

# Computer Vision - Lecture 15

## Deep Learning for Object Categorization

21.12.2016

Bastian Leibe

RWTH Aachen

<http://www.vision.rwth-aachen.de>

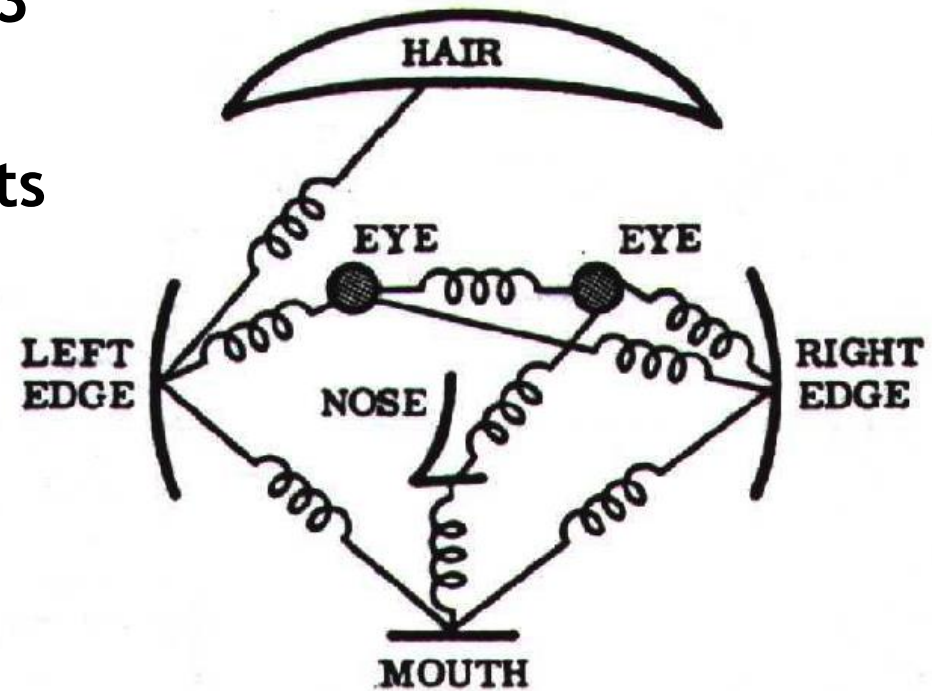
[leibe@vision.rwth-aachen.de](mailto:leibe@vision.rwth-aachen.de)

# Course Outline

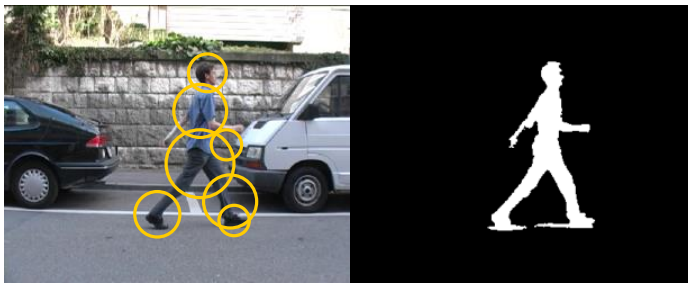
- Image Processing Basics
- Segmentation & Grouping
- Object Recognition
- Object Categorization I
  - Sliding Window based Object Detection
- Local Features & Matching
  - Local Features - Detection and Description
  - Recognition with Local Features
  - Indexing & Visual Vocabularies
- Object Categorization II
  - Bag-of-Words Approaches & Part-based Approaches
  - **Deep Learning Methods**
- 3D Reconstruction

# Recap: Part-Based Models

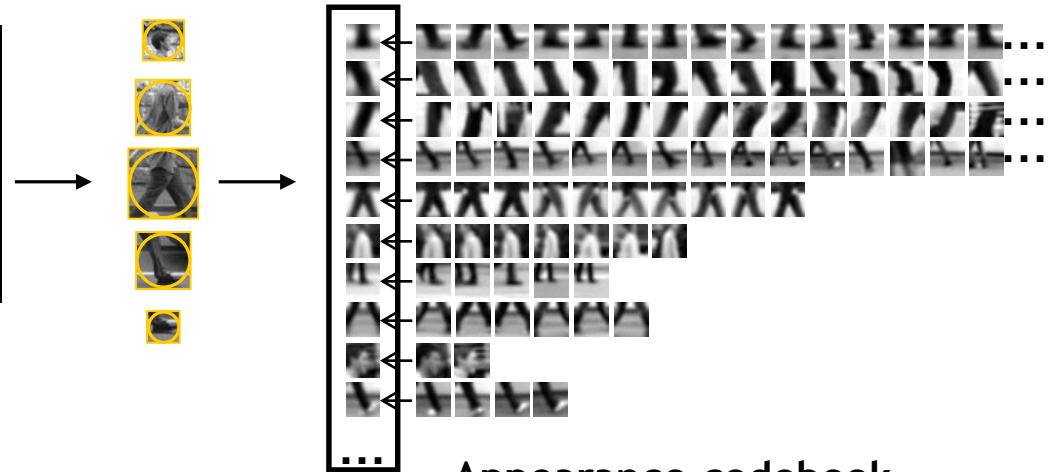
- Fischler & Elschlager 1973
- Model has two components
  - parts  
(2D image fragments)
  - structure  
(configuration of parts)



# Recap: Implicit Shape Model - Representation

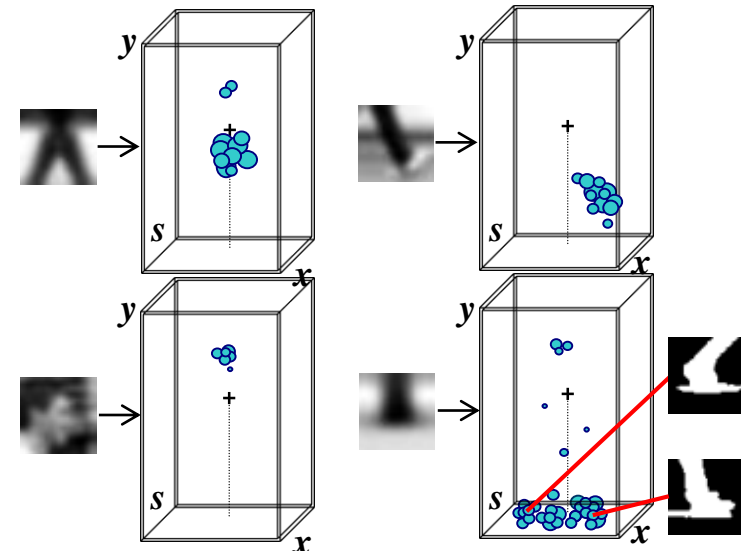


Training images  
(+reference segmentation)



Appearance codebook

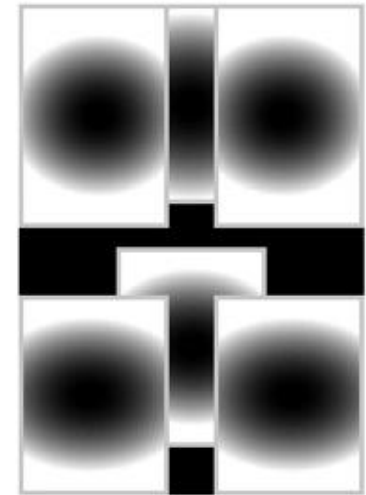
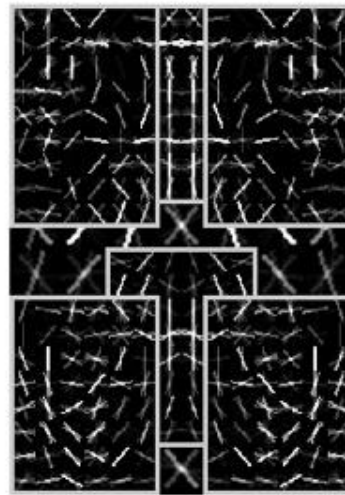
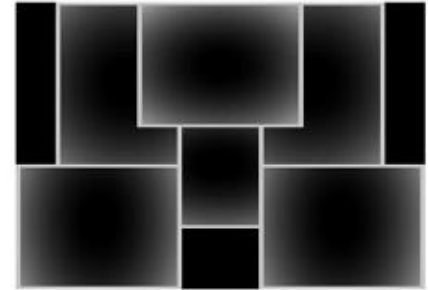
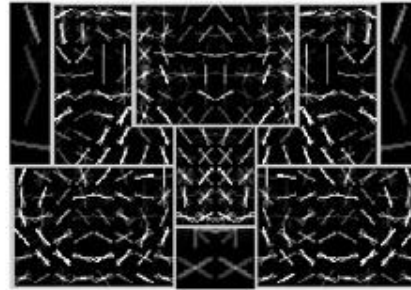
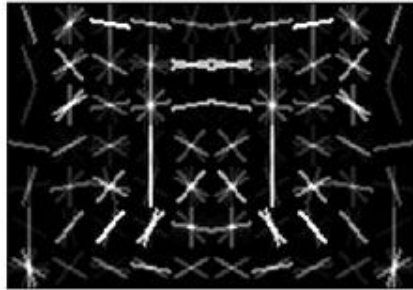
- Learn appearance codebook
  - Extract local features at interest points
  - Clustering  $\Rightarrow$  appearance codebook
- Learn spatial distributions
  - Match codebook to training images
  - Record matching positions on object



Spatial occurrence distributions

+ local figure-ground labels 5

# Recap: Deformable Part-Based Model

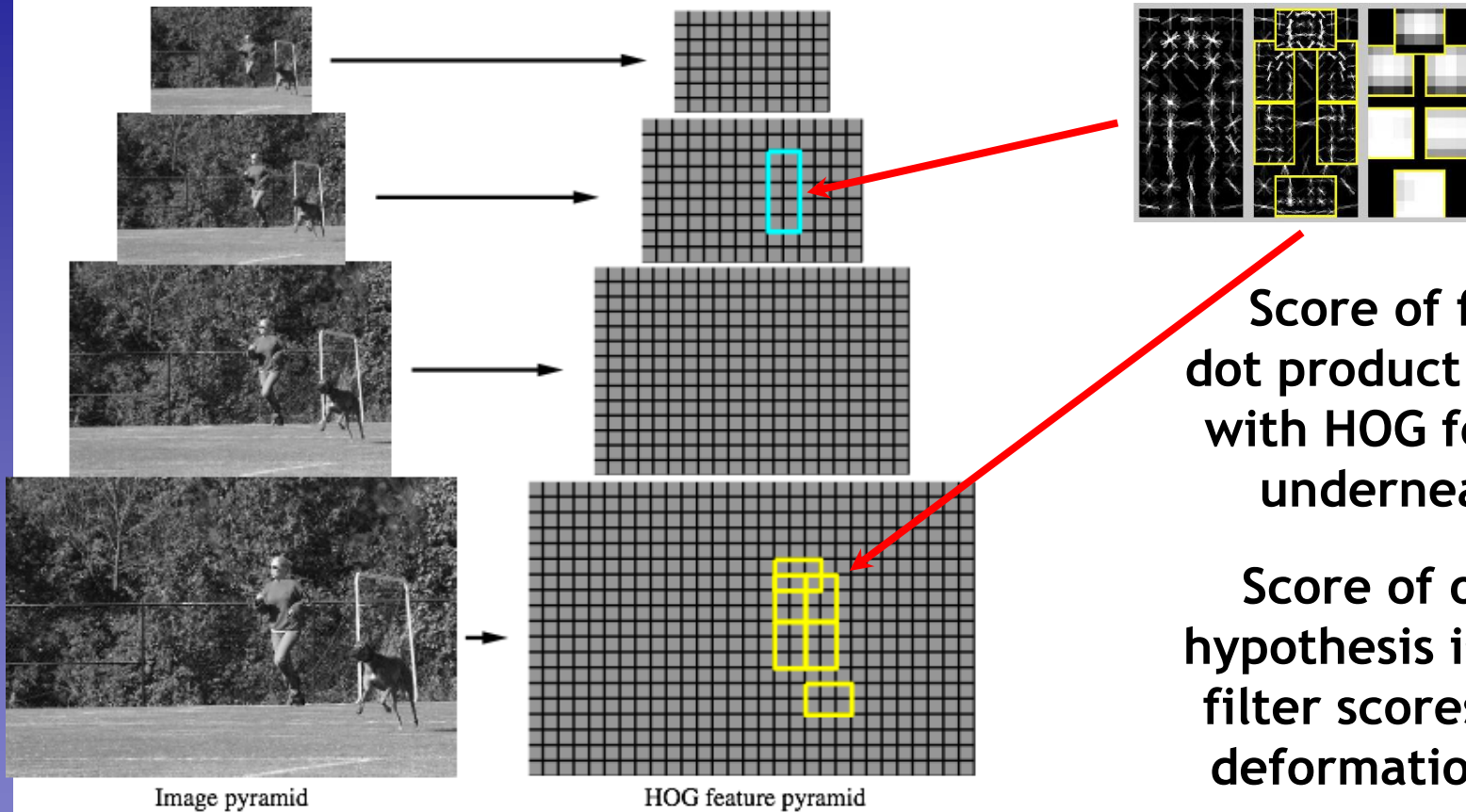


**Root filters**  
coarse resolution

**Part filters**  
finer resolution

**Deformation**  
models

# Recap: Object Hypothesis



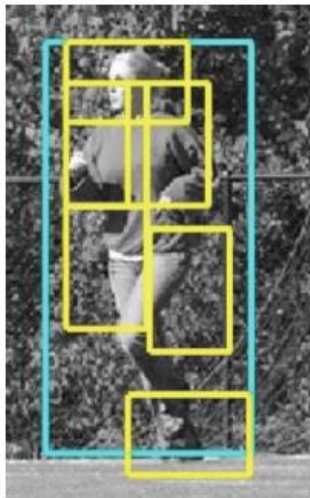
- **Multiscale model captures features at two resolutions**

# Recap: Score of a Hypothesis

$$\text{score}(p_0, \dots, p_n) = \sum_{i=0}^n F_i \cdot \phi(H, p_i) - \sum_{i=1}^n d_i \cdot (dx_i^2, dy_i^2)$$

“data term”  
 $\sum_{i=0}^n F_i \cdot \phi(H, p_i)$   
 ↑ filters
 

 “spatial prior”  
 $\sum_{i=1}^n d_i \cdot (dx_i^2, dy_i^2)$   
 ↑ displacements  
 deformation parameters



$$\text{score}(z) = \beta \cdot \Psi(H, z)$$

concatenation filters and deformation parameters

concatenation of HOG features and part displacement features



# Topics of This Lecture

- **Deep Learning**
  - Motivation
- **Convolutional Neural Networks**
  - Convolutional Layers
  - Pooling Layers
  - Nonlinearities
- **CNN Architectures**
  - LeNet
  - AlexNet
  - VGGNet
  - GoogLeNet
- **Applications**

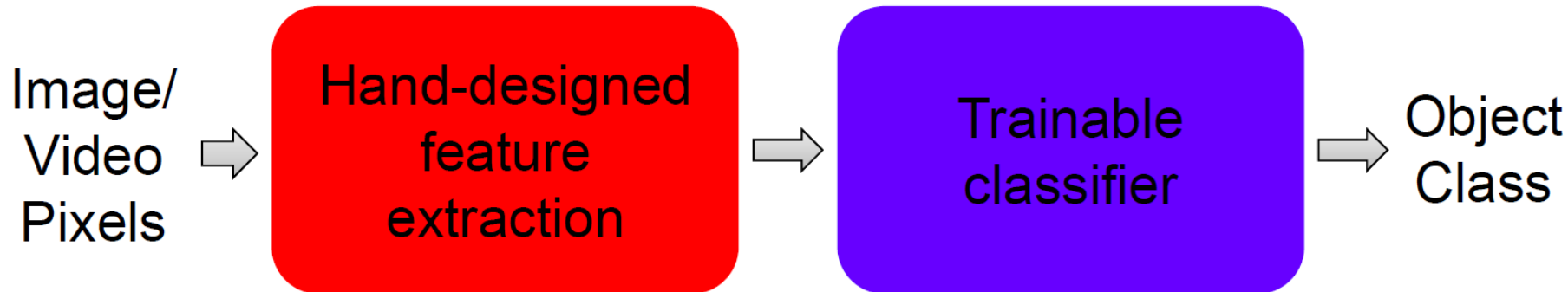


**We've finally got there!**



**Deep Learning**

# Traditional Recognition Approach

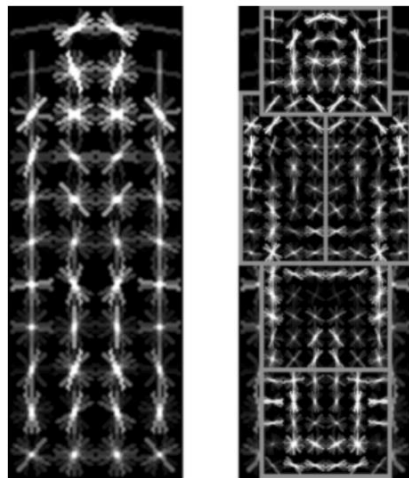


- **Characteristics**

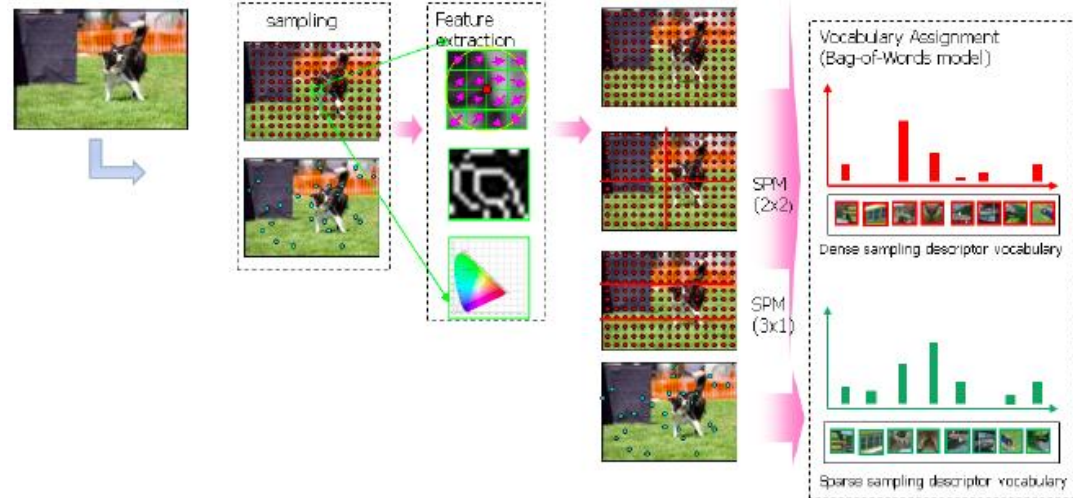
- Features are not learned, but engineered
  - Trainable classifier is often generic (e.g., SVM)
- ⇒ Many successes in 2000-2010.

