

# Supervised Deep Learning

Marc'Aurelio Ranzato



Facebook A.I. Research

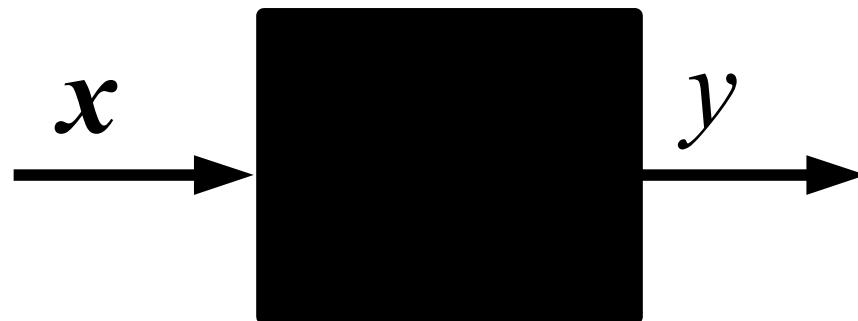
# Supervised Learning

$\{(x^i, y^i), i=1 \dots P\}$  training dataset

$x^i$  i-th input training example

$y^i$  i-th target label

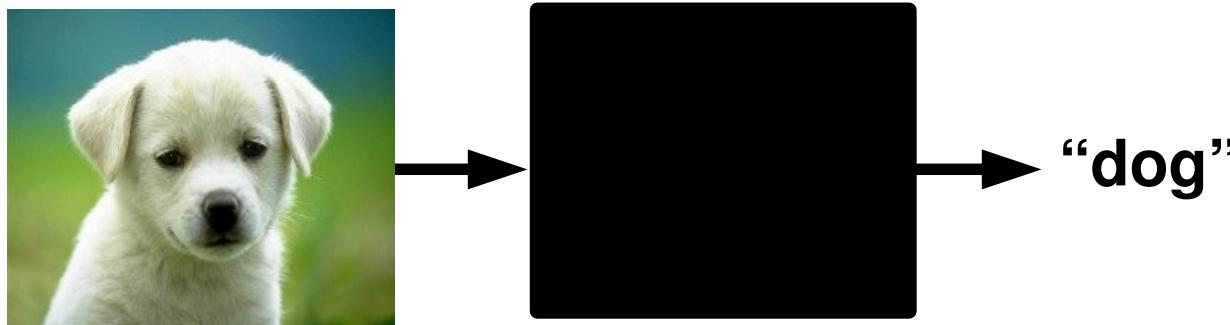
$P$  number of training examples



Goal: predict the target label of unseen inputs.

# Supervised Learning: Examples

## Classification



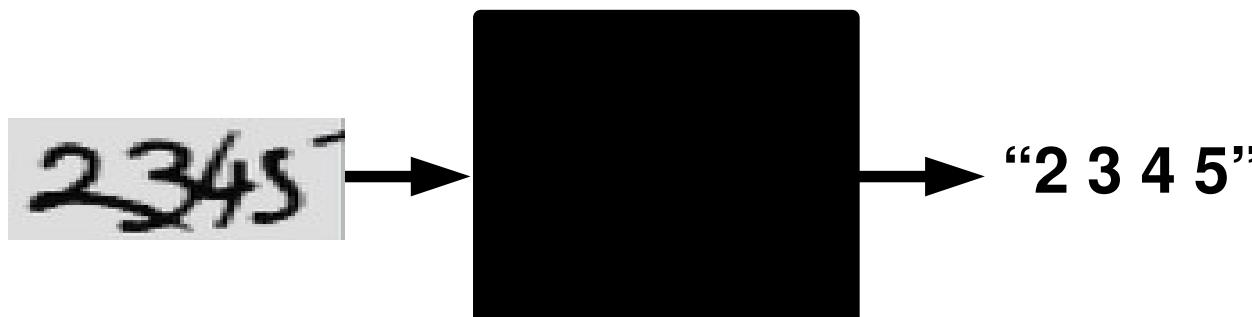
classification

## Denoising



regression

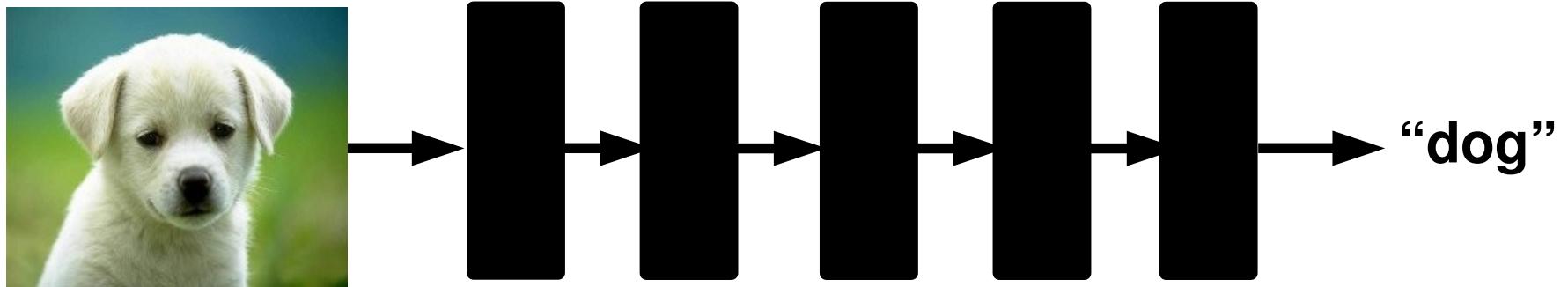
## OCR



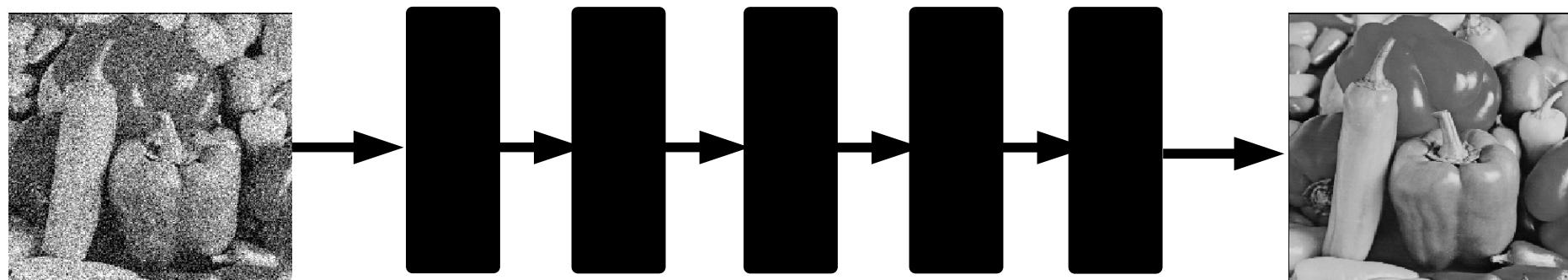
structured prediction

# Supervised Deep Learning

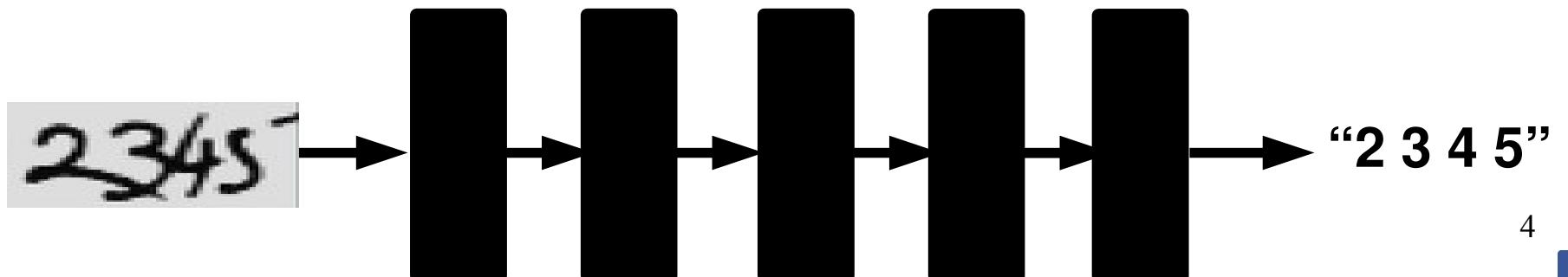
## Classification



## Denoising



## OCR



# Outline

- Supervised Neural Networks
- Convolutional Neural Networks
- Examples
- Tips

# Neural Networks

Assumptions (for the next few slides):

- The input image is vectorized (disregard the spatial layout of pixels)
- The target label is discrete (classification)

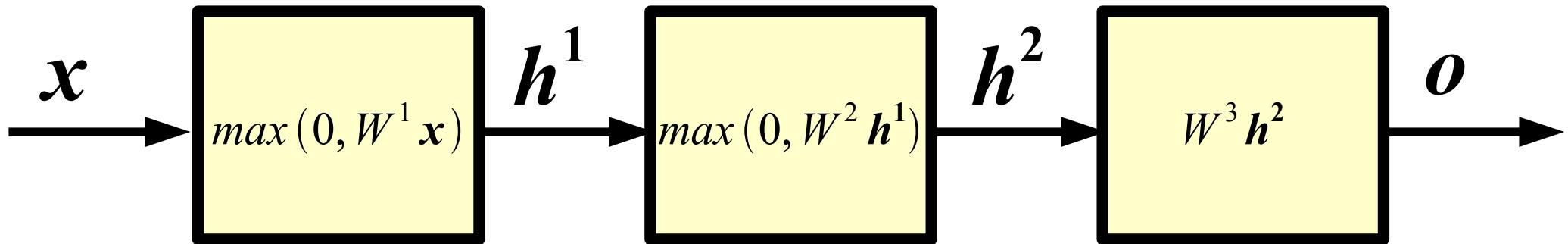
**Question:** what class of functions shall we consider to map the input into the output?

**Answer:** composition of simpler functions.

**Follow-up questions:** Why not a linear combination? What are the “simpler” functions? What is the interpretation?

**Answer:** later...

# Neural Networks: example



$x$  input

$h^1$  1-st layer hidden units

$h^2$  2-nd layer hidden units

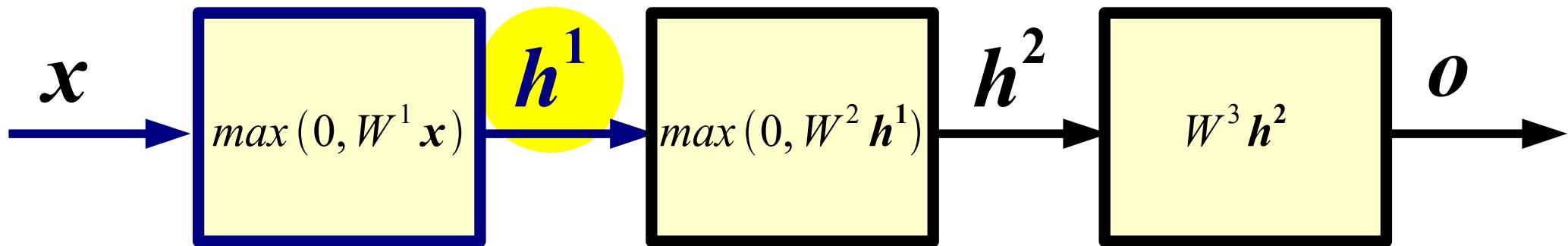
$o$  output

Example of a 2 hidden layer neural network (or 4 layer network,  
counting also input and output).

# Forward Propagation

**Def.:** Forward propagation is the process of computing the output of the network given its input.

# Forward Propagation



$$x \in R^D \quad W^1 \in R^{N_1 \times D} \quad b^1 \in R^{N_1} \quad h^1 \in R^{N_1}$$

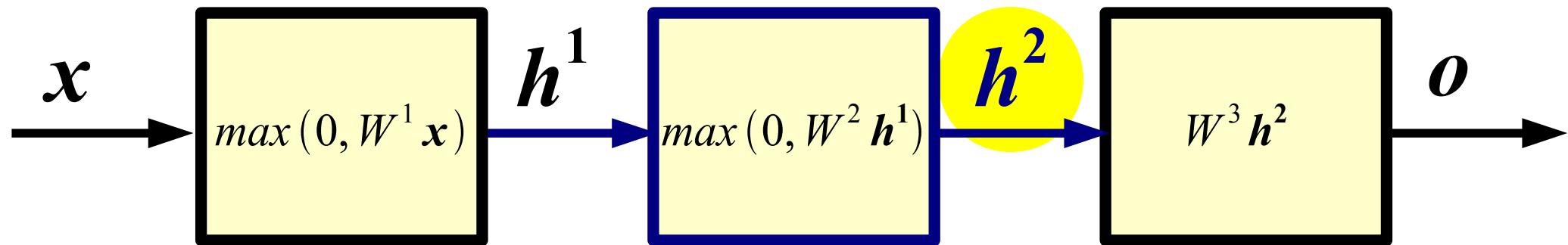
$$h^1 = \max(0, W^1 x + b^1)$$

$W^1$  1-st layer weight matrix or weights

$b^1$  1-st layer biases

The non-linearity  $u = \max(0, v)$  is called **ReLU** in the DL literature. Each output hidden unit takes as input all the units at the previous layer: each such layer is called “**fully connected**”.

# Forward Propagation



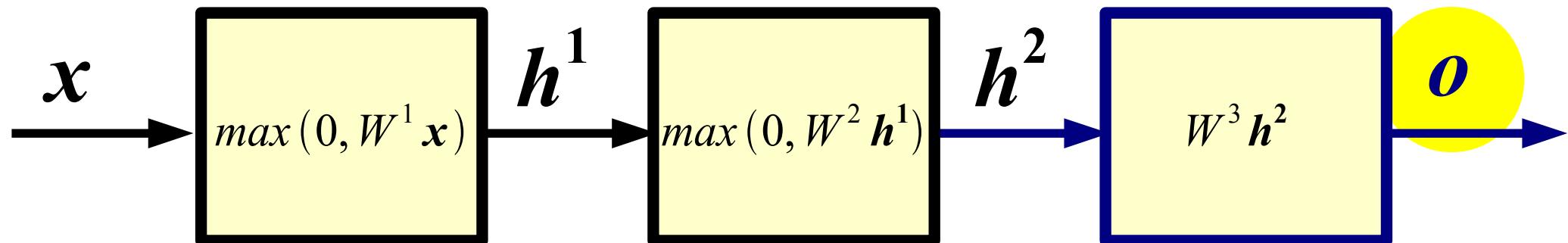
$$h^1 \in R^{N_1} \quad W^2 \in R^{N_2 \times N_1} \quad b^2 \in R^{N_2} \quad h^2 \in R^{N_2}$$

$$h^2 = \max(0, W^2 h^1 + b^2)$$

$W^2$  2-nd layer weight matrix or weights

$b^2$  2-nd layer biases

# Forward Propagation

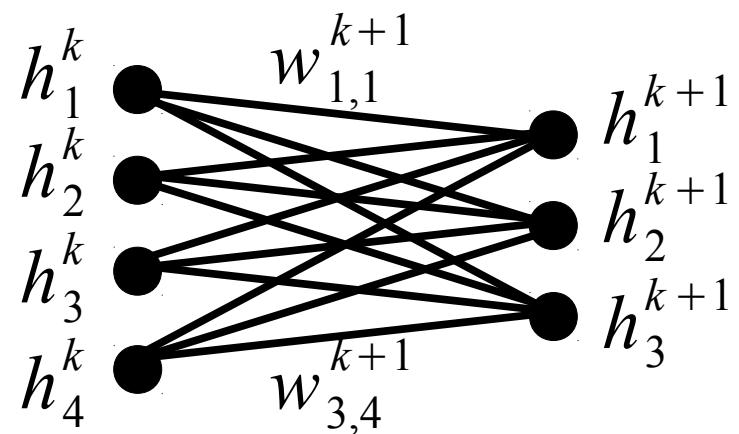
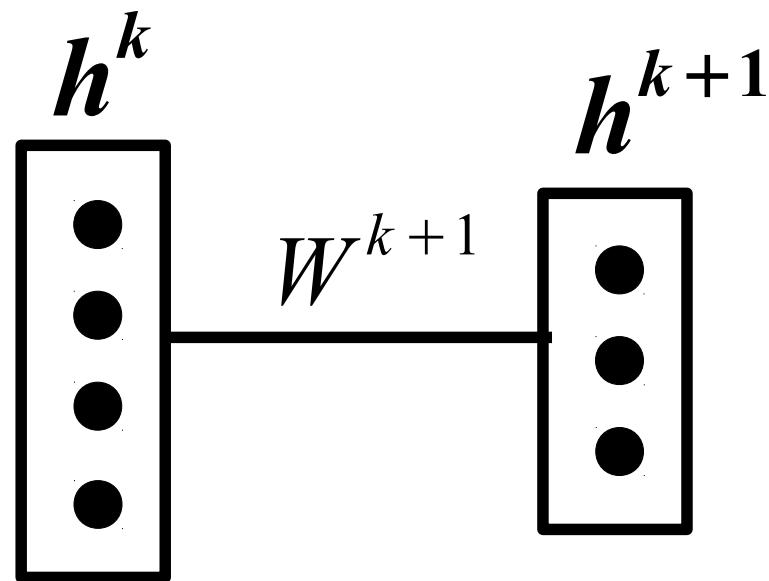
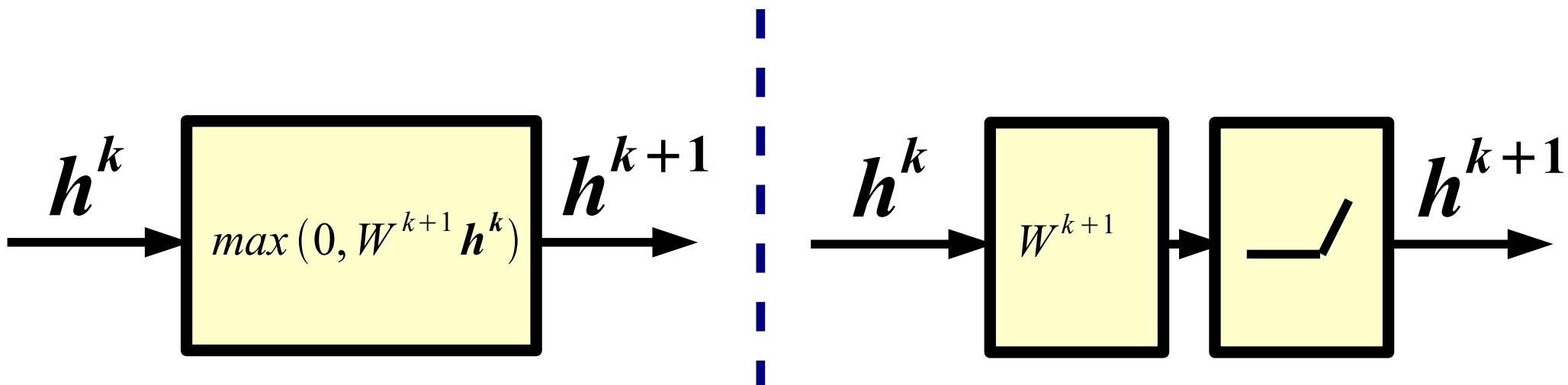


$$h^2 \in R^{N_2} \quad W^3 \in R^{N_3 \times N_2} \quad b^3 \in R^{N_3} \quad o \in R^{N_3}$$

$$o = \max(0, W^3 h^2 + b^3)$$

$W^3$  3-rd layer weight matrix or weights  
 $b^3$  3-rd layer biases

# Alternative Graphical Representation



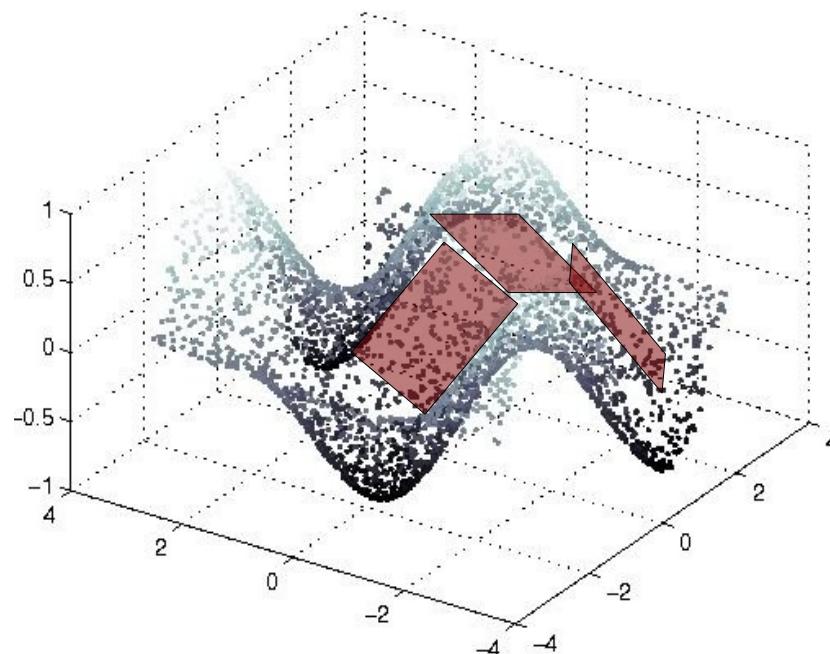
# Interpretation

**Question:** Why can't the mapping between layers be linear?

**Answer:** Because composition of linear functions is a linear function. Neural network would reduce to (1 layer) logistic regression.

**Question:** What do ReLU layers accomplish?

**Answer:** Piece-wise linear tiling: mapping is locally linear.

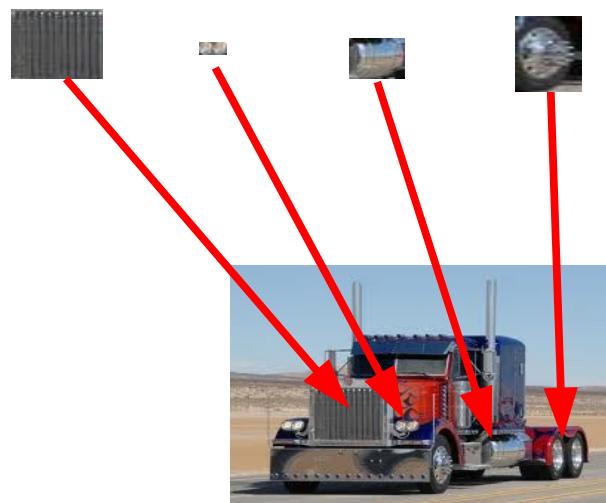


# Interpretation

**Question:** Why do we need many layers?

**Answer:** When input has hierarchical structure, the use of a hierarchical architecture is potentially more efficient because intermediate computations can be re-used. DL architectures are efficient also because they use **distributed representations** which are shared across classes.

[0 0 1 0 0 0 0 1 0 0 1 1 0 0 1 0 ... ] truck feature

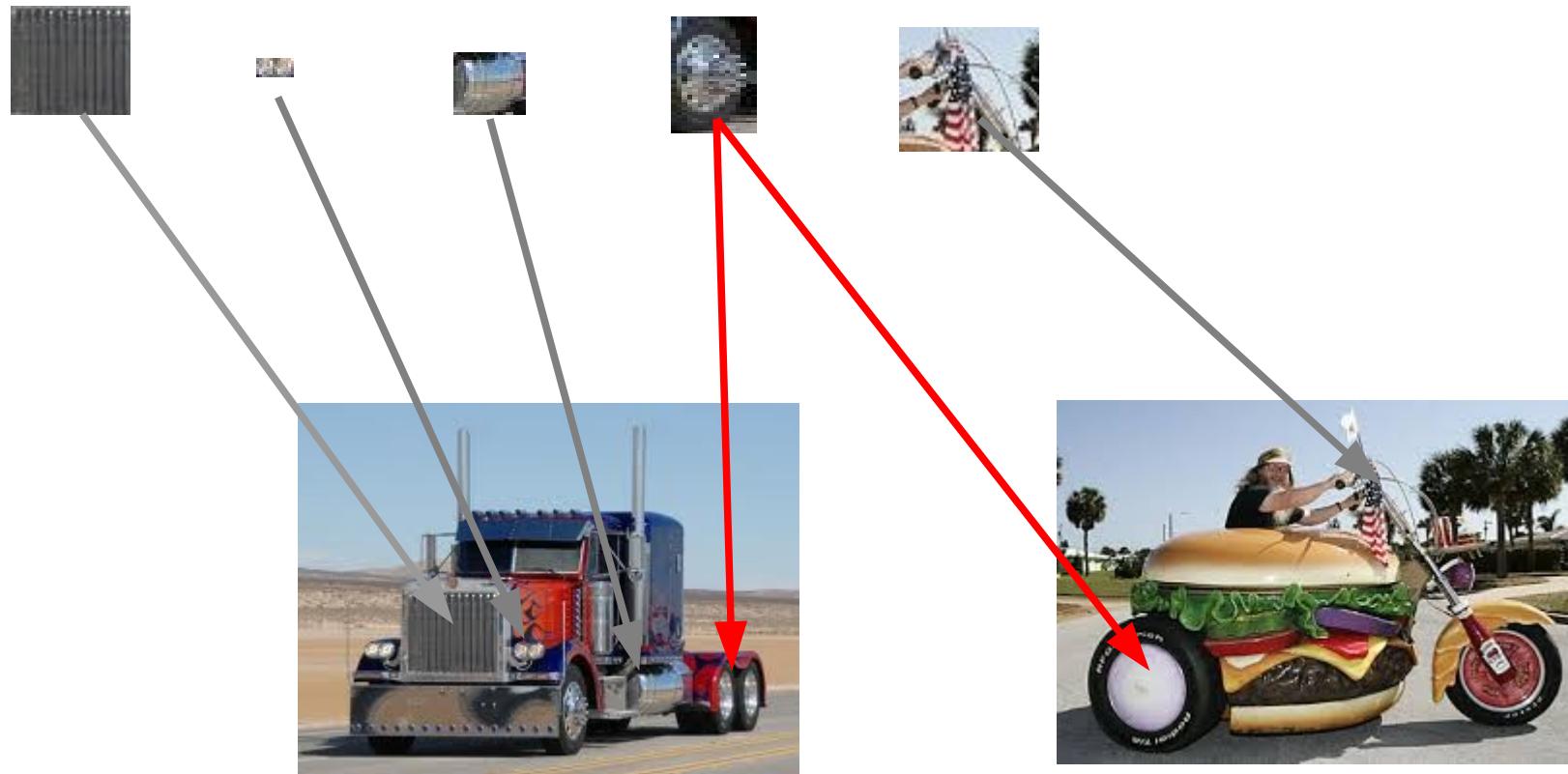


Exponentially more efficient than a 1-of-N representation (a la k-means)

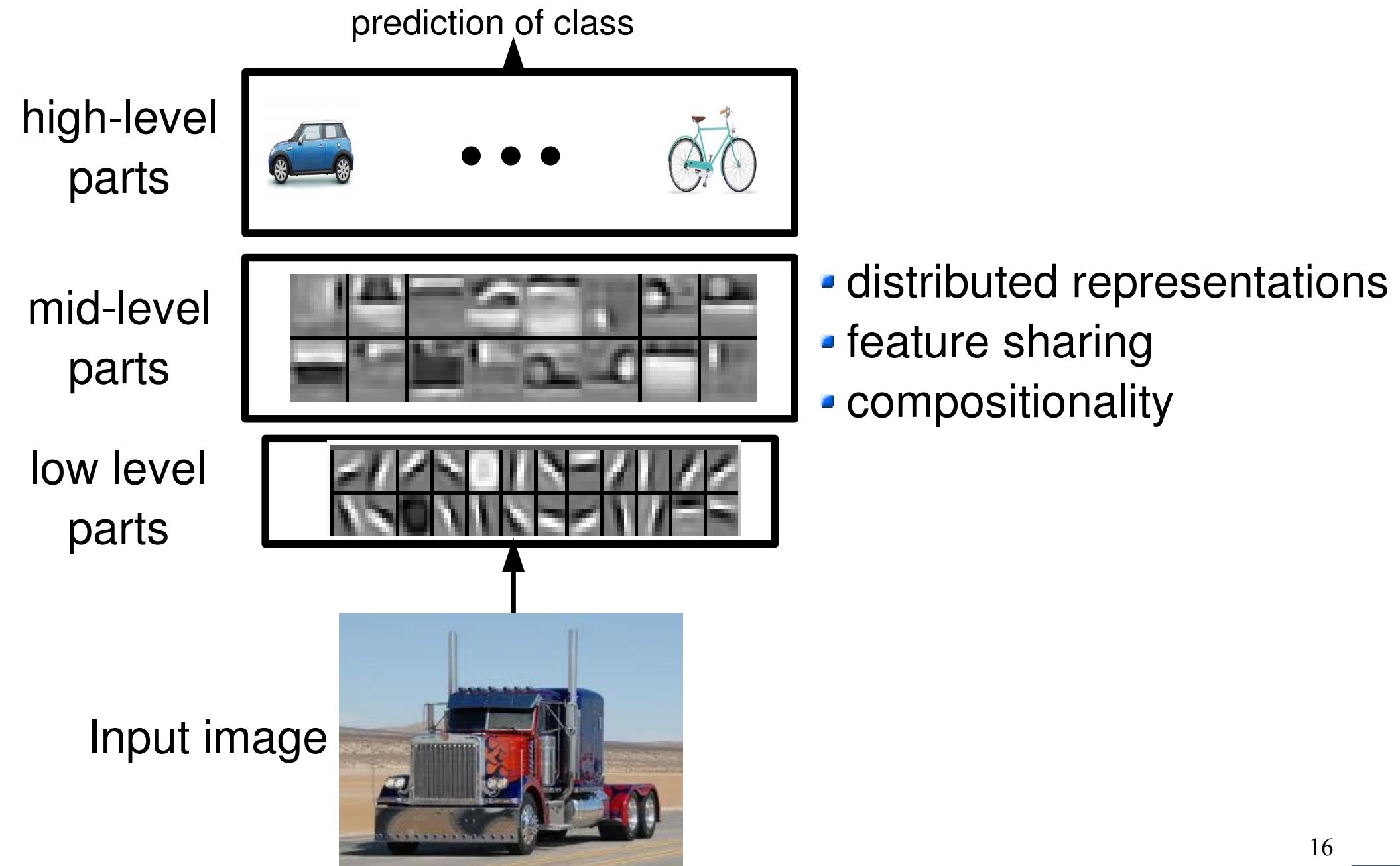
# Interpretation

[1 1 0 0 0 1 0 **1** 0 0 0 0 1 1 0 1 ... ]    motorbike

[0 0 1 0 0 0 0 **1** 0 0 1 1 0 0 1 0 ... ]    truck



# Interpretation



# Interpretation

**Question:** What does a hidden unit do?

**Answer:** It can be thought of as a classifier or feature detector.

**Question:** How many layers? How many hidden units?

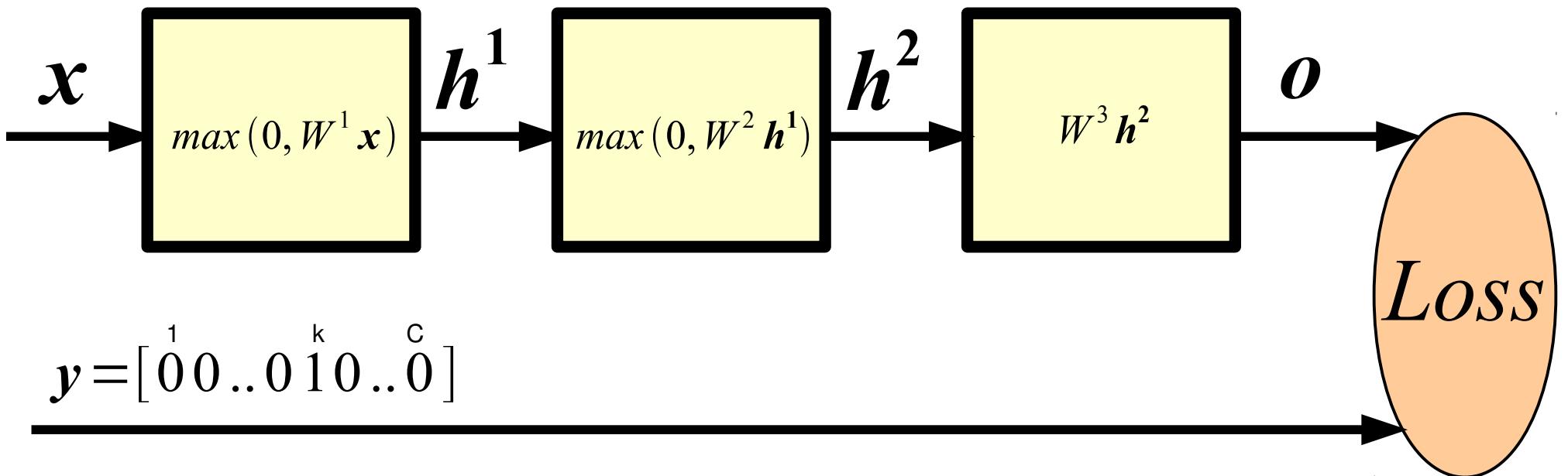
**Answer:** Cross-validation or hyper-parameter search methods are the answer. In general, the wider and the deeper the network the more complicated the mapping.

**Question:** How do I set the weight matrices?

**Answer:** Weight matrices and biases are learned.

First, we need to define a measure of quality of the current mapping. Then, we need to define a procedure to adjust the parameters.

# How Good is a Network?



Probability of class  $k$  given input (softmax):

$$p(c_k=1|x) = \frac{e^{o_k}}{\sum_{j=1}^C e^{o_j}}$$

(Per-sample) **Loss**; e.g., negative log-likelihood (good for classification of small number of classes):

$$L(x, y; \theta) = -\sum_j y_j \log p(c_j|x)$$

# Training

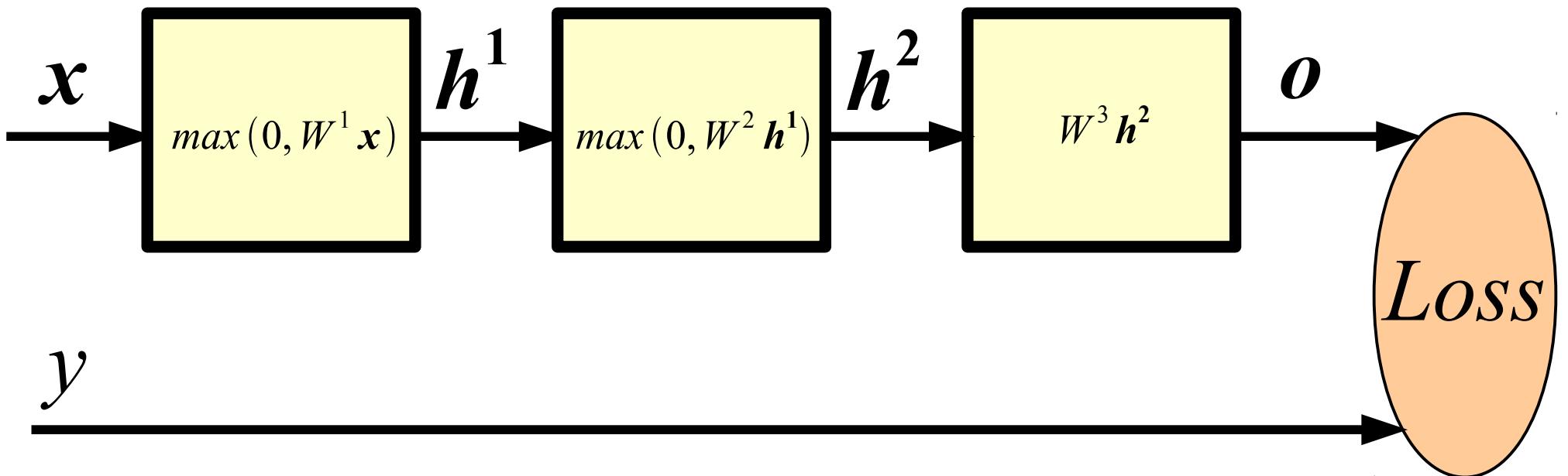
**Learning** consists of minimizing the loss (plus some regularization term) w.r.t. parameters over the whole training set.

$$\theta^* = \arg \min_{\theta} \sum_{n=1}^P L(\mathbf{x}^n, y^n; \theta)$$

**Question:** How to minimize a complicated function of the parameters?

**Answer:** Chain rule, a.k.a. **Backpropagation**! That is the procedure to compute gradients of the loss w.r.t. parameters in a multi-layer neural network.

