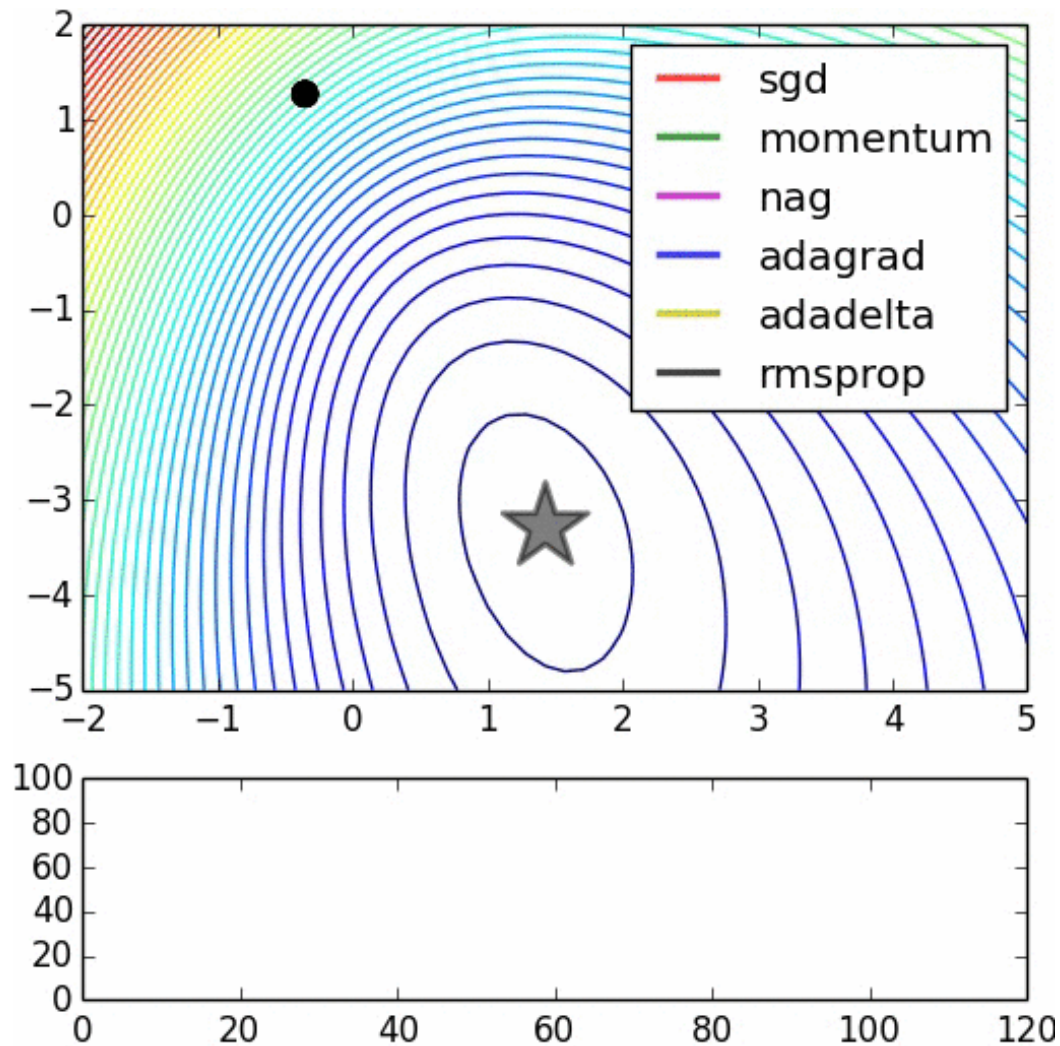




CSE 176 Introduction to Machine Learning

Lecture 12: Convolutional Neural Network

From last lecture: Optimization methods



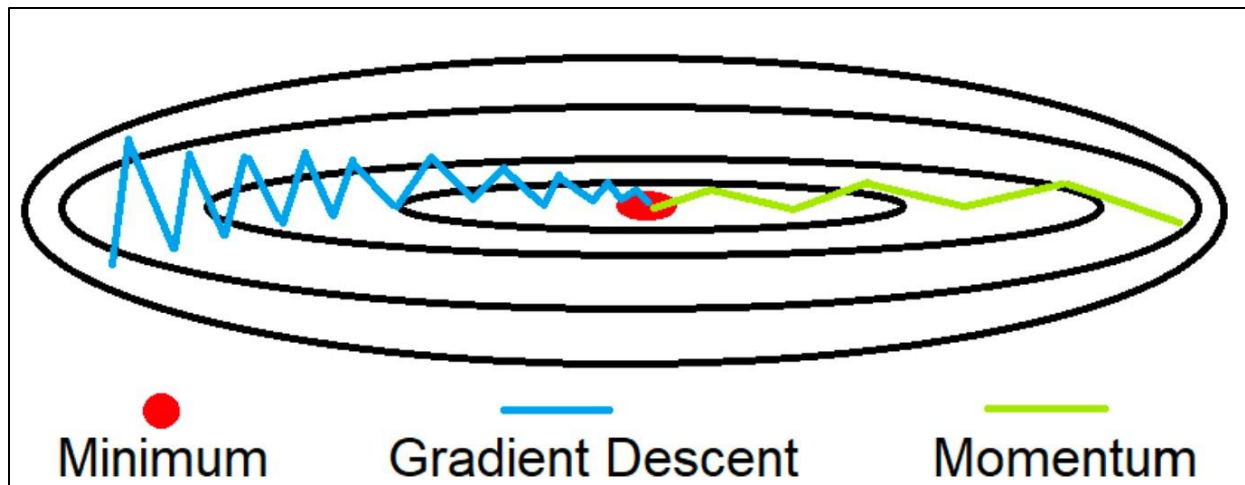
Quiz

Which of the following optimization algorithm(s) uses adaptive learning rate for each network parameter? Select all that apply.

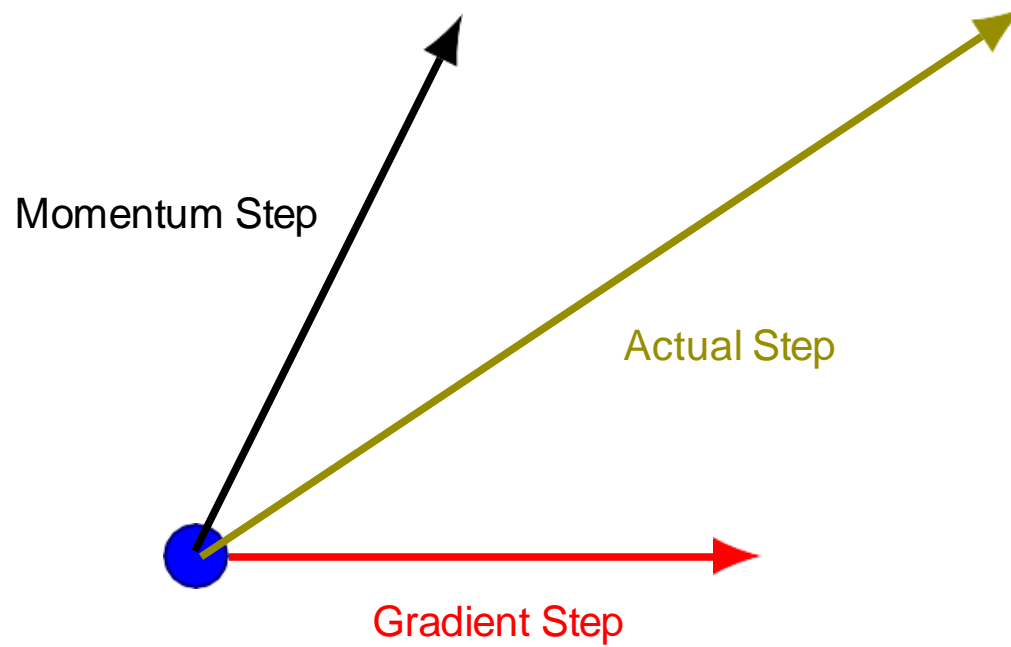
- ☐ Stochastic Gradient Descent (SGD)
- ☐ SGD with momentum
- ☐ AdaGrad

Gradient Descent with Momentum

- Momentum reduces oscillation in w_2 direction



Momentum



SGD:

Compute gradient estimate: $\hat{\mathbf{g}} \leftarrow +\frac{1}{m} \nabla_{\boldsymbol{\theta}} \sum_i L(f(\mathbf{x}^{(i)}; \boldsymbol{\theta}), \mathbf{y}^{(i)})$

Apply update: $\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} - \epsilon \hat{\mathbf{g}}$

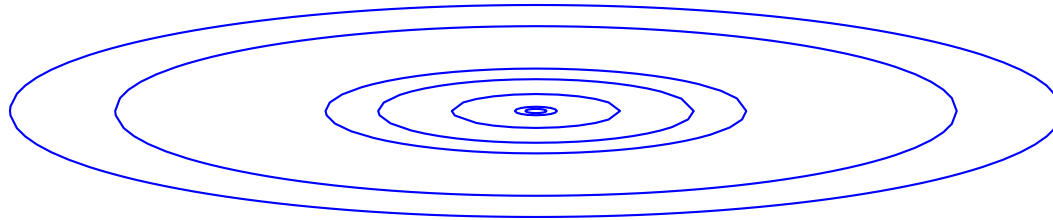
SGD with momentum:

Compute gradient estimate: $\mathbf{g} \leftarrow \frac{1}{m} \nabla_{\boldsymbol{\theta}} \sum_i L(f(\mathbf{x}^{(i)}; \boldsymbol{\theta}), \mathbf{y}^{(i)})$

Compute velocity update: $\mathbf{v} \leftarrow \alpha \mathbf{v} - \epsilon \mathbf{g}$

Apply update: $\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} + \mathbf{v}$

Motivation for Adaptive Learning Rate for each Parameter



Harder

AdaGrad

Require: Global learning rate ϵ

Require: Initial parameter θ

Require: Small constant δ , perhaps 10^{-7} , for numerical stability

Initialize gradient accumulation variable $\mathbf{r} = \mathbf{0}$

while stopping criterion not met **do**

Sample a minibatch of m examples from the training set $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\}$ with corresponding targets $\mathbf{y}^{(i)}$.

Compute gradient: $\mathbf{g} \leftarrow \frac{1}{m} \nabla_{\theta} \sum_i L(f(\mathbf{x}^{(i)}; \theta), \mathbf{y}^{(i)})$

Accumulate squared gradient: $\mathbf{r} \leftarrow \mathbf{r} + \mathbf{g} \odot \mathbf{g}$

Compute update: $\Delta\theta \leftarrow -\frac{\epsilon}{\delta + \sqrt{\mathbf{r}}} \odot \mathbf{g}$. (Division and square root applied element-wise)

Apply update: $\theta \leftarrow \theta + \Delta\theta$

end while

Topics today

☐ Convolutional Neural Networks

- ☐ Convolution layer

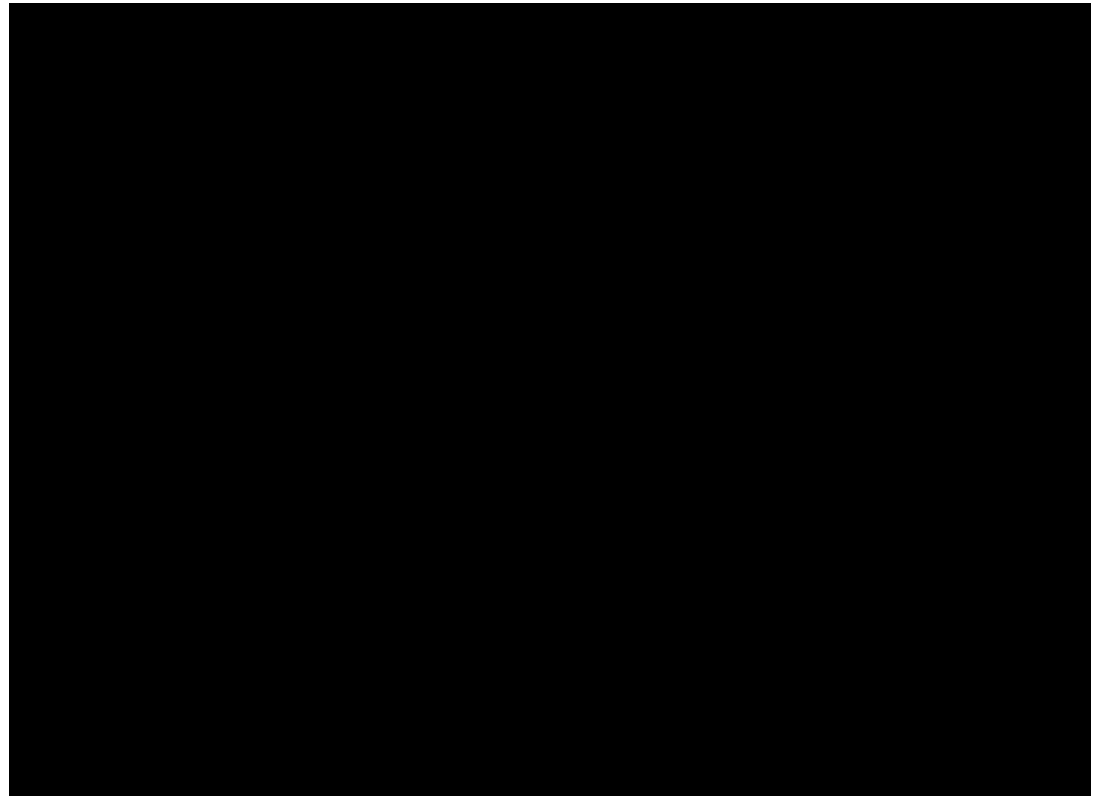
- ☐ Pooling layer

- ☐ Fully connected layers

First CNN architectures for classification

- **first CNNs** (1982-89) *Neocognitron: A new algorithm for pattern recognition tolerant of deformations and shifts in position*
(a.k.a. *convNets*)
K. Fukushima, S. Miyake - Pattern Recognition 1982

- **LeNet** (1998)



https://youtu.be/FwFduRA_L6Q

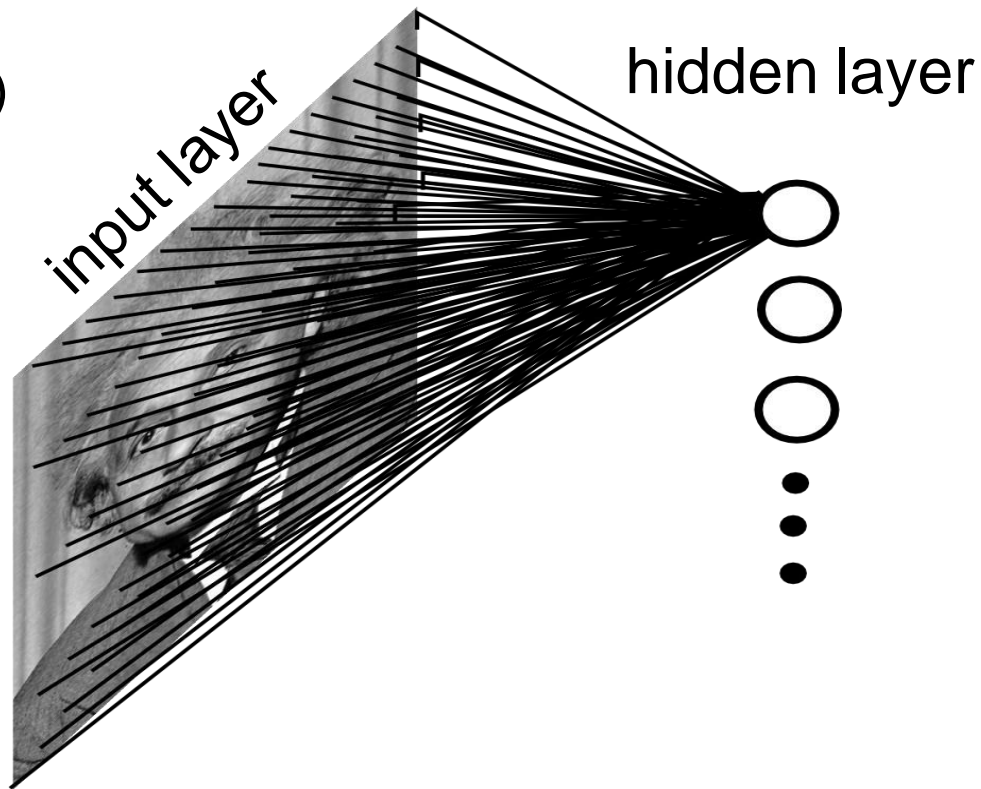
Handwritten digit recognition with a back-propagation network
Y. LeCun et al - NIPS 1989

