



# Statistical Methods in AI (CS7.403)

Introduction to CNN/Deep Learning

Ravi Kiran ([ravi.kiran@iiit.ac.in](mailto:ravi.kiran@iiit.ac.in))

<https://ravika.github.io>



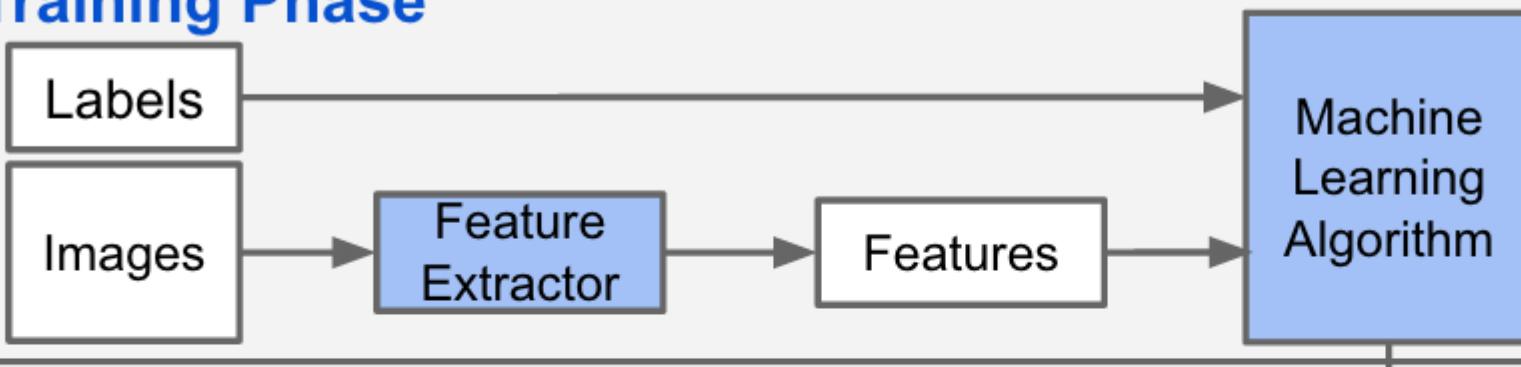
Center for Visual Information Technology (CVIT)  
IIIT Hyderabad





# Classification

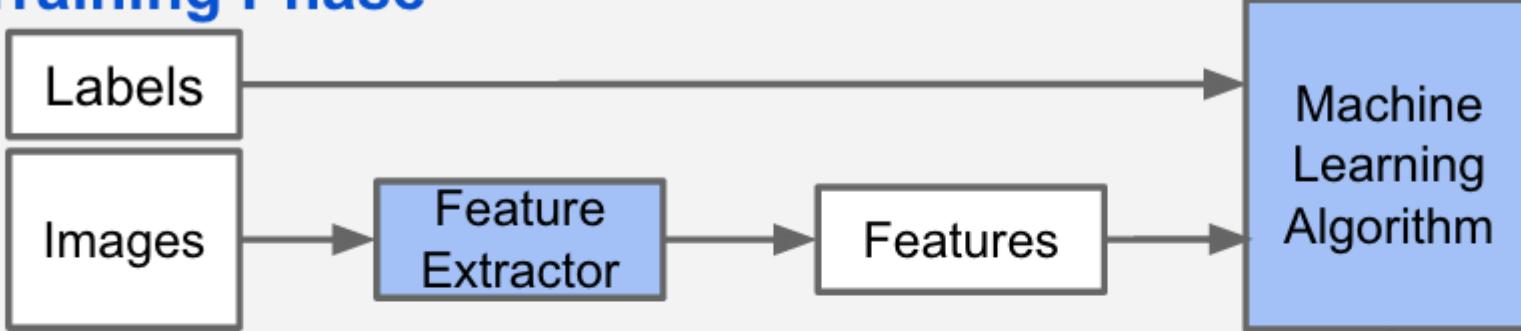
## Training Phase





# Classification

## Training Phase



## Prediction Phase





# Feature Vectors

- Usually a single object can be represented using several features, e.g.

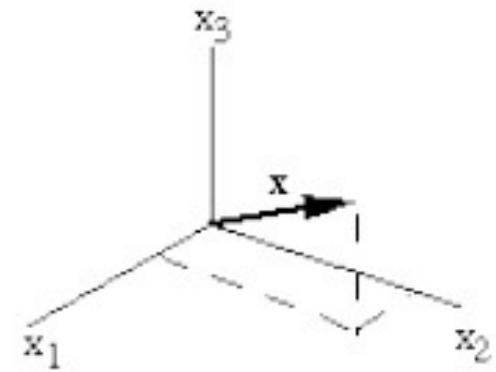
- $x_1$  = shape (e.g. nr of sides)
- $x_2$  = size (e.g. some numeric value)
- $x_3$  = color (e.g. rgb values)
- ...
- $x_d$  = some other (numeric) feature.

- **X** becomes a feature vector
  - $x$  is a point in a d-dimensional **feature space**.



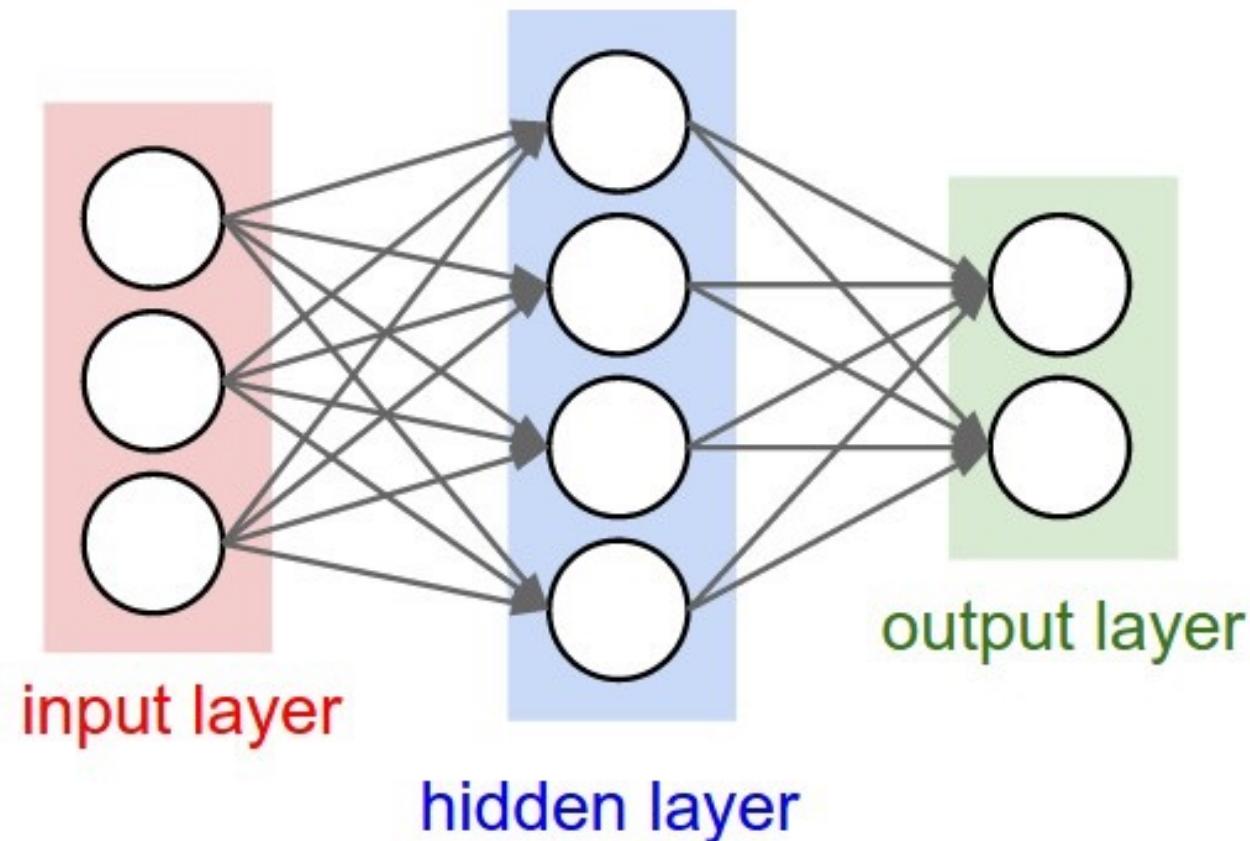
$$\begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \mathbf{x}$$

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_d \end{bmatrix}$$





# Traditional approach: Extract features, feed to NN





---

**Wouldn't it be nice if we could feed the image directly ?**

**... and let the “learning process” figure out which features to extract ?**



---

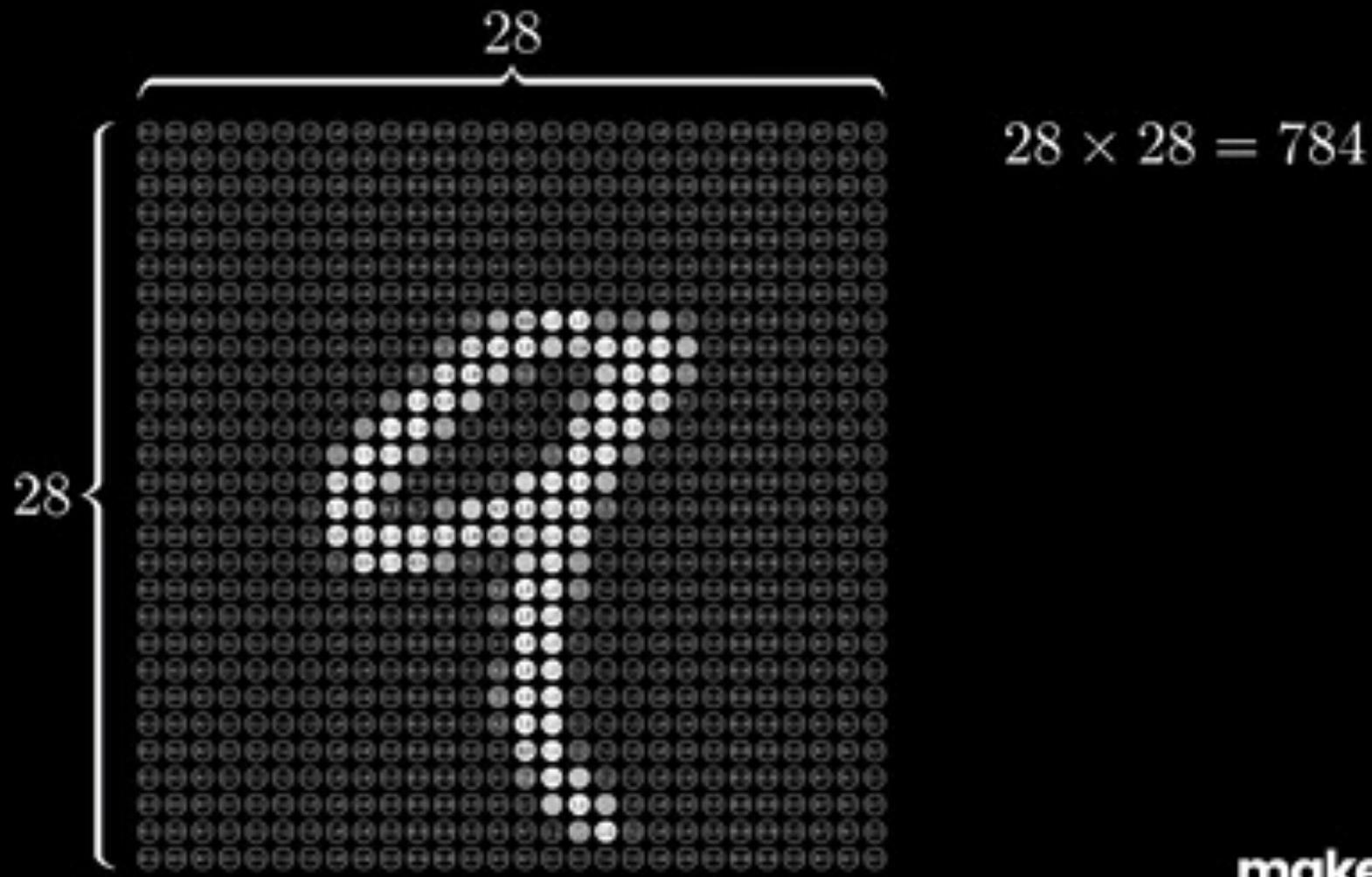
# MNIST Handwritten Digits Dataset



0000000000000000  
1111111111111111  
2222222222222222  
3333333333333333  
4444444444444444  
5555555555555555  
6666666666666666  
7777777777777777  
8888888888888888  
9999999999999999



