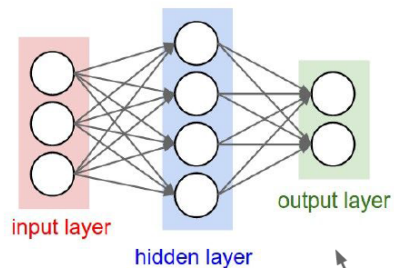# ELEG404/604: Imaging & Deep Learning

Gonzalo R. Arce

**Department of Electrical and Computer Engineering**
**University of Delaware**

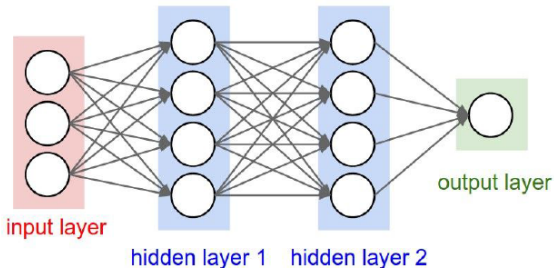Convolutional Neural Networks

# Neural Networks Architectures



"2-layer Neural Net", or
"1-hidden-layer Neural Net"

**"Fully-connected" layers**

"3-layer Neural Net", or
"2-hidden-layer Neural Net"

$4 + 2 = 6$ neurons.
$[3 \times 4] + [4 \times 2] = 20$ weights
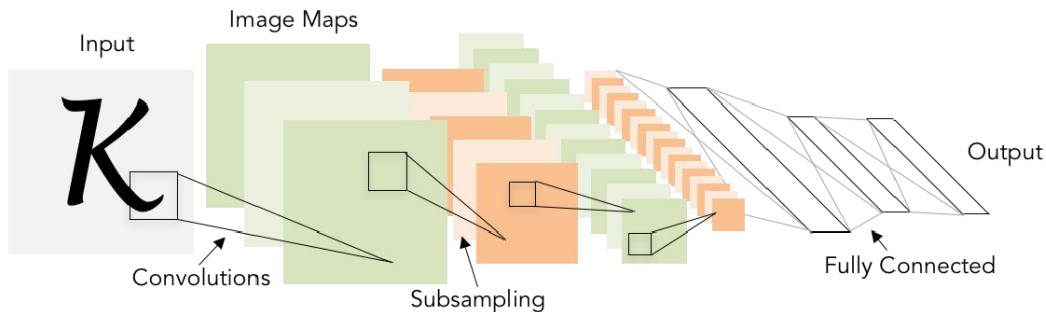$4 + 2 = 6$ biases.

$4 + 4 + 1 = 9$ neurons.
$[3 \times 4] + [4 \times 4] + [4 \times 1] = 32$ weights
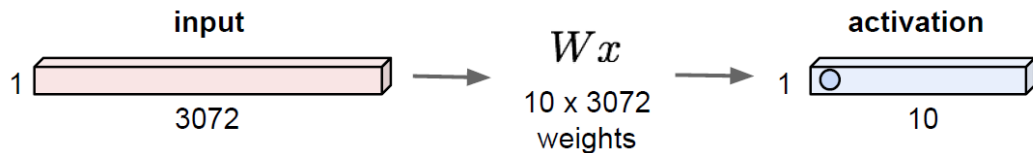$4 + 4 + 1 = 9$ biases.

# Convolutional Neural Networks Architectures

- ▶ Very similar to ordinary Neural Networks.

- ▶ Add convolutional layers. Neurons with 3 dimensions: width, height and depth.

