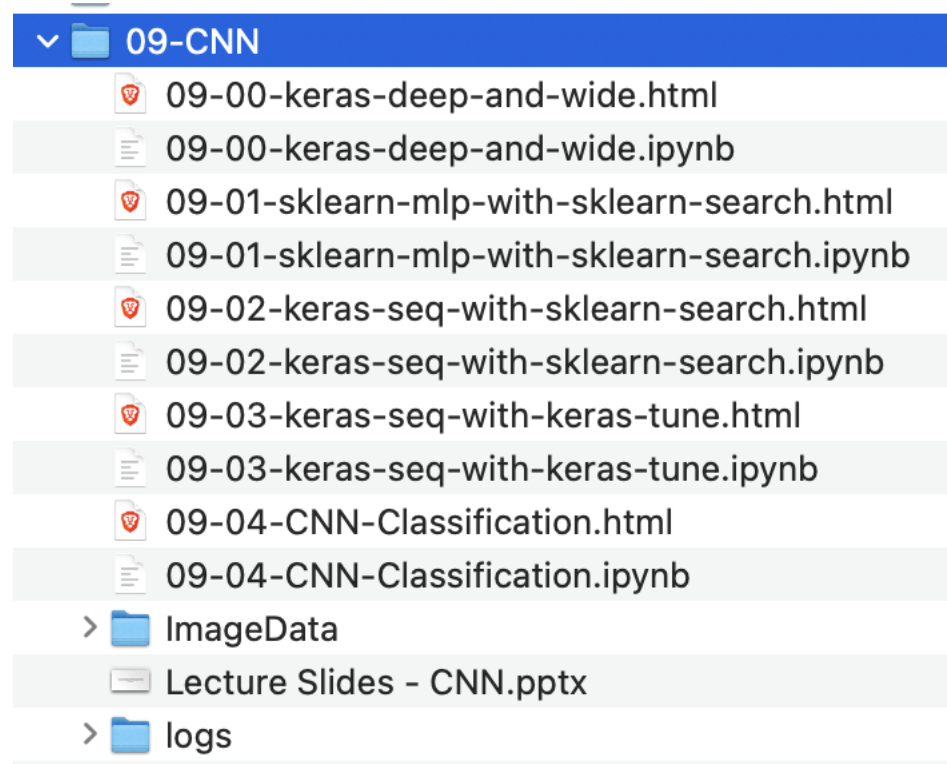# Convolutional Neural Nets

Tim Smith, PhD

# Class09

- Review DNN
  - Demonstrate various deep and wide modeling techniques
- Discuss and elaborate on hyperparameter tuning
  - Sklearn model with Sklearn tuning
  - Keras model with Sklearn tuning
  - Keras model with Keras tuning
- Introduce the concept of convolutional neural networks
  - Build an image recognition model using CNN
- Demonstrate debugging models using Tensorflow

# For the first part of the lecture:



Download the files from the following link...

https://usf.box.com/s/7bn9cox3sqyyw32lqqps2cb4nmcq0buv

# Convolutional Neural Networks

- "Discrete convolution": *a mathematical operation on two functions (f and g) to produce a third function that expresses how the shape of one is modified by the other (Wikipedia)*

- Input data: **images (usually)**

- Used in
  - Image search
  - Image classification
  - Self-driving cars

# Review on Images

- If black and white, there is 1 channel
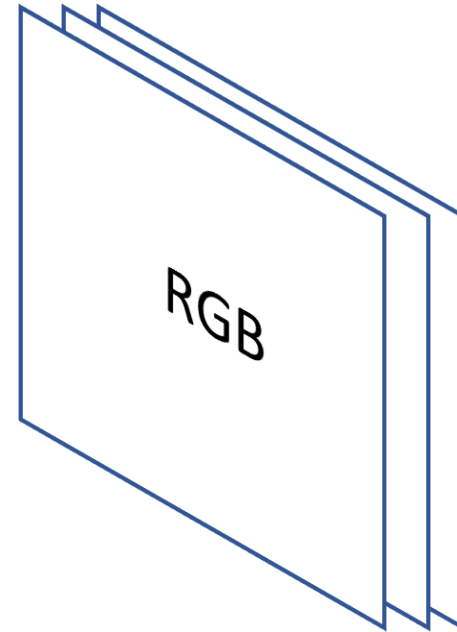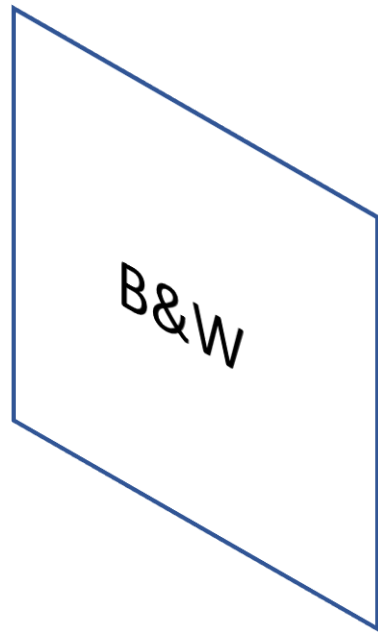- If color, there are 3 channels: Red, Green, Blue (RGB)

B&W

RGB

# Image Data

- Image size is measured in pixels.
- Ex: 480 x 480 means:

480

480

B&W

RGB

Size = 480 x 480 x 1

Size = 480 x 480 x 3

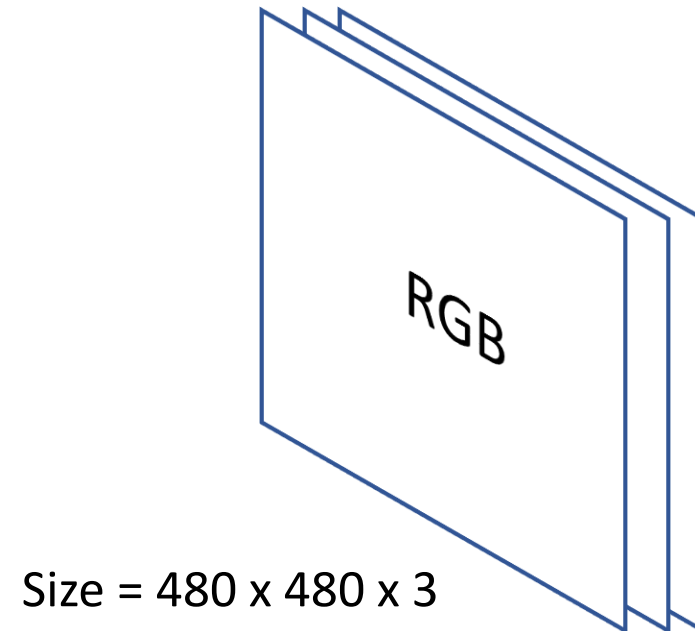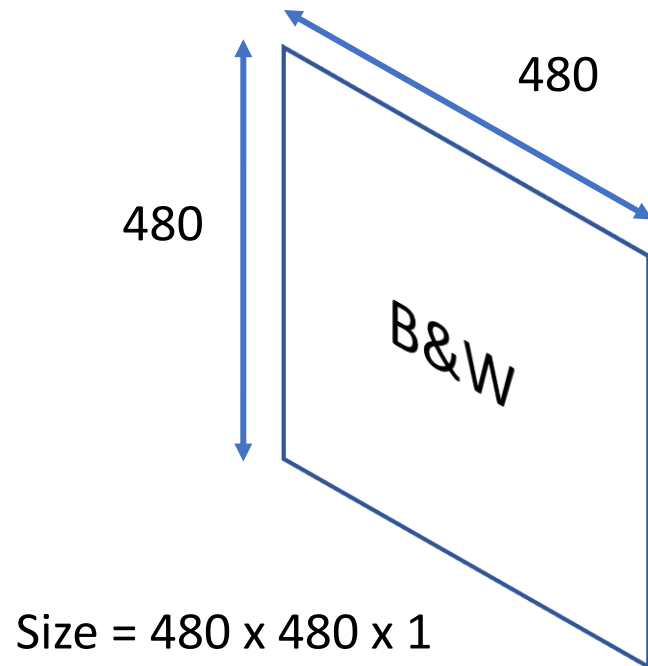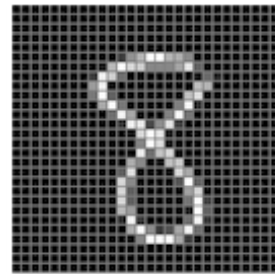# Image Data

- Each pixel has a value between 0-255



28 x 28
784 pixels

# Image Data

- Let's think of a 1M (megapixel) image
- 1 megapixel = 1,000,000 pixels (1,000 x 1,000)



1,000

1,000

Resize

100

100

10,000 pixels

10K image

# Image Data



**Problem???**

10 million connections in the first layer alone!

