

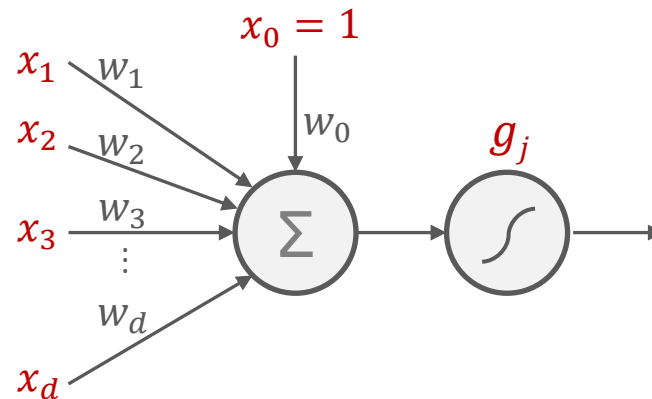
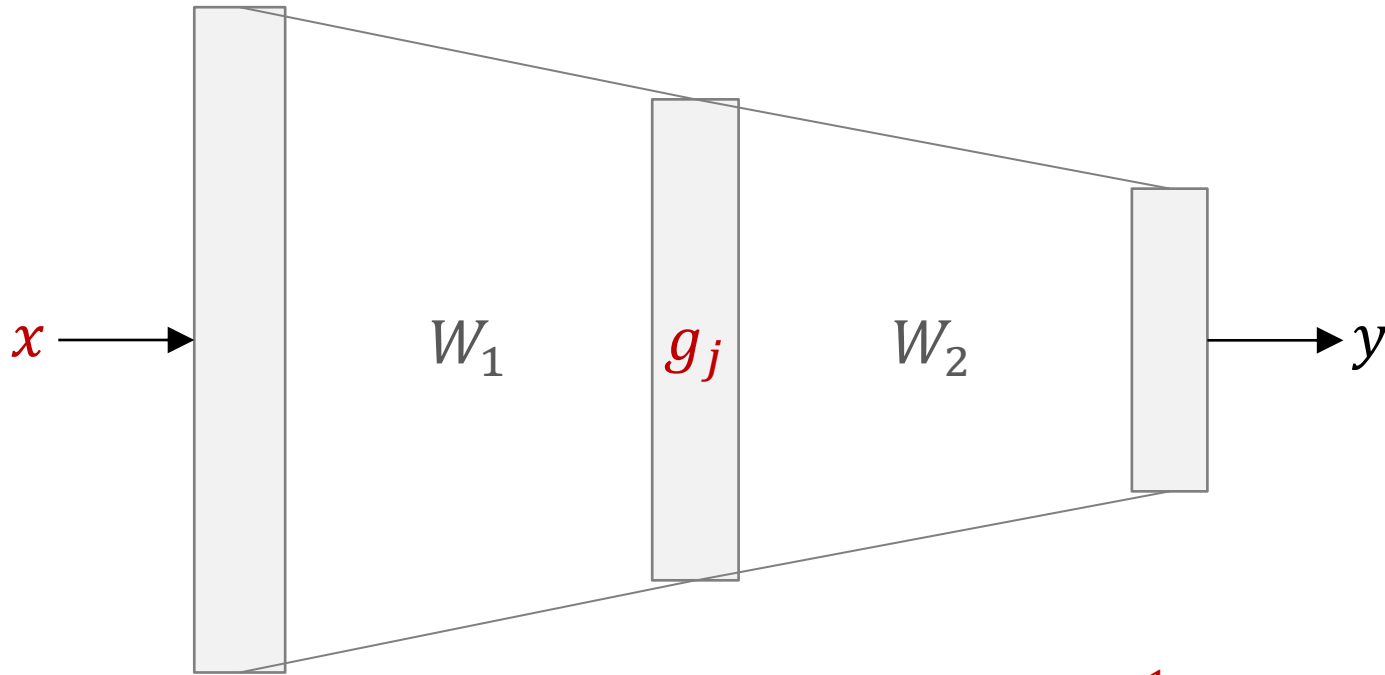


MACHINE LEARNING
University Master's Degree in Intelligent Systems

convolutional neural networks

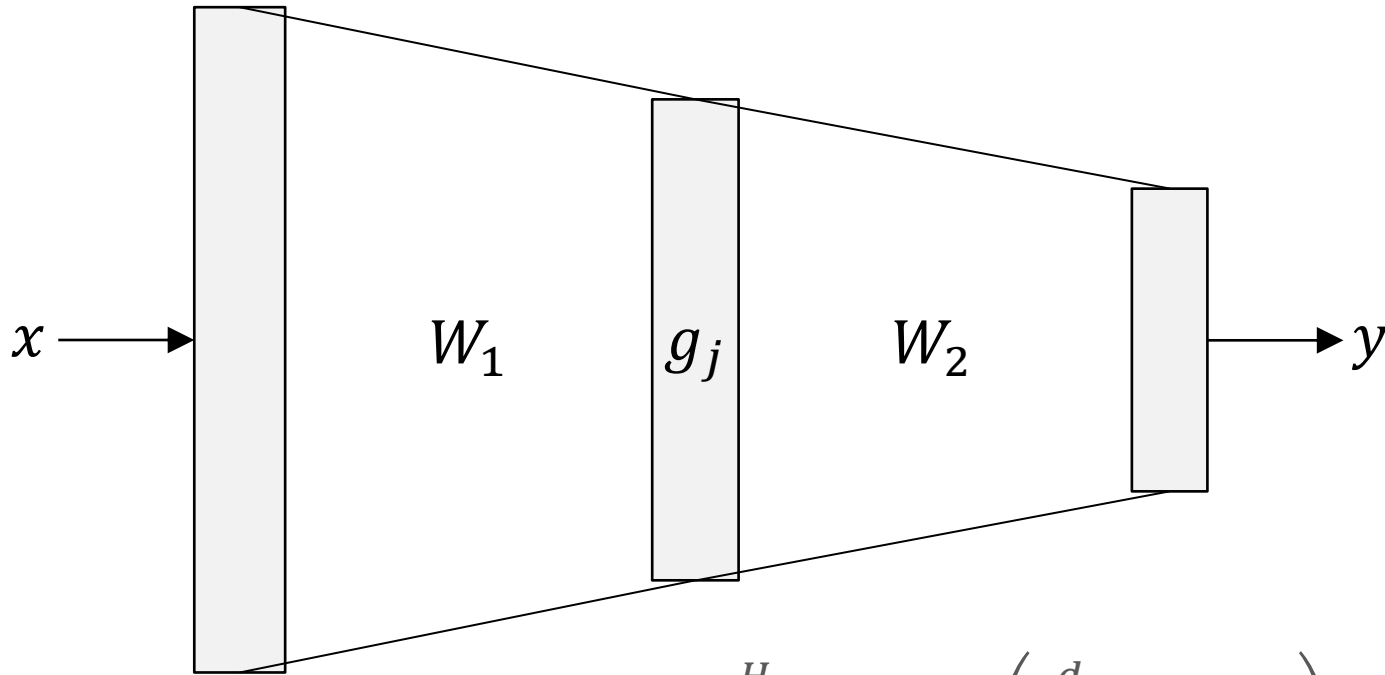
Ramón A. Mollineda Cárdenas

what we already know
fully-connected neural network



what we already know

fully-connected neural network



$$y_k = \sum_{j=1}^H w_2^{jk} \cdot g \left(\sum_{i=1}^d w_1^{ij} \cdot x_i \right) = W_2 \cdot g(W_1 \cdot x)$$

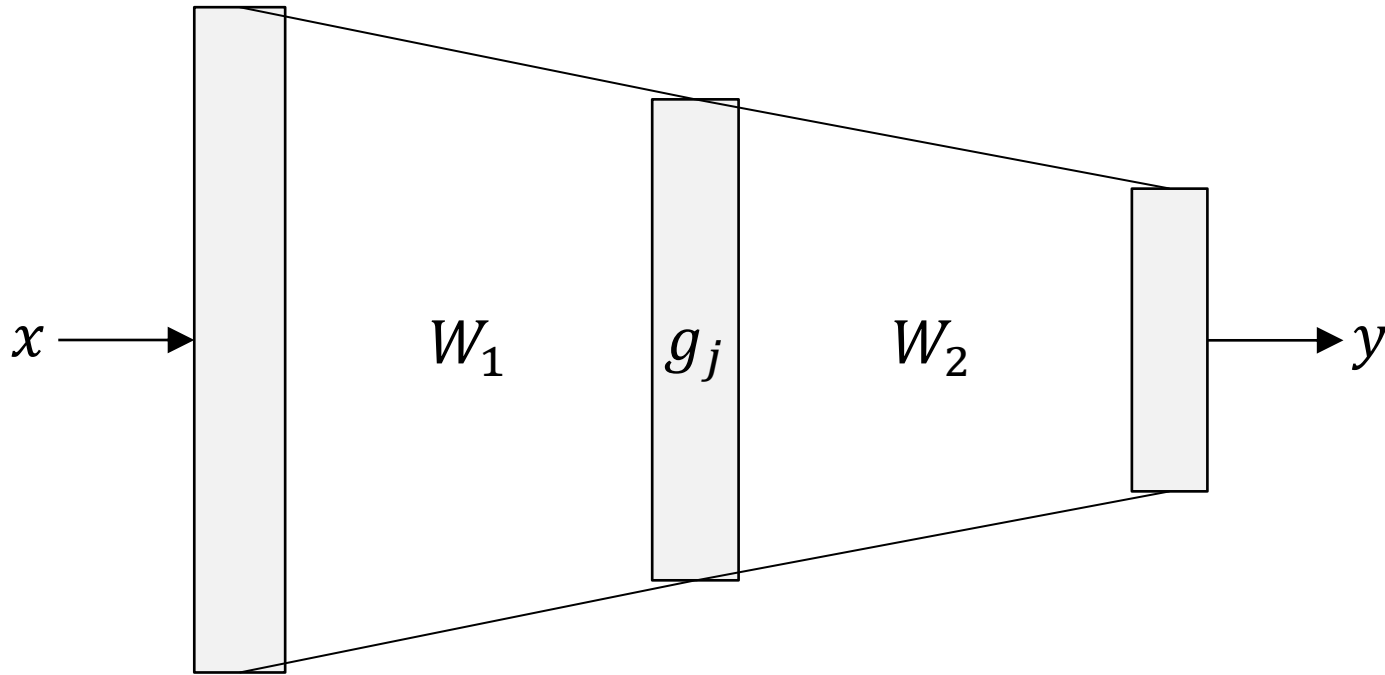
d , input space dimension

H , number of hidden layer units

g , nonlinear activation function

what we already know

fully-connected neural network

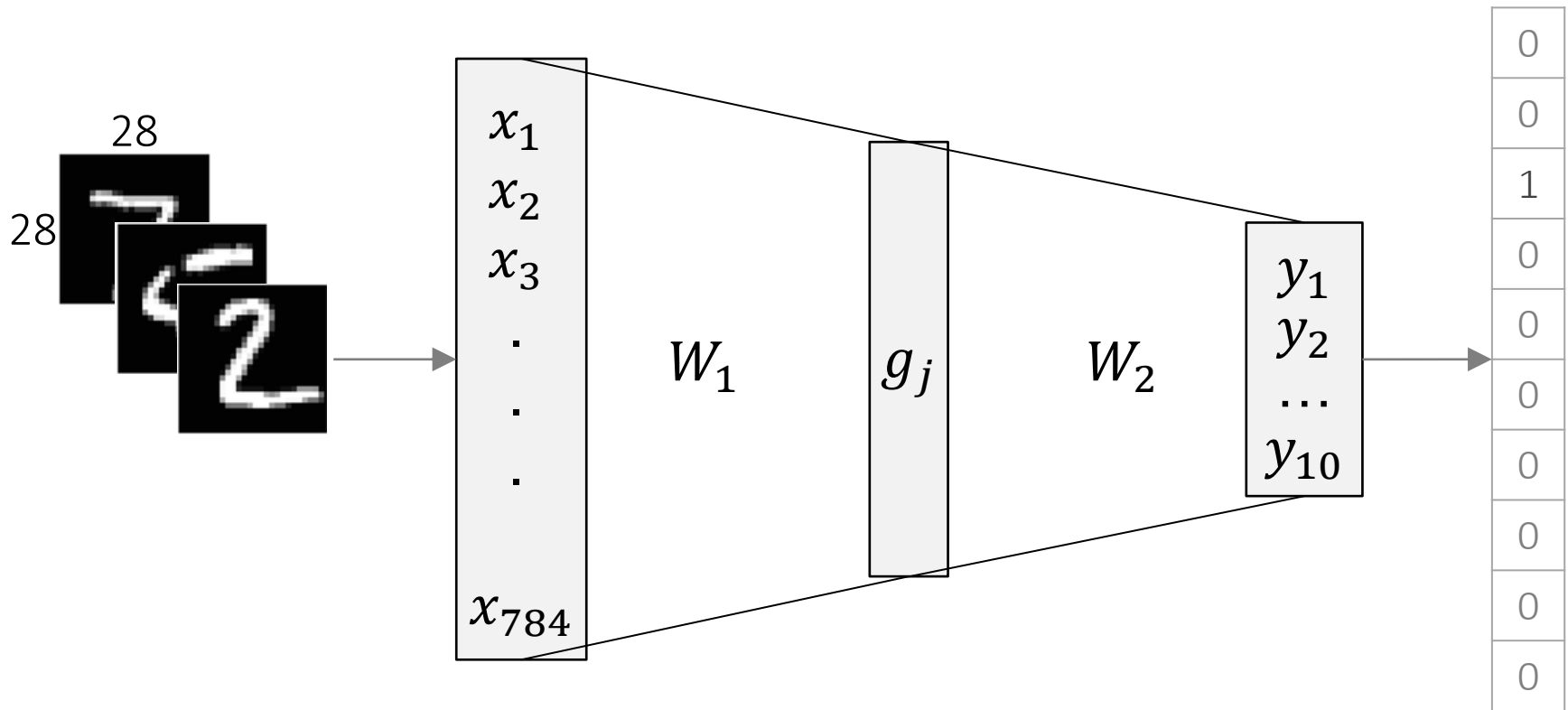


- input: $x \in \mathbb{R}^N$
- two-layer neural network
- g non-linear functions (e.g. sigmoide, relu)
- output: linear combination of g_j

what we already know

fully-connected neural network

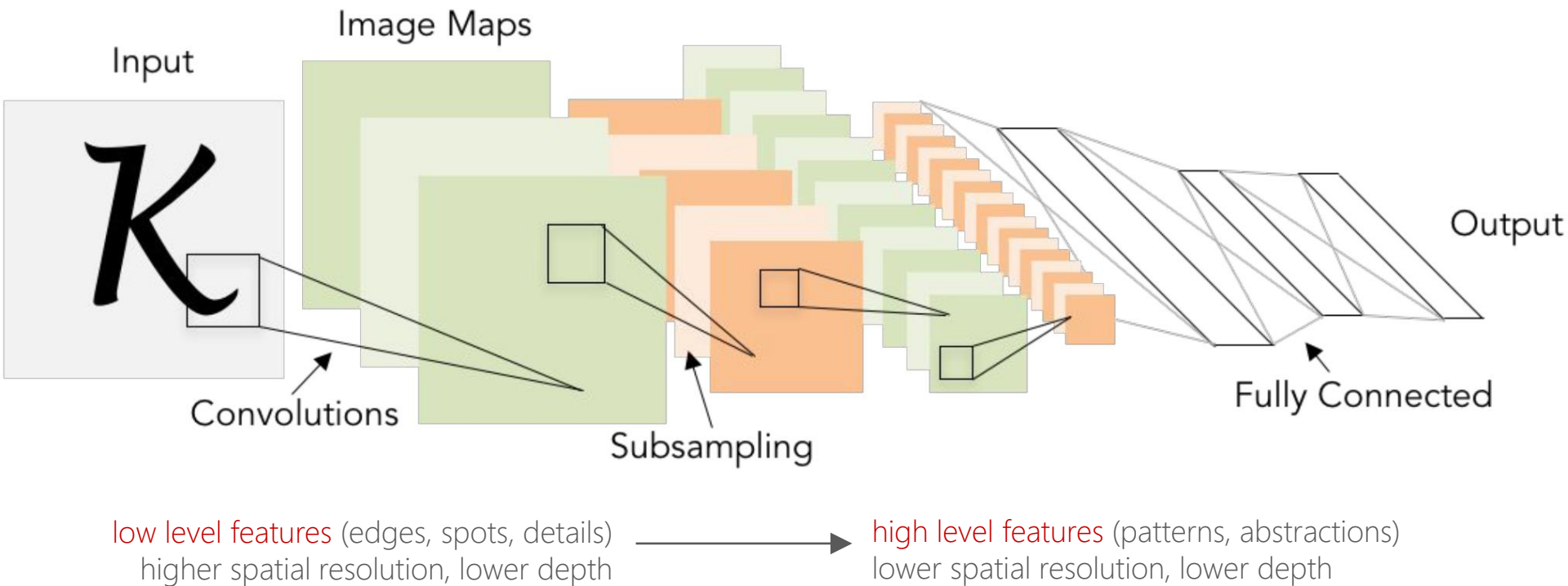
MNIST digits case study



convolutional neural network (CNN)

