

# Artificial Intelligence II

Part 2: Lecture 7  
Yalda Mohsenzadeh

# Convolutional Neural Networks

# Today

- How to use networks for images
- Why “C”-NNs
- Standard building blocks of CNNs
- Some important networks & their tricks
- Some debugging tools

# Image classification

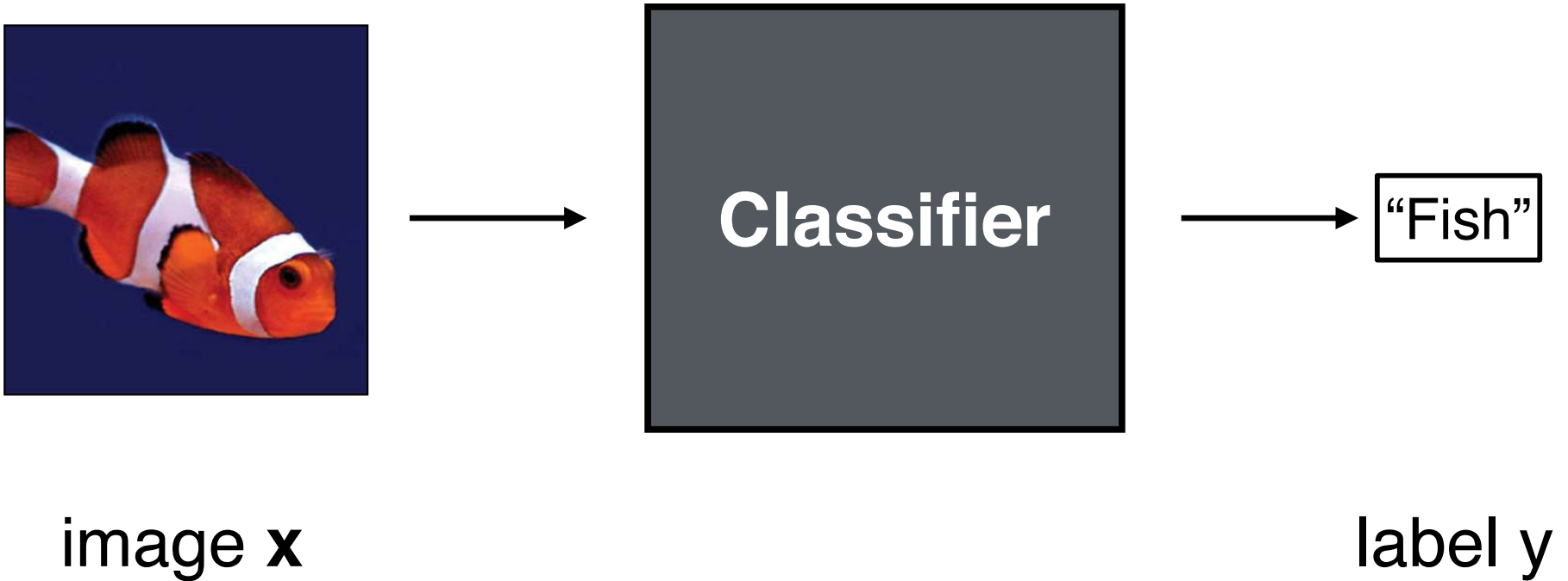
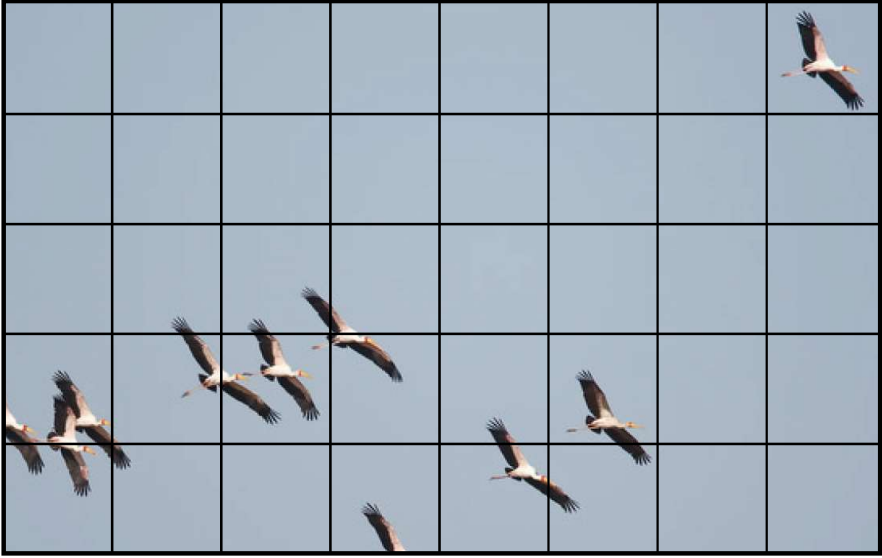
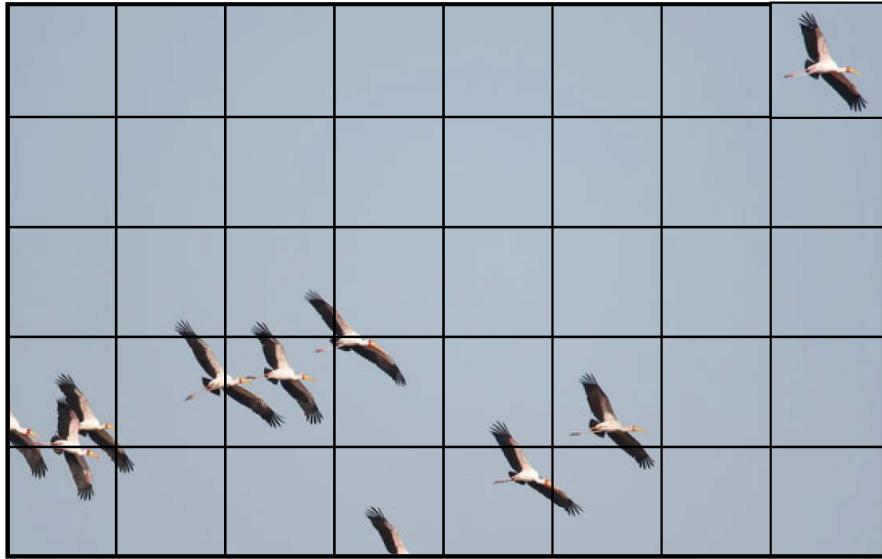
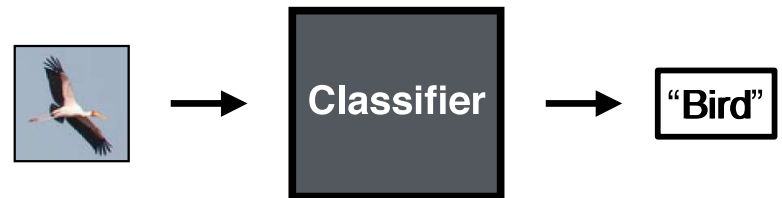
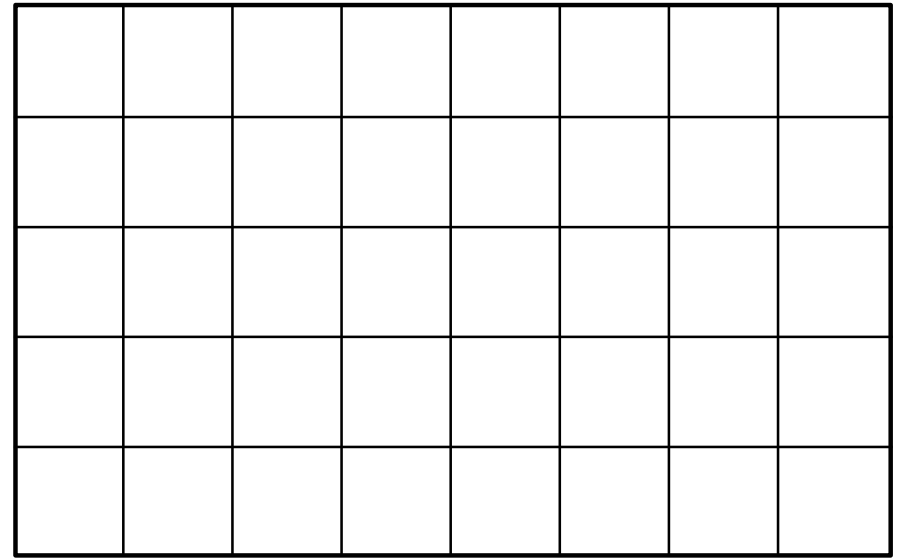
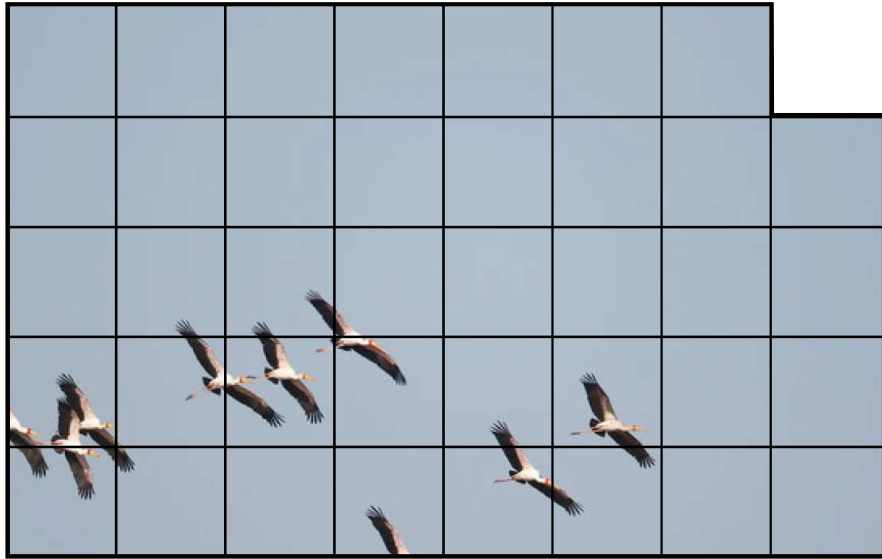




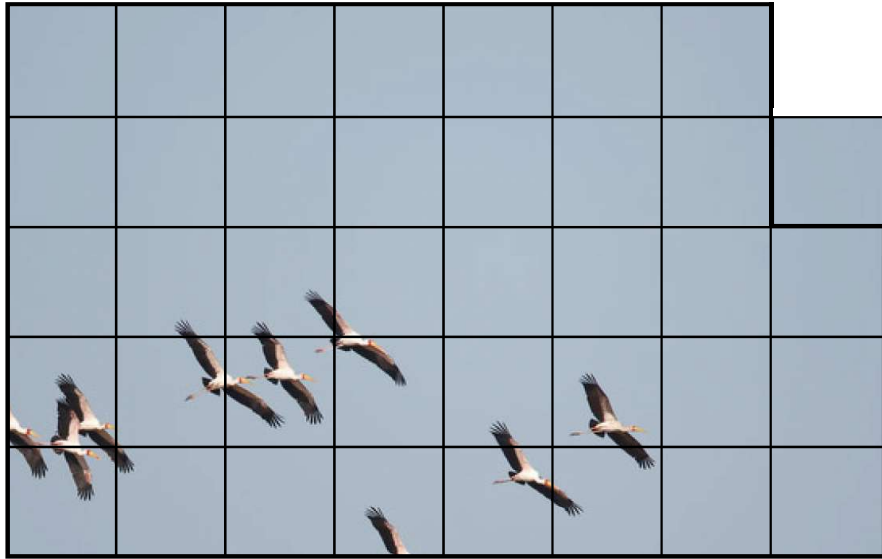
Photo credit: Fredo Durand





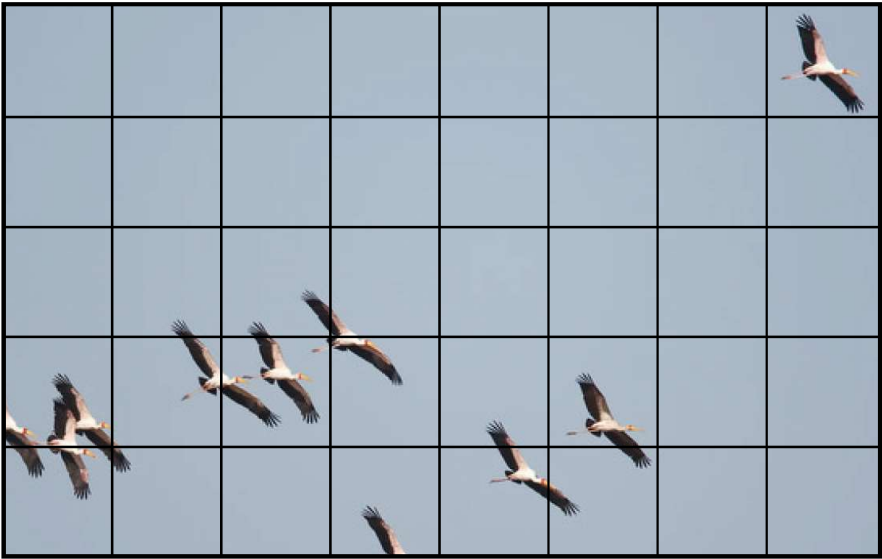




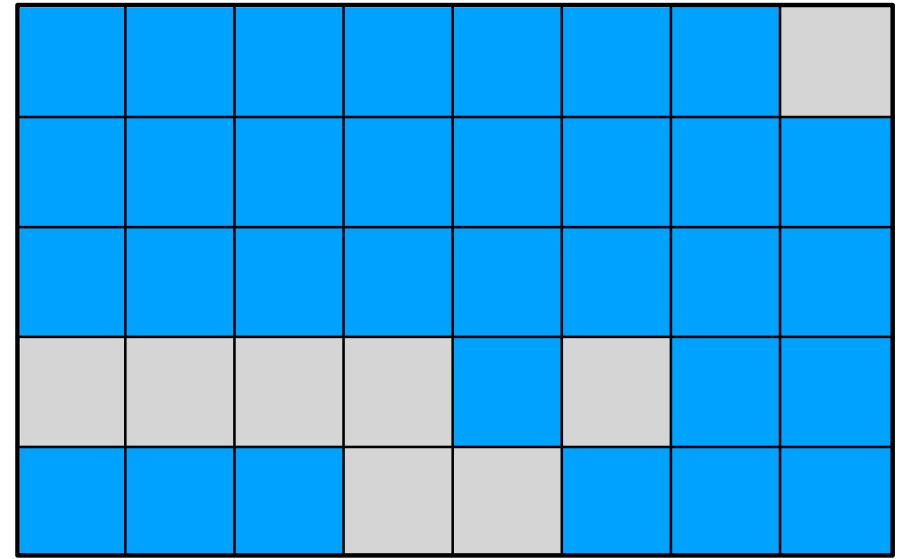
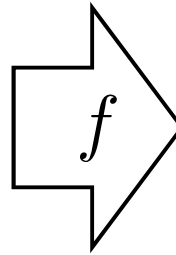
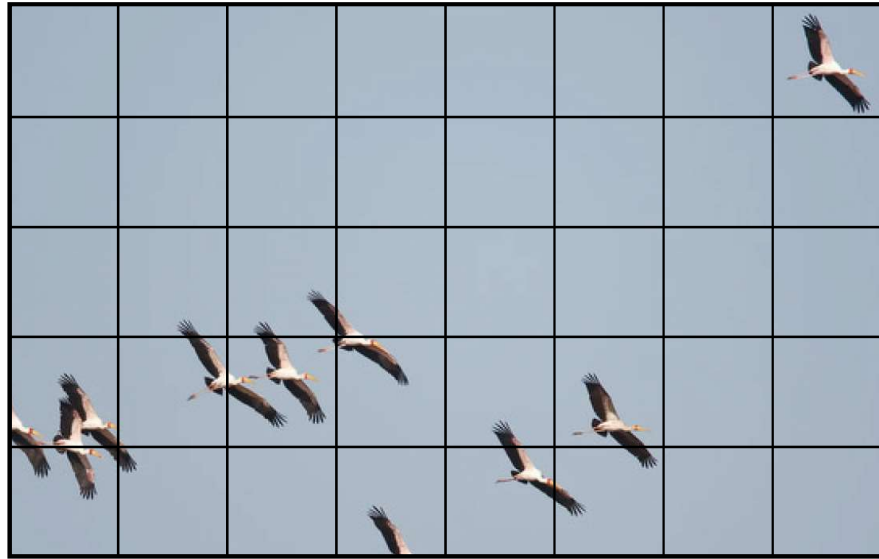


							Bird





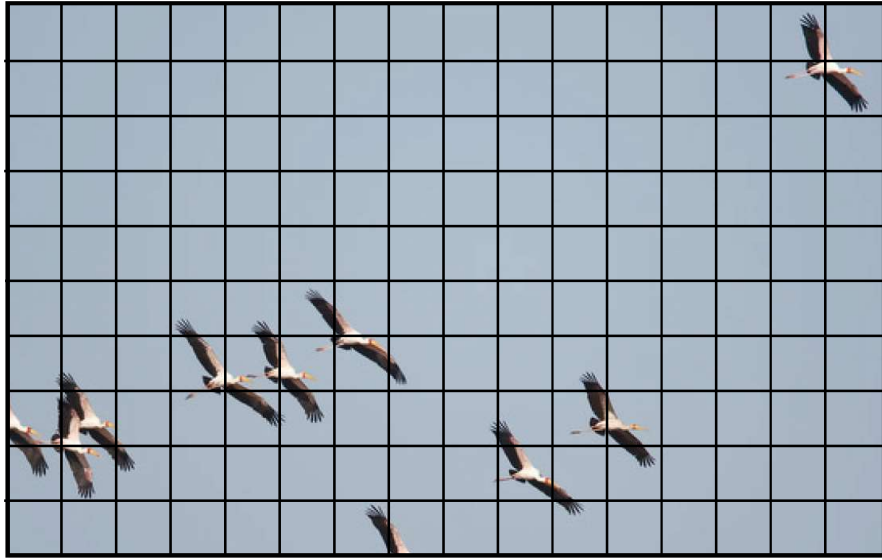
Sky	Sky	Sky	Sky	Sky	Sky	Sky	Bird
Sky	Sky	Sky	Sky	Sky	Sky	Sky	Sky
Sky	Sky	Sky	Sky	Sky	Sky	Sky	Sky
Bird	Bird	Bird	Sky	Bird	Sky	Sky	Sky
Sky	Sky	Sky	Bird	Sky	Sky	Sky	Sky



## Problem:

What happens to objects that are bigger?

What if an object crosses multiple cells?



“Cell”-based approach is limited.

What can we do instead?





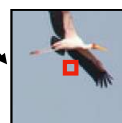
What's the object class of the center pixel?



"Bird"



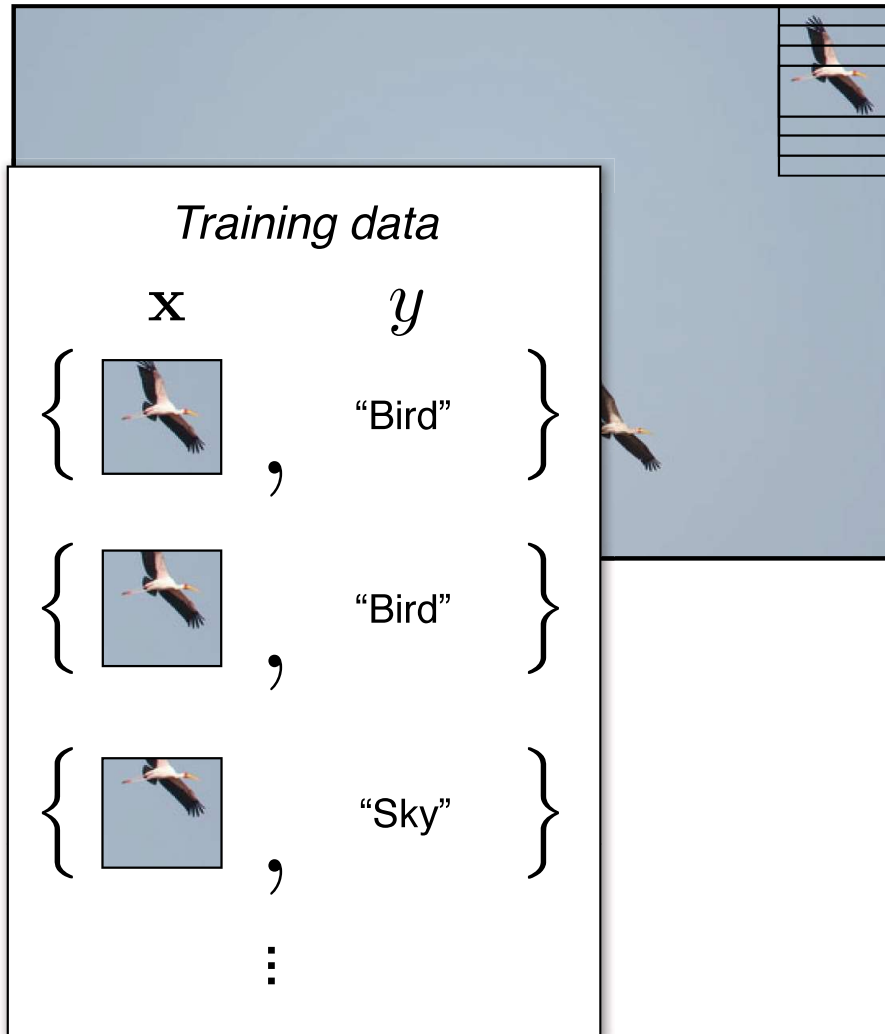
"Bird"



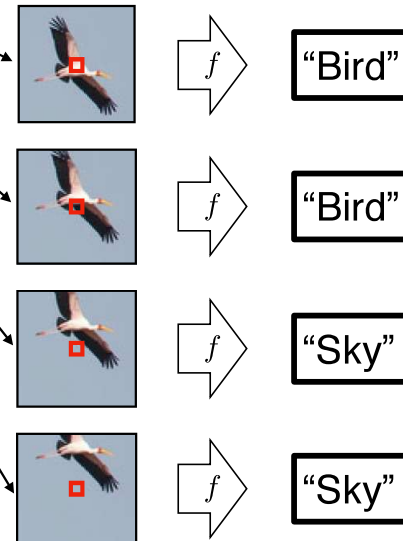
"Sky"

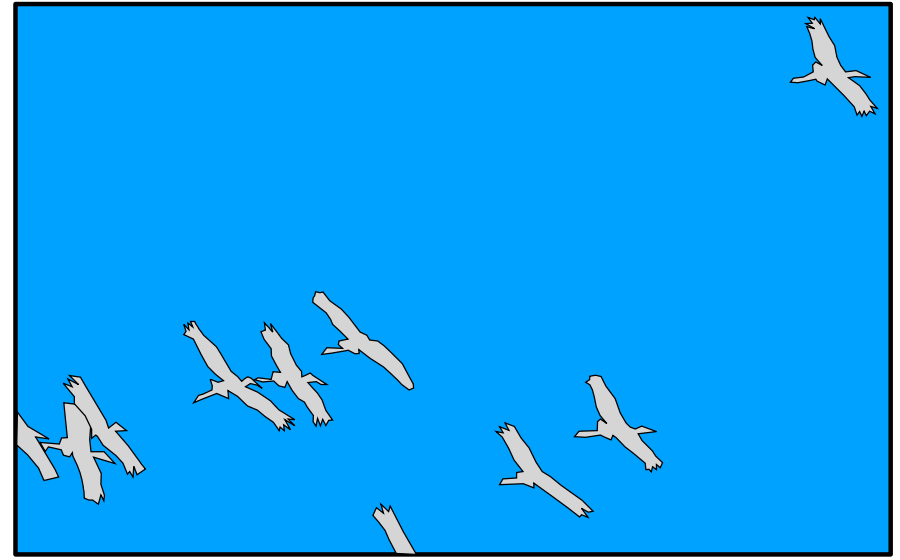
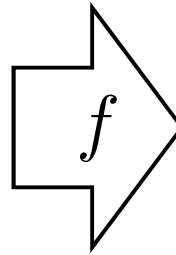


"Sky"



What's the object class of the center pixel?