100

100
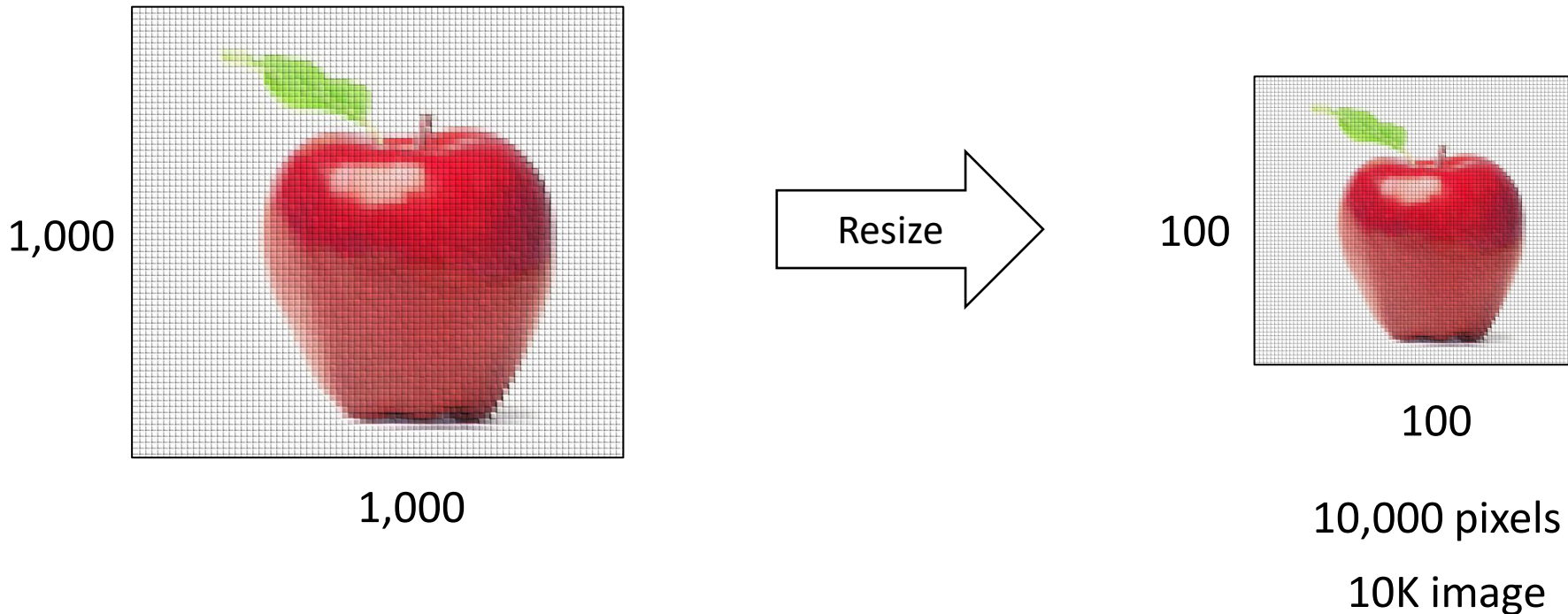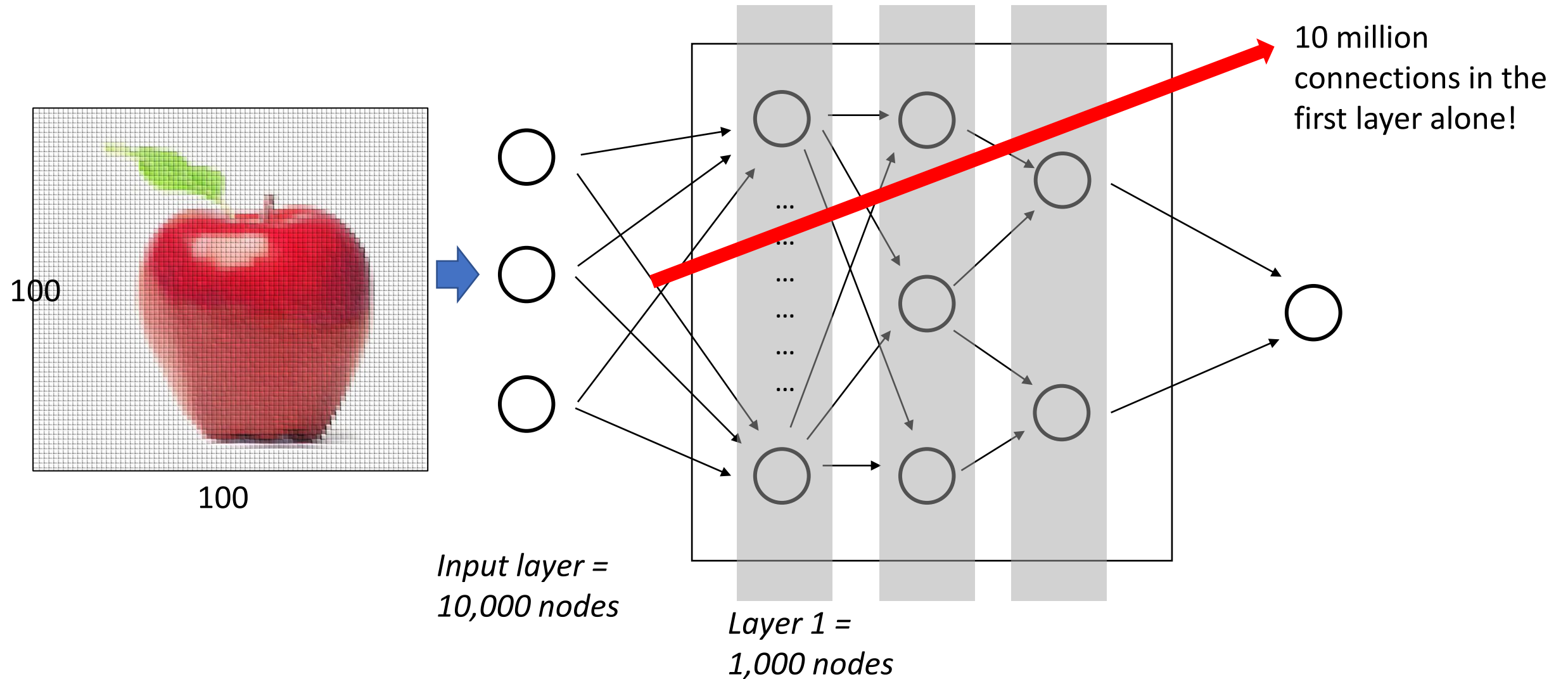
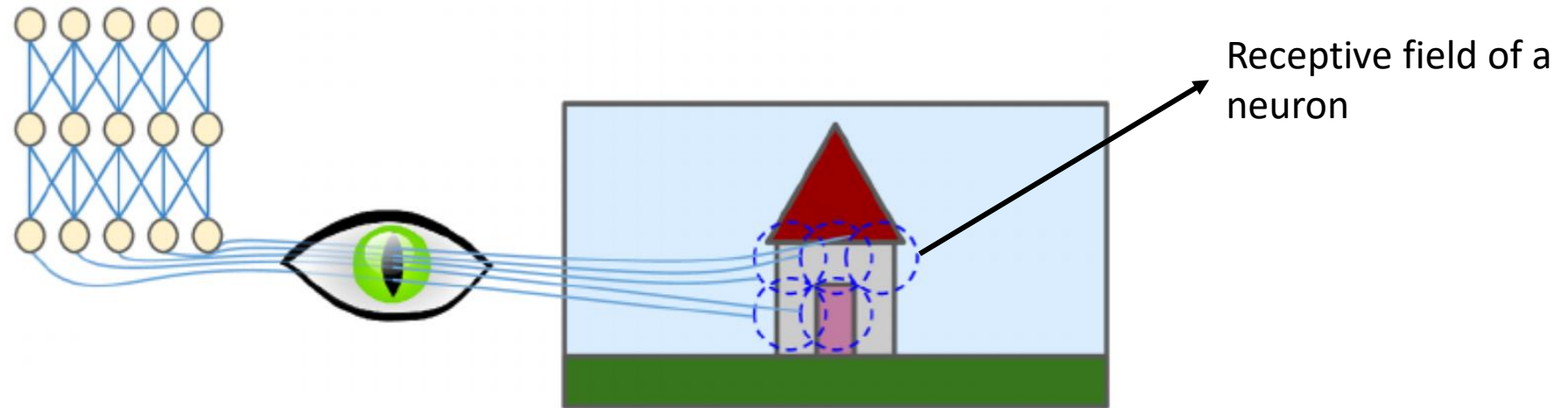*Input layer = 10,000 nodes*

*Layer 1 = 1,000 nodes*

# Solution:

- **STEP 1 (convolution):** Extract features from an image:
  - Horizontal lines,
  - Vertical lines,
  - Curvature,
  - Colors,
  - Etc.
- **STEP 2 (pooling):** Reduce the variables (without losing information)

