# Separation between NN and kernel

- For approximation and optimization, neural network has no advantage over kernel. Why NN gives better performance: generalization.

- [Allen-Zhu and Li '20] Construct a class of functions $\mathscr{F}$ such that $y = f(x)$ for some $f \in \mathscr{F}$:
  - no kernel is sample-efficient;
  - Exists a neural network that is sample-efficient.

# Separation between NN and kernel

# Separation between NN and kernel

# Separation between NN and kernel

# Separation between NN and kernel

# Separation between NN and kernel

# Separation between NN and kernel

# Separation between NN and kernel

# Convolutional Neural Networks

# Multi-layer Neural Network

$$a^{(1)} = x$$

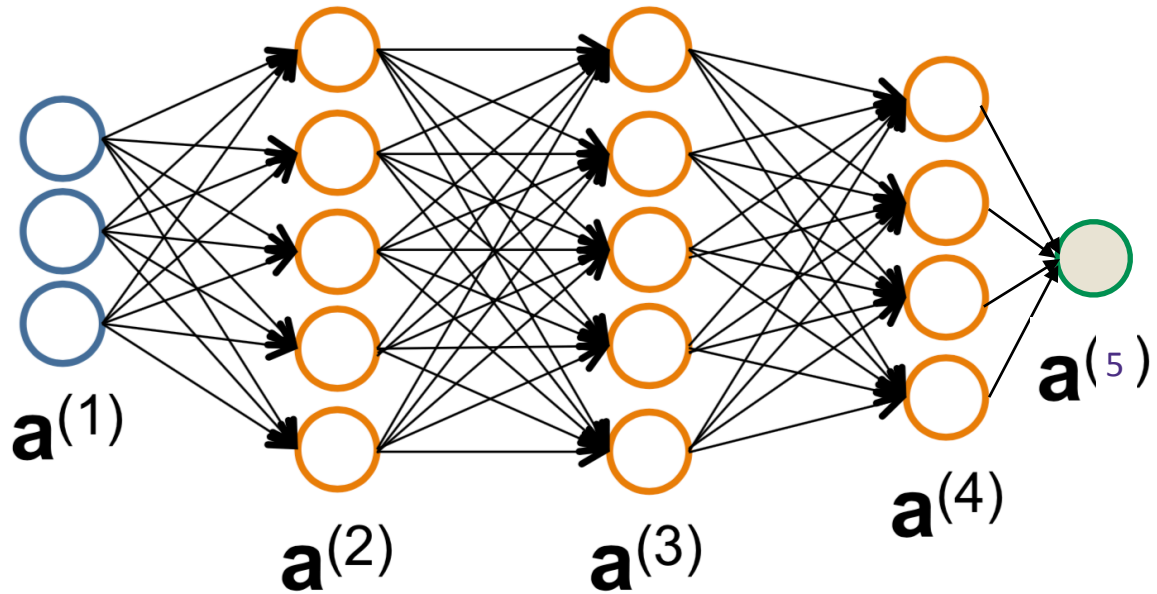$$z^{(2)} = \Theta^{(1)} a^{(1)}$$

$$a^{(2)} = g\left(z^{(2)}\right)$$

$$\vdots$$

$$z^{(l+1)} = \Theta^{(l)} a^{(l)}$$

$$a^{(l+1)} = g\left(z^{(l+1)}\right)$$

$$\vdots$$

$$\widehat{y} = a^{(L+1)}$$



$\mathbf{a}^{(1)}$    $\mathbf{a}^{(2)}$    $\mathbf{a}^{(3)}$    $\mathbf{a}^{(4)}$    $\mathbf{a}^{(5)}$
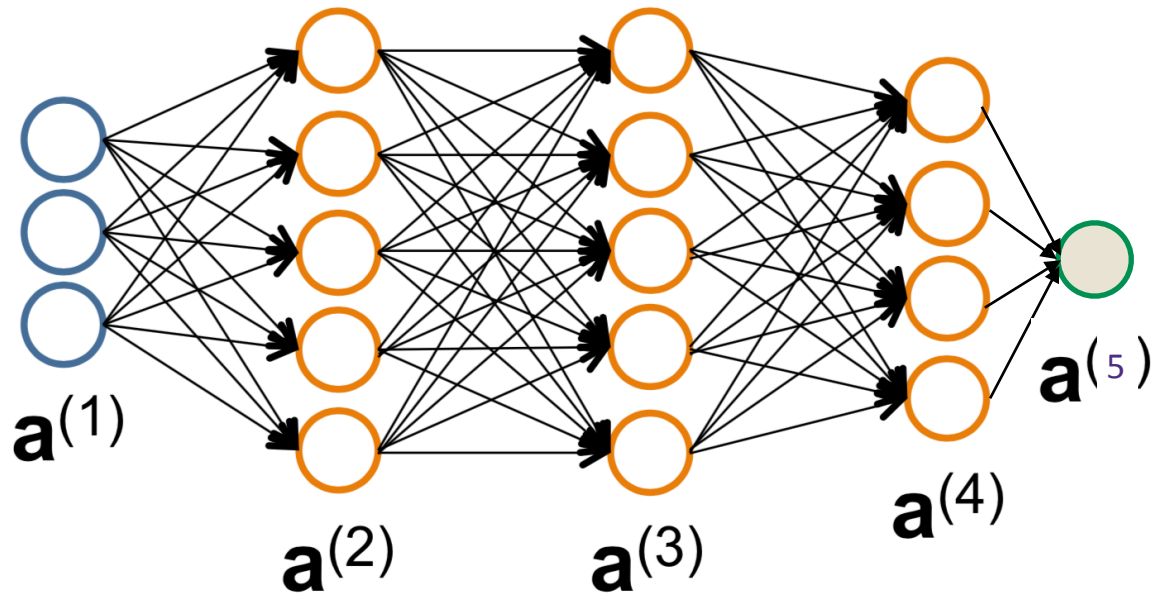
$$L(y, \widehat{y}) = y \log(\widehat{y}) + (1 - y)\log(1 - \widehat{y})$$