# Flatten image , Feed to NN

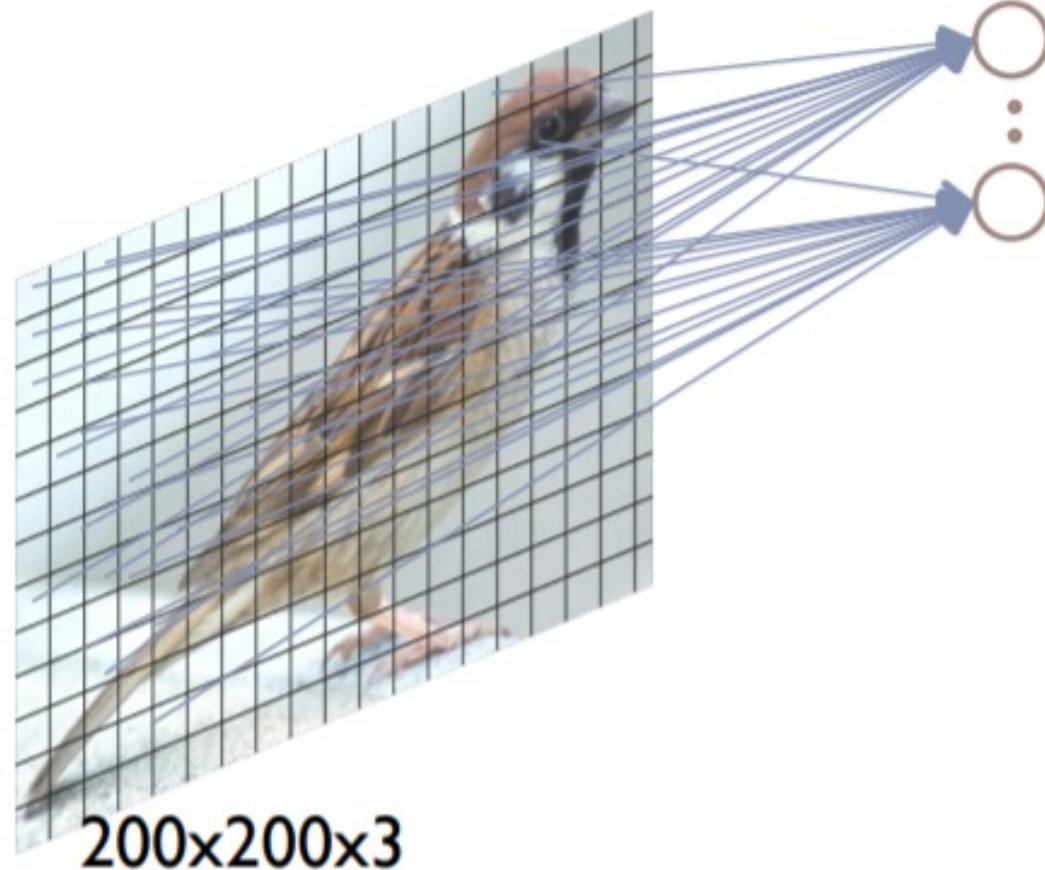


makeagif.com



---

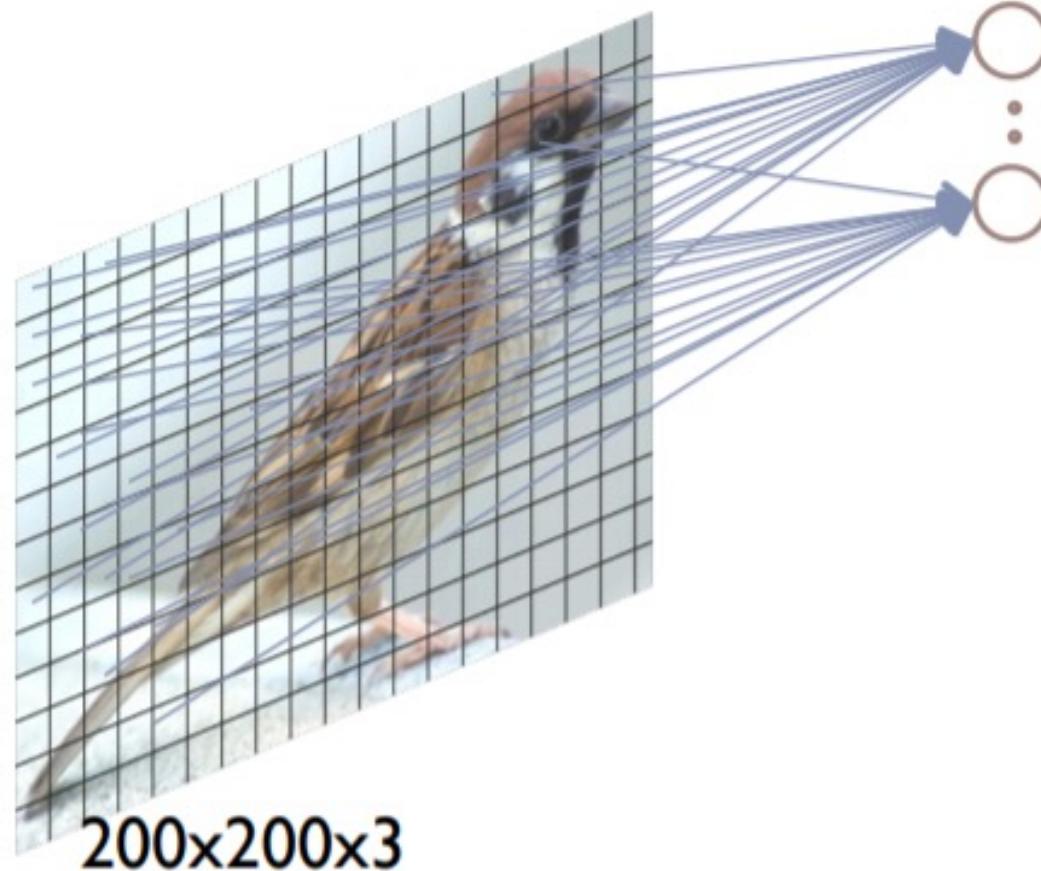
# For color images ?



- #Hidden Units: 120,000
- #Params: 14.4 billion
- Need huge training data to prevent over-fitting!



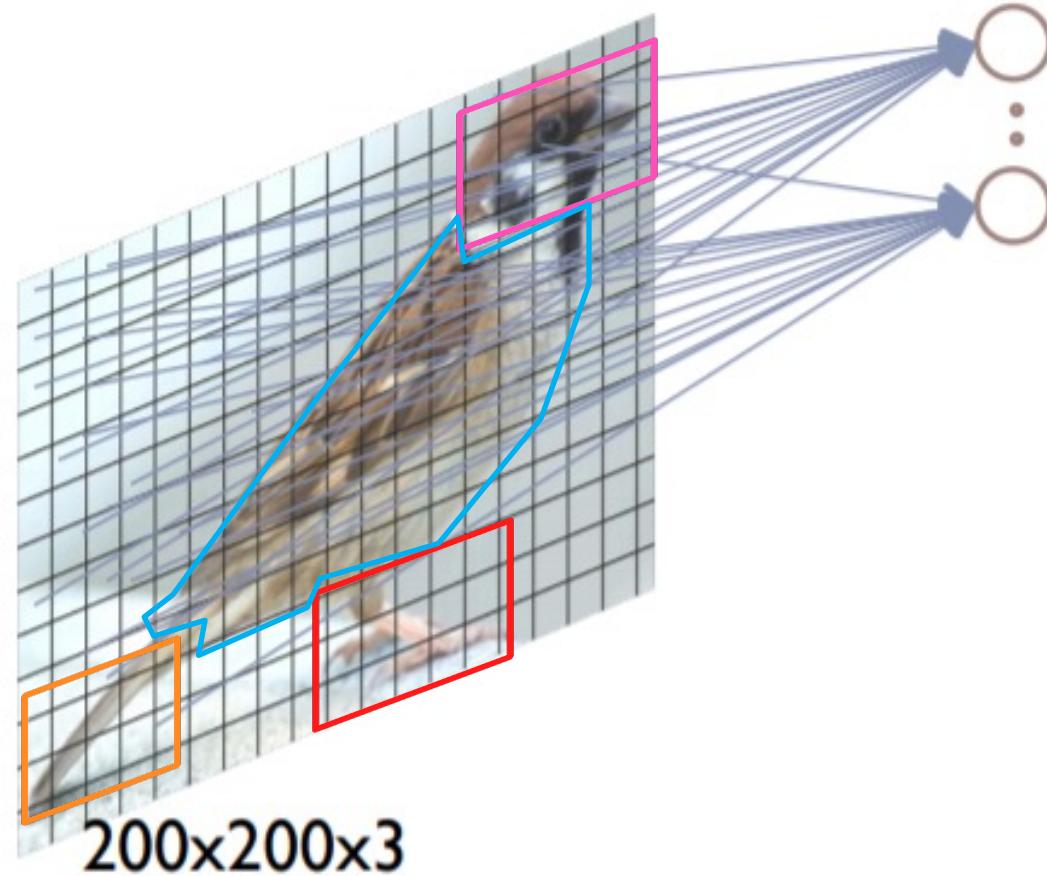
# Do we really need full connectivity ?



- #Hidden Units: 120,000
- #Params: 14.4 billion
- Need huge training data to prevent over-fitting!



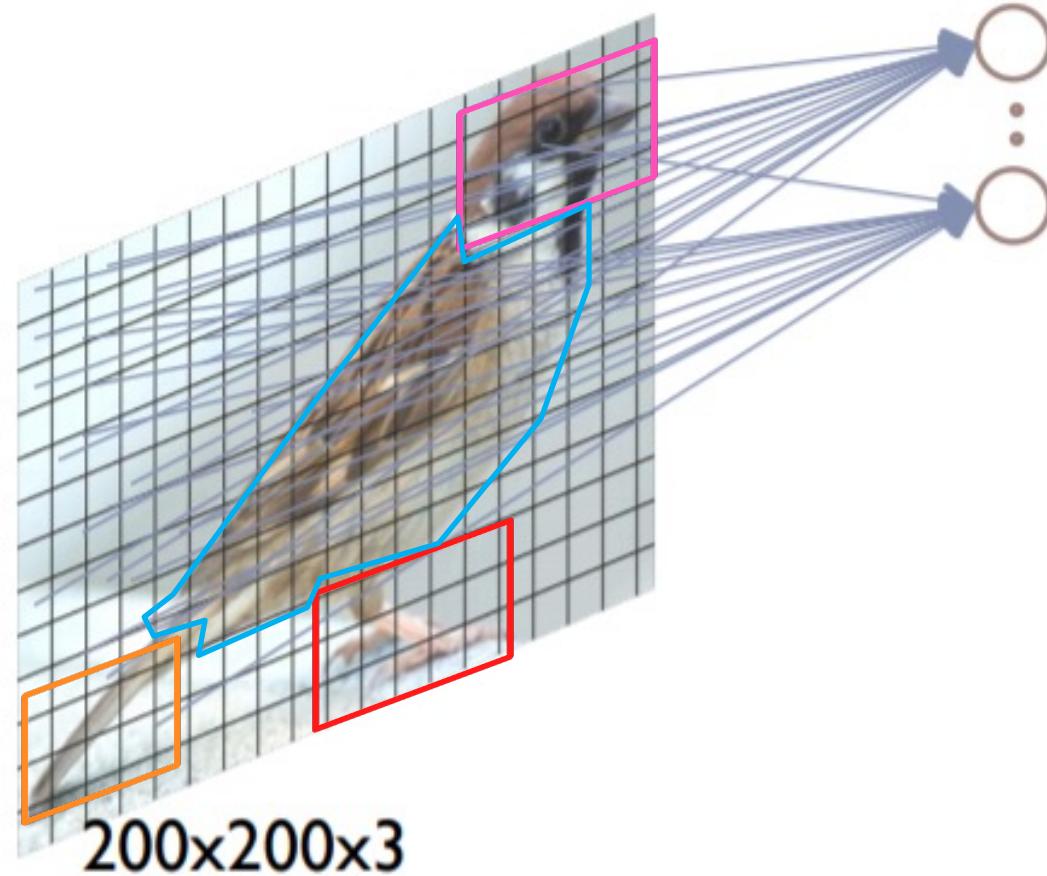
# Do we really need full connectivity ?



- Head + Body + Tail + Legs = sparrow



# Do we really need full connectivity ?

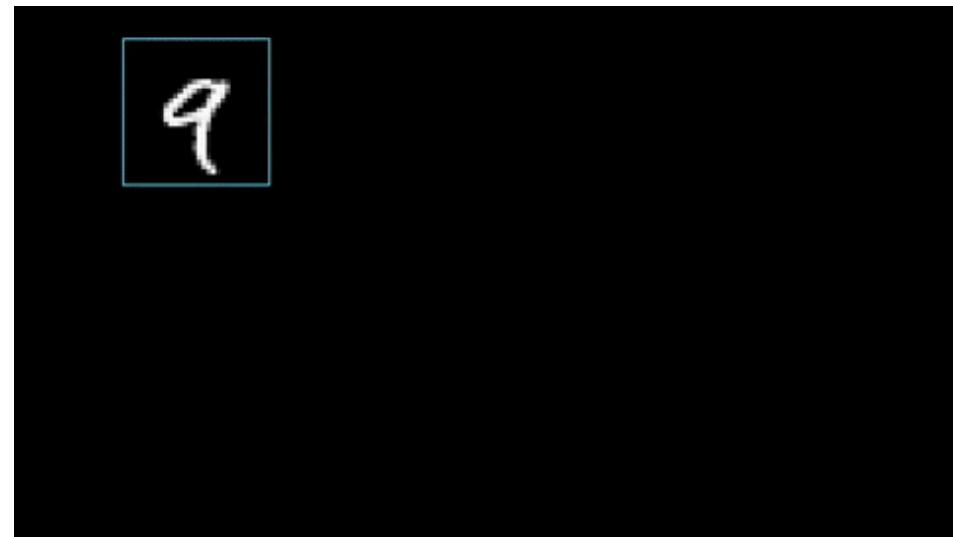


- Head + Body + Tail + Legs = sparrow



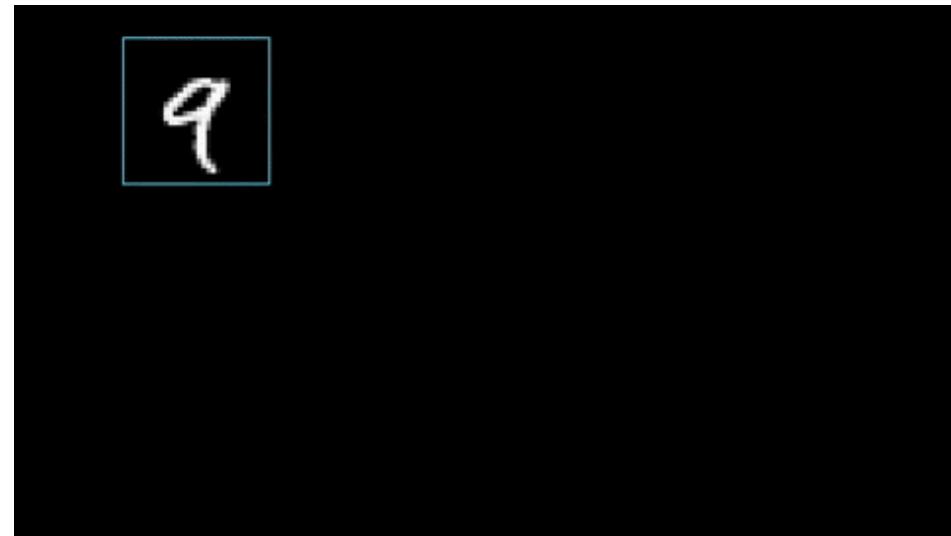


- 
- Images are 2D.
  - Assumption: Object image = combination of **2D image patterns**



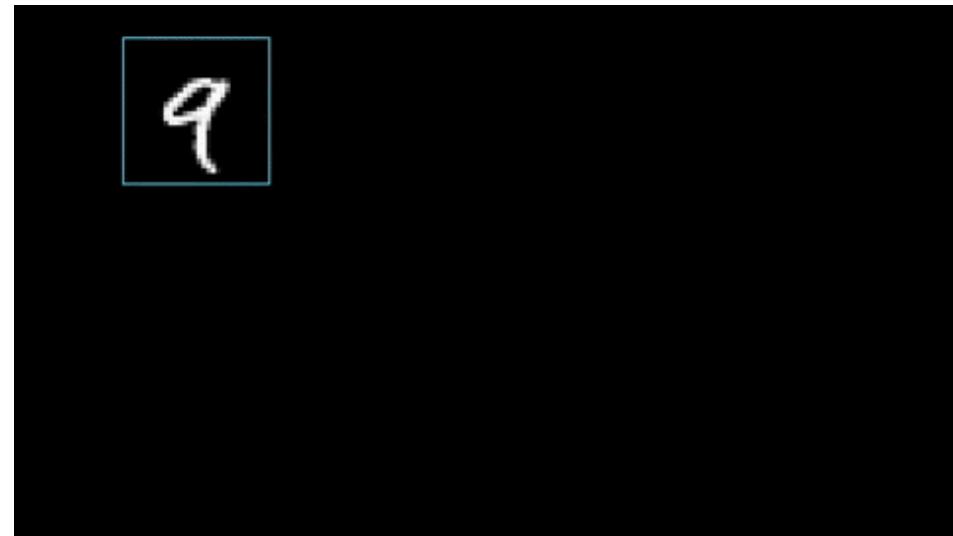


- Images are 2D.
- Assumption: Object image = combination of **2D image patterns**
- Machine Learning Strategy:
  - a) [Pre-final layer] Determine 2D image patterns
  - b) Map 2D image patterns → Target label



