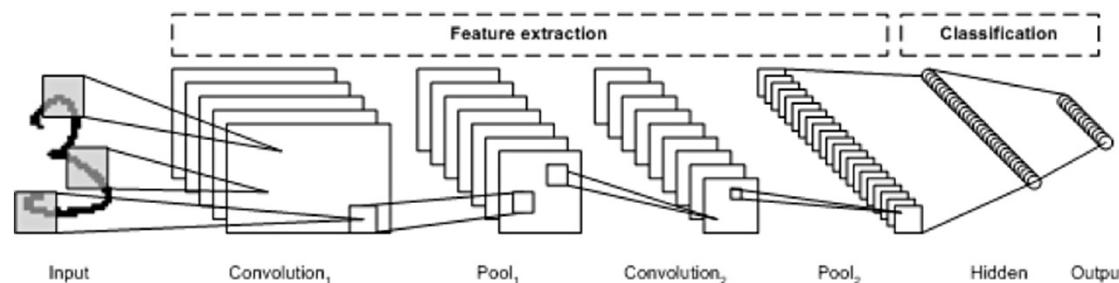




Computational vision:

Convolutional Neural Networks

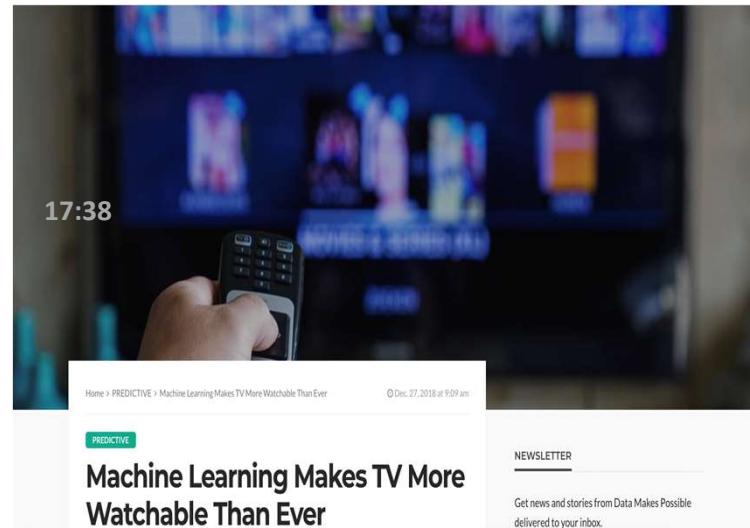
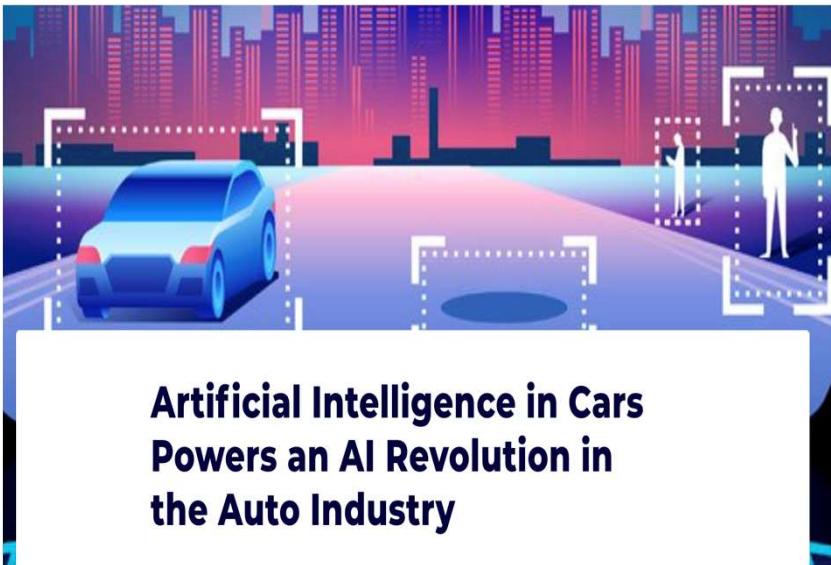
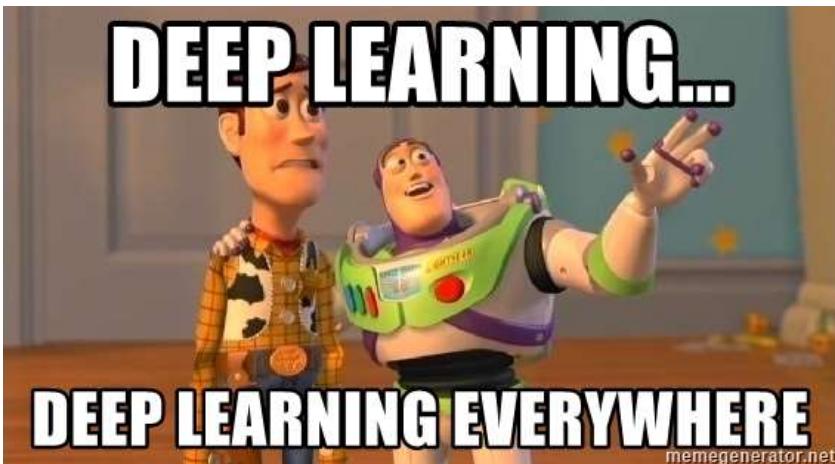
Class 6: Computational Vision





- ↗ AI, Machine learning & Deep learning
- ↗ What is a Convolutional Neural Network?
 - ↗ Layers
 - ↗ Optimization
- ↗ Applications

Deep learning everywhere



Artificial Intelligence

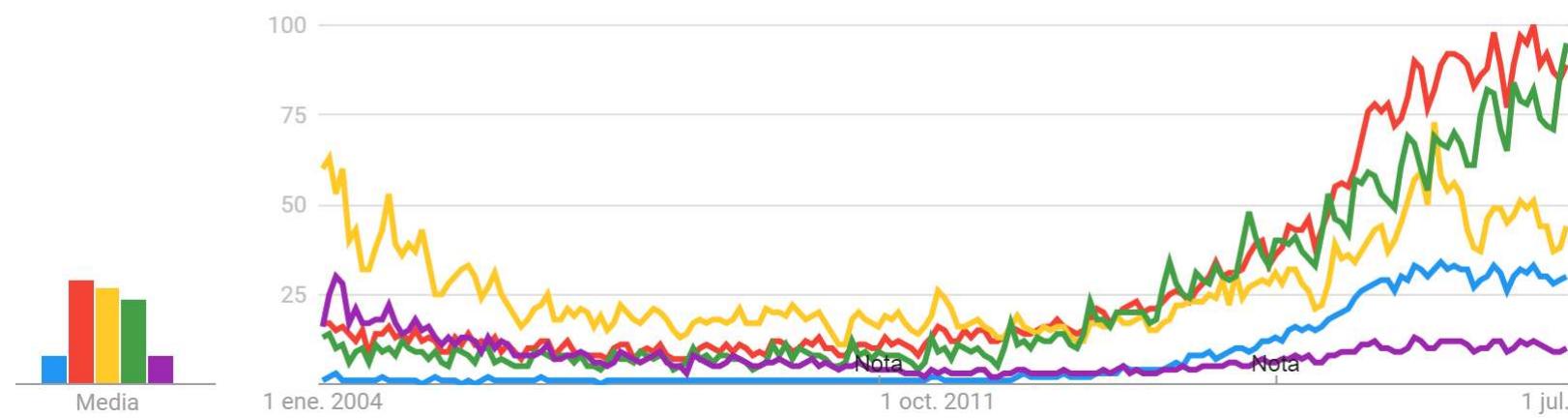
↗ Artificial Intelligence

↗ Machine Learning

↗ Deep learning

↗ Neural networks

↗ Data science



17:38

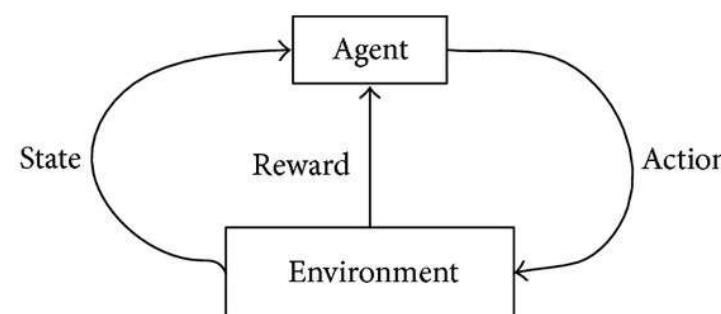
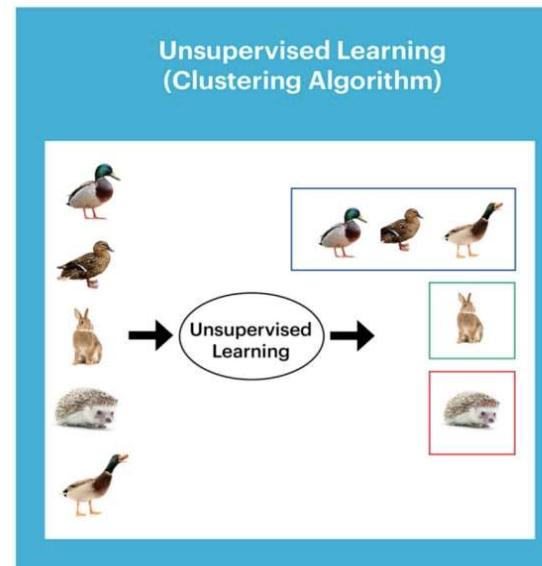
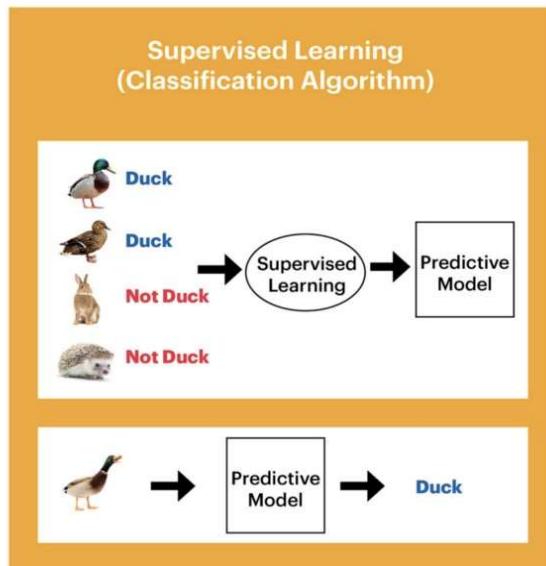
Google Scholar reveals its most influential papers

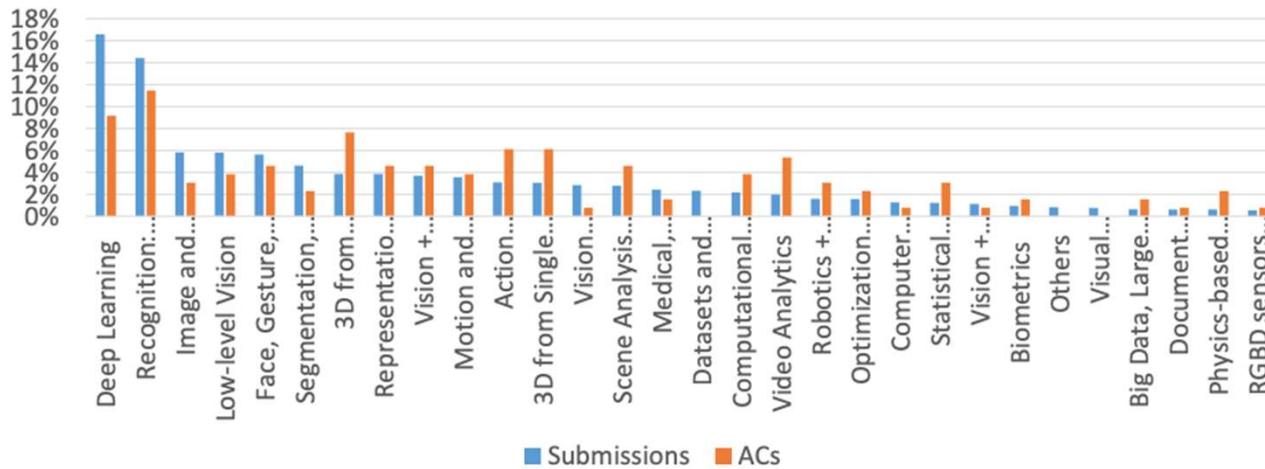


1. **"Deep Residual Learning for Image Recognition" (2016)** *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition* 25,256 citations
2. **"Deep learning" (2015)** *Nature* 16,750 citations
3. **"Going Deeper with Convolutions" (2015)** *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition* 14,424 citations
4. **"Fully Convolutional Networks for Semantic Segmentation" (2015)** *Proceedings of the IEEE Conf. on Computer Vision and Pattern Recognition* 10,153 citations
5. **"Prevalence of Childhood and Adult Obesity in the United States, 2011-2012" (2014)** *JAMA* 8,057 citations
6. **"Global, regional, and national prevalence of overweight and obesity in children and adults during 1980–2013: a systematic analysis for the Global Burden of Disease Study 2013" (2014)** *Lancet* 7,371 citations
7. **"Observation of Gravitational Waves from a Binary Black Hole Merger" (2016)** *Physical Review Letters* 6,009 citations



Supervised vs. Unsupervised vs Reinforcement Learning





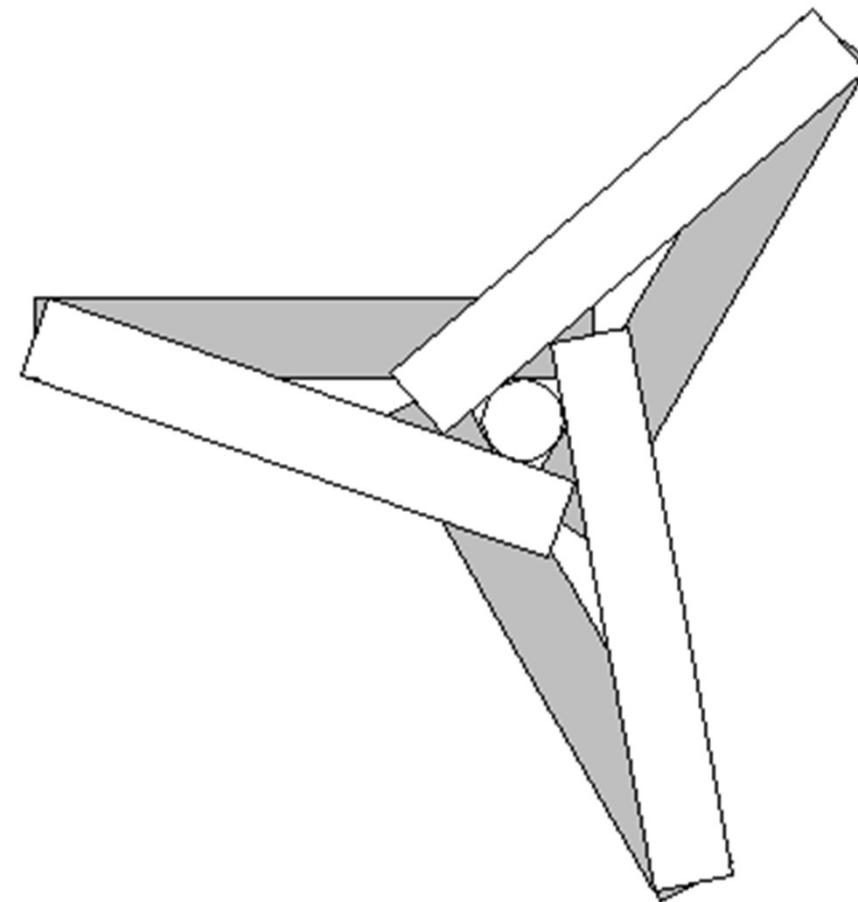
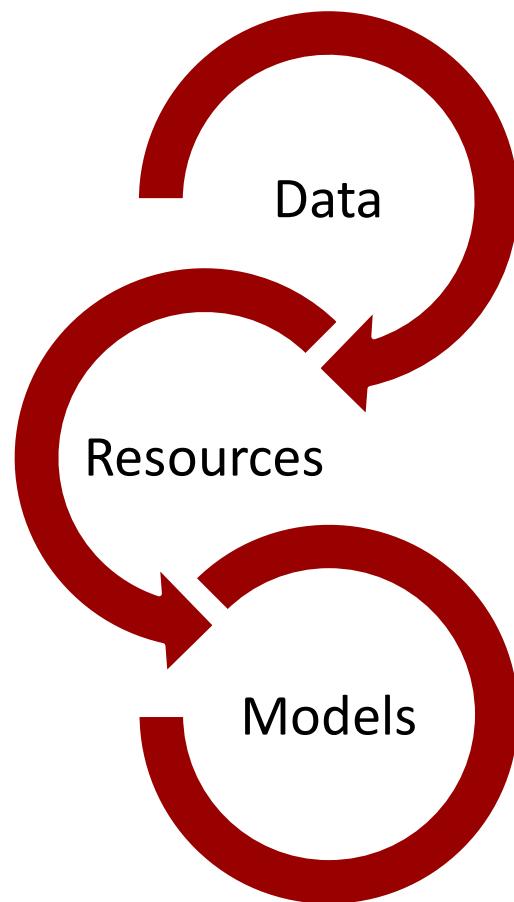
action adaptation adversarial attention based clouds convolutional
 data deep detection domain efficient estimation face
 feature generative graph human image instance joint
learning local matching model motion network
 neural object person point pose prediction recognition reconstruction
 representation robust scene segmentation semantic shape
 single structure supervised tracking transfer unsupervised video visual

Deep learning – where does it come from?

- ↗ **1943:** neurophysiologist Warren McCulloch and mathematician Walter Pitts – a **neuronal model** using an electrical circuit
- ↗ **1950:** Alan Turing created the world-famous **Turing Test**: a computer to pass, has to be able to convince a human that it is a human and not a computer.
- ↗ **1952:** Arthur Samuel created the first computer **program which could learn** as it ran. It was a game which played checkers.
- ↗ **1958:** Frank Rosenblatt designed the first artificial neural network, called **Perceptron**. The main goal of this was pattern and shape recognition.
- ↗ **1959**, Bernard Widrow and Marcian Hoff created
 - ↗ ADELINe to detect binary patterns (in a stream of bits, to predict the next one).
 - ↗ MADELINe to eliminate echo on phone lines, still in use today.
- ↗ **1986:** Neural networks use **back propagation** allowing multiple layers to be used in a neural network, creating what are known as 'slow learners'.
- ↗ Late **1980s** and **1990s** did not bring much to the field.
- ↗ **1997**, the IBM computer Deep Blue (a chess-playing computer) beat the world chess champion.
- ↗ **1998** AT&T Bell Laboratories on digit recognition resulted in good accuracy in detecting handwritten postcodes from the US Postal Service.

What changed today?

The magic triangle



Data

90% of all digital data were generated last 2 years.