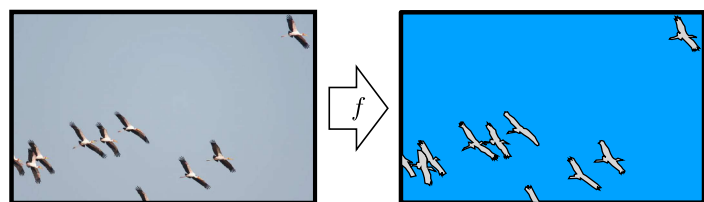
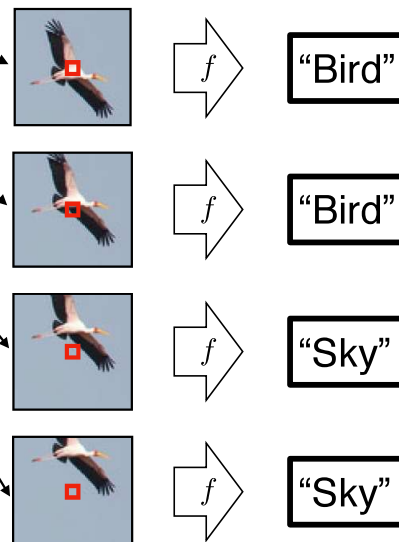
(Colors represent one-hot codes)

This problem is called **semantic segmentation**





What's the object class of the center pixel?

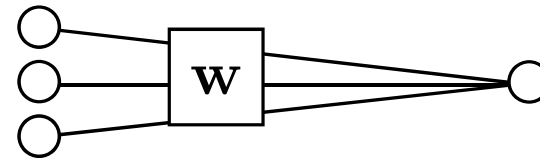


Translation invariance: process each patch in the same way.

An *equivariant* mapping:

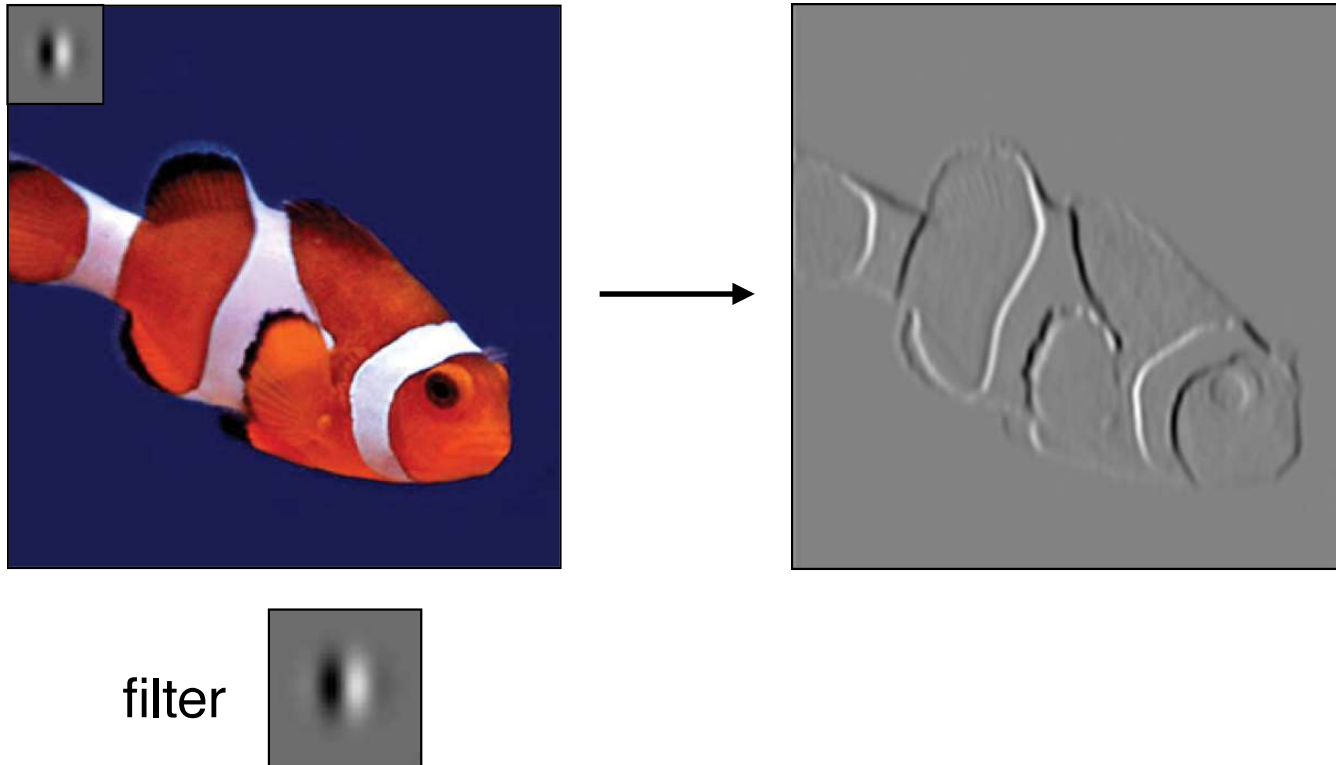
$$f(\text{translate}(x)) = \text{translate}(f(x))$$

**W** computes a weighted sum of all pixels in the patch



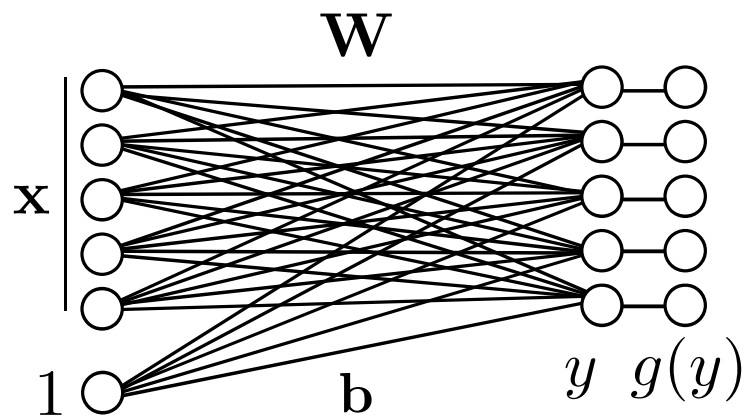
**W** is a convolutional kernel applied to the full image!

# Convolution

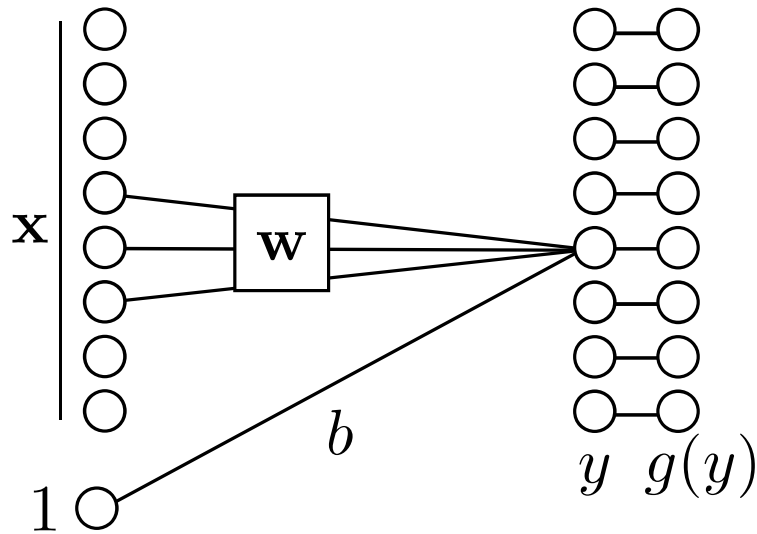


# Fully-connected network

## Fully-connected (fc) layer



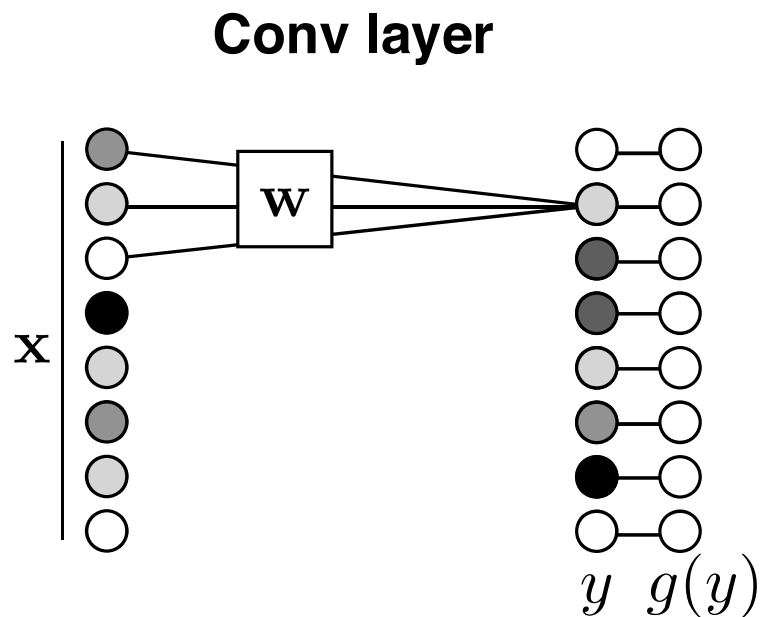
# Locally connected network



Often, we assume output is a **local** function of input.

If we use the same weights (**weight sharing**) to compute each local function, we get a convolutional neural network.

# Convolutional neural network

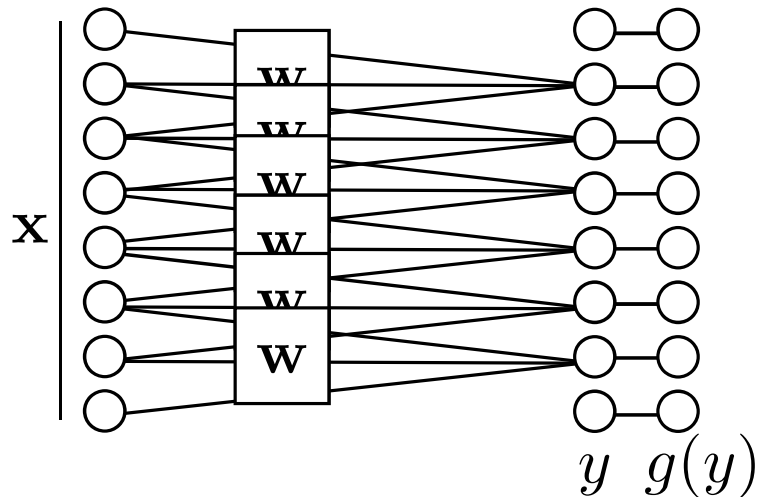


Often, we assume output is a **local** function of input.

If we use the same weights (**weight sharing**) to compute each local function, we get a convolutional neural network.

# Weight sharing

## Conv layer



Often, we assume output is a **local** function of input.

If we use the same weights (**weight sharing**) to compute each local function, we get a convolutional neural network.

**Toeplitz matrix**

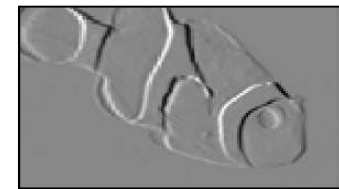
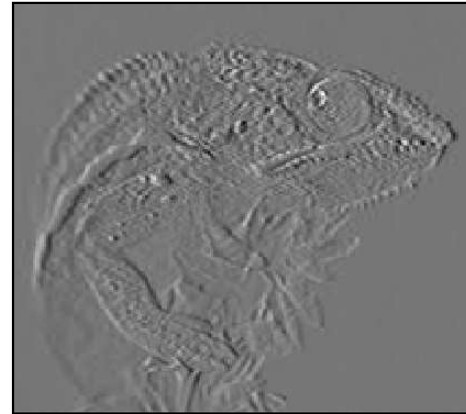
$$\begin{pmatrix} a & b & c & d & e \\ f & a & b & c & d \\ g & f & a & b & c \\ h & g & f & a & b \\ i & h & g & f & a \end{pmatrix}$$

$$\mathbf{x}^{(l+1)} = \begin{matrix} \text{[Toeplitz Matrix]} \end{matrix} * \mathbf{x}^{(l)}$$

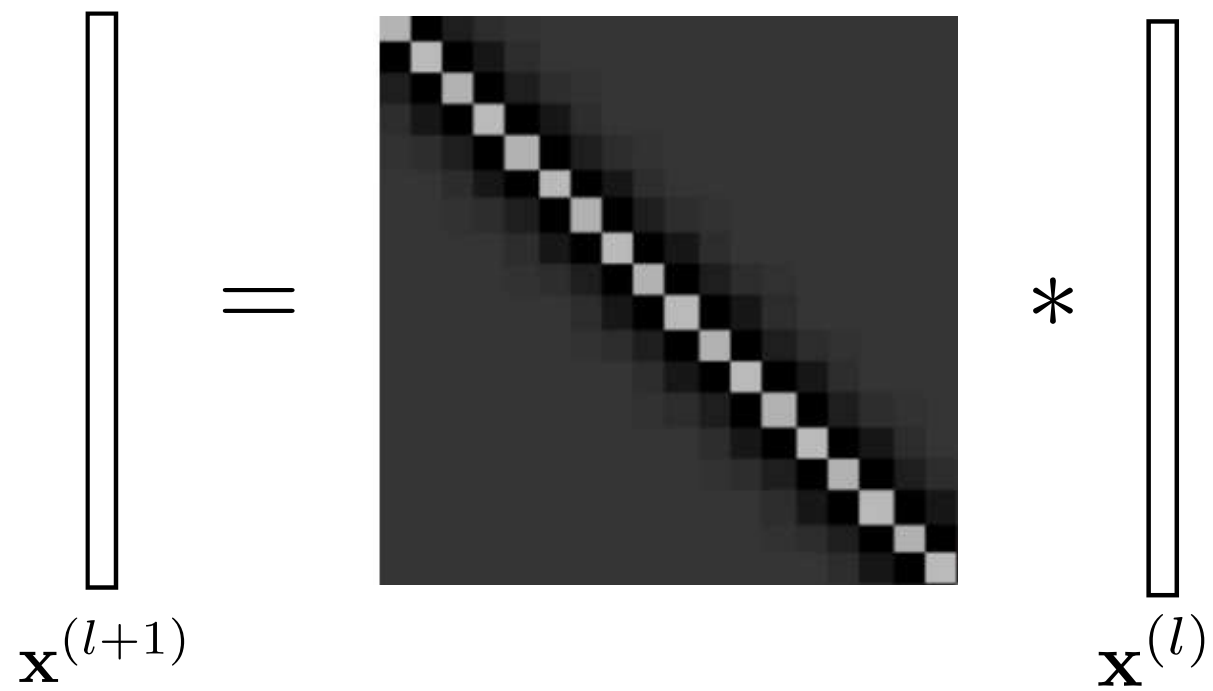
e.g., pixel image

- Constrained linear layer
- Fewer parameters → easier to learn, less overfitting





Conv layers can be applied to arbitrarily-sized inputs

$$\mathbf{x}^{(l+1)} = \mathbf{W} * \mathbf{x}^{(l)}$$


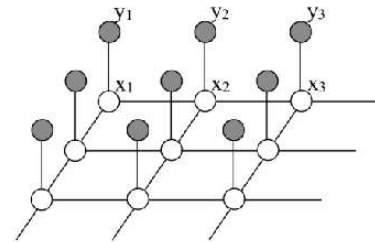
$$\mathbf{x}^{(l+1)} = \begin{bmatrix} \text{diagonal matrix} \end{bmatrix} * \mathbf{x}^{(l)}$$

The diagram illustrates a vector-matrix multiplication. On the left, a vertical rectangle represents the vector  $\mathbf{x}^{(l+1)}$ . In the center, a square matrix is shown with a prominent diagonal of alternating black and white squares, indicating a diagonal matrix. To the right of the matrix is an asterisk  $*$ , representing multiplication. On the far right, another vertical rectangle represents the vector  $\mathbf{x}^{(l)}$ .

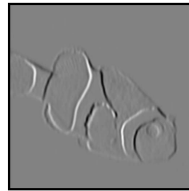
# Five views on convolutional layers

1. Equivariant with translation (stationarity)  $f(\text{translate}(x)) = \text{translate}(f(x))$

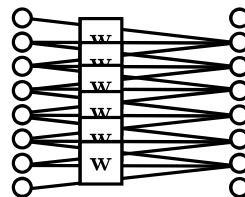
2. Patch processing (Markov assumption)



3. Image filter



4. Parameter sharing



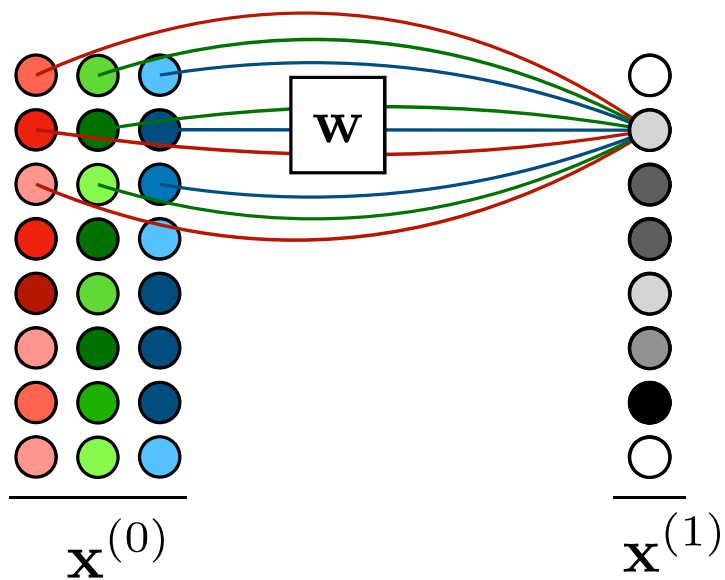
5. A way to process variable-sized tensors

# What if we have color?

(aka multiple input channels?)

