

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



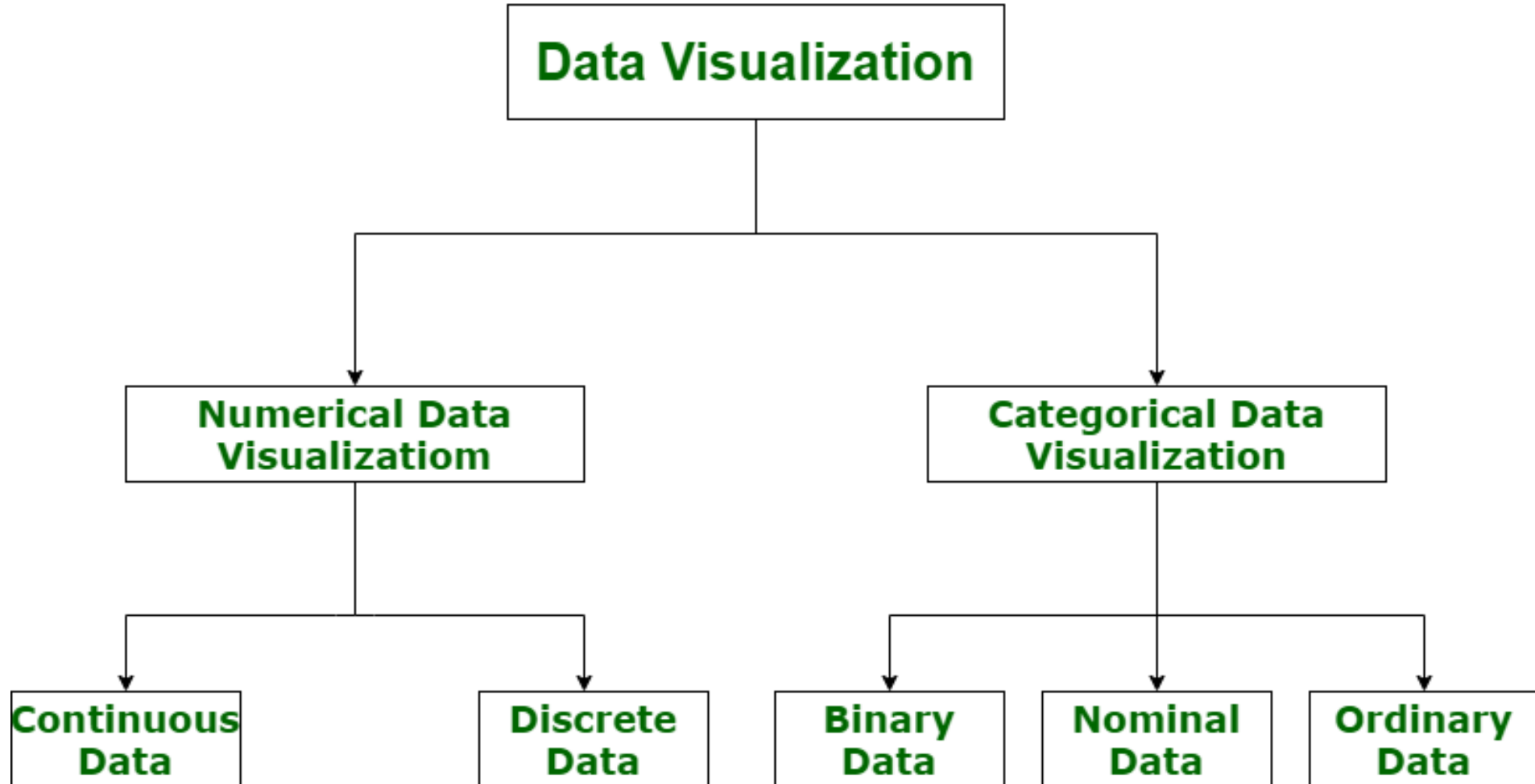
images by Mark P. Witton

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

A network graph with nodes and edges, overlaid with a white banner containing the word 'Classification'. The graph consists of numerous nodes connected by thin lines, forming a complex web. The nodes are colored in shades of blue and green. The white banner is positioned horizontally across the center of the image, with the word 'Classification' written in bold black text.

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?

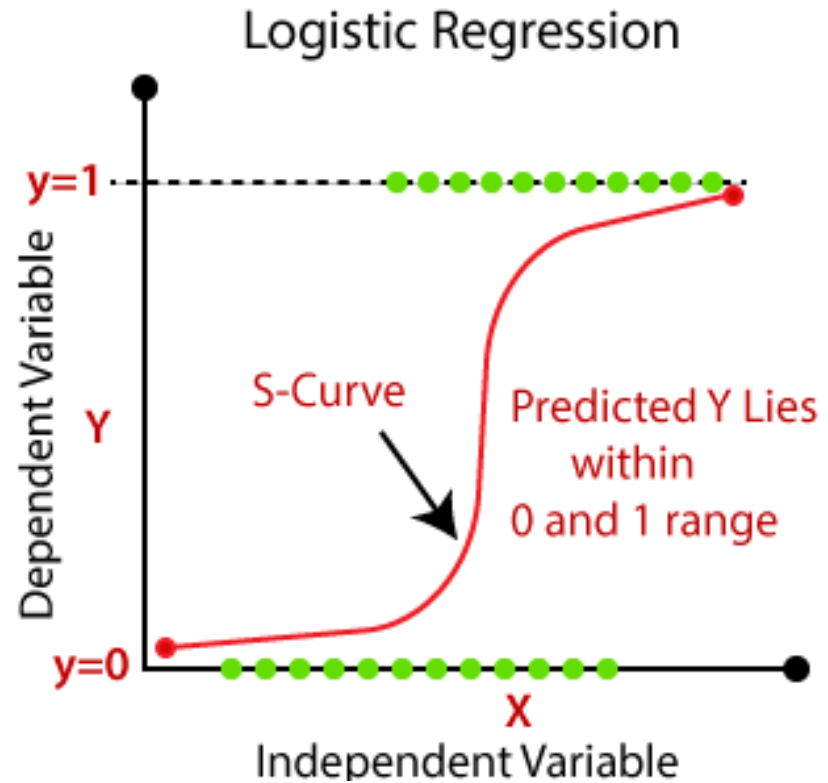
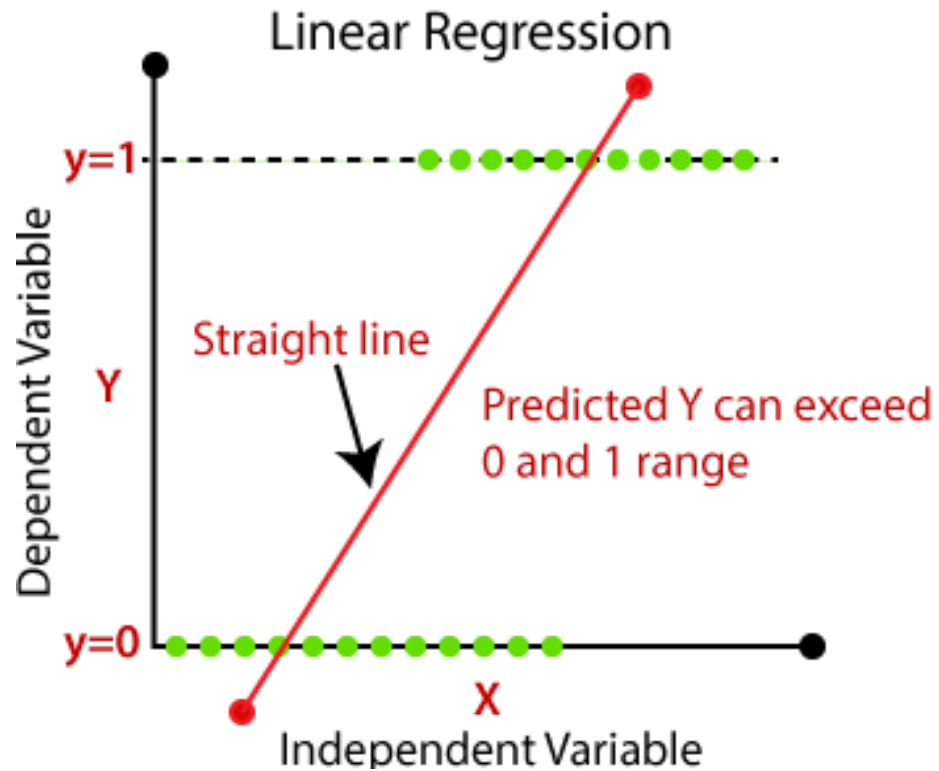


- Is this an apple?



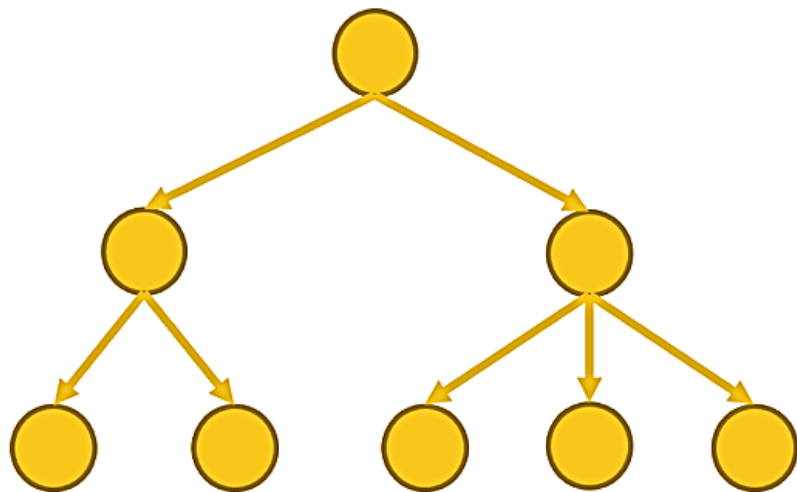
Given a number of examples, identify to which class a given observation belongs to.

Regression

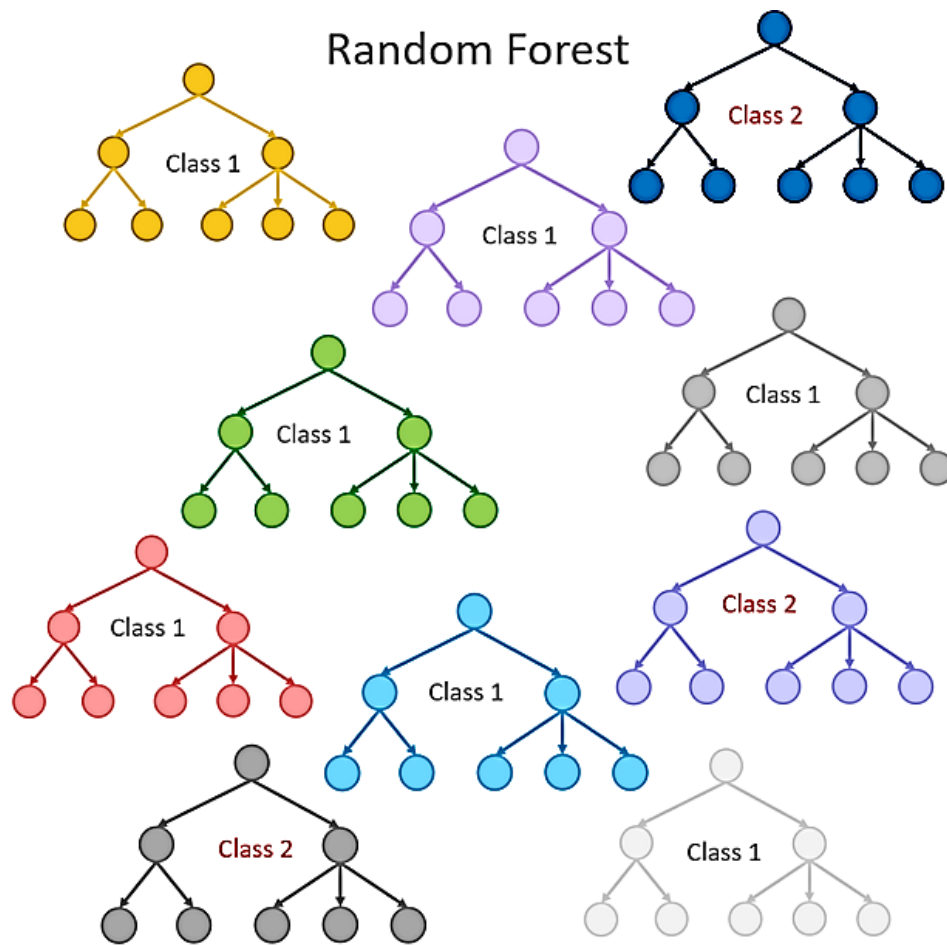


Tree

Single Decision Tree



Random Forest



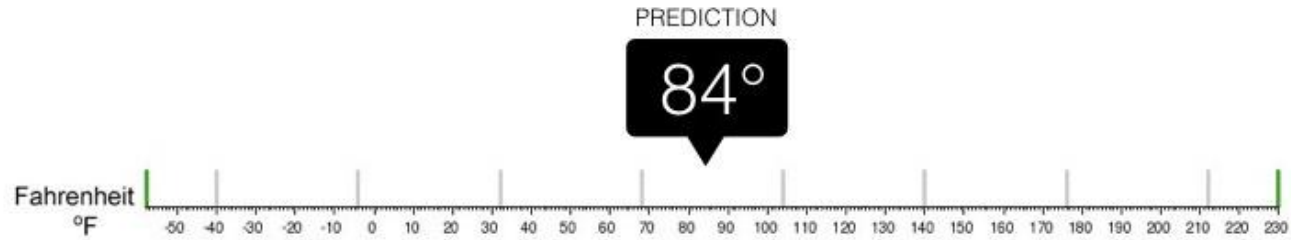
Regression vs. Classification

- **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

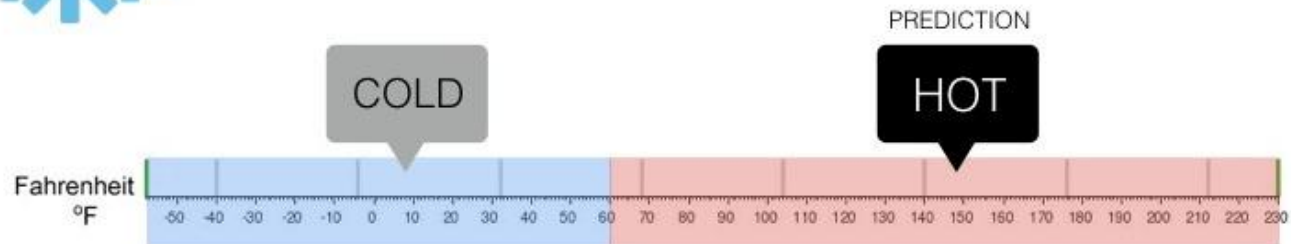
Regression vs. Classification



What is the temperature going to be tomorrow?



Will it be Cold or Hot tomorrow?

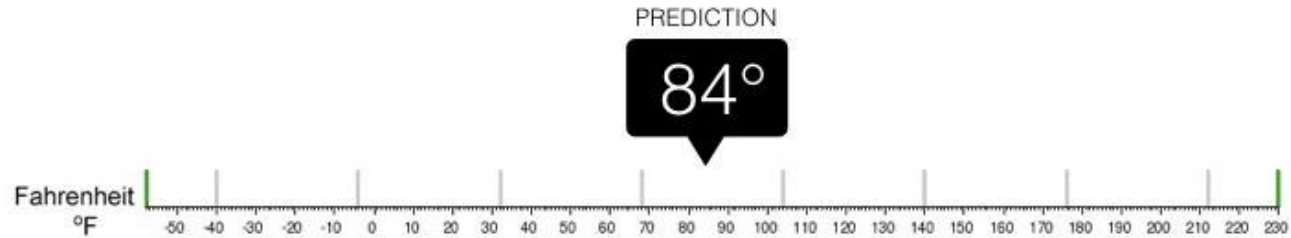


Regression vs. Classification



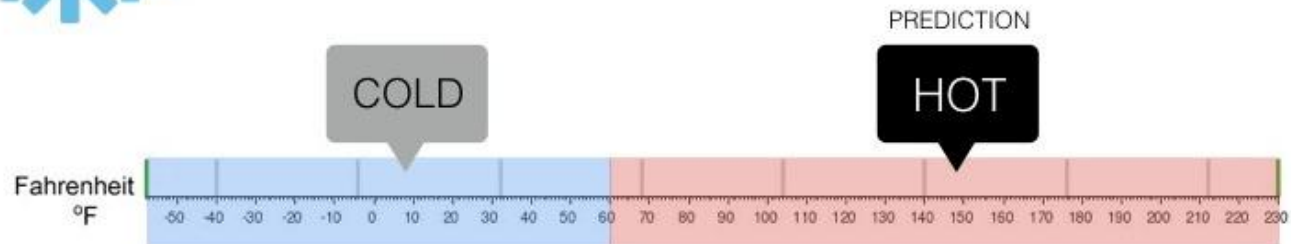
Regression

What is the temperature going to be tomorrow?

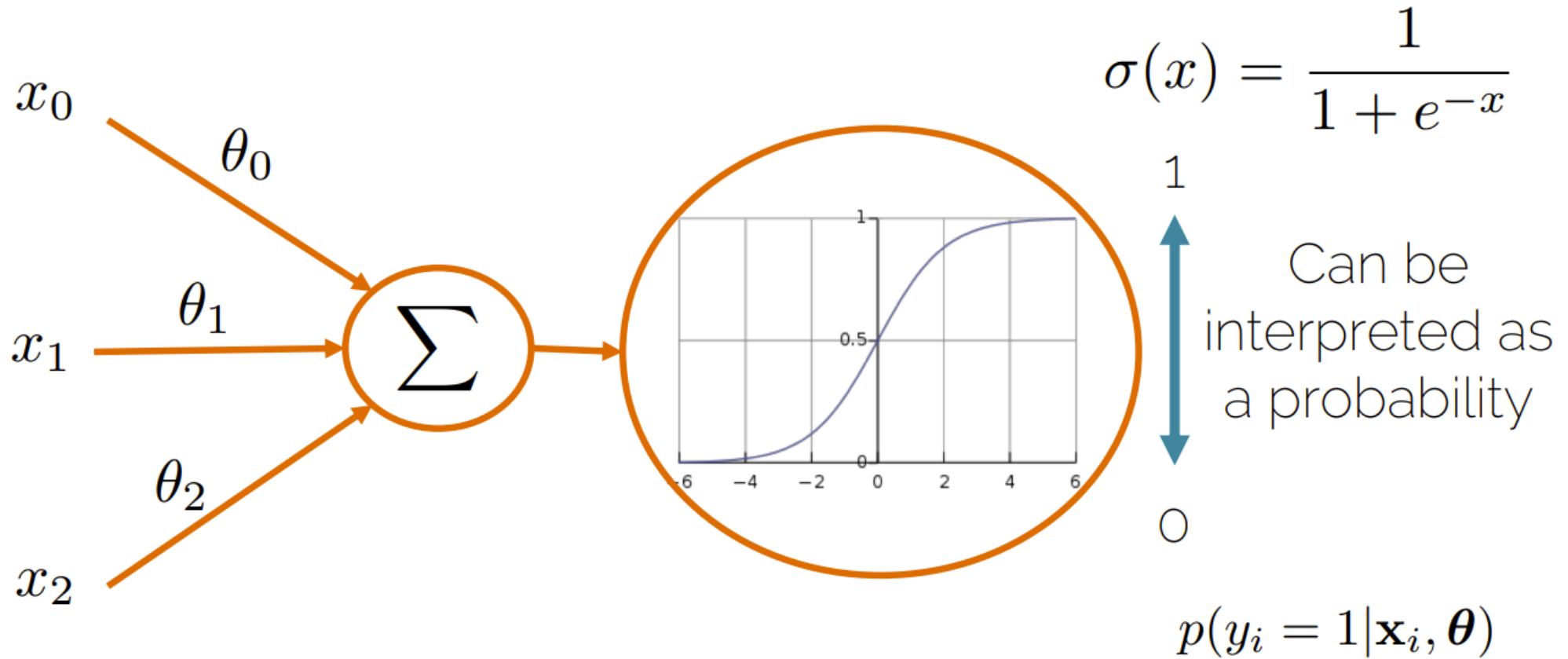


Classification

Will it be Cold or Hot tomorrow?