Convolutional Network: Motivation

Consider a **fully connected network** (most weights $W[i,j] \neq 0$)

Example: 200 by 200 image,
 4×10^4 connections to one
hidden unit

For 10^5 hidden units $\rightarrow 4 \times 10^9$
connections



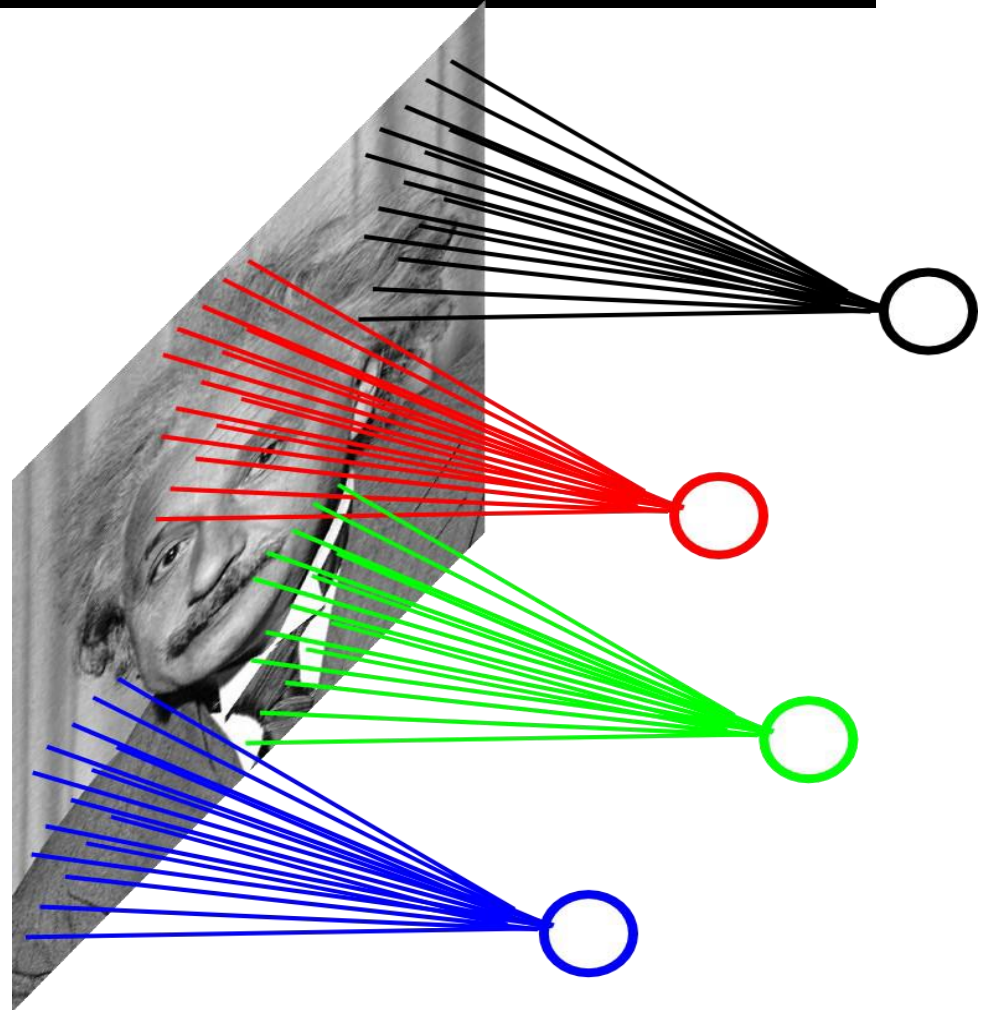
Motivation for local connections

Connect only pixels in a local patch, say 10×10

For 200 by 200 image, 10^2 connections to one hidden unit

For 10^5 hidden units $\rightarrow 10^7$ connections

- contrast with 4×10^9 for fully connected layer
- factor of 400 decrease

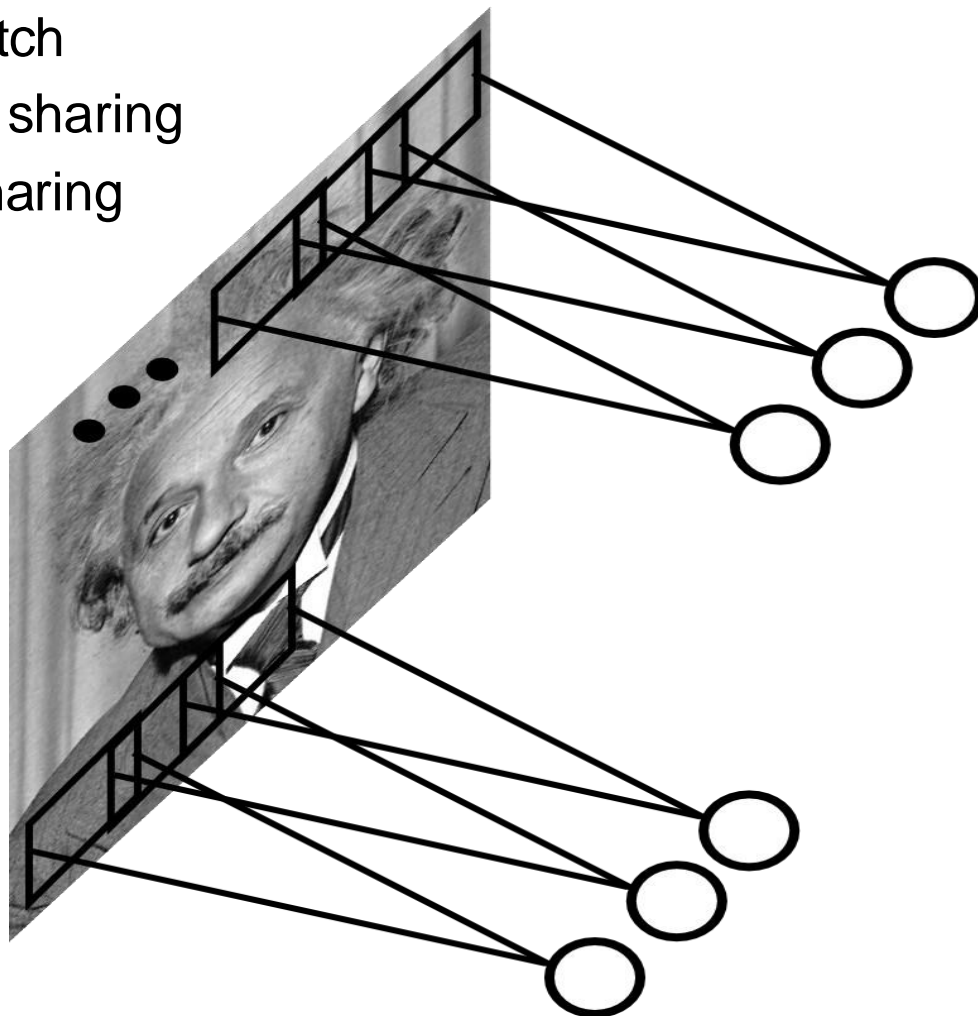


Motivation for Weight Sharing

Much fewer parameters to learn

For 10^5 hidden units and 10×10 patch

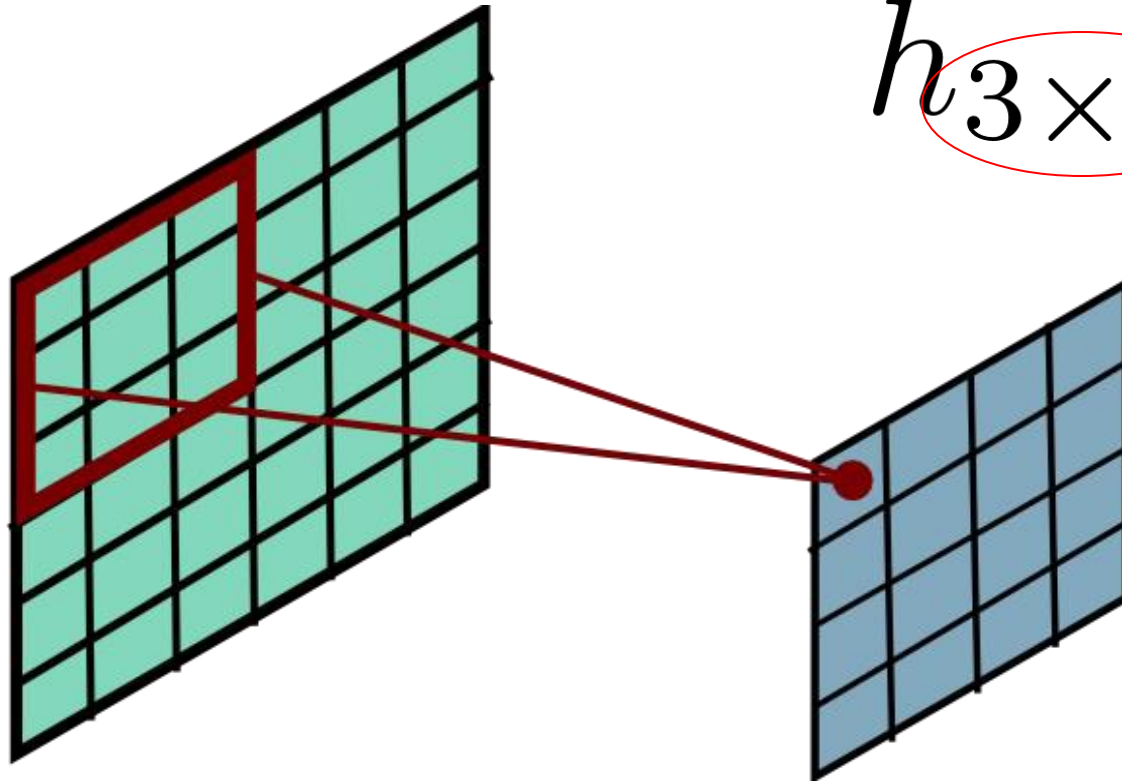
- 10^7 parameters to learn without sharing
- 10^2 parameters to learn with sharing



