



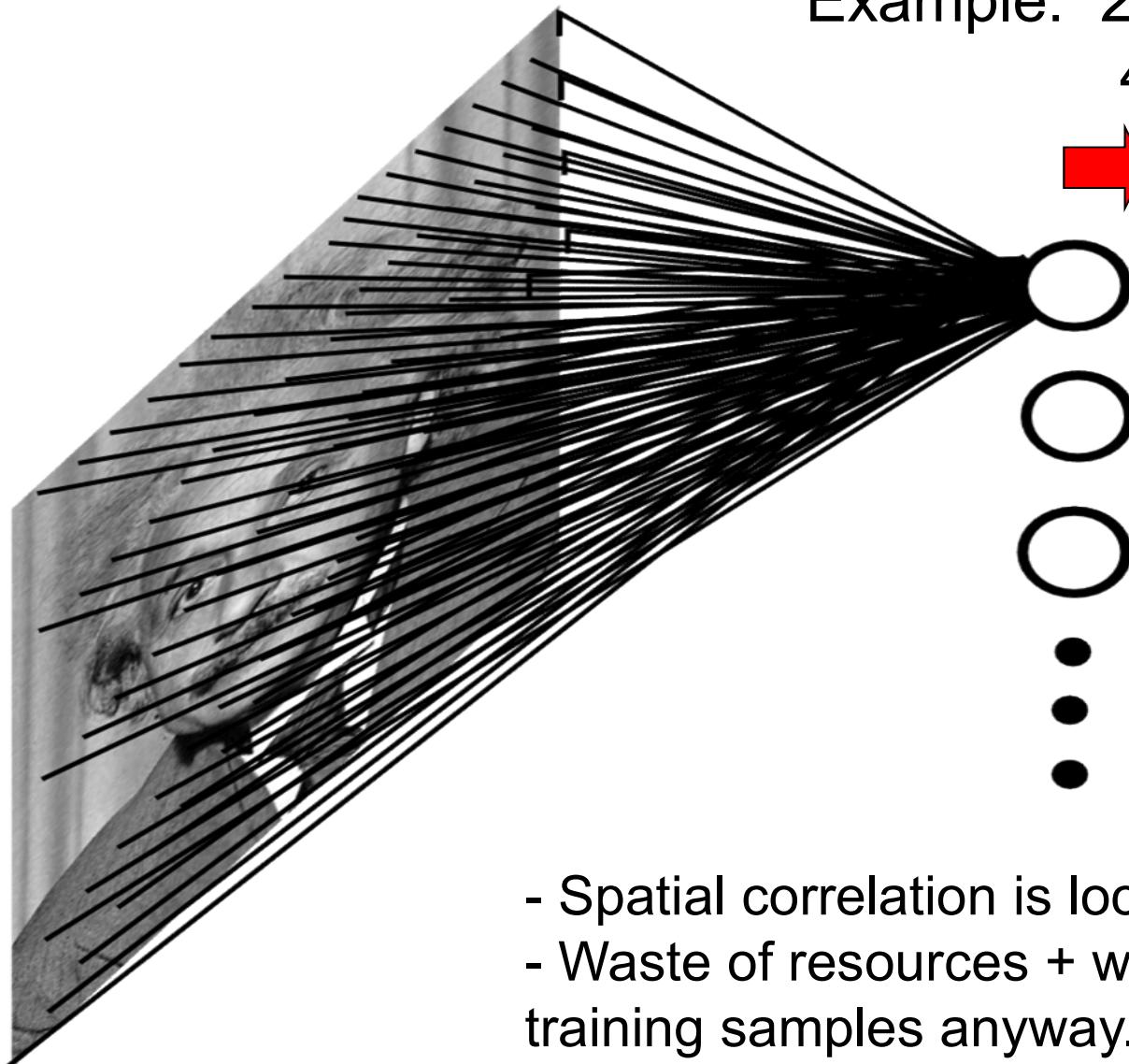
# 고급 소프트웨어 실습(CSE4152)

---

## Convolutional Neural Network

(CNN)

# Fully Connected Layer



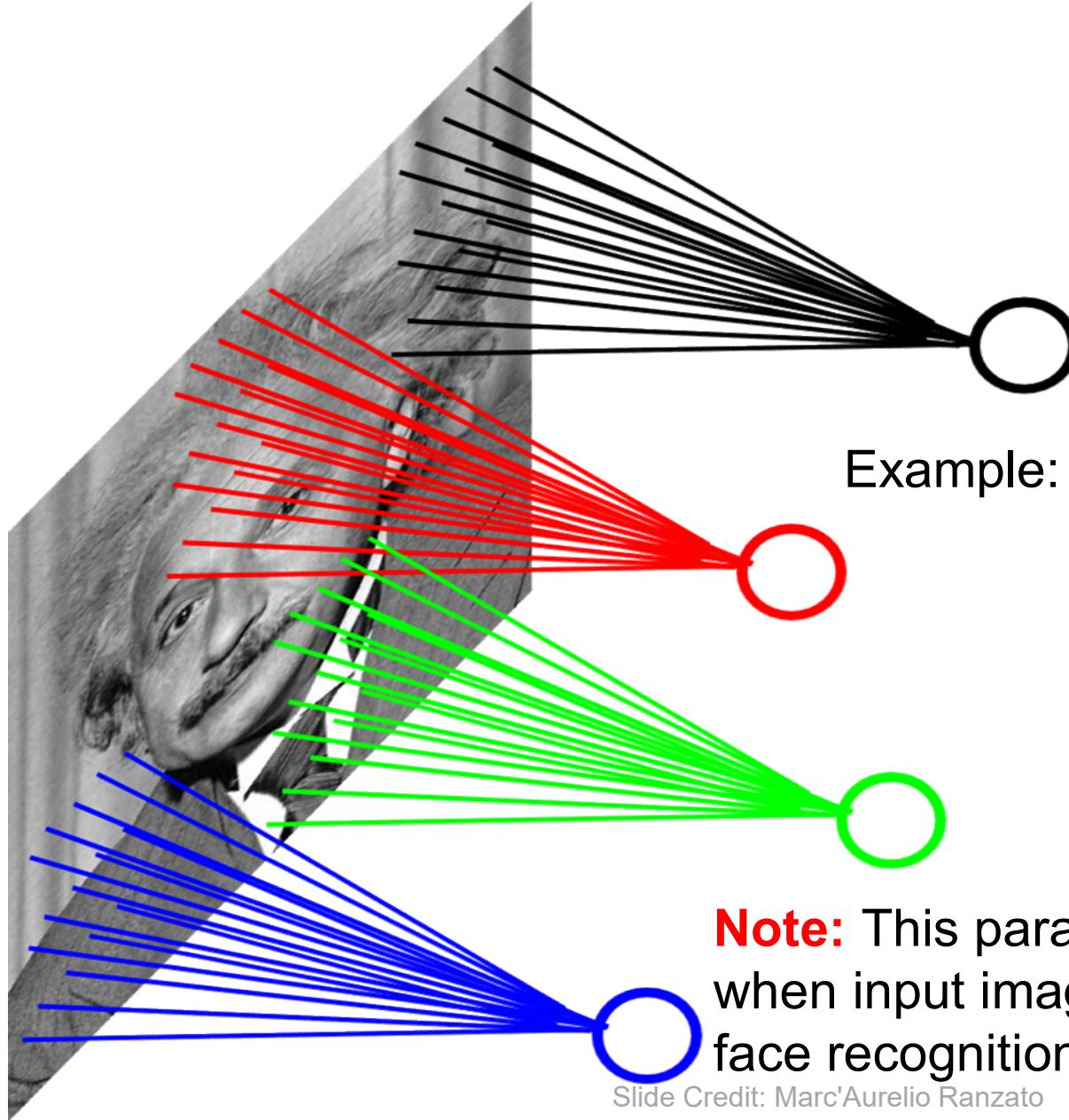
Example: 200x200 image  
40K hidden units

**~2B parameters!!!**

$$40K \times 40K = 1600M = 1.6B$$

- Spatial correlation is local
- Waste of resources + we have not enough training samples anyway..

# Locally Connected Layer

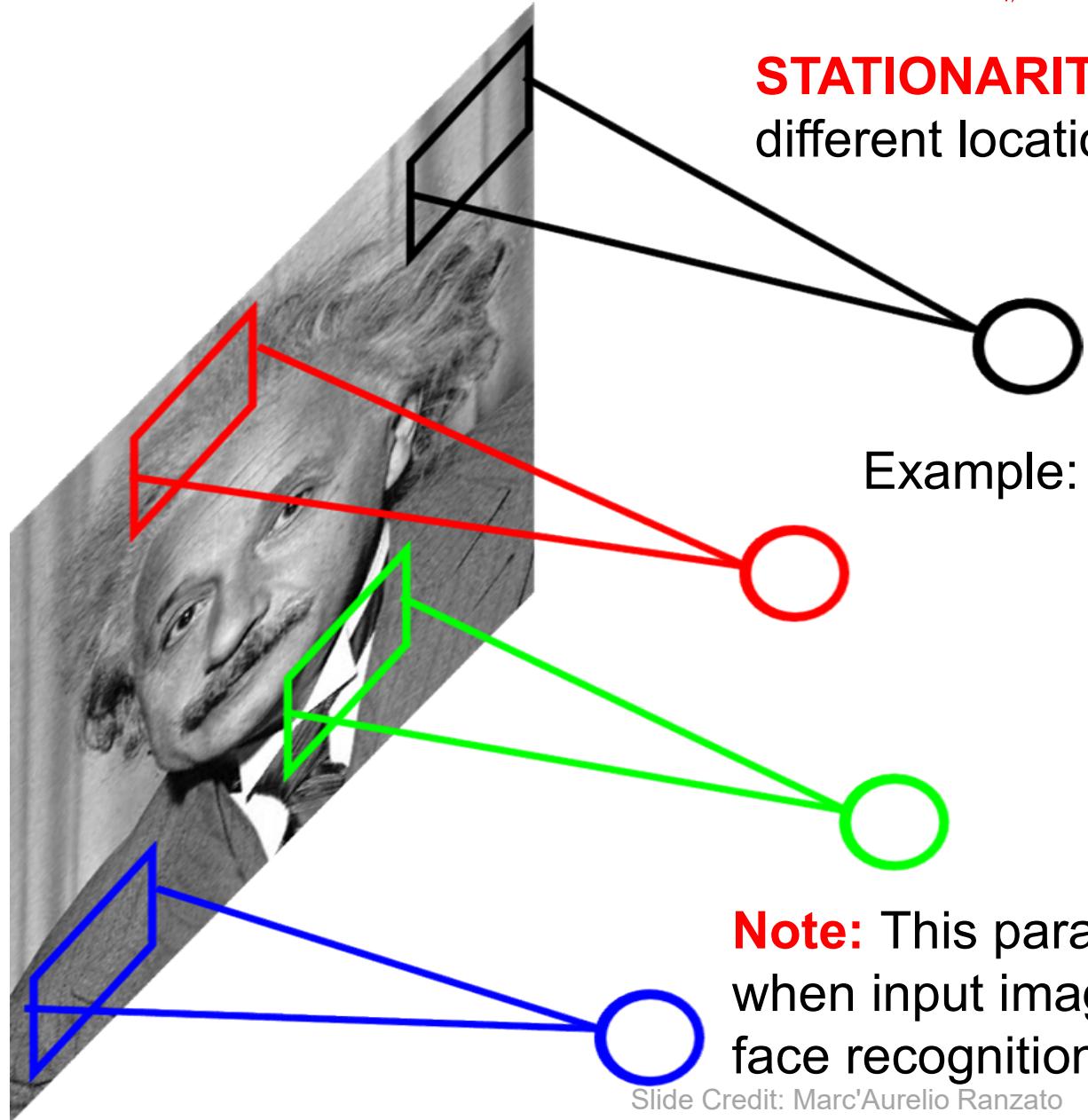


Example:  
200x200 image  
40K hidden units  
Filter size: 10x10  
4M parameters

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

Slide Credit: Marc'Aurelio Ranzato

# Locally Connected Layer



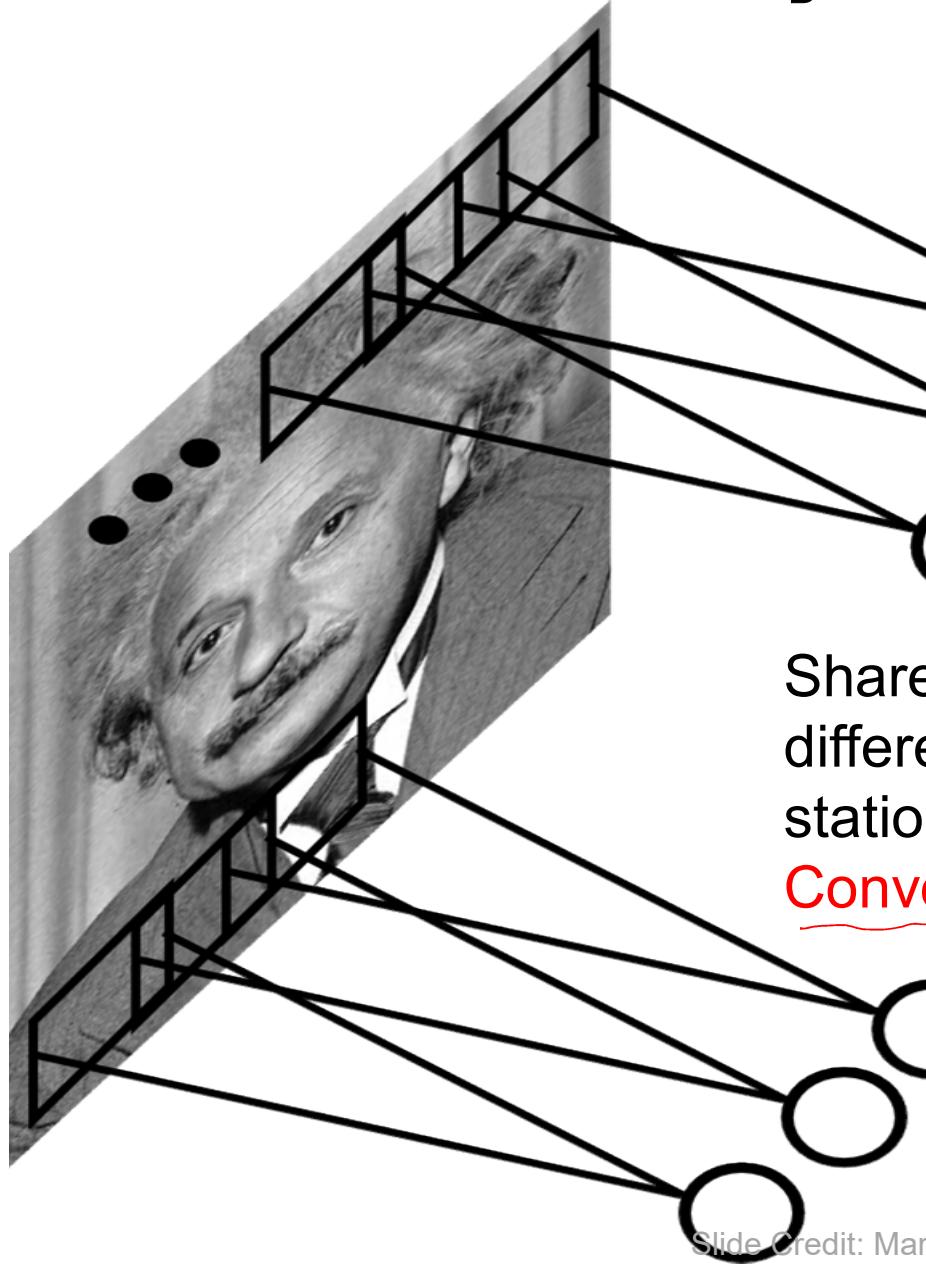
**STATIONARITY?** Statistics is similar at different locations

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).

