

Lecture 12

CNN

A bit of history:

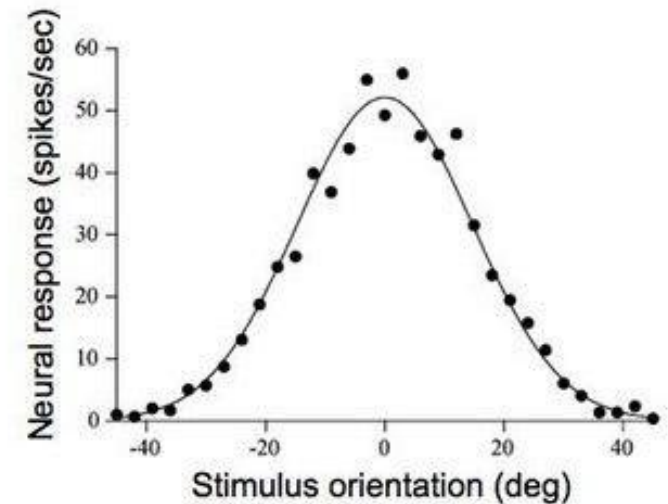
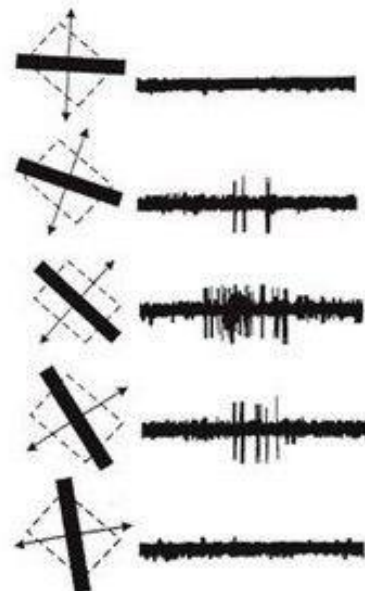
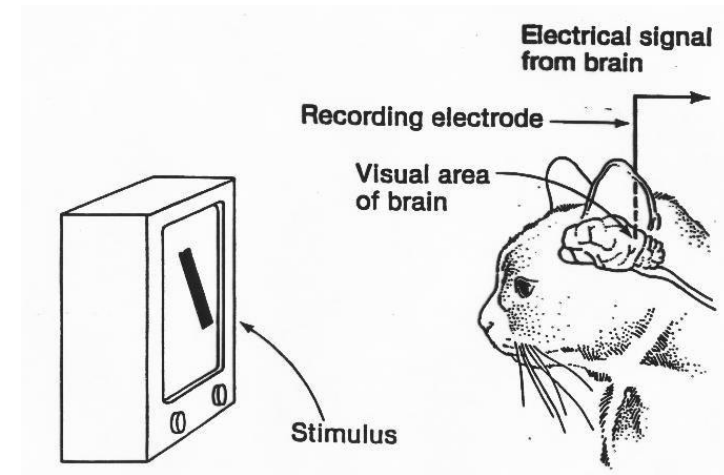
Hubel & Wiesel, 1959

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

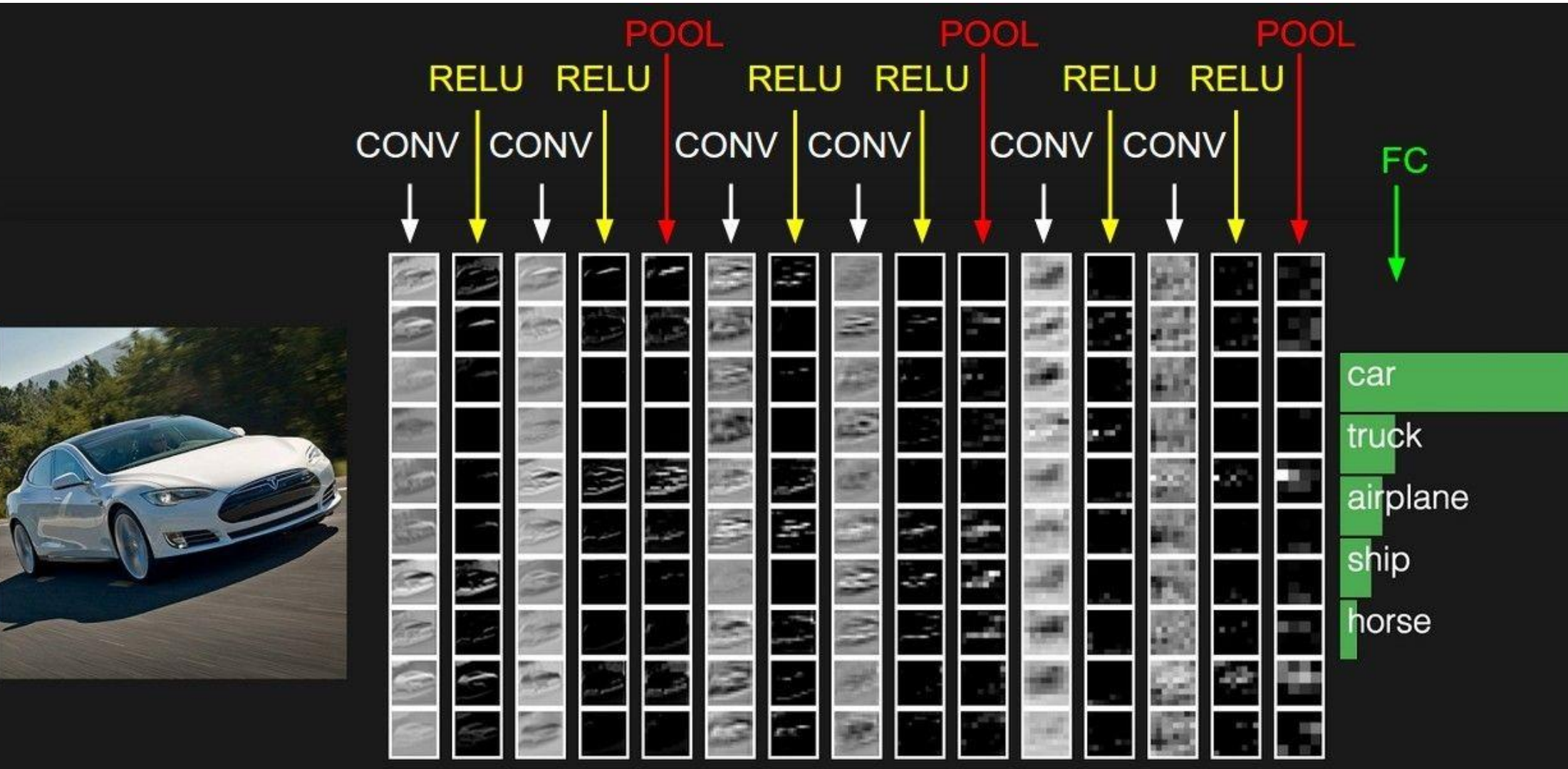
1962

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

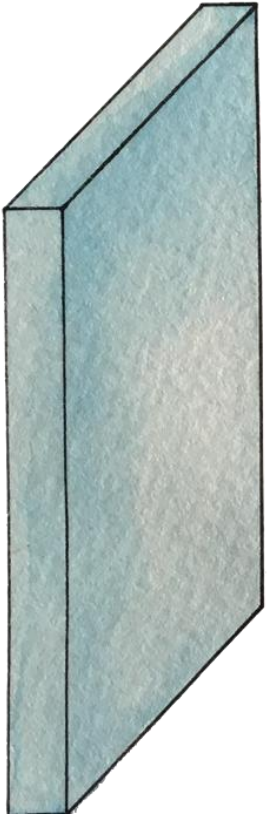
1968...



preview:

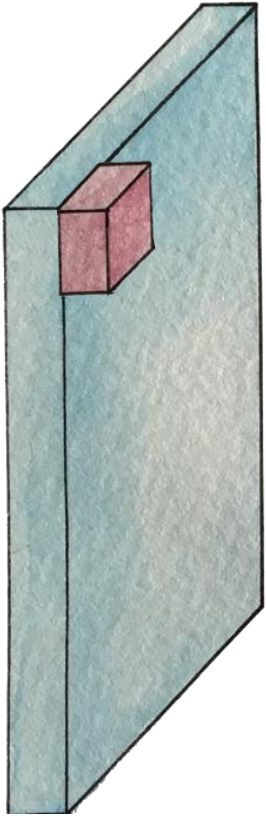


Start with an image (width x height x depth)



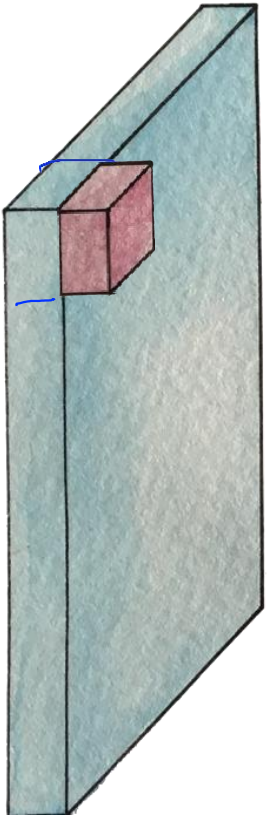
32x32x3 image

Let's focus on a small area only



32x32x3 image

Let's focus on a small area only (5x5x3)

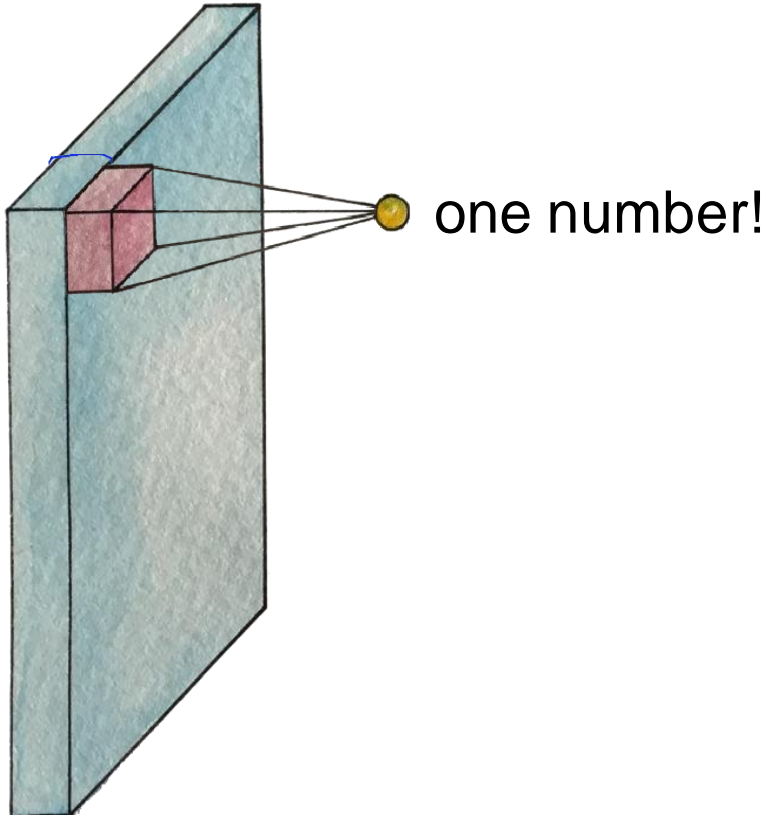


32x32x3 image



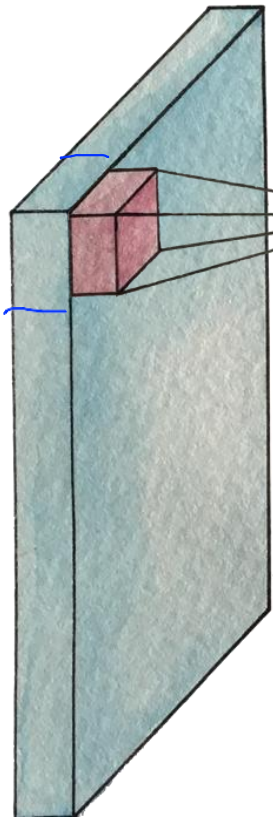
5x5x3 filter

Get one number using the filter



32x32x3 image

Get one number using the filter

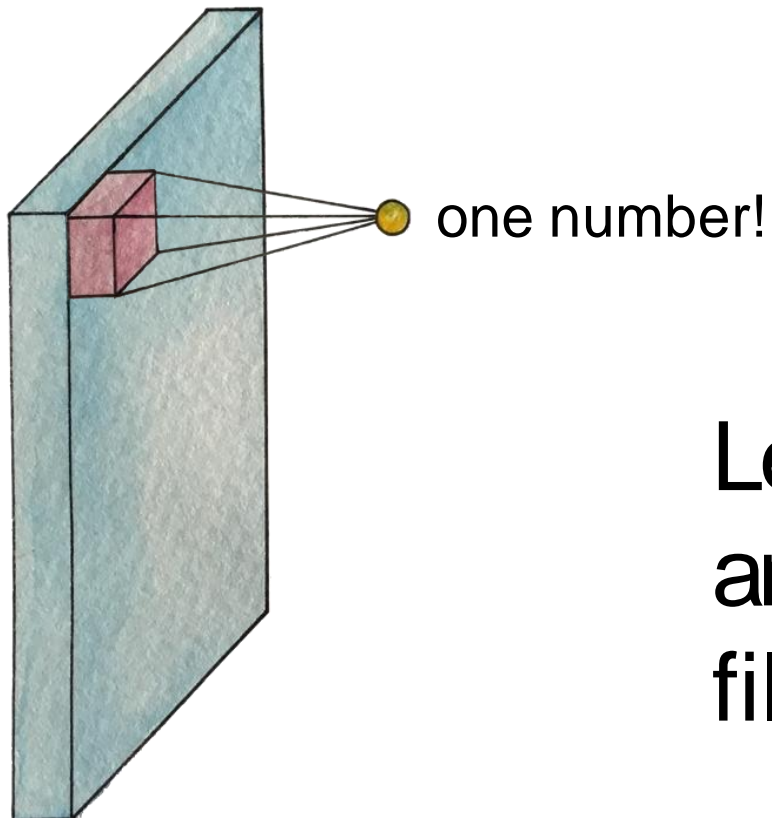


one number! $=Wx+b$

$=\text{ReLU}(Wx+b)$

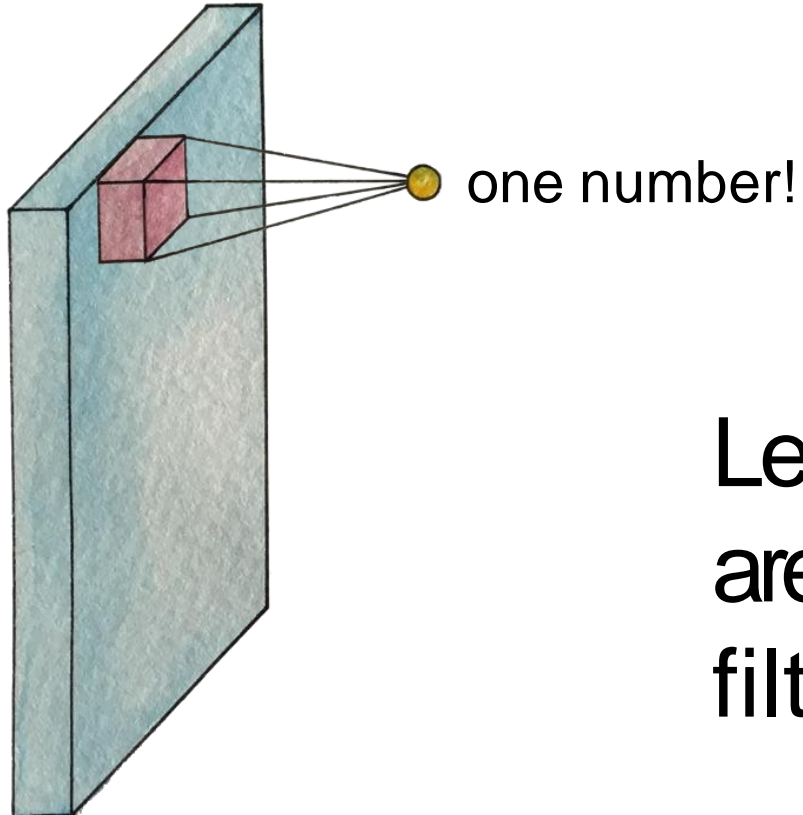
5x5x3 filter

32x32x3 image



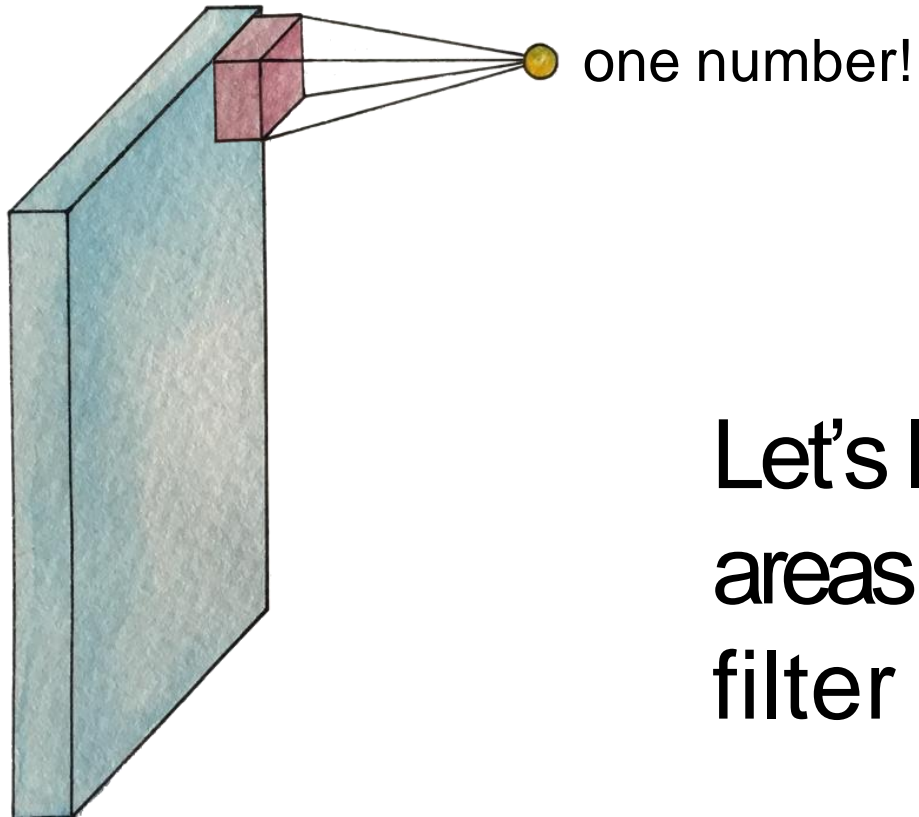
Let's look at other
areas with the same
filter (w)

32x32x3 image



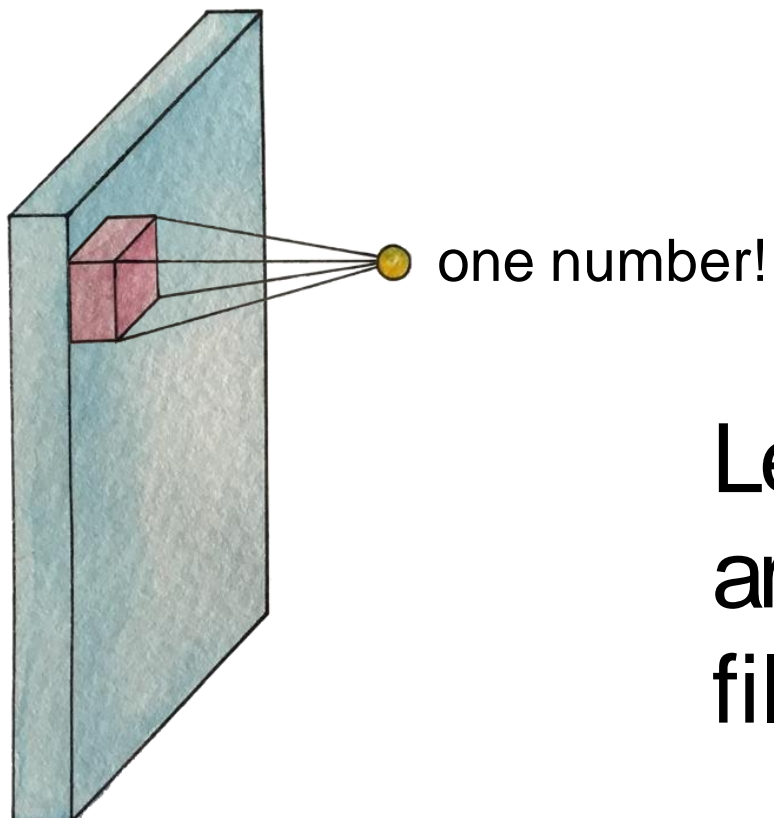
Let's look at other
areas with the same
filter (w)

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)

32x32x3 image

A closer look at spatial dimensions:

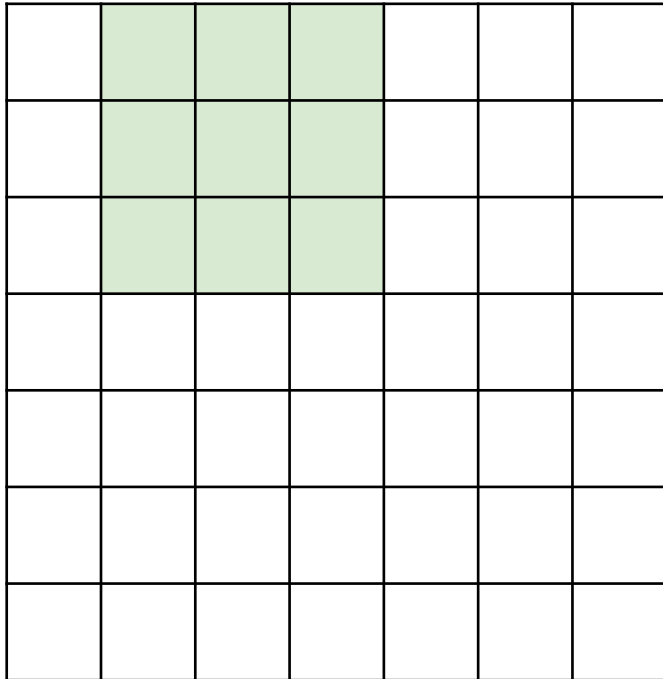
7

7

7x7 input (spatially)
assume 3x3 filter

A closer look at spatial dimensions:

7

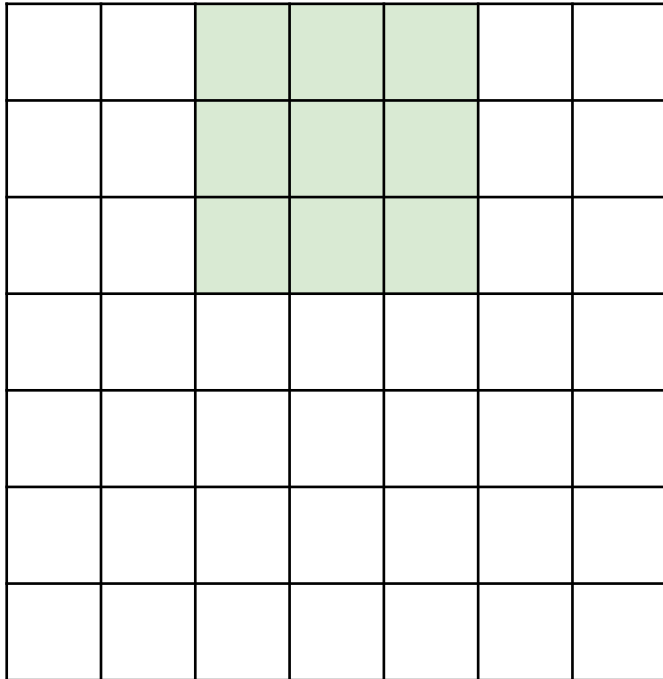


7x7 input (spatially)
assume 3x3 filter

7

A closer look at spatial dimensions:

7

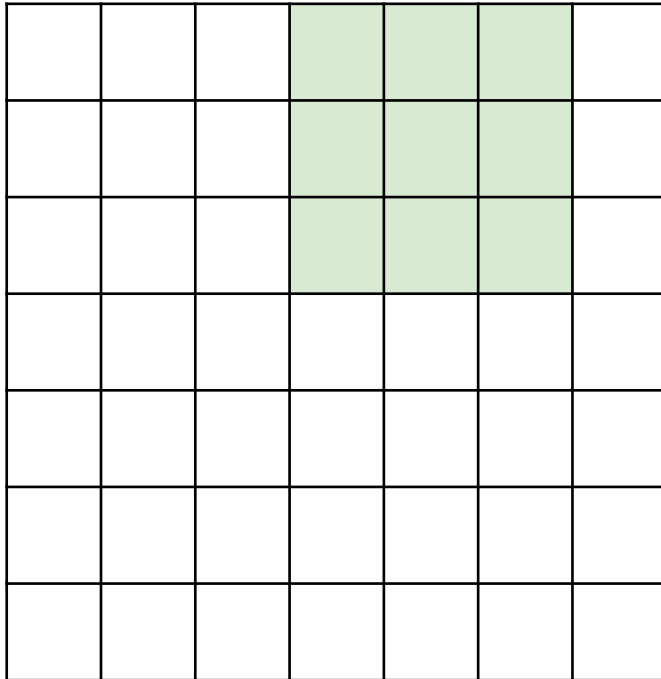


7x7 input (spatially)
assume 3x3 filter

7

A closer look at spatial dimensions:

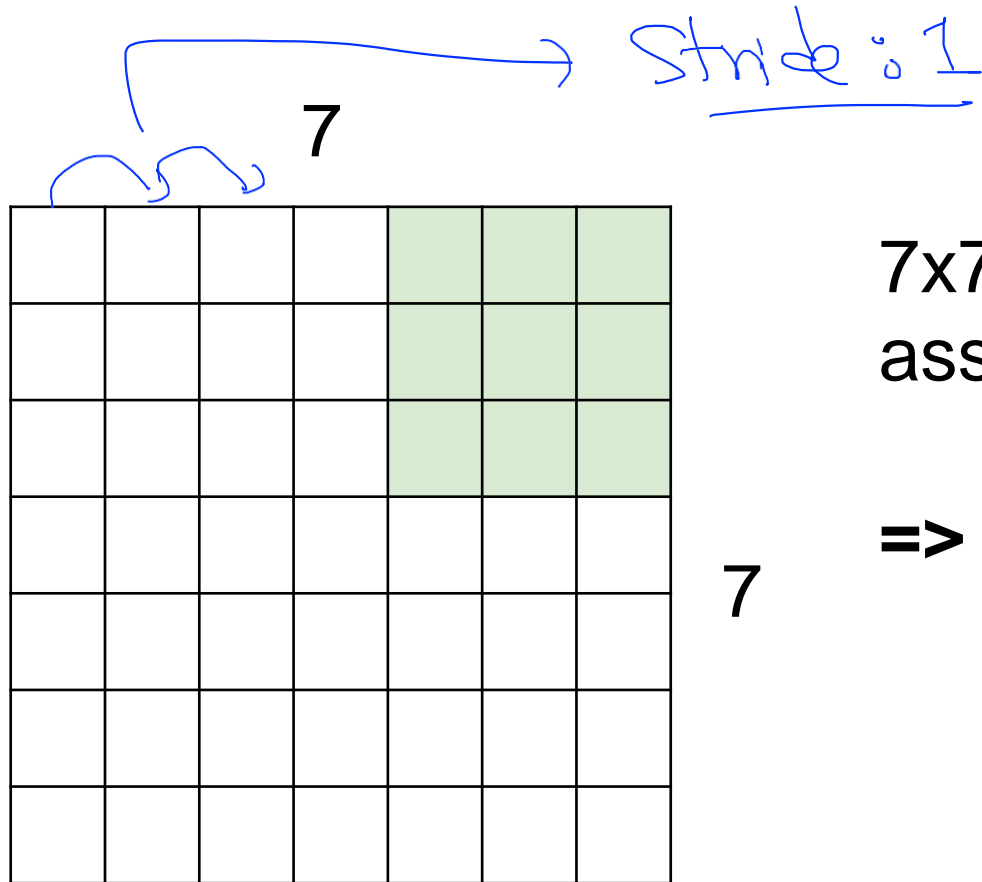
7



7

7x7 input (spatially)
assume 3x3 filter

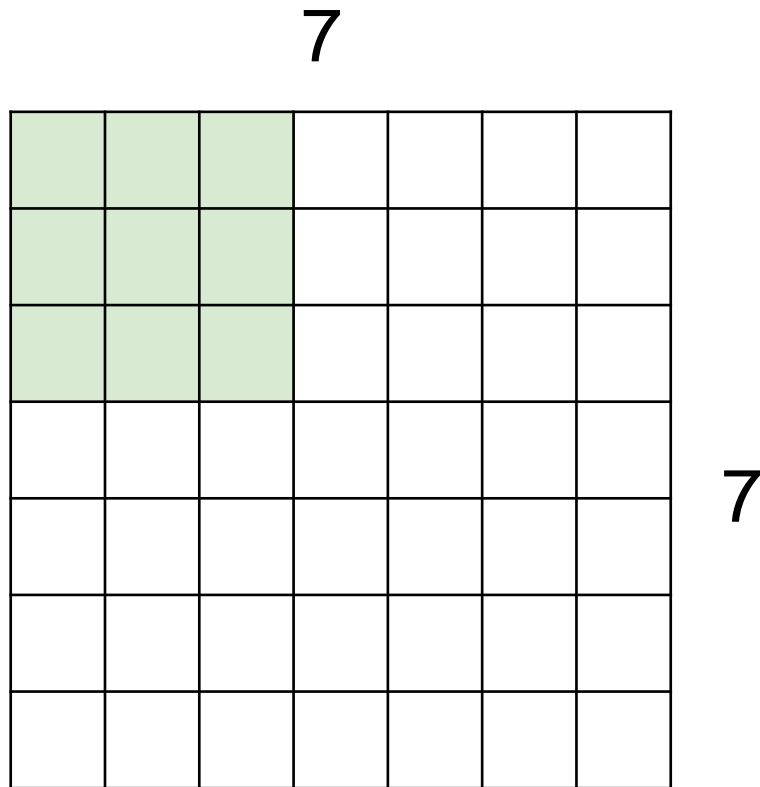
A closer look at spatial dimensions:



7x7 input (spatially)
assume 3x3 filter

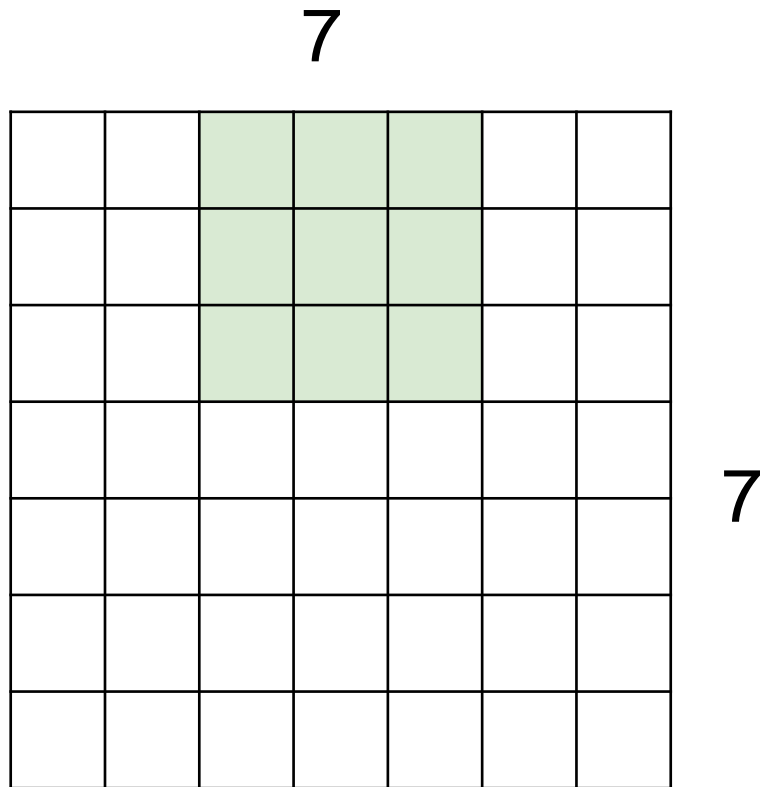
=> **5x5 output**

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:

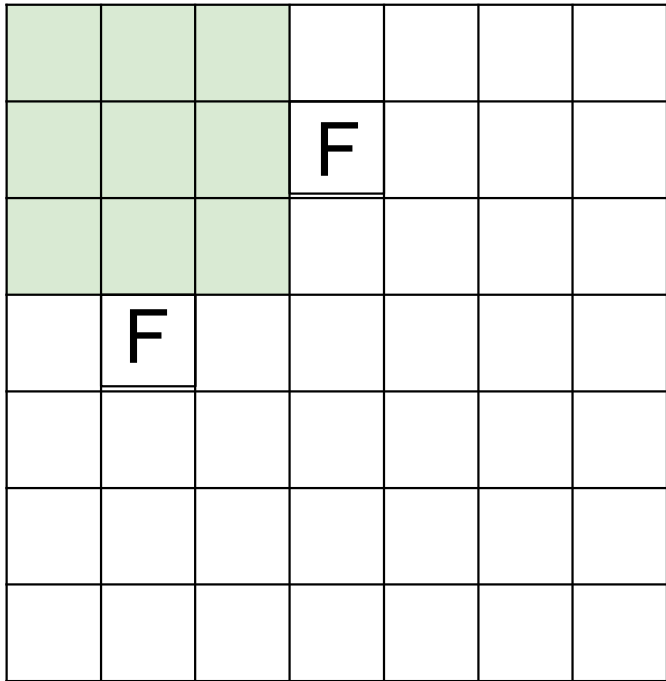
7

7

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

=> 3x3 output!

N



N

Output size:

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

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

$$\text{stride } 1 \Rightarrow (7 - 3) / 1 + 1 = 5$$

$$\text{stride } 2 \Rightarrow (7 - 3) / 2 + 1 = 3$$

$$\text{stride } 3 \Rightarrow (7 - 3) / 3 + 1 = 2.33 \text{ :}\backslash$$

In practice: Common to zero pad the border

0	0	0	0	0	0			
0								
0								
0								
0								

e.g. input 7x7

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

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

(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 7x7

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

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

7x7 output!

In practice: Common to zero pad the border

0	0	0	0	0	0			
0								
0								
0								
0								

e.g. input 7x7

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

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

7x7 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)

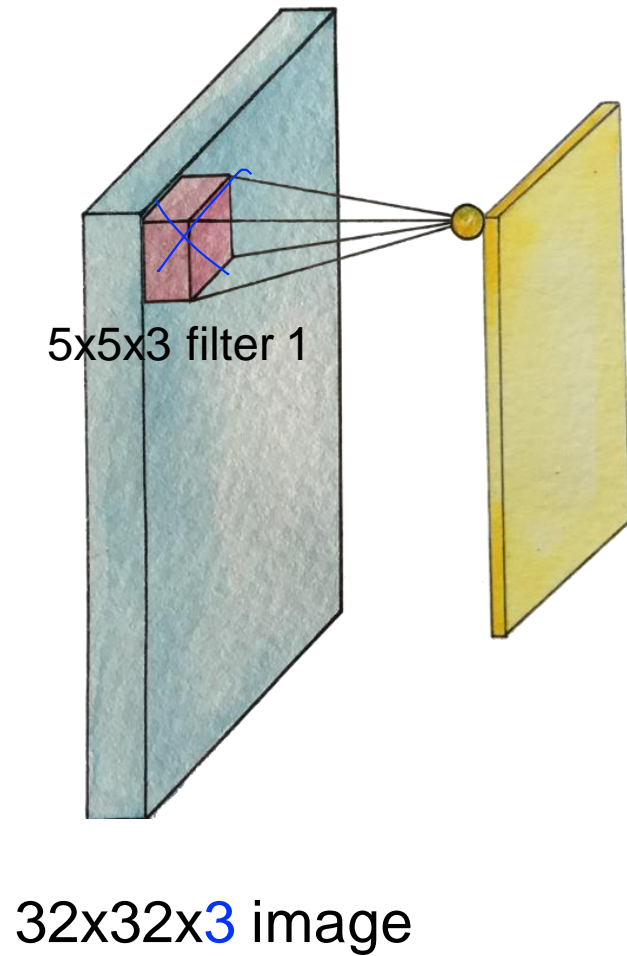
e.g. $F = 3 \Rightarrow$ zero pad with 1

$F = 5 \Rightarrow$ zero pad

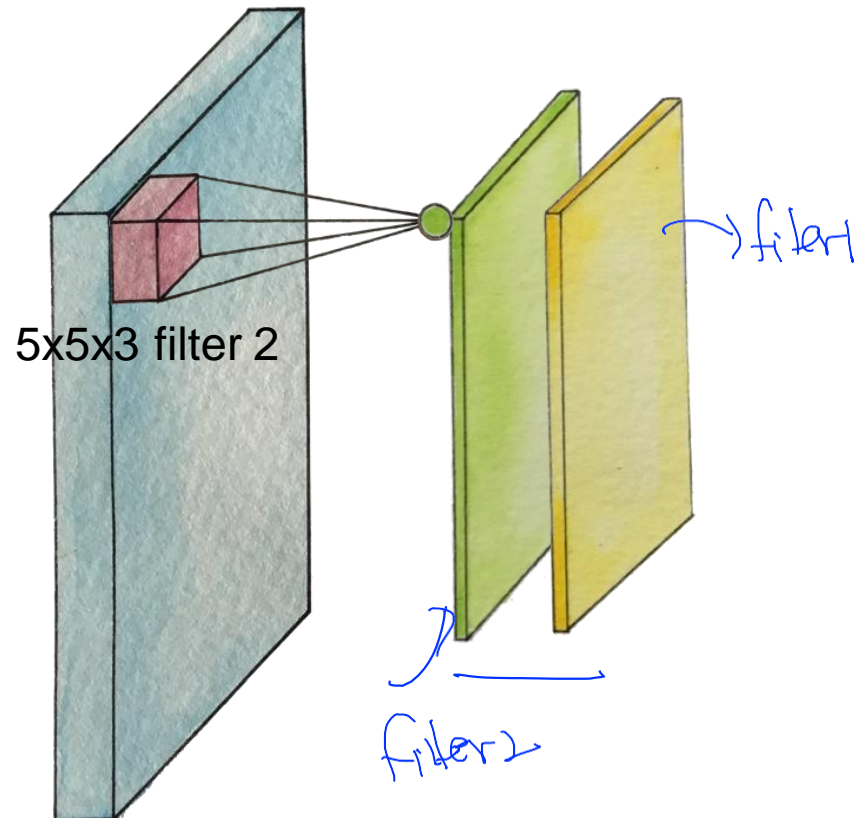
with 2 $F = 7 \Rightarrow$ zero

pad with 3

Swiping the entire image

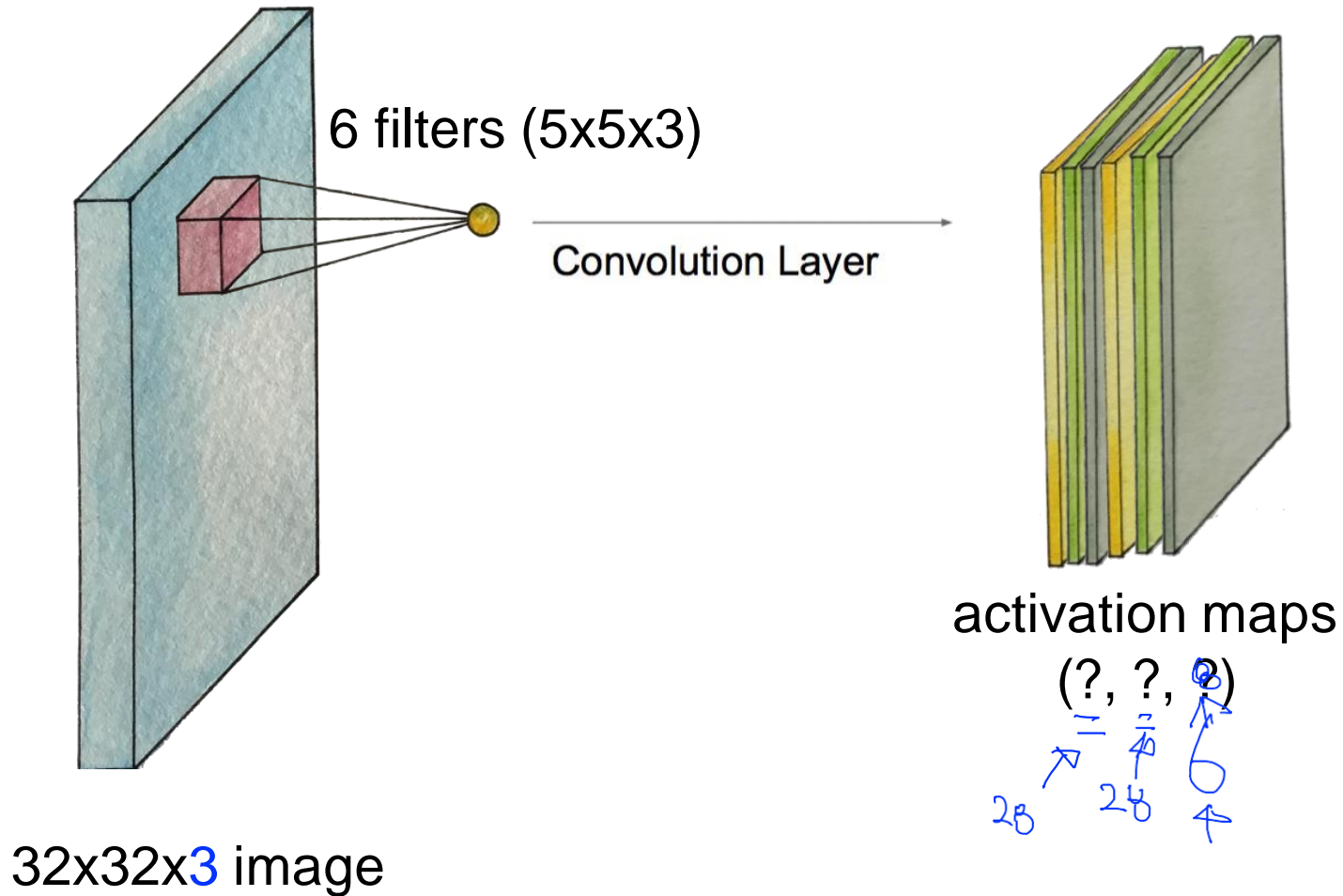


Swiping the entire image

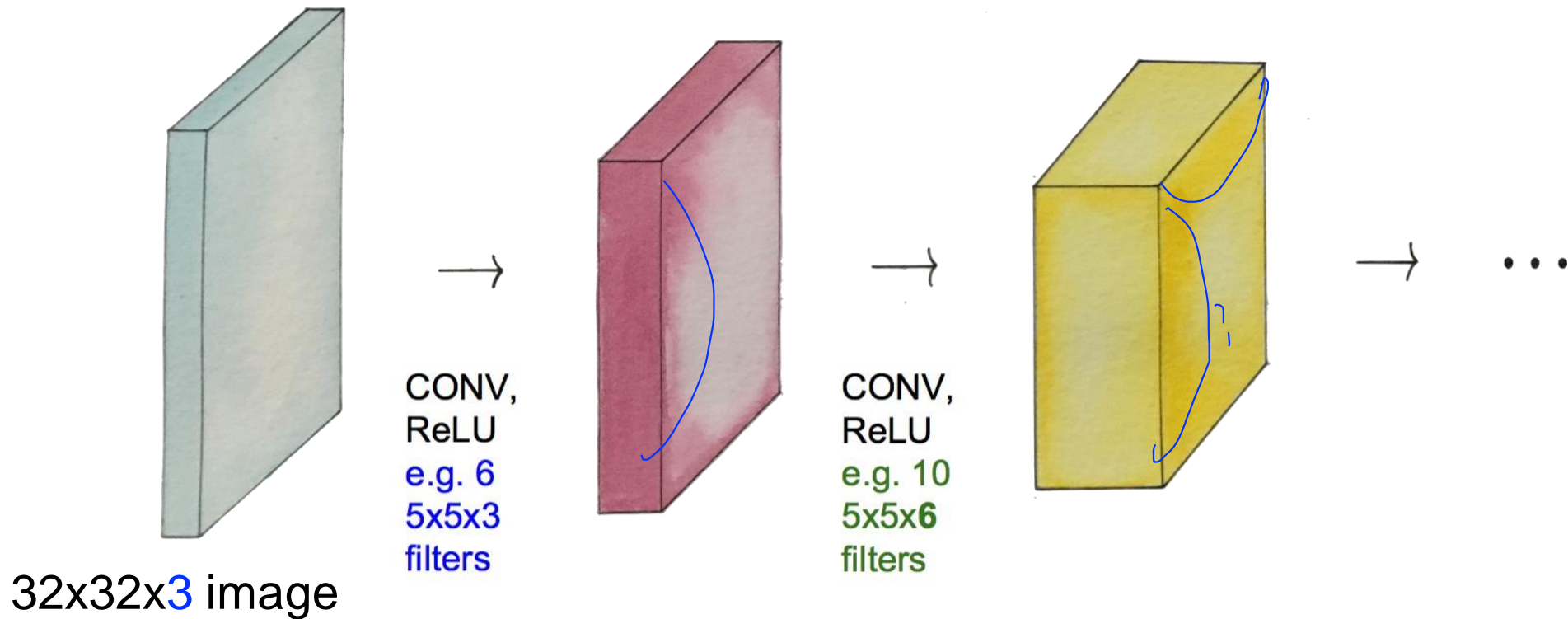


32x32x3 image

Swiping the entire image

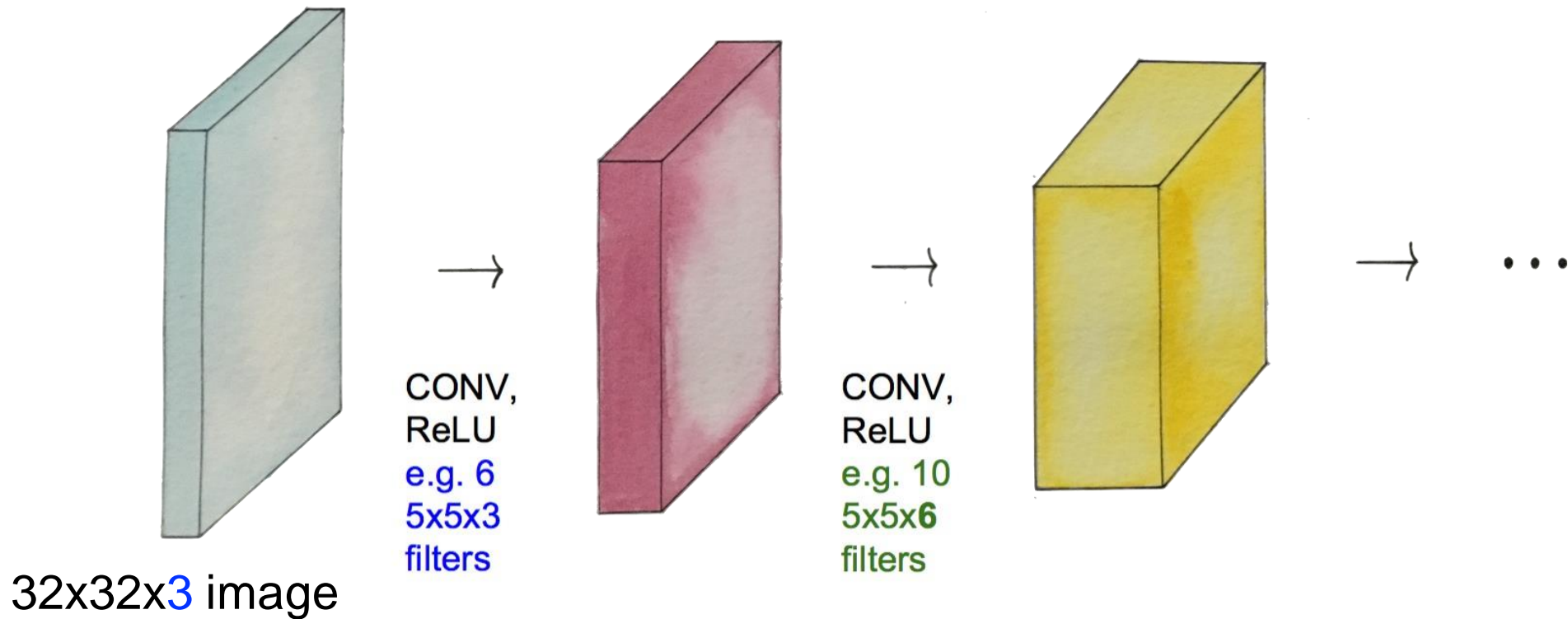


Convolution layers

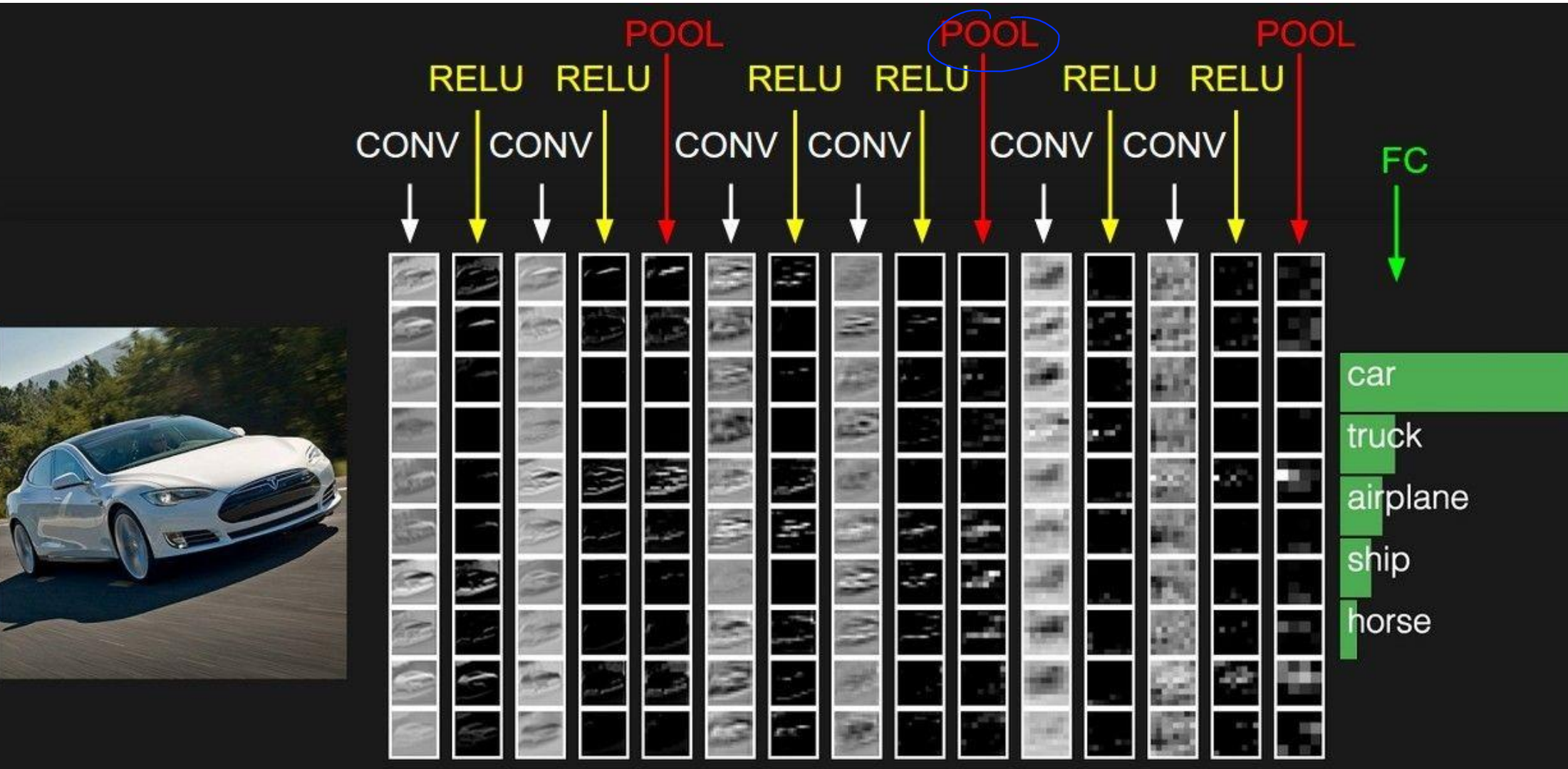


Convolution layers

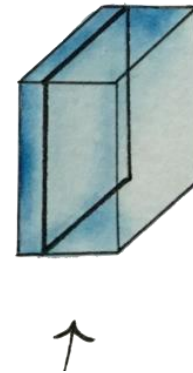
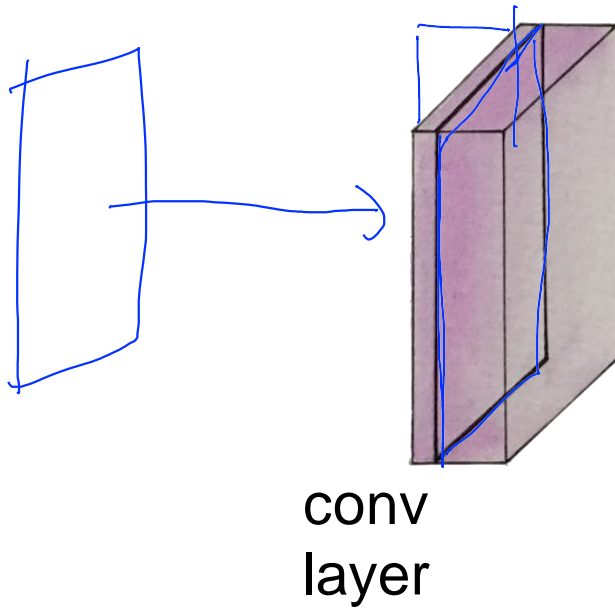
How many weight variables? How to set them?



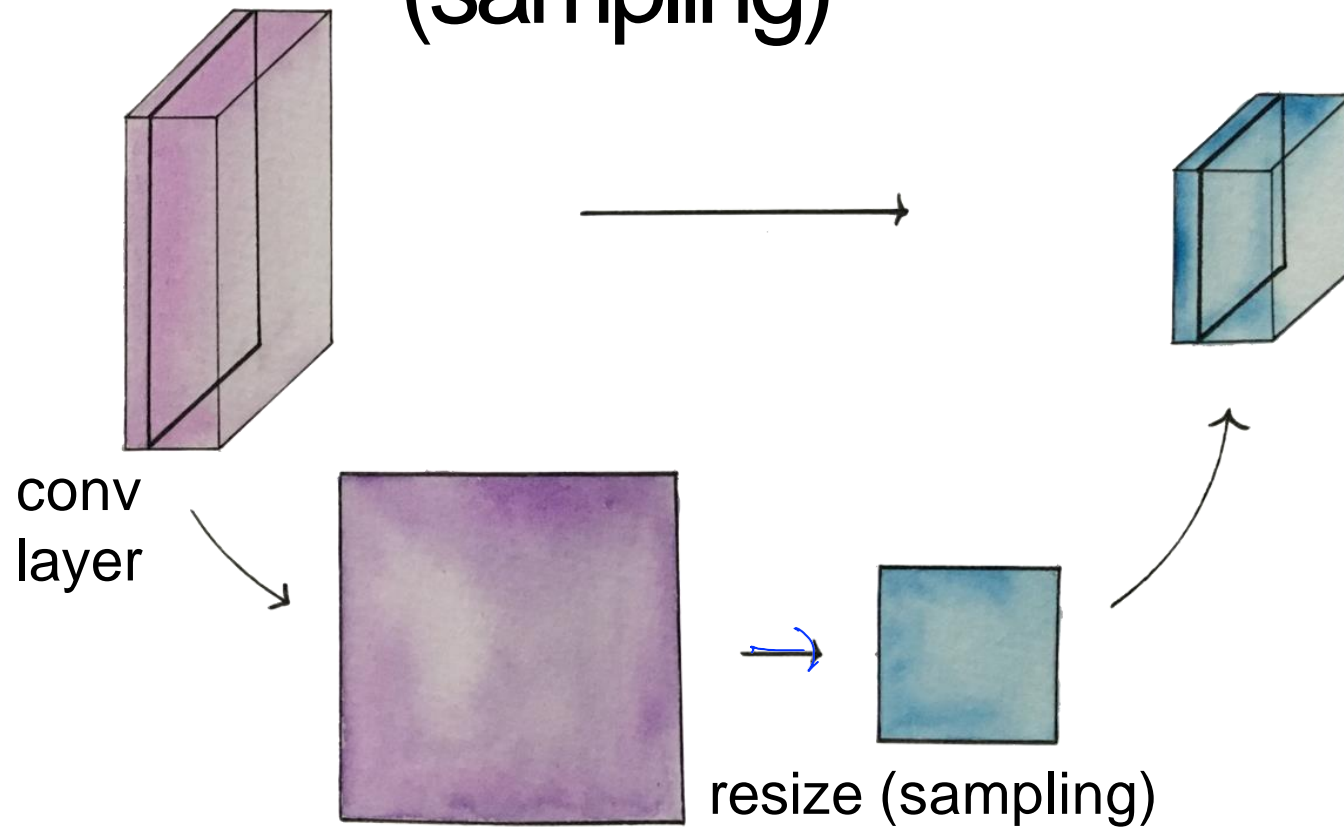
preview:



Pooling layer (sampling)



Pooling layer (sampling)



MAX POOLING

Single depth slice

1	1	2	4
5	6	7	8
3	2	1	0
1	2	3	4

max pool with 2x2 filters
and stride 2



2x2

6	8
3	4

MAX POOLING

Single depth slice

x

1	1	2	4
5	6	7	8
3	2	1	0
1	2	3	4

y

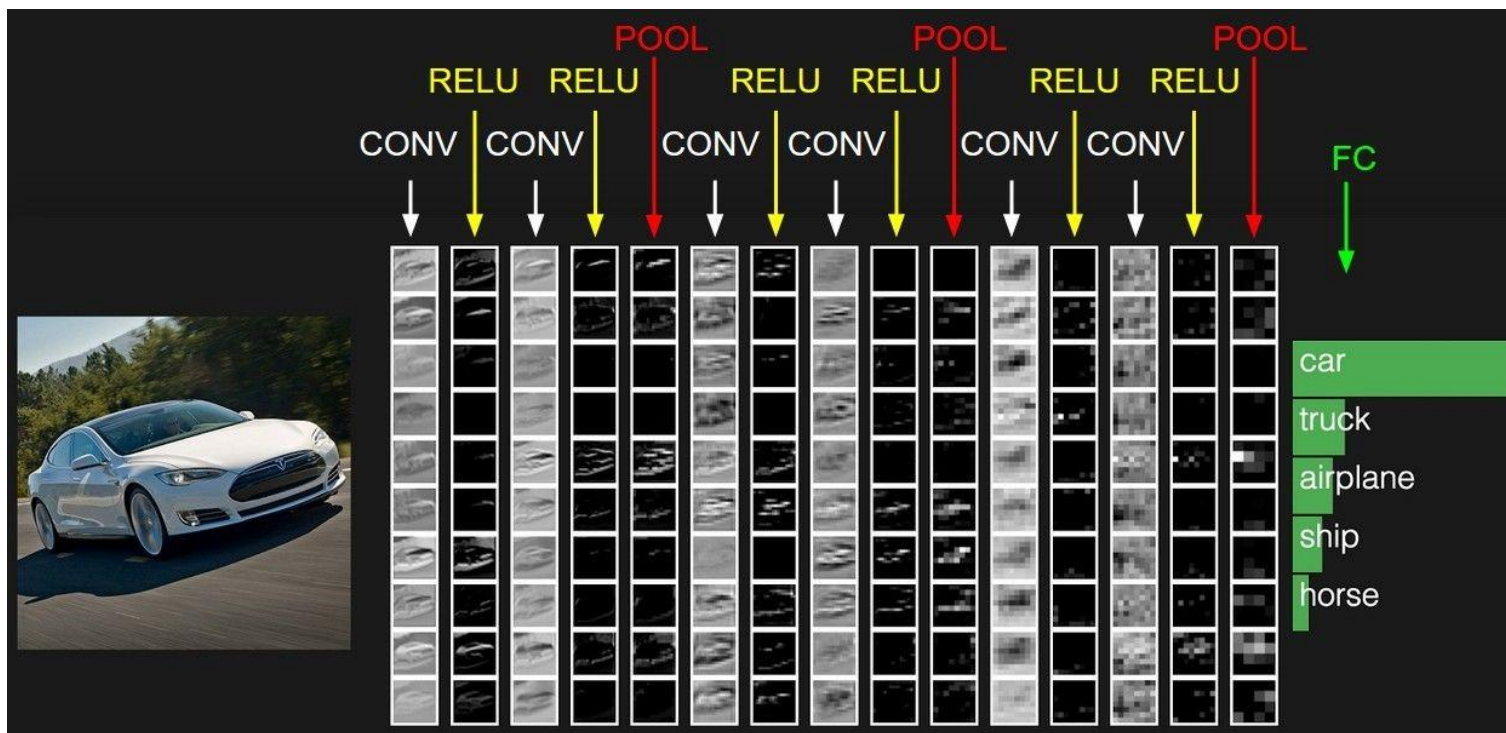
max pool with 2x2 filters
and stride 2



6	8
3	4

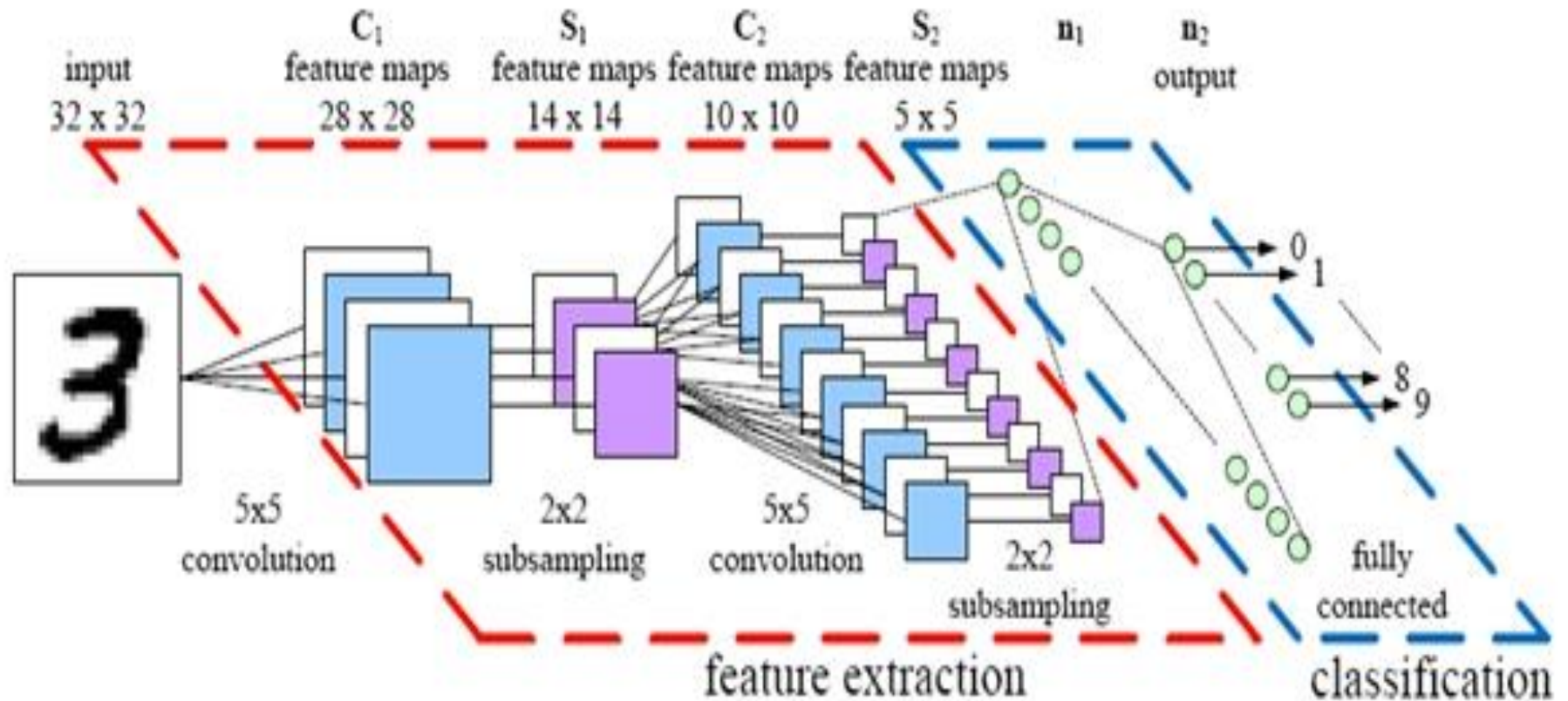
Fully Connected Layer (FC layer)

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

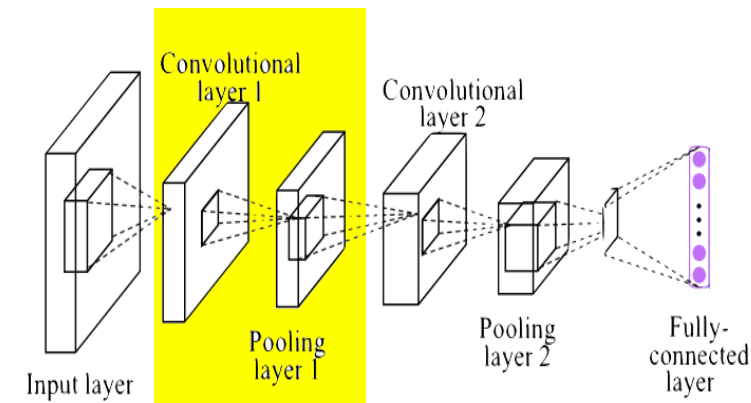


<https://cs.stanford.edu/people/karpathy/convnetjs/demo/cifar10.html>

CNN



Conv layer 1



```
# input placeholders
```

```
X = tf.placeholder(tf.float32, [None, 784])
```

```
X_img = tf.reshape(X, [-1, 28, 28, 1]) # img 28x28x1 (black/white)
```

```
Y = tf.placeholder(tf.float32, [None, 10])
```

```
# L1 ImgIn shape=(?, 28, 28, 1)
```

```
W1 = tf.Variable(tf.random_normal([3, 3, 1, 32], stddev=0.01))
```

```
# Conv -> (?, 28, 28, 32)
```

```
# Pool -> (?, 14, 14, 32)
```

```
L1 = tf.nn.conv2d(X_img, W1, strides=[1, 1, 1, 1], padding='SAME')
```

```
L1 = tf.nn.relu(L1)
```

```
L1 = tf.nn.max_pool(L1, ksize=[1, 2, 2, 1],  
                    strides=[1, 2, 2, 1], padding='SAME')
```

```
...
```

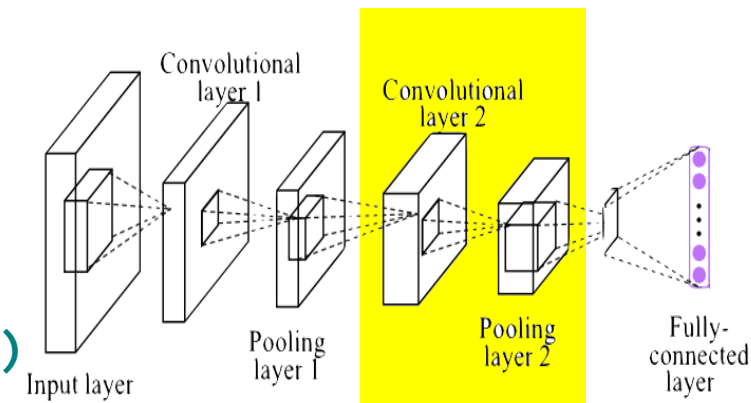
```
Tensor("Conv2D:0", shape=(?, 28, 28, 32), dtype=float32)
```

```
Tensor("Relu:0", shape=(?, 28, 28, 32), dtype=float32)
```

```
Tensor("MaxPool:0", shape=(?, 14, 14, 32), dtype=float32)
```

```
...
```

Conv layer 2



```
...
```

```
Tensor("Conv2D:0", shape=(?, 28, 28, 32), dtype=float32)
Tensor("Relu:0", shape=(?, 28, 28, 32), dtype=float32)
Tensor("MaxPool:0", shape=(?, 14, 14, 32), dtype=float32)
...
```

```
# L2 ImgIn shape=(?, 14, 14, 32)
W2 = tf.Variable(tf.random_normal([3, 3, 32, 64], stddev=0.01))
# Conv ->(?, 14, 14, 64)
# Pool ->(?, 7, 7, 64)
L2 = tf.nn.conv2d(L1, W2, strides=[1, 1, 1, 1], padding='SAME')
L2 = tf.nn.relu(L2)
L2 = tf.nn.max_pool(L2, ksize=[1, 2, 2, 1], strides=[1, 2, 2, 1], padding='SAME')
L2 = tf.reshape(L2, [-1, 7 * 7 * 64])
...
```

```
Tensor("Conv2D_1:0", shape=(?, 14, 14, 64), dtype=float32)
Tensor("Relu_1:0", shape=(?, 14, 14, 64), dtype=float32)
Tensor("MaxPool_1:0", shape=(?, 7, 7, 64), dtype=float32)
Tensor("Reshape_1:0", shape=(?, 3136), dtype=float32)
```

Fully Connected (FC, Dense) layer

```
...
```

```
Tensor("Conv2D_1:0", shape=(?, 14, 14, 64), dtype=float32)  
Tensor("Relu_1:0", shape=(?, 14, 14, 64), dtype=float32)  
Tensor("MaxPool_1:0", shape=(?, 7, 7, 64), dtype=float32)  
Tensor("Reshape_1:0", shape=(?, 3136), dtype=float32)  
...
```

```
L2 = tf.reshape(L2, [-1, 7 * 7 * 64])
```

```
# Final FC 7x7x64 inputs -> 10 outputs
```

```
W3 = tf.get_variable("W3", shape=[7 * 7 * 64, 10],  
initializer=tf.contrib.layers.xavier_initializer())
```

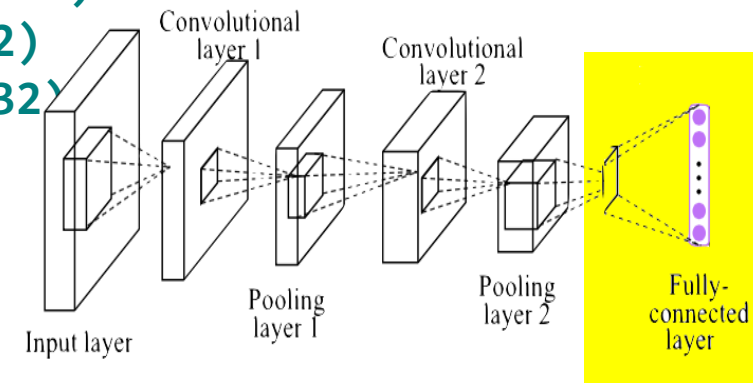
```
b = tf.Variable(tf.random_normal([10]))
```

```
hypothesis = tf.matmul(L2, W3) + b
```

```
# define cost/loss & optimizer
```

```
cost = tf.reduce_mean(tf.nn.softmax_cross_entropy_with_logits(logits=hypothesis,  
labels=Y))
```

```
optimizer = tf.train.AdamOptimizer(learning_rate=learning_rate).minimize(cost)
```



Training and Evaluation

```
# initialize
sess = tf.Session()
sess.run(tf.global_variables_initializer())

# train my model
print('Learning started. It takes sometime.')
for epoch in range(training_epochs):
    avg_cost = 0
    total_batch = int(mnist.train.num_examples / batch_size)
    for i in range(total_batch):
        batch_xs, batch_ys = mnist.train.next_batch(batch_size)
        feed_dict = {X: batch_xs, Y: batch_ys}
        c, _, = sess.run([cost, optimizer], feed_dict=feed_dict)
        avg_cost += c / total_batch
    print('Epoch:', '%04d' % (epoch + 1), 'cost =', '{:.9f}'.format(avg_cost))

print('Learning Finished!')

# Test model and check accuracy
correct_prediction = tf.equal(tf.argmax(hypothesis, 1), tf.argmax(Y, 1))
accuracy = tf.reduce_mean(tf.cast(correct_prediction, tf.float32))
print('Accuracy:', sess.run(accuracy, feed_dict={X: mnist.test.images, Y:
mnist.test.labels}))
```

Epoch: 0001 cost =
0.340291267

Epoch: 0002 cost =
0.090731326

Epoch: 0003 cost =
0.064477619

Epoch: 0004 cost =
0.050683064

...

Epoch: 0011 cost =
0.017758641

Epoch: 0012 cost =
0.014156652

Epoch: 0013 cost =
0.012397016

Epoch: 0014 cost =
0.010693789

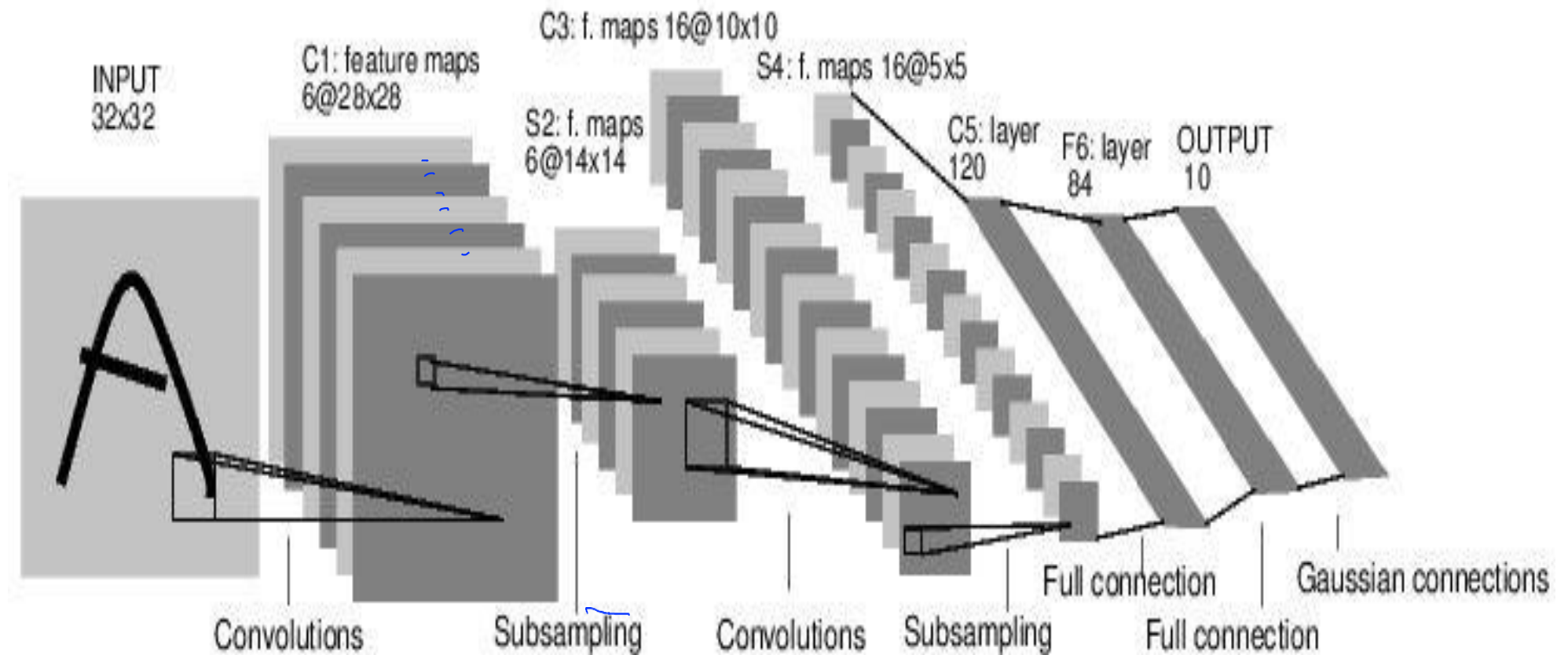
Epoch: 0015 cost =
0.009469977

Learning Finished!

Accuracy: **0.9885**

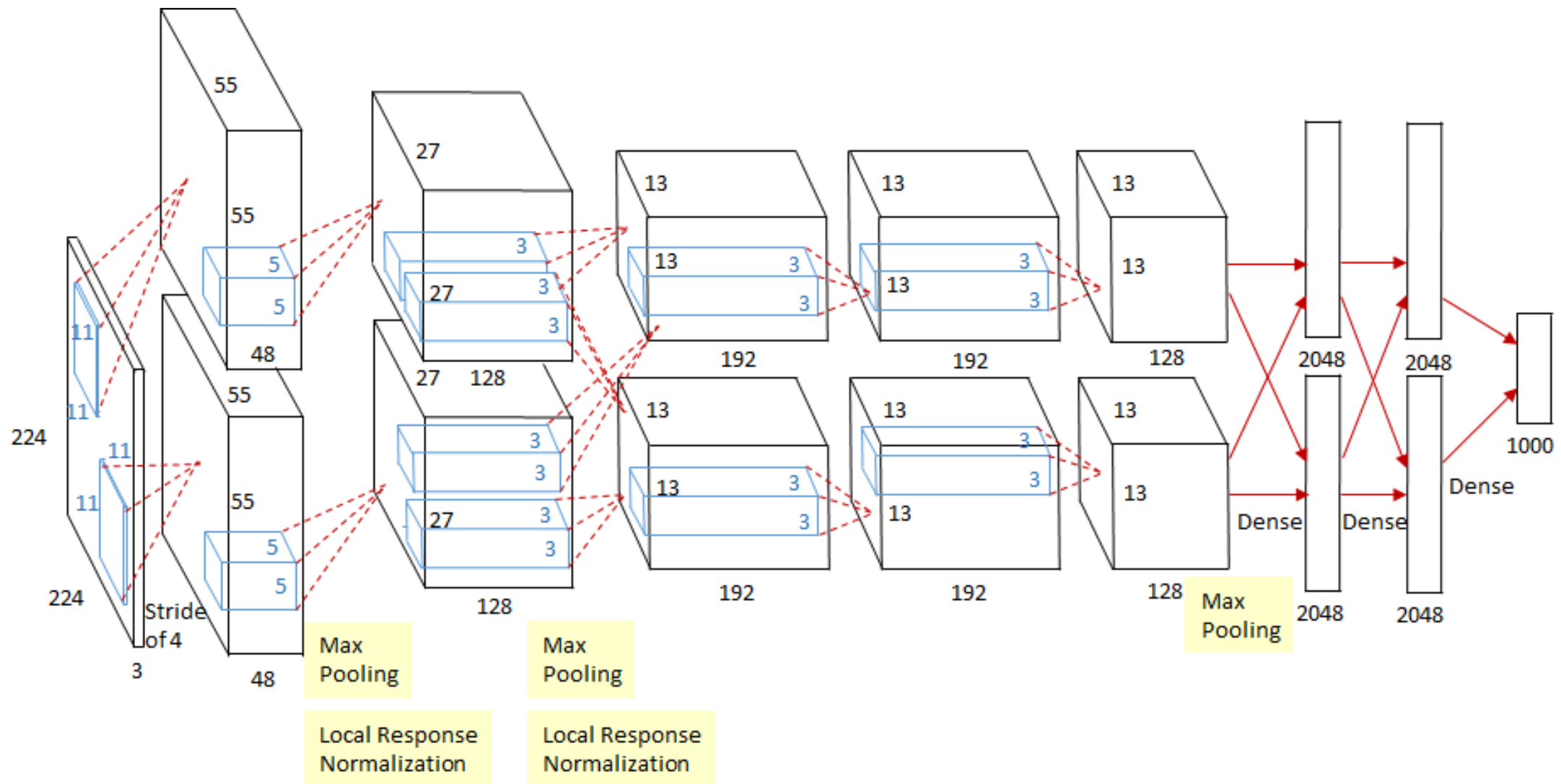
Case Study: LeNet-5 (60K param)

[LeCun et al., 1998]



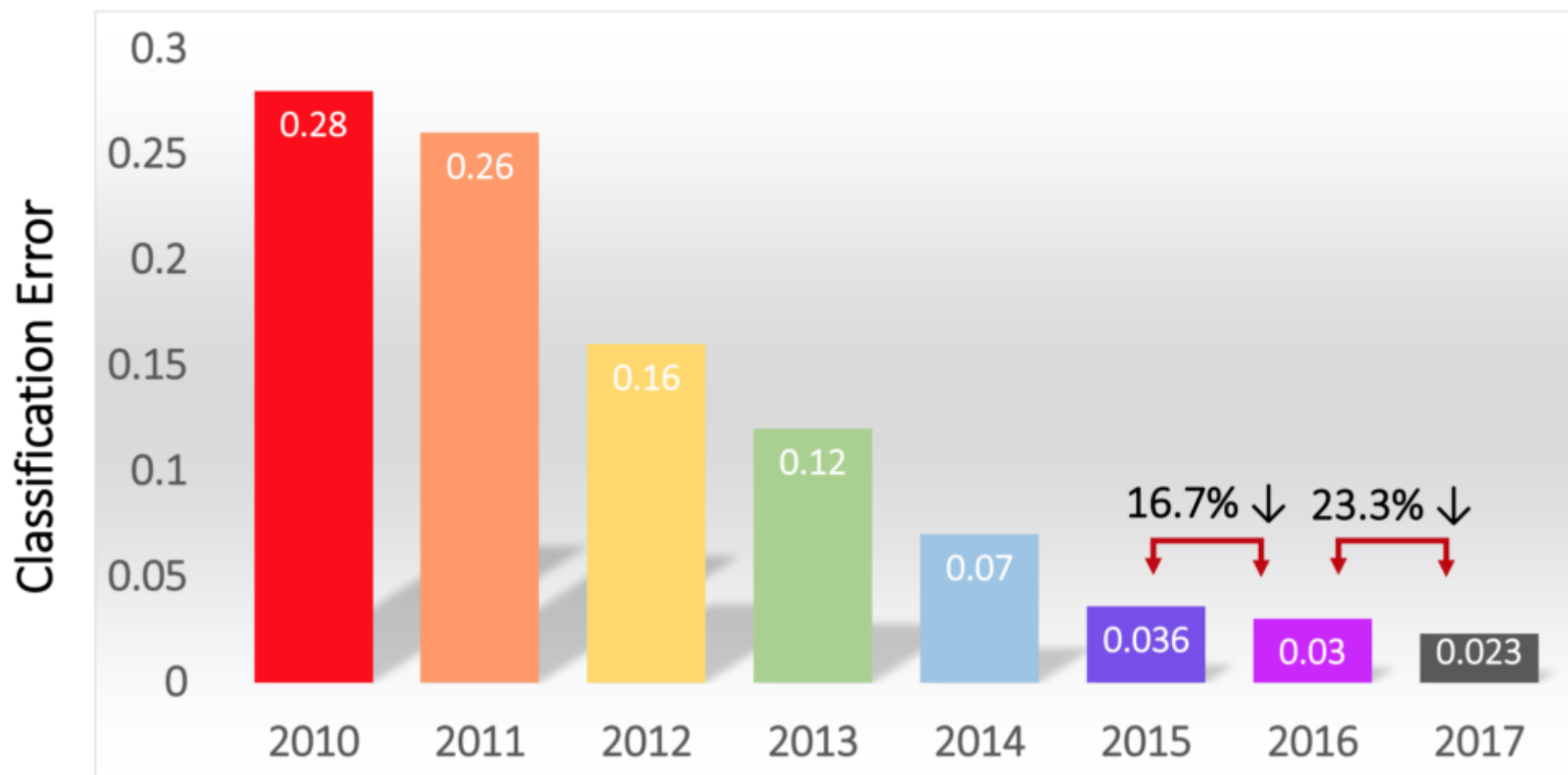
Conv filters were 5x5, applied at stride 1
Subsampling (Pooling) layers were 2x2 applied at stride 2
i.e. architecture is [CONV-POOL-CONV-POOL-CONV-FC]

AlexNet

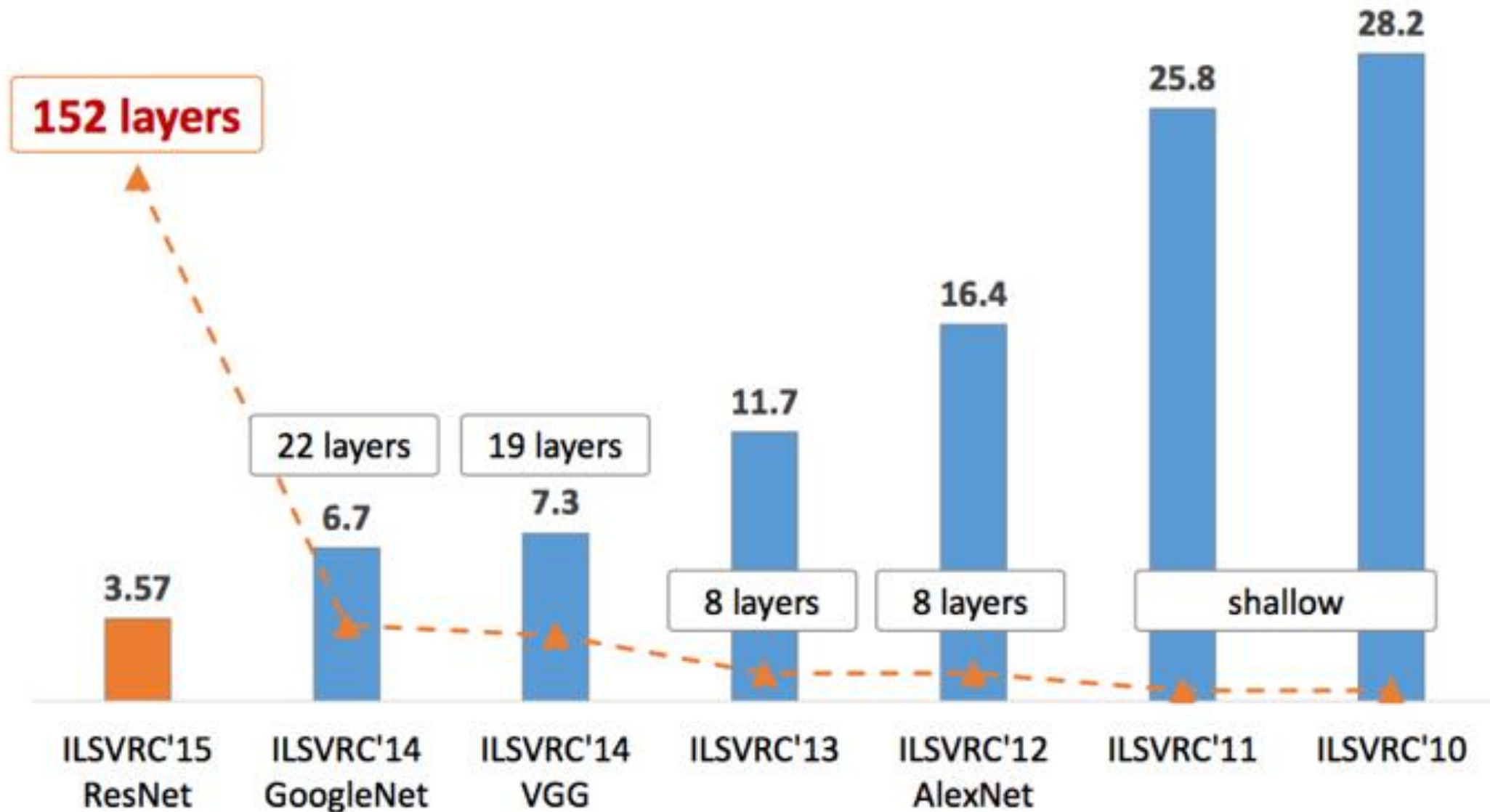




Classification Results (CLS)



CNN Architectures: LeNet, AlexNet, VGG, GoogLeNet, ResNet



Deep Visualization Toolbox

<https://github.com/yosinski/deep-visualization-toolbox>

<https://www.youtube.com/watch?v=AgkfIQ4IGaM>