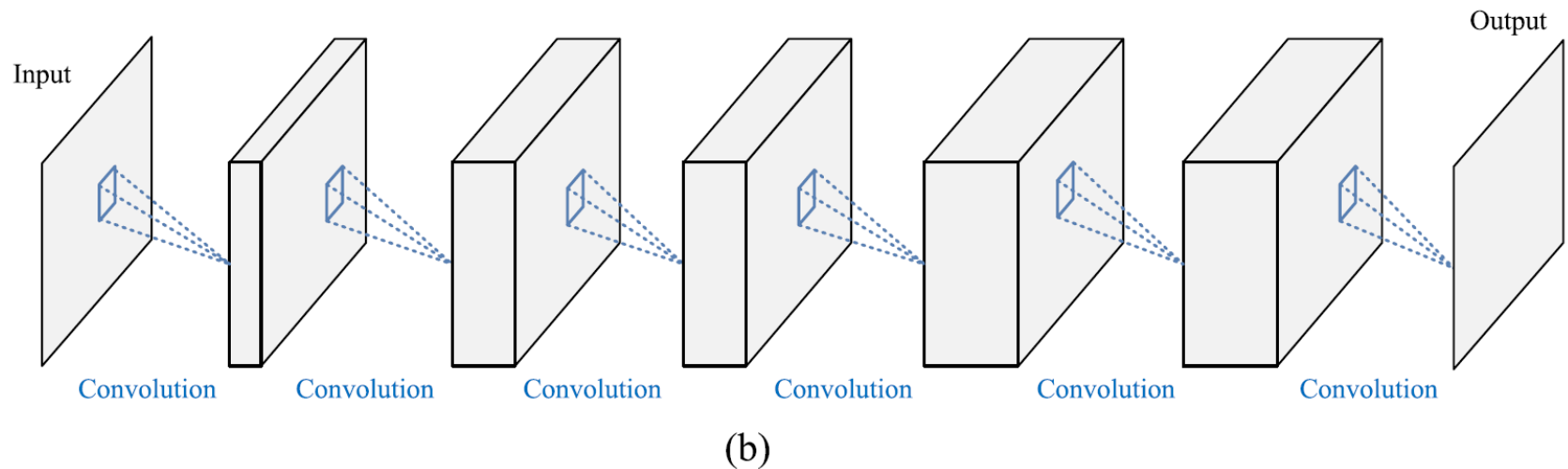
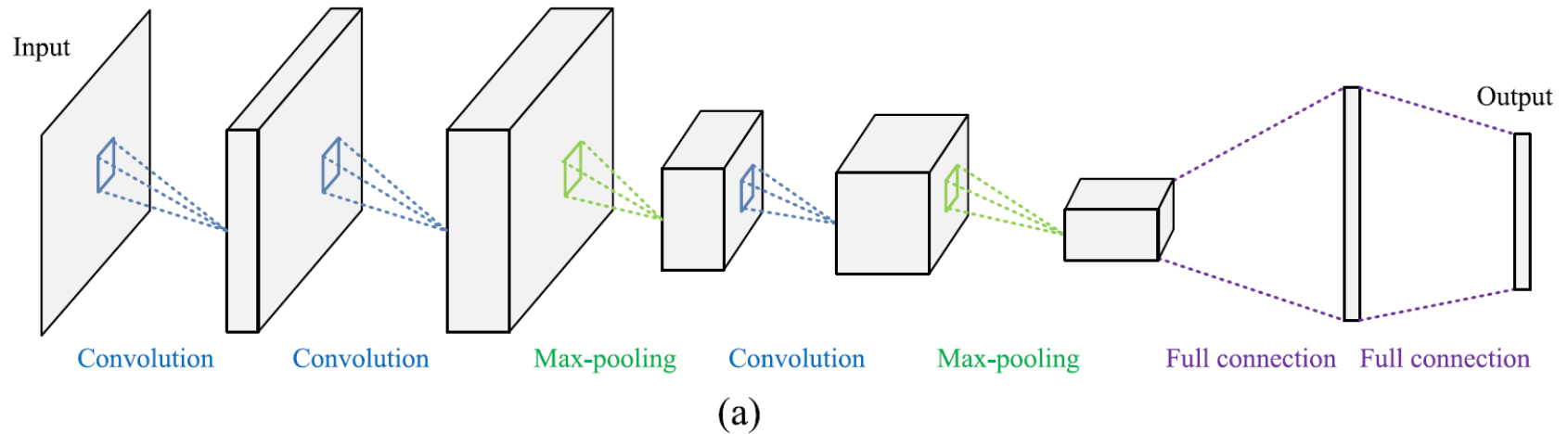
overview

CNN: end-to-end solutions transforms primary representations into categories



convolutional neural network (CNN)

architectural patterns

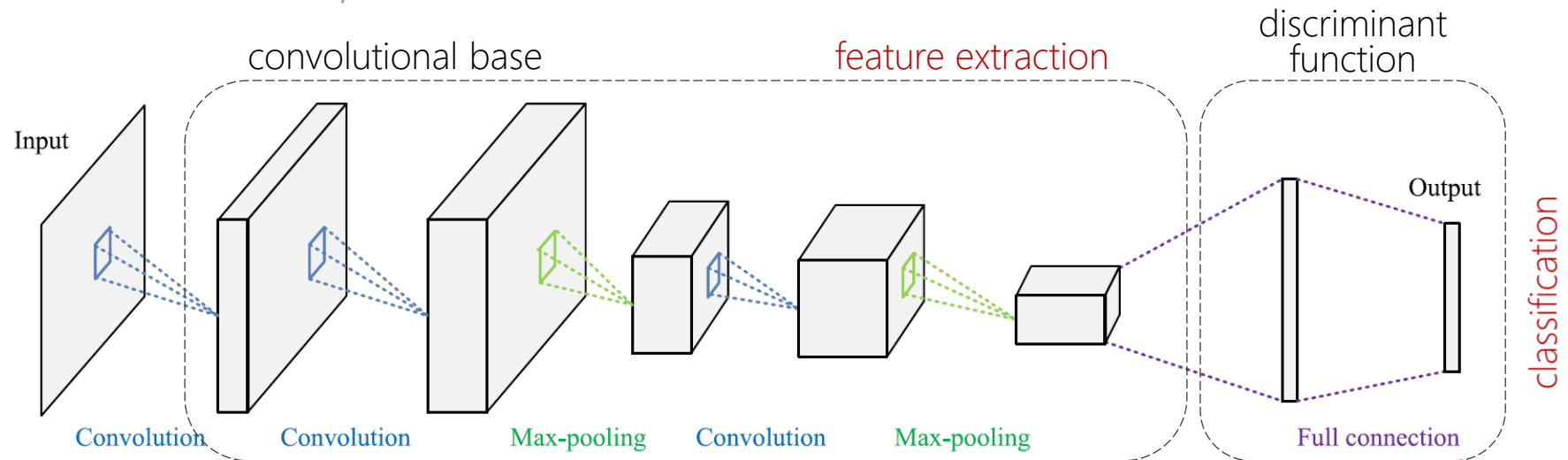


(a) classification

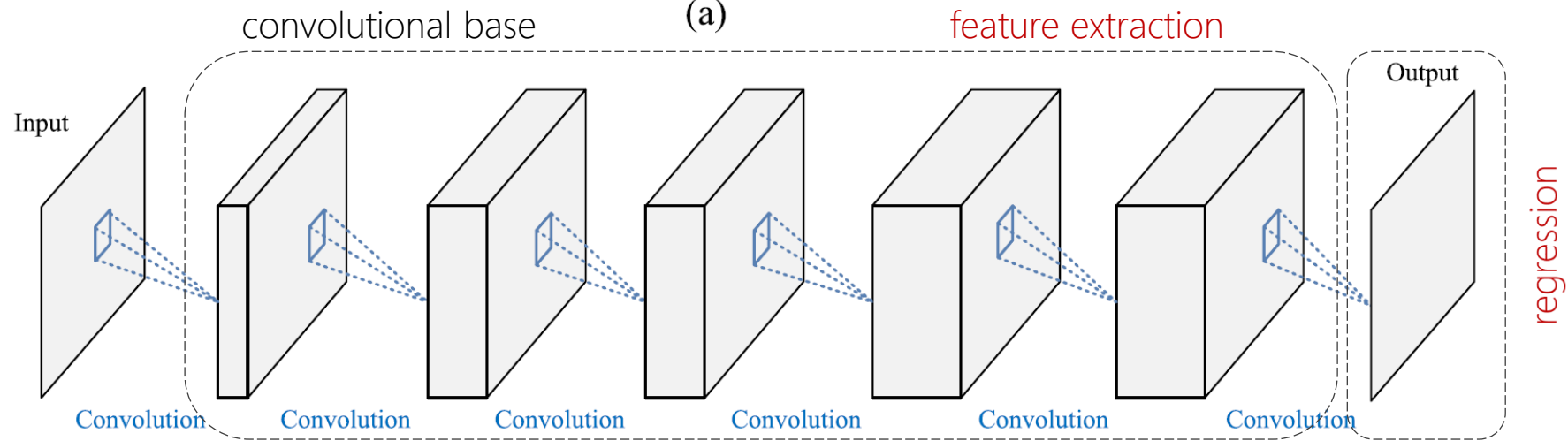
(b) regression (image reconstruction)

convolutional neural network (CNN)

architectural patterns



(a)



(b)

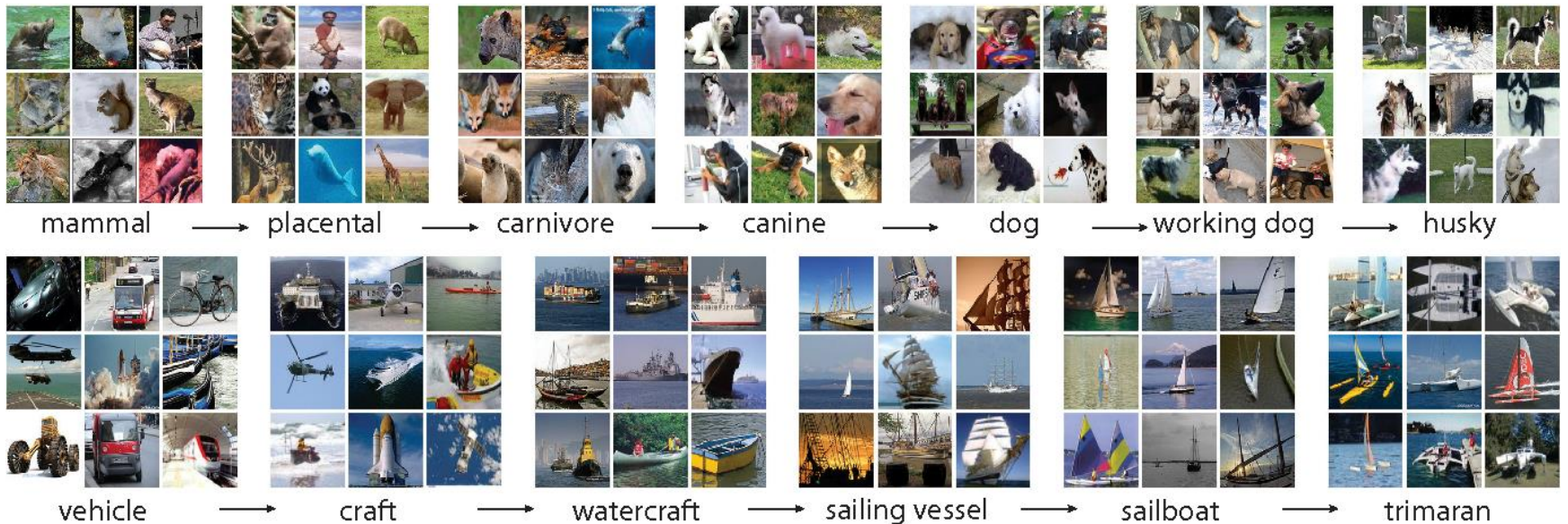
(a) **classification**

(b) **regression** (image reconstruction)

ImageNet challenge

Large Scale Visual Recognition Challenge (ILSVRC)

IMAGENET (<http://image-net.org>)



ImageNet challenge

Large Scale Visual Recognition Challenge (ILSVRC)

task: given an image, identify the main object in the image

training data:

- 1,200,000 labeled images
- 1,000 final classes/categories (ground truth)
- one class label per image (identifies the main object)

validation and test data:

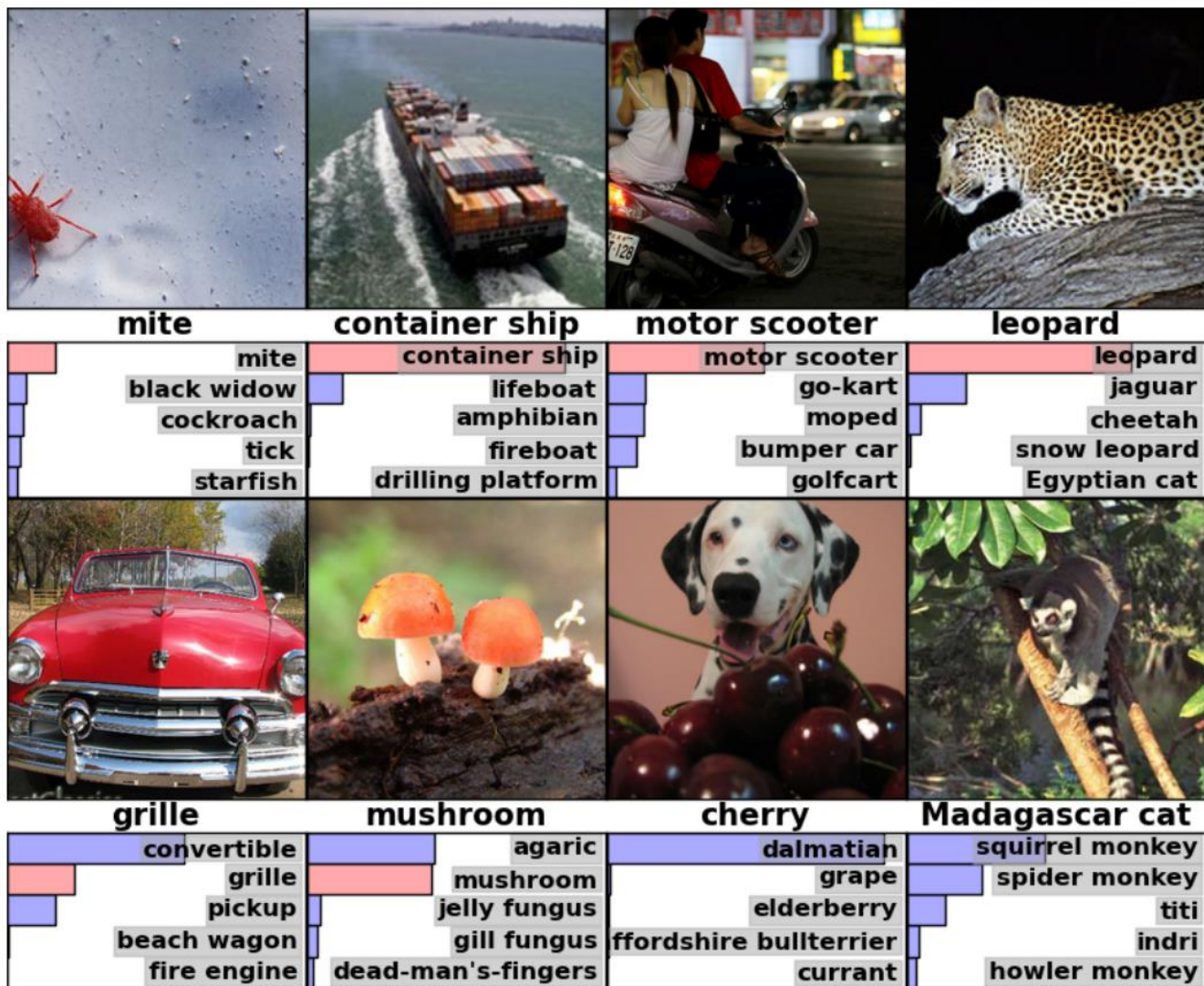
- 150,000 pictures (from Flickr and other search engines)
- 50,000 labeled validation images
- 100,000 unlabeled test data

success/hit: the correct class (ground truth) is one of the 5 most probable classes found by the model (top-5 error)

ImageNet challenge

Large Scale Visual Recognition Challenge (ILSVRC)

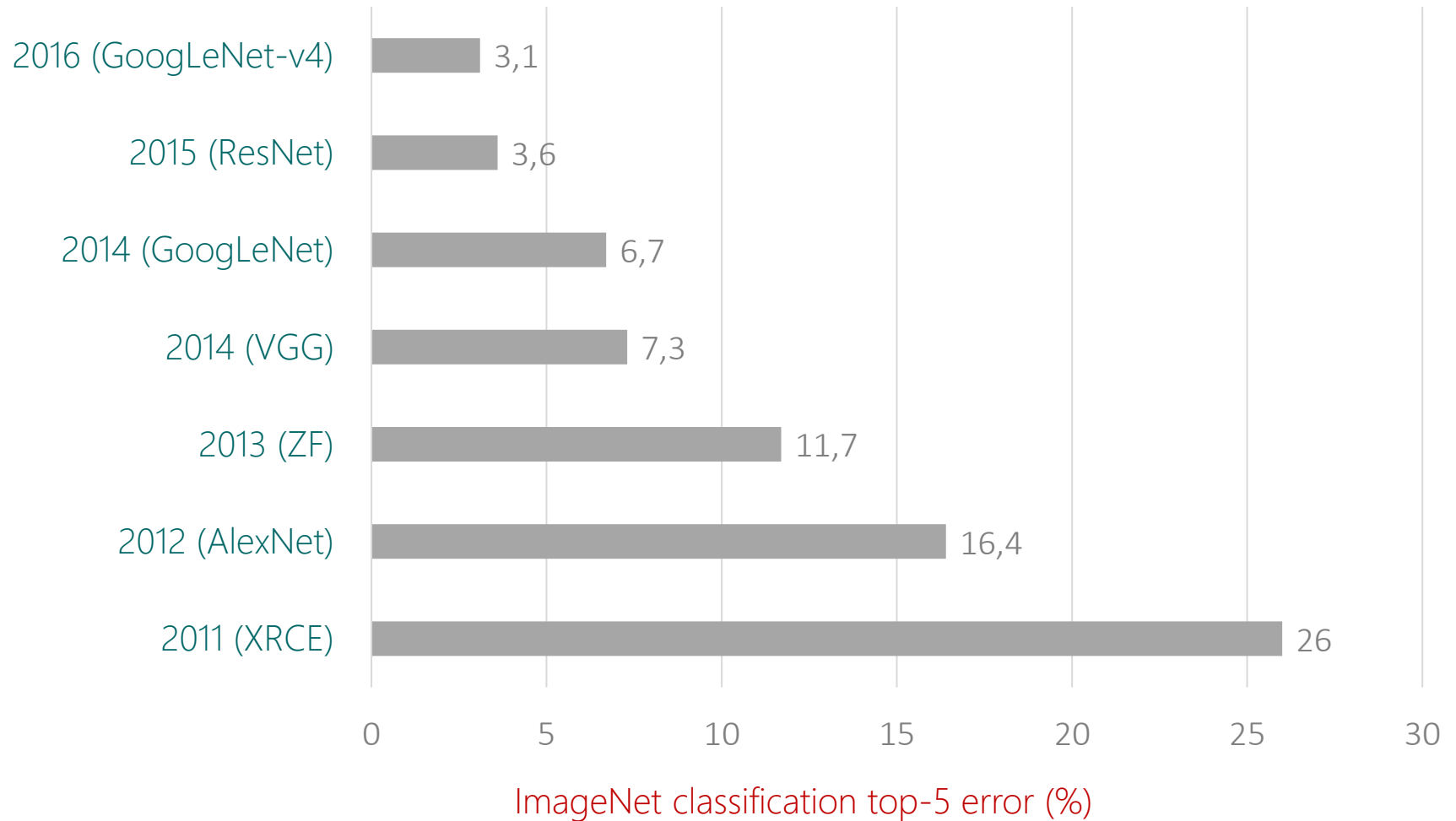
classification (top-5 error)



Krizhevsky, Alex, Ilya Sutskever, and Geoffrey E. Hinton. "Imagenet classification with deep convolutional neural networks." *Advances in neural information processing systems*. 2012.