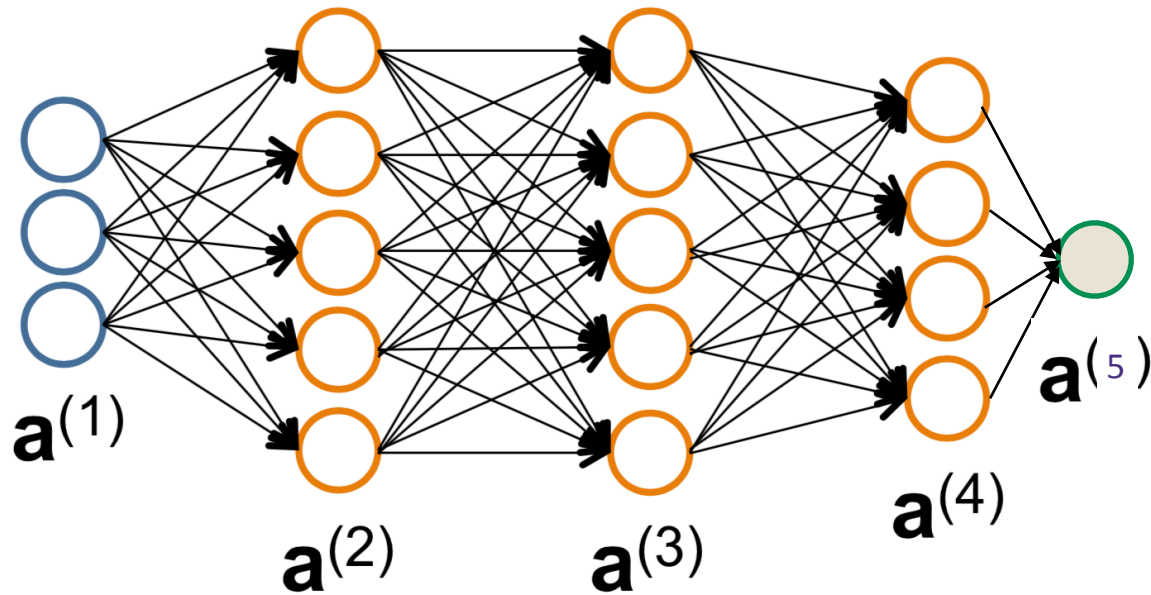$$g(z) = \frac{1}{1 + e^{-z}}$$

Binary Logistic Regression

# Neural Network Architecture

The neural network architecture is defined by the number of layers, and the number of nodes in each layer, but also by **allowable edges**.

# Neural Network Architecture

The neural network architecture is defined by the number of layers, and the number of nodes in each layer, but also by **allowable edges**.



$\mathbf{a}^{(1)}$ $\mathbf{a}^{(2)}$ $\mathbf{a}^{(3)}$ $\mathbf{a}^{(4)}$ $\mathbf{a}^{(5)}$
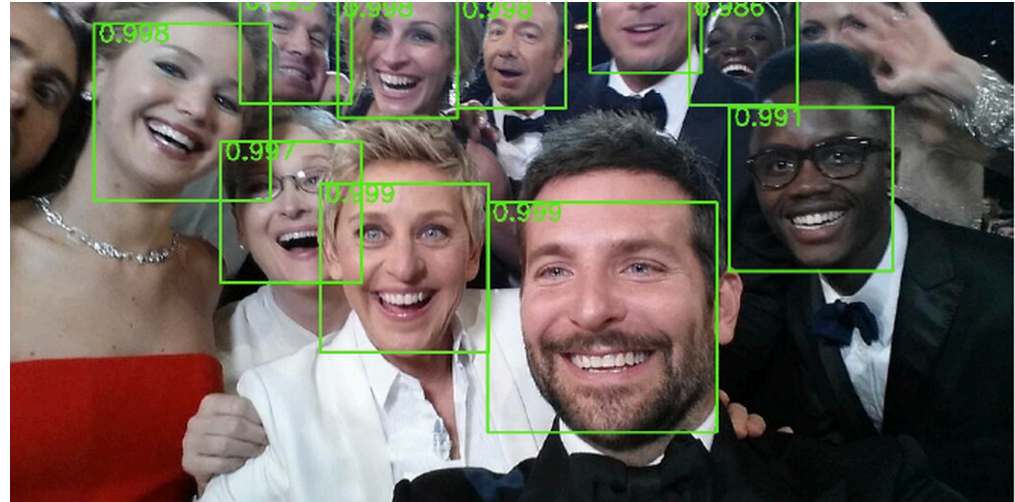
We say a layer is **Fully Connected (FC)** if all linear mappings from the current layer to the next layer are permissible.

$$\mathbf{a}^{(k+1)} = g(\Theta \mathbf{a}^{(k)}) \quad \text{for any } \Theta \in \mathbb{R}^{n_{k+1} \times n_k}$$

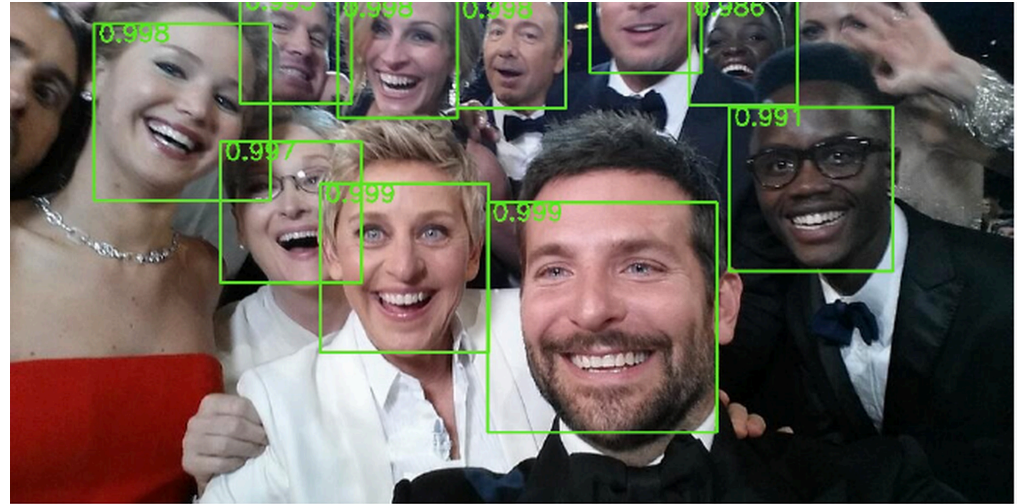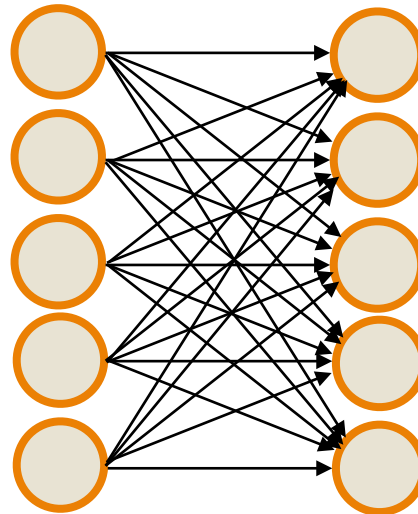A lot of parameters!! $\quad n_1 n_2 + n_2 n_3 + \cdots + n_L n_{L+1}$

# Neural Network Architecture

Objects are often **localized in space** so to find the faces in an image, not every pixel is important for classification—makes sense to drag a window across an image.