Sigmoid for binary predictions



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

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: $\frac{-(\ln(0.3) + \ln(0.8) + \ln(0.8))}{3} \approx 0.91$

MSE: $\left. \begin{array}{l} (0.3 - 1)^2 + (0.4 - 0)^2 + (0.3 - 0)^2 = 0.74 \\ (0.1 - 0)^2 + (0.8 - 1)^2 + (0.1 - 0)^2 = 0.06 \\ (0.1 - 0)^2 + (0.1 - 0)^2 + (0.8 - 1)^2 = 0.06 \end{array} \right\} \frac{0.74 + 0.06 + 0.06}{3} = 0.286$

A network diagram with nodes and edges, rendered in light blue and teal colors, serves as a background for the slide. The nodes are small circles, and the edges are thin lines connecting them, creating a complex web-like structure.

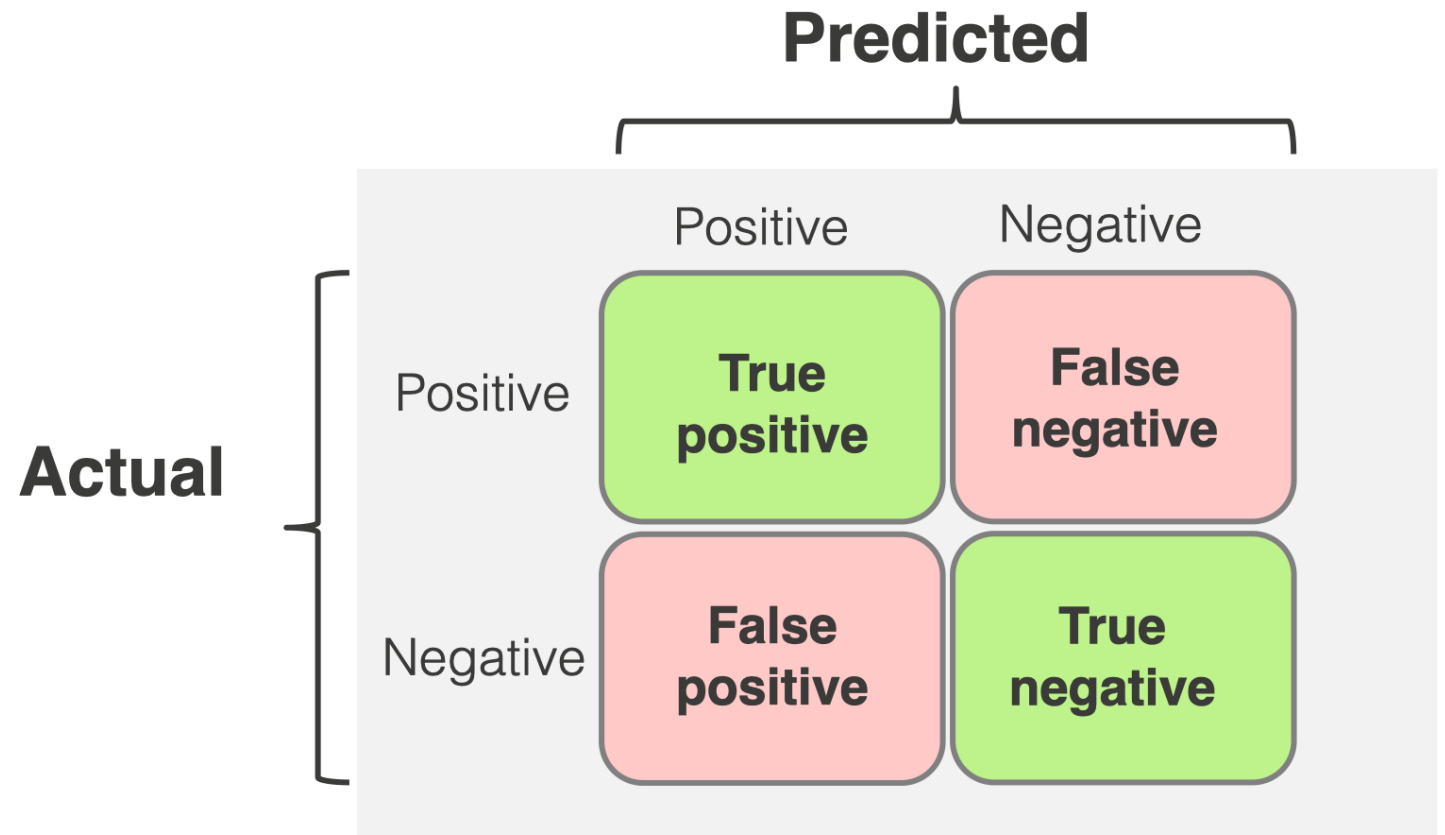
Evaluation methods

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:

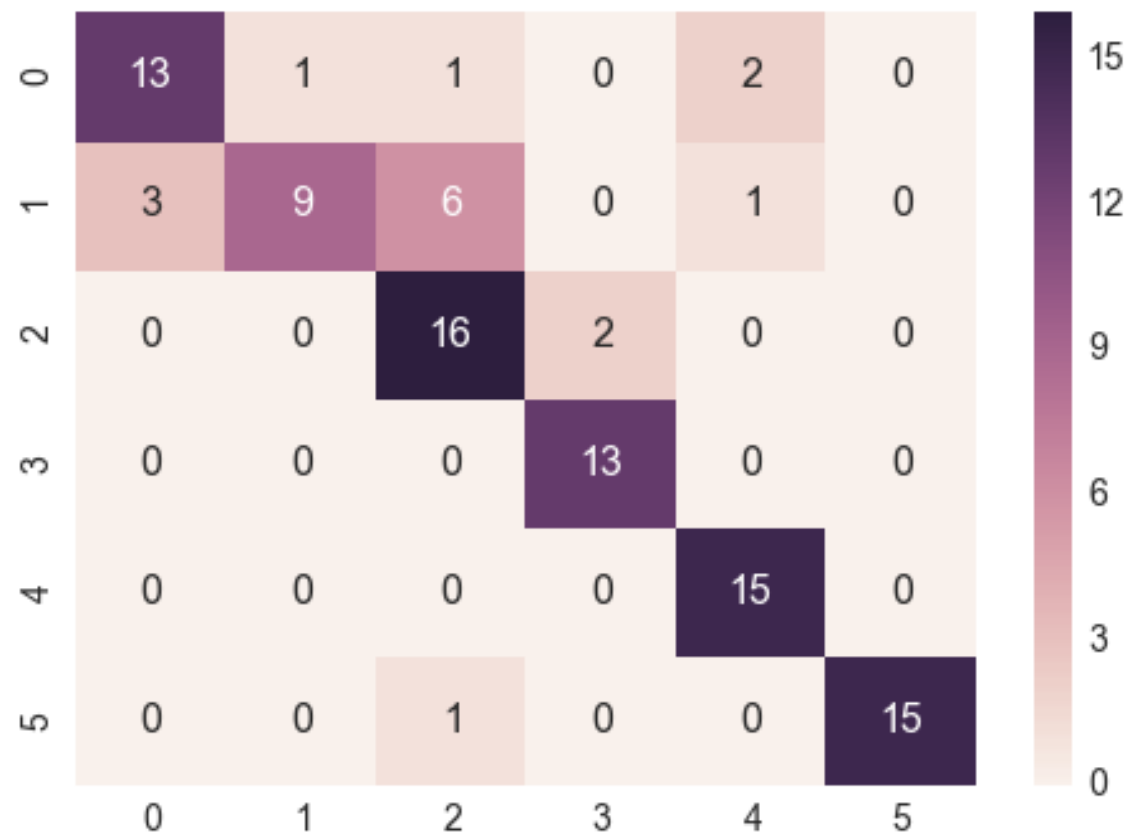
		Y_predicted		
		Class	Rocket	Car
Y_target	Rocket	15	3	8
	Car	4	23	1
	Airplane	9	2	19

Confusion matrix: binary classification

Rocket		Y_predicted	
		Positive	Negative
Y_target	Condition		
	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

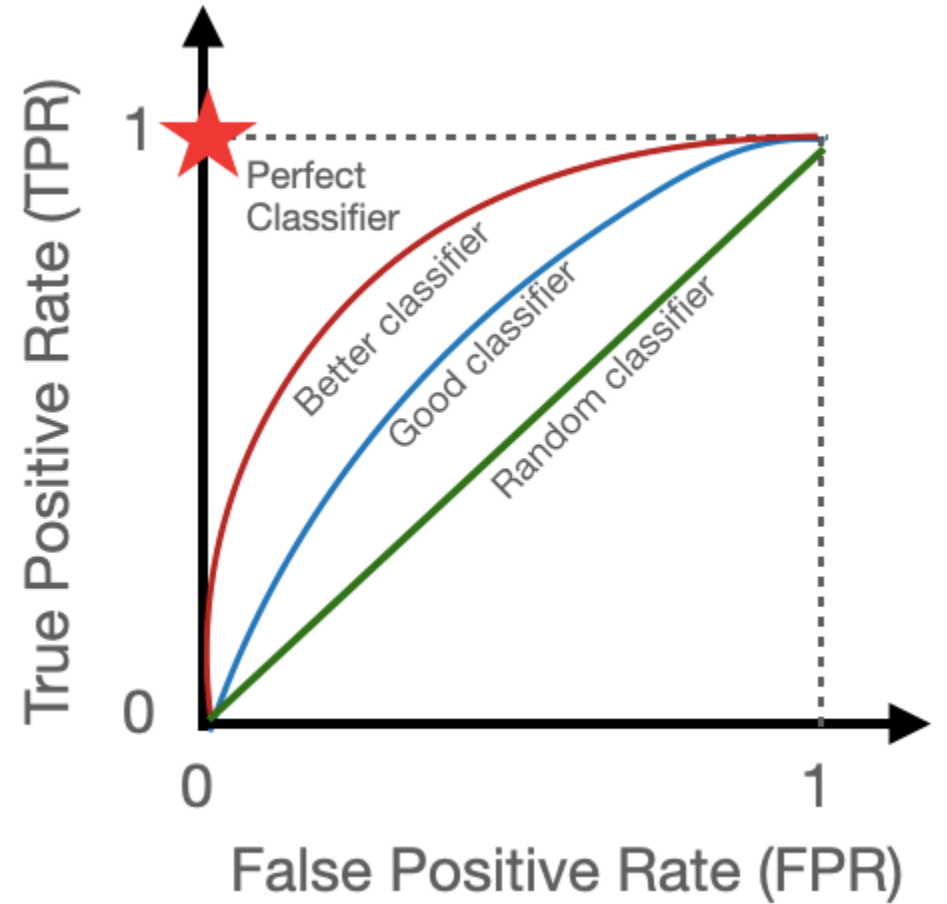
- **Macro:** calculate the value for all classes separately, and take the average.
- **Micro:** calculate the values for all classes, and calculate precision/recall/F1 from that.
- **Weighted:** weights the classes based as a function of the imbalances. Those classes which have more values have higher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

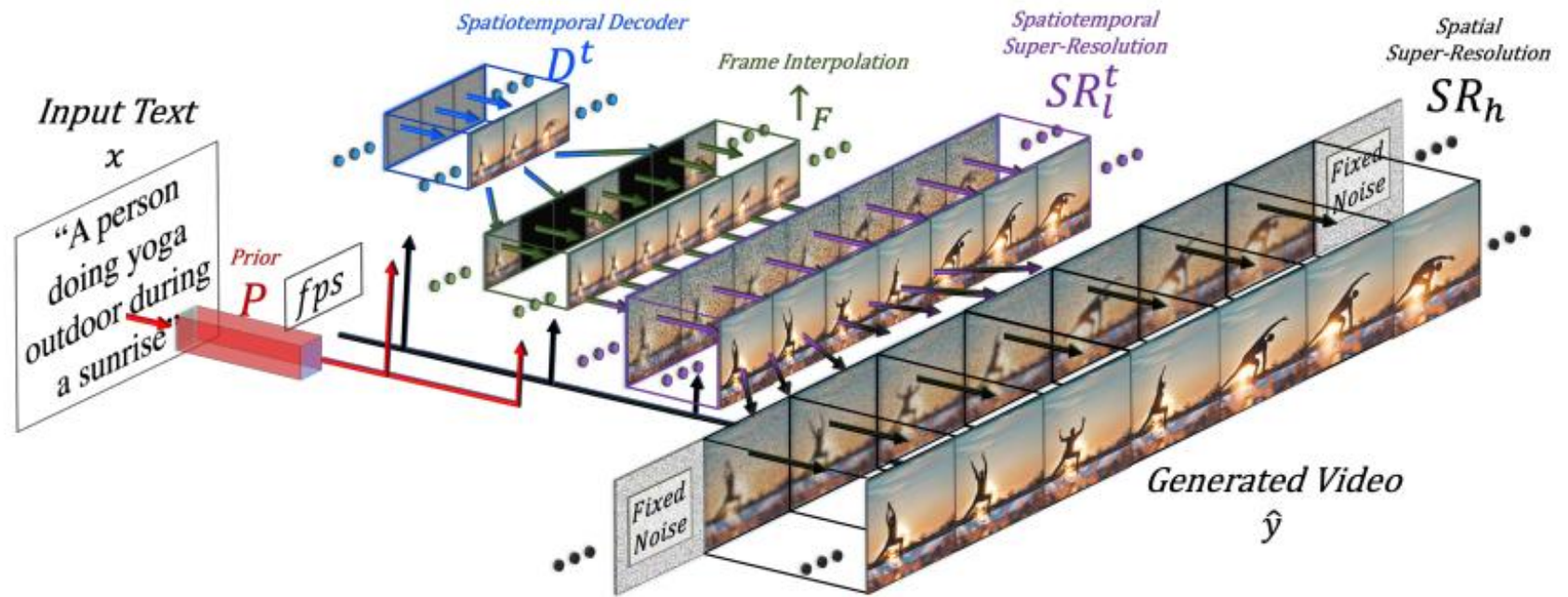
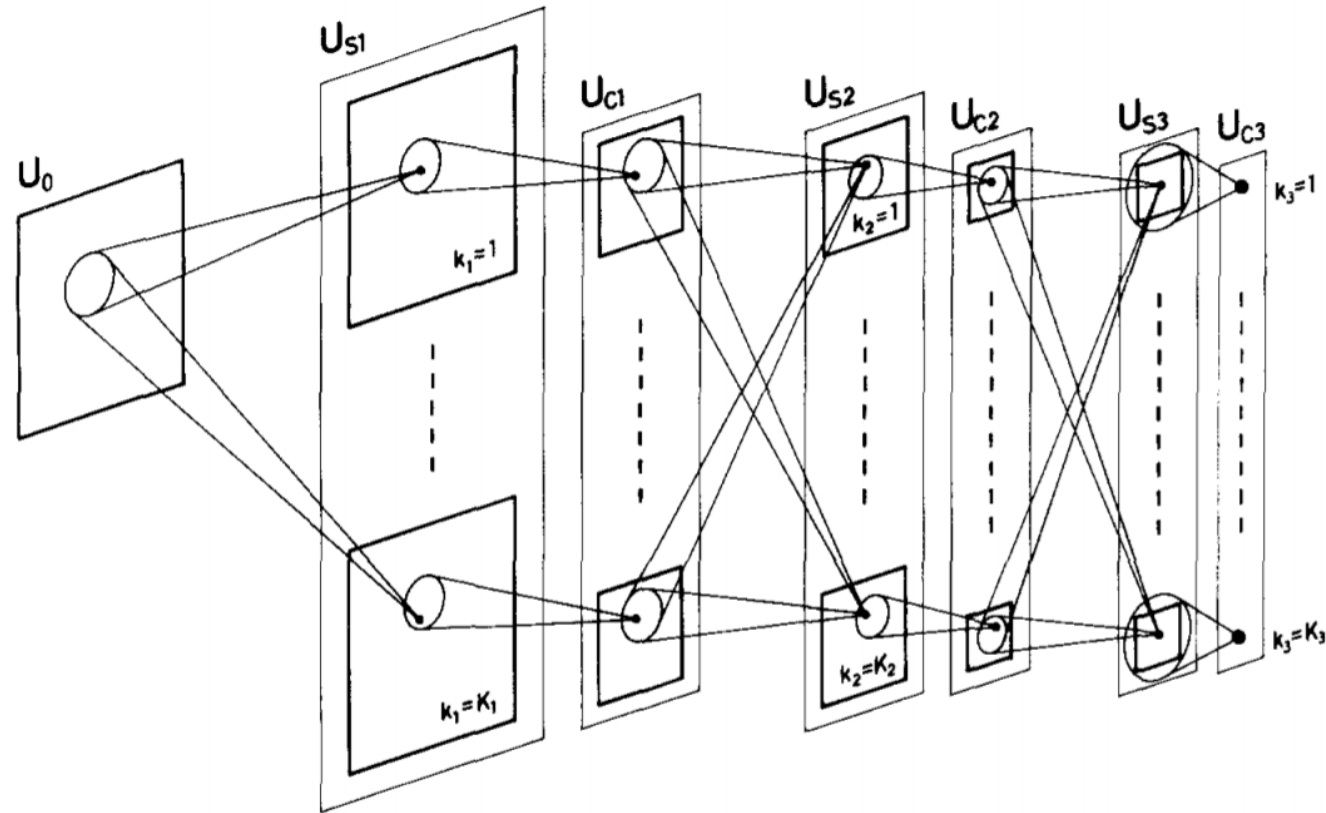


Figure 2: **Make-A-Video high-level architecture.** Given input text x translated by the prior P into an image embedding, and a desired frame rate fps , the decoder D^t generates 16 64×64 frames, which are then interpolated to a higher frame rate by \uparrow_F , and increased in resolution to 256×256 by SR_i^t and 768×768 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.