Slide Credit: Marc'Aurelio Ranzato

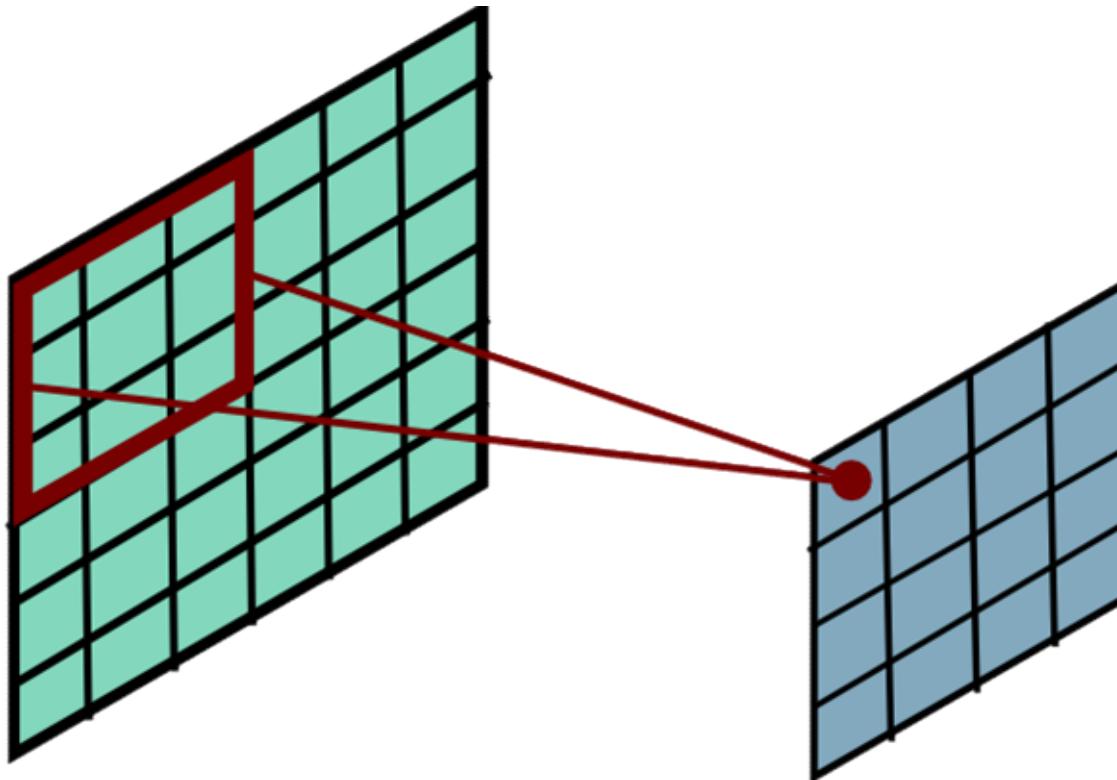
# Convolutional Layer



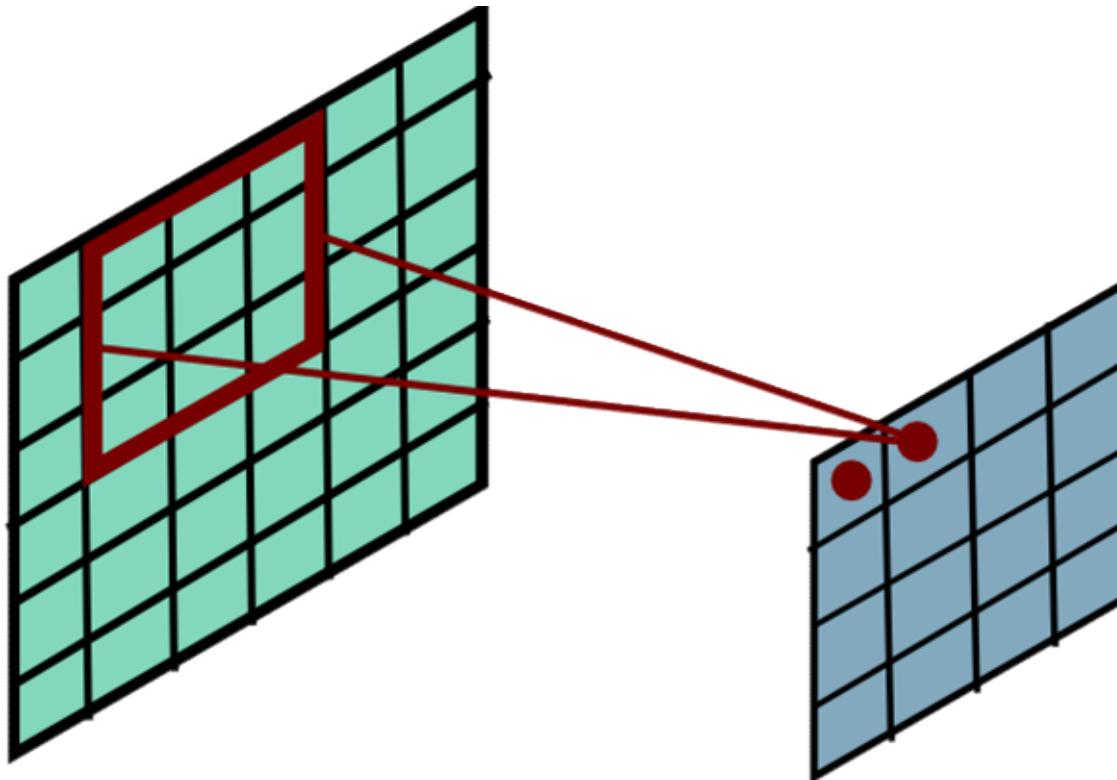
$$\begin{matrix} w_1 & w_2 & w_{13} \\ \cdots & & \cdots \\ & & w_{37} \end{matrix}$$

Share the same parameters across different locations (assuming input is stationary):  
Convolutions with learned kernels