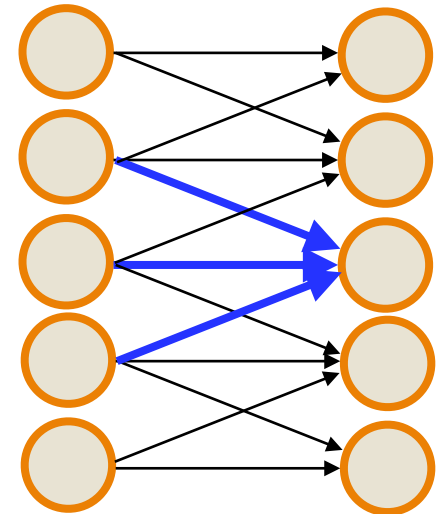# Neural Network Architecture

Objects are often **localized in space** so to find the faces in an image, not every pixel is important for classification—makes sense to drag a window across an image.
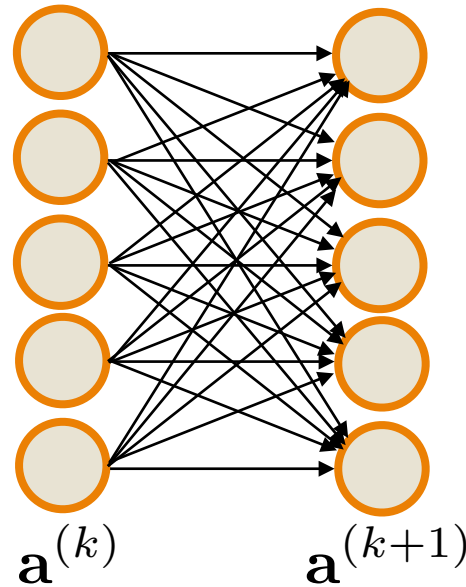
Similarly, to identify edges or other local structure, it makes sense to only look at **local information**
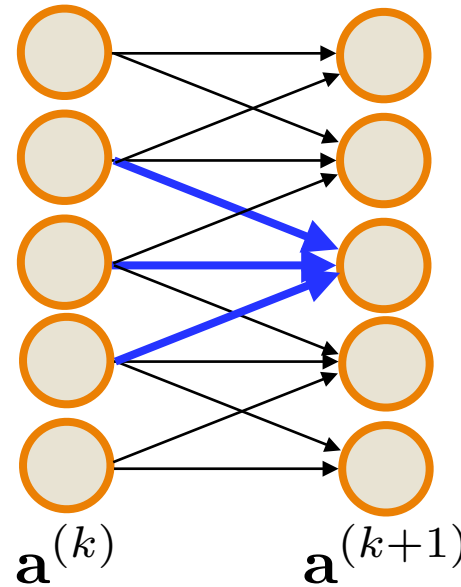
vs.

# Neural Network Architecture



vs.

$$\mathbf{a}^{(k)} \qquad \mathbf{a}^{(k+1)} \qquad\qquad \mathbf{a}^{(k)} \qquad \mathbf{a}^{(k+1)}$$

$$\begin{bmatrix} \Theta_{0,0} & \Theta_{0,1} & \Theta_{0,2} & \Theta_{0,3} & \Theta_{0,4} \\ \Theta_{1,0} & \Theta_{1,1} & \Theta_{1,2} & \Theta_{1,3} & \Theta_{1,4} \\ \Theta_{2,0} & \Theta_{2,1} & \Theta_{2,2} & \Theta_{2,3} & \Theta_{2,4} \\ \Theta_{3,0} & \Theta_{3,1} & \Theta_{3,2} & \Theta_{3,3} & \Theta_{3,4} \\ \Theta_{4,0} & \Theta_{4,1} & \T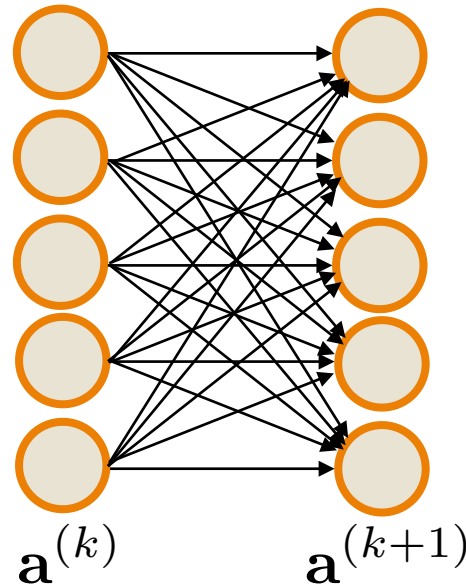heta_{4,2} & \Theta_{4,3} & \Theta_{4,4} \end{bmatrix} \qquad \begin{bmatrix} \Theta_{0,0} & \Theta_{0,1} & 0 & 0 & 0 \\ \Theta_{1,0} & \Theta_{1,1} & \Theta_{1,2} & 0 & 0 \\ 0 & \Theta_{2,1} & \Theta_{2,2} & \Theta_{2,3} & 0 \\ 0 & 0 & \Theta_{3,2} & \Theta_{3,3} & \Theta_{3,4} \\ 0 & 0 & 0 & \Theta_{4,3} & \Theta_{4,4} \end{bmatrix}$$
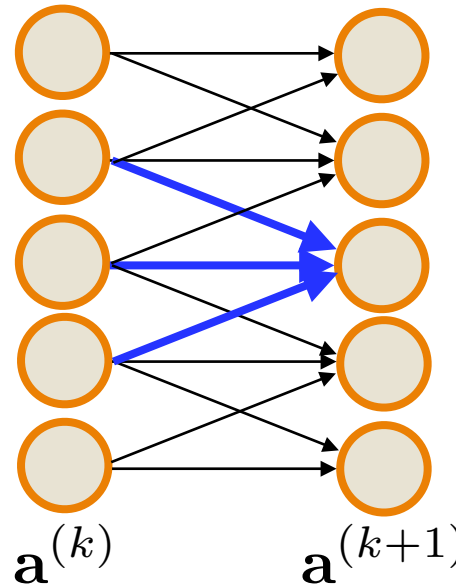
Parameters: $\qquad n^2 \qquad\qquad\qquad 3n-2$

$$\mathbf{a}_i^{(k+1)} = g\left( \sum_{j=0}^{n-1} \Theta_{i,j} \mathbf{a}_j^{(k)} \right)$$

# Neural Network Architecture



vs.

**Mirror/share local weights everywhere (e.g., structure equally likely to be anywhere in image)**

$$\mathbf{a}^{(k)} \qquad \mathbf{a}^{(k+1)} \qquad\qquad \mathbf{a}^{(k)} \qquad \mathbf{a}^{(k+1)}$$

$$\begin{bmatrix} \Theta_{0,0} & \Theta_{0,1} & \Theta_{0,2} & \Theta_{0,3} & \Theta_{0,4} \\ \Theta_{1,0} & \Theta_{1,1} & \Theta_{1,2} & \Theta_{1,3} & \Theta_{1,4} \\ \Theta_{2,0} & \Theta_{2,1} & \Theta_{2,2} & \Theta_{2,3} & \Theta_{2,4} \\ \Theta_{3,0} & \Theta_{3,1} & \Theta_{3,2} & \Theta_{3,3} & \Theta_{3,4} \\ \Theta_{4,0} & \Theta_{4,1} & \Theta_{4,2} & \Theta_{4,3} & \Theta_{4,4} \end{bmatrix} \quad \begin{bmatrix} \Theta_{0,0} & \Theta_{0,1} & 0 & 0 & 0 \\ \Theta_{1,0} & \Theta_{1,1} & \Theta_{1,2} & 0 & 0 \\ 0 & \Theta_{2,1} & \Theta_{2,2} & \Theta_{2,3} & 0 \\ 0 & 0 & \Theta_{3,2} & \Theta_{3,3} & \Theta_{3,4} \\ 0 & 0 & 0 & \Theta_{4,3} & \Theta_{4,4} \end{bmatrix} \quad \begin{bmatrix} \theta_1 & \theta_2 & 0 & 0 & 0 \\ \theta_0 & \theta_1 & \theta_2 & 0 & 0 \\ 0 & \theta_0 & \theta_1 & \theta_2 & 0 \\ 0 & 0 & \theta_0 & \theta_1 & \theta_2 \\ 0 & 0 & 0 & \theta_0 & \theta_1 \end{bmatrix}$$

