

# I 202: INFORMATION ORGANIZATION & RETRIEVAL FALL 2025

---

Class 16: Image Classification

# Today's Outline

Image Classification

Data Creation & Augmentation

Image Convolutions

Convolutional Neural Nets

# IMAGE PROCESSING ALGORITHMS

- Today: Image recognition brief intro, CNNs
- Later: Transformer-based image models

# IMAGE CLASSIFICATION

Image classification is the task of taking an input image and outputting a class or a **probability of classes** that best describes the image



input image

classification →

“dog”



input image

classification →

“cat”

# Image Classification Application: Autonomous Driving



# IMAGE ANALYSIS CHALLENGE: VARIABLE VIEWPOINT



Michelangelo 1475-1564

# CHALLENGE: VARIABLE ILLUMINATION DIFFERENCES

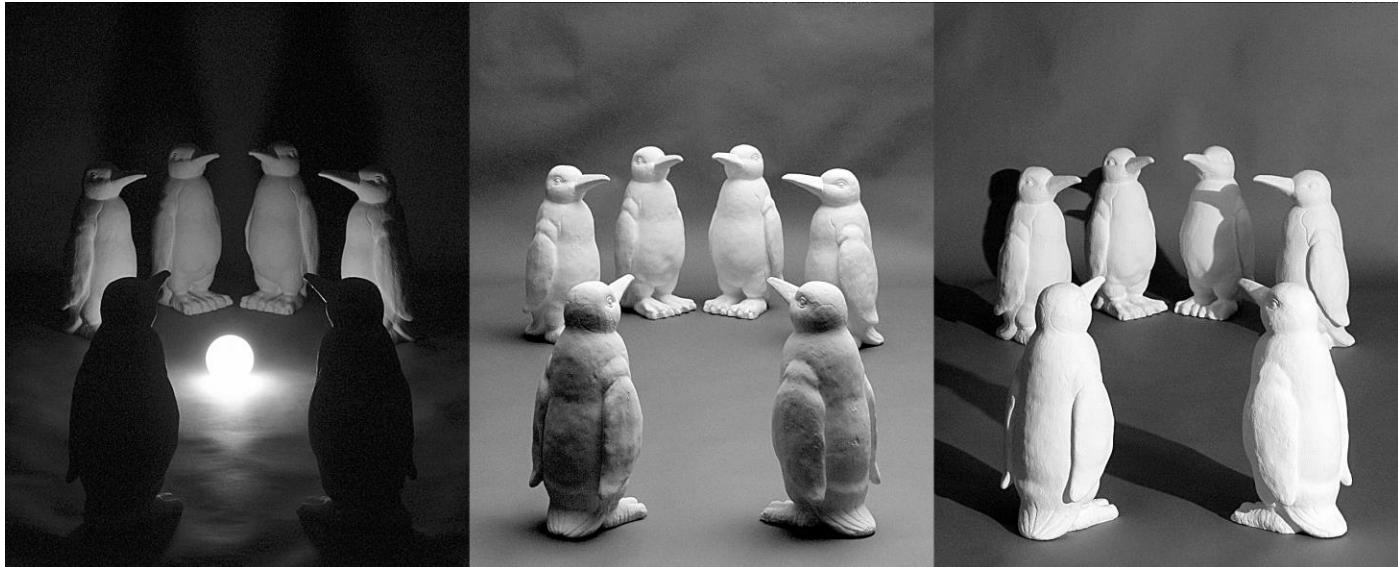


image credit: J. Koenderink

# CHALLENGE: DEFORMATION

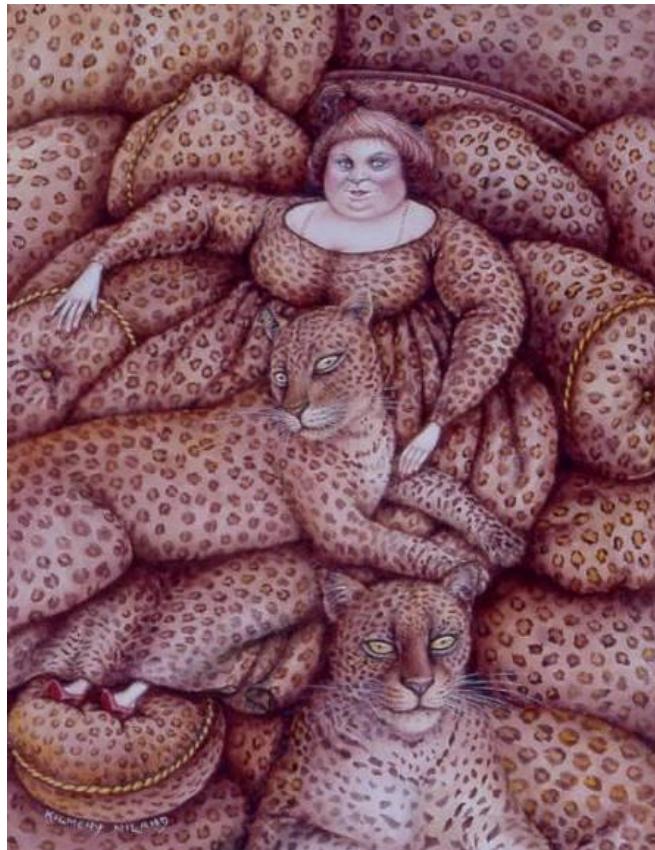


# CHALLENGE: OCCLUSION



Magritte, 1957

# CHALLENGE: BACKGROUND CLUTTER



Kilmeny Niland. 1995

# Challenge: Intra-Class Variations

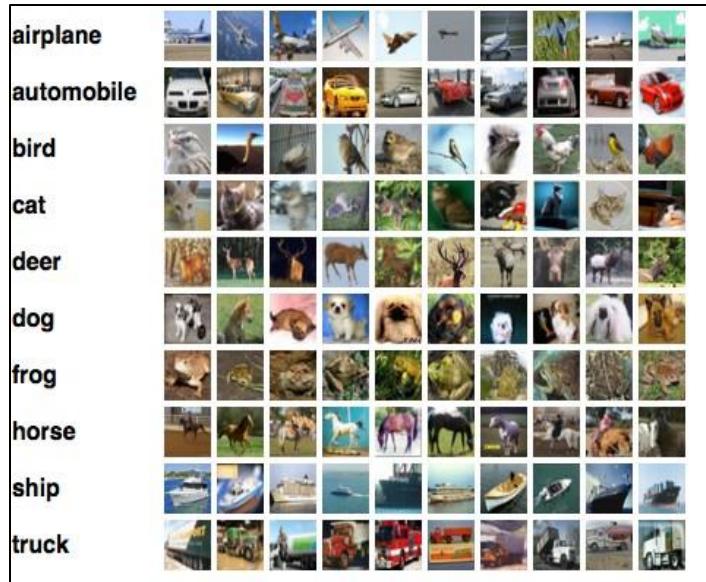


# **DATA GATHERING & AUGMENTATION**

# IMAGE CATEGORY DATASET

## CIFAR-10 Dataset (<https://www.cs.toronto.edu/~kriz/cifar.html>)

- Consists of 60,000 32x32 **color** images in 10 classes, with 6,000 images per class. There are 50,000 training images and 10,000 test images.



# IMAGENET: A DATA GATHERING BREAKTHROUGH > 14 M IMAGES AND 20,000 CATEGORIES!



# IMAGENET AND WORDNET

WordNet inspired the creation of ImageNet and provides its structure

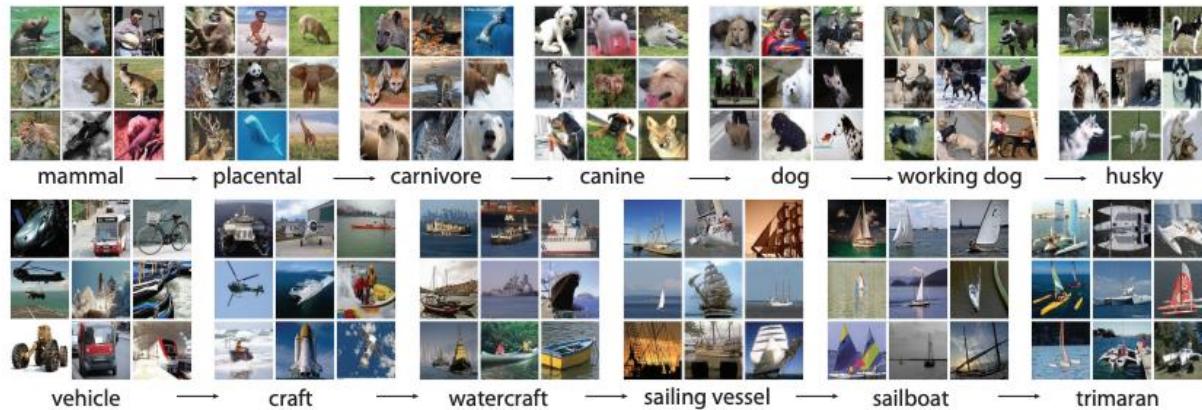
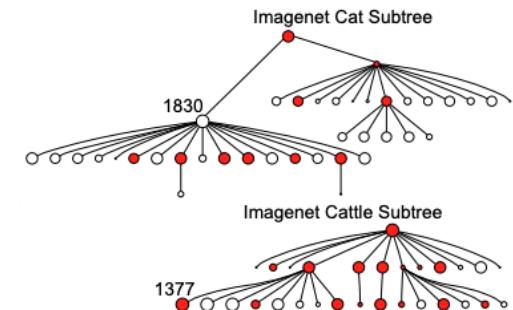


Figure 1: A snapshot of two root-to-leaf branches of ImageNet: the **top** row is from the mammal subtree; the **bottom** row is from the vehicle subtree. For each synset, 9 randomly sampled images are presented.