ImageNet challenge

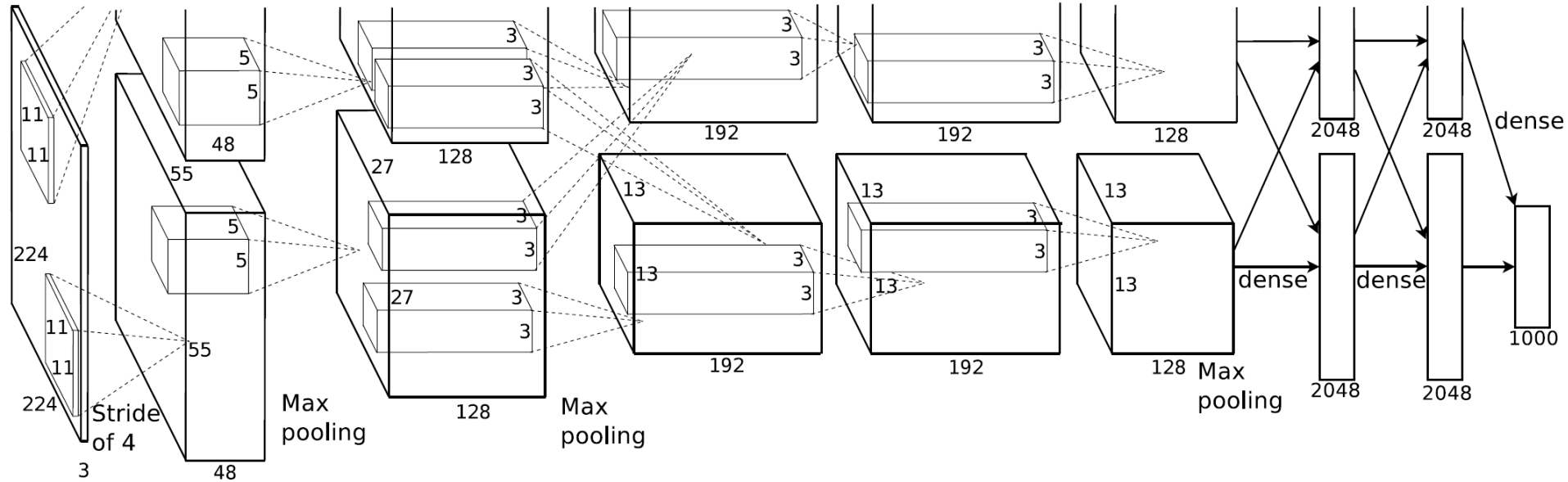
winning models



ImageNet challenge

Large Scale Visual Recognition Challenge (ILSVRC)

AlexNet (winning model in 2012)



Krizhevsky, Alex, Ilya Sutskever, and Geoffrey E. Hinton. "Imagenet classification with deep convolutional neural networks." *Advances in neural information processing systems*. 2012 ([ver](#)).

top-1 error: 38,1%

top-5 error: 16,4%

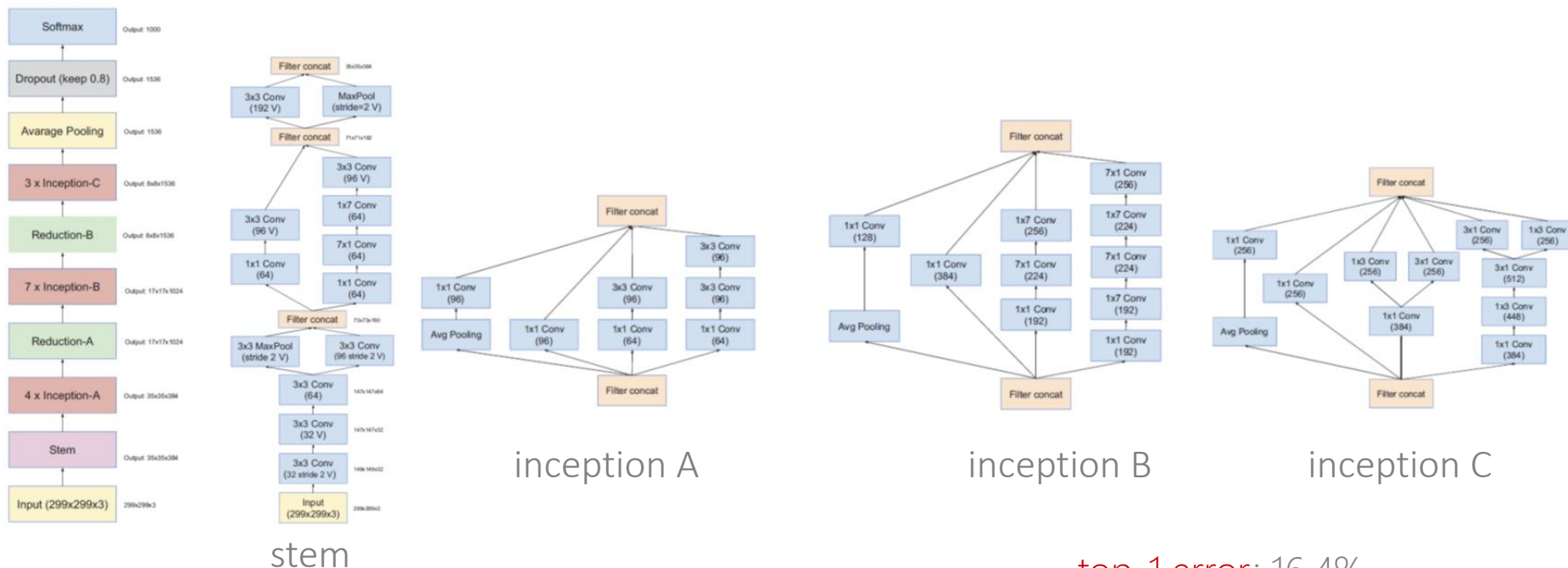
layers: 8

¡60 millones de parámetros!

ImageNet challenge

Large Scale Visual Recognition Challenge (ILSVRC)

Inception v4 (winning model in 2016)



Szegedy, Christian, et al. "Inception-v4, inception-resnet and the impact of residual connections on learning." AAAI. Vol. 4. 2017 ([ver](#)).

top-1 error: 16,4%
top-5 error: 3,1%
layers \approx 170 (trainable)
43 million parameters!

convolution

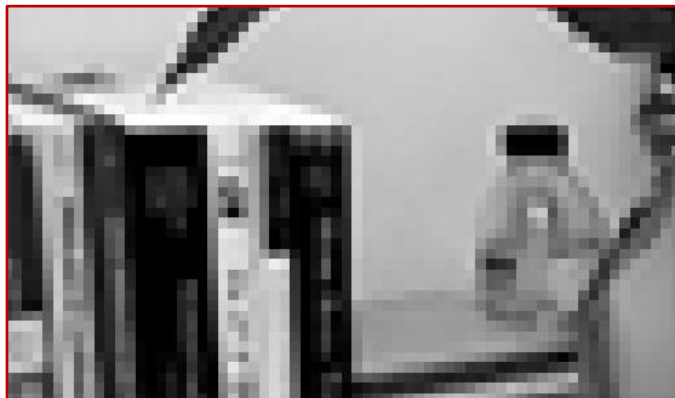
Convolutional networks are simply **neural networks** that use **convolution** in place of general matrix multiplication in at least one of their layers.

convolution

digital representation of images

grayscale images

- array of elements called **image pixels**
- matrix dimensions are called the **image resolution**
- each cell/element/pixel stores a value: its intensity or **gray level**
- intensity/gray levels take integer values between 0 and 255
- dark values are close to 0 (black); light greys, close to 255 (white)

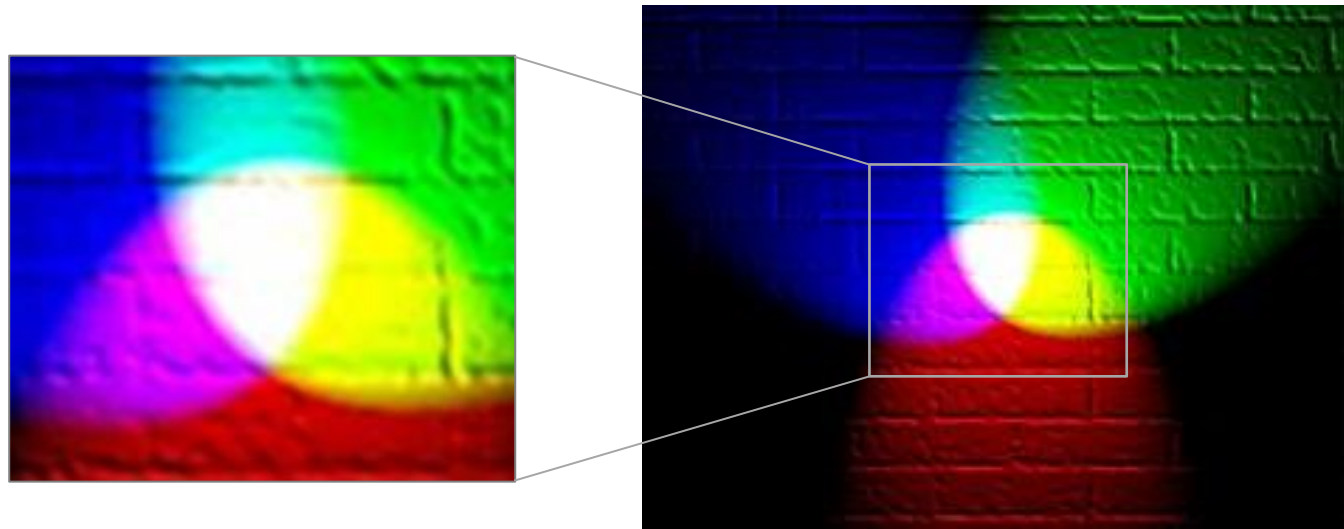


convolution

digital representation of images

color images

- array of elements called **image pixels**
- matrix dimensions are called the **image resolution**
- each pixel stores its intensity in 3 channels/values: Red, Green, Blue
- each R, G or B intensity takes value between 0 and 255
- each final color is the result of combining the R, G or B values



convolution

digital representation of images

let...

- $p_{i,j,k}$ be the intensity value of channel $k \in \{R, G, B\}$ for pixel (i, j)
- for grayscale images, $p_{i,j,k}$ reduces to $p_{i,j}$

normalization [0, 1]

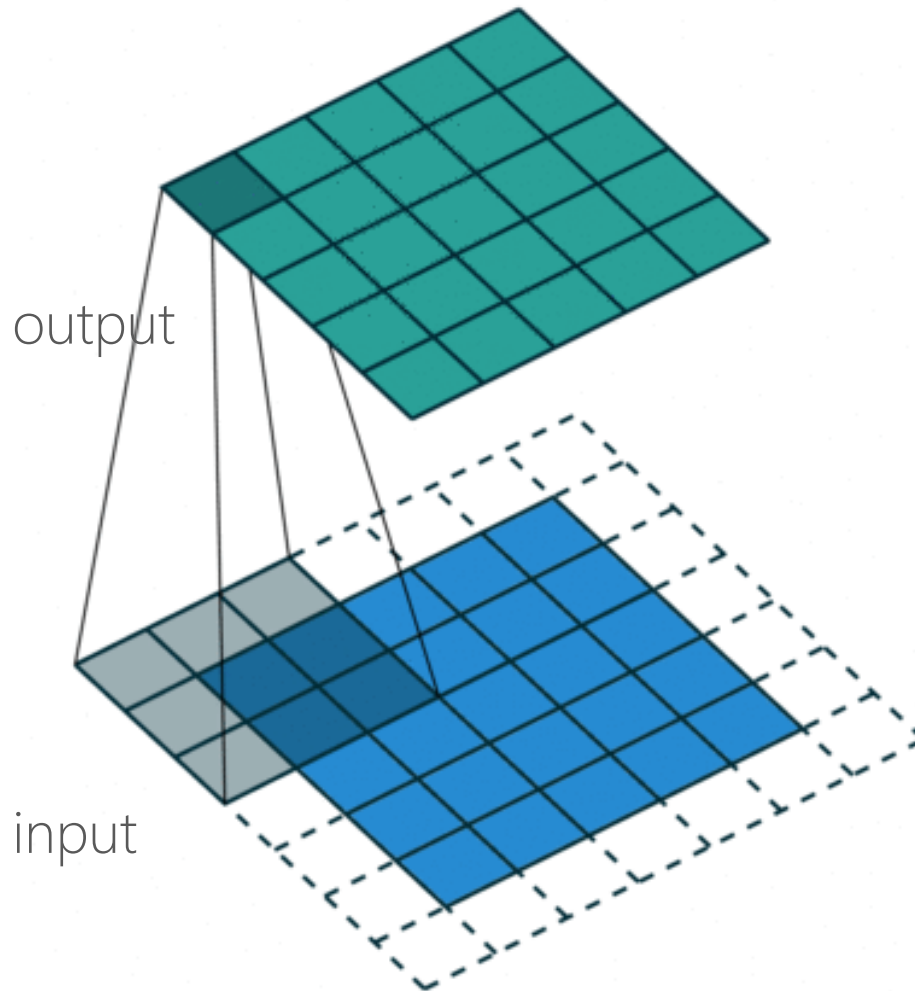
$$\tilde{p}_{i,j,k} = \frac{p_{i,j,k}}{255}$$

normalization [-0.5, 0.5]

$$\tilde{p}_{i,j,k} = \frac{p_{i,j,k}}{255} - 0.5$$

convolution

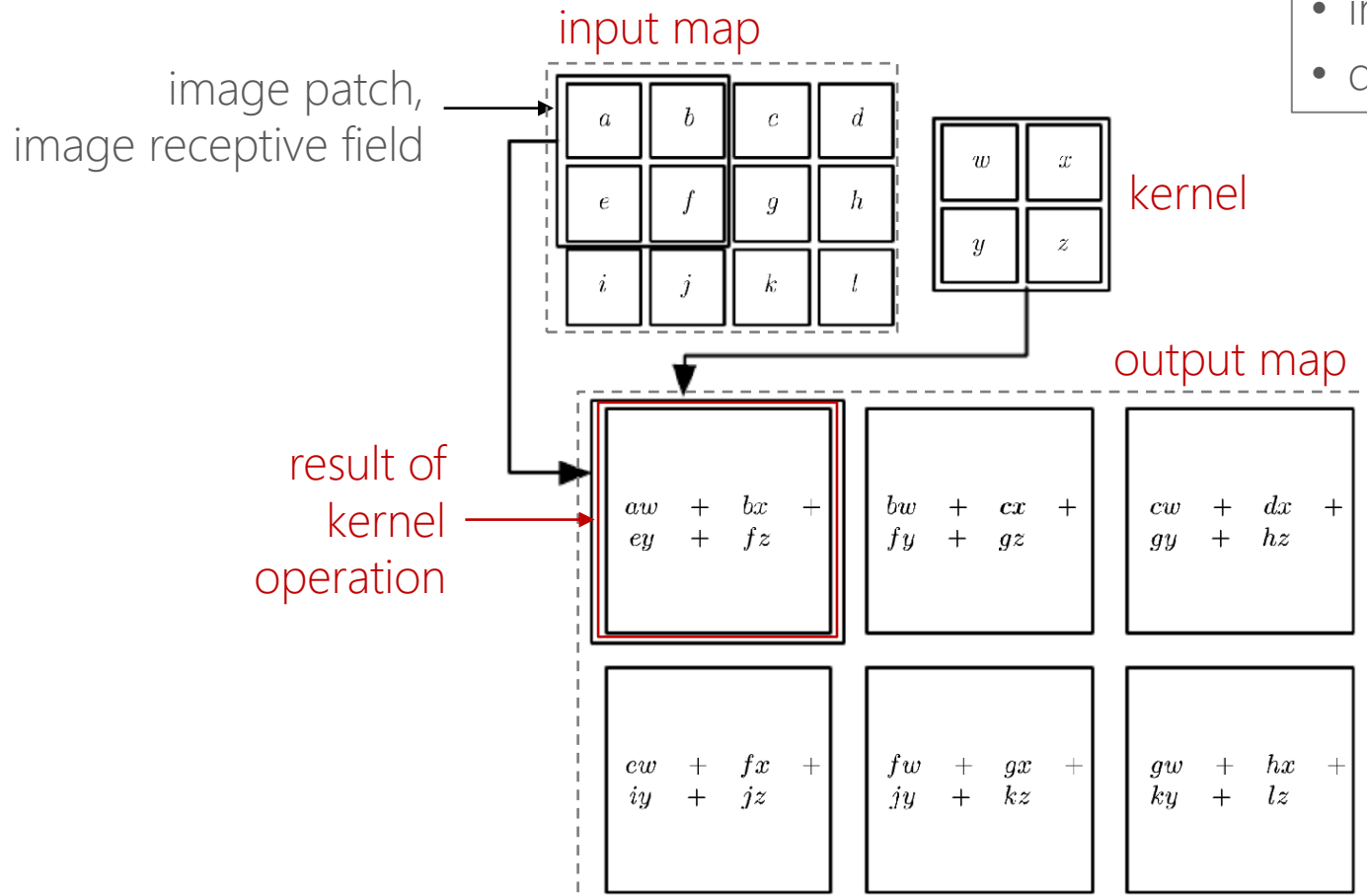
2D convolution



stride = 1
padding = 1

convolution

2D convolution



convolution operation

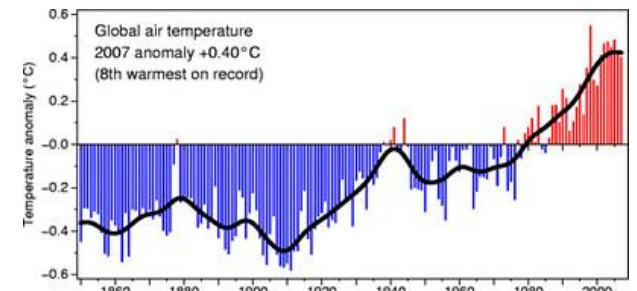
- scalar product
- inner product
- dot product

convolution

2D convolution

a convolution operates on “grid-shaped” data:

- a time series, interpretable as a 1D grid
- an image, interpretable as a 2D/3D grid



