# Tensorflow CNN turorial

2017/03/10

# Lenet-5

[LeCun et al., 1998]



INPUT 32x32 — C1: feature maps 6@28x28 — S2: f. maps 6@14x14 — C3: f. maps 16@10x10 — S4: f. maps 16@5x5 — C5: layer 120 — F6: layer 84 — OUTPUT 10

Convolutions — Subsampling — Convolutions — Subsampling — Full connection — Full connection — Gaussian connections

# Today's example



input     conv1     pool1     conv2     pool2   hidden4   output

Convolution     Subsample     Convolution     Subsample
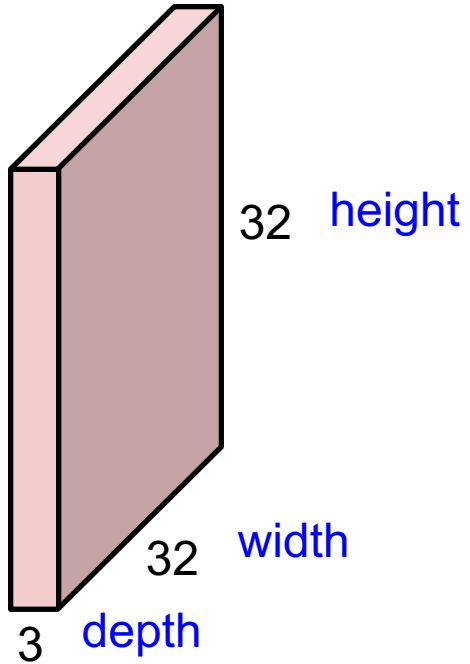
Full Connection

**Full Connection**

**The slides are from
1. "Lecture 13: Neural networks for machine vision, Dr. Richard E. Turner"
2. Lecture 7 & 12 in Stanford CS231n**

# Convolution Layer
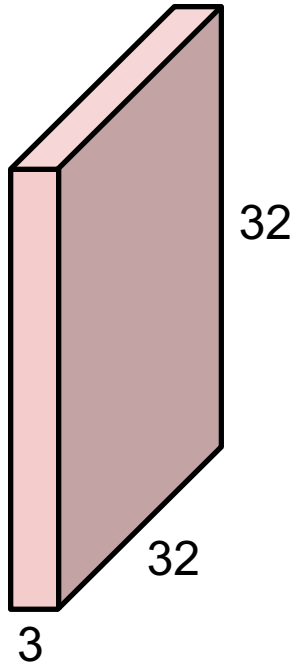
32x32x3 image

32 height

32 width

3 depth

# Convolution Layer

**32x32x3 image**

32

32

3

**5x5x3 filter**

**Convolve** the filter with the image
i.e. "slide over the image spatially,
computing dot products"
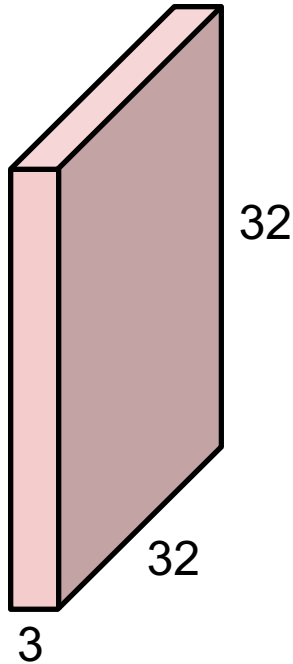
# Convolution Layer

32x32x3 image

32

32

3

5x5x3 filter

**Convolve** the filter with the image i.e. "slide over the image spatially, computing dot products"
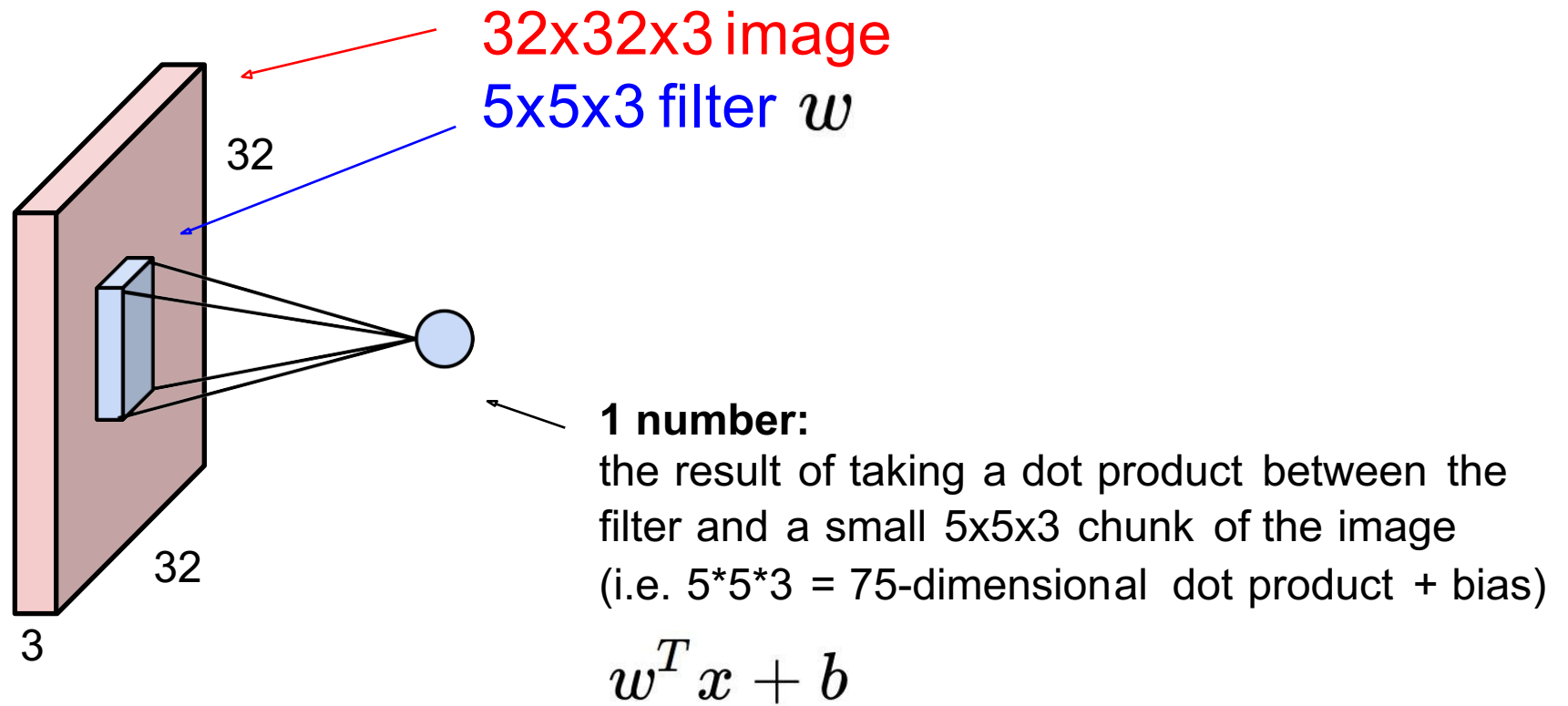
# Convolution Layer



32x32x3 image
5x5x3 filter $w$

**1 number:**
the result of taking a dot product between the filter and a small 5x5x3 chunk of the image (i.e. 5*5*3 = 75-dimensional dot product + bias)
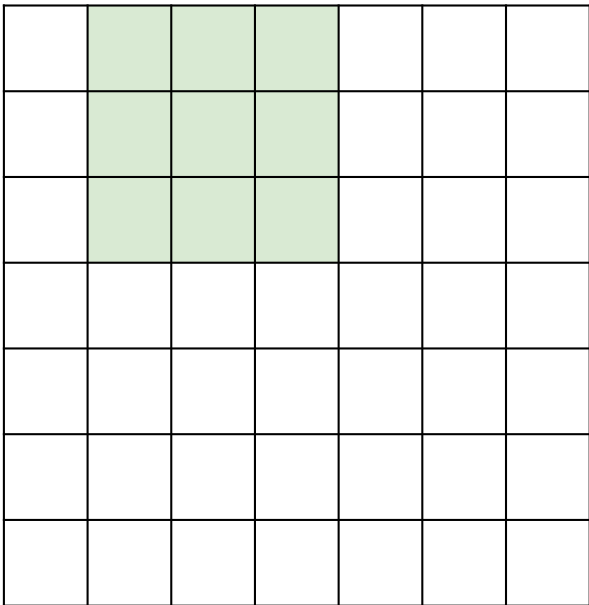
$$w^T x + b$$

**7**

**7**
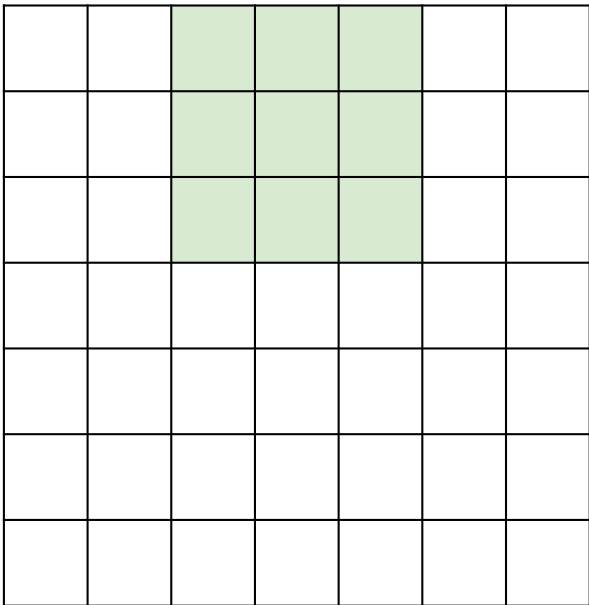
7x7 input (spatially)
assume 3x3 filter

7

7x7 input (spatially)
assume 3x3 filter

7

7x7 input (spatially)
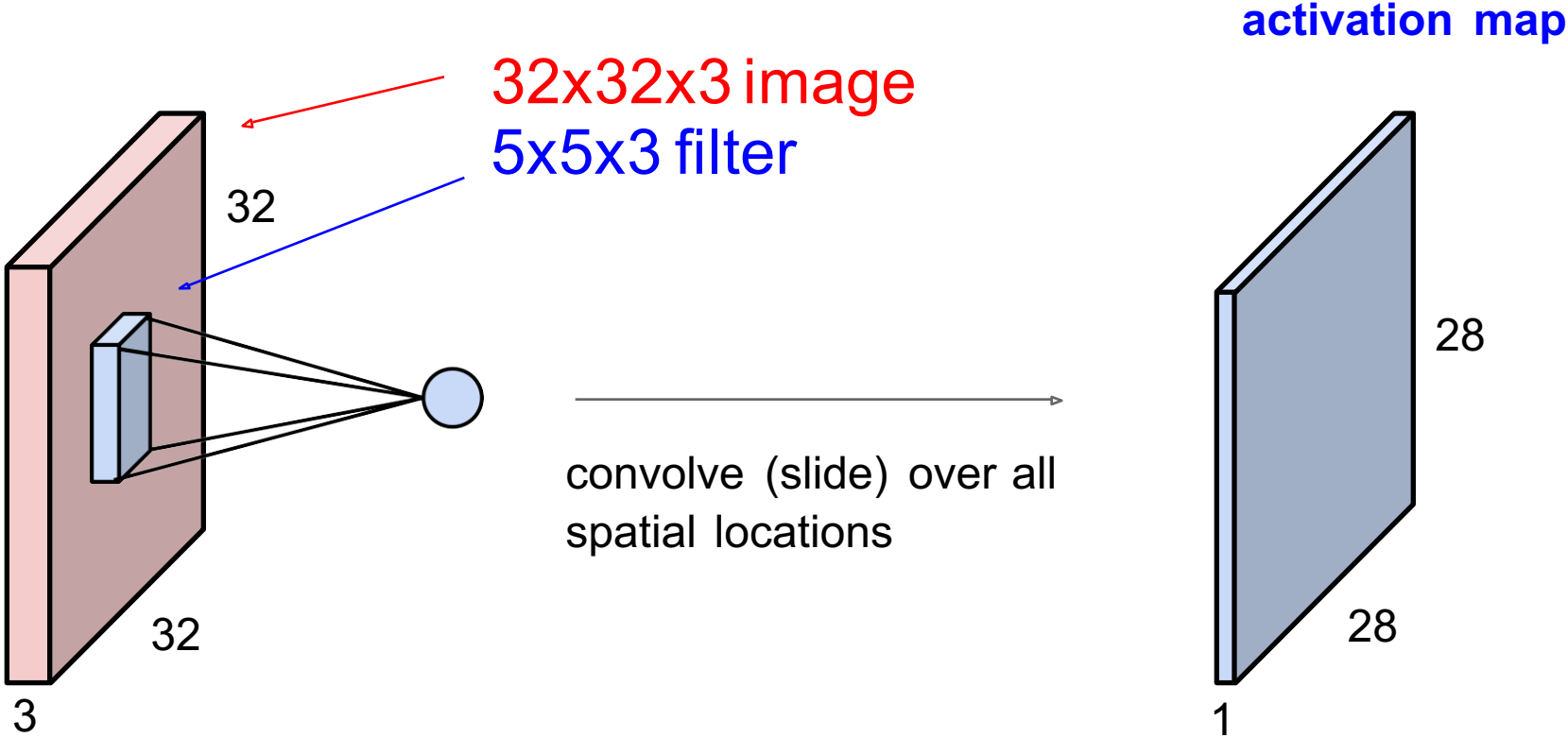assume 3x3 filter

7

7x7 input (spatially)
assume 3x3 filter

7
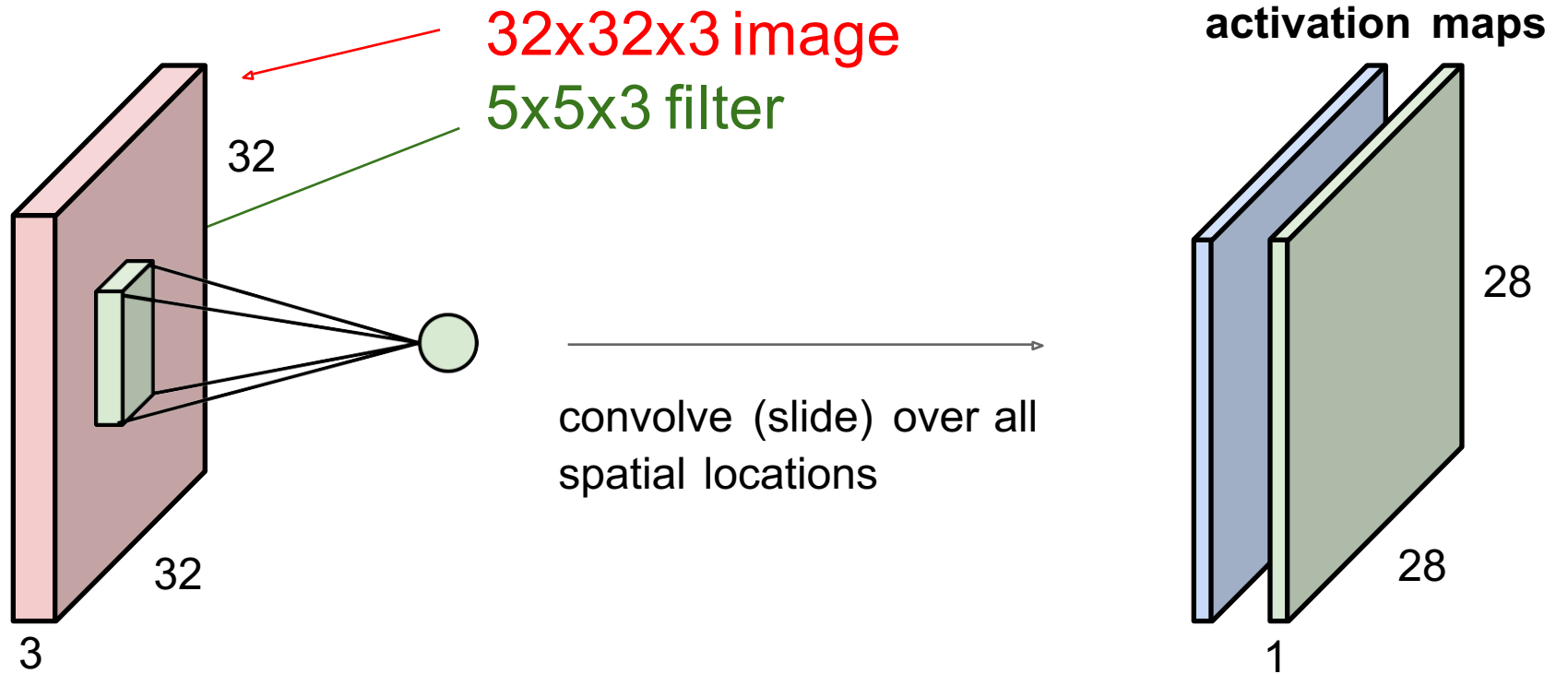
**7**



**7**

7x7 input (spatially)
assume 3x3 filter

**=> 5x5 output**

# Convolution Layer



**32x32x3 image**

**5x5x3 filter**

32

32

3

convolve (slide) over all spatial locations

**activation map**

28

28

1

# Convolution Layer

consider a second, green filter

32x32x3 image
5x5x3 filter

32

32

3

convolve (slide) over all
spatial locations

**activation maps**

28

28

1

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

**activation maps**

32

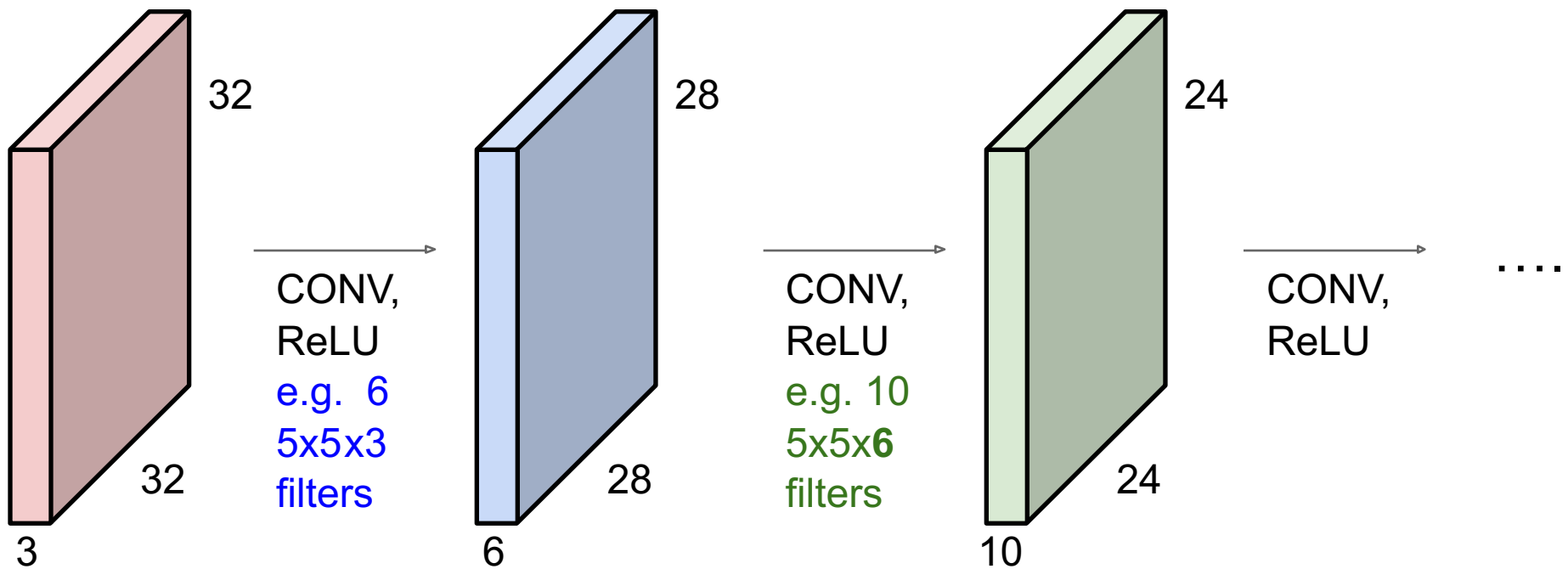32

3

Convolution Layer

28

28

6

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

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



32

32

3

CONV,
ReLU
e.g. 6
5x5x3
filters

28

28

6

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

32

32

3

CONV,
ReLU
e.g. 6
5x5x3
filters

28

28

6

CONV,
ReLU
e.g. 10
5x5x**6**
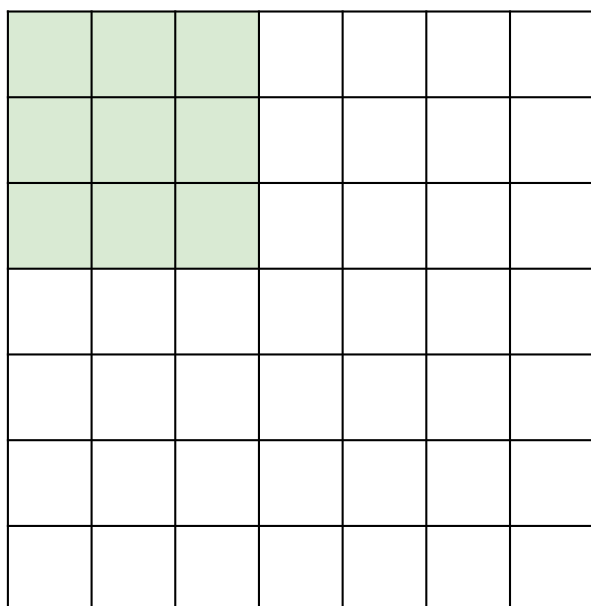filters

24

24

10

CONV,
ReLU

....

32x32 input convolved repeatedly with 5x5 filters shrinks volumes spatially! (32 -> 28 -> 24 ...). Shrinking too fast is not good, doesn't work well.

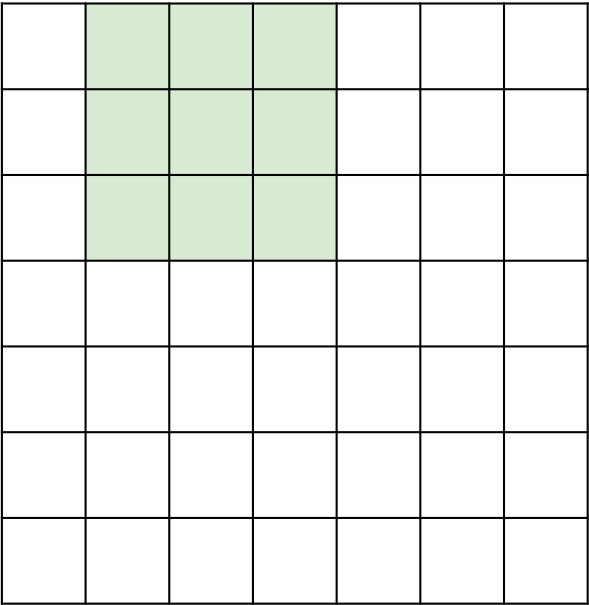A closer look at spatial dimensions:

7



7

7x7 input (spatially)
assume 3x3 filter

A closer look at spatial dimensions:

7



7

7x7 input (spatially)
assume 3x3 filter
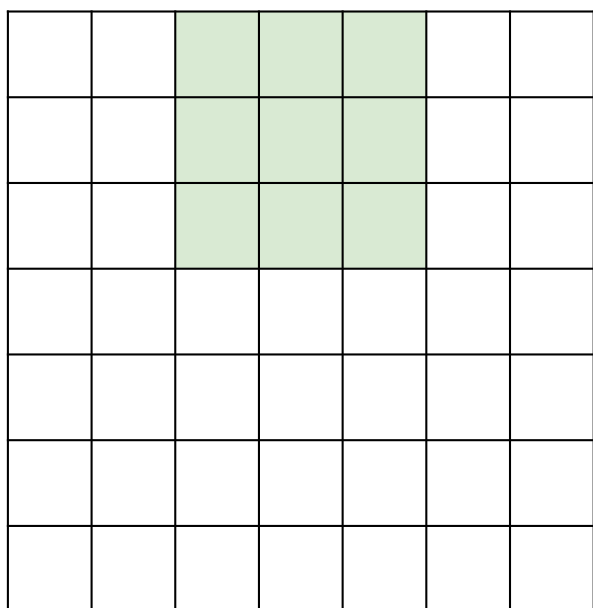
A closer look at spatial dimensions:

7



7
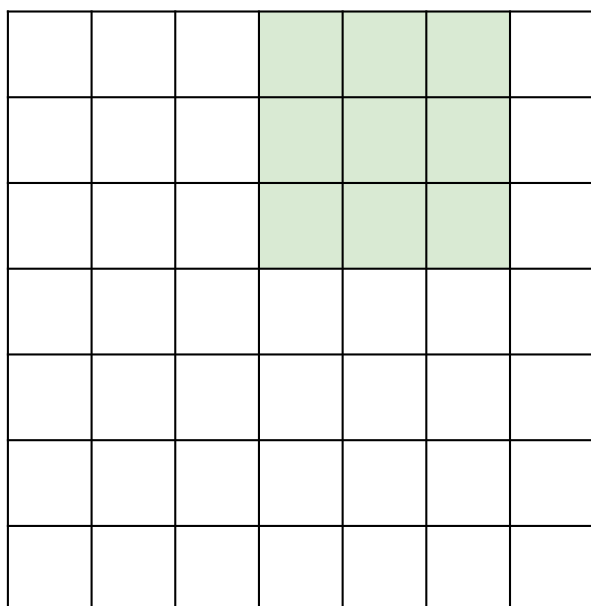
7x7 input (spatially)
assume 3x3 filter

A closer look at spatial dimensions:

7



7x7 input (spatially)
assume 3x3 filter

7

A closer look at spatial dimensions:

7



7

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:

7



7

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!**

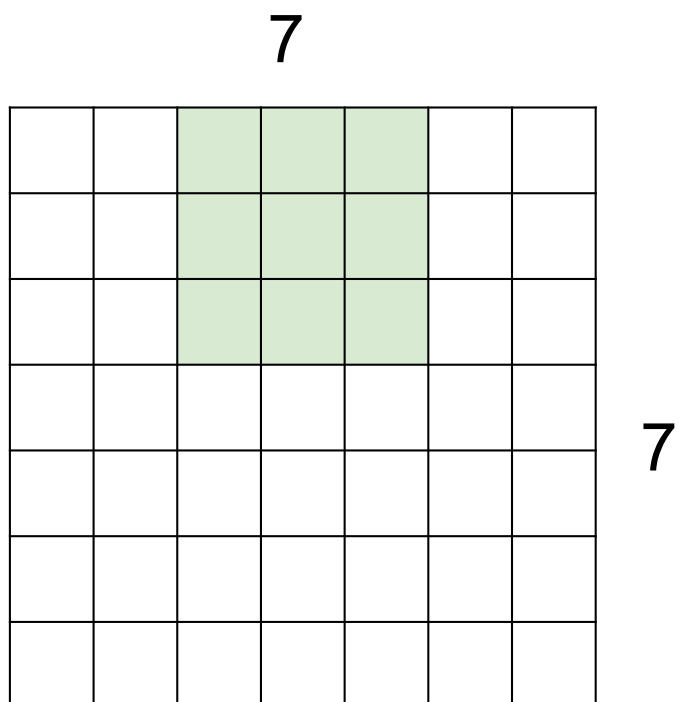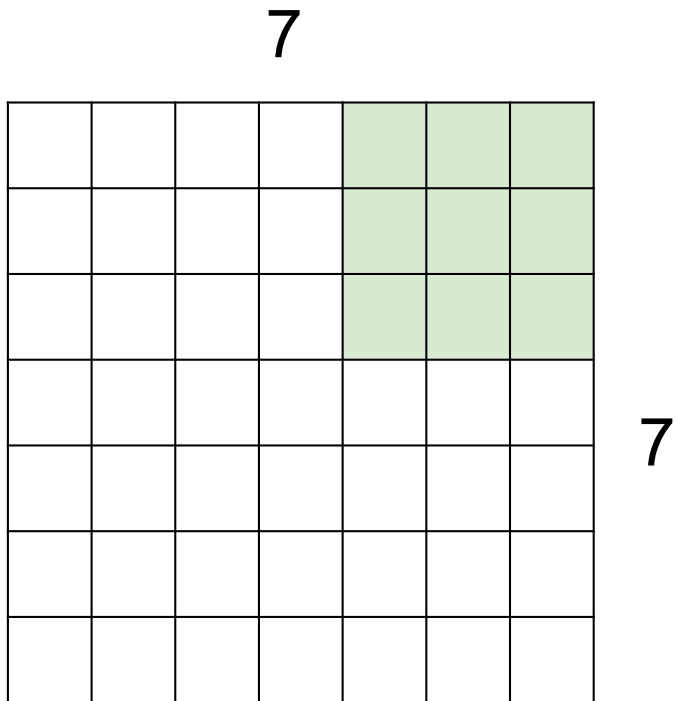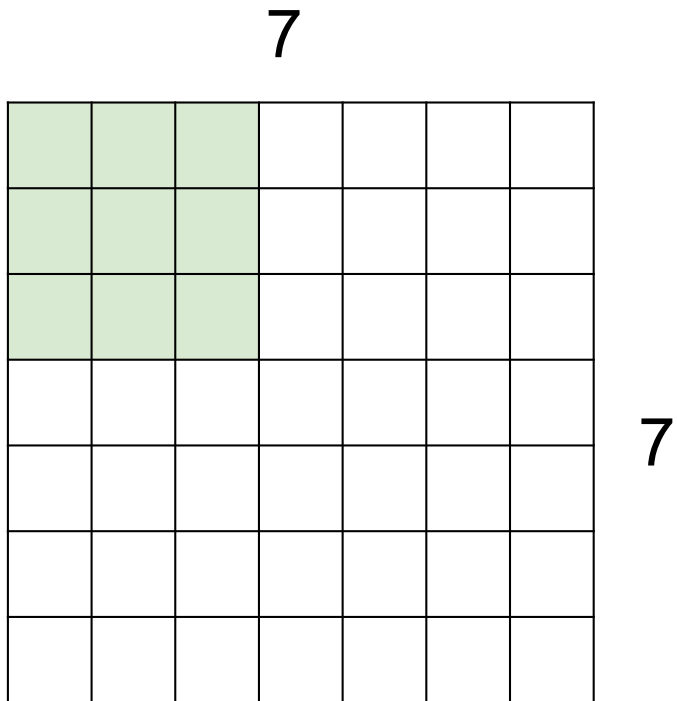A closer look at spatial dimensions:



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

A closer look at spatial dimensions:



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

**doesn't fit!**
cannot apply 3x3 filter on
7x7 input with stride 3.

# 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



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 FxF, and zero-padding with (F-1)/2. (will preserve size spatially)

e.g. F = 3 => zero pad with 1
      F = 5 => zero pad with 2
      F = 7 => zero pad with 3

# Pooling layer

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

# MAX POOLING

Single depth slice

# Tensorflow implementation

- Weight Initialization

- Convolution and Pooling

- Convolution layer

- Fully connected layer

- Readout Layer


- Reference and image source: https://www.tensorflow.org/get_started/mnist/pros (See section 'Build a Multilayer Convolutional Network')

# Input (placeholder)

```
x = tf.placeholder(tf.float32, shape=[None, input_size])
y = tf.placeholder(tf.float32, shape=[None, classes_num])
```

x is placeholder for input image.
y is label with one-hot representation, so second dimension of y is equal to number of classes.

None indicates that the first dimension, corresponding to the batch size, which can be any size.

tf.placeholder

# Weight Initialization

```python
def weight_variable(shape):
    initial = tf.truncated_normal(shape, stddev=0.1)
    return tf.Variable(initial)

def bias_variable(shape):
    initial = tf.constant(0.1, shape=shape)
    return tf.Variable(initial)
```

tf.truncated_normal

These variable will be initialized when user run 'tf.global_variables_initializer'.
Now they are just nodes in a graph without any value.

# Convolution and Pooling

```python
def conv2d(x, W):
    return tf.nn.conv2d(x, W, strides=[1, 1, 1, 1], padding='SAME')

def max_pool_2x2(x):
    return tf.nn.max_pool(x, ksize=[1, 2, 2, 1],
                          strides=[1, 2, 2, 1], padding='SAME')
```

Strides is 4-d, following NHWC format.
(Num_samples x Height x Width x Channels)

Recall strides and padding.
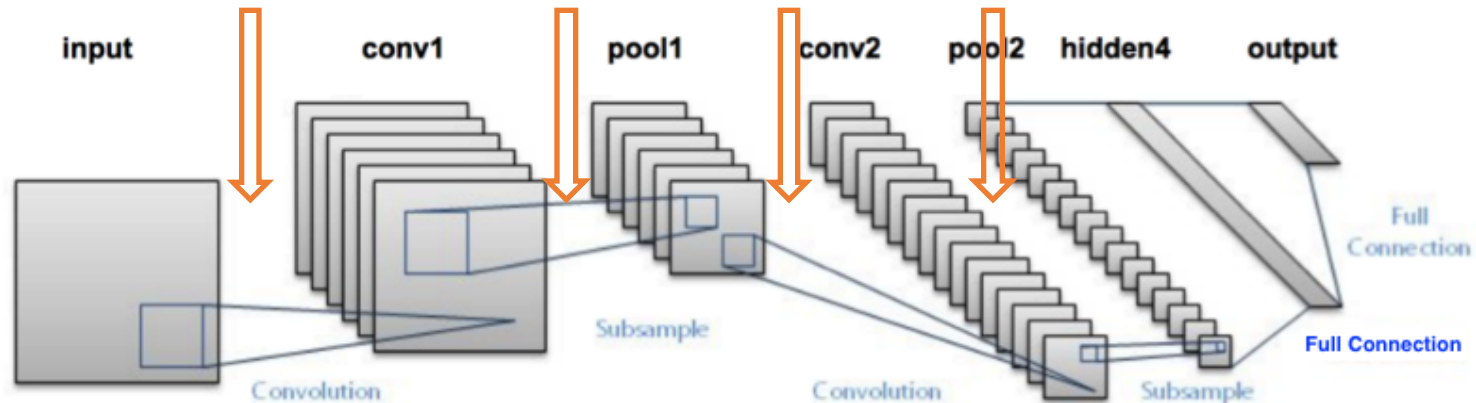padding = 'SAME' means apply padding to keep output size as same as input size.

Conv2d pads with zeros and max_pool pads with –inf.

tf.nn.conv2d
tf.nn.max_pool

# Convolution layer

```
W_conv1 = weight_variable([5, 5, 1, 32])
b_conv1 = bias_variable([32])
x_image = tf.reshape(x, [-1,28,28,1])
h_conv1 = tf.nn.relu(conv2d(x_image, W_conv1) + b_conv1)
h_pool1 = max_pool_2x2(h_conv1)
W_conv2 = weight_variable([5, 5, 32, 64])
b_conv2 = bias_variable([64])

h_conv2 = tf.nn.relu(conv2d(h_pool1, W_conv2) + b_conv2)
h_pool2 = max_pool_2x2(h_conv2)
```



tf.reshape

# Convolution layer

```
W_conv1 = weight_variable([5, 5, 1, 32])
b_conv1 = bias_variable([32])
x_image = tf.reshape(x, [-1,28,28,1])
h_conv1 = tf.nn.relu(conv2d(x_image, W_conv1) + b_conv1)
h_pool1 = max_pool_2x2(h_conv1)
W_conv2 = weight_variable([5, 5, 32, 64])
b_conv2 = bias_variable([64])

h_conv2 = tf.nn.relu(conv2d(h_pool1, W_conv2) + b_conv2)
h_pool2 = max_pool_2x2(h_conv2)
```

See how the code creates a model by wrapping layers.
Be care of shape of each layer.
-1 means match the size of that dimension is computed
so that the total size remains constant.

tf.reshape

# Reshape

For example:

tensor 't' is [[1, 2], [3, 4], [5, 6], [7, 8]] , so t has shape [4, 2]

(1) reshape(t, [2,4]) ➔ [[1, 2, 3, 4], [5, 6, 7, 8]]
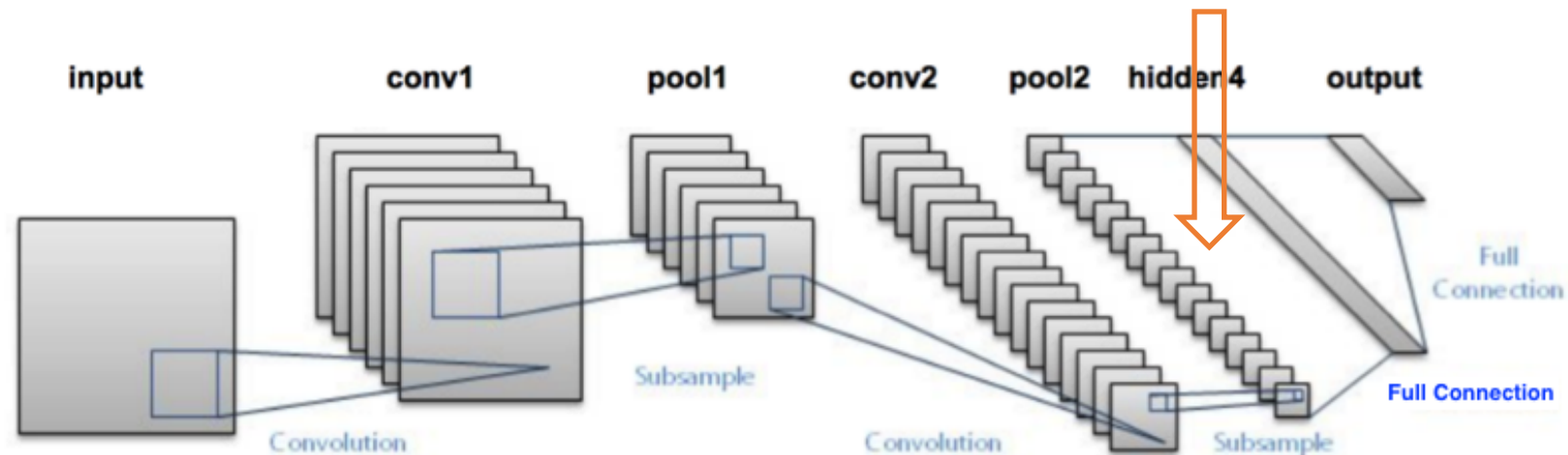(2) reshape(t, [-1, 4]) ➔ [[1, 2, 3, 4], [5, 6, 7, 8]]

-1 would be computed and becomes '2'

tf.reshape

# Fully connected layer

```
W_fc1 = weight_variable([7 * 7 * 64, 1024])
b_fc1 = bias_variable([1024])

h_pool2_flat = tf.reshape(h_pool2, [-1, 7*7*64])
h_fc1 = tf.nn.relu(tf.matmul(h_pool2_flat, W_fc1) + b_fc1)
```
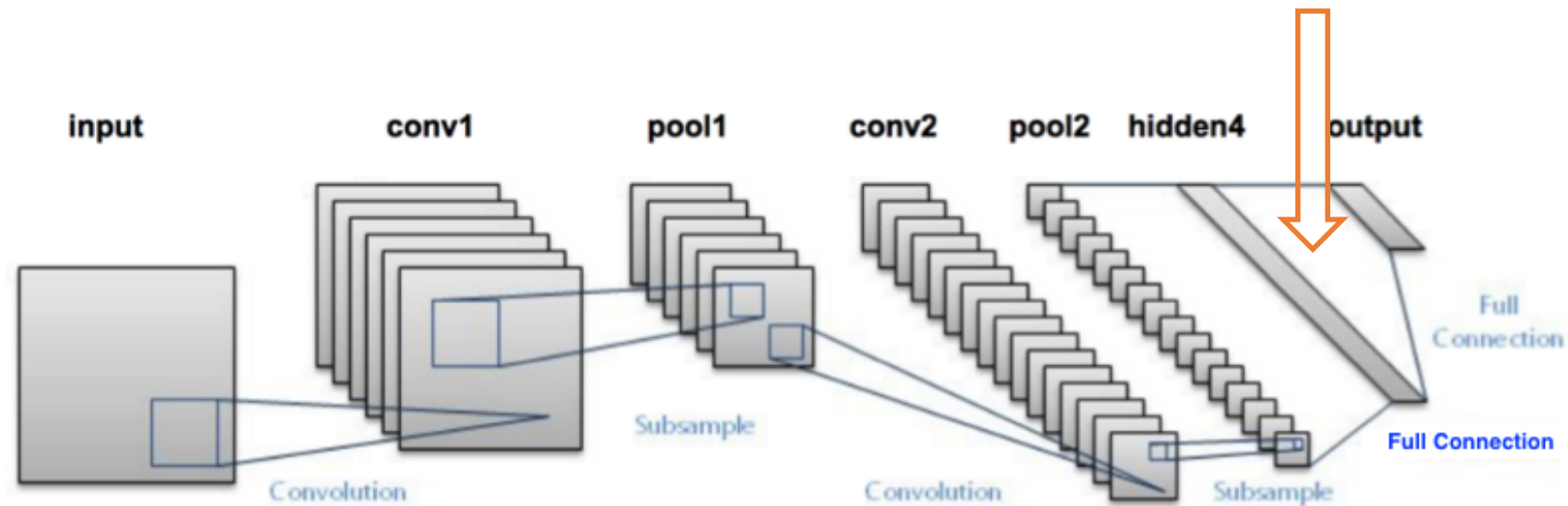
Flatten all the maps and connect them with fully connected layer.
Again, be care of shape.

# Readout Layer

```
W_fc2 = weight_variable([1024, 10])
b_fc2 = bias_variable([10])

y_conv = tf.matmul(h_fc1_drop, W_fc2) + b_fc2
```

Use a layer to match output size.
Done!

# Training and Evaluation (optional)

```python
cross_entropy = tf.reduce_mean(tf.nn.softmax_cross_entropy_with_logits(y_conv, y_))
train_step = tf.train.AdamOptimizer(1e-4).minimize(cross_entropy)
correct_prediction = tf.equal(tf.argmax(y_conv,1), tf.argmax(y_,1))
accuracy = tf.reduce_mean(tf.cast(correct_prediction, tf.float32))
sess.run(tf.global_variables_initializer())
for i in range(20000):
    batch = mnist.train.next_batch(50)
    if i%100 == 0:
        train_accuracy = accuracy.eval(feed_dict={
            x:batch[0], y_: batch[1], keep_prob: 1.0})
        print("step %d, training accuracy %g"%(i, train_accuracy))
    train_step.run(feed_dict={x: batch[0], y_: batch[1], keep_prob: 0.5})

print("test accuracy %g"%accuracy.eval(feed_dict={
    x: mnist.test.images, y_: mnist.test.labels, keep_prob: 1.0}))
```

# Recommendation

- Search for each function, and you'll what's everything going on.