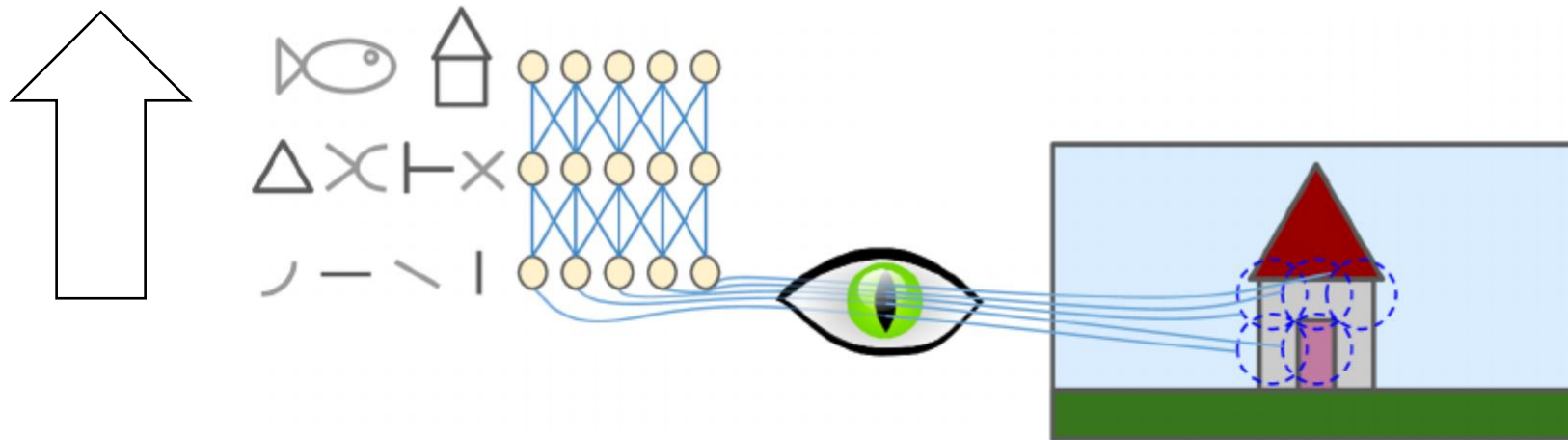- **GET HELP FROM HUMAN VISUAL CORTEX**

# Convolutional Neural Networks

- The visual cortex is made up of neurons

- Each neuron has a specific receptive field (and do not necessarily see the things outside of that field)

Receptive field of a neuron

# Convolutional Neural Networks

- Further, some neurons react only to horizontal lines
- Some react to different orientations
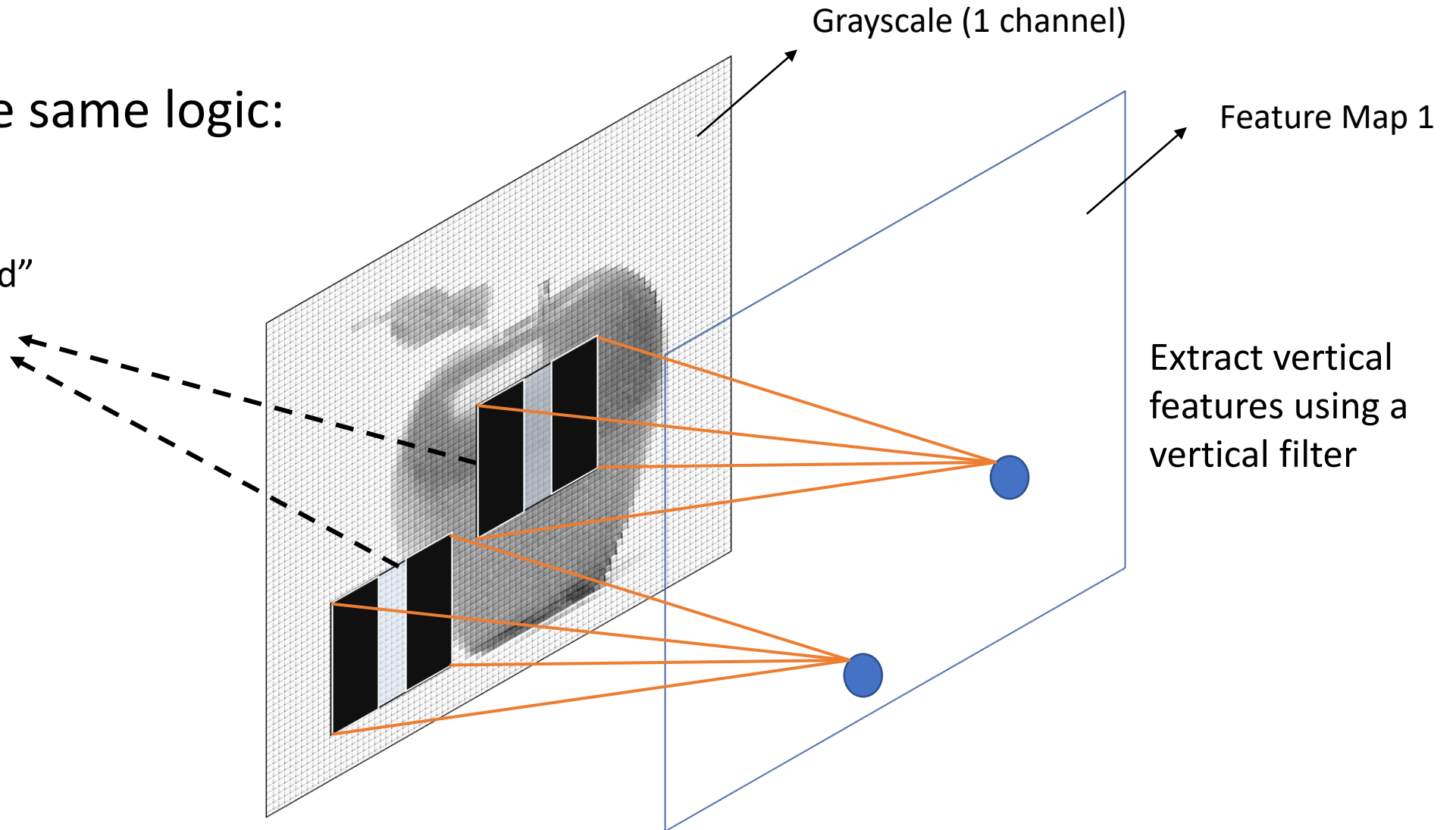- Higher level neurons make sense of the pattern detected at lower levels

# Convolutional Layer

- Apply the same logic:

Grayscale (1 channel)

Feature Map 1

"Receptive field"
Also called:
- "**Filter**"
- "Kernel"

Extract vertical features using a vertical filter

# Convolutional Layer

- Apply the same logic:

Grayscale (1 channel)

Feature Map 2

"Receptive field"
Also called:
- "**Filter**"
- "Kernel"

Extract horizontal features using a horizontal filter

# Convolutional Layer



Feature Map 1

Feature Map 2

Feature Map n

Conv Layer 1

# Convolutional Layer



Stride = 1

$f_h = 3$

$f_w = 3$

Zero padding

Makes the size of second layer the same as first

# Convolutional Layer

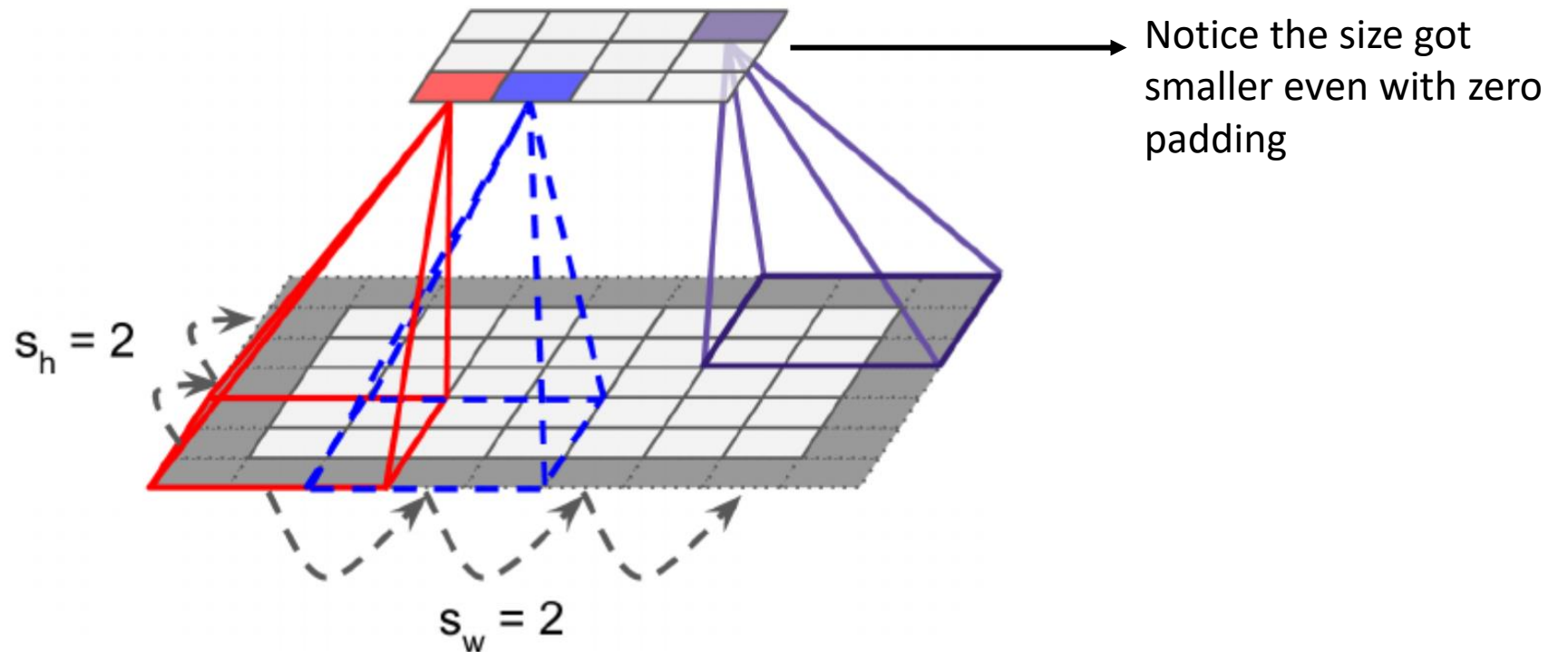- You can have varying strides (both horizontally, and vertically)