Every minute of the day:

- 4M YouTube videos watched
 - 456K tweets on Twitter
 - 46K photos posted in Instagram
 - 16M text messages sent
 - 103M spam emails sent



DL Datasets



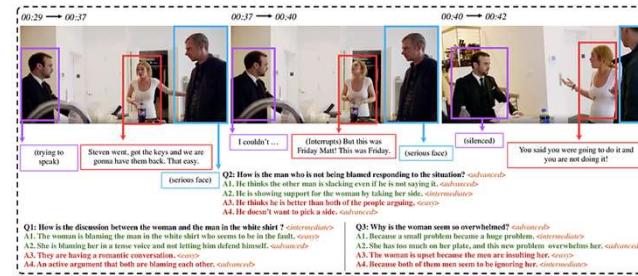
Lyft Level 5



LVIS Challenge: 2.2M masks, 16K images



Places2: 10M images



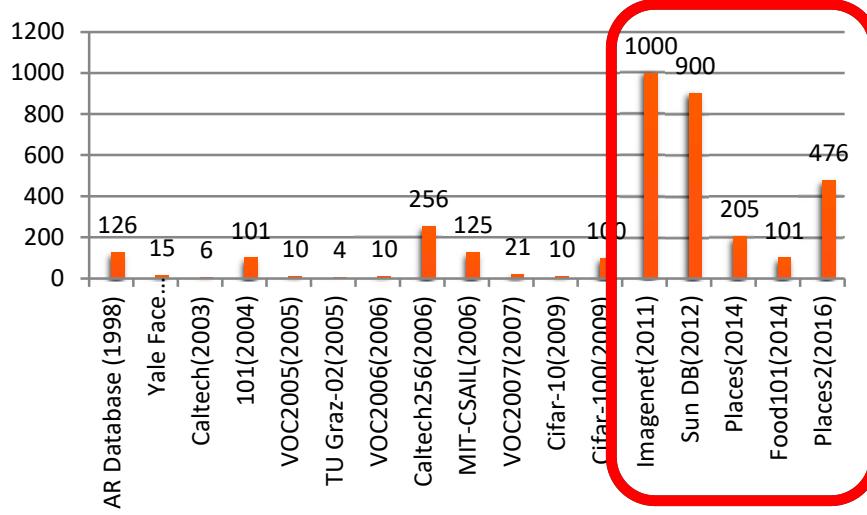
SocialIQ



TACO: Waste in the wild

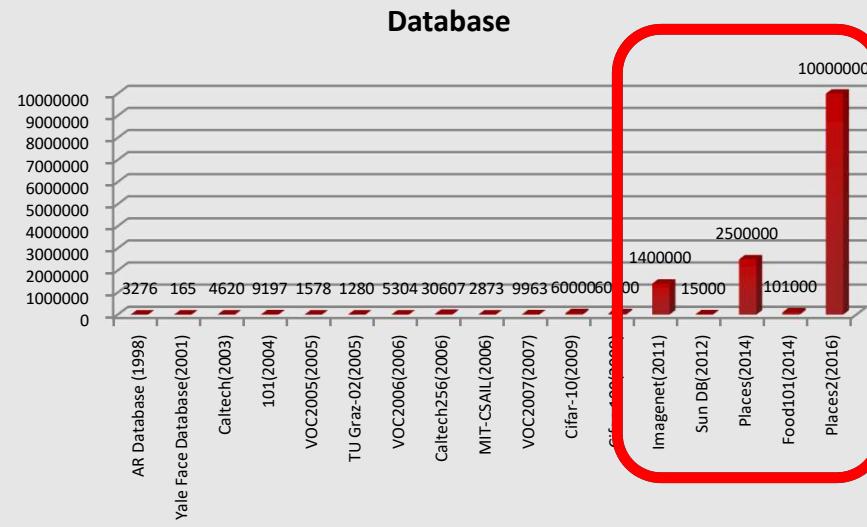
Image databases evolution

Number of objects/Database



ImageNet &
Deep learning

Number of images/Database



Imagenet



IM[▲]GENET

www.image-net.org

22K categories and **14M** images

- Animals
 - Bird
 - Fish
 - Mammal
 - Invertebrate
- Plants
 - Tree
 - Flower
 - Food
 - Materials
- Structures
 - Artifact
 - Tools
 - Appliances
 - Structures
- Person
 - Scenes
 - Indoor
 - Geological Formations
 - Sport Activities



Deng, Dong, Socher, Li, Li, & Fei-Fei, 2009

Back **TED** Ideas worth spreading

WATCH DISCOVER ATT

Fei-Fei Li:

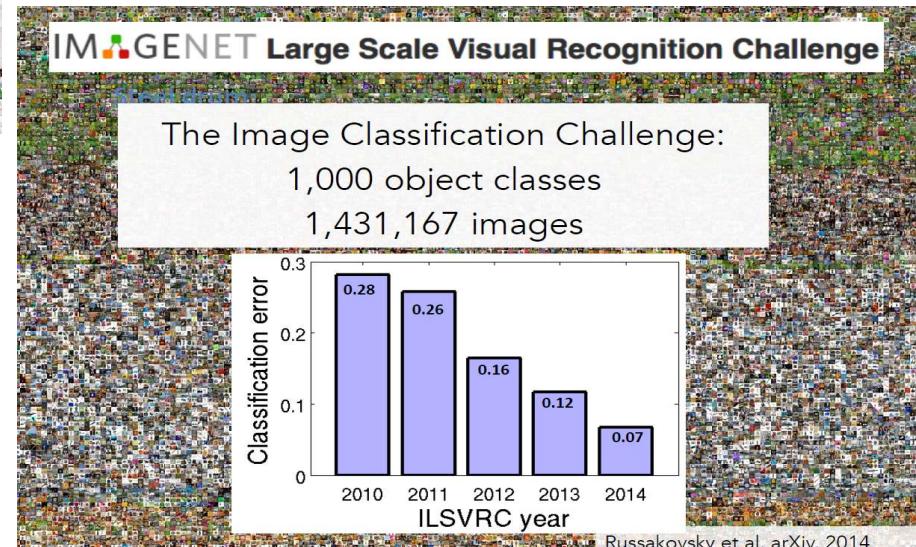
How we're teaching computers to understand pictures

TED2015 · 17:58 · Filmed Mar 2015

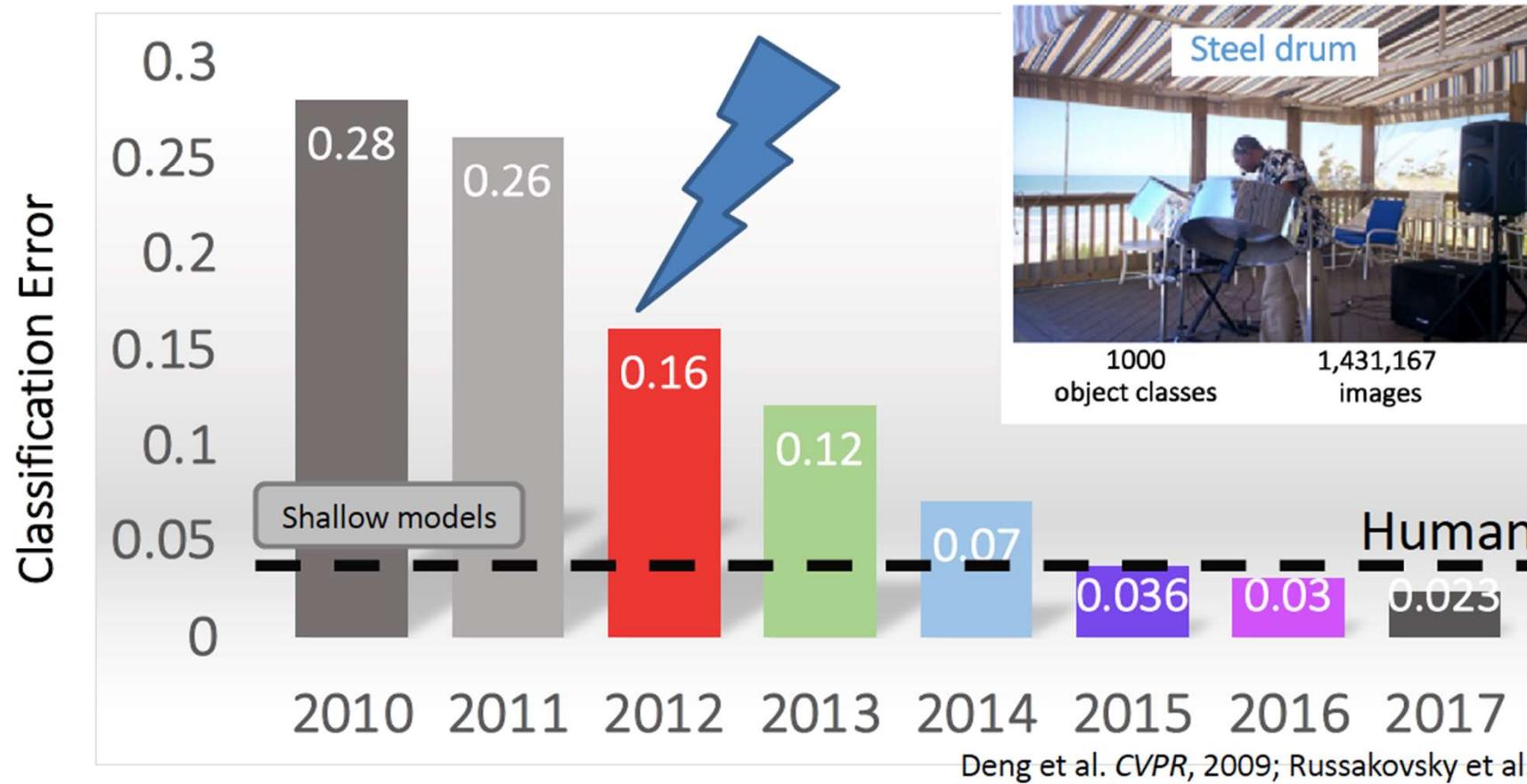
26 subtitle languages ⓘ View interactive transcript



Later Download Rate share 1,607,730 Total views

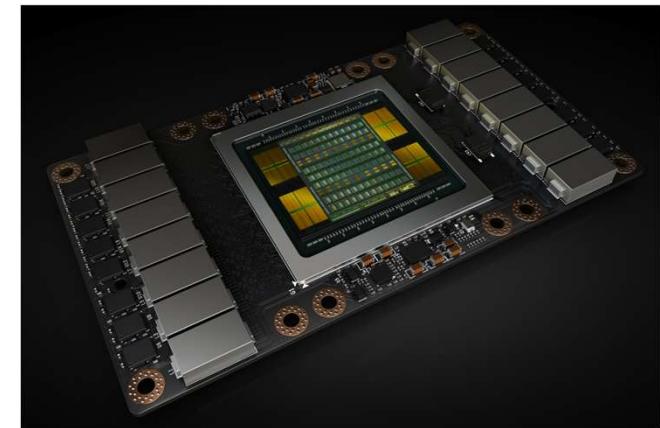


IMAGENET Classification Task



The Importance of GPUs

- ↗ Nvidia Tensor Cores - 2017
- ↗ Google Tensor Processing Unit (TPU) - 2016
- ↗ Intel - Nervana Neural Processor - 2017
- ↗ GPUs in Cloud Computing (Google, 2017)



$$D = \left(\begin{array}{cccc} A_{0,0} & A_{0,1} & A_{0,2} & A_{0,3} \\ A_{1,0} & A_{1,1} & A_{1,2} & A_{1,3} \\ A_{2,0} & A_{2,1} & A_{2,2} & A_{2,3} \\ A_{3,0} & A_{3,1} & A_{3,2} & A_{3,3} \end{array} \right) \text{FP16 or FP32} \quad \left(\begin{array}{cccc} B_{0,0} & B_{0,1} & B_{0,2} & B_{0,3} \\ B_{1,0} & B_{1,1} & B_{1,2} & B_{1,3} \\ B_{2,0} & B_{2,1} & B_{2,2} & B_{2,3} \\ B_{3,0} & B_{3,1} & B_{3,2} & B_{3,3} \end{array} \right) \text{FP16} + \left(\begin{array}{cccc} C_{0,0} & C_{0,1} & C_{0,2} & C_{0,3} \\ C_{1,0} & C_{1,1} & C_{1,2} & C_{1,3} \\ C_{2,0} & C_{2,1} & C_{2,2} & C_{2,3} \\ C_{3,0} & C_{3,1} & C_{3,2} & C_{3,3} \end{array} \right) \text{FP16 or FP32}$$

GPU cores is based on matrix multiplication

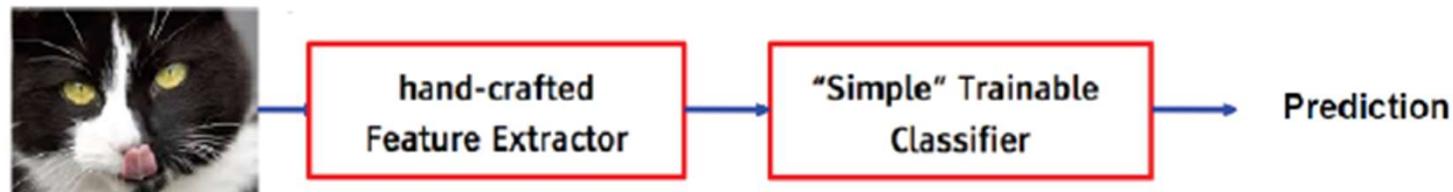


What is a Neural Network?

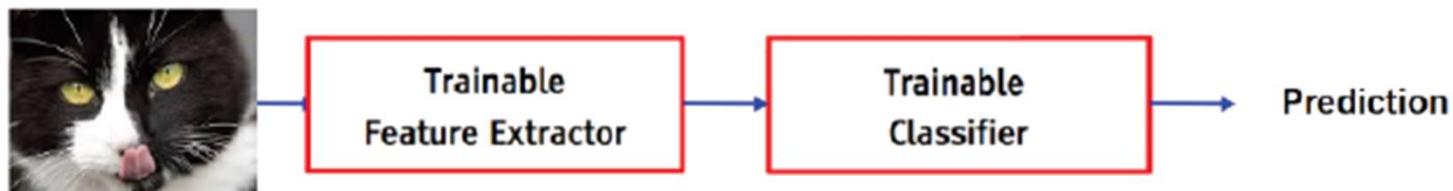
- ↗ AI, Machine learning & Deep learning
- ↗ What is a Convolutional Neural Network?
 - ↗ Layers
 - ↗ Optimization
- ↗ Applications

Why Deep learning?

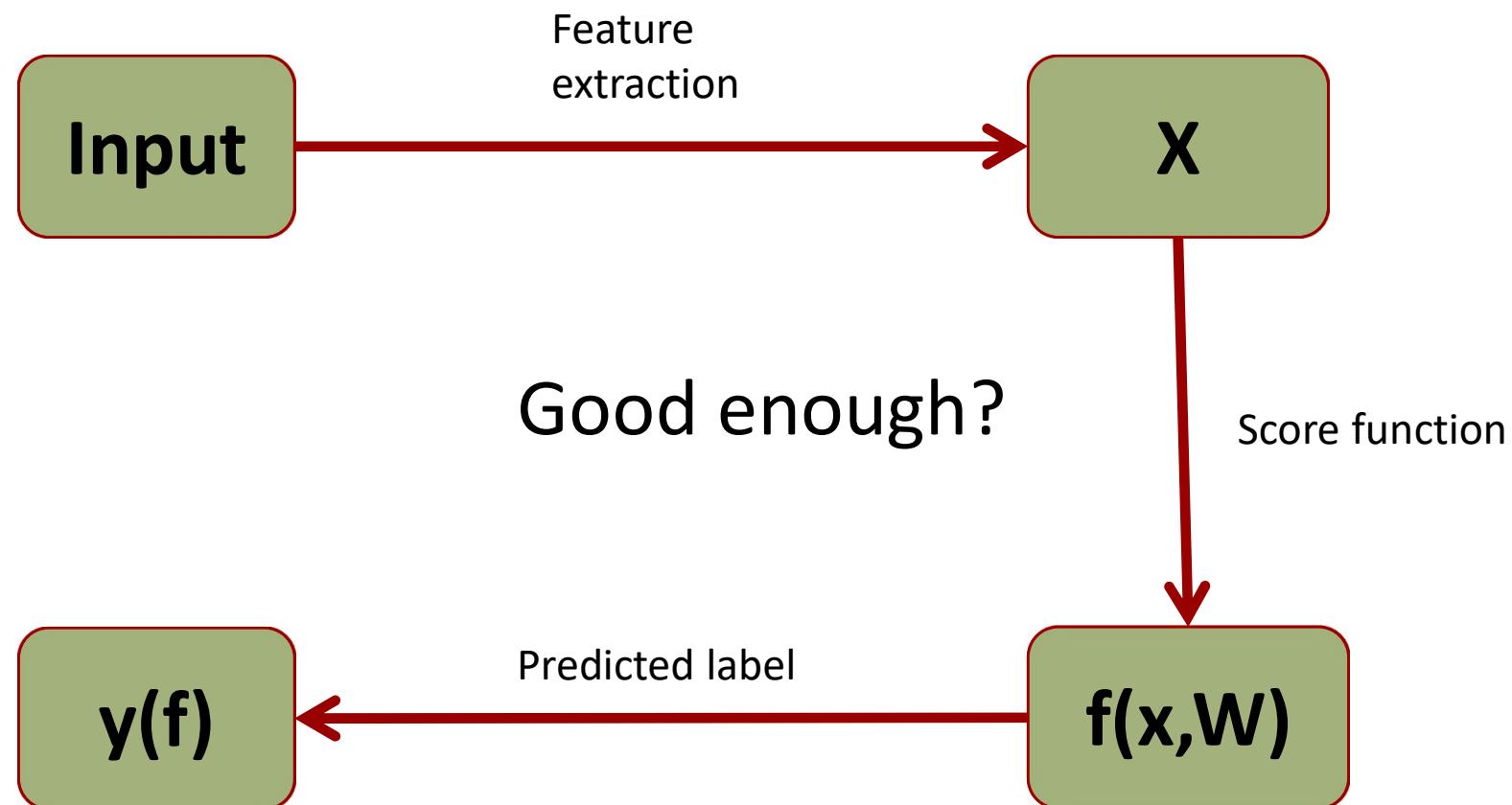
Classical way of solving Computer Vision problems



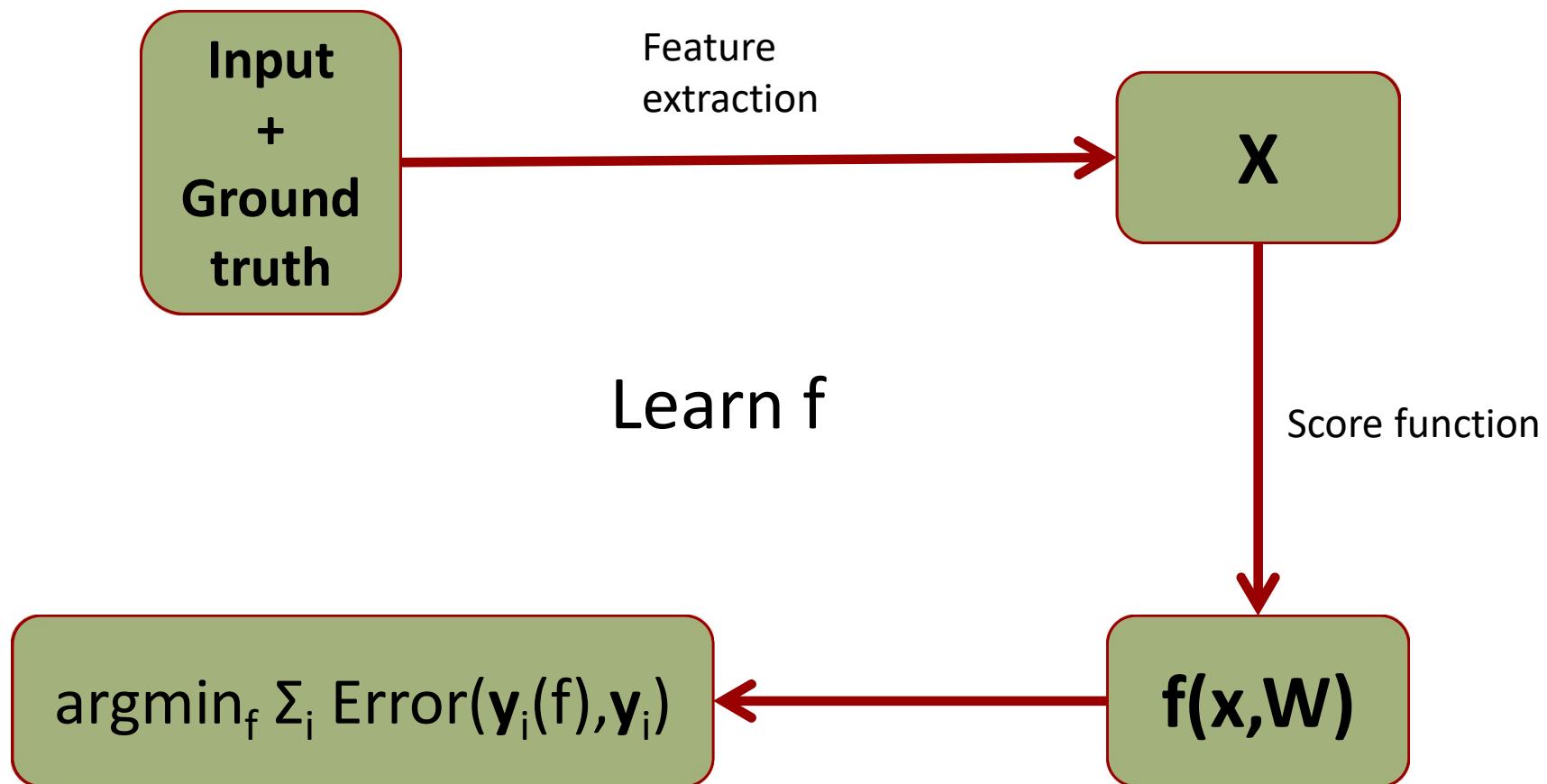
How Computer Vision problems are solved by Deep Learning