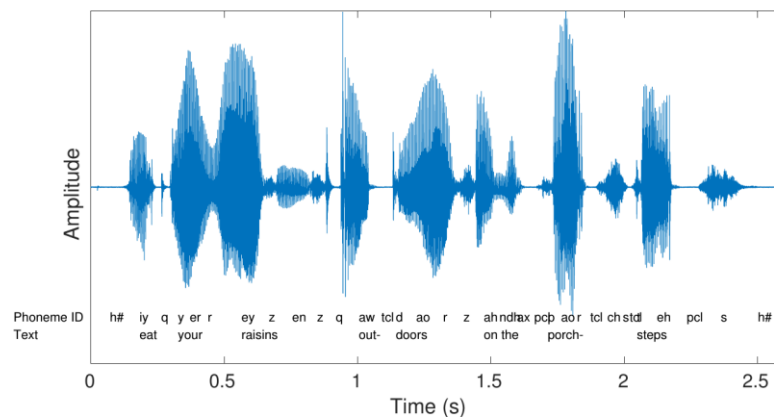
Major milestones

- **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:

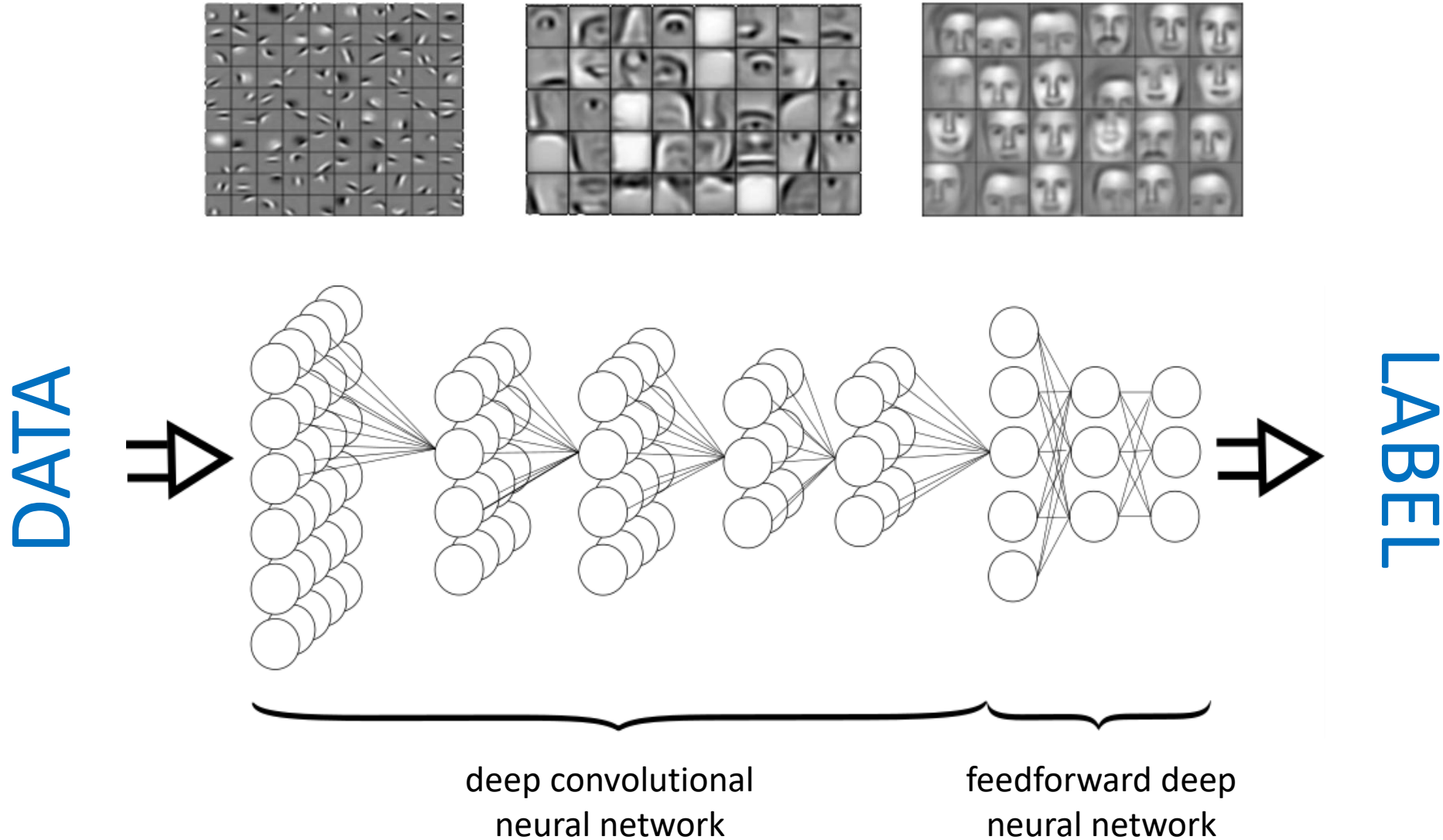
- 1D signals (e.g., speech waveforms)

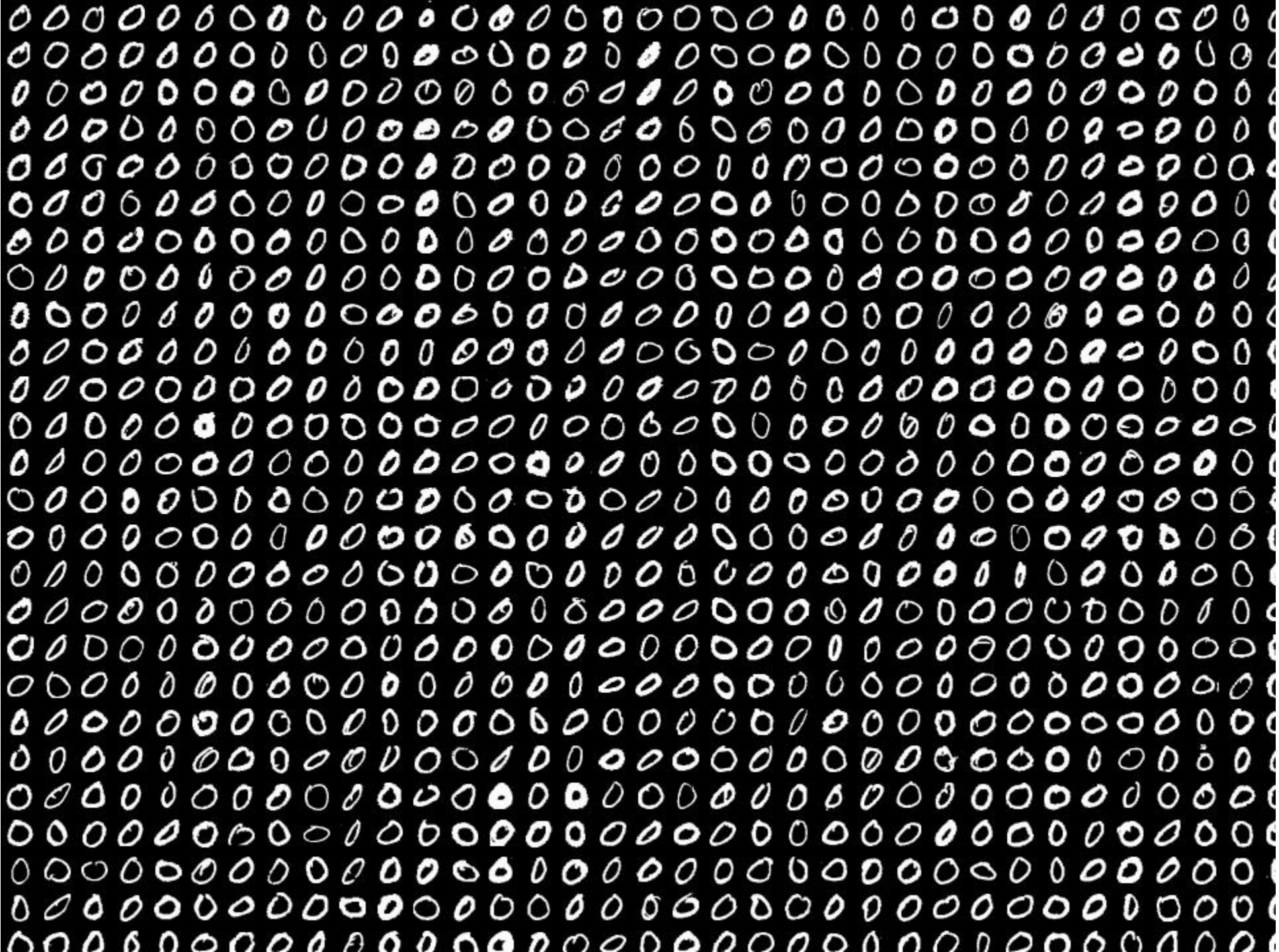


- 2D signals (e.g., images)

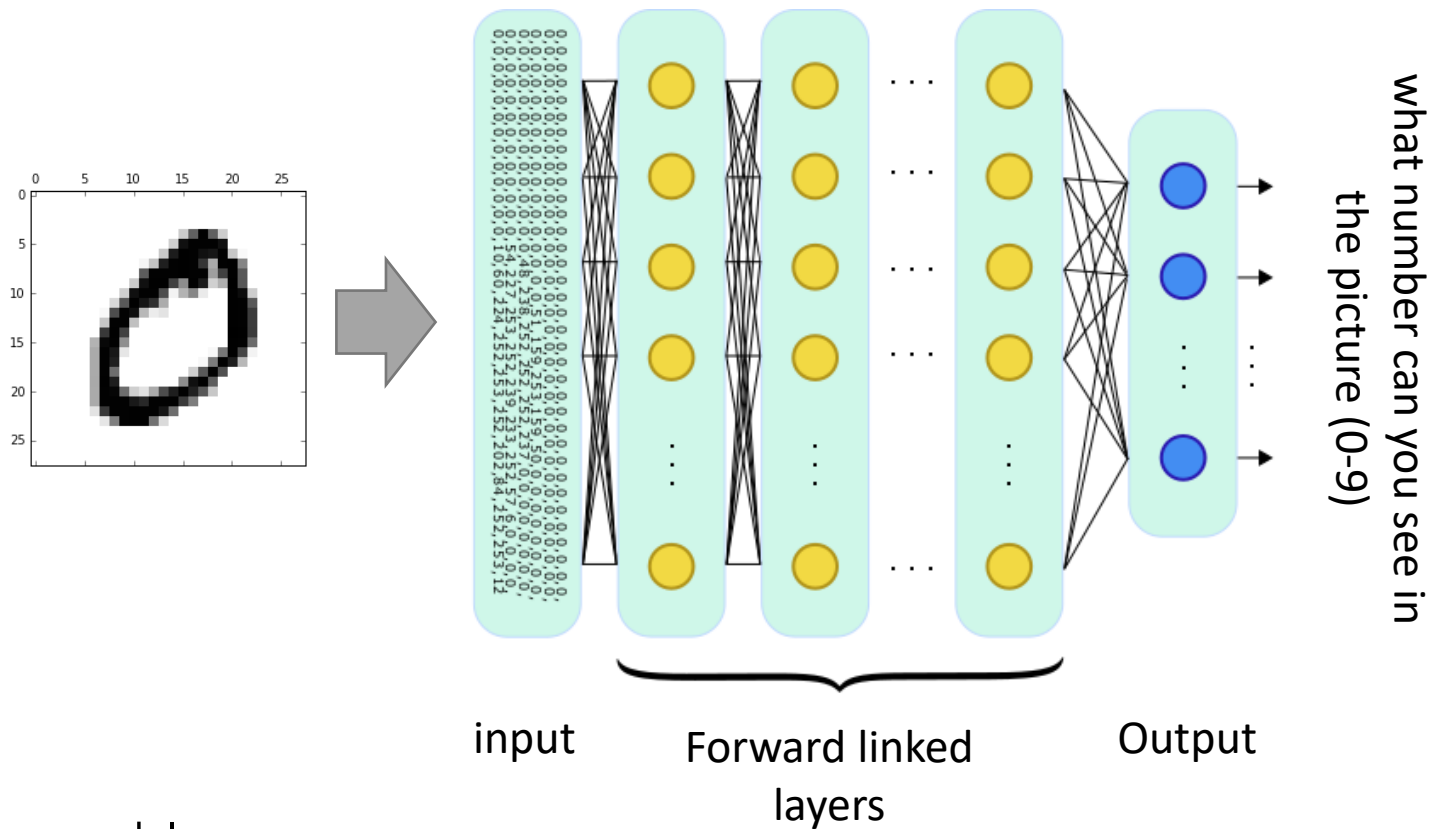


General structure of a CNN network





Feedforward network approach



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 _{x1}	1 _{x0}	1 _{x1}	0	0
0 _{x0}	1 _{x1}	1 _{x0}	1	0
0 _{x1}	0 _{x0}	1 _{x1}	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

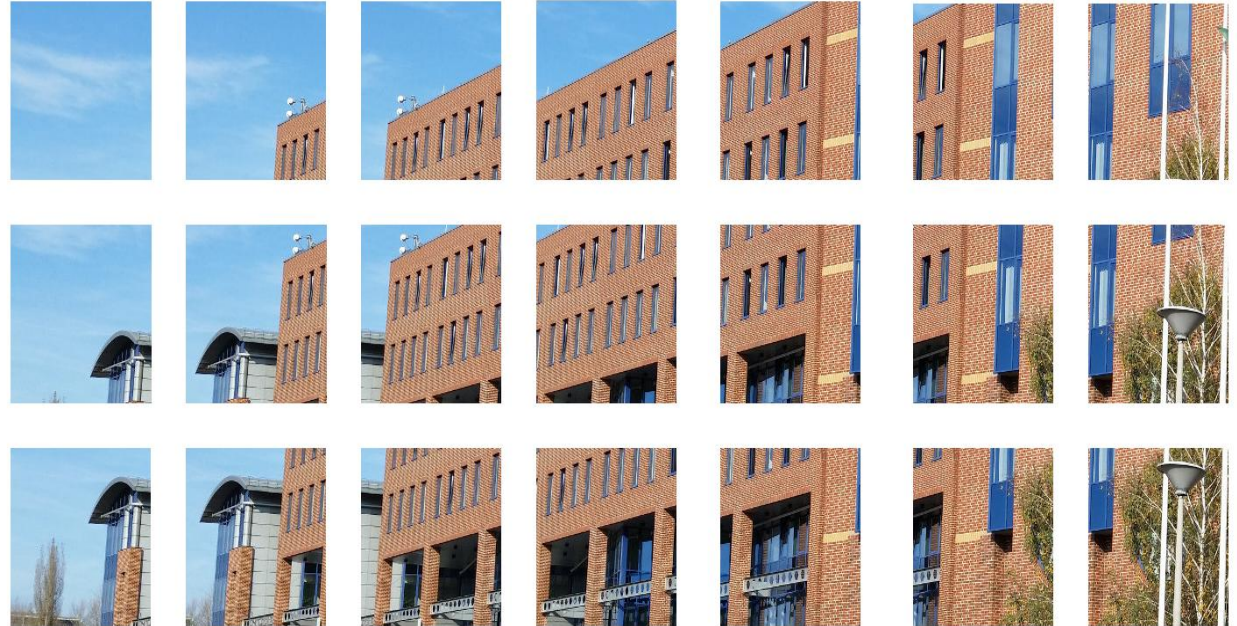
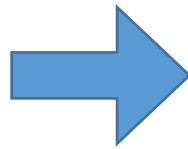
Dependence on space
and/or time

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.

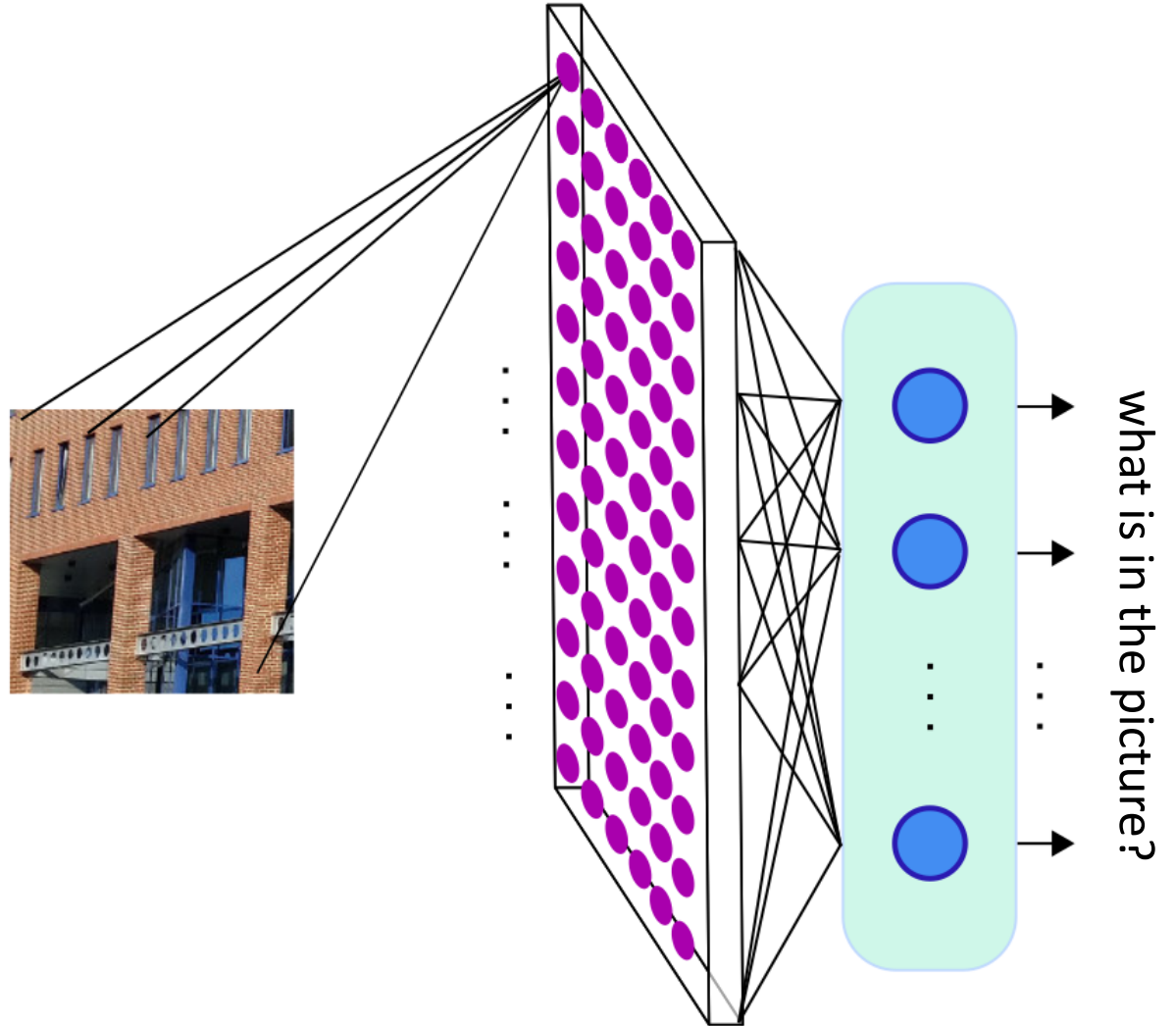


Building blocks: convolutional layer 1

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

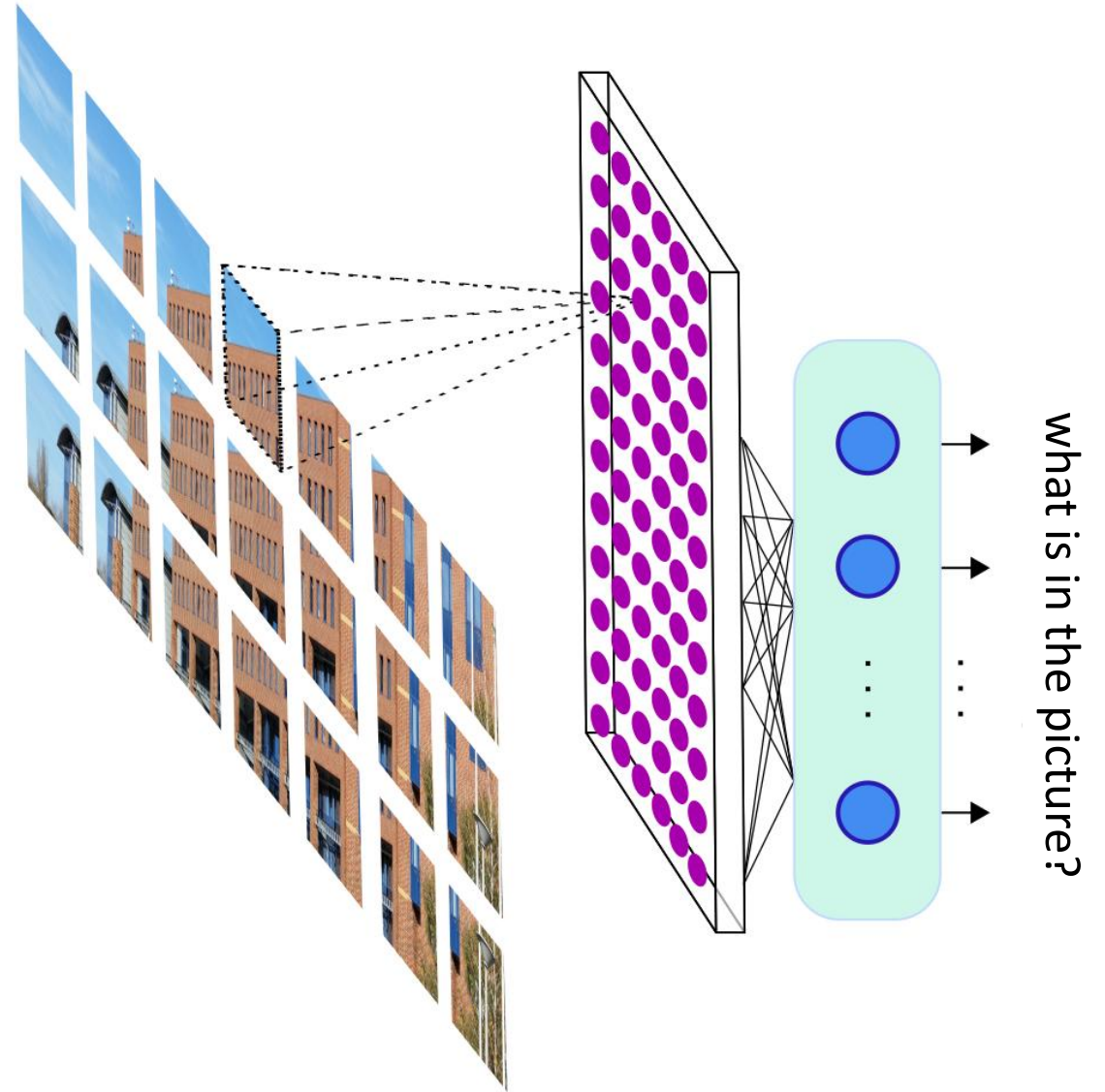
Shared weights:

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: `{dilataion_x}x{dilataion_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

0	1	2	1	3	1	4
2	3	1	0	0	1	2
3	4	2	1	1	0	1
.
.
.
.
.

