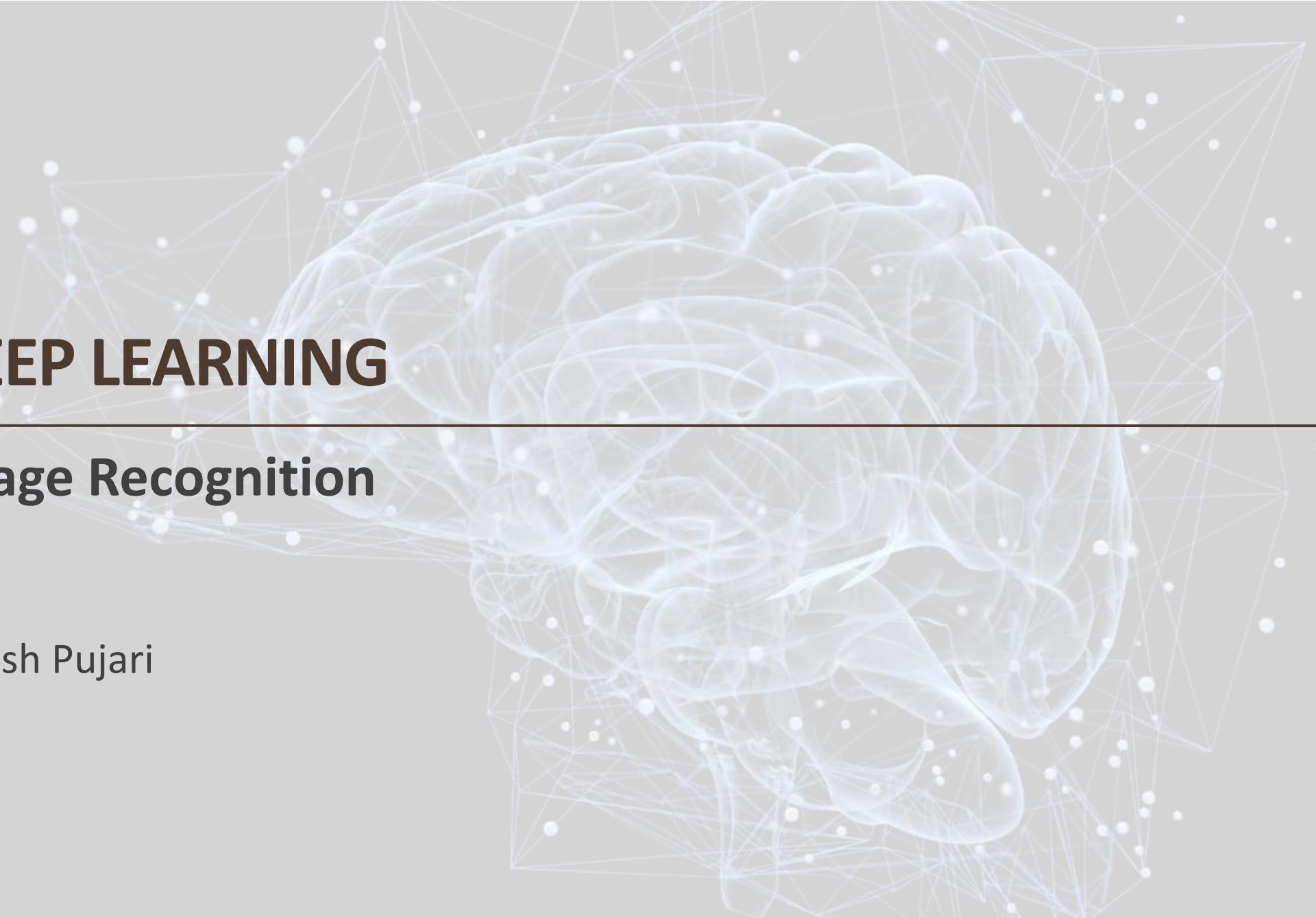


# DEEP LEARNING

---

## Image Recognition

Ashish Pujari



# Image Recognition

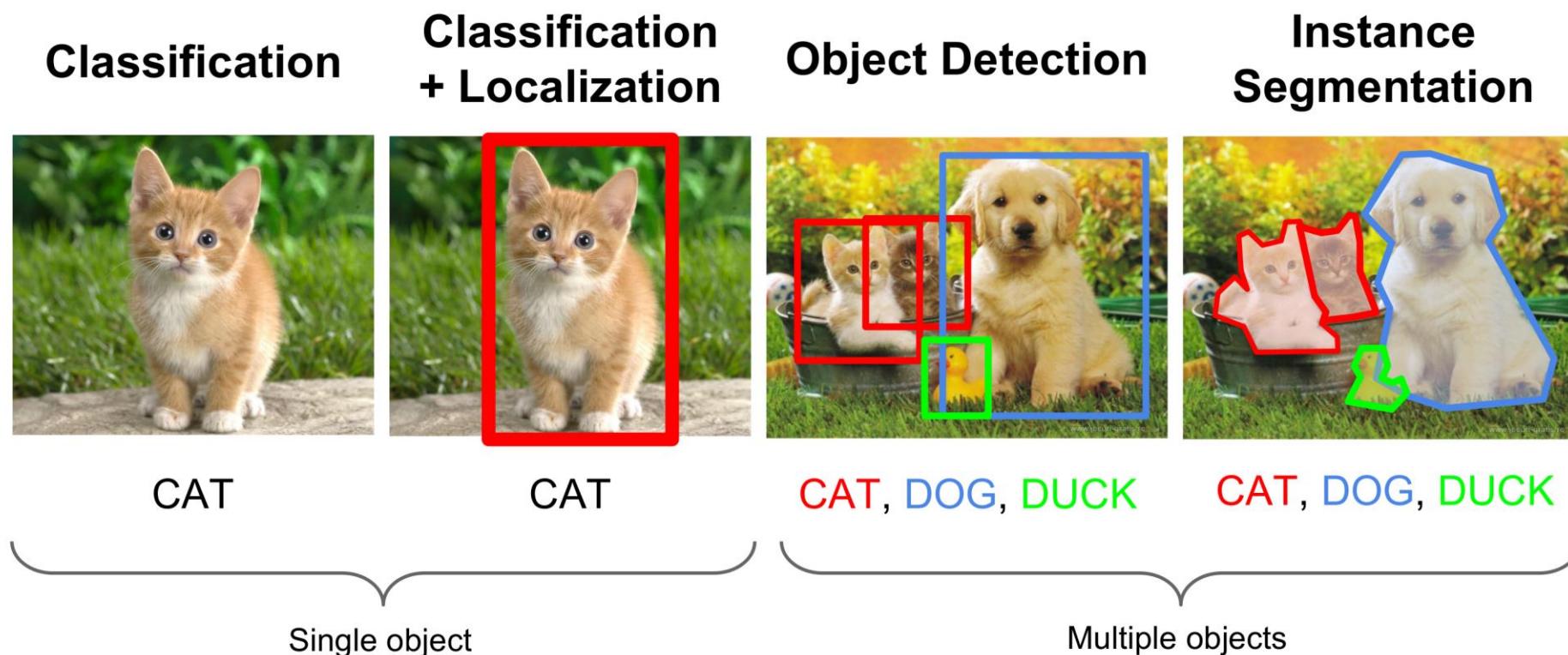
1. Computer Vision
2. Convolutional Neural Networks
3. CNN Training
4. CNN Applications

# COMPUTER VISION

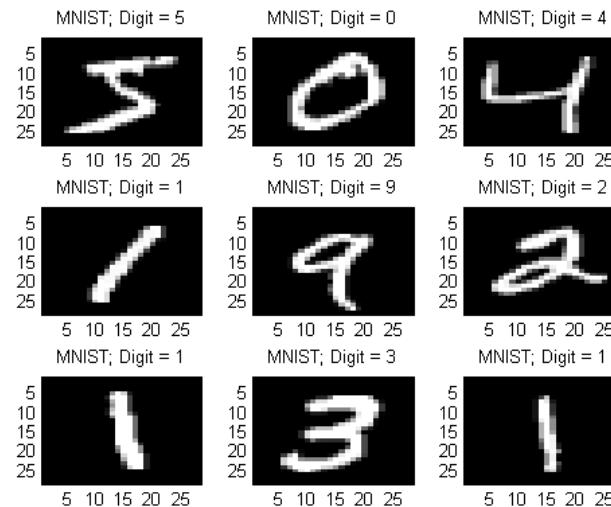
---

Image Datasets, Computer Vision Problems, Invariance

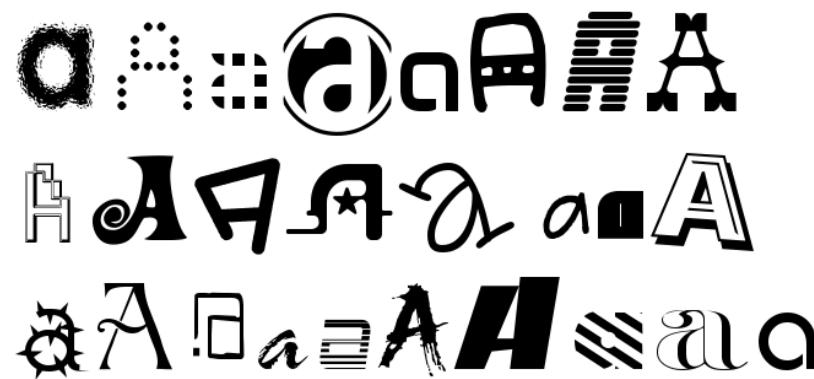
# Computer Vision



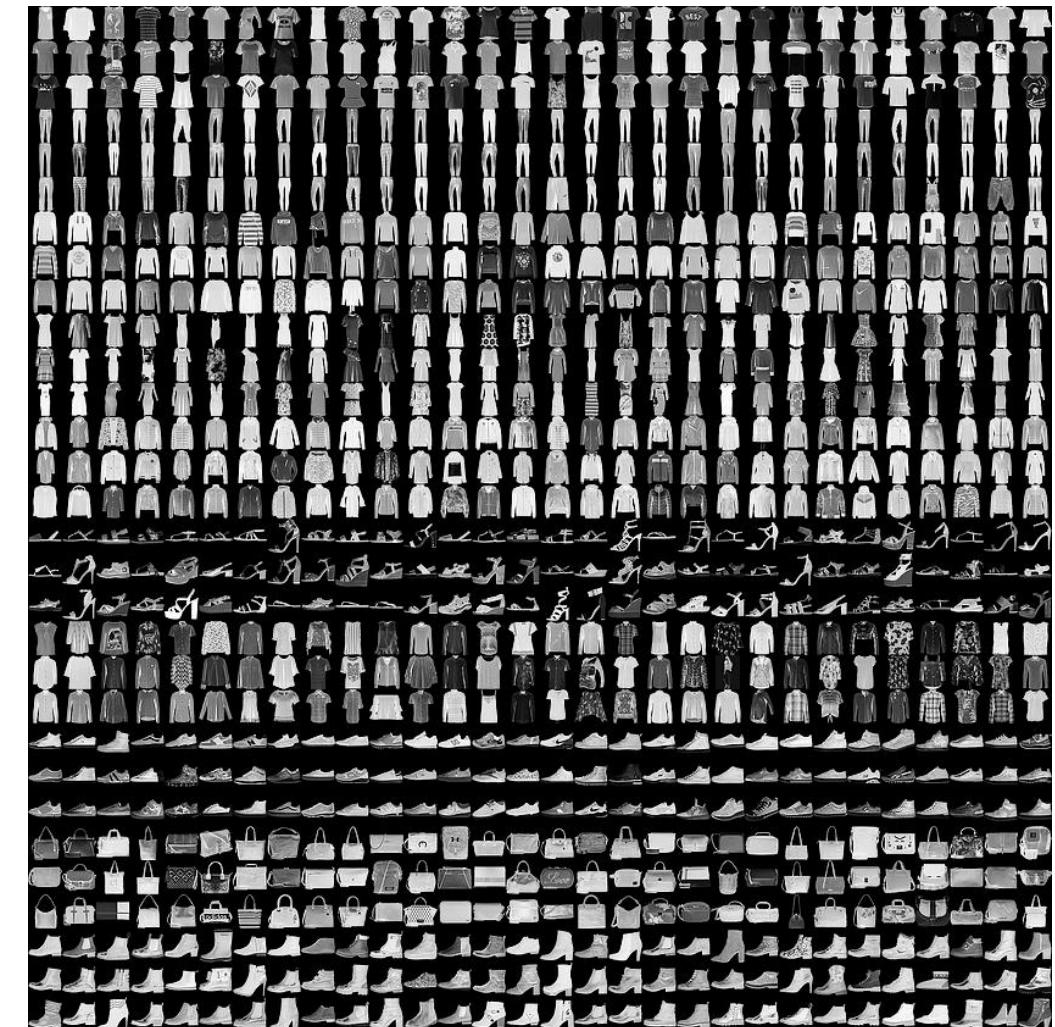
# Images - Benchmarking Datasets



<http://yann.lecun.com/exdb/mnist/>



<http://yaroslavvb.blogspot.com/2011/09/notmnist-dataset.html>



<https://github.com/zalandoresearch/fashion-mnist>

# Images - Benchmarking Datasets

**airplane****automobile****bird****cat****deer****dog****frog****horse****ship****truck**

<https://www.cs.toronto.edu/~kriz/cifar.html>

## COCO Explorer

COCO 2017 train/val browser (123,287 images, 886,284 instances). Crowd labels not shown.



bicycle ✕ person ✕

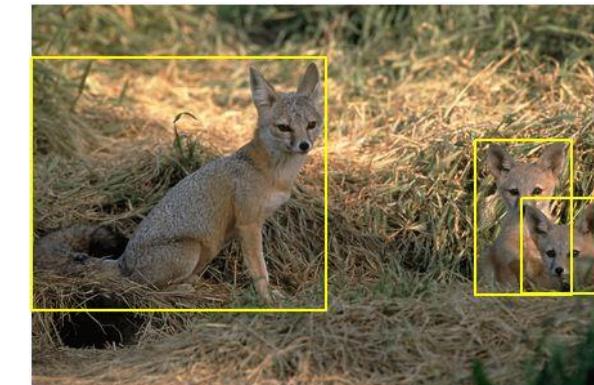
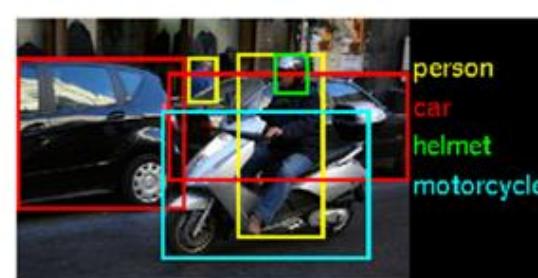
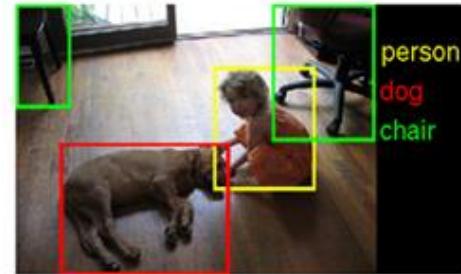
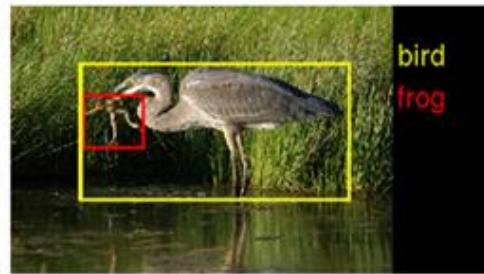
search

3401 results



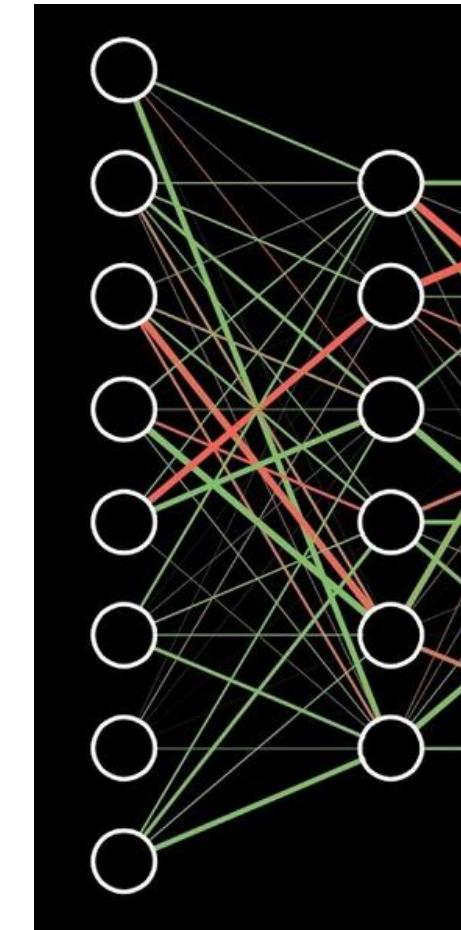
<http://cocodataset.org/#explore>

# Images - Competition Datasets



# MLP Limitations

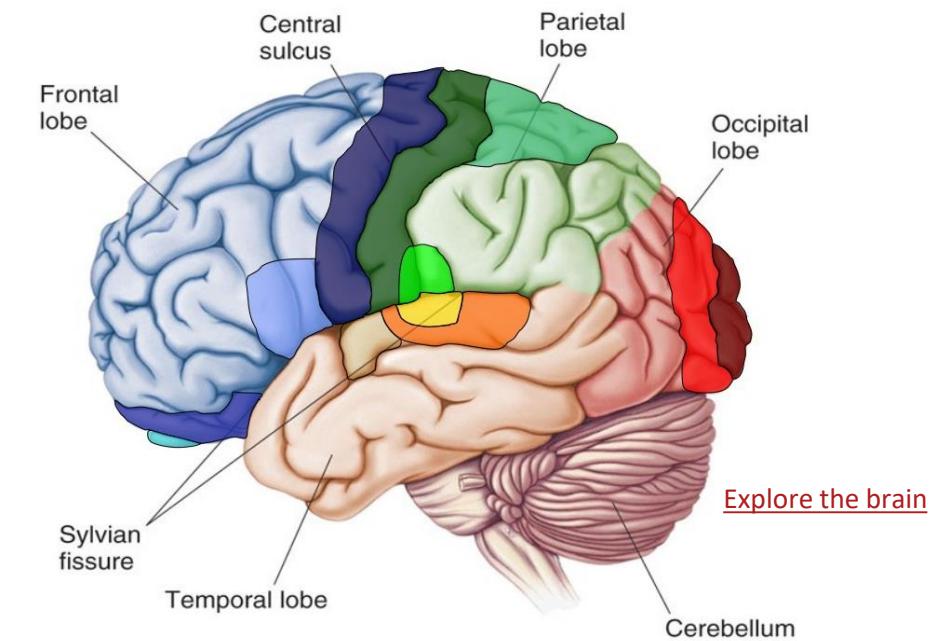
- High Dimensionality
  - $150\text{px} \times 150\text{px} = 22500$
  - Color Images:  $22500 \times 3 \text{ RGB} = 67500$
  - 3D images, video
- Computational Cost
  - MLPs suffer from high dimensionality due to full matrix multiplication between layers
- Spatial Information
  - MLPs are insensitive to order of elements in input vector. A images have intrinsic spatial information that is lost.



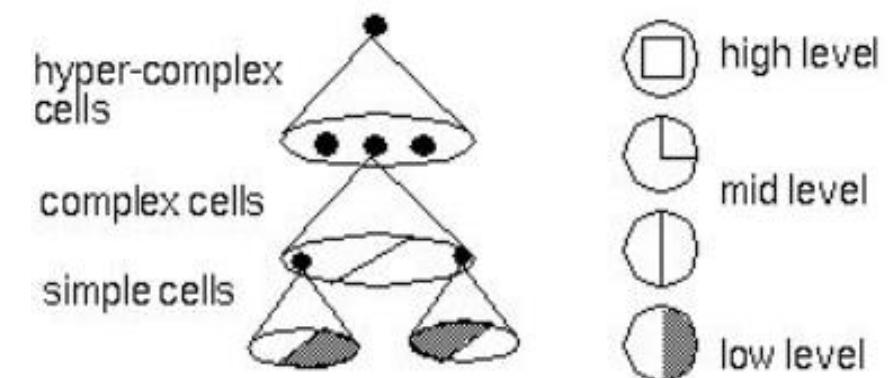
MLP – dense connections

# Inspiration - Physiology of Vision

- Animal Visual Cortex
  - Contains complex arrangements of cells
  - Simple cells responsible for detecting edges
  - Complex cells have larger receptive fields with a degree of spatial invariance or contrast-insensitive manner
  - Higher level features respond to patterns of activation in lower level cells, and propagating activation upwards to still higher level cells
  - Experiments on cats showed that different cells responded to only certain kinds of stimuli such as an edge or a particular color

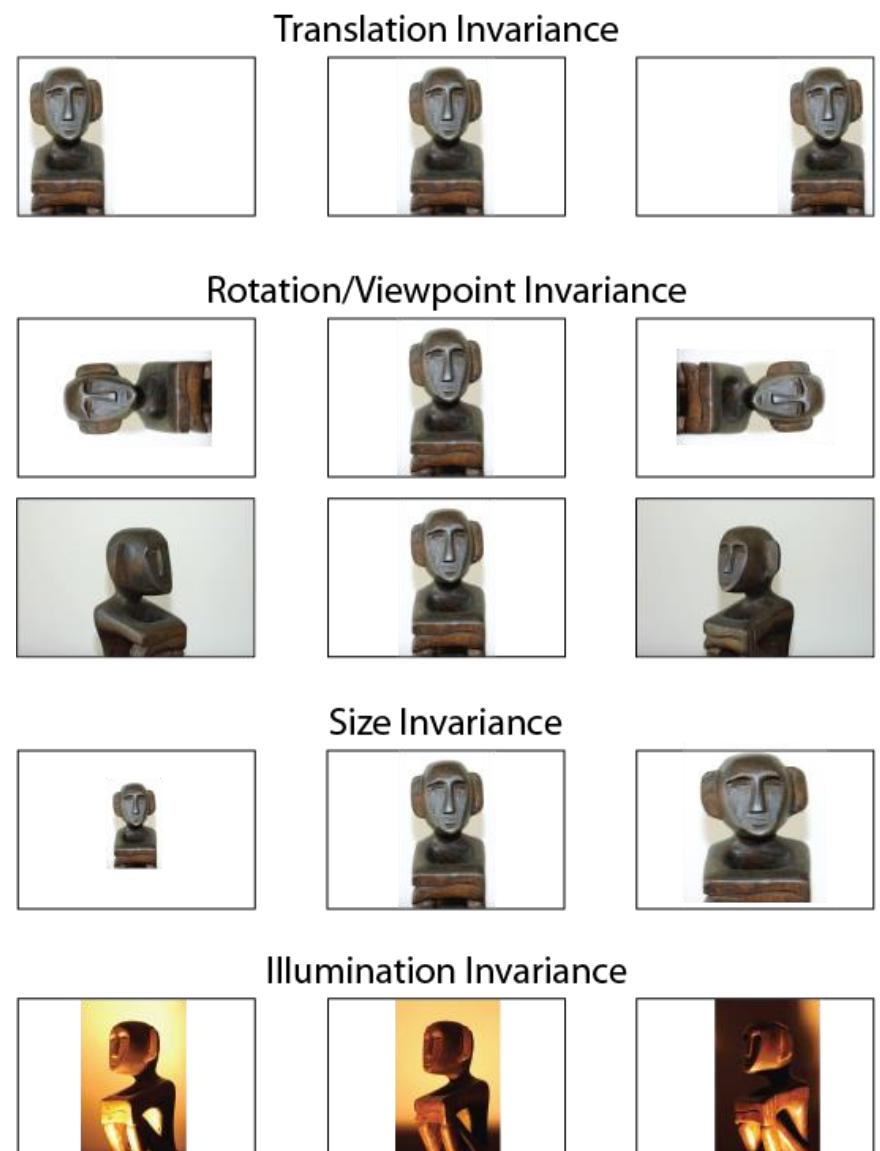


## featural hierarchy



# Statistical Invariance

- Invariance
  - Translation/Spatial invariance
  - Orientation/Rotation invariance
  - Size/Scale invariance
  - Illumination invariance
- Challenge
  - Recognition in the presence of common transformations such as translation, rotation, etc.
- Solution
  - Since images have high statistical invariance, we could squeeze out the spatial information and map parameters that focus on content

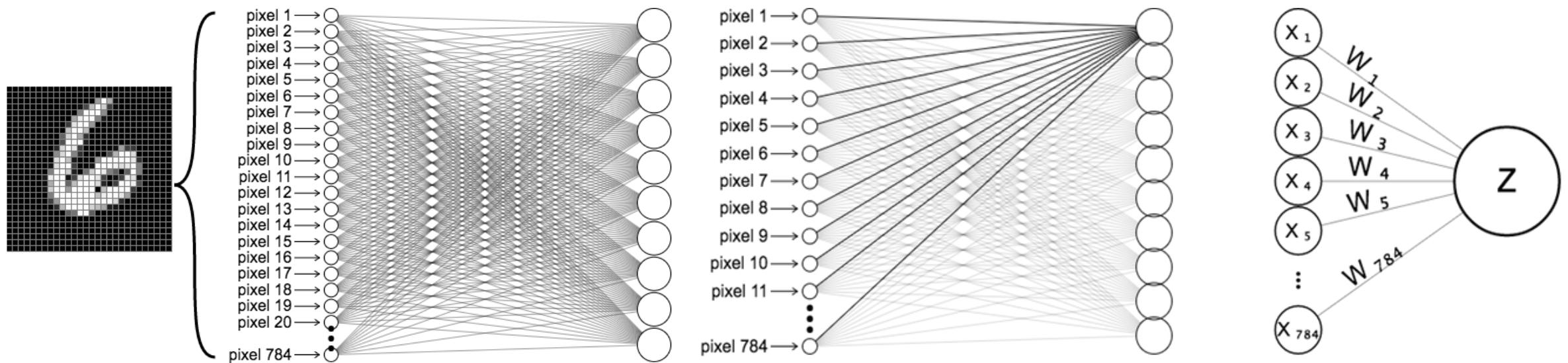


# CONVOLUTIONAL NEURAL NETS (CNN)

---

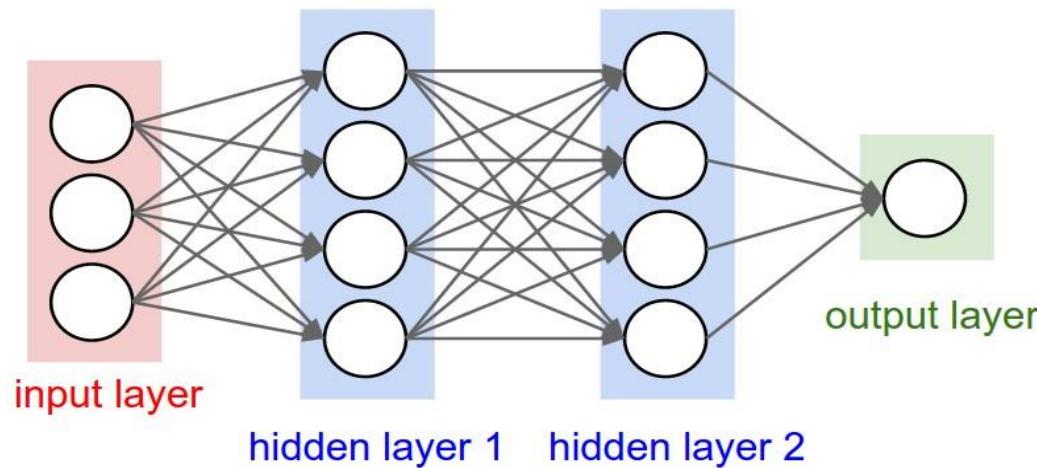
# MLP Challenges

- When we input an image into an MLP, we need to unroll the pixels into a single column of neurons E.g. For a 32x32 MNIST image the number of input pixels will be 784 serially fed into an MLP through the input layer

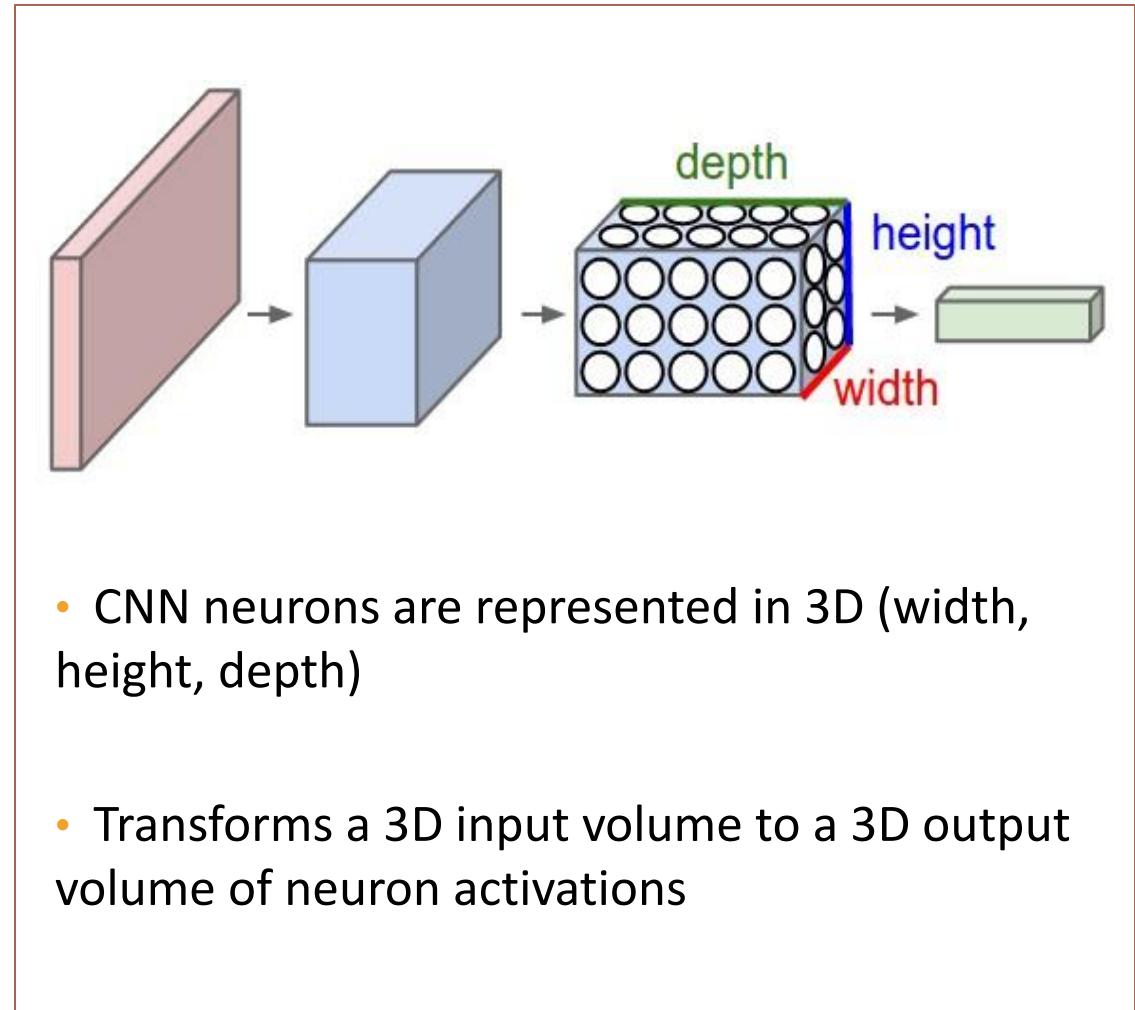


1-layer neural network for MNIST. The 10 output neurons correspond to our classes, the 10 digits from 0 to 9.

# MLP vs CNN

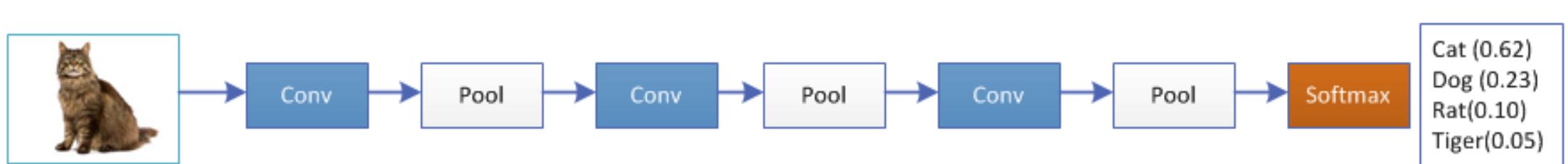


- MLPs suffer from curse of dimensionality due to full connectivity between nodes.
- MLPs do not scale well to higher resolution images



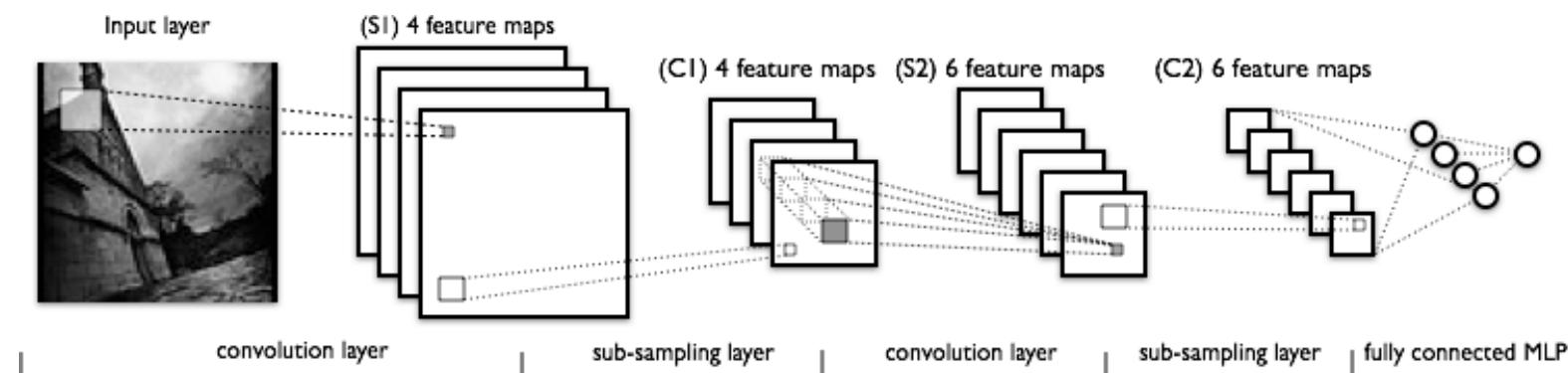
# Convolutional Neural Networks (CNN or ConvNets)

- Deep neural networks inspired by a hierarchy of feature detectors in the visual cortex.
- Each layer progressively extracts higher and higher-level features of the image
- Leverage 3 important principles:
  - Sparse Connections
  - Parameter Sharing
  - Equivariant Representations



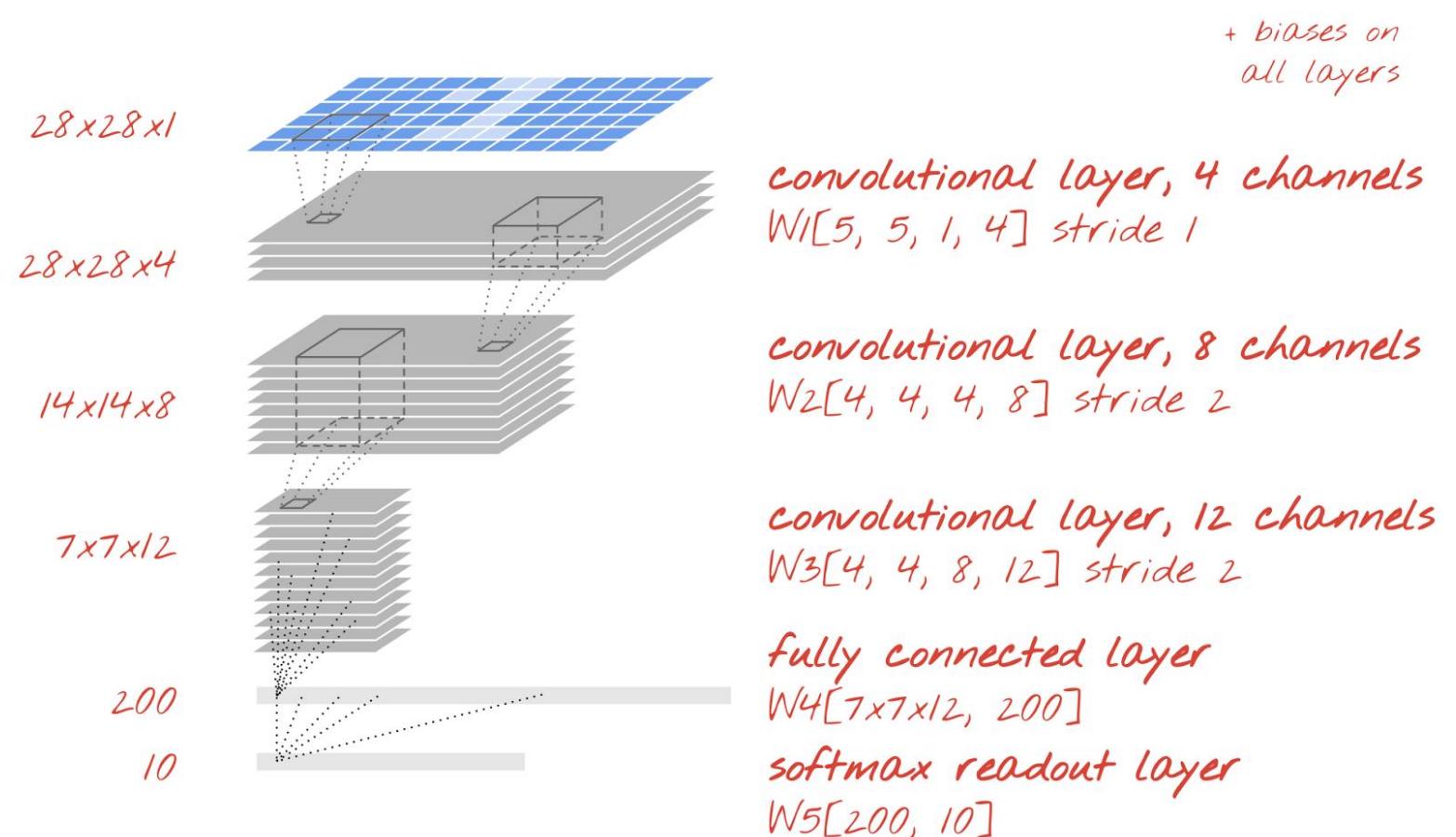
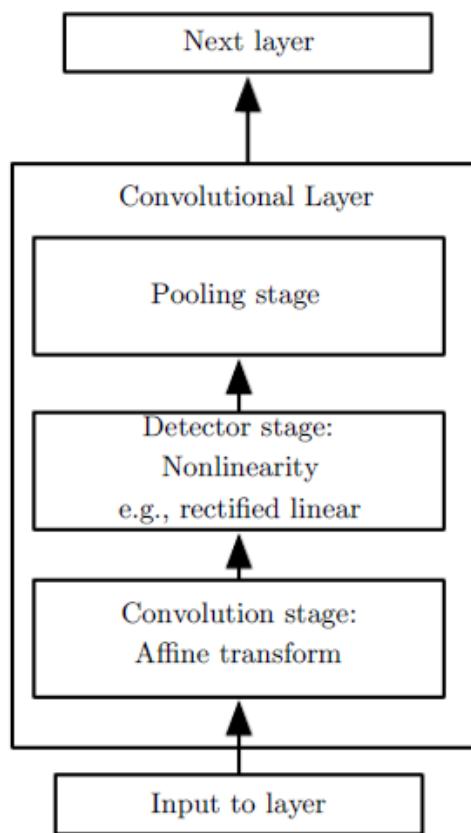
# CNN Advantages

- Composability
  - The output of one convolutional layer becomes an input into another. With each layer, the network can detect higher-level, more abstract features.
- Reduced Dimensionality
  - Reduce dimensionality by orders of magnitude over ANNs by learning invariant features by using weight sharing
  - Smaller representations are achieved by shared parameters

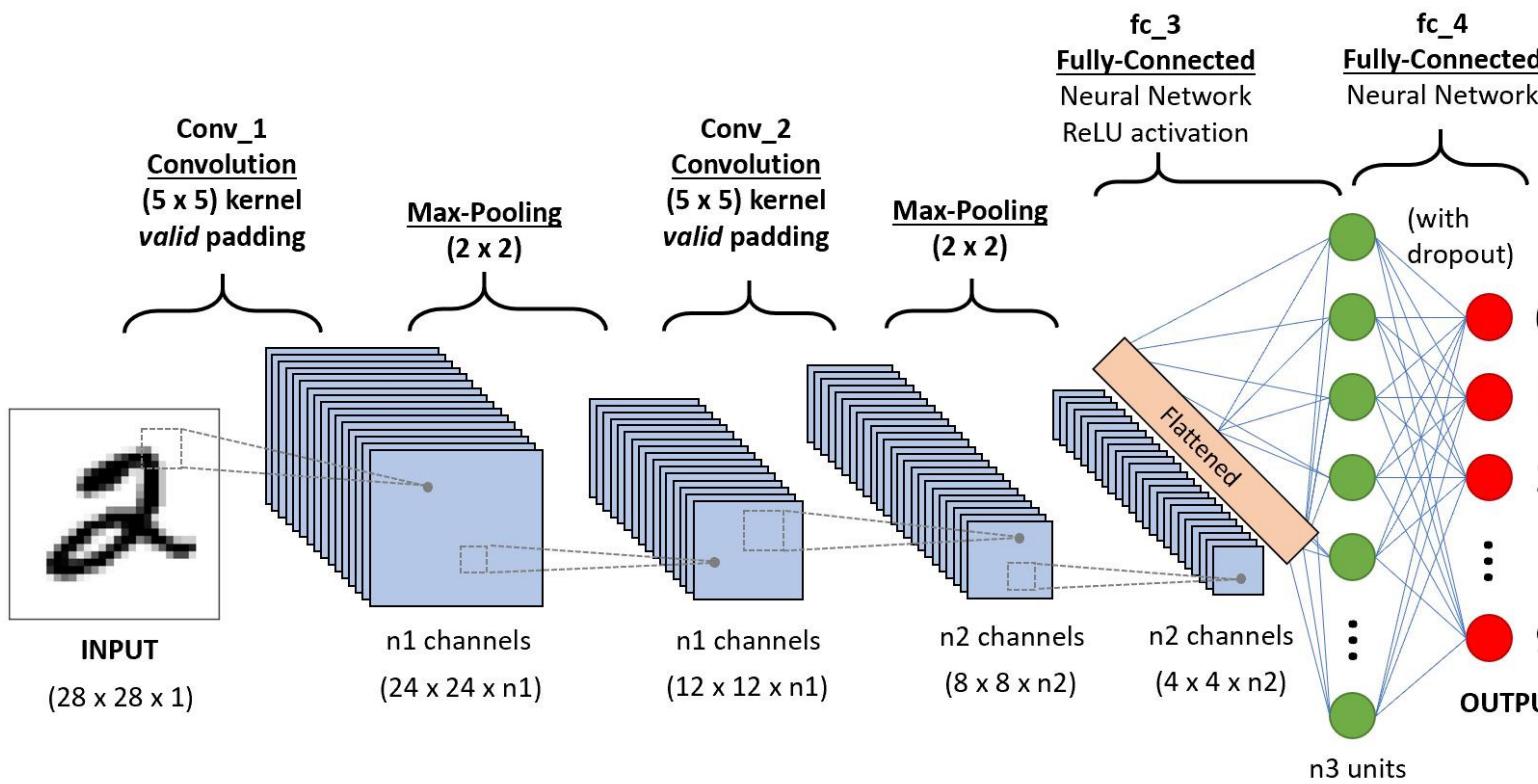
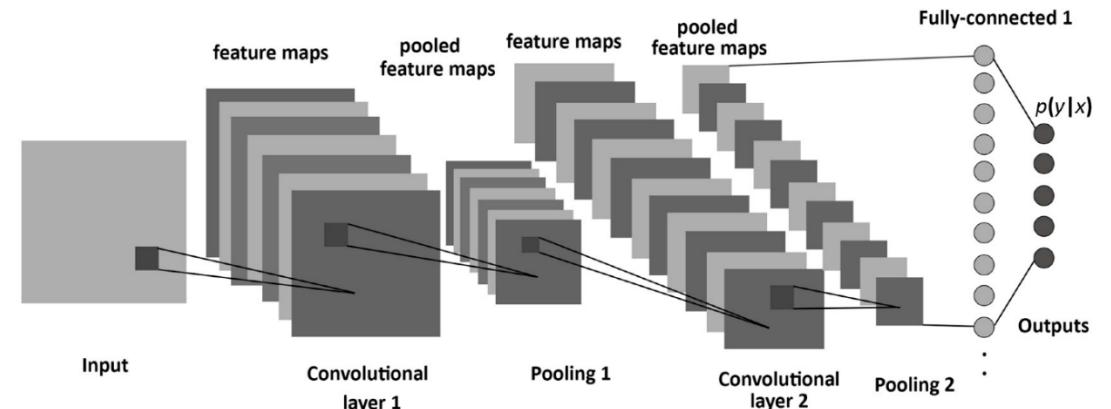


# CNN Layers

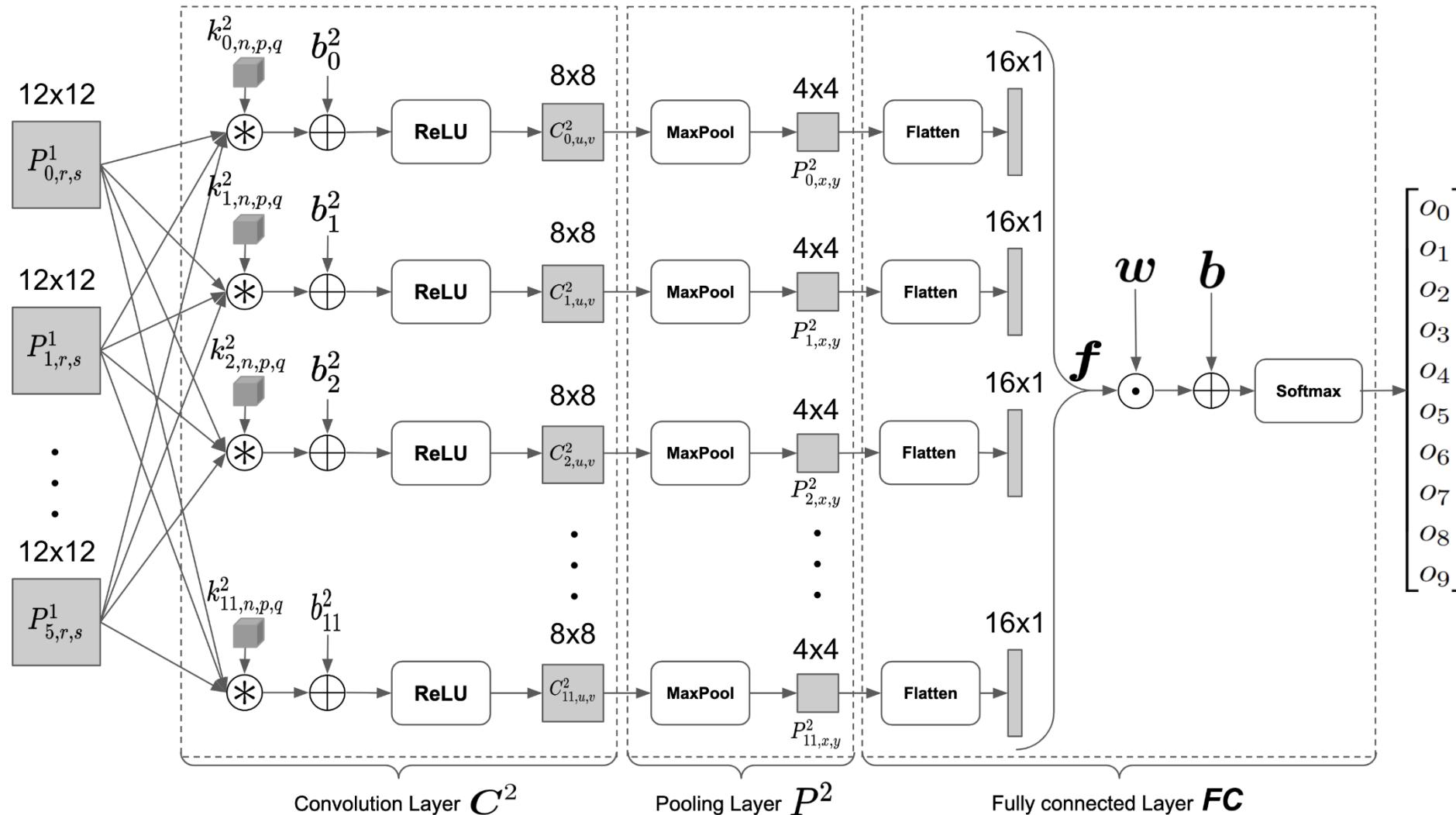
Complex layer terminology



# CNN Design



# CNN Classifier

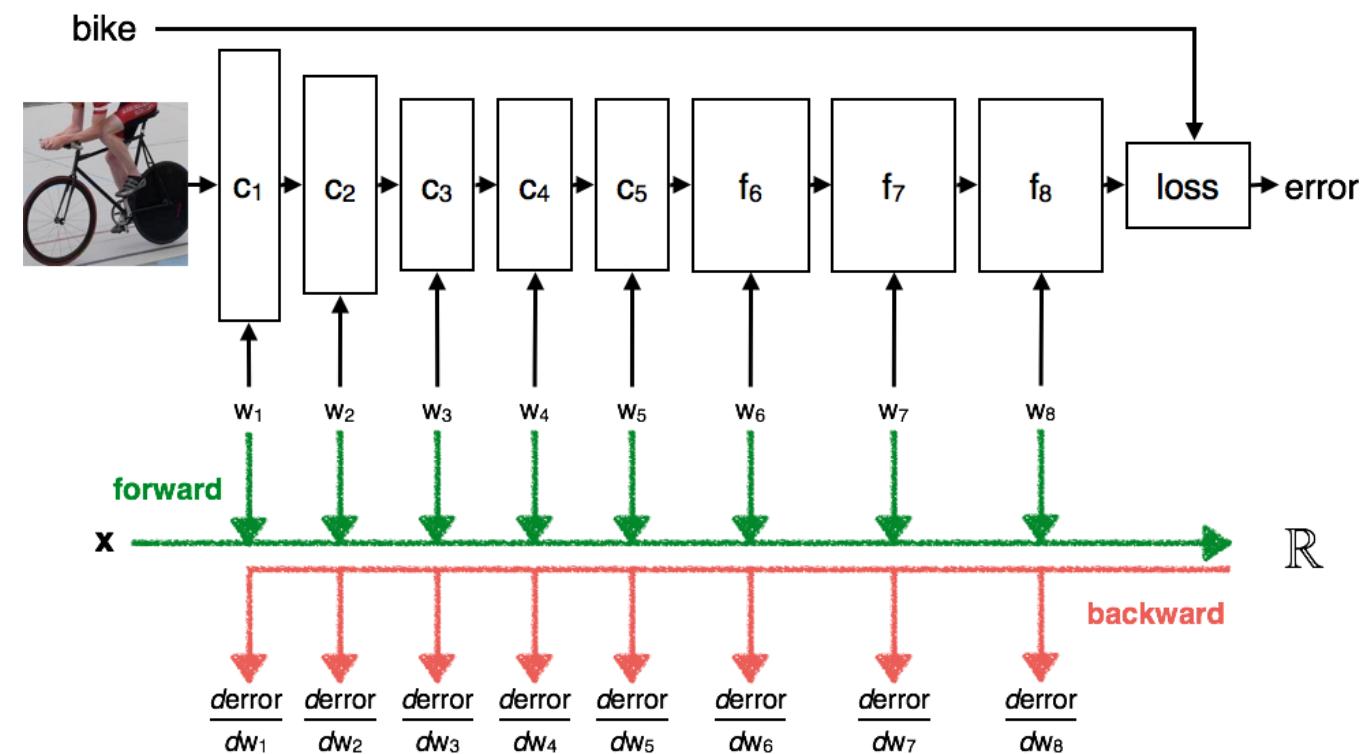


# CNN TRAINING

---

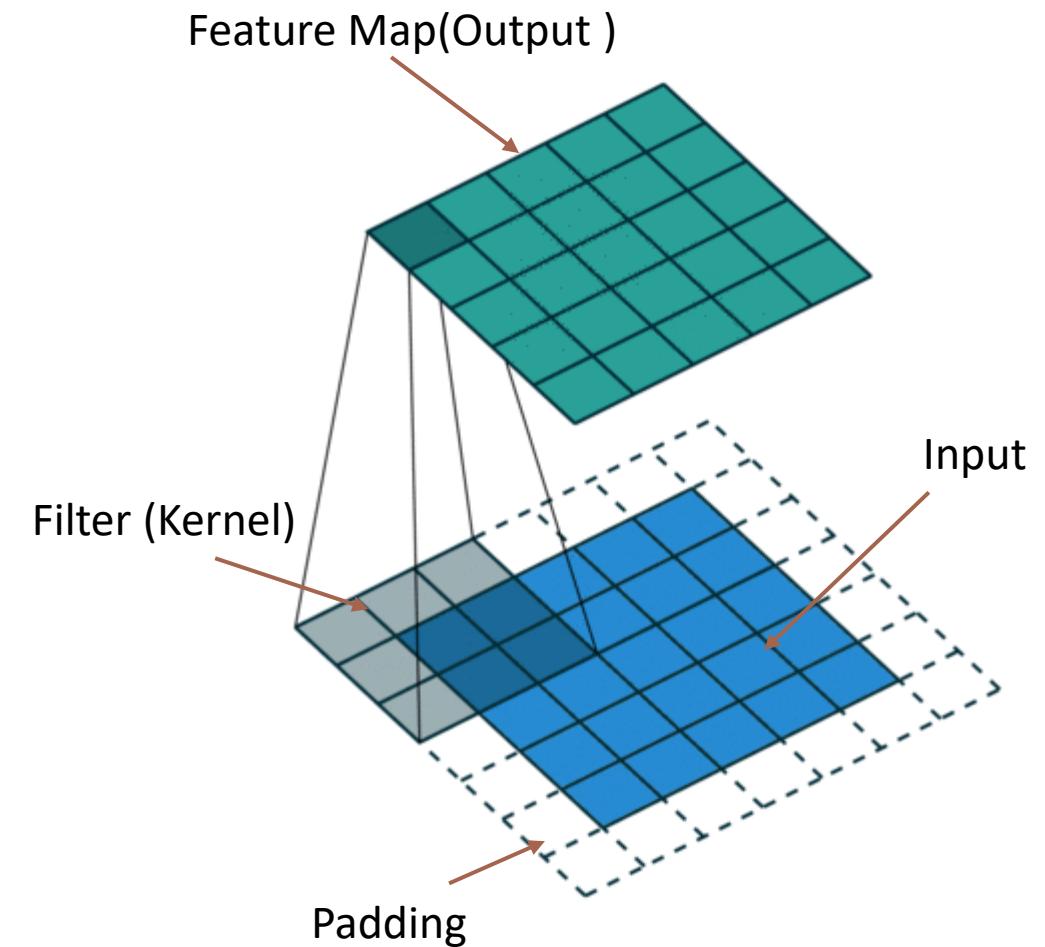
# CNN Training

- CNNs train using Backpropagation
- Weights of the CNN are just the filters
- Filters are not known but learned through backpropagation
  - Instead of fixed numbers in the kernel, assign parameters to kernels that train on the data
  - With training, kernels get better at filtering a given image for relevant information



# Convolution Parameters

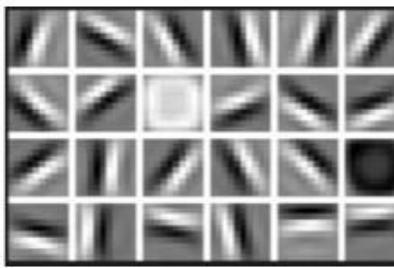
- Filter Size (F)
  - Size of convolution filter
- Padding (P)
  - Allows filter to slide through the image
  - Generates output feature map of desired size.
  - Padding with 0 is most common
- Stride (S)
  - No of steps the filter slides
  - $S = 1$ , 1 pixel at a time
  - $S \geq 2$  produces smaller output volumes



# Convolution Mapping

- The size of the output may shrink and the number of channels may grow as more convolutions are applied to the original input

$$o = \left\lfloor \frac{i + 2p - k}{s} \right\rfloor + 1.$$



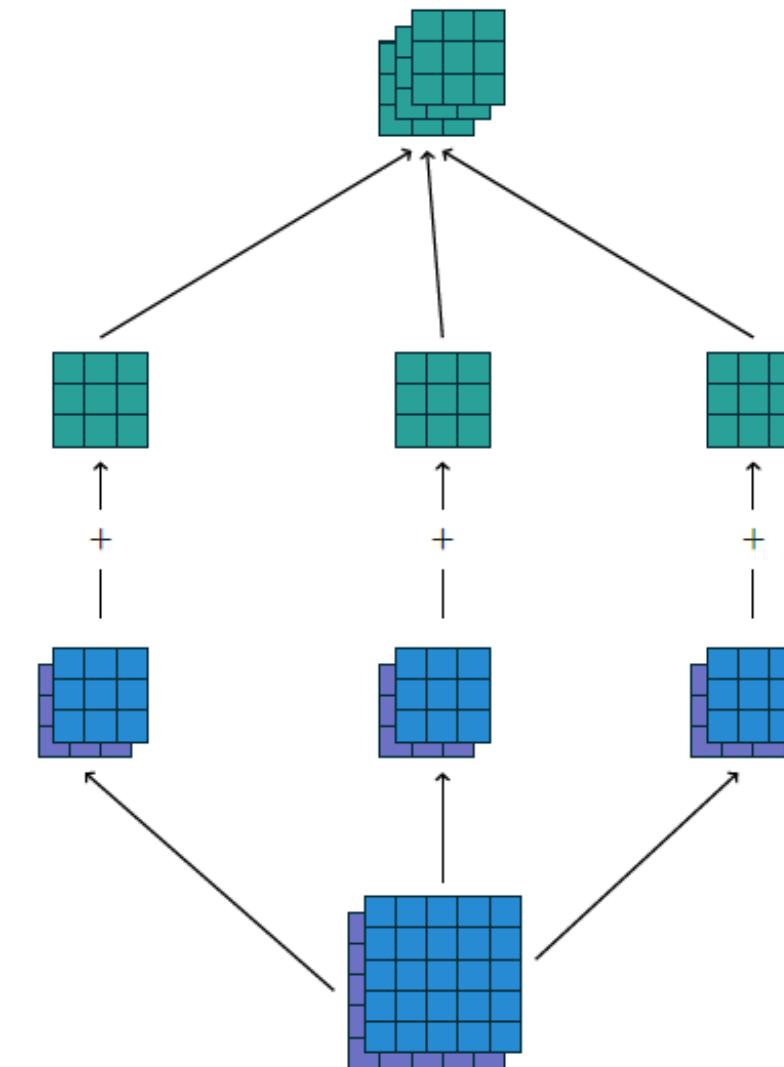
First Layer Representation



Second Layer Representation

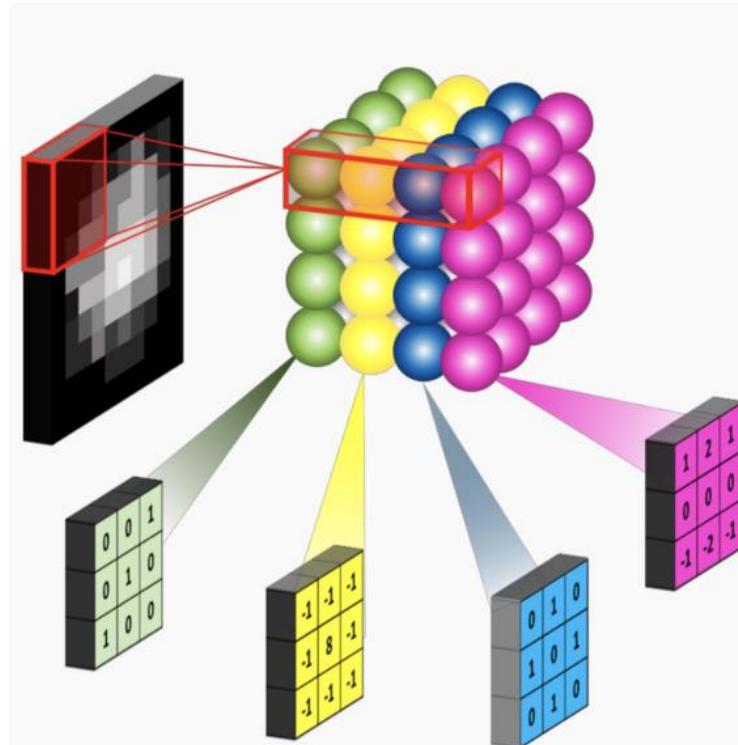


Third Layer Representation

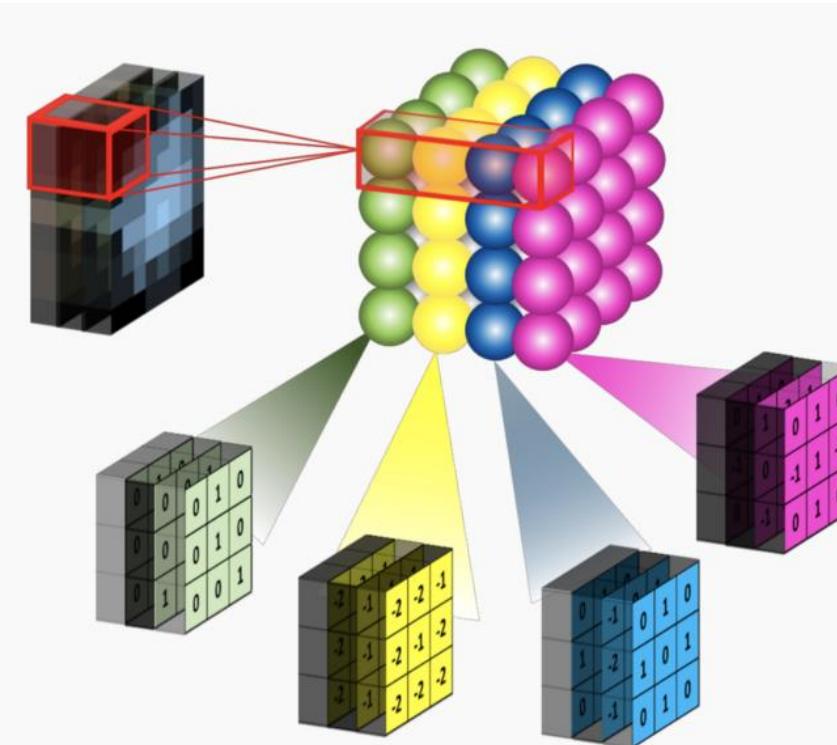


Mapping from two input feature maps to three output feature maps using a  $3 \times 2 \times 3 \times 3$  collection of kernels  $w$

# Multi-channel Convolution



Convolutional layer with four 3x3 filters on a **black and white image** (just one channel)



Convolutional layer with four 3x3 filters on an **RGB image**. As you can see, the filters are now cubes, and they are applied on the full depth of the image..

# Multi-channel Convolution

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	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	...
...	...	...	...	...	...	...

Input Channel #2 (Green)

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

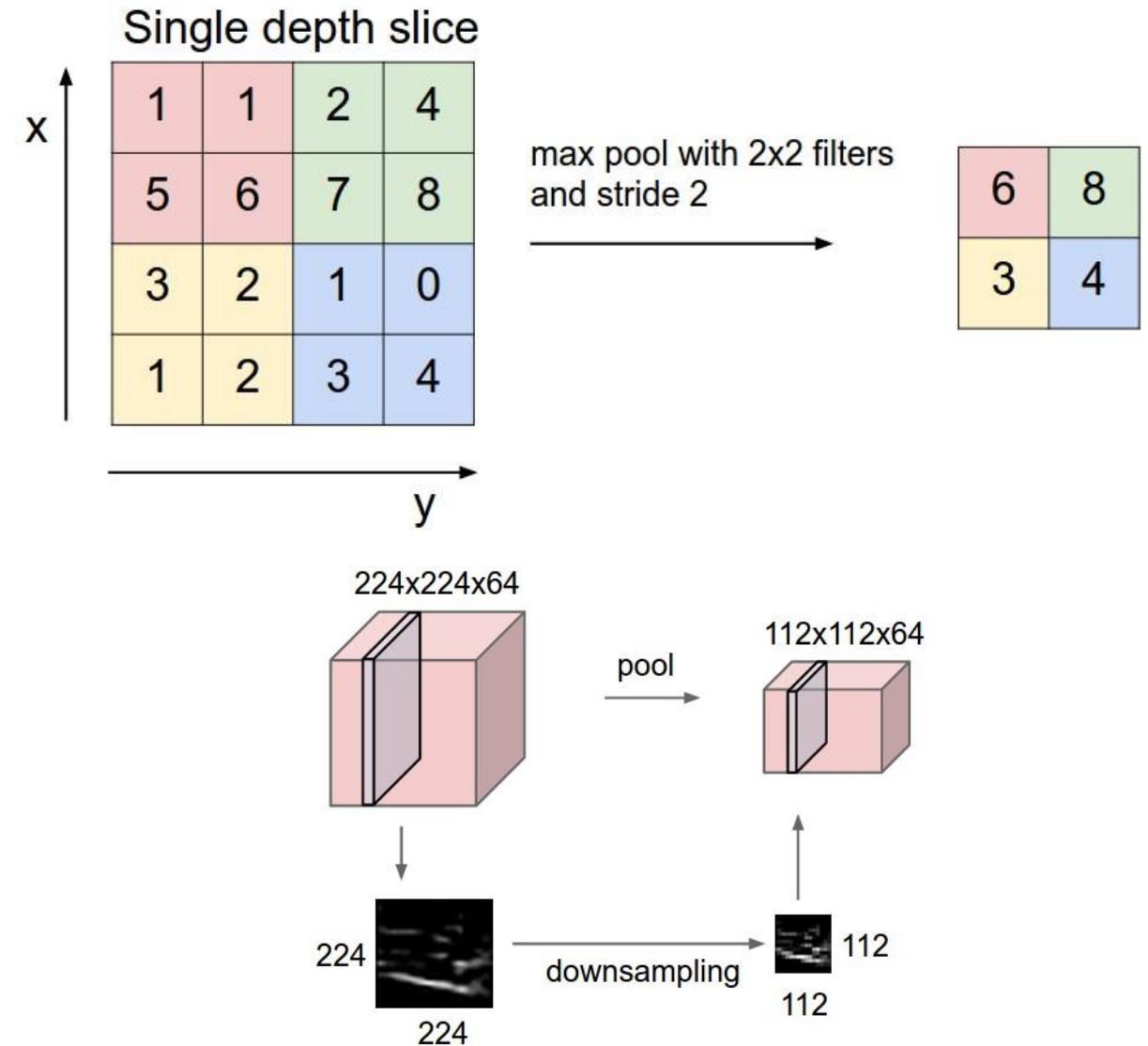
$$+ \quad 164 + 1 = -25$$

$$\begin{array}{c} \uparrow \\ \text{Bias} = 1 \end{array}$$

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

# Pooling

- Pooling is a sub-sampling operation
- Reduces size of feature maps by summarizing sub-regions, such as average or the maximum value (max-pooling)
- Max-pooling reduces the dependency of the feature vectors on their exact placement in an image
- Max-pooling layers do not actually do any learning themselves. Instead, they reduce the size of the problem by introducing sparseness



# Pooling

Pooling operations reduce the size of feature maps by using some function to summarize sub-regions, such as taking the average or the maximum value

$$o = \left\lfloor \frac{i - k}{s} \right\rfloor + 1$$

K should be set as powers of 2 for computational efficiency

F is generally taken as odd number

F=1 might sometimes be used and it makes sense because there is a depth component involved

3	3	2	1	0
0	0	1	3	1
3	1	2	2	3
2	0	0	2	2
2	0	0	0	1

1.7	1.7	1.7
1.0	1.2	1.8
1.1	0.8	1.3

3	3	2	1	0
0	0	1	3	1
3	1	2	2	3
2	0	0	2	2
2	0	0	0	1

1.7	1.7	1.7
1.0	1.2	1.8
1.1	0.8	1.3

3	3	2	1	0
0	0	1	3	1
3	1	2	2	3
2	0	0	2	2
2	0	0	0	1

1.7	1.7	1.7
1.0	1.2	1.8
1.1	0.8	1.3

3	3	2	1	0
0	0	1	3	1
3	1	2	2	3
2	0	0	2	2
2	0	0	0	1

1.7	1.7	1.7
1.0	1.2	1.8
1.1	0.8	1.3

3	3	2	1	0
0	0	1	3	1
3	1	2	2	3
2	0	0	2	2
2	0	0	0	1

1.7	1.7	1.7
1.0	1.2	1.8
1.1	0.8	1.3

3	3	2	1	0
0	0	1	3	1
3	1	2	2	3
2	0	0	2	2
2	0	0	0	1

1.7	1.7	1.7
1.0	1.2	1.8
1.1	0.8	1.3

3	3	2	1	0
0	0	1	3	1
3	1	2	2	3
2	0	0	2	2
2	0	0	0	1

1.7	1.7	1.7
1.0	1.2	1.8
1.1	0.8	1.3

3	3	2	1	0
0	0	1	3	1
3	1	2	2	3
2	0	0	2	2
2	0	0	0	1

1.7	1.7	1.7
1.0	1.2	1.8
1.1	0.8	1.3

3	3	2	1	0
0	0	1	3	1
3	1	2	2	3
2	0	0	2	2
2	0	0	0	1

1.7	1.7	1.7
1.0	1.2	1.8
1.1	0.8	1.3

Computing the output values of a 3x3 average pooling operation on a 5x5 input using 1x1 strides

# Convolution

$$\begin{array}{|c|c|} \hline O_{11} & O_{12} \\ \hline O_{21} & O_{22} \\ \hline \end{array} = \text{Convolution} \left( \begin{array}{|c|c|c|} \hline X_{11} & X_{12} & X_{13} \\ \hline X_{21} & X_{22} & X_{23} \\ \hline X_{31} & X_{32} & X_{33} \\ \hline \end{array}, \begin{array}{|c|c|} \hline F_{11} & F_{12} \\ \hline F_{21} & F_{22} \\ \hline \end{array} \right)$$

Input  $X$       Filter  $F$

$\otimes$

$X_{11}$	$X_{12}$	$X_{13}$
$X_{21}$	$X_{22}$	$X_{23}$
$X_{31}$	$X_{32}$	$X_{33}$

Input  $X$

$F_{11}$	$F_{12}$
$F_{21}$	$F_{22}$

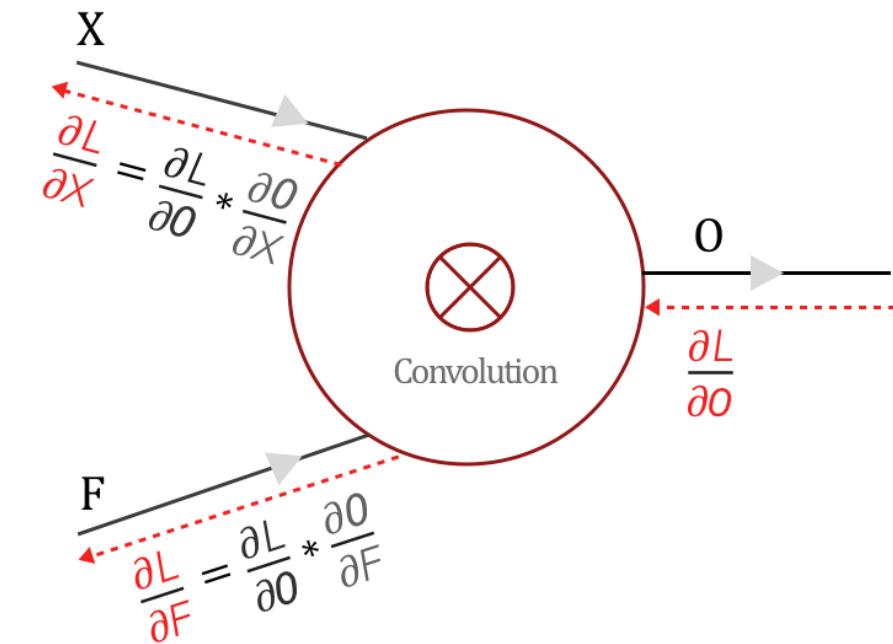
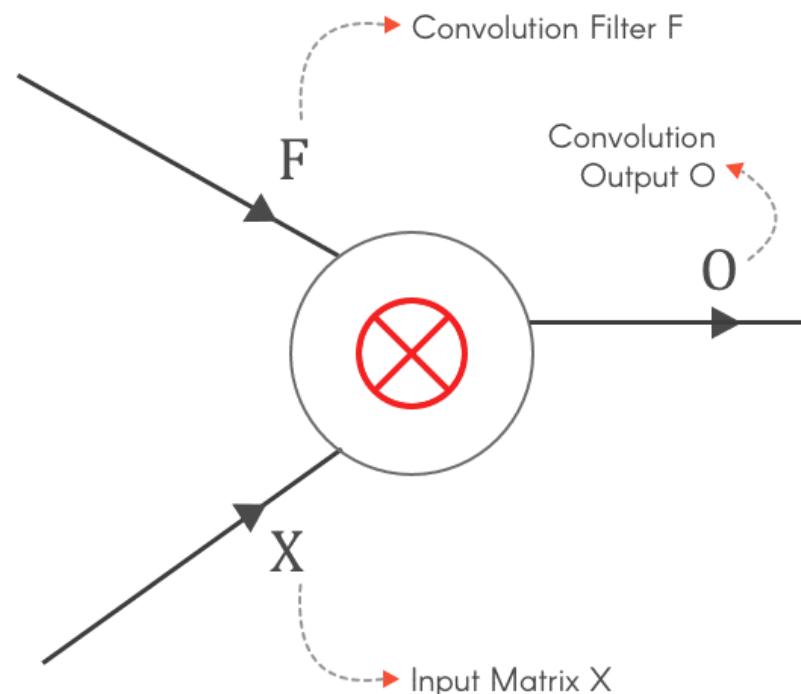
Filter  $F$

$X_{11}F_{11}$	$X_{12}F_{12}$	$X_{13}$
$X_{21}F_{21}$	$X_{22}F_{22}$	$X_{23}$
$X_{31}$	$X_{32}$	$X_{33}$

$$O_{11} = X_{11}F_{11} + X_{12}F_{12} + X_{21}F_{21} + X_{22}F_{22}$$

# CNN Backpropagation

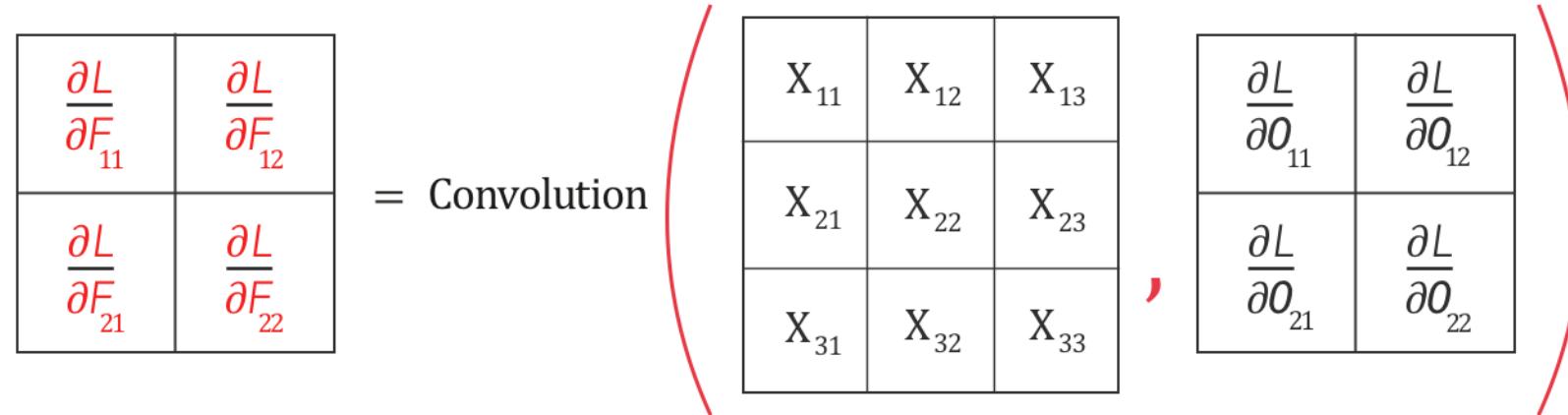
We can find the local gradients  $\partial O / \partial X$  and  $\partial O / \partial F$  with respect to Output  $O$ . And with loss gradient from previous layers —  $\partial L / \partial O$  and using chain rule, we can calculate  $\partial L / \partial X$  and  $\partial L / \partial F$



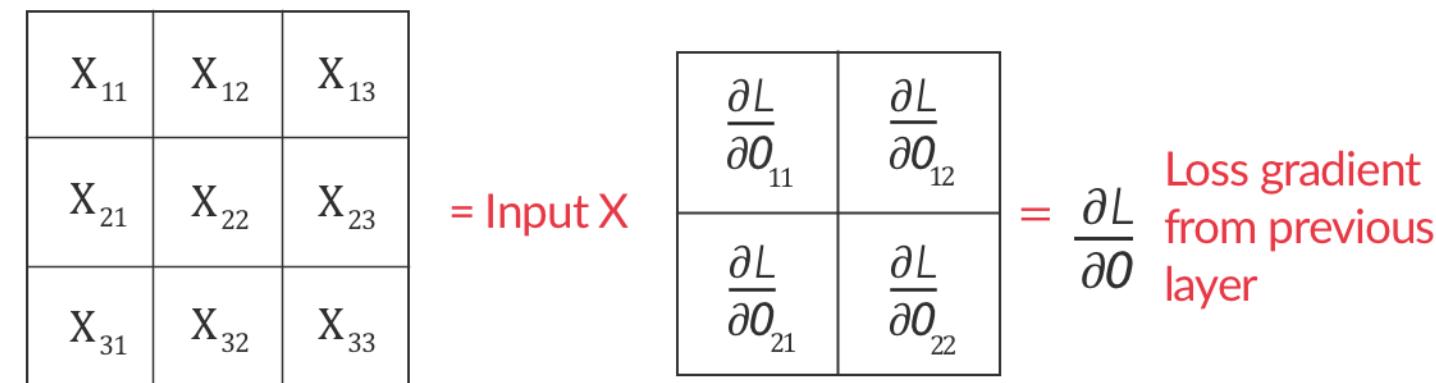
$\frac{\partial O}{\partial X}$  &  $\frac{\partial O}{\partial F}$  are local gradients

$\frac{\partial L}{\partial Z}$  is the loss from the previous layer which has to be backpropagated to other layers

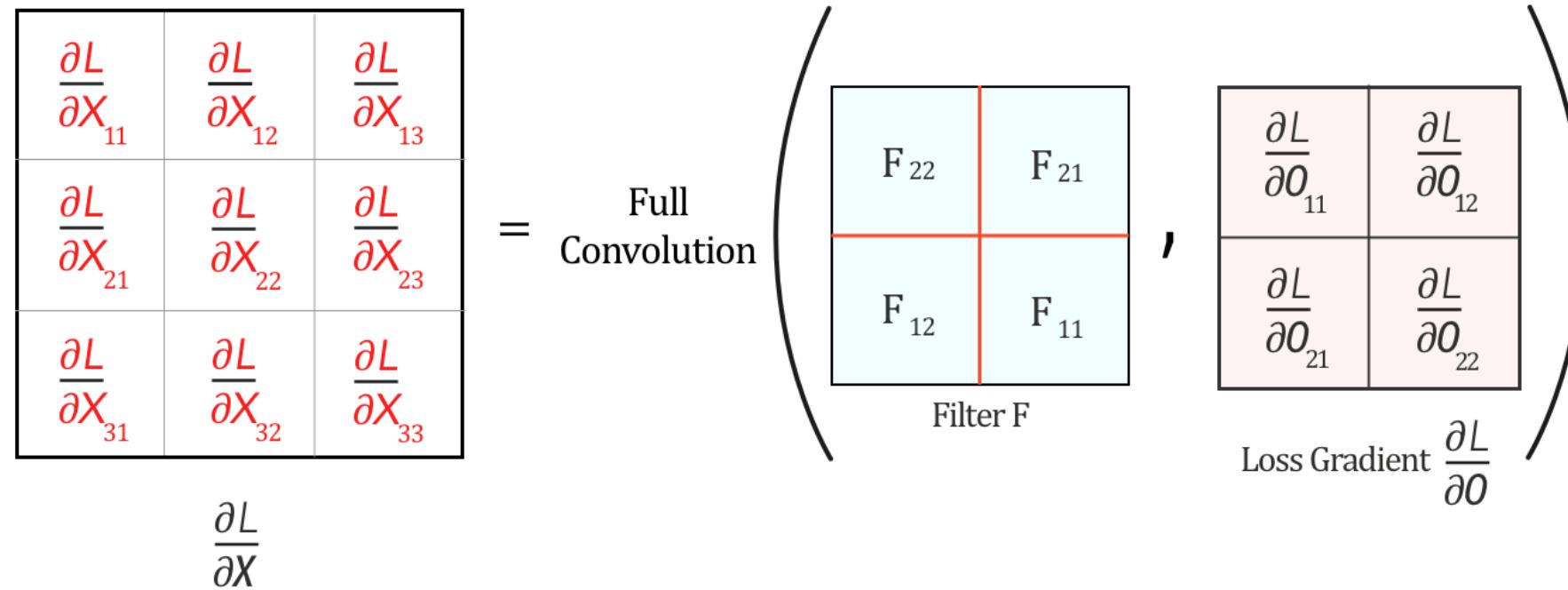
# CNN Backpropagation



where



# CNN Backpropagation



$\frac{\partial L}{\partial X}$  can be represented as 'full' convolution between a 180-degree rotated Filter  $F$  and loss gradient  $\frac{\partial L}{\partial O}$

# CNN Backpropagation

- Both the Forward pass and the Backpropagation of a Convolutional layer are Convolutions

$$\frac{\partial L}{\partial F} = \text{Convolution} \left( \text{Input } X, \text{ Loss gradient } \frac{\partial L}{\partial O} \right)$$

$$\frac{\partial L}{\partial X} = \text{Full Convolution} \left( 180^\circ \text{rotated Filter } F, \text{ Loss Gradient } \frac{\partial L}{\partial O} \right)$$

How to calculate  $\partial L / \partial X$  and  $\partial L / \partial F$

# CNN Architecture Considerations & Best Practices

- Architecture
  - Number of channels generally increases or stays the same while we progress through layers
  - Dropout layer is placed after the Max-Pool layer
  - Batch Normalization layer is placed after the activation function.
  - Avoid negative dimensions due to high values of stride length and filter-sizes
  - Choose the architecture(s) for which you wish to do hyper-parameter optimization
  - Track dimensions, activation shapes, and sizes for the architecture

Layer	Number of Filters	Padding	Activation Shape	Activation Size
Input Image	-	-	(28,28,1)	784
Conv2d(f=3,s=1)	8	Same	(28,28,8)	6,272
MaxPool(f=2,s=2)	-	Valid	(14,14,8)	1568
Conv2d(f=5,s=1)	16	Valid	(10,10,16)	1,600
MaxPool(f=2,s=2)	-	Valid	(5,5,16)	400
Flatten	-	-	(400,1)	400
Flatten	-	-	(120,0)	120
Dense	-	-	(64,1)	64
Softmax	-	-	(10,1)	10

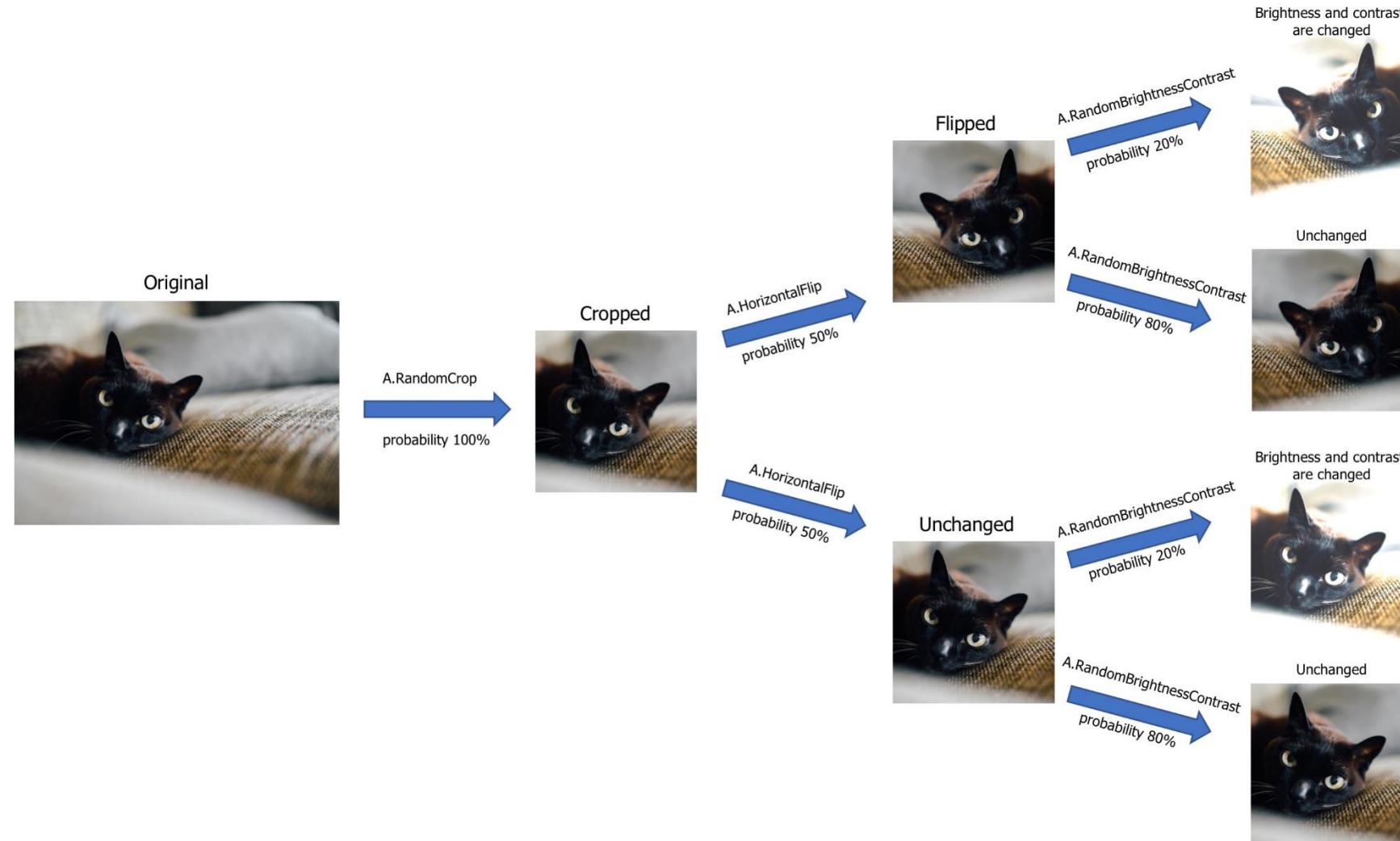
# CNN Training Considerations & Best Practices

- Training Data
  - CNNs require large training samples
  - Perform image augmentation - multiple image orientations and noise
- Vanishing & Exploding Gradients
  - Lowering the learning rate can help keep ReLU units from dying
  - Add Batch Normalization to prevent overfitting and exploding gradients
- Regularization
  - Pick dropout values between 0.0 (no dropout) and 1.0 (dropout everything)
- Training Speed
  - GPU/TPUs
  - BLAS libraries

<https://developers.google.com/machine-learning/crash-course/training-neural-networks/best-practices>

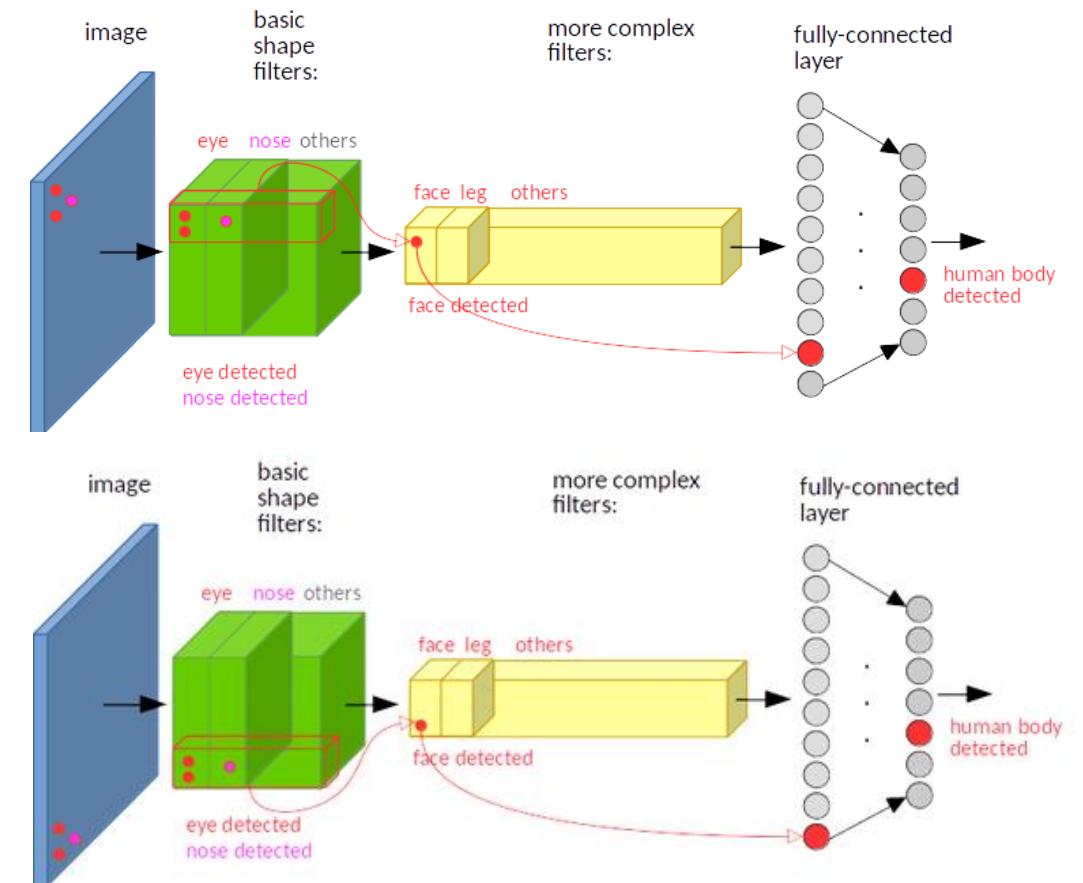
<https://www.mosaicml.com/blog/5-best-practices-for-efficient-model-training>

# Image Augmentation



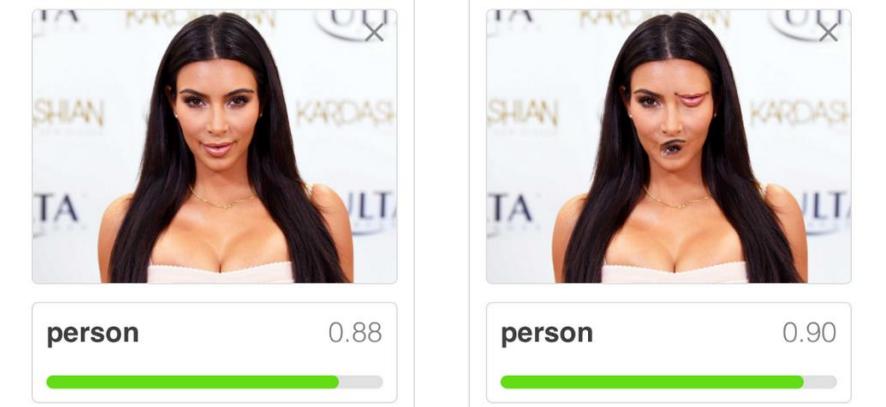
# Invariance

- CNNs are invariant to object translation due to convolution filters
- Regular CNNs are not rotation invariant



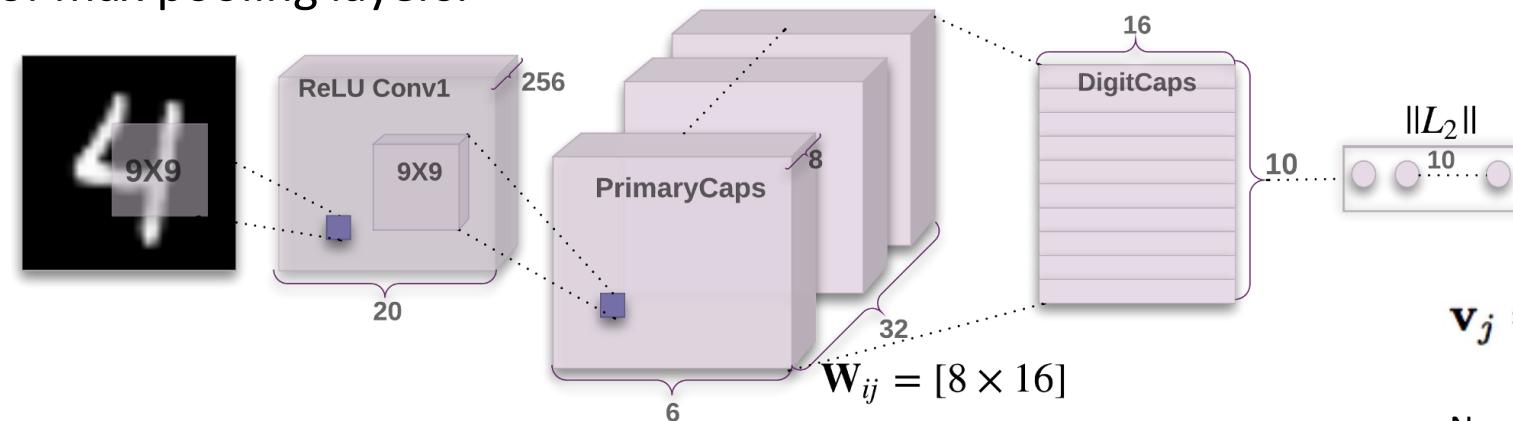
# CNN Issues

- Fooled by images with features in the wrong place
- Easily confused when viewing an image in a different orientation
  - One way to combat this is with excessive training of all possible angles, but this takes a lot of time
- CNNs can be susceptible to white box adversarial attacks
  - Embedding a secret pattern into an object to make it look like something else



# Capsule Networks

- Capsule is a nested set of neural layers - more layers inside a single layer
- Invariance
  - Capsule takes into account the existence of the specific feature like mouth or nose.
- Equivariance
  - Instead of making the feature translation invariant, capsule will make it **translation-equivariant** or viewpoint-equivariant. As the feature moves and changes its position in the image, feature vector representation also changes thus making it equivariant. This property of capsules tries to solve the drawback of max pooling layers.



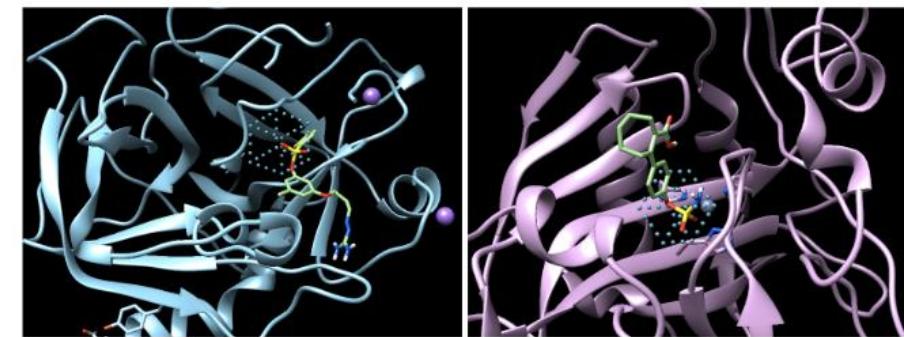
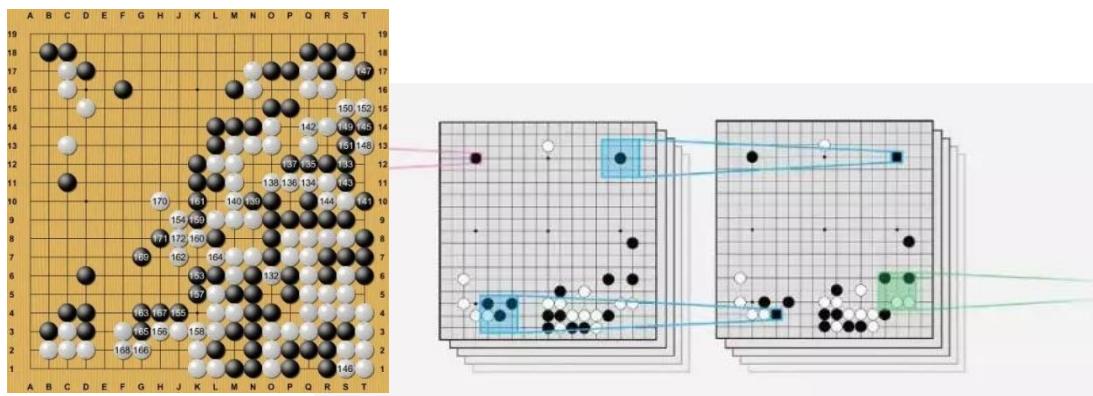
$$\mathbf{v}_j = \frac{\|\mathbf{s}_j\|^2}{1 + \|\mathbf{s}_j\|^2} \frac{\mathbf{s}_j}{\|\mathbf{s}_j\|}$$

Novel Squashing Function

# CNN APPLICATIONS

---

# CNN Applications

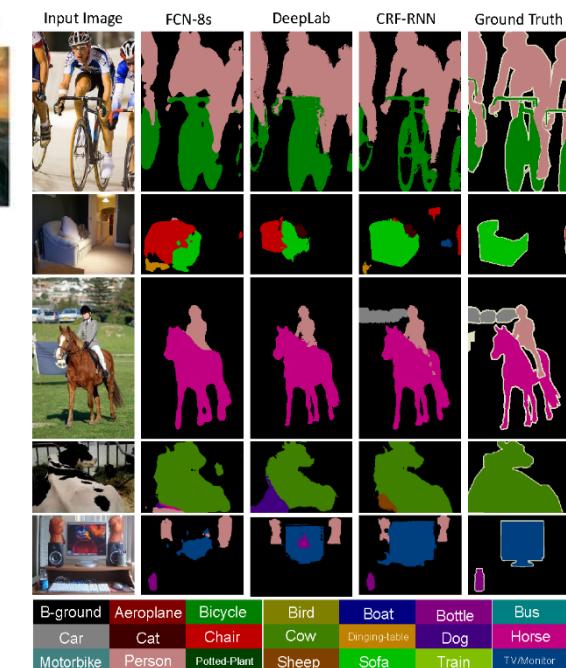


<https://arxiv.org/pdf/1510.02855.pdf>

Figure 5: Sulfonyl/sulfonamide detection with autonomously trained convolutional filters.



<https://github.com/jcjohnson/neural-style/>



<http://www.robots.ox.ac.uk/~szheng/CRFasRNN.html>

# Yoga Pose Detection

- 10 Yoga poses detected
- ~500+ images per pose
- Model: CNN Resnet-50 with Transfer Learning
- Tensorflow on iPhone



<https://grahamschool.uchicago.edu/academic-programs/masters-degrees/analytics/yoga>

# Lab

1. MNIST MLP
2. MNIST CNN - Keras
3. CIFAR CNN - Keras
4. <https://transcranial.github.io/keras-js/#/>

