



# Convolutional Neural Network Architecture

2023-2 KUBIG 방학세션  
DL



1. Convolution layer
2. Max Pooling layer
3. AlexNet- Top Network from ImageNet Classification
4. Q & A

# 0. Paper Review

# 0. Paper Review

## Normalization in Pytorch

$$y = \frac{x - E[x]}{\sqrt{\text{Var}[x] + \epsilon}} * \gamma + \beta$$

`nn.BatchNorm2d(num_features, eps=1e-5, affine=True)` `nn.LayerNorm(normalized_shape, eps=1e-5, elementwise_affine=True)`

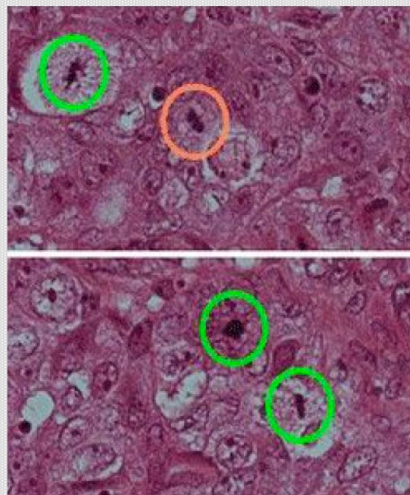
```
# With Learnable Parameters
m = nn.BatchNorm2d(100)
# Without Learnable Parameters
m = nn.BatchNorm2d(100, affine=False)
input = torch.randn(20, 100, 35, 45)
output = m(input)
```

`nn.BatchNorm3d(num_features, ...)`

```
# With Learnable Parameters
m = nn.BatchNorm3d(100)
# Without Learnable Parameters
m = nn.BatchNorm3d(100, affine=False)
input = torch.randn(20, 100, 35, 45, 10)
output = m(input)
```

```
# NLP Example
batch, sentence_length, embedding_dim = 20, 5, 10
embedding = torch.randn(batch, sentence_length, embedding_dim)
layer_norm = nn.LayerNorm(embedding_dim)
# Activate module
layer_norm(embedding)
```

# 0. CNNs are Everywhere



[Ciresan et al., 2013]

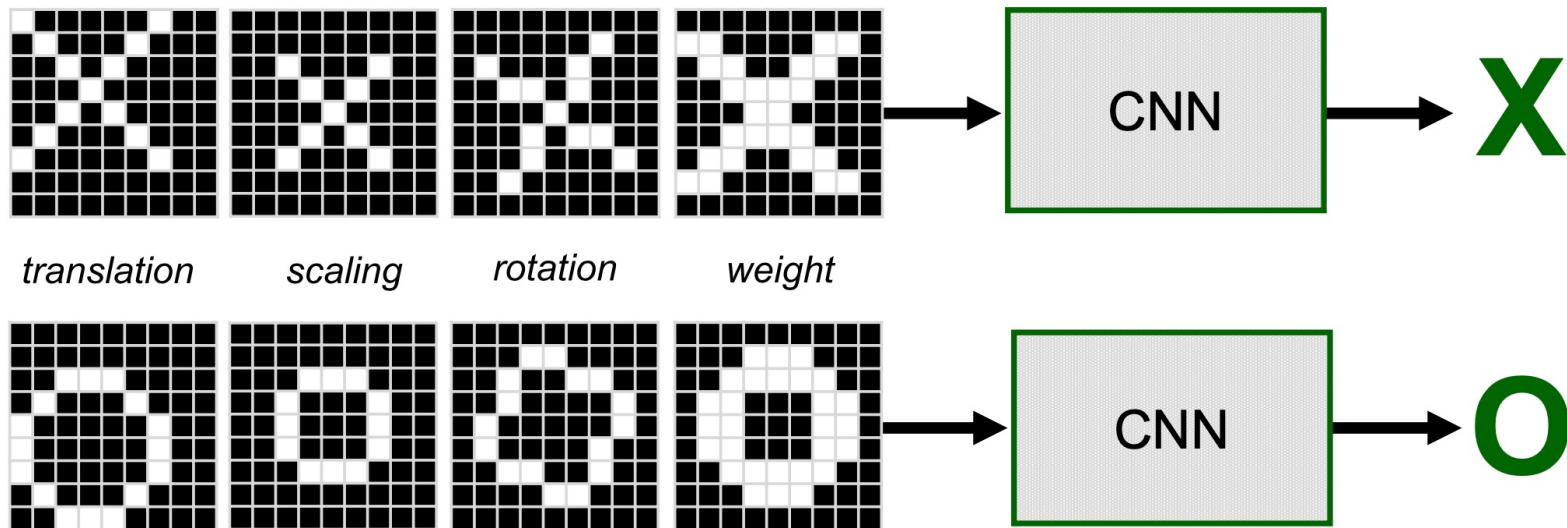


[Sermanet et al., 2011]



# 0. CNNs are Everywhere

## ► Trickier cases



[Citation] Lecture 13: Convolutional Neural Networks -1, 데이터  
학습과 지능, 정태수, 고려대학교

# 0. CNNs are Everywhere

## ► Computers are literal

-1	-1	-1	-1	-1	-1	-1	-1	-1
-1	1	-1	-1	-1	-1	-1	1	-1
-1	-1	1	-1	-1	-1	1	-1	-1
-1	-1	-1	1	-1	1	-1	-1	-1
-1	-1	-1	-1	1	-1	-1	-1	-1
-1	-1	-1	1	-1	1	-1	-1	-1
-1	-1	1	-1	-1	-1	1	-1	-1
-1	1	-1	-1	-1	-1	-1	1	-1
-1	-1	-1	-1	-1	-1	-1	-1	-1



-1	-1	-1	-1	-1	-1	-1	-1	-1
-1	-1	-1	-1	-1	-1	1	-1	-1
-1	1	-1	-1	-1	1	-1	-1	-1
-1	-1	1	1	-1	1	-1	-1	-1
-1	-1	-1	-1	1	-1	-1	-1	-1
-1	-1	-1	1	-1	1	1	-1	-1
-1	-1	-1	1	-1	-1	-1	1	-1
-1	-1	1	-1	-1	-1	-1	-1	-1
-1	-1	-1	-1	-1	-1	-1	-1	-1

[Citation] Lecture 13: Convolutional Neural Networks -1, 데이터  
학습과 지능, 정태수, 고려대학교

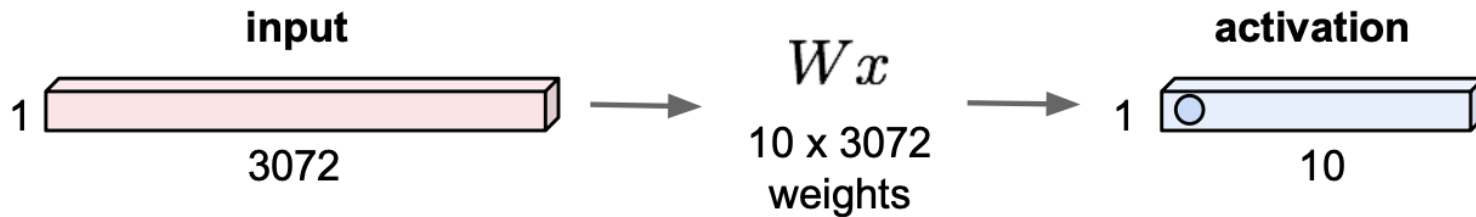


# 1. Convolution Layer

---

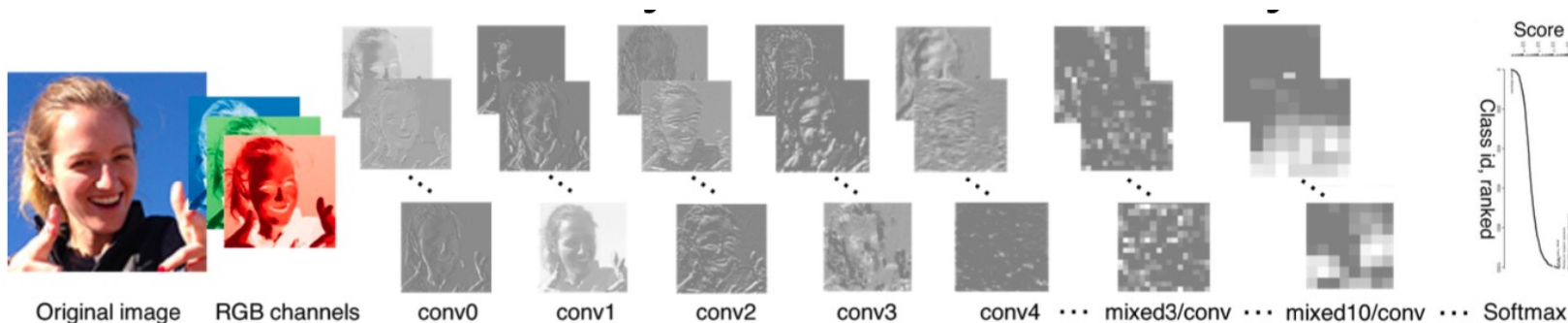
\* Recap: FC layer

32x32x3 image -> stretch to 3072 x 1



# 1. Convolution Layer

- \* Convolution = "합성곱"
- \* Similar to fully connected layer but somewhat different
- \* Use "filter" to "convolve" input data

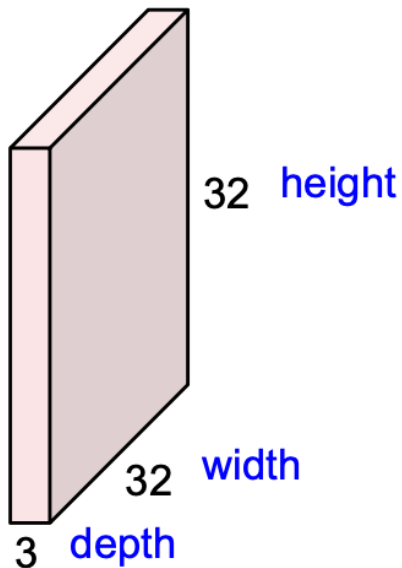


[Citation] Lecture 5: Convolutional Neural Networks, Fei-Fei-Li & Justin Johnson & Serena Yeung

# 1. Convolution Layer

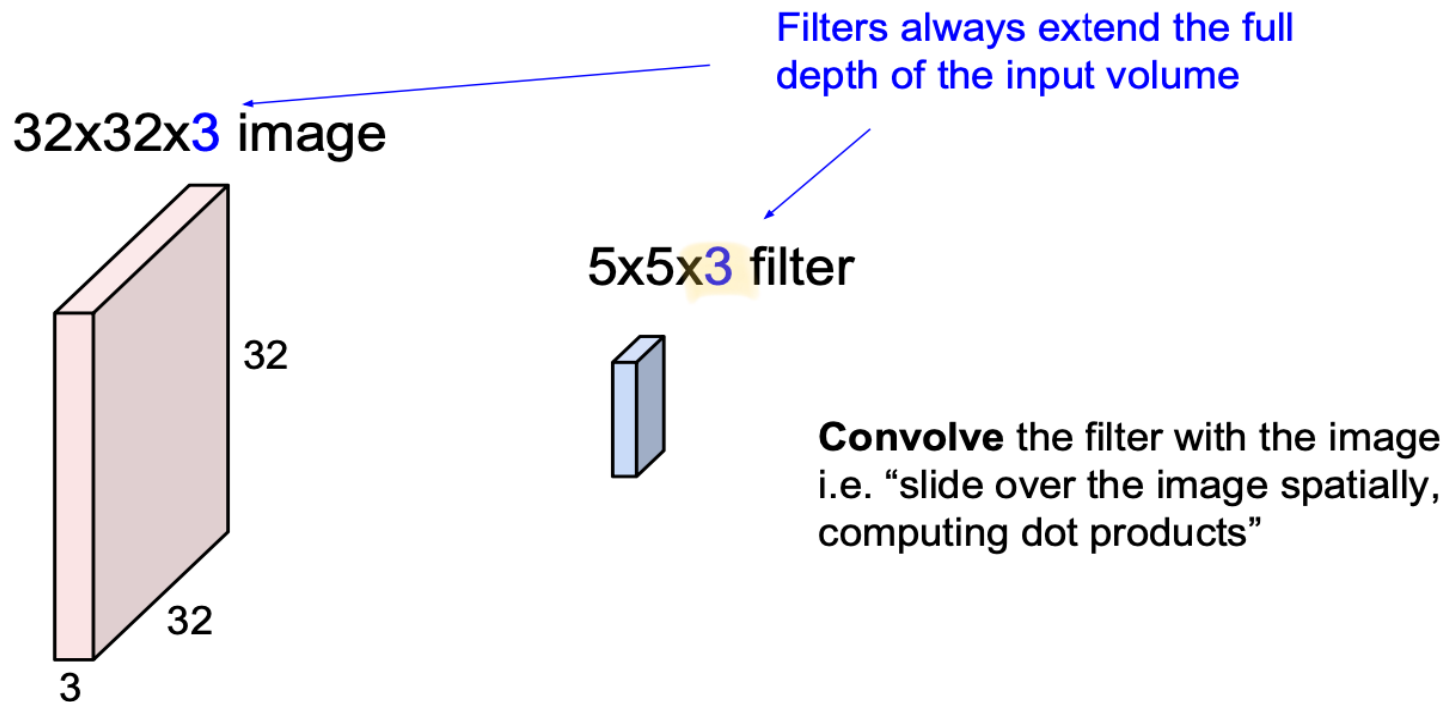
---

32x32x3 image -> preserve spatial structure

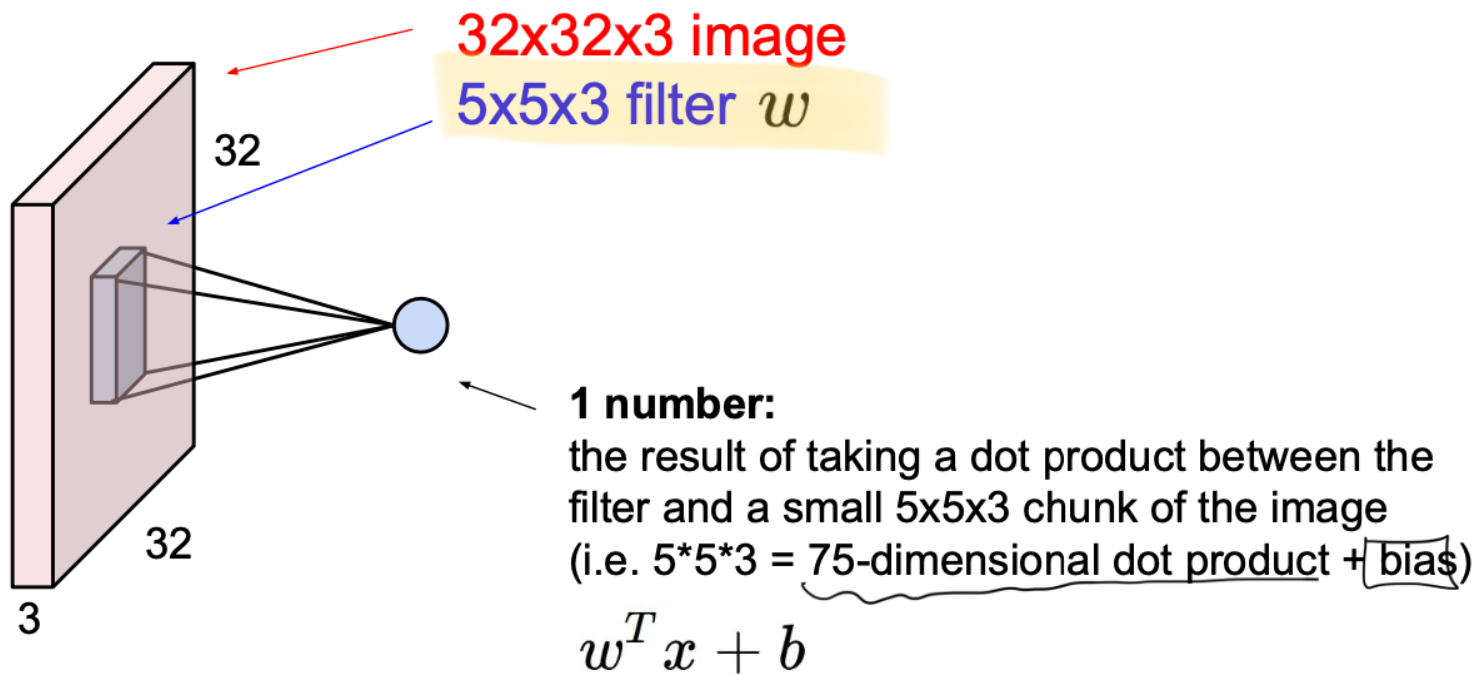


[Citation] Lecture 5: Convolutional Neural Networks, Fei-Fei-Li & Justin Johnson & Serena Yeung

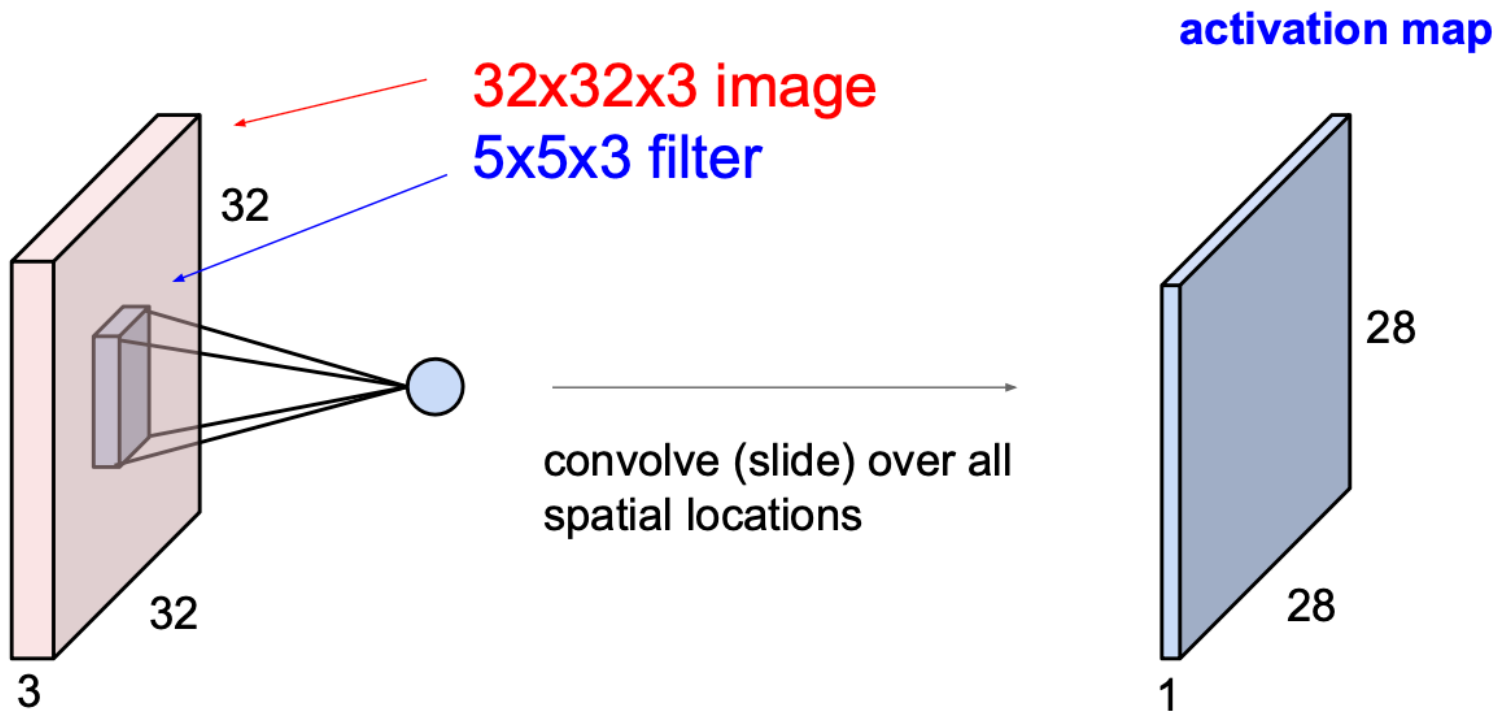
# 1. Convolution Layer



# 1. Convolution Layer

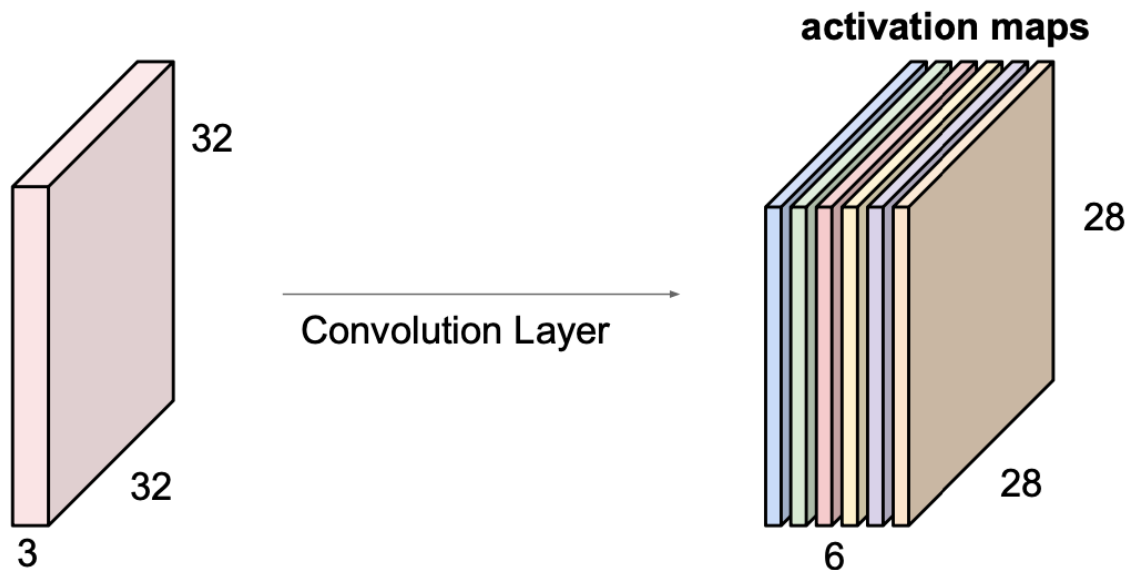


# 1. Convolution Layer



# 1. Convolution Layer

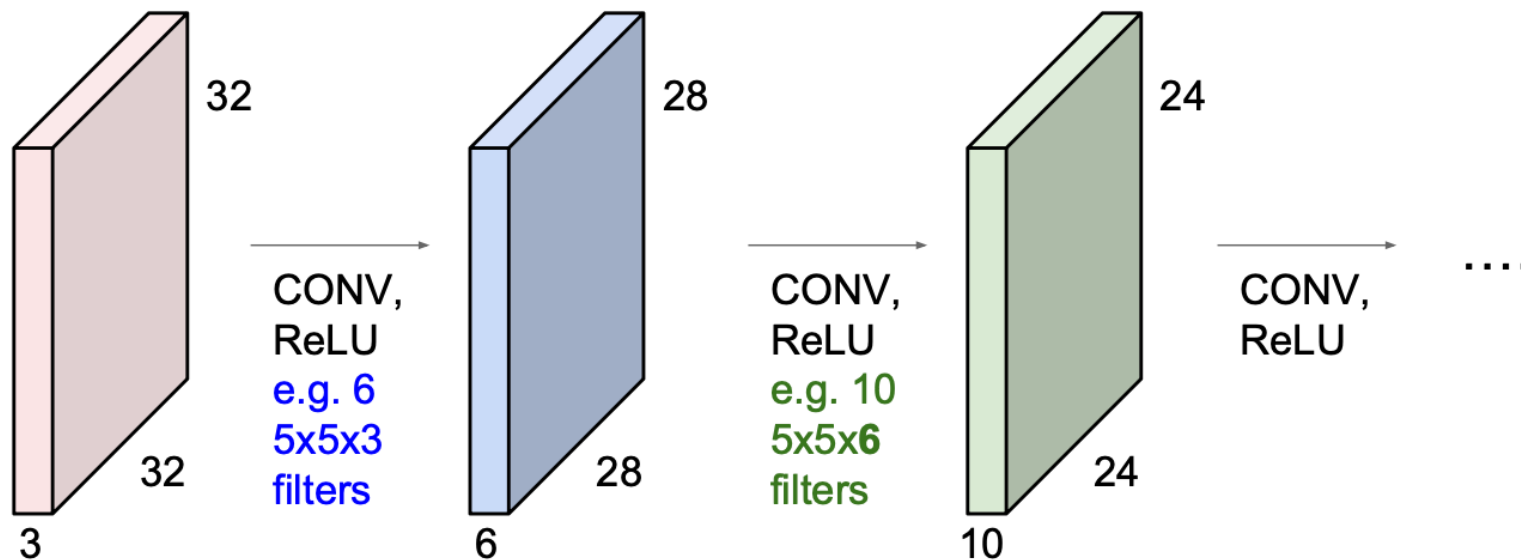
For example, if we had 6 5x5 filters, we'll get 6 separate activation maps:



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

# 1. Convolution Layer

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

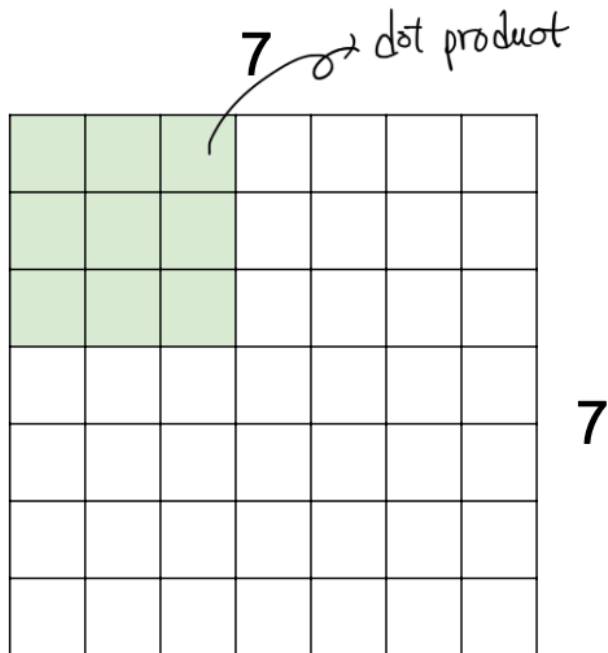


[Citation] Lecture 5: Convolutional Neural Networks, Fei-Fei-Li & Justin Johnson & Serena Yeung



# 1. Convolution Layer

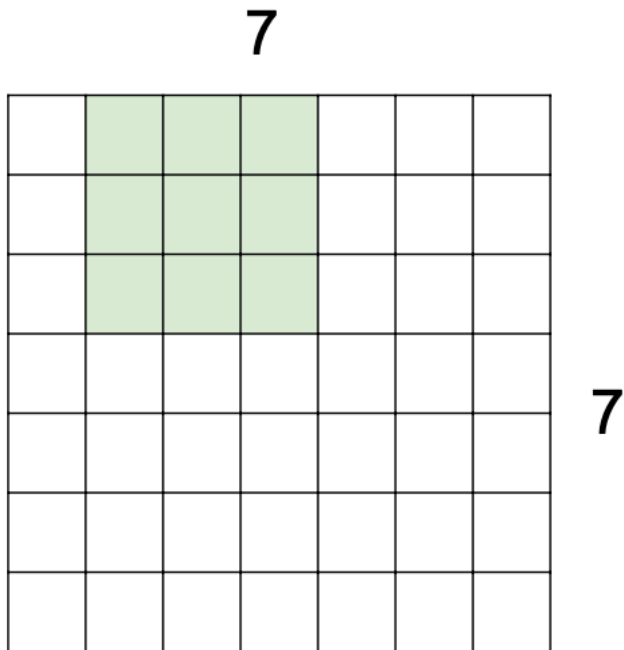
---



7x7 input (spatially)  
assume 3x3 filter

# 1. Convolution Layer

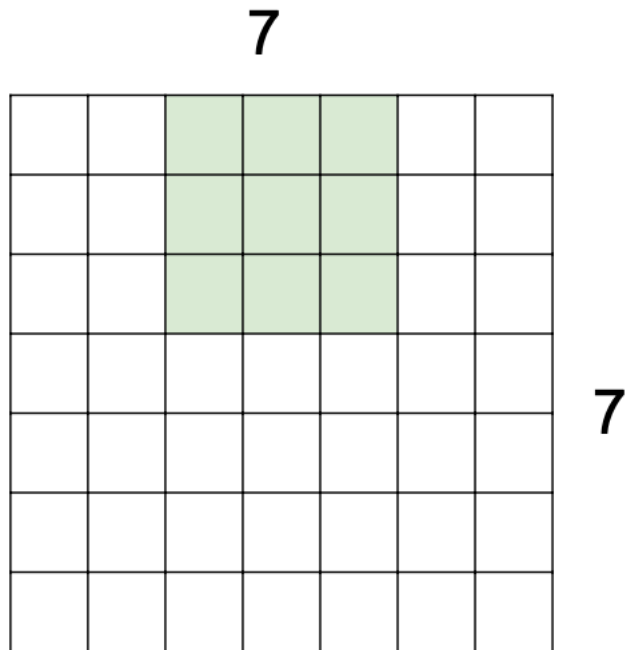
---



7x7 input (spatially)  
assume 3x3 filter

# 1. Convolution Layer

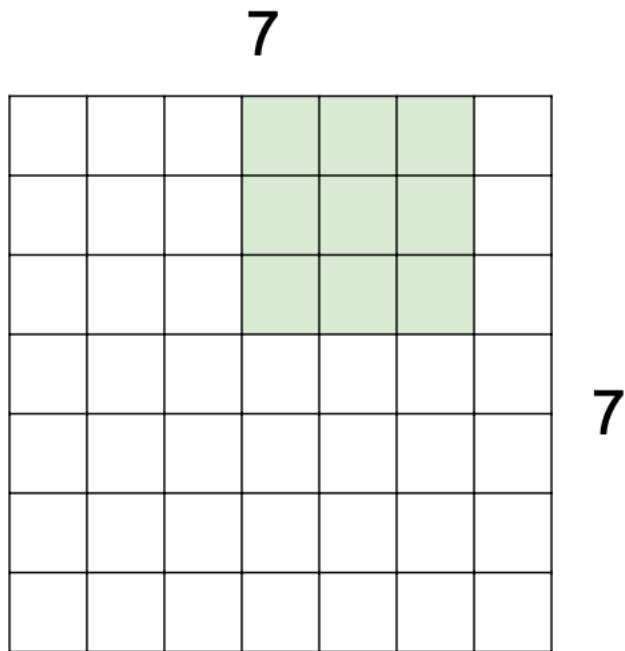
---



7x7 input (spatially)  
assume 3x3 filter

# 1. Convolution Layer

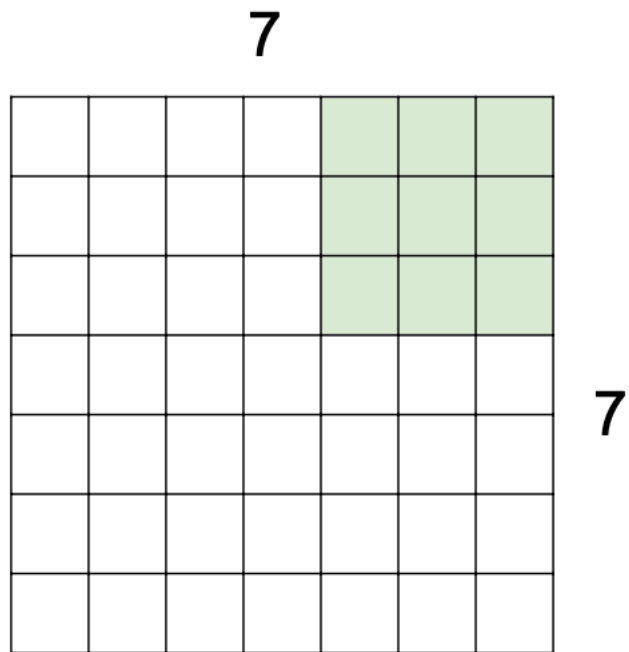
---



7x7 input (spatially)  
assume 3x3 filter

# 1. Convolution Layer

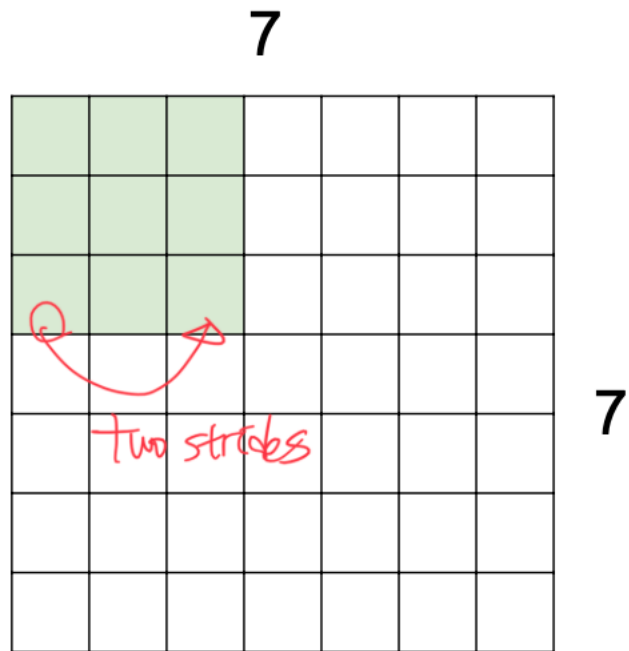
---



7x7 input (spatially)  
assume 3x3 filter

=> **5x5 output**

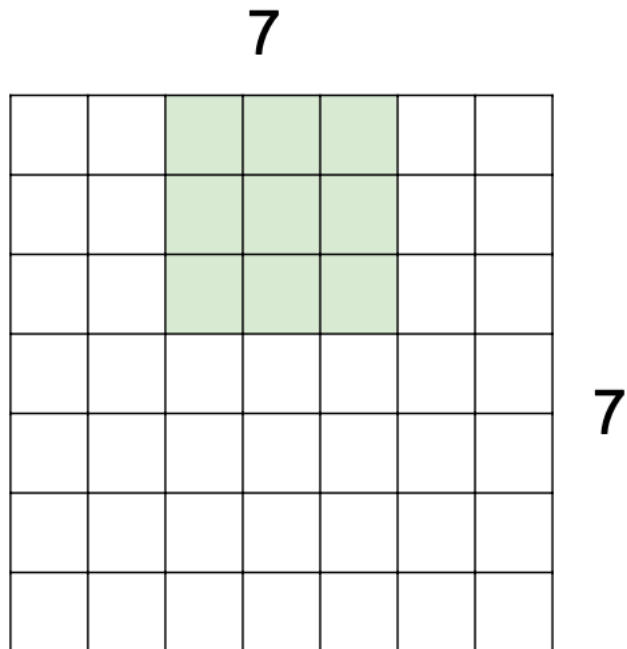
# 1. Convolution Layer



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

# 1. Convolution Layer

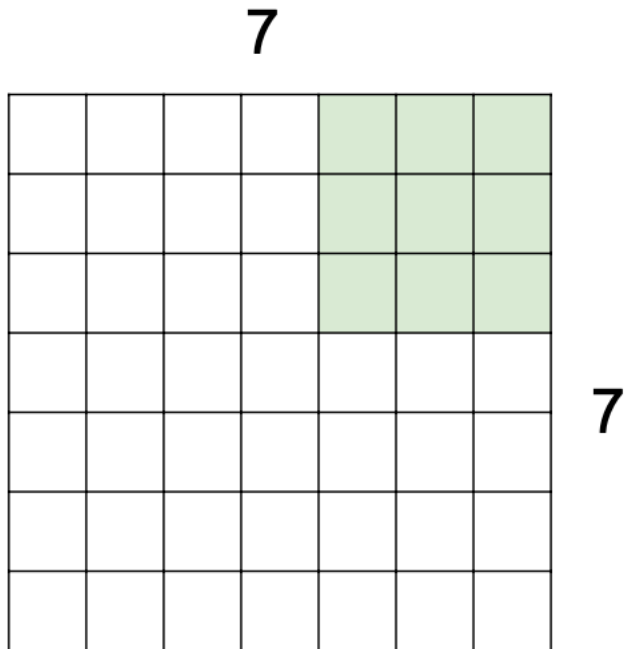
---



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

# 1. Convolution Layer

---

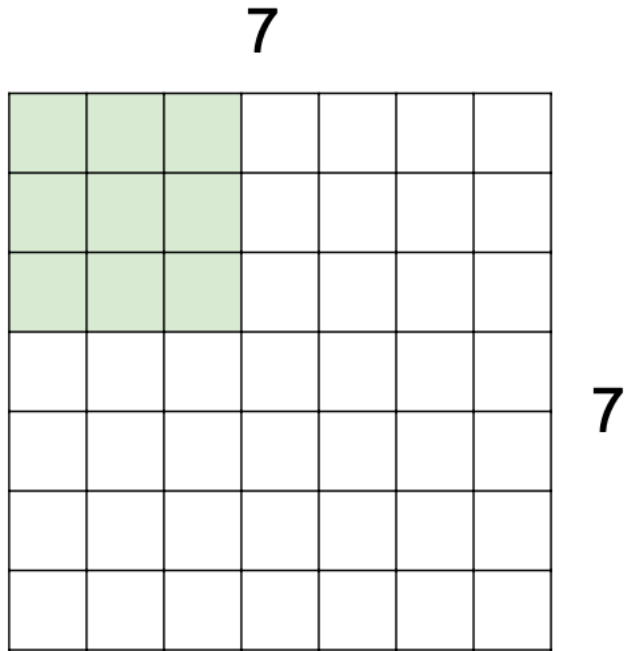


7x7 input (spatially)  
assume 3x3 filter  
applied **with stride 2**  
=> **3x3 output!**



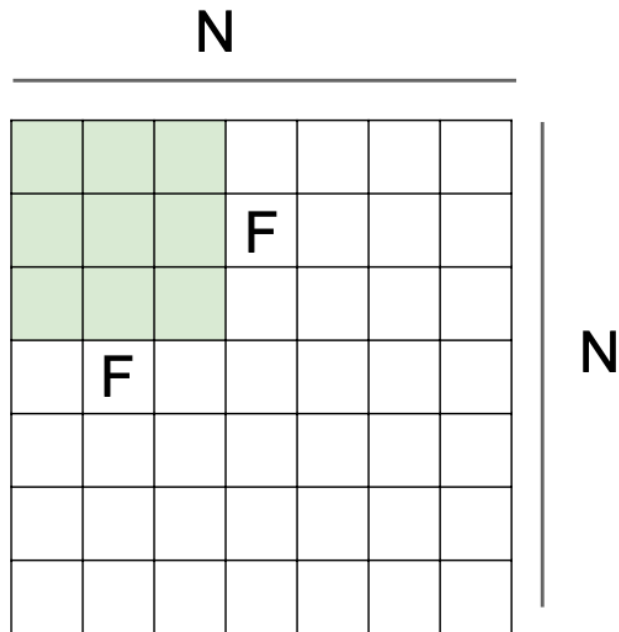
# 1. Convolution Layer

---



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

# 1. Convolution Layer



Output size:

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

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

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

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

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

# 1. Convolution Layer

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!**

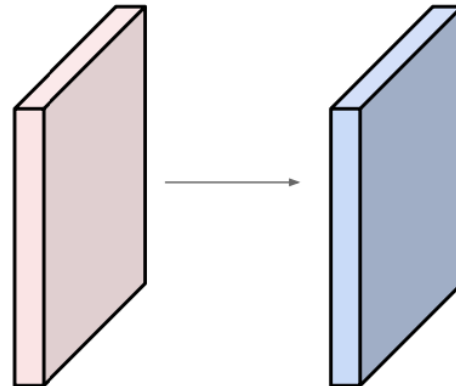
# 1. Convolution Layer

---

Examples time:

Input volume: **32x32x3**  
10 5x5 filters with stride 1, pad 2

Output volume size: ?



# 1. Convolution Layer

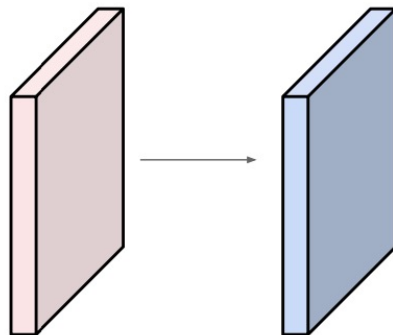
Examples time:

Input volume: **32x32x3**

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

Output volume size:

$$\frac{(\text{Input } 32 + 2 \times \text{padding } 2 - \text{Filter } 5)}{\text{stride } 1} + 1 = 32 \text{ spatially, so}$$
**32x32x10**



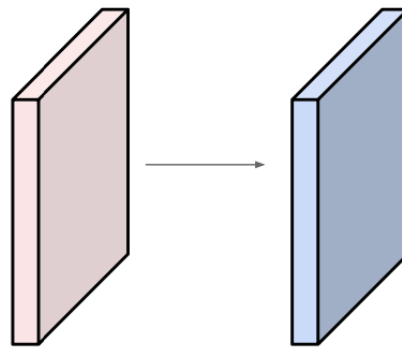
# 1. Convolution Layer

---

Examples time:

Input volume: **32x32x3**

10 5x5 filters with stride 1, pad 2

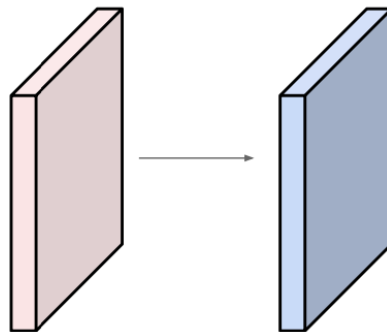


Number of parameters in this layer?

# 1. Convolution Layer

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

$$\Rightarrow 76 * 10 = 760$$

ten filters

$$W^T x + b$$

One

(+1 for bias)

# 1. Convolution Layer(Summary)

Common settings:

**Summary.** To summarize, the Conv Layer:

- Accepts a volume of size  $W_1 \times H_1 \times D_1$
- Requires four hyperparameters:

- Number of filters  $K$ ,
- their spatial extent  $F$ ,
- the stride  $S$ ,
- the amount of zero padding  $P$ .

- Produces a volume of size  $W_2 \times H_2 \times D_2$  where:
  - $W_2 = (W_1 - F + 2P)/S + 1$
  - $H_2 = (H_1 - F + 2P)/S + 1$  (i.e. width and height are computed equally by symmetry)
  - $D_2 = K$
- With parameter sharing, it introduces  $F \cdot F \cdot D_1$  weights per filter, for a total of  $(F \cdot F \cdot D_1) \cdot K$  weights and  $K$  biases.
- In the output volume, the  $d$ -th depth slice (of size  $W_2 \times H_2$ ) is the result of performing a valid convolution of the  $d$ -th filter over the input volume with a stride of  $S$ , and then offset by  $d$ -th bias.

$K =$  (powers of 2, e.g. 32, 64, 128, 512)

- $F = 3, S = 1, P = 1$
- $F = 5, S = 1, P = 2$
- $F = 5, S = 2, P = ?$  (whatever fits)
- $F = 1, S = 1, P = 0$



# 1. Convolution Layer

```
# In[11]:
```

```
conv = nn.Conv2d(3, 16, kernel_size=3) ←  
conv
```

Instead of the shortcut `kernel_size=3`, we could equivalently pass in the tuple that we see in the output: `kernel_size=(3, 3)`.

```
# Out[11]:
```

```
Conv2d(3, 16, kernel_size=(3, 3), stride=(1, 1))
```

```
# In[12]:
```

```
conv.weight.shape, conv.bias.shape
```

```
# Out[12]:
```

```
(torch.Size([16, 3, 3, 3]), torch.Size([16]))
```

# 1. Convolution Layer

```
# In[16]:  
conv = nn.Conv2d(3, 1, kernel_size=3, padding=1)      ← Now with padding  
output = conv(img.unsqueeze(0))  
img.unsqueeze(0).shape, output.shape  
  
# Out[16]:  
(torch.Size([1, 3, 32, 32]), torch.Size([1, 1, 32, 32]))
```

`nn.ConstantPad1d(padding, value)`  
`nn.ConstantPad2d(padding, value)`  
`nn.ConstantPad3d(padding, value)`

# 1. Convolution Layer

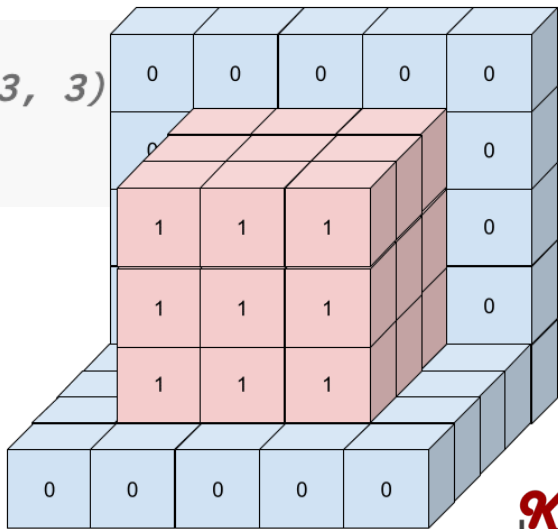
`torch.nn.functional.pad`(input, pad, mode='constant', value=None)

-> padding starts from the last dimension

-> pad: int or tuple\_(left,right) or (left,right, top, bottom)  
or (left, right, top, bottom, front back)

```
t4d = torch.empty(3, 3, 4, 2)
p3d = (0, 1, 2, 1, 3, 3) # pad by (0, 1), (2, 1), and (3, 3)
out = F.pad(t4d, p3d, "constant", 0)
print(out.size())
```

\*How about Size?

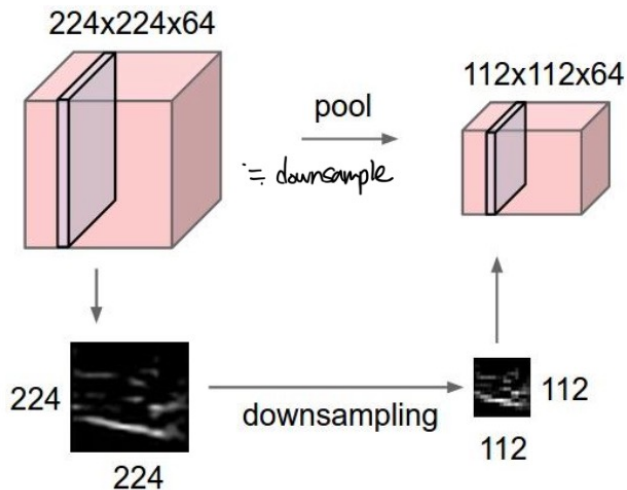


## 2. Max Pooling Layer

## 2. Max Pooling Layer

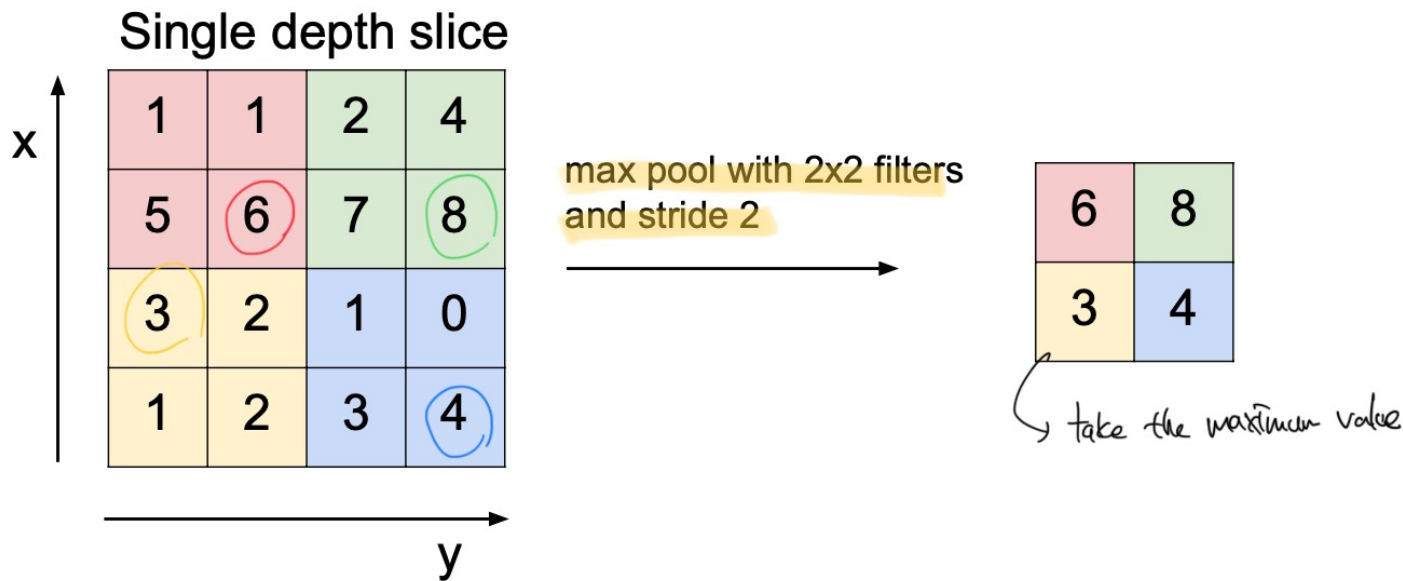
### Pooling layer

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



## 2. Max Pooling Layer

### MAX POOLING



## 2. Max Pooling Layer

---

- Accepts a volume of size  $W_1 \times H_1 \times D_1$
- Requires three hyperparameters:
  - their spatial extent  $F$ , : FilterSize
  - the stride  $S$ ,
- Produces a volume of size  $W_2 \times H_2 \times D_2$  where:
  - $W_2 = (W_1 - F)/S + 1$
  - $H_2 = (H_1 - F)/S + 1$
  - $D_2 = D_1$
- Introduces zero parameters since it computes a fixed function of the input
- Note that it is not common to use zero-padding for Pooling layers

## 2. Max Pooling Layer

```
# In[21]:
pool = nn.MaxPool2d(2)
output = pool(img.unsqueeze(0))

img.unsqueeze(0).shape, output.shape

# Out[21]:
(torch.Size([1, 3, 32, 32]), torch.Size([1, 3, 16, 16]))
```

```
>>> # pool of square window of size=3, stride=2
>>> m = nn.MaxPool2d(3, stride=2)
>>> # pool of non-square window
>>> m = nn.MaxPool2d((3, 2), stride=(2, 1))
>>> input = torch.randn(20, 16, 50, 32)
>>> output = m(input)
```

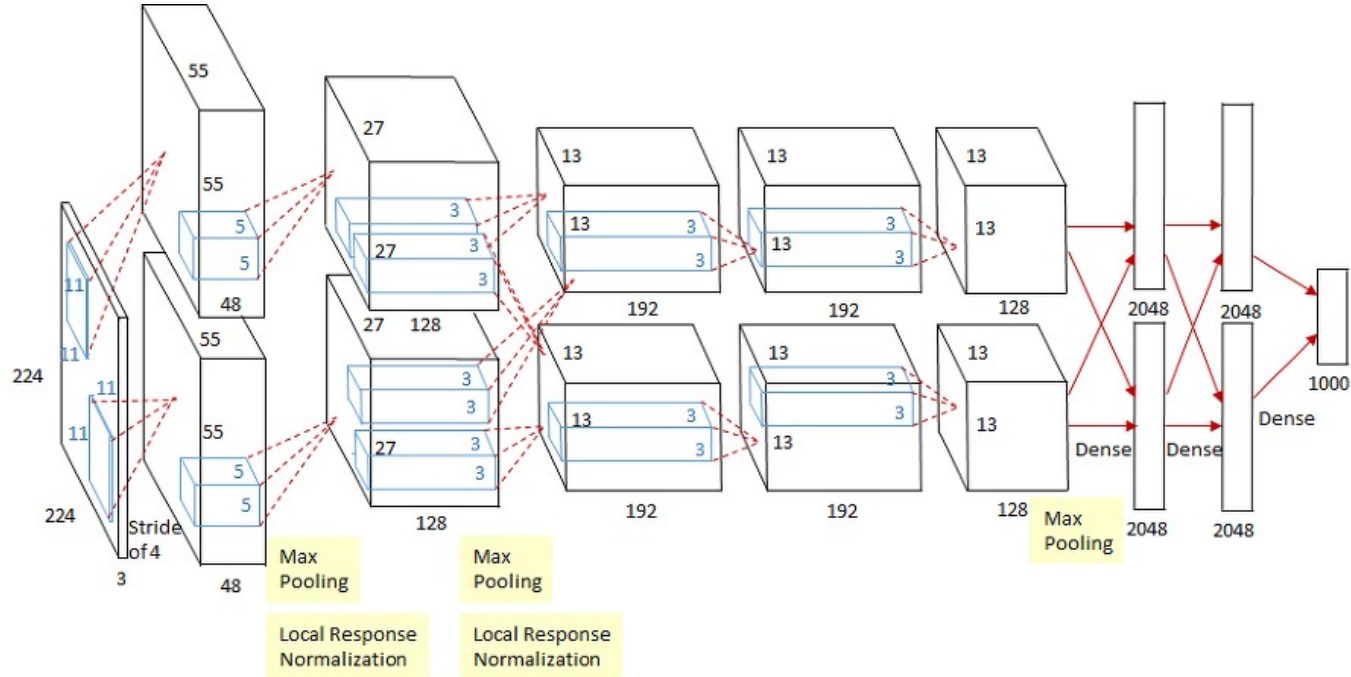




### 3. AlexNet

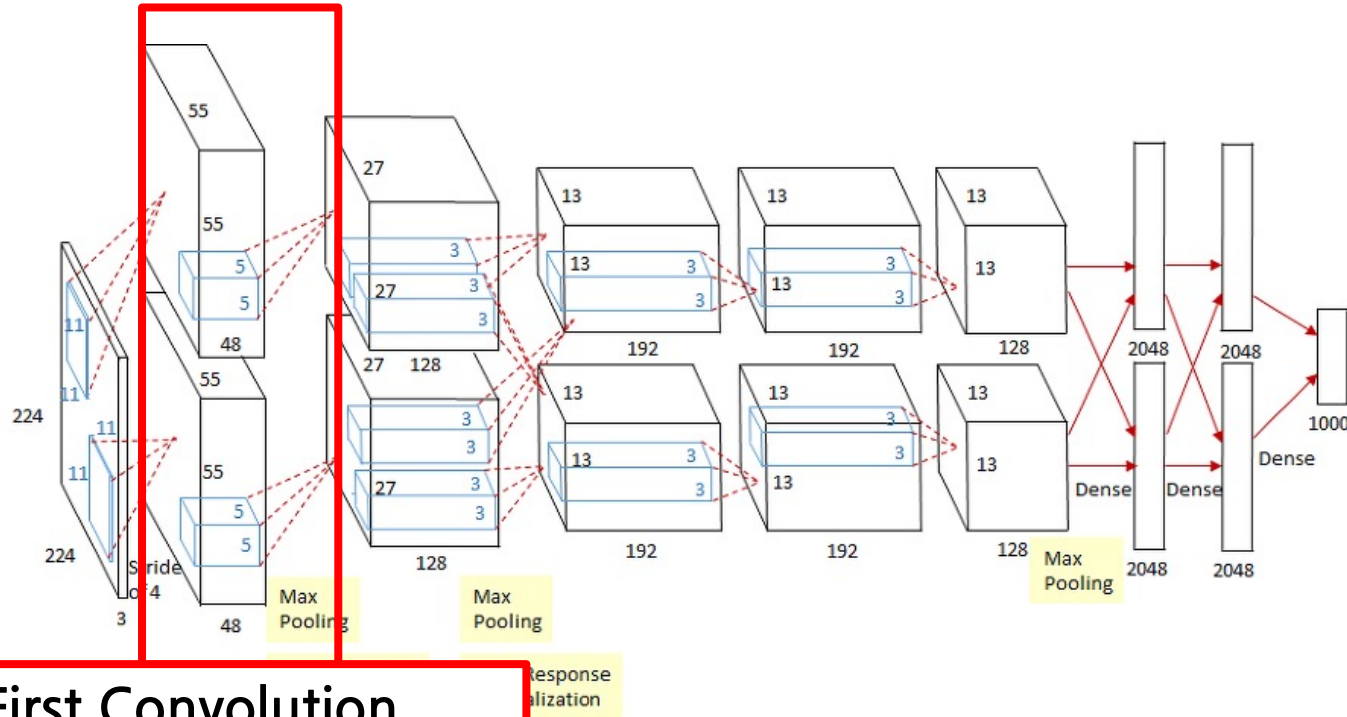
- 2012 Top Network from ImageNet Classification

# 3. AlexNet - 2012 Top Network from ImageNet Classification



[Citation] Lecture 5: Convolutional Neural Networks, Fei-Fei-Li & Justin Johnson & Serena Yeung

### 3. AlexNet - 2012 Top Network from ImageNet Classification

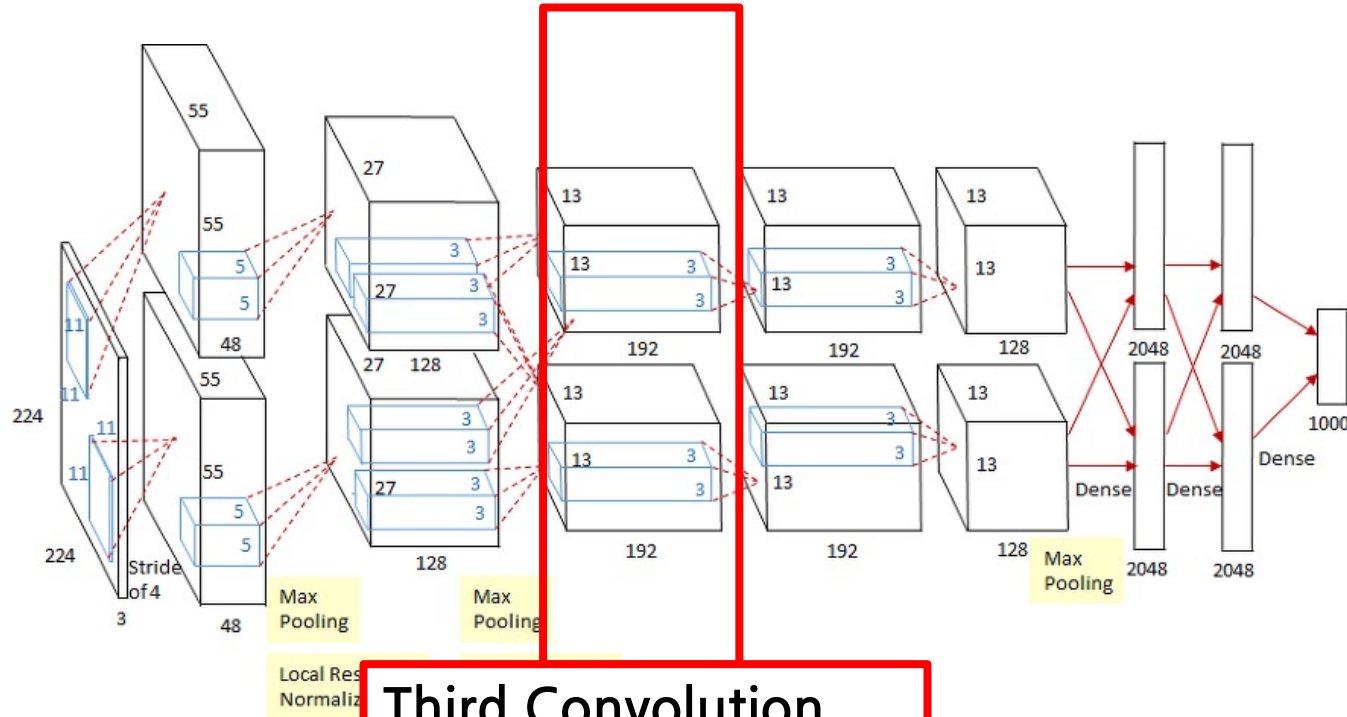


### 3. AlexNet - 2012 Top Network from ImageNet Classification



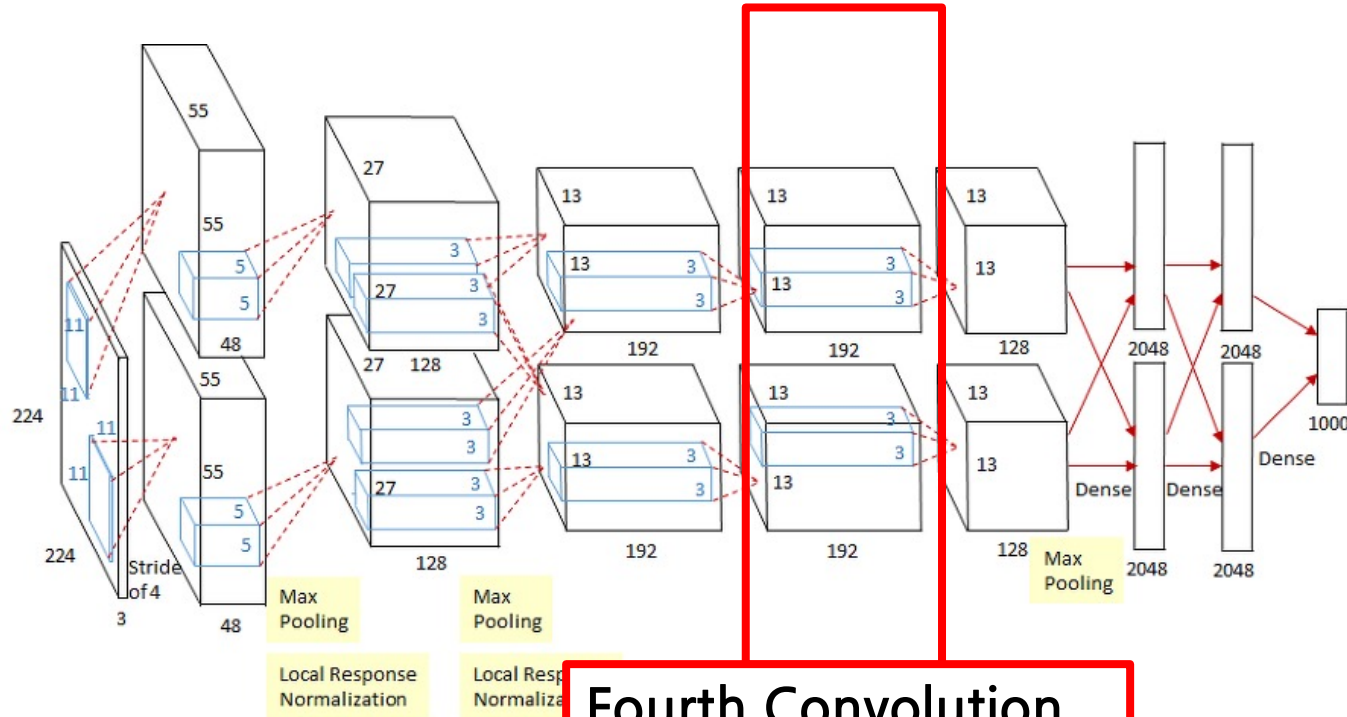
Second Convolution  
with 256 Filters 5\*5

### 3. AlexNet - 2012 Top Network from ImageNet Classification



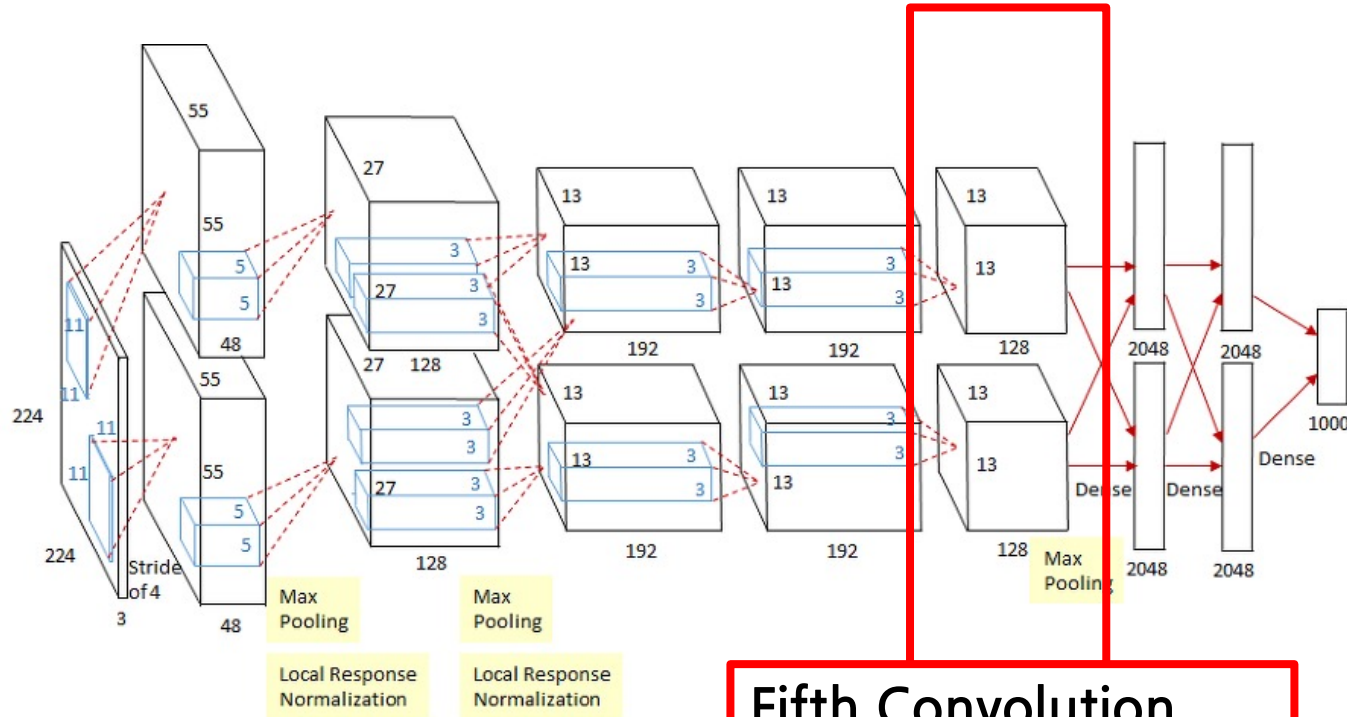
Third Convolution  
with 384 Filters  $3 \times 3$

### 3. AlexNet - 2012 Top Network from ImageNet Classification



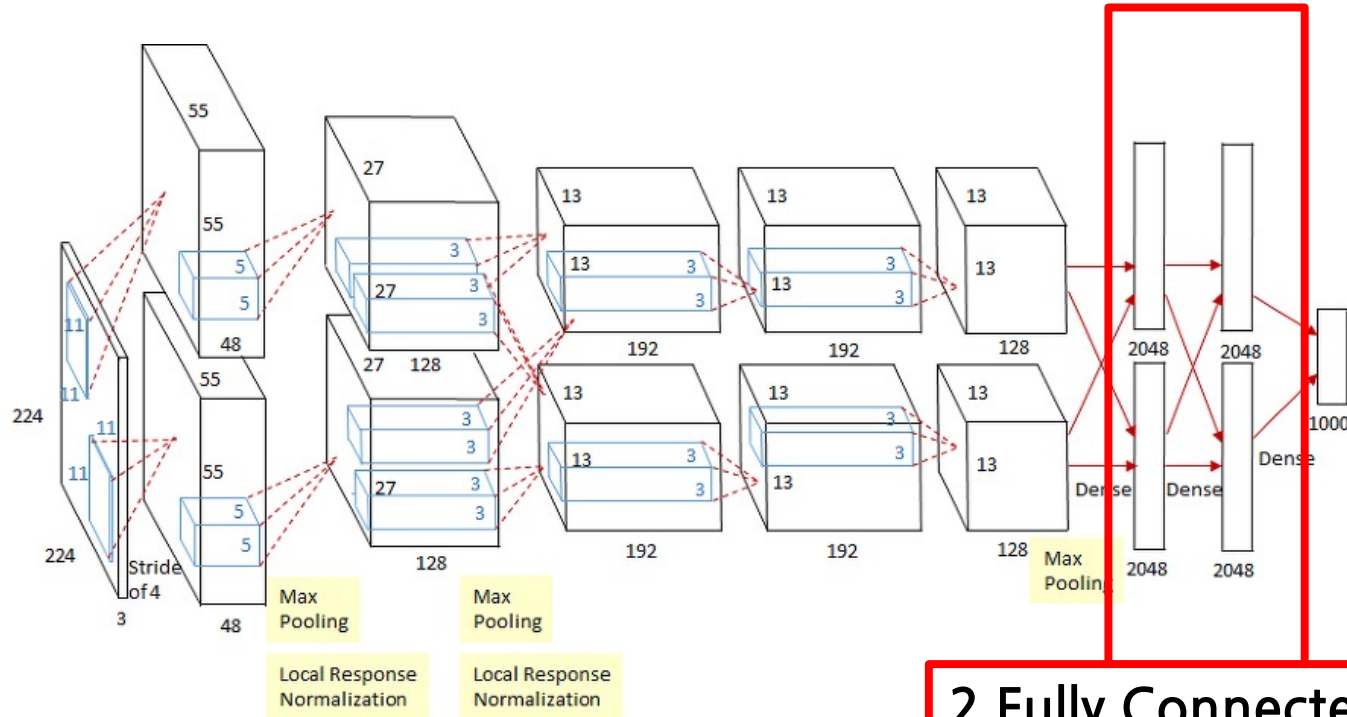
Fourth Convolution  
with 384 Filters 3\*3

### 3. AlexNet - 2012 Top Network from ImageNet Classification



Fifth Convolution  
with 256 Filters 3\*3

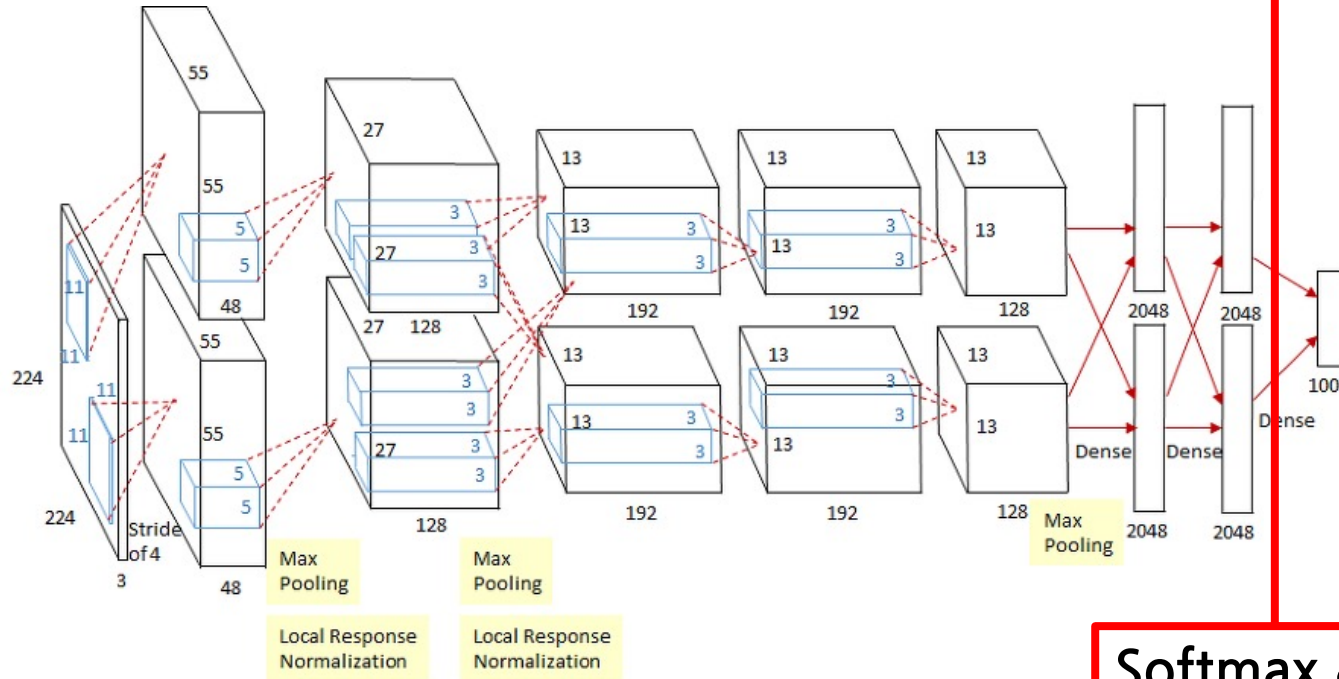
### 3. AlexNet - 2012 Top Network from ImageNet Classification



2 Fully Connected  
layer with 4096-dim



### 3. AlexNet - 2012 Top Network from ImageNet Classification



**Softmax &  
Backpropagation !**

## 4. Other Remarks

---

ConvBlock = stacking Conv layer, Pooling layer, FC layer

- \* General Structure of ConvNet

- (Conv-ReLU)\*N  $\rightarrow$  Pool(or not)  $\rightarrow$  FC layer  $\rightarrow$  ReLU  $\rightarrow$  FC or Softmax

- \* Training Technique

- Batch Normalization, LR Scheduler, Optimizer, Dropout, Data Augmentation
- SOTA: Transformer,,

- \* Limitation

- Deeper Conv layers don't guarantee better performance

## 4. Other Remarks

```
# In[17]:  
with torch.no_grad():  
    conv.bias.zero_()  
  
with torch.no_grad():  
    conv.weight.fill_(1.0 / 9.0)
```

Blurred image due to convolution

```
# In[18]:  
output = conv(img.unsqueeze(0))  
plt.imshow(output[0, 0].detach(), cmap='gray')  
plt.show()
```

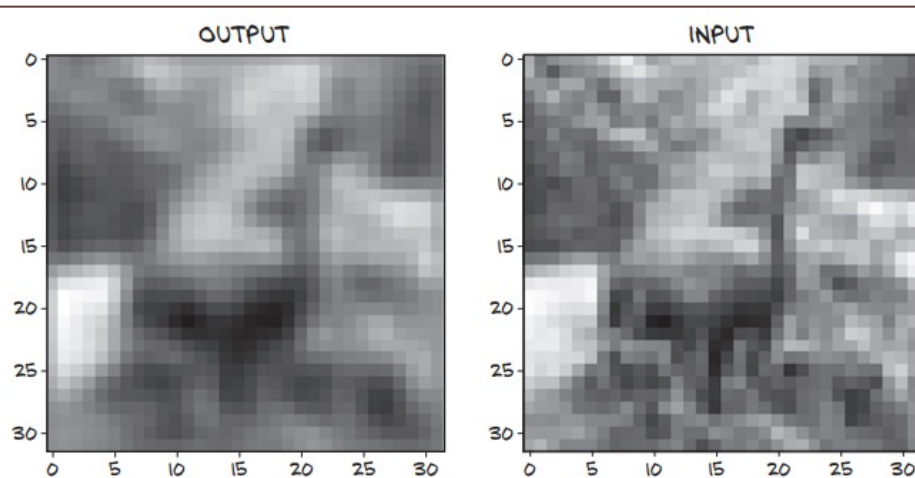


Figure 8.4 Our bird, this time blurred thanks to a constant convolution kernel

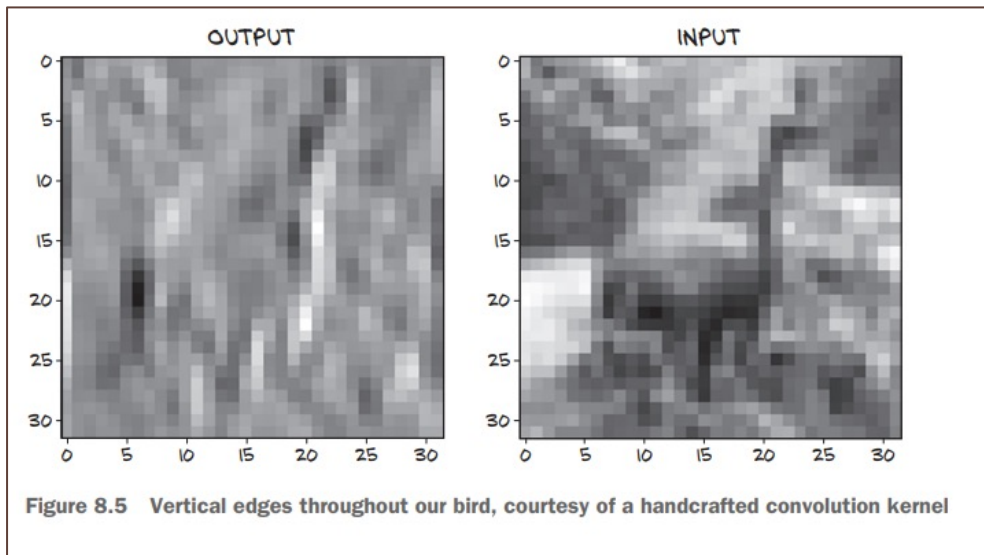
## 4. Other Remarks

```
# In[19]:
conv = nn.Conv2d(3, 1, kernel_size=3, padding=1)

with torch.no_grad():
    conv.weight[:] = torch.tensor([[-1.0, 0.0, 1.0],
                                   [-1.0, 0.0, 1.0],
                                   [-1.0, 0.0, 1.0]])

    conv.bias.zero_()
```

### Conv layer with Vertical Edge



Q&A

Have a nice week 🥰