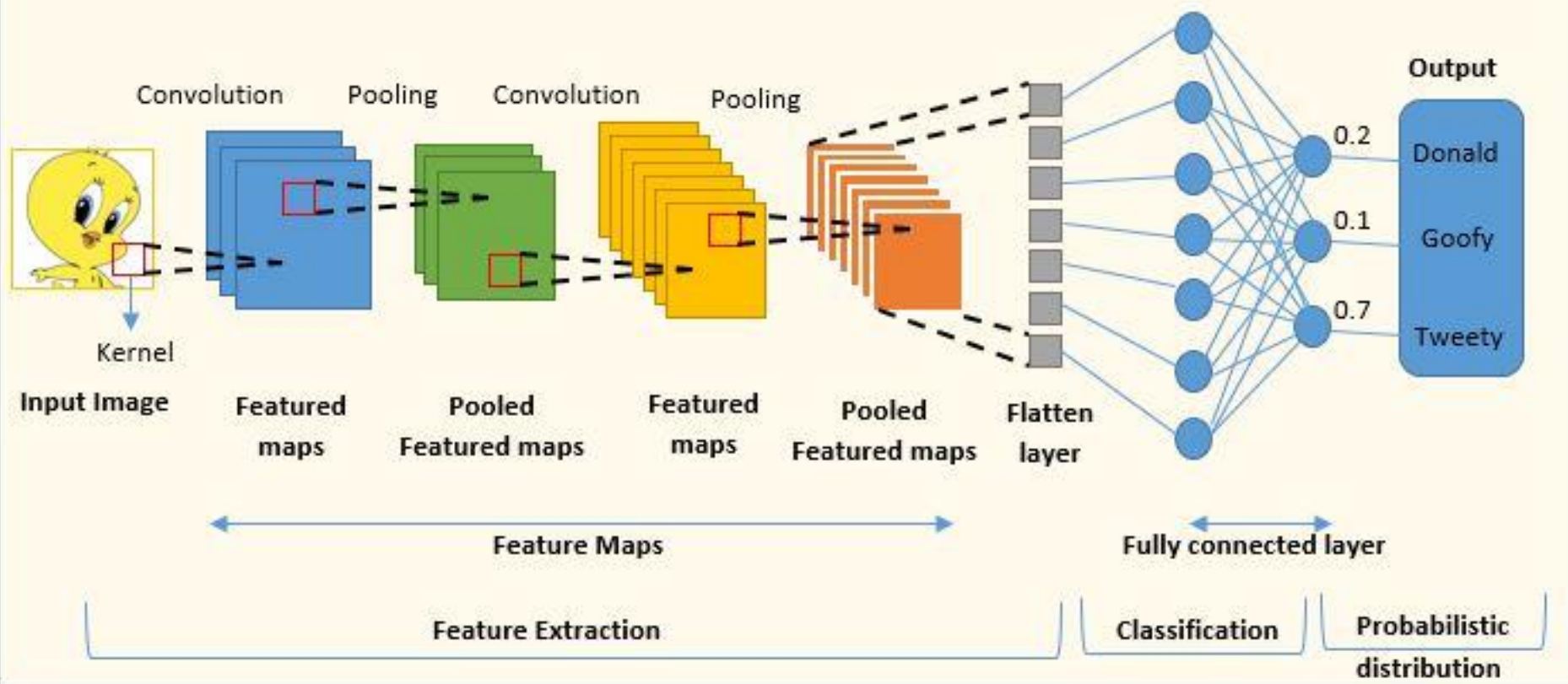


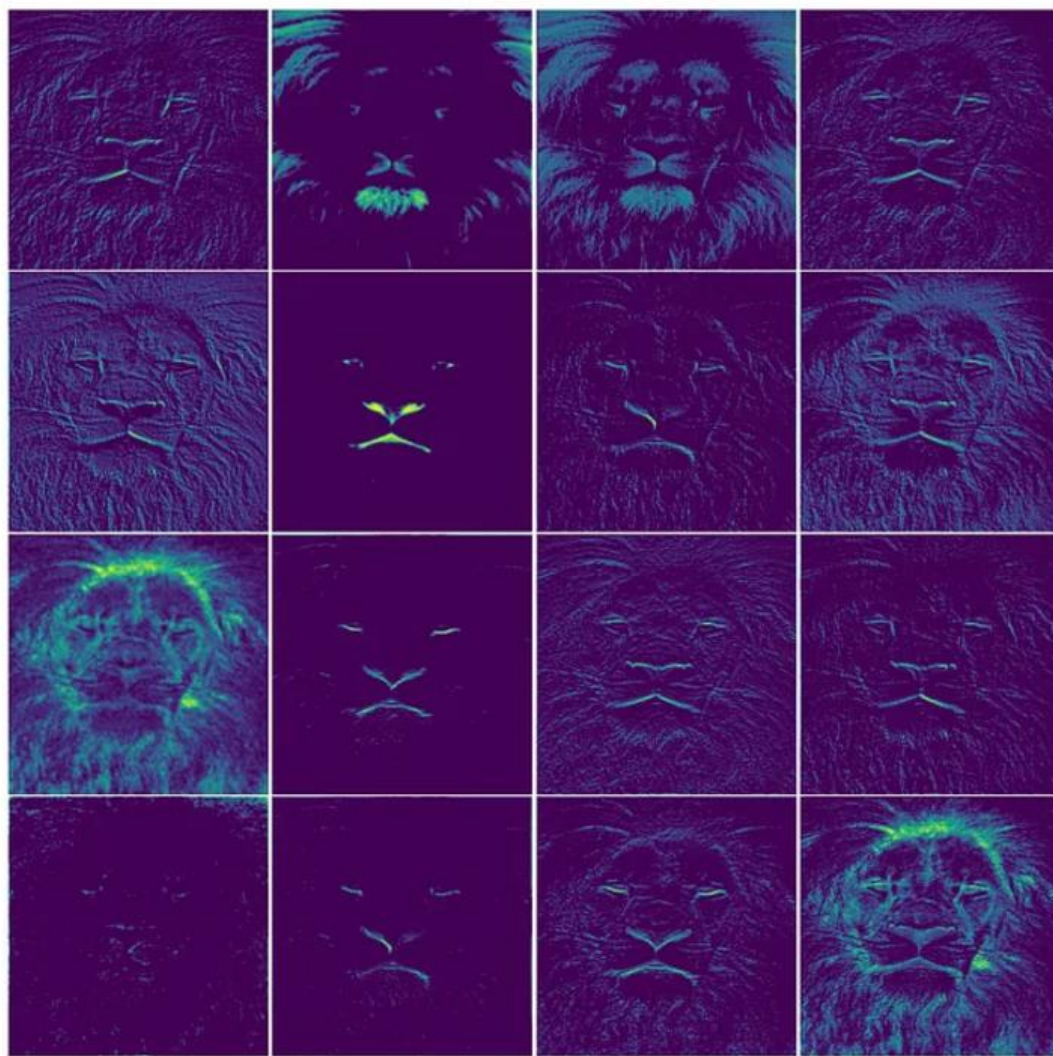
После получения карт активаций, мы **развернем все карты в векторы, сконкатенируем** и подадим на вход полносвязной сети



Вместо ручного конструирования признаков изображения с помощью операции свертки НС сама обучается извлечению признаков. Но для комплексных изображений необходимо много больше промежуточных представлений

### *A Typical Convolutional Neural Network (CNN)*





**(Left)** Feature extraction performed over the image of a lion using vgg19 CNN architecture (image by author). **(Right)** Original picture of the lion (public domain, available at [Pexels](https://www.pexels.com/photo/young-lion-cub-lying-down/)).



## Свертка:

- Выявляет наличие на изображении паттерна, который задается ядром свертки
- Результат операции свертки – промежуточное представление (новое изображение)
- Чем сильнее на исходном изображении представлен паттерн, тем выше будут значения свертки
- Хотим много различных паттернов – используем много сверток

## А что вообще обучать?

- Ядро свертки – подбирается в процессе обучения
- Веса полносвязного слоя – подбираются в процессе обучения



- Первые свертки выделяют какие-то низкоуровневые паттерны (изгибы, край, линии)
- Последние сверточные слои выделяют части изображений. Выход последнего слоя - признаковое описание изображения
- Полученные представления (признаковые описания) подаются на вход полносвязной сети. Например, для задачи многоклассовой классификации: размер выходного слоя - количество классов, функция активации - `softmax`

