

# CAP 5415: Computer Vision

## Deep (Neural) Networks



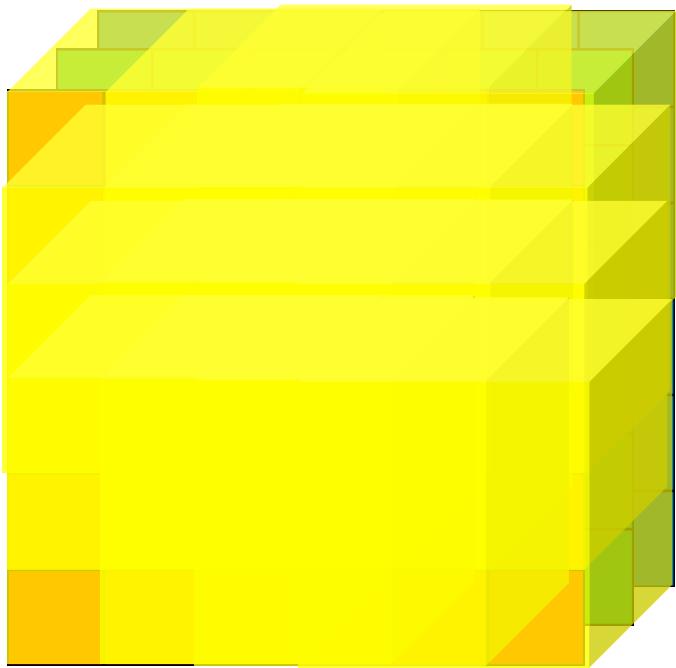
Dr. Sedat Ozer



# Announcements:

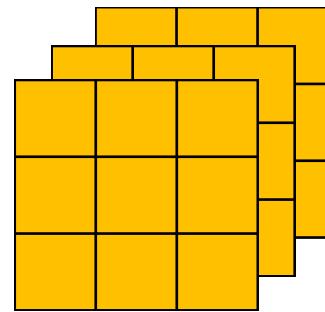
- Midterm:
  - October 4, 2018, Thursday 3pm-4:15pm
    - Classical: Written, in-class exam.
    - Can include content from all the covered material both in class and in the assignments.
    - May or may not have multiple choice questions.
    - You can bring one page hand written cheat-sheet.
- **Non-Existing** Final exam schedule:
  - **Non-existing** Final exam time: **Thursday, December 6, 2018 1:00 PM – 3:50 PM**
  - (scheduled automatically :)

# Convolutions on RGB image

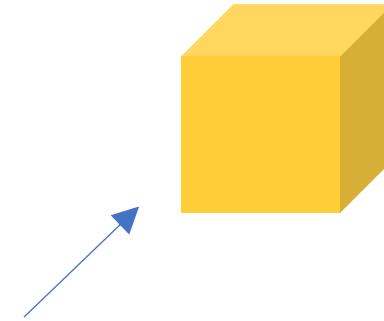


$6 \times 6 \times 3$

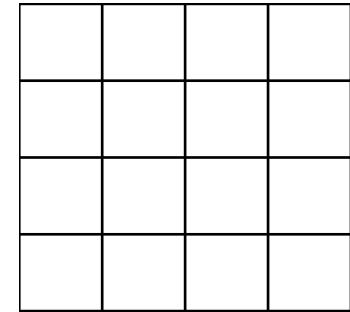
\*



$3 \times 3 \times 3$

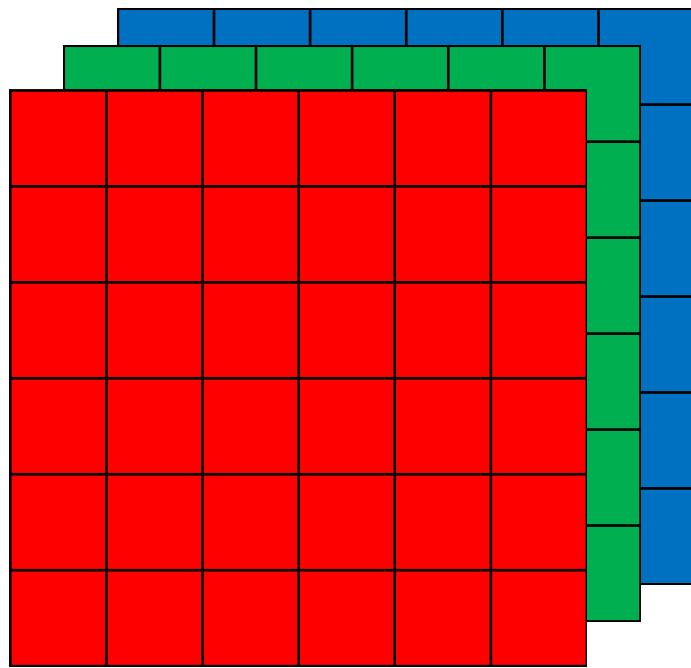


=



$4 \times 4$

# Multiple filters

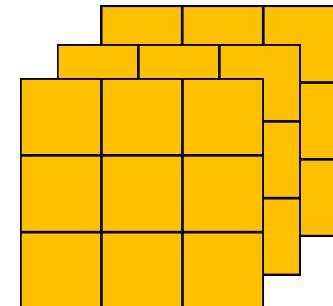


$6 \times 6 \times 3$

Number of channels

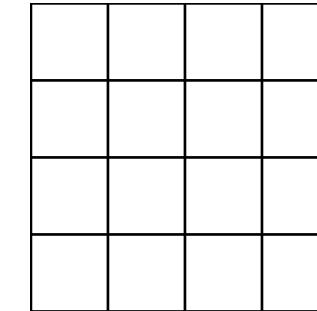
$n \times n \times n_c$

\*



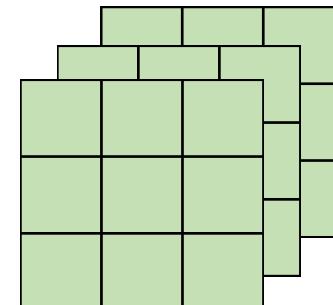
$3 \times 3 \times 3$

=



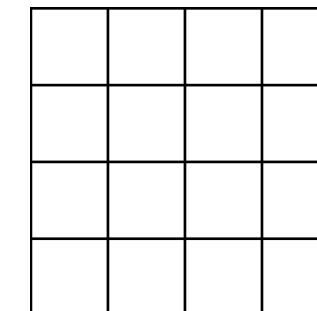
$4 \times 4$

\*



$3 \times 3 \times 3$

=



$4 \times 4$

Number of filters

\*

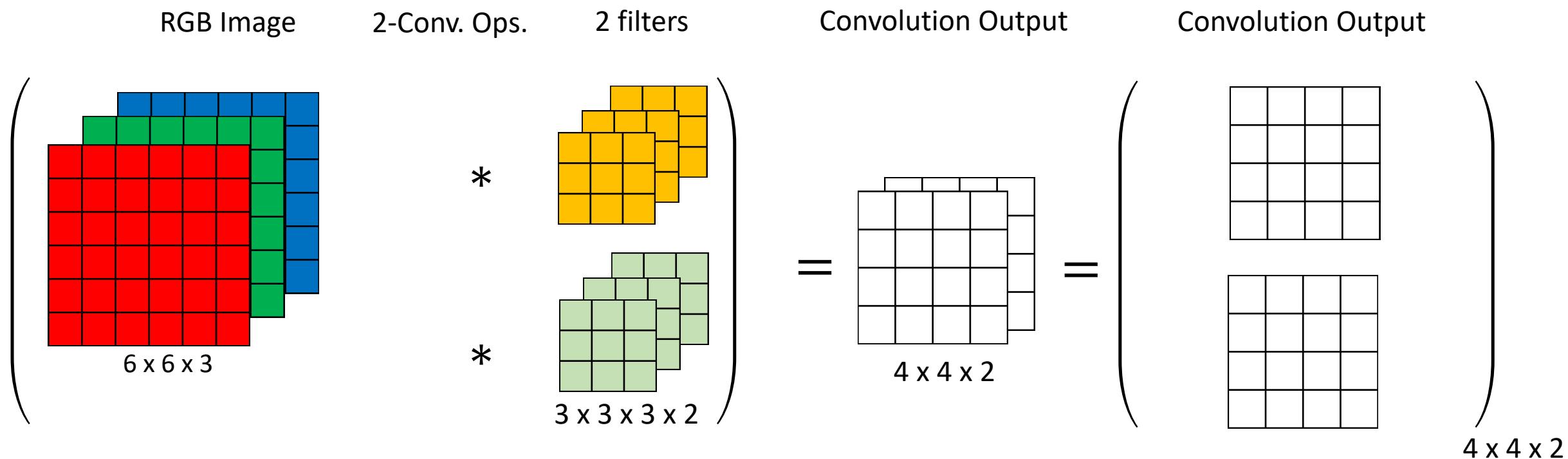
$f \times f \times n_c \times 2$

Number of filters

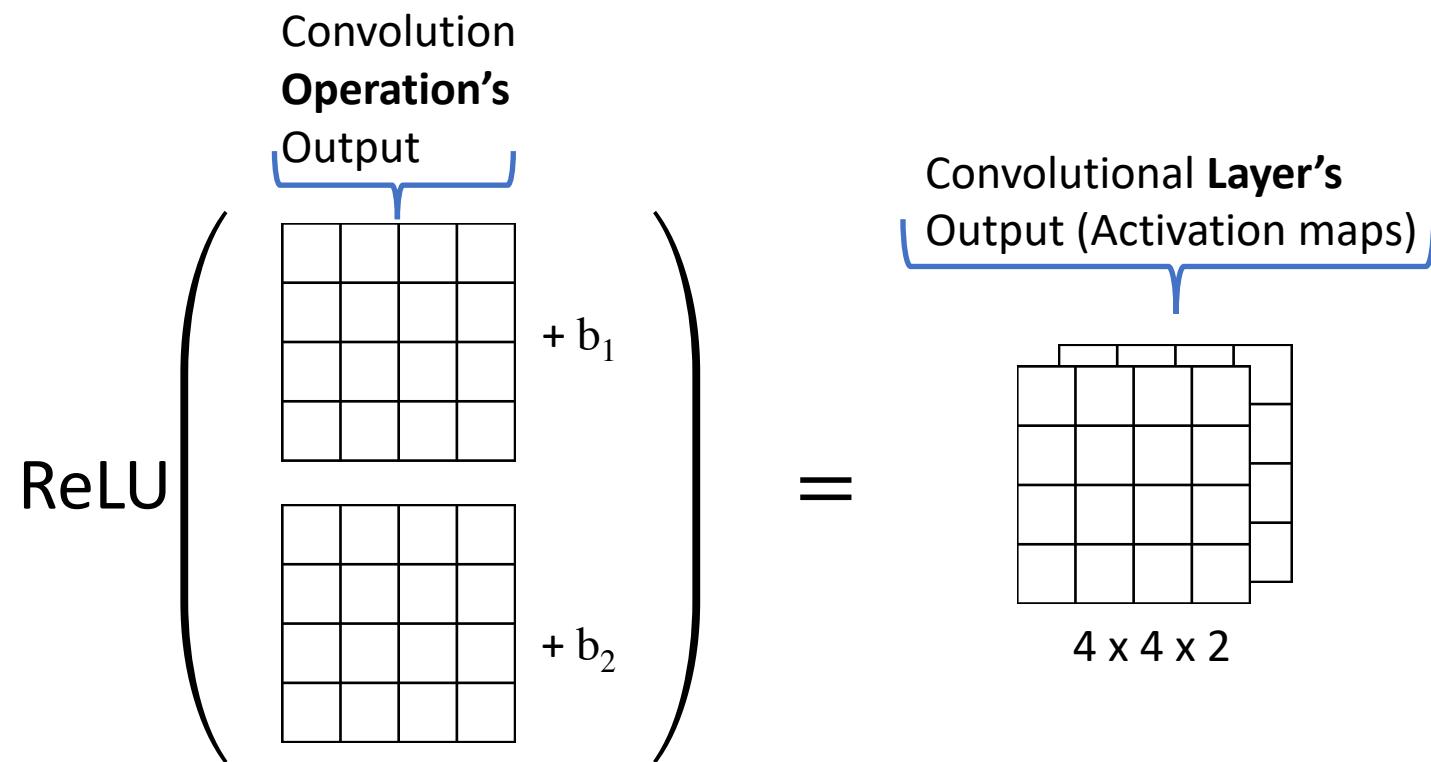
=

$(n - f + 1) \times (n - f + 1) \times 2$

# Multiple filters



# A Convolutional Layer of CNN (ConvNet)



Convolutional Layer Output: Activation\_map =  $\text{ReLU} ((RGBImage * filter) + b)$

Output of a neuron in a FC NN:  
NN

$a = \text{ReLU} (\mathbf{w}^T \mathbf{x} + b)$  in FC

Similar notation!  
(but not equal)

# Input – Output Dimensions for $l^{th}$ Convolutional Layer

Input (One image):  $n_H^{[l-1]} \times n_W^{[l-1]} \times n_C^{[l-1]}$



Output:  $n_H^{[l]} \times n_W^{[l]} \times n_C^{[l]}$

Padding size:  $p^{[l]} \times p^{[l]}$

Stride size:  $s^{[l]} \times s^{[l]}$

Filter size:  $f^{[l]} \times f^{[l]} \times n_C^{[l-1]} \times n_C^{[l]} = \text{One\_FilterDims} \times n_C^{[l]}$

(Number of Weights = Filter size, Bias size:  $n_C^{[l]}$ )

Output dims (for one image) :  $\left\lfloor \frac{n_H^{[l-1]} + 2p^{[l]} - f^{[l]}}{s^{[l]}} + 1 \right\rfloor \times \left\lfloor \frac{n_W^{[l-1]} + 2p^{[l]} - f^{[l]}}{s^{[l]}} + 1 \right\rfloor \times n_C^{[l]}$

$n_H^{[l]}$       x       $n_W^{[l]}$       x       $n_C^{[l]}$

The output dimensions for  $m$  input samples  $A^{[l]}$ :  $m \times n_H^{[l]} \times n_W^{[l]} \times n_C^{[l]}$

# Common Layer Types in a ConvNet

- Convolutional Layer (CONV)
- Fully Connected (FC, Dense)
- Locally Connected (not so commonly used!)

These three layers have parameters to be learned!

- Pooling (POOL)
- Activation Layer
  - (Sometimes, activation function is considered as a separate layer)

No learned parameters for these layers!

...

We have already seen both CONV and FC layers!

For now, we will use only those two layer types mentioned above.

# Pooling layer: Max pooling

Input			
1	3	2	1
2	7	1	0
2	3	4	3
5	6	1	2

Filter size:  $f = 2$   
Stride size:  $s = 2$  

Output	
7	2
6	4

Pooling is typically applied on each channel separately.

# Pooling layer: Max pooling

Input

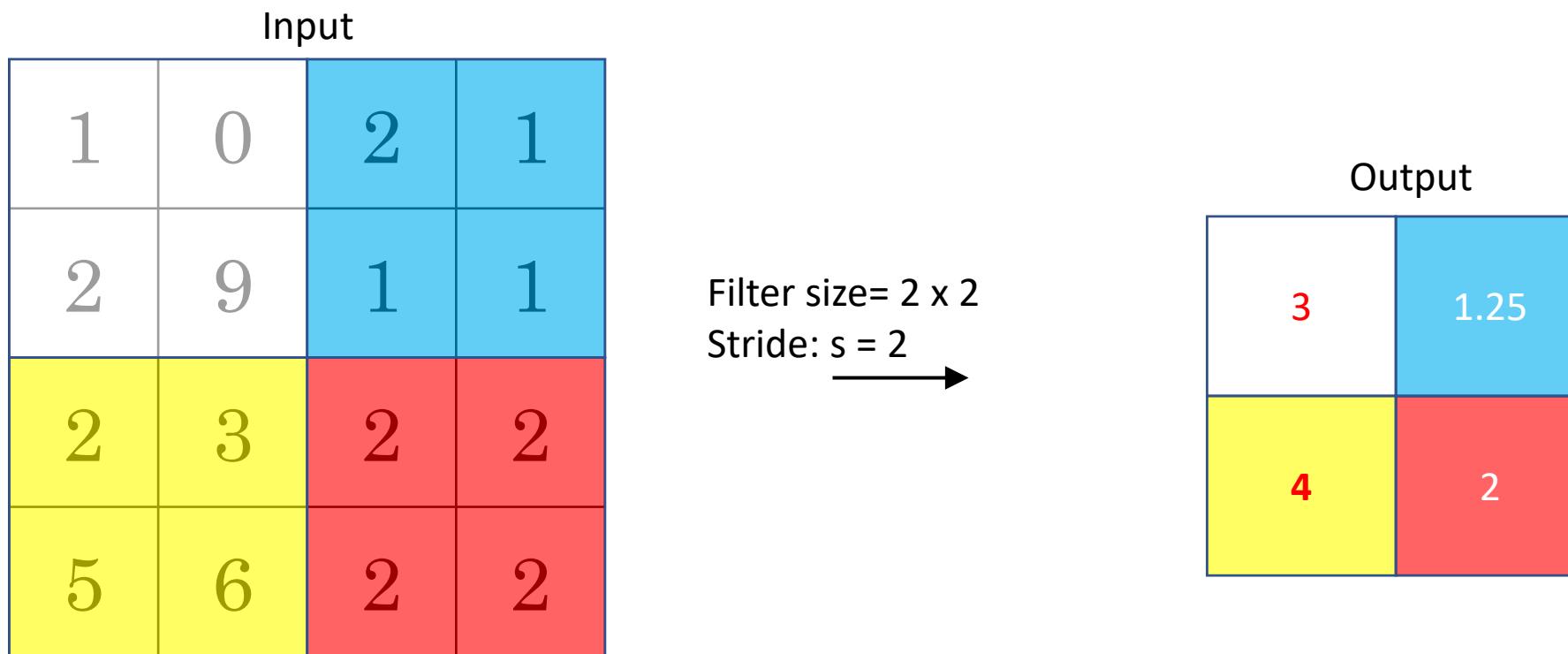
1	2	2	1	4
2	7		3	5
1				2
6	3		1	0
5	6	0	2	7

$f=3$   
 $s=1$



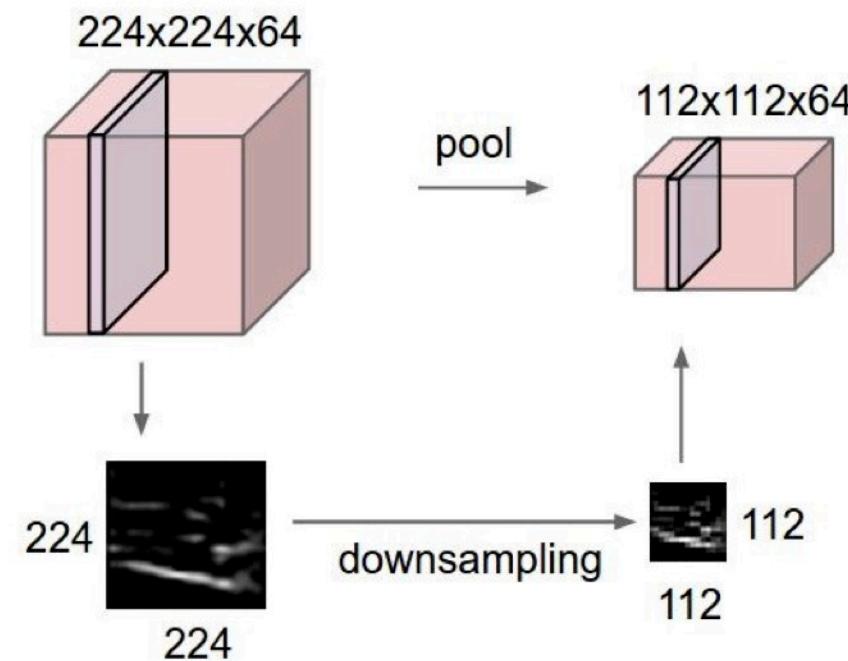
Output


# Pooling layer: Average pooling



# Pooling layer

- makes the representations smaller and more manageable
- operates over each activation map independently:



# Summary of pooling

Hyperparameters:

f : filter size

f=2, s=2

s : stride

f=3, s=2

p = 0

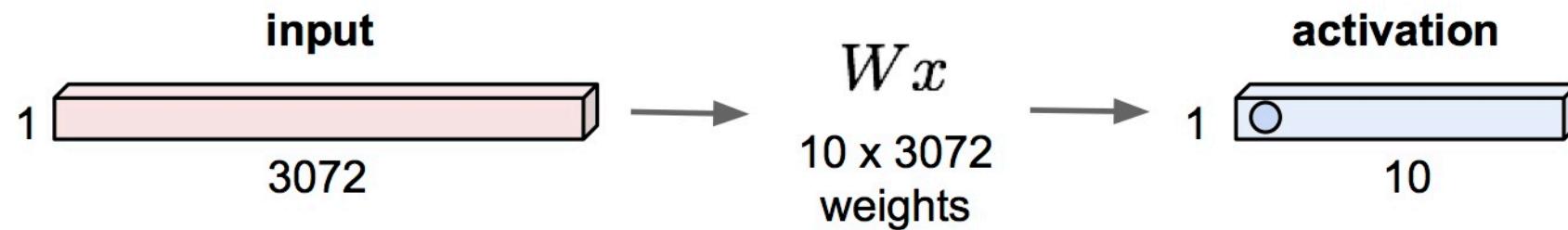
**Max or Average** pooling

Input:  $n_H \times n_W \times n_C$

Output:  $\left\lfloor \frac{n_H - f}{s} + 1 \right\rfloor \times \left\lfloor \frac{n_W - f}{s} + 1 \right\rfloor \times n_C$

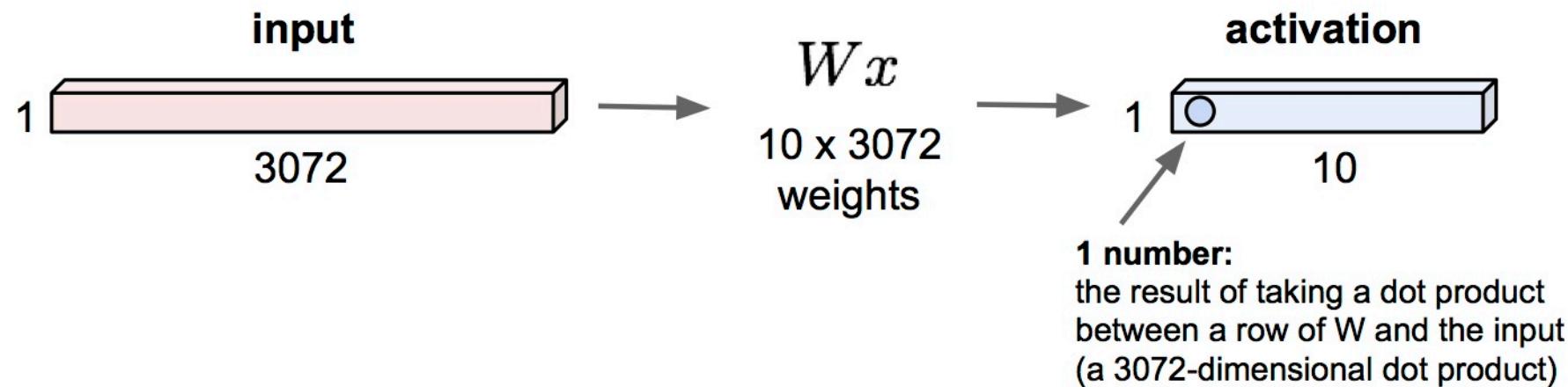
# Fully Connected Layer

32x32x3 image -> stretch to 3072 x 1



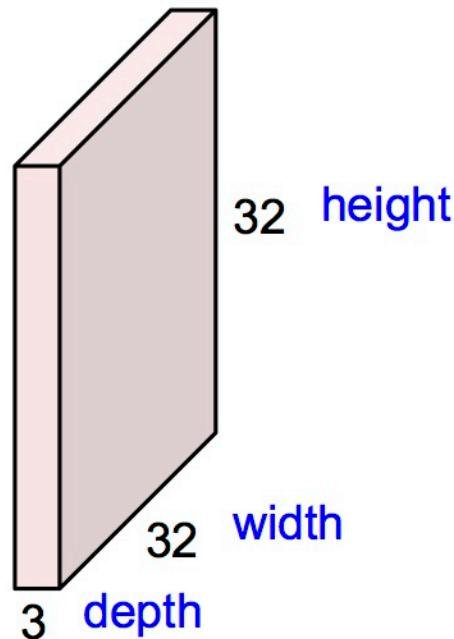
# Fully Connected Layer

32x32x3 image -> stretch to  $3072 \times 1$



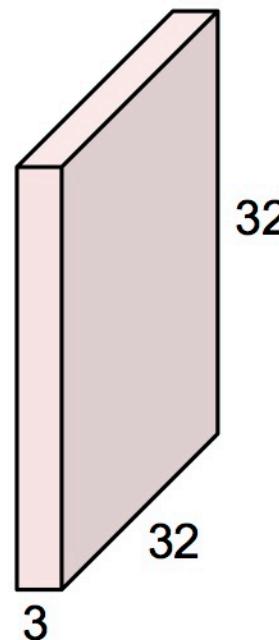
# Convolution Layer

32x32x3 image -> preserve spatial structure

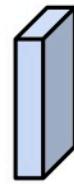


# Convolution Layer

32x32x3 image



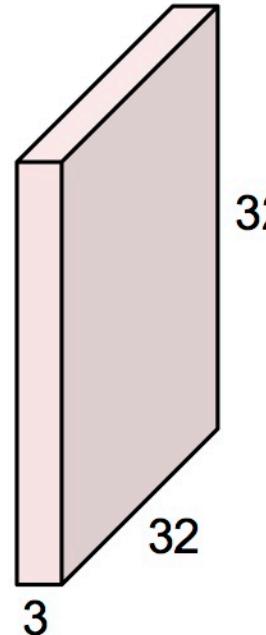
5x5x3 filter



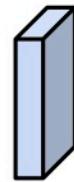
**Convolve** the filter with the image  
i.e. “slide over the image spatially,  
computing dot products”

# Convolution Layer

32x32x3 image



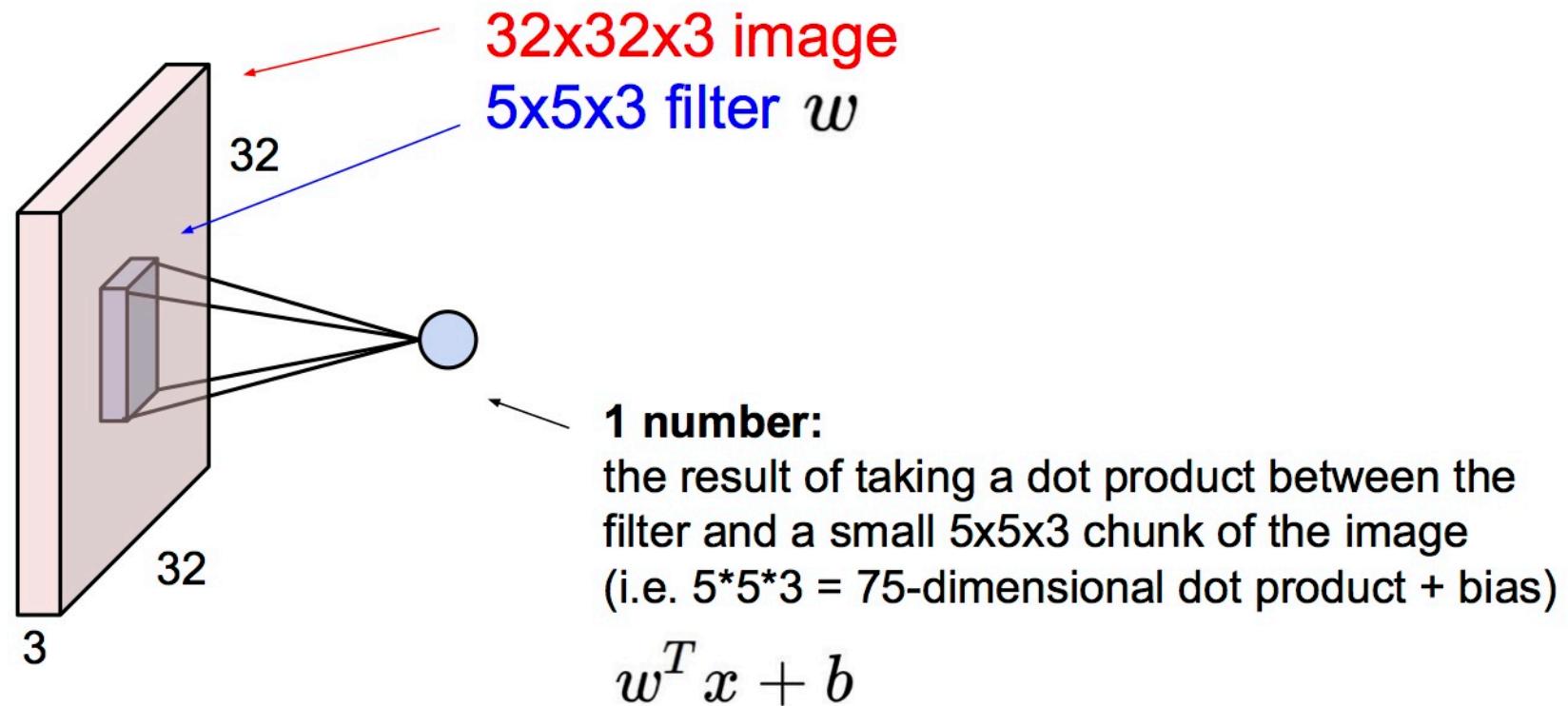
5x5x3 filter



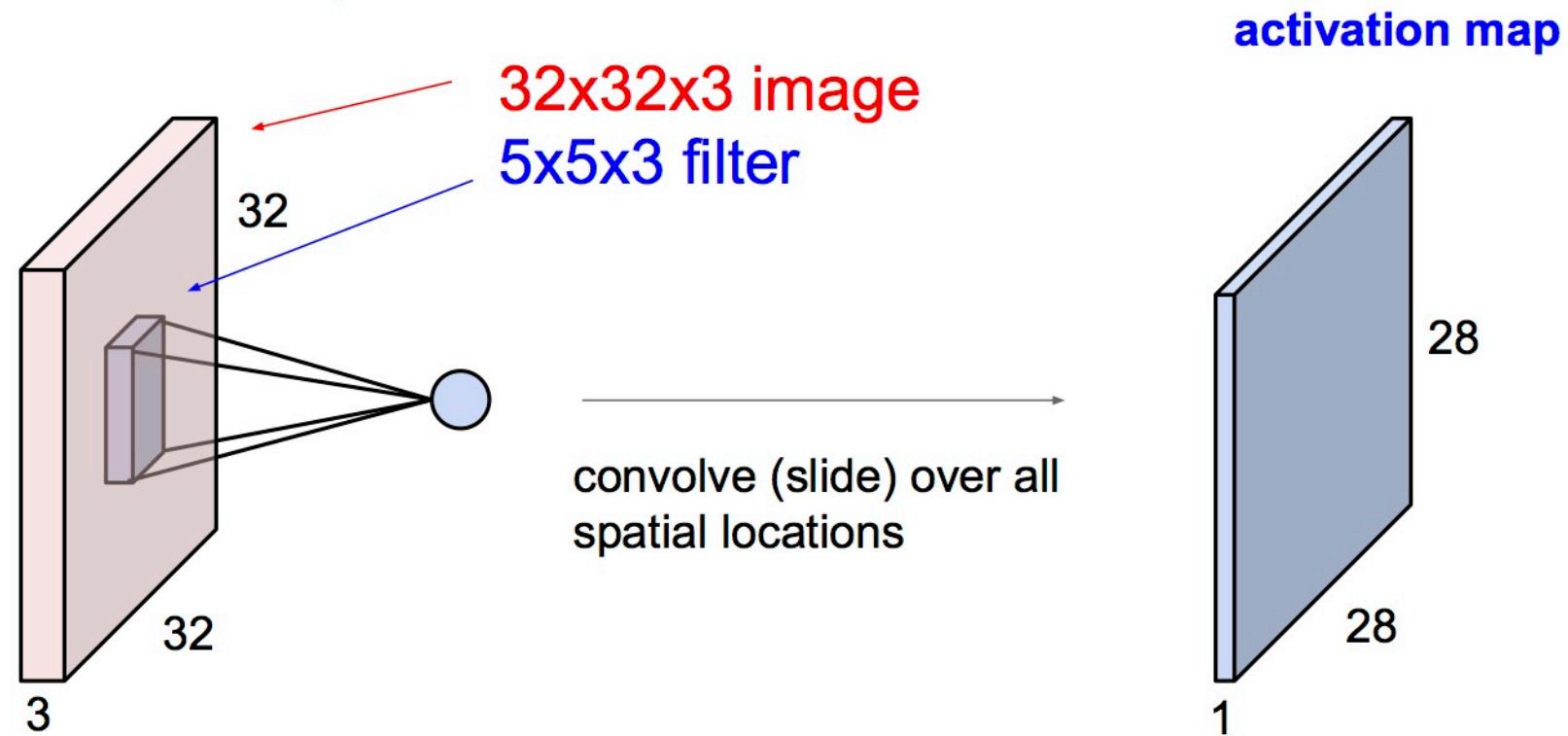
Filters always extend the full depth of the input volume

**Convolve** the filter with the image  
i.e. “slide over the image spatially,  
computing dot products”

# Convolution Layer

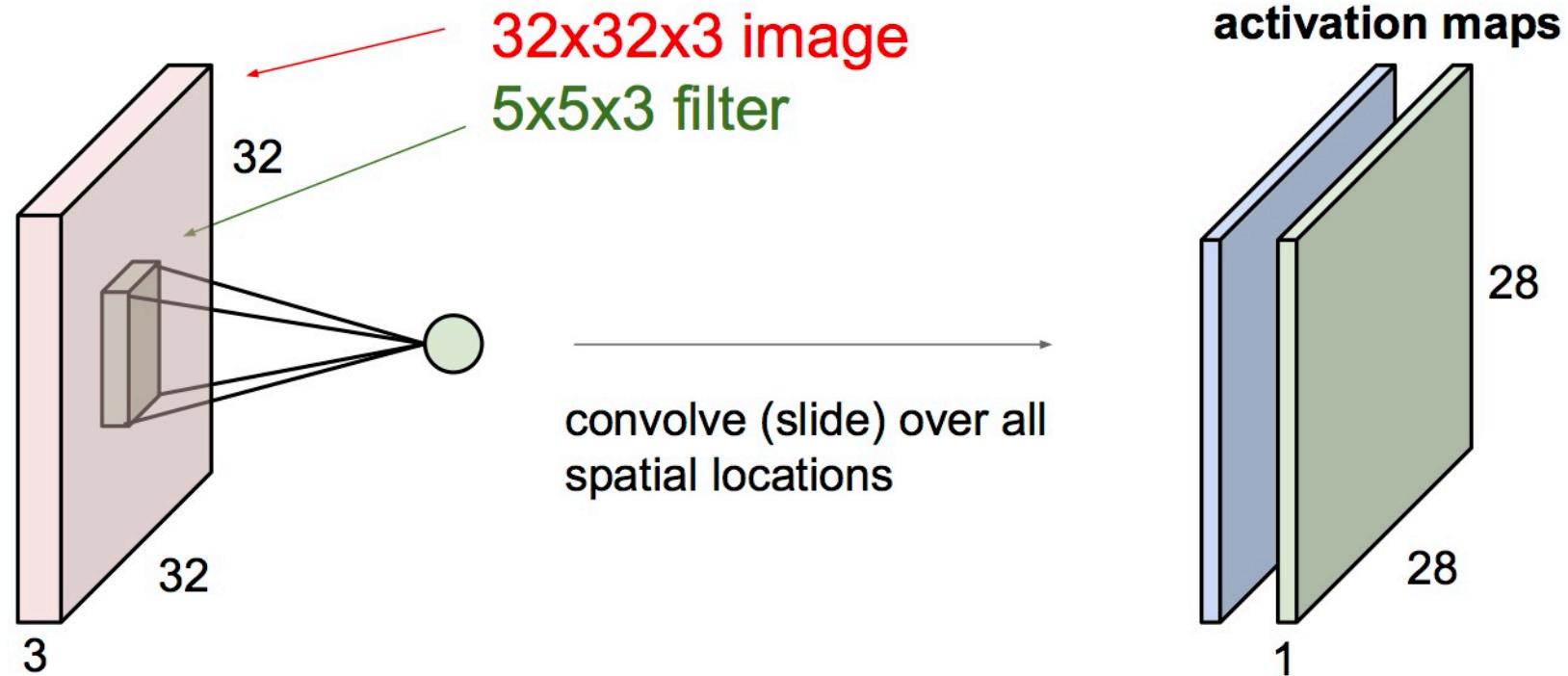


# Convolution Layer

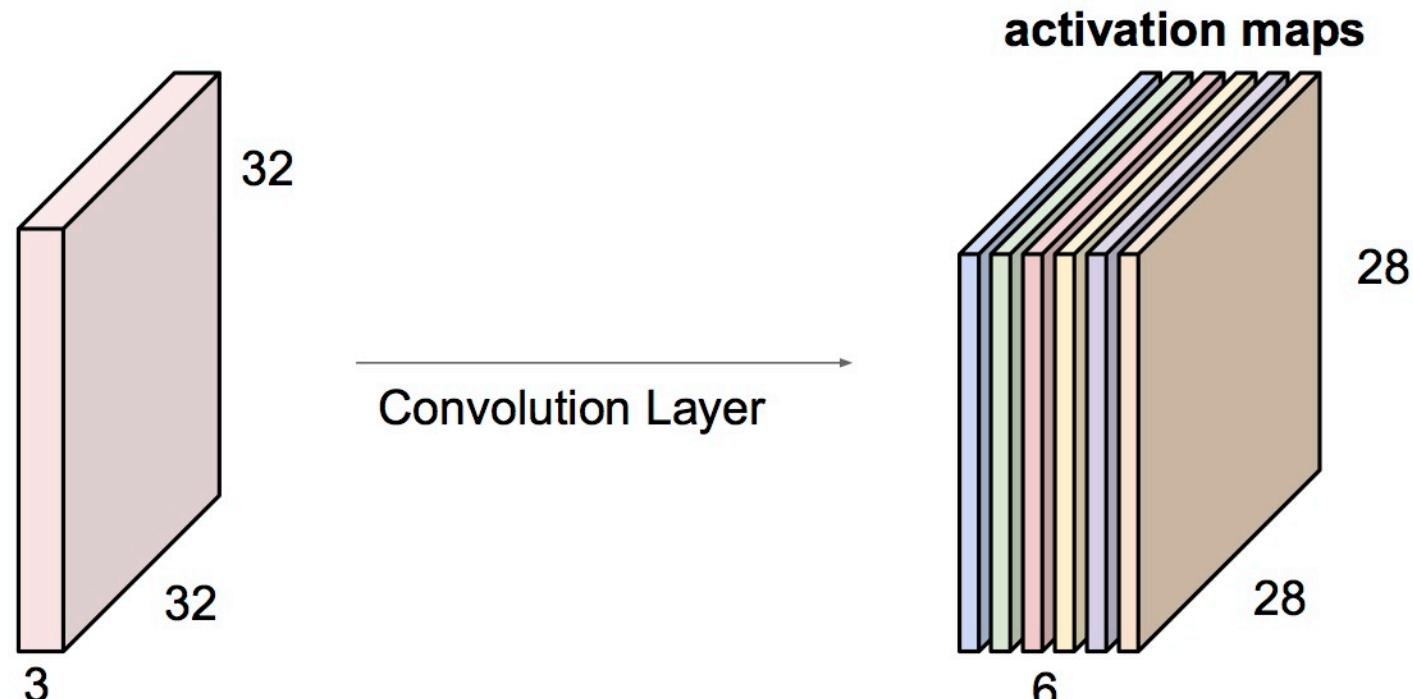


# Convolution Layer

consider a second, **green** filter

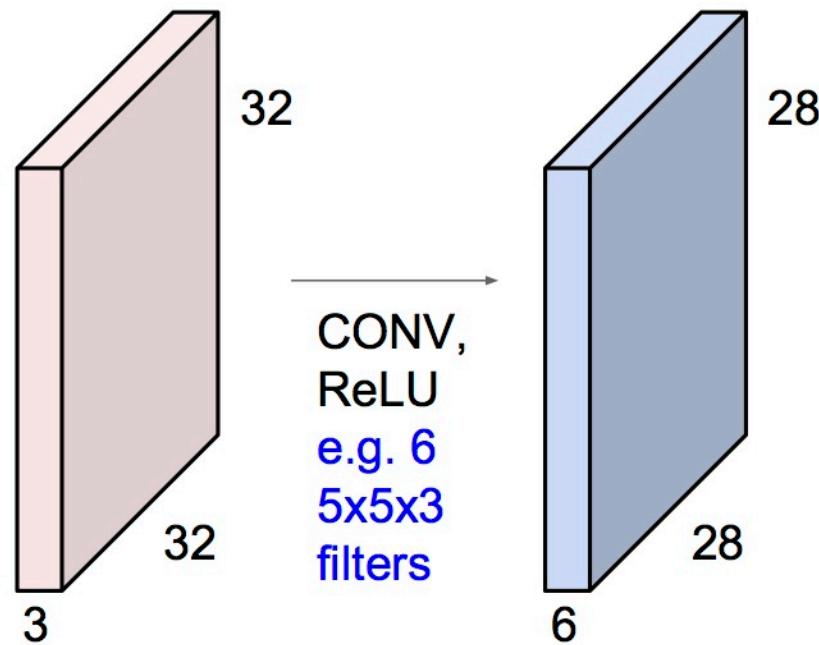


For example, if we had 6 5x5 filters, we'll get 6 separate activation maps:

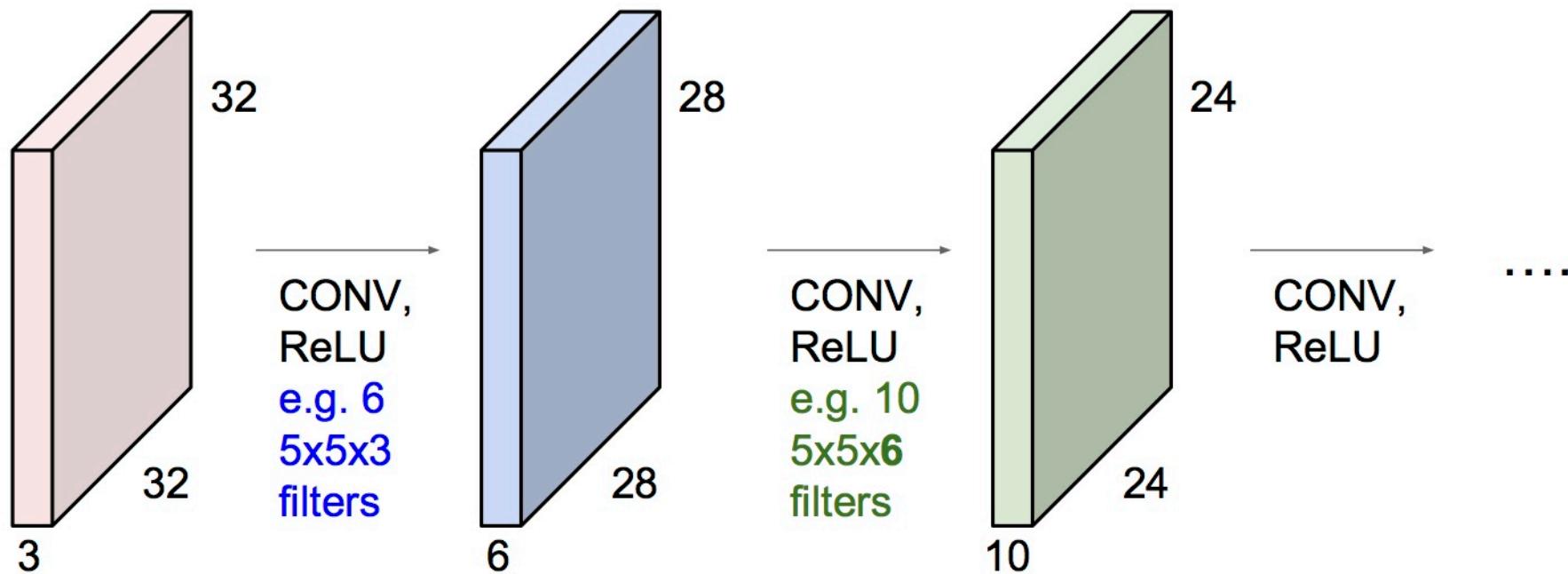


We stack these up to get a “new image” of size 28x28x6!

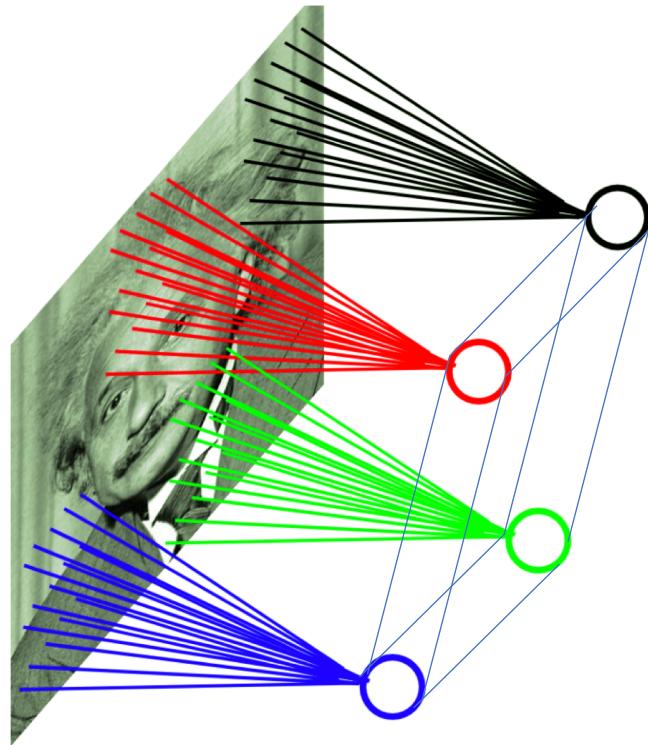
**Preview:** ConvNet is a sequence of Convolution Layers, interspersed with activation functions



**Preview:** ConvNet is a sequence of Convolutional Layers, interspersed with activation functions



# Parameter Sharing



One filter to rule them all!



- Image Classification & Recognition: A main task in computer vision
  - PASCAL
  - IMAGENET

# Pop Quiz

- How many parameters do we have in a conv layer with **100** filters that are  $3 \times 3 \times 3$  dimensional each?
- For one filter: we have  $3 \times 3 \times 3 = 27$  filter weights; and + 1 bias = **28** parameters.
- For all **100** filters: we have  **$28 \times 100 = 2800$  parameters** to learn for that layer!

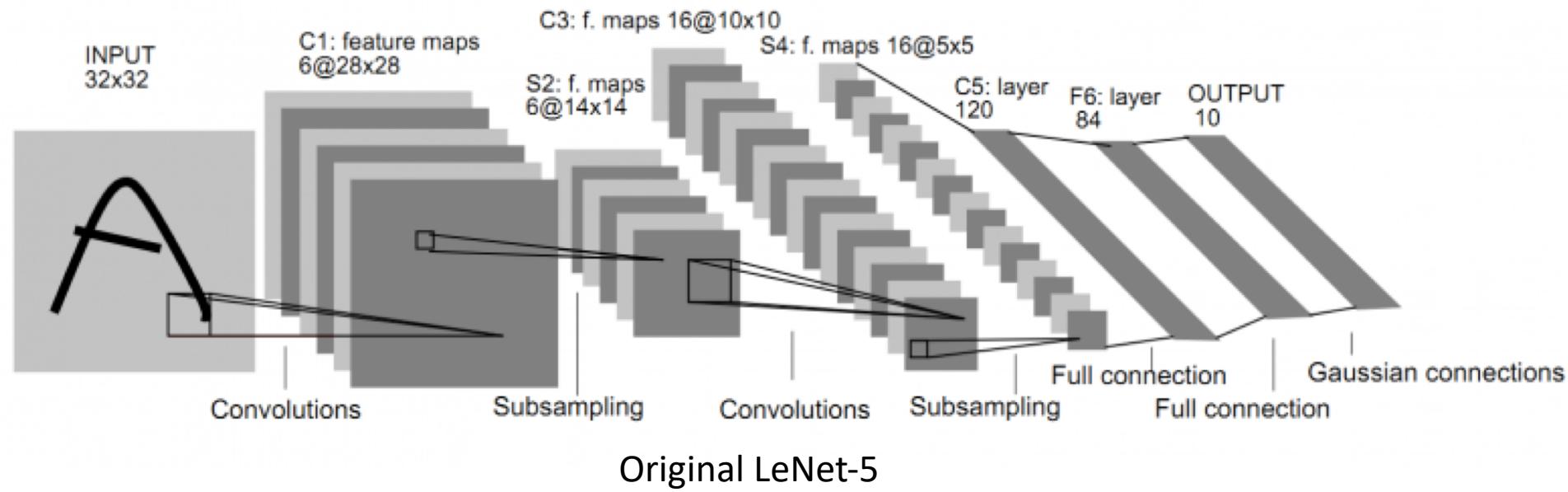
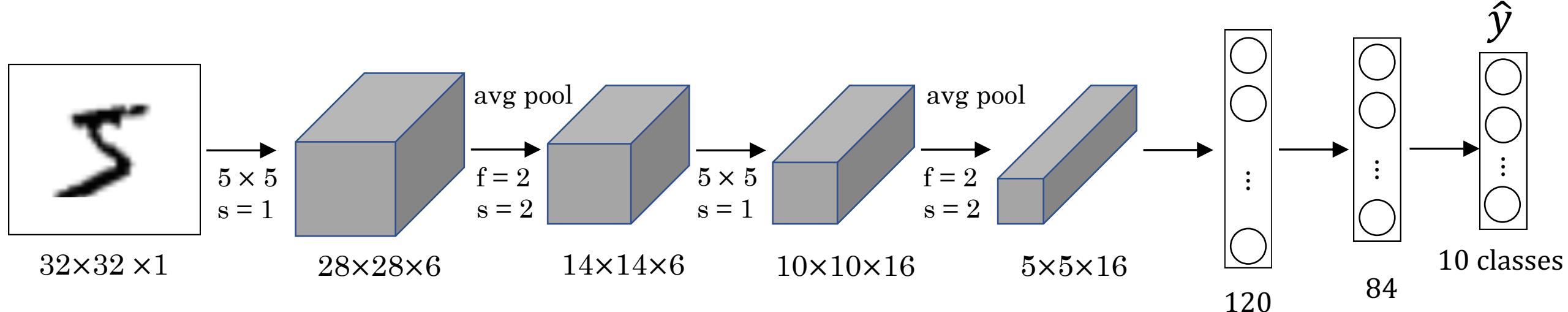
# Pop Quiz

- Imagine that you are designing a CNN with 5 layers. The first 4 layers are convolutional and the last layer is FC layer with one neuron (using logistic reg).
  - The input dims are: 600x600x3.
  - In each of the 4 conv. layers:
    - We have 10 filters (for each,  $f = 3$ )
    - Stride is 2
    - Padding is “valid”.

What is the output dims at the end of the 5<sup>th</sup> layer?

**30 seconds to return your answers! ☺**

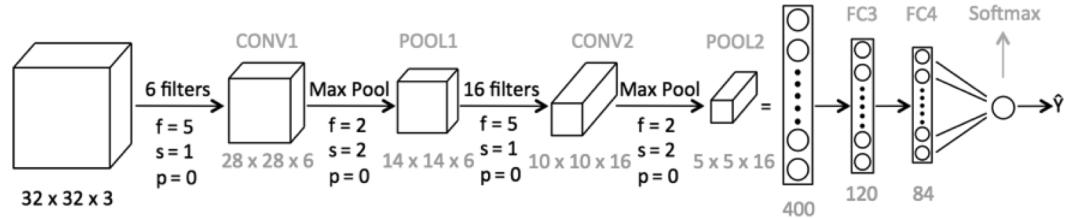
# (Slightly Modified) LeNet-5



# (Slightly Modified) LeNet-5

For: Grayscale image (single channel)

	Activation shape	Activation Size	# parameters
Input:	(32,32,1)	1,024	0
CONV1 (f=5, s=1)	(28,28,6)	4,704	150 + 6 bias = <b>156</b>
POOL1	(14,14,6)	1,176	0
CONV2 (f=5, s=1)	(10,10,16)	1,600	(5x5x6) x 16 = 2400 2400 + 16 = <b>2416</b>
POOL2	(5,5,16)	400	0
FC3	(120,1)	120	48,000 + 120 = <b>48,120</b>
FC4	(84,1)	84	10,080 + 84 = <b>10,164</b>
Softmax	(10,1)	10	84x10 + 10 = <b>850</b>



For: Colored (RGB) image

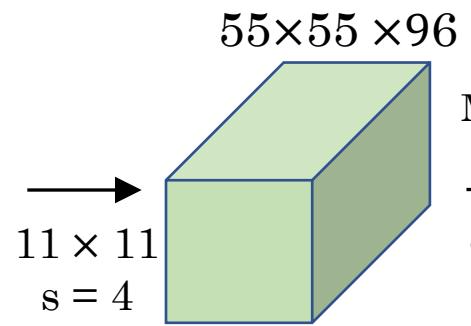
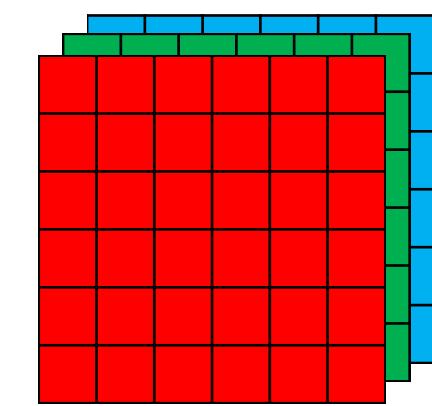
	Activation shape	Activation Size	# parameters
Input:	(32,32,3)	3,072	0
CONV1 (f=5, s=1)	(28,28,6)	4,704	450 + 6 bias = <b>456</b>
POOL1	(14,14,6)	1,176	0
CONV2 (f=5, s=1)	(10,10,16)	1,600	(5x5x6) x 16 = 2400 2400 + 16 = <b>2416</b>
POOL2	(5,5,16)	400	0
FC3	(120,1)	120	48,000 + 120 = <b>48,120</b>
FC4	(84,1)	84	10,080 + 84 = <b>10,164</b>
Softmax	(10,1)	10	84x10 + 10 = <b>850</b>

Tensorflow output for the (grayscale) model summarizing the model parameters:

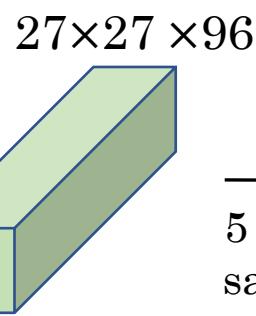
```
Variables: name (type shape) [size]
-----
Variable:0 (float32_ref 5x5x1x6) [150, bytes: 600]
conv1_biases:0 (float32_ref 6) [6, bytes: 24]
Variable_1:0 (float32_ref 5x5x6x16) [2400, bytes: 9600]
conv2_biases:0 (float32_ref 16) [16, bytes: 64]
Variable_2:0 (float32_ref 400x120) [48000, bytes: 192000]
fc1_biases:0 (float32_ref 120) [120, bytes: 480]
Variable_3:0 (float32_ref 120x84) [10080, bytes: 40320]
fc2_biases:0 (float32_ref 84) [84, bytes: 336]
Variable_4:0 (float32_ref 84x10) [840, bytes: 3360]
fc3_biases:0 (float32_ref 10) [10, bytes: 40]
Total size of variables: 61706
Total bytes of variables: 246824
sedat@sedat-ThinkPad-P50:~/
```

# (Slightly modified) AlexNet

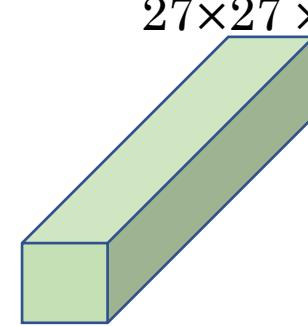
$227 \times 227 \times 3$



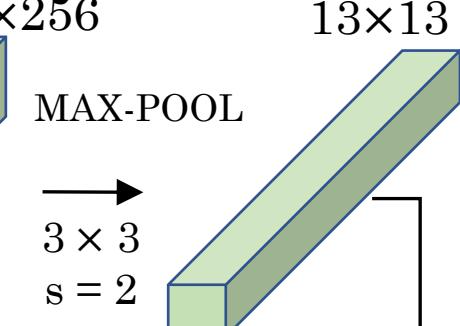
MAX-POOL  
 $3 \times 3$   
 $s = 2$



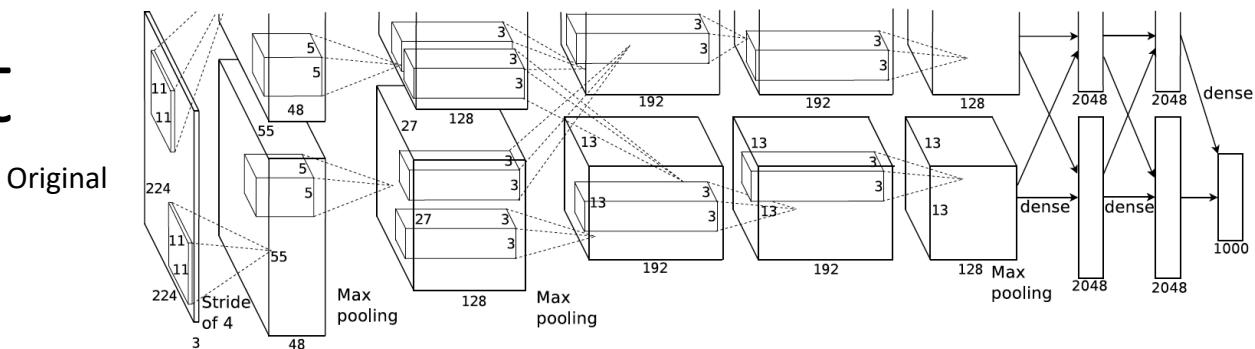
MAX-POOL  
 $3 \times 3$   
 $s = 2$



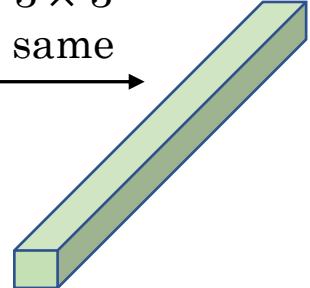
MAX-POOL  
 $3 \times 3$   
 $s = 2$



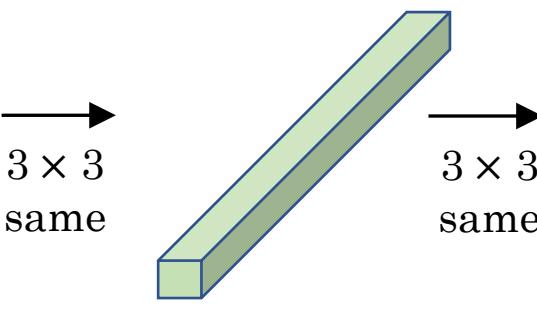
MAX-POOL  
 $3 \times 3$   
 $s = 2$



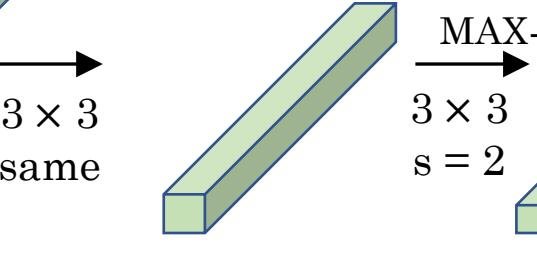
$3 \times 3$   
same



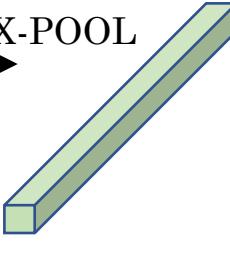
$3 \times 3$   
same



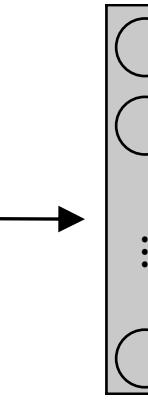
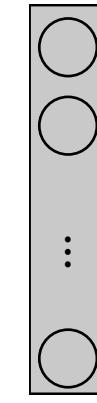
$3 \times 3$   
same



MAX-POOL  
 $3 \times 3$   
 $s = 2$



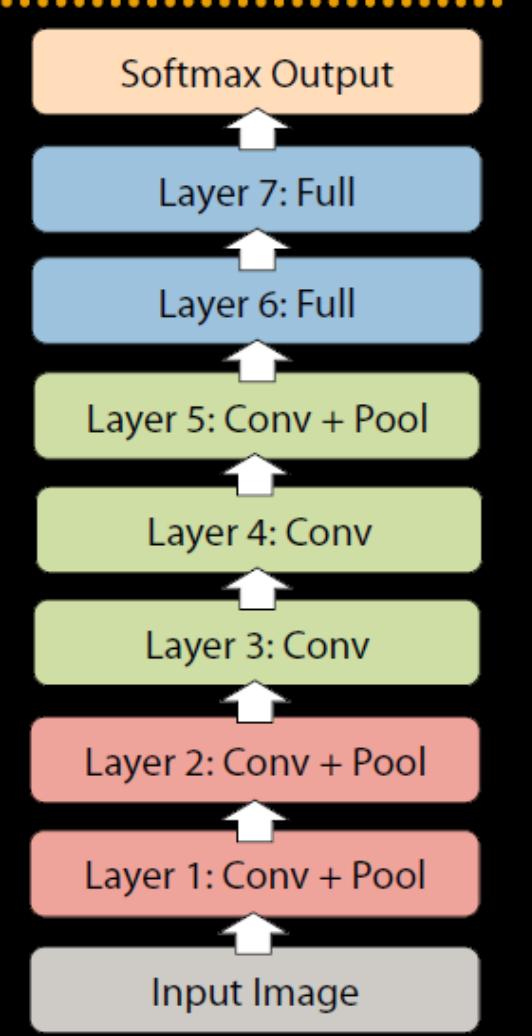
=



Softmax  
1000

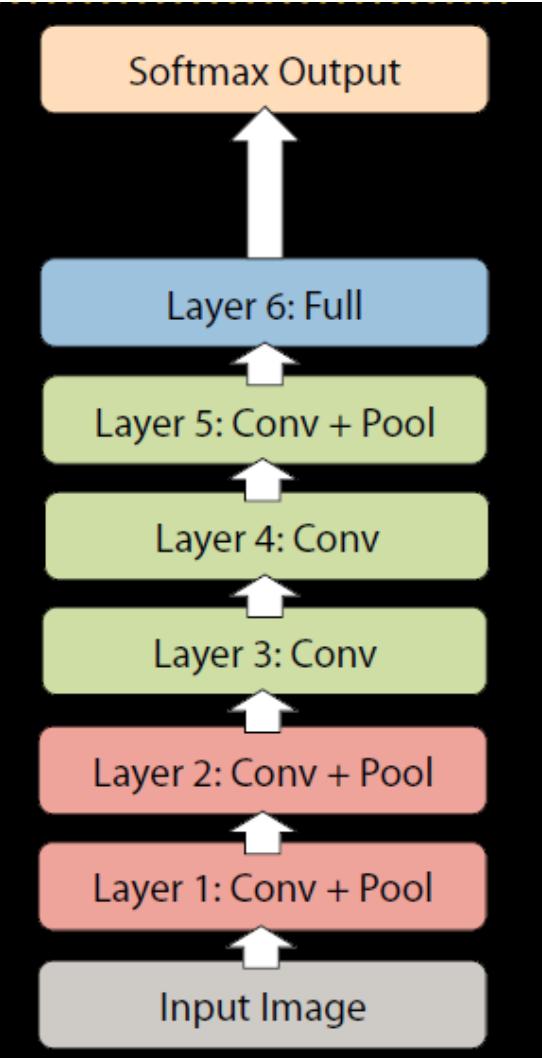
# Why ConvNet should be Deep?

- 8 layers total
- Trained on Imagenet dataset [Deng et al. CVPR'09]
- 18.2% top-5 error
- Our reimplementation:  
18.1% top-5 error



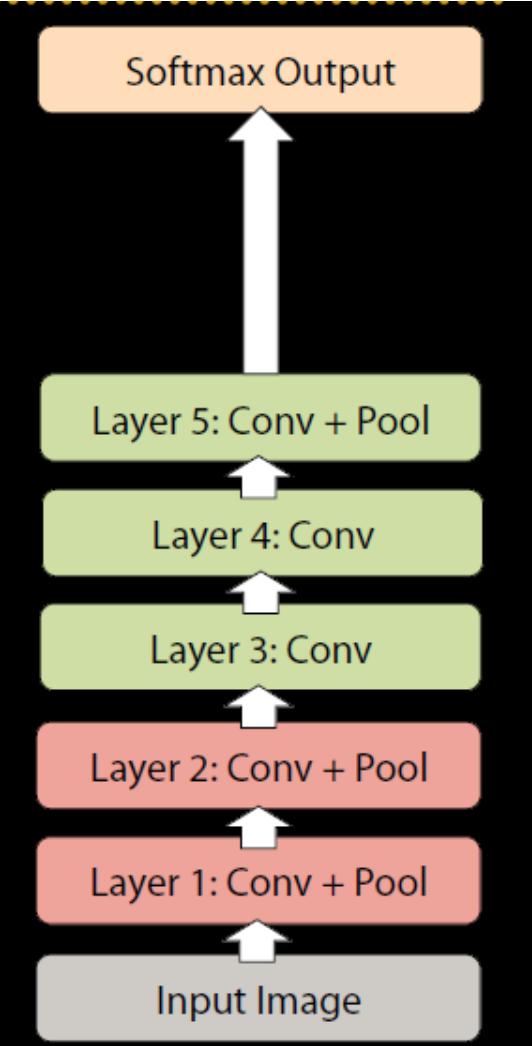
# Why ConvNet should be Deep?

- Remove top fully connected layer
  - Layer 7
- Drop 16 million parameters
- Only 1.1% drop in performance!



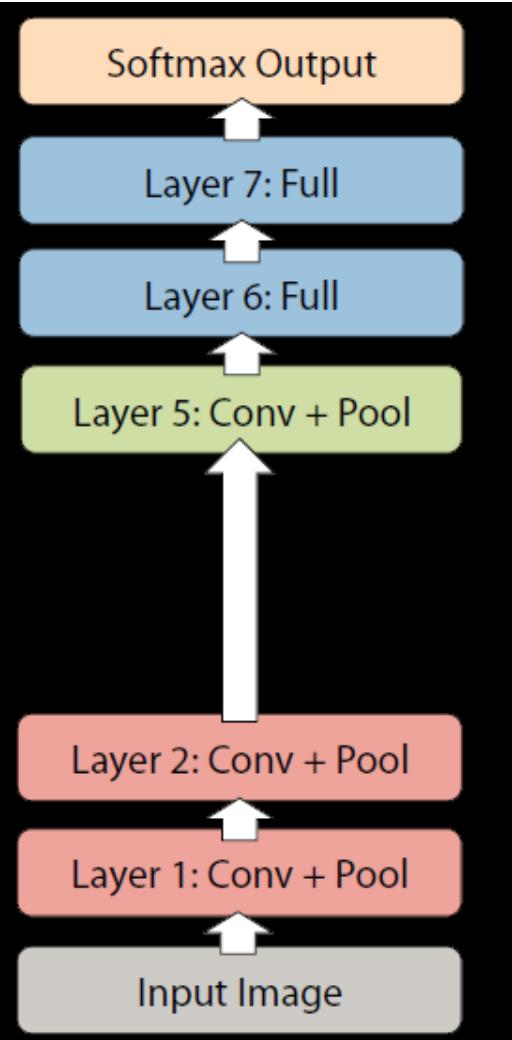
# Why ConvNet should be Deep?

- Remove both fully connected layers
  - Layer 6 & 7
- Drop ~50 million parameters
- 5.7% drop in performance



# Why ConvNet should be Deep?

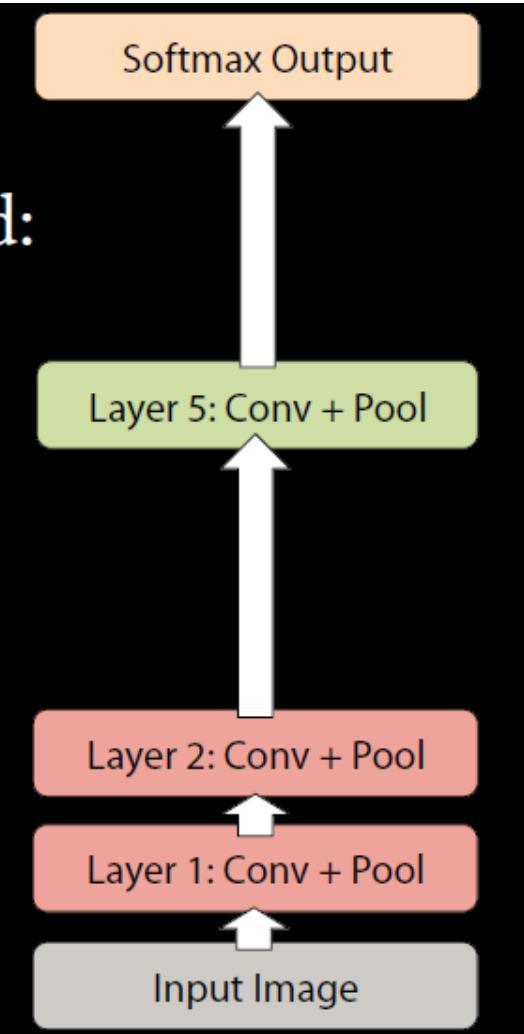
- Now try removing upper feature extractor layers:
  - Layers 3 & 4
- Drop ~1 million parameters
- 3.0% drop in performance



# Why ConvNet should be Deep?

- Now try removing upper feature extractor layers & fully connected:
  - Layers 3, 4, 6 ,7
- Now only 4 layers
- 33.5% drop in performance

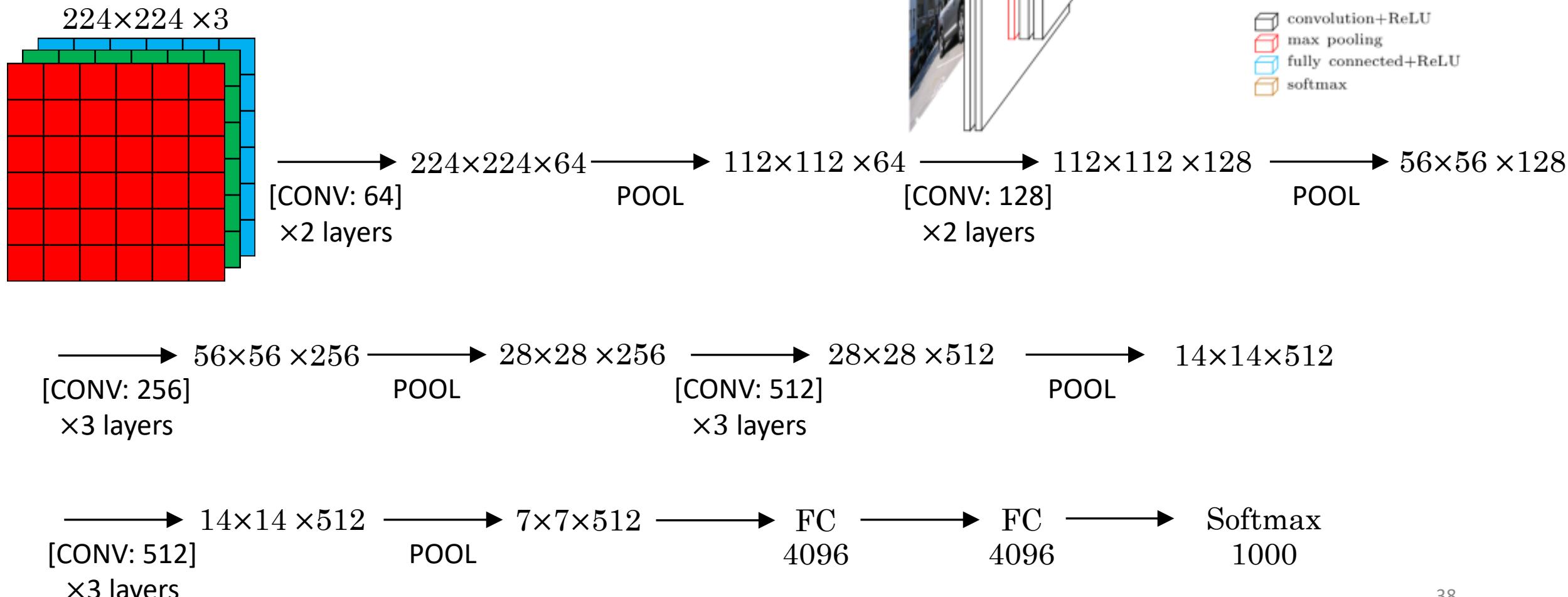
→ Depth of network is key



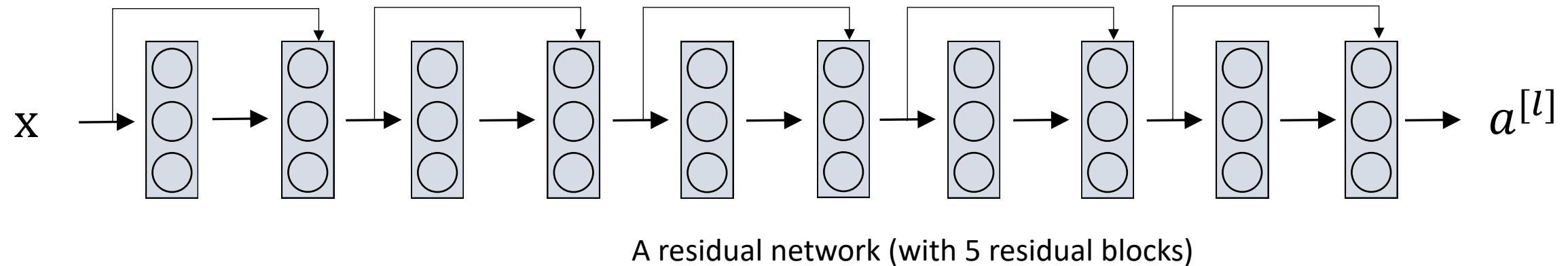
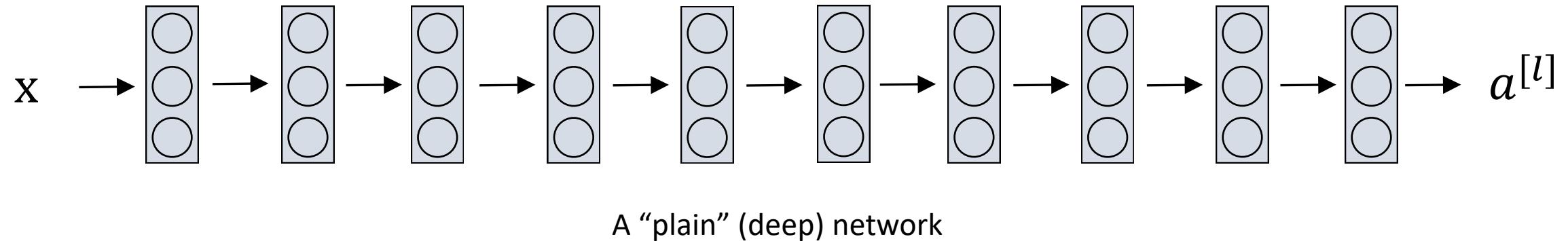
# VGG – 16 (16 layers)

CONV =  $3 \times 3$  filter, s = 1, “same” padding

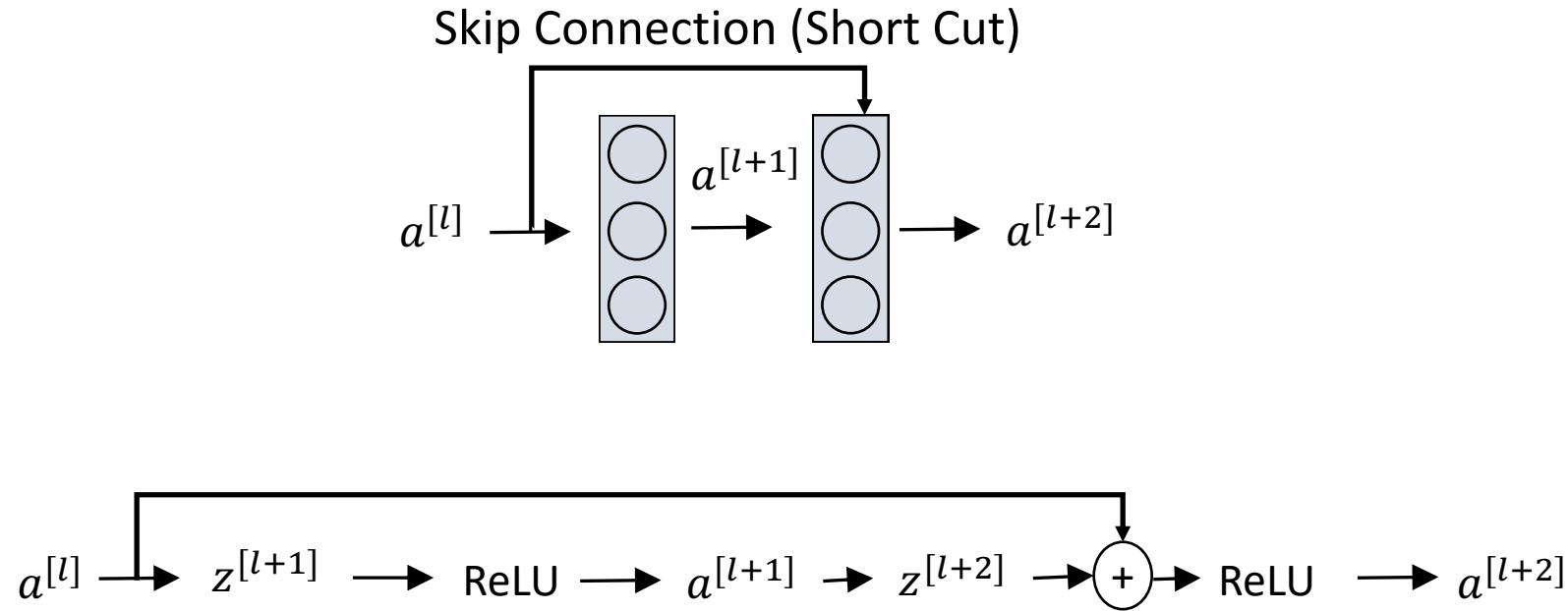
MAX-POOL =  $2 \times 2$ , s = 2



# Residual Network - ResNet (Microsoft)



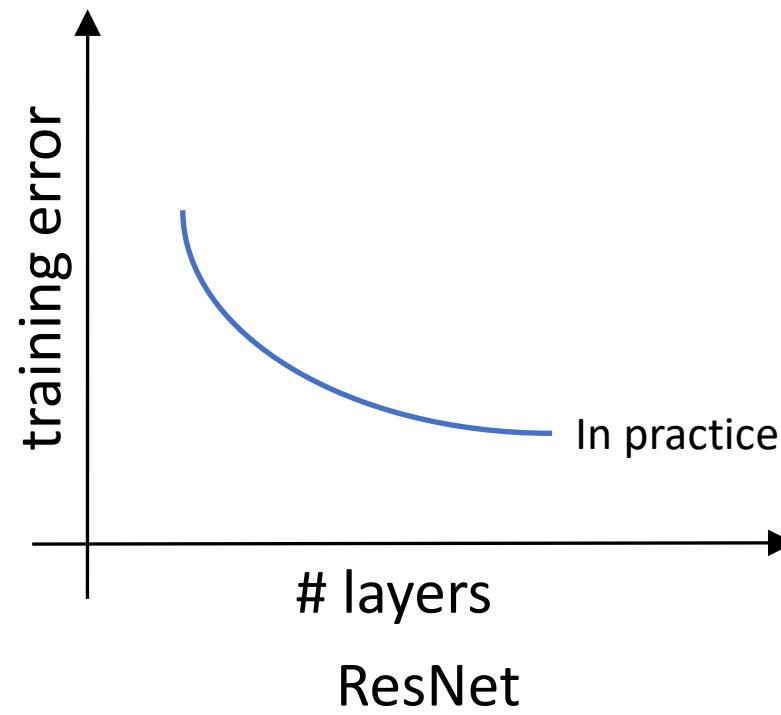
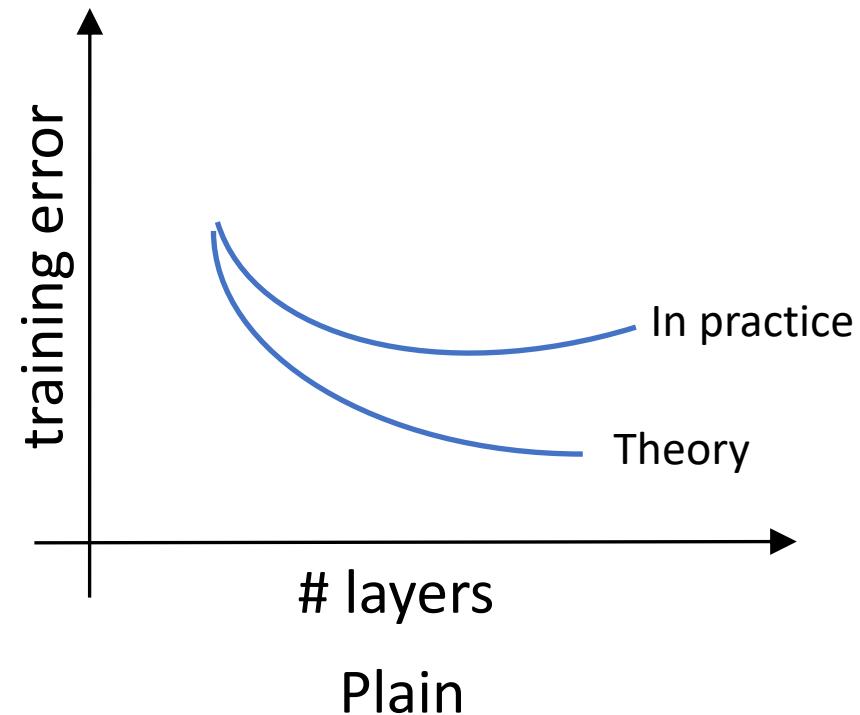
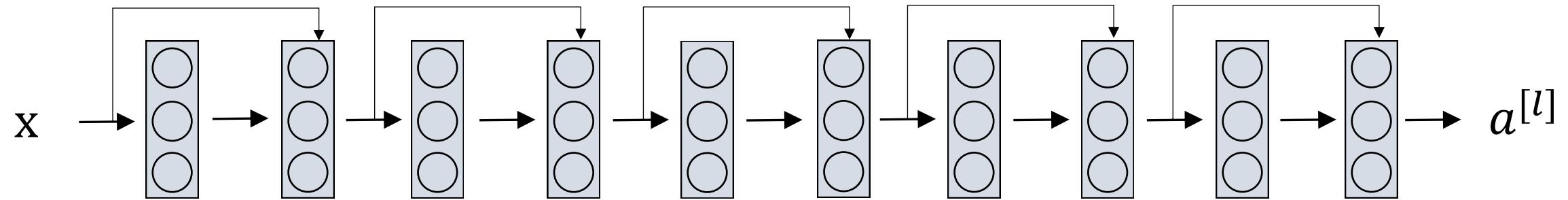
# Residual Building Block



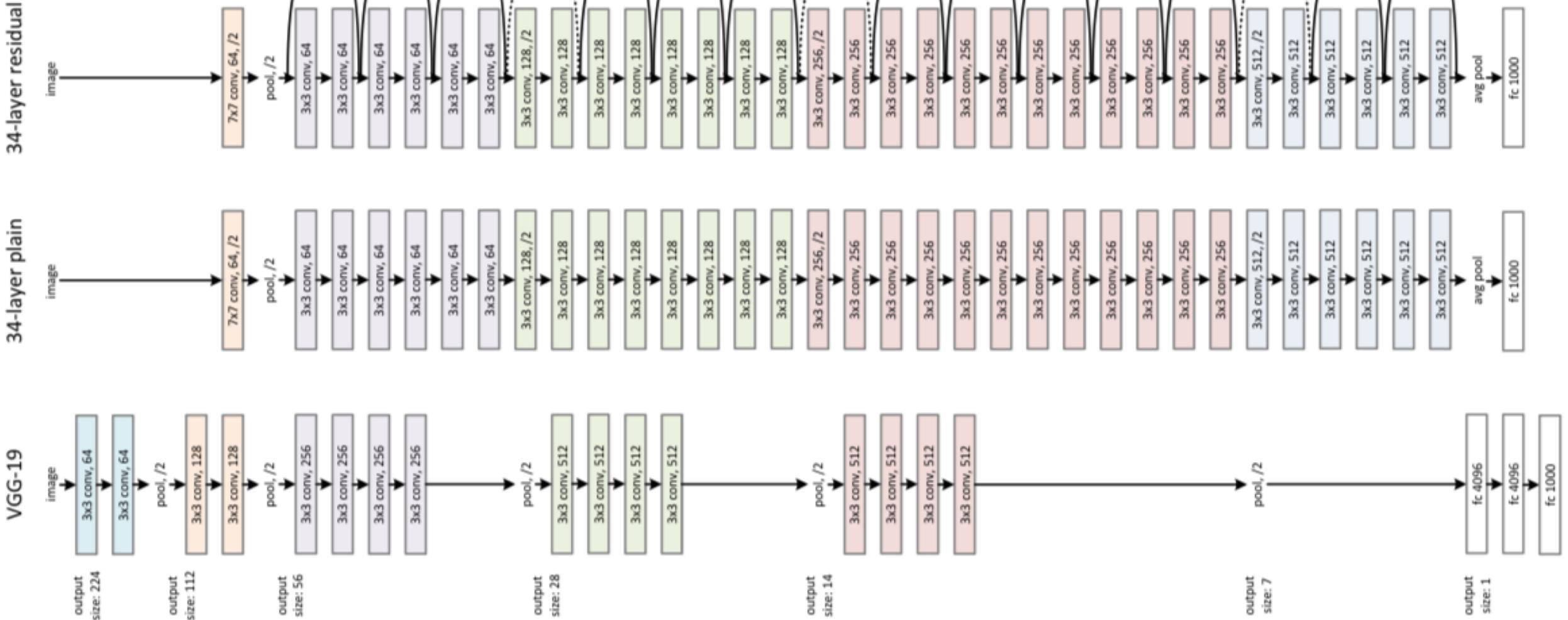
$$z^{[l+1]} = W^{[l+1]} a^{[l]} + b^{[l+1]} \quad a^{[l+1]} = g(z^{[l+1]}) \quad z^{[l+2]} = W^{[l+2]} a^{[l+1]} + b^{[l+2]} \quad a^{[l+2]} = g(z^{[l+2]})$$

With skip connection:  $a^{[l+2]} = g(z^{[l+2]} + a^{[l]})$

# Residual Network



## ResNet



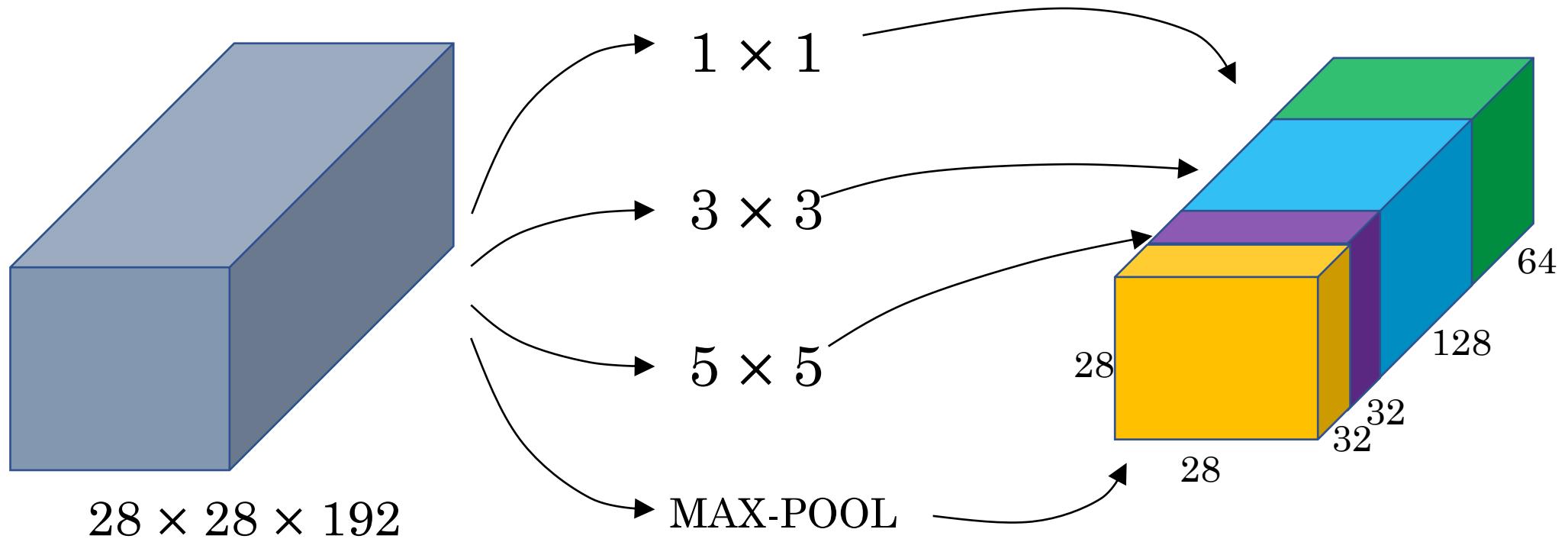
[Source: He et al., 2015. Deep residual networks for image recognition]

# Inception Network



(GoogleNet)

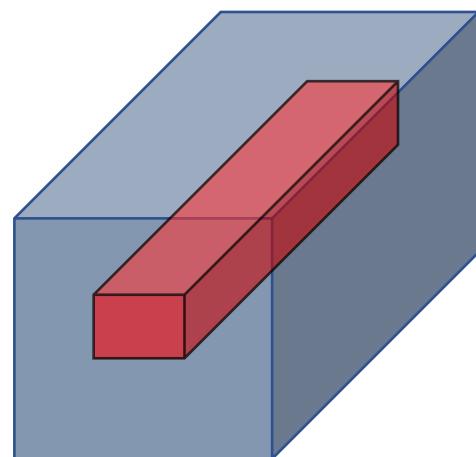
# Main motivation for inception network



# What does a $1 \times 1$ convolution do?

3	1	3	6	5	8
3	2	4	1	3	4
2	1	3	4	9	2
4	5	8	1	7	9
1	5	3	3	4	3
5	4	4	8	3	1

$6 \times 6$



$6 \times 6 \times 32$

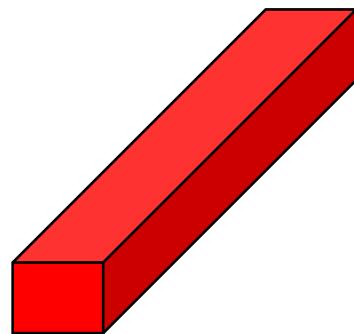
\*

2

=

6	2	6	.	.	.
6	4	8	.	.	.
.	.	.	.	.	.
.	.	.	.	.	.
.	.	.	.	.	.
.	.	.	.	.	.

\*

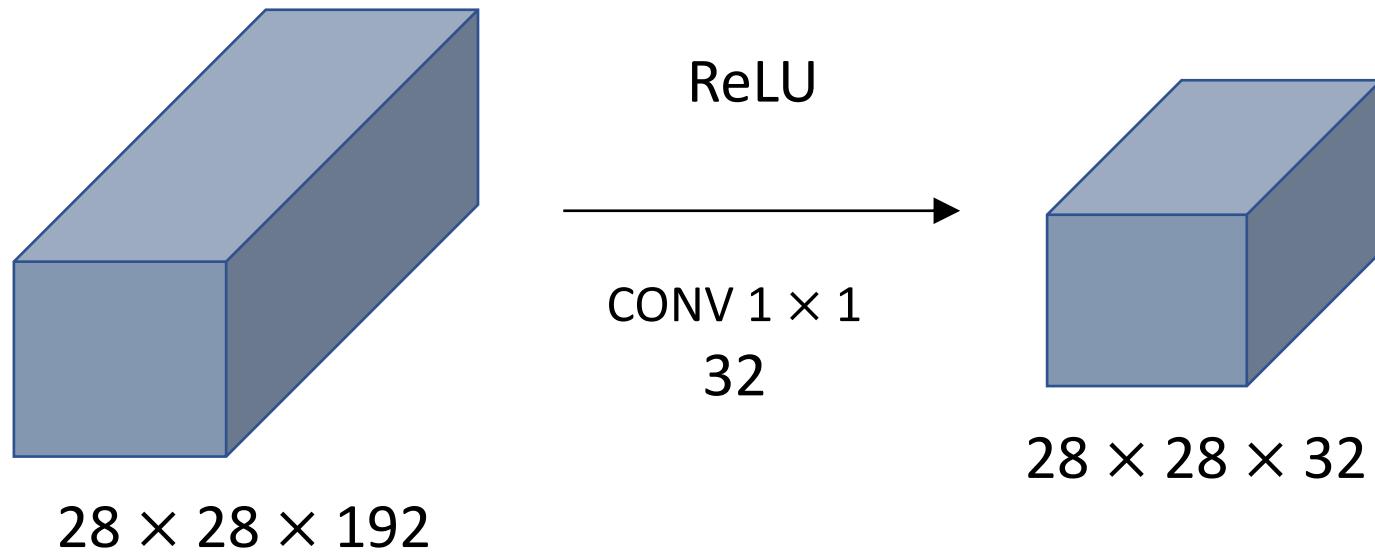


=


$1 \times 1 \times 32$

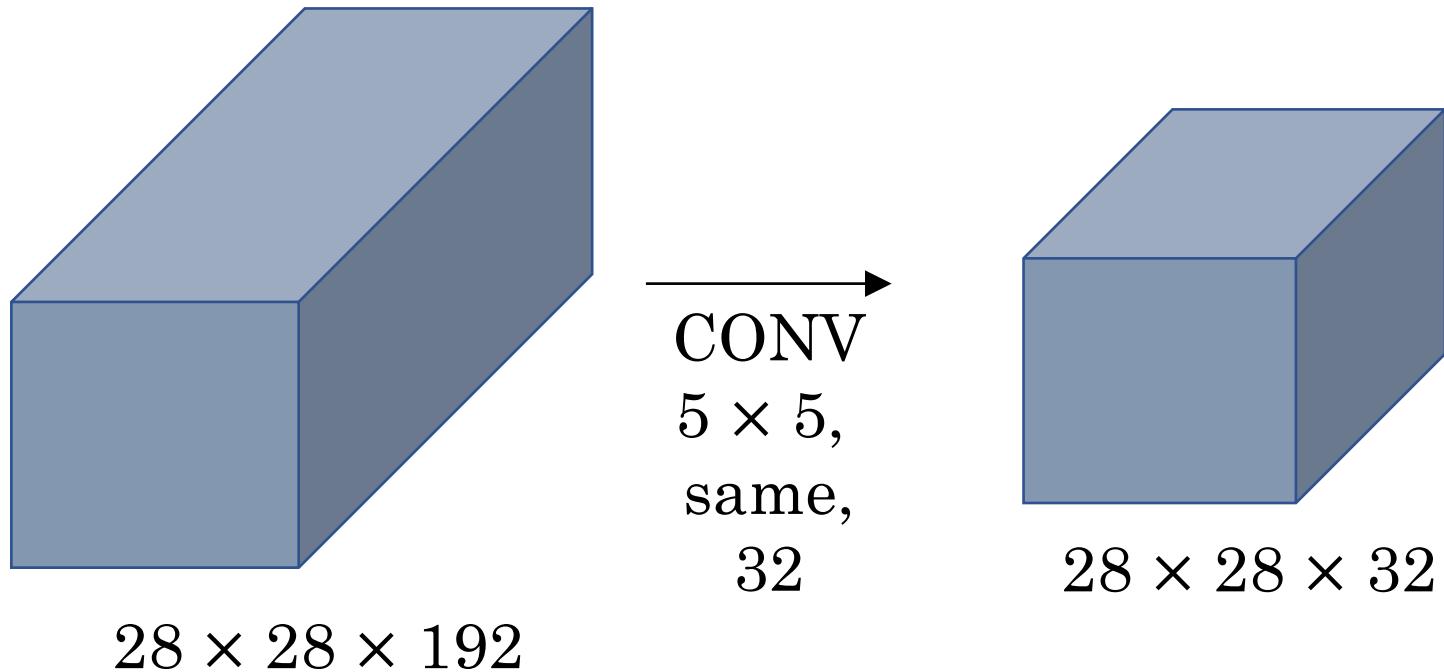
$6 \times 6 \times \# \text{ filters}$

# $1 \times 1$ convolutions can shrink the layer dims



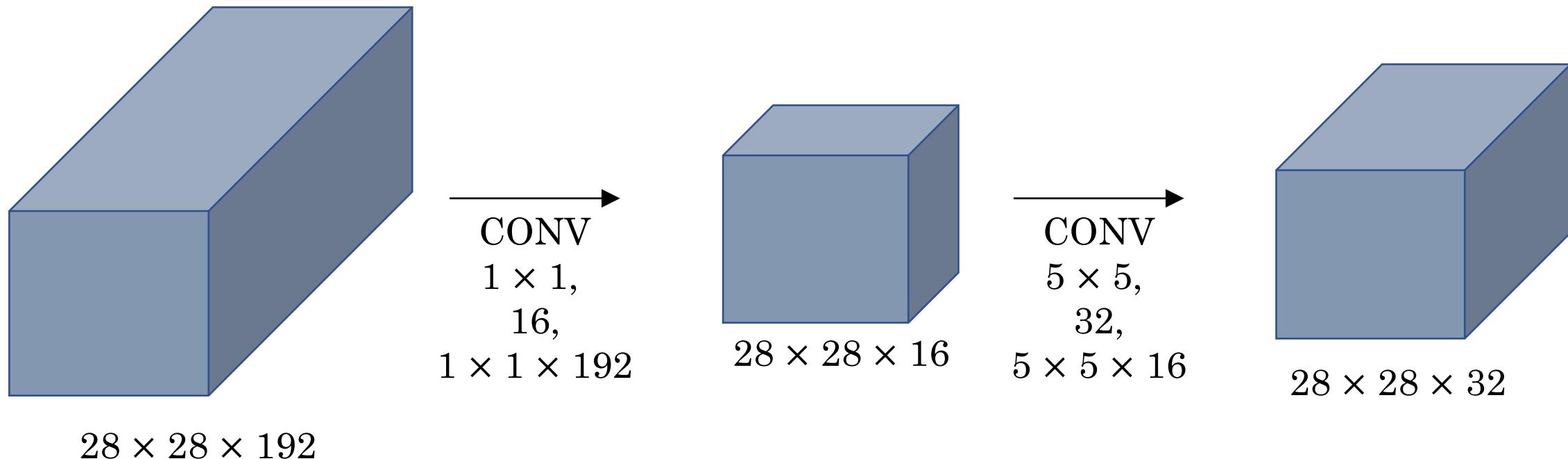
Helps reducing the number of parameters!

# The problem of computational cost



- One filter volume:  $5 \times 5 \times 192$  for one voxel in the output volume (in the activation map).
- Since we have  $28 \times 28 \times 32$  dimensional volume at the output, we need:  
 $(5 \times 5 \times 192) \times (28 \times 28 \times 32)$  multiplications! ( $\sim 120$ Million)

# Using $1 \times 1$ convolution for reducing the computation

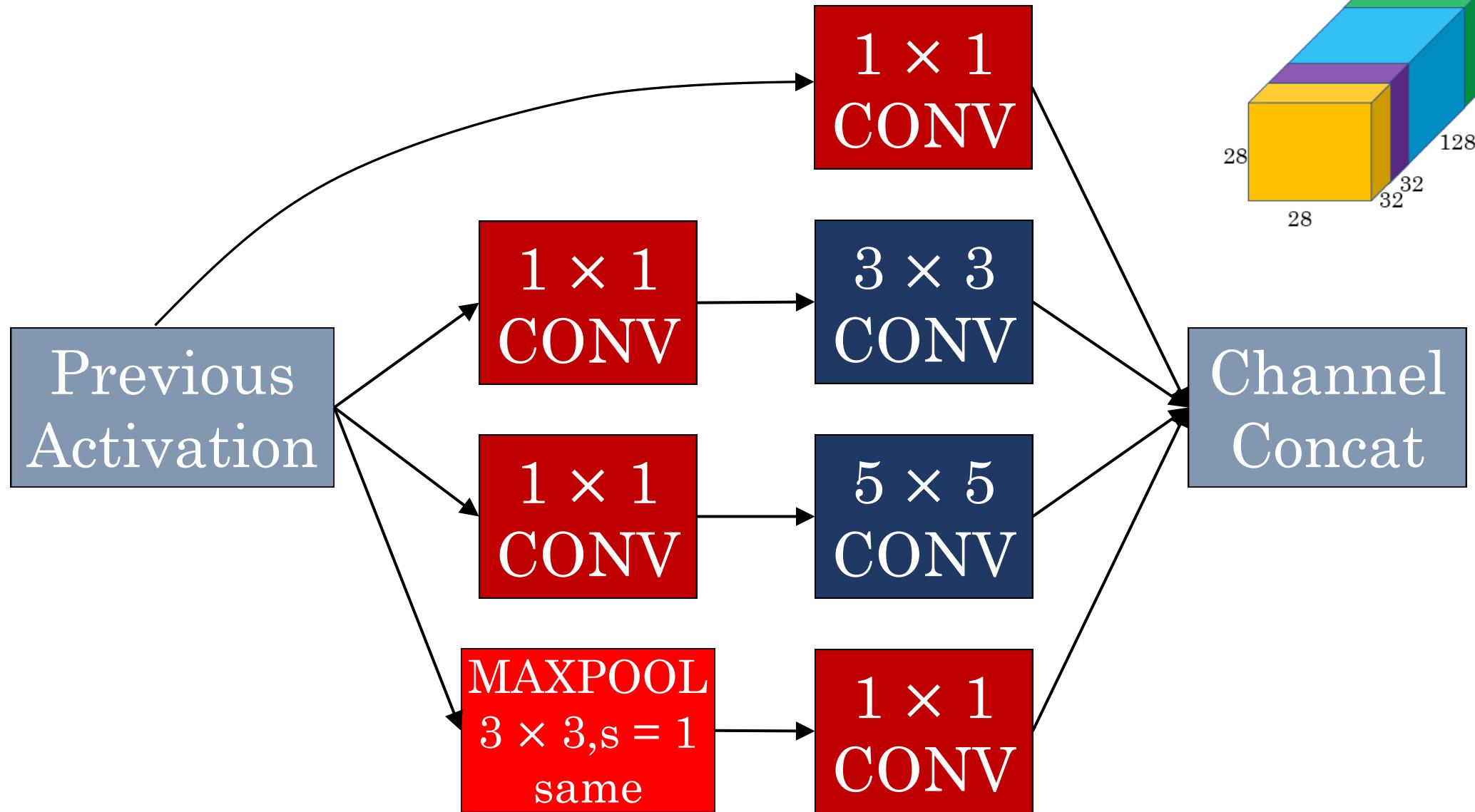


Multiplications needed:  $(28 \times 28 \times 16) \times (1 \times 1 \times 192) = \sim 2.4 \text{ Million}$

$(28 \times 28 \times 32) \times (5 \times 5 \times 16) = \sim 10 \text{ Million}$

**Total:** 12.4 Million << 120 Million

# Inception module

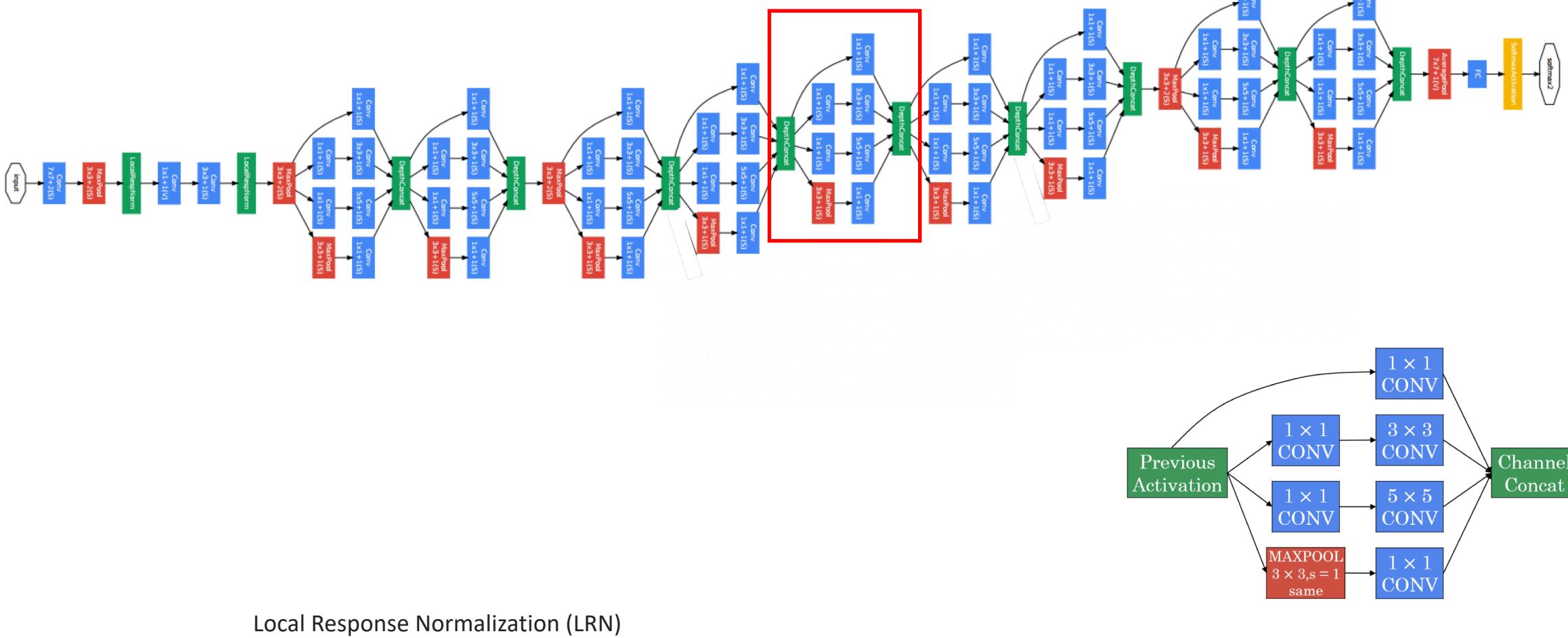


# Inception Network (Google)

GooLeNet

or

GoogLeNet



- This lecture has materials from CS231n (Stanford), Andrew Ng and Ulas Bagci.