# Traditional Recognition Approach

- Features are key to recent progress in recognition
    - Multitude of hand-designed features currently in use
    - SIFT, HOG, .....
- ⇒ *Where next? Better classifiers? Or keep building more features?*



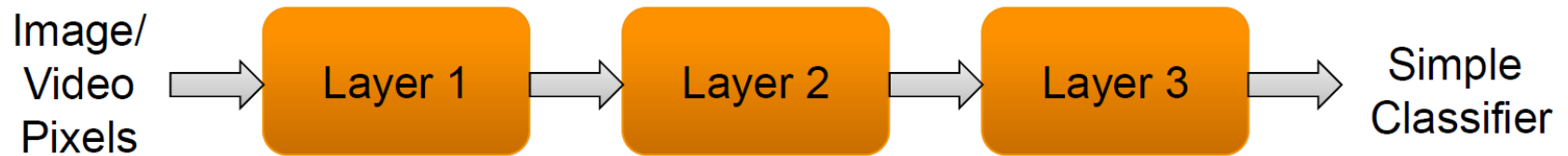
**DPM**  
[Felzenszwalb  
et al., PAMI'07]



**Dense SIFT+LBP+HOG → BOW → Classifier**  
[Yan & Huan '10]  
(Winner of PASCAL 2010 Challenge)

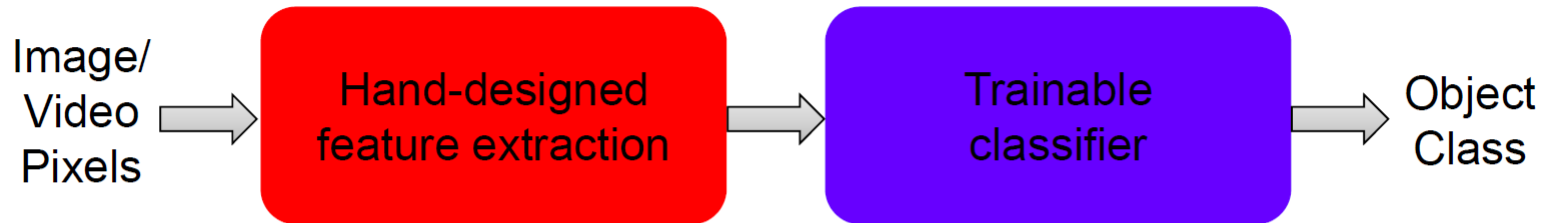
# What About Learning the Features?

- Learn a *feature hierarchy* all the way from pixels to classifier
  - Each layer extracts features from the output of previous layer
  - Train all layers jointly



# “Shallow” vs. “Deep” Architectures

## Traditional recognition: “Shallow” architecture

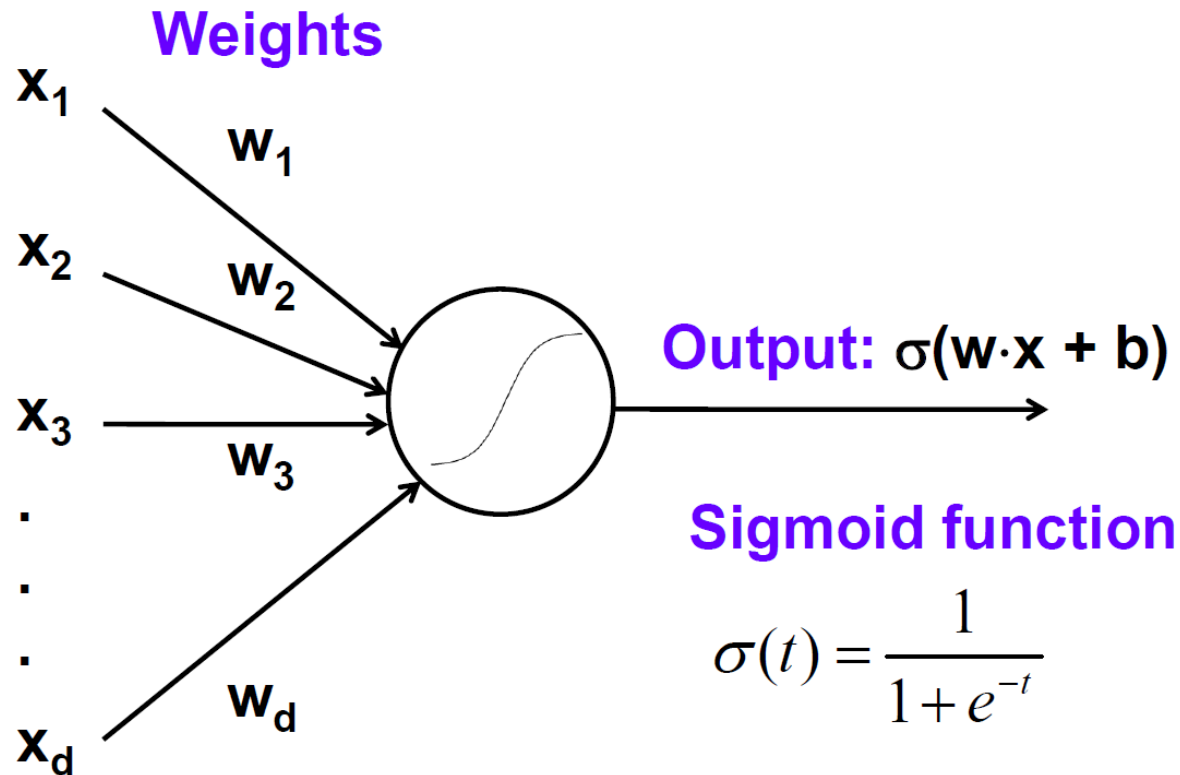


## Deep learning: “Deep” architecture

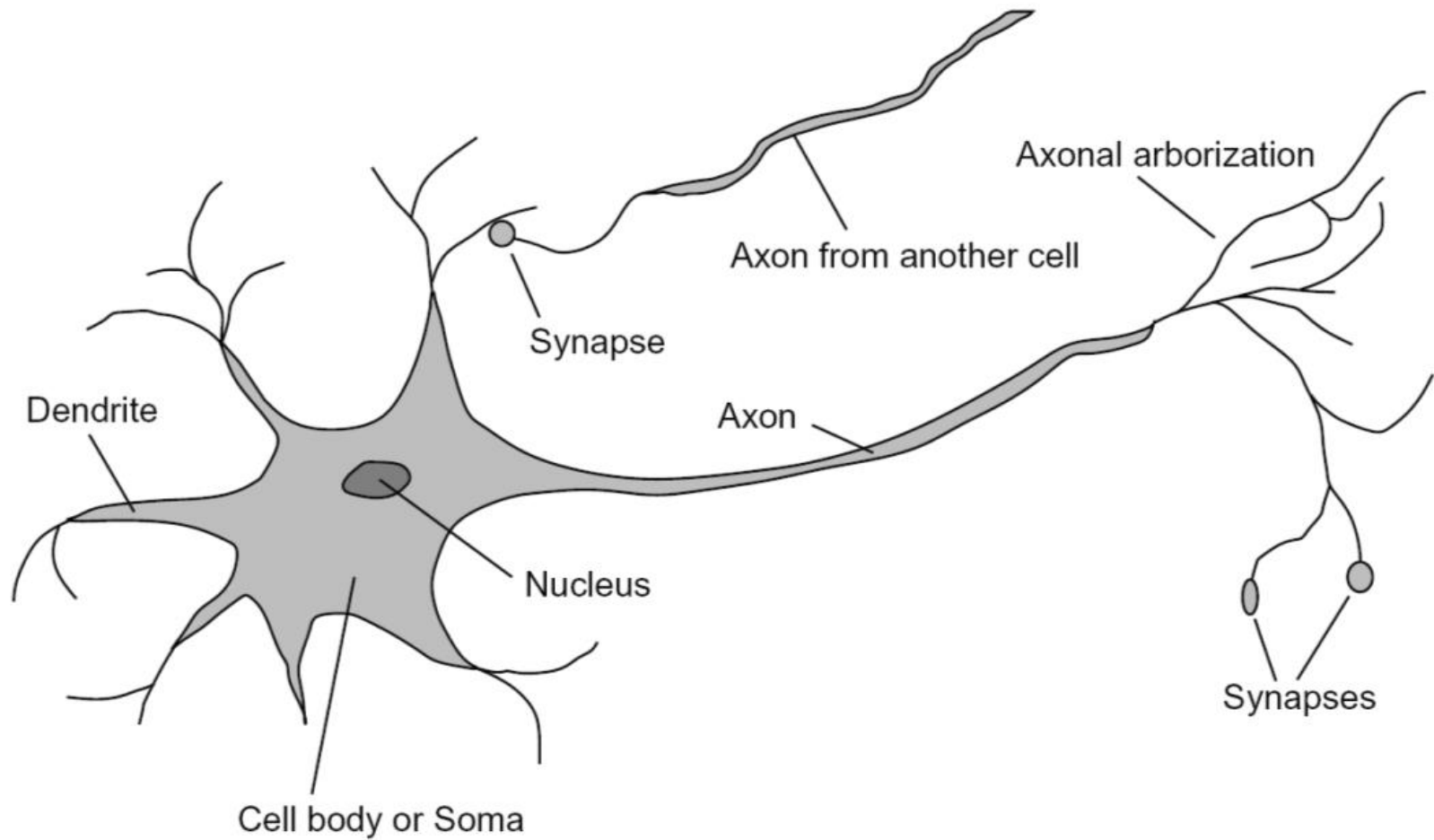


# Background: Perceptrons

Input

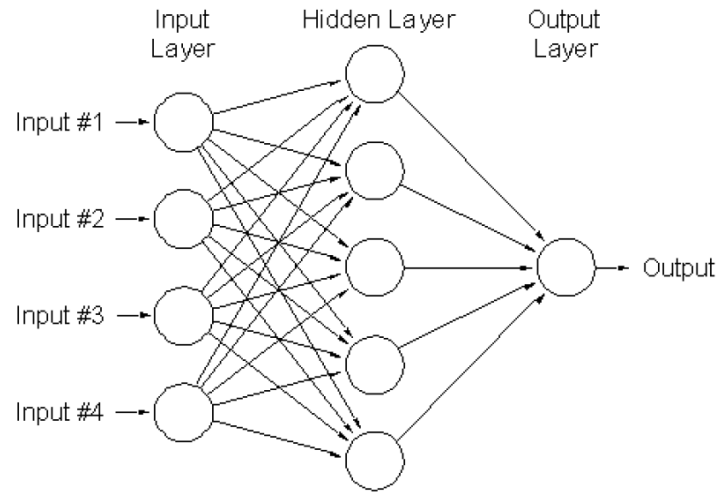


# Inspiration: Neuron Cells





# Background: Multi-Layer Neural Networks



- **Nonlinear classifier**

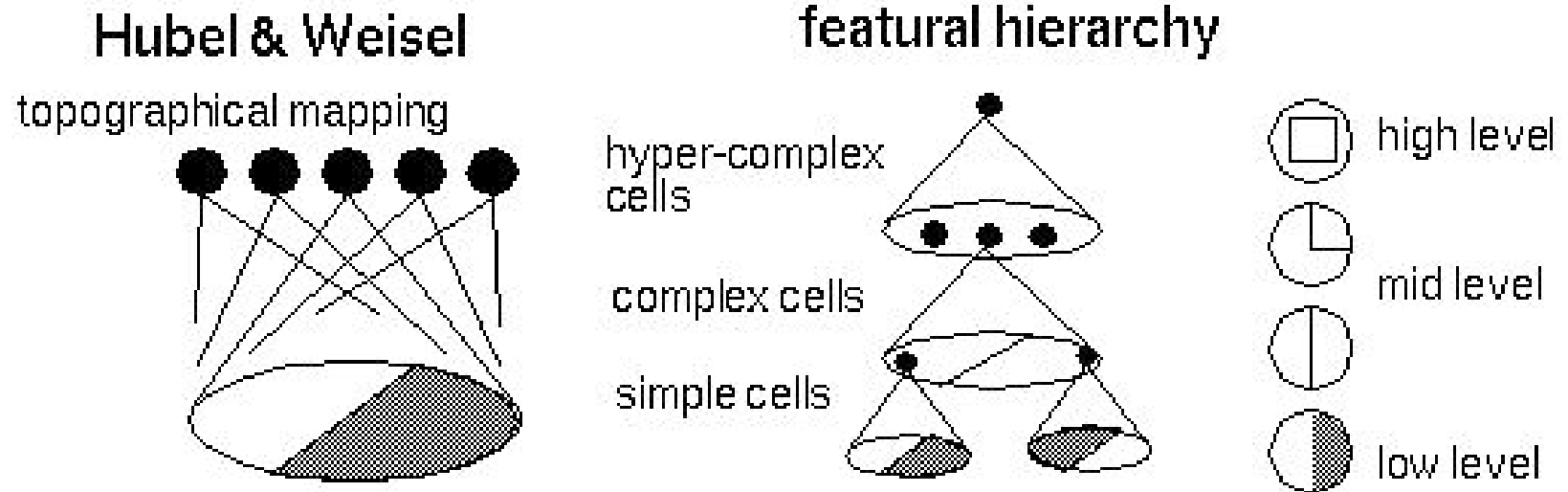
- **Training:** find network weights  $\mathbf{w}$  to minimize the error between true training labels  $t_n$  and estimated labels  $f_{\mathbf{w}}(x_n)$ :

$$E(\mathbf{W}) = \sum_n L(t_n, f(\mathbf{x}_n; \mathbf{W}))$$

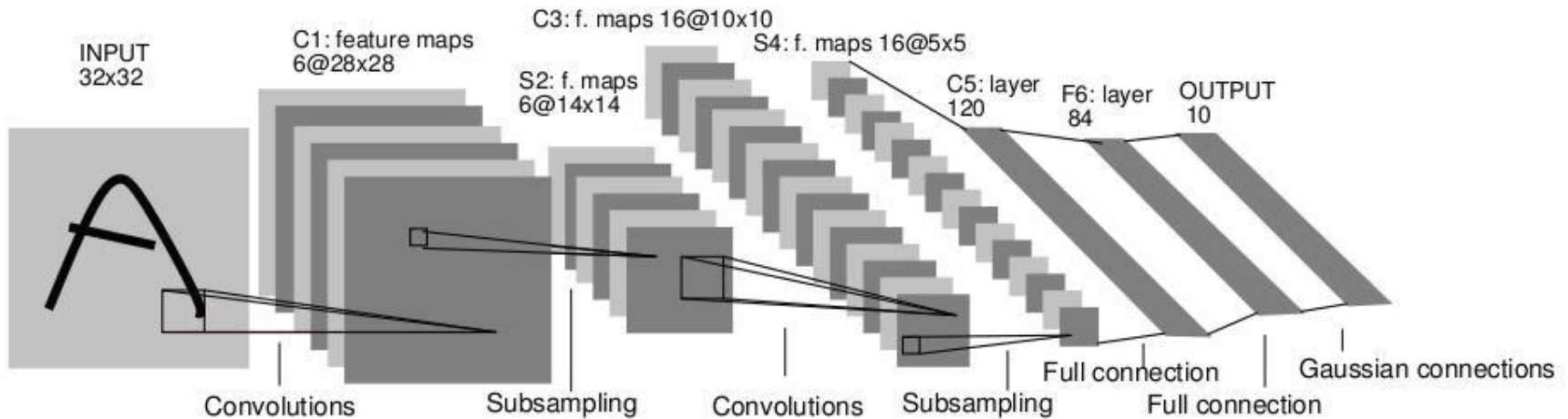
- Minimization can be done by gradient descent, provided  $f$  is differentiable
  - Training method: **error backpropagation.**

# Hubel/Wiesel Architecture

- D. Hubel, T. Wiesel (1959, 1962, Nobel Prize 1981)
  - Visual cortex consists of a hierarchy of *simple*, *complex*, and *hyper-complex* cells



# Convolutional Neural Networks (CNN, ConvNet)



- Neural network with specialized connectivity structure
  - Stack multiple stages of feature extractors
  - Higher stages compute more global, more invariant features
  - Classification layer at the end

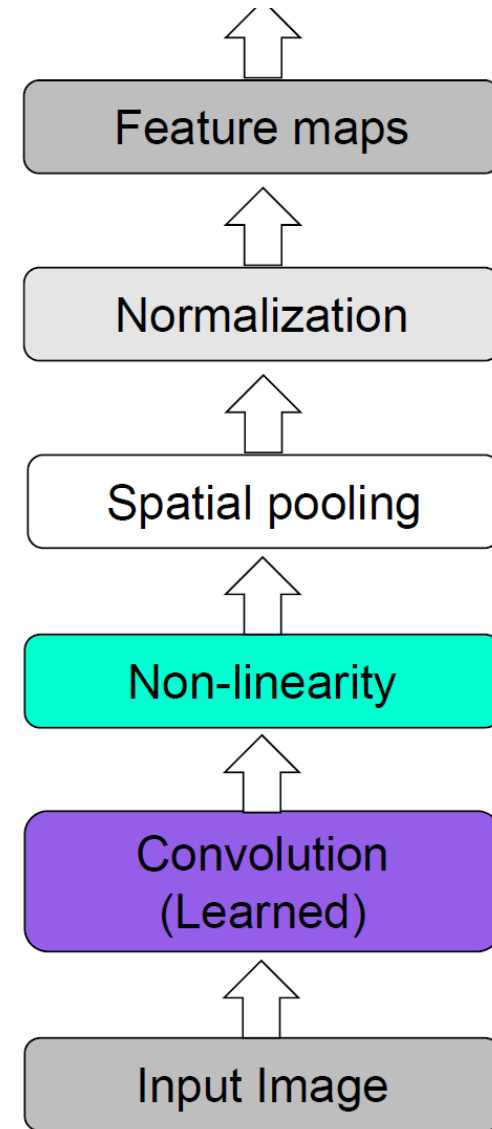
Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner, [Gradient-based learning applied to document recognition](#), Proceedings of the IEEE 86(11): 2278-2324, 1998.

