

Local descriptors

Return of DL

YOUNG & YANDEX

Радослав Нейчев
Выпускник и преподаватель ШАД и МФТИ,
руководитель группы ML-разработки в Яндексе,
основатель  girafe

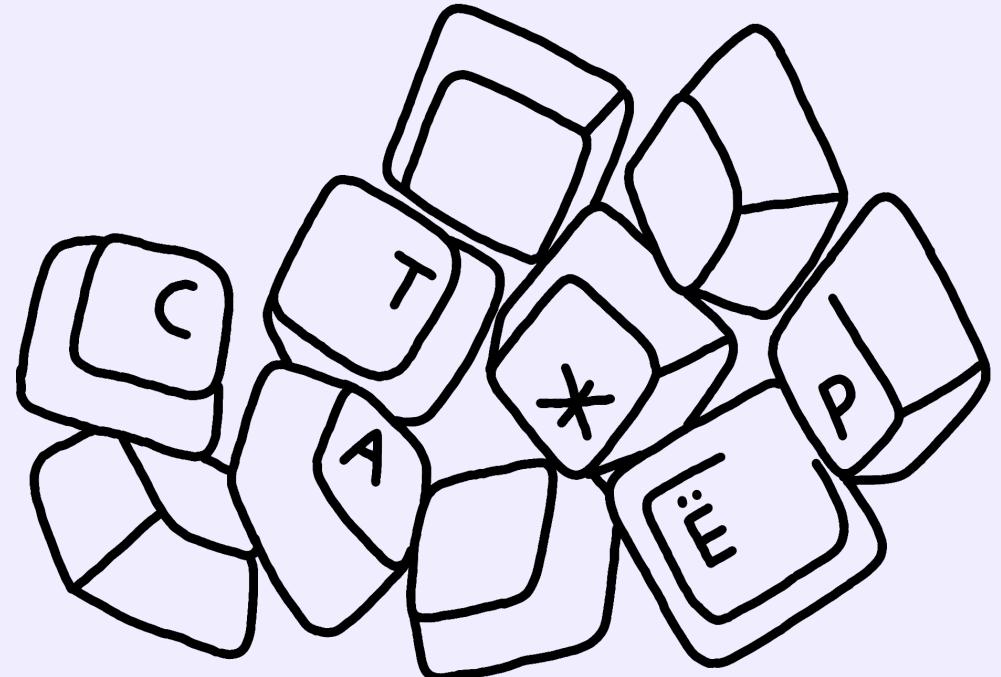


Outline

- 01** Local image descriptors
- 02** Backpropagation recap
- 03** Classification architectures overview

Local image descriptors

01



How to match images?

image 1



image 2



Source: [CSC420 slides by David Lindell, original slides credits: Babak Taati](#)
[←Ahmed Ashraf ← Sanja Fidler](#)

How to match images?

image 1



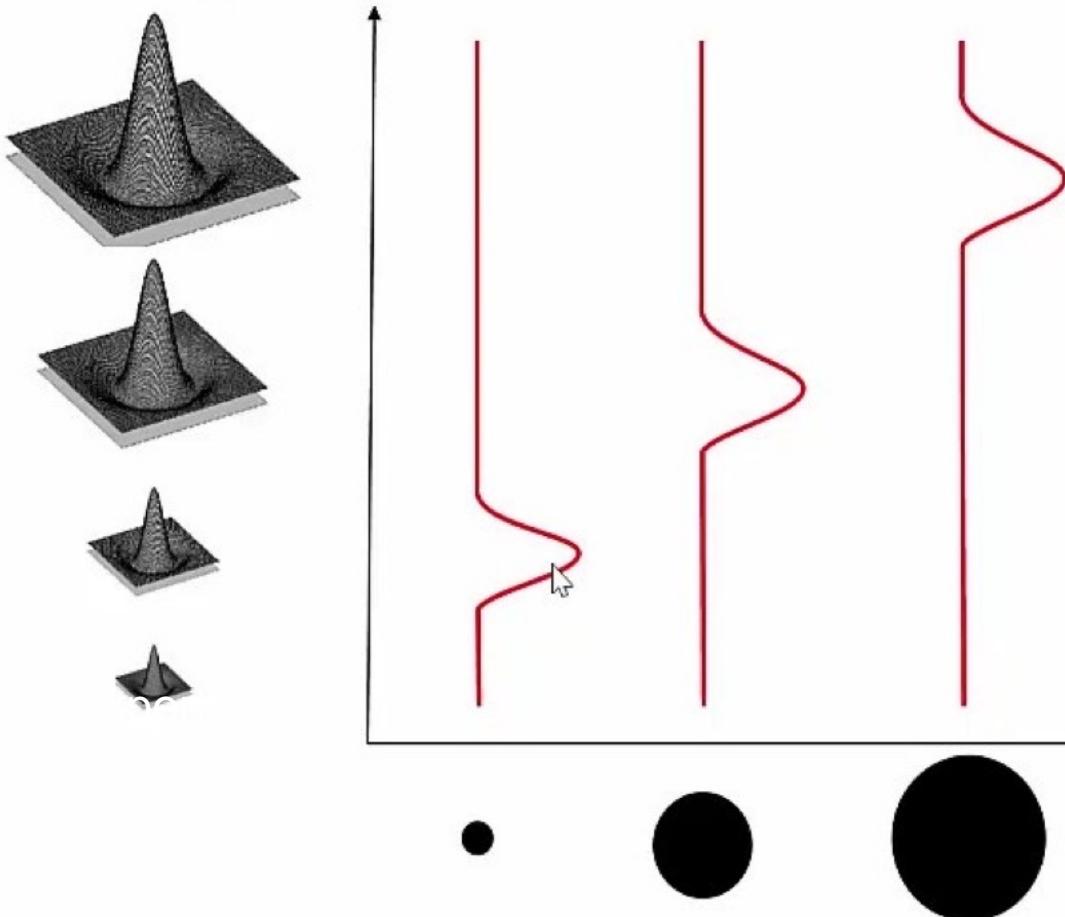
image 2



Source: [CSC420 slides by David Lindell, original slides credits: Babak Taati](#)
[←Ahmed Ashraf ← Sanja Fidler](#)

Blob detection

- Laplacian-of-Gaussian = “blob” detector

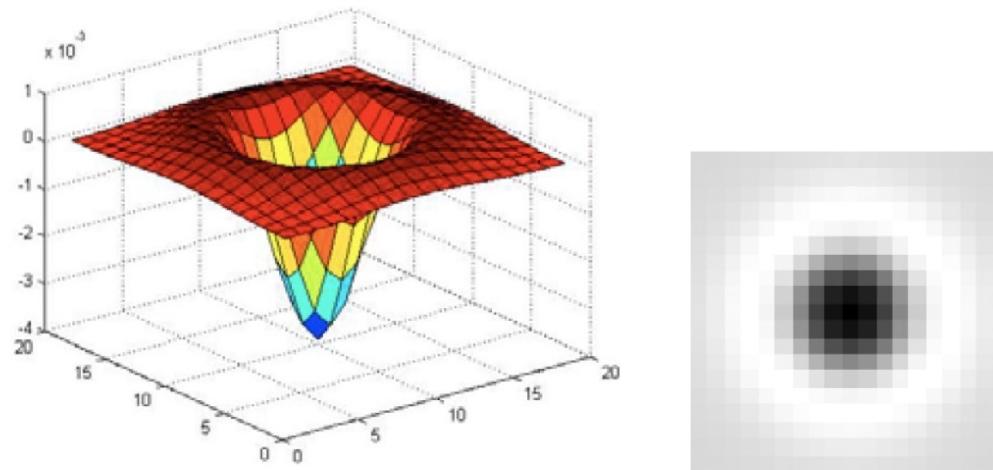


Source: [CSC420 slides by David Lindell, original slides credits: Babak Taati](#)
[←Ahmed Ashraf ← Sanja Fidler](#)

Blob detection

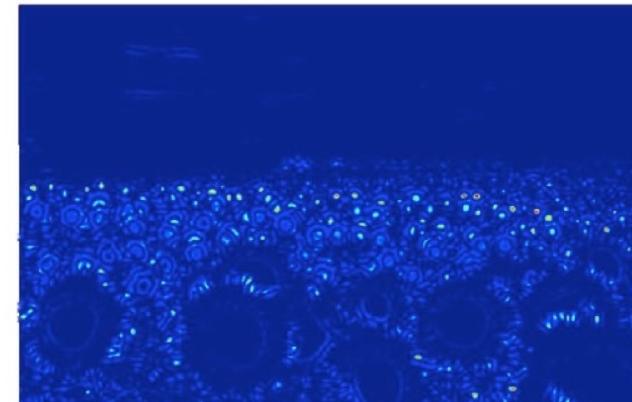
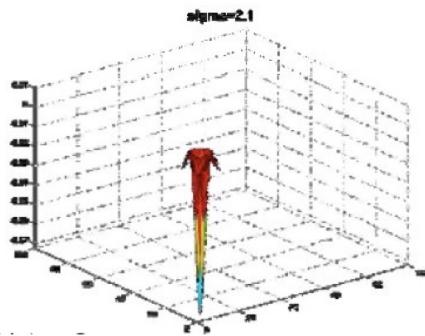
- Laplacian of Gaussian:

$$\begin{aligned} \bullet \nabla_g^2(x, y, \sigma) &= \frac{\partial^2 g(x, y, \sigma)}{\partial x^2} + \frac{\partial^2 g(x, y, \sigma)}{\partial y^2} \text{ where } G \text{ is Gaussian} \\ \bullet \nabla_g^2(x, y, \sigma) &= -\frac{1}{\pi\sigma^4} \left(1 - \frac{x^2+y^2}{2\sigma^2}\right) \exp -\frac{x^2+y^2}{2\sigma^2} \end{aligned}$$



[Source: K. Grauman]

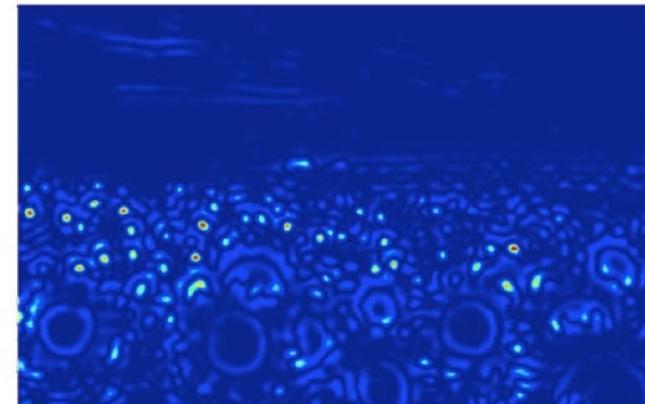
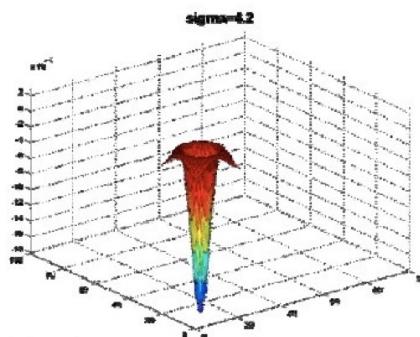
Blob detection



[Source: K. Grauman]

Source: [CSC420 slides by David Lindell, original slides credits: Babak Taati](#)
[←Ahmed Ashraf ← Sanja Fidler](#)

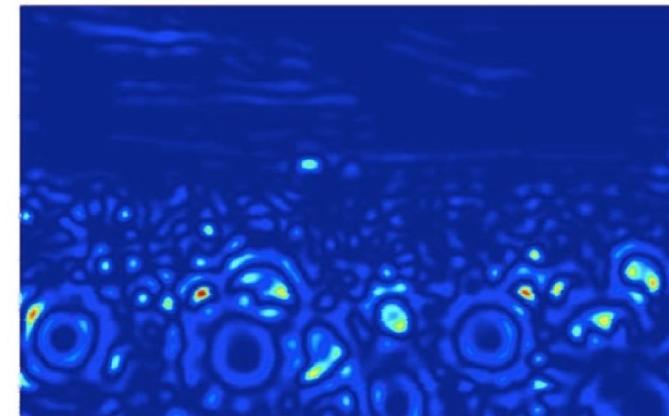
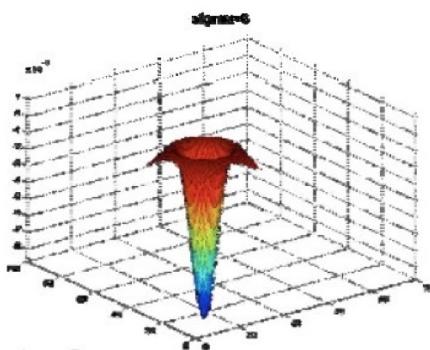
Blob detection



[Source: K. Grauman]

Source: [CSC420 slides by David Lindell, original slides credits: Babak Taati](#)
[←Ahmed Ashraf ← Sanja Fidler](#)

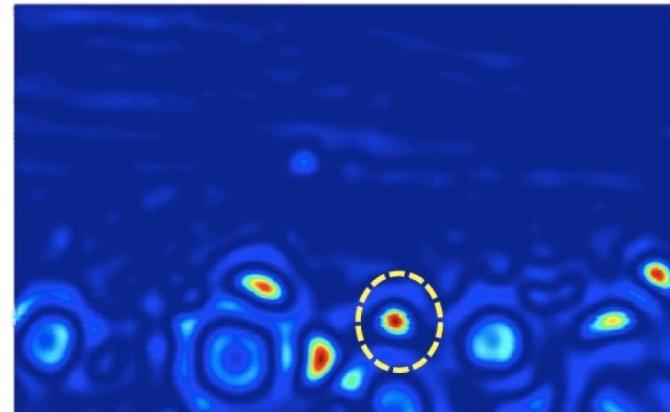
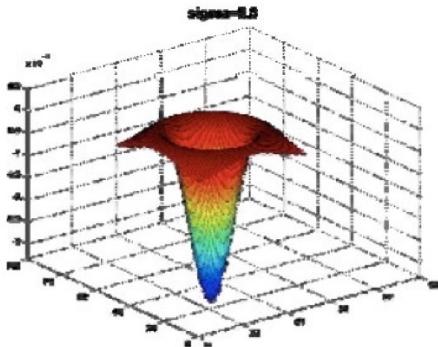
Blob detection



[Source: K. Grauman]

Source: [CSC420 slides by David Lindell, original slides credits: Babak Taati](#)
[←Ahmed Ashraf ← Sanja Fidler](#)

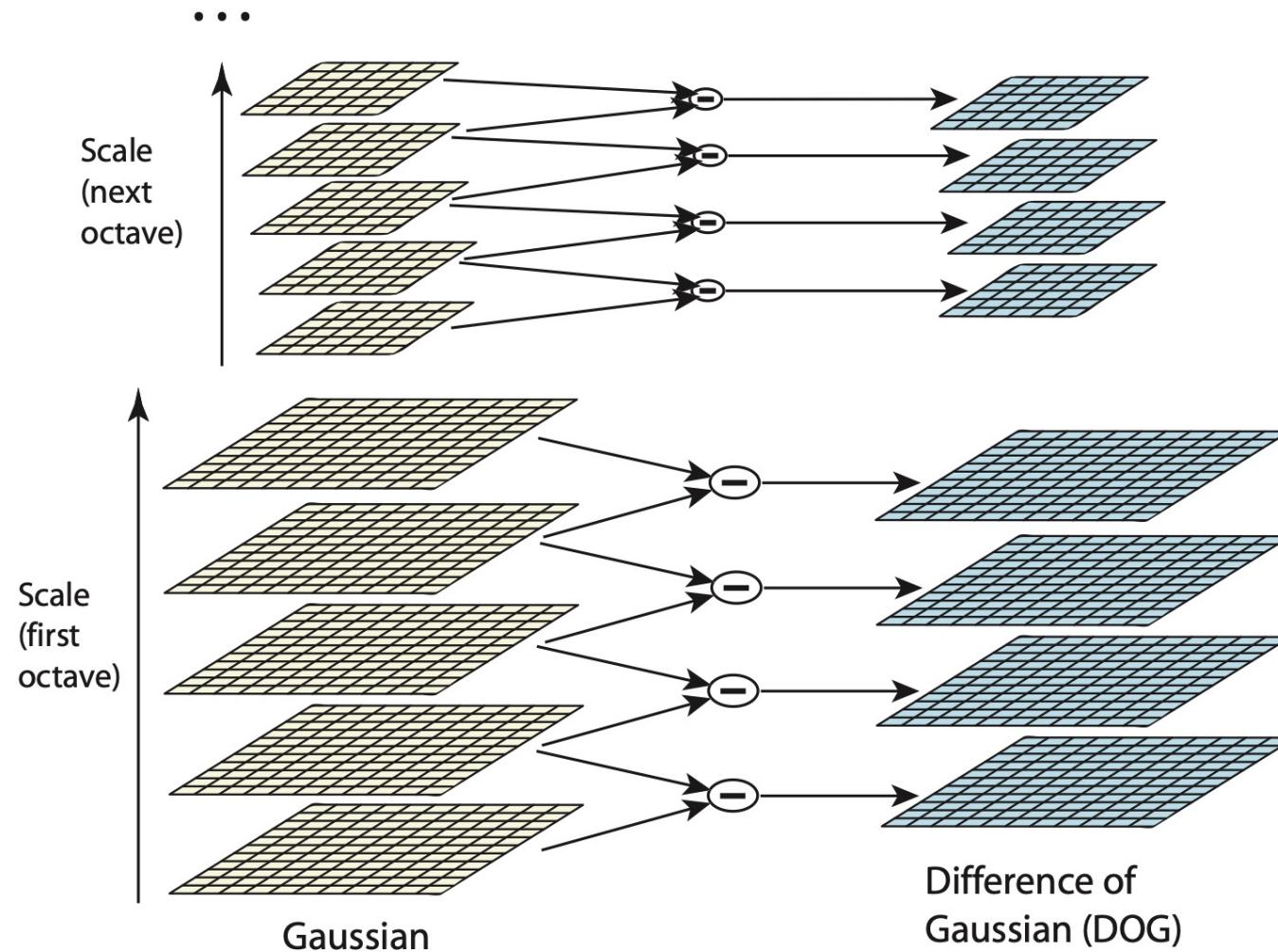
Blob detection



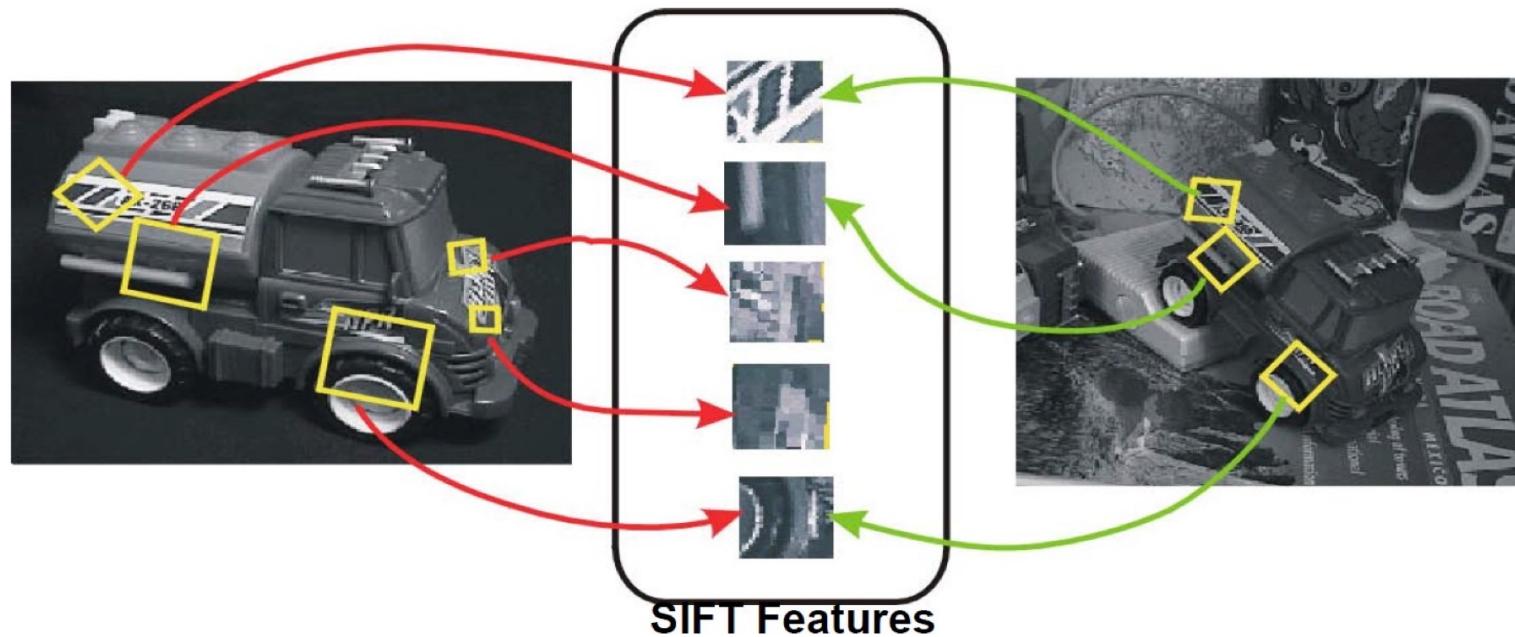
[Source: K. Grauman]

Source: [CSC420 slides by David Lindell, original slides credits: Babak Taati](#)
[←Ahmed Ashraf ← Sanja Fidler](#)

How to get invariant to scale change?



