



UNIVERSITAT DE  
BARCELONA



DATA SCIENCE @ UB

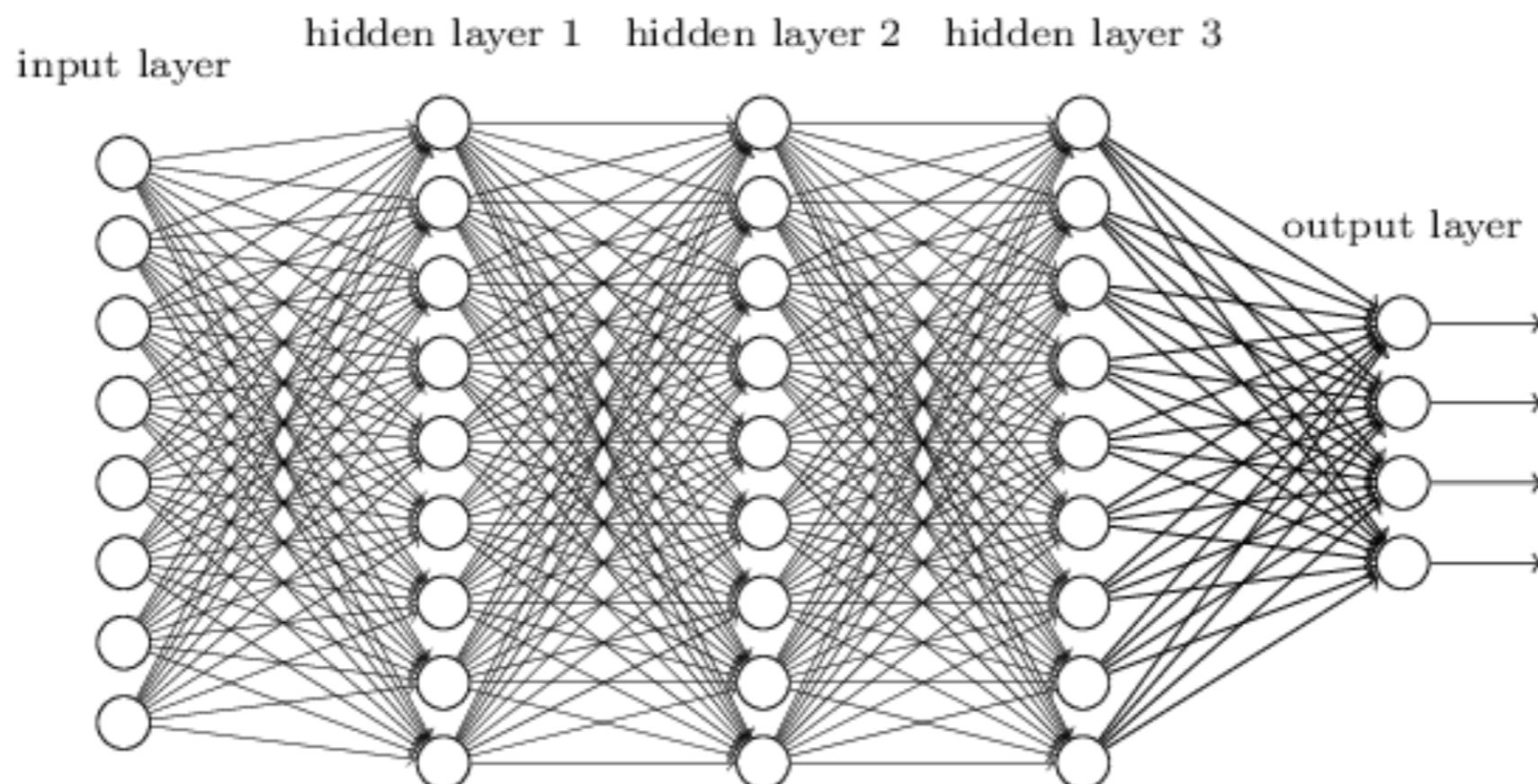


Deep Learning From Scratch  
**Convolutional Neural Networks**

Santi Seguí

# Neural Networks for Images

## Multi Layer Perceptron



How many parameter?

$$8 * 9 + 9 * 9 + 9 * 9 + 9 * 4 = 570$$

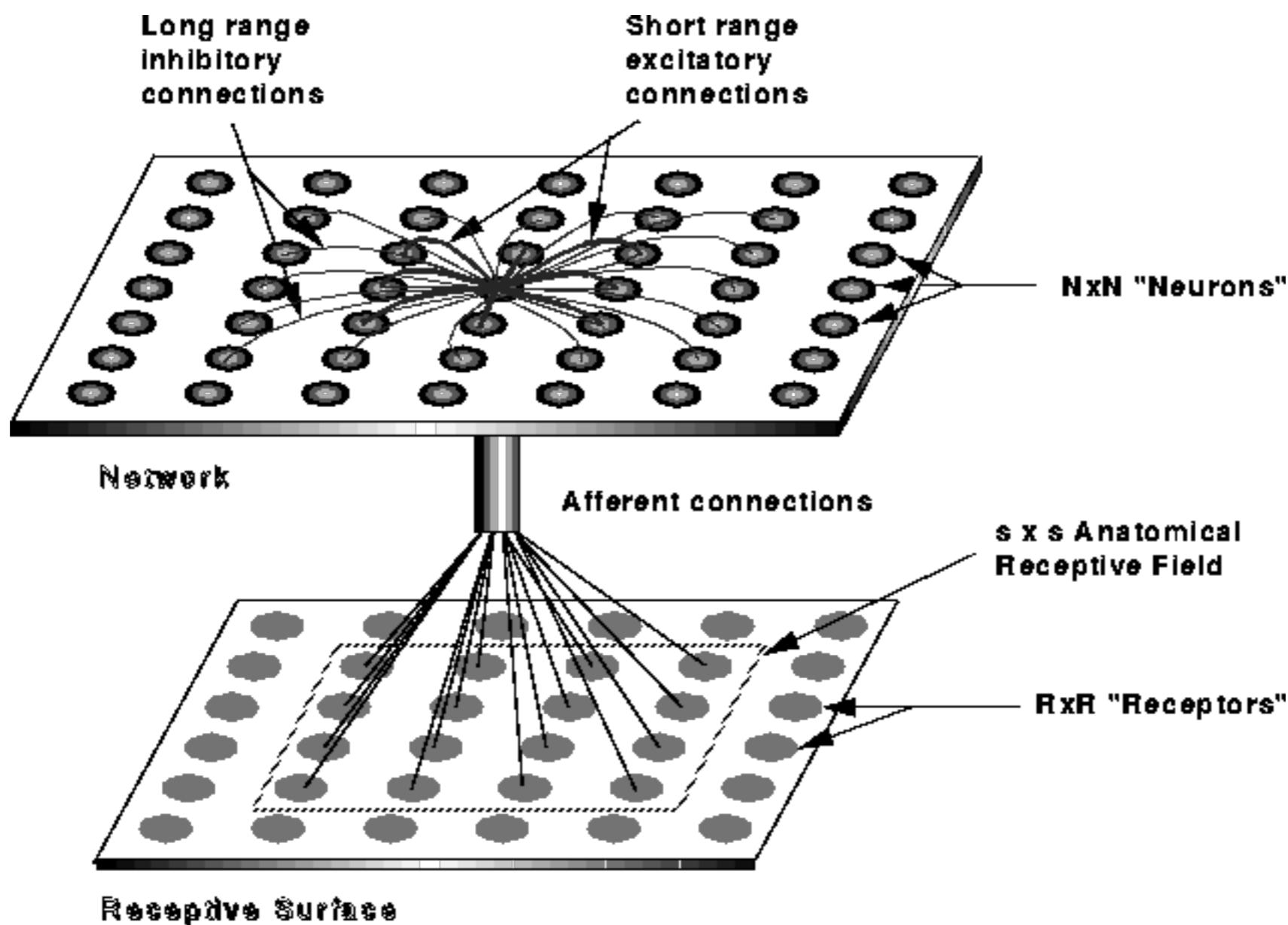
# Neural Networks for Images



# Neural Networks for Images

- An image is just a Matrix of  $N \times M$  pixels.
- So, input is a standard vector of size  $N \times M$ 
  - Imagine an medium resolution color image of 256x256 pixels
  - If we think on a Multi Layer Perceptron with just one hidden layer of 256 neurons + an output layer of 1 neuron it will have more than **48 million** parameters.
  - **Does it make sense? Can we do it better?**

# Local Receptive Fields



But, in an image:

A dog can appear **anywhere** in the image!



