

---

# CAP 5516

# Medical Image Computing

# (Spring 2025)

Dr. Chen Chen  
Associate Professor

Center for Research in Computer Vision (CRCV)  
University of Central Florida  
Office: HEC 221

Email: [chen.chen@crcv.ucf.edu](mailto:chen.chen@crcv.ucf.edu)

Web: <https://www.crcv.ucf.edu/chenchen/>

---

# Lecture 6: Introduction to Deep Learning (1)

---

## ARTIFICIAL INTELLIGENCE

A program that can sense, reason,  
act, and adapt

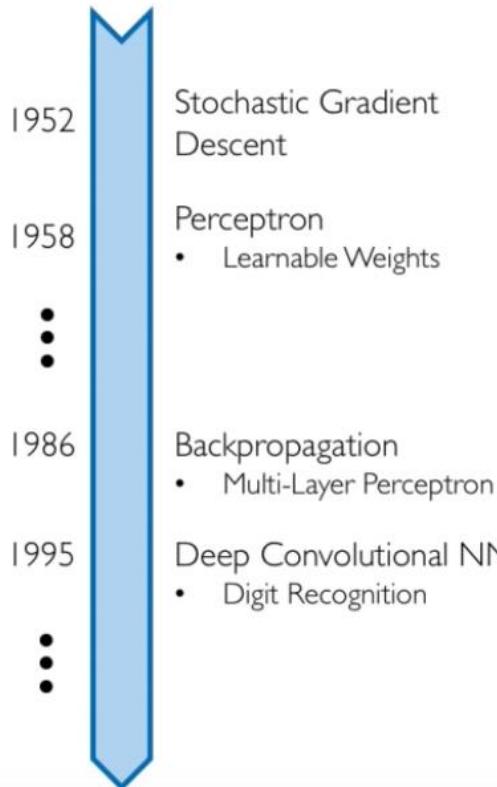
## MACHINE LEARNING

Algorithms whose performance improve  
as they are exposed to more data over time

## DEEP LEARNING

Subset of machine learning in  
which multilayered neural  
networks learn from  
vast amounts of data

# Why Now?



Neural Networks date back decades, so why the resurgence?

## I. Big Data

- Larger Datasets
- Easier Collection & Storage

IM<sup>2</sup>GENET



## 2. Hardware

- Graphics Processing Units (GPUs)
- Massively Parallelizable



## 3. Software

- Improved Techniques
- New Models
- Toolboxes



PYTORCH

Caffe2

Chainer

mxnet

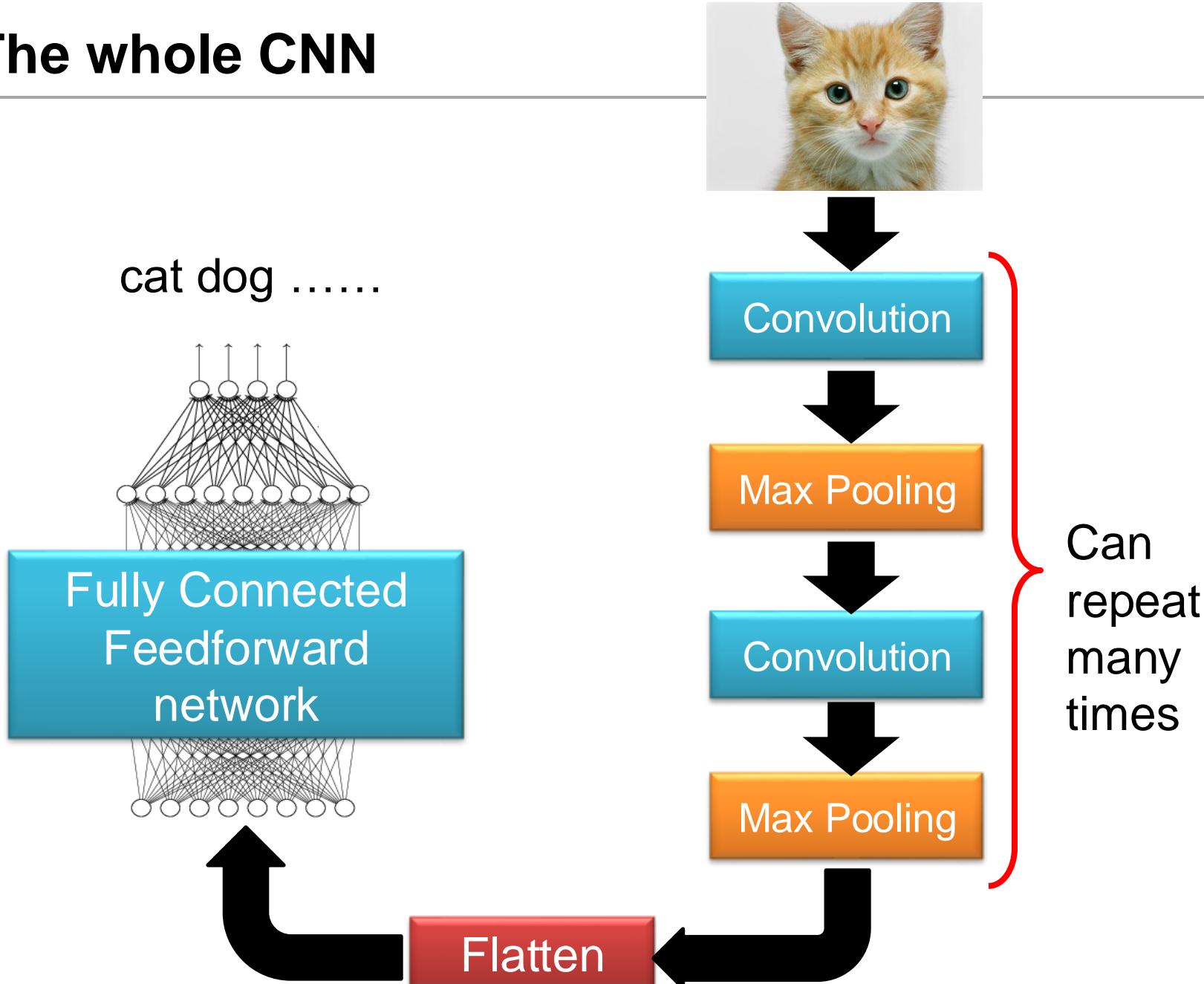
## 4. Open resource (GitHub)

Credit: Lex Fridman

# Convolutional Neural Network (CNN)

Widely used in image  
processing and computer vision

# The whole CNN



# The whole CNN



## Property 1

- Some patterns are much smaller than the whole image

## Property 2

- The same patterns appear in different regions.

## Property 3

- Subsampling the pixels will not change the object

Convolution

Max Pooling

Convolution

Max Pooling

Can  
repeat  
many  
times

Flatten



# CNN – Convolution

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

6 x 6 image

Those are the network parameters to be learned.

1	-1	-1
-1	1	-1
-1	-1	1

Filter/kernel 1  
Matrix

-1	1	-1
-1	1	-1
-1	1	-1

Filter 2  
Matrix

: :

Each filter detects a small pattern (3 x 3).

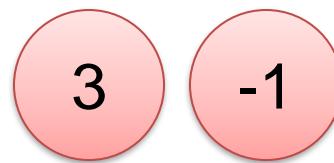
# CNN – Convolution

1	-1	-1
-1	1	-1
-1	-1	1

Filter/Kernel 1

stride=1

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



6 x 6 image

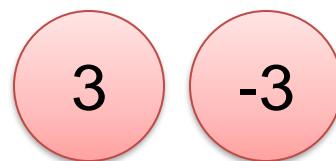
# CNN – Convolution

1	-1	-1
-1	1	-1
-1	-1	1

Filter 1

If stride=2

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

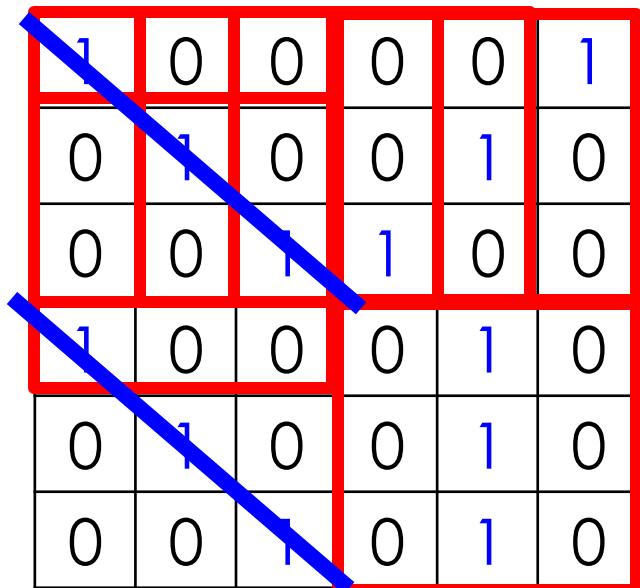


We set stride=1 below

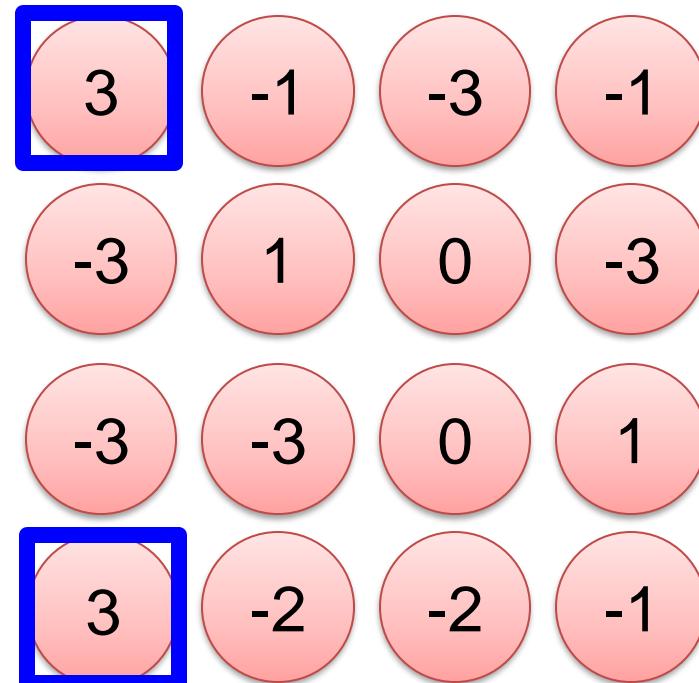
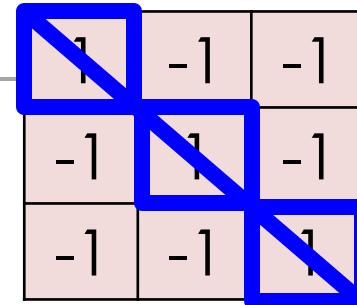
6 x 6 image

# CNN – Convolution

stride=1



6 x 6 image



# CNN – Convolution

-1	1	-1
-1	1	-1
-1	1	-1

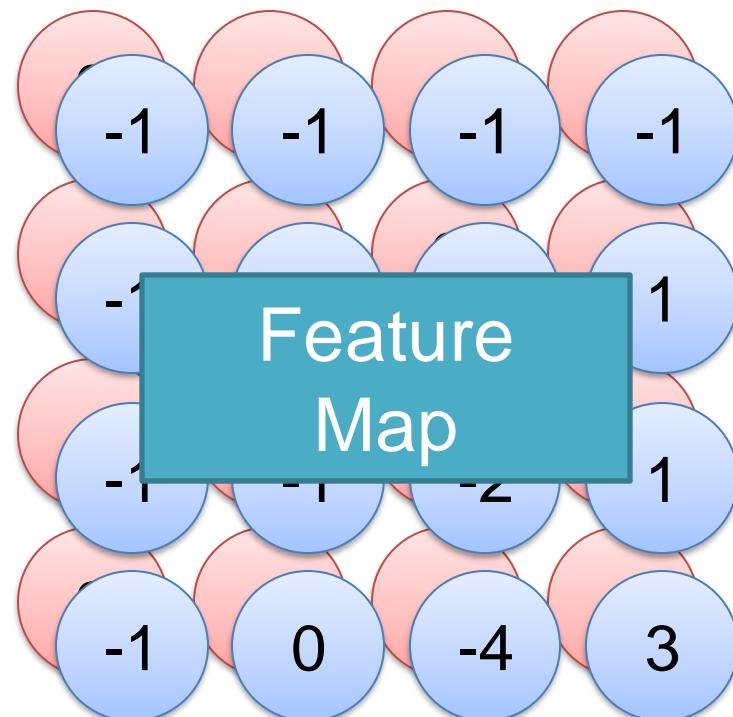
Filter 2

stride=1

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

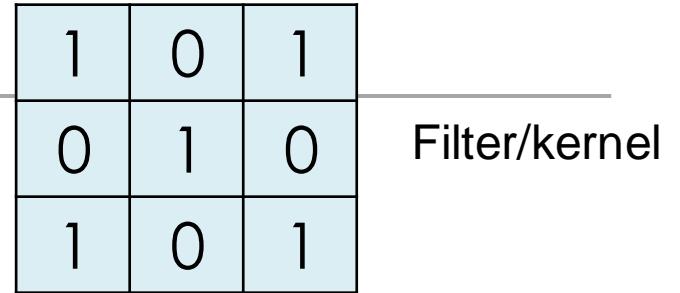
6 x 6 image

Do the same process  
for every filter



4 x 4 image (reduced size)

# Demo



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

Image

4		

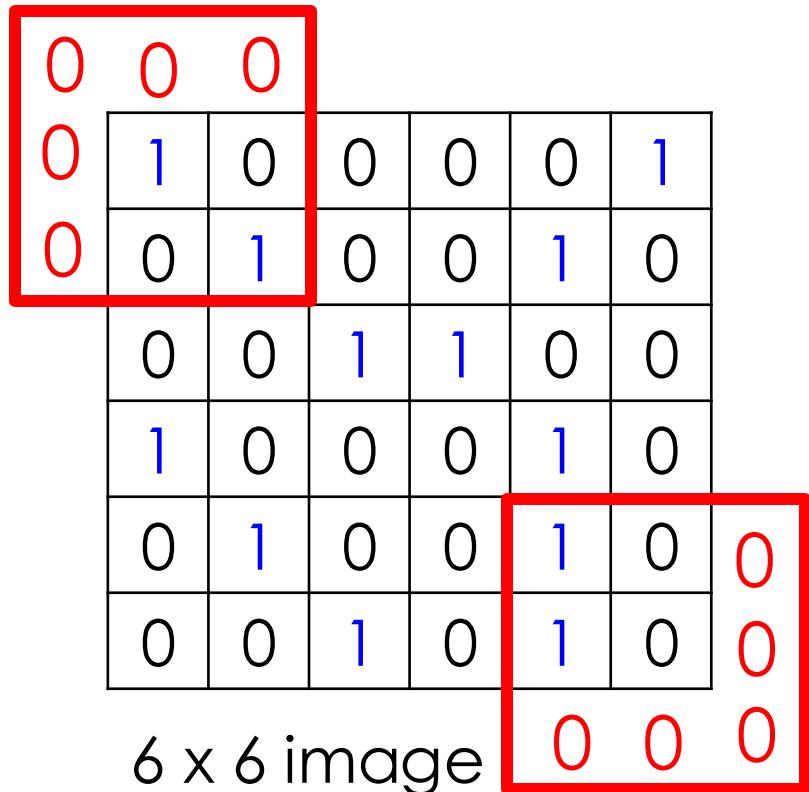
Convolved  
Feature

<https://ujjwalkarn.me/2016/08/11/intuitive-explanation-convnets/>

# CNN – Zero Padding

1	-1	-1
-1	1	-1
-1	-1	1

Filter 1

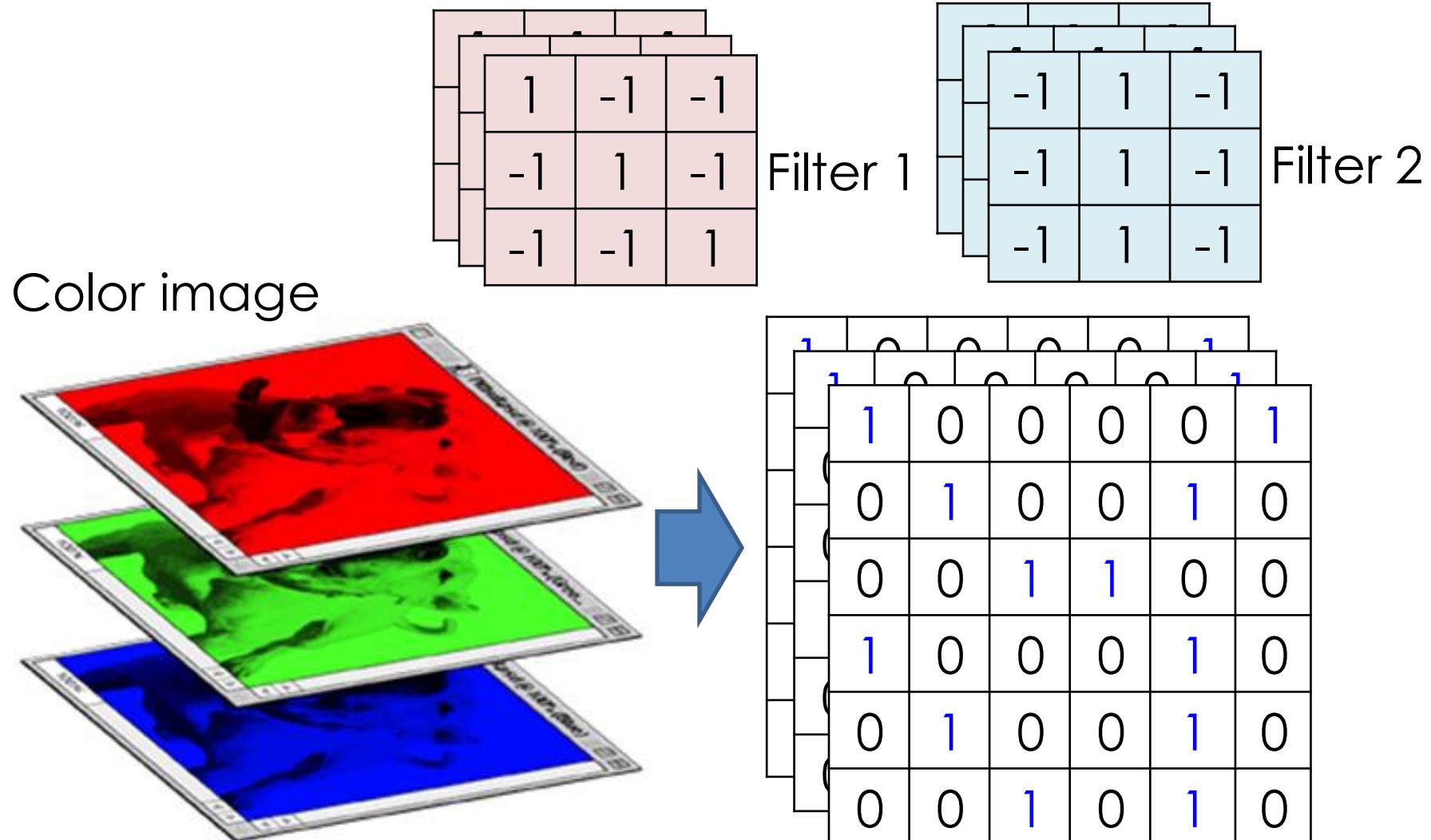


You will get another  $6 \times 6$  image in this way



Zero  
padding

# CNN – Color image



Input Volume (+pad 1) (7x7x3)