The learning pipeline



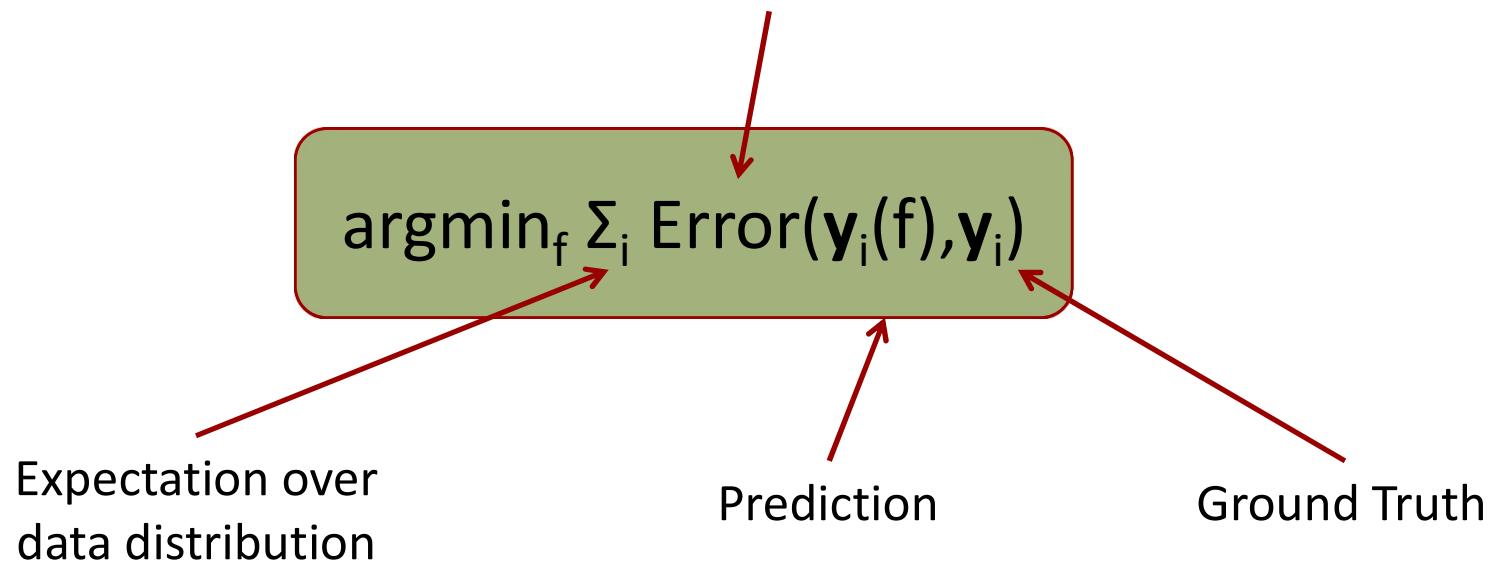
The traning process



The learning process

Training data $\{(x_i, y_i), i = 1, 2, \dots, n\}$

Measure of prediction quality (error, loss)

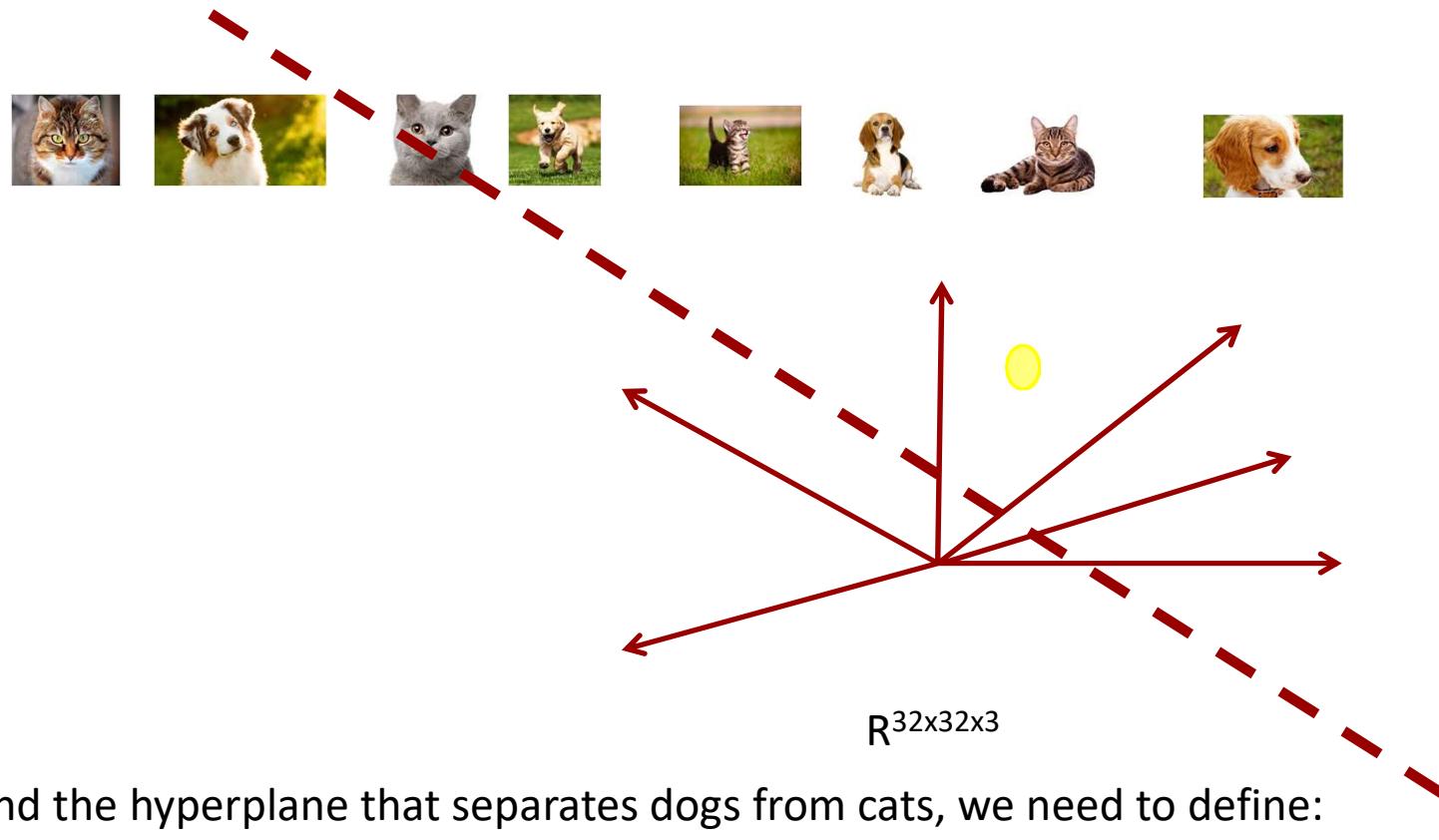


Loss function the negative conditional log-likelihood, with the interpretation that $f_i(X)$ estimates $P(Y=i|X)$:

$$L(f(x), y) = -\log f_i(x), \text{ where } f_i(x) \geq 0, \quad \sum_i f_i(x) = 1.$$

Linear classification

Given two classes how to learn a hyperplane to separate them?



To find the hyperplane that separates dogs from cats, we need to define:

- The score function
- The loss function
- And the optimization process.

Linear classification

How to project data in the feature space:

3x1

3072x1

$$f(x) = W x + b$$

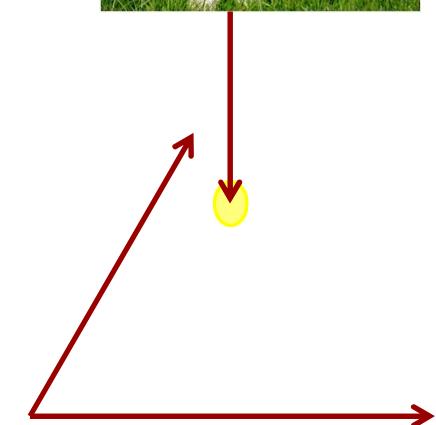
3x3072

3x1

If x is an image of $(32 \times 32 \times 3)$, $\rightarrow x$ in R^{3072} ,

The matrix W is (3×3072) .

The bias vector b is 3-dimensional.



Linear classification

How to project data in the feature space:

$$f(x) = Wx + b$$

3×1 3072×1
 3×3072 3×1

If we have 3 classes, $f(x)$ will give 3 scores.

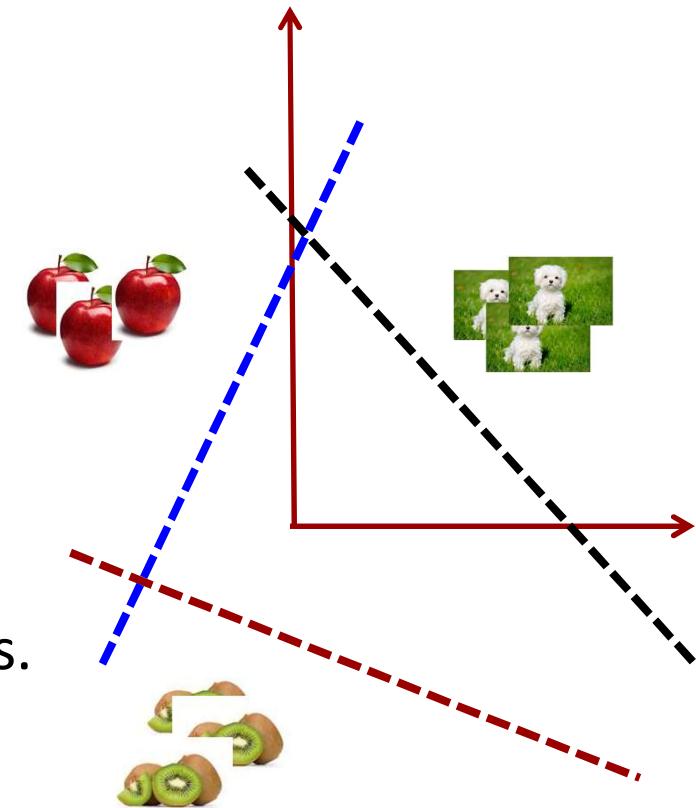
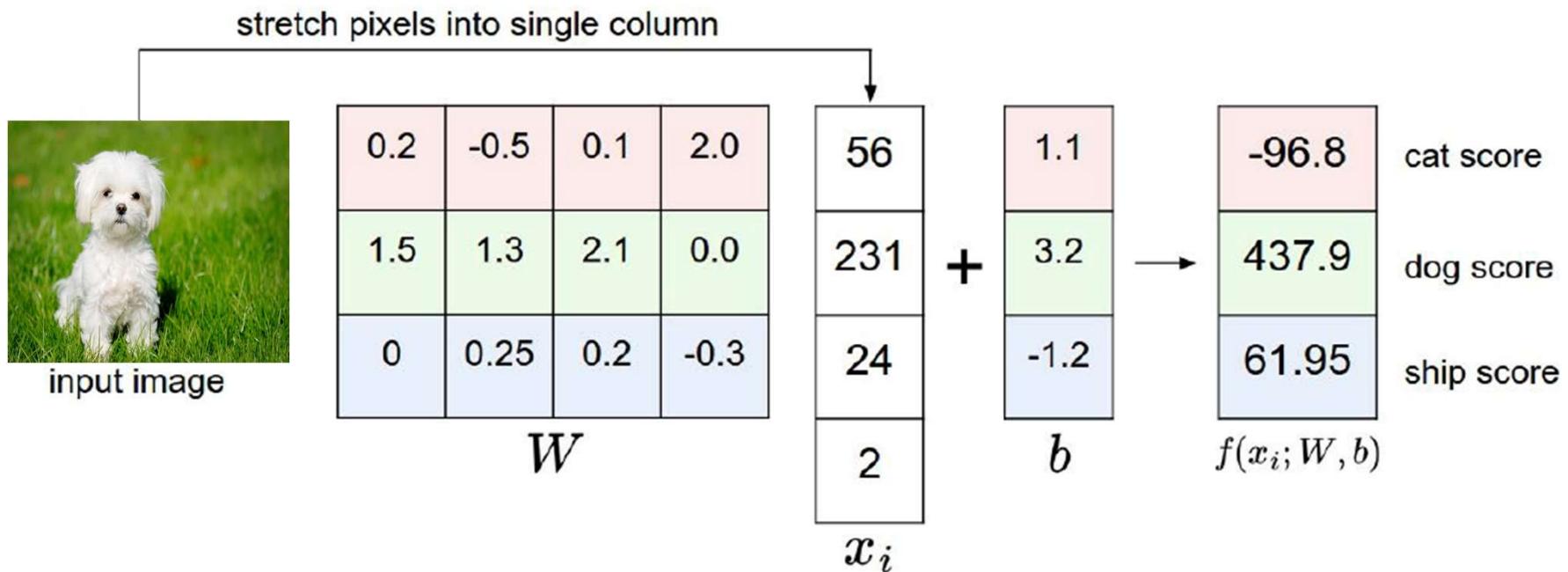


Image classification

Example with an image with 4 pixels, and 3 classes (**cat/dog/ship**)



Adapted from: Fei-Fei Li & Andrej Karpathy & Justin Johnson

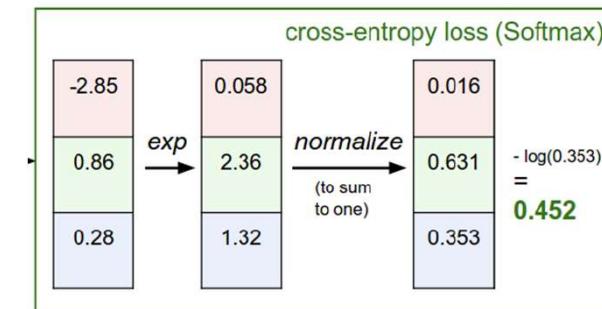
17:38

Loss function and optimisation

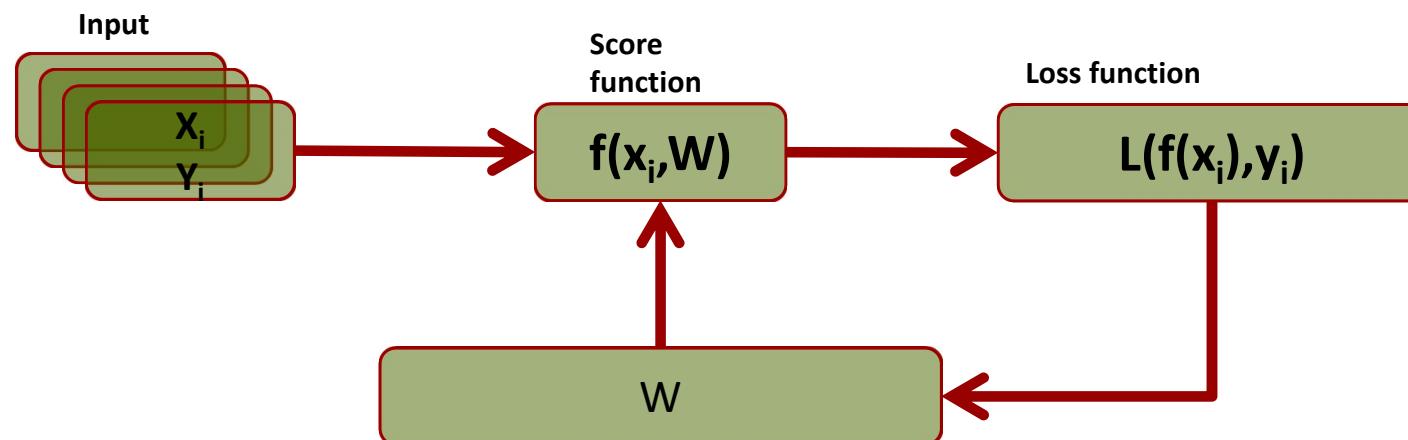
- ↗ **Question:** if you were to assign a single number to how unhappy you are with these scores, what would you do?

$$L_i = -\log \left(\frac{e^{f_{y_i}}}{\sum_j e^{f_j}} \right)$$

softmax function



Question : Given the score and the loss function, how to find the parameters W?

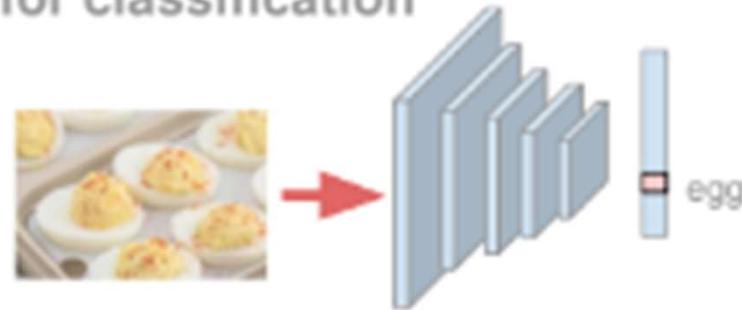


Single-label classification

- ↗ Change the lost function to the Binary cross-entropy function

L_b :

Conventional approach:
CNN for classification



Softmax

$$P(y_i|x) = \frac{\exp^{f(x)_i}}{\sum_i \exp^{f(x)_j}}$$

loss
←

Categorical cross-entropy

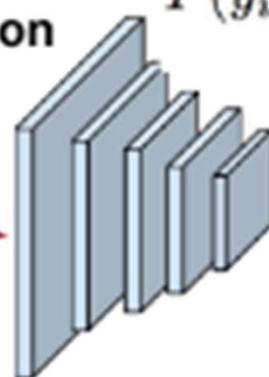
$$L_c = - \sum_x \log(P(\hat{y}_x|x))$$

Multi-label classification

- ↗ Change the lost function to the Binary cross-entropy function

L_b :

Our proposal:
Adaptation for
multi-label recognition



Softmax

Sigmoid

$$P(y_i|x) = \frac{1}{1 + \exp^{-f(x)_i}}$$



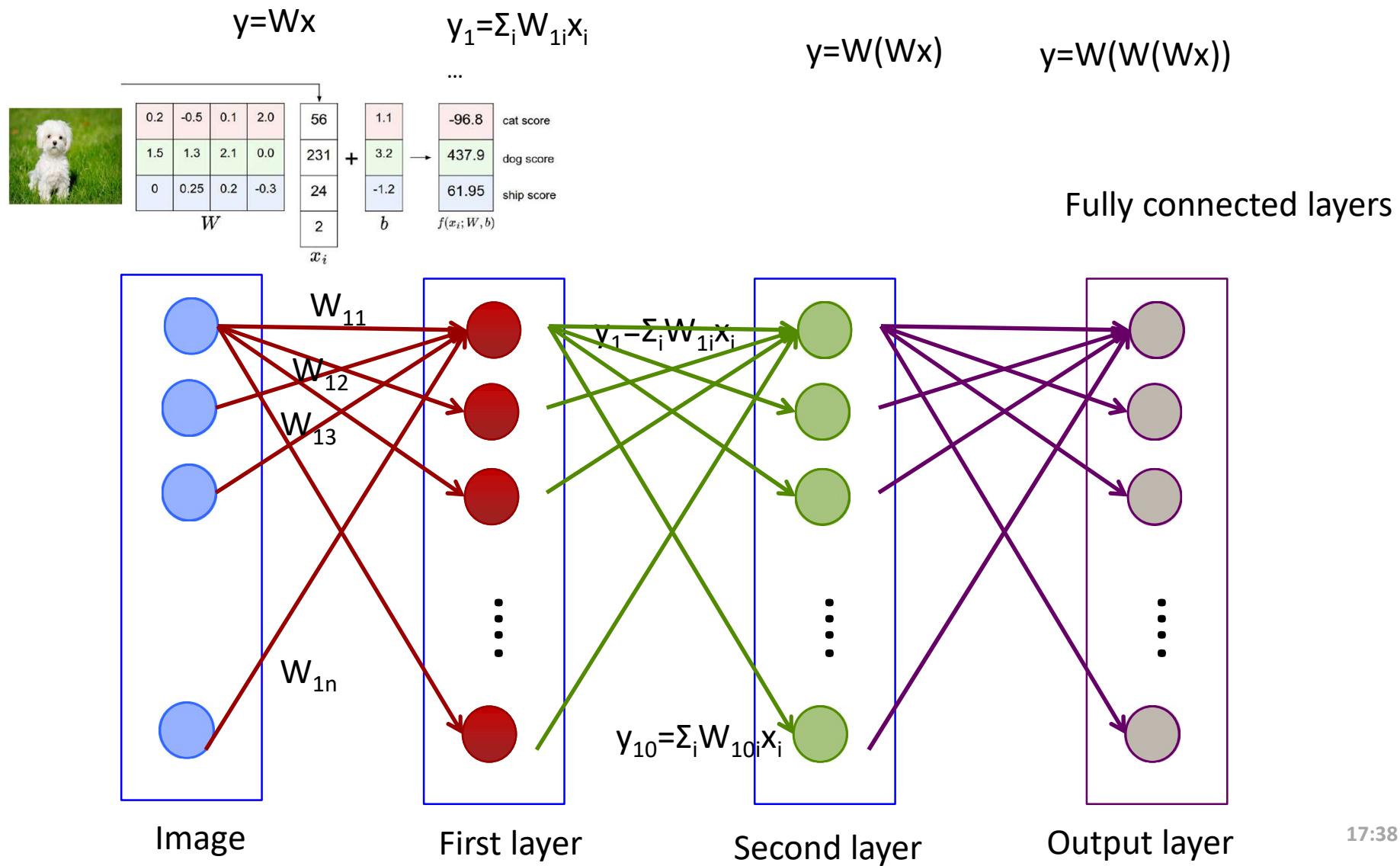
loss

Binary cross-entropy [2]