- Зная параметры сети (размер ядра свертки, паддинг, пуллинг и страйд) можно посчитать размер выходной матрицы признаков:

$$n_{out} = \left\lfloor \frac{n_{in} + 2p - k}{s} \right\rfloor + 1$$

$n_{in}$ : number of input features

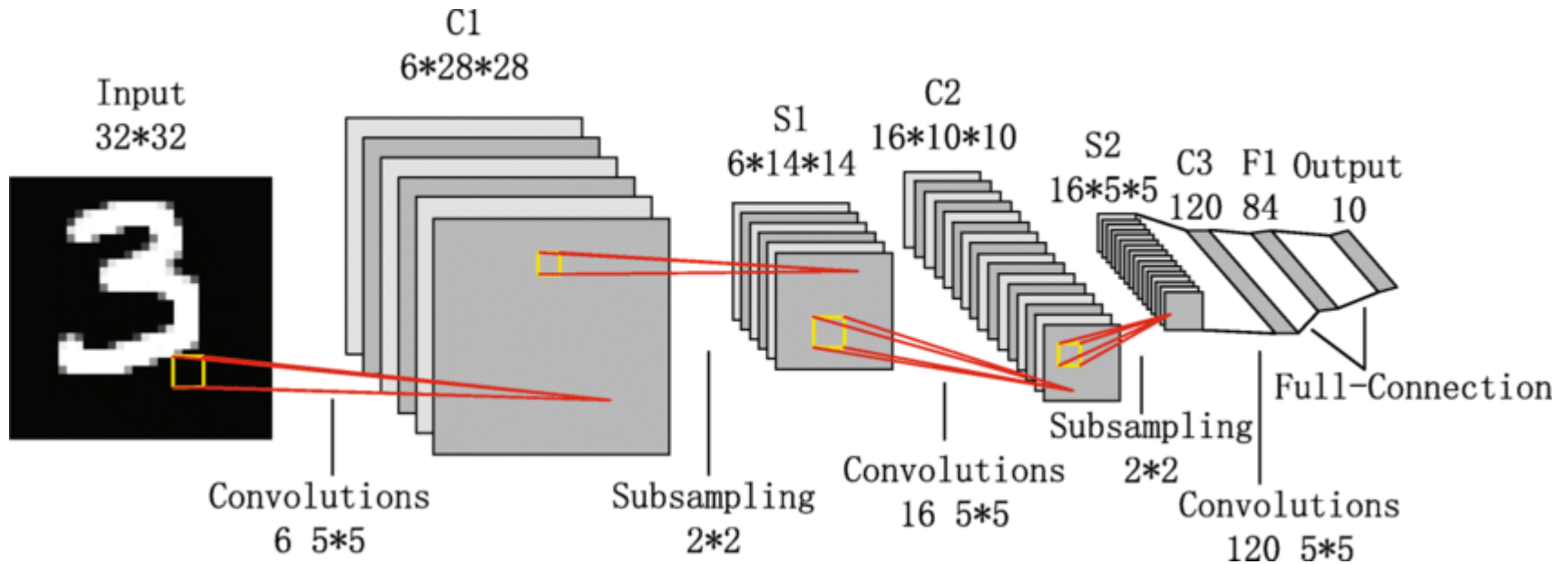
$n_{out}$ : number of output features

$k$ : convolution kernel size

$p$ : convolution padding size

$s$ : convolution stride size

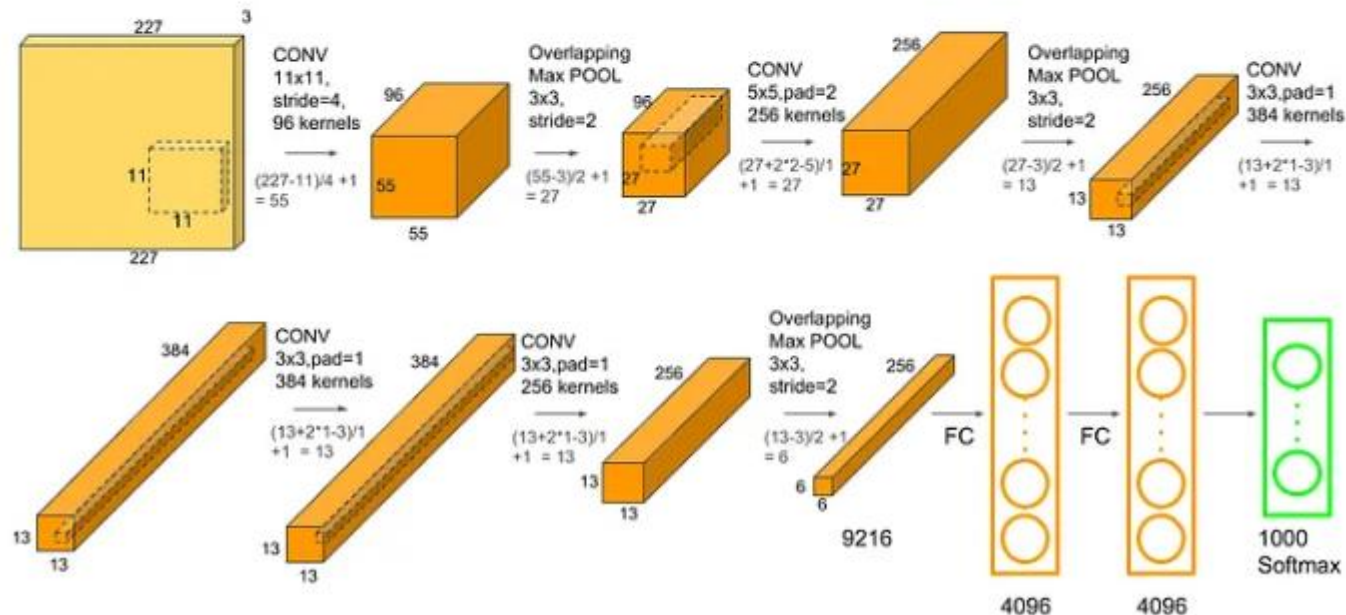
# Насколько глубокие должны быть сети?



Архитектура нейронной сети LeNet-5, 1998 год.  
Первая версия LeNet-1 была представлена в 1989  
году

# Насколько глубокие должны быть сети?

- AlexNet — 2012. Uses max pooling instead of average pooling. Uses ReLU activation instead of tanh. Takes in 3 channels (RGB) as input. Obtained significantly better results in ImageNet Large Scale Visual Recognition Challenge compared to other models at the time.



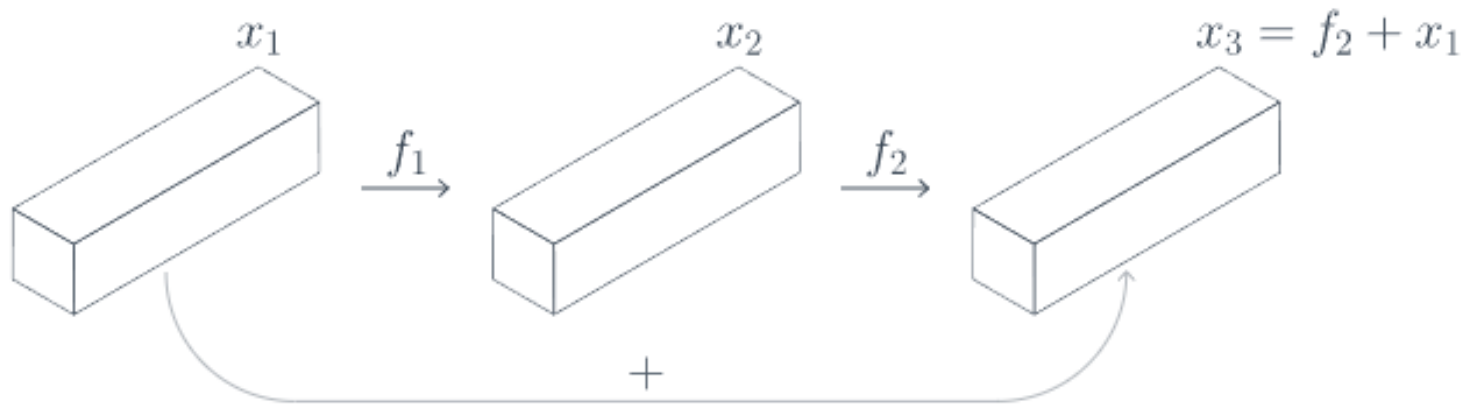
Layer	Type	Maps	Size	Kernel size	Stride	Padding	Activation
out	Fully connected	__	1000	__	__	__	Softmax
F9	Fully connected	__	4096	__	__	__	ReLU
F8	Fully connected	__	4096	__	__	__	ReLU
C7	Convolution	256	13*13	3*3	1	Same	ReLU
C6	Convolution	384	13*13	3*3	1	Same	ReLU
C5	Convolution	384	13*13	3*3	1	Same	ReLU
S4	Max Pooling	256	13*13	3*3	2	Valid	__
C3	Convolution	256	27*27	5*5	1	Same	ReLU
S2	Max Pooling	96	27*27	3*3	2	Valid	__
C1	Convolution	96	55*55	11*11	4	Valid	ReLU
In	Input	3(RGB)	227*227	__	__	__	__

*AlexNet Summary*

- Чем больше сверток, тем лучше?