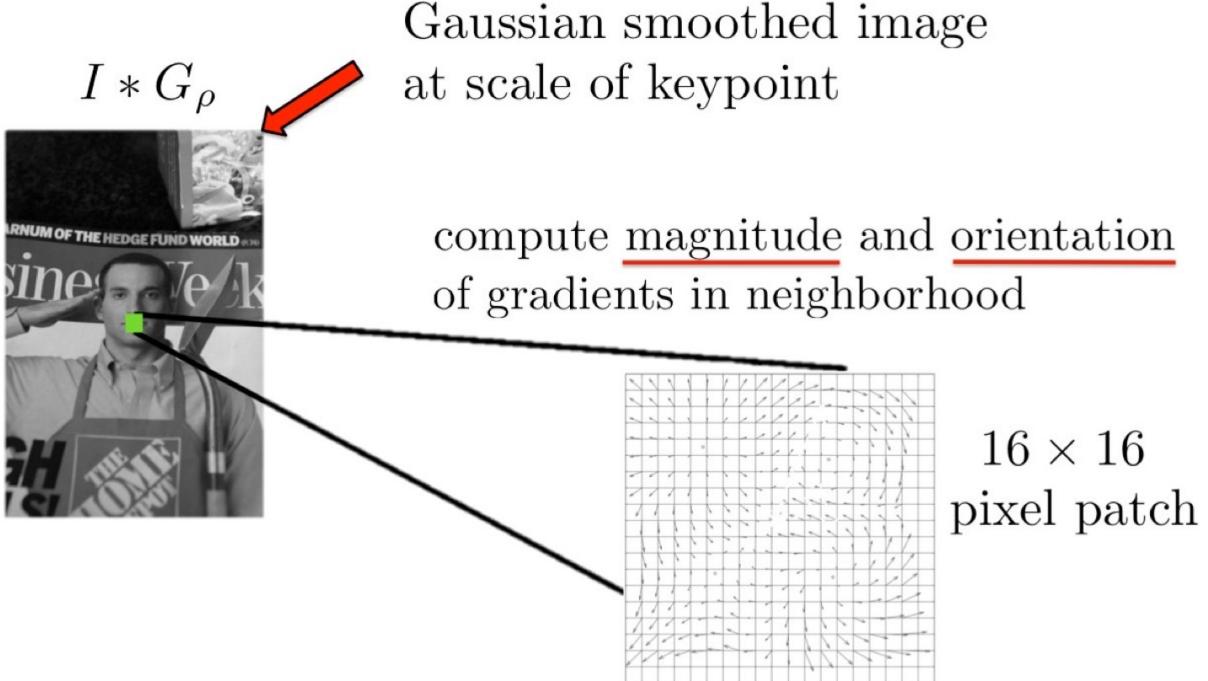
SIFT detector



Source: [CSC420 slides by David Lindell, original slides credits: Babak Taati](#)
[←Ahmed Ashraf ← Sanja Fidler](#)

SIFT detector

- Compute the gradient magnitude and orientation in neighborhood of each keypoint proportional to the detected scale



[Adopted from: F. Flores-Mangas]

SIFT detector

- Compute the gradient magnitude and orientation in neighborhood of each keypoint proportional to the detected scale

magnitude of gradient:

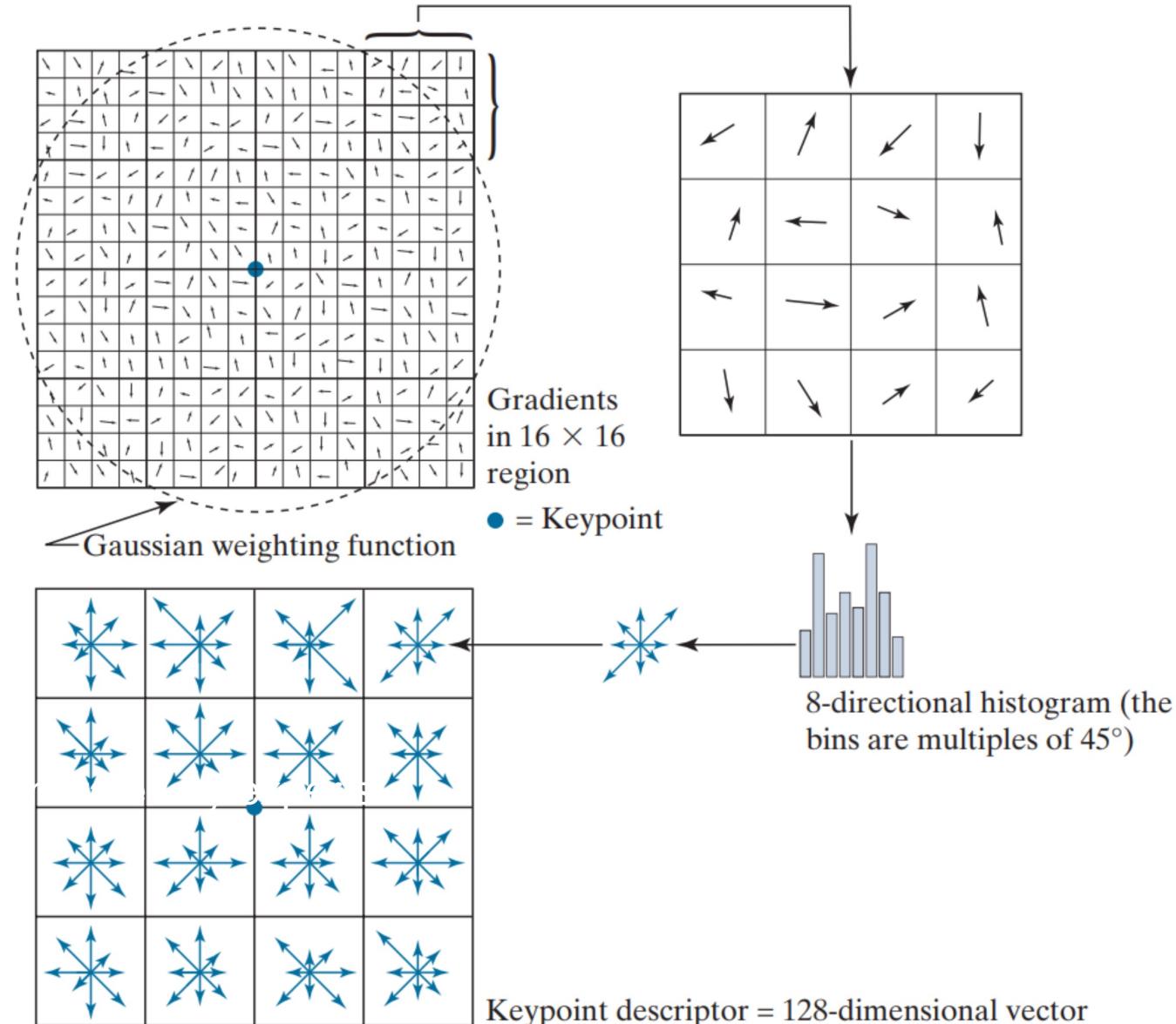
$$|\nabla I(x, y)| = \sqrt{\left(\frac{\partial(I(x, y) * G_\rho)}{\partial x}\right)^2 + \left(\frac{\partial(I(x, y) * G_\rho)}{\partial y}\right)^2}$$

gradient orientation:

$$\theta(x, y) = \arctan\left(\frac{\partial I * G_\rho}{\partial y} / \frac{\partial I * G_\rho}{\partial x}\right)$$

(in case you forgot ;))

[Adopted from: F. Flores-Mangas]



Source: <https://blog.roboflow.com/sift/>

best 50 matches



best 50 matches

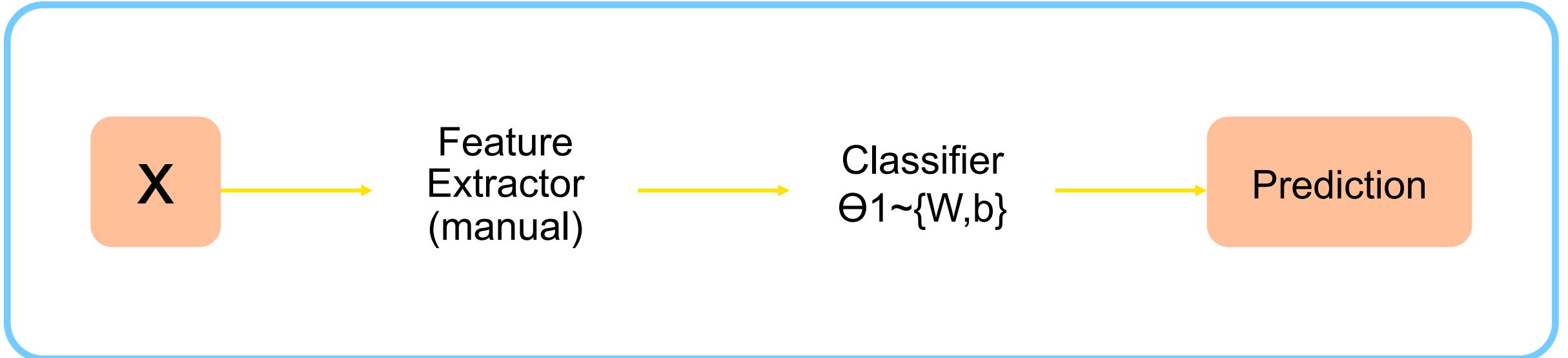


Return to Deep Learning

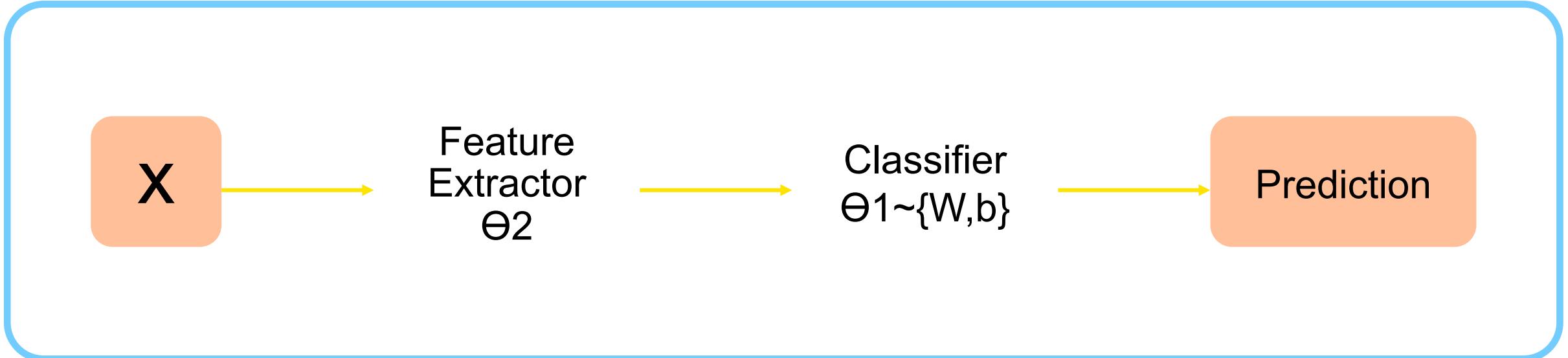
02



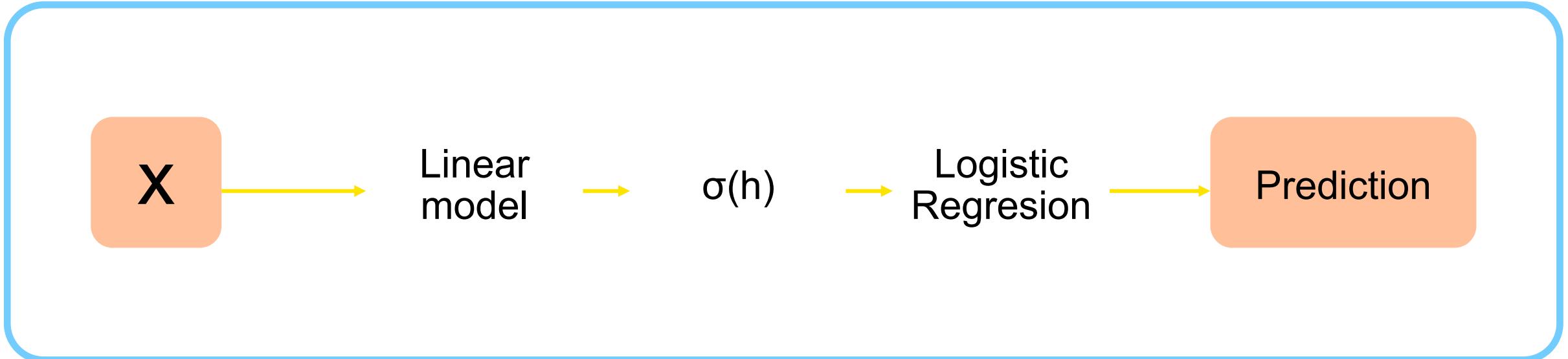
Classic pipeline



NN pipeline



NN pipeline: example



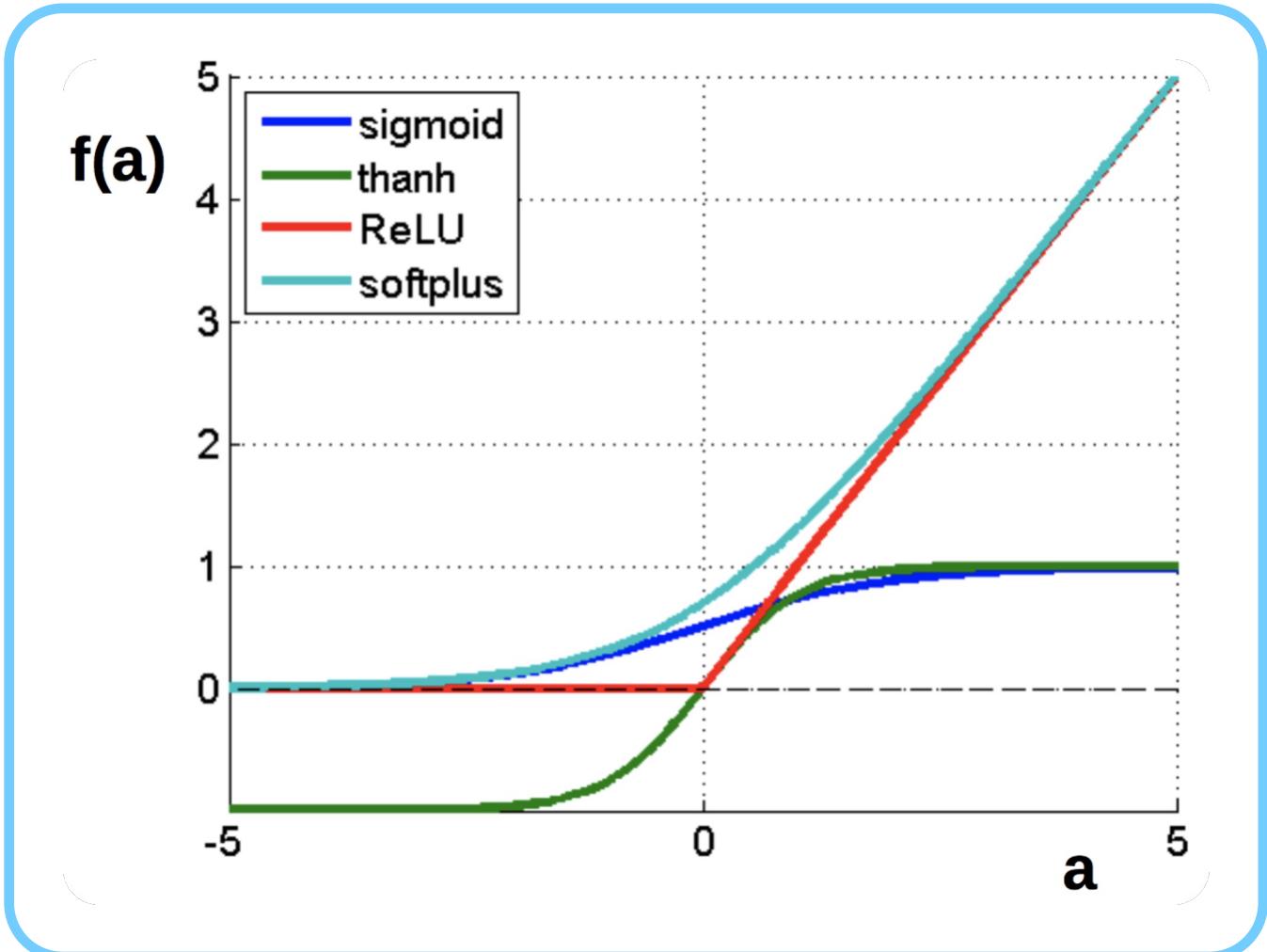
Activation functions: nonlinearities

$$f(a) = \frac{1}{1 + e^{-a}}$$

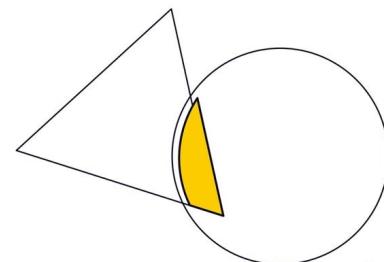
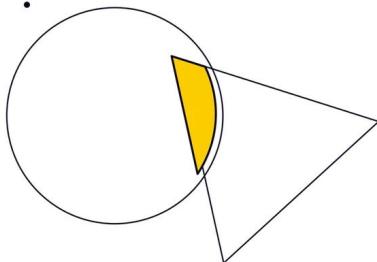
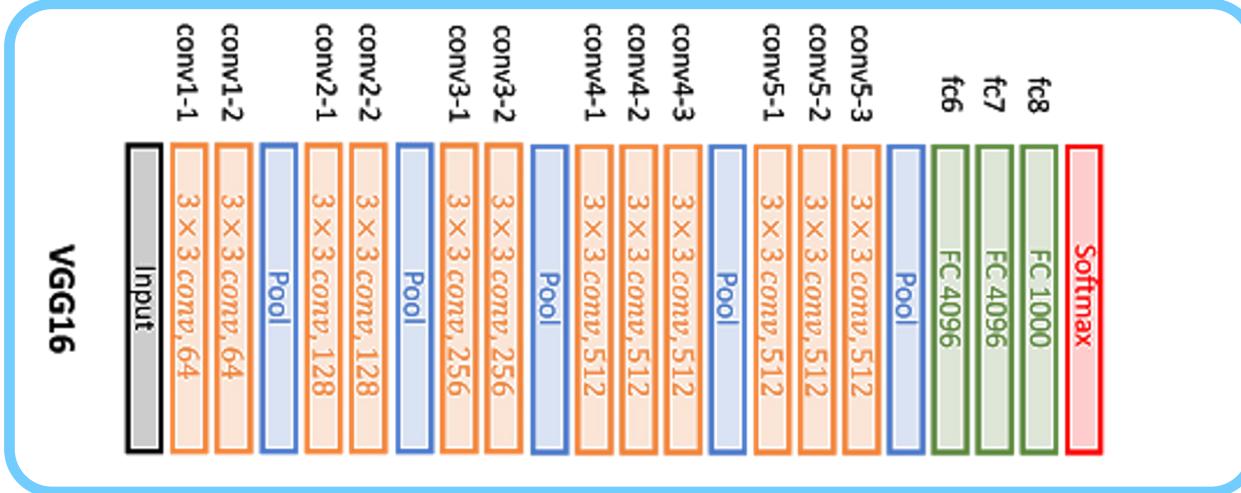
$$f(a) = \tanh(a)$$