# Key Idea: Wiggle To Decrease Loss



Let's say we want to decrease the loss by adjusting  $W_{i,j}^1$ .  
We could consider a very small  $\epsilon = 1e-6$  and compute:

$$L(\mathbf{x}, y; \boldsymbol{\theta})$$

$$L(\mathbf{x}, y; \boldsymbol{\theta} \setminus W_{i,j}^1, W_{i,j}^1 + \epsilon)$$

Then, update:

$$W_{i,j}^1 \leftarrow W_{i,j}^1 + \epsilon \operatorname{sgn}(L(\mathbf{x}, y; \boldsymbol{\theta}) - L(\mathbf{x}, y; \boldsymbol{\theta} \setminus W_{i,j}^1, W_{i,j}^1 + \epsilon))$$

# Derivative w.r.t. Input of Softmax

$$p(c_k=1|\mathbf{x}) = \frac{e^{o_k}}{\sum_j e^{o_j}}$$

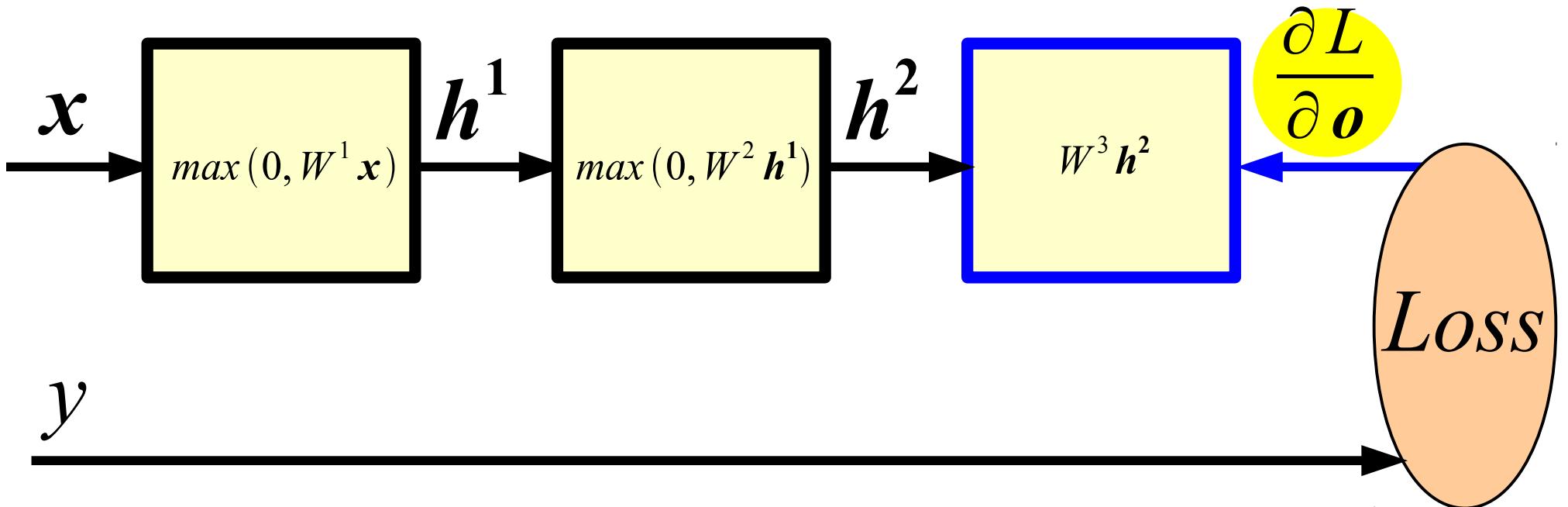
$$L(\mathbf{x}, y; \boldsymbol{\theta}) = -\sum_j y_j \log p(c_j|\mathbf{x}) \quad y = [0^1 0..0^k 1^0 .. 0^c]$$

By substituting the first formula in the second, and taking the derivative w.r.t.  $o$  we get:

$$\frac{\partial L}{\partial o} = p(c|\mathbf{x}) - y$$

**HOMEWORK: prove it!**

# Backward Propagation

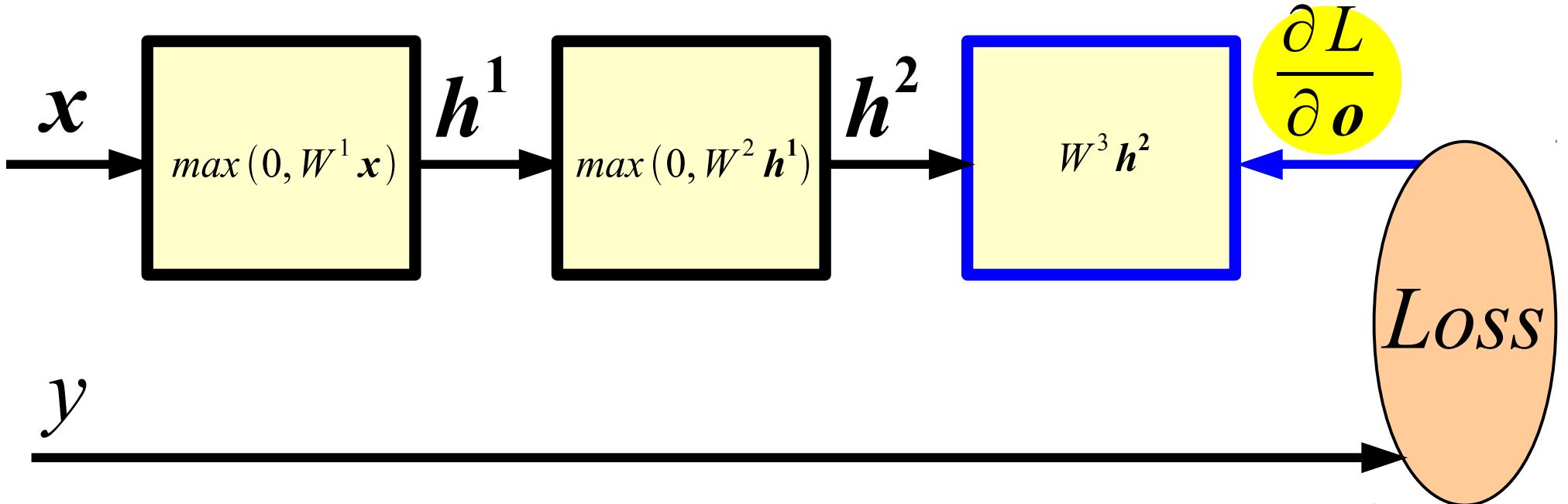


Given  $\frac{\partial L}{\partial \mathbf{o}}$  and assuming we can easily compute the Jacobian of each module, we have:

$$\frac{\partial L}{\partial W^3} = \frac{\partial L}{\partial \mathbf{o}} \frac{\partial \mathbf{o}}{\partial W^3}$$

$$\frac{\partial L}{\partial h^2} = \frac{\partial L}{\partial \mathbf{o}} \frac{\partial \mathbf{o}}{\partial h^2}$$

# Backward Propagation



Given  $\frac{\partial L}{\partial \mathbf{o}}$  and assuming we can easily compute the Jacobian of each module, we have:

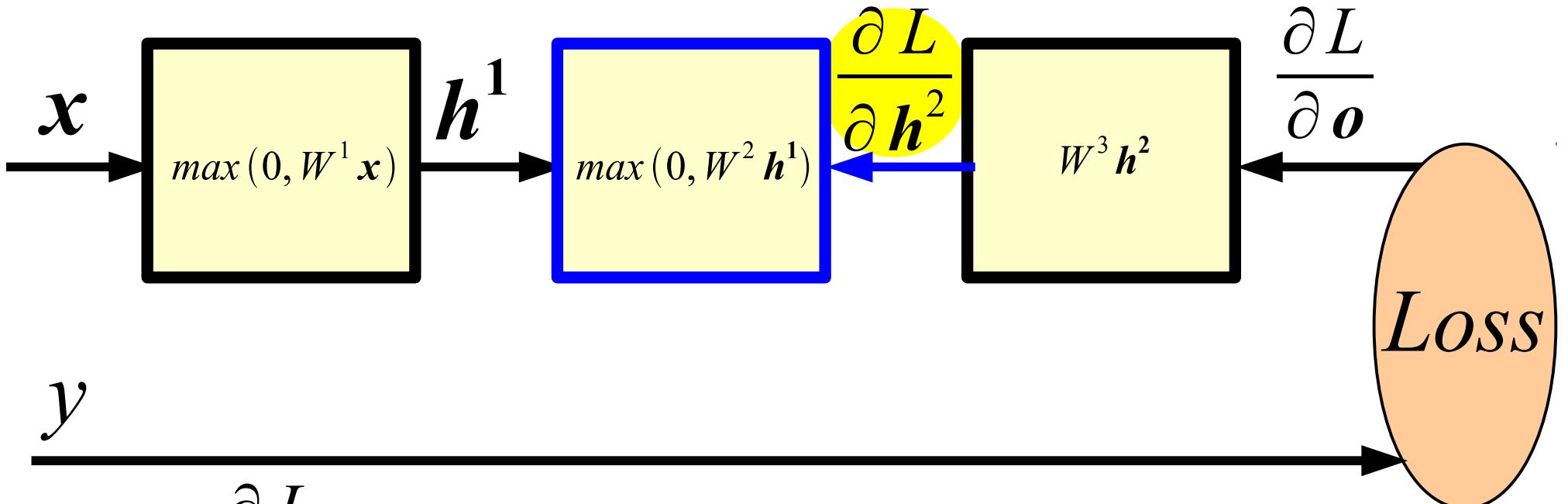
$$\frac{\partial L}{\partial W^3} = \frac{\partial L}{\partial \mathbf{o}} \frac{\partial \mathbf{o}}{\partial W^3}$$

$$\frac{\partial L}{\partial h^2} = \frac{\partial L}{\partial \mathbf{o}} \frac{\partial \mathbf{o}}{\partial h^2}$$

$$\frac{\partial L}{\partial W^3} = (p(c|x) - y) h^{2T}$$

$$\frac{\partial L}{\partial h^2} = W^{3T} (p(c|x) - y)$$

# Backward Propagation

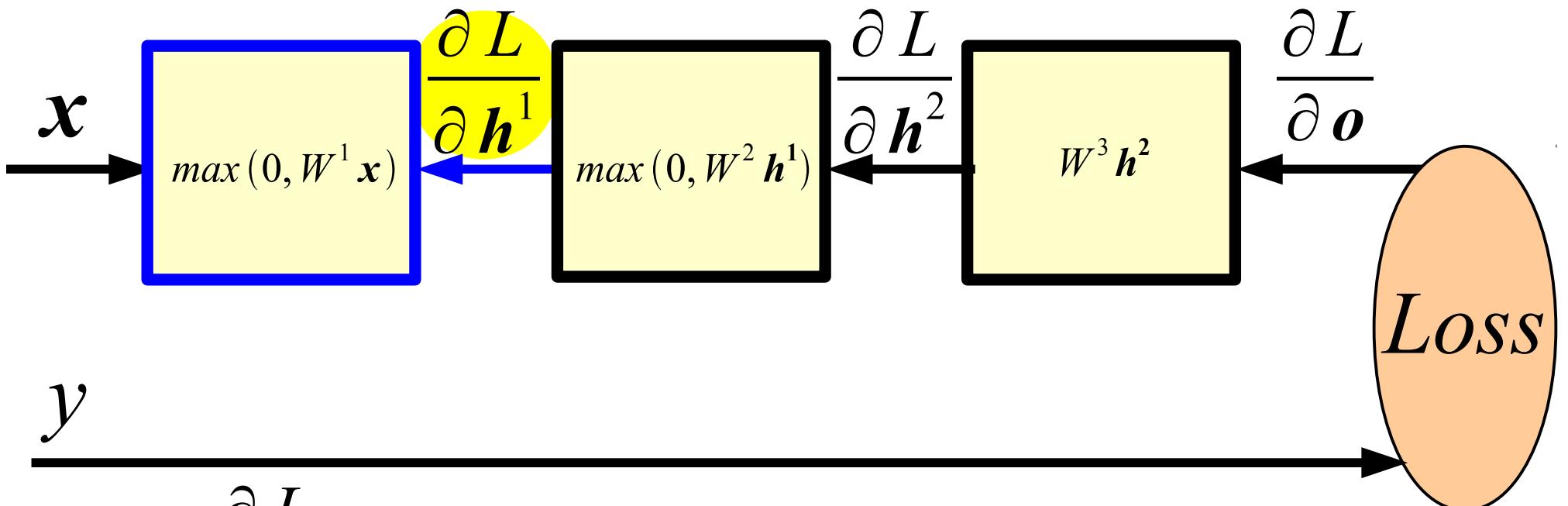


Given  $\frac{\partial L}{\partial h^2}$  we can compute now:

$$\frac{\partial L}{\partial W^2} = \frac{\partial L}{\partial h^2} \frac{\partial h^2}{\partial W^2}$$

$$\frac{\partial L}{\partial h^1} = \frac{\partial L}{\partial h^2} \frac{\partial h^2}{\partial h^1}$$

# Backward Propagation



Given  $\frac{\partial L}{\partial h^1}$  we can compute now:

$$\frac{\partial L}{\partial W^1} = \frac{\partial L}{\partial h^1} \frac{\partial h^1}{\partial W^1}$$

# Backward Propagation

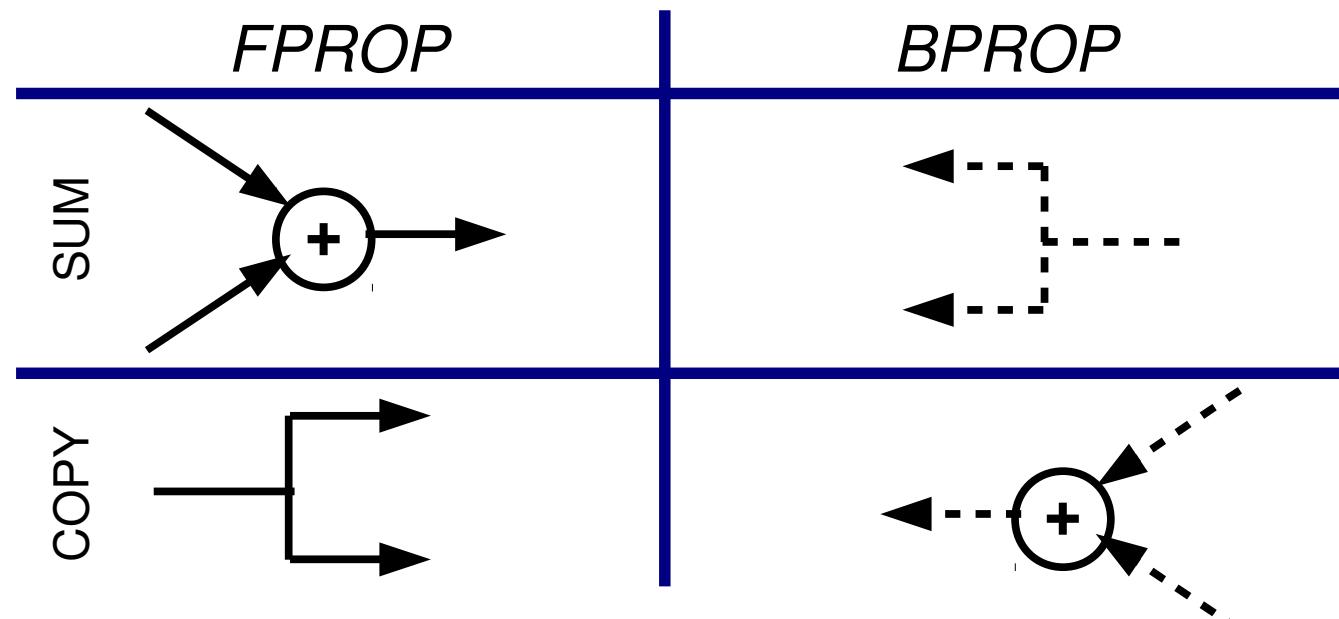
**Question:** Does BPROP work with ReLU layers only?

**Answer:** Nope, any a.e. differentiable transformation works.

**Question:** What's the computational cost of BPROP?

**Answer:** About twice FPROP (need to compute gradients w.r.t. input and parameters at every layer).

**Note:** FPROP and BPROP are dual of each other. E.g.,:



# Optimization

**Stochastic Gradient Descent (on mini-batches):**

$$\theta \leftarrow \theta - \eta \frac{\partial L}{\partial \theta}, \eta \in (0, 1)$$

**Stochastic Gradient Descent with Momentum:**

$$\theta \leftarrow \theta - \eta \Delta$$

$$\Delta \leftarrow 0.9 \Delta + \frac{\partial L}{\partial \theta}$$

**Note: there are many other variants...**

# Toy Code (Matlab): Neural Net Trainer

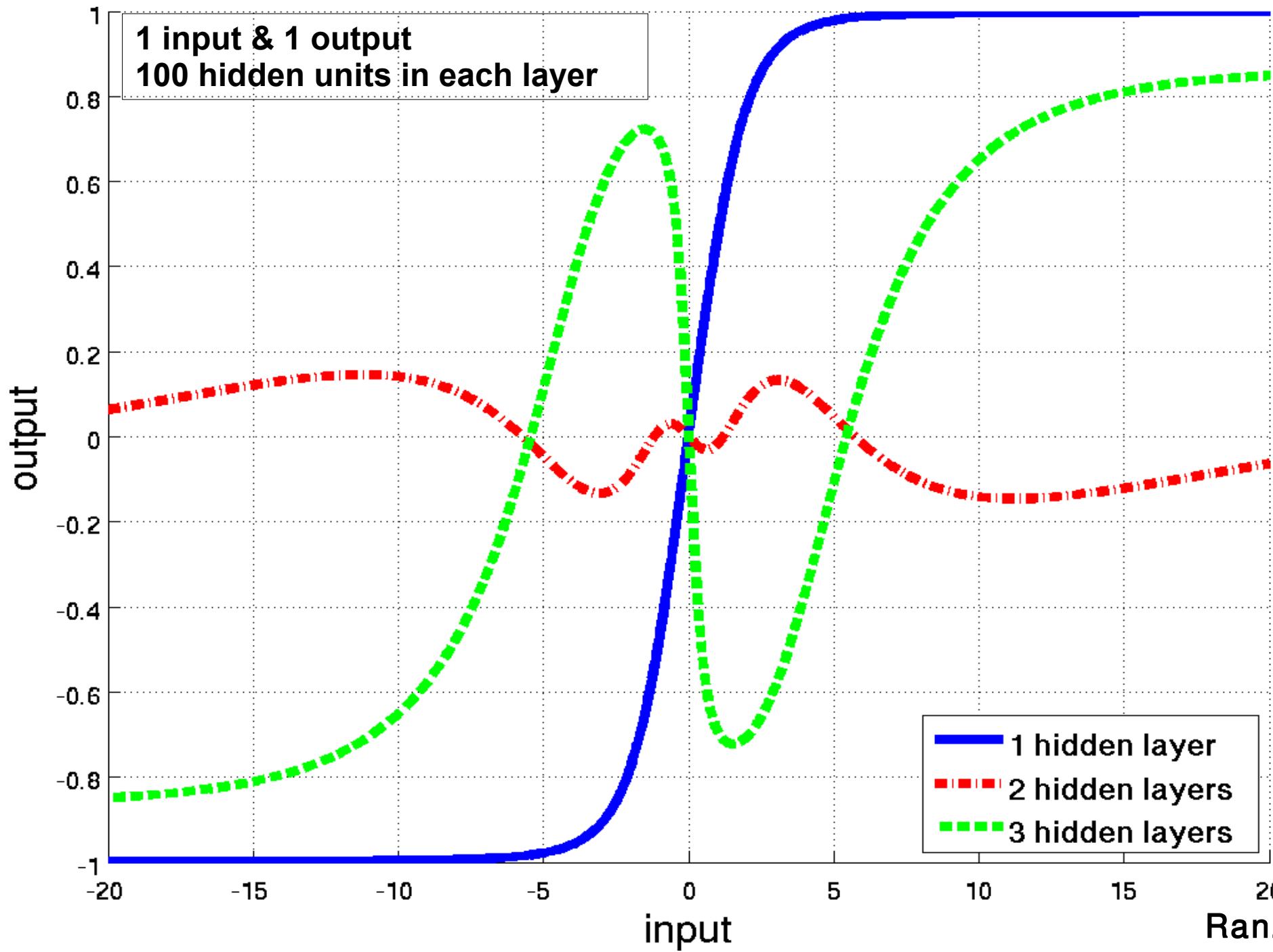
```
% F-PROP
for i = 1 : nr_layers - 1
    [h{i} jac{i}] = nonlinearity(W{i} * h{i-1} + b{i});
end
h{nr_layers-1} = W{nr_layers-1} * h{nr_layers-2} + b{nr_layers-1};
prediction = softmax(h{l-1});

% CROSS ENTROPY LOSS
loss = - sum(sum(log(prediction) .* target)) / batch_size;

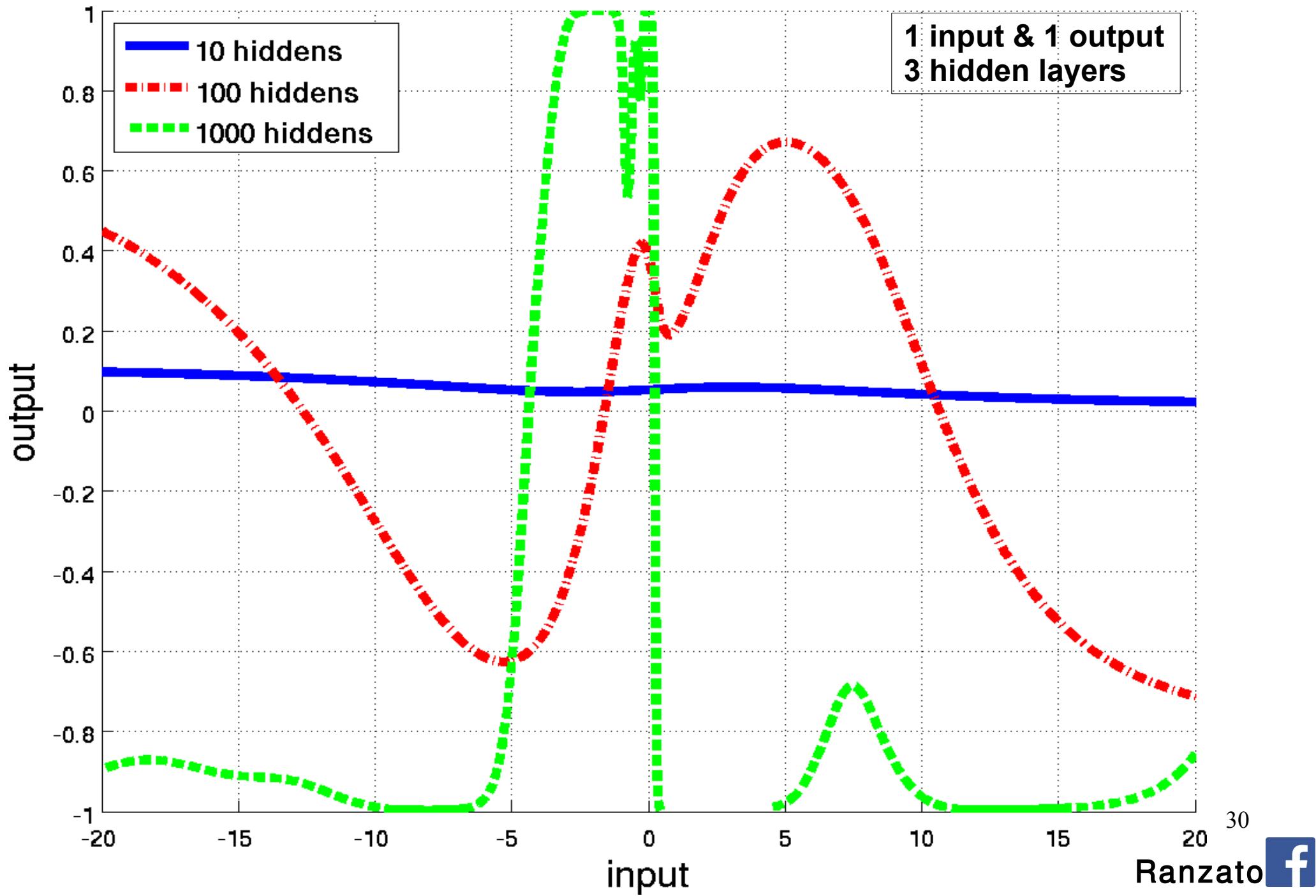
% B-PROP
dh{l-1} = prediction - target;
for i = nr_layers - 1 : -1 : 1
    Wgrad{i} = dh{i} * h{i-1}';
    bgrad{i} = sum(dh{i}, 2);
    dh{i-1} = (W{i}' * dh{i}) .* jac{i-1};
end

% UPDATE
for i = 1 : nr_layers - 1
    W{i} = W{i} - (lr / batch_size) * Wgrad{i};
    b{i} = b{i} - (lr / batch_size) * bgrad{i};
end
```

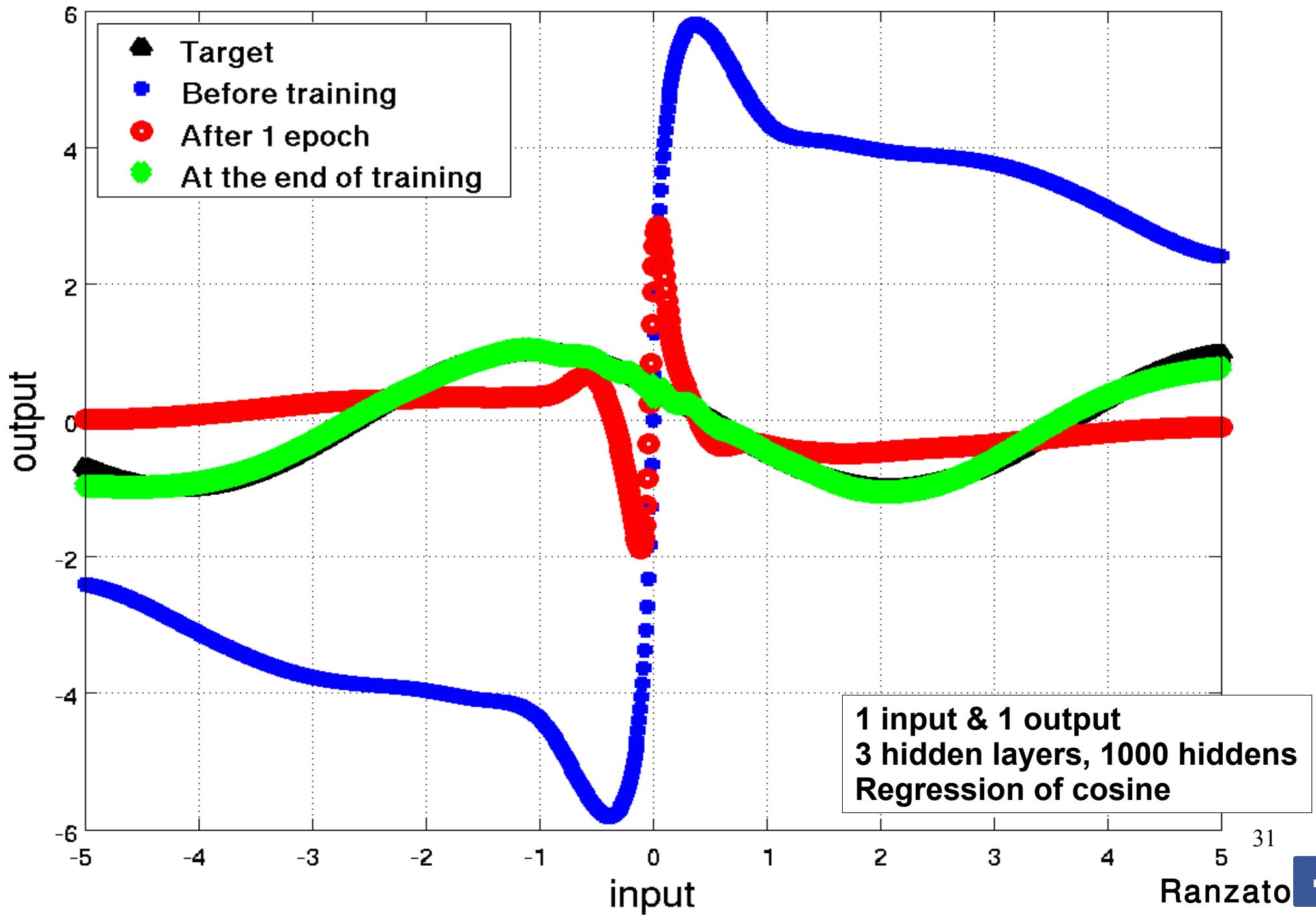
# Toy Example: Synthetic Data



# Toy Example: Synthetic Data



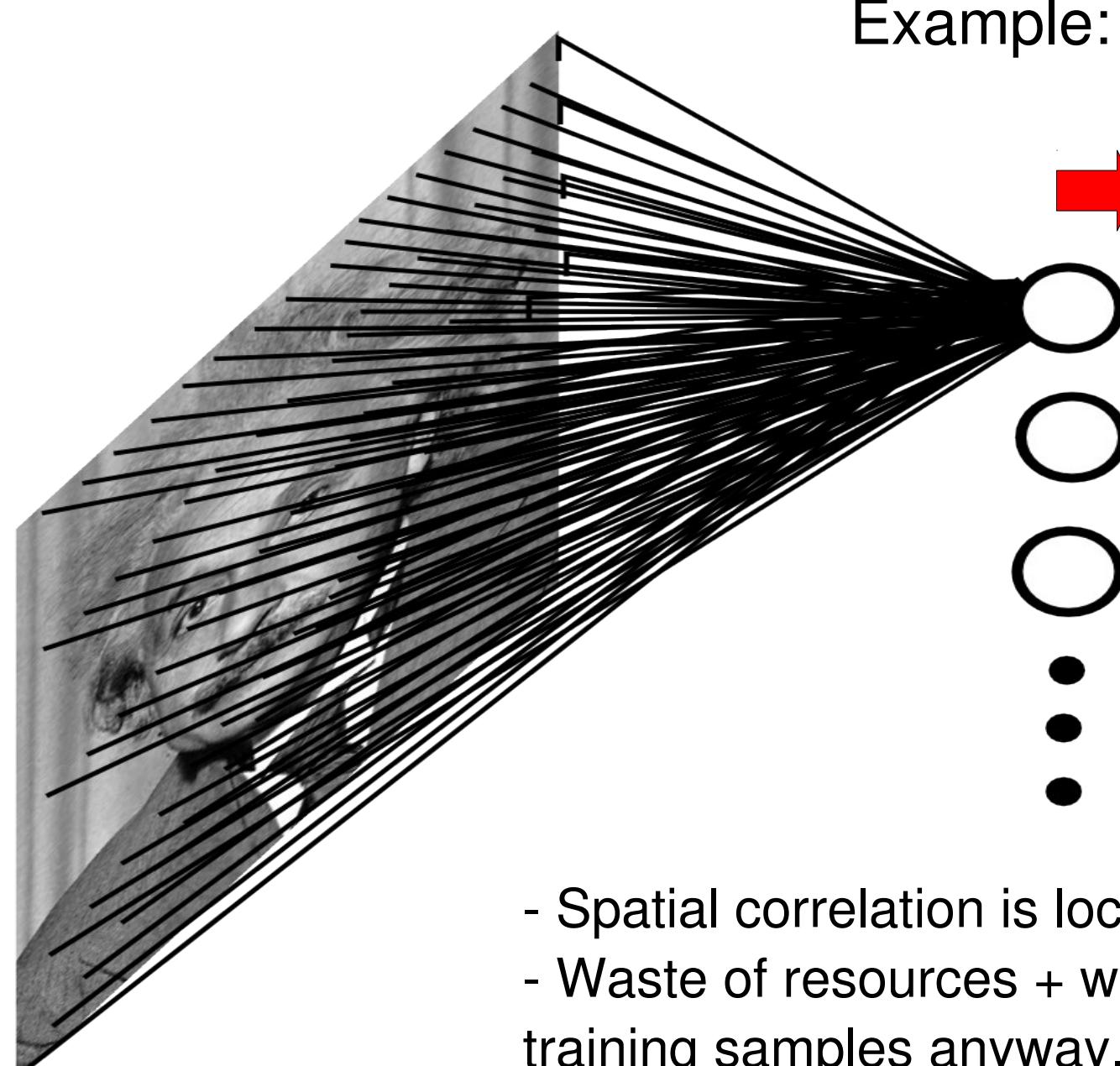
# Toy Example: Synthetic Data



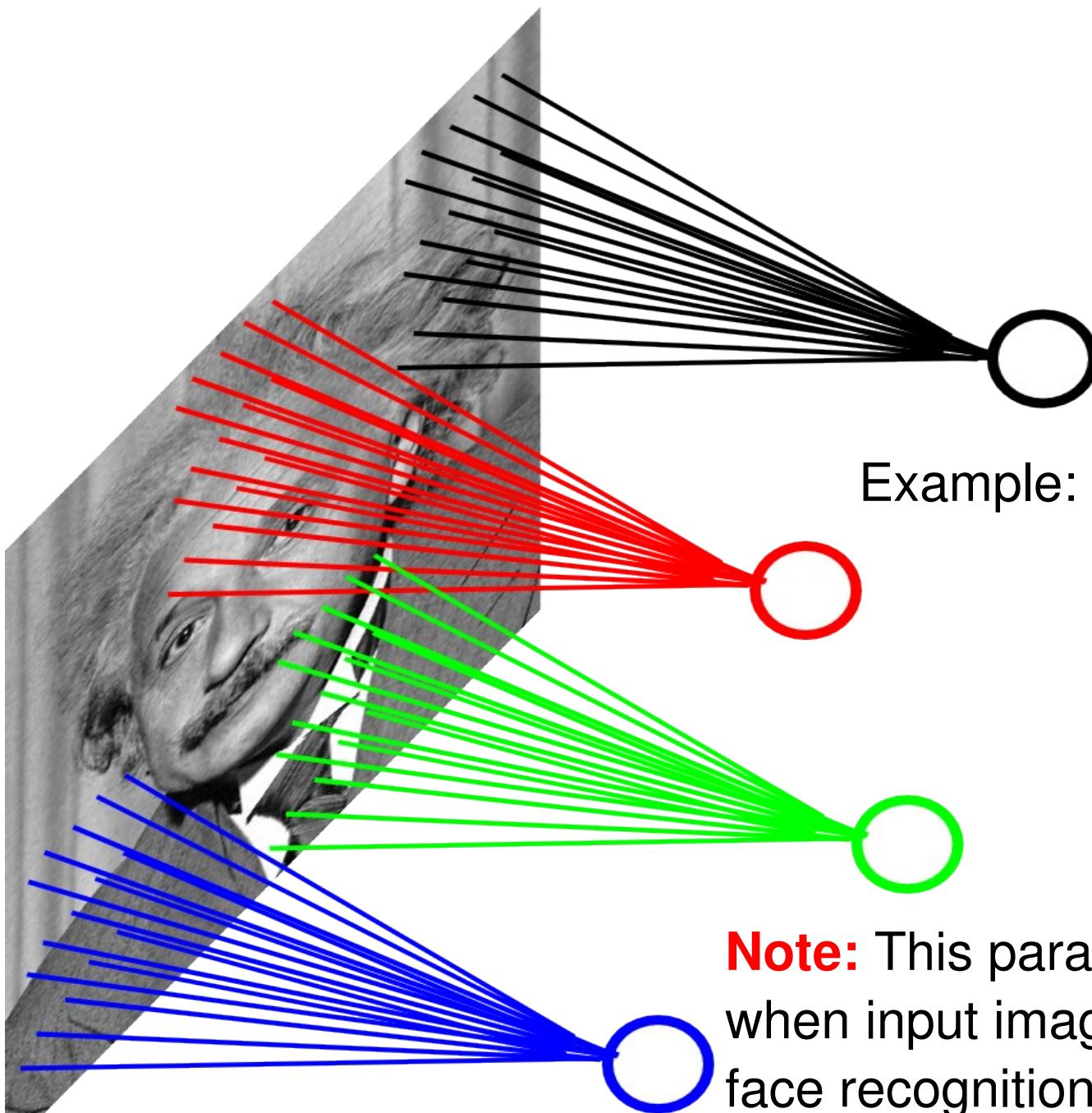
# Outline

- Supervised Neural Networks
- Convolutional Neural Networks
- Examples
- Tips

# Fully Connected Layer



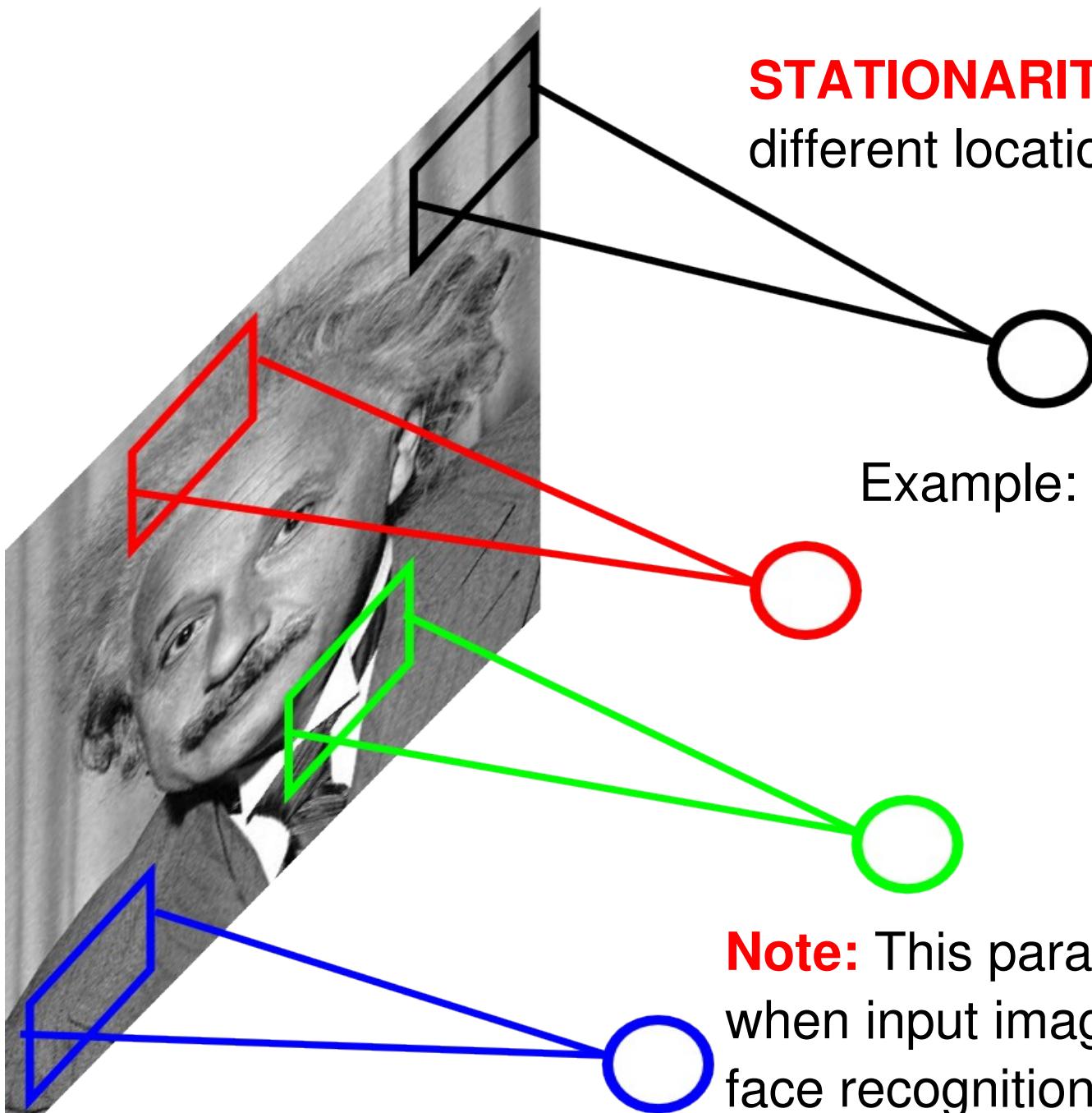
# Locally Connected Layer



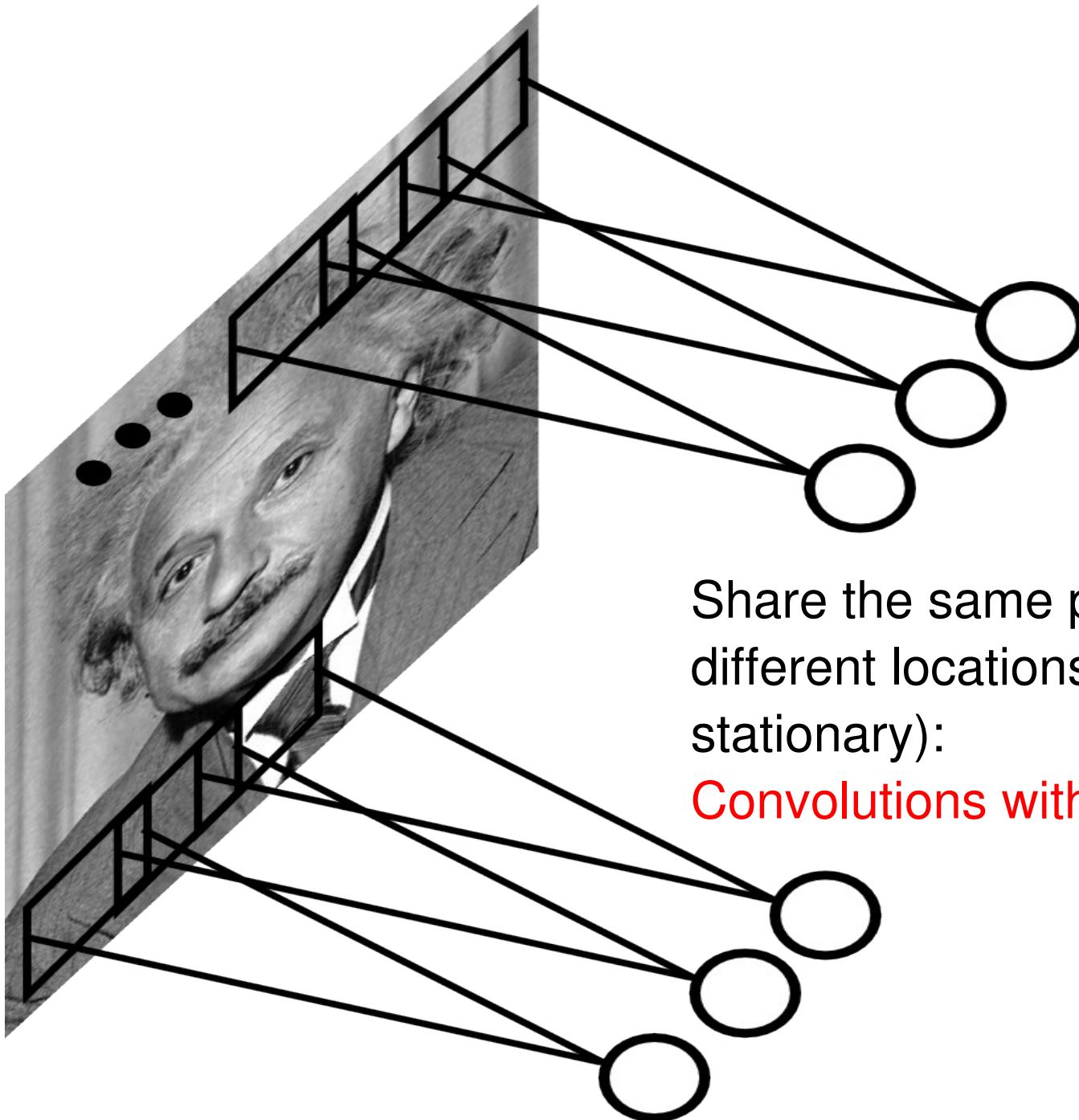
Example: 200x200 image  
40K hidden units  
Filter size: 10x10  
4M parameters