Если мы будем бесконтрольно добавлять сверточные слои, то, несмотря на использование ReLU и batchNorm, градиенты все равно будут затухать и на первых слоях будут близки к нулю

- Residual connection (архитектура ResNet, Residual NN) – прокидываем признаки на предыдущем слое мимо сверток на следующем



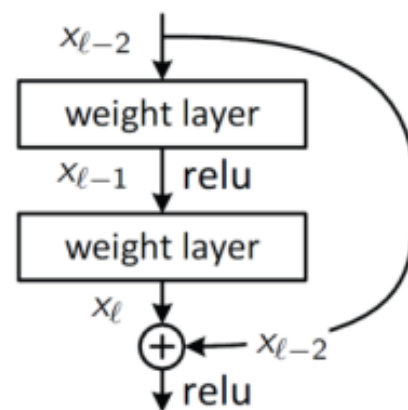
- DenseNet – вместо суммы конкатенация, результаты лучше, но вычислительно сложно



Сквозная связь (skip connection) слоя  $\ell$  с предшествующим слоем  $\ell - d$ :

$$x_\ell = \sigma(Wx_{\ell-1}) + x_{\ell-d}$$

Слой  $\ell$  выучивает не новое векторное представление  $x_\ell$ , а его приращение  $x_\ell - x_{\ell-d}$



- Приращения более устойчивы  $\Rightarrow$  улучшается сходимость
- Появляется возможность увеличивать число слоёв
- Обобщение — Highway Networks:

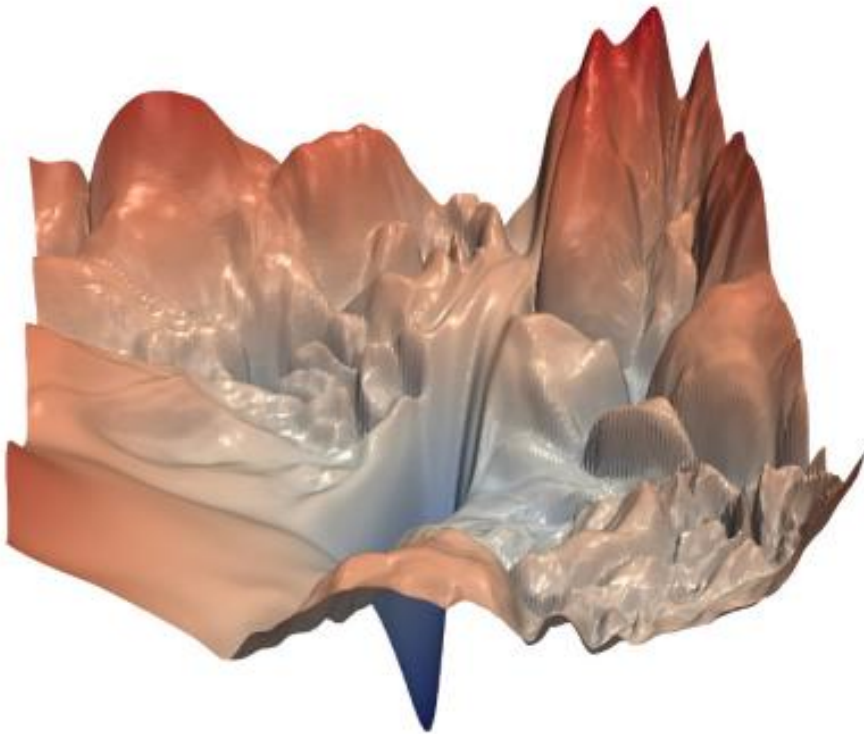
$$x_\ell = \sigma(Wx_{\ell-1}) \underbrace{\tau(W'x_{\ell-1})}_{\text{transform gate}} + x_{\ell-d} \underbrace{(1 - \tau(W'x_{\ell-1}))}_{\text{carry gate}}$$