Source: flickr.com

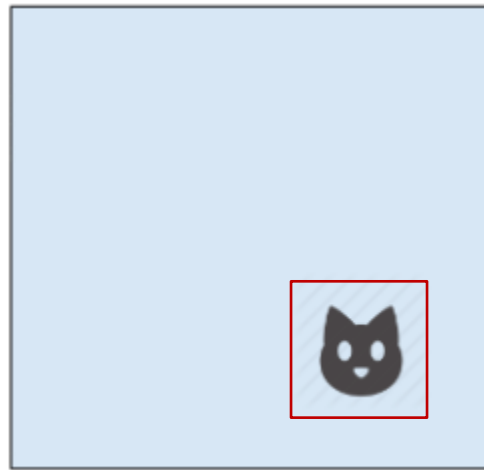


convolution

2D convolution

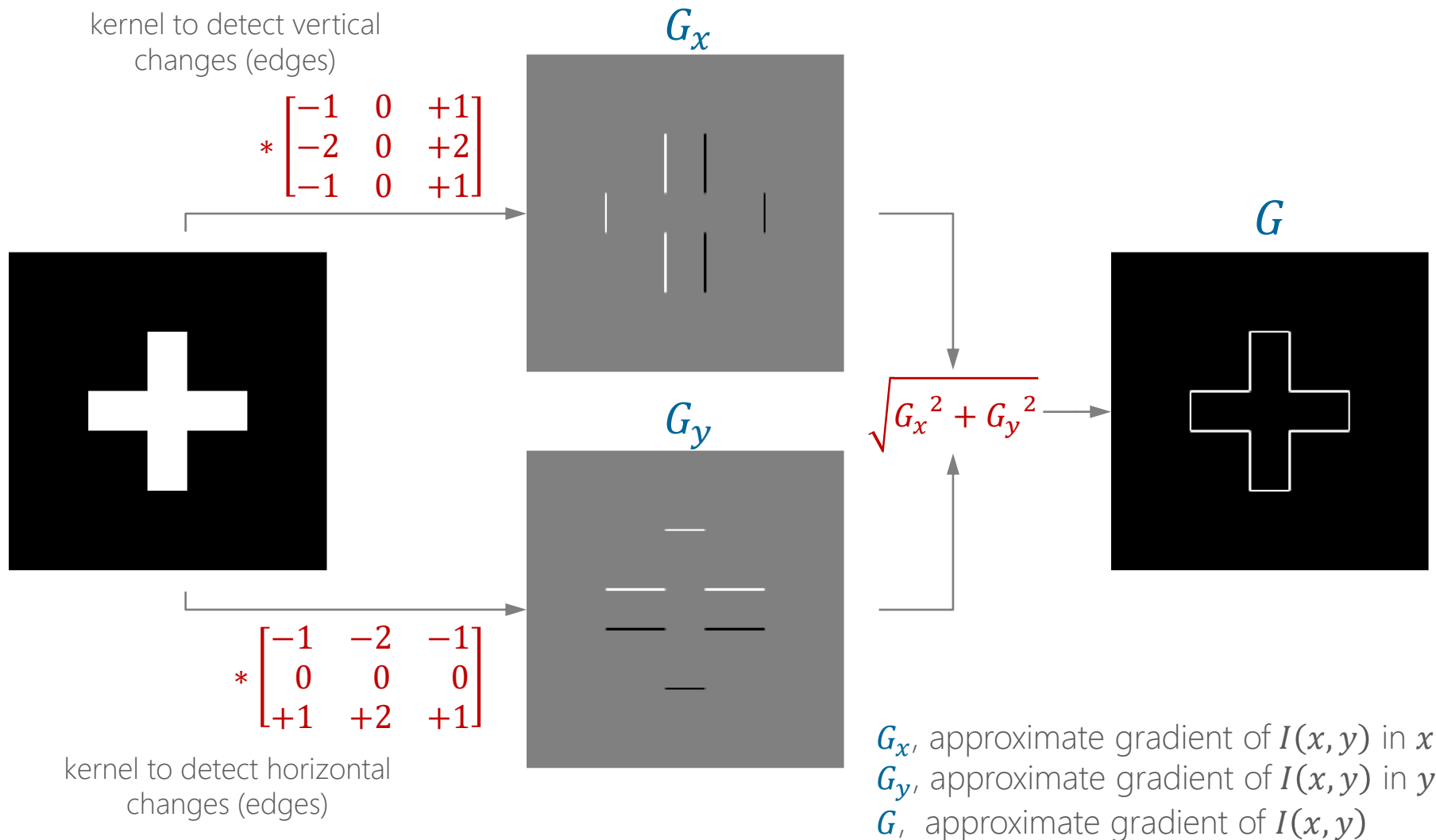
convolution operation properties

- measures spatial correlations in the image receptive field
- translation invariance: detects pattern at any position in the image



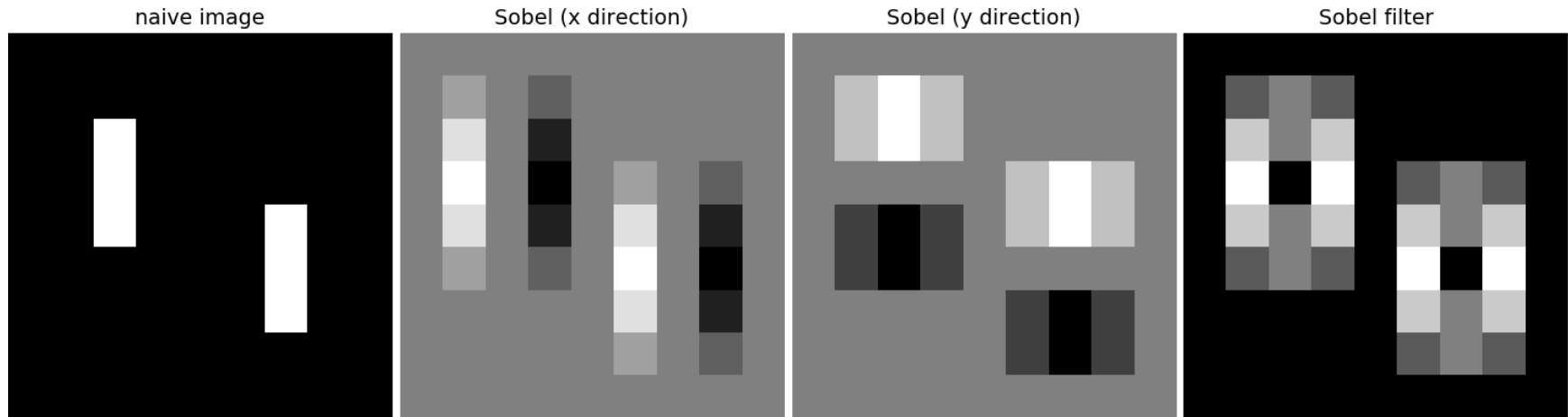
convolution

2D convolution: Sobel-Feldman (edge detection)



convolution

2D convolution: Sobel-Feldman (edge detection)


$$G_x = \begin{bmatrix} 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & -1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 3 & 0 & -3 & 0 & 0 & 0 & 0 & 0 \\ 0 & 4 & 0 & -4 & 0 & 1 & 0 & -1 & 0 \\ 0 & 3 & 0 & -3 & 0 & 3 & 0 & -3 & 0 \\ 0 & 1 & 0 & -1 & 0 & 4 & 0 & -4 & 0 \\ 0 & 0 & 0 & 0 & 0 & 3 & 0 & -3 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 & 0 & -1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \end{bmatrix} \quad G_y = \begin{bmatrix} 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 2 & 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 2 & 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 & 2 & 1 & 0 \\ 0 & -1 & -2 & -1 & 0 & 1 & 2 & 1 & 0 \\ 0 & -1 & -2 & -1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & -1 & -2 & -1 & 0 \\ 0 & 0 & 0 & 0 & 0 & -1 & -2 & -1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \end{bmatrix}$$

image

G

$$\begin{bmatrix} 0. & 0. & 0. & 0. & 0. & 0. & 0. & 0. & 0. \\ 0. & 0. & 0. & 0. & 0. & 0. & 0. & 0. & 0. \\ 0. & 0. & 1. & 0. & 0. & 0. & 0. & 0. & 0. \\ 0. & 0. & 1. & 0. & 0. & 0. & 0. & 0. & 0. \\ 0. & 0. & 1. & 0. & 0. & 0. & 1. & 0. & 0. \\ 0. & 0. & 0. & 0. & 0. & 0. & 1. & 0. & 0. \\ 0. & 0. & 0. & 0. & 0. & 0. & 1. & 0. & 0. \\ 0. & 0. & 0. & 0. & 0. & 0. & 0. & 0. & 0. \\ 0. & 0. & 0. & 0. & 0. & 0. & 0. & 0. & 0. \end{bmatrix}$$

numeric output

Do you notice the invariance to the position of the pattern?

```
[ [0. 0. 0. 0. 0. 0. 0. 0. 0. ]
  [0. 1.4 2. 1.4 0. 0. 0. 0. 0. ]
  [0. 3.2 2. 3.2 0. 0. 0. 0. 0. ]
  [0. 4. 0. 4. 0. 1.4 2. 1.4 0. ]
  [0. 3.2 2. 3.2 0. 3.2 2. 3.2 0. ]
  [0. 1.4 2. 1.4 0. 4. 0. 4. 0. ]
  [0. 0. 0. 0. 0. 3.2 2. 3.2 0. ]
  [0. 0. 0. 0. 0. 1.4 2. 1.4 0. ]
  [0. 0. 0. 0. 0. 0. 0. 0. 0. ] ]
```

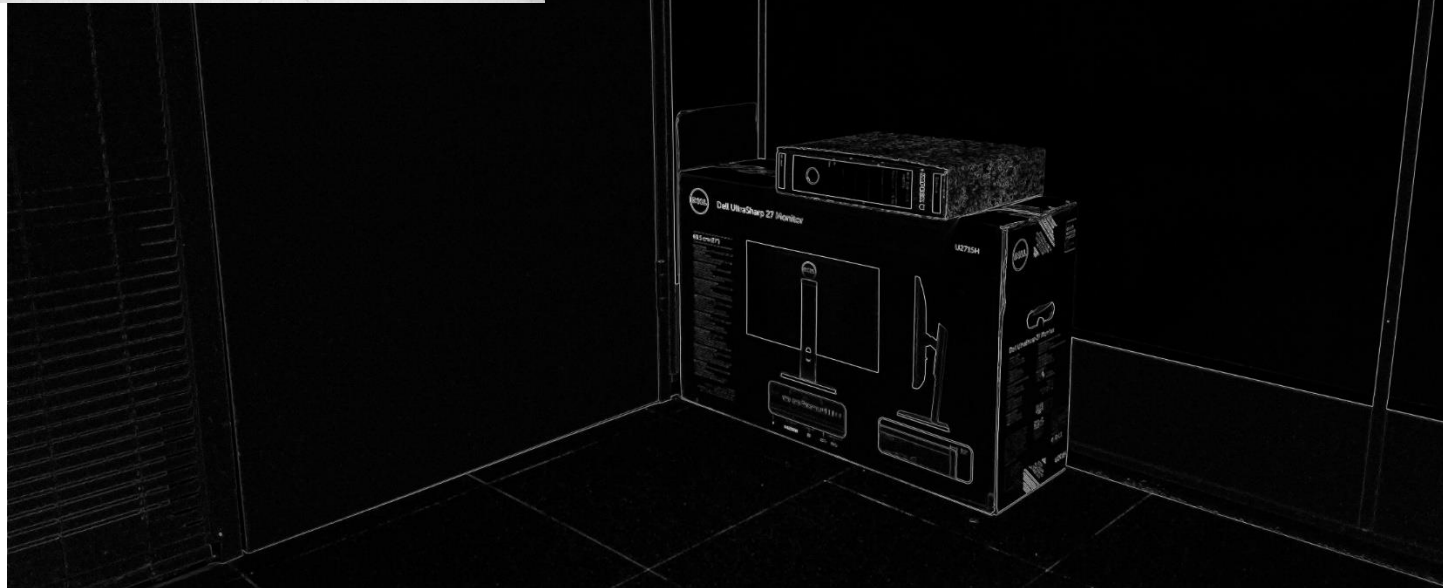

convolution

2D convolution: Sobel-Feldman (edge detection)



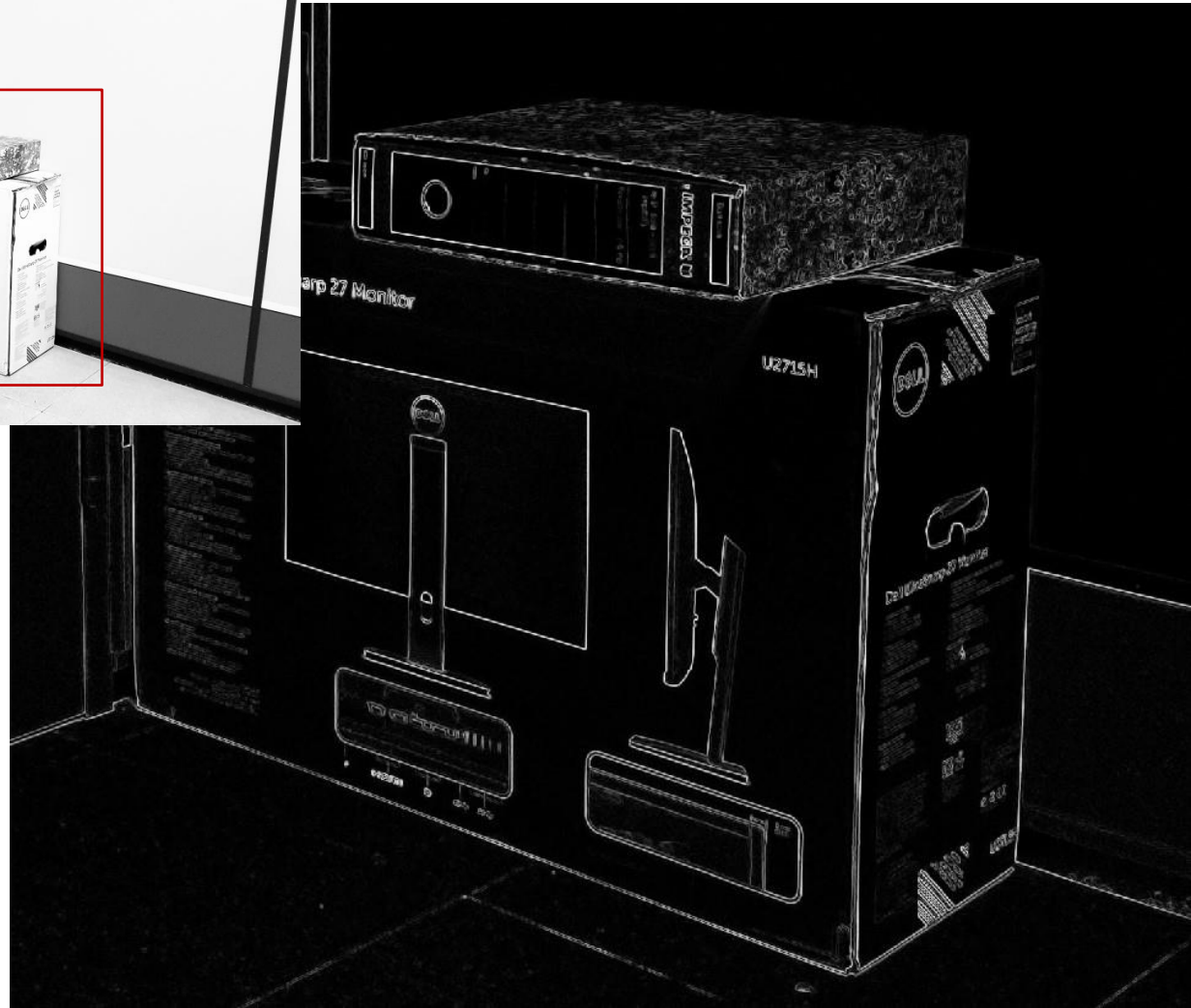
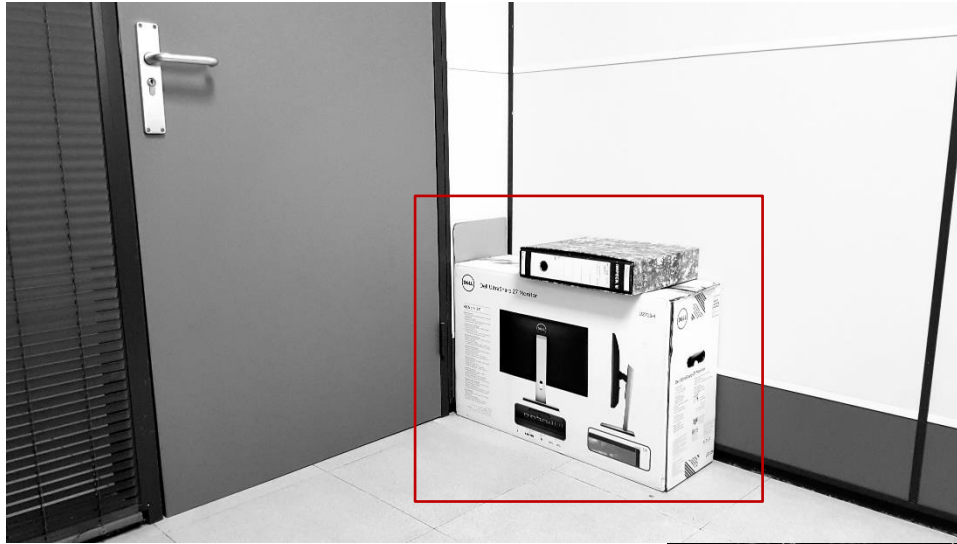
easy task

corner of my office with mostly straight edges and a few textures



convolution

2D convolution: Sobel-Feldman (edge detection)



convolution

2D convolution: Sobel-Feldman (edge detection)