 $x[:, :, 0]$ 

0	0	0	0	0	0	0	0
0	2	0	0	2	1	0	
0	2	1	0	2	2	0	

0	2	1	0	2	0	0
0	2	2	1	2	2	0
0	0	1	1	0	1	0
0	0	0	0	0	0	0

0	0	0	0	0	0	0	0
0	1	1	2	2	0	0	
0	1	2	1	0	1	0	

0	0	2	0	2	2	0
0	0	2	0	0	1	0
0	2	1	2	2	0	0

0	0	0	0	0	0	0	0
0	0	1	2	1	0	0	
0	0	1	0	2	2	0	
0	0	0	0	0	0	0	0

Filter W0 (3x3x3)

 $w0[:, :, 0]$ 

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

 $w0[:, :, 1]$ 

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

 $w0[:, :, 2]$ 

0	1	1
1	1	0
1	-1	0

Bias  $b_0$  (1x1x1)  
 $b_0[:, :, 0]$ 

1

Filter W1 (3x3x3)

 $w1[:, :, 0]$ 

0	-1	0
1	1	1
0	1	1

 $w1[:, :, 1]$ 

-1	-1	-1
-1	1	-1
-1	1	-1

 $w1[:, :, 2]$ 

-1	-1	1
-1	-1	1
1	-1	-1

Bias  $b_1$  (1x1x1)  
 $b_1[:, :, 0]$ 

0

Output Volume (3x3x2)

 $o[:, :, 0]$ 

-2	2	-1
3	6	-3
4	1	-2

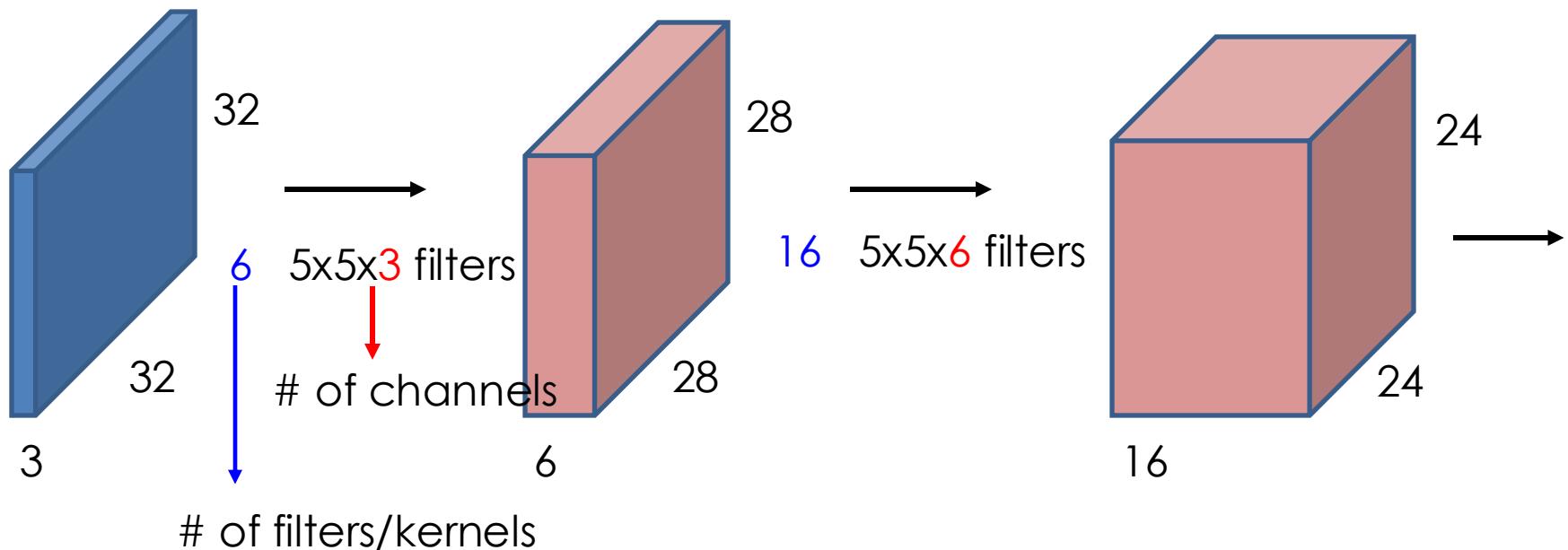
 $o[:, :, 1]$ 

4	-1	3
-4	-4	-3
-2	-3	-8

toggle movement

# Convolutional Network

- Convolution network is a sequence of these layers



# Convolution - Intuition

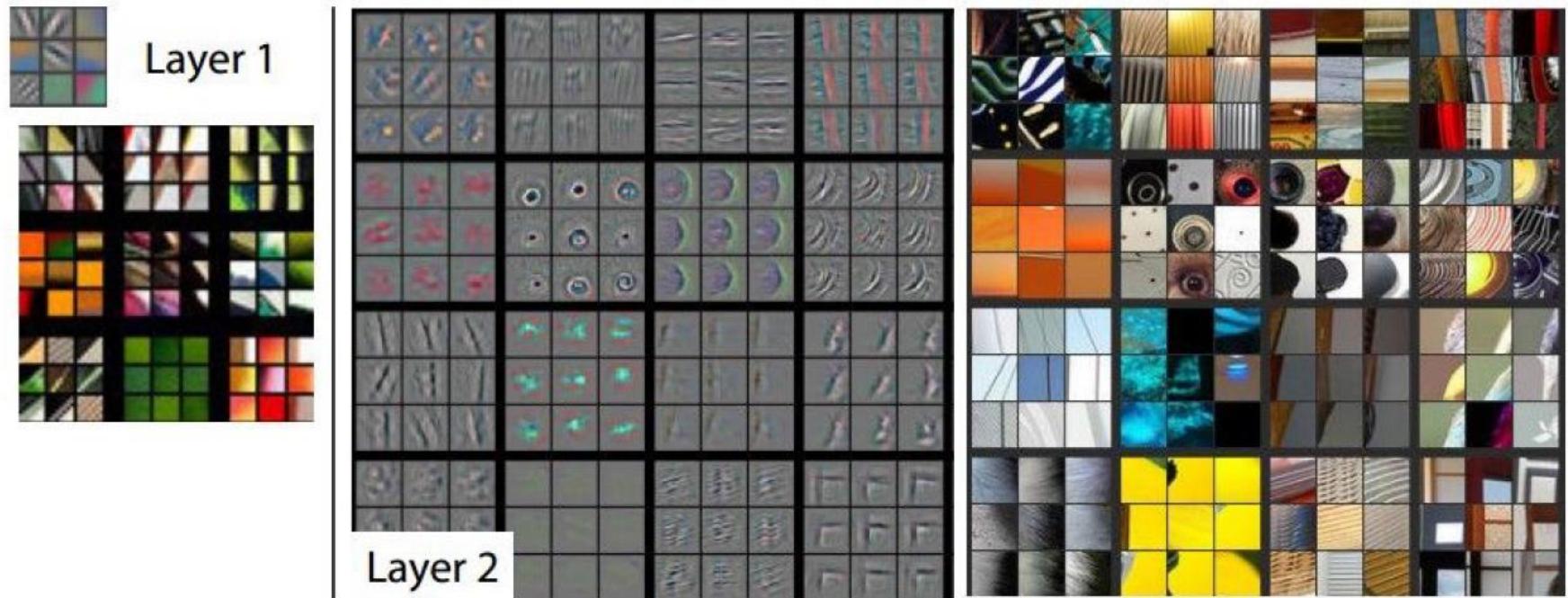
---



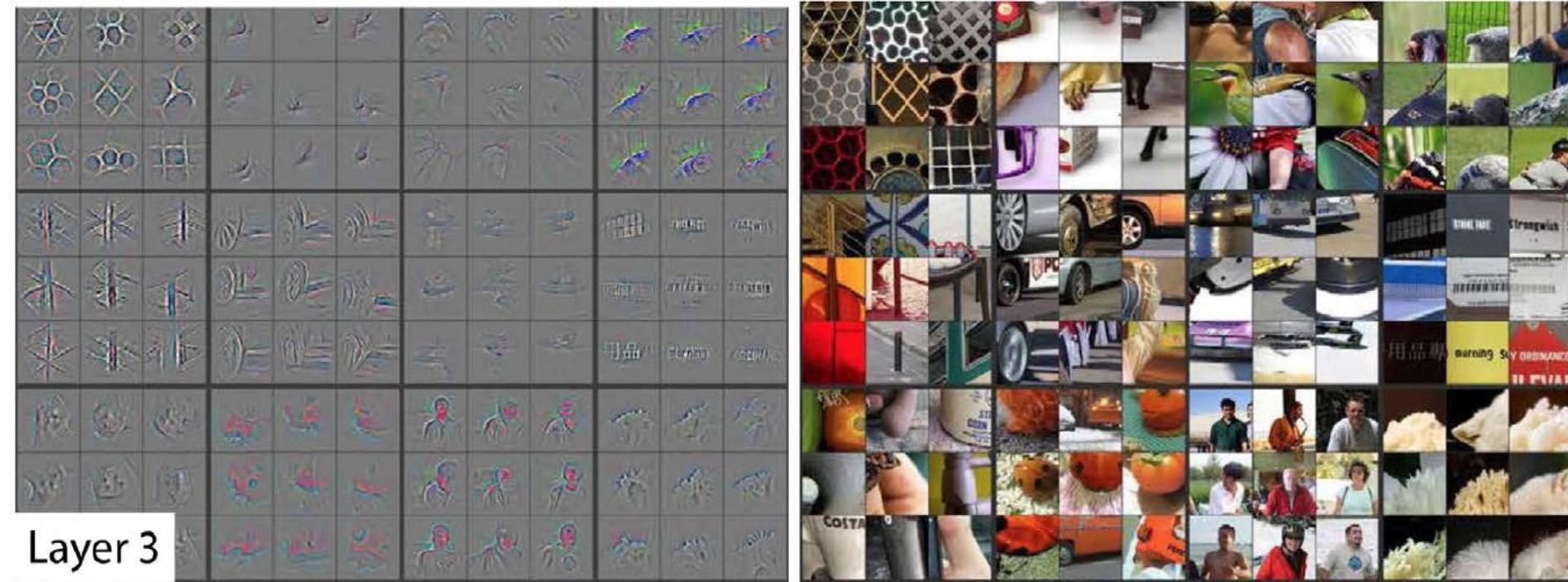
Source : [https://cs.nyu.edu/~fergus/tutorials/deep\\_learning\\_cvpr12/](https://cs.nyu.edu/~fergus/tutorials/deep_learning_cvpr12/)

# Visualizing CNN

---



Zeiler, Matthew D., and Rob Fergus. "Visualizing and understanding convolutional networks." *European conference on computer vision*. Springer, Cham, 2014.



Zeiler, Matthew D., and Rob Fergus. "Visualizing and understanding convolutional networks." European conference on computer vision. Springer, Cham, 2014.



Layer 4



Layer 5



# Variants of Convolution Operation

---

- Dilated convolution [1]
- Depth-wise separable convolution [2]
- Grouped convolution [3]

[1] Yu, Fisher, and Vladlen Koltun. "Multi-scale context aggregation by dilated convolutions." arXiv preprint arXiv:1511.07122 (2015).

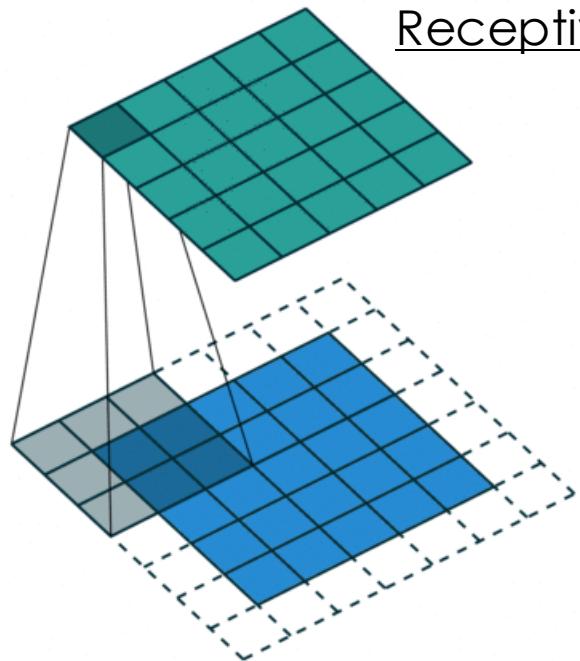
[2] Howard, Andrew G., et al. "Mobilenets: Efficient convolutional neural networks for mobile vision applications." arXiv preprint arXiv:1704.04861 (2017).

[3] Zhang, Xiangyu, et al. "Shufflenet: An extremely efficient convolutional neural network for mobile devices." Proceedings of the IEEE conference on computer vision and pattern recognition. 2018.

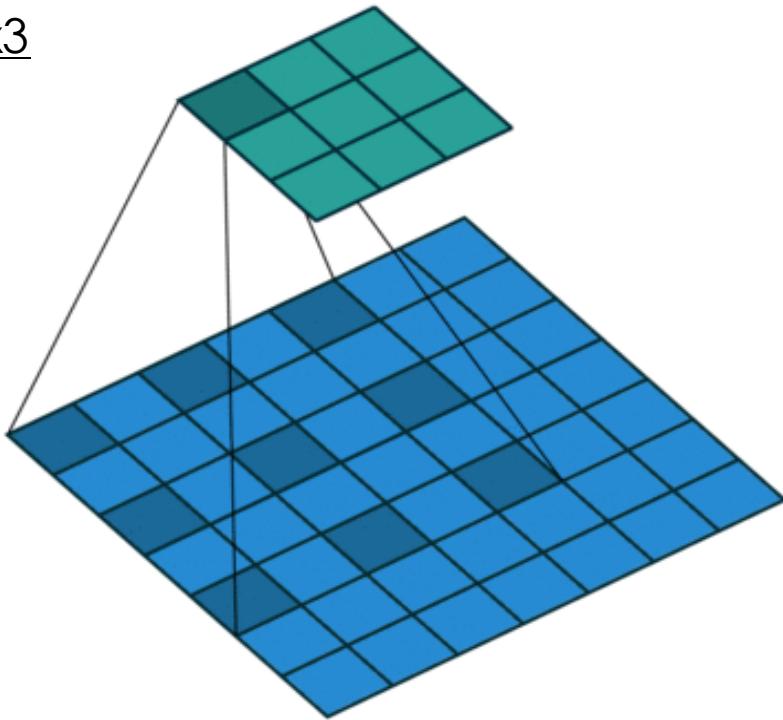
# Variants of Convolution Operation

- Dilated/Atrous Convolution

Dilation rate = 2  
Receptive field =  $5 \times 5$   
With only 9 parameters



Normal convolution

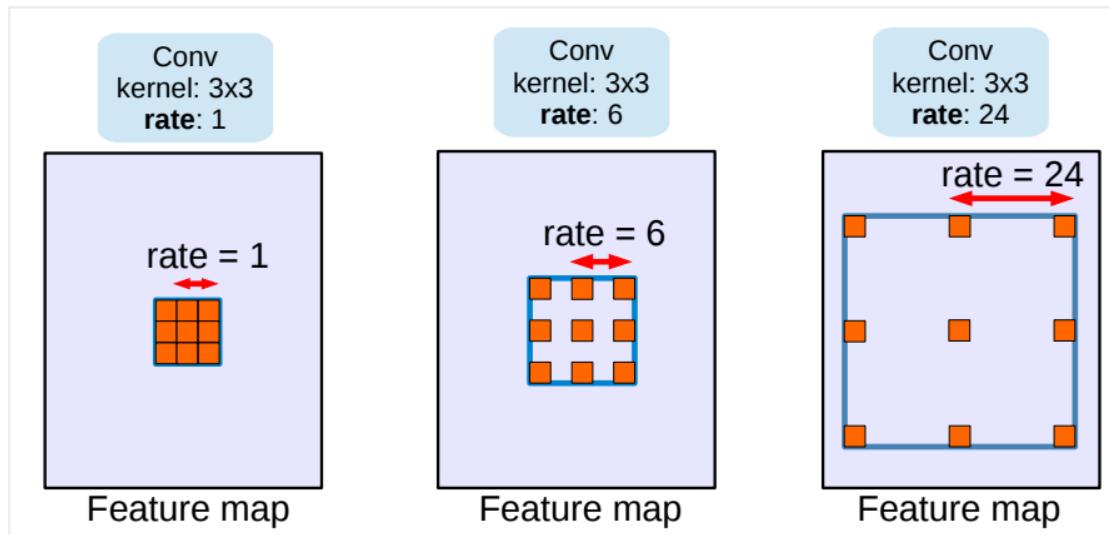


Effective in segmentation task

<https://towardsdatascience.com/types-of-convolutions-in-deep-learning-717013397f4d>

# Variants of Convolution Operation

- Dilated Convolution



- Large value of atrous/dilation rate enlarges the model's field-of-view
- Enabling object encoding at multiple scales

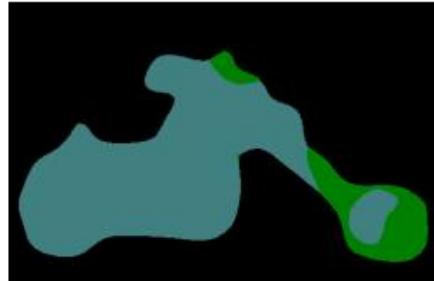
Chen, Liang-Chieh, et al. "Rethinking atrous convolution for semantic image segmentation." *arXiv preprint arXiv:1706.05587* (2017).