$$f(a) = \max(0, a)$$

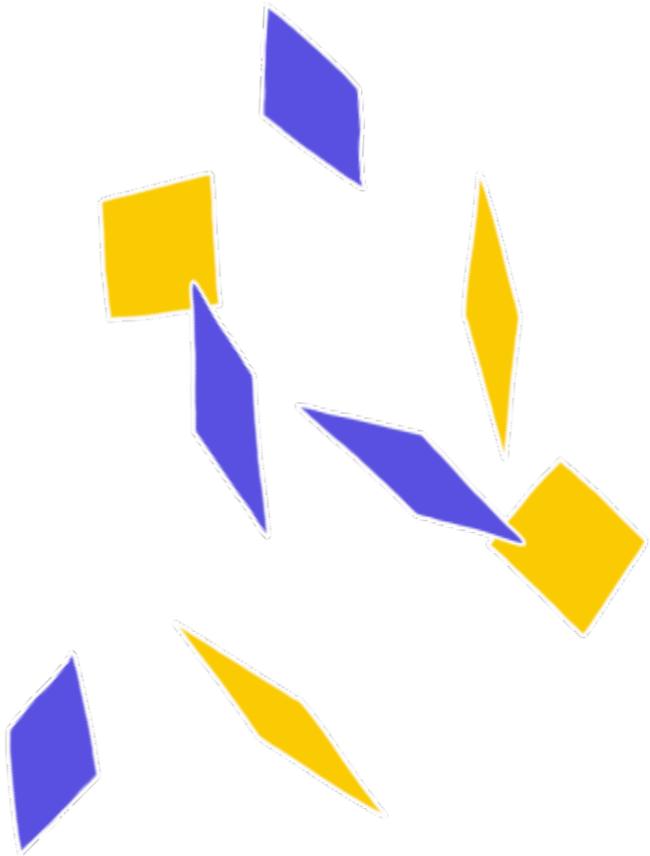
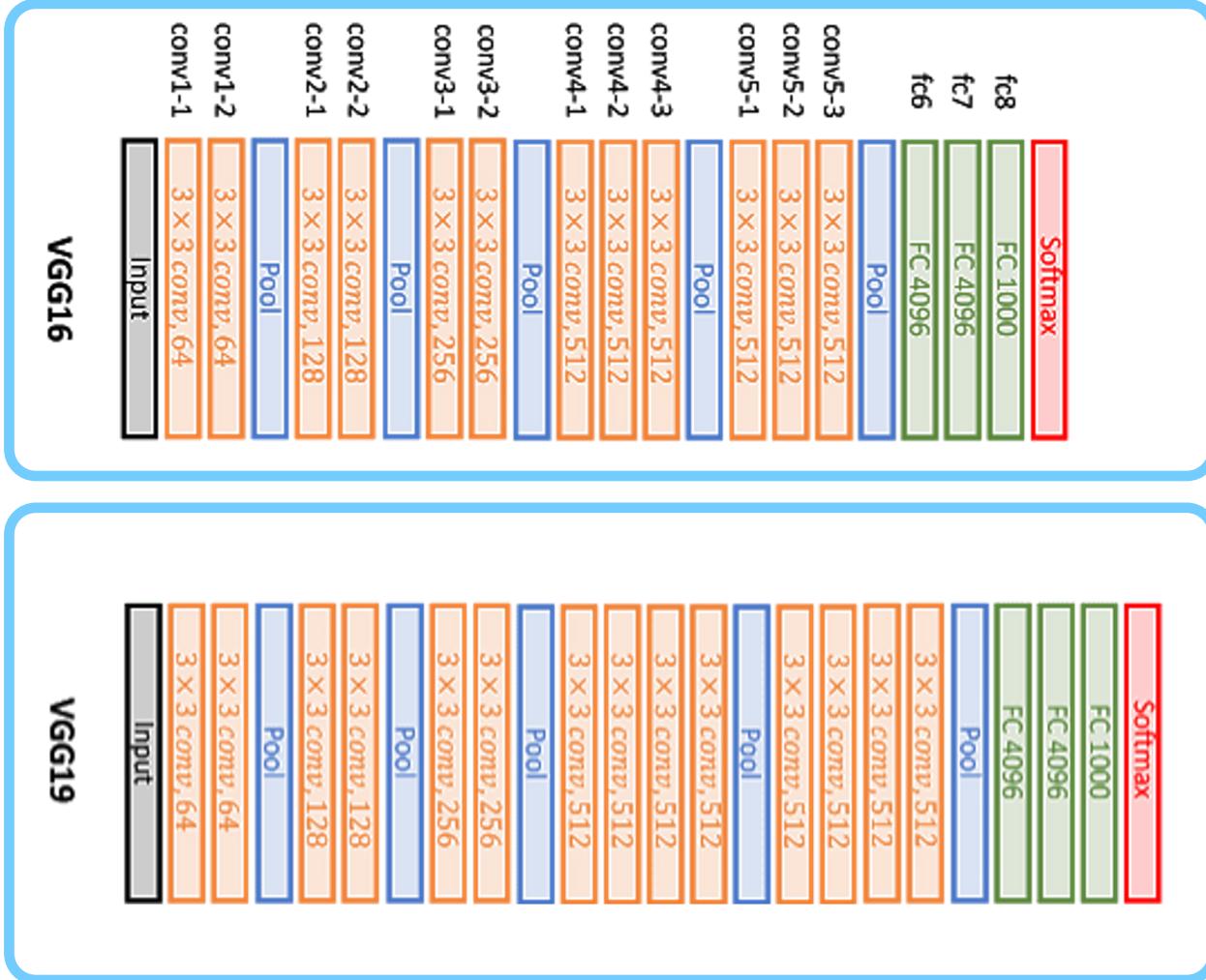
$$f(a) = \log(1 + e^a)$$



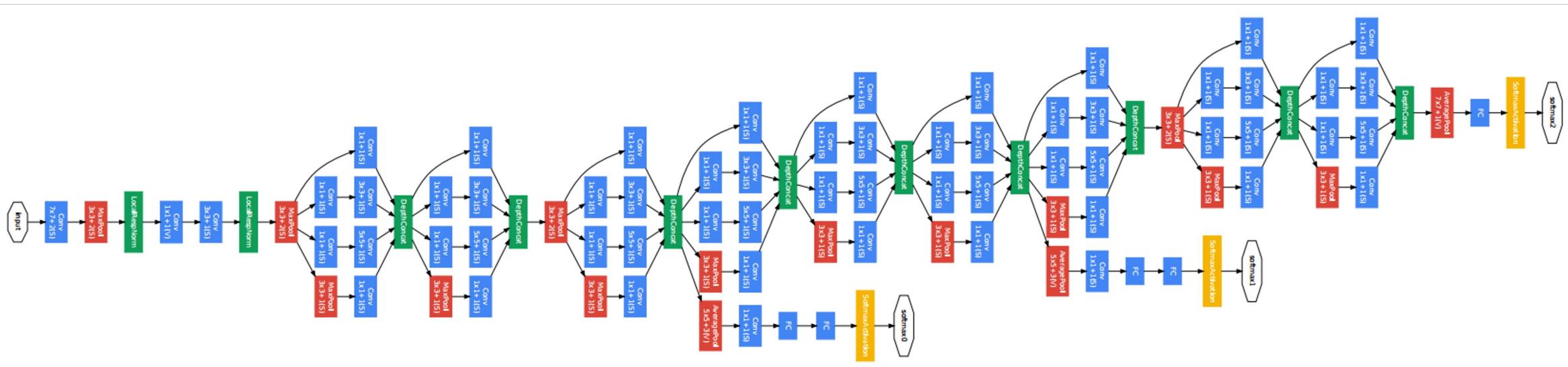
Actually, networks can be deep



And deeper...



Much deeper...

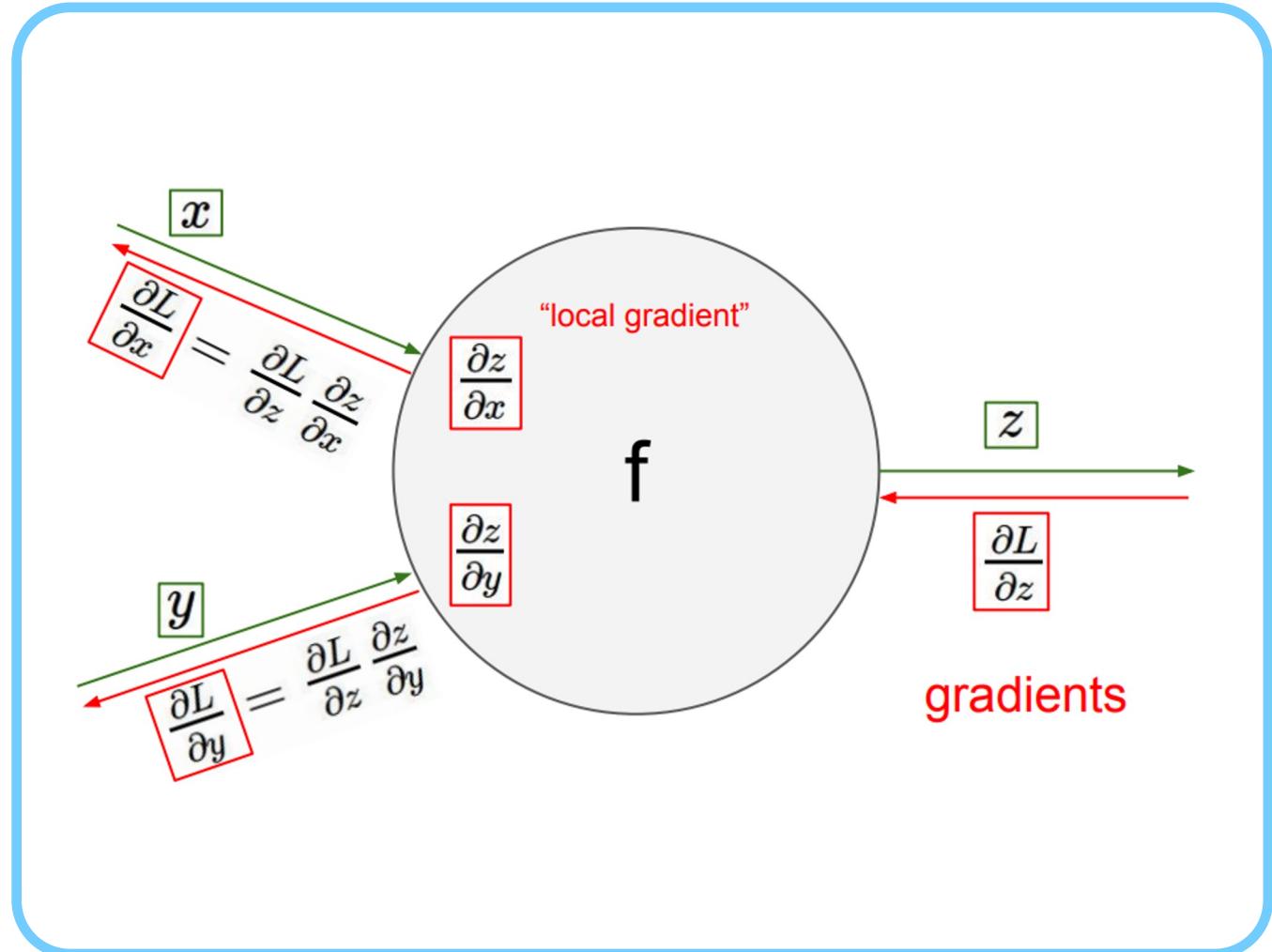


Backpropagation and chain rule

Chain rule is just simple math:

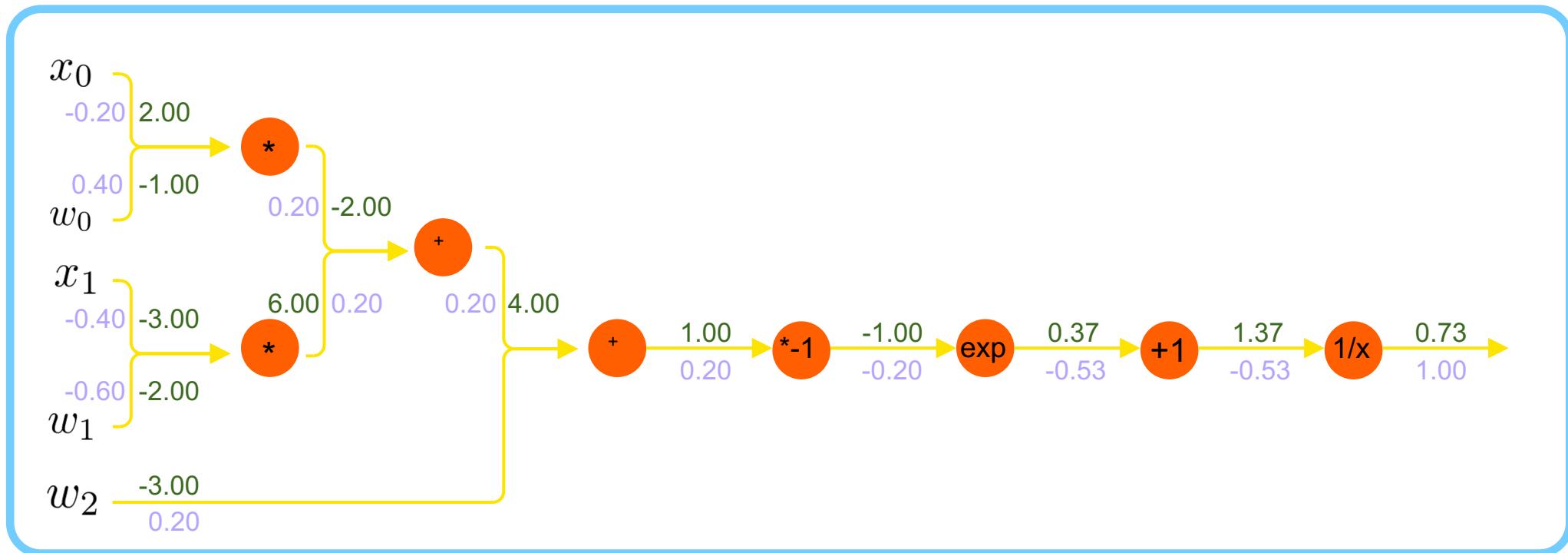
$$\frac{\partial L}{\partial x} = \frac{\partial L}{\partial z} \frac{\partial z}{\partial x}$$

Backprop is just way
to use it in NN training



Backpropagation example

$$L(\mathbf{w}, \mathbf{x}) = \frac{1}{1 + \exp(-(\mathbf{x}_0 w_0 + \mathbf{x}_1 w_1 + w_2))}$$



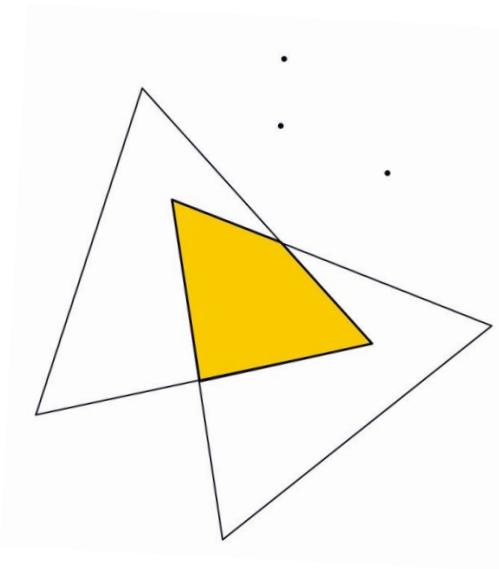
Backpropagation: matrix form

$$y_1 = f_1(\mathbf{x}) = x_1$$

$$y_2 = f_2(\mathbf{x}) = x_2$$

⋮

$$y_n = f_n(\mathbf{x}) = x_n$$



$$\frac{\partial \mathbf{y}}{\partial \mathbf{x}} = \begin{bmatrix} \nabla f_1(\mathbf{x}) \\ \nabla f_2(\mathbf{x}) \\ \dots \\ \nabla f_m(\mathbf{x}) \end{bmatrix} = \begin{bmatrix} \frac{\partial}{\partial \mathbf{x}} f_1(\mathbf{x}) \\ \frac{\partial}{\partial \mathbf{x}} f_2(\mathbf{x}) \\ \dots \\ \frac{\partial}{\partial \mathbf{x}} f_m(\mathbf{x}) \end{bmatrix} = \begin{bmatrix} \frac{\partial}{\partial x_1} f_1(\mathbf{x}) & \frac{\partial}{\partial x_2} f_1(\mathbf{x}) & \dots & \frac{\partial}{\partial x_n} f_1(\mathbf{x}) \\ \frac{\partial}{\partial x_1} f_2(\mathbf{x}) & \frac{\partial}{\partial x_2} f_2(\mathbf{x}) & \dots & \frac{\partial}{\partial x_n} f_2(\mathbf{x}) \\ \dots \\ \frac{\partial}{\partial x_1} f_m(\mathbf{x}) & \frac{\partial}{\partial x_2} f_m(\mathbf{x}) & \dots & \frac{\partial}{\partial x_n} f_m(\mathbf{x}) \end{bmatrix}$$

Backpropagation: matrix form

	scalar	vector
scalar	x	\mathbf{x}
vector	f	$\frac{\partial f}{\partial \mathbf{x}}$
	$\frac{\partial f}{\partial x}$	$\frac{\partial \mathbf{f}}{\partial \mathbf{x}}$

Backpropagation: matrix form

$$\frac{\partial \mathbf{y}}{\partial \mathbf{x}} = \begin{bmatrix} \frac{\partial}{\partial \mathbf{x}} f_1(\mathbf{x}) \\ \frac{\partial}{\partial \mathbf{x}} f_2(\mathbf{x}) \\ \vdots \\ \frac{\partial}{\partial \mathbf{x}} f_m(\mathbf{x}) \end{bmatrix} = \begin{bmatrix} \frac{\partial}{\partial x_1} f_1(\mathbf{x}) & \frac{\partial}{\partial x_2} f_1(\mathbf{x}) & \dots & \frac{\partial}{\partial x_n} f_1(\mathbf{x}) \\ \frac{\partial}{\partial x_1} f_2(\mathbf{x}) & \frac{\partial}{\partial x_2} f_2(\mathbf{x}) & \dots & \frac{\partial}{\partial x_n} f_2(\mathbf{x}) \\ \vdots \\ \frac{\partial}{\partial x_1} f_m(\mathbf{x}) & \frac{\partial}{\partial x_2} f_m(\mathbf{x}) & \dots & \frac{\partial}{\partial x_n} f_m(\mathbf{x}) \end{bmatrix}$$

$$= \begin{bmatrix} \frac{\partial}{\partial x_1} x_1 & \frac{\partial}{\partial x_2} x_1 & \dots & \frac{\partial}{\partial x_n} x_1 \\ \frac{\partial}{\partial x_1} x_2 & \frac{\partial}{\partial x_2} x_2 & \dots & \frac{\partial}{\partial x_n} x_2 \\ \vdots \\ \frac{\partial}{\partial x_1} x_n & \frac{\partial}{\partial x_2} x_n & \dots & \frac{\partial}{\partial x_n} x_n \end{bmatrix}$$

(and since $\frac{\partial}{\partial x_j} x_i = 0$ for $j \neq i$)

$$= \begin{bmatrix} \frac{\partial}{\partial x_1} x_1 & 0 & \dots & 0 \\ 0 & \frac{\partial}{\partial x_2} x_2 & \dots & 0 \\ \ddots \\ 0 & 0 & \dots & \frac{\partial}{\partial x_n} x_n \end{bmatrix}$$

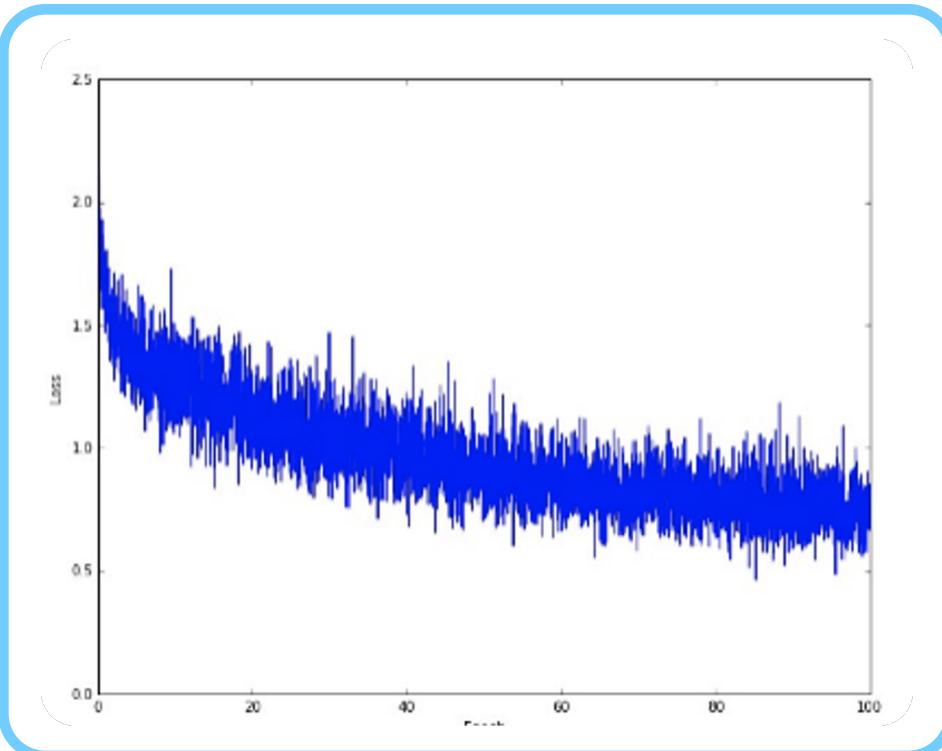
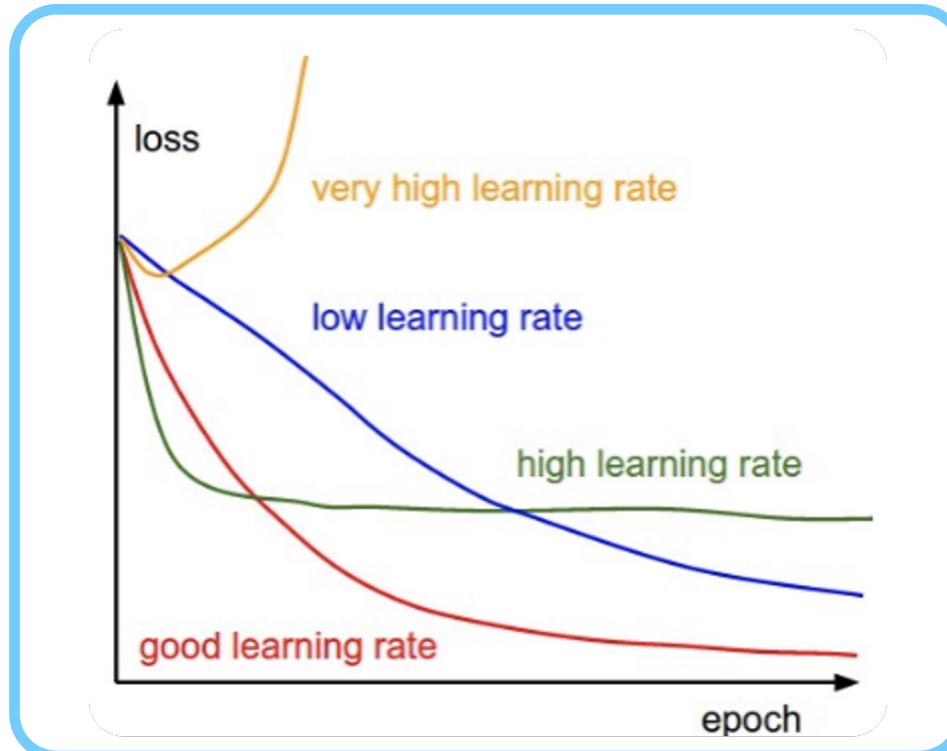
$$= \begin{bmatrix} 1 & 0 & \dots & 0 \\ 0 & 1 & \dots & 0 \\ \ddots \\ 0 & 0 & \dots & 1 \end{bmatrix}$$

