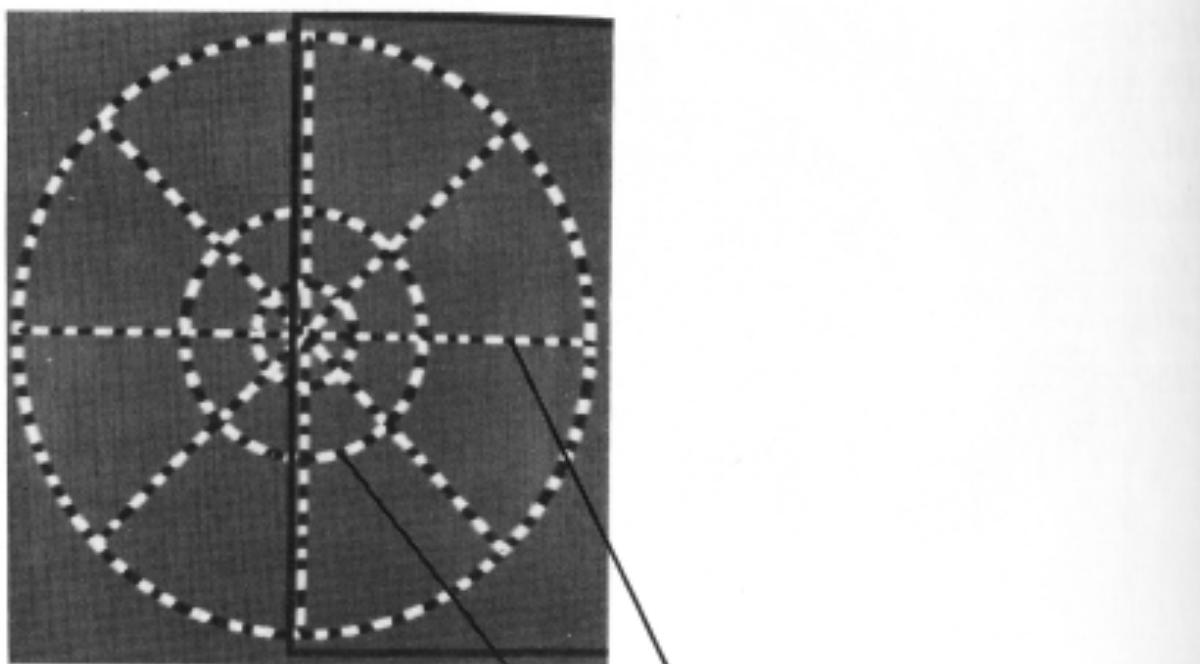


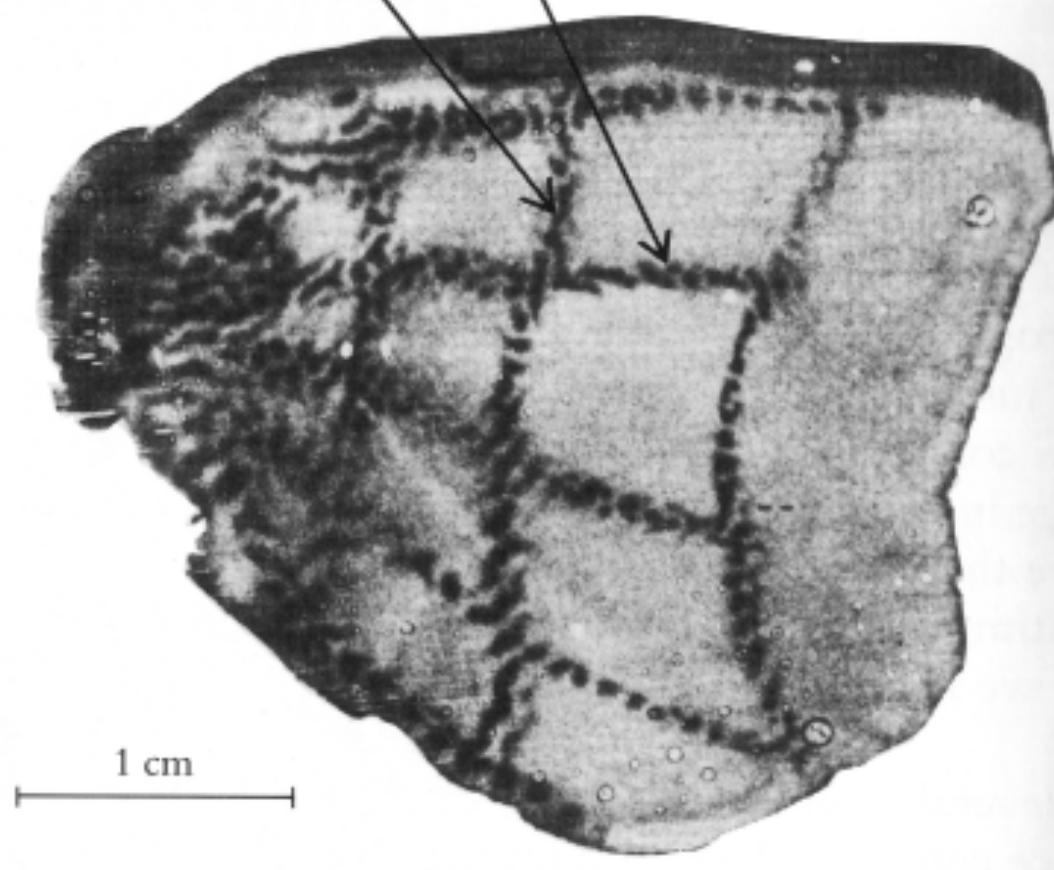
# Learning in Convolutional Neural Networks (CNNs)



## Retinotopic Map

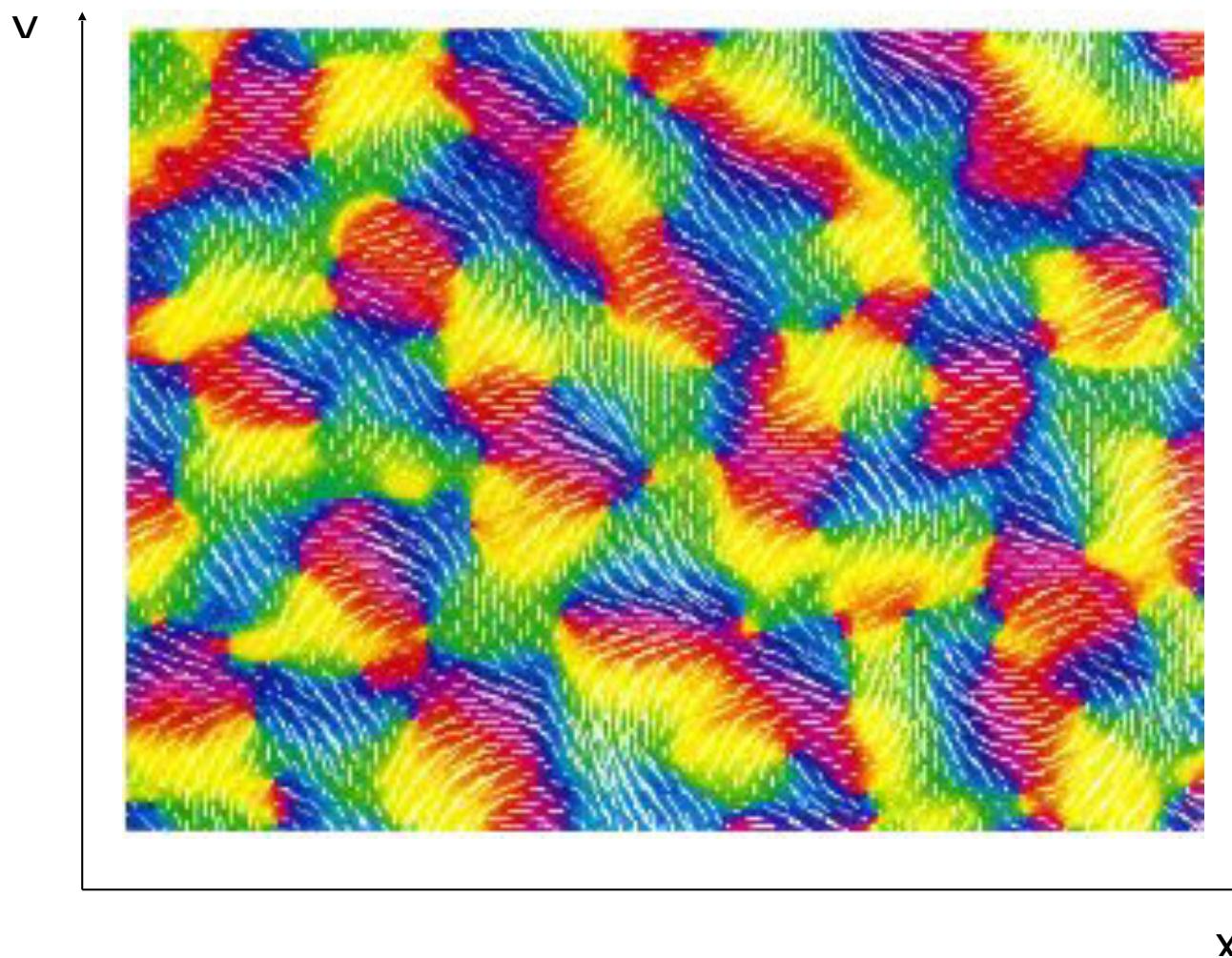
A flickering visual pattern excites cells across the visual field