# IMAGENET: A BREAKTHROUGH IN DATA COLLECTION

“We went to the Internet, the biggest treasure trove of pictures that humans have ever created. We downloaded nearly a billion images and used crowdsourcing technology like the Amazon Mechanical Turk platform to help us to label these images. At its peak, ImageNet was one of the biggest employers of the Amazon Mechanical Turk workers: together, almost 50,000 workers from 167 countries around the world helped us to clean, sort and label nearly a billion candidate images. That was how much effort it took to capture even a fraction of the imagery a child's mind takes in in the early developmental years”. – Fei Fei Li

# Image Data Gathering Tasks



Single object, plain background

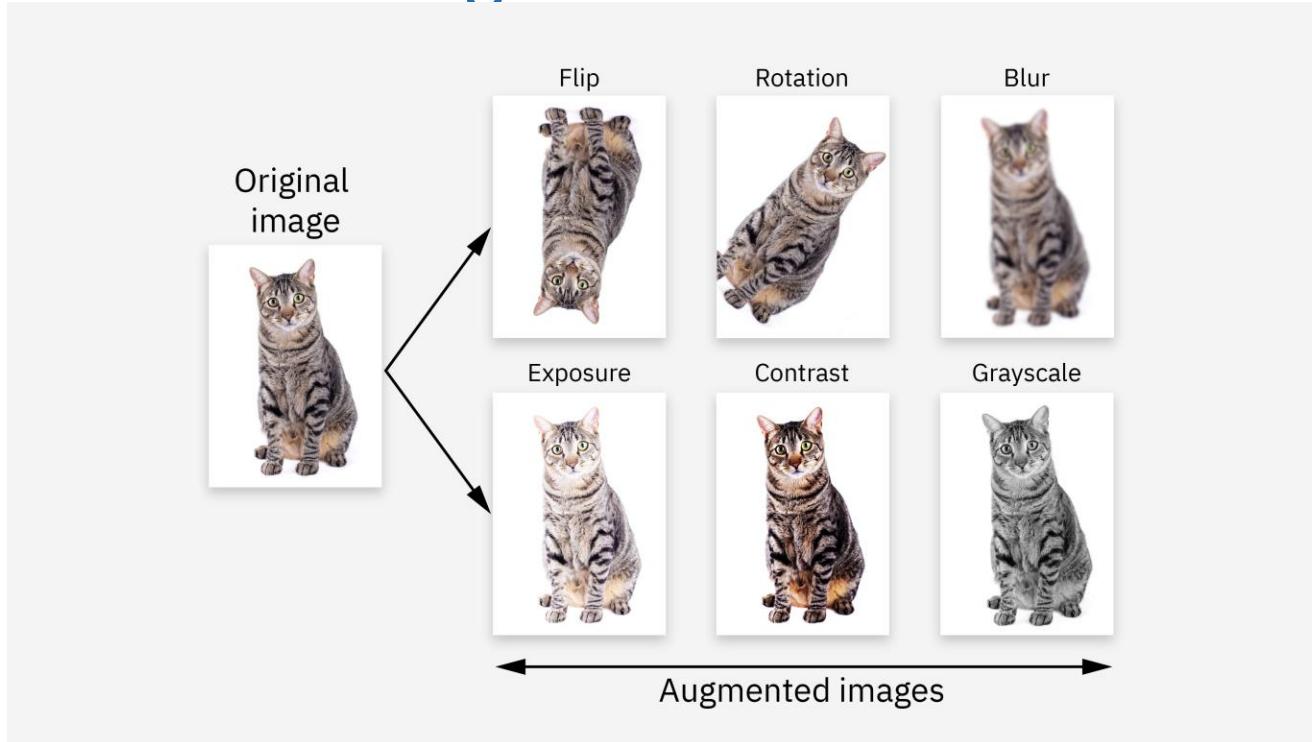


Choose a more general category when uncertain

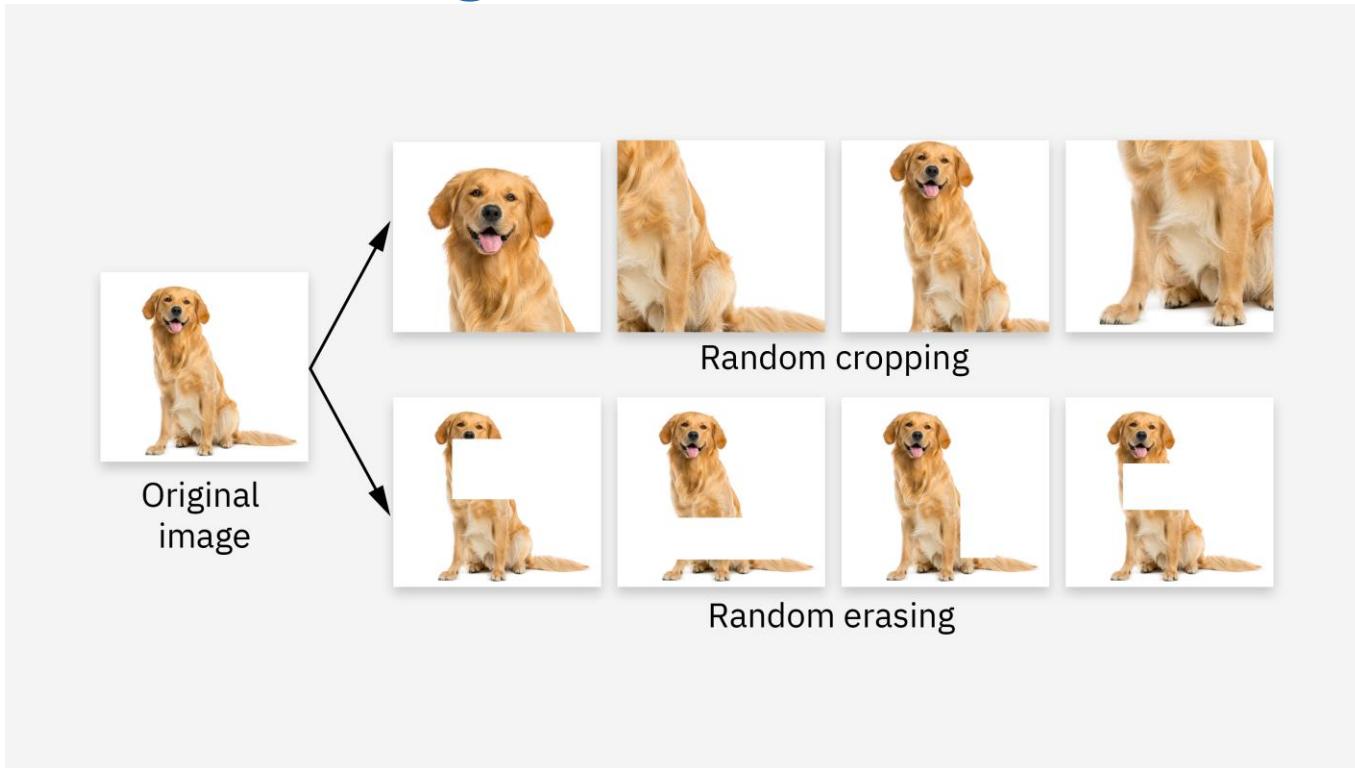


Multiple objects

# Image Data Gathering Task: Data Augmentation



# Image Data Gathering Task: Data Augmentation



## APPLICATION: DATA GATHERING VIA IMAGE RECOGNITION

Got millions of Google street view images, and classified the makes and models of cars