Parameters: $\qquad n^2 \qquad\qquad\qquad 3n-2 \qquad\qquad\qquad 3$

$$\mathbf{a}_i^{(k+1)} = g\left( \sum_{j=0}^{n-1} \Theta_{i,j} \mathbf{a}_j^{(k)} \right) \qquad\qquad \mathbf{a}_i^{(k+1)} = g\left( \sum_{j=0}^{m-1} \theta_j \mathbf{a}_{i+j}^{(k)} \right)$$

# Neural Network Architecture

**Fully Connected (FC) Layer**

$$\begin{bmatrix} \Theta_{0,0} & \Theta_{0,1} & \Theta_{0,2} & \Theta_{0,3} & \Theta_{0,4} \\ \Theta_{1,0} & \Theta_{1,1} & \Theta_{1,2} & \Theta_{1,3} & \Theta_{1,4} \\ \Theta_{2,0} & \Theta_{2,1} & \Theta_{2,2} & \Theta_{2,3} & \Theta_{2,4} \\ \Theta_{3,0} & \Theta_{3,1} & \Theta_{3,2} & \Theta_{3,3} & \Theta_{3,4} \\ \Theta_{4,0} & \Theta_{4,1} & \Theta_{4,2} & \Theta_{4,3} & \Theta_{4,4} \end{bmatrix}$$

**Convolutional (CONV) Layer (1 filter)**

$$\begin{bmatrix} \theta_1 & \theta_2 & 0 & 0 & 0 \\ \theta_0 & \theta_1 & \theta_2 & 0 & 0 \\ 0 & \theta_0 & \theta_1 & \theta_2 & 0 \\ 0 & 0 & \theta_0 & \theta_1 & \theta_2 \\ 0 & 0 & 0 & \theta_0 & \theta_1 \end{bmatrix} \quad \text{m=3}$$

$$\mathbf{a}_i^{(k+1)} = g\left( \sum_{j=0}^{n-1} \Theta_{i,j} \mathbf{a}_j^{(k)} \right)$$

$$\mathbf{a}_i^{(k+1)} = g\left( \sum_{j=0}^{m-1} \theta_j \mathbf{a}_{i+j}^{(k)} \right) = g([\theta * \mathbf{a}^{(k)}]_i)$$

Convolution*

$$\theta = (\theta_0, \ldots, \theta_{m-1}) \in \mathbb{R}^m \text{ is referred to as a "filter"}$$

# Example (1d convolution)

$$(\theta * x)_i = \sum_{j=0}^{m-1} \theta_j x_{i+j}$$

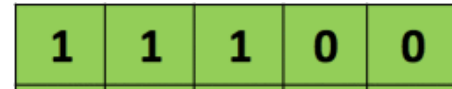| 1 | 1 | 1 | 0 | 0 |
|---|---|---|---|---|

Input $x \in \mathbb{R}^n$

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

Filter $\theta \in \mathbb{R}^m$

|  |  |  |
|---|---|---|

Output $\theta * x$

# Example (1d convolution)

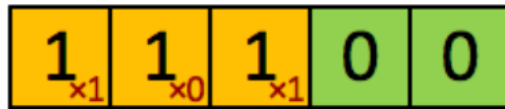$$(\theta * x)_i = \sum_{j=0}^{m-1} \theta_j x_{i+j}$$

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

Input $x \in \mathbb{R}^n$

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

Filter $\theta \in \mathbb{R}^m$

| $1_{\times 1}$ | $1_{\times 0}$ | $1_{\times 1}$ | 0 | 0 |
|---|---|---|---|---|

| 2 | | |
|---|---|---|

Output $\theta * x$

# Example (1d convolution)

$$(\theta * x)_i = \sum_{j=0}^{m-1} \theta_j x_{i+j}$$

| 1 | 1 | 1 | 0 | 0 |

Input $x \in \mathbb{R}^n$

| 1 | 0 | 1 |

Filter $\theta \in \mathbb{R}^m$

| 1 | 1$_{\times 1}$ | 1$_{\times 0}$ | 0$_{\times 1}$ | 0 |

| 2 | 1 | |

Output $\theta * x$

# Example (1d convolution)

$$(\theta * x)_i = \sum_{j=0}^{m-1} \theta_j x_{i+j}$$

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

Input $x \in \mathbb{R}^n$

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

Filter $\theta \in \mathbb{R}^m$

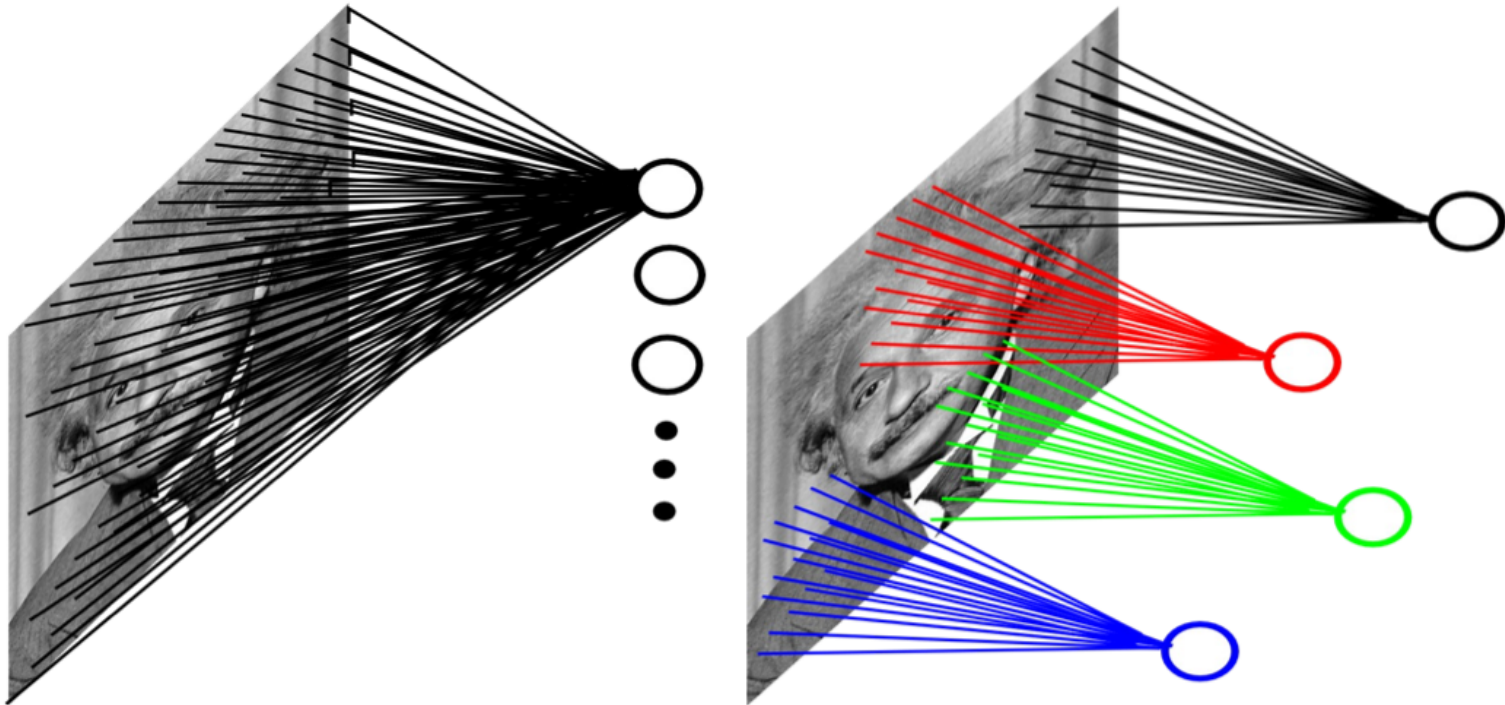| 1 | 1 | 1 ×1 | 0 ×0 | 0 ×1 |
|---|---|---|---|---|