Convolutional Layer

convolution kernel

h 3×3 size



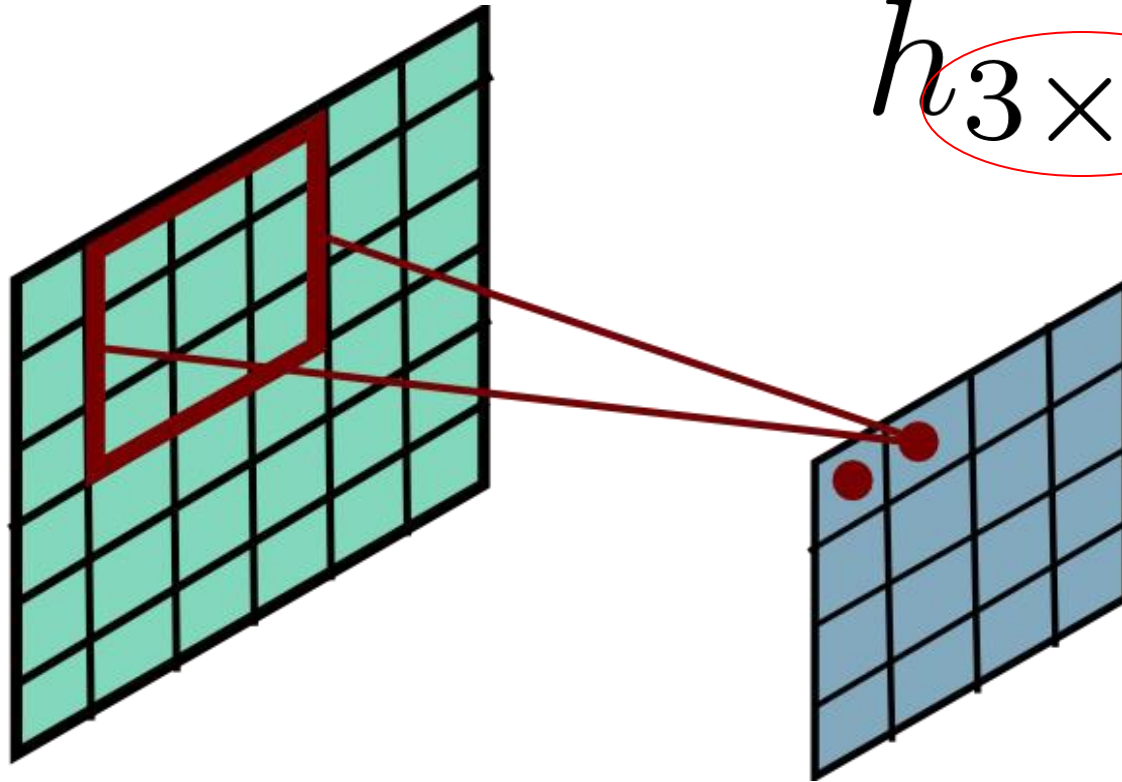
input

output

Convolutional Layer

convolution kernel

h 3×3 size



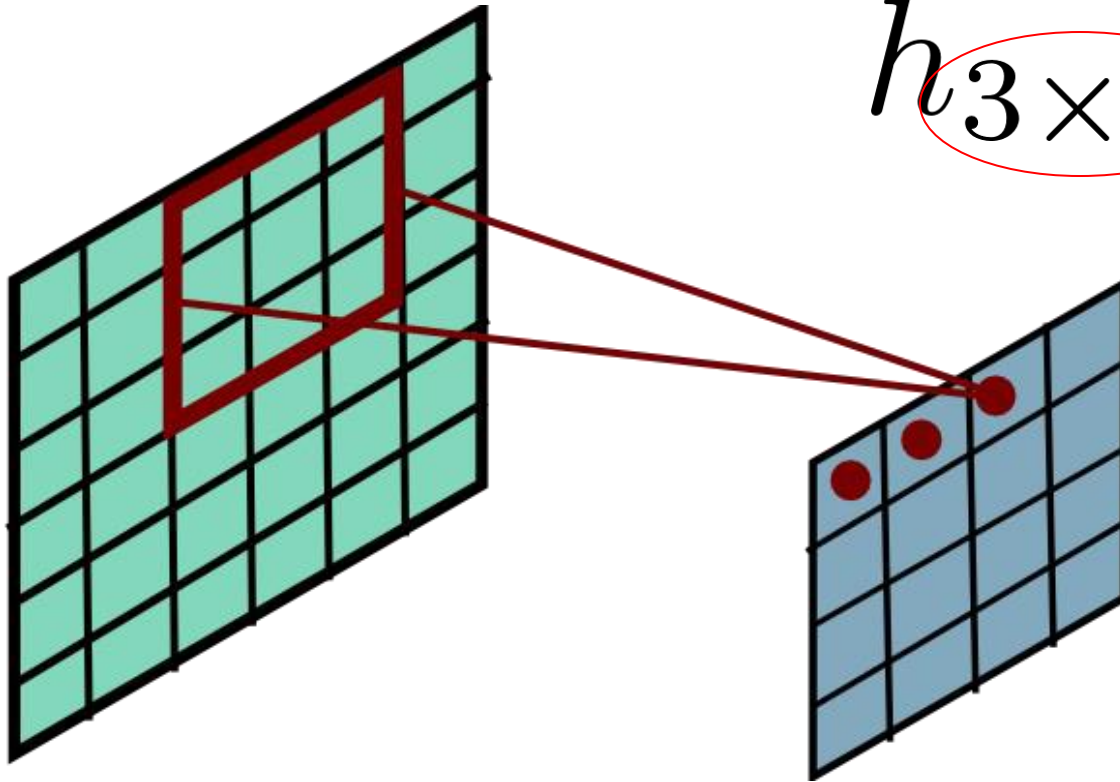
input

output

Convolutional Layer

convolution kernel

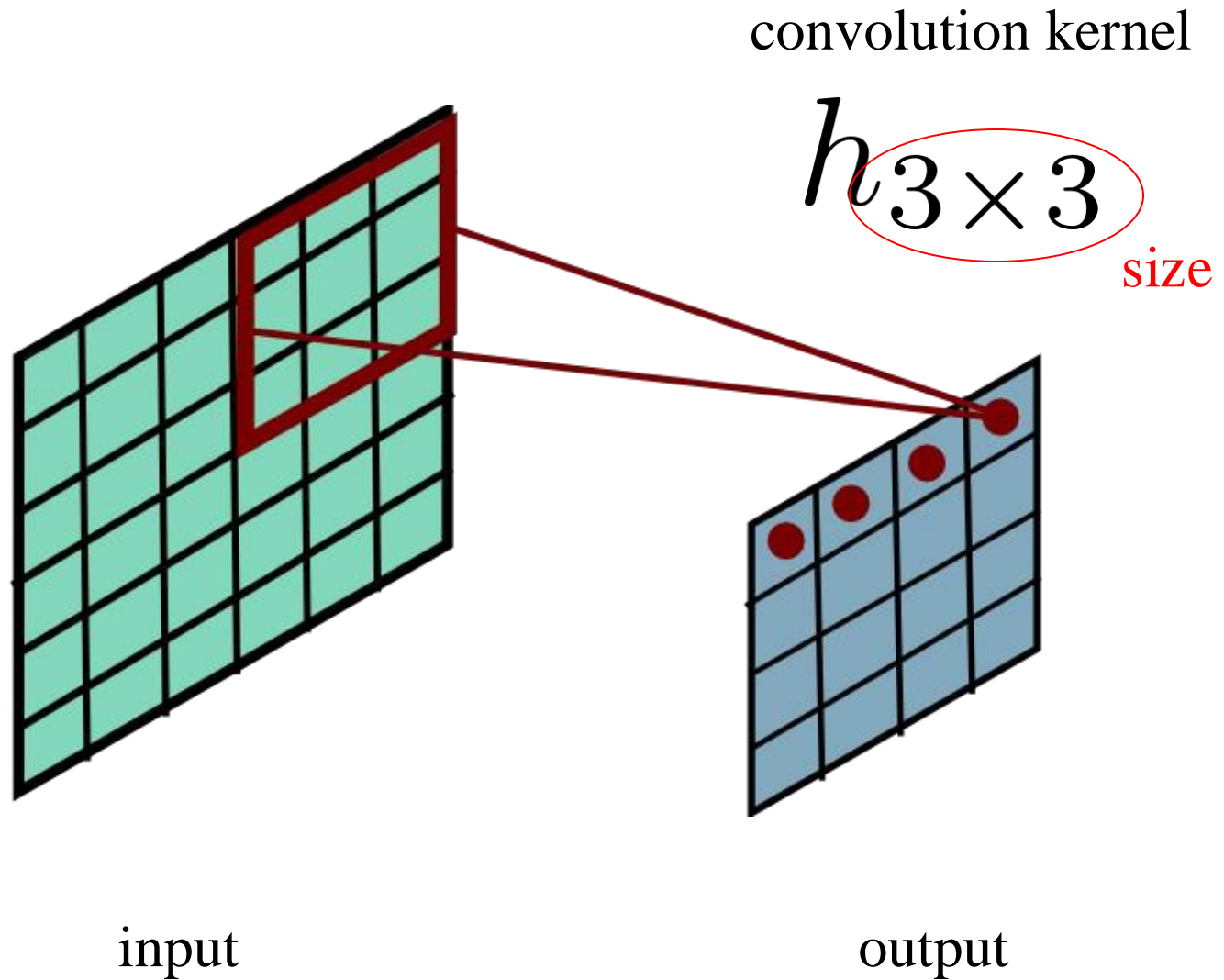
$h_{3 \times 3}$
size



input

output

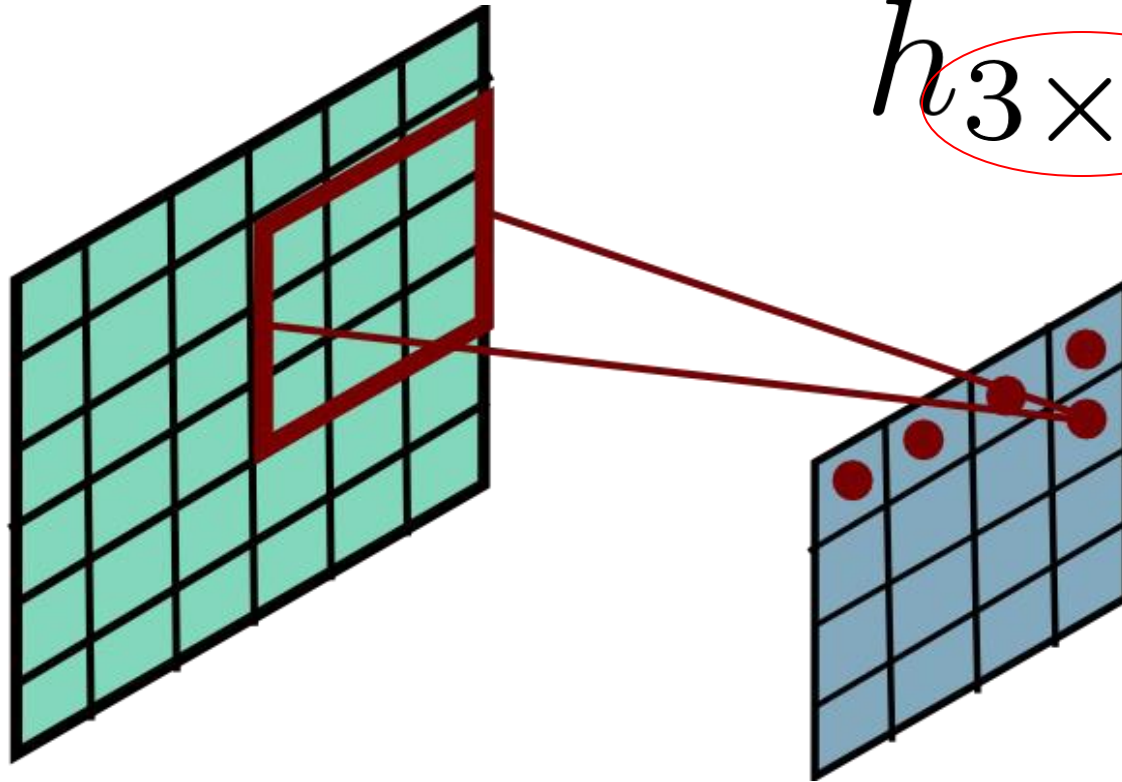
Convolutional Layer



Convolutional Layer

convolution kernel

h 3×3 size



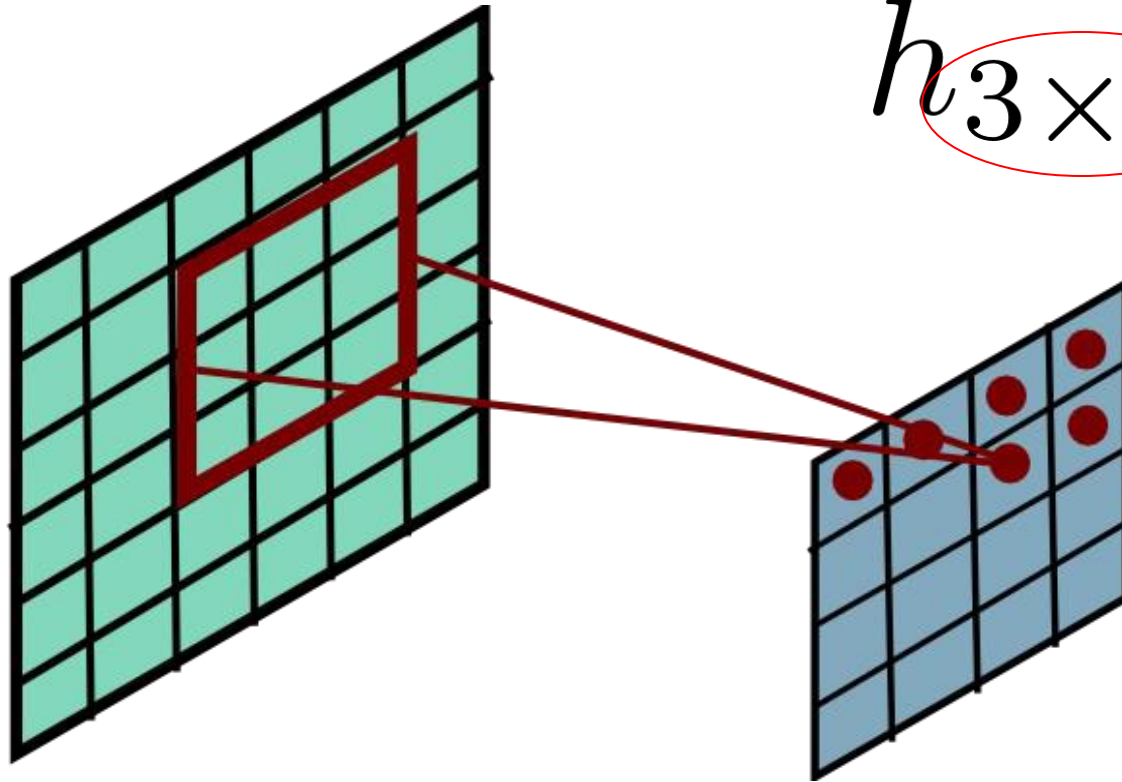
input

output

Convolutional Layer

convolution kernel

h 3×3 size



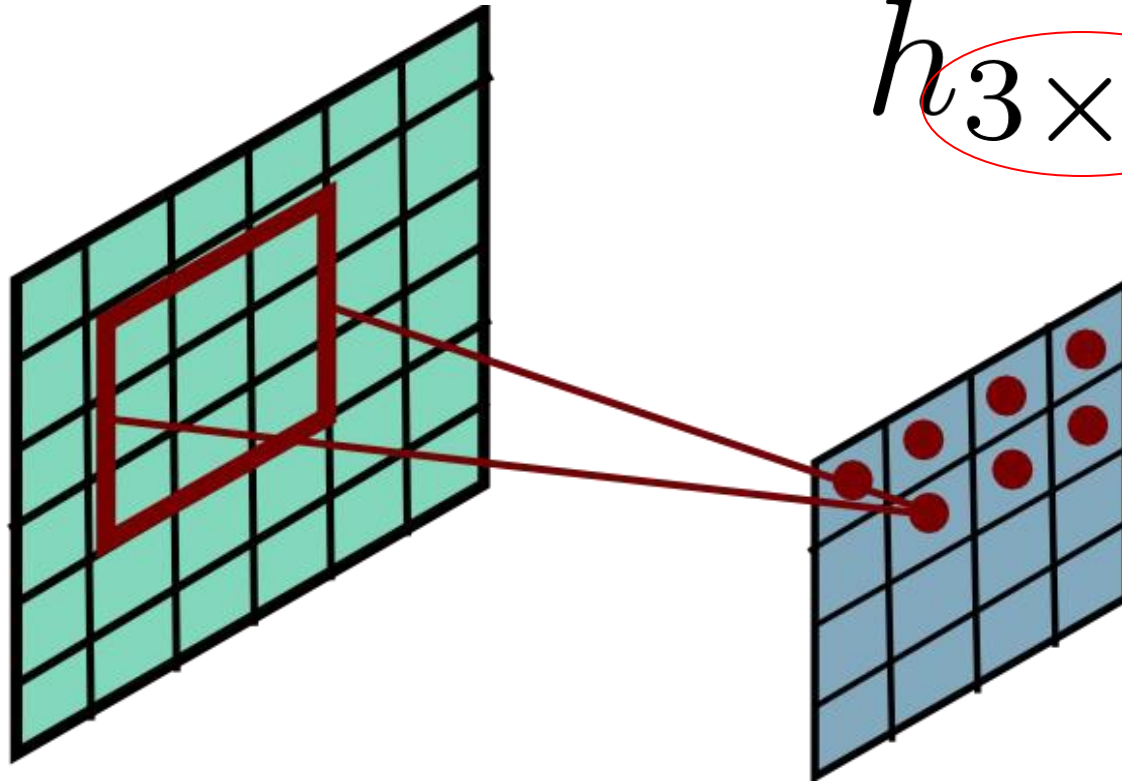
input

output

Convolutional Layer

convolution kernel

h 3×3 size



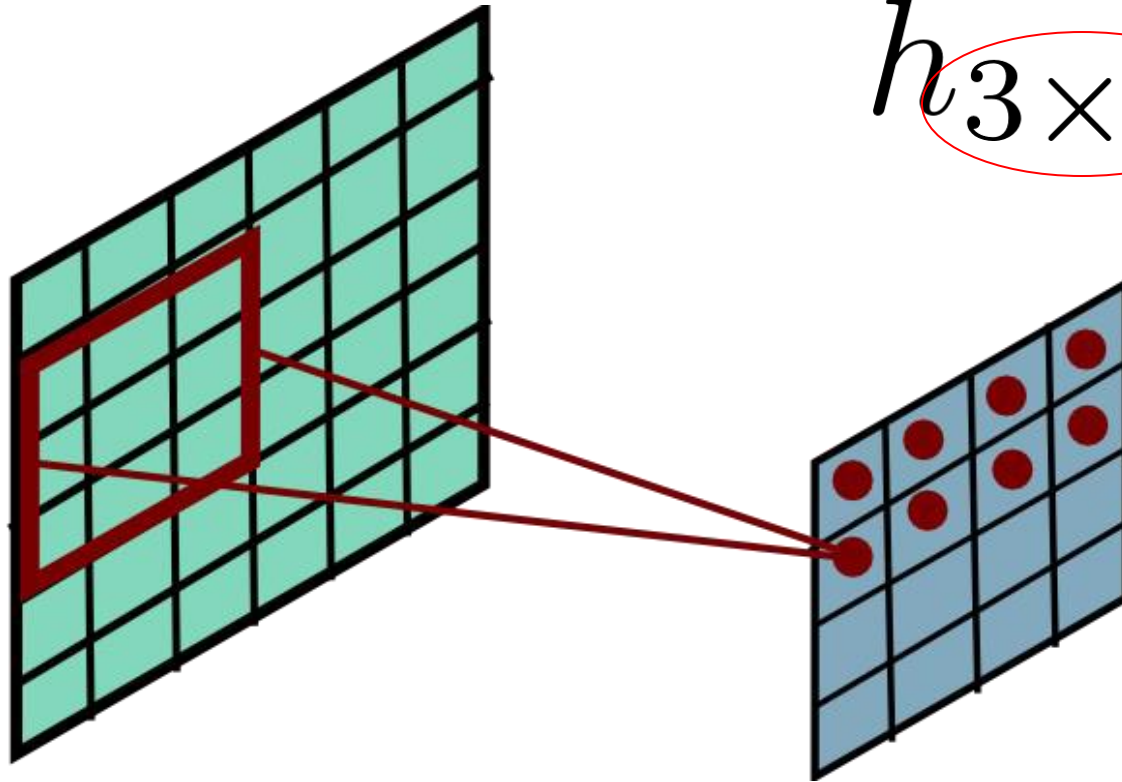
input

output

Convolutional Layer

convolution kernel

h 3×3 size



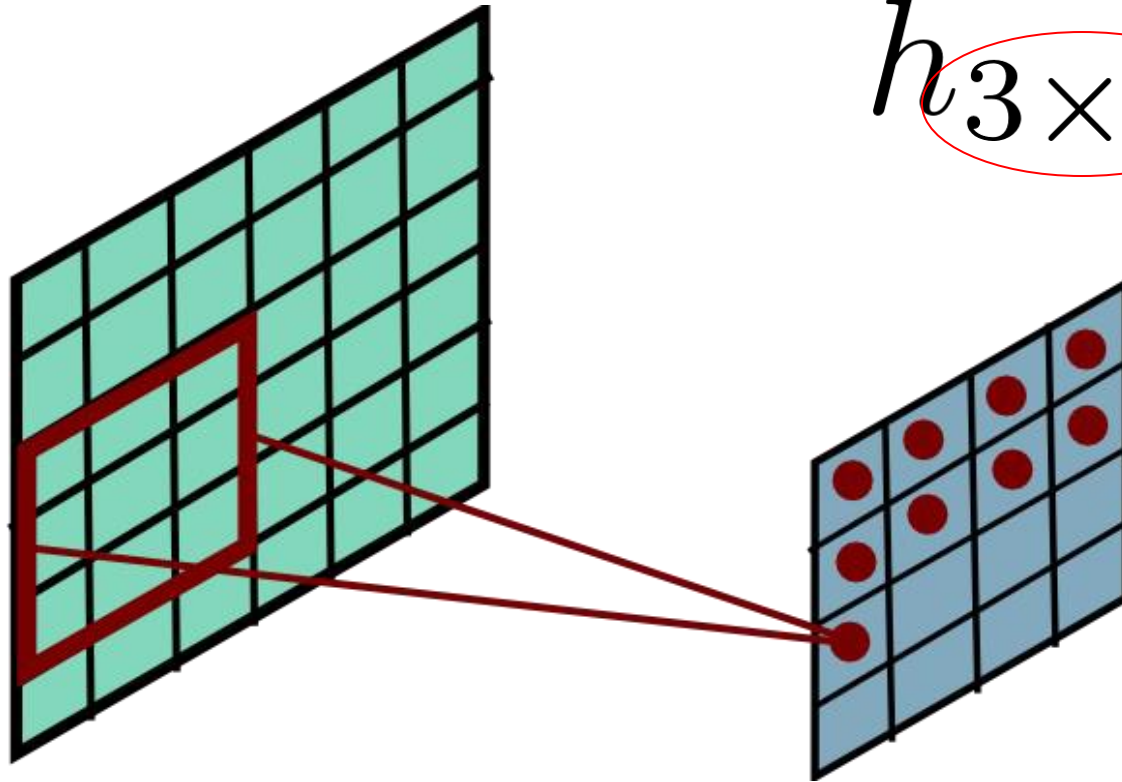
input

output

Convolutional Layer

convolution kernel

$h_{3 \times 3}$ size



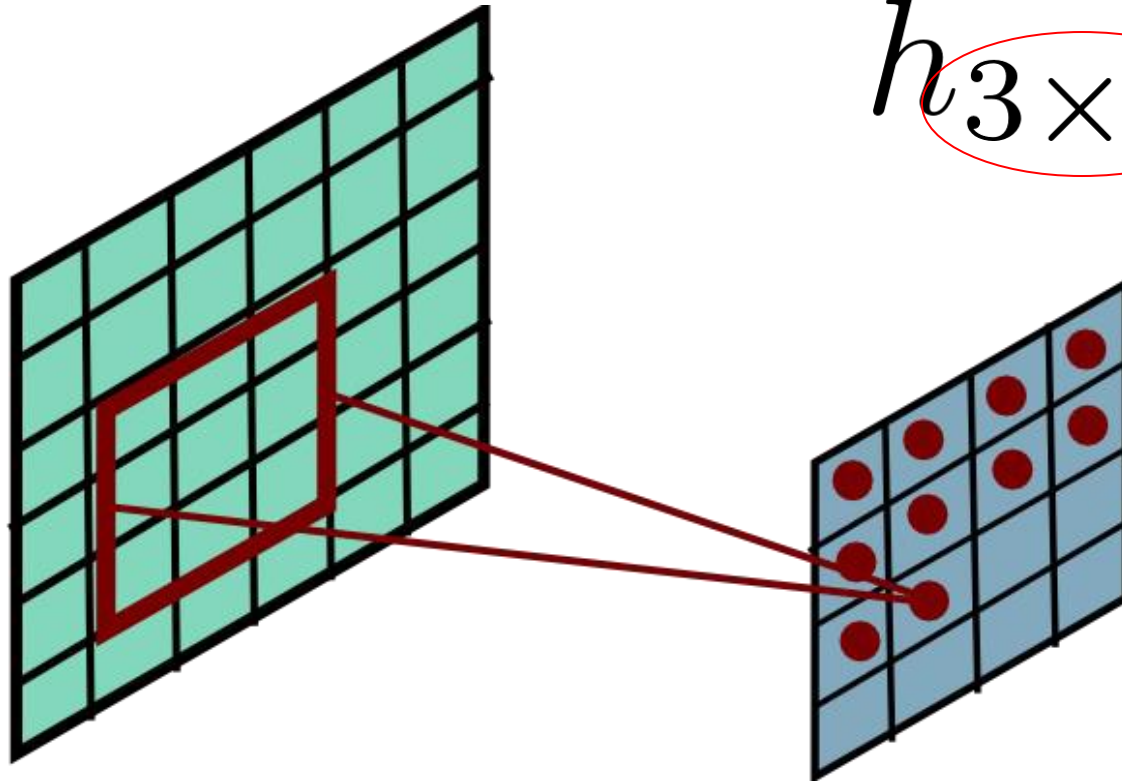
input

output

Convolutional Layer

convolution kernel

h 3×3 size



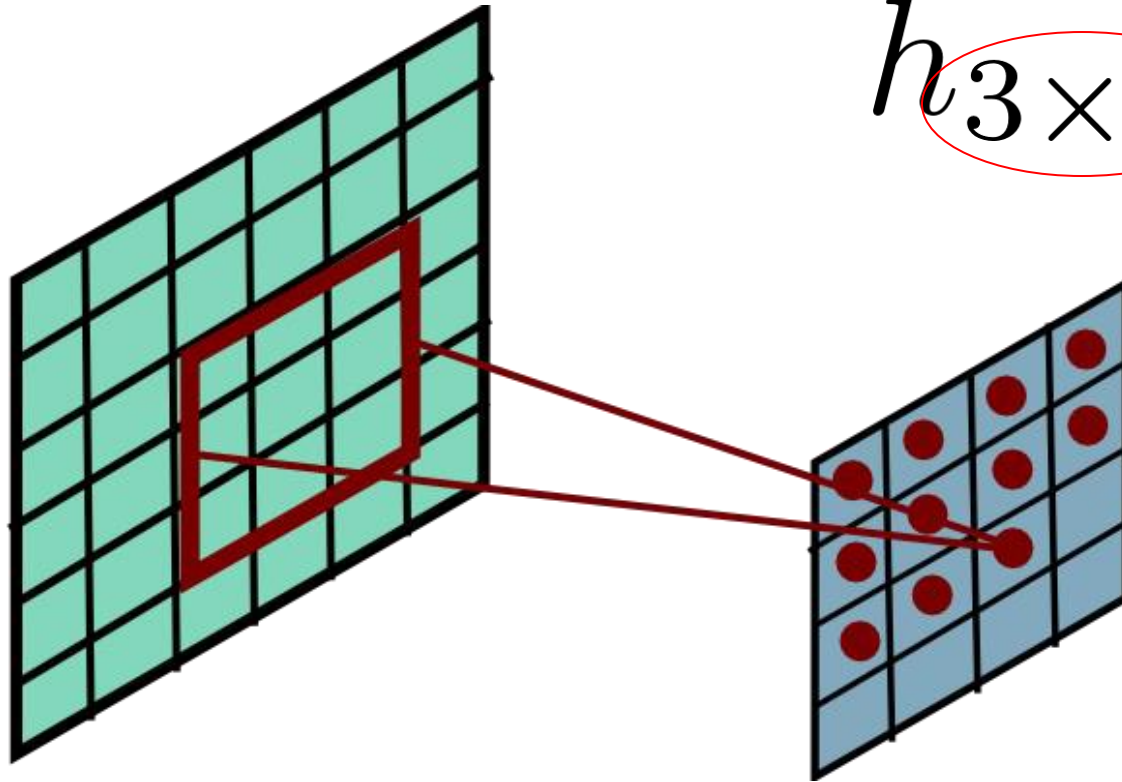
input

output

Convolutional Layer

convolution kernel

h 3×3 size



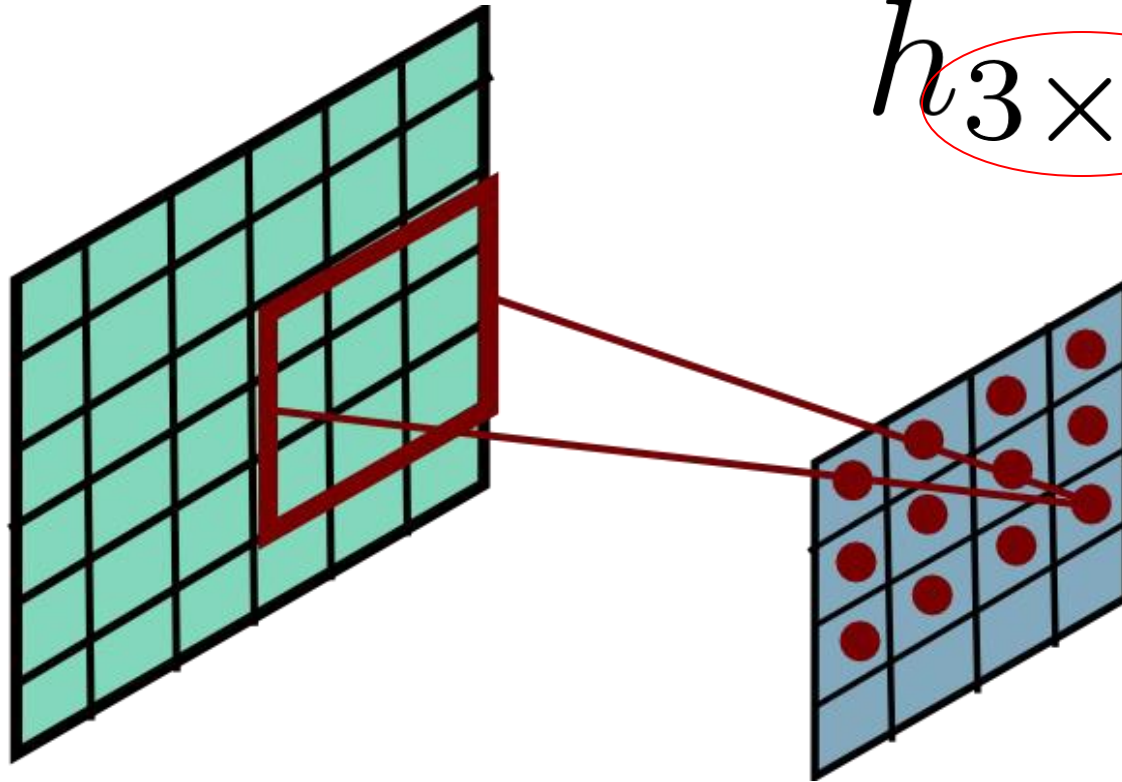
input

output

Convolutional Layer

convolution kernel

h 3×3 size



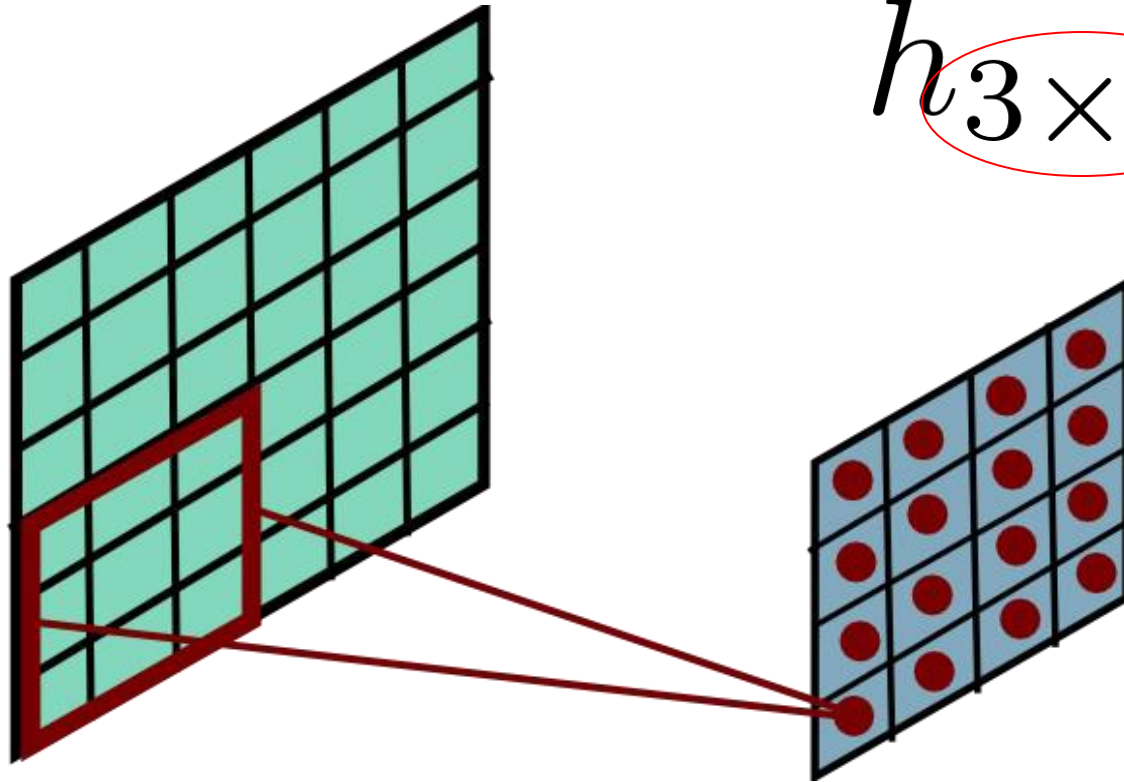
input

output

Convolutional Layer

convolution kernel

h 3×3 size



input

output

2D Convolution

A 2D image $f[i,j]$ can be filtered by a **2D kernel** $h[u,v]$ to produce an output image $g[i,j]$:

$$g[i, j] = \sum_{u=-k}^k \sum_{v=-k}^k h[u, v] \cdot f[i + u, j + v]$$

This is called a **convolution** operation and written:

$$g = h \circ f$$

h is called “**kernel**” or “**mask**” or “**filter**” which representing a given “window function”

Mean filtering

0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0
0	0	0	90	90	90	90	90	0	0
0	0	0	90	90	90	90	90	0	0
0	0	0	90	90	90	90	90	0	0
0	0	0	90	0	90	90	90	0	0
0	0	0	90	90	90	90	90	0	0
0	0	0	0	0	0	0	0	0	0
0	0	90	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0

$f[x, y]$

		10							

$g[x, y]$

Mean filtering

[illegible]
$$f[x, y]$$

A 10x10 grid with a red square containing the number 80 at the intersection of the 4th column and 6th row, and the number 10 at the intersection of the 3rd column and 8th row.

$$g[x, y]$$

Mean filtering

side effect of mean filtering: **blurring**

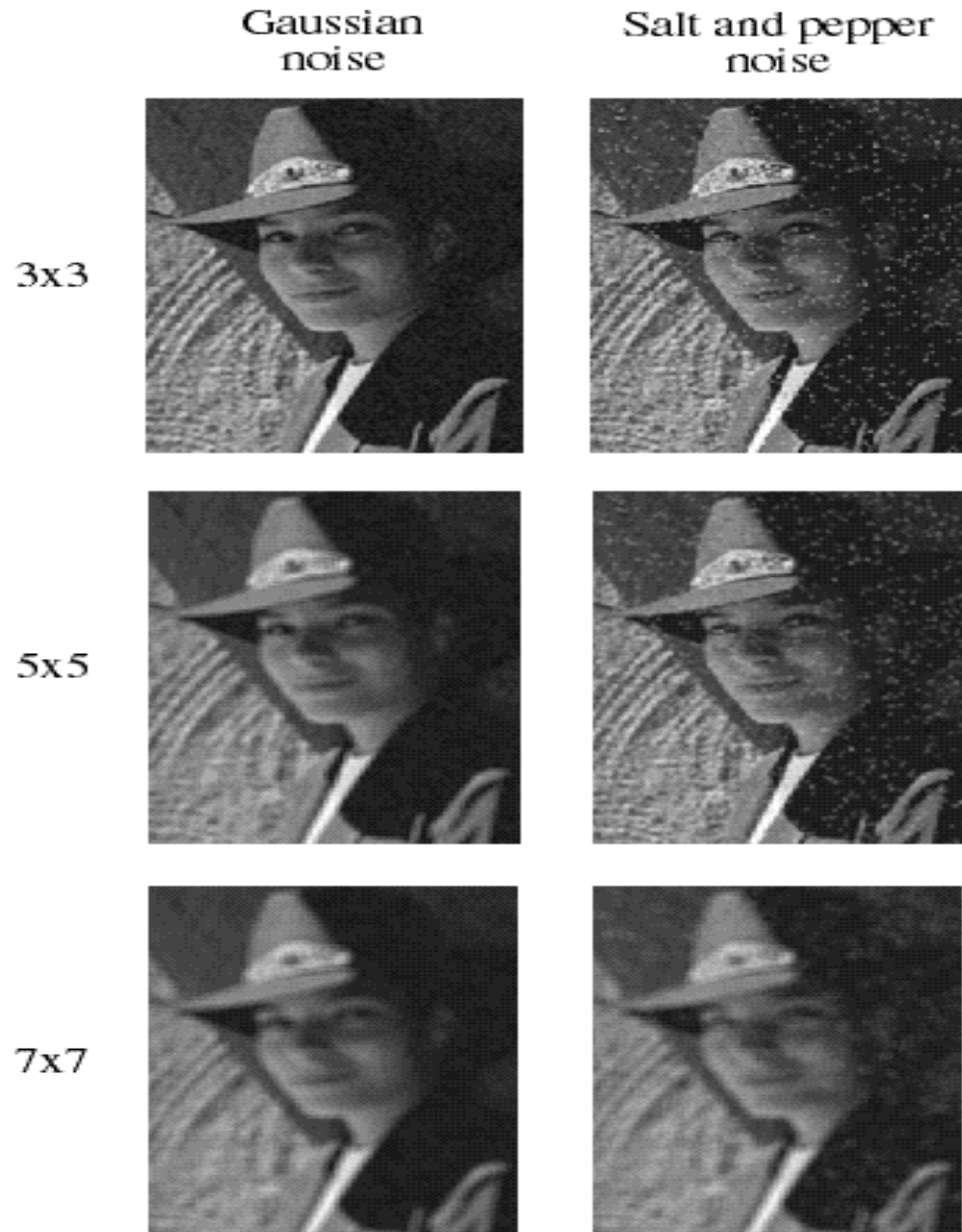
0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0
0	0	0	90	90	90	90	90	0	0
0	0	0	90	90	90	90	90	0	0
0	0	0	90	90	90	90	90	0	0
0	0	0	90	0	90	90	90	0	0
0	0	0	90	90	90	90	90	0	0
0	0	0	0	0	0	0	0	0	0
0	0	90	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0

$f[x, y]$

	0	10	20	30	30	30	20	10	
	0	20	40	60	60	60	40	20	
	0	30	60	90	90	90	60	30	
	0	30	50	80	80	90	60	30	
	0	30	50	80	80	90	60	30	
	0	20	30	50	50	60	40	20	
	10	20	30	30	30	30	20	10	
	10	10	10	0	0	0	0	0	

$g[x, y]$

Mean filtering



Mean kernel

□ What's the kernel for a 3x3 mean filter?