$= I$ (I is the identity matrix with ones down the diagonal)

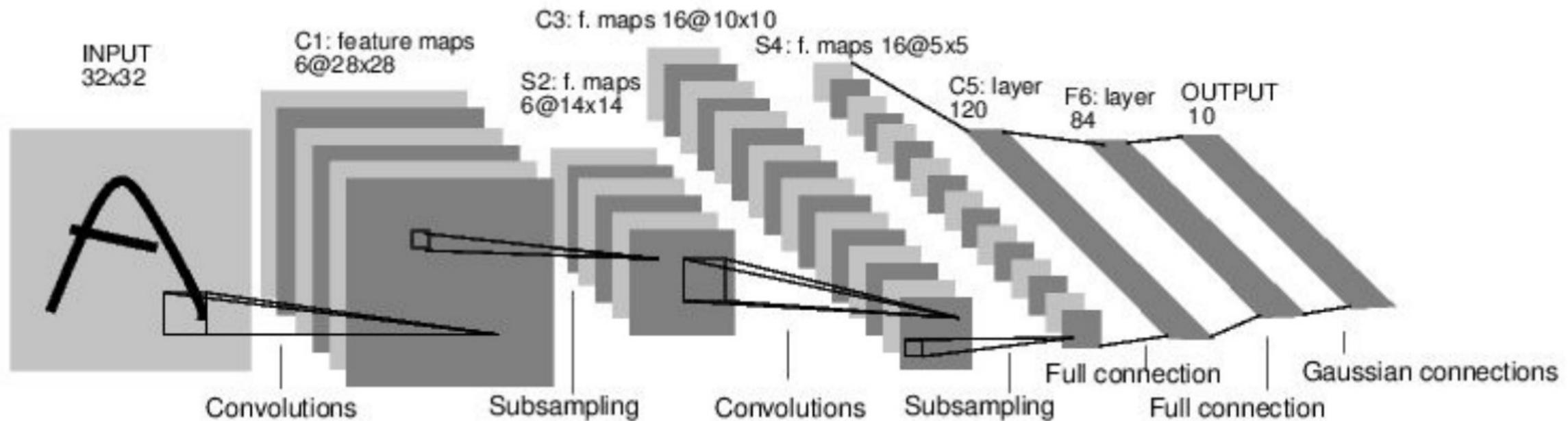
Gradient optimization

Stochastic gradient descent is used to optimize NN parameters

$$\boldsymbol{w}_{t+1} = \boldsymbol{w}_t - \eta \frac{dL}{d\boldsymbol{w}}$$



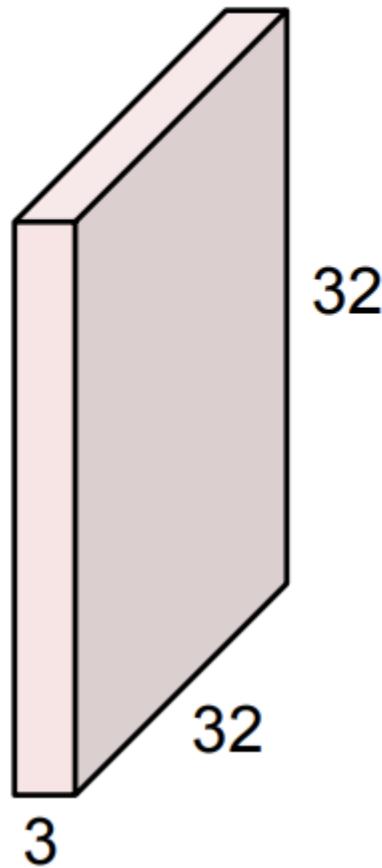
Convolutions in Neural Networks



[LeNet-5, LeCun 1998]

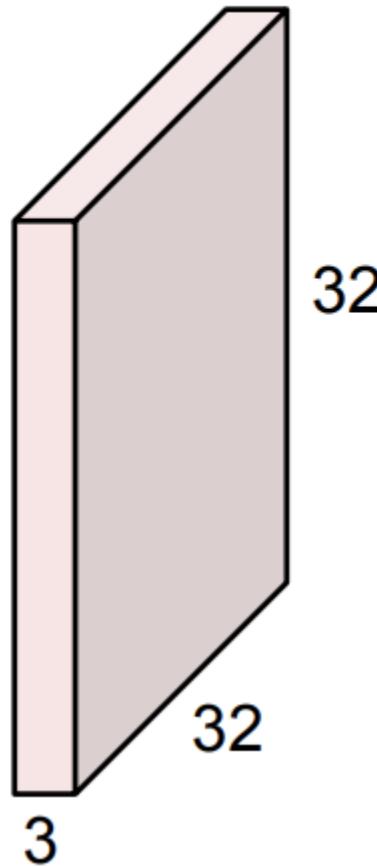
Convolutional layer

32x32x3 image



Convolutional layer

32x32x3 image



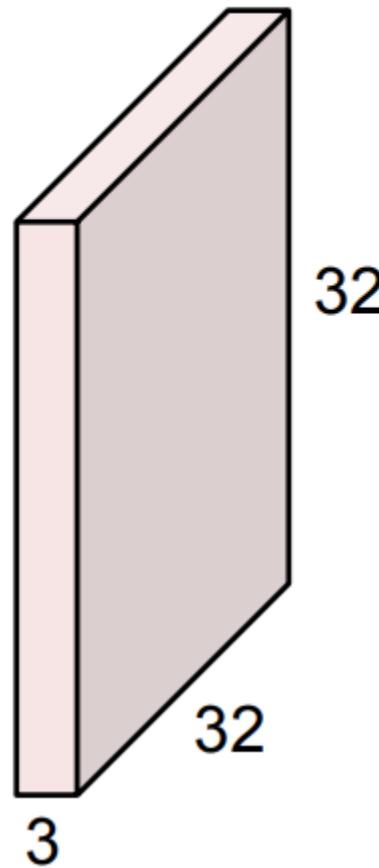
5x5x3 filter



Convolve the filter with the image
i.e. “slide over the image spatially,
computing dot products”

Convolutional layer

32x32x3 image



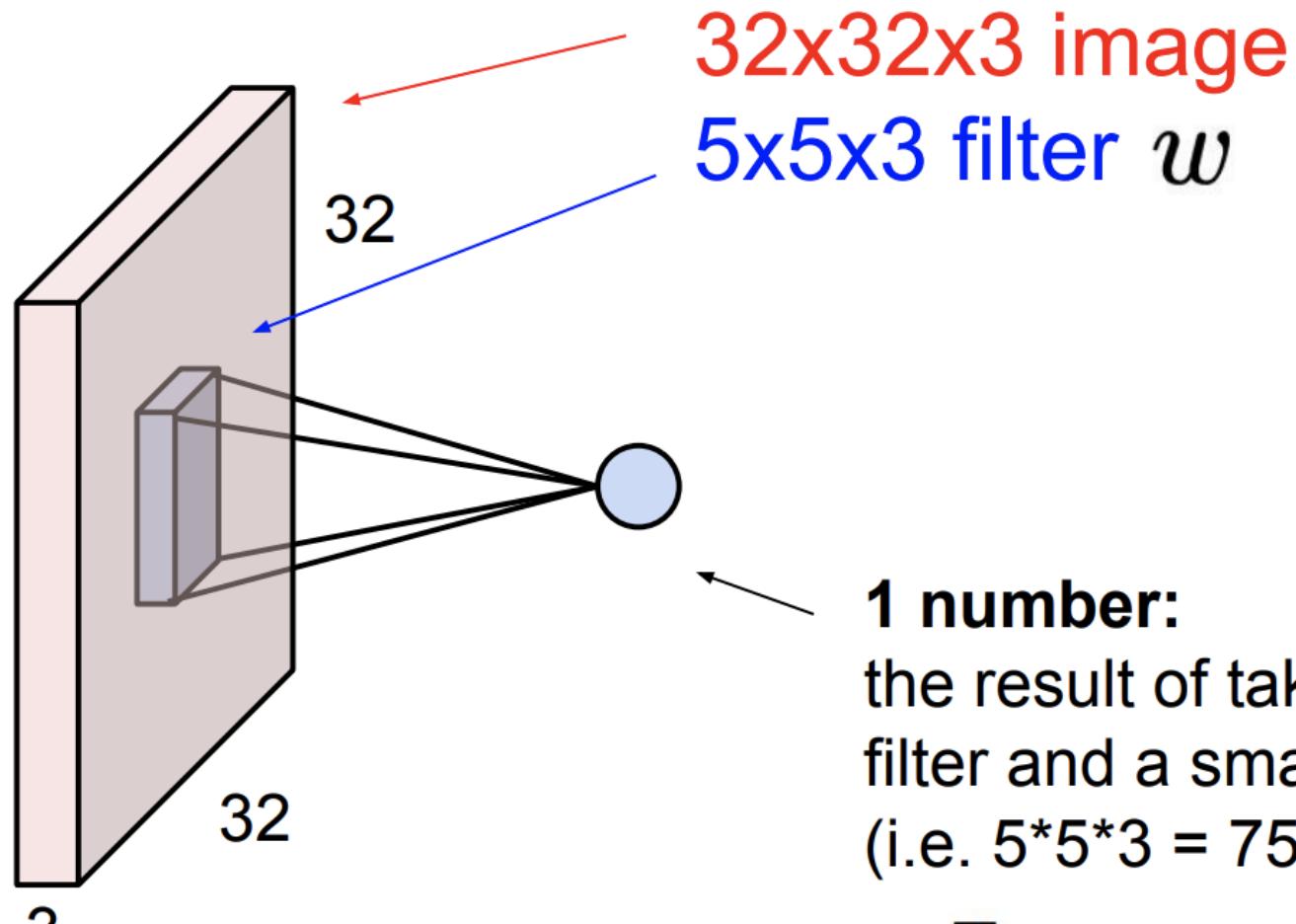
Filters extend the depth of the original image

5x5x3 filter



Convolve the filter with the image
i.e. “slide over the image spatially,
computing dot products”

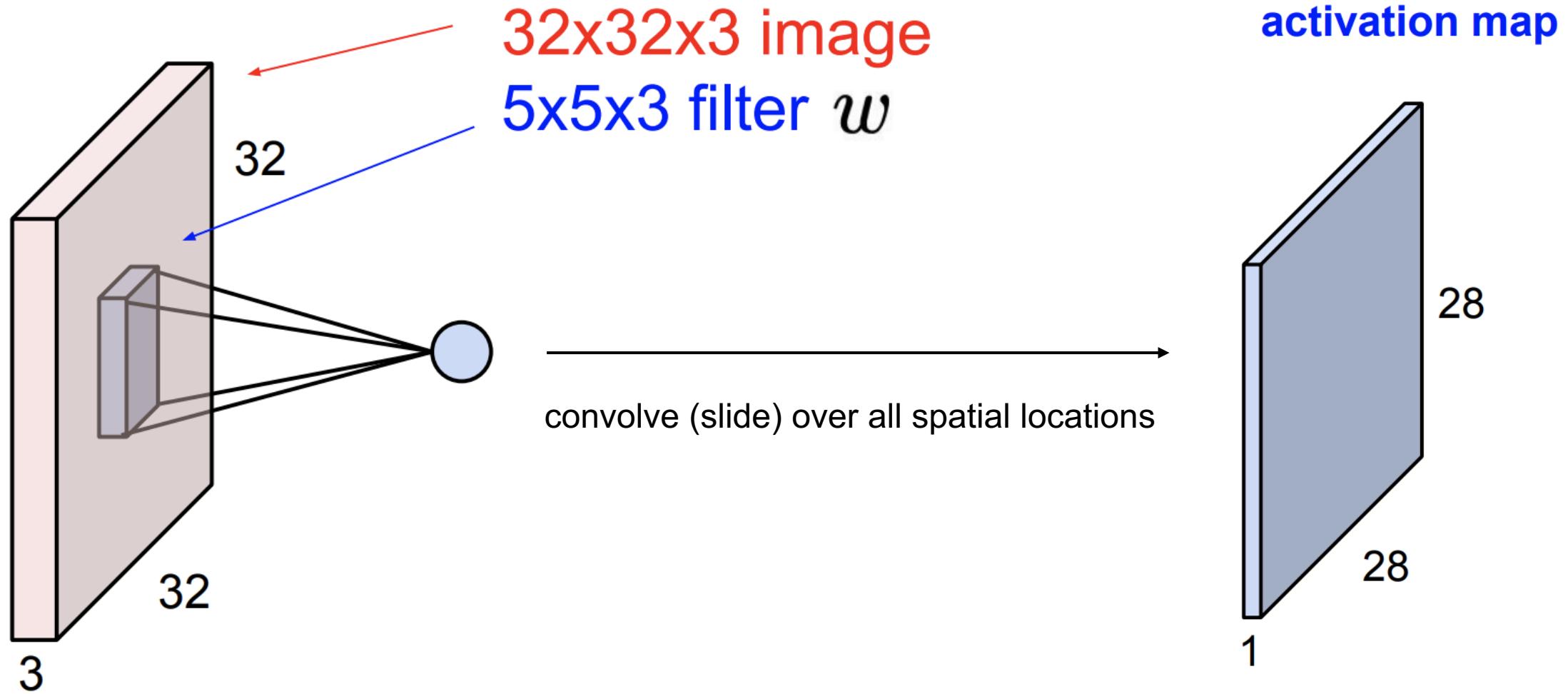
Convolutional layer



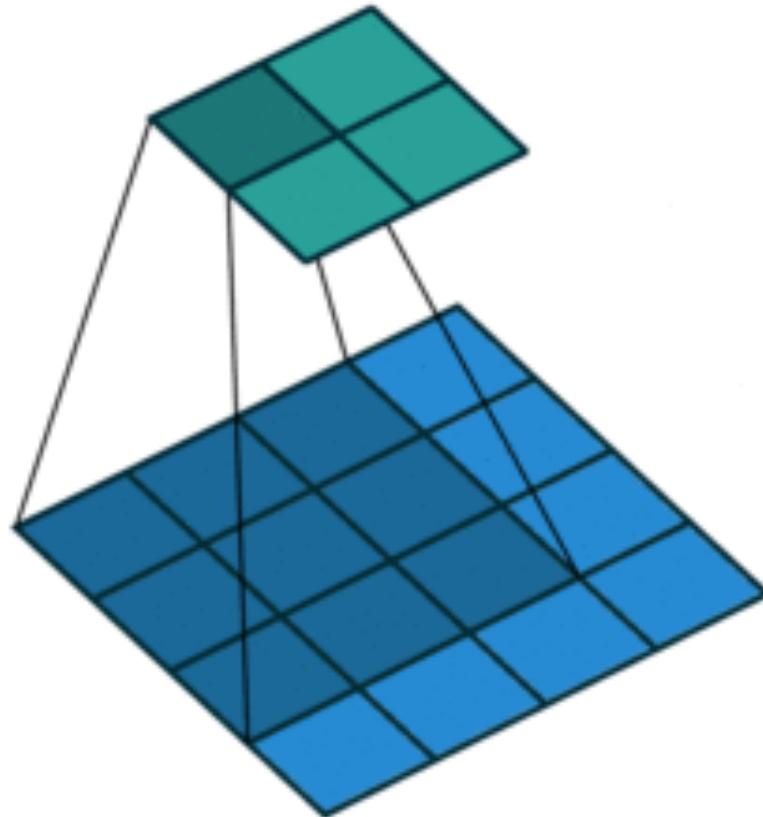
1 number:
the result of taking a dot product between the
filter and a small $5 \times 5 \times 3$ chunk of the image
(i.e. $5 \times 5 \times 3 = 75$ -dimensional dot product + bias)