- ▶ Inputs are also volumes.
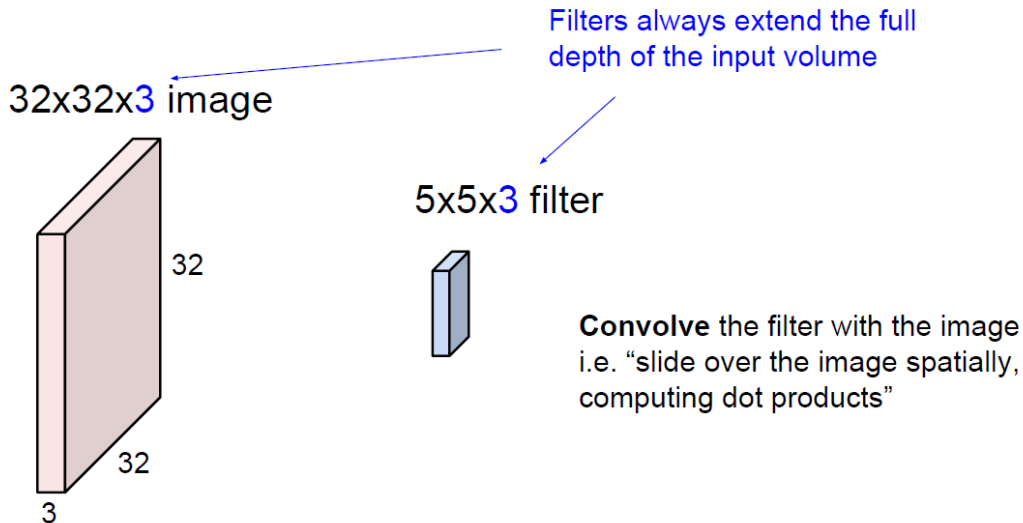
# Neural Network - Fully Connected (FC) Layer

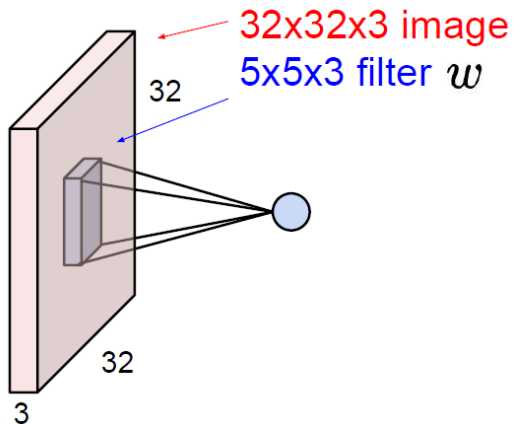Consider a $32 \times 32 \times 3$ image $\rightarrow$ stretch to $3072 \times 1$



Each output is the result of a dot product between a row of **W** and the input **x**. 10 neurons outputs.

# Convolutional Layer

Consider a $32 \times 32 \times 3$ image $\rightarrow$ preserve spatial structure.

Filters always extend the full depth of the input volume

32x32x3 image

5x5x3 filter

32

32

3

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

# Convolutional Layer



32x32x3 image
5x5x3 filter $w$

**Result:** dot product between the filter and a small $5 \times 5 \times 3$ chunk of the image.

Volume convolution at $(x, y)$, for **all** maps of the input volume:

$$\text{conv}_{x,y} = \sum_i w_i v_i$$

where $w$s are kernel weights, $v$s chuck of the image.
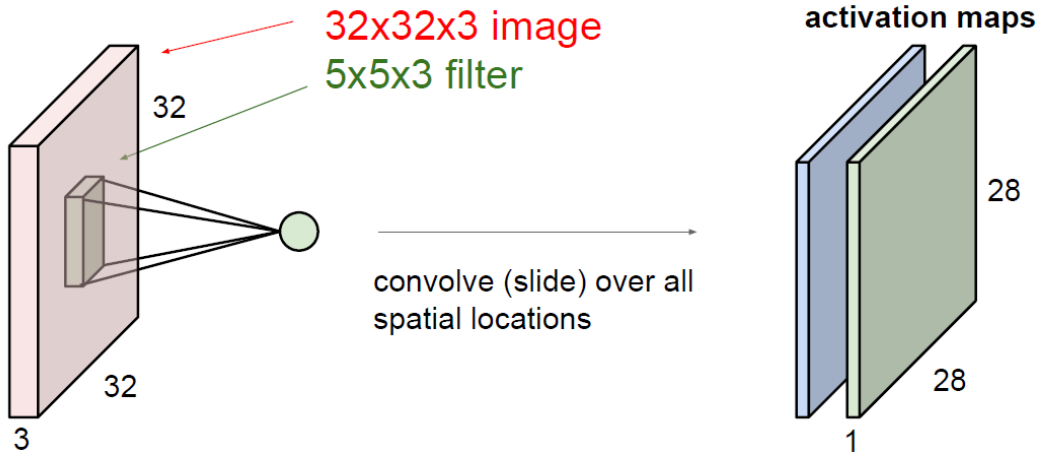
Adding scalar bias $b$:

$$z_{x,y} = \sum_i w_i v_i + b$$

# Convolutional Layer



32x32x3 image
5x5x3 filter

32

32

3

convolve (slide) over all
spatial locations

**activation map**

28

28

1

# Convolutional Layer

Consider a second, green filter:



32x32x3 image
5x5x3 filter

convolve (slide) over all spatial locations

activation maps

# Convolutional Layer

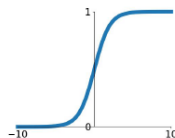Consider 6 filters ($5 \times 5$), we get 6 separate activation maps:



We stack these up to get a "new image volume" of size $28 \times 28 \times 6$

# Activation Functions

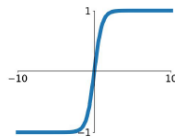Pass every element of each activation map through a nonlinearity:
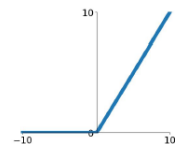
**Sigmoid**
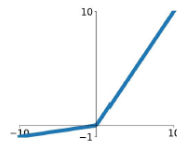$\sigma(x) = \frac{1}{1+e^{-x}}$

**tanh**
$\tanh(x)$

**ReLU**
$\max(0, x)$

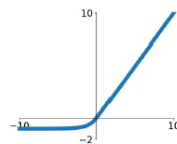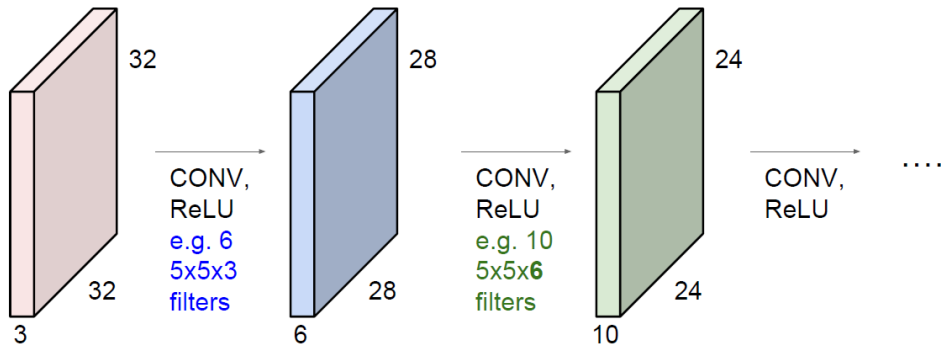**Leaky ReLU**
$\max(0.1x, x)$

**Maxout**
$\max(w_1^T x + b_1, w_2^T x + b_2)$

**ELU**
$$\begin{cases} x & x \geq 0 \\ \alpha(e^x - 1) & x < 0 \end{cases}$$
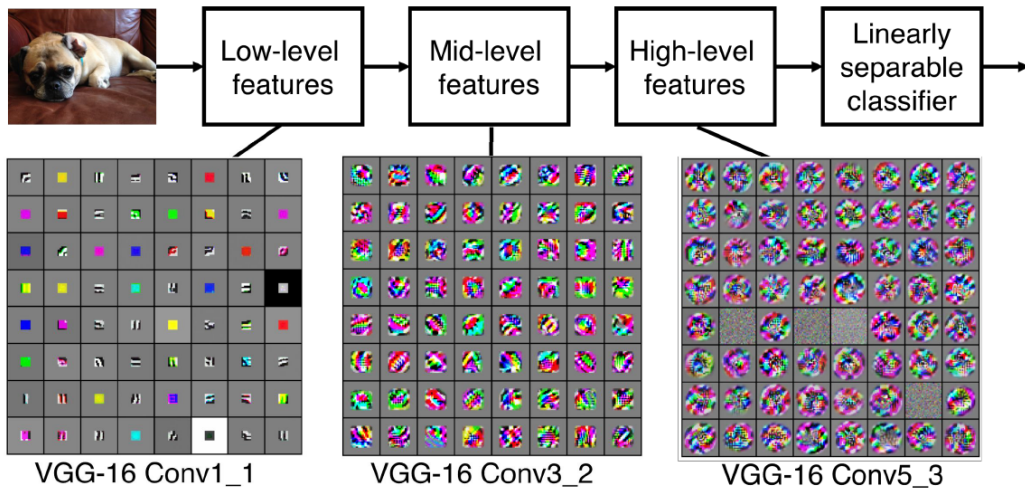
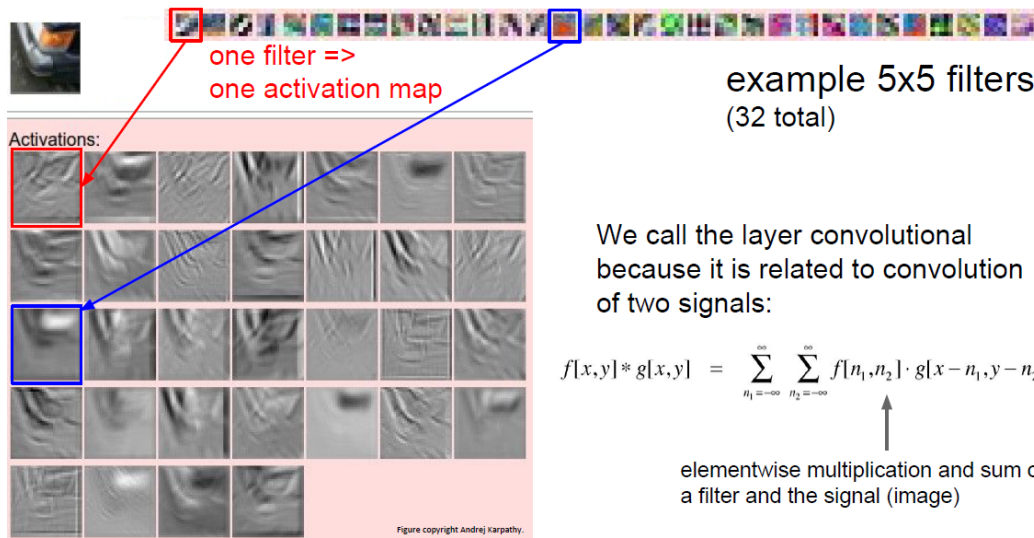ConvNet is a sequence of Convolutional Layers, interspersed with activation functions:



Notice how the activation maps get smaller, this can be solved by zero padding.

# Interpretation

Filters Learned:



VGG-16 Conv1_1          VGG-16 Conv3_2          VGG-16 Conv5_3

# Interpretation



one filter =>
one activation map

example 5x5 filters
(32 total)

We call the layer convolutional
because it is related to convolution
of two signals:

$$f[x,y] * g[x,y] = \sum_{n_1=-\infty}^{\infty} \sum_{n_2=-\infty}^{\infty} f[n_1, n_2] \cdot g[x-n_1, y-n_2]$$

elementwise multiplication and sum of
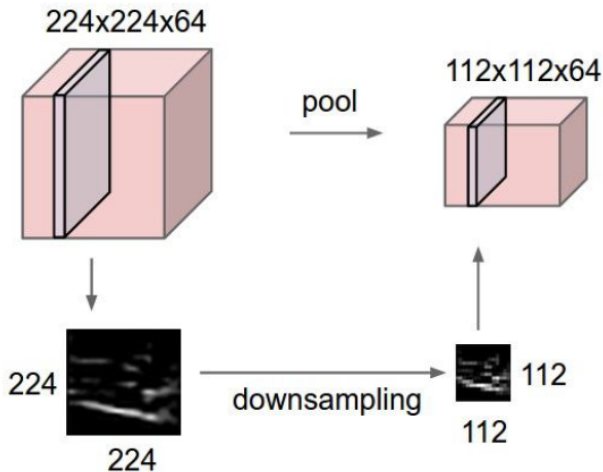a filter and the signal (image)

Activations:
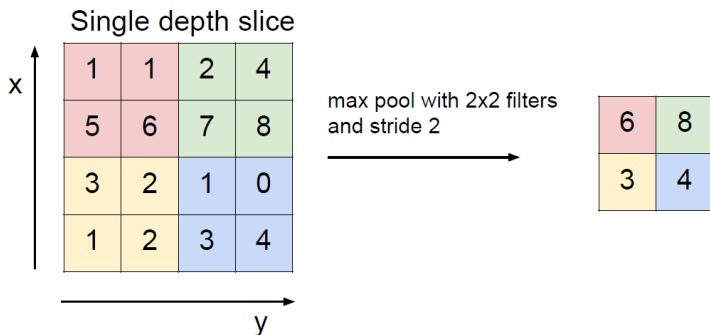
Figure copyright Andrej Karpathy.

# Pooling Layer

▶ Makes the representations smaller and more manageable.

