

# Lecture7

# Convolutional Neural Networks

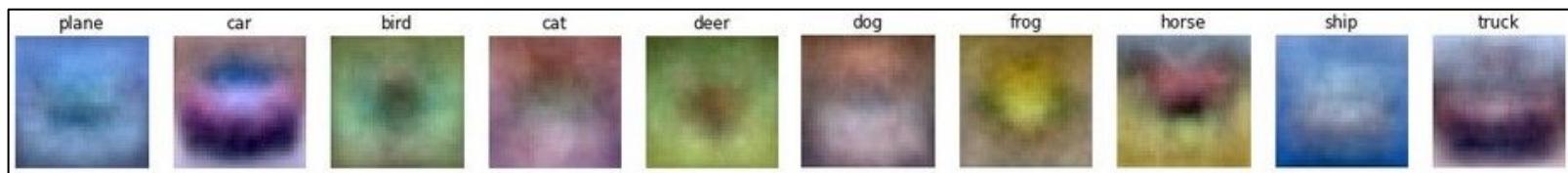
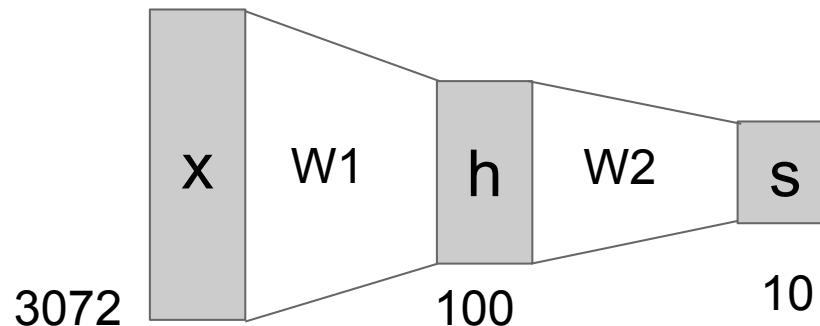
# Last time: Neural Networks

Linear score function:

$$f = Wx$$

2-layer Neural Network

$$f = W_2 \max(0, W_1 x)$$



# Next: Convolutional Neural Networks

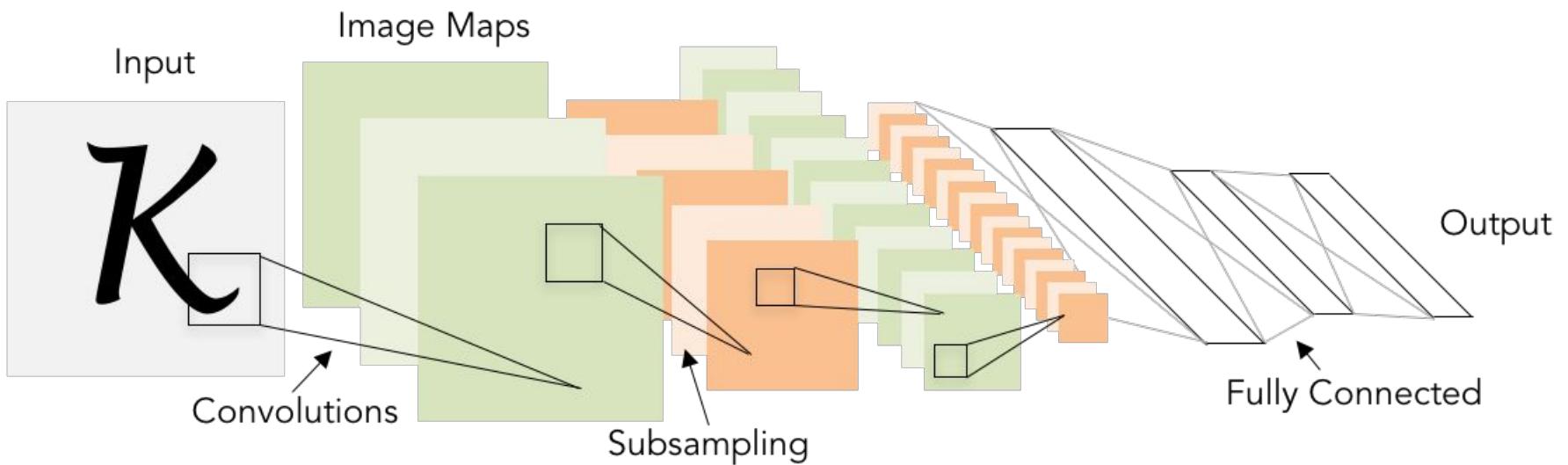
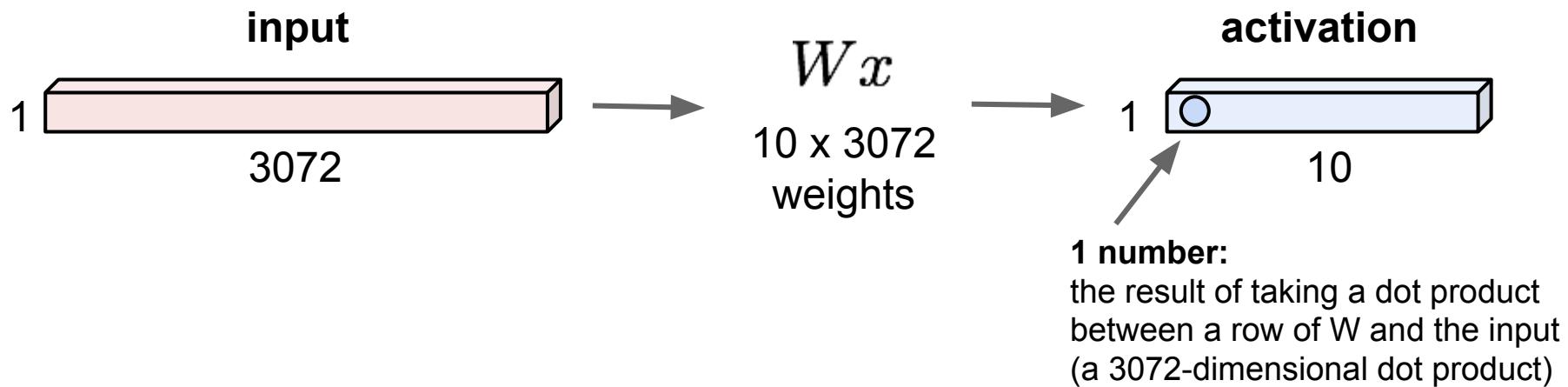


Illustration of LeCun et al. 1998 from CS231n 2017 Lecture 1

# Fully connected layer

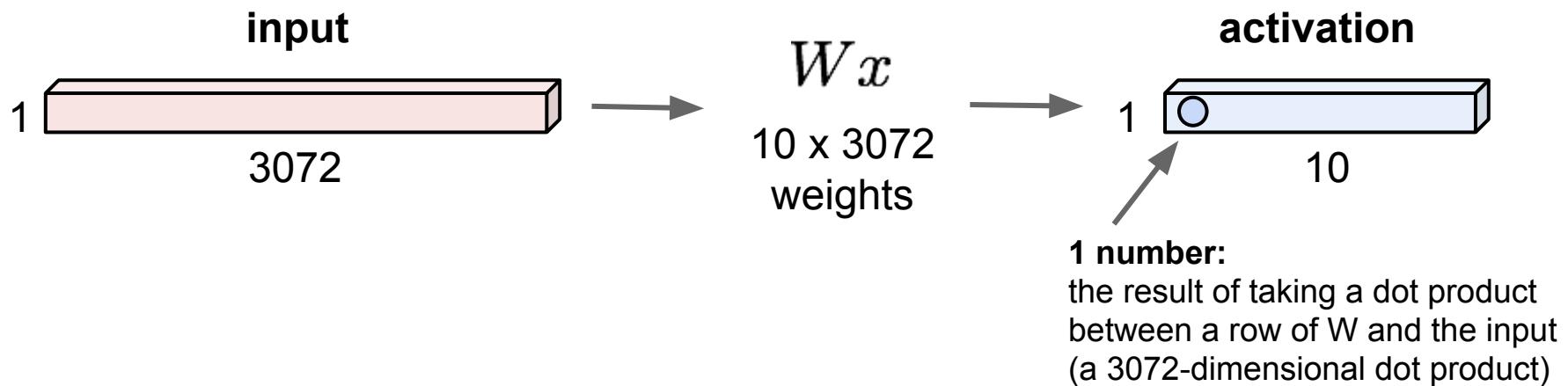
32x32x3 image -> stretch to 3072 x 1



# Convolutional Layer

# Recall fully connected layer

32x32x3 image -> stretch to 3072 x 1

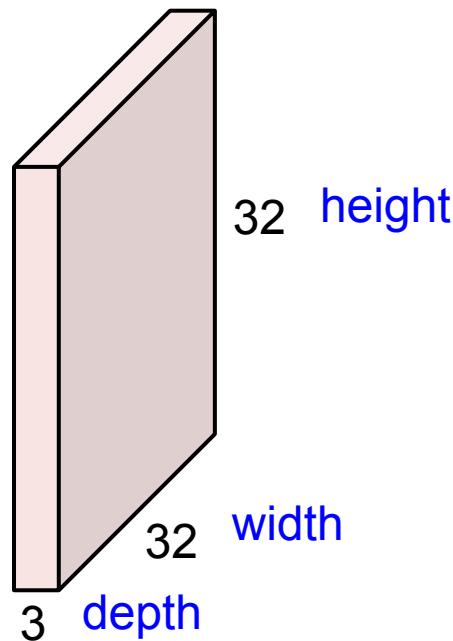


## Problem:

The spatial structure of images is broken. Spatially adjacent pixels are no longer close to each other after stretching

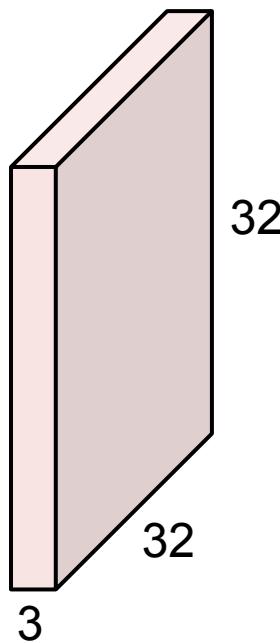
# Convolutional layer

32x32x3 image -> preserve spatial structure



# Convolutional Layer

32x32x3 image



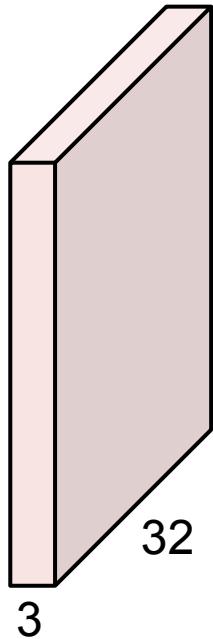
5x5x3 filter



**Convolve** the filter with the image  
i.e. “slide over the image spatially,  
computing dot products”

# Convolutional Layer

32x32x3 image



5x5x3 filter



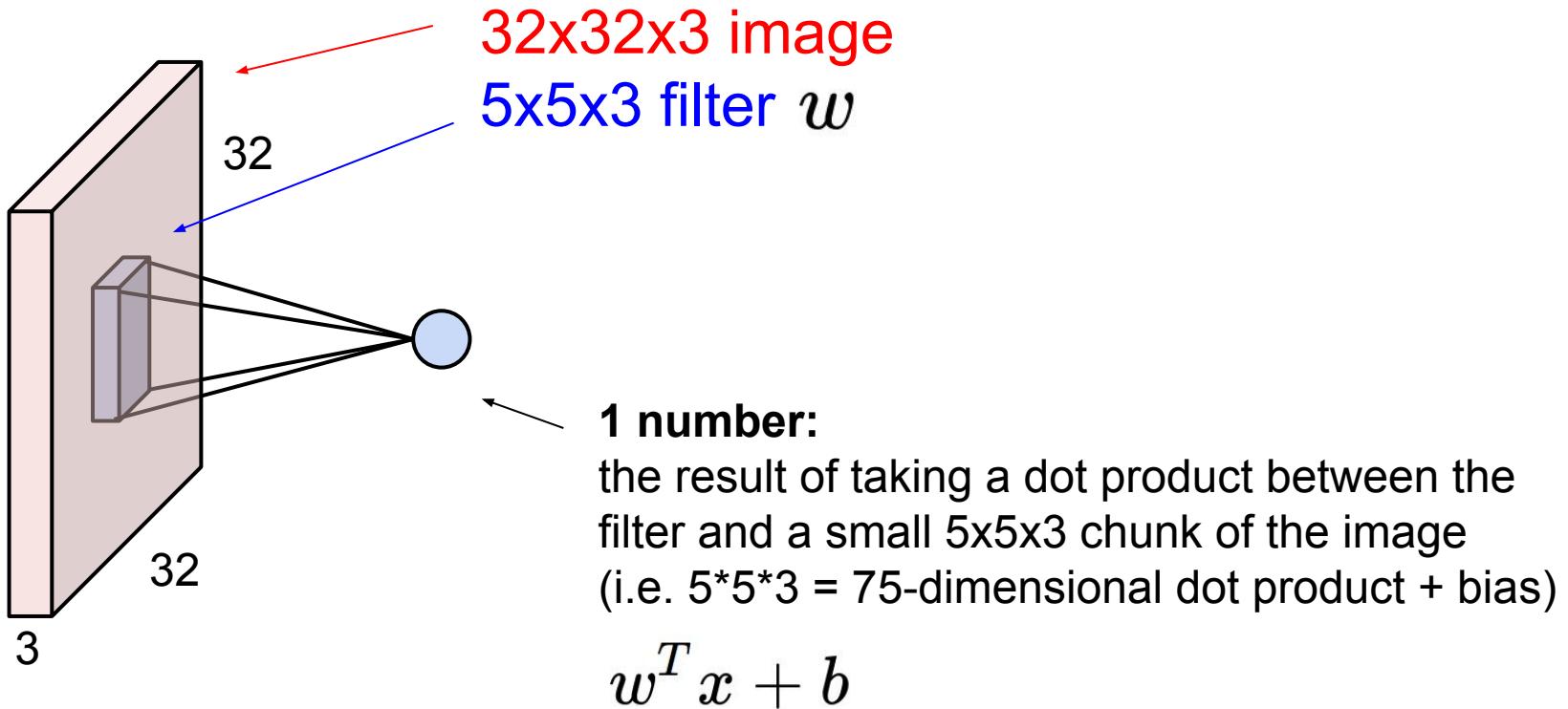
Filters always extend the full depth of the input volume

**Convolve** the filter with the image  
i.e. “slide over the image spatially,  
computing dot products”

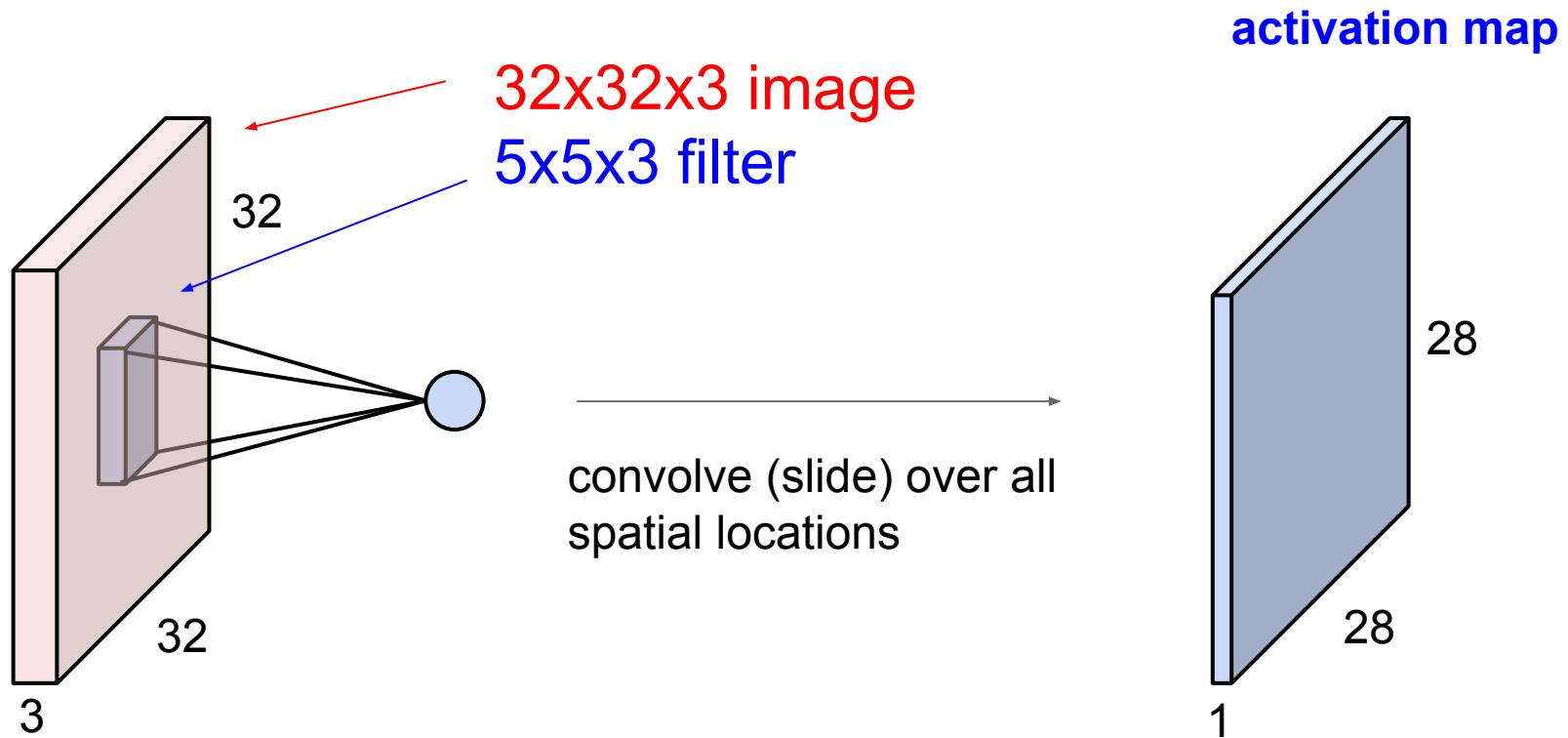
Filters move spatial-wise, not depth-wise

Q: Will the results be different if we move filters along different directions?