Input

*

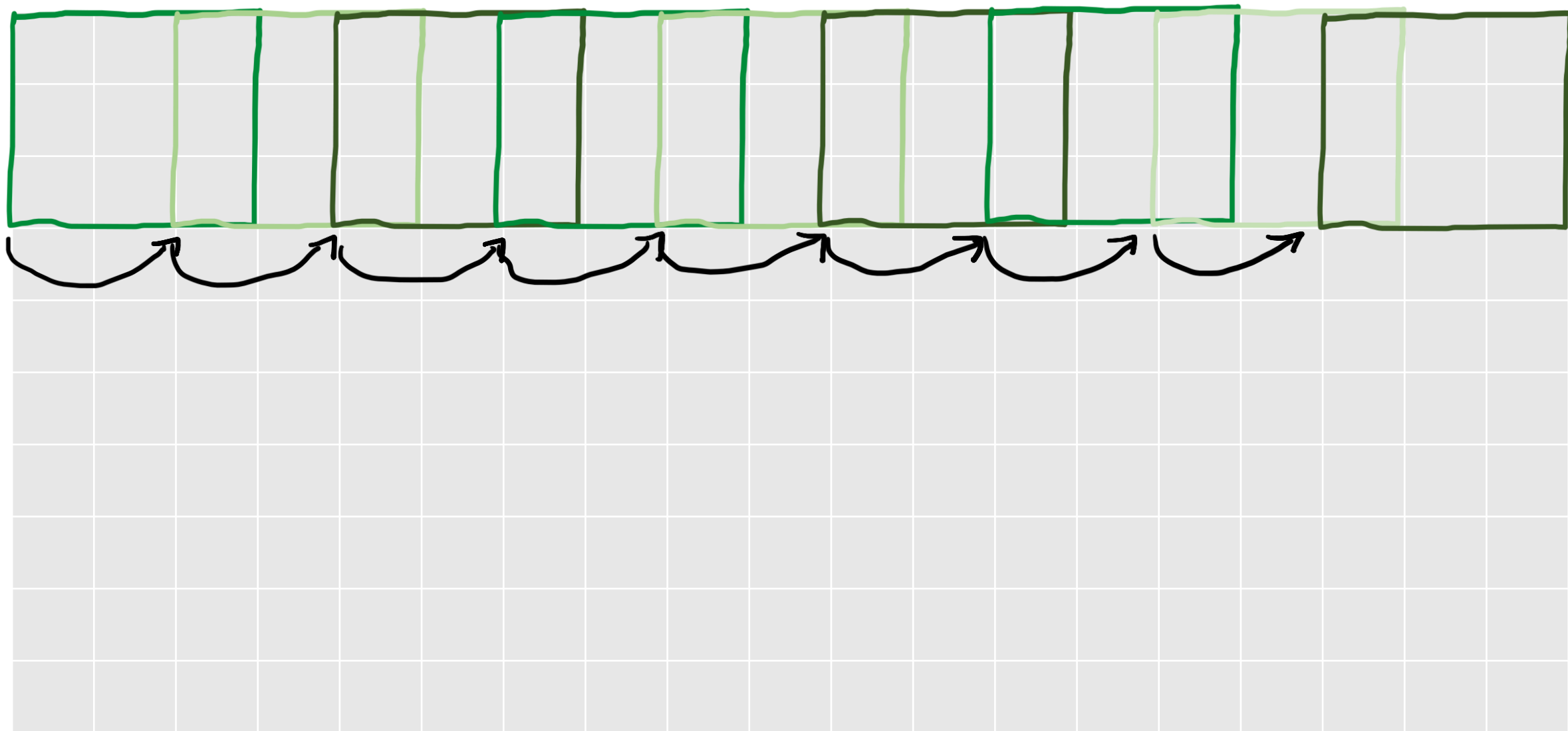
2	1
1	3

=

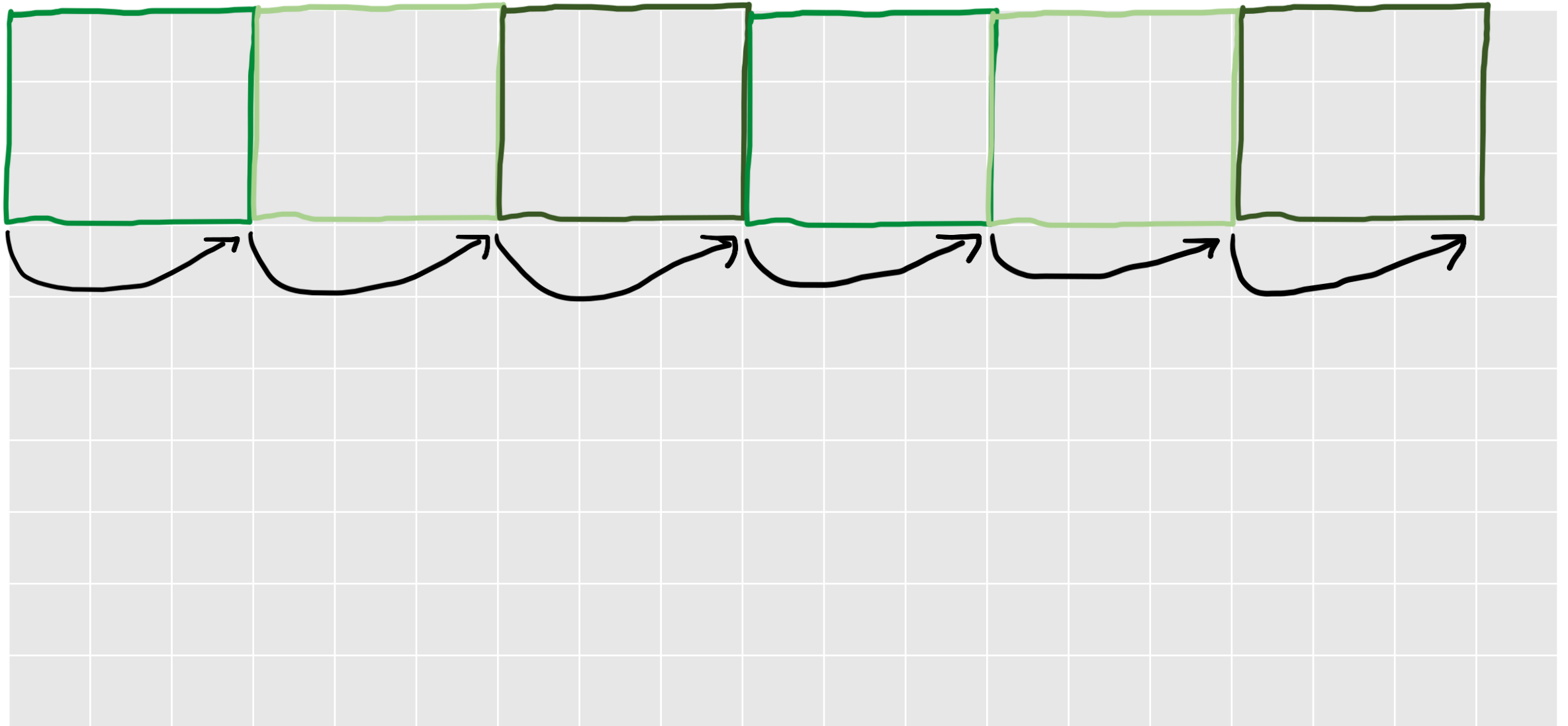
12	10	6	5	10	13
22	17	7	4	2	7
.
.
.
.
.
.
.

Output

Stride = 2

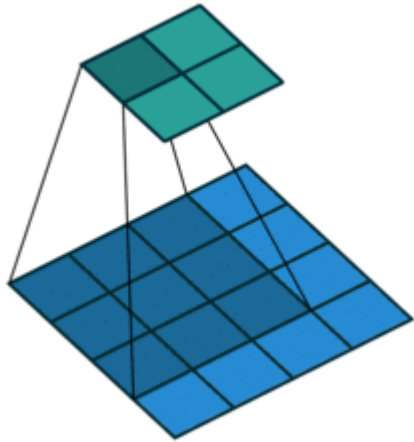


Stride = 3

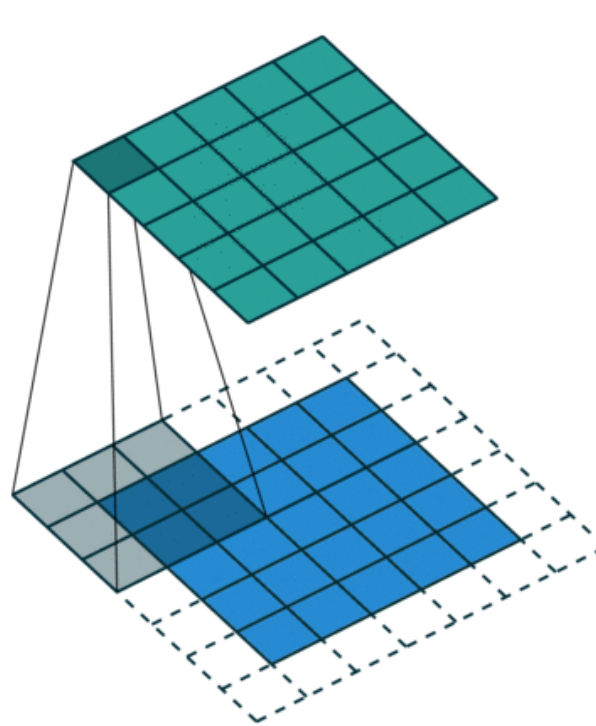


Zeropadding

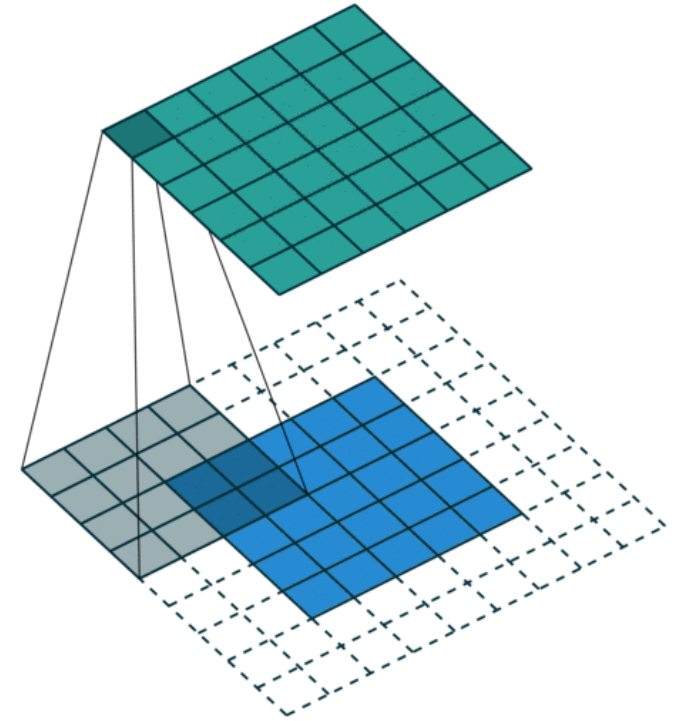
No zero padding, unit strides



Zero padding, unit strides



1 x 1



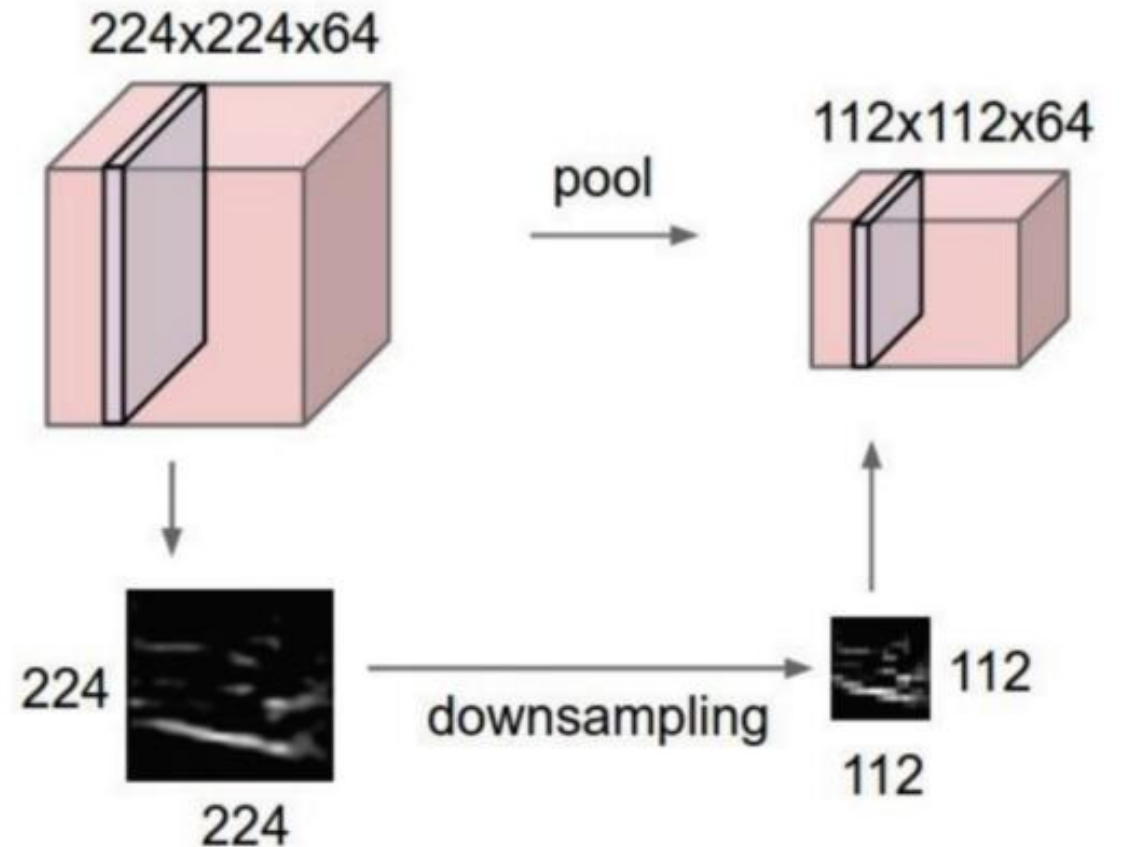
2 x 2

Pooling

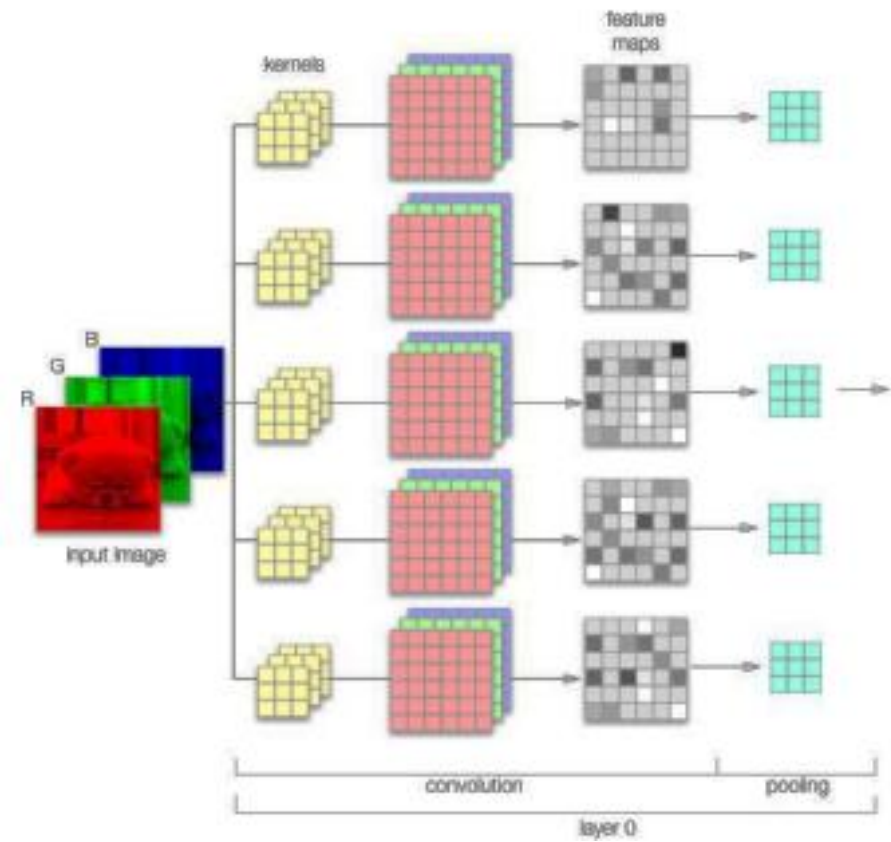
- makes the representations smaller and more manageable
- operates over each activation map independently