difficult task

corner of my office with edges of
diverse geometry and rich textures



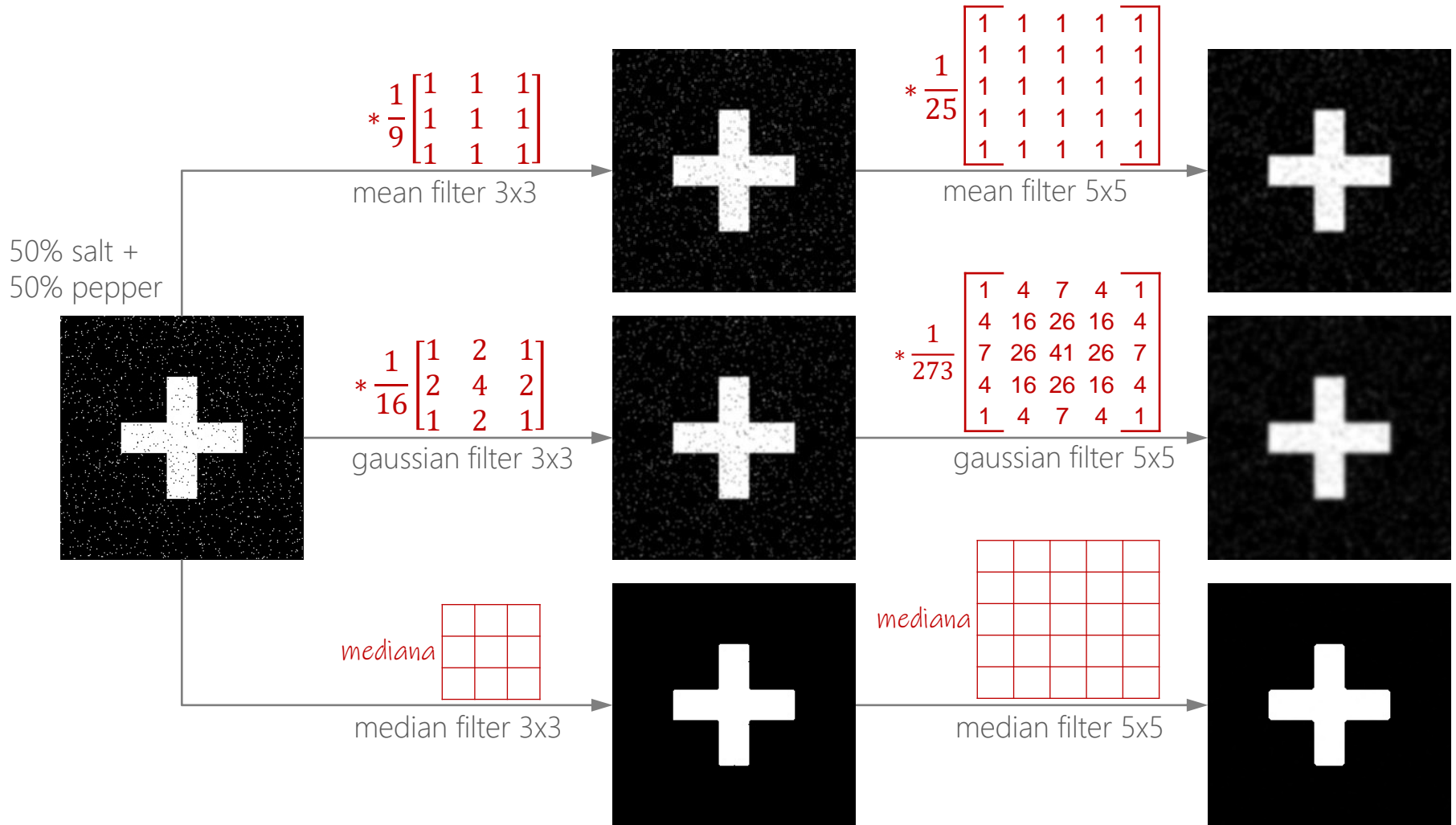
convolution

2D convolution: Sobel-Feldman (edge detection)



convolution

2D convolution: smoothing images



convolution

2D convolution: smoothing images



50% salt +
50% pepper

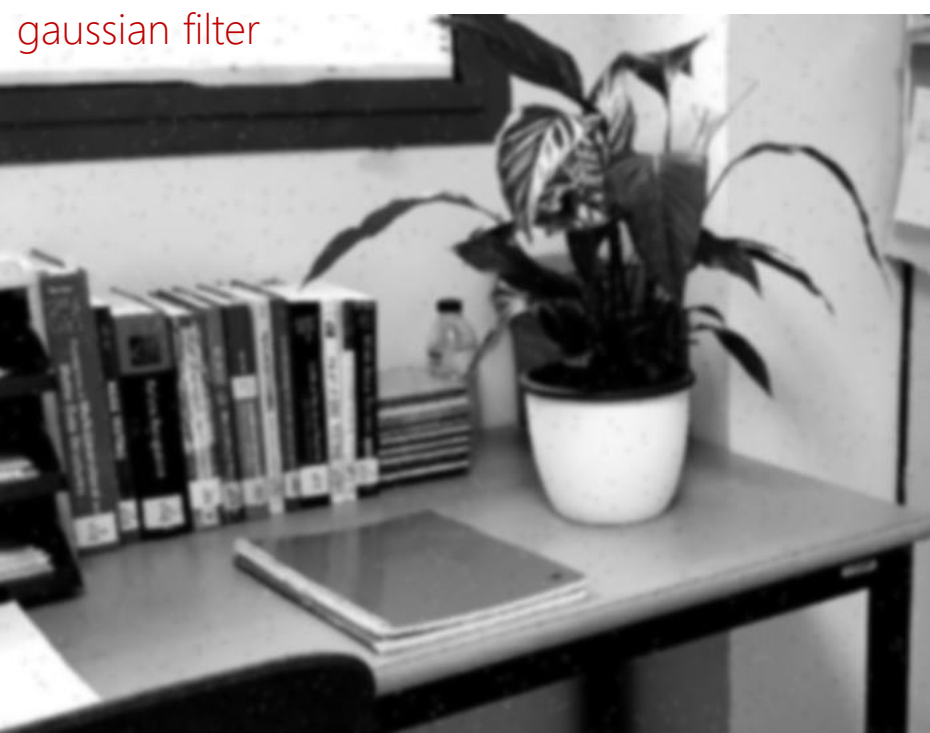


median filter



convolution
2D convolution:
smoothing images

gaussian filter

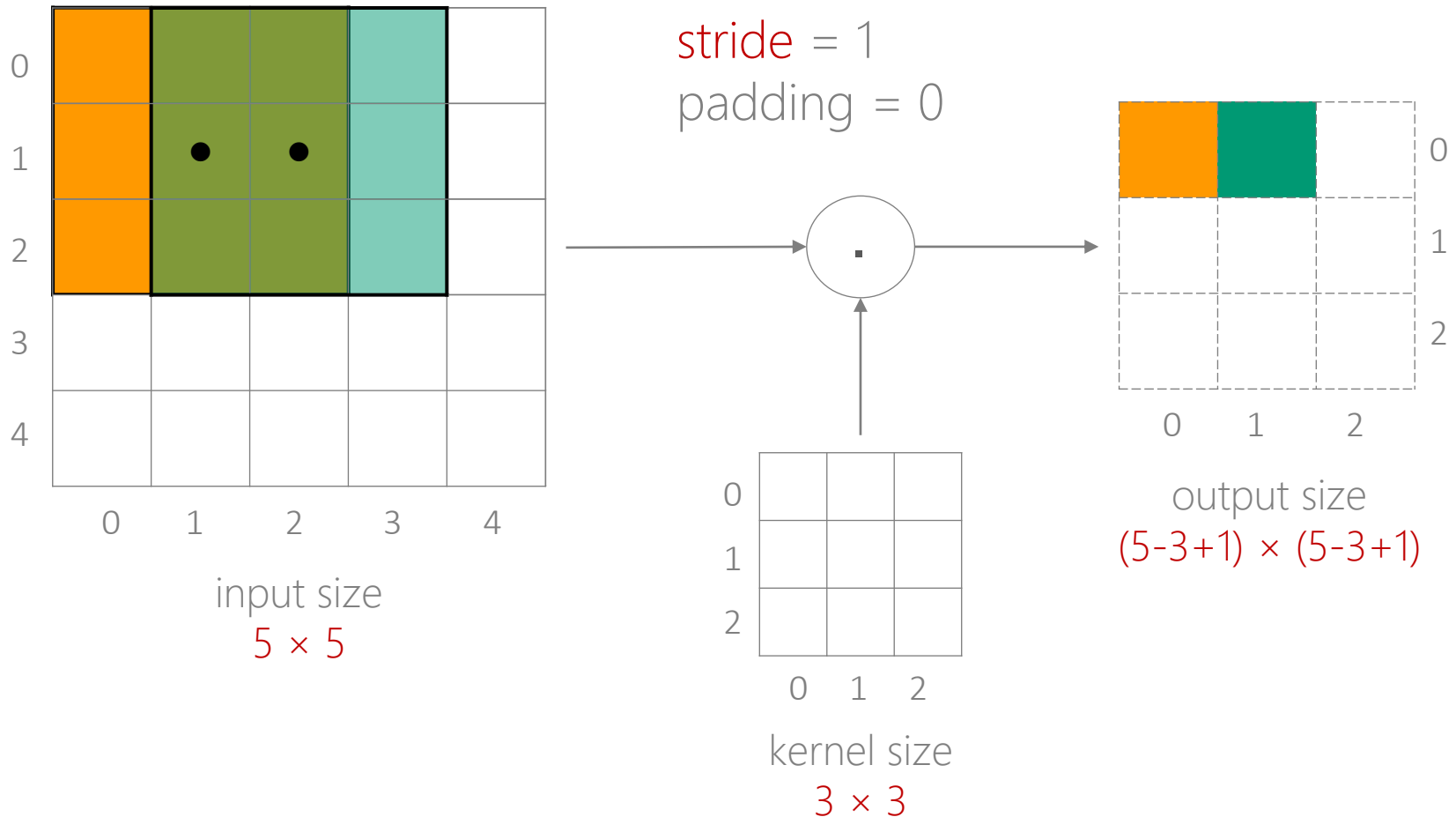


noisy source image



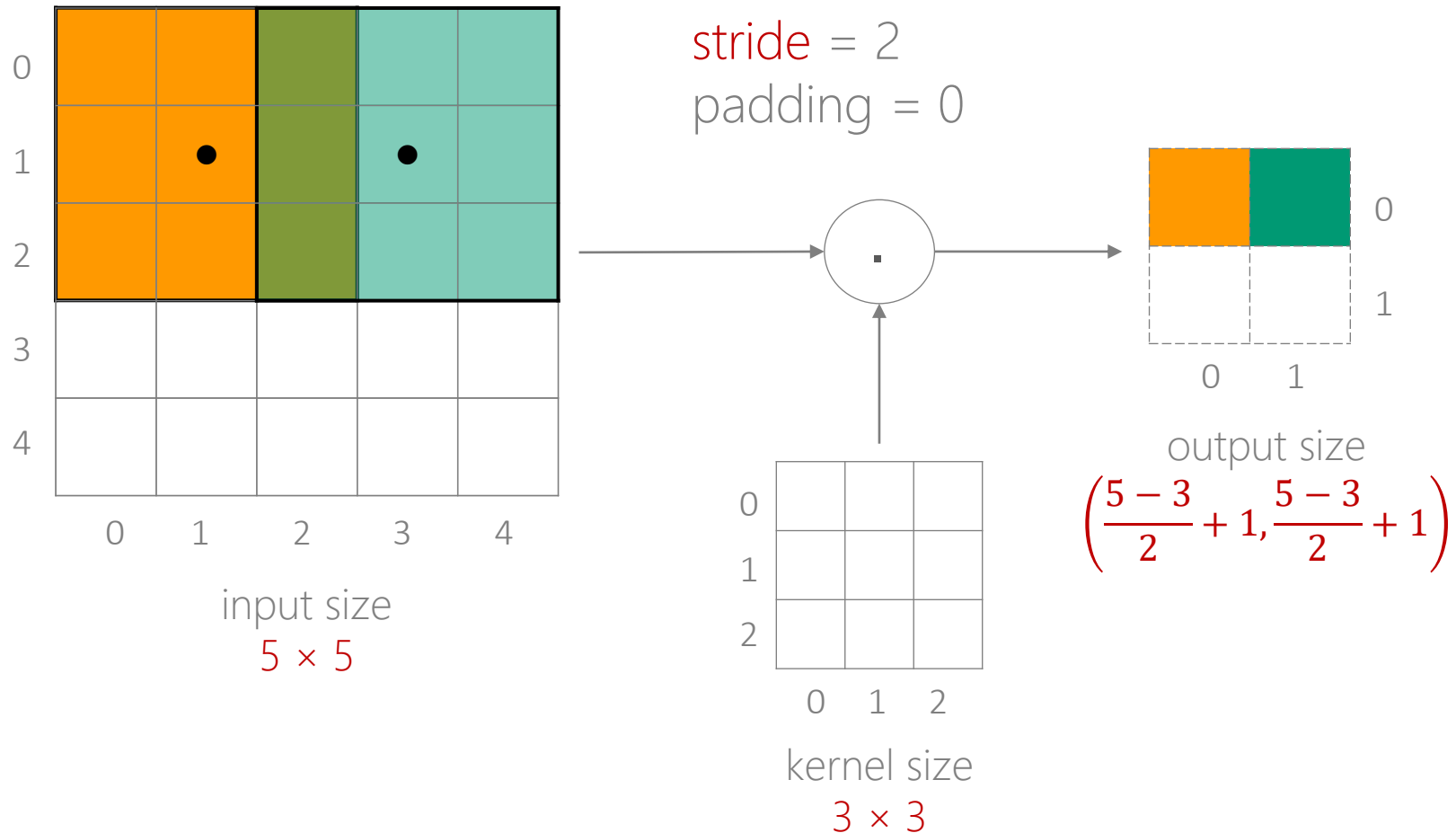
convolution

2D convolution: *stride, padding*



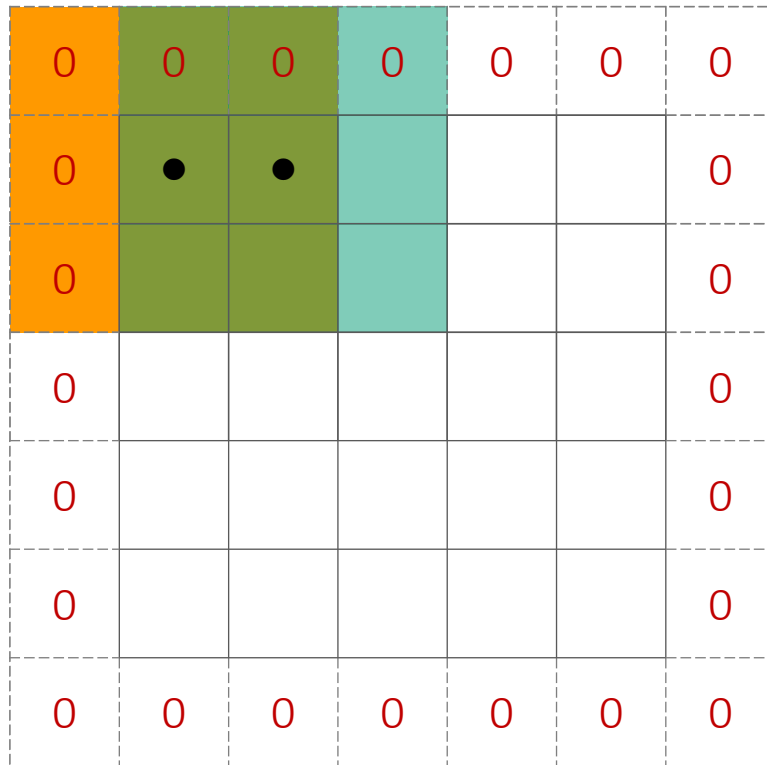
convolution

2D convolution: *stride, padding*

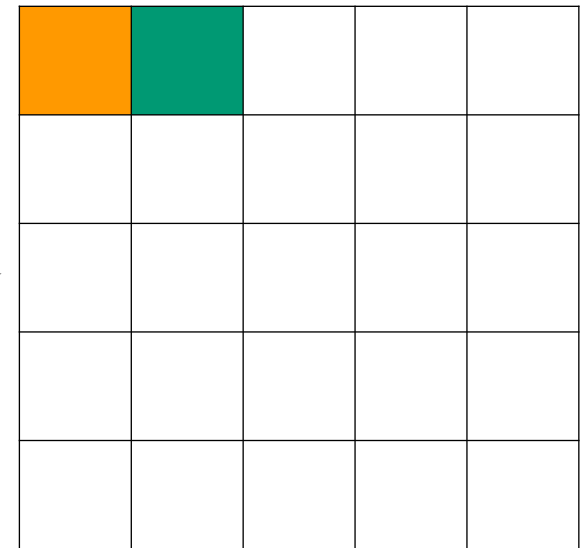
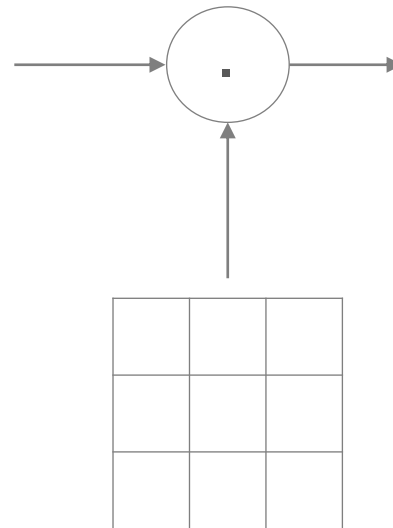