# Topics of This Lecture

- Deep Learning
  - Motivation
- Convolutional Neural Networks
  - Convolutional Layers
  - Pooling Layers
  - Nonlinearities
- CNN Architectures
  - LeNet
  - AlexNet
  - VGGNet
  - GoogLeNet
- Applications

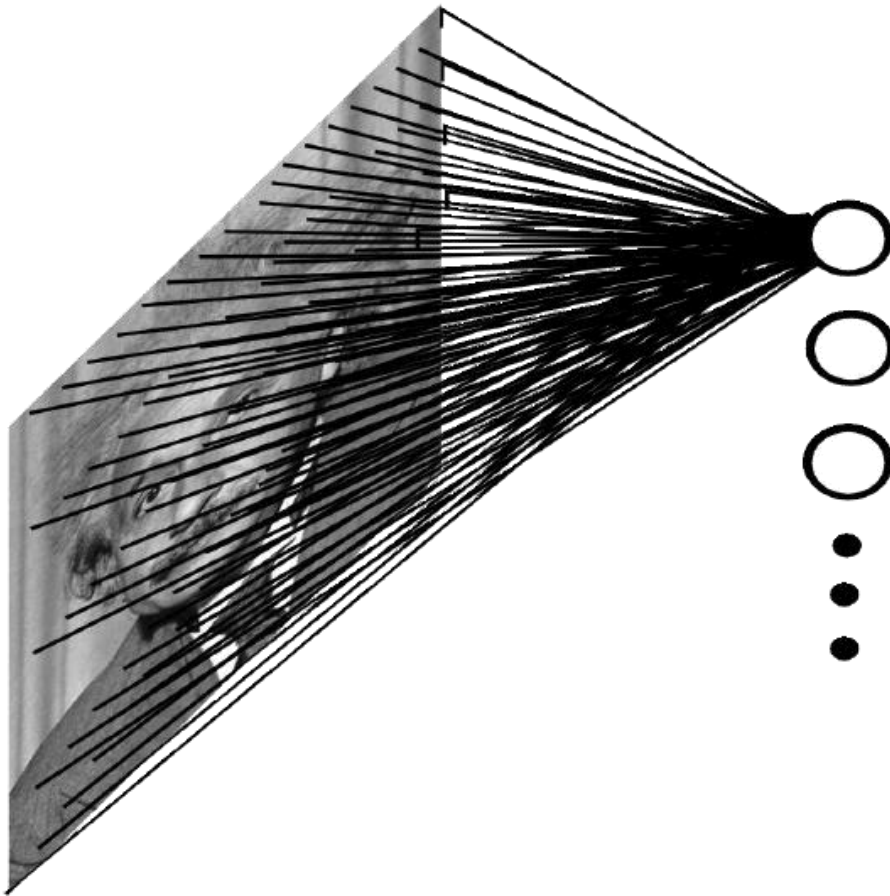
# Convolutional Networks: Structure

- **Feed-forward feature extraction**
  1. Convolve input with learned filters
  2. Non-linearity
  3. Spatial pooling
  4. (Normalization)
- **Supervised training of convolutional filters by back-propagating classification error**



# Convolutional Networks: Intuition

- Fully connected network
  - E.g.  $1000 \times 1000$  image  
1M hidden units  
⇒ 1T parameters!

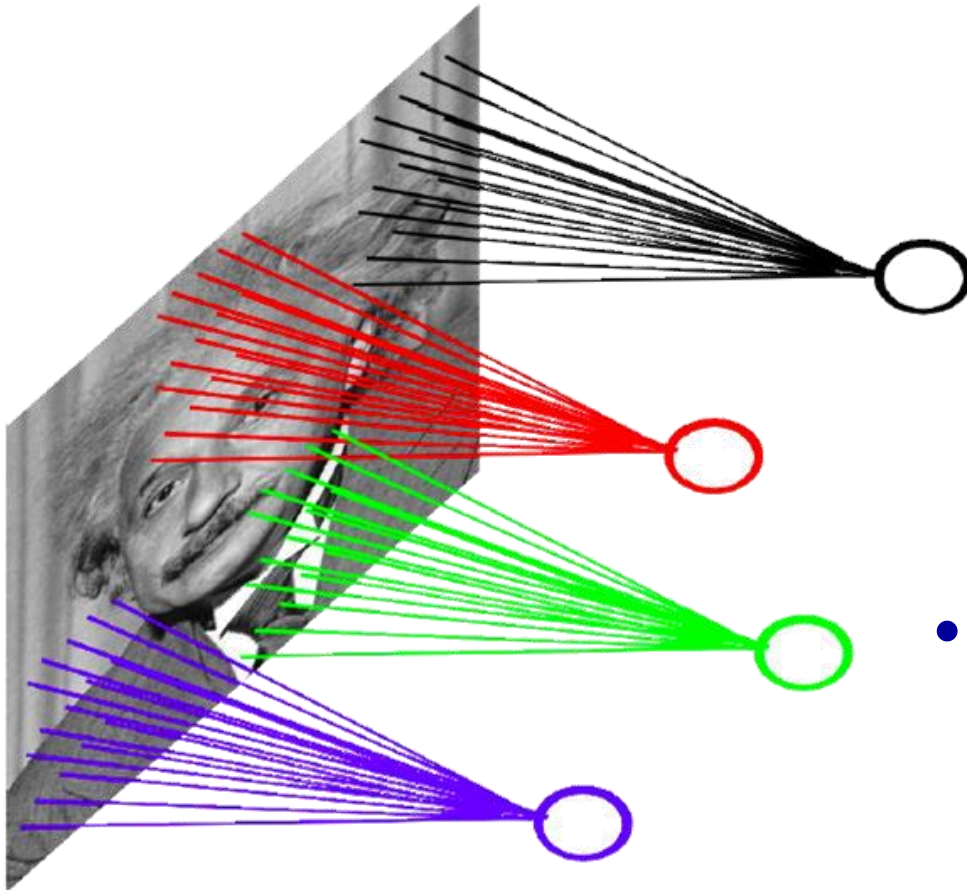


- Ideas to improve this
  - Spatial correlation is local

# Convolutional Networks: Intuition

- **Locally connected net**

- E.g.  $1000 \times 1000$  image  
1M hidden units  
 $10 \times 10$  receptive fields  
⇒ 100M parameters!



- **Ideas to improve this**

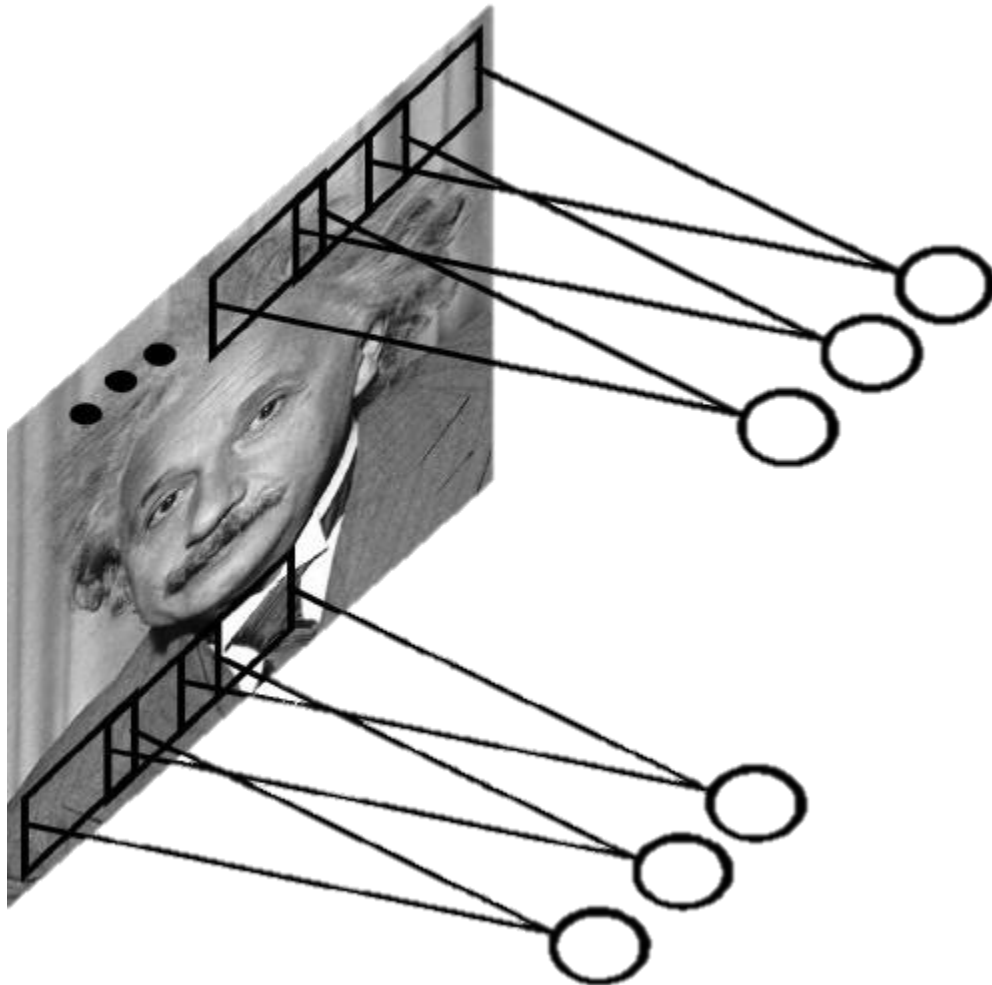
- Spatial correlation is local
- Want translation invariance



# Convolutional Networks: Intuition

- Convolutional net

- Share the same parameters across different locations
- Convolutions with learned kernels



# Convolutional Networks: Intuition

- Convolutional net

- Share the same parameters across different locations
- Convolutions with learned kernels

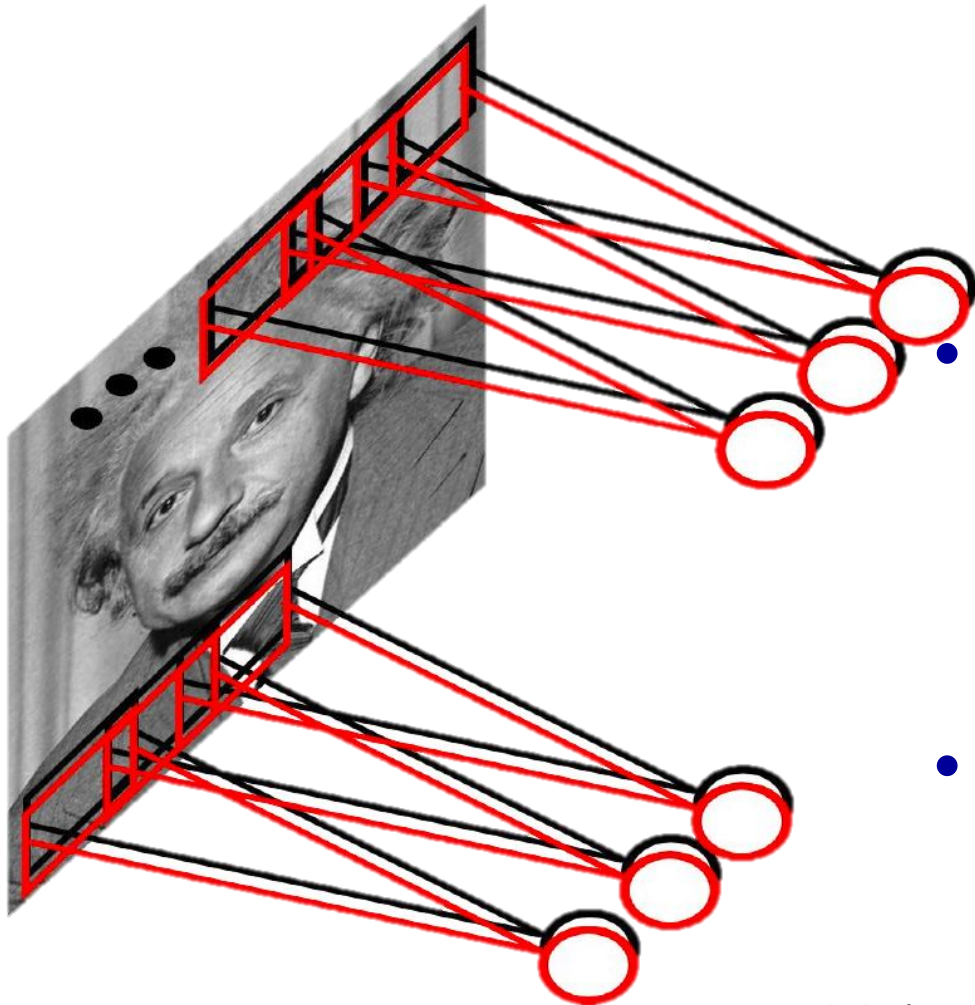
- Learn *multiple* filters

- E.g.  $1000 \times 1000$  image  
100 filters  
 $10 \times 10$  filter size

⇒ 10k parameters

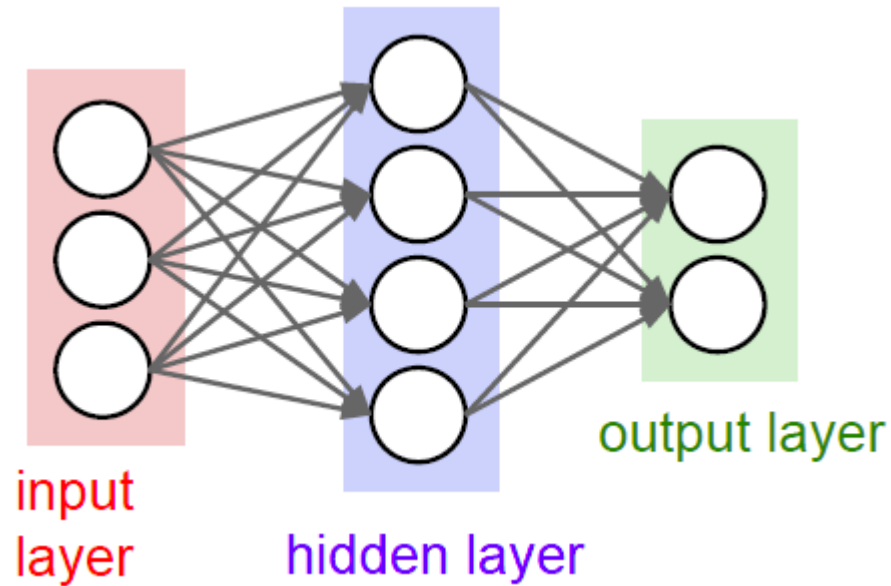
- Result: Response map

- size:  $1000 \times 1000 \times 100$
- Only memory, not params!

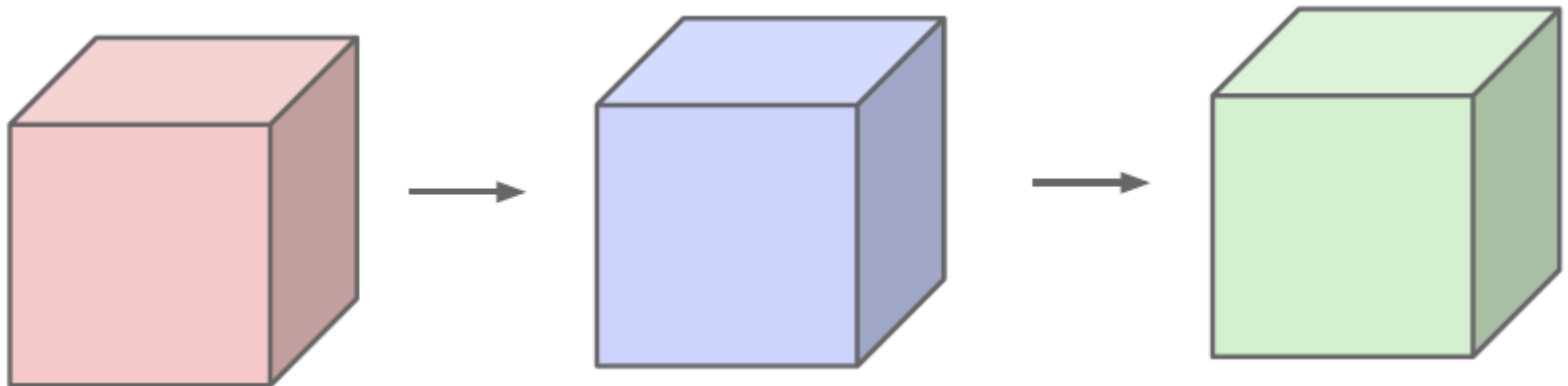


# Important Conceptual Shift

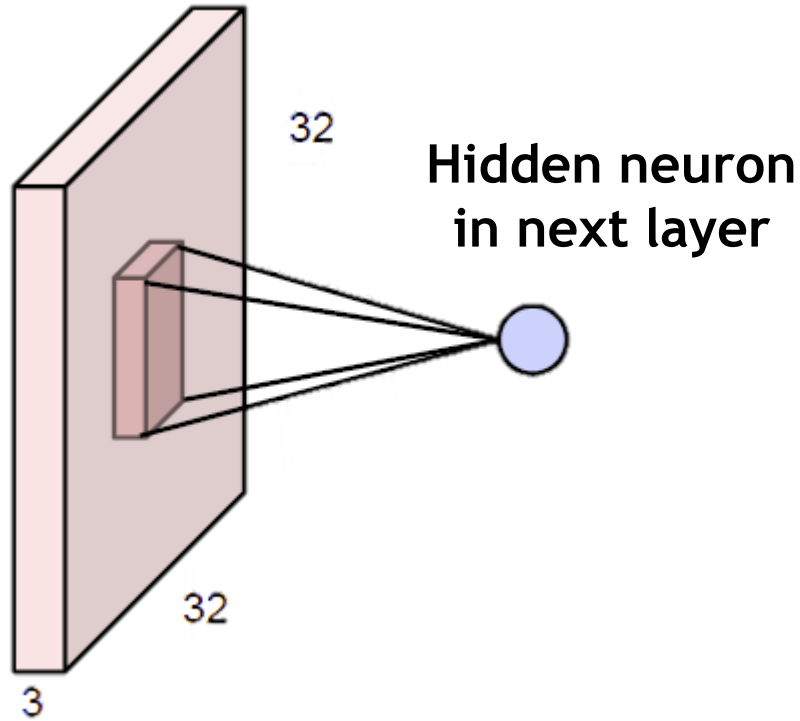
- Before



- Now:



# Convolution Layers



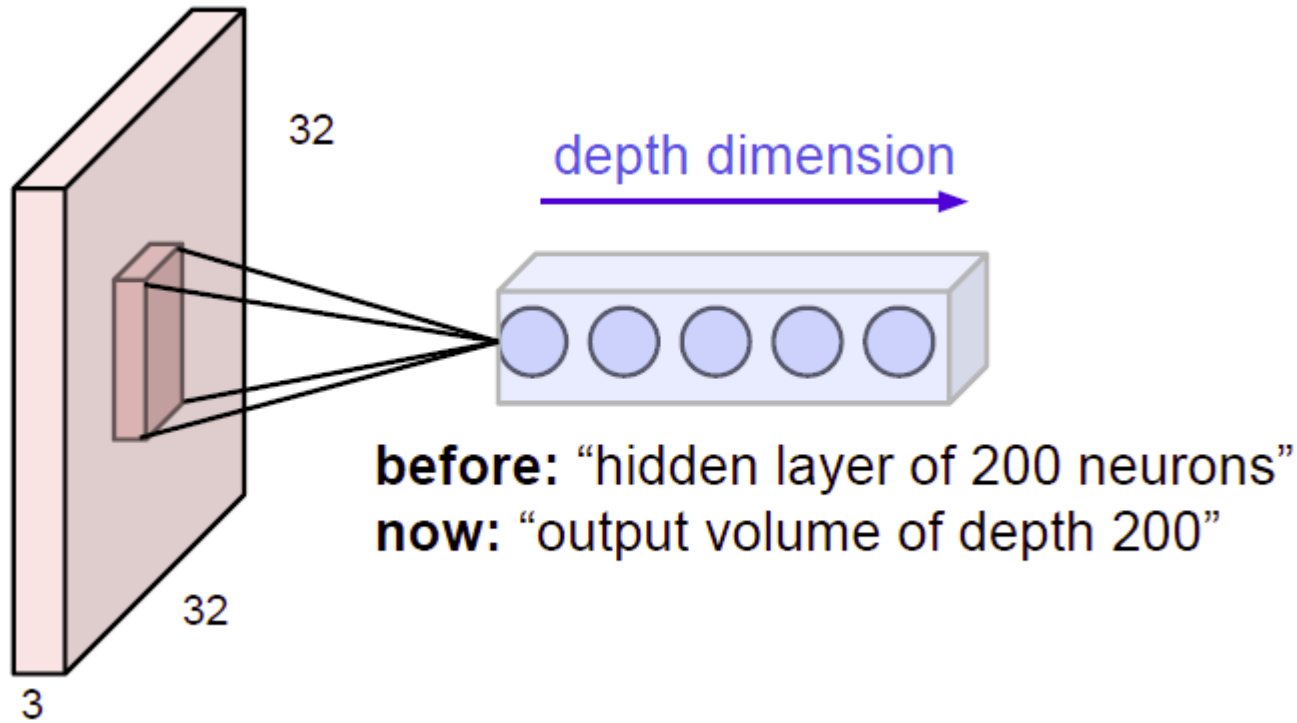
Example  
image:  $32 \times 32 \times 3$  volume

**Before:** Full connectivity  
 $32 \times 32 \times 3$  weights

**Now:** Local connectivity  
One neuron connects to, e.g.,  
 $5 \times 5 \times 3$  region.  
 $\Rightarrow$  Only  $5 \times 5 \times 3$  **shared weights.**

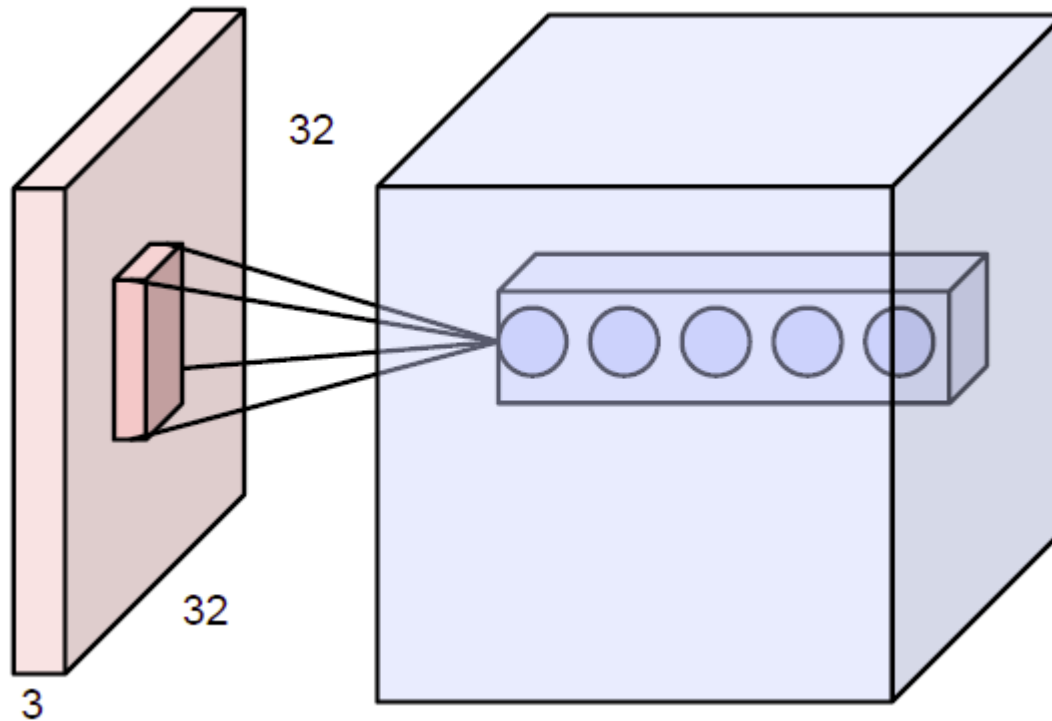
- **Note: Connectivity is**
  - Local in space ( $5 \times 5$  inside  $32 \times 32$ )
  - But full in depth (all 3 depth channels)

# Convolution Layers

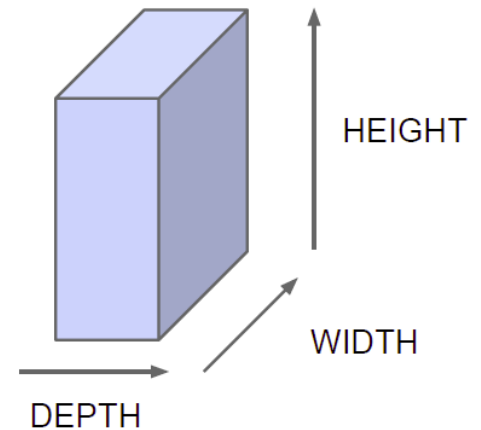


- All Neural Net activations arranged in 3 dimensions
  - Multiple neurons all looking at the same input region, stacked in depth

# Convolution Layers

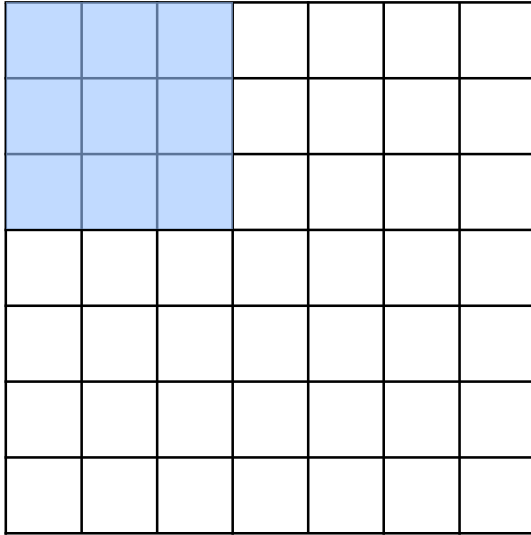


Naming convention:



- All Neural Net activations arranged in 3 dimensions
  - Multiple neurons all looking at the same input region, stacked in depth
  - Form a single  $[1 \times 1 \times \text{depth}]$  depth column in output volume.

# Convolution Layers

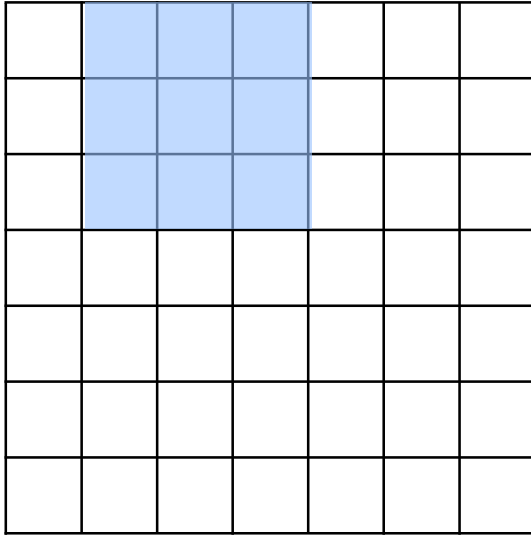


Example:  
 $7 \times 7$  input  
assume  $3 \times 3$  connectivity  
stride 1

- Replicate this column of hidden neurons across space, with some **stride**.



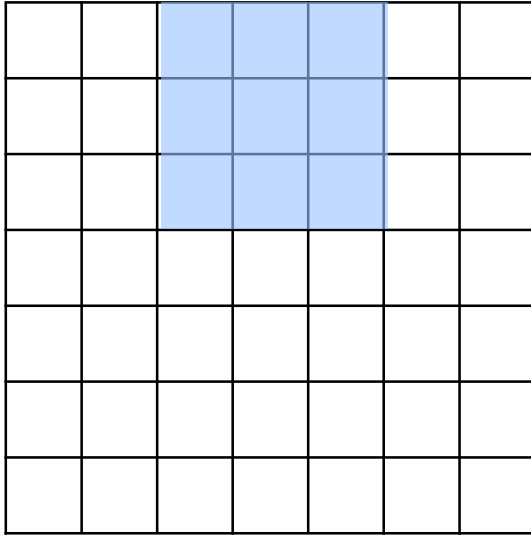
# Convolution Layers



Example:  
 $7 \times 7$  input  
assume  $3 \times 3$  connectivity  
stride 1

- Replicate this column of hidden neurons across space, with some **stride**.

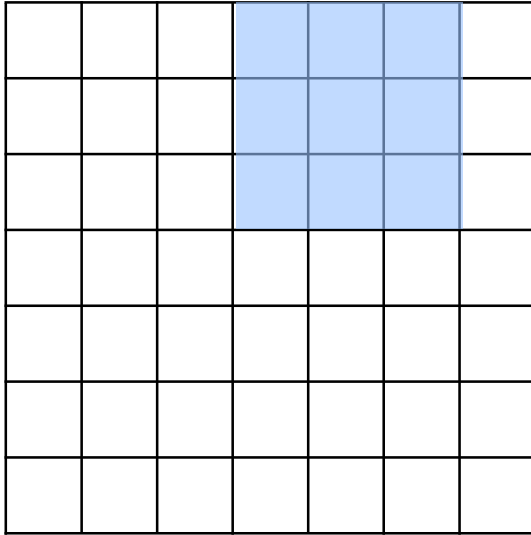
# Convolution Layers



Example:  
 $7 \times 7$  input  
assume  $3 \times 3$  connectivity  
stride 1

- Replicate this column of hidden neurons across space, with some **stride**.

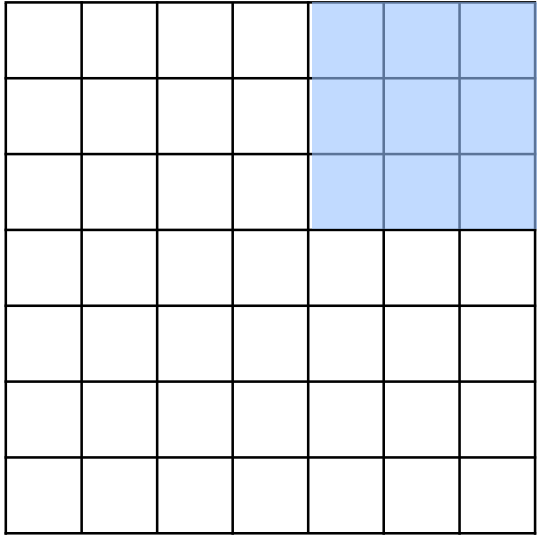
# Convolution Layers



Example:  
 $7 \times 7$  input  
assume  $3 \times 3$  connectivity  
stride 1

- Replicate this column of hidden neurons across space, with some **stride**.

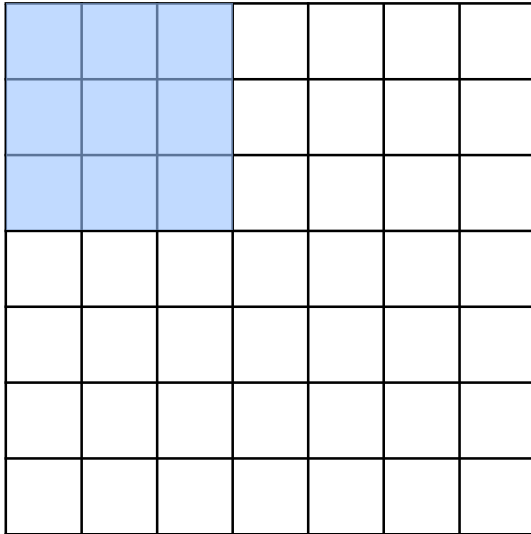
# Convolution Layers



Example:  
 $7 \times 7$  input  
assume  $3 \times 3$  connectivity  
stride 1  
 $\Rightarrow 5 \times 5$  output

- Replicate this column of hidden neurons across space, with some **stride**.

# Convolution Layers

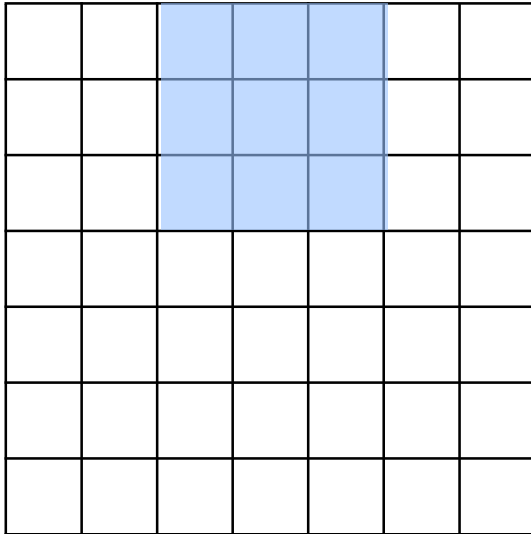


Example:  
 $7 \times 7$  input  
assume  $3 \times 3$  connectivity  
stride 1  
 $\Rightarrow 5 \times 5$  output

What about stride 2?

- Replicate this column of hidden neurons across space, with some **stride**.

# Convolution Layers

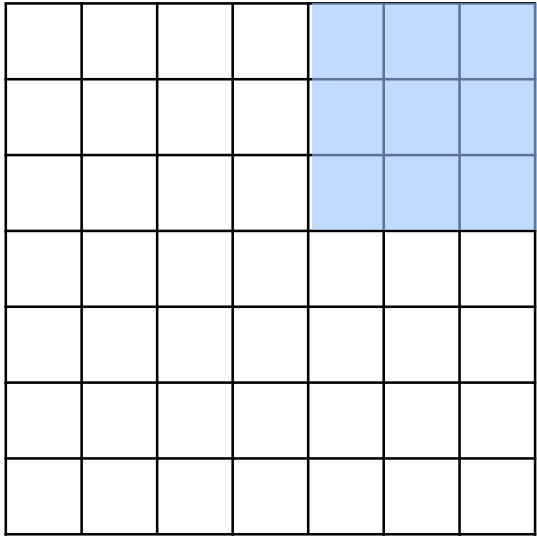


Example:  
 $7 \times 7$  input  
assume  $3 \times 3$  connectivity  
stride 1  
 $\Rightarrow 5 \times 5$  output

What about stride 2?

- Replicate this column of hidden neurons across space, with some **stride**.

# Convolution Layers



Example:

$7 \times 7$  input

assume  $3 \times 3$  connectivity

stride 1

$\Rightarrow 5 \times 5$  output

What about stride 2?

$\Rightarrow 3 \times 3$  output

- Replicate this column of hidden neurons across space, with some **stride**.



# Convolution Layers

0	0	0	0	0				
0								
0								
0								
0								

Example:

$7 \times 7$  input

assume  $3 \times 3$  connectivity

stride 1

$\Rightarrow 5 \times 5$  output

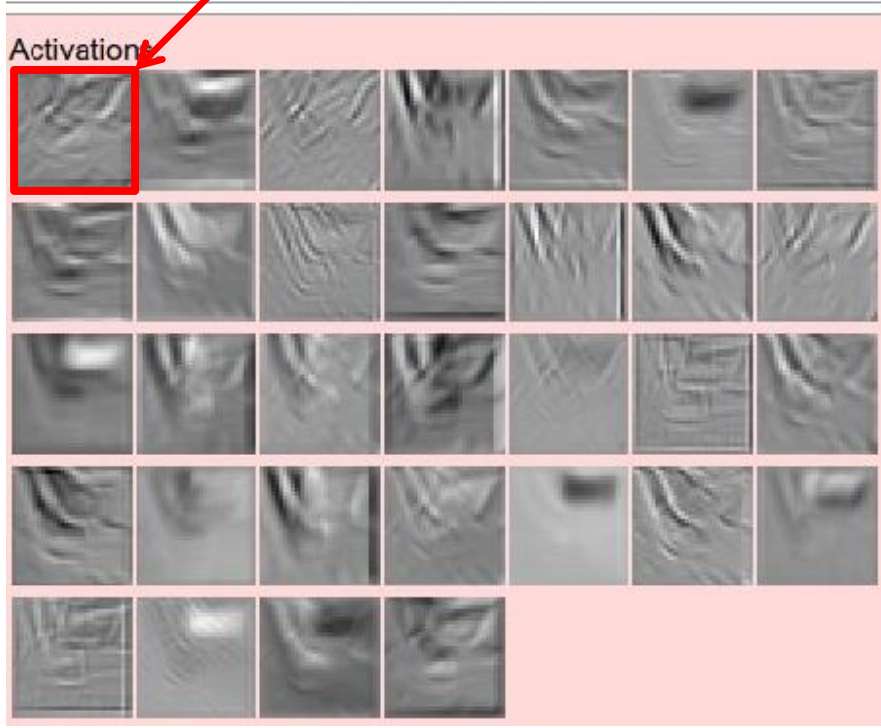
What about stride 2?

$\Rightarrow 3 \times 3$  output

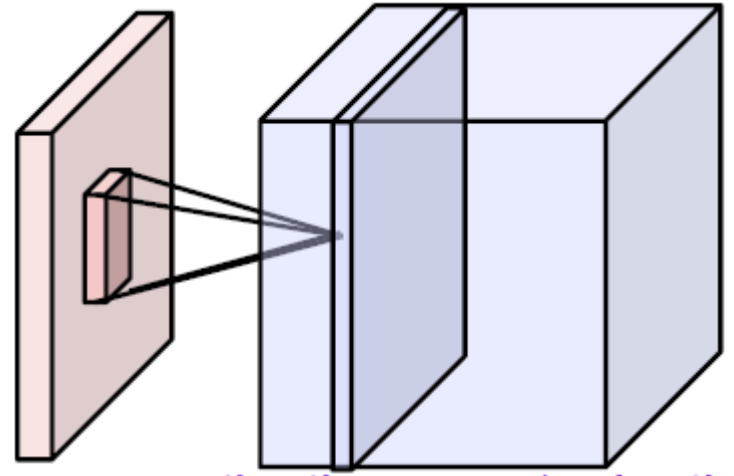
- Replicate this column of hidden neurons across space, with some **stride**.
- In practice, common to zero-pad the border.
  - Preserves the size of the input spatially.