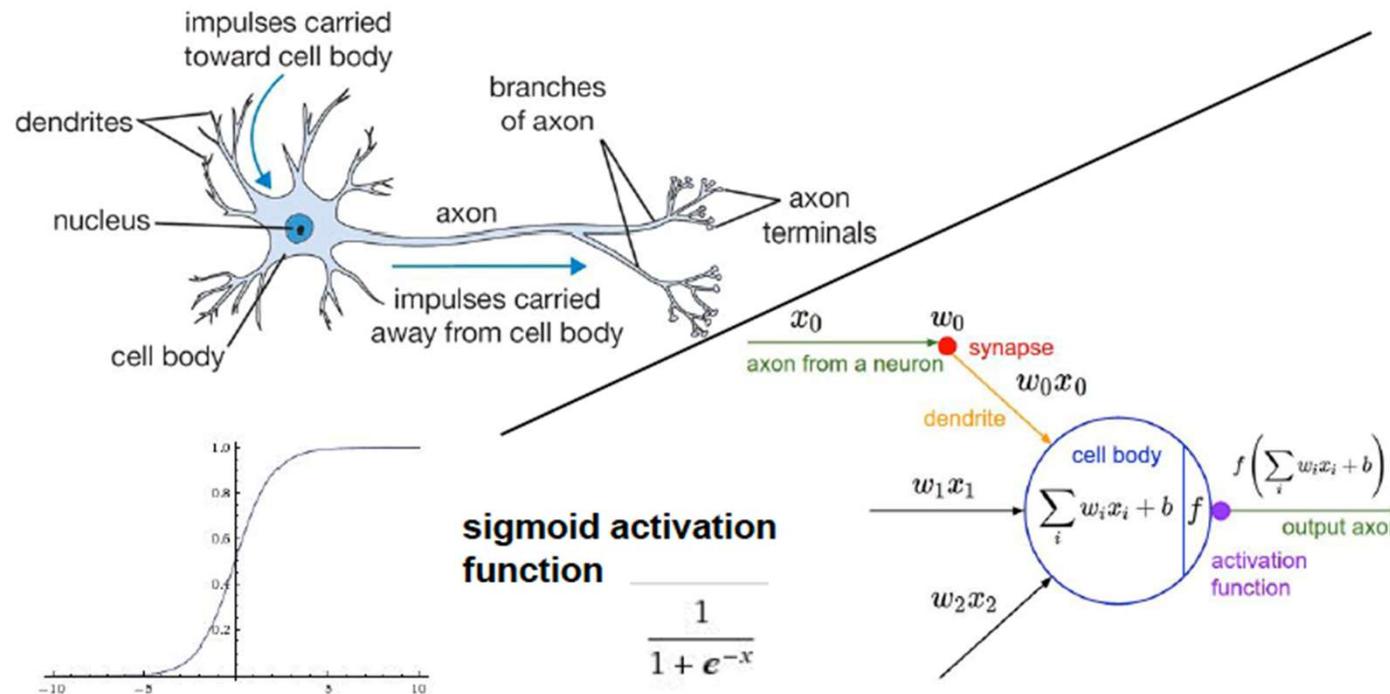
$$L_b = - \sum_x \sum_i^N (\hat{y}_x^i \cdot \log(P(y_i|x)) + (1-\hat{y}_x^i) \cdot \log(1-P(y_i|x)))$$

- ↗ AI, Machine learning & Deep learning
- ↗ What is a Convolutional Neural Network?
 - ↗ Layers
 - ↗ Optimization
- ↗ Applications

How is a CNN doing deep learning?



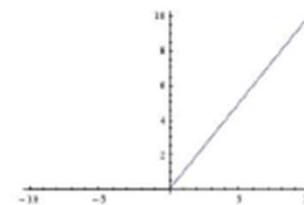
Why a CNN is a neural network?



sigmoid activation function

$$\frac{1}{1 + e^{-x}}$$

ReLU $\max(0, x)$

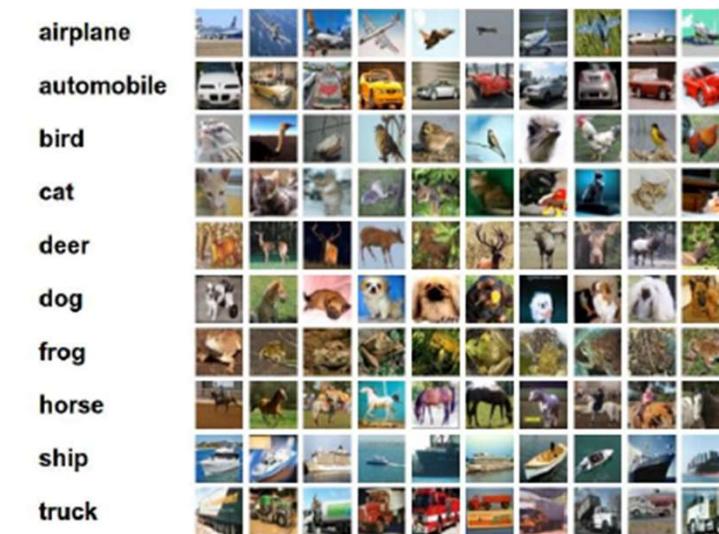


Modern CNNs – 10M neurons

Human CNNs – 5B of neurons.

From: Fei-Fei Li & Andrej Karpathy & Justin Johnson

Why is it convolutional?



$$f(x_i, W, b) = Wx_i + b$$

$$\begin{array}{c} \begin{array}{cccc} 0.2 & -0.5 & 0.1 & 2.0 \\ 1.5 & 1.3 & 2.1 & 0.0 \\ 0 & 0.25 & 0.2 & -0.3 \end{array} & \downarrow & \begin{array}{c} 56 \\ 231 \\ 24 \\ 2 \end{array} & + & \begin{array}{c} 1.1 \\ 3.2 \\ -1.2 \end{array} & \longrightarrow & \begin{array}{c} -96.8 \\ 437.9 \\ 61.95 \end{array} \\ W & & b & & & & f(x_i; W, b) \end{array}$$

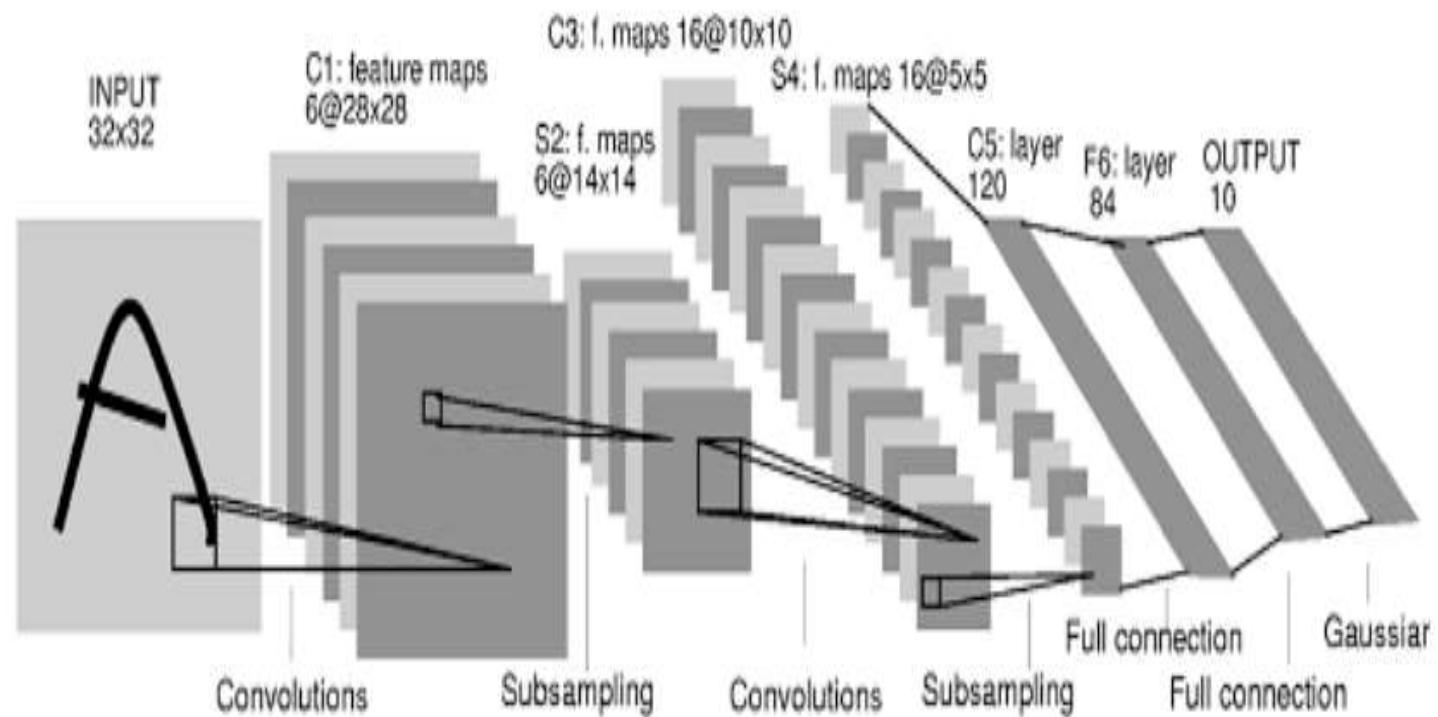


Adapted from: Fei-Fei Li & Andrej Karpathy & Justin Johnson

What is new in the Convolutional Neural Network?

1998

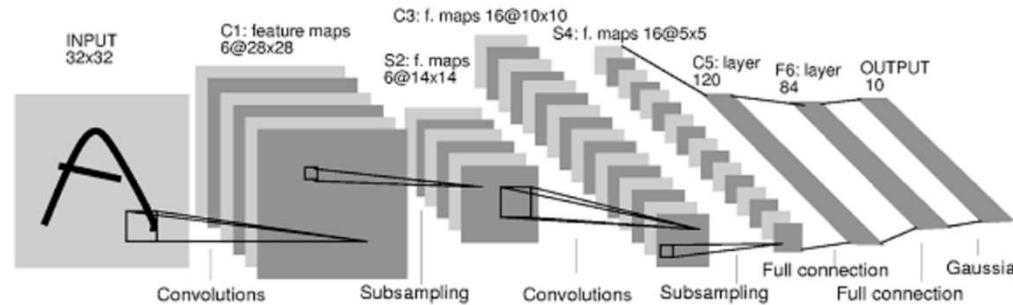
LeCun et al.



CNN evolution

1998

LeCun et al.

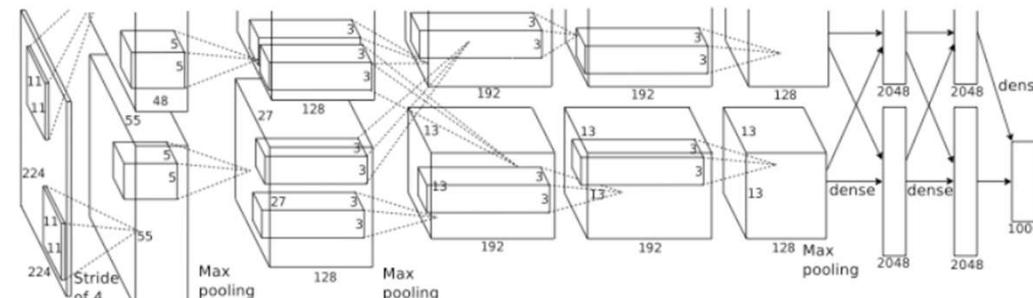


of transistors
 10^6
pentium® II

of pixels used in training
 10^7

2012

Krizhevsky
et al.

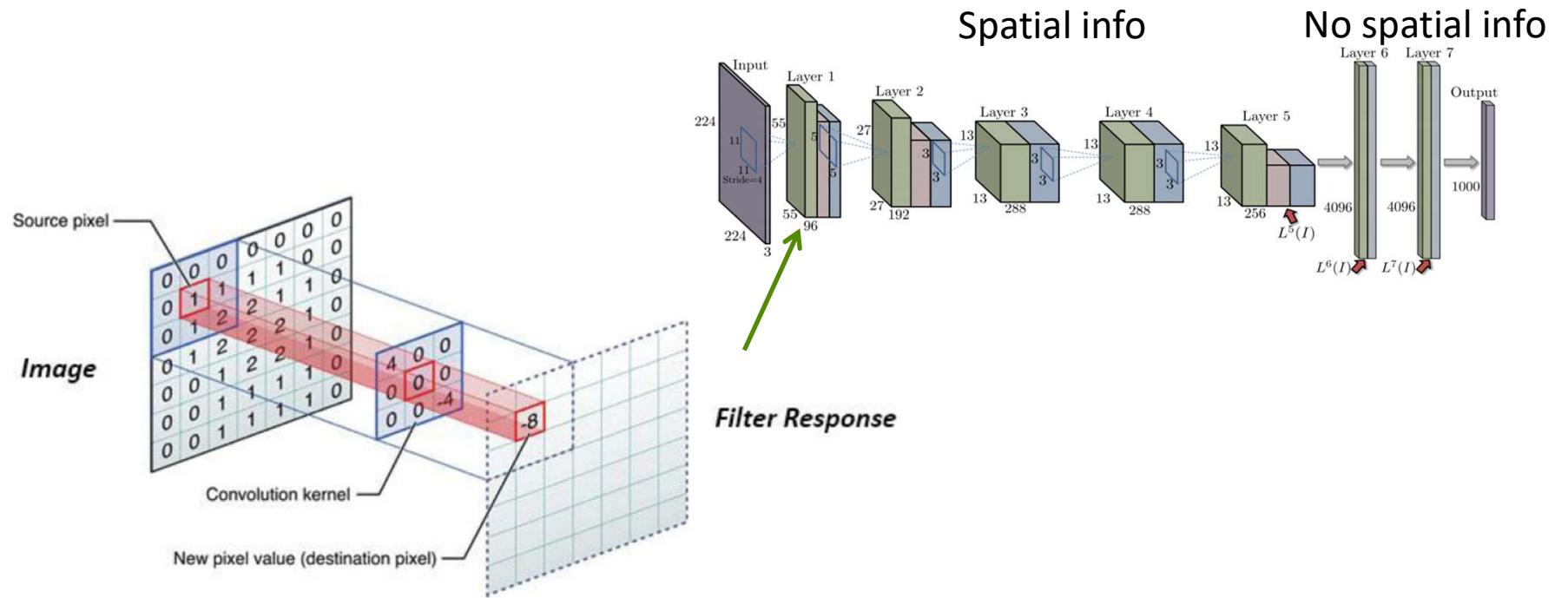


of transistors
 10^9



of pixels used in training
 10^{14}

Convolutional and Max-pooling layer



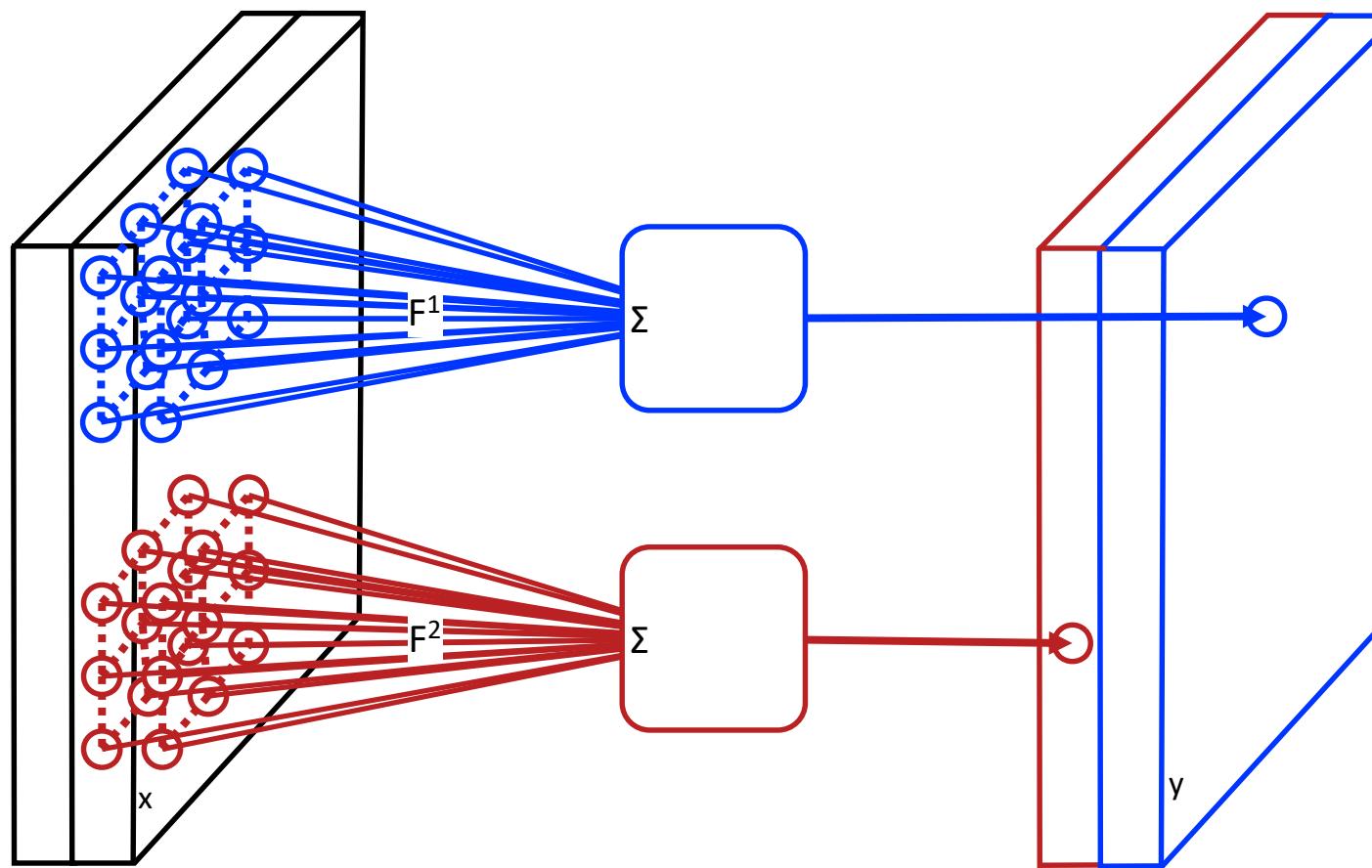
Convolutional layer

Why is it convolutional?

input features

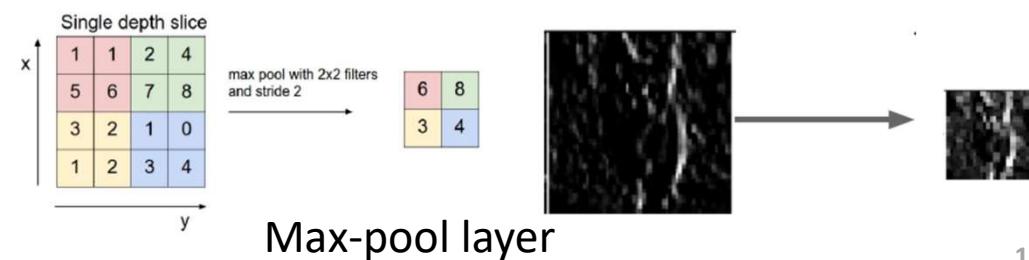
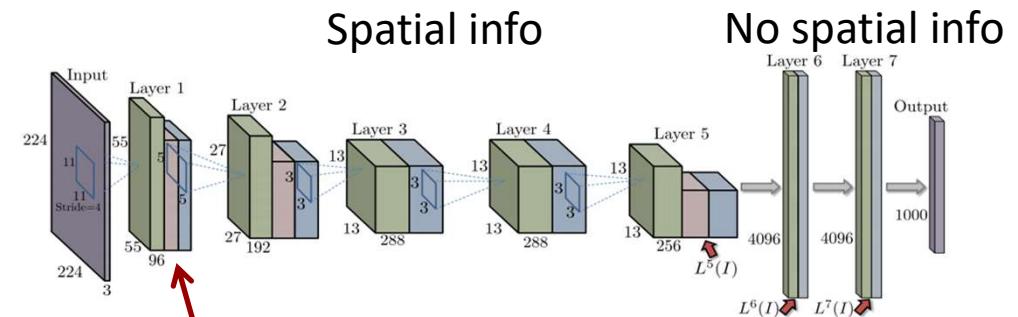
a bank of 2 filters

