

# Deep Learning Classification & Convolutional Neural Networks

Dr. Mohammed Salah Al-Radhi  
(slides by: Dr. Bálint Gyires-Tóth )



# Copyright

Copyright © **Mohammed Salah Al-Radhi**, All Rights Reserved.

This presentation and its contents are protected by copyright law. The intellectual property contained herein, including but not limited to text, images, graphics, and design elements, are the exclusive property of the copyright holder identified above. Any unauthorized use, reproduction, distribution, or modification of this presentation or its contents is strictly prohibited without prior written consent from the copyright holder.

**No Recordings or Reproductions:** Attendees, viewers, and recipients of this presentation are expressly prohibited from making any audio, video, or photographic recordings, as well as screen captures, screenshots, or any form of reproduction, of this presentation, its content, or any related materials, whether during its live presentation or subsequent access. Violation of this prohibition may result in legal action.

For permissions, inquiries, or licensing requests, please contact: **malradhi@tmit.bme.hu**

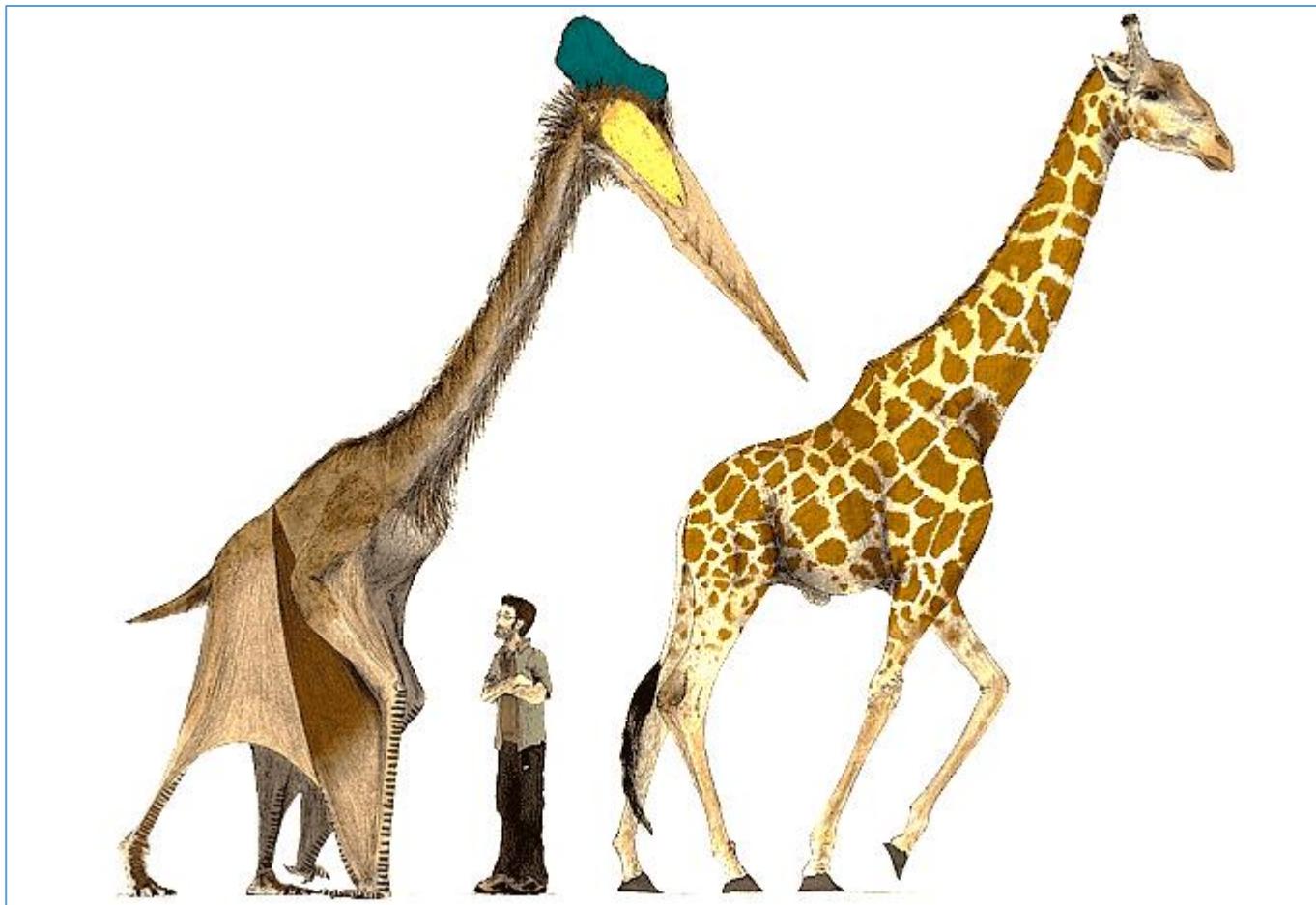
Unauthorized use, distribution, or reproduction of this presentation may result in civil and criminal penalties. Thank you for respecting the intellectual property rights of the copyright holder.

# Lecture Overview

1. Data Visualization
2. Classification
  - a. Cost function
  - b. Confusion matrix
  - c. ROC & AUC curves
3. Convolution Neural Network
  - a. Stride
  - b. Zero padding
  - c. Pooling
  - d. Dilation
  - e. Layers
  - f. General architecture

# Data Visualization

# Data Visualization



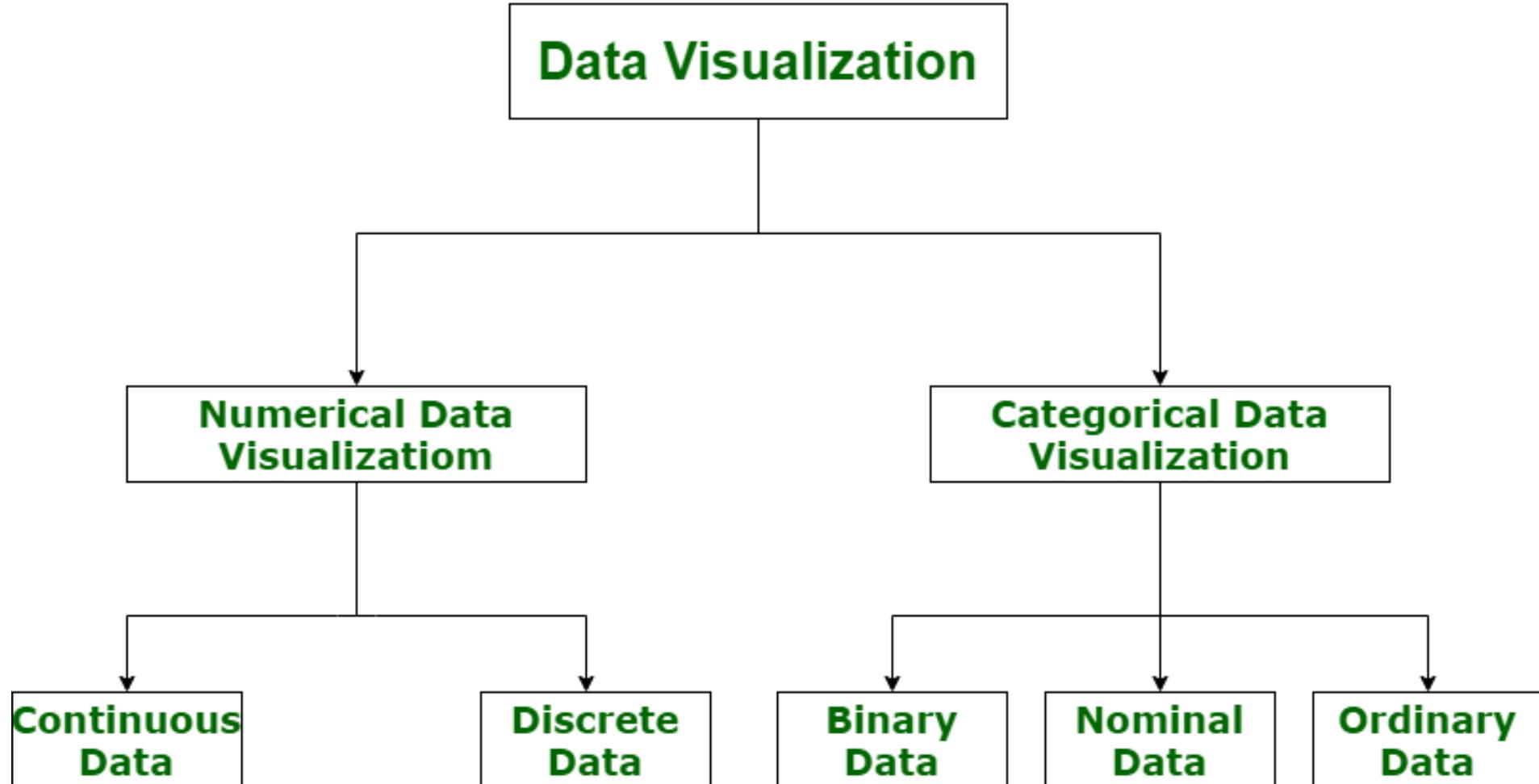
images by Mark P. Witton

# What is Data Visualization?

- presentation of data in a graphic format.
- helps people see patterns and trends.
- used from simple graphs to complex network diagrams.
- puts the Data into the correct context
- Saves Time
- Tells a Data Story and making it more accessible, understandable, and usable.



# Data Visualization



# Data Visualization

## Top Data Visualization Tools:

- Tableau
- Looker
- Zoho Analytics
- Sisense
- IBM Cognos Analytics
- Qlik Sense
- Microsoft Power BI
- SAP Analytics Cloud

## Top Data Visualization Libraries Available in Python Libraries:

- Matplotlib
- Plotly
- Seaborn
- Altair
- Bokeh
- Scikit Learn
- Pandas
- YellowBrick

# Classification

# Is this Classification?

- Is this bank transfer fraudulent?



- Is this patient healthy?



- Will you vote for me, for X, or for Y?



- Will these two people fit to each other?

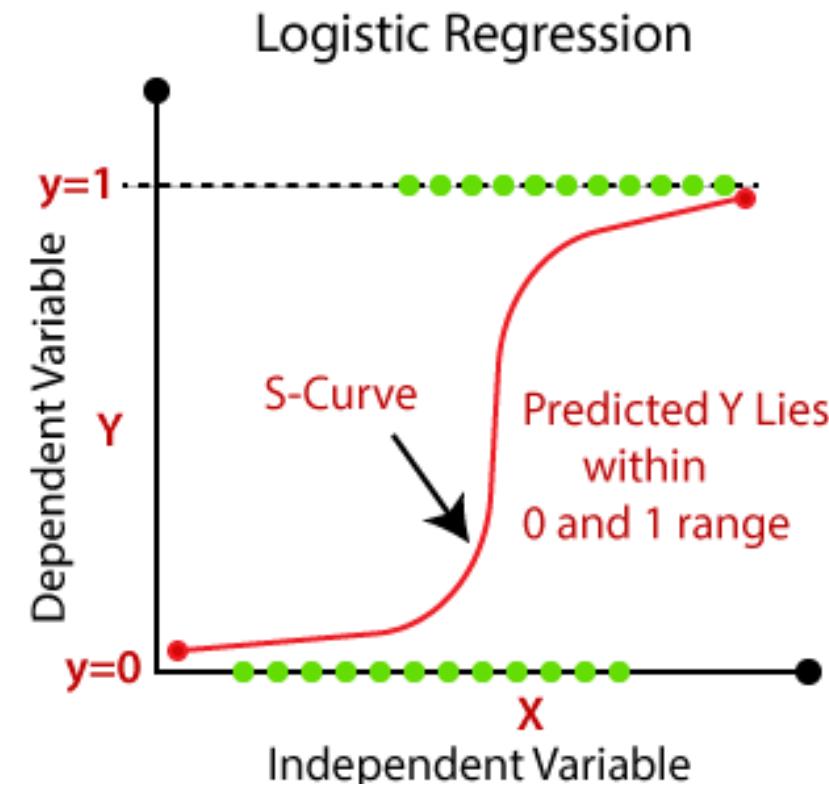
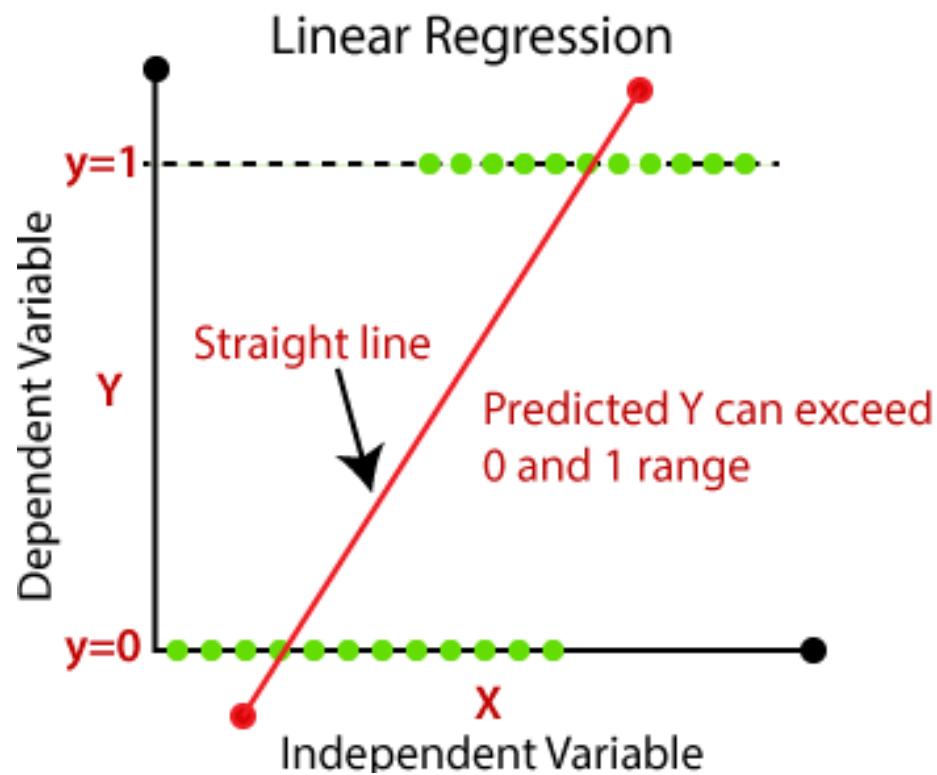


- Is this an apple?