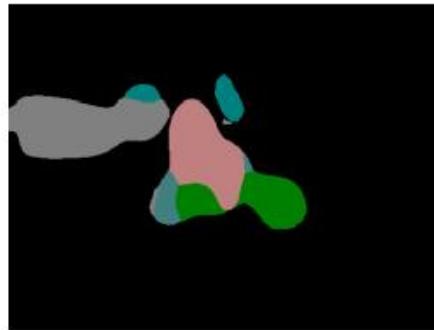
# Variants of Convolution Operation



Input



LargeFOV

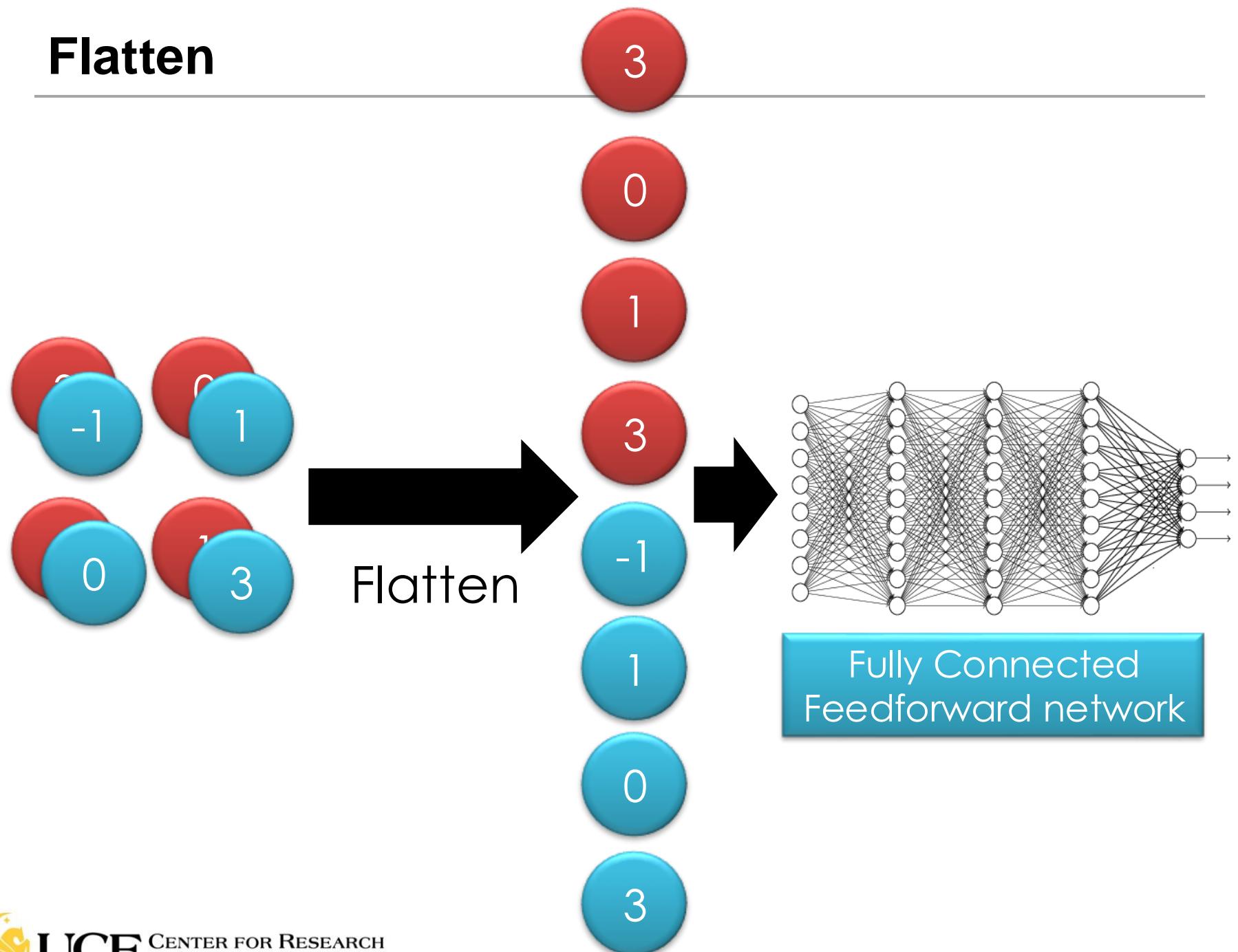


Atrous Spatial Pyramid Pooling (ASPP)

Chen, Liang-Chieh, et al. "Deeplab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected crfs." *IEEE transactions on pattern analysis and machine intelligence* 40.4 (2018): 834-848.

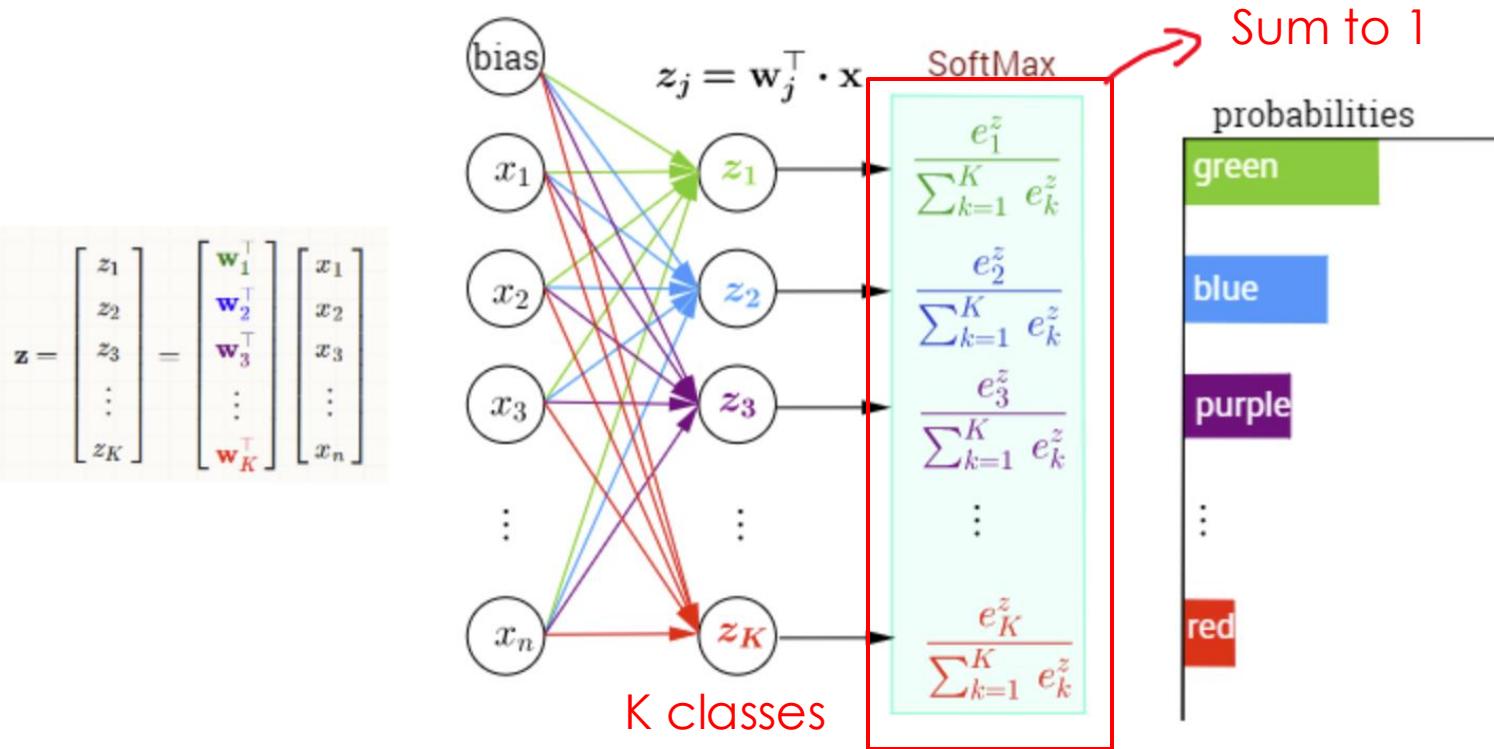
# Flatten

---



# Softmax

## Multi-Class Classification with NN and SoftMax Function



The softmax as

$$\sigma(j) = \frac{\exp(\mathbf{w}_j^\top \mathbf{x})}{\sum_{k=1}^K \exp(\mathbf{w}_k^\top \mathbf{x})} = \frac{\exp(z_j)}{\sum_{k=1}^K \exp(z_k)}$$

# Loss Function

---

- Way to define how good the network is performing
  - In terms of prediction
- Network training (Optimization)
  - Find the best network parameters to minimize the loss

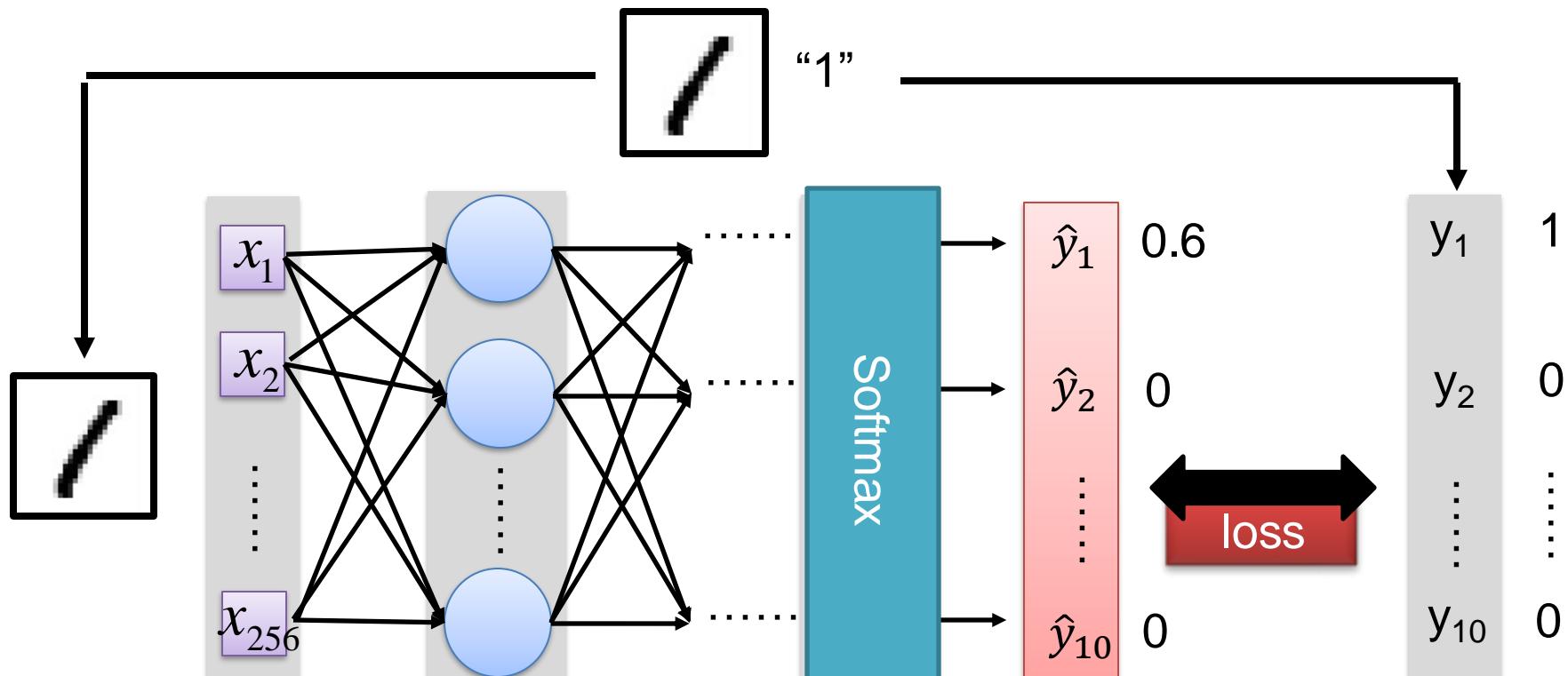
$$L(W) = \frac{1}{N} \sum_{i=1}^N L_i(f(\mathbf{x}_i, W), \mathbf{y}_i)$$

Total loss for a training set  
N samples

The diagram shows the loss function equation with red arrows pointing from text labels to specific parts of the equation:

- An arrow points from "input" to the variable  $\mathbf{x}_i$ .
- An arrow points from "Ground truth" to the variable  $\mathbf{y}_i$ .
- An arrow points from "Network parameters" to the variable  $W$ .
- An arrow points from "network" to the term  $f(\mathbf{x}_i, W)$ .
- An arrow points from "Loss function" to the summation part  $L_i(f(\mathbf{x}_i, W), \mathbf{y}_i)$ .
- An arrow points from "N samples" to the denominator  $N$ .
- An arrow points from "Total loss for a training set" to the overall expression  $L(W)$ .

# Choosing Proper Loss

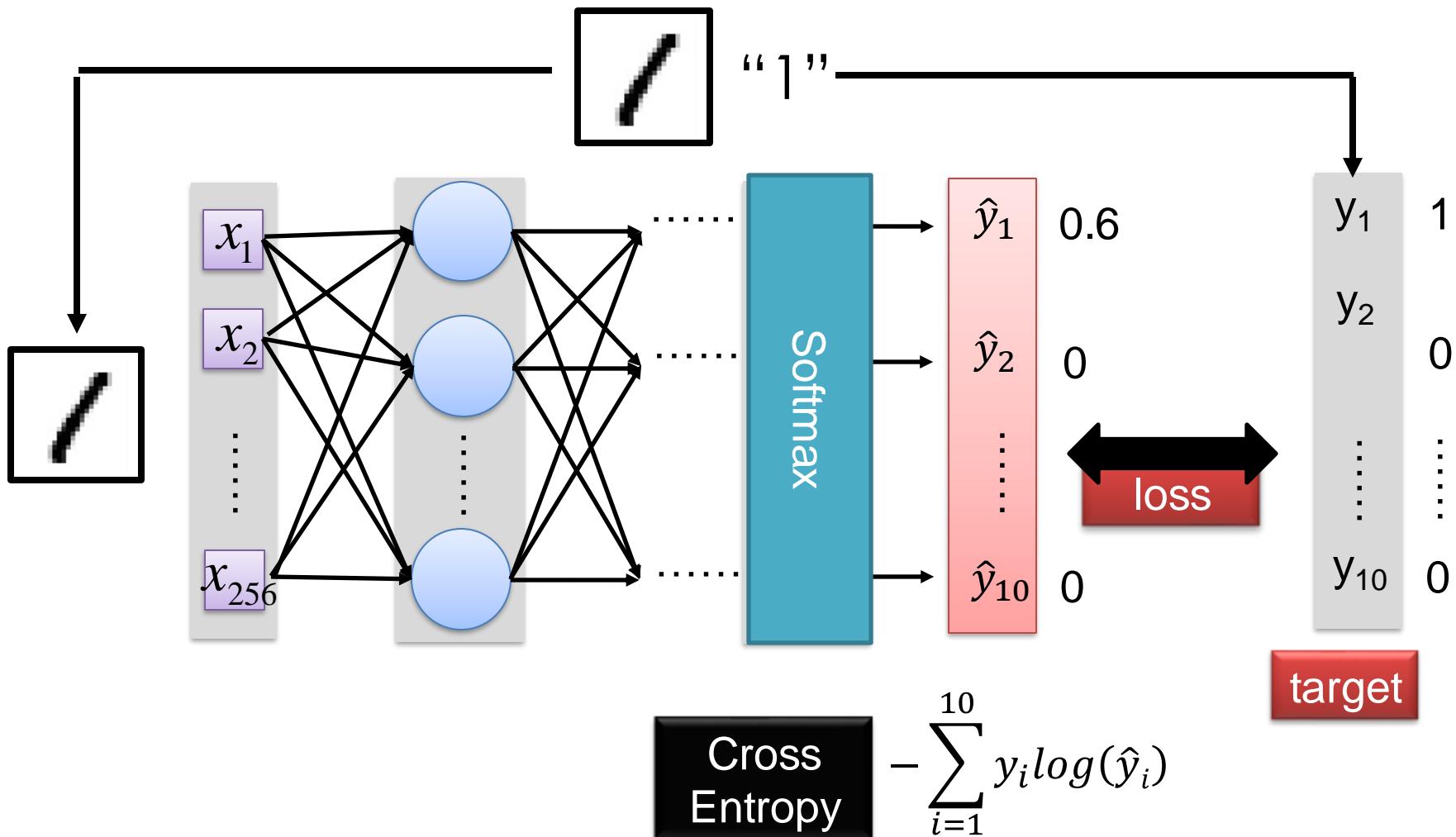


Square  
Error

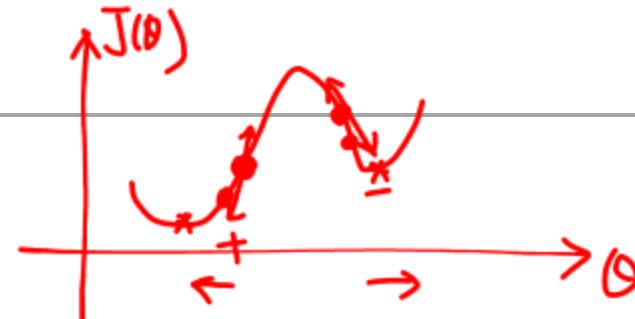
$$\sum_{i=1}^{10} (y_i - \hat{y}_i)^2$$

Target  
(one hot encoding)

