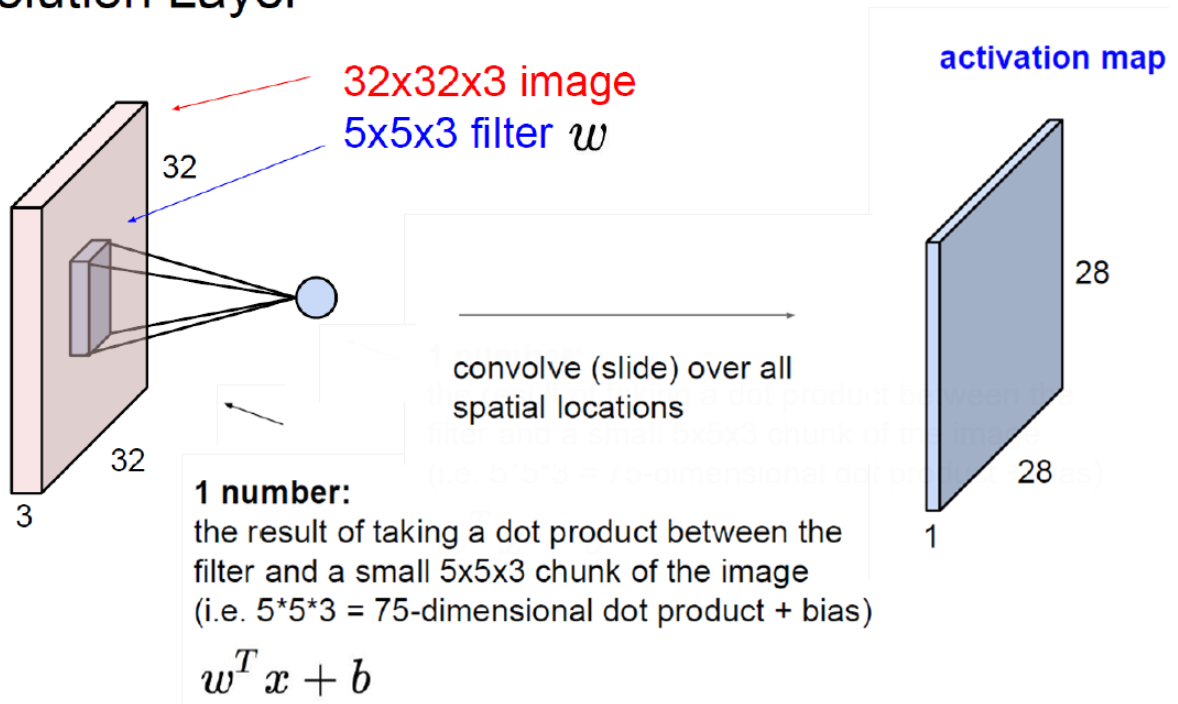


# Convolutional Neural Network

## 卷积

### Convolution Layer

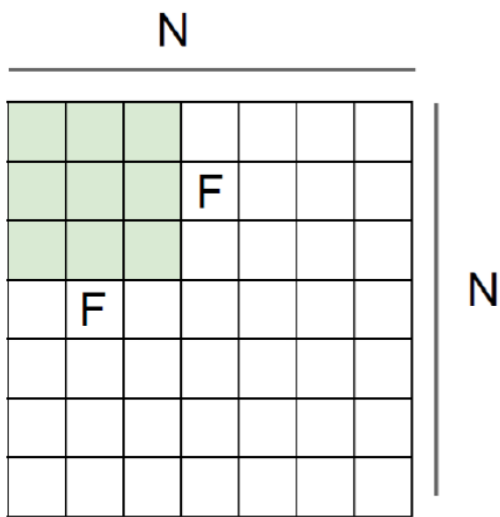


$$f[x,y] * g[x,y] = \sum_{n_1=-\infty}^{\infty} \sum_{n_2=-\infty}^{\infty} f[n_1,n_2] \cdot g[x-n_1,y-n_2]$$

↑  
elementwise multiplication and sum of a filter and the signal (image)

注：上面是卷积的计算公式，n 是卷积核的尺寸，x 和 y 是图片的尺寸。

### stride and padding

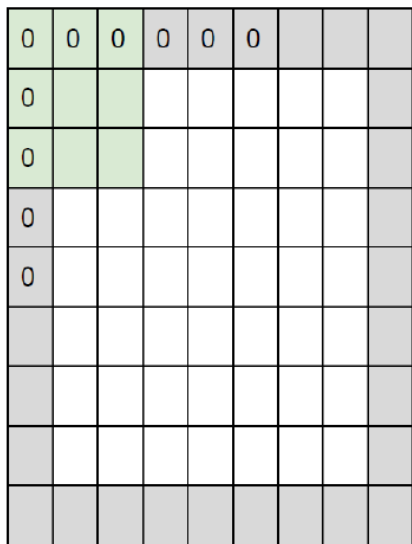


Output size:  
 $(N - F) / \text{stride} + 1$

e.g.  $N = 7, F = 3$ :  
 stride 1  $\Rightarrow (7 - 3) / 1 + 1 = 5$   
 stride 2  $\Rightarrow (7 - 3) / 2 + 1 = 3$   
 stride 3  $\Rightarrow (7 - 3) / 3 + 1 = 2.33 \therefore \backslash$

注：对于 7x7 的图像，不能使用 3x3 的且步长为 3 的卷积核，因为不能整除。

## In practice: Common to zero pad the border



e.g. input 7x7  
**3x3** filter, applied with **stride 1**  
**pad with 1 pixel** border  $\Rightarrow$  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)

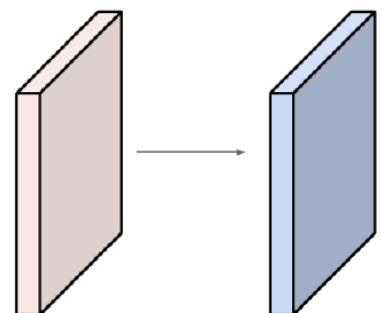
e.g.  $F = 3 \Rightarrow$  zero pad with 1  
 $F = 5 \Rightarrow$  zero pad with 2  
 $F = 7 \Rightarrow$  zero pad with 3

对于 padding 后的尺寸：(图片尺寸 + 2 \* padding - kernel 尺寸) / stride + 1。

Examples time:

Input volume: **32x32x3**

**10** **5x5** filters with stride **1**, pad **2**

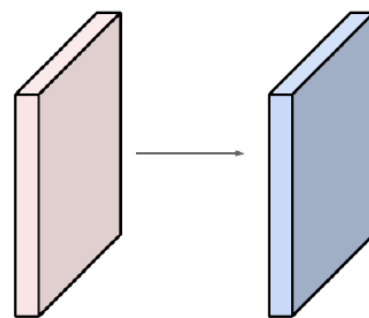


Output volume size:

$(32 + 2 * 2 - 5) / 1 + 1 = 32$  spatially, so

**32x32x10**

Examples time:



Input volume: **32x32x3**

**10** **5x5** filters with stride 1, pad 2

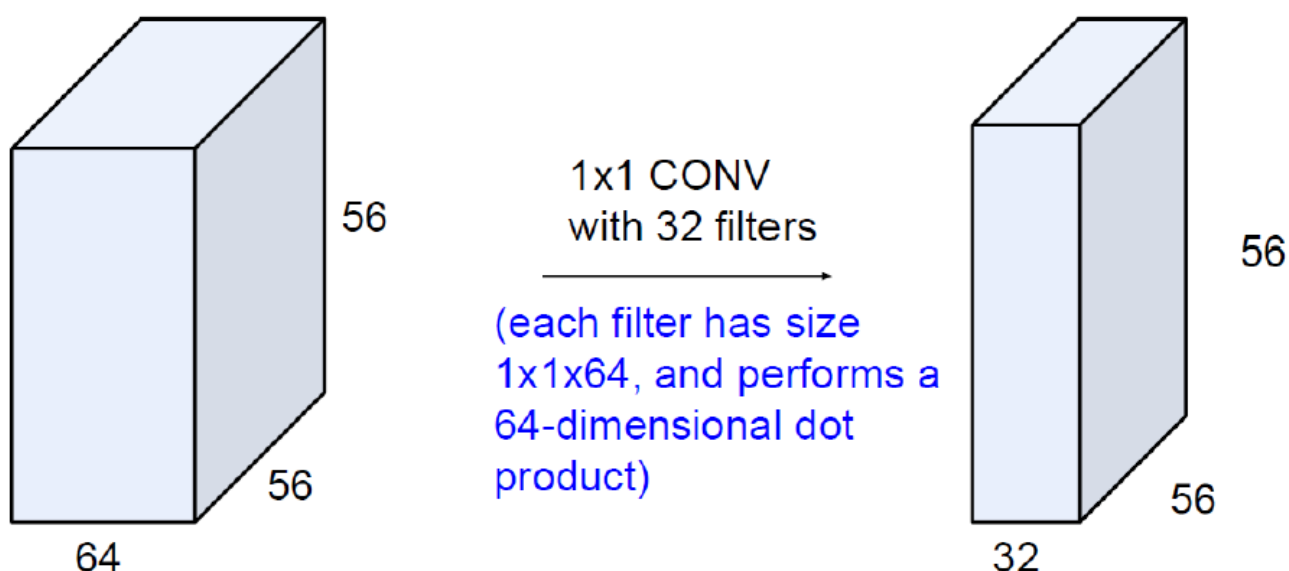
Number of parameters in this layer?

each filter has  $5*5*3 + 1 = 76$  params (+1 for bias)

=>  $76*10 = 760$

注：记得参数中 bias 的那个 1。

(btw, 1x1 convolution layers make perfect sense)

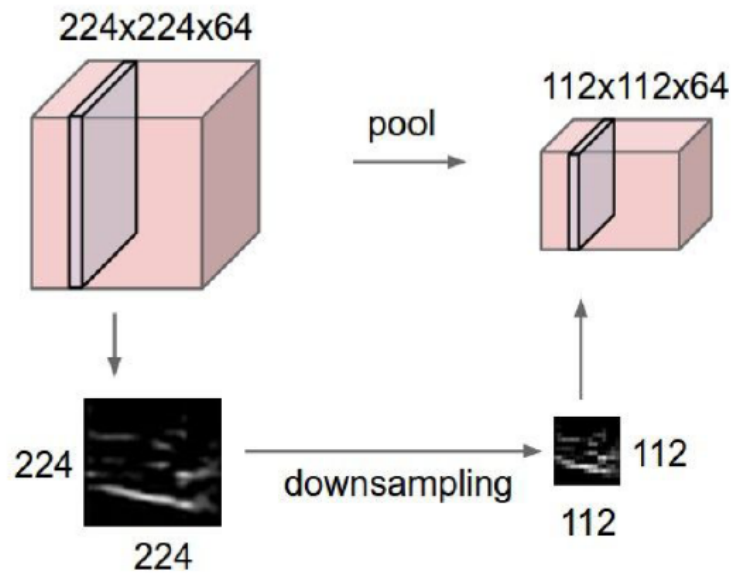


注：1x1 卷积核用于降维。

**pooling**

# Pooling layer

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



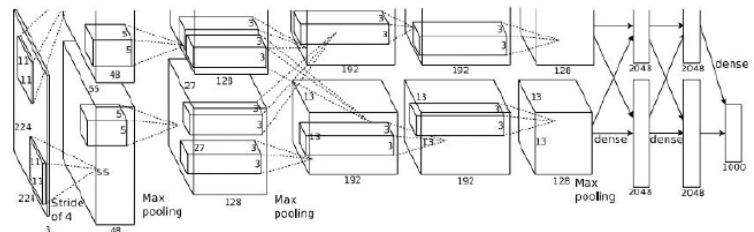
pooling 的参数一般为 0。

## Case study

### AlexNet

#### Case Study: AlexNet

[Krizhevsky et al. 2012]



Input:  $227 \times 227 \times 3$  images

**First layer** (CONV1): 96  $11 \times 11$  filters applied at stride 4

=>

Output volume  **$55 \times 55 \times 96$**

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

### VGGnet

# Case Study: VGGNet

[Simonyan and Zisserman, 2014]

Only 3x3 CONV stride 1, pad 1  
and 2x2 MAX POOL stride 2

best model

11.2% top 5 error in ILSVRC 2013

->

7.3% top 5 error

ConvNet Configuration					
A	A-LRN	B	C	D	E
11 weight layers	11 weight layers	13 weight layers	16 weight layers	16 weight layers	19 weight layers
input (224 × 224 RGB image)					
conv3-64	conv3-64 LRN	conv3-64	conv3-64	conv3-64	conv3-64
maxpool					
conv3-128	conv3-128	conv3-128	conv3-128	conv3-128	conv3-128
maxpool					
conv3-256	conv3-256	conv3-256	conv3-256	conv3-256	conv3-256
conv3-256	conv3-256	conv3-256	conv3-256	conv3-256	conv3-256
maxpool					
conv3-512	conv3-512	conv3-512	conv3-512	conv3-512	conv3-512
conv3-512	conv3-512	conv3-512	conv3-512	conv3-512	conv3-512
maxpool					
conv3-512	conv3-512	conv3-512	conv3-512	conv3-512	conv3-512
conv3-512	conv3-512	conv3-512	conv3-512	conv3-512	conv3-512
maxpool					
FC-4096					
FC-4096					
FC-1000					
soft-max					

Table 2: Number of parameters (in millions).

Network	A,A-LRN	B	C	D	E
Number of parameters	133	133	134	138	144

INPUT: [224x224x3] memory:  $224*224*3=150K$  params: 0 (not counting biases)

CONV3-64: [224x224x64] memory:  $224*224*64=3.2M$  params:  $(3*3*3)*64 = 1,728$

CONV3-64: [224x224x64] memory:  $224*224*64=3.2M$  params:  $(3*3*64)*64 = 36,864$

POOL2: [112x112x64] memory:  $112*112*64=800K$  params: 0

CONV3-128: [112x112x128] memory:  $112*112*128=1.6M$  params:  $(3*3*64)*128 = 73,728$

CONV3-128: [112x112x128] memory:  $112*112*128=1.6M$  params:  $(3*3*128)*128 = 147,456$

POOL2: [56x56x128] memory:  $56*56*128=400K$  params: 0

CONV3-256: [56x56x256] memory:  $56*56*256=800K$  params:  $(3*3*128)*256 = 294,912$

CONV3-256: [56x56x256] memory:  $56*56*256=800K$  params:  $(3*3*256)*256 = 589,824$

CONV3-256: [56x56x256] memory:  $56*56*256=800K$  params:  $(3*3*256)*256 = 589,824$

POOL2: [28x28x256] memory:  $28*28*256=200K$  params: 0

CONV3-512: [28x28x512] memory:  $28*28*512=400K$  params:  $(3*3*256)*512 = 1,179,648$

CONV3-512: [28x28x512] memory:  $28*28*512=400K$  params:  $(3*3*512)*512 = 2,359,296$

CONV3-512: [28x28x512] memory:  $28*28*512=400K$  params:  $(3*3*512)*512 = 2,359,296$

POOL2: [14x14x512] memory:  $14*14*512=100K$  params: 0

CONV3-512: [14x14x512] memory:  $14*14*512=100K$  params:  $(3*3*512)*512 = 2,359,296$

CONV3-512: [14x14x512] memory:  $14*14*512=100K$  params:  $(3*3*512)*512 = 2,359,296$

CONV3-512: [14x14x512] memory:  $14*14*512=100K$  params:  $(3*3*512)*512 = 2,359,296$

POOL2: [7x7x512] memory:  $7*7*512=25K$  params: 0

FC: [1x1x4096] memory: 4096 params:  $7*7*512*4096 = 102,760,448$

FC: [1x1x4096] memory: 4096 params:  $4096*4096 = 16,777,216$

FC: [1x1x1000] memory: 1000 params:  $4096*1000 = 4,096,000$

Note:

Most memory is in early CONV

Most params are in late FC

TOTAL memory:  $24M * 4 \text{ bytes} \sim 93MB / \text{image}$  (only forward!  $\sim 2$  for bwd)

TOTAL params: 138M parameters

## Summary

Typical architectures look like

**[(CONV-RELU)\*N-POOL?]\*M-(FC-RELU)\*K,SOFTMAX**

where N is usually up to  $\sim 5$ , M is large,  $0 \leq K \leq 2$ .

- but recent advances such as ResNet/GoogLeNet challenge this paradigm