

# DEEP LEARNING - VISION 1

## INTRODUCTION TO CNNs

Sumandeep Banerjee  
Module 3: Deep Learning  
<https://shala2020.github.io/>

# LEARNING OBJECTIVES

- Applications of Computer Vision
- Motivation from Digital Image Processing Filters
- Implementation of Convolution
- Building a Convolutional Neural Network
- Suitability of Convolutions to Computer Vision
- Case Study - Image Classification using CNNs
- Advanced CNN Architectures
- Practical Training Techniques

# APPLICATIONS OF COMPUTER VISION

# CLASSIFICATION VS LOCALIZATION VS DETECTION



Car or No Car ?

Image Classification

Where is the car ?

Object Localization

Where are the cars ?

Object Detection

What all is there ?

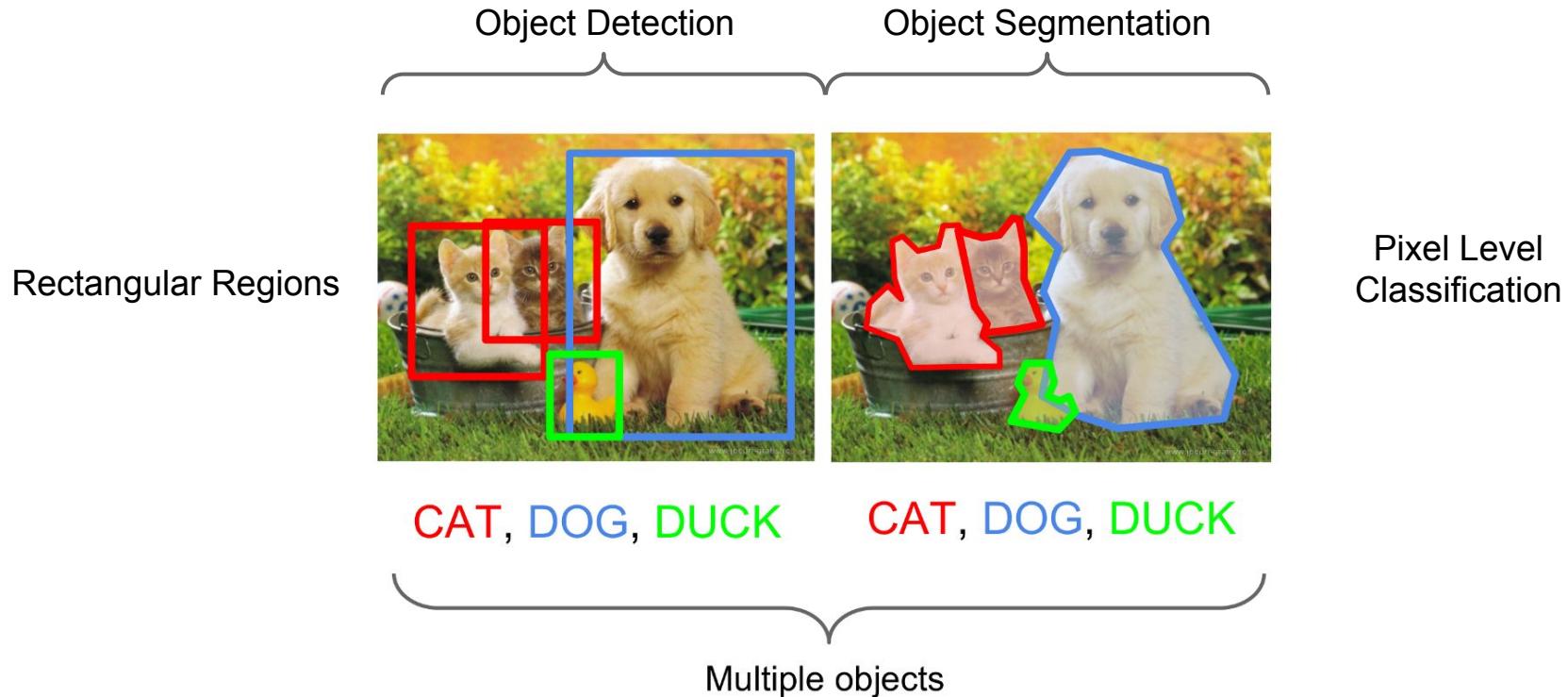
Multi-Class Detection

Special Case: Recognition

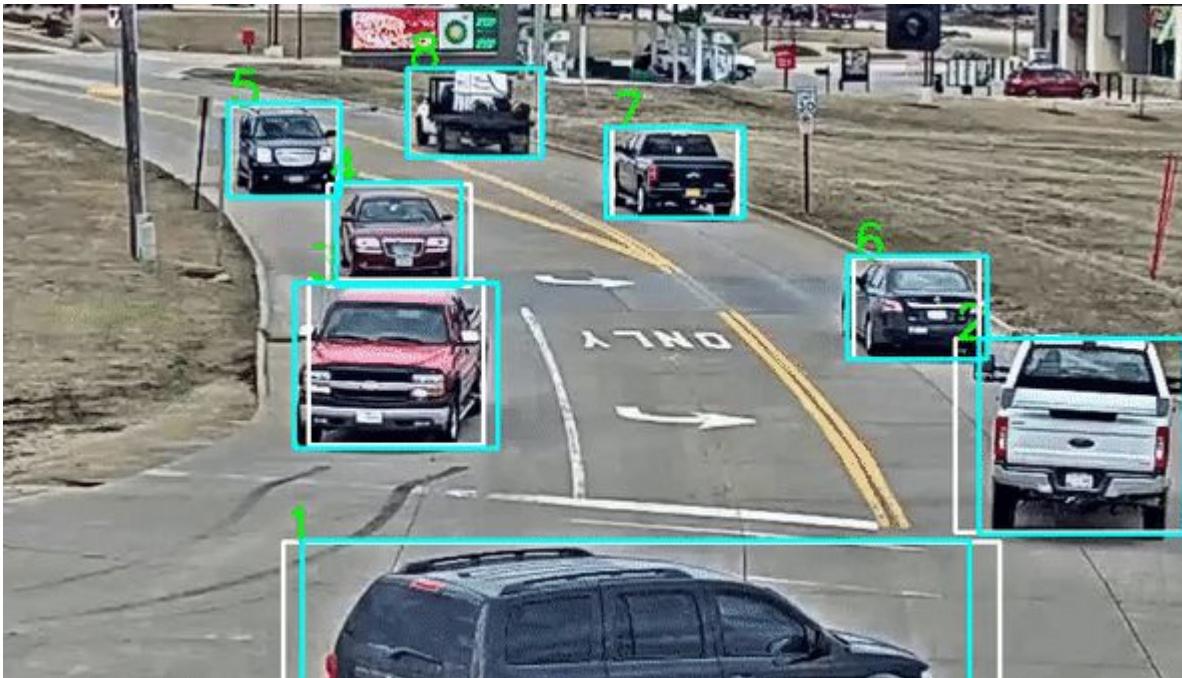
Single Object

Multiple Objects

# DETECTION VS SEGMENTATION



# OBJECT TRACKING



- Where did it go in the next frame ?
- Prediction of changing positions, mapping the correct object identity

# COMPUTER VISION USING DEEP LEARNING - CHALLENGES



- Classification - Dog or Not Dog ?
- Input Image
  - Size - 1024 X 1024 pixels
  - Colors - 3 channels (RGB)
- Input Layer Dimension
  - 3M X 1 vector
- First Layer Units
  - 1000 - reasonable considering complexity
- Dimension of Weight Matrix
  - 3M X 1000 matrix
  - 3 Billion elements - not sparse  
!!!
  - 32 bit float => 12GB MEM just the first layer
  - Not feasible
- Solution - Convolution

# MOTIVATION FROM DIGITAL IMAGE PROCESSING FILTERS

# WHAT IS CONVOLUTION OPERATION ?

0	0	0	0	0	0	
0	105	102	100	97	96	
0	103	99	103	101	102	
0	101	98	104	102	100	
0	99	101	106	104	99	
0	104	104	104	100	98	

Image Matrix

0	-1	0
-1	5	-1
0	-1	0

Kernel Matrix

320					

Output Matrix

$$\begin{aligned}
 & 0 * 0 + 0 * -1 + 0 * 0 \\
 & + 0 * -1 + 105 * 5 + 102 * -1 \\
 & + 0 * 0 + 103 * -1 + 99 * 0 = 320
 \end{aligned}$$

# EDGE DETECTION - SOBEL



-1	0	+1
-2	0	+2
-1	0	+1



+1	+2	+1
0	0	0
-1	-2	-1



# NOISE REMOVAL - GAUSSIAN LPF

 $\frac{1}{16}$ 

1	2	1
2	4	2
1	2	1



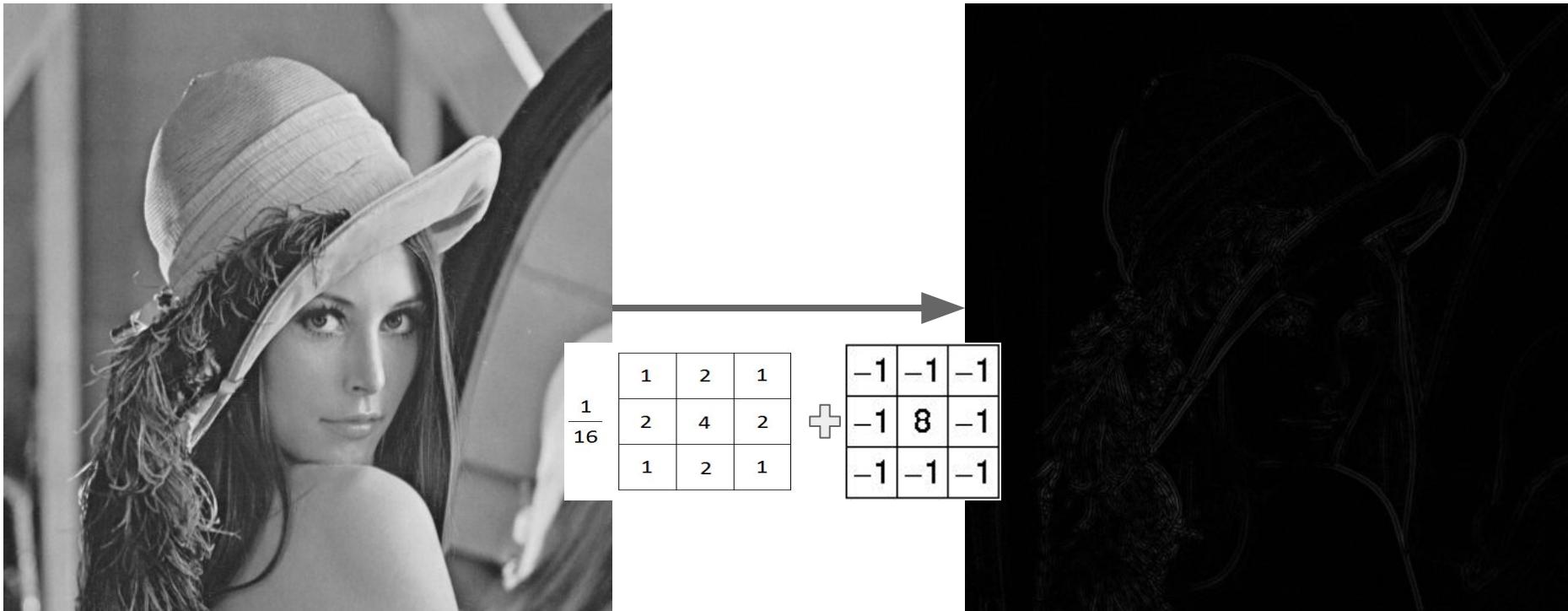
# SHARP FEATURE DETECTION - LAPLACIAN HPF