$$w^T x + b$$

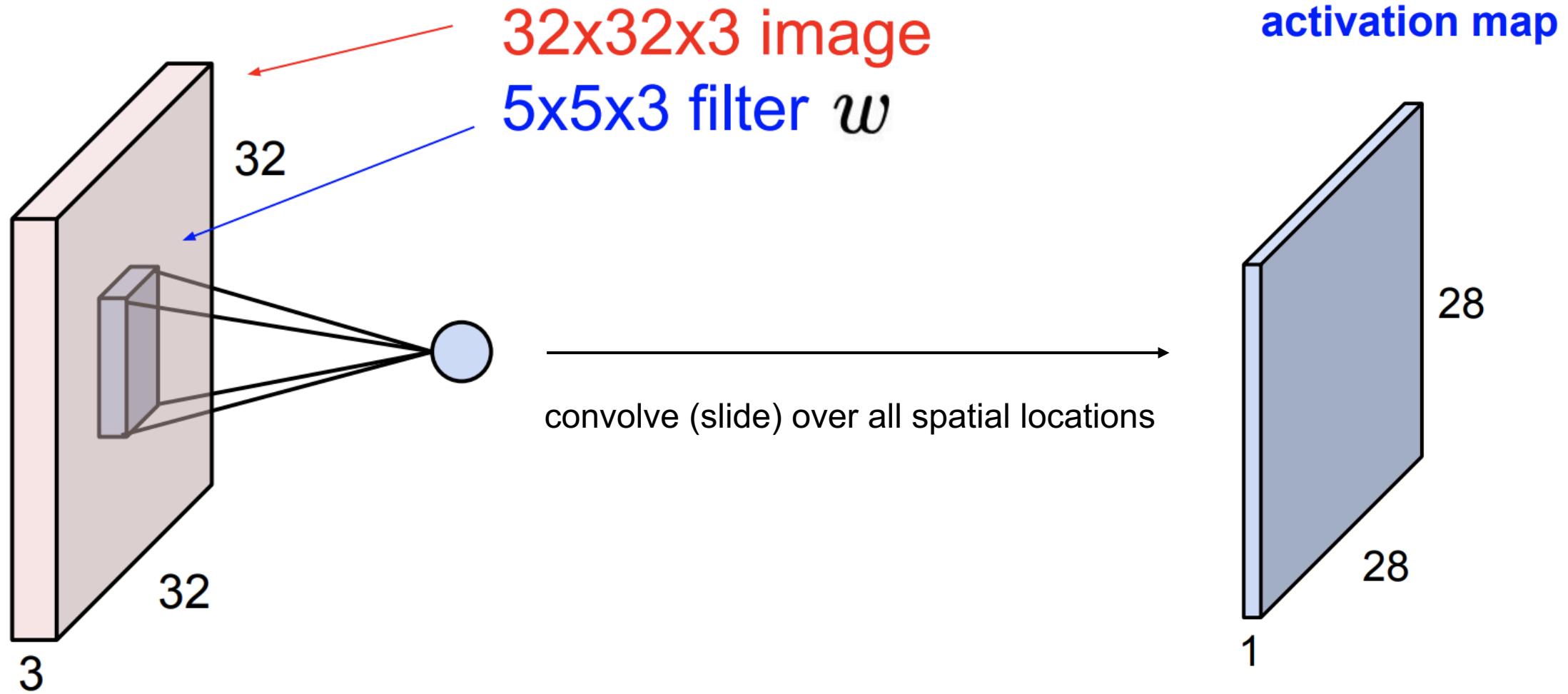
Convolutional layer



Convolutional layer

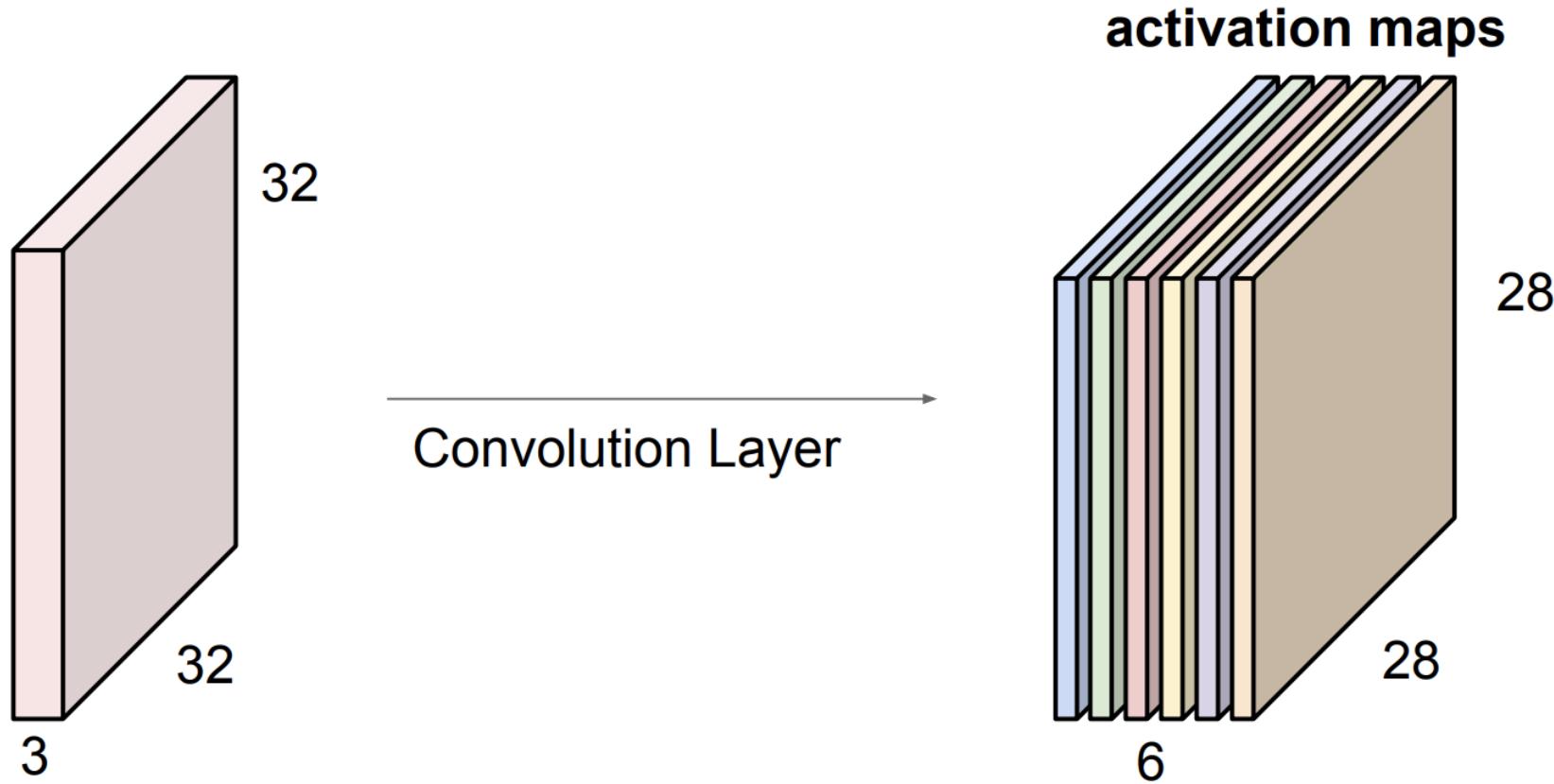


Convolutional layer



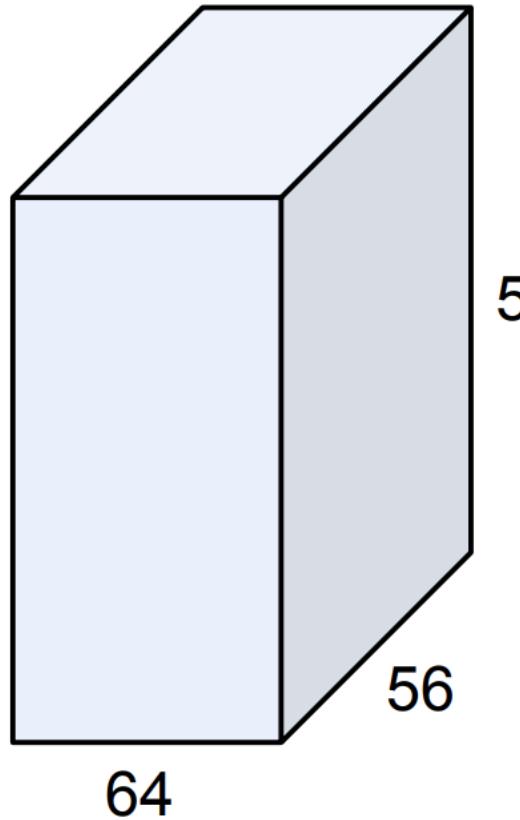
Convolutional layer

For example, if we had 6 5x5 filters, we'll get 6 separate activation maps:



We stack these up to get a “new image” of size 28x28x6!

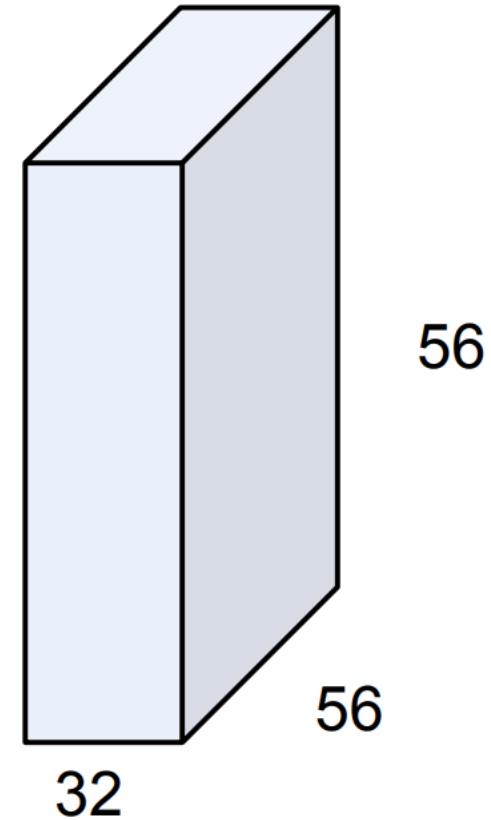
Pointwise (1x1) Convolution



1x1 CONV
with 32 filters

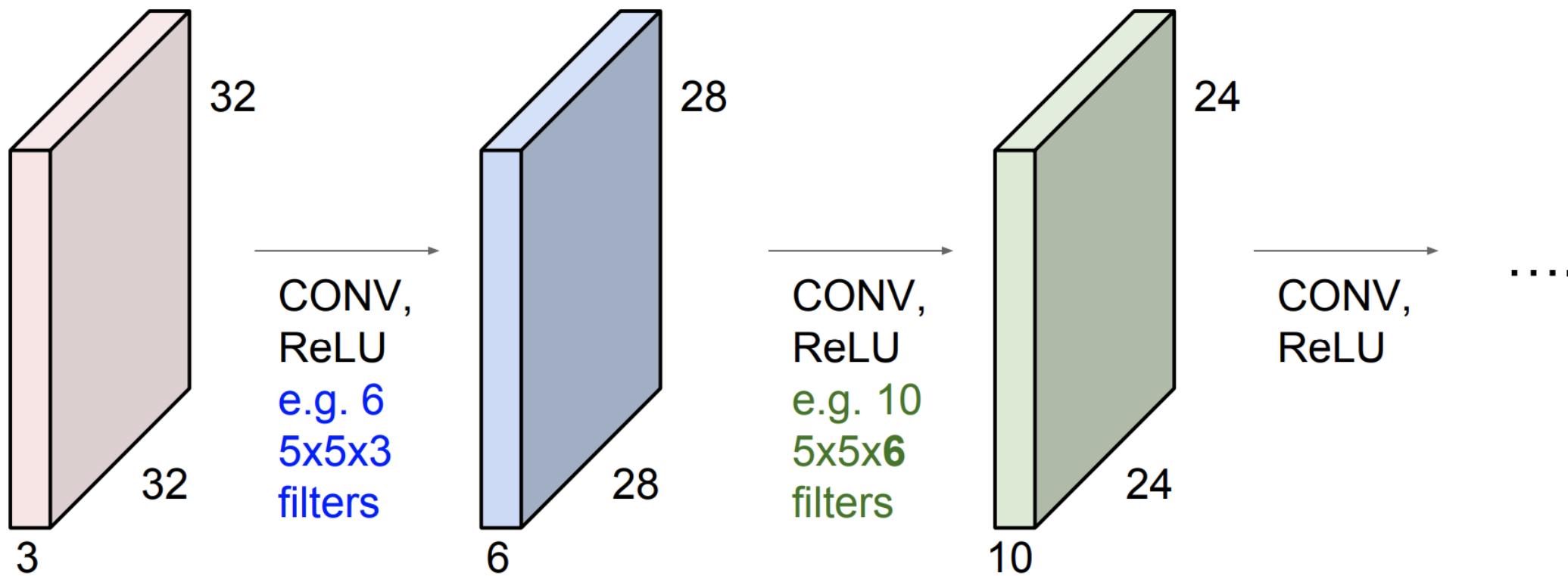
→

(each filter has size
 $1 \times 1 \times 64$, and performs a
64-dimensional dot
product)

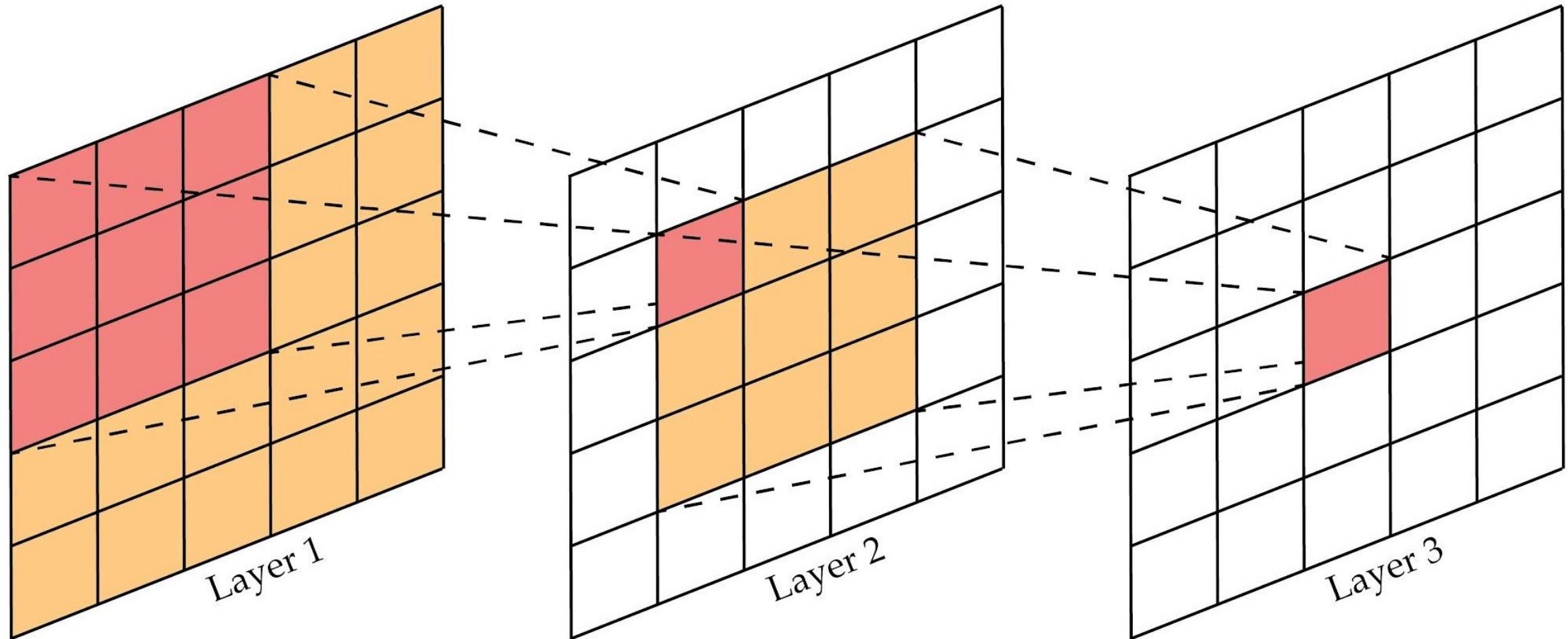


Convolutional layer

Preview: ConvNet is a sequence of Convolutional Layers, interspersed with activation functions



Receptive field

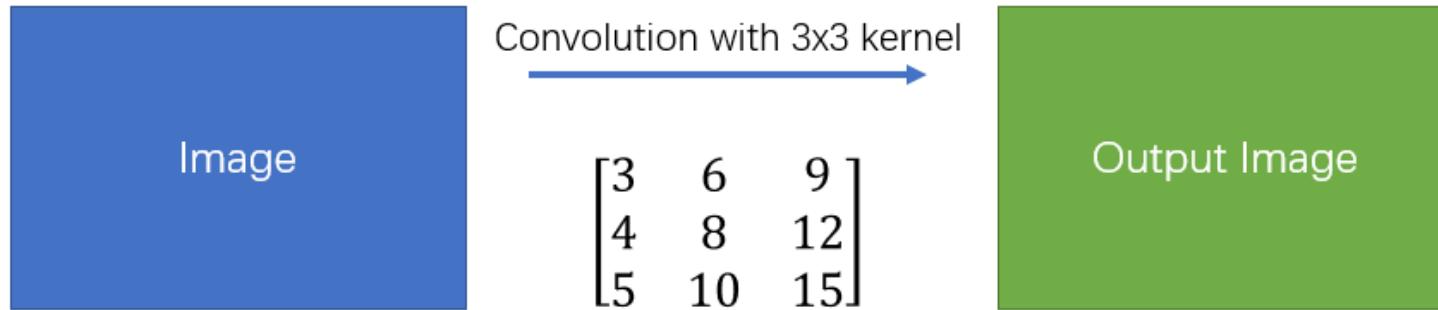


Spatial Separable Convolution

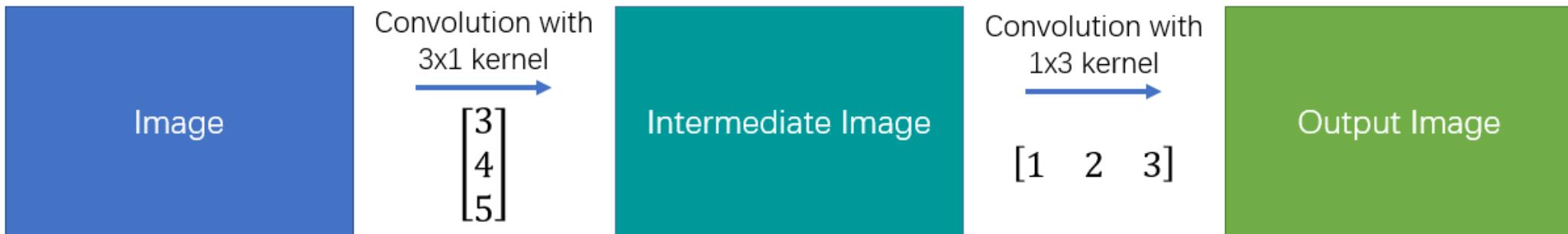
$$\begin{bmatrix} 3 & 6 & 9 \\ 4 & 8 & 12 \\ 5 & 10 & 15 \end{bmatrix} = \begin{bmatrix} 3 \\ 4 \\ 5 \end{bmatrix} \times [1 \quad 2 \quad 3]$$

Spatial Separable Convolution

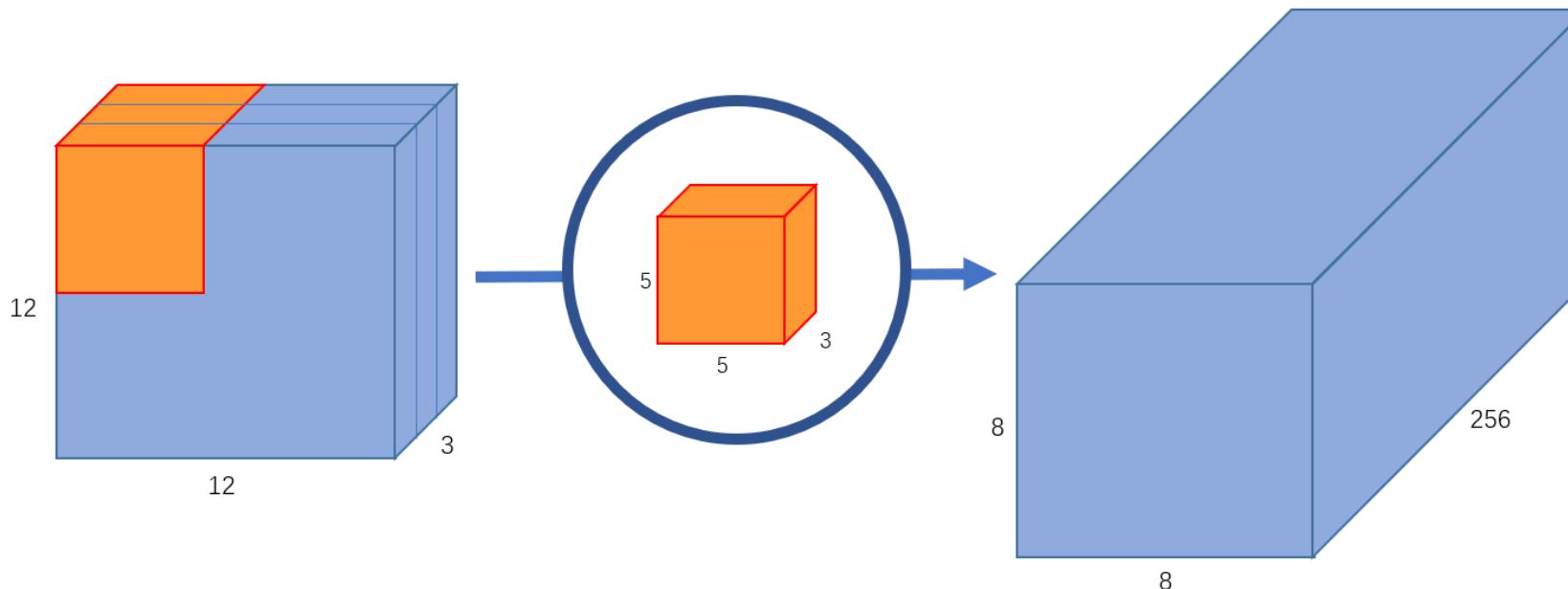
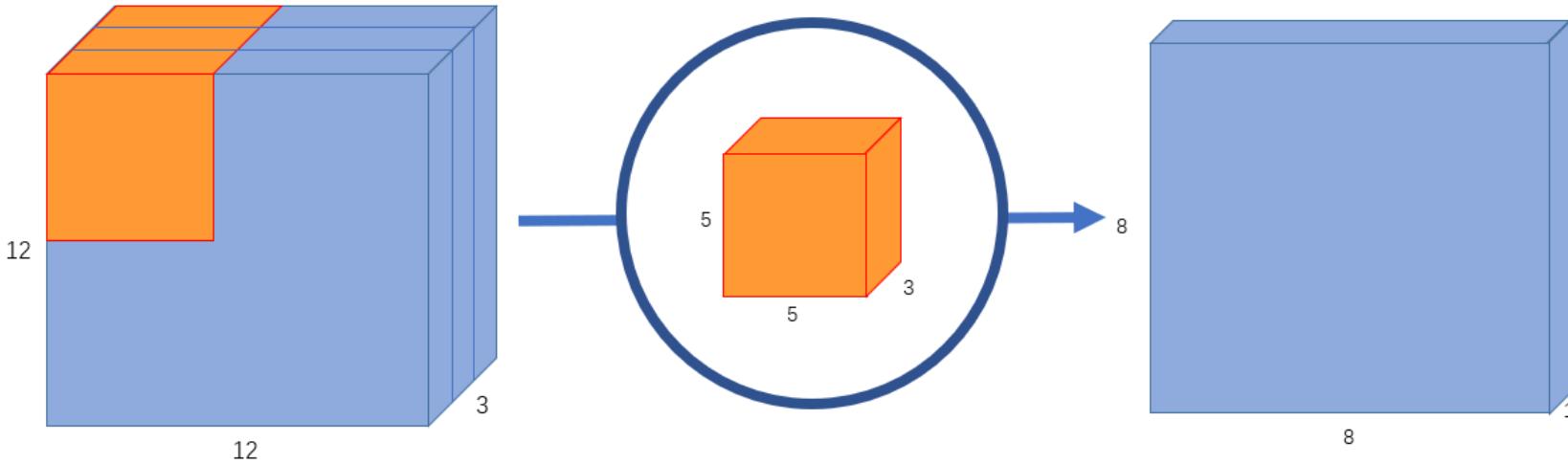
Simple Convolution



