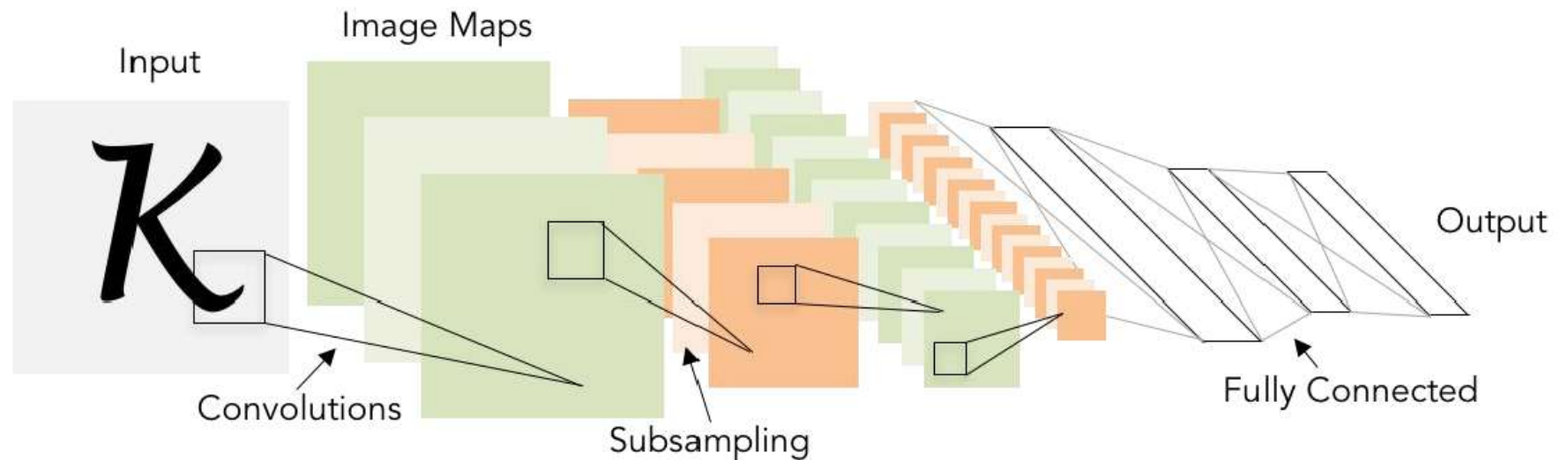


# CNN Architectures

Some slides were adapted/taken from various sources, including Andrew Ng's Coursera Lectures, CS231n: Convolutional Neural Networks for Visual Recognition lectures, Stanford University CSWaterloo Canada lectures, Aykut Erdem, et.al. tutorial on Deep Learning in Computer Vision, Ismini Lourentzou's lecture slide on "Introduction to Deep Learning", Ramprasaath's lecture slides, and many more. We thankfully acknowledge them. Students are requested to use this material for their study only and **NOT** to distribute it.

# LeNet-5

[LeCun et al., 1998]

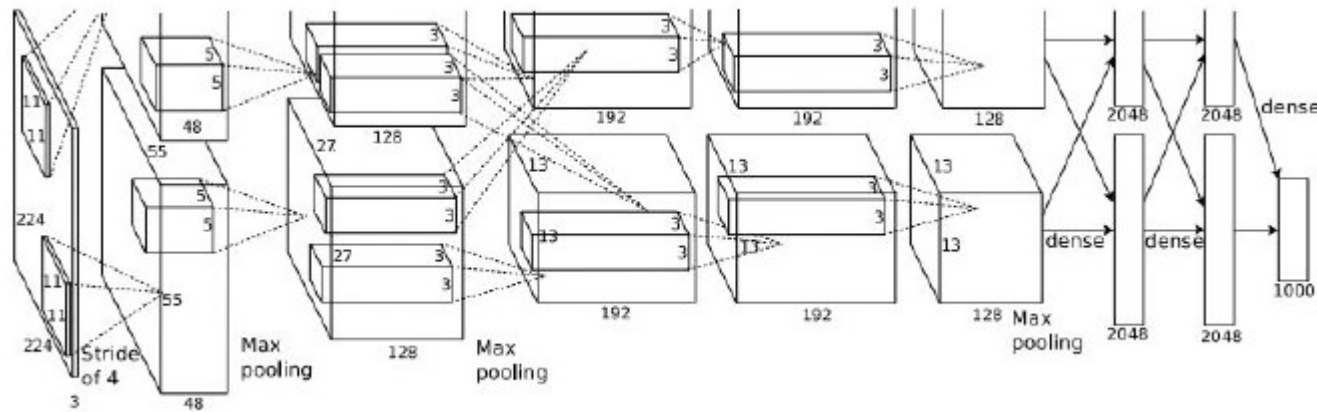


Conv filters were 5x5, applied at stride 1  
Subsampling (Pooling) layers were 2x2 applied at stride 2  
i.e. architecture is [CONV-POOL-CONV-POOL-FC-FC]

Slide Credit: Fei Fei Li et. al.

# AlexNet

Krizhevsky et. al. 2012



**Architecture:**

CONV1

MAX POOL1

NORM1

CONV2

MAX POOL2

NORM2

CONV3

CONV4

CONV5

Max POOL3

FC6

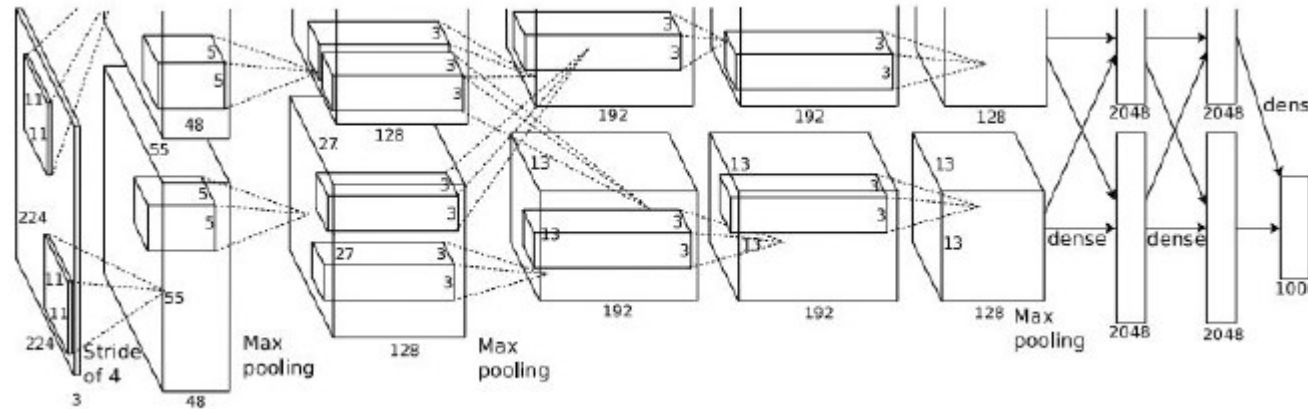
FC7

FC8

Slide Credit: Fei Fei Li et. al.

# AlexNet

Krizhevsky et. al. 2012



Architecture:

CONV1

MAX POOL1

NORM1

CONV2

MAX POOL2

NORM2

CONV3

CONV4

CONV5

Max POOL3

FC6

FC7

FC8

Input: 227x227x3 images

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

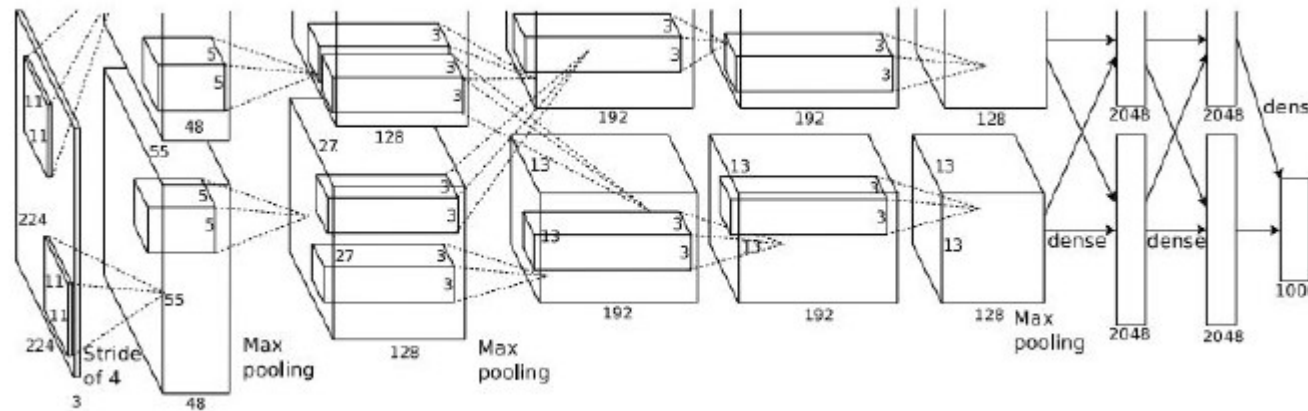
=>

Q: what is the output volume size? Hint:  $(227-11)/4+1 = 55$

Slide Credit: Fei Fei Li et. al.

# AlexNet

Krizhevsky et. al. 2012



Input: 227x227x3 images

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

=>

Output volume **[55x55x96]**

Q: What is the total number of parameters in this layer?

Architecture:

CONV1

MAX POOL1

NORM1

CONV2

MAX POOL2

NORM2

CONV3

CONV4

CONV5

Max POOL3

FC6

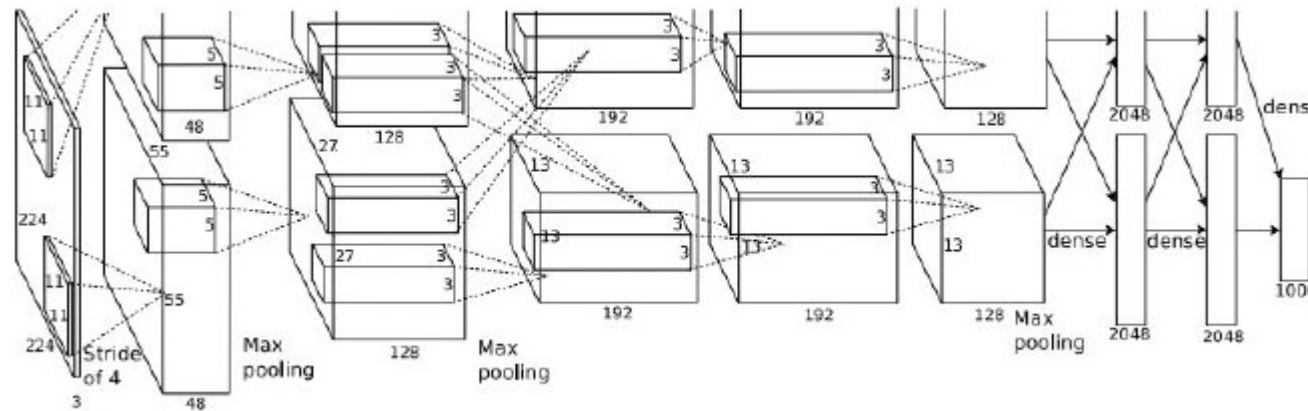
FC7

FC8

Slide Credit: Fei Fei Li et. al.

# AlexNet

Krizhevsky et. al. 2012



Architecture:

CONV1

MAX POOL1

NORM1

CONV2

MAX POOL2

NORM2

CONV3

CONV4

CONV5

Max POOL3

FC6

FC7

FC8

Input: 227x227x3 images

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

=>

Output volume **[55x55x96]**

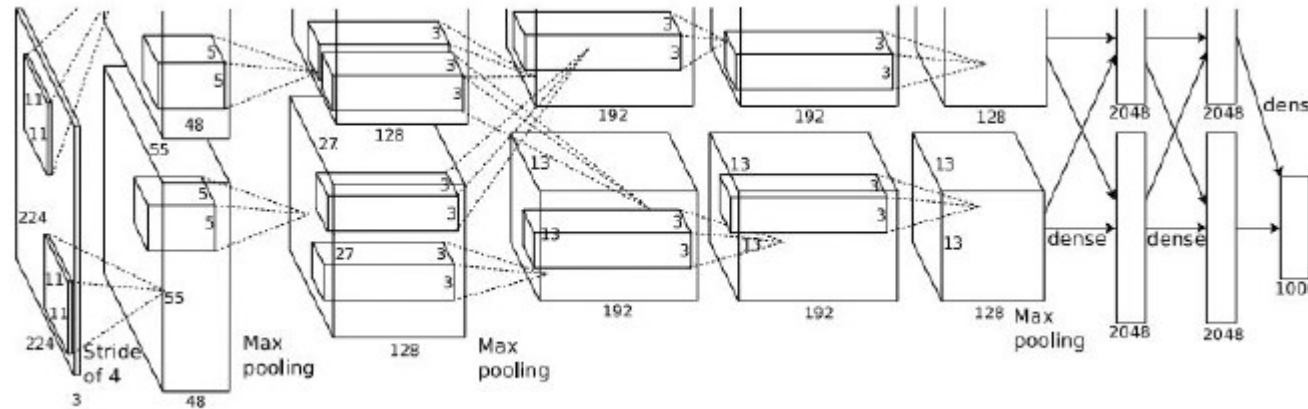
Parameters:  $(11 \times 11 \times 3) \times 96 = 35K$

Slide Credit: Fei Fei Li et. al.