# Choosing Proper Loss



# Gradient Descent



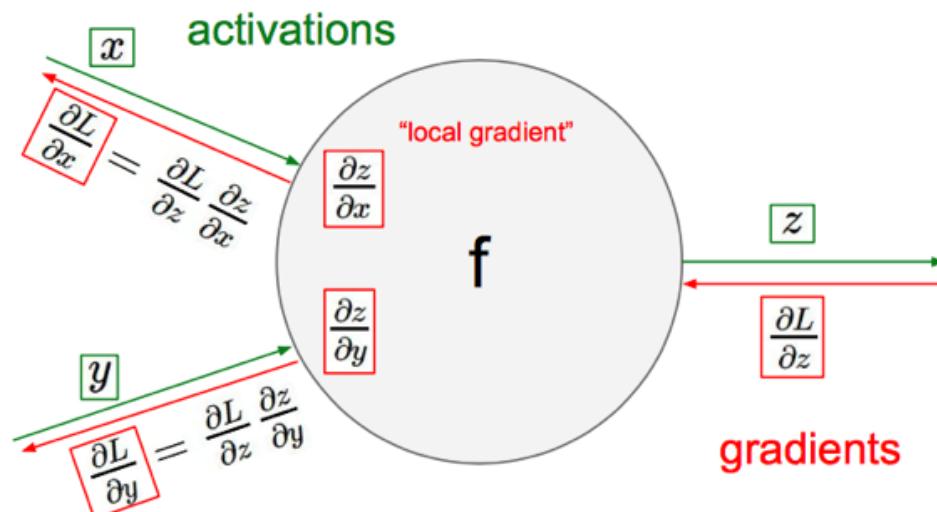
objective/cost function  $J(\theta)$

$$\theta_j^{new} = \theta_j^{old} - \alpha \frac{d}{d\theta_j^{old}} J(\theta) \quad \text{Update each element of } \theta$$

$$\theta^{new} = \theta^{old} - \alpha \underline{\nabla_{\theta} J(\theta)}$$

learning rate

Matrix notation for all parameters

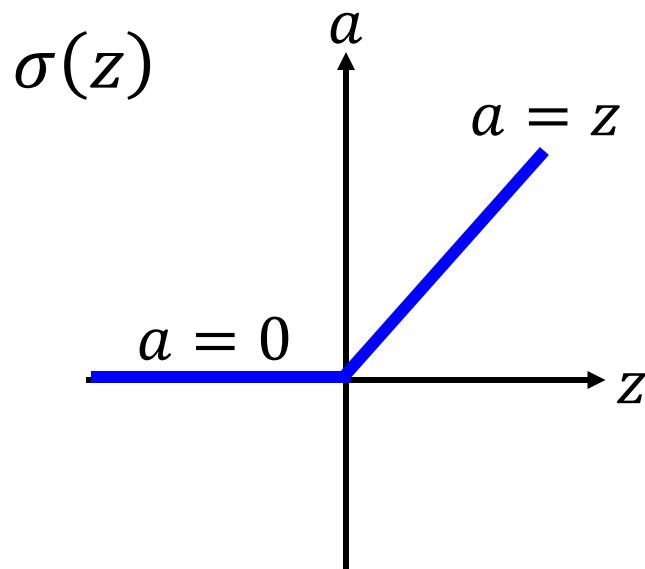


Recursively apply **chain rule** though each node

# ReLU

---

- Rectified Linear Unit (ReLU)



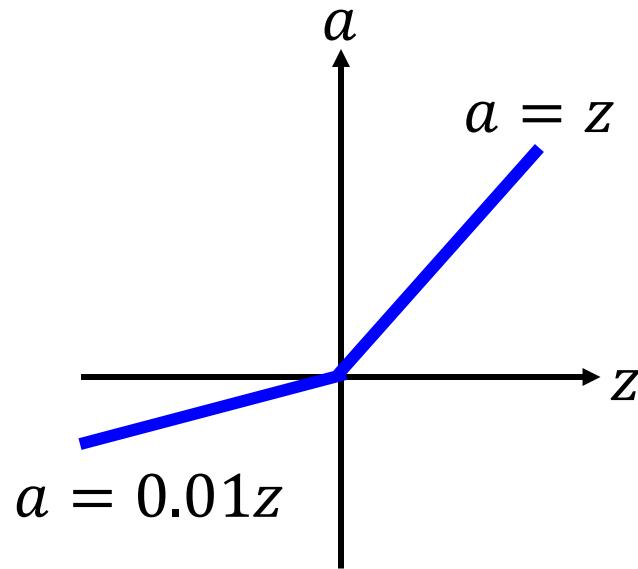
[Xavier Glorot, AISTATS'11]  
[Andrew L. Maas, ICML'13]  
[Kaiming He, arXiv'15]

Vanishing  
gradient problem

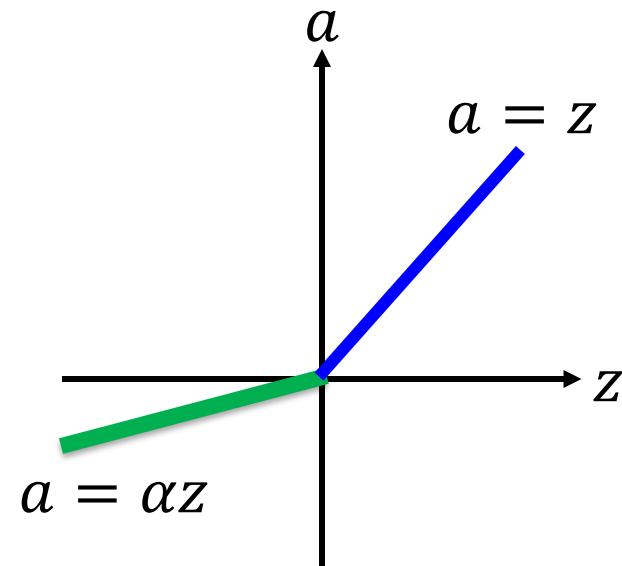
# ReLU - variant

---

*Leaky ReLU*



*Parametric ReLU*

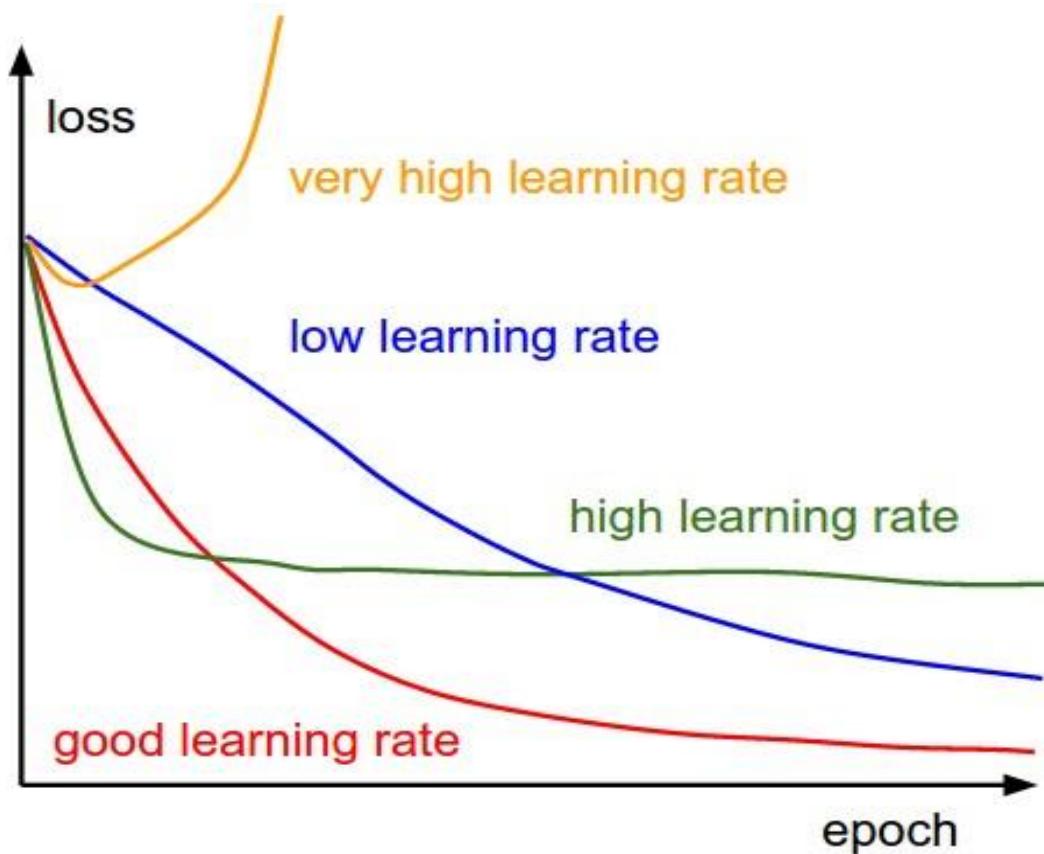


$\alpha$  also learned by  
gradient descent

# How to pick the learning rate?

---

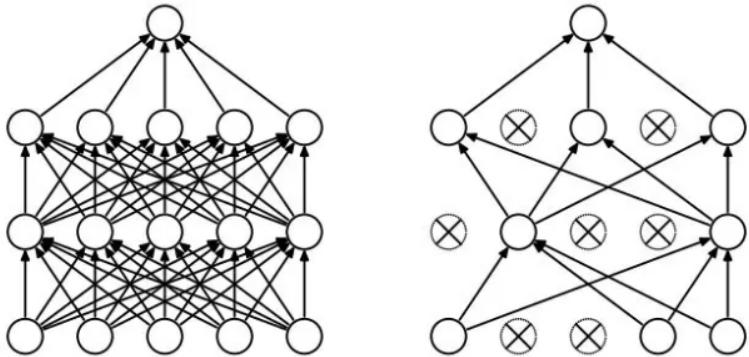
- Too big = diverge, too small = slow convergence
- No “one learning rate to rule them all”
- Start from a high value and keep cutting by half if model diverges
- Learning rate schedule: decay learning rate over time



<http://cs231n.github.io/assets/nn3/learningrates.jpeg>

# Regularization

---



## Dropout

- Randomly drop units (along with their connections) during training
- Each unit retained with fixed probability  $p$ , independent of other units
- Hyper-parameter  $p$  to be chosen (tuned)

Srivastava, Nitish, et al. [\*"Dropout: a simple way to prevent neural networks from overfitting."\*](#) Journal of machine learning research (2014)

## L2 = weight decay

- Regularization term that penalizes big weights, added to the objective – reduce overfitting
- Weight decay value determines how dominant regularization is during gradient computation
- Big weight decay coefficient → big penalty for big weights

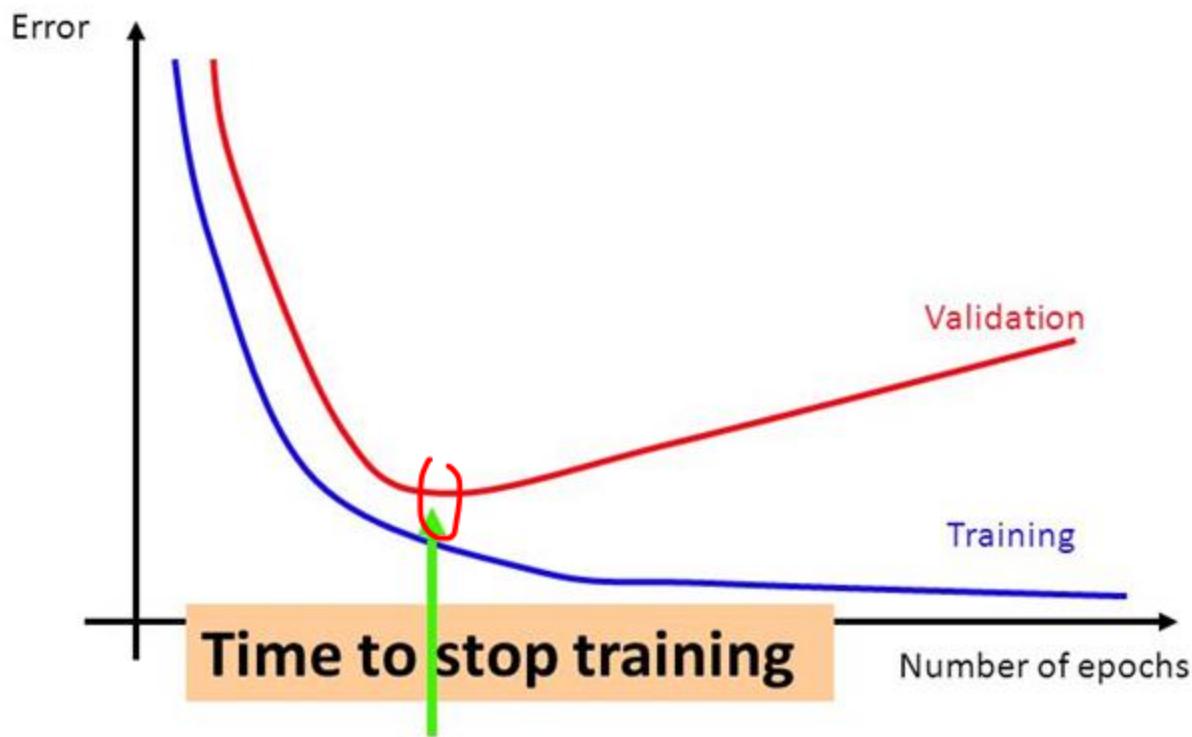
$$J_{reg}(\theta) = J(\theta) + \lambda \sum_k \theta_k^2$$

---

## Early-stopping

- Use validation error to decide when to stop training
- Stop when monitored quantity has not improved after  $n$  subsequent epochs

# Early Stopping

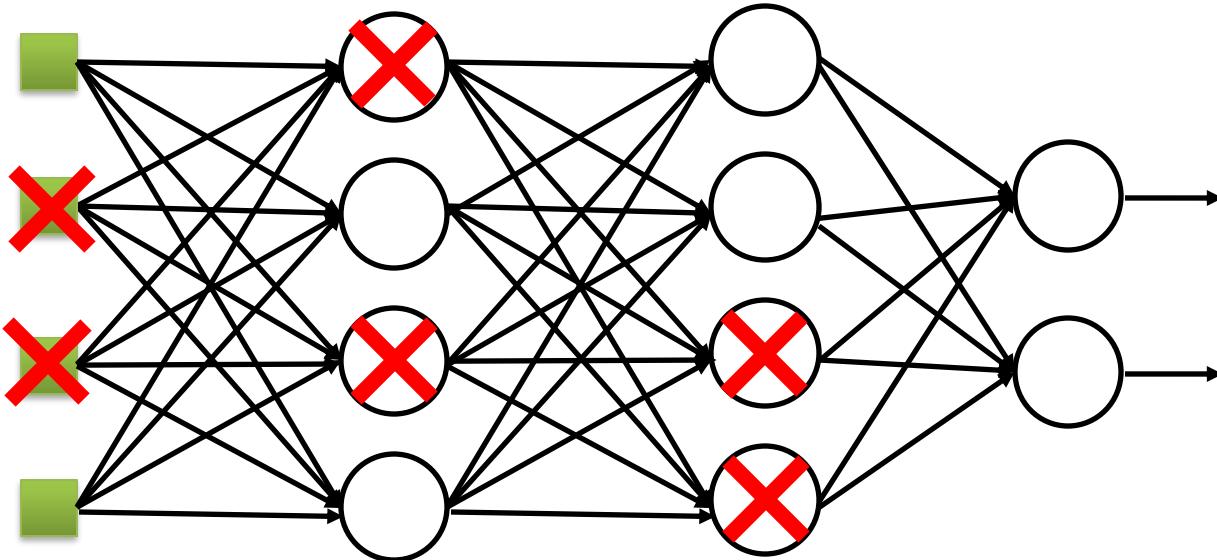


Credit: Stephen Marsland

# Dropout

---

## Training:

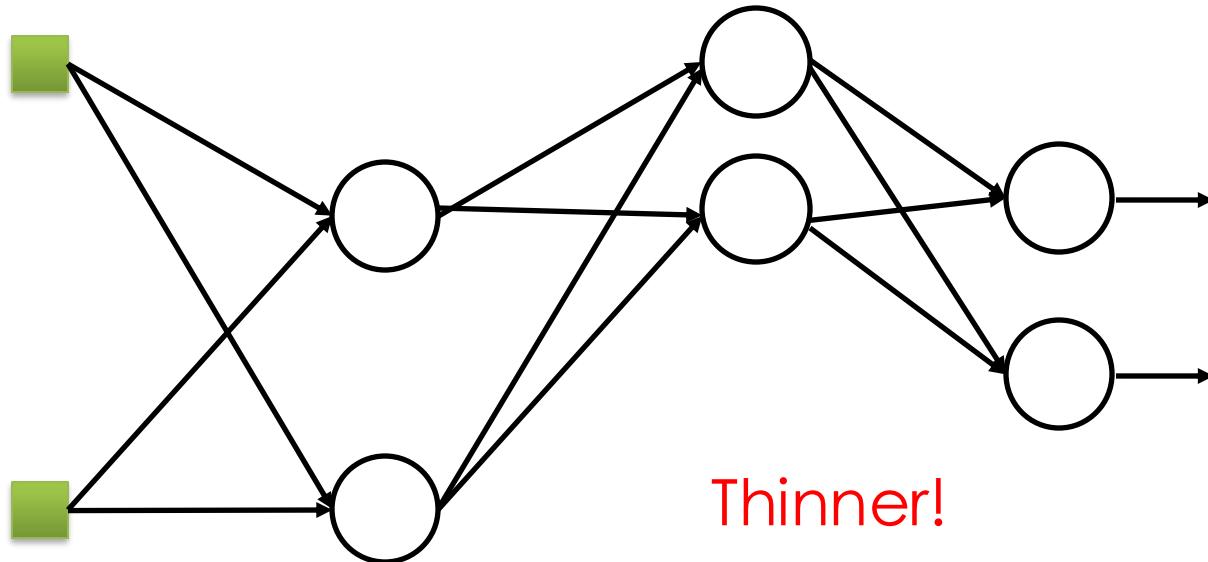


- **Each time before updating the parameters**
  - Each neuron has  $p\%$  to dropout

# Dropout

---

## Training:

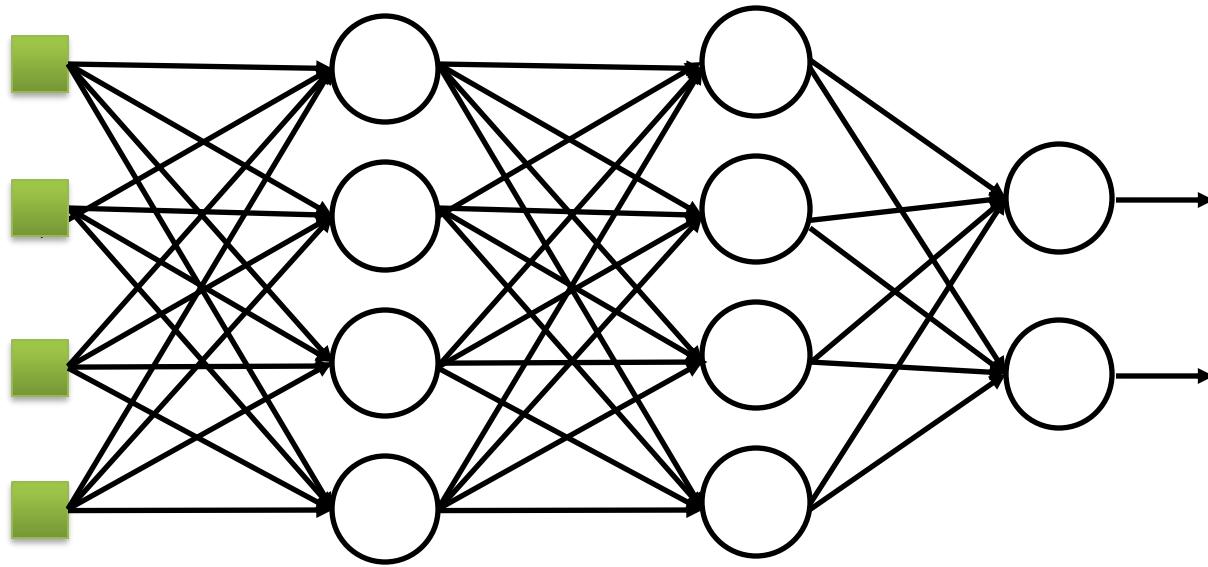


- **Each time before updating the parameters**
    - Each neuron has  $p\%$  to dropout
    - Using the new network for training
- The structure of the network is changed.**