**Doesn't matter where** they appear,  
**they look similar anywhere!**

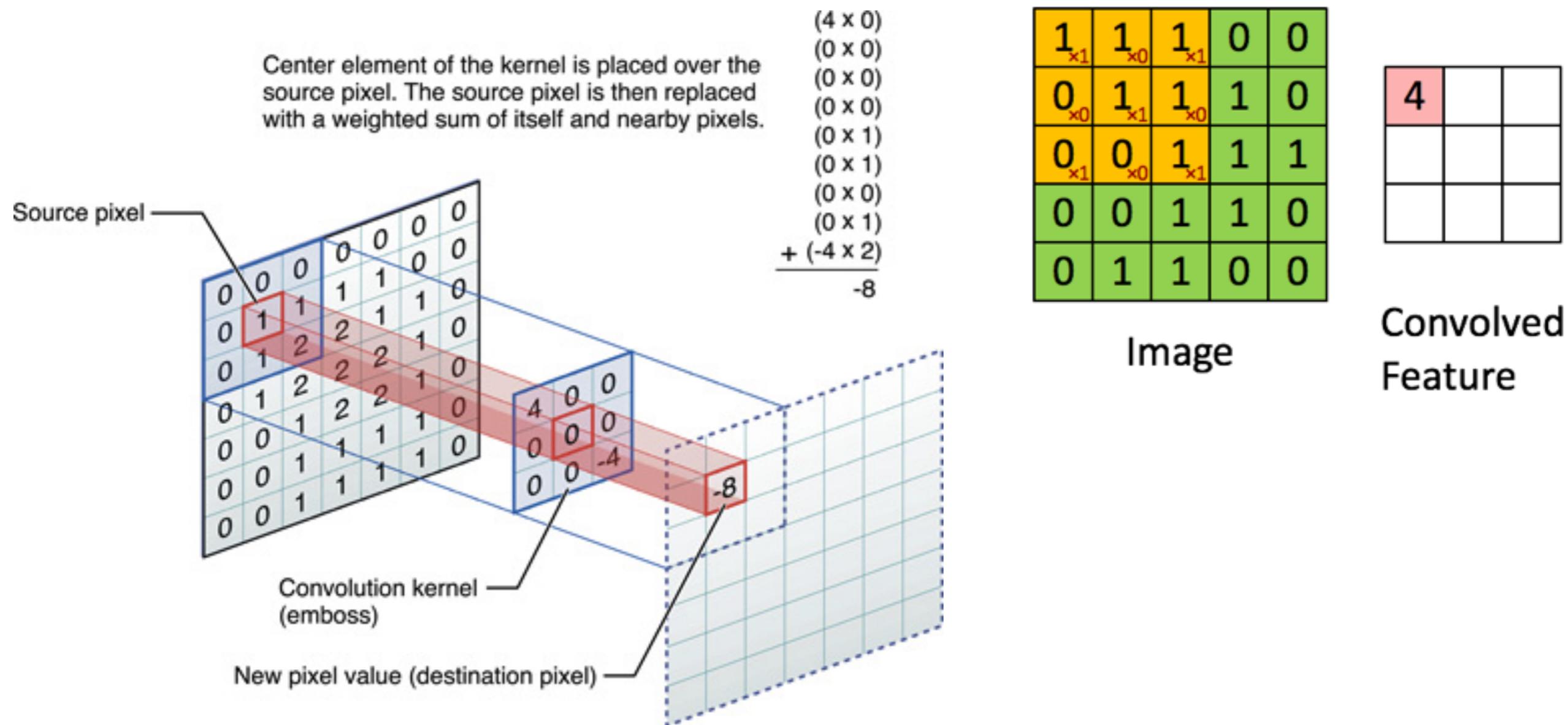
# Convolutional Neural Networks (CNNs)

- Three main ideas:
  1. **local receptive fields,**
  2. **shared weights,**
  3. **sub-sampling**

# Convolutional Neural Networks (CNNs)

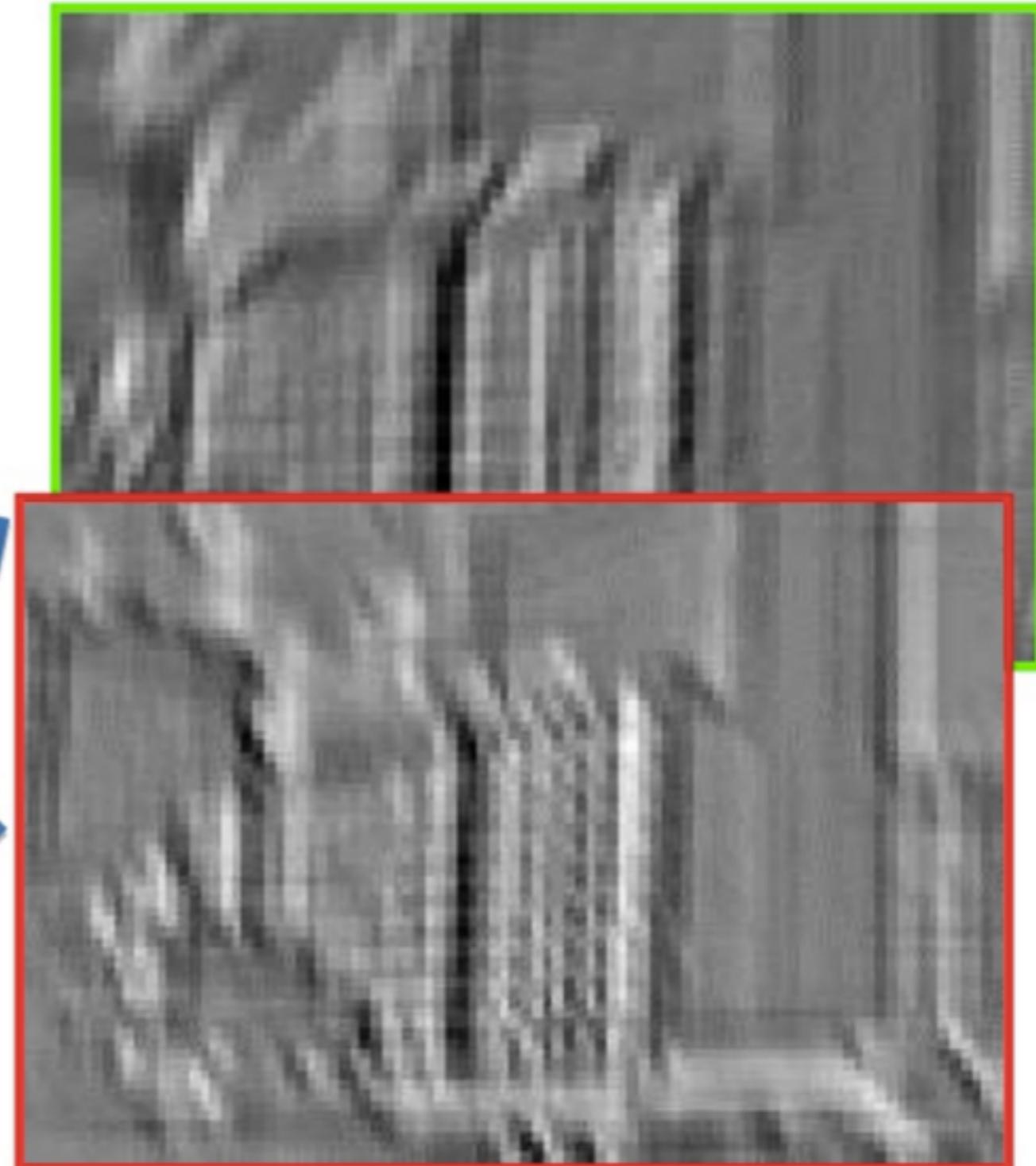
- **Repetitive** blocks of neurons that are **applied across space** (for images) or time (for audio signals etc).
- For **images**, these blocks of neurons can be interpreted as **2D convolutional kernels**, repeatedly applied over each patch of the image.
- For **speech**, they can be seen as the **1D convolutional kernels** applied across time-windows.
- At training time, the **weights** for these repeated blocks are '**shared**', i.e. the weight gradients learned over various image patches are averaged.