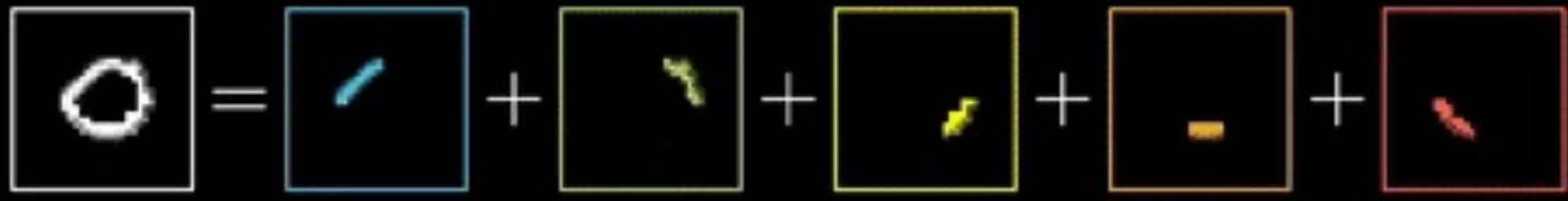


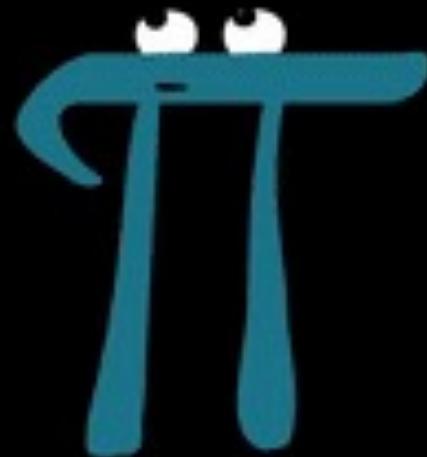
- Images are 2D.
- Assumption: Object image = combination of **2D image patterns**
- Machine Learning Strategy:
  - [Pre-final layer] Determine 2D image patterns
  - Map 2D image patterns → Target label
- NOTE: 2D image patterns are smaller than image





# Hierarchical decomposition of input/features

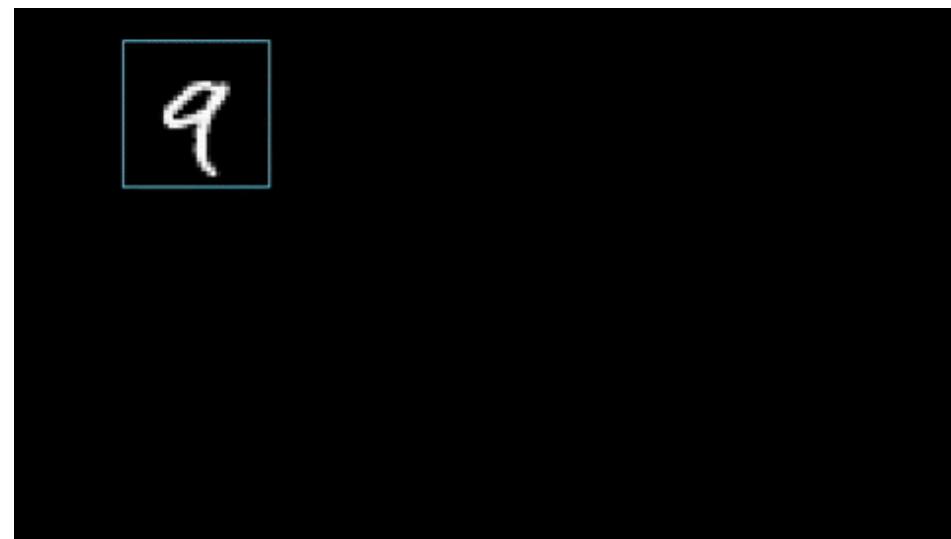
$$\text{input} = \text{stroke} + \text{background} + \text{foreground} + \text{noise} + \text{color}$$




makeagif.com



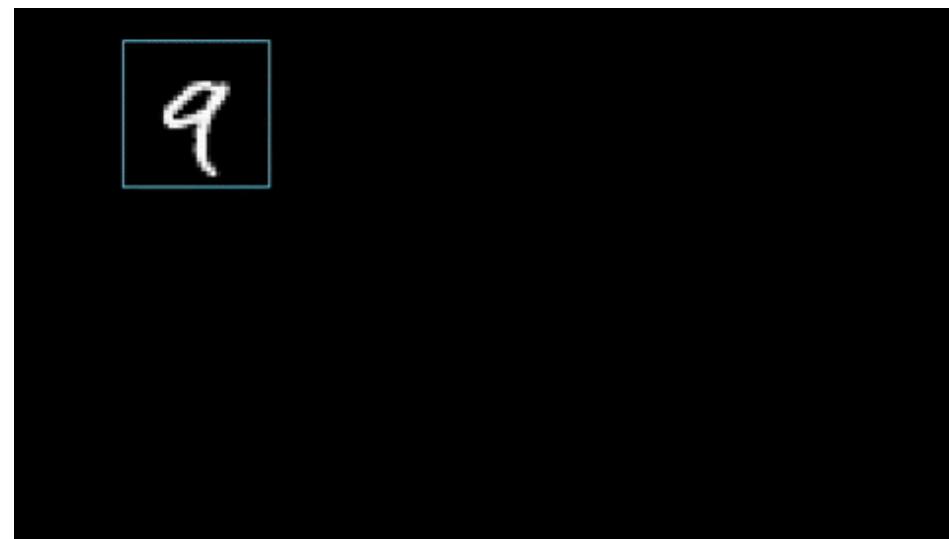
- Images are 2D.
- Assumption: Object image = **hierarchical** combination of **2D image patterns**
- **Machine Learning Strategy:**
  - [All except pre-final] Determine 2D image patterns
  - [Pre-final layer] Map 2D image patterns → Target label
- NOTE: 2D image patterns are smaller than image





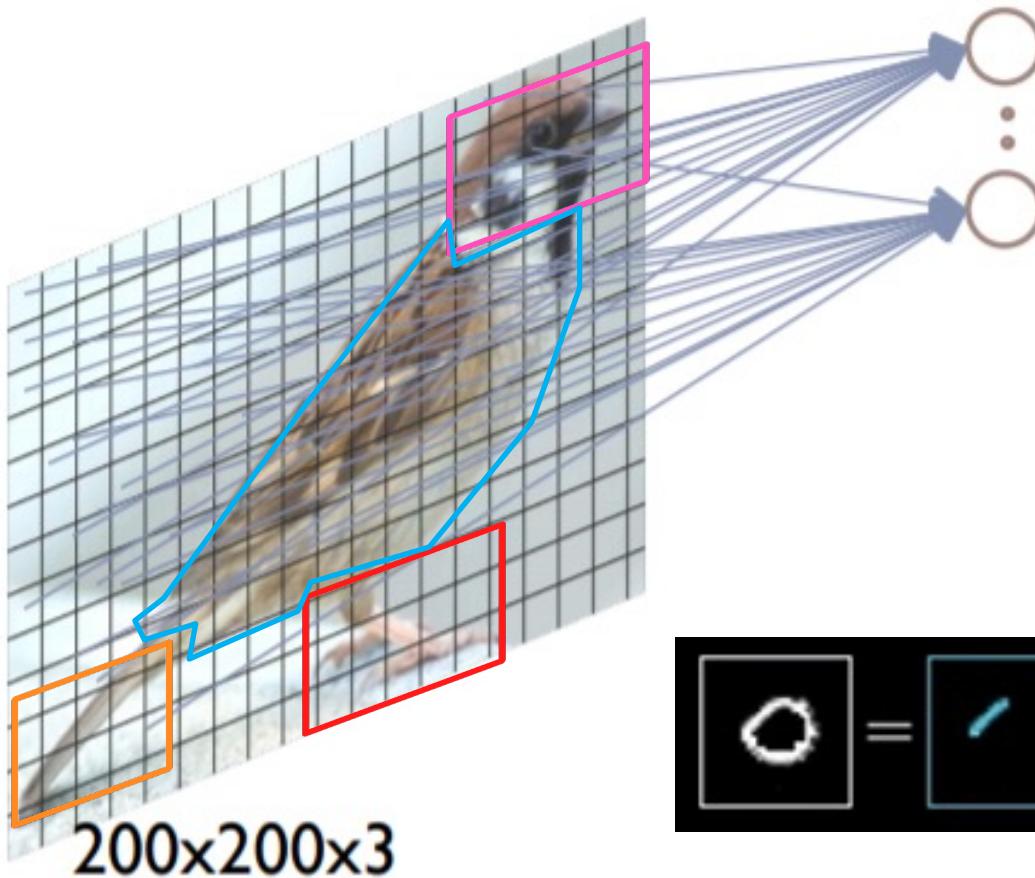
*Efficient !*

- Images are 2D.
- Assumption: Object image = **hierarchical** combination of **2D image patterns**
- **Machine Learning Strategy:**
  - [All except pre-final] Determine 2D image patterns
  - [Pre-final layer] Map 2D image patterns → Target label
- NOTE: 2D image patterns are smaller than image

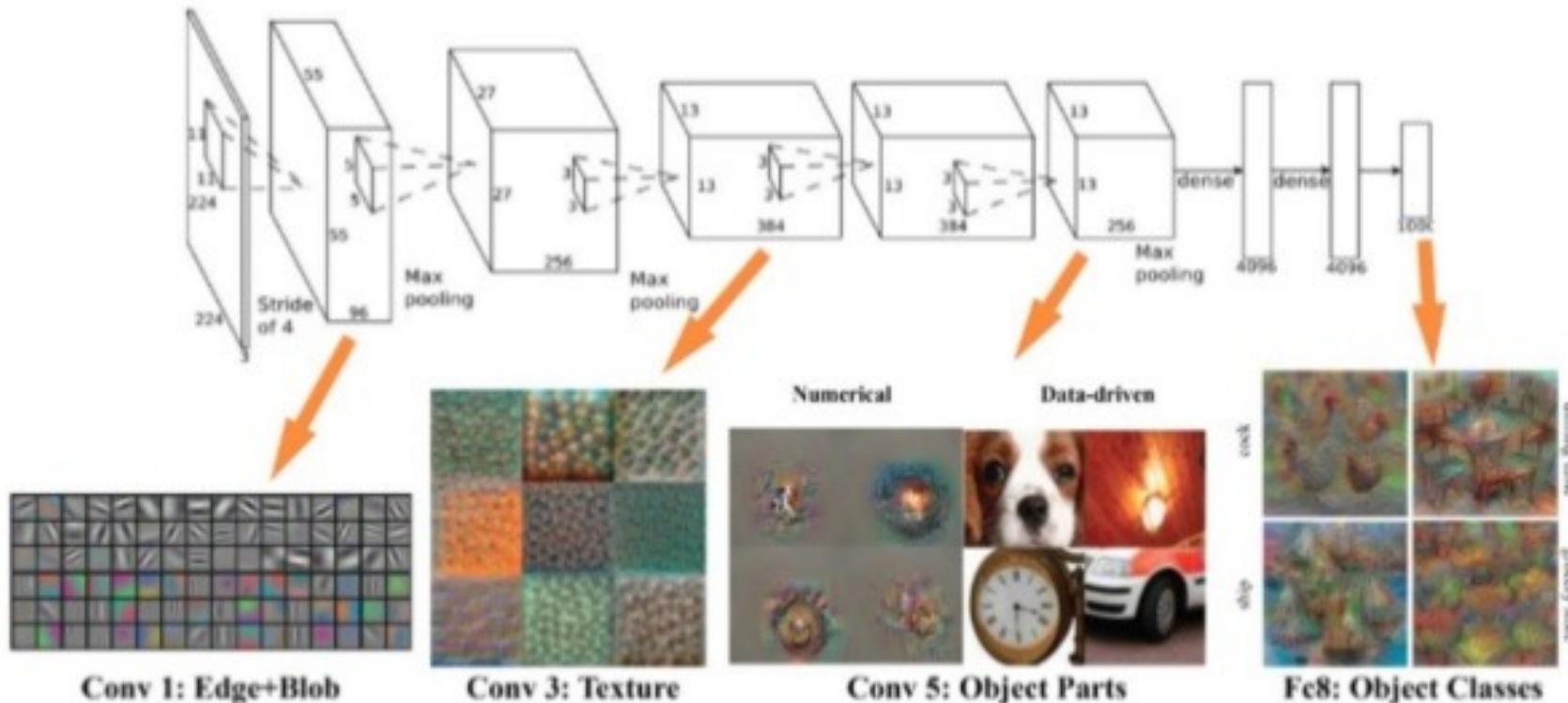




**SPOILER ALERT!**



Neural Networks don't really learn such conveniently interpretable Representations – this is just a toy illustration



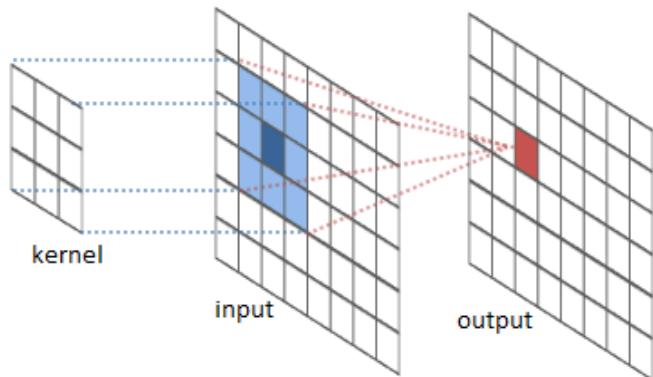
From: [mNeuron: A Matlab Plugin to Visualize Neurons from Deep Models](#), Donglai Wei et. al.

- Pattern detectors = Filters
- What should be the size of filters ?
- How many filters do we need at each level ?
- How many levels ?