# Multiple channels

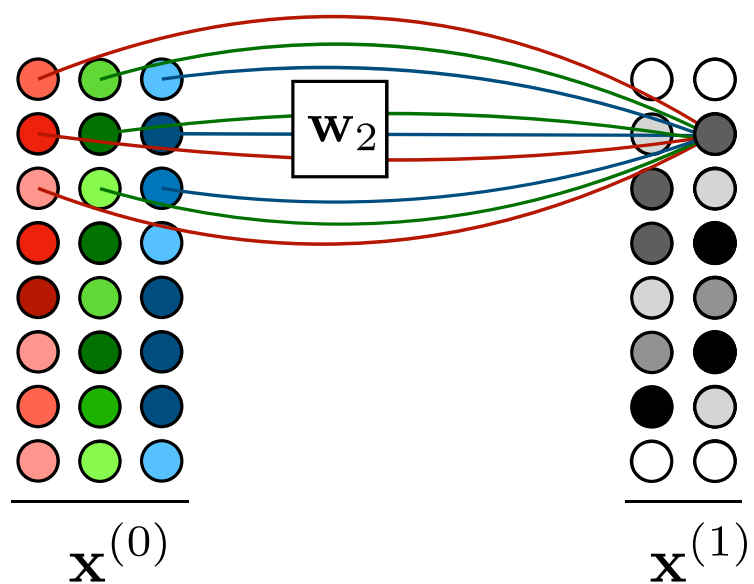
**Conv layer**



$$\mathbb{R}^{N \times C} \rightarrow \mathbb{R}^{N \times 1}$$

# Multiple channels

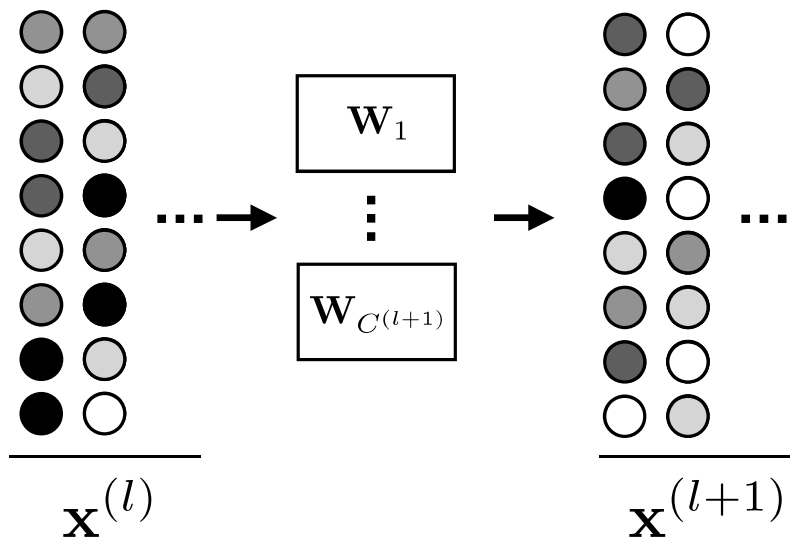
**Conv layer**



$$\mathbb{R}^{N \times C^{(0)}} \rightarrow \mathbb{R}^{N \times C^{(1)}}$$

# Multiple channels

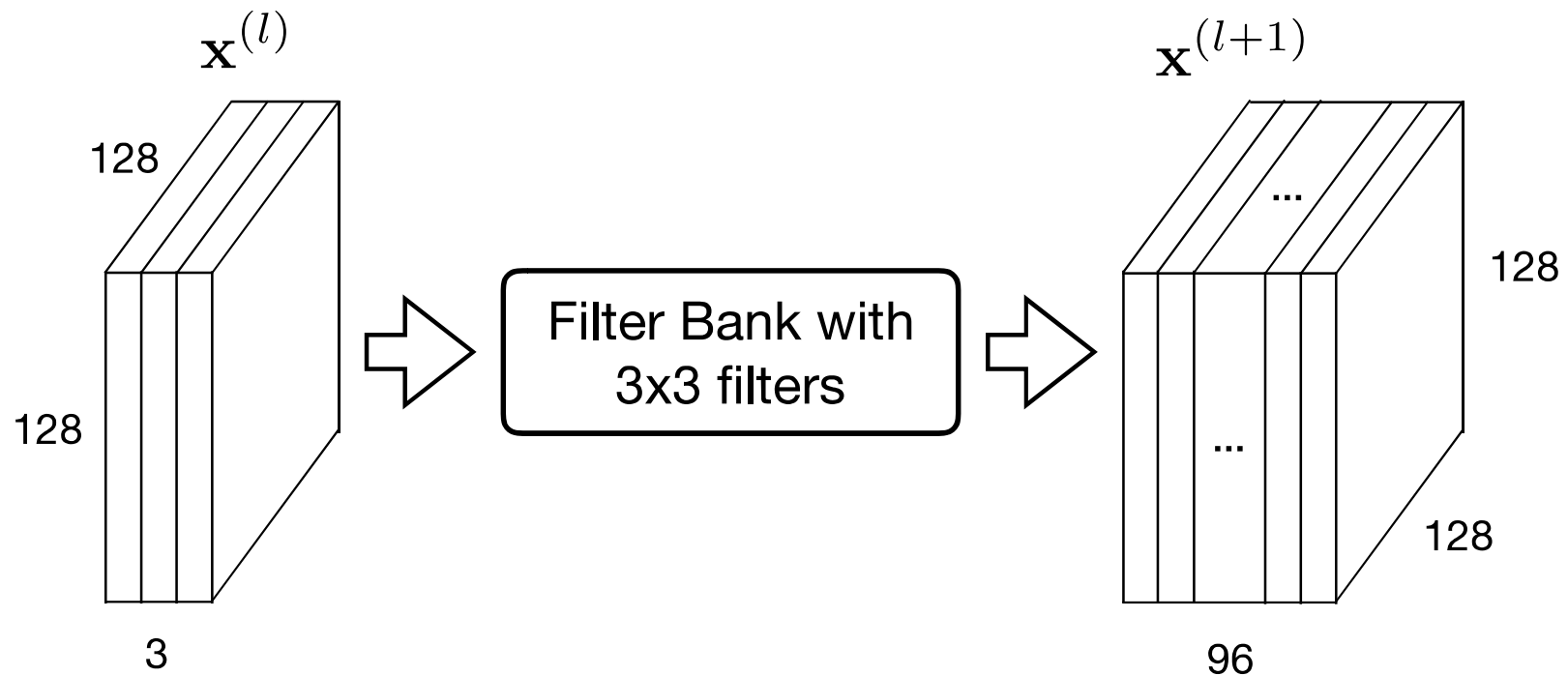
## Conv layer



$$\mathbb{R}^{N \times C^{(l)}} \rightarrow \mathbb{R}^{N \times C^{(l+1)}}$$



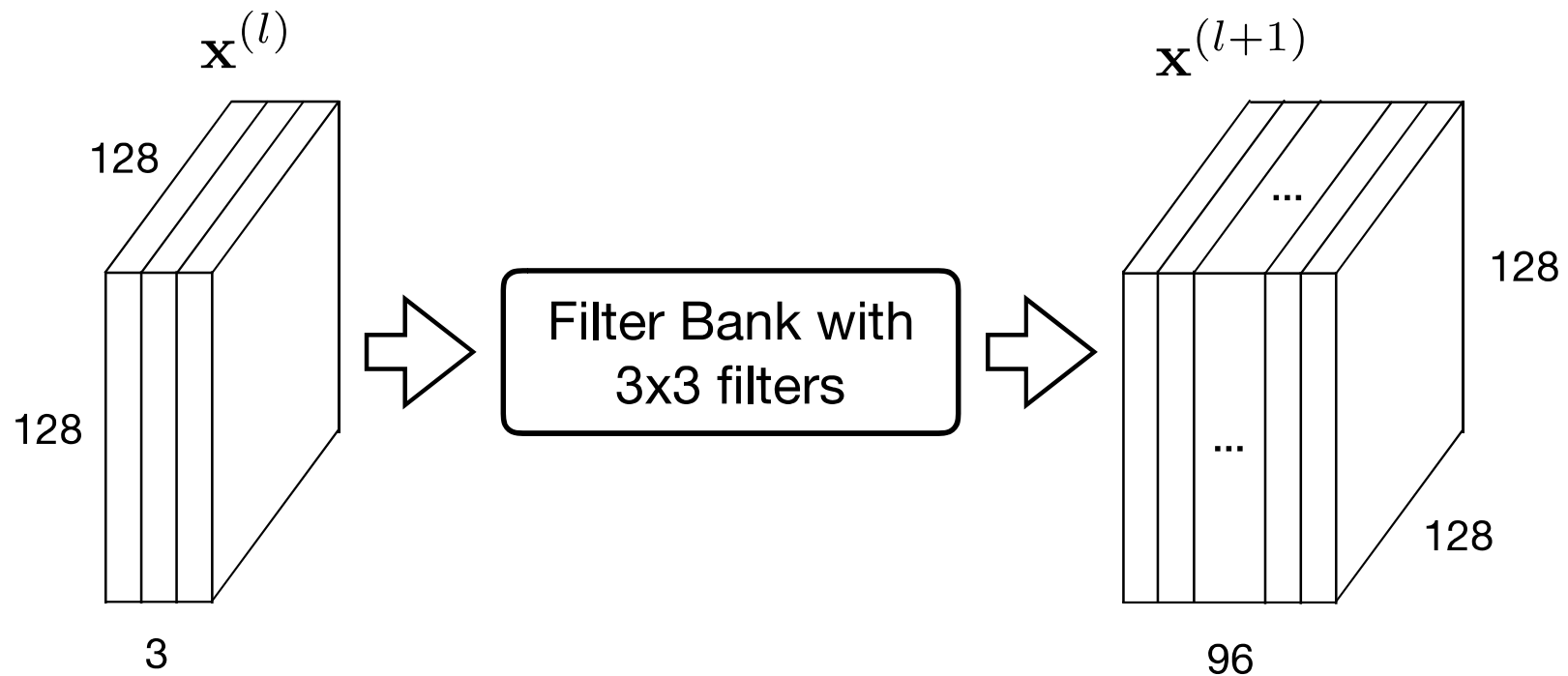
## Multiple channels: Example



How many parameters does *each filter* have?

- (a) 9      (b) 27      (c) 96      (d) 864

## Multiple channels: Example



How many filters are in the bank?

- (a) 3      (b) 27      (c) 96      (d) can't say

# Filter sizes

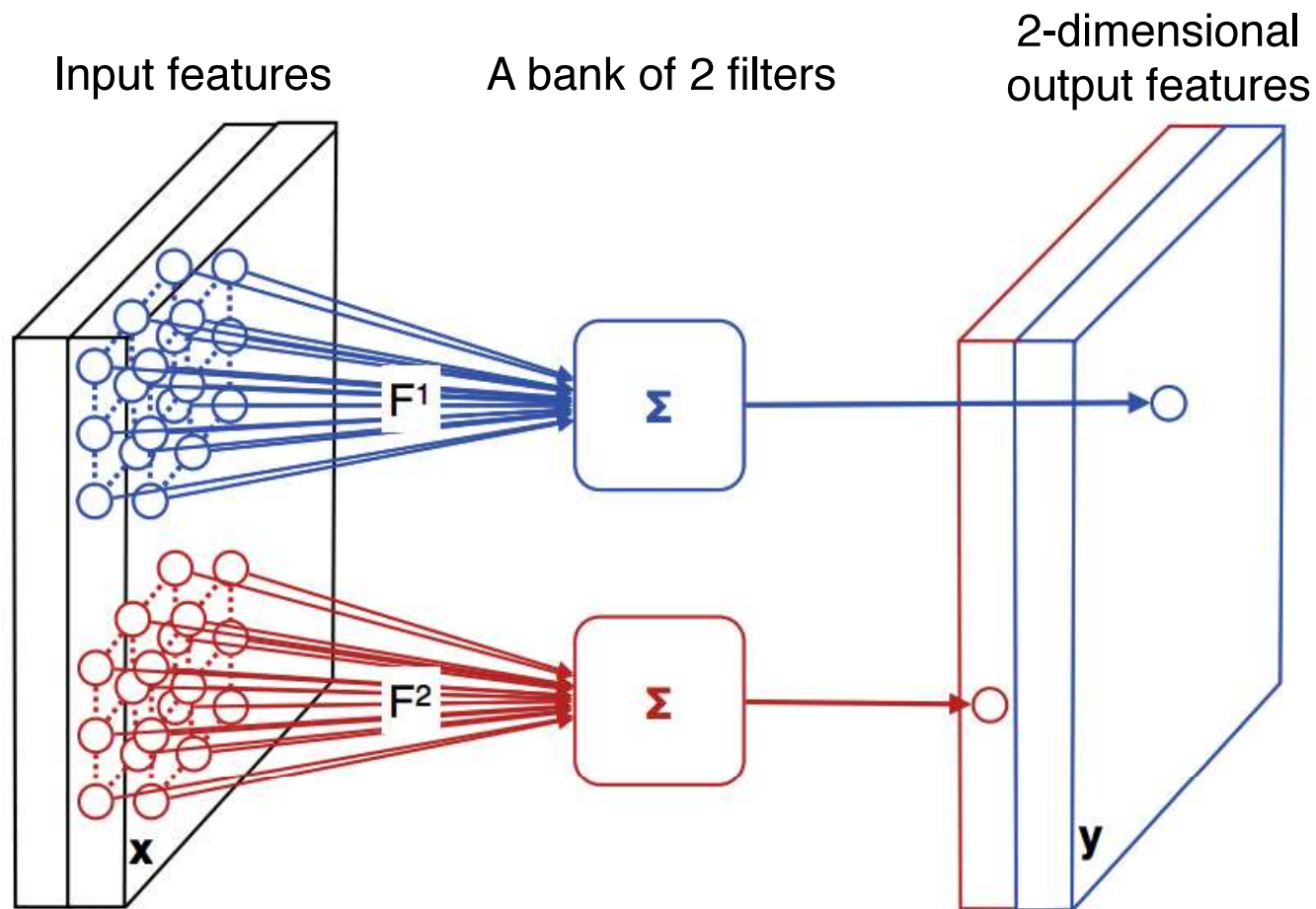
When mapping from

$$\mathbf{x}^{(l)} \in \mathbb{R}^{H \times W \times C^{(l)}} \rightarrow \mathbf{x}^{(l+1)} \in \mathbb{R}^{H \times W \times C^{(l+1)}}$$

using an filter of spatial extent  $M \times N$

Number of parameters per filter:  $M \times N \times C^{(l)}$

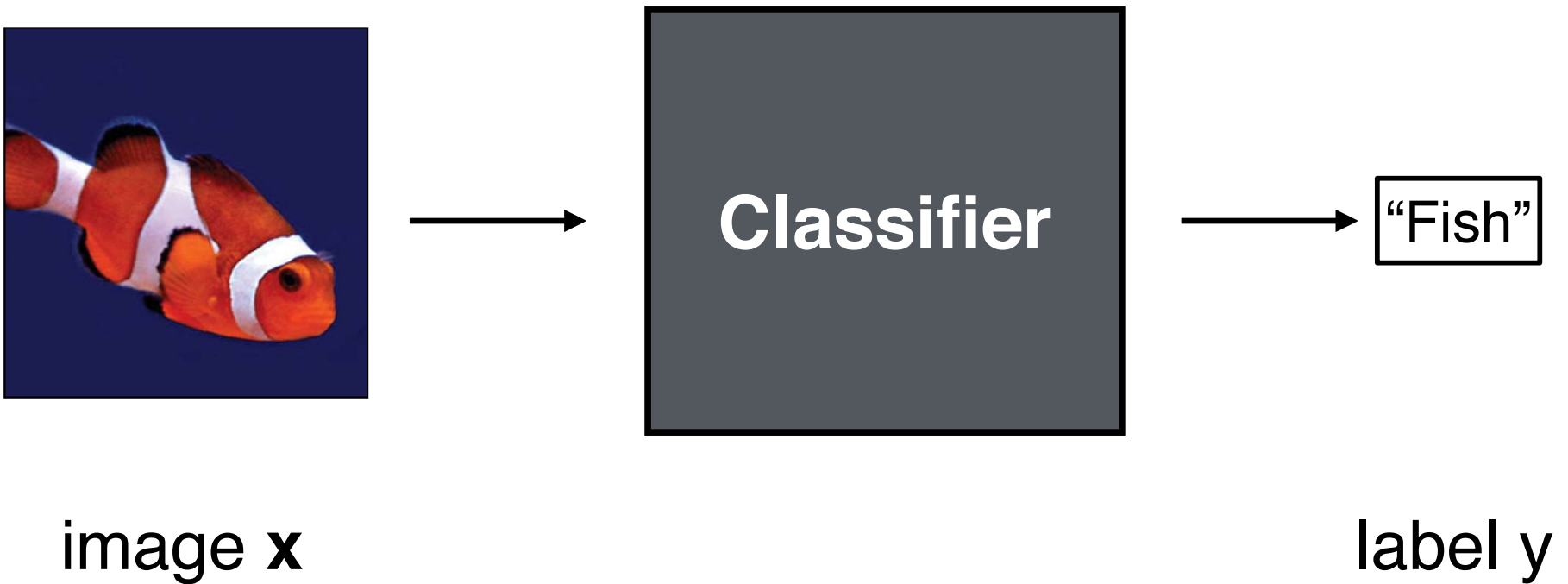
Number of filters:  $C^{(l+1)}$



$$\mathbb{R}^{H \times W \times C^{(l)}} \rightarrow \mathbb{R}^{H \times W \times C^{(l+1)}}$$

[Figure from Andrea Vedaldi]

# Image classification



# Image classification

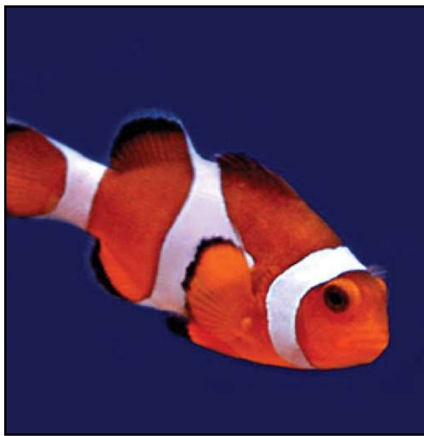
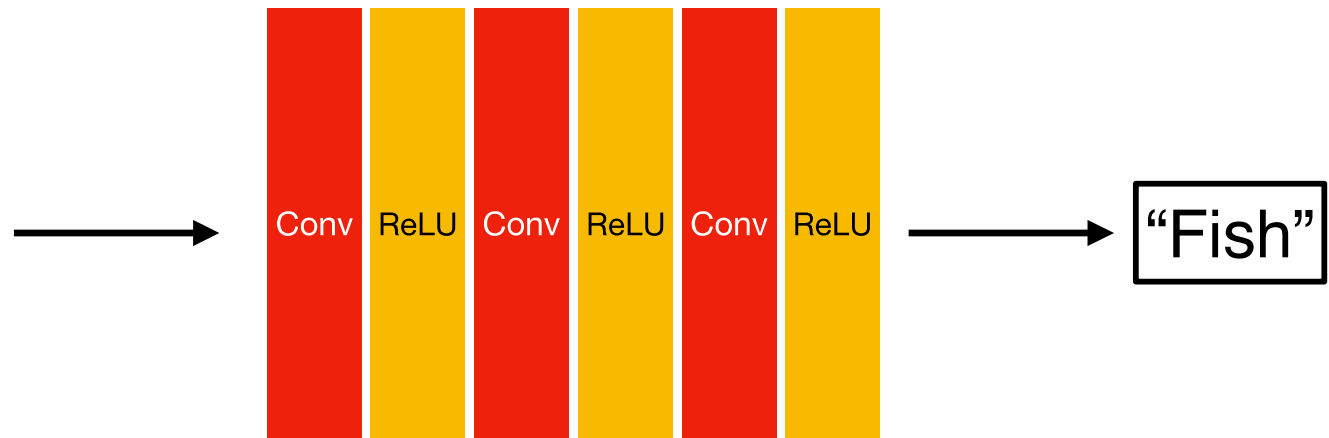
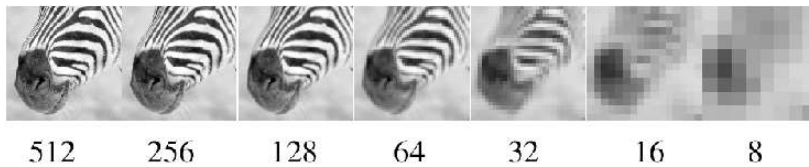


image  $x$

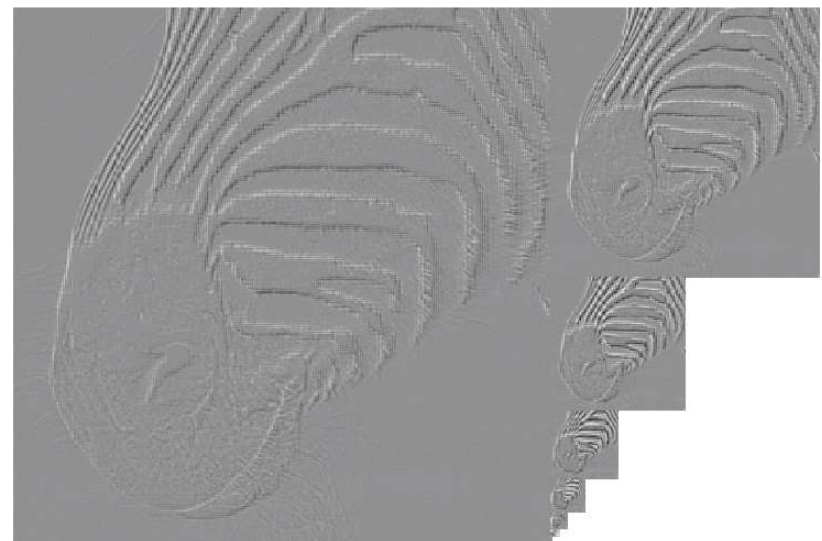
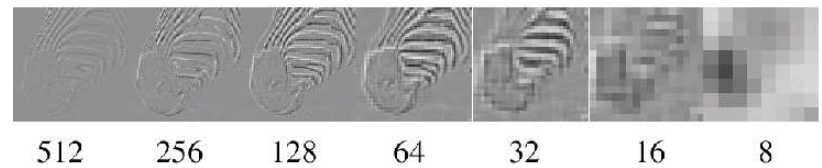


label  $y$

## Multiscale representations are great!



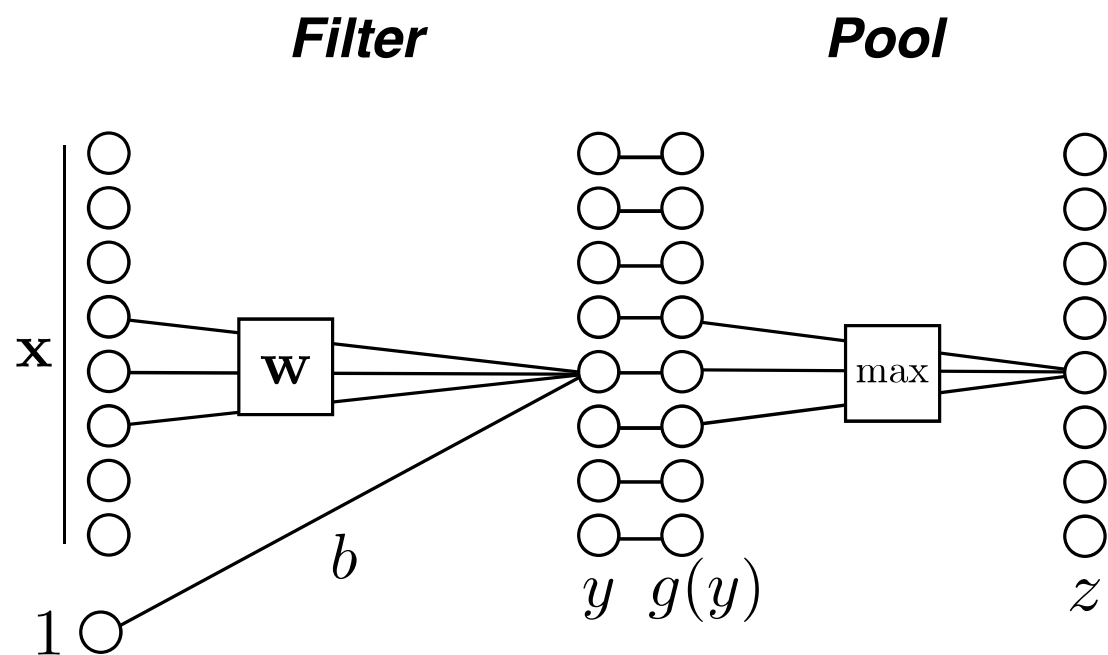
Gaussian Pyr



Laplacian Pyr

**How can we use multi-scale modeling in Convnets?**

# Pooling

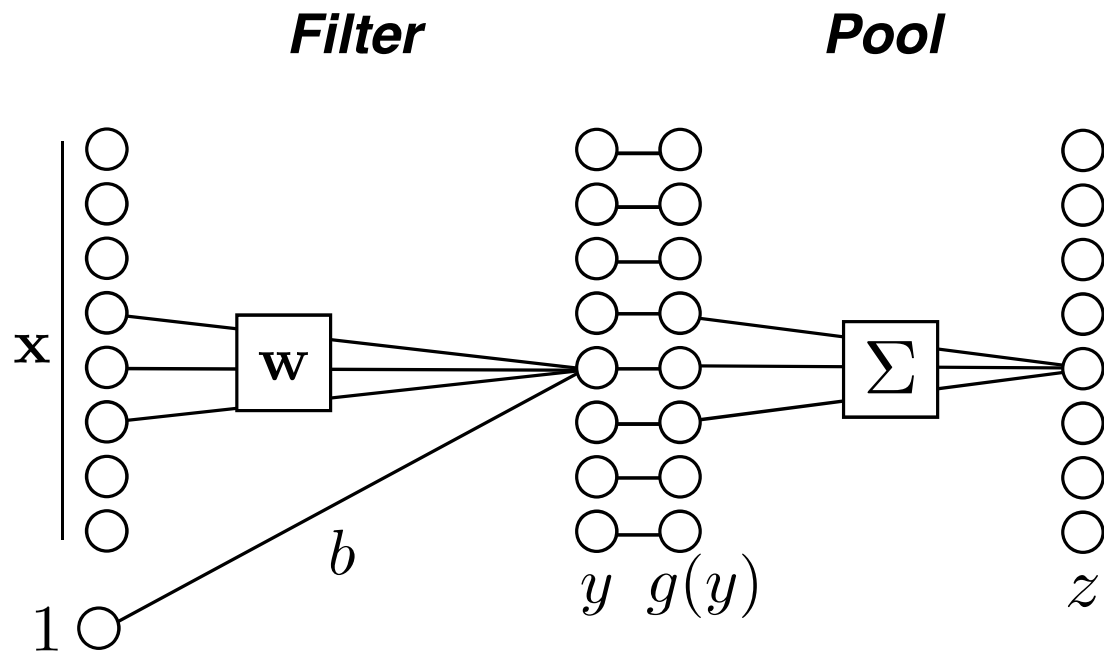


## Max pooling

$$z_k = \max_{j \in \mathcal{N}(k)} g(y_j)$$



# Pooling



## Max pooling

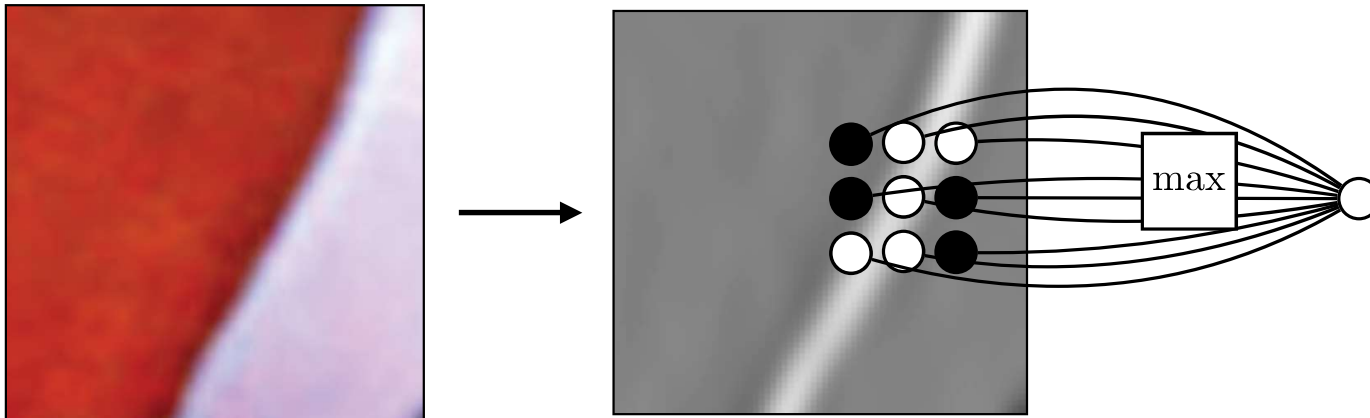
$$z_k = \max_{j \in \mathcal{N}(k)} g(y_j)$$