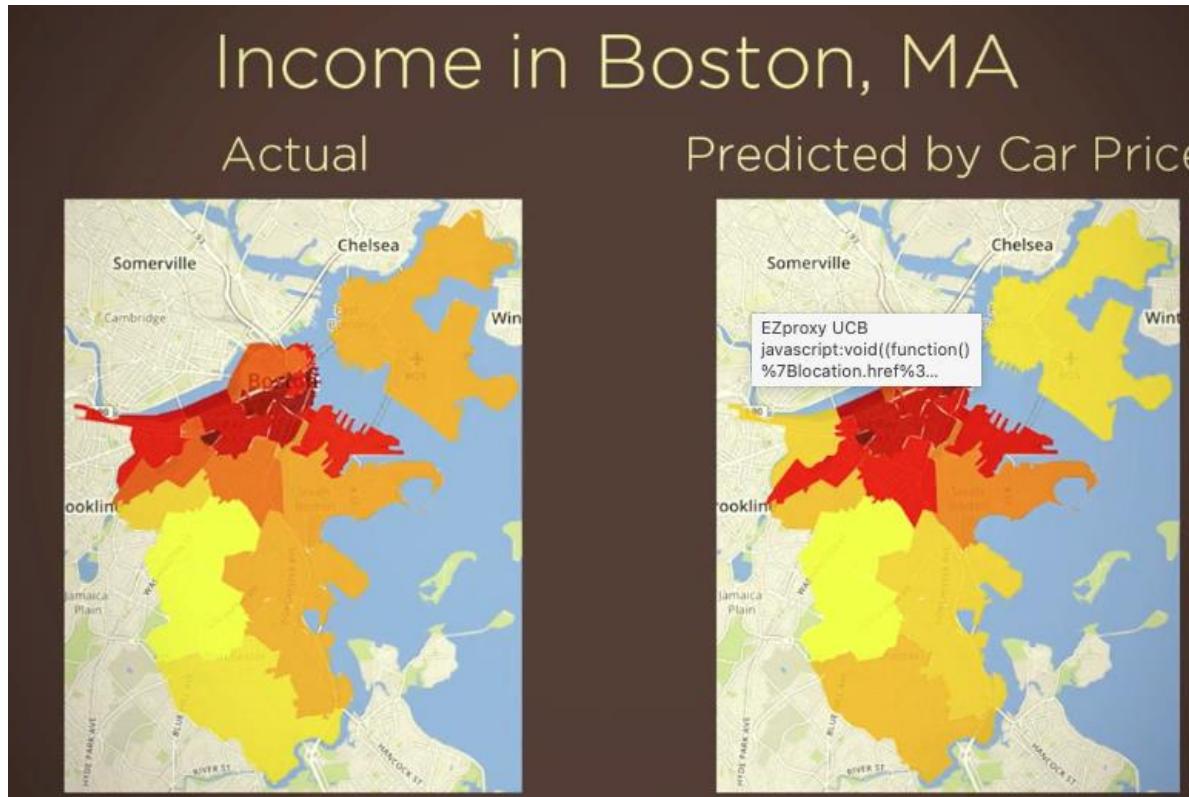


2005 Nissan  
Sentra 1.8 S

2000 Toyota  
4Runner Limited

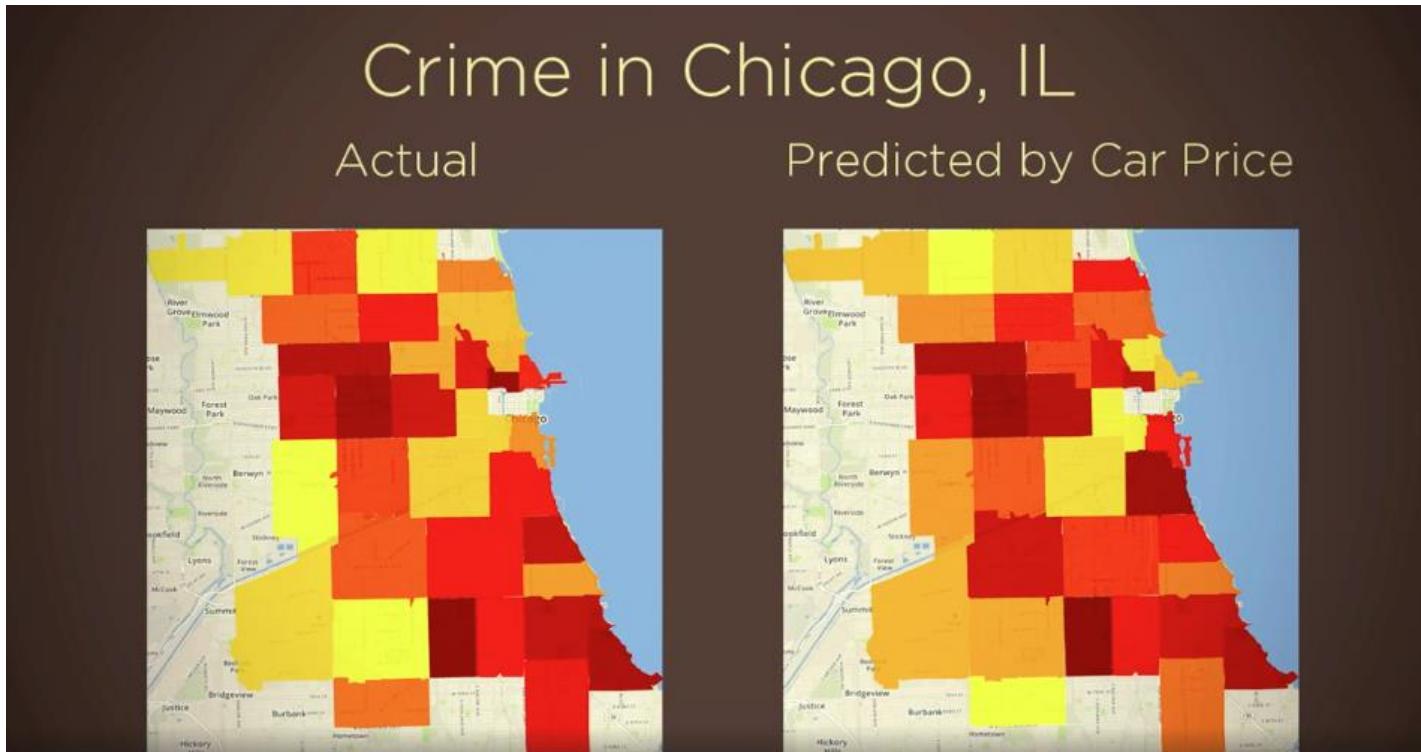
## APPLICATION: DATA GATHERING VIA IMAGE RECOGNITION

Then showed the expected correlation between car model and average income of a region



## APPLICATION: DATA GATHERING VIA IMAGE RECOGNITION

But also found unexpected correlation between car model and crime rate



# ENSURING THE CATEGORY LABELS ARE VERY ACCURATE

- “I ended up training on 500 validation images and then switched to the test set of 1500 images. The labeling happened at a rate of about 1 per minute, but this decreased over time. I only enjoyed the first ~200, and the rest I only did *#forscience*. ....
- The labeling time distribution was strongly bimodal: Some images are easily recognized, while some images (such as those of fine-grained breeds of dogs, birds, or monkeys) can require multiple minutes of concentrated effort. I became very good at identifying breeds of dogs.”

# ENSURING THE CATEGORY LABELS ARE VERY ACCURATE

## Labeling Interface

I developed a labeling interface that would help us evaluate the human performance. It looked similar to, but not identical, to the screenshot below:



# EXERCISE: TRY THE LABELING INTERFACE

## SCROLL DOWN TO SEE THE CATEGORY CHOICES

<https://cs.stanford.edu/people/karpathy/ilsvrc/>



Show answer

Show google prediction

Tibetan mastiff

GoogLeNet predictions:

Tibetan mastiff

Bernese mountain dog

Newfoundland, Newfoundland dog

Gordon setter

Appenzeller



mastiff

Tibetan mastiff



bulldog, English bulldog

French bulldog



spitz

Pomeranian



keeshond

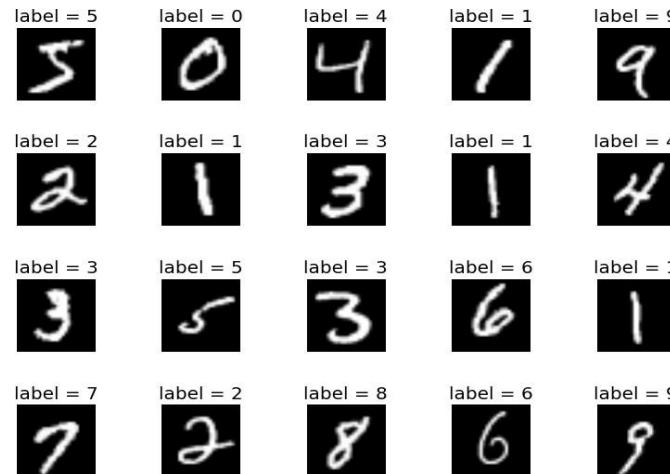


# **IMAGE REPRESENTATION & CONVOLUTIONS**

# MNIST: A Classic Image Dataset

## MNIST Dataset

- 60,000 training examples
- 10,000 test examples
- Rank of best Classifiers and Errors



# How Computers Represent Images

What we see vs. What computers see

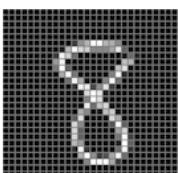


05	02	22	97	38	15	00	40	00	75	04	05	07	78	52	12	50	77	01
49	49	99	40	17	81	18	57	60	87	17	40	98	43	69	41	03	36	62
81	49	31	73	55	79	14	29	93	71	40	67	57	83	30	03	49	13	36
92	70	95	23	04	60	11	42	69	17	65	56	01	32	56	71	37	02	36
22	31	16	71	51	63	13	59	41	92	36	54	22	40	40	23	66	33	13
24	47	33	60	99	03	45	02	44	75	33	53	78	36	24	23	35	17	12
32	98	81	28	64	23	67	10	26	38	40	67	59	54	70	66	18	38	64
70	67	26	20	68	02	62	12	20	95	63	94	39	63	08	40	91	66	49
21	24	55	58	05	66	73	99	26	97	17	78	78	96	83	14	88	34	89
72	21	36	23	09	75	00	76	44	20	45	35	14	00	61	33	97	34	31
95	78	17	53	28	22	75	31	67	15	94	03	80	04	62	16	14	09	53
92	16	39	05	42	96	35	31	47	55	58	88	24	00	17	54	24	36	29
57	86	56	00	48	35	71	89	07	05	44	44	37	44	60	21	54	51	54
17	58	19	80	81	68	05	94	47	69	28	73	92	13	86	52	17	77	04
40	40	04	52	08	83	97	35	99	16	07	97	57	32	16	26	26	79	33
66	66	04	46	68	87	57	62	20	72	03	46	33	67	46	55	12	32	65
69	93	53	04	42	16	73	36	85	39	11	24	94	72	18	08	46	29	32
56	56	20	69	36	41	72	30	23	88	34	61	89	69	82	67	59	85	74
16	16	20	73	35	29	72	31	90	01	74	31	49	71	48	81	41	26	23
57	57	05	01	70	54	71	83	51	54	69	16	92	33	48	61	43	32	01
48	48	48	02	02	02	02	02	02	02	02	02	02	02	02	02	02	02	02