**Note:** This parameterization is good  
when input image is registered (e.g.,  
face recognition).

# Locally Connected Layer



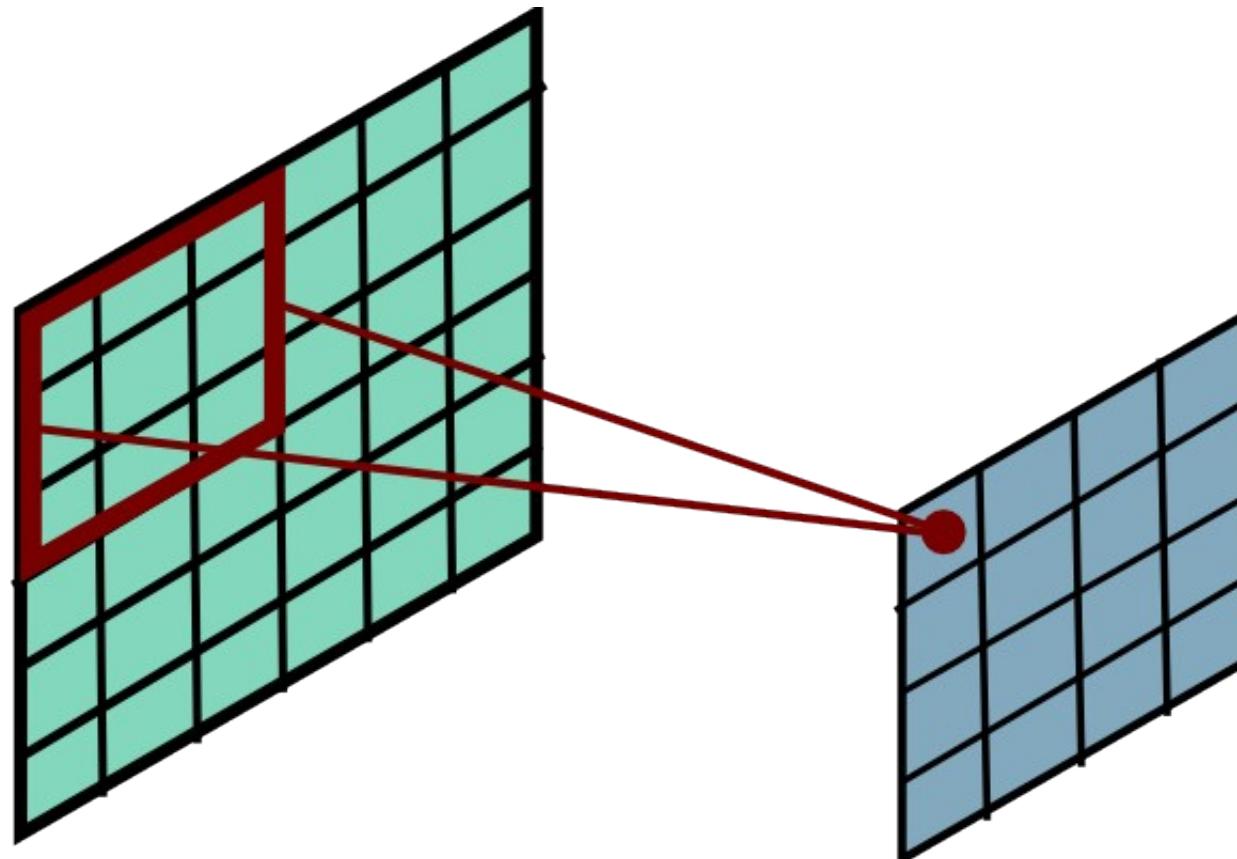
# Convolutional Layer



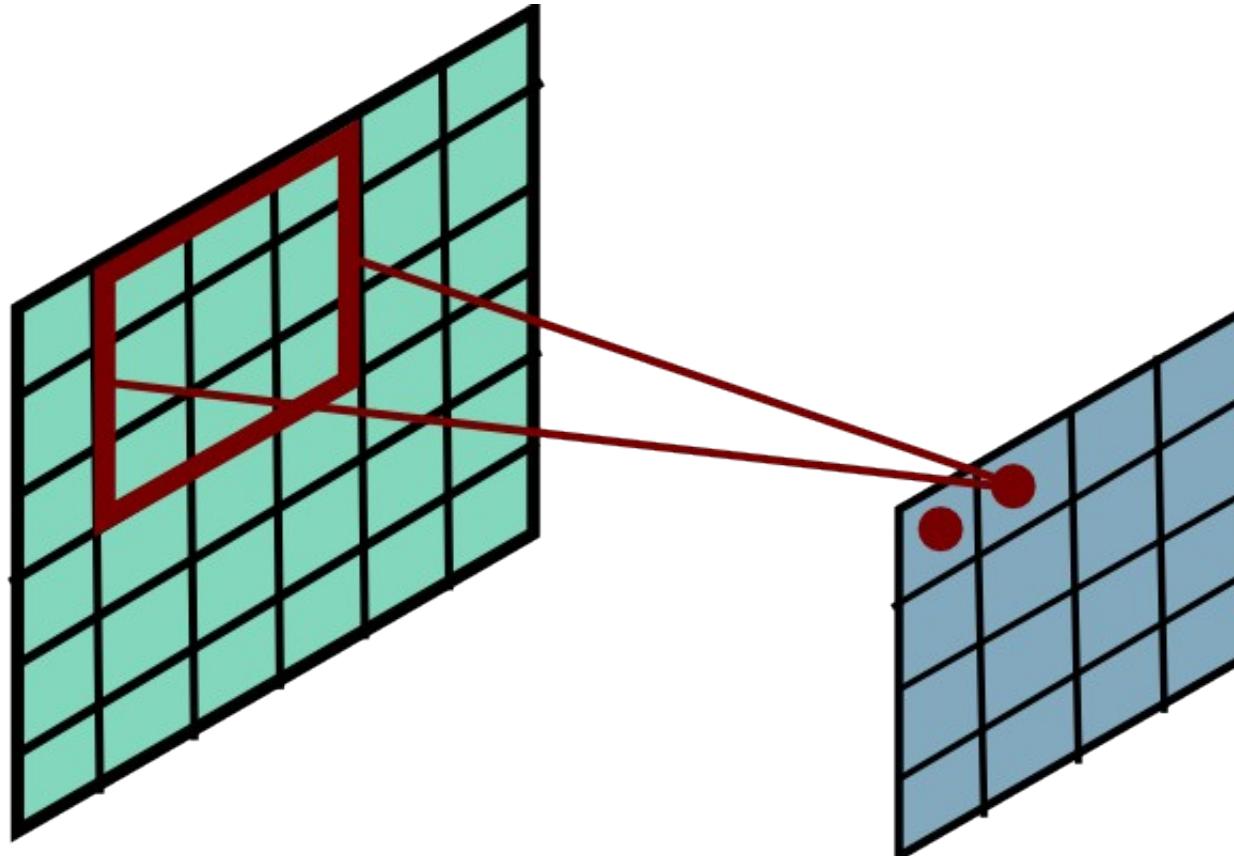
Share the same parameters across  
different locations (assuming input is  
stationary):