# Convolutional Layer



# Convolutional Layer



An activation map is a sheet of neuron outputs

- 1) Each is connected to a small region (receptive field, 感受野) in the input
- 2) All of them share parameters → translation invariant

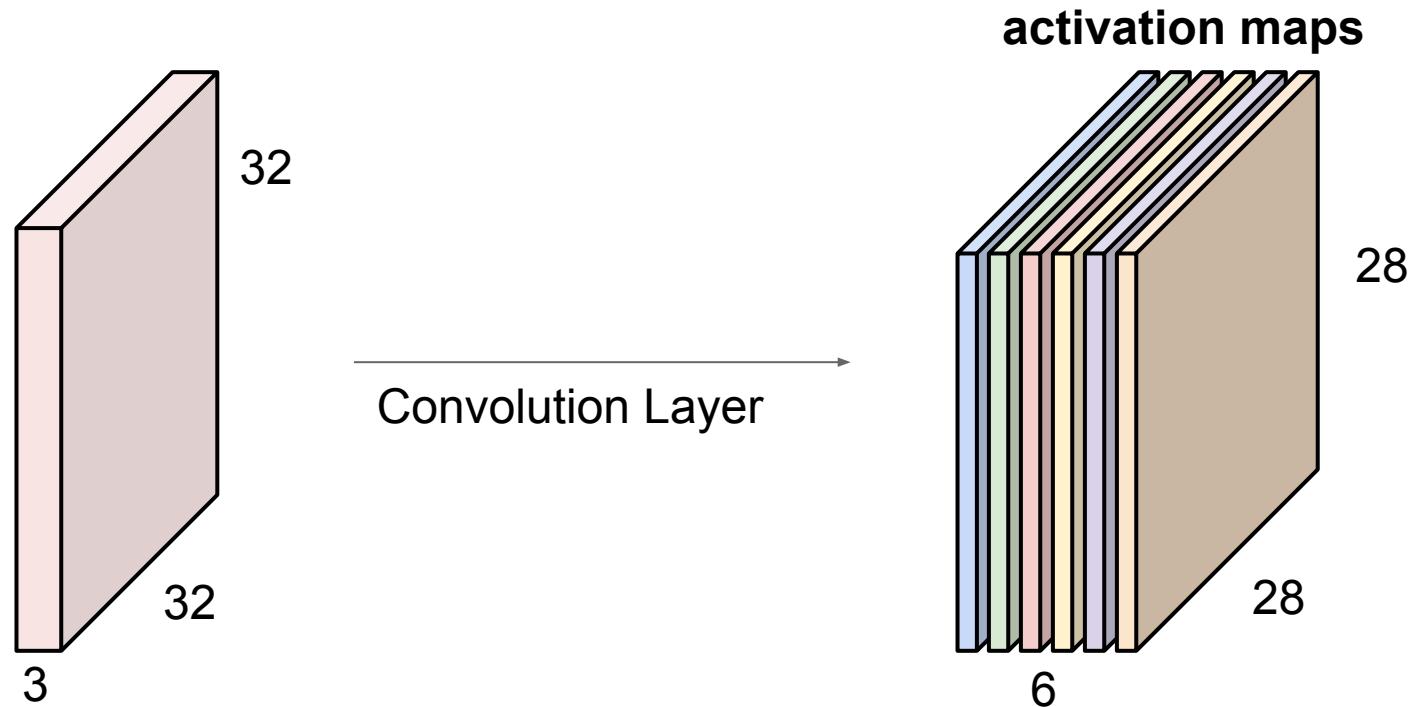
# Convolutional Layer



Consider another **green** filter

# Convolutional Layer

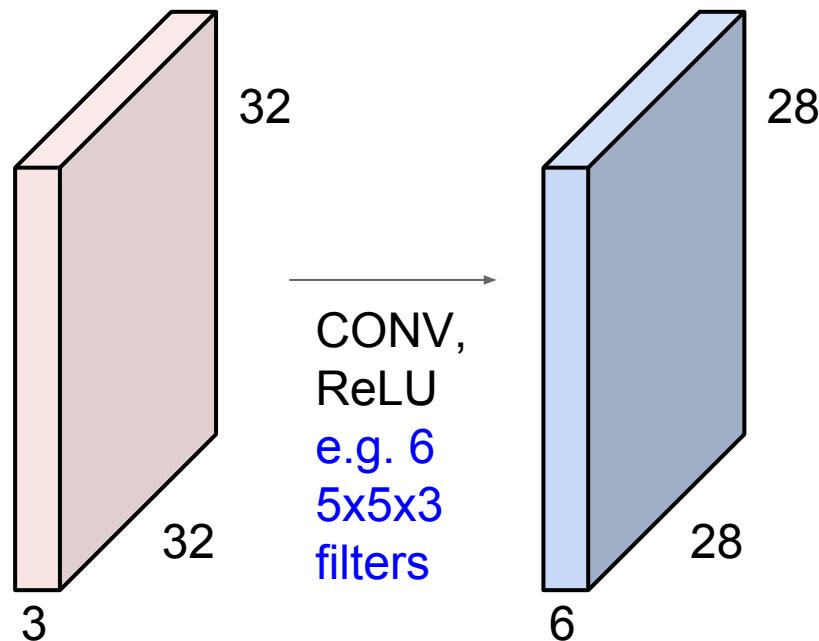
For example, if we had 6  $5 \times 5$  filters, we'll get 6 separate activation maps:



We stack these up to get a “new image” of size  $28 \times 28 \times 6$ !

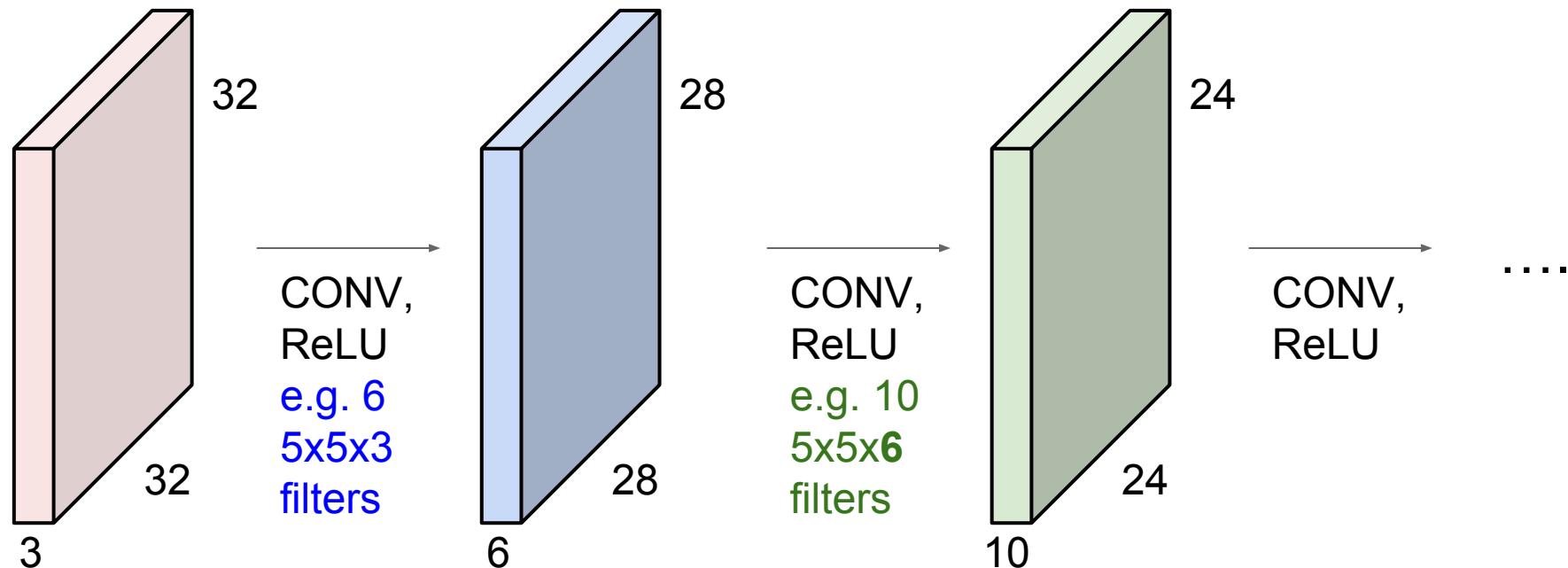
# ConvNet

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



# ConvNet

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



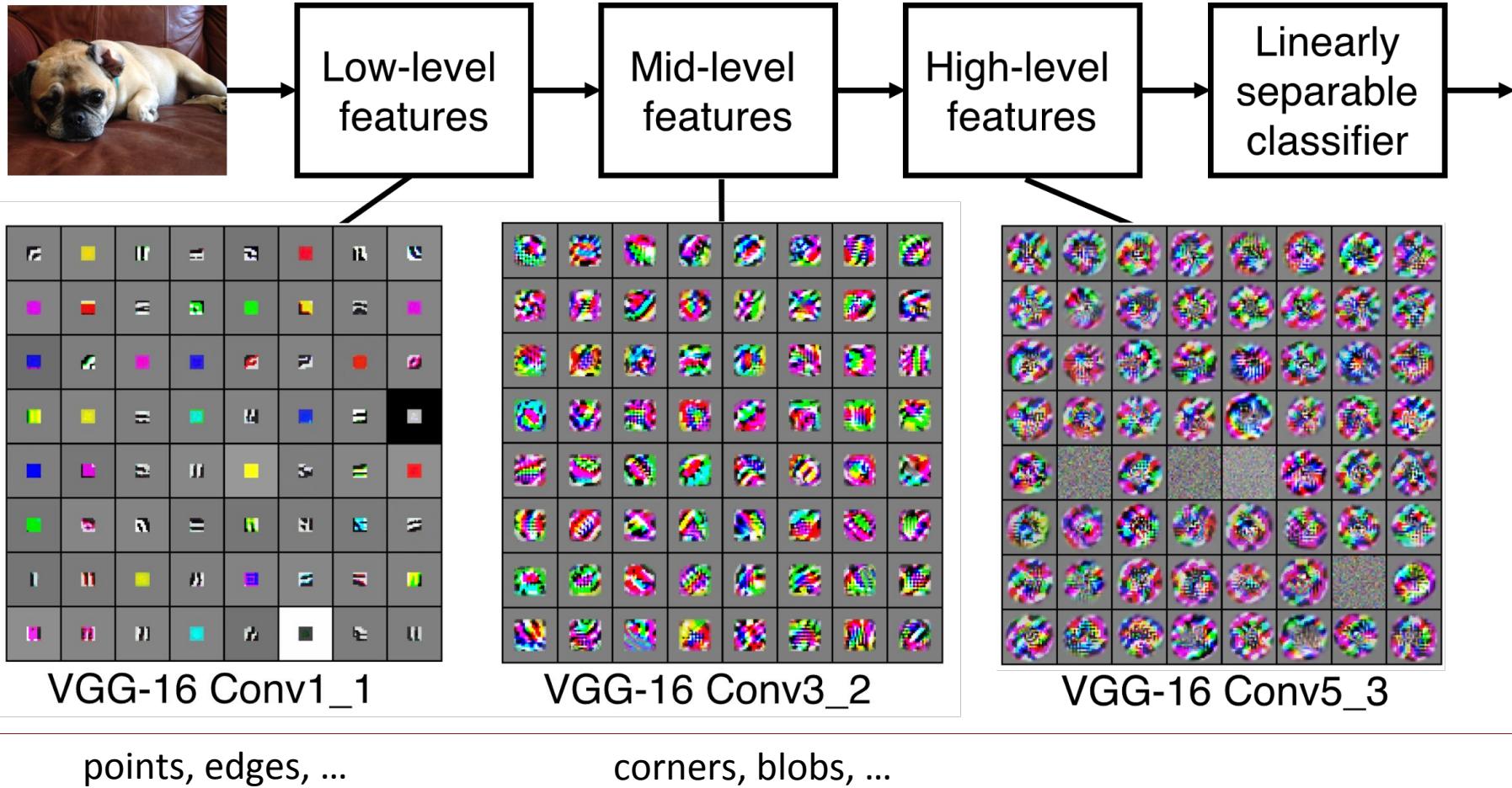
# Visualization

Let's see what the Inputs maximizing the activations of conv neurons look like

## Preview

[Zeiler and Fergus 2013]

Visualization of VGG-16 by Lane McIntosh. VGG-16 architecture from [Simonyan and Zisserman 2014].



# 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

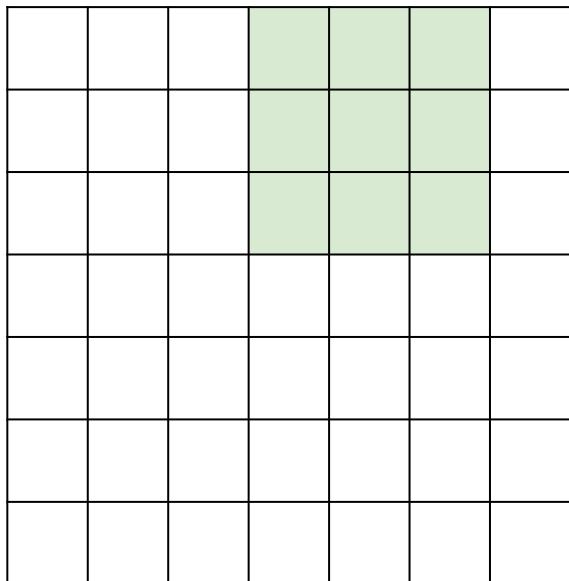
7


7x7 input (spatially)  
assume 3x3 filter

7

# A closer look at spatial dimensions

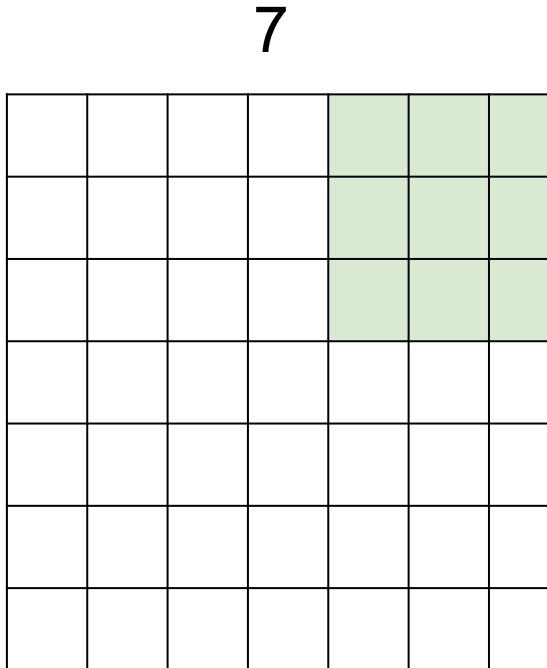
7



7x7 input (spatially)  
assume 3x3 filter

7

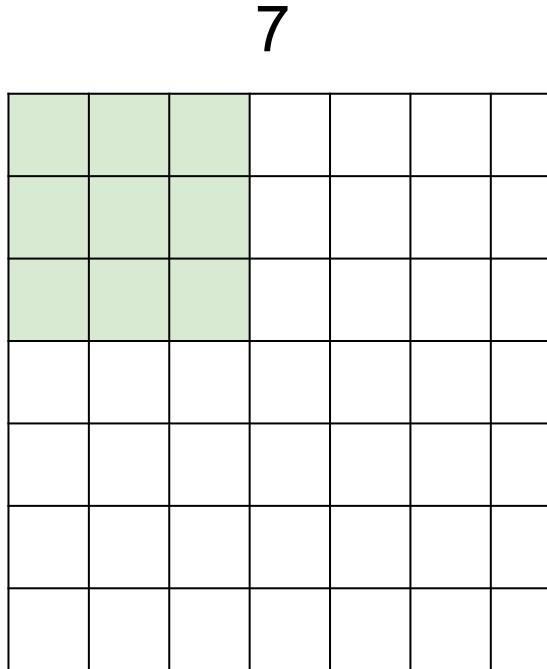
# A closer look at spatial dimensions



7x7 input (spatially)  
assume 3x3 filter

**=> 5x5 output**

# A closer look at spatial dimensions

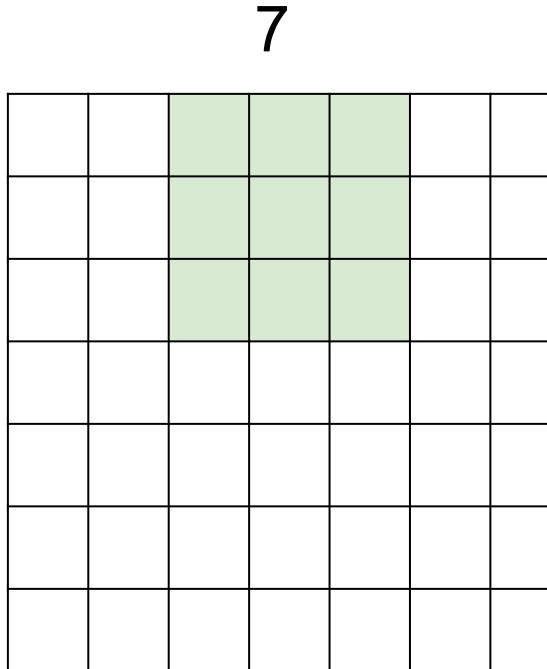


7x7 input (spatially)  
assume 3x3 filter  
applied **with stride 2**

7

Stride : the step length  
(in number of pixels)  
when you move a filter.

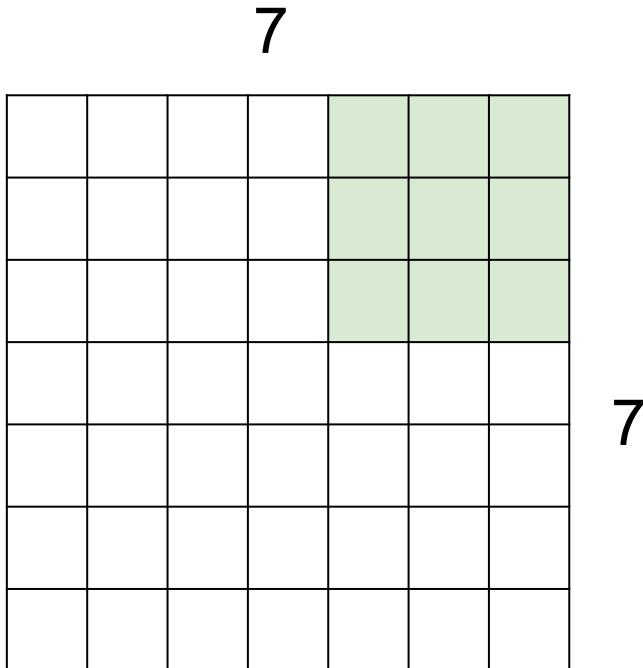
# A closer look at spatial dimensions



7

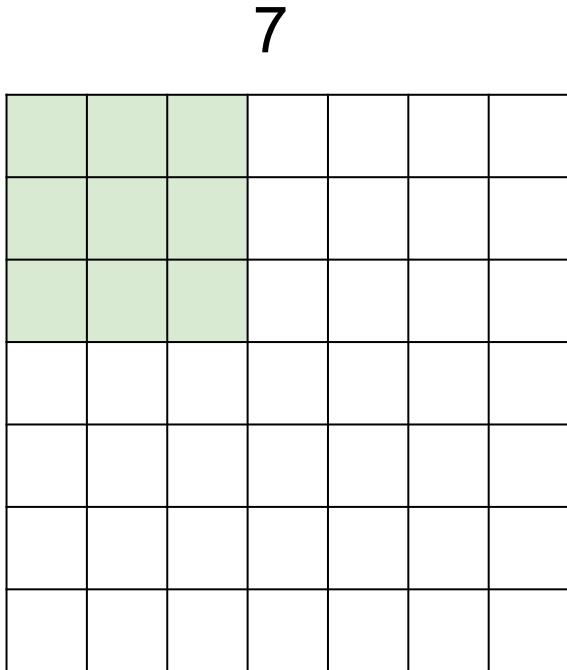
7x7 input (spatially)  
assume 3x3 filter  
applied **with stride 2**

# A closer look at spatial dimensions



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

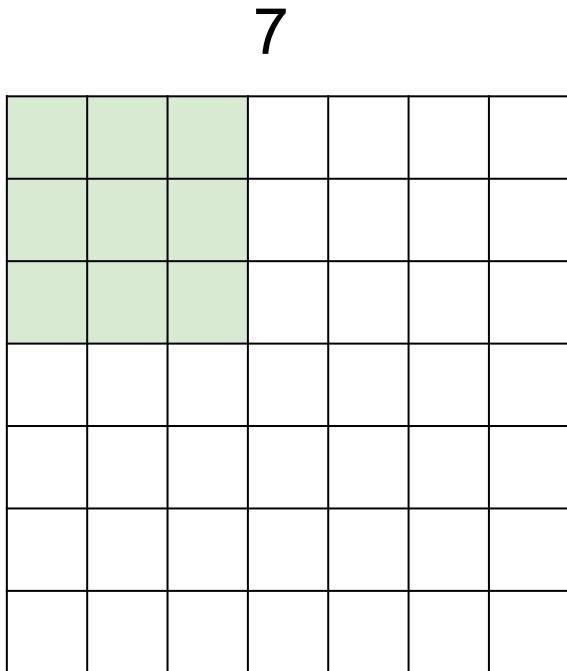
# A closer look at spatial dimensions



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

7

# A closer look at spatial dimensions



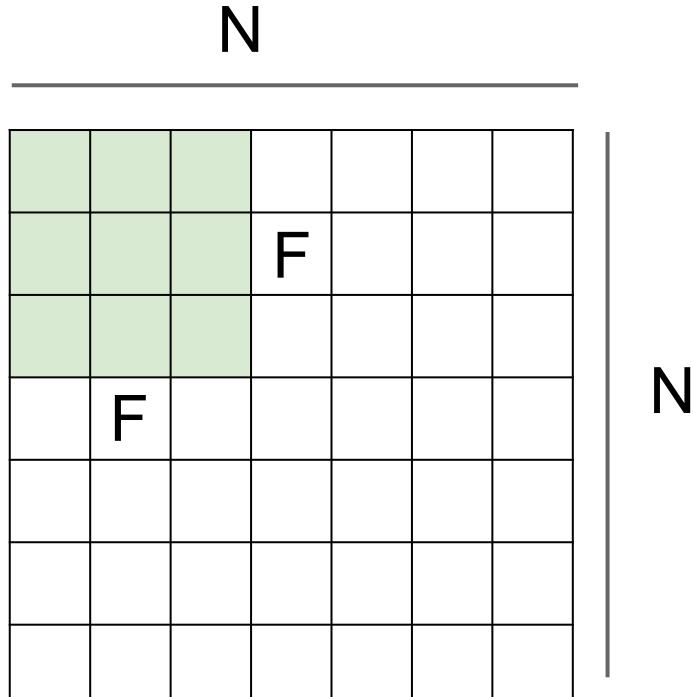
7

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

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

Leads to asymmetric  
outputs and avoid doing  
that!