2-dimensional
output features

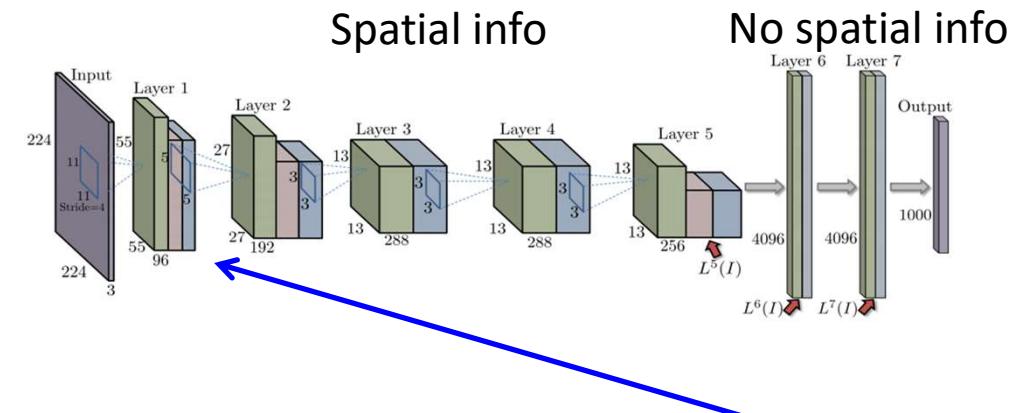


Convolution

Convolutional and Max-pooling layer

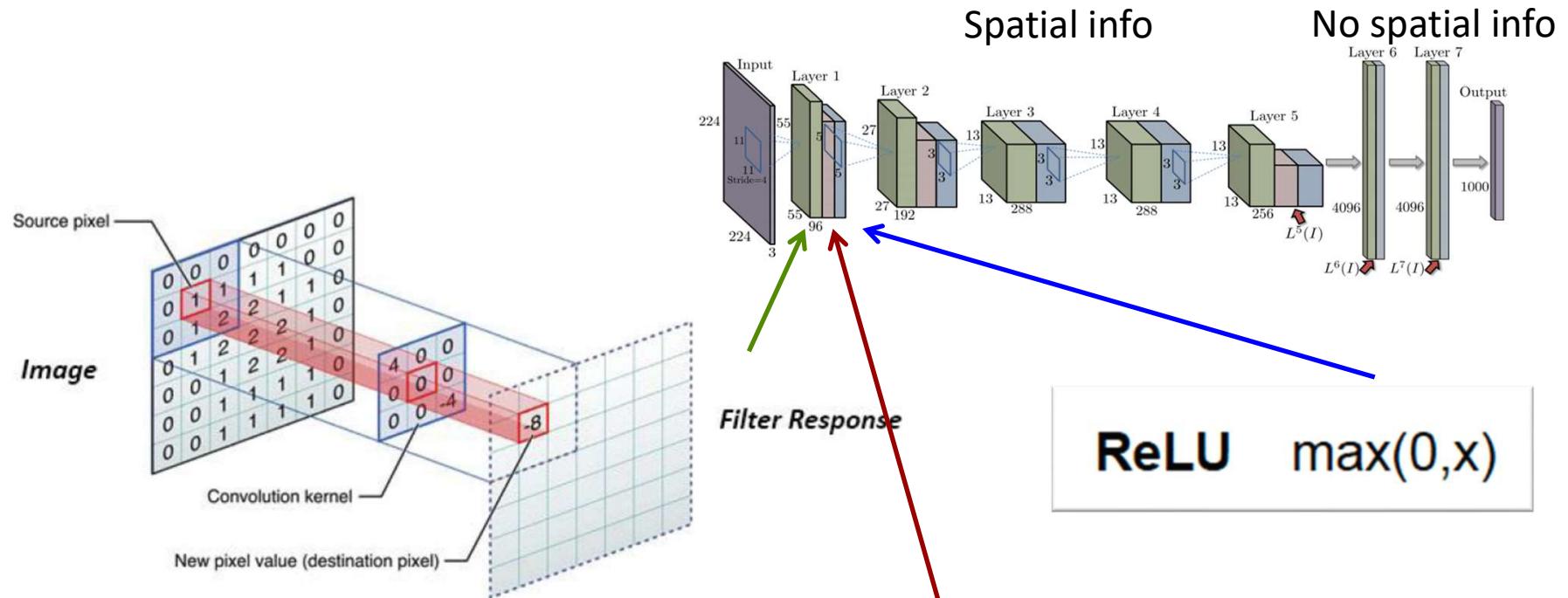


Convolutional and Max-pooling layer



ReLU $\max(0, x)$

Convolutional and Max-pooling layer



Convolutional layer

Single depth slice			
x	1	1	2
	5	6	7
y	3	2	1
	1	2	3

max pool with 2x2 filters and stride 2

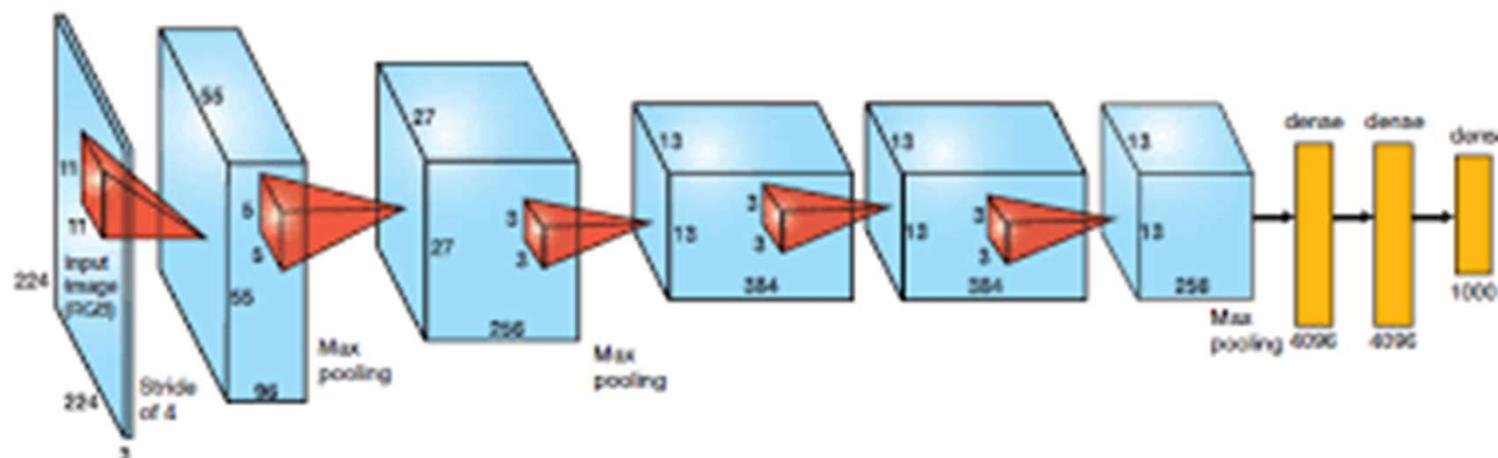
6	8
3	4

Max-pool layer



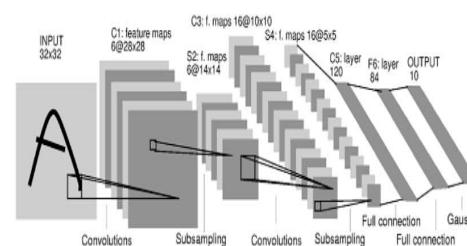
17:38

Why is it convolutional?



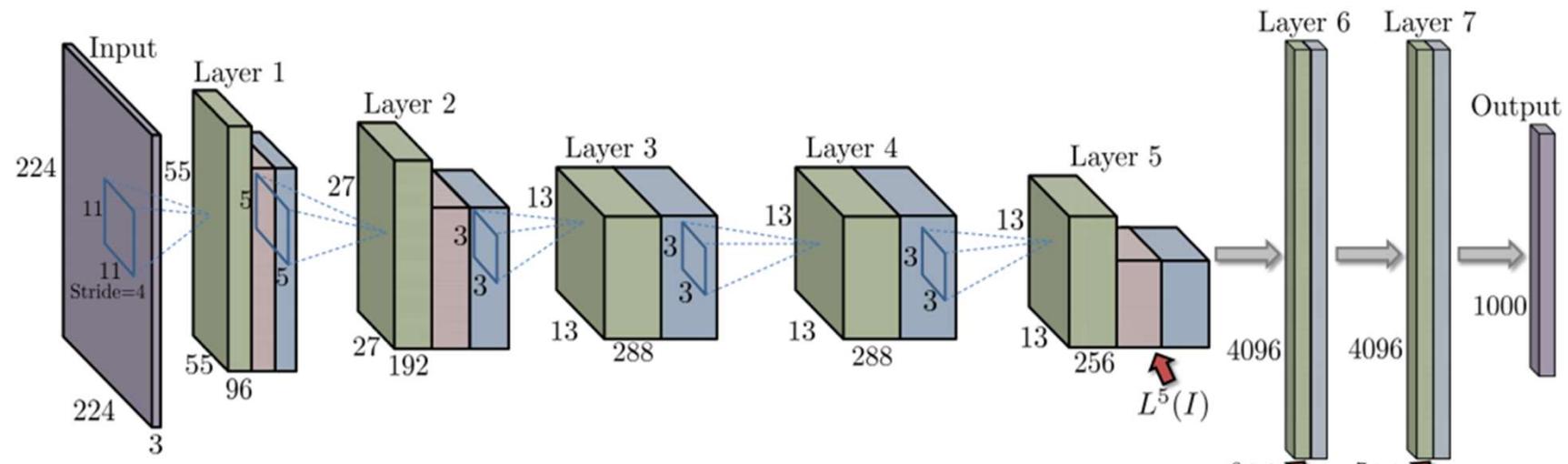
1998

LeCun et al.



LeCun, Chief AI Scientist for Facebook AI Research (FAIR), and a Silver Professor at New York University

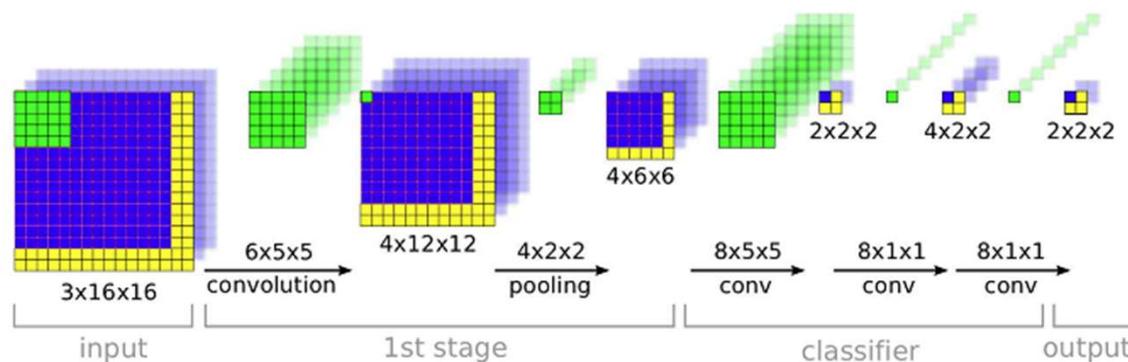
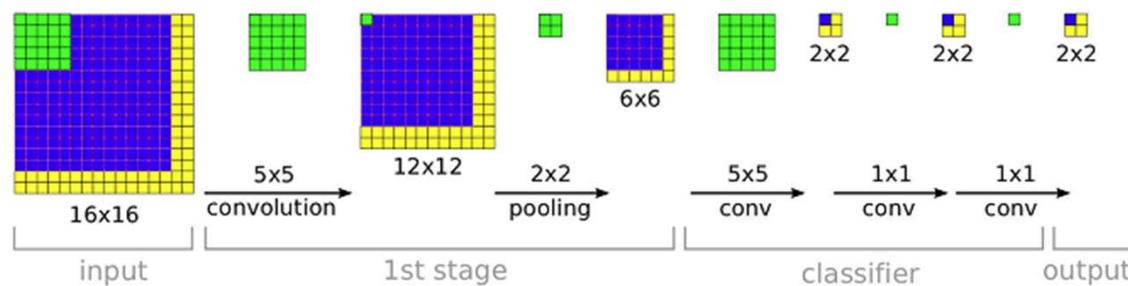
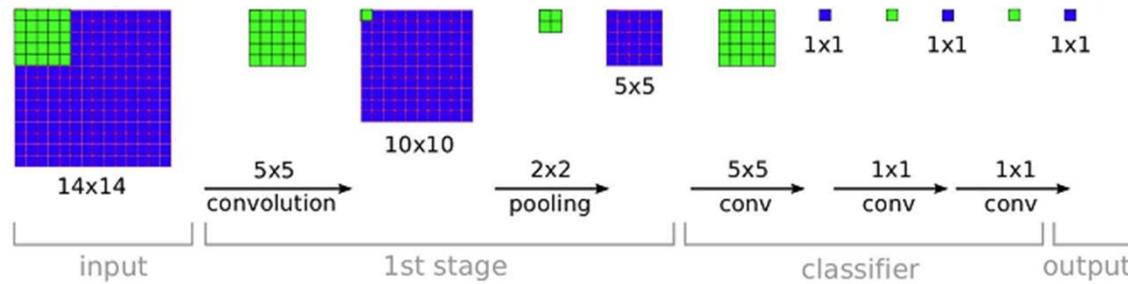
Training a CNN



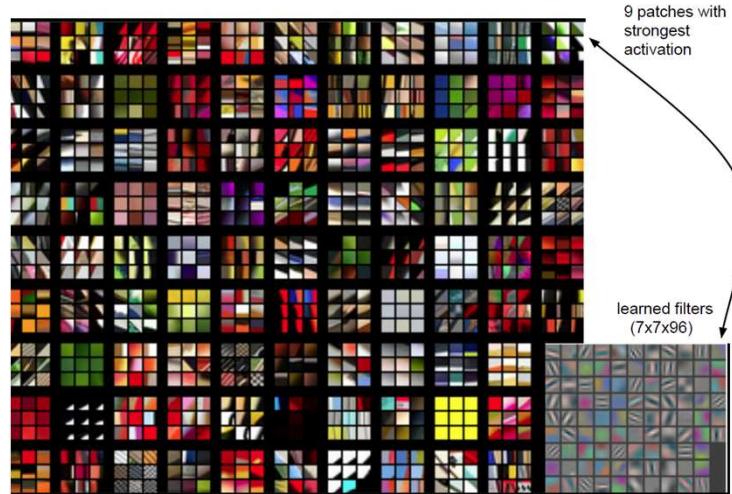
The process of training a CNN consists of training all hyperparameters: convolutional matrices and weights of the fully connected layers.

- **Several millions of parameters!!!**

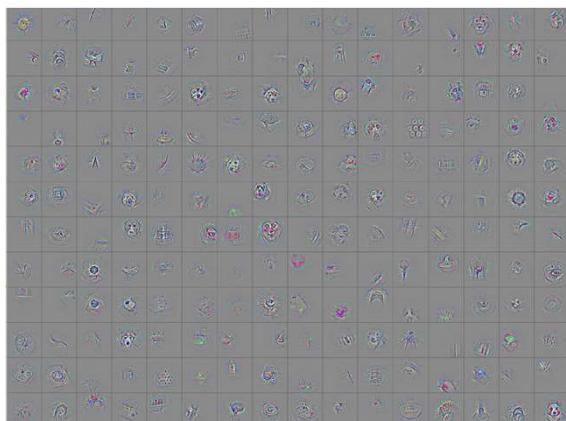
How does the CNN work?



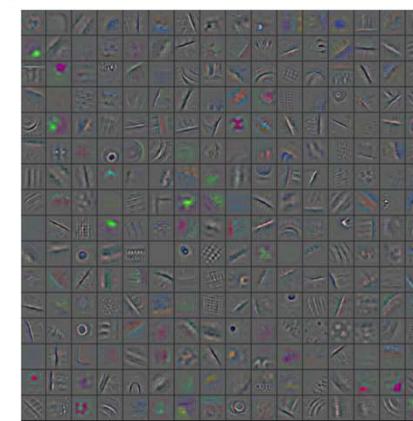
Learned convolutional filters



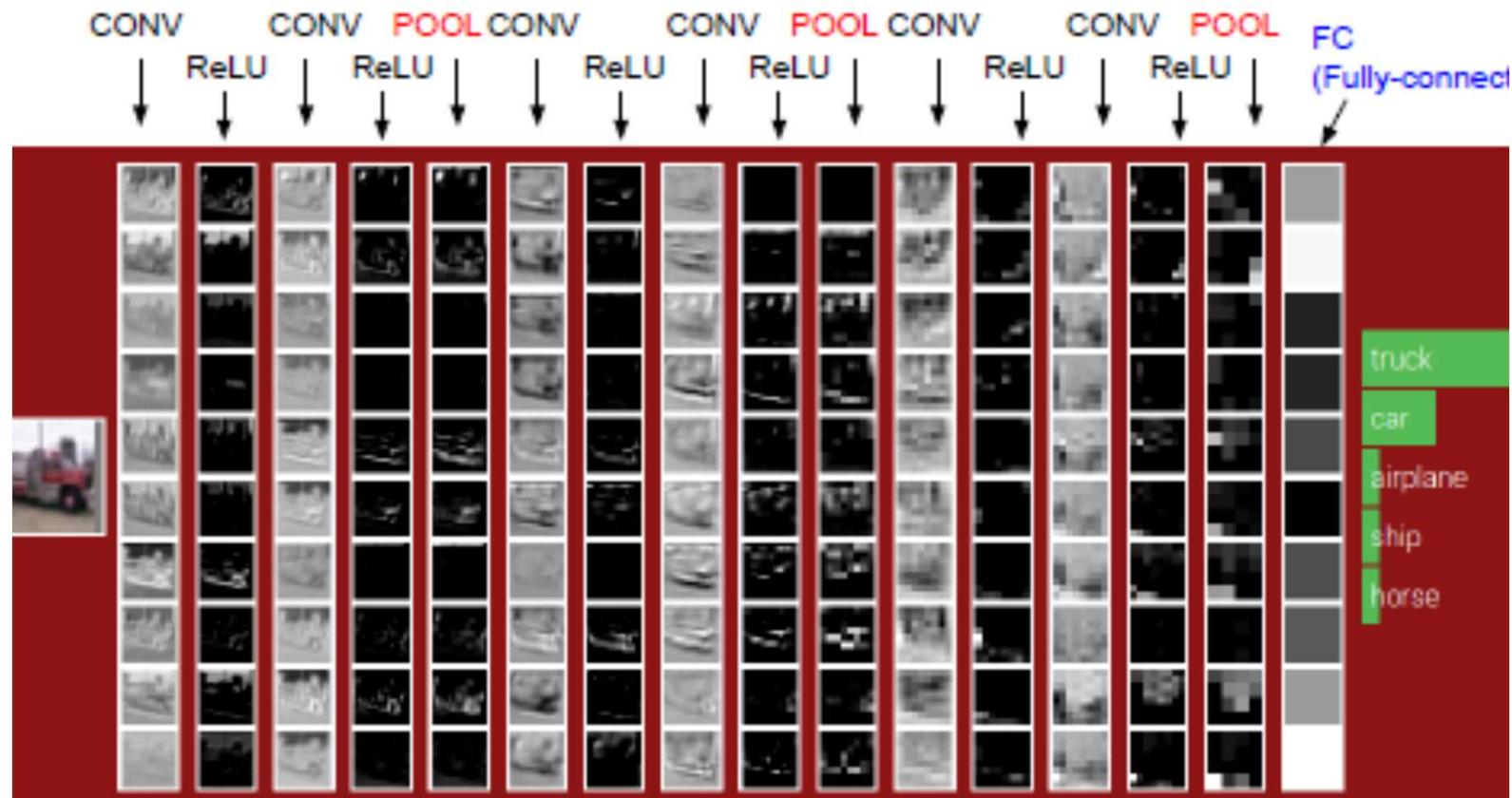
Strongest activations: Stage 5



Strongest activations: Stage 2



Example architecture



The trick is to train the weights such that when the network sees a picture of a truck, the last layer will say “truck”.