# Dropout

---

**Testing:**

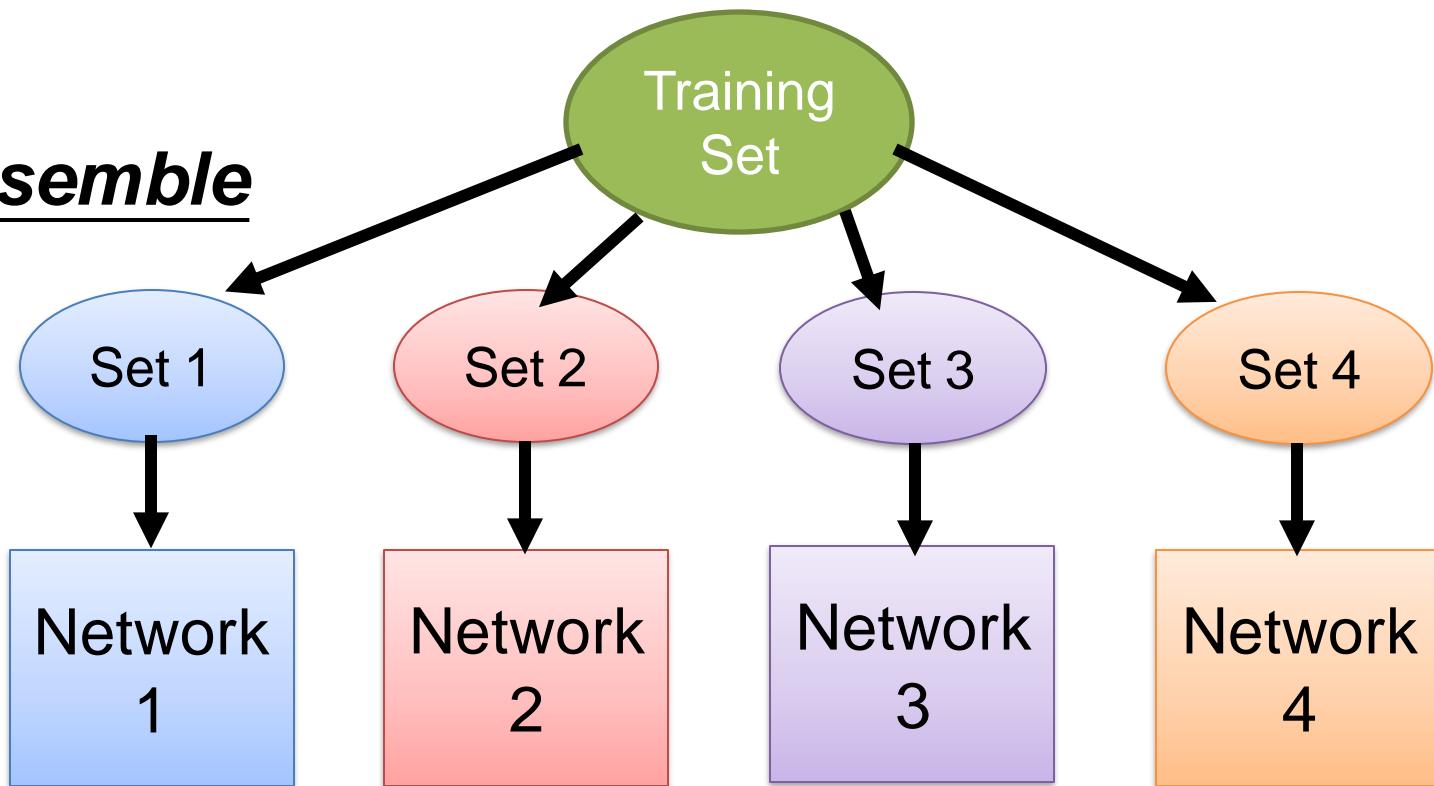


➤ **No dropout**

# Dropout is a kind of ensemble

---

**Ensemble**



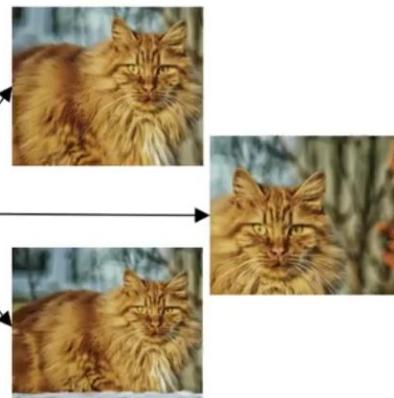
Train a bunch of networks with different structures

# Data Augmentation

- Why data augmentation -> increase training data



Y



{  
Rotation  
Shearing  
Local warping  
...

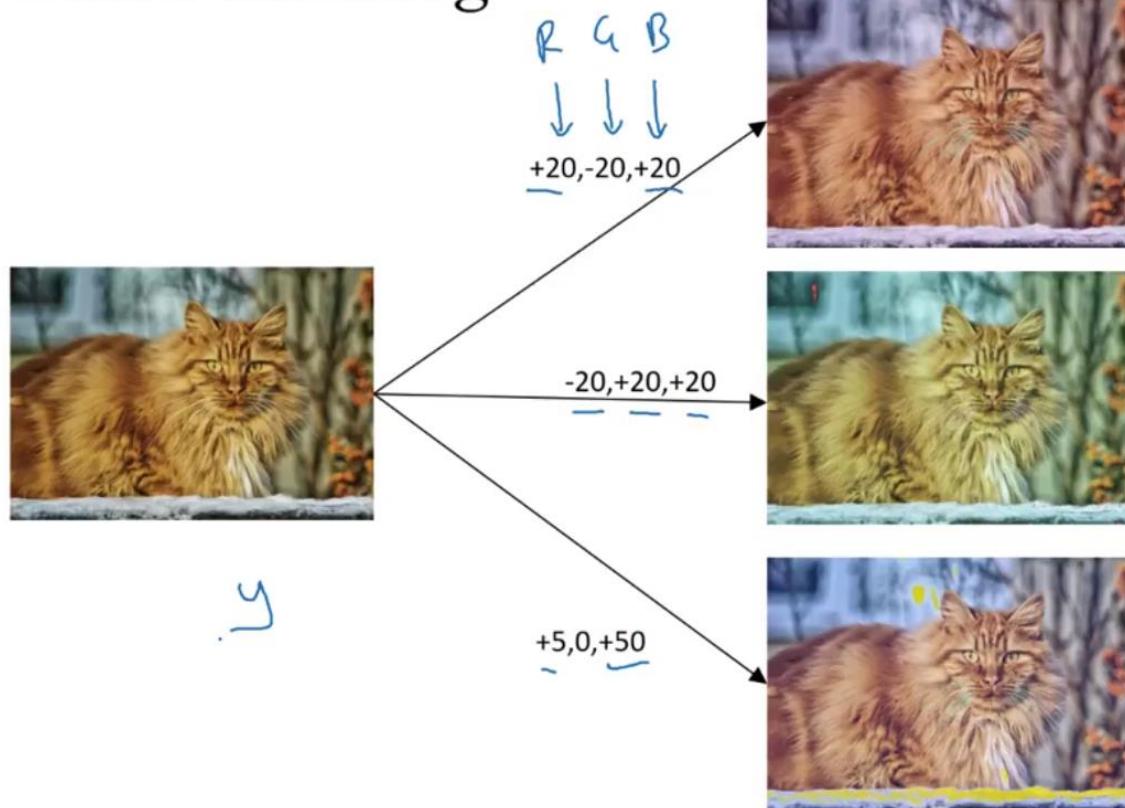


Credit: Andrew Ng



# Data Augmentation

## Color shifting



Credit: Andrew Ng

# Data Augmentation

---

- Color space conversion

Comparison of results for different color spaces on CIFAR-10 with simple CNN

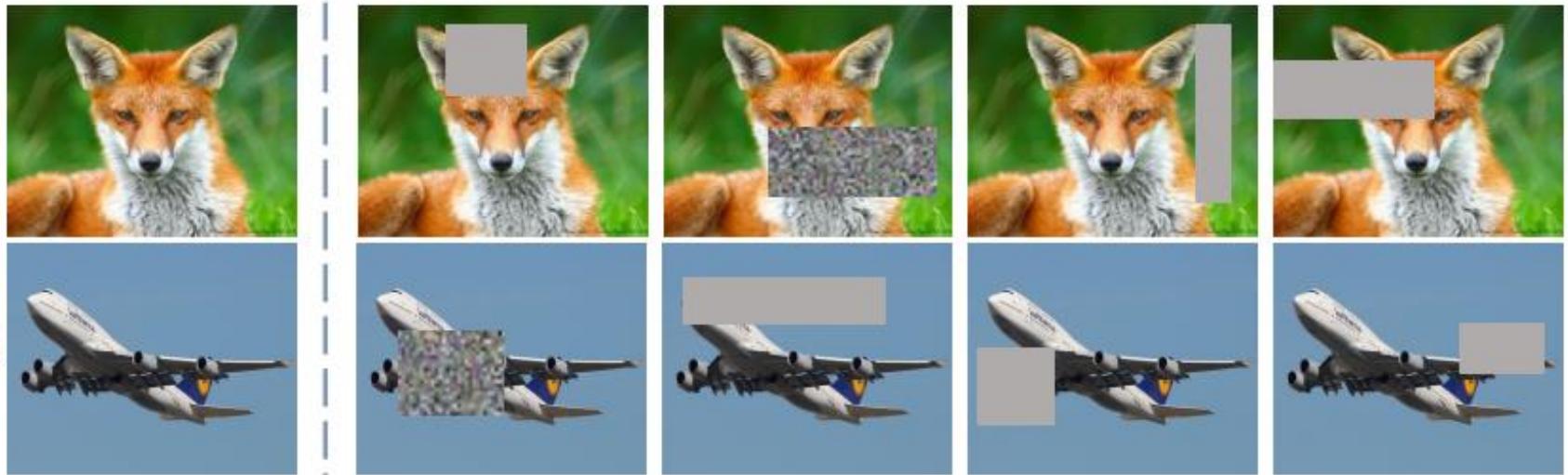
Color Space	Accuracy	Time
RGB	78.89	26 secs
HSV	78.57	26 secs
YUV	78.89	26 secs
LAB	<b>80.43</b>	26 secs
YIQ	78.79	26 secs
XYZ	78.72	26 secs
YPbPr	78.78	26 secs
YCbCr	78.81	26 secs
HED	78.98	26 secs
LCH	78.82	26 secs

Gowda, Shreyank N., and Chun Yuan. "ColorNet: Investigating the importance of color spaces for image classification." *arXiv preprint arXiv:1902.00267* (2019).

# Data Augmentation

- Random erasing

image classification



Reduces the risk of over-fitting and makes the model robust to  
occlusion  
Improve the generalization ability of CNNs

# Data Augmentation

- Learning Augmentation Policies from Data

Search space of operations: ShearX/Y, TranslateX/Y, Rotate, AutoContrast, Invert, Equalize, Solarize, Posterize, Contrast, Color, Brightness, Sharpness,

	Original	Sub-policy 1	Sub-policy 2	Sub-policy 3	Sub-policy 4	Sub-policy 5
Batch 1						
Batch 2						
Batch 3						
	Equalize, 0.4, 4 Rotate, 0.8, 8	Solarize, 0.6, 3 Equalize, 0.6, 7	Posterize, 0.8, 5 Equalize, 1.0, 2	Rotate, 0.2, 3 Solarize, 0.6, 8	Equalize, 0.6, 8 Posterize, 0.4, 6	

Cubuk, Ekin D., et al. "Autoaugment: Learning augmentation policies from data." *arXiv preprint arXiv:1805.09501* (2018).

# Data Augmentation

- Cutmix

Image	ResNet-50	Mixup [48]	Cutout [3]	CutMix
Label	Dog 1.0	Dog 0.5 Cat 0.5	Dog 1.0	Dog 0.6 Cat 0.4
ImageNet	76.3	77.4	77.1	<b>78.6</b>
Cls (%)	(+0.0)	(+1.1)	(+0.8)	(+2.3)
ImageNet	46.3	45.8	46.7	<b>47.3</b>
Loc (%)	(+0.0)	(-0.5)	(+0.4)	(+1.0)
Pascal VOC	75.6	73.9	75.1	<b>76.7</b>
Det (mAP)	(+0.0)	(-1.7)	(-0.5)	(+1.1)

Yun, Sangdoo, et al. "Cutmix: Regularization strategy to train strong classifiers with localizable features." *Proceedings of the IEEE International Conference on Computer Vision*. 2019.

- 
- More on data augmentation
    - (Pytorch) torchvision.transforms:  
<https://pytorch.org/vision/master/transforms.html>
    - <https://github.com/AgaMiko/data-augmentation-review>

# Reading Material

---

- Training a Classifier — PyTorch Tutorials:  
[https://pytorch.org/tutorials/beginner/blitz/cifar10\\_tutorial.html](https://pytorch.org/tutorials/beginner/blitz/cifar10_tutorial.html)
- CNNs, Part 1: An Introduction to Convolutional Neural Networks (Keras implementation)
  - <https://victorzhou.com/blog/intro-to-cnns-part-1/>
- CNNs, Part 2: Training a Convolutional Neural Network (Keras implementation)
  - <https://victorzhou.com/blog/intro-to-cnns-part-2/>
- **MNIST digit classification (Keras implementation)**
  - <https://victorzhou.com/blog/keras-cnn-tutorial/#the-full-code>
- Mnist classification (CNN Keras)
  - <https://www.kaggle.com/elcaiseri/mnist-simple-cnn-keras-accuracy-0-99-top-1>
- Bag of Tricks for Image Classification with Convolutional Neural Networks
  - [https://openaccess.thecvf.com/content\\_CVPR\\_2019/papers/He\\_Bag\\_of\\_Tricks\\_for\\_Image\\_Classification\\_with\\_Convolutional\\_Neural\\_Networks\\_CVPR\\_2019\\_paper.pdf](https://openaccess.thecvf.com/content_CVPR_2019/papers/He_Bag_of_Tricks_for_Image_Classification_with_Convolutional_Neural_Networks_CVPR_2019_paper.pdf)

# Reading Material

---

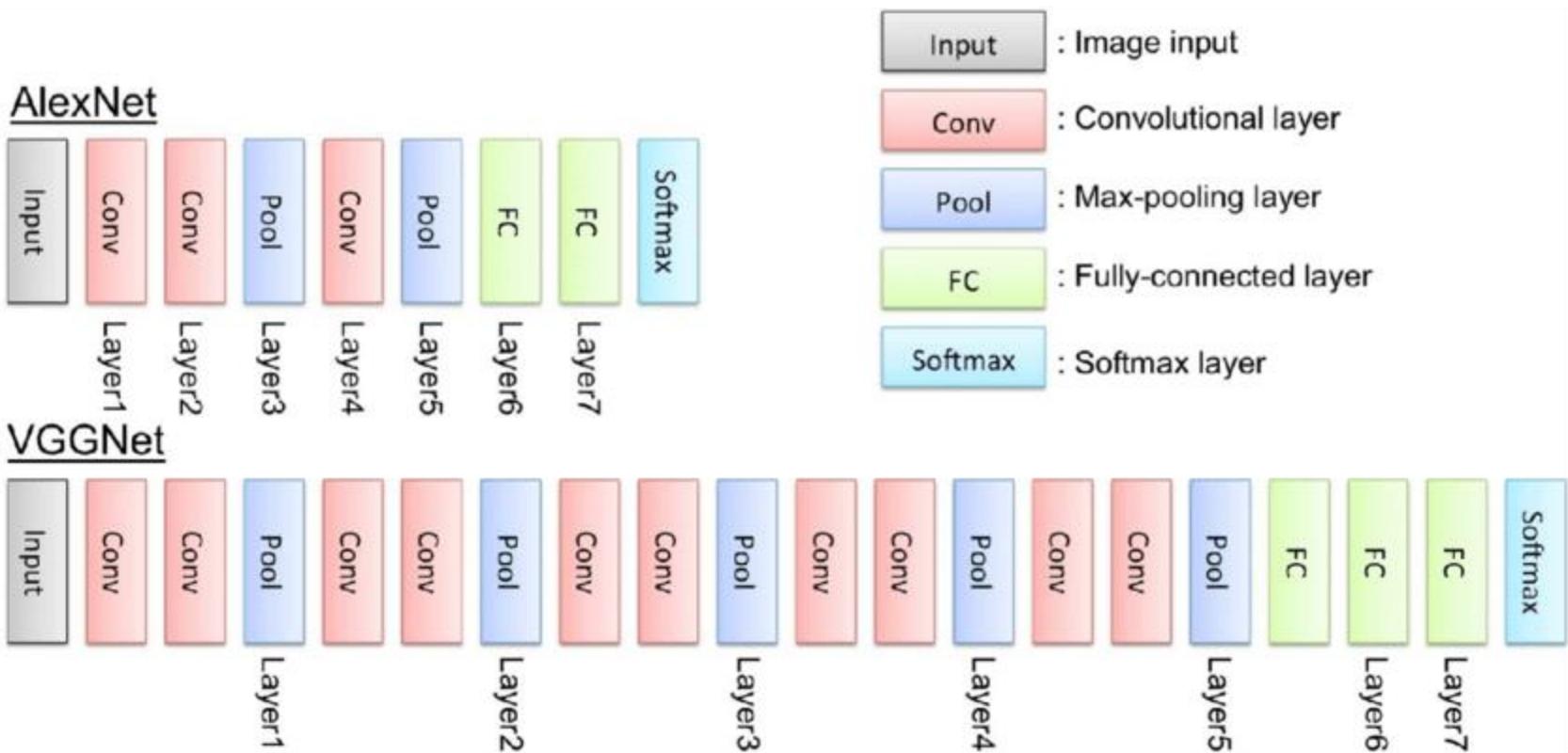
- Deep Learning Tuning Playbook
- [https://github.com/google-research/tuning\\_playbook](https://github.com/google-research/tuning_playbook)

## Who is this document for?

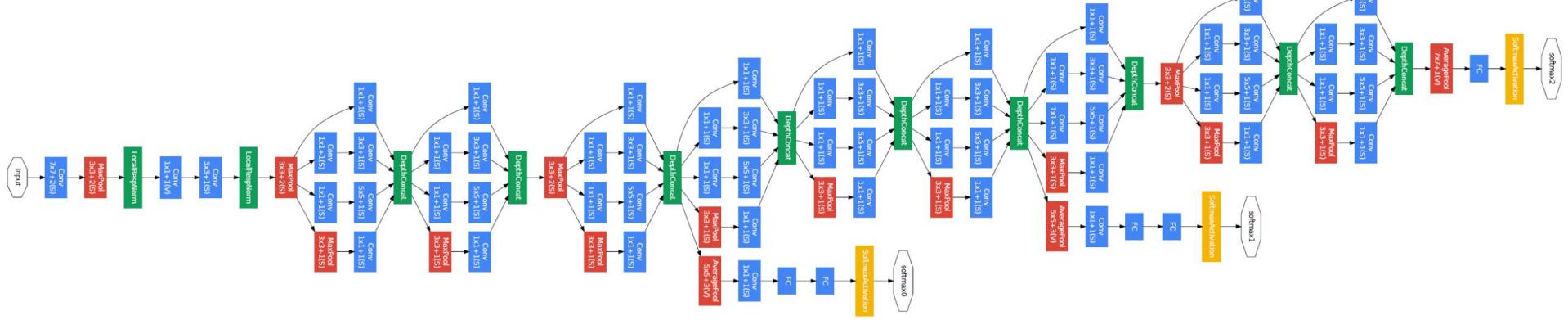
This document is for engineers and researchers (both individuals and teams) interested in **maximizing the performance of deep learning models**. We assume basic knowledge of machine learning and deep learning concepts.