# How do filters work ?

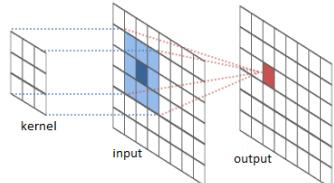
- Determine the filter
- ‘Scan’ the filter on input image / representation
- Output = A record of all input locations that ‘match’ the filter



$$\begin{pmatrix} 0 & 1 & 2 \\ 2 & 2 & 0 \\ 0 & 1 & 2 \end{pmatrix}$$



# How do filters work ?



$$\begin{matrix} 0 & 0 & 0 \\ 0 & -1 & 1 \\ 0 & 0 & 0 \end{matrix}$$



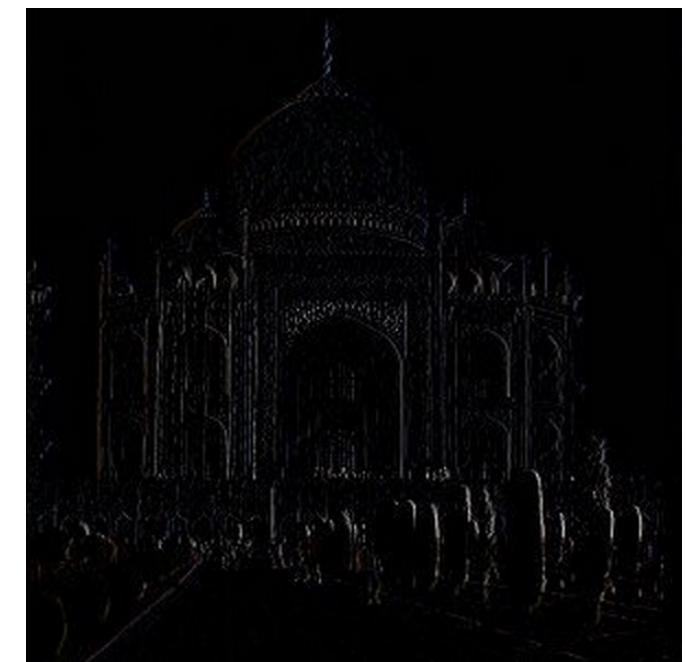
Input



Kernel Convolution

-1	-1	-1
-1	8	-1
-1	-1	-1

Feature Map



Features



---

**Wouldn't it be nice if we could feed the image directly ?**

**... and let the “learning process” figure out which features to extract ?**



**Wouldn't it be nice if we could feed the image directly ?**

**... and let the “learning process” figure out which features to extract ?**

**(which filters to construct)**



---

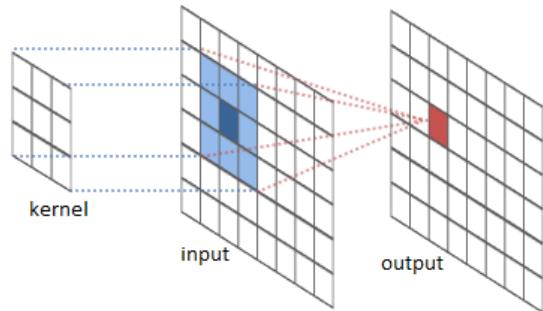
# Convolution (in general)

3 <sub>0</sub>	3 <sub>1</sub>	2 <sub>2</sub>	1	0
0 <sub>2</sub>	0 <sub>2</sub>	1 <sub>0</sub>	3	1
3 <sub>0</sub>	1 <sub>1</sub>	2 <sub>2</sub>	2	3
2	0	0	2	2
2	0	0	0	1

12	12	17
10	17	19
9	6	14

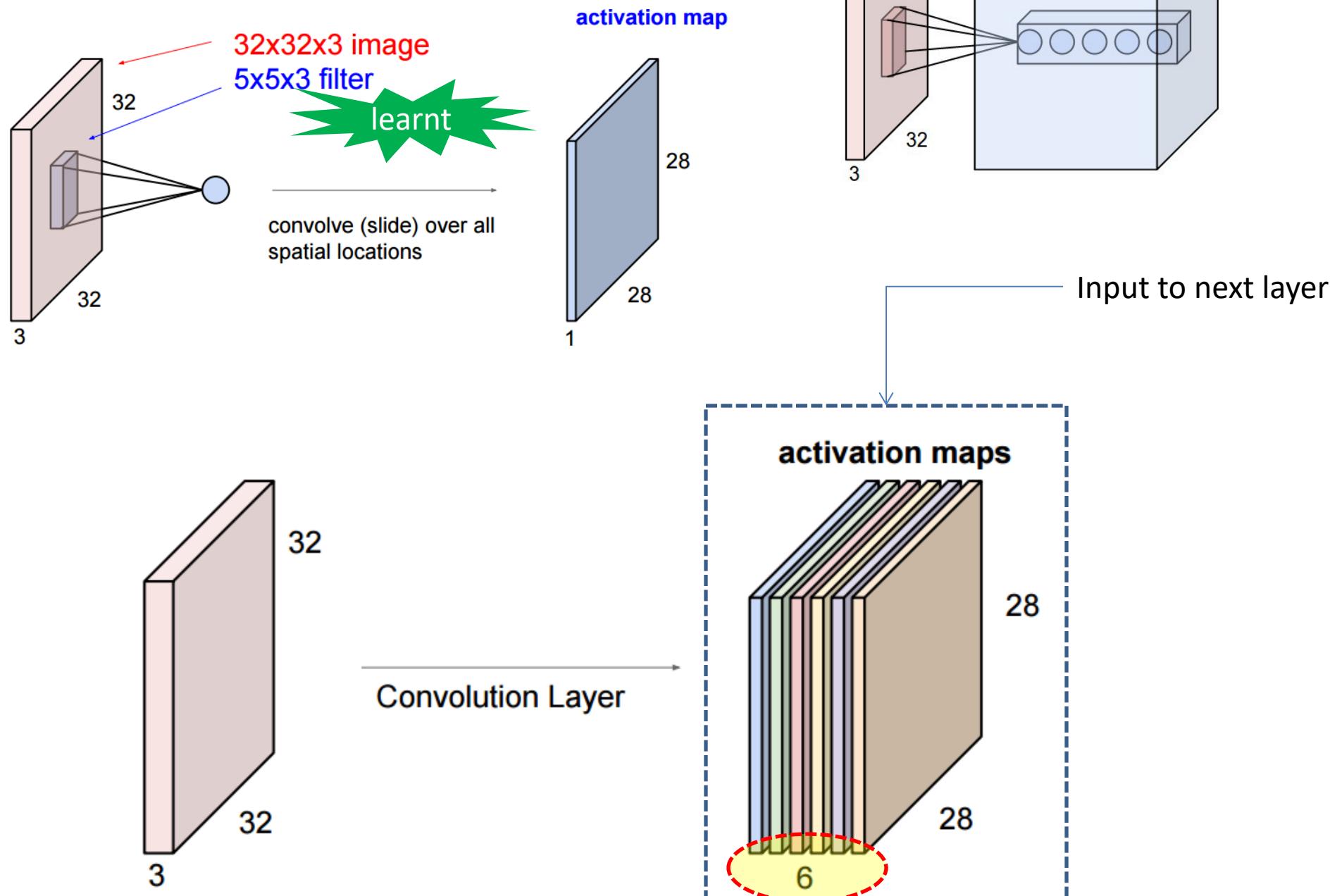


# Multiple filters, multiple ‘feature maps’



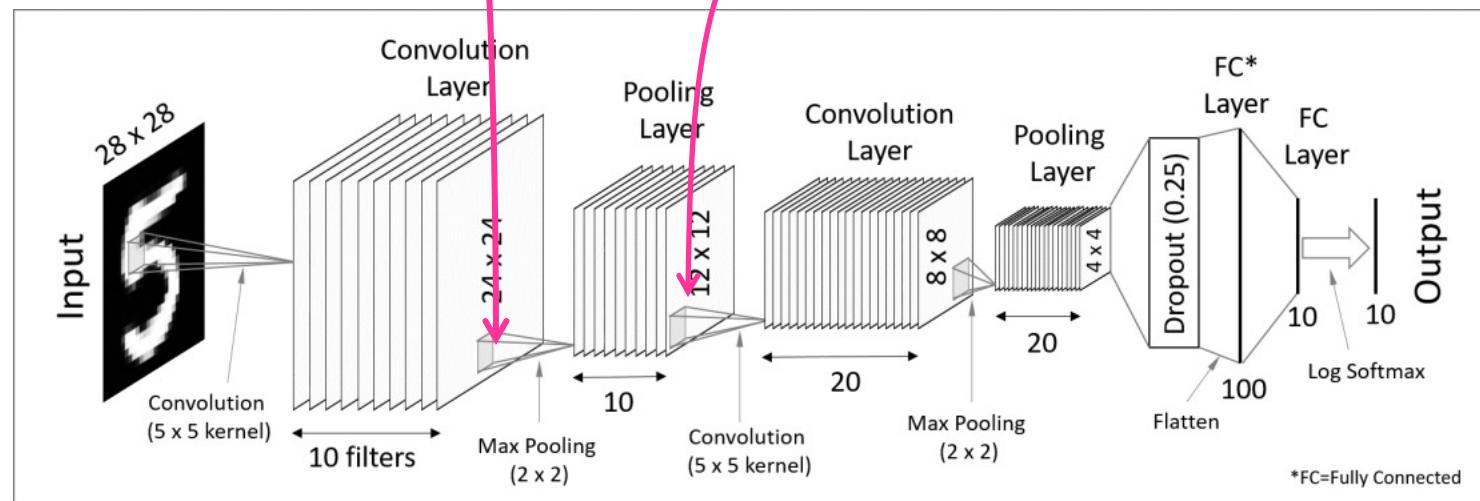
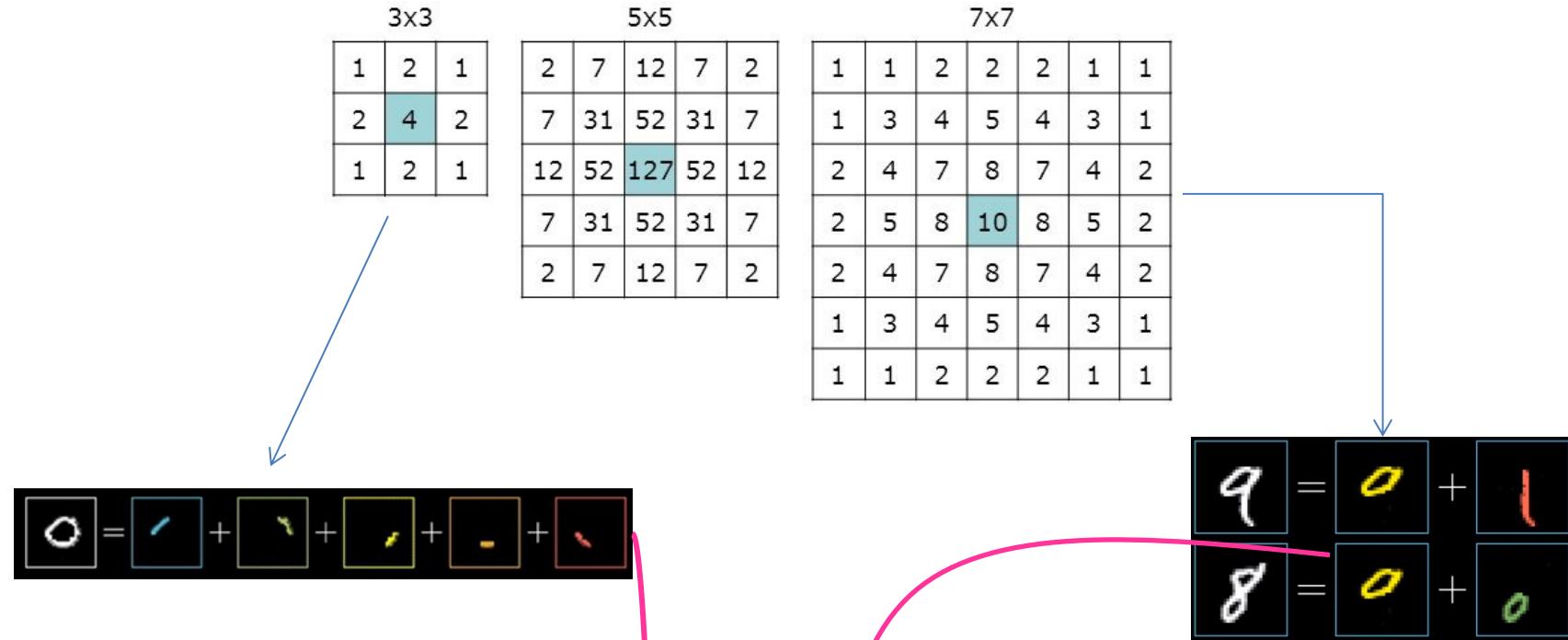


# Convolution Layer



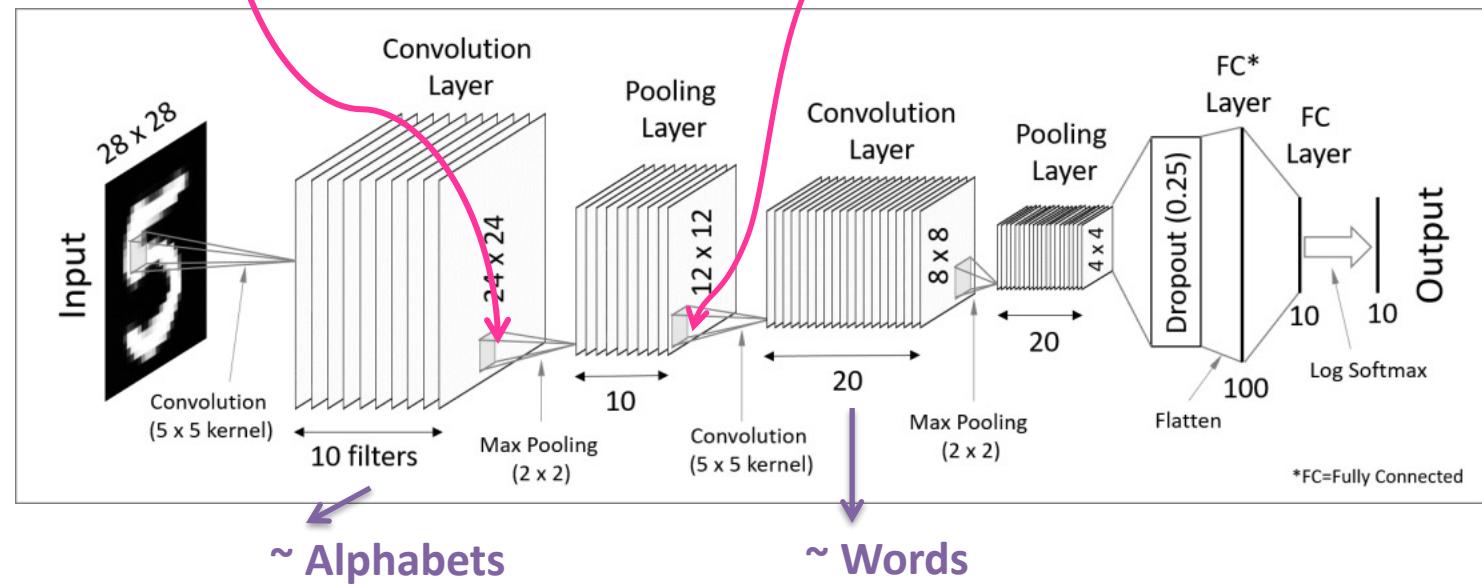
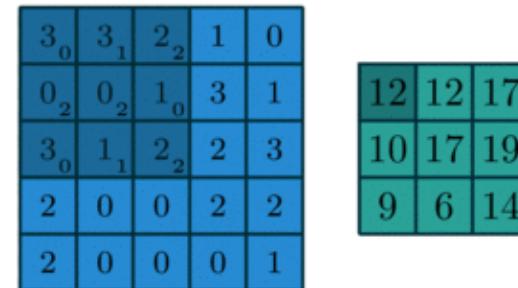