What the computer sees

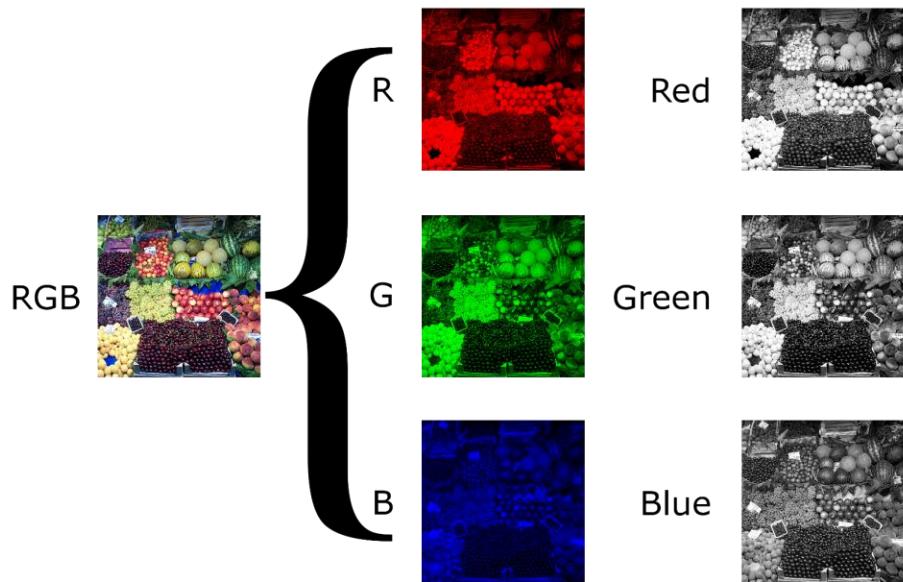
# MNIST Dataset

label = 5	label = 0	label = 4	label = 1	label = 9
				
label = 2	label = 1	label = 3	label = 1	label = 4
				
label = 3	label = 5	label = 3	label = 6	label = 1
				
label = 7	label = 2	label = 8	label = 6	label = 9
				

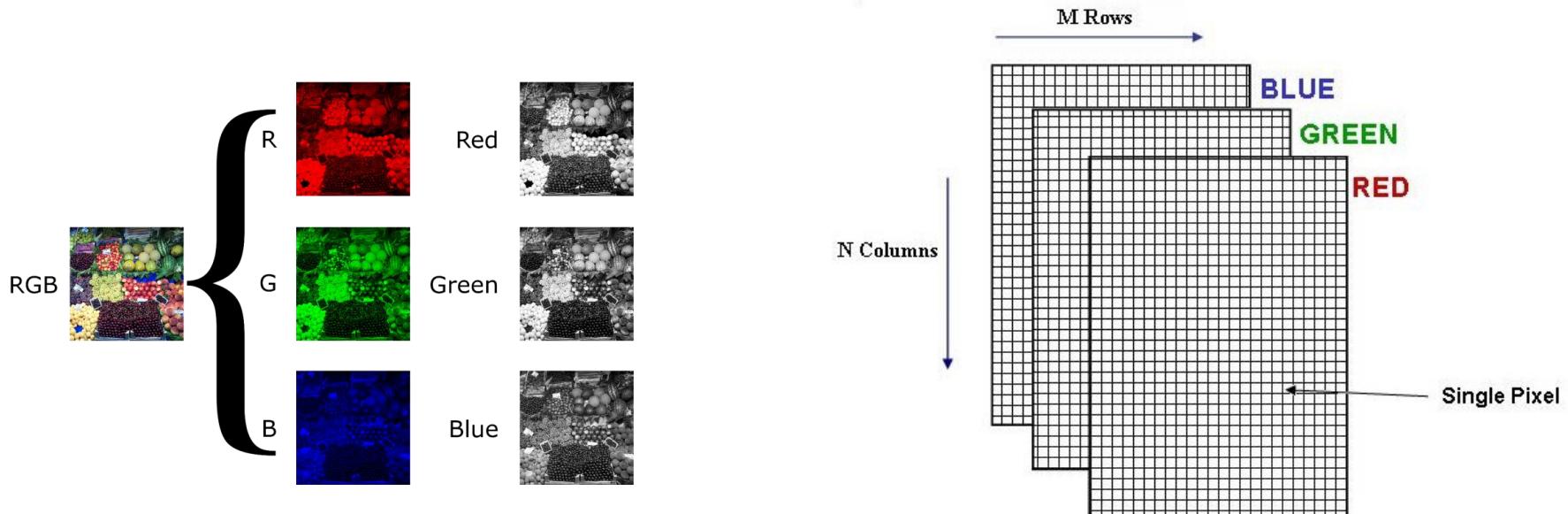


**28 x 28  
784 pixels**

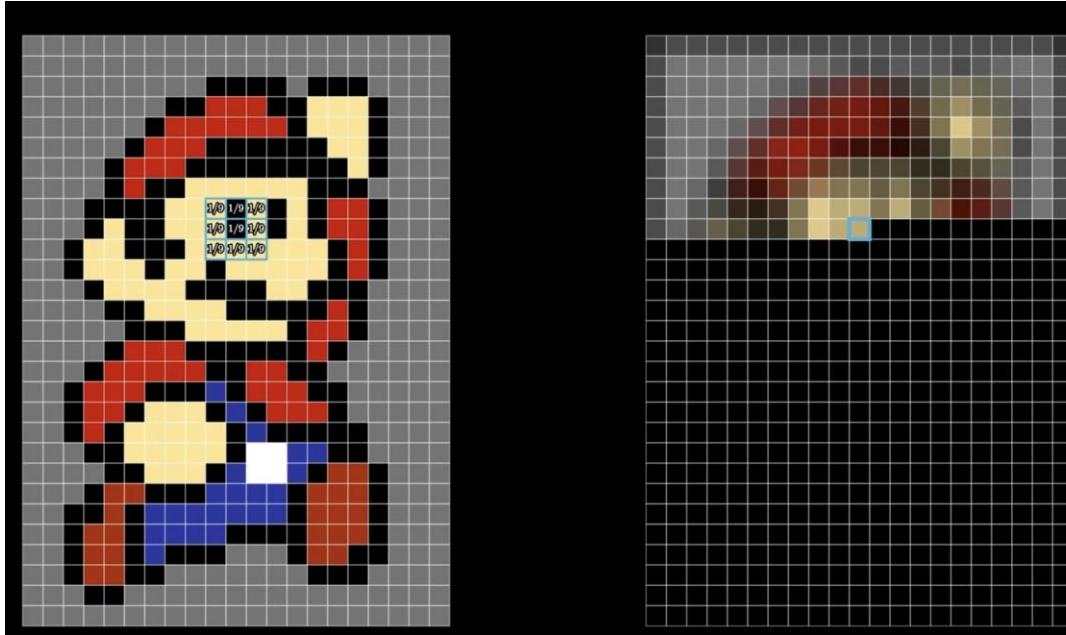
# Color Images



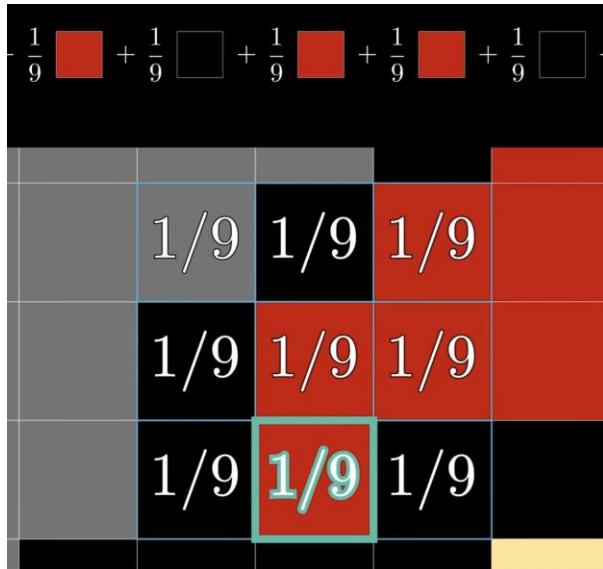
# Color Images



# CONVOLUTIONS TO “AVERAGE” OVER AN IMAGE



# CONVOLUTIONS OVER AN IMAGE



Create a grid of weights  
In this case, 3x3 grid  
Equally weighted with 1/9

Move the grid along the image

At each position in the image, multiply  
the grid values by the value of the color  
of the underlying pixel

Place that value in a new image

That value is the average color of the 9  
pixels

# Convolutions Filters to “Average” over an Image

1	1	1	0	0
0	0	1	1	0
0	1	0	1	1
1	1	0	0	0
1	0	0	1	1

“image”

\*

0	1	0
1	0	1
0	1	0

convolutional filter

# Convolutions Filters to “Average” over an Image

1 x0	1 x1	1 x0	0	0
0 x1	0 x0	1 x1	1	0
0 x0	1 x1	0 x0	1	1
1	1	0	0	0
1	0	0	1	1

“image”

$$\begin{matrix} & \begin{matrix} * & = \end{matrix} & \begin{matrix} 3 \\ \hline \end{matrix} \\ \begin{matrix} 0 & 1 & 0 \\ 1 & 0 & 1 \\ 0 & 1 & 0 \end{matrix} & & \begin{matrix} \hline & \hline & \hline \end{matrix} \end{matrix}$$

convolutional filter

# Convolutions Filters to “Average” over an Image

1	1 $x_0$	1 $x_1$	0 $x_0$	0
0	0 $x_1$	1 $x_0$	1 $x_1$	0
0	1 $x_0$	0 $x_1$	1 $x_0$	1
1	1	0	0	0
1	0	0	1	1

“image”

