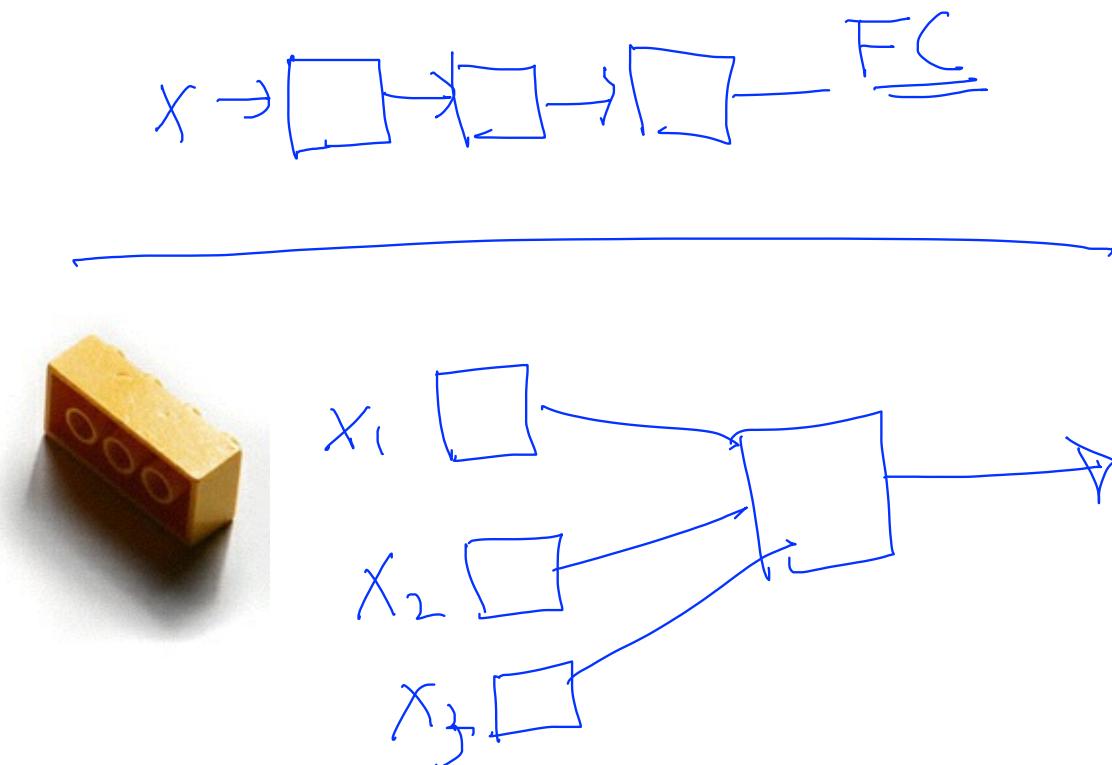
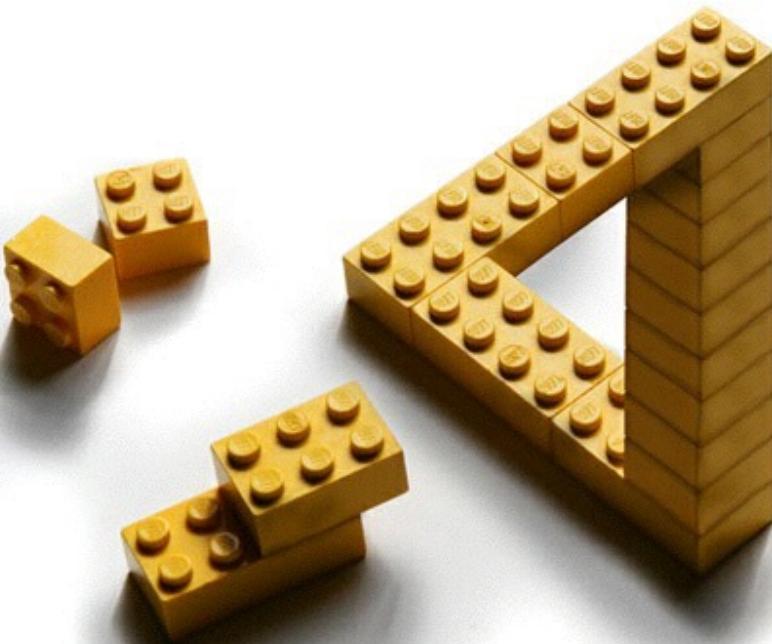


# Lecture II - I

## CNN introduction

Sung Kim <hunkim+mr@gmail.com>

# 'The only limit is your imagination'



# Convolutional Neural Networks

<http://cs231n.stanford.edu/>

# A bit of history:

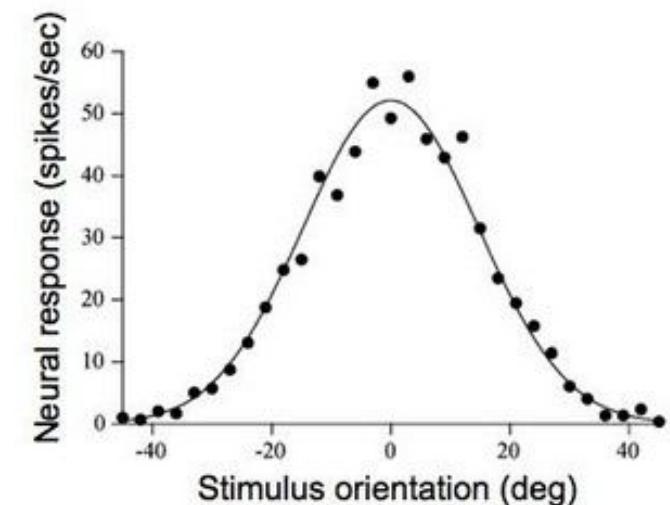
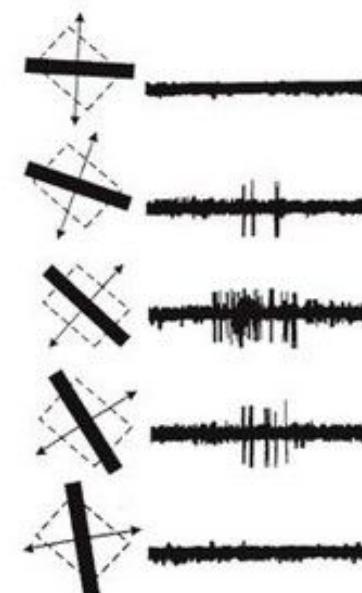
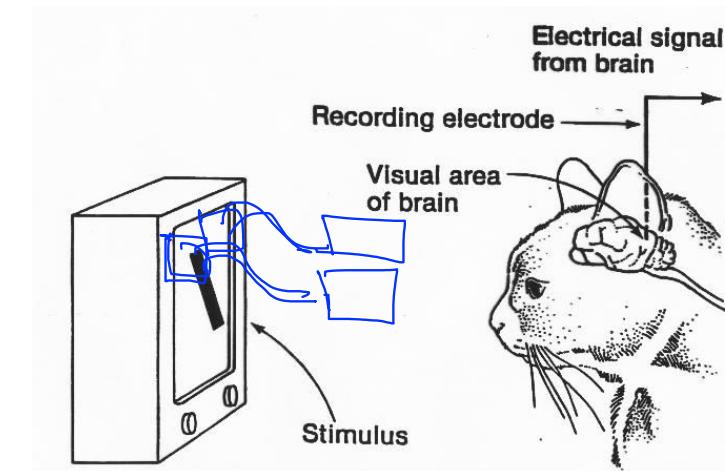
## Hubel & Wiesel, 1959

RECEPTIVE FIELDS OF SINGLE  
NEURONES IN  
THE CAT'S STRIATE CORTEX

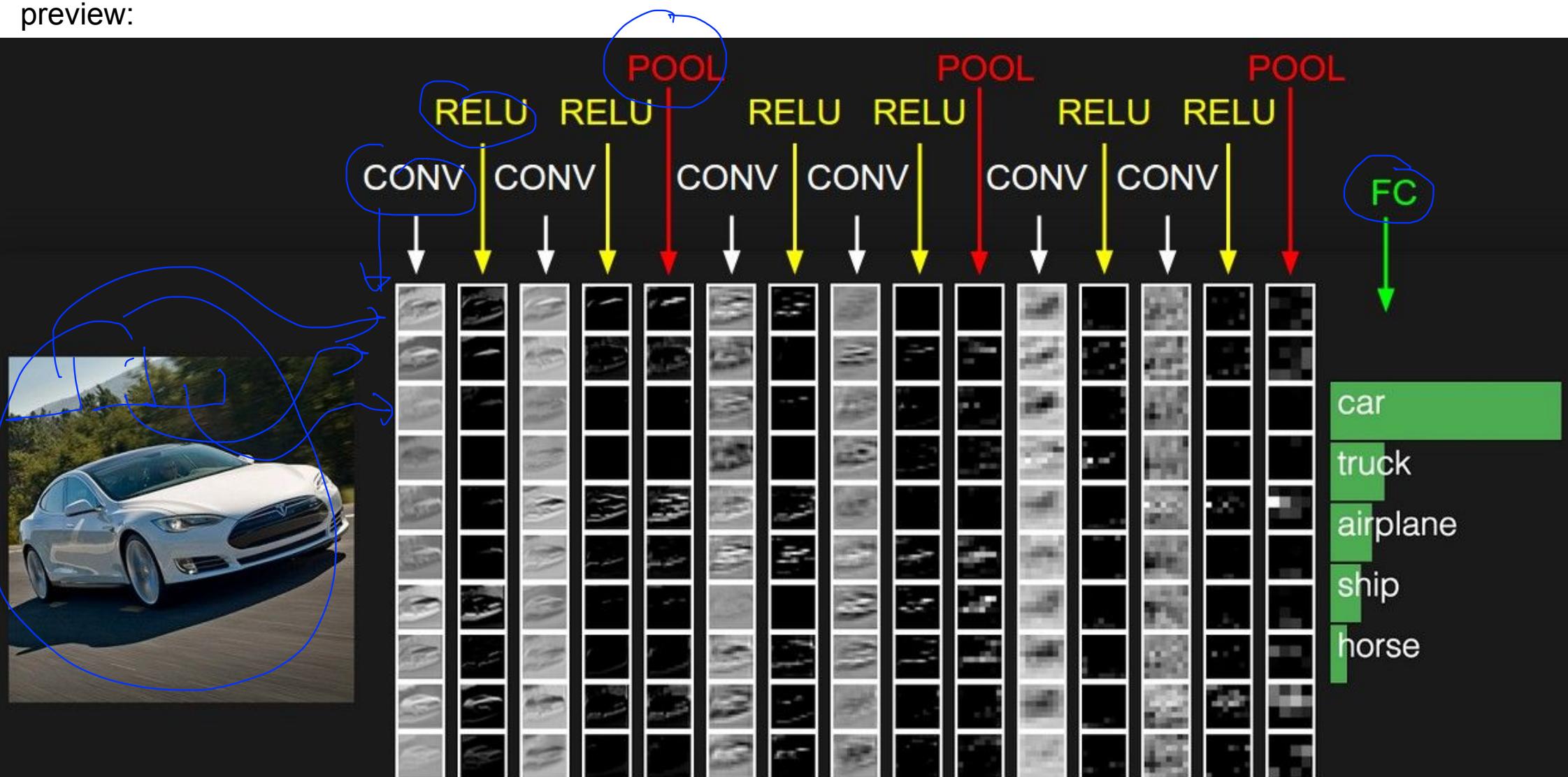
1962

RECEPTIVE FIELDS, BINOCULAR  
INTERACTION  
AND FUNCTIONAL ARCHITECTURE IN  
THE CAT'S VISUAL CORTEX

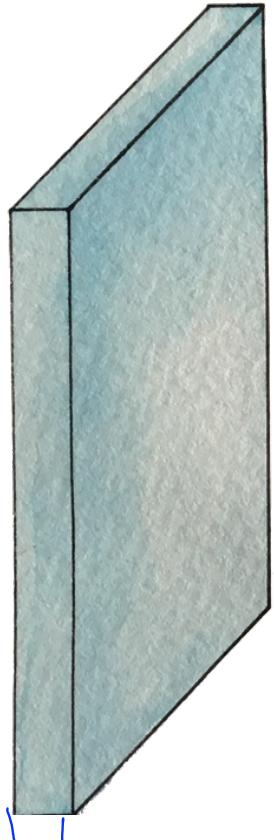
1968...



preview:

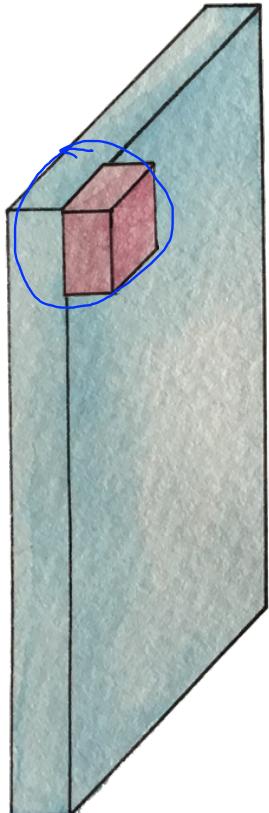


Start with an image (width x height x depth)



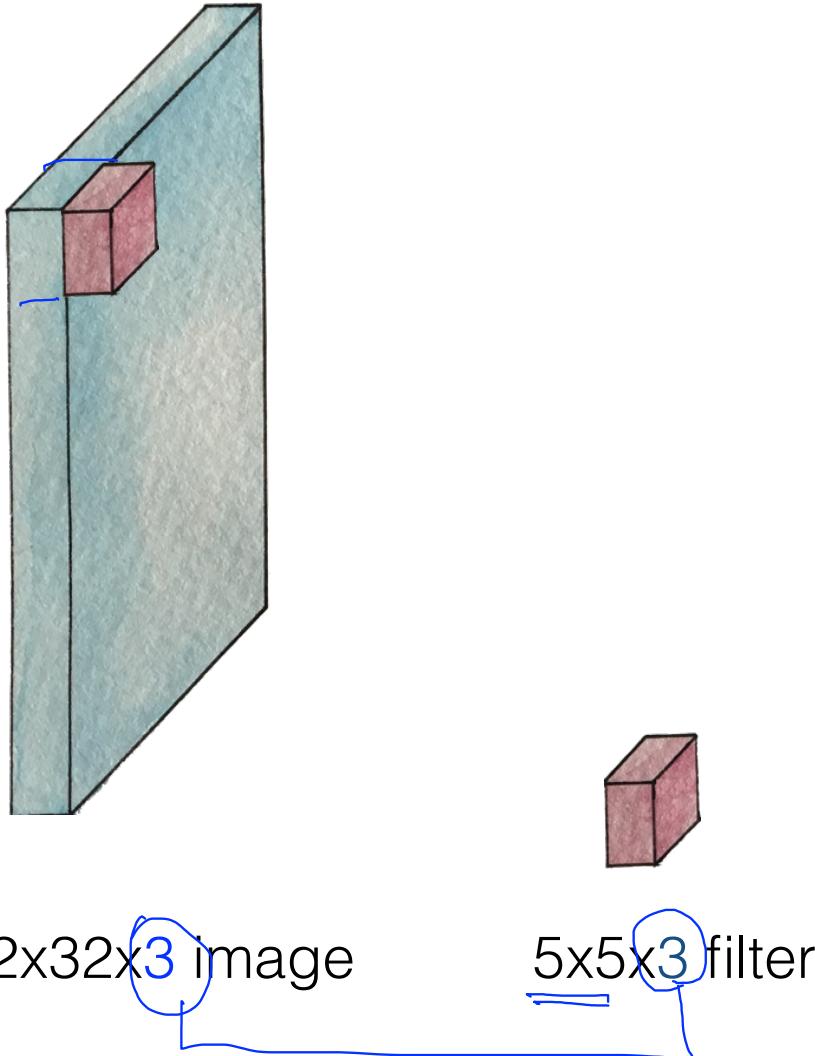
32x32x3 image

# Let's focus on a small area only

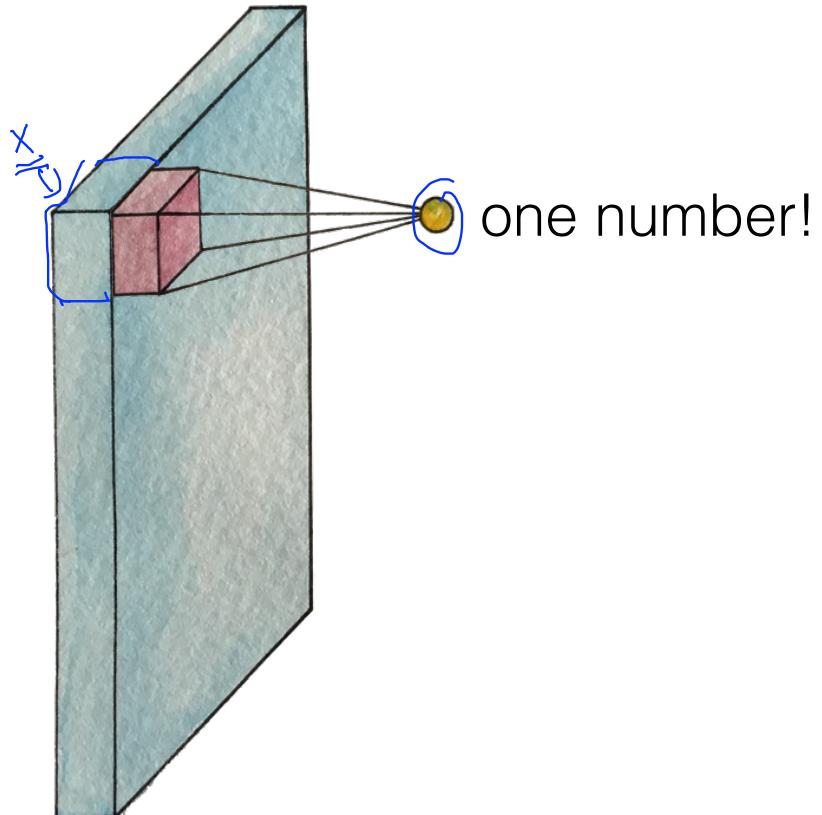


32x32x3 image

# Let's focus on a small area only ( $5 \times 5 \times 3$ )

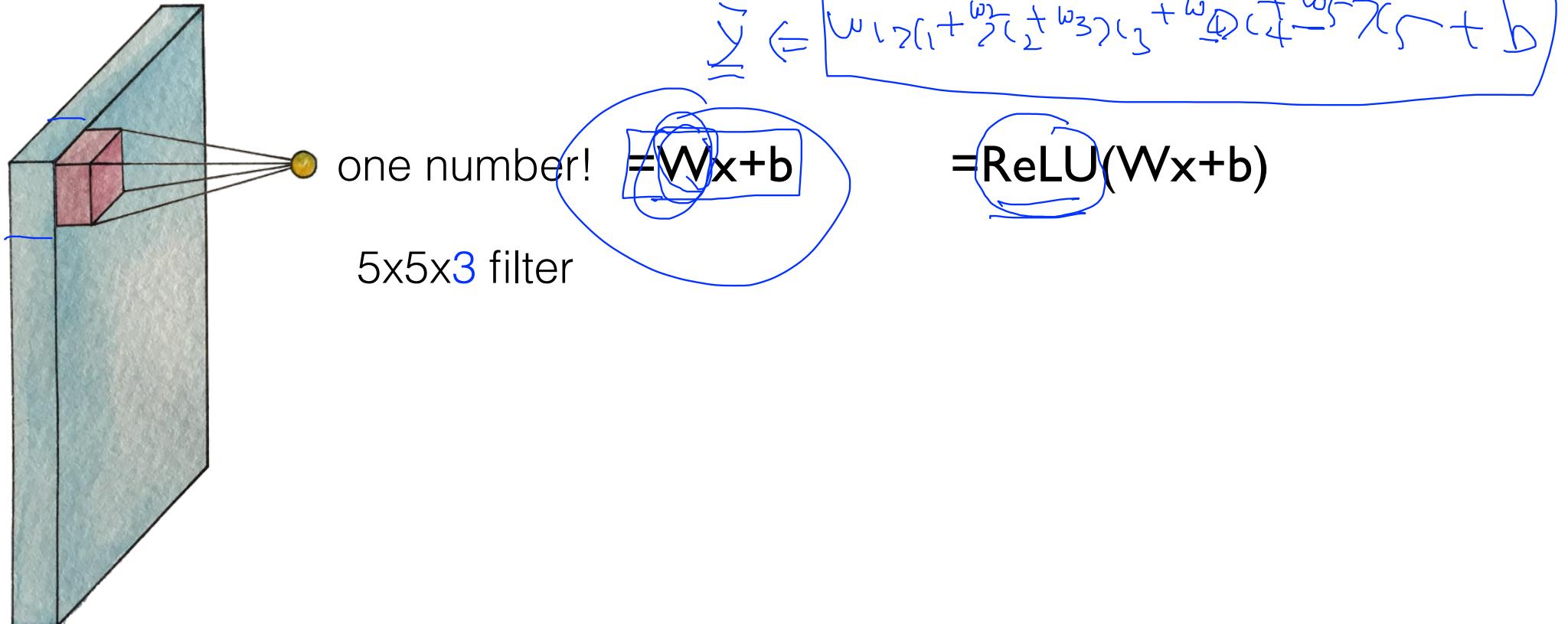


# Get one number using the filter



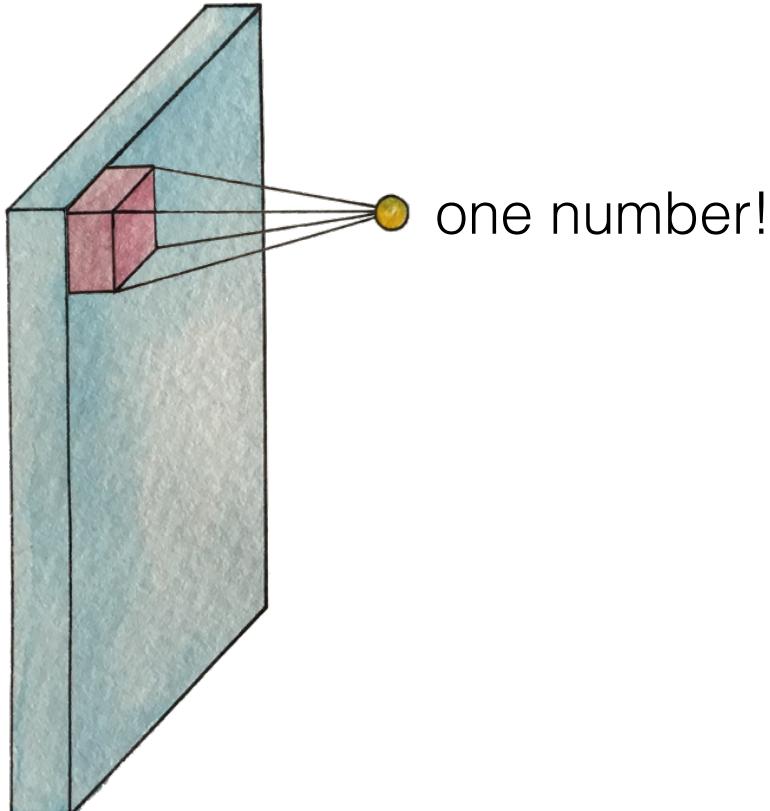
32x32x $\color{blue}{3}$  image

# Get one number using the filter



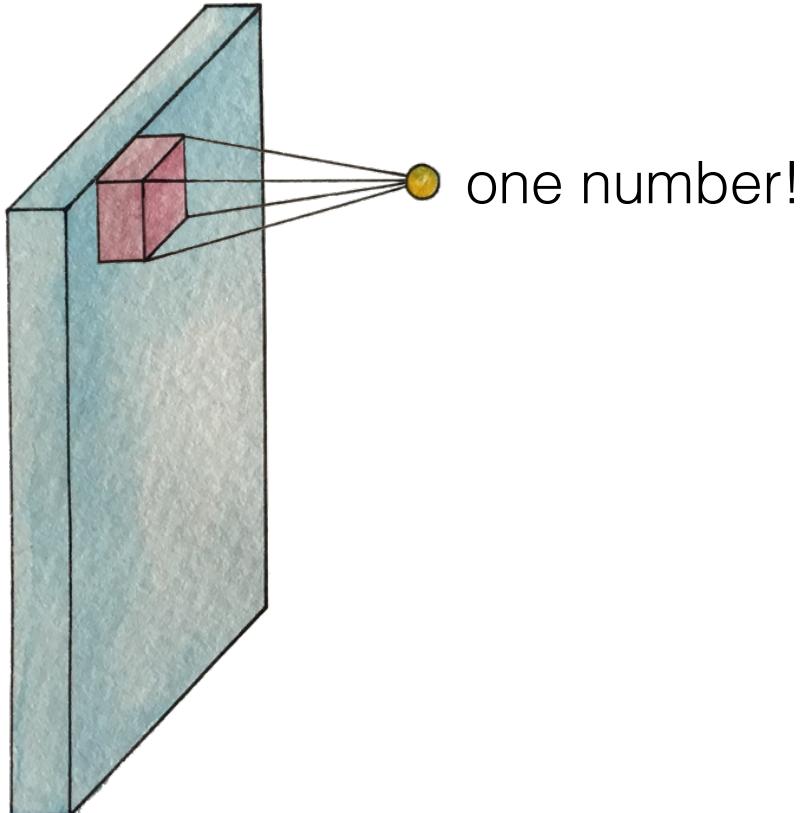
$32 \times 32 \times 3$  image

# Let's look at other areas with the same filter (w)



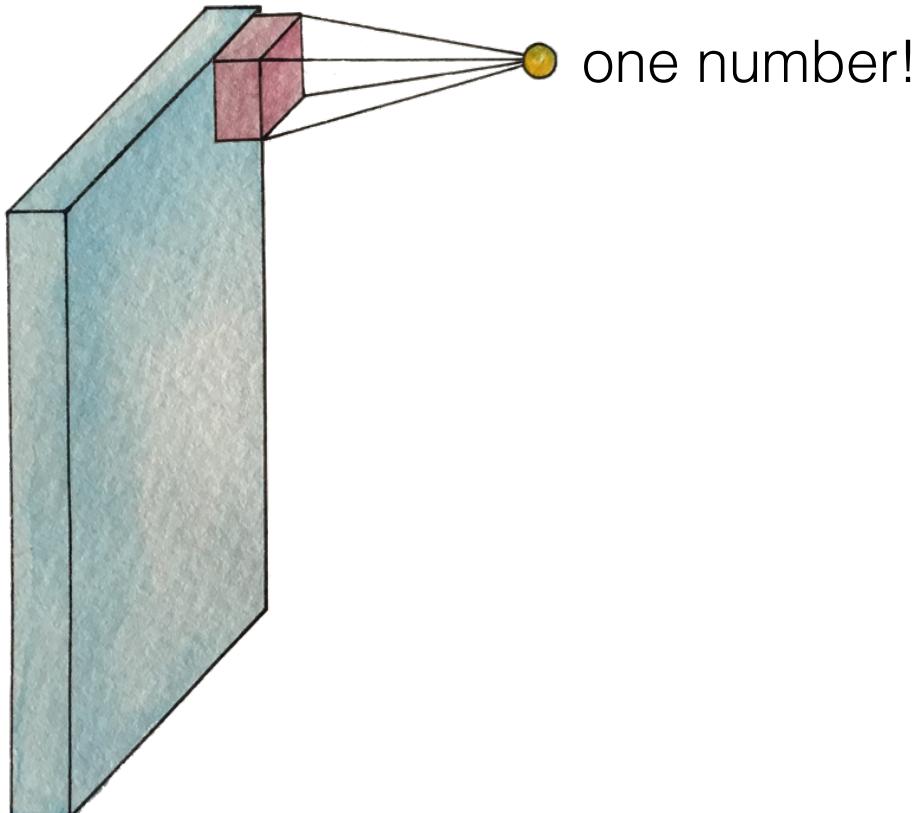
32x32x3 image

# Let's look at other areas with the same filter ( $w$ )



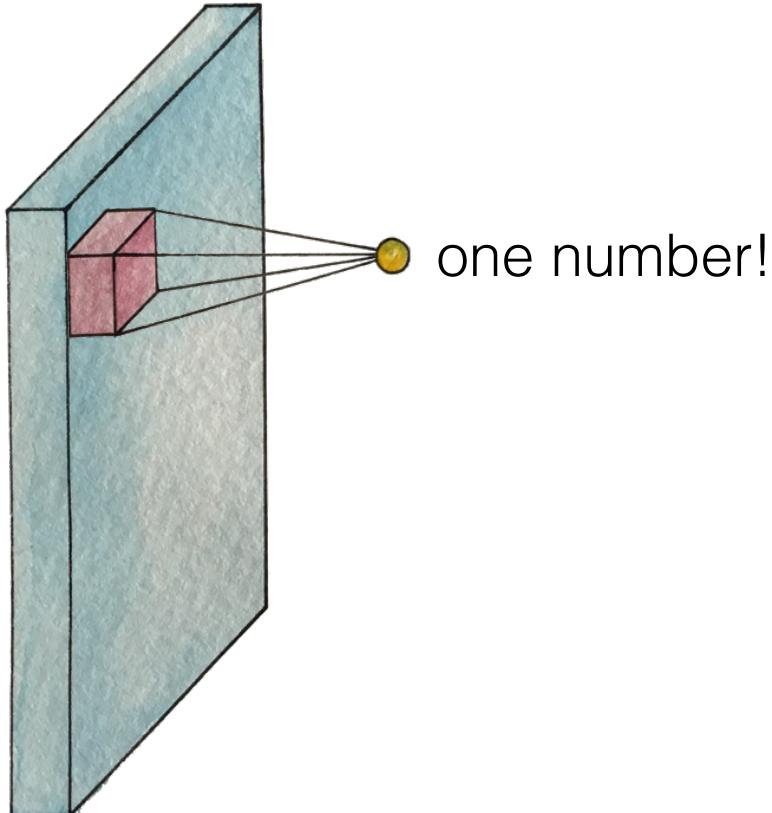
32x32x3 image

# Let's look at other areas with the same filter ( $w$ )



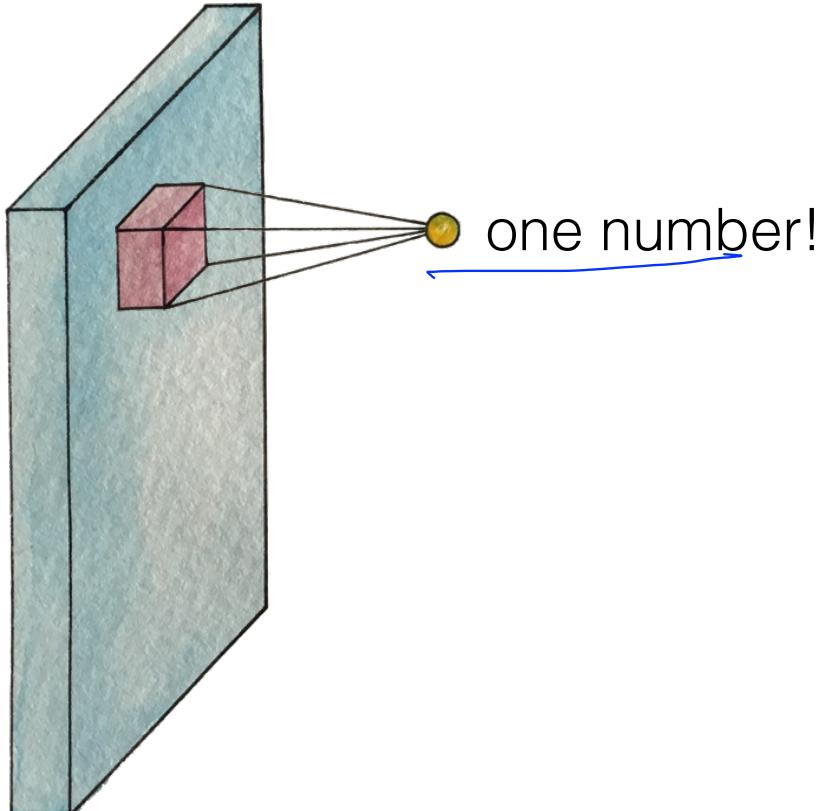
32x32x3 image

# Let's look at other areas with the same filter ( $w$ )



32x32x3 image

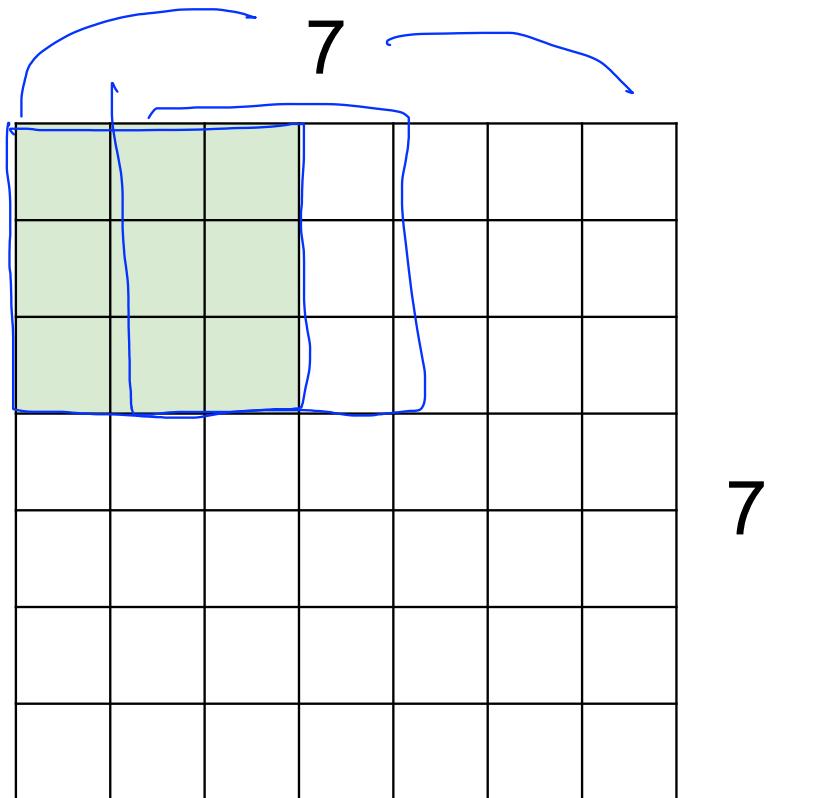
# Let's look at other areas with the same filter (w)



How many numbers  
can we get?

32x32x3 image

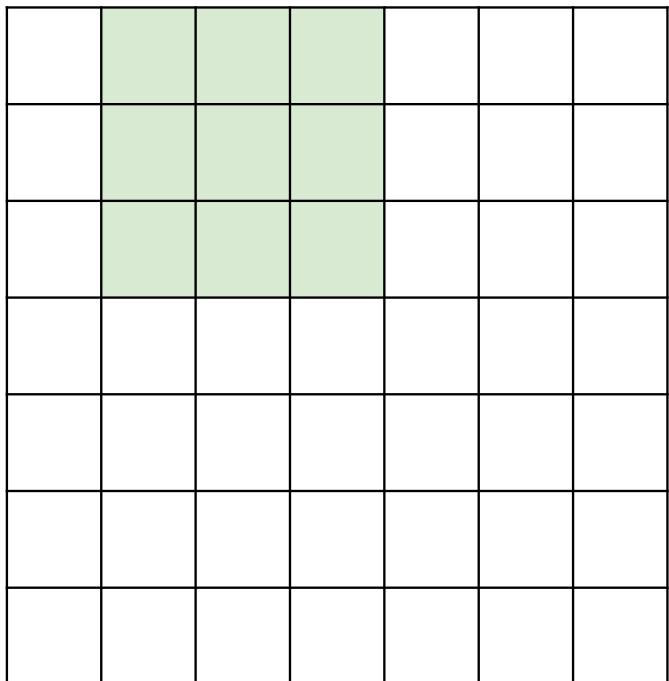
## A closer look at spatial dimensions:



7x7 input (spatially)  
assume 3x3 filter

## A closer look at spatial dimensions:

7

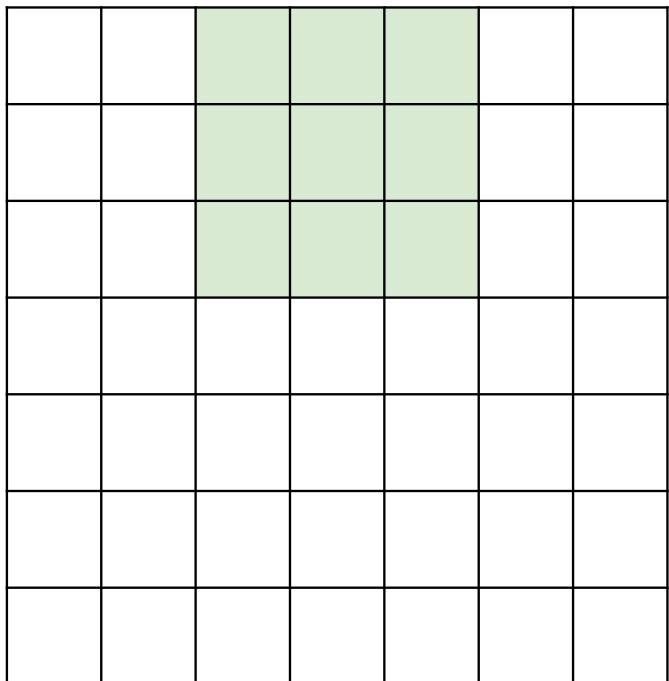


7x7 input (spatially)  
assume 3x3 filter

7

## A closer look at spatial dimensions:

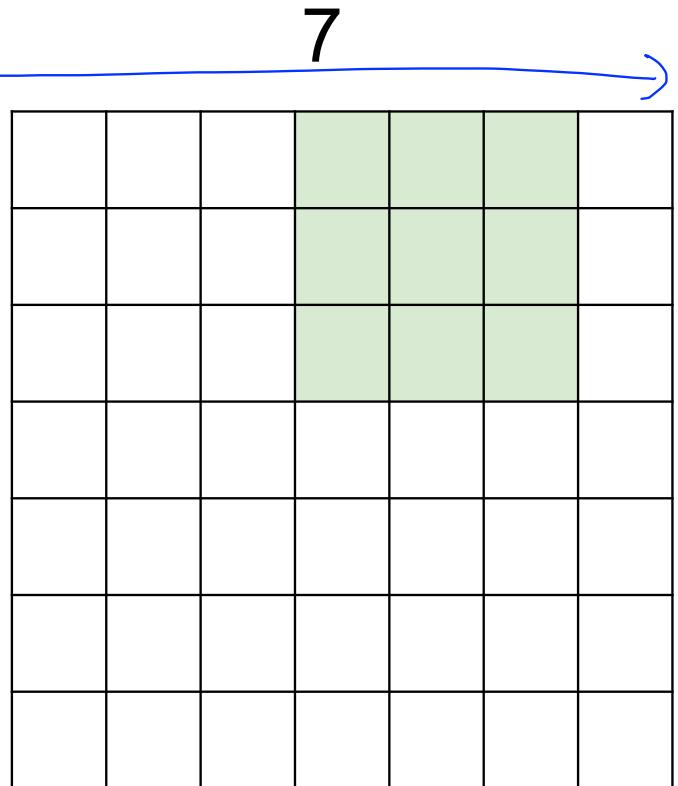
7



7x7 input (spatially)  
assume 3x3 filter

7

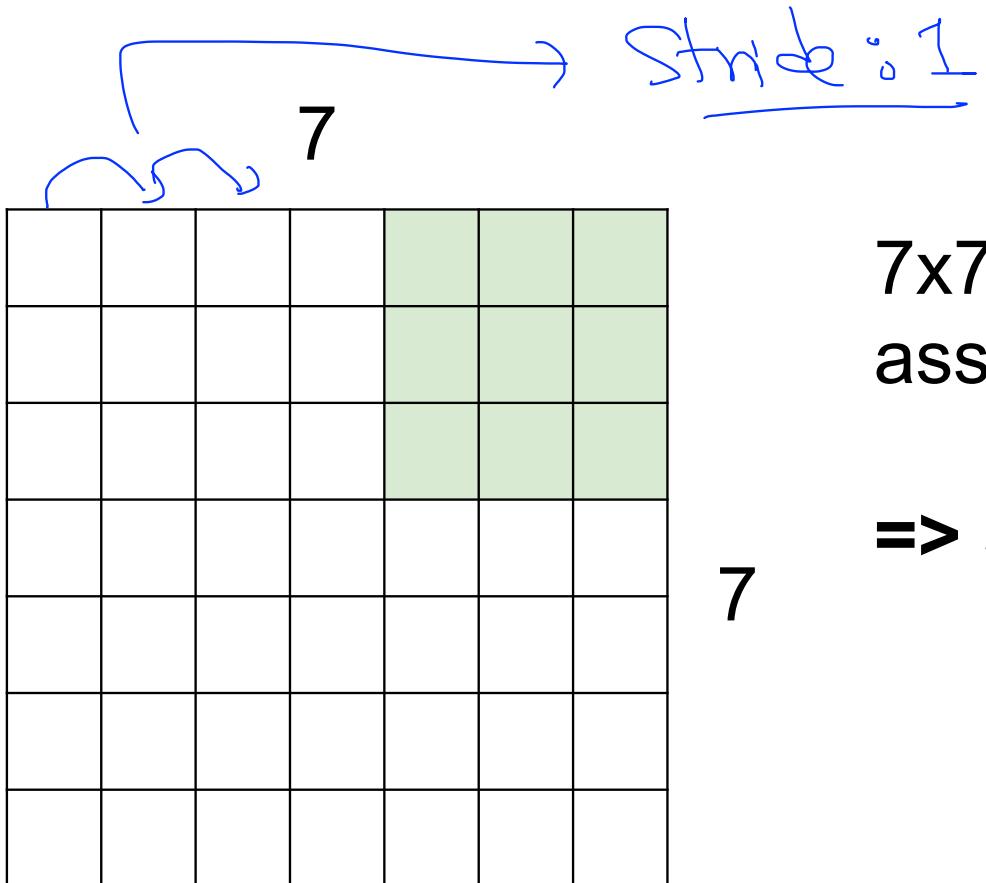
## A closer look at spatial dimensions:



7x7 input (spatially)  
assume 3x3 filter

5x5

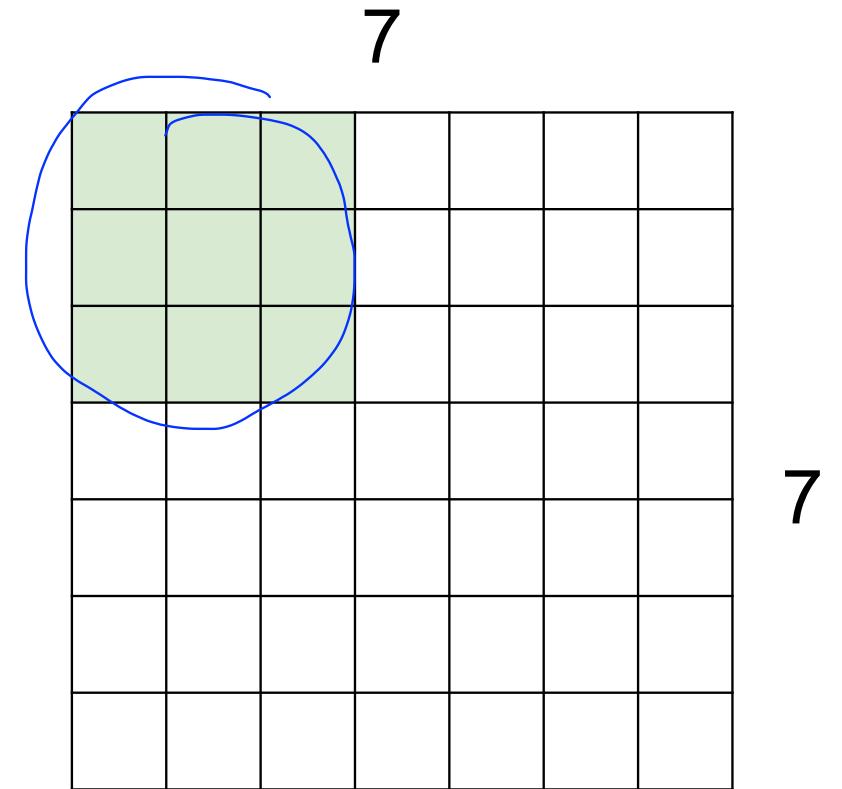
## A closer look at spatial dimensions:



7x7 input (spatially)  
assume 3x3 filter

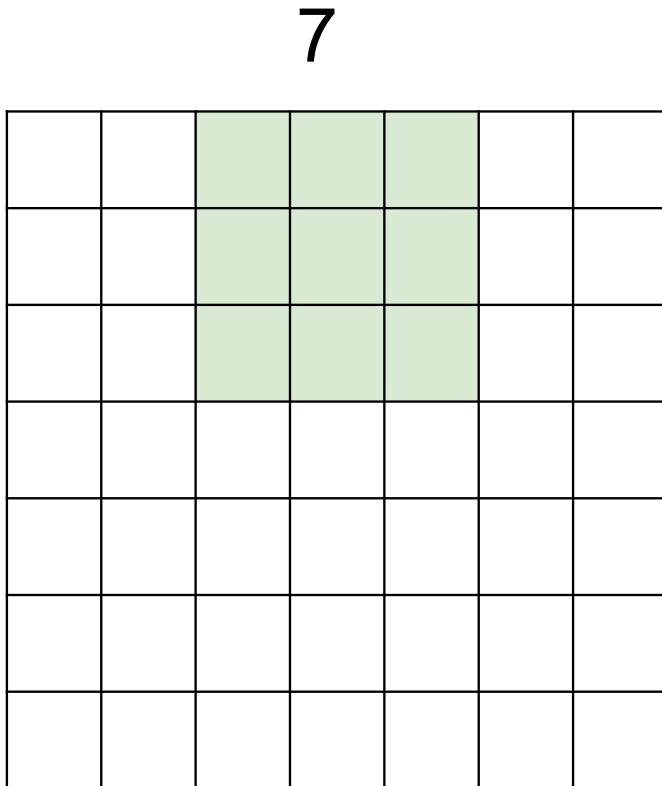
**=> 5x5 output**

## A closer look at spatial dimensions:



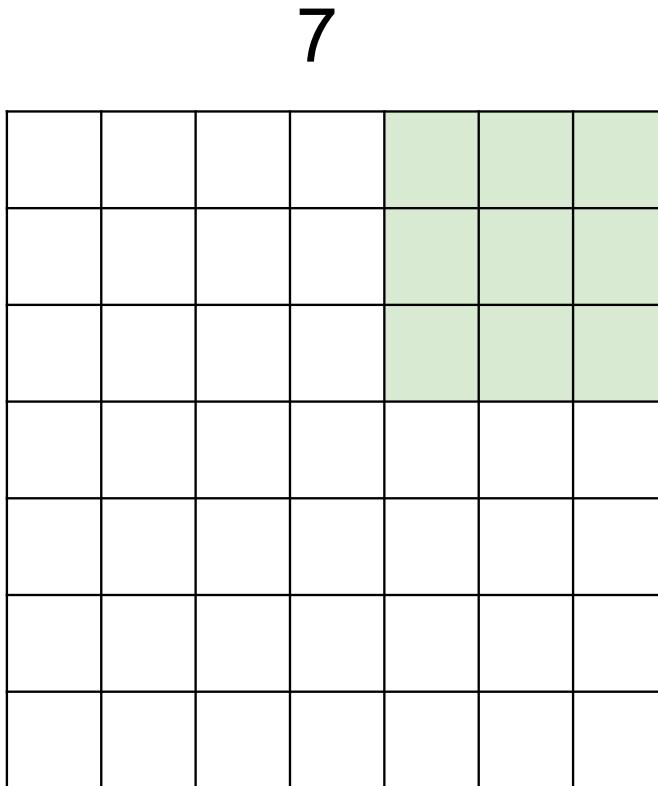
7x7 input (spatially)  
assume 3x3 filter  
applied **with stride 2**

## A closer look at spatial dimensions:

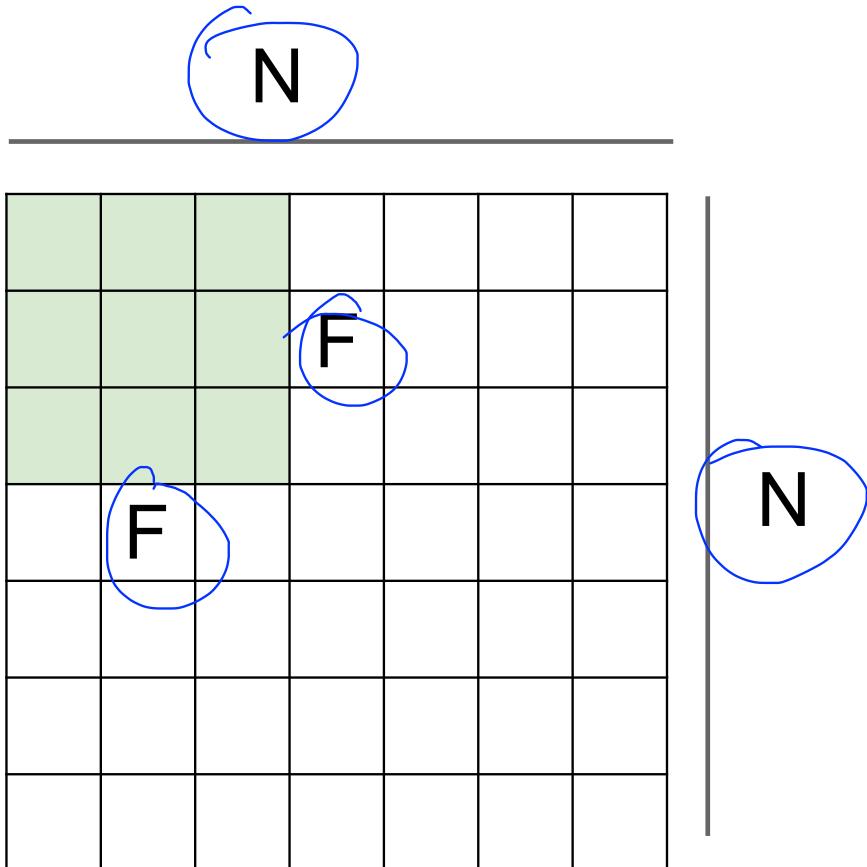


7x7 input (spatially)  
assume 3x3 filter  
applied **with stride 2**

## A closer look at spatial dimensions:



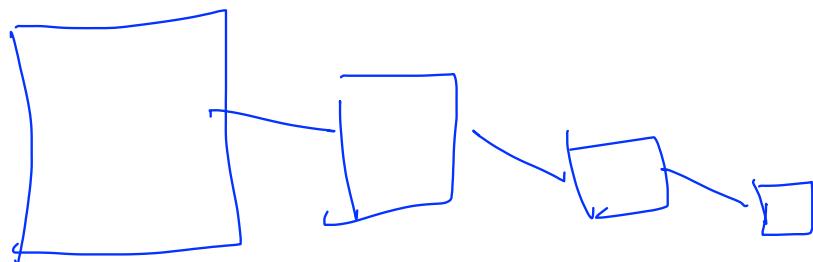
7x7 input (spatially)  
assume 3x3 filter  
applied **with stride 2**  
=> 3x3 output!



Output size:  
 $(N - F) / \text{stride} + 1$

e.g.  $N = 7, F = 3$ :

$$\begin{aligned} \text{stride } 1 &\Rightarrow (7 - 3)/1 + 1 = 5 \\ \text{stride } 2 &\Rightarrow (7 - 3)/2 + 1 = 3 \\ \text{stride } 3 &\Rightarrow (7 - 3)/3 + 1 = 2.33 : \end{aligned}$$



# In practice: Common to zero pad the border

0	0	0	0	0	0		
0							
0							
0							
0							
0							

e.g. input  $7 \times 7$

$\underline{3 \times 3}$  filter, applied with **stride 1**

**pad with 1 pixel border** => what is the output?

$$7 \times 7 \xrightarrow[\text{pad}]{} 9 \times 9$$

(recall:)

$$(N - F) / \text{stride} + 1$$

# In practice: Common to zero pad the border

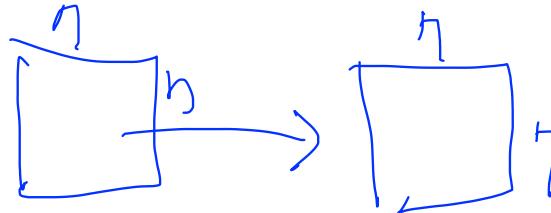
0	0	0	0	0	0			
0								
0								
0								
0								

e.g. input  $7 \times 7$

$3 \times 3$  filter, applied with **stride 1**

**pad with 1 pixel border** => what is the output?

**$7 \times 7$  output!**



# In practice: Common to zero pad the border

0	0	0	0	0	0			
0								
0								
0								
0								

e.g. input  $7 \times 7$

**3x3 filter, applied with stride 1**

**pad with 1 pixel border => what is the output?**

**$7 \times 7$  output!**

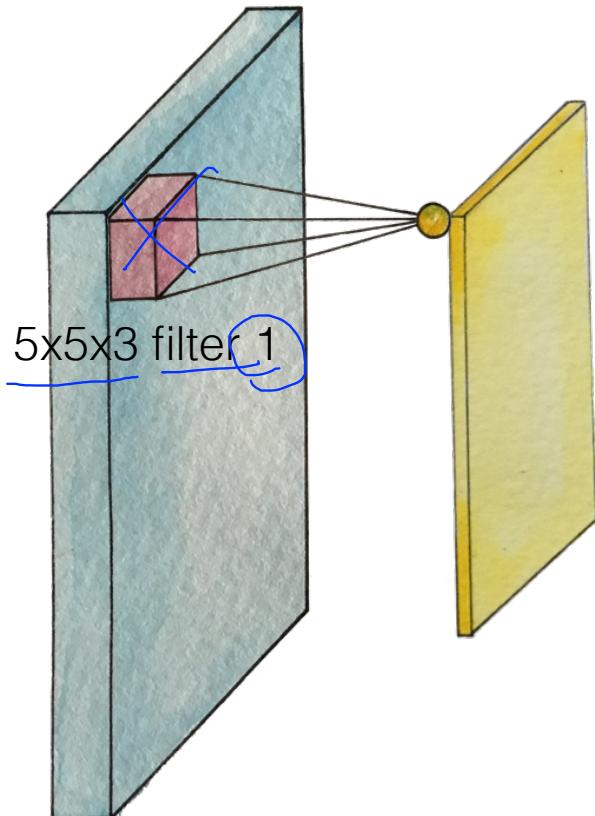
in general, common to see CONV layers with stride 1, filters of size  $F \times F$ , and zero-padding with  $(F-1)/2$ . (will preserve size spatially)

e.g.  $F = 3 \Rightarrow$  zero pad with 1

$F = 5 \Rightarrow$  zero pad with 2

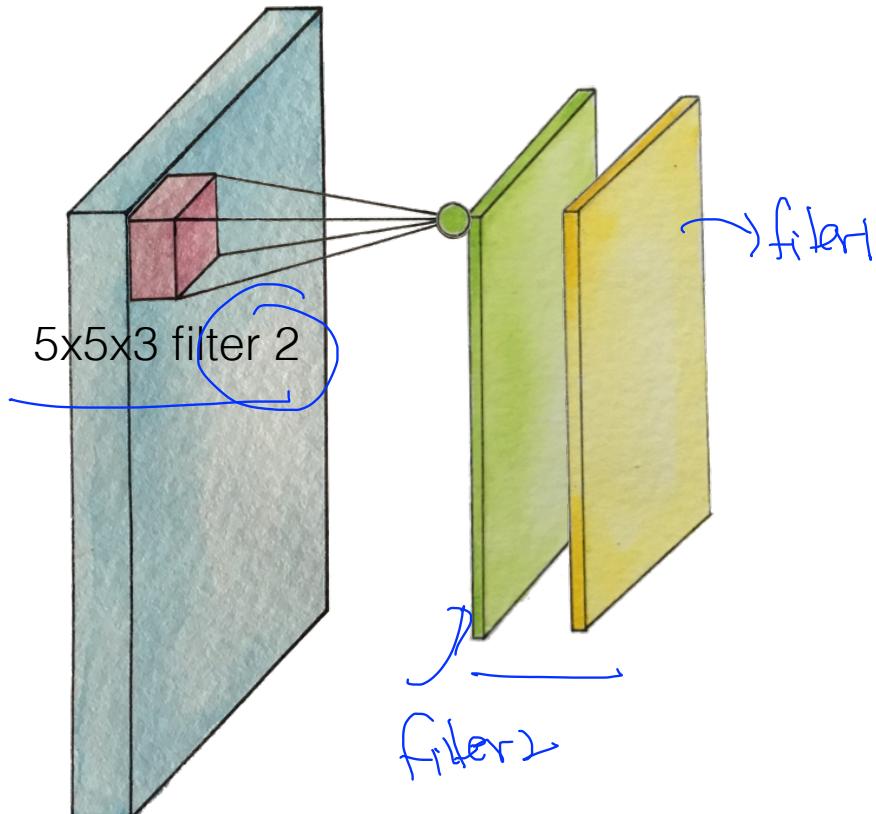
$F = 7 \Rightarrow$  zero pad with 3

# Swiping the entire image



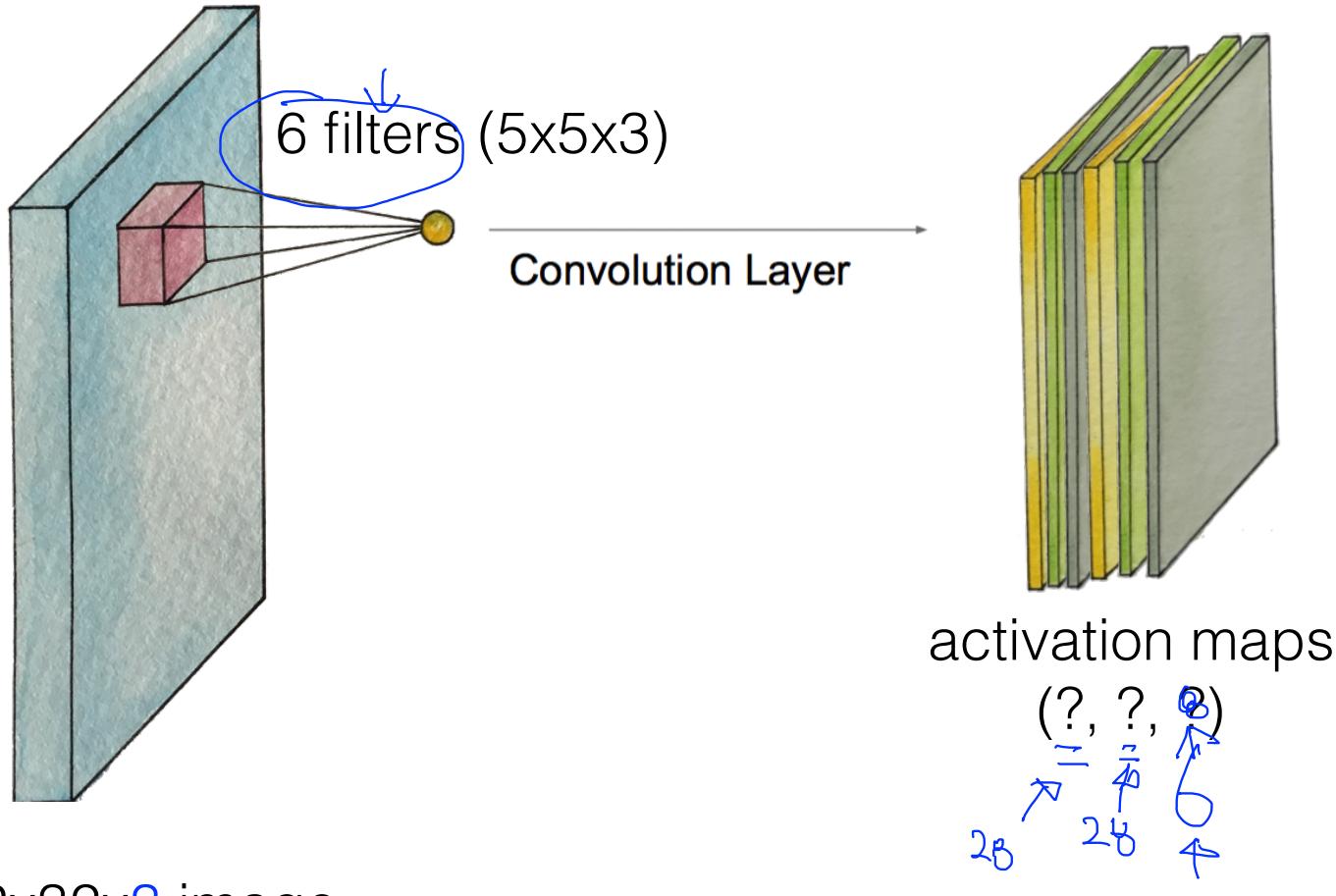
$32 \times 32 \times 3$  image

# Swiping the entire image

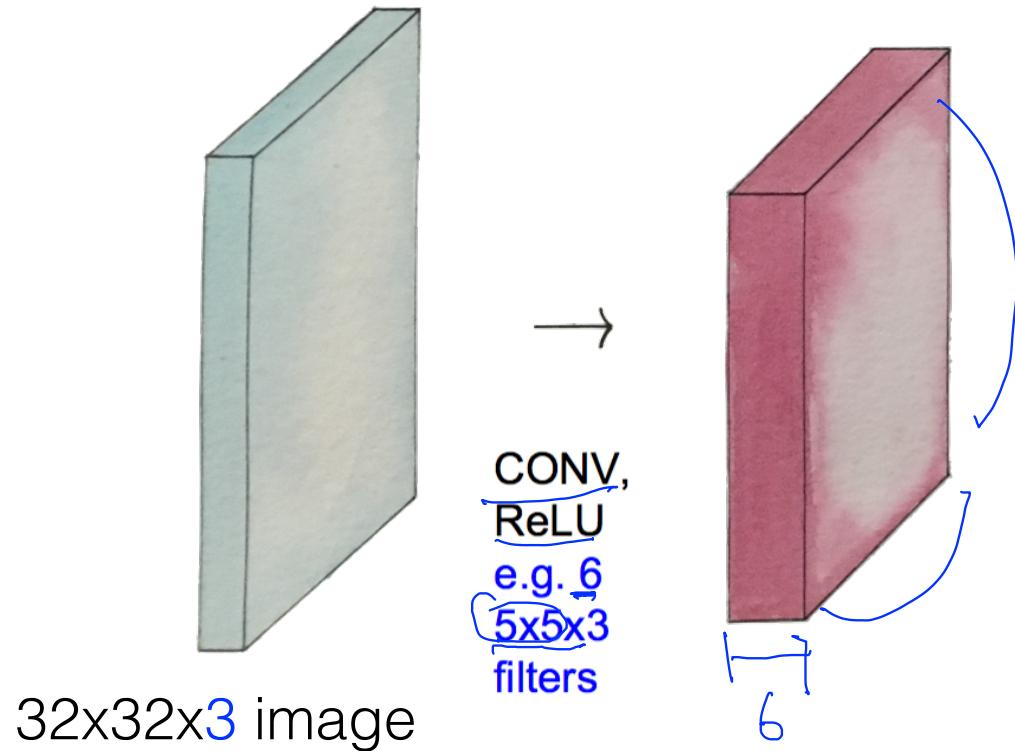


$32 \times 32 \times 3$  image

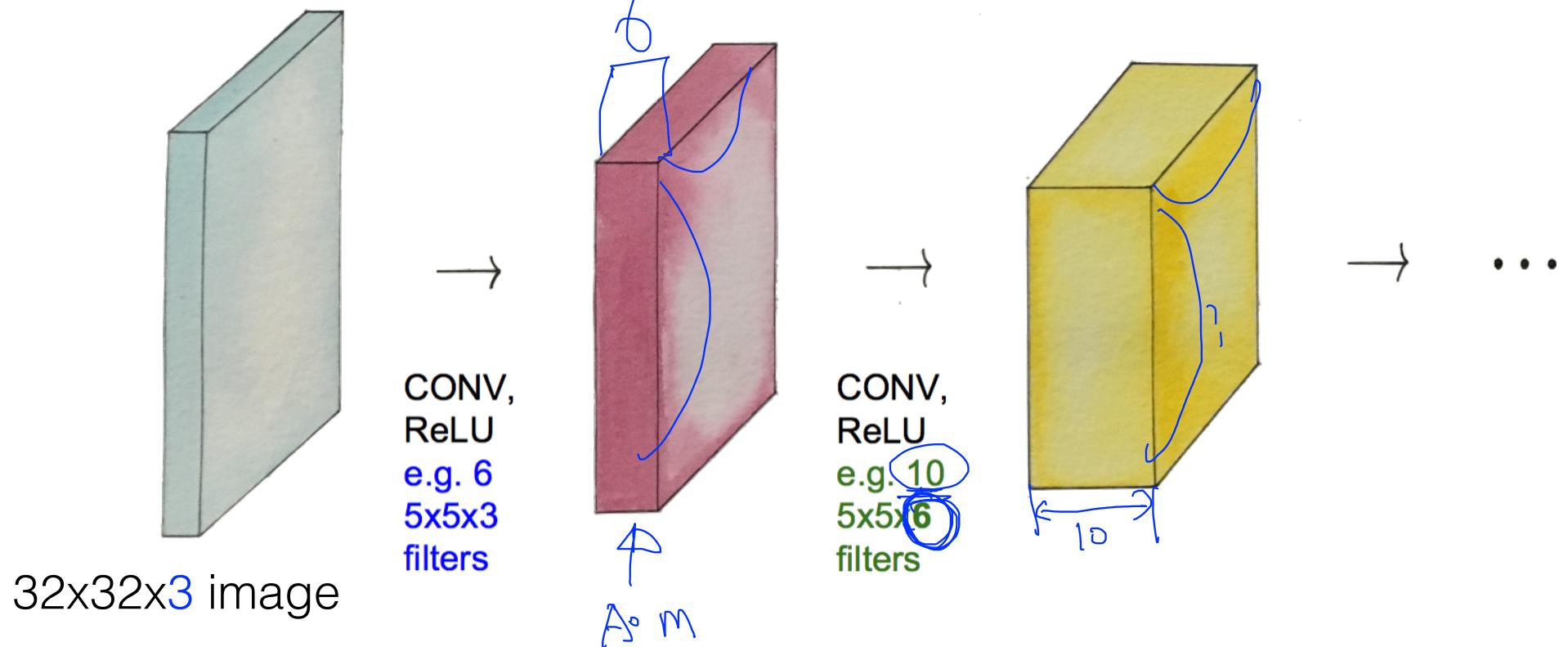
# Swiping the entire image



# Convolution layers

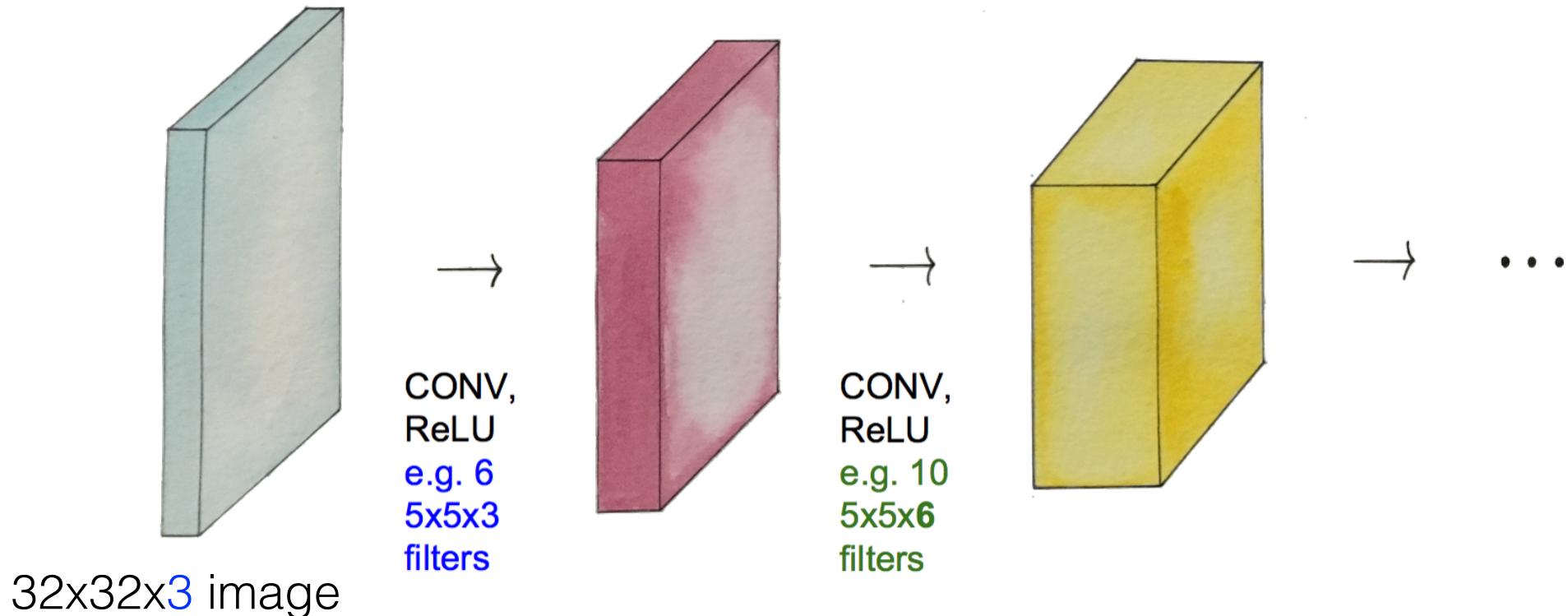


# Convolution layers



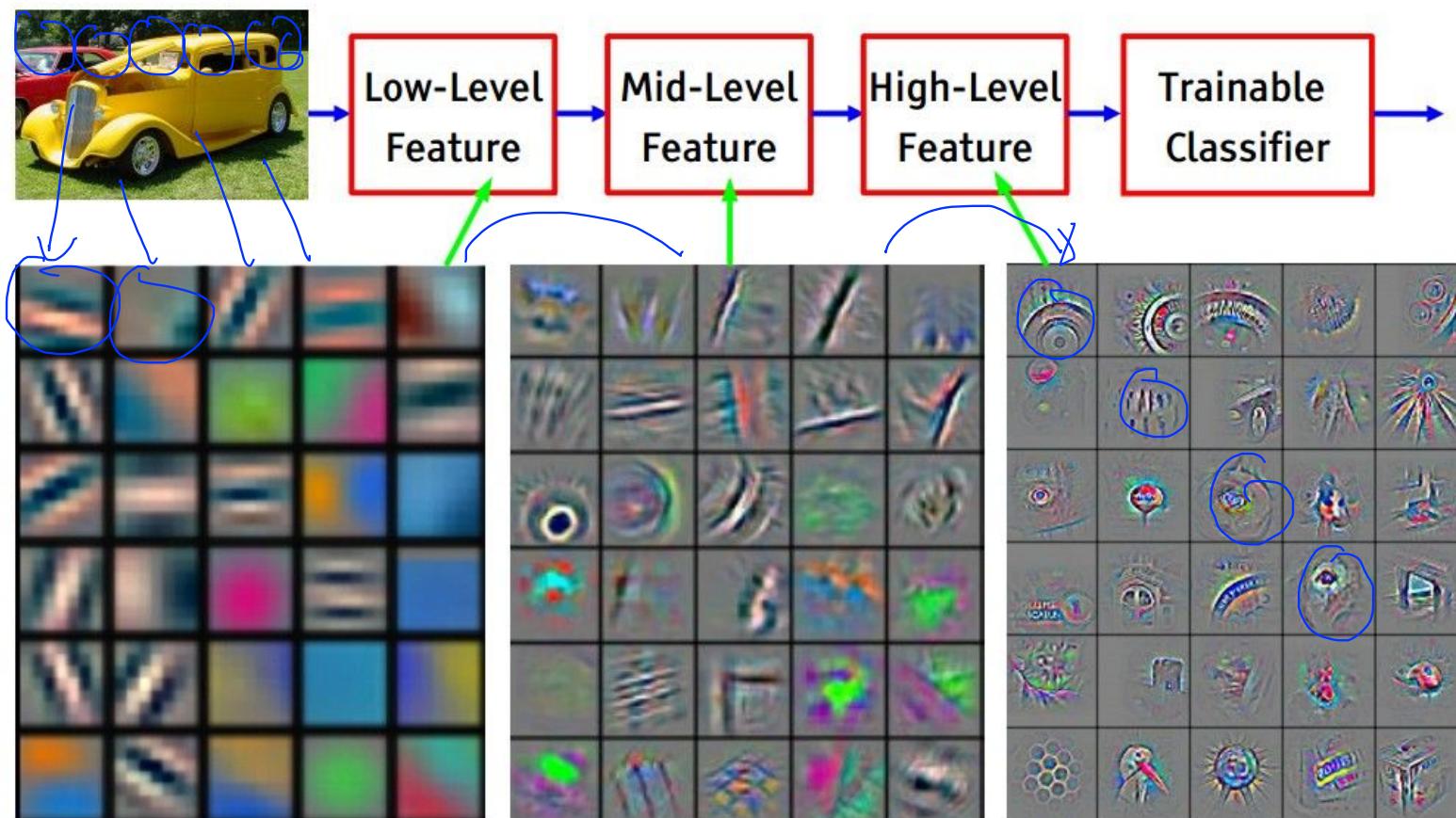
# Convolution layers

How many weight variables? How to set them?



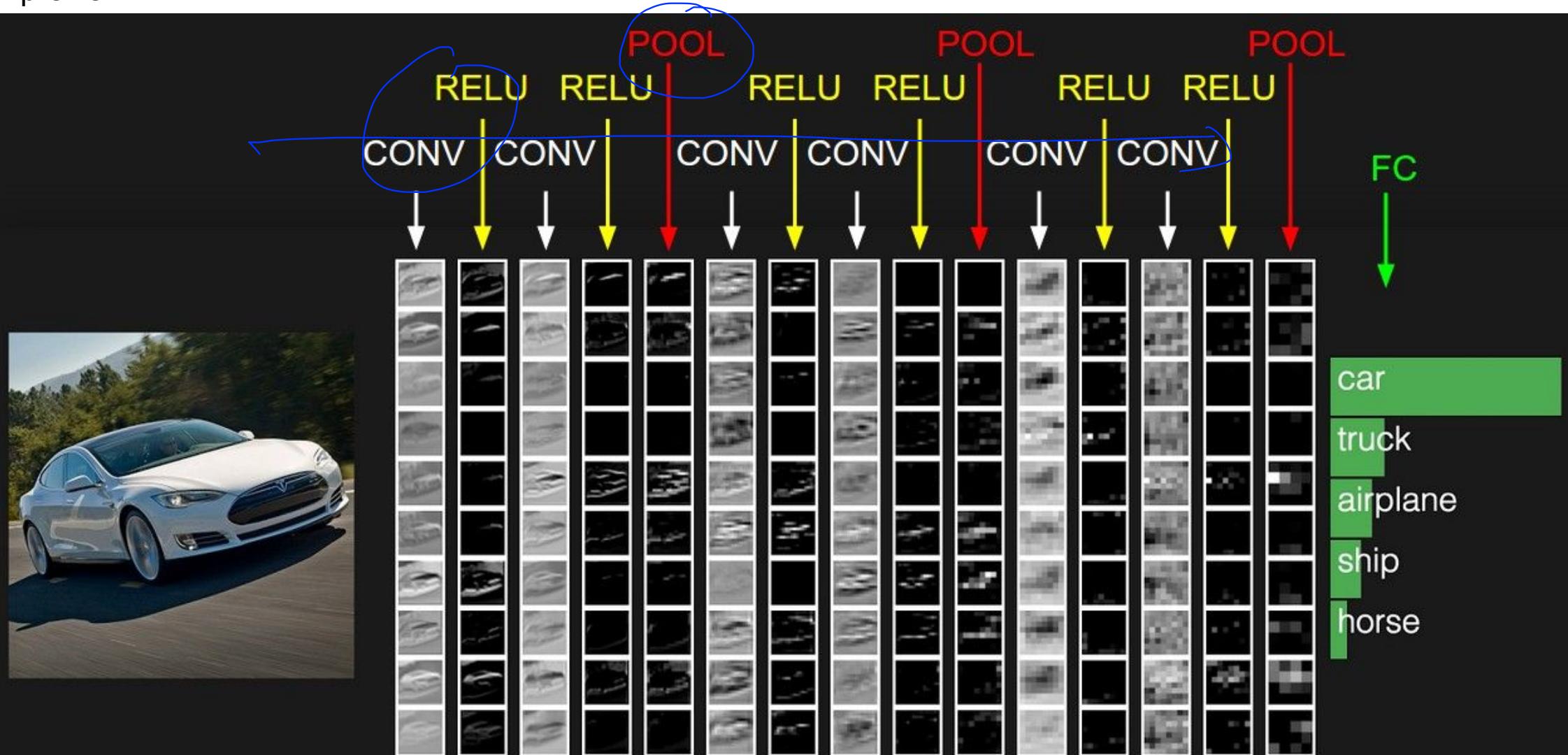
## Preview

[From recent Yann LeCun slides]



Feature visualization of convolutional net trained on ImageNet from [Zeiler & Fergus 2013]

preview:

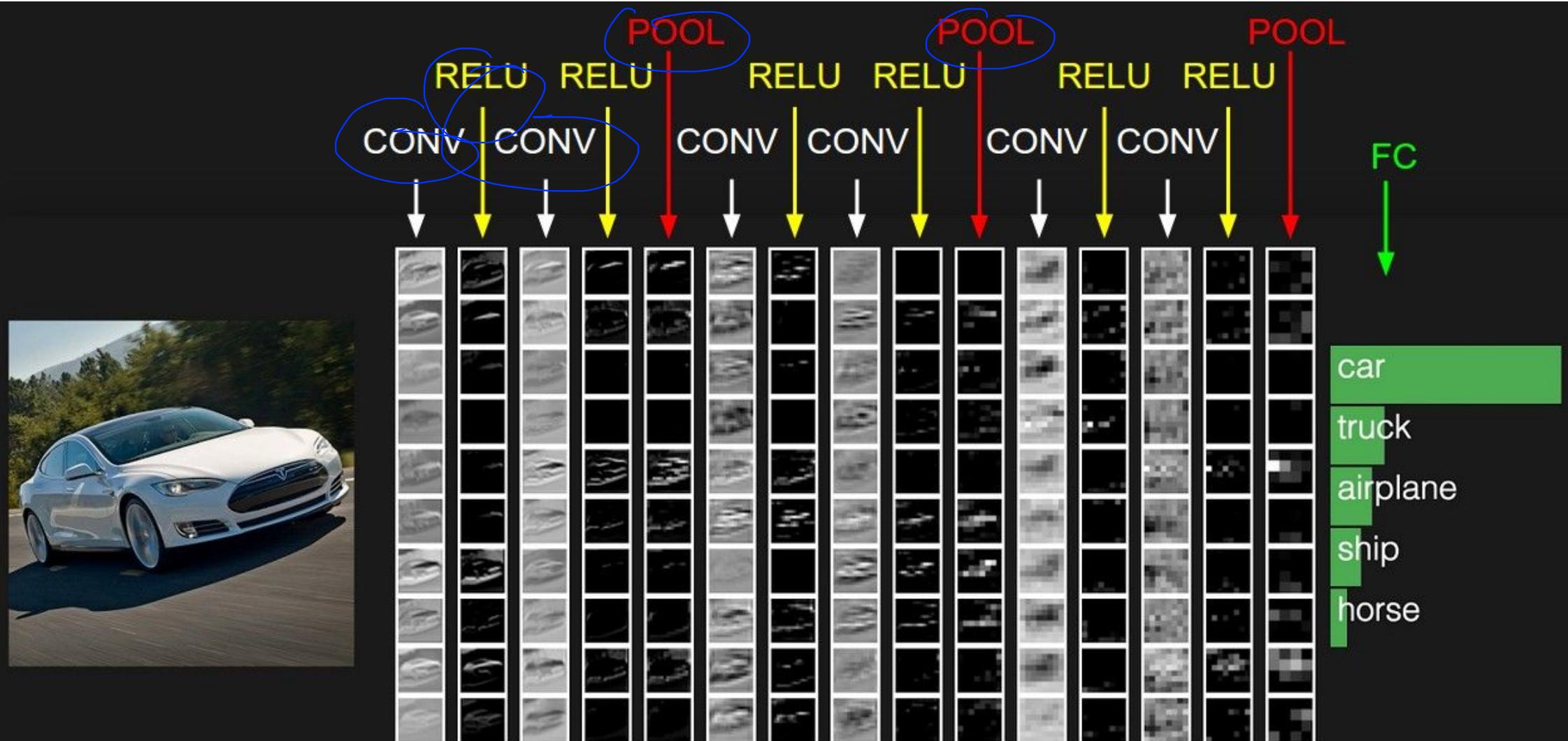


# Lecture II-2

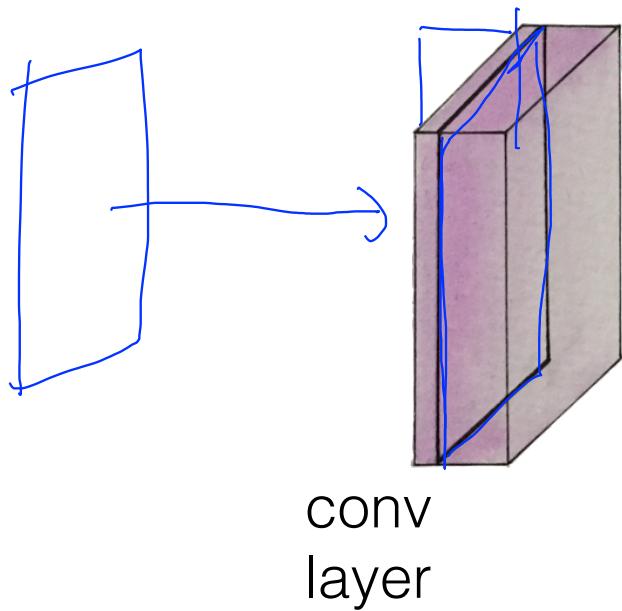
CNN introduction: Max pooling and others

Sung Kim <hunkim+mr@gmail.com>

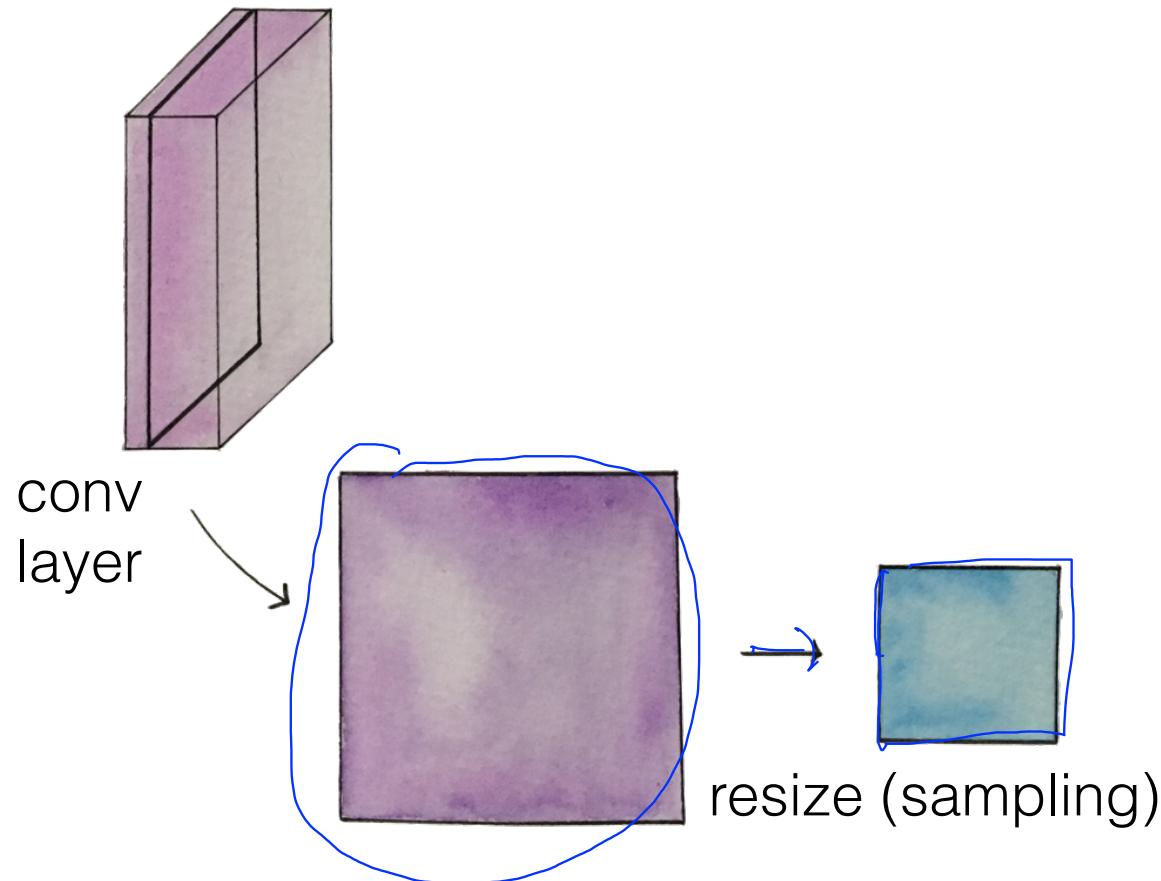
preview:



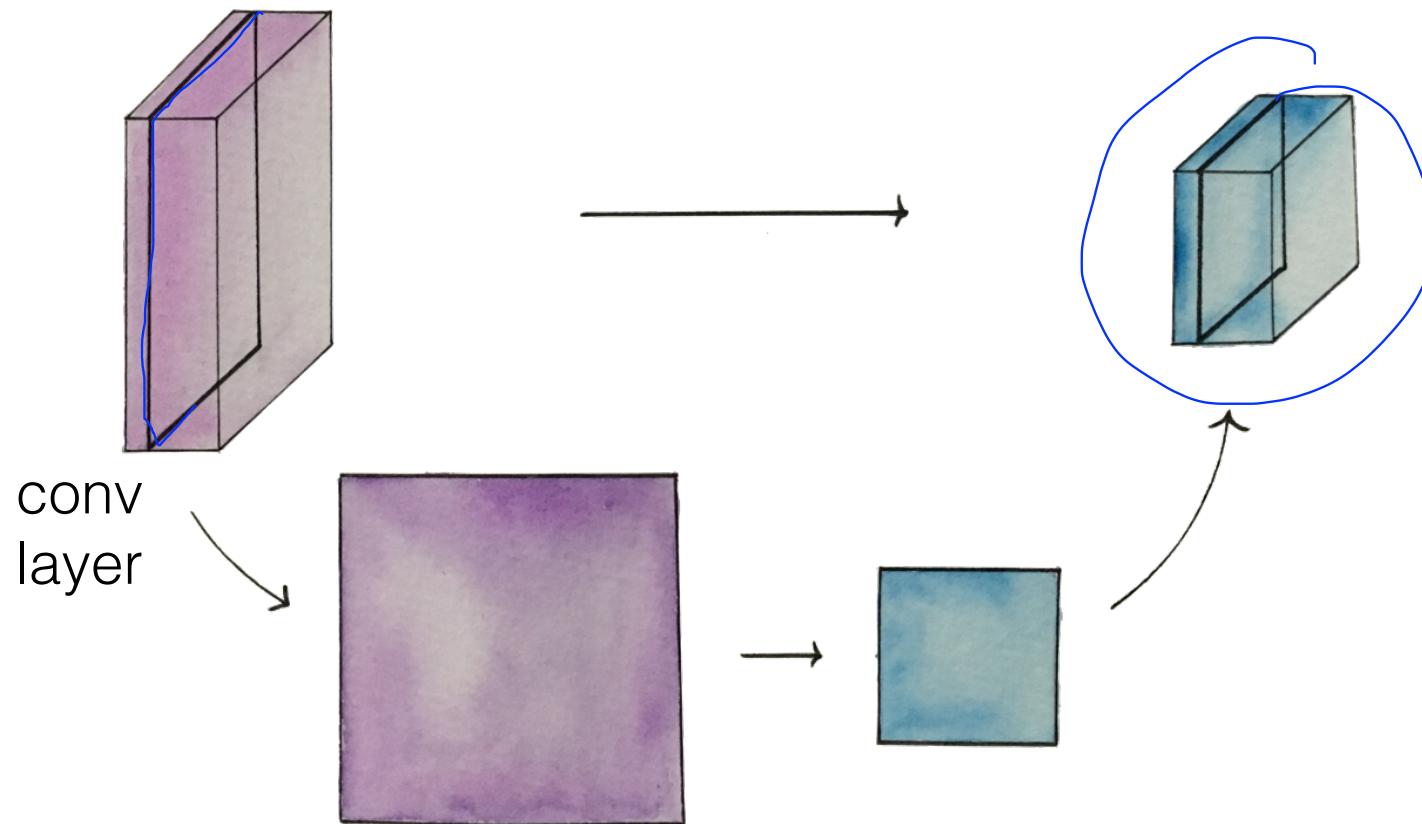
# Pooling layer (sampling)



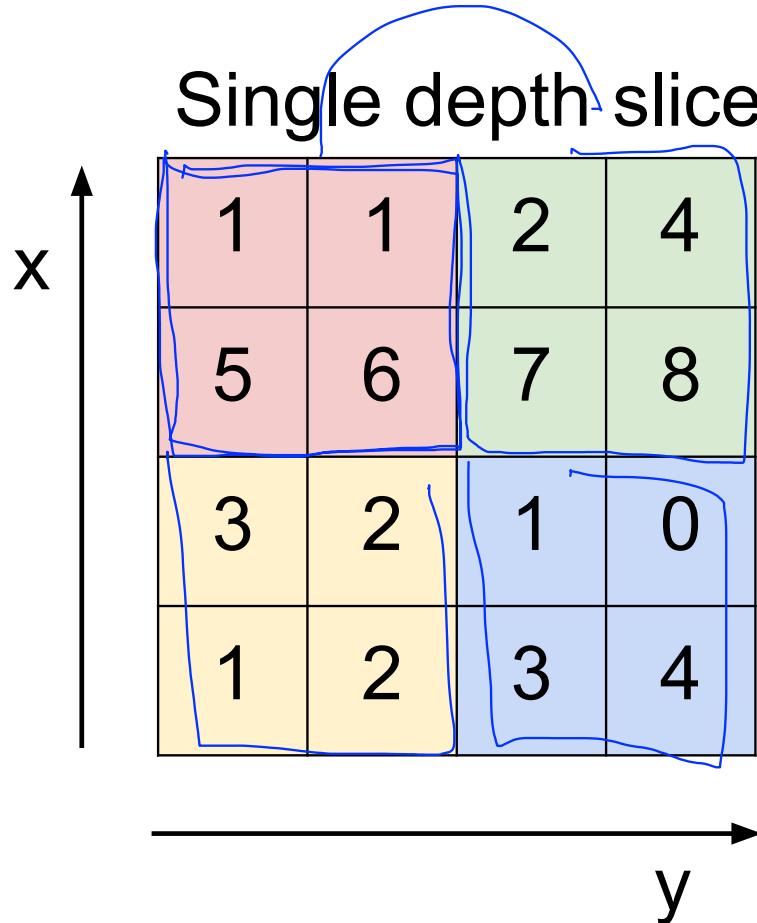
# Pooling layer (sampling)



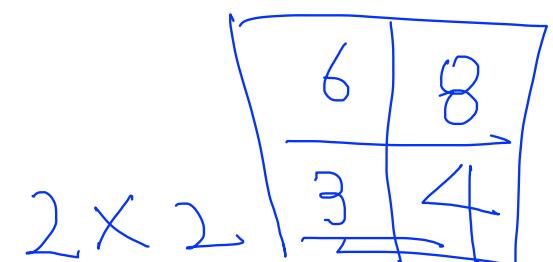
# Pooling layer (sampling)



# MAX POOLING

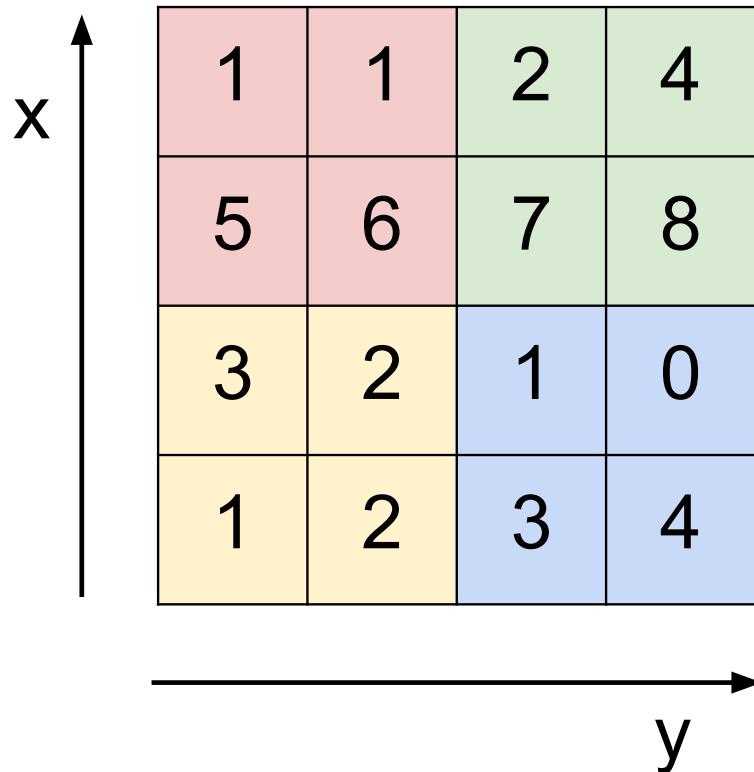


max pool with  $2 \times 2$  filters  
and stride 2

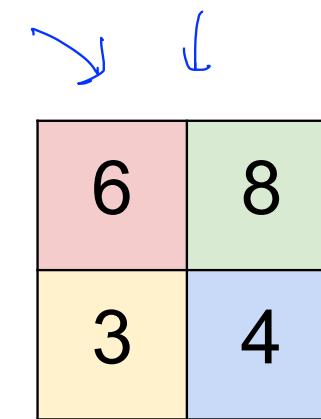


# MAX POOLING

Single depth slice

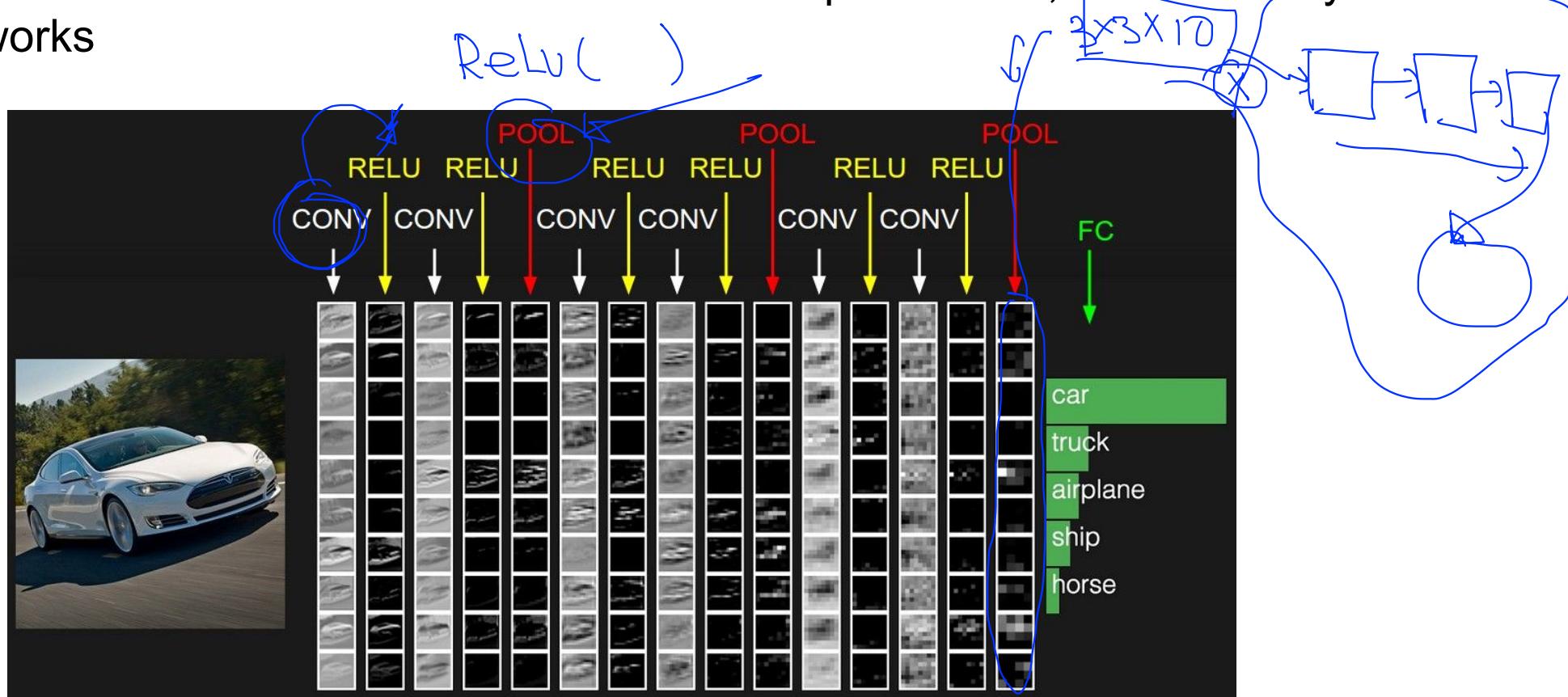


max pool with 2x2 filters  
and stride 2



# Fully Connected Layer (FC layer)

- Contains neurons that connect to the entire input volume, as in ordinary Neural Networks



## [ConvNetJS demo: training on CIFAR-10]



<http://cs.stanford.edu/people/karpathy/convnetjs/demo/cifar10.html>

# Lecture II-3

## CNN case study

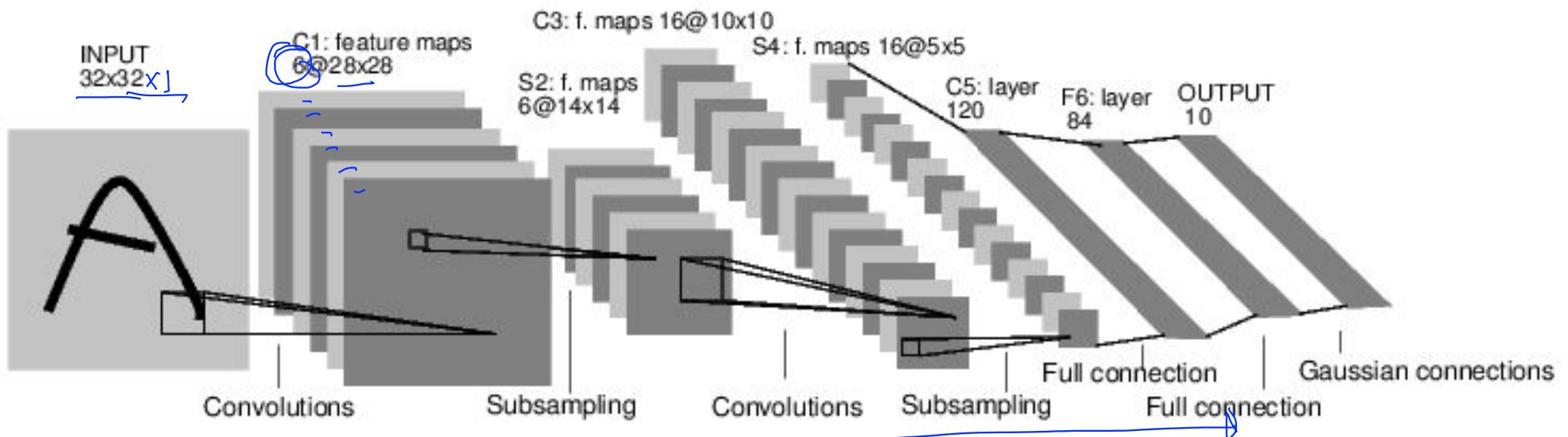
Sung Kim <hunkim+mr@gmail.com>  
<http://hunkim.github.io/ml/>

# Convolutional Neural Networks

<http://cs231n.stanford.edu/>

# Case Study: LeNet-5

[LeCun et al., 1998]

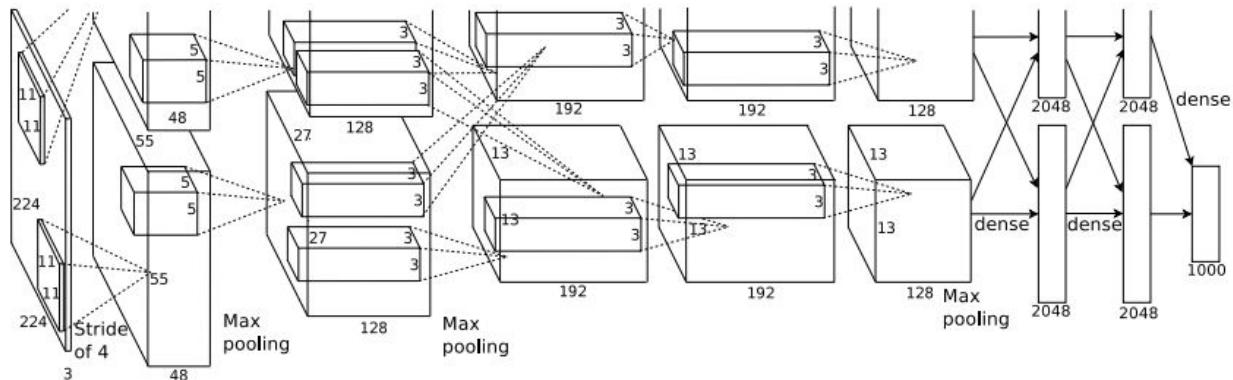


Conv filters were 5x5, applied at stride 1

Subsampling (Pooling) layers were 2x2 applied at stride 2  
i.e. architecture is [CONV-POOL-CONV-POOL-CONV-FC]

# Case Study: AlexNet

[Krizhevsky et al. 2012]



Input: 227x227x3 images

First layer (CONV1): 96 11x11 filters applied at stride 4

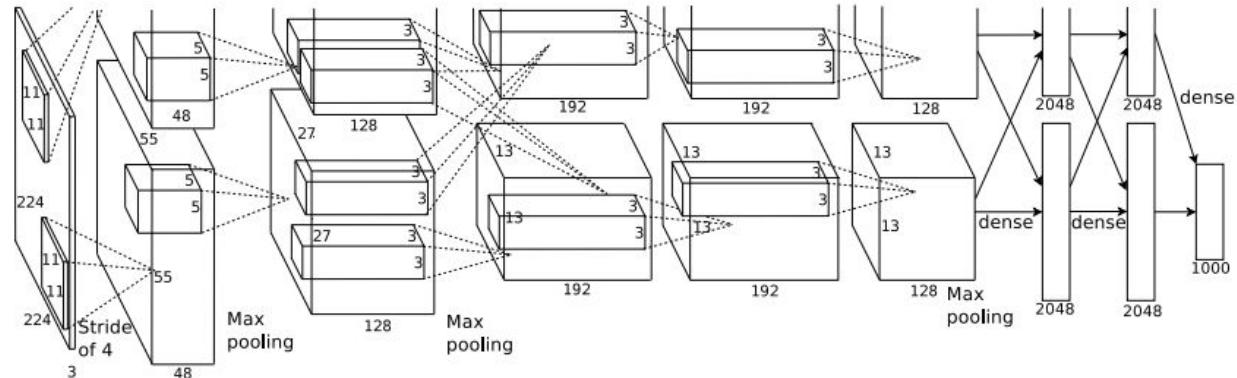
=>

Output volume [55x55x96]

Parameters:  $(11 \times 11 \times 3) \times 96 = 35K$

# Case Study: AlexNet

[Krizhevsky et al. 2012]



Input: 227x227x3 images

After CONV1: 55x55x96

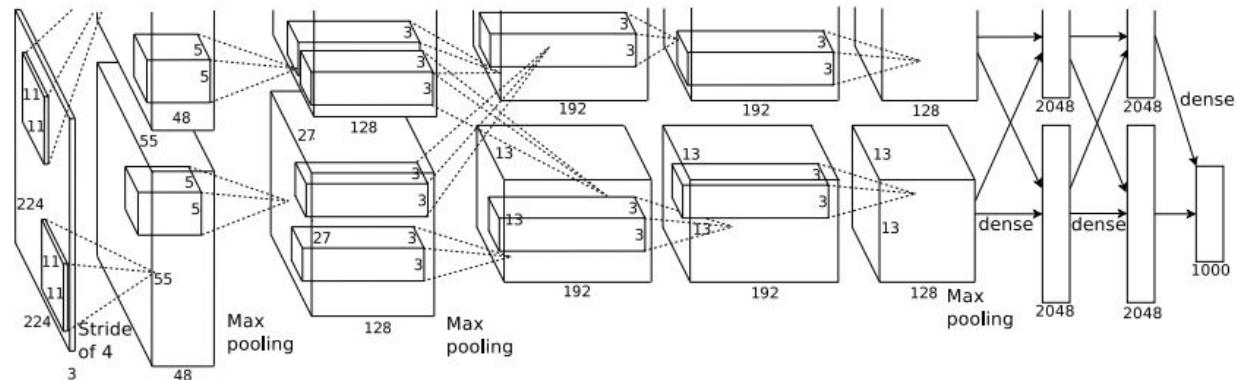
**Second layer (POOL1): 3x3 filters applied at stride 2**

Output volume: 27x27x96

Parameters: 0!

# Case Study: AlexNet

[Krizhevsky et al. 2012]



Input: 227x227x3 images

After CONV1: 55x55x96

After POOL1: 27x27x96

...

# Case Study: AlexNet

[Krizhevsky et al. 2012]

Full (simplified) AlexNet architecture:

[227x227x3] INPUT

[55x55x96] CONV1: 96 11x11 filters at stride 4, pad 0

[27x27x96] MAX POOL1: 3x3 filters at stride 2

[27x27x96] NORM1: Normalization layer

[27x27x256] CONV2: 256 5x5 filters at stride 1, pad 2

[13x13x256] MAX POOL2: 3x3 filters at stride 2

[13x13x256] NORM2: Normalization layer

[13x13x384] CONV3: 384 3x3 filters at stride 1, pad 1

[13x13x384] CONV4: 384 3x3 filters at stride 1, pad 1

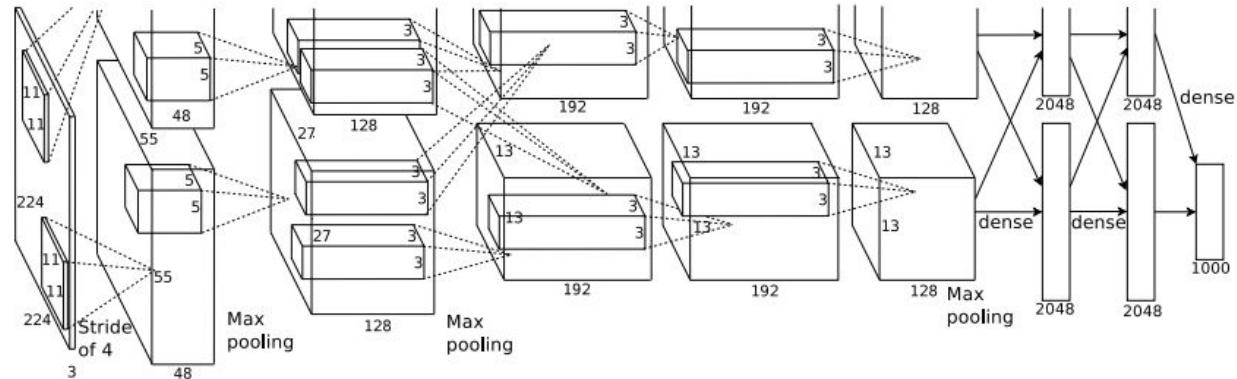
[13x13x256] CONV5: 256 3x3 filters at stride 1, pad 1

[6x6x256] MAX POOL3: 3x3 filters at stride 2

[4096] FC6: 4096 neurons

[4096] FC7: 4096 neurons

[1000] FC8: 1000 neurons (class scores)



# Case Study: AlexNet

[Krizhevsky et al. 2012]

Full (simplified) AlexNet architecture:

[227x227x3] INPUT

[55x55x96] CONV1: 96 11x11 filters at stride 4, pad 0

[27x27x96] MAX POOL1: 3x3 filters at stride 2

[27x27x96] NORM1: Normalization layer

[27x27x256] CONV2: 256 5x5 filters at stride 1, pad 2

[13x13x256] MAX POOL2: 3x3 filters at stride 2

[13x13x256] NORM2: Normalization layer

[13x13x384] CONV3: 384 3x3 filters at stride 1, pad 1

[13x13x384] CONV4: 384 3x3 filters at stride 1, pad 1

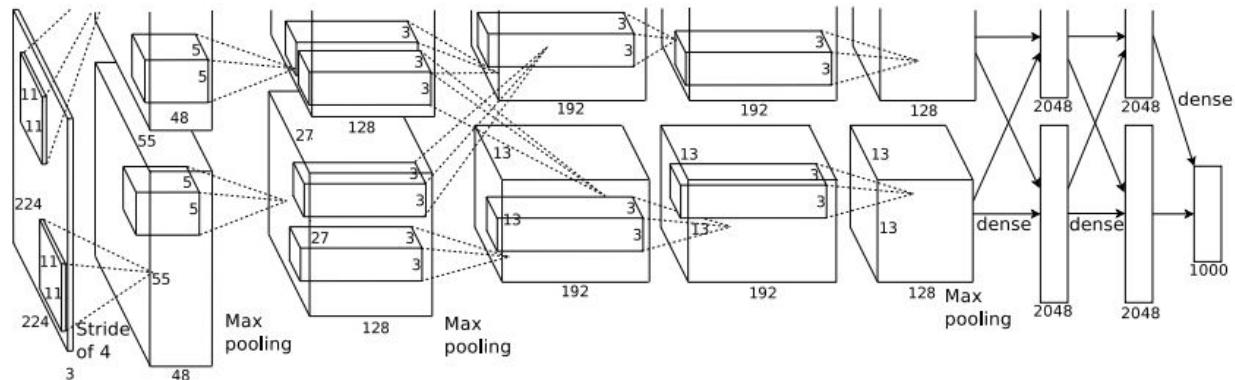
[13x13x256] CONV5: 256 3x3 filters at stride 1, pad 1

[6x6x256] MAX POOL3: 3x3 filters at stride 2

[4096] FC6: 4096 neurons

[4096] FC7: 4096 neurons

[1000] FC8: 1000 neurons (class scores)

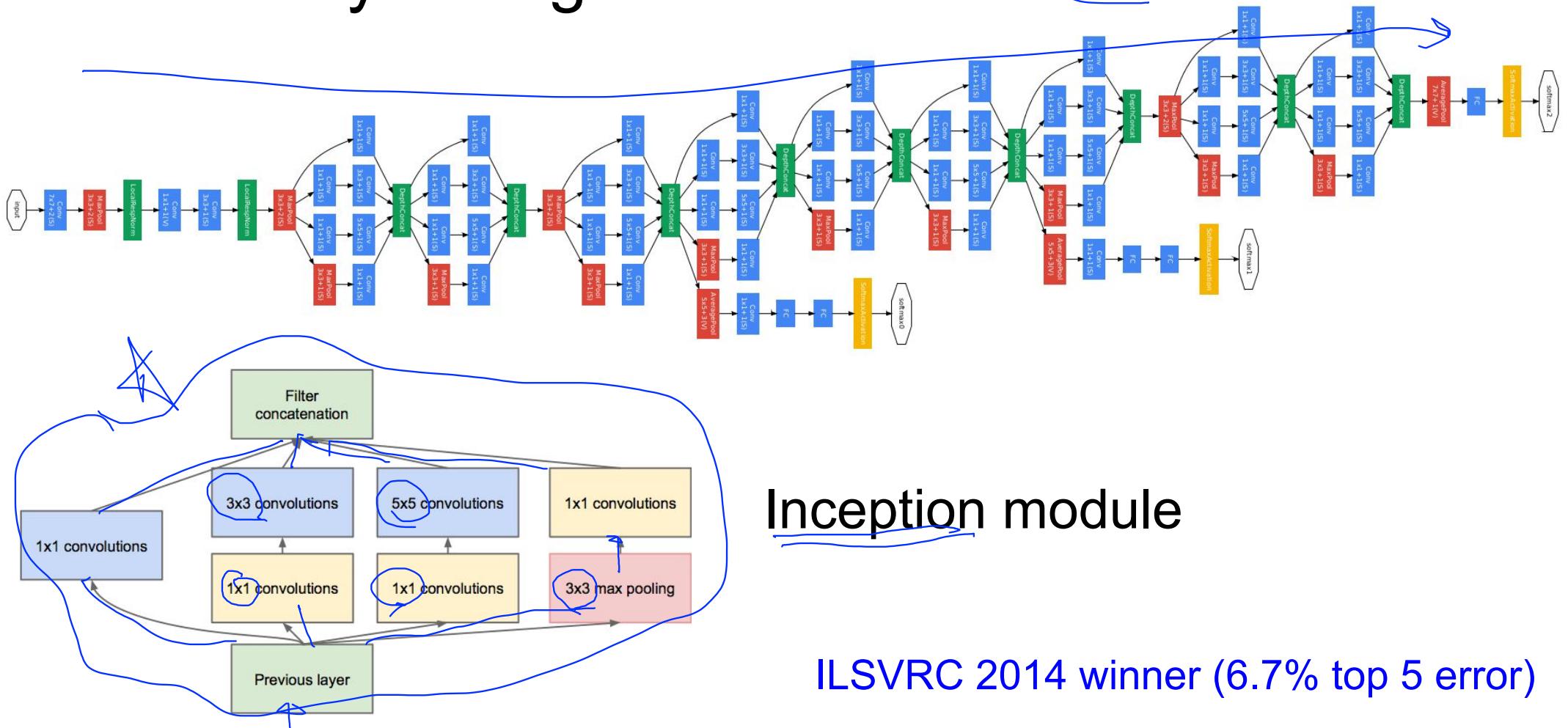


## Details/Retrospectives:

- first use of ReLU
- used Norm layers (not common anymore)
- heavy data augmentation
- dropout 0.5
- batch size 128
- SGD Momentum 0.9
- Learning rate 1e-2, reduced by 10 manually when val accuracy plateaus
- L2 weight decay 5e-4
- 7 CNN ensemble > 18.2% -> 15.4%

# Case Study: GoogLeNet

[Szegedy et al., 2014]



# Case Study: ResNet

[He et al., 2015]



ILSVRC 2015 winner (3.6% top 5 error)

MSRA @ ILSVRC & COCO 2015 Competitions

Microsoft Research

• **1st places** in all five main tracks

- ImageNet Classification: “Ultra-deep” (quote Yann) **152-layer** nets
- ImageNet Detection: **16%** better than 2nd
- ImageNet Localization: **27%** better than 2nd
- COCO Detection: **11%** better than 2nd
- COCO Segmentation: **12%** better than 2nd

\*improvements are relative numbers

ICCV15  
International Conference on Computer Vision

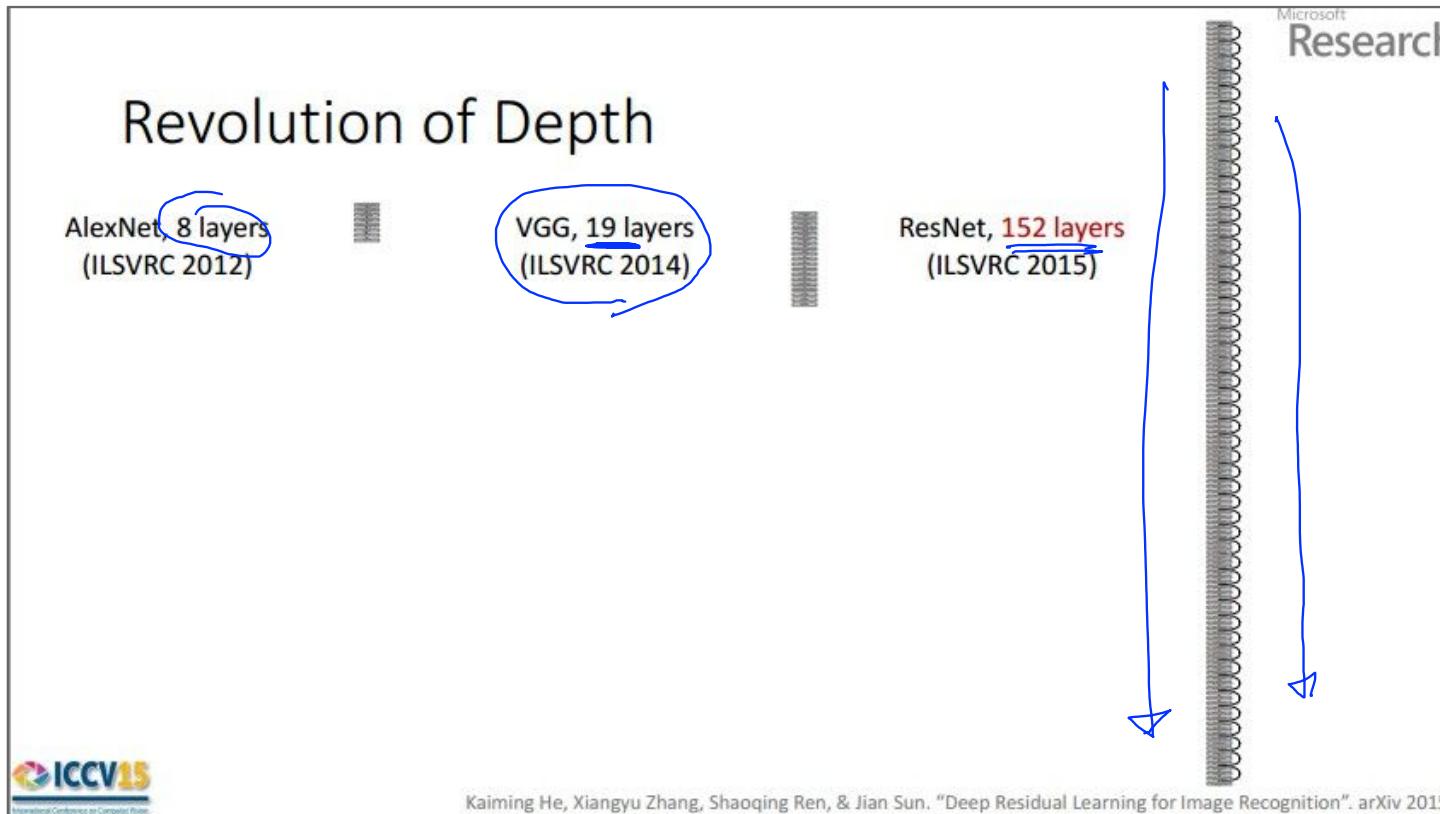
Kaiming He, Xiangyu Zhang, Shaoqing Ren, & Jian Sun. "Deep Residual Learning for Image Recognition". arXiv 2015.

Slide from Kaiming He's recent presentation <https://www.youtube.com/watch?v=1PGLj-uKT1w>

# Case Study: ResNet

[He et al., 2015]

ILSVRC 2015 winner (3.6% top 5 error)



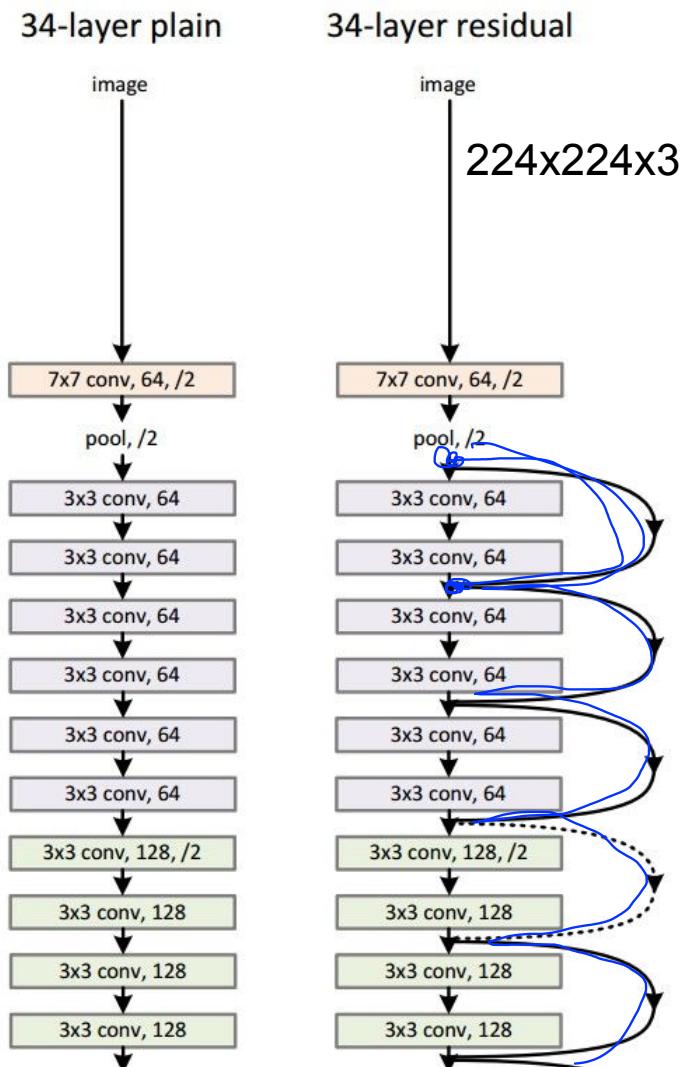
(slide from Kaiming He's recent presentation)

2-3 weeks of training  
on 8 GPU machine

at runtime: faster  
than a VGGNet!  
(even though it has  
8x more layers)

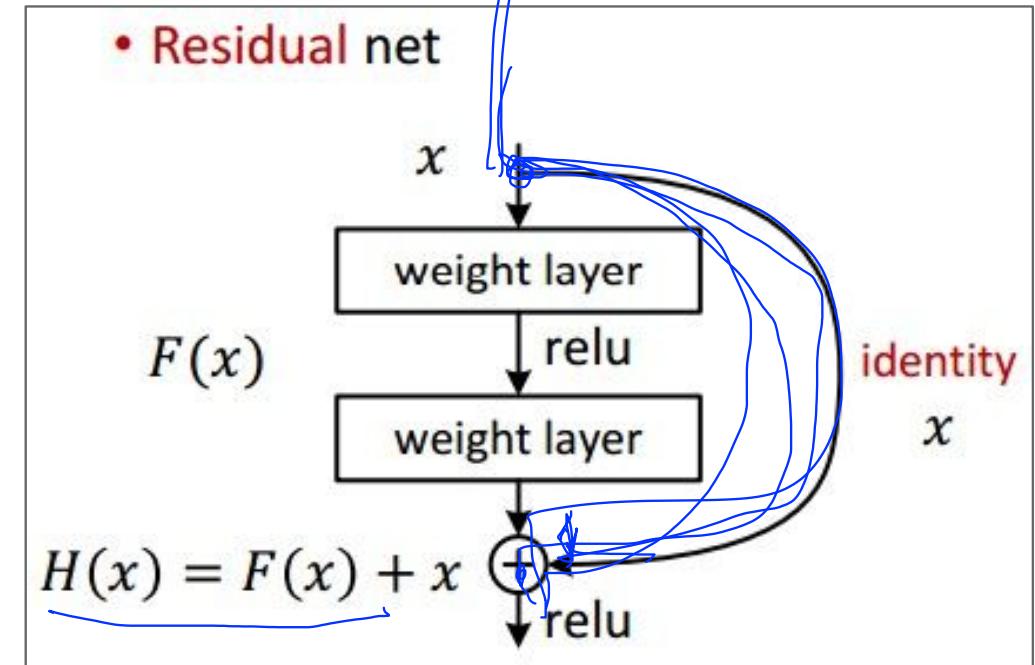
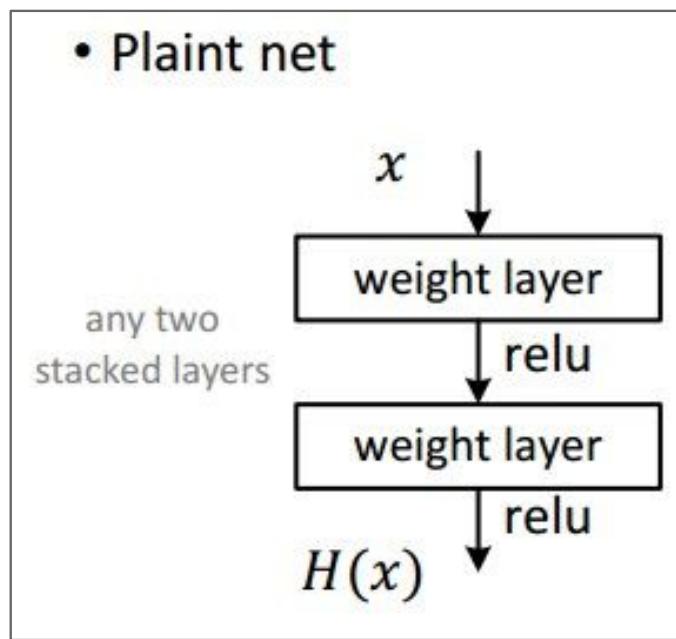
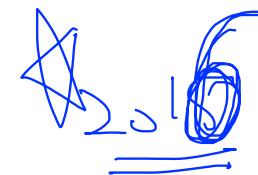
# Case Study: ResNet

[He et al., 2015]



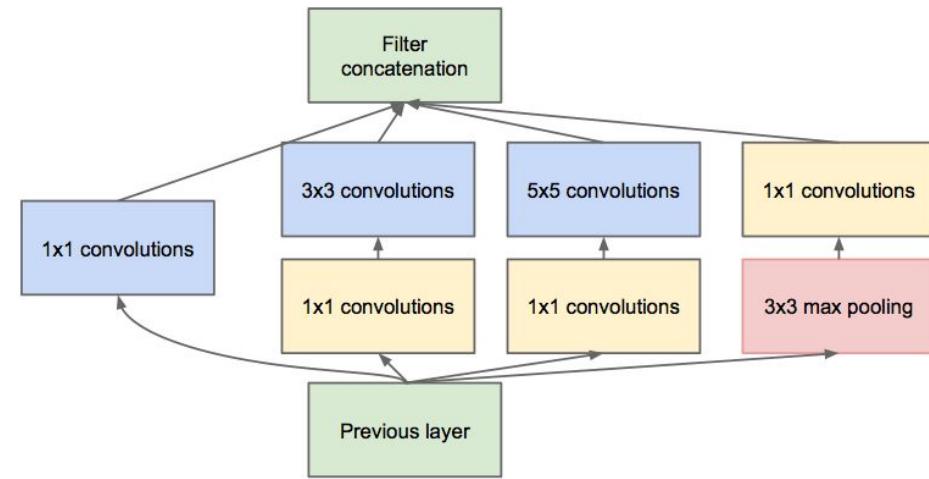
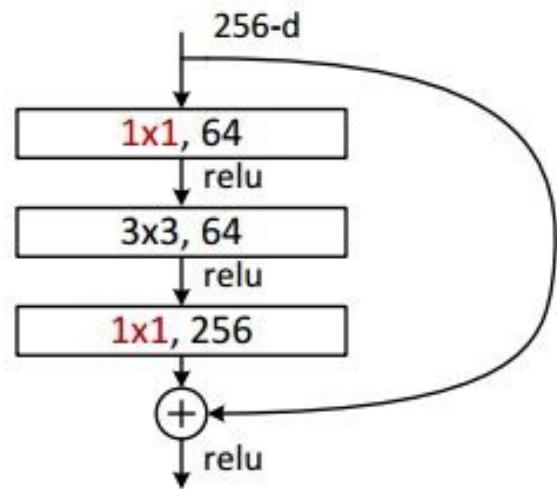
# Case Study: ResNet

[He et al., 2015]



# Case Study: ResNet

[He et al., 2015]



# Convolutional Neural Networks for Sentence Classification

[Yoon Kim, 2014]

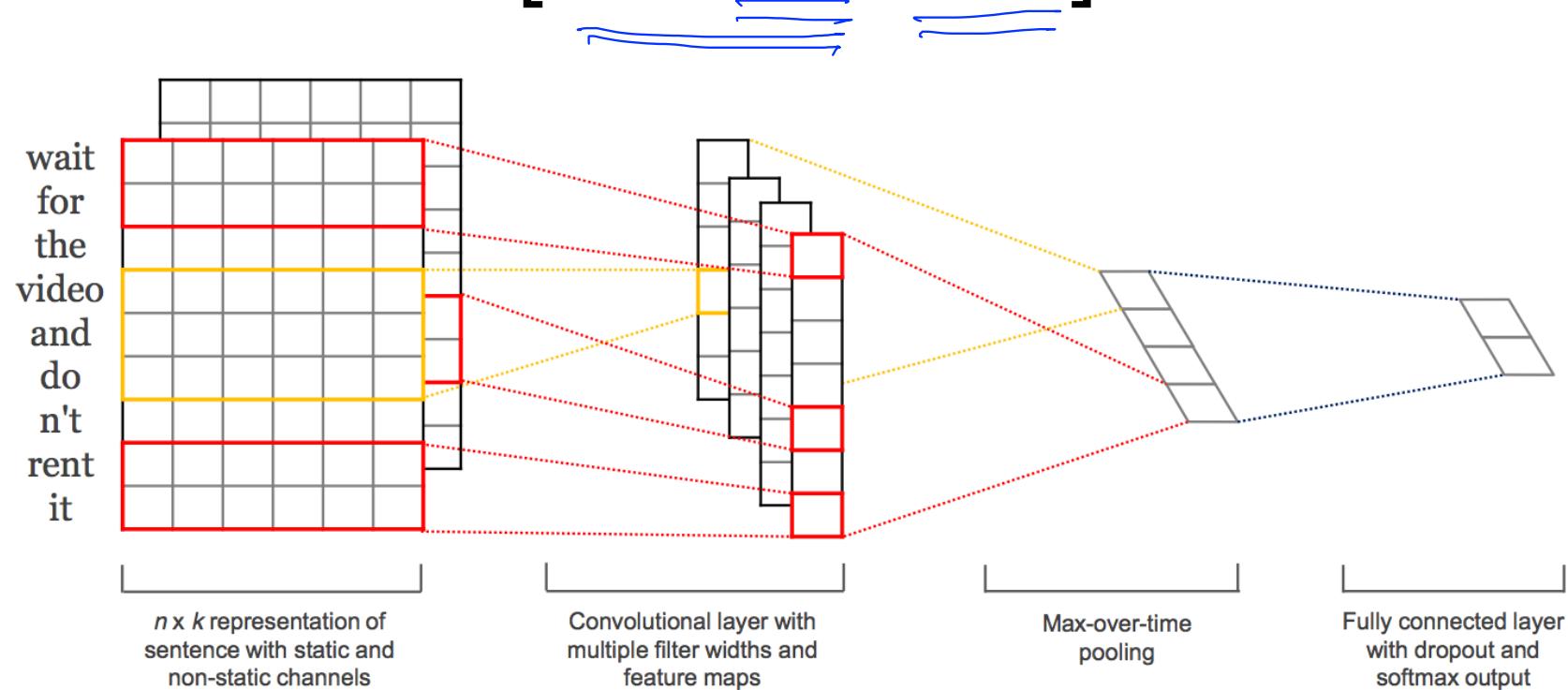
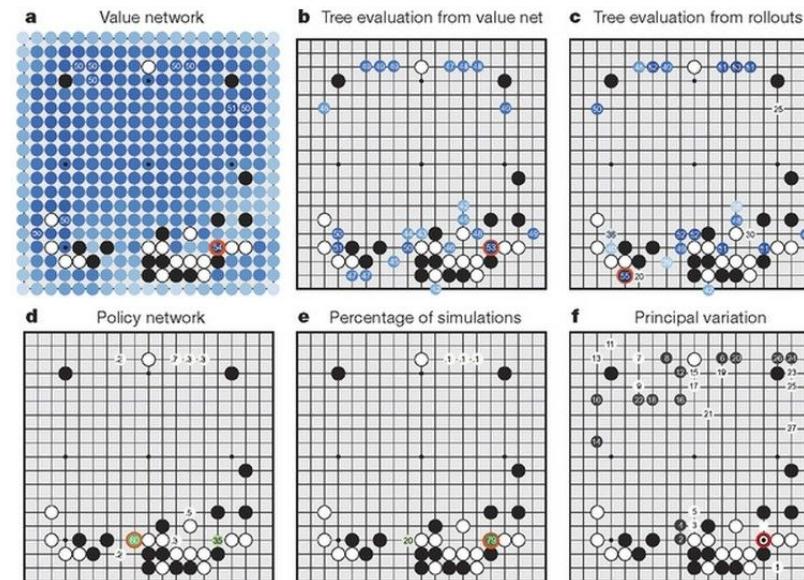


Figure 1: Model architecture with two channels for an example sentence.

# Case Study Bonus: DeepMind's AlphaGo



The input to the policy network is a  $19 \times 19 \times 48$  image stack consisting of 48 feature planes. The first hidden layer zero pads the input into a  $23 \times 23$  image, then convolves  $k$  filters of kernel size  $5 \times 5$  with stride 1 with the input image and applies a rectifier nonlinearity. Each of the subsequent hidden layers 2 to 12 zero pads the respective previous hidden layer into a  $21 \times 21$  image, then convolves  $k$  filters of kernel size  $3 \times 3$  with stride 1, again followed by a rectifier nonlinearity. The final layer convolves 1 filter of kernel size  $1 \times 1$  with stride 1, with a different bias for each position, and applies a softmax function. The match version of AlphaGo used  $k = 192$  filters; Fig. 2b and Extended Data Table 3 additionally show the results of training with  $k = 128, 256$  and  $384$  filters.

## policy network:

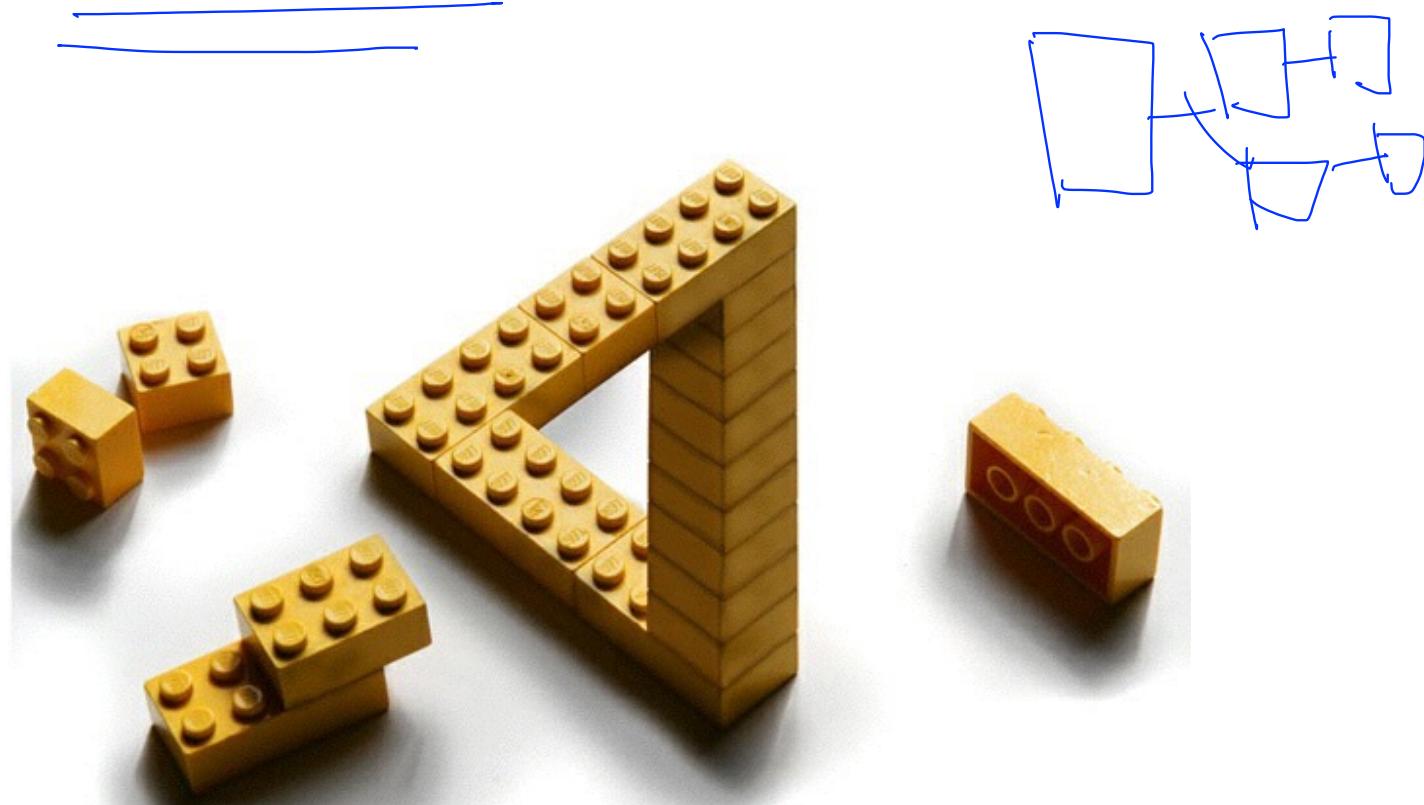
[ $19 \times 19 \times 48$ ] Input

CONV1: 192  $5 \times 5$  filters , stride 1, pad 2 => [ $19 \times 19 \times 192$ ]

CONV2..12: 192  $3 \times 3$  filters, stride 1, pad 1 => [ $19 \times 19 \times 192$ ]

CONV: 1  $1 \times 1$  filter, stride 1, pad 0 => [ $19 \times 19$ ] (*probability map of promising moves*)

# 'The only limit is your imagination'



<http://itchyi.squarespace.com/thelatest/2012/5/17/the-only-limit-is-your-imagination.html>

Next  
Recurrent  
Neural Nets (RNN)

