

CAP 5415: Computer Vision

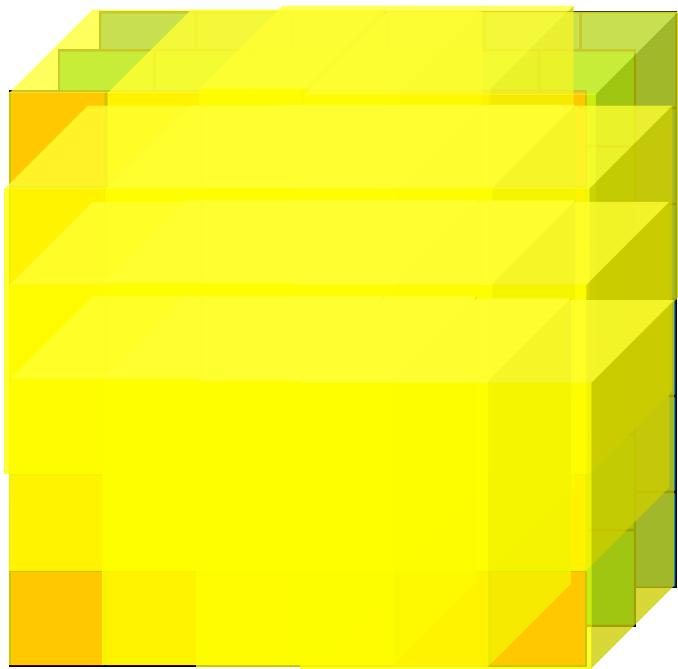
Deep (Neural) Networks



Dr. Sedat Ozer

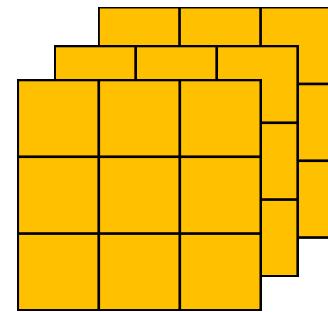


Convolutions on RGB image

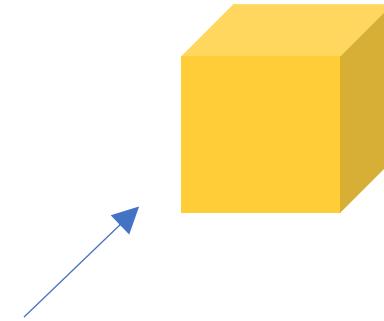


$6 \times 6 \times 3$

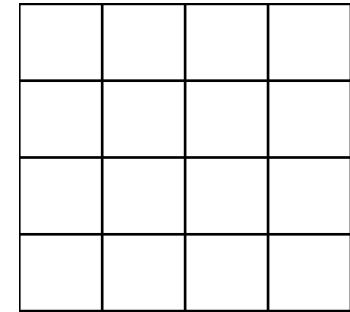
*



$3 \times 3 \times 3$

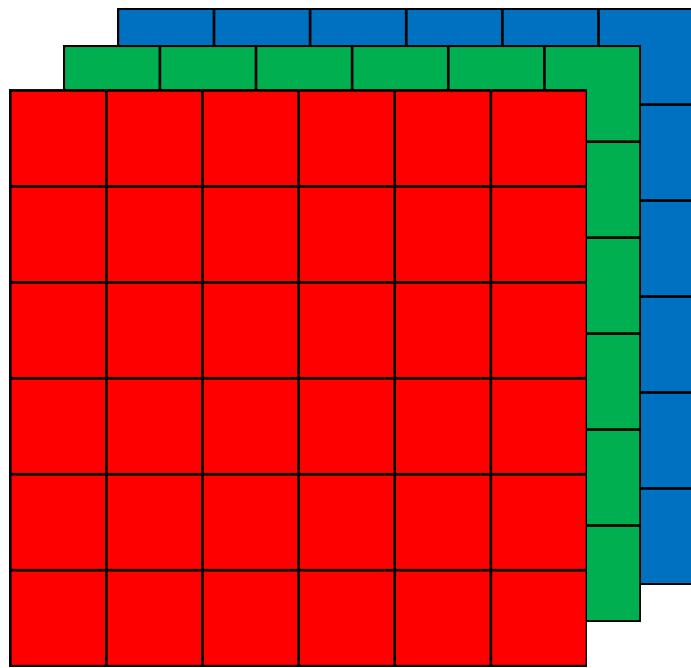


=



4×4

Multiple filters

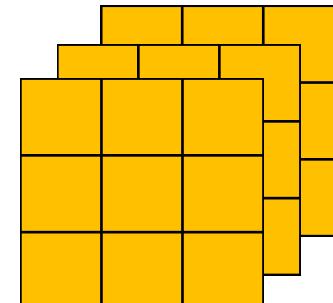


$6 \times 6 \times 3$

Number of channels

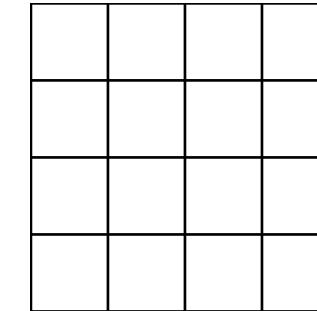
$n \times n \times n_c$

*



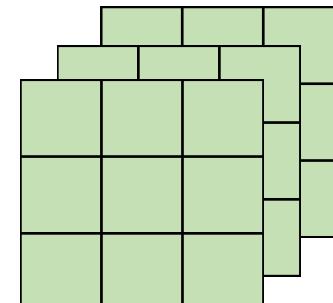
$3 \times 3 \times 3$

=



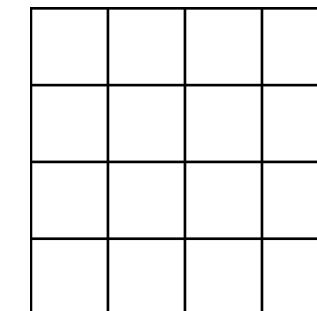
4×4

*



$3 \times 3 \times 3$

=



4×4

Number of filters

*

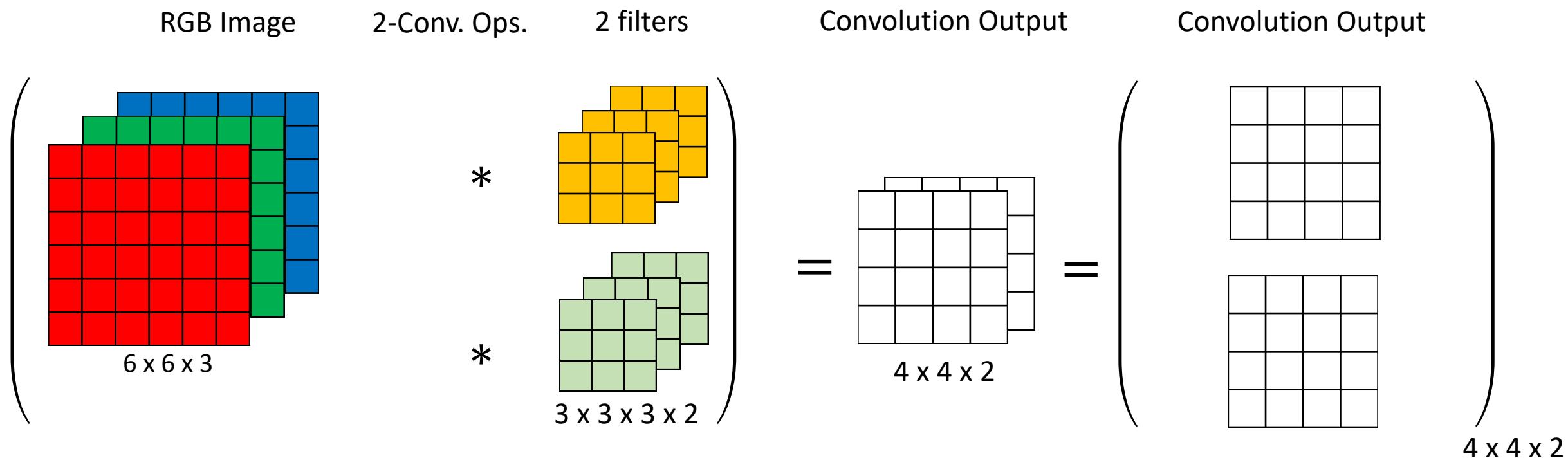
$f \times f \times n_c \times 2$

Number of filters

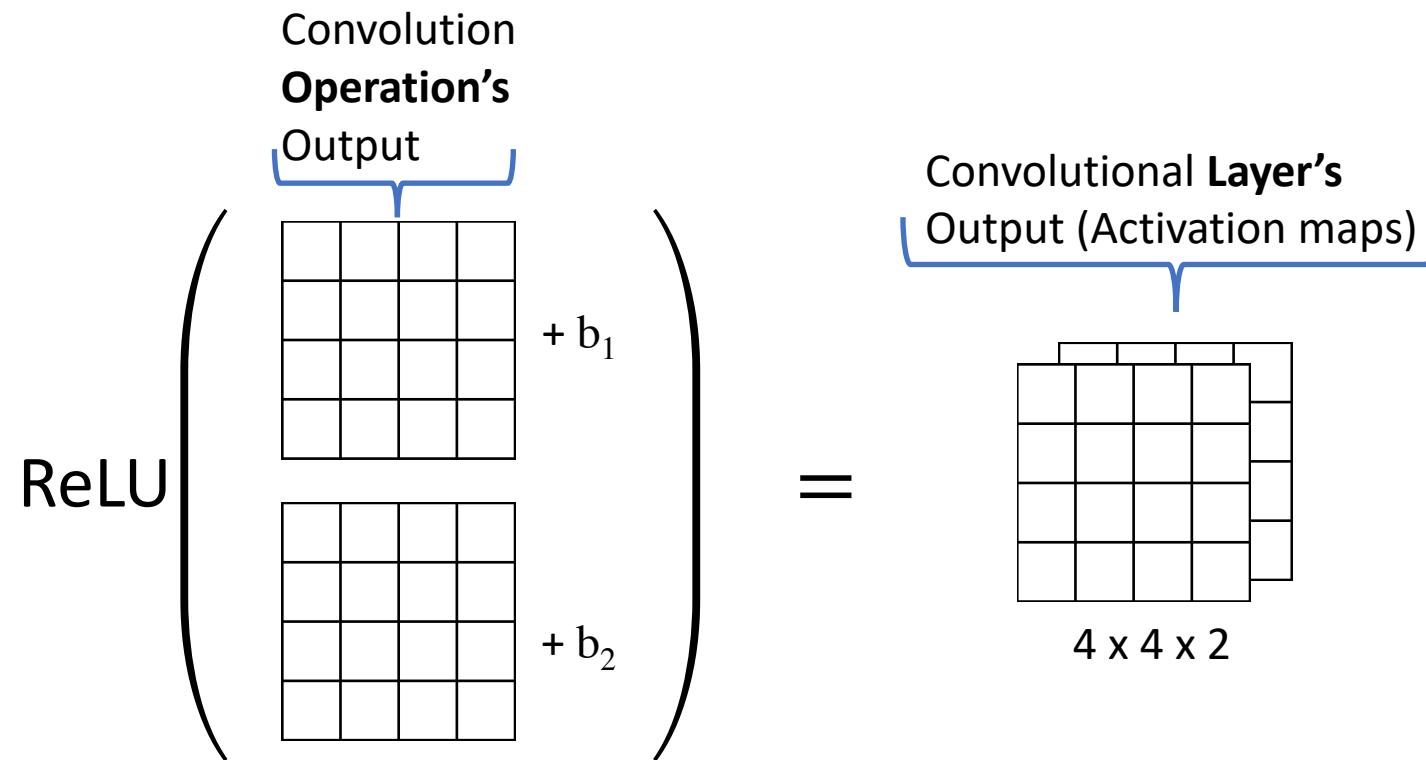
=

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

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

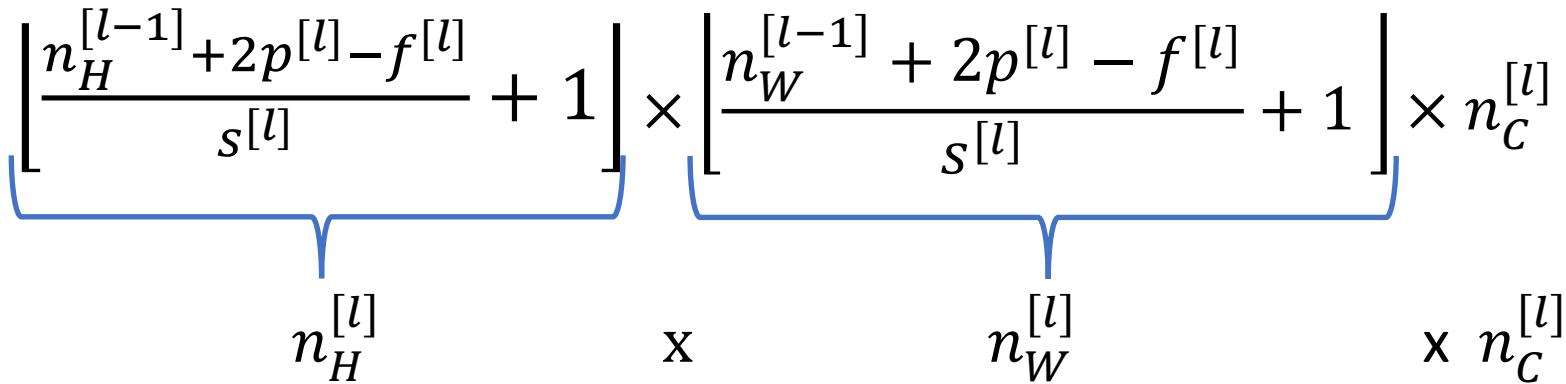
Padding size: $p^{[l]} \times p^{[l]}$

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



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

Pop Quiz

- How many parameters do we have in a conv layer with 10 filters that are 3x3x20 dimensional each?
- For one filter: we have $3 \times 3 \times 20 = 180$ filter weights + 1 bias = **181** parameters.
- For all **10** filters: we have **$181 \times 10 = 1810$ 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.
 - 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 dims of the output at the end of the 5th layer?

30 seconds to return your answers! 😊

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!

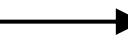
...

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

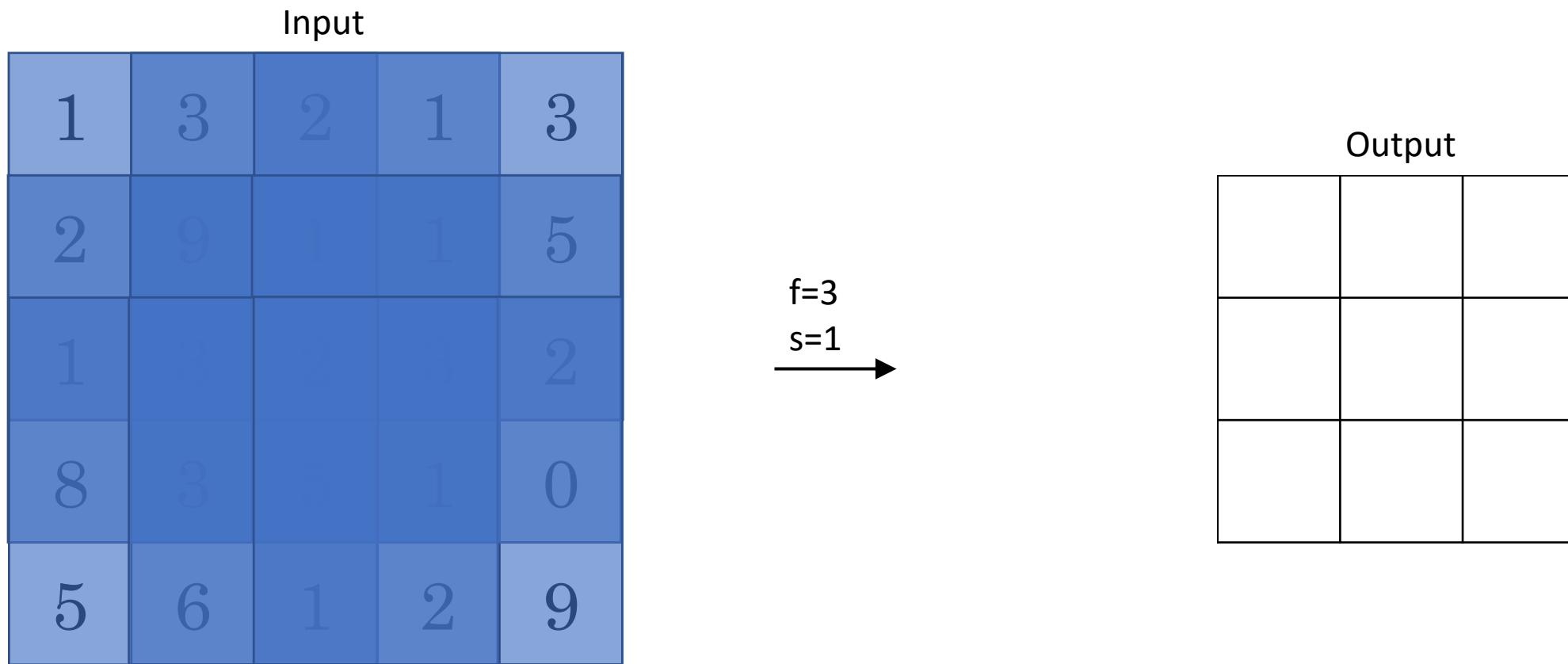
Input			
1	3	2	1
2	9	1	1
1	3	2	3
5	6	1	2

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

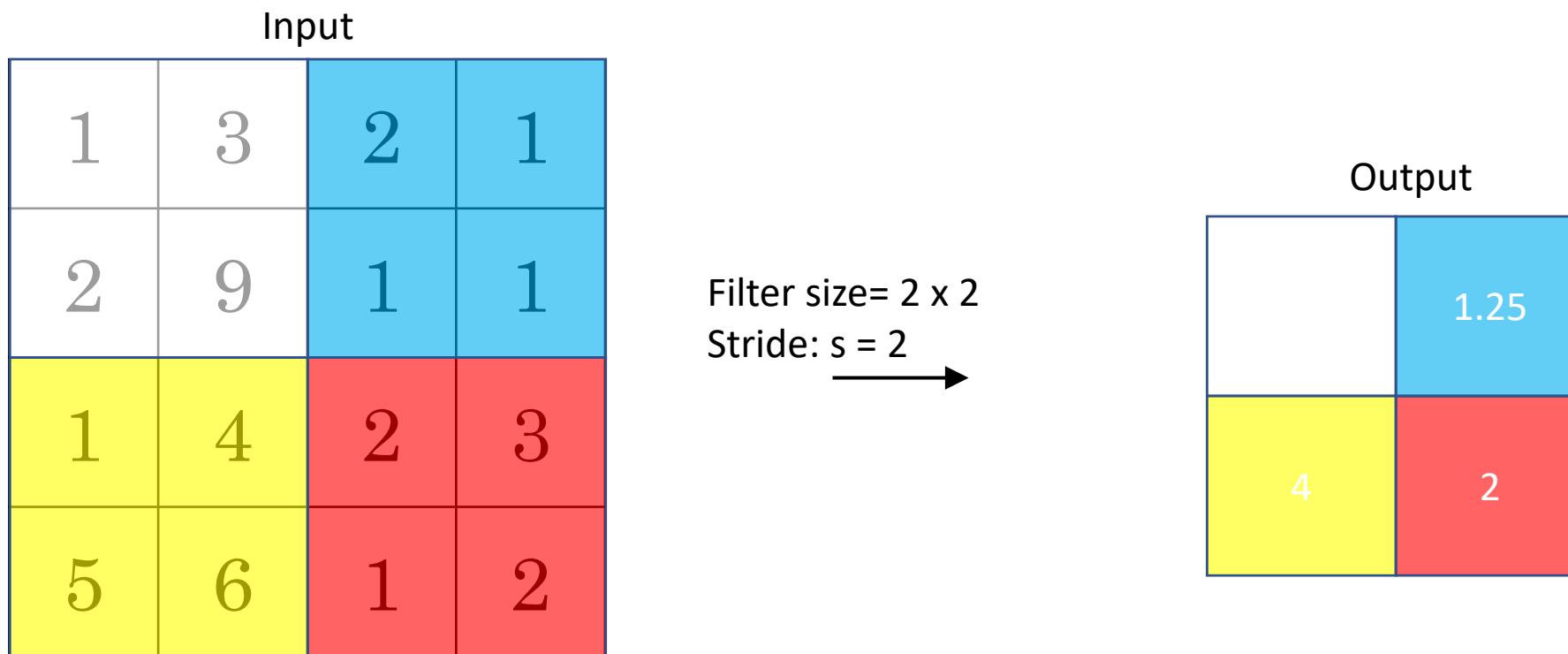
Output	

Pooling is typically applied on each channel separately.

Pooling layer: Max pooling

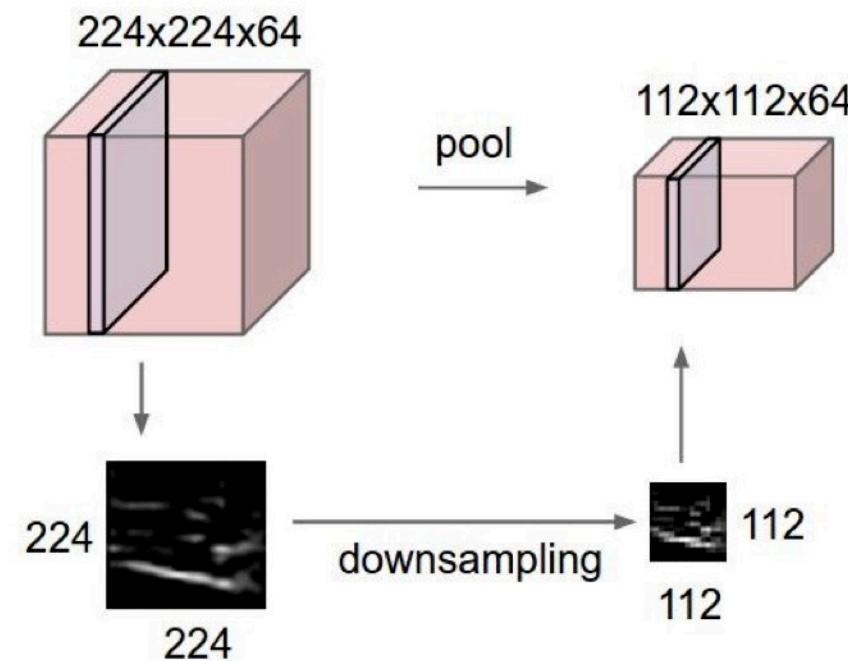


Pooling layer: Average pooling



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

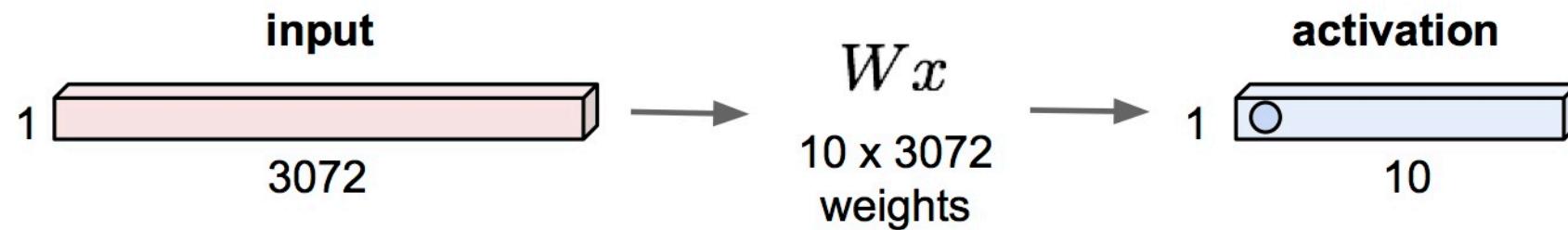
Max or Average pooling

Input: $n_H \times n_W \times n_C$

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

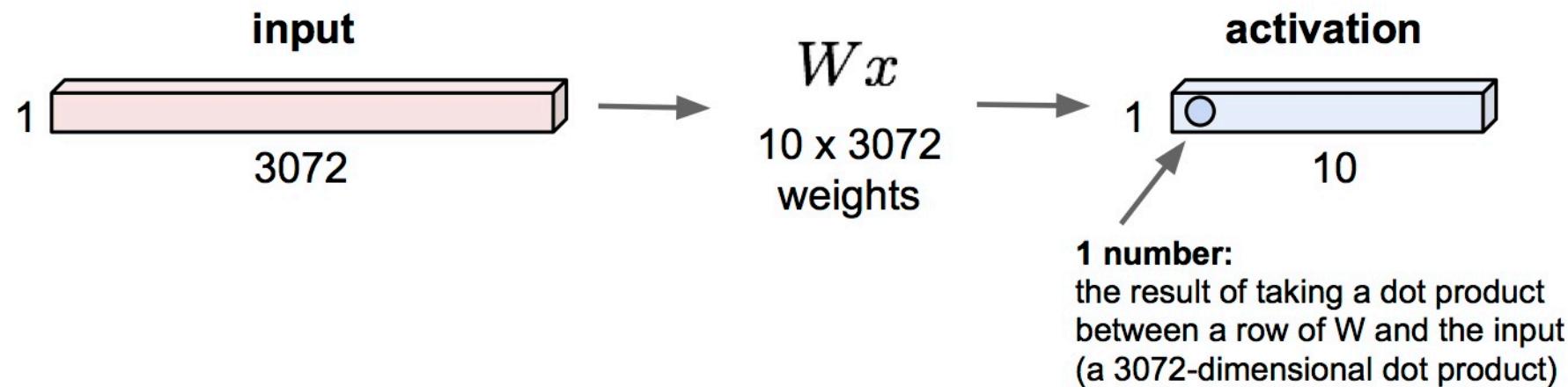
Fully Connected Layer

32x32x3 image -> stretch to 3072 x 1



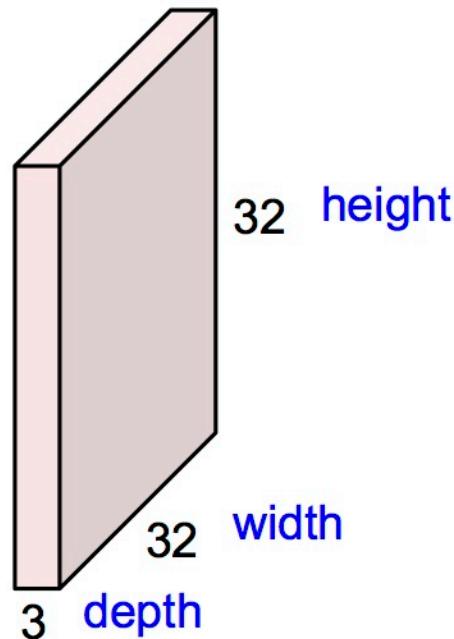
Fully Connected Layer

32x32x3 image -> stretch to 3072×1



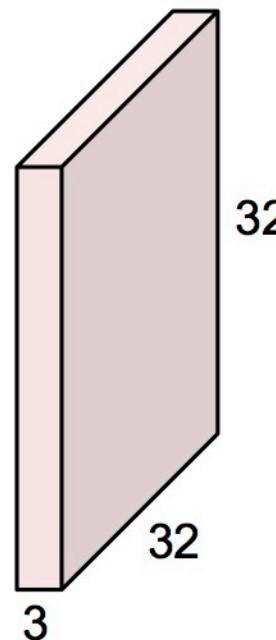
Convolution Layer

32x32x3 image -> preserve spatial structure

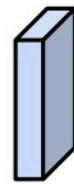


Convolution Layer

32x32x3 image



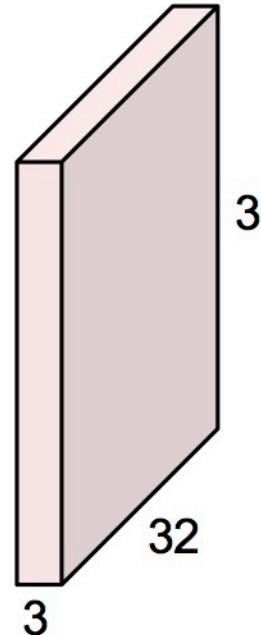
5x5x3 filter



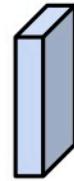
Convolve the filter with the image
i.e. “slide over the image spatially,
computing dot products”

Convolution Layer

32x32x3 image



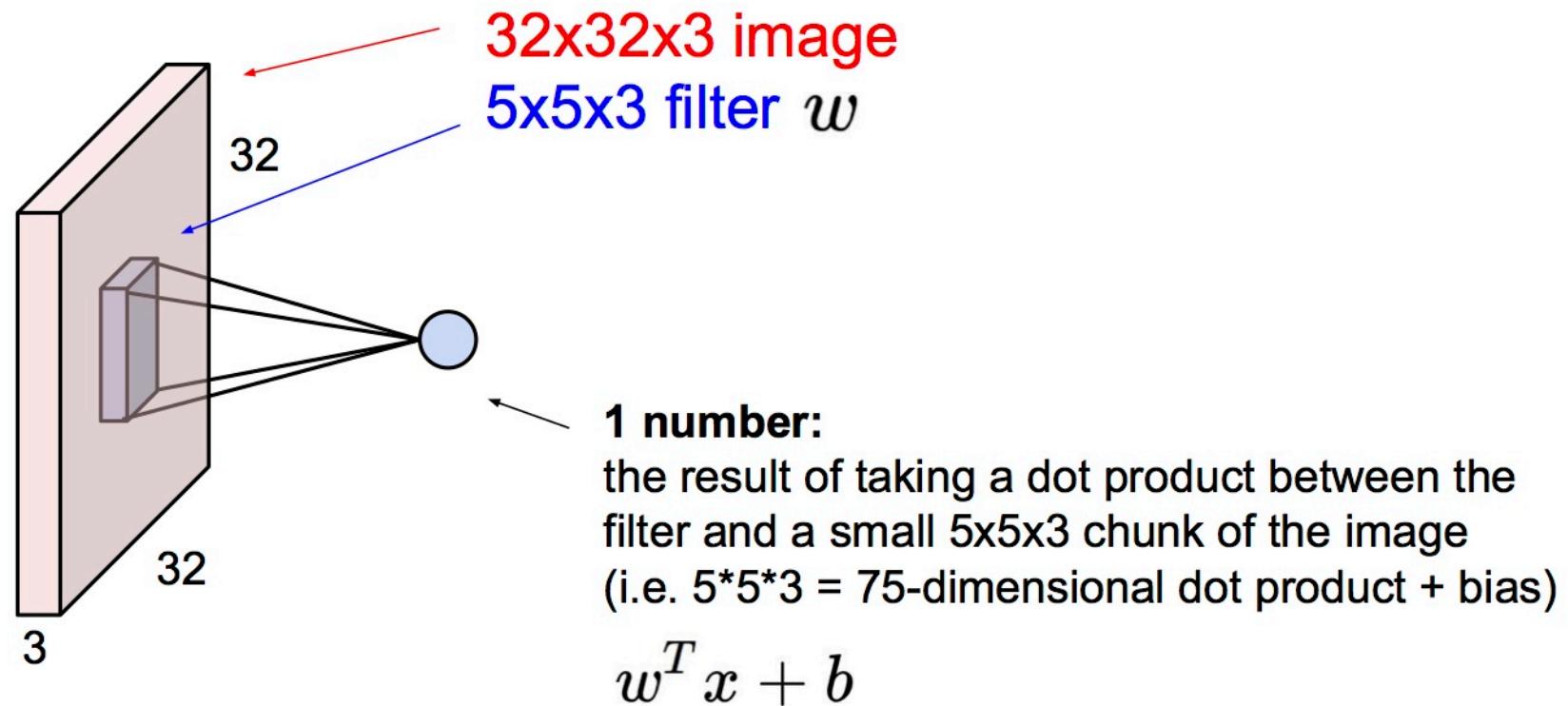
5x5x3 filter



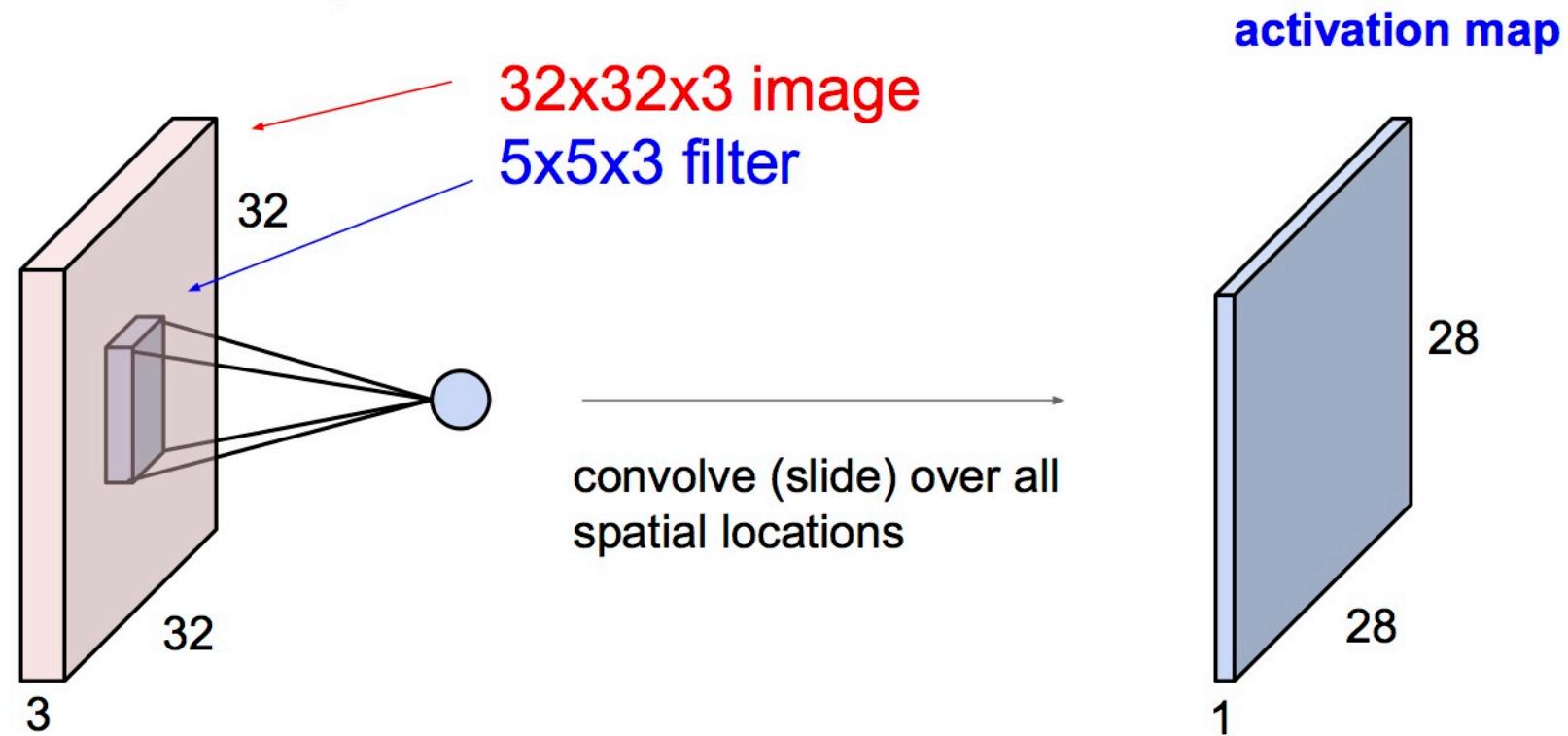
Filters always extend the full depth of the input volume

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

Convolution Layer

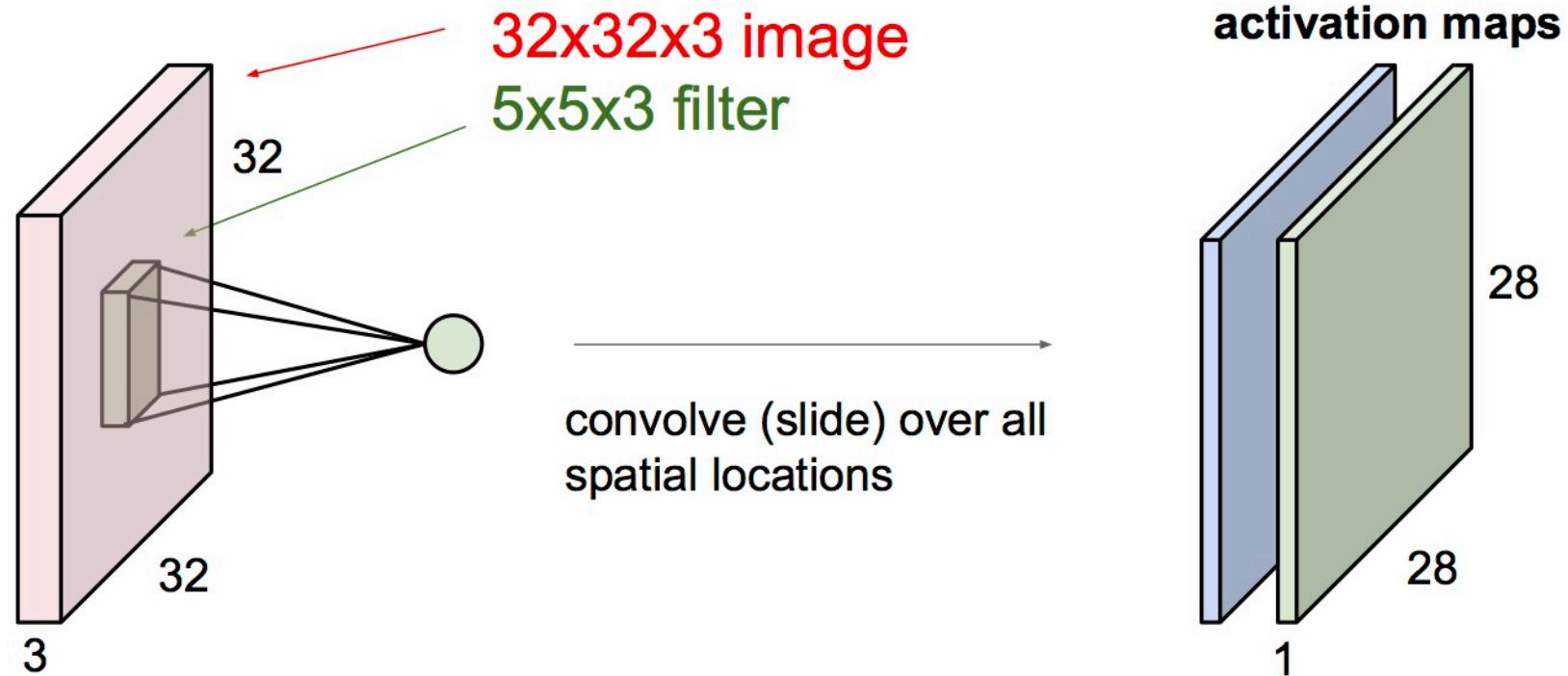


Convolution Layer

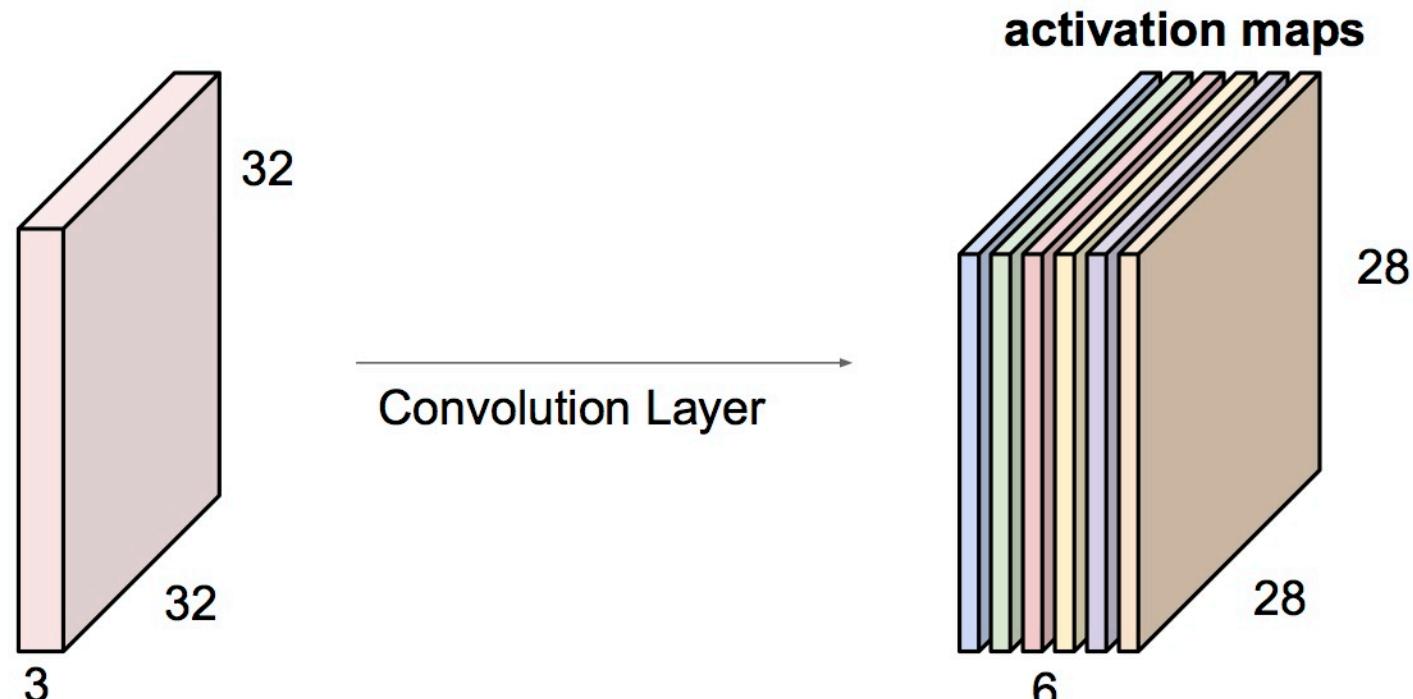


Convolution Layer

consider a second, **green** filter

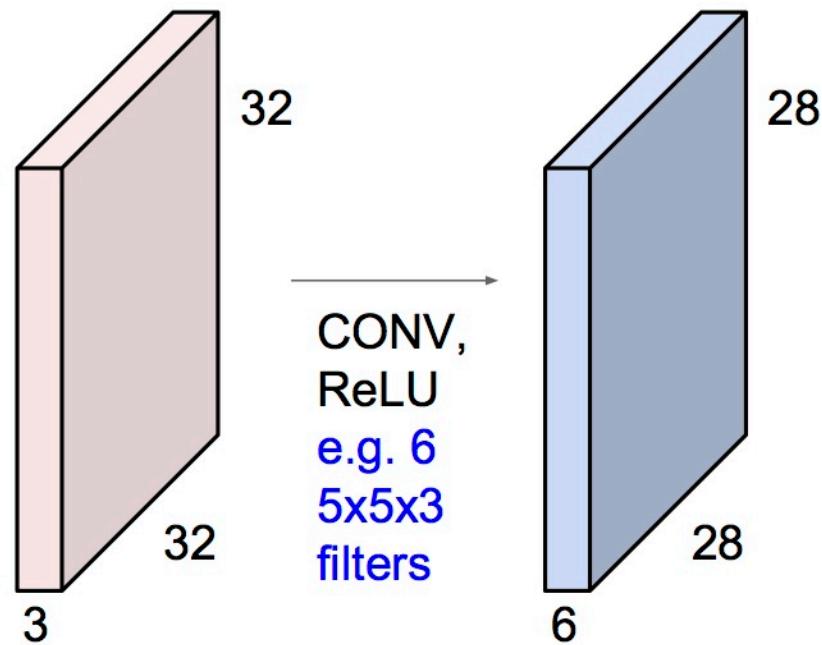


For example, if we had 6 5x5 filters, we'll get 6 separate activation maps:

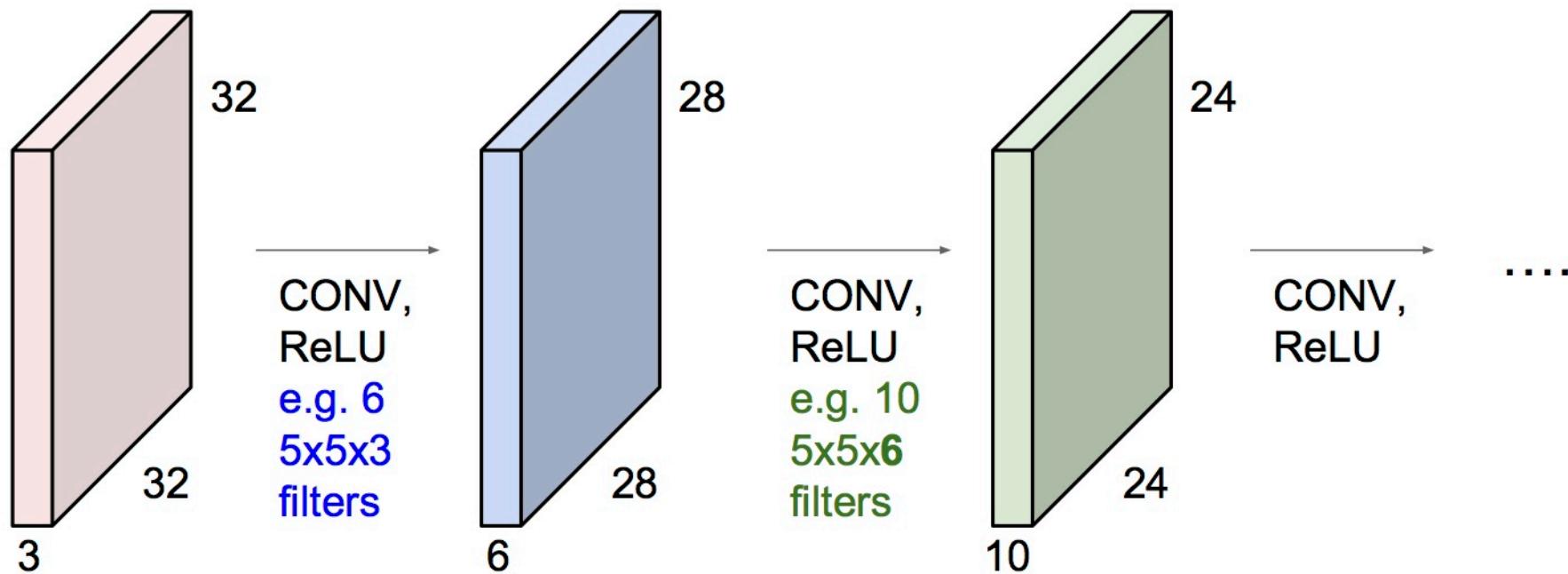


We stack these up to get a “new image” of size 28x28x6!

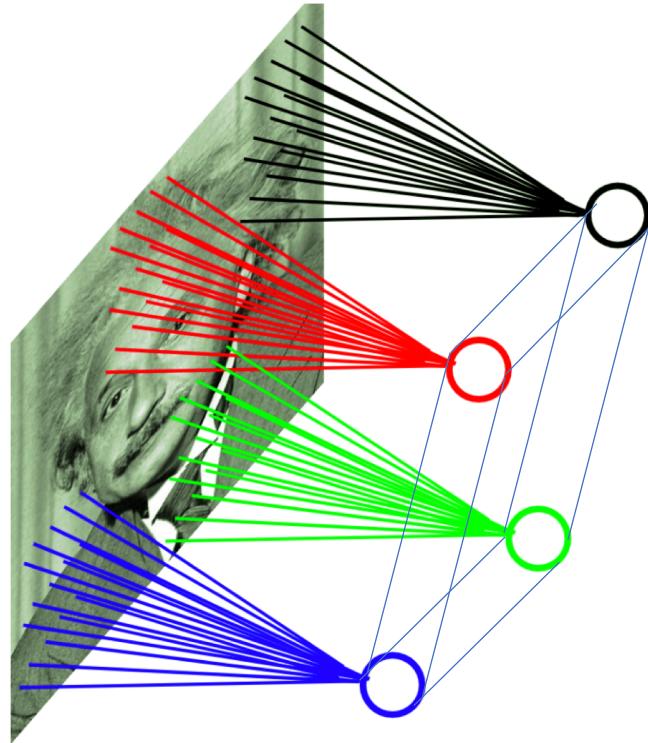
Preview: ConvNet is a sequence of Convolution Layers, interspersed with activation functions



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



Parameter Sharing

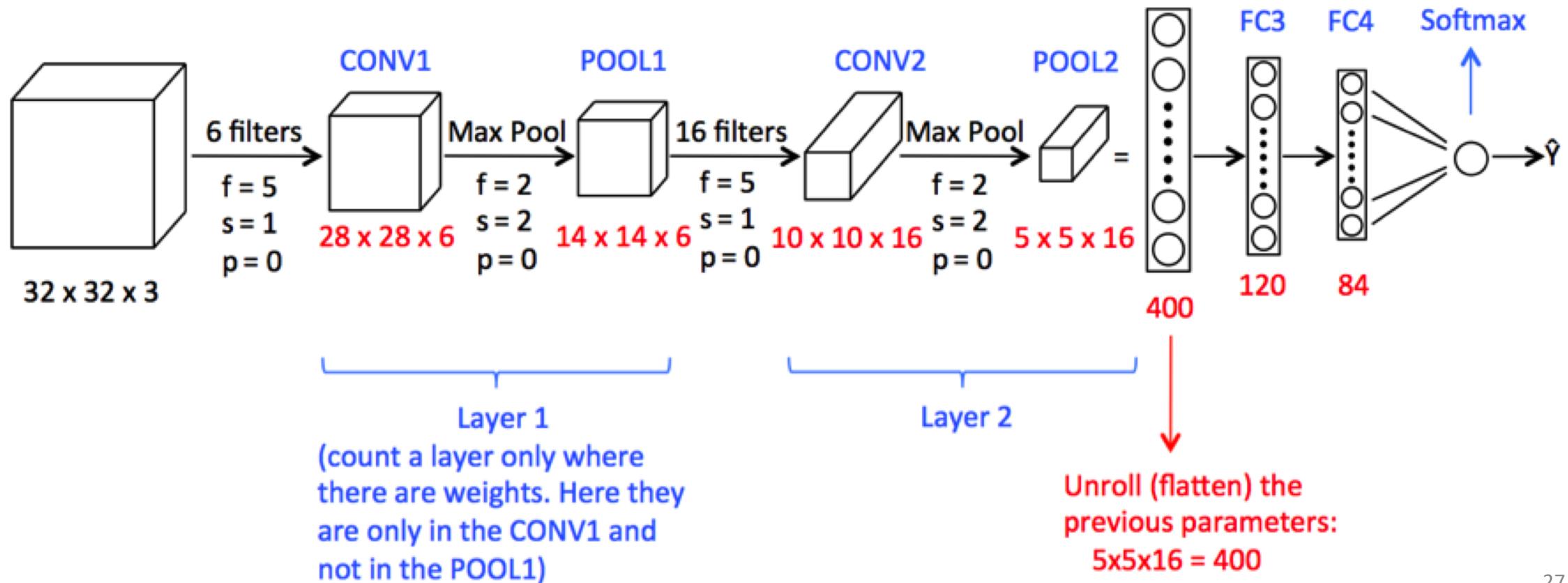


One filter to rule them all!



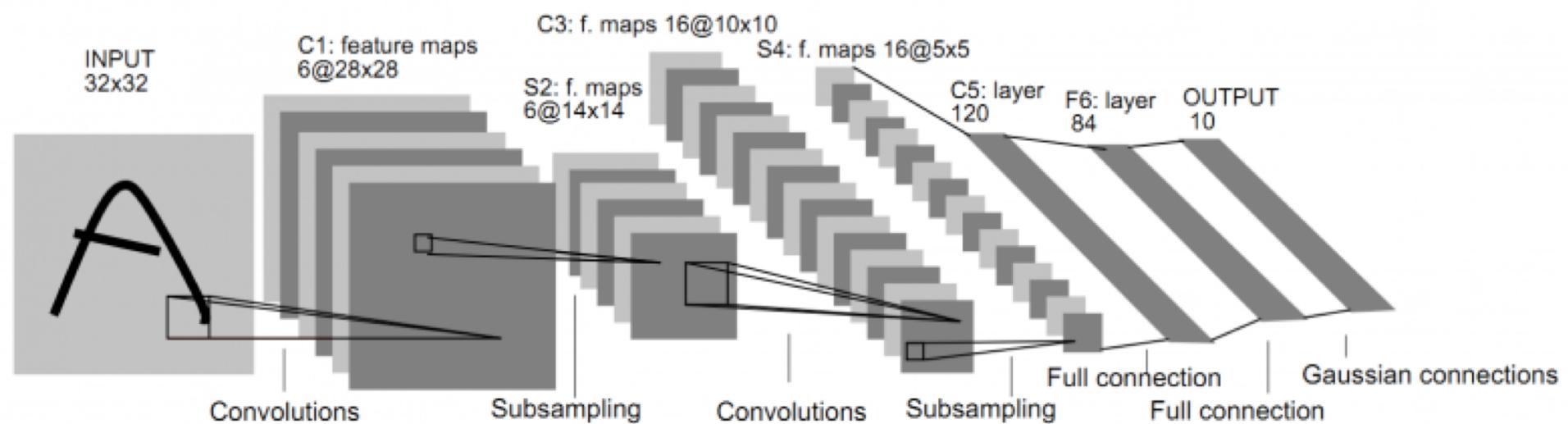
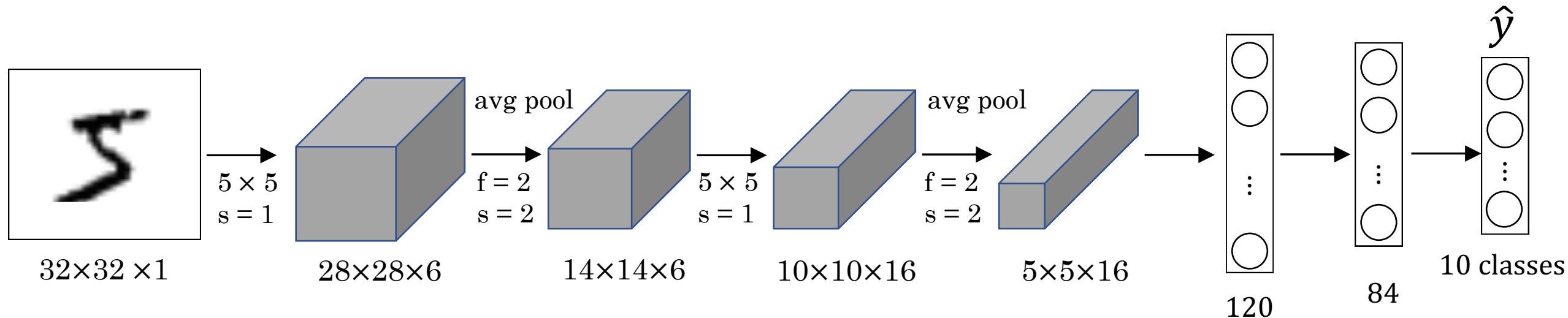
A CNN Example (LeNet-5)

LeNet-5

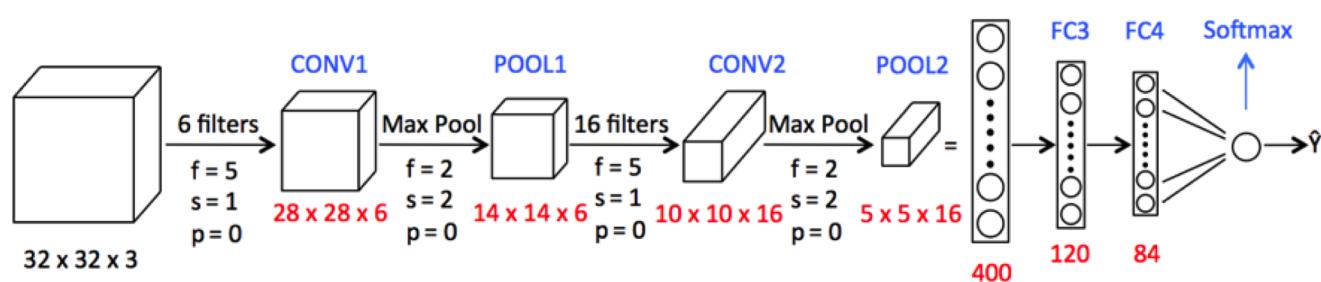


- Image Classification: A main task in computer vision
 - PASCAL
 - IMAGENET

LeNet - 5



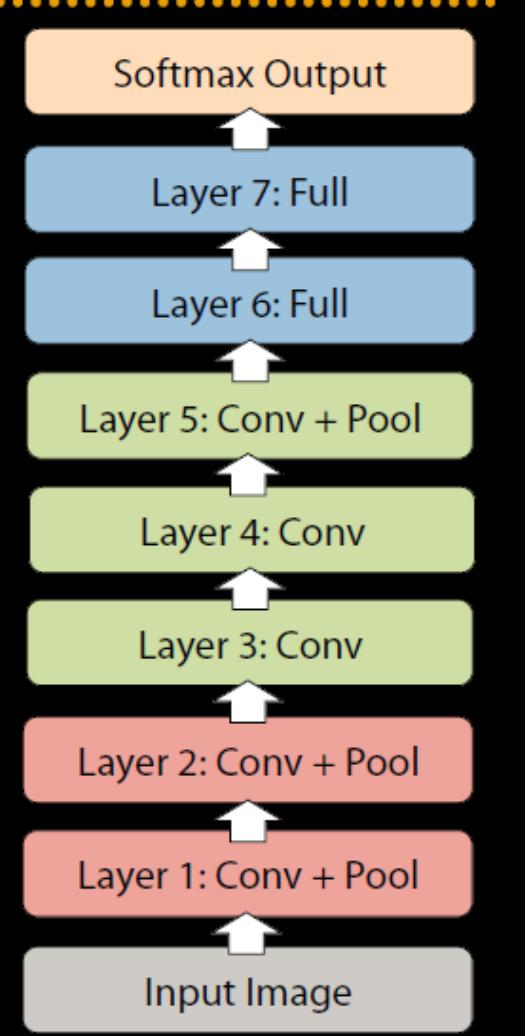
(Modified) LeNet-5



	Activation shape	Activation Size	# parameters
Input:	$(32, 32, 3)$	3,072	0

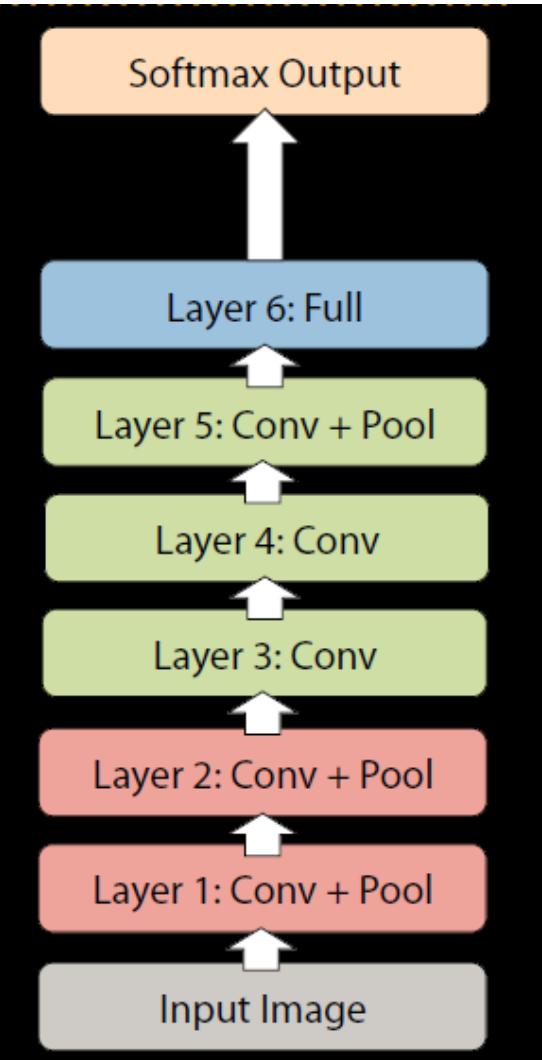
Why ConvNet should be Deep?

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



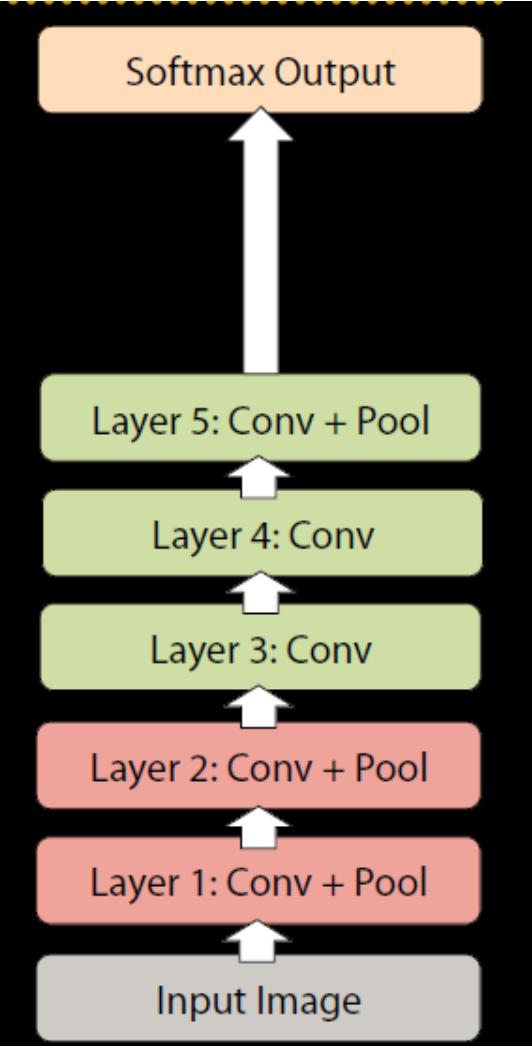
Why ConvNet should be Deep?

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



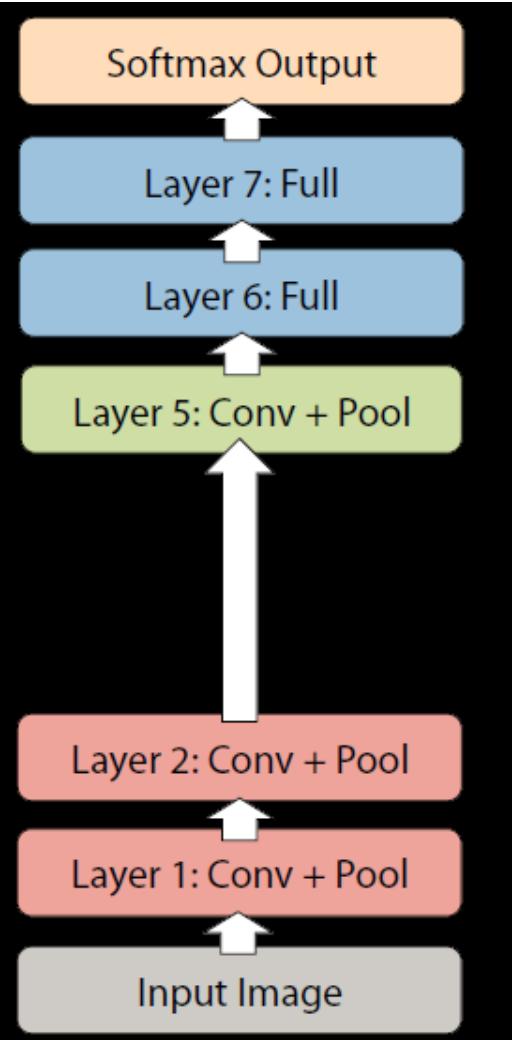
Why ConvNet should be Deep?

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



Why ConvNet should be Deep?

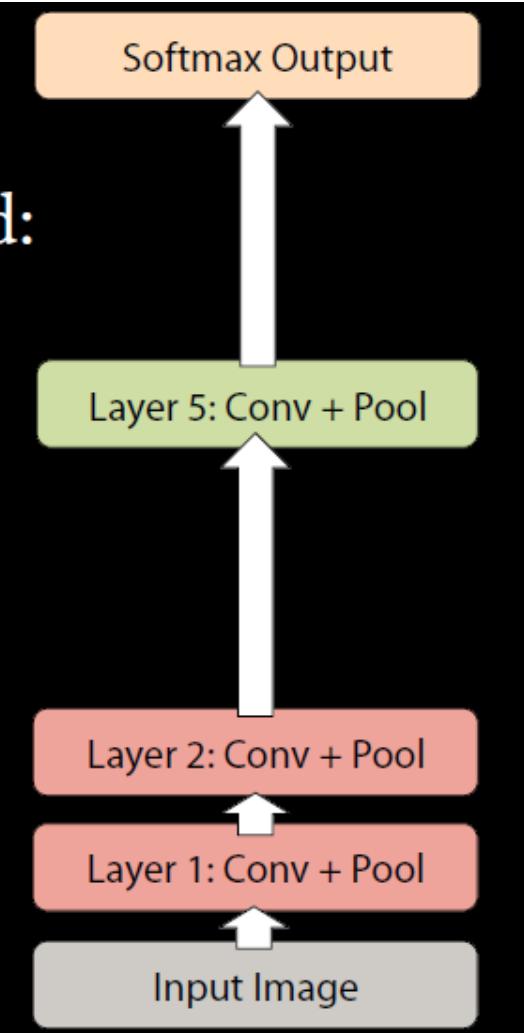
- Now try removing upper feature extractor layers:
 - Layers 3 & 4
- Drop ~1 million parameters
- 3.0% drop in performance



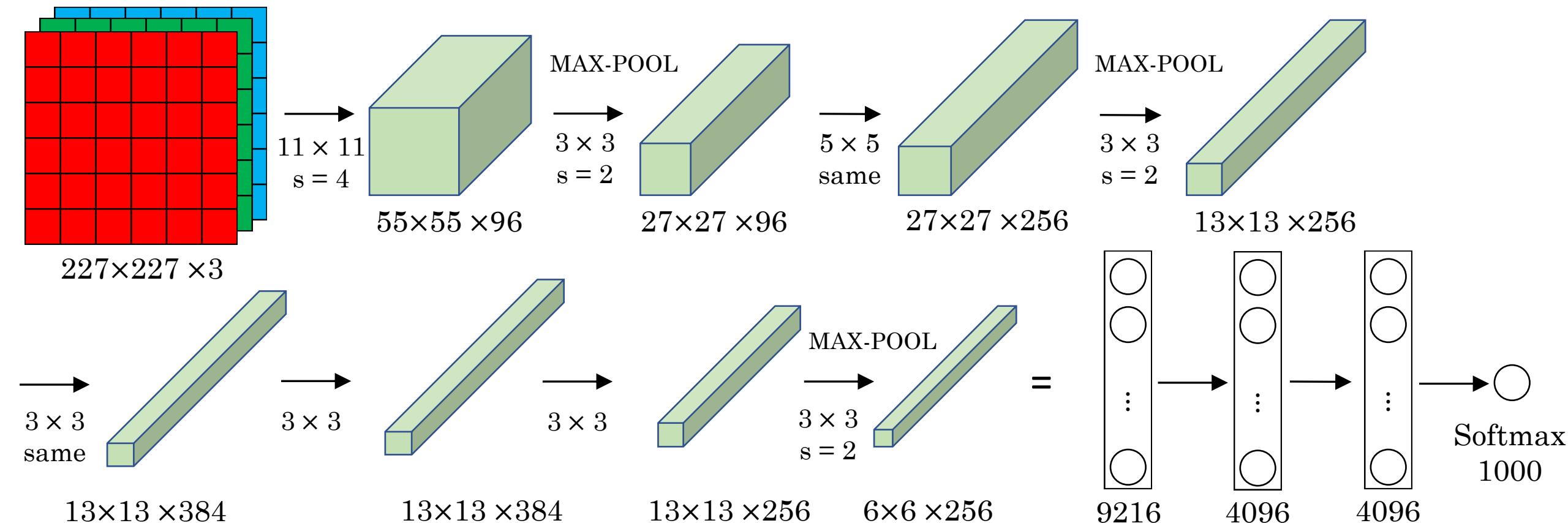
Why ConvNet should be Deep?

- 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



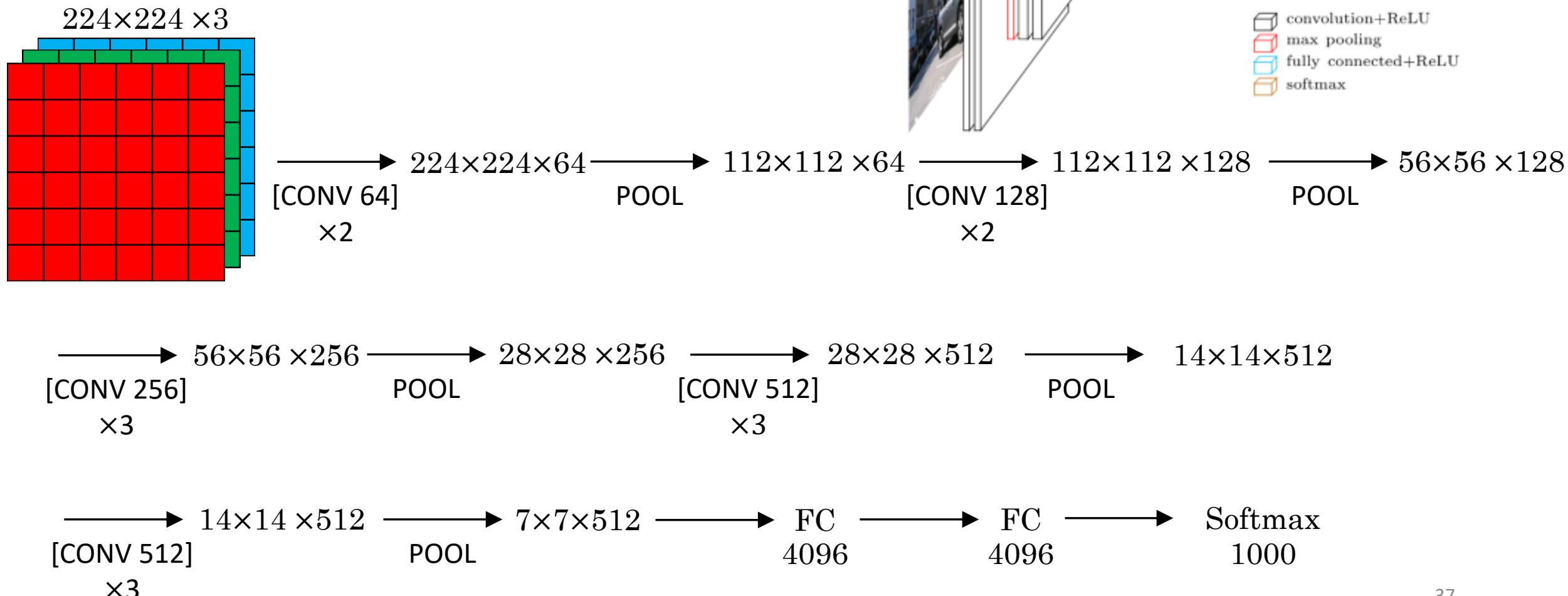
AlexNet



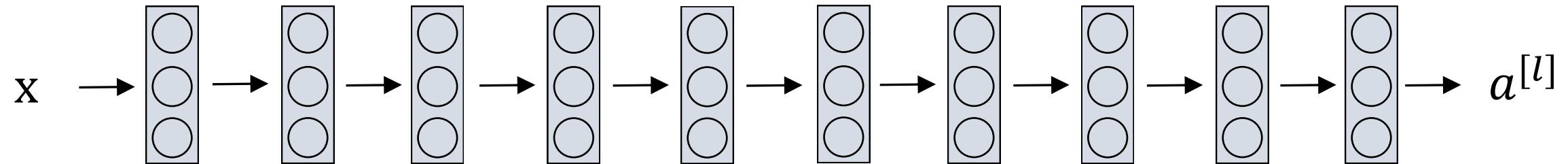
VGG - 16

CONV = 3×3 filter, s = 1, same

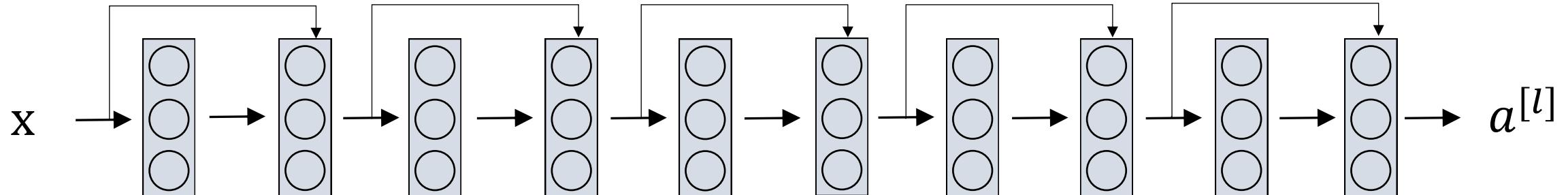
MAX-POOL = 2×2 , s = 2



Residual Network

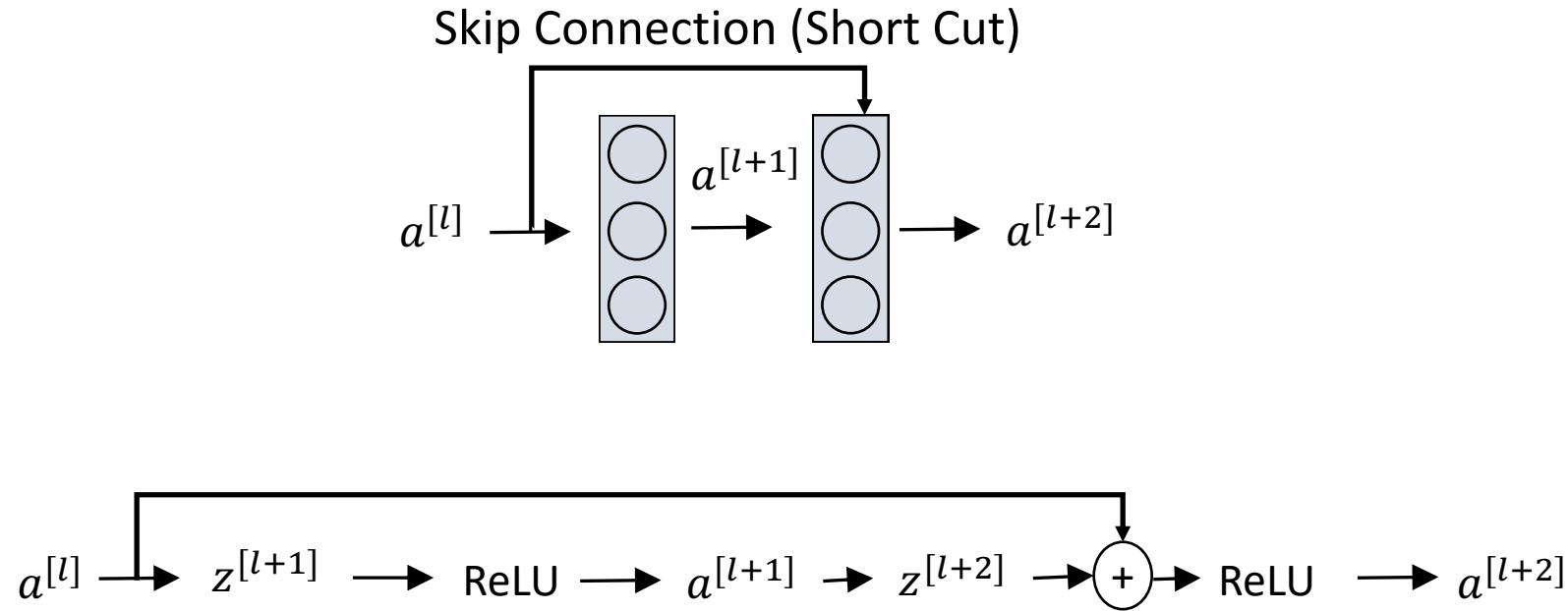


A “plain” (deep) network



A residual network (with 5 residual blocks)

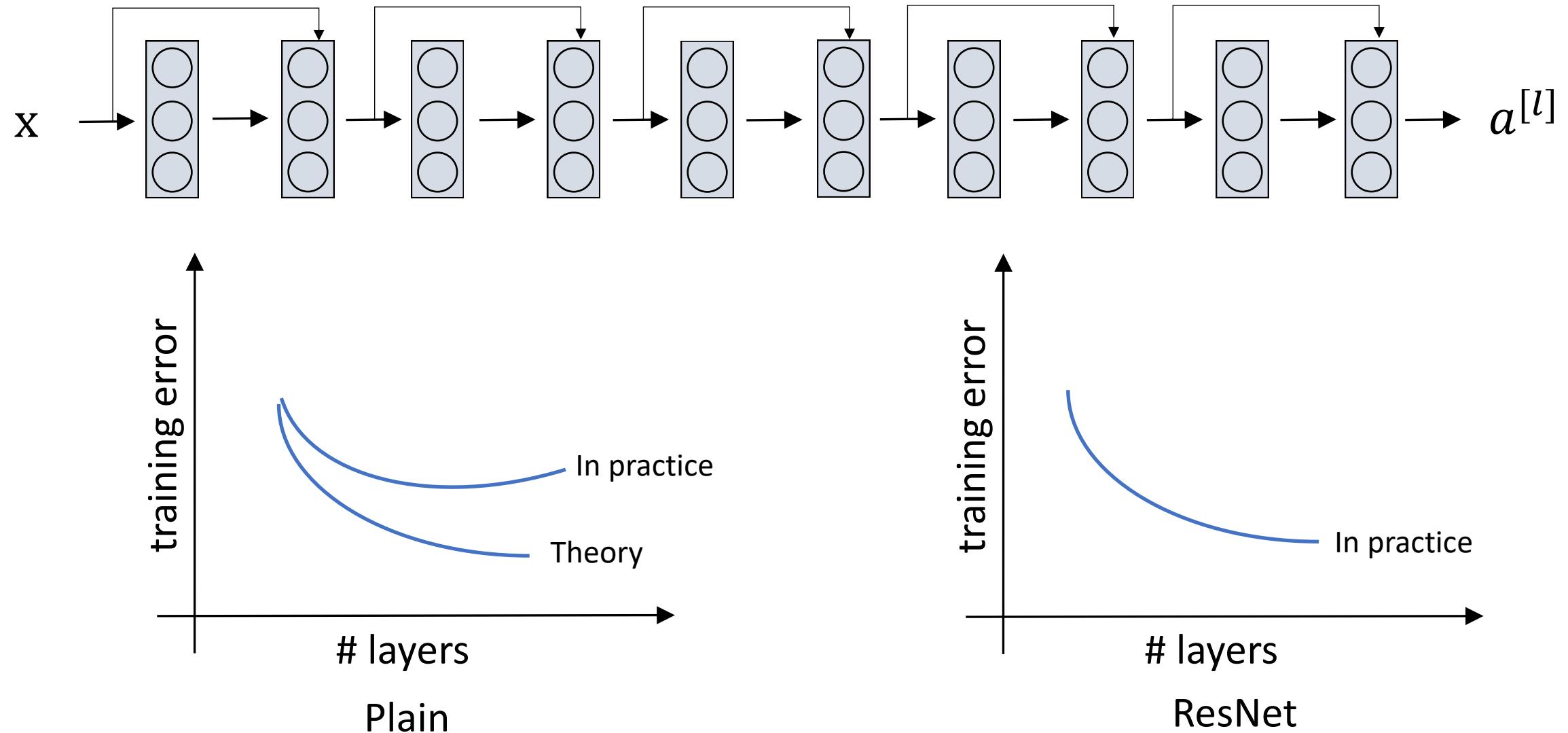
Residual Building Block



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

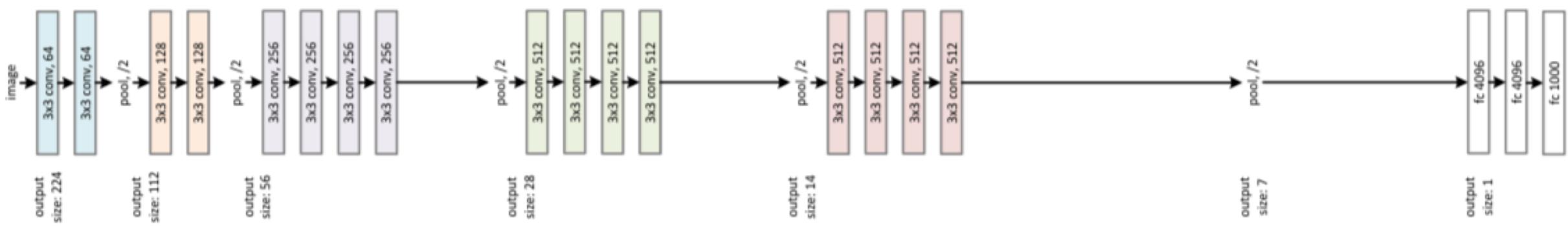
With skip connection: $a^{[l+2]} = g(z^{[l+2]} + a^{[l]})$

Residual Network

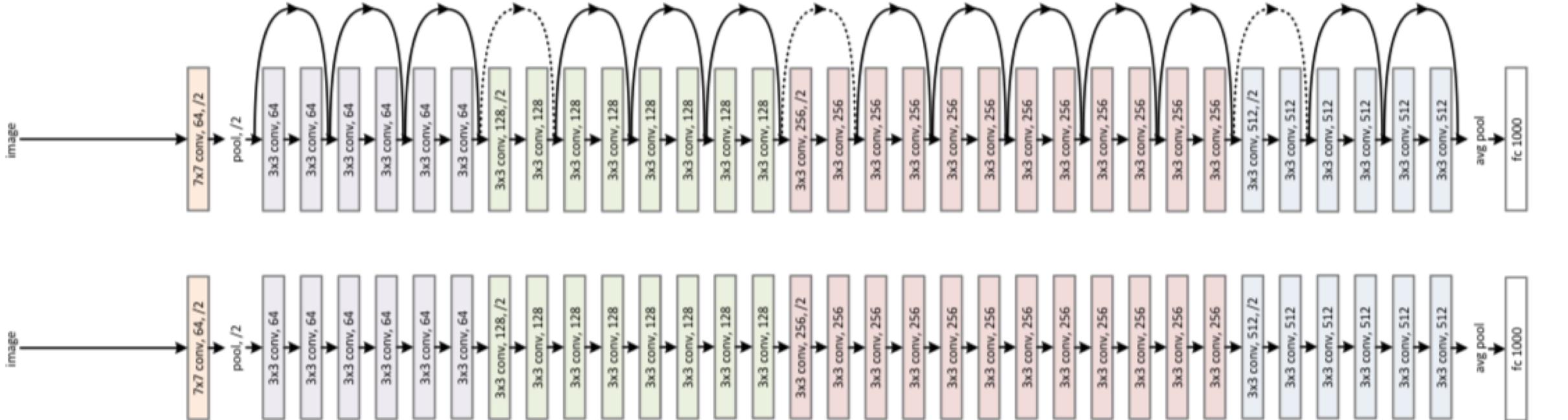


ResNet

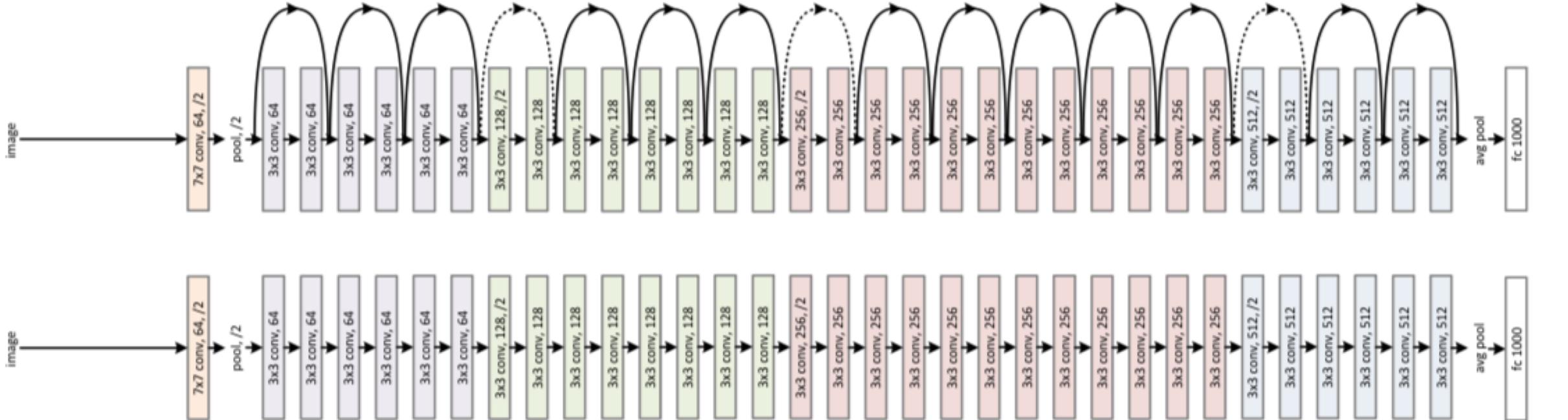
VGG-19



34-layer plain



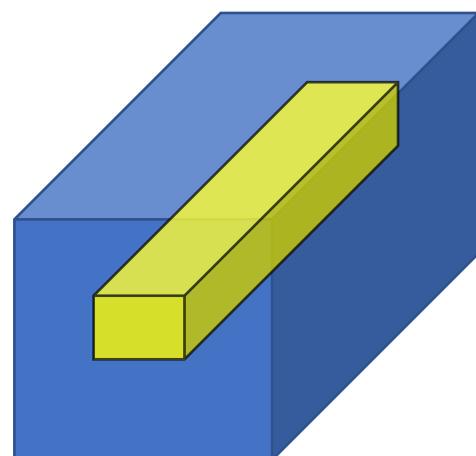
34-layer residual



What does a 1×1 convolution do?

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

6×6



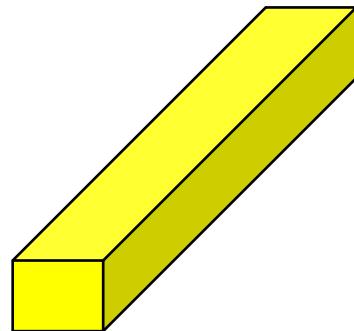
$6 \times 6 \times 32$

*

2

=

*

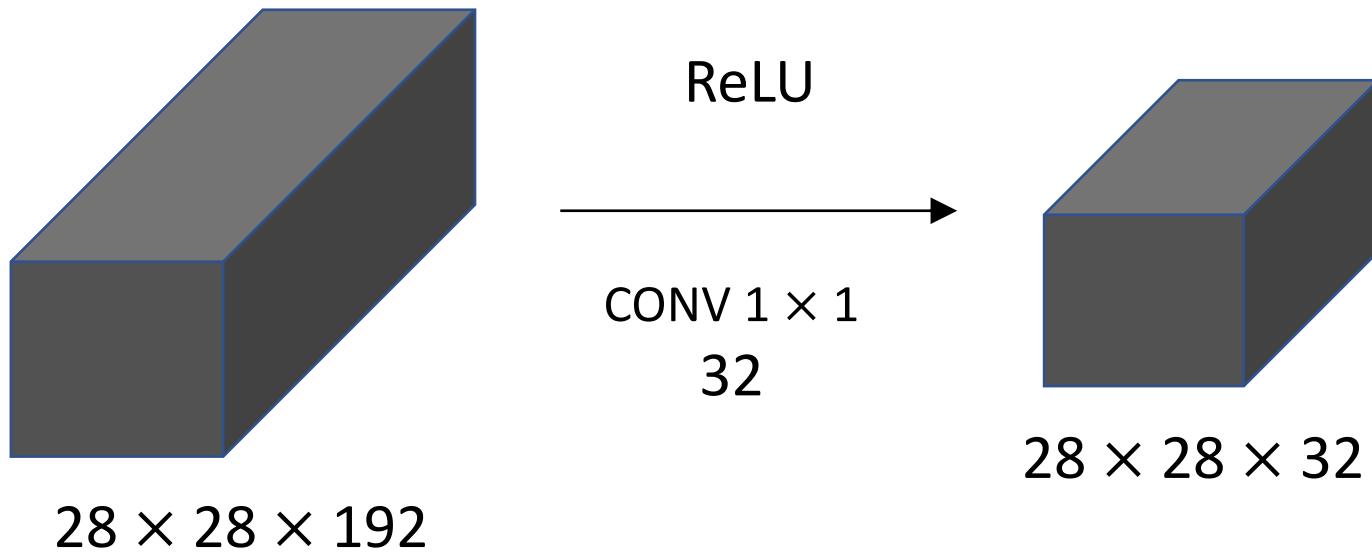


=

$1 \times 1 \times 32$

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

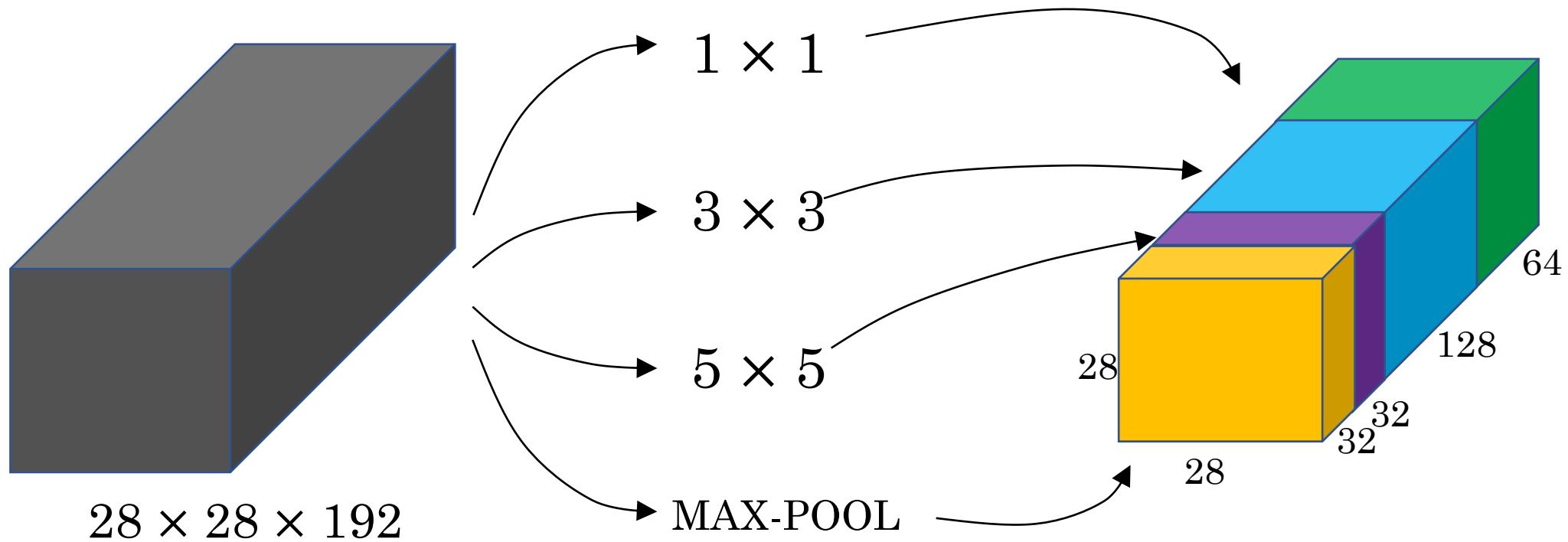
Using 1×1 convolutions



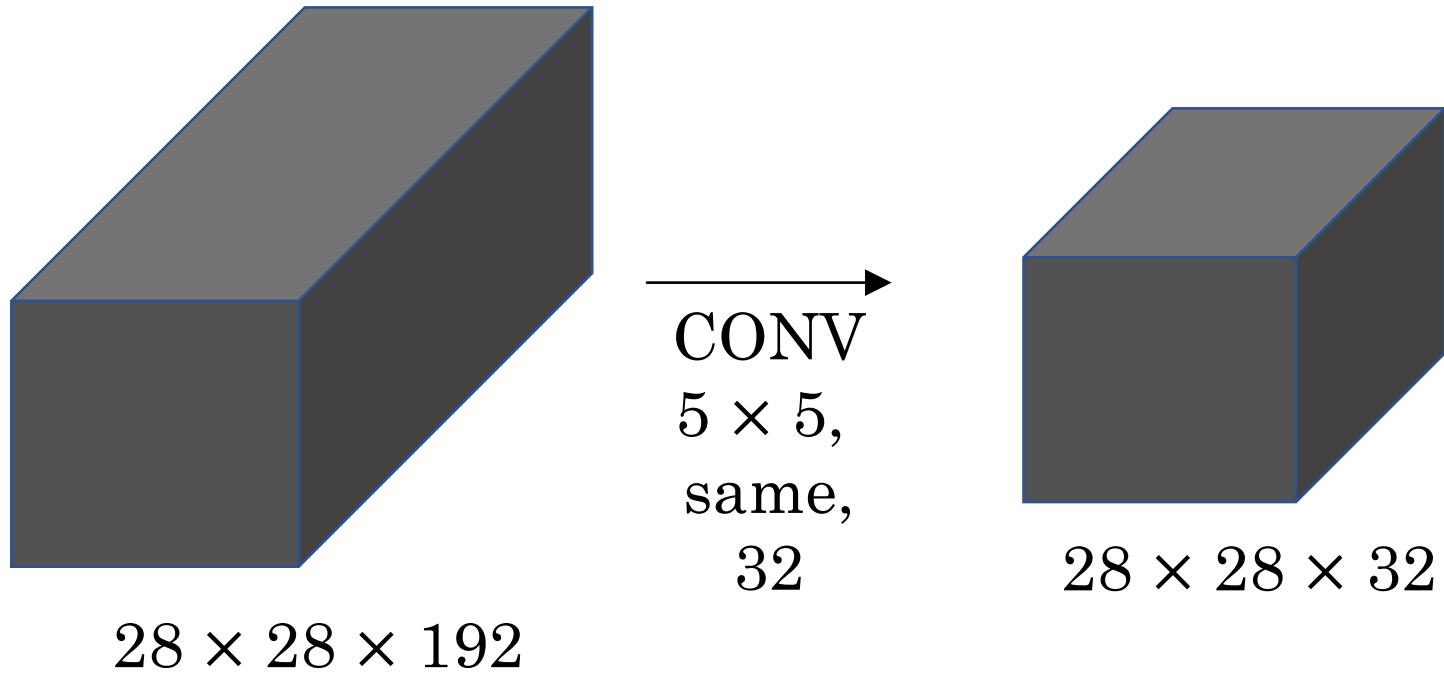
Inception Network



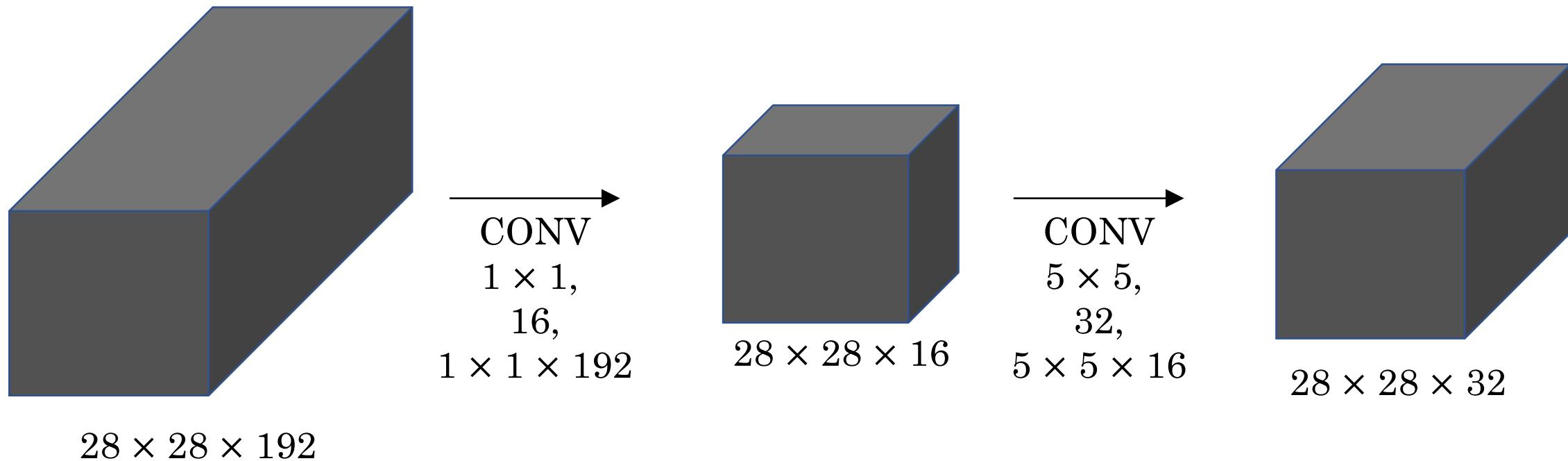
Motivation for inception network



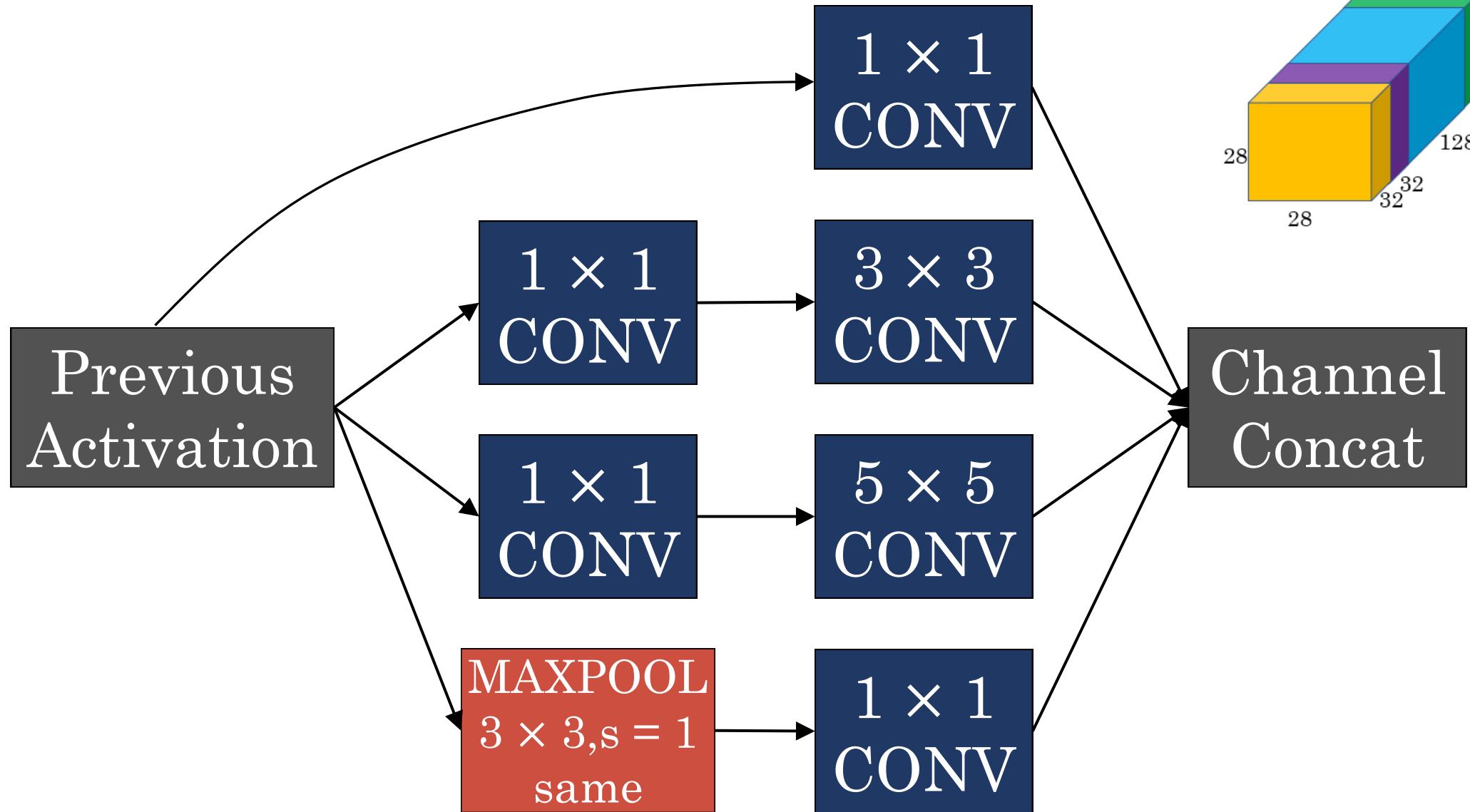
The problem of computational cost



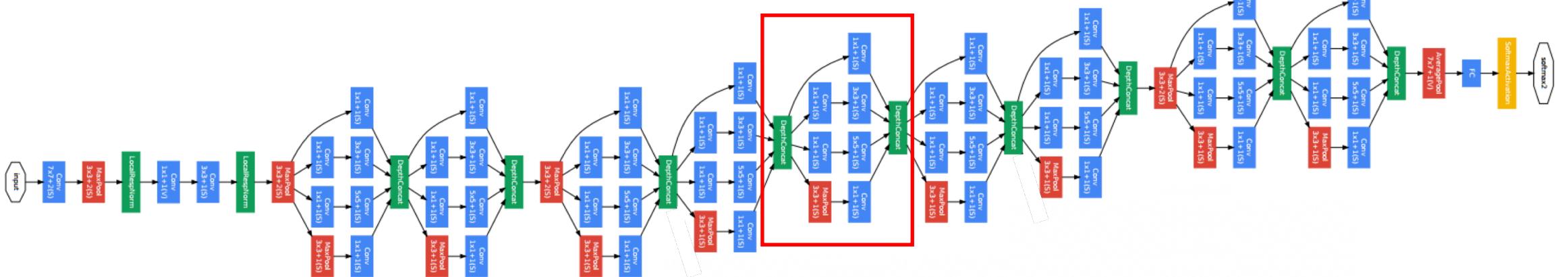
Using 1×1 convolution



Inception module



Inception network



[Szegedy et al., 2014, Going Deeper with Convolutions]

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