# A closer look at spatial dimensions



Output size:  
**(N - F) / stride + 1**

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

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

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

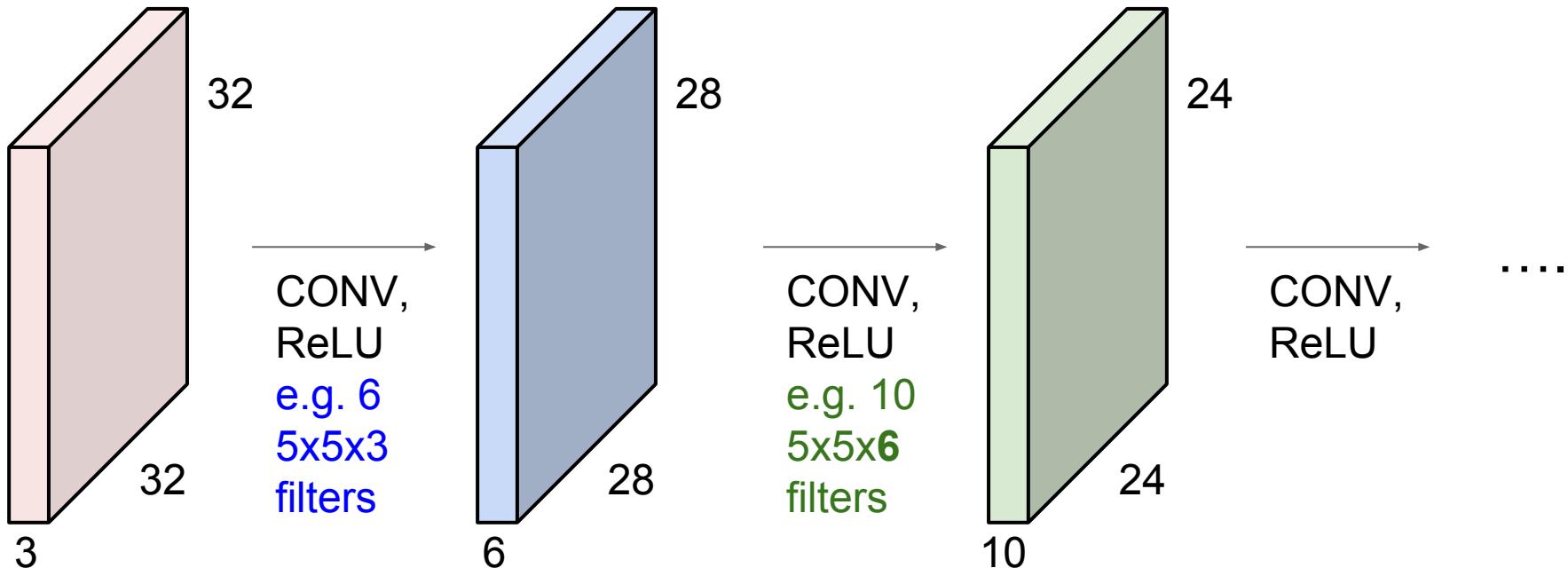
$F = 5 \Rightarrow$  zero pad with 2

$F = 7 \Rightarrow$  zero pad with 3

# In practice: common to zero pad the border

**Remember back to...**

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



Now it is possible to build a very deep ConvNet

# Example

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

Output volume size: ?

# Example

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

Output volume size:  
 $(32+2*2-5)/1+1 = 32$  spatially, so  
**32x32x10**

# Example

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

Number of parameters in this layer?

# Example

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

$$\Rightarrow \mathbf{76 * 10 = 760}$$

# Example

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

Number of add&mul operations in this layer?

# Example

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

Number of add&mul operations in this layer?  
each output neuron has 5x5x3 ops, and there  
are 32x32x10 output neurons =>  
 $32 \times 32 \times 10 \times 5 \times 5 \times 3$

# Summary

**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.

# Summary

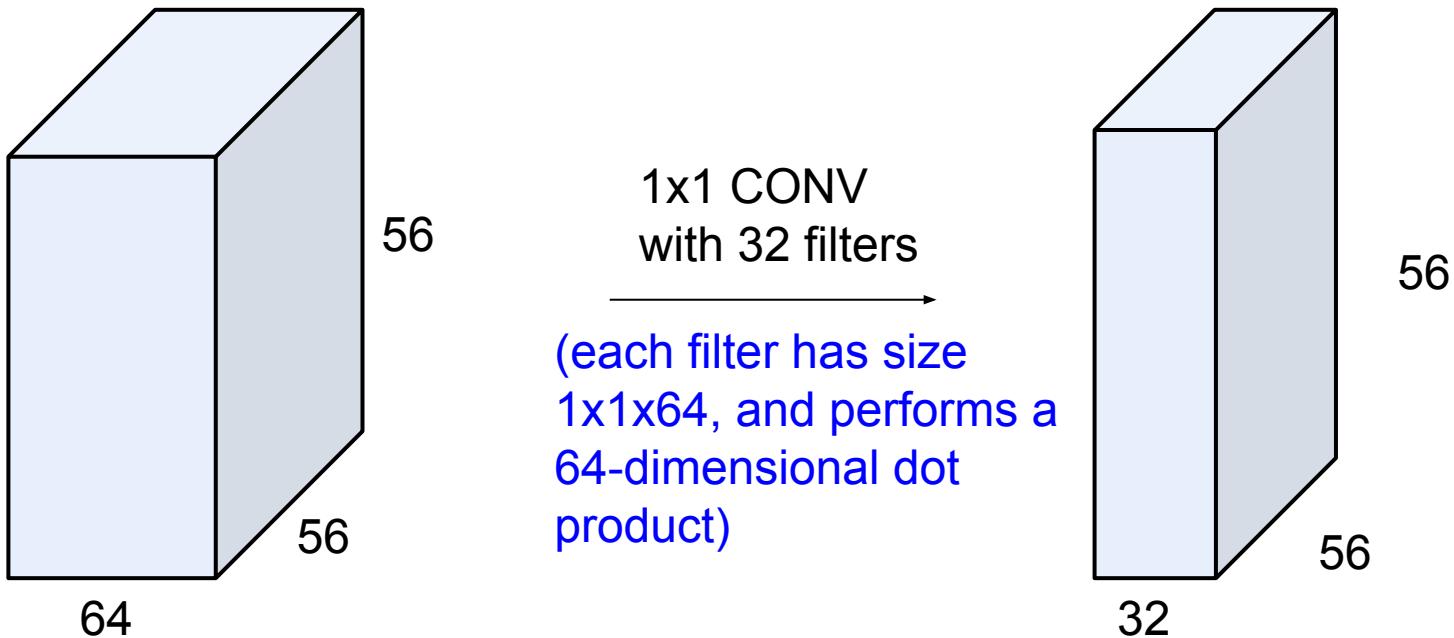
**Summary.** To summarize, the Conv Layer:

- 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.

Common settings:

- $K = (\text{powers of 2, e.g. } 32, 64, 128, 512)$
- $F = 3, S = 1, P = 1$
  - $F = 5, S = 1, P = 2$
  - $F = 5, S = 2, P = ?$  (whatever fits)
  - $F = 1, S = 1, P = 0$

# 1x1 convolution

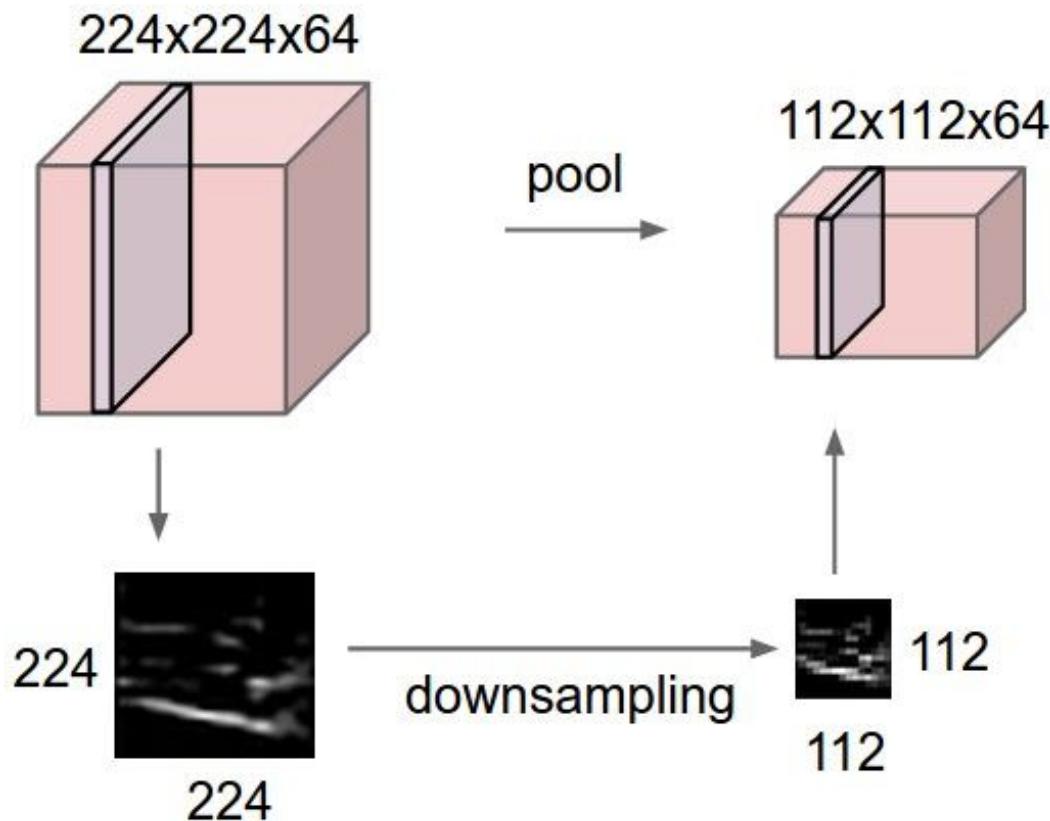


Q: Is a fully-connected layer a convolutional layer?

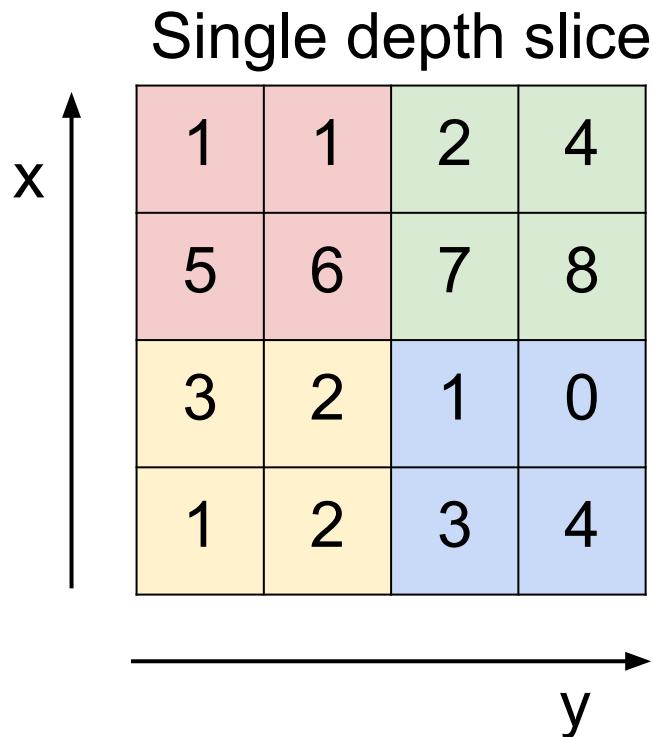
# Pooling Layer

# Pooling layer

- To make the representations smaller and more manageable
- Operates over each activation map independently: **output-depth = input-depth**



# Max Pooling



max pool with 2x2 filters  
and stride 2

6	8
3	4

# Summary

- Accepts a volume of size  $W_1 \times H_1 \times D_1$
- Requires three hyperparameters:
  - their spatial extent  $F$ ,
  - 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

# Summary

- Accepts a volume of size  $W_1 \times H_1 \times D_1$
- Requires three hyperparameters:
  - their spatial extent  $F$ ,
  - 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

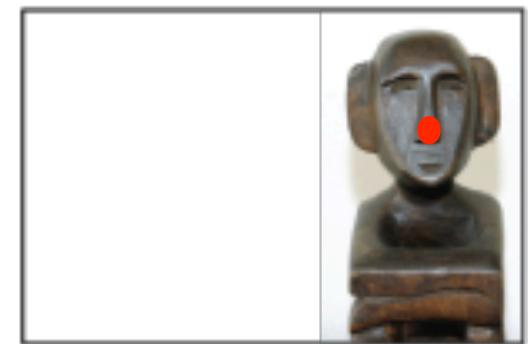
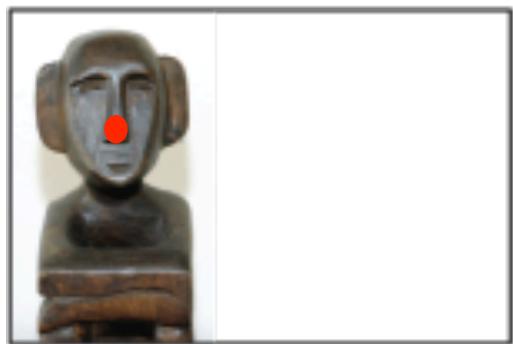
Common settings:

$F = 2, S = 2$

$F = 3, S = 2$

# Convolution + Pooling $\approx$ Translation invariance

- Convolution: local feature detector
- Pooling: invariant to small translation

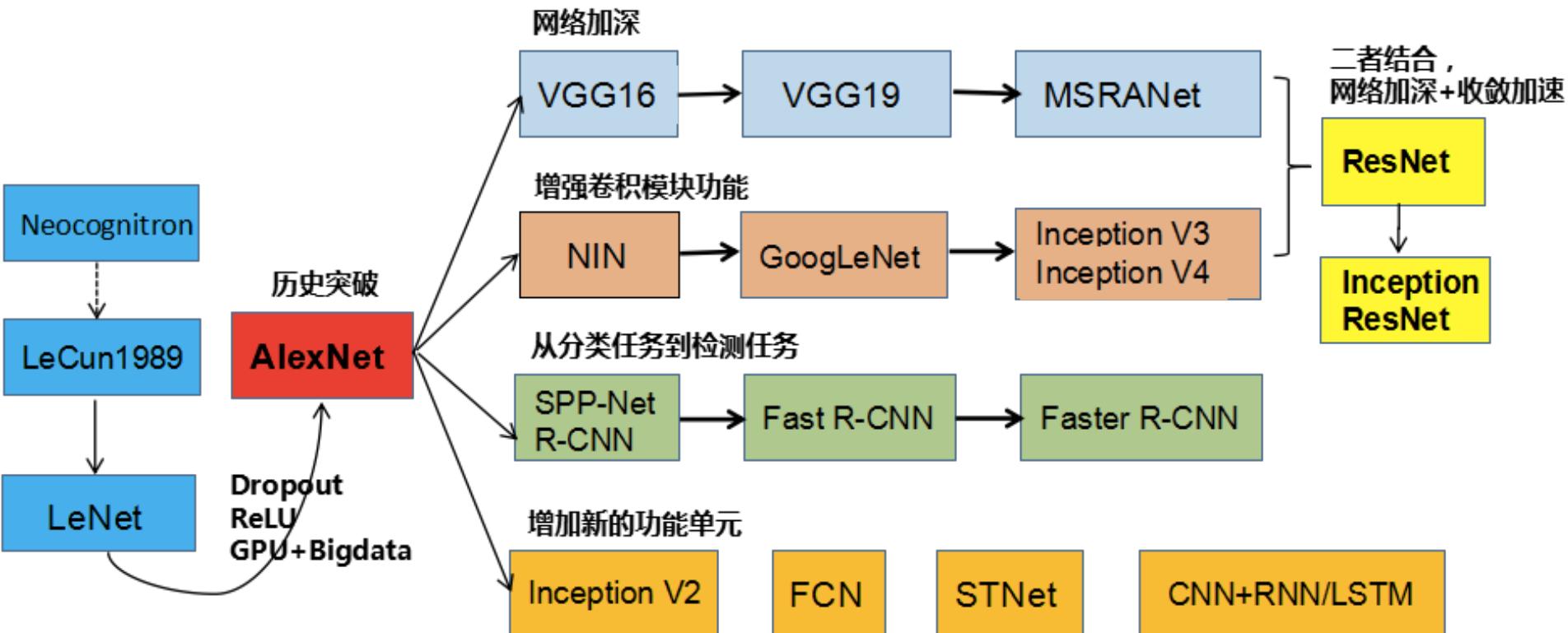


# Taking-home messages

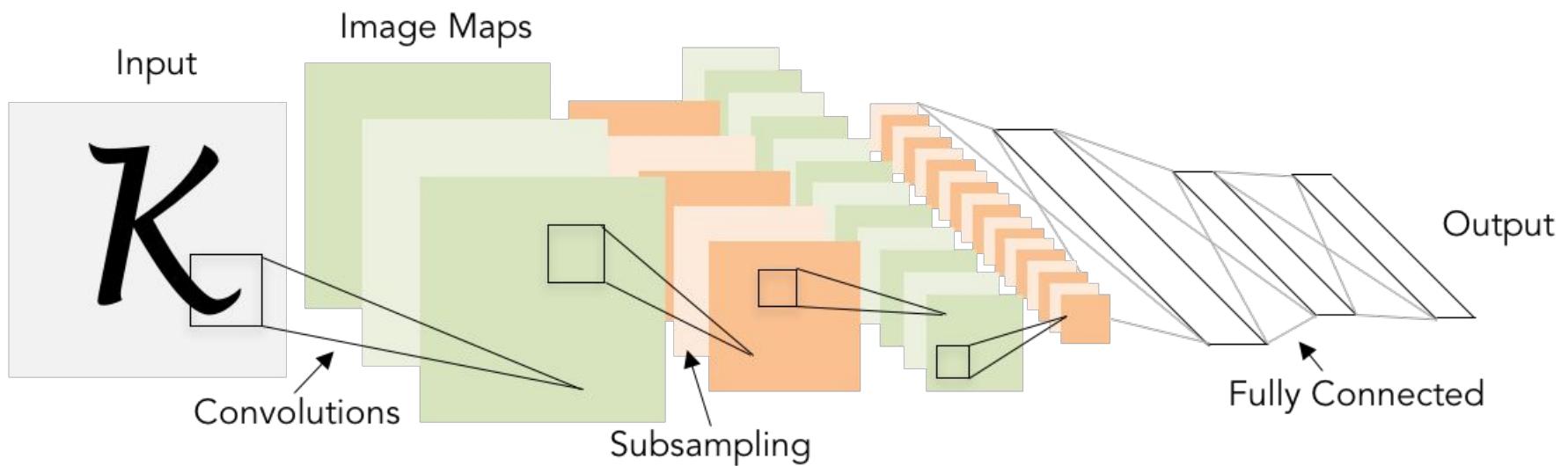
- ConvNets stack CONV, POOL, FC layers
- Trend towards smaller filters and deeper architectures
- Trend towards getting rid of POOL/FC layers (just CONV)
- Typical architectures look like:  
$$[(\text{CONV-RELU})^*N - \text{POOL?}]^*M - (\text{FC-RELU})^*K - \text{SOFTMAX}$$
where N is usually up to  $\sim 5$ , M is large,  $0 \leq K \leq 2$ . But recent advances such as ResNet challenge this paradigm
- Convolution + Pooling  $\approx$  Translation invariance