# Activation Maps of Convolutional Filters

Activations:

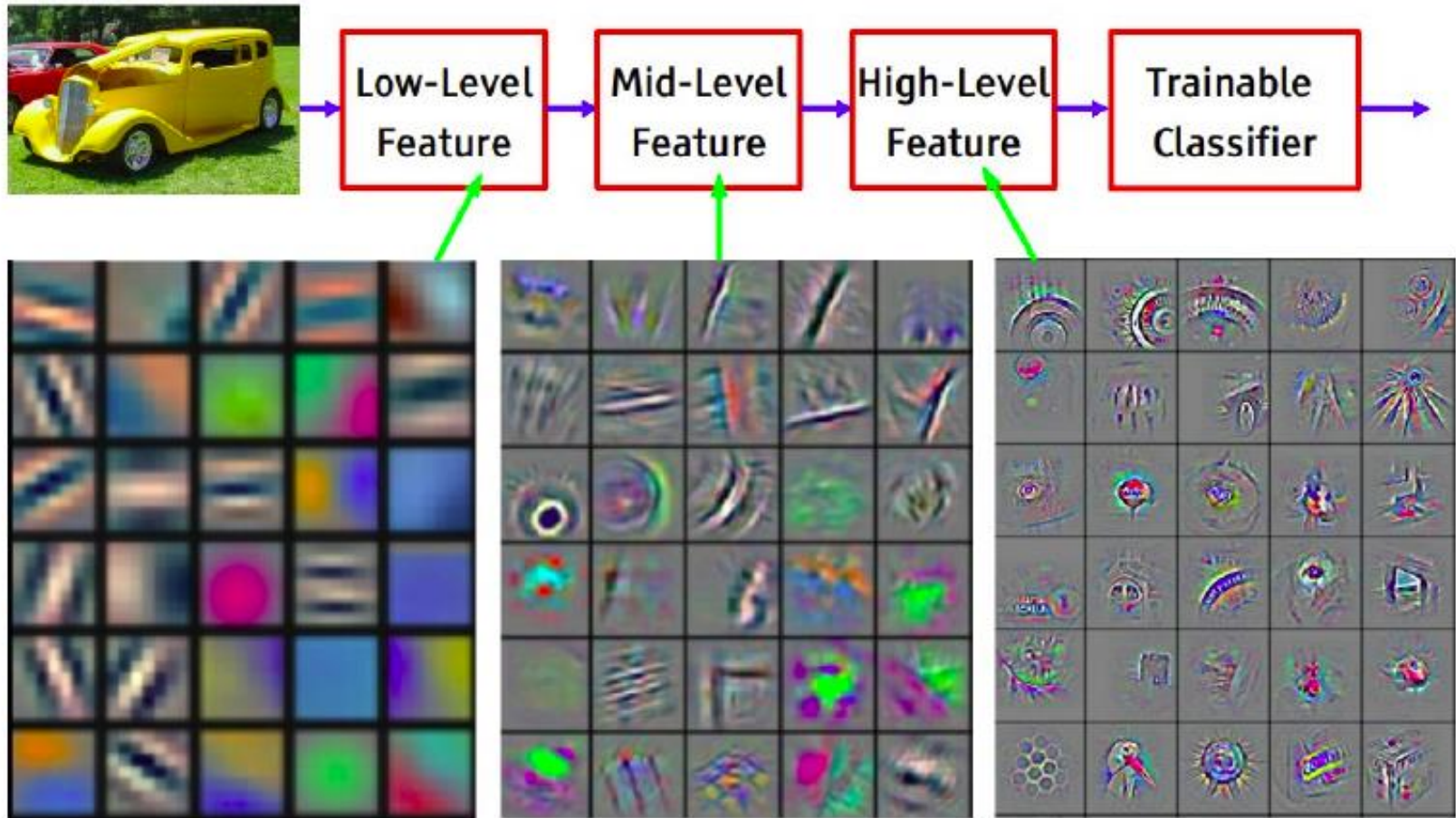


Activation maps



Each activation map is a depth slice through the output volume.

# Effect of Multiple Convolution Layers



Feature visualization of convolutional net trained on ImageNet from [Zeiler & Fergus 2013]

# Commonly Used Nonlinearities

- **Sigmoid**

$$\begin{aligned}g(a) &= \sigma(a) \\ &= \frac{1}{1 + \exp\{-a\}}\end{aligned}$$

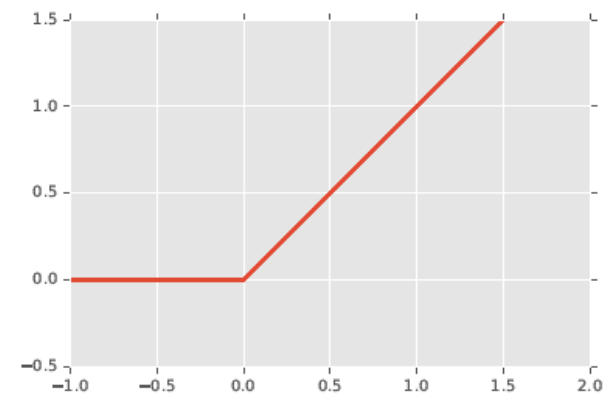
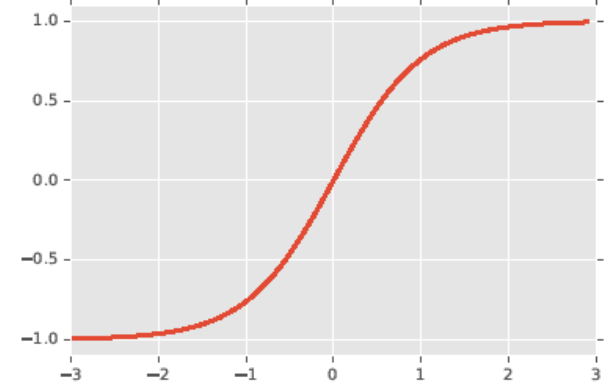
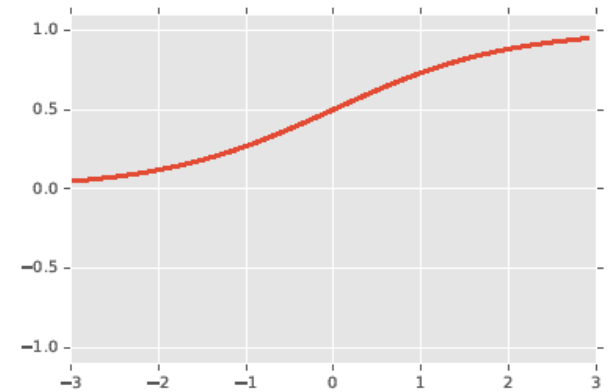
- **Hyperbolic tangent**

$$\begin{aligned}g(a) &= \tanh(a) \\ &= 2\sigma(2a) - 1\end{aligned}$$

- **Rectified linear unit (ReLU)**

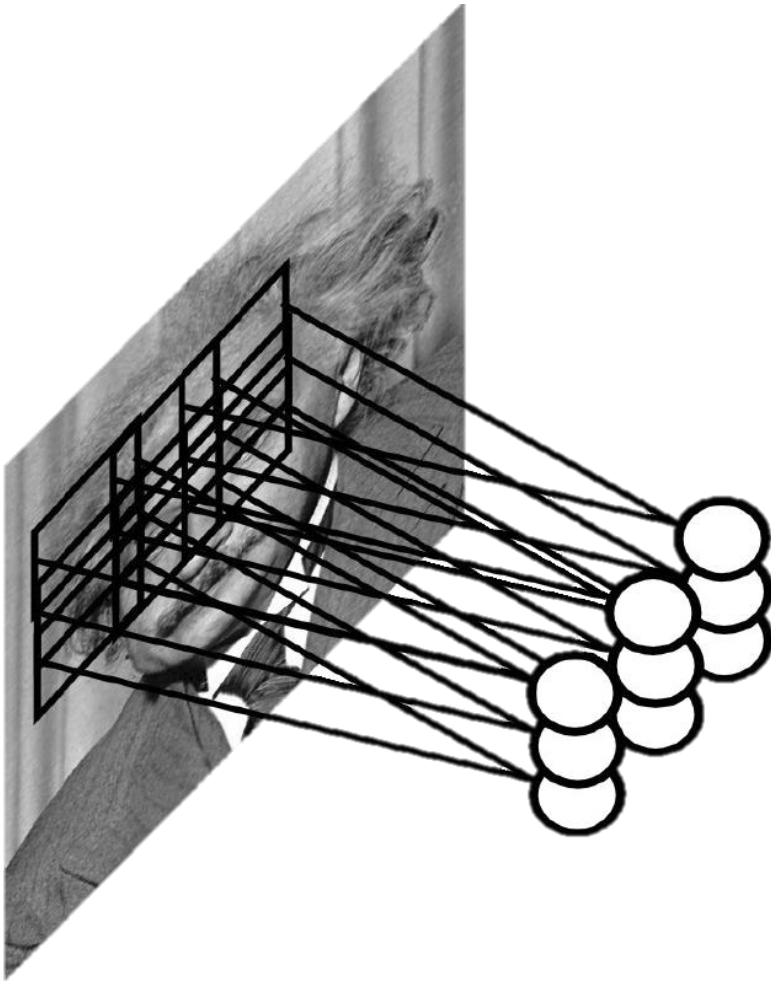
$$g(a) = \max\{0, a\}$$

**Preferred option for deep networks**



# Convolutional Networks: Intuition

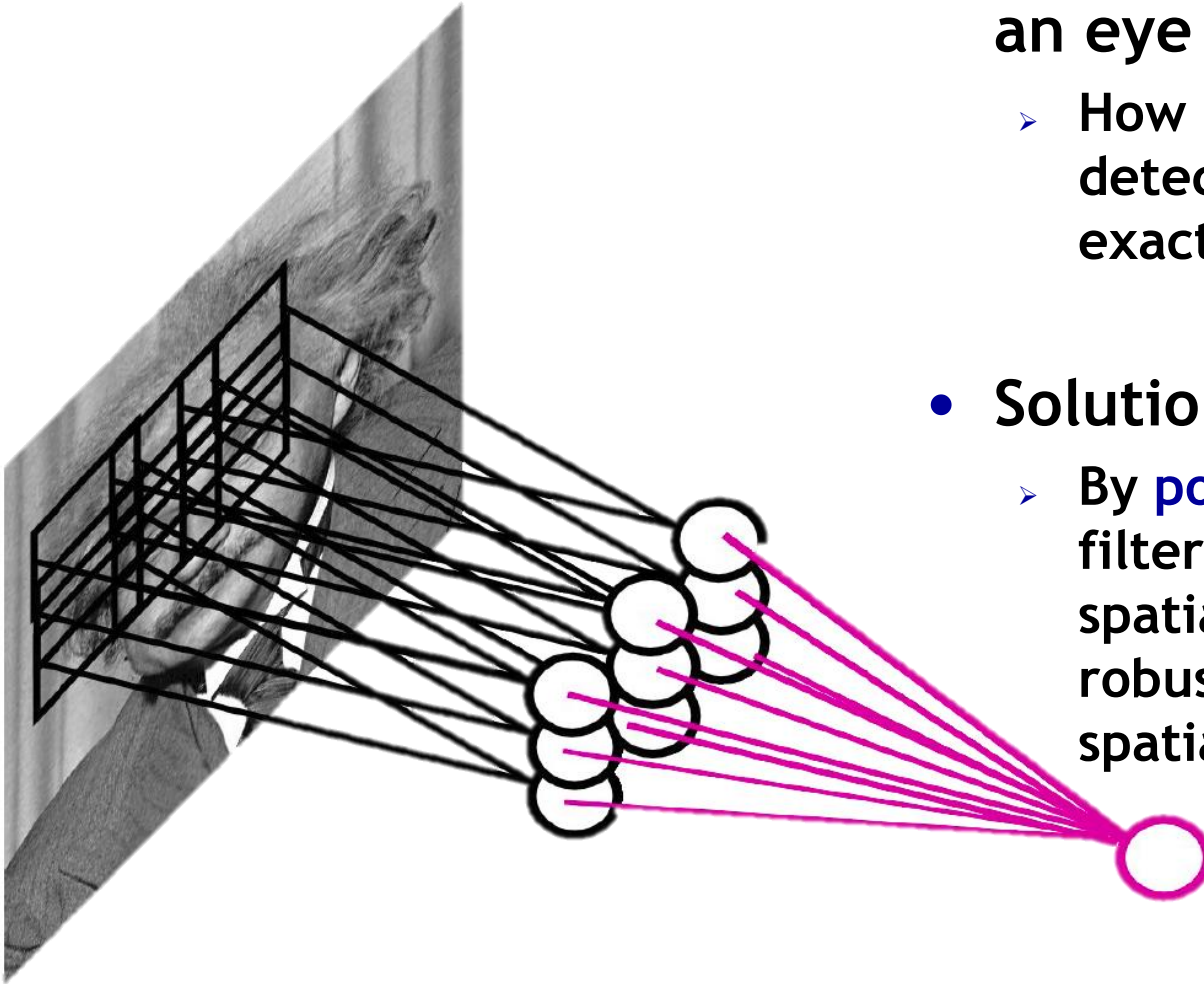
- Let's assume the filter is an eye detector
  - How can we make the detection robust to the exact location of the eye?



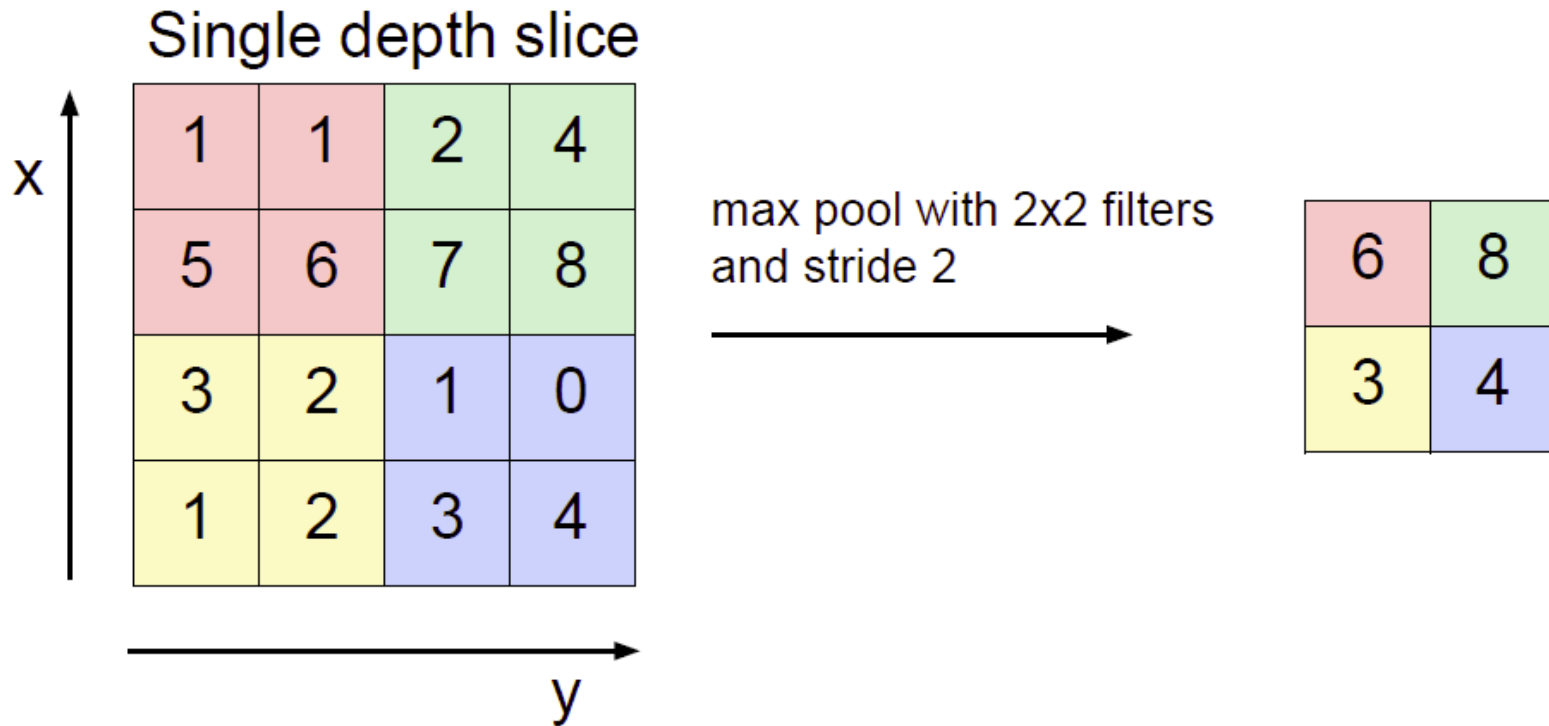


# Convolutional Networks: Intuition

- Let's assume the filter is an eye detector
  - How can we make the detection robust to the exact location of the eye?
- Solution:
  - By **pooling** (e.g., max or avg) filter responses at different spatial locations, we gain robustness to the exact spatial location of features.



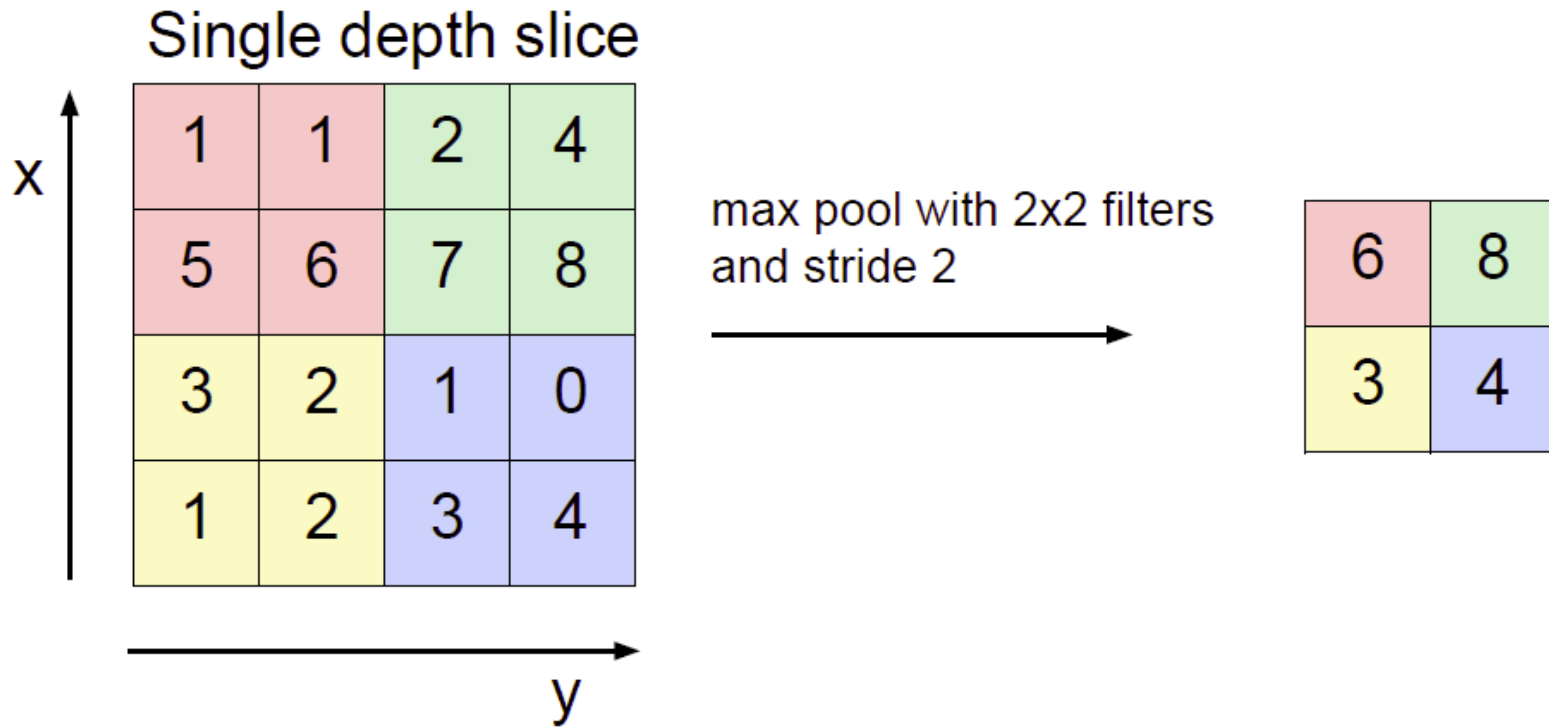
# Max Pooling



- **Effect:**

- Make the representation smaller without losing too much information
- Achieve robustness to translations

# Max Pooling



- **Note**

- Pooling happens independently across each slice, preserving the number of slices.