convolution

2D convolution: *stride, padding*



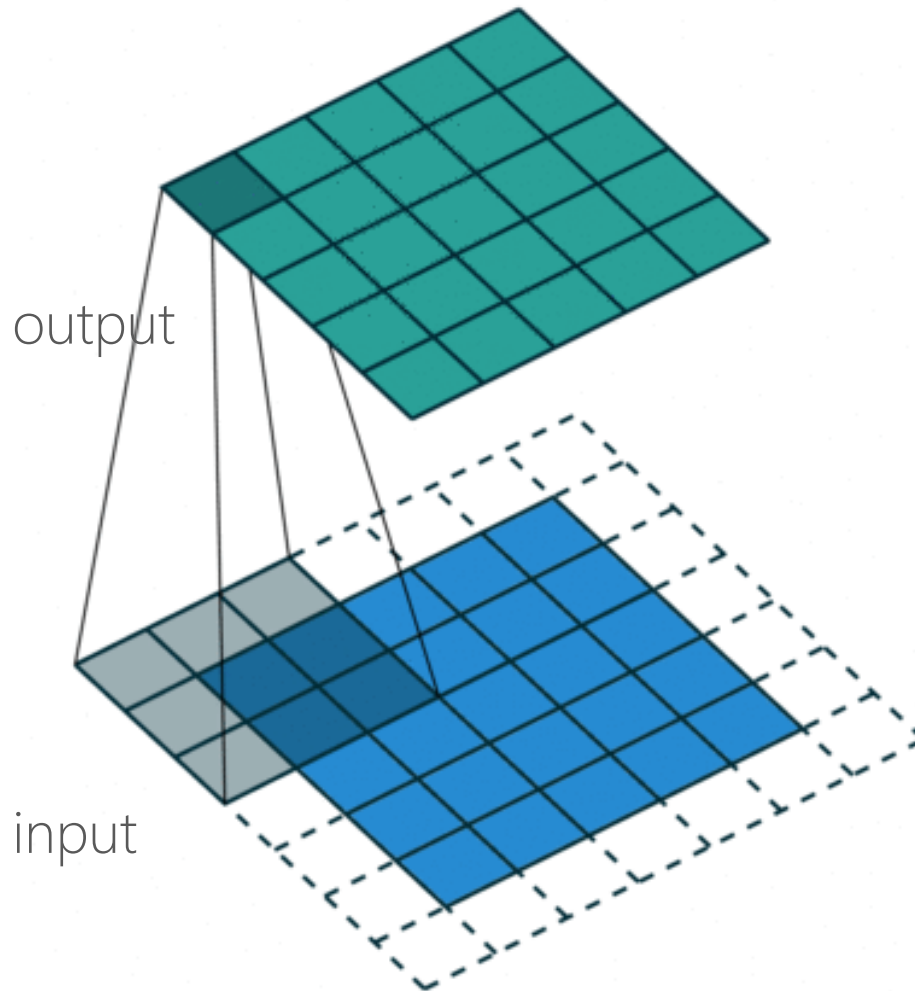
stride = 1
padding = 1



$$padding \leq \left\lceil \frac{\text{tamaño kernel}}{2} \right\rceil$$

convolution

2D convolution: *stride, padding*

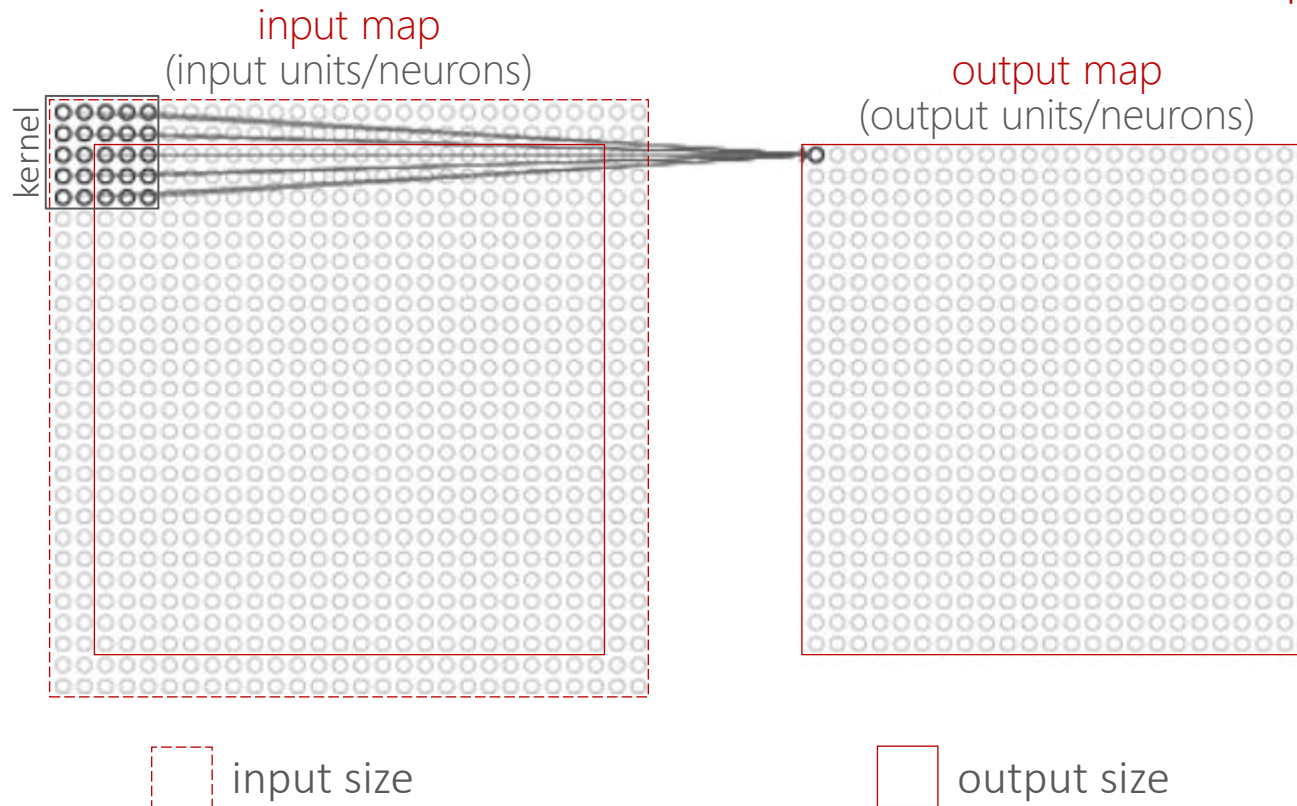


stride = 1
padding = 1

convolution

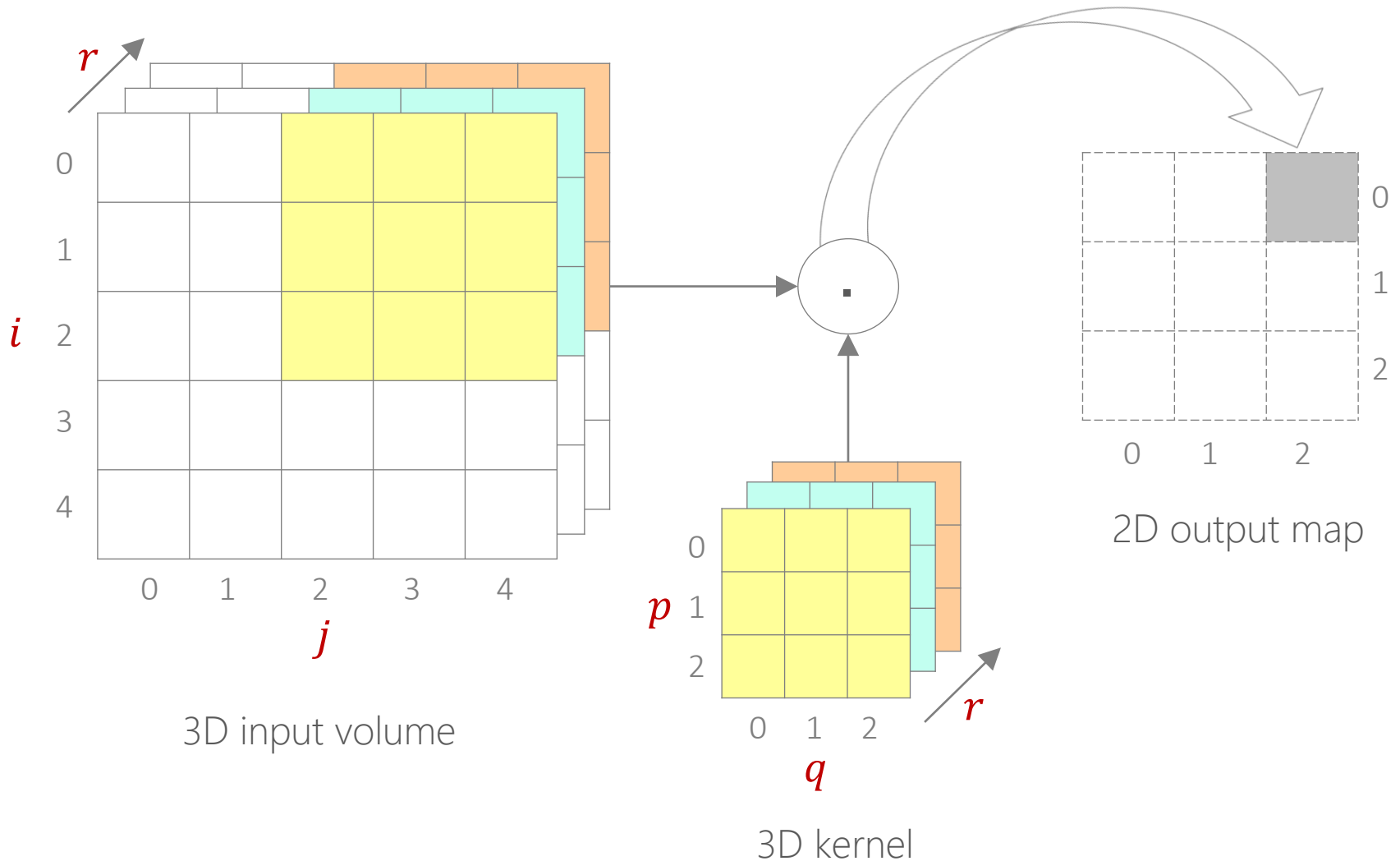
2D convolution

stride = 1
padding = 2



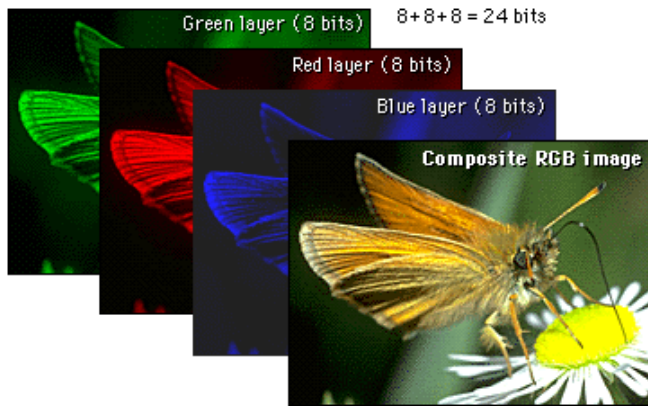
convolution

2D convolution (on 3D maps)



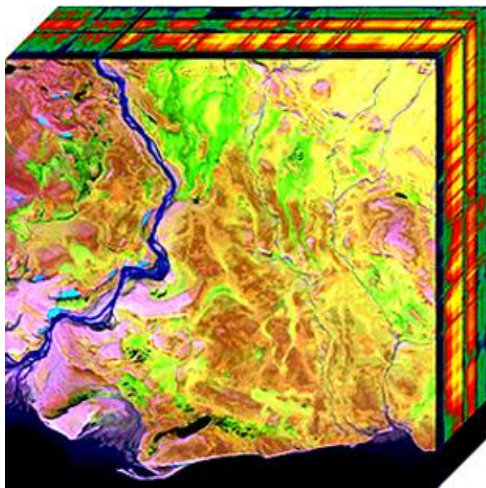
convolution

2D convolution (on multi-channel images)



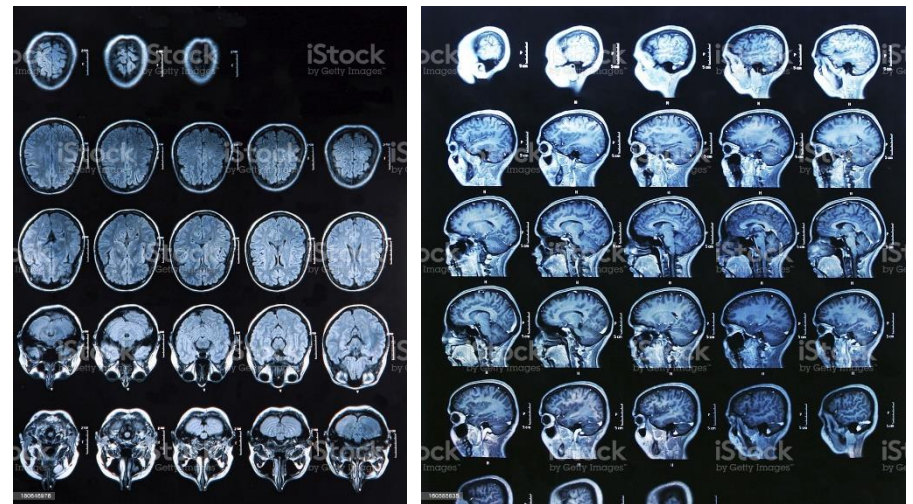
Source: https://webstyleguide.com/wsg1/graphics/display_primer.html

multi-channel images



multi-spectral image

Source: wikipedia.org



Magnetic Resonance Image (MRI)

Source: istockphoto.com

convolution

2D convolution

Let...

I , be a 3D input volume

W , be a 3D kernel (filter, operator)

O , be the 2D output map

$$O(i, j) = (W \cdot I)(i, j) = \sum_p \sum_q \sum_r I_{i+p, j+q, r} \cdot W_{p, q, r} + b$$

convolution

2D convolution

Let...

I , be a 3D input volume

W^c , the c -th 3D kernel of a set of C kernels

O , be the 3D output volume

$$O(i, j, c) = (W^c \cdot I)(i, j, c) = \sum_p \sum_q \sum_r I_{i+p, j+q, r} \cdot W_{p, q, r}^c + b^c$$

convolution

2D convolution: practical exercise

example above: with a $5 \times 5 \times 3$ input map, $3 \times 3 \times 3$ kernel, padding 0, and stride 1, the output map size was $3 \times 3 \times 1$.

what would have been the dimension of the output map with a kernel...

- $2 \times 2 \times 3 \rightarrow ?$
- $4 \times 4 \times 3 \rightarrow ?$
- $5 \times 5 \times 3 \rightarrow ?$

generalization: given an $H \times W \times C$ input, and $K \times K \times C$ kernel, padding P , and stride S , what would be the size of the output map?

- $K \times K \times C \rightarrow ?$

convolution

2D convolution: practical exercise

example above: with a $5 \times 5 \times 3$ input map, $3 \times 3 \times 3$ kernel, padding 0, and stride 1, the output map size was $3 \times 3 \times 1$.

what would have been the dimension of the output map with a kernel...

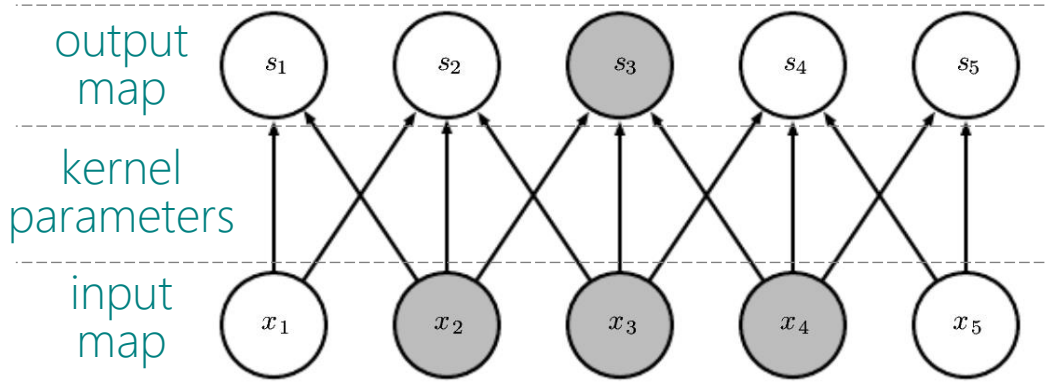
- $2 \times 2 \times 3 \rightarrow 4 \times 4 \times 1$
- $4 \times 4 \times 3 \rightarrow 2 \times 2 \times 1$
- $5 \times 5 \times 3 \rightarrow 1 \times 1 \times 1$

generalization: given an $H \times W \times C$ input, and $K \times K \times C$ kernel, padding P , and stride S , what would be the size of the output map?

- $K \times K \times C \rightarrow \left(\frac{H - K + 2P}{S} + 1, \frac{W - K + 2P}{S} + 1 \right)$

convolution

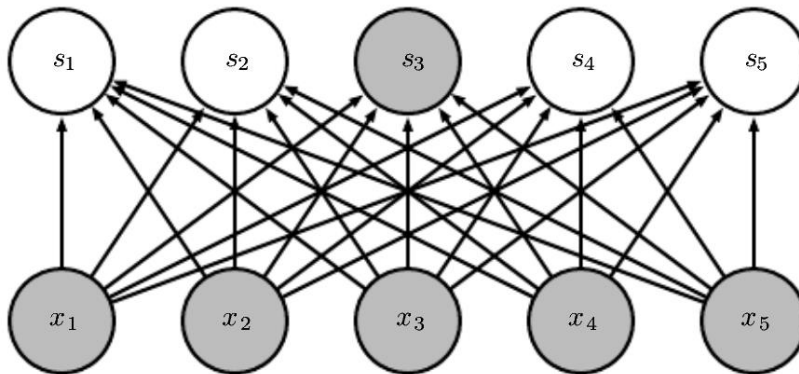
convolutional connection pattern vs dense pattern



CNN

Convolutional
Neural
Network

sparse connectivity
(when kernel is smaller than input)



FCN

Fully
Connected
Network

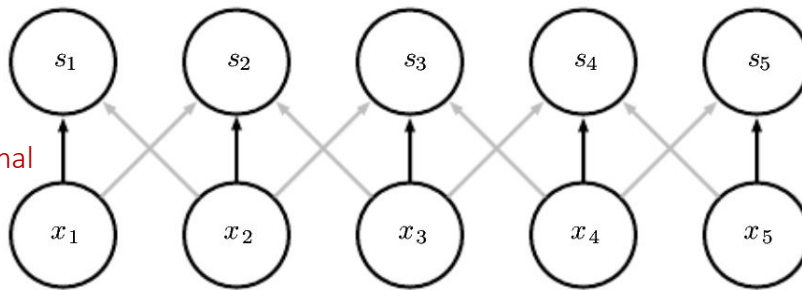
dense connectivity

convolution

convolutional connection pattern vs dense pattern

CNN

Convolutional
Neural
Network

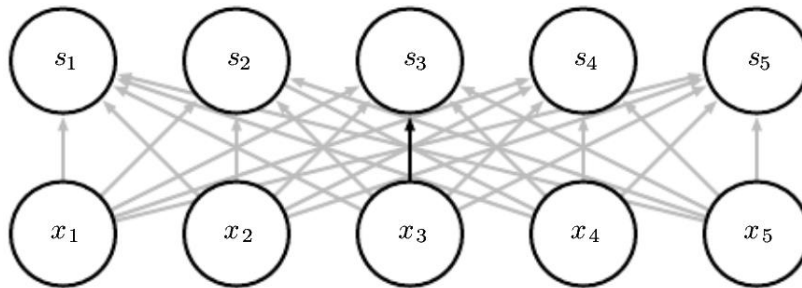


shared parameters

- the no. param. depends on the kernel
- fewer parameters are stored
- fewer calculations are made

FCN

Fully
Connected
Network



one parameter per connection

- the number of parameters depends on the units of the connected layers

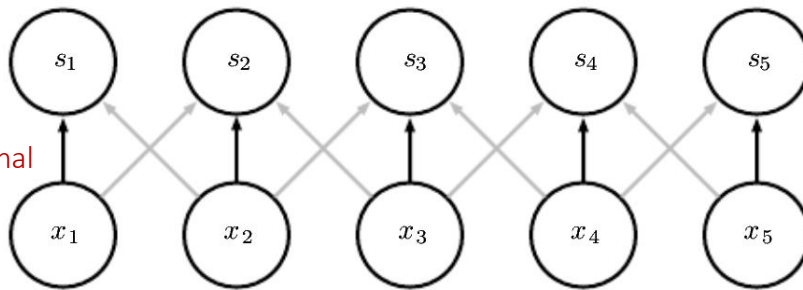
convolution

convolutional connection pattern vs dense pattern

What if you change the size of the network input?