---

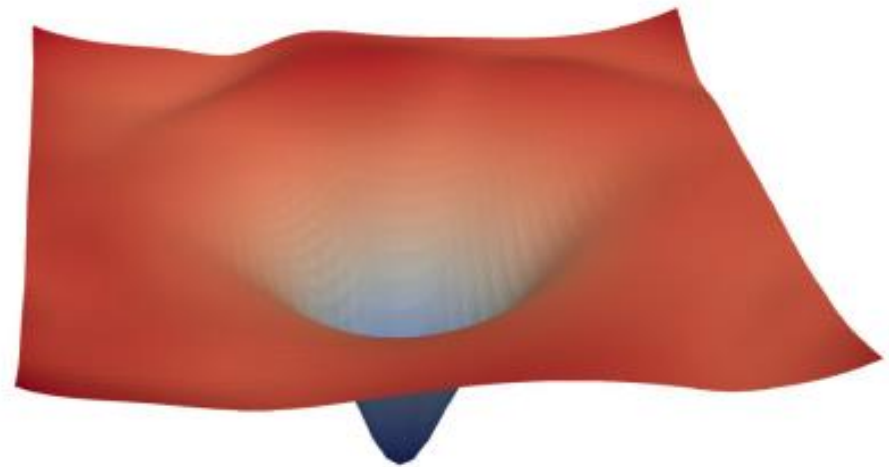
Kaiming He, Xiangyu Zhang, Shaoqing Ren, Jian Sun. Deep Residual Learning for Image Recognition. 2015

R.K.Srivastava, K.Greff, J.Schmidhuber. Highway Networks. 2015

Сквозные связи (skip connection) упрощают оптимизируемый критерий, устраняя локальные экстремумы и седловые точки:



without skip connections



with skip connections

---

*Hao Li et al. Visualizing the Loss Landscape of Neural Nets. 2018*

# Знаковые архитектуры в мире сверточных сетей для задач классификации изображений:

1. [LeNet - 1998](#)
2. [AlexNet - 2012](#)
3. [Network in network - 2013](#)
4. [VGG - 2014](#)
5. [GoogLeNet \(inception\) - 2014](#)
6. [ResNet - 2015](#)
7. [MobileNet - 2017](#)
8. [EfficientNet - 2019](#)

## Что учитываем при обучении CNN:

- Функции активации без горизонтальных асимптот (ReLU)
- Адаптивные градиентные методы
- Dropout
- Batch normalization
- Residual NN
- Подбираем число слоев и размеры слоев
- Используем аугментацию

# Практическое применение CNN

**Semantic Segmentation**



GRASS, CAT,  
TREE, SKY

No objects, just pixels

**Classification  
+ Localization**



CAT

Single Object

**Object  
Detection**



DOG, DOG, CAT

Multiple Object

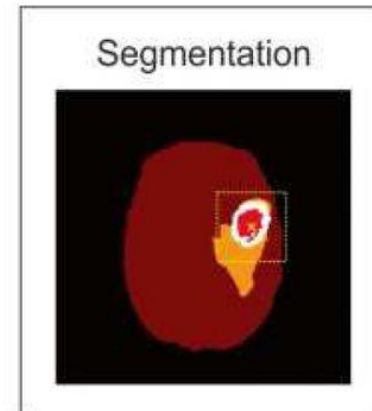
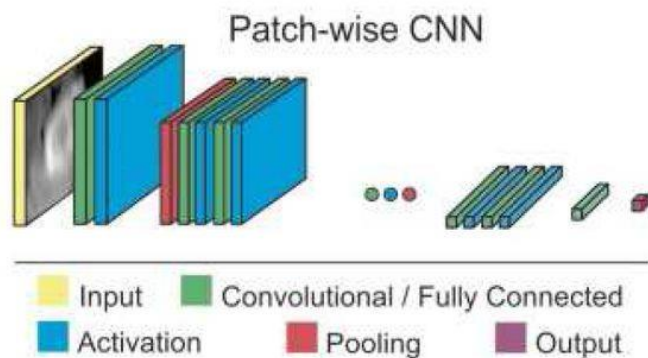
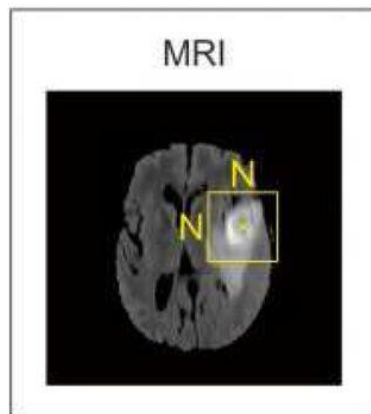
**Instance  
Segmentation**