# Design choices : Filter size, # of filters





# Design choices : Filter size, # of filters





# Pooling Layer

1	1	2	4
5	6	7	8
3	2	1	0
1	2	3	4

max pool with 2x2 filters  
and stride 2



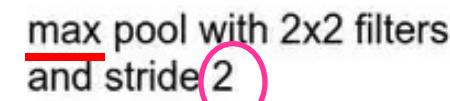
6	8
3	4



# Pooling Layer

1	1	2	4
5	6	7	8
3	2	1	0
1	2	3	4

max pool with 2x2 filters  
and stride 2



6	8
3	4



# Pooling Layer

1	1	2	4
5	6	7	8
3	2	1	0
1	2	3	4

max pool with 2x2 filters  
and stride 2



6	8
3	4

- **Motivation: We care about presence of features, not their exact location !**
- Dimensionality Reduction
- Prevents overfitting



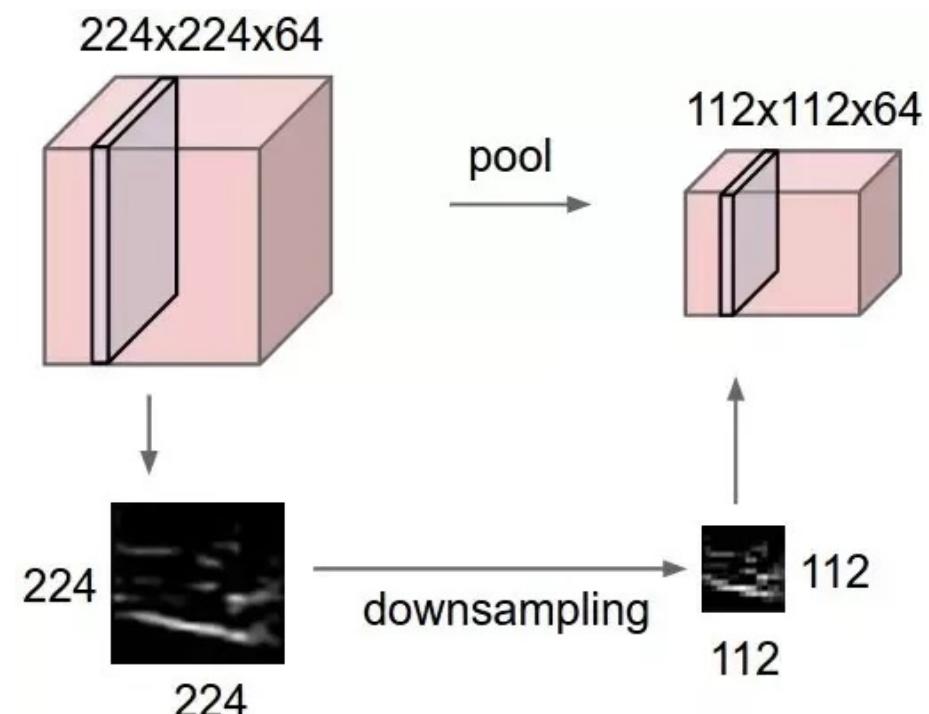
# Pooling Layer

1	1	2	4
5	6	7	8
3	2	1	0
1	2	3	4

max pool with 2x2 filters  
and stride 2

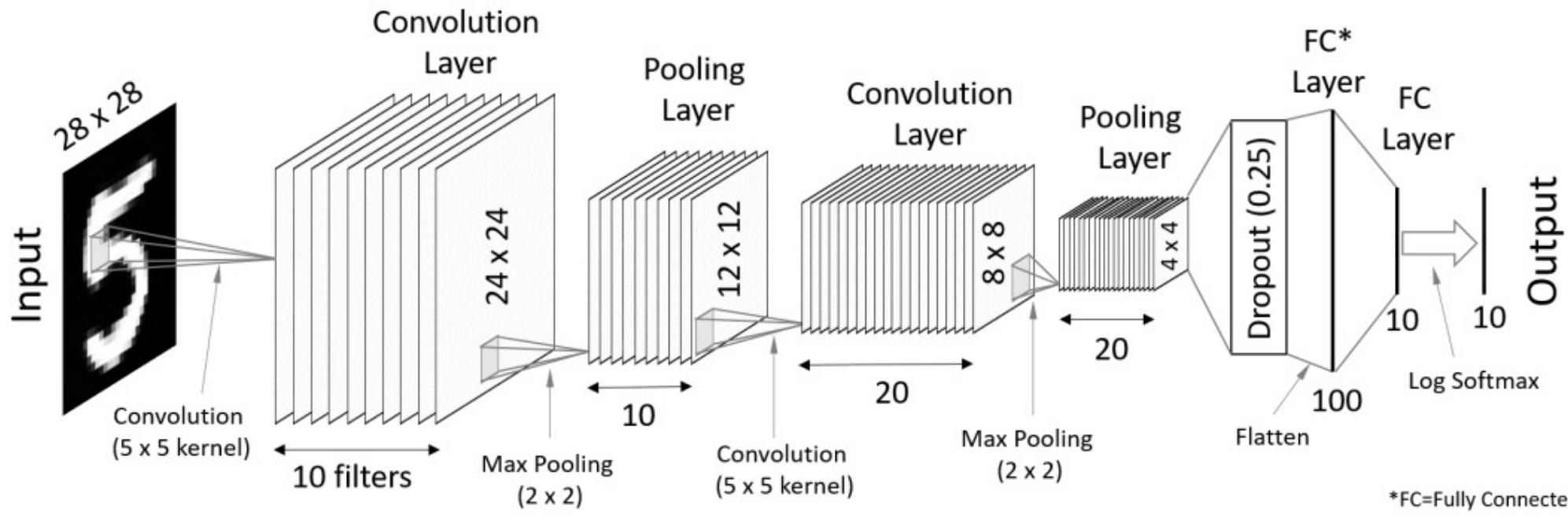
6	8
3	4

- **Motivation: We care about presence of features, not their exact location !**
- Dimensionality Reduction
- Prevents overfitting