- **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.MaxPool2D(  
    pool_size=(2,2),  
    strides=2
```



Pooling



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

Max Pool
Filter - (2 x 2)
Stride - (2, 2)

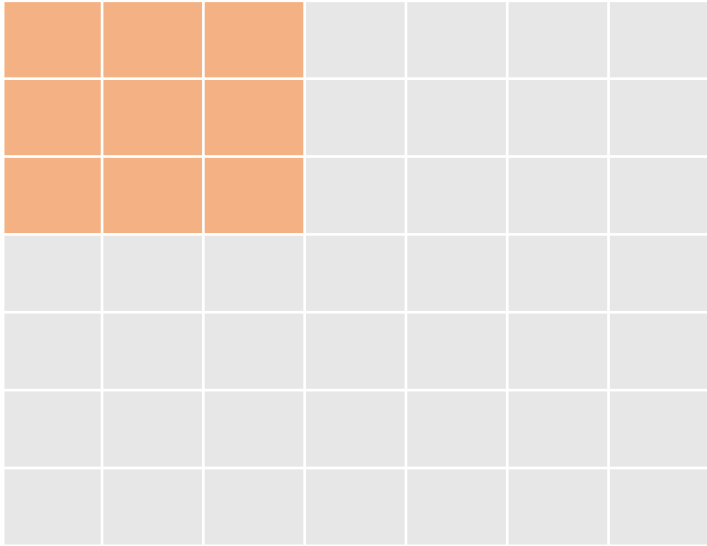
9	7
8	6

Average Pool
Filter - (2 x 2)
Stride - (2, 2)

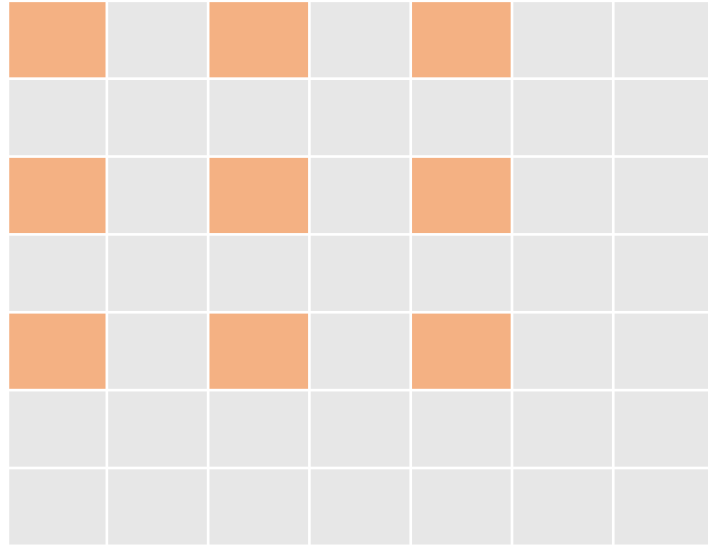
4.25	4.25
4.25	3.5

Reduce dimensionality

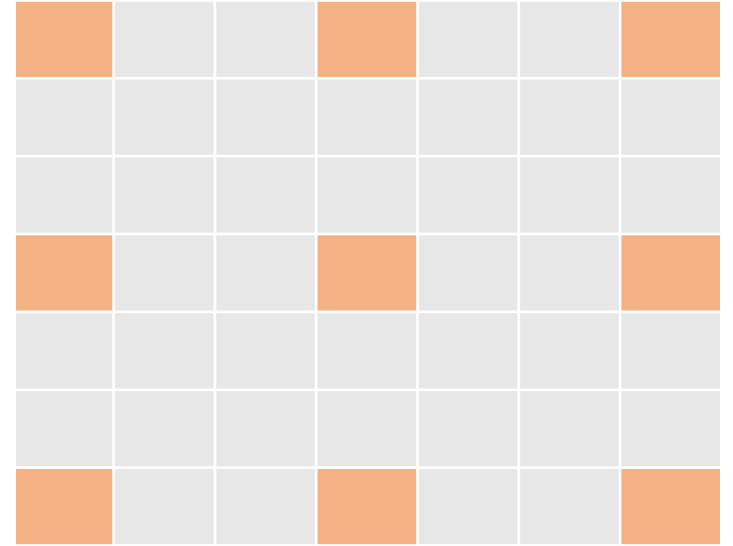
Dilation



Kernel: 3x3
Dilation: 1x1

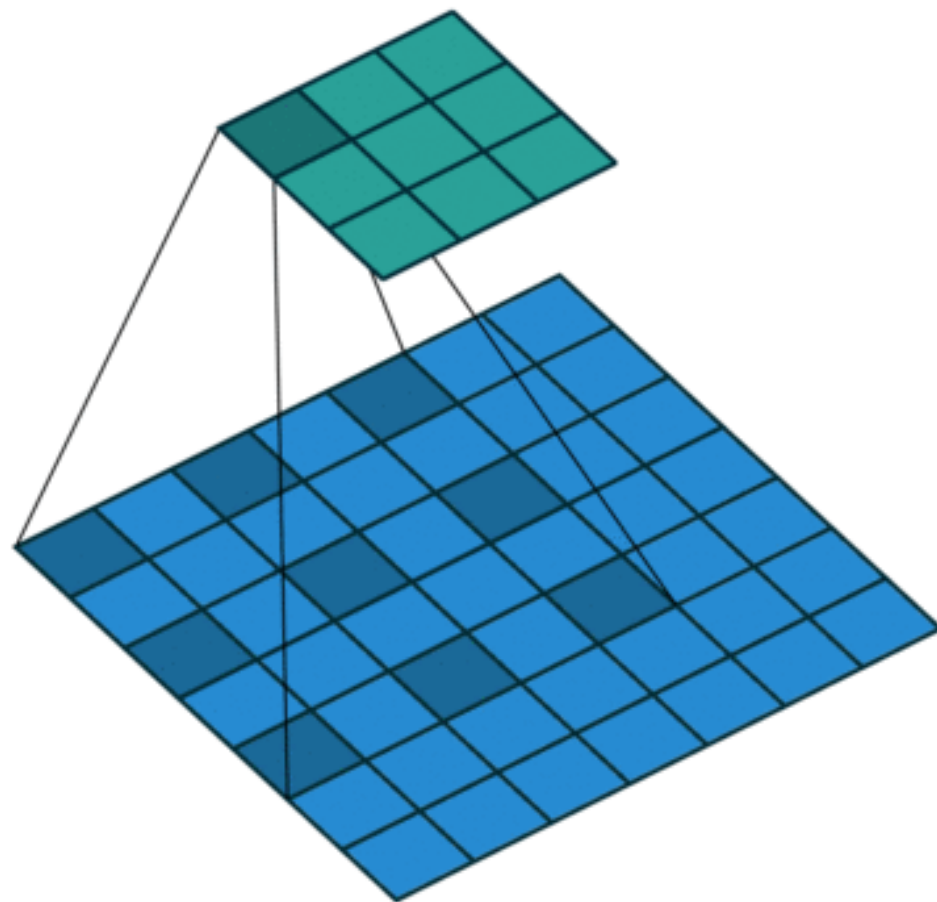


Kernel: 3x3
Dilation: 2x2



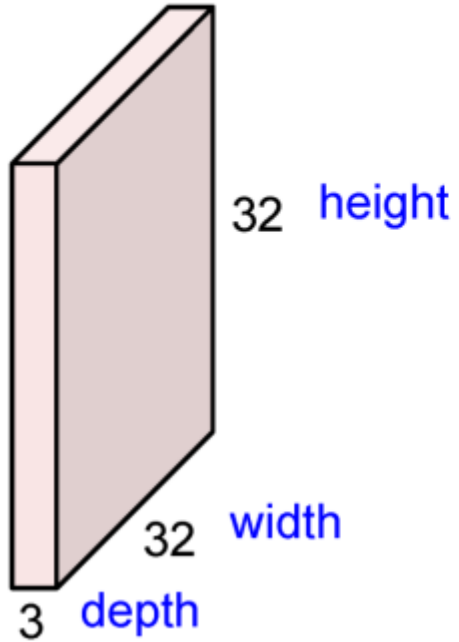
Kernel: 3x3
Dilation: 3x3

Dilation

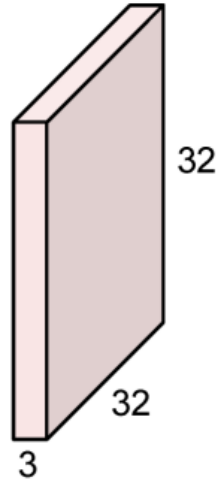


Convolution Layer

32x32x3 image -> preserve spatial structure

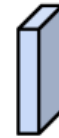


32x32x3 image



Filters always extend the full depth of the input volume

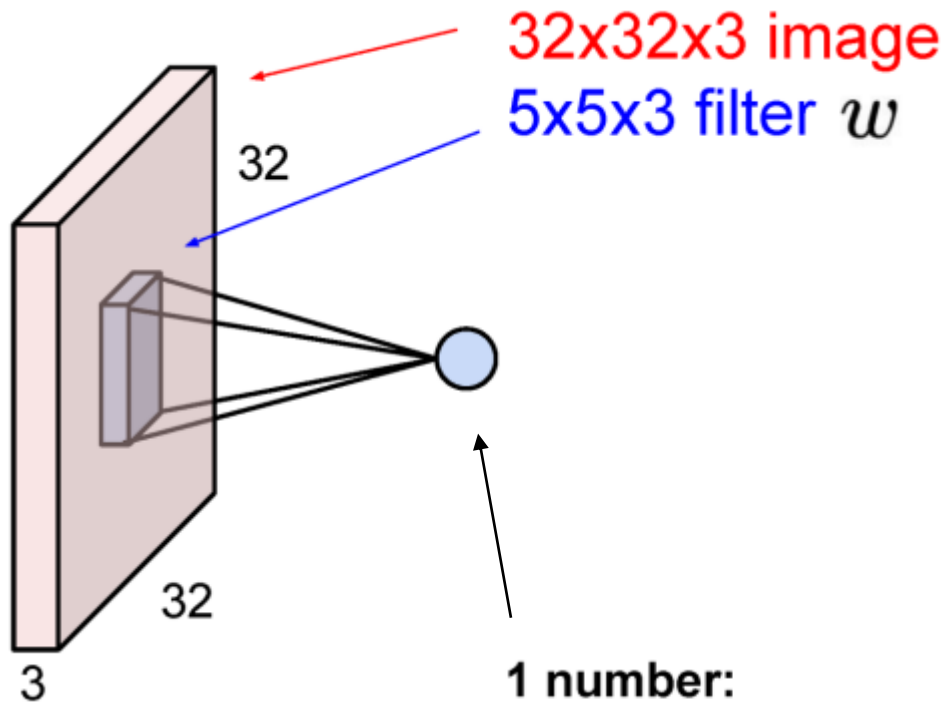
5x5x3 filter



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

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

Convolution Layer



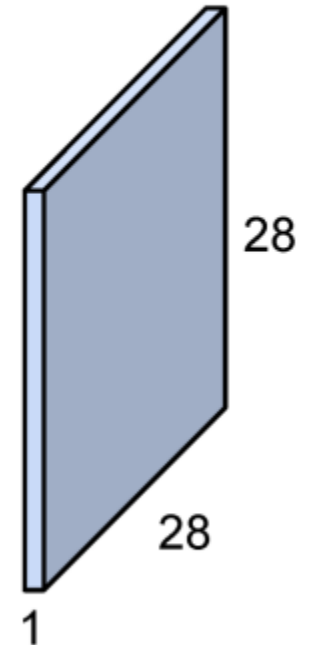
1 number:

the result of taking a dot product between the filter and a small 5x5x3 chunk of the image (i.e. $5*5*3 = 75$ -dimensional dot product + bias)

$$w^T x + b$$

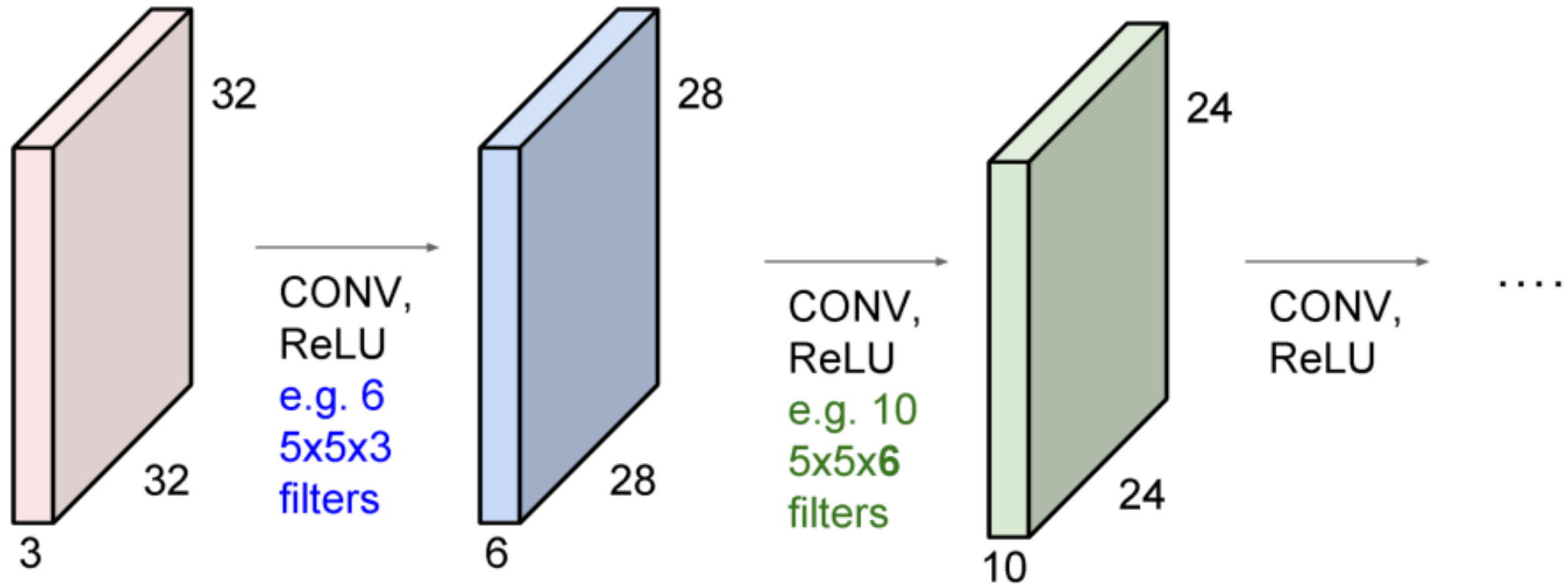
convolve (slide) over all spatial locations

activation map



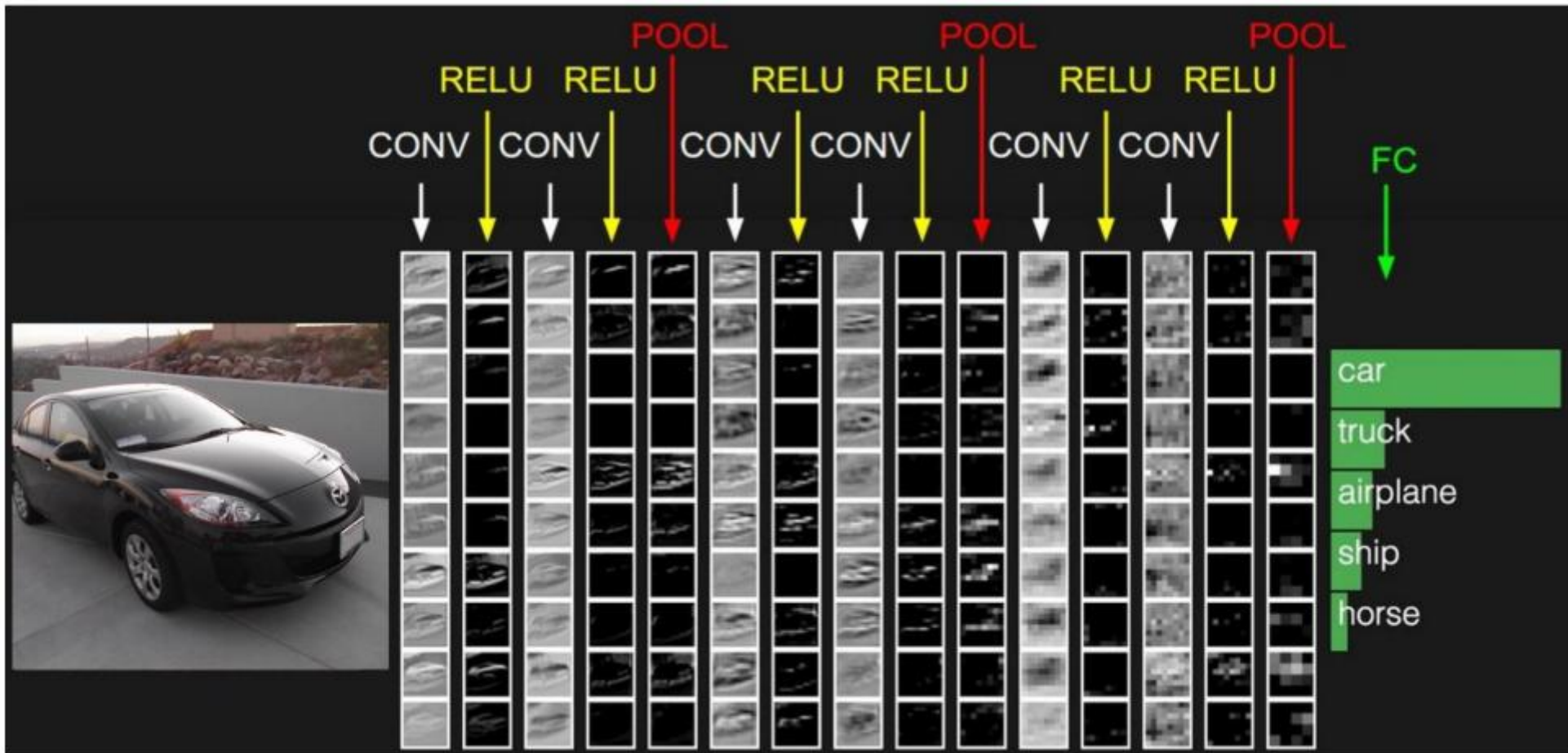
Convolution Layer

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

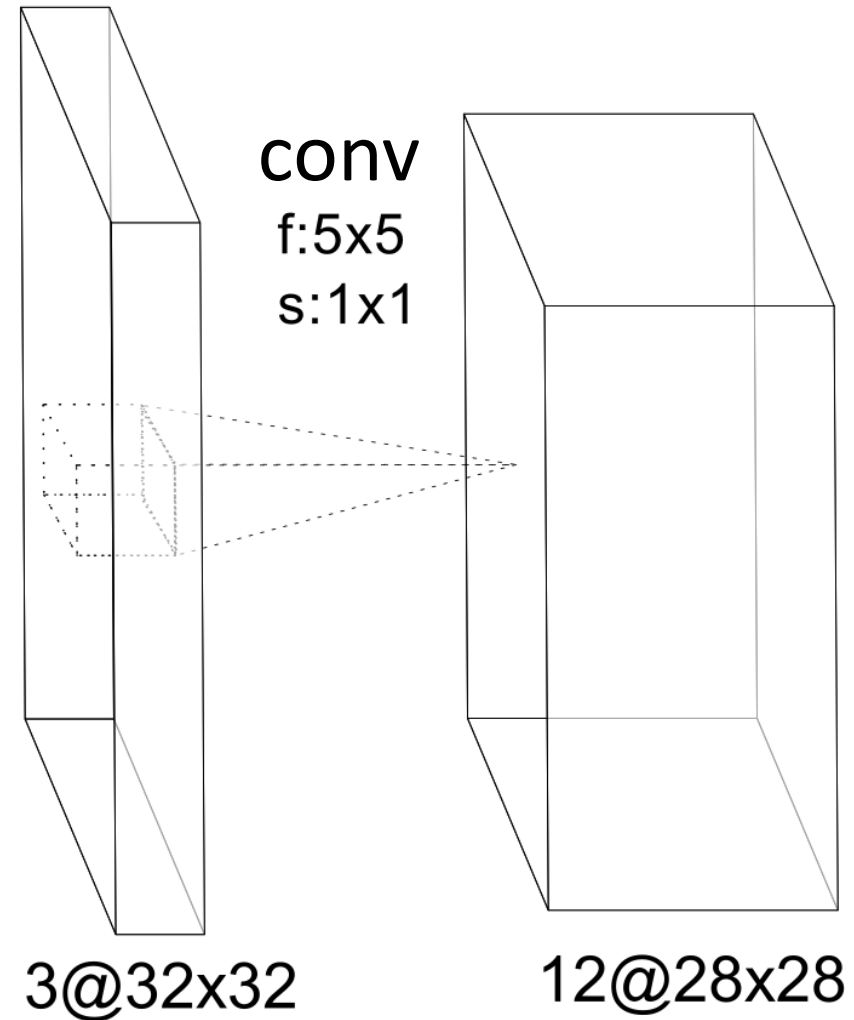
Output: 12@28x28

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:

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

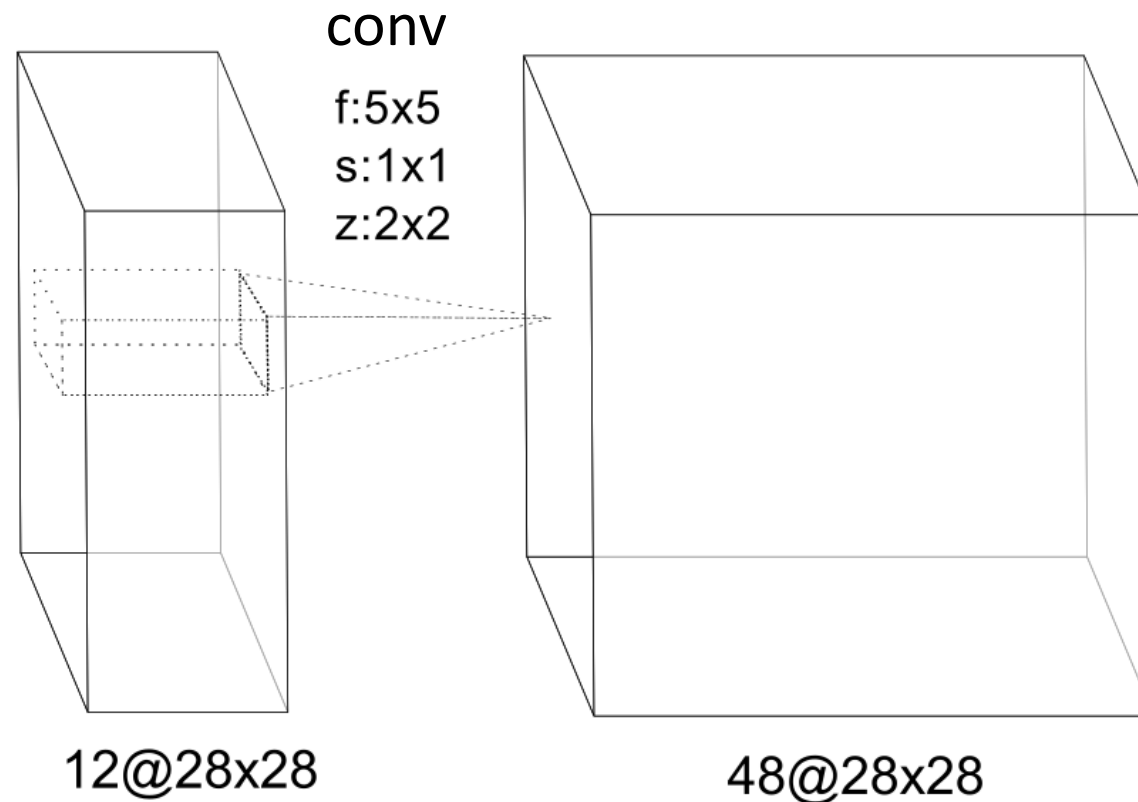


Example: 12@28x28, f:5x5 @48, s:1x1, z:2x2

Output: 48@28x28

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

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

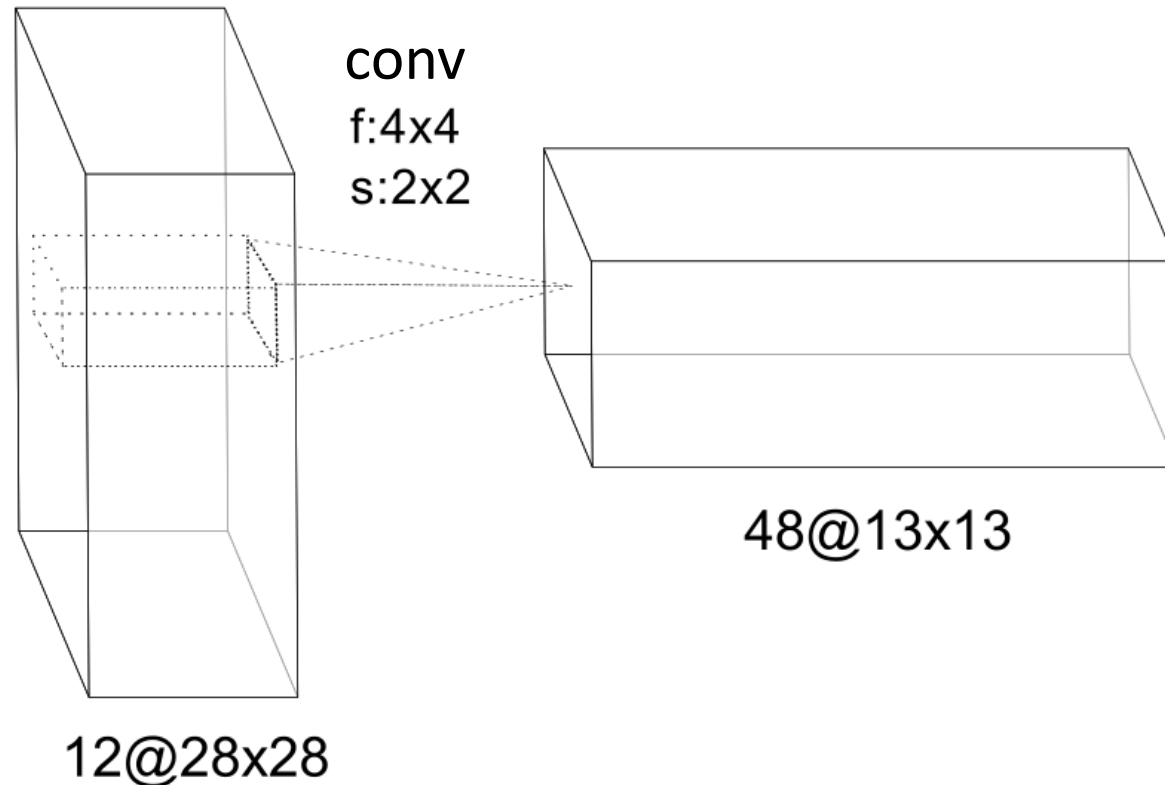


Example: 12@28x28, f:4x4@48, s:2x2, z:0x0

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.

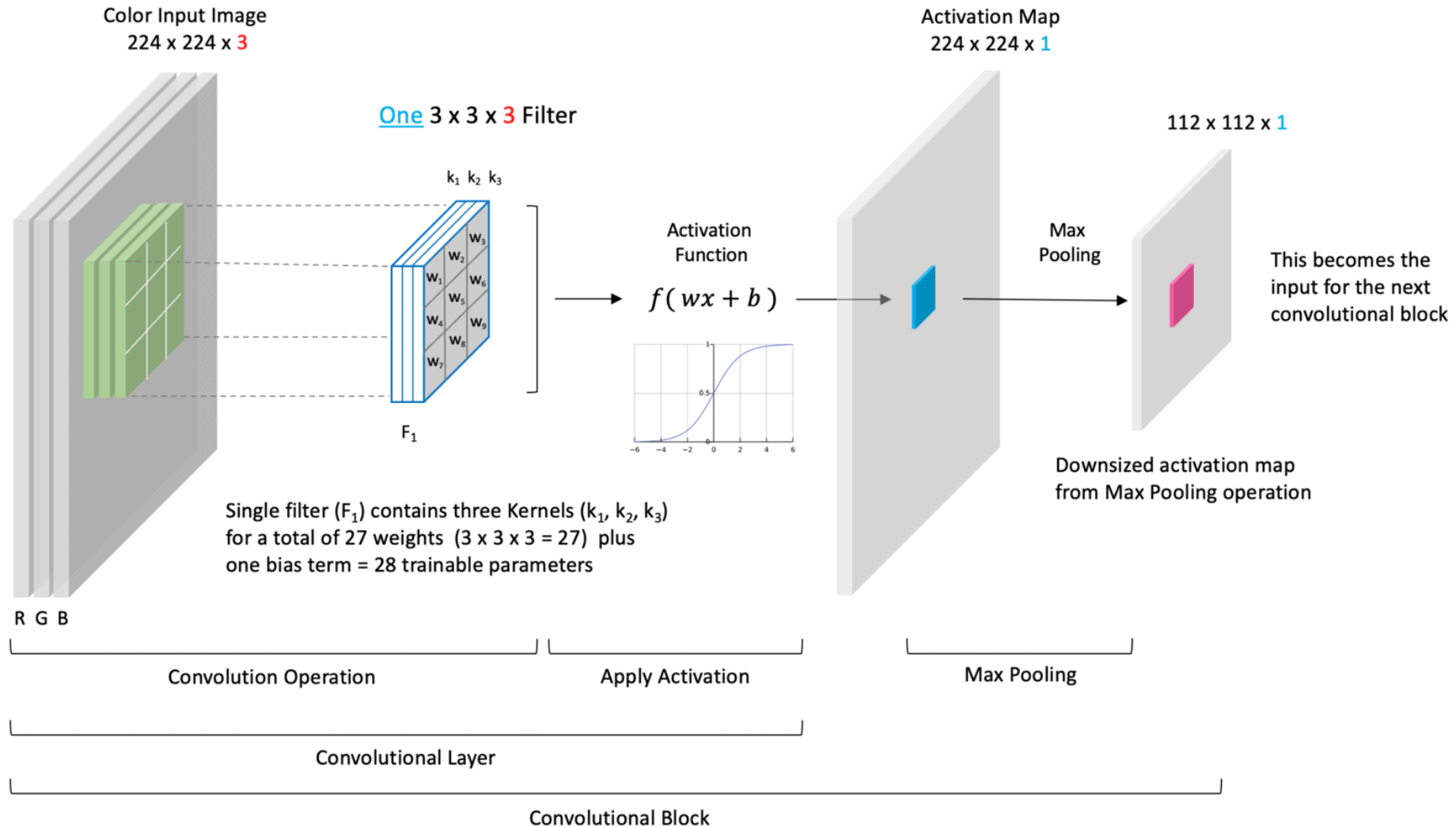
By default, in this demo we're using Adadelta which is one of per-parameter adaptive step size methods, so we don't have to worry about changing learning rates or momentum over time. However, I still included the text fields for changing these if you'd like to play around with SGD+Momentum trainer.

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

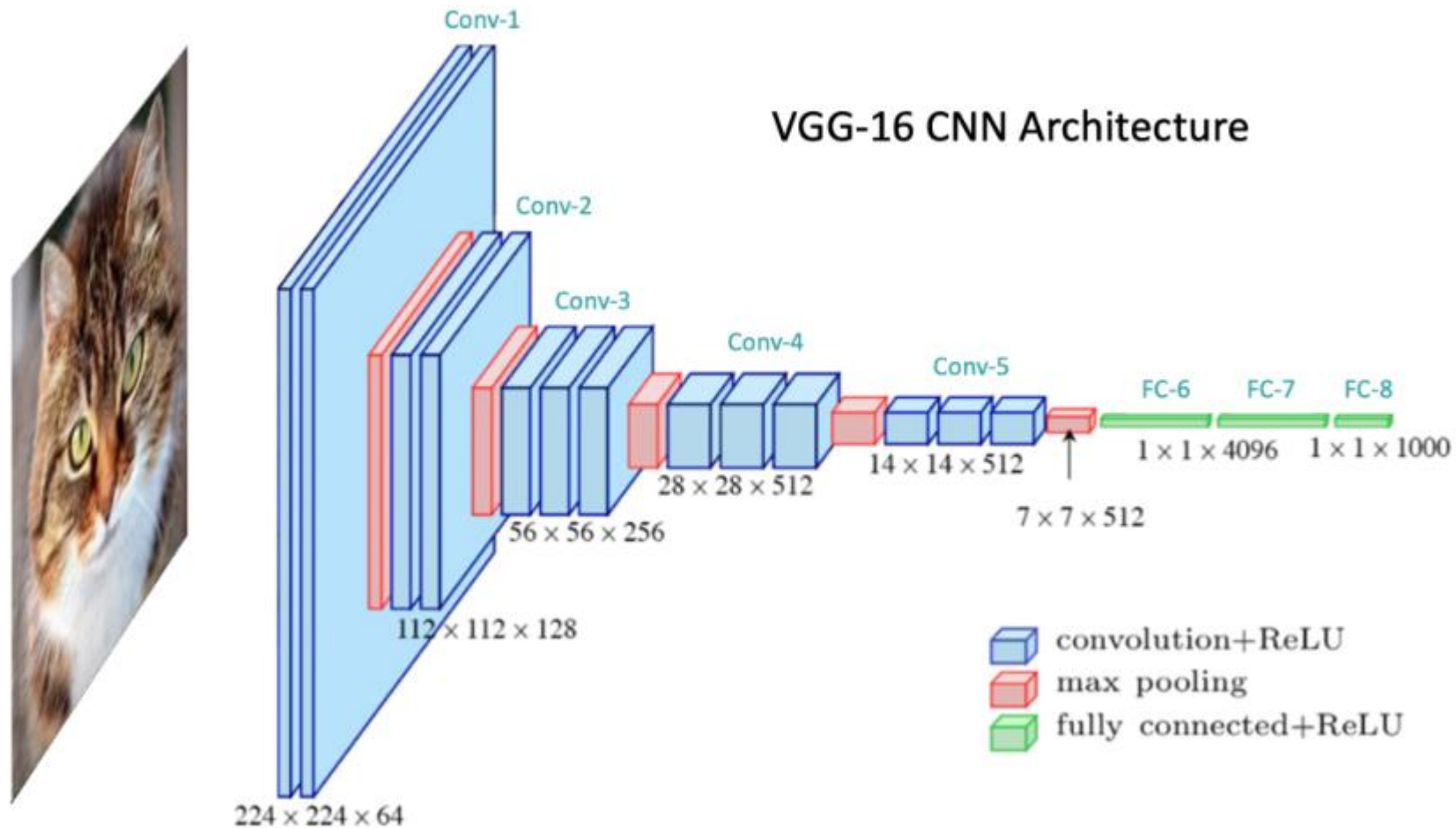


<https://cs.stanford.edu/people/karpathy/convnetjs/demo/cifar10.html>

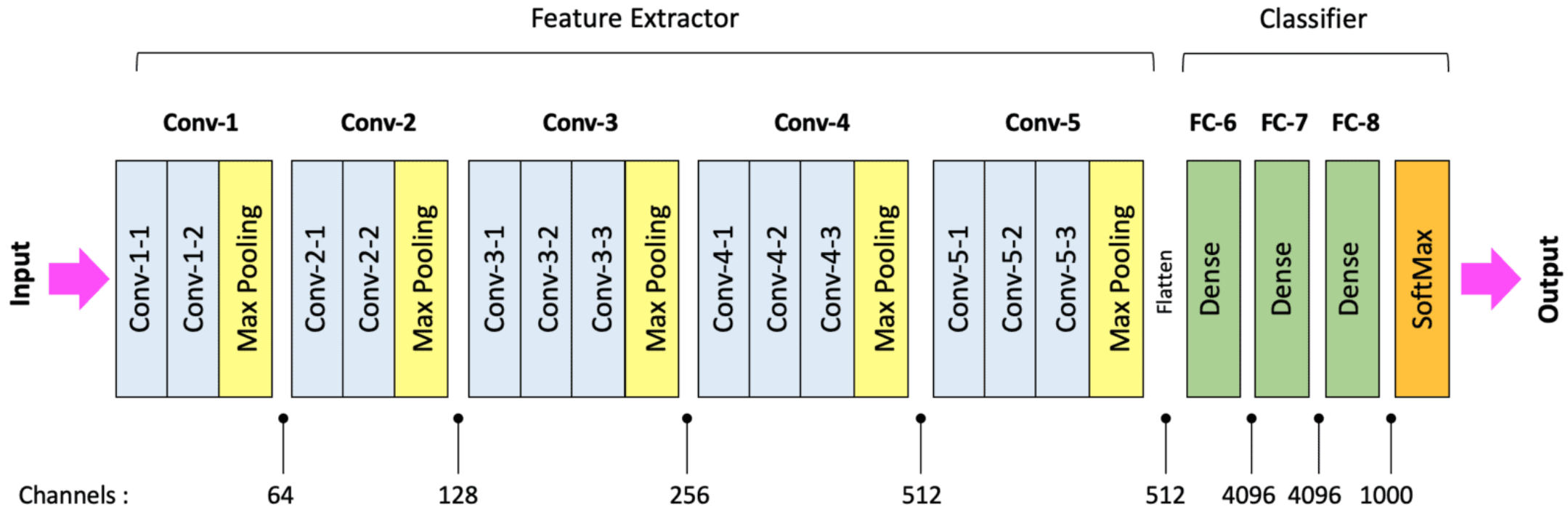
General Convolution Net Architecture



General Convolution Net Architecture



General Convolution Net Architecture



- The fully connected (dense) layers in a CNN architecture transform features into class probabilities.

The 1D convolution Operator

Discrete Form

The diagram illustrates the discrete form of the 1D convolution operator. It features the equation $s_i = (\mathbf{k} * \mathbf{x})_i = \sum_{m=0}^{M-1} k_m x_{i-m+(M-1)}$ centered on the page. Five labels with arrows point to specific parts of the equation: 'FEATURE MAP' points to s_i ; 'KERNEL (or FILTER)' points to \mathbf{k} ; 'KERNEL SIZE' points to the upper limit $M-1$ of the summation; '(TIME) INDEX' points to the subscript i ; and 'INPUT' points to the input variable x .

Labels and arrows:

- FEATURE MAP points to s_i
- KERNEL (or FILTER) points to \mathbf{k}
- KERNEL SIZE points to $M-1$
- (TIME) INDEX points to i
- INPUT points to x

$$s_i = (\mathbf{k} * \mathbf{x})_i = \sum_{m=0}^{M-1} k_m x_{i-m+(M-1)}$$

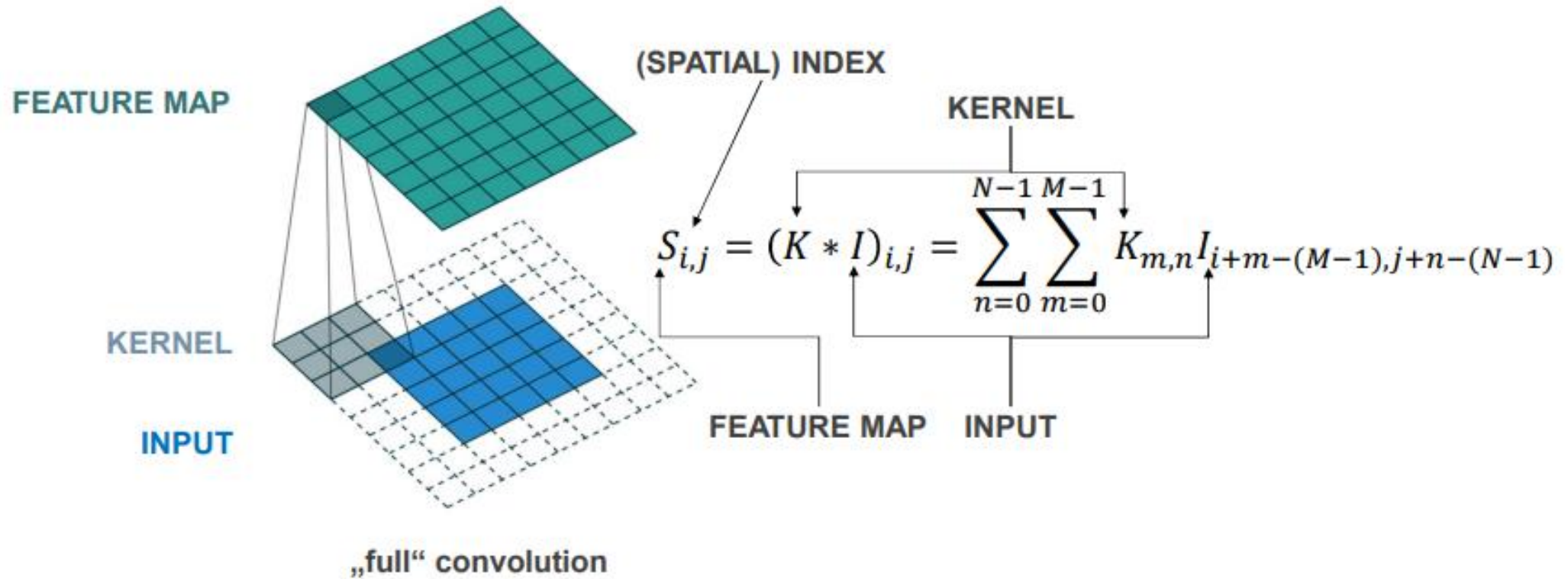
1D convolution implementation

Keras

```
model.add(  
    Convolution1D(filter_depth, filter, border_mode='same',  
    input_shape=(depth, x))  
)
```

Input/output data: $(batch_size, steps, input_dim)$

The 2D convolution Operator



Size After Convolution

Feature map size:

$$O = \frac{W - K + 2P}{S} + 1$$

output width

input width

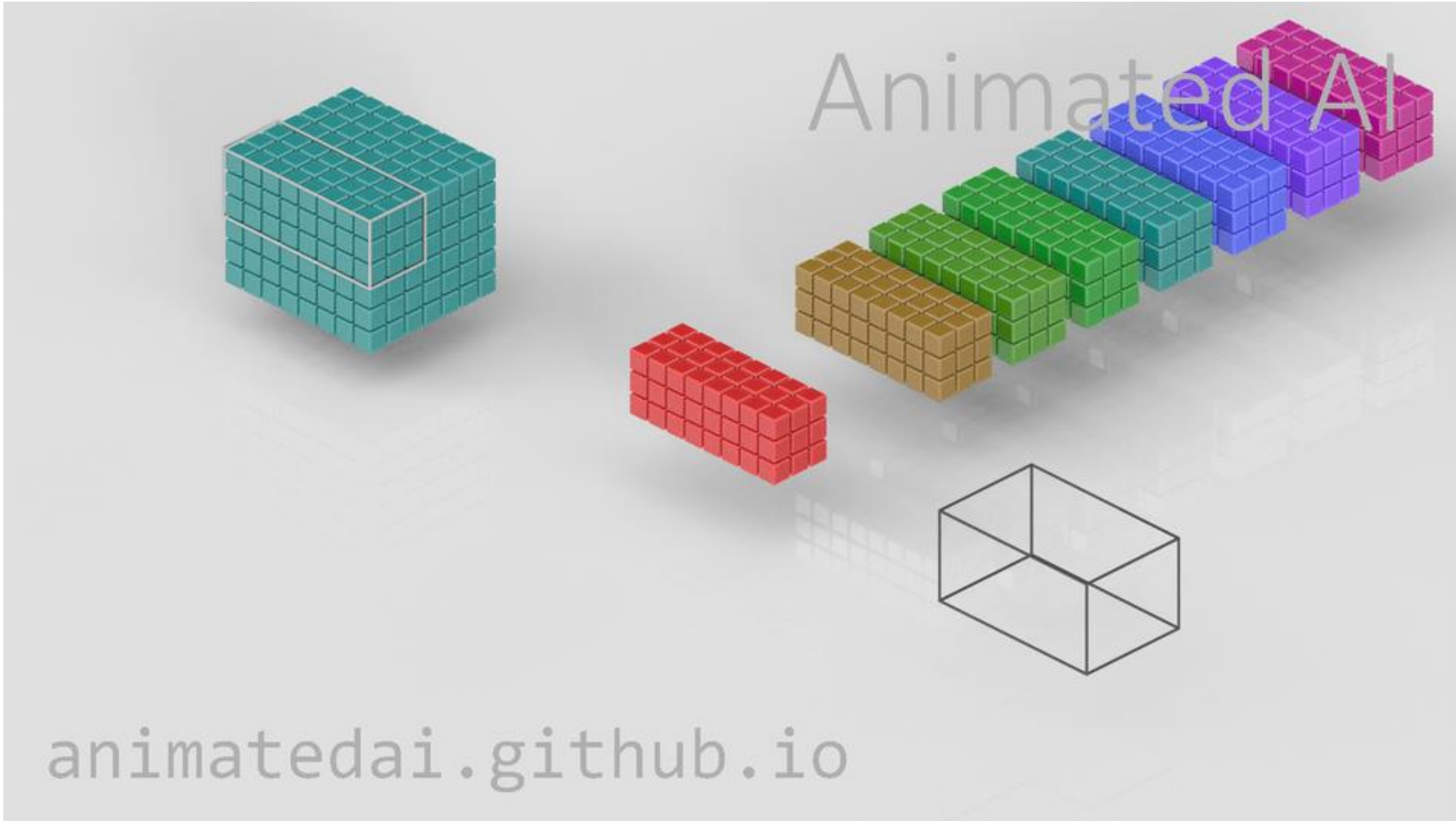
kernel width

padding

stride

The diagram shows the formula for the output width of a convolution. The variables are labeled with arrows: 'input width' points to W, 'kernel width' points to K, 'padding' points to P, and 'stride' points to S. The output width O is indicated by an arrow from the label 'output width'.

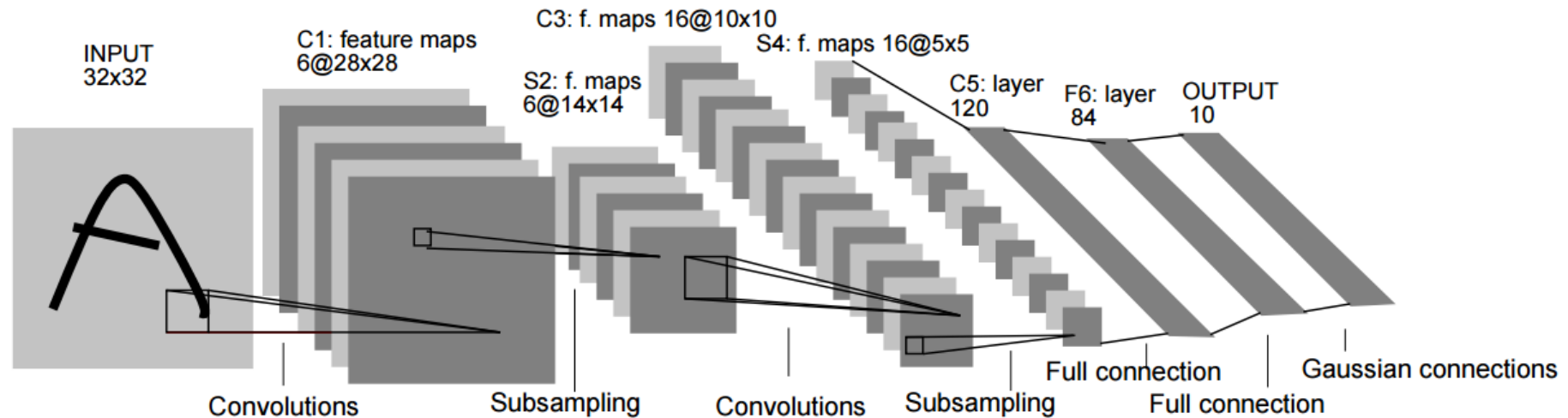
Convolution



<https://animatedai.github.io/>

LeNet5

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.



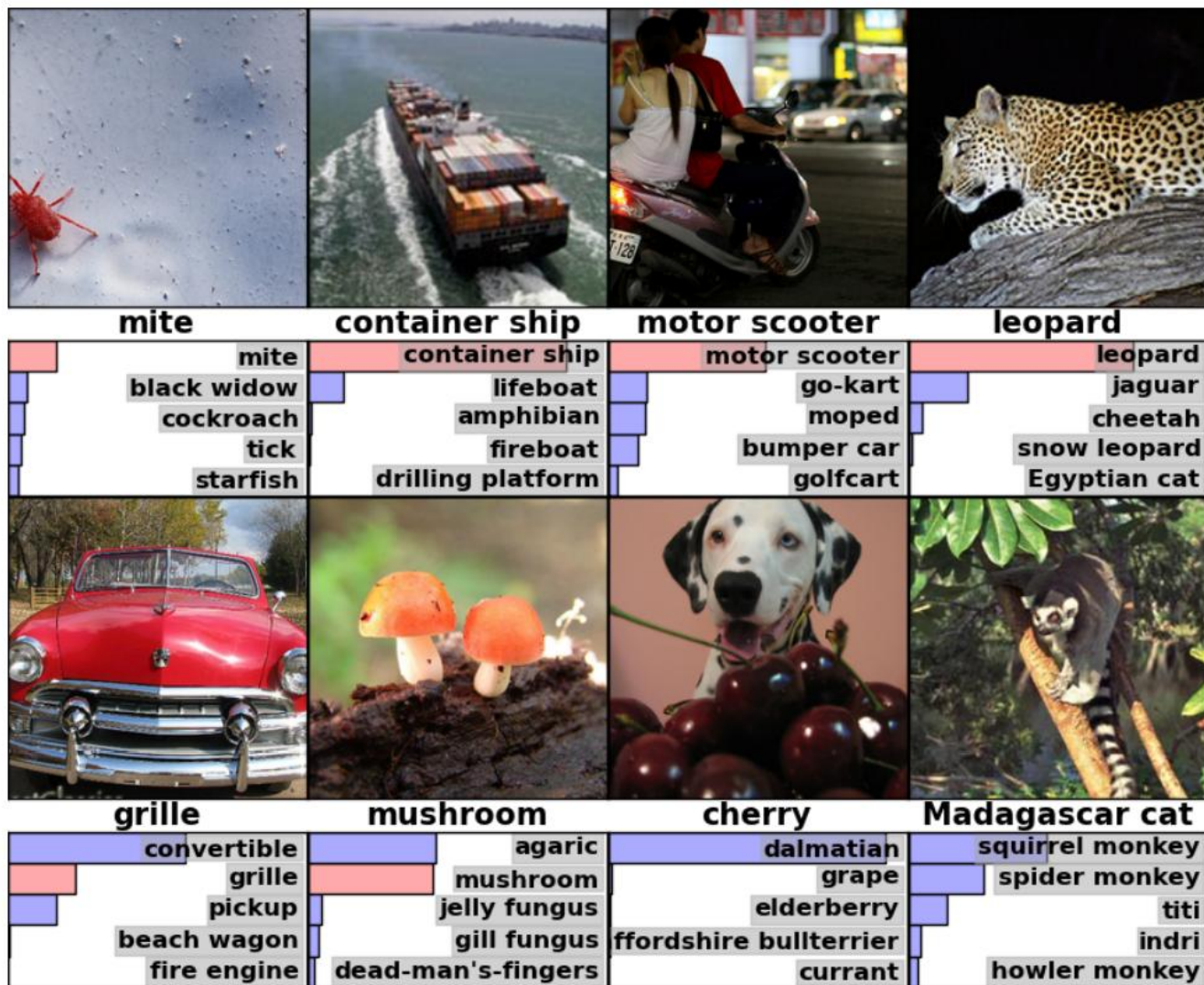
Visualisation: <http://scs.ryerson.ca/~aharley/vis/>

LeNet5 – Keras

```
model = Sequential()
model.add(Conv2D(32, kernel_size=(3, 3),
activation='relu', input_shape=input_shape))
model.add(Conv2D(64, (3, 3), activation='relu'))
model.add(MaxPooling2D(pool_size=(2, 2)))
model.add(Dropout(0.25))
model.add(Flatten())
model.add(Dense(128, activation='relu'))
model.add(Dropout(0.5))
model.add(Dense(num_classes, activation='softmax'))

model.compile(loss=keras.losses.categorical_crossentropy,
optimizer=keras.optimizers.Adadelta(),
metrics=['accuracy'])
```

Image – AlexNet 2012



AlexNet – Keras

```
model = Sequential()
model.add(Conv2D(filters=96, input_shape=(224, 224, 3),
kernel_size=(11, 11), strides=(4, 4), padding='valid'))
model.add(Activation('relu'))
model.add(MaxPooling2D(pool_size=(2, 2), strides=(2, 2),
padding='valid'))
model.add(Conv2D(filters=256, kernel_size=(11, 11),
strides=(1, 1), padding='valid'))
model.add(Activation('relu'))
model.add(MaxPooling2D(pool_size=(2, 2), strides=(2, 2),
padding='valid'))
```

AlexNet – Keras

```
model.add(Conv2D(filters=384, kernel_size=(3,3),  
strides=(1,1), padding='valid'))
```

```
model.add(Activation('relu'))
```

```
model.add(Conv2D(filters=384, kernel_size=(3,3),  
strides=(1,1), padding='valid'))
```

```
model.add(Activation('relu'))
```

```
model.add(Conv2D(filters=256, kernel_size=(3,3),  
strides=(1,1), padding='valid'))
```

```
model.add(Activation('relu'))
```

AlexNet – Keras

```
model.add(MaxPooling2D(pool_size=(2,2), strides=(2,2),  
padding='valid'))
```

```
model.add(Flatten())
```

```
model.add(Dense(4096))
```

```
model.add(Activation('relu'))
```

```
model.add(Dropout(0.4))
```

```
model.add(Dense(4096))
```

```
model.add(Activation('relu'))
```

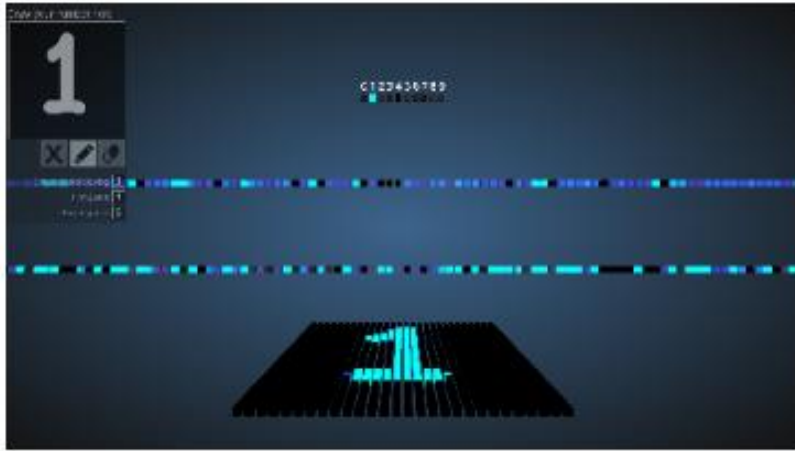
```
model.add(Dropout(0.4))
```