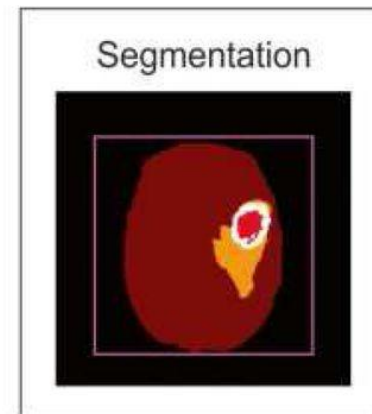
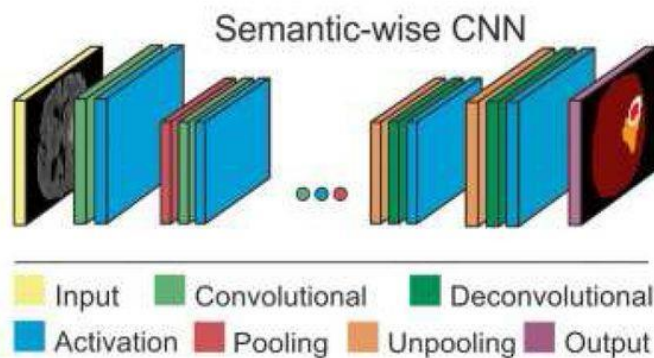
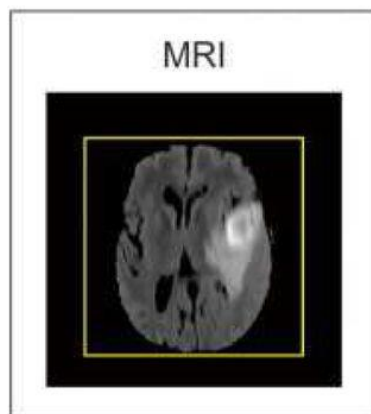
DOG, DOG, CAT