CNN

Convolutional
Neural
Network

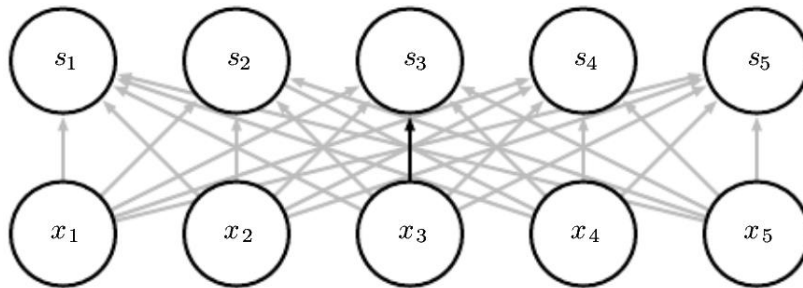


it keeps working

connections/parameters do not depend on input size

FCN

Fully
Connected
Network



it stops working

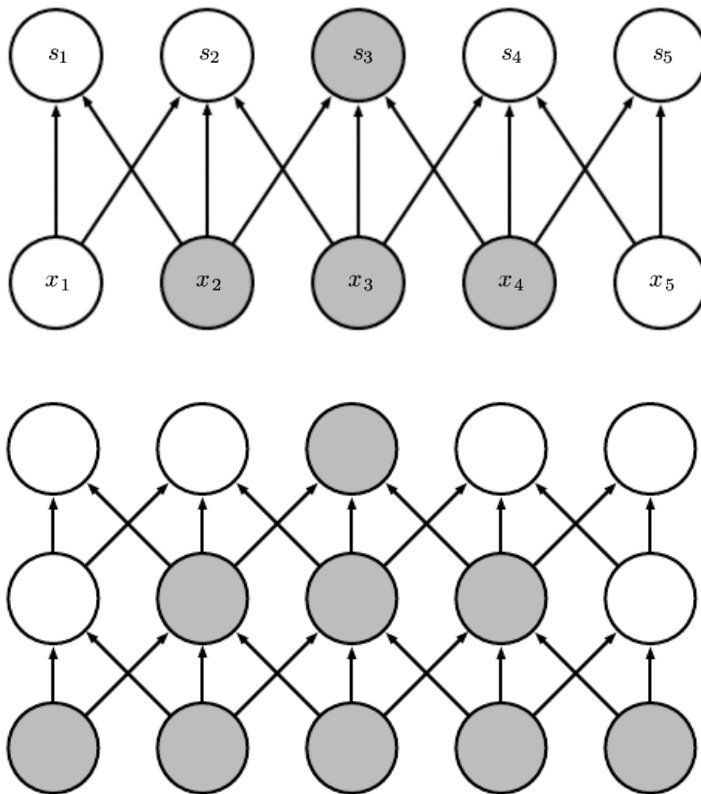
connections/parameters DO depend on input size

convolution

convolutional connection pattern vs dense pattern

CNN

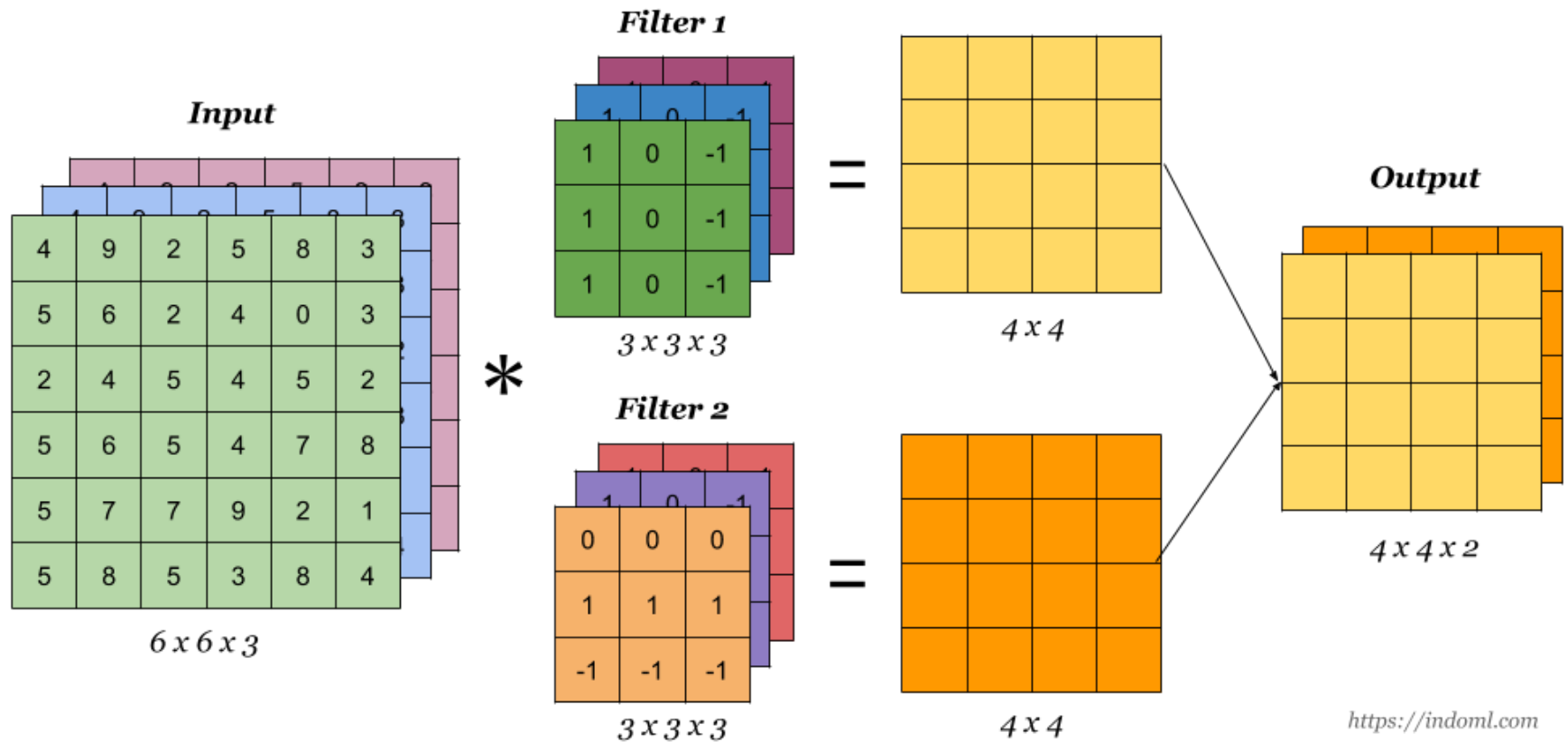
Convolutional
Neural
Network



- x_2, x_3, x_4 (input units) are known as the **receptive field** of s_3
- deeper layer units are related to larger **receptive fields**
- units in very deep layers may be indirectly connected to all or most of the input image

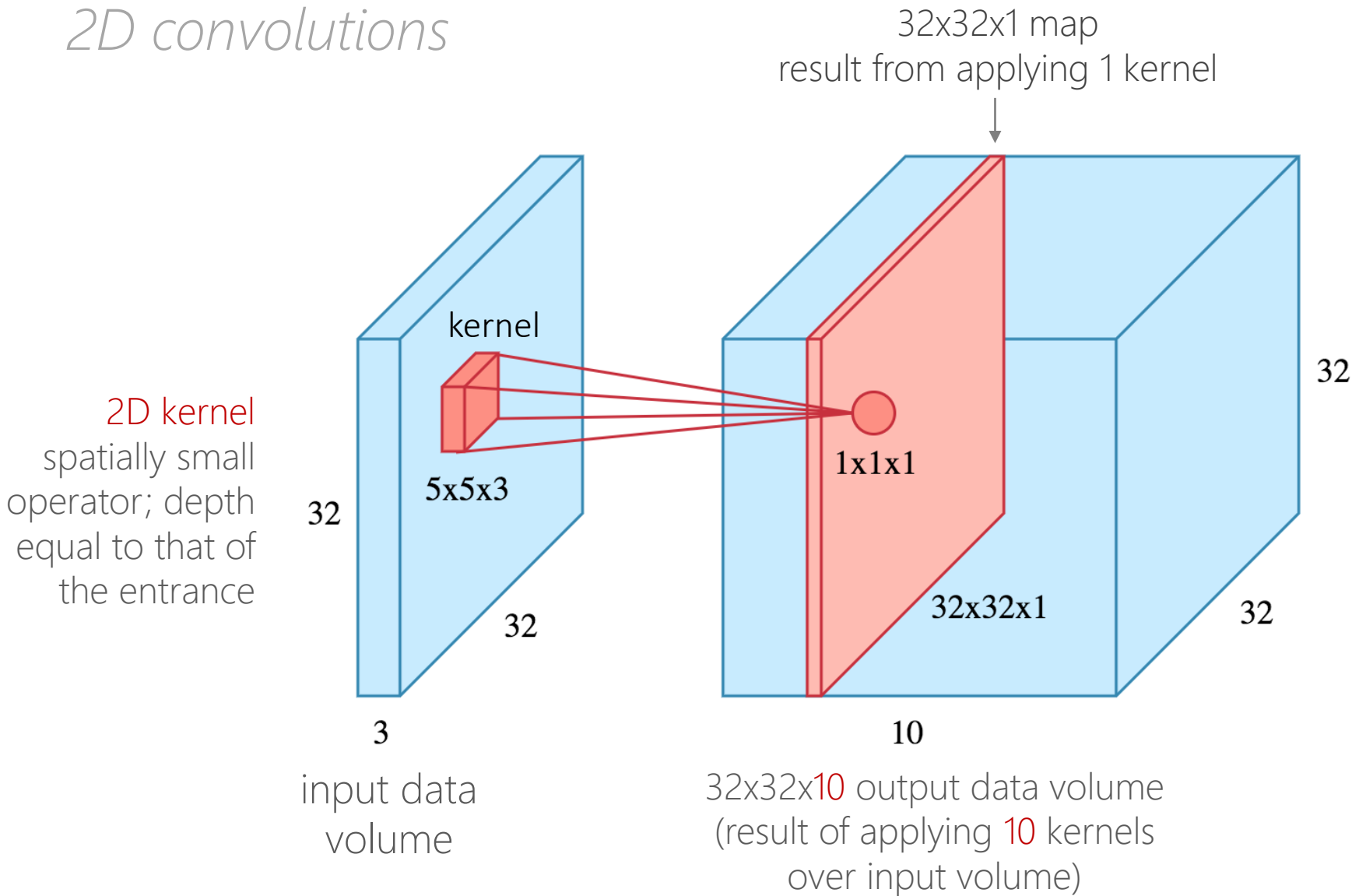
convolutional layer

2D convolutions



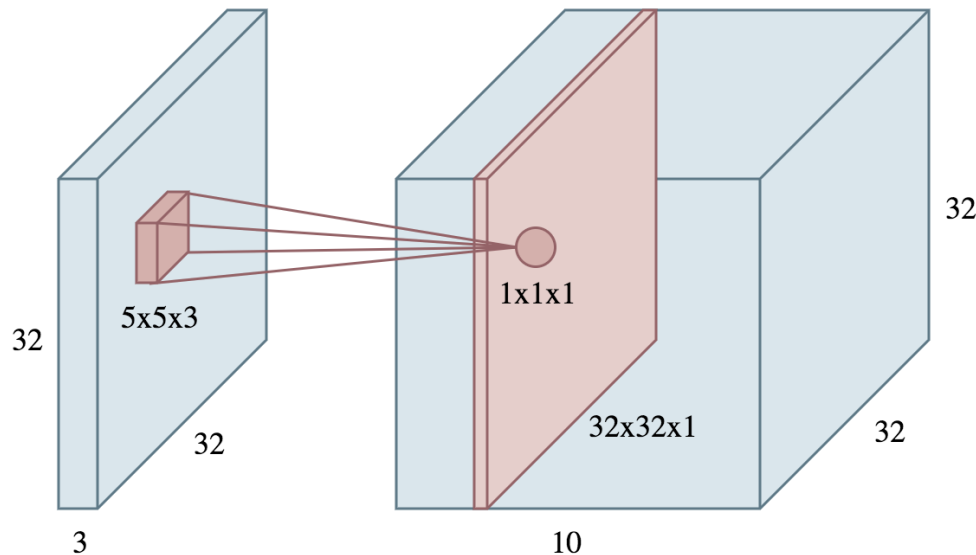
convolutional layer

2D convolutions



convolutional layer

2D convolutions

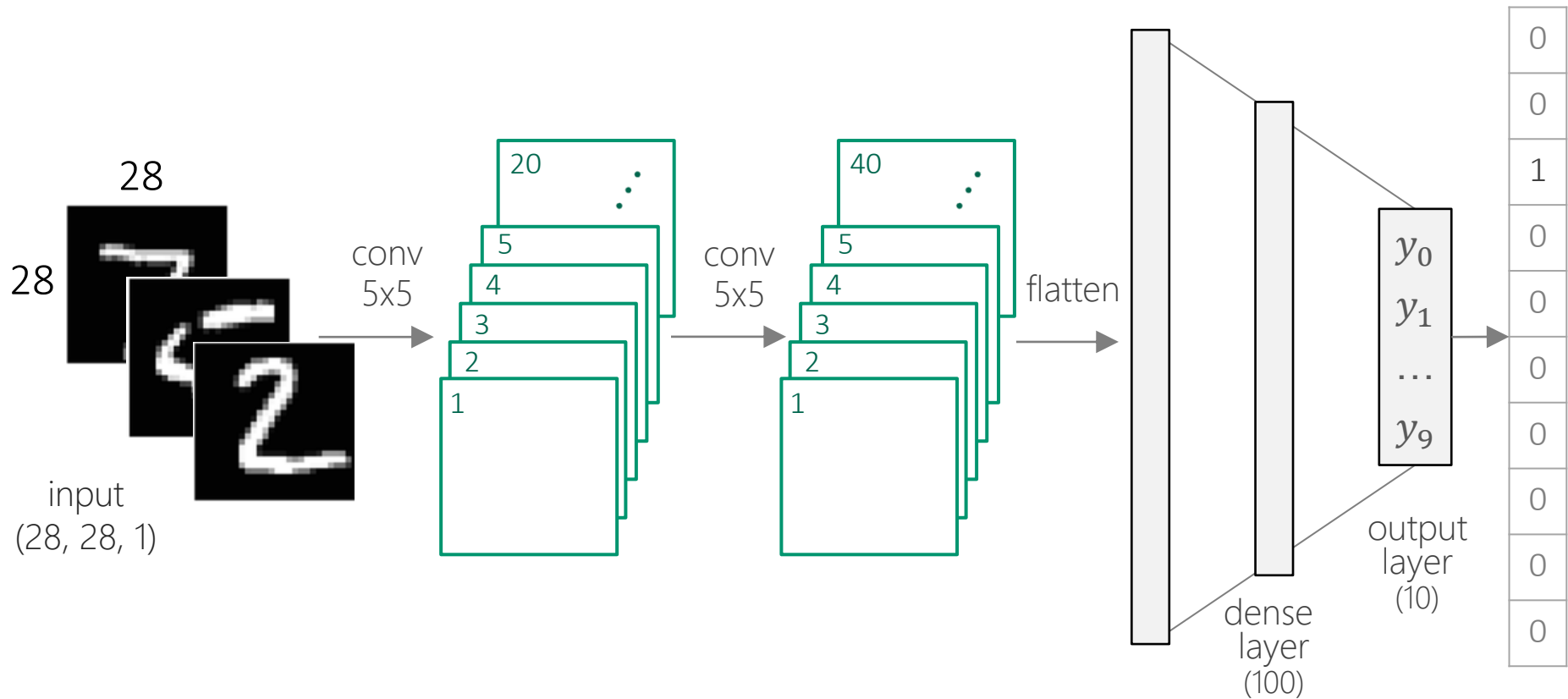


the kernel moves over the input volume
and convolves with each position

- **operator**: $5 \times 5 \times 3$ kernel
- **input** (to the convolution operation): $5 \times 5 \times 3$ block of the input volume
- **operation**: scalar product
- **result**: scalar ($1 \times 1 \times 1$)
- **kernel trainable params** = $75 + 1$
- **kernel depth** = input volumen depth
- **output volumen depth** = number of kernels