# CNN Architectures

# The evolution of CNN architectures



# LeNet-5 [LeCun et al., 1998]



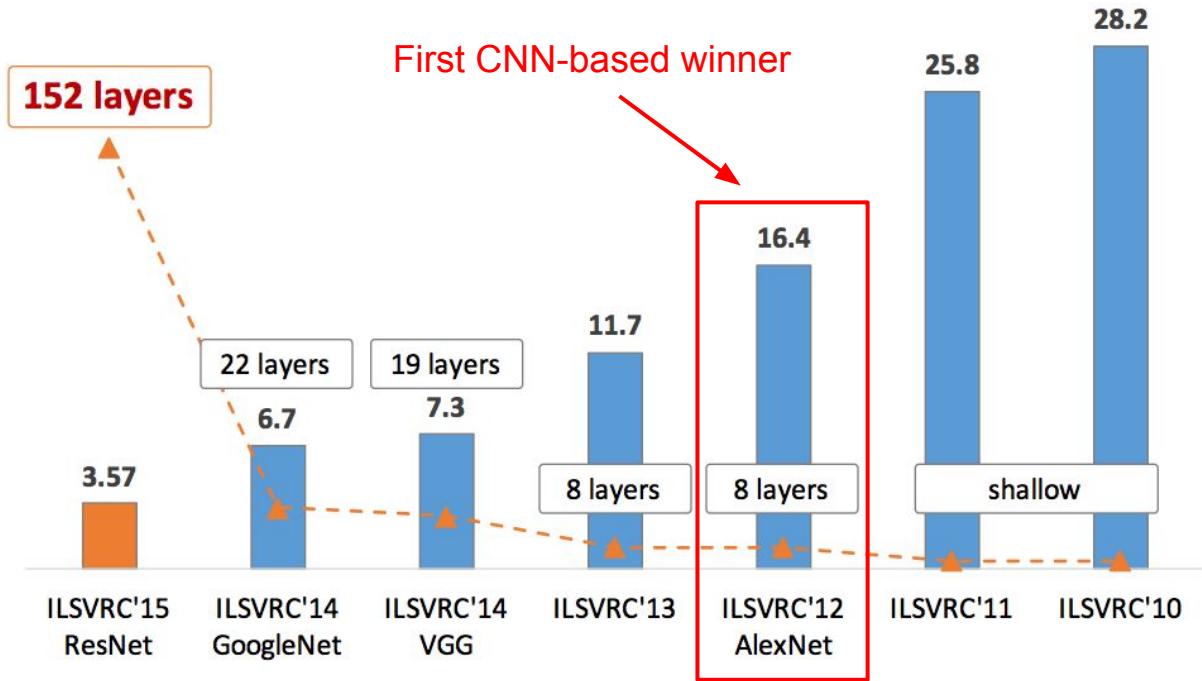
Conv filters were 5x5, applied at stride 1

Subsampling (Pooling) layers were 2x2 applied at stride 2

i.e. architecture is [CONV-POOL-CONV-POOL-FC-FC]

# AlexNet [Krizhevsky et al. 2012]

ImageNet Large Scale Visual Recognition Challenge (ILSVRC) winners



# AlexNet [Krizhevsky et al. 2012]

Architecture:

CONV1

MAX POOL1

NORM1

CONV2

MAX POOL2

NORM2

CONV3

CONV4

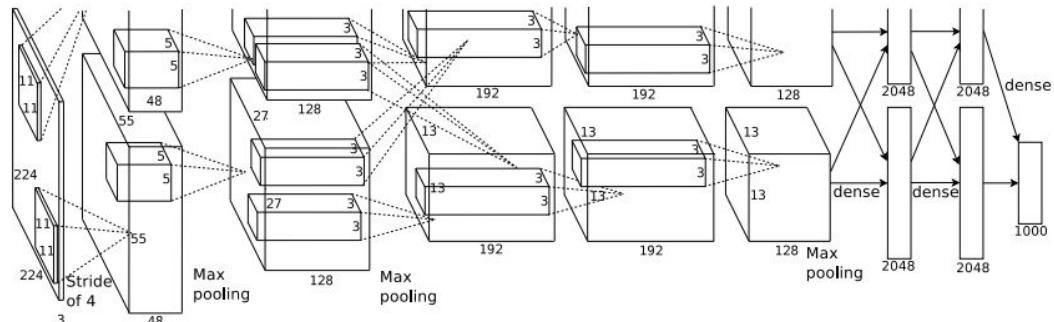
CONV5

Max POOL3

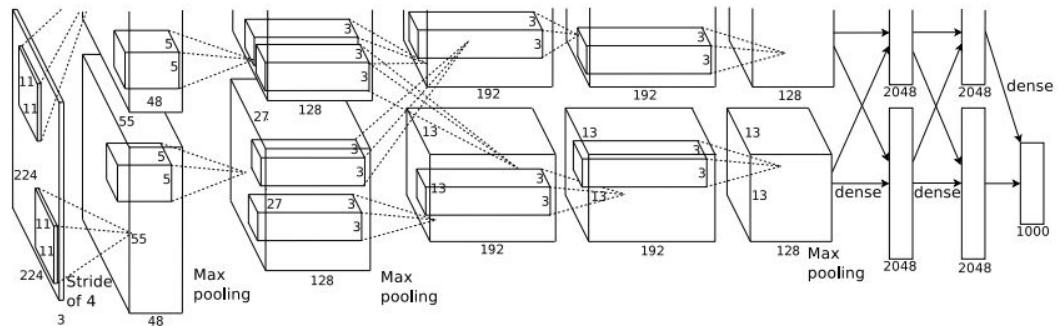
FC6

FC7

FC8



# AlexNet [Krizhevsky et al. 2012]



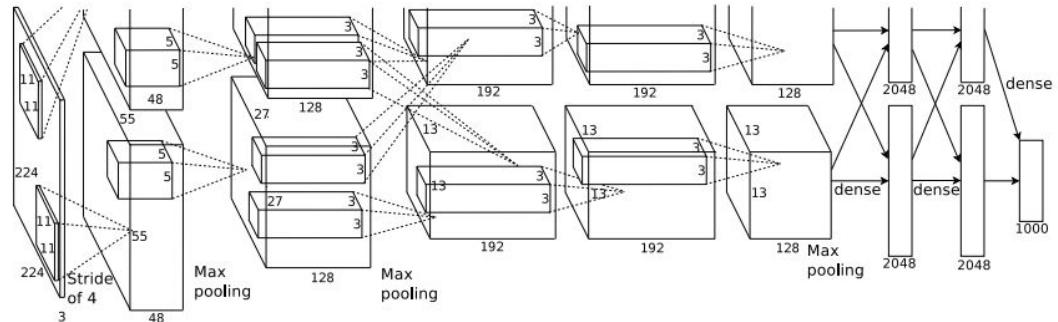
Input: 227x227x3 images

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

=>

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

# AlexNet [Krizhevsky et al. 2012]



Input: 227x227x3 images

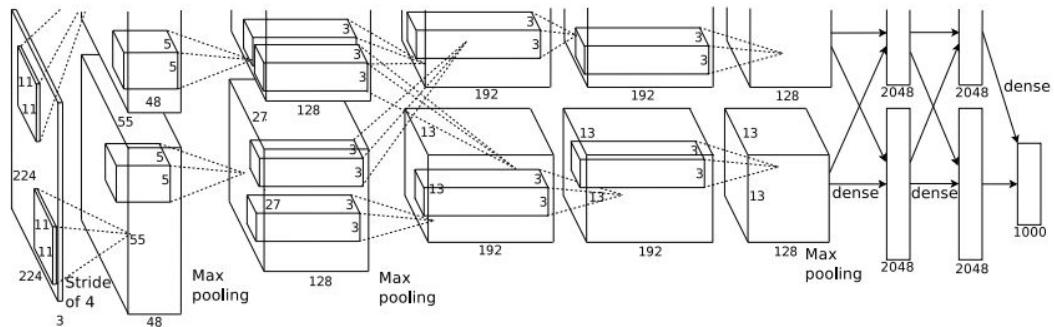
**First layer (CONV1):** 96 11x11 filters applied at stride 4

=>

Output volume **[55x55x96]**

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

# AlexNet [Krizhevsky et al. 2012]



Input: 227x227x3 images

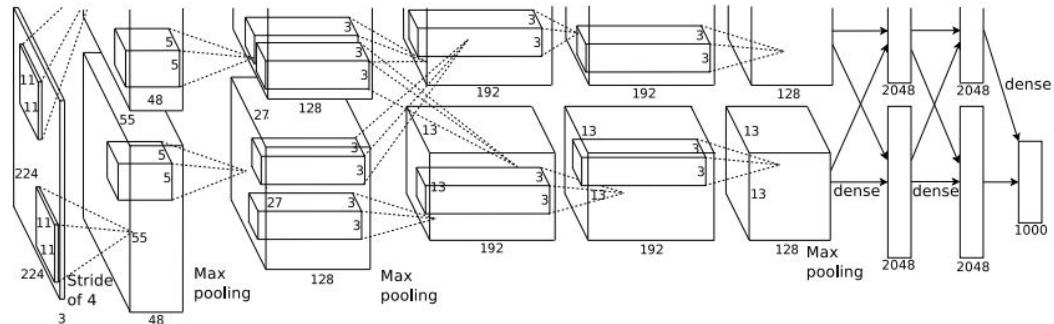
**First layer (CONV1):** 96 11x11 filters applied at stride 4

=>

Output volume **[55x55x96]**

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

# AlexNet [Krizhevsky et al. 2012]



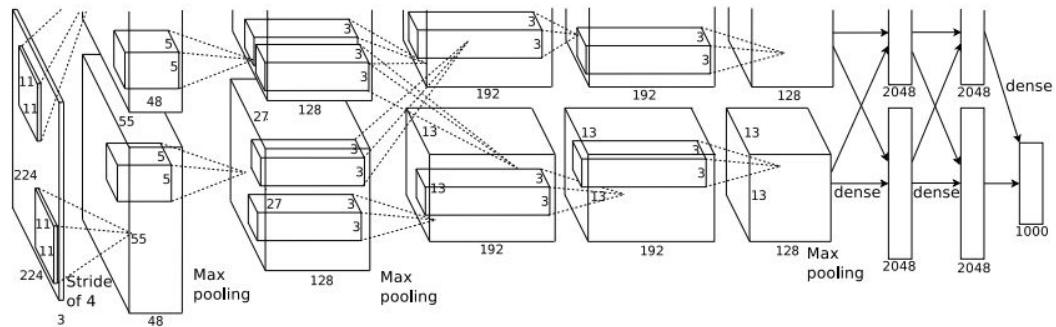
Input: 227x227x3 images

After CONV1: 55x55x96

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

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

# AlexNet [Krizhevsky et al. 2012]



Input: 227x227x3 images

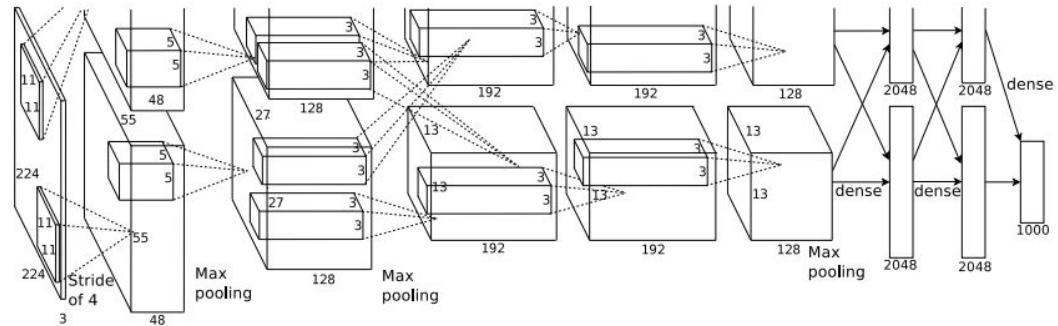
After CONV1: 55x55x96

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

Output volume: 27x27x96

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

# AlexNet [Krizhevsky et al. 2012]



Input: 227x227x3 images

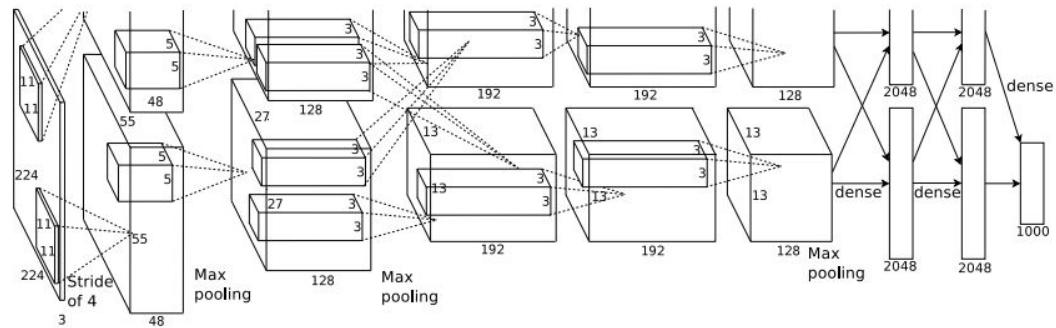
After CONV1: 55x55x96

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

Output volume: 27x27x96

Parameters: 0!

# AlexNet [Krizhevsky et al. 2012]



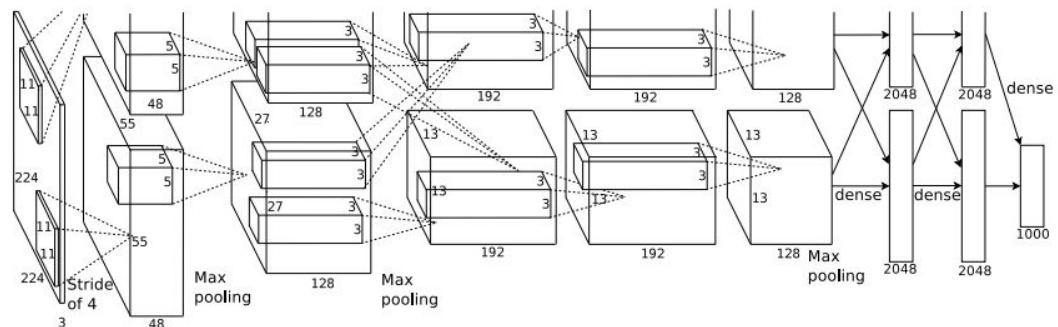
Input: 227x227x3 images

After CONV1: 55x55x96

After POOL1: 27x27x96

...

# 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)

# 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

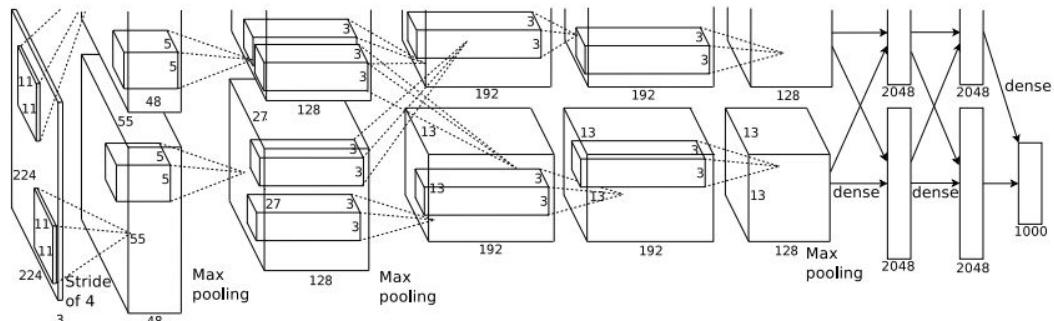
[13x13x256] CONV5: 256 3x3 filters at stride 1, pad 1

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

[4096] FC6: 4096 neurons

[4096] FC7: 4096 neurons

[1000] FC8: 1000 neurons (class scores)



## Details/Retrospectives:

- first use of ReLU
- used Norm layers (not common anymore)
- heavy data augmentation
- dropout 0.5
- batch size 128
- SGD Momentum 0.9
- Learning rate 1e-2, reduced by 10 manually when val accuracy plateaus
- L2 weight decay 5e-4
- 7 CNN ensemble: 18.2% → 15.4%

# 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

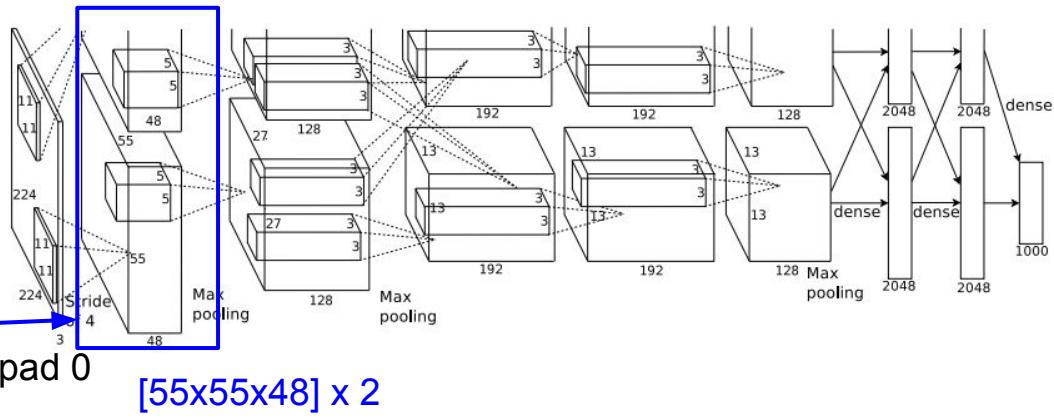
[13x13x256] CONV5: 256 3x3 filters at stride 1, pad 1

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

[4096] FC6: 4096 neurons

[4096] FC7: 4096 neurons

[1000] FC8: 1000 neurons (class scores)



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.