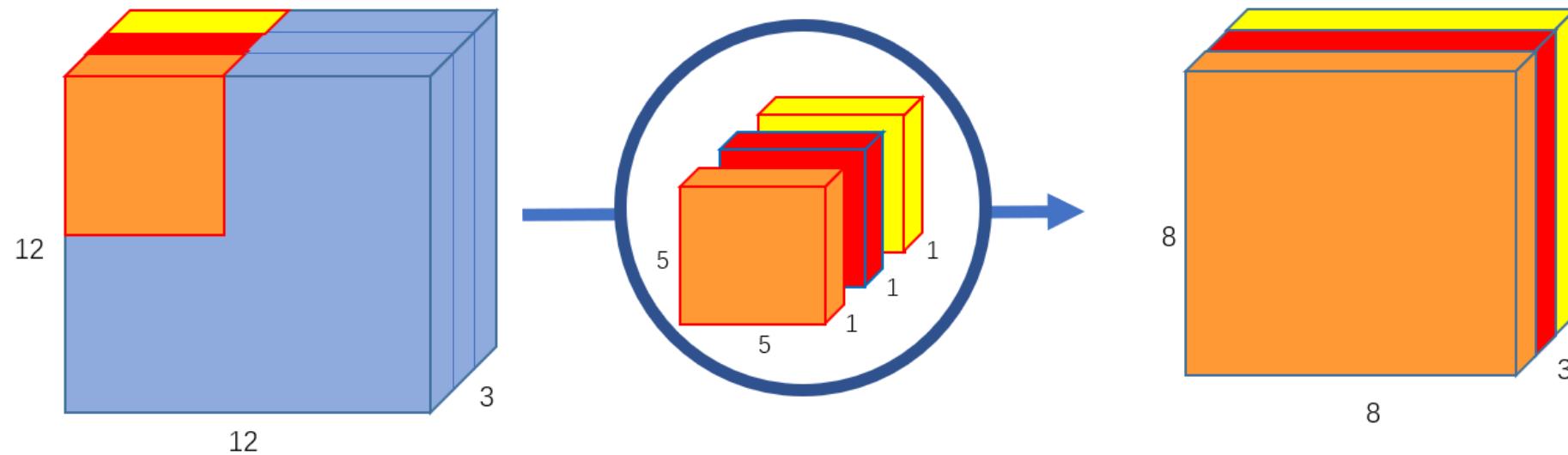
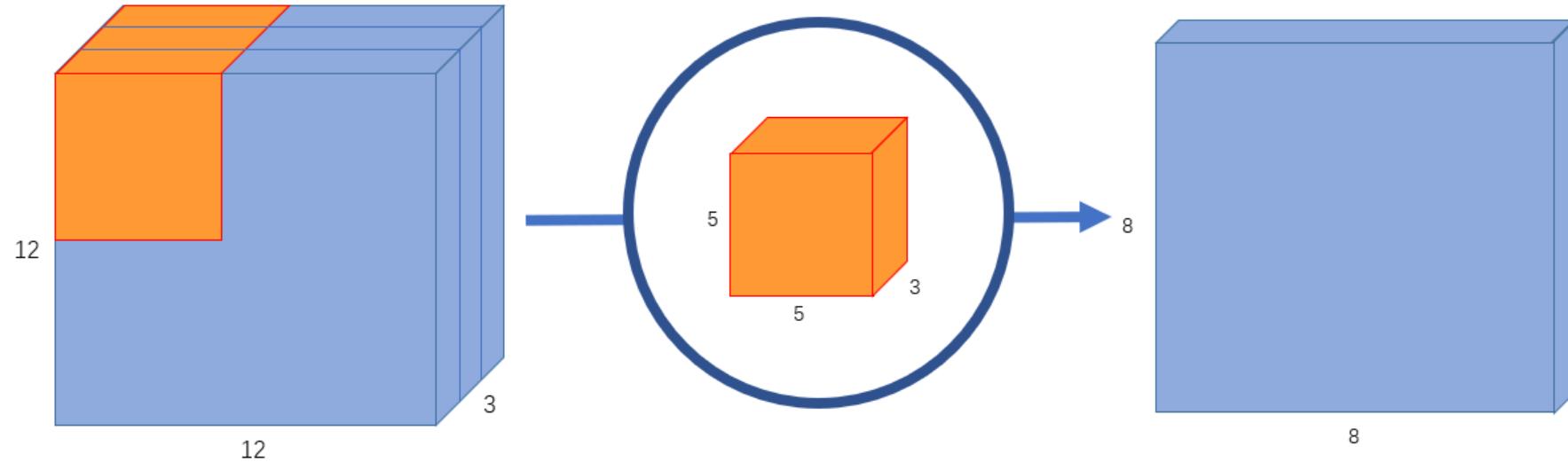
Spatial Separable Convolution



Normal Convolution

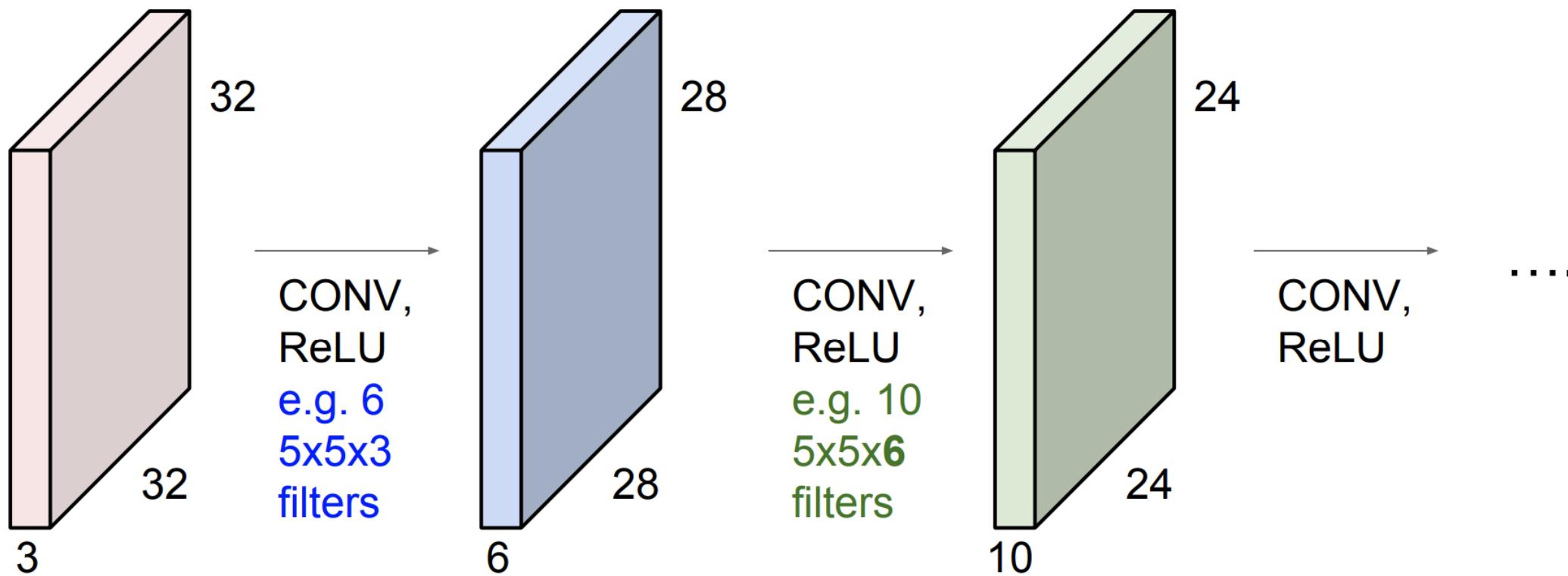


Depthwise Convolution



Convolutional layer

Preview: ConvNet is a sequence of Convolutional Layers, interspersed with activation functions



Vanishing gradient

Vanishing gradient is present in **all** deep neural network architectures.

- Due to chain rule / choice of nonlinearity function, gradient can become vanishingly small during backpropagation
- Lower levels are hard to train and are trained slower
- Solutions
 - skip-connections (ResNet)

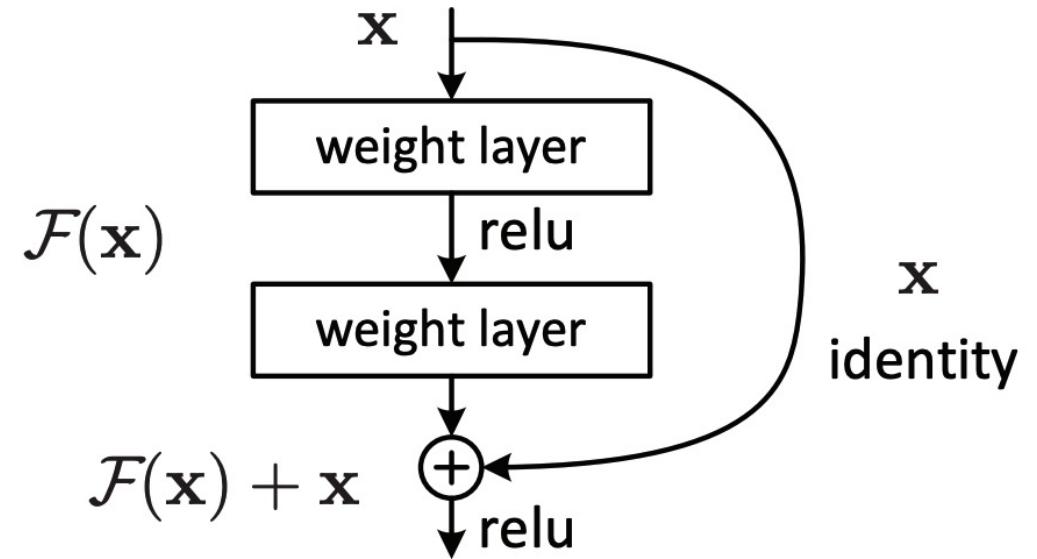
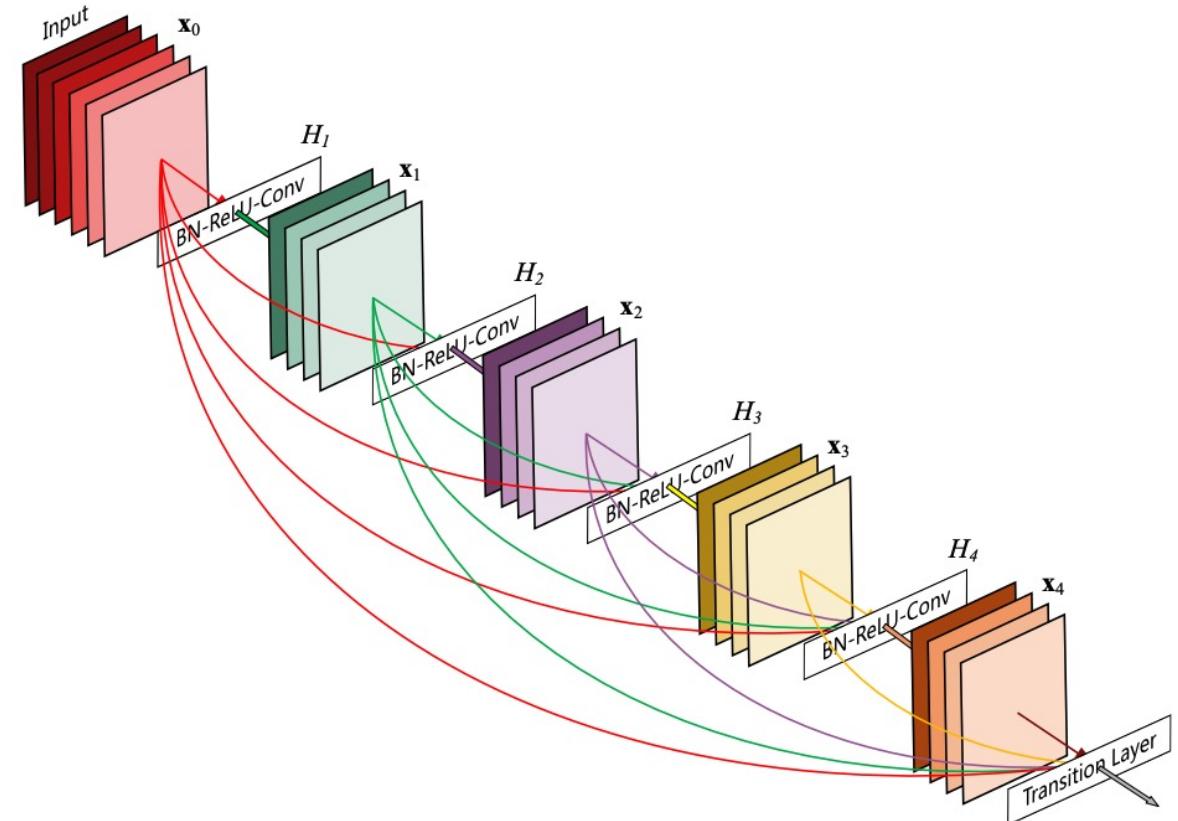


Figure 2. Residual learning: a building block.

Vanishing gradient

Vanishing gradient is present in all deep neural network architectures.

- Due to chain rule / choice of nonlinearity function, gradient can become vanishingly small during backpropagation
- Lower levels are hard to train and are trained slower
- Solutions
 - skip-connections (ResNet)
 - dense connections (DenseNet)



Classification architectures overview

03

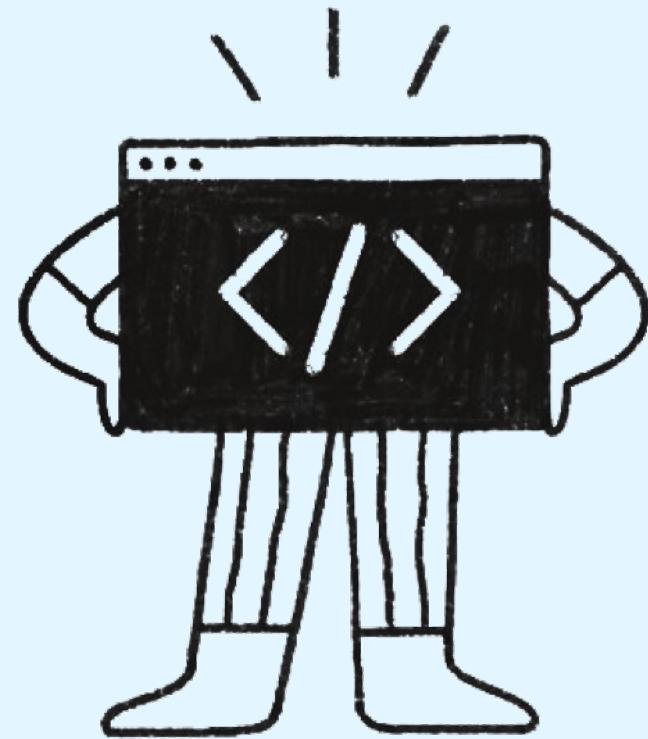


Image classification then

2007: Asirra

CAPTCHA that asks users to identify cats out of a set of 12 photographs of both cats and dogs.

Solved by **humans** in **99.6%** of cases.

Barring a major advance in machine vision, we expect **computers** will have no better than a **1/54,000 chance** of solving it

Please click on all the images that show cats:



[Elson et al](#)

Image classification now

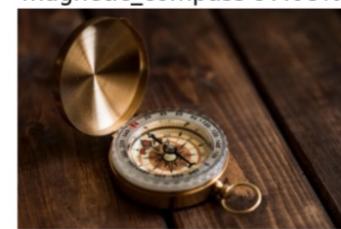
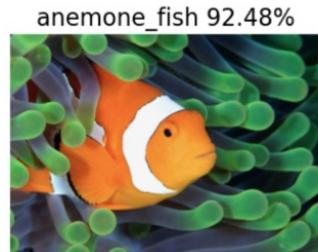
2014: Asirra on kaggle

[Competition page](#)

#	△	Team	Members	Score	Entries	Last	Solution
1	—	Pierre Sermanet		0.98914	5	10y	
2	▲ 4	orchid		0.98308	17	10y	Last Submission: 02/02/2014, 1:43:19 GMT+4
3	—	Owen		0.98171	15	10y	
4	—	Paul Covington		0.98171	3	10y	
5	▼ 3	Maxim Milakov		0.98137	24	10y	
6	▼ 1	we've been in KAIST	 	0.98102	8	10y	
7	▲ 1	Doug Koch		0.98057	6	10y	
8	▲ 2	fastml.com/cats-and-dogs		0.98000	6	10y	

IMAGENET competition

- ImageNet Large Scale Visual Recognition Challenge (ILSVRC)
- since 2010
- 14.2m images
- 1000 classes
- <http://image-net.org/explore>



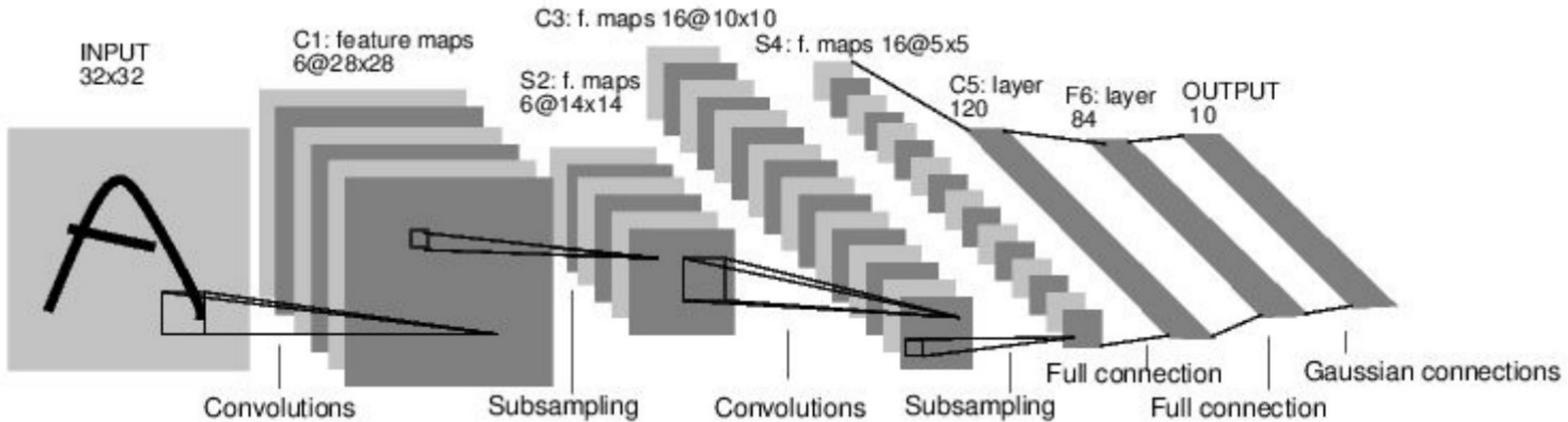
IMAGENET dataset

- Most of classes are dogs
- Most correlated classes are chihuahua and cupcake

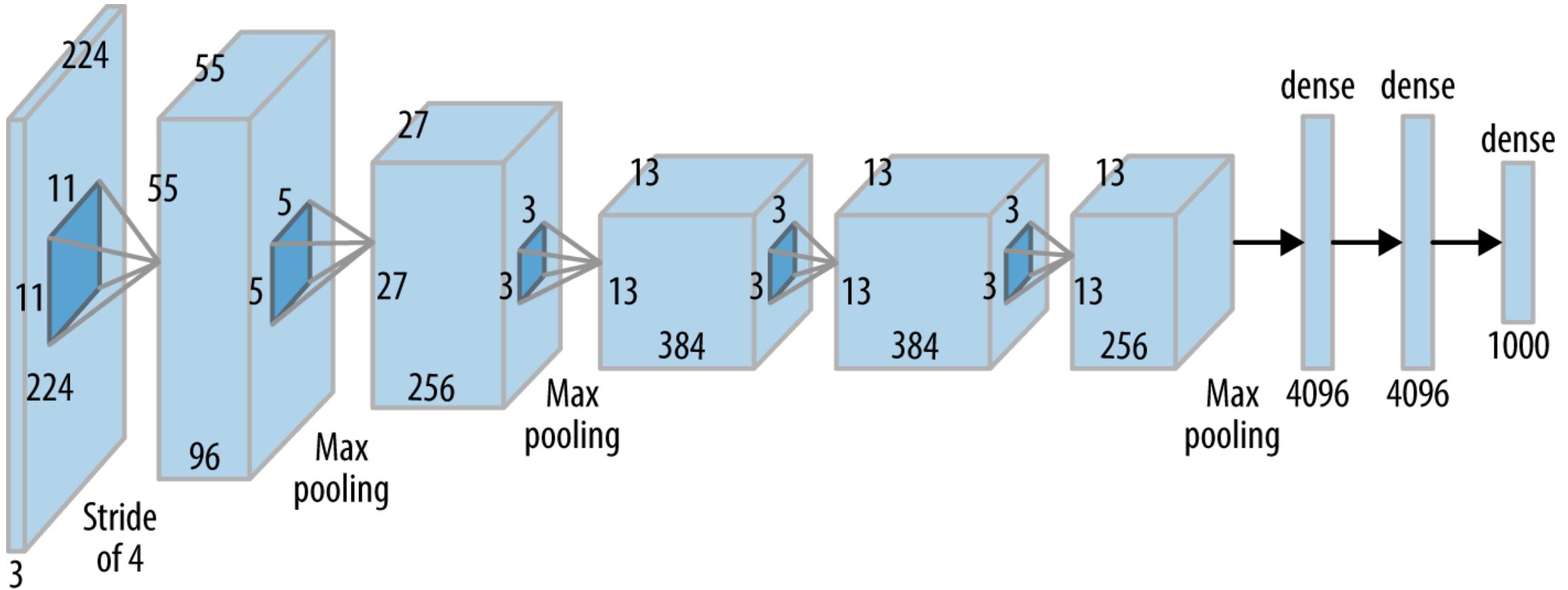
Чтобы подтвердить, что вы не робот, отметьте
КЕКСЫ