# What is an image convolution?

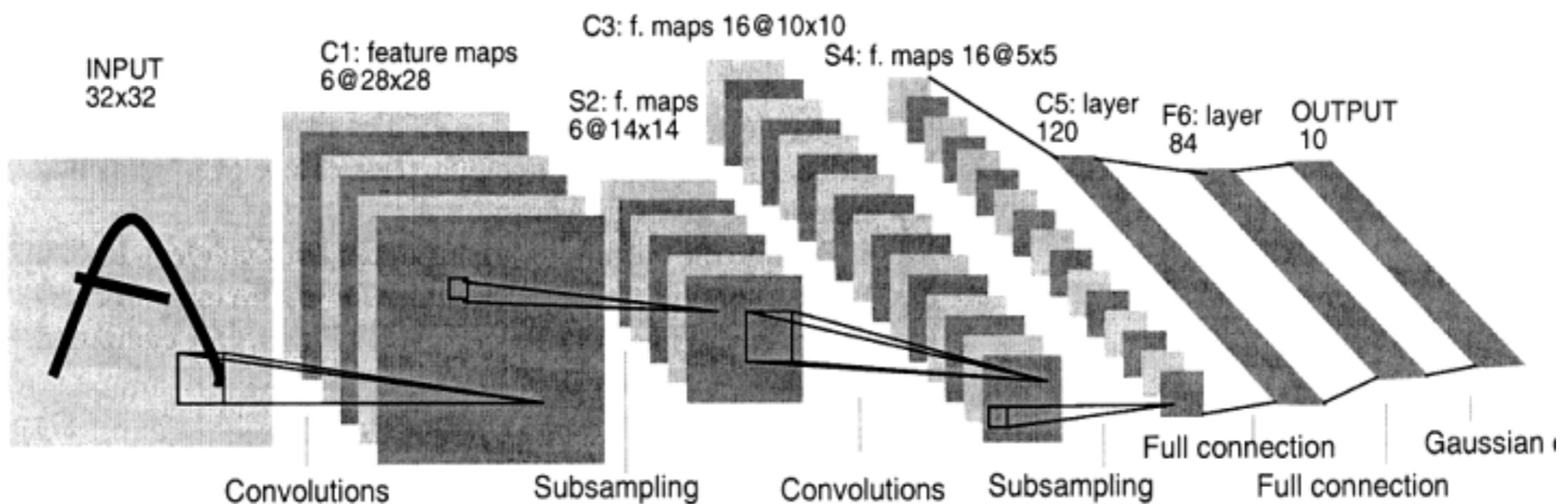


# What is an image convolution?

Weighted moving sum



# “Nothing New”

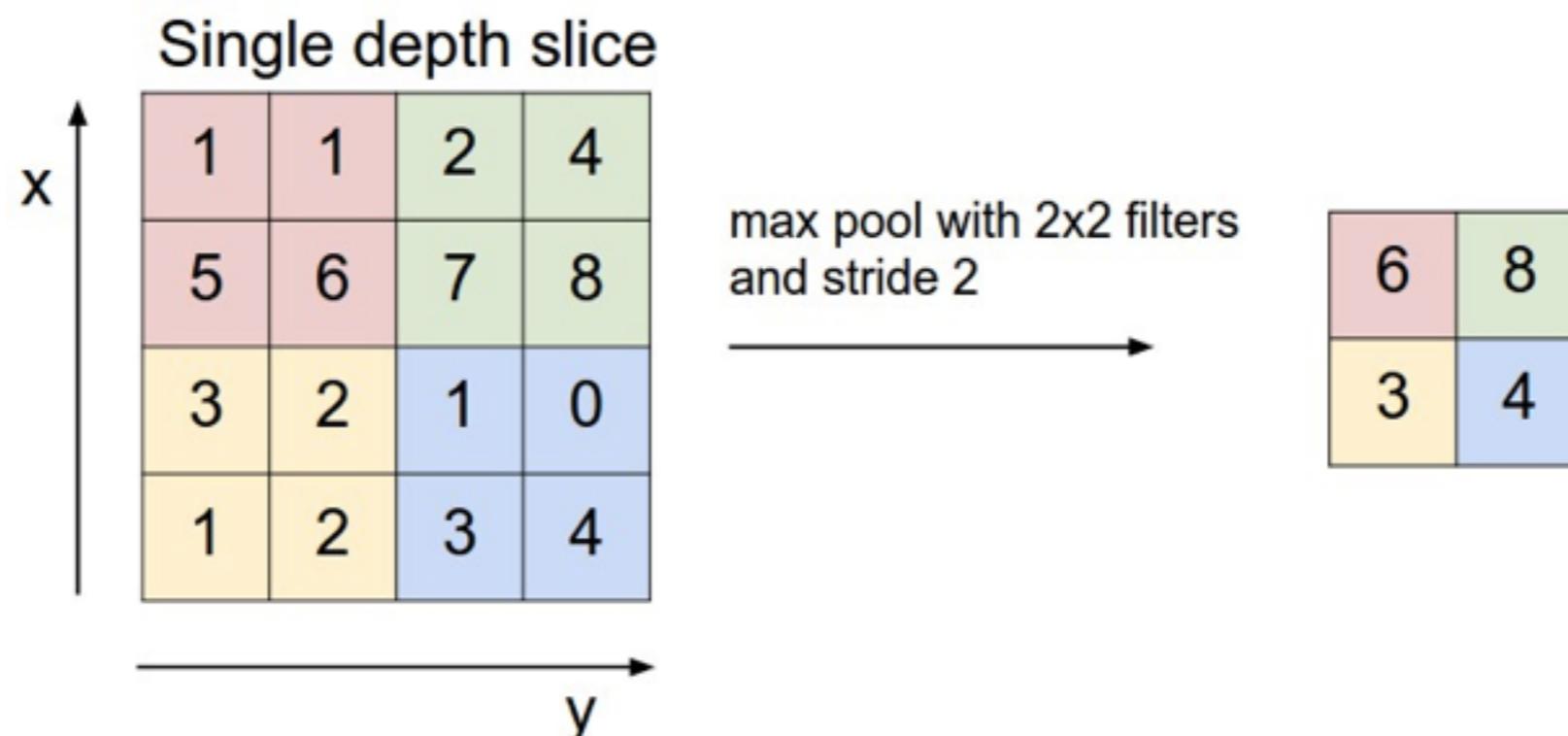


LeCun et al. 1992

# Max pooling

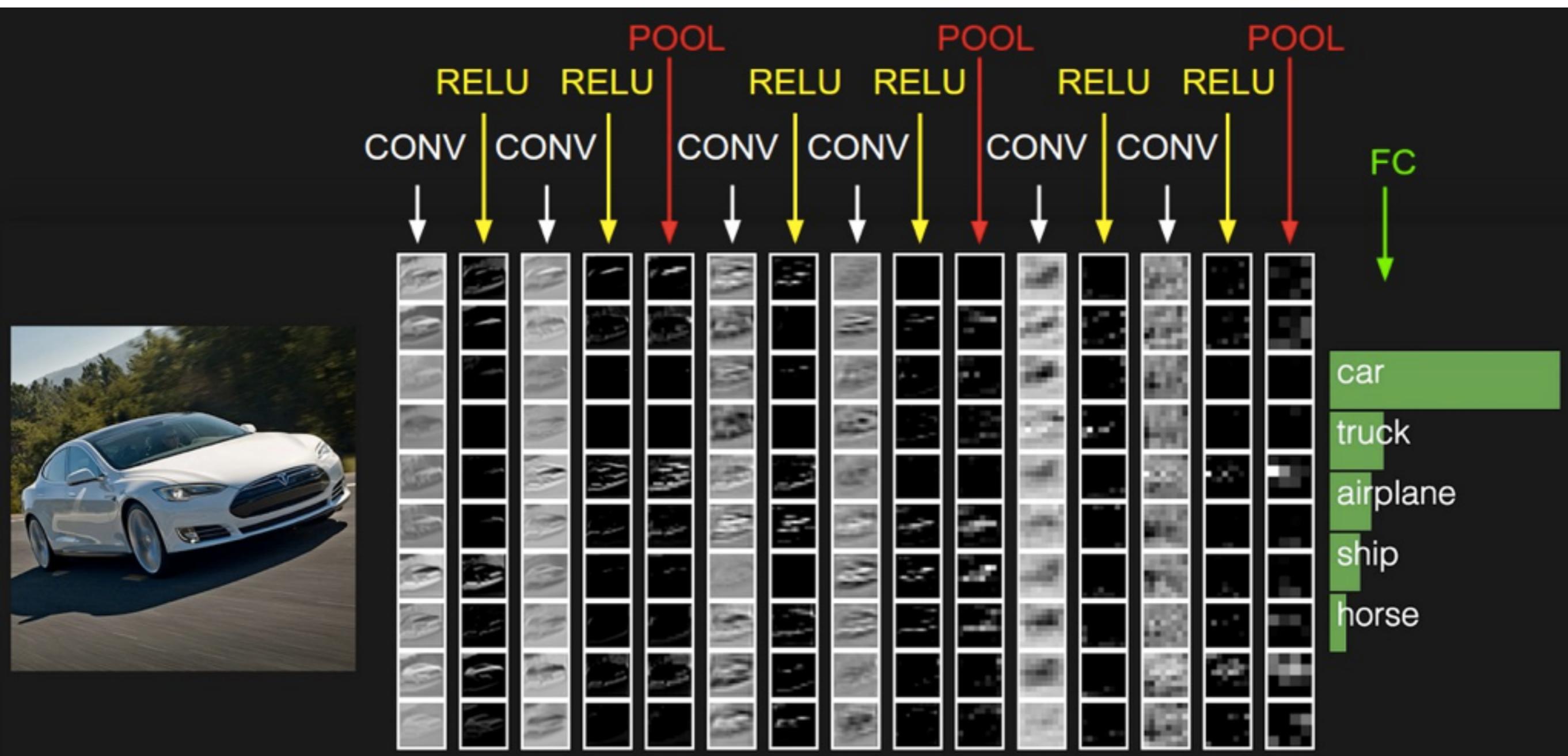
Pooling is a way of sub-sampling, i.e. reducing the dimension of the input (or at some hidden layer).

It is usually done after some of the convolutional layers