Notice the size got smaller even with zero padding

$s_h = 2$

$s_w = 2$

# Convolutional Layer

- If you have RGB (3-channel) image, the filter has 3 channels too



32x32x3 image
5x5x3 filter

# Convolutional Layer

- Eventually, you have multiple feature maps (per convolutional layer)

# Convolutional Layer

- How to process the values inside the "filter"
- Multiply them with "weights" then add them up
- Example: Apply "vertical" filter (to highlight vertical lines)

5 x 5 image

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

3 x 3 filter

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

(Actually, this sum is passed through an activation function (ReLU) to introduce non-linearity

Output is called "feature map"

# Convolutional Layer

- Then, apply "horizontal" filter (to highlight horizontal lines)

5 x 5 image

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

3 x 3 filter

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

- There are many more filters…

# Convolutional Layer
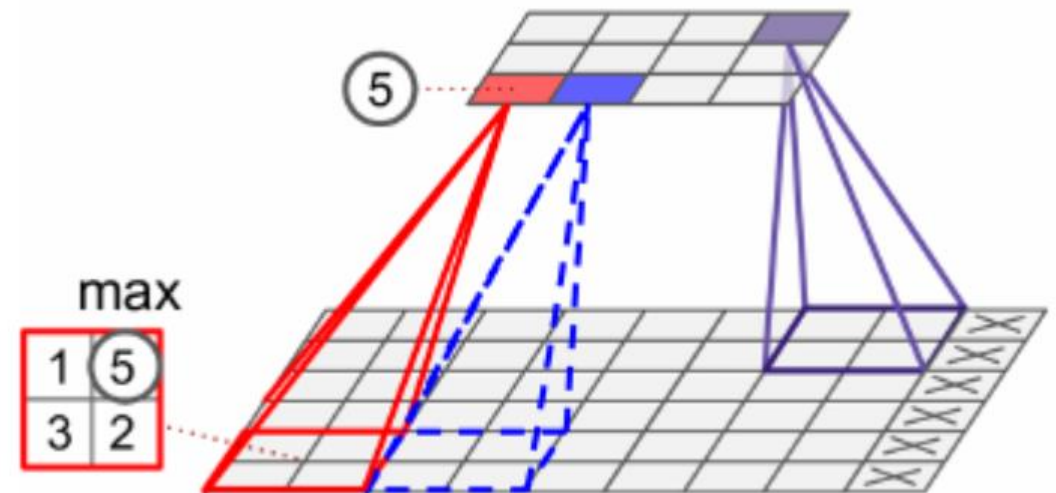
# Memory Requirements

- Imagine a single image with 150 x 100 pixels
  - Filter of 5 x 5 with 1 stride
  - 200 feature maps
  - Padding (add when necessary)
- Calculations:
  - ( 5 x 5 x 3 + 1 ) x 200 =  15,200 terms
  - Each feature map has 150 x 100 neurons that need to compute 5 x 5 x 3 inputs. Makes 225 million multiplications
- 1 image needs 12 MB RAM.
- 100 images 1.2 GB RAM.

# Pooling Layer

- "Shrinks" the convolutional layer (by sampling)
- An "aggregation" technique!
- Why:
  - Reduce computational load
  - Reduce number of parameters
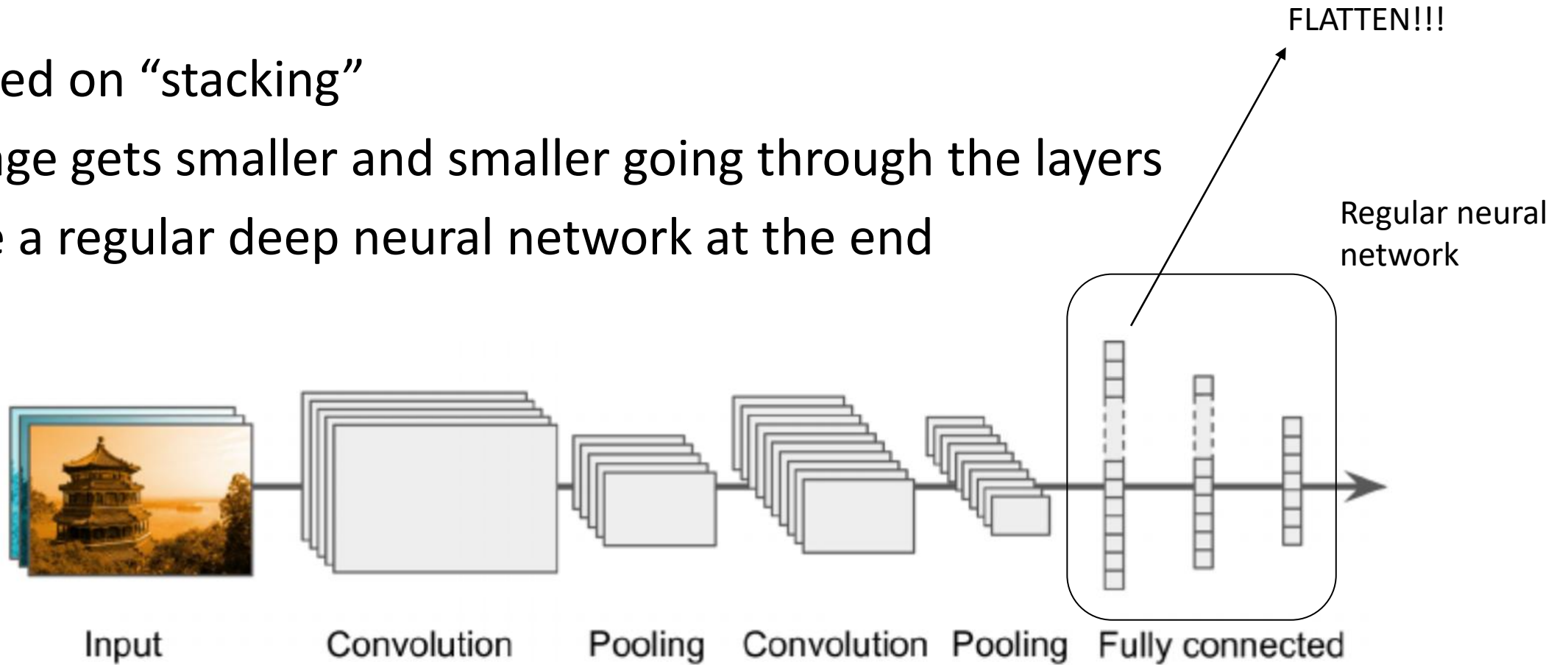  - Reduce overfitting

# Pooling Layer

- Similar to earlier: define the size, stride, padding type
- It AGGREGATES the values
  - Example: use min, max, average, etc.
- A 2 x 2 kernel with a stride of 2 reduces the image (inputs) by 75%
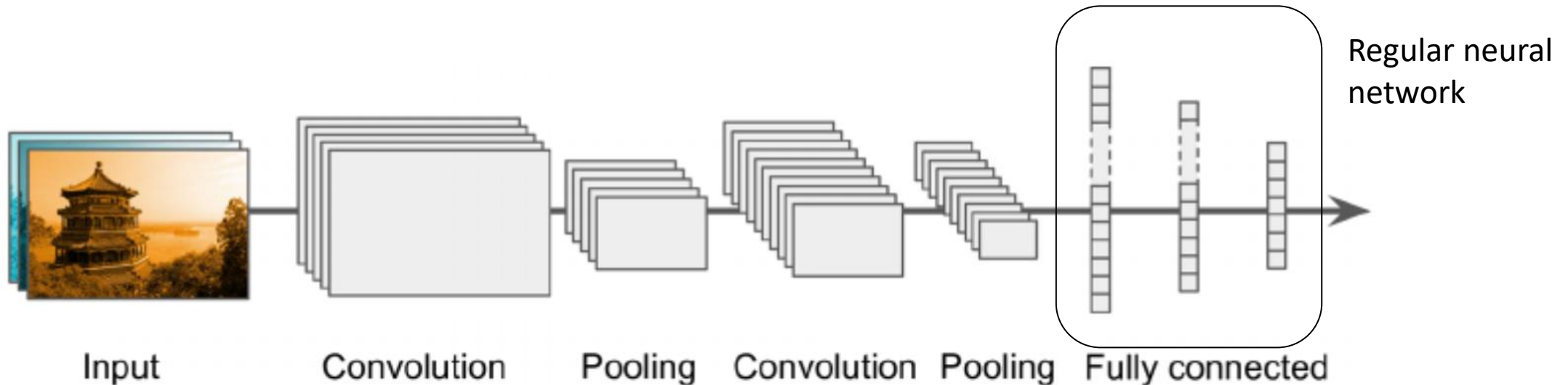- You can set it to shrink channels too

# CNN Architectures

- Based on "stacking"
- Image gets smaller and smaller going through the layers
- Use a regular deep neural network at the end

FLATTEN!!!

Regular neural network



Input     Convolution     Pooling     Convolution   Pooling   Fully connected

# CNN Architectures

- Can also stack two convolutional layers one after another

- Instead of having a 5 x 5 convo layer, have two
  3 x 3 stacked convo layers. Might perform better!

- Double the filter size as you go deeper in CNN



Regular neural network

Input    Convolution    Pooling    Convolution    Pooling    Fully connected

# LeNET-5

- Developed in 1998
- Used the MNIST data set (handwritten digits)



- **Labeled data**
- Black and white
- 28x28
- 60,000 images for training
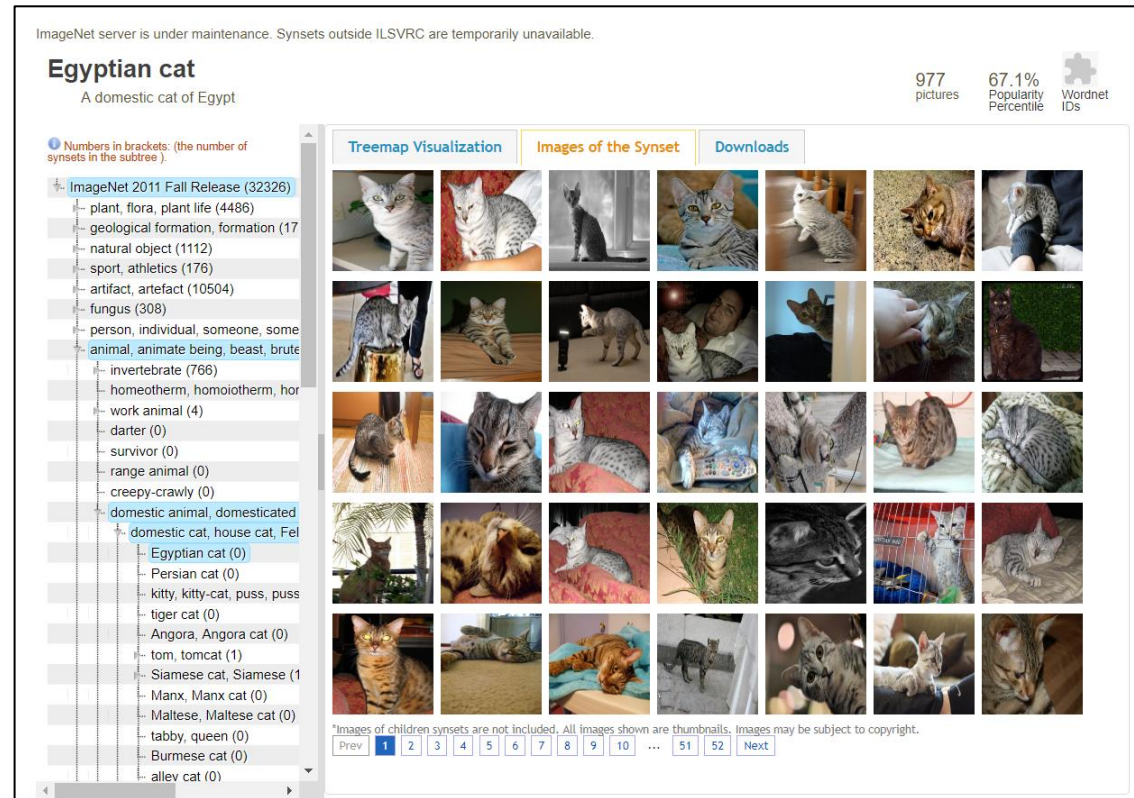- 10,000 images for testing

# LeNET-5

- Architecture:

| Layer | Type | Maps | Size | Kernel size | Stride | Activation |
|-------|------|------|------|-------------|--------|------------|
| Out | Fully connected | – | 10 | – | – | RBF |
| F6 | Fully connected | – | 84 | – | – | tanh |
| C5 | Convolution | 120 | $1 \times 1$ | $5 \times 5$ | 1 | tanh |
| S4 | Avg pooling | 16 | $5 \times 5$ | $2 \times 2$ | 2 | tanh |
| C3 | Convolution | 16 | $10 \times 10$ | $5 \times 5$ | 1 | tanh |
| S2 | Avg pooling | 6 | $14 \times 14$ | $2 \times 2$ | 2 | tanh |
| C1 | Convolution | 6 | $28 \times 28$ | $5 \times 5$ | 1 | tanh |
| In | Input | 1 | $32 \times 32$ | – | – | – |

Error: less than 1%

# AlexNET

- Developed in 2012
- Achieved 17% error rate in ILSVRC imagenet challenge
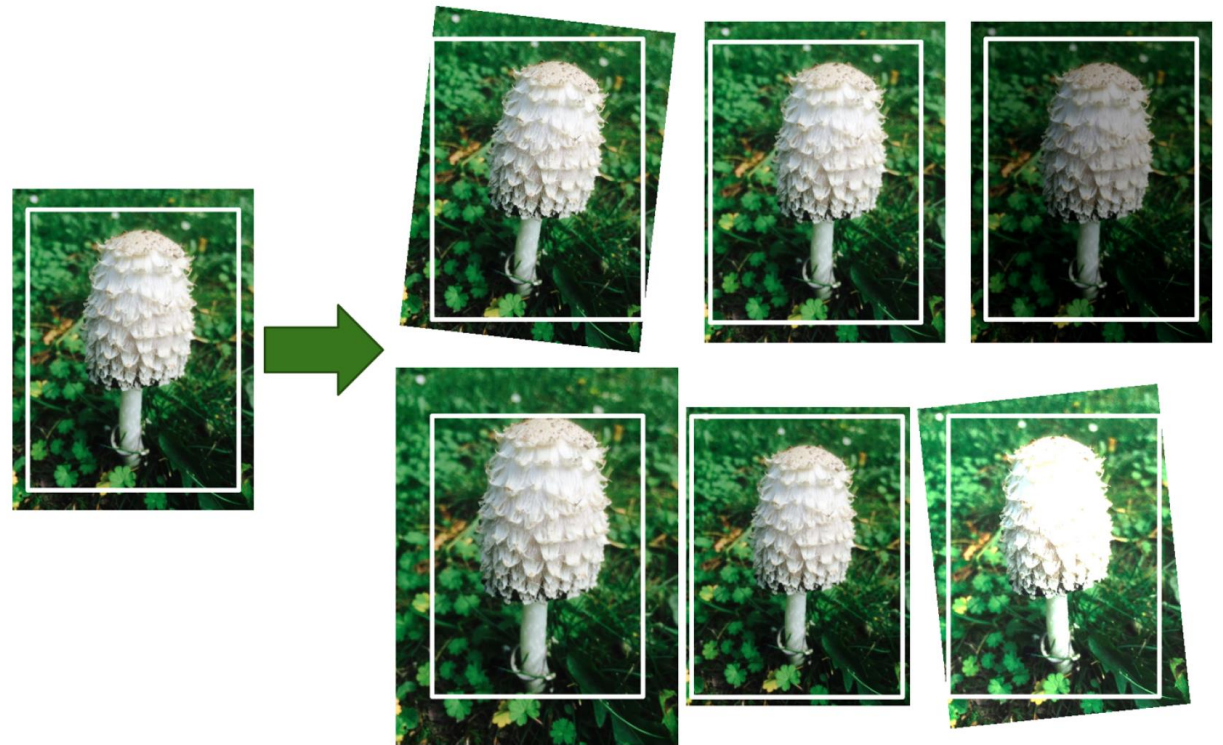  - Labeled data (1,000 classes)
  - 150,000 images in total

# AlexNET

- Architecture:

| Layer | Type | Maps | Size | Kernel size | Stride | Padding | Activation |
|-------|------|------|------|-------------|--------|---------|------------|
| Out | Fully connected | – | 1,000 | – | – | – | Softmax |
| F10 | Fully connected | – | 4,096 | – | – | – | ReLU |
| F9 | Fully connected | – | 4,096 | – | – | – | ReLU |
| S8 | Max pooling | 256 | $6 \times 6$ | $3 \times 3$ | 2 | valid | – |
| C7 | Convolution | 256 | $13 \times 13$ | $3 \times 3$ | 1 | same | ReLU |
| C6 | Convolution | 384 | $13 \times 13$ | $3 \times 3$ | 1 | same | ReLU |
| C5 | Convolution | 384 | $13 \times 13$ | $3 \times 3$ | 1 | same | ReLU |
| S4 | Max pooling | 256 | $13 \times 13$ | $3 \times 3$ | 2 | valid | – |
| C3 | Convolution | 256 | $27 \times 27$ | $5 \times 5$ | 1 | same | ReLU |
| S2 | Max pooling | 96 | $27 \times 27$ | $3 \times 3$ | 2 | valid | – |
| C1 | Convolution | 96 | $55 \times 55$ | $11 \times 11$ | 4 | valid | ReLU |
| In | Input | 3 (RGB) | $227 \times 227$ | – | – | – | – |

- Used 50% dropout at F9 and 10
- Performed data augmentation
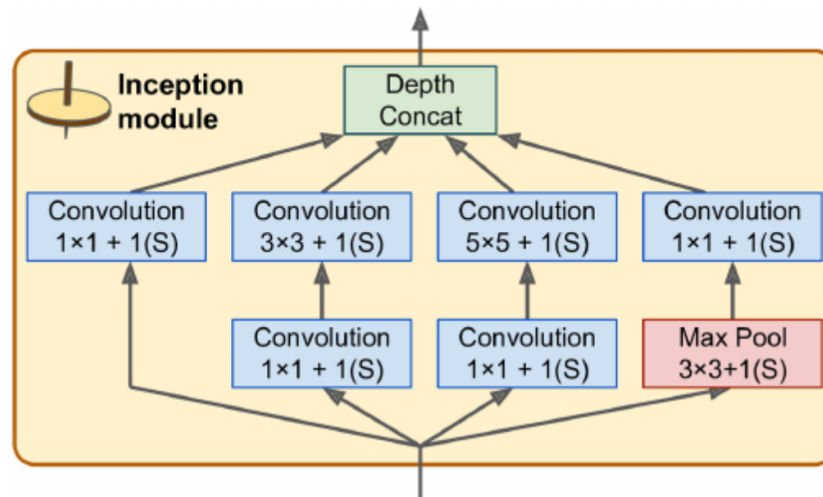- Changed lighting conditions
- Local response normalization

# Data Augmentation

- Generate new images from existing training images

- Artificially increase the size
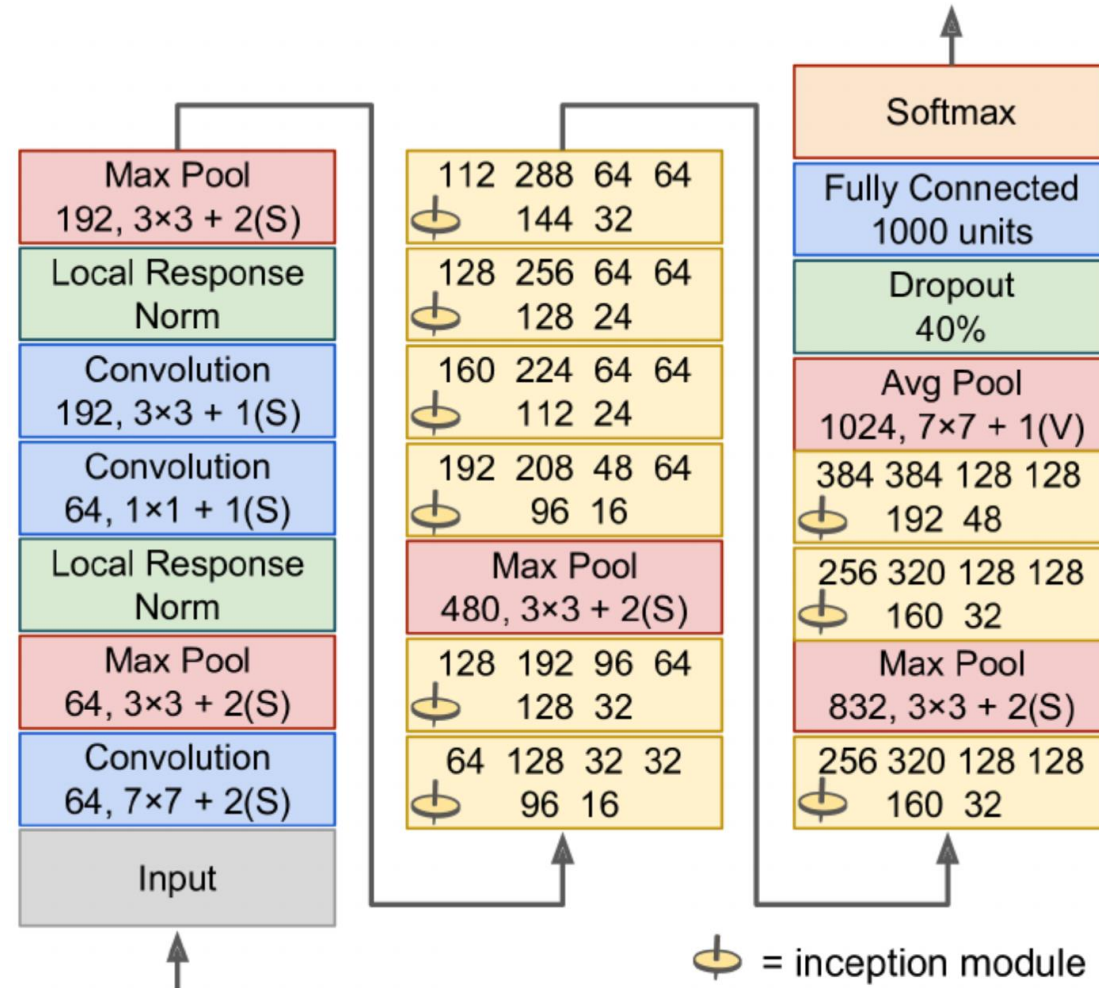
- Reduces overfitting

# GoogLeNET

- Achieved an error rate of below 7% (ILSVRC imagenet challenge)

- Deeper network

- Reduced AlexNET's parameters by 10 times (from 60M to 6M)

- Used sub-networks called inception
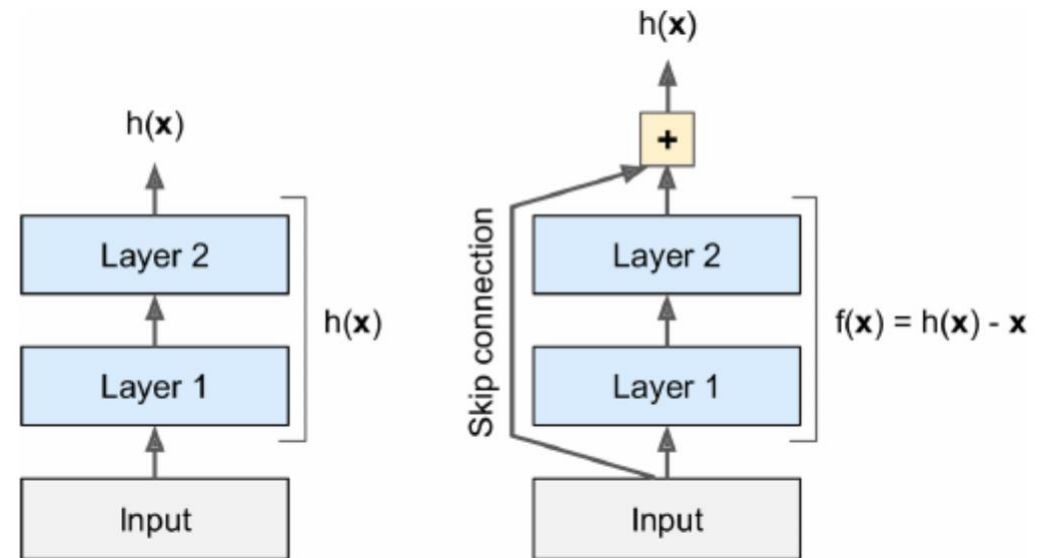


Why use 1x1 filters?
- Capture patterns along depth dimension
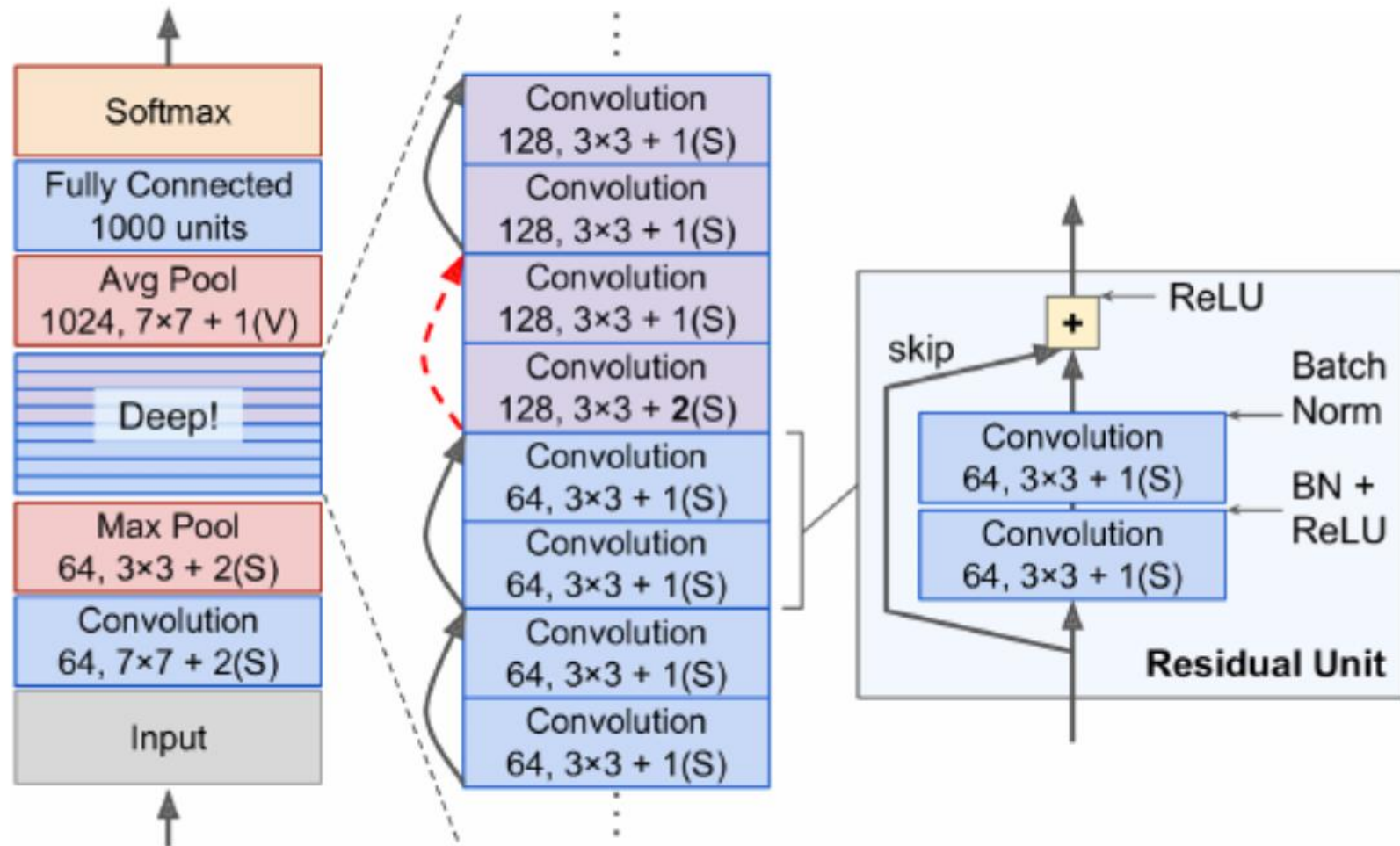- Reduce dimensionality (bottleneck layers)

# GoogLeNET

# ResNET

- Reduced the error rate below 3.6%

- A total of 152 layers

- Used skip (shortcut) connections
  - Helps perform residual learning

- Skipping connections helps train faster

# ResNET

# Other Models

- VGGNet

- Xception (a variant of GoogLeNet)

- SENet (reduced the error rate to 2.25% in 2017)

- Inception v4
  - Combination of GoogLeNet and Resnet

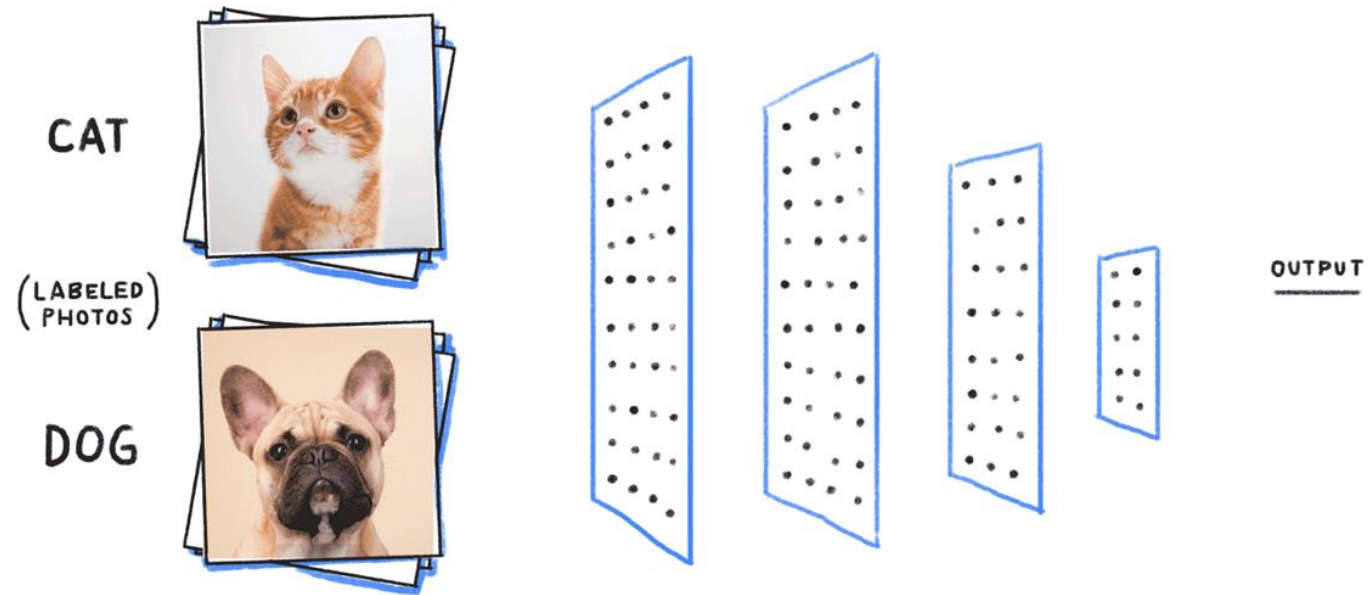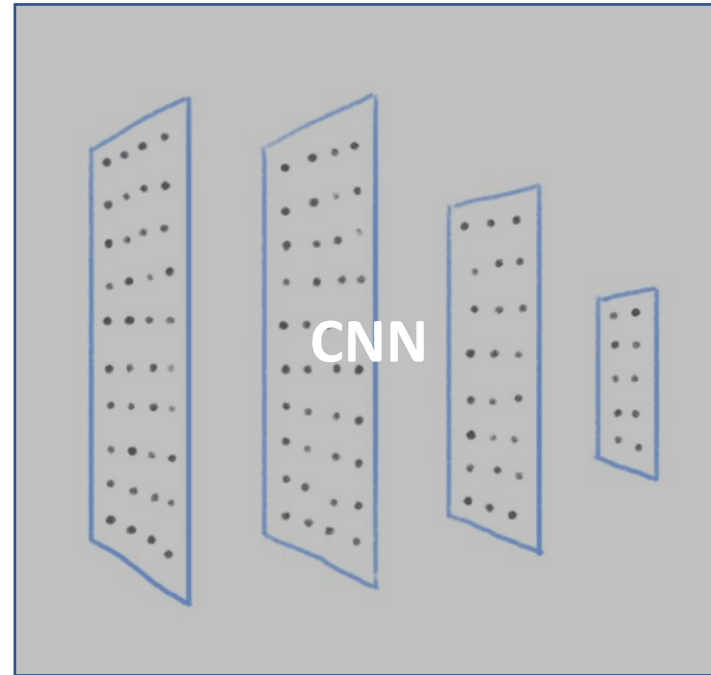# Data Requirements

- We need **labeled** images (most of the time)

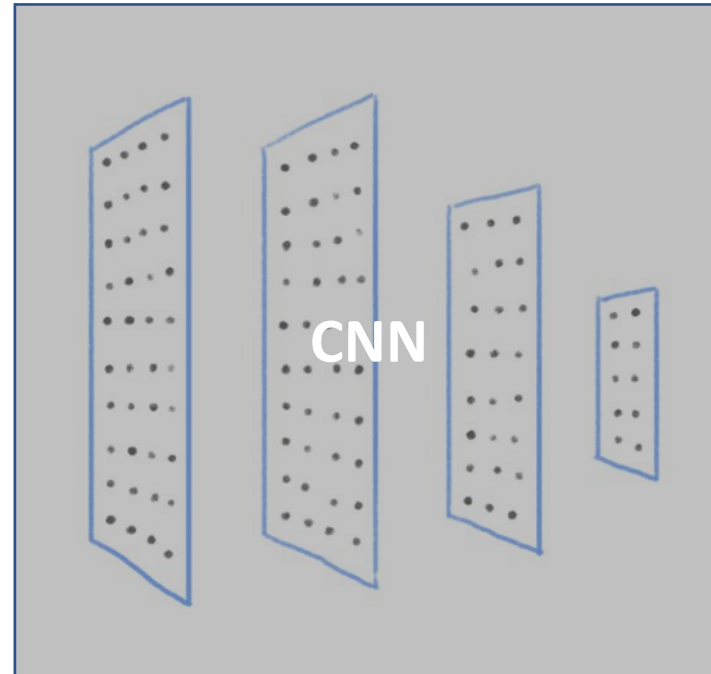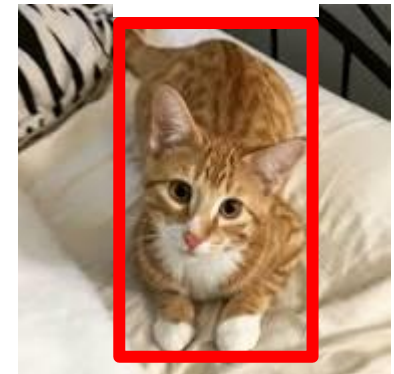# Image classification