```
model.add(Dense(1000))
```

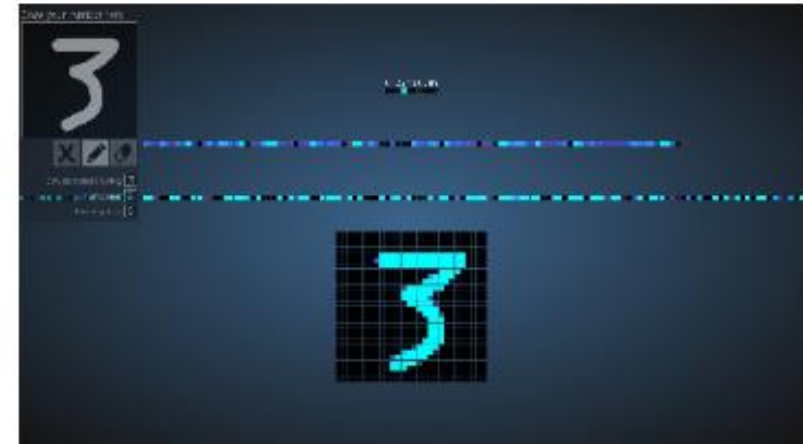
```
model.add(Activation('softmax'))
```

Node-Link Visualization of CNNs

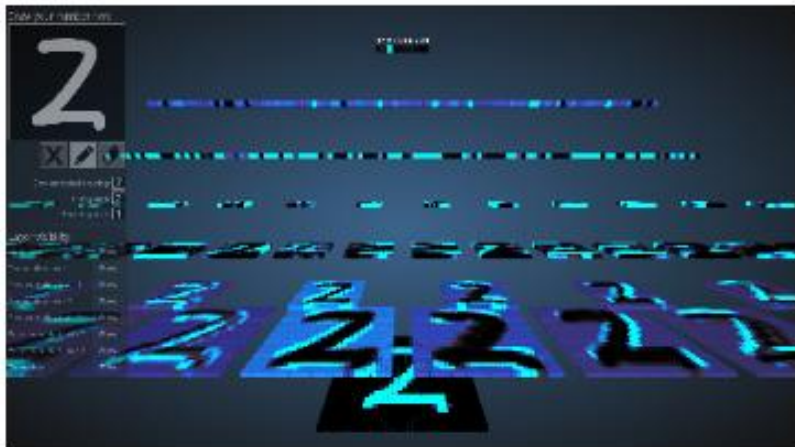
1. 3d visualization of a multi-layer perceptron:



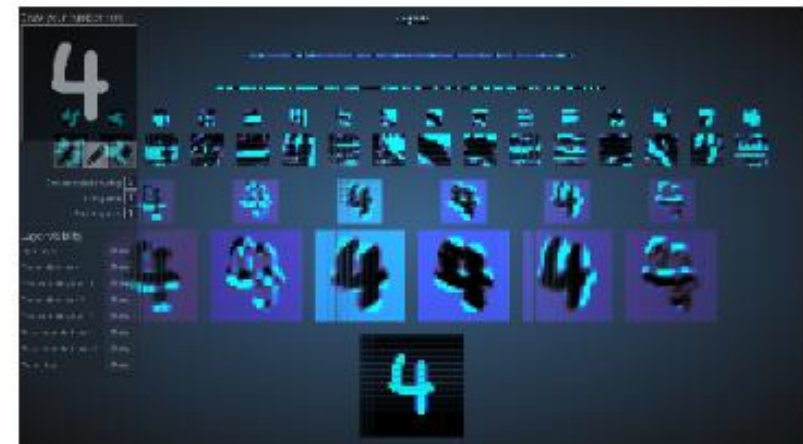
3. 2d visualization of a multi-layer perceptron:



2. 3d visualization of a convolutional network:

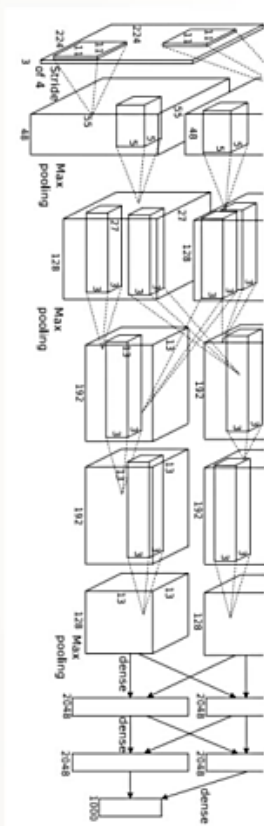


4. 2d visualization of a convolutional network:



Popular nets

“AlexNet”



[Krizhevsky et al. NIPS 2012]

“GoogLeNet”



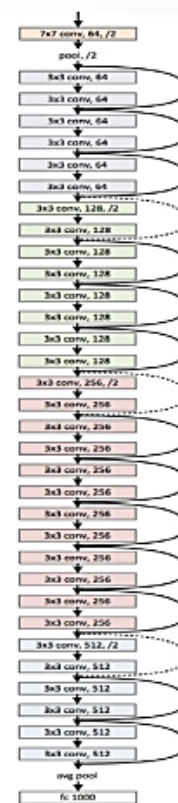
[Szegedy et al. CVPR 2015]

“VGG Net”



[Simonyan & Zisserman, ICLR 2015]

“ResNet”



[He et al. CVPR 2016]

Popular image databases

- MNIST
- NIST
- CIFAR10/CIFAR100
- SLT-10
- SVHN
- IMAGENET (ILSVRC2012)
- PascalVOC
- Google Open Image Dataset
- Stb.

Reformatting data 1

Mini-batch compilation for a parameter flow:

batch (b) x steps (m) x 1

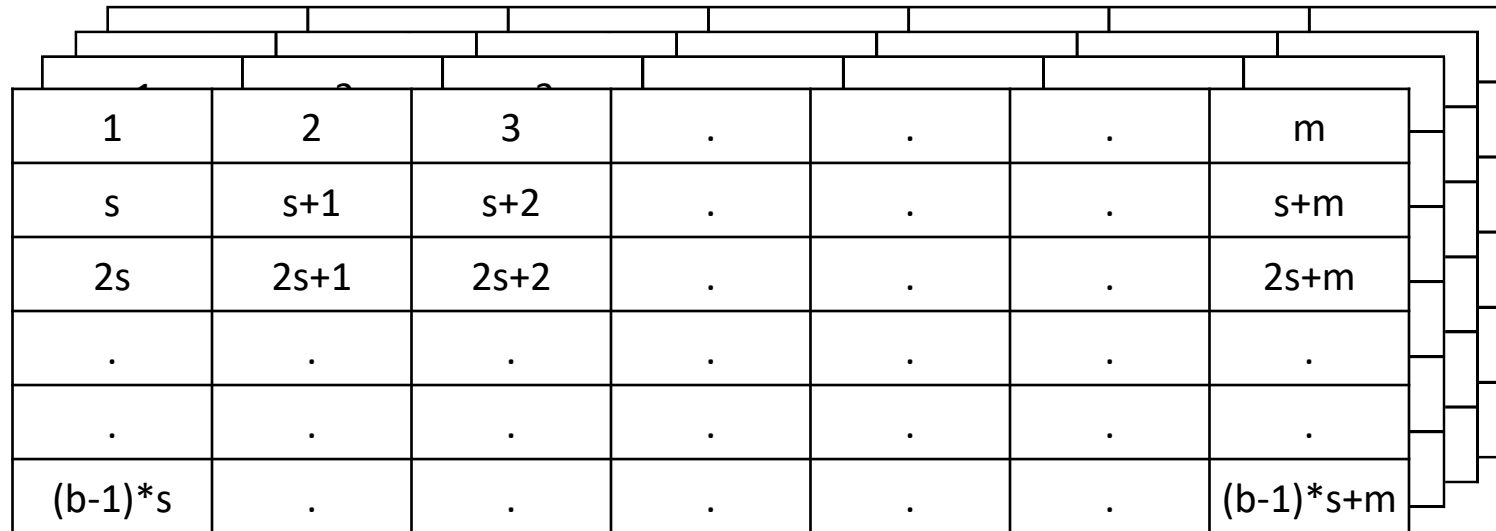
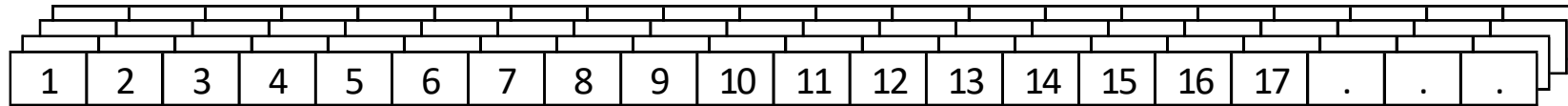
s: step interval between two samples

1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	.	.	.
---	---	---	---	---	---	---	---	---	----	----	----	----	----	----	----	----	---	---	---

1	2	3	.	.	.	m
s	s+1	s+2	.	.	.	s+m
2s	2s+1	2s+2	.	.	.	2s+m
.
.
(b-1)*s	(b-1)*s+m

Reformatting data 2

Mini-batch compilation for multiple parameter streams:
batch (b) x steps (m) x input_dim



Example of convolution

- Texts
 - Translation
 - Sentiment analysis
 - IMDB
 - <https://github.com/bhaveshoswal/CNN-text-classification-keras>
 - Plagiarism detection
 - Natural language processing
 - Text synthesis
 - Speech synthesis
- Time series
- Sensor data
- etc.

Tricks, methods

- Data augmentation
- Standardization
- Object detection, segmentation
- Batch Normalization (faster convergence, no need for DropOut)
- Dilated convolution (PixelCNN, WaveNet)
- Transfer learning

Tricks, methods

Object Detection



Redmon, J., Divvala, S., Girshick, R., & Farhadi, A. (2016). You only look once: Unified, real-time object detection. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition* (pp. 779-788).

Tricks, methods

Object Segmentation

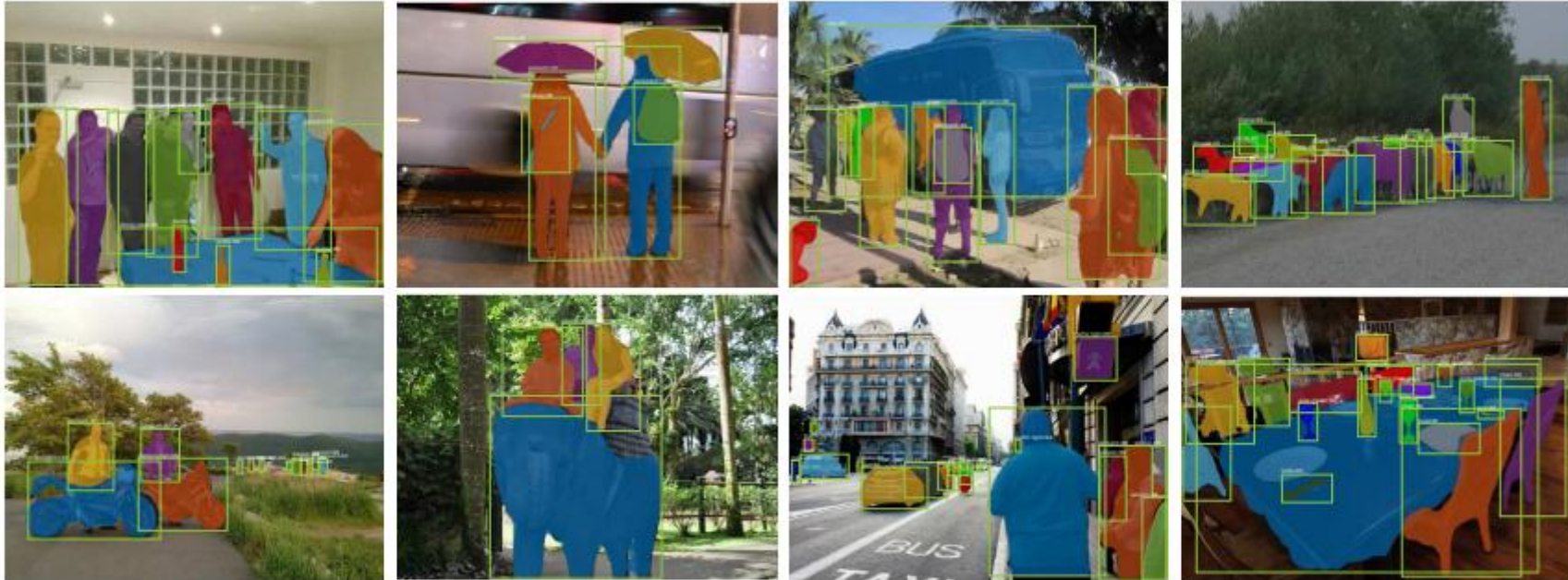
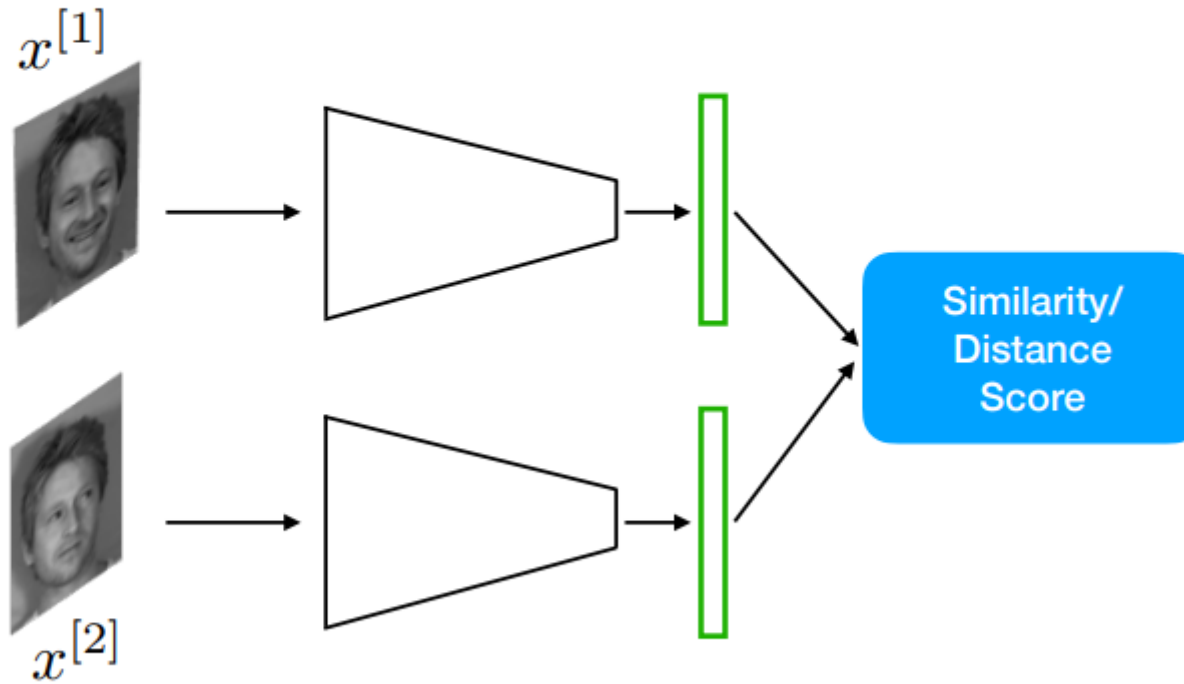


Figure 2. **Mask R-CNN** results on the COCO test set. These results are based on ResNet-101 [15], achieving a *mask AP* of 35.7 and running at 5 fps. Masks are shown in color, and bounding box, category, and confidences are also shown.

He, Kaiming, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. "Mask R-CNN." In *Proceedings of the IEEE International Conference on Computer Vision*, pp. 2961-2969. 2017.

Tricks, methods

Face Recognition



Why image Classification is hard

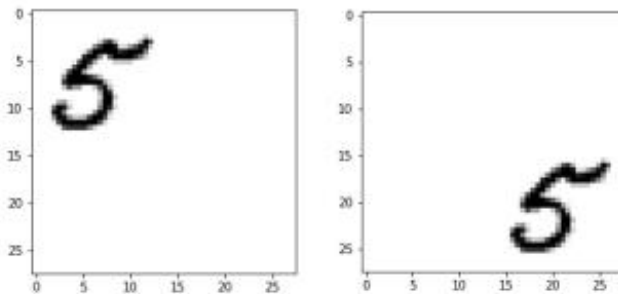
Different lighting, contrast, viewpoints, etc.



Image Source:
twitter.com/%2Fcats&psig=AOvWaw30_o-PCM-K21DIMAjQimQ4&ust=1553887775741551

Image Source: https://www.123rf.com/photo_76714328_side-view-of-tabby-cat-face-over-white.html

Or even simple translation

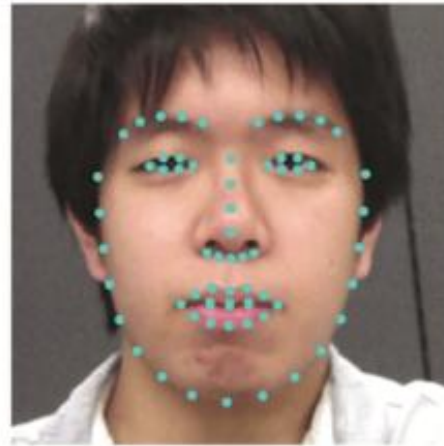


This is hard for traditional methods like multi-layer perceptrons, because the prediction is basically based on a sum of pixel intensities

Why image Classification is hard

➤ Traditional Approaches

a) Use hand-engineered features



(a) Detected facial keypoints



(b) Facial organ keypoints

Why image Classification is hard

➤ Traditional Approaches

b) Preprocess images (centering, cropping, etc.)



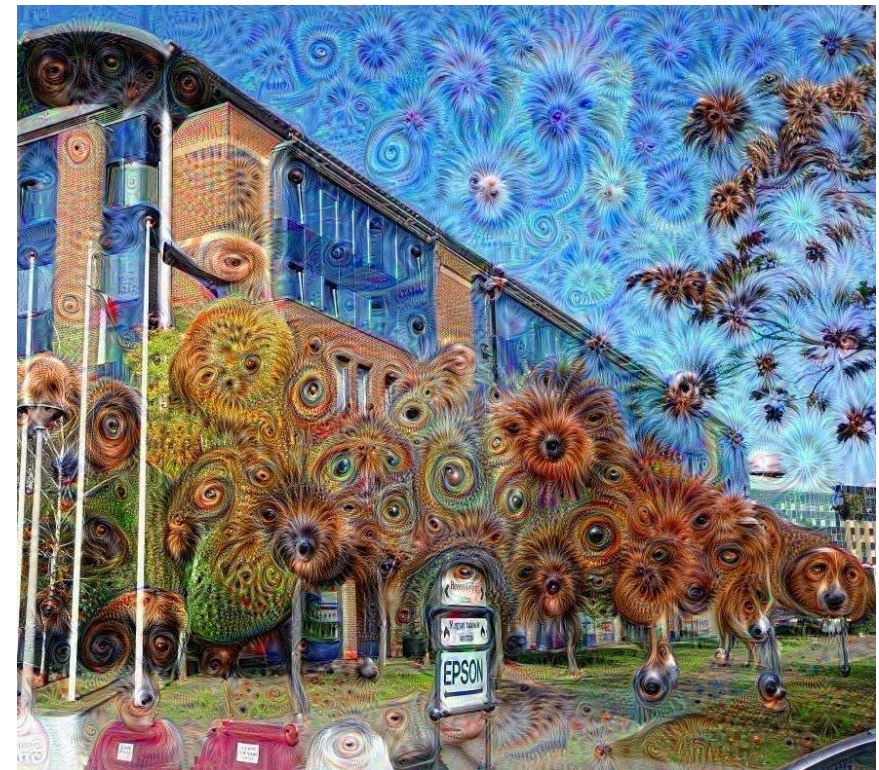
Related links

Further description, information

<http://cs231n.github.io/convolutional-networks/>



Deep dream generator

<https://deepdreamgenerator.com/>




Tutorial


TensorFlow > Learn > TensorFlow Core > Tutorials

Was this helpful?  

Convolutional Neural Network (CNN)

 [Run in Google Colab](#)

 [View source on GitHub](#)

 [Download notebook](#)

This tutorial demonstrates training a simple [Convolutional Neural Network \(CNN\)](#) to classify [CIFAR images](#). Because this tutorial uses the [Keras Sequential API](#), creating and training your model will take just a few lines of code.

Import TensorFlow

```
import tensorflow as tf

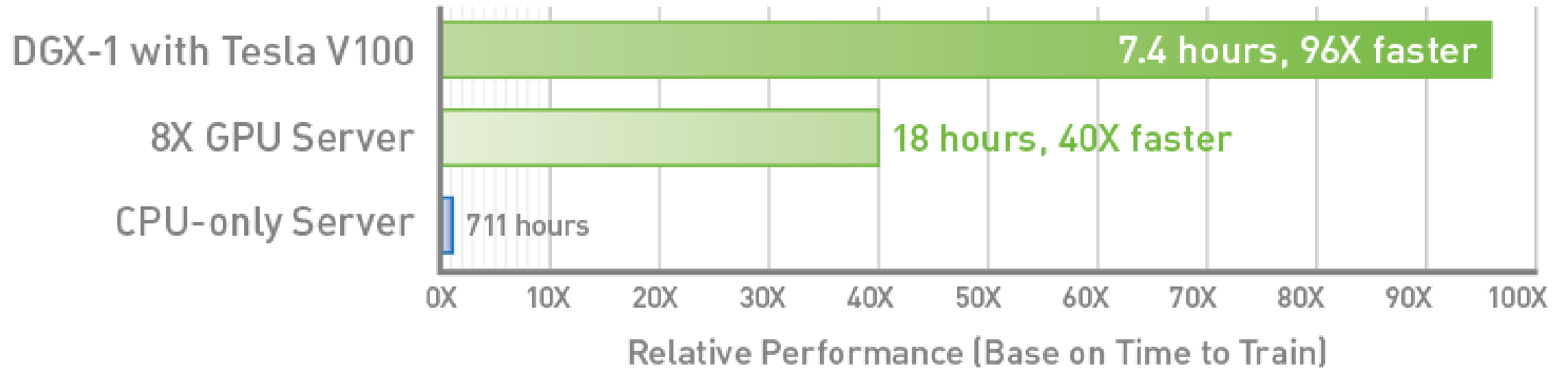
from tensorflow.keras import datasets, layers, models
import matplotlib.pyplot as plt
```

Download and prepare the CIFAR10 dataset

The CIFAR10 dataset contains 60,000 color images in 10 classes, with 6,000 images in each class. The dataset is divided into 50,000 training images and 10,000 testing images. The classes are mutually exclusive and there is no overlap between them.

- <https://github.com/tensorflow/docs/blob/master/site/en/tutorials/images/cnn.ipynb>
- <https://github.com/tensorflow/docs/blob/master/site/en/tutorials/images/classification.ipynb>
- <https://colab.research.google.com/github/tensorflow/docs/blob/master/site/en/tutorials/images/cnn.ipynb>

NVIDIA DGX-1 Delivers 96X Faster Training



Workload: ResNet50, 90 epochs to solution | CPU Server: Dual Xeon E5-2699 v4, 2.6GHz

November 13 result: ResNet50 90 epoch training in 3.7 minutes on 2176 V100 GPUs.

Source: Hiroaki Mikami, Hisahiro Suganuma, Pongsakorn U-chupala, Yoshiki Tanaka, Yuichi Kageyama: ImageNet/ResNet-50 Training in 224 Seconds

DGX SATURNV



summit

IBM

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

(slides by: Dr. Bálint Gyires-Tóth)

15 October 2024