Our emphasis is on the **process of hyperparameter tuning**. We touch on other aspects of deep learning training, such as pipeline implementation and optimization, but our treatment of those aspects is not intended to be complete.

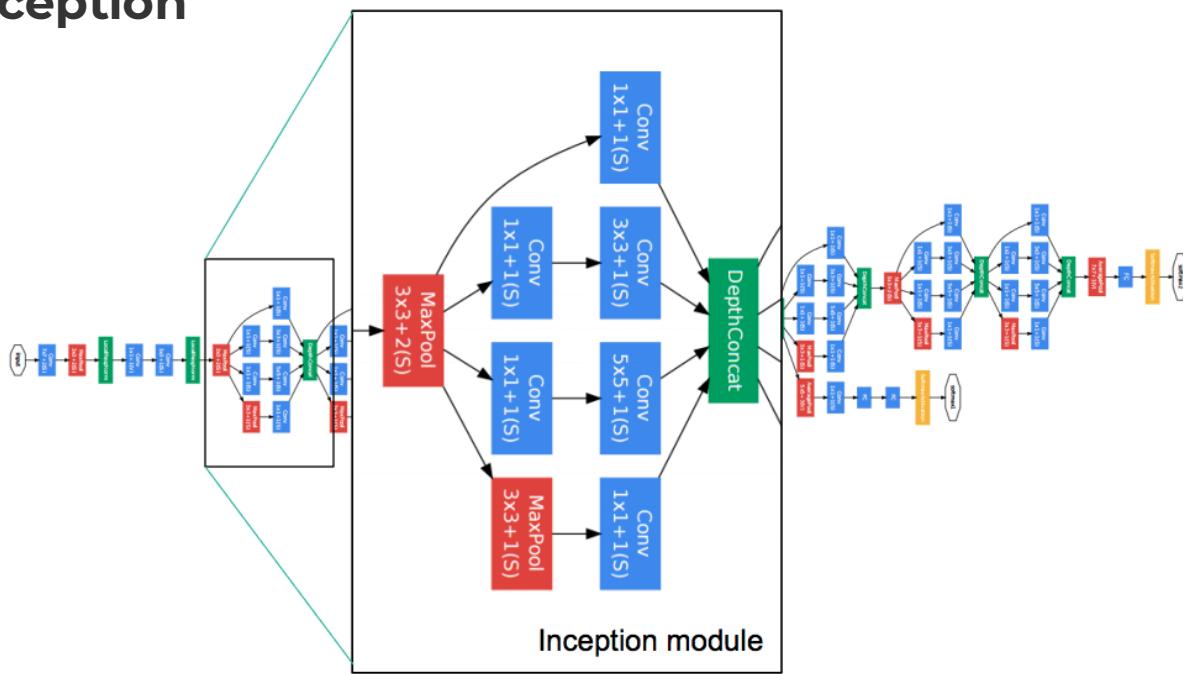
# Case study



# Case study



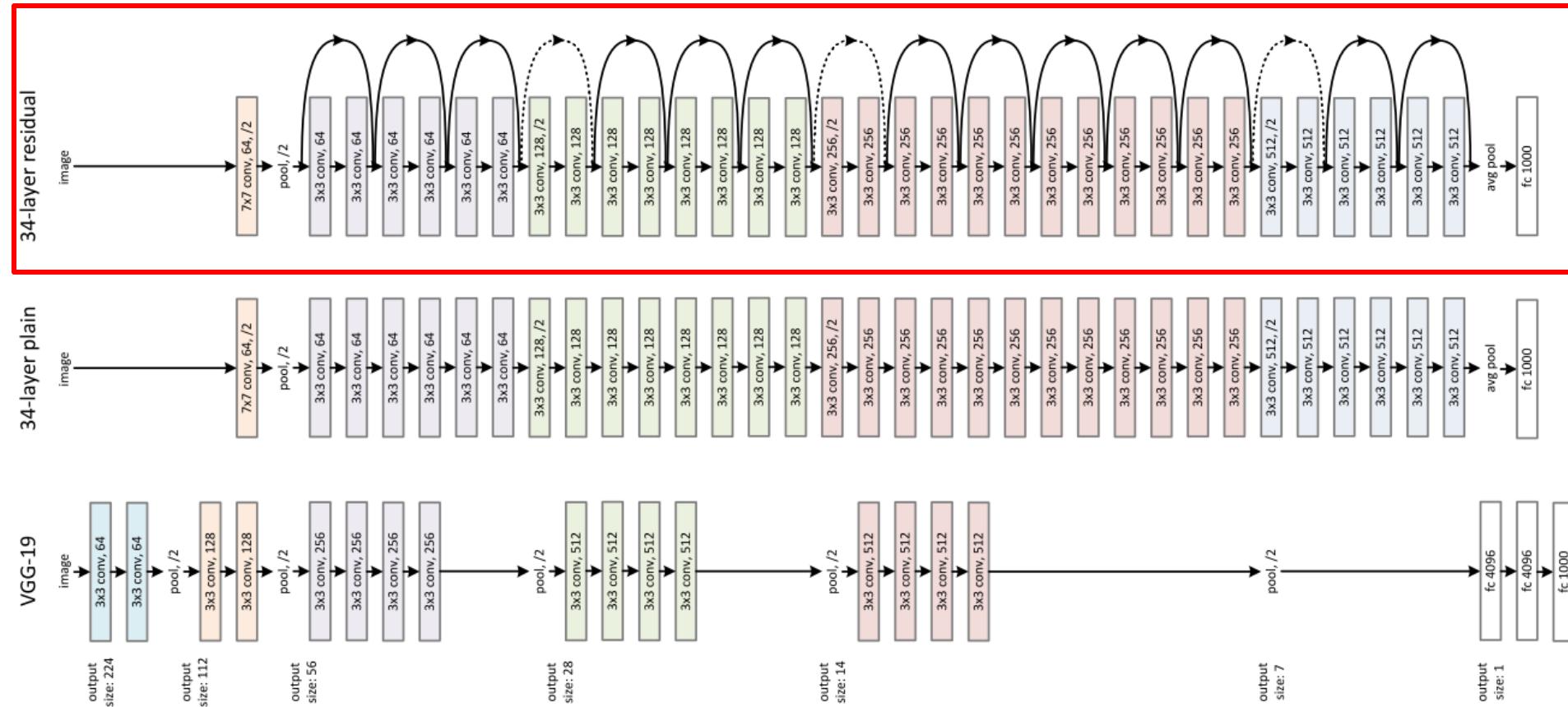
## GoogLeNet/Inception



# Case study

Skip connections:  
Help in backpropagation

Residual Networks



He et. al. Deep Residual Learning for Image Recognition, 2015



# Common network models

---

- Pytorch

```
import torchvision.models as models
resnet18 = models.resnet18()
alexnet = models.alexnet()
vgg16 = models.vgg16()
squeezenet = models.squeezenet1_0()
densenet = models.densenet161()
inception = models.inception_v3()
googlenet = models.googlenet()
shufflenet = models.shufflenet_v2_x1_0()
mobilenet = models.mobilenet_v2()
resnext50_32x4d = models.resnext50_32x4d()
wide_resnet50_2 = models.wide_resnet50_2()
mnasnet = models.mnasnet1_0()
```

## TORCHVISION.MODELS

The models subpackage contains definitions of models for addressing different semantic segmentation, object detection, instance segmentation, person keypoint estimation, and other computer vision tasks.

### Classification

The models subpackage contains definitions for the following model architectures:

- AlexNet
- VGG
- ResNet
- SqueezeNet
- DenseNet
- Inception v3
- GoogLeNet
- ShuffleNet v2
- MobileNet v2
- ResNeXt
- Wide ResNet
- MNASNet

<https://pytorch.org/docs/stable/torchvision/models.html>

# Common network models

---

- Keras

## Available models

Model	Size	Top-1 Accuracy	Top-5 Accuracy	F
Xception	88 MB	0.790	0.945	
VGG16	528 MB	0.713	0.901	
VGG19	549 MB	0.713	0.900	
ResNet50	98 MB	0.749	0.921	
ResNet101	171 MB	0.764	0.928	
ResNet152	232 MB	0.766	0.931	
ResNet50V2	98 MB	0.760	0.930	
ResNet101V2	171 MB	0.772	0.938	
ResNet152V2	232 MB	0.780	0.942	
InceptionV3	92 MB	0.779	0.937	
InceptionResNetV2	215 MB	0.803	0.953	
MobileNet	16 MB	0.704	0.895	
MobileNetV2	14 MB	0.713	0.901	
DenseNet121	33 MB	0.750	0.923	
DenseNet169	57 MB	0.762	0.932	
DenseNet201	80 MB	0.773	0.936	
NASNetMobile	23 MB	0.744	0.919	

<https://keras.io/api/applications/>

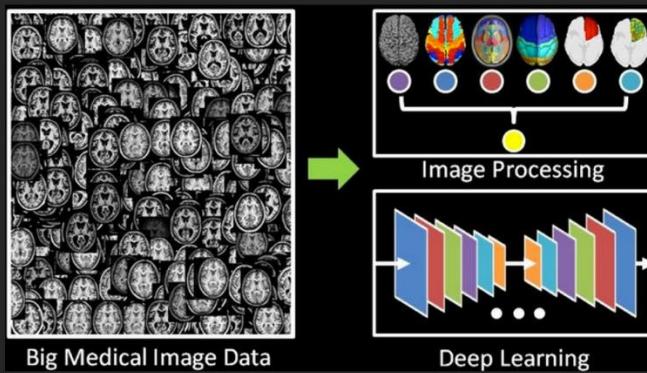
# Interpretability of Deep Neural Networks

## Safety of AI models



Autonomous Driving

## Trust of AI decision



## Policy and Regulation



Right to the explanation  
for algorithmic decisions

# Interpretability of Deep Neural Networks

---

Debate on if interpretability is necessary

For: Rich Caruana, Patrice Simard

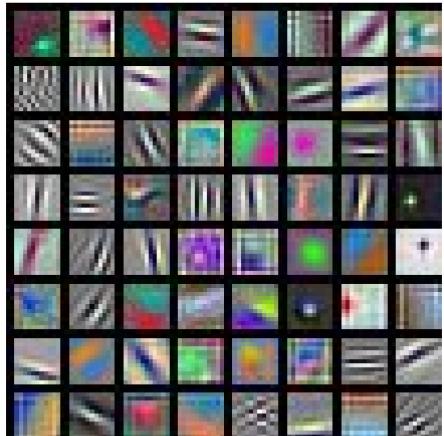
Against: Kilian Weinberger, Yann LeCun



Interpretable Machine Learning Symposium

<https://www.youtube.com/watch?v=93Xv8vJ2acl>

# First Layer: Visualize Filters



AlexNet:  
 $64 \times 3 \times 11 \times 11$



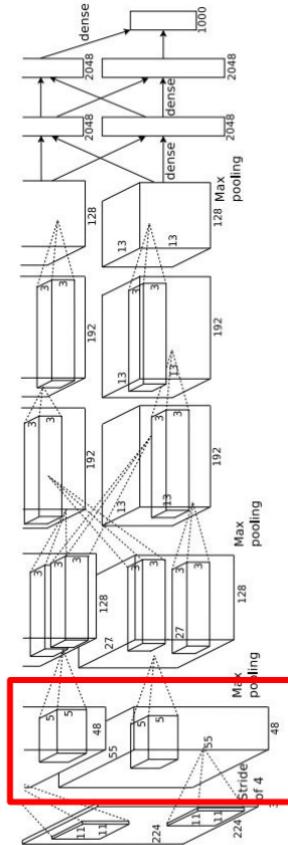
ResNet-18:  
 $64 \times 3 \times 7 \times 7$



ResNet-101:  
 $64 \times 3 \times 7 \times 7$



DenseNet-121:  
 $64 \times 3 \times 7 \times 7$



Krizhevsky, "One weird trick for parallelizing convolutional neural networks", arXiv 2014  
He et al, "Deep Residual Learning for Image Recognition", CVPR 2016  
Huang et al, "Densely Connected Convolutional Networks", CVPR 2017

# Visualize the filters/kernels (raw weights)

We can visualize filters at higher layers, but not that interesting

(these are taken from ConvNetJS  
CIFAR-10 demo)

Weights:  

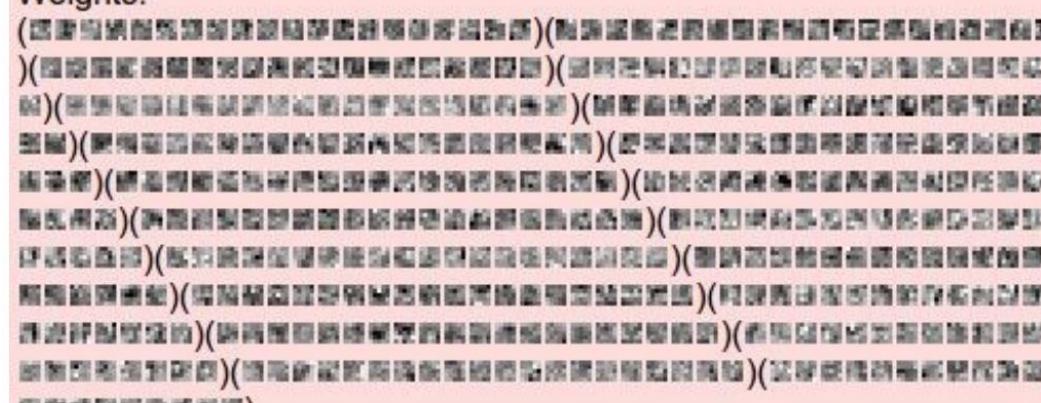

layer 1 weights

$16 \times 3 \times 7 \times 7$

Weights:  
()

layer 2 weights

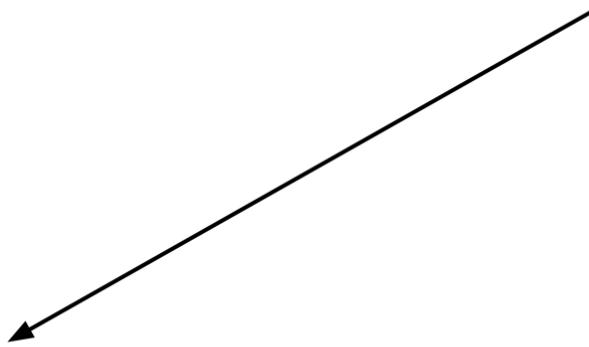
$20 \times 16 \times 7 \times 7$

Weights:  
()

layer 3 weights

$20 \times 20 \times 7 \times 7$

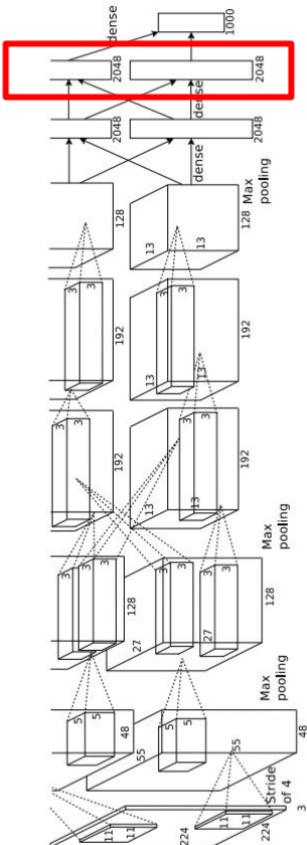
# Last Layer



FC7 layer

4096-dimensional feature vector for an image  
(layer immediately before the classifier)

Run the network on many images, collect the  
feature vectors



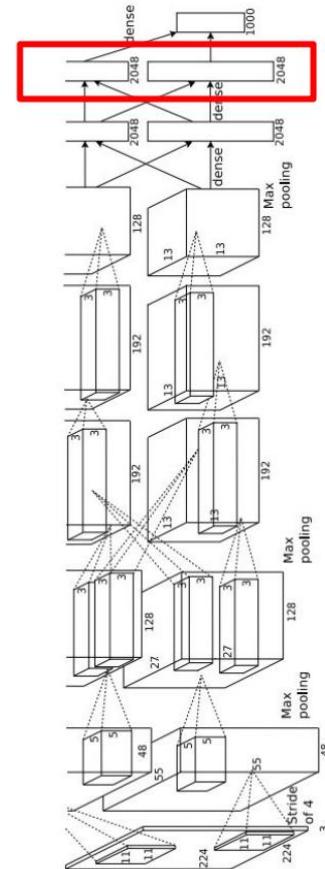
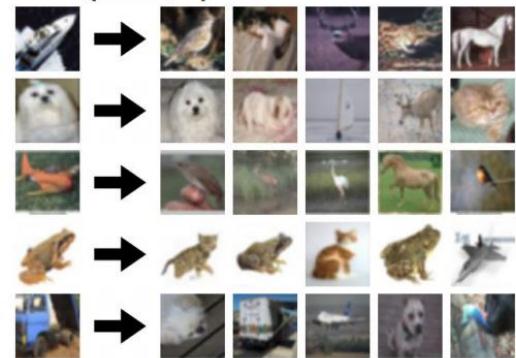
# Last Layer: Nearest Neighbors