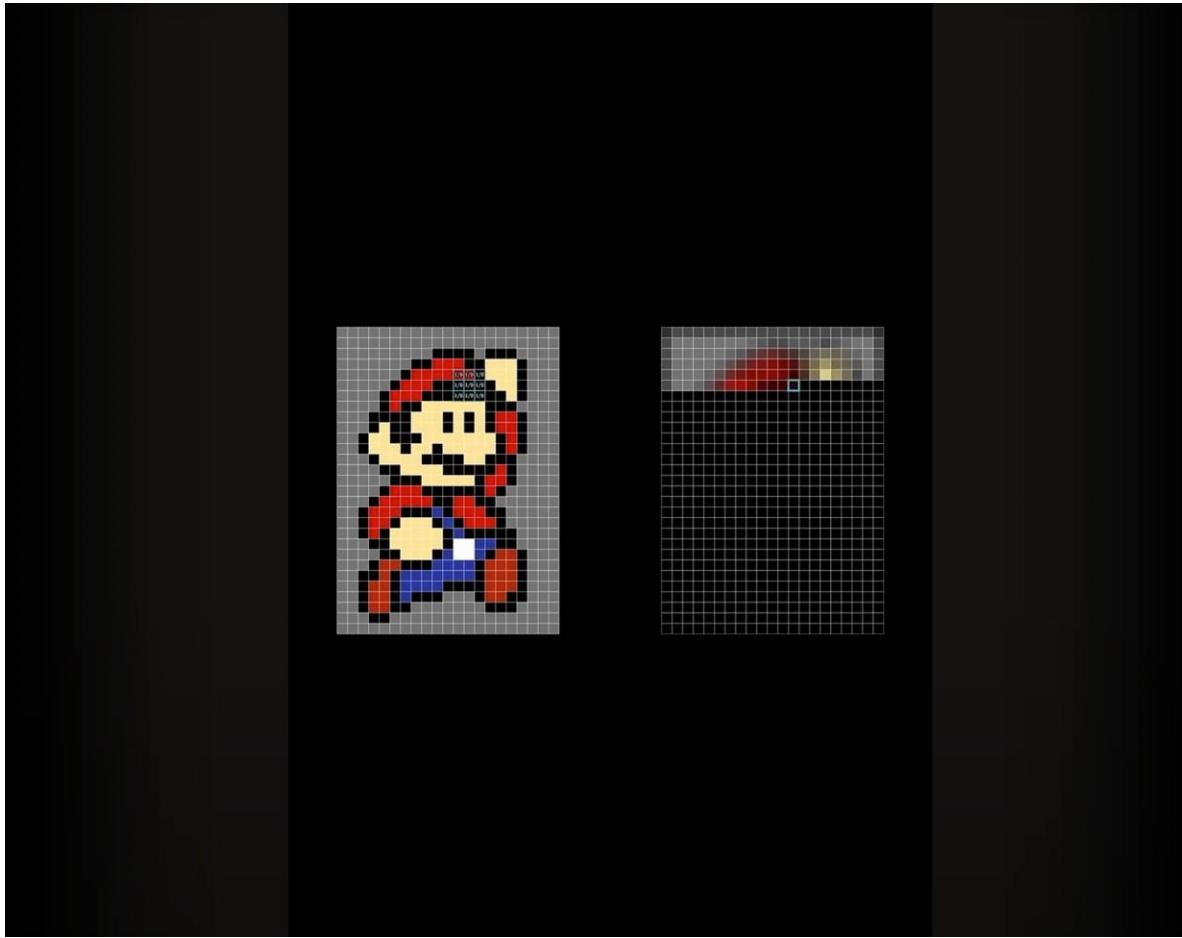
\*

convolutional filter

=

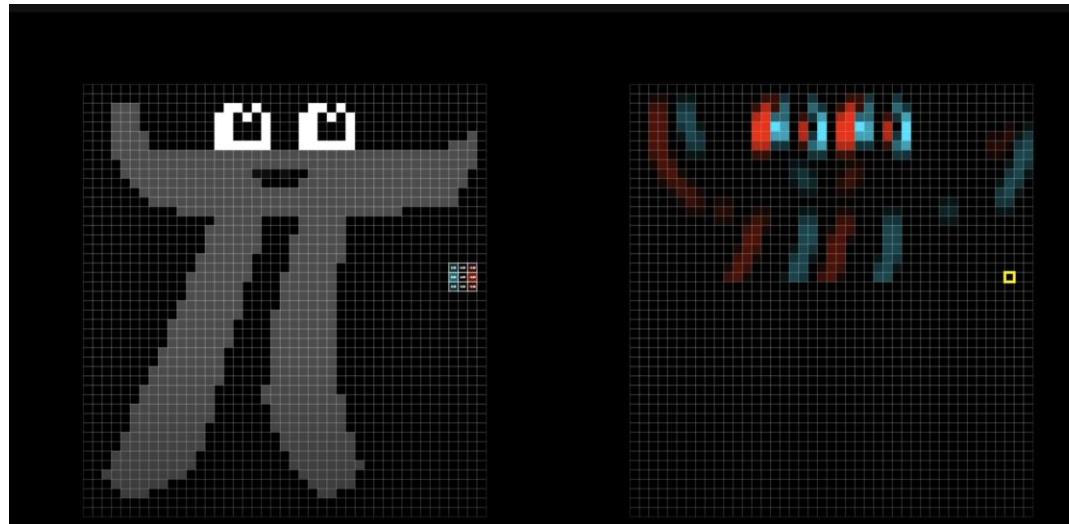
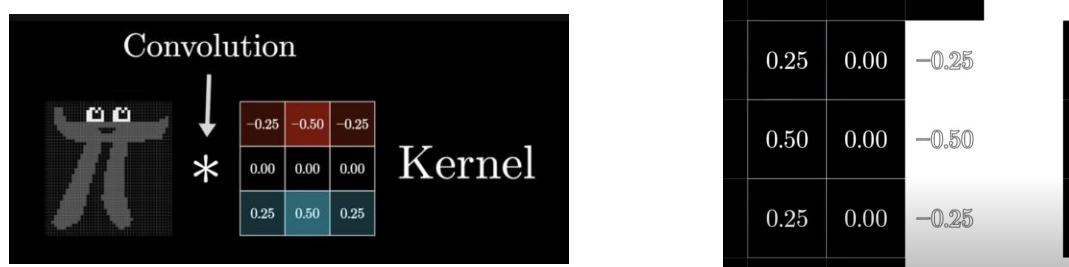
0	1	0
1	0	1
0	1	0

3	2	



But What is a Convolution? 3Blue1Brown <https://www.youtube.com/watch?v=KuXjwB4LzSA>

# CONVOLUTIONS FOR EDGE DETECTION



# CONVOLUTIONS FOR EDGE DETECTION

0.25	0.00	-0.25
0.50	0.00	-0.50
0.25	0.00	-0.25

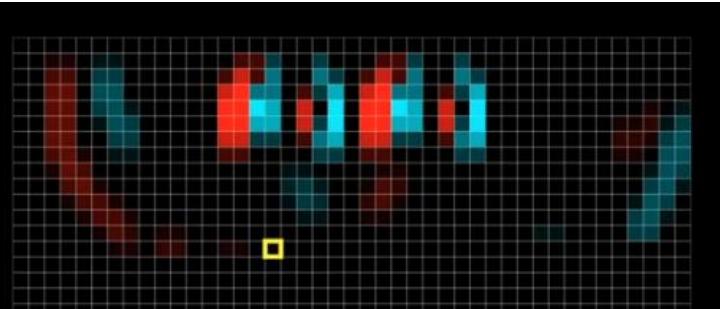
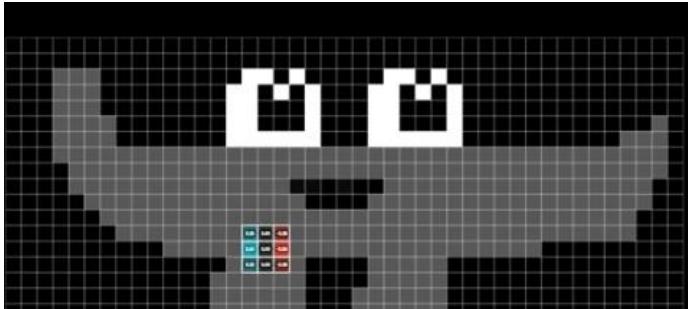
Here, black pixels are 0  
So left & middle column are 0  
(multiply by 0 = 0)

Here, white columns are 1  
So we will get negative results in the  
right column (-.25, -.50, -.25)

If we color the new image:

- negative as red,
- positive as blue
- zero as black

Then we start to see red edges  
on the left and blue edges on  
the right



# CONVOLUTIONS FOR EDGE DETECTION

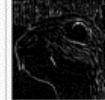
0.25	0.00	-0.25
0.50	0.00	-0.50
0.25	0.00	-0.25

# Varying Convolution Filters

Different filters will produce different feature maps for the same input image

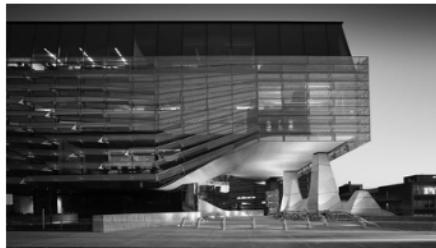


Input image

Operation	Filter	Convolved Image
Identity	$\begin{bmatrix} 0 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \end{bmatrix}$	
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)	$\frac{1}{16} \begin{bmatrix} 1 & 2 & 1 \\ 2 & 4 & 2 \\ 1 & 2 & 1 \end{bmatrix}$	

# Exercise: Match the Convolutions with the Output They Produce

Match the following convolutional filters with the output they produce.