| 2 | 1 | 1 |
|---|---|---|

Output $\theta * x$

# 2d Convolution Layer

Example: 200x200 image
- Fully-connected, 400,000 hidden units = 16 billion parameters
- Locally-connected, 400,000 hidden units 10x10 fields = 40 million params
- Local connections capture local dependencies

# Convolution of images (2d convolution)

$$(I * K)(i, j) = \sum_m \sum_n I(i + m, j + n) K(m, n)$$

| 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 $I$

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

Filter $K$

| 1 $\times 1$ | 1 $\times 0$ | 1 $\times 1$ | 0 | 0 |
|---|---|---|---|---|
| 0 $\times 0$ | 1 $\times 1$ | 1 $\times 0$ | 1 | 0 |
| 0 $\times 1$ | 0 $\times 0$ | 1 $\times 1$ | 1 | 1 |
| 0 | 0 | 1 | 1 | 0 |
| 0 | 1 | 1 | 0 | 0 |

Image

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

Convolved
Feature
$I * K$

# Convolution of images

$$(I * K)(i, j) = \sum_m \sum_n I(i + m, j + n) K(m, n)$$

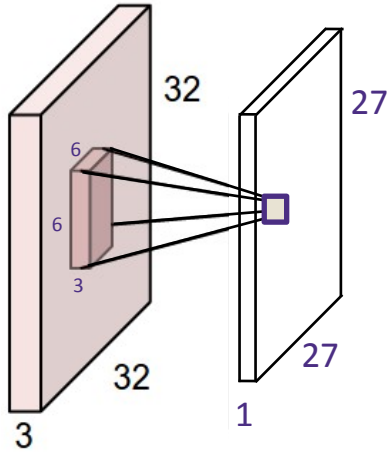Image $I$



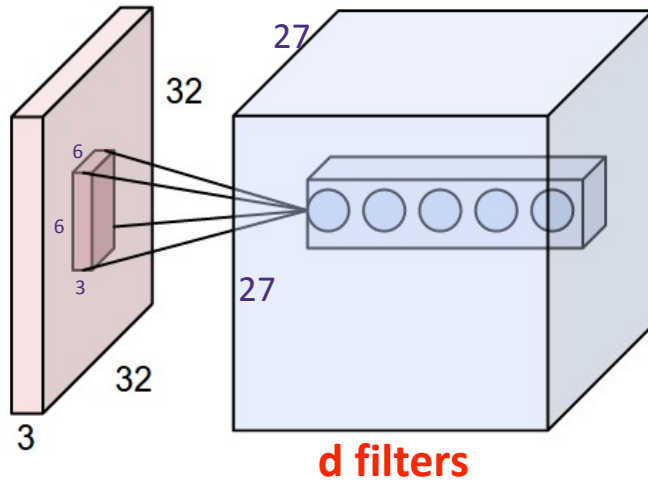| Operation | Filter $K$ | Convolved Image $I * K$ |
|---|---|---|
| Edge detection | $\begin{bmatrix} 1 & 0 & -1 \\ 0 & 0 & 0 \\ -1 & 0 & 1 \end{bmatrix}$ | |
| | $\begin{bmatrix} 0 & 1 & 0 \\ 1 & -4 & 1 \\ 0 & 1 & 0 \end{bmatrix}$ | |
| | $\begin{bmatrix} -1 & -1 & -1 \\ -1 & 8 & -1 \\ -1 & -1 & -1 \end{bmatrix}$ | |
| Sharpen | $\begin{bmatrix} 0 & -1 & 0 \\ -1 & 5 & -1 \\ 0 & -1 & 0 \end{bmatrix}$ | |
| Box blur (normalized) | $\frac{1}{9}\begin{bmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{bmatrix}$ | |
| Gaussian blur (approximation) | $\frac{1}{16}\begin{bmatrix} 1 & 2 & 1 \\ 2 & 4 & 2 \\ 1 & 2 & 1 \end{bmatrix}$ | |