1001 benefits of CNN



Transfer learning: Fine tuning for object recognition

Replace and retrain the classifier on top of the ConvNet

Fine-tune the weights of the pre-trained network by continuing the backpropagation

Feature extraction by CNN

Object detection

Object segmentation

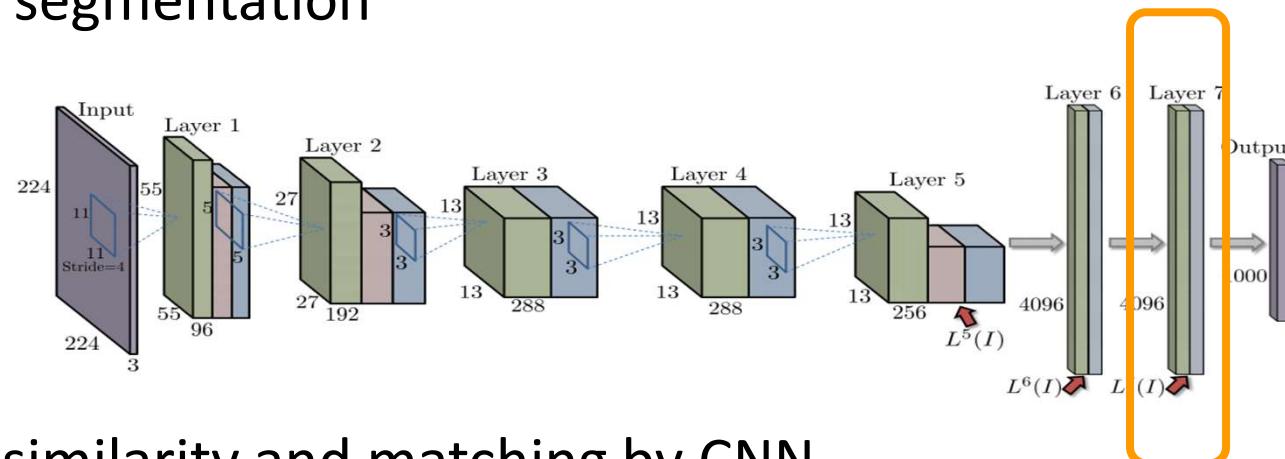


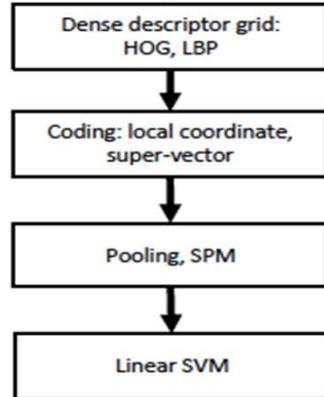
Image similarity and matching by CNN

Convolutional Neural Networks (4096 Features)

IMAGENET Large Scale Visual Recognition Challenge

Year 2010

NEC-UIUC

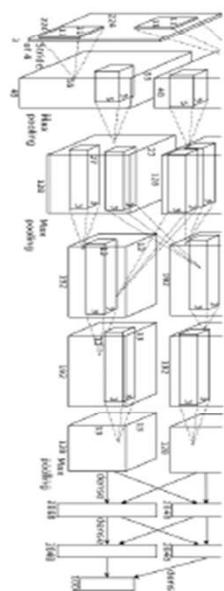


[Lin CVPR 2011]

Lion image by Swissfrog
is
licensed under CC BY 3.0.

Year 2012

SuperVision



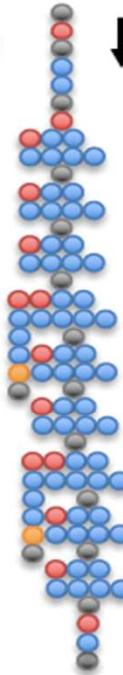
[Krizhevsky NIPS 2012]

Figure copyright Alex Krizhevsky, Ilya Sutskever, and Geoffrey Hinton, 2012.
Reproduced with permission.

Year 2014

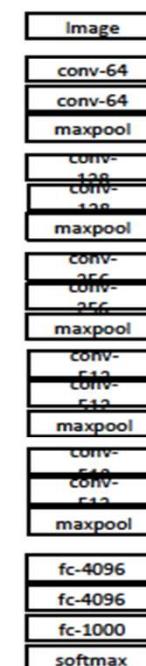
GoogLeNet

- Pooling
- Convolution
- Softmax
- Other



[Szegedy arxiv 2014]

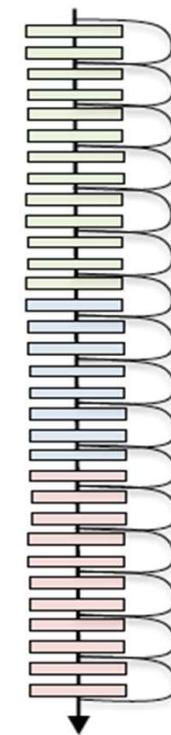
VGG



[Simonyan arxiv 2014]

Year 2015

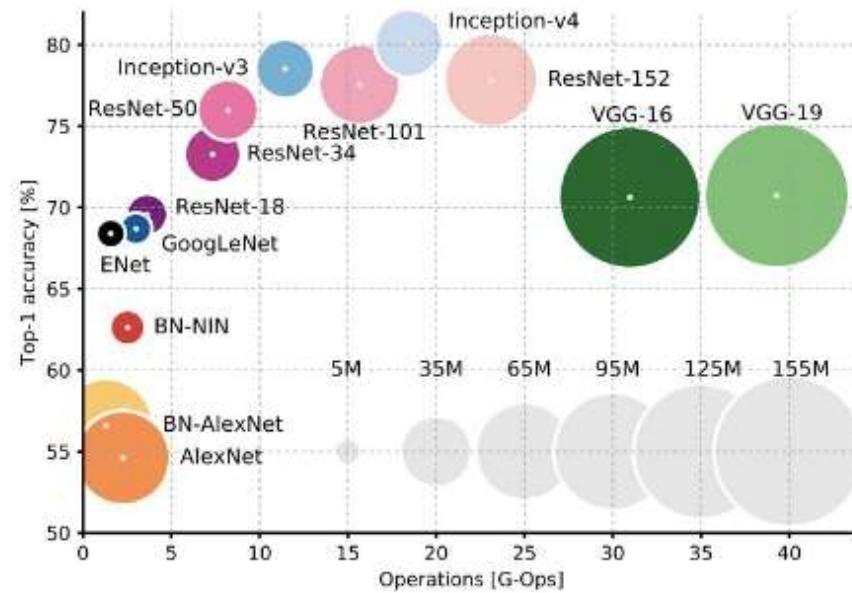
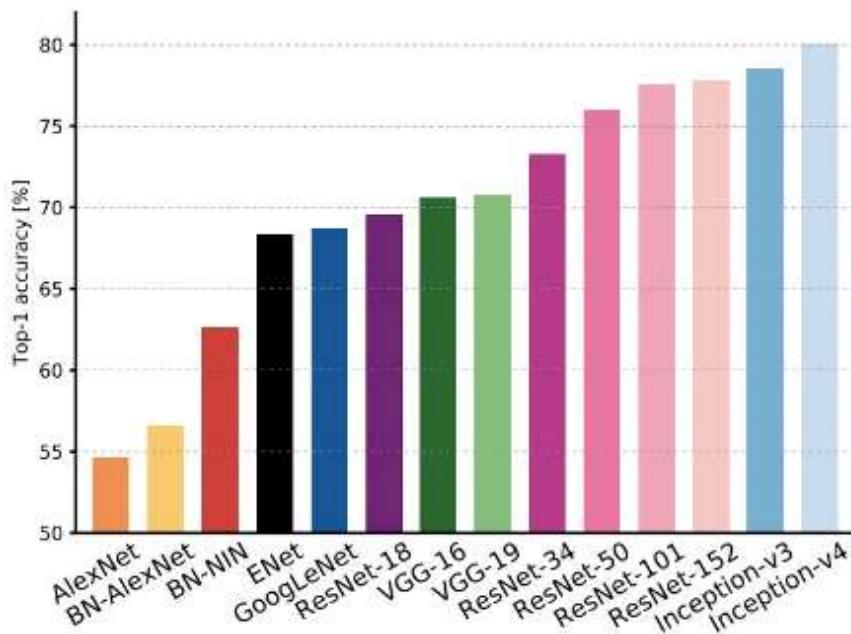
MSRA



[He ICCV 2015]

$$o = f_1(f_2(f_3(\dots f_n(x, \theta_n))))$$

Analysis of CNNs



- Millions of parameters!!!

The process of training a CNN consists of training all hyperparameters: convolutional matrices and weights of the fully connected layers.

- ↗ AI, Machine learning & Deep learning

- ↗ What is a Convolutional Neural Network?

- ↗ Layers

- ↗ Loss function and CNN Optimization

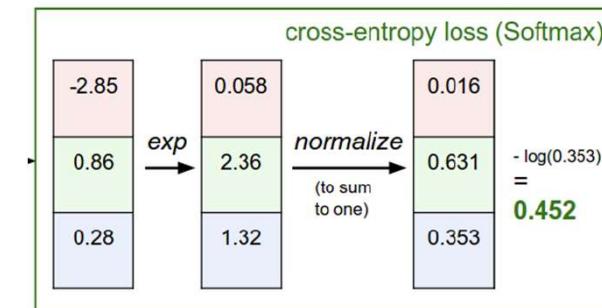
- ↗ Applications

Loss function and optimisation

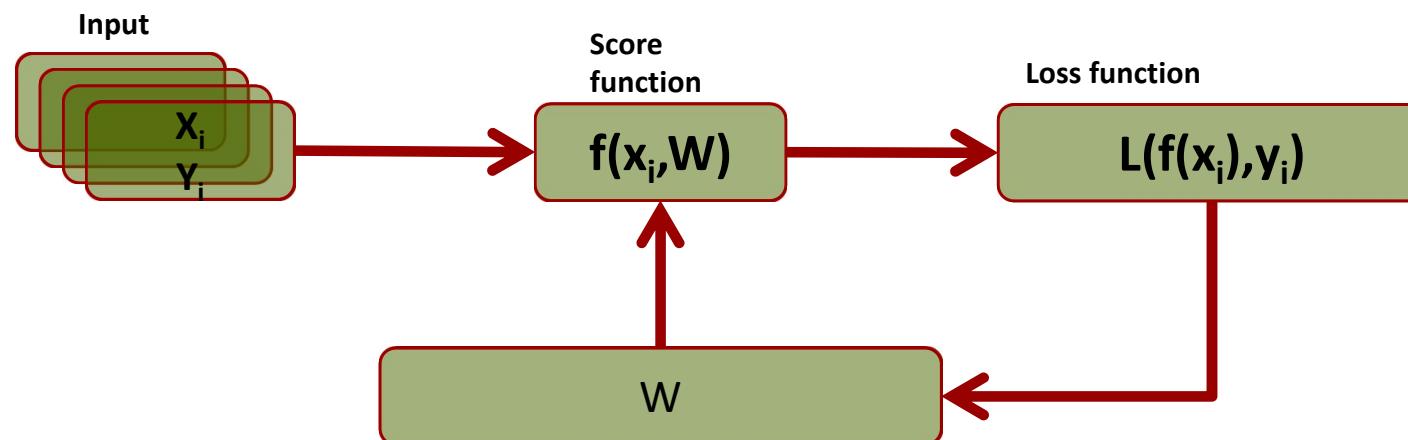
↗ **Question:** if you were to assign a single number to how unhappy you are with these scores, what would you do?

$$L_i = -\log \left(\frac{e^{f_{y_i}}}{\sum_j e^{f_j}} \right)$$

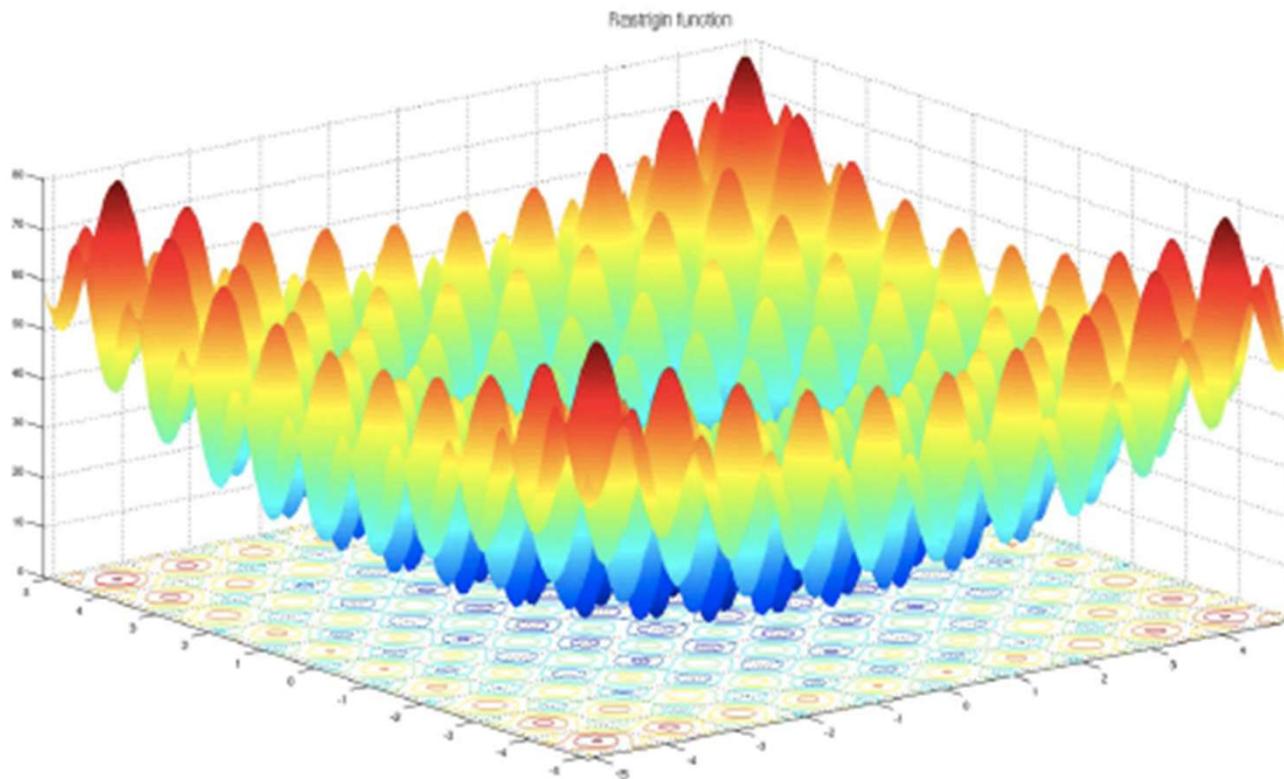
softmax function



Question : Given the score and the loss function, how to find the parameters W?



Optimization

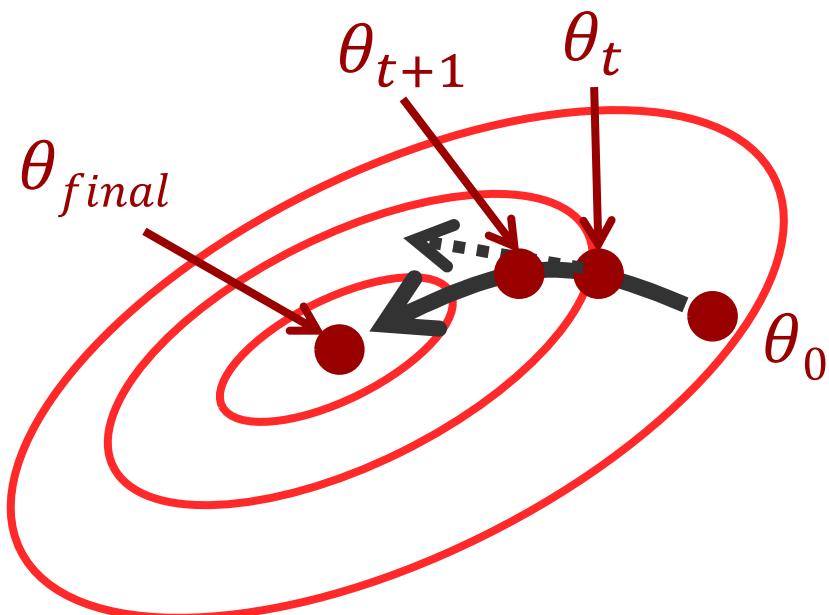


The process of training a CNN consists of training all hyperparameters: convolutional matrices and weights of the fully connected layers.

- Millions of parameters!!!

Gradient descent

$$o = f_1(f_2(f_3(\dots, f_n(x, \theta_n))))$$

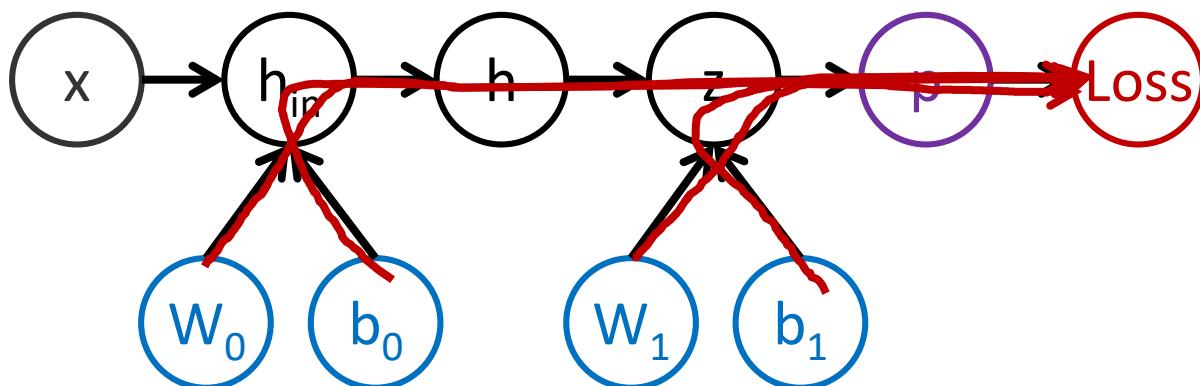


- ↗ Initialize θ_0 randomly
 - ↗ For t in $0, \dots, T_{\text{maxiter}}$
 $\theta^{t+1} = \theta^t - \eta_t \cdot \nabla L(\theta^t)$
- Gradient of the objective
- Learning rate (step size)

- Computation of $\nabla L(\theta^t)$ requires a full sweep over the training data
- Per-iteration comp. cost = $O(n)$

Chain rule

- Identify how each variable influence the loss



$$\frac{\partial \text{Loss}}{\partial W_1} = \frac{\partial \text{Loss}}{\partial p} \cdot \dots \cdot \frac{\partial z}{\partial W_1}$$

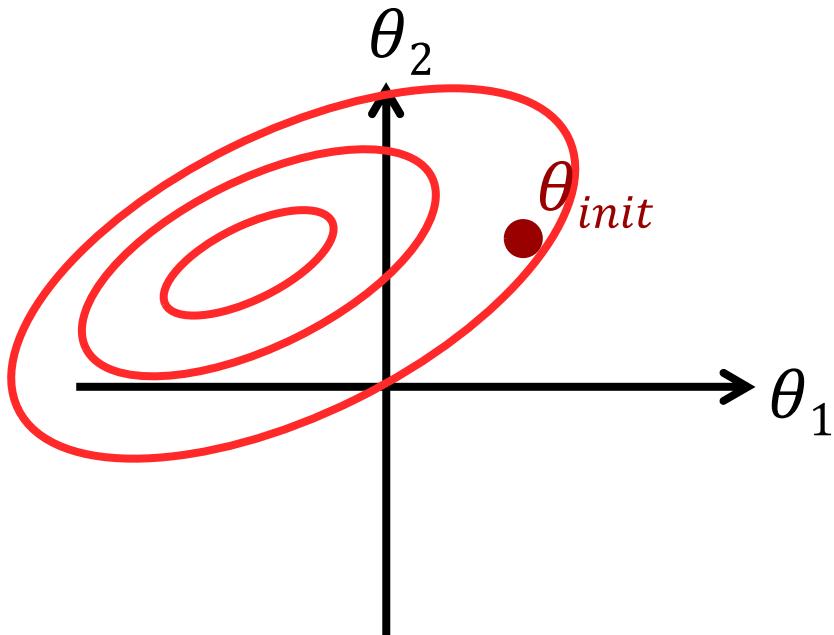
$$\frac{\partial \text{Loss}}{\partial b_1} = \frac{\partial \text{Loss}}{\partial p} \cdot \dots \cdot \frac{\partial z}{\partial b_1}$$

$$\frac{\partial \text{Loss}}{\partial W_0} = \frac{\partial \text{Loss}}{\partial p} \cdot \dots \cdot \frac{\partial h_{in}}{\partial W_0}$$