LeNet-5



AlexNet

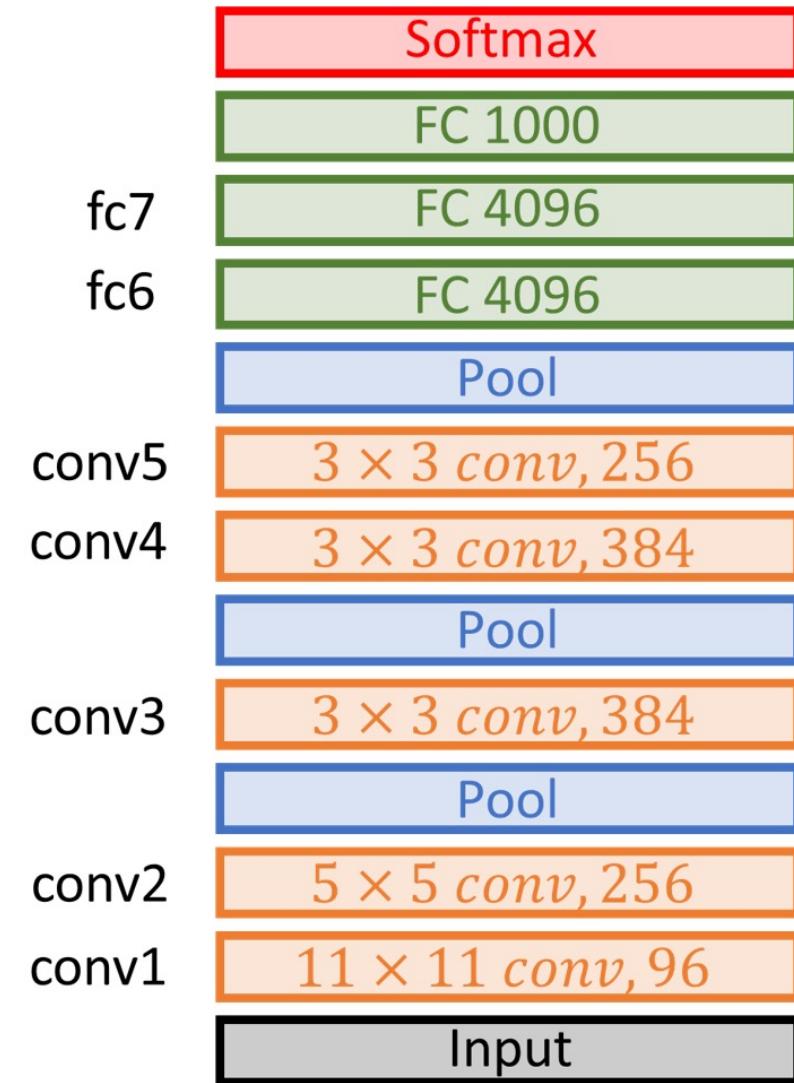


AlexNet

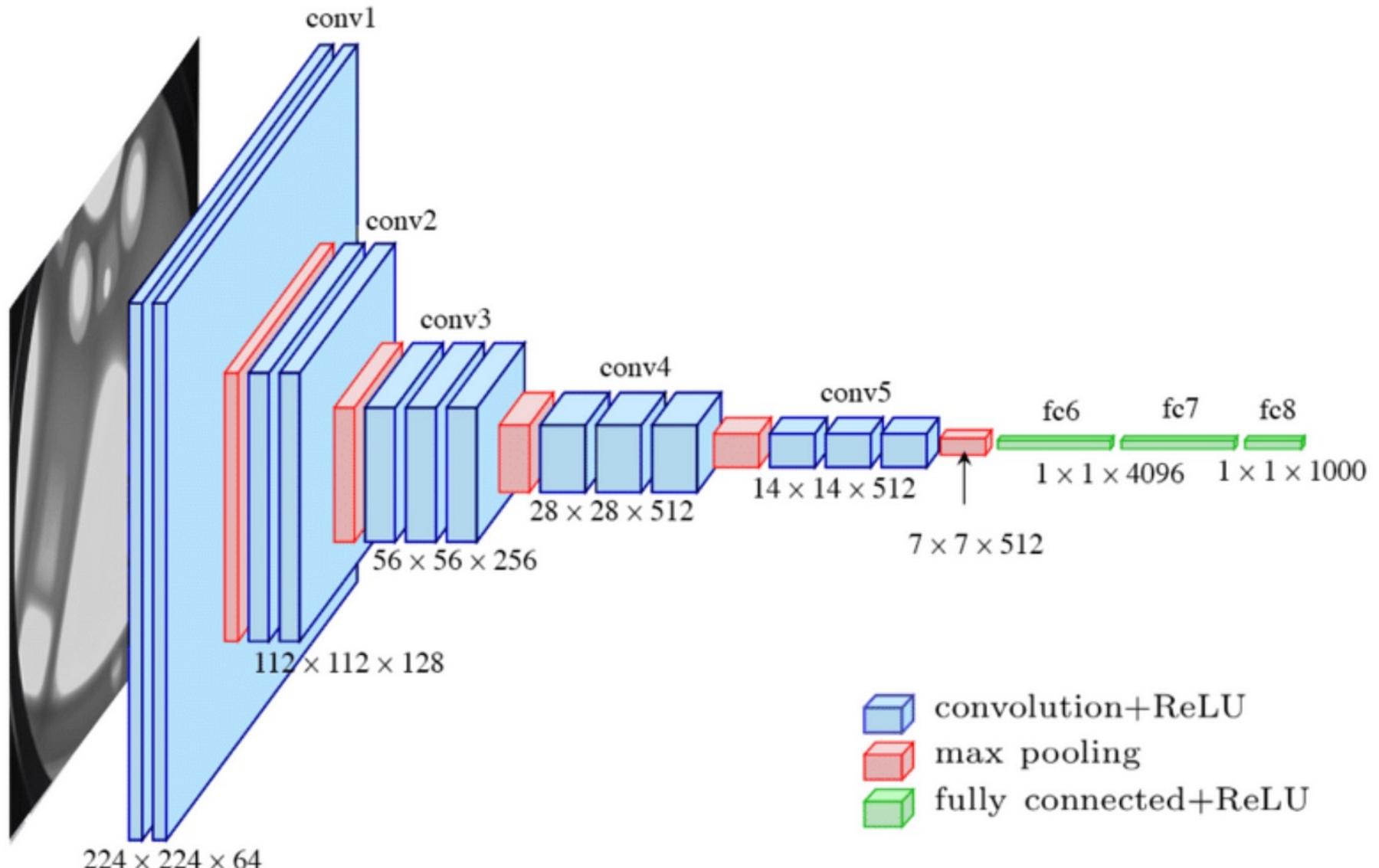
- Krizhevsky et al. 2012
- Used 2 GPUs for training
- First win of ImageNet competition with Deep Neural Networks
- Used big convolutions
- Fully connected at the end
- 7 layers

Takeouts:

- First use of ReLU
- Increasing number of channels
- 16.4%

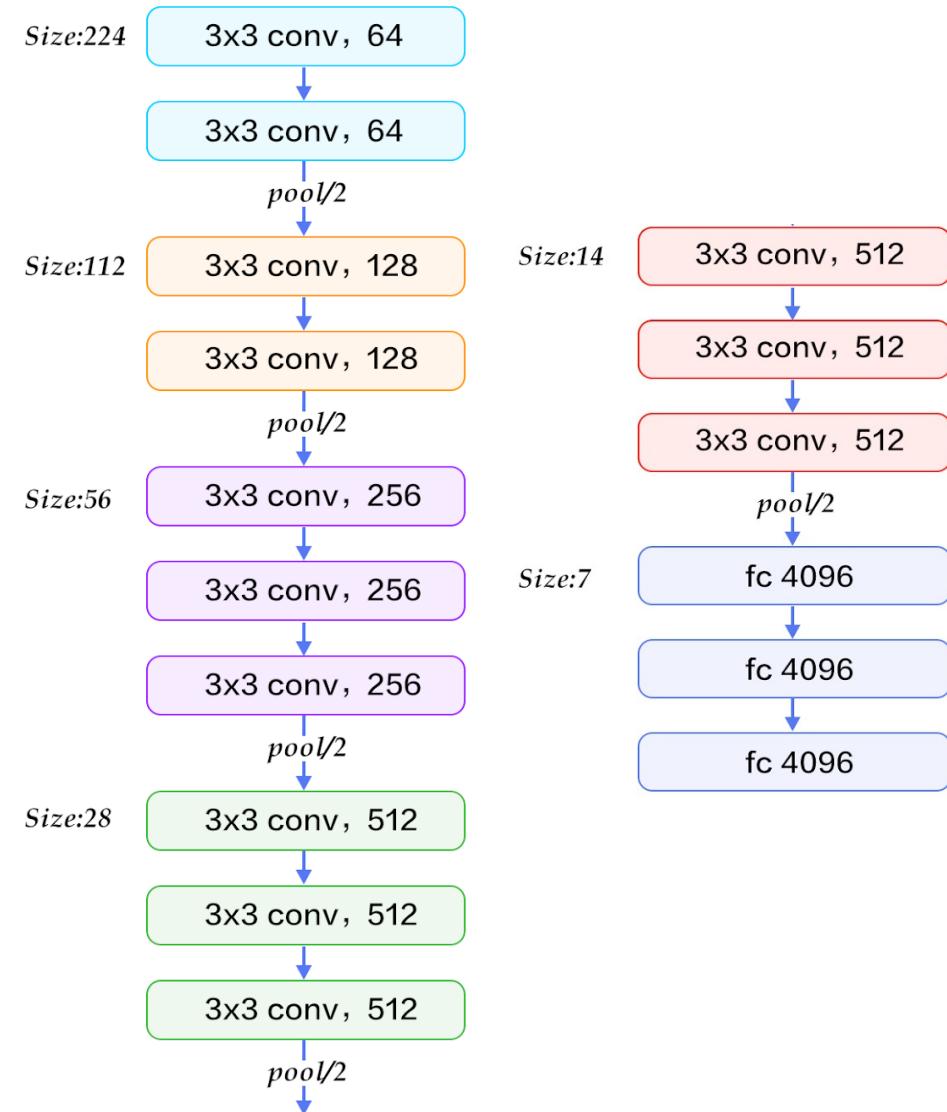


VGGNet

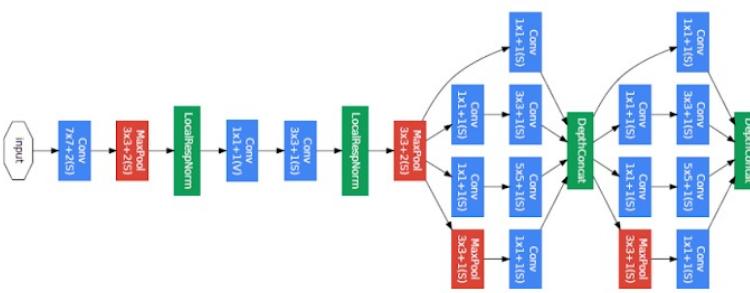


VGGNet

- Simonyan, Zisserman 2014
- All convolutions are 3x3
- Exponentially increasing channels
- 19 layers
- Requires a lot of memory
- Gradient vanishes =(



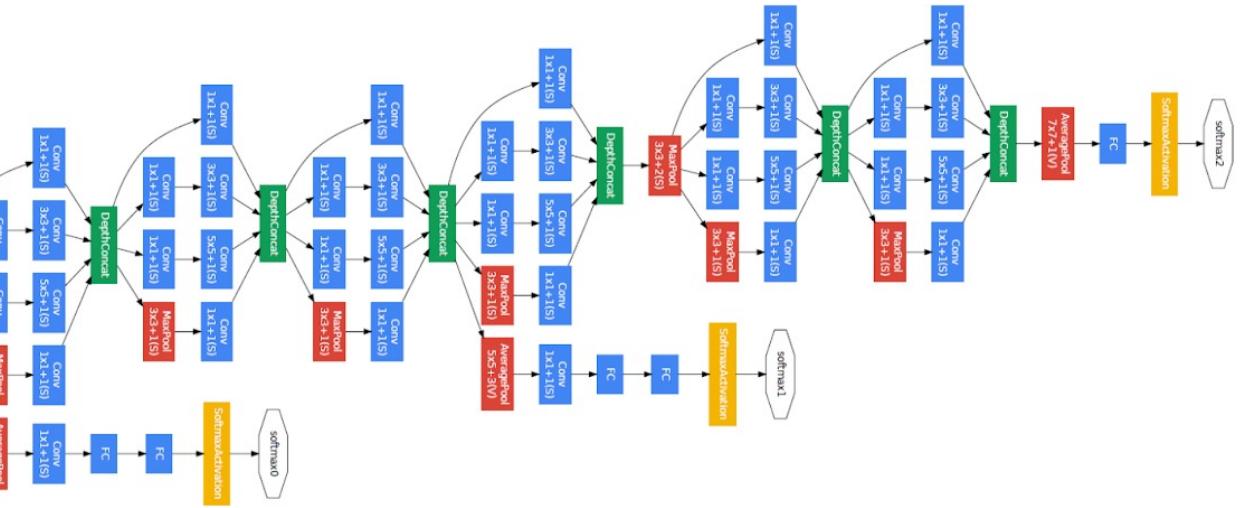
GoogLeNet



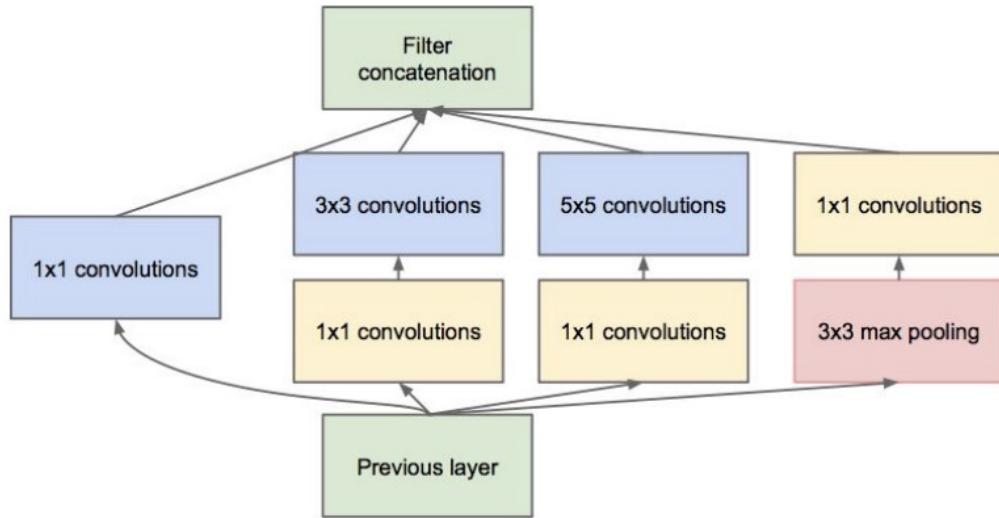
Inception module



[Szegedy et al., 2014]



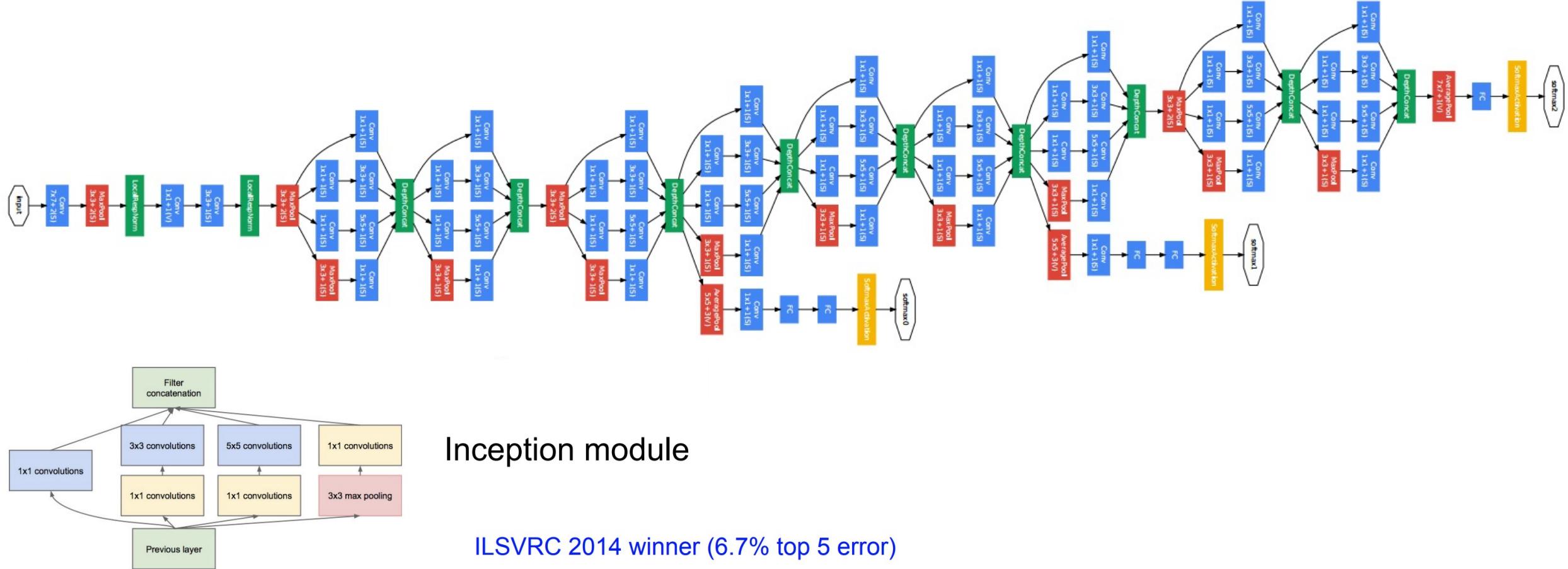
GoogLeNet



Inception module

ILSVRC 2014 winner (6.7% top 5 error)

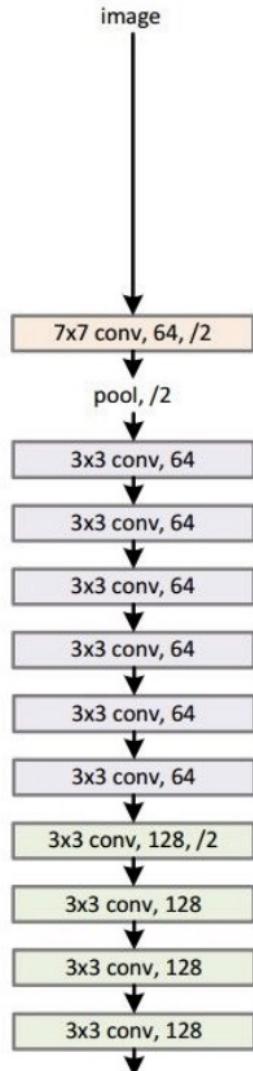
GoogLeNet



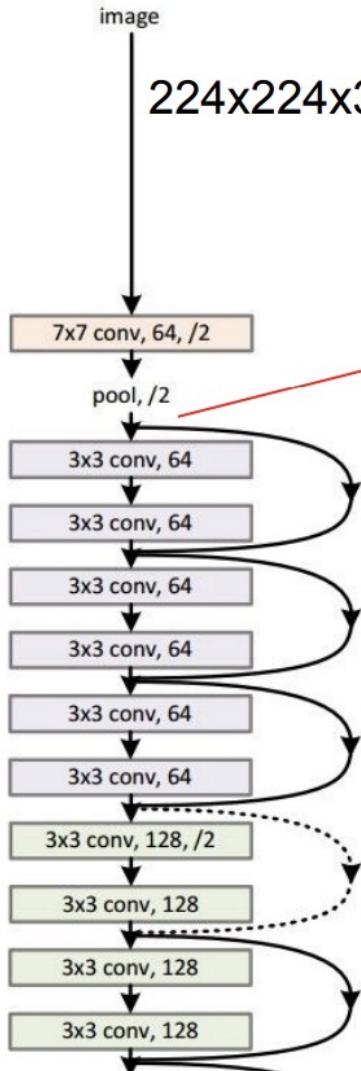
[Szegedy et al., 2014]

ResNet

34-layer plain

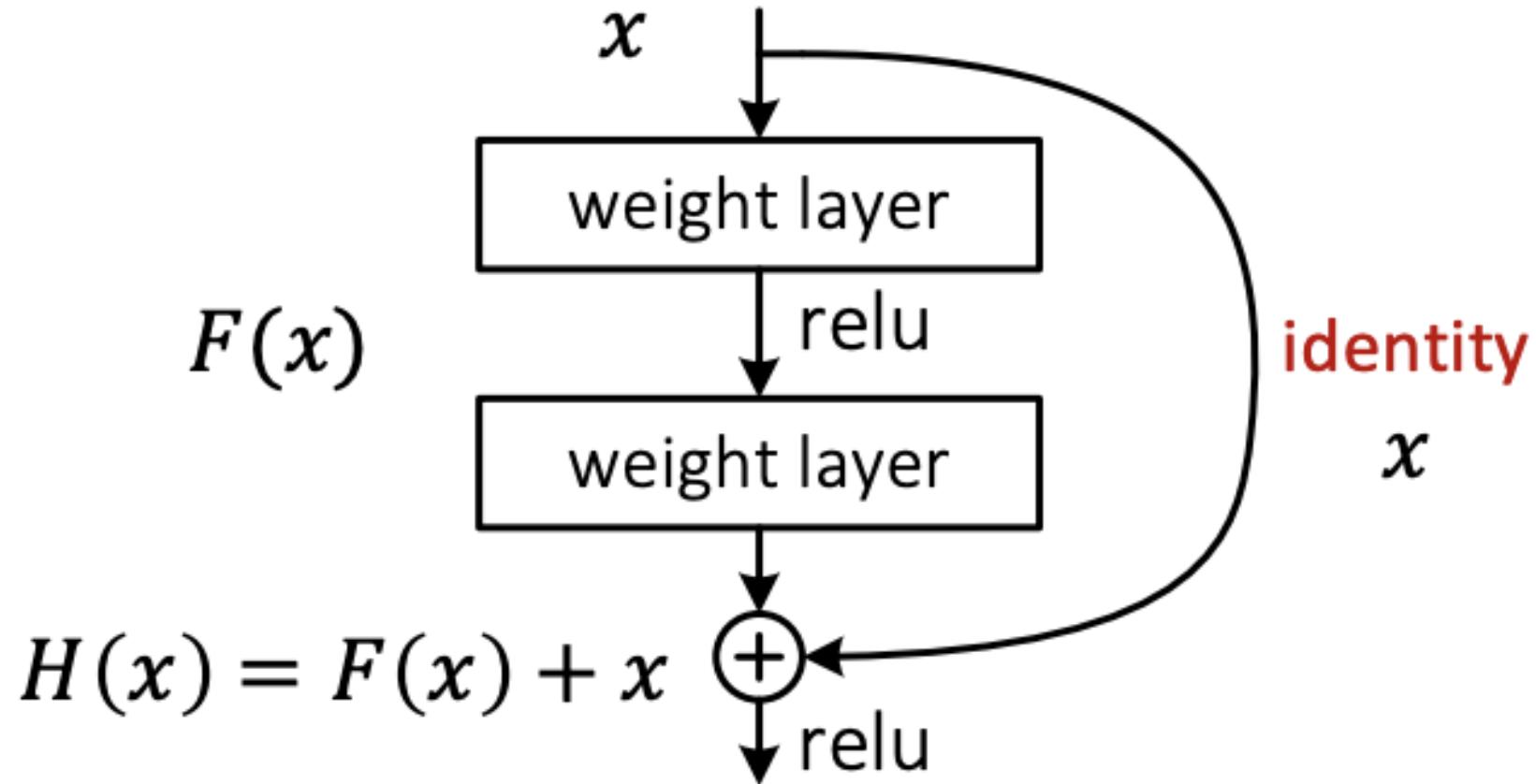


34-layer residual

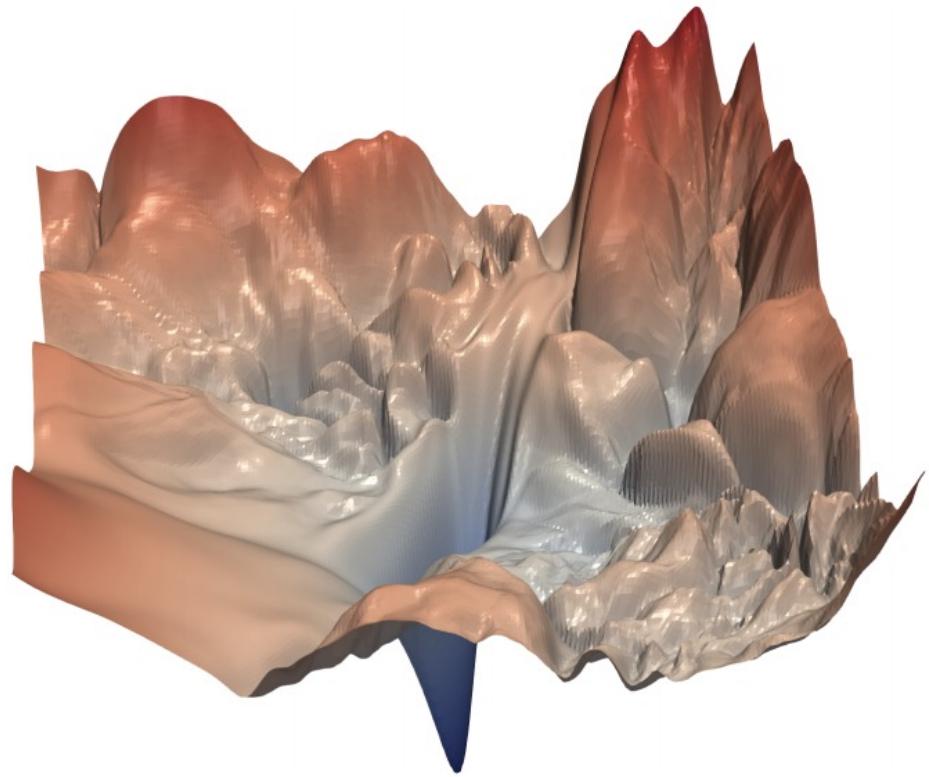


spatial dimension
only 56×56 !

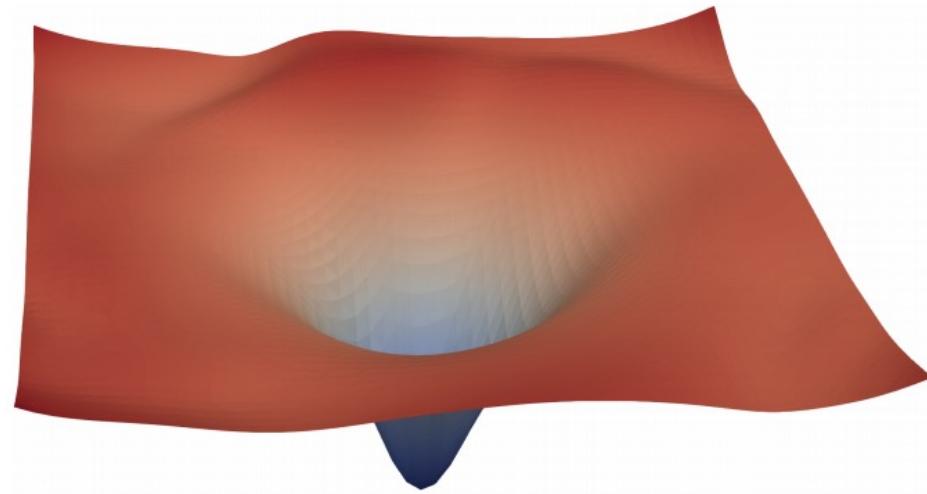
Residual Block



Residual Block



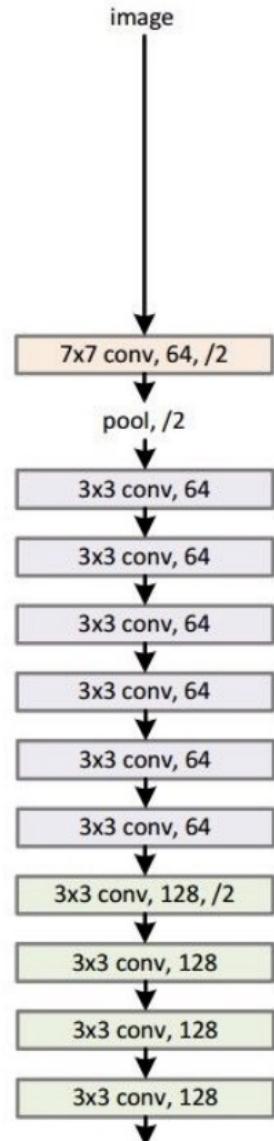
(a) without skip connections



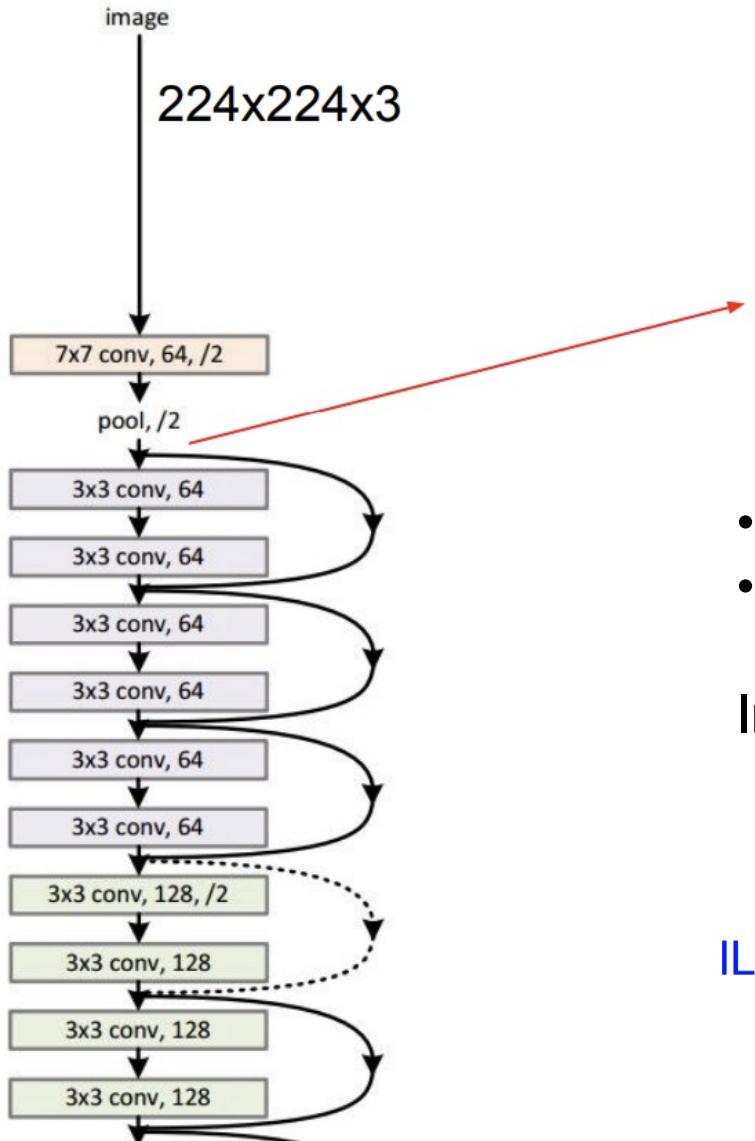
(b) with skip connections

Figure 1: The loss surfaces of ResNet-56 with/without skip connections. The proposed filter normalization scheme is used to enable comparisons of sharpness/flatness between the two figures.

34-layer plain



34-layer residual



spatial dimension
only 56×56 !

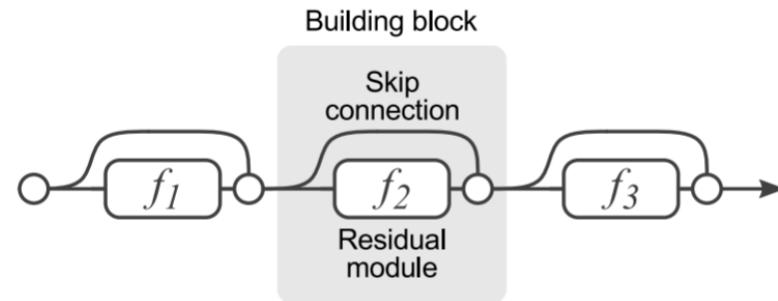
- Batch Normalization after every conv layer
- No dropout used

ImageNet 2015 winner

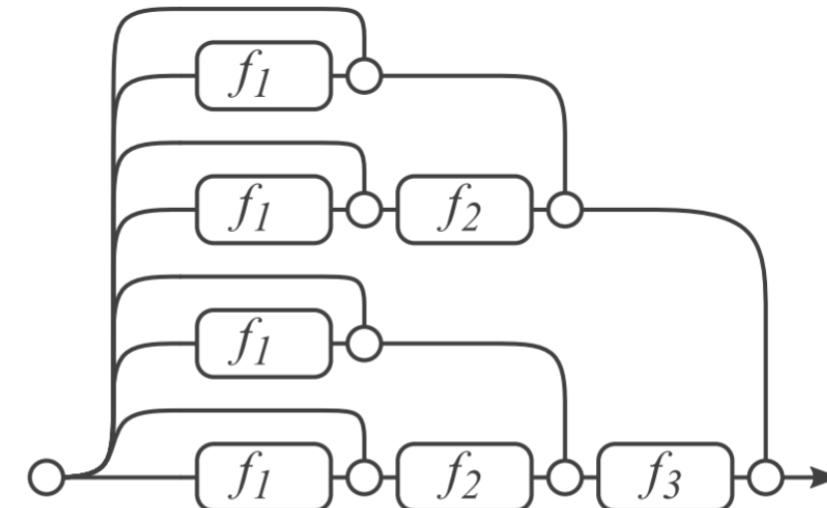
ILSVRC 2015 winner (3.6% top 5 error)

ResNet is implicit ensemble

Residual Networks Behave Like Ensembles of Relatively Shallow Networks



(a) Conventional 3-block residual network



(b) Unraveled view of (a)

Paper: <https://arxiv.org/pdf/1605.06431.pdf>

ResNeXt

Same as ResNet, but scalable

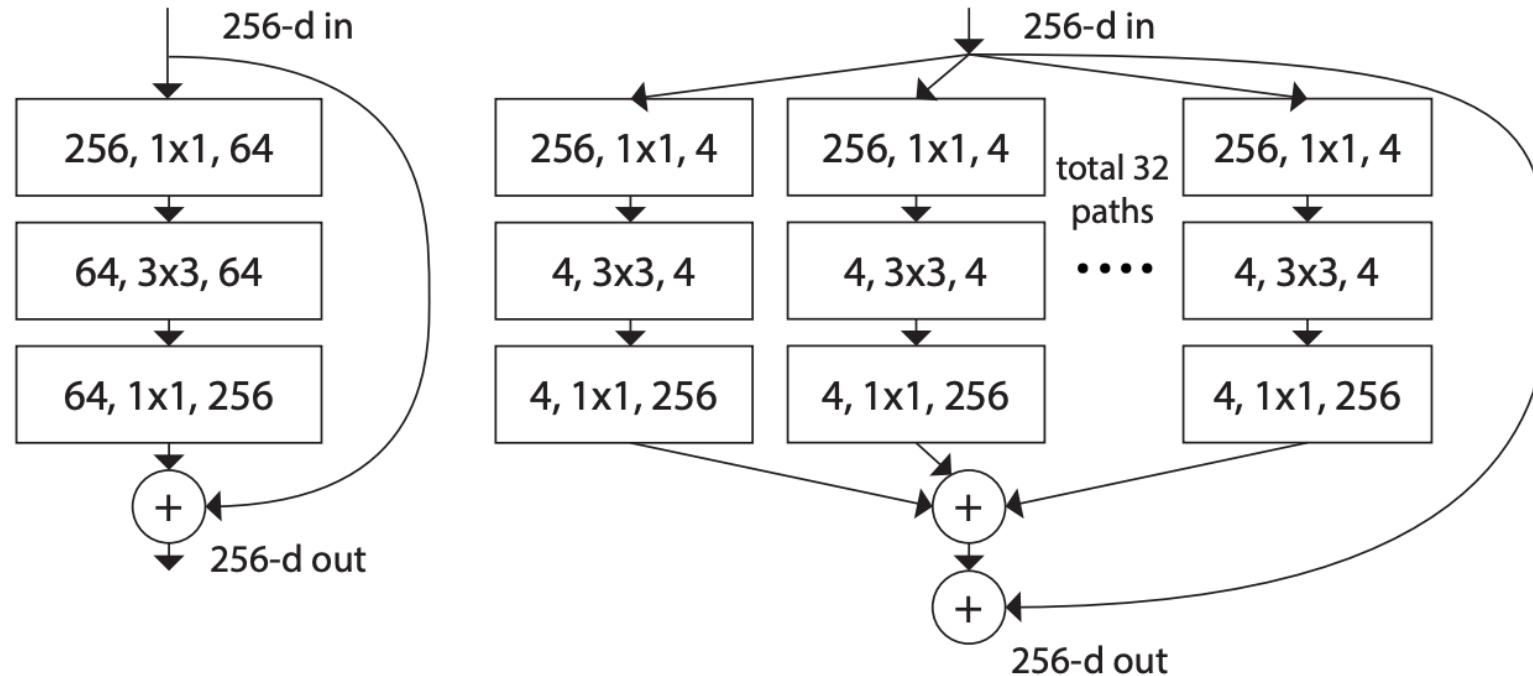


Figure 1. **Left:** A block of ResNet [14]. **Right:** A block of ResNeXt with cardinality = 32, with roughly the same complexity. A layer is shown as (# in channels, filter size, # out channels).

SENet

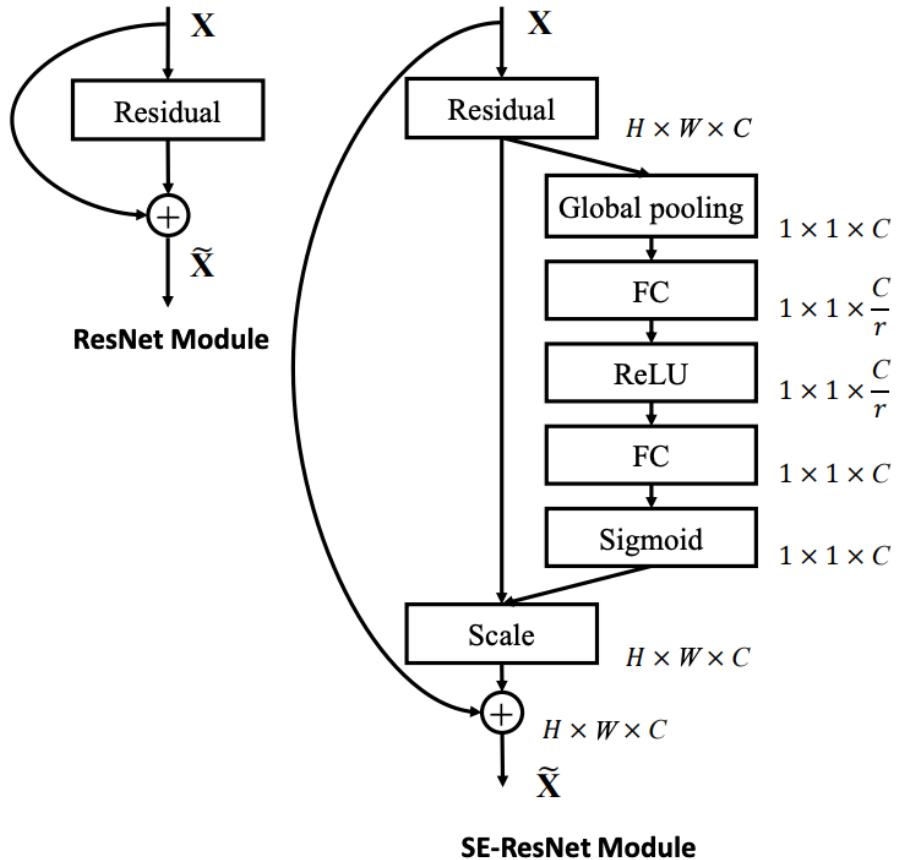
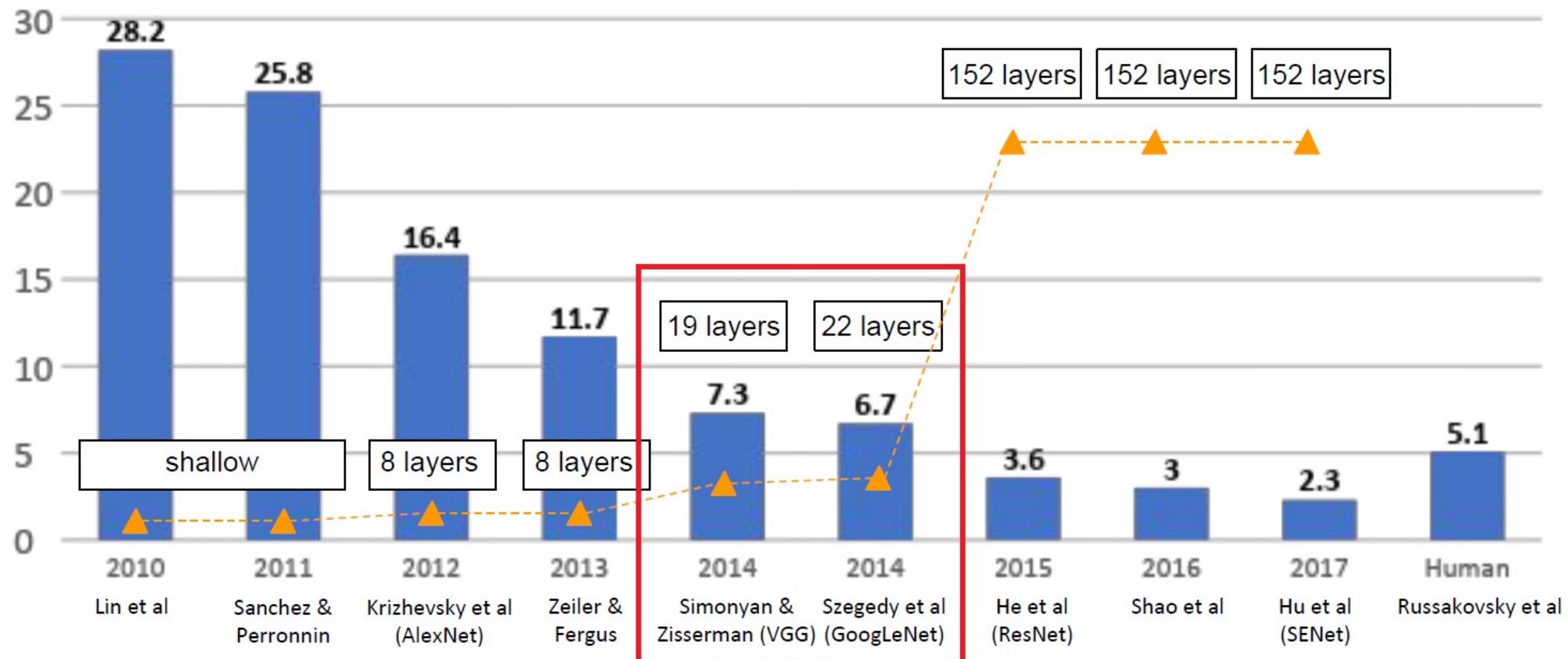


Fig. 3. The schema of the original Residual module (left) and the SE-ResNet module (right).

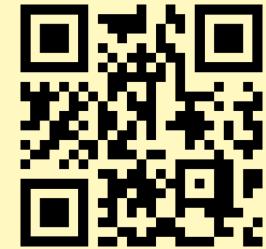
ImageNet summary

ImageNet Large Scale Visual Recognition Challenge (ILSVRC) winners



Канал стажировок Яндекса

Спасибо за внимание



Доп. материалы



YANDEX