# Stacking convolved images
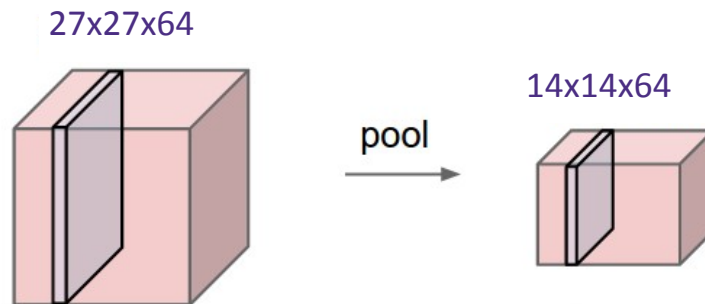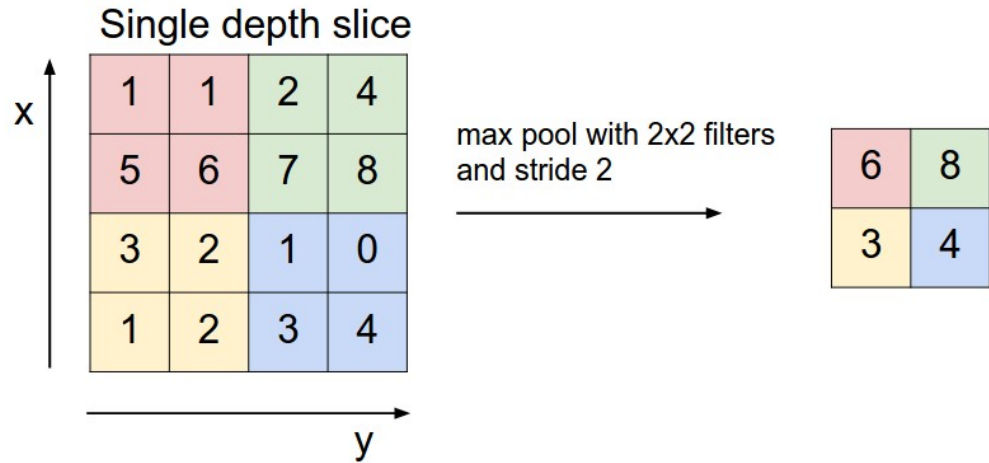


$$x \in \mathbb{R}^{n \times n \times r}$$

# Stacking convolved images



**Repeat with d filters!**
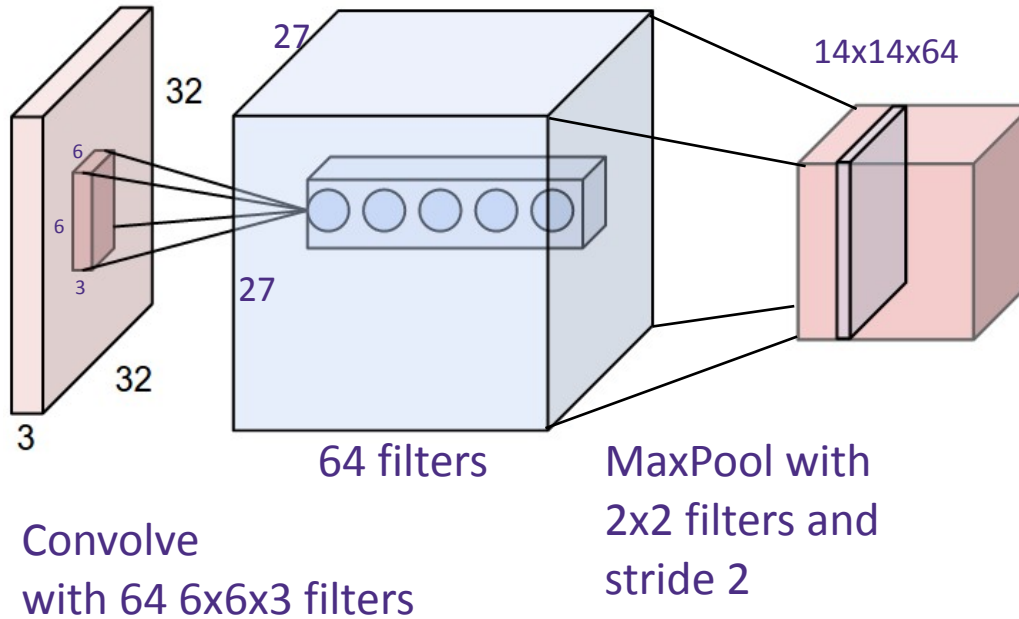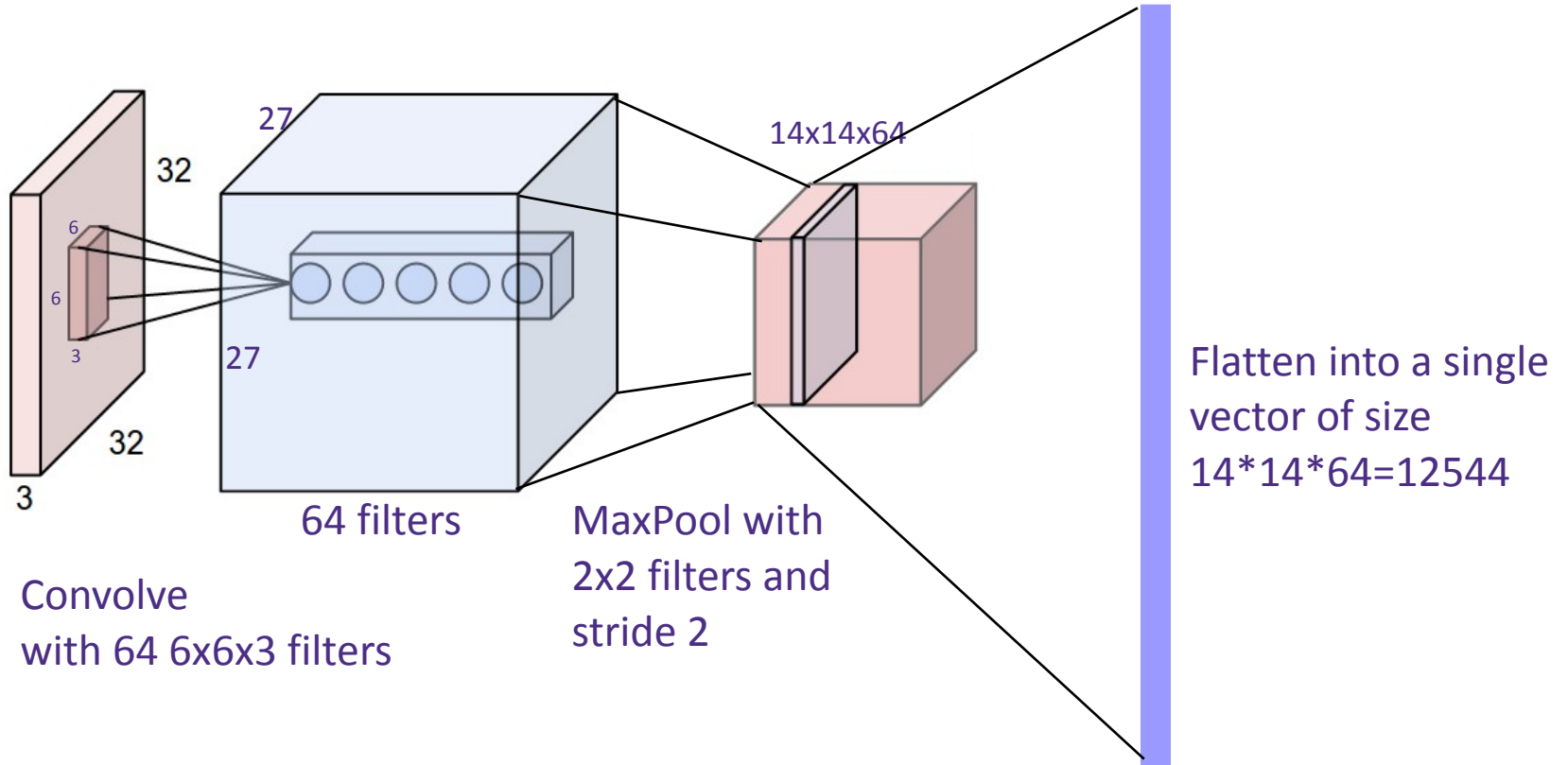
**d filters**

# Pooling

Pooling reduces the dimension and can be interpreted as "This filter had a high response in this general region"