4096-dim vector

Test image    L2 Nearest neighbors in feature space



Recall: Nearest neighbors  
in pixel space



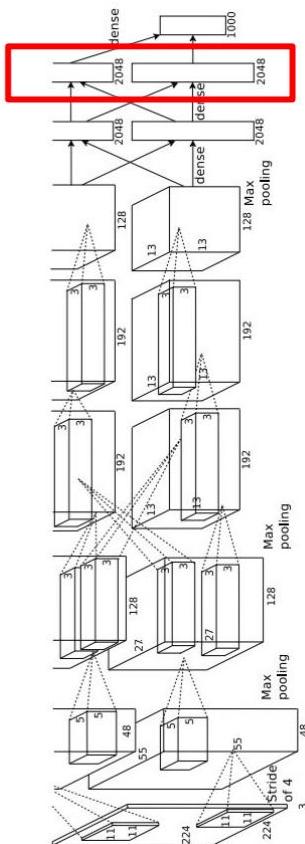
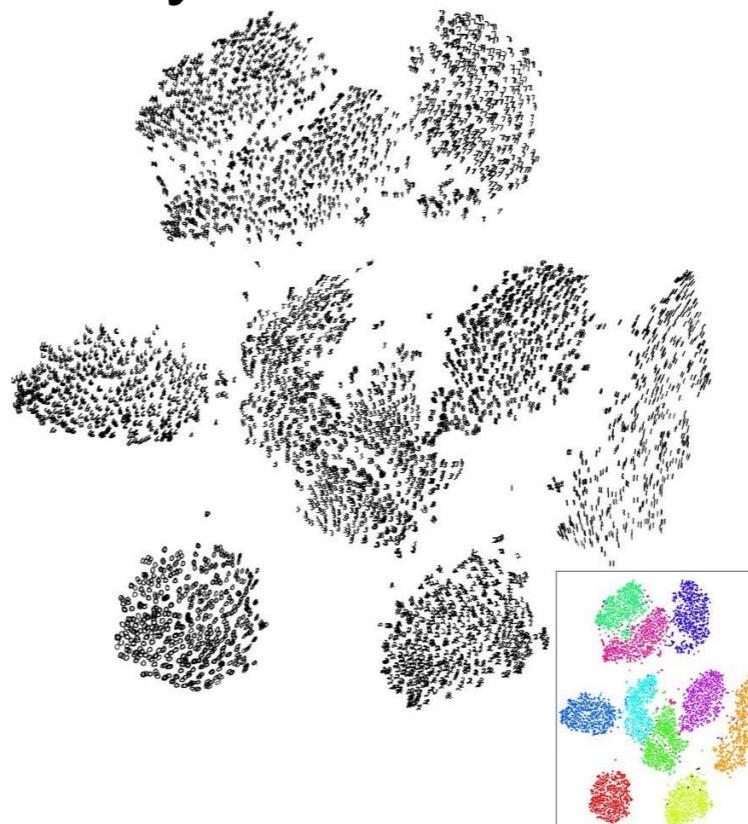
Krizhevsky et al, "ImageNet Classification with Deep Convolutional Neural Networks", NIPS 2012.  
Figures reproduced with permission.

# Last Layer: Dimensionality Reduction

Visualize the “space” of FC7 feature vectors by reducing dimensionality of vectors from 4096 to 2 dimensions

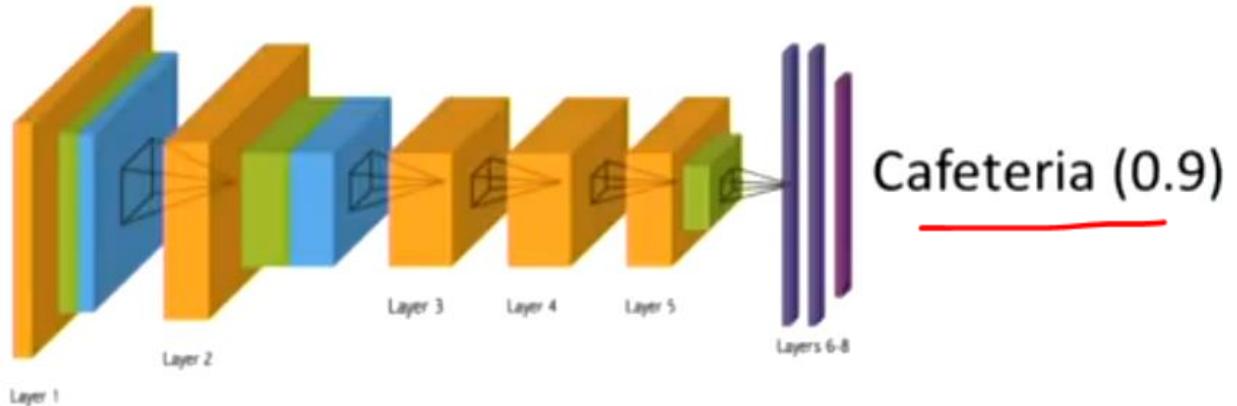
Simple algorithm: Principle Component Analysis (PCA)

More complex: t-SNE



Van der Maaten and Hinton, "Visualizing Data using t-SNE", JMLR 2008  
Figure copyright Laurens van der Maaten and Geoff Hinton, 2008. Reproduced with permission.

# Why the network gives such prediction?



Credit: Bolei Zhou

# Class Activation Map

- Explain Prediction of Deep Neural Network

Prediction: Conference Center



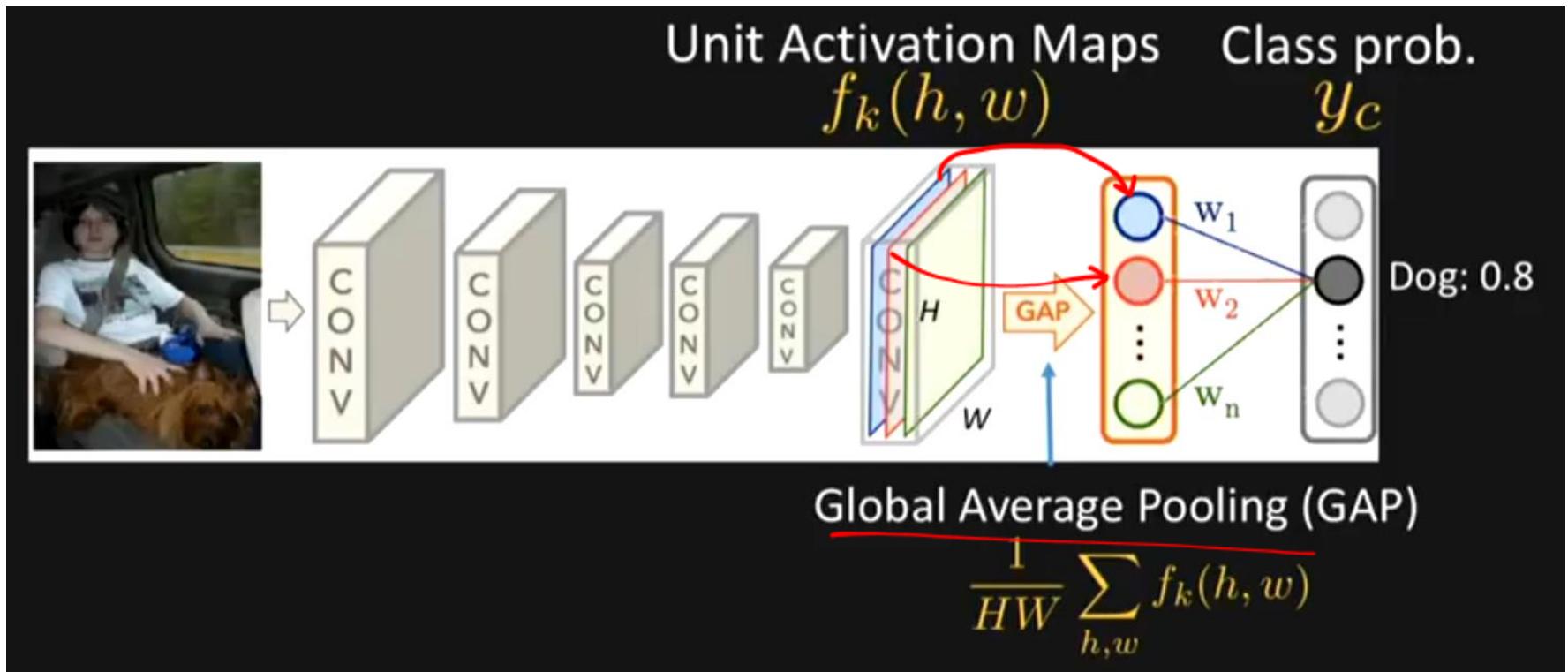
Prediction: Indoor Booth



Credit: Bolei Zhou

Zhou, Bolei, et al. "Learning deep features for discriminative localization." *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2016.

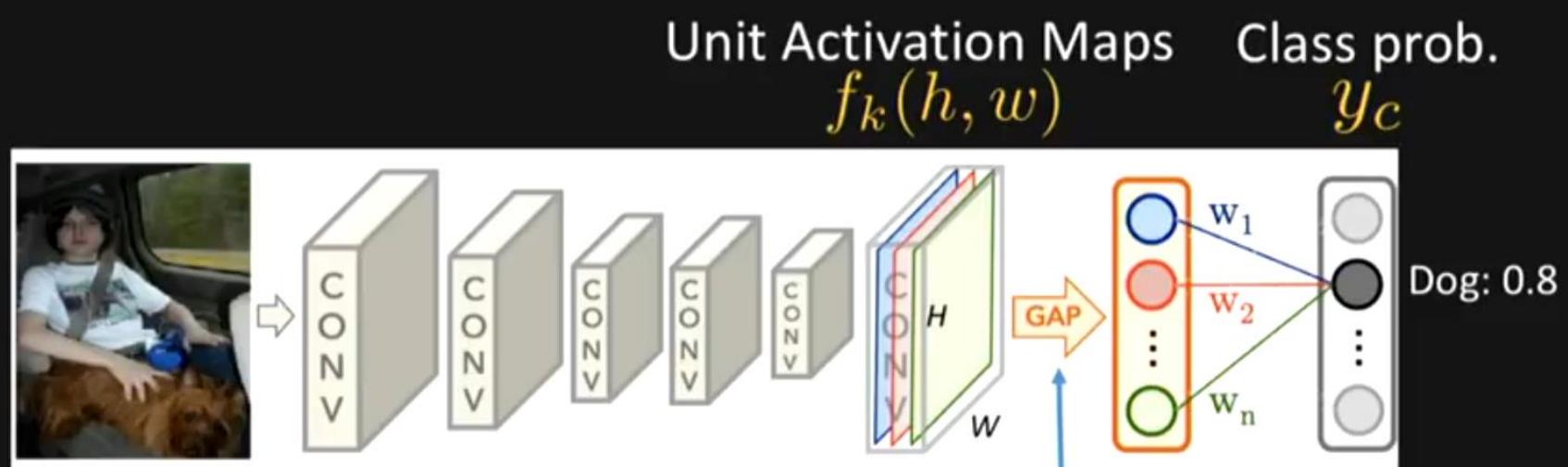
# What the CNN is looking



Credit: Bolei Zhou

Zhou, Bolei, et al. "Learning deep features for discriminative localization." *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2016.

# What the CNN is looking



Global Average Pooling (GAP)

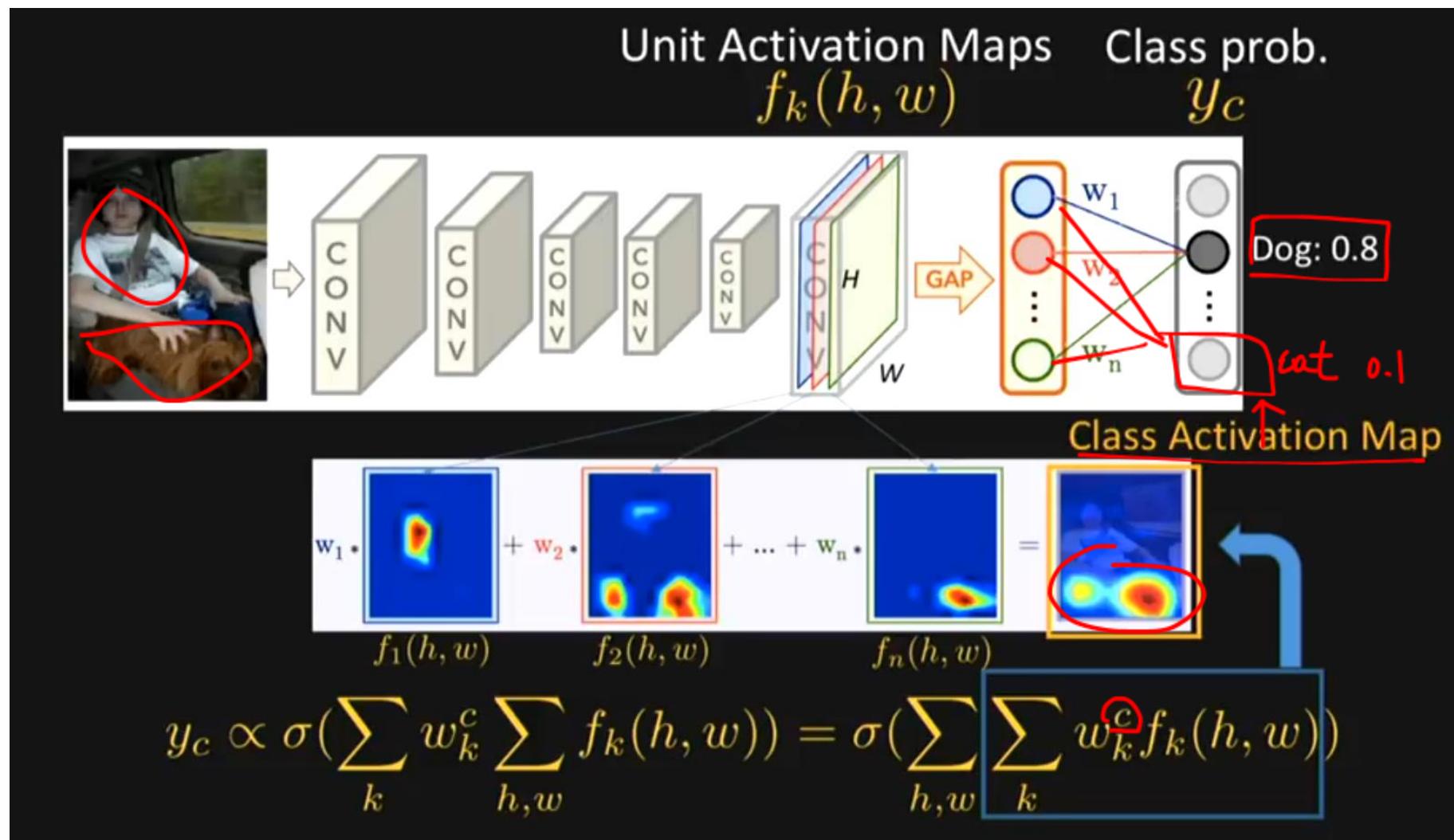
$$\frac{1}{HW} \sum_{h,w} f_k(h, w)$$

$$y_c \propto \sigma \left( \sum_k w_k^c \sum_{h,w} f_k(h, w) \right) = \sigma \left( \sum_{h,w} \sum_k w_k^c f_k(h, w) \right)$$

Zhou, Bolei, et al. "Learning deep features for discriminative localization." *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2016.

Credit: Bolei Zhou

# What the CNN is looking

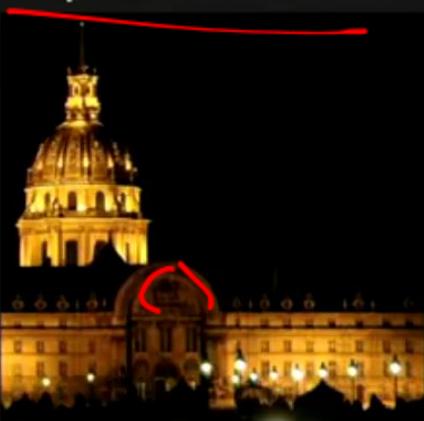


Zhou, Bolei, et al. "Learning deep features for discriminative localization." *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2016.