**Convolutions with learned kernels**

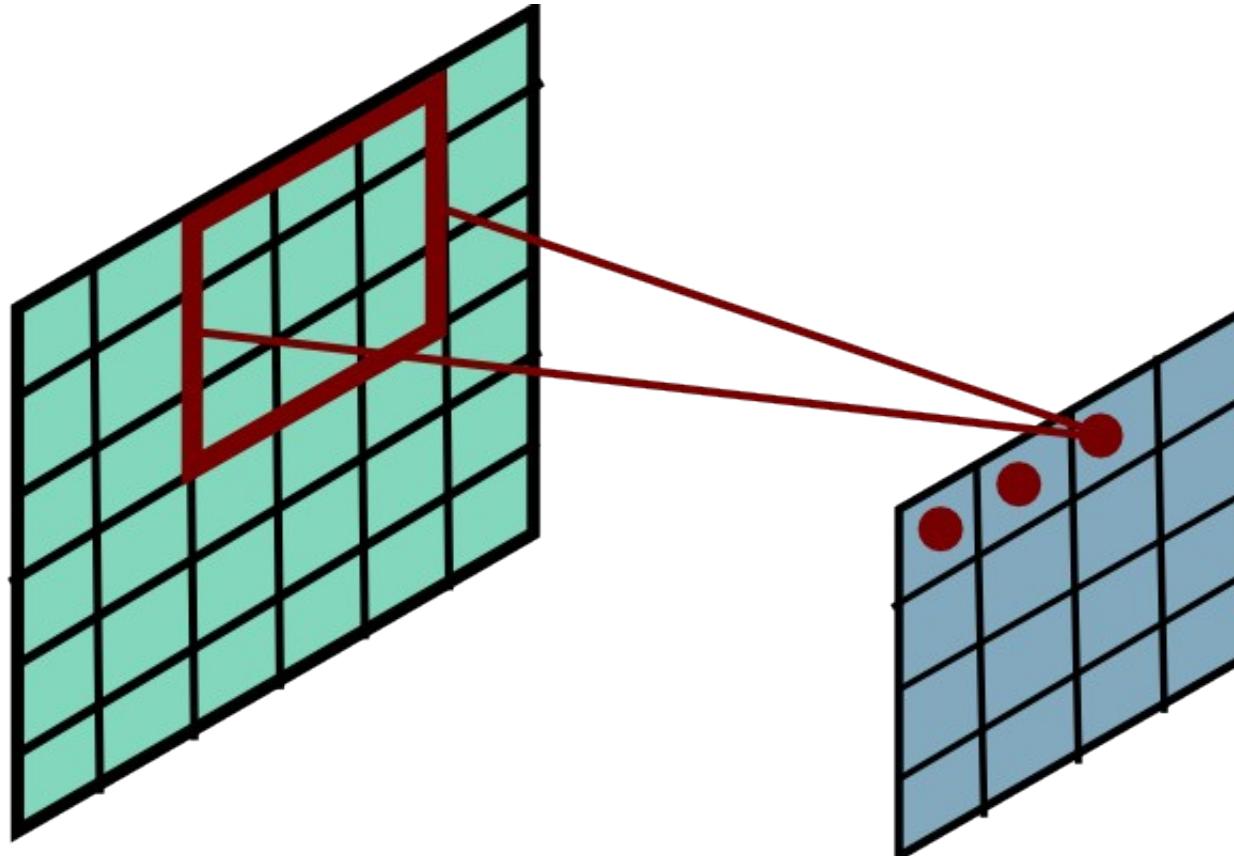
# Convolutional Layer



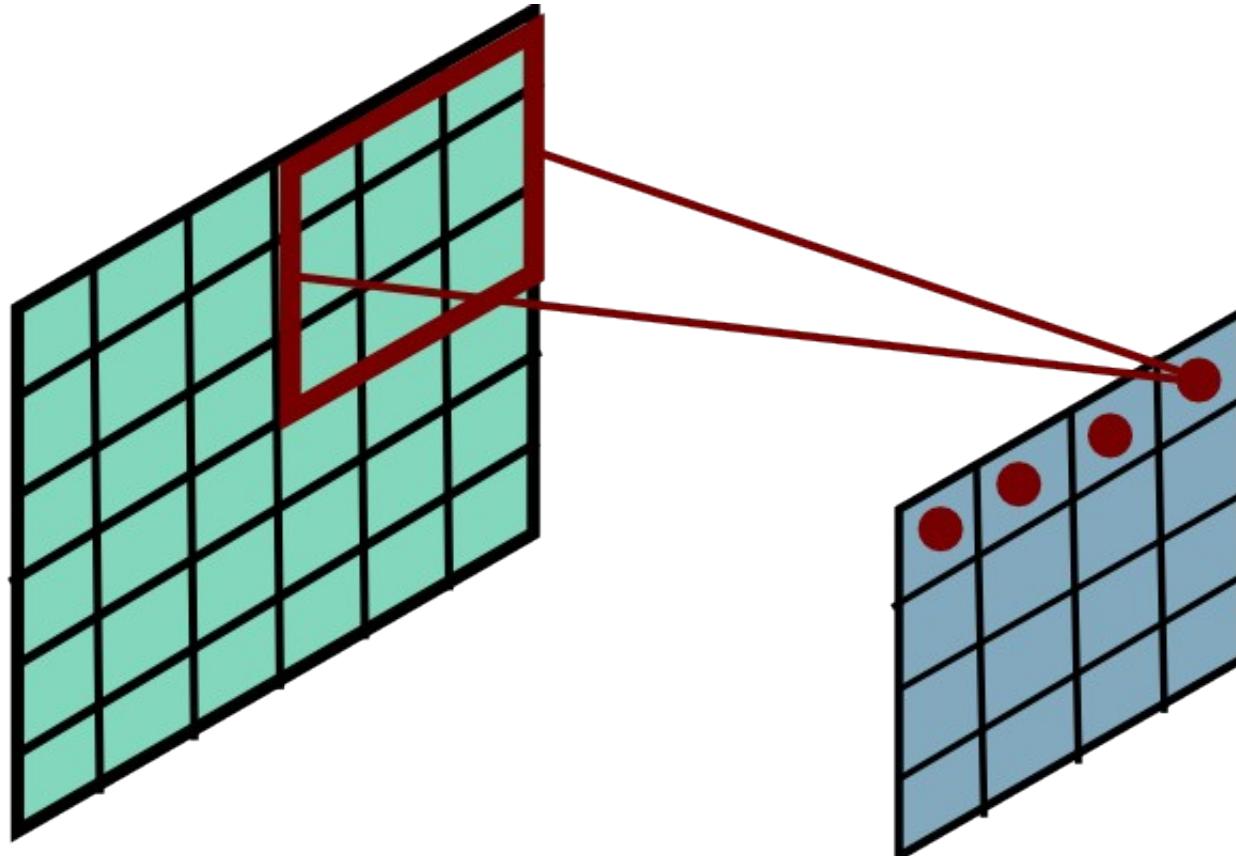
# Convolutional Layer



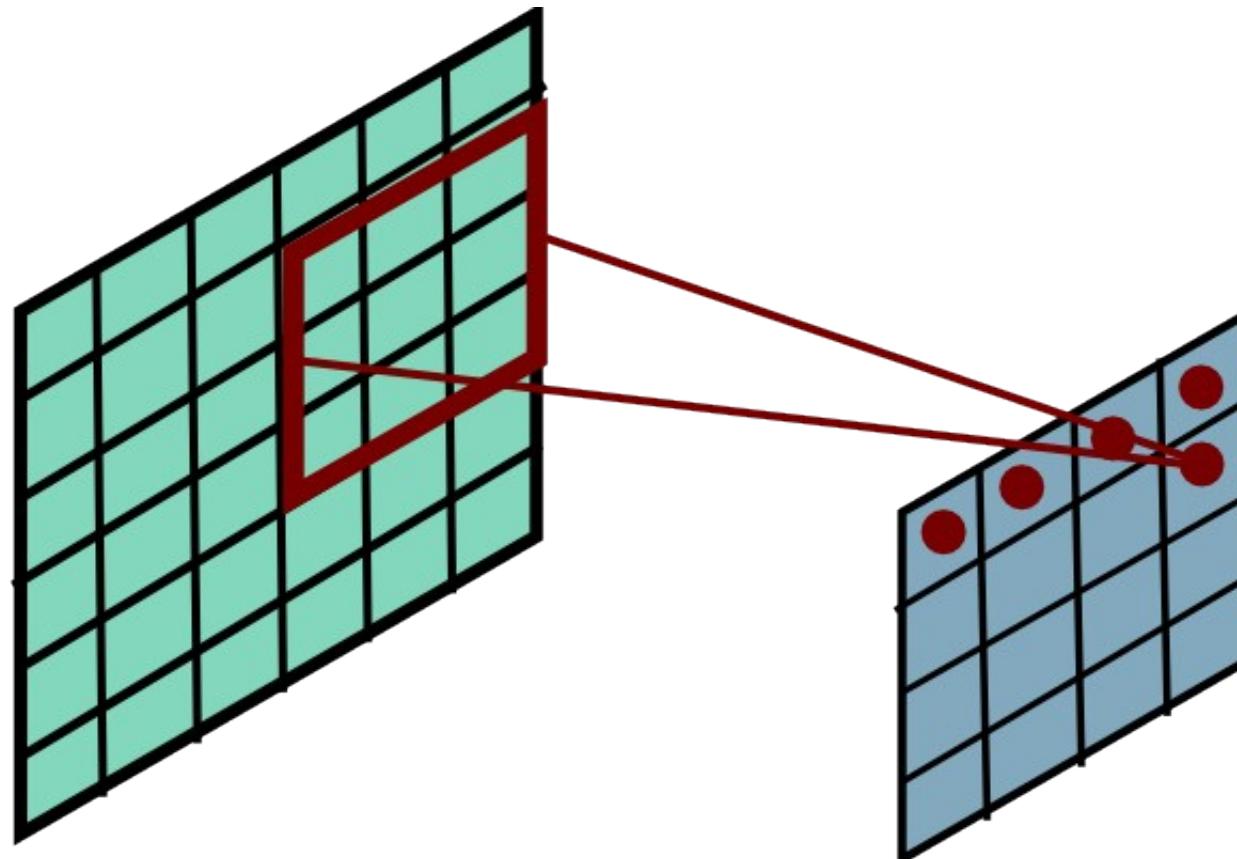
# Convolutional Layer



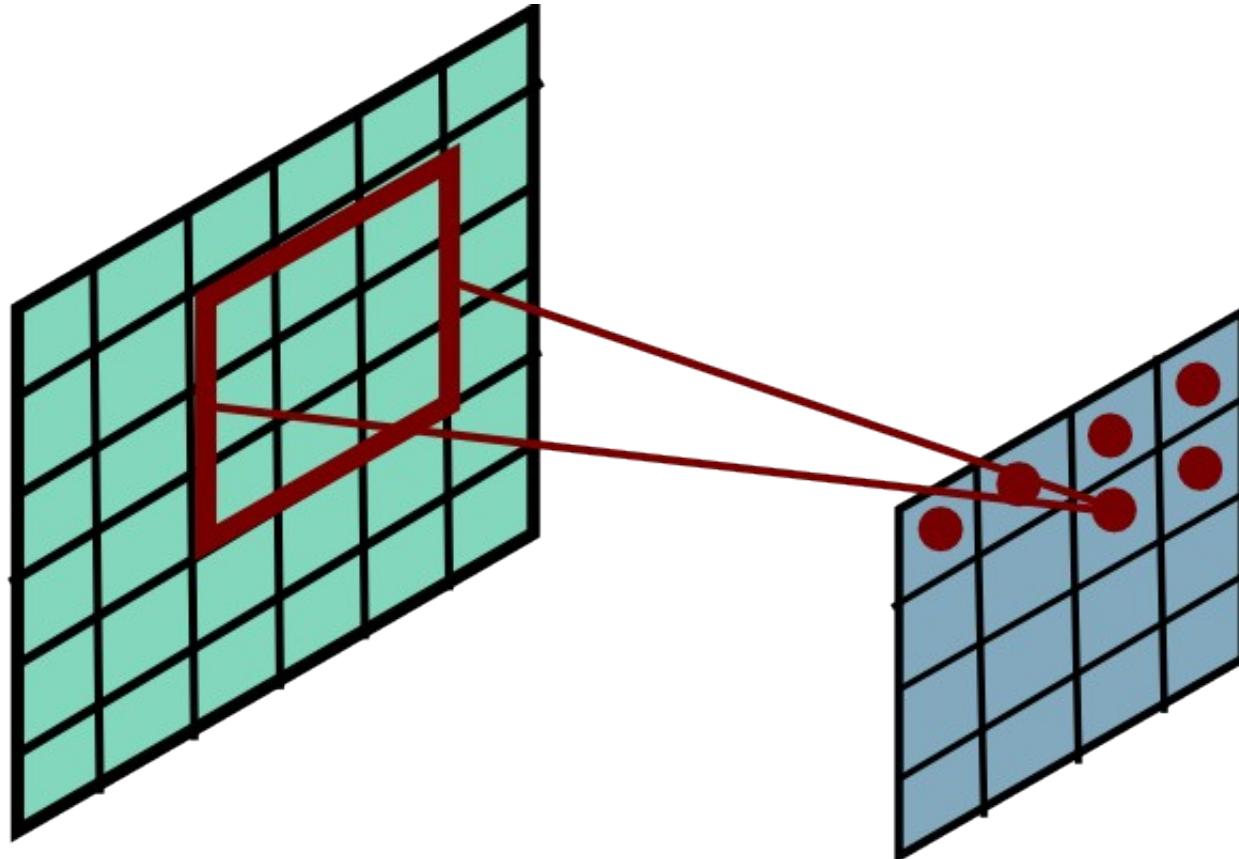
# Convolutional Layer



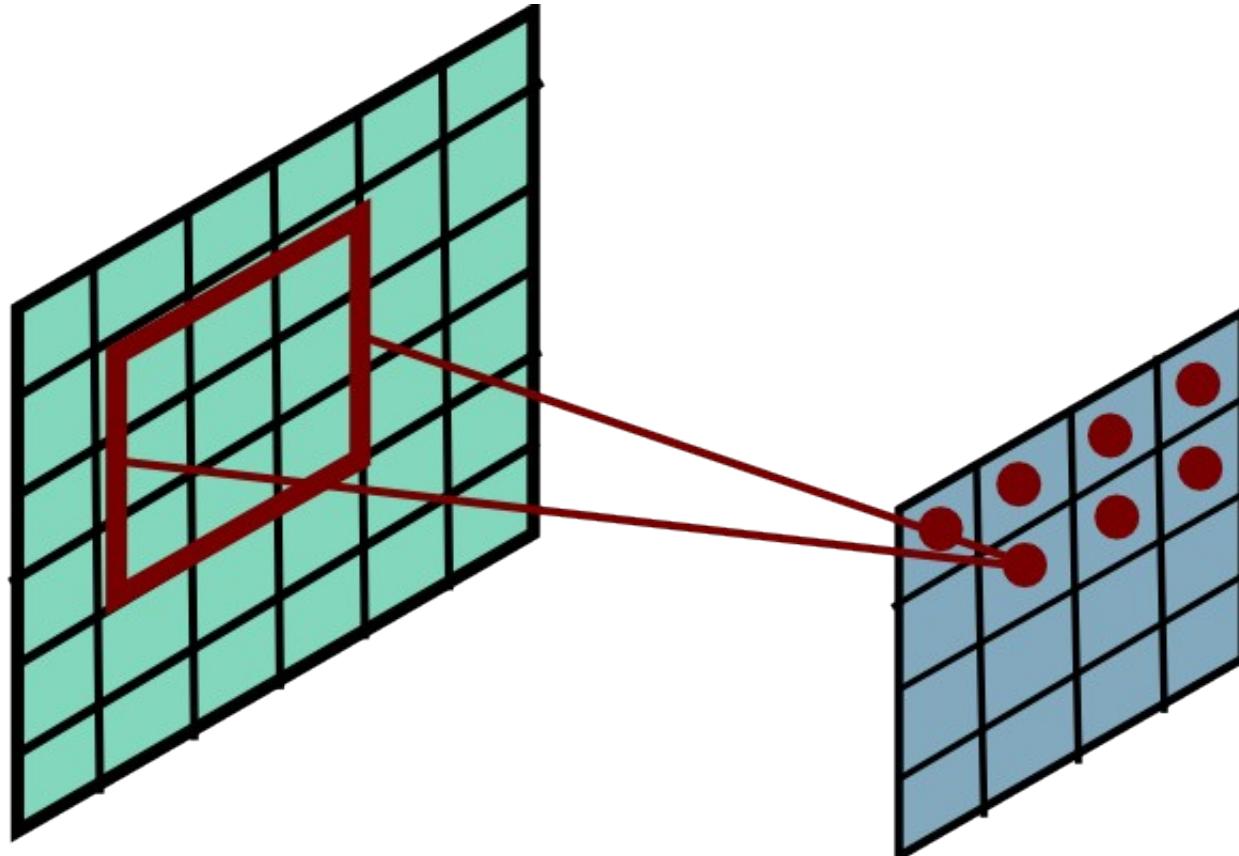
# Convolutional Layer



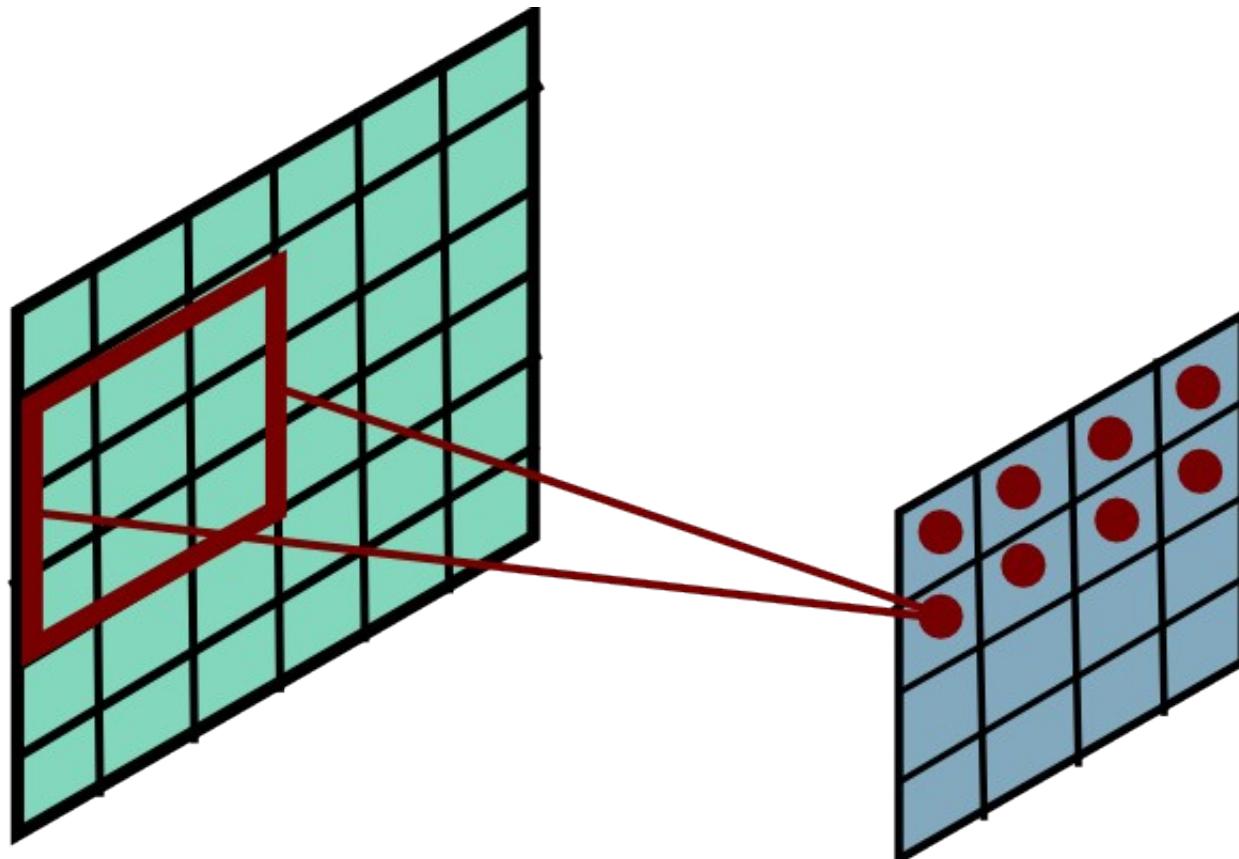
# Convolutional Layer



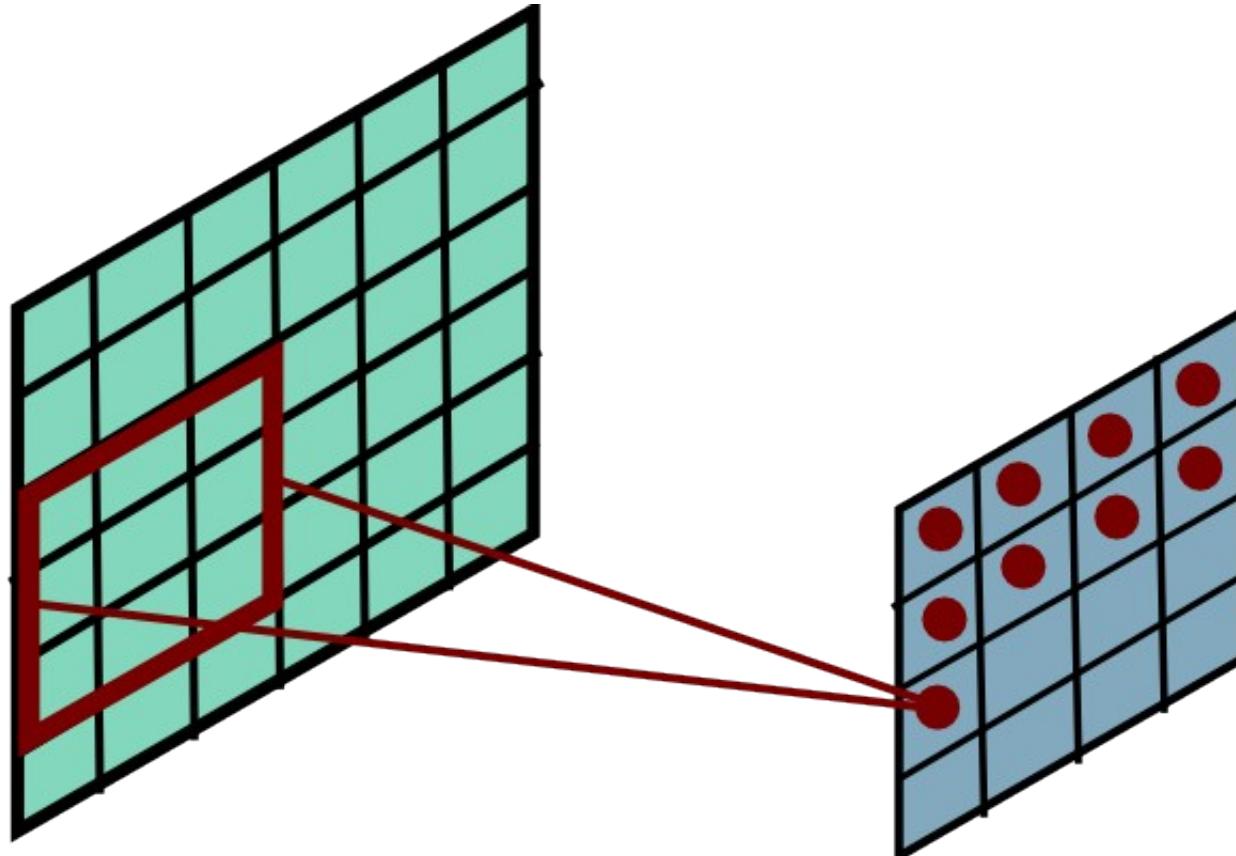
# Convolutional Layer



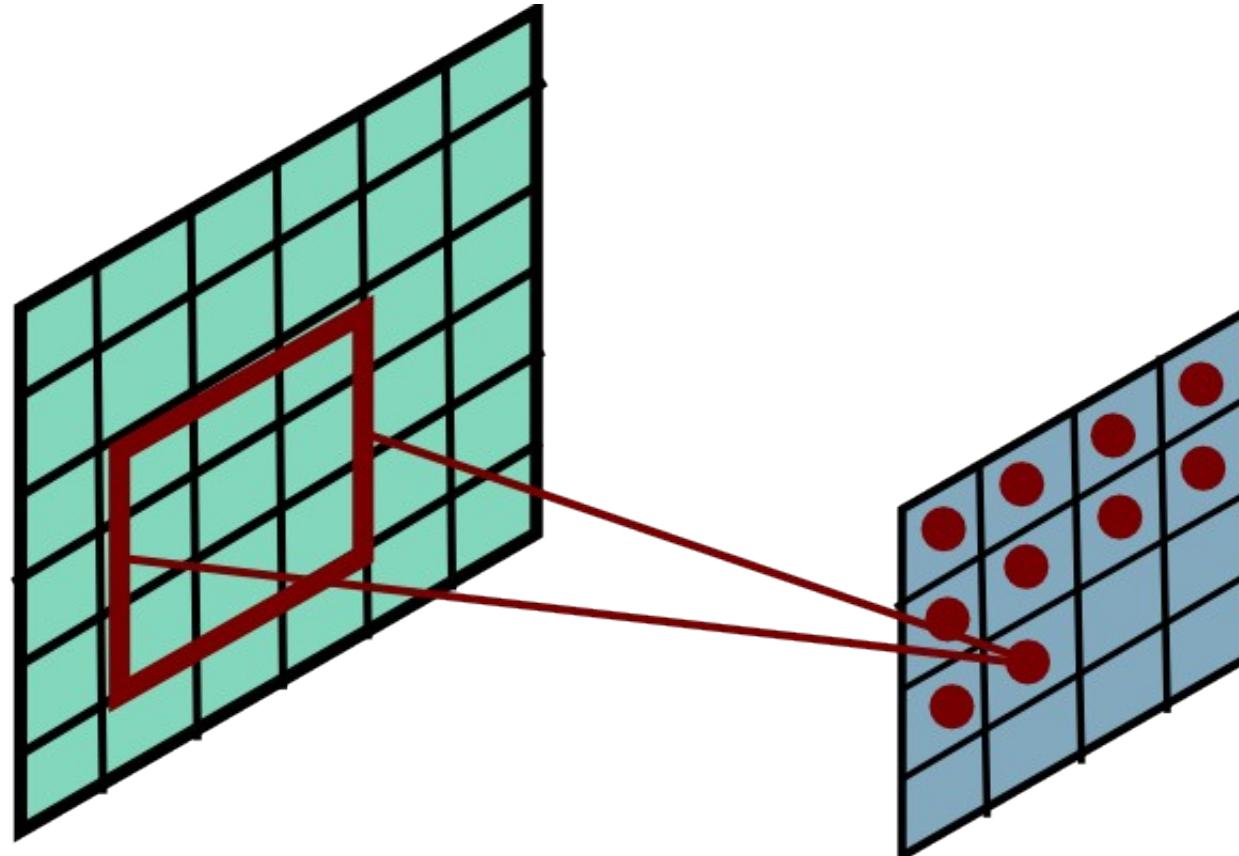
# Convolutional Layer



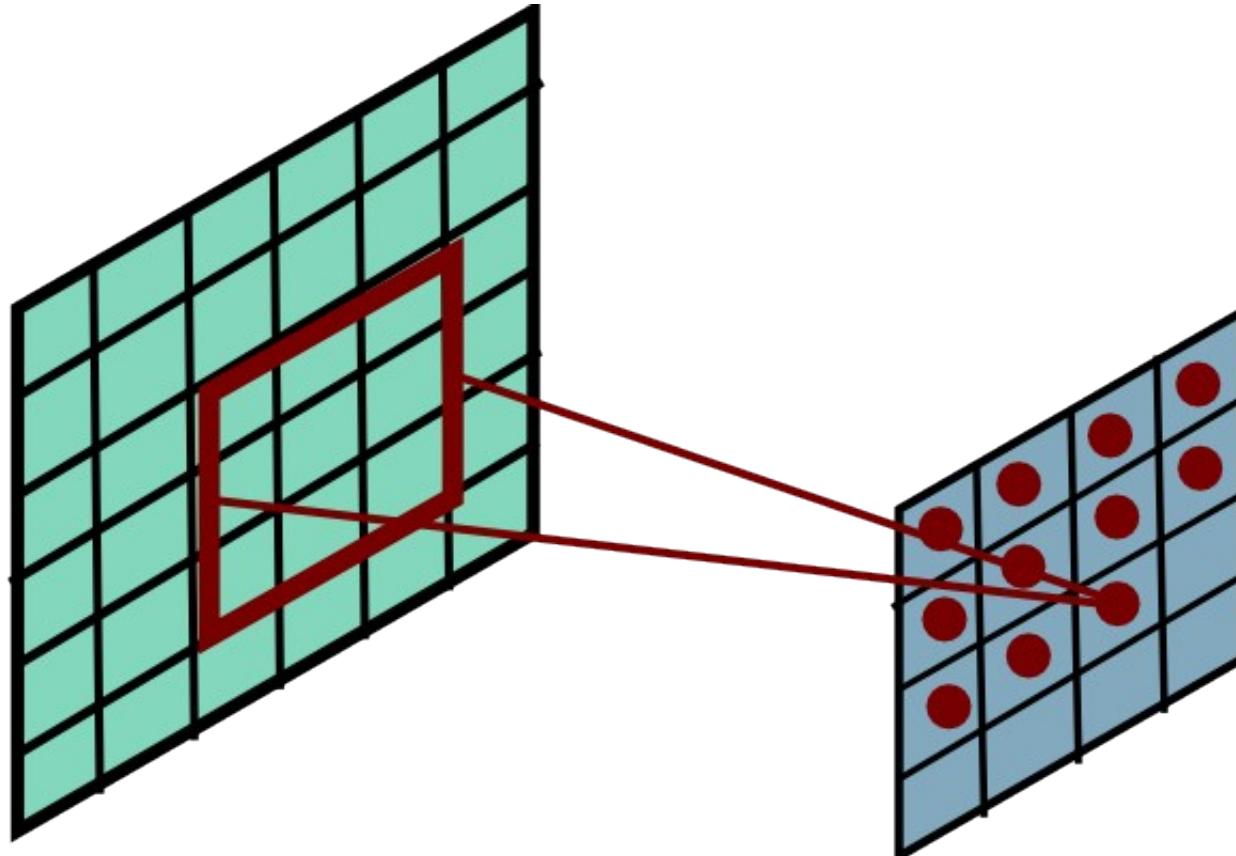
# Convolutional Layer



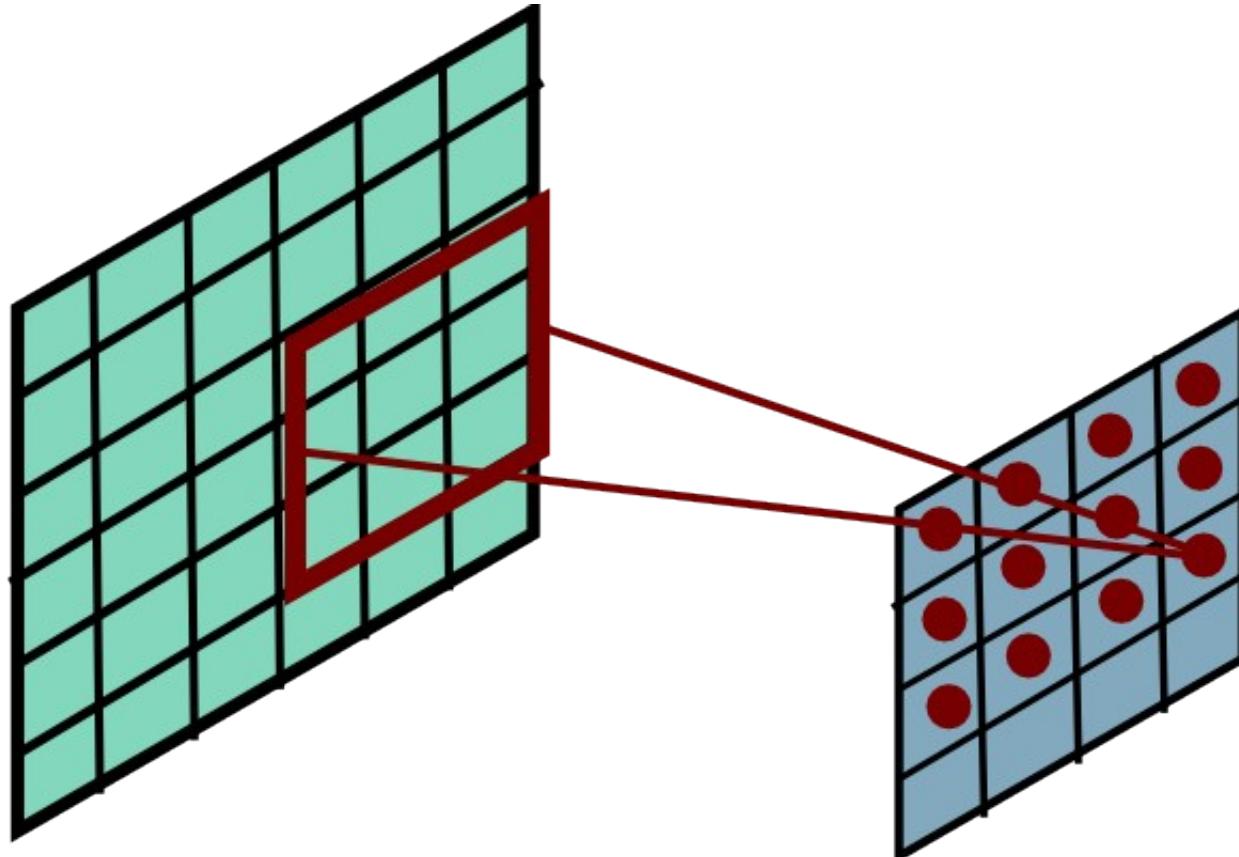
# Convolutional Layer



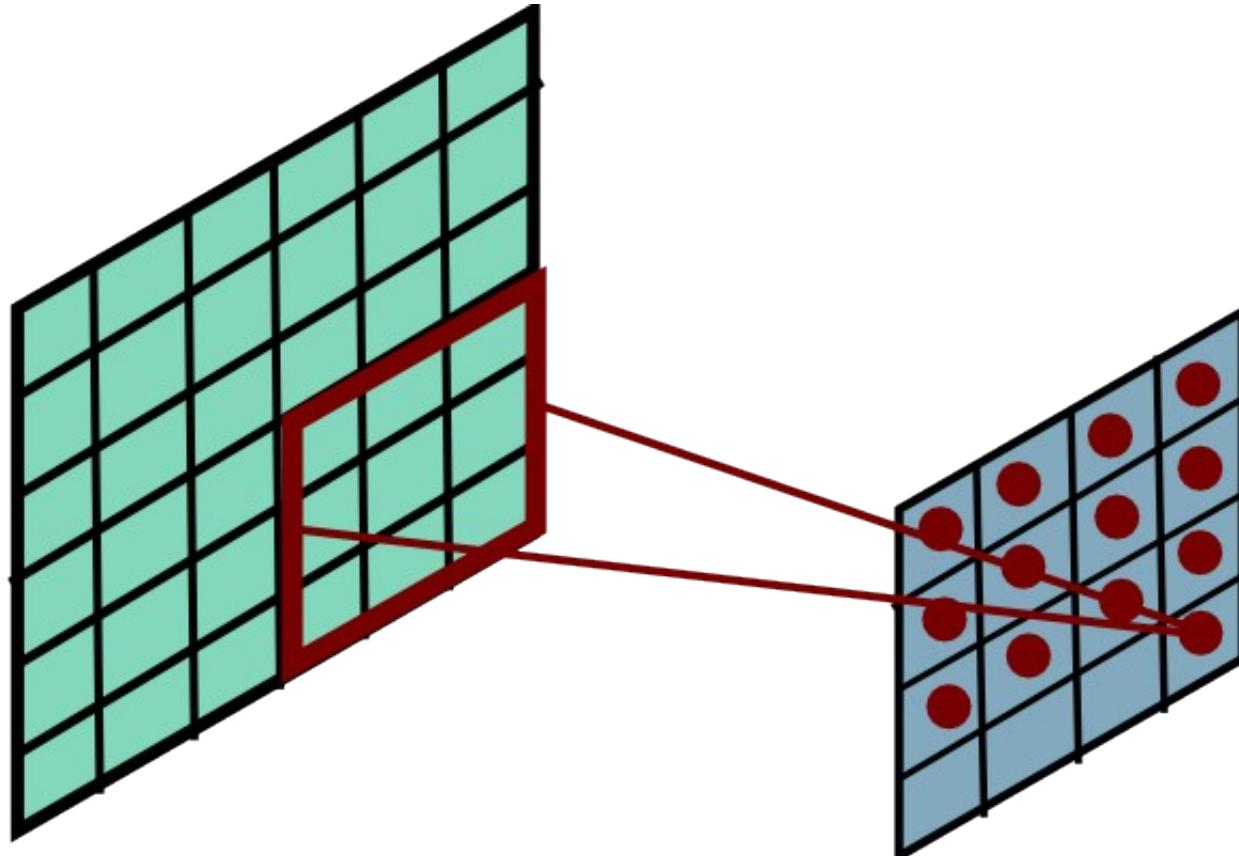
# Convolutional Layer



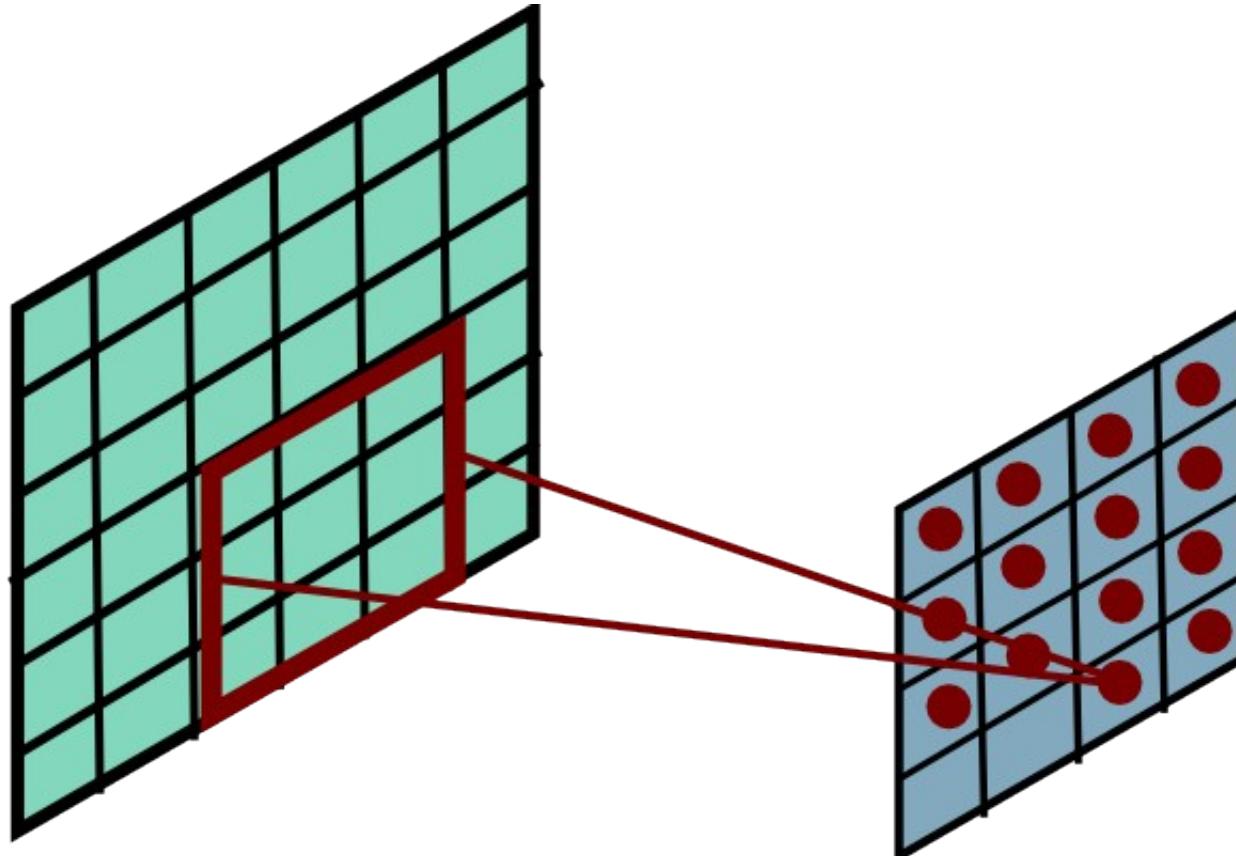
# Convolutional Layer



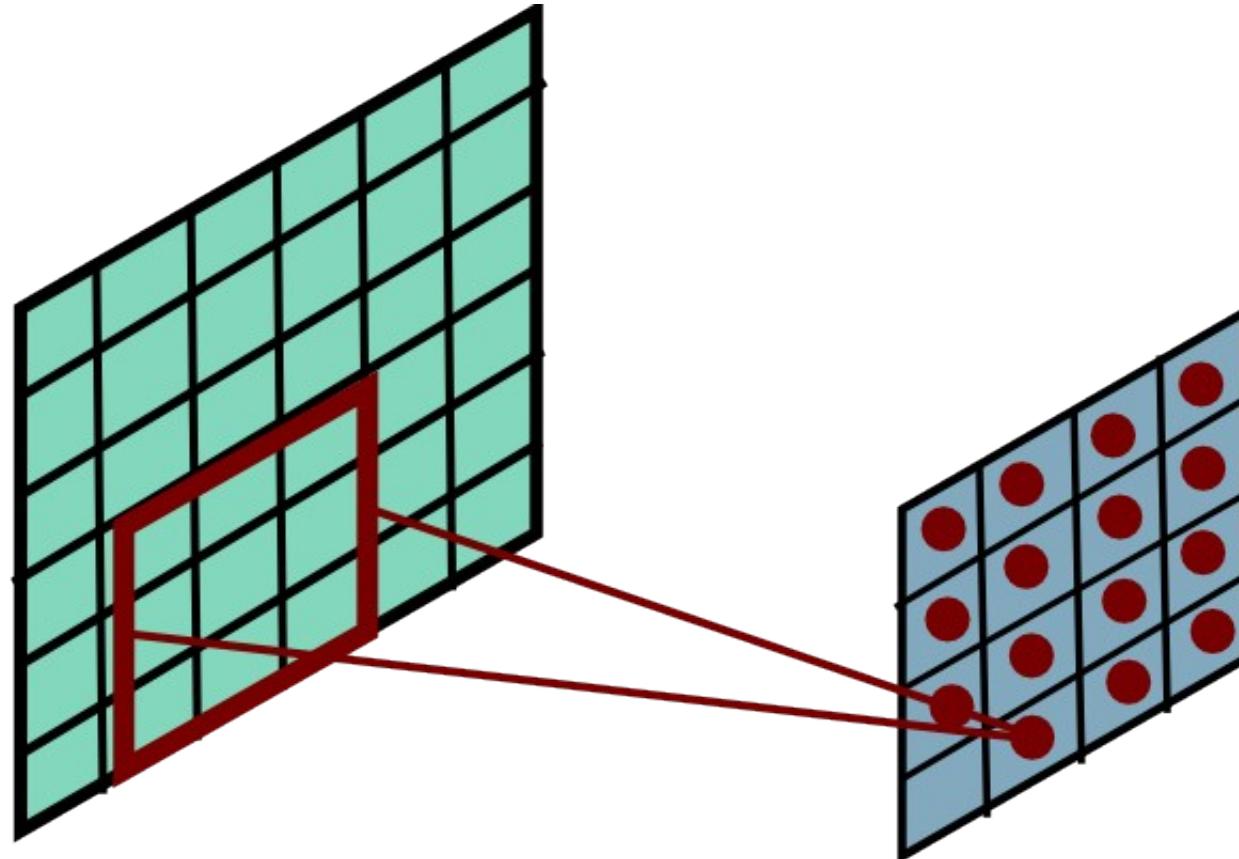
# Convolutional Layer



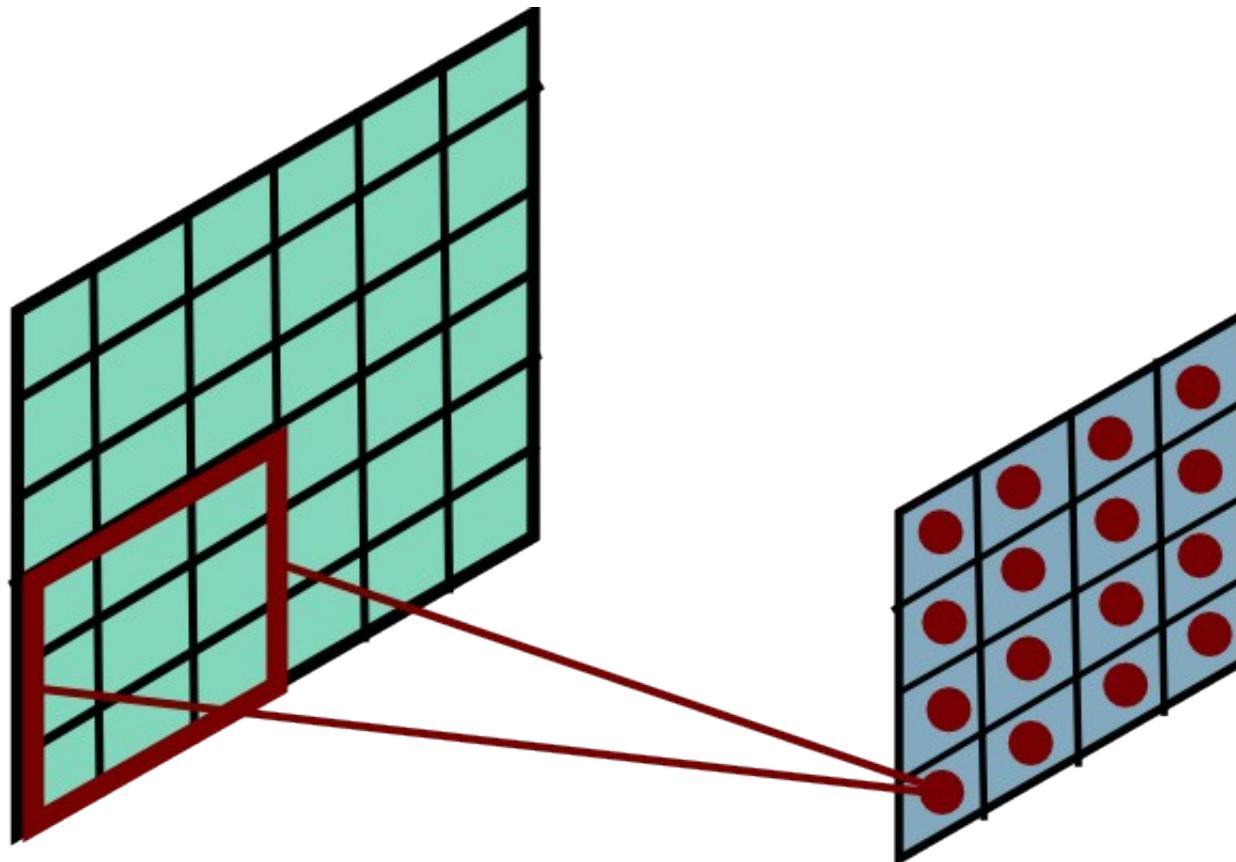
# Convolutional Layer



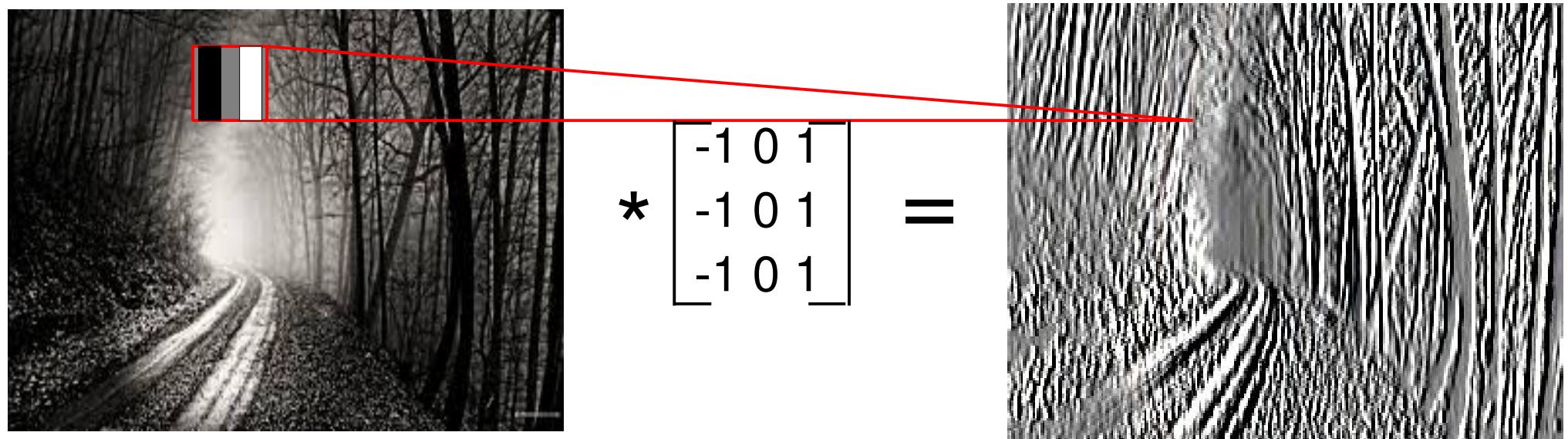
# Convolutional Layer



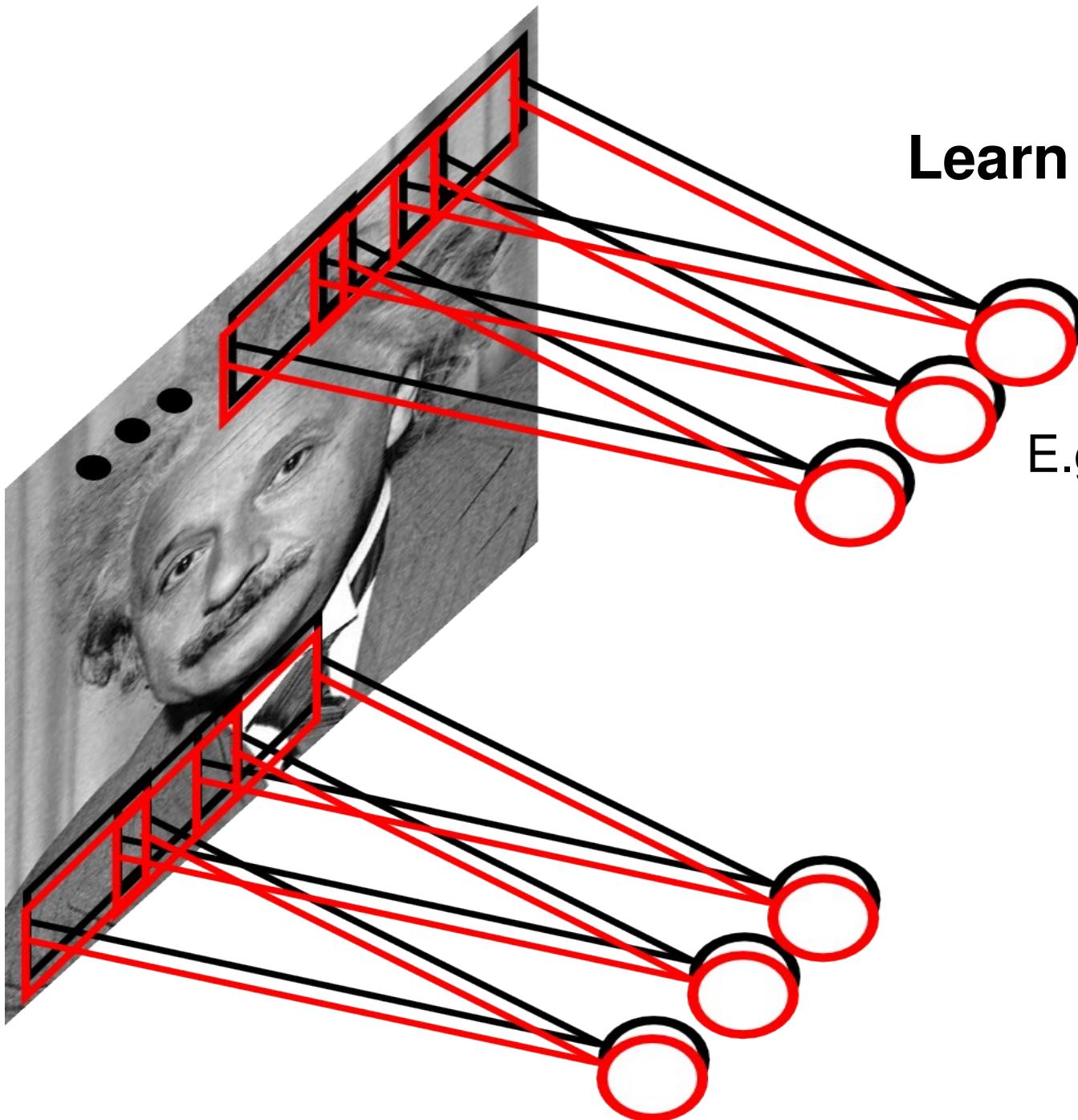
# Convolutional Layer



# Convolutional Layer



# Convolutional Layer



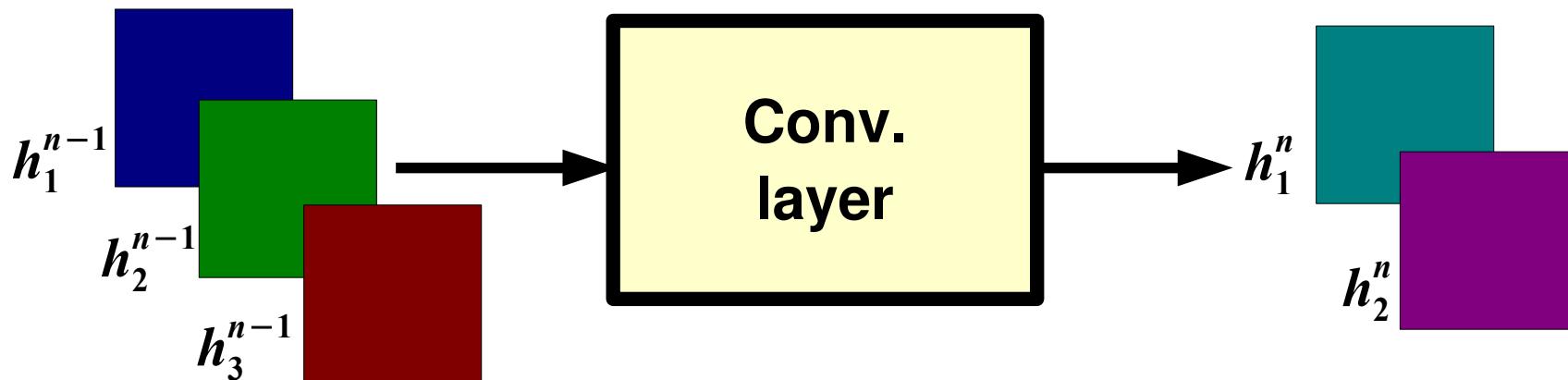
**Learn multiple filters.**

E.g.: 200x200 image  
100 Filters  
Filter size: 10x10  
10K parameters

# Convolutional Layer

$$h_j^n = \max \left( 0, \sum_{k=1}^K h_k^{n-1} * w_{kj}^n \right)$$

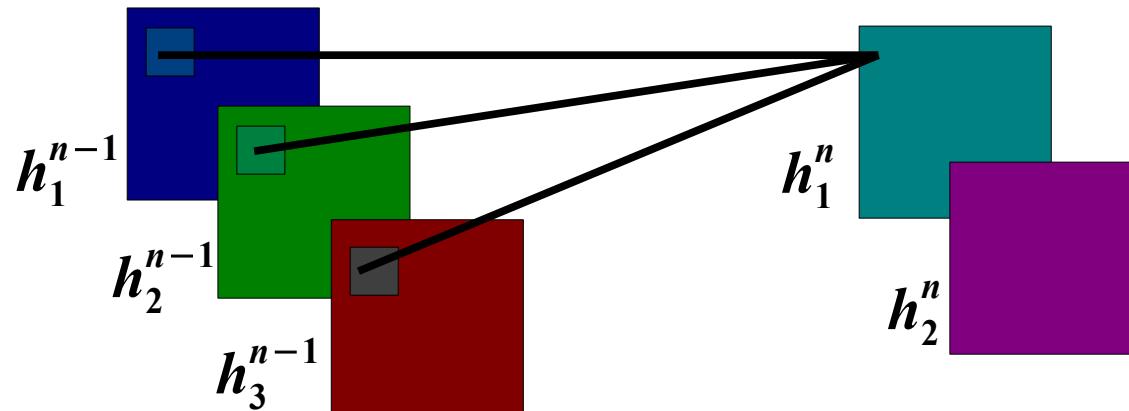
**output feature map**      **input feature map**      **kernel**



# Convolutional Layer

$$h_j^n = \max(0, \sum_{k=1}^K h_k^{n-1} * w_{kj}^n)$$

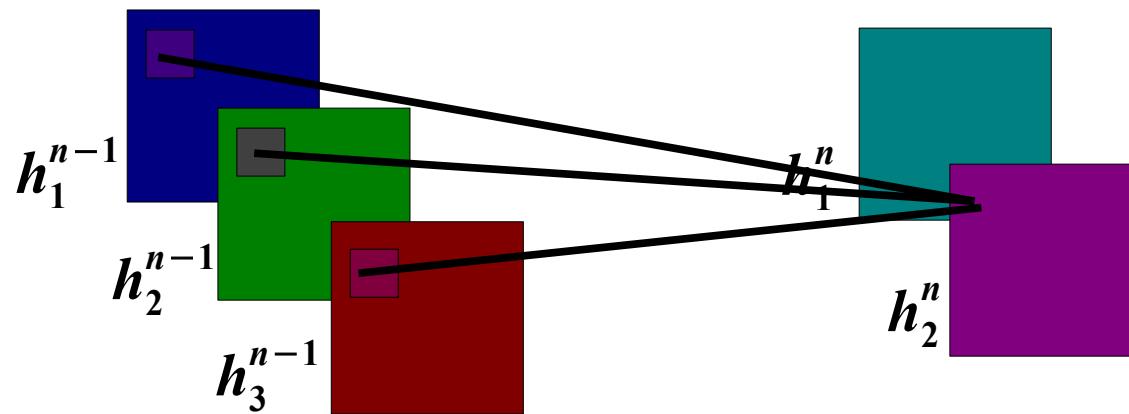
**output feature map**      **input feature map**      **kernel**



# Convolutional Layer

$$h_j^n = \max(0, \sum_{k=1}^K h_k^{n-1} * w_{kj}^n)$$

**output feature map**      **input feature map**      **kernel**



# Convolutional Layer

**Question:** What is the size of the output? What's the computational cost?

**Answer:** It is proportional to the number of filters and depends on the stride. If kernels have size  $K \times K$ , input has size  $D \times D$ , stride is 1, and there are  $M$  input feature maps and  $N$  output feature maps then:

- the input has size  $M @ D \times D$
- the output has size  $N @ (D - K + 1) \times (D - K + 1)$
- the kernels have  $M \times N \times K \times K$  coefficients (which have to be learned)
- cost:  $M \cdot K^2 \cdot N \cdot (D - K + 1)^2$

**Question:** How many feature maps? What's the size of the filters?

**Answer:** Usually, there are more output feature maps than input feature maps. Convolutional layers can increase the number of hidden units by big factors (and are expensive to compute). The size of the filters has to match the size/scale of the patterns we want to detect (task dependent).

# Key Ideas

A standard neural net applied to images:

- scales quadratically with the size of the input
- does not leverage stationarity

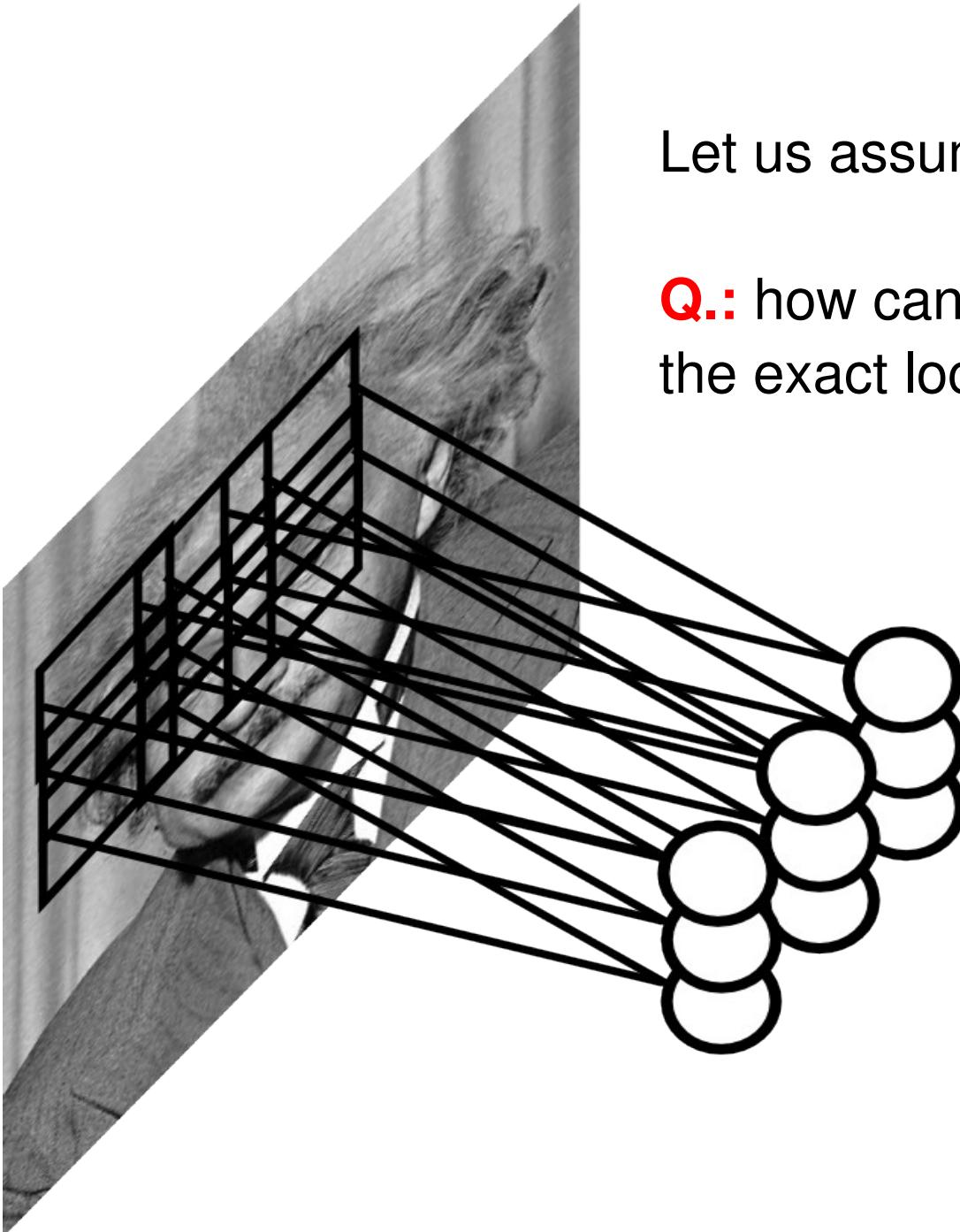
Solution:

- connect each hidden unit to a small patch of the input
- share the weight across space

This is called: **convolutional layer**.

A network with convolutional layers is called **convolutional network**.

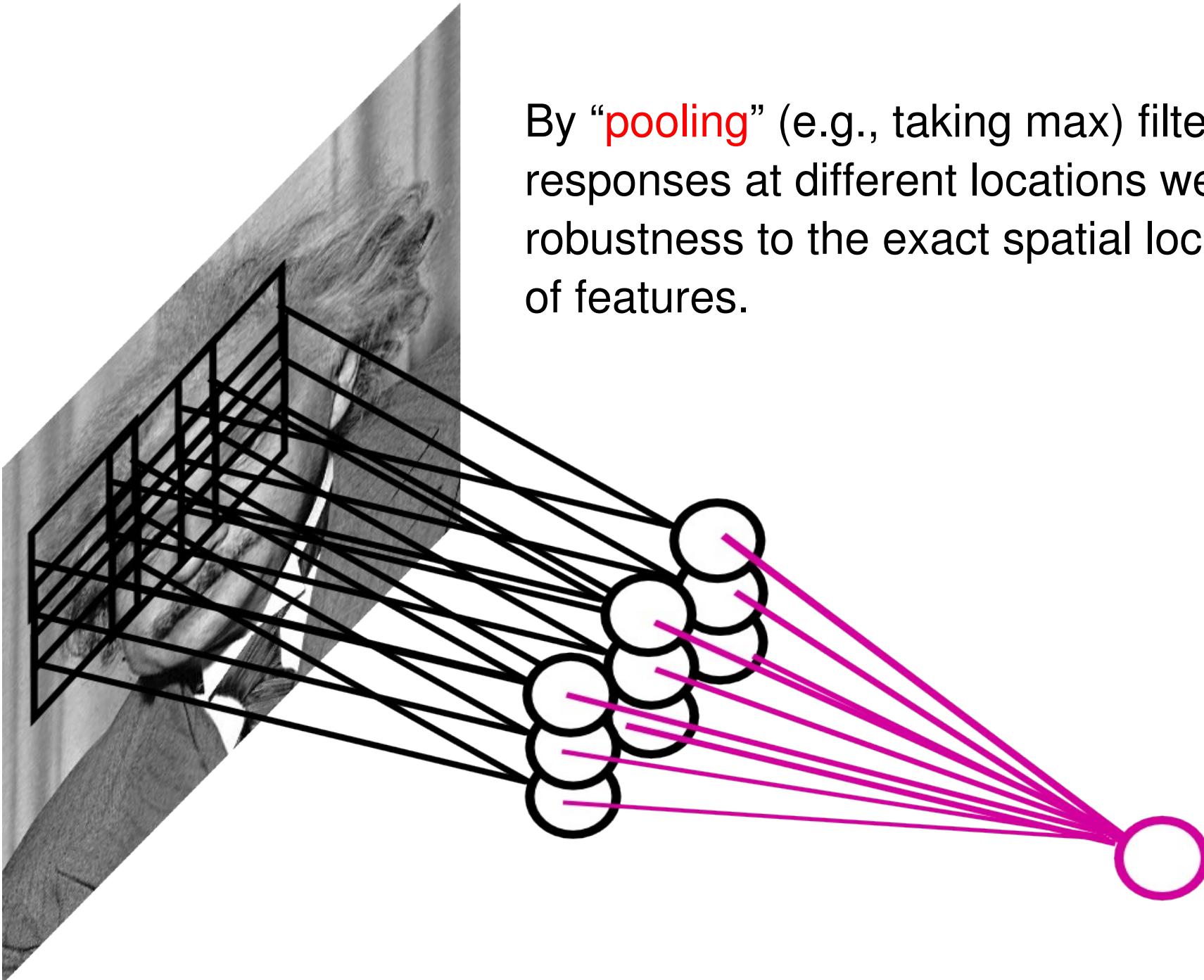
# Pooling Layer



Let us assume filter is an “eye” detector.

**Q.:** how can we make the detection robust to the exact location of the eye?

# Pooling Layer



# Pooling Layer: Examples

Max-pooling:

$$h_j^n(x, y) = \max_{\bar{x} \in N(x), \bar{y} \in N(y)} h_j^{n-1}(\bar{x}, \bar{y})$$

Average-pooling:

$$h_j^n(x, y) = 1/K \sum_{\bar{x} \in N(x), \bar{y} \in N(y)} h_j^{n-1}(\bar{x}, \bar{y})$$

L2-pooling:

$$h_j^n(x, y) = \sqrt{\sum_{\bar{x} \in N(x), \bar{y} \in N(y)} h_j^{n-1}(\bar{x}, \bar{y})^2}$$

L2-pooling over features:

$$h_j^n(x, y) = \sqrt{\sum_{k \in N(j)} h_k^{n-1}(x, y)^2}$$

# Pooling Layer

**Question:** What is the size of the output? What's the computational cost?

**Answer:** The size of the output depends on the stride between the pools. For instance, if pools do not overlap and have size  $K \times K$ , and the input has size  $D \times D$  with  $M$  input feature maps, then:

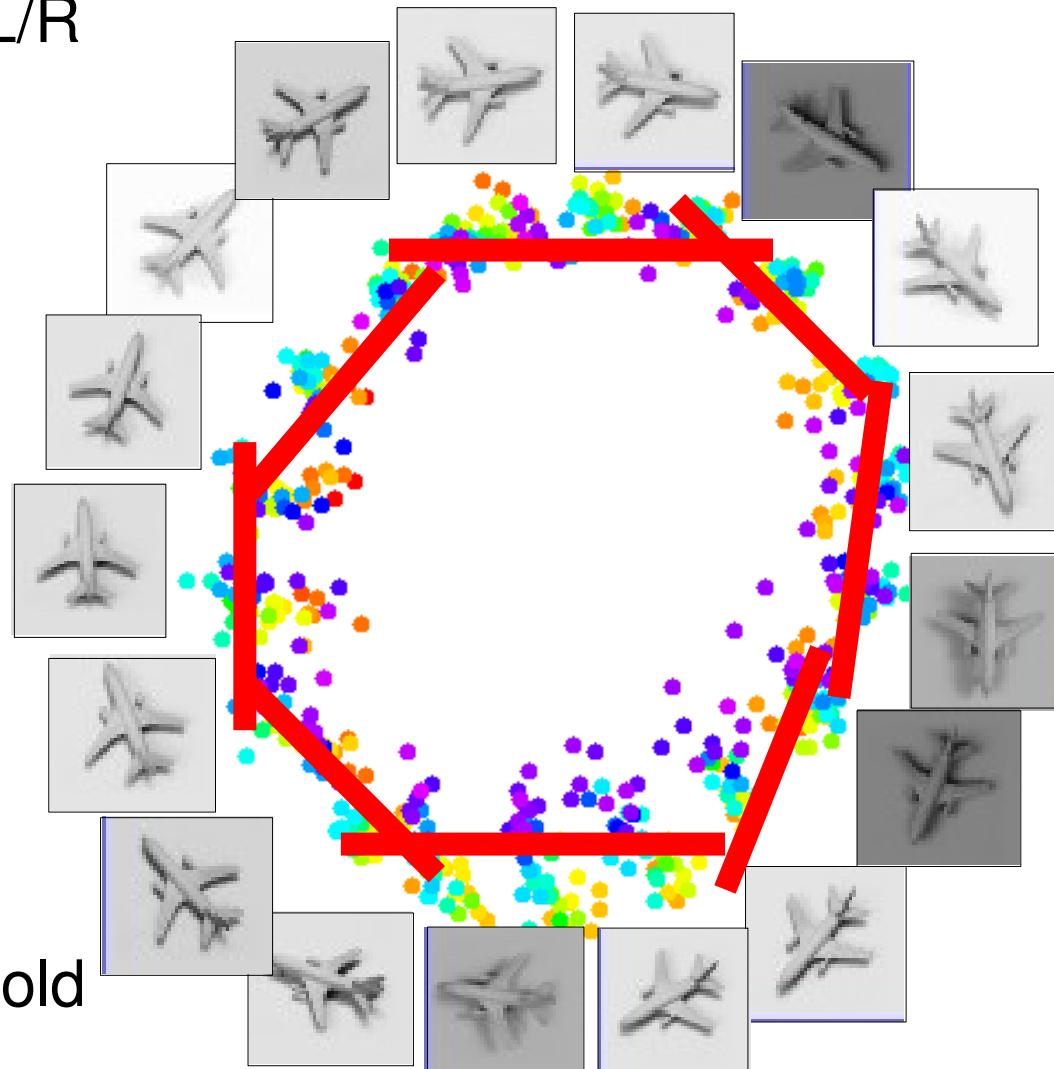
- output is  $M @ (D/K) \times (D/K)$
- the computational cost is proportional to the size of the input (negligible compared to a convolutional layer)

**Question:** How should I set the size of the pools?

**Answer:** It depends on how much “invariant” or robust to distortions we want the representation to be. It is best to pool slowly (via a few stacks of conv-pooling layers).

# Pooling Layer: Interpretation

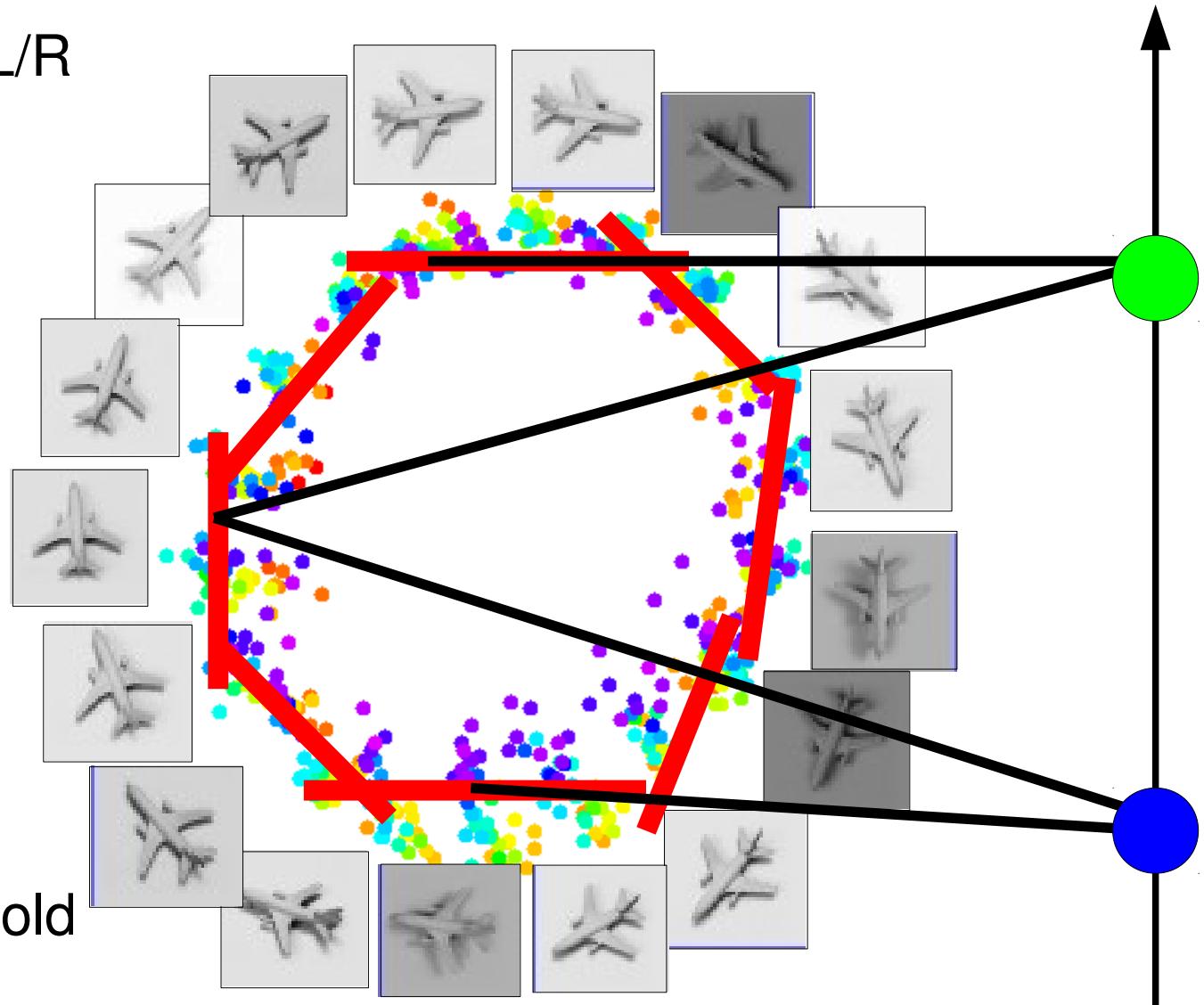
Task: detect orientation L/R



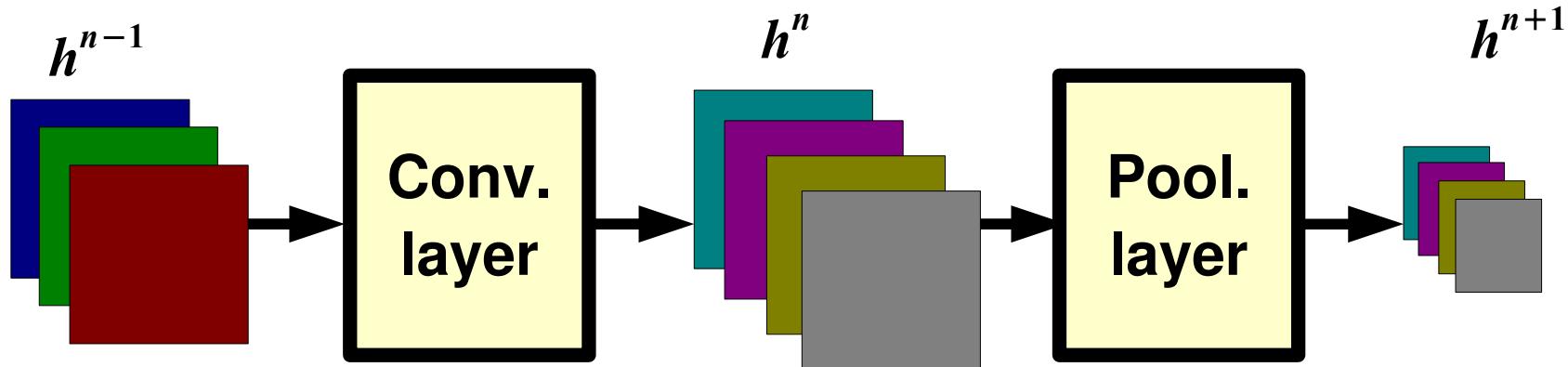
Conv layer:  
linearizes manifold

# Pooling Layer: Interpretation

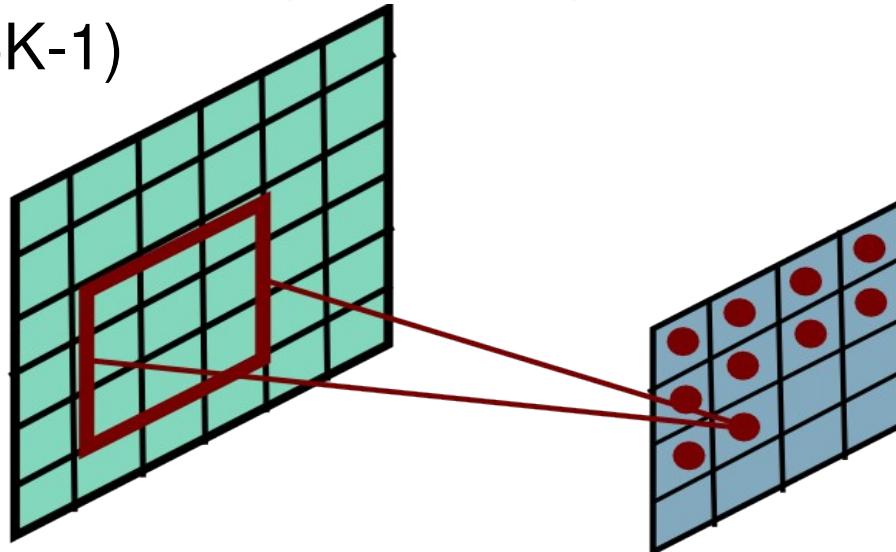
Task: detect orientation L/R



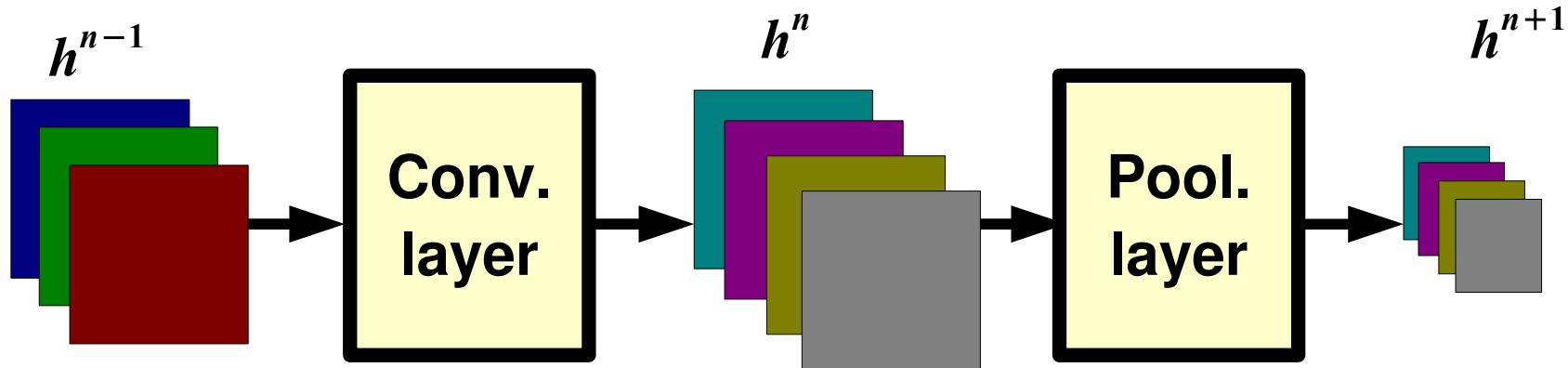
# Pooling Layer: Receptive Field Size



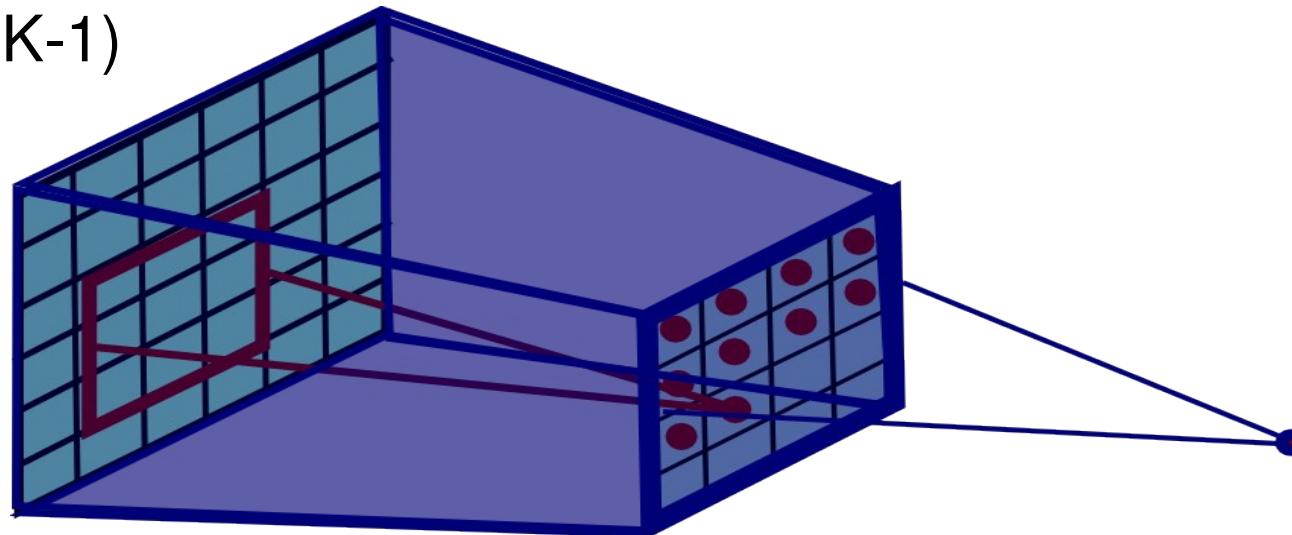
If convolutional filters have size  $K \times K$  and stride 1, and pooling layer has pools of size  $P \times P$ , then each unit in the pooling layer depends upon a patch (at the input of the preceding conv. layer) of size:  
 $(P+K-1) \times (P+K-1)$



# Pooling Layer: Receptive Field Size

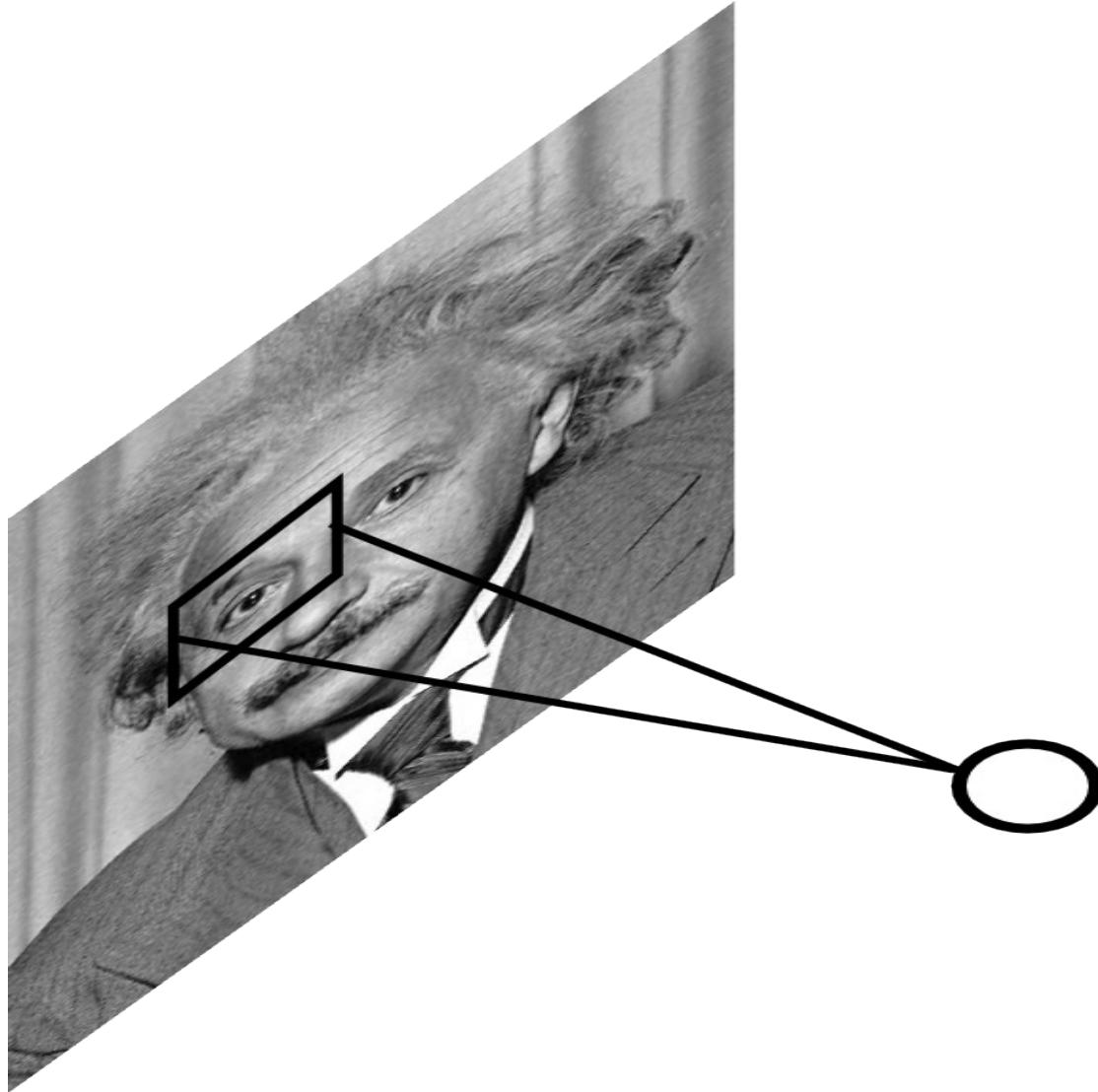


If convolutional filters have size  $K \times K$  and stride 1, and pooling layer has pools of size  $P \times P$ , then each unit in the pooling layer depends upon a patch (at the input of the preceding conv. layer) of size:  
 $(P+K-1) \times (P+K-1)$



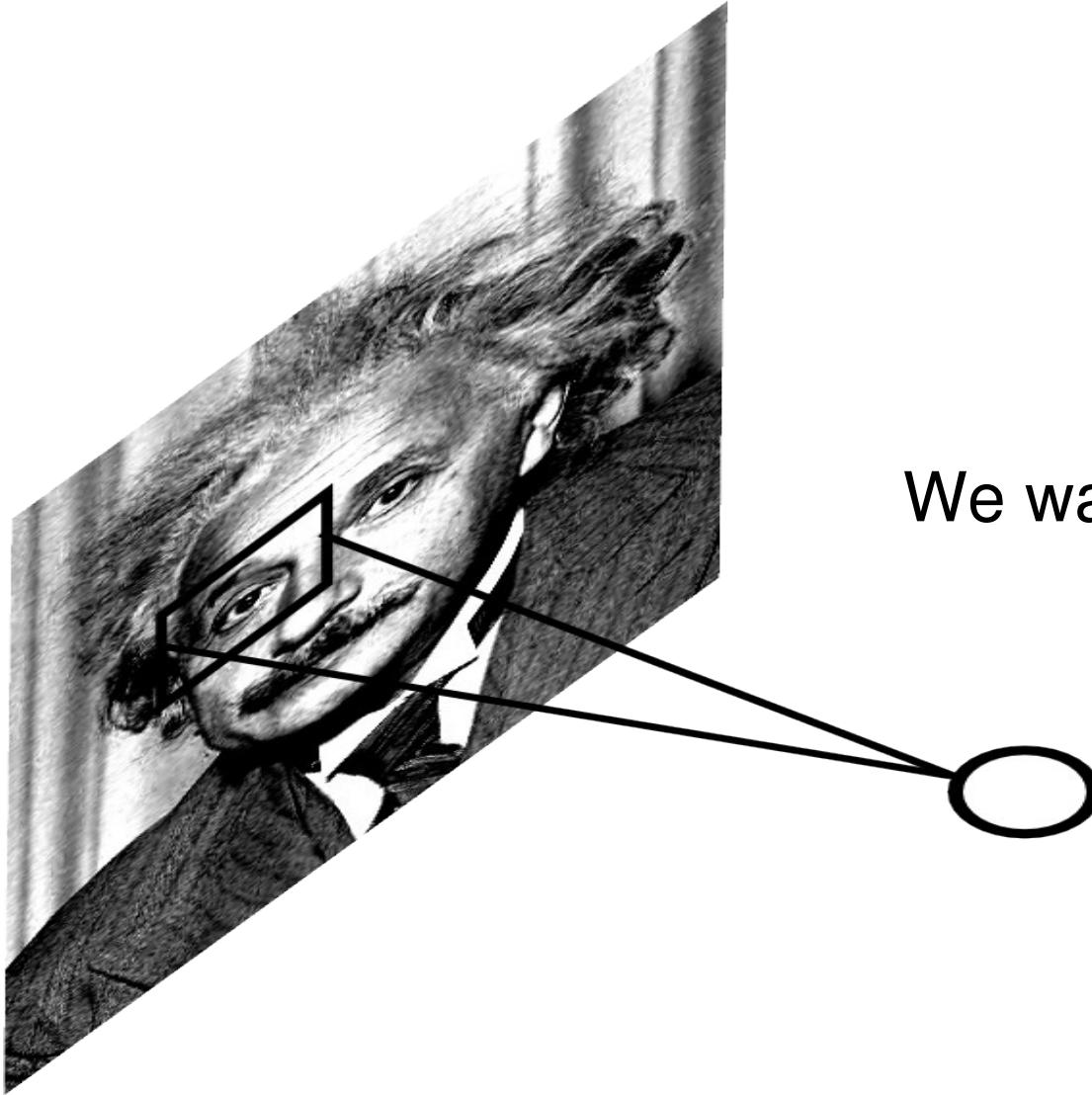
# Local Contrast Normalization

$$h^{i+1}(x, y) = \frac{h^i(x, y) - m^i(N(x, y))}{\sigma^i(N(x, y))}$$



# Local Contrast Normalization

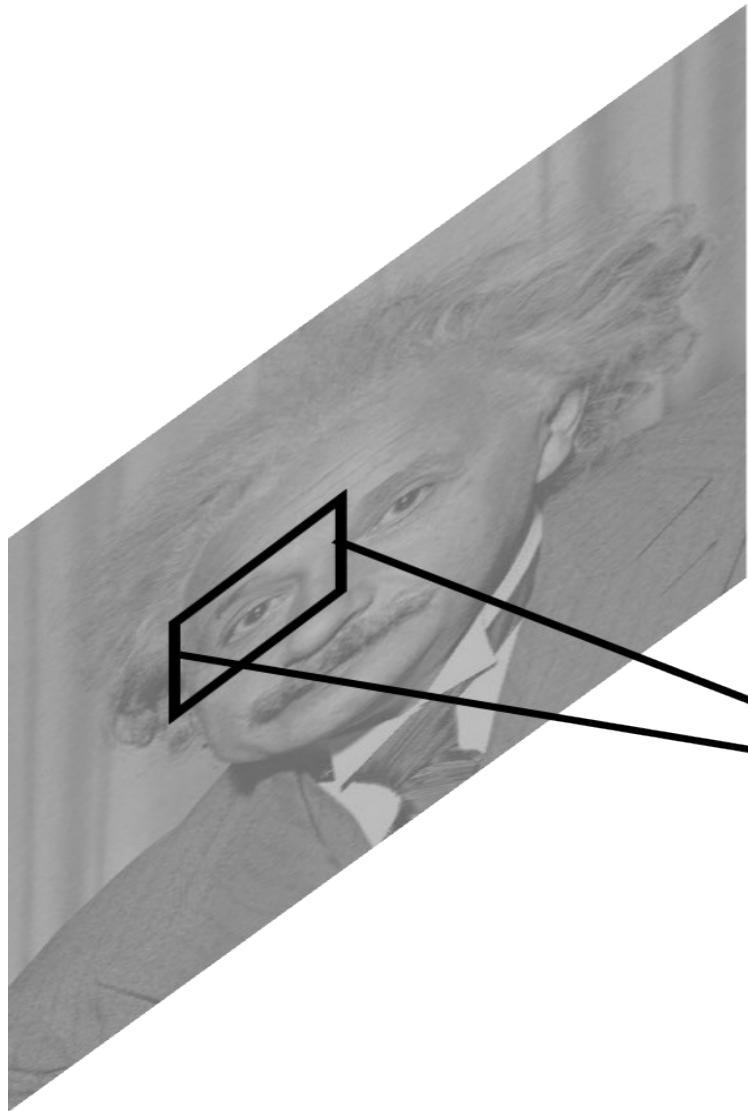
$$h^{i+1}(x, y) = \frac{h^i(x, y) - m^i(N(x, y))}{\sigma^i(N(x, y))}$$



We want the same response.

# Local Contrast Normalization

$$h^{i+1}(x, y) = \frac{h^i(x, y) - m^i(N(x, y))}{\sigma^i(N(x, y))}$$



Performed also across features  
and in the higher layers..

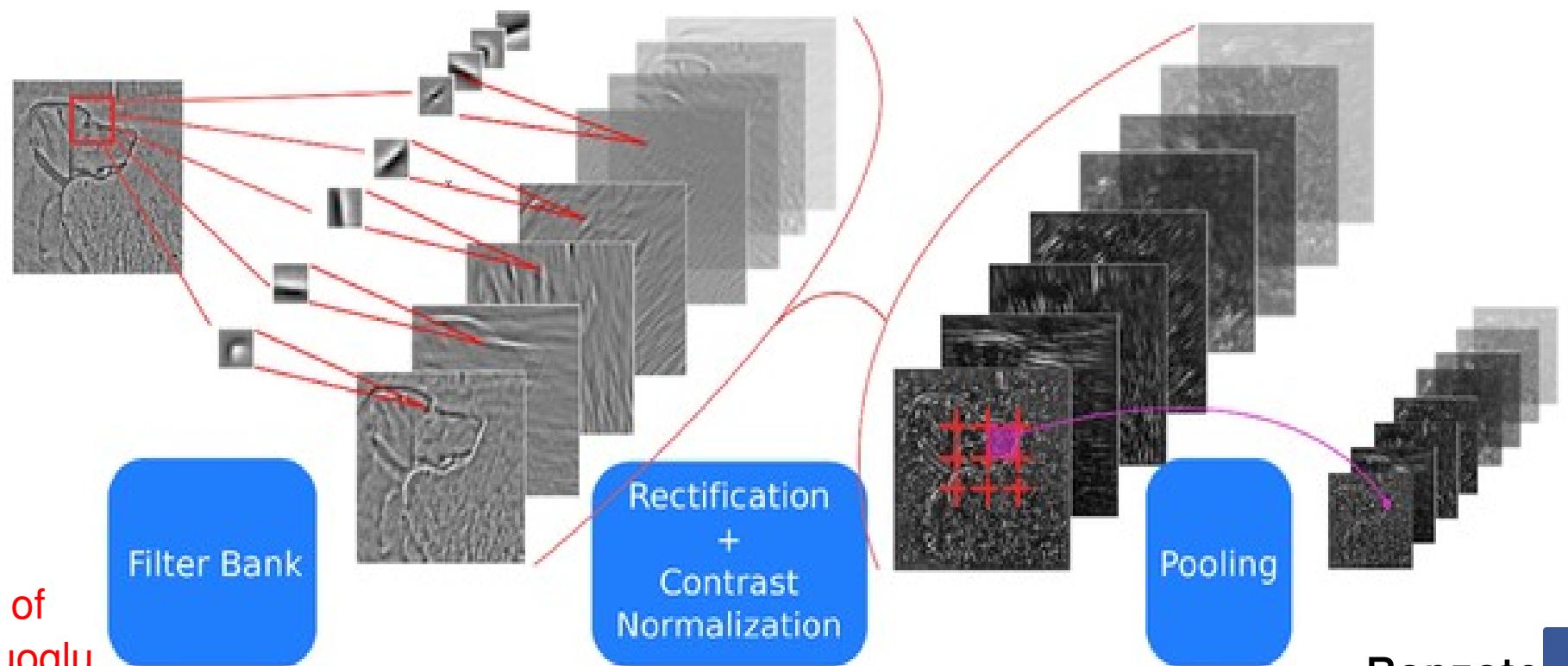
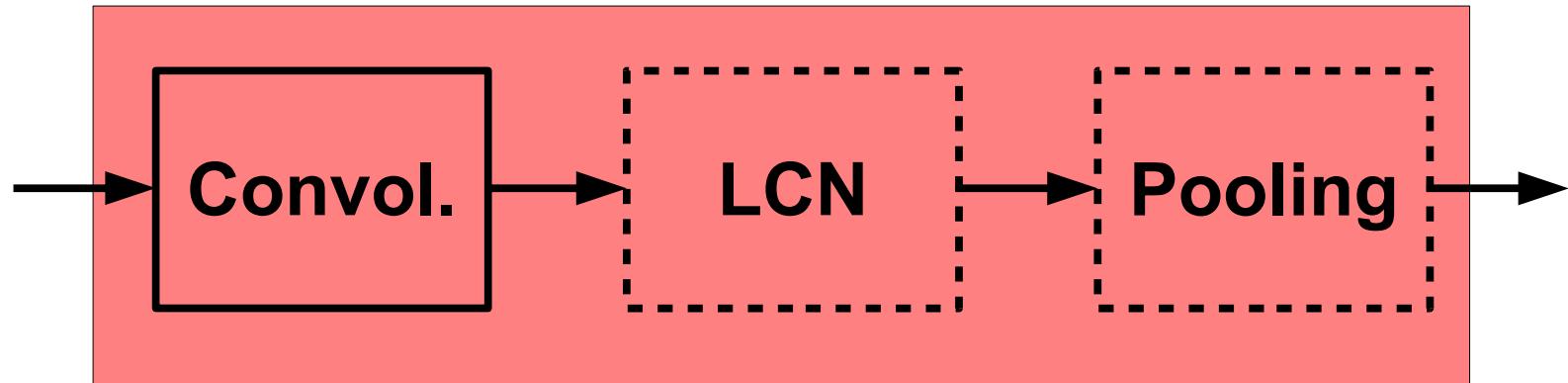
Effects:

- improves invariance
- improves optimization
- increases sparsity

**Note:** computational cost is negligible w.r.t. conv. layer.

# ConvNets: Typical Stage

One stage (zoom)

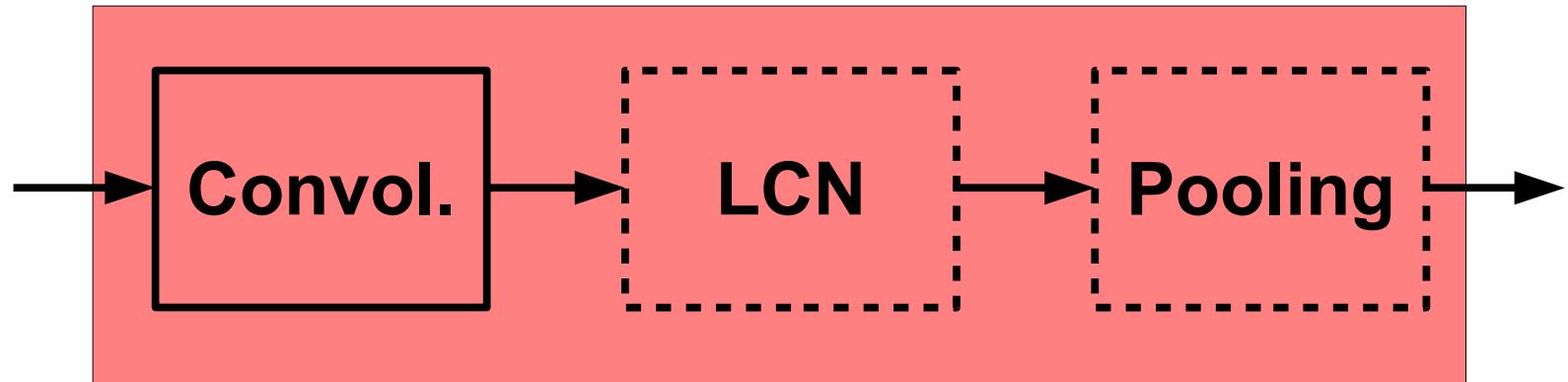


courtesy of  
K. Kavukcuoglu

Ranzato

# ConvNets: Typical Stage

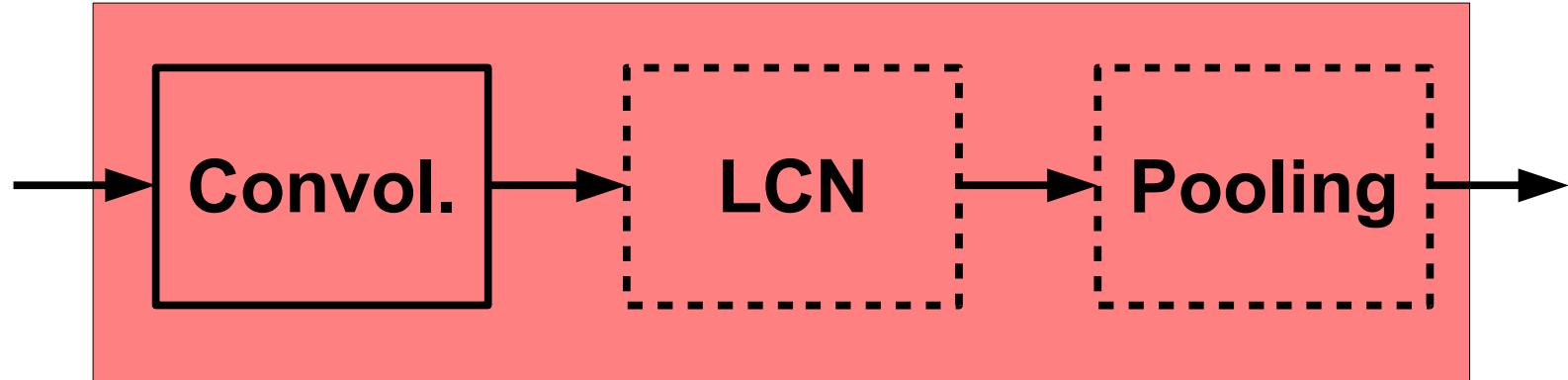
One stage (zoom)



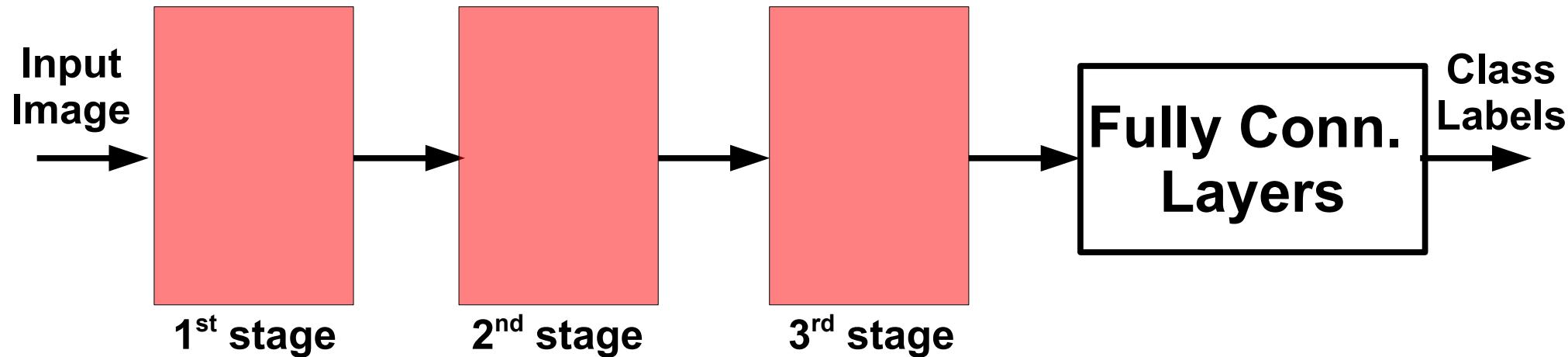
Conceptually similar to: SIFT, HoG, etc.

# ConvNets: Typical Architecture

One stage (zoom)

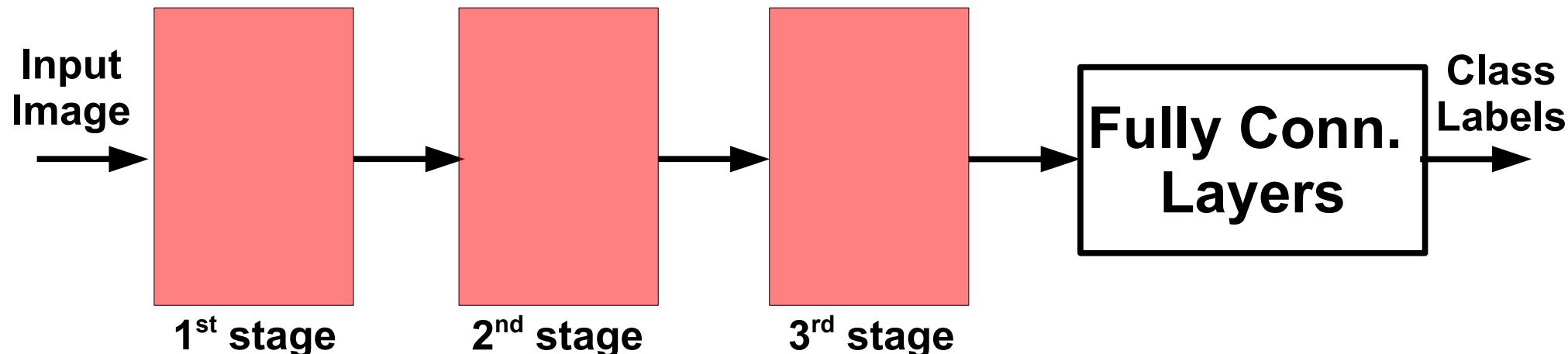


Whole system



# ConvNets: Typical Architecture

Whole system



Conceptually similar to:

SIFT → K-Means → Pyramid Pooling → SVM

Lazebnik et al. "...Spatial Pyramid Matching..." CVPR 2006

SIFT → Fisher Vect. → Pooling → SVM

Sanchez et al. "Image classification with F.V.: Theory and practice" IJCV 2012

# ConvNets: Training

All layers are differentiable (a.e.).

We can use standard back-propagation.

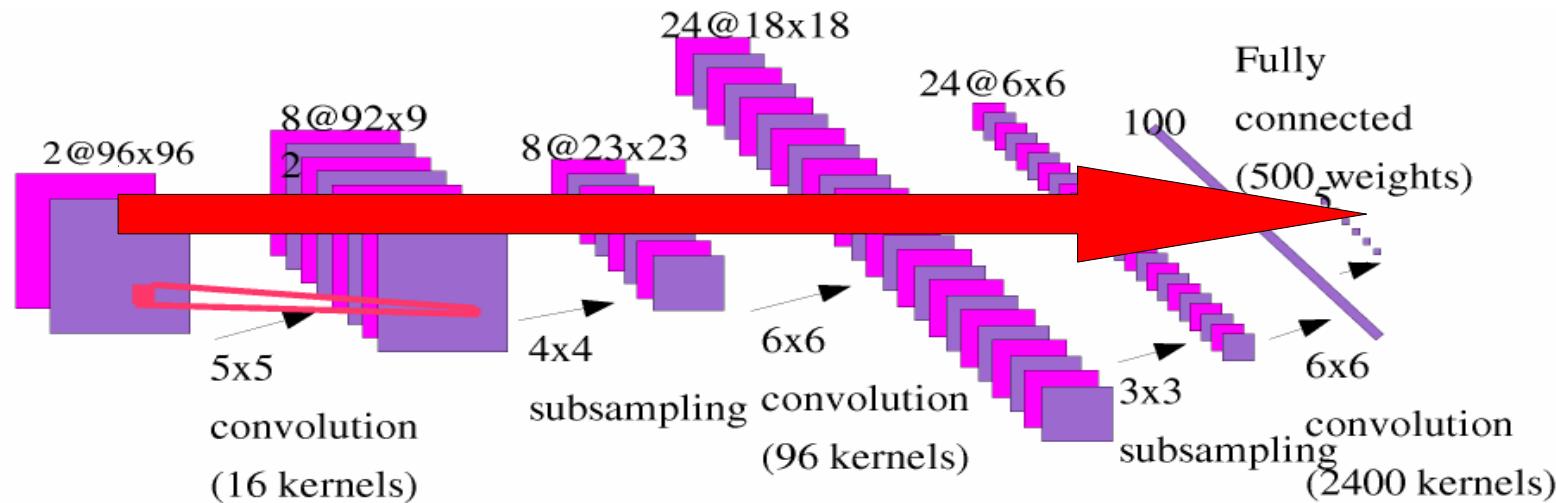
## Algorithm:

**Given a small mini-batch**

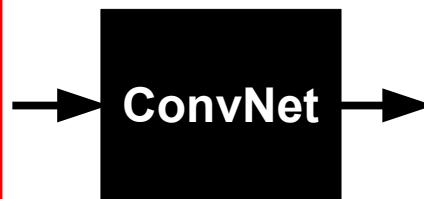
- **F-PROP**
- **B-PROP**
- **PARAMETER UPDATE**

# ConvNets: Test

At test time, run only is forward mode (FPROP).



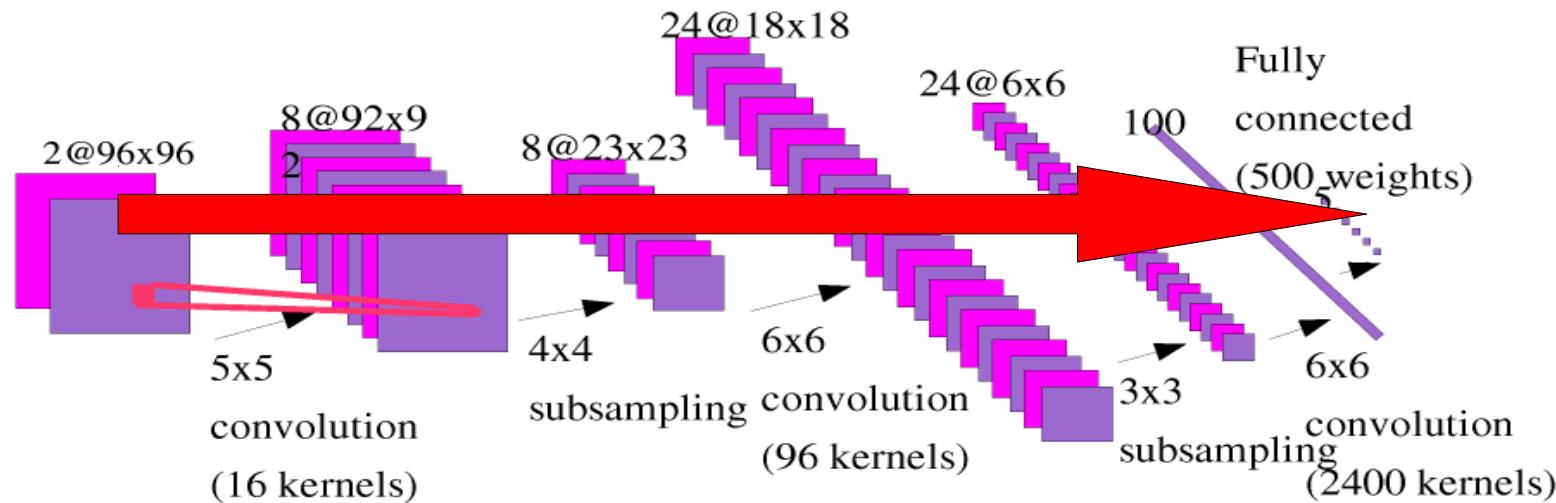
Naturally, convnet can process larger images at little cost.



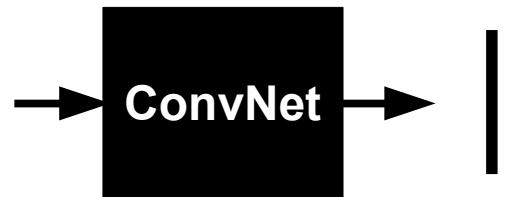
Traditional methods use inefficient sliding windows.

# ConvNets: Test

At test time, run only is forward mode (FPROP).



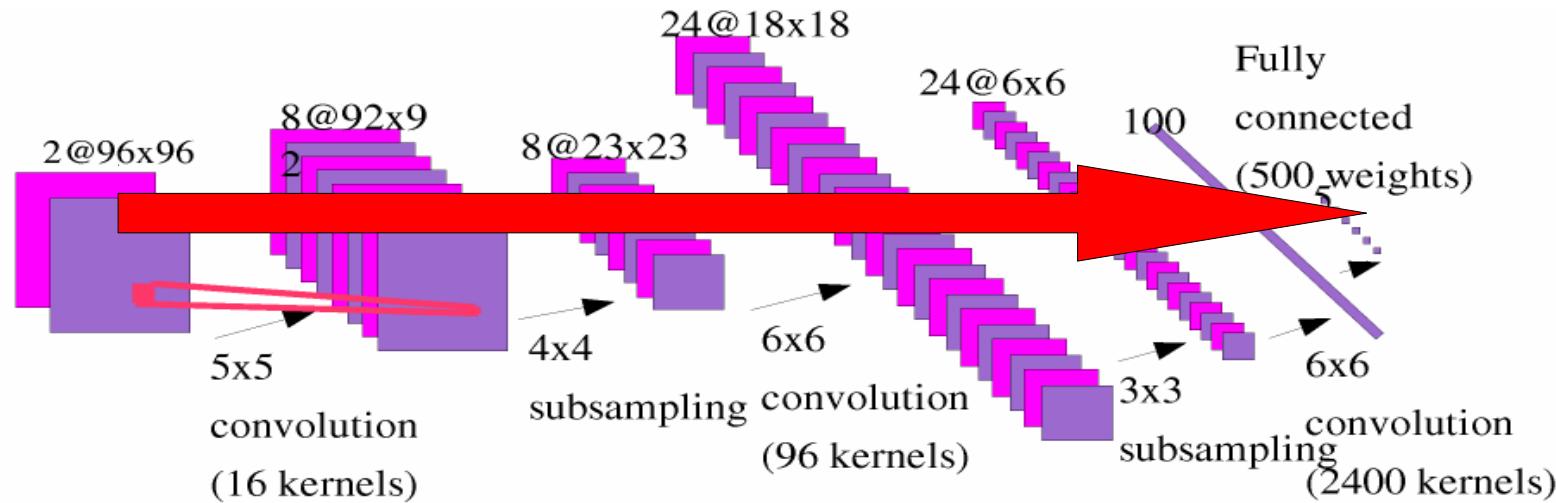
Naturally, convnet can process larger images at little cost.



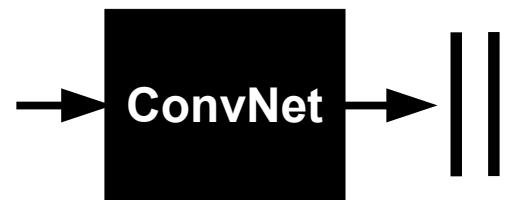
Traditional methods use inefficient sliding windows.

# ConvNets: Test

At test time, run only is forward mode (FPROP).



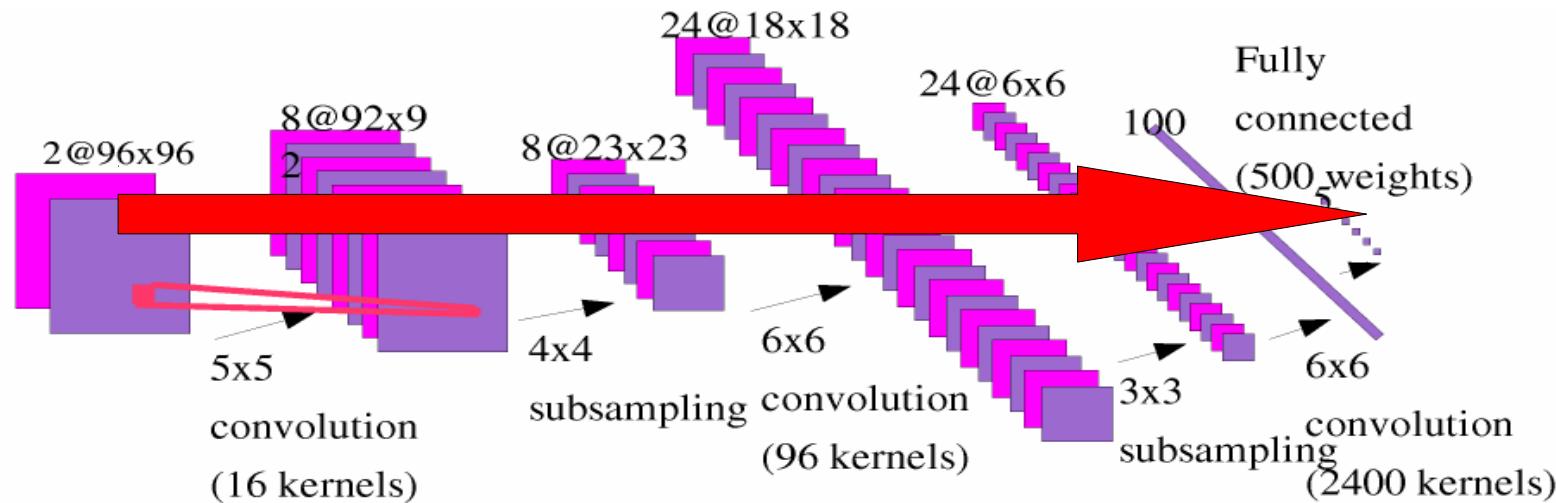
Naturally, convnet can process larger images at little cost.



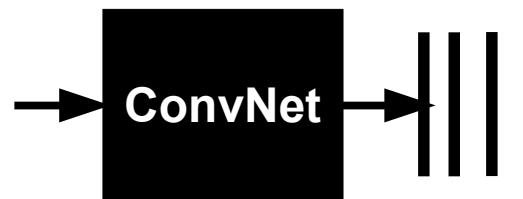
Traditional methods use inefficient sliding windows.

# ConvNets: Test

At test time, run only is forward mode (FPROP).



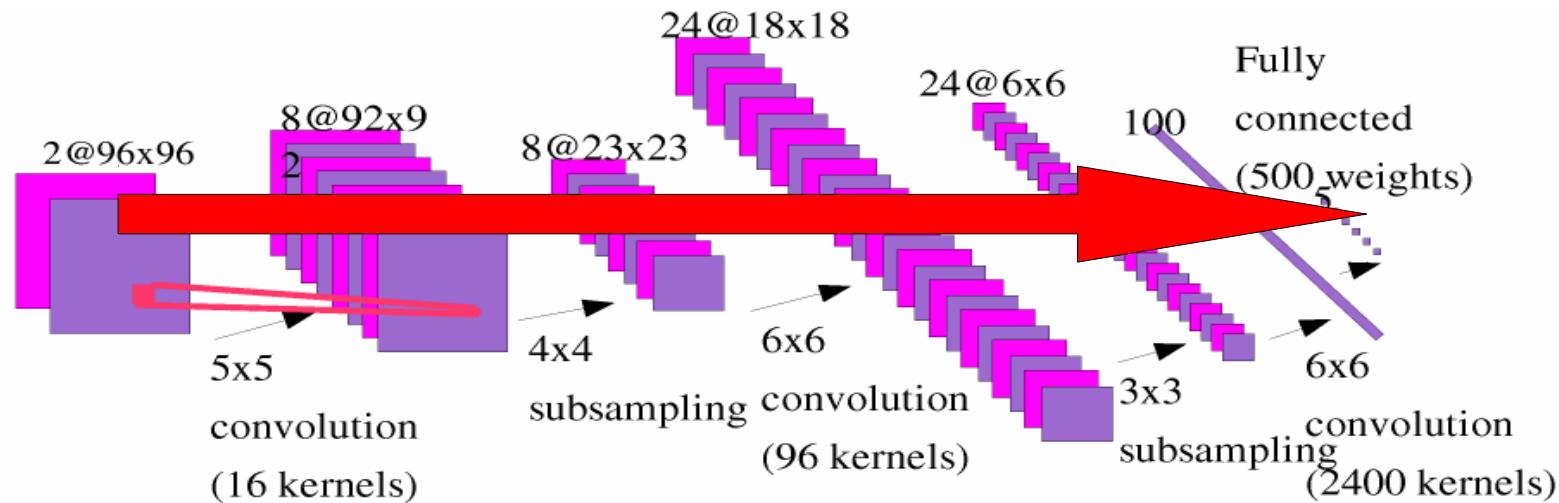
Naturally, convnet can process larger images at little cost.



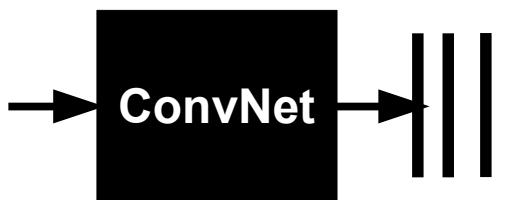
Traditional methods use inefficient sliding windows.

# ConvNets: Test

At test time, run only is forward mode (FPROP).



Naturally, convnet can process larger images at little cost.



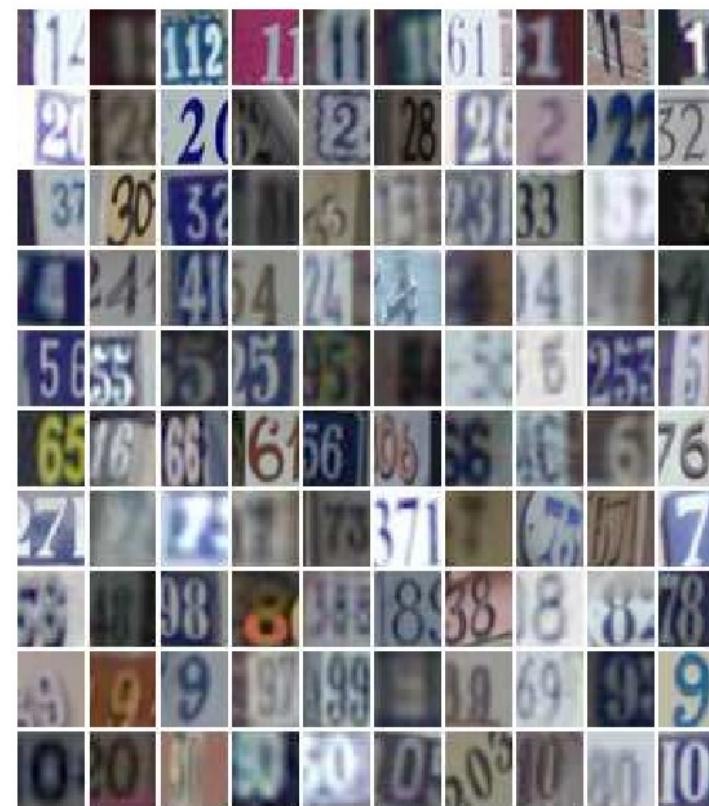
ConvNet: unrolls convolutions over bigger images and produces outputs at several locations.

# Outline

- Supervised Neural Networks
- Convolutional Neural Networks
- Examples
- Tips

# CONV NETS: EXAMPLES

- OCR / House number & Traffic sign classification



Ciresan et al. "MCDNN for image classification" CVPR 2012

Wan et al. "Regularization of neural networks using dropconnect" ICML 2013

82

Jaderberg et al. "Synthetic data and ANN for natural scene text recognition" arXiv 2014

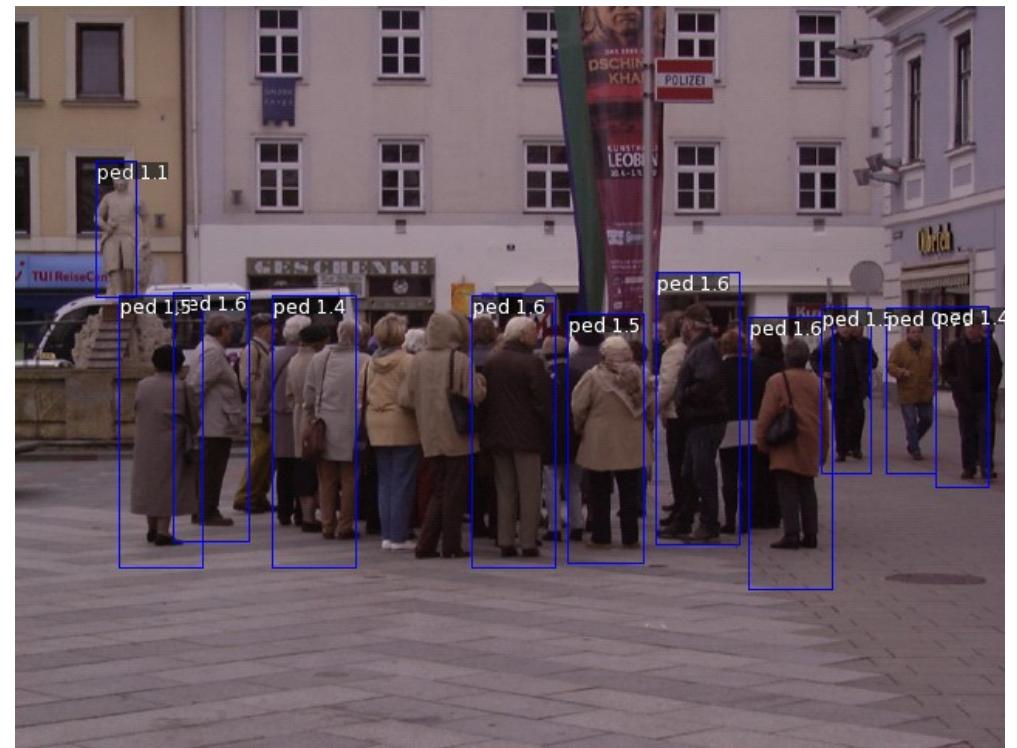
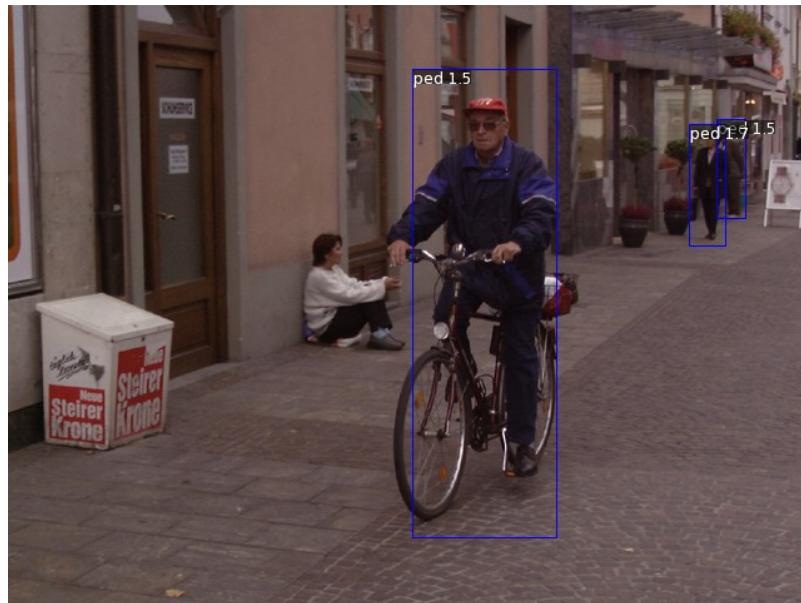
# CONV NETS: EXAMPLES

## - Texture classification



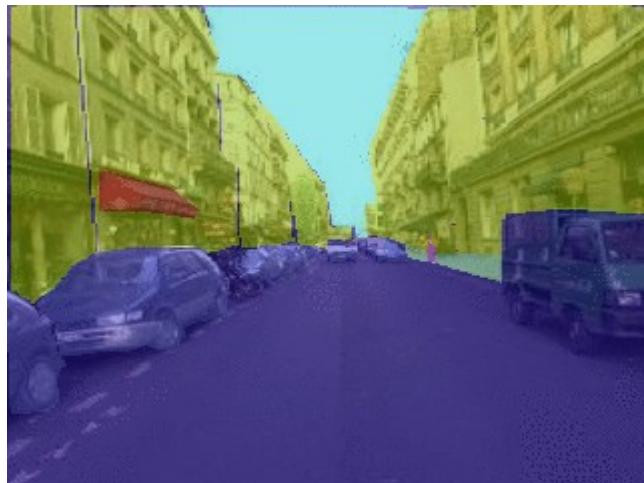
# CONV NETS: EXAMPLES

## - Pedestrian detection



# CONV NETS: EXAMPLES

## - Scene Parsing



A vibrant collage featuring a variety of objects and scenes. In the upper left, a yellow building is labeled 'building'. To its right, a blue sky is labeled 'sky'. The center of the image shows a red building labeled 'building' and a green awning labeled 'awning'. Below these, several cars are labeled 'car' in different colors. A road is labeled 'road' in the lower center. On the far right, there's a bridge labeled 'bridge' and a tree labeled 'tree'. The bottom left contains a green rock labeled 'rock'. The overall composition is a dense, multi-colored patchwork of everyday items.

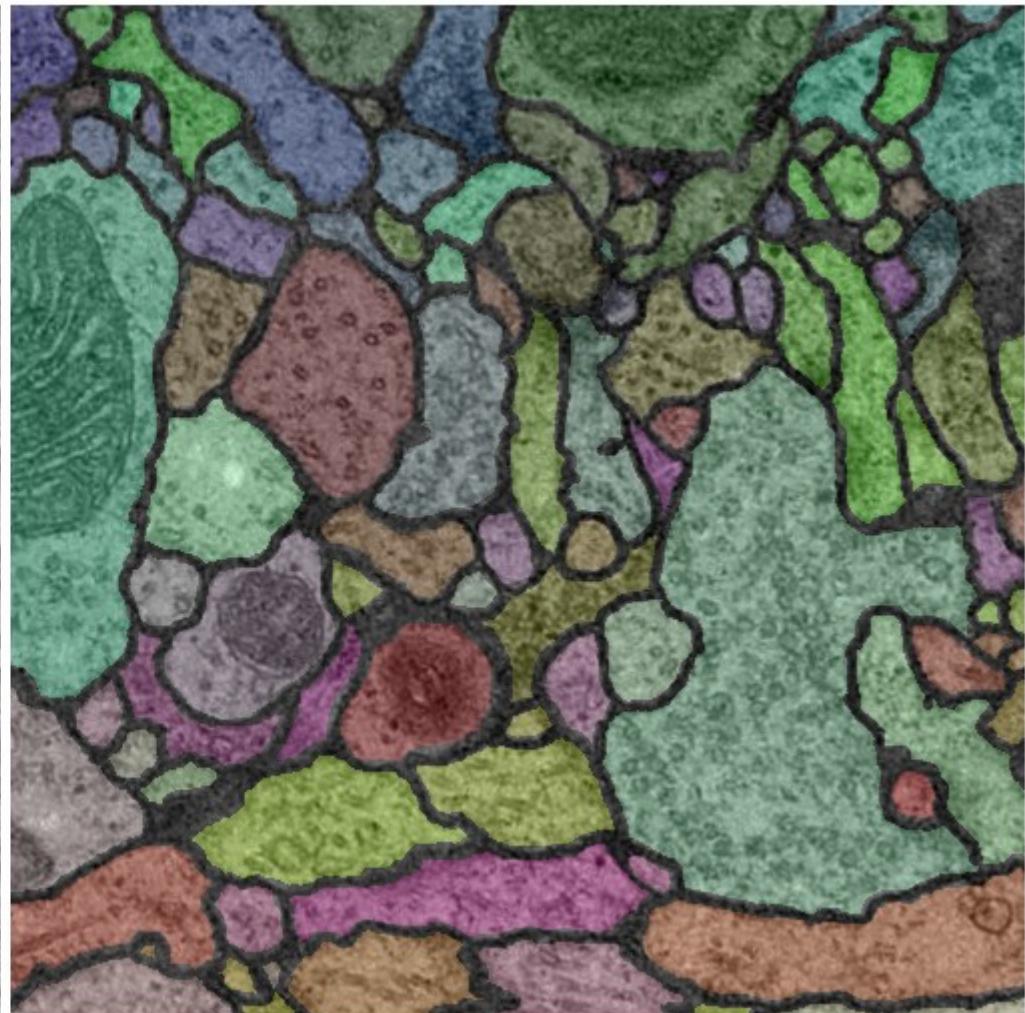
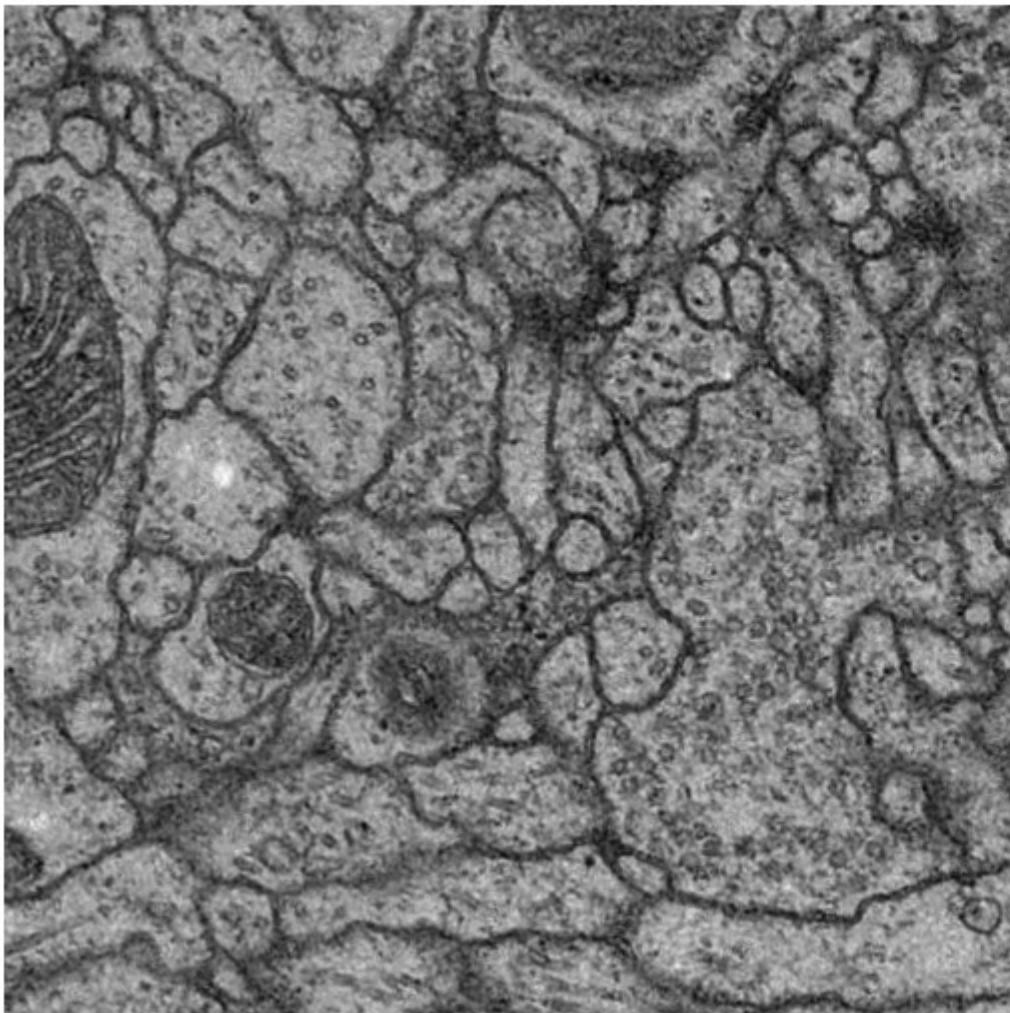
A screenshot from a video game showing a brown building with multiple windows and doors. The building is labeled with various objects: "building" (multiple times), "pole" (multiple times), "window" (multiple times), "door" (multiple times), and "sidewalk" (multiple times). The building is surrounded by a green lawn and a blue sky.

An abstract illustration composed of various colored shapes (green, blue, red, yellow) representing different elements of a landscape. Overlaid on these shapes are numerous white, sans-serif labels identifying specific items. The labels are repeated multiple times and include: 'sky' (top left), 'tree' (multiple instances throughout), 'plant' (multiple instances), 'building' (multiple instances), 'grass' (multiple instances), 'sign' (multiple instances), 'car' (multiple instances), 'road' (multiple instances), and 'sidewalk' (multiple instances). The overall effect is a dense, colorful collage where each label corresponds to a similar object or shape within the scene.

Farabet et al. "Learning hierarchical features for scene labeling" PAMI 2013  
Pinheiro et al. "Recurrent CNN for scene parsing" arxiv 2013

# CONV NETS: EXAMPLES

- Segmentation 3D volumetric images



Ciresan et al. "DNN segment neuronal membranes..." NIPS 2012

Turaga et al. "Maximin learning of image segmentation" NIPS 2009

# CONV NETS: EXAMPLES

- Action recognition from videos



Taylor et al. "Convolutional learning of spatio-temporal features" ECCV 2010

Karpathy et al. "Large-scale video classification with CNNs" CVPR 2014

Simonyan et al. "Two-stream CNNs for action recognition in videos" arXiv 2014

# CONV NETS: EXAMPLES

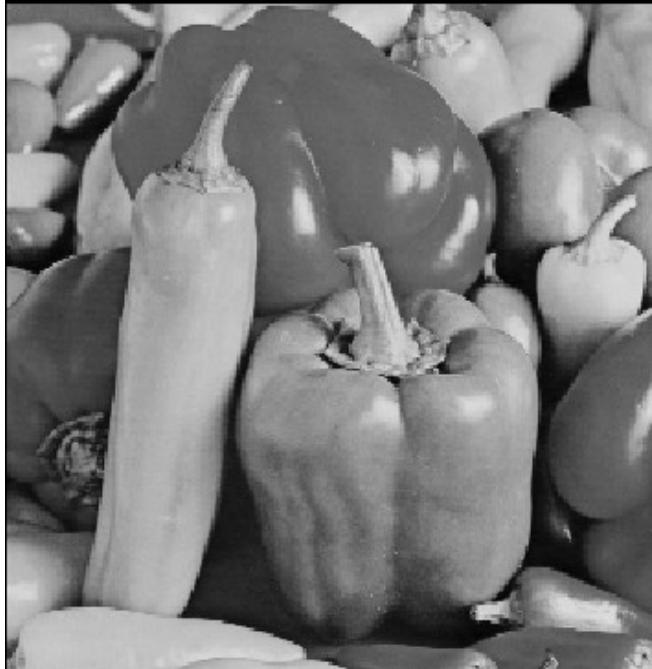
- Robotics



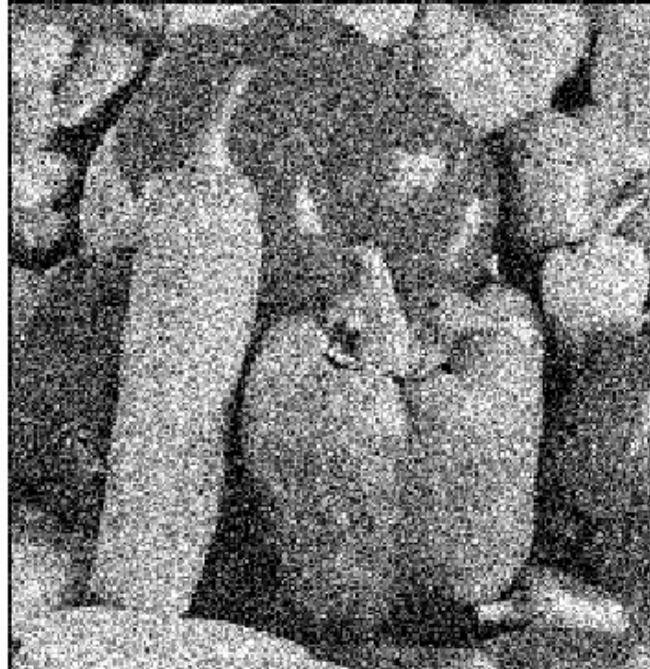
# CONV NETS: EXAMPLES

## - Denoising

original



noised

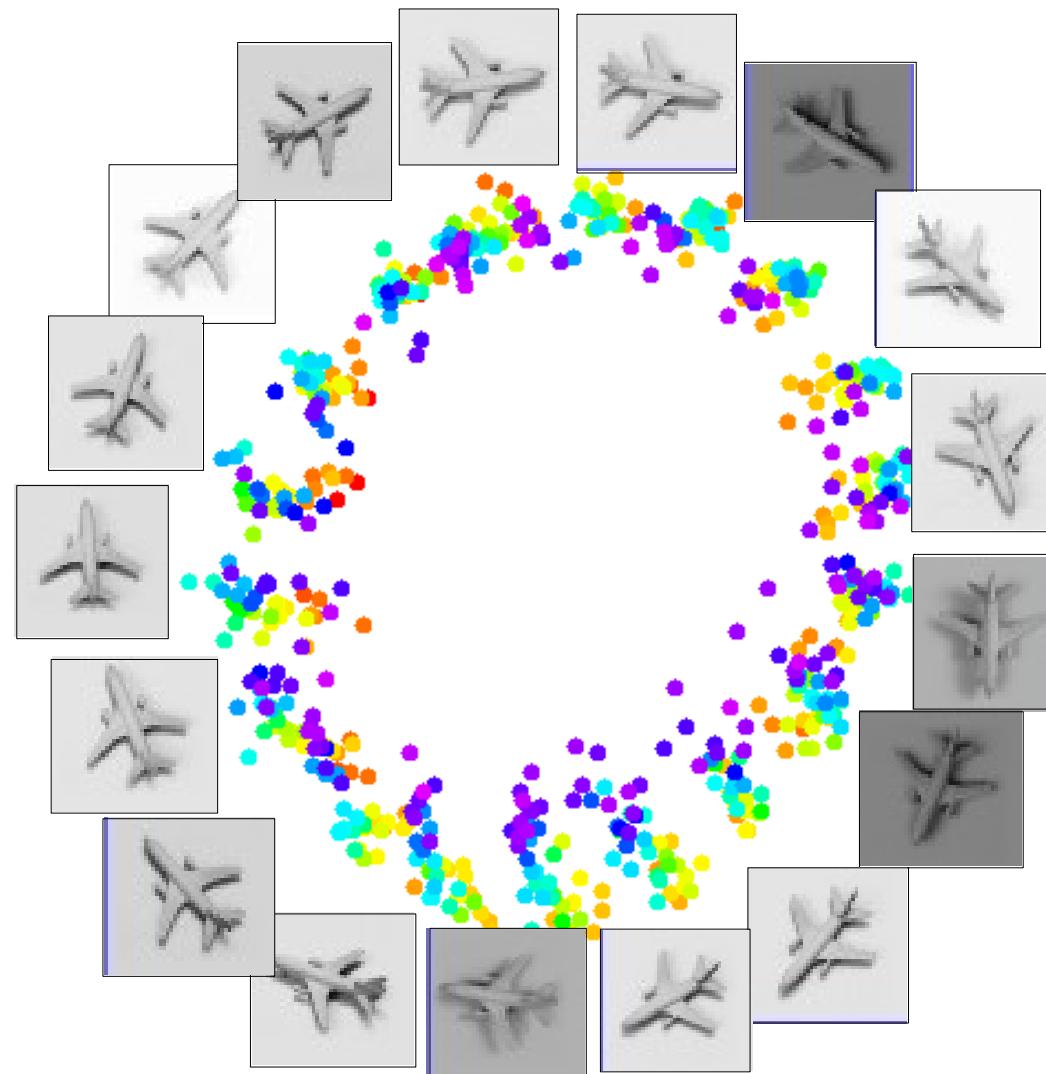


denoised



# CONV NETS: EXAMPLES

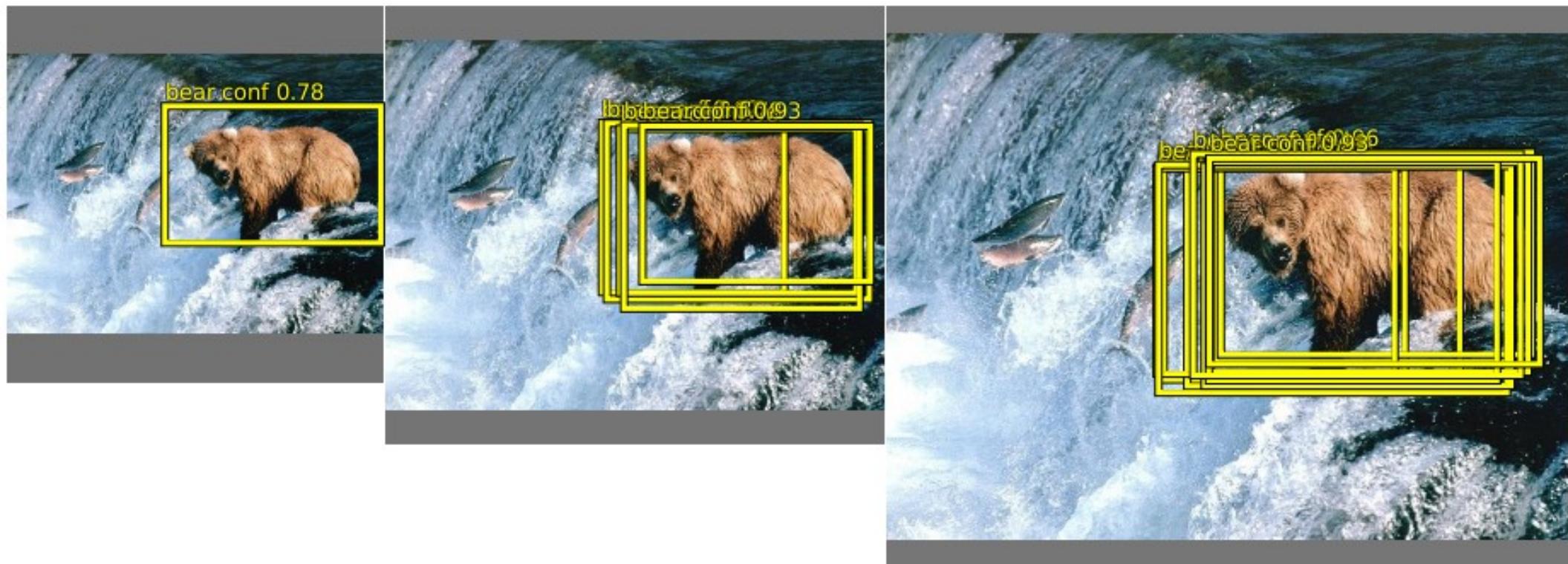
- Dimensionality reduction / learning embeddings



90

# CONV NETS: EXAMPLES

## - Object detection



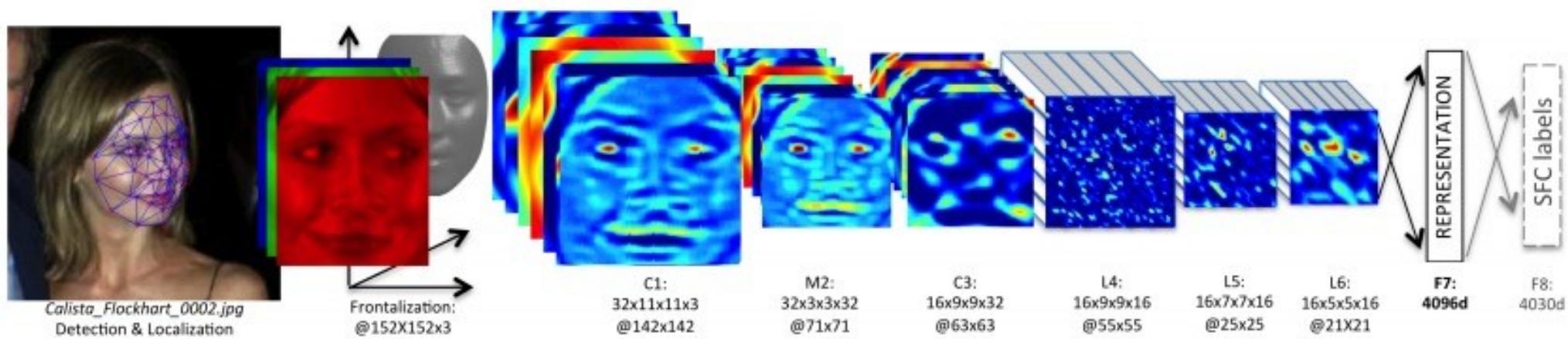
Sermanet et al. “OverFeat: Integrated recognition, localization, ...” arxiv 2013

Girshick et al. “Rich feature hierarchies for accurate object detection...” arxiv 2013 91

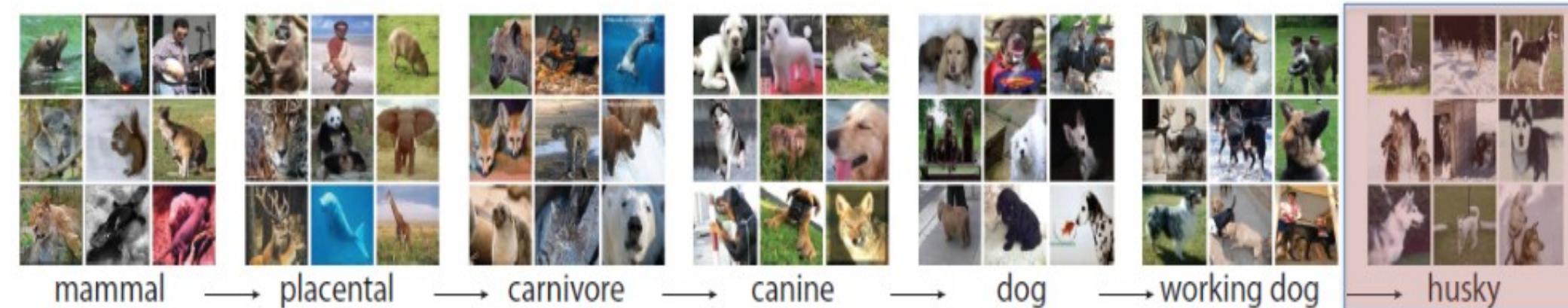
Szegedy et al. “DNN for object detection” NIPS 2013

# CONV NETS: EXAMPLES

## - Face Verification & Identification



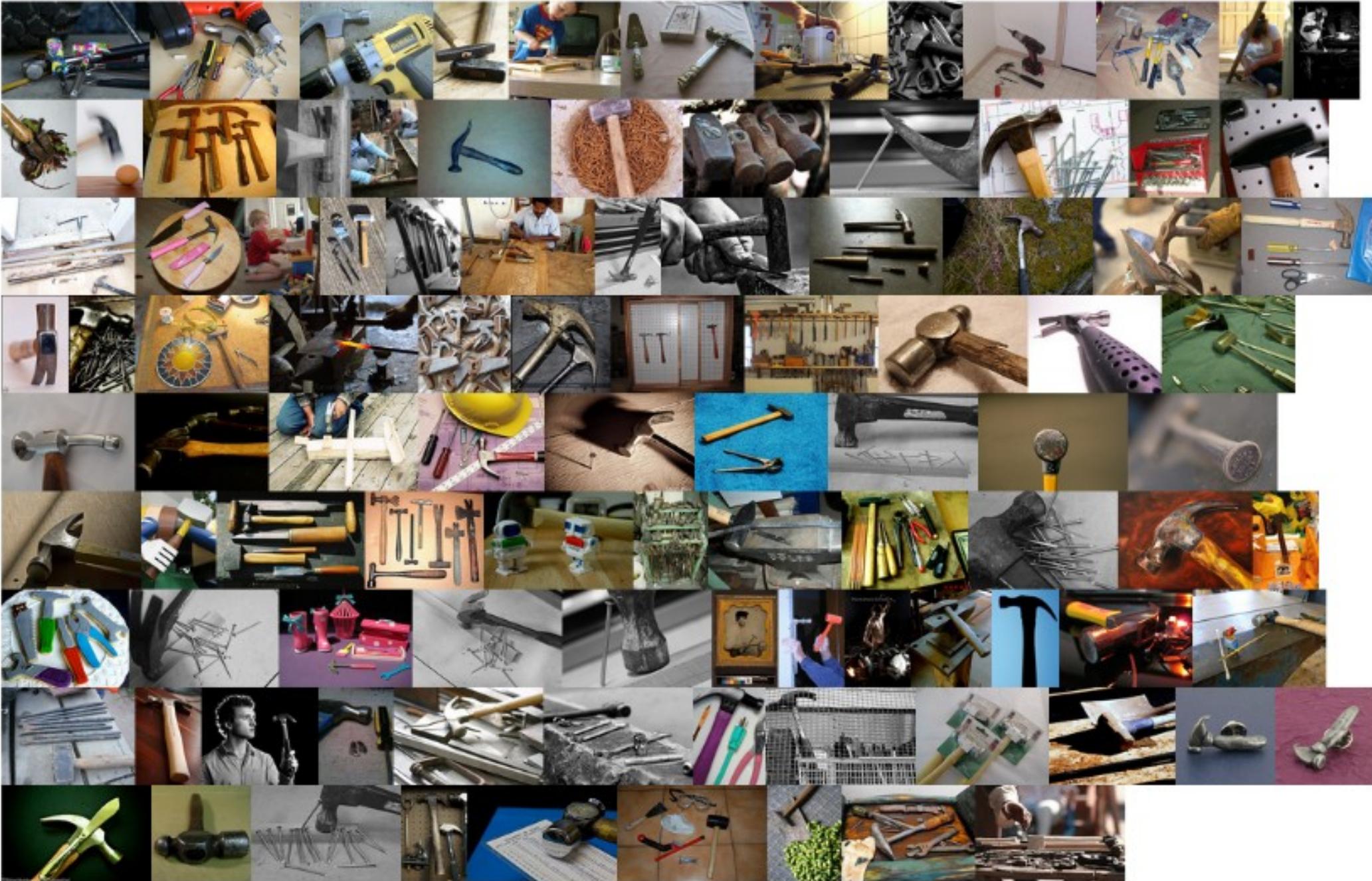
# Dataset: ImageNet 2012



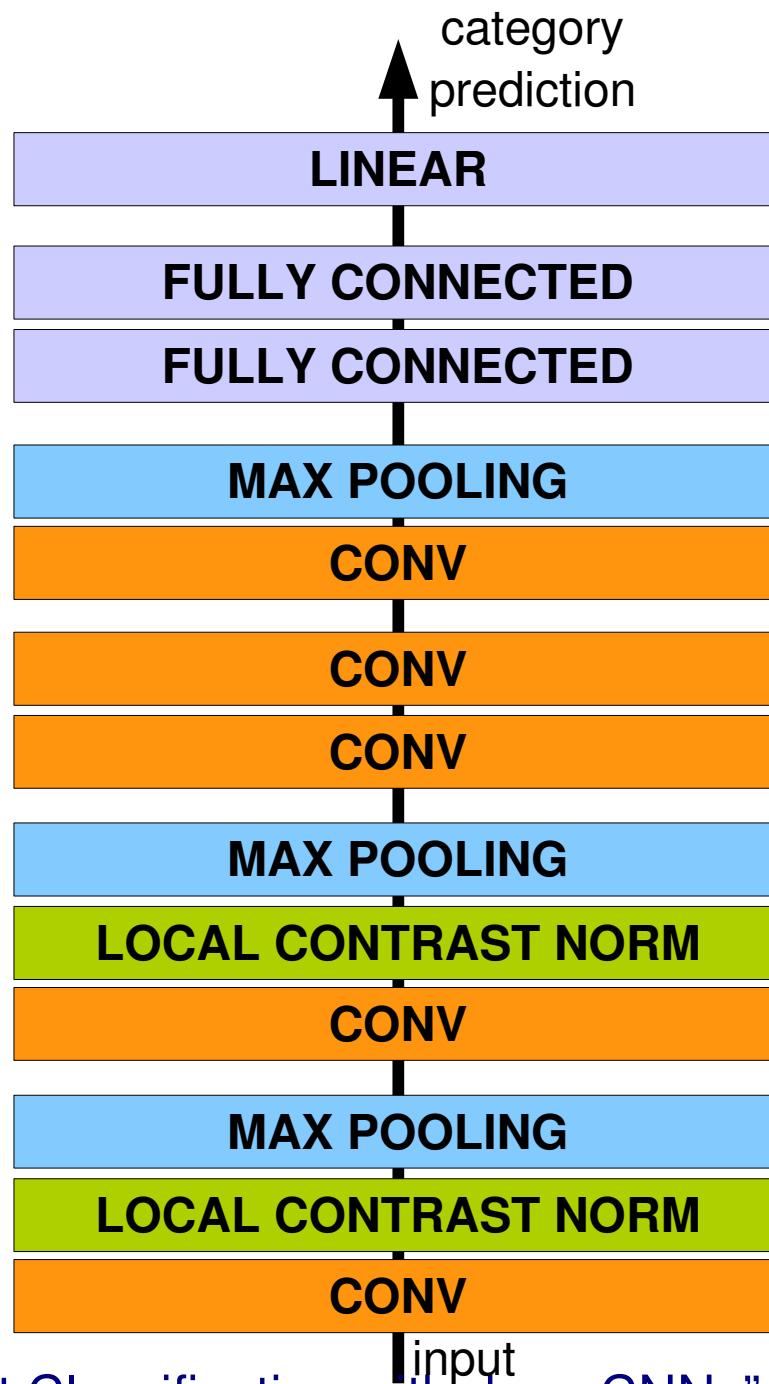
- S: (n) [Eskimo dog](#), [husky](#) (breed of heavy-coated Arctic sled dog)
  - [direct hypernym](#) / [inherited hypernym](#) / [sister term](#)
- S: (n) [working dog](#) (any of several breeds of usually large powerful dogs bred to work as draft animals and guard and guide dogs)
  - S: (n) [dog](#), [domestic dog](#), [Canis familiaris](#) (a member of the genus *Canis* (probably descended from the common wolf) that has been domesticated by man since prehistoric times; occurs in many breeds) "the dog barked all night"
  - S: (n) [canine](#), [canid](#) (any of various fissiped mammals with nonretractile claws and typically long muzzles)
  - S: (n) [carnivore](#) (a terrestrial or aquatic flesh-eating mammal) "terrestrial carnivores have four or five clawed digits on each limb"
  - S: (n) [placental](#), [placental mammal](#), [eutherian](#), [eutherian mammal](#) (mammals having a placenta; all mammals except monotremes and marsupials)
  - S: (n) [mammal](#), [mammalian](#) (any warm-blooded vertebrate having the skin more or less covered with hair; young are born alive except for the small subclass of monotremes and nourished with milk)
    - S: (n) [vertebrate](#), [craniate](#) (animals having a bony or cartilaginous skeleton with a segmented spinal column and a large brain enclosed in a skull or cranium)
    - S: (n) [chordate](#) (any animal of the phylum Chordata having a notochord or spinal column)
      - S: (n) [animal](#), [animate being](#), [beast](#), [brute](#), [creature](#), [fauna](#) (a living organism characterized by voluntary movement)
      - S: (n) [organism](#), [being](#) (a living thing that has (or can develop) the ability to act or function independently)
      - S: (n) [living thing](#), [animate thing](#) (a living (or once living) entity)
    - S: (n) [whole](#), [unit](#) (an assemblage of parts that is regarded as a single entity) "how big is that part compared to the whole?"; "the team is a unit"
    - S: (n) [object](#), [physical object](#) (a tangible and visible entity; an entity that can cast a shadow) "it was full of rackets, balls and other objects"
    - S: (n) [physical entity](#) (an entity that has physical existence)
    - S: (n) [entity](#) (that which is perceived or known or inferred to have its own distinct existence (living or nonliving))

# ImageNet

Examples of hammer:

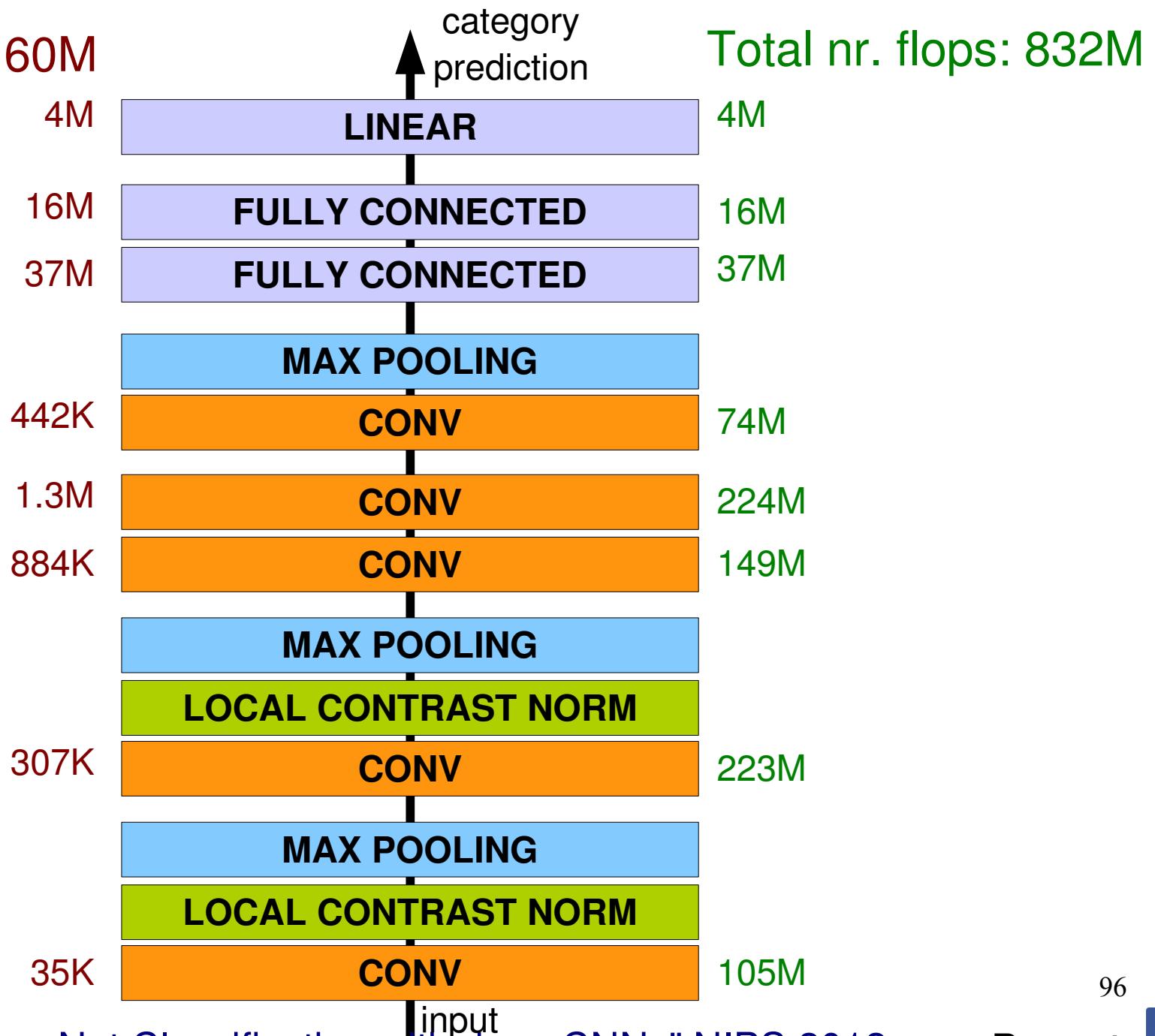


# Architecture for Classification



# Architecture for Classification

Total nr. params: 60M



# Optimization

## SGD with momentum:

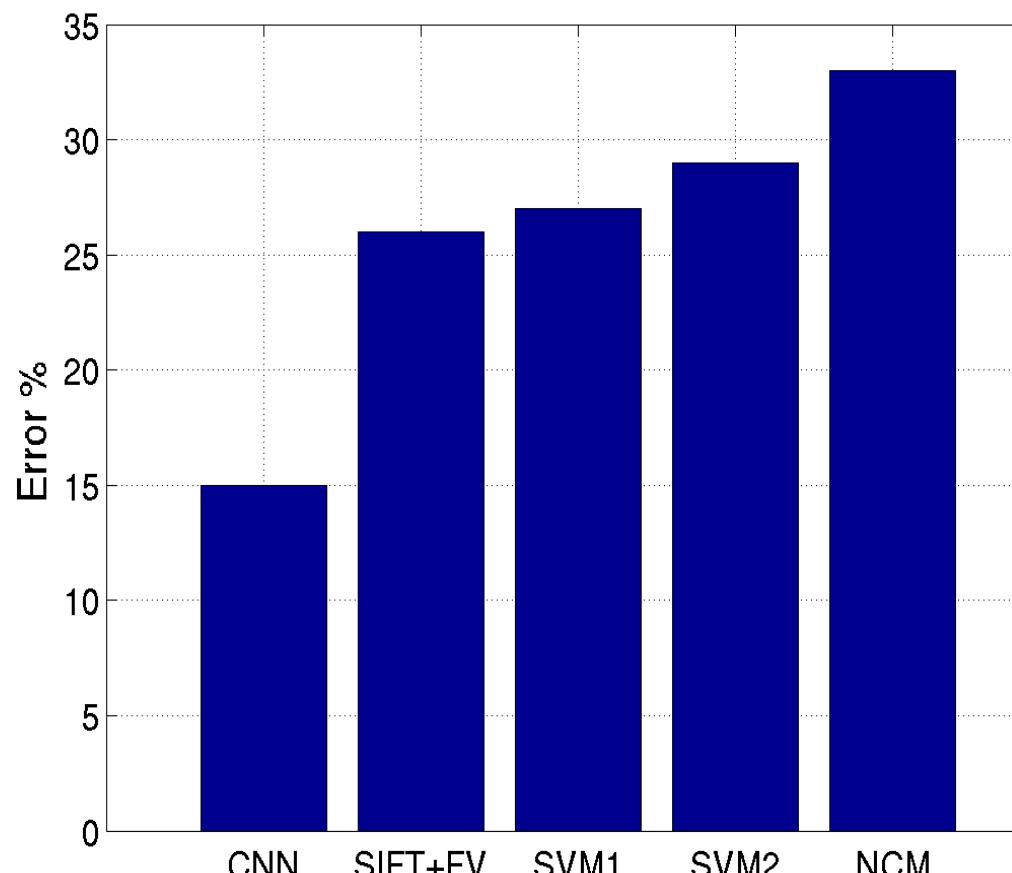
- Learning rate = 0.01
- Momentum = 0.9

## Improving generalization by:

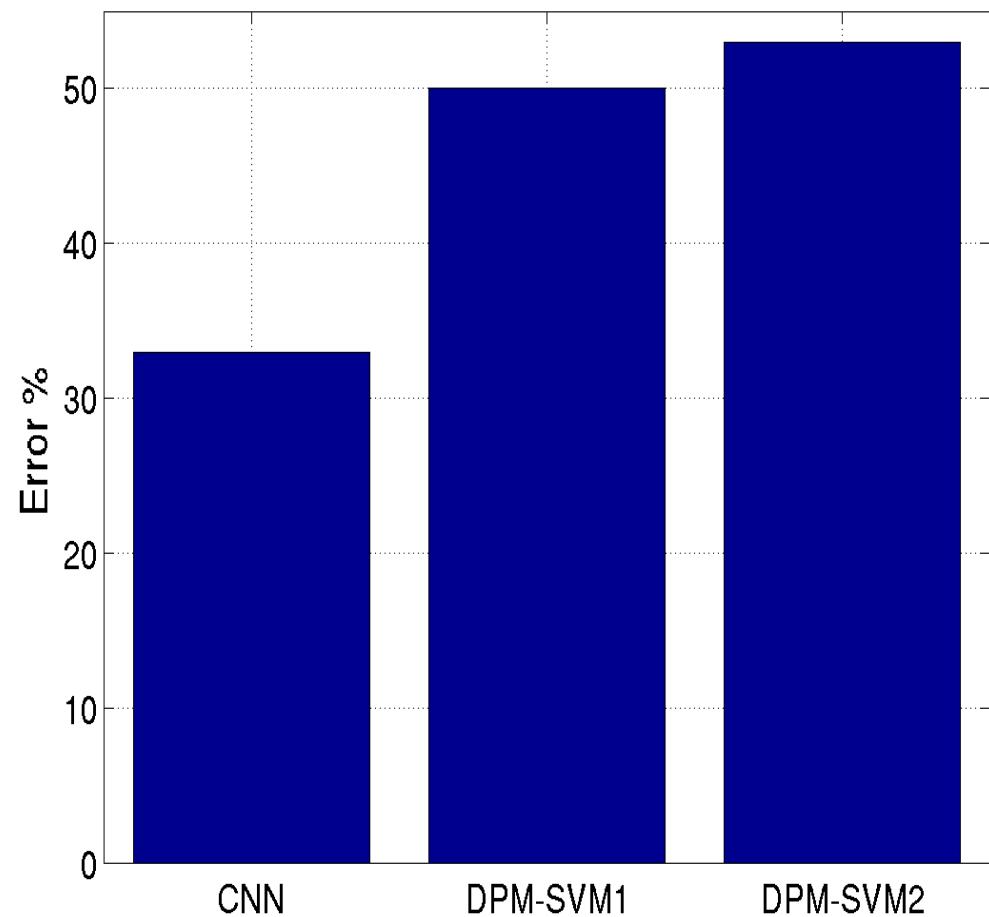
- Weight sharing (convolution)
- Input distortions
- Dropout = 0.5
- Weight decay = 0.0005

# Results: ILSVRC 2012

TASK 1 - CLASSIFICATION



TASK 2 - DETECTION





mite

mite  
black widow  
cockroach  
tick  
starfish

container ship

container ship  
lifeboat  
amphibian  
fireboat  
drilling platform

motor scooter

motor scooter  
go-kart  
moped  
bumper car  
golfcart

leopard

leopard  
jaguar  
cheetah  
snow leopard  
Egyptian cat



grille

convertible  
grille  
pickup  
beach wagon  
fire engine

mushroom

agaric  
mushroom  
jelly fungus  
gill fungus  
dead-man's-fingers

cherry

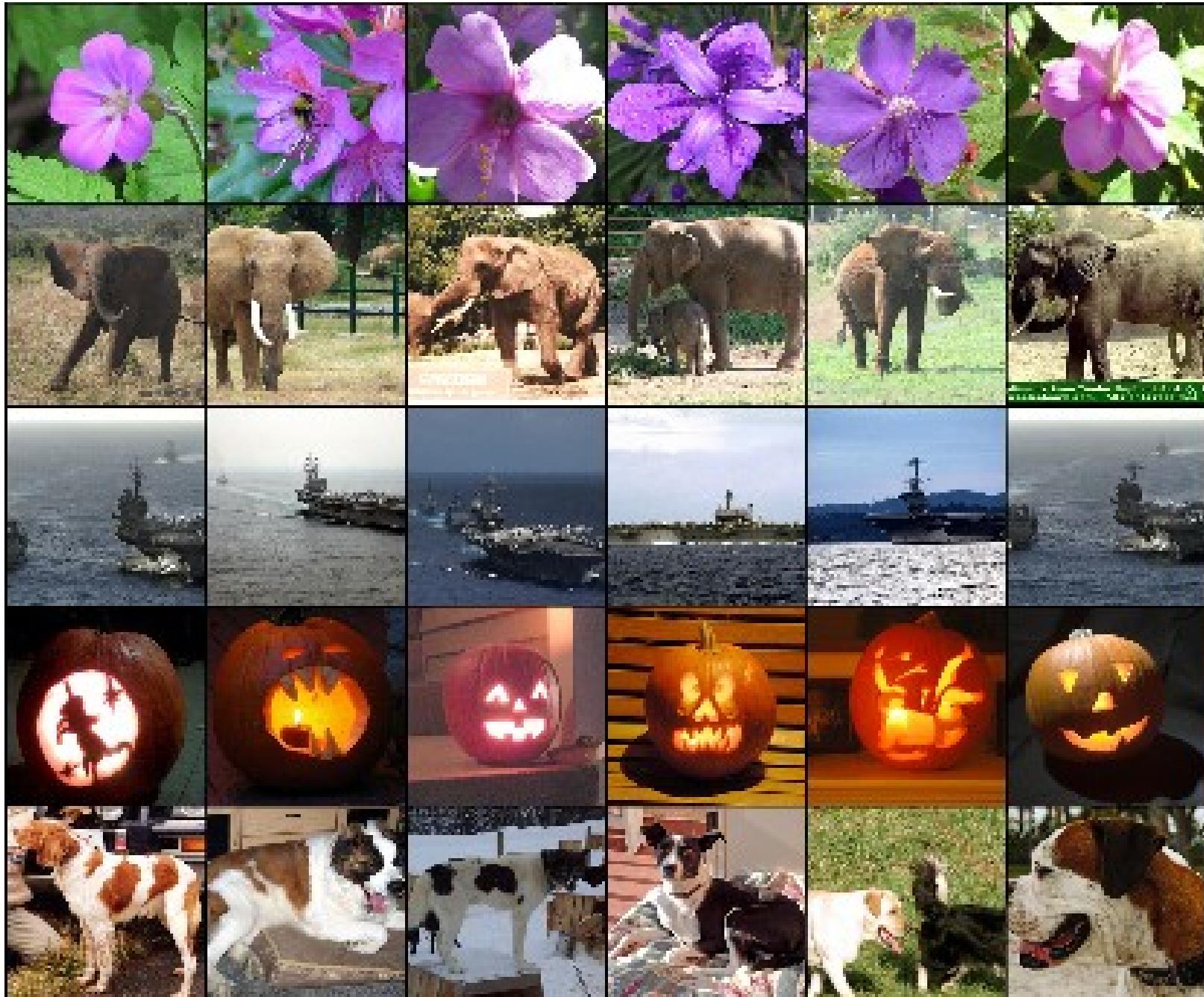
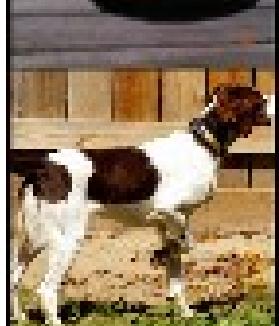
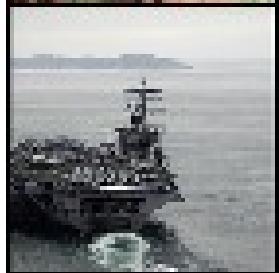
dalmatian  
grape  
elderberry  
ffordshire bullterrier  
currant

Madagascar cat

squirrel monkey  
spider monkey  
titi  
indri  
howler monkey

# TEST IMAGE

# RETRIEVED IMAGES



# Demo of classifier by Matt Zeiler & Rob Fergus:

<http://horatio.cs.nyu.edu/>

The screenshot shows a web browser window with the following details:

- Header:** File Edit View History Bookmarks Tools Help
- Title Bar:** Image Classifier Demo
- Address Bar:** horatio.cs.nyu.edu/index.html
- Toolbar:** Back, Forward, Stop, Refresh, Home, etc.
- Page Content:**
  - Section Header:** Image Classifier Demo
  - Text:** Upload your images to have them classified by a machine! Upload multiple images using the button below or dropping them on this page. The predicted objects will be refreshed automatically. Images are resized such that the smallest dimension becomes 256, then the center 256x256 crop is used. More about the demo can be found [here](#).
  - Buttons:** + Upload Images, Remove All, Show help tips
  - Agreement:** I agree to the Terms of Use (checkbox checked)
  - Predicted objects:**
    1. ✓ ✗ ⓘ Lion, King Of Beasts, Panthera Leo (0.34)
    2. ✓ ✗ ⓘ Hartebeest (0.19)
    3. ✓ ✗ ⓘ Hyena, Hyaena (0.16)
    4. ✓ ✗ ⓘ Arabian Camel, Dromedary, Camelus Dromedarius (0.06)
    5. ✓ ✗ ⓘ Dingo, Warrigal, Warragal, Canis Dingo (0.04)

A cursor is hovering over the fifth item in the list.

**Other objects:** (empty input field)

**NYU Logo:** NYU

# Demo of classifier by Yangqing Jia & Trevor Darrell:

***<http://decafberkeleyvision.org/>***

File Edit View History Bookmarks Tools Help

Yangqing Jia Decaf Demo DeCAF: A Deep Convolution... +

← decaf.berkeleyvision.org/classify\_url?imageurl=http%3A%2F%2Fwww.careerealism.com%2Fhome%2Fjtodonnell%2Fcareerealism. star ↗ g decaf berkeley

Most Visited Getting Started Latest Headlines Click Here!

## Decaf Image Classifier

New: get the [software](#) and [tech report](#) that we have released!

[\[About this demo\]](#) [\[Sign up for Updates!\]](#)

Provide an image and have it classified by decaf. [Click for a Quick Example](#)



Flat Prediction	Maximize Infogain
lion	0.87769
dhole	0.01987
red fox	0.01606
coyote	0.01503
red wolf	0.01321

CNN took 0.462 seconds.

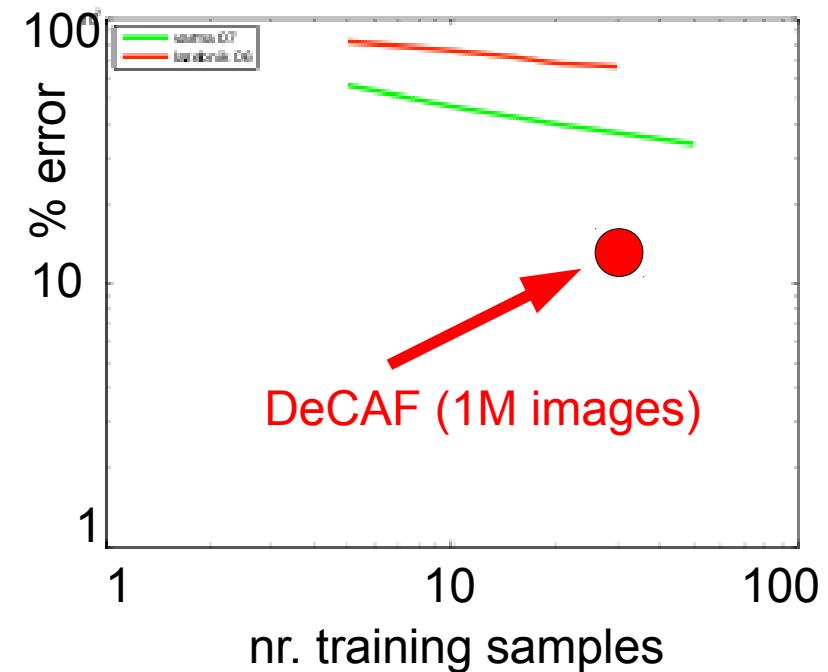
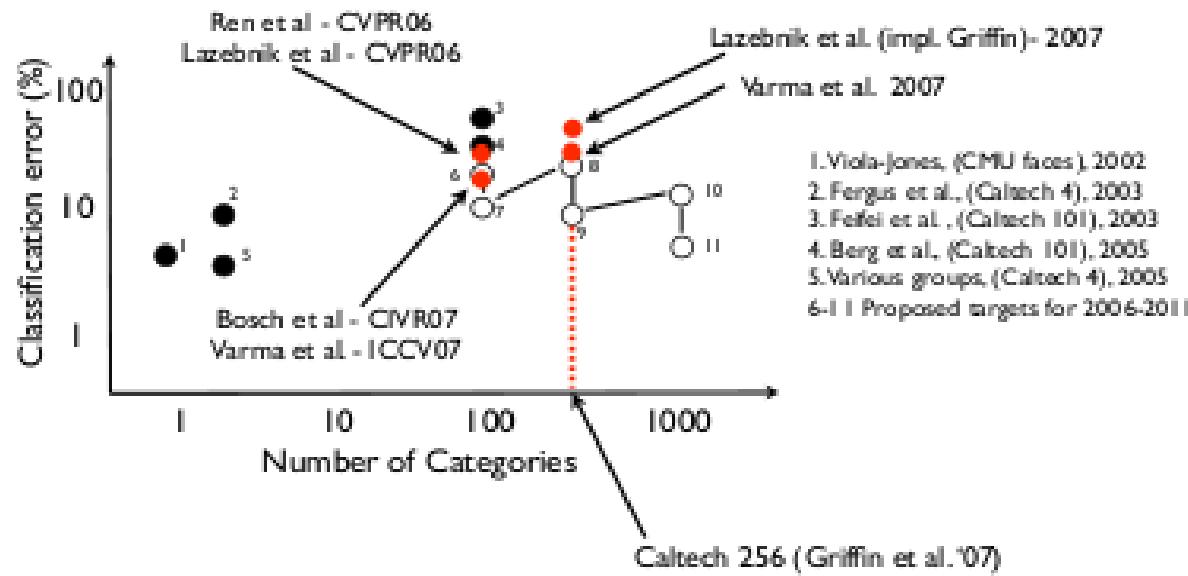


Figure 3: How well are we doing? (Left) Classification performance has seen steady improvement in the last few years, both in the number of categories on which algorithms are tested and in classification error rates. (Right) Performance of the best 2006 [Lazebnik et al., 2006] and the best 2007 algorithm [Varma, 2007] are compared here (classification error rates vs number of training examples). One may notice the significant year-on-year progress (see also left panel). Extrapolation enthusiasts may calculate that  $10^8$  training examples would be sufficient to achieve 1% error rates with current algorithms. Furthermore, if the pace of year-on-year progress is constant on this log scale chart, 1% error rates with 30 training examples will be achieved in 8-10 years.