▶ Operates over each activation map independently.
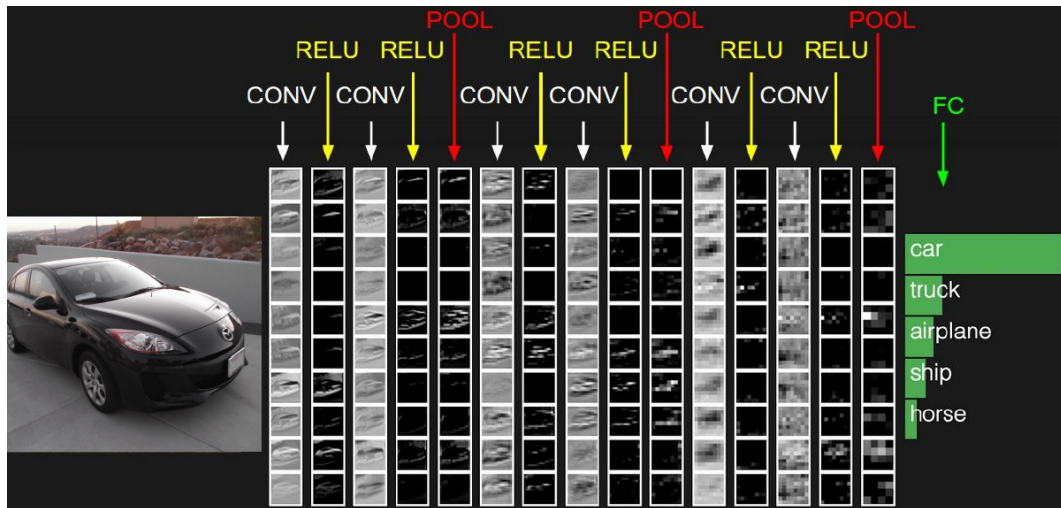
▶ Neighborhood of $2 \times 2$ is replaced by the average.

# Max Pooling

▶ Neighborhood of $2 \times 2$ is replaced by the maximum value.

▶ Effective in classifying large image databases.

▶ Simple and fast.



Single depth slice

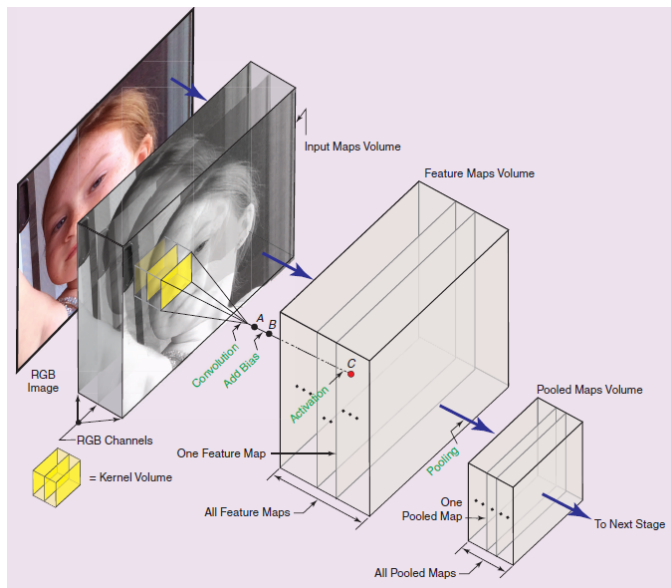max pool with 2x2 filters and stride 2

▶ $L_2$ pooling is also used. Neighborhood of $2 \times 2$ is replaced by the squared root of the sum of their squared values.
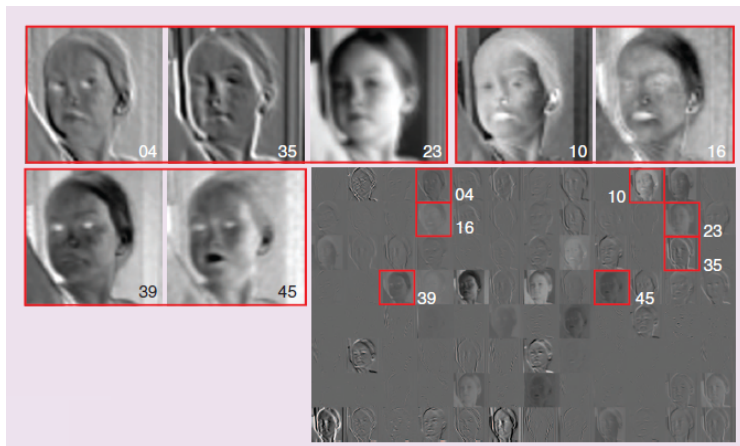
# Example - Image classification