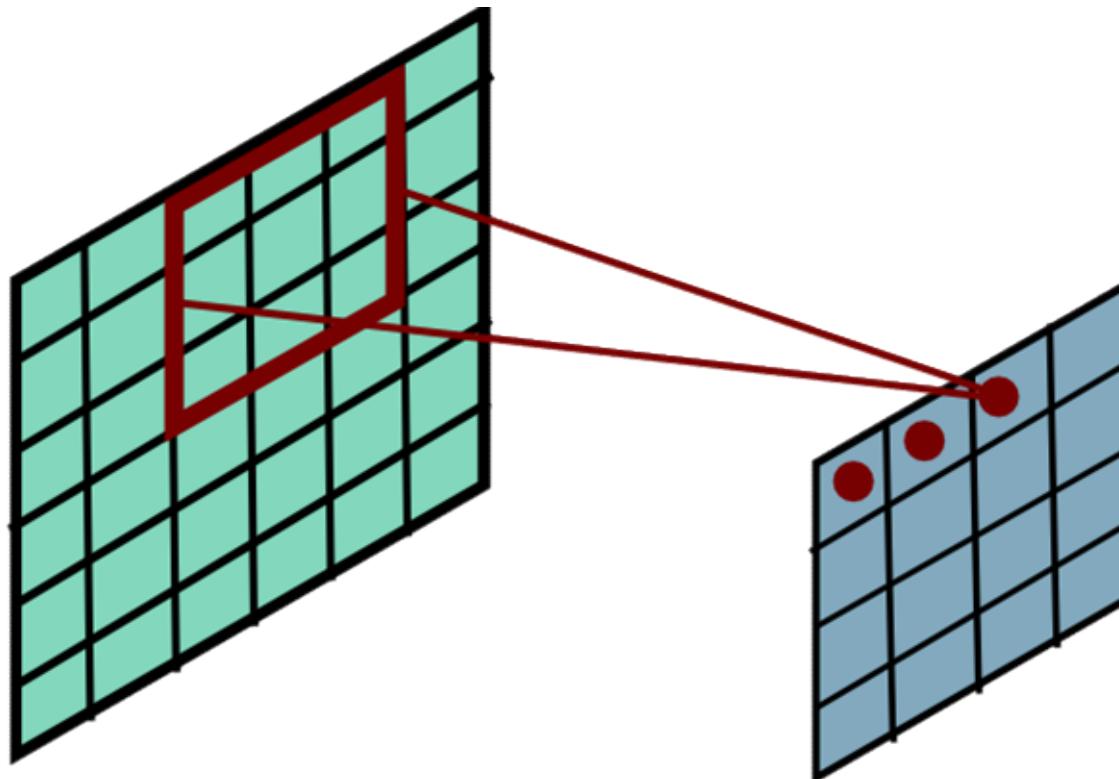
# Convolutional Layer



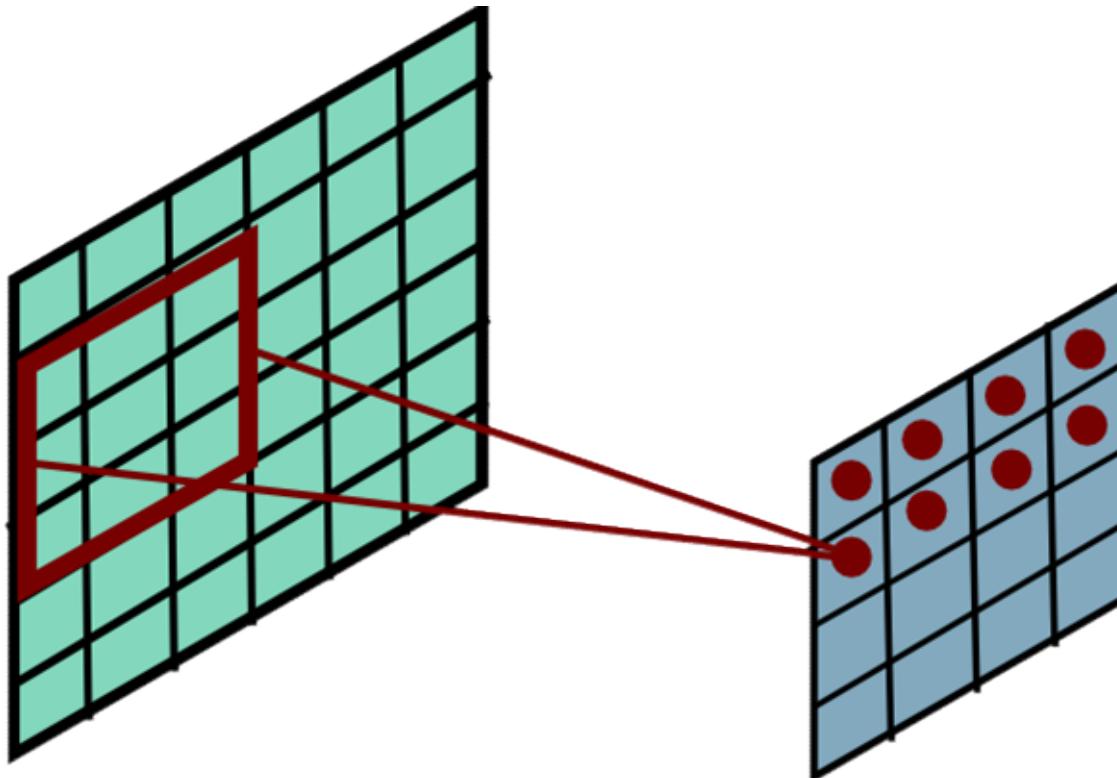
# Convolutional Layer



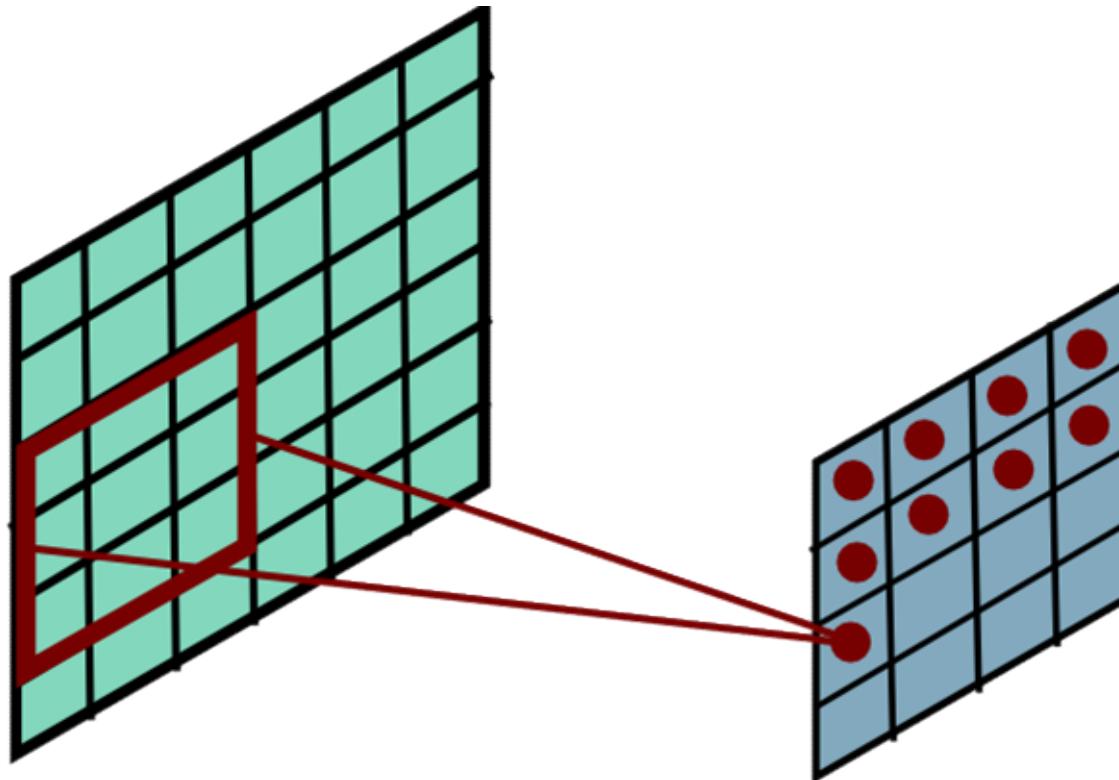
# Convolutional Layer



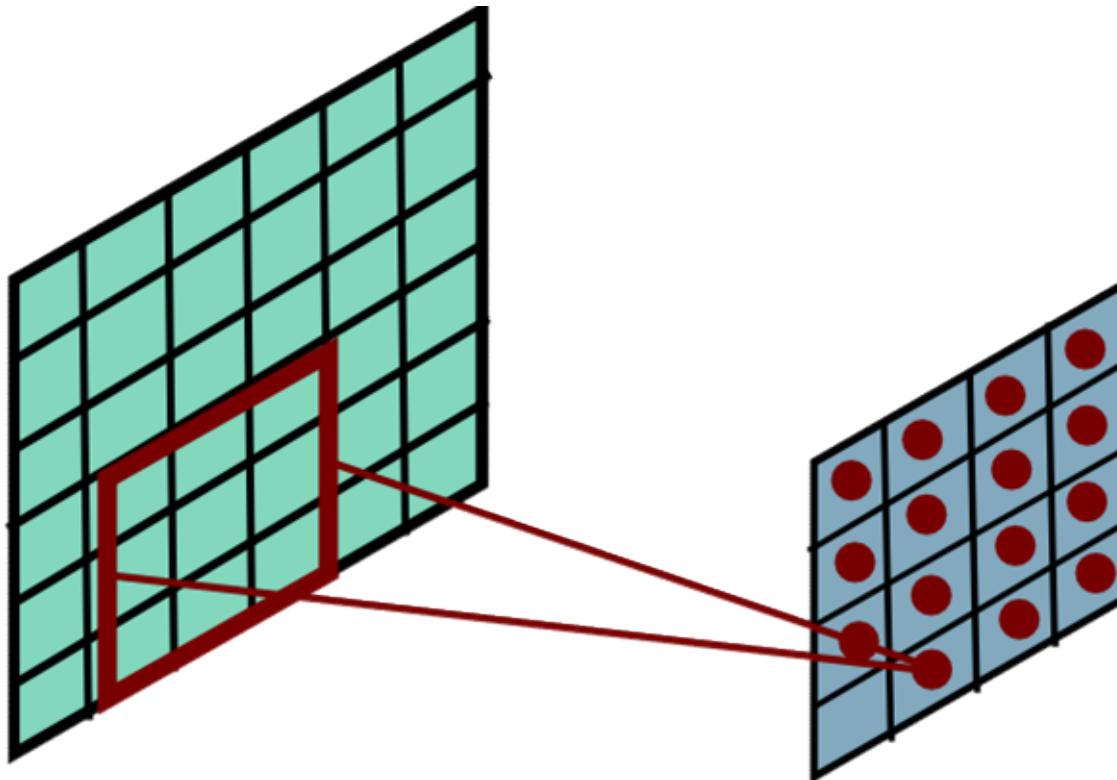
# Convolutional Layer



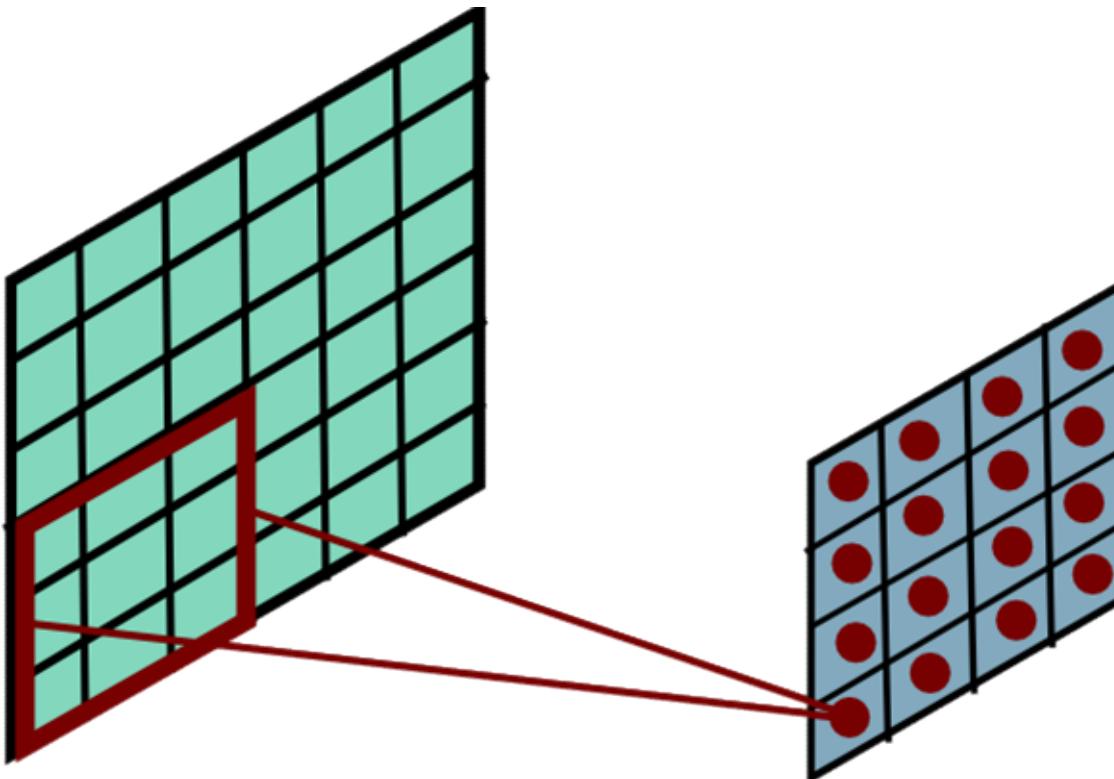
# Convolutional Layer



# Convolutional Layer



# Convolutional Layer

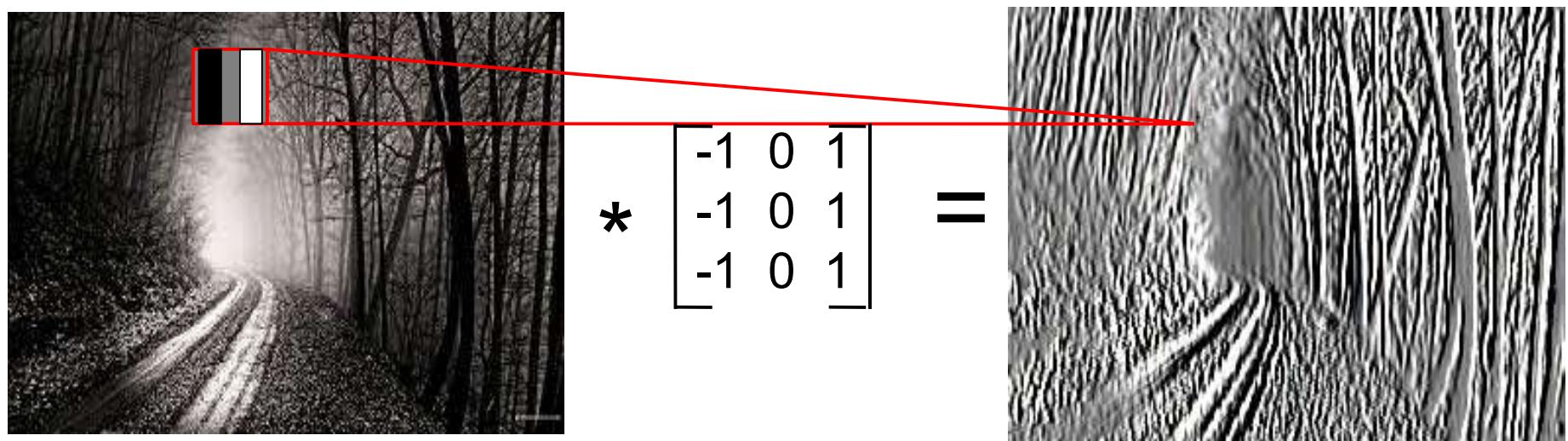


Mathieu et al. "Fast training of CNNs through FFTs" ICLR 2014

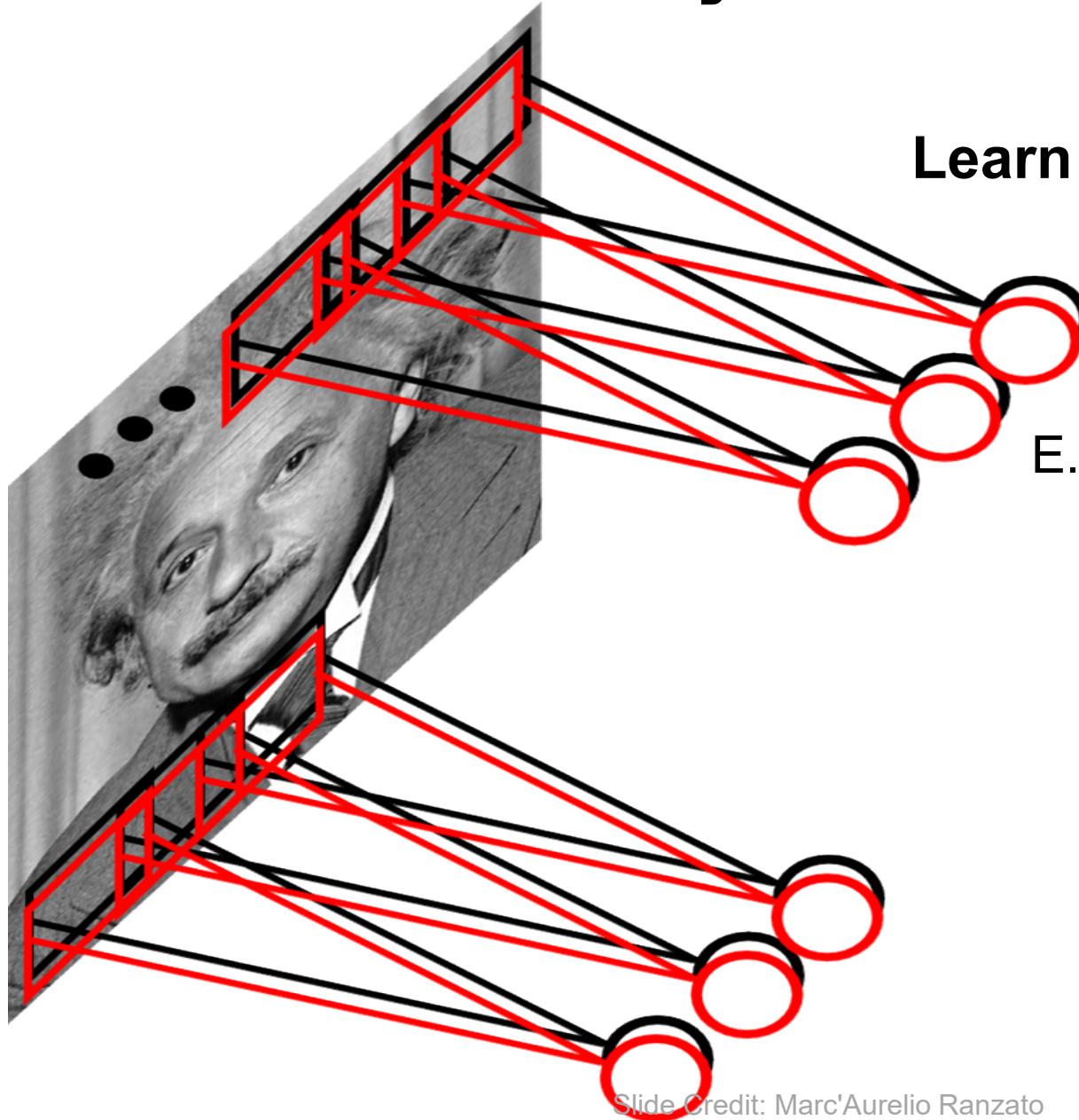
Slide Credit: Marc'Aurelio Ranzato

12

# Convolutional Layer



# Convolutional Layer



**Learn multiple filters.**

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