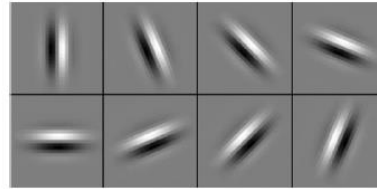
# Compare: SIFT Descriptor

Lowé  
[IJCV 2004]

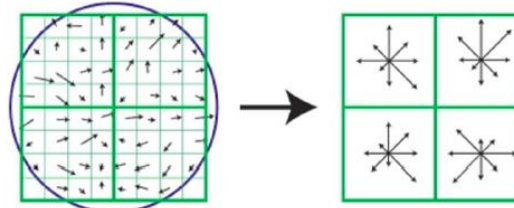
Image  
Pixels



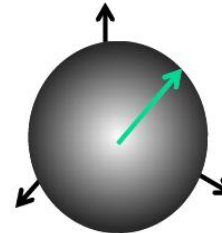
Apply  
oriented filters



Spatial pool  
(Sum)



Normalize to  
unit length

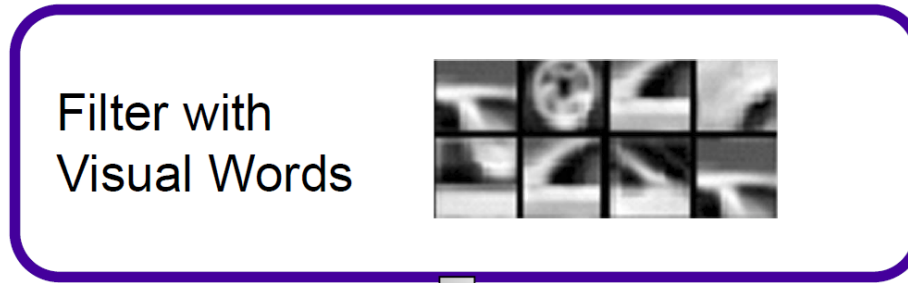


Feature  
Vector

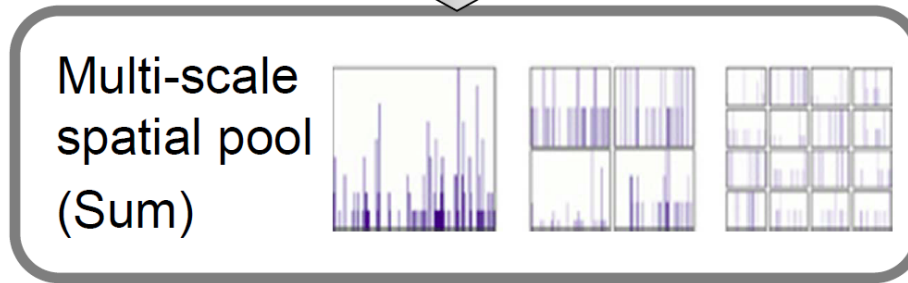
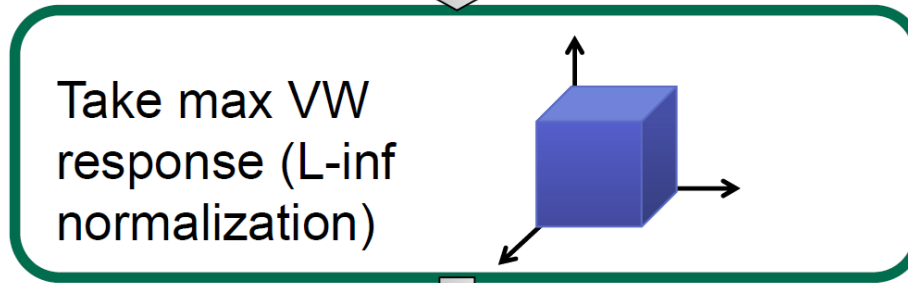


# Compare: Spatial Pyramid Matching

SIFT features →



Lazebnik,  
Schmid,  
Ponce  
[CVPR 2006]

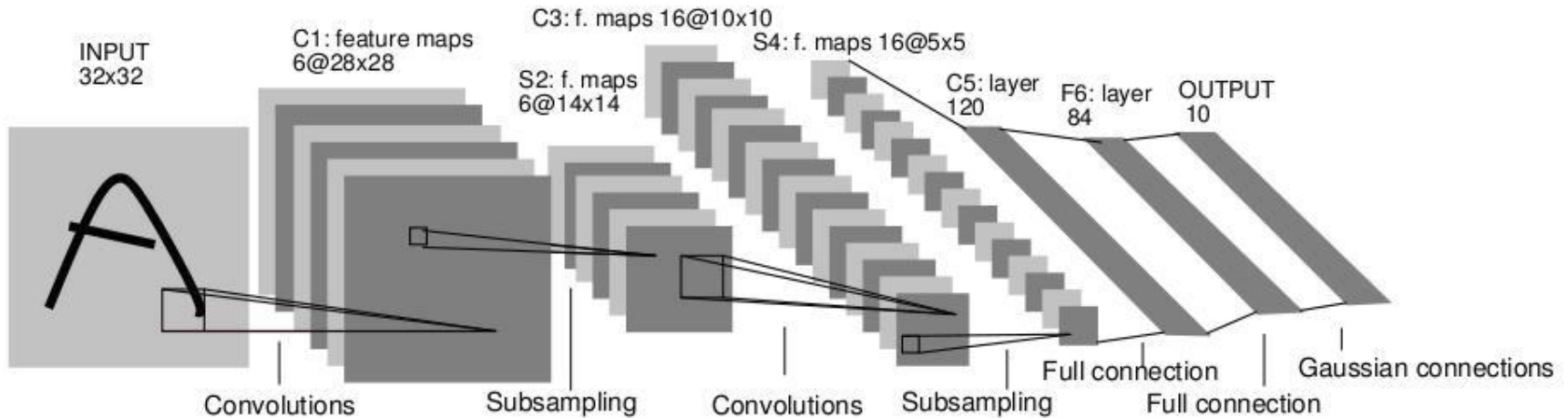


→ Global image descriptor

# Topics of This Lecture

- Deep Learning
  - Motivation
- Convolutional Neural Networks
  - Convolutional Layers
  - Pooling Layers
  - Nonlinearities
- **CNN Architectures**
  - **LeNet**
  - **AlexNet**
  - **VGGNet**
  - **GoogLeNet**
- Applications

# CNN Architectures: LeNet (1998)



- Early convolutional architecture
  - 2 Convolutional layers, 2 pooling layers
  - Fully-connected NN layers for classification
  - Successfully used for handwritten digit recognition (MNIST)

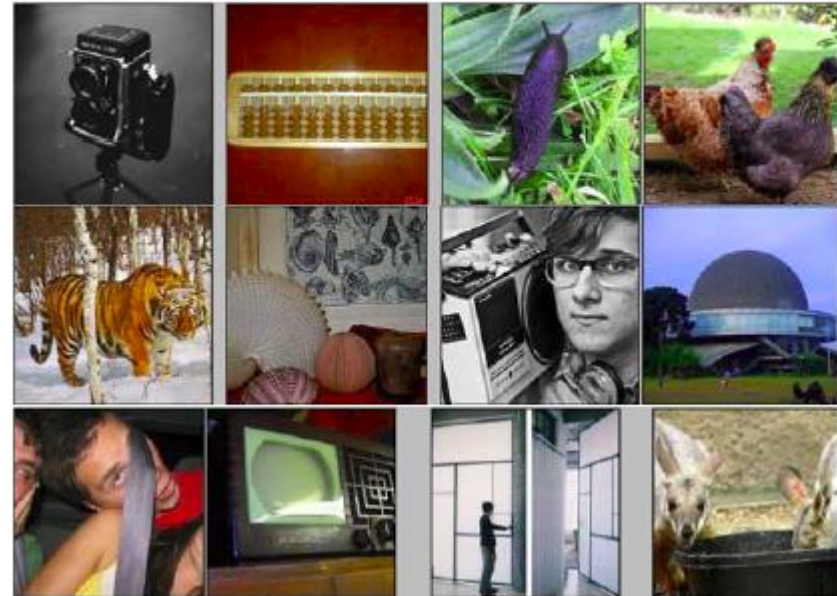
Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner, [Gradient-based learning applied to document recognition](#), Proceedings of the IEEE 86(11): 2278-2324, 1998.

# ImageNet Challenge 2012

- ImageNet

- ~14M labeled internet images
- 20k classes
- Human labels via Amazon Mechanical Turk

IM  GENET

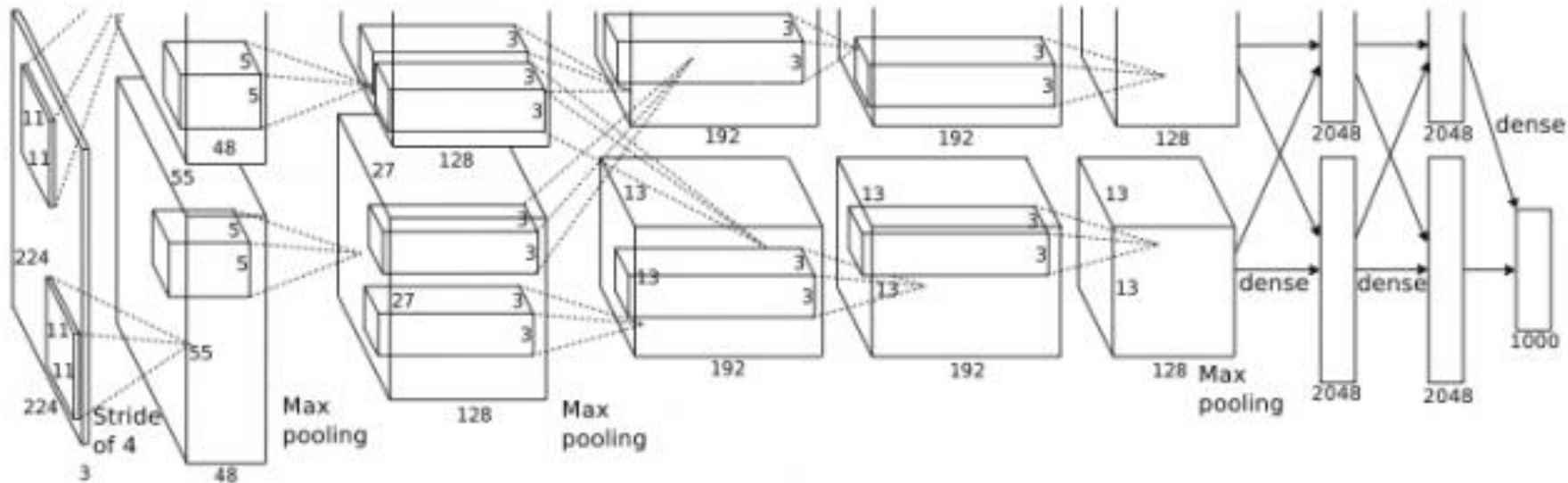


- Challenge (ILSVRC)

- 1.2 million training images
- 1000 classes
- Goal: Predict ground-truth class within top-5 responses
- Currently one of the top benchmarks in Computer Vision

[Deng et al., CVPR'09]

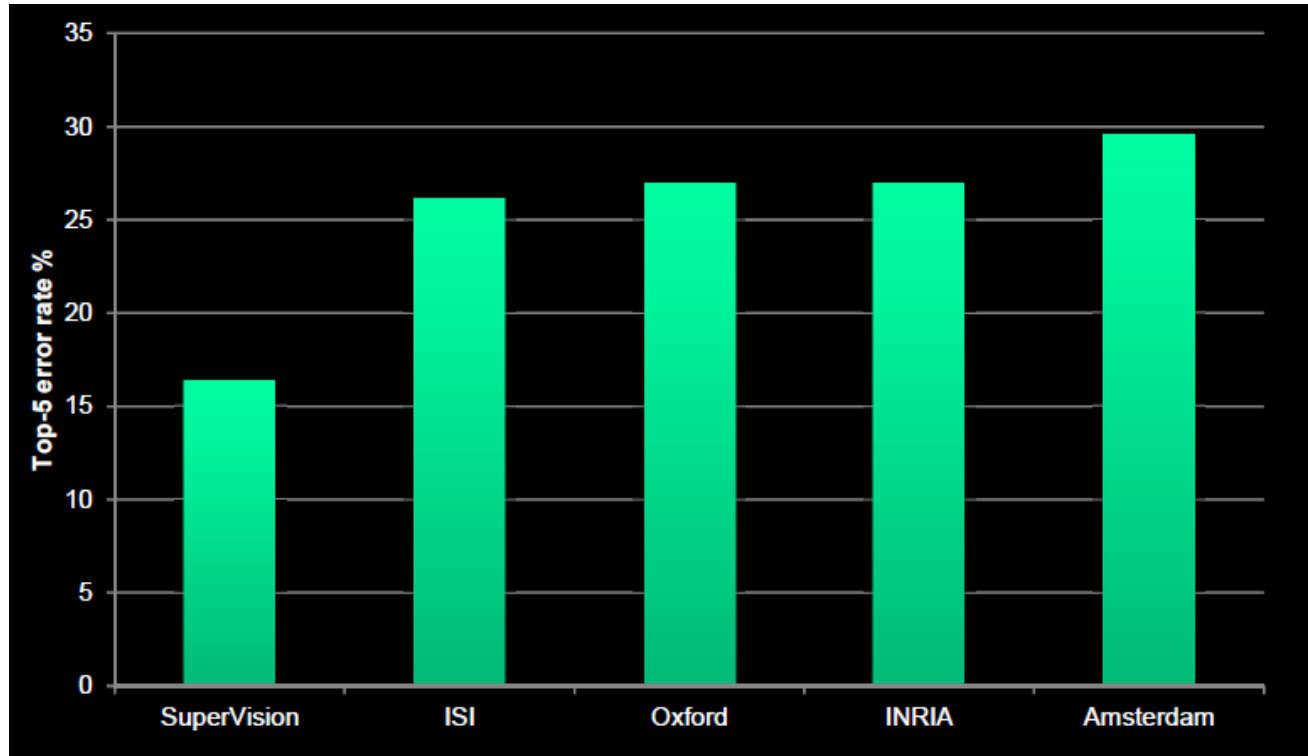
# CNN Architectures: AlexNet (2012)



- **Similar framework as LeNet, but**
  - **Bigger model (7 hidden layers, 650k units, 60M parameters)**
  - **More data ( $10^6$  images instead of  $10^3$ )**
  - **GPU implementation**
  - **Better regularization and up-to-date tricks for training (Dropout)**

A. Krizhevsky, I. Sutskever, and G. Hinton, [ImageNet Classification with Deep Convolutional Neural Networks](#), NIPS 2012.

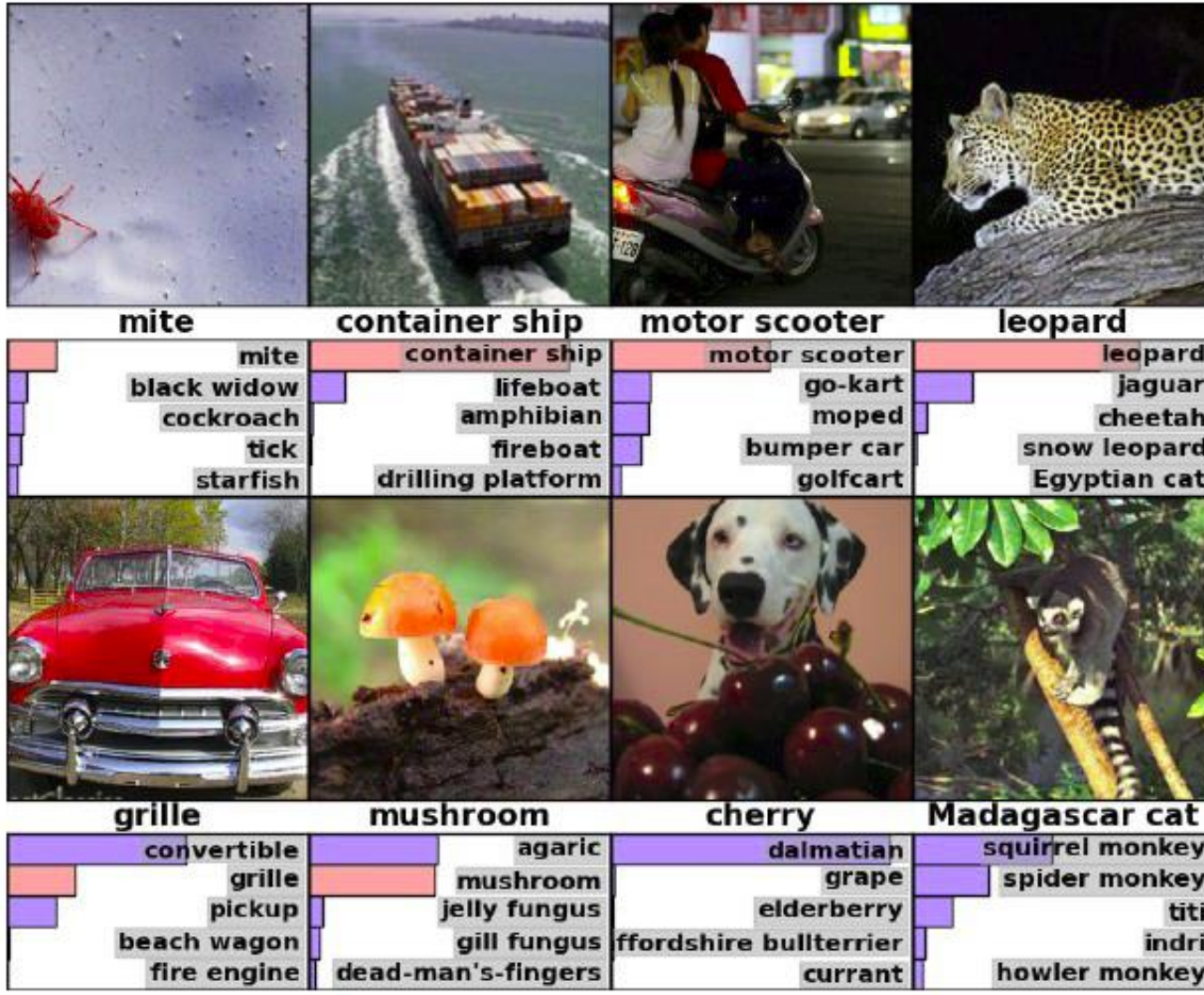
# ILSVRC 2012 Results



- AlexNet almost halved the error rate
  - 16.4% error (top-5) vs. 26.2% for the next best approach
  - ⇒ A revolution in Computer Vision
  - Acquired by Google in Jan '13, deployed in Google+ in May '13



# AlexNet Results





# AlexNet Results

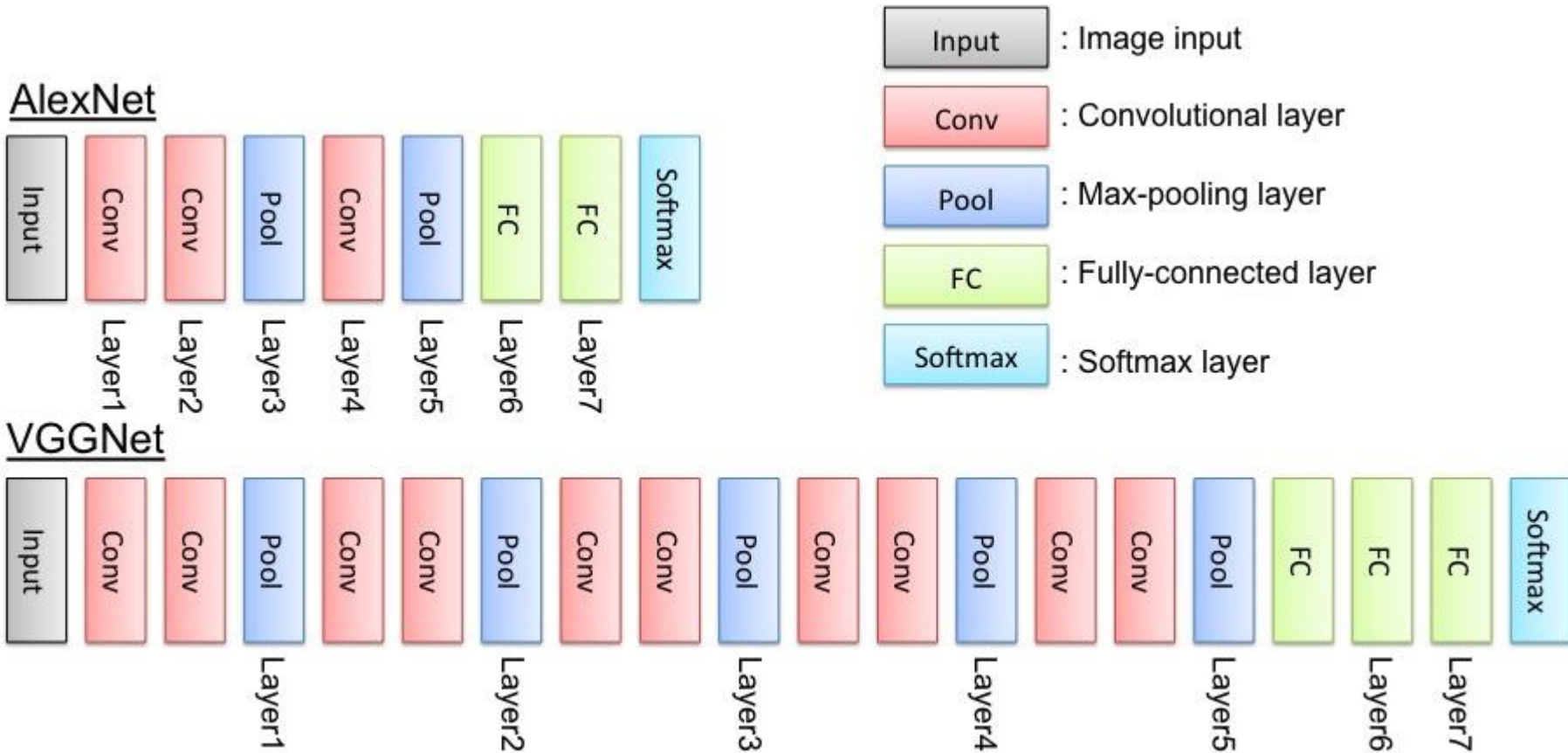


Test image



Retrieved images

# CNN Architectures: VGGNet (2014/15)



K. Simonyan, A. Zisserman, [Very Deep Convolutional Networks for Large-Scale Image Recognition](#), ICLR 2015

# CNN Architectures: VGGNet (2014/15)

- Main ideas

- Deeper network
- Stacked convolutional layers with smaller filters (+ nonlinearity)
- Detailed evaluation of all components

- Results

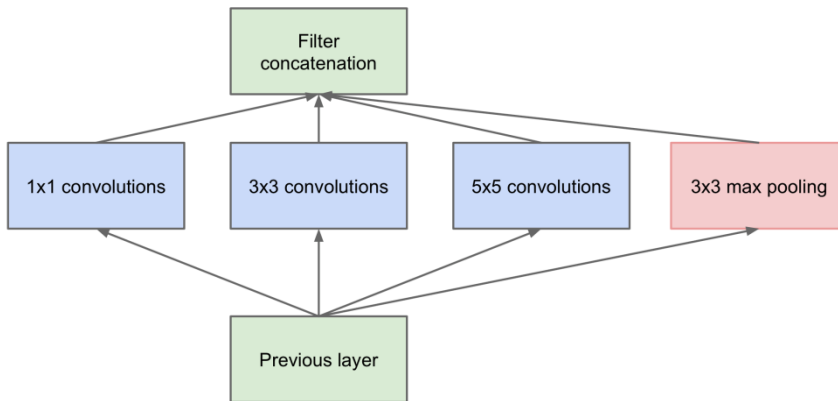
- Improved ILSVRC top-5 error rate to 6.7%.

ConvNet Configuration					
A	A-LRN	B	C	D	E
11 weight layers	11 weight layers	13 weight layers	16 weight layers	16 weight layers	19 weight layers
input (224 × 224 RGB image)					
conv3-64	conv3-64 <b>LRN</b>	conv3-64 <b>conv3-64</b>	conv3-64 conv3-64	conv3-64 conv3-64	conv3-64 conv3-64
maxpool					
conv3-128	conv3-128	conv3-128 <b>conv3-128</b>	conv3-128 conv3-128	conv3-128 conv3-128	conv3-128 conv3-128
maxpool					
conv3-256 conv3-256	conv3-256 conv3-256	conv3-256 conv3-256	conv3-256 conv3-256 <b>conv1-256</b>	conv3-256 conv3-256 <b>conv3-256</b>	conv3-256 conv3-256 conv3-256 <b>conv3-256</b>
maxpool					
conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512 <b>conv1-512</b>	conv3-512 conv3-512 <b>conv3-512</b>	conv3-512 conv3-512 conv3-512 <b>conv3-512</b>
maxpool					
conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512 <b>conv1-512</b>	conv3-512 conv3-512 <b>conv3-512</b>	conv3-512 conv3-512 conv3-512 <b>conv3-512</b>
maxpool					
				<b>Mainly used</b>	
FC-4096					
FC-4096					
FC-1000					
soft-max					

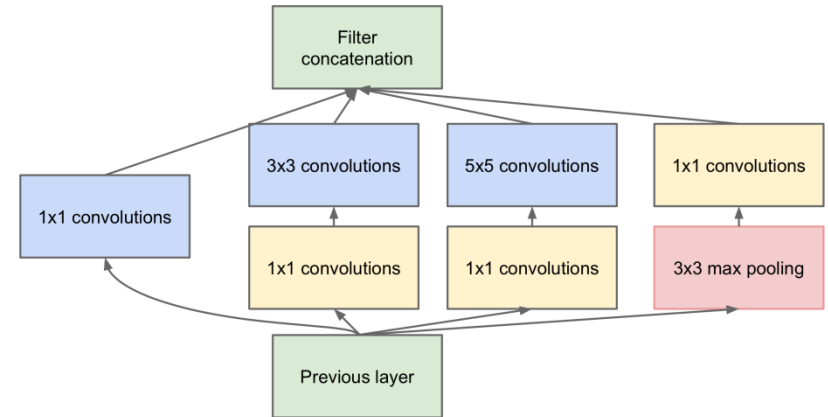
# Comparison: AlexNet vs. VGGNet

- **Receptive fields in the first layer**
  - AlexNet:  $11 \times 11$ , stride 4
  - Zeiler & Fergus:  $7 \times 7$ , stride 2
  - VGGNet:  $3 \times 3$ , stride 1
- **Why that?**
  - If you stack three  $3 \times 3$  on top of another  $3 \times 3$  layer, you effectively get a  $5 \times 5$  receptive field.
  - With three  $3 \times 3$  layers, the receptive field is already  $7 \times 7$ .
  - But much fewer parameters:  $3 \cdot 3^2 = 27$  instead of  $7^2 = 49$ .
  - In addition, non-linearities in-between  $3 \times 3$  layers for additional discriminativity.

# CNN Architectures: GoogLeNet (2014)



(a) Inception module, naïve version



(b) Inception module with dimension reductions

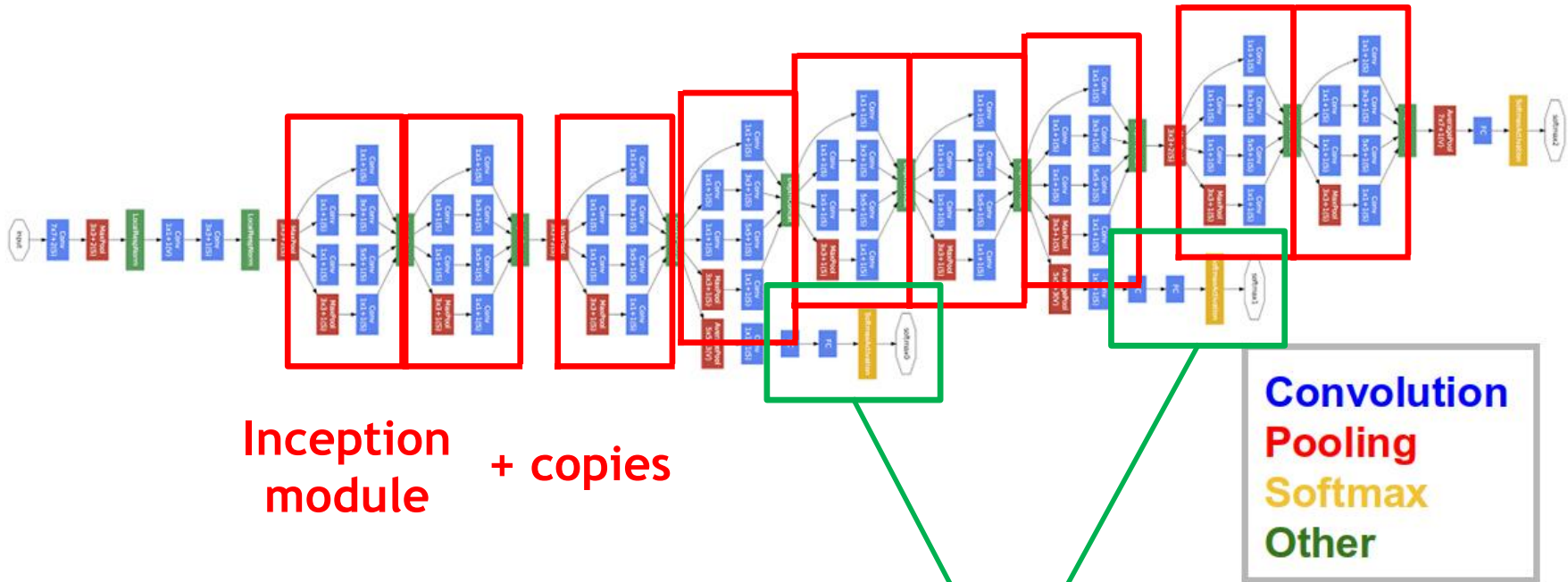
- **Main ideas**

- “Inception” module as modular component
- Learns filters at several scales within each module

C. Szegedy, W. Liu, Y. Jia, et al, [Going Deeper with Convolutions](#), arXiv:1409.4842, 2014.



# GoogLeNet Visualization



Inception module + copies

Auxiliary classification outputs for training the lower layers (deprecated)

Convolution  
Pooling  
Softmax  
Other

# Results on ILSVRC

Method	top-1 val. error (%)	top-5 val. error (%)	top-5 test error (%)
VGG (2 nets, multi-crop & dense eval.)	<b>23.7</b>	<b>6.8</b>	<b>6.8</b>
VGG (1 net, multi-crop & dense eval.)	24.4	7.1	7.0
VGG (ILSVRC submission, 7 nets, dense eval.)	24.7	7.5	7.3
GoogLeNet (Szegedy et al., 2014) (1 net)	-	7.9	
GoogLeNet (Szegedy et al., 2014) (7 nets)	-	<b>6.7</b>	
MSRA (He et al., 2014) (11 nets)	-	-	8.1
MSRA (He et al., 2014) (1 net)	27.9	9.1	9.1
Clarifai (Russakovsky et al., 2014) (multiple nets)	-	-	11.7
Clarifai (Russakovsky et al., 2014) (1 net)	-	-	12.5
Zeiler & Fergus (Zeiler & Fergus, 2013) (6 nets)	36.0	14.7	14.8
Zeiler & Fergus (Zeiler & Fergus, 2013) (1 net)	37.5	16.0	16.1
OverFeat (Sermanet et al., 2014) (7 nets)	34.0	13.2	13.6
OverFeat (Sermanet et al., 2014) (1 net)	35.7	14.2	-
Krizhevsky et al. (Krizhevsky et al., 2012) (5 nets)	38.1	16.4	16.4
Krizhevsky et al. (Krizhevsky et al., 2012) (1 net)	40.7	18.2	-

- **VGGNet and GoogLeNet perform at similar level**
  - **Comparison: human performance ~5% [Karpathy]**

<http://karpathy.github.io/2014/09/02/what-i-learned-from-competing-against-a-convnet-on-imagenet/>





# Newest Development: Residual Networks

AlexNet, 8 layers  
(ILSVRC 2012)



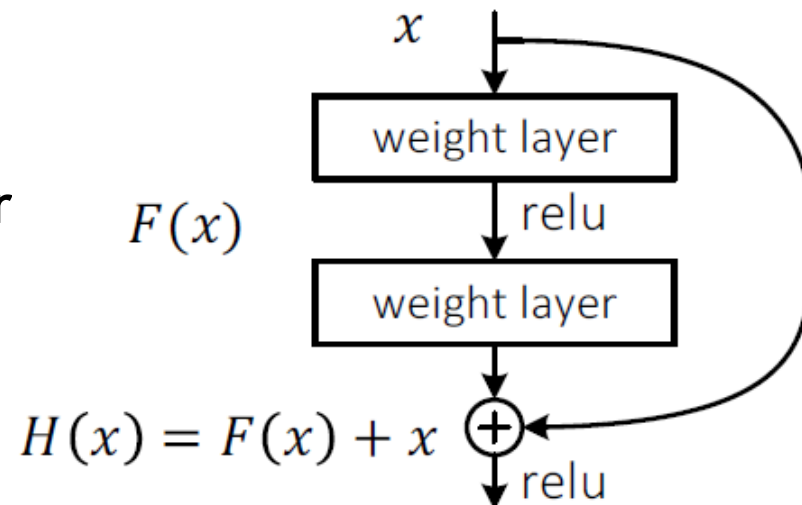
VGG, 19 layers  
(ILSVRC 2014)



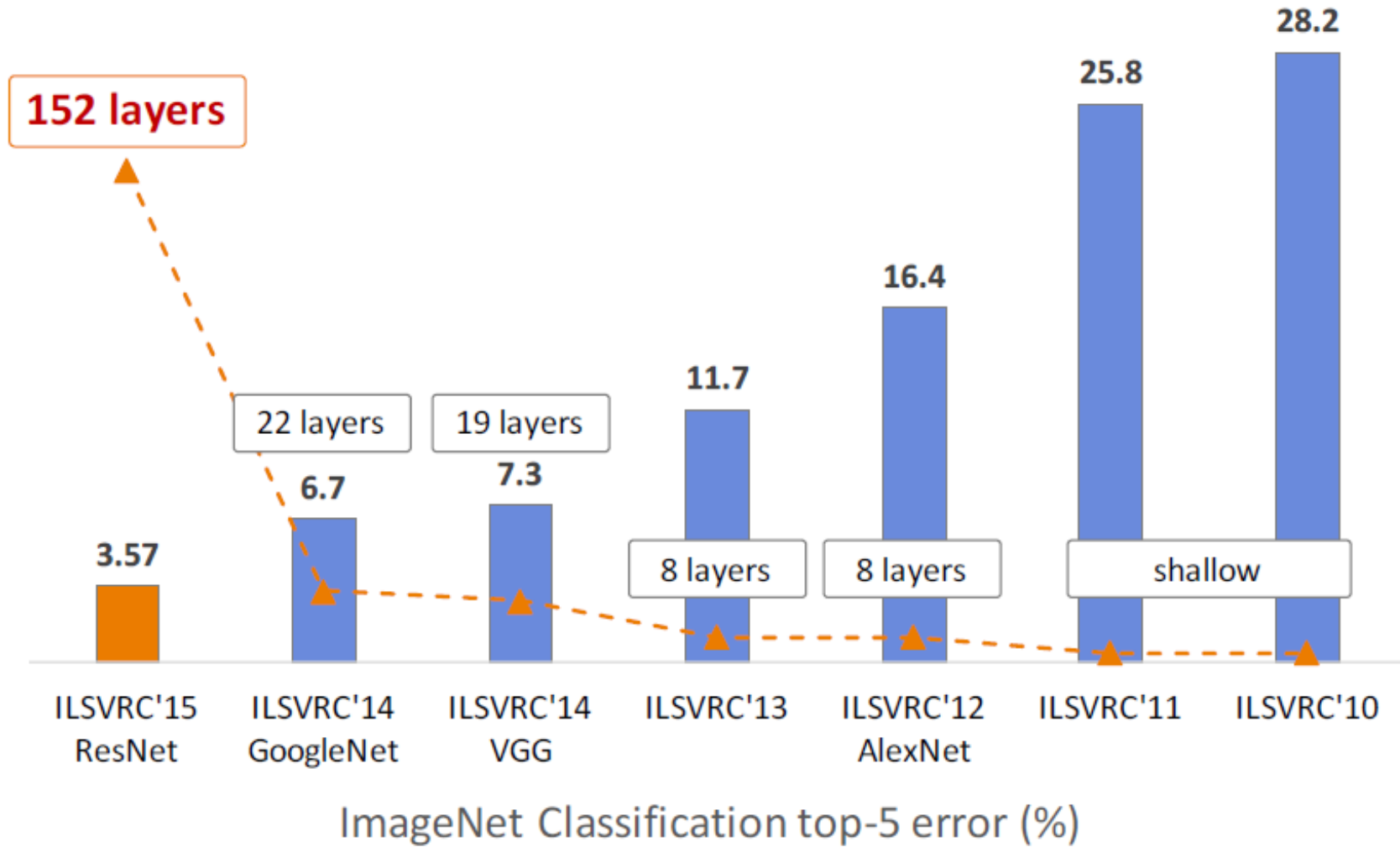
ResNet, 152 layers  
(ILSVRC 2015)

- Core component

- Skip connections bypassing each layer
- Better propagation of gradients to the deeper layers



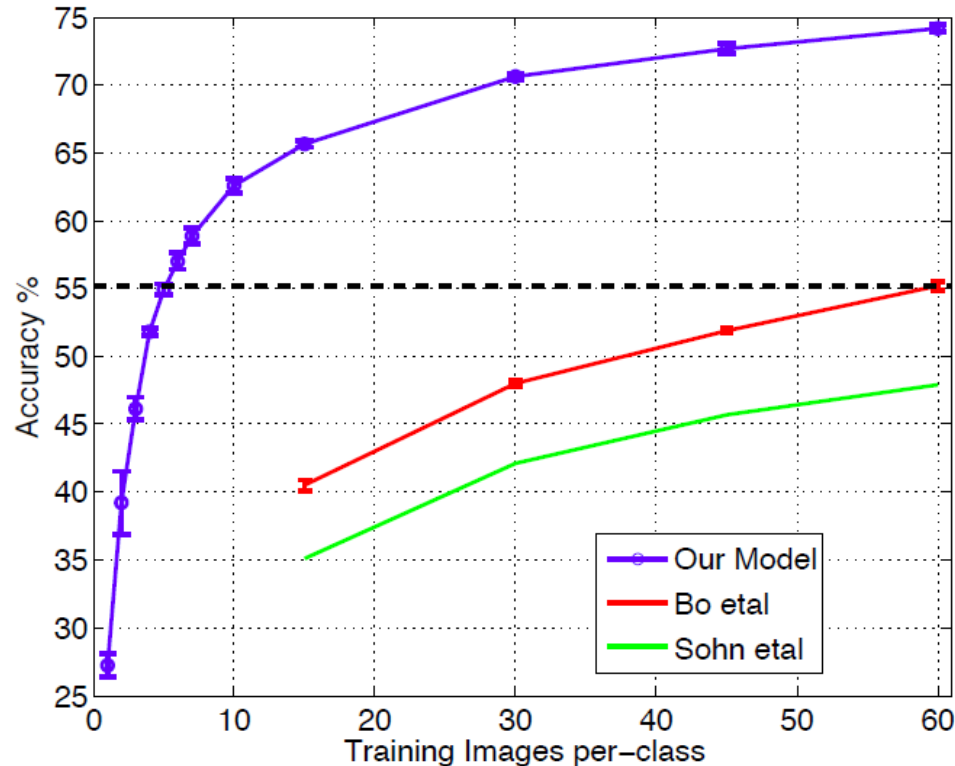
# ImageNet Performance



# Topics of This Lecture

- Deep Learning
  - Motivation
- Convolutional Neural Networks
  - Convolutional Layers
  - Pooling Layers
  - Nonlinearities
- CNN Architectures
  - LeNet
  - AlexNet
  - VGGNet
  - GoogLeNet
- **Applications**

# The Learned Features are Generic



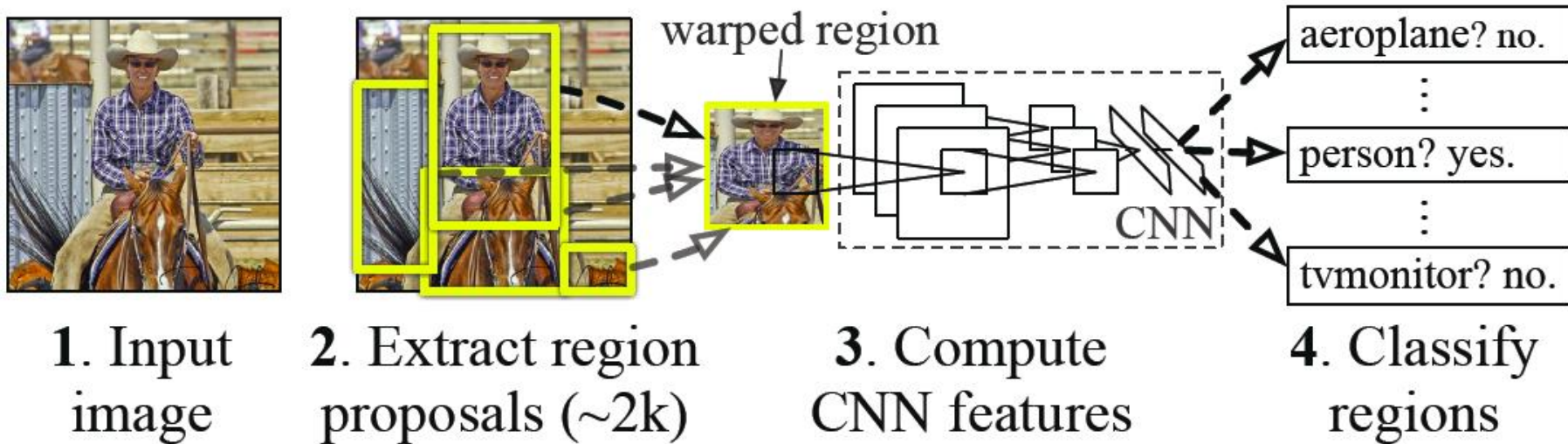
state of the art  
level (pre-CNN)

- **Experiment: feature transfer**

- Train network on ImageNet
  - Chop off last layer and train classification layer on CalTech256
- ⇒ State of the art accuracy already with only 6 training images

# Other Tasks: Detection

## R-CNN: *Regions with CNN features*



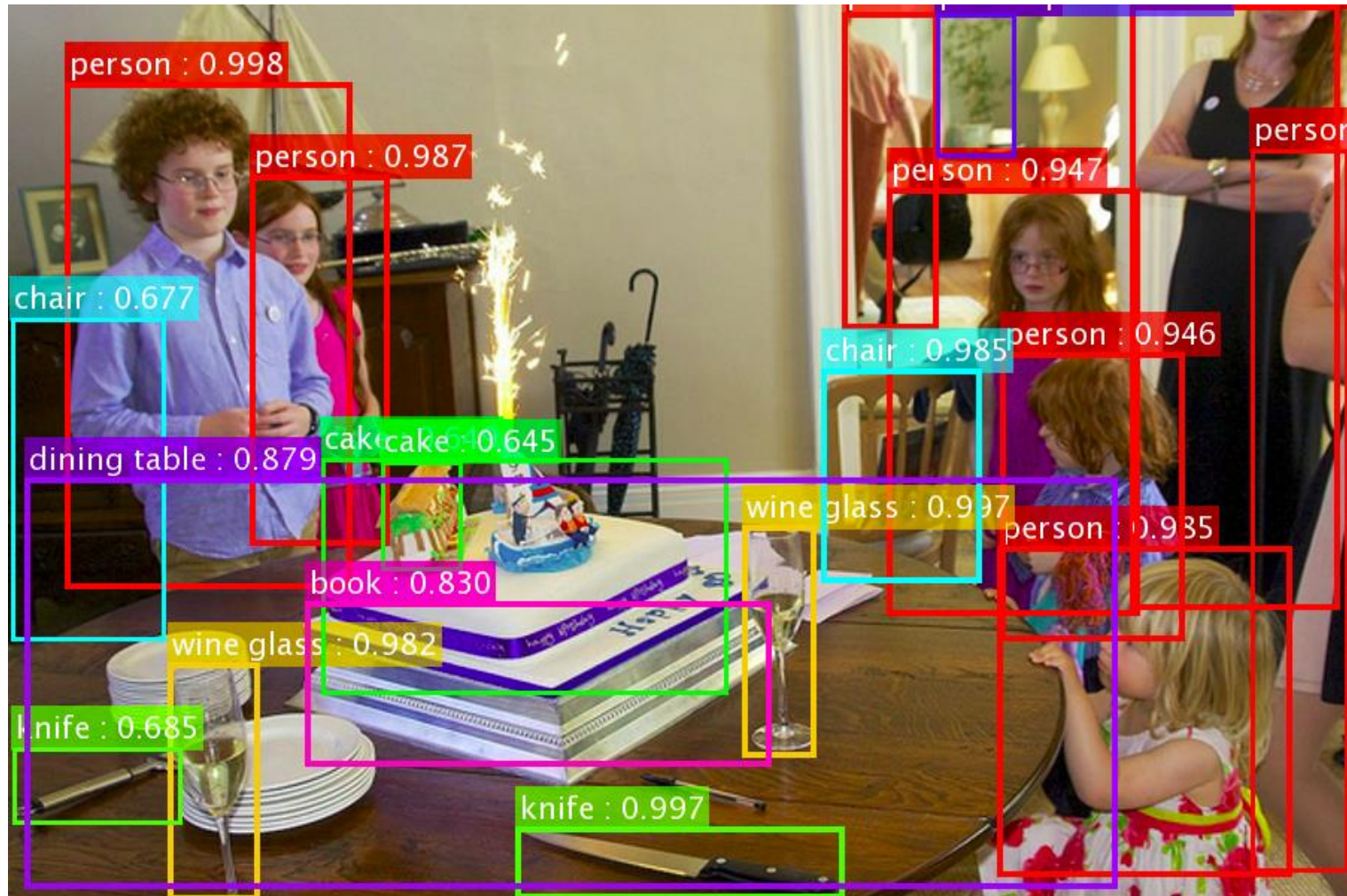
- **Results on PASCAL VOC Detection benchmark**

- **Pre-CNN state of the art: 35.1% mAP [Uijlings et al., 2013]**
  - 33.4% mAP DPM
  - **R-CNN: 53.7% mAP**

R. Girshick, J. Donahue, T. Darrell, and J. Malik, [Rich Feature Hierarchies for Accurate Object Detection and Semantic Segmentation](#), CVPR 2014

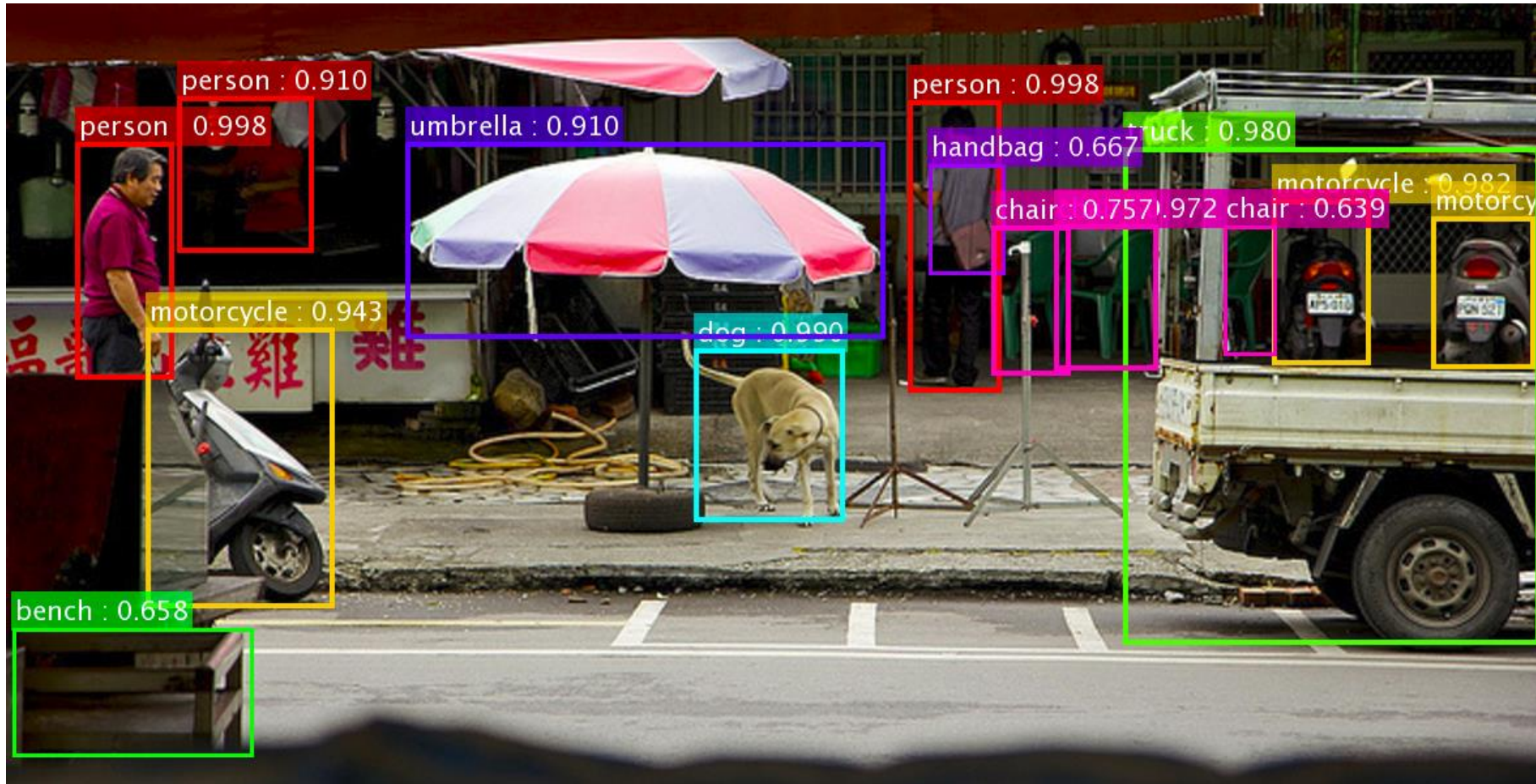


# Faster R-CNN (based on ResNets)



K. He, X. Zhang, S. Ren, J. Sun, [Deep Residual Learning for Image Recognition](#), CVPR 2016.

# Faster R-CNN (based on ResNets)



K. He, X. Zhang, S. Ren, J. Sun, [Deep Residual Learning for Image Recognition](#), CVPR 2016.







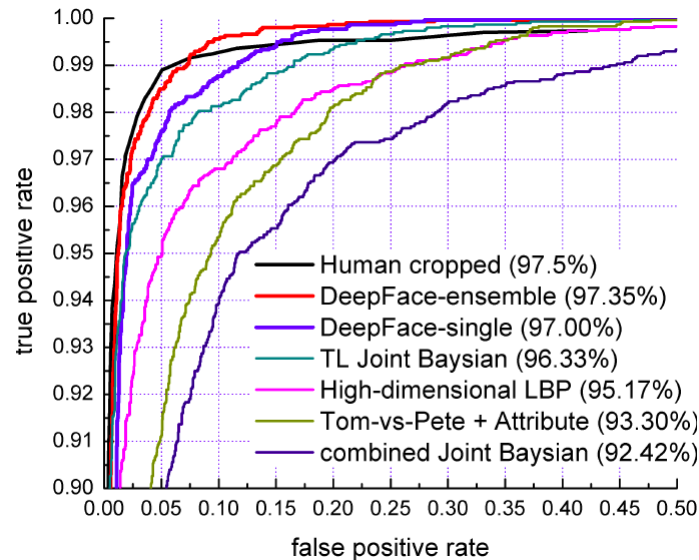
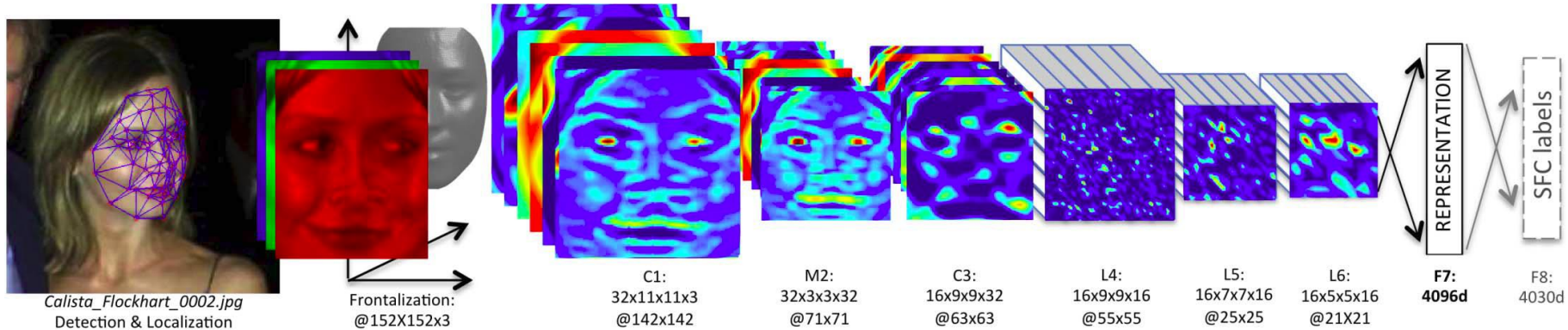
# Semantic Segmentation



[Pohlen, Hermans, Mathias, Leibe, arXiv 2016]

- **More recent results**
  - Based on an extension of ResNets

# Other Tasks: Face Verification



Y. Taigman, M. Yang, M. Ranzato, L. Wolf, [DeepFace: Closing the Gap to Human-Level Performance in Face Verification](#), CVPR 2014

# Commercial Recognition Services

- E.g., **clarifai**



Try it out with your own media

Upload an image or video file under 100mb or give us a direct link to a file on the web.

Paste a url here... ENGLISH ▼

USE THE URL CHOOSE A FILE INSTEAD

\*By using the demo you agree to our terms of service

- **Be careful when taking test images from Google Search**
  - Chances are they may have been seen in the training set...

# Commercial Recognition Services



coffee croissant beverage  
morning breakfast food



night bridge city  
suspension bridge river



winter snow cold mammal  
dog arctic

clarifai

# References and Further Reading

- **LeNet**

- Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner, [Gradient-based learning applied to document recognition](#), Proceedings of the IEEE 86(11): 2278-2324, 1998.

- **AlexNet**

- A. Krizhevsky, I. Sutskever, and G. Hinton, [ImageNet Classification with Deep Convolutional Neural Networks](#), NIPS 2012.

- **VGGNet**

- K. Simonyan, A. Zisserman, [Very Deep Convolutional Networks for Large-Scale Image Recognition](#), ICLR 2015

- **GoogLeNet**

- C. Szegedy, W. Liu, Y. Jia, et al, [Going Deeper with Convolutions](#), arXiv:1409.4842, 2014.