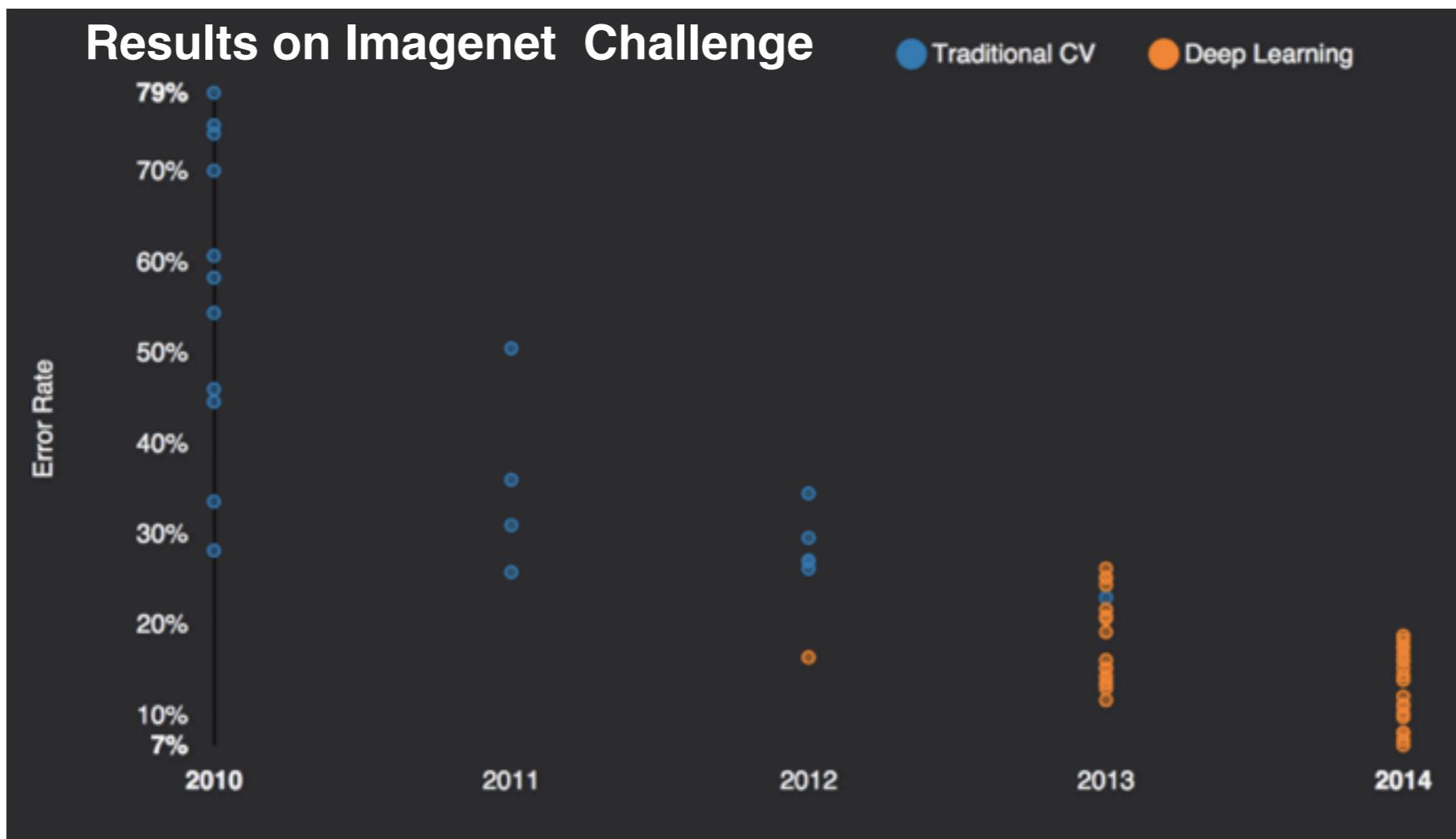
But it is also useful since it provides a form of translation **invariance**

# Finally..



# Convolutional Neural Networks (CNNs)

In computer Vision the breakthrough resulted in 2011 when Ciresan et.al introduced an algorithm to train these networks by using graphical cards (GPUs)



# AlexNet

Similar framework to LeCun'98 but:

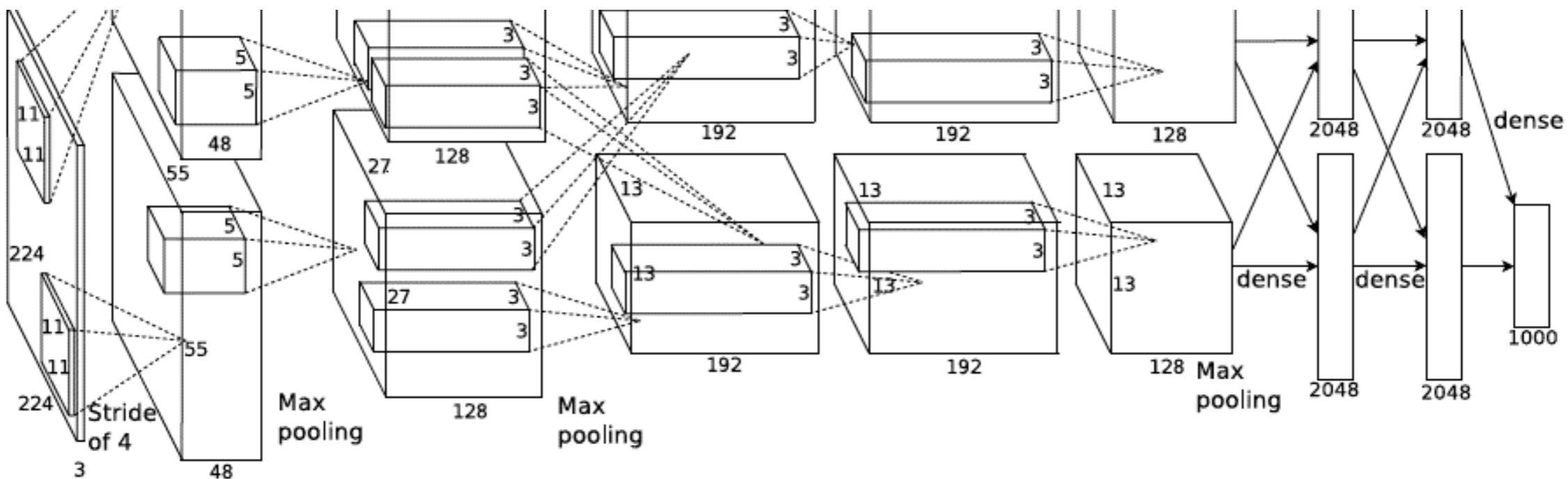
**Bigger model:**

7 hidden layers, 650.000 units, 60 million parameter

**More Data:**

$10^6$  vs  $10^3$  images

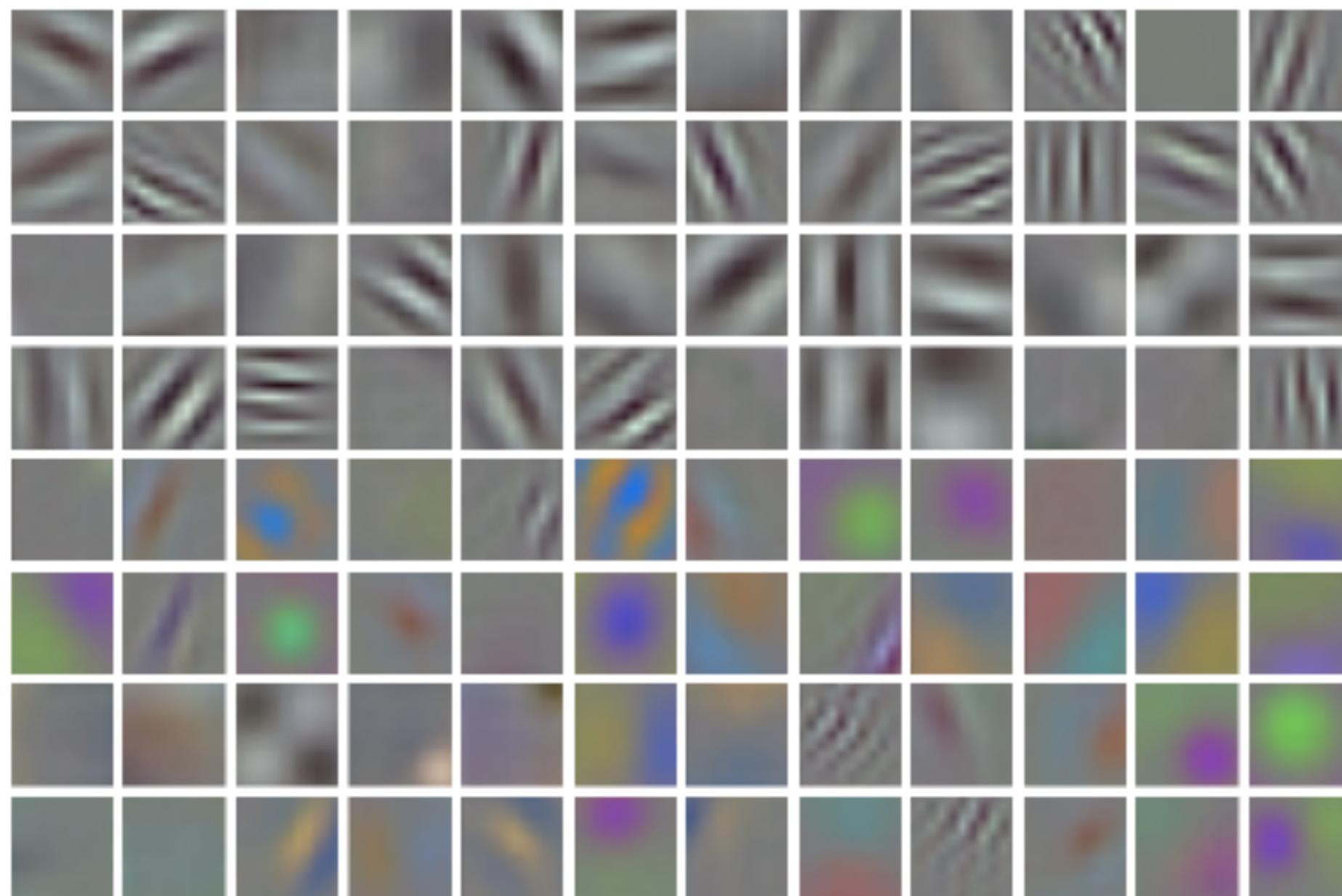
GPU implementation (50x speedup over CPU)



# AlexNet

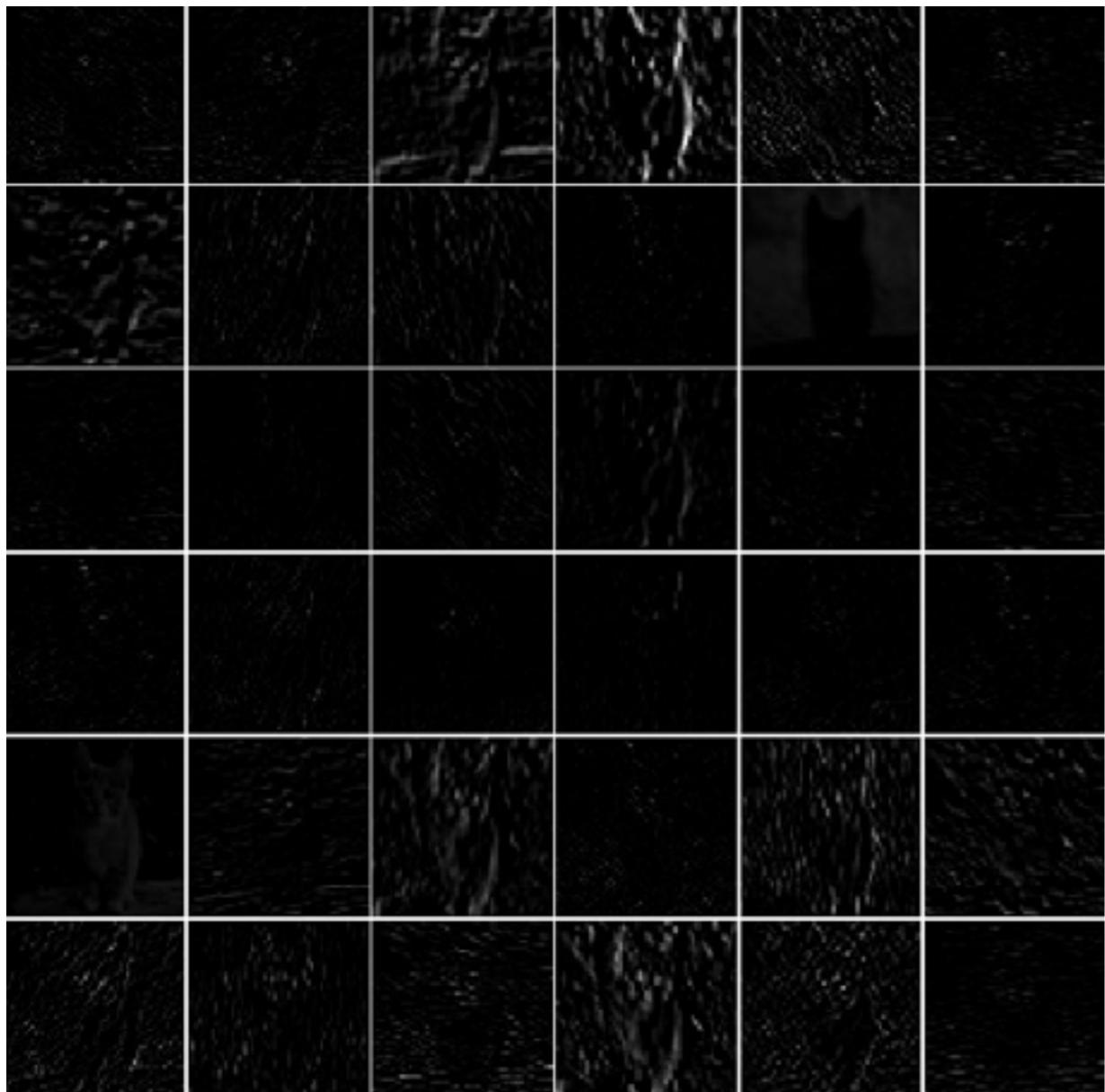
- 1st Layer: 96 conv filters. Size 11x11 (step 4)
  - Pooling + norm
- 2nd Layer: 256 conv filters. Size 5x5
  - Pooling + norm
- 3rd Layer: 384 conv filters. Size 3x3
- 4th Layer: 385 conv filters. Size 3x3
- 5th Layer: 256 conv filters. Size 3x3
  - Pooling
- 6th Layers: Fully Connect. 4096 Neurons
- 7th Layers: Fully Connect. 4096 Neurons
- Output Layer: Fully Connect. 1000 Neurons