# AlexNet

Krizhevsky et. al. 2012



Input: 227x227x3 images

After CONV1: 55x55x96

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

Q: what is the output volume size? Hint:  $(55-3)/2+1 = 27$

Architecture:

CONV1

MAX POOL1

NORM1

CONV2

MAX POOL2

NORM2

CONV3

CONV4

CONV5

Max POOL3

FC6

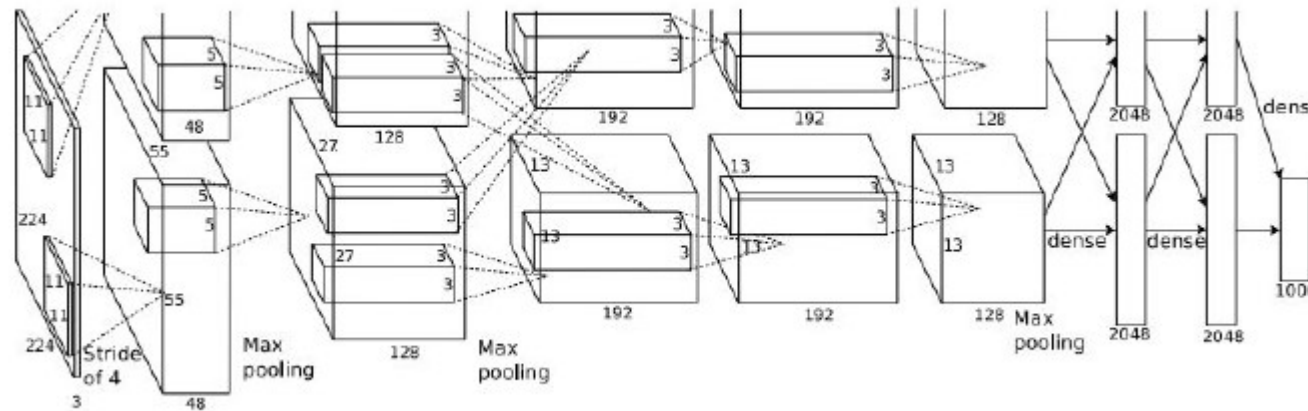
FC7

FC8

Slide Credit: Fei Fei Li et. al.

# AlexNet

Krizhevsky et. al. 2012



Input: 227x227x3 images

After CONV1: 55x55x96

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

Output volume: 27x27x96

Q: what is the number of parameters in this layer?

Architecture:

CONV1

MAX POOL1

NORM1

CONV2

MAX POOL2

NORM2

CONV3

CONV4

CONV5

Max POOL3

FC6

FC7

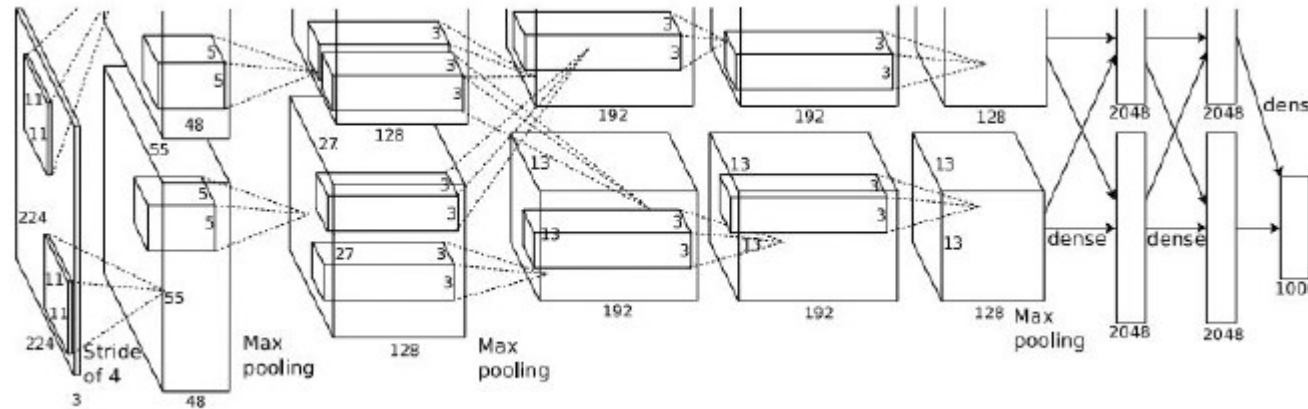
FC8

Slide Credit: Fei Fei Li et. al.



# 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!

Architecture:

CONV1

MAX POOL1

NORM1

CONV2

MAX POOL2

NORM2

CONV3

CONV4

CONV5

Max POOL3

FC6

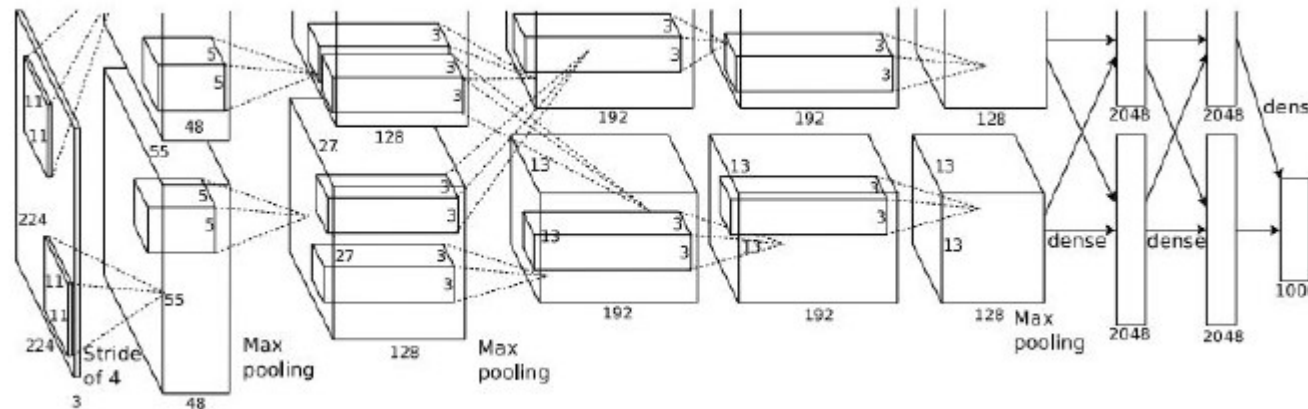
FC7

FC8

Slide Credit: Fei Fei Li et. al.

# 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)

Architecture:

CONV1

MAX POOL1

NORM1

CONV2

MAX POOL2

NORM2

CONV3

CONV4

CONV5

Max POOL3

FC6

FC7

FC8

Slide Credit: Fei Fei Li et. al.

# 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

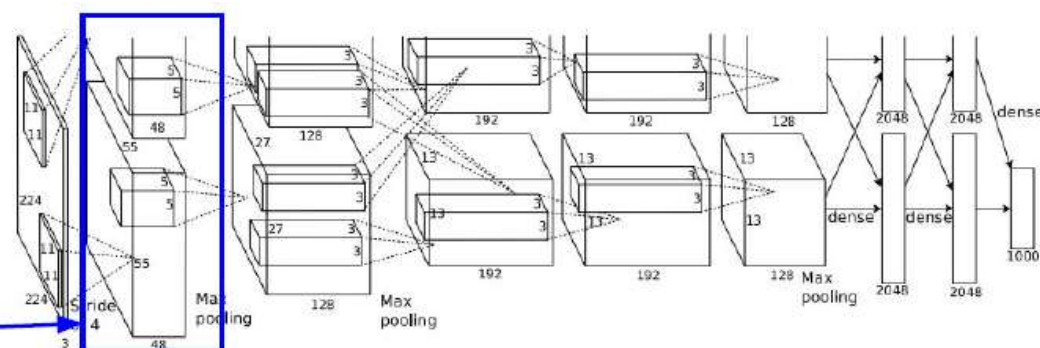
[13x13x256] CONV5: 256 3x3 filters at stride 1, pad 1

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

[4096] FC6: 4096 neurons

[4096] FC7: 4096 neurons

[1000] FC8: 1000 neurons (class scores)



[55x55x48] x 2

Historical note: Trained on GTX 580 GPU with only 3 GB of memory. Network spread across 2 GPUs, half the neurons (feature maps) on each GPU.

Figure copyright Alex Krizhevsky, Ilya Sutskever, and Geoffrey Hinton, 2012. Reproduced with permission.

Slide Credit: Fei Fei Li et. al.



# 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

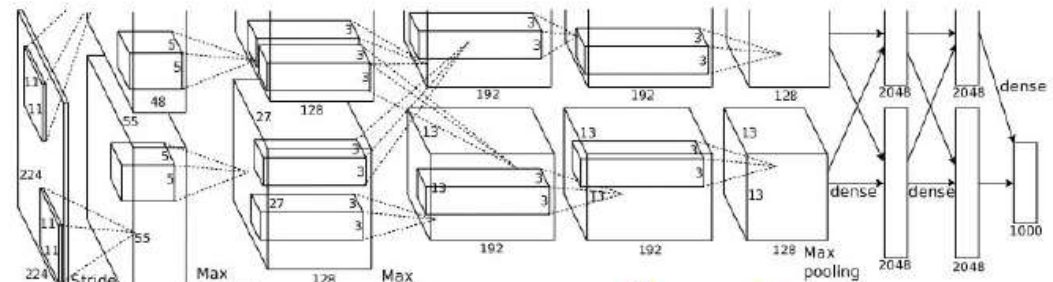
[13x13x256] **CONV5**: 256 3x3 filters at stride 1, pad 1

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

[4096] **FC6**: 4096 neurons

[4096] **FC7**: 4096 neurons

[1000] **FC8**: 1000 neurons (class scores)



**CONV1, CONV2, CONV4, CONV5:**  
Connections only with feature maps  
on same GPU

Figure copyright Alex Krizhevsky, Ilya Sutskever, and Geoffrey Hinton, 2012. Reproduced with permission.

Slide Credit: Fei Fei Li et. al.

# 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

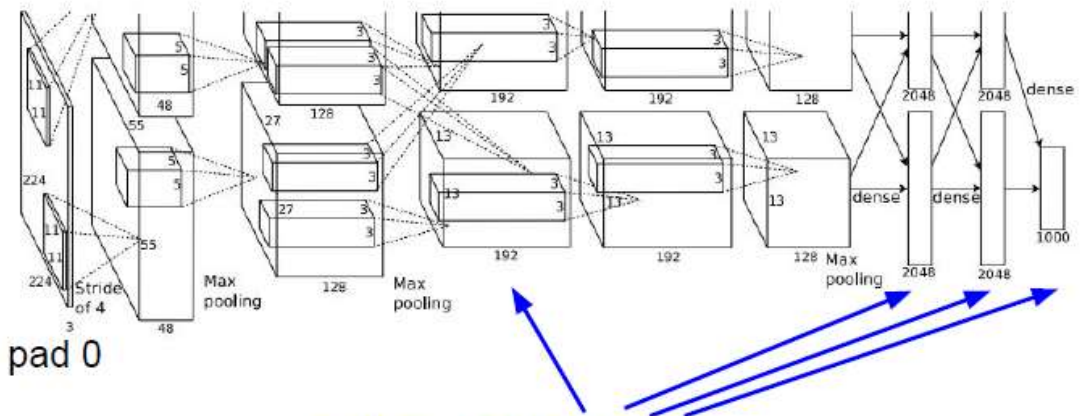
[13x13x256] **CONV5**: 256 3x3 filters at stride 1, pad 1

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

[4096] **FC6**: 4096 neurons

[4096] **FC7**: 4096 neurons

[1000] **FC8**: 1000 neurons (class scores)



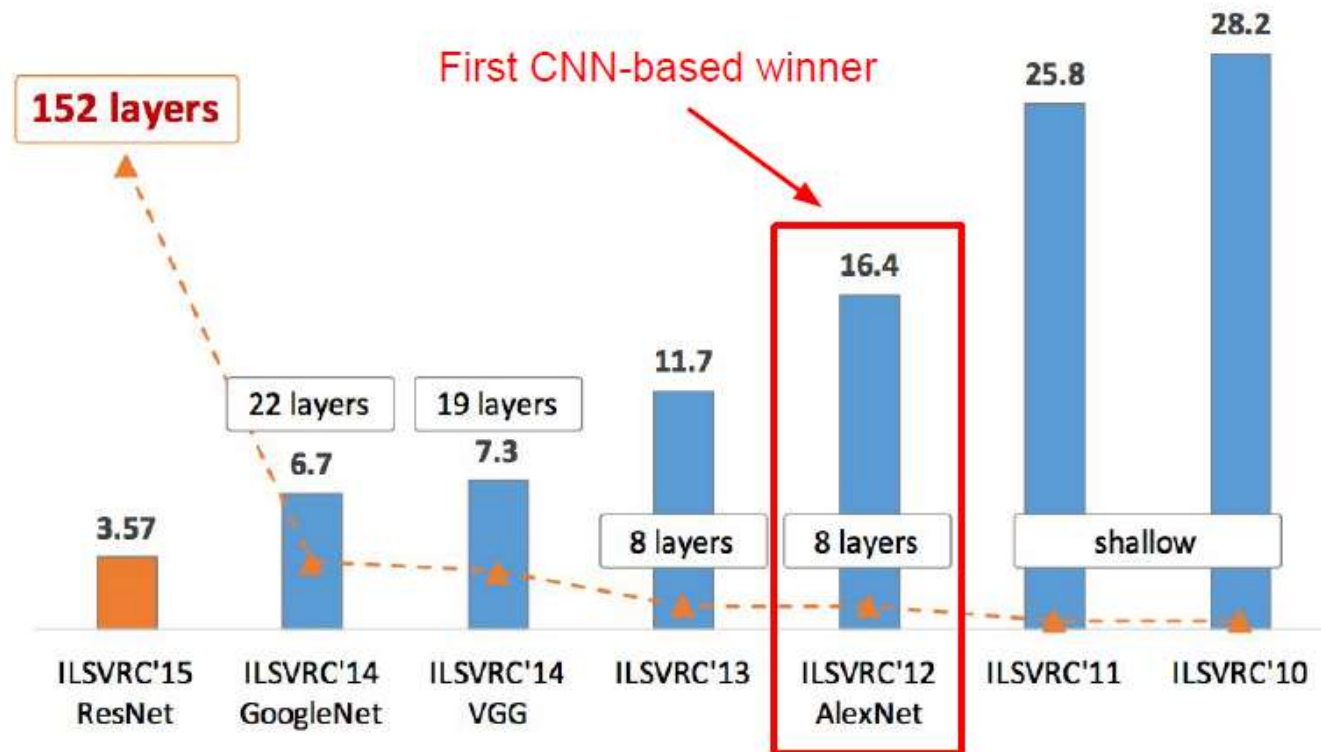
**CONV3, FC6, FC7, FC8:**  
Connections with all feature maps in preceding layer, communication across GPUs

Figure copyright Alex Krizhevsky, Ilya Sutskever, and Geoffrey Hinton, 2012. Reproduced with permission.