Each mini dark region represents 150,000 cells



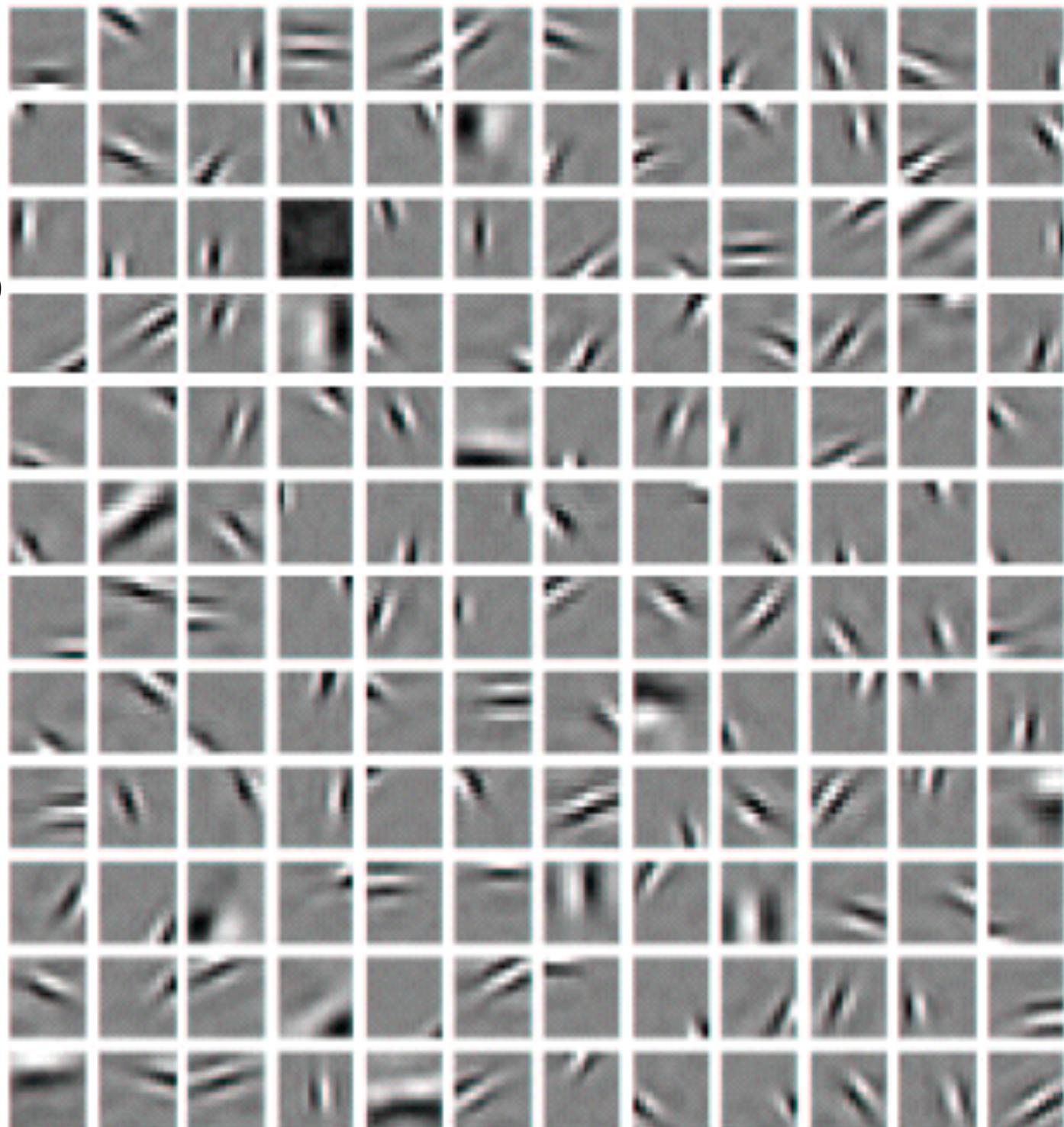
## Orientation Maps in Striate Cortex

Across visual coordinates, cells responding to local photometric “edges” at different orientation (here denoted by different colors) are clumped into groups



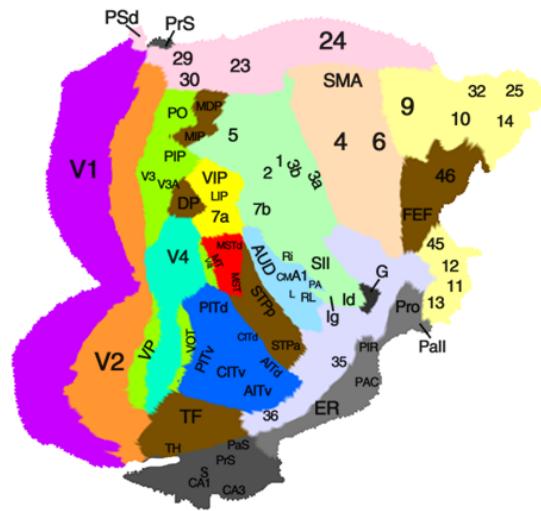
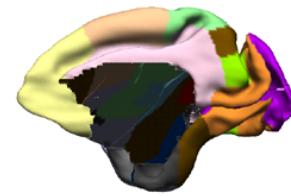
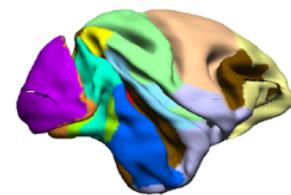
## Different features

Receptive fields (here 144) of cells at a given location can be learned

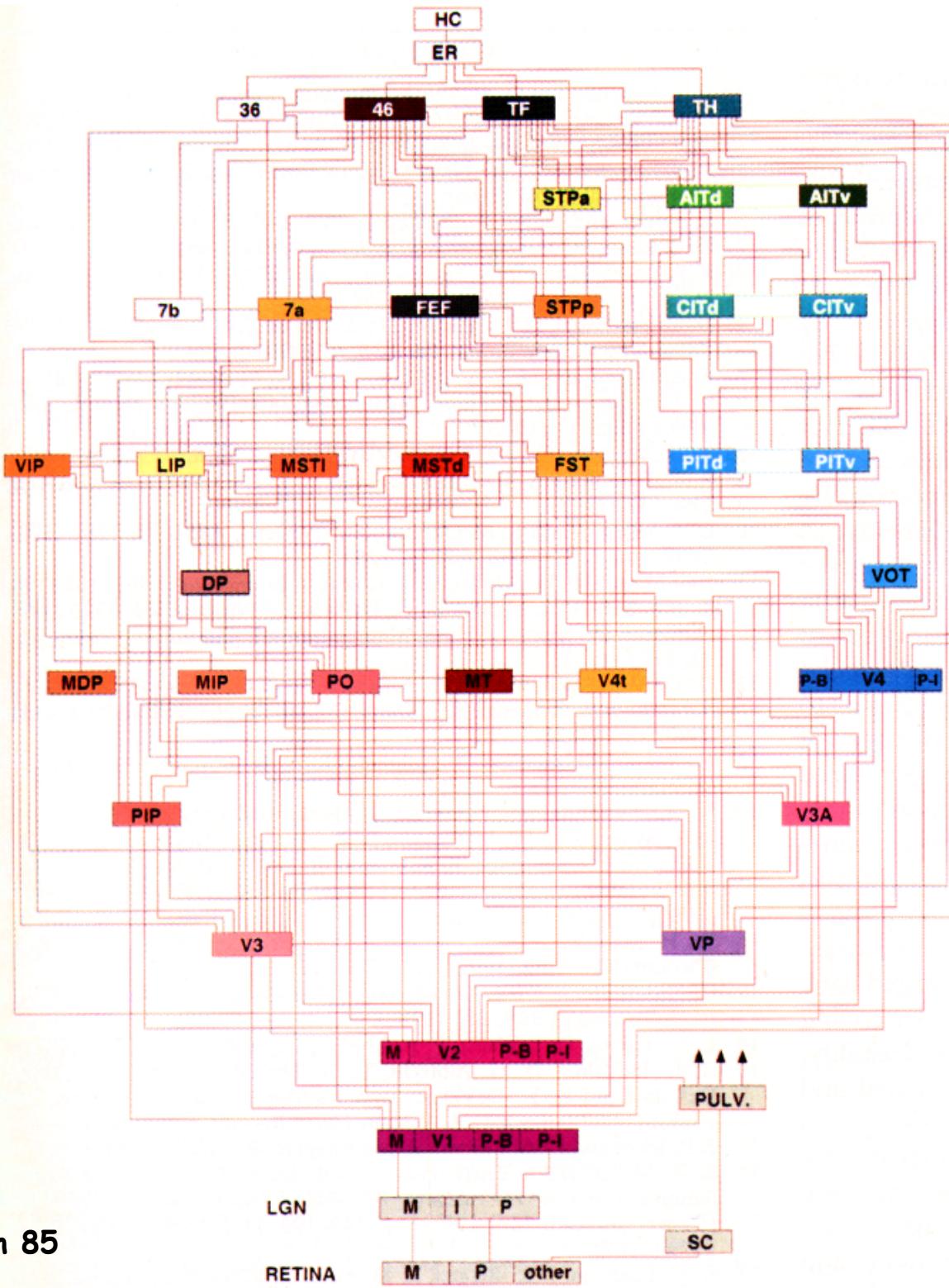
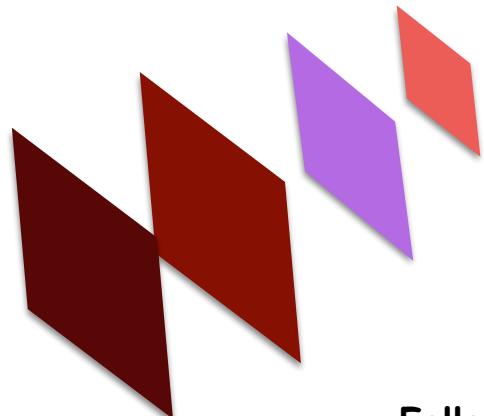