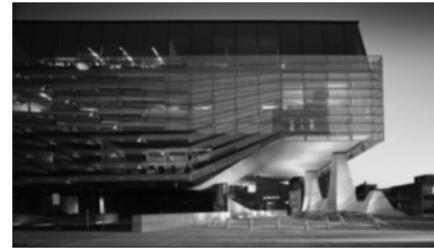


input image

-1	-1	-1
0	0	0
1	1	1

-1	0	1
-1	0	1
-1	0	1

1/9	1/9	1/9
1/9	1/9	1/9
1/9	1/9	1/9

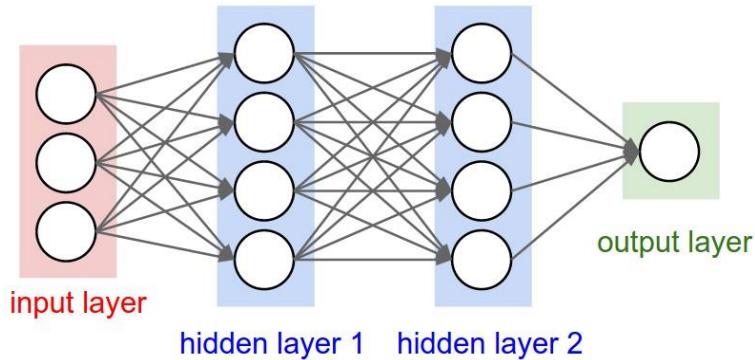


# CONVOLUTIONAL NEURAL NETWORKS (CNNs)

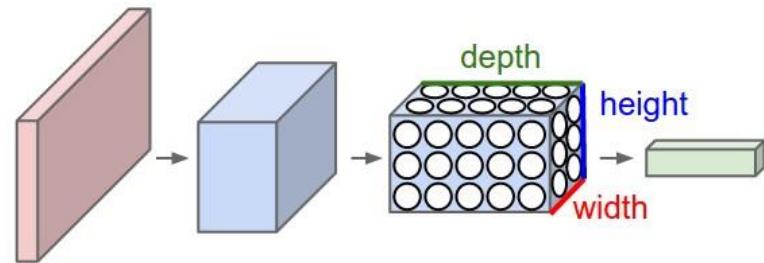
# PRE-DEEP LEARNING IMAGE CLASSIFICATION

- Hand-Crafted Features
- **Texture Features:** Histogram based, Entropy, Haralick features (Co-occurrence matrix), Gray-level run length metrics, Local Binary Pattern, Fractal, etc.
- **Morphological Features:** Hu's moments, Shape features, Granulometry, Bending Energy, Roundness ratio, etc

# Standard vs Convolutional Neural Nets



A standard 3-layer Neural Network.

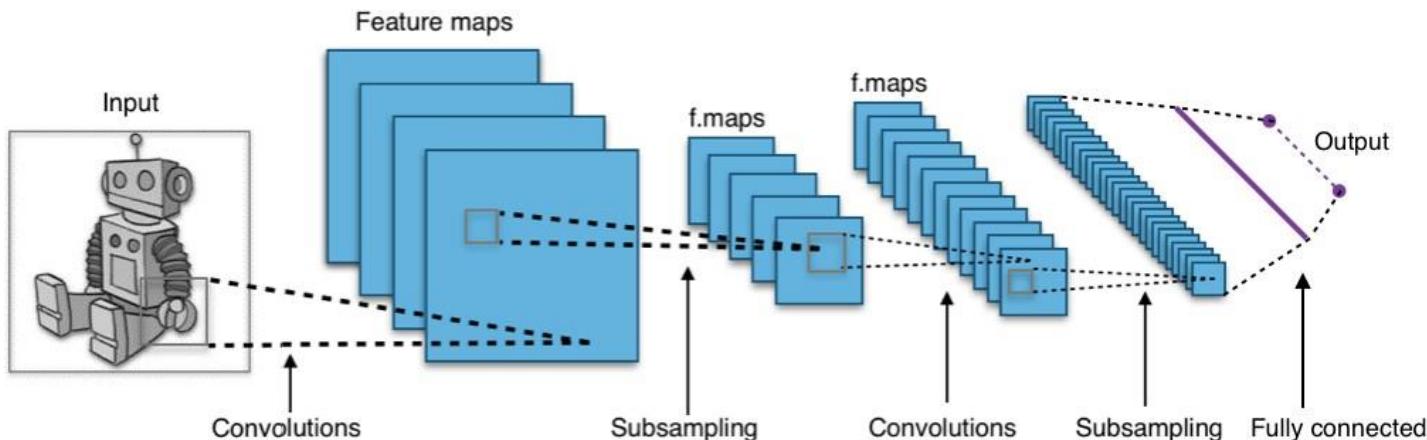


A ConvNet arranges its neurons in three dimensions (width, height, depth), as visualized in one of the layers. Every layer of a ConvNet transforms the 3D input volume to a 3D output volume of neuron activations.

# Convolutional Neural Networks (CNNs)

## Network Architecture

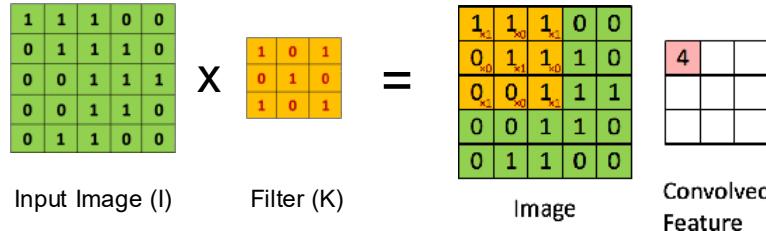
Convolutional Layer, Pooling Layer, Fully Connected Layer



# CNN TERMINOLOGY

## Convolution Operator

$$(I * K)_{xy} = \sum_{i=1}^h \sum_{j=1}^w K_{ij} \cdot I_{x+i-1, y+j-1}$$



- The  $3 \times 3$  matrix ( $K$ ) is called a '**filter**' or '**kernel**' or '**feature detector**' and the matrix formed by sliding the filter over the image and computing the dot product is called the '**Convolved Feature**' or '**Activation Map**' or the '**Feature Map**'.

# Make a convolution filter for each RGB color channel

$$\text{conv}(I, K)_{xy} = \sigma \left( b + \sum_{i=1}^h \sum_{j=1}^w \sum_{k=1}^d K_{ijk} \cdot I_{x+i-1, y+j-1, k} \right)$$

0	0	0	0	0	0	0	...
0	156	155	156	158	158	...	
0	153	154	157	159	159	...	
0	149	151	155	158	159	...	
0	146	146	149	153	158	...	
0	145	143	143	148	158	...	
...	...	...	...	...	...	...	

Input Channel #1 (Red)

0	0	0	0	0	0	0	...
0	167	166	167	169	169	...	
0	164	165	168	170	170	...	
0	160	162	166	169	170	...	
0	156	156	159	163	168	...	
0	155	153	153	158	168	...	
0	154	152	152	157	167	...	
...	...	...	...	...	...	...	

Input Channel #2 (Green)

0	0	0	0	0	0	0	...
0	163	162	163	165	165	...	
0	160	161	164	166	166	...	
0	156	158	162	165	166	...	
0	155	155	158	162	167	...	
0	154	152	152	157	167	...	
...	...	...	...	...	...	...	

Input Channel #3 (Blue)

-1	-1	1
0	1	-1
0	1	1

Kernel Channel #1



308

1	0	0
1	-1	-1
1	0	-1

Kernel Channel #2



-498

0	1	1
0	1	0
1	-1	1

Kernel Channel #3



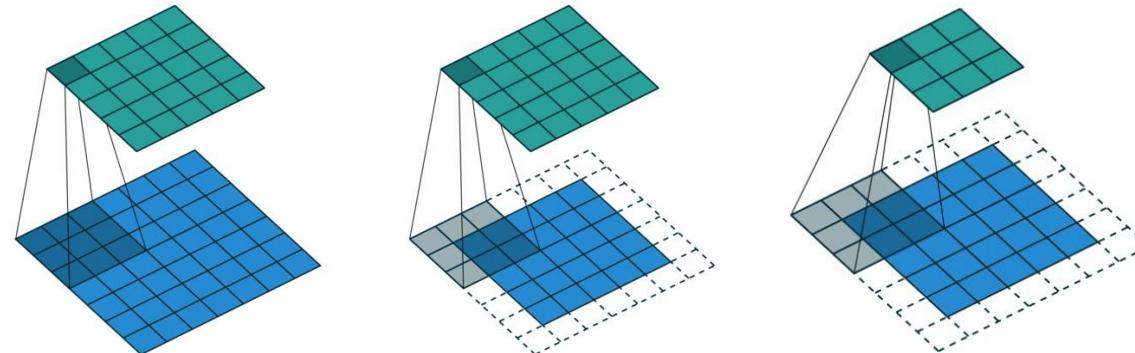
164 + 1 = -25

-25				...
				...
				...
				...
...	...	...	...	...