Slide Credit: Fei Fei Li et. al.

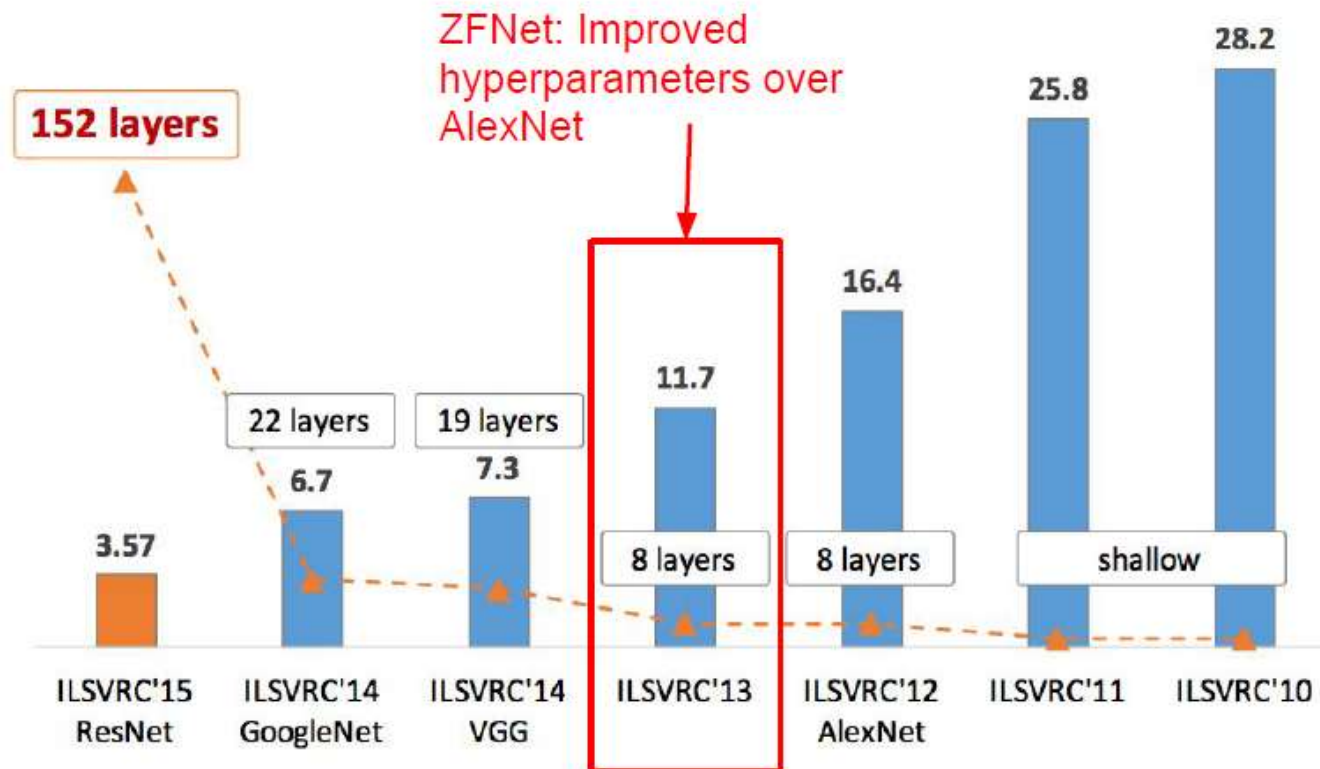


## ImageNet Large Scale Visual Recognition Challenge (ILSVRC) winners



Slide Credit: Fei Fei Li et. al.

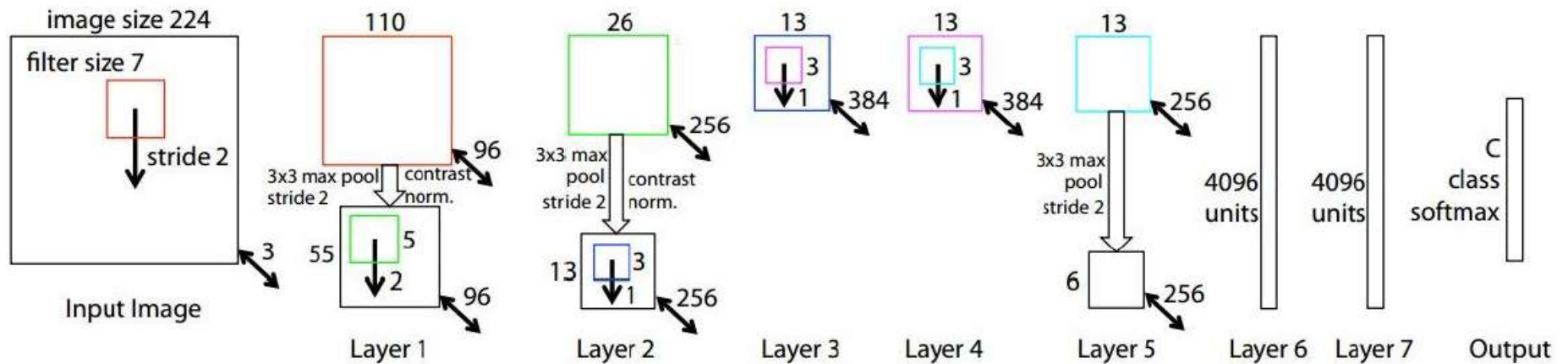
# ImageNet Large Scale Visual Recognition Challenge (ILSVRC) winners



Slide Credit: Fei Fei Li et. al.

# ZFNet

[Zeiler and Fergus, 2013]



TODO: remake figure

AlexNet but:

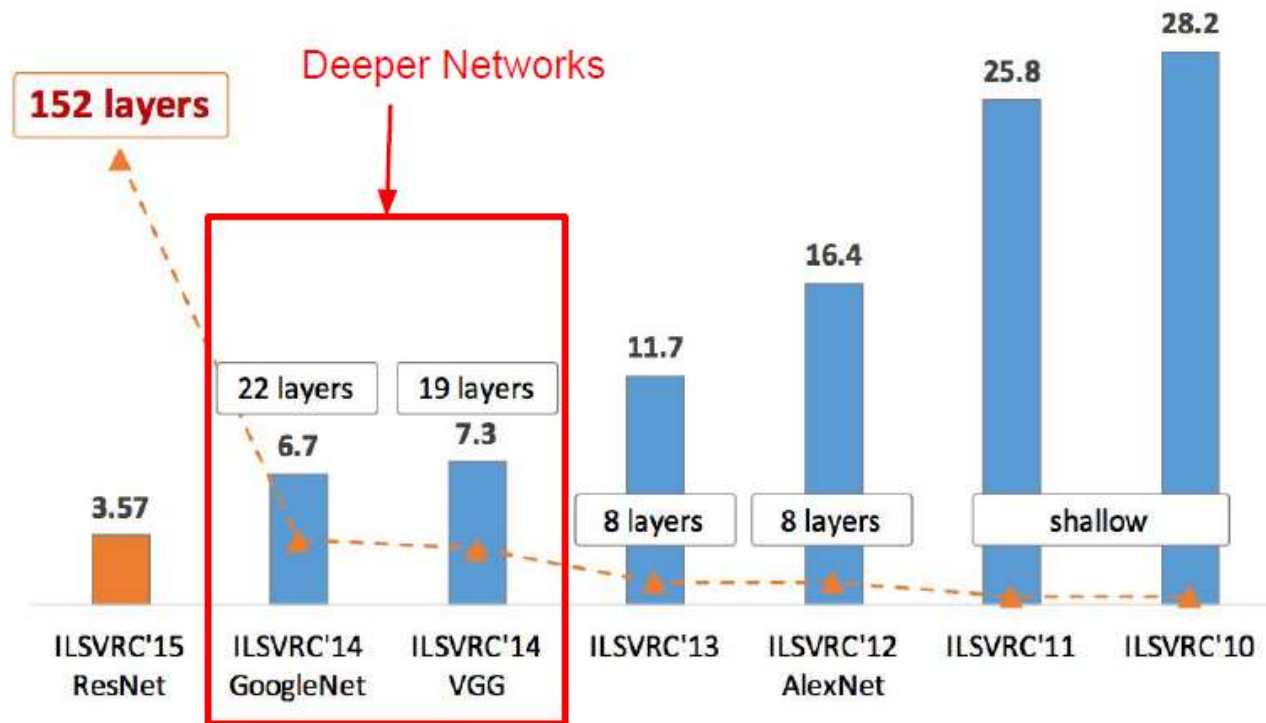
CONV1: change from (11x11 stride 4) to (7x7 stride 2)

CONV3,4,5: instead of 384, 384, 256 filters use 512, 1024, 512

ImageNet top 5 error: 16.4% -> 11.7%

Slide Credit: Fei Fei Li et. al.

## ImageNet Large Scale Visual Recognition Challenge (ILSVRC) winners



Slide Credit: Fei Fei Li et. al.

# VGG Net

Simonyan and Zisserman 2014

Small filters, Deeper networks

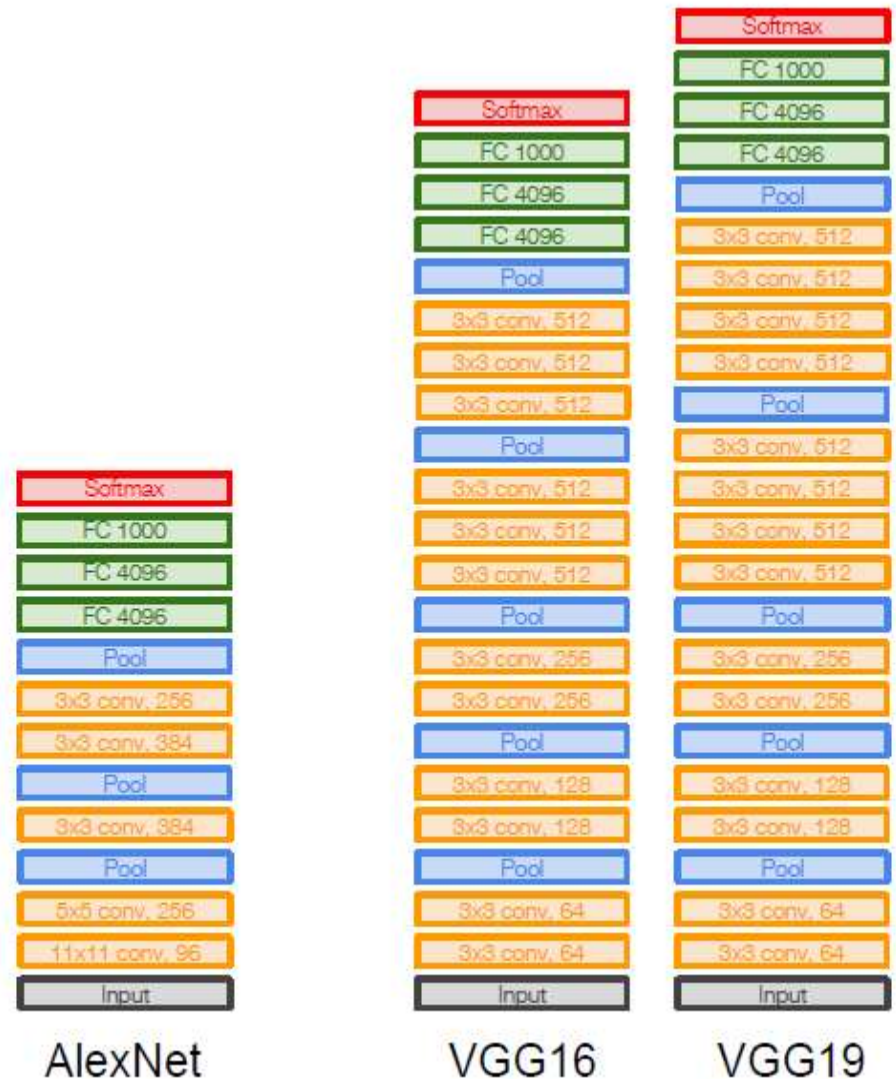
8 layers (AlexNet)

-> 16 - 19 layers (VGG16Net)

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

11.7% top 5 error in ILSVRC'13  
(ZFNet)

-> 7.3% top 5 error in ILSVRC'14



Slide Credit: Fei Fei Li et. al.



# VGG Net

Simonyan and Zisserman 2014

Q: Why use smaller filters? (3x3 conv)



Slide Credit: Fei Fei Li et. al.

# VGG Net

Simonyan and Zisserman 2014

Q: Why use smaller filters? (3x3 conv)

Stack of three 3x3 conv (stride 1) layers has same **effective receptive field** as one 7x7 conv layer

Q: What is the effective receptive field of three 3x3 conv (stride 1) layers?



Slide Credit: Fei Fei Li et. al.

# VGG Net

Simonyan and Zisserman 2014

Q: Why use smaller filters? (3x3 conv)

Stack of three 3x3 conv (stride 1) layers has same **effective receptive field** as one 7x7 conv layer

[7x7]



Slide Credit: Fei Fei Li et. al.



# VGG Net

Simonyan and Zisserman 2014

Q: Why use smaller filters? (3x3 conv)

Stack of three 3x3 conv (stride 1) layers has same **effective receptive field** as one 7x7 conv layer

But deeper, more non-linearities

And fewer parameters:  $3 * (3^2 C^2)$  vs.  $7^2 C^2$  for C channels per layer



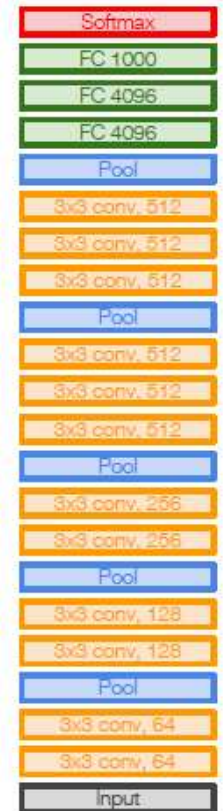
Slide Credit: Fei Fei Li et. al.

# Simonyan and Zisserman 2014

```

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

```



VGG16

TOTAL memory:  $24M * 4 \text{ bytes} \approx 96MB$  / image (only forward!  $\sim *2$  for bwd)

TOTAL params: 138M parameters

Slide Credit: Fei Fei Li et. al.



# VGG Net

Simonyan and Zisserman 2014

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:

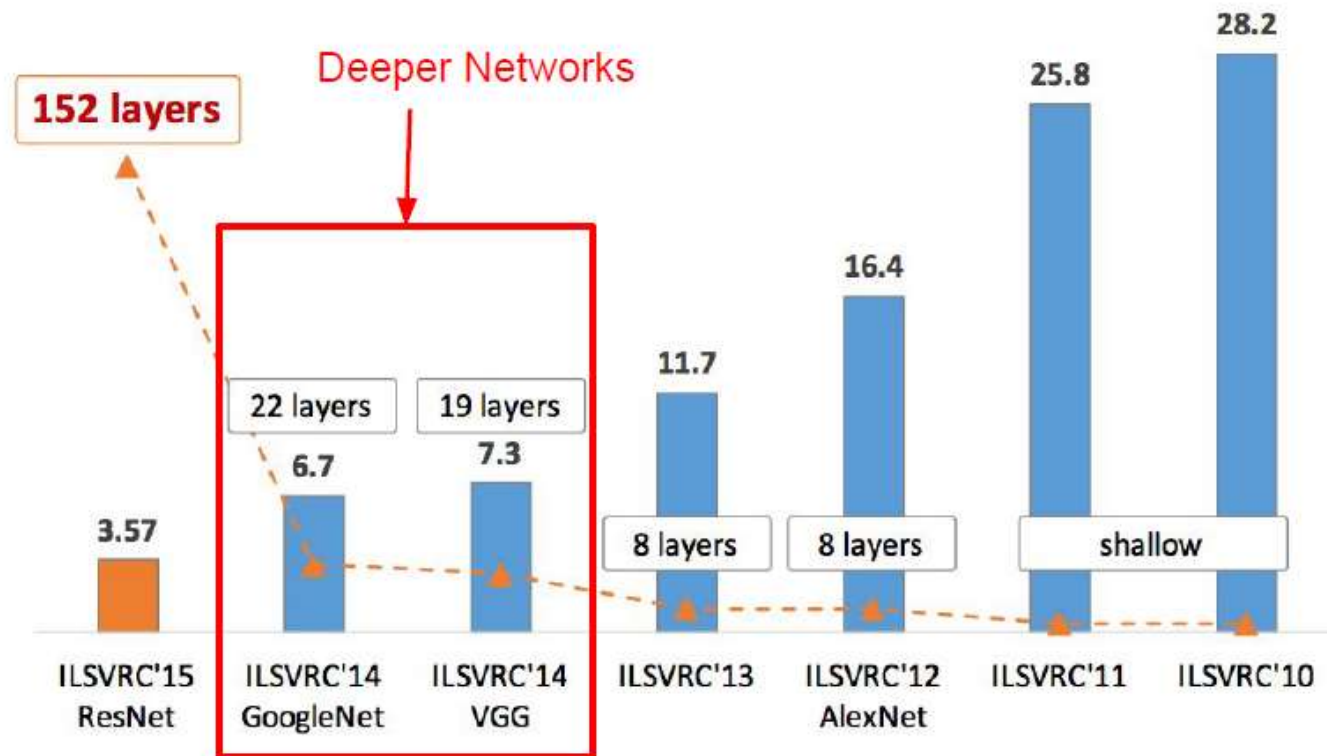
Most memory is in early CONV

Most params are in late FC

TOTAL memory:  $24M * 4 \text{ bytes} \sim 96MB / \text{image}$  (only forward!  $\sim 2$  for bwd)  
TOTAL params: 138M parameters

Slide Credit: Fei Fei Li et. al.

## ImageNet Large Scale Visual Recognition Challenge (ILSVRC) winners



Slide Credit: Fei Fei Li et. al.

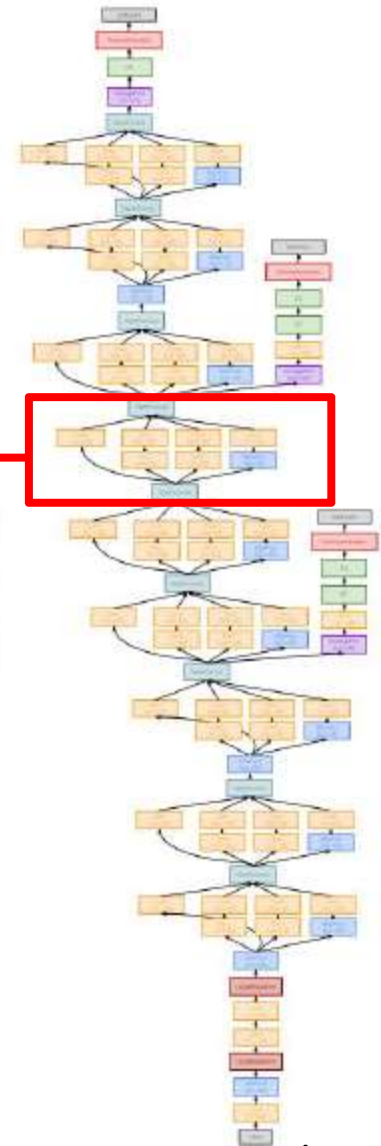
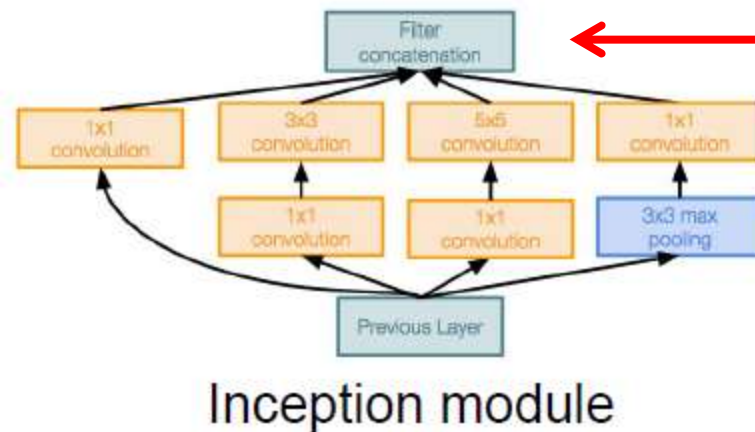
# Google Net

Szegedy et. al. 2014

Deeper networks, with computational efficiency

- 22 layers
- Efficient “Inception” module
- No FC layers
- Only 5 million parameters!  
12x less than AlexNet
- ILSVRC’14 classification winner  
(6.7% top 5 error)

“Inception module”: design a good local network topology (network within a network) and then stack these modules on top of each other

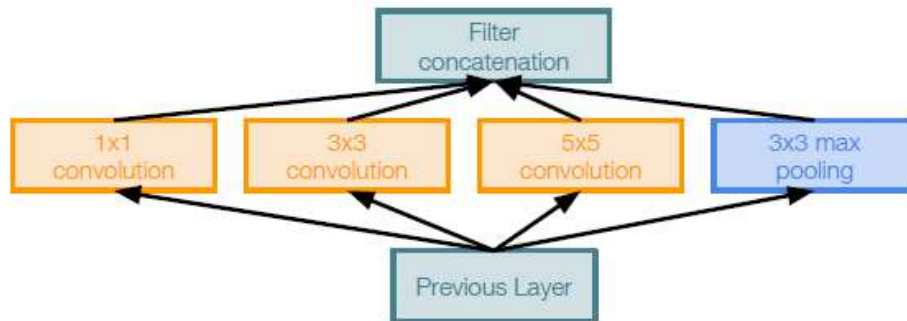


Slide Credit: Fei Fei Li et. al.



# Google Net

Szegedy et. al. 2014



Naive Inception module

Apply parallel filter operations on the input from previous layer:

- Multiple receptive field sizes for convolution (1x1, 3x3, 5x5)
- Pooling operation (3x3)

Concatenate all filter outputs together depth-wise

Q: What is the problem with this?  
[Hint: Computational complexity]

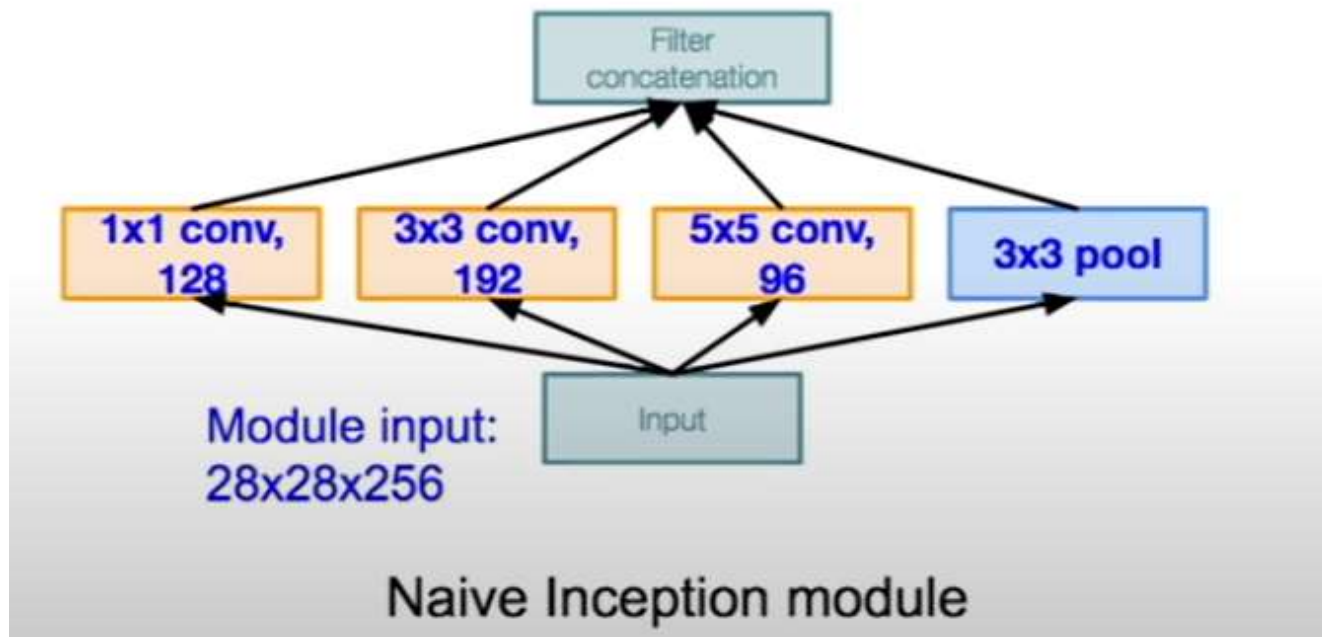
# Google Net

Szegedy et. al. 2014

Q: What is the problem with this?  
[Hint: Computational complexity]

Example:

Q1: What is the output size of the  
1x1 conv, with 128 filters?





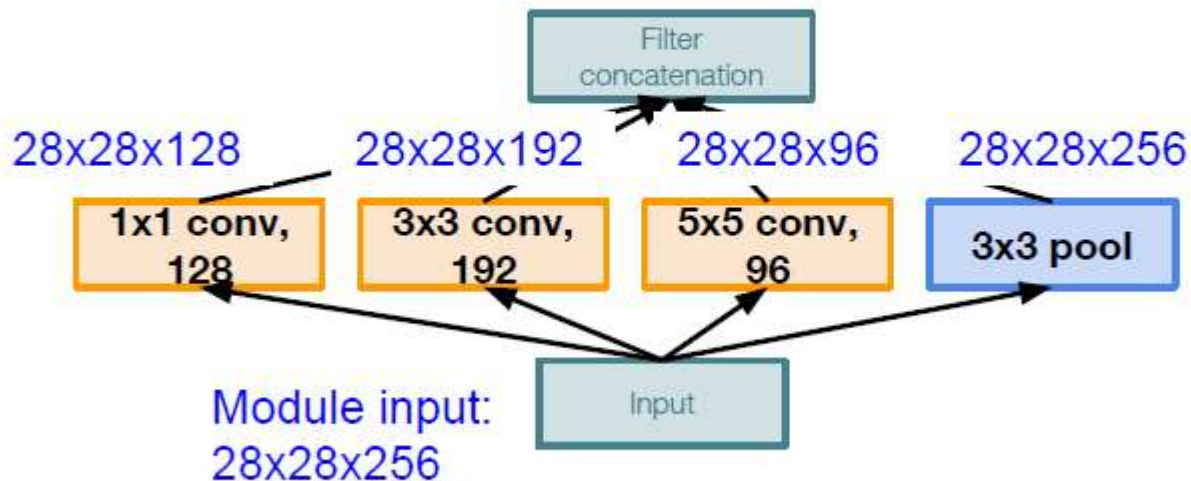
# Google Net

Szegedy et. al. 2014

Q: What is the problem with this?  
[Hint: Computational complexity]

Example:

Q2: What are the output sizes of  
all different filter operations?



Naive Inception module

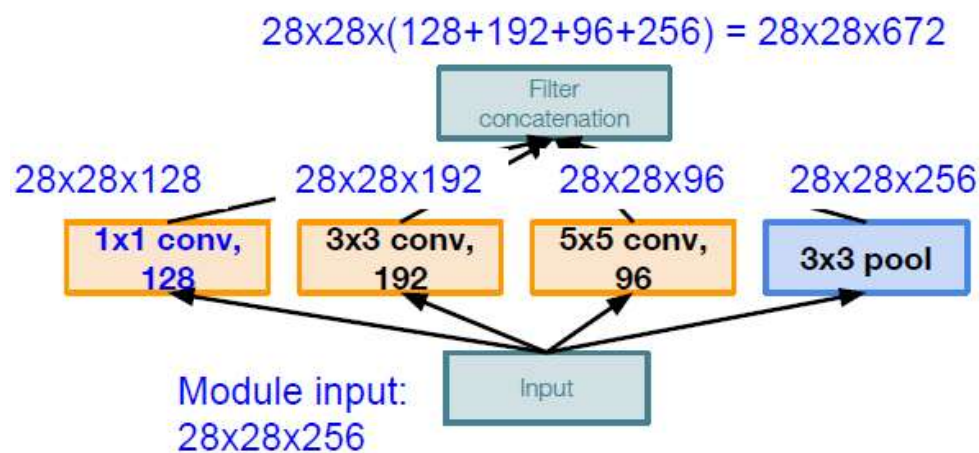
# Google Net

Szegedy et. al. 2014

Q: What is the problem with this?  
[Hint: Computational complexity]

Example:

Q3: What is output size after filter concatenation?



Naive Inception module

## Conv Ops:

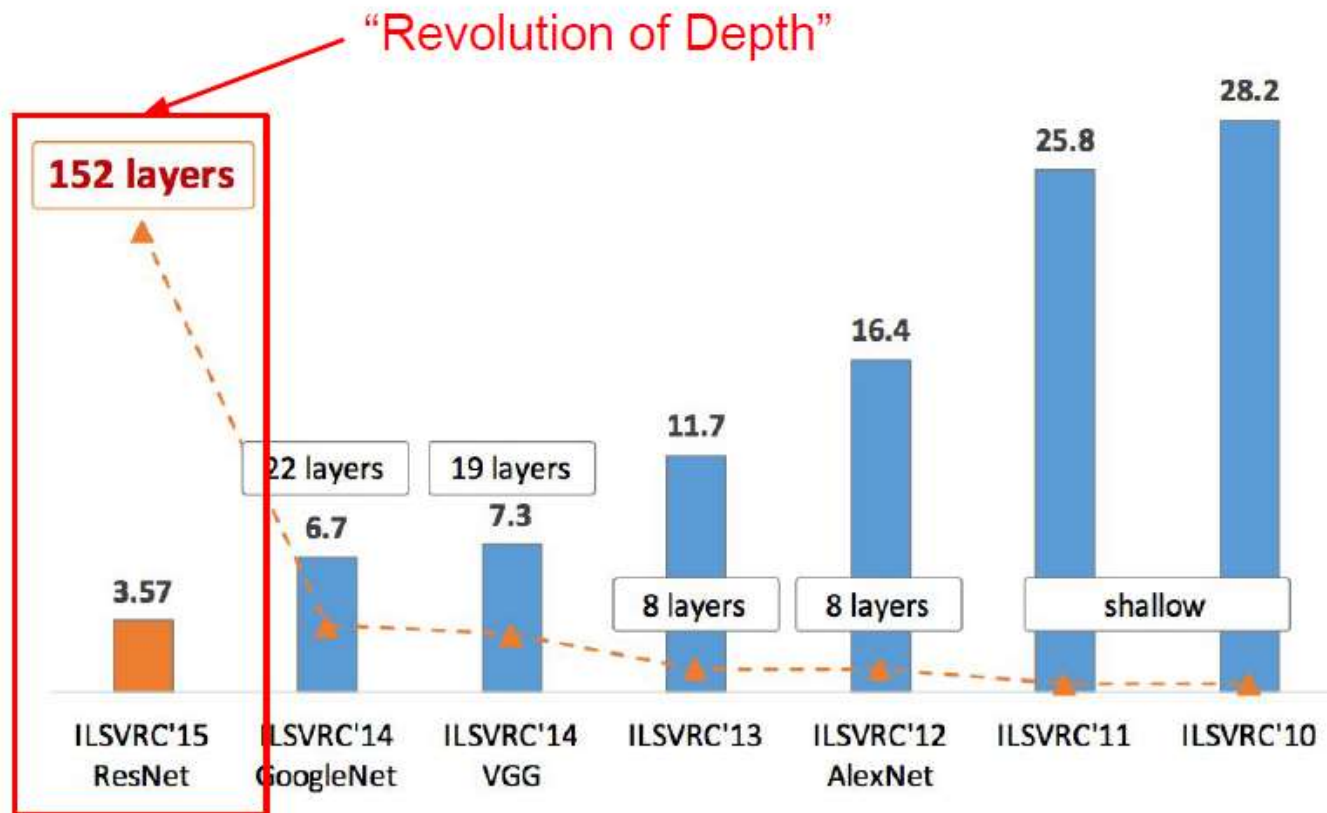
[1x1 conv, 128]  $28 \times 28 \times 128 \times 1 \times 1 \times 256$   
[3x3 conv, 192]  $28 \times 28 \times 192 \times 3 \times 3 \times 256$   
[5x5 conv, 96]  $28 \times 28 \times 96 \times 5 \times 5 \times 256$

**Total: 854M ops**

Very expensive compute

Pooling layer also preserves feature depth, which means total depth after concatenation can only grow at every layer!

## ImageNet Large Scale Visual Recognition Challenge (ILSVRC) winners



# What does 1x1 convolution do?

1	2	3	6	5	8
3	5	5	1	3	4
2	1	3	4	9	3
4	7	8	5	7	9
1	5	3	7	4	8
5	4	9	8	3	5

6 × 6

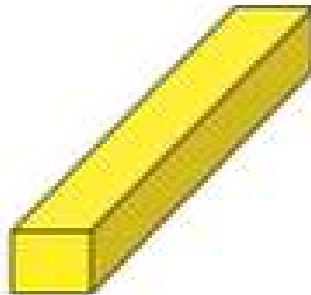
\*

2

=




\*

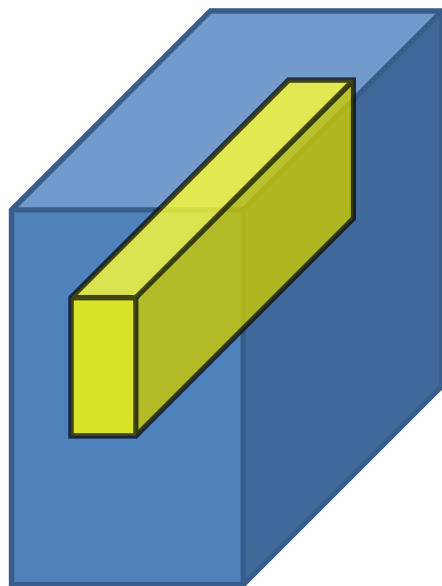


=


6 × 6 × # filters

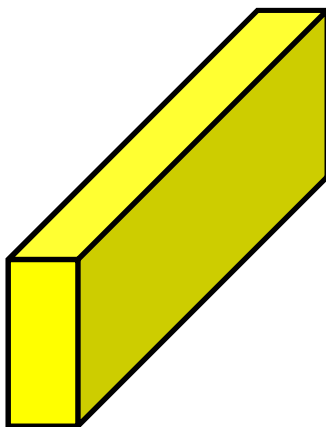


# 1X1 Convolutions

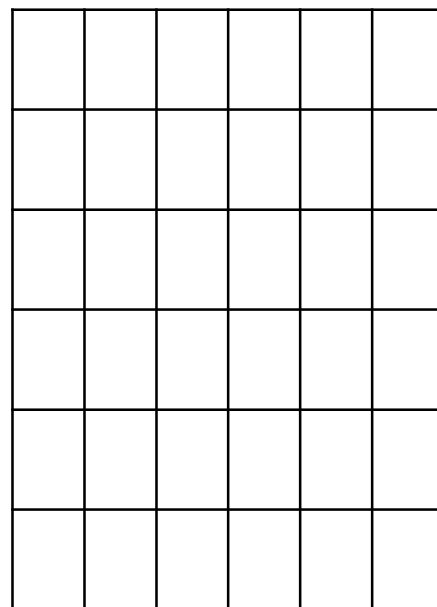


$6 \times 6 \times 32$

\*

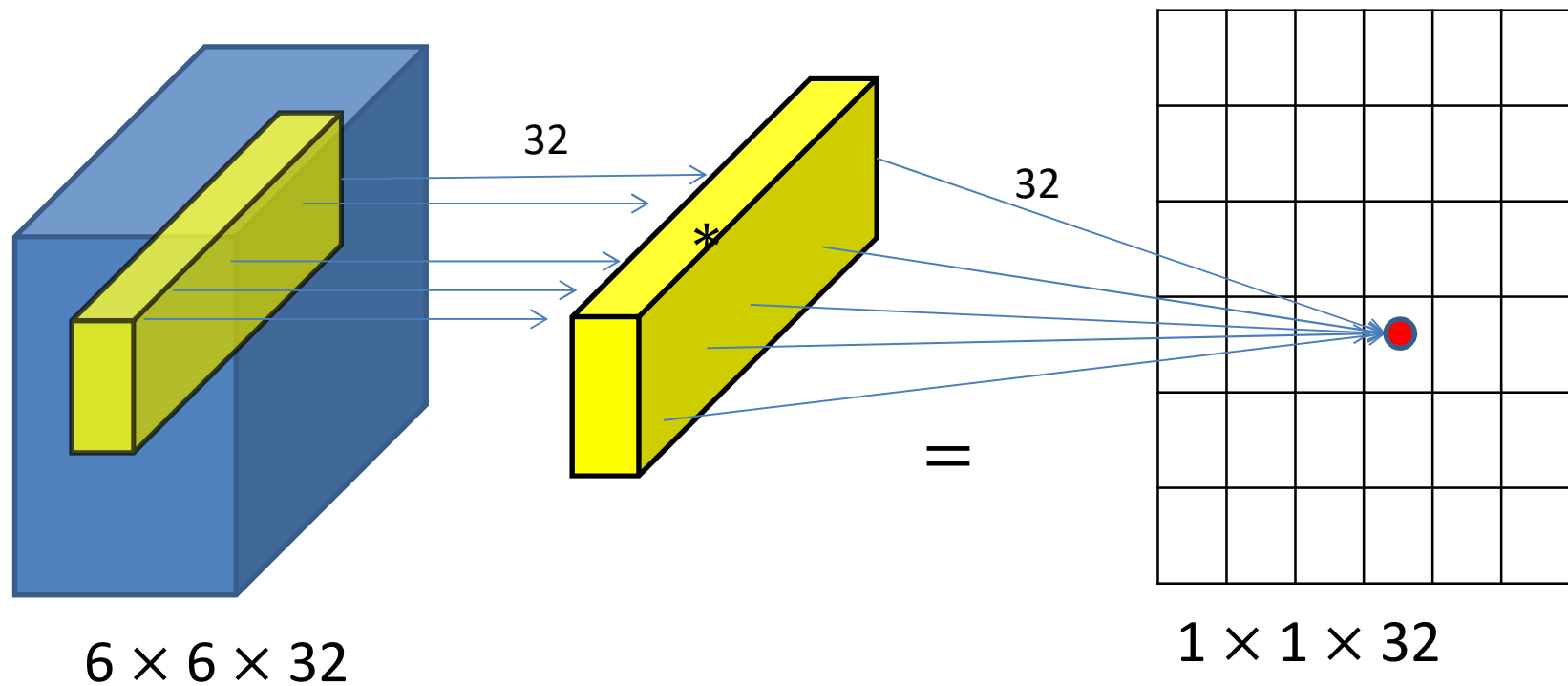


=



$1 \times 1 \times 32$

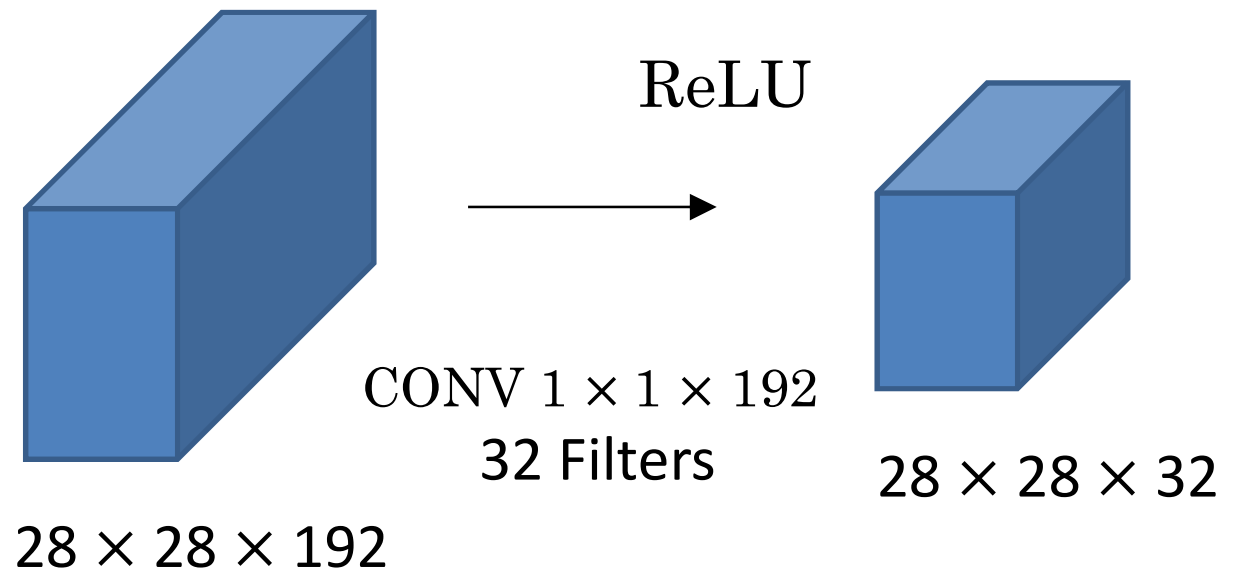
# 1X1 Convolutions



Pixel values from each of the 32 channels are multiplied with corresponding convolution kernel coefficients and are series summed to get one pixel of response matrix.

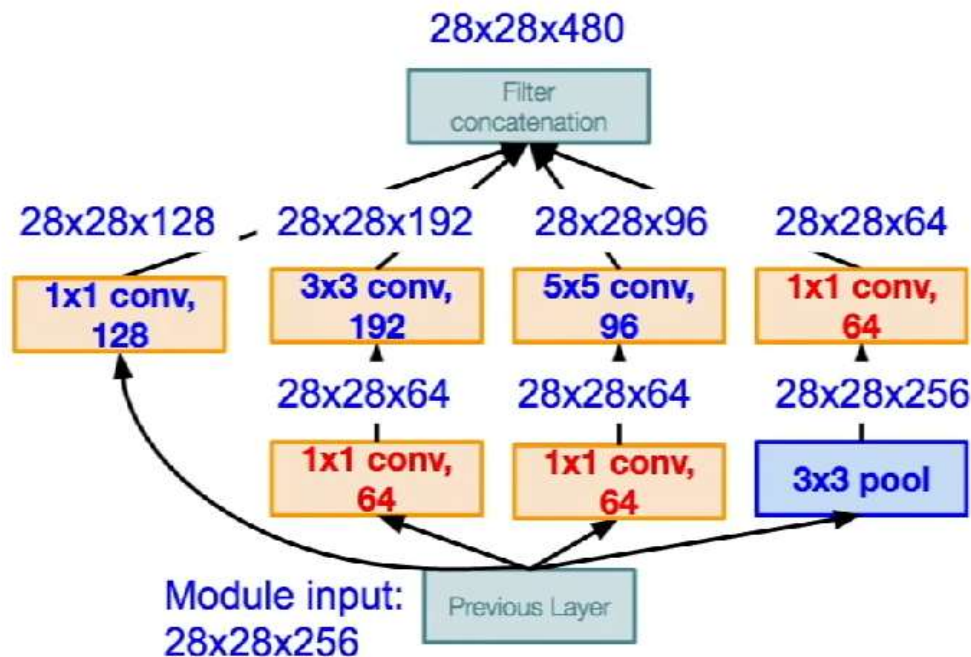
# 1X1 Convolutions

1X1 filters are often used to reduce the dimensionality of a layer



# Case Study: GoogLeNet

[Szegedy et al., 2014]



Inception module with dimension reduction

Using same parallel layers as naive example, and adding “1x1 conv, 64 filter” bottlenecks:

## Conv Ops:

[1x1 conv, 64] 28x28x64x1x1x256  
[1x1 conv, 64] 28x28x64x1x1x256  
[1x1 conv, 128] 28x28x128x1x1x256  
[3x3 conv, 192] 28x28x192x3x3x64  
[5x5 conv, 96] 28x28x96x5x5x64  
[1x1 conv, 64] 28x28x64x1x1x256

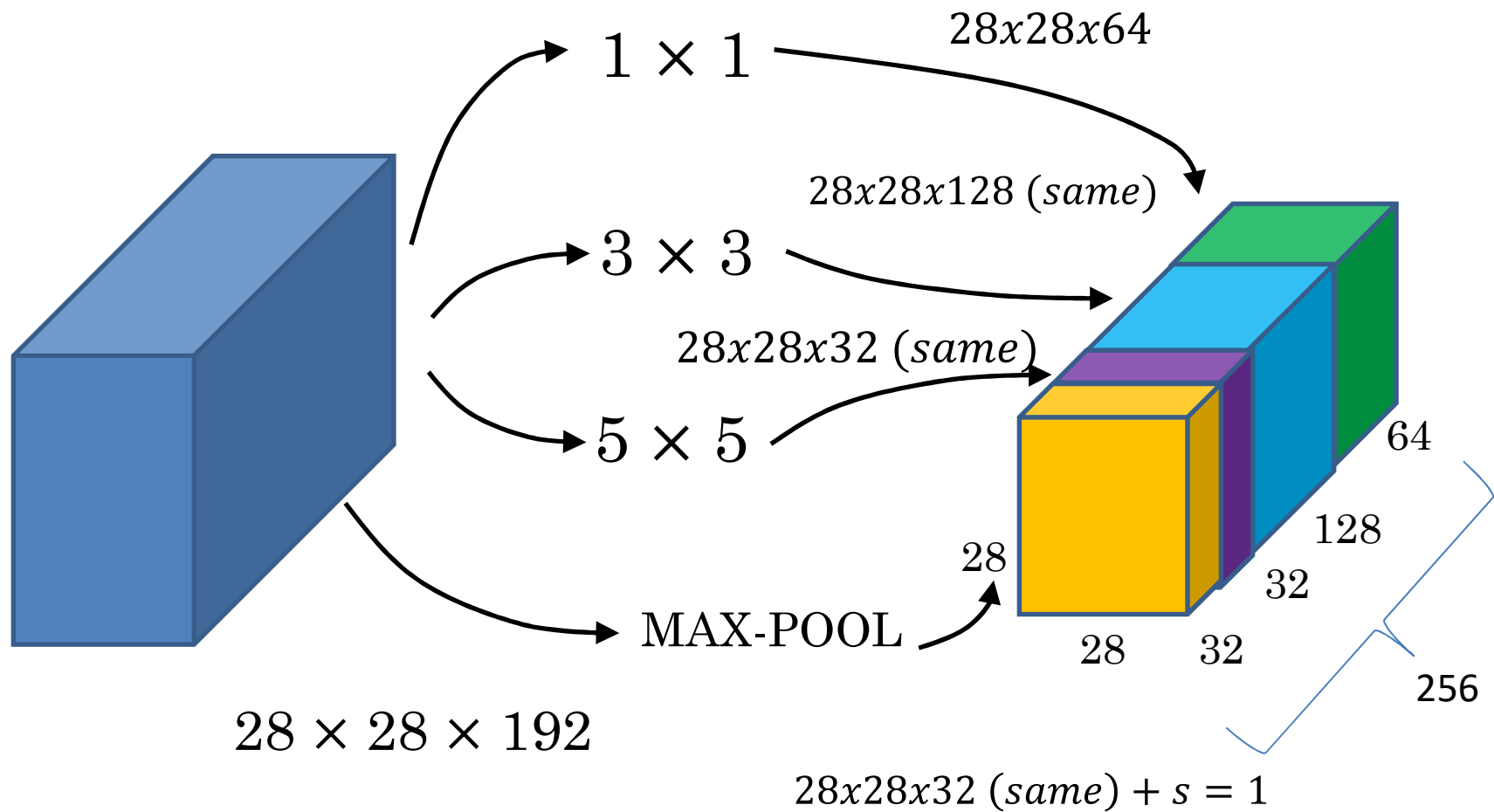
**Total: 358M ops**

Compared to 854M ops for naive version

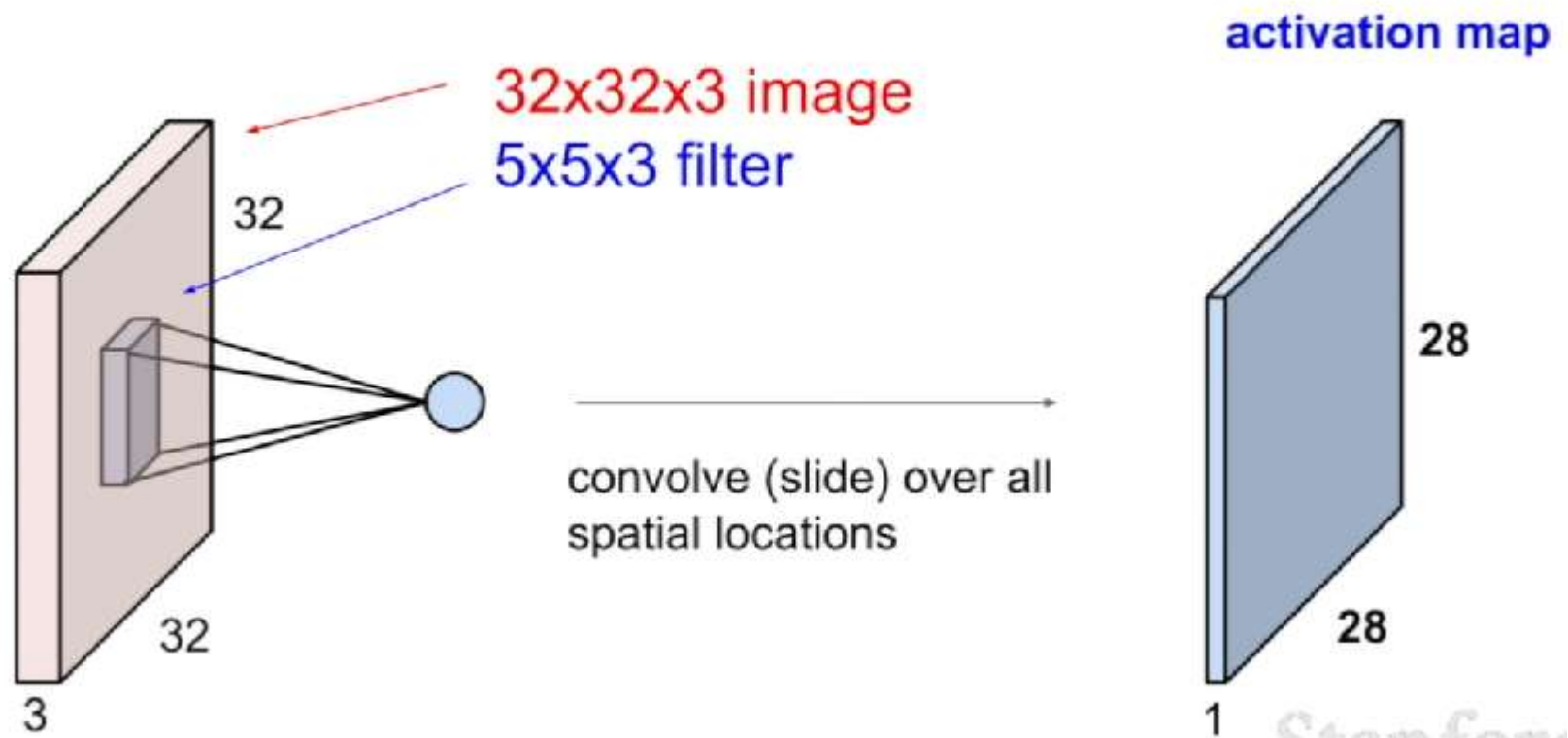
Bottleneck can also reduce depth after pooling layer



# Inception Network



[Szegedy et al. 2014. Going deeper with convolutions]



7


7

7x7 input (spatially)  
assume 3x3 filter

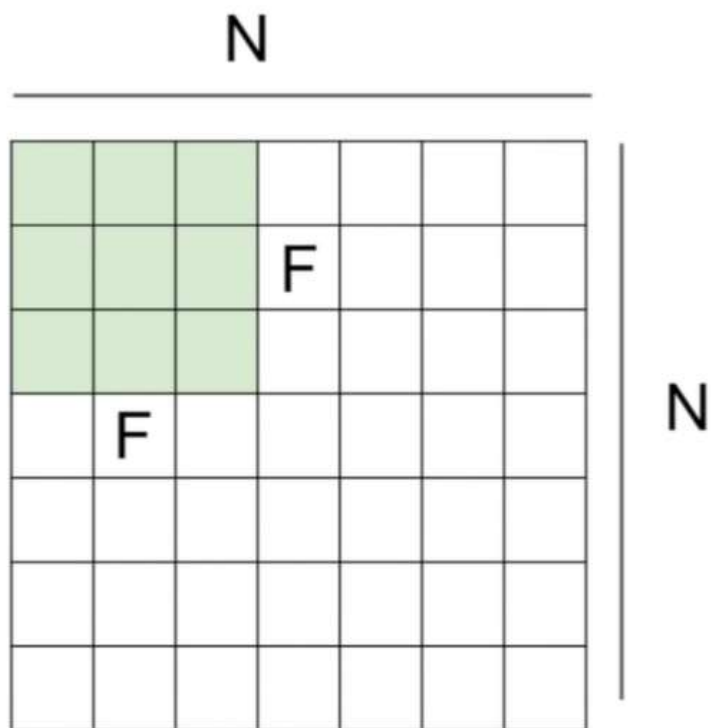
**=> 5x5 output**

7


7

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





Output size:

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

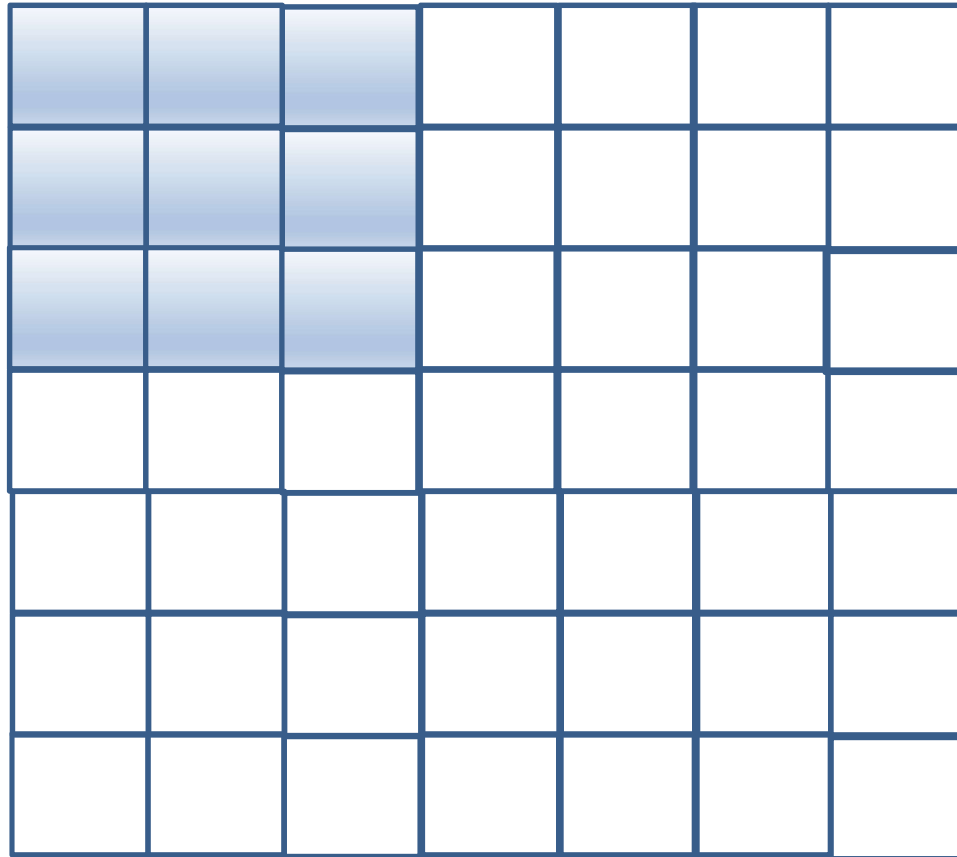
e.g.  $N = 7$ ,  $F = 3$ :

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

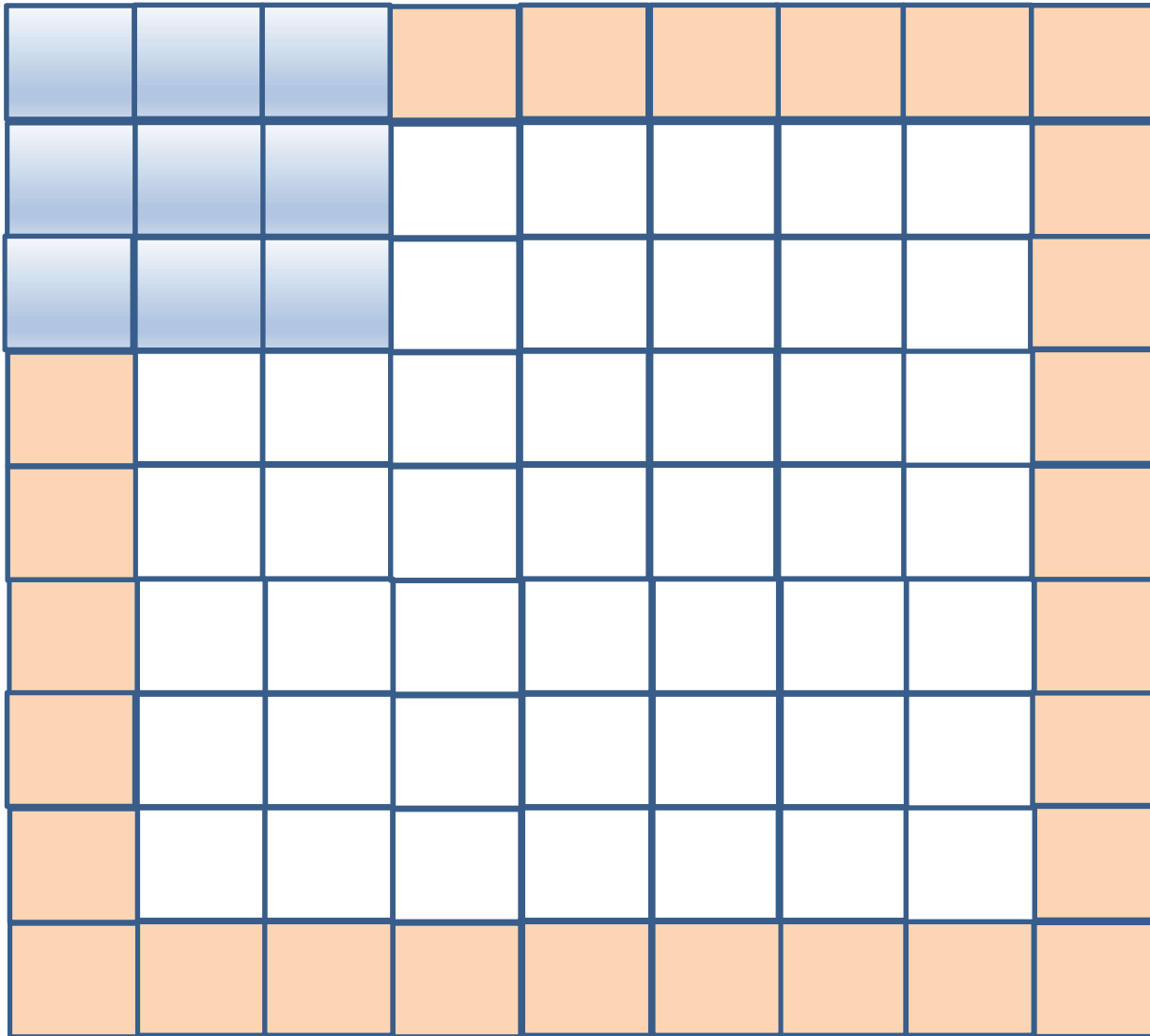
stride 2  $\Rightarrow (7 - 3)/2 + 1 = 3$

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

# Convolution Spatial Dimension



Padding:  $(9-3)/1+1=7$



Padding 1+1

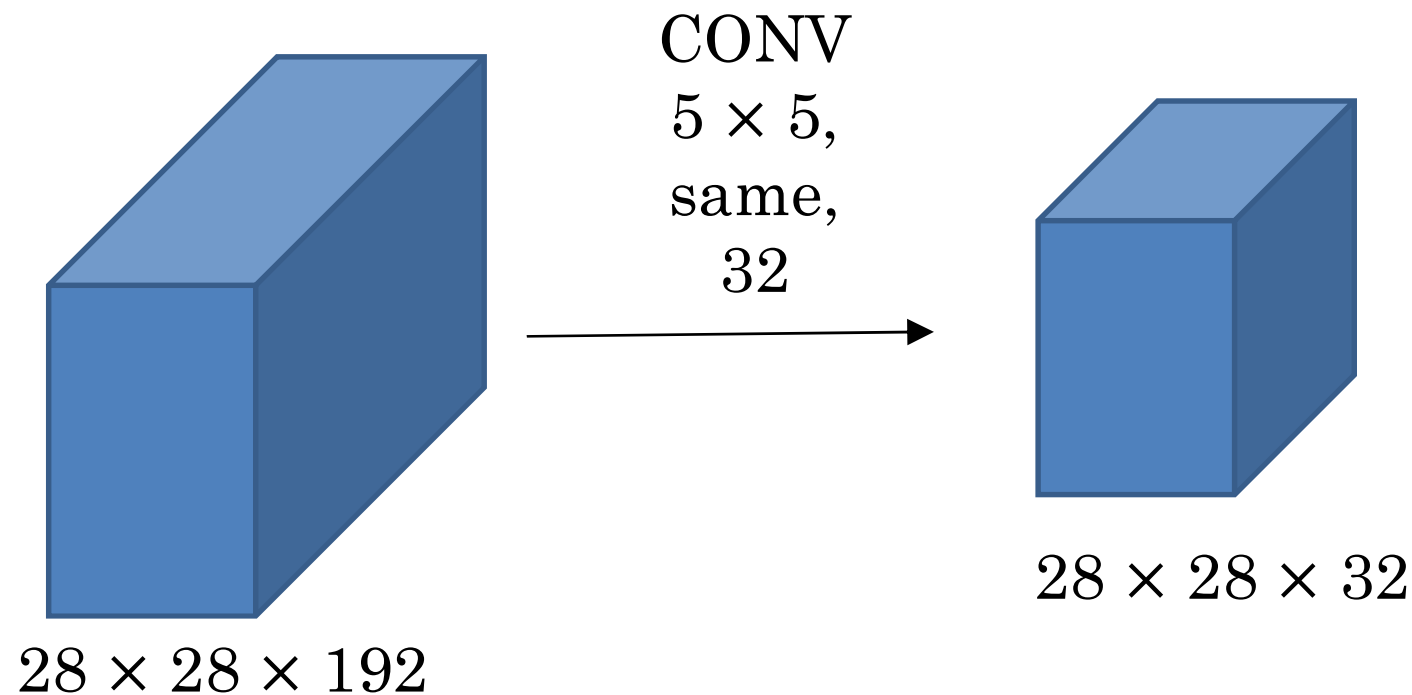
$I = 7 \times 7$

With Padding =  $9 \times 9$

Stride = 1

Spatial dimension  
remain same

# Inception Networks: Computational Cost



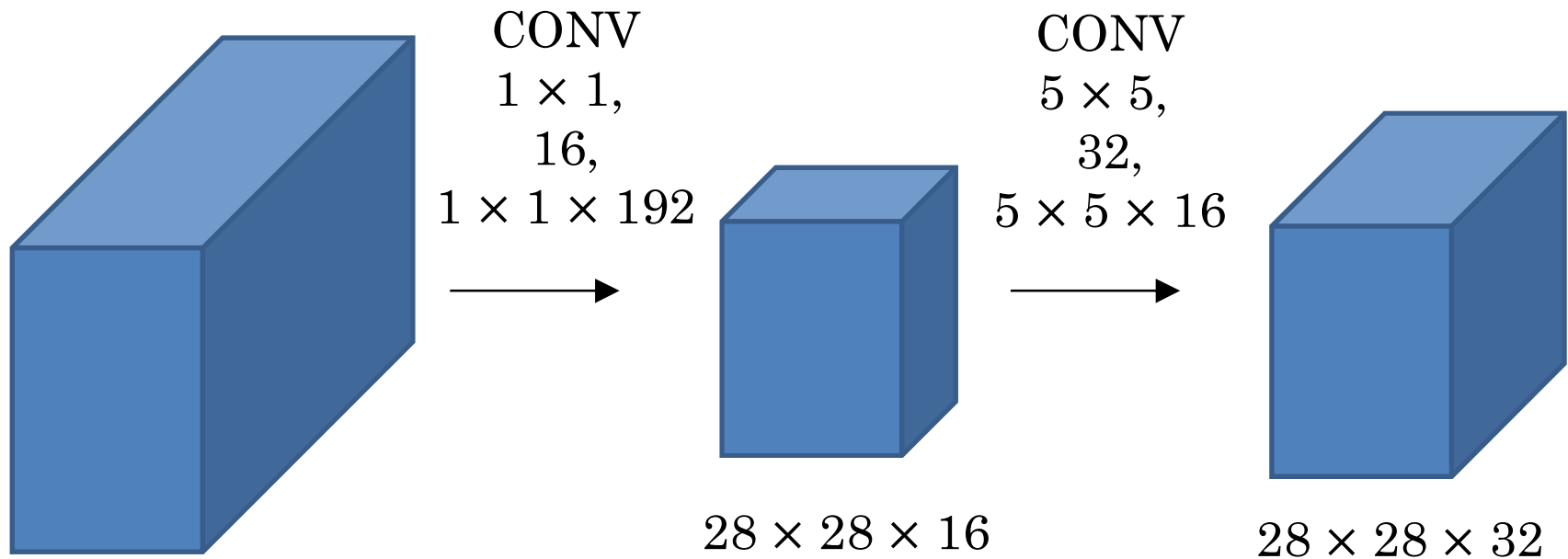
One output pixel requires  $5 \times 5 \times 192$  multiplications

There are total  $28 \times 28 \times 32$  output pixels

Total # of multiplications required =  $(5 \times 5 \times 192) \times (28 \times 28 \times 32) = \mathbf{120M}$



## How Computational cost can be reduced with 1X1 Convolutions



$28 \times 28 \times 192$

Multiplications for first layer  
 $= 28 \times 28 \times 6 \times 192 = 2.4 \text{ M}$

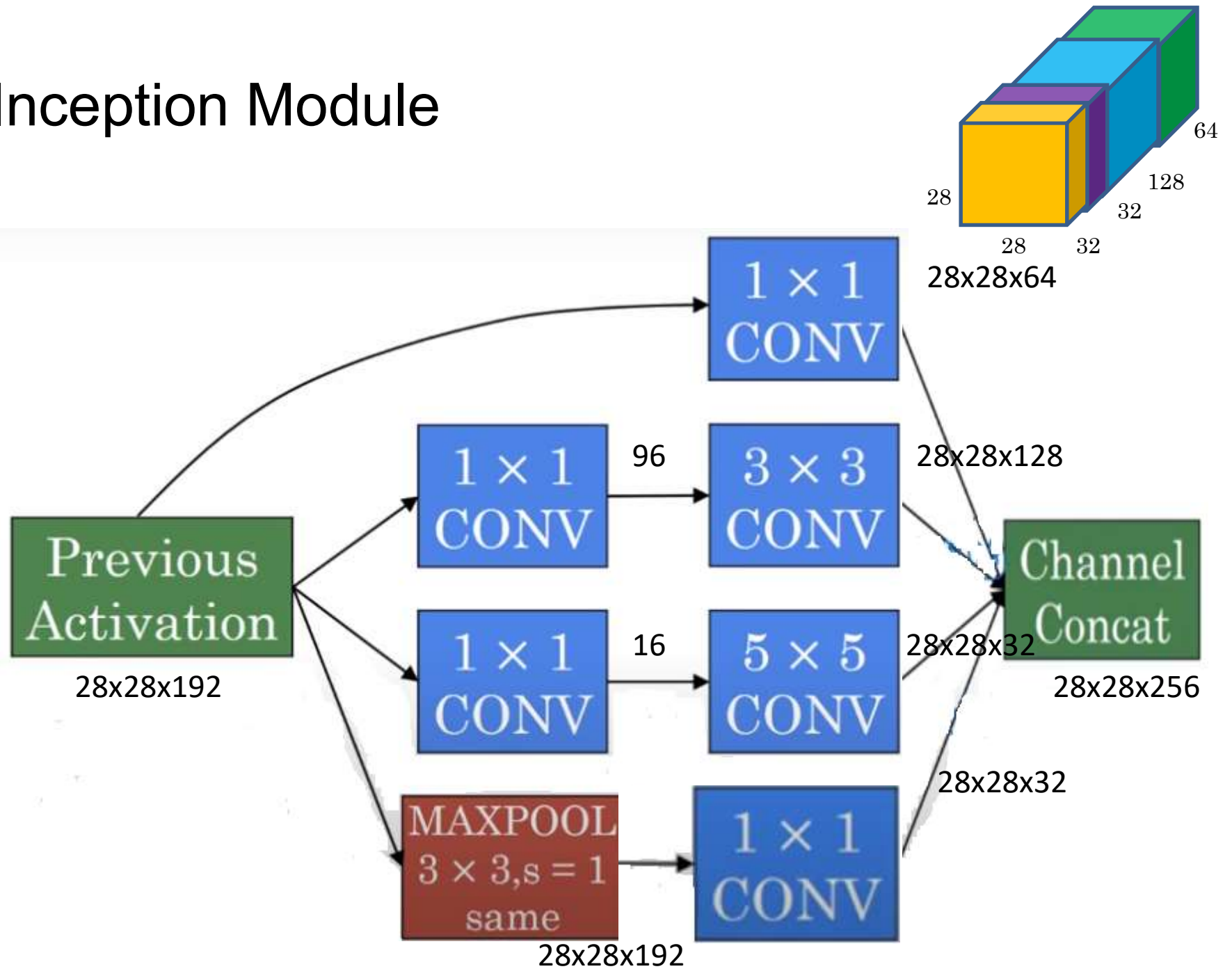
$28 \times 28 \times 16$

Multiplications for second layer  
 $= (28 \times 28 \times 32) \times (5 \times 5 \times 16) = 10 \text{ M}$

$28 \times 28 \times 32$

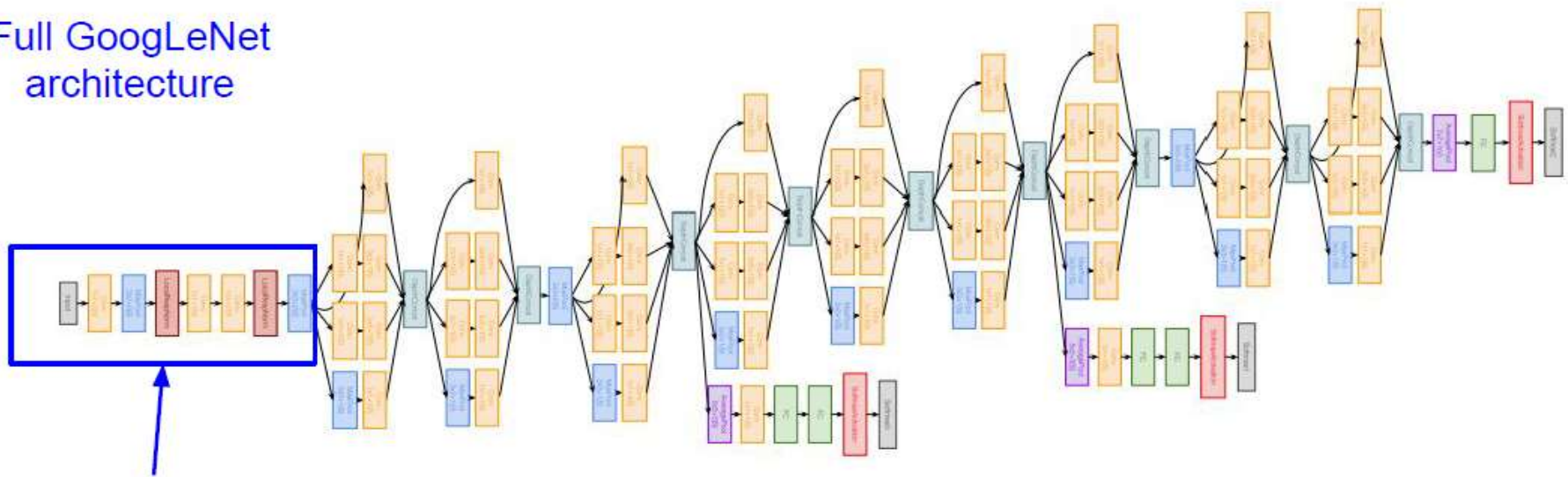
Total # of Multiplications  $= 2.4 + 10 = 12.4 \text{ M}$  which is approximately **one tenth** of the previous inception model (120 M)

# Inception Module



# Inception Network: Google Net

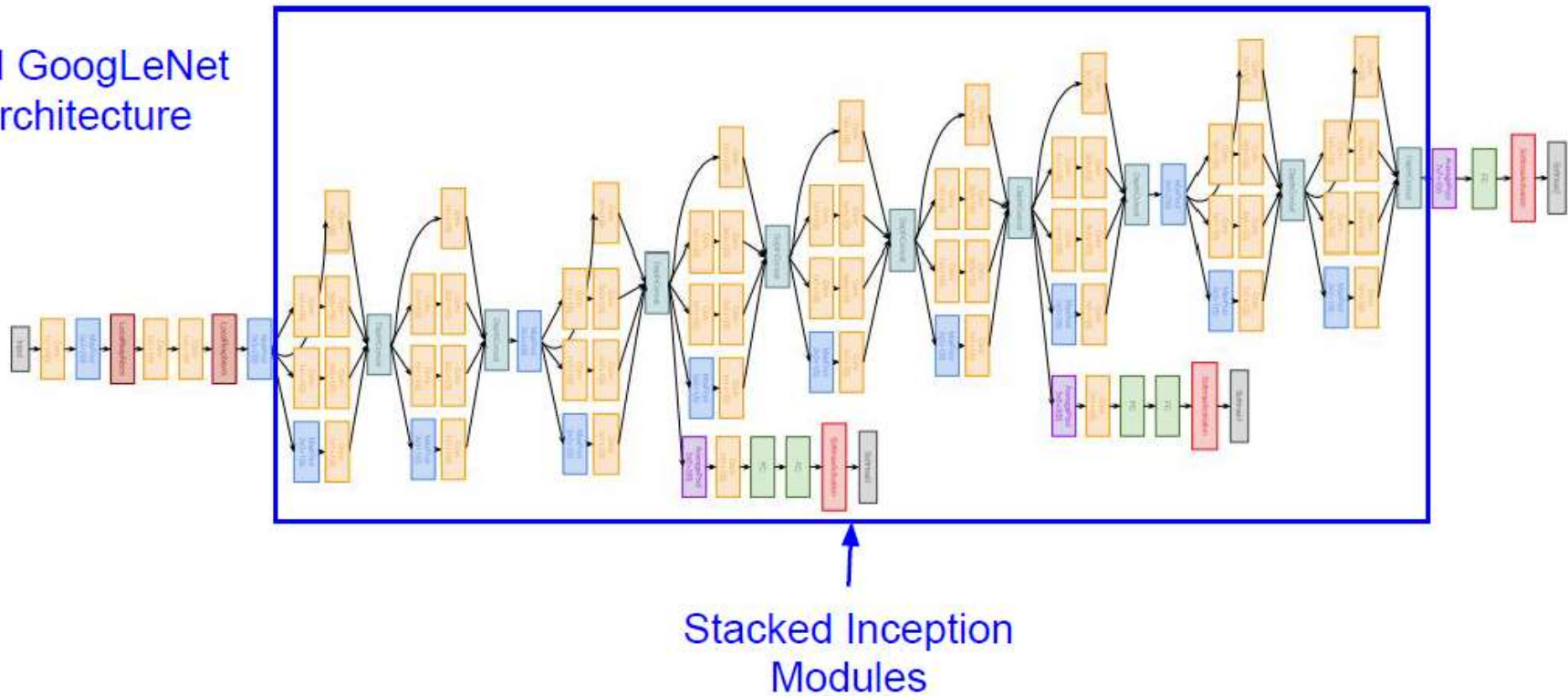
## Full GoogLeNet architecture



Stem Network:  
Conv-Pool-  
2x Conv-Pool

# Inception Network: Google Net

## Full GoogLeNet architecture

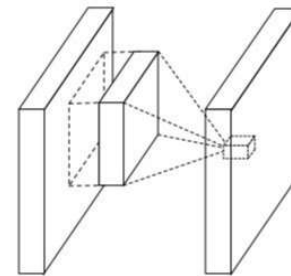




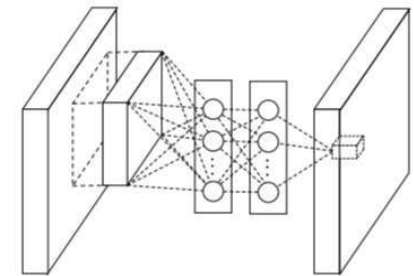
# Network in Network (NiN)

[Lin et al. 2014]

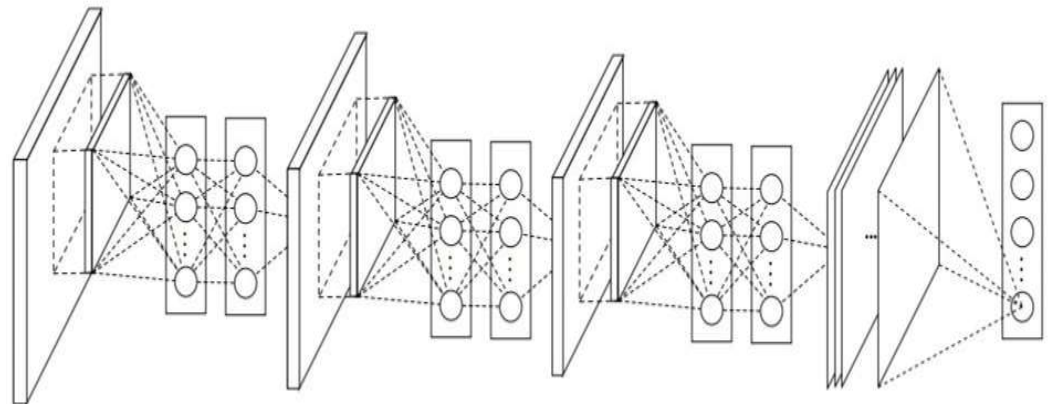
- Mlpconv layer with “micronetwork” within each conv layer to compute more abstract features for local patches
- Micronetwork uses multilayer perceptron (FC, i.e. 1x1 conv layers)
- Precursor to GoogLeNet and ResNet “bottleneck” layers
- Philosophical inspiration for GoogLeNet



(a) Linear convolution layer



(b) Mlpconv layer



Figures copyright Lin et al., 2014.

To continue...