Excerpt from Perona Visual Recognition 2007

103

Donahue, Jia et al. DeCAF arXiv 1310.1531 2013

# Outline

- Supervised Neural Networks
- Convolutional Neural Networks
- Examples
- Tips

# Choosing The Architecture

- Task dependent
- Cross-validation
- [Convolution → LCN → pooling]\* + fully connected layer
- The more data: the more layers and the more kernels
  - Look at the number of parameters at each layer
  - Look at the number of flops at each layer
- Computational resources
- Be creative :)

# How To Optimize

- SGD (with momentum) usually works very well
- Pick learning rate by running on a subset of the data  
Bottou “Stochastic Gradient Tricks” Neural Networks 2012
  - Start with large learning rate and divide by 2 until loss does not diverge
  - Decay learning rate by a factor of ~1000 or more by the end of training
- Use  non-linearity
- Initialize parameters so that each feature across layers has similar variance. Avoid units in saturation.

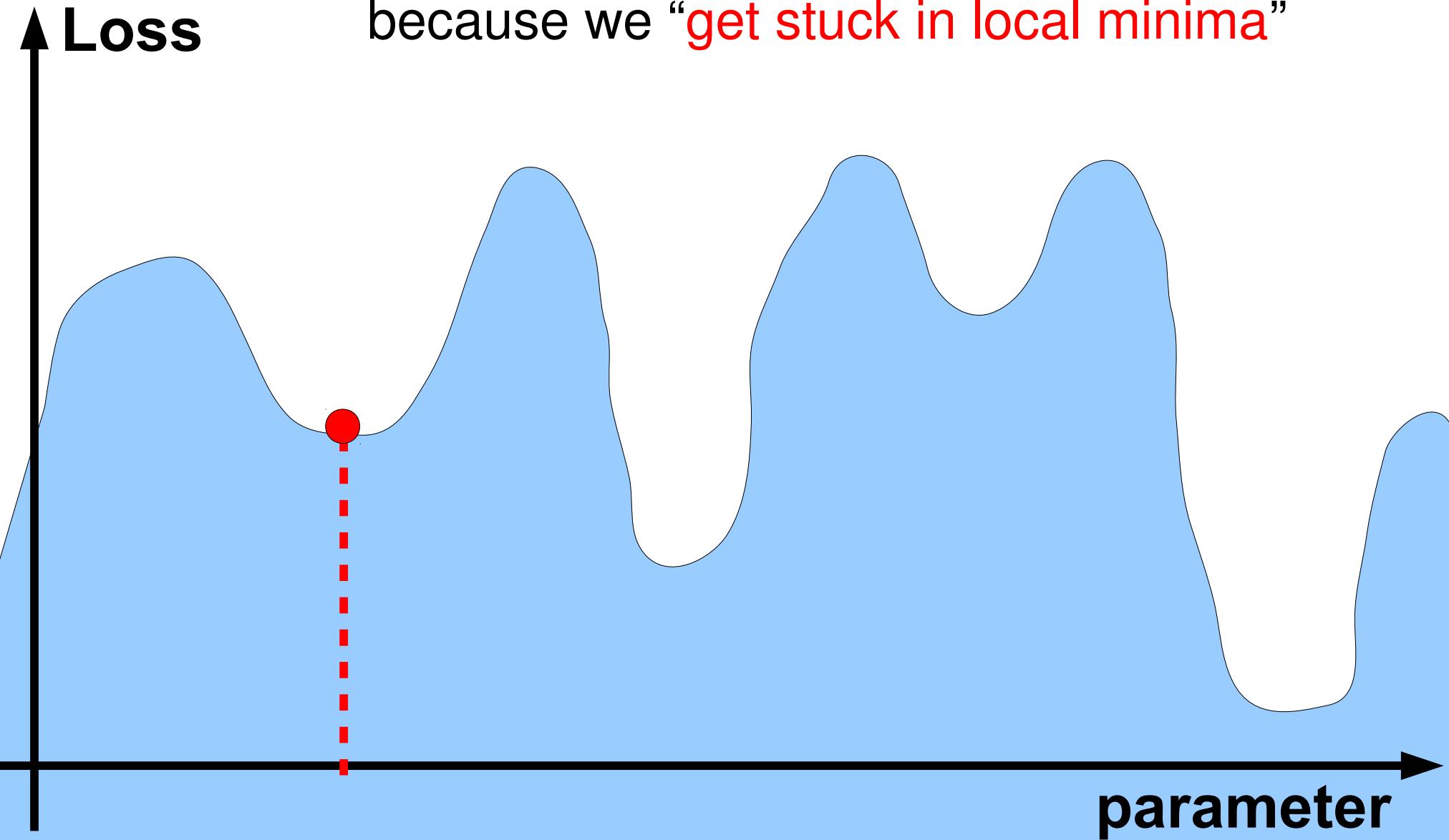
# Improving Generalization

- Weight sharing (greatly reduce the number of parameters)
- Data augmentation (e.g., jittering, noise injection, etc.)
- Dropout

Hinton et al. “Improving Nns by preventing co-adaptation of feature detectors”  
arxiv 2012
- Weight decay (L2, L1)
- Sparsity in the hidden units
- Multi-task (unsupervised learning)

# ConvNets: till 2012

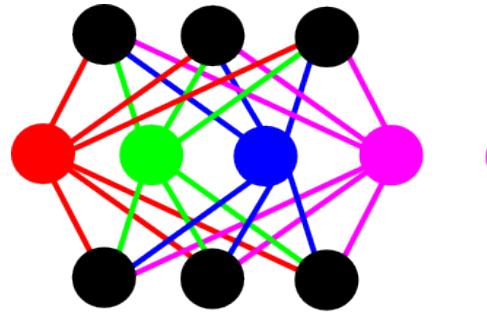
Common wisdom: training does not work because we “get stuck in local minima”



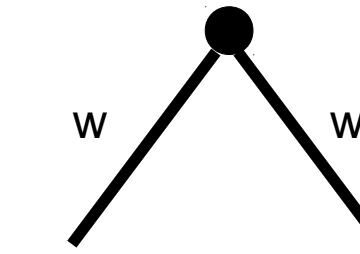
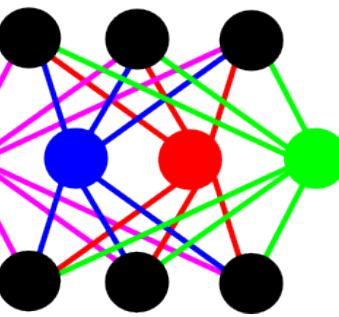
# ConvNets: today

Local minima are all similar, there are long plateaus,  
it can take long time to break symmetries.

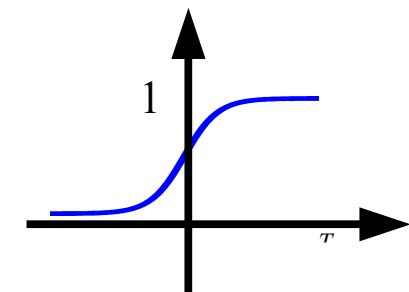
**Loss**



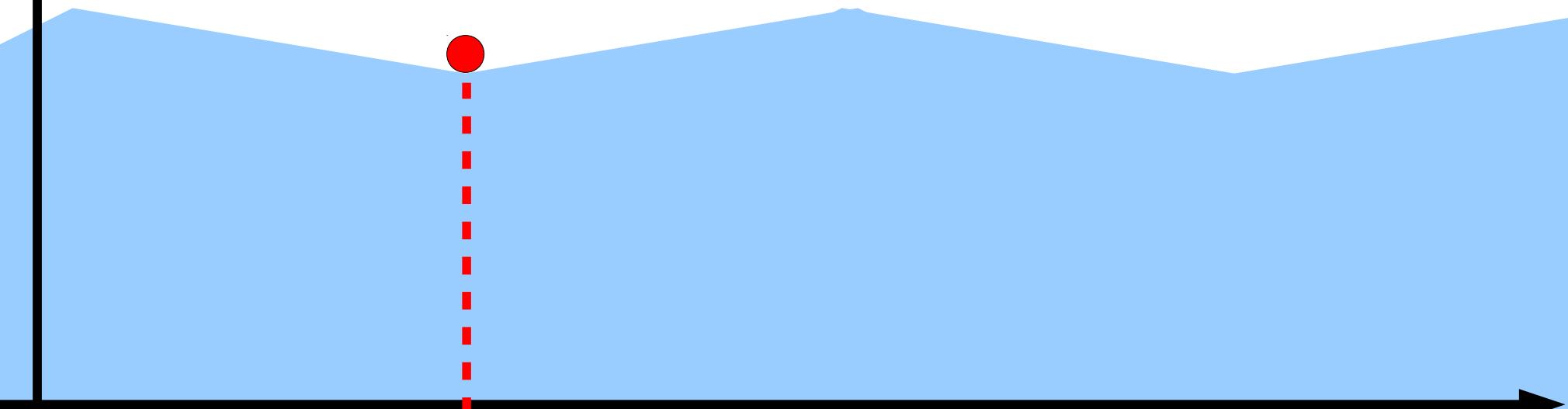
input/output invariant to permutations



breaking ties  
between parameters



Saturating units



# Neural Net Optimization is...

Like walking on a ridge between valleys

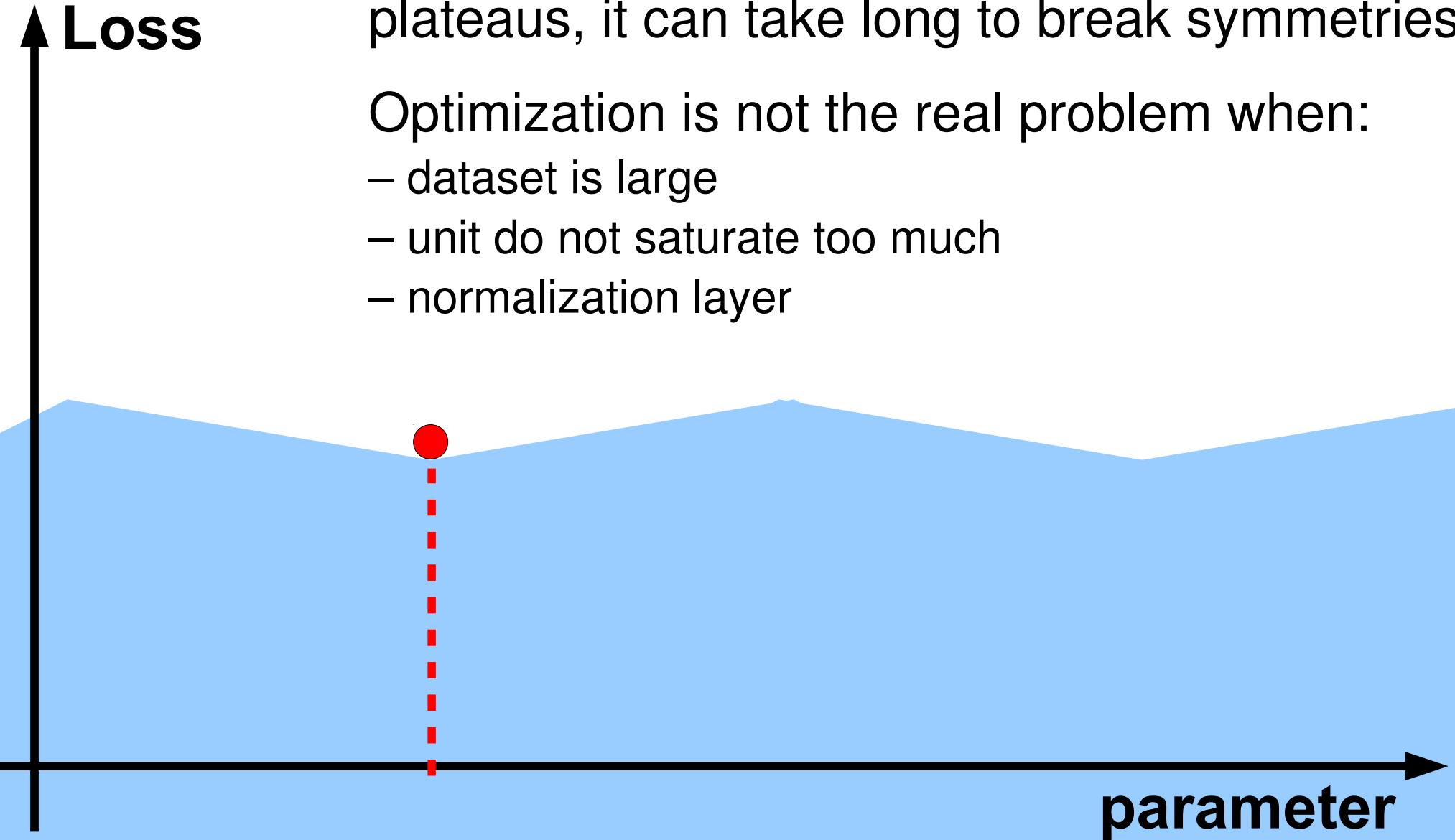


# ConvNets: today

Local minima are all similar, there are long plateaus, it can take long to break symmetries.

Optimization is not the real problem when:

- dataset is large
- unit do not saturate too much
- normalization layer



# ConvNets: today

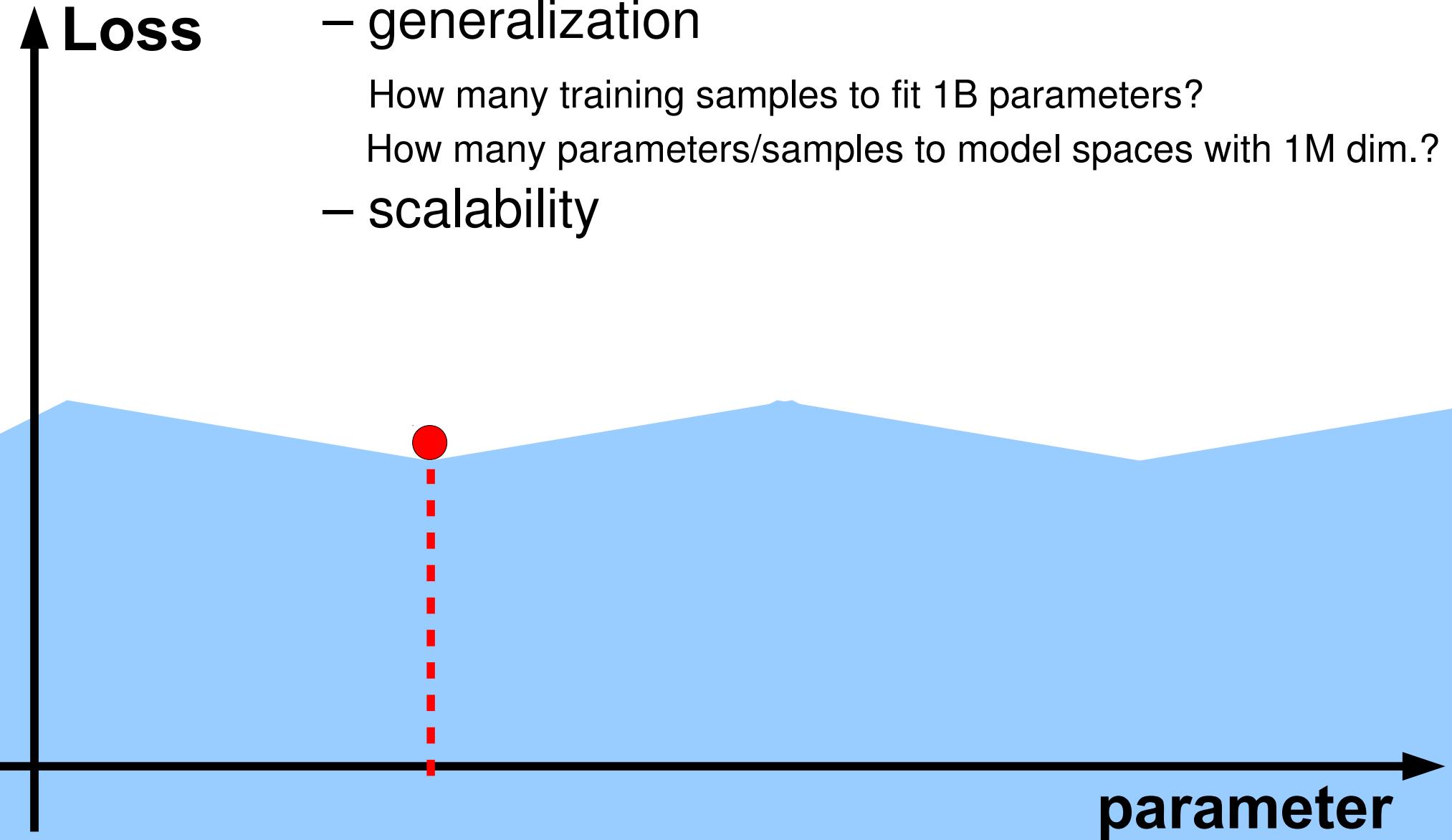
Today's belief is that the challenge is about:

- generalization

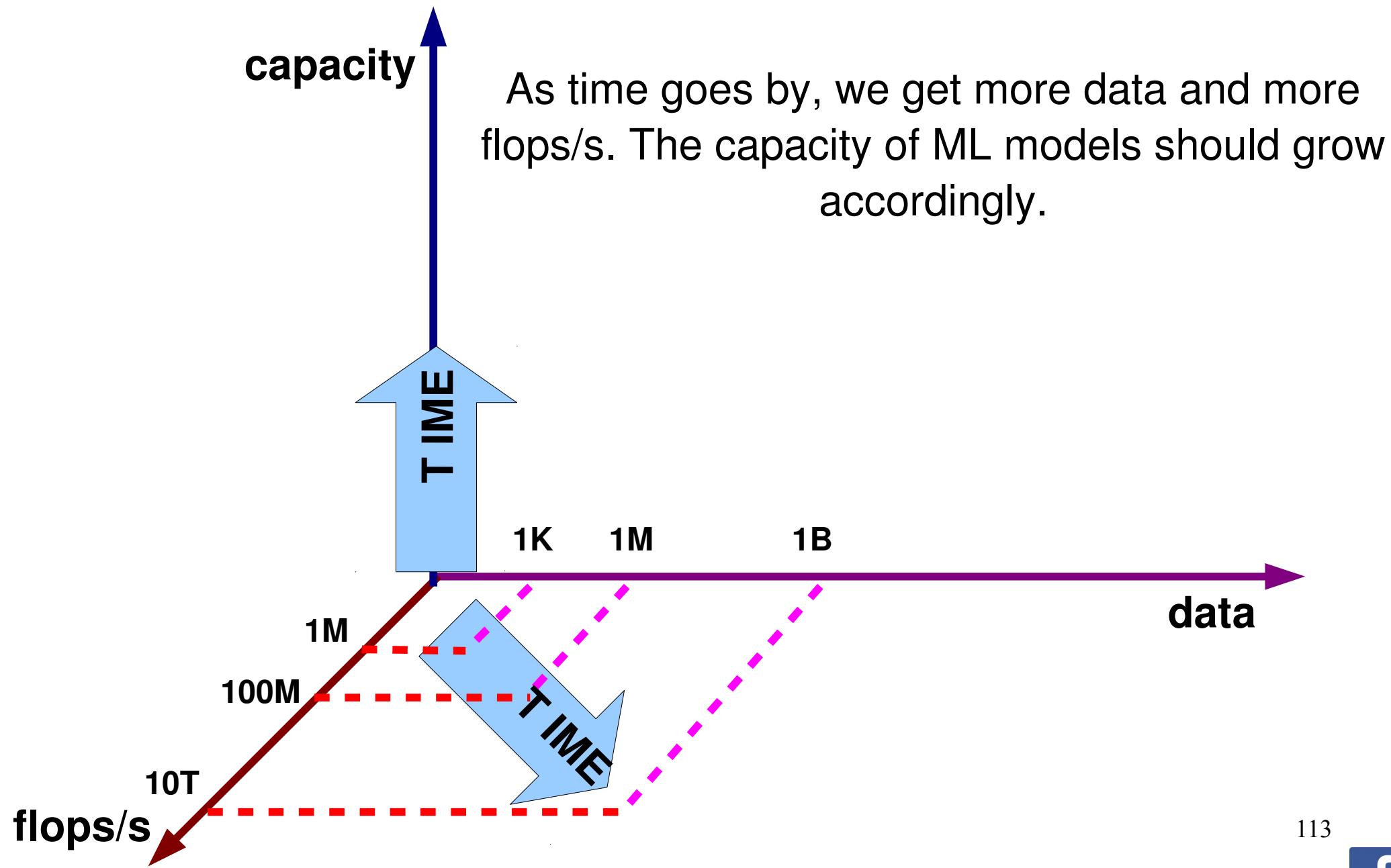
- How many training samples to fit 1B parameters?

- How many parameters/samples to model spaces with 1M dim.?

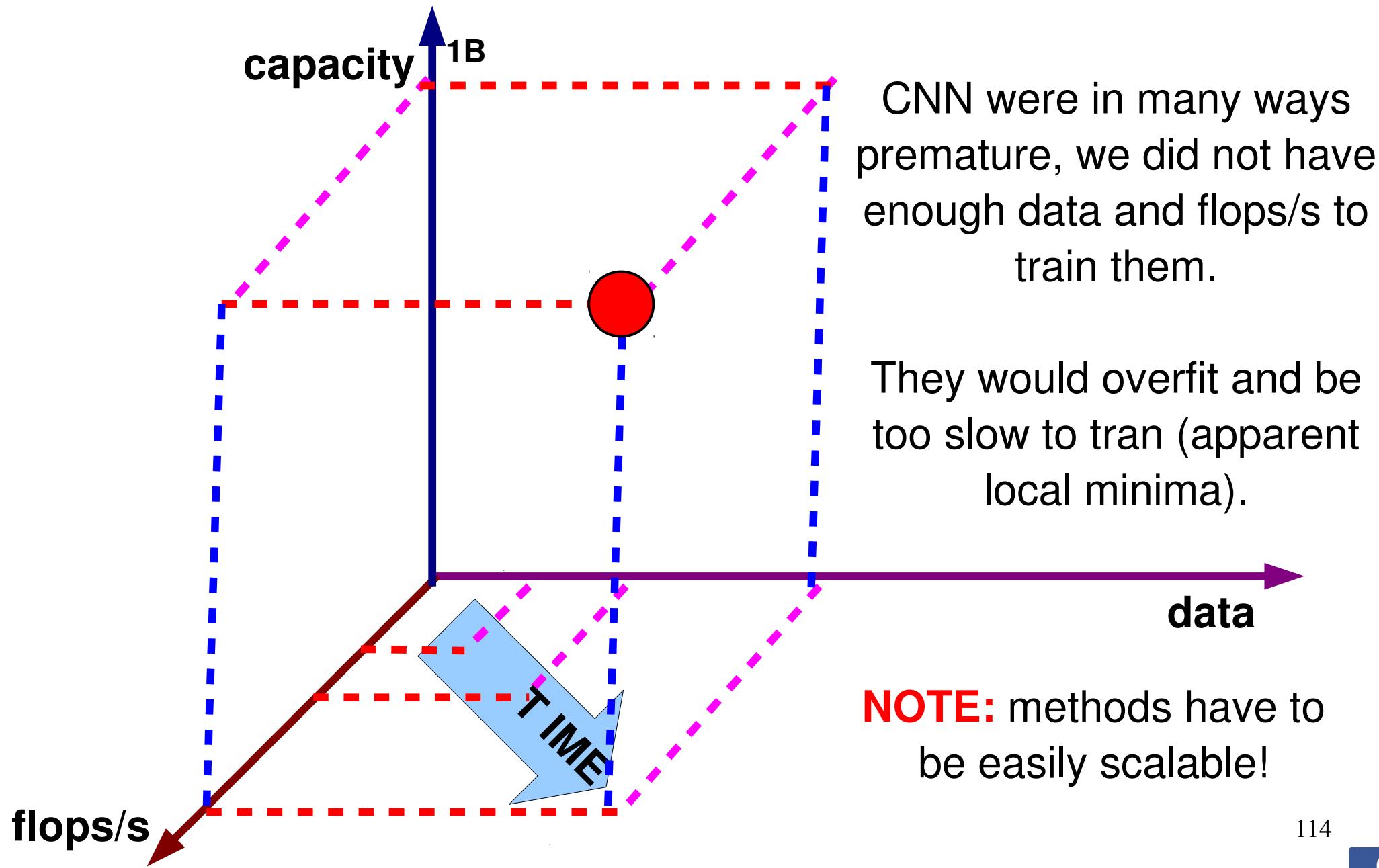
- scalability



# ConvNets: Why so successful today?

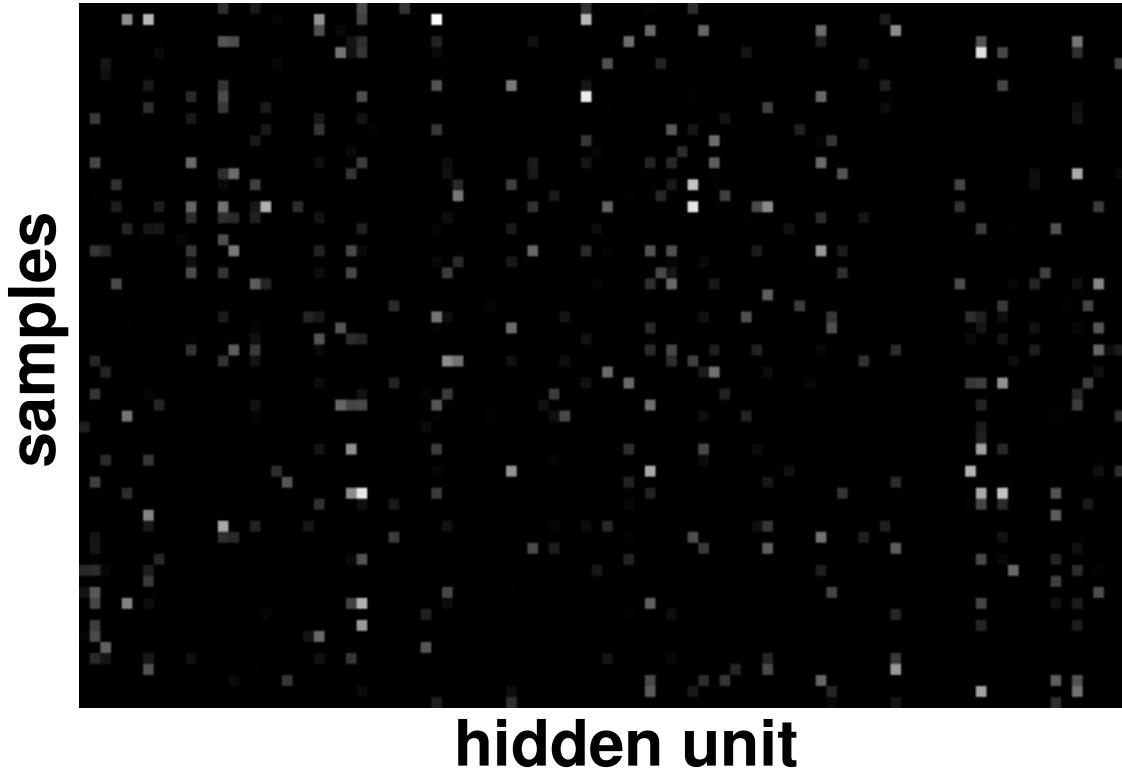


# ConvNets: Why so successful today?



# Good To Know

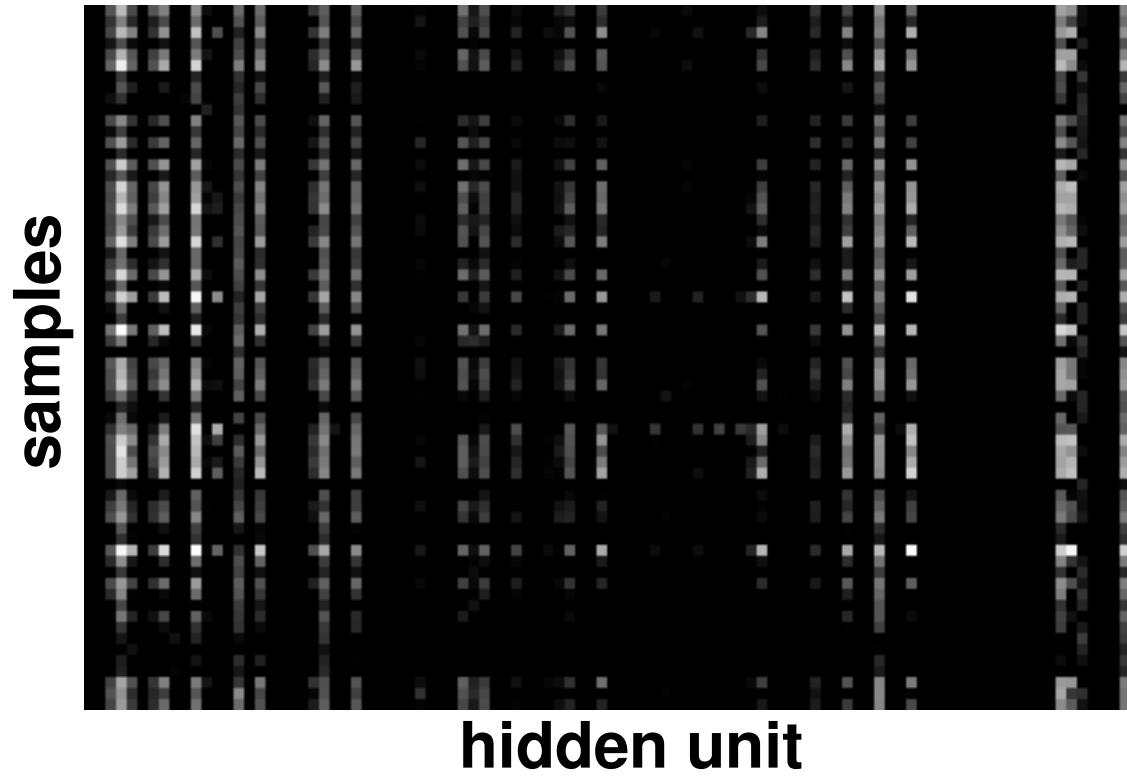
- Check gradients numerically by finite differences
- Visualize features (feature maps need to be uncorrelated) and have high variance.



**Good training:** hidden units are sparse across samples  
and across features.

# Good To Know

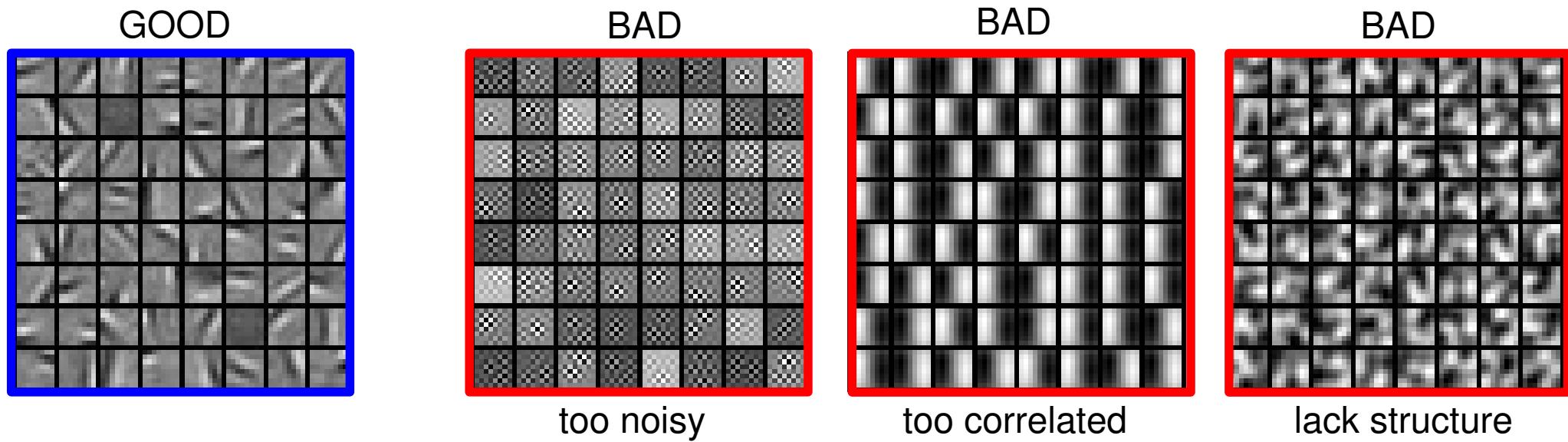
- Check gradients numerically by finite differences
- Visualize features (feature maps need to be uncorrelated) and have high variance.



**Bad training:** many hidden units ignore the input and/or exhibit strong correlations.

# Good To Know

- Check gradients numerically by finite differences
- Visualize features (feature maps need to be uncorrelated) and have high variance.
- Visualize parameters



**Good training:** learned filters exhibit structure and are uncorrelated.

117

Zeiler, Fergus "Visualizing and understanding CNNs" arXiv 2013

Simonyan, Vedaldi, Zisserman "Deep inside CNNs: visualizing image classification models.." ICLR 2014

# Good To Know

- Check gradients numerically by finite differences
- Visualize features (feature maps need to be uncorrelated) and have high variance.
- Visualize parameters
- Measure error on both training and validation set.
- Test on a small subset of the data and check the error → 0.

# What If It Does Not Work?

- Training diverges:
  - Learning rate may be too large → decrease learning rate
  - BPROP is buggy → numerical gradient checking
- Parameters collapse / loss is minimized but accuracy is low
  - Check loss function:
    - Is it appropriate for the task you want to solve?
    - Does it have degenerate solutions? Check “pull-up” term.
- Network is underperforming
  - Compute flops and nr. params. → if too small, make net larger
  - Visualize hidden units/params → fix optimization
- Network is too slow
  - Compute flops and nr. params. → GPU,distrib. framework, make net smaller

# Summary

- Supervised learning: today it is the most successful set up.  
ConvNets are used for a great variety of tasks.
- Optimization
  - Don't we get stuck in local minima? No, they are all the same!
  - In large scale applications, local minima are even less of an issue.
- Scaling
  - GPUs
  - Distributed framework (Google)
  - Better optimization techniques
- Generalization on small datasets (curse of dimensionality):
  - data augmentation
  - weight decay
  - dropout
  - unsupervised learning
  - multi-task learning

# SOFTWARE

## Torch7: learning library that supports neural net training

---

<http://www.torch.ch>

<http://code.cogbits.com/wiki/doku.php> (tutorial with demos by C. Farabet)

<https://github.com/sermanet/OverFeat>

## Python-based learning library (U. Montreal)

---

- <http://deeplearning.net/software/theano/> (does automatic differentiation)

## Efficient CUDA kernels for ConvNets (Krizhevsky)

---

– [code.google.com/p/cuda-convnet](http://code.google.com/p/cuda-convnet)

## Caffe (Yangqing Jia)

---

– <http://caffe.berkeleyvision.org>

# REFERENCES

## Convolutional Nets

---

- LeCun, Bottou, Bengio and Haffner: Gradient-Based Learning Applied to Document Recognition, Proceedings of the IEEE, 86(11):2278-2324, November 1998
- Krizhevsky, Sutskever, Hinton “ImageNet Classification with deep convolutional neural networks” NIPS 2012
- Jarrett, Kavukcuoglu, Ranzato, LeCun: What is the Best Multi-Stage Architecture for Object Recognition?, Proc. International Conference on Computer Vision (ICCV'09), IEEE, 2009
- Kavukcuoglu, Sermanet, Boureau, Gregor, Mathieu, LeCun: Learning Convolutional Feature Hierachies for Visual Recognition, Advances in Neural Information Processing Systems (NIPS 2010), 23, 2010
- see [yann.lecun.com/exdb/publis](http://yann.lecun.com/exdb/publis) for references on many different kinds of convnets.
- see <http://www.cmap.polytechnique.fr/scattering/> for scattering networks (similar to convnets but with less learning and stronger mathematical foundations)
- see <http://www.idsia.ch/~juergen/> for other references to ConvNets and LSTMs.

# REFERENCES

## Applications of Convolutional Nets

---

- Farabet, Couprie, Najman, LeCun. Scene Parsing with Multiscale Feature Learning, Purity Trees, and Optimal Covers”, ICML 2012
- Pierre Sermanet, Koray Kavukcuoglu, Soumith Chintala and Yann LeCun: Pedestrian Detection with Unsupervised Multi-Stage Feature Learning, CVPR 2013
- D. Ciresan, A. Giusti, L. Gambardella, J. Schmidhuber. Deep Neural Networks Segment Neuronal Membranes in Electron Microscopy Images. NIPS 2012
- Raia Hadsell, Pierre Sermanet, Marco Scoffier, Ayse Erkan, Koray Kavackuoglu, Urs Muller and Yann LeCun. Learning Long-Range Vision for Autonomous Off-Road Driving, Journal of Field Robotics, 26(2):120-144, 2009
- Burger, Schuler, Harmeling. Image Denoising: Can Plain Neural Networks Compete with BM3D?, CVPR 2012
- Hadsell, Chopra, LeCun. Dimensionality reduction by learning an invariant mapping, CVPR 2006
- Bergstra et al. Making a science of model search: hyperparameter optimization in hundred of dimensions for vision architectures, ICML 2013

# REFERENCES

## Latest and Greatest Convolutional Nets

---

- Girshick, Donahue, Darrell, Malick. “Rich feature hierarchies for accurate object detection and semantic segmentation”, arXiv 2014
- Karpathy, Toderici, Shetty, Leung, Sukthankar, FeiFei “Large-scale video classification with convolutional neural networks, CVPR 2014
- Cadieu, Hong, Yamins, Pinto, Ardila, Solomon, Majaj, DiCarlo. “DNN rival in representation of primate IT cortex for core visual object recognition”. arXiv 2014
- Erhan, Szegedy, Toshev, Anguelov “Scalable object detection using DNN” CVPR 2014
- Dauphin, Pascanu, Gulcehre, Cho, Ganguli, Bengio “Identifying and attacking the saddle point problem in high-dimensional non-convex optimization” arXiv 2014
- Razavian, Azizpour, Sullivan, Carlsson “CNN features off-the-shelf: and astounding baseline for recognition” arXiv 2014

# REFERENCES

## Deep Learning in general

---

- deep learning tutorial slides at ICML 2013
- Yoshua Bengio, Learning Deep Architectures for AI, Foundations and Trends in Machine Learning, 2(1), pp.1-127, 2009.
- LeCun, Chopra, Hadsell, Ranzato, Huang: A Tutorial on Energy-Based Learning, in Bakir, G. and Hofman, T. and Schölkopf, B. and Smola, A. and Taskar, B. (Eds), Predicting Structured Data, MIT Press, 2006

# THANK YOU