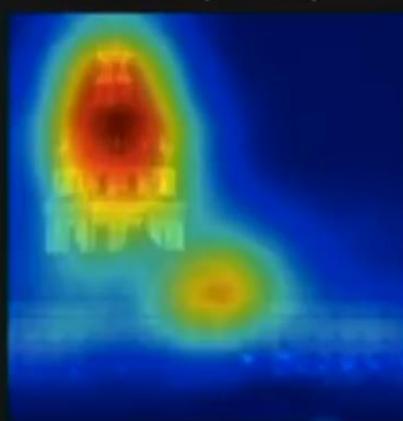
Credit: Bolei Zhou

# Class Activation Mapping: Explain Prediction of Deep Neural Network

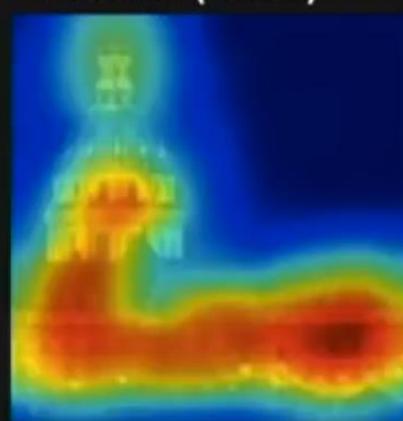
Top3 Predictions:



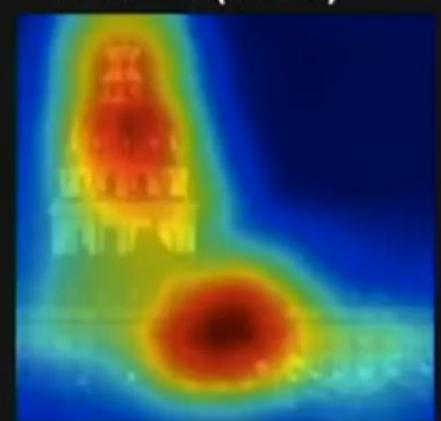
Dome (0.45)



Palace (0.21)



Church (0.10)



Credit: Bolei Zhou



Credit: Bolei Zhou

# Explain the failure cases

---



GT: House  
Prediction: Sushi bar

Credit: Bolei Zhou



# Explain the failure cases

---

Prediction: Martial Arts Gym (0.21)



Prediction: Martial Arts Gym (0.21)

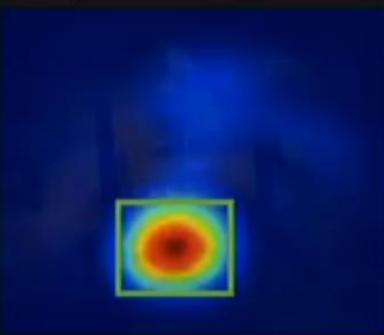
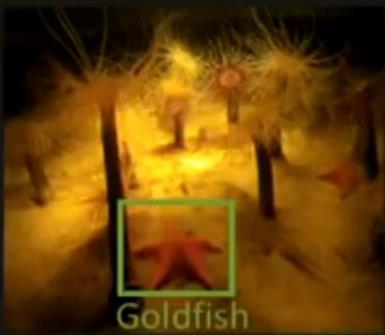


Credit: Bolei Zhou

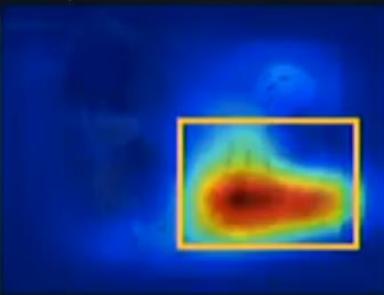
# Using CAM for localization

## Evaluation on Weakly-Supervised Localization

Prediction: Starfish (0.83)



Prediction: Tricycle (0.92)



Method	Supervision	Localization Accuracy(%)
Backpropagation	weakly	53.6
Our method	weakly	62.9
AlexNet	full	65.8

Result on ImageNet Localization Benchmark

Credit: Bolei Zhou

# Limitation of CAM

---

- To apply CAM, any CNN-based network must change its architecture, where GAP is a must before the output layer
  - i.e., architectural changes and hence re-training is needed

# Grad-CAM (Gradient-weighted Class Activation Mapping)

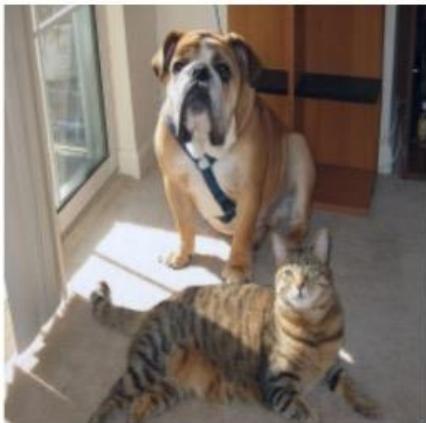
---

- Grad-CAM (Gradient-weighted Class Activation Mapping) generalizes CAM for a wide variety of CNN-based architectures
  - i.e., without requiring architectural changes or re-training
- Characteristics
  - Without GAP layer, we need a way to define weights -  $w_k^c$
  - Grad-CAM uses the gradients of any target class (e.g., “dog” in a classification network) flowing into the final convolutional layer, and derive summary statistics out of it to represent the weights (importance)

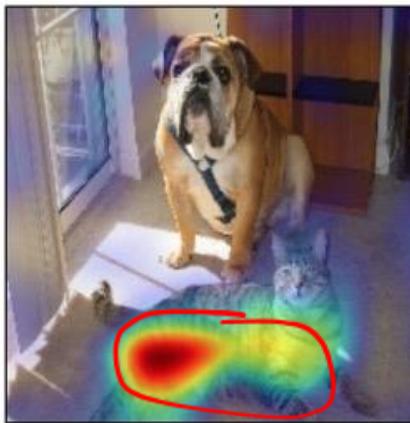
Selvaraju, Ramprasaath R., et al. "Grad-cam: Visual explanations from deep networks via gradient-based localization." Proceedings of the IEEE International Conference on Computer Vision. 2017.

# Grad-CAM

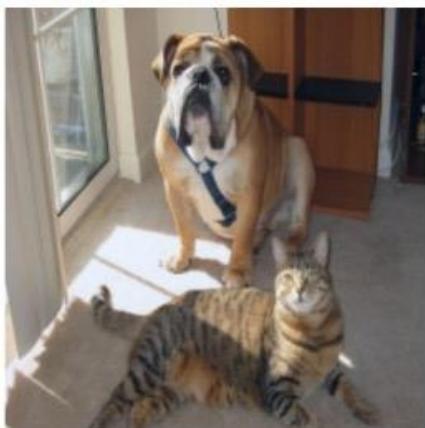
---



(a) Original Image



(c) Grad-CAM 'Cat'



(g) Original Image



(i) Grad-CAM 'Dog'

Selvaraju, Ramprasaath R., et al. "Grad-cam: Visual explanations from deep networks via gradient-based localization." Proceedings of the IEEE International Conference on Computer Vision. 2017.

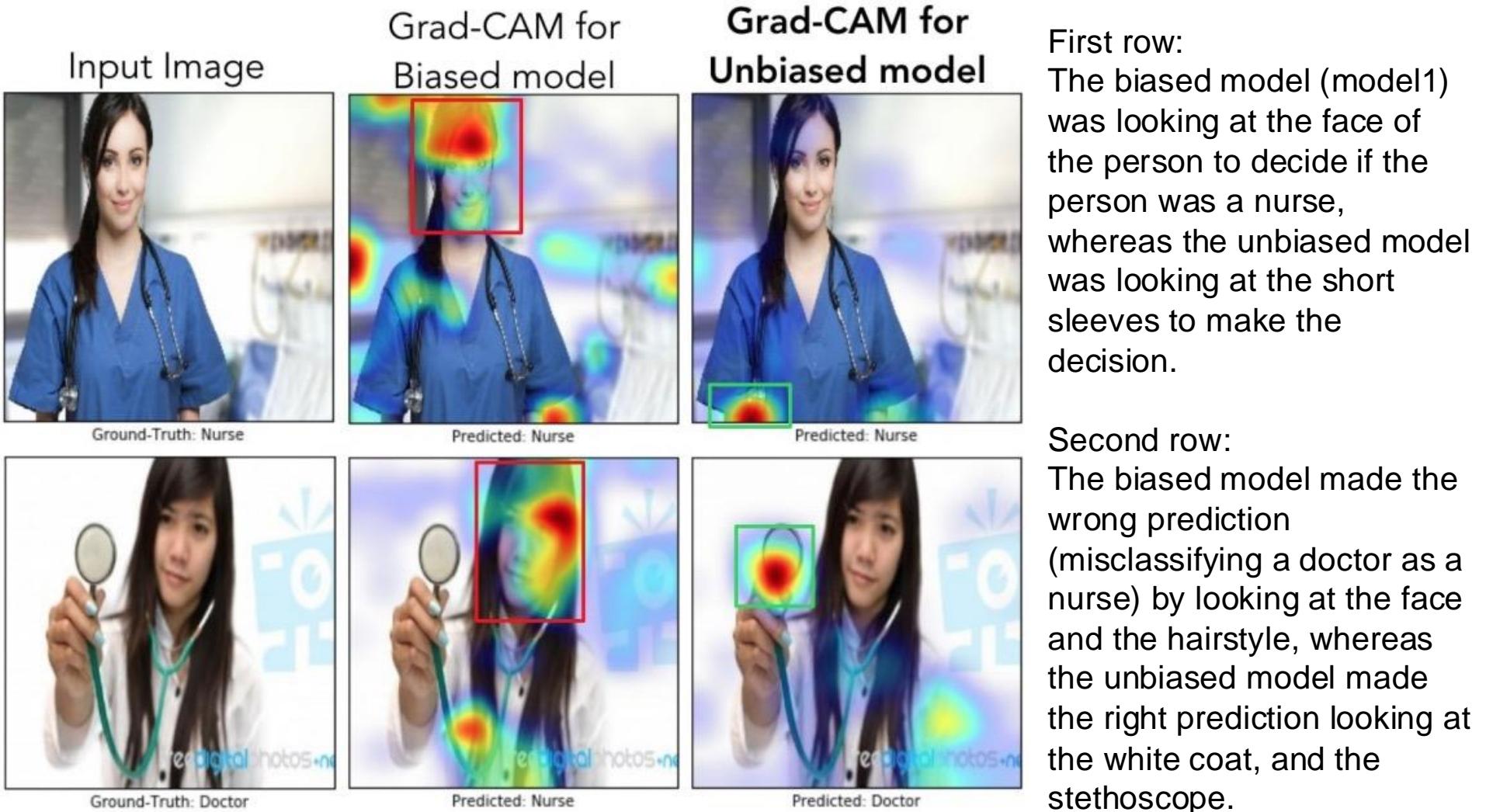
# Grad-CAM for Weakly-supervised Localization

		Classification		Localization	
		Top-1	Top-5	Top-1	Top-5
VGG-16	Backprop [51]	30.38	10.89	61.12	51.46
	c-MWP [58]	30.38	10.89	70.92	63.04
	Grad-CAM (ours)	30.38	10.89	<b>56.51</b>	46.41
AlexNet	CAM [59]	33.40	12.20	57.20	<b>45.14</b>
	c-MWP [58]	44.2	20.8	92.6	89.2
	Grad-CAM (ours)	44.2	20.8	68.3	56.6
GoogleNet	Grad-CAM (ours)	31.9	11.3	60.09	49.34
	CAM [59]	31.9	11.3	60.09	49.34

Table 1: Classification and localization error % on ILSVRC-15 val (lower is better) for VGG-16, AlexNet and GoogleNet. We see that Grad-CAM achieves superior localization errors without compromising on classification performance.

Selvaraju, Ramprasaath R., et al. "Grad-cam: Visual explanations from deep networks via gradient-based localization." Proceedings of the IEEE International Conference on Computer Vision. 2017.

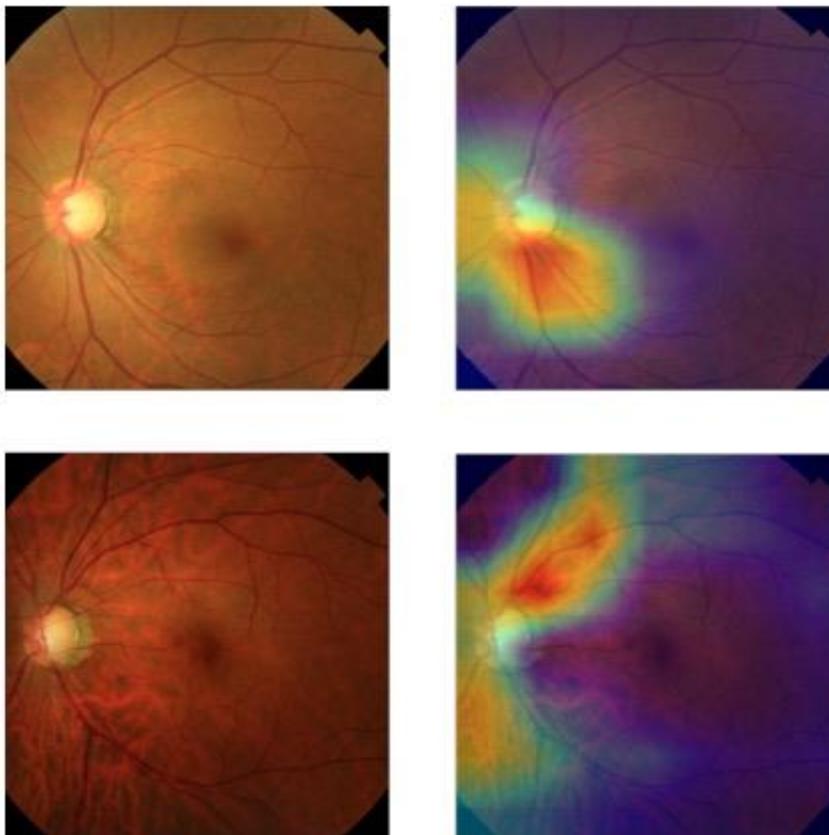
# Grad-CAM for Identifying Bias in Dataset



Selvaraju, Ramprasaath R., et al. "Grad-cam: Visual explanations from deep networks via gradient-based localization." Proceedings of the IEEE International Conference on Computer Vision. 2017.

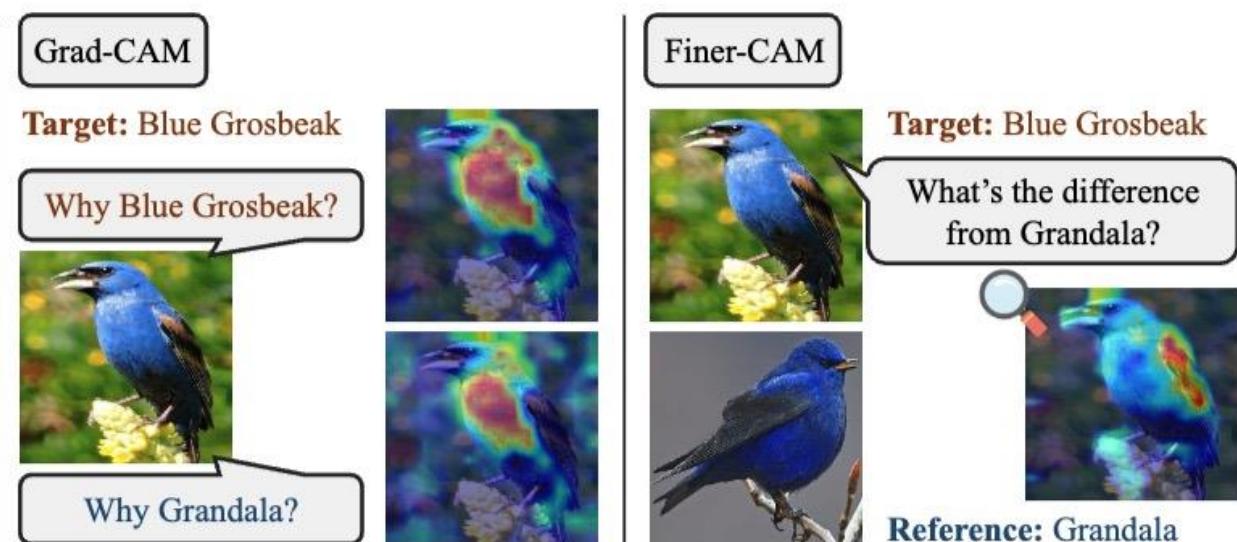
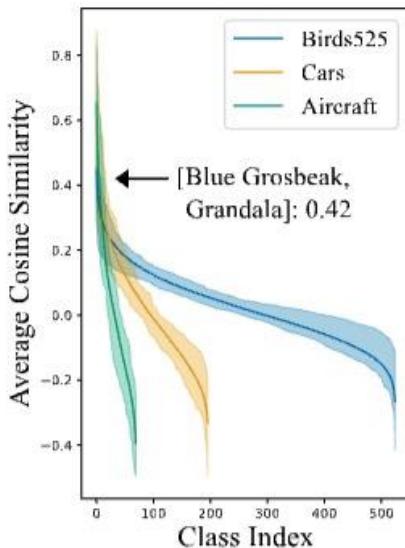
# Grad-CAM for Medical Image Analysis

- Grad-CAM for glaucoma diagnosis and localization



(a) Correctly localized suspicious areas.

Kim, Mijung, et al. "Web applicable computer-aided diagnosis of glaucoma using deep learning." arXiv preprint arXiv:1812.02405 (2018).



**Figure 1. Illustration of Finer-CAM.** **Left:** Sorted cosine similarity between the linear classifier weights of one class and the other classes, averaged across all classes. Many pairs of classes are highly similar, yet neural networks can effectively distinguish them to achieve high fine-grained classification accuracy. **Middle:** Standard CAM methods highlight main regions contributing to the target class's logit value, inadvertently including regions predictive of similar classes and overshadowing fine discriminative details. **Right:** We propose Finer-CAM to explicitly compare the target class with similar classes and spot the difference, enabling accurate localization of discriminative details.

Zhang et al. Finer-CAM: Fine-grained Visual Interpretability through Class-Specific Gradient Refinements  
<https://arxiv.org/pdf/2501.11309>

# More reading

---

- Bau, David, et al. "Network dissection: Quantifying interpretability of deep visual representations." *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 2017.
- Selvaraju, Ramprasaath R., et al. "Grad-cam: Visual explanations from deep networks via gradient-based localization." *Proceedings of the IEEE International Conference on Computer Vision*. 2017.
- Chattopadhyay, Aditya, et al. "Grad-cam++: Generalized gradient-based visual explanations for deep convolutional networks." *2018 IEEE Winter Conference on Applications of Computer Vision (WACV)*. IEEE, 2018.
- Zhang et al. Finer-CAM: Fine-grained Visual Interpretability through Class-Specific Gradient Refinements <https://arxiv.org/pdf/2501.11309>

# Implementation

---

- CAM
  - <https://github.com/zhoubolei/CAM>
- Grad-CAM
  - <https://github.com/ramprs/grad-cam/>
  - Demo: <http://gradcam.cloudcv.org/>
- TorchCAM: class activation explorer
  - <https://github.com/frgfm/torch-cam>

---

Thank you!

Question?

# References and Slide Credits

---

- Many slides are adapted from the existing teaching or tutorial slides by Hung-yi Lee, Andrew Ng, Alexander Amini, Lex Fridman, and Stanford course - CS231n: Convolutional Neural Networks for Visual Recognition
- Special thanks to Dr. Hung-yi Lee for making his machine learning course slides and materials available
- Alexander Amini, MIT 6.S191 Introduction to Deep Learning:  
<http://introtodeeplearning.com/>
  - Youtube videos:  
[https://www.youtube.com/watch?v=5tvmMX8r\\_OM&list=PLtBw6njQRU-rwp5\\_7C0oIVt26ZgjG9NI&index=1](https://www.youtube.com/watch?v=5tvmMX8r_OM&list=PLtBw6njQRU-rwp5_7C0oIVt26ZgjG9NI&index=1)
- Lex Fridman, MIT Deep Learning and Artificial Intelligence Lectures: <https://deeplearning.mit.edu/>  
<https://www.youtube.com/watch?v=O5xeyoRL95U>