## Mean pooling

$$z_k = \frac{1}{|\mathcal{N}(k)|} \sum_{j \in \mathcal{N}(k)} g(y_j)$$

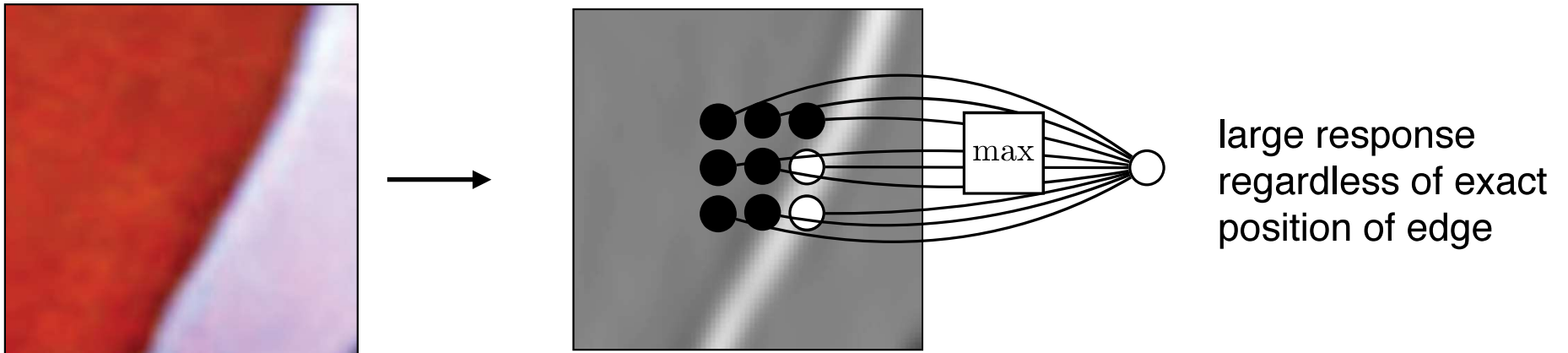
# Pooling — Why?

Pooling across spatial locations achieves stability w.r.t. small translations:



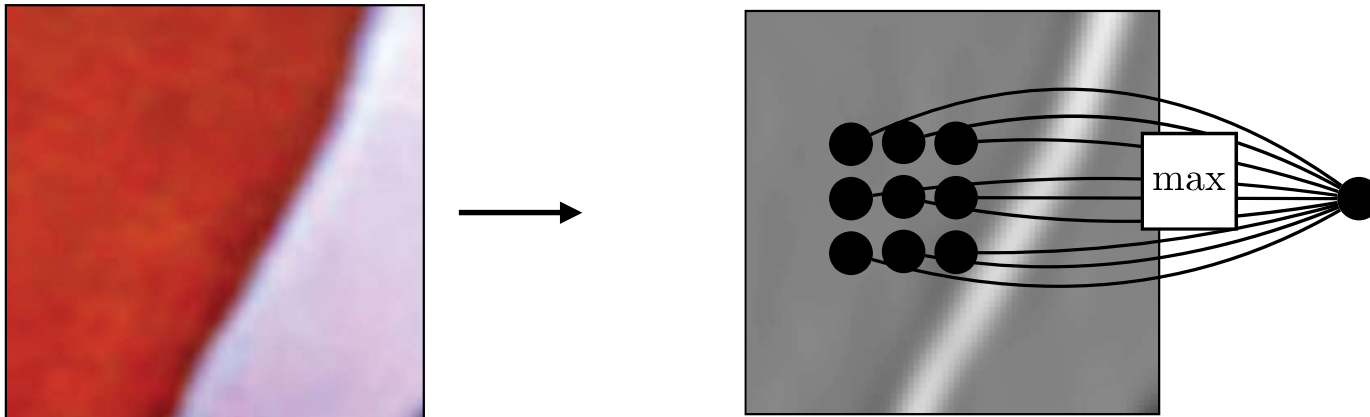
# Pooling — Why?

Pooling across spatial locations achieves stability w.r.t. small translations:



# Pooling — Why?

Pooling across spatial locations achieves stability w.r.t. small translations:



CNNs are stable w.r.t. diffeomorphisms



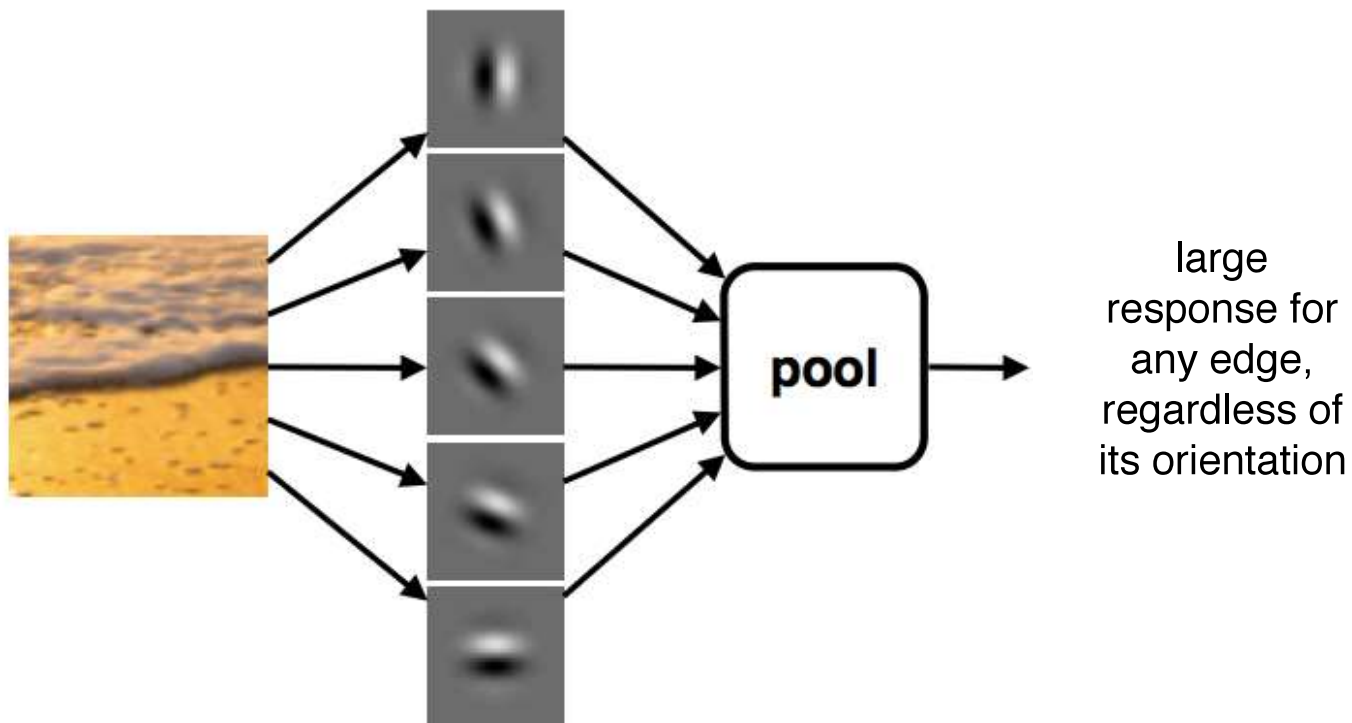
$\approx$



[“Unreasonable effectiveness of Deep Features as a Perceptual Metric”, Zhang et al. 2018]

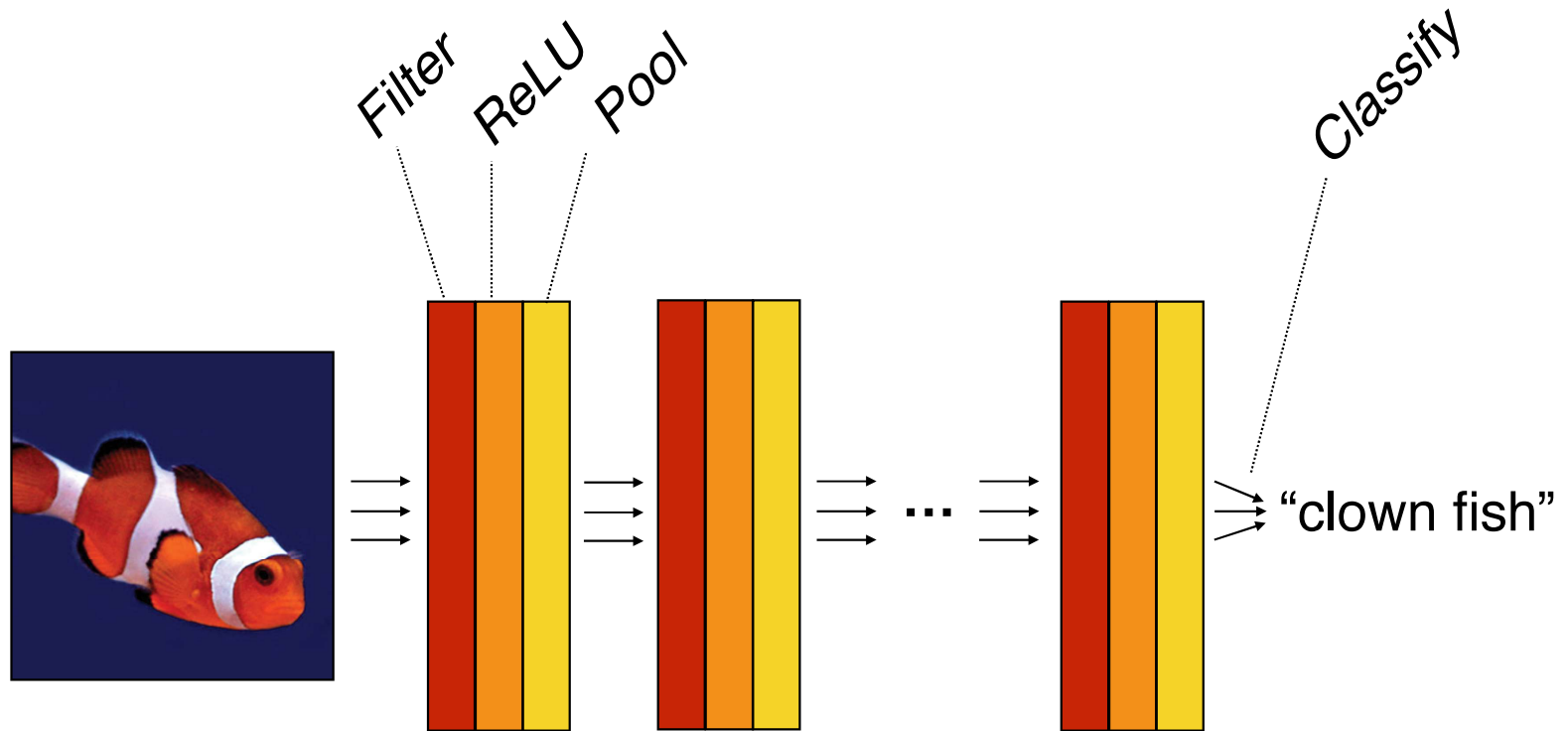
# Pooling — Why?

Pooling across feature channels (filter outputs) can achieve other kinds of invariances:



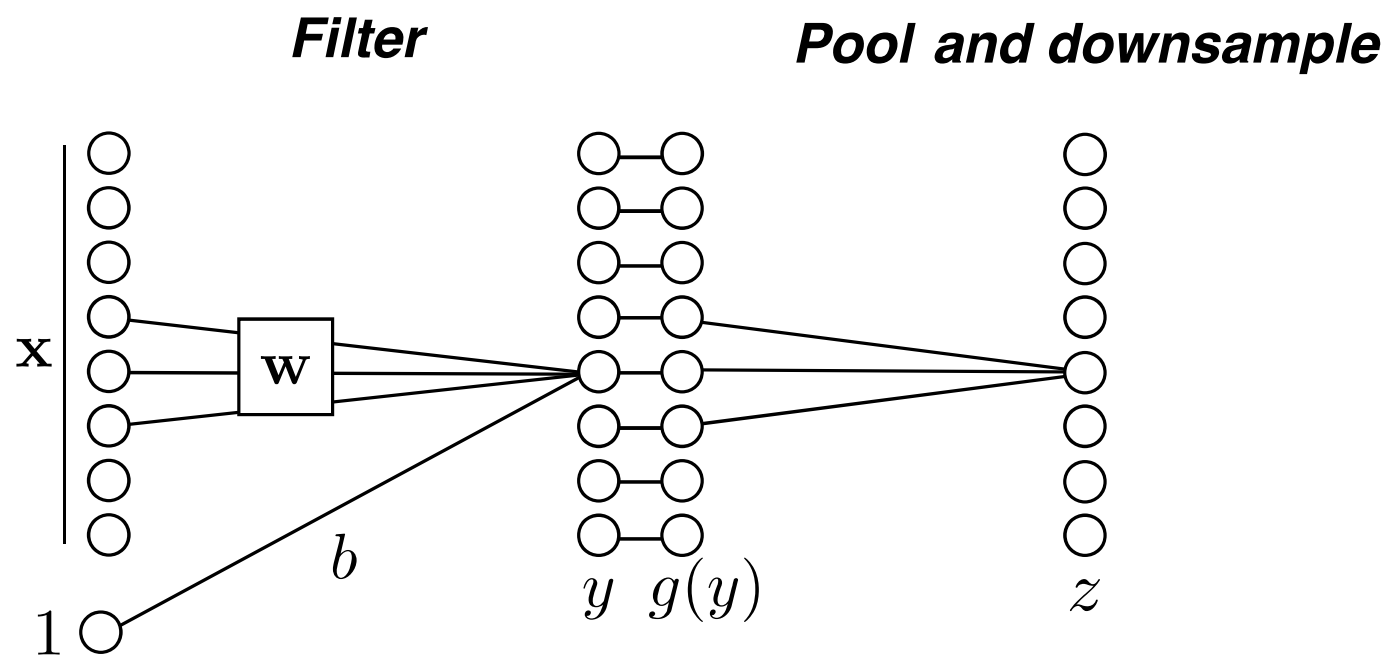
[Derived from slide by Andrea Vedaldi]

# Computation in a neural net



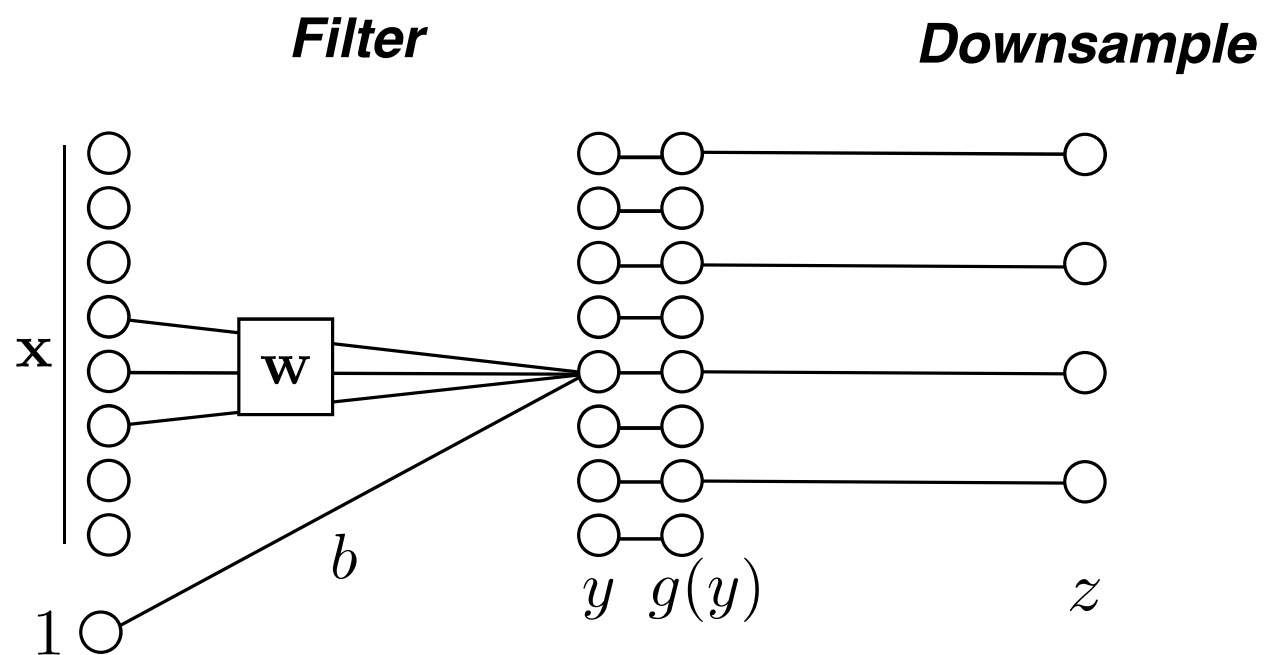
$$f(\mathbf{x}) = f_L(\dots f_2(f_1(\mathbf{x})))$$

# Downsampling



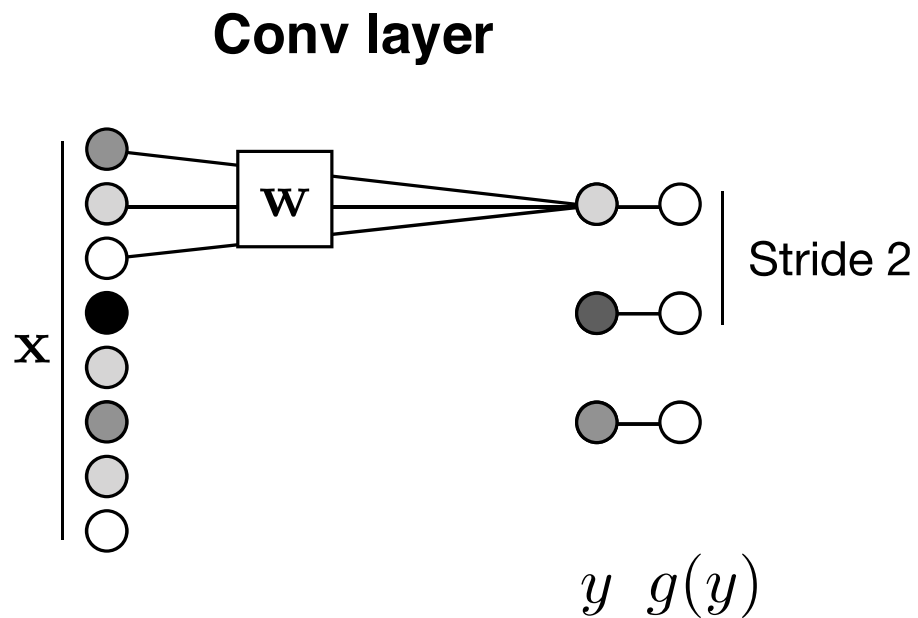


# Downsampling



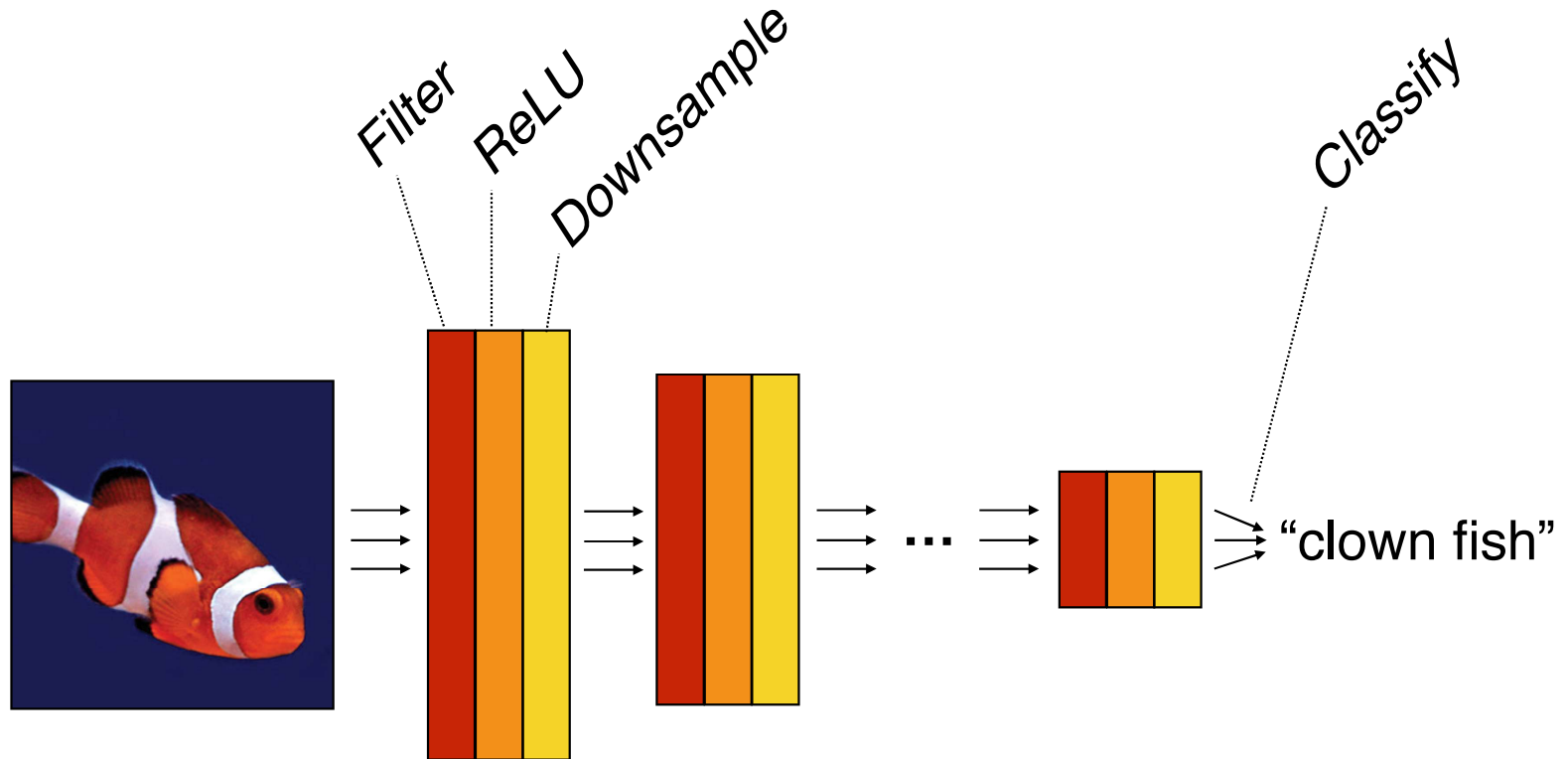
$$\mathbb{R}^{H^{(l)} \times W^{(l)} \times C^{(l)}} \rightarrow \mathbb{R}^{H^{(l+1)} \times W^{(l+1)} \times C^{(l+1)}}$$