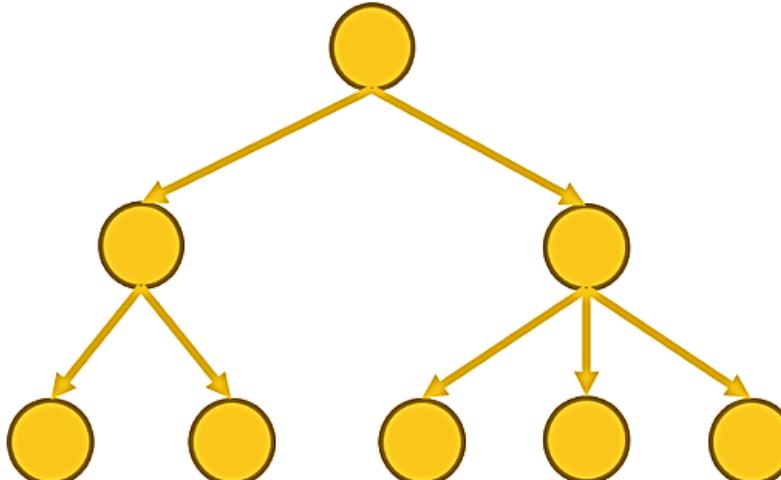
Given a number of examples, identify to which class a given observation belongs to.

# Regression

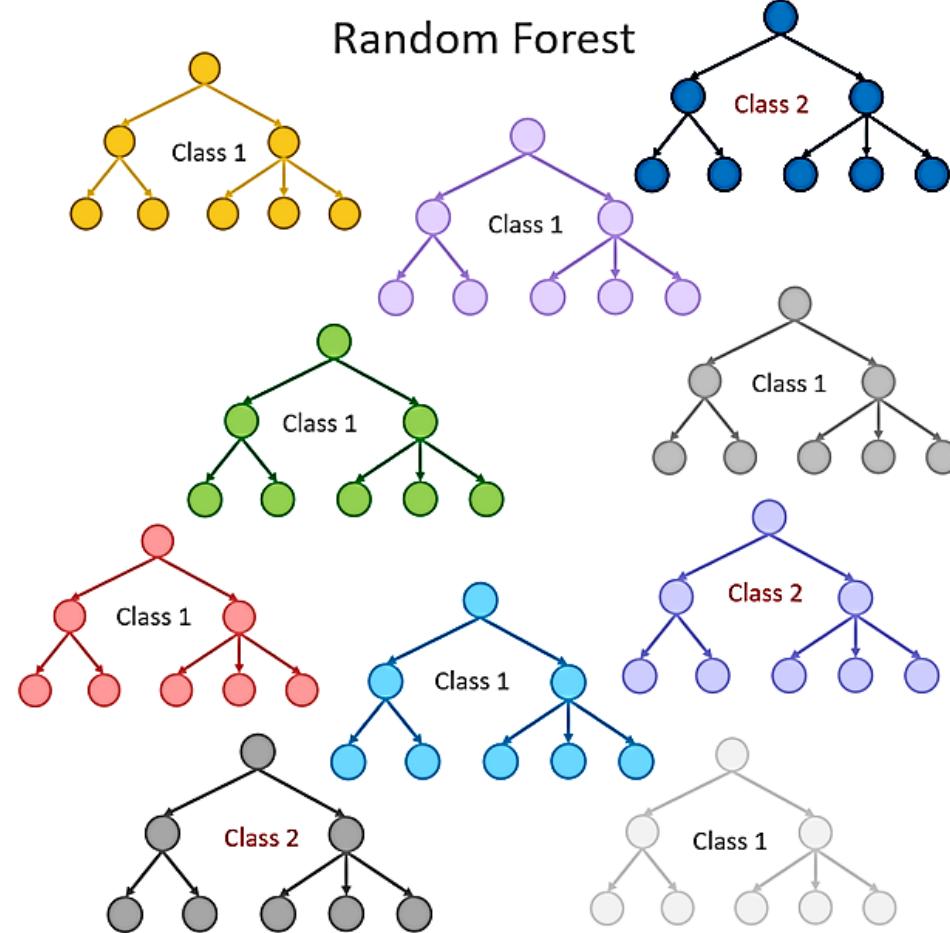


# Tree

Single Decision Tree



Random Forest



# Regression vs. Classification

- **Regression**: predict a continuous output value
  - e.g., temperature of a room
- **Classification**: predicts a discrete value
  - Binary classification: output is either 0 or 1
  - Multi-class classification: set of N classes

# Regression vs. Classification



What is the temperature going to be tomorrow?

PREDICTION

84°



Will it be Cold or Hot tomorrow?

PREDICTION

COLD

HOT



# Regression vs. Classification



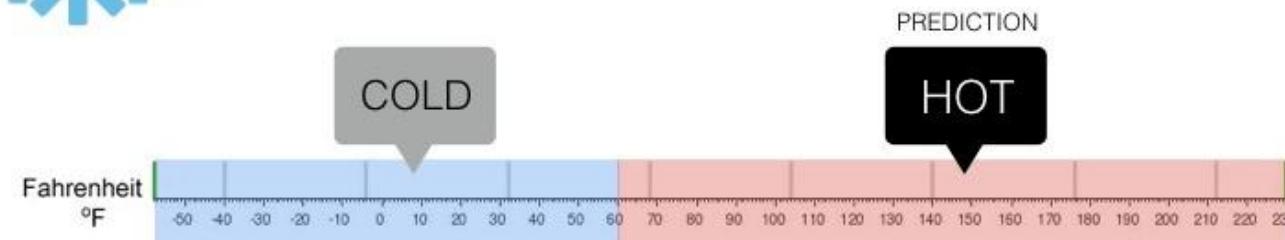
## Regression

What is the temperature going to be tomorrow?

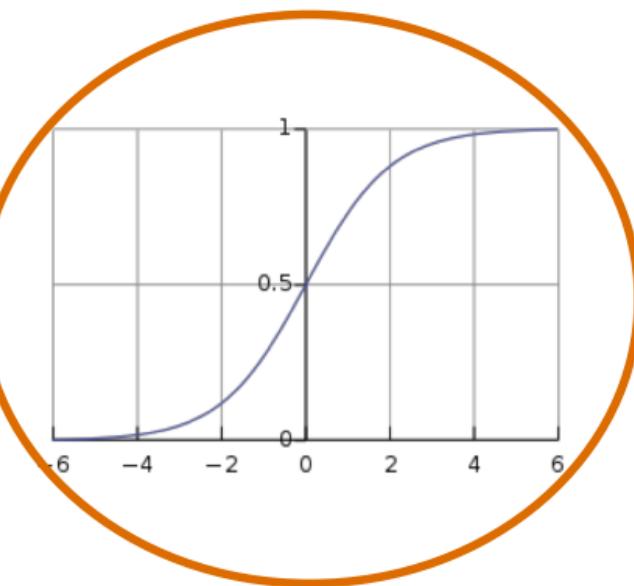
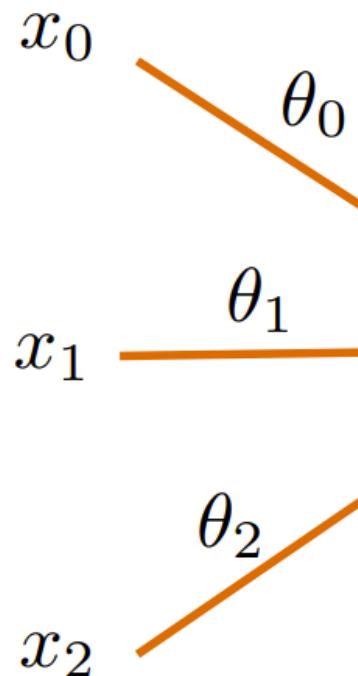


## Classification

Will it be Cold or Hot tomorrow?



# Sigmoid for binary predictions



$$\sigma(x) = \frac{1}{1 + e^{-x}}$$

1  
0

Can be interpreted as a probability

$$p(y_i = 1 | \mathbf{x}_i, \boldsymbol{\theta})$$

# Cost functions, output layer types

- Regression
  - Sigmoid / tanh
  - Data scaled to: between 0.01..0.99 or -0.99 .. 0.99
  - Loss function: mean\_squared\_error
- Classification
  - Two classes: sigmoid (one output) + binary\_crossentropy
  - Multiple (k) labels: sigmoid (k outputs, multihot encoding) + binary\_crossentropy
  - Multiple (n) classes: softmax (n outputs, onehot encoding) + categorical\_crossentropy

# Classification vs regression

- Regression: sigmoid / tanh / linear + MSE
- Classification: softmax + cross-entropy
- Cost function: cross-entropy

$$C = - \sum_{n=1}^M \sum_{i=1}^K y_i^{(n)} \cdot \ln (\hat{y}_i^{(n)})$$

M: number of samples  
K: number of classes

- SoftMax: converts numbers to „probabilities” (with sum of 1)

$$y_i(x) = \frac{e^{x_i}}{\sum_{k=1}^K e^{x_k}} \quad i=1 \dots K$$

# Cross-entropy for classification 1

Y_target (OneHot)			Y_predicted (probability)			Class	Correct?
1	0	0	0.1	0.1	0.8	car	no
0	1	0	0.3	0.4	0.3	rocket	yes
0	0	1	0.3	0.3	0.4	airplane	yes

Classification accuracy:  $\frac{2}{3} = 0.6$

Cross-entropy:  $-(\ln(0.1) \cdot 1 + \ln(0.1) \cdot 0 + \ln(0.8) \cdot 0) = -\ln(0.1)$   
 $-(\ln(0.3) \cdot 0 + \ln(0.4) \cdot 1 + \ln(0.3) \cdot 0) = -\ln(0.4)$   
 $-(\ln(0.3) \cdot 0 + \ln(0.3) \cdot 0 + \ln(0.4) \cdot 1) = -\ln(0.4)$

Average cross-entropy:  $\frac{-(\ln(0.1) + \ln(0.4) + \ln(0.4))}{3} \approx 1.38$

MSE:  $\left\{ \begin{array}{l} (0.1-1)^2 + (0.1-0)^2 + (0.8-0)^2 = 1.46 \\ (0.3-0)^2 + (0.4-1)^2 + (0.3-0)^2 = 0.54 \\ (0.3-0)^2 + (0.3-0)^2 + (0.4-1)^2 = 0.54 \end{array} \right\} \frac{1.46 + 0.54 + 0.54}{3} \approx 0.847$

# Cross-entropy for classification 2

Y_target (OneHot)			Y_predicted (probability)			Class	Correct?
1	0	0	0.3	0.4	0.3	Car	no
0	1	0	0.1	0.8	0.1	Rocket	yes
0	0	1	0.1	0.1	0.8	Airplane	yes

Classification accuracy:

$$\cup \quad \cup \quad \frac{2}{3} = 0.6$$

Cross-entropy:

$$- \{ (\ln(0.3) \cdot 1) + (\ln(0.4) \cdot 0) + (\ln(0.3) \cdot 0) \} = -\ln(0.3)$$

$$- \{ (\ln(0.1) \cdot 0) + (\ln(0.8) \cdot 1) + (\ln(0.1) \cdot 0) \} = -\ln(0.8)$$

$$- \{ (\ln(0.1) \cdot 0) + (\ln(0.1) \cdot 0) + (\ln(0.8) \cdot 1) \} = -\ln(0.8)$$

Average cross-entropy:

$$\frac{-(\ln(0.1) + \ln(0.8) + \ln(0.8))}{3} \approx 0.91$$

$$\left. \begin{array}{l} (0.3-1)^2 + (0.4-0)^2 + (0.3-0)^2 = 0.74 \\ (0.1-0)^2 + (0.8-1)^2 + (0.1-0)^2 = 0.06 \\ (0.1-0)^2 + (0.1-0)^2 + (0.8-1)^2 = 0.06 \end{array} \right\}$$

$$\frac{0.74 + 0.06 + 0.06}{3} = 0.286$$

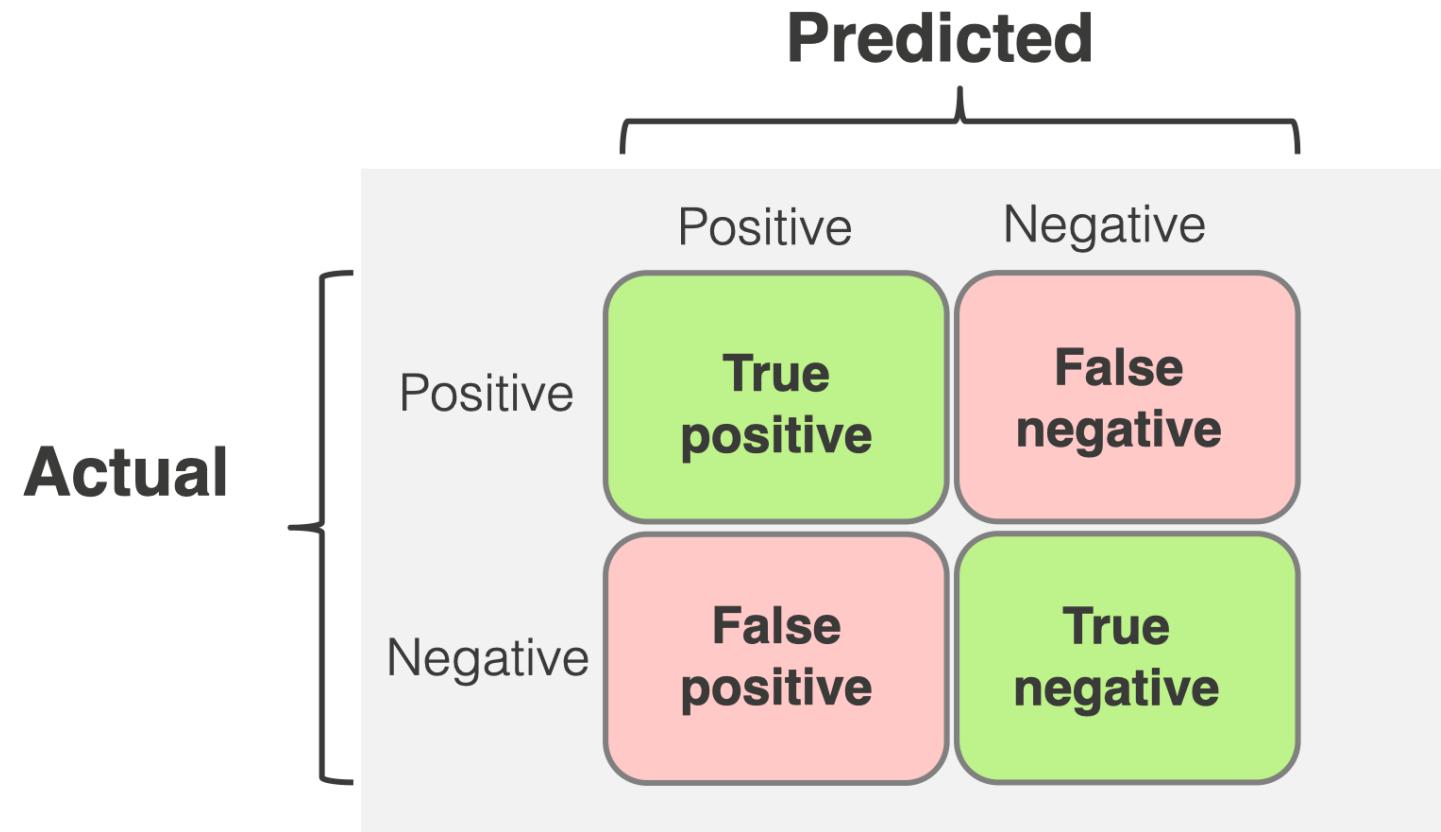
# Evaluation methods

# Classification: evaluation methods

- Confusion matrix (TP, TN, FP, FN)
- Classification metrics
  - Recall
  - Precision
  - Accuracy
  - F1
- ROC
- AUC

# Confusion matrix

- Easy to read, includes all results.