# Alexnet 1st Conv Filters

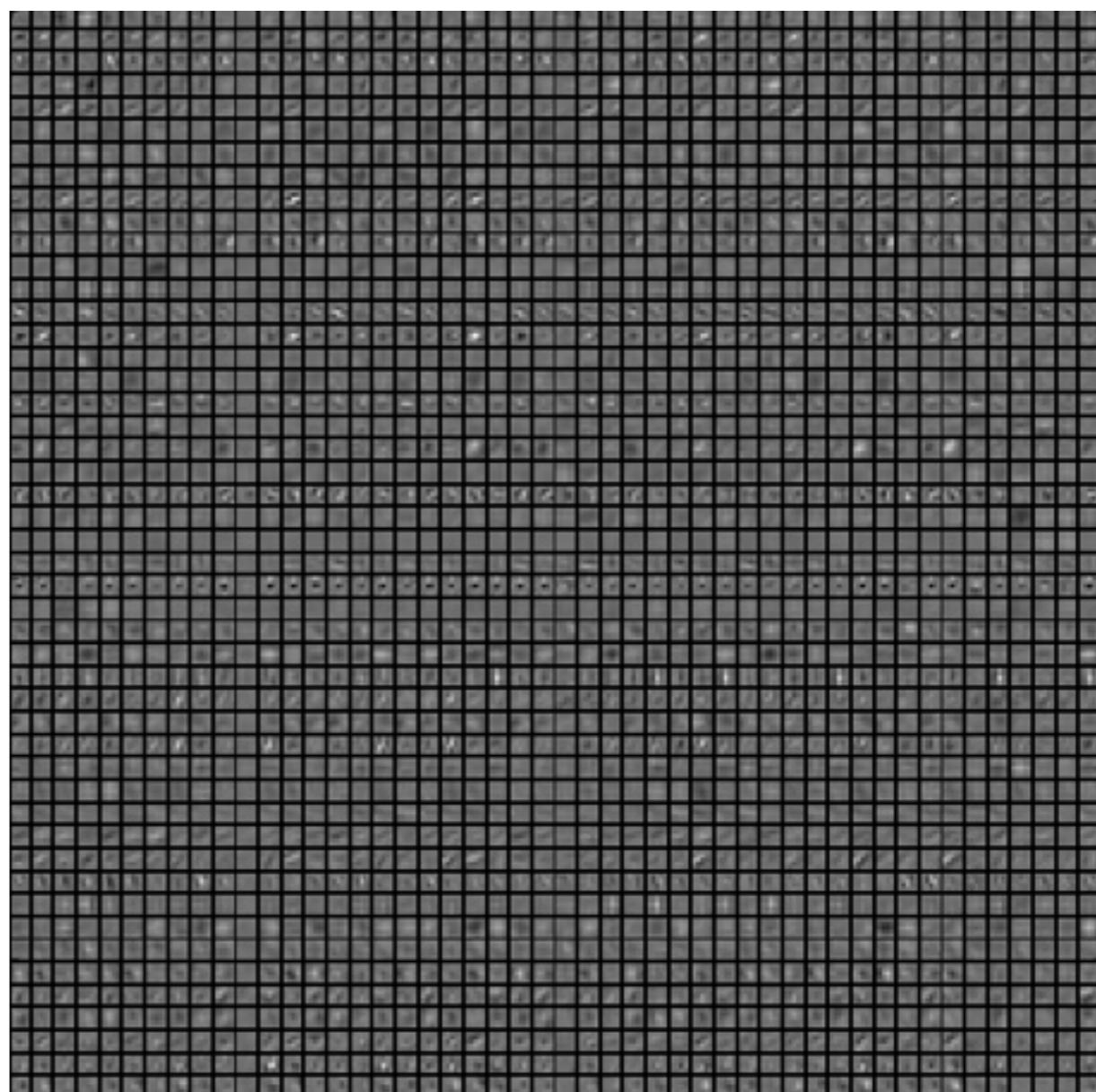


# Alexnet

Feature Map Conv1



Conv2



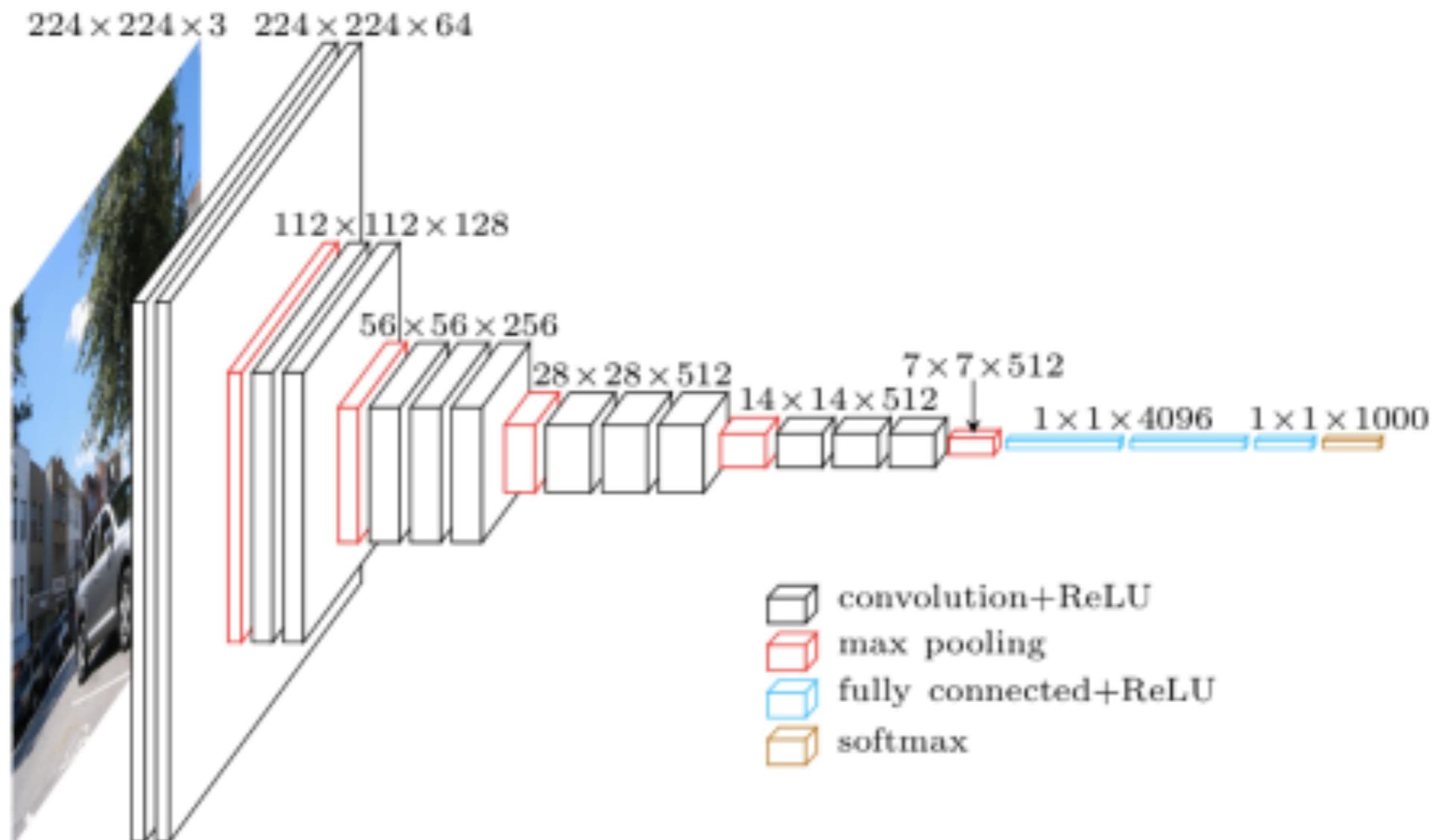
A dramatic, close-up shot of two men in dark suits. The man on the left has light-colored hair and is looking directly at the camera with a serious expression. The man on the right has dark hair and is looking slightly away from the camera. The lighting is low-key, creating strong shadows on their faces.