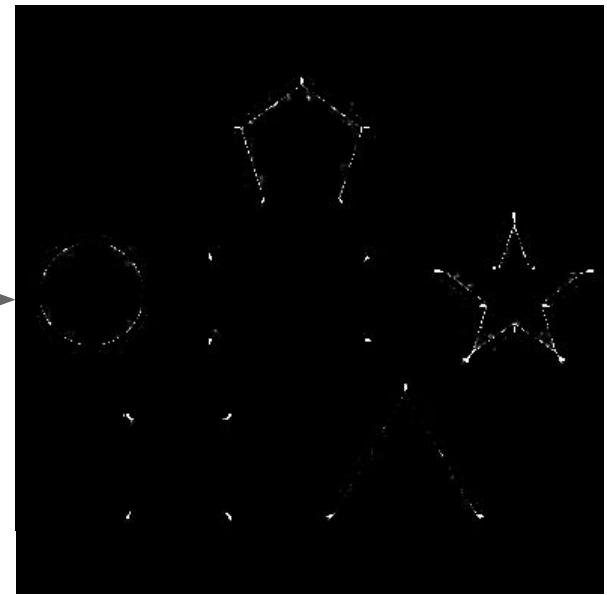
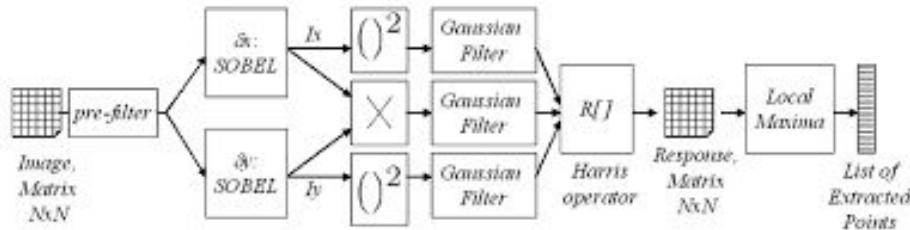
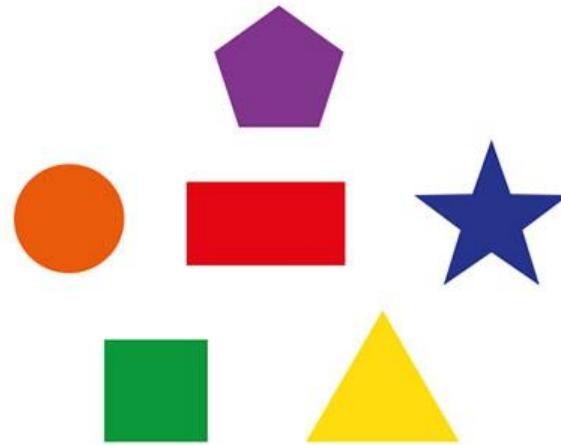
$$\begin{array}{ccc} -1 & -1 & -1 \\ -1 & 8 & -1 \\ -1 & -1 & -1 \end{array}$$



# TWO STAGE EDGE DETECTION - LAPLACIAN OF GAUSSIAN



# HARRIS CORNER DETECTION (6 STAGE FILTER)

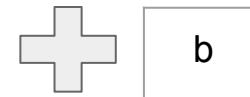


# HAND PICKED VS DATA DRIVEN FILTER DESIGN

- What features are be computed ?
- What should be the best filter weights ?

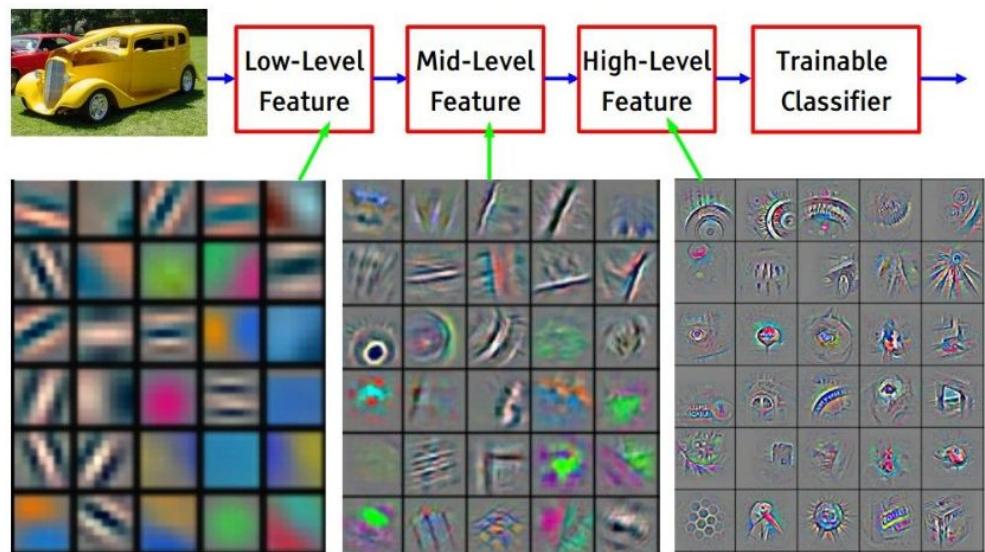
Use training to learn the filter coefficients at each stage

$W_{11}$	$W_{12}$	$W_{13}$	....	$W_{1n}$
$W_{21}$	$W_{22}$	$W_{23}$	....	$W_{2n}$
$W_{31}$	$W_{32}$	$W_{33}$	....	$W_{3n}$
⋮	⋮	⋮		⋮
$W_{n1}$	$W_{n2}$	$W_{n3}$	....	$W_{nn}$



# WHAT DOES CNN LAYERS DO ?

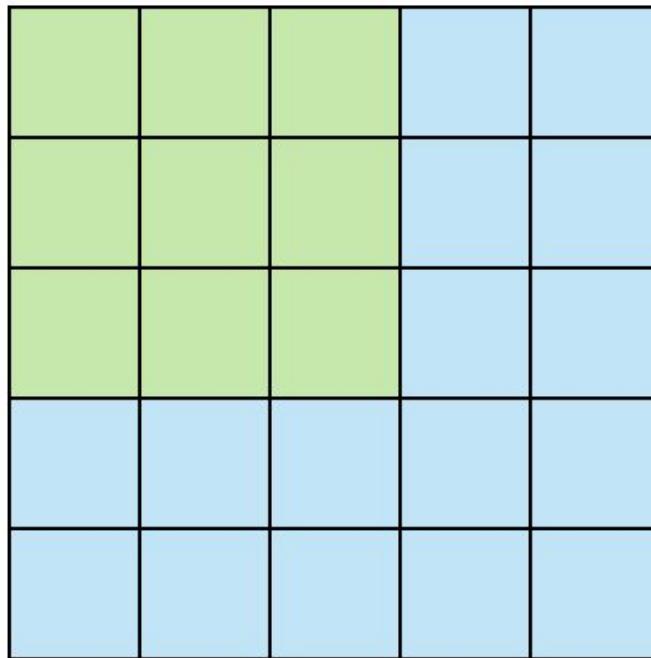
- Early Layers
  - Detects edges, corners, blobs etc
- Middle Layers
  - Detect parts of object
- Later Layers
  - Detects near complete sub-objects
- Last Layer
  - Final classification output



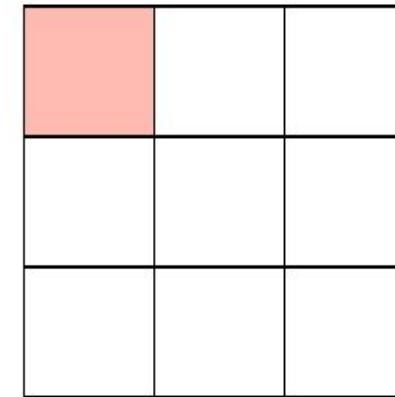
# IMPLEMENTATION OF CONVOLUTION

# VALID CONVOLUTION

Image:  $n \times n$   
Filter:  $f \times f$



Stride 1

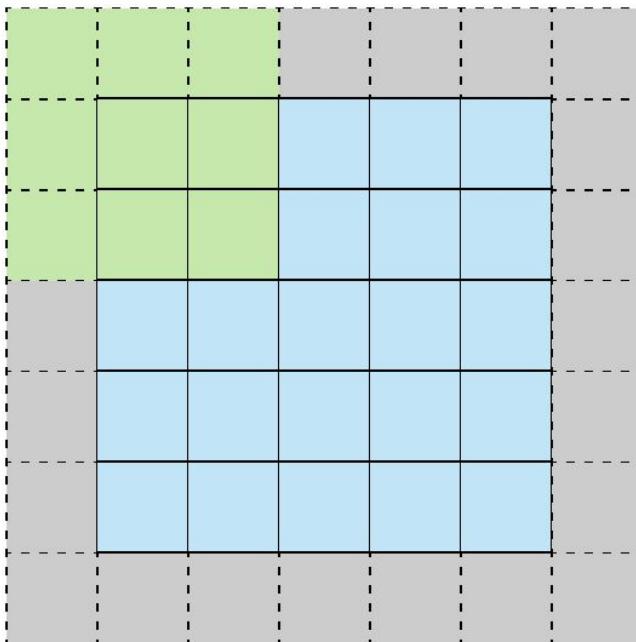


Output:  $(n-f+1) \times (n-f+1)$

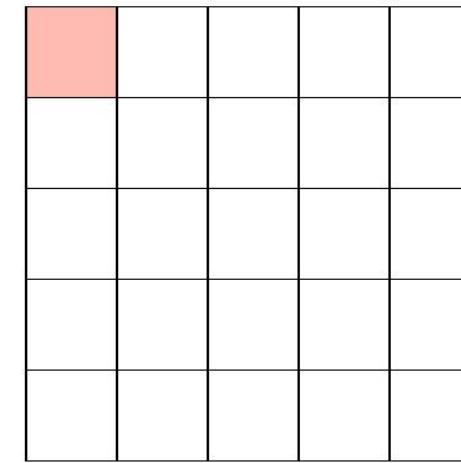
Feature Map

# SAME CONVOLUTION - PADDING

Image:  $n \times n$   
Filter:  $f \times f$   
Padding:  $p$



Stride 1 with Padding



Output:  $(n+2p-f+1) \times (n+2p-f+1)$

To maintain size:  $p = (f-1)/2$

Feature Map

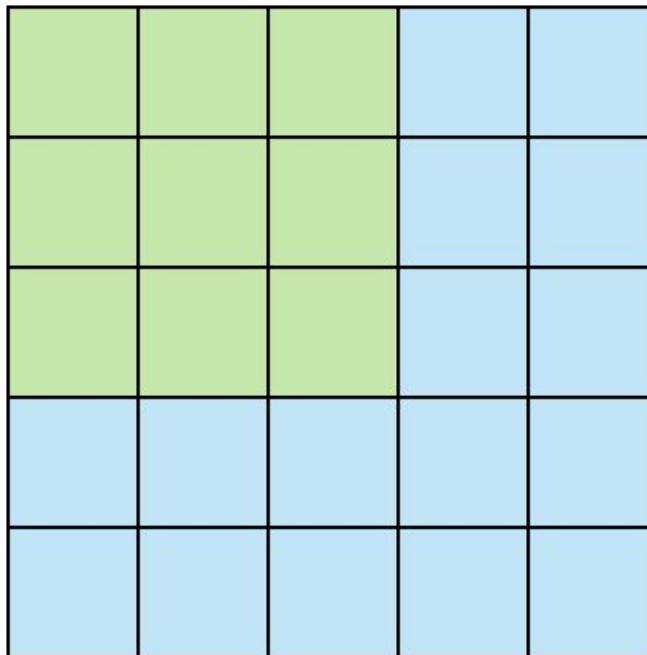
# STRIDES

Image:  $n \times n$

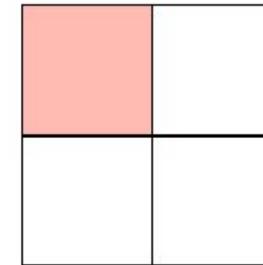
Filter:  $f \times f$

Padding:  $p$

Stride:  $s$



Stride 2

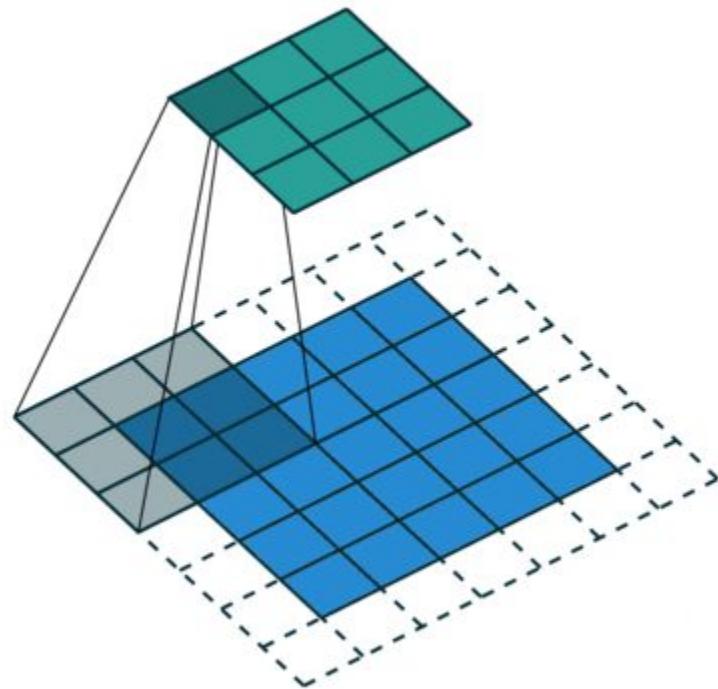


Output:  $((n + 2p - f) / s + 1) \times ((n + 2p - f) / s + 1)$

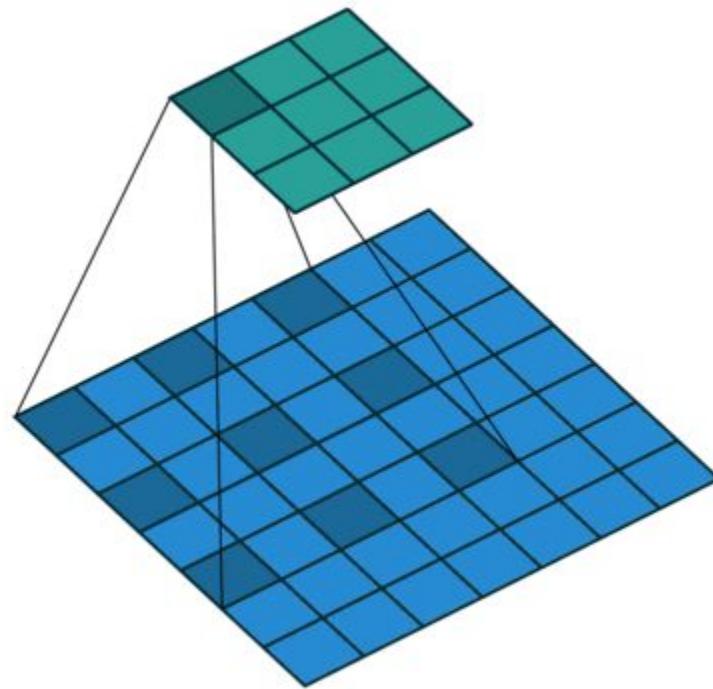
Avoid fractions, use floor function if cannot

Feature Map

# DILATED CONVOLUTION

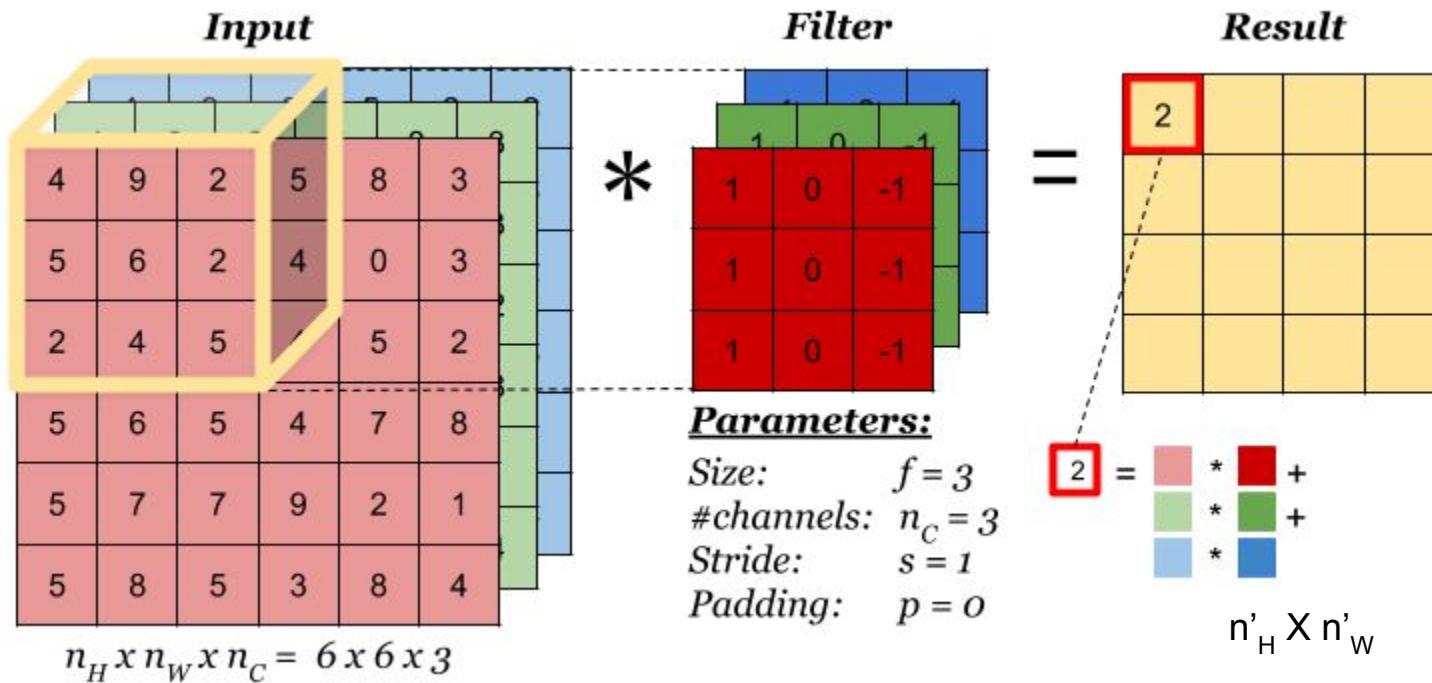


Standard Convolution ( $l=1$ )



Dilated Convolution ( $l=2$ )

# MULTI-CHANNEL CONVOLUTION



# MULTI-CHANNEL CONVOLUTION

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

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

+

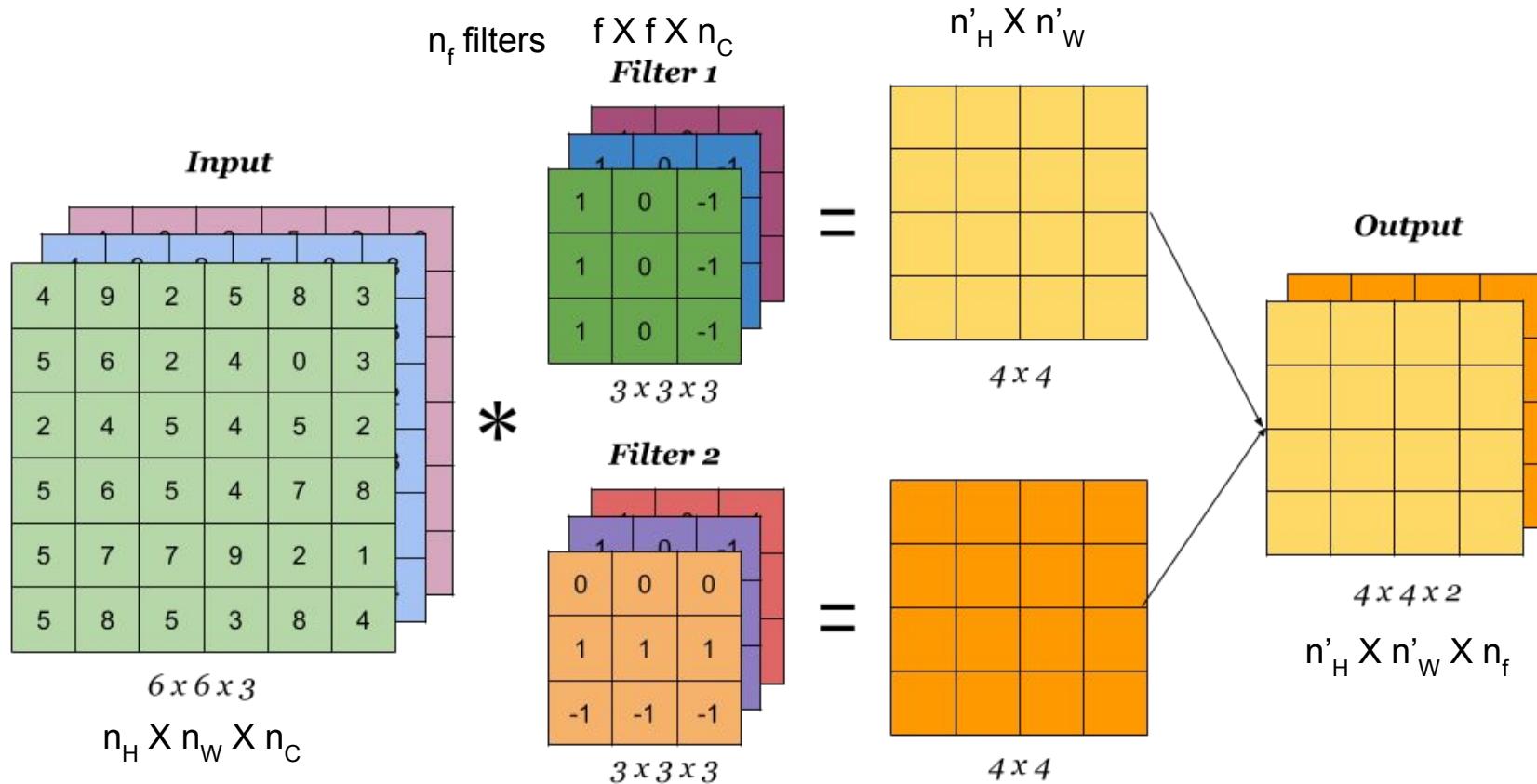
$$+ 1 = -25$$

Bias = 1

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

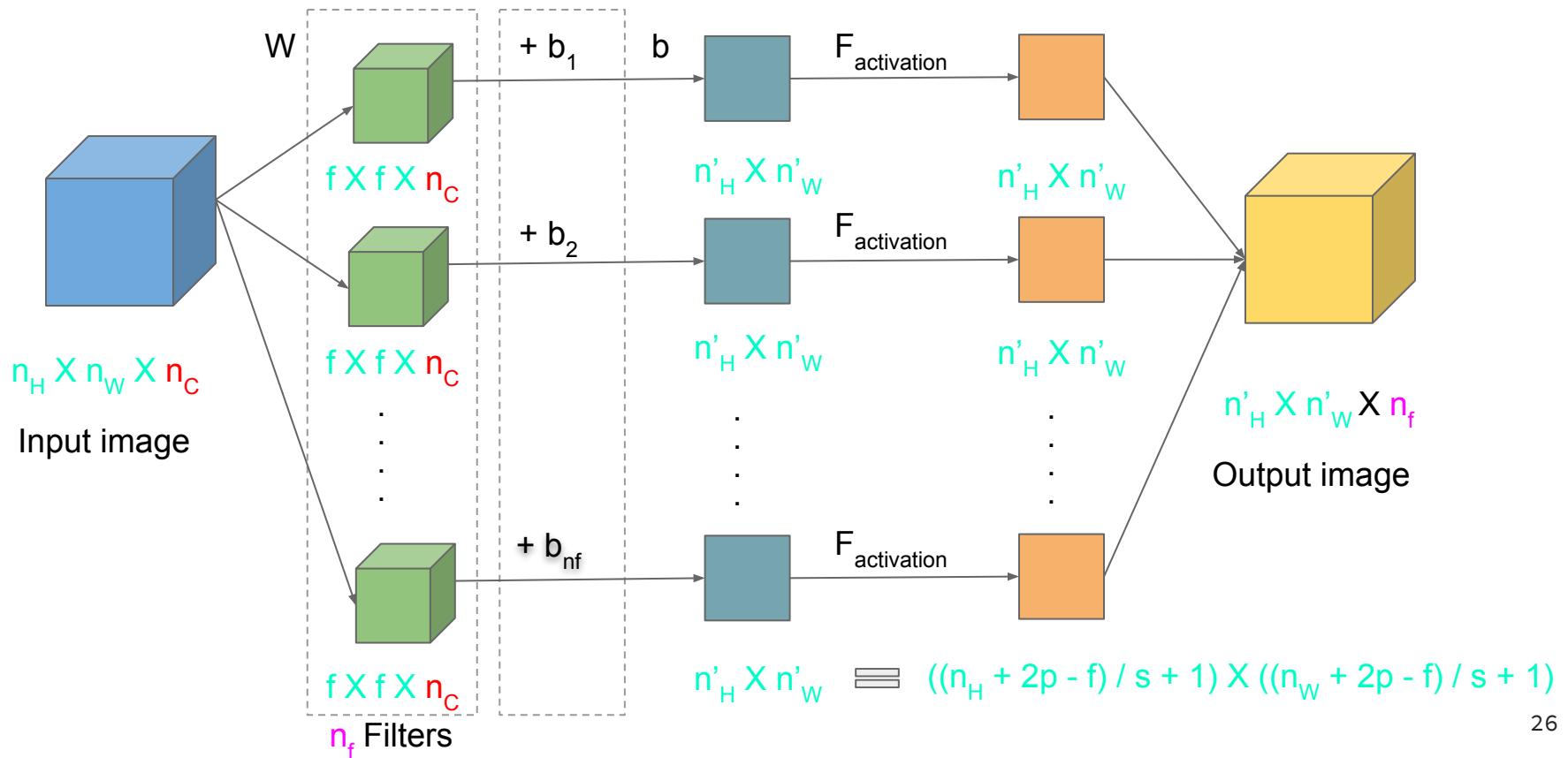
Output

# MULTI-FILTER CONVOLUTION



# BUILDING A CONVOLUTIONAL NEURAL NETWORK

# SINGLE CONVOLUTIONAL LAYER

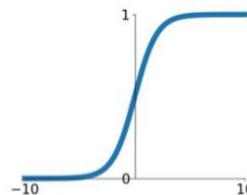


# ACTIVATION FUNCTIONS

Adds non-linearity to the network

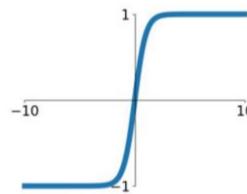
## Sigmoid

$$\sigma(x) = \frac{1}{1+e^{-x}}$$



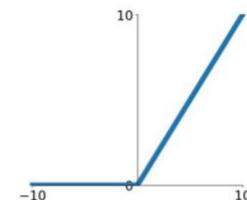
## tanh

$$\tanh(x)$$



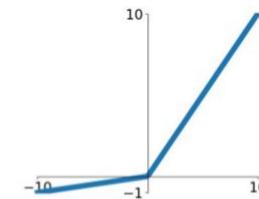
## ReLU

$$\max(0, x)$$



## Leaky ReLU

$$\max(0.1x, x)$$

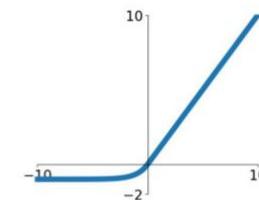


## Maxout

$$\max(w_1^T x + b_1, w_2^T x + b_2)$$

## ELU

$$\begin{cases} x & x \geq 0 \\ \alpha(e^x - 1) & x < 0 \end{cases}$$



# POOLING LAYERS

Enlarges the receptive field of subsequent layers

Max Pool

2	3	1	9
4	7	3	5
8	2	2	2
1	3	4	5



7	9
8	5

Max-Pool with a  
2 by 2 filter and  
stride 2.

Select the most important feature  
(almost exclusively used nowadays)

Apply on each channel independently

Average Pool

2	3	1	9
4	7	3	5
8	2	2	2
1	3	4	5



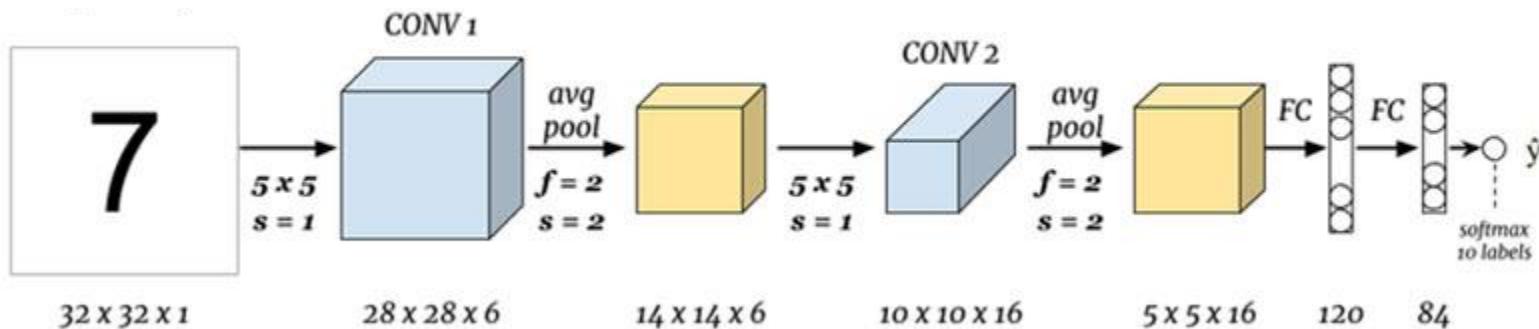
4	4.5
3.25	3.25

Average Pool with  
a 2 by 2 filter and  
stride 2.

Select the average of incoming features  
(earlier research)

# LENET-5 MODEL

$n_H$ ,  $n_W$  goes down,  $n_C$  goes up along the depth. Why?

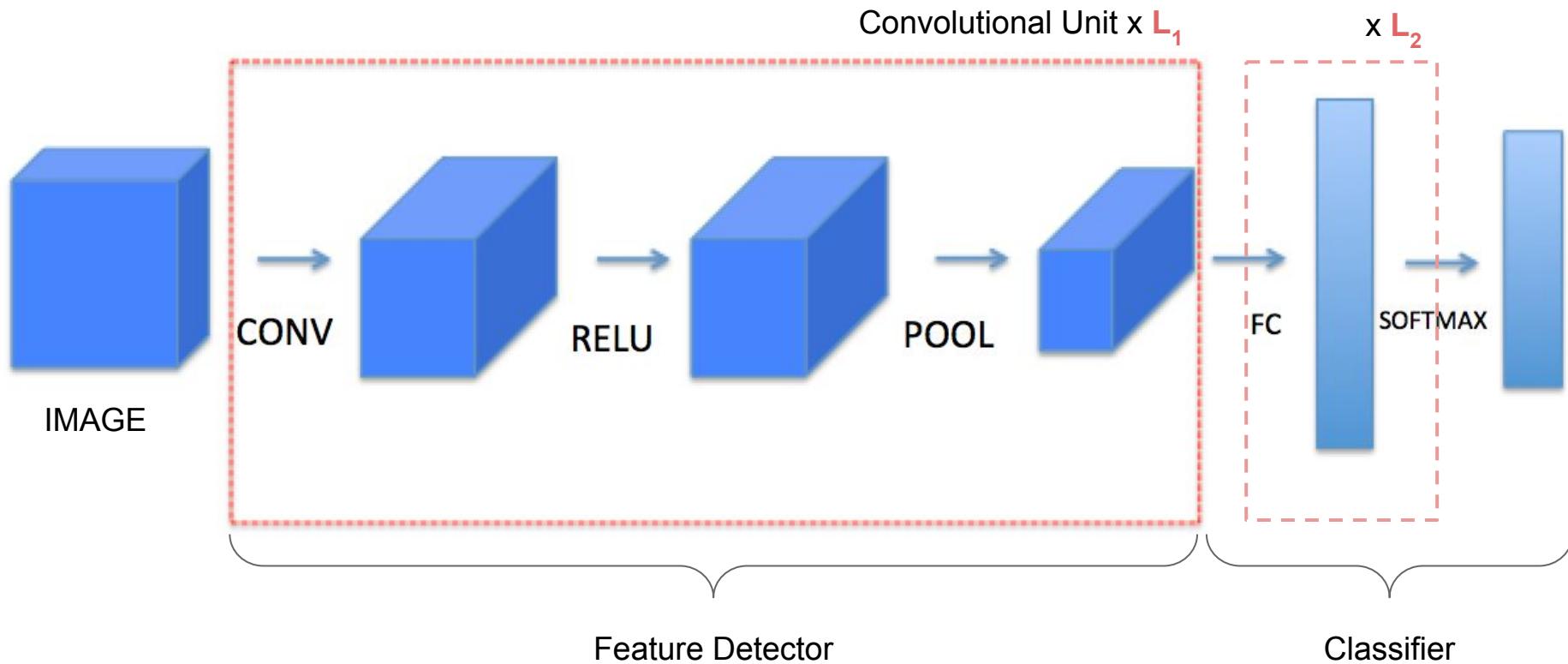


tanh activation for all layers  
softmax activation for last (output) layer

~ 61K parameters

[LeCun et al., 1998, Gradient-based learning applied to document recognition]

# A COMPLETE CONVOLUTIONAL NEURAL NETWORK



# HYPERPARAMETERS OF A CLASSIC CNN

- Input image size ( $n_h$ ,  $n_w$ ,  $n_c$ )
- Number of Convolutional Units ( $L_1$ )
  - Size of filters in each layer ( $f$ )
  - Padding ( $p$ )
  - Stride ( $s$ )
  - Number of filters in each layer ( $n_f$ )
  - Size of pooling filter ( $f_p$ )
  - Stride of pooling filter ( $s_p$ )
- Number of fully connected layers ( $L_2$ )
  - Number of nodes in each FC layer ( $n_u$ )

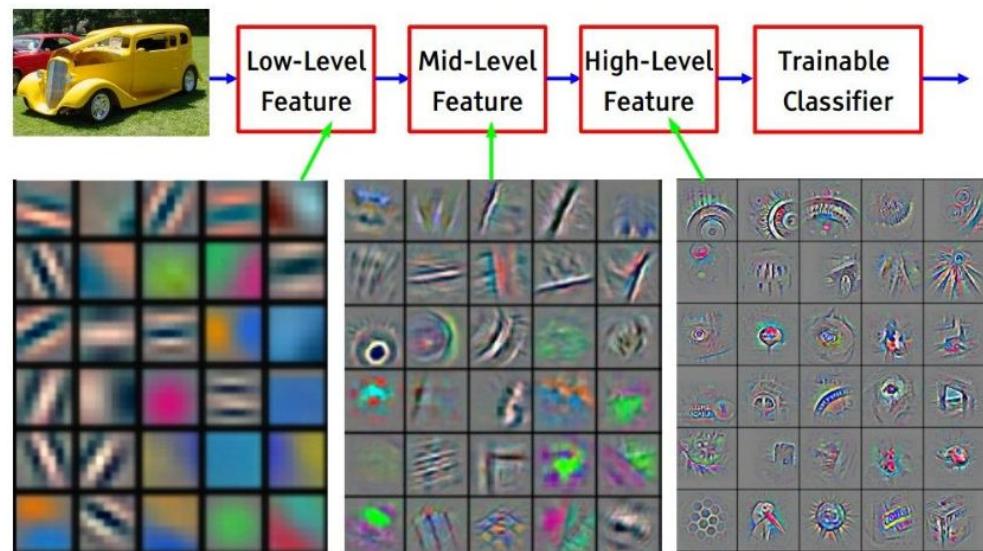
# SUITABILITY OF CONVOLUTIONS TO COMPUTER VISION

# WHY CONVOLUTIONS OVER DENSE NEURAL NETWORKS ?

- Parameter Sharing => **Less Computations**
  - Each filter detects a different feature
  - Same feature detector useful in one part of image most likely useful in another part as well
- Sparsity of Connections => **Robustness to Localization**
  - Each output pixel is only connected to a limited number of input values
  - Relevant content may be contained only in a part of an image and not span the whole image

# WHAT DOES CNN LAYERS DO ?

- Early Layers
  - Detects edges, corners, blobs etc
- Middle Layers
  - Detect parts of object
- Later Layers
  - Detects near complete sub-objects
- Last Layer
  - Final classification output



# FEATURE VISUALIZATION OF CNN LAYERS



**Edges** (layer conv2d0)

**Textures** (layer mixed3a)

**Patterns** (layer mixed4a)

**Parts** (layers mixed4b & mixed4c)

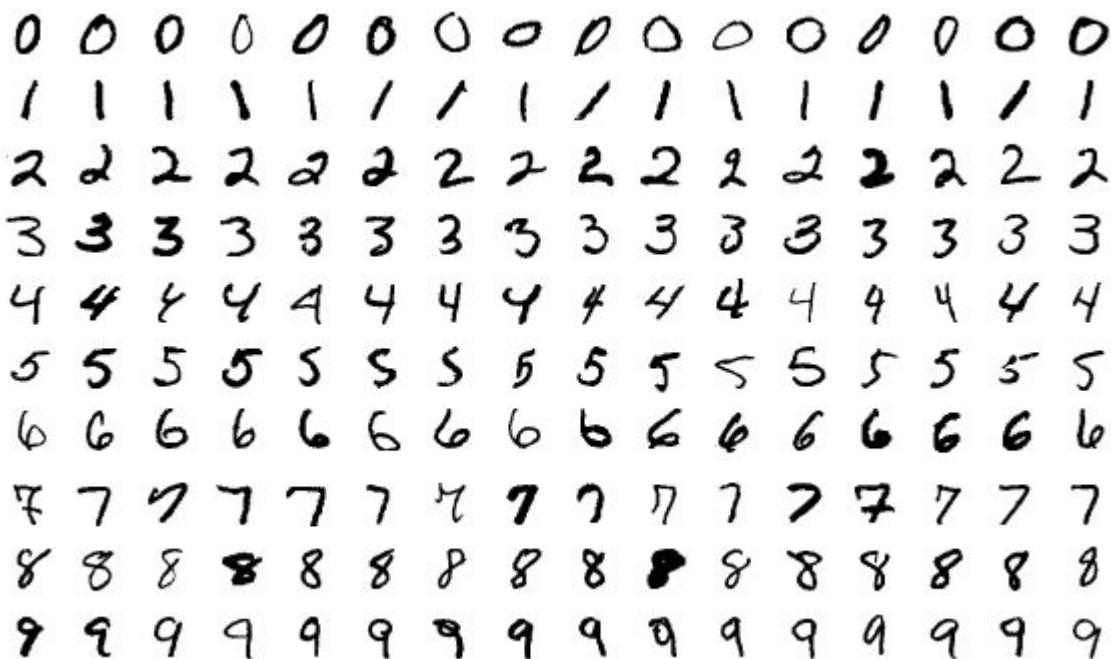
**Objects** (layers mixed4d & mixed4e)

Feature visualization allows us to see how GoogLeNet [1], trained on the ImageNet [2] dataset, builds up its understanding of images over many layers. Visualizations of all channels are available in the [appendix](#).

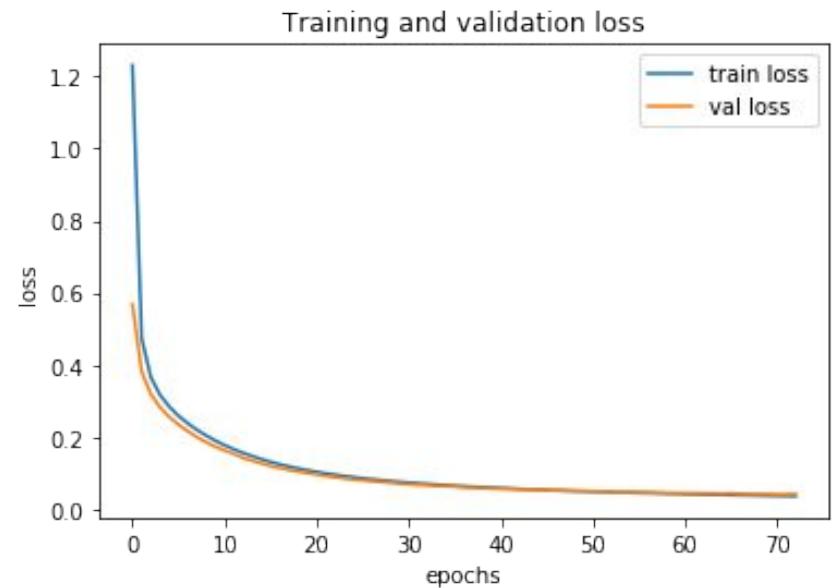
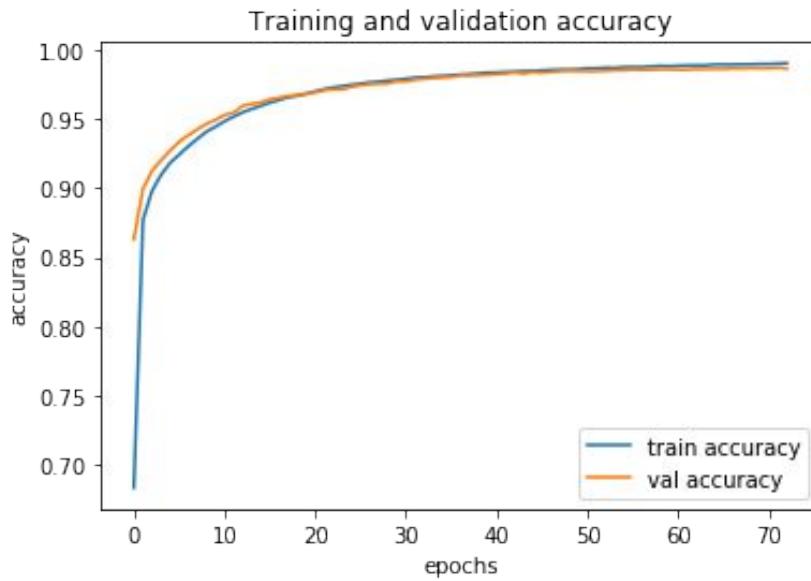
# CASE STUDY - IMAGE CLASSIFICATION USING CNNs

# MNIST

- MNIST
  - handwritten number classification set (0-9)
  - 60K train set, 10K test set - 28X28 grayscale



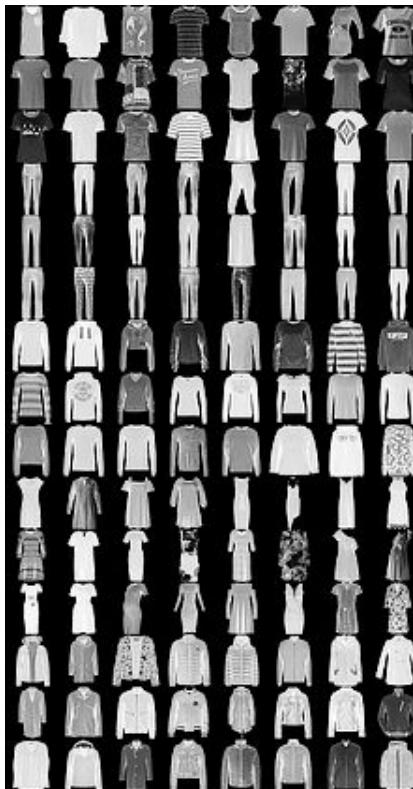
# MNIST DATASET CLASSIFICATION USING LENET5 CNN



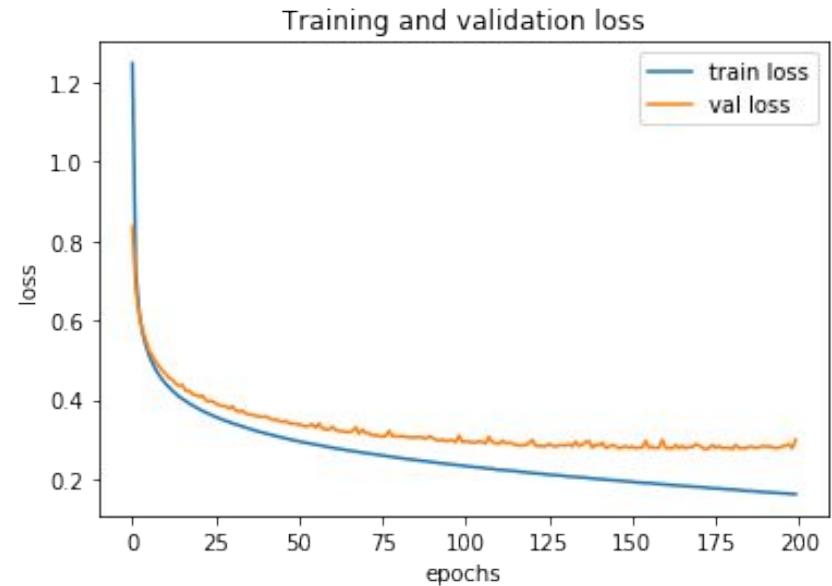
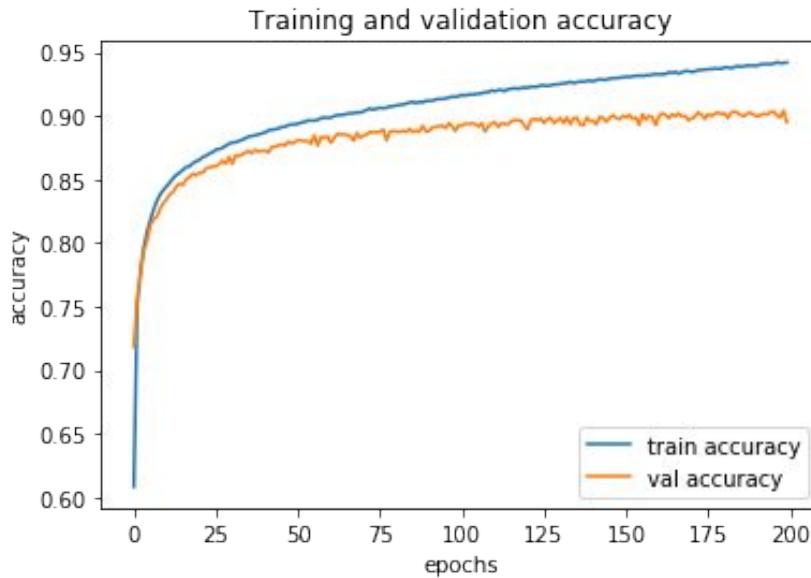
- Training accuracy reached **99%**. Validation accuracy reached **98.5%**
- Model is a **good fit**
  - Very good training and validation accuracy
  - Validation loss tracks training loss

# FASHION-MNIST

- Fashion-MNIST
  - image object classification set - **t-shirt/top, trouser, pullover, dress, coat, sandal, shirt, sneaker, bag, ankle boot**
  - 60K train set, 10K test set - 28X28 grayscale



# FASHION MNIST DATASET CLASSIFICATION USING LENET5 CNN



- Training accuracy reached **95%**. Validation accuracy reached **90%**
- Model is **average with a tendency of slight overfitting**
  - Good training accuracy and acceptable validation accuracy
  - Validation loss diverges from training loss after initial few epochs

# CIFAR10

- CIFAR10
  - image object classification set - **airplane, car, bird, cat, deer, dog, frog, horse, ship, truck**
  - 50K train set, 10K test set - 32X32 rgb color

airplane



automobile



bird



cat



deer



dog



frog



horse



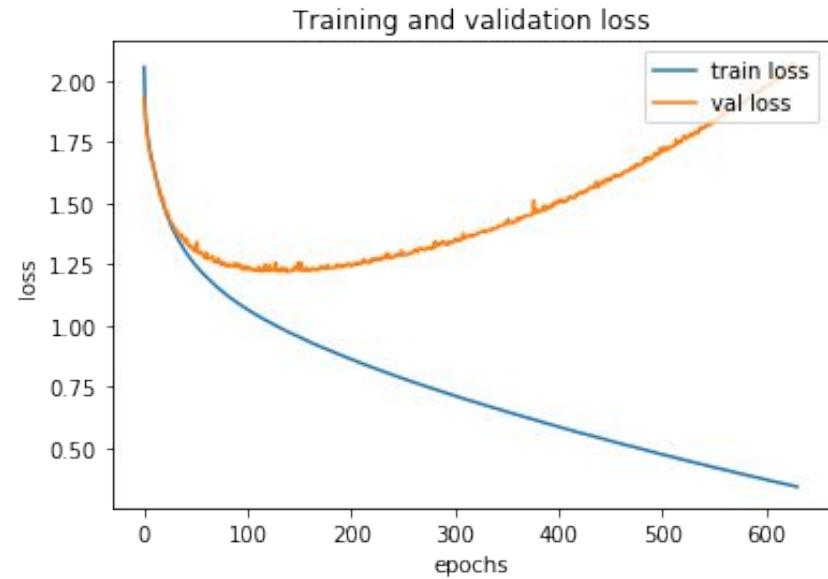
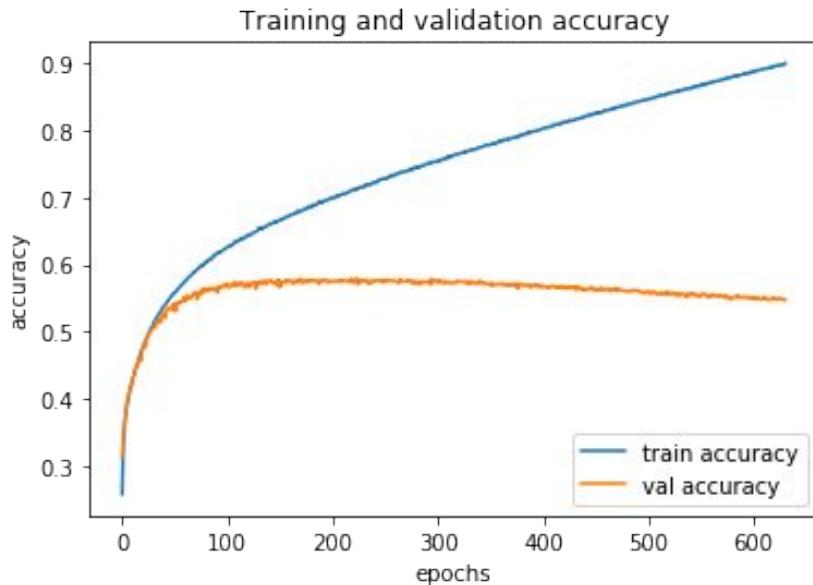
ship



truck



# CIFAR-10 DATASET CLASSIFICATION USING LENET5 CNN

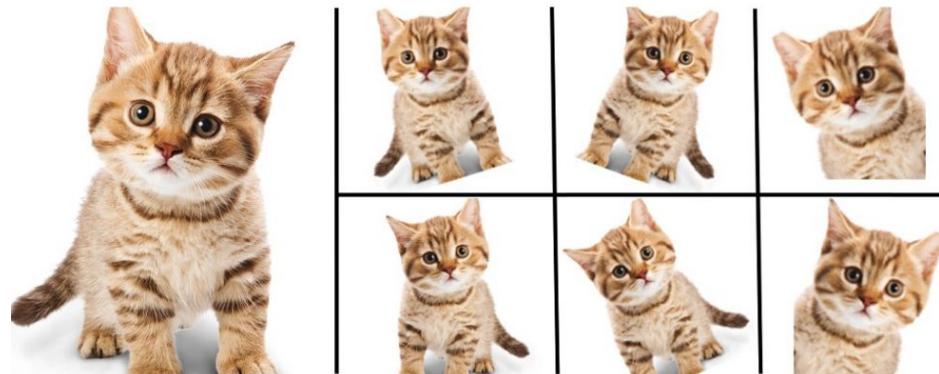


- Training accuracy reached **90%**. Validation accuracy reached **55%**
- Model is **poor with extreme overfitting**
  - Good training accuracy and very bad validation accuracy
  - Validation loss hugely diverges from training loss after a number of epochs

# PRACTICAL TRAINING TECHNIQUE I - DATA AUGMENTATION

- Overfitting
  - Insufficient data
    - Trains well
  - Training data does not contain all potential variations
    - Hence poor validation/test results
- Gathering more data is expensive and time consuming
  - Labeling lots of new data is extremely tedious task
- Solution: Simulate variations
  - Modify training images and add to the training

# DATA AUGMENTATION



- Mirroring
  - Horizontal flip
  - Vertical flip (rare)
- Zoom +/-
- Translation
- Rotation
- Shearing
- Distortions
  - Local Warping
- Brightness
- Color Shifting
  - Change tint

# ADVANCED CNN ARCHITECTURES

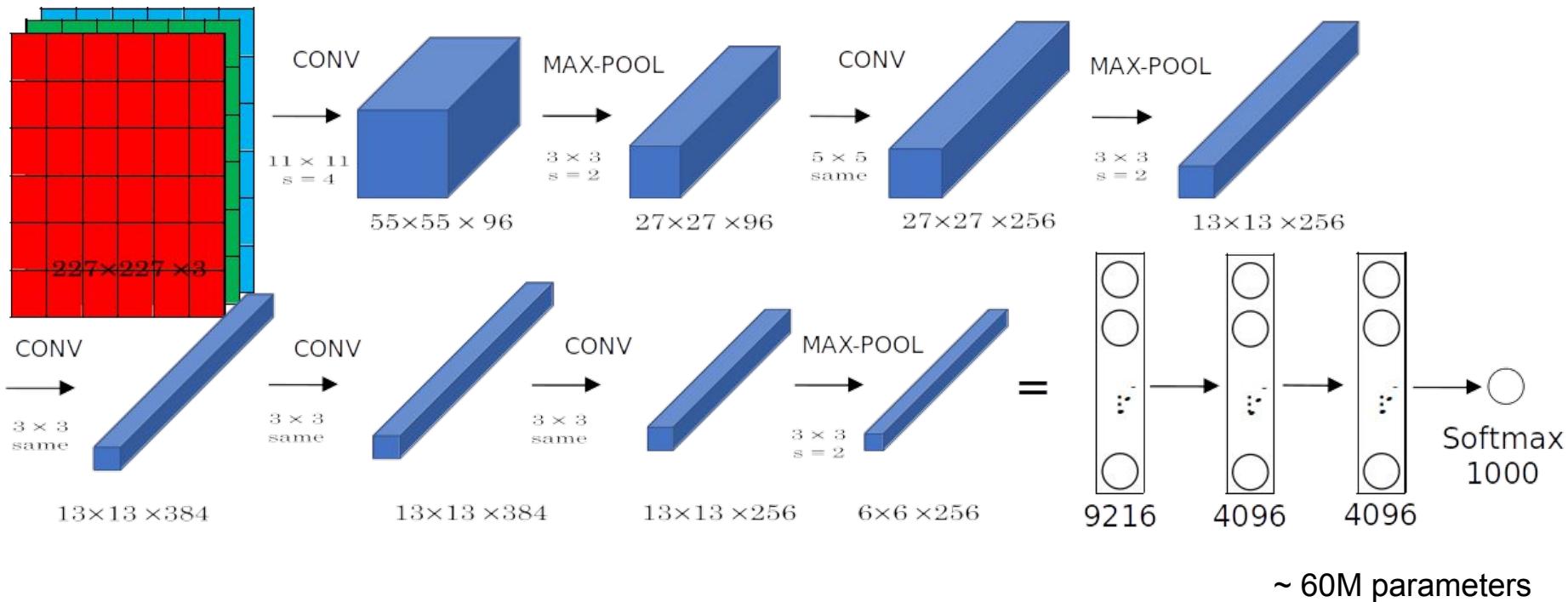
# PROGRESSION OF CNN ARCHITECTURES

Story of Image Net Challenge Winners since 2012

- 2012 - AlexNet, University of Toronto - **83.6%**
  - 5 convolutional with/without max pooling layers, 3 fully connected layers with dropout, **first use of ReLU activation**
- 2013 - VGG-16/19, University of Oxford - **92.7%**
  - 16 -19 layers with **only small 3X3 receptive fields**
- 2014 - GoogleLeNet, Google - **93.3%**
  - 22 layers utilizing **Network-in-Network** inception module
- 2016 - ResNet, Microsoft - **96.4%**
  - 152 layers having **residual blocks** with identity skip connections

PS: Human Top 5 Accuracy ~ **95%**

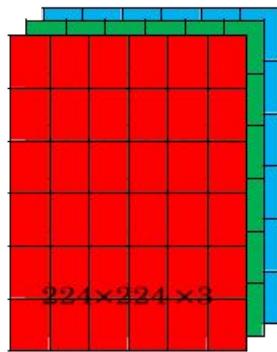
# ALEXNET ARCHITECTURE (9 LAYERS)



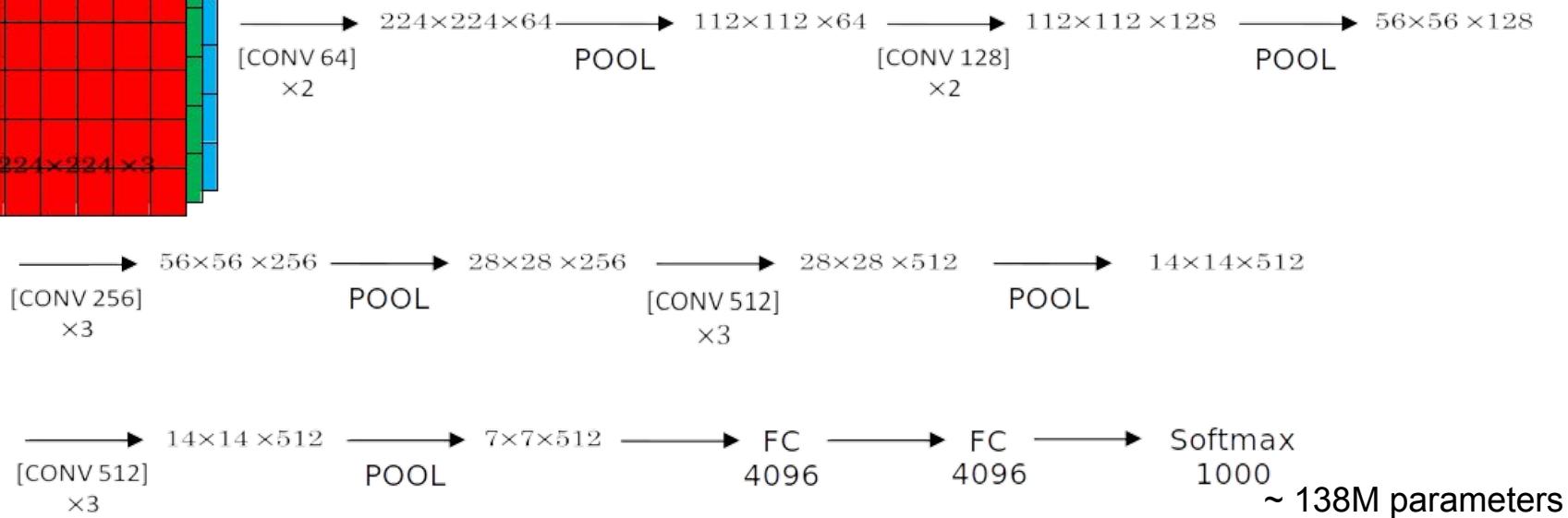
[Krizhevsky et al., 2012, ImageNet classification with deep convolutional neural networks]

# VGG-16

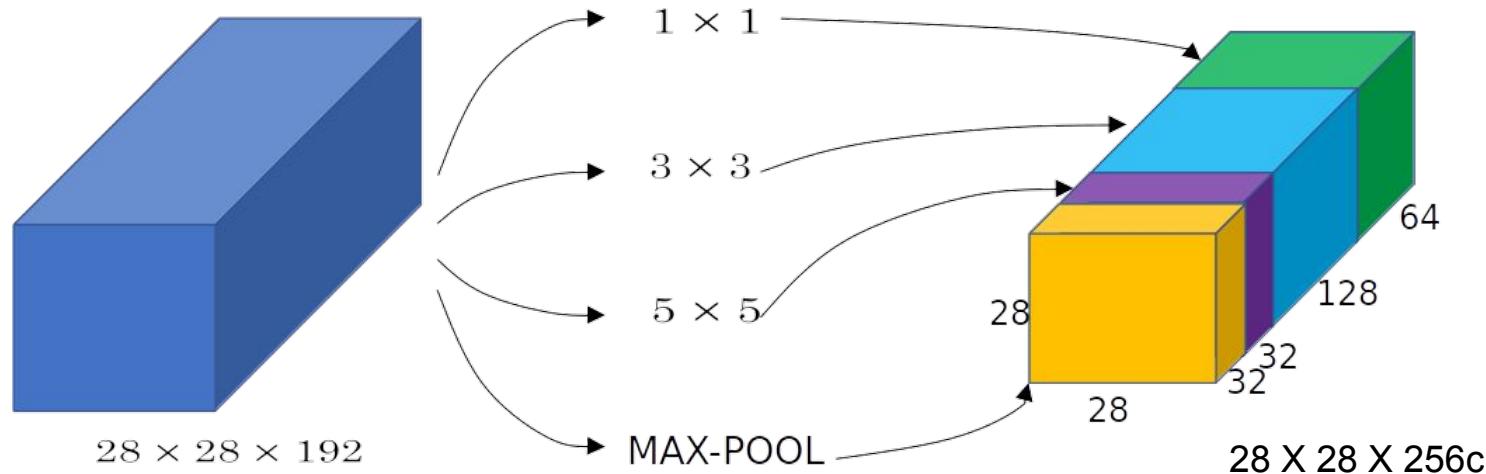
CONV =  $3 \times 3$  filter,  $s = 1$ , same



MAX-POOL =  $2 \times 2$ ,  $s = 2$



# FILTER CONCATENATION AKA NETWORK-IN-NETWORK

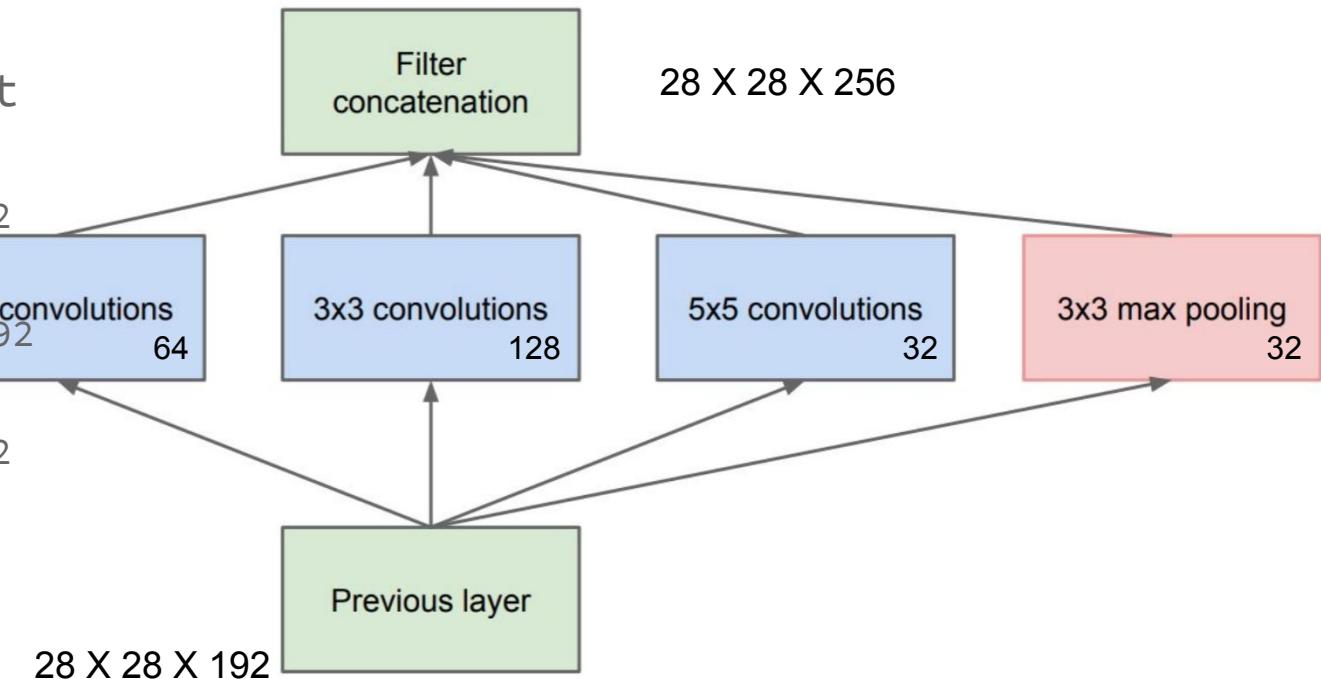


- 1X1 Conv: Effectively a fully connected layer within sub-data (weighted sum across channels with activation)
  - Non-trivially alter  $n_c$  just as Pooling layers shrink  $n_H$ ,  $n_w$

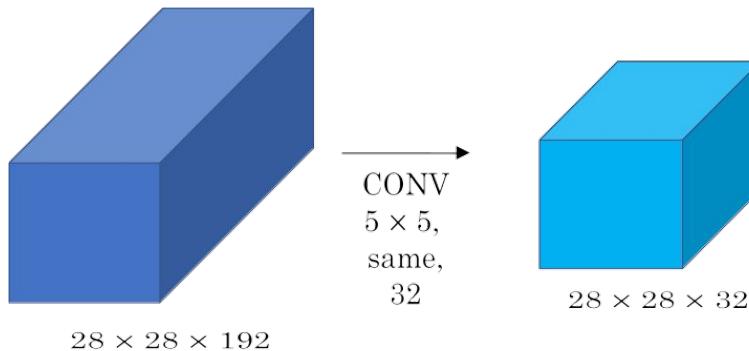
[Lin et al., 2013. Network in Network]

# INCEPTION MODULE

- Computation Cost
- [1X1 Conv]
  - $28 \times 28 \times 64 \times 1 \times 1 \times 192$
- [3X3 Conv]
  - $28 \times 28 \times 128 \times 3 \times 3 \times 192$
- [5X5 Conv]
  - $28 \times 28 \times 32 \times 5 \times 5 \times 192$
- 303M operations



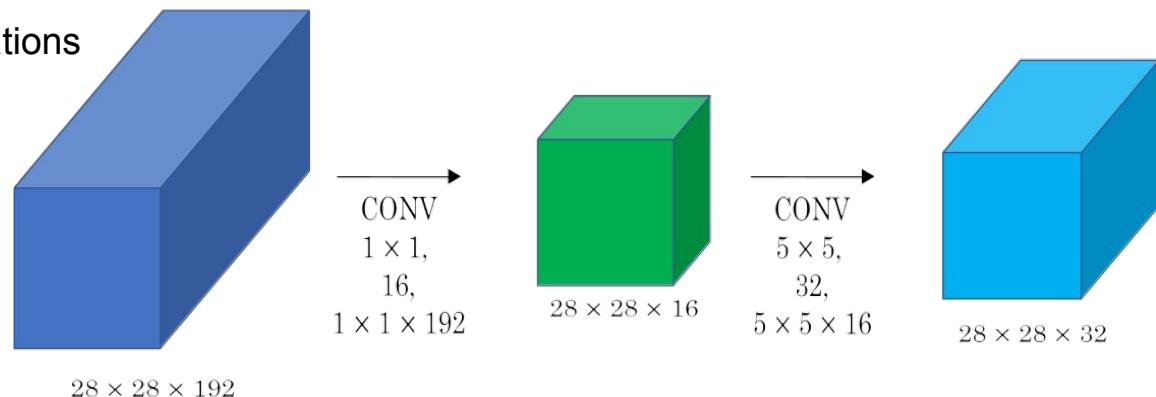
# HIGH COMPUTATIONAL COST - SOLN. IS BOTTLENECK LAYER



Channel lowering within reason doesn't show any noticeable drop in accuracy

$28 \times 28 \times 32 \times 5 \times 5 \times 192 \sim 120M$  operations

10X Speed-up

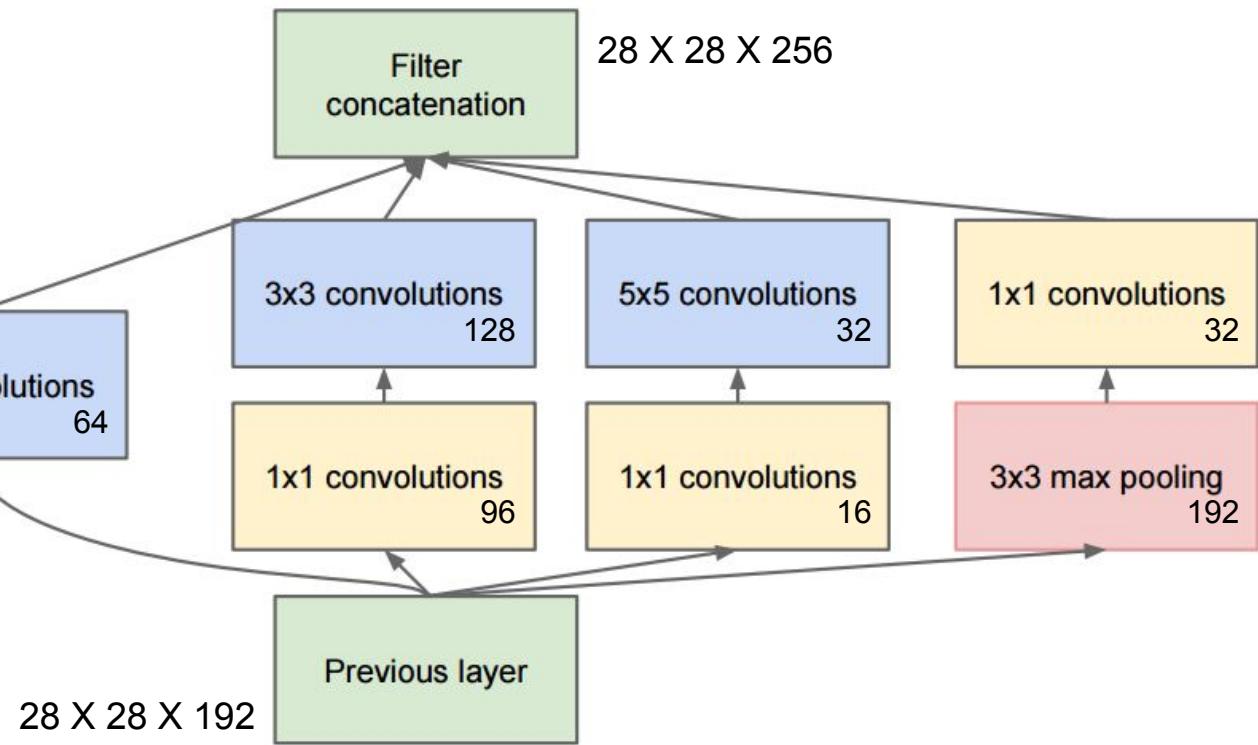


Total  $\sim 12.4M$  Only  $\longrightarrow 28 \times 28 \times 16 \times 192 \sim 2.4M$

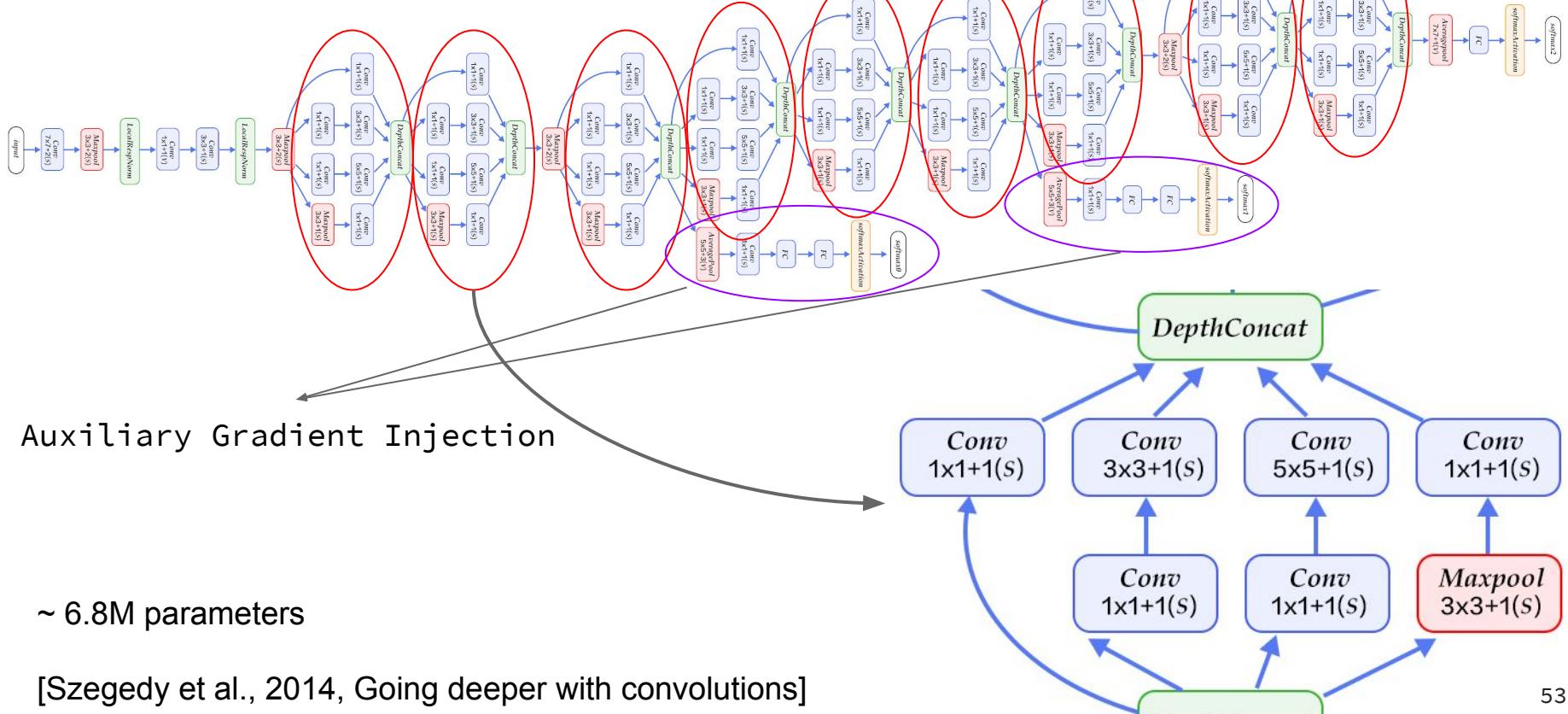
$28 \times 28 \times 32 \times 5 \times 5 \times 16 \sim 10M$

# INCEPTION MODULE - OPTIMIZED

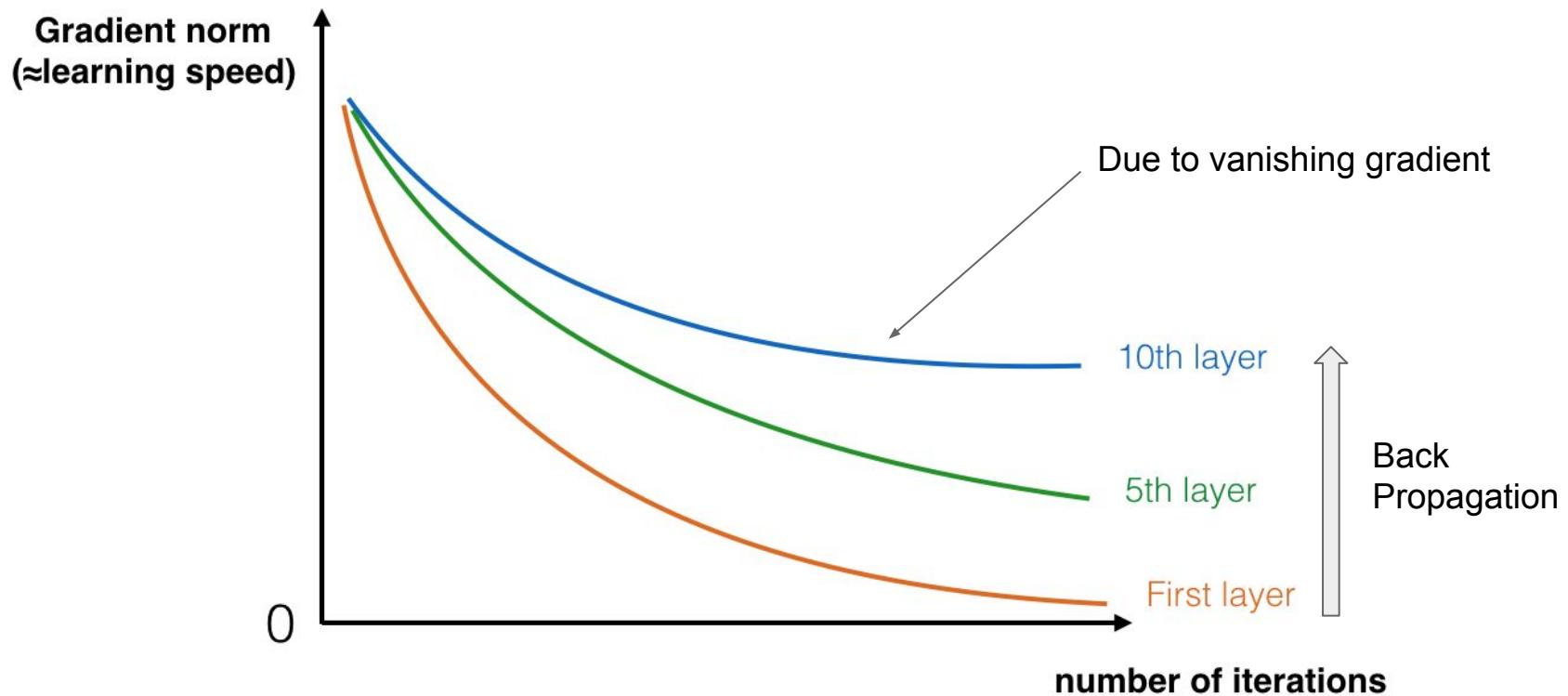
- Computation Cost
- [1X1 Conv]
  - $28 \times 28 \times 64 \times 1 \times 1 \times 192$
- [3X3 Conv]
  - $28 \times 28 \times 96 \times 1 \times 1 \times 192$
  - $28 \times 28 \times 128 \times 3 \times 3 \times 96$
- [5X5 Conv]
  - $28 \times 28 \times 16 \times 1 \times 1 \times 192$
  - $28 \times 28 \times 32 \times 5 \times 5 \times 16$
- [1X1 Conv]
  - $28 \times 28 \times 32 \times 1 \times 1 \times 192$
- 128M operations



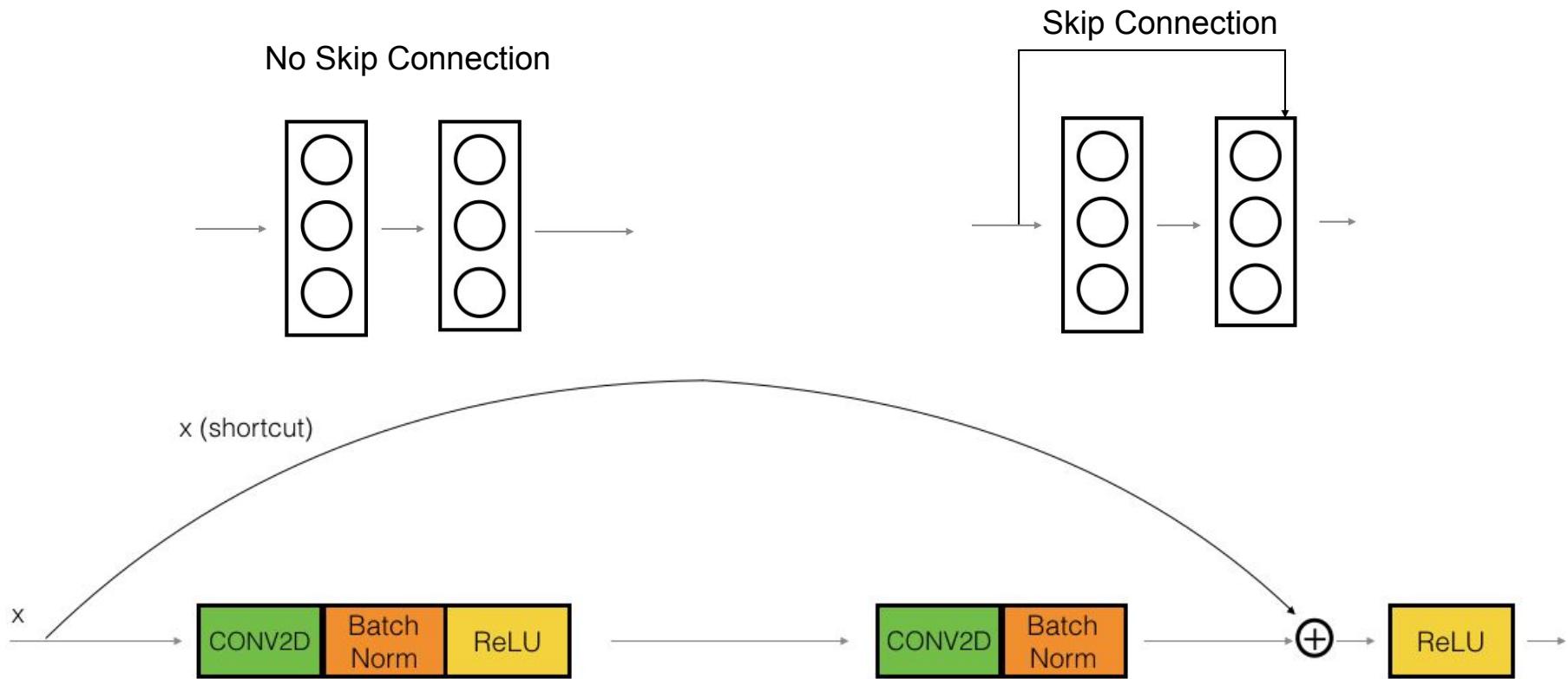
# GOOGLELENET (22 LAYERS)



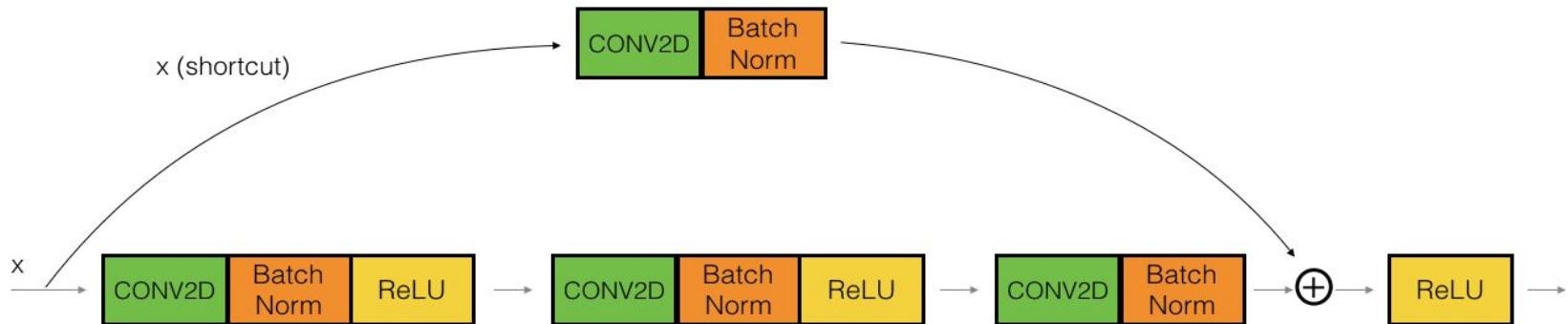
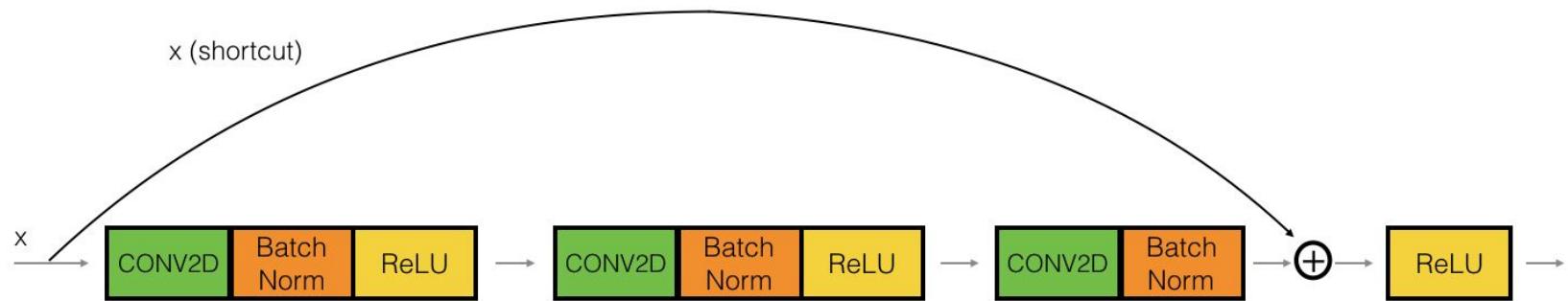
# PROBLEM WITH VERY DEEP NEURAL NETWORKS



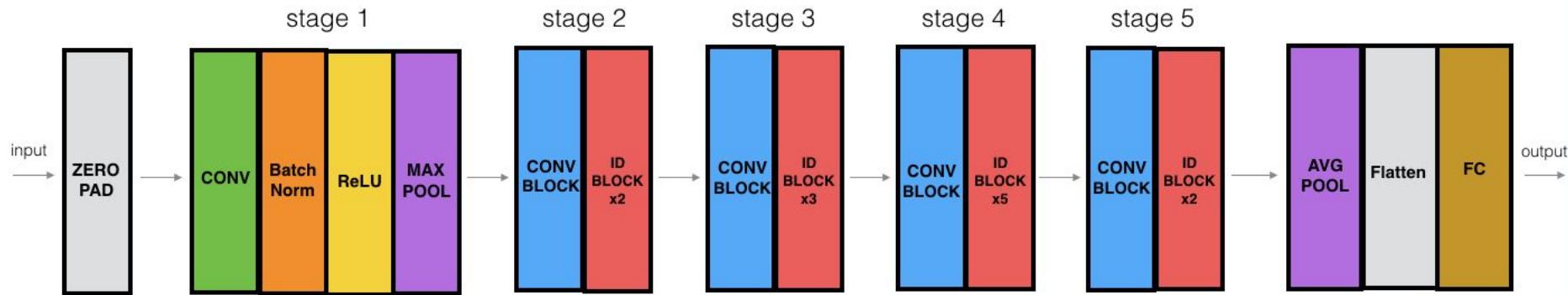
# RESNET IDENTITY BLOCK



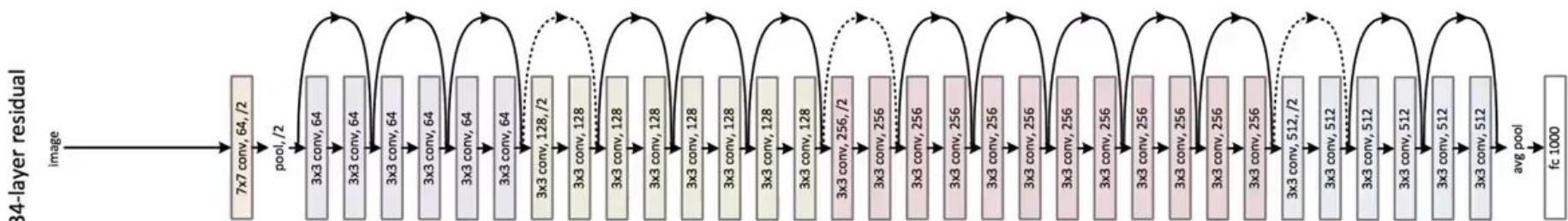
# MORE TYPES OF RESNET IDENTITY BLOCKS



# RESNET 50 MODEL



ResNet



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

# PRACTICAL TRAINING TECHNIQUE II - TRANSFER LEARNING

- Use open source model implementation
  - Architecture of networks often published by authors
  - Can use the whole model or a part of it
- Use pre-trained model weights
  - Use them as it is or as a starting point instead of random initialization then fine tune
- Utilizes prior efforts on time consuming steps
  - Creating extensive training data
  - Training using expensive computation resources such as multiple GPUs
  - Hyperparameter tuning

# TRANSFER LEARNING

- Transfer learning works well in Computer Vision
- Early/mid layers of CNN perform generic image processing and feature extraction
  - Basic constituent building blocks such as edges, corners, blobs, tend to apply across multiple applications
  - ‘Freeze’ these layers, train rest of the network
- Train a small alternative network with the output of ‘frozen’ layers for your own needs
  - Usually Fully Connected
  - May also be Convolutional
- How many layers we ‘unfreeze’ ?
  - Similarities / differences with new application
  - Amount of new data available

THANK YOU