# Strided operations



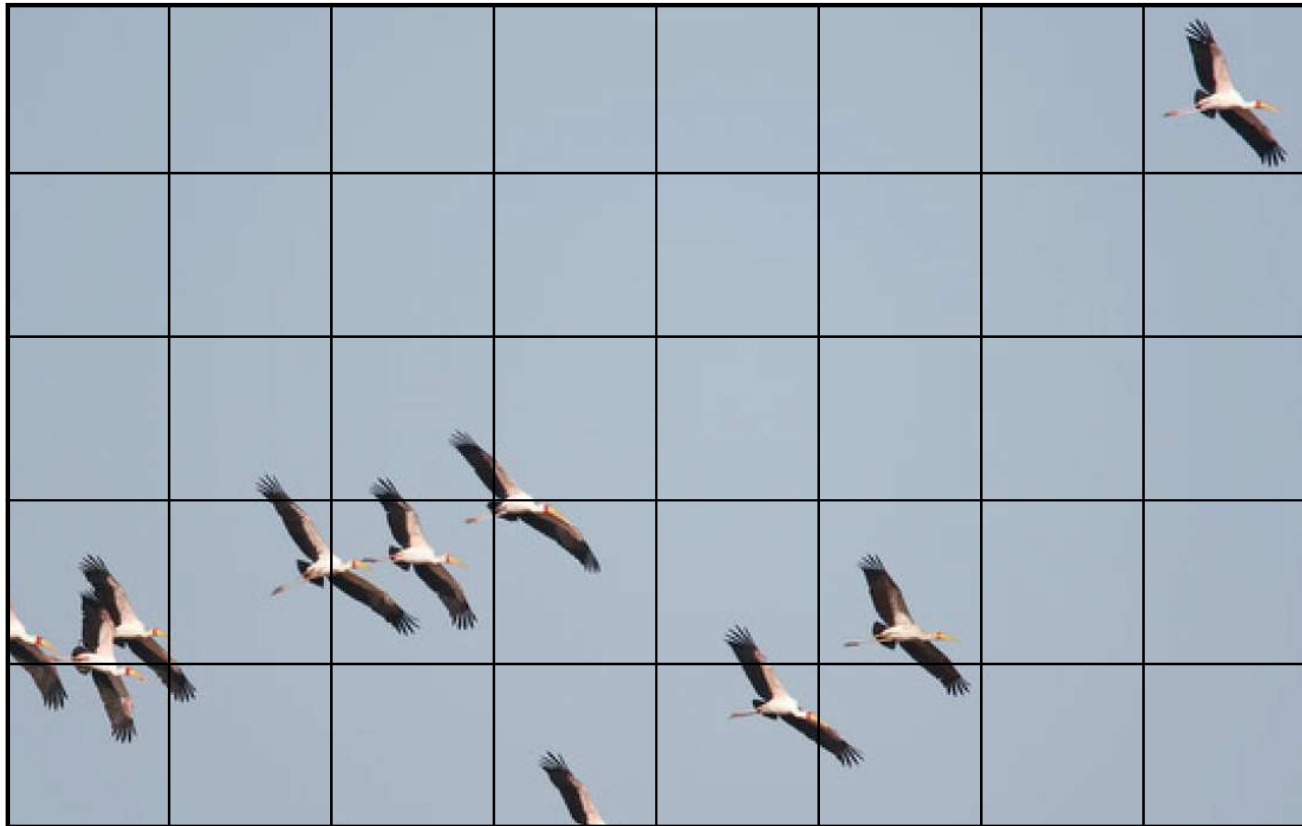
**Strided operations** combine a given operation (convolution or pooling) and downsampling into a single operation.

# Computation in a neural net

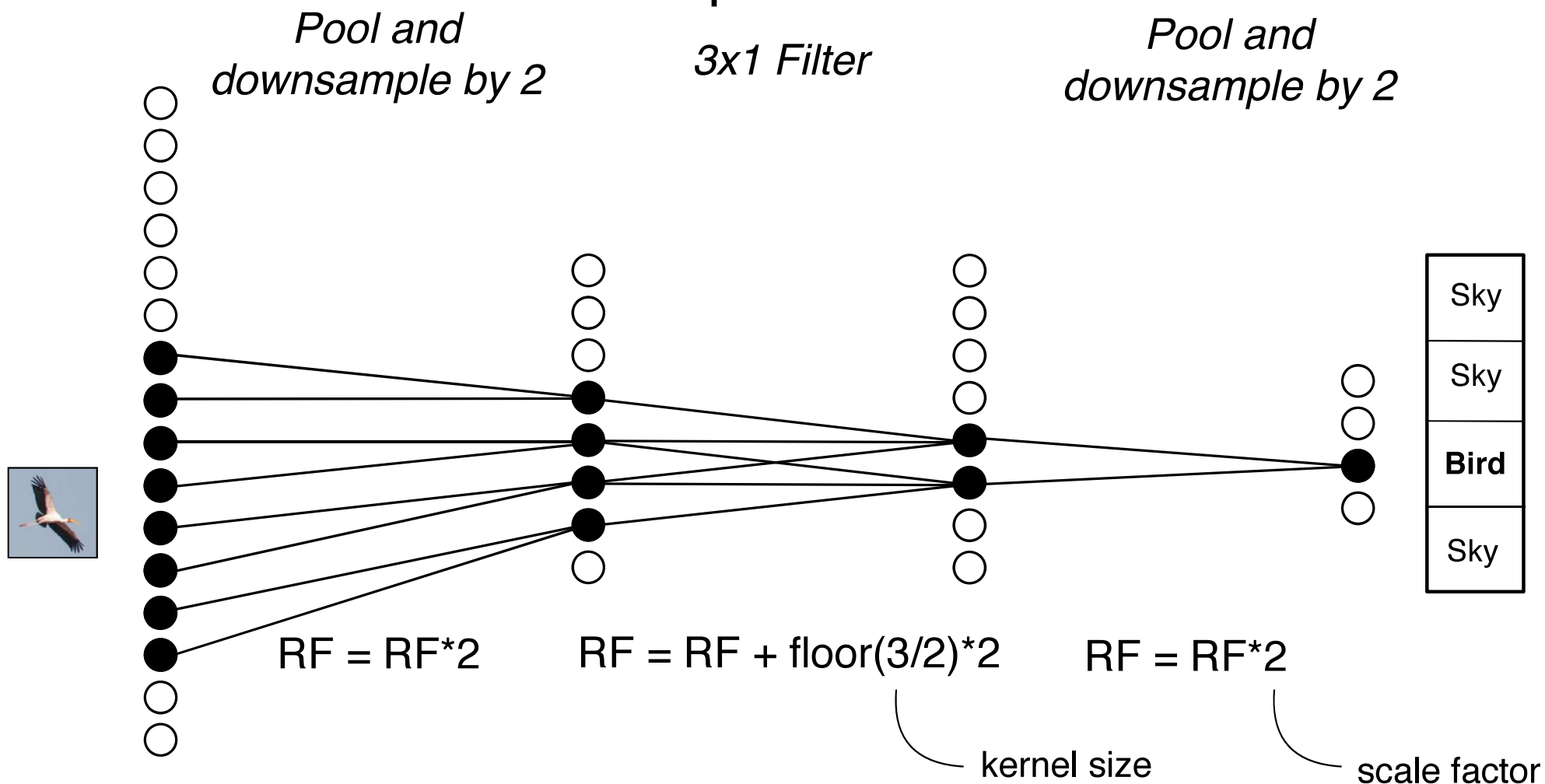


$$f(\mathbf{x}) = f_L(\dots f_2(f_1(\mathbf{x})))$$

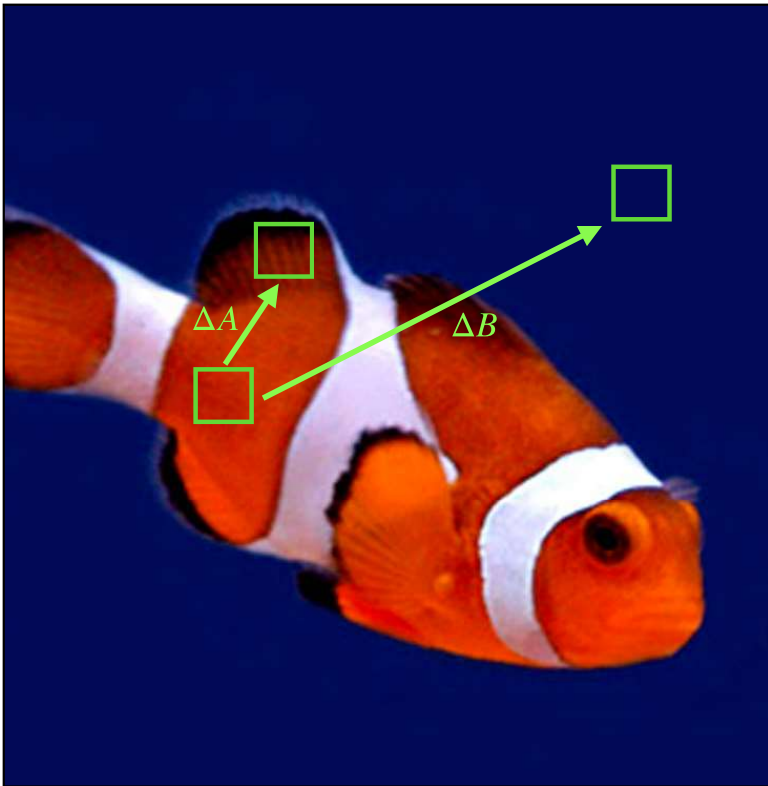
# Receptive fields



# Receptive fields



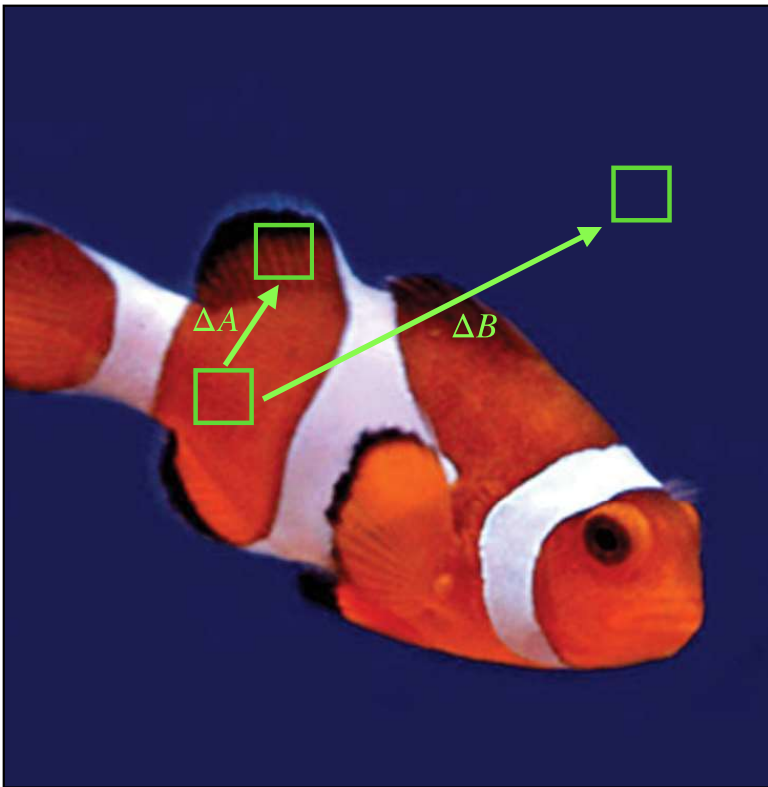
# Local vs. global processing



**Across all images, which is higher:**

- (1) correlation between points with distance  $\Delta A$
- (2) correlation between points with distance  $\Delta B$
- (3) can't say

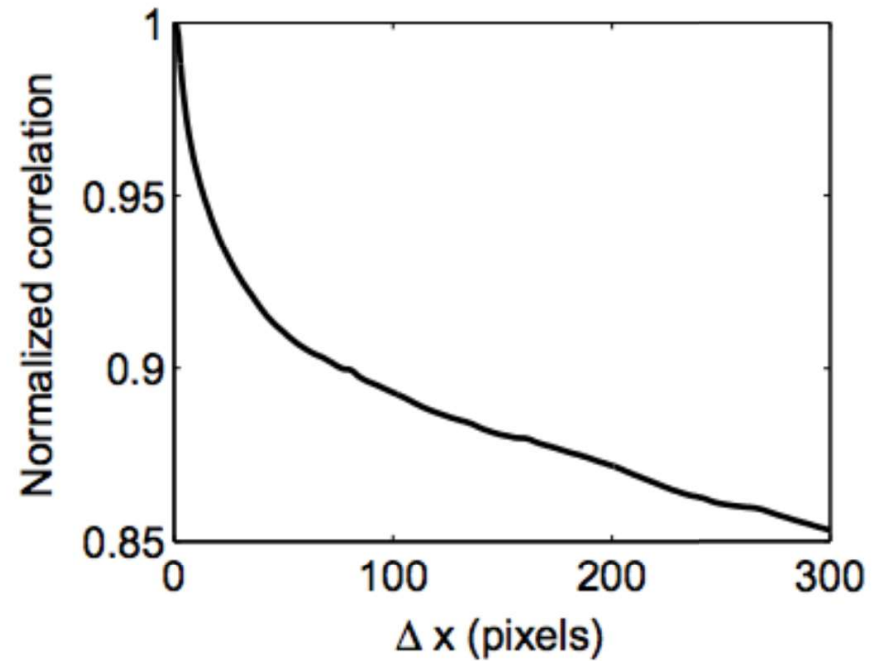
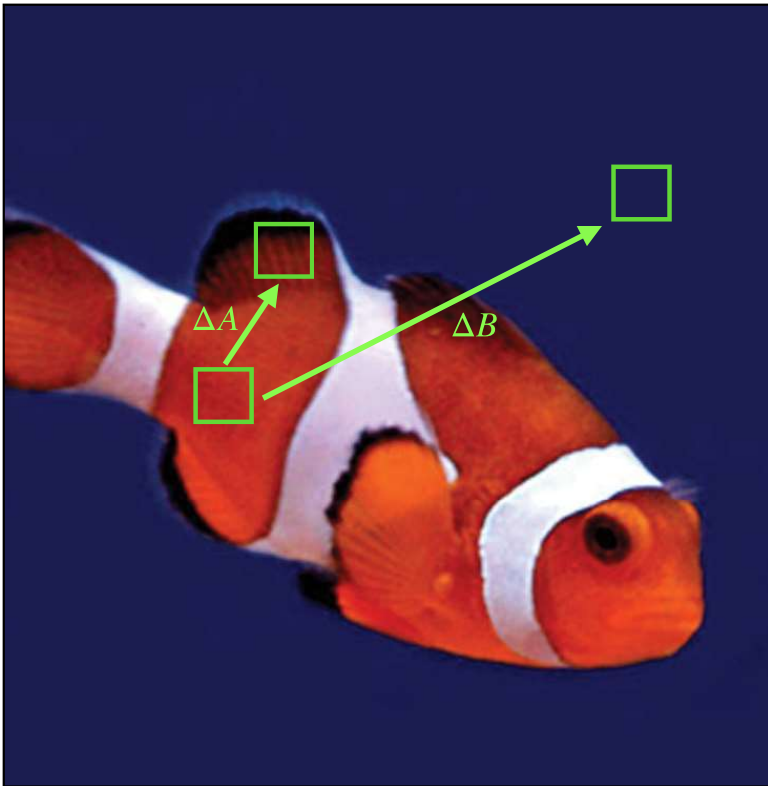
# Local vs. global processing



**Across all images, which is higher:**

- (1) correlation between points with distance  $\Delta A$
- (2) correlation between points with distance  $\Delta B$
- (3) can't say

# Local vs. global processing



*[Simoncelli: Statistical Modelling of Photographic Images, 2005]*



# CNNs — Why?

Statistical dependences between pixels decay as a power law of distance between the pixels.

It is therefore often sufficient to model local dependences only. —> **Convolution**

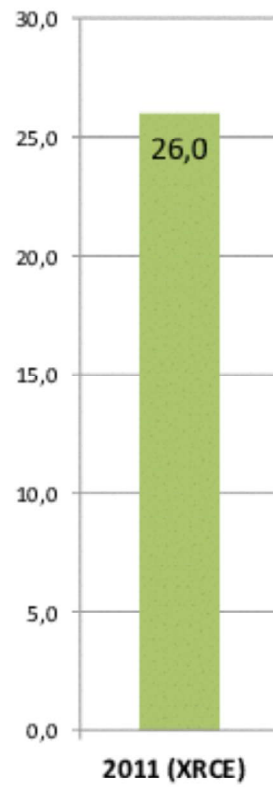
More generally, we should allocate parameters that model dependencies in proportion to the strength of those dependences. —> **Multiscale, hierarchical representations**

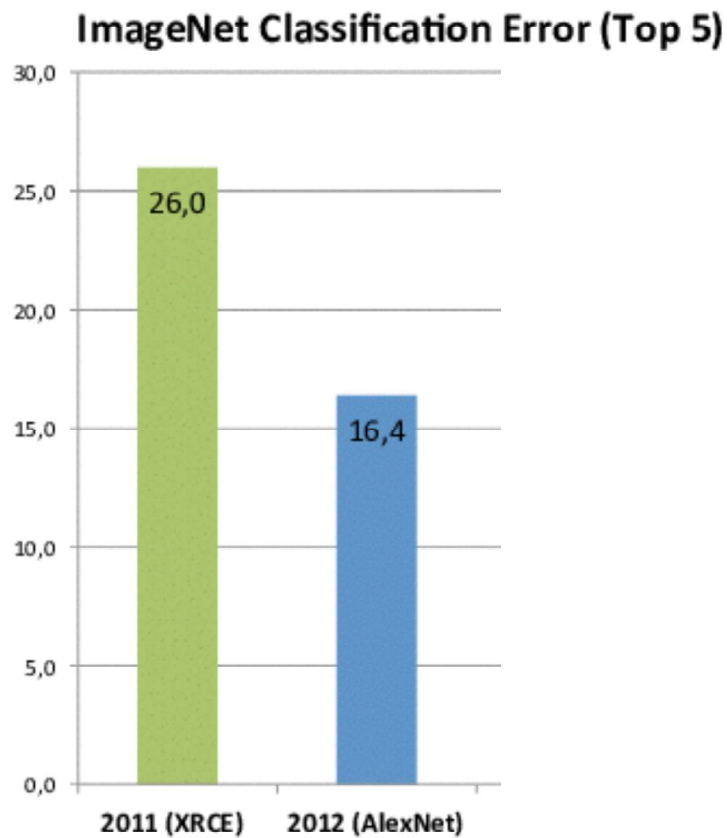
[For more discussion, see “Why does Deep and Cheap Learning Work So Well?”, Lin et al. 2017]

# Some networks

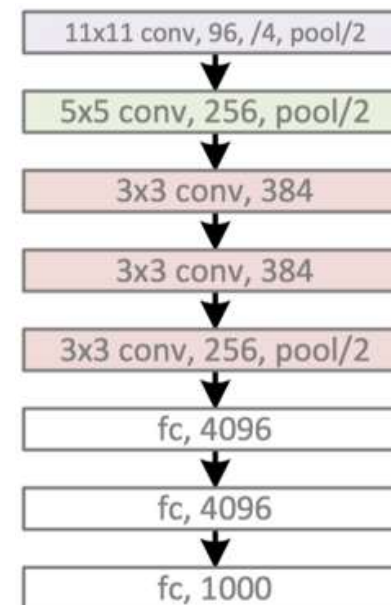
... and what makes them work

**ImageNet Classification Error (Top 5)**





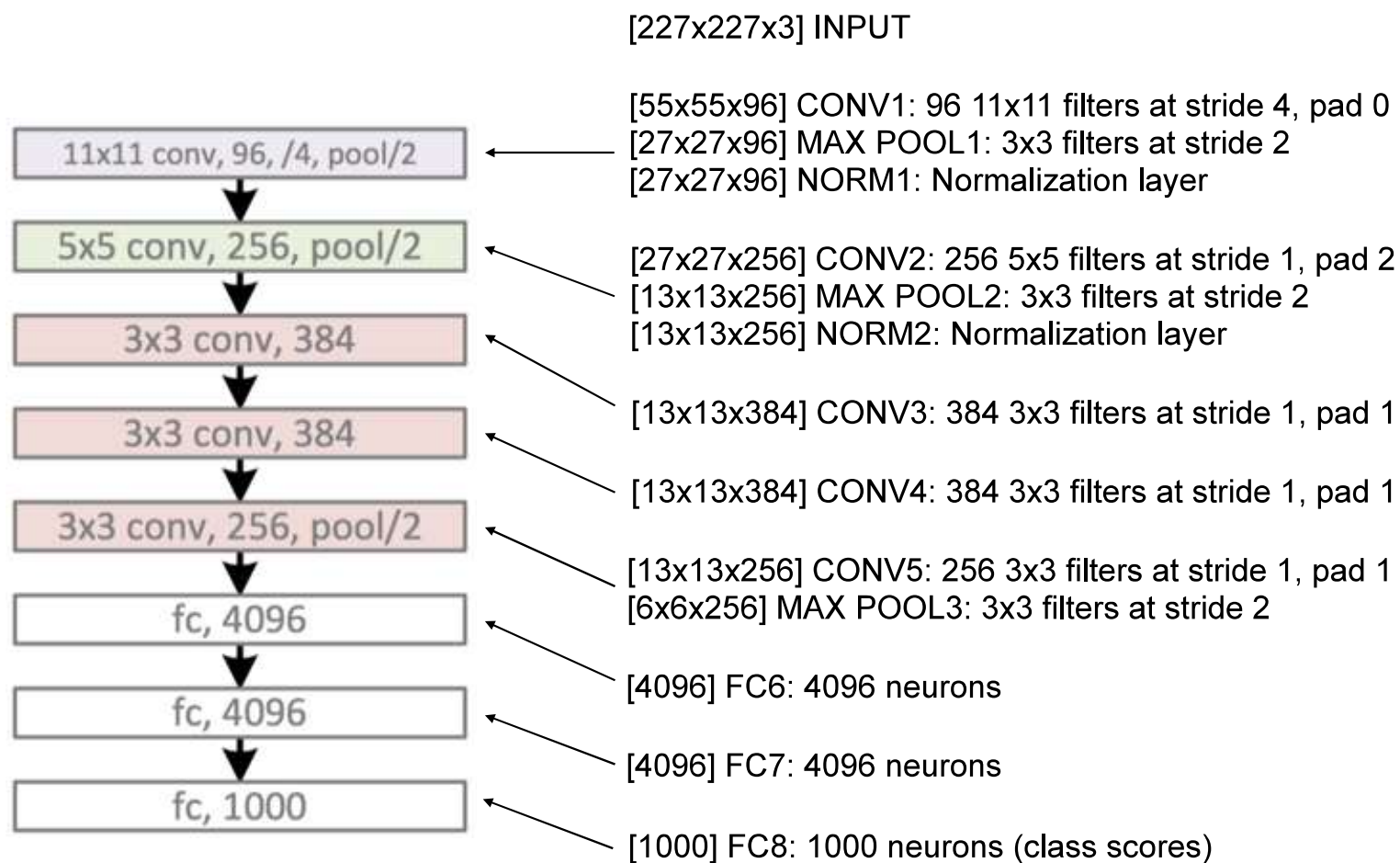
2012: AlexNet  
5 conv. layers

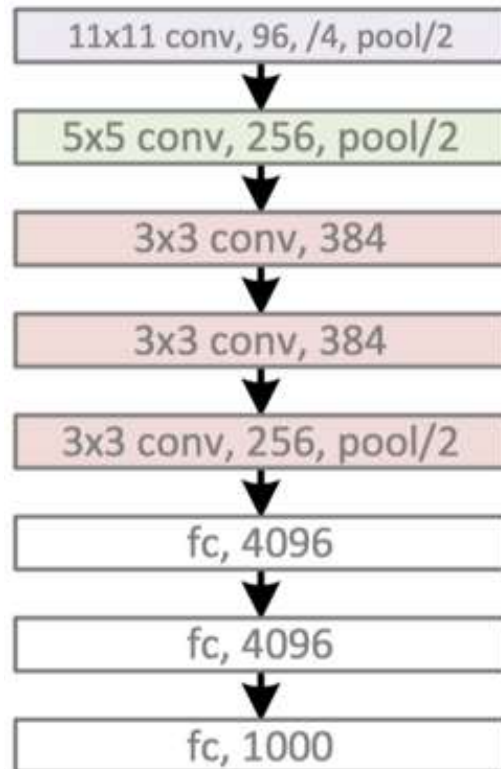


Error: 16.4%

*[Krizhevsky et al: ImageNet Classification with Deep Convolutional Neural Networks, NIPS 2012]*