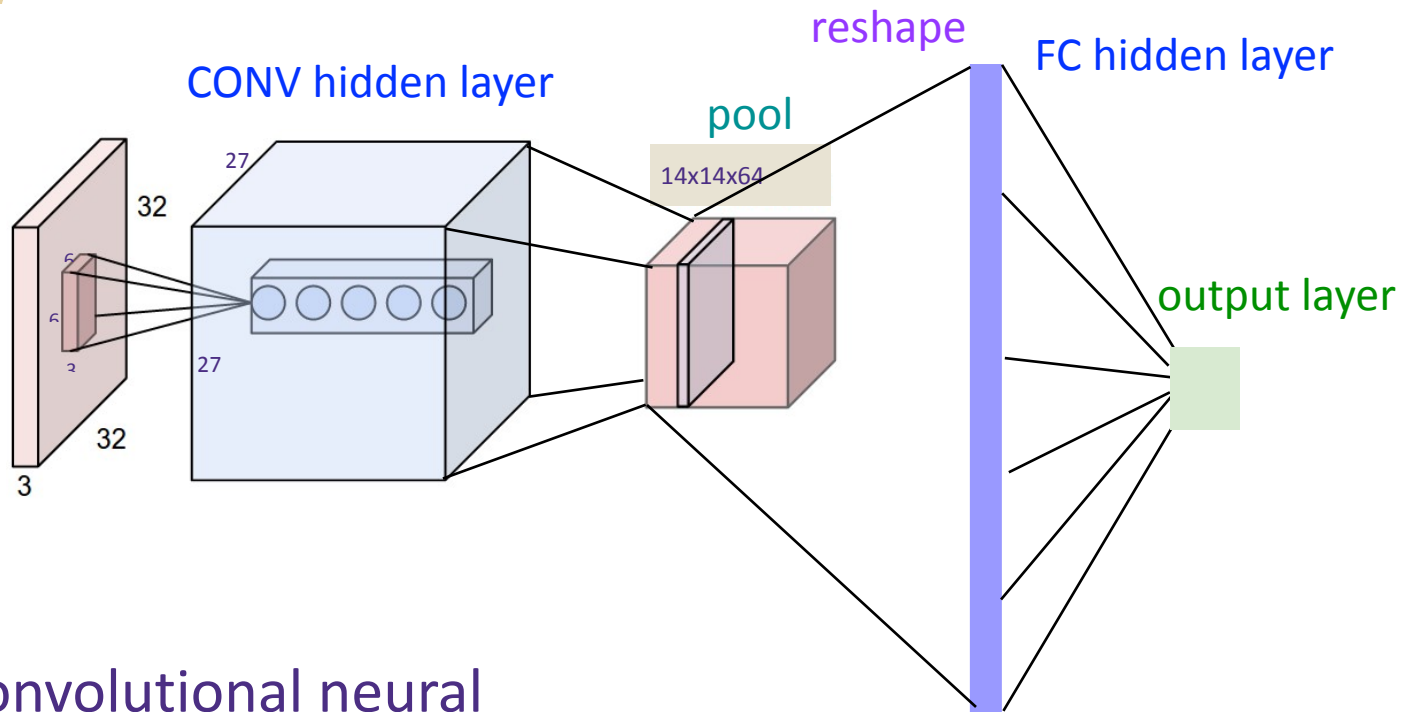Single depth slice

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

x

y

max pool with 2x2 filters and stride 2

| 6 | 8 |
|---|---|
| 3 | 4 |

27x27x64

pool

14x14x64

# Pooling Convolution layer



27

32

6

6

3

32

3

27

64 filters

27

14x14x64

Convolve
with 64 6x6x3 filters

MaxPool with
2x2 filters and
stride 2

# Flattening



27

32

6

6

3

27

32

3

64 filters

Convolve
with 64 6x6x3 filters

MaxPool with
2x2 filters and
stride 2

14x14x64

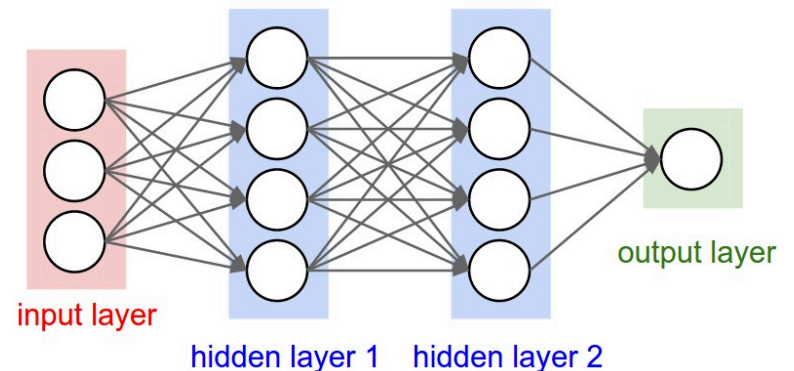Flatten into a single
vector of size
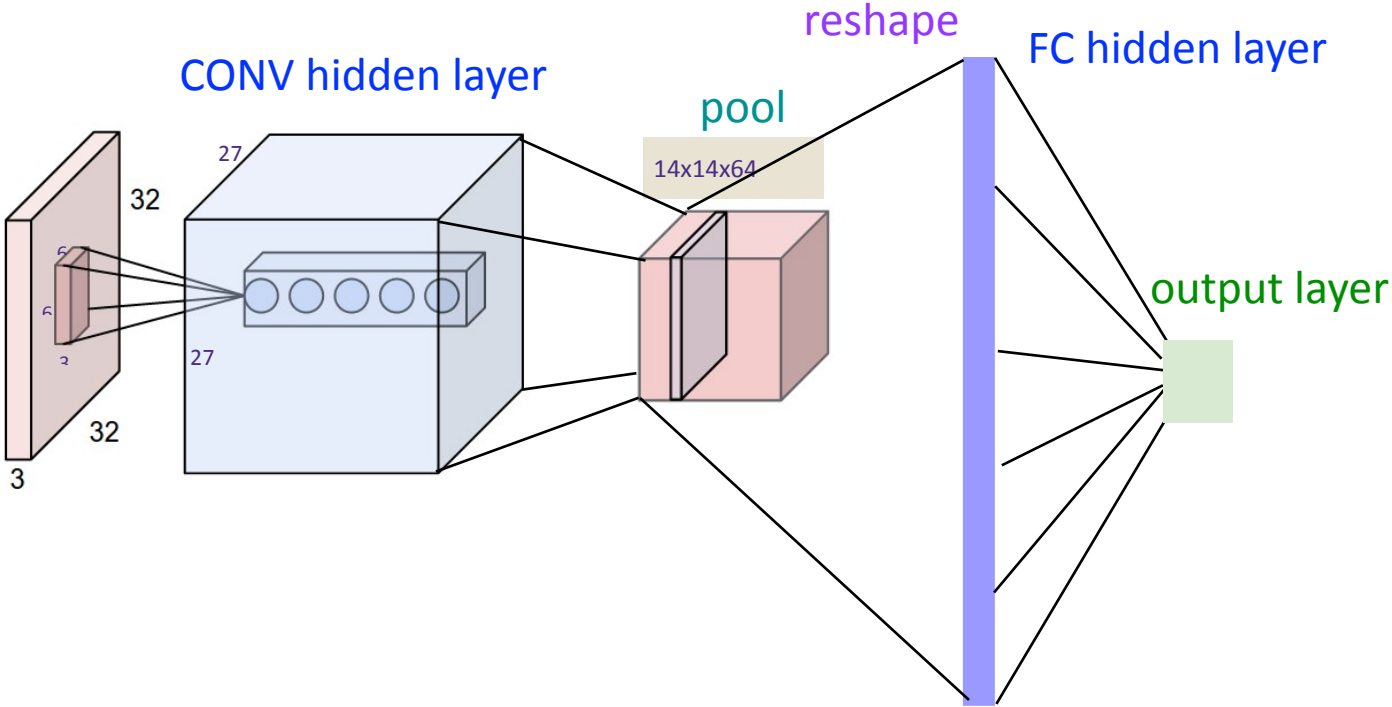14*14*64=12544