**WE NEED TO GO**

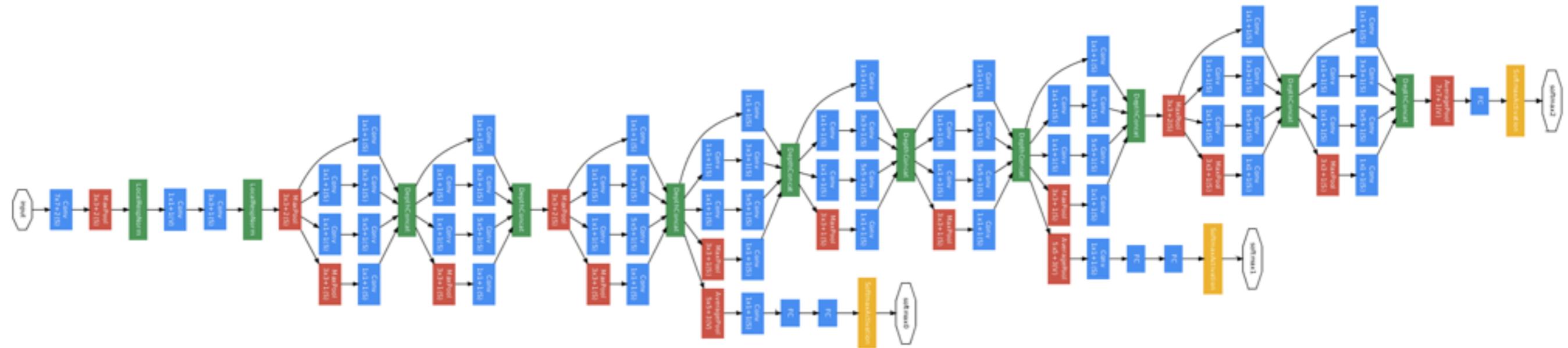
**DEEPER**

troll-face

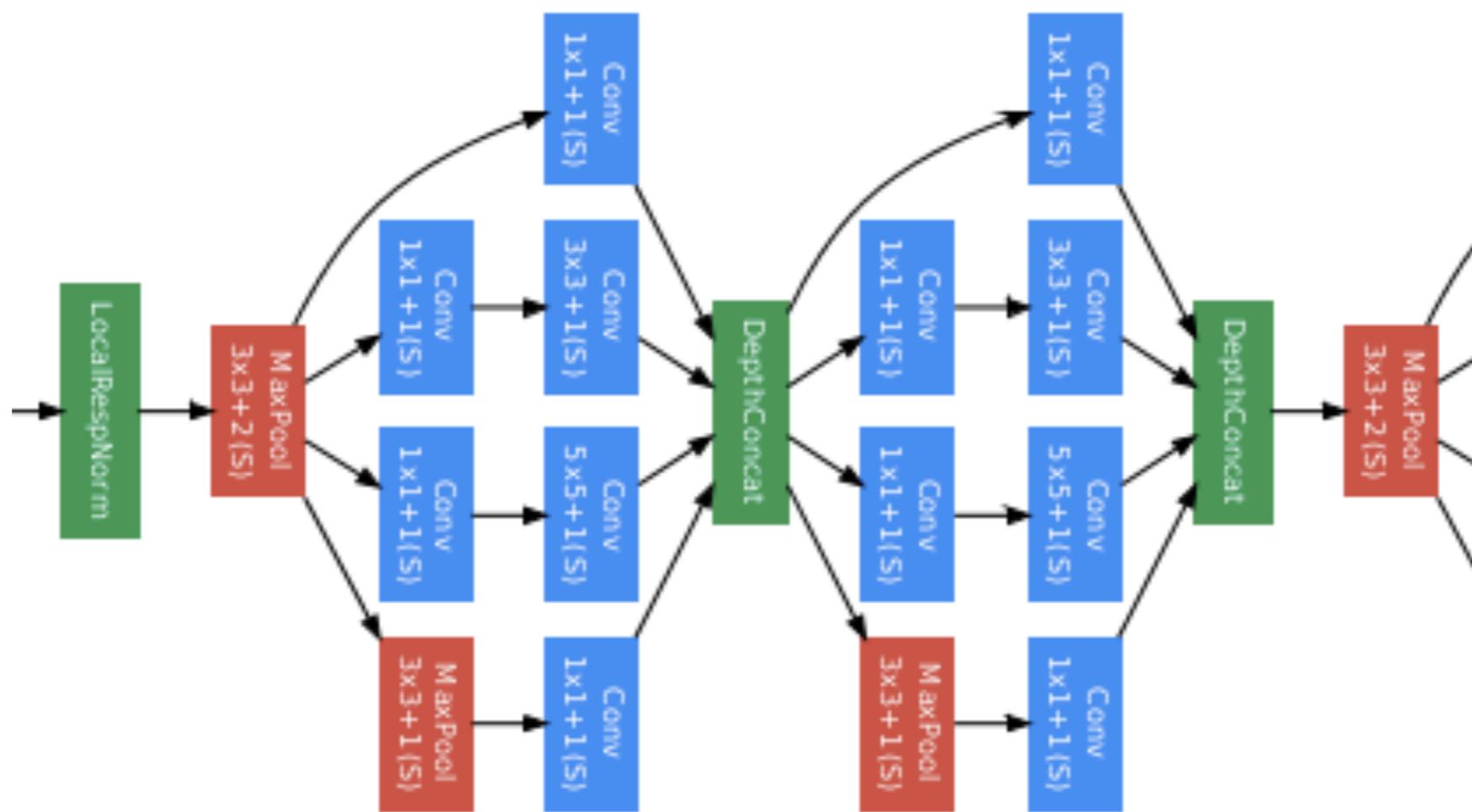
# VGG Net



# GoogleNet



# GoogleNet



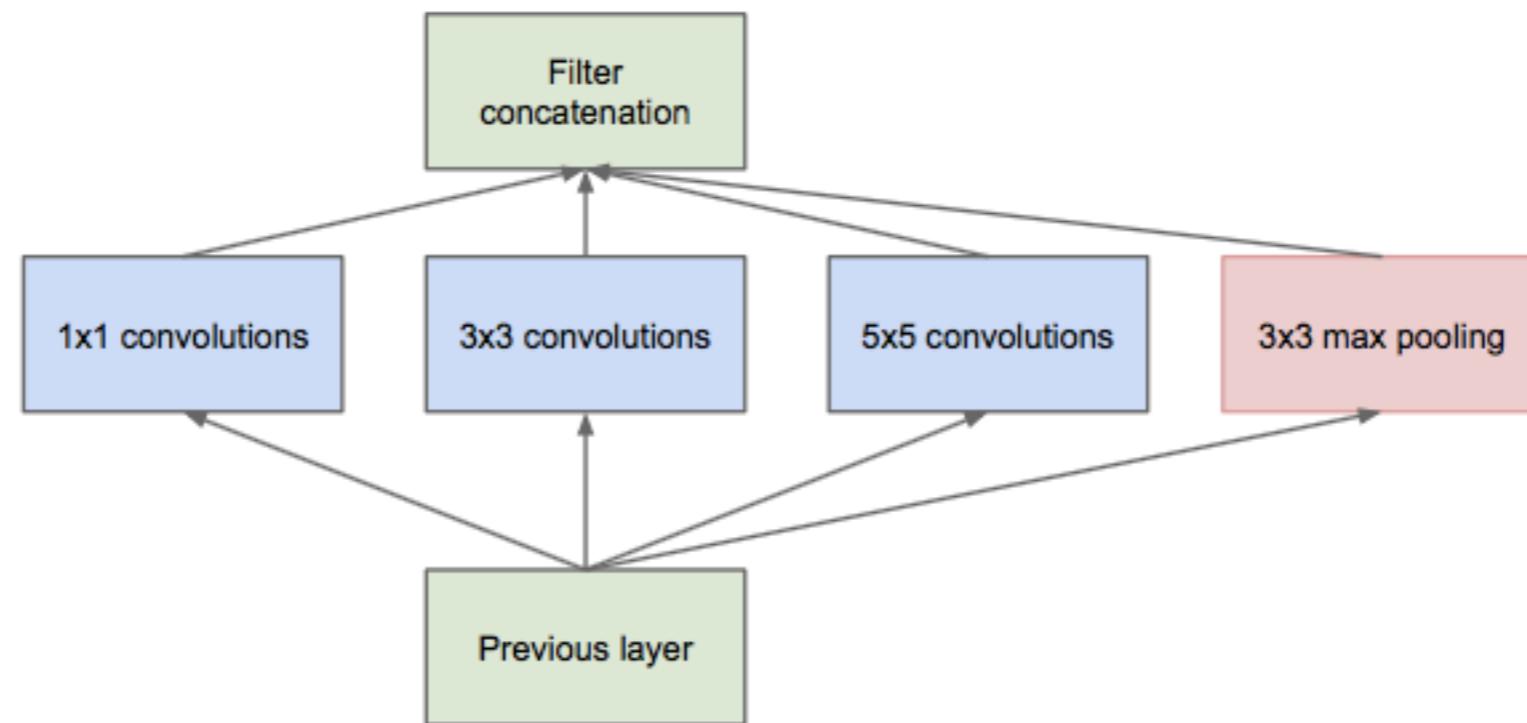
# Inception Module: Naive Version

Convolutional filters with different sizes can cover different clusters of information.

By finding the optimal local construction and repeating it spatially, they approximate the optimal sparse structure with dense components.

For convenience of computation, they use  $1 \times 1$ ,  $3 \times 3$  and  $5 \times 5$  filters + pooling.

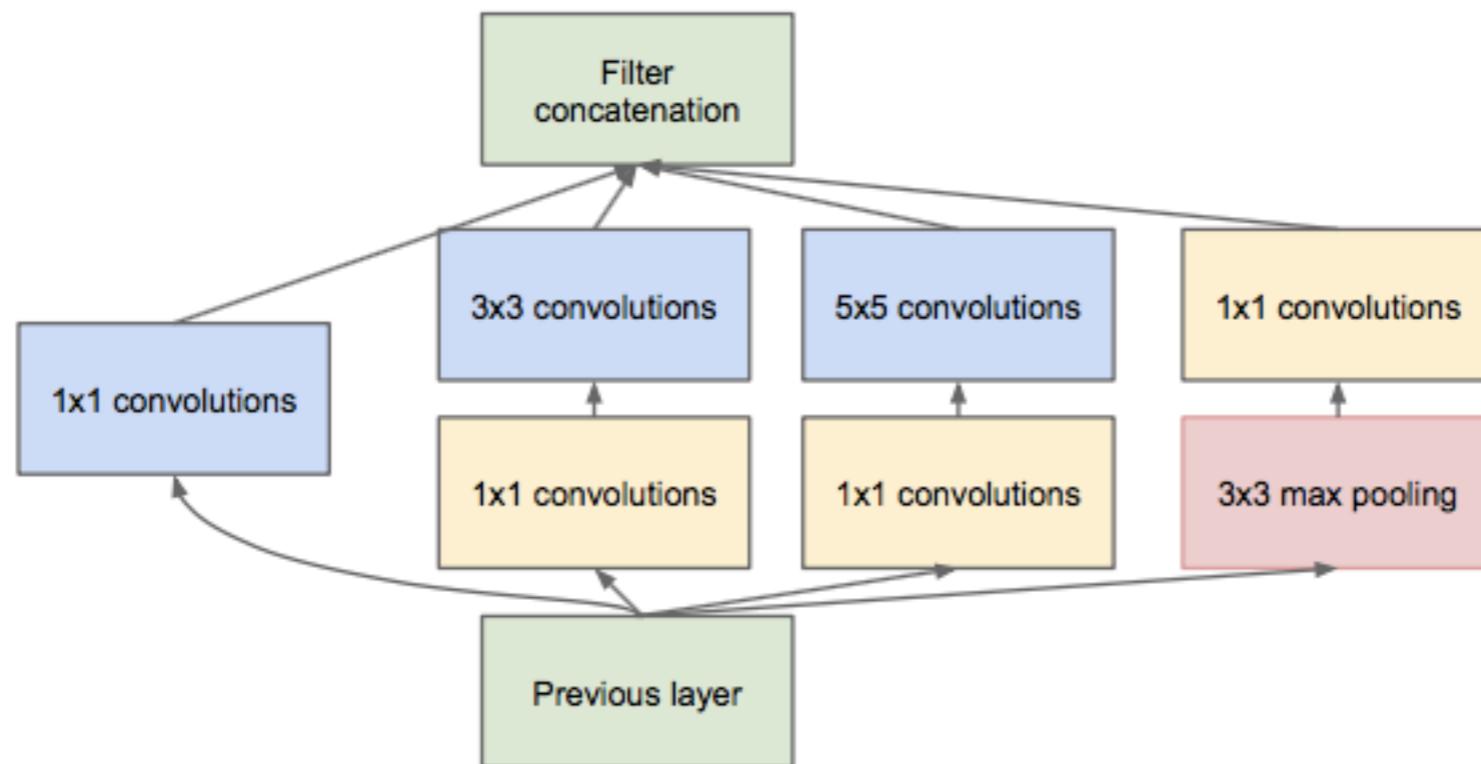
Together these made up the naive Inception module.



# Inception Module: Naive Version

Stacking these inception modules on top of each would lead to an exploding number of outputs

Solution: inspired by "Network in Network" add 1x1 convolutions for dimensionality reduction