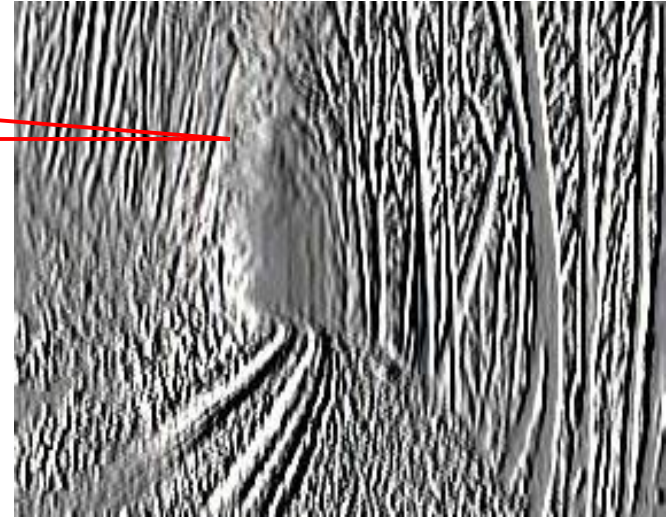
0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0
0	0	0	90	90	90	90	90	0	0
0	0	0	90	90	90	90	90	0	0
0	0	0	90	90	90	90	90	0	0
0	0	0	90	0	90	90	90	0	0
0	0	0	90	90	90	90	90	0	0
0	0	0	0	0	0	0	0	0	0
0	0	90	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0

$$\frac{1}{9} \cdot \begin{bmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{bmatrix}$$

Image Filtering by Other Kernel



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



2D filtering for noise reduction

- ❑ Common types of noise:
 - ❑ **Salt and pepper noise:** random occurrences of black and white pixels
 - ❑ **Gaussian noise:** variations in intensity drawn from a Gaussian normal distribution



Original



Salt and pepper noise



Impulse noise



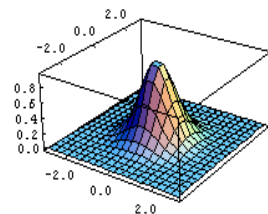
Gaussian noise

Gaussian filtering

□ A Gaussian kernel gives less weight to pixels further from the center

0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0
0	0	0	90	90	90	90	90	0	0
0	0	0	90	90	90	90	90	0	0
0	0	0	90	90	90	90	90	0	0
0	0	0	90	0	90	90	90	0	0
0	0	0	90	90	90	90	90	0	0
0	0	0	0	0	0	0	0	0	0
0	0	90	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0

$$\frac{1}{16} \cdot \begin{bmatrix} 1 & 2 & 1 \\ 2 & 4 & 2 \\ 1 & 2 & 1 \end{bmatrix}$$



discrete approximation of
a Gaussian (density) function

$$h(u, v) = \frac{1}{2\pi\sigma^2} e^{-\frac{u^2+v^2}{\sigma^2}}$$

Median filter

- ☐ A **Median Filter** operates over a window by selecting the median intensity in the window.
- ☐ What advantage does a median filter have over a mean filter?
- ☐ Is a median filter a kind of convolution?
 - ☐ - No, median filter is non-linear

Comparison: salt and pepper noise

3x3



5x5



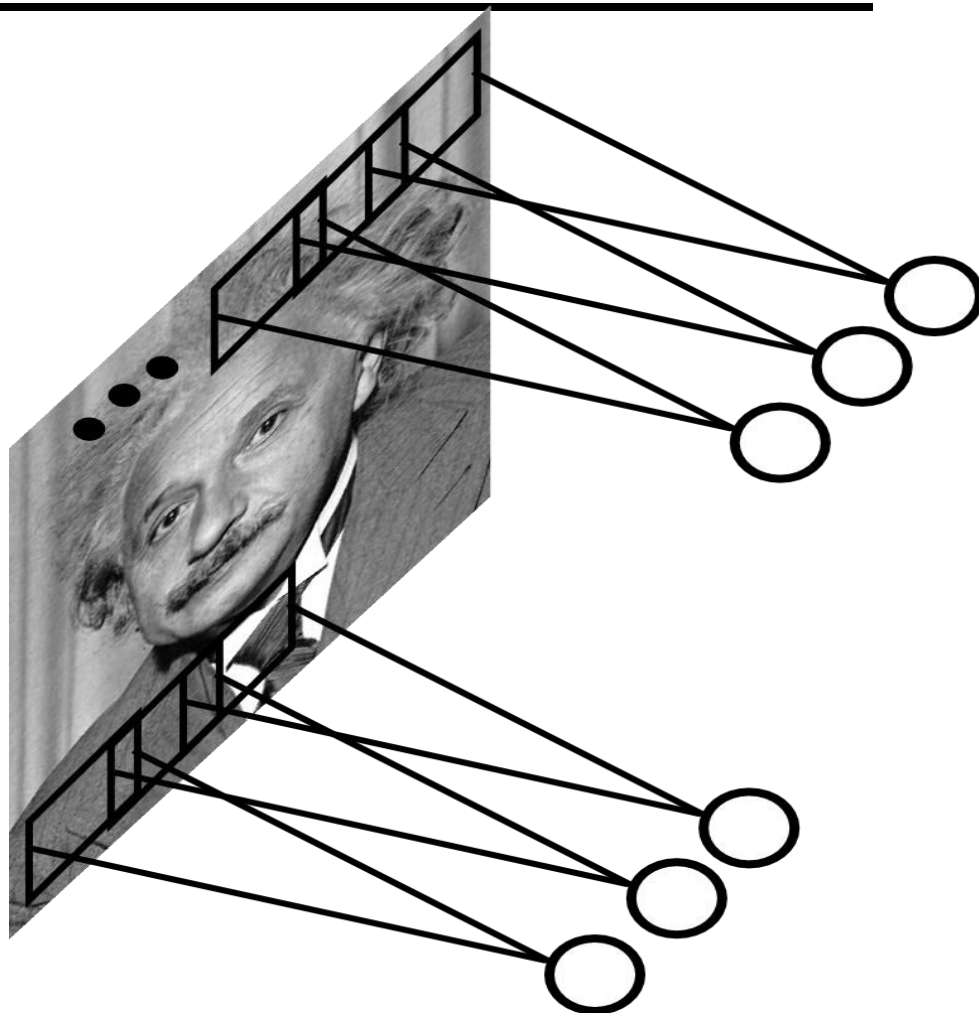
7x7



Convolutional Layer

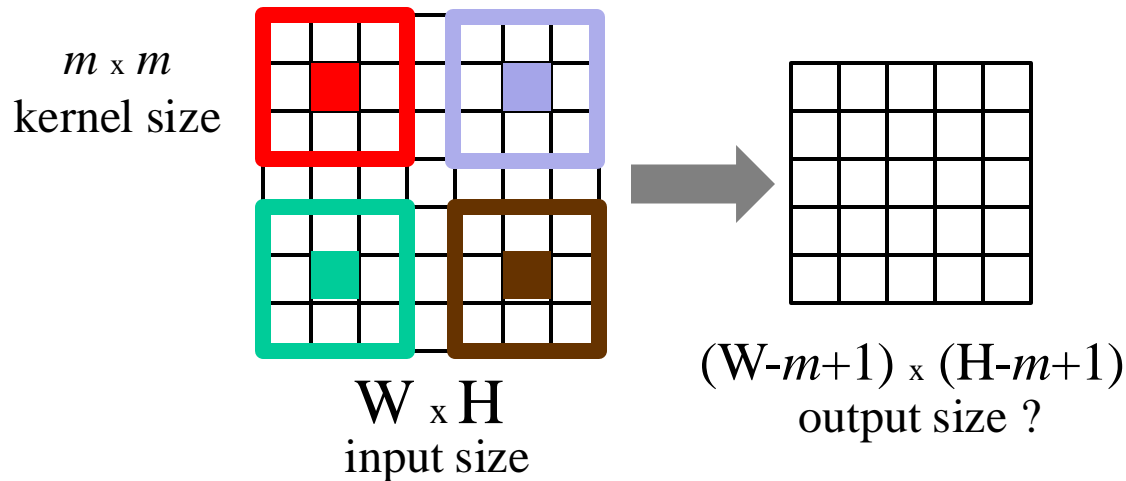
Same as convolution with
some fixed filter

But here the filter
parameters
will be learned

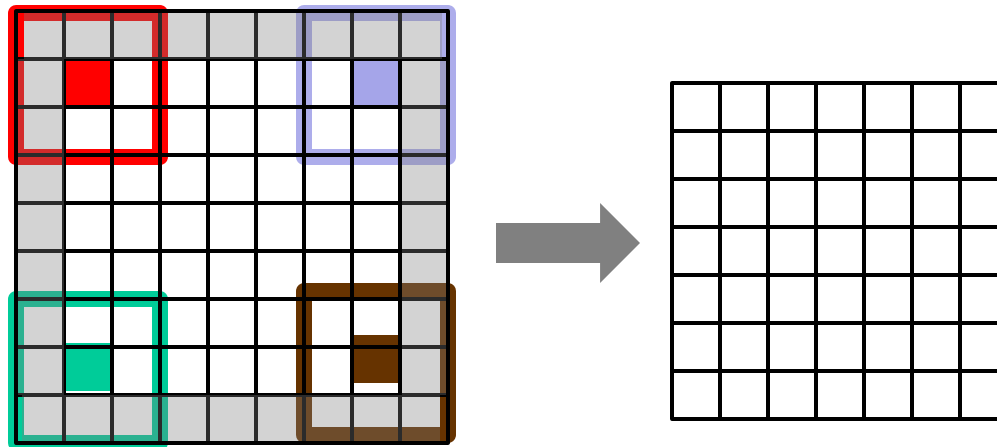


Convolutional Layer - Size Change

Output is usually slightly smaller because the borders of the image are left out



If want output to be the same size, zero-pad the input



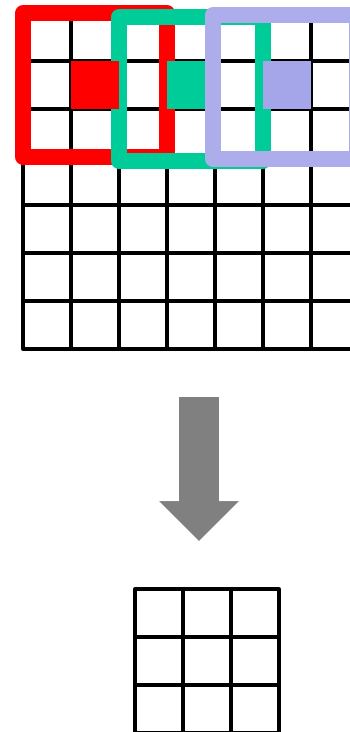
Convolutional Layer - Stride

Can apply convolution only to some pixels (say every second)

- output layer is smaller

Example

- stride = 2 means apply convolution every second pixel
- makes output image approximately **twice smaller** in each dimension
 - image not zero-padded in this example

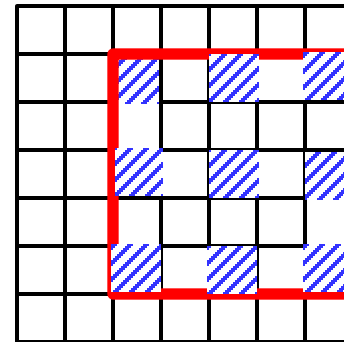


strided convolution

Convolutional Layer - Dilation

It may be helpful to **increase kernel size**
to **enlarge “*receptive field*”**
for each element of the output

But larger kernels could be expensive...



Use only **subset of points** within the kernel's window

atrours convolution

(Fr. *à trous* – hole)

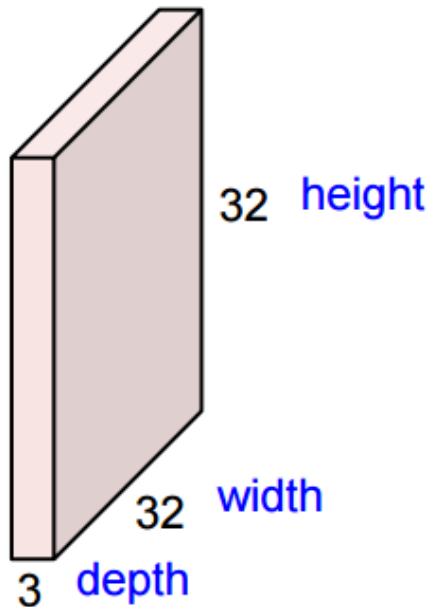
a.k.a. ***dilated convolution***

larger *receptive field* (5x5) for output elements
while effectively using smaller kernels (3x3)

Convolutional Layer – Feature Depth

Input image is usually color, has 3 channels or depth 3

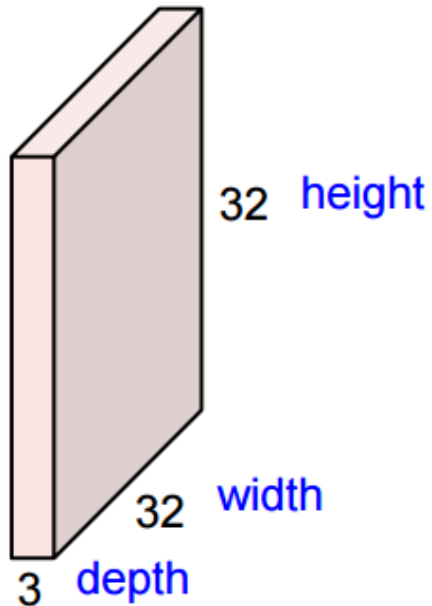
32x32x3 image



Convolutional Layer – Feature Depth

Convolve 3D image with 3D filter

32x32x3 image



5x5x3 filter



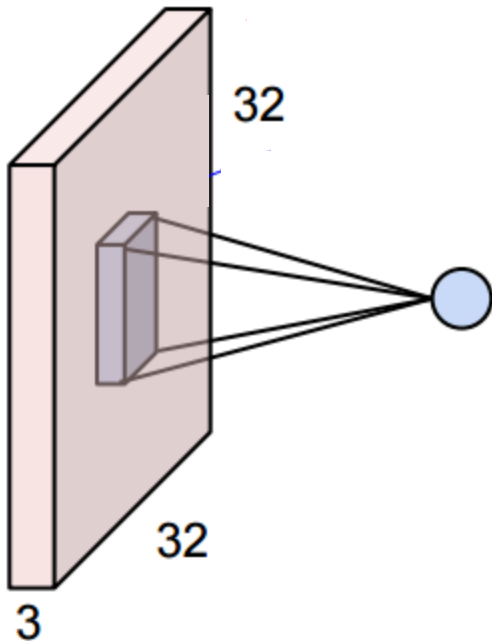
75 parameters

Convolutional Layer – Feature Depth

Each convolution step is a 75 dimensional dot product between the $5 \times 5 \times 3$ filter and a piece of image of size $5 \times 5 \times 3$