# 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

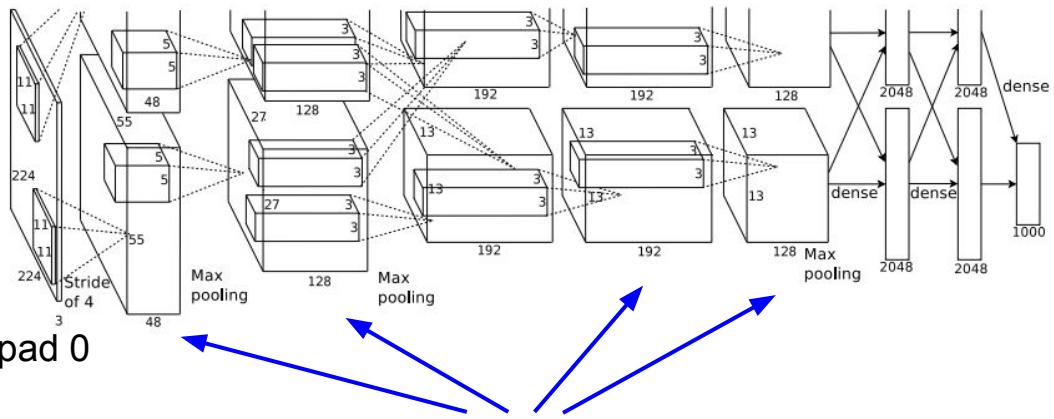
[13x13x256] CONV5: 256 3x3 filters at stride 1, pad 1

[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

# 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

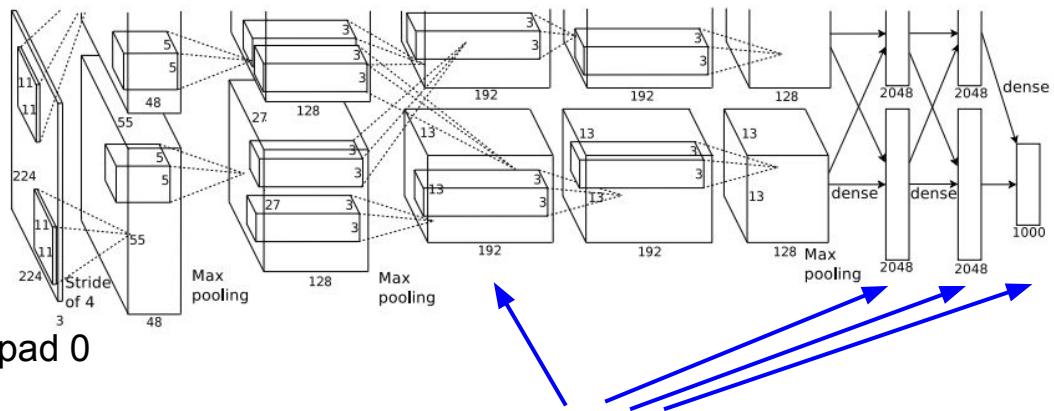
[13x13x256] CONV5: 256 3x3 filters at stride 1, pad 1

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

[4096] FC6: 4096 neurons

[4096] FC7: 4096 neurons

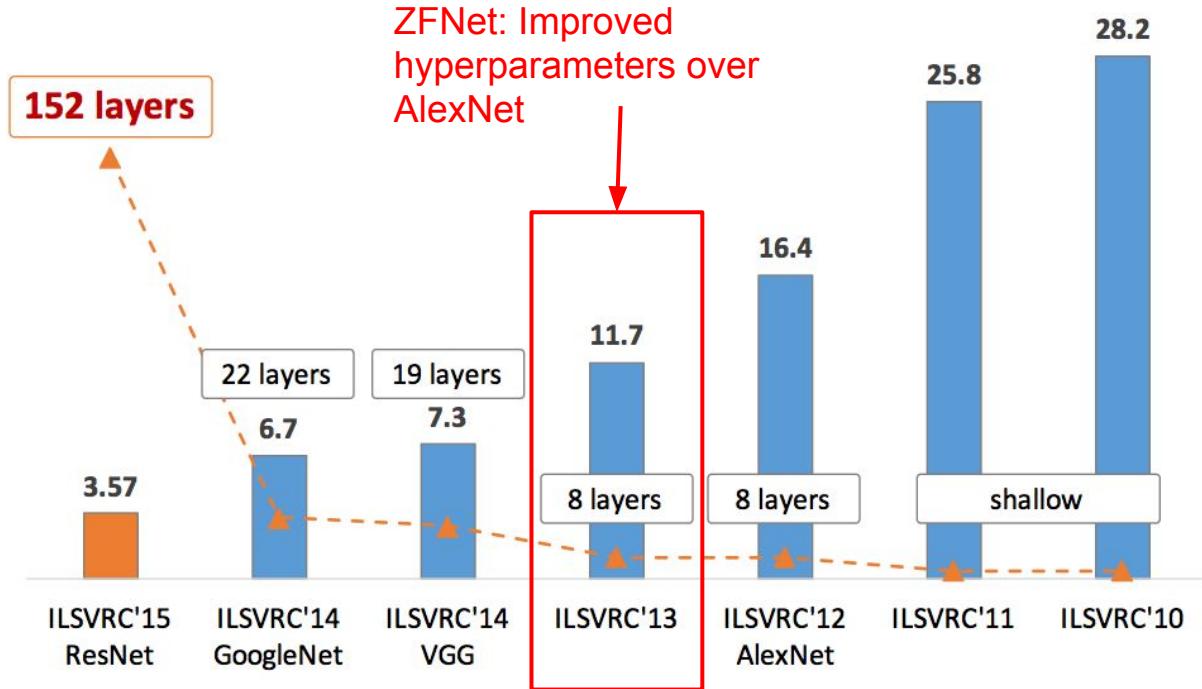
[1000] FC8: 1000 neurons (class scores)



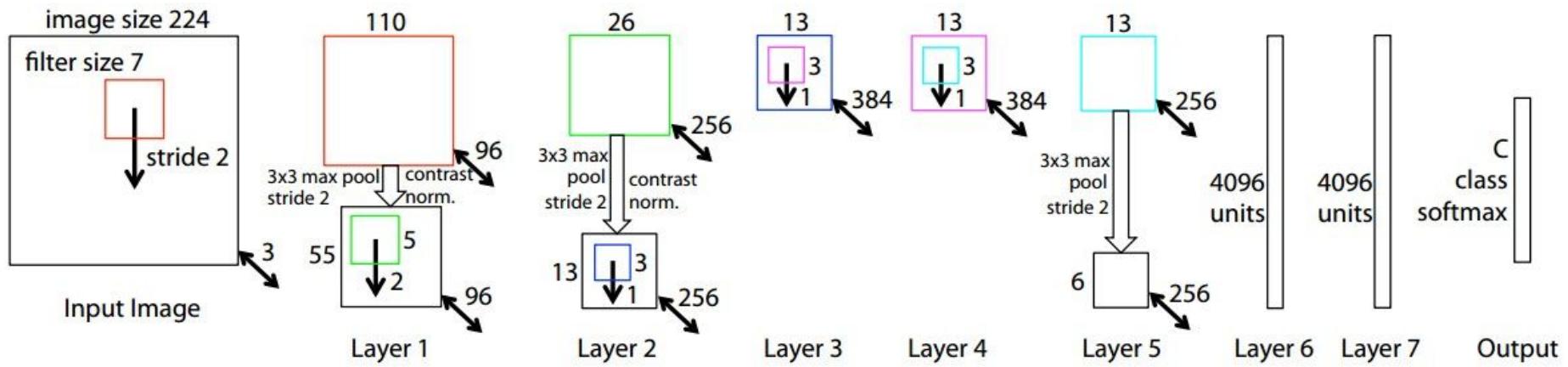
CONV3, FC6, FC7, FC8:  
Connections with all feature maps in  
preceding layer, communication  
across GPUs

# ZFNet [Zeiler and Fergus, 2013]

ImageNet Large Scale Visual Recognition Challenge (ILSVRC) winners



# ZFNet [Zeiler and Fergus, 2013]



AlexNet but:

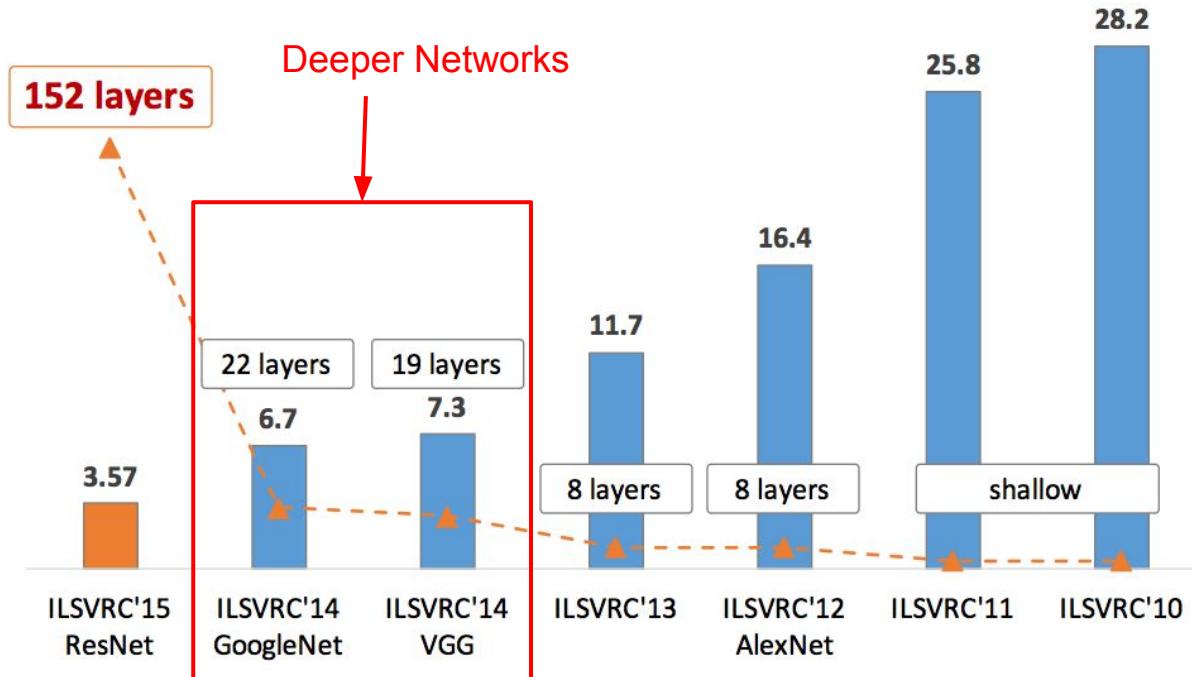
CONV1: change from (11x11 stride 4) to (7x7 stride 2)

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

ImageNet top 5 error: 16.4% -> 11.7%

# VGGNet [Simonyan and Zisserman, 2014]

ImageNet Large Scale Visual Recognition Challenge (ILSVRC) winners



# VGGNet [Simonyan and Zisserman, 2014]

Small filters, Deeper networks

8 layers (AlexNet)

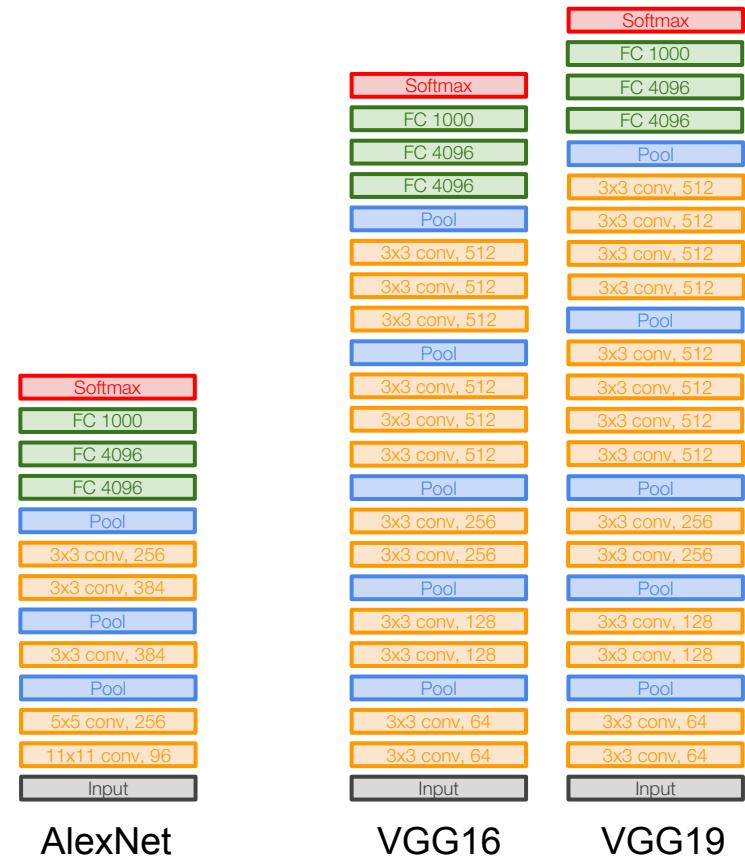
-> 16 - 19 layers (VGG16Net)

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

11.7% top 5 error in ILSVRC'13

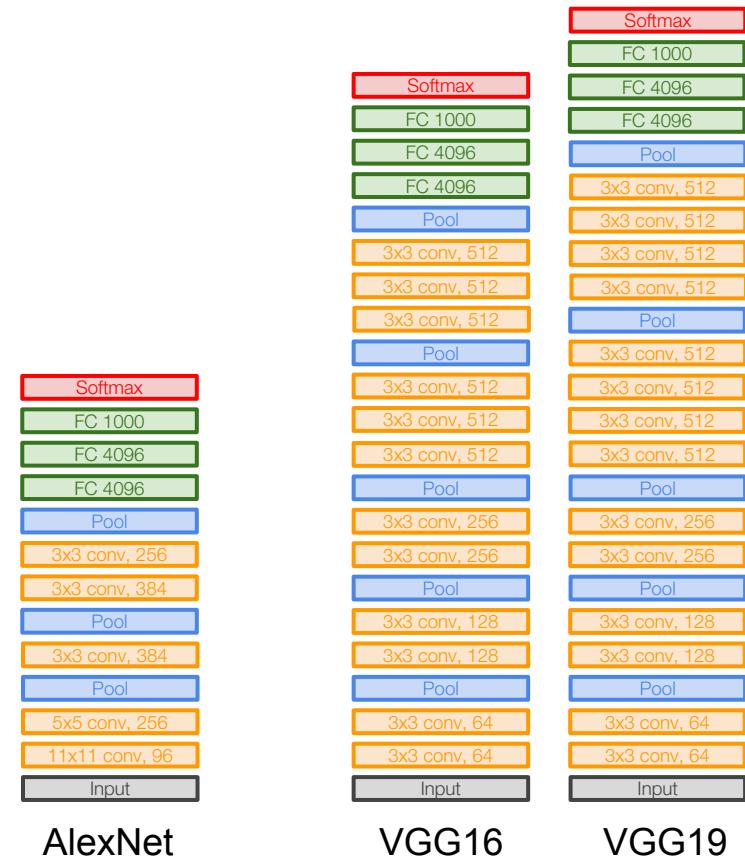
(ZFNet)

-> 7.3% top 5 error in ILSVRC'14



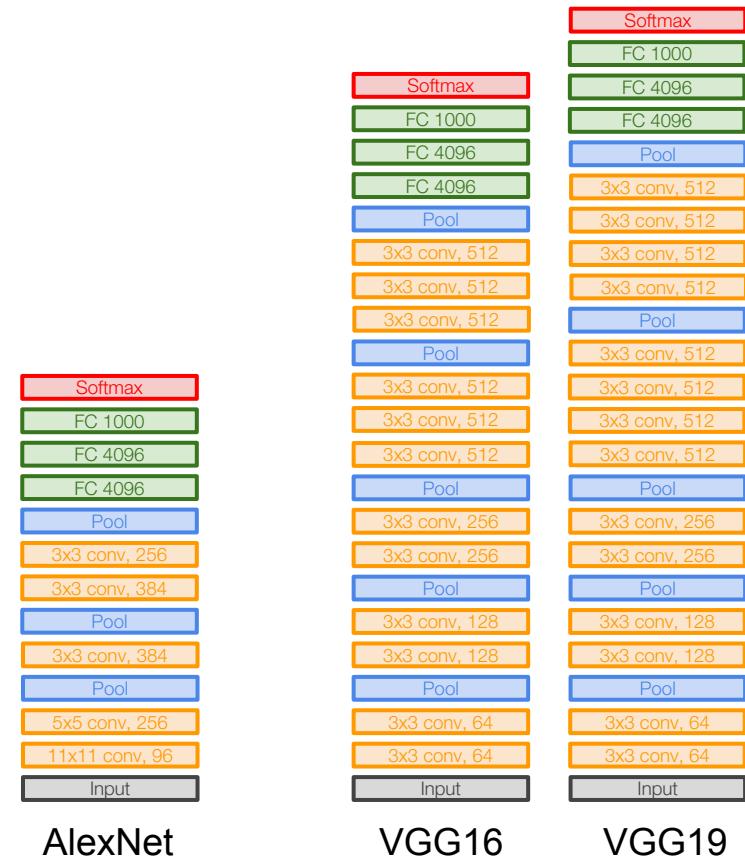
# VGGNet [Simonyan and Zisserman, 2014]

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



# VGGNet [Simonyan and Zisserman, 2014]

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

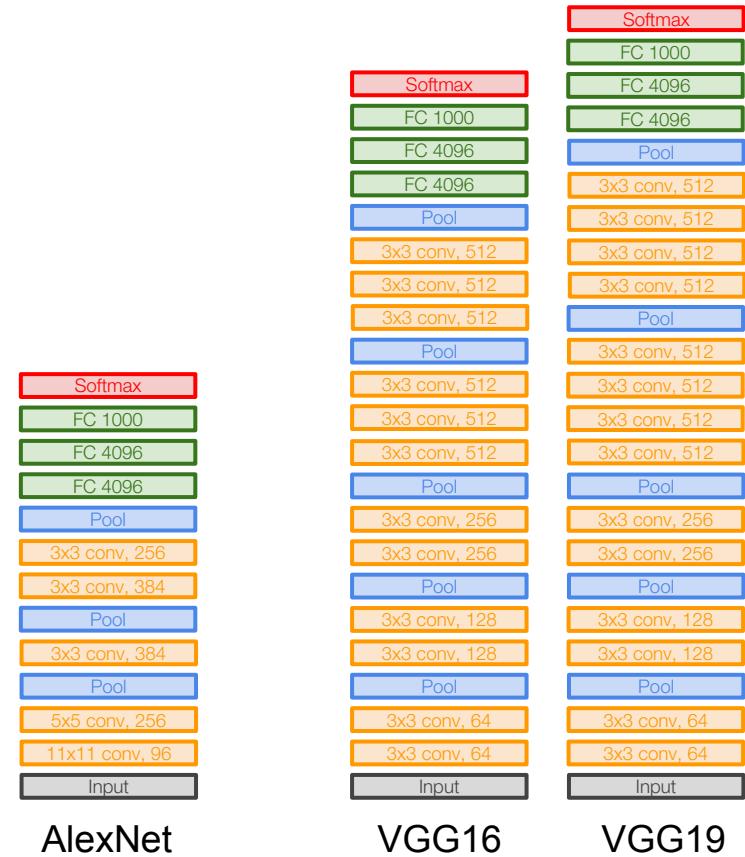


# VGGNet [Simonyan and Zisserman, 2014]

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]



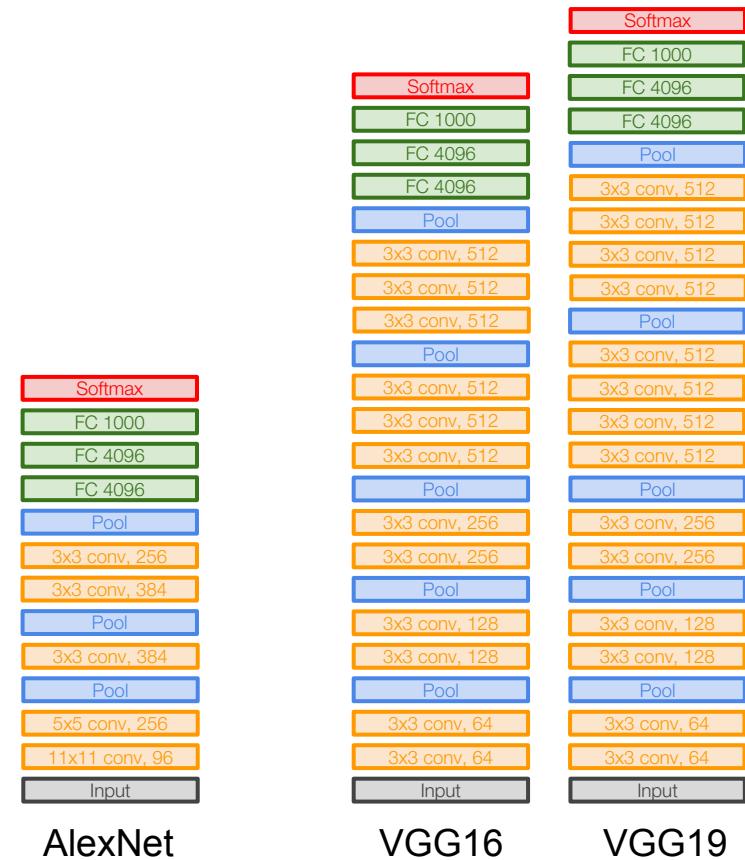
# VGGNet [Simonyan and Zisserman, 2014]

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



# VGGNet [Simonyan and Zisserman, 2014]

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

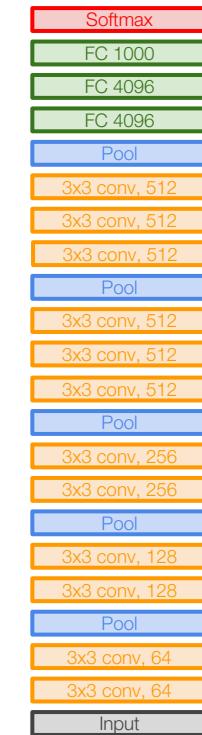
CONV3-512: [14x14x512]      memory:  $14 \times 14 \times 512 = 100K$     params:  $(3 \times 3 \times 512) \times 512 = 2,359,296$