Bias = 1

# CONVOLUTIONAL NEURAL NETWORKS (CNNs)

- In practice, a CNN learns the values of these **filters** on its own during the **training process**
- Although we still need to specify parameters such as **number of filters**, **filter size**, **padding**, and **stride** before the training process

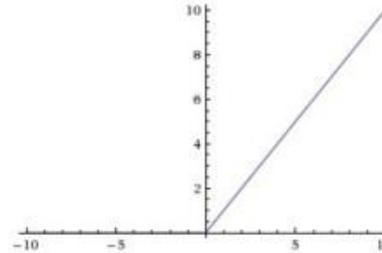


# CNN ACTIVATION LAYERS

## Activation Layer (ReLU)

- An additional operation called Rectified Linear Unit (ReLU) which can be used after every Convolution operation

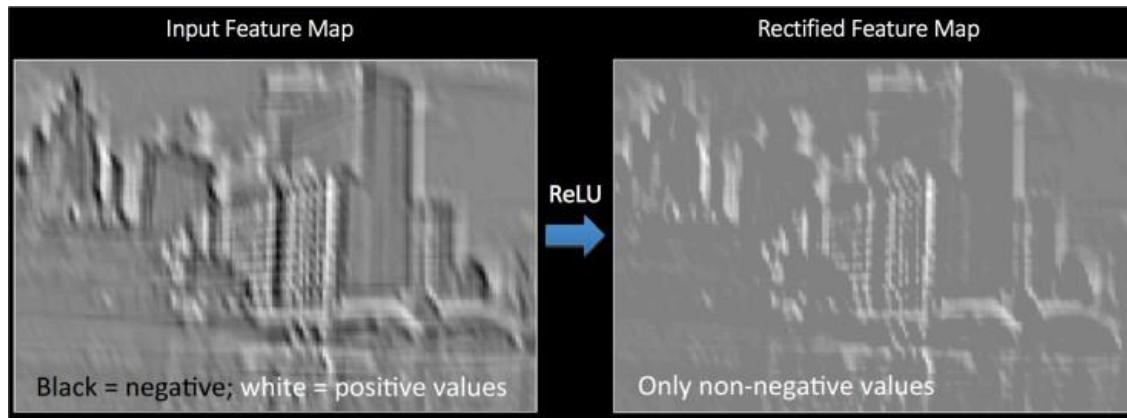
$$\text{Output} = \text{Max(zero, Input)}$$



- Basically, ReLU is an element wise operation (applied per pixel) and replaces all negative pixel values in the feature map by zero
- The purpose of ReLU is to introduce non-linearity to the network

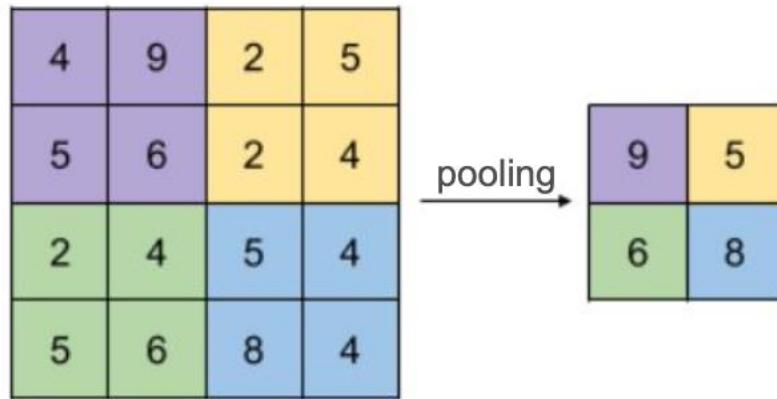
# CNN ACTIVATION LAYERS

## Activation Layer (ReLU)



- Other nonlinear functions such as **tanh** or **sigmoid** can also be used instead of ReLU, but ReLU has been found to perform better in most situations.

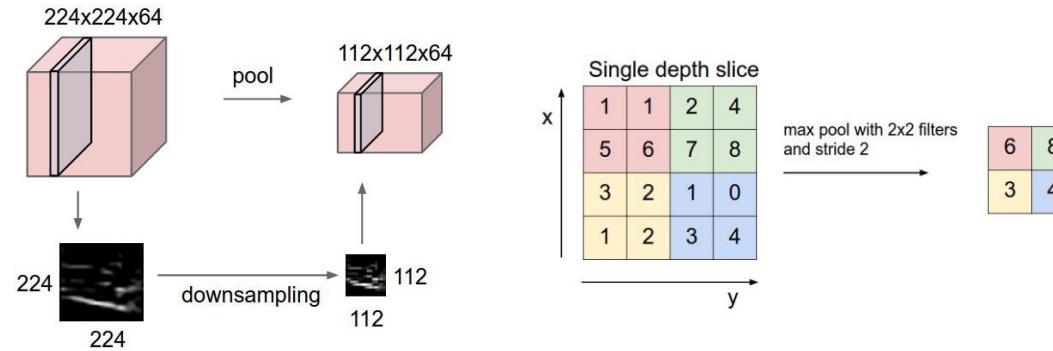
# CNN MAX POOLING



# CNN POOLING LAYER

## Pooling Layer

Pooling layer **downsamples** the volume spatially, **independently** in **each depth** slice of the input

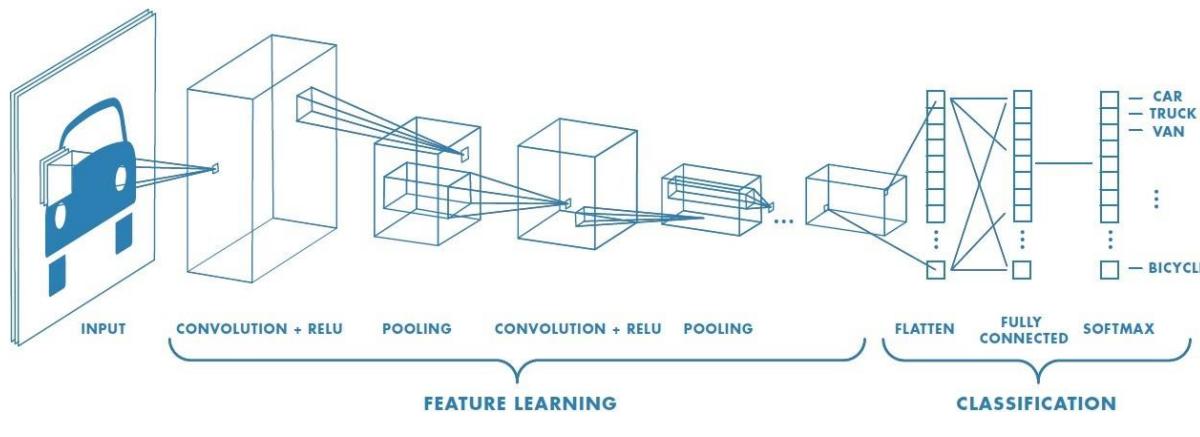


The most common downsampling operation is **max**, giving rise to **max pooling**, here shown with a stride of 2

# CNN FULLY CONNECTED LAYER FOR CLASSIFICATION

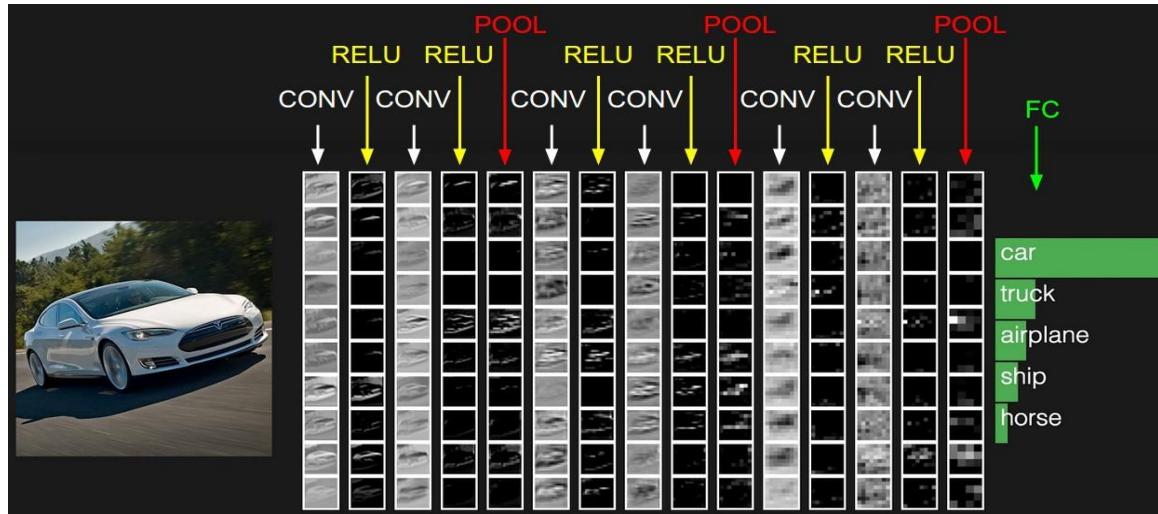
## Fully Connected Layer

- Neurons in a fully connected layer have full connections to all activations in the previous layer