Can be expressed as $\mathbf{w}^t \mathbf{x}$, 75 parameters to learn (\mathbf{w})

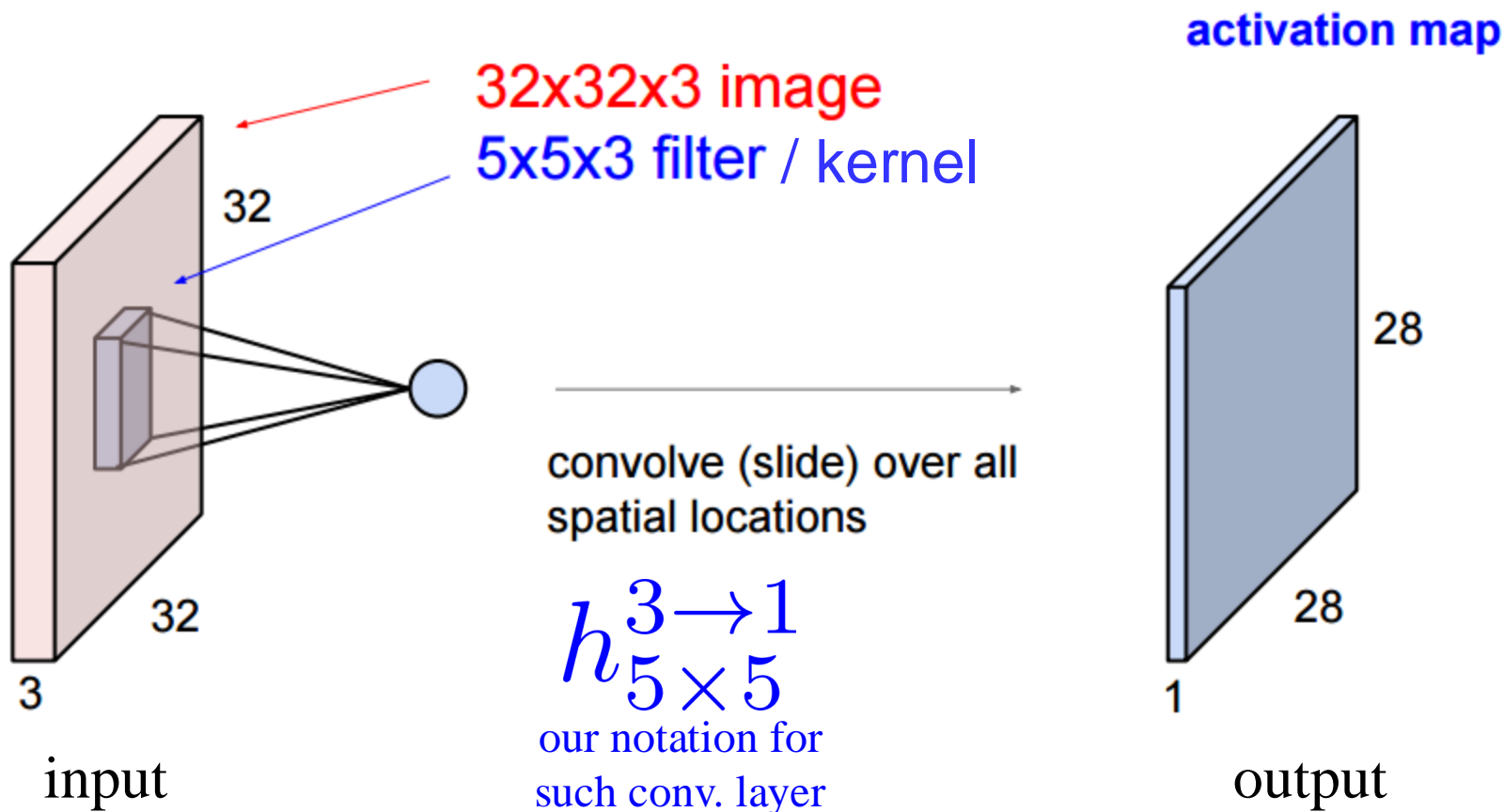
Can add bias $\mathbf{w}^t \mathbf{x} + b$, 76 parameters to learn (\mathbf{w}, b)



Convolutional Layer

Convolve 3D image with 3D filter

- result is a 28x28x1 activation map, no zero padding used



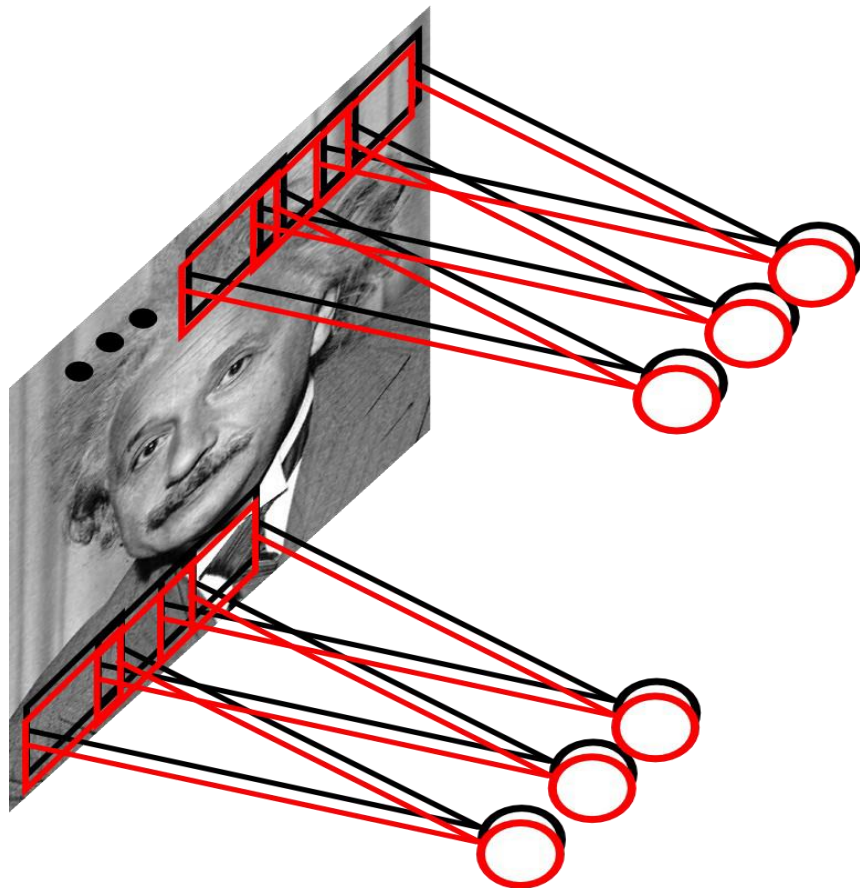
Convolutional Layer

One filter is responsible for
one feature type

Learn multiple filters

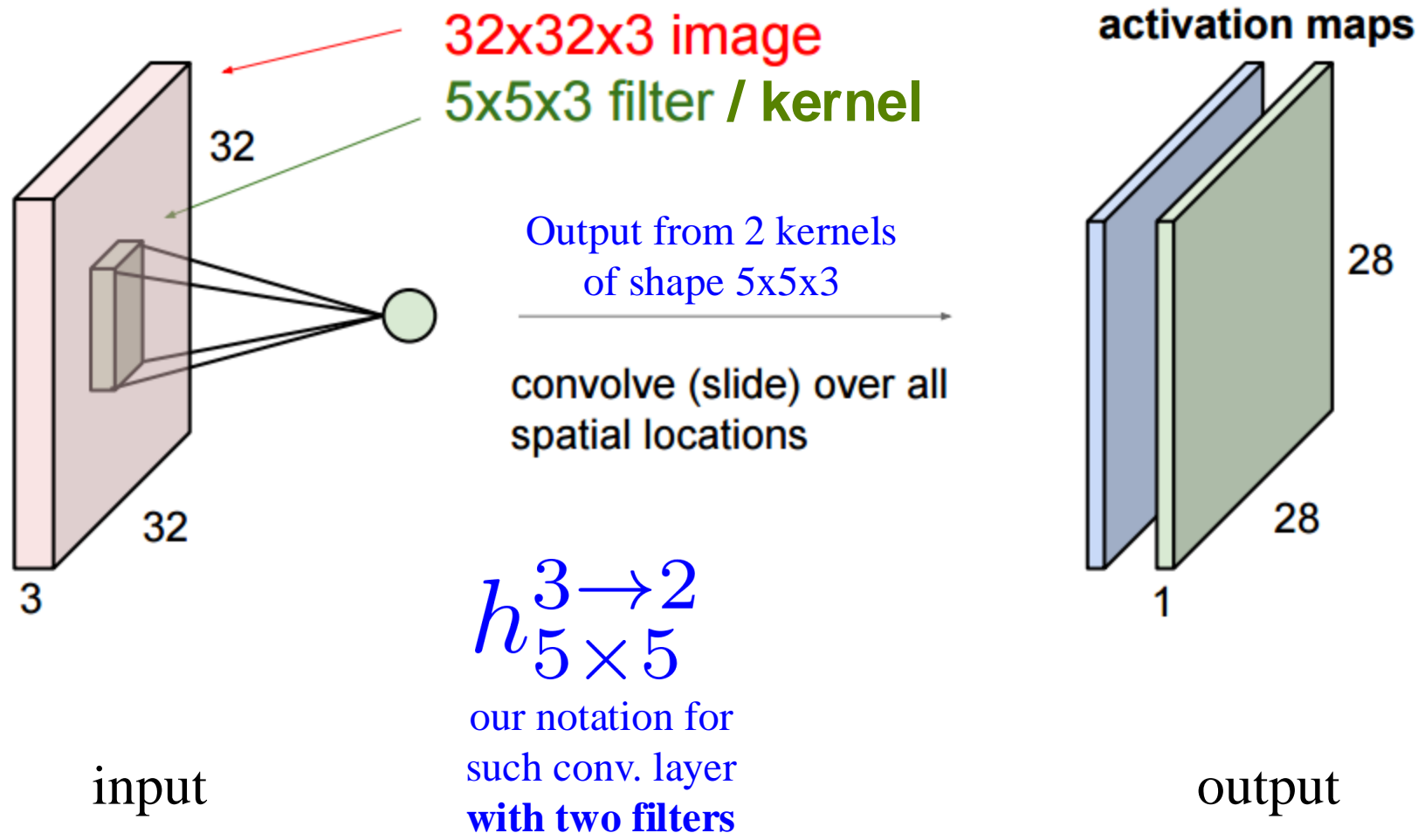
Example:

- 10x10 patch
- 100 filters
- only 10^4 parameters to learn



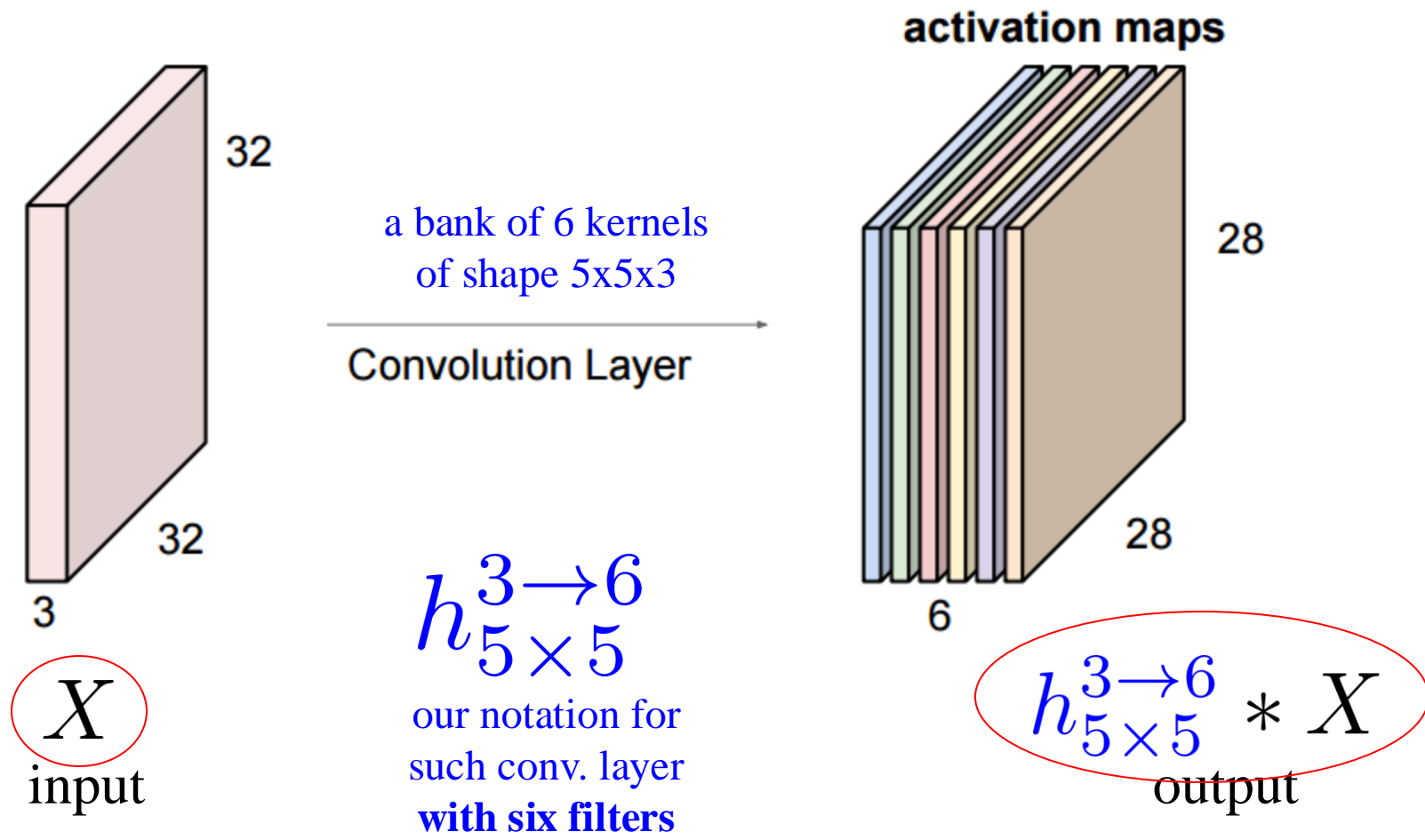
Convolutional Layer

Consider one **extra filter**



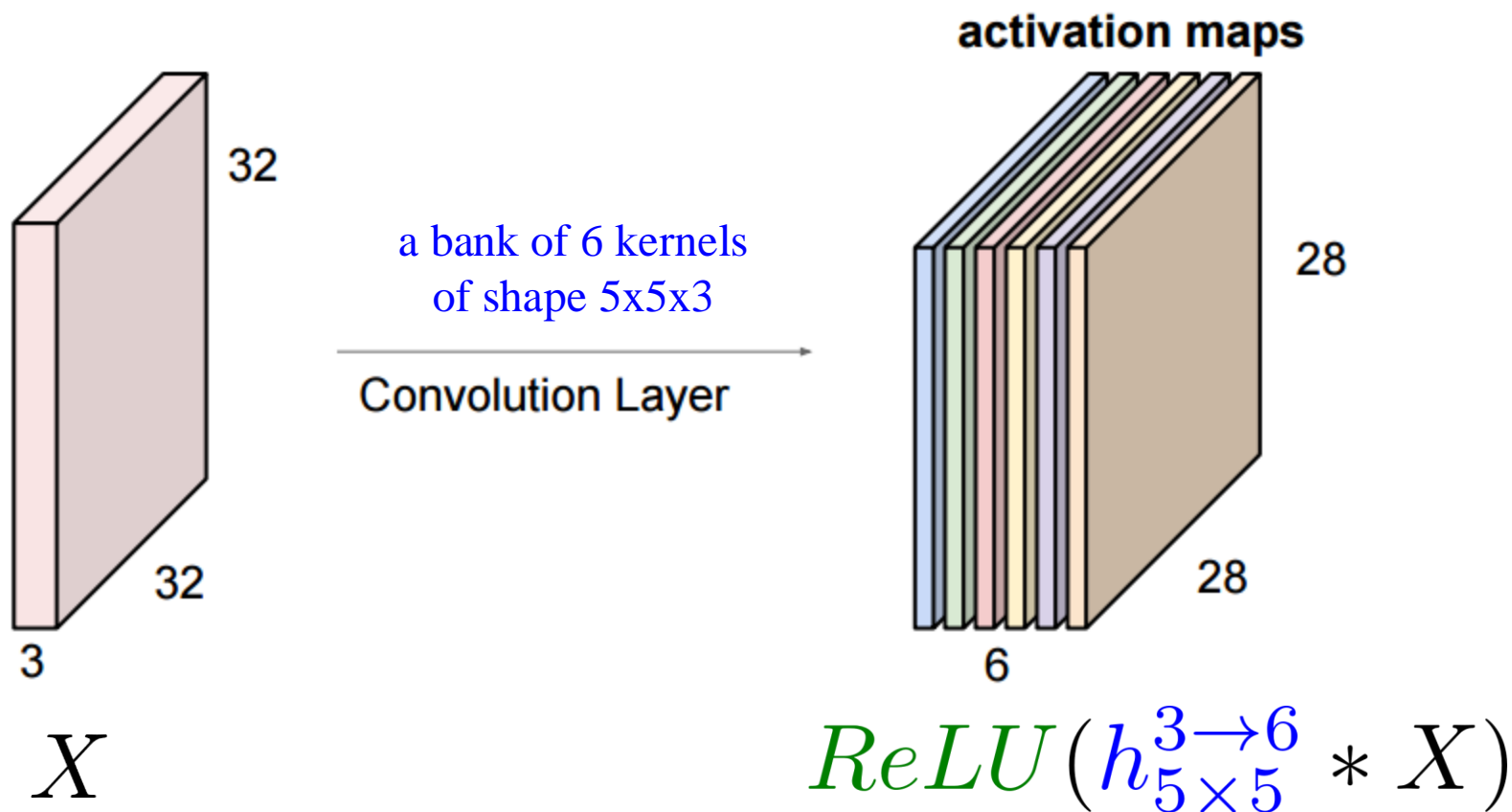
Convolutional Layer

- If have 6 filters (each of size 5x5x3) get 6 activation maps, 28x28 each
- Stack them to get new 28x28x6 “image”



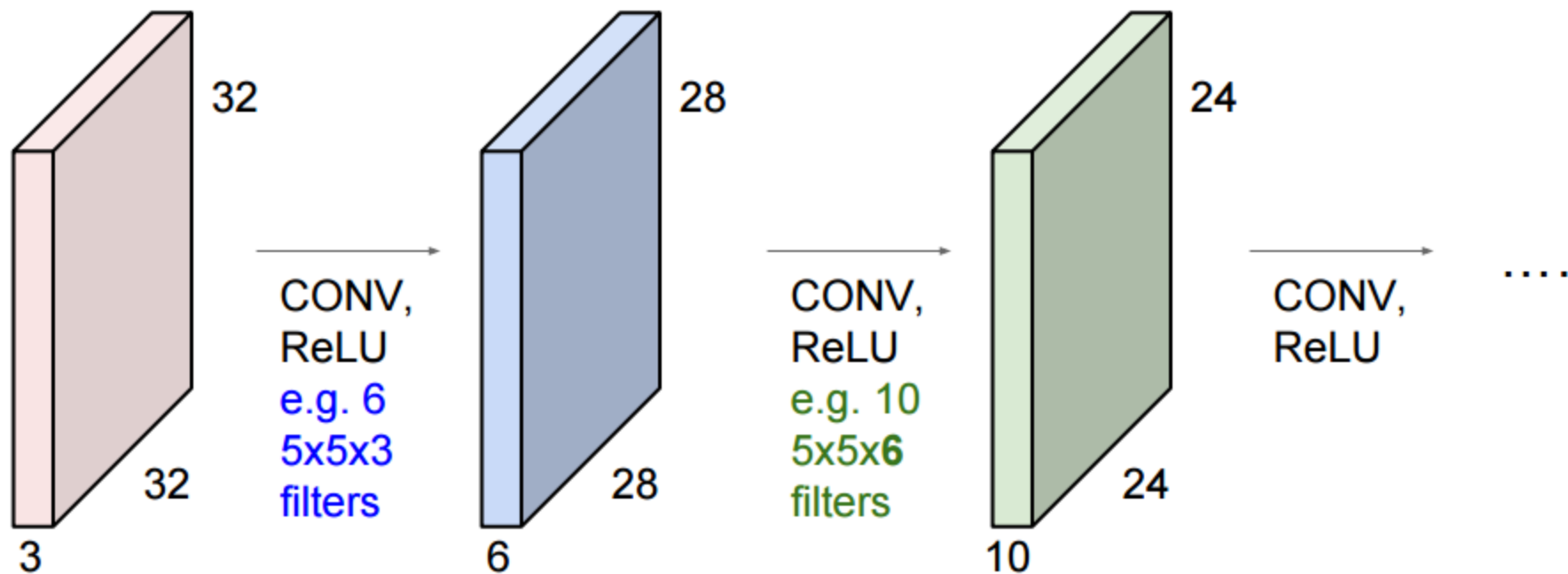
Convolutional Layer

Apply activation function (say **ReLU**) to the activation map



Several Convolution Layers

Construct a sequence of convolution layers interspersed with activation functions



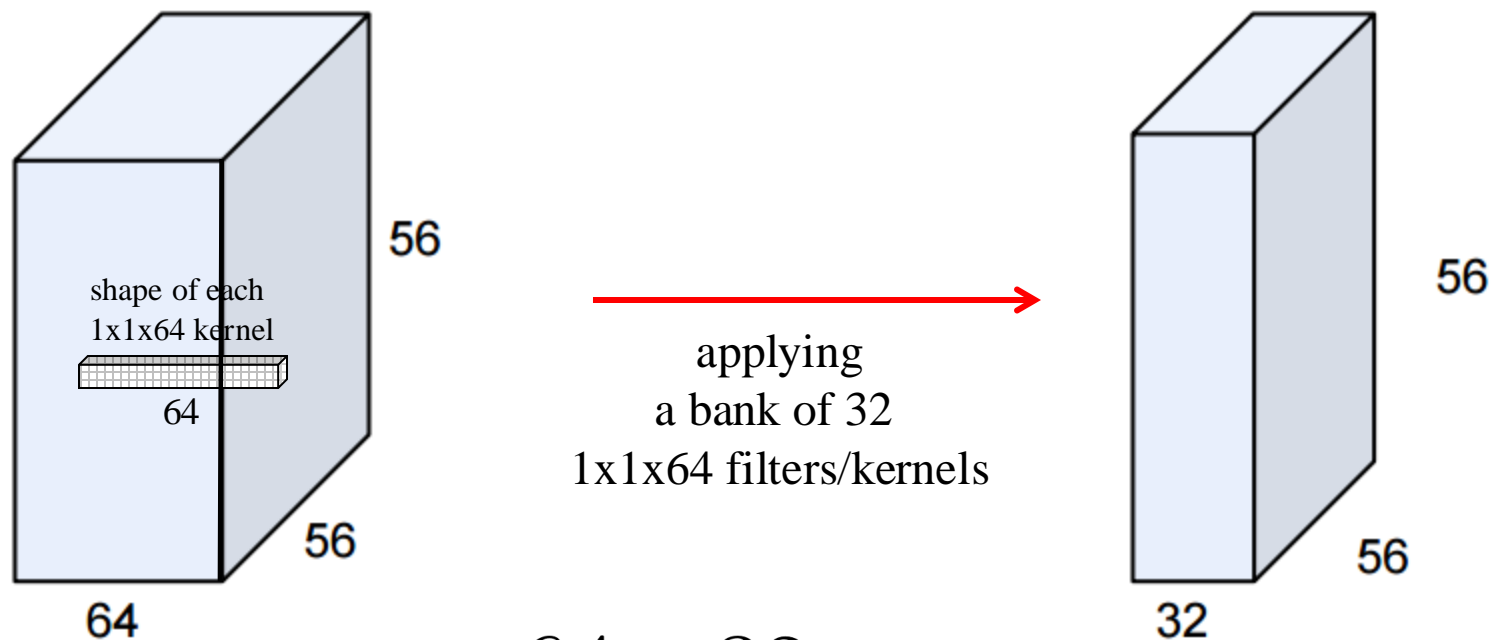
$$ReLU(h_{5 \times 5}^{3 \rightarrow 6} * X) \quad ReLU(h_{5 \times 5}^{6 \rightarrow 10} * X)$$

Convolutional Layer

1x1 convolutions make perfect sense

Example

- Input image of size 56x56x64
- Convolve with 32 filters, each of size 1x1x64

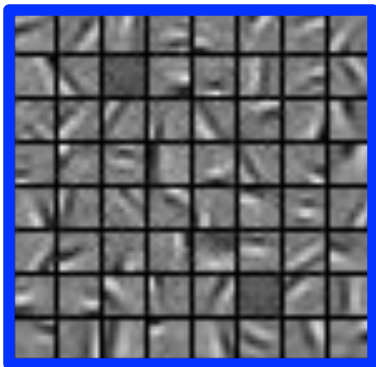


$$h_{1 \times 1}^{64 \rightarrow 32} * X$$

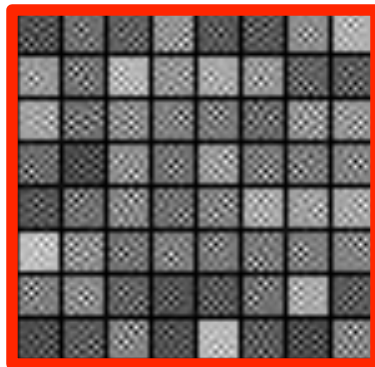
Check Learned Convolutions

- Good training: learned filters exhibit structure and are uncorrelated

GOOD

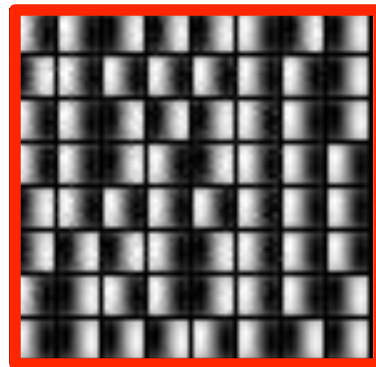


BAD



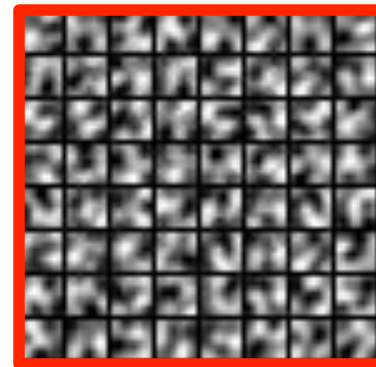
too noisy

BAD



too
correlated

BAD



lack
structure

Convolutional Layer Summary

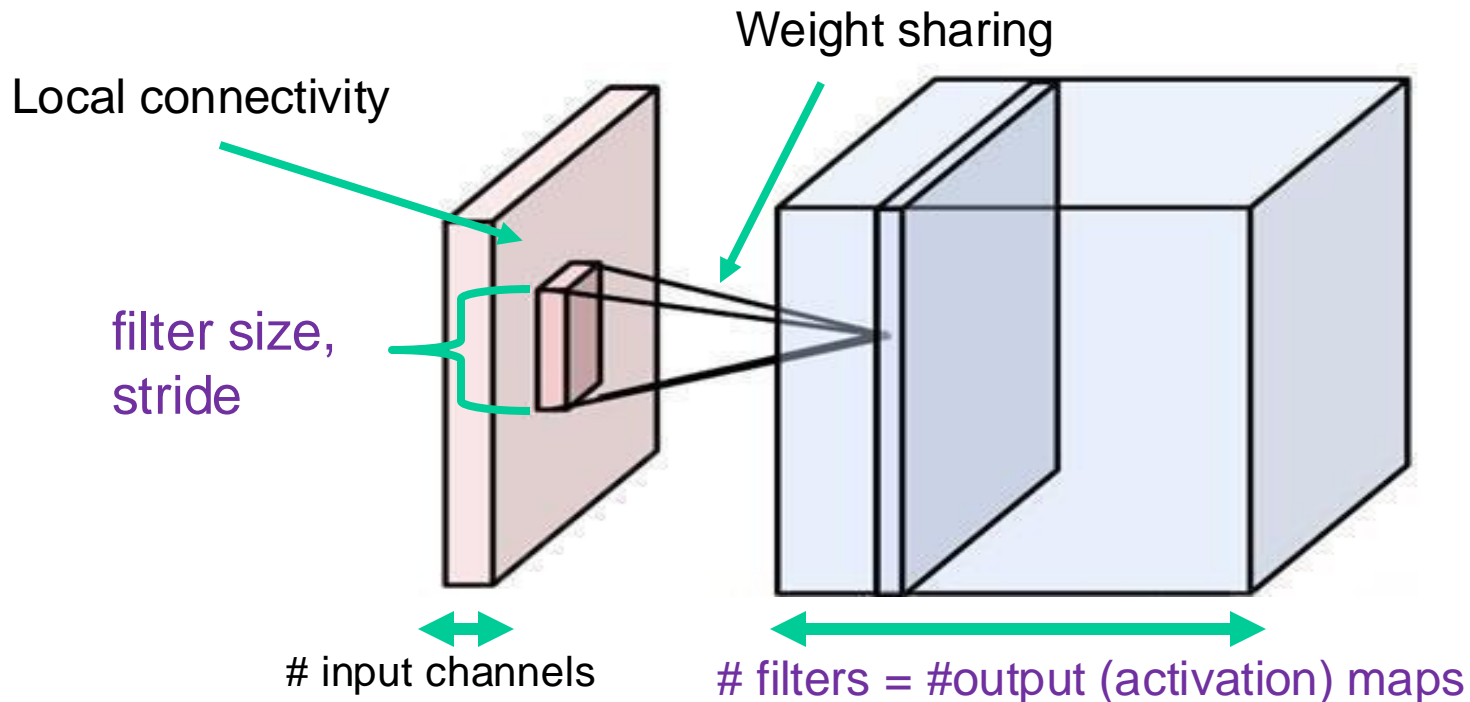
Local connectivity

Weight sharing

Handling multiple input/output channels

Retains location associations

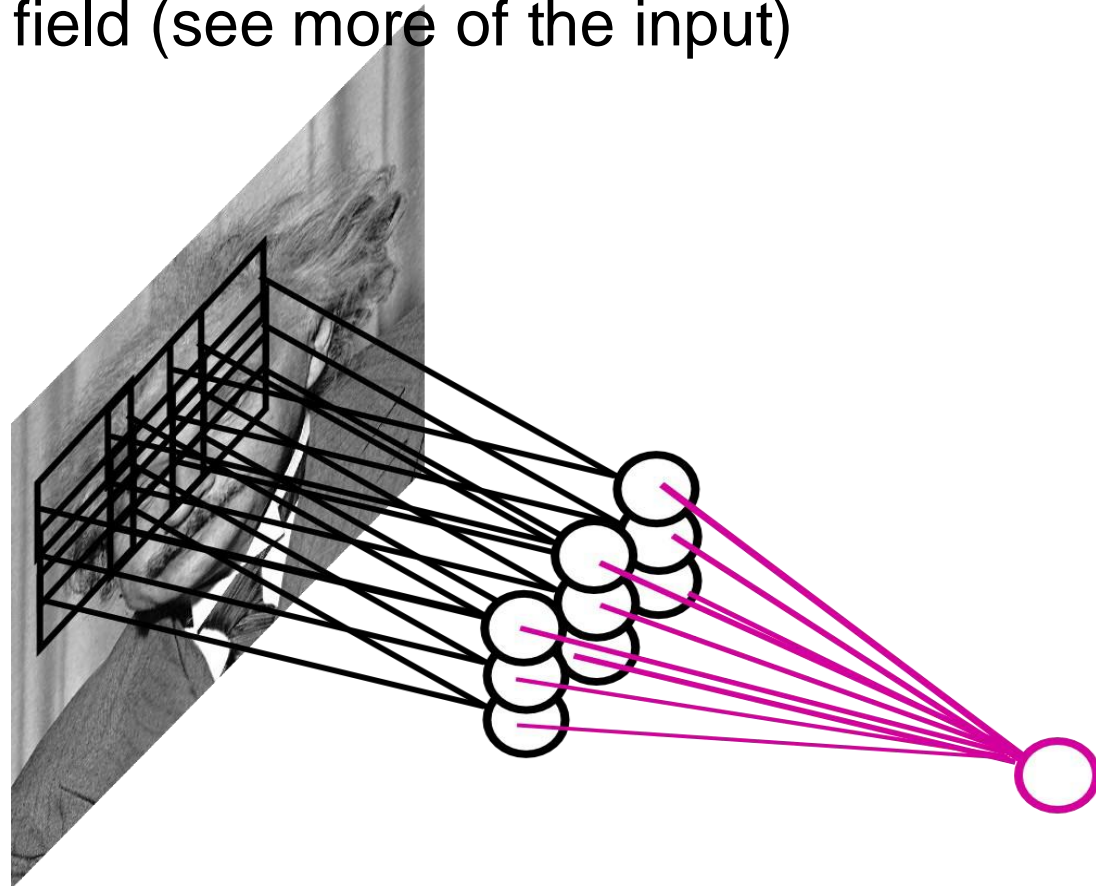
Transforms 3D tensor into 3D tensor (**tensor flow**)