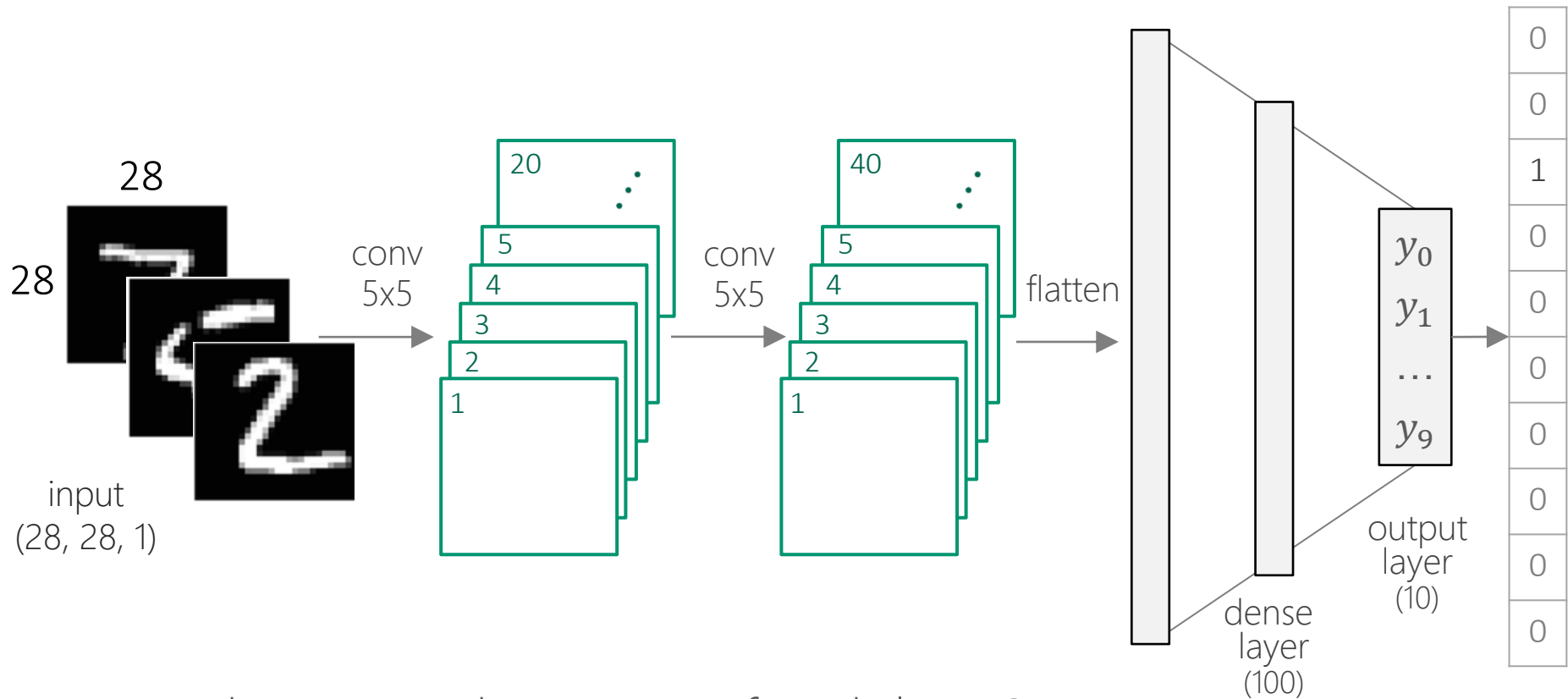
case study

convolutional neural network for the MNIST task



case study

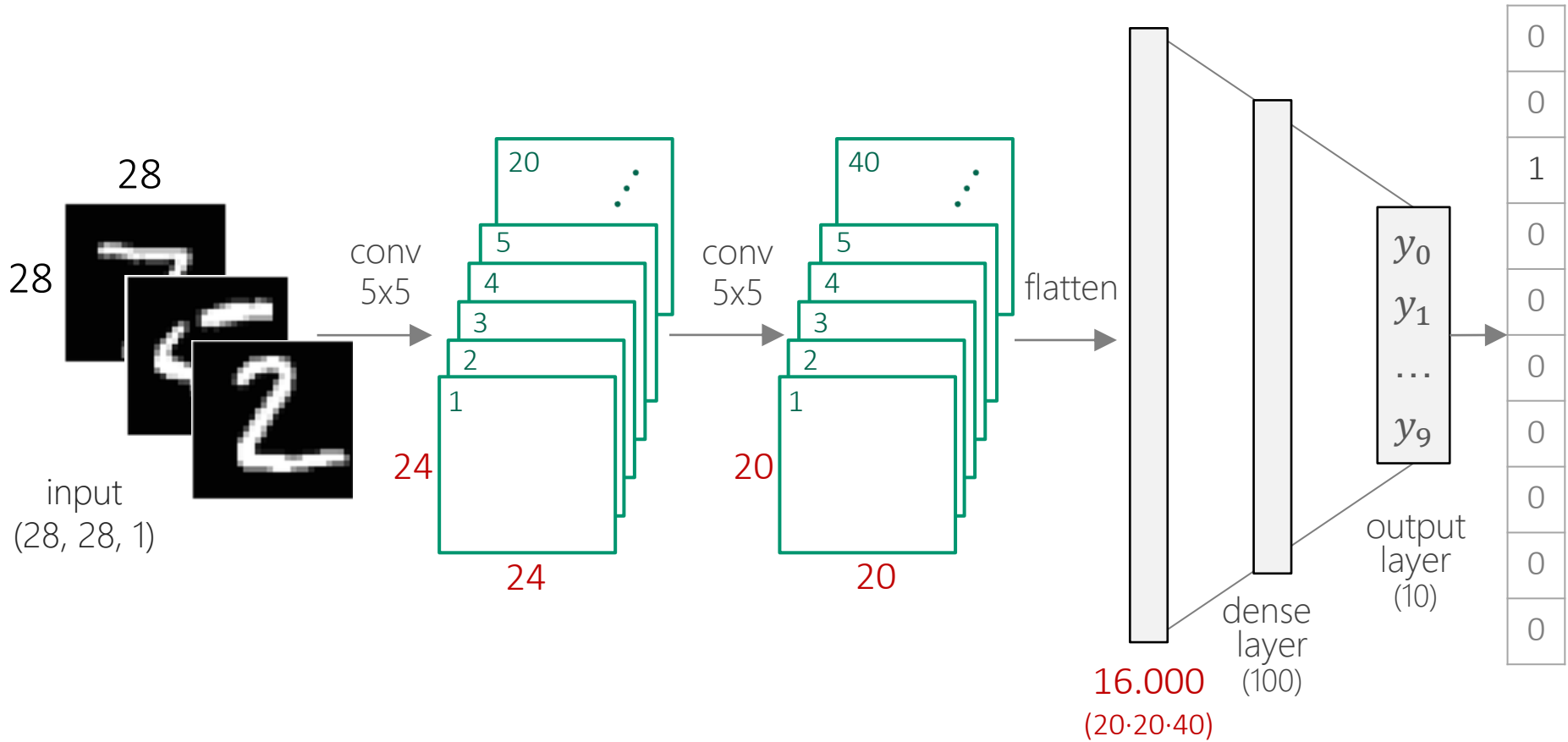
convolutional neural network for the MNIST task



What **size** is the **output** of each layer?

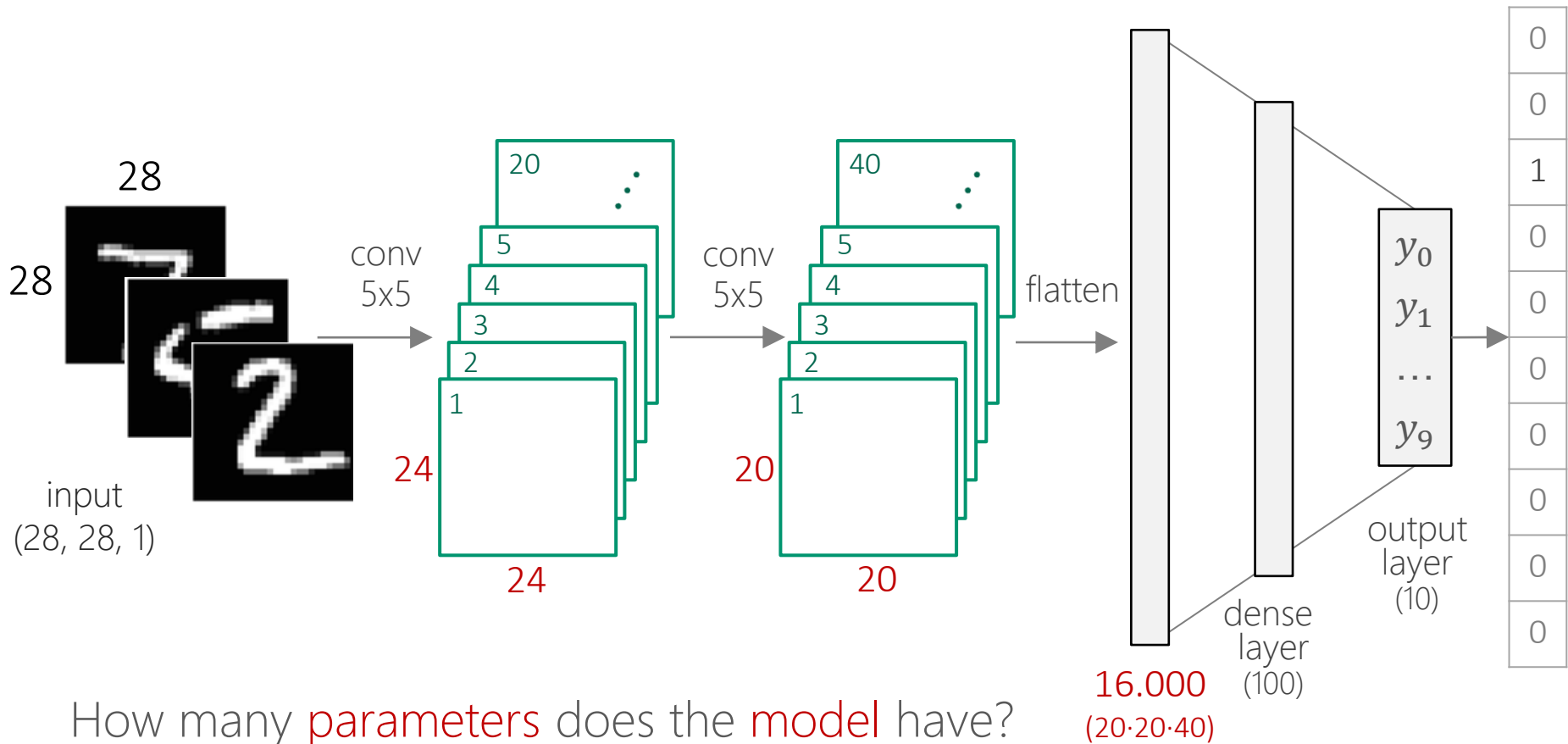
case study

convolutional neural network for the MNIST task



case study

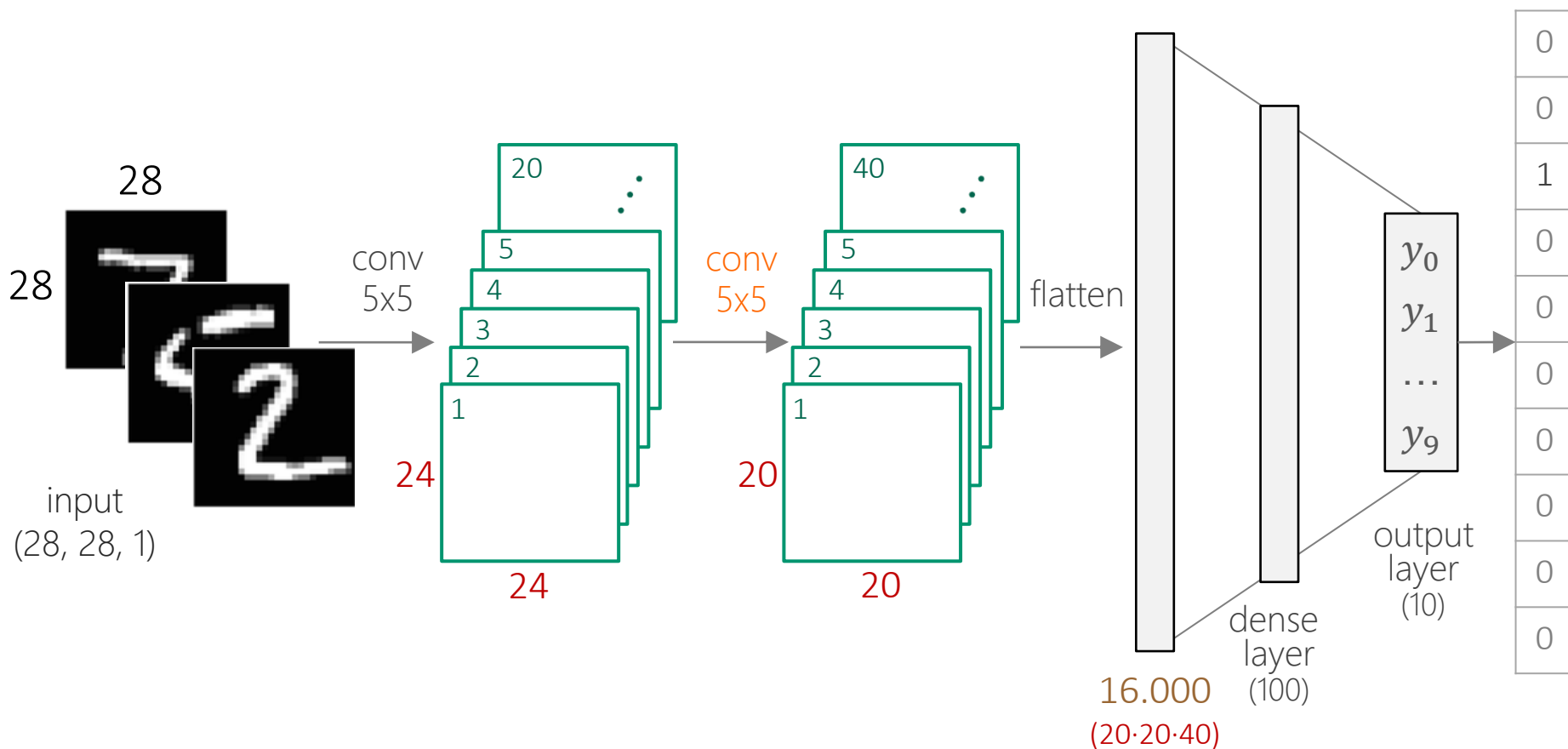
convolutional neural network for the MNIST task



How many **parameters** does the **model** have?

case study

convolutional neural network for the MNIST task



$20 \times 5 \times 5 + 20$	+	$40 \times 5 \times 5 \times 20 + 40$	+	$16.000 \times 100 + 100$	+	$100 \times 10 + 10$	=
520	+	20.040	+	1.600.100	+	1.010	=
1.621.670 trainable parameters							

case study

convolutional neural network for the MNIST task

implementation in Keras...

- **Keras** is a Python library for creating neural networks (deep learning framework)
- **interface to high-level modules** based on Tensorflow, Theano, or the Microsoft Cognitive Toolkit
- **it allows you to combine predefined pieces** common in neural networks such as layers, objective functions, activation functions, optimizers, etc.

case study

convolutional neural network for the MNIST task

simplified script

```
(X_train,y_train), (X_test, y_test) = mnist.load_data()
X_train = X_train.reshape(60000, 28, 28, 1).astype('float32')/255
X_test = X_test.reshape(10000, 28, 28, 1).astype('float32')/255
y_train = to_categorical(y_train)
y_test = to_categorical(y_test)

model = Sequential()
model.add(Conv2D(20, kernel_size=5, activation='relu', input_shape=(28, 28, 1)))
model.add(Conv2D(40, kernel_size=5, activation='relu'))
model.add(Flatten())
model.add(Dense(100, activation='relu'))
model.add(Dense(10, activation='softmax'))
model.compile(optimizer='adam', loss='categorical_crossentropy', metrics=['accuracy'])
model.fit(X_train, y_train, validation_data=(X_test, y_test), epochs=20)
```

case study

convolutional neural network for the MNIST task

hyperparameters

```
(X_train,y_train), (X_test, y_test) = mnist.load_data()
X_train = X_train.reshape(60000, 28, 28, 1).astype('float32')/255
X_test = X_test.reshape(10000, 28, 28, 1).astype('float32')/255
y_train = to_categorical(y_train)
y_test = to_categorical(y_test)

model = Sequential()
model.add(Conv2D(20, kernel_size=5, activation='relu', input_shape=(28, 28, 1)))
model.add(Conv2D(40, kernel_size=5, activation='relu'))
model.add(Flatten())
model.add(Dense(100, activation='relu'))
model.add(Dense(10, activation='softmax'))
model.compile(optimizer='adam', loss='categorical_crossentropy', metrics=['accuracy'])
model.fit(X_train, y_train, validation_data=(X_test, y_test), epochs=20)
```

case study

convolutional neural network for the MNIST task

activation 'relu'

- Rectified Linear Unit (ReLU)
- activation function of hidden layers

activation 'softmax'

- output layer activation function
- probability distribution over K outputs

loss 'categorical_crossentropy'

- one-hot vector + softmax + cross entropy
- loss function
- measures discrepancy between two distributions

optimizer 'adam'

- adam = adaptive moment estimation
- parameter optimization method

optimizer

adam: Adaptive Moment Estimation

gradient descent

$$\theta_{t+1} = \theta_t - \gamma \frac{\partial L}{\partial \theta}(\theta_t)$$

adam:

$$\theta_{t+1} = \theta_t - \frac{\eta}{\sqrt{\hat{v}_t + \epsilon}} \hat{m}_t$$

$$\hat{m}_t = \frac{m_t}{1 - \beta_1^t} \quad \hat{v}_t = \frac{v_t}{1 - \beta_2^t}$$

$$m_t = \beta_1 m_{t-1} + (1 - \beta_1) g_t$$

$$v_t = \beta_2 v_{t-1} + (1 - \beta_2) g_t^2$$

$$g_t = \frac{\partial L}{\partial \theta}(\theta_t)$$

typical values

$$\beta_1 = 0.9$$

$$\beta_2 = 0.999$$

$$\epsilon = 10^{-8}$$

optimizer

adam: considerations

- Adam adjusts learning rates individually for each parameter
- it is an evolution of Stochastic Gradient Descent with momentum
- \mathbf{m}_t is a first-order moment (moving average of the gradient)
- \mathbf{m}_t is a mean estimate of the last $1/(1 - \beta_1)$ gradients
- \mathbf{v}_t is a second-order moment (it characterizes the variance of \mathbf{g})
- \mathbf{v}_t scales the learning rates adaptively for the parameter θ
- the adapted coefficient is inversely proportional to the variance of the gradient
- The epsilon (ϵ) value is a small constant added for numerical stability.
- $\mathbf{m}_0 = \mathbf{0}$ y $\mathbf{v}_0 = \mathbf{0}$; $\hat{\mathbf{m}}_t$ y $\hat{\mathbf{v}}_t$ are unbiased estimators of the mean and variance of \mathbf{g}
- Adam has become a go-to choice for many machine learning practitioners

hyperparameters

convolutional neural network for the MNIST task

CNN hyperparameters (MNIST case study)

- 4 trainable layers: 2 convolutional + 2 dense
- layer 1: 20 x 5x5 filters, ReLU activation
- layer 2: 40 x 5x5 filters, ReLU activation
- layer 3: dense with 100 units, ReLU activation
- layer 4: output layer, softmax activation
- padding: 0
- stride: 1

ConvNetJS CIFAR-10 demo

Training Stats

pause

Forward time per example: 53ms

Backprop time per example: 77ms

Classification loss: 0.63136

L2 Weight decay loss: 0.00187

Training accuracy: 0.82

Validation accuracy: 0.79

Examples seen: 2965

Learning rate: 0.0001

change

Momentum: 0.9

change

Batch size: 2

change

Weight decay: 0.00001

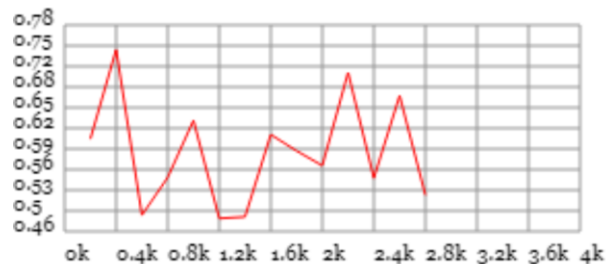
change

save network snapshot as JSON

init network from JSON snapshot

load a pretrained network (achieves ~80% accuracy)

Loss:



clear graph

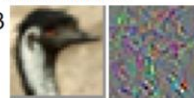
Network Visualization

input (32x32x3)

max activation: 0.33921, min: -0.46863

max gradient: 0.36864, min: -0.31779

Activations:



conv (32x32x16)

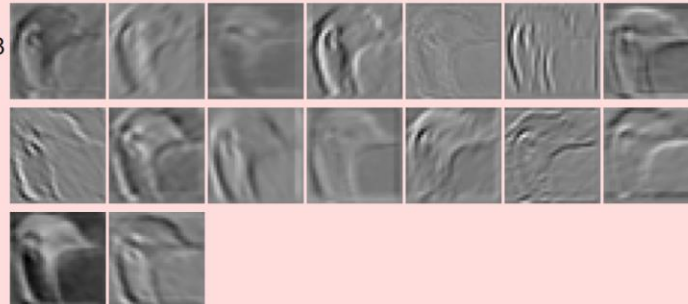
filter size 5x5x3, stride 1

max activation: 2.70593, min: -3.20258

max gradient: 0.11371, min: -0.08803

parameters: $16 \times 5 \times 5 \times 3 + 16 = 1216$

Activations:



CNN

application scope

application scope

- CNN is able to detect and characterize local spatial patterns
- typical applications on CNN-based images:
 - image classification
 - detection (and location) of objects in images
 - image reconstruction, noise removal
 - image fusion
 - super-resolution
 - ...

limitations

- CNN might not be useful if data has no spatial or temporal order