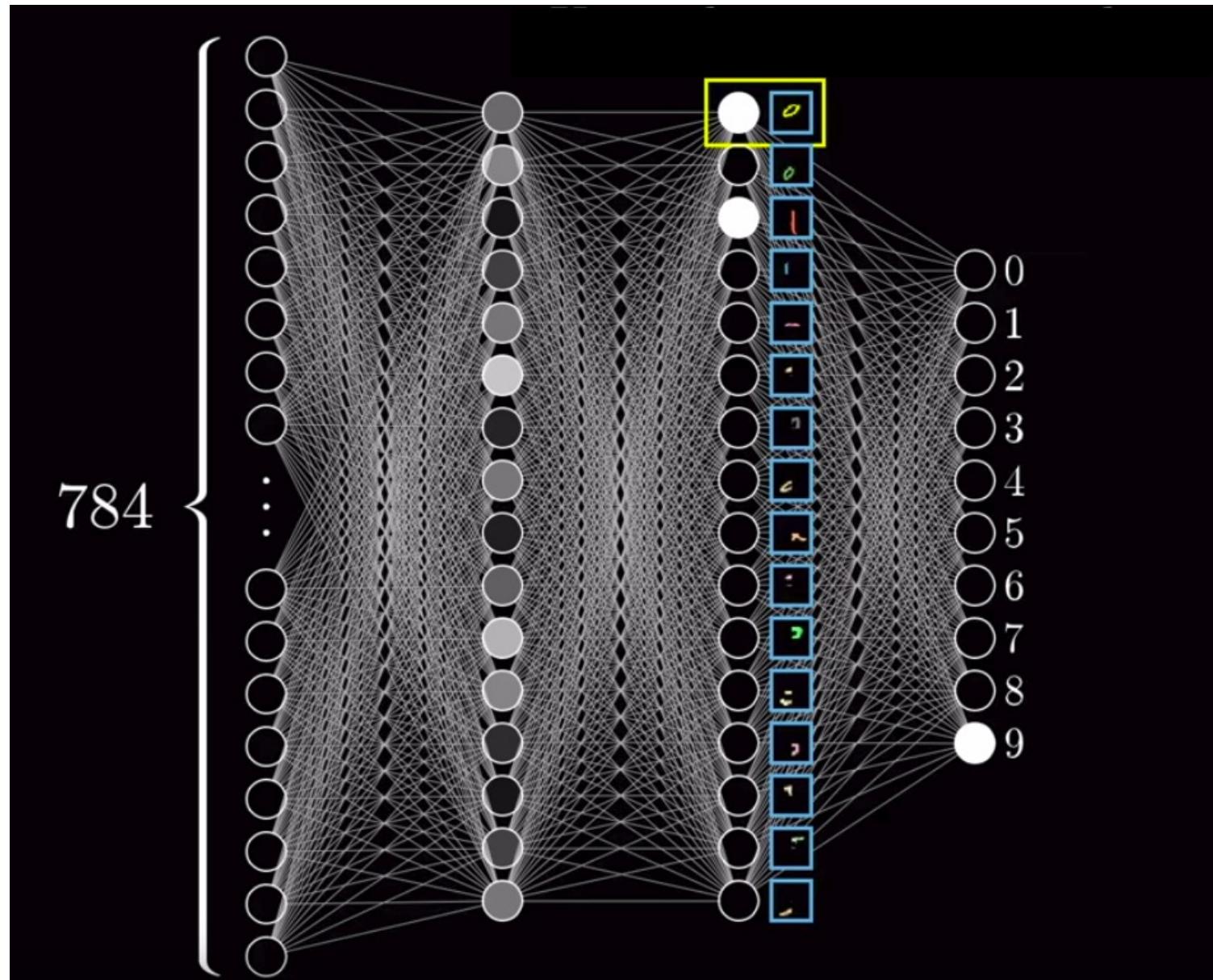


# Fully-connected layers





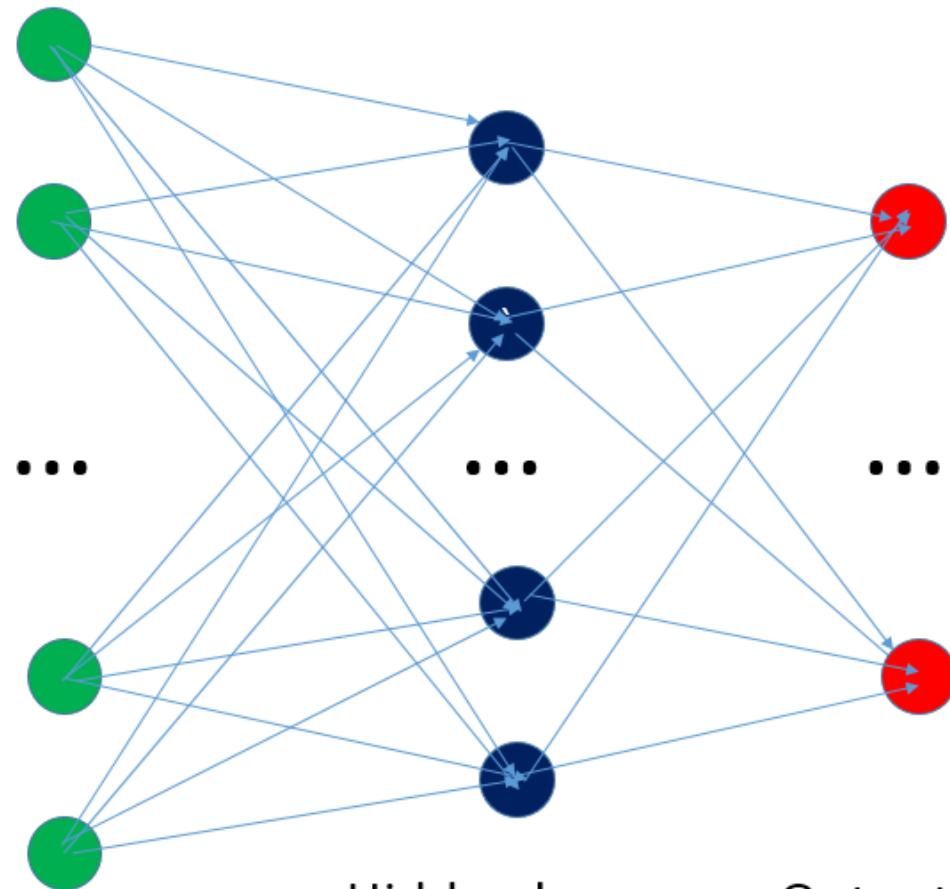
# Fully Connected Layers





7

Input image:  
28x28 pixels



Input layer:  
784 neurons,  
one per pixel

Hidden layer:  
100 neurons

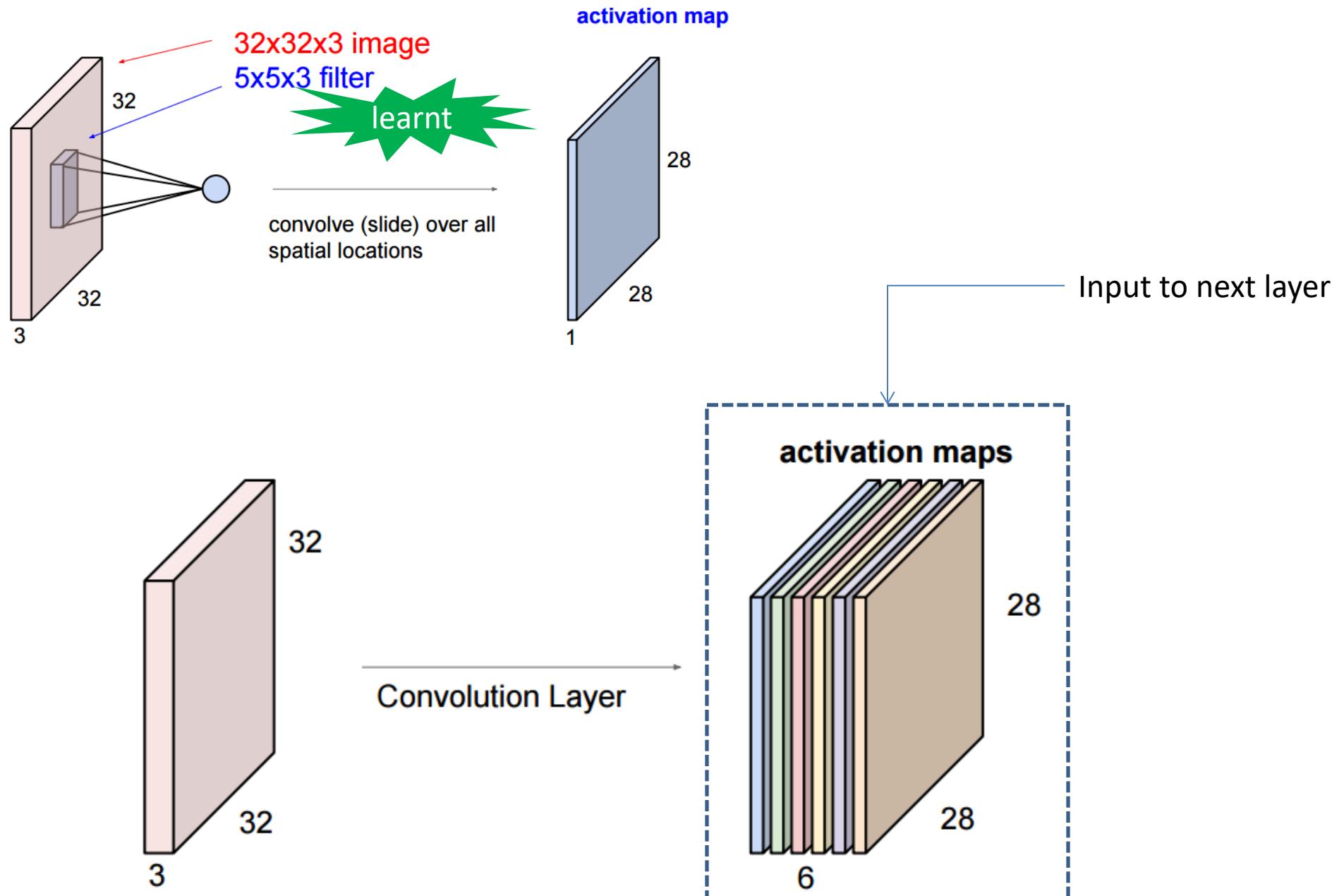
Output layer:  
10 neurons



Output:  
predicted  
digit value

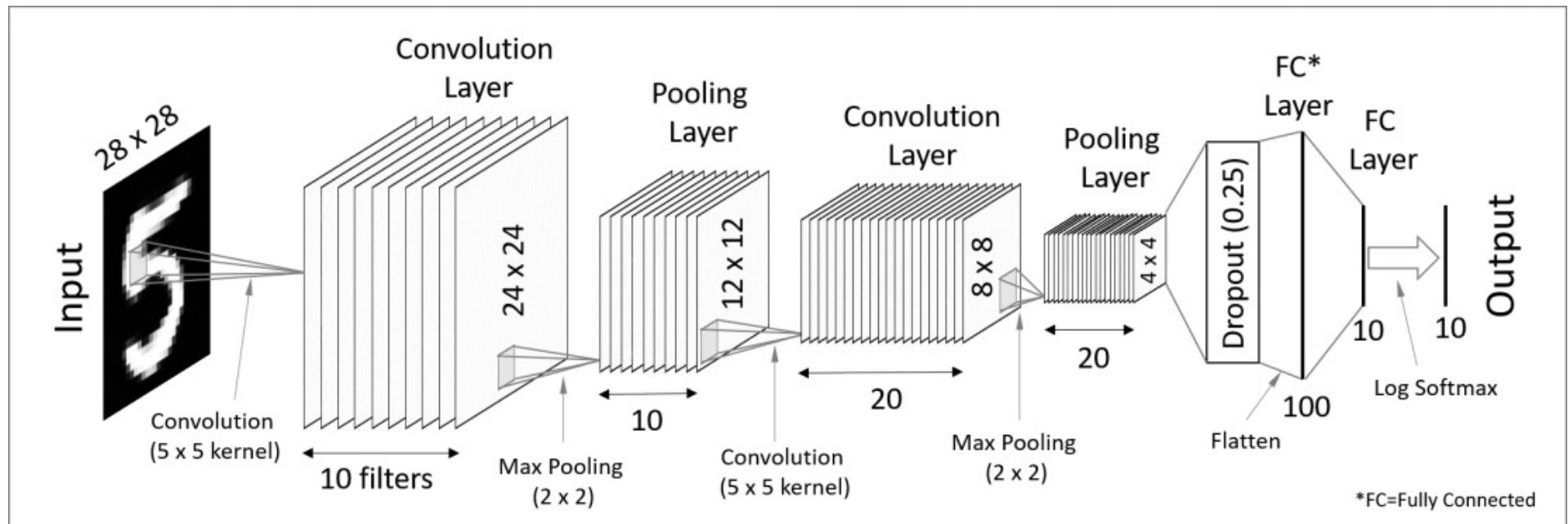


# Convolution Layer





# Convolutional Neural Network





# Pooling Layer

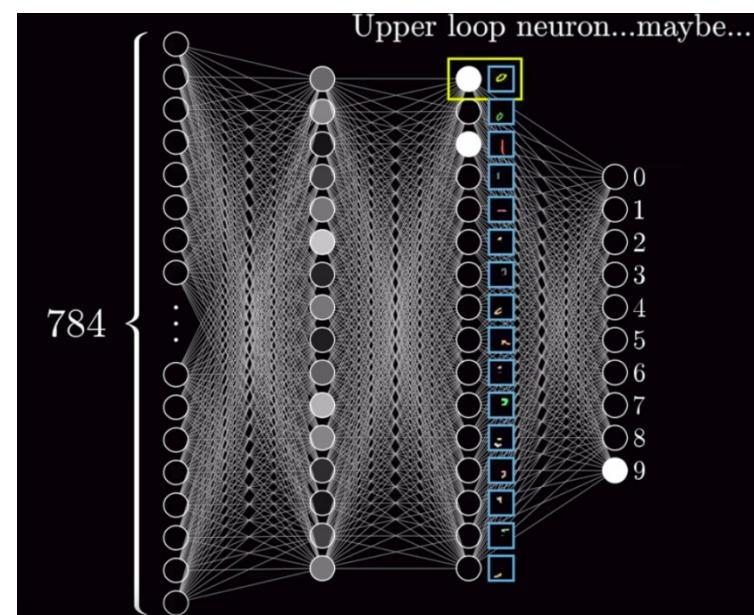
1	1	2	4
5	6	7	8
3	2	1	0
1	2	3	4

max pool with 2x2 filters  
and stride 2



6	8
3	4

# Fully-connected Layer



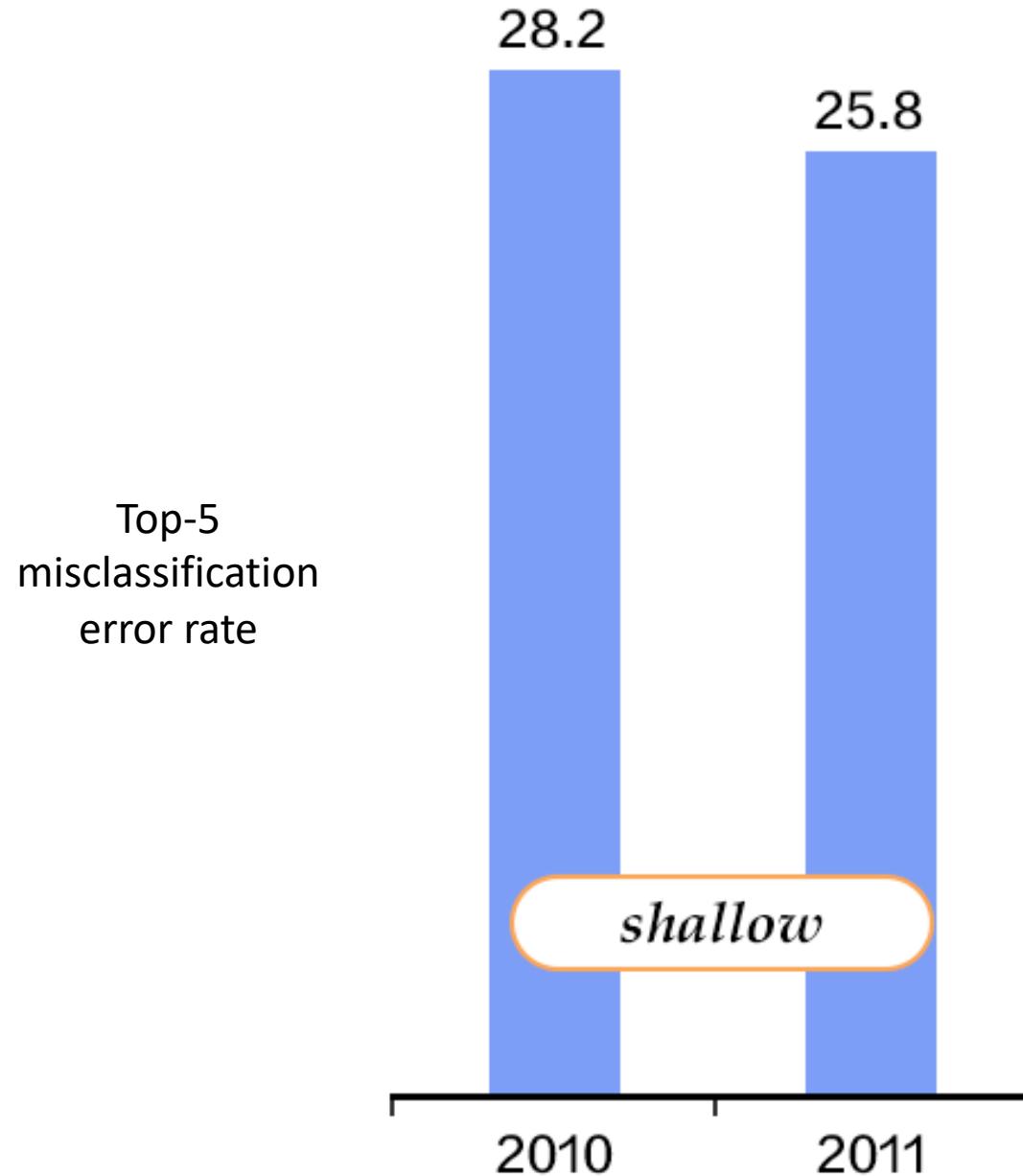


# ImageNet Challenge

## IMAGENET

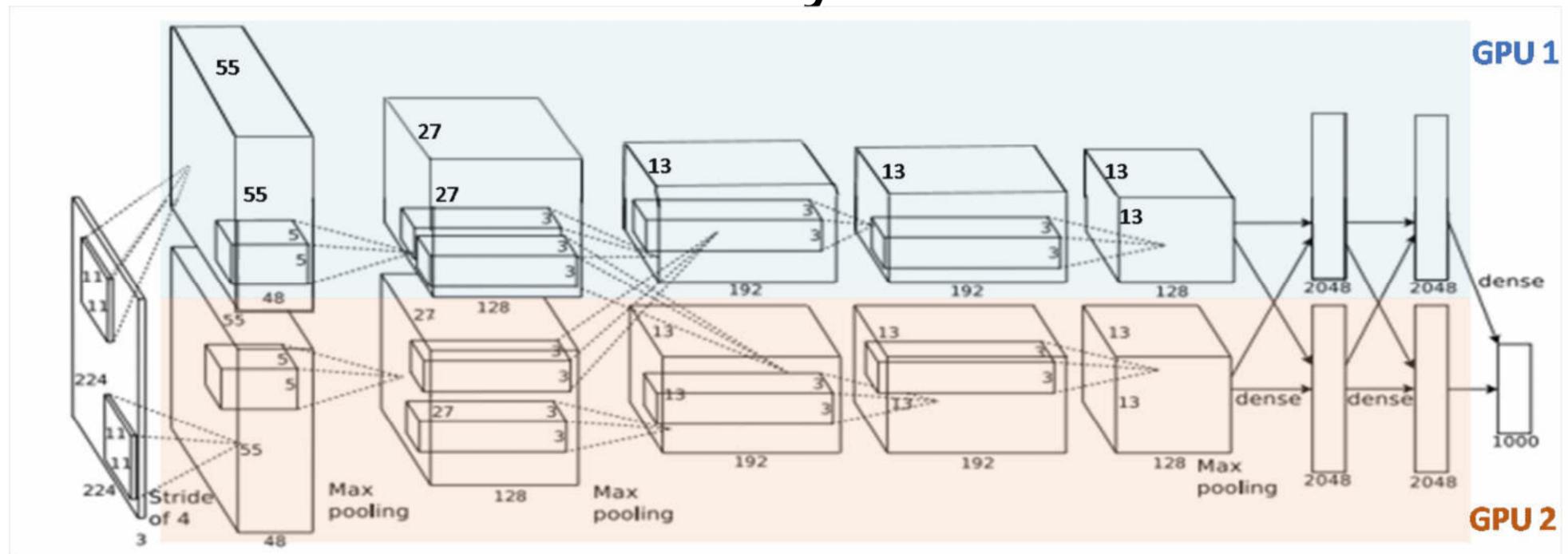
- 1,000 object classes (categories).
- Images:
  - 1.2 M train
  - 100k test.







# Case Study: AlexNet

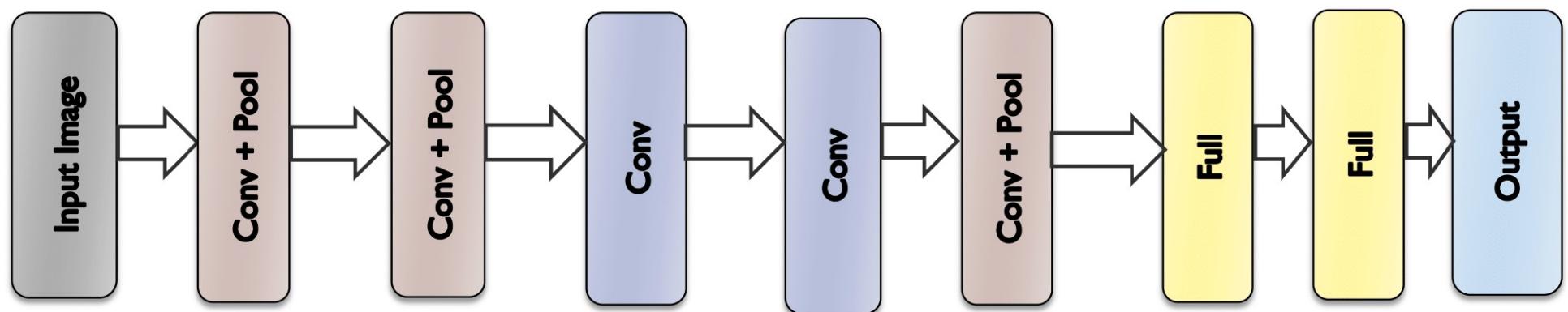
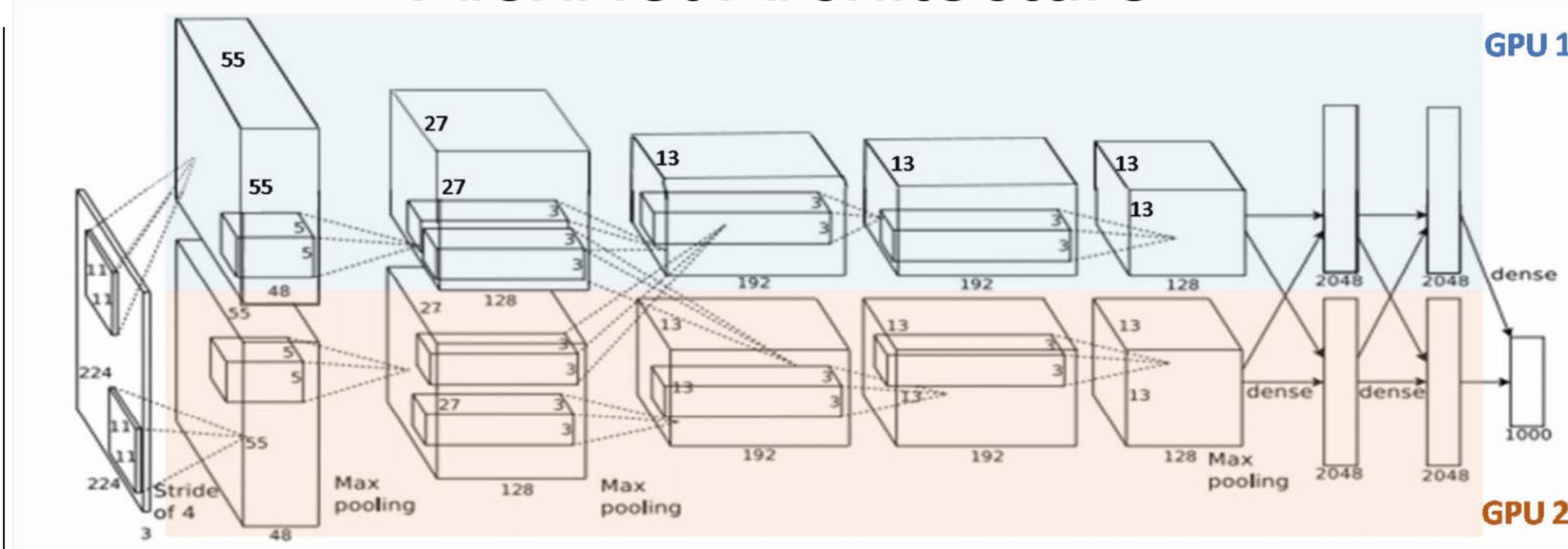


- Winner of ImageNet LSVRC-2010.
- Trained over 1.2M images using SGD with regularization.
- Deep architecture (60M parameters.)
- Optimized GPU implementation (cuda-convnet)

Krizhevsky, Alex, Ilya Sutskever, and Geoffrey E. Hinton. "Imagenet classification with deep convolutional neural networks." *NIPS* 2012.

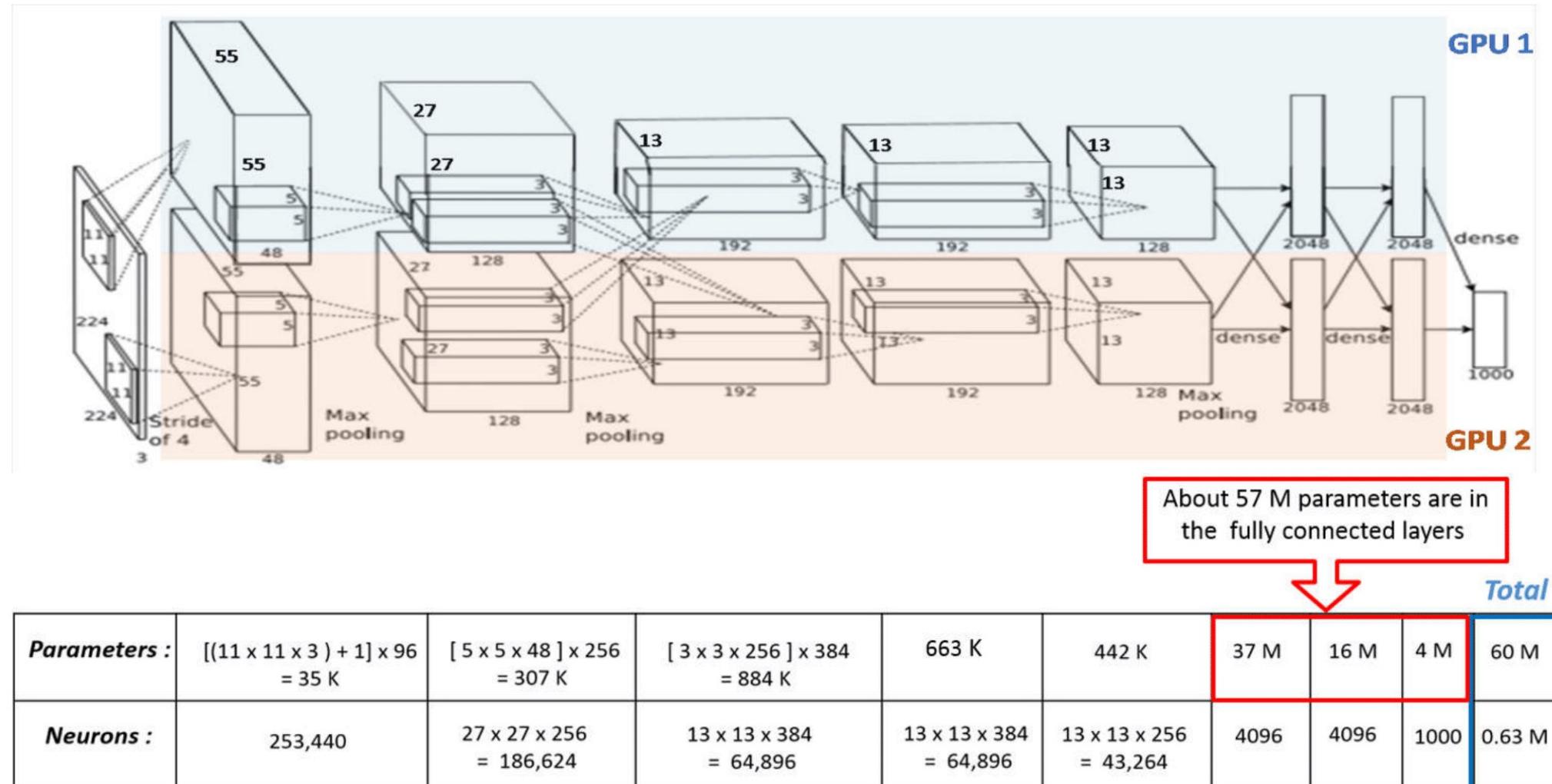


# AlexNet Architecture

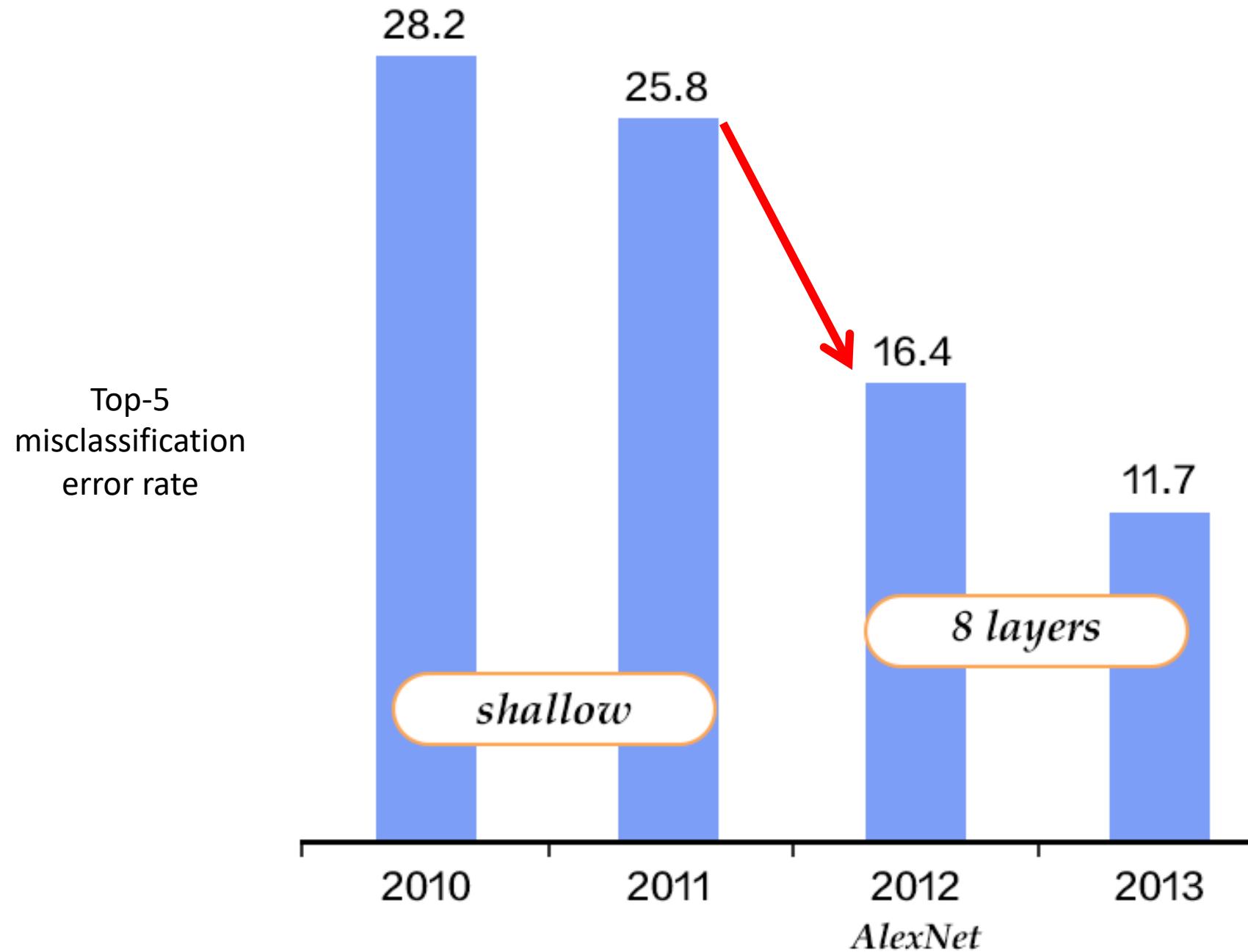




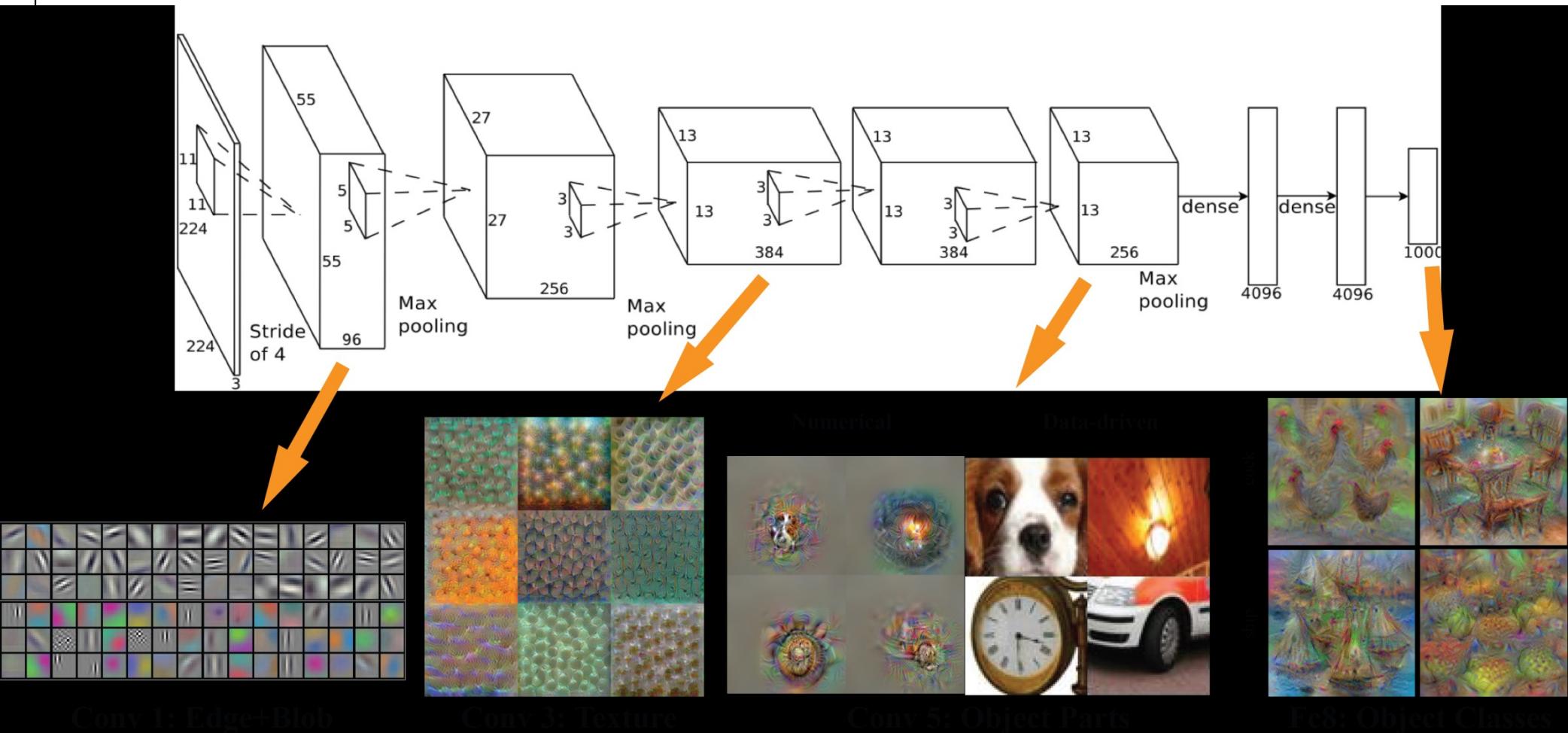
# AlexNet Architecture



- Convolutional layers cumulatively contain about 90-95% of computation, only about 5% of the parameters
- Fully-connected layers contain about 95% of parameters.

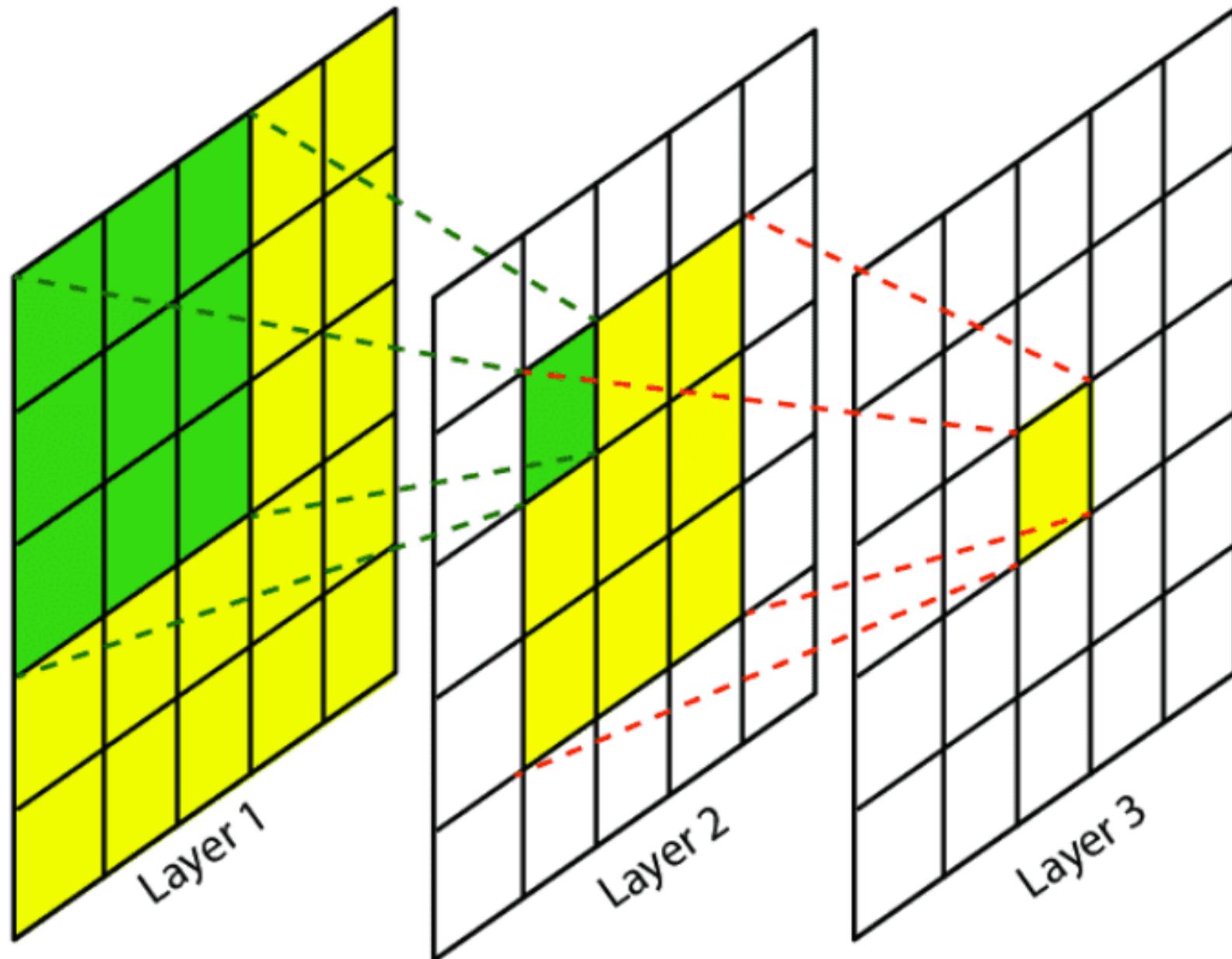


# Filters learnt by AlexNet



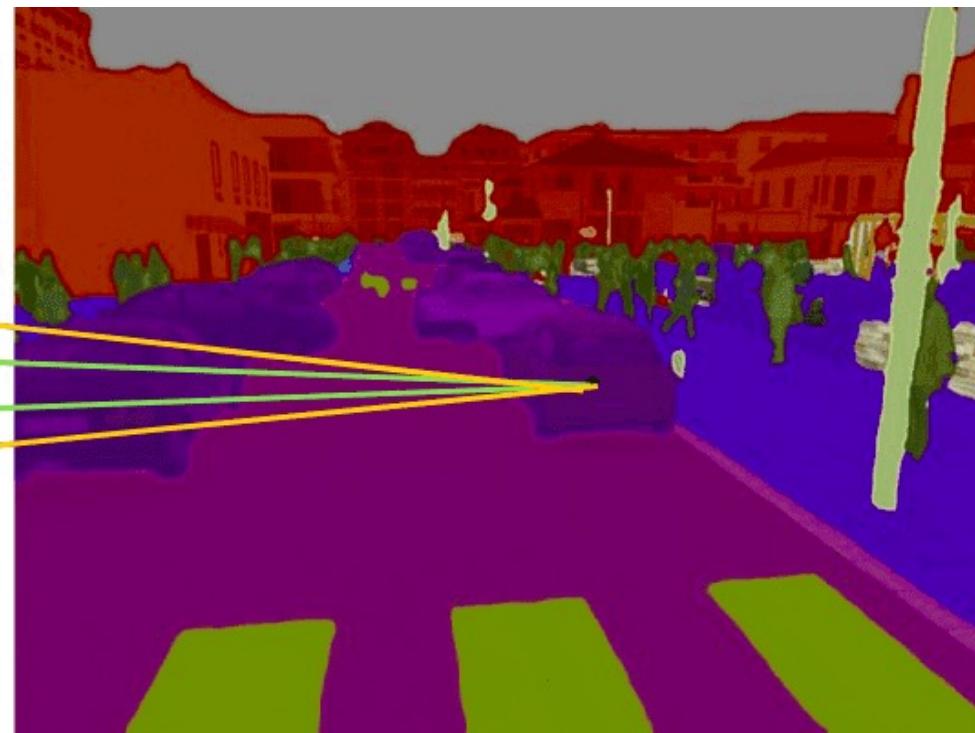


# Receptive Field



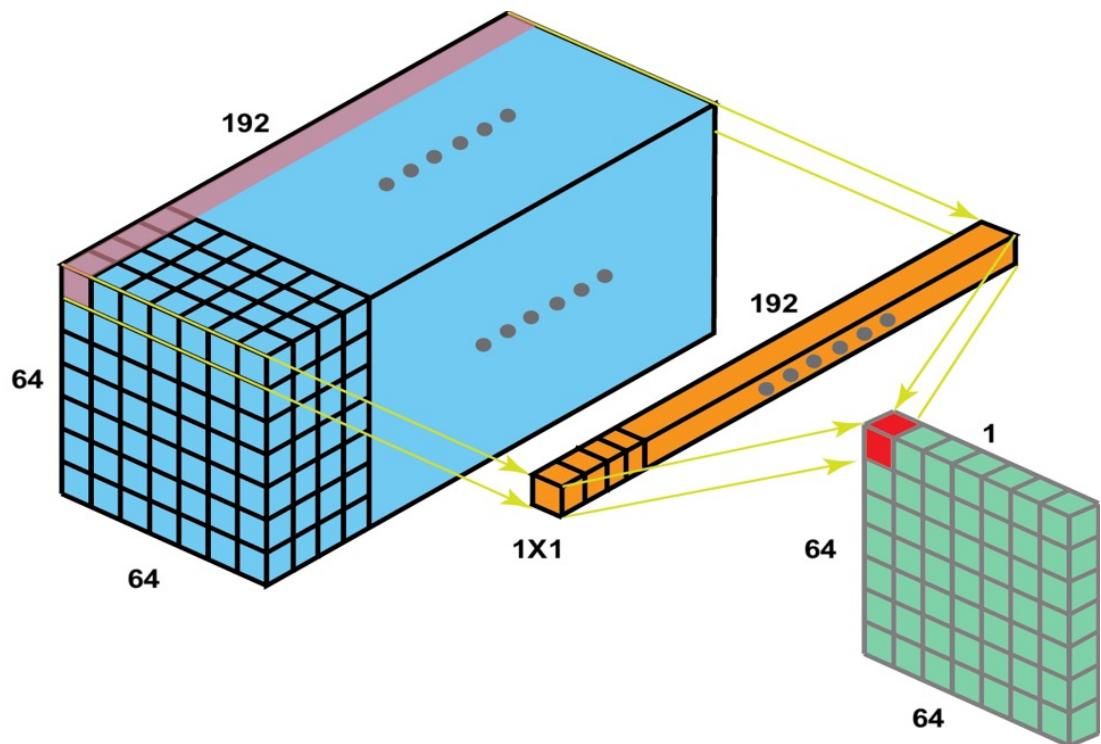


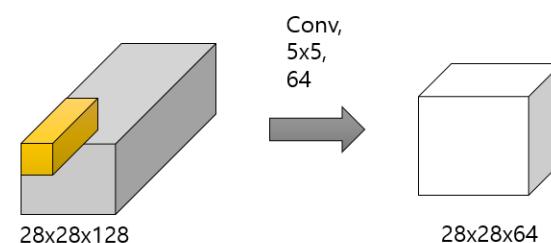
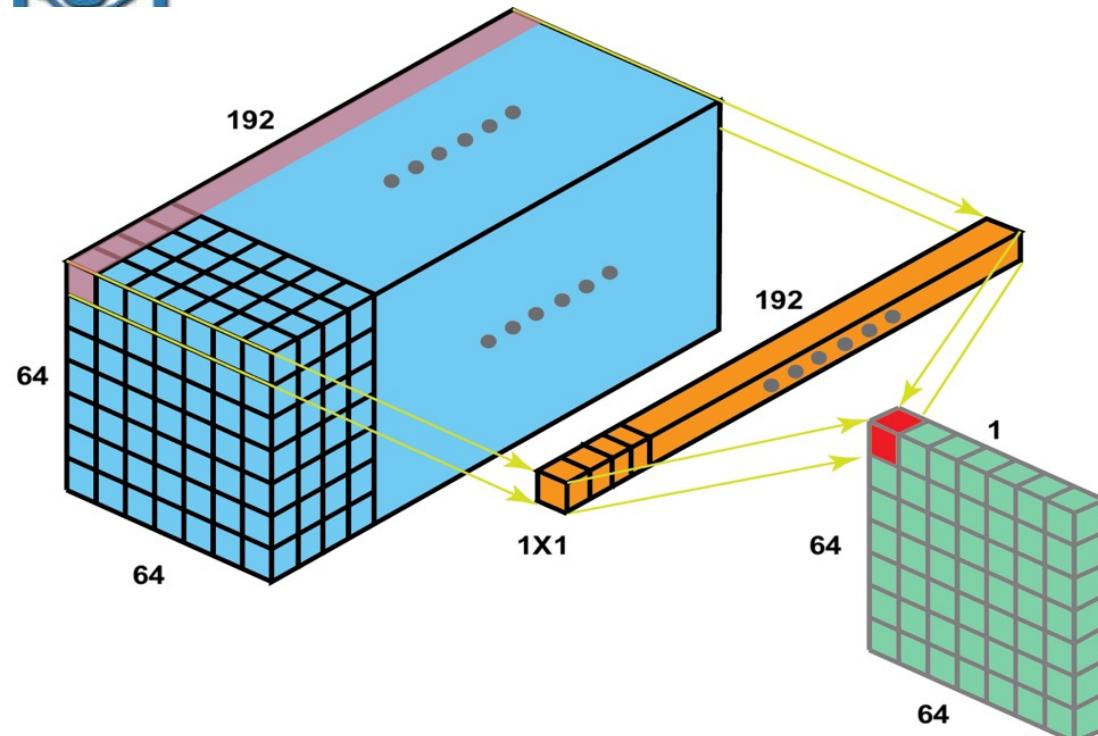
Opacity 69%  
Mask 0%



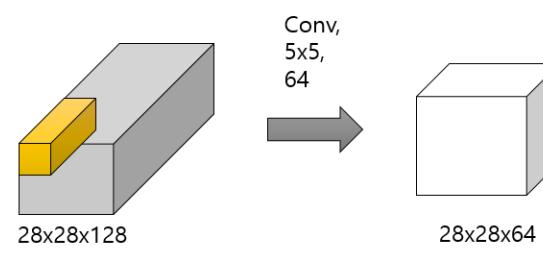
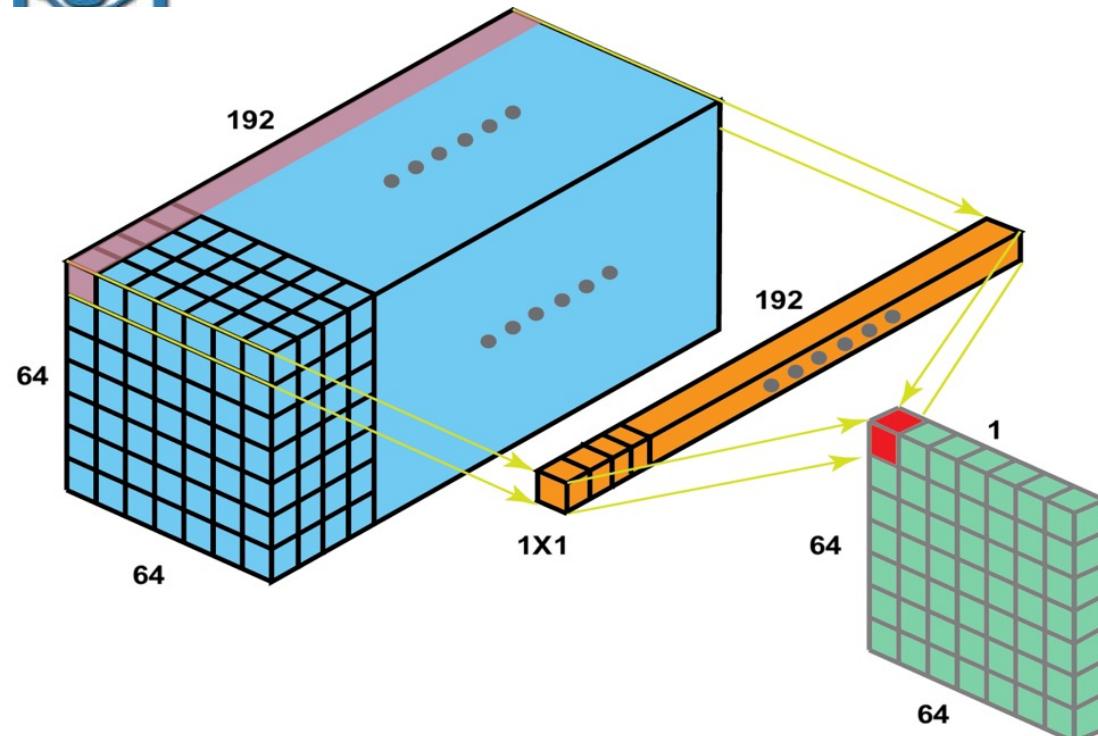


# 1x1 Convolutions

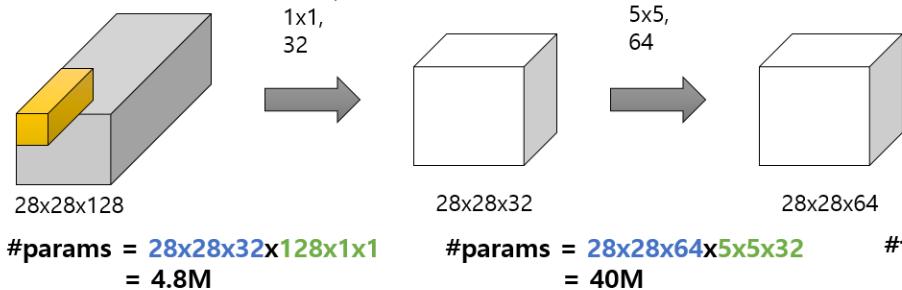




#params =  $28 \times 28 \times 64 \times 5 \times 5 \times 128 = 160M$



#params =  $28 \times 28 \times 64 \times 5 \times 5 \times 128 = 160M$



#params =  $28 \times 28 \times 32 \times 128 \times 1 \times 1 = 4.8M$

#params =  $28 \times 28 \times 64 \times 5 \times 5 \times 32 = 40M$  #total = 44.8M



---

# Dilated Convolutions

[https://github.com/vdumoulin/conv\\_arithmetic](https://github.com/vdumoulin/conv_arithmetic)



---

# Resources

- <https://www.youtube.com/@far1din619/videos>
- What is a convolution: <https://www.youtube.com/watch?v=KuXjwB4LzSA>
- [https://github.com/vdumoulin/conv\\_arithmetic](https://github.com/vdumoulin/conv_arithmetic)
- BP in CNN:  
[https://deeplearning.cs.cmu.edu/F21/document/recitation/Recitation5/CNN\\_Backprop\\_Recitation\\_5\\_F21.pdf](https://deeplearning.cs.cmu.edu/F21/document/recitation/Recitation5/CNN_Backprop_Recitation_5_F21.pdf)