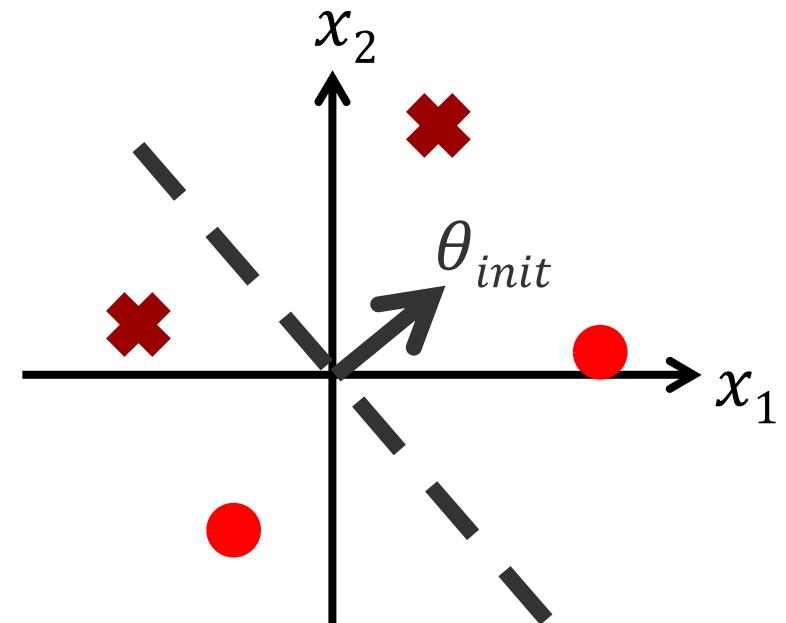
$$\frac{\partial \text{Loss}}{\partial b_0} = \frac{\partial \text{Loss}}{\partial p} \cdot \dots \cdot \frac{\partial h_{in}}{\partial b_0}$$

Landscape of training objective

Parameter space

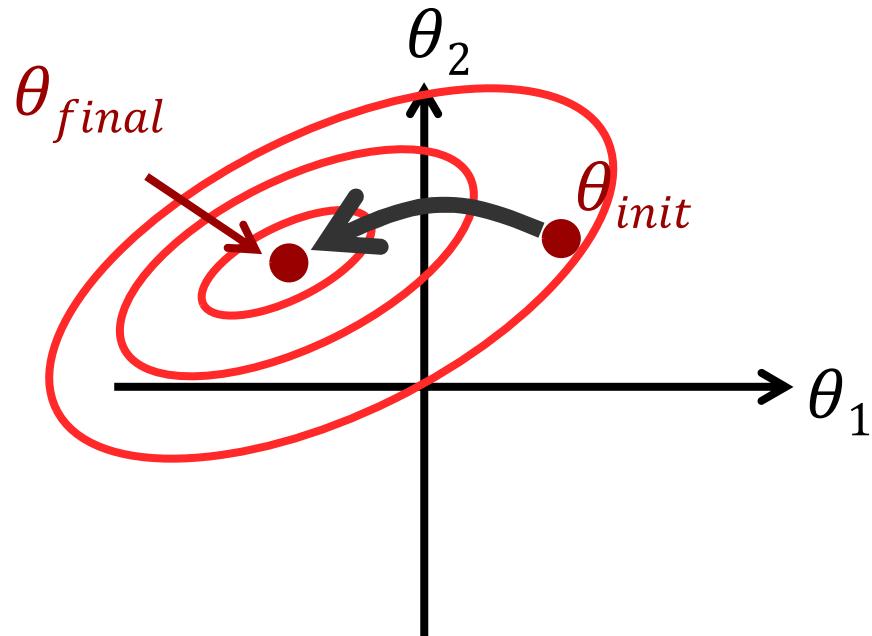


Example space

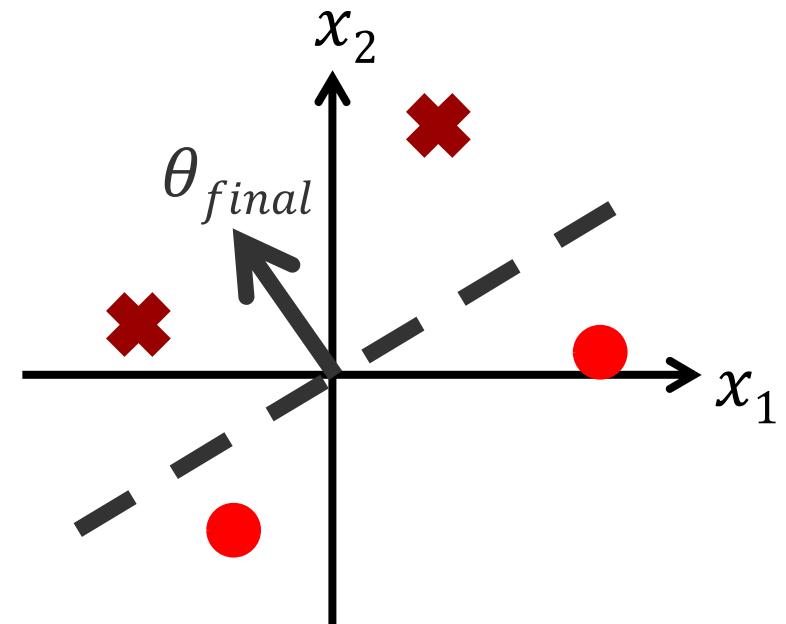


Landscape of training objective

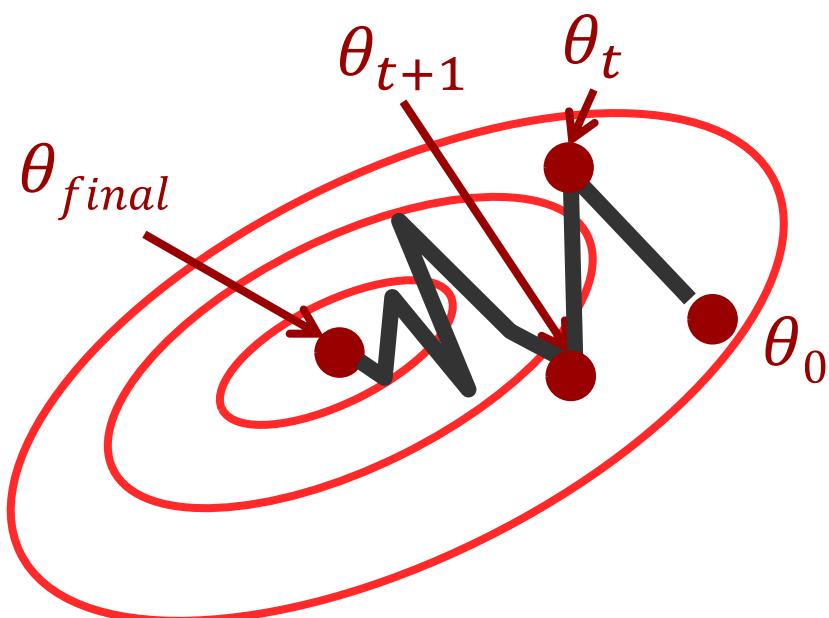
Parameter space



Example space



Stochastic gradient descent (SGD)



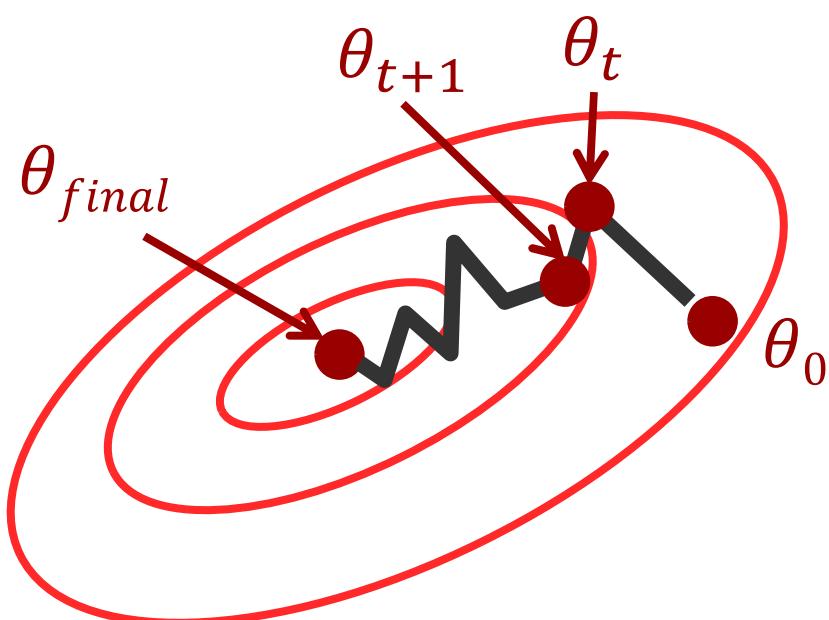
- ↗ Initialize θ_0 randomly
- ↗ For t in $0, \dots, T_{\text{maxiter}}$
$$\theta^{t+1} = \theta^t - \eta_t \cdot \nabla \text{Loss}(f_\theta(x_i), y_i)$$

Stochastic gradient

where index i is chosen randomly

- computation of $\nabla \text{Loss}(\dots)$ requires only one training example
- Per-iteration comp. cost = $O(1)$

Mini-batch stochastic gradient descent



↗ Initialize θ_0 randomly
↗ For t in $0, \dots, T_{\text{maxiter}}$
$$\theta^{t+1} = \theta^t - \eta_t \cdot \tilde{\nabla}_B L(\theta)$$

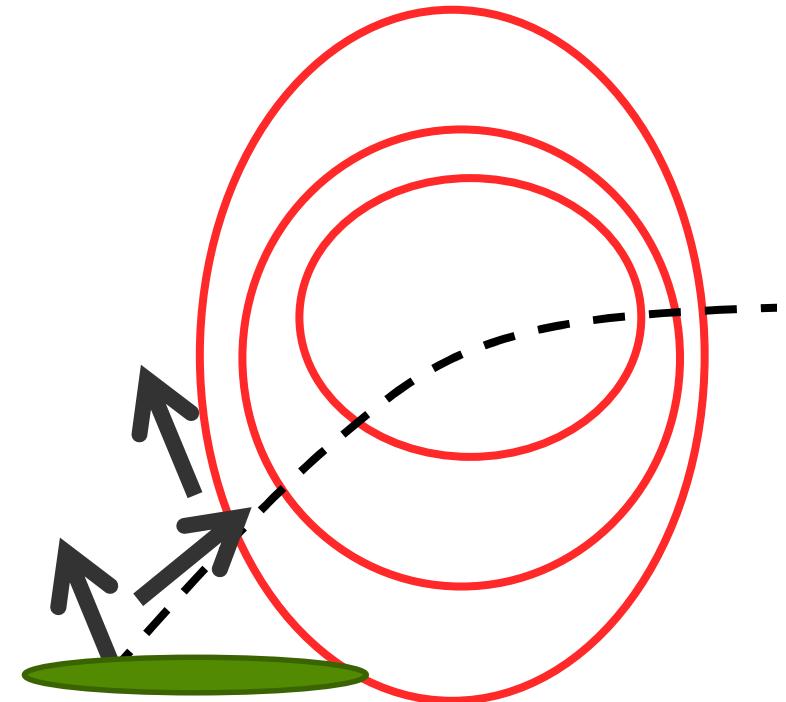
minibatch gradient

where minibatch B is chosen randomly

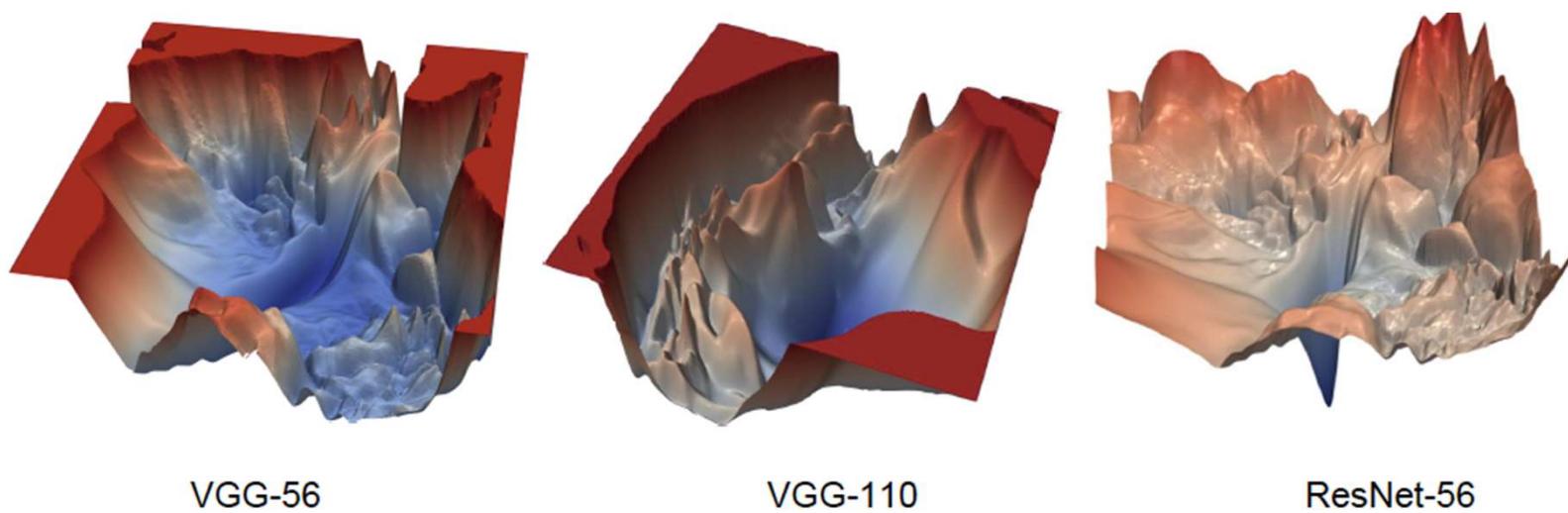
- $\tilde{\nabla}L(\theta)$ is average gradient over random subset of data of size B
- Per-iteration comp. cost = $O(B)$

More optimization algorithms

- ↗ **Momentum SGD**: improves SGD by incorporating “momentum”
- ↗ **Adam** [Kingma & Ba 2015]: uses first and second order statistics of the gradients so that gradients are normalized
- ↗ **Benefit**: prevents the vanishing/exploding gradient problem

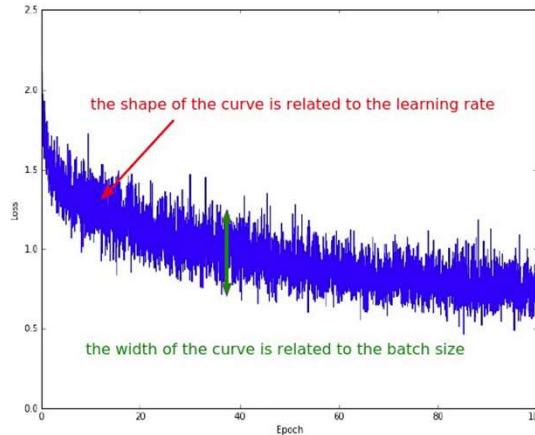
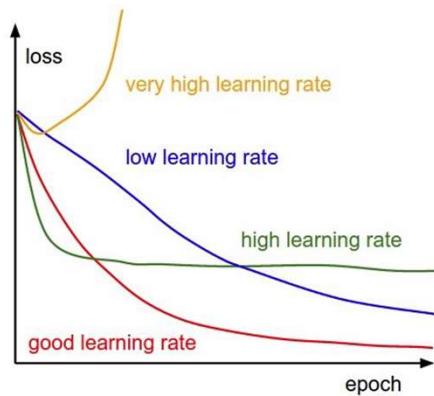


Stochastic Gradient Descent

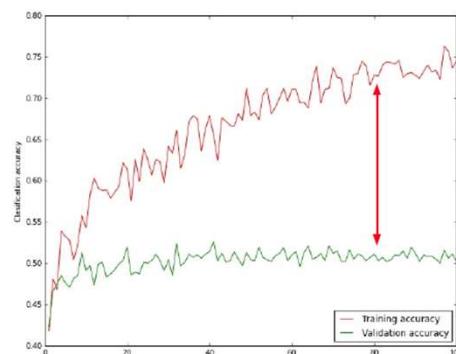
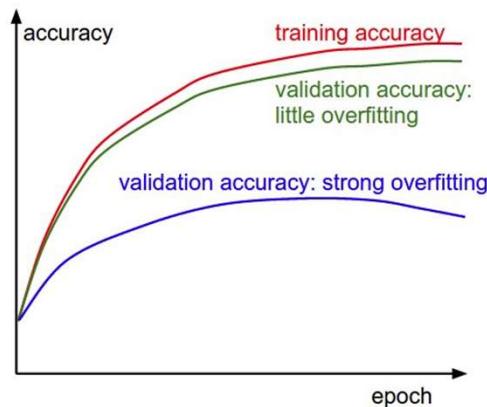


Hao Li et al., NIPS, 2017

Monitoring loss and accuracy



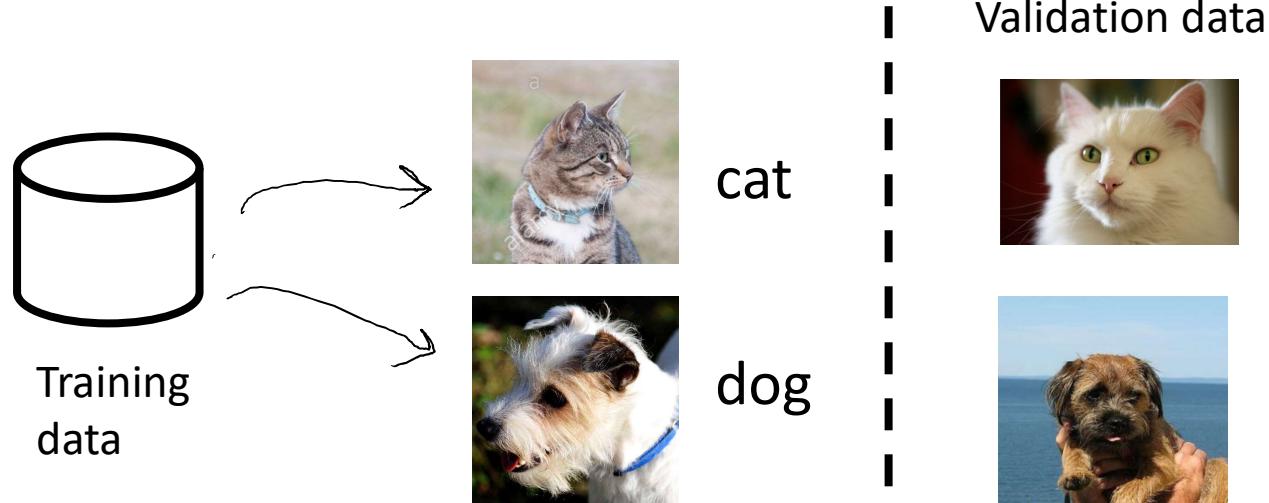
Looks linear?
Learning rate too low!
Decreases too slowly?
Learning rate too high.
Looks too noisy?
Increases the batch size.



Big gap?
- you're overfitting,
increase
regularization!

Overfitting – what is signal vs noise?

↗ Imagine:



- ↗ Powerful models are more likely to overfit
- ↗ We need validation data: leave out some portion of the training data to validate the generalizability of the model

Overfitting



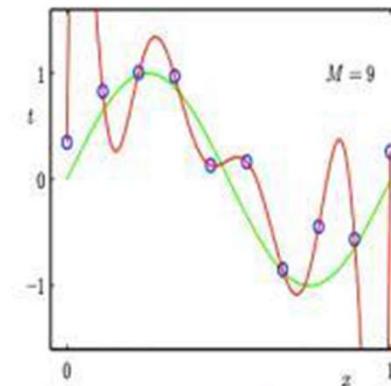
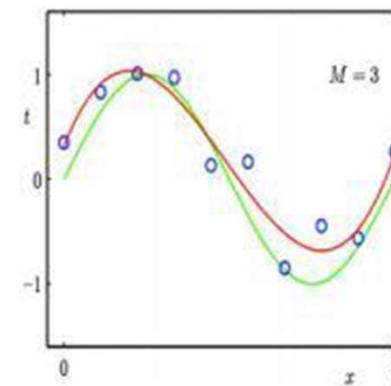
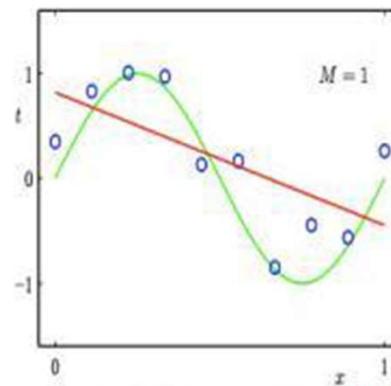
With four parameters I can fit an elephant, and with five I can make him wiggle his trunk.

— *John von Neumann* —

AZ QUOTES

Under- and Over-fitting

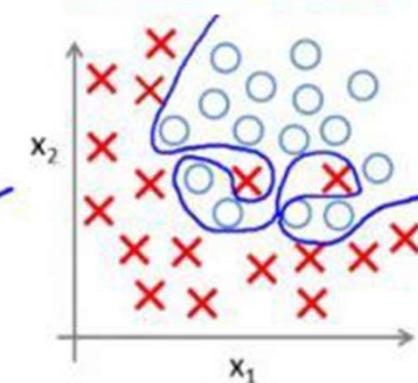
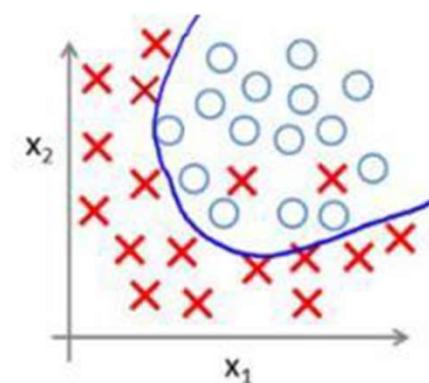
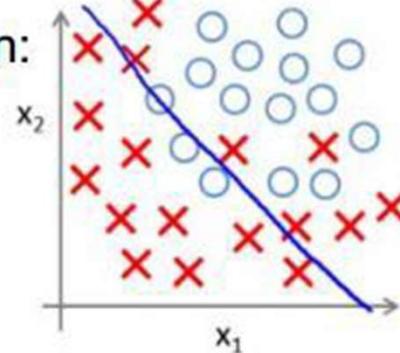
Regression:



predictor too inflexible:
cannot capture pattern

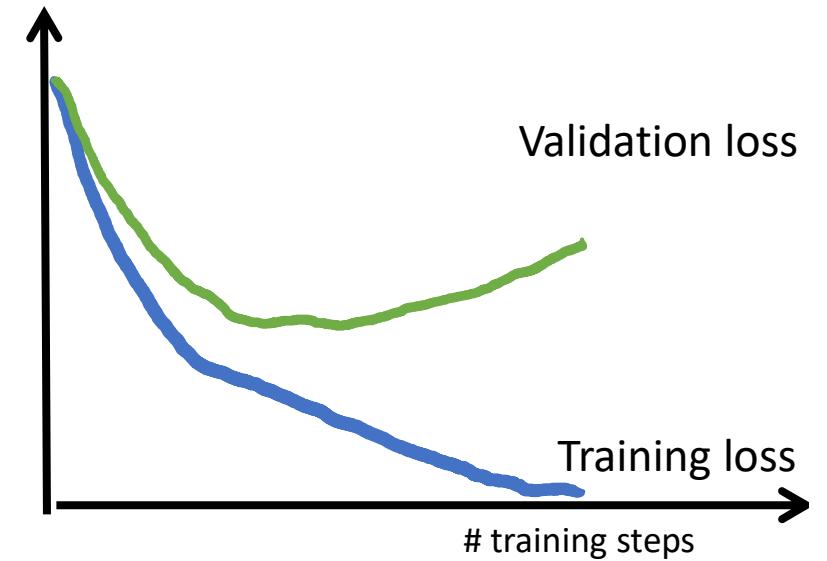
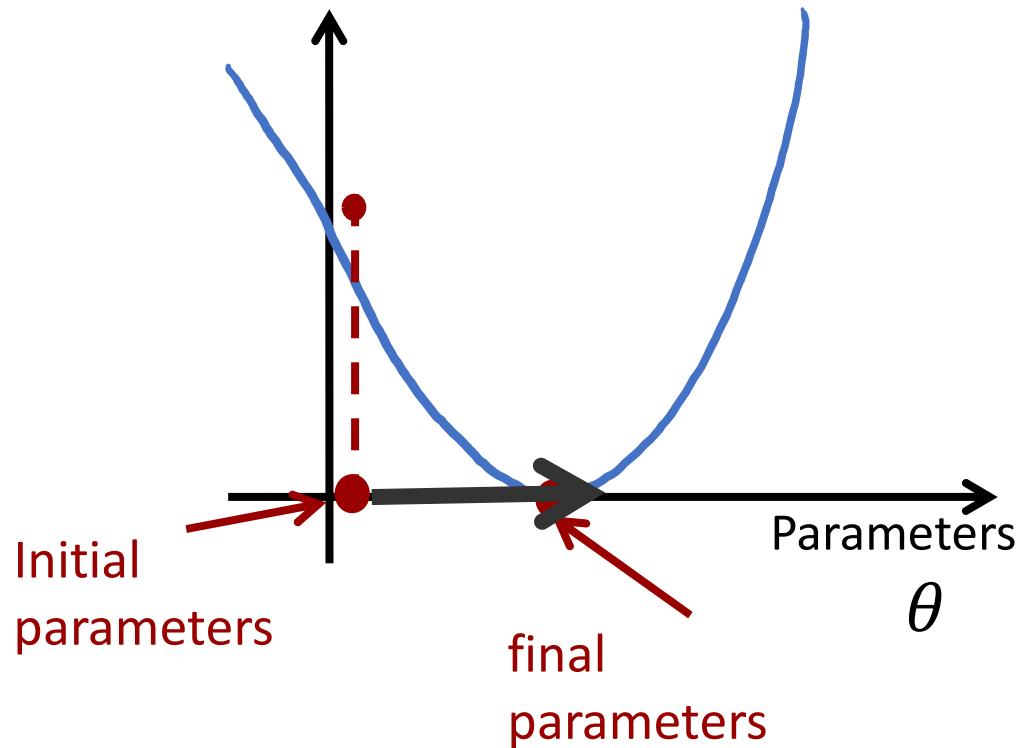
predictor too flexible:
fits noise in the data

Classification:



Landscape of training objective

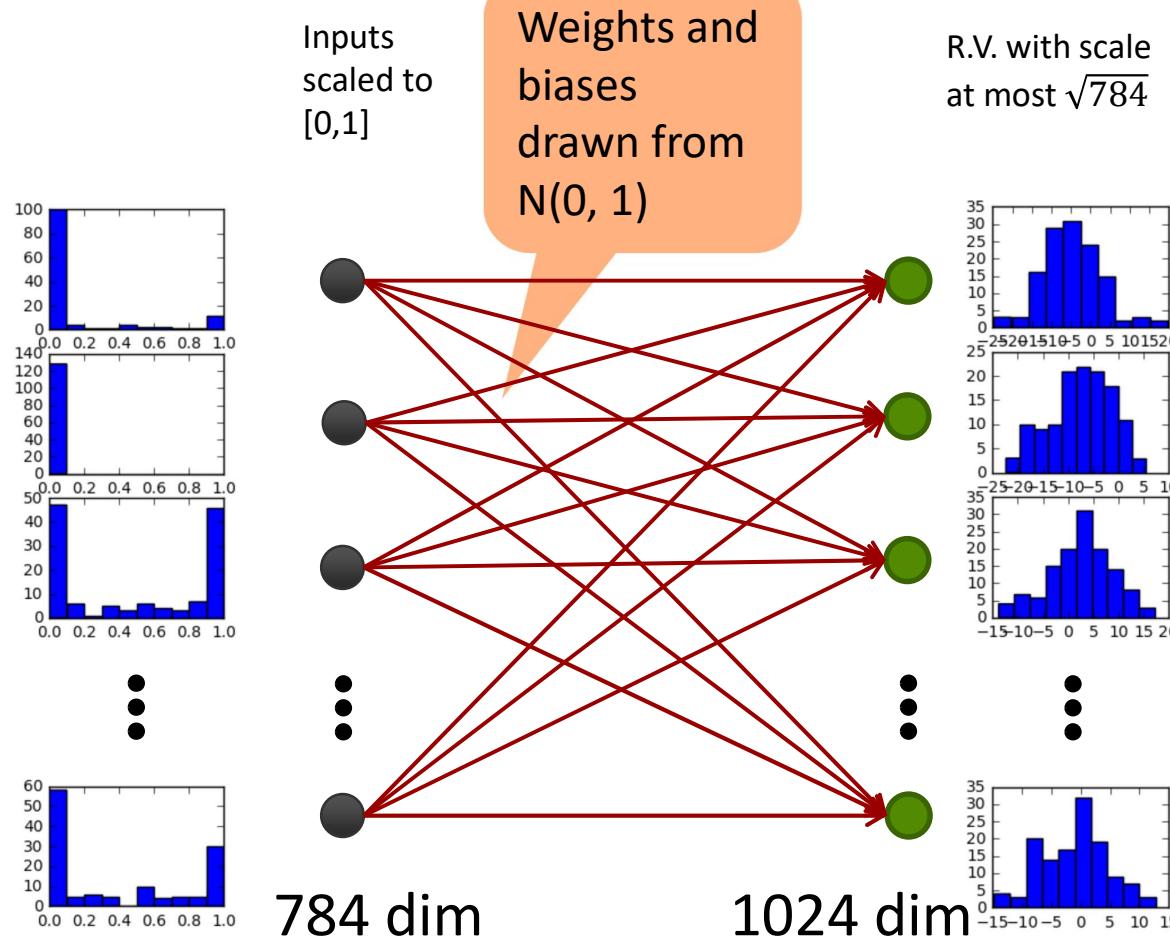
Objective function: $L(\theta) = \frac{1}{n} \sum_{i=1}^n Loss(f_\theta(x_i), y_i)$



Techniques to reduce overfitting

- ↗ Reduce the number of parameters
 - ↗ Parameter sharing (convnets, recurrent neural nets)
- ↗ Weight decay (aka L2 regularization)
 - ↗ Penalizes the magnitude of the parameters $\sum_{j=1}^d w_j^2$
- ↗ Early stopping
 - ↗ Indirectly controls the magnitude of the parameters
- ↗ Alternatives
 - ↗ Dropout, batch normalization

Batch normalization



Against overfitting: data augmentation

- Resizing images keeping the aspect ratio.
- Enhancing images using random distortions (color, contrast, brightness and sharpness).
- Applying random crops using the same dimension for the width and height.
- Applying random horizontal flips.



- ↗ AI, Machine learning & Deep learning
- ↗ What is a Convolutional Neural Network?
 - ↗ Layers
 - ↗ Optimization
- ↗ Applications

Everybody dances now



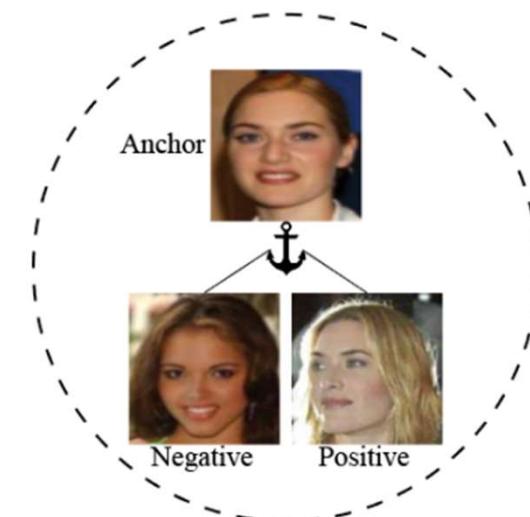
[More links](#)

17:38

Convolutional NN

CNN beat humans in many tasks as:

- ↗ object recognition,
- ↗ lip reading,
- ↗ high-end surveillance,
- ↗ facial recognition,
- ↗ object-based searches,
- ↗ license plate readers,
- ↗ traffic violations detection,
- ↗ breast tomosynthesis diagnosis,
- ↗ etc., etc.



Computer Vision & Deep Learning ↗



Social media



Autodriving (Tesla)



Security (Airports)



Shopping (Mango, Amazon)

Biometric data for cheating detection

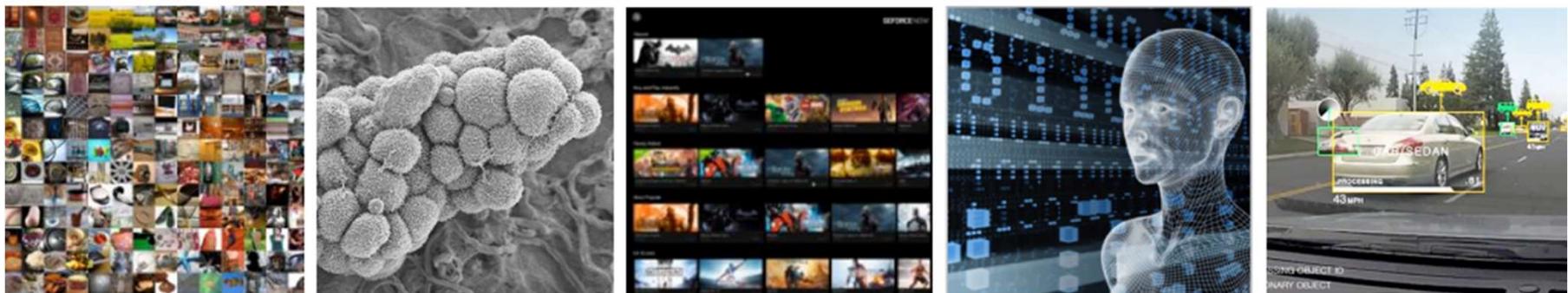


Analizes up to 38 faciales
facial micro-gestures, digital
prints and scans hand veins

iBorderCtrl will reduce the cost of
frontiers access of near 700M
crossing European countries



Deep learning everywhere



INTERNET & CLOUD

Image Classification
Speech Recognition
Language Translation
Language Processing
Sentiment Analysis
Recommendation

MEDICINE & BIOLOGY

Cancer Cell Detection
Diabetic Grading
Drug Discovery

MEDIA & ENTERTAINMENT

Video Captioning
Video Search
Real Time Translation

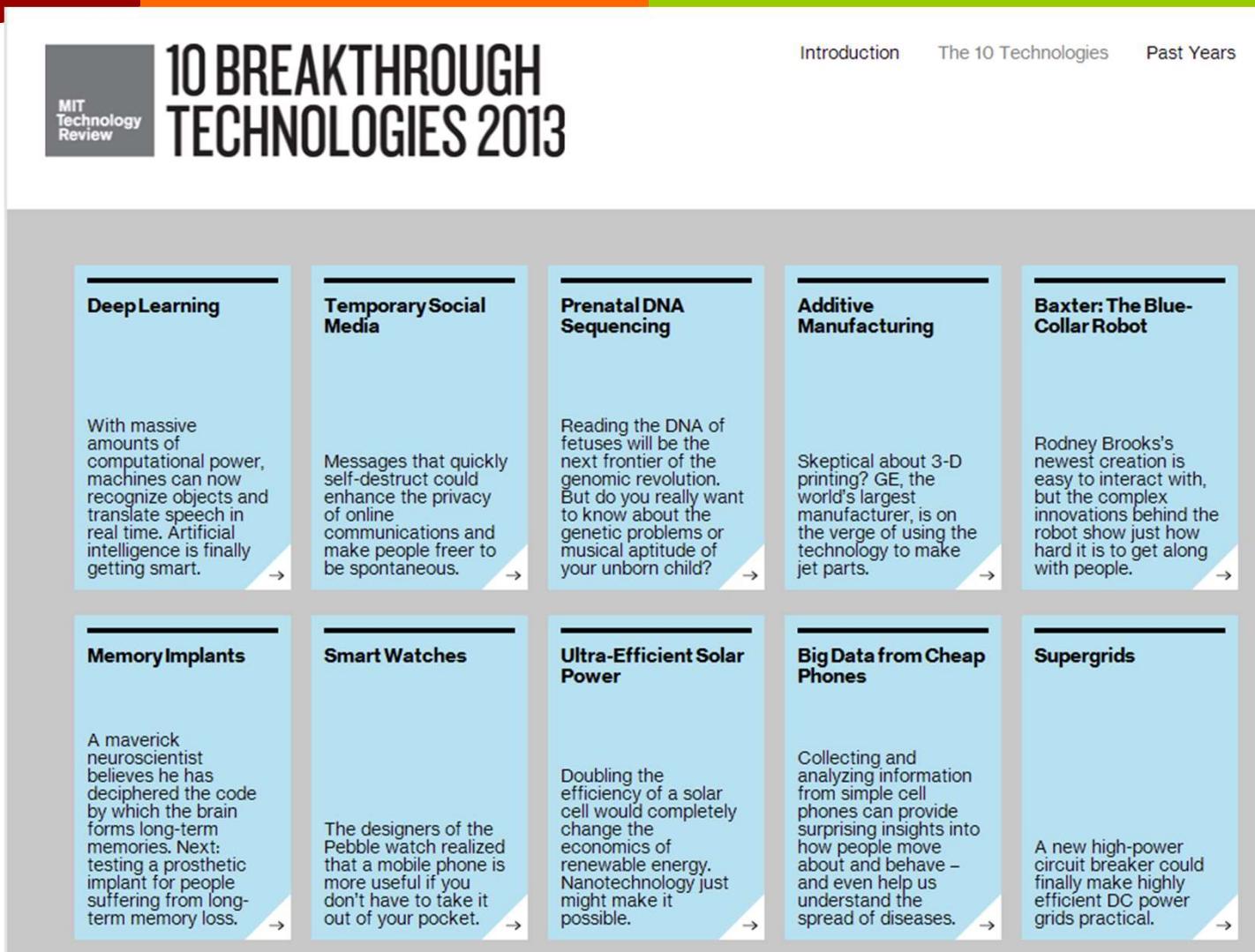
SECURITY & DEFENSE

Face Detection
Video Surveillance
Satellite Imagery

AUTONOMOUS MACHINES

Pedestrian Detection
Lane Tracking
Recognize Traffic Sign

Deep learning - one of the 10 breakthrough technologies



The screenshot shows the MIT Technology Review website for the "10 Breakthrough Technologies 2013" list. The header features the MIT Technology Review logo and the title "10 BREAKTHROUGH TECHNOLOGIES 2013". Below the header, there are ten cards, each representing a technology. The technologies listed are: Deep Learning, Temporary Social Media, Prenatal DNA Sequencing, Additive Manufacturing, Baxter: The Blue-Collar Robot, Memory Implants, Smart Watches, Ultra-Efficient Solar Power, Big Data from Cheap Phones, and Supergrids. Each card contains a brief description and a right-pointing arrow.

- Deep Learning**
With massive amounts of computational power, machines can now recognize objects and translate speech in real time. Artificial intelligence is finally getting smart. →
- Temporary Social Media**
Messages that quickly self-destruct could enhance the privacy of online communications and make people freer to be spontaneous. →
- Prenatal DNA Sequencing**
Reading the DNA of fetuses will be the next frontier of the genomic revolution. But do you really want to know about the genetic problems or musical aptitude of your unborn child? →
- Additive Manufacturing**
Skeptical about 3-D printing? GE, the world's largest manufacturer, is on the verge of using the technology to make jet parts. →
- Baxter: The Blue-Collar Robot**
Rodney Brooks's newest creation is easy to interact with, but the complex innovations behind the robot show just how hard it is to get along with people. →
- Memory Implants**
A maverick neuroscientist believes he has deciphered the code by which the brain forms long-term memories. Next: testing a prosthetic implant for people suffering from long-term memory loss. →
- Smart Watches**
The designers of the Pebble watch realized that a mobile phone is more useful if you don't have to take it out of your pocket. →
- Ultra-Efficient Solar Power**
Doubling the efficiency of a solar cell would completely change the economics of renewable energy. Nanotechnology just might make it possible. →
- Big Data from Cheap Phones**
Collecting and analyzing information from simple cell phones can provide surprising insights into how people move about and behave – and even help us understand the spread of diseases. →
- Supergrids**
A new high-power circuit breaker could finally make highly efficient DC power grids practical. →

Deep learning everywhere

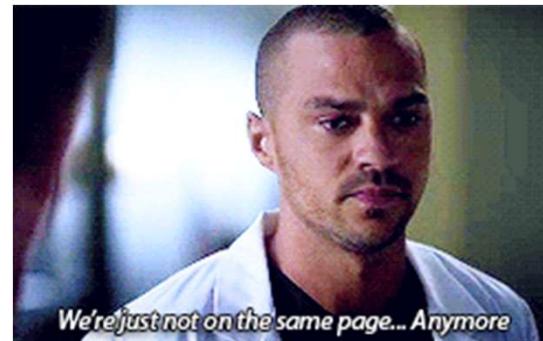
[“Preparing for the Future of Artificial Intelligence.”](#) ★¹



This 58-page report outlines a number of important topics related to artificial intelligence.

Concerns about DL

- Different development approach.



We're just not on the same page... Anymore

- No clear view on how insight is generated.



- A system is only as good as the data it learns from.

Conclusions



Deep learning – a technology that came to stay

A new technological trend that is affecting directly our environment

CNN applicable to different Computer Vision problems

Different CNN models can be found. No optimal one exists.