# Confusion matrix

- E.g. 3 classes:

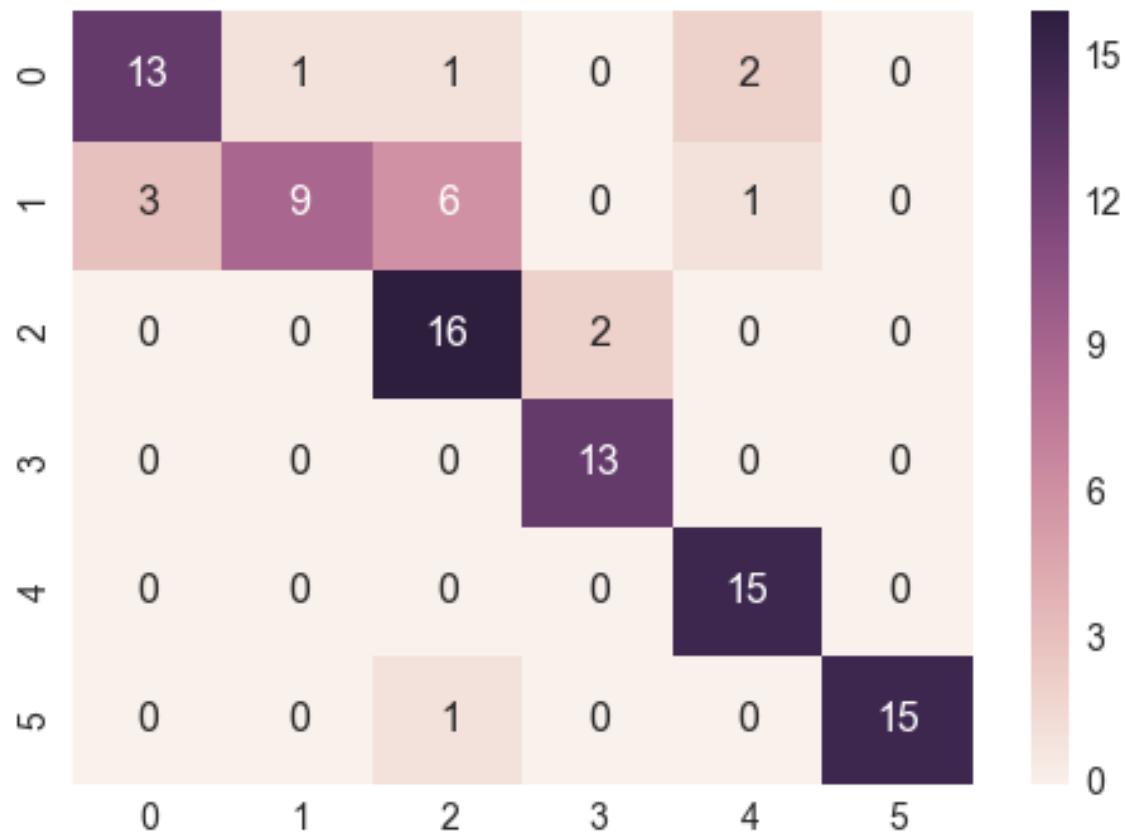
		Y_predicted		
		Class	Rocket	Car
Y_target	Rocket	15	3	8
	Car	4	23	1
	Airplane	9	2	19

# Confusion matrix: binary classification

Rocket		<i>Condition</i>	Y_predicted	
			Positive	Negative
Y_target	Positive	15 True Positive (TP)		3+8=11 False Negative (FN)
	Negative	4+9=13 False Positive (FP)		23+19+1+2=42 True Negative (TN)

Pl. shoot, if rocket! (to protect the city):

- **True Positive (TP):** We correctly identified it as a rocket (✓).
- **True Negative (TN):** We correctly identified it as not a rocket (✓).
- **False Positive (FP):** We incorrectly identified it as a rocket (X).
- **False Negative (FN):** We incorrectly identified it as not a rocket (X).



```
sns.heatmap(confusion_matrix, annot=True)
```

# Classification metrics

$$\text{Accuracy} = \frac{\text{TP} + \text{TN}}{\text{TP} + \text{TN} + \text{FP} + \text{FN}}$$

$$\text{Precision or } = \frac{\text{TP}}{\text{TP} + \text{FP}}$$

$$\text{Recall or TPR} = \frac{\text{TP}}{\text{TP} + \text{FN}}$$

$$\text{Specificity or TNR} = \frac{\text{TN}}{\text{TN} + \text{FP}}$$

$$\text{F1 Score} = 2 * \frac{\text{Recall} \times \text{Precision}}{\text{Recall} + \text{Precision}}$$

$$\text{False Positive Ratio(FPR)} = \frac{\text{FP}}{\text{FP} + \text{FN}}$$

More: MAP, ROC, AUROC

# SkLearn Precision/Recall/F1 means

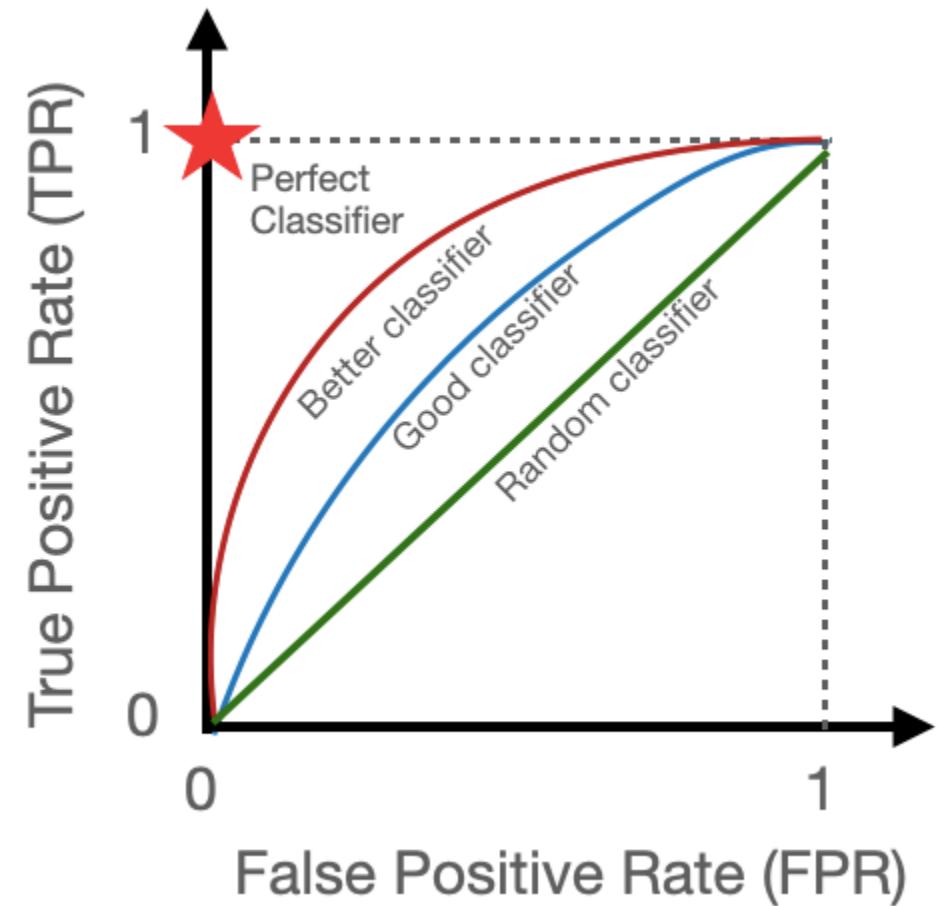
- **Macro:** calculate the value for all classes separately, and take the average.
- **Micro:** calculate the values for all classes, and calculate precision/recall/F1 from that.
- **Weighted:** weights the classes based as a function of the imbalances. Those classes which have more values have higher weight.

# ROC Curve, AUC

Receiver Operating Characteristic

TP and FP rates under various  
Decision thresholds

AUC: Area Under Curve



# Convolution Neural Network

# “CNN” vs “ConvNet”

- There are many papers that use either phrases, but:
- “ConvNet” is the preferred term, since “CNN” clashes with other things called CNN ☺

# Deep learning news station

- Video generation based on text: <https://makeavideo.studio/>

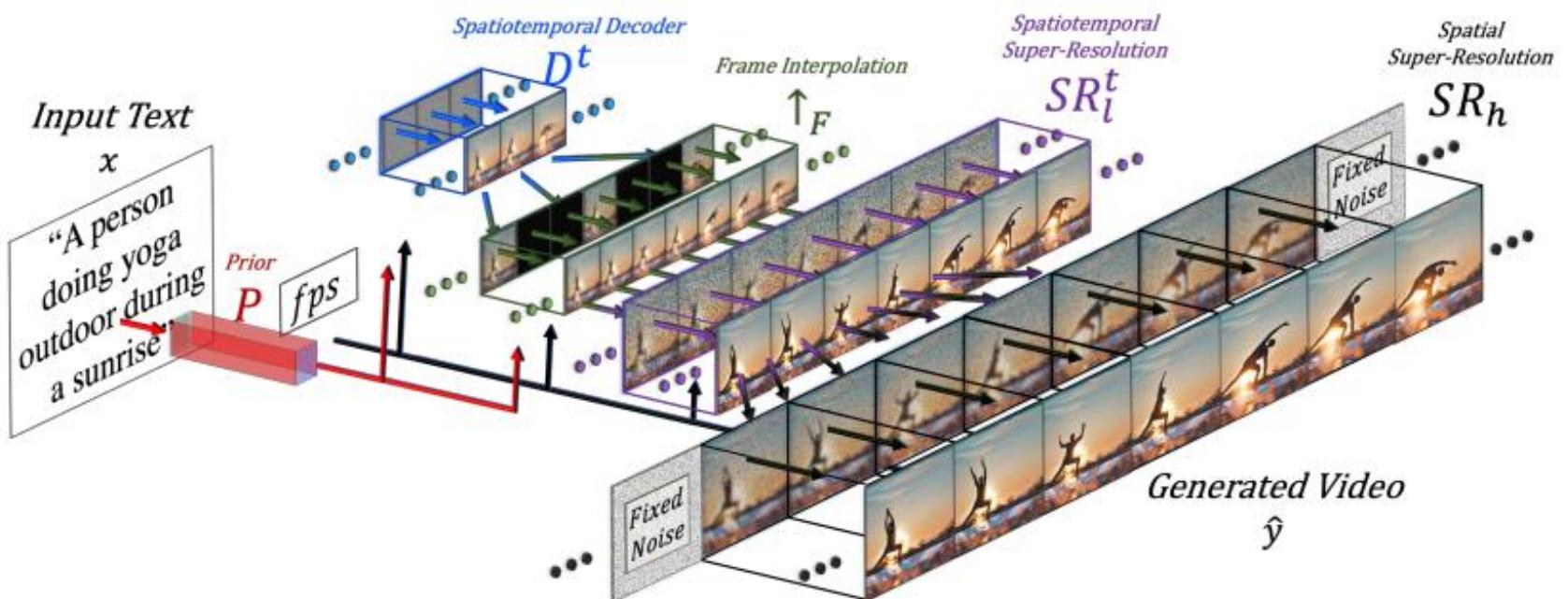
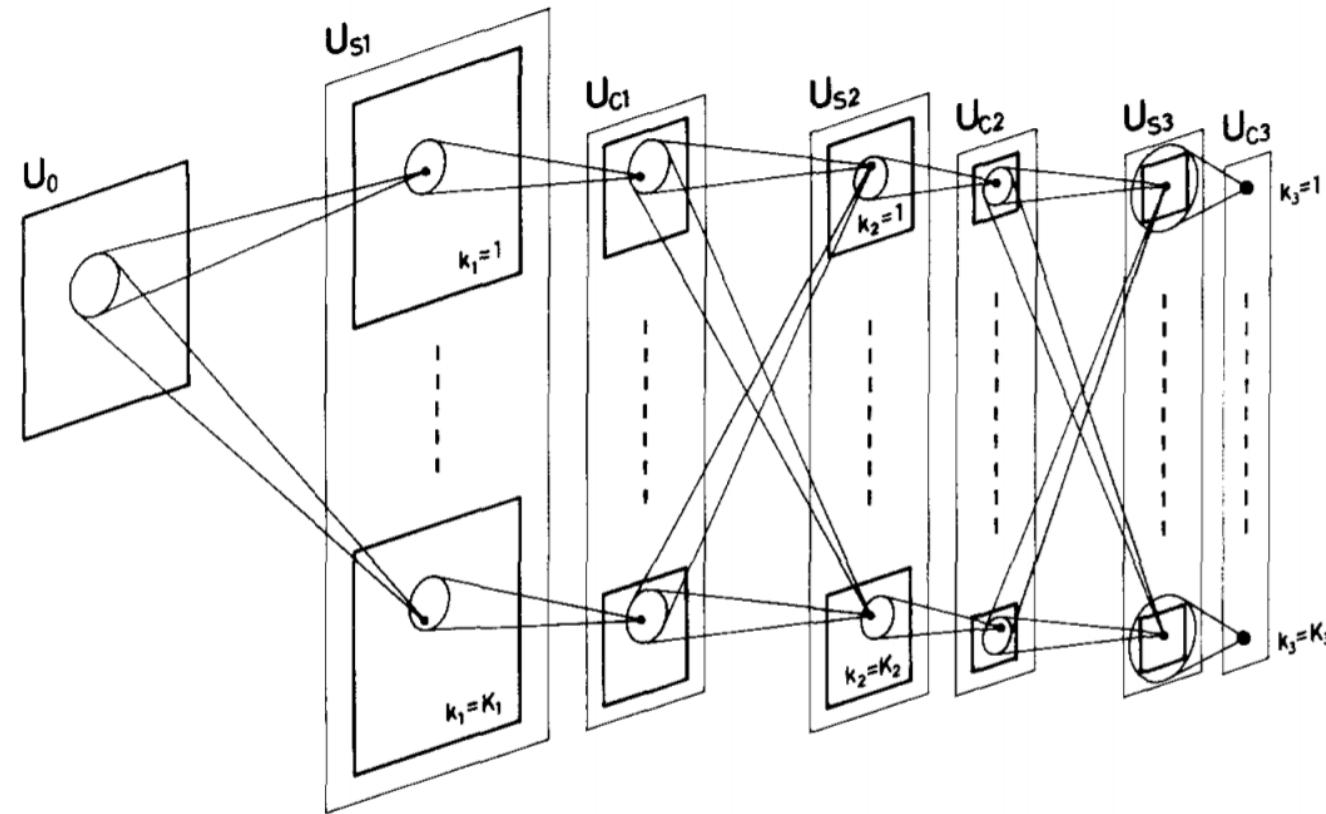


Figure 2: **Make-A-Video high-level architecture.** Given input text  $x$  translated by the prior  $P$  into an image embedding, and a desired frame rate  $fps$ , the decoder  $D^t$  generates 16  $64 \times 64$  frames, which are then interpolated to a higher frame rate by  $\uparrow_F$ , and increased in resolution to  $256 \times 256$  by  $SR_l^t$  and  $768 \times 768$  by  $SR_h$ , resulting in a high-spatiotemporal-resolution generated video  $\hat{y}$ .

# Major milestones

- Neocognitron: Fukushima, Kunihiko, and Sei Miyake. "Neocognitron: A self-organizing neural network model for a mechanism of visual pattern recognition." In Competition and cooperation in neural nets, pp. 267-285. Springer, Berlin, Heidelberg, 1982.



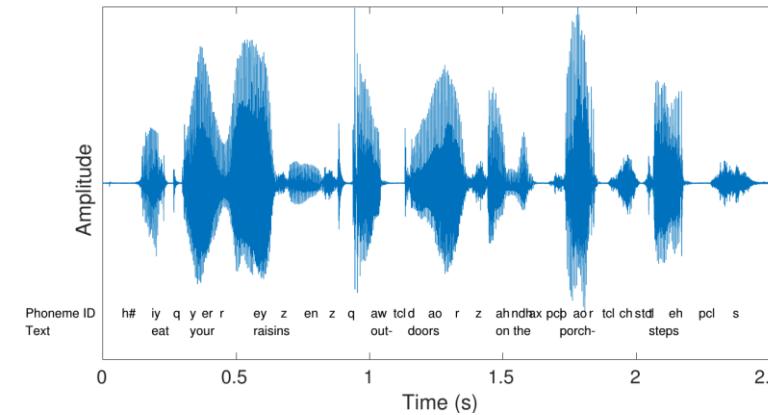
# Major milestones

- **LeNet**: LeCun, Y., & Bengio, Y. (1995). Convolutional networks for images, speech, and time series. *The handbook of brain theory and neural networks*, 3361(10), 1995.
- **LeNet5**: LeCun, Y., Bottou, L., Bengio, Y., & Haffner, P. (1998). Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11), 2278-2324.
- **AlexNet**: Krizhevsky, A., Sutskever, I., & Hinton, G. E. (2012). Imagenet classification with deep convolutional neural networks. In *Advances in neural information processing systems* (pp. 1097-1105).

# But what is CNN?

Version of deep neural networks designed for signals:

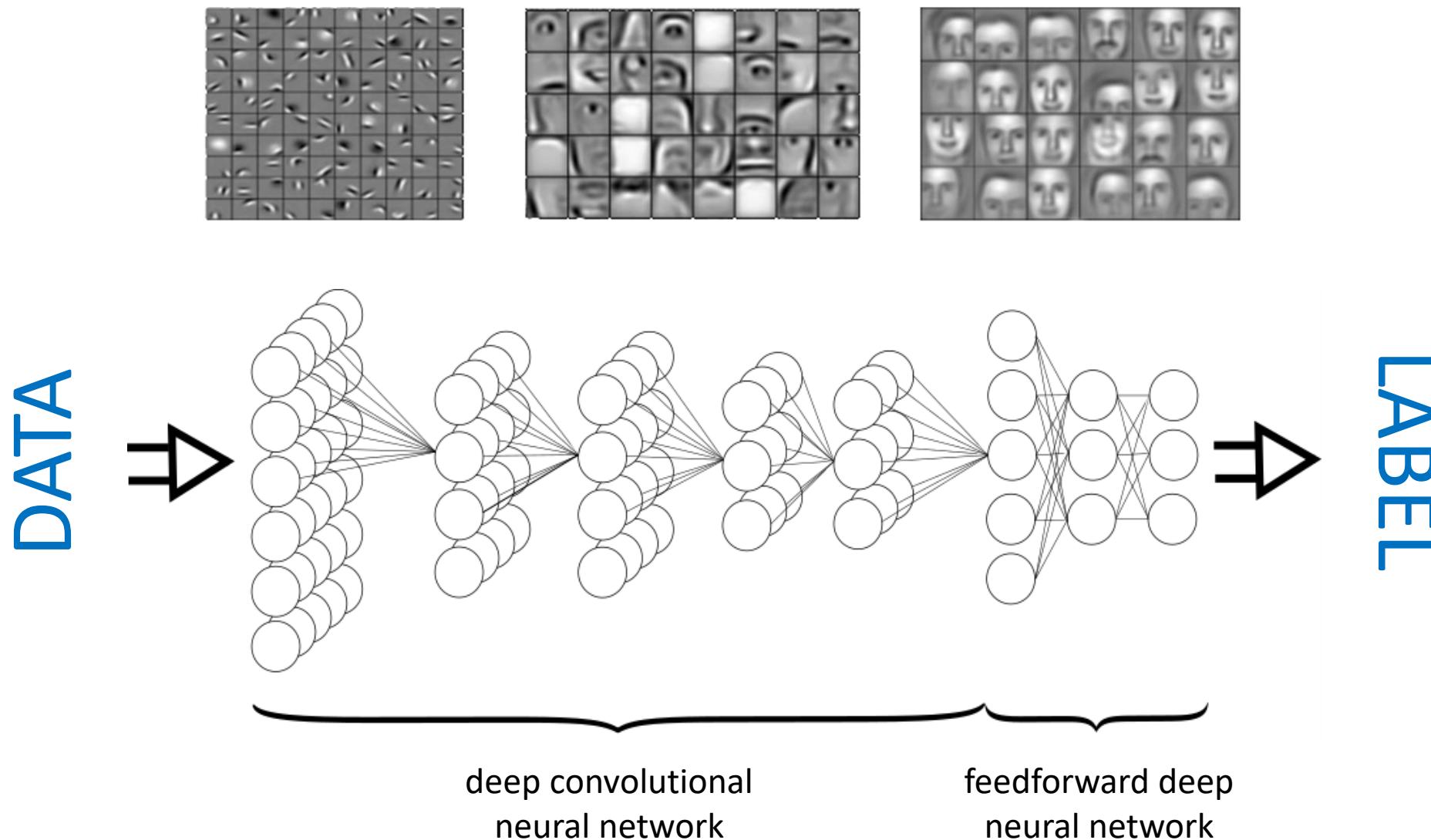
- 1D signals (e.g., speech waveforms)



- 2D signals (e.g., images)



# General structure of a CNN network

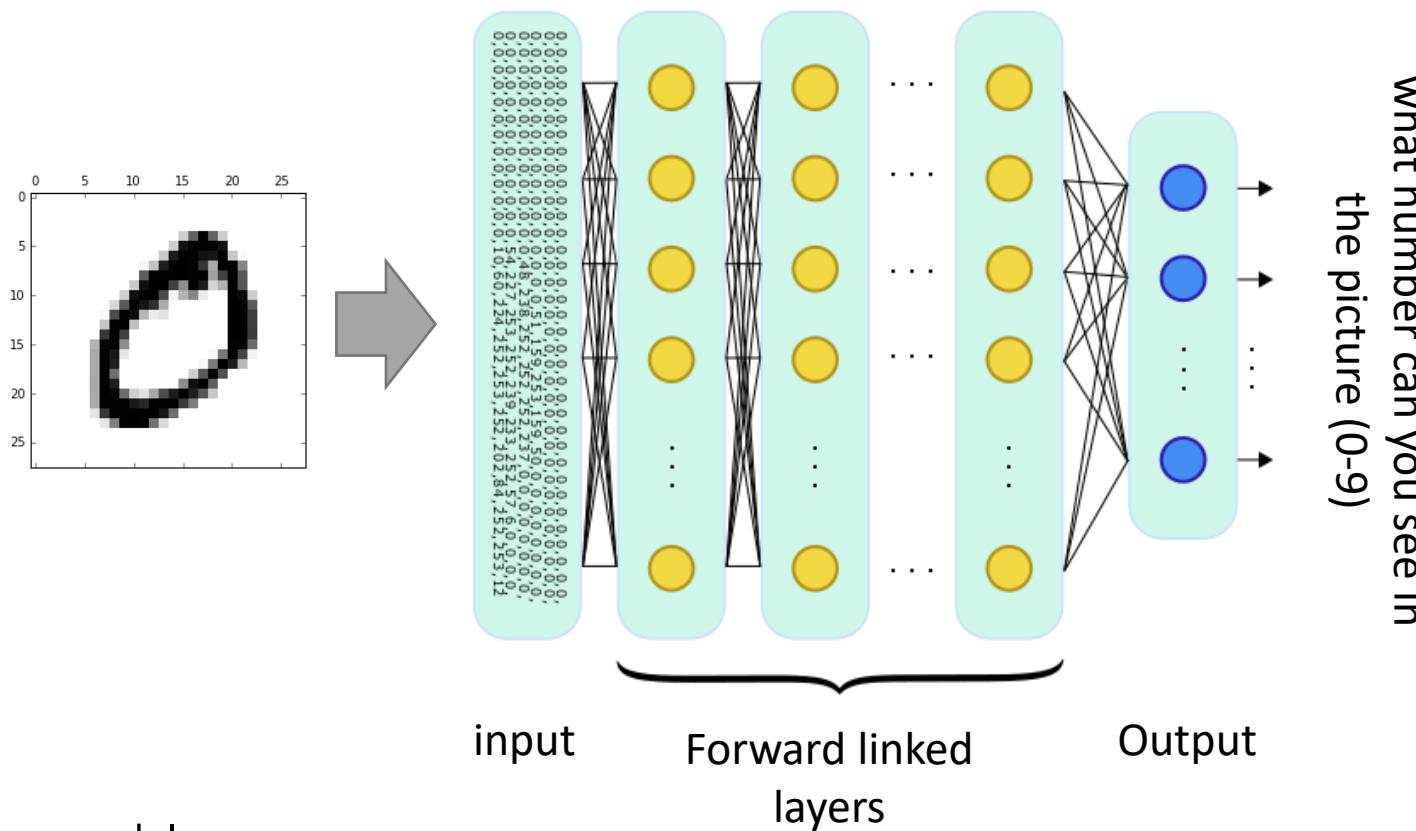


# MNIST: Yann LeCun handwriting database

# MNIST



# Feedforward network approach



what number can you see in  
the picture (0-9)

Main problem

- Parameter space too large
  - Pl.  $200 \times 200 \times 3 = 120000$  input, additional connections ...)
- It recognizes data, not patterns

# Convolutional neural network approach

- The convolution operation involves element-wise **multiplication** of a filter (also known as a kernel) with a corresponding region of the input, followed by the **summation** of these products to produce a single value in the output feature map.

1	0	1
0	1	0
1	0	1

Filter / Kernel

1 <sub>x1</sub>	1 <sub>x0</sub>	1 <sub>x1</sub>	0	0
0 <sub>x0</sub>	1 <sub>x1</sub>	1 <sub>x0</sub>	1	0
0 <sub>x1</sub>	0 <sub>x0</sub>	1 <sub>x1</sub>	1	1
0	0	1	1	0
0	1	1	0	0

Image

4		

Convolved Feature

# CNN – Producing Feature Maps



Original



Sharpen



Edge Detect



“Strong” Edge  
Detect

## CNN main ideas

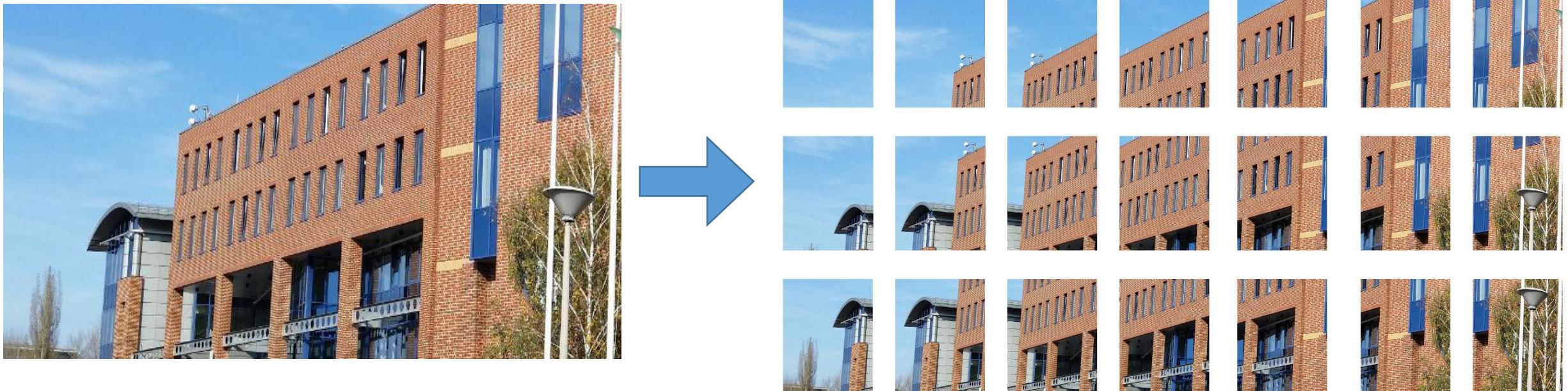
Dependence on space  
and/or time

Shared weights

Multiple (shared) weight  
matrix per layer

# Spatial dependence

- By breaking down the data into smaller pieces, we train the network at different levels.

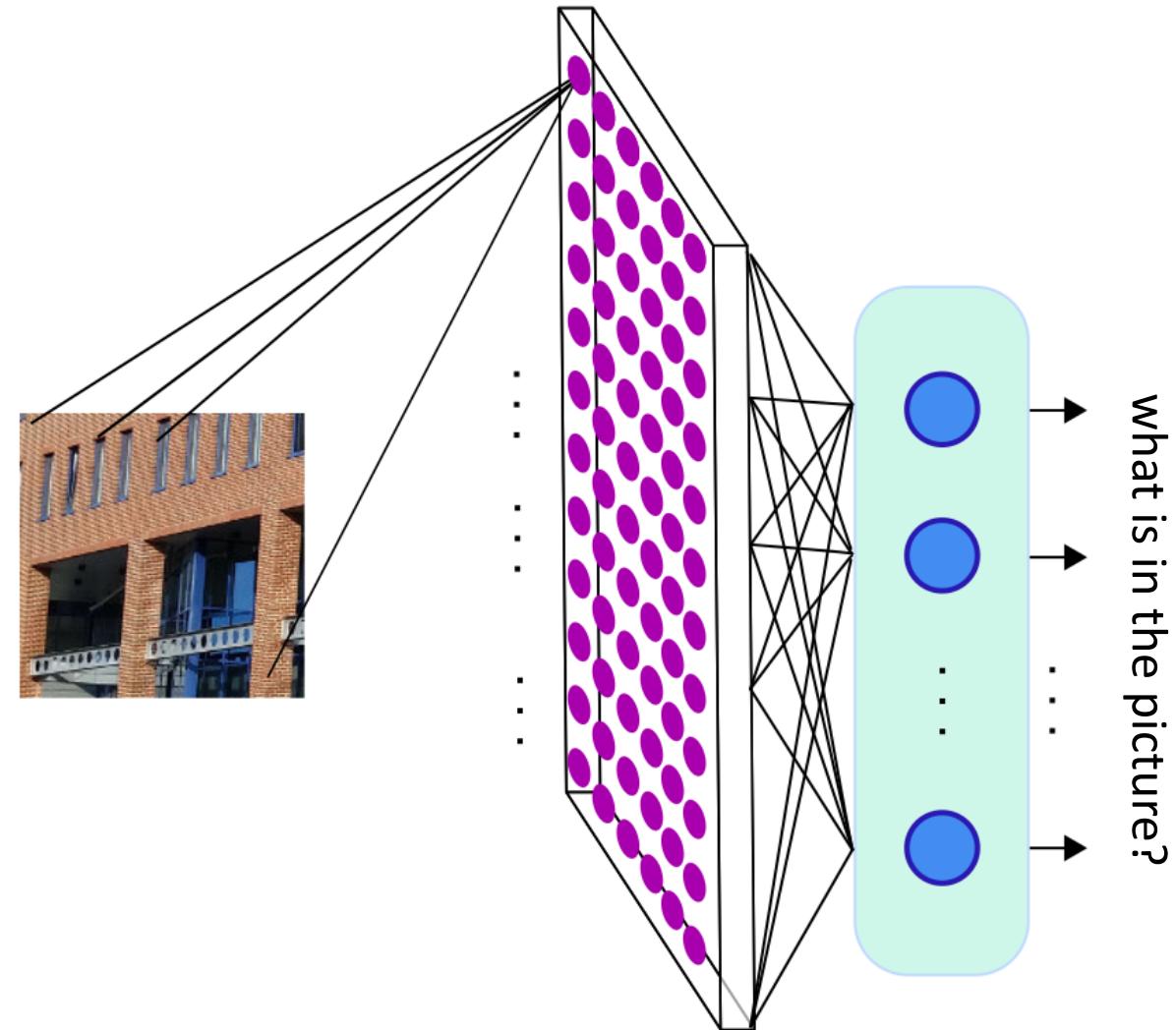


# Building blocks: convolutional layer 1

Training the net with pieces of the original image.

Pl. 20x20 px piece, then 400 weight + 1 bias.

Now the neurons of the feed-forward network are arranged in 2D.



# Building blocks: convolutional layer 2

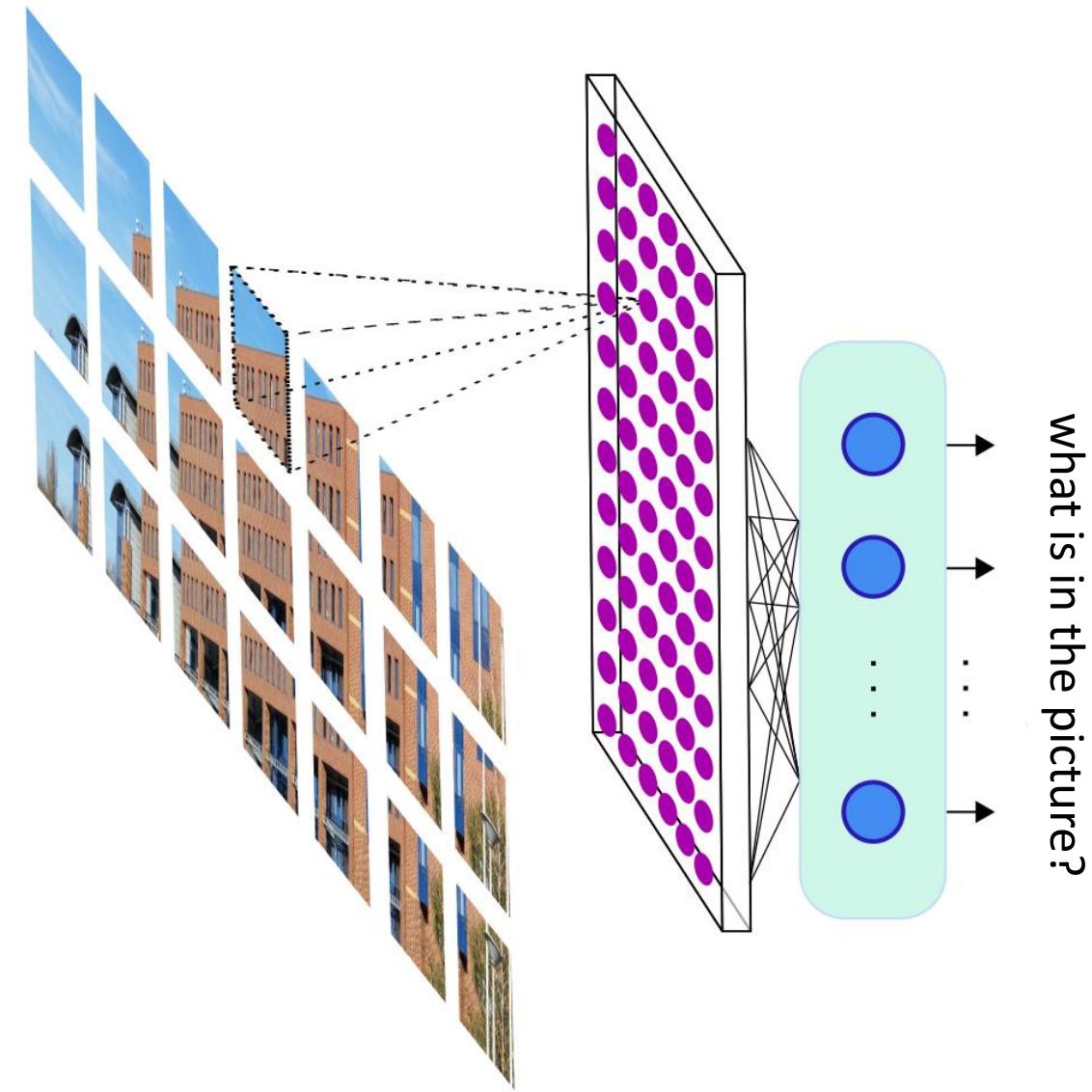
Shared weights:

We use the same weights for each piece of the original image.

If a feature is important at one point in the picture, it will be important elsewhere!

Why we call it CNN?

Convolution of input and shared weights (filters).



# Hyperparameters (conv2d)

INPUT:

{depth}@{input\_x}x{input\_y}

FILTER / KERNEL:

{filter\_x}x{filter\_y}@{filter\_depth}

STRIDE:

{stride\_x}x{stride\_y}

ZERO PADDING:

{zeropadding\_x}x{zeropadding\_y}

DILATATION:

{dilatation\_x}x{dilatation\_y}

Output:

{filter\_depth}@  
{((input\_x-filter\_x+2\*zeropadding\_x)/stride\_x)+1}x  
{((input\_y-filter\_y+2\*zeropadding\_y)/stride\_y)+1}

# Convolution operation

## Input

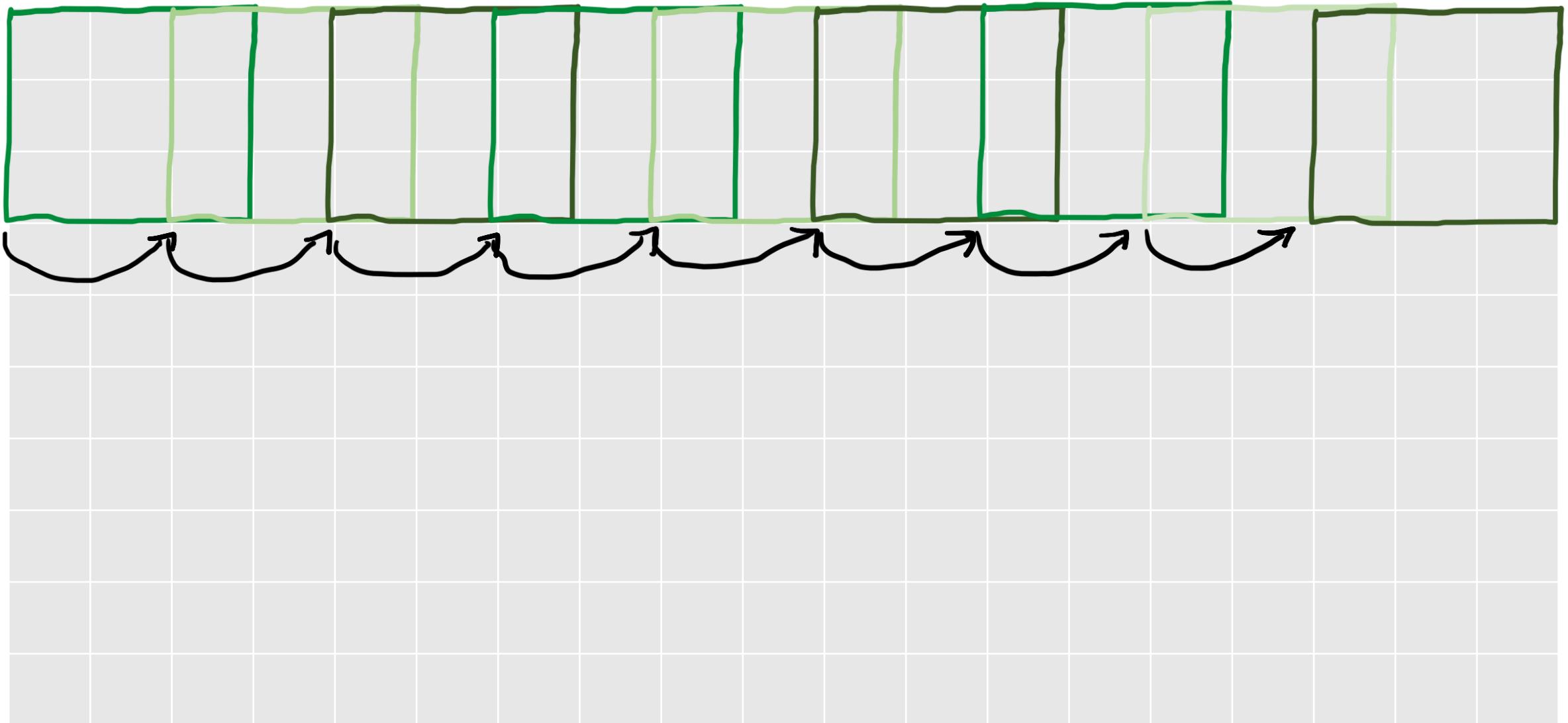
\*

## Kernel

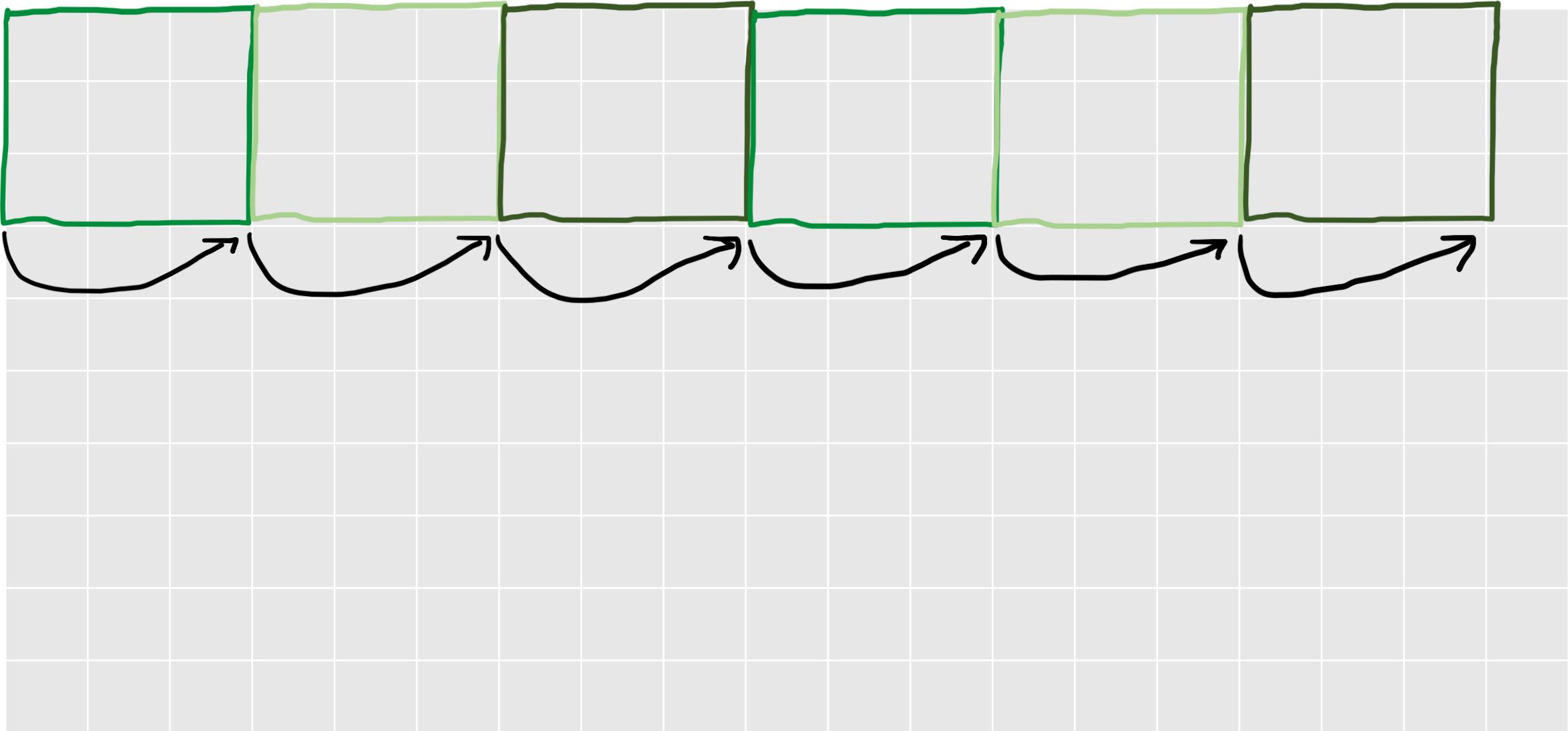
2

## Output

Stride = 2



Stride = 3

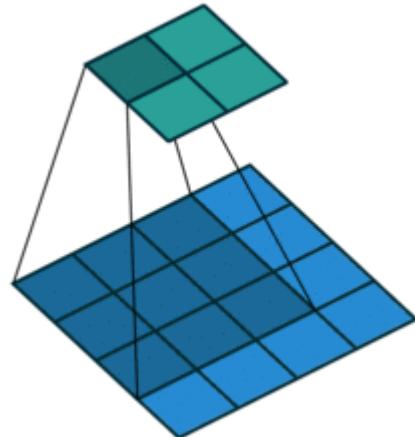


# Zeropadding = 1x1

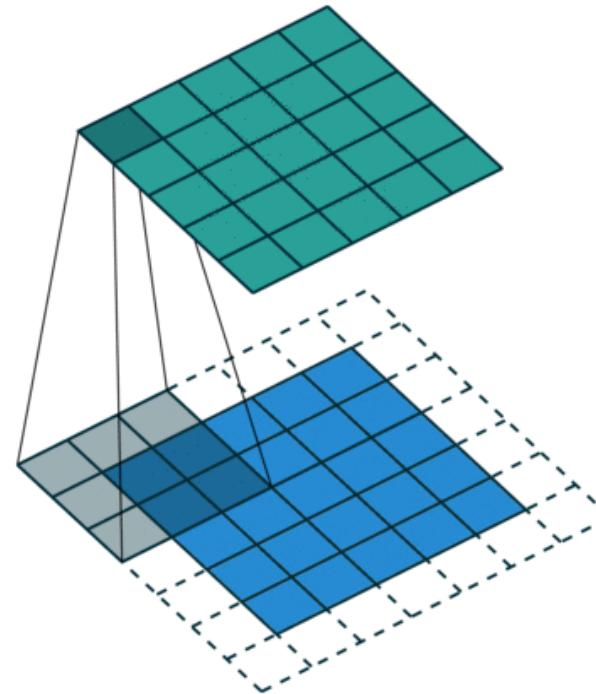
# Zeropadding = 2x2

# Zeropadding

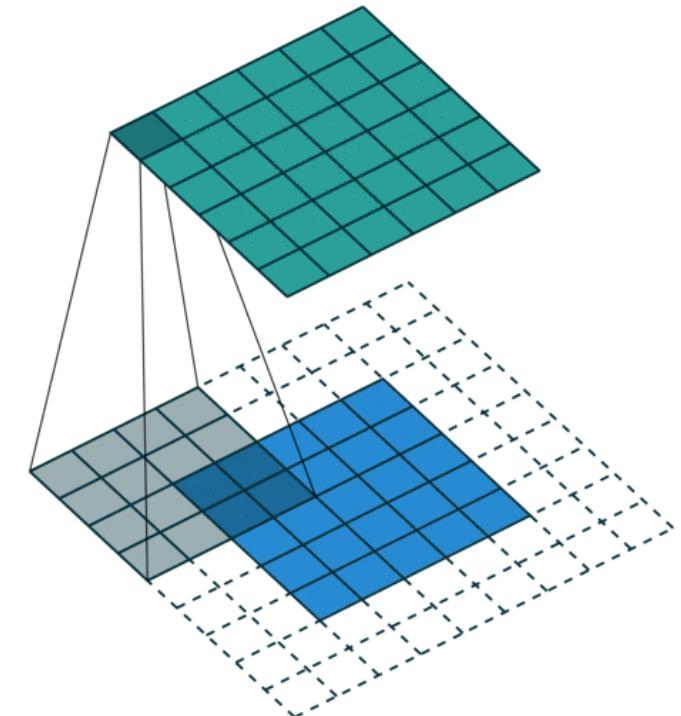
No zero padding, unit strides



Zero padding, unit strides



**1 x 1**



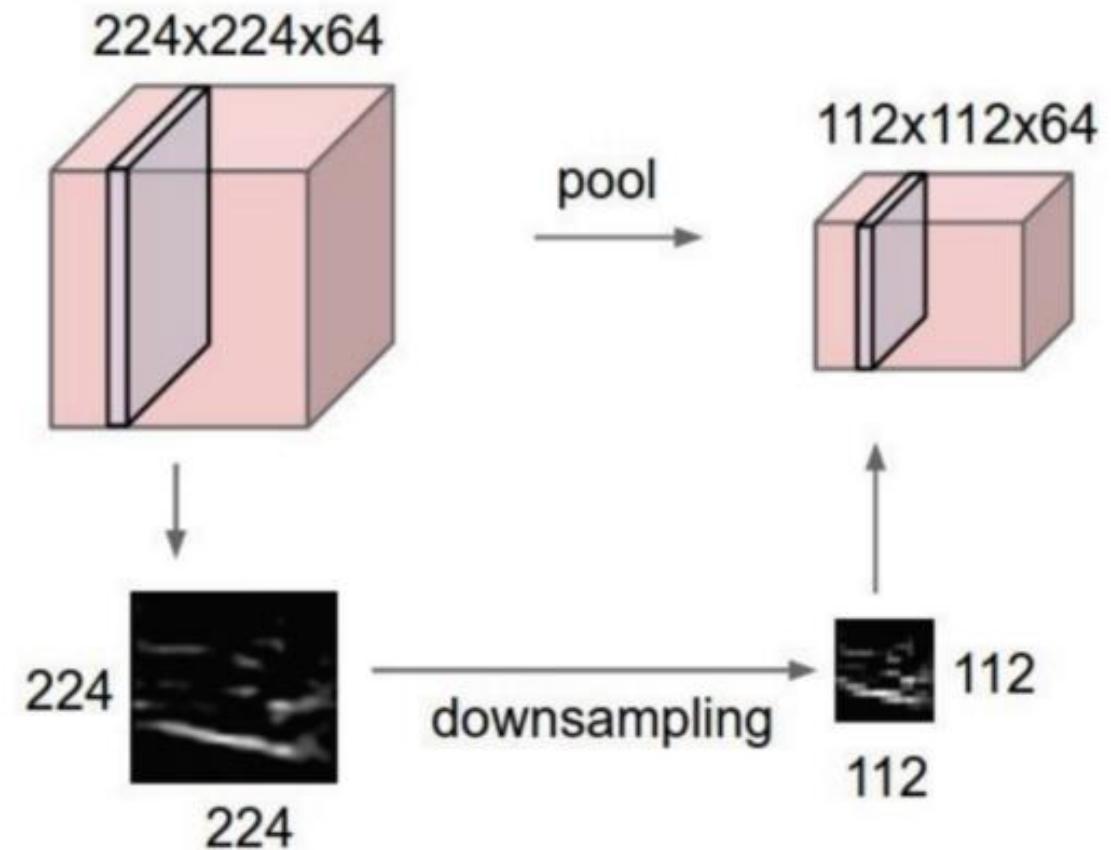
**2 x 2**

# Pooling

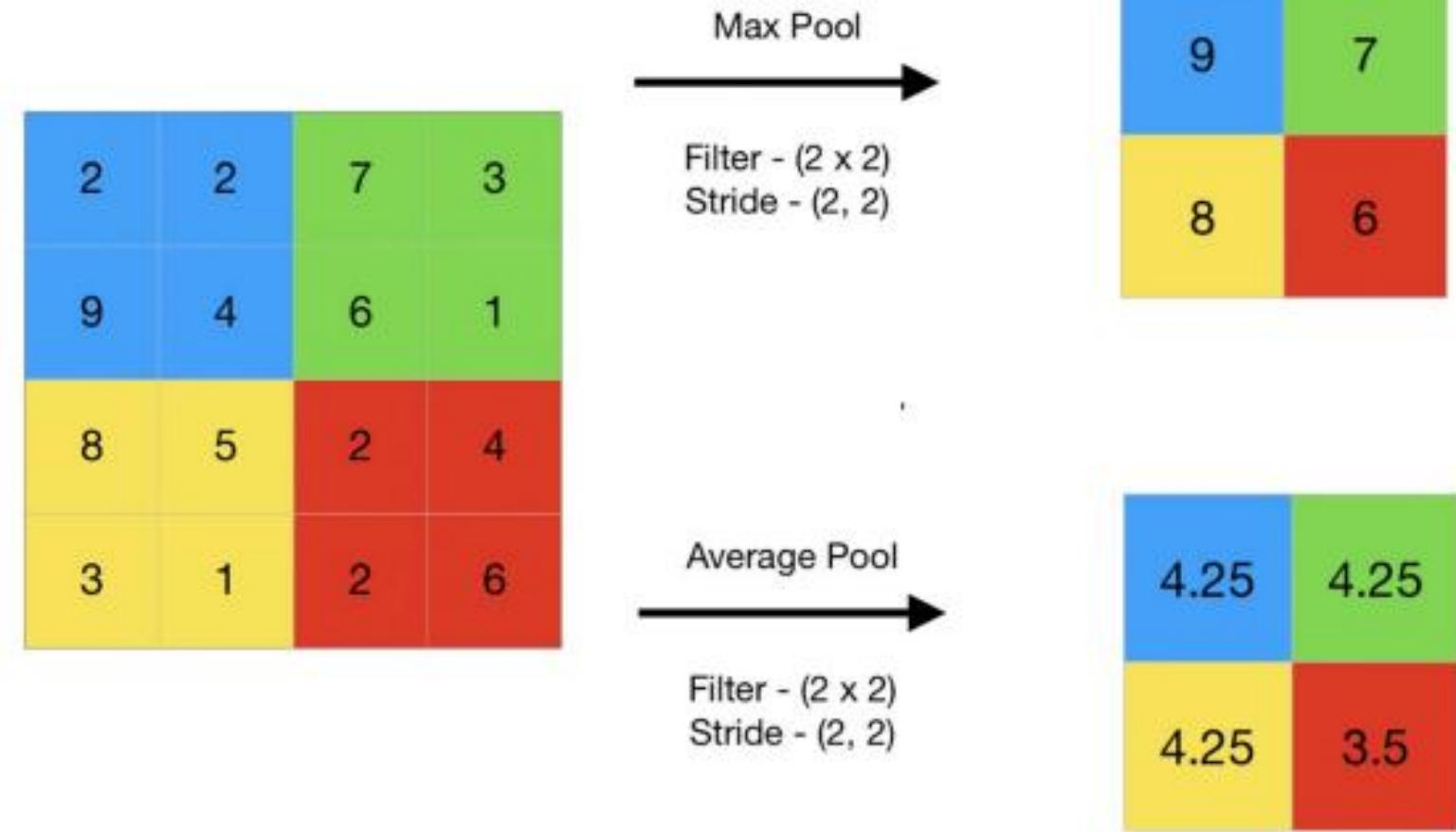
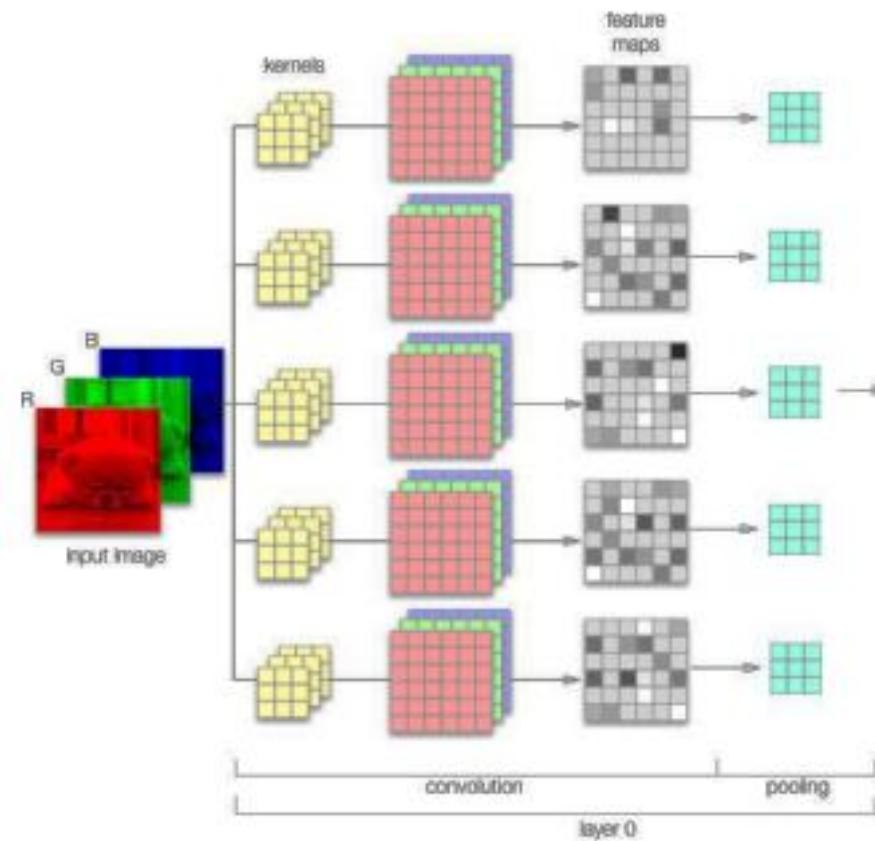
- makes the representations smaller and more manageable
- operates over each activation map independently

- **Pooling** is a form of downsizing that uses a 2D sliding filter. The filter passes over the input slice according to a configure parameter called the stride.
- **Stride** is the number pixels the filter moves across the input slice from one position to the next.
- There are two types of pooling operations: **average pooling** and **max pooling**. However, max pooling is the most common.

```
tf.keras.layers.Max  
Pool2D(  
pool_size=(2,2),  
strides=2  
)
```

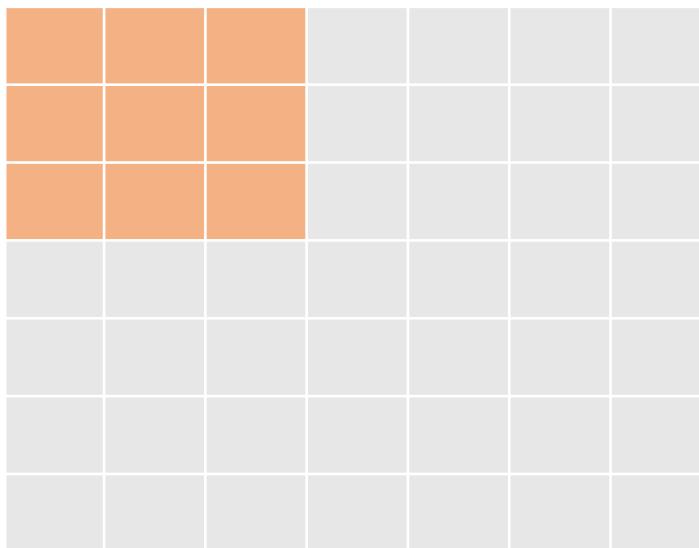


# Pooling

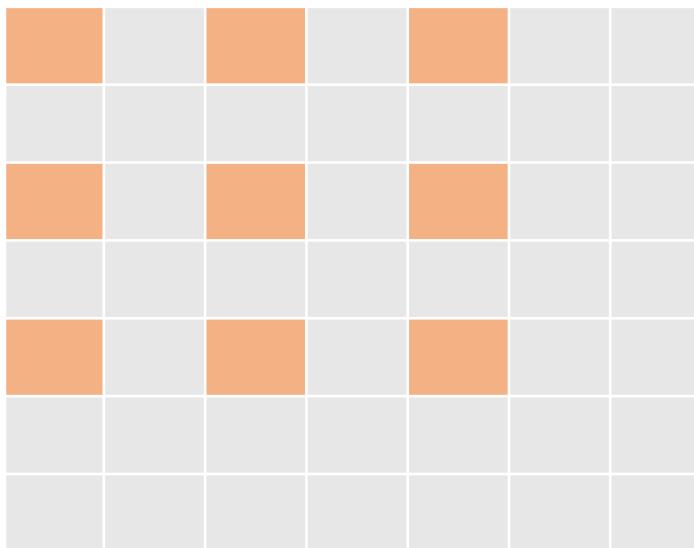


**Reduce dimensionality**

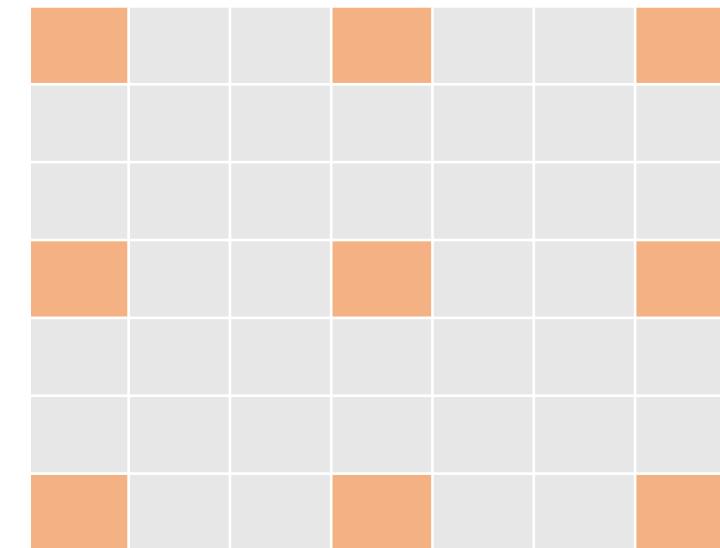
# Dilation



Kernel: 3x3  
Dilation: 1x1

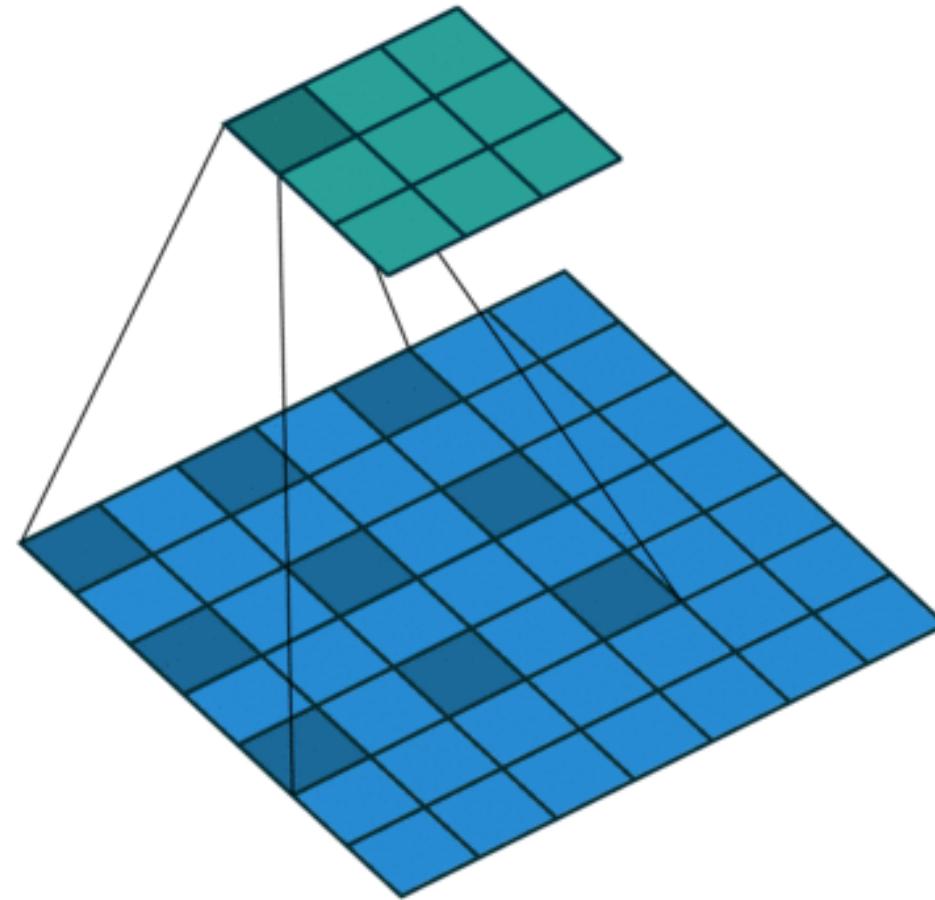


Kernel: 3x3  
Dilation: 2x2



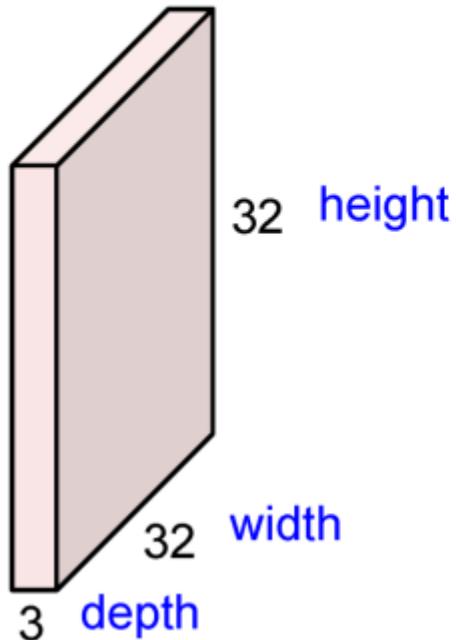
Kernel: 3x3  
Dilation: 3x3

# Dilation

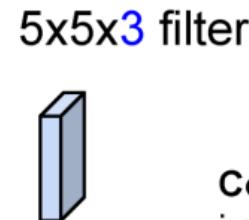
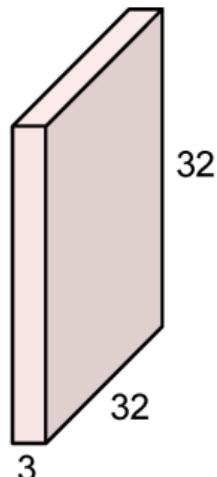


# Convolution Layer

32x32x3 image -> preserve spatial structure



32x32x3 image

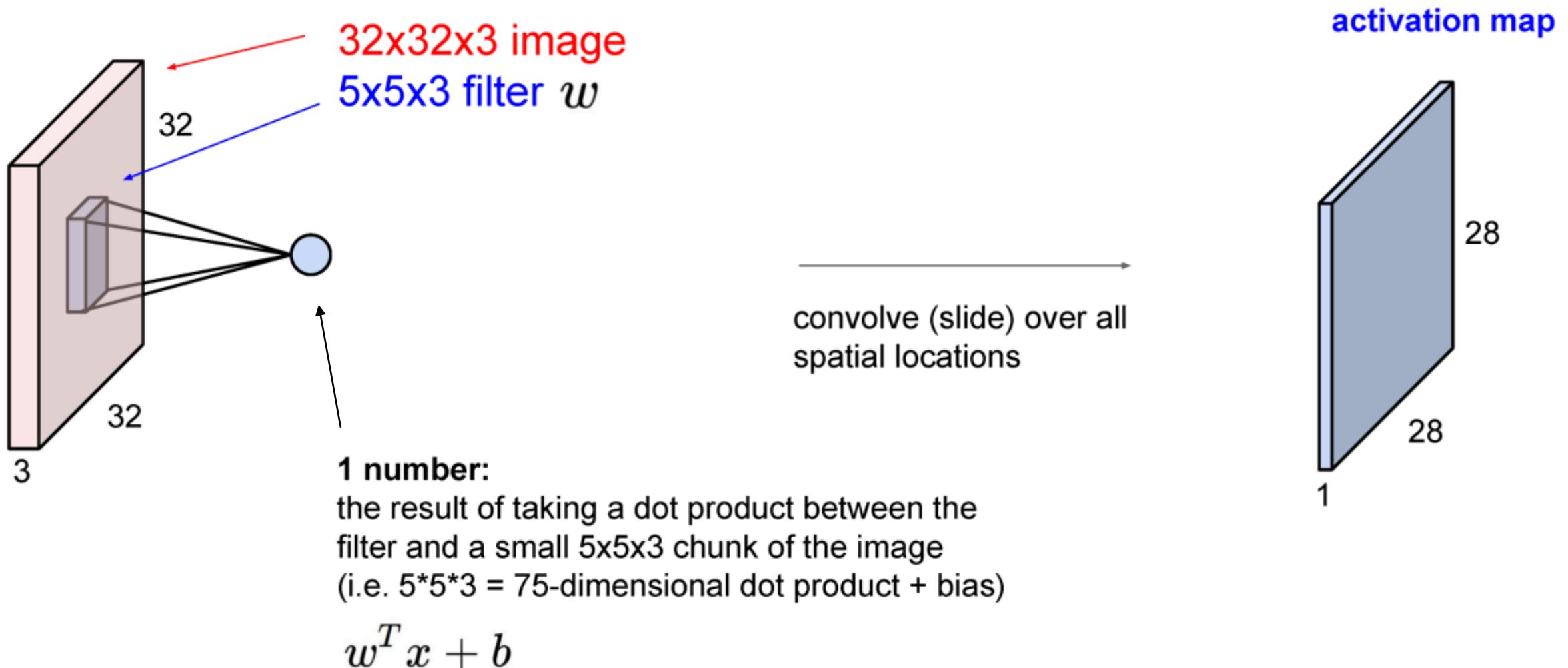


Filters always extend the full depth of the input volume

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

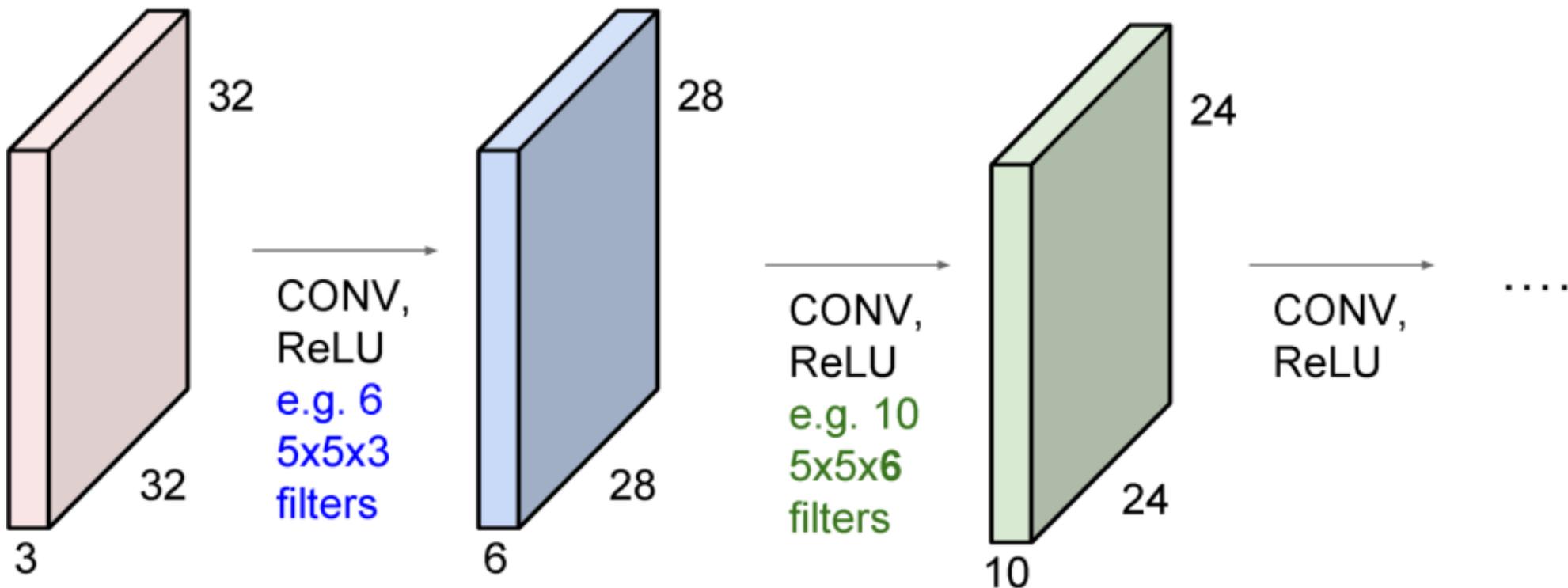
Number of weights:  $5 \times 5 \times 3 + 1 = 76$   
(vs. 3072 for a fully-connected layer)  
(+1 for bias)

# Convolution Layer



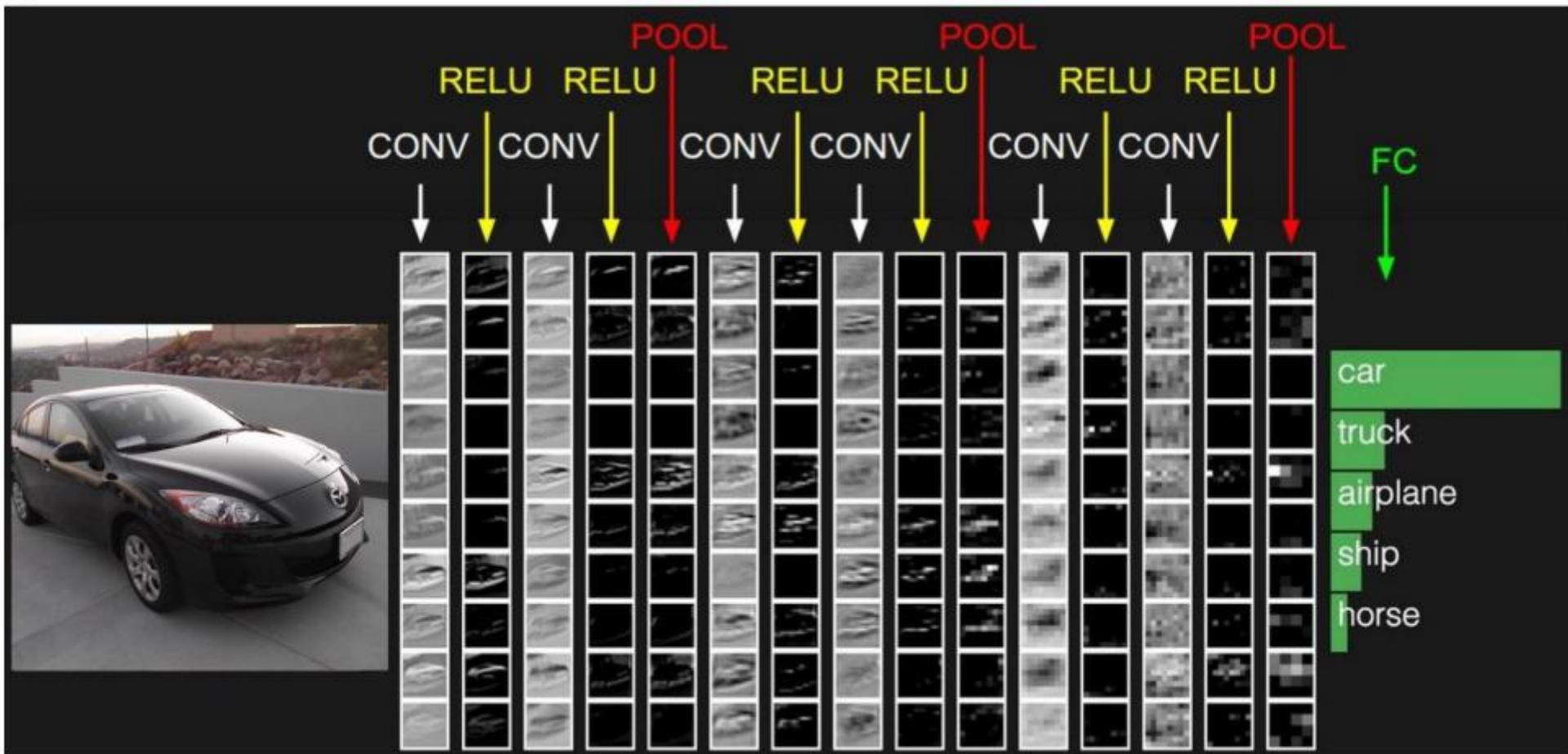
# Convolution Layer

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



# Fully Connected Layer

- Contains neurons that connect to the entire input volume, as in ordinary Neural Networks



# Example: 3@32x32,f:5x5@12,s:1x1,z:0x0

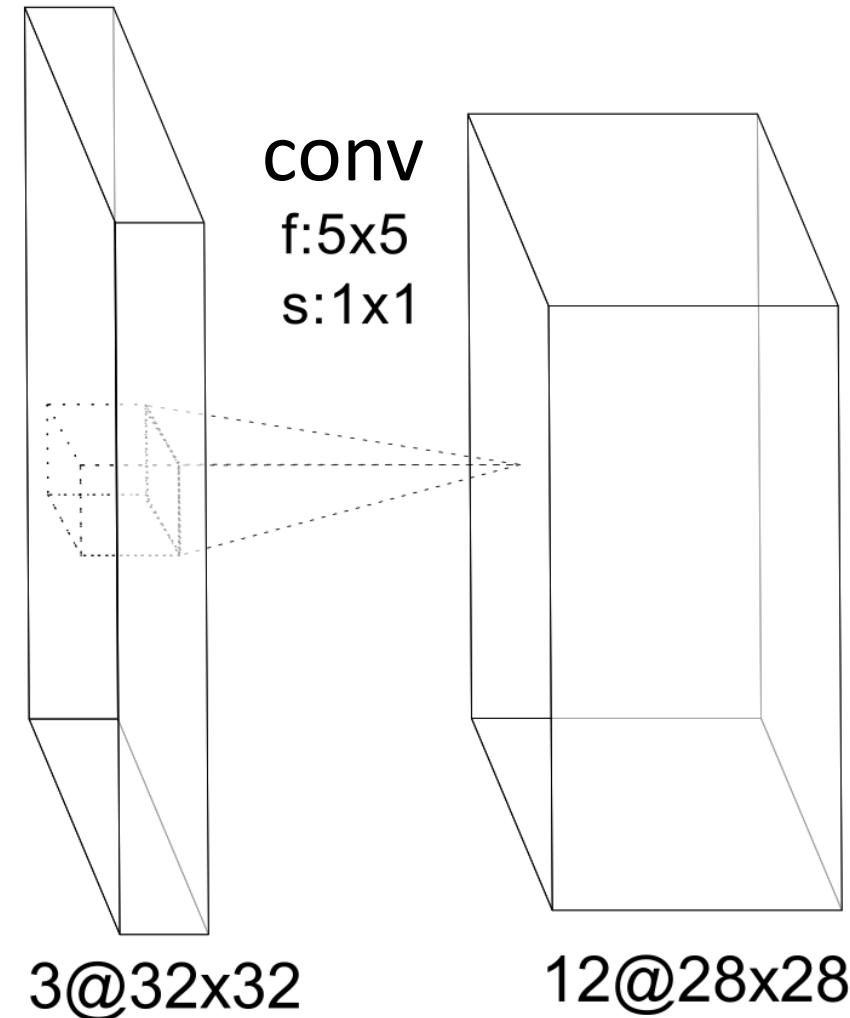
Output: 12@28x28

Number of weights between a first layer and a plane in a second layer:

$$5*5*3 + 1 \text{ (bias)} = 76$$

Total number of weights:

$$(5*5*3 + 1)*12 = 912$$

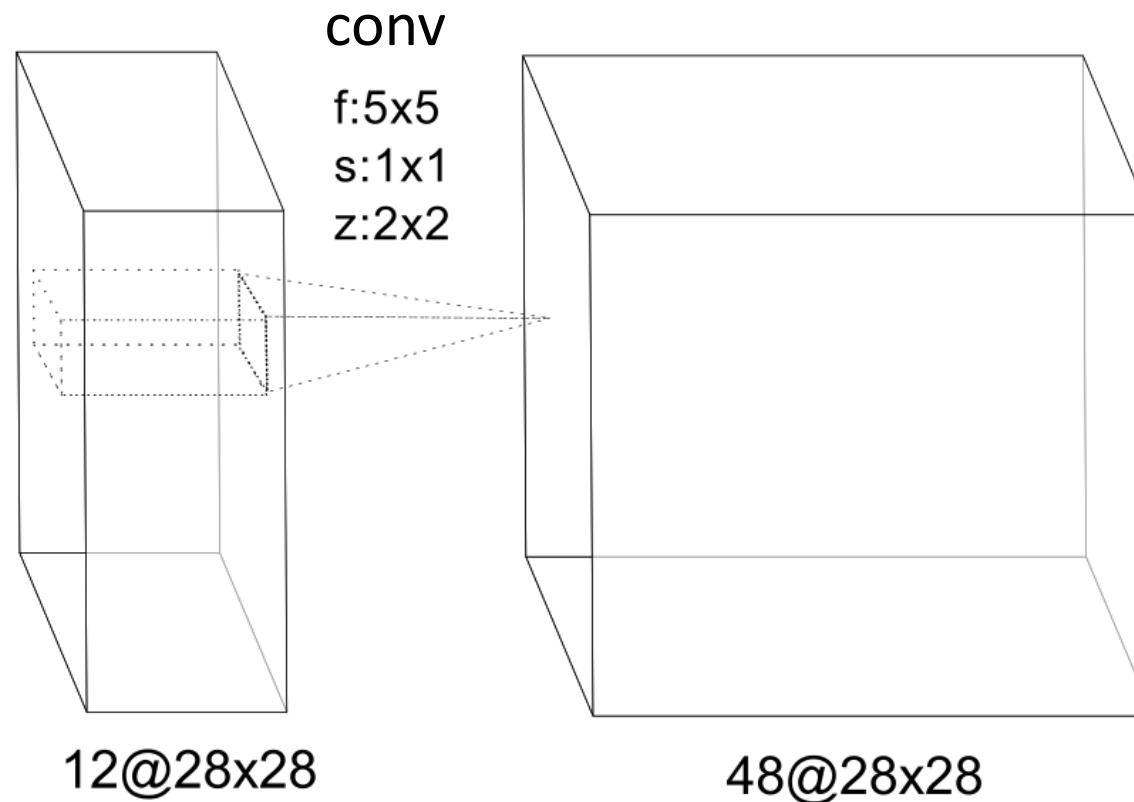


# Example: 12@28x28,f:5x5@48,s:1x1,z:2x2

Output: 48@28x28

Number of weights per plane:  $5*5*12 + 1$  (bias) = 301

Number of weights:  $(5*5*12+1)*48 = 14448$

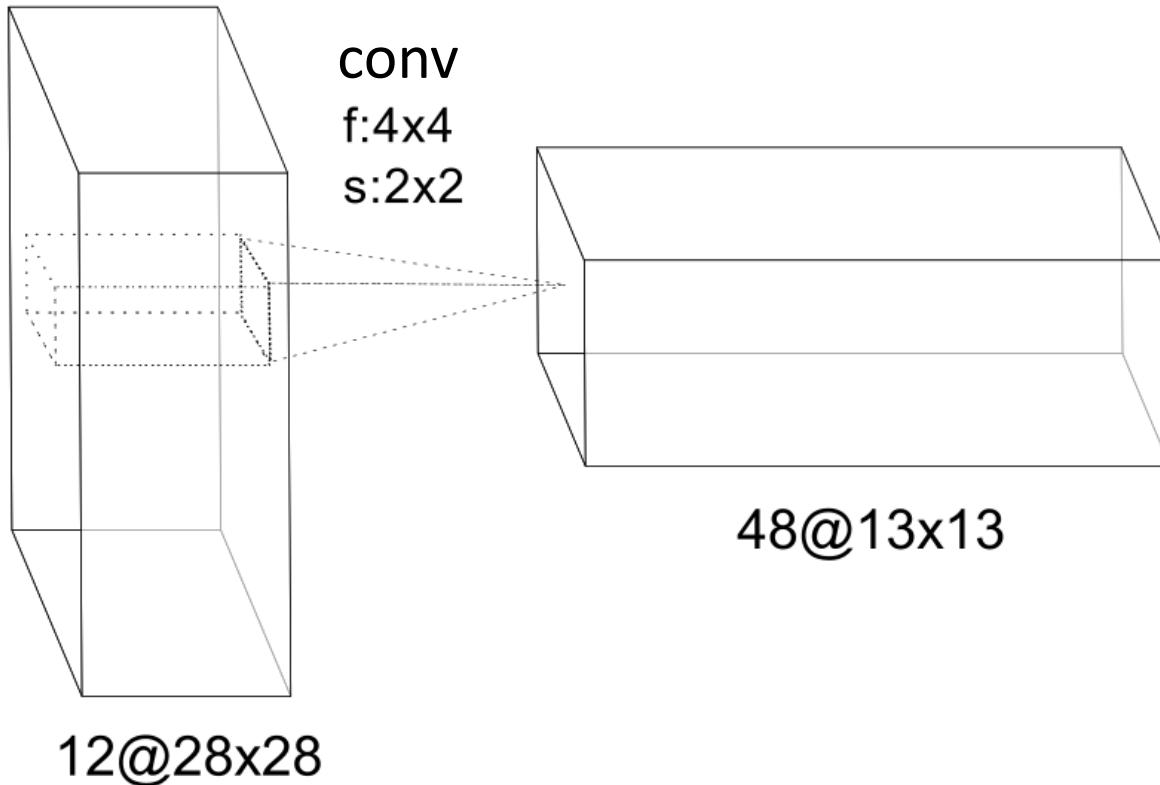


# Example: 12@28x28,f:4x4@48,s:2x2,z:0x0

Output: 48@13x13

Number of weights per plane:  $4*4*12+1=193$

Number of weights:  $(4*4*12+1)*48=9264$



# [ConvNetJS demo: training on CIFAR-10]

**ConvNetJS CIFAR-10 demo**

**Description**

This demo trains a Convolutional Neural Network on the [CIFAR-10 dataset](#) in your browser, with nothing but Javascript. The state of the art on this dataset is about 90% accuracy and human performance is at about 94% (not perfect as the dataset can be a bit ambiguous). I used [this python script](#) to parse the [original files](#) (python version) into batches of images that can be easily loaded into page DOM with img tags.

This dataset is more difficult and it takes longer to train a network. Data augmentation includes random flipping and random image shifts by up to 2px horizontally and vertically.

By default, in this demo we're using Adadelta which is one of per-parameter adaptive step size methods, so we don't have to worry about changing learning rates or momentum over time. However, I still included the text fields for changing these if you'd like to play around with SGD+Momentum trainer.

Report questions/bugs/suggestions to [@karpathy](#).

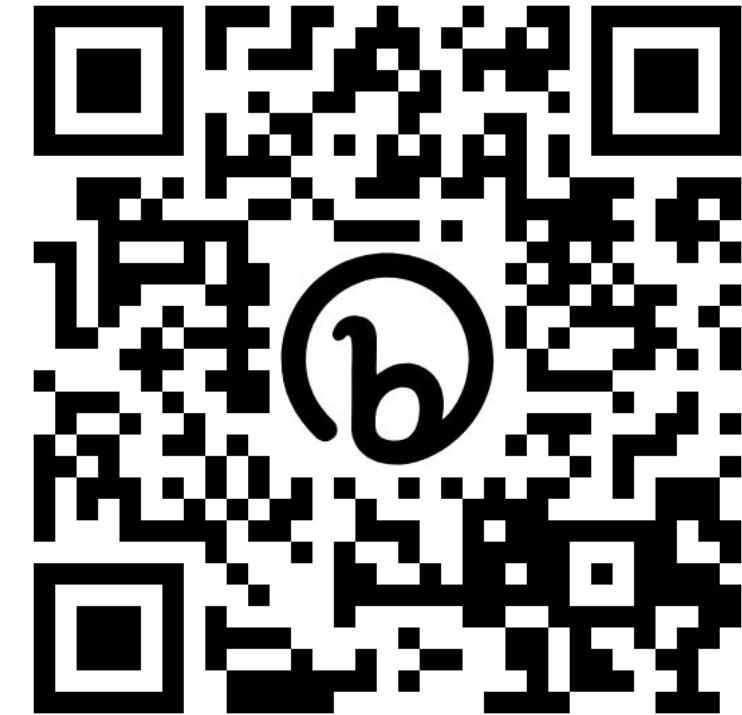


<https://cs.stanford.edu/people/karpathy/convnetjs/demo/cifar10.html>

To be continued next week

Please, don't forget  
to send feedback:

<https://bit.ly/bme-dl>



# Thank you for your attention

Dr. Mohammed Salah Al-Radhi  
[malradhi@tmit.bme.hu](mailto:malradhi@tmit.bme.hu)

(slides by: Dr. Bálint Gyires-Tóth)

08 October 2024