# Convolutional Neural Networks

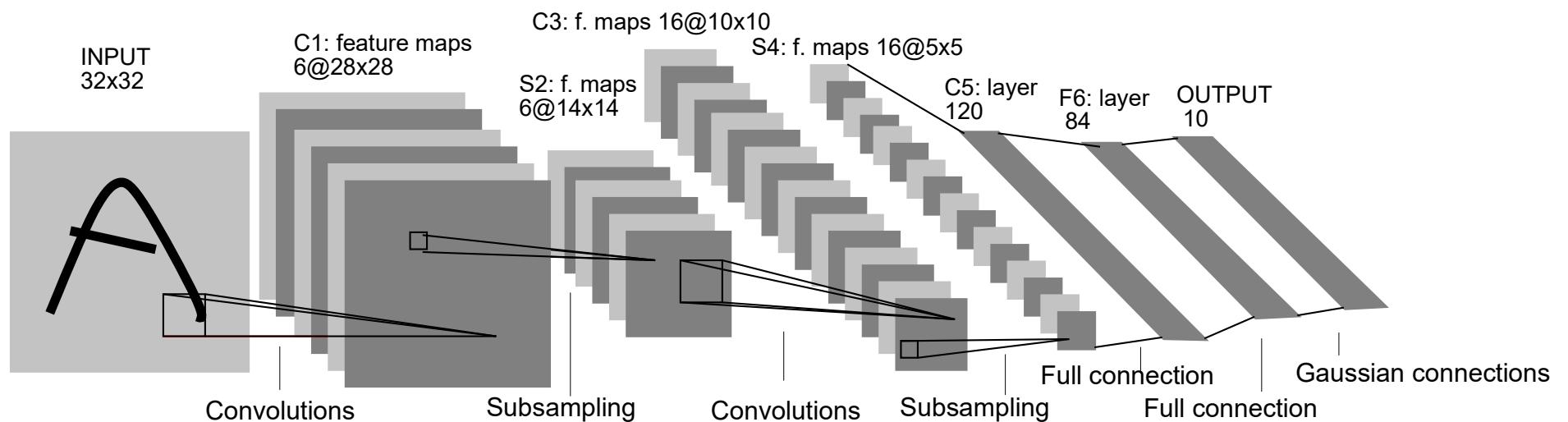
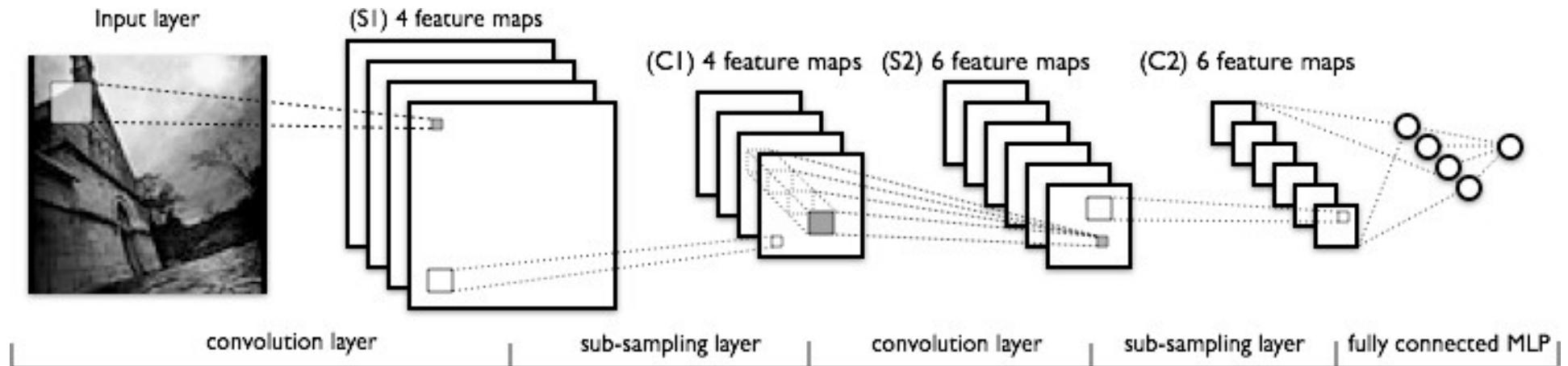
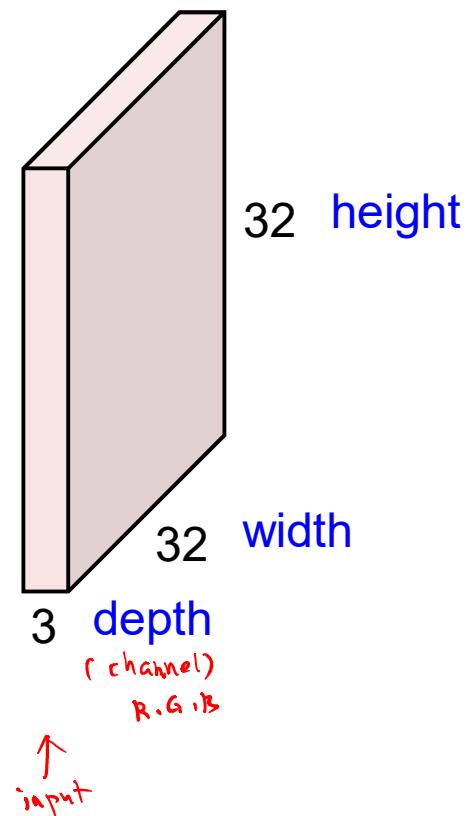


Image Credit: Yann LeCun, Kevin Murphy

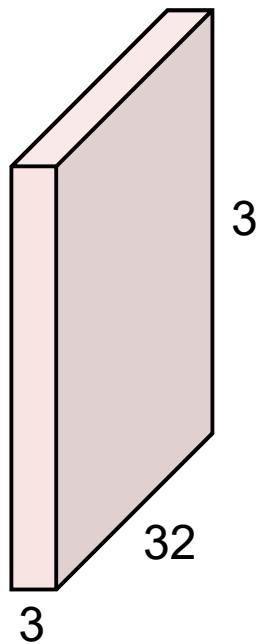
# Convolution Layer

32x32x3 image -> preserve spatial structure

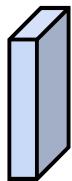


# Convolution Layer

32x32x3 image



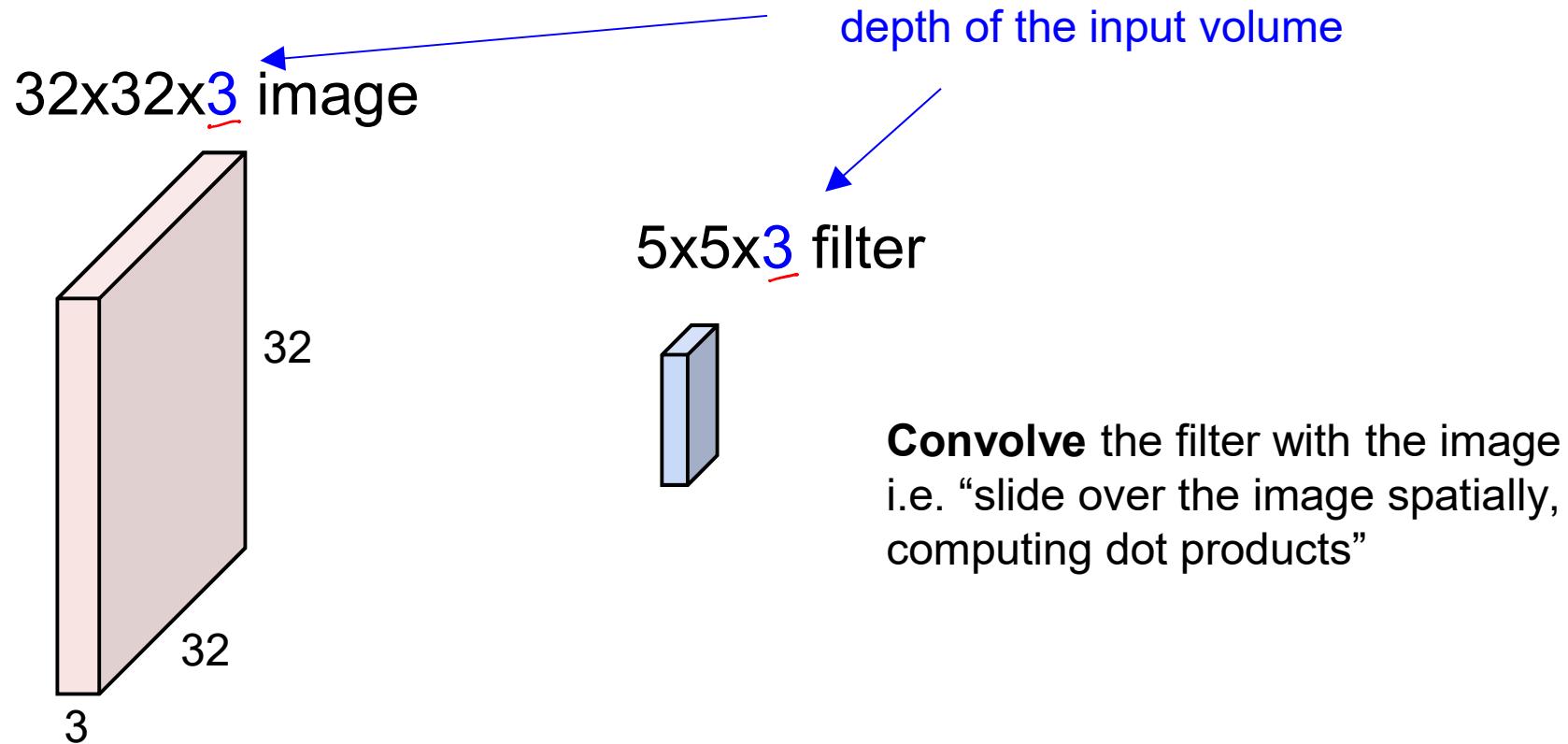
5x5x3 filter



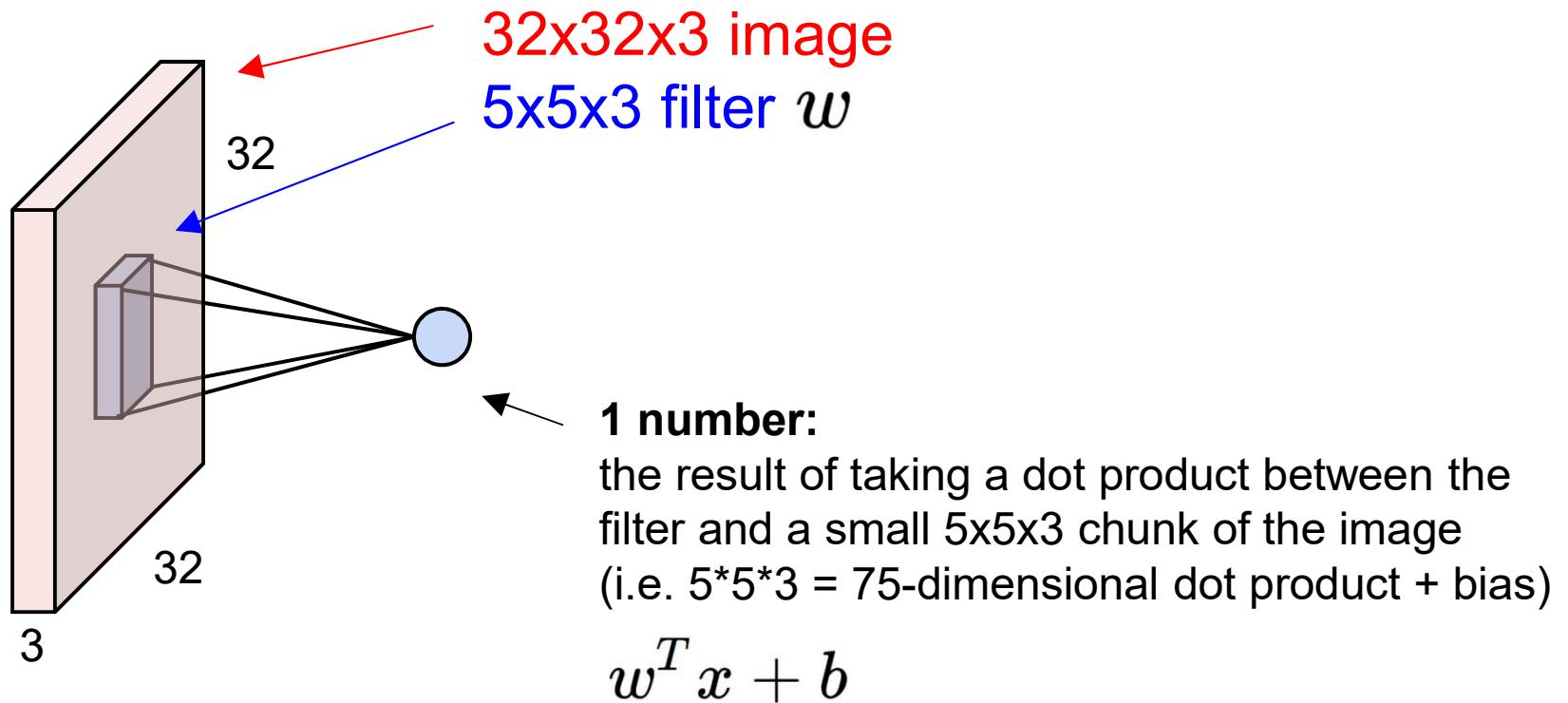
**Convolve** the filter with the image  
i.e. “slide over the image spatially,  
computing dot products”

2D convolution  
3D convolution ~~~ 3D convolution in NN  
conv3D( ) ~~~ w1|w2|w3