# Training Convolutional Networks



Recall: Convolutional neural networks (CNN) are just regular fully connected (FC) neural networks with some connections removed.
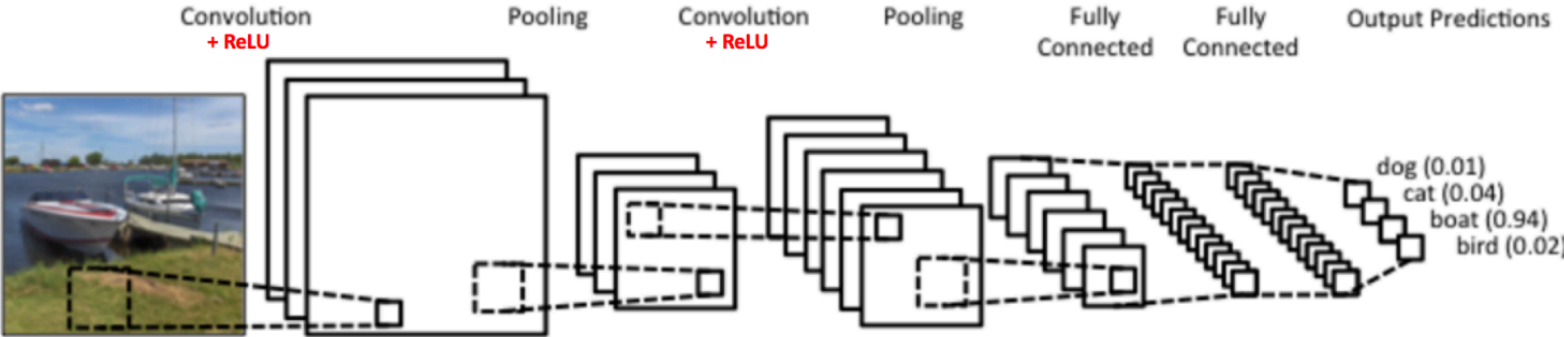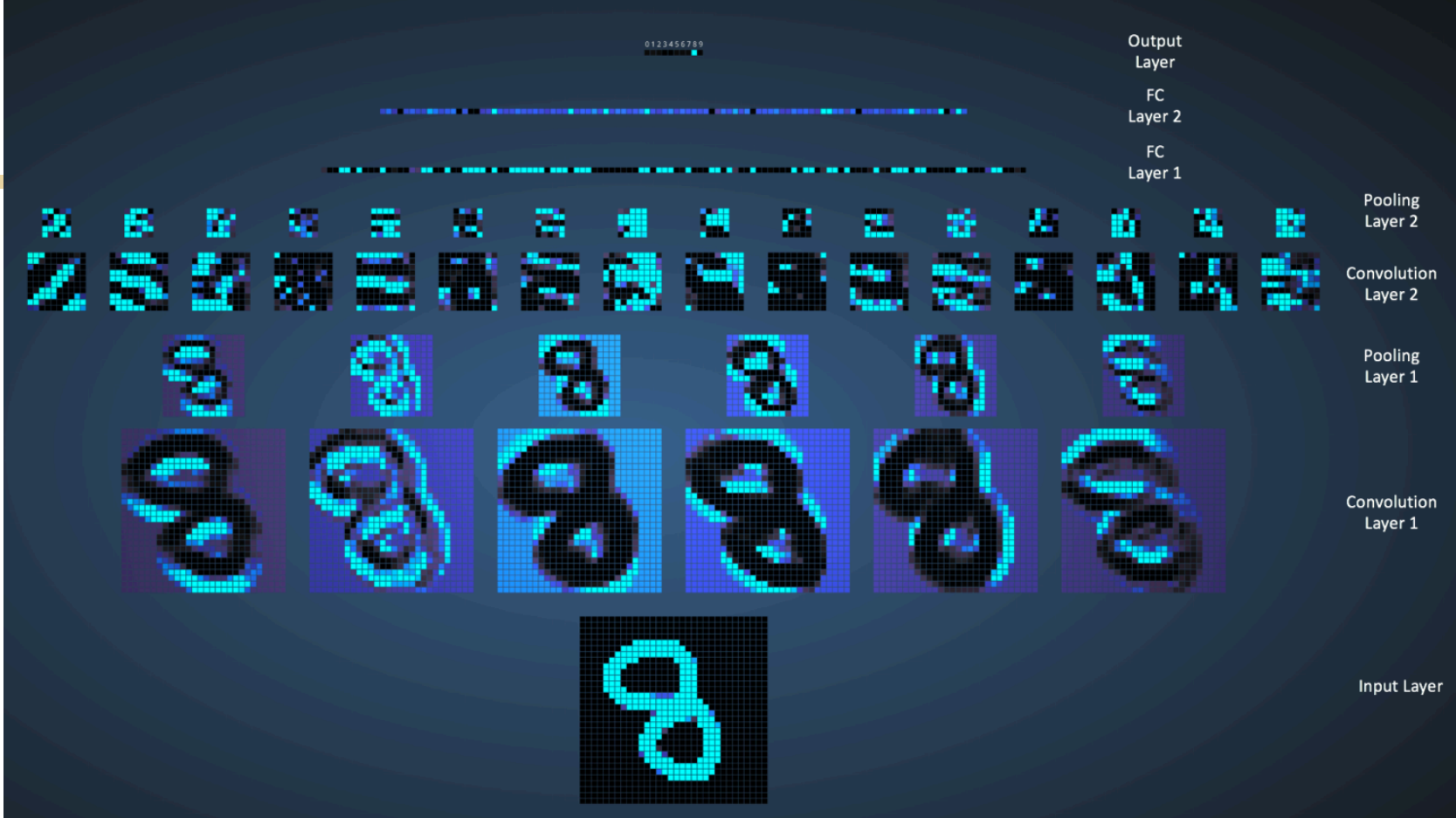**Train with SGD!**

# Training Convolutional Networks



CONV hidden layer

reshape

FC hidden layer

pool

14x14x64

output layer

32

32

3

27

27

Real example network: LeNet



Convolution + ReLU    Pooling    Convolution + ReLU    Pooling    Fully Connected    Fully Connected    Output Predictions

dog (0.01)
cat (0.04)
boat (0.94)
bird (0.02)

Real example network: LeNet