Pooling Layer

Pool responses at different locations

- by taking **max**, **average**, etc.
- robustness to exact spatial location
- also larger receptive field (see more of the input)
- Usually pooling applied with stride > 1
- This reduces resolution of output map



Pooling Layer: Max Pooling Example

Single depth slice

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

max pool with 2x2 filters
and stride 2

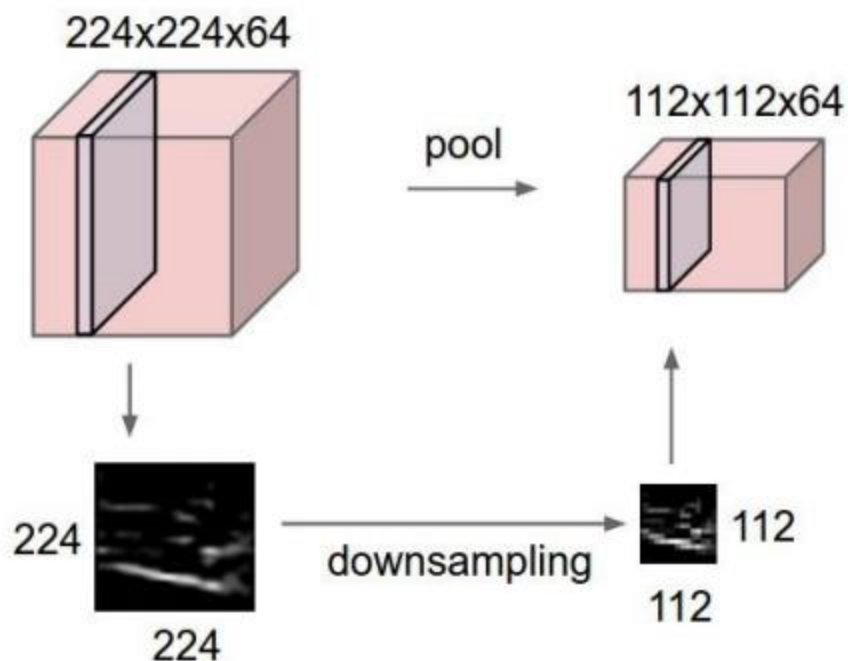


6	8
3	4

- pooling can be interpreted as *downsampling*

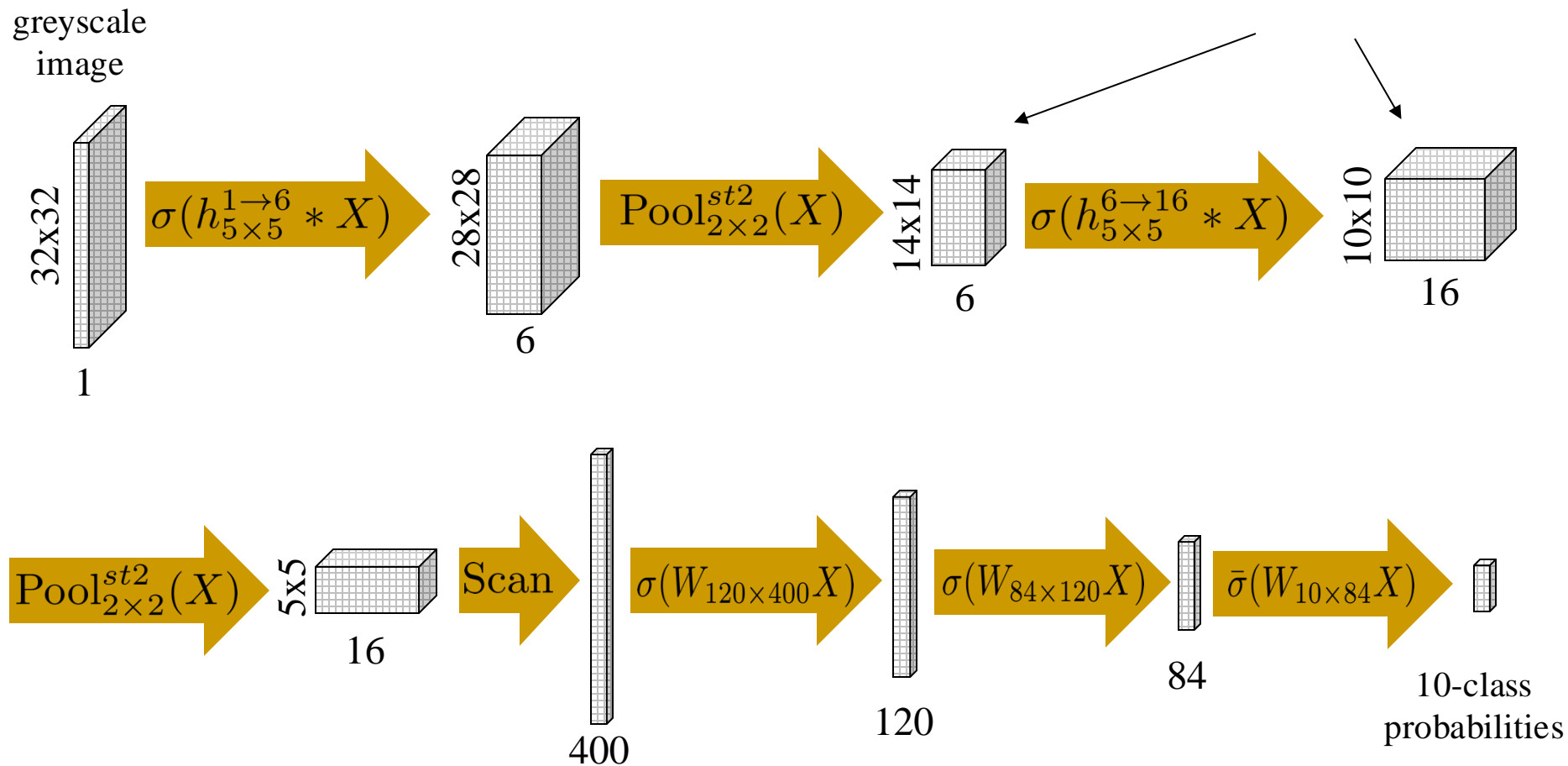
Pooling Layer

Pooling usually applied to each activation map separately



Basic CNN example (à la *LeNet* -1998)

NOTE: transformation of multi-dimensional arrays (**tensors**)

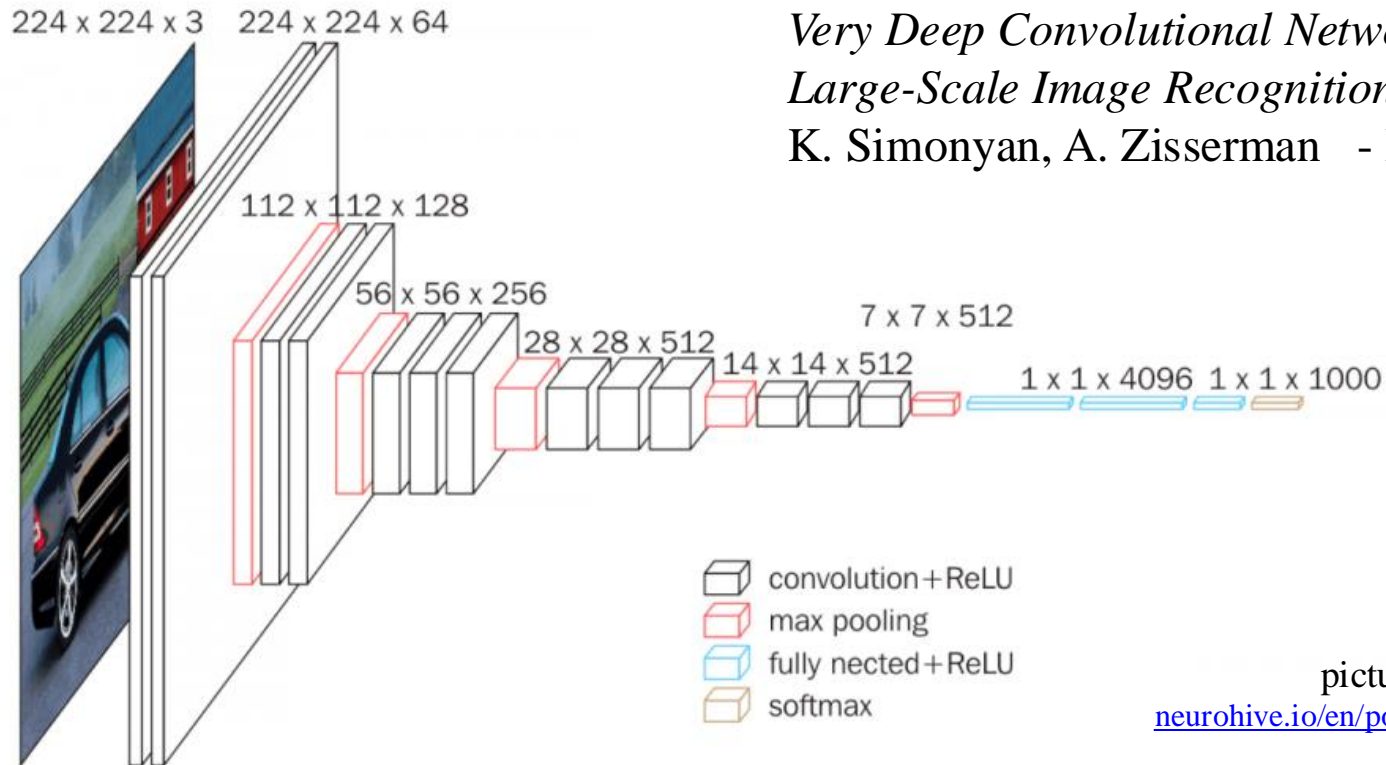


Deep CNN architectures for classification

- **AlexNet** (2012) *ImageNet classification with deep convolutional neural networks*
Alex Krizhevsky, Ilya Sutskever, Geoffrey Hinton - NIPS 2012.
- **VGG** (2014) *Very Deep Convolutional Networks for Large-Scale Image Recognition*
K. Simonyan, A. Zisserman - ICLR 2015
<http://www.robots.ox.ac.uk/~vgg/practicals/cnn/index.html>
- **ResNet** (2016) *Deep residual learning for image recognition*
K. He, X. Zhang, S. Ren, J. Sun. - CVPR 2016

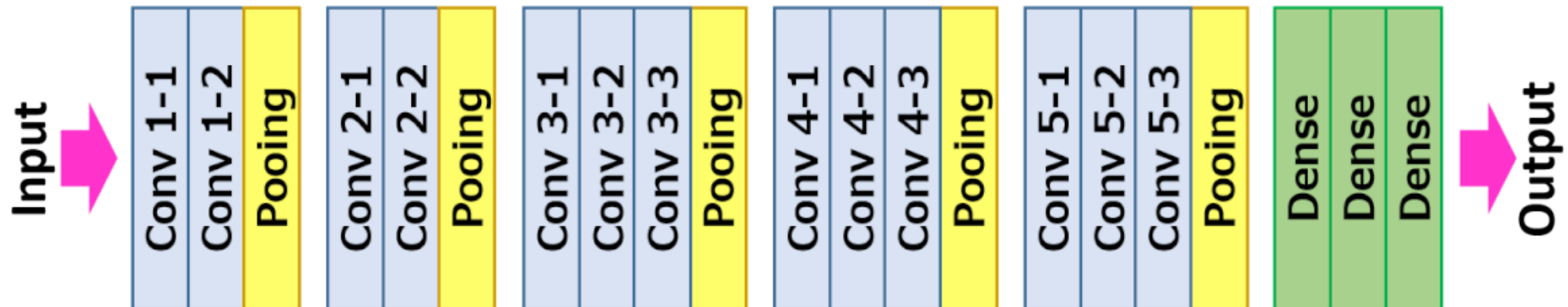
VGG -16

*Very Deep Convolutional Networks for
Large-Scale Image Recognition*
K. Simonyan, A. Zisserman - ICLR 2015



picture credits

neurohive.io/en/popular-networks/vgg16/



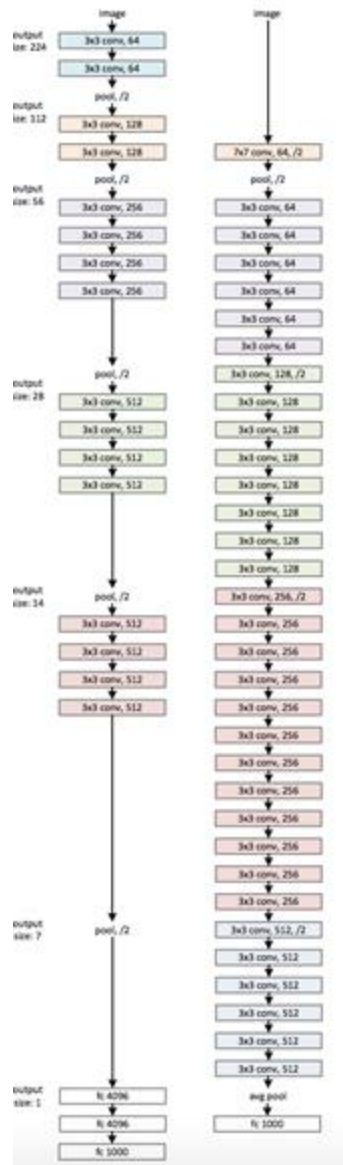
ResNet

very deep 😊

one of the
state of the art
on *image net*

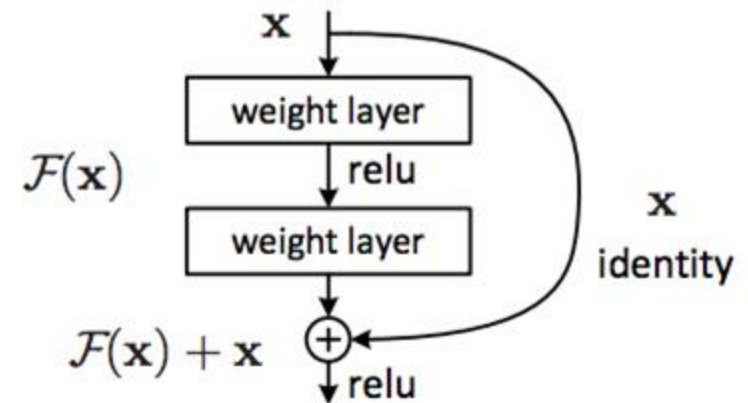
www.image-net.org

- very large dataset
of labeled images
>14,000,000



Deep residual learning for image recognition. K. He, X. Zhang, S. Ren, and J. Sun. CVPR 2016

key technical trick



resnet block

(residual link helps gradient descent)

FashionMNIST classification example