## Flattened Representation of Monkey Cerebral Cortex

Van Essen and Drury 1998

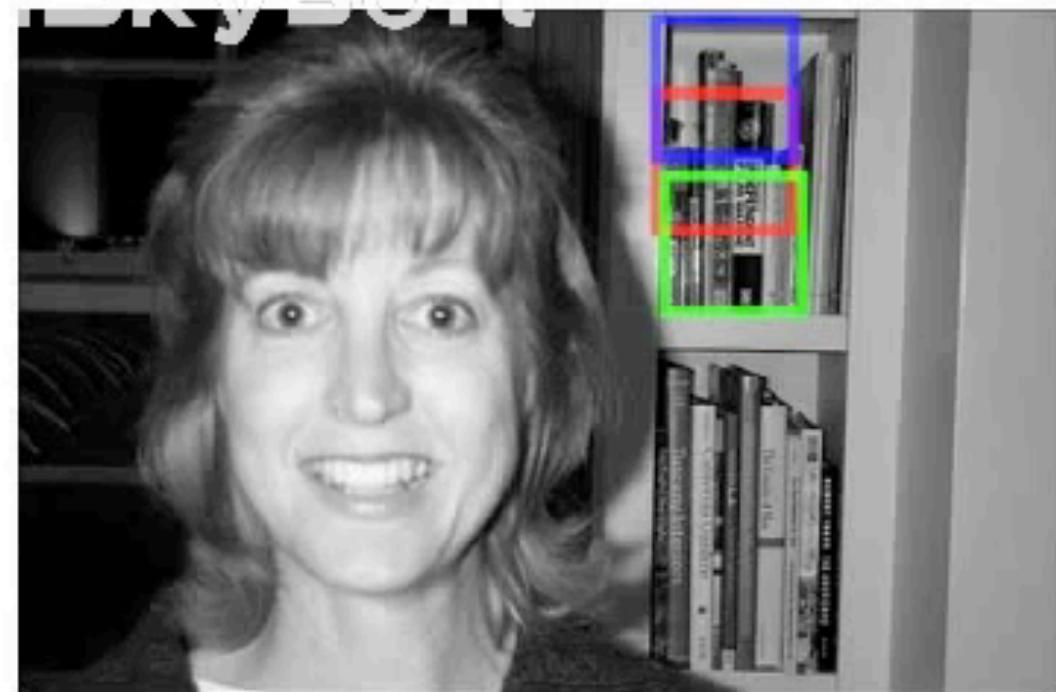


Successive maps are connected to each other



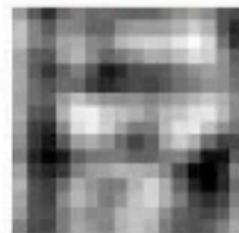


image\_0182.pgm

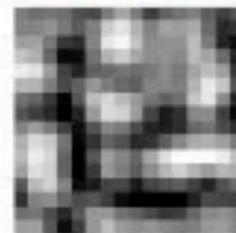


1. Templates tile the image
2. Best match template adds its input into itself and all its copies

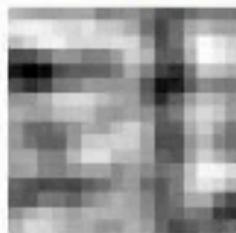
1 spikes

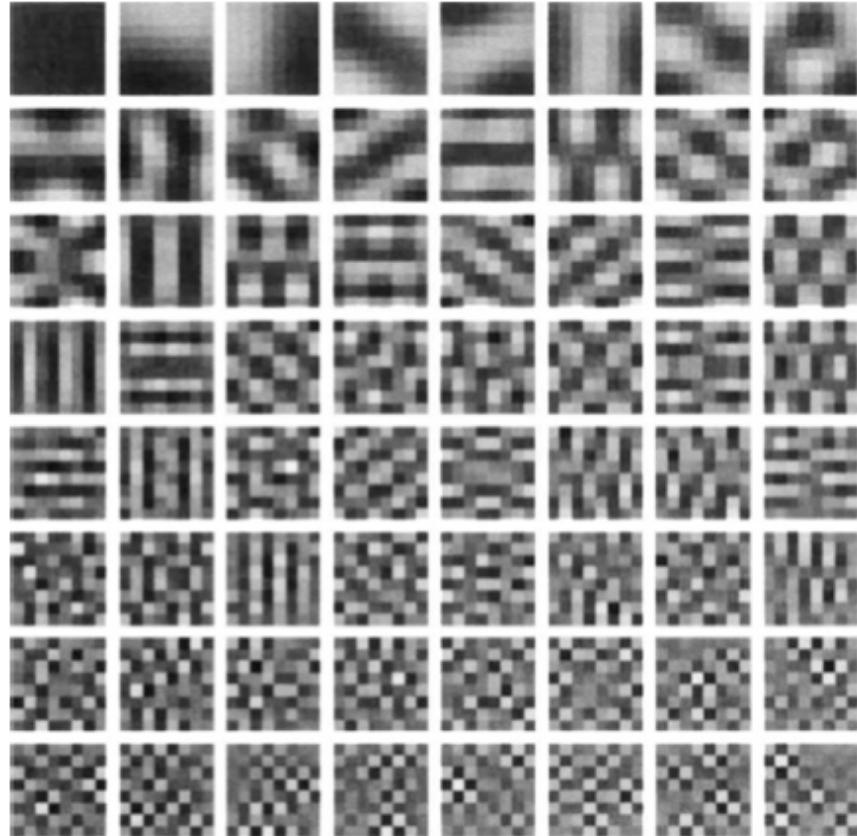