This image is CC0 public domain

# Применение в медицине



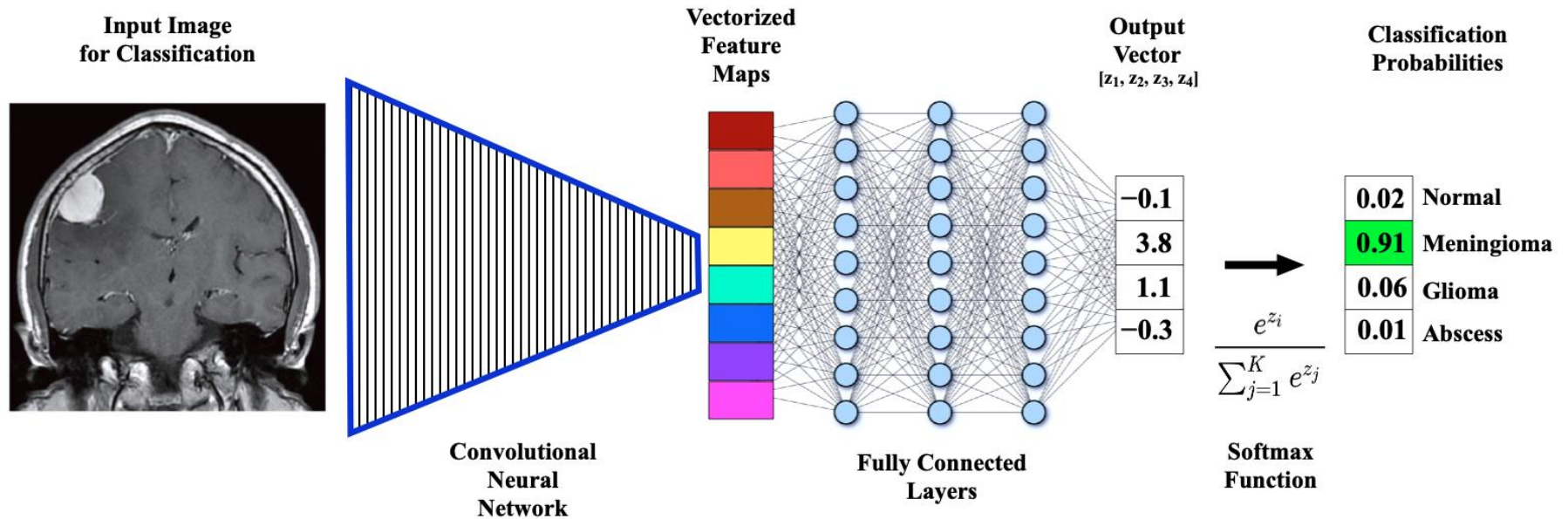
(a) Voxel-wise CNN architecture



(b) Semantic-wise CNN architecture



# Применение в медицине



# Применение в беспилотном транспорте

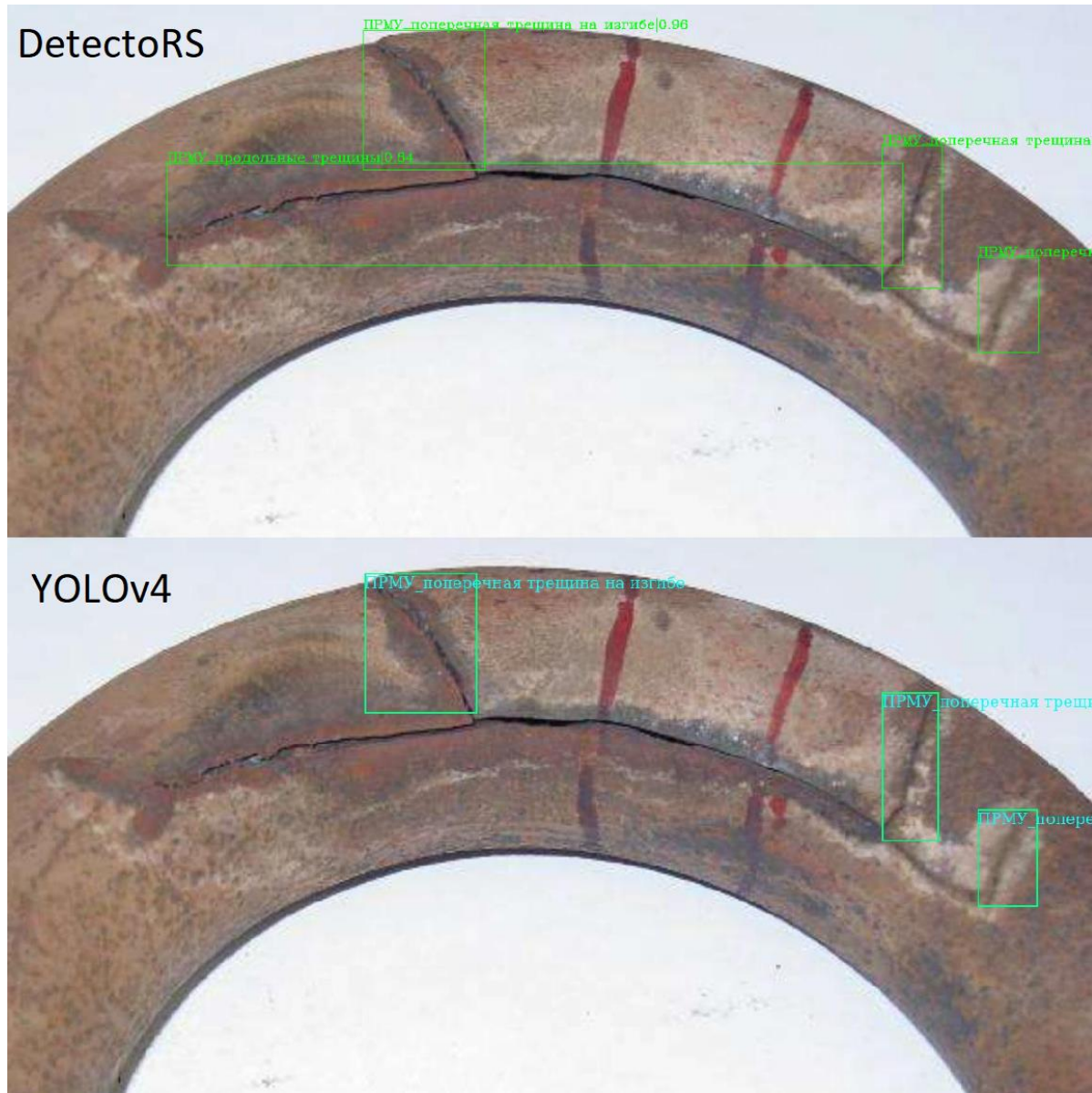


# Применение в беспилотном транспорте





# Применение в промышленности (поиск дефектов)



# CNN explainer — наглядное средство визуализации CNN