# GoogleNet

# ImageNet Challenge

## 2012-2014

<b>Team</b>	<b>Year</b>	<b>Place</b>	<b>Error (top-5)</b>	<b>External data</b>
SuperVision – Toronto (7 layers)	2012	-	16.4%	no
SuperVision	2012	1st	15.3%	ImageNet 22k
Clarifai – NYU (7 layers)	2013	-	11.7%	no
Clarifai	2013	1st	11.2%	ImageNet 22k
VGG – Oxford (16 layers)	2014	2nd	7.32%	no
GoogLeNet (19 layers)	2014	1st	6.67%	no
<u>Human expert*</u>			5.1%	

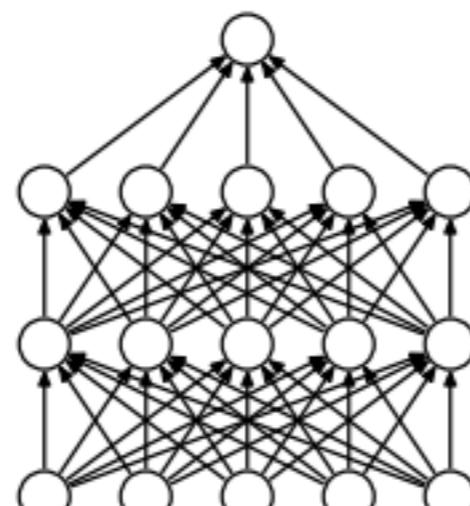
# Training a CNN

- Backpropagation + stochastic gradient descent with momentum
- Dropout
- Data Augmentation
- Batch Normalization
- Initialization
  - Transfer Learning

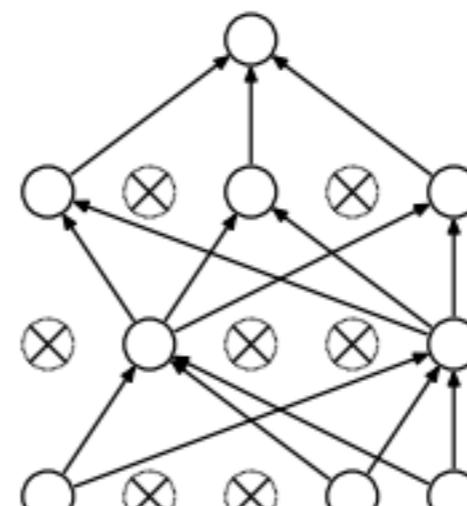
# Dropout

**Dropout:** A Simple Way to Prevent Neural Networks from Overfitting

*Journal of Machine Learning Research 15 (2014) 1929-1958*



(a) Standard Neural Net



(b) After applying dropout.