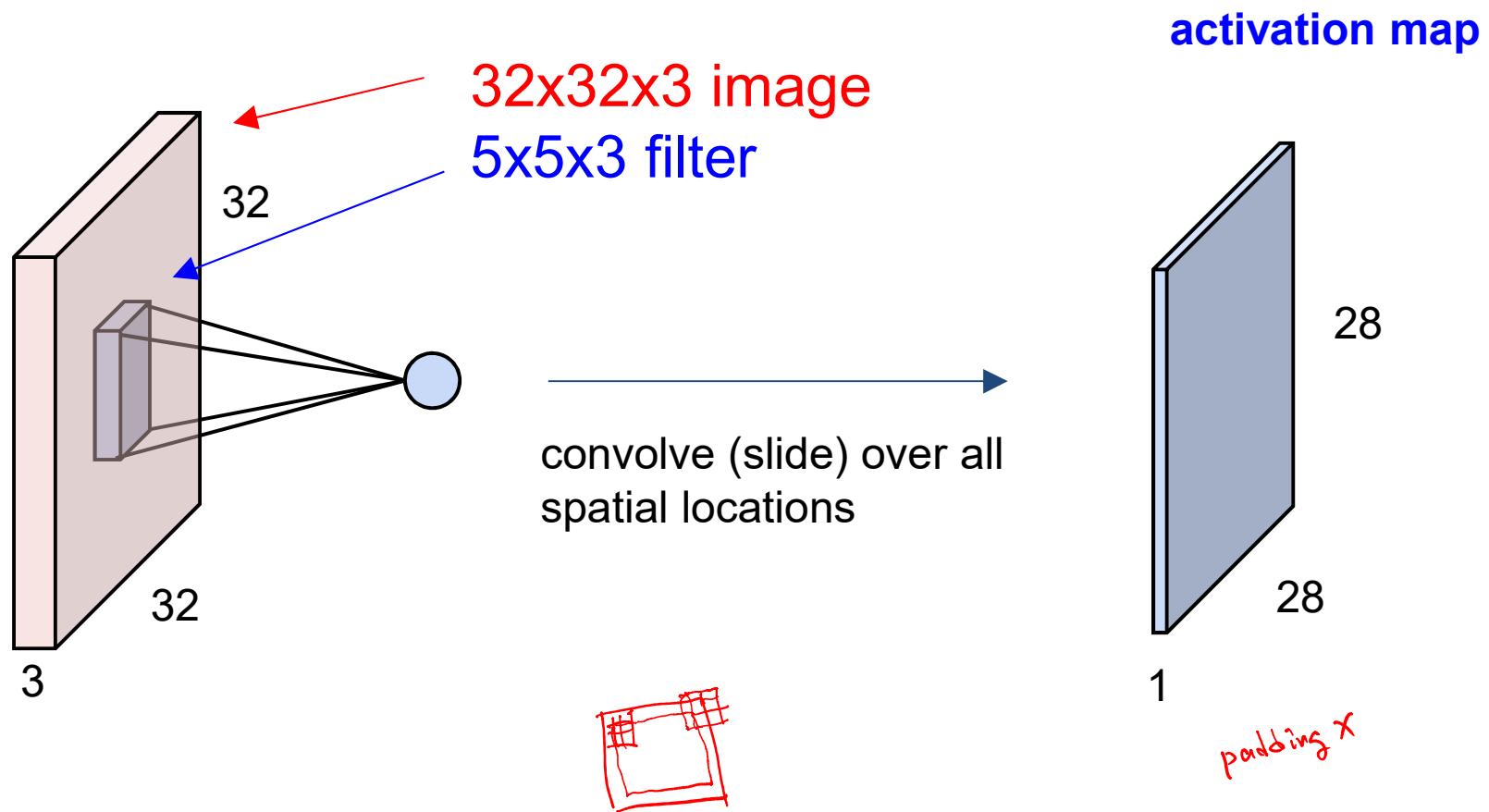
# Convolution Layer



# Convolution Layer

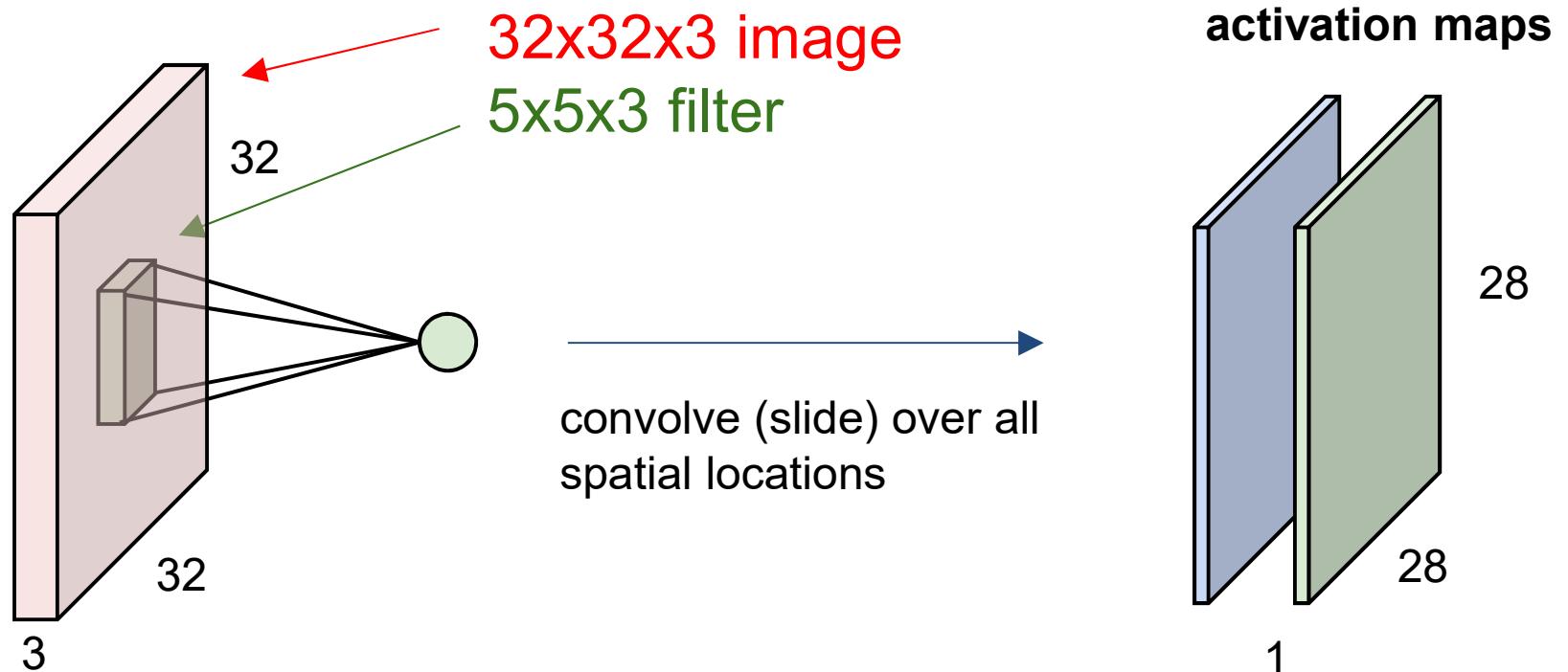


# Convolution Layer



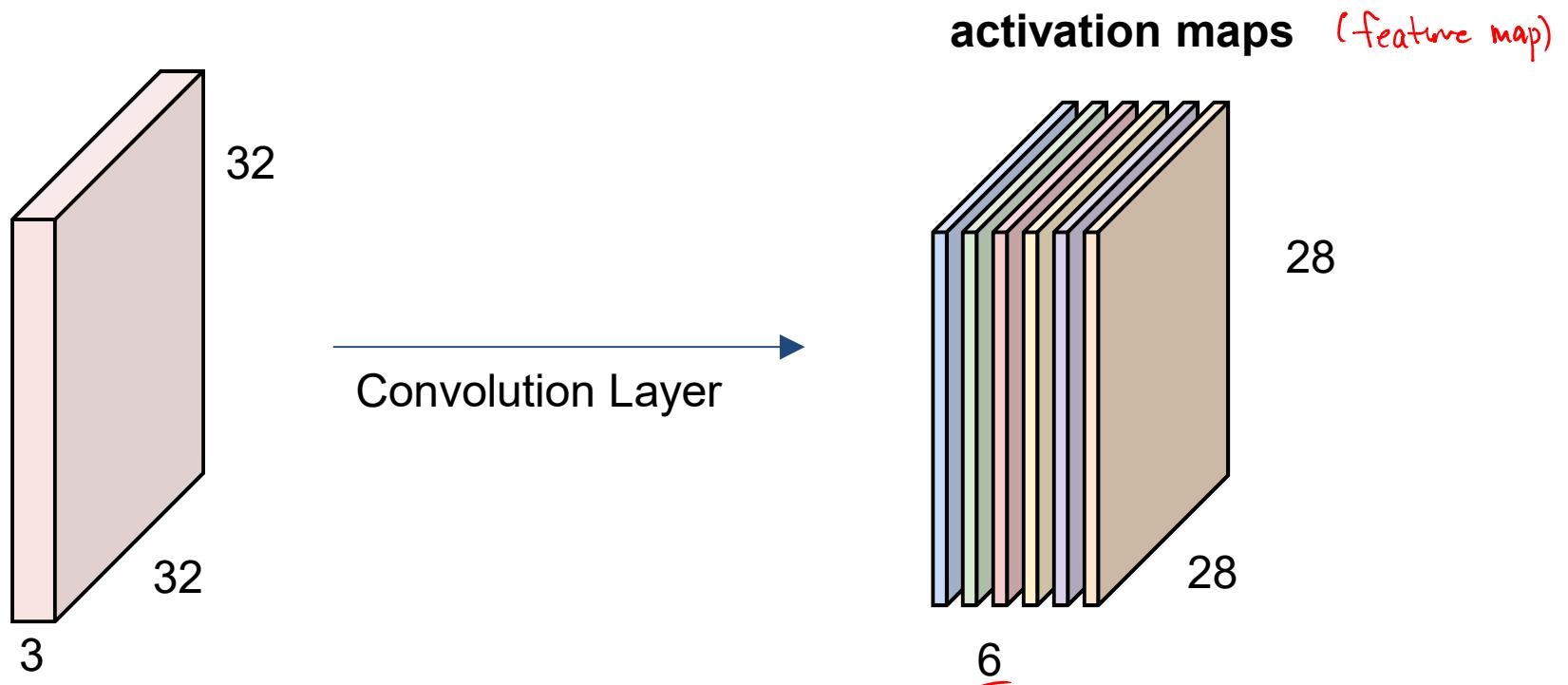
# Convolution Layer

consider a second, green filter



# Convolution Layer

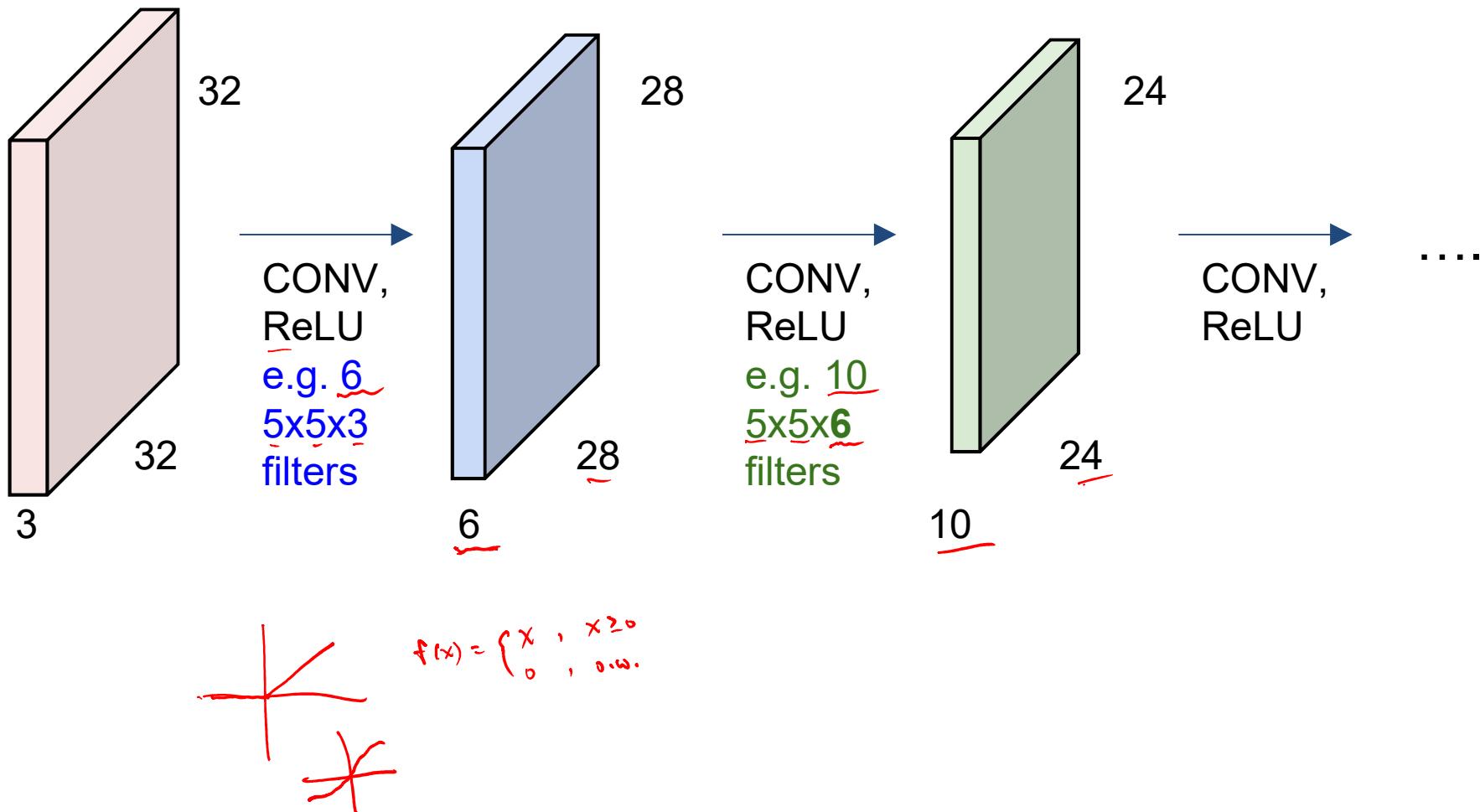
For example, if we had 6  $5 \times 5 \times 3$  filters, we'll get 6 separate activation maps:



We stack these up to get a “new image” of size  $28 \times 28 \times 6$ !

# Convolution Layer

**Preview:** ConvNet is a sequence of Convolutional Layers, interspersed with activation functions

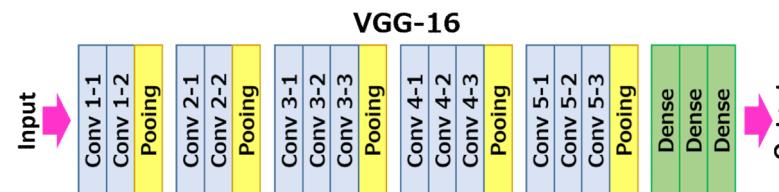
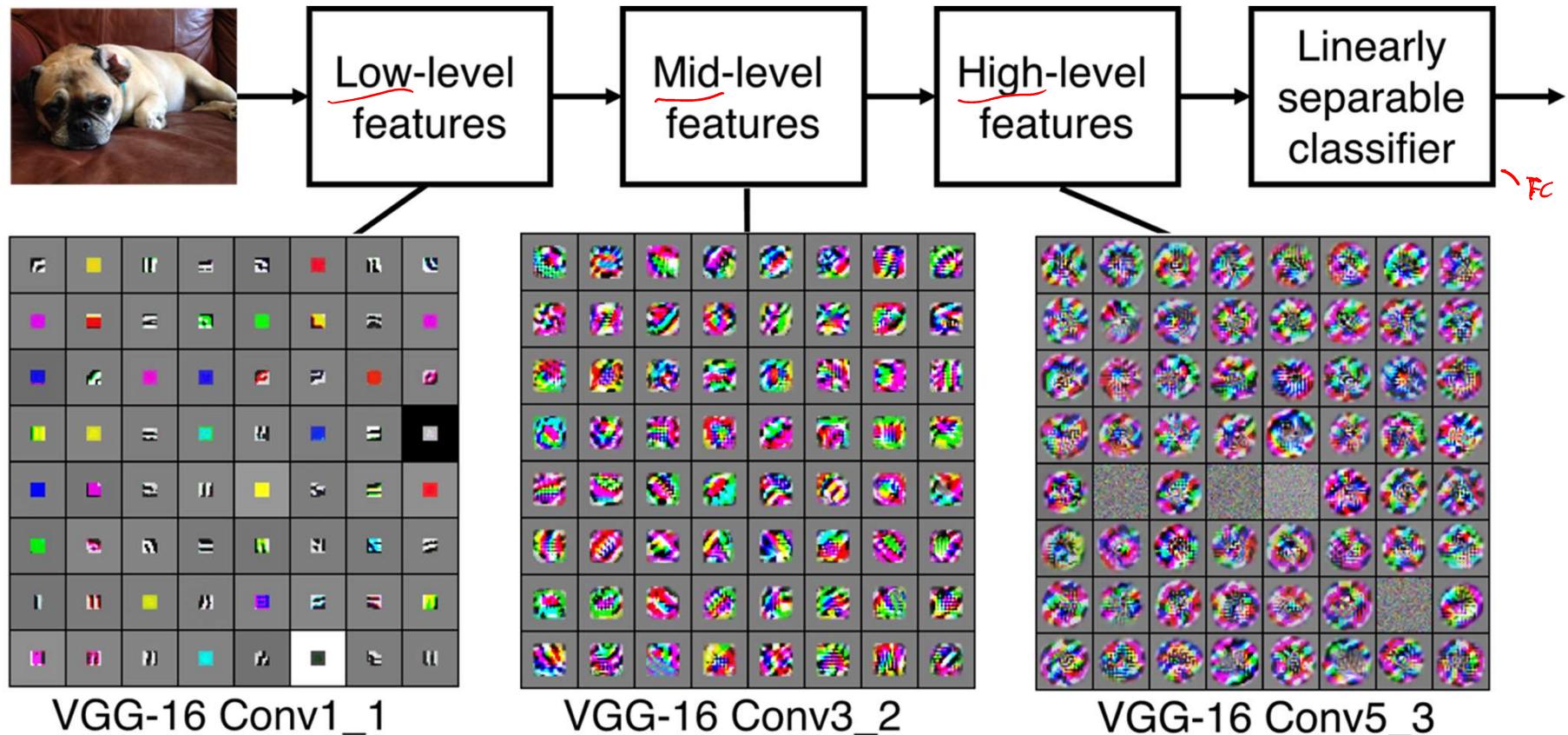


# Convolutional Layer

## Preview

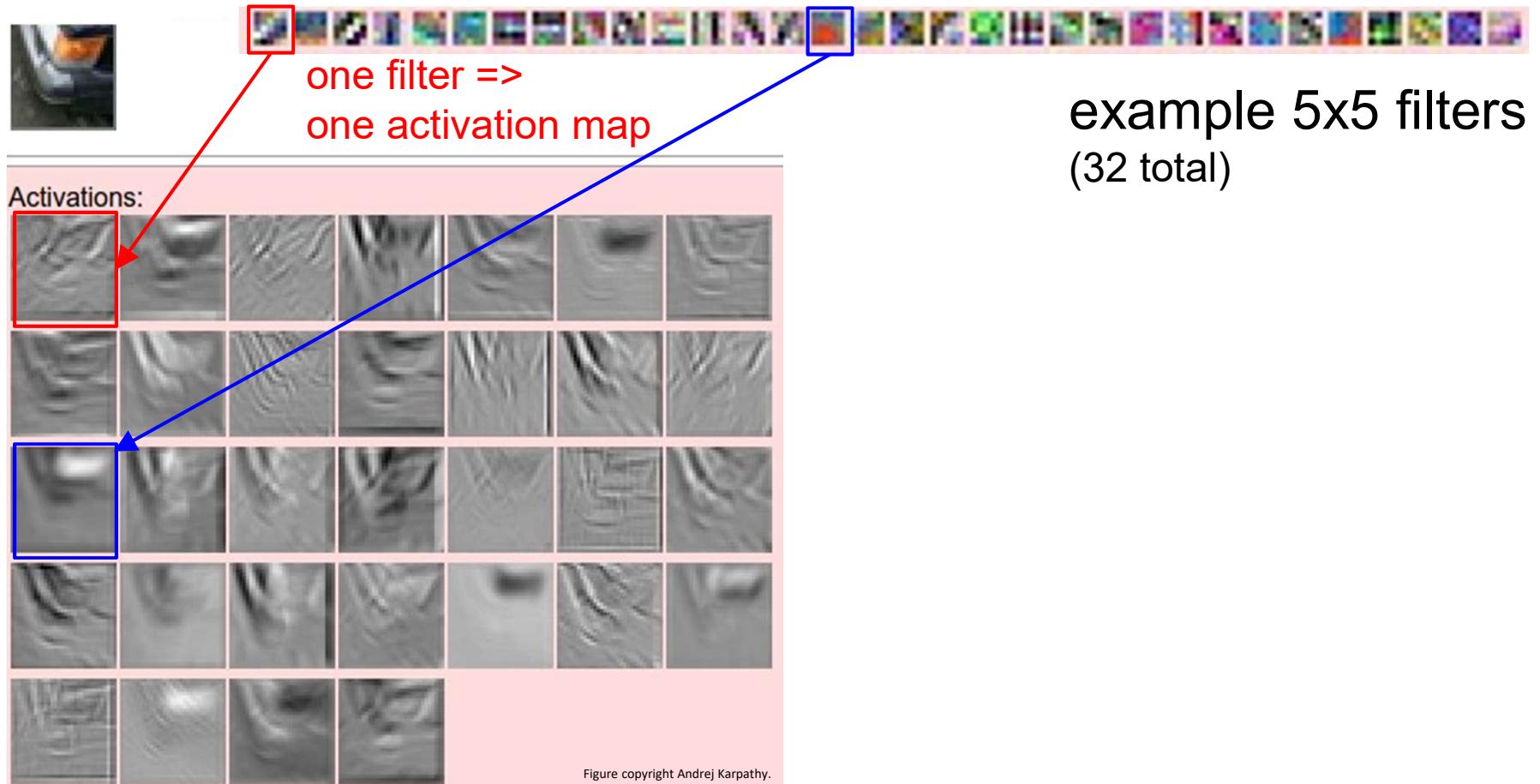
[Zeiler and Fergus 2013]

Visualization of VGG-16 by Lane McIntosh. VGG-16 architecture from [Simonyan and Zisserman 2014].



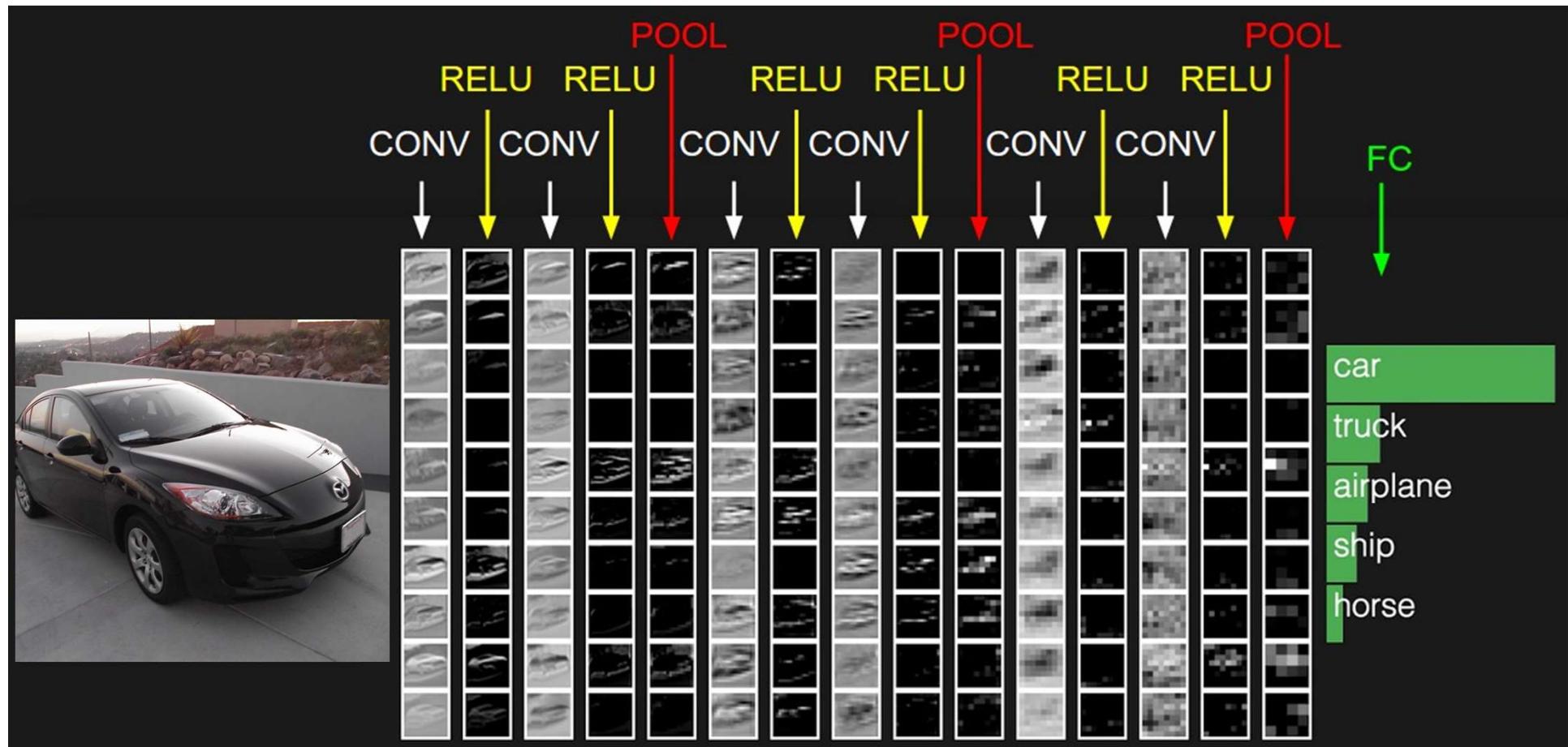
Slide Credit: Fei-Fei Li, Justin Johnson, Serena Yeung, CS 231n

# Convolutional Layer



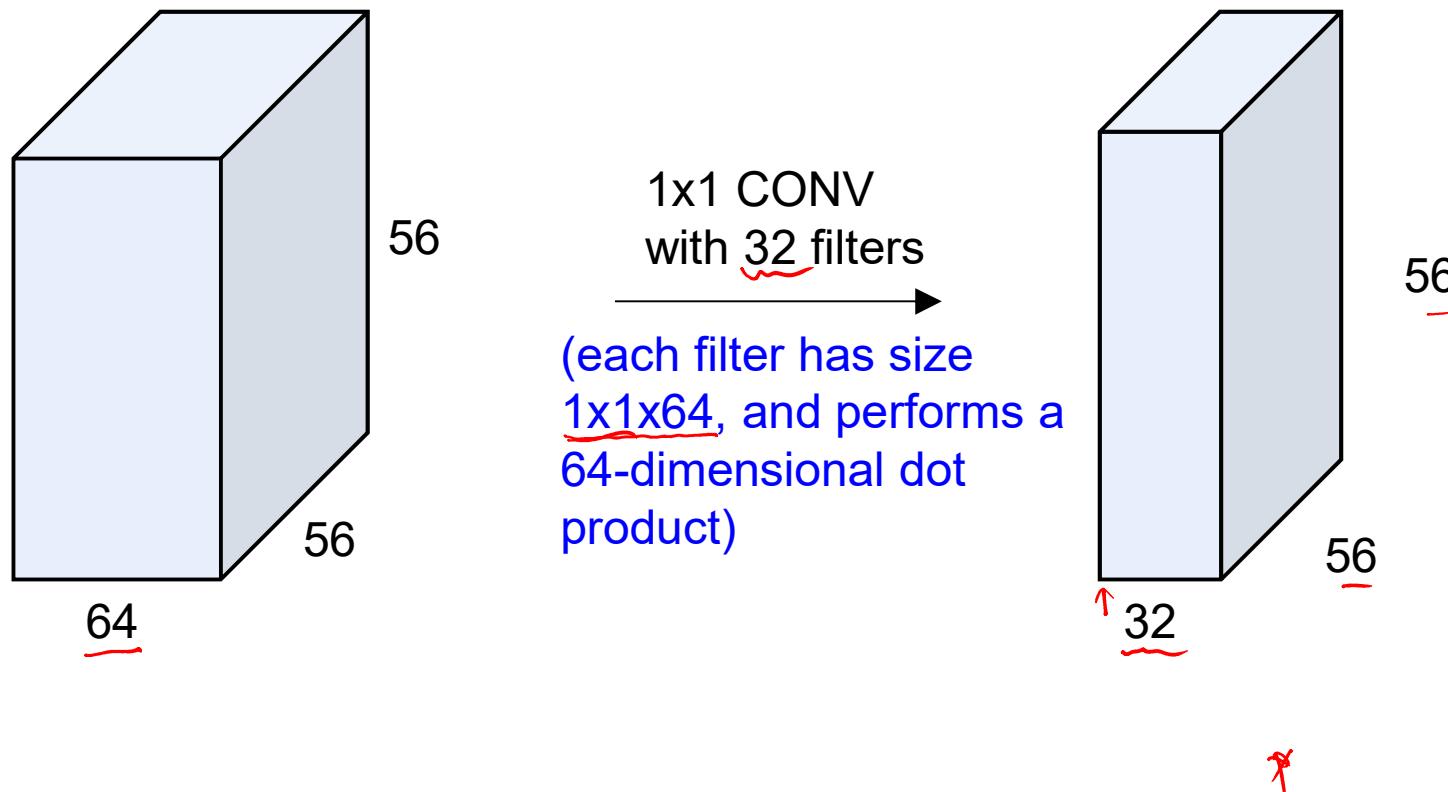
# Convolutional Layer

preview:



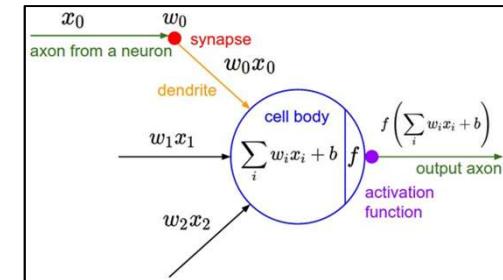
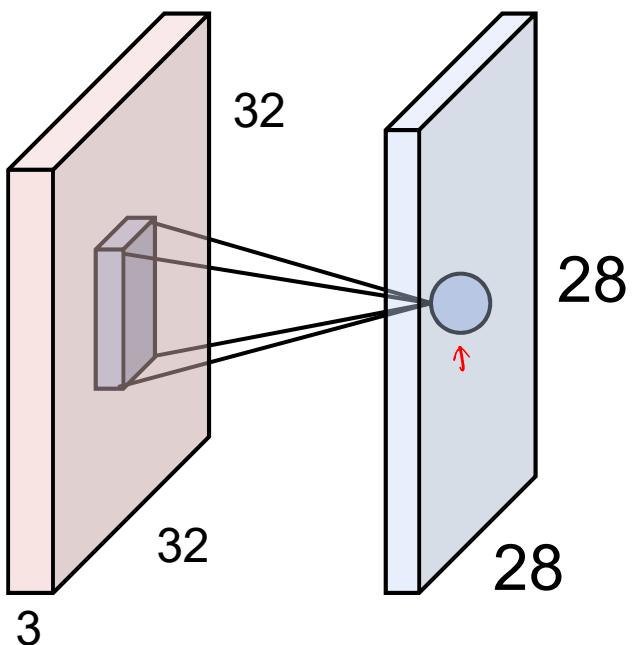
# Convolutional Layer

(btw, 1x1 convolution layers make perfect sense)



# Convolutional Layer

The brain/neuron view of CONV Layer

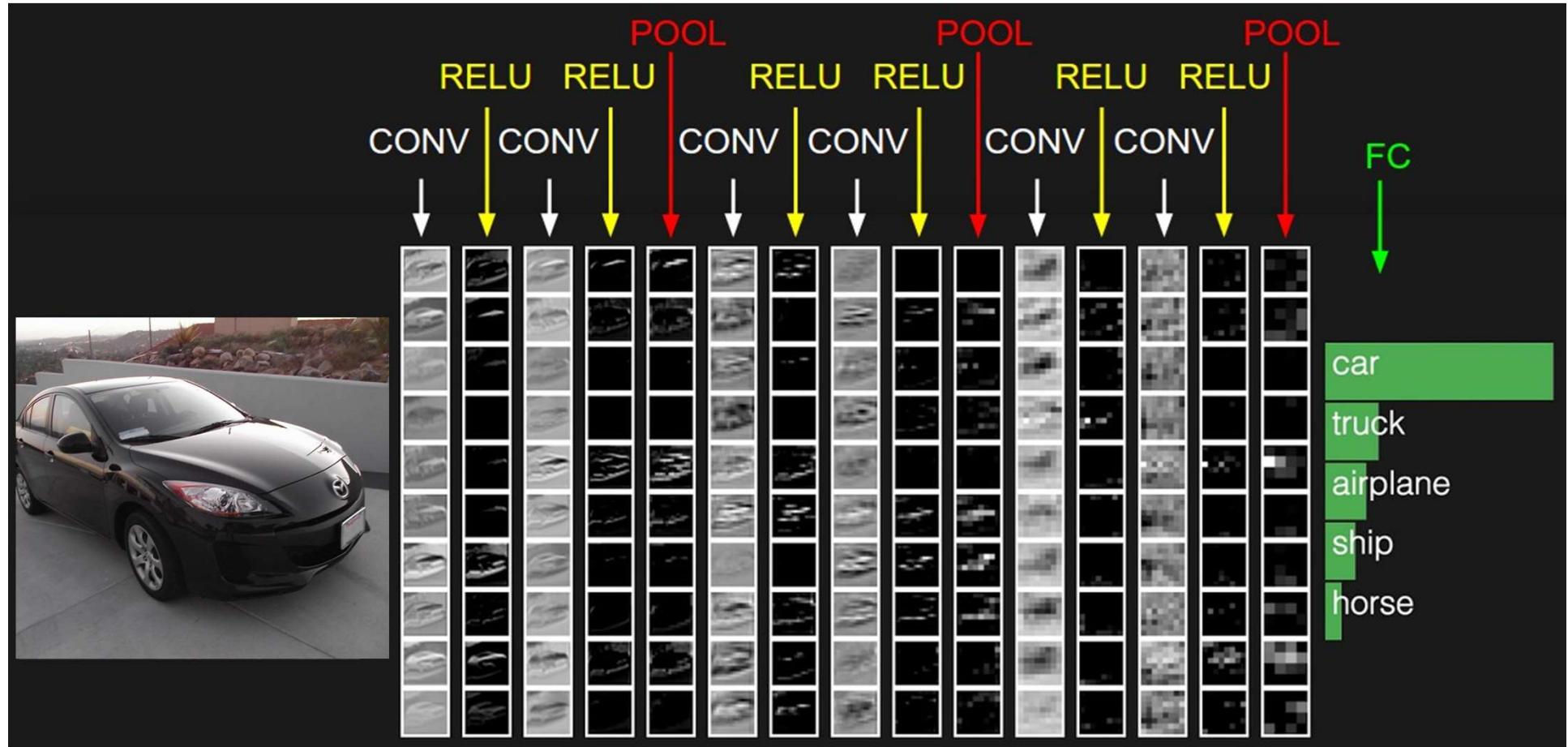


An activation map is a 28x28 sheet of neuron outputs:

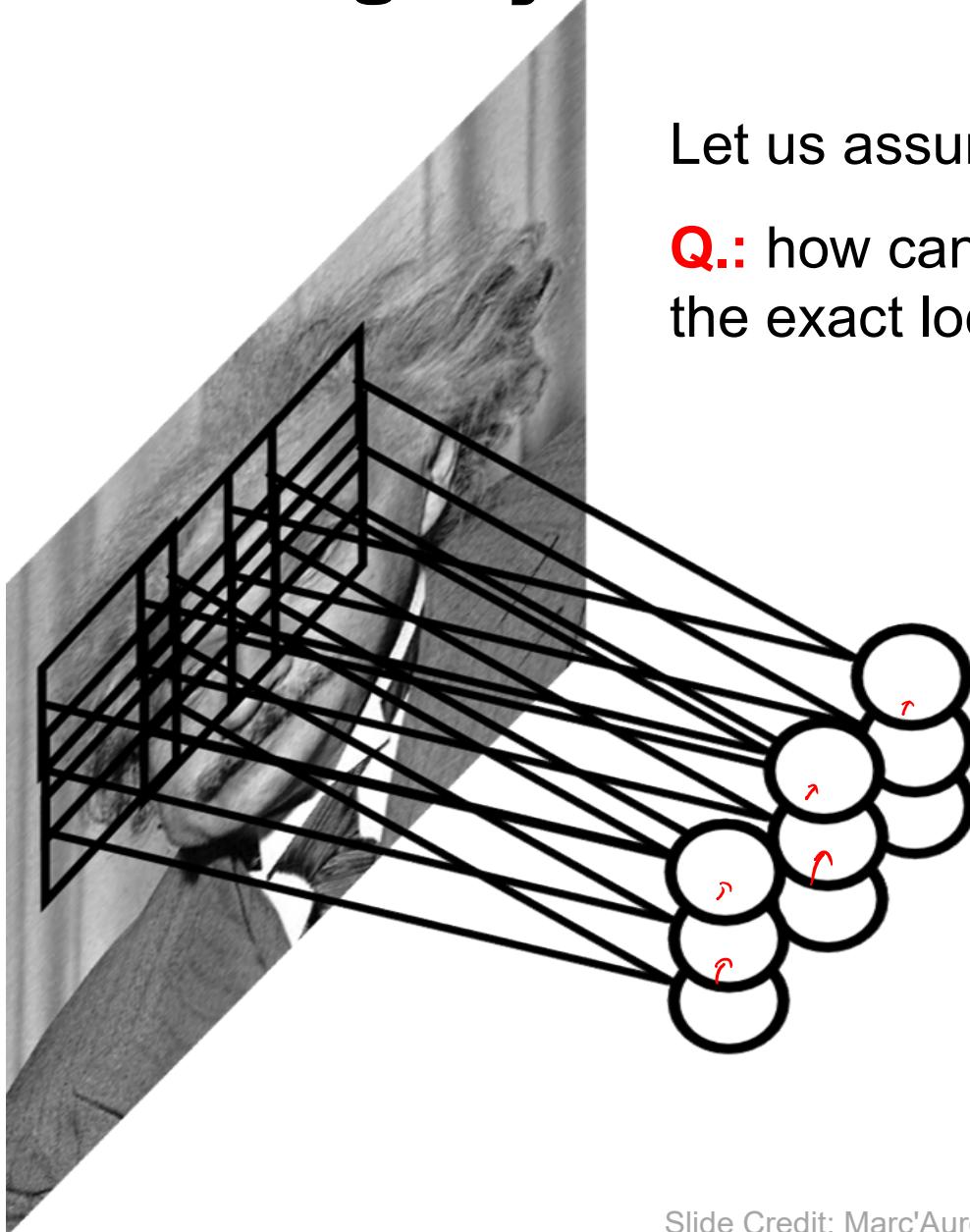
1. Each is connected to a small region in the input
2. All of them share parameters

“5x5 filter” -> “5x5 receptive field for each neuron”

# Two more layers to go: POOL/FC



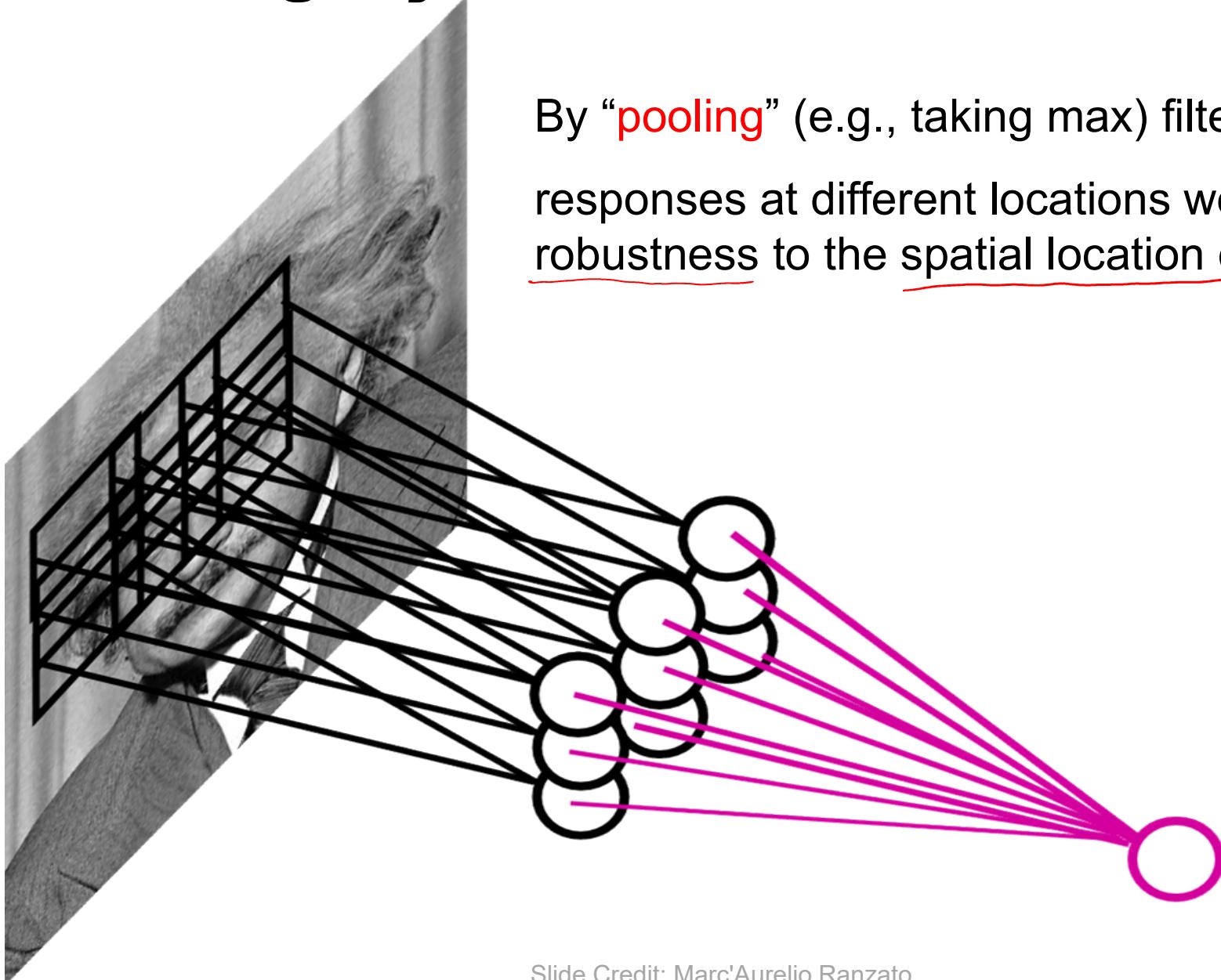
# Pooling Layer



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

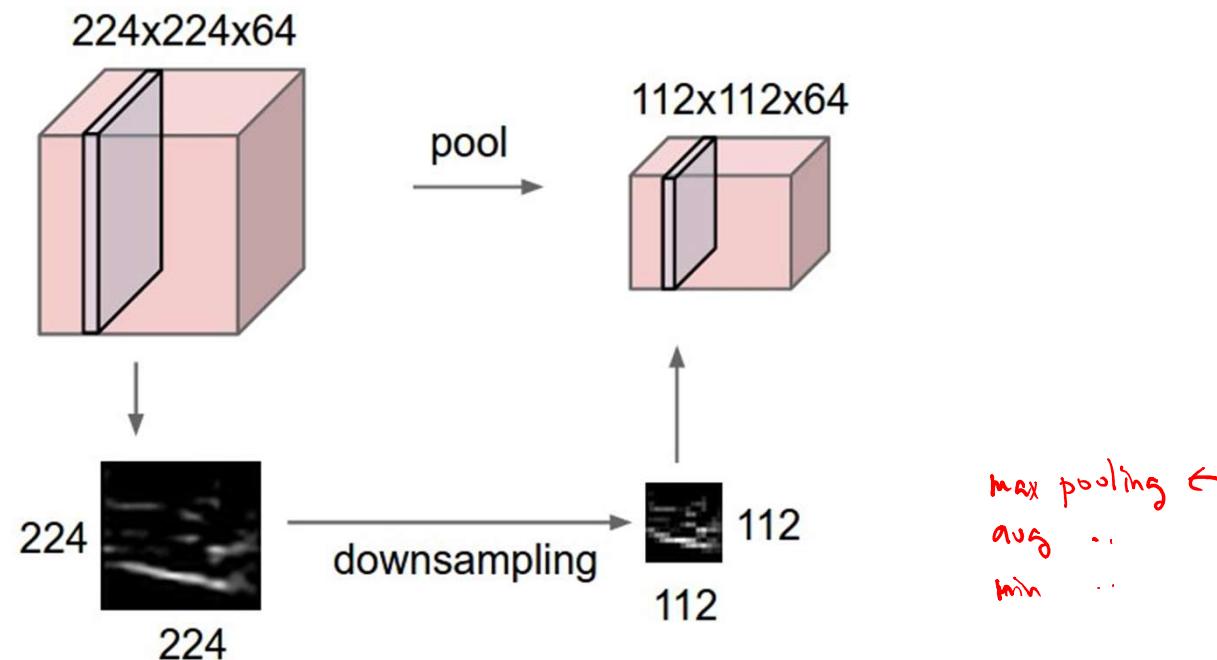


Slide Credit: Marc'Aurelio Ranzato

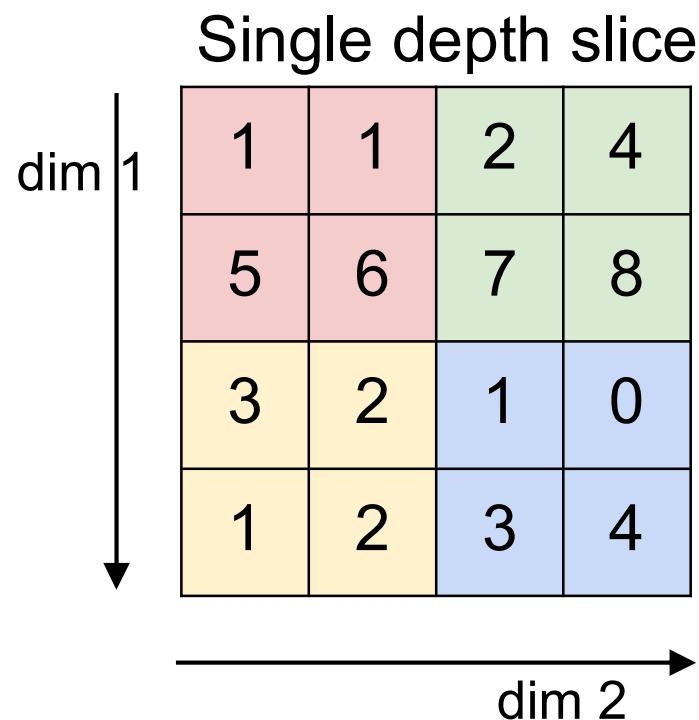
31

# Pooling layer

- makes the representations smaller and more manageable
- operates over each activation map independently:



# MAX POOLING

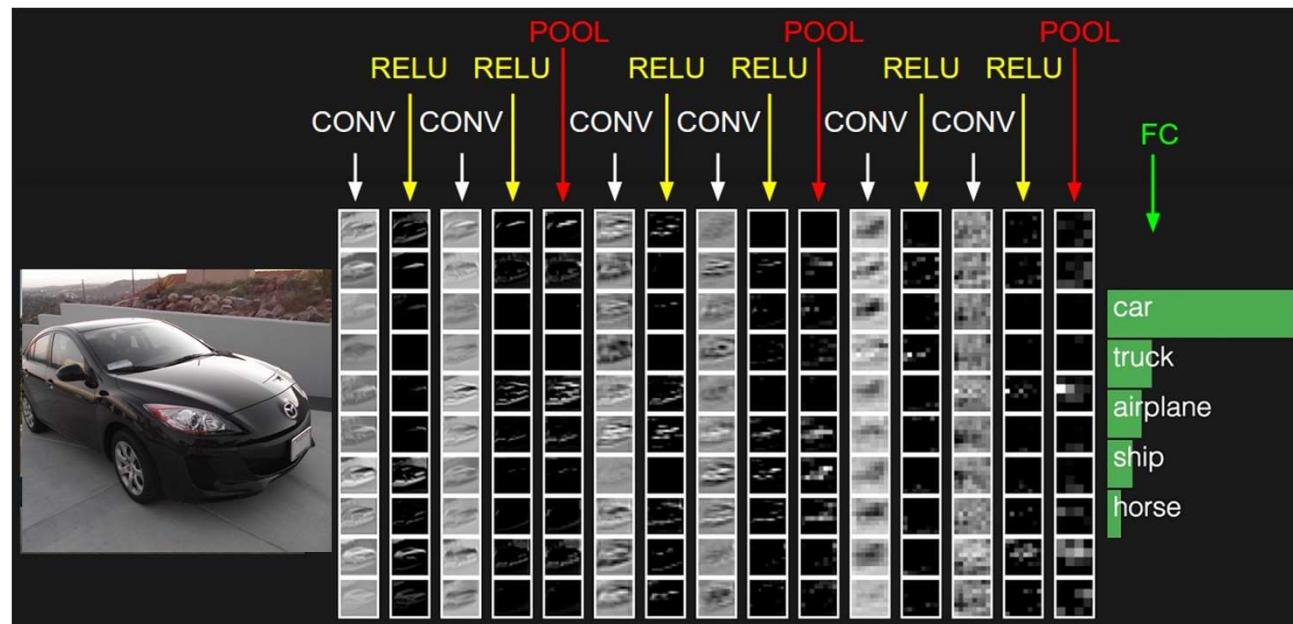


max pool with 2x2 filters  
and stride 2

6	8
3	4

# Fully Connected Layer (FC layer)

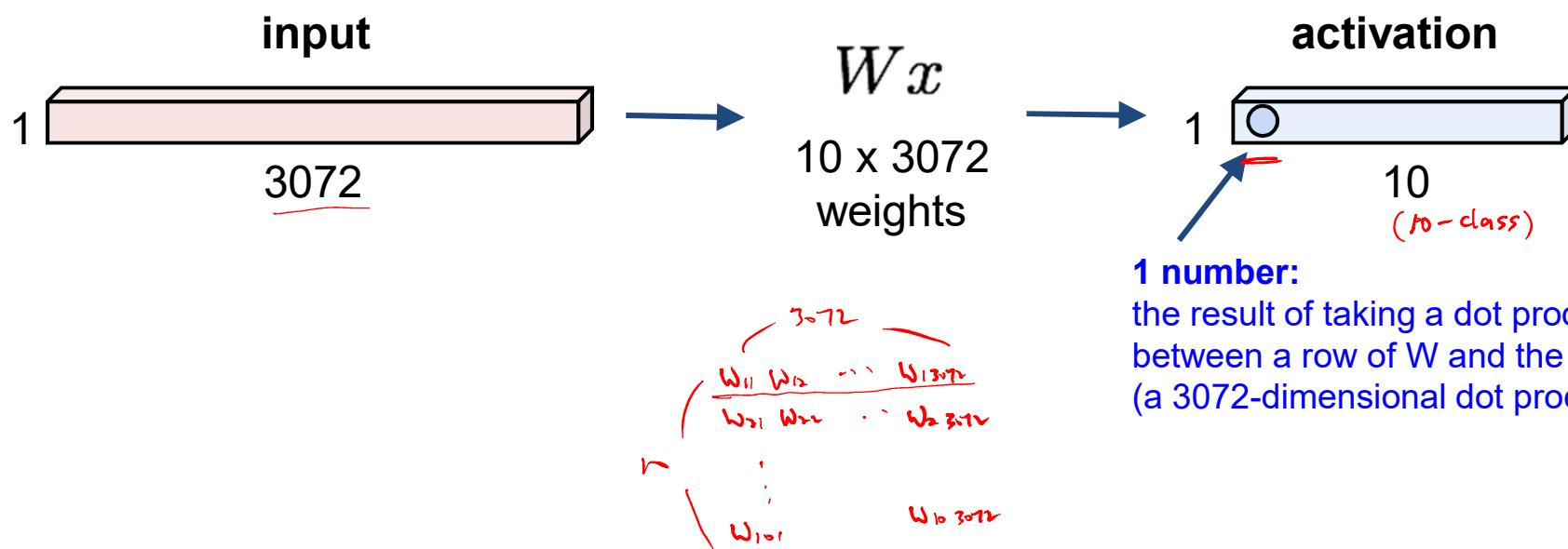
- Contains neurons that connect to the entire input volume, as in ordinary Neural Networks



# Fully Connected Layer (FC layer)

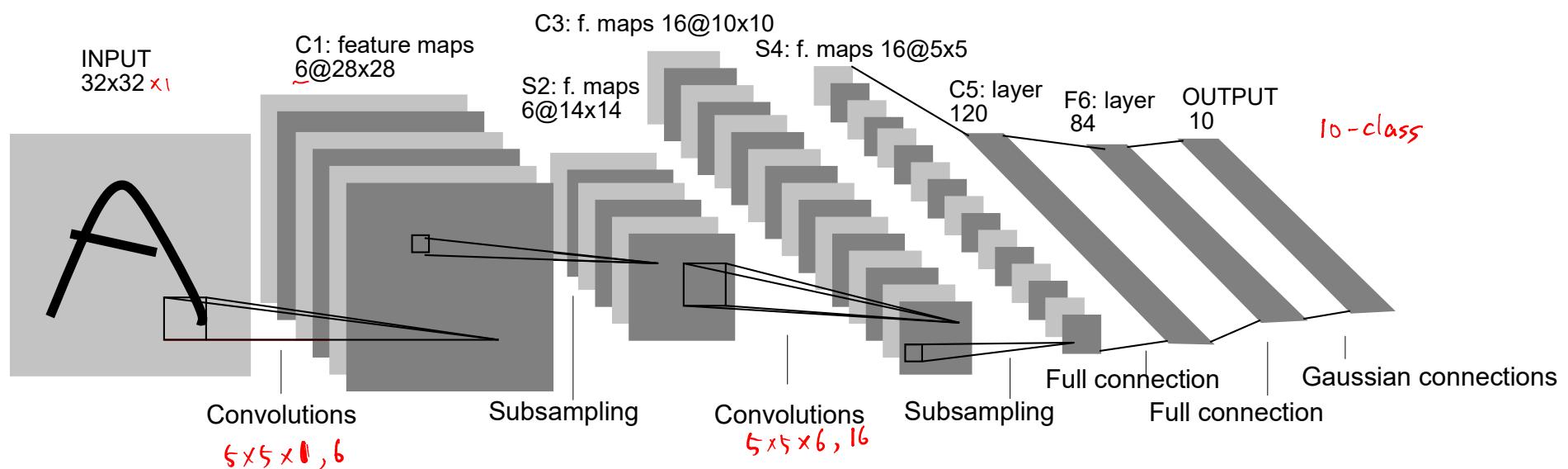
32x32x3 image -> stretch to 3072 x 1

Each neuron  
looks at the full  
input volume



# Convolutional Nets

- Example:
  - <http://yann.lecun.com/exdb/lenet/index.html>



# Backpropagation Algorithm

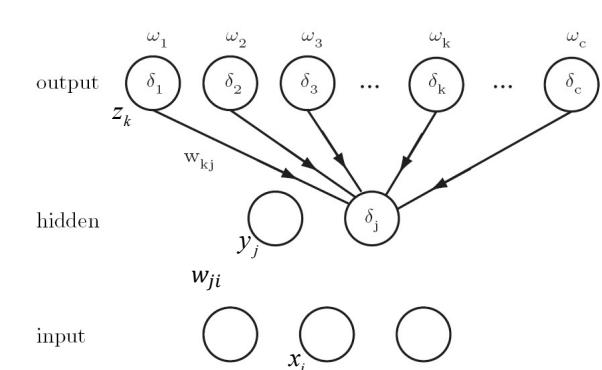
$$w(m+1) = w(m) + \Delta w(m) \quad \Delta w = -\eta \frac{\partial J}{\partial w} \quad J(w) = \frac{1}{2} \sum_{k=1}^c (t_k - z_k)^2 = \frac{1}{2} \|t - z\|^2$$

$$net_k = w_k^T y \quad net_j = w_j^T x \quad z_k = f(net_k) \quad y_j = f(net_j) \quad \frac{\partial net_k}{\partial w_{kj}} = y_j$$

$$\frac{\partial J}{\partial w_{kj}} = \frac{\partial J}{\partial z_k} \frac{\partial z_k}{\partial net_k} \frac{\partial net_k}{\partial w_{kj}} = -(t_k - z_k) f'(net_k) y_j$$

$$\delta_k = -\frac{\partial J}{\partial net_k} = -\frac{\partial J}{\partial z_k} \frac{\partial z_k}{\partial net_k} = (t_k - z_k) f'(net_k)$$

$$\Delta w_{kj} = -\eta \frac{\partial J}{\partial w_{kj}} = -\eta \frac{\partial J}{\partial net_k} \frac{\partial net_k}{\partial w_{kj}} = \eta (t_k - z_k) f'(net_k) y_j = \eta \delta_k y_j$$



$$\frac{\partial J}{\partial w_{ji}} = \frac{\partial J}{\partial y_j} \frac{\partial y_j}{\partial net_j} \frac{\partial net_j}{\partial w_{ji}}$$

$$\begin{aligned} \frac{\partial J}{\partial y_j} &= \frac{\partial}{\partial y_j} \left[ \frac{1}{2} \sum_{k=1}^c (t_k - z_k)^2 \right] = -\sum_{k=1}^c (t_k - z_k) \frac{\partial z_k}{\partial y_j} \\ &= -\sum_{k=1}^c (t_k - z_k) \frac{\partial z_k}{\partial net_k} \frac{\partial net_k}{\partial y_j} = -\sum_{k=1}^c (t_k - z_k) f'(net_k) w_{kj} \end{aligned}$$

$$\begin{aligned} \frac{\partial J}{\partial w_{ji}} &= \frac{\partial J}{\partial y_j} \frac{\partial y_j}{\partial net_j} \frac{\partial net_j}{\partial w_{ji}} \\ &= -\sum_{k=1}^c (t_k - z_k) f'(net_k) w_{kj} f'(net_j) x_i \end{aligned}$$

$$\begin{aligned} \delta_j &= -\frac{\partial J}{\partial net_j} = -\frac{\partial J}{\partial y_j} \frac{\partial y_j}{\partial net_j} \\ &= \sum_{k=1}^c (t_k - z_k) f'(net_k) w_{kj} f'(net_j) = f'(net_j) \sum_{k=1}^c w_{kj} \delta_k \end{aligned}$$

$$\Delta w_{ji} = \eta \delta_j x_i = \eta \left[ \sum_{k=1}^c w_{kj} \delta_k \right] f'(net_j) x_i$$

# Backpropagation in CNN

- Can be considered as a special case of MLP
  - Local connection – non connecting edges are disregarded
  - Weight sharing – kernels (filters) are updated
- Convolution based forward passing makes the backward passing also like convolution operations
- Pooling
  - E.g., max pooling – only the max element is updated