1 spikes



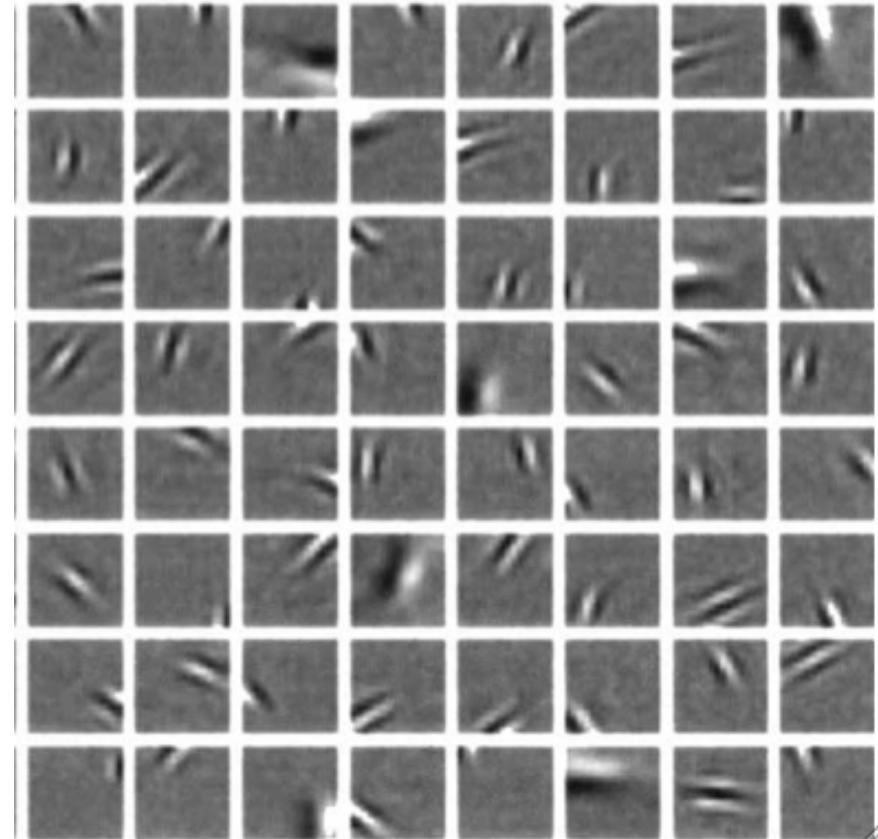
1 spikes





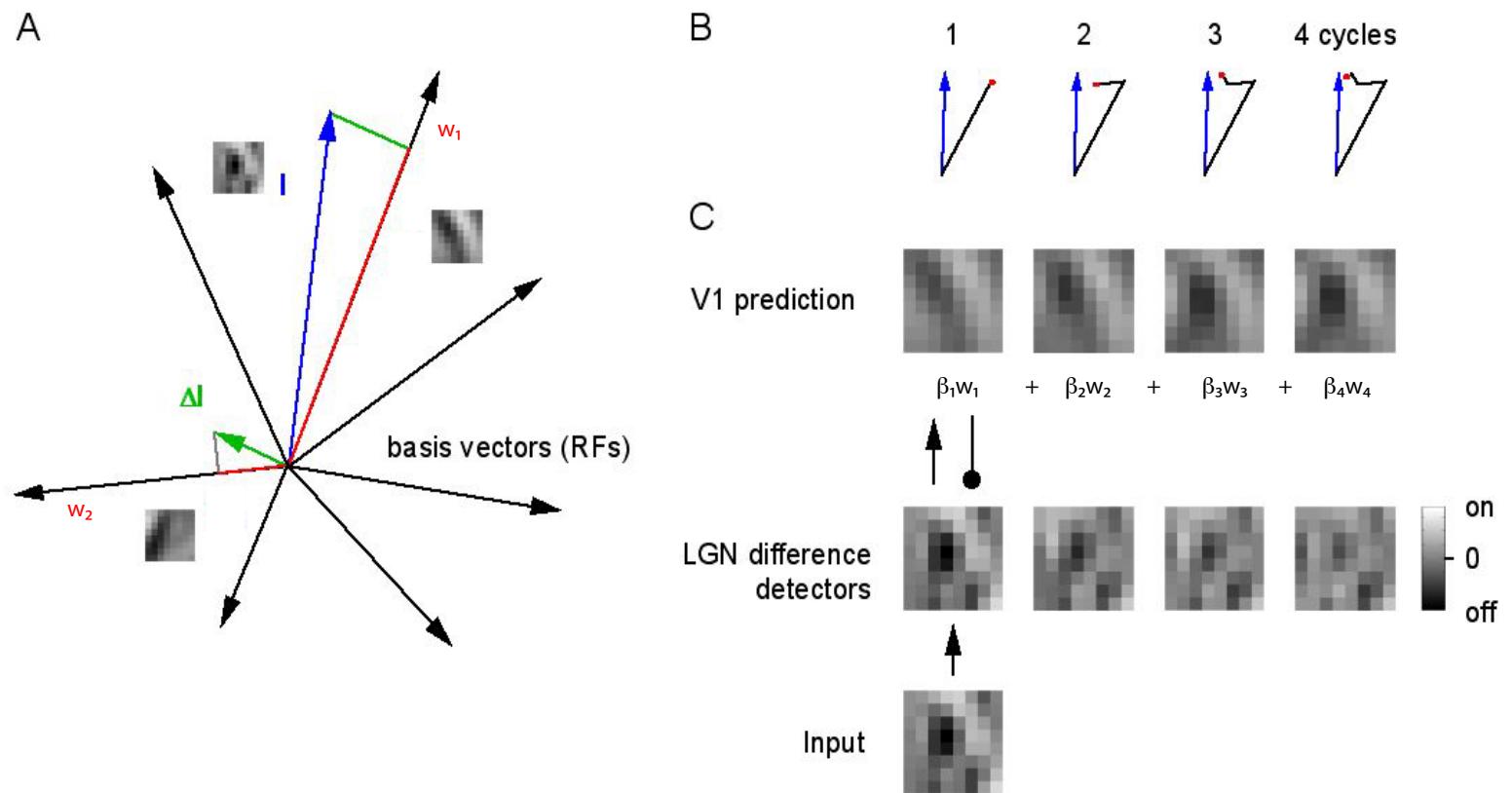
$$\min_{\mathbf{r}} \quad \|\mathbf{I} - U\mathbf{r}\|_2$$

Learn receptive fields  
by minimizing the  $\ell_2$  fit only.

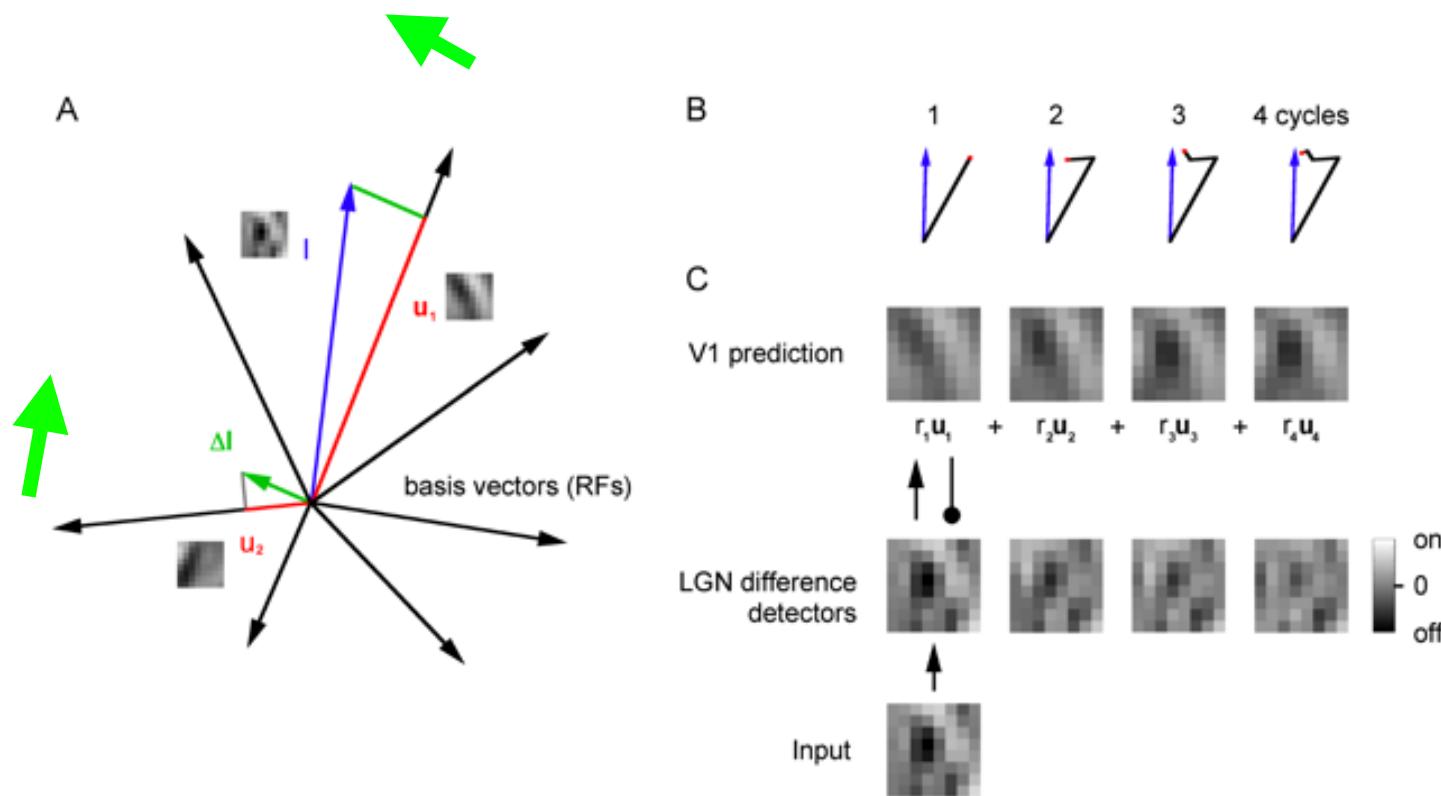


$$\min_{\mathbf{r}} \quad \|\mathbf{I} - U\mathbf{r}\|_2 + \lambda \|\mathbf{r}\|_1$$

Learn receptive fields  
by adding a term that penalizes large synapses  
 $\|\cdot\|_1$  means absolute value

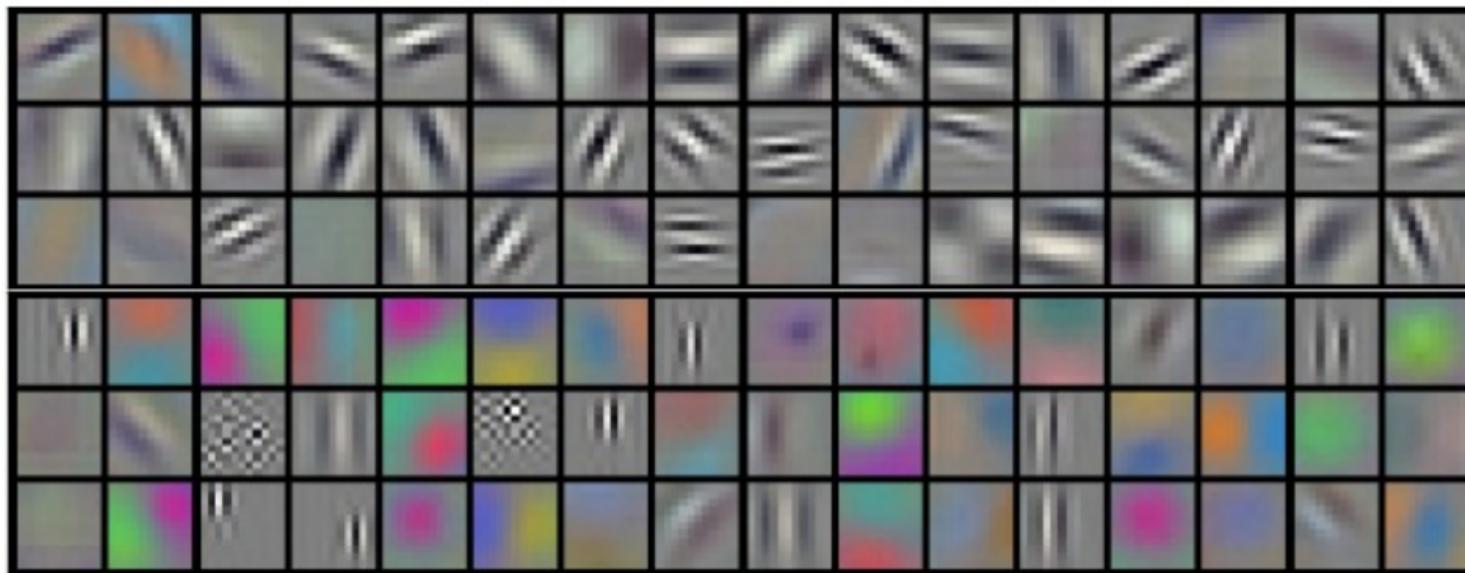


# A standard Algorithm: Matching Pursuit

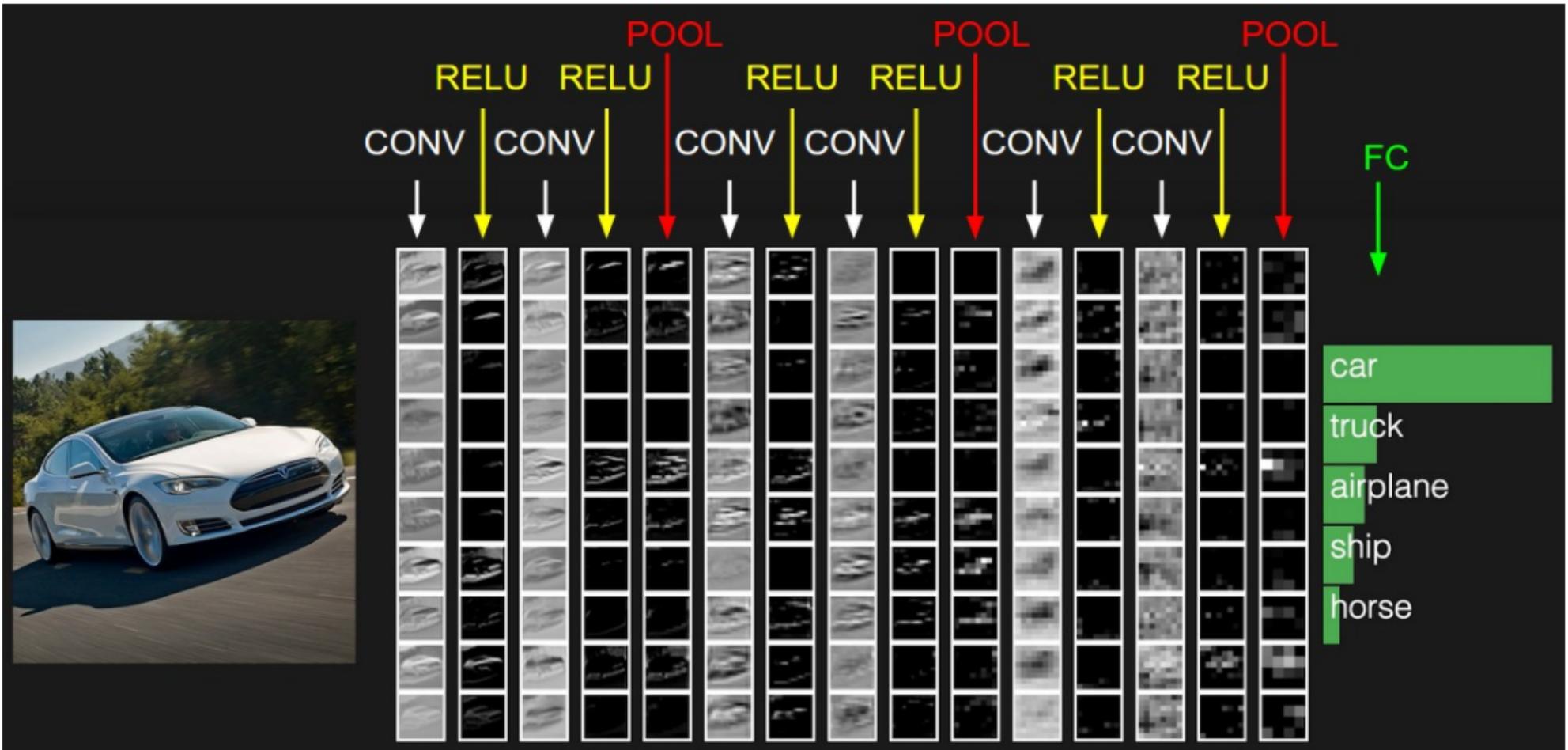


**Learning algorithm: move RFs  
of winning neurons towards inputs**

# Learned filters



Example filters learned by Krizhevsky et al. Each of the 96 filters shown here is of size [11x11x3], and each one is shared by the 55\*55 neurons in one depth slice. Notice that the parameter sharing assumption is relatively reasonable: If detecting a horizontal edge is important at some location in the image, it should intuitively be useful at some other location as well due to the translationally-invariant structure of images. There is therefore no need to relearn to detect a horizontal edge at every one of the 55\*55 distinct locations in the Conv layer output volume.



|   |   |   |
|---|---|---|
| 1 | 0 | 1 |
| 0 | 1 | 0 |
| 1 | 0 | 1 |

Convolution  
Template

How convolution operator works

|   |   |   |   |   |
|---|---|---|---|---|
| 1 | 1 | 1 | 0 | 0 |
| 0 | 1 | 1 | 1 | 0 |
| 0 | 0 | 1 | 1 | 1 |
| 0 | 0 | 1 | 1 | 0 |
| 0 | 1 | 1 | 0 | 0 |

Image

|   |  |  |
|---|--|--|
| 4 |  |  |
|   |  |  |
|   |  |  |
|   |  |  |

Convolved  
Feature

Input Volume (+pad 1) (7x7x3)

| x[::,::,0]    |
|---------------|
| 0 0 0 0 0 0 0 |
| 0 1 1 1 2 2 0 |
| 0 2 0 2 1 1 0 |
| 0 1 2 2 2 1 0 |
| 0 1 0 2 1 1 0 |
| 0 0 1 1 2 0 0 |
| 0 0 0 0 0 0 0 |

Filter W0 (3x3x3)

| w0[::,::,0] |
|-------------|
| -1 0 1      |
| 1 1 1       |
| 0 -1 0      |

| w0[::,::,1] |
|-------------|
| 1 1 0       |
| 0 -1 -1     |
| -1 -1 -1    |

| w0[::,::,2] |
|-------------|
| 0 0 0       |
| -1 1 1      |
| -1 -1 -1    |

Bias b0 (1x1x1)  
b0[::,::,0]

1

| x[::,::,1]    |
|---------------|
| 0 0 0 0 0 0 0 |
| 0 2 0 2 2 0 0 |
| 0 0 0 1 2 2 0 |
| 0 1 0 2 2 1 0 |
| 0 0 0 2 0 2 0 |
| 0 0 2 2 2 2 0 |
| 0 0 0 0 0 0 0 |

| x[::,::,2]    |
|---------------|
| 0 0 0 0 0 0 0 |
| 0 1 2 0 0 2 0 |
| 0 2 0 0 1 0 0 |
| 0 1 0 0 2 2 0 |
| 0 2 1 1 0 0 0 |
| 0 2 1 1 0 2 0 |
| 0 0 0 0 0 0 0 |

Filter W1 (3x3x3)

| w1[::,::,0] |
|-------------|
| -1 0 -1     |
| 1 -1 1      |
| -1 1 -1     |

| w1[::,::,1] |
|-------------|
| 1 0 1       |
| 1 1 -1      |
| -1 1 0      |

| w1[::,::,2] |
|-------------|
| -1 0 0      |
| 1 1 0       |
| -1 0 -1     |

Bias b1 (1x1x1)

| b1[::,::,0] |
|-------------|
| 0           |

Output Volume (3x3x2)

| o[::,::,0] |
|------------|
| 0 -7 1     |
| 0 1 3      |
| 3 3 5      |

| o[::,::,1] |
|------------|
| 5 5 3      |
| 3 5 10     |
| 1 3 4      |

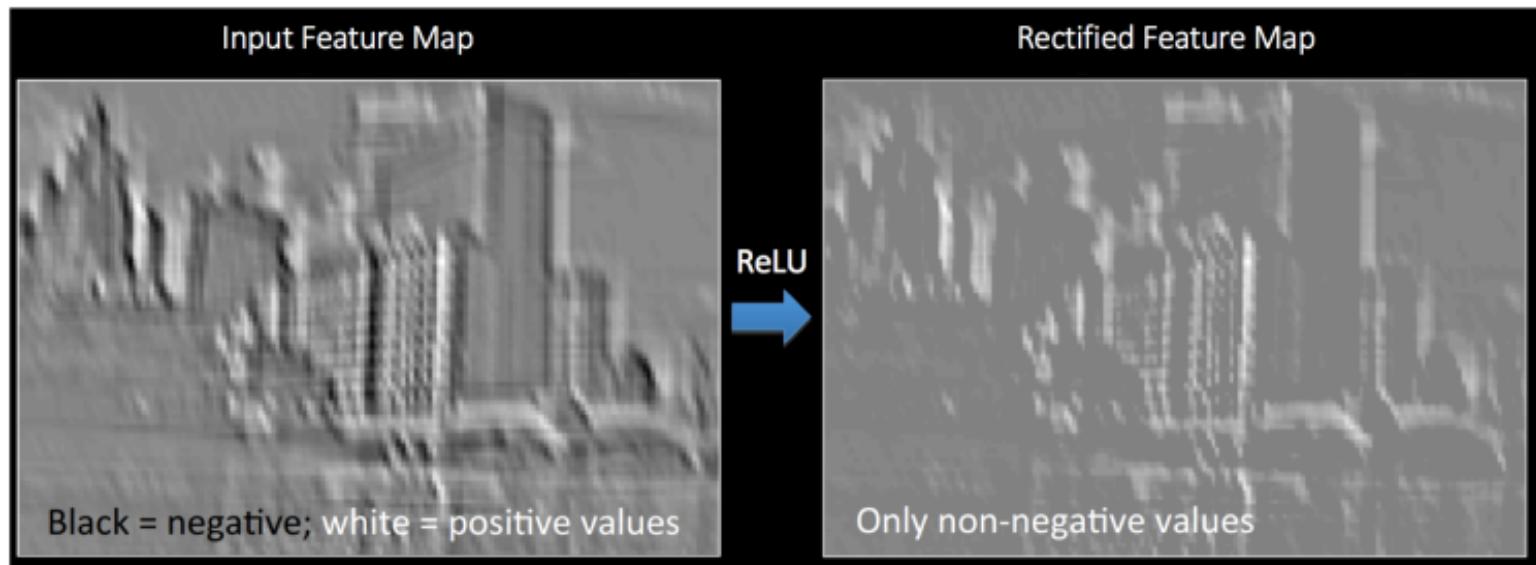
toggle movement

Example for a larger image with two feature templates

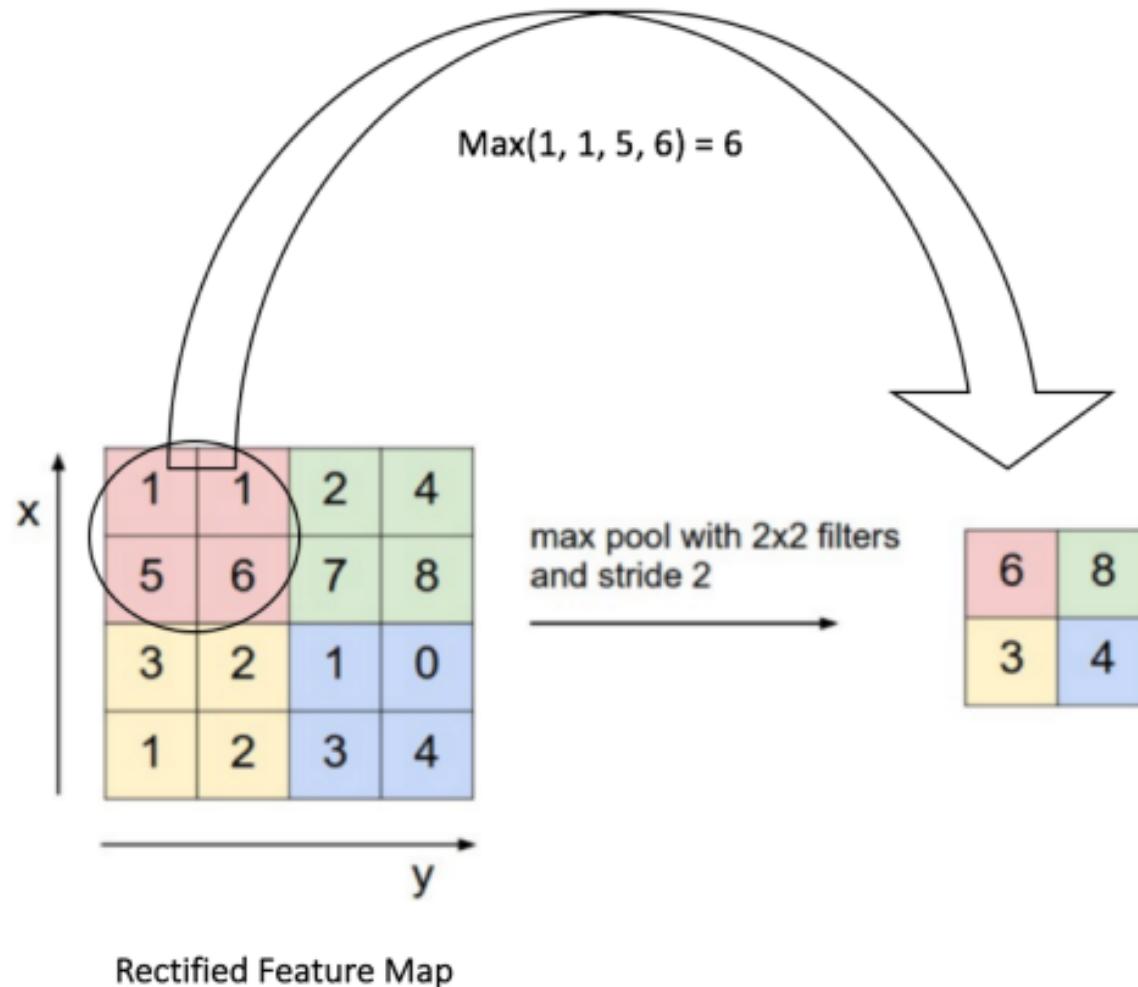


Input

RELUs [Rectified Linear Units] introduce nonlinearity



# Stride and Max Pooling



## Backpropagation

To keep things simple, let us just work with one pattern. In that case the objective function is defined to be

$$E = \frac{1}{2} \|\boldsymbol{x}^K - \boldsymbol{d}\|^2 \quad (1)$$

where  $K$  is an index denoting the last layer in the network and  $d$  is the desired output.

The equation for updating the states is given by

$$\boldsymbol{x}^{k+1} = g(W^k \boldsymbol{x}^k) \quad (2)$$

where  $W^k$  is a matrix used to store the weights between layer  $k$  and layer  $k + 1$ ,

$$W^k = \begin{bmatrix} w_{11}^k & w_{12}^k & \dots \\ w_{21}^k & \ddots & \\ \vdots & & w_{nn}^k \end{bmatrix} = \begin{pmatrix} \boldsymbol{w}_1^k \\ \boldsymbol{w}_2^k \\ \vdots \\ \boldsymbol{w}_n^k \end{pmatrix}$$

and the special understanding we shall have is that the function  $g$  applied to a vector is just that function applied to its elements; that is,

$$g(W\boldsymbol{x}) = \begin{pmatrix} g(\boldsymbol{w}_1 \cdot \boldsymbol{x}) \\ g(\boldsymbol{w}_2 \cdot \boldsymbol{x}) \\ \vdots \end{pmatrix}$$

This is just a version of the optimal control problem where Equation 1 is the optimand and Equation 2 is the dynamics. Here instead of the control  $\boldsymbol{u}$ , the matrices  $W^k, k = 1, \dots, K$  represent the “control” variables. One difference, of course, is the dimensionality: In networks the dimension of the state space may range into the thousands! But the important thing to realize is that the mathematical treatment is the same.

## Generating the Solution

To generate the solution, assume an initial  $W$  matrix. Then solve the dynamic equations going forward in levels in the network. This solution provides a value for  $\mathbf{x}^K$ . This value allows the adjoint system to be solved backward for  $k = K$  to  $k = 0$ :

$$\begin{aligned}\boldsymbol{\lambda}^K &= \mathbf{x}^K - \mathbf{d} \quad \text{for } k = K \\ \boldsymbol{\lambda}^k &= -W^T \Lambda^{k+1} g'(W\mathbf{x}^k) \quad \text{for } k = 0, \dots, K-1\end{aligned}$$

where  $g'$  denotes differentiation with respect to its argument. Now improve the estimate for  $W$  using gradient descent. To do so, calculate the adjustment for the weights as  $\frac{\partial H}{\partial w_{ij}}$ :

$$\frac{\partial H}{\partial w_{ij}^k} = \lambda_i^{k+1} x_j^k g'(\mathbf{w}_i^k \cdot \mathbf{x}^k) \quad (5)$$

Equation 5 is very compact and so is worth unpacking in order to understand it. The multiplier  $\lambda_i^{k+1}$  can be understood by realizing that at level  $K$  it just keeps track of the output error. This is its role for the internal units also; it appropriately keeps track of how the internal weight change affects output error. The derivative  $g'$  is also simple, especially when using

$$g(u) = \frac{1}{1 + e^{-u}}$$

since it can be easily verified that

$$g'(u) = g(1 - g)$$

Thus to calculate  $g'$  for a given unit, just determine how much input the  $x_j^k$ <sup>th</sup> unit receives, and apply that as an argument to  $g'()$ .

Biological features that have inspired CNN architecture are mostly based on how the brain's cortex handles vision

Retinopathy: visual features are represented across the visual field *explicitly*

Specific properties: visual features are *repeated* across the visual field explicitly

Feature maps are connected to a layer above and down sampled to create a hierarchy of maps with properties of increasing abstraction

All these properties are incorporated into CNNs

# Now back to Backpropagation

For the fully connected network ...

For the weights:

$$\Delta w_{ij}^k \propto \lambda_i^{k+1} x_j^k g'(\mathbf{w}_i^k \cdot \mathbf{x}^k)$$

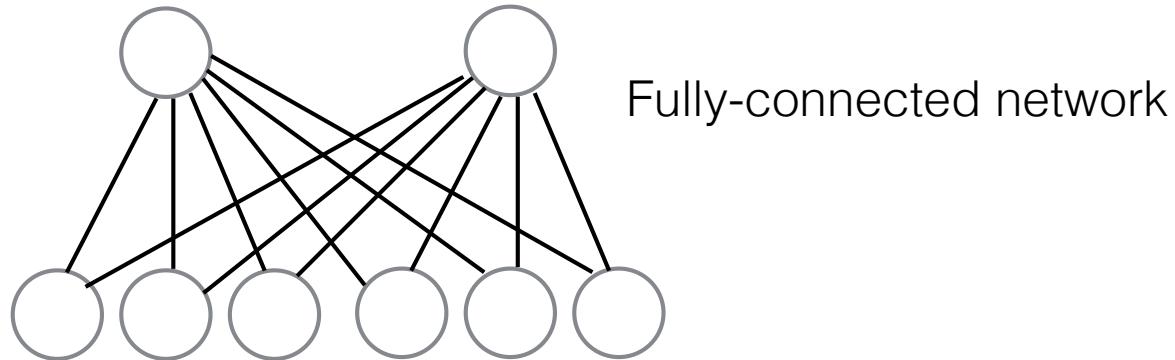
$$\boldsymbol{\lambda}^K = \mathbf{x}^K - \mathbf{d} \text{ for } k = K$$

For the errors:

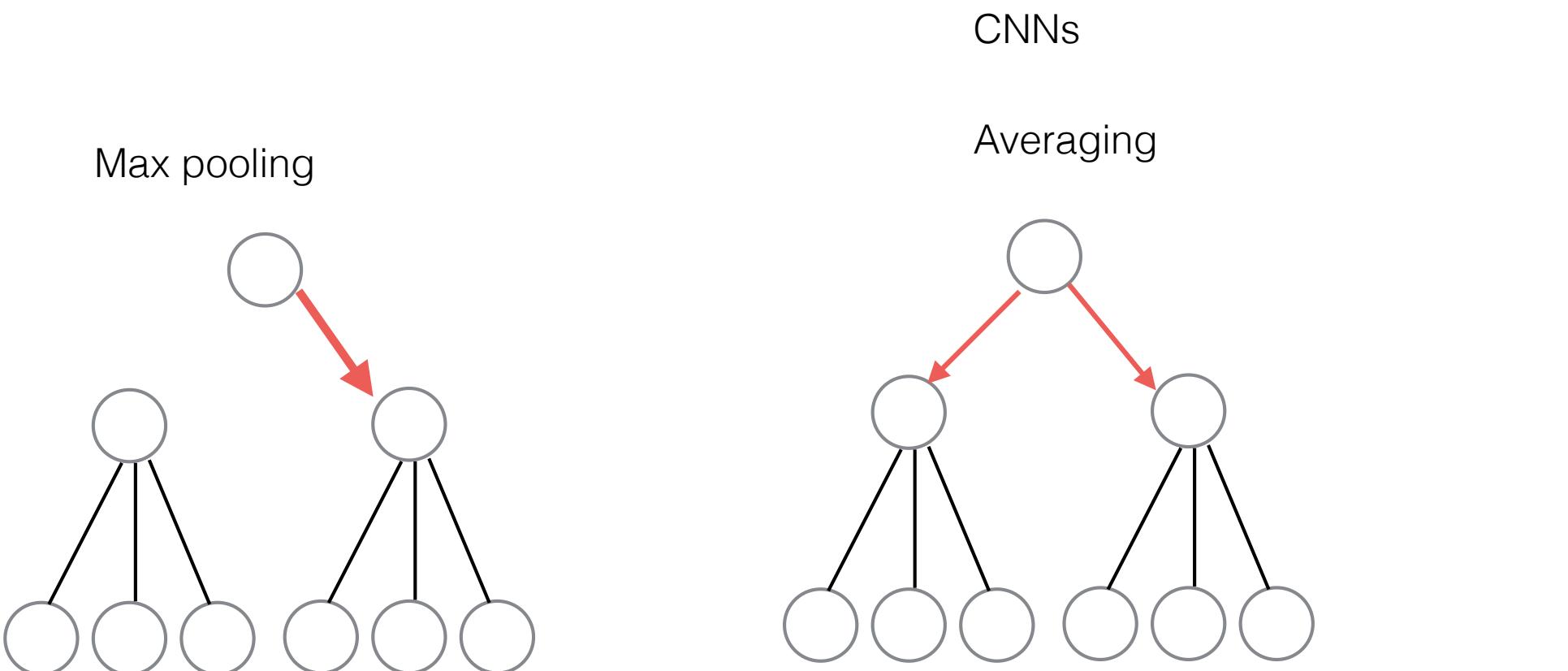
$$\boldsymbol{\lambda}^k = -W^T \Lambda^{k+1} g'(W\mathbf{x}^k) \text{ for } k = 0, \dots, K-1$$

For the CNN ...

$$\boldsymbol{\lambda}^k = \text{upsample} ( -W^T \Lambda^{k+1} g'(W\mathbf{x}^k) )$$



Fully-connected network

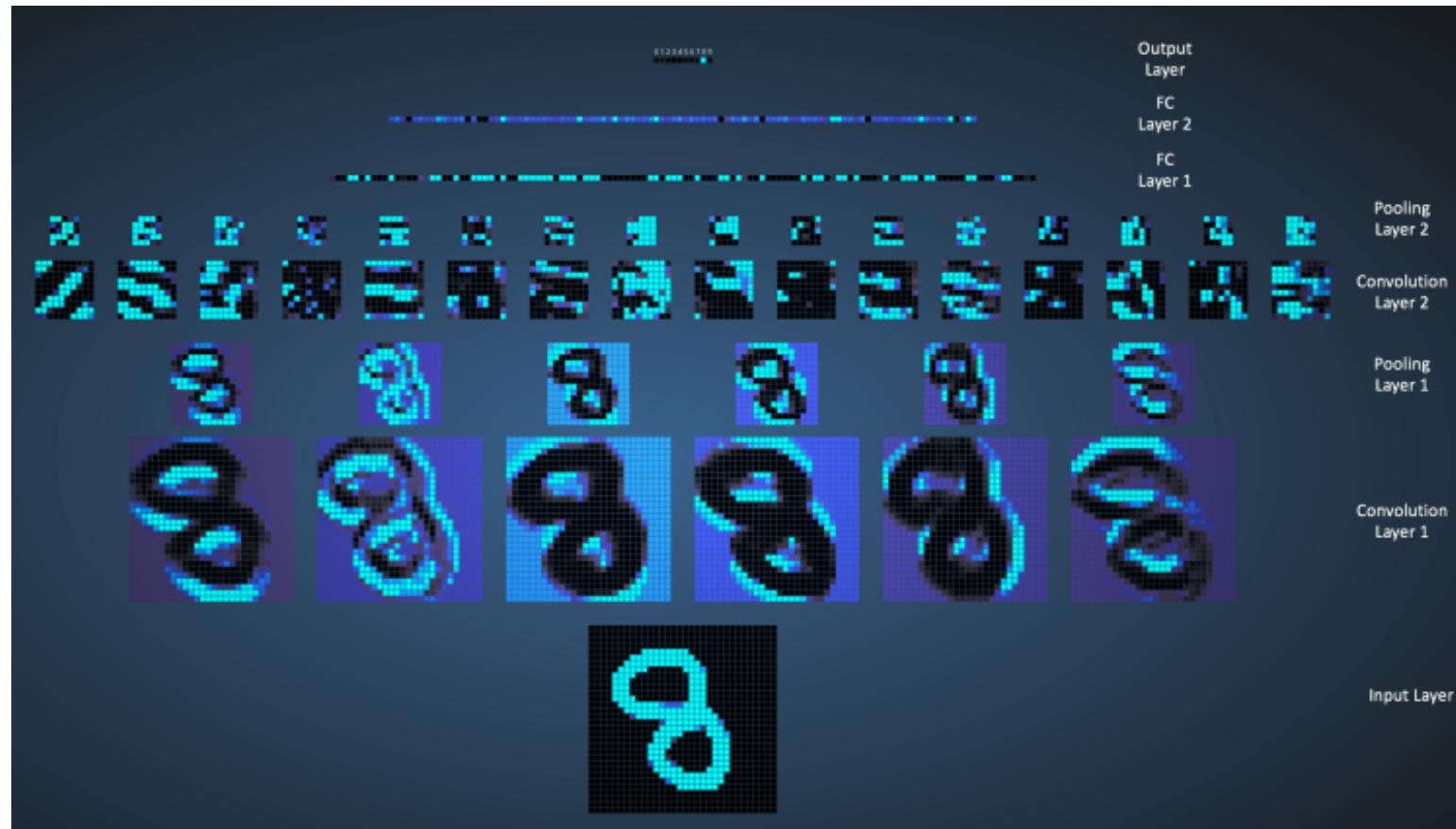


Max pooling

CNNs

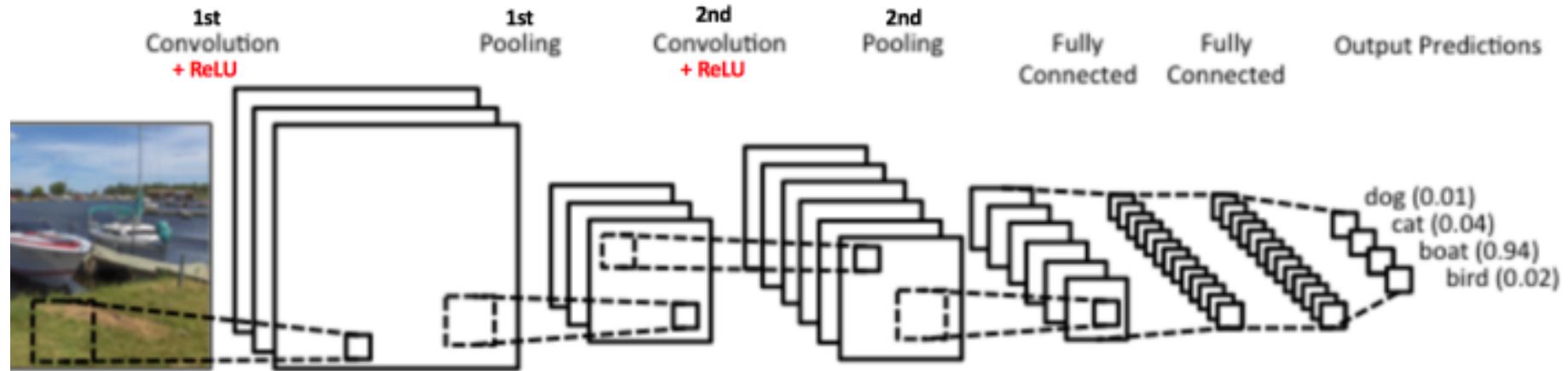
Averaging

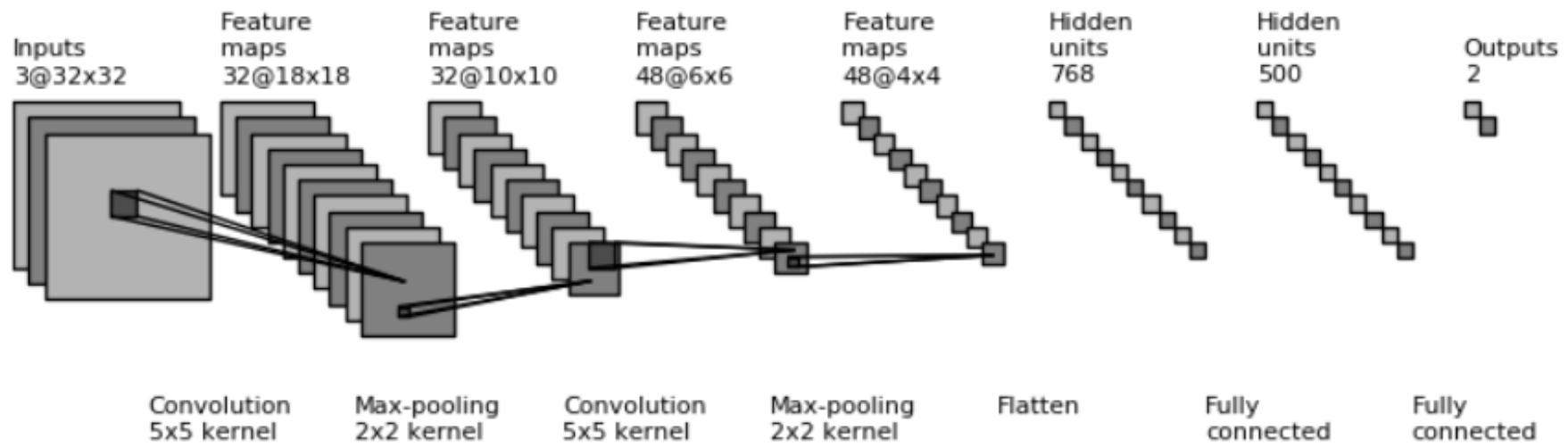
# Visualizations of a deep neural network with CNNs



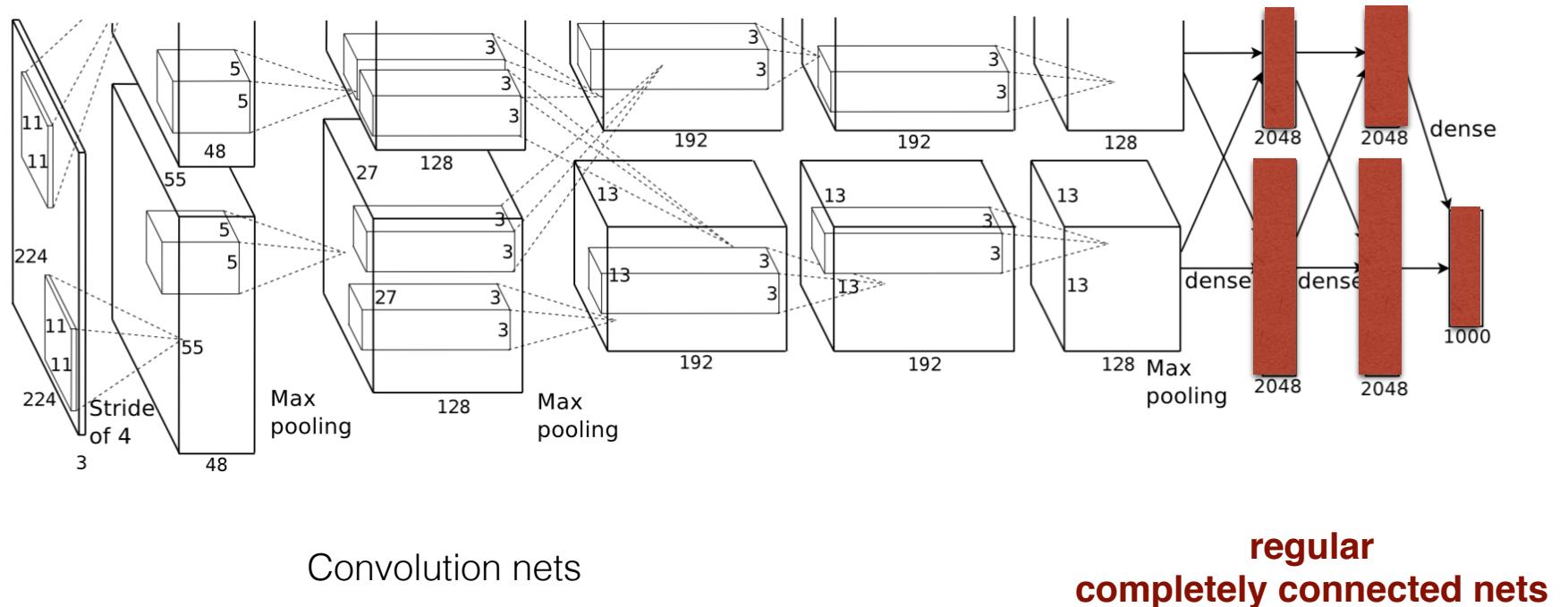
A. W. Harley, "An Interactive Node-Link Visualization of Convolutional Neural Networks," in ISVC, pages 867-877, 2015

## Summary: An example with two CNNs and two fully-connected layers





# A Deep Learning Network handcrafted for image classification



Krizhevsky, A., Sutskever, I. and Hinton, G. E.  
**ImageNet Classification with Deep Convolutional Neural Networks**  
Advances in Neural Information Processing 25, MIT Press, Cambridge, MA