# Alexnet — [Krizhevsky et al. NIPS 2012]

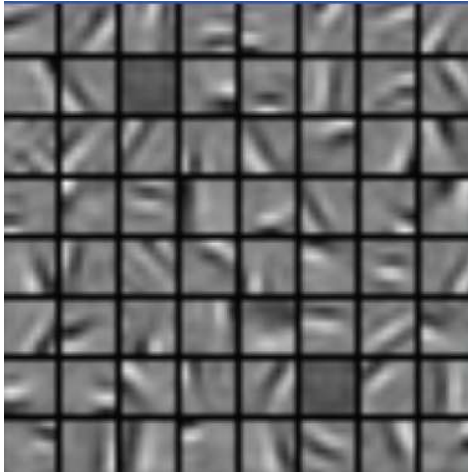




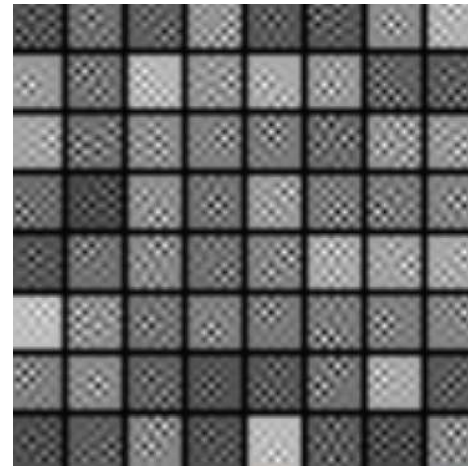
What filters are learned?

# What filters are learned?

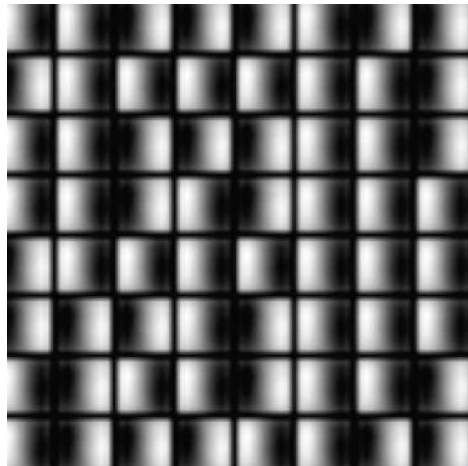
A



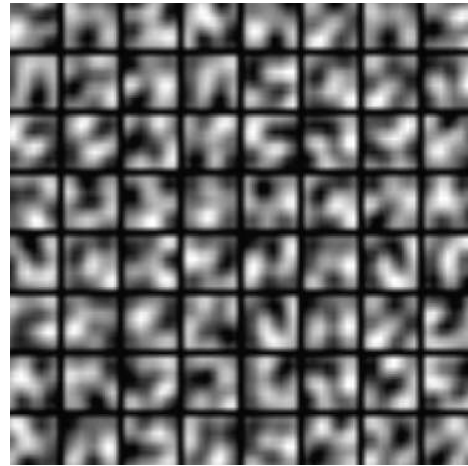
B



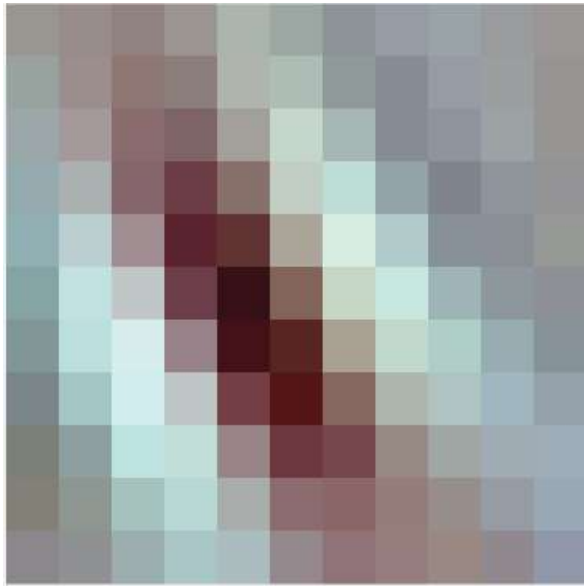
C



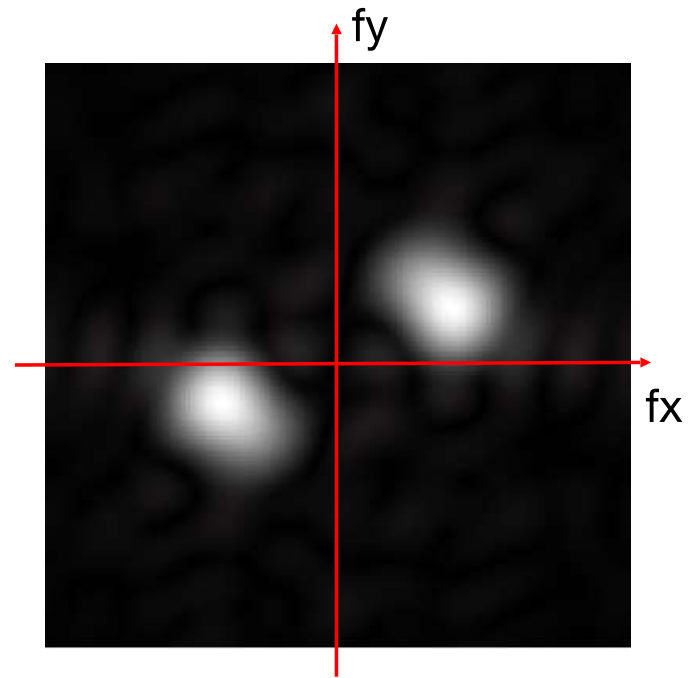
D



# Get to know your units

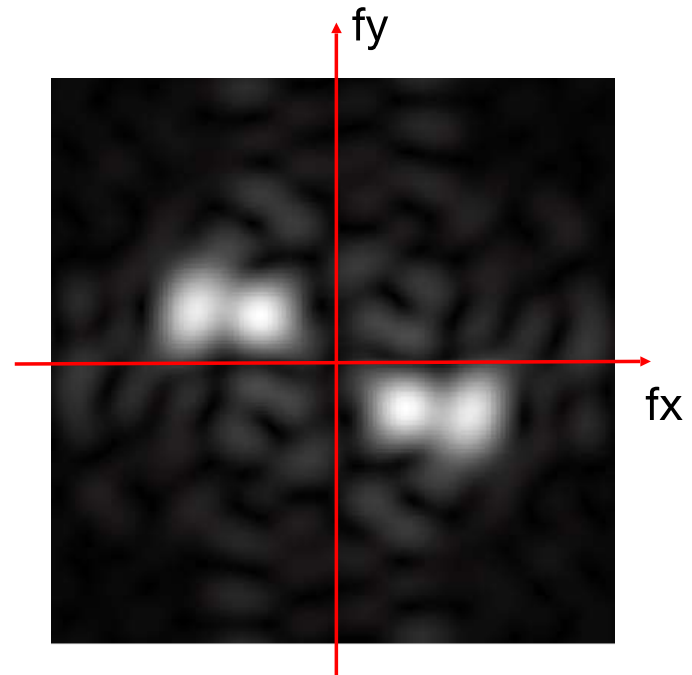
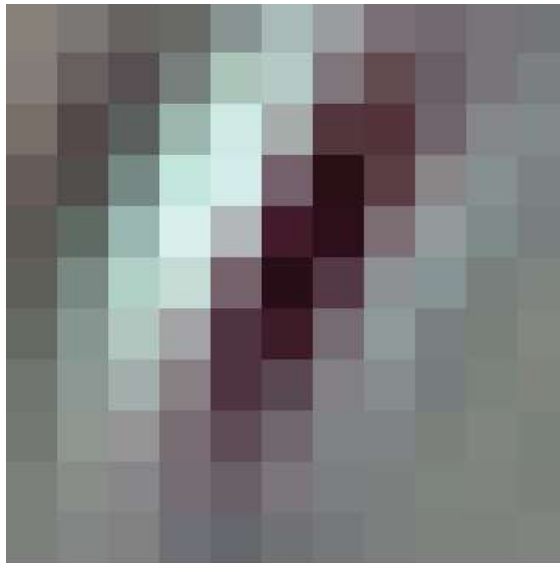


11x11 convolution kernel  
(3 color channels)

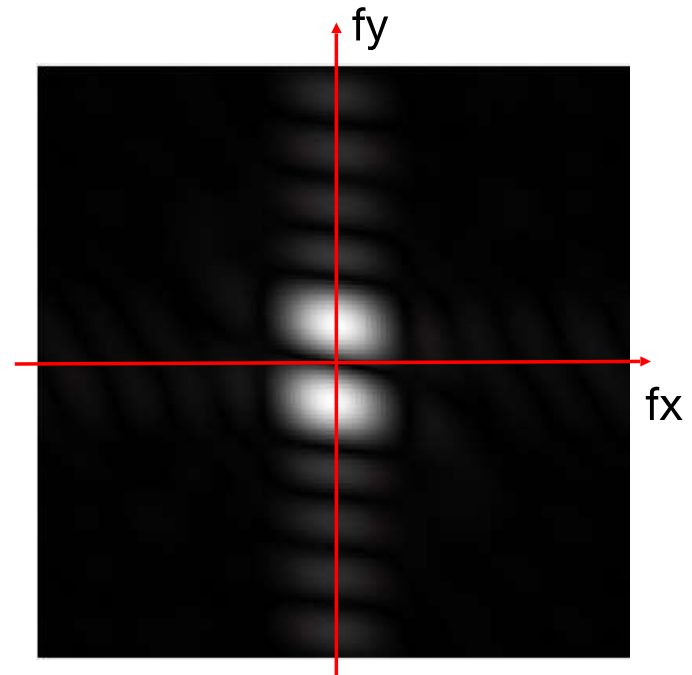
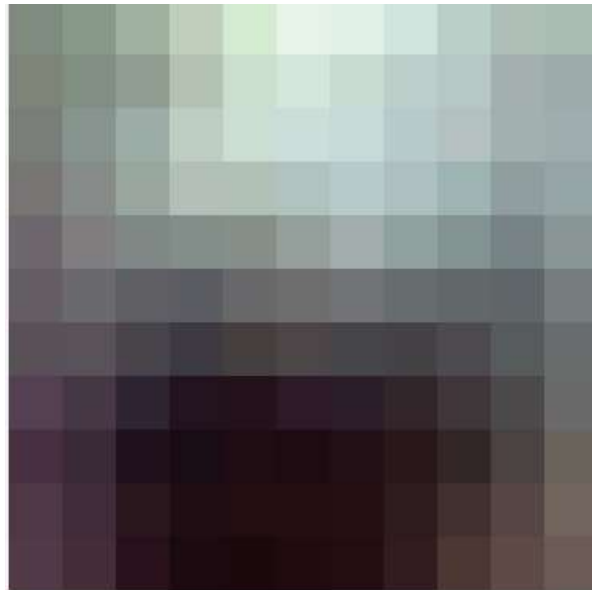




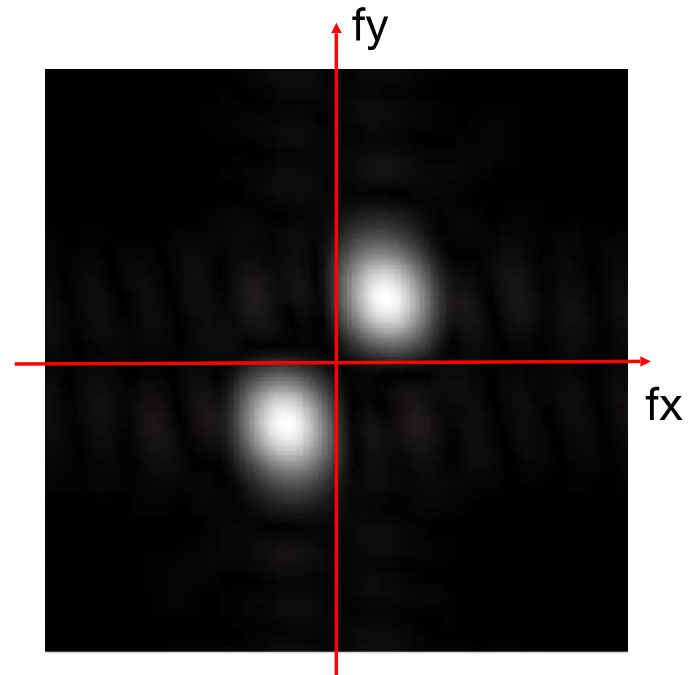
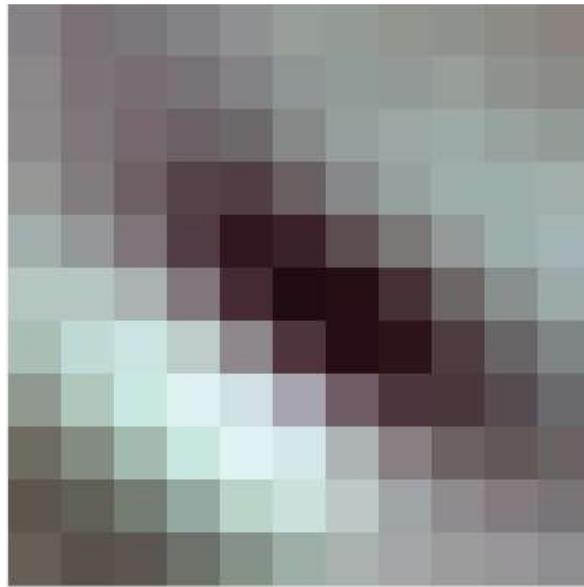
# Get to know your units



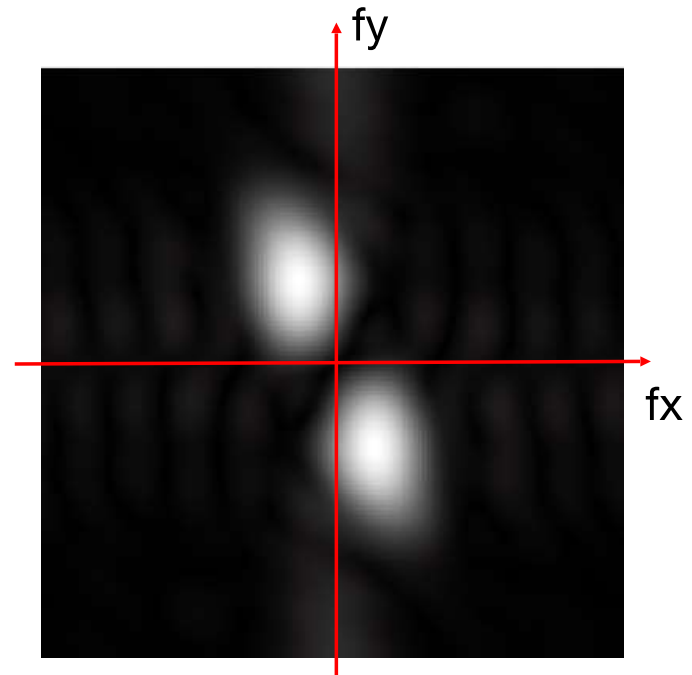
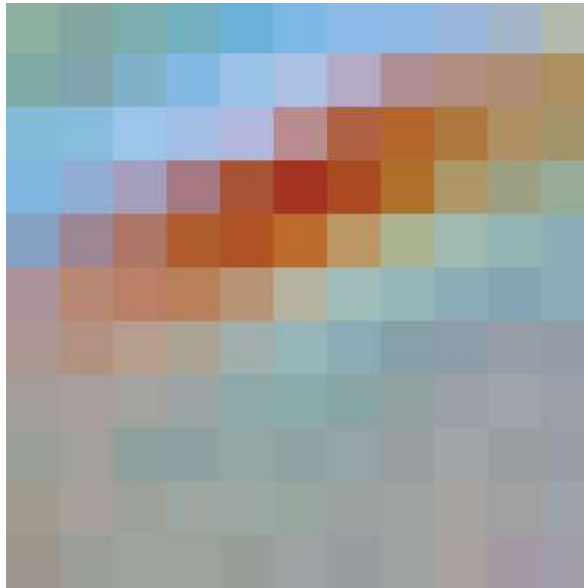
# Get to know your units



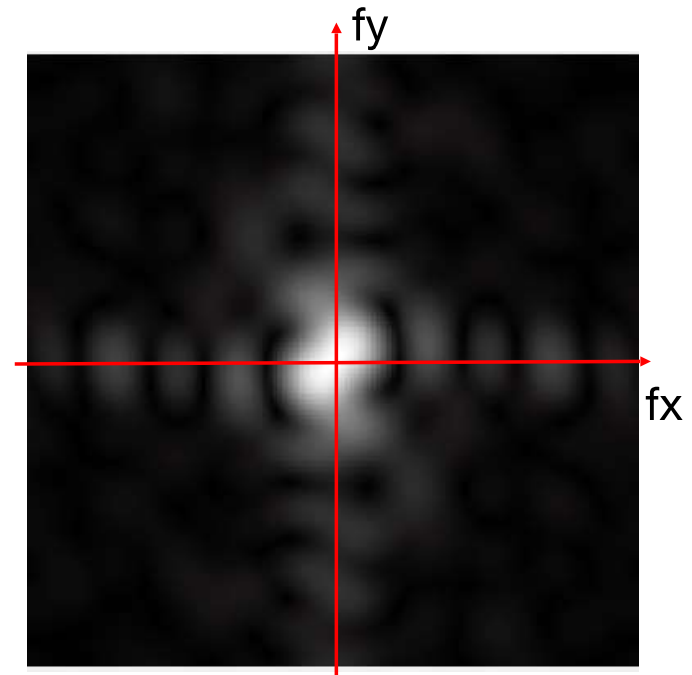
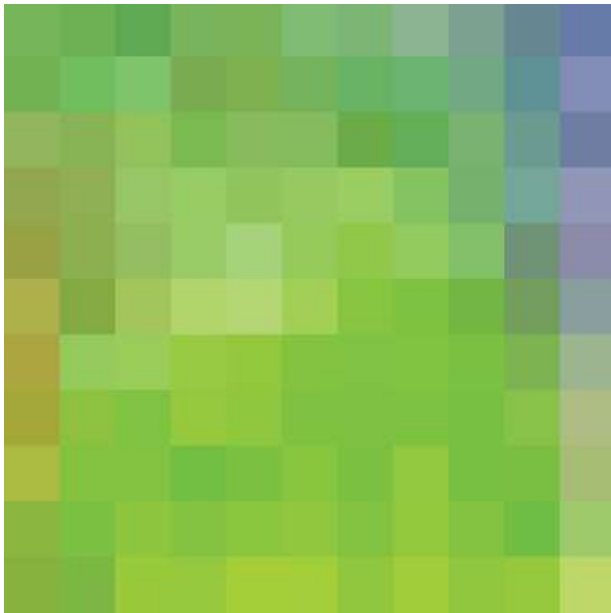
# Get to know your units



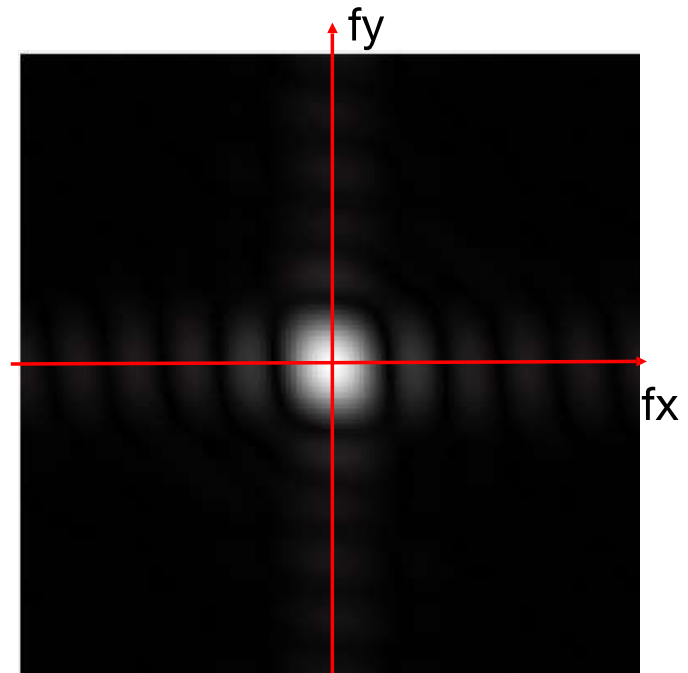
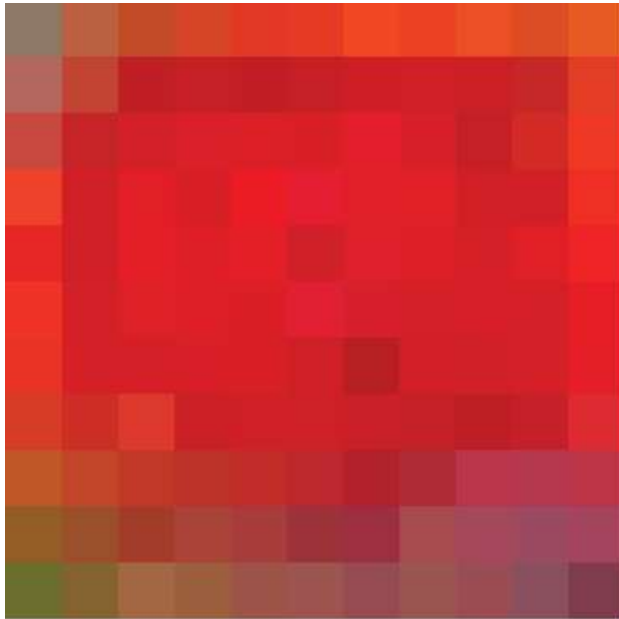
# Get to know your units