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

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

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

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



VGG16

# VGGNet [Simonyan and Zisserman, 2014]

INPUT: [224x224x3] memory:  $224 \times 224 \times 3 = 150\text{K}$  params: 0 (not counting biases)

CONV3-64: [224x224x64] memory:  $224 \times 224 \times 64 = 3.2\text{M}$  params:  $(3 \times 3 \times 3) \times 64 = 1,728$

CONV3-64: [224x224x64] memory:  $224 \times 224 \times 64 = 3.2\text{M}$  params:  $(3 \times 3 \times 64) \times 64 = 36,864$

POOL2: [112x112x64] memory:  $112 \times 112 \times 64 = 800\text{K}$  params: 0

CONV3-128: [112x112x128] memory:  $112 \times 112 \times 128 = 1.6\text{M}$  params:  $(3 \times 3 \times 64) \times 128 = 73,728$

CONV3-128: [112x112x128] memory:  $112 \times 112 \times 128 = 1.6\text{M}$  params:  $(3 \times 3 \times 128) \times 128 = 147,456$

POOL2: [56x56x128] memory:  $56 \times 56 \times 128 = 400\text{K}$  params: 0

CONV3-256: [56x56x256] memory:  $56 \times 56 \times 256 = 800\text{K}$  params:  $(3 \times 3 \times 128) \times 256 = 294,912$

CONV3-256: [56x56x256] memory:  $56 \times 56 \times 256 = 800\text{K}$  params:  $(3 \times 3 \times 256) \times 256 = 589,824$

CONV3-256: [56x56x256] memory:  $56 \times 56 \times 256 = 800\text{K}$  params:  $(3 \times 3 \times 256) \times 256 = 589,824$

POOL2: [28x28x256] memory:  $28 \times 28 \times 256 = 200\text{K}$  params: 0

CONV3-512: [28x28x512] memory:  $28 \times 28 \times 512 = 400\text{K}$  params:  $(3 \times 3 \times 256) \times 512 = 1,179,648$

CONV3-512: [28x28x512] memory:  $28 \times 28 \times 512 = 400\text{K}$  params:  $(3 \times 3 \times 512) \times 512 = 2,359,296$

CONV3-512: [28x28x512] memory:  $28 \times 28 \times 512 = 400\text{K}$  params:  $(3 \times 3 \times 512) \times 512 = 2,359,296$

POOL2: [14x14x512] memory:  $14 \times 14 \times 512 = 100\text{K}$  params: 0

CONV3-512: [14x14x512] memory:  $14 \times 14 \times 512 = 100\text{K}$  params:  $(3 \times 3 \times 512) \times 512 = 2,359,296$

CONV3-512: [14x14x512] memory:  $14 \times 14 \times 512 = 100\text{K}$  params:  $(3 \times 3 \times 512) \times 512 = 2,359,296$

CONV3-512: [14x14x512] memory:  $14 \times 14 \times 512 = 100\text{K}$  params:  $(3 \times 3 \times 512) \times 512 = 2,359,296$

POOL2: [7x7x512] memory:  $7 \times 7 \times 512 = 25\text{K}$  params: 0

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

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

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

Note:

Most memory is in early CONV

Most params are in late FC

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

TOTAL params: 138M parameters

# VGGNet [Simonyan and Zisserman, 2014]

INPUT: [224x224x3]      memory:  $224 \times 224 \times 3 = 150\text{K}$     params: 0      (not counting biases)

CONV3-64: [224x224x64]      memory:  $224 \times 224 \times 64 = 3.2\text{M}$     params:  $(3 \times 3 \times 3) \times 64 = 1,728$

CONV3-64: [224x224x64]      memory:  $224 \times 224 \times 64 = 3.2\text{M}$     params:  $(3 \times 3 \times 64) \times 64 = 36,864$

POOL2: [112x112x64]      memory:  $112 \times 112 \times 64 = 800\text{K}$     params: 0

CONV3-128: [112x112x128]      memory:  $112 \times 112 \times 128 = 1.6\text{M}$     params:  $(3 \times 3 \times 64) \times 128 = 73,728$

CONV3-128: [112x112x128]      memory:  $112 \times 112 \times 128 = 1.6\text{M}$     params:  $(3 \times 3 \times 128) \times 128 = 147,456$

POOL2: [56x56x128]      memory:  $56 \times 56 \times 128 = 400\text{K}$     params: 0

CONV3-256: [56x56x256]      memory:  $56 \times 56 \times 256 = 800\text{K}$     params:  $(3 \times 3 \times 128) \times 256 = 294,912$

CONV3-256: [56x56x256]      memory:  $56 \times 56 \times 256 = 800\text{K}$     params:  $(3 \times 3 \times 256) \times 256 = 589,824$

CONV3-256: [56x56x256]      memory:  $56 \times 56 \times 256 = 800\text{K}$     params:  $(3 \times 3 \times 256) \times 256 = 589,824$

POOL2: [28x28x256]      memory:  $28 \times 28 \times 256 = 200\text{K}$     params: 0

CONV3-512: [28x28x512]      memory:  $28 \times 28 \times 512 = 400\text{K}$     params:  $(3 \times 3 \times 256) \times 512 = 1,179,648$

CONV3-512: [28x28x512]      memory:  $28 \times 28 \times 512 = 400\text{K}$     params:  $(3 \times 3 \times 512) \times 512 = 2,359,296$

CONV3-512: [28x28x512]      memory:  $28 \times 28 \times 512 = 400\text{K}$     params:  $(3 \times 3 \times 512) \times 512 = 2,359,296$

POOL2: [14x14x512]      memory:  $14 \times 14 \times 512 = 100\text{K}$     params: 0

CONV3-512: [14x14x512]      memory:  $14 \times 14 \times 512 = 100\text{K}$     params:  $(3 \times 3 \times 512) \times 512 = 2,359,296$

CONV3-512: [14x14x512]      memory:  $14 \times 14 \times 512 = 100\text{K}$     params:  $(3 \times 3 \times 512) \times 512 = 2,359,296$

CONV3-512: [14x14x512]      memory:  $14 \times 14 \times 512 = 100\text{K}$     params:  $(3 \times 3 \times 512) \times 512 = 2,359,296$

POOL2: [7x7x512]      memory:  $7 \times 7 \times 512 = 25\text{K}$     params: 0

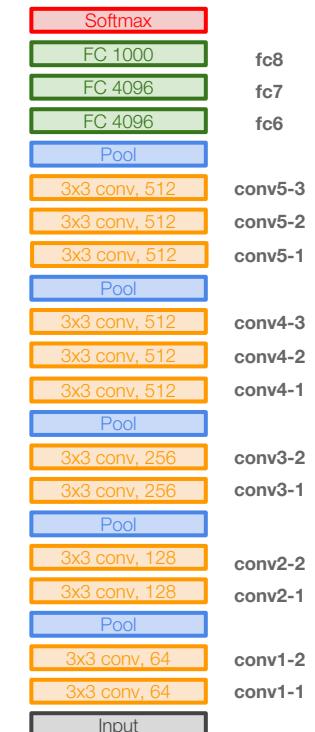
FC: [1x1x4096]      memory: 4096    params:  $7 \times 7 \times 512 \times 4096 = 102,760,448$

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

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

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

**TOTAL** params: 138M parameters



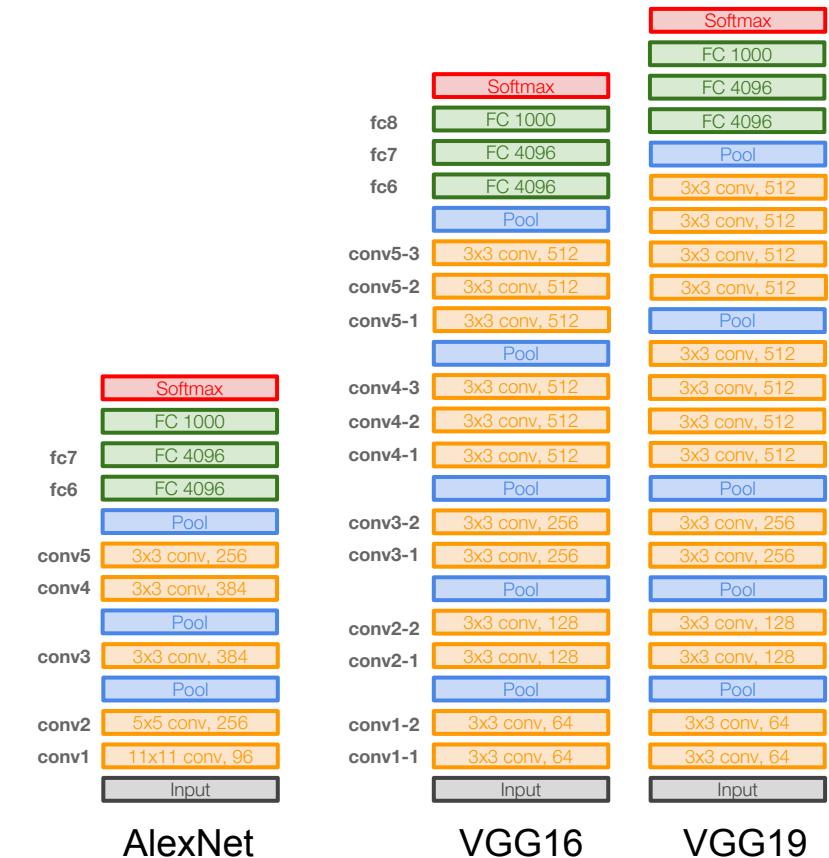
VGG16

Common names

# VGGNet [Simonyan and Zisserman, 2014]

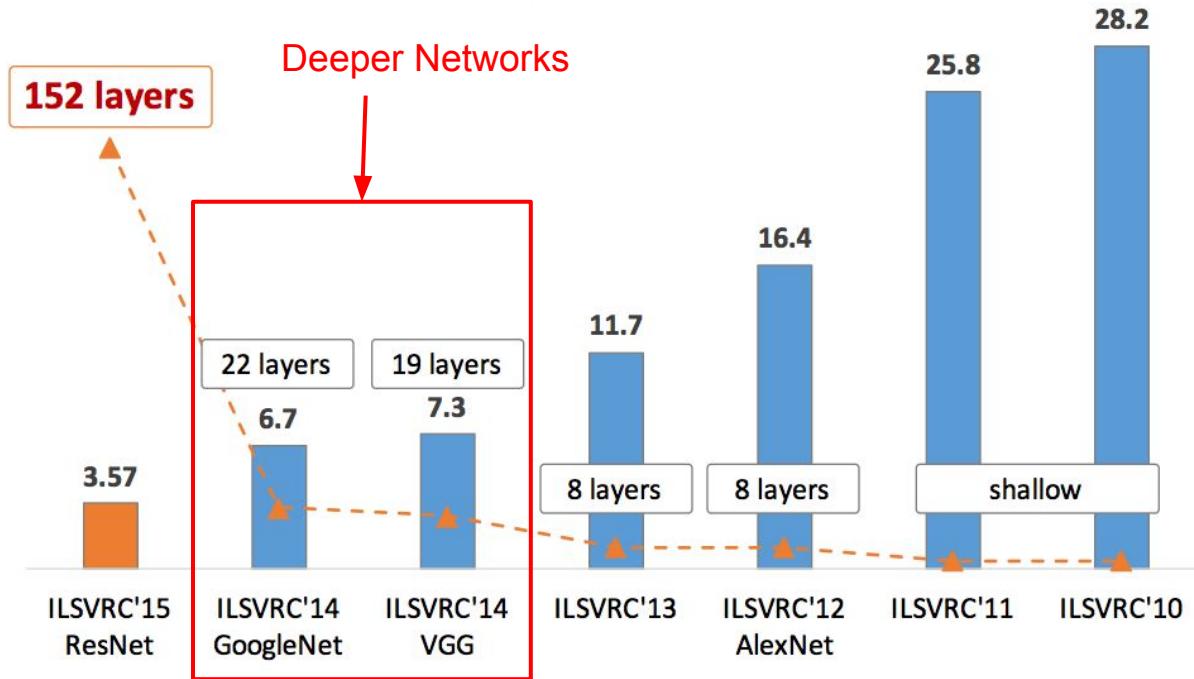
## Details:

- ILSVRC'14 2nd in classification, 1st in localization
- Similar training procedure as Krizhevsky 2012
- No Local Response Normalisation (LRN)
- Use VGG16 or VGG19 (VGG19 only slightly better, more memory)
- Use ensembles for best results
- FC7 features generalize well to other tasks



# GoogLeNet [Szegedy et al., 2014]

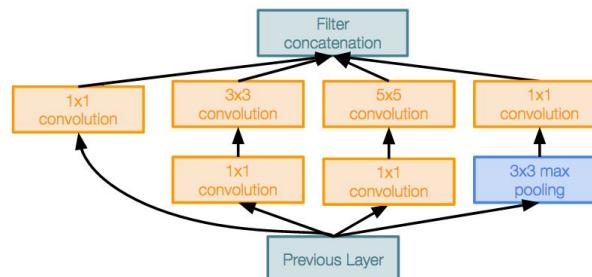
ImageNet Large Scale Visual Recognition Challenge (ILSVRC) winners



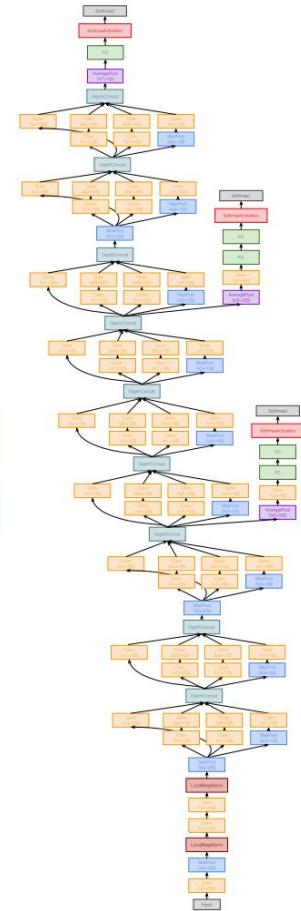
# GoogLeNet [Szegedy et al., 2014]

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)

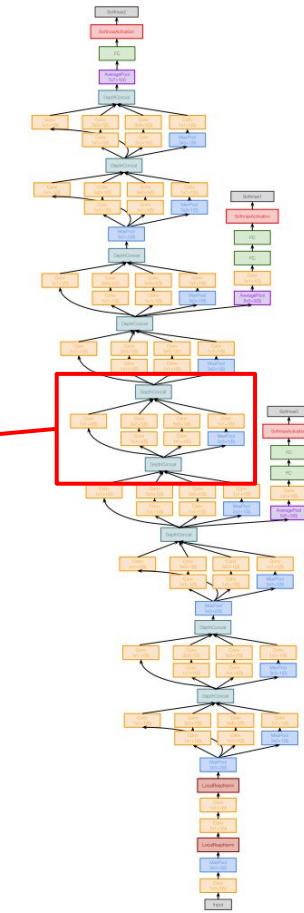
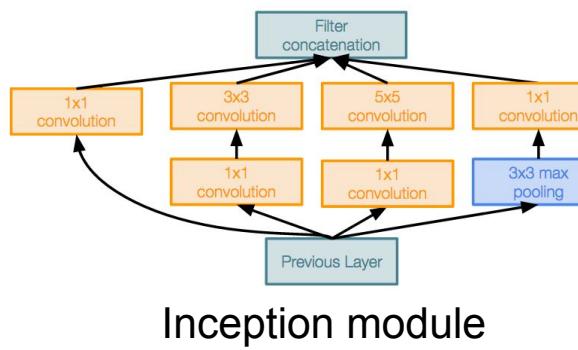


Inception module

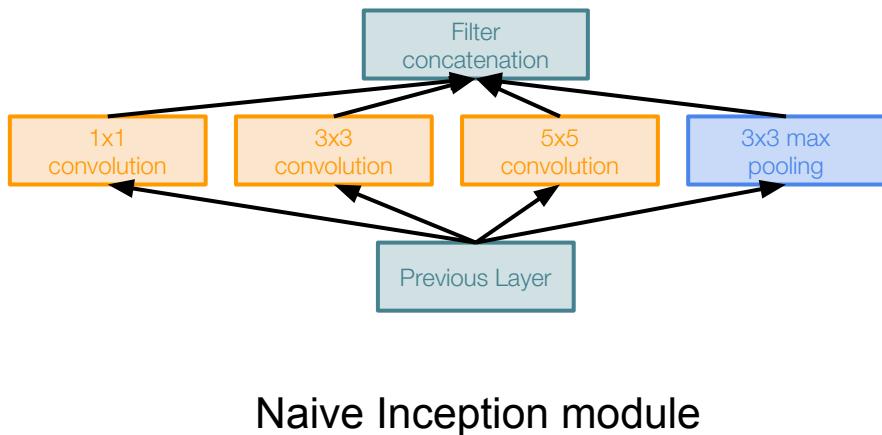


# GoogLeNet [Szegedy et al., 2014]

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



# GoogLeNet [Szegedy et al., 2014]

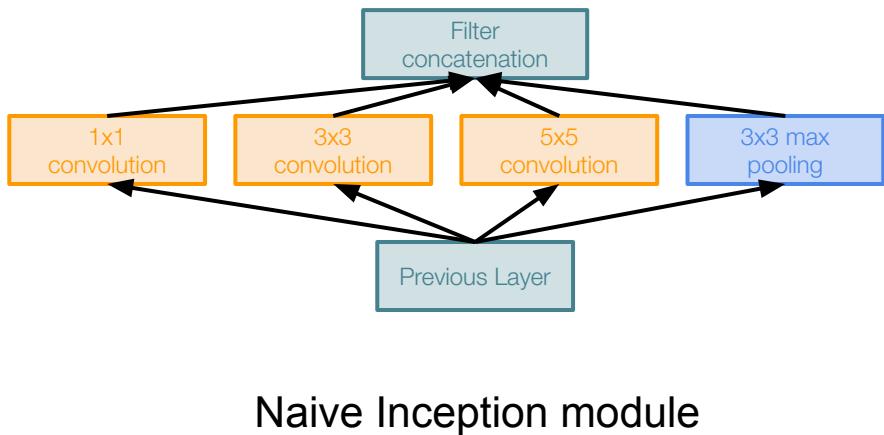


Apply parallel filter operations on the input from previous layer:

- Multiple receptive field sizes for convolution ( $1 \times 1$ ,  $3 \times 3$ ,  $5 \times 5$ )
- Pooling operation ( $3 \times 3$ )

Concatenate all filter outputs together depth-wise

# GoogLeNet [Szegedy et al., 2014]



Apply parallel filter operations on the input from previous layer:

- Multiple receptive field sizes for convolution ( $1 \times 1$ ,  $3 \times 3$ ,  $5 \times 5$ )
- Pooling operation ( $3 \times 3$ )

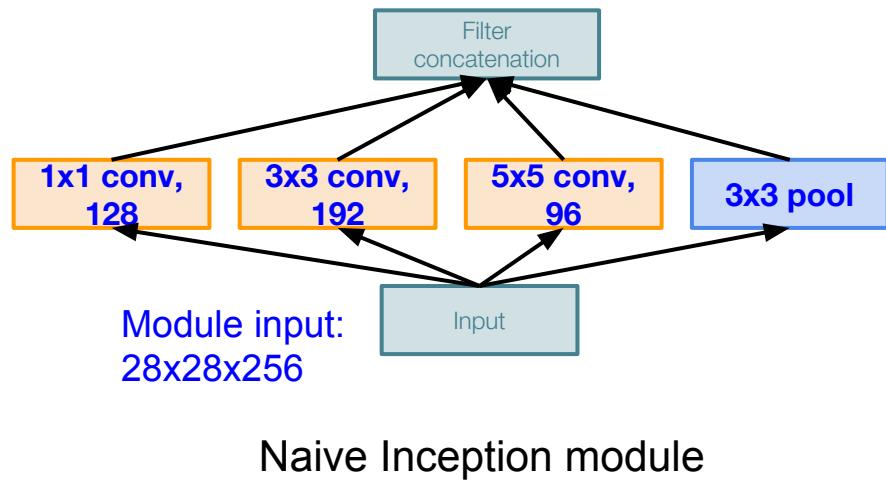
Concatenate all filter outputs together depth-wise

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

# GoogLeNet [Szegedy et al., 2014]

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

Example:

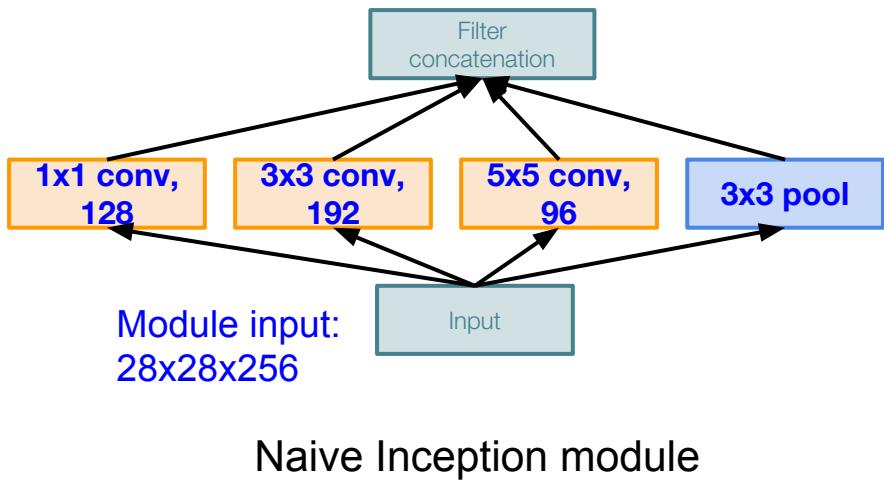


# GoogLeNet [Szegedy et al., 2014]

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

Example:

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

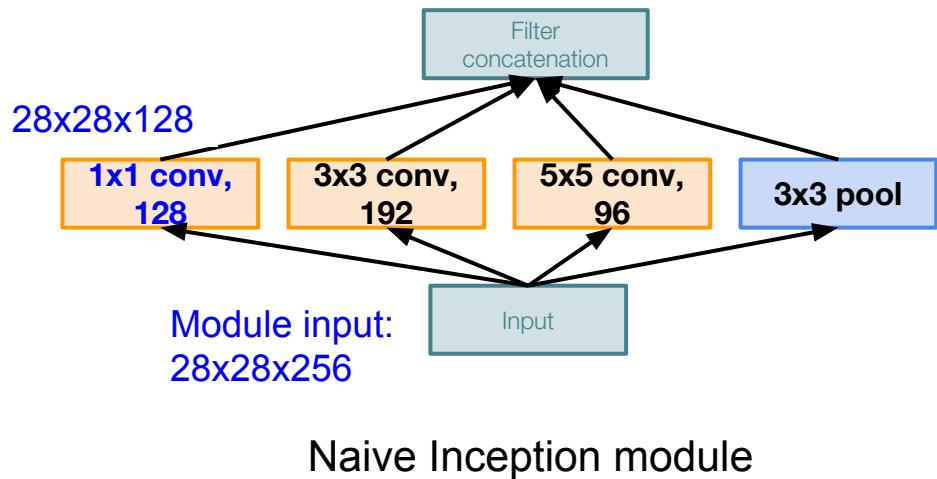


# GoogLeNet [Szegedy et al., 2014]

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

Example:

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

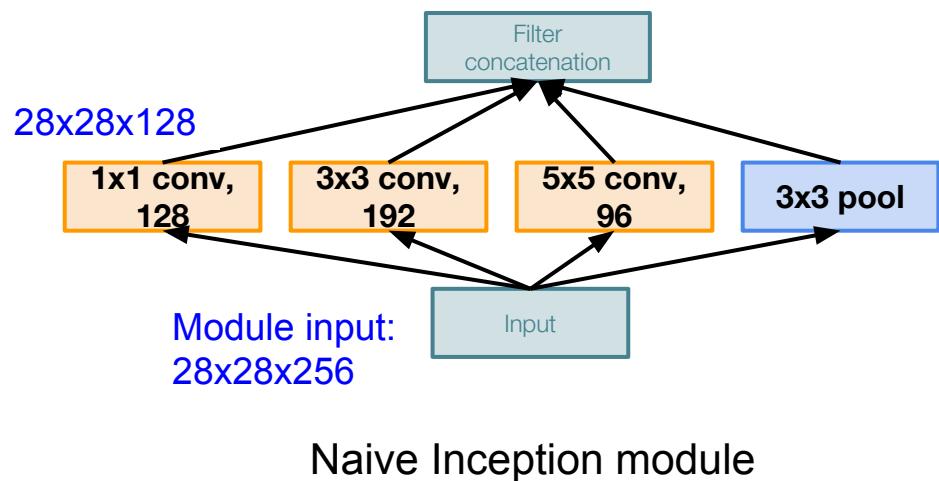


# GoogLeNet [Szegedy et al., 2014]

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

Example:

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

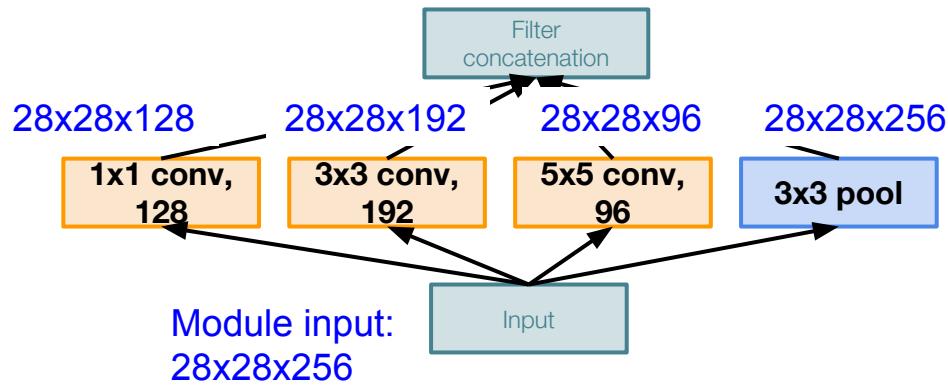


# GoogLeNet [Szegedy et al., 2014]

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

Example:

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



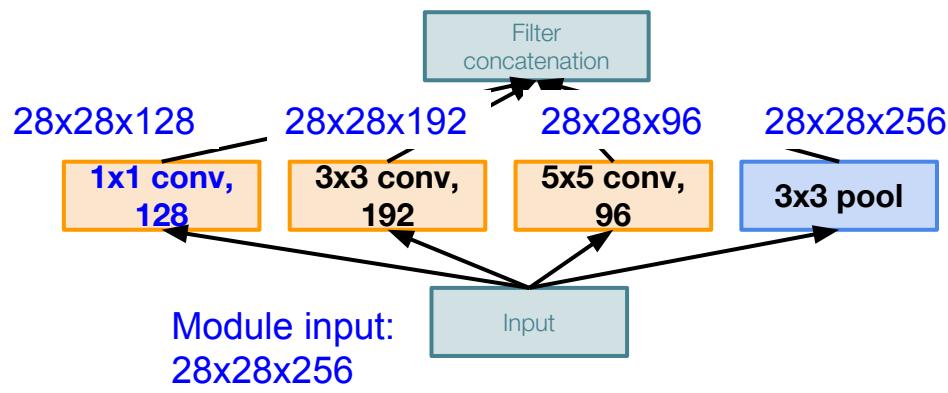
Naive Inception module

# GoogLeNet [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

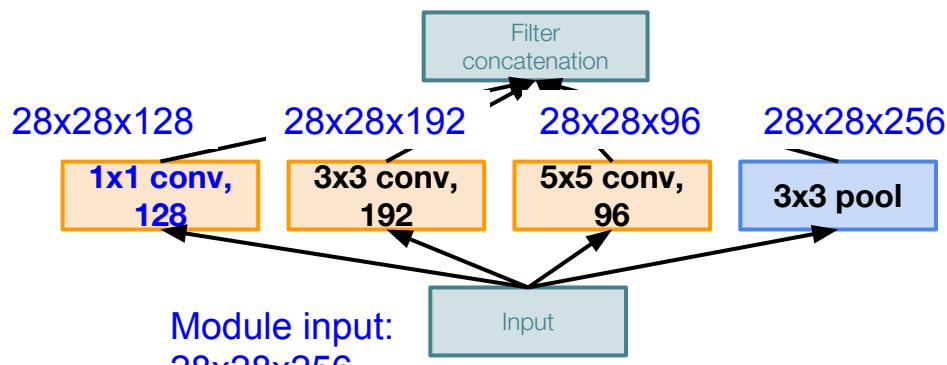
# GoogLeNet [Szegedy et al., 2014]

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

Example:

Q3:What is output size after  
filter concatenation?

$$28 \times 28 \times (128 + 192 + 96 + 256) = 28 \times 28 \times 672$$



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

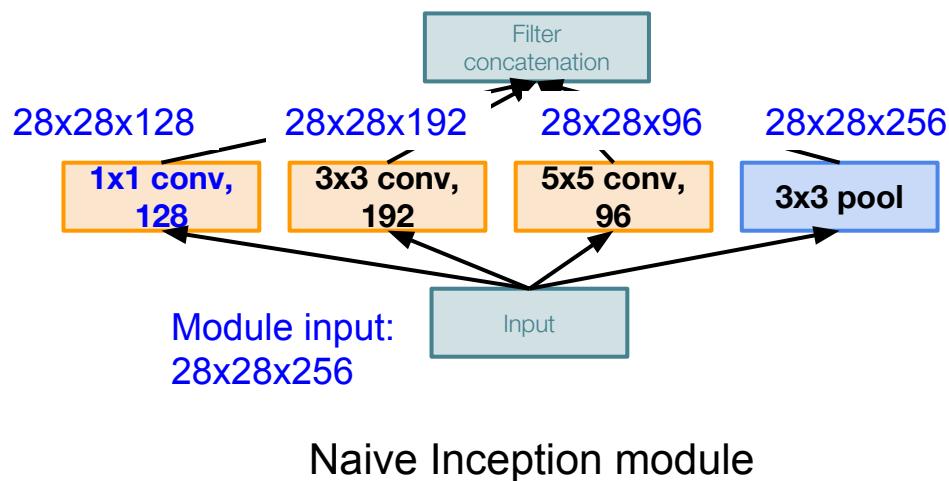
# GoogLeNet [Szegedy et al., 2014]

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

Example:

Q3:What is output size after  
filter concatenation?

$$28 \times 28 \times (128 + 192 + 96 + 256) = 28 \times 28 \times 672$$



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!

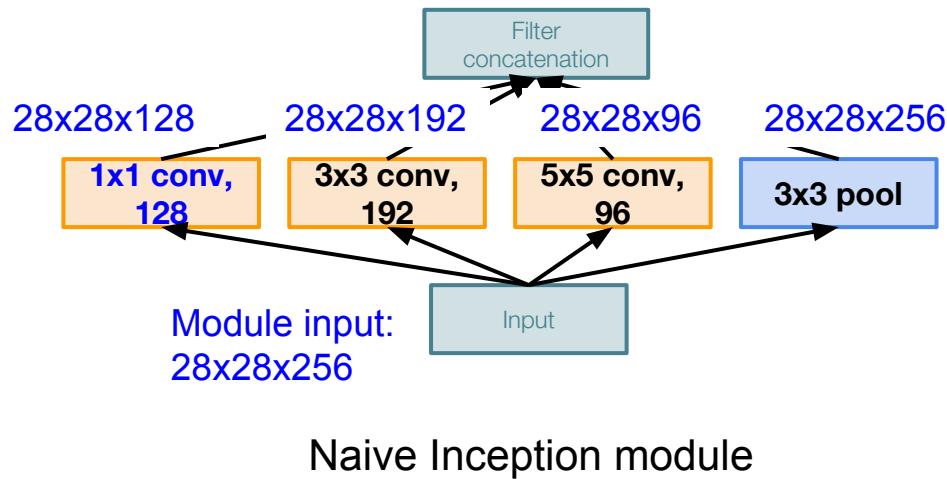
# GoogLeNet [Szegedy et al., 2014]

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

Example:

Q3:What is output size after  
filter concatenation?

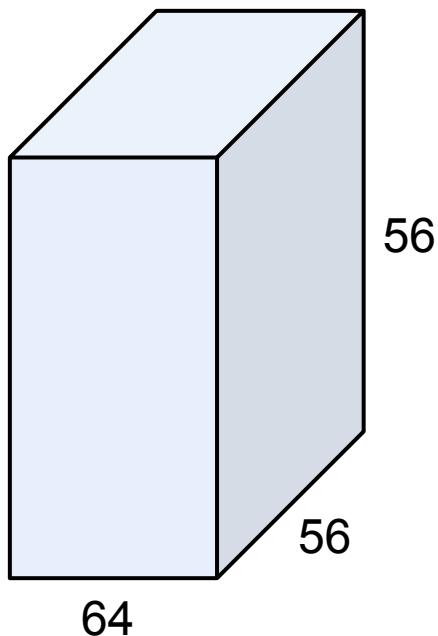
$$28 \times 28 \times (128 + 192 + 96 + 256) = 529k$$



Solution: “bottleneck” layers that use 1x1 convolutions to reduce feature depth

# GoogLeNet [Szegedy et al., 2014]

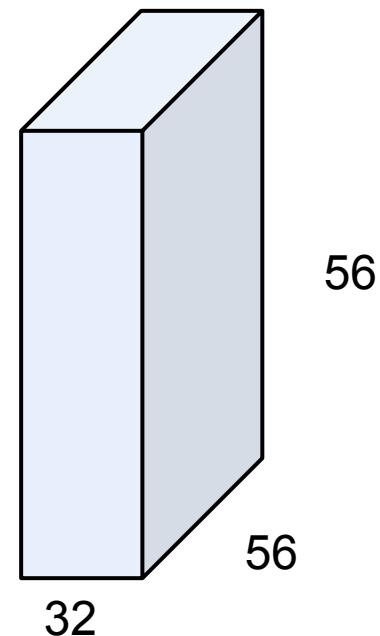
Reminder: 1x1 convolutions



1x1 CONV  
with 32 filters

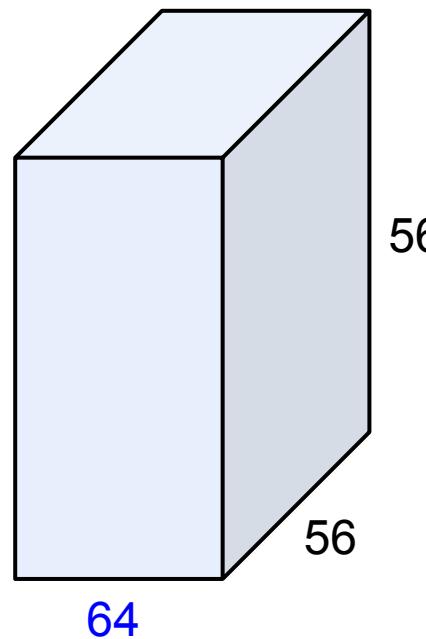
---

(each filter has size  
1x1x64, and performs a  
64-dimensional dot  
product)

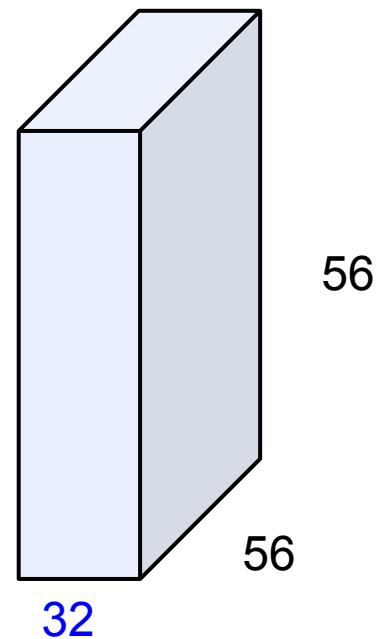


# GoogLeNet [Szegedy et al., 2014]

Reminder: 1x1 convolutions



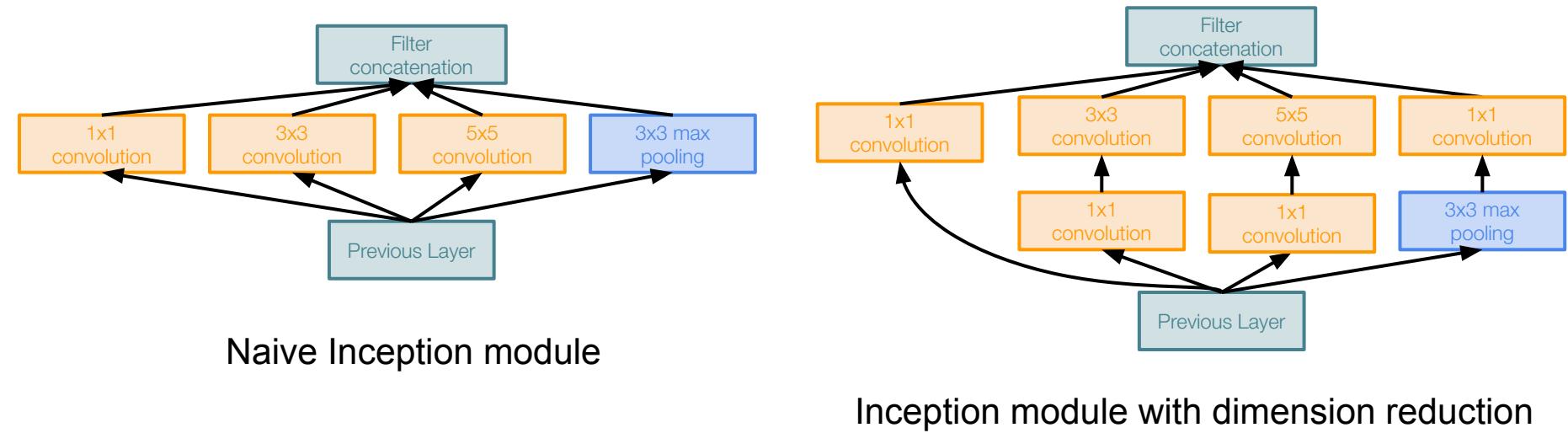
1x1 CONV  
with 32 filters



preserves spatial  
dimensions, reduces depth!

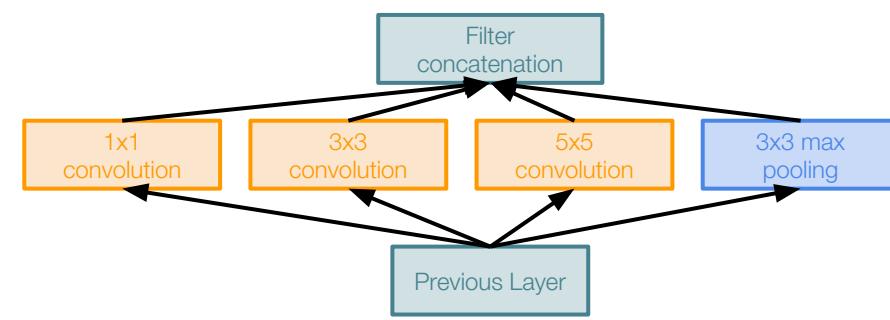
Projects depth to lower  
dimension (combination of  
feature maps)

# GoogLeNet [Szegedy et al., 2014]

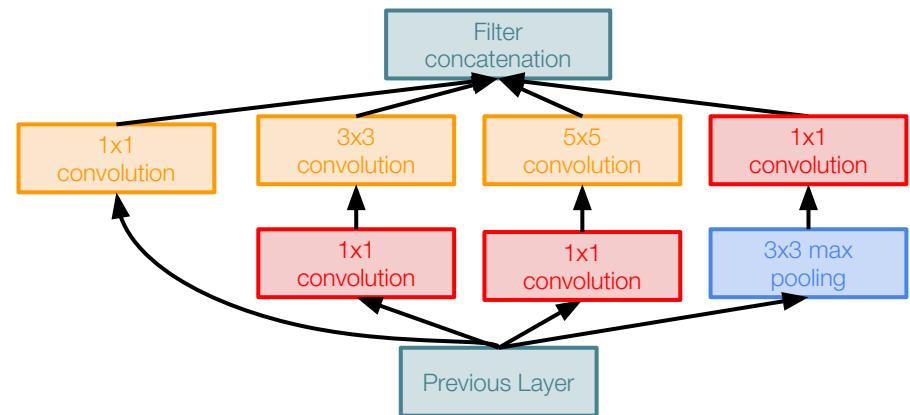


# GoogLeNet [Szegedy et al., 2014]

1x1 conv “bottleneck”  
layers

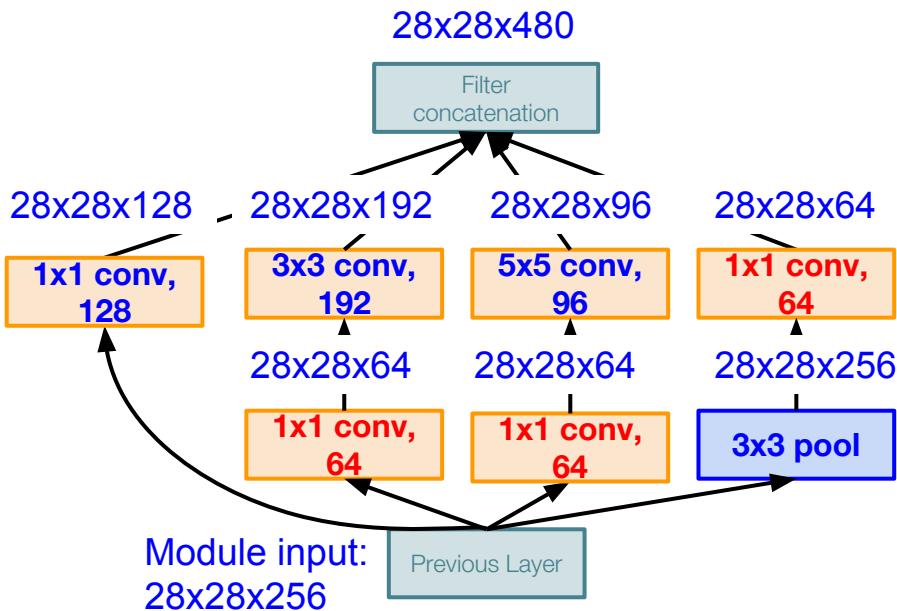


Naive Inception module



Inception module with dimension reduction

# 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:

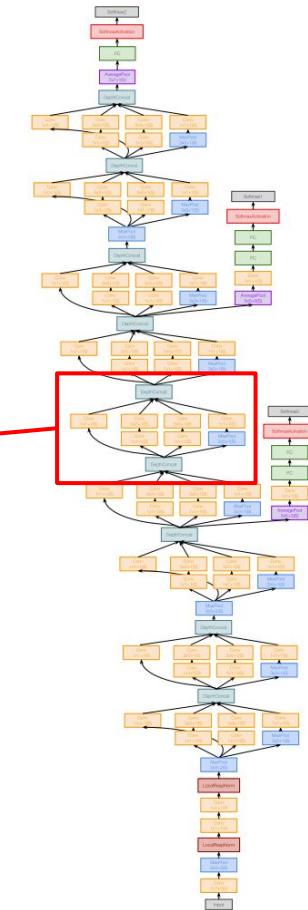
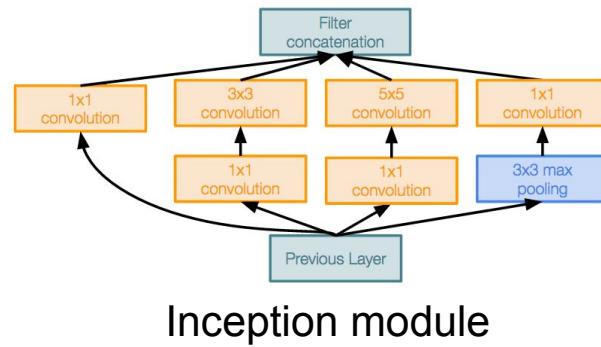
- [1x1 conv, 64] 28x28x64x1x1x256
- [1x1 conv, 64] 28x28x64x1x1x256
- [1x1 conv, 128] 28x28x128x1x1x256
- [3x3 conv, 192] 28x28x192x3x3x64
- [5x5 conv, 96] 28x28x96x5x5x64
- [1x1 conv, 64] 28x28x64x1x1x256

Total: 358M ops

Compared to 854M ops for naive version  
Bottleneck can also reduce depth after pooling layer

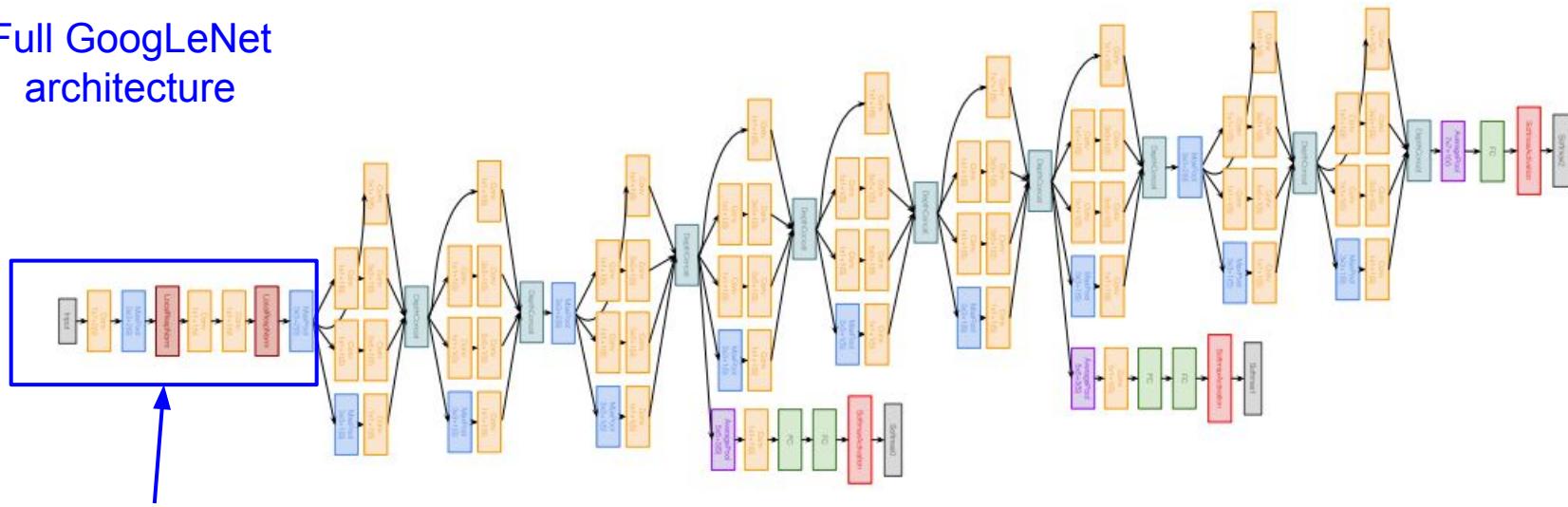
# GoogLeNet [Szegedy et al., 2014]

Stack Inception modules  
with dimension reduction  
on top of each other



# GoogLeNet [Szegedy et al., 2014]

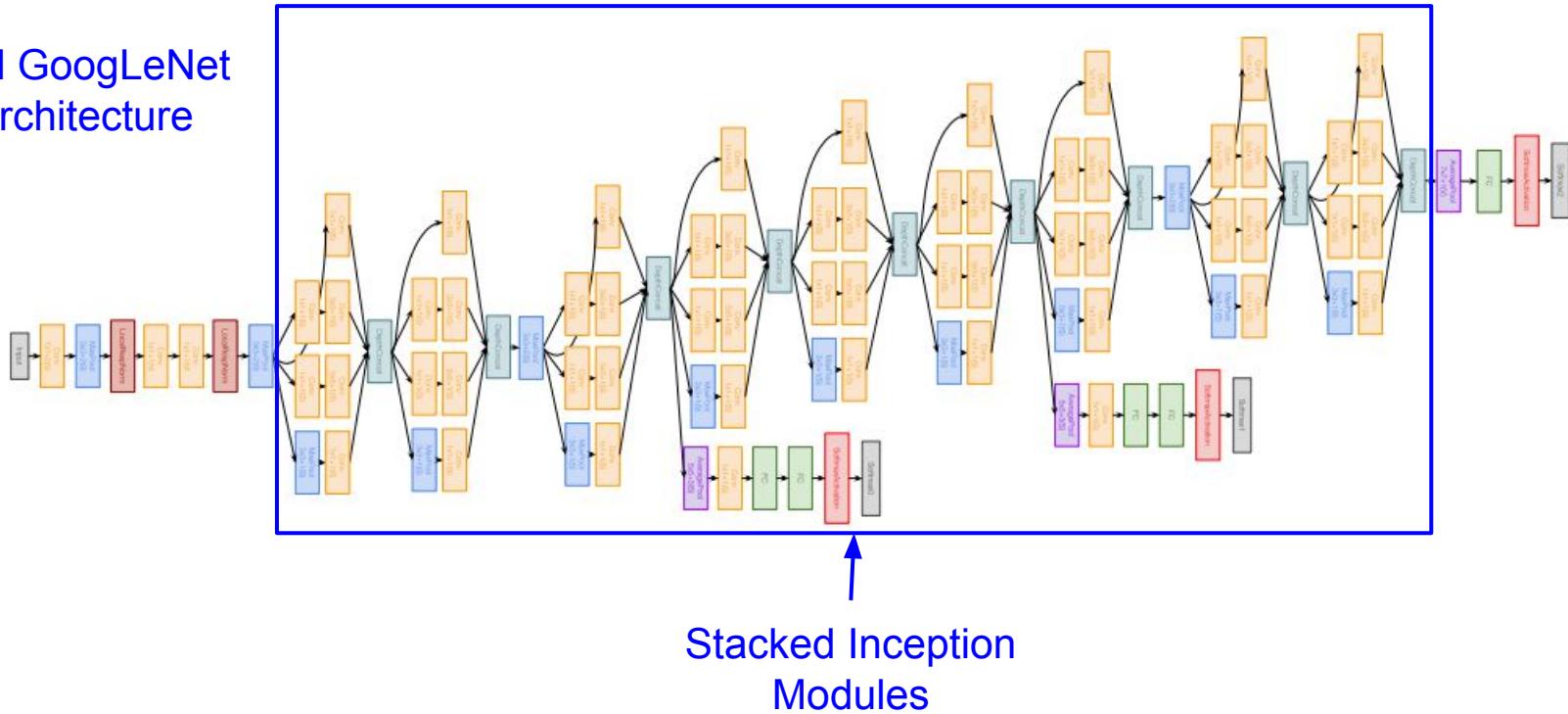
Full GoogLeNet  
architecture



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

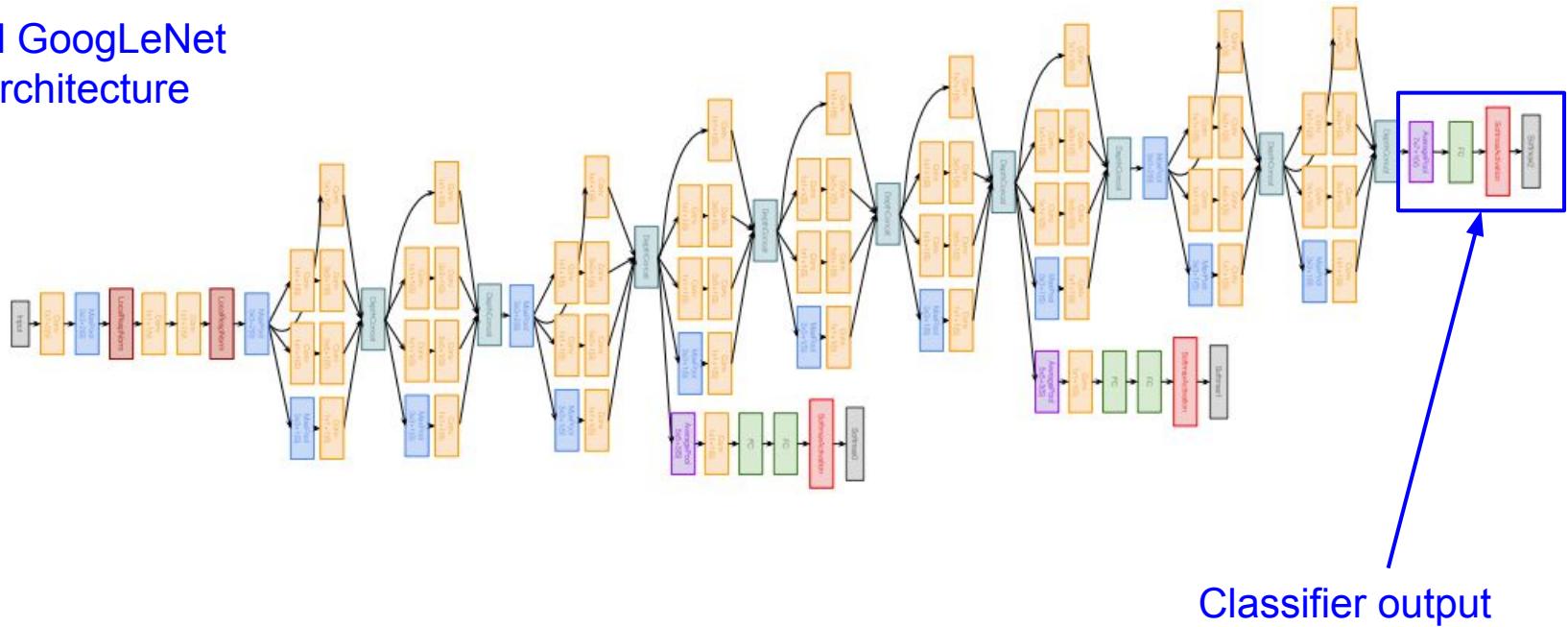
# GoogLeNet [Szegedy et al., 2014]

Full GoogLeNet  
architecture



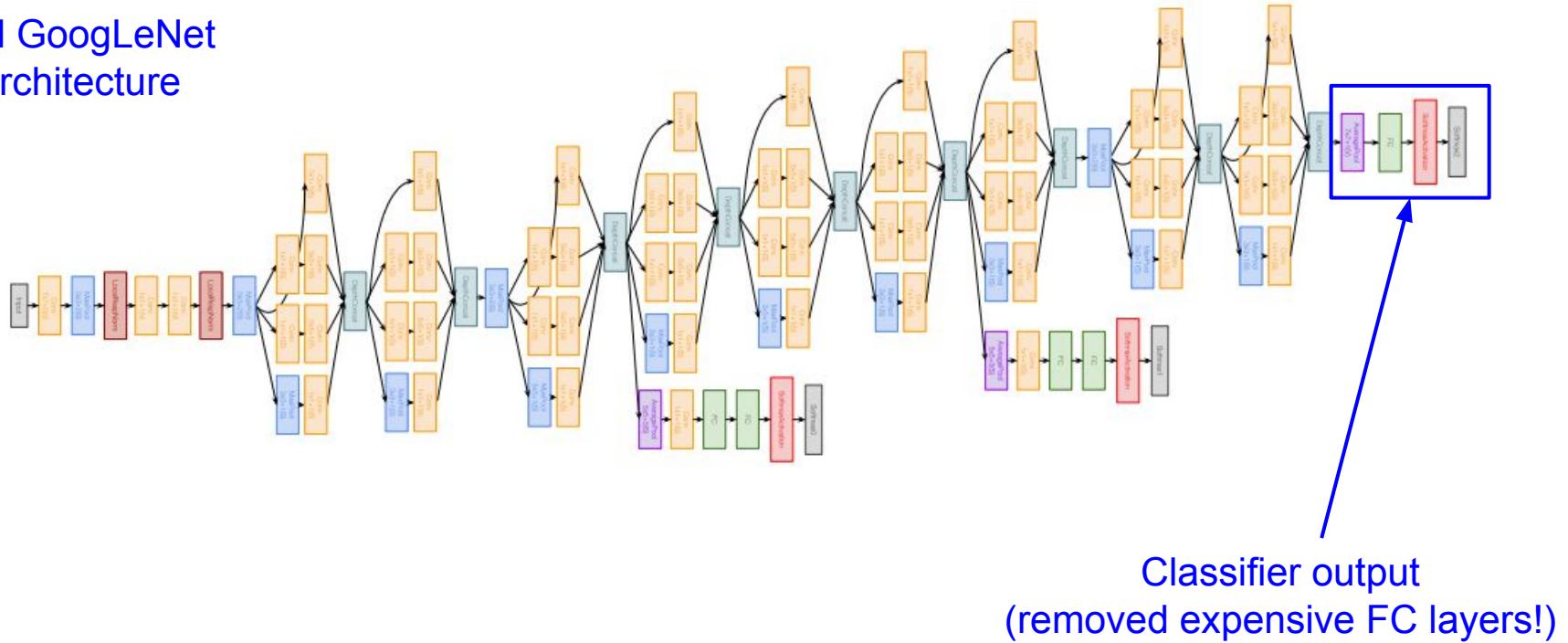
# GoogLeNet [Szegedy et al., 2014]

Full GoogLeNet  
architecture



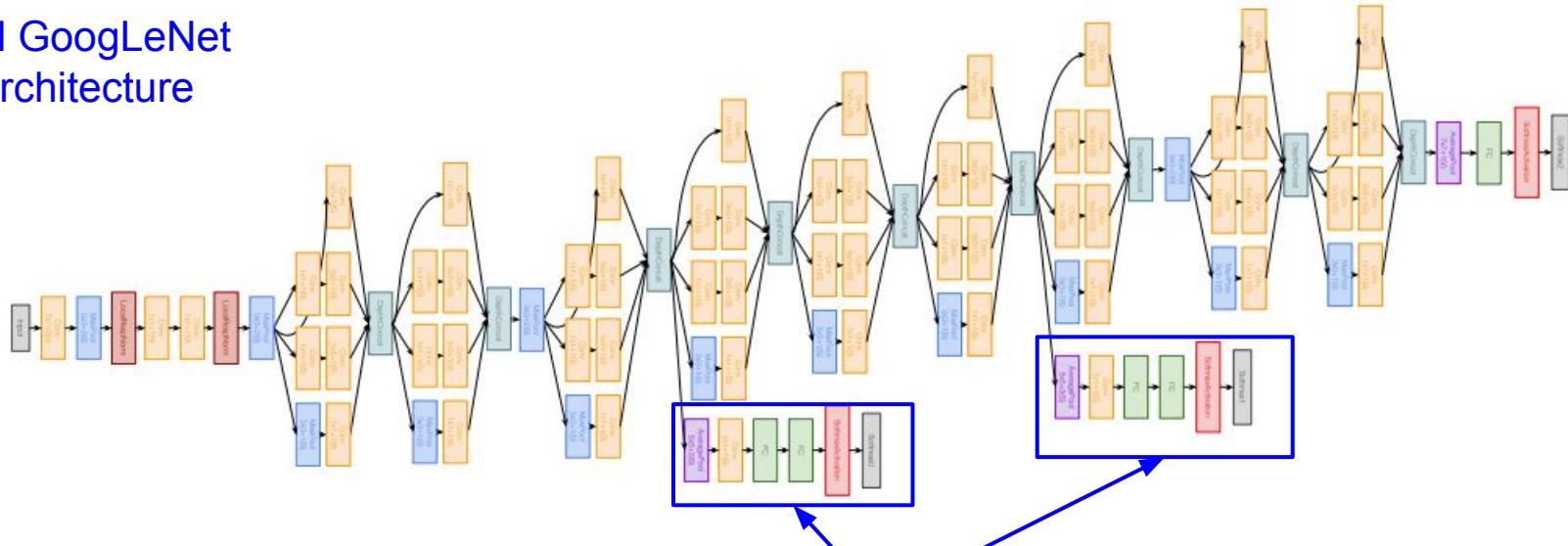
# GoogLeNet [Szegedy et al., 2014]

Full GoogLeNet  
architecture



# GoogLeNet [Szegedy et al., 2014]