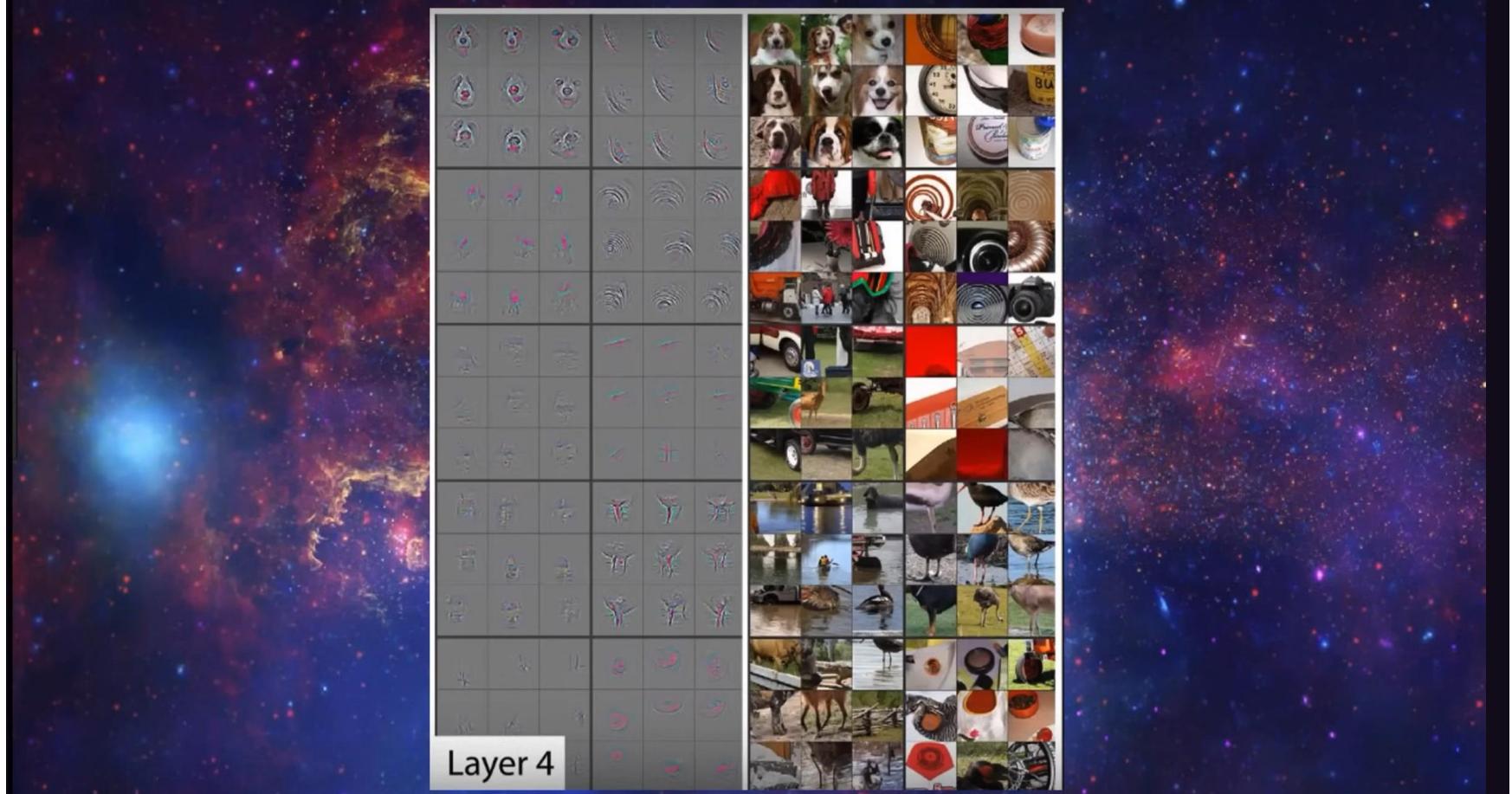


# CNN Output at Each Layer

**Example:** Input >> [ [ Conv >> ReLU ] \* 2 >> Pool ] \* 3 >> FC



The leftmost column stores the raw image pixels and the rightmost stores the class scores (right). Since it's difficult to visualize 3D volumes, we show image slices in rows. The last layer volume holds the scores for each class.

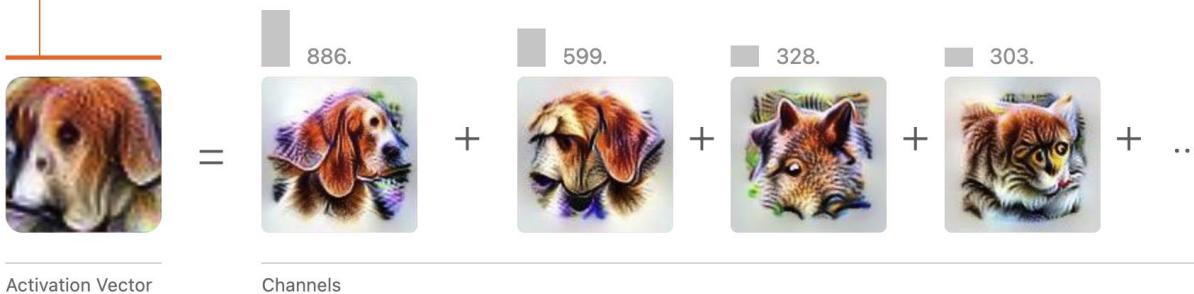


[https://www.youtube.com/watch?v=XRhxdVk\\_sls](https://www.youtube.com/watch?v=XRhxdVk_sls)

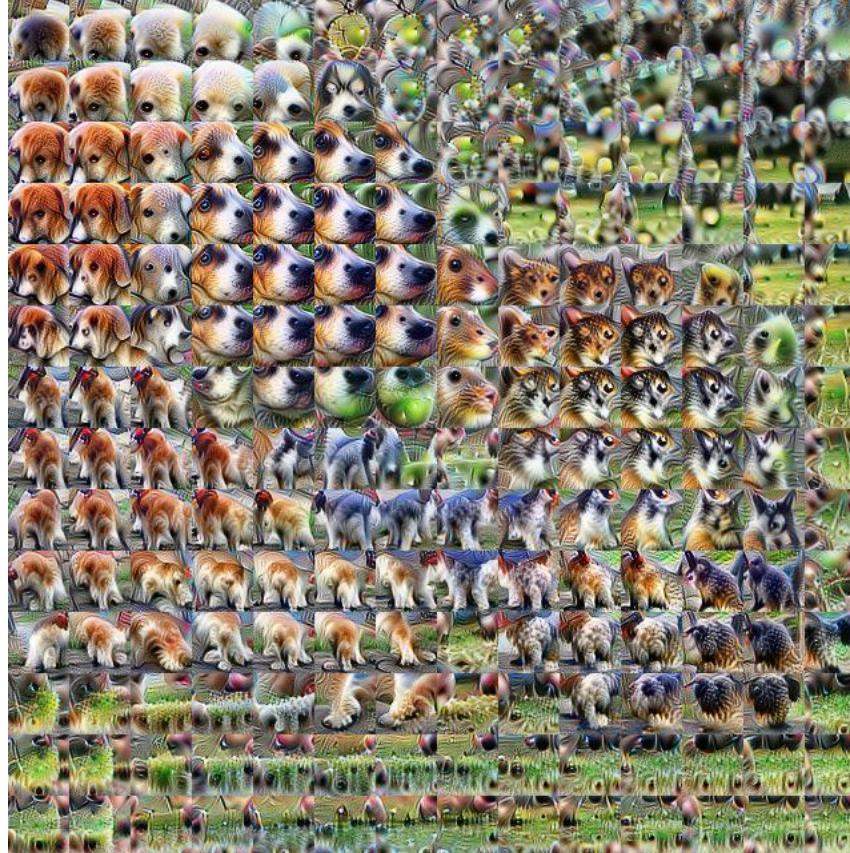
# What Does the Network See?



Semantic dictionaries give us a fine-grained look at an activation: what does each single neuron detect? Building off this representation, we can also consider an activation vector as a whole. Instead of visualizing individual neurons, we can instead visualize the *combination* of neurons that fire at a given spatial location. (Concretely, we optimize the image to maximize the dot product of its activations with the original activation vector.)



# Explore CNN Features



<https://distill.pub/2018/building-blocks/>

# Advantages of CNNs for Image Classification

## Convolutions

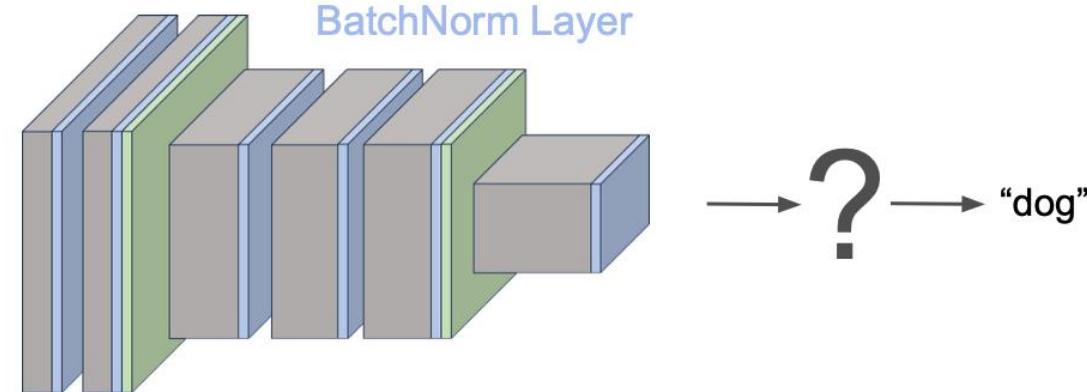
Maintain spatial relation between pixels  
Reduce number of parameters through weight sharing

## Pooling

Captures key information from across different areas of the feature maps  
Together with convolutions allows for translational invariance

## BatchNorm

Increases speed and stability of training



# CNNs IN SUMMARY

- CNNs are primarily designed to process and analyze visual data, such as images and videos.
- Key components: convolution layers, pooling layers, activation functions, normalization layers
- Advantages:
  - *Translational Invariance*
  - *Parameter sharing*
  - *Feature learning*
- Can be trained with backprop
- Used for tasks such as segmentation, classification, object detection, etc.
- Achieved results around 86% on ImageNet
- Have been superseded today by Vision Transformers

# To Learn More...

To learn about computer vision in detail:



## Info 290T Computer Vision **3 units**

---

**Course Description**

This course introduces the theoretical and practical aspects of computer vision, covering both classical and state of the art deep-learning based approaches. This course covers everything from the basics of the image formation process in digital cameras and biological systems, through a mathematical and practical treatment of basic image processing, space/frequency representations, classical computer vision techniques for making 3D measurements from images, and modern deep-learning based techniques for image classification and recognition.

**Prerequisites**

Linear algebra and Python (INFO 206A & B or equivalent).



Prof Hany Farid: