

## Using Deep Neural Networks for Discriminative Feature Localization

Javier Sánchez Rois  
Daniel González Jiménez



Àgata Lapedriza García



**Javier Sánchez Rois**  
Research Software Engineer



**Daniel González**  
R&D Director, MMI Area



**SUnAI** SCENE UNDERSTANDING  
& ARTIFICIAL INTELLIGENCE  
Computer Vision & AI Recognition for Objects and Gestures



**Àgata Lapedriza**  
Associate Professor (UOC)  
Visiting Researcher (MIT)



**Àgata Lapedriza**  
Associate Professor

# Deep {Neural Nets, Learning}

43 y.  
1.70 m  
Iowa  
woman

$$\begin{bmatrix} 43 \\ 170 \\ 29 \\ 1 \end{bmatrix}$$


**Deep Learning**

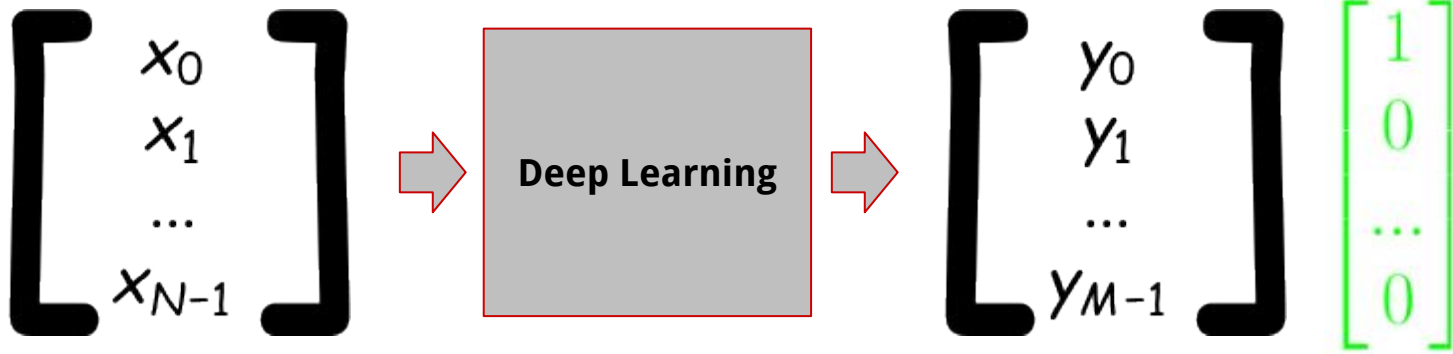

$$\begin{bmatrix} 0 \\ 1 \end{bmatrix}$$


26 y.  
1.68 m  
Louisiana  
man

$$\begin{bmatrix} 26 \\ 168 \\ 18 \\ 0 \end{bmatrix}$$


**Deep Learning**


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



# A Linear Model

$$\mathbf{w} \cdot \begin{bmatrix} x_0 \\ x_1 \\ \dots \\ x_{N-1} \end{bmatrix} + \mathbf{b} = \begin{bmatrix} y_0 \\ y_1 \\ \dots \\ y_{M-1} \end{bmatrix}$$

$$\mathbf{w} \cdot \begin{bmatrix} x_0 \\ x_1 \\ \dots \\ x_{N-1} \end{bmatrix} + \mathbf{b} = \begin{bmatrix} 13.2 \\ -1.3 \\ \dots \\ 0.01 \end{bmatrix} \quad (\text{raw scores})$$



$$\text{Softmax} \left( \mathbf{w} \cdot \begin{bmatrix} x_0 \\ x_1 \\ \dots \\ x_{N-1} \end{bmatrix} + \mathbf{b} \right) = \begin{bmatrix} 1 \\ 0 \\ \dots \\ 0 \end{bmatrix} \quad (\text{one-hot encoding})$$

$$\text{Softmax} \left( \mathbf{W} \cdot \begin{bmatrix} x_0 \\ x_1 \\ \dots \\ x_{N-1} \end{bmatrix} + \mathbf{b} \right) = \begin{bmatrix} 0.9 \\ \sim 0 \\ \dots \\ \sim 0 \end{bmatrix} \quad (... \text{ almost})$$

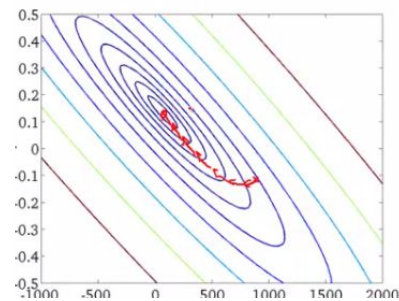
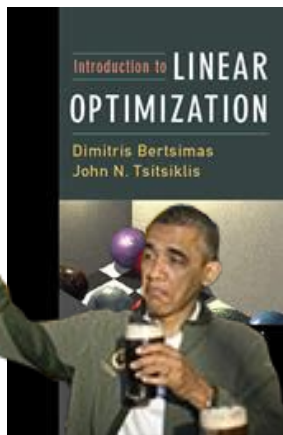
***Logistic Regression***

$$\text{Softmax} \left( \underset{\uparrow}{\mathbf{W}} \cdot \begin{bmatrix} x_0 \\ x_1 \\ \dots \\ x_{N-1} \end{bmatrix} + \underset{\uparrow}{\mathbf{b}} \right)$$

$$Error_i = f\left(\begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 0.3 \\ 0.6 \\ 0.05 \\ 0.05 \end{bmatrix}\right)$$

$$E_i = CE(\mathbf{l}_i, \text{Softmax}(\mathbf{W} \cdot \mathbf{x}_i + \mathbf{b}))$$

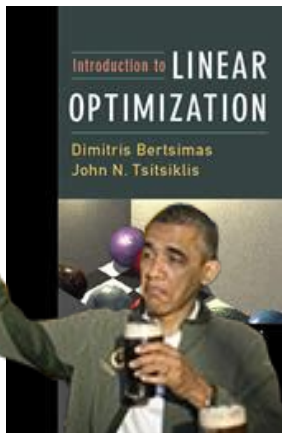
$$\min_{W,b} \sum_i E_i \quad ?$$



### Gradient Descent

$$W_{i+1} \leftarrow W_i - \alpha \frac{\partial E}{\partial W}$$

# Linear Models: Pros

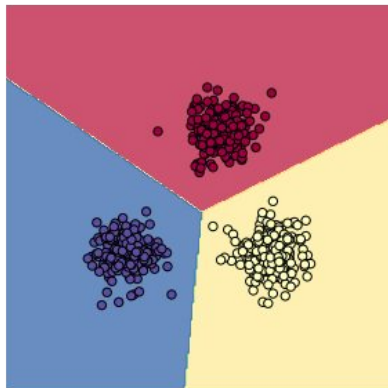


Easy to solve for minimal loss

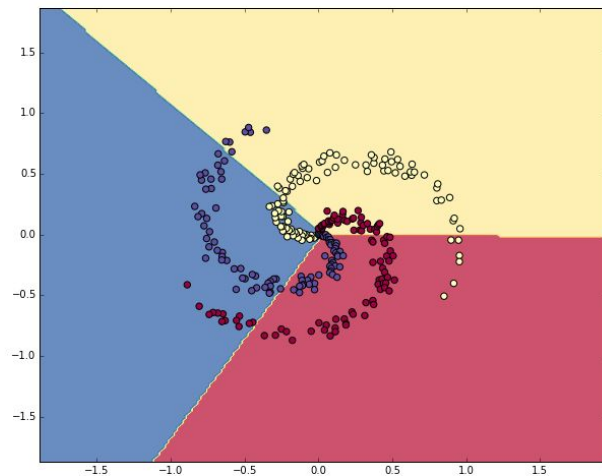
GPUs for optimal implementation



# Linear Models: Cons



OK for easy scenarios



Underfitting complex data

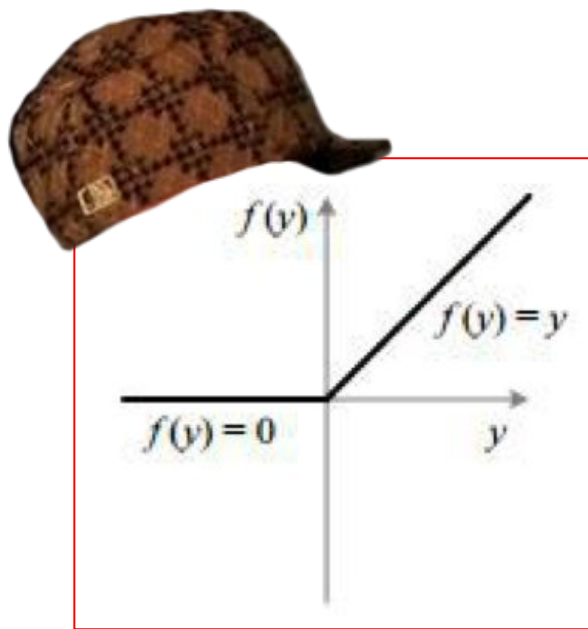




We need to go *nonlinear*

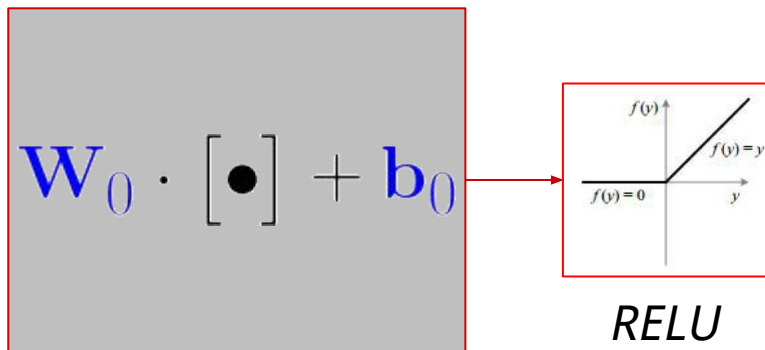


# RELU: rectified linear unit



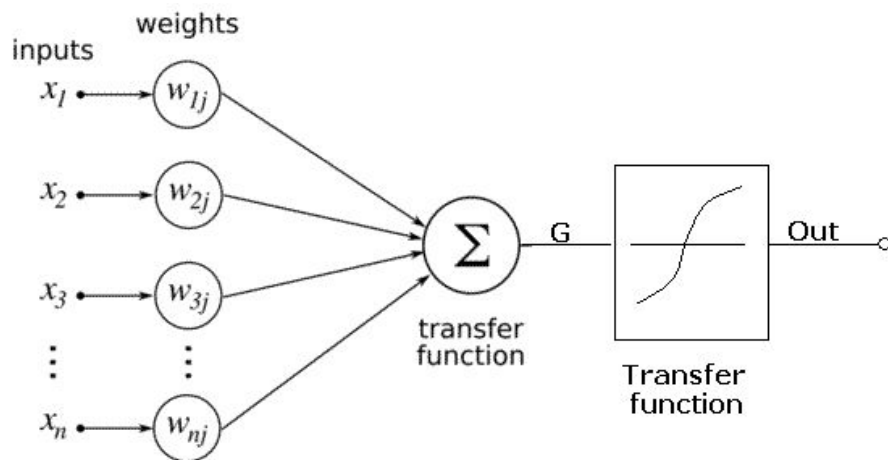
*(scumbag nonlinear function)*

# The *perceptron*

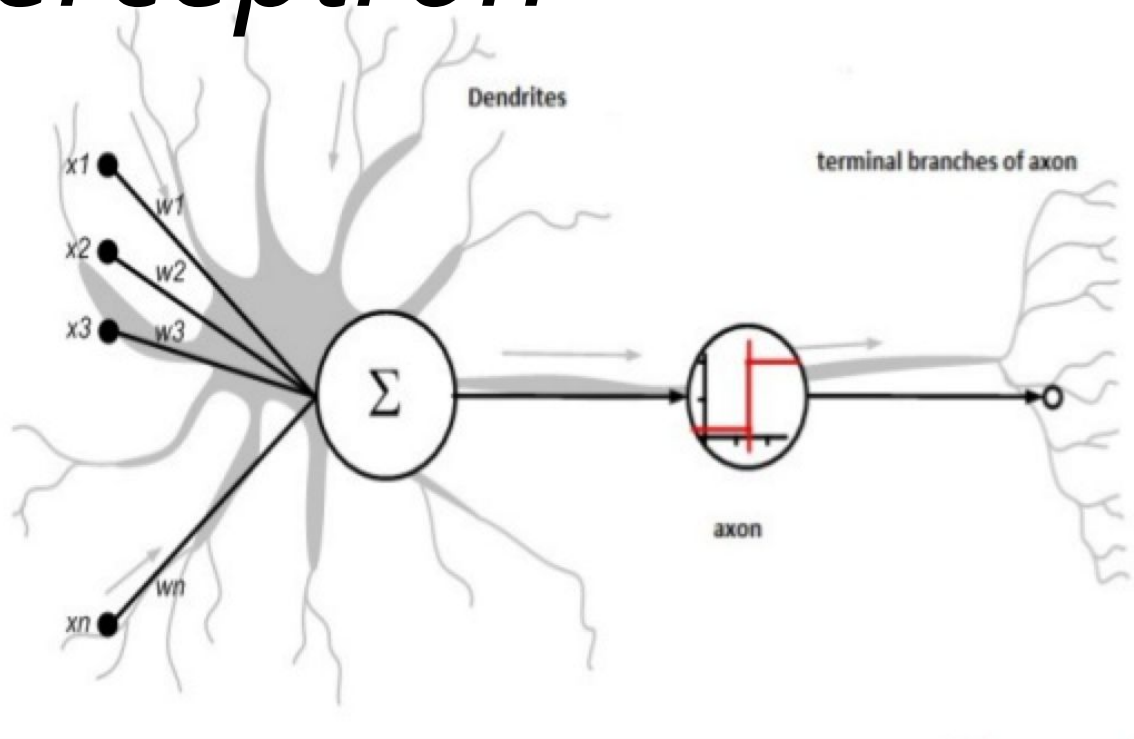


***Fully-connected layer***

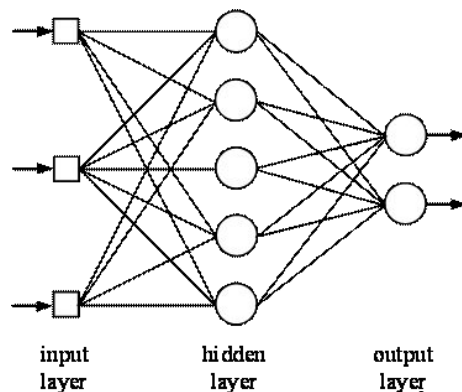
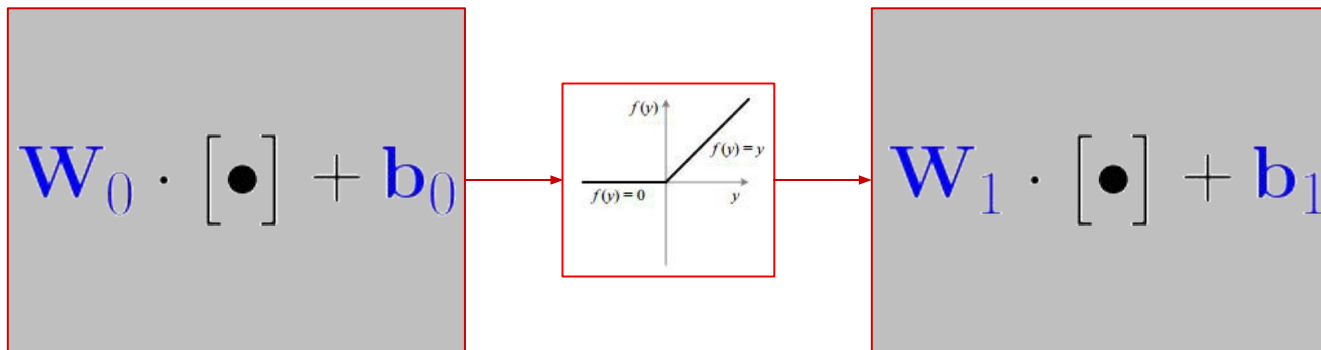
# The *perceptron*



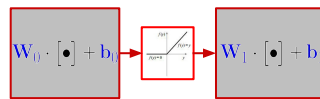
# The *perceptron*



# The *multilayer perceptron*



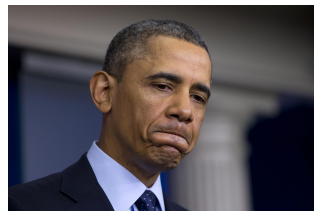
$$\min_{W_0, b_0, W_1, b_1} \sum_i E_i$$



too hard

**Gradient Descent**

$$W_{i+1} \leftarrow W_i - \alpha \frac{\partial E}{\partial W}$$





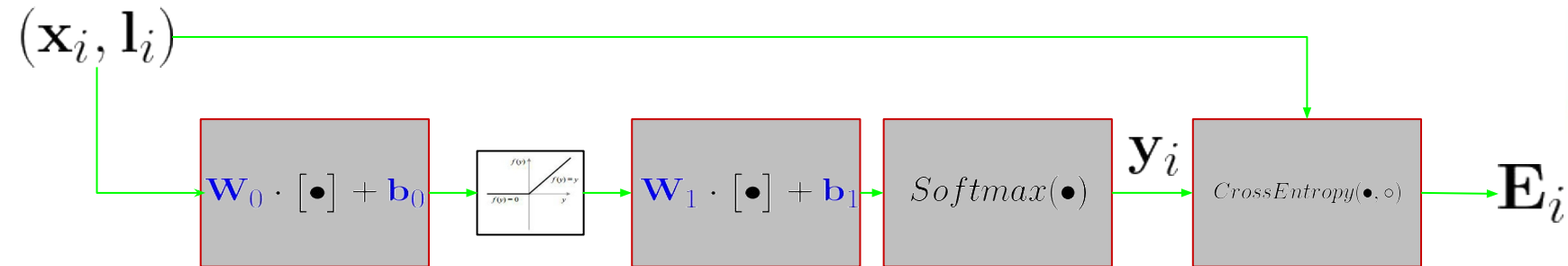


# Backpropagation\*

Backpropagation = Rule of Chain

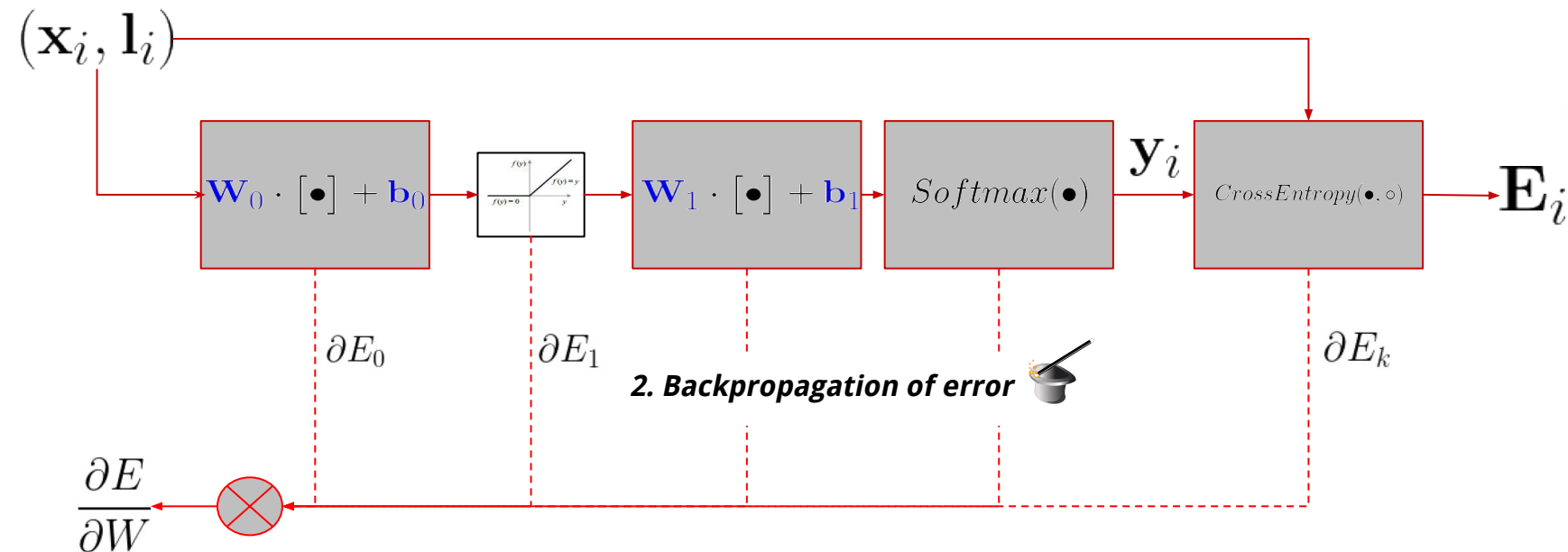
$$\frac{\partial}{\partial W} Error \left( \boxed{W_0 \cdot [\bullet] + b_0} \rightarrow \boxed{\frac{\exp}{1 + \exp}} \rightarrow \boxed{W_1 \cdot [\bullet] + b_1} \right) = \partial E_0 \cdot \partial E_1 \dots \cdot \partial E_k$$

# Training

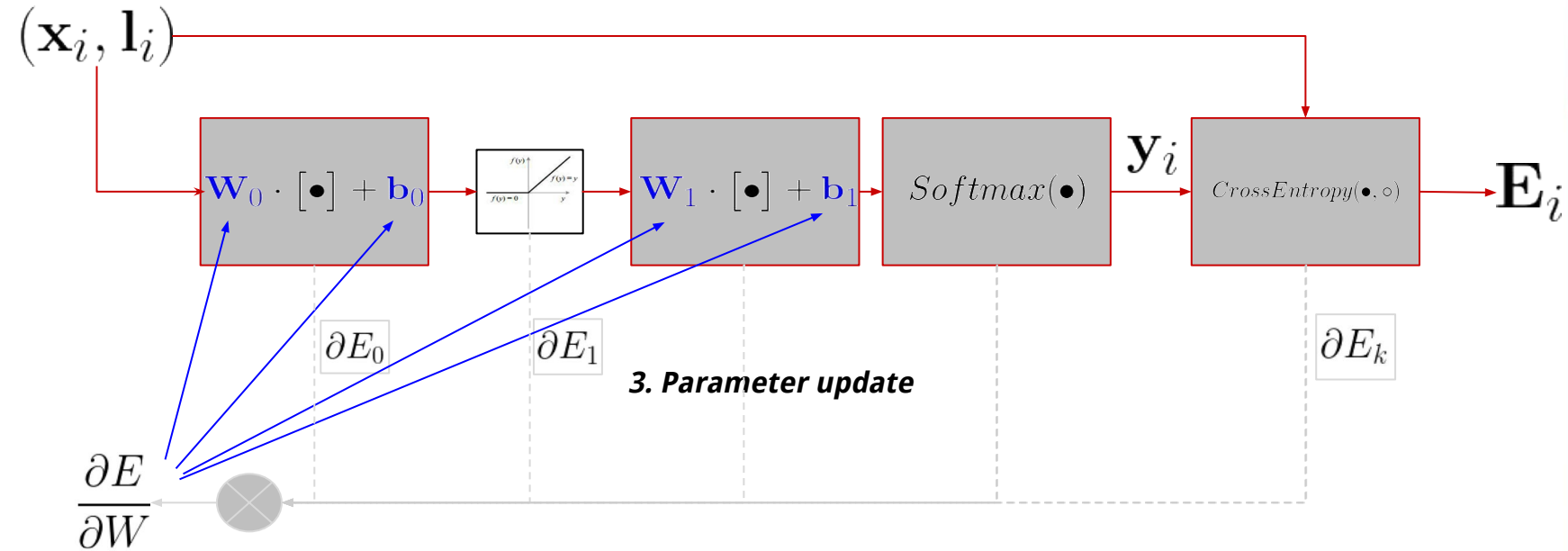


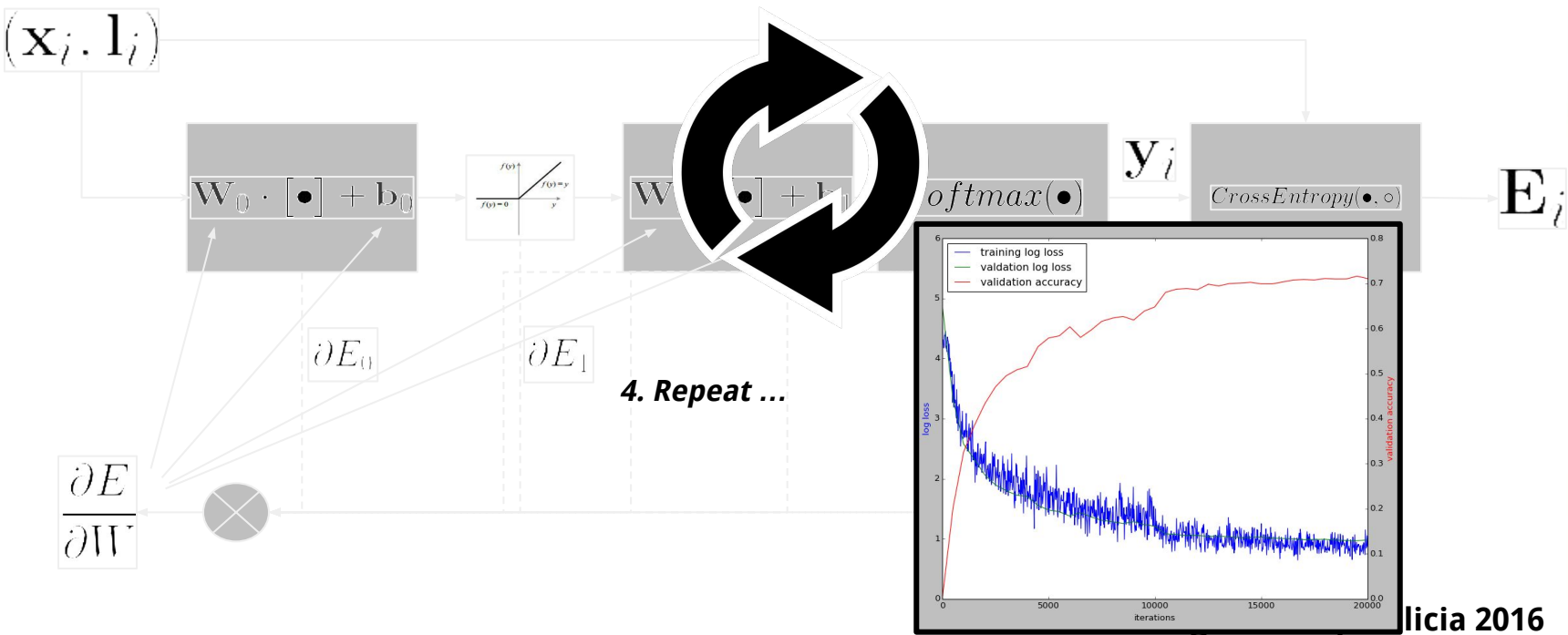
## 1. Forward Propagation

# Training

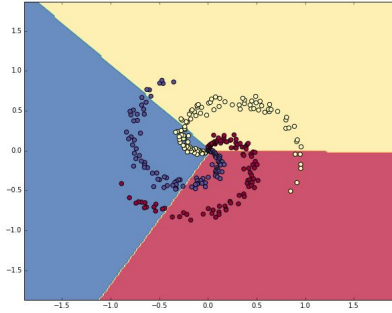


# Training

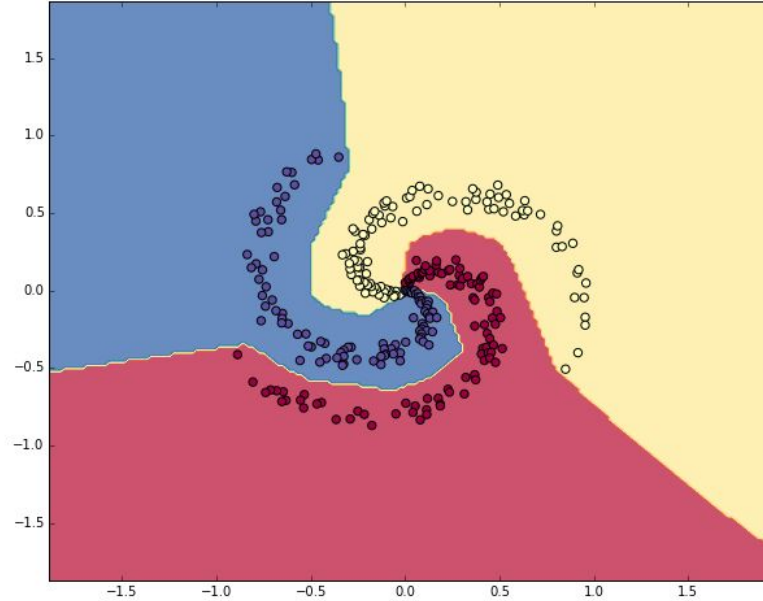




$$W_0 \cdot \begin{bmatrix} \bullet \\ \bullet \end{bmatrix} + b_0$$



$$W_0 \cdot \begin{bmatrix} \bullet \\ \bullet \end{bmatrix} + b_0 \xrightarrow{\frac{\max(0, x)}{x}} W_1 \cdot \begin{bmatrix} \bullet \\ \bullet \end{bmatrix} + b_1$$



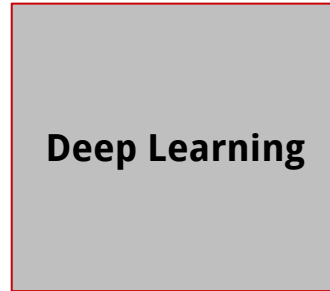
*Andrej Karpathy - CS231n: Convolutional Neural Networks for Visual Recognition.*

**Machine Learning Workshop Galicia 2016**

**Neural Nets are machine learning network characterized by**

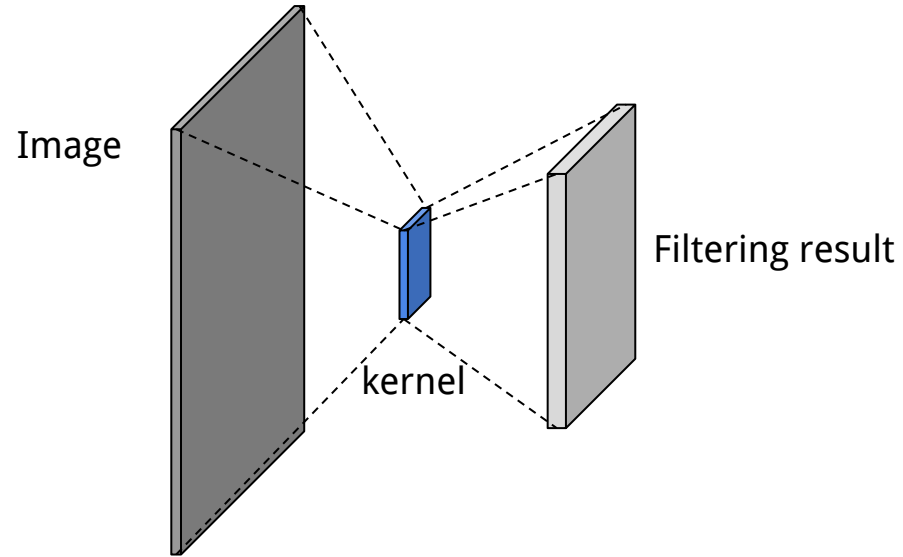
- Use of **nonlinear** models
- **Backpropagation** algorithm for training

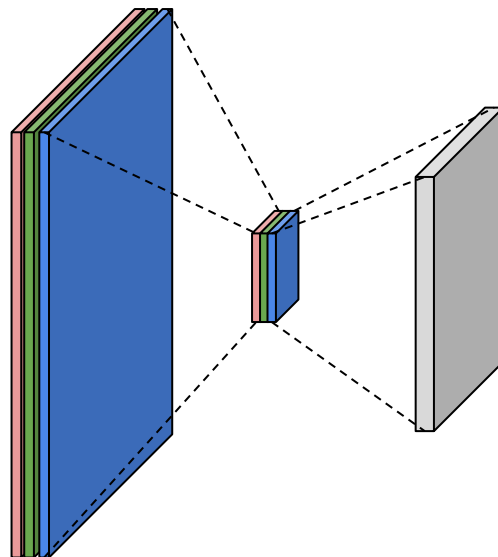
# DNN'ing images

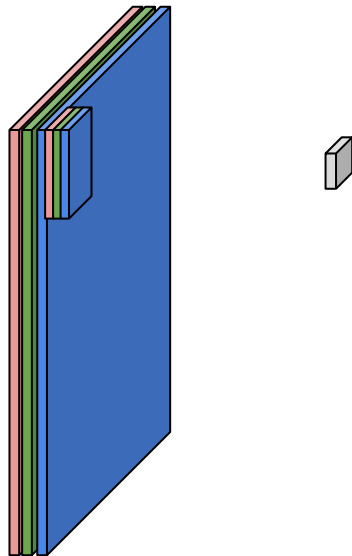


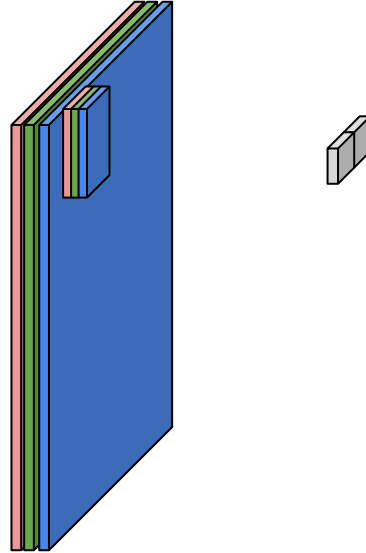
$$\begin{bmatrix} y_0 \\ y_1 \\ \dots \\ y_{M-1} \end{bmatrix}$$

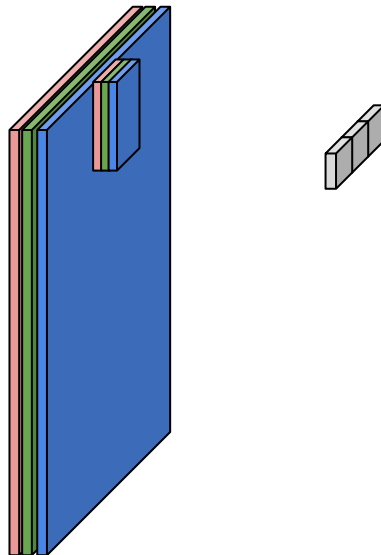


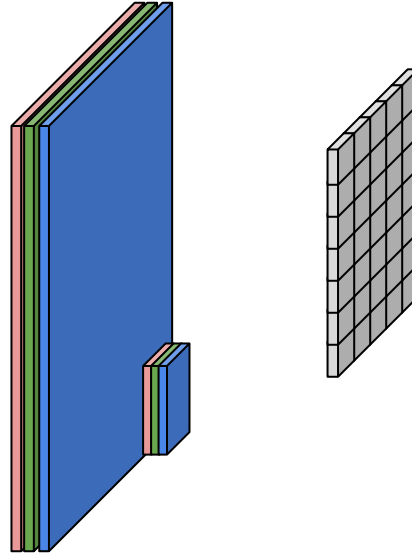


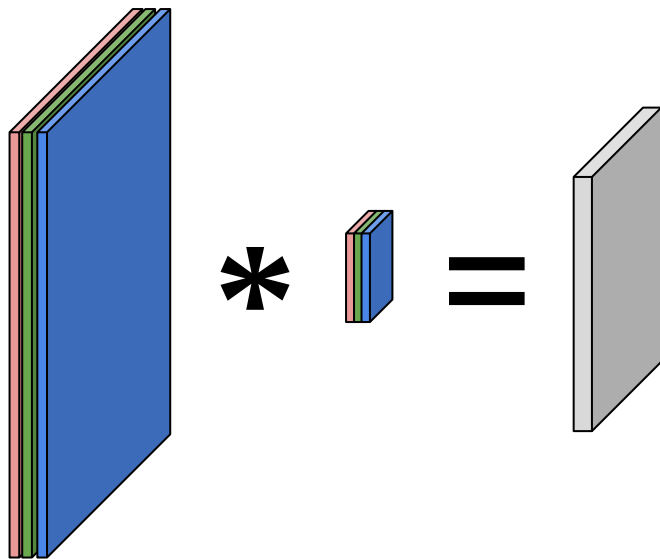


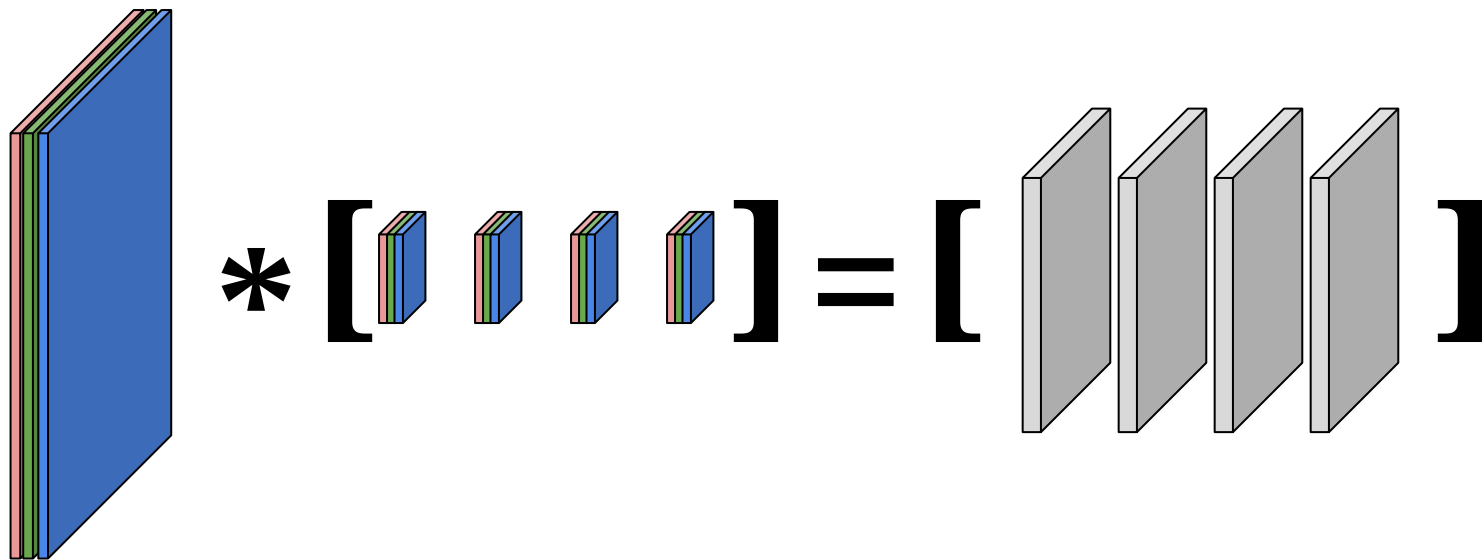




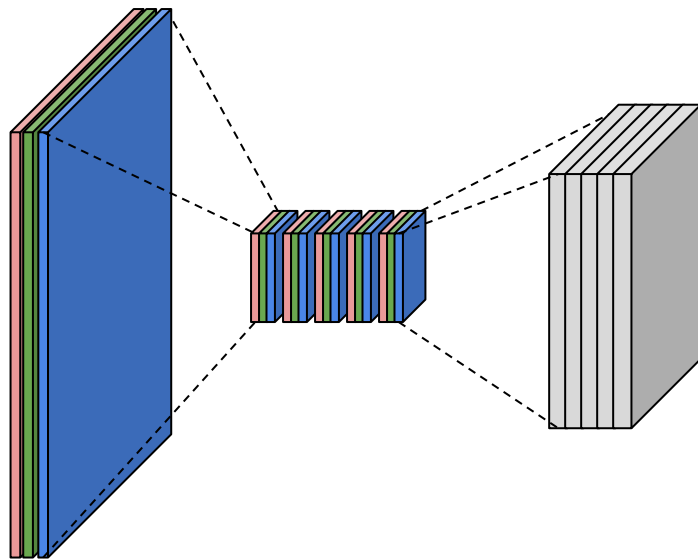


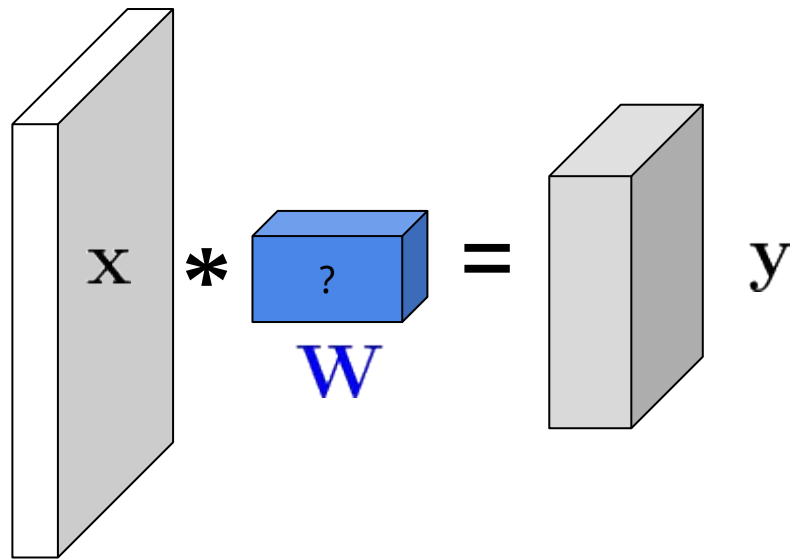


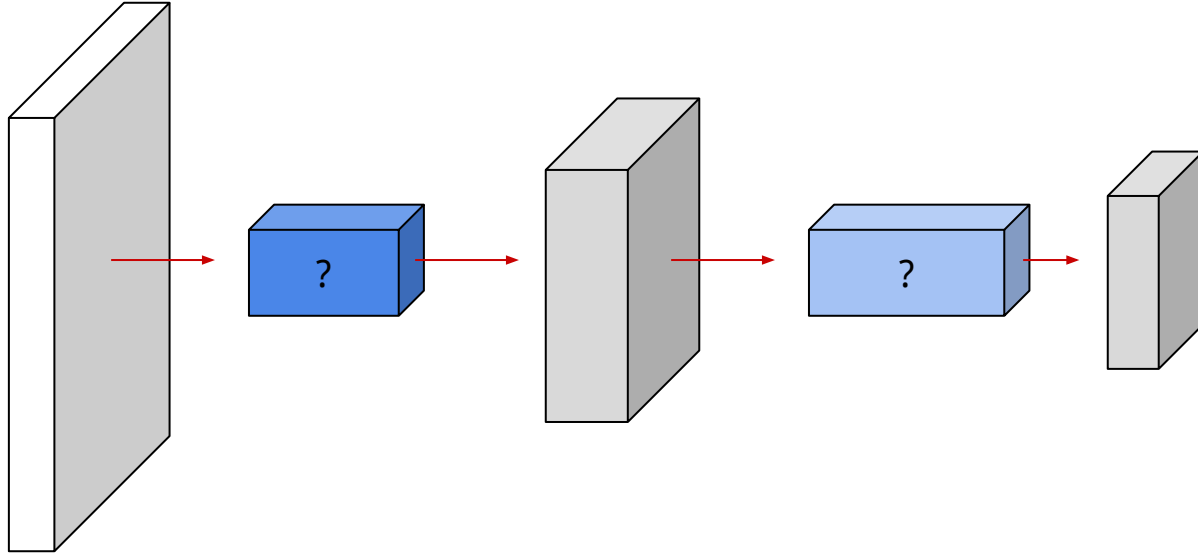


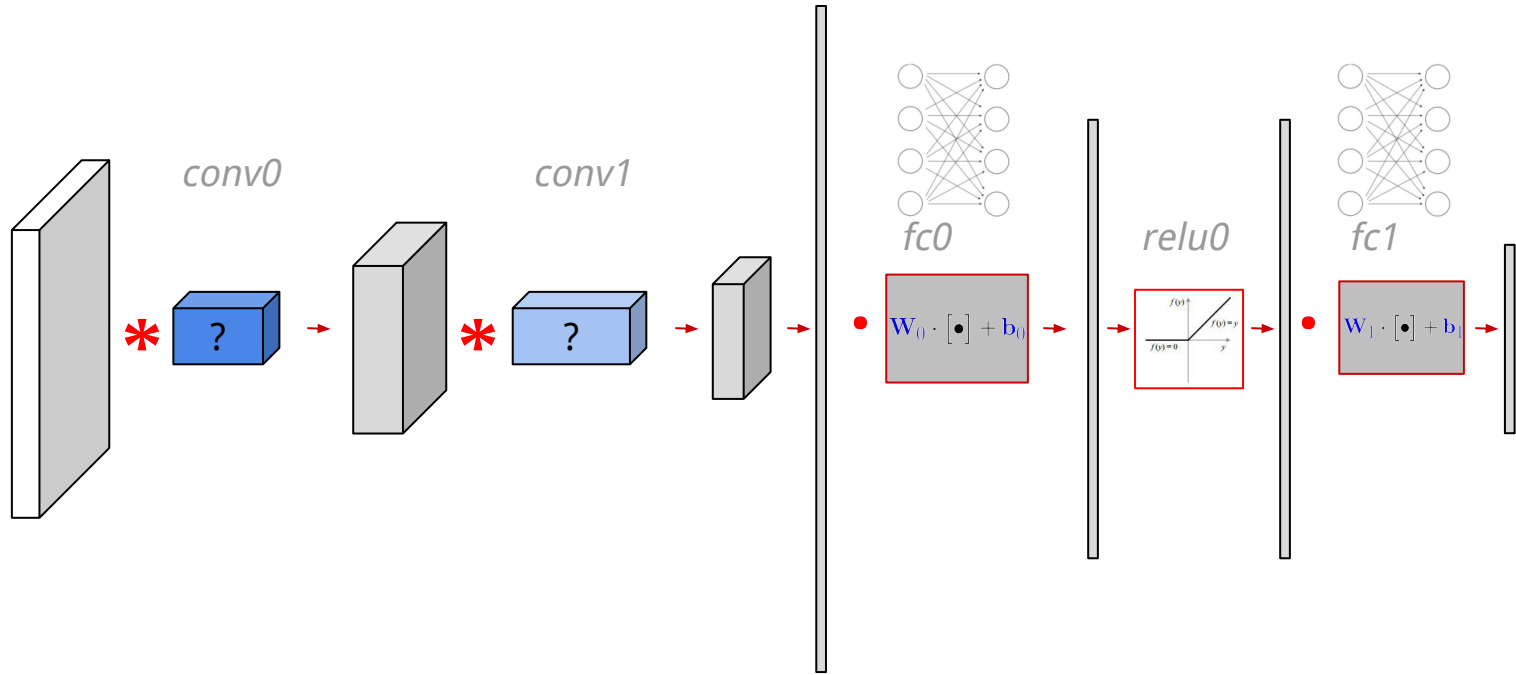












# Network Architecture

How many convolutional layers?  
How many fully connected layers?  
Parameters?  
...

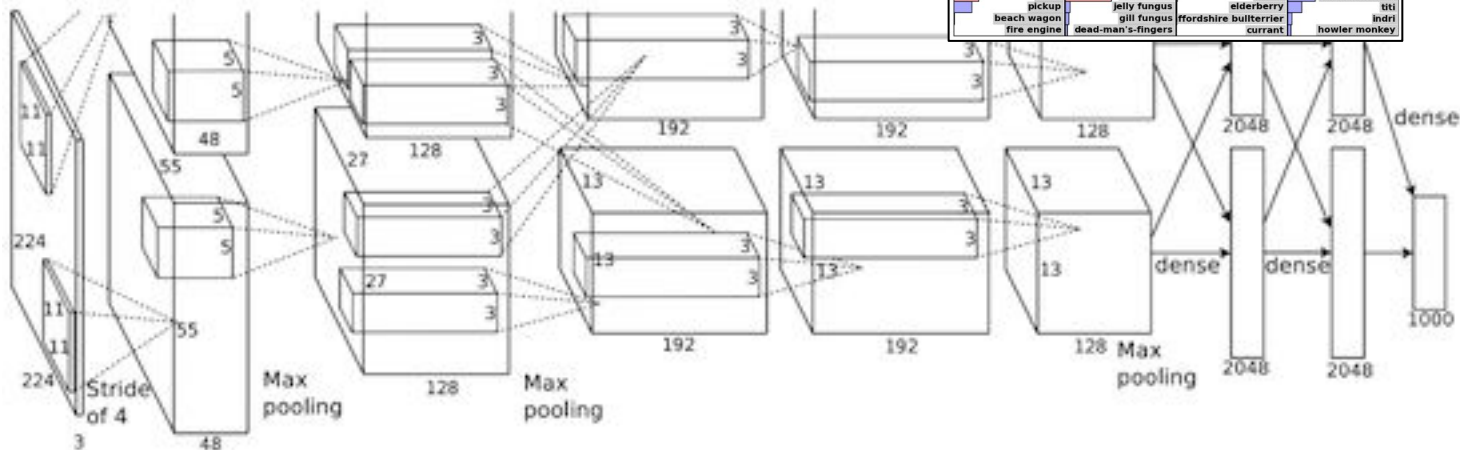
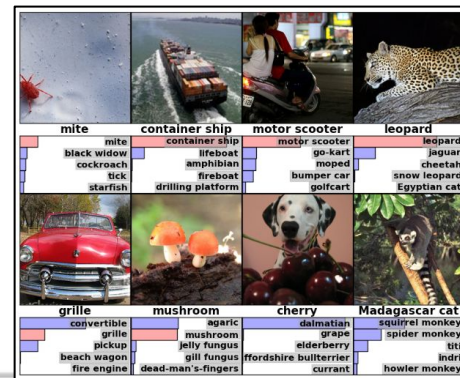
# CNN architectures: AlexNet (2012)

## ImageNet Classification with Deep Convolutional Neural Networks

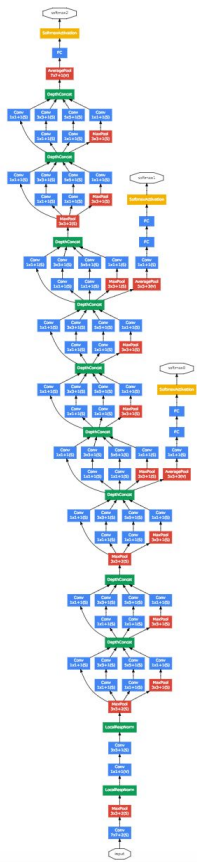
Alex Krizhevsky  
University of Toronto  
kriz@cs.utoronto.ca

Ilya Sutskever  
University of Toronto  
ilya@cs.utoronto.ca

Geoffrey E. Hinton  
University of Toronto  
hinton@cs.utoronto.ca



# CNN architectures: GoogLeNet (2014)



## Going Deeper with Convolutions

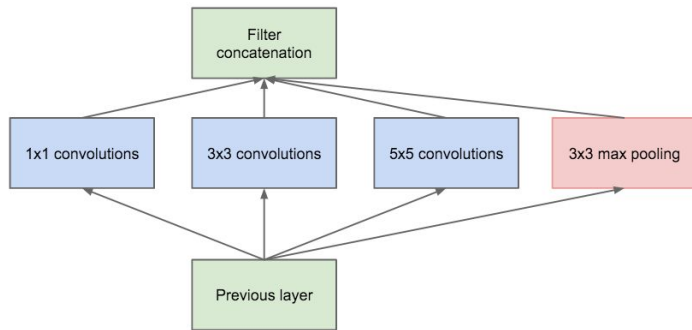
Christian Szegedy<sup>1</sup>, Wei Liu<sup>2</sup>, Yangqing Jia<sup>1</sup>, Pierre Sermanet<sup>1</sup>, Scott Reed<sup>3</sup>,  
Dragomir Anguelov<sup>1</sup>, Dumitru Erhan<sup>1</sup>, Vincent Vanhoucke<sup>1</sup>, Andrew Rabinovich<sup>4</sup>

<sup>1</sup>Google Inc. <sup>2</sup>University of North Carolina, Chapel Hill

<sup>3</sup>University of Michigan, Ann Arbor <sup>4</sup>Magic Leap Inc.

<sup>1</sup>{szegedy, jia, q, sermanet, dragomir, dimitru, vanhoucke}@google.com

<sup>2</sup>wliu@cs.unc.edu, <sup>3</sup>reedscott@umich.edu, <sup>4</sup>arabinovich@magic Leap Inc.



# Caffe

## So what is Caffe?

- Pure C++ / CUDA architecture for deep learning
  - command line, Python, MATLAB interfaces
- Fast, well-tested code
- Tools, reference models, demos, and recipes
- Seamless switch between CPU and GPU
  - `Caffe::set_mode(Caffe::GPU);`



Prototype



Training



Deployment

All with essentially the same code!

## Caffe Model Zoo

```
layer {
  name: "conv1"
  type: "Convolution"
  bottom: "data"
  top: "conv1"
  param {
    lr_mult: 1
    decay_mult: 1
  }
  param {
    lr_mult: 2
    decay_mult: 0
  }
  convolution_param {
    num_output: 96
    kernel_size: 11
    stride: 4
    weight_filler {
      type: "gaussian"
      std: 0.01
    }
    bias_filler {
      type: "constant"
      value: 0
    }
  }
}
layer {
  name: "relu1"
```

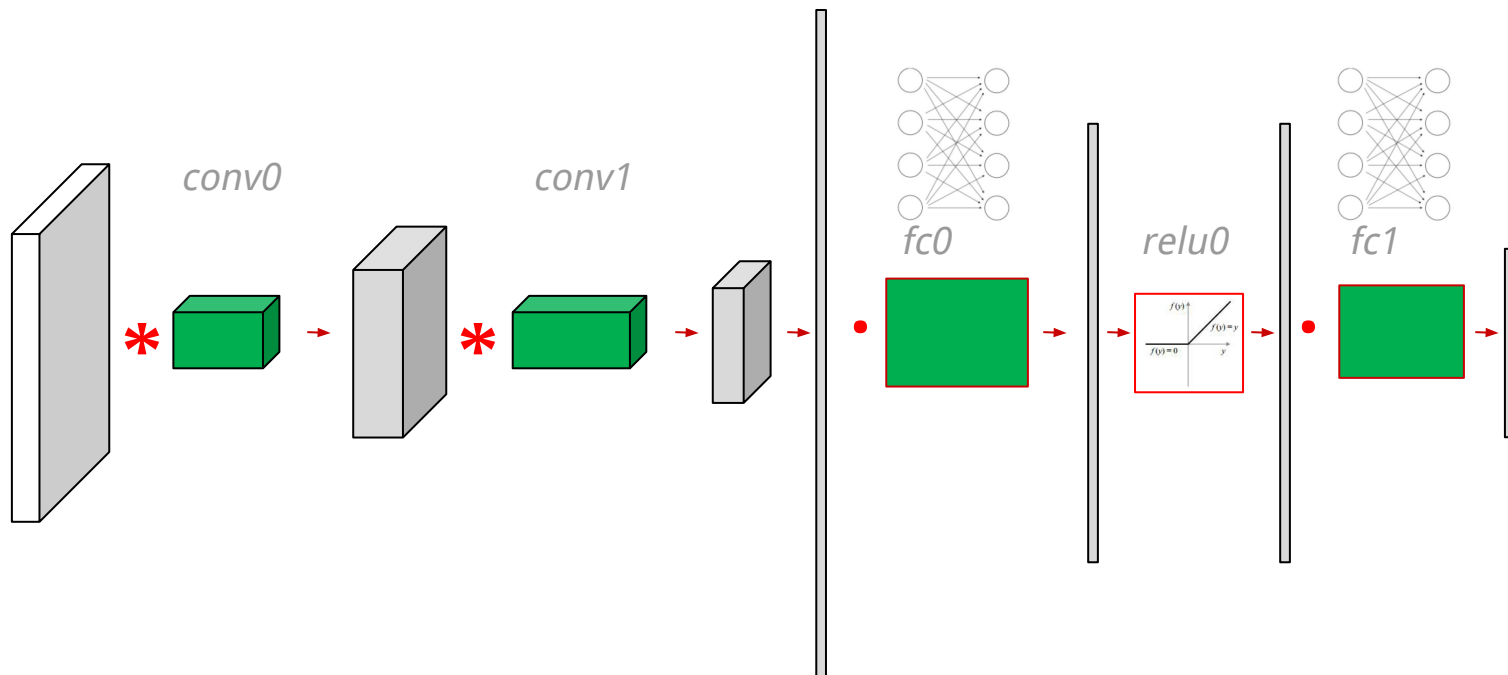


# Discriminative Feature Localization

# Fine-tuning

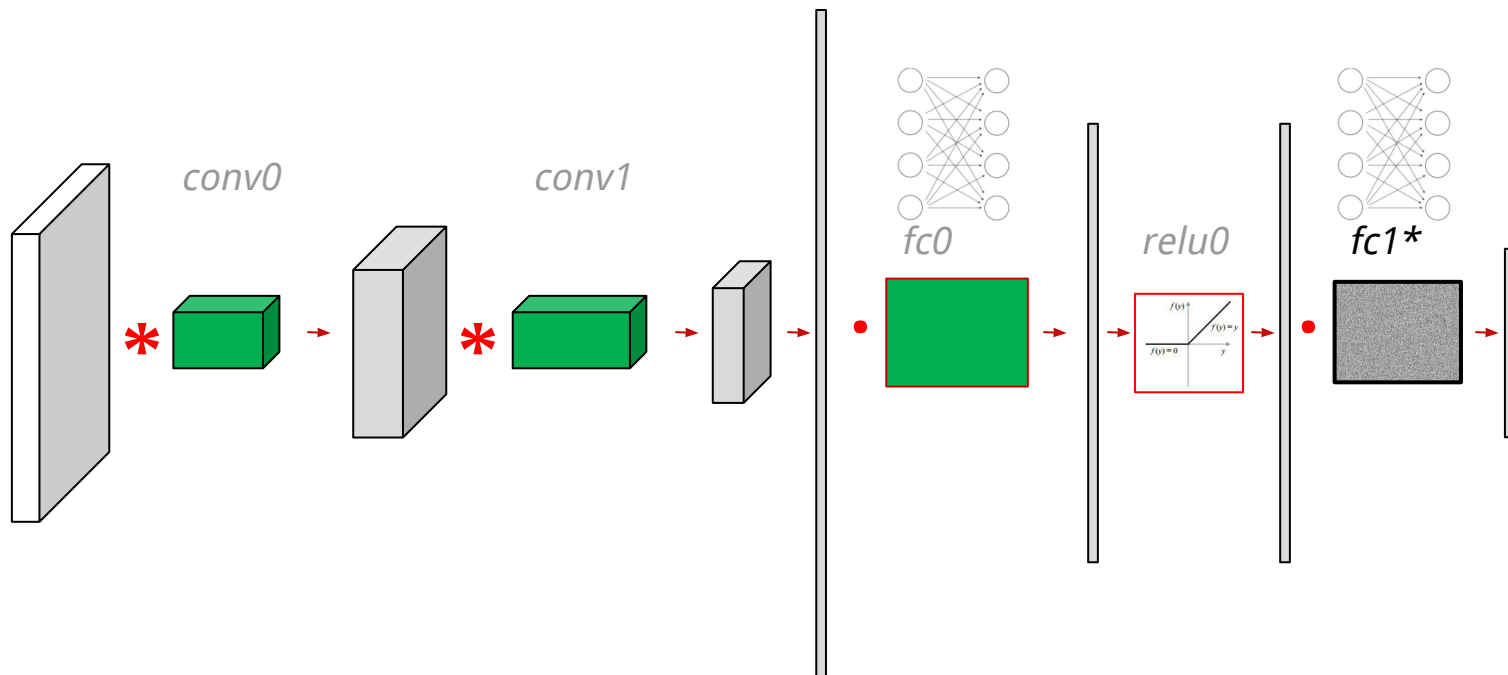
*Adjusting the parameters of an existing, pre-trained network to solve a different problem with “few” images from the new problem domain.*

# Fine-tuning



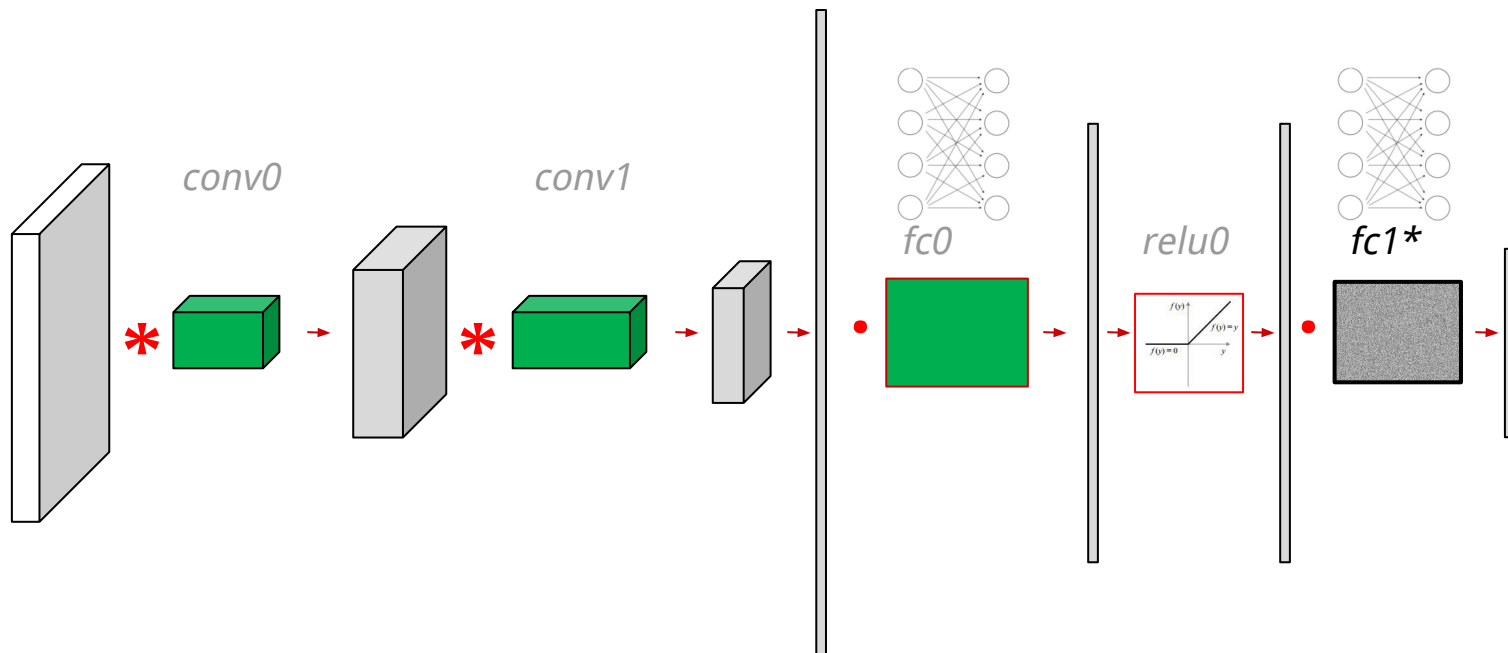
Pretrained network for problem **P** Machine Learning Workshop Galicia 2016

# Fine-tuning



Substitute layer(s)

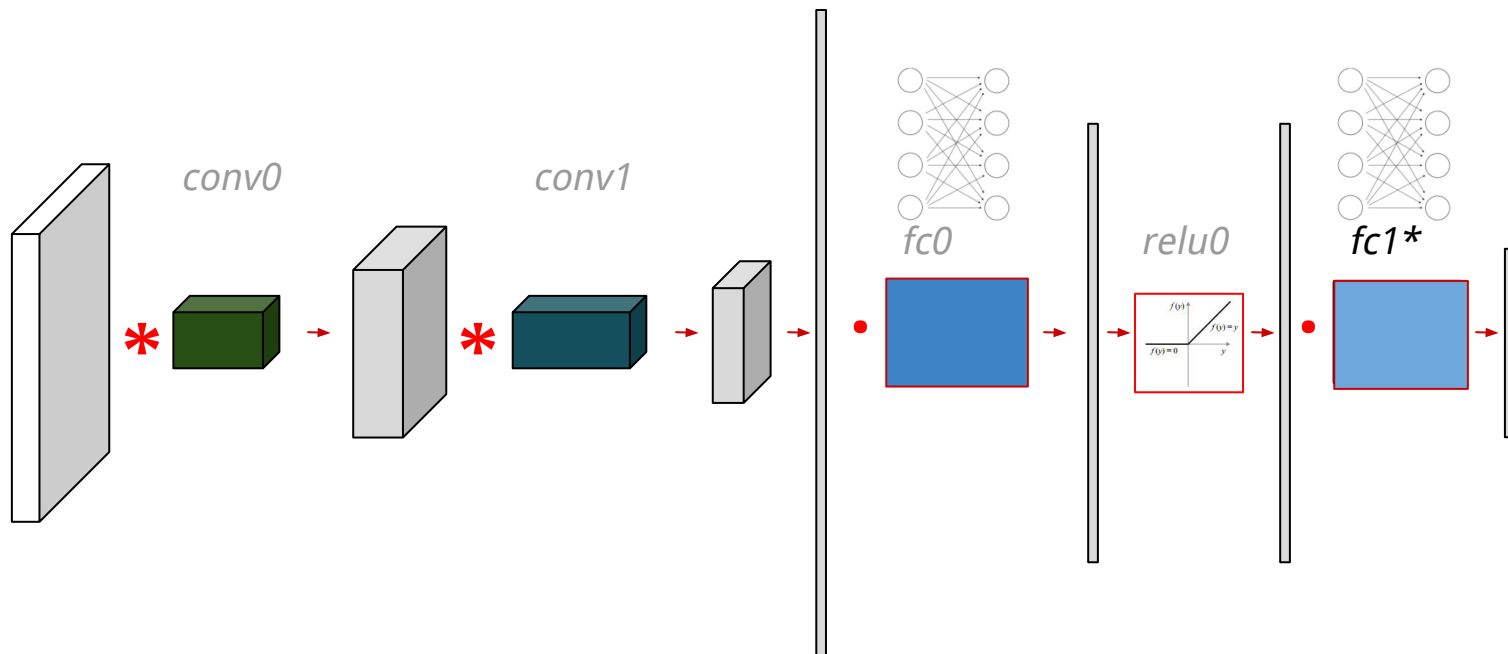
# Fine-tuning



Train with new samples for problem Q

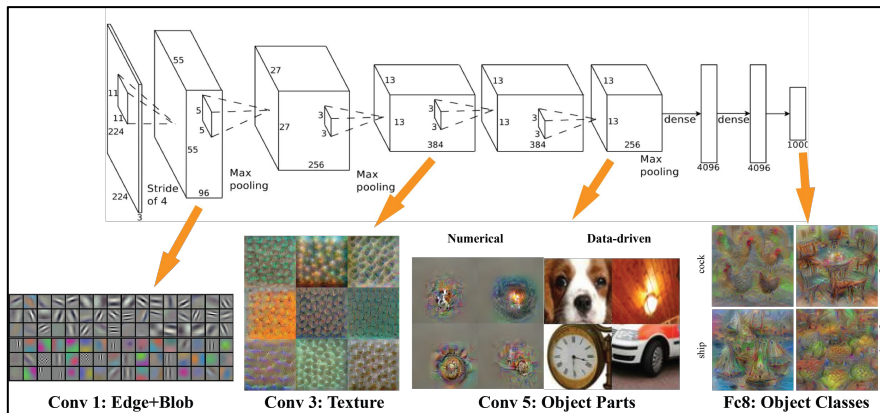
Learning Workshop Galicia 2016

# Fine-tuning

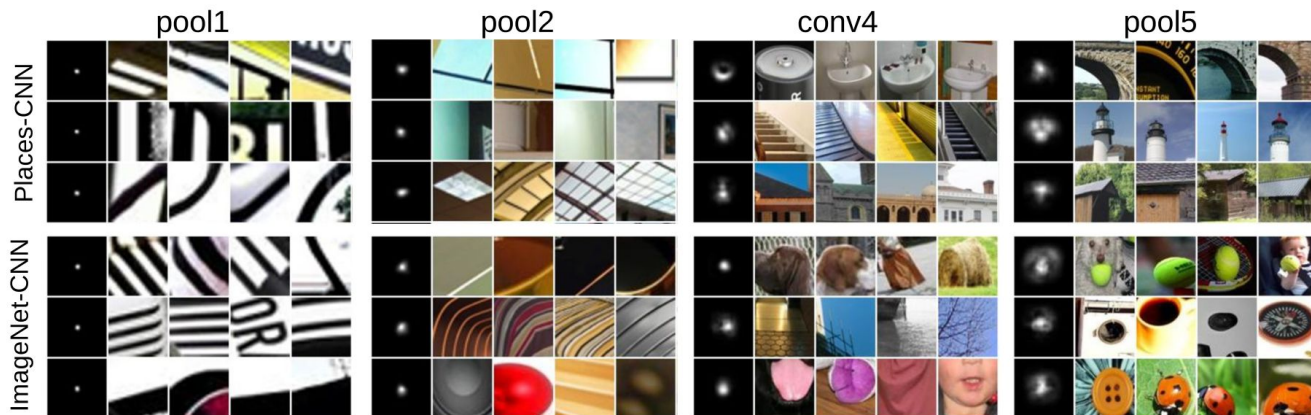


Train with new samples for problem **Q** ing Workhop Galicia 2016

# Image Semantics



# Image Semantics



Published as a conference paper at ICLR 2015

## OBJECT DETECTORS EMERGE IN DEEP SCENE CNNs

**Bolei Zhou, Aditya Khosla, Agata Lapedriza, Aude Oliva, Antonio Torralba**  
Computer Science and Artificial Intelligence Laboratory, MIT  
{bolei, khosla, agata, oliva, torralba}@mit.edu



# Learning Deep Features for Discriminative Localization

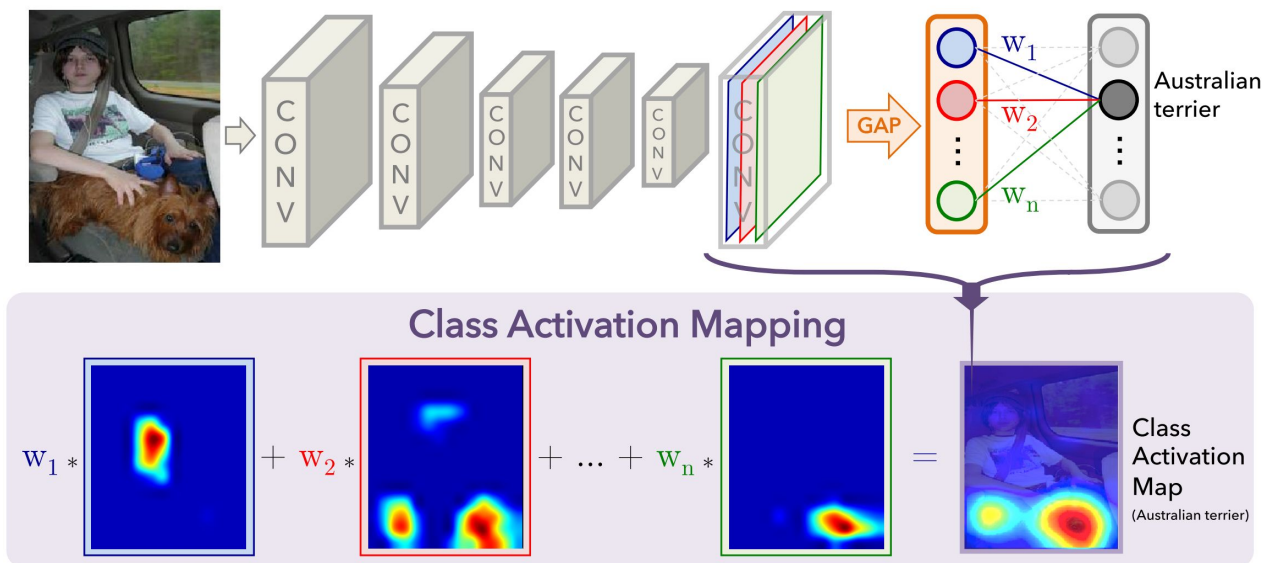
Bolei Zhou, Aditya Khosla, Agata Lapedriza, Aude Oliva, Antonio Torralba  
Computer Science and Artificial Intelligence Laboratory, MIT  
{bzhou, khosla, agata, oliva, torralba}@csail.mit.edu

## Abstract

*In this work, we revisit the global average pooling layer proposed in [13], and shed light on how it explicitly enables the convolutional neural network to have remarkable localization ability despite being trained on image-level labels. While this technique was previously proposed as a means for regularizing training, we find that it actually builds a generic localizable deep representation that can be applied to a variety of tasks. Despite the apparent simplicity of global average pooling, we are able to achieve 37.1% top-5*



Figure 1. A simple modification of the global average pool-

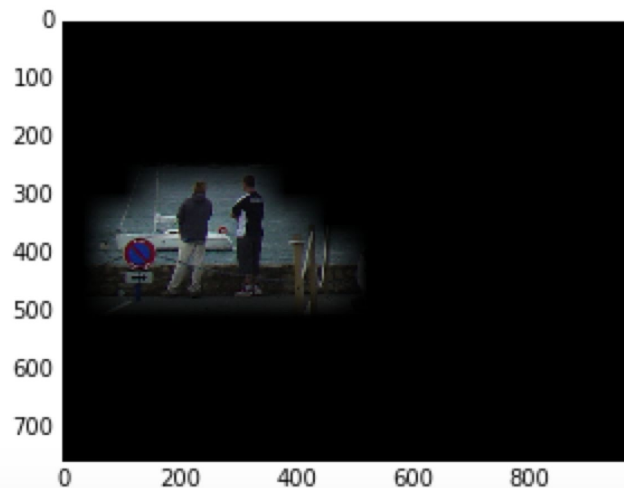
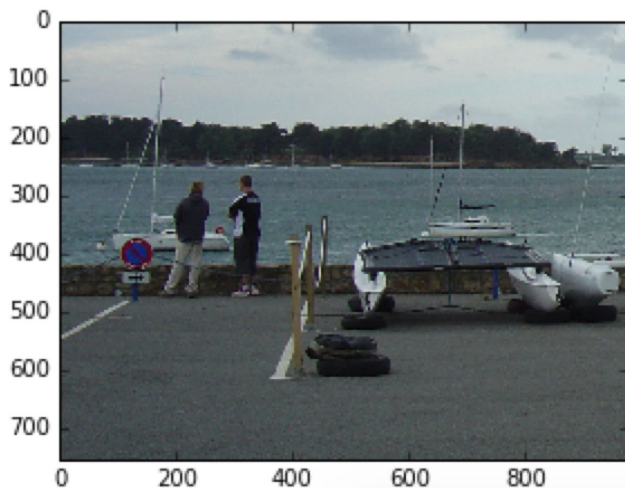


# Experiment: Localization using CAM

2 classes : {person, not person}

Image classified as 1 with probability 0.999978

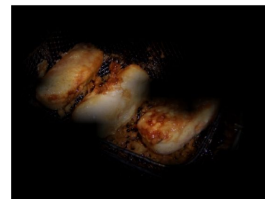
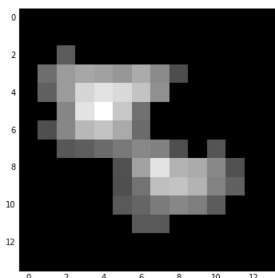
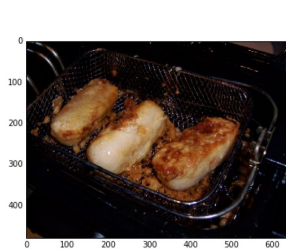
Class no. 1 Activation Mapping



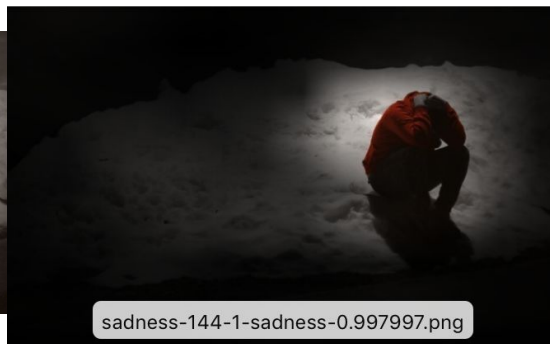
# Experiment: Sentiment Analysis + CAM

6 classes : {anger, disgust, fear, joy, sadness, surprise}

classified as "disgust" (0.999673)

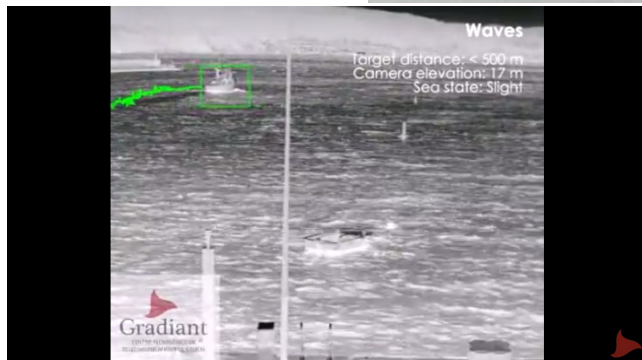
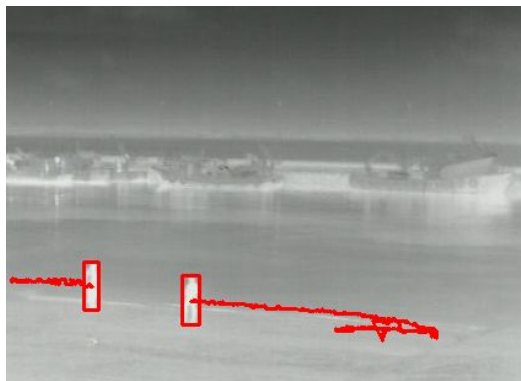


sadness-144-0.png



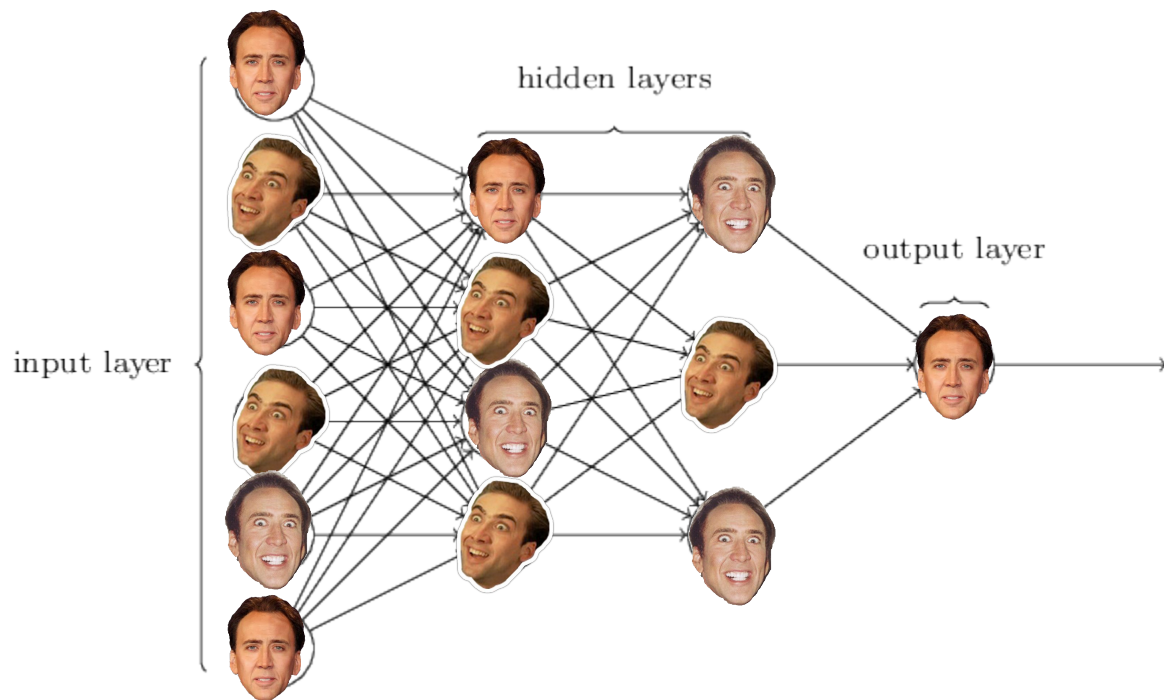
sadness-144-1-sadness-0.997997.png

# New Challenges: Stable Object Detection



**SEERS**  
Snapshot Spectral Imager for IR Surveillance

# THANK YOU!



## Questions?

Machine Learning Workshop Galicia 2016