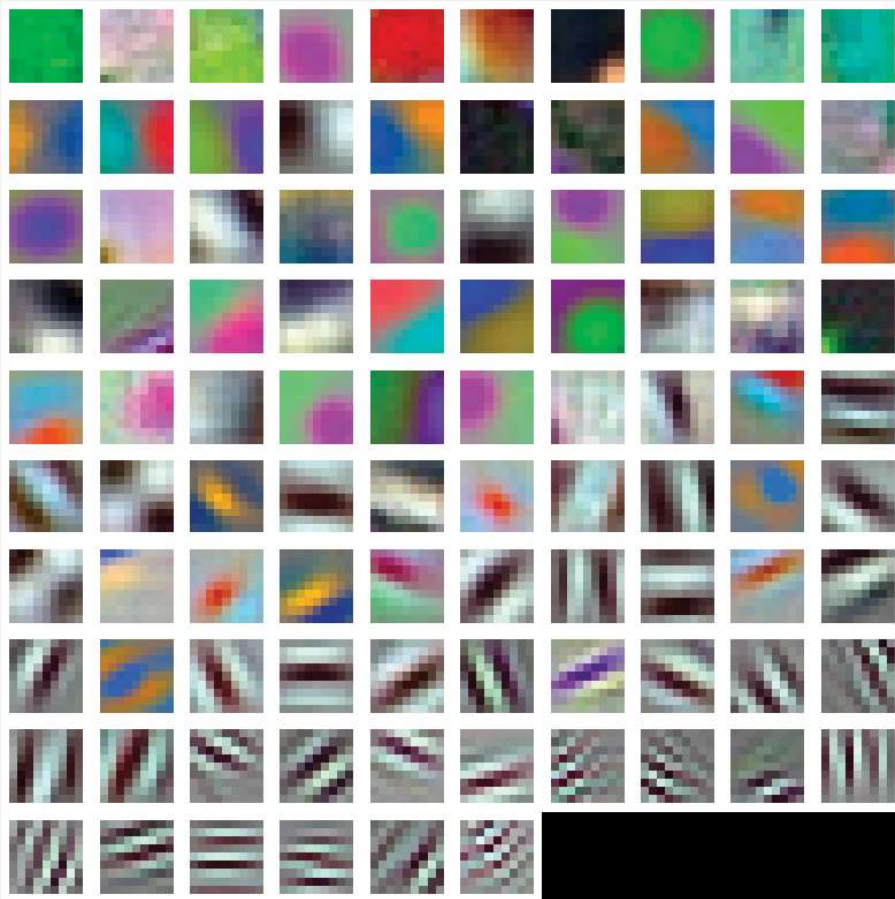
# Get to know your units



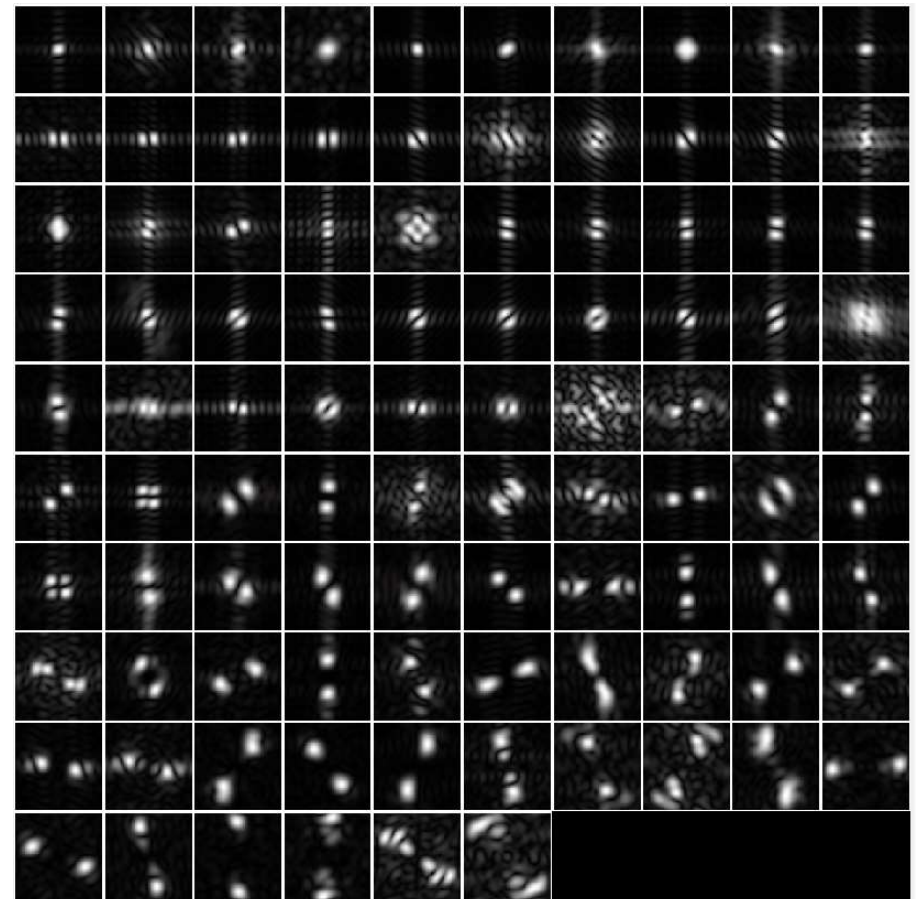
# Get to know your units



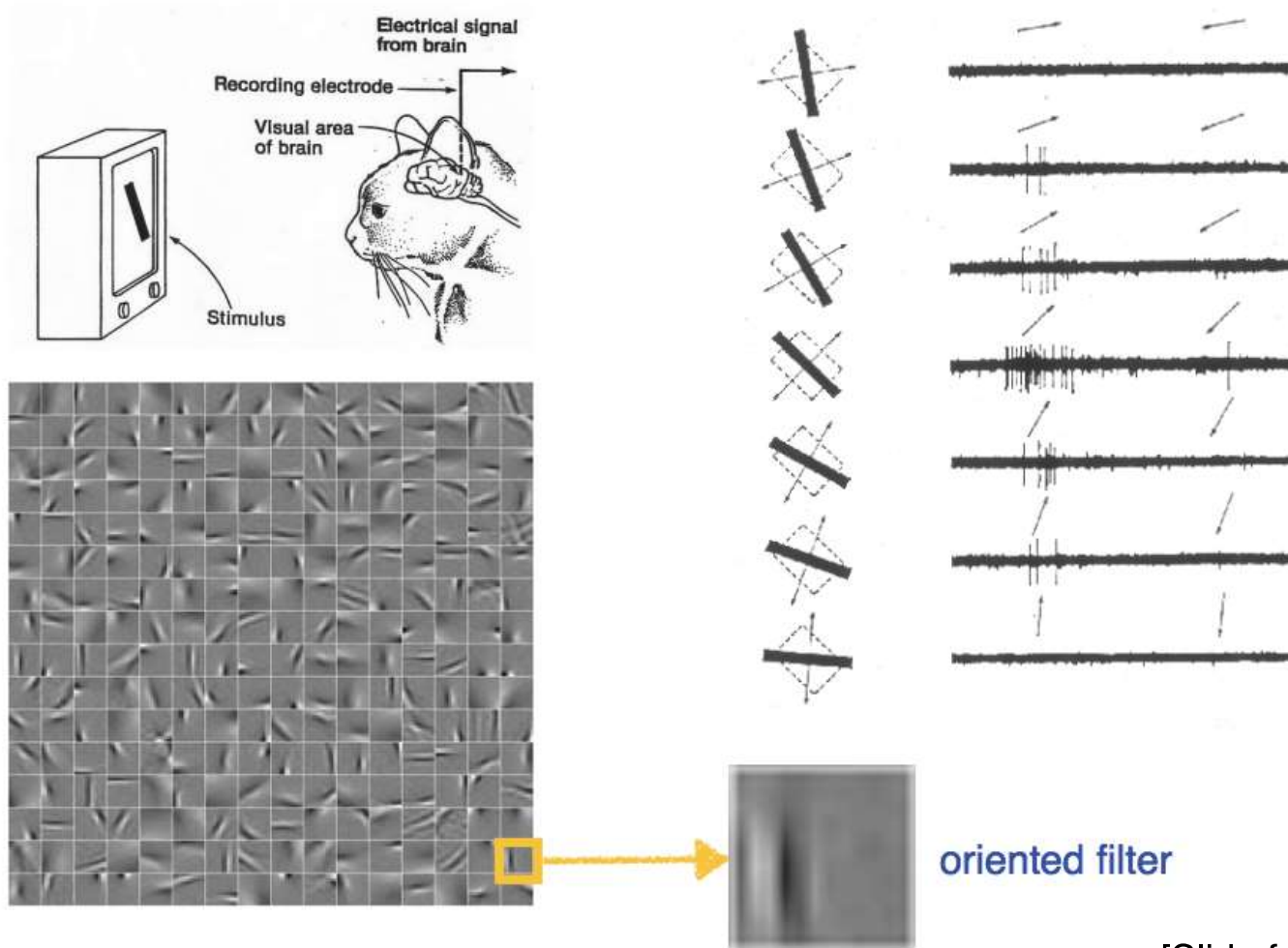
# Get to know your units



96 Units in conv1



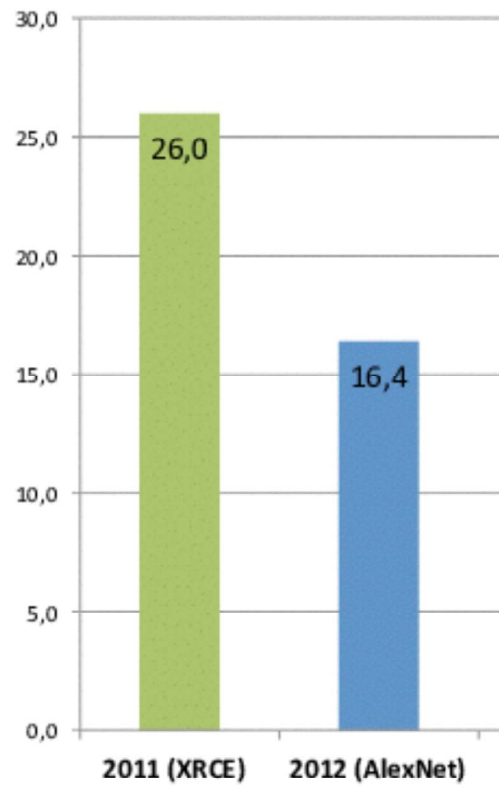
[Hubel and Wiesel 59]

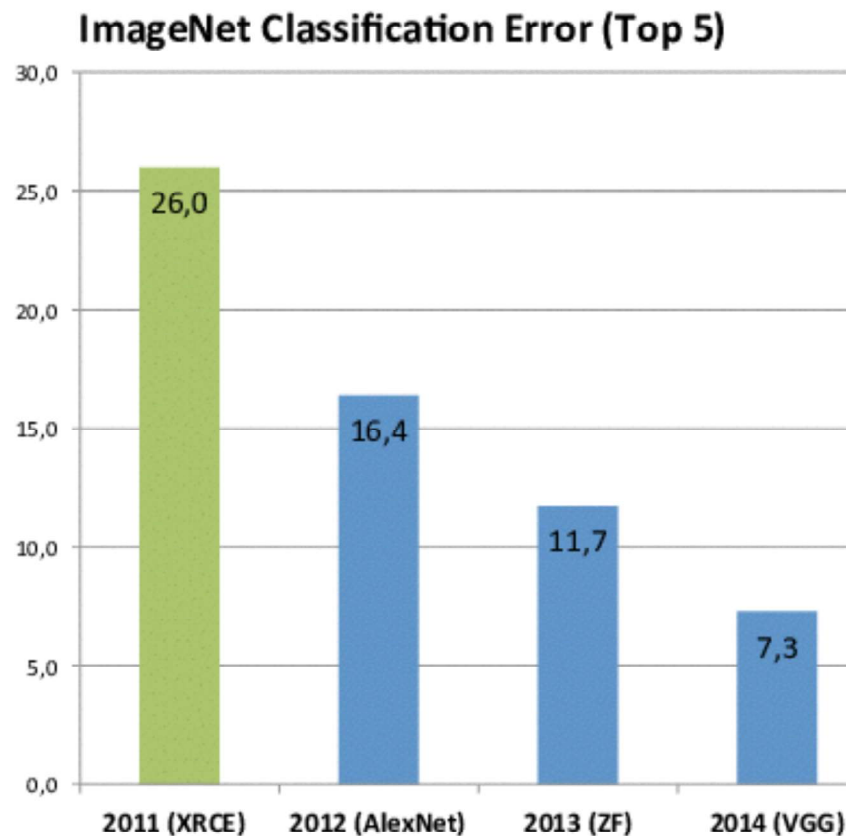


[Slide from Andrea Vedaldi]

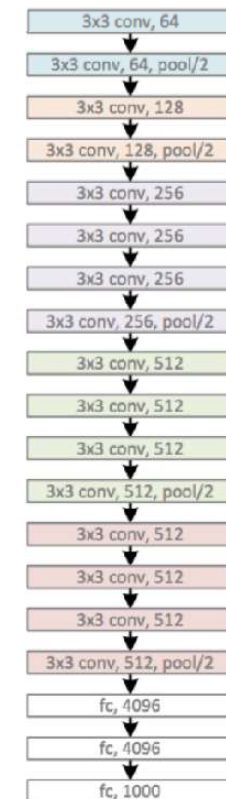


**ImageNet Classification Error (Top 5)**





2014: VGG  
16 conv. layers

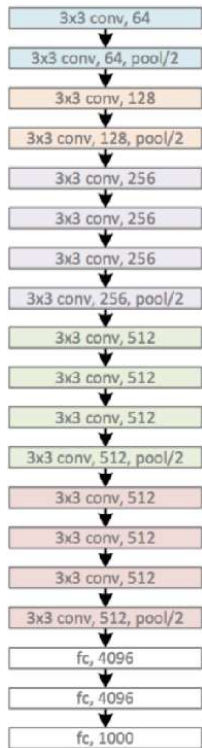


Error: 7.3%

*[Simonyan & Zisserman: Very Deep Convolutional Networks for Large-Scale Image Recognition, ICLR 2015]*

# VGG-Net [Simonyan & Zisserman, 2015]

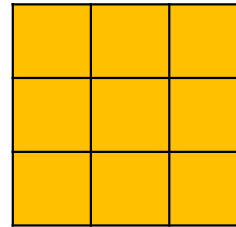
2014: VGG  
16 conv. layers



Error: 7.3%

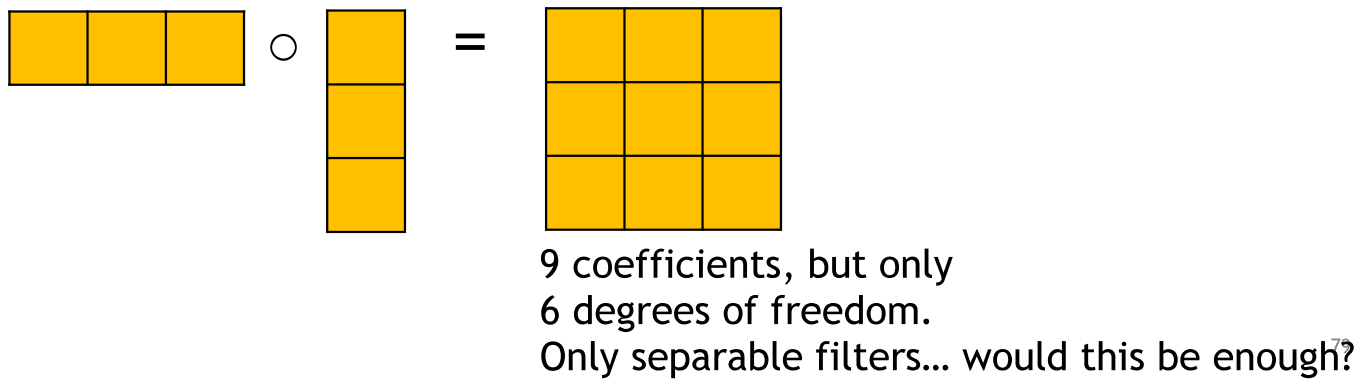
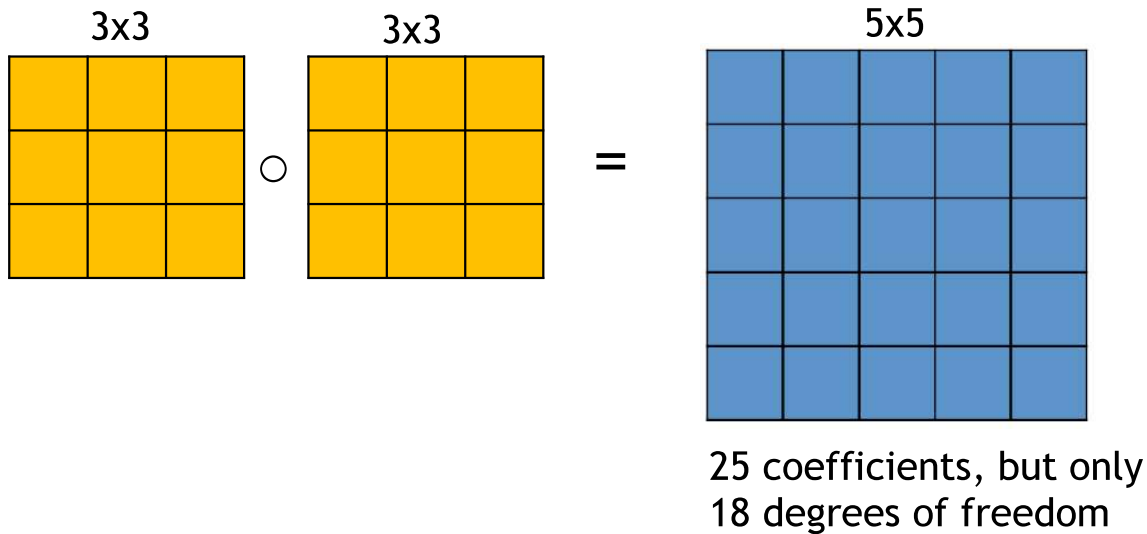
## Main developments

- Small convolutional kernels: only 3x3

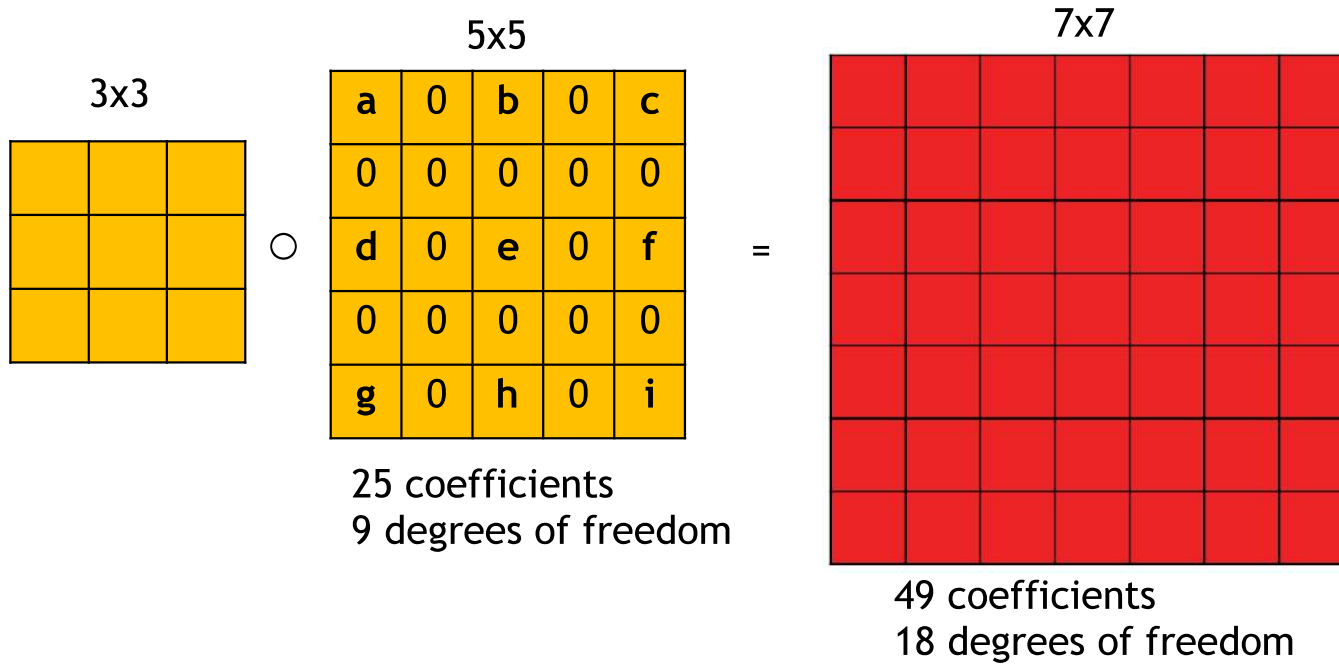


- Increased depth (5 -> 16/19 layers)

# Chaining convolutions



# Dilated convolutions



What is lost?

[<https://arxiv.org/pdf/1511.07122.pdf>]

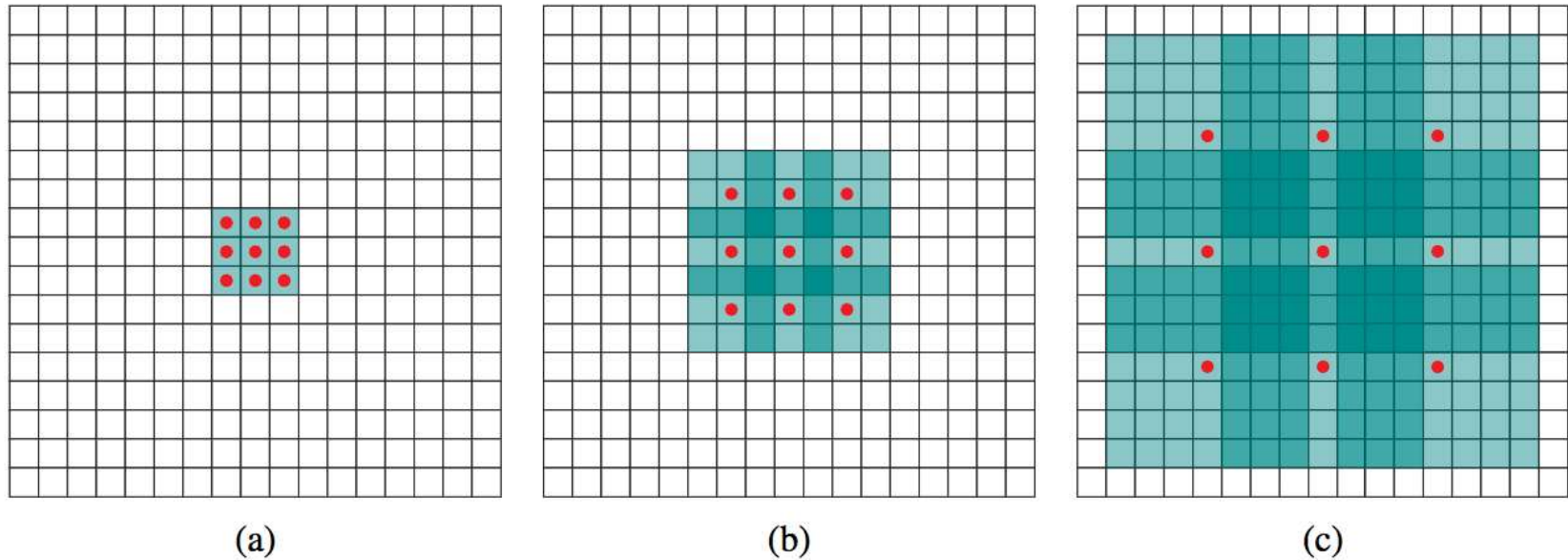
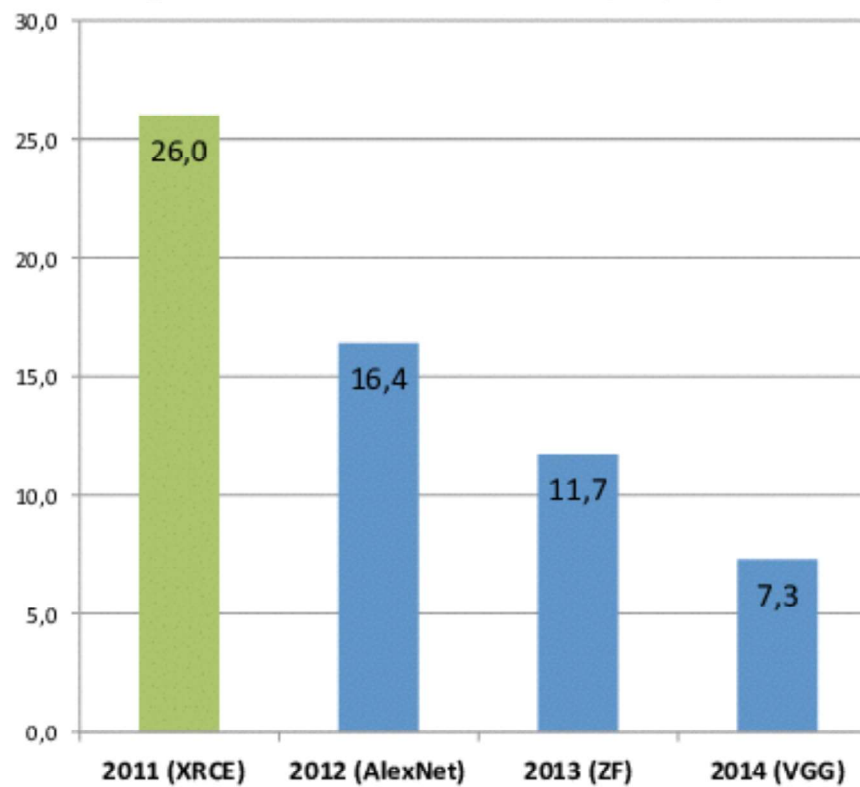
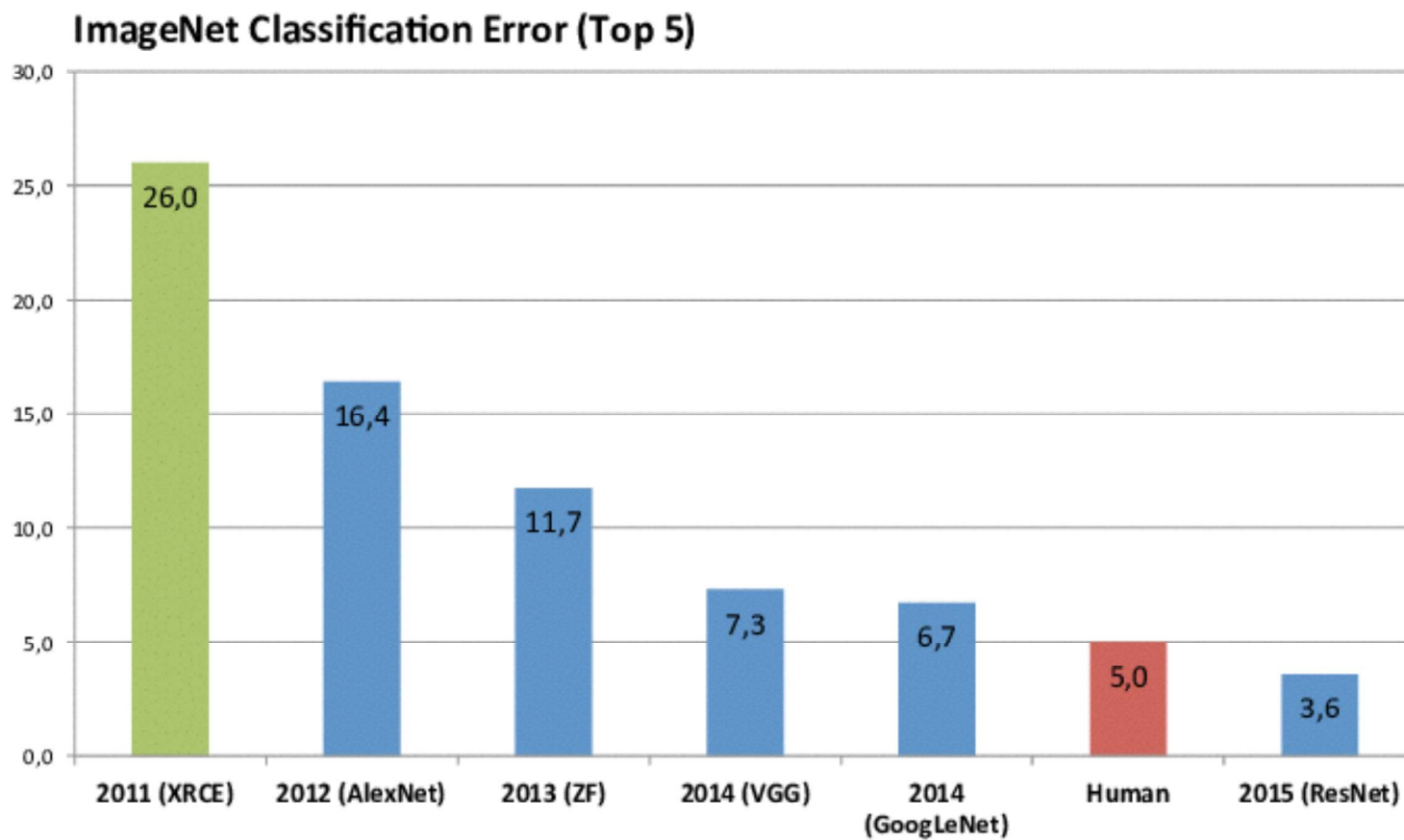


Figure 1: Systematic dilation supports exponential expansion of the receptive field without loss of resolution or coverage. (a)  $F_1$  is produced from  $F_0$  by a 1-dilated convolution; each element in  $F_1$  has a receptive field of  $3 \times 3$ . (b)  $F_2$  is produced from  $F_1$  by a 2-dilated convolution; each element in  $F_2$  has a receptive field of  $7 \times 7$ . (c)  $F_3$  is produced from  $F_2$  by a 4-dilated convolution; each element in  $F_2$  has a receptive field of  $15 \times 15$ . The number of parameters associated with each layer is identical. The receptive field grows exponentially while the number of parameters grows linearly.

**ImageNet Classification Error (Top 5)**





Error: 3.6%

2016: ResNet  
>100 conv. layers

*[He et al: Deep Residual Learning for Image Recognition, CVPR 2016]*

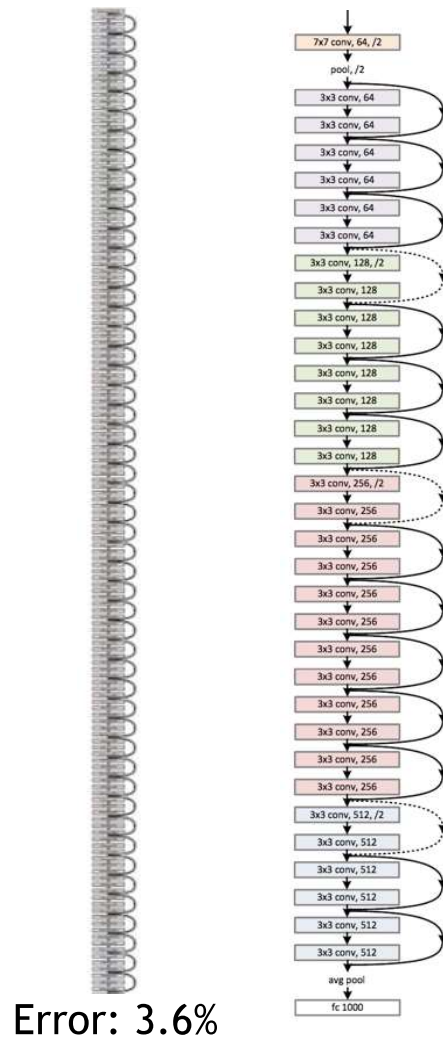


2016: ResNet  
>100 conv. layers

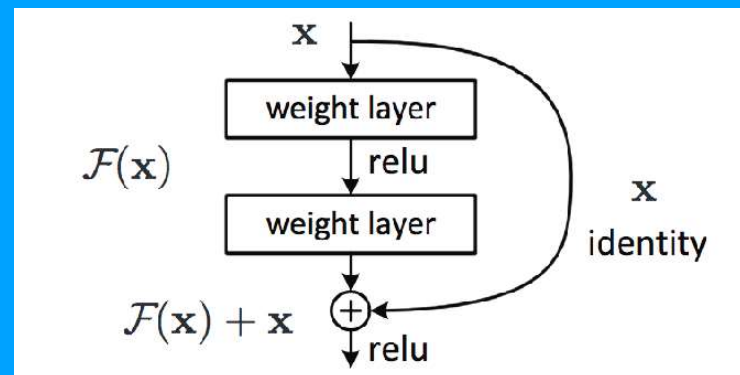
# ResNet [He et al, 2016]

## Main developments

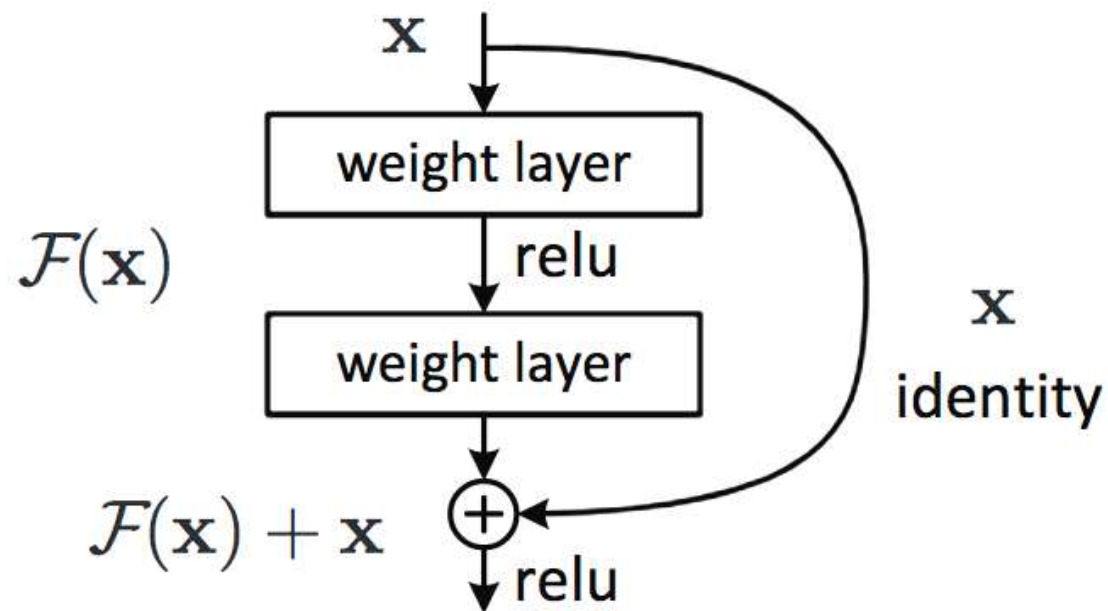
- Increased depth possible through residual blocks



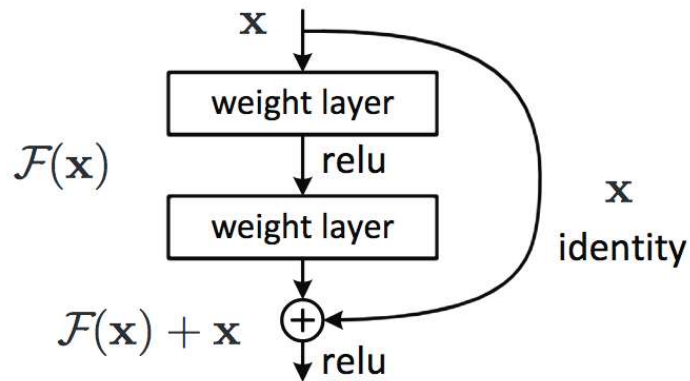
Error: 3.6%



# Residual Blocks



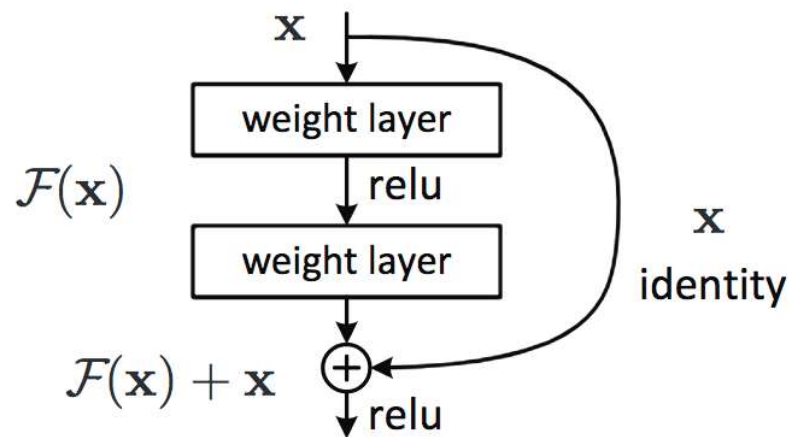
# Residual Blocks



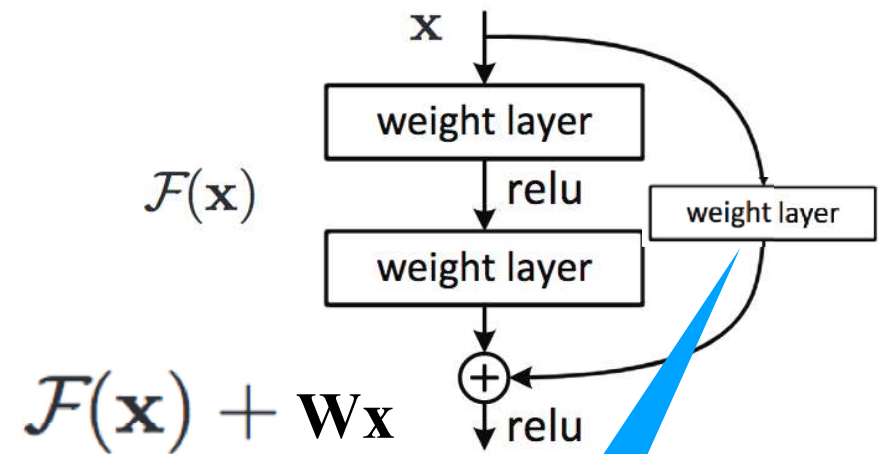
## Why do they work?

- Gradients can propagate faster (via the identity mapping)
- Within each block, only small residuals have to be learned

If output has same size as input:



If output has a different size:

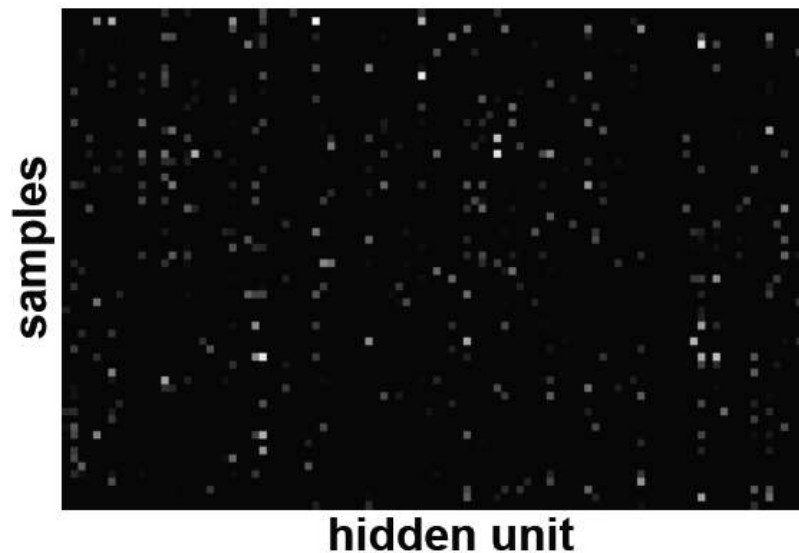


Projects into the right  
dimensionality:  
 $\dim(\mathcal{F}(\mathbf{x})) = \dim(\mathbf{W}\mathbf{x})$

Some debugging advice

# Other good things to know

- Check gradients numerically by finite differences
- Visualize hidden activations — should be uncorrelated and high variance

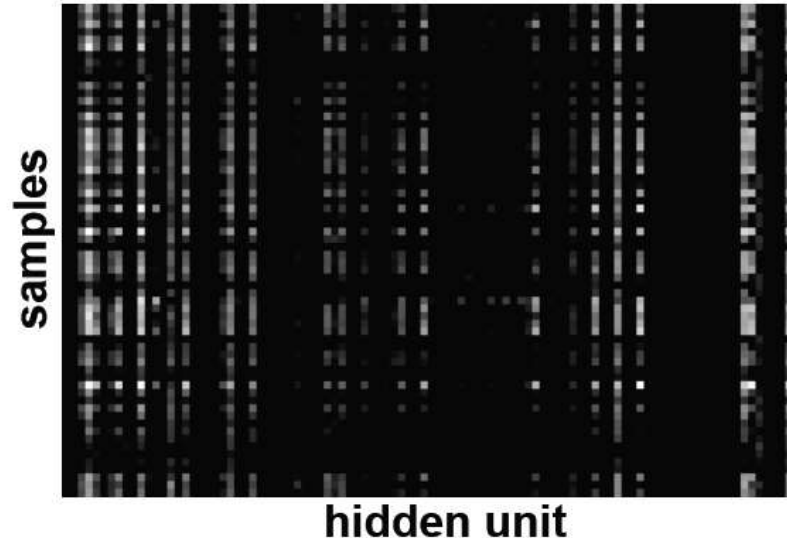


**Good training:** hidden units are sparse across samples and across features.

[Derived from slide by Marc'Aurelio Ranzato]

# Other good things to know

- Check gradients numerically by finite differences
- Visualize hidden activations — should be uncorrelated and high variance

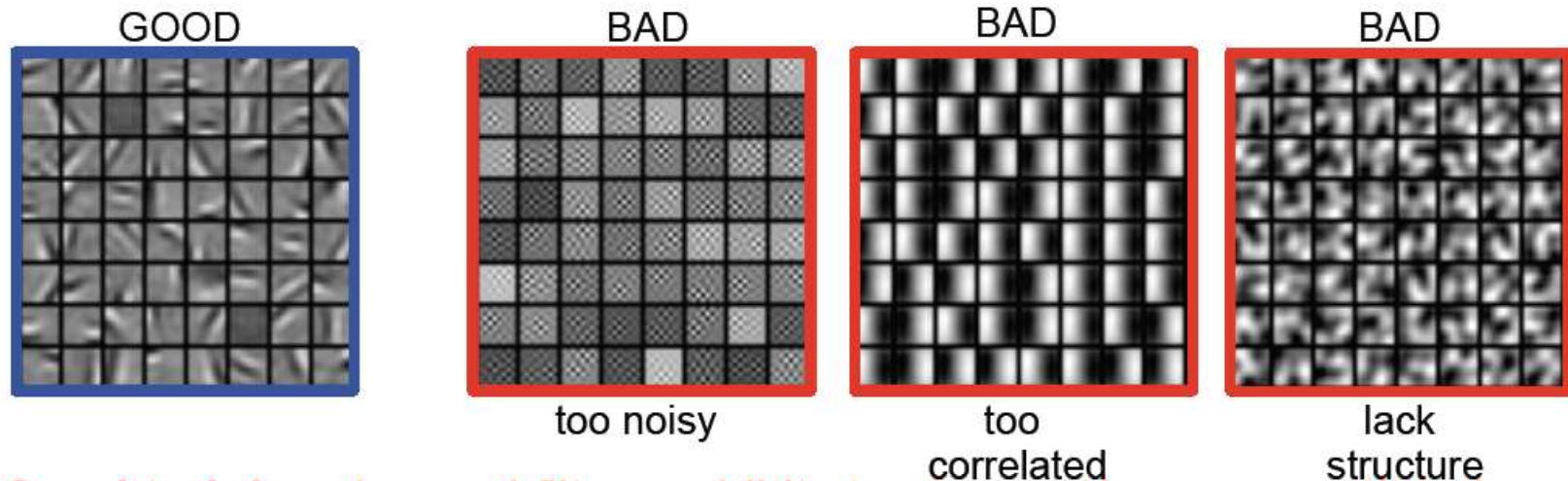


**Bad training:** many hidden units ignore the input and/or exhibit strong correlations.

[Derived from slide by Marc'Aurelio Ranzato]

# Other good things to know

- Check gradients numerically by finite differences
- Visualize hidden activations — should be uncorrelated and high variance
- Visualize filters



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

[Derived from slide by Marc'Aurelio Ranzato]