New (unseen) image



CNN

OUTPUT  **Cat**

# Object detection

New (unseen) image



CNN

OUTPUT

Cat 0.95

# Object detection

- Example

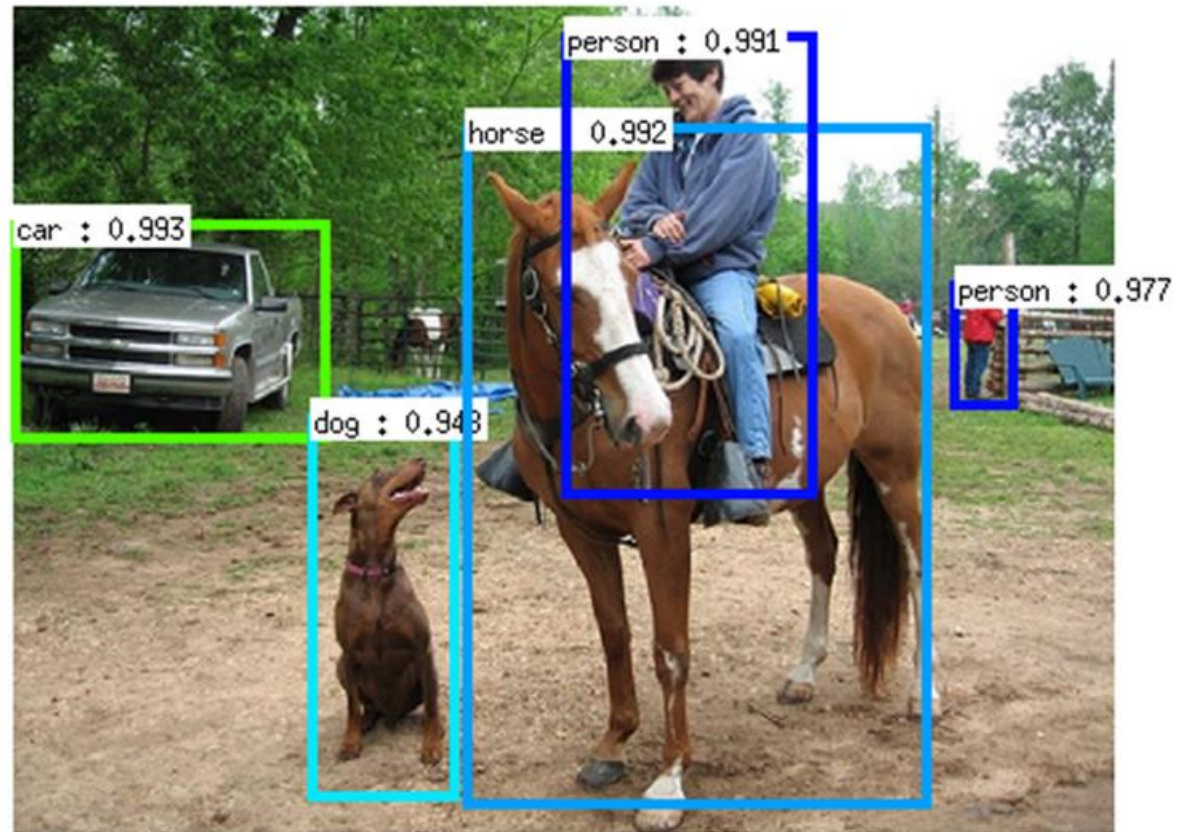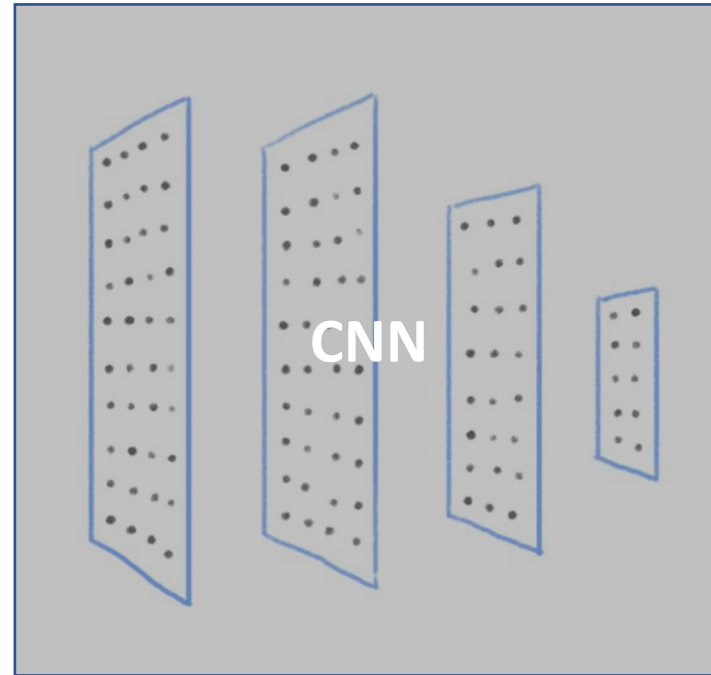A popular model:
You Only Look Once (YOLOv3)

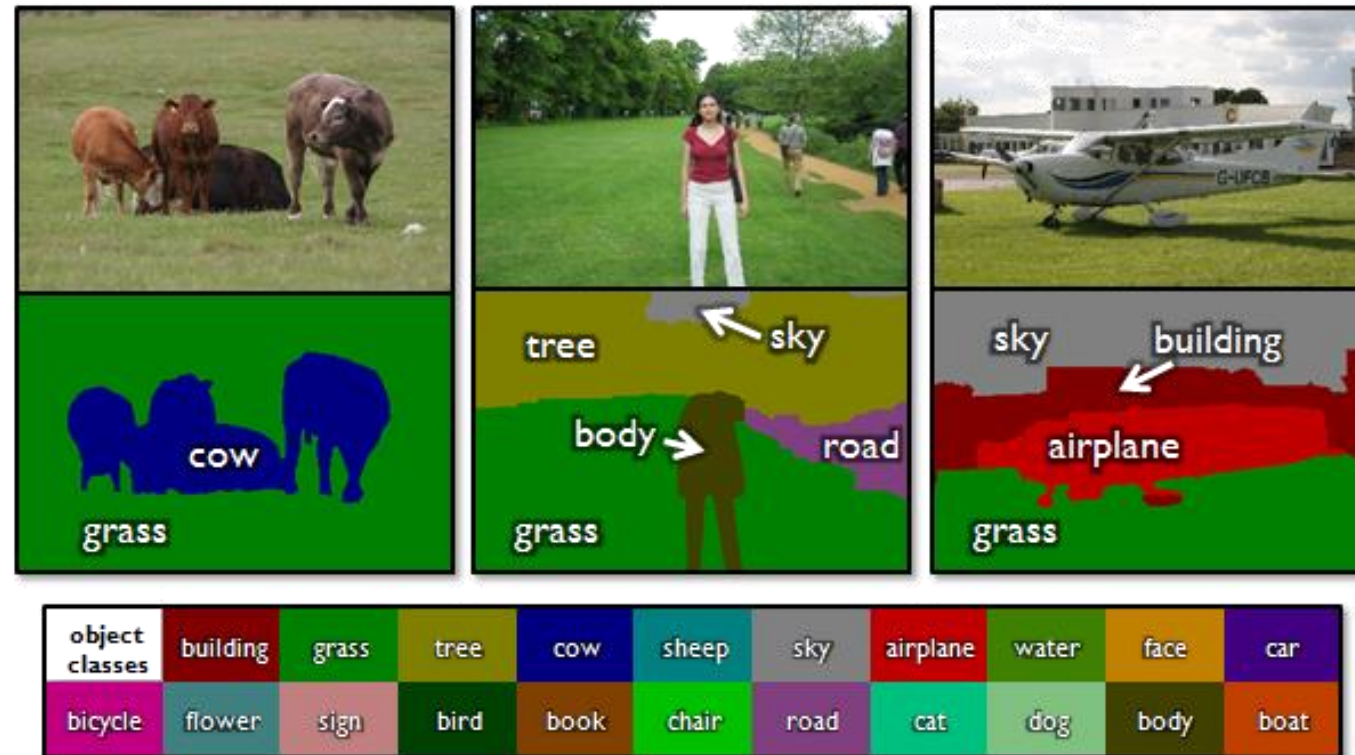# Image segmentation

New (unseen) image

CNN

OUTPUT

Cat

# What Kind of Output Does it Generate?

- Image segmentation example

# How Are We Going to Use CNNs?

- Image classification