# Convolutional Neural Networks Complete Scheme
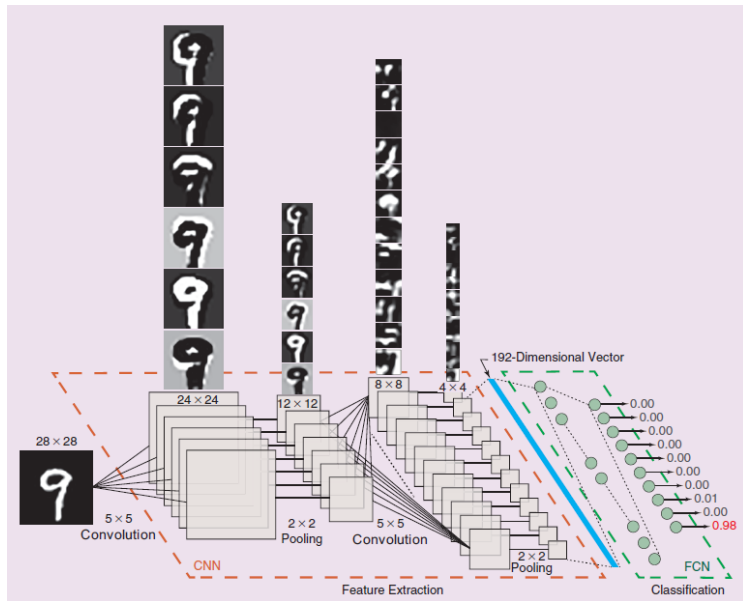


- $277 \times 277$ pixels RGB image.
- 96 feature maps.
- 96 kernels volumes of size $11 \times 11 \times 3$
- This weights came from AlexNet: CNN trained using more than 1 million images belonging to 1,000 object categories.
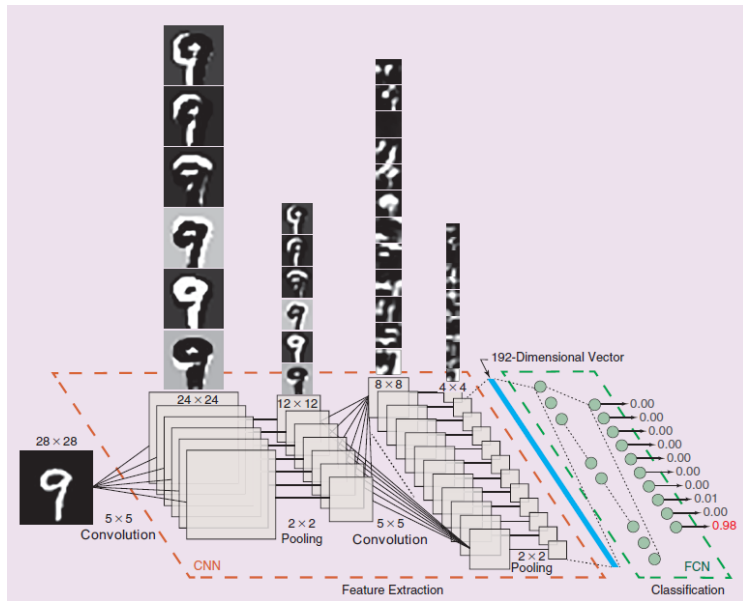
# Result Feature Maps



(4) and (35) emphasize edge content. (23) is a blurred version of the input. (10) and (16) capture complementary shades of gray (hair). (39) emphasizes eyes and dress (blue). (45) blue and red tones (lips, hair, skin).

# Example - Handwritten Numerals Classification



- ▶ Training: 60,000 grayscale images.
- ▶ Testing: 10,000 grayscale images.
- ▶ Network trained for 200 epochs.
- ▶ Performance: 99.4% in training set.
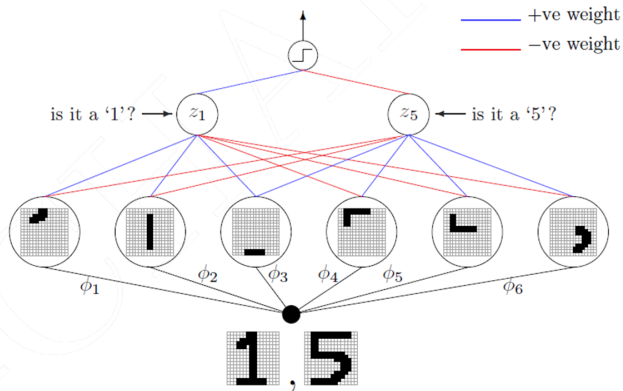- ▶ Performance: 99.1% in testing set.

# Example - Handwritten Numerals Classification



- ▶ First stage: 6 features maps.
- ▶ Second stage: 12 features maps.
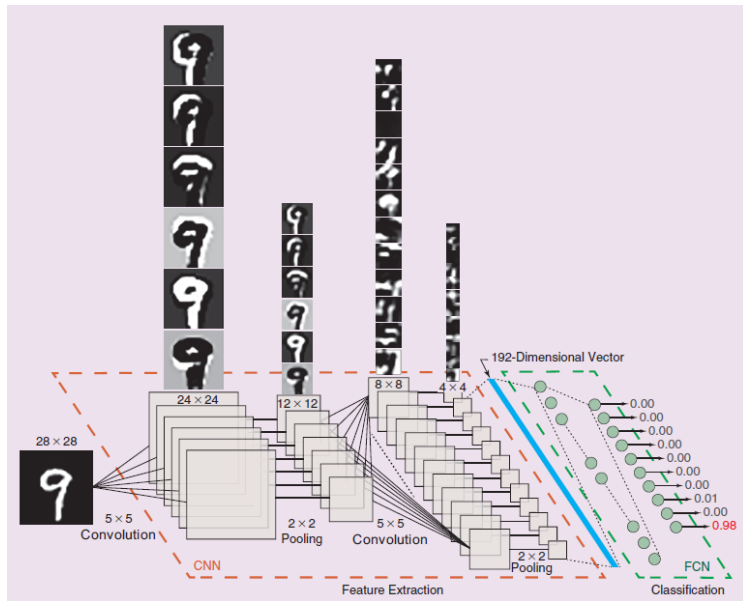- ▶ Kernels of size $5 \times 5$.
- ▶ Fully Connected Layer without hidden layers.

# Remember: Networks with many layers - Example

$\phi_i$ is feature function which computes the presence $(+1)$ and absence $(-1)$ of the corresponding feature.



If we feed in '1', $\phi_1, \phi_2, \phi_3$ compute $+1$ and $\phi_4, \phi_5, \phi_6$ compute $-1$. Combining with the signs of the weights, $z_1$ will be positive and $z_5$ will be negative.
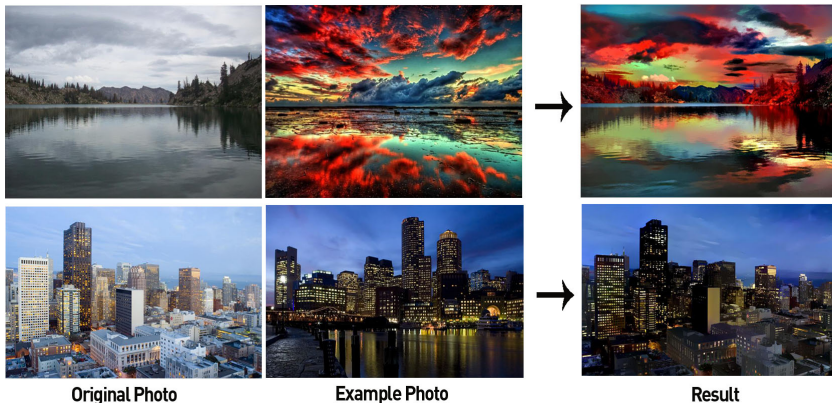
# Features Map Interpretation



- First feature map: strong vertical components on the left.
- Second: strong components in the northwest area of the top of the character and the left vertical lower area.
- Third: strong horizontal components.

# Style Transfer

- ▶ Goal: Rendering the semantic content of an image in different styles.
- ▶ Challenge: separate image content from style.

A Neural Algorithm of Artistic Style can separate and recombine the image content and style of natural images.



Original Photo      Example Photo      Result

## Deep Image Representations

VGG-19 is a convolutional neural network that is trained on more than a million images from the ImageNet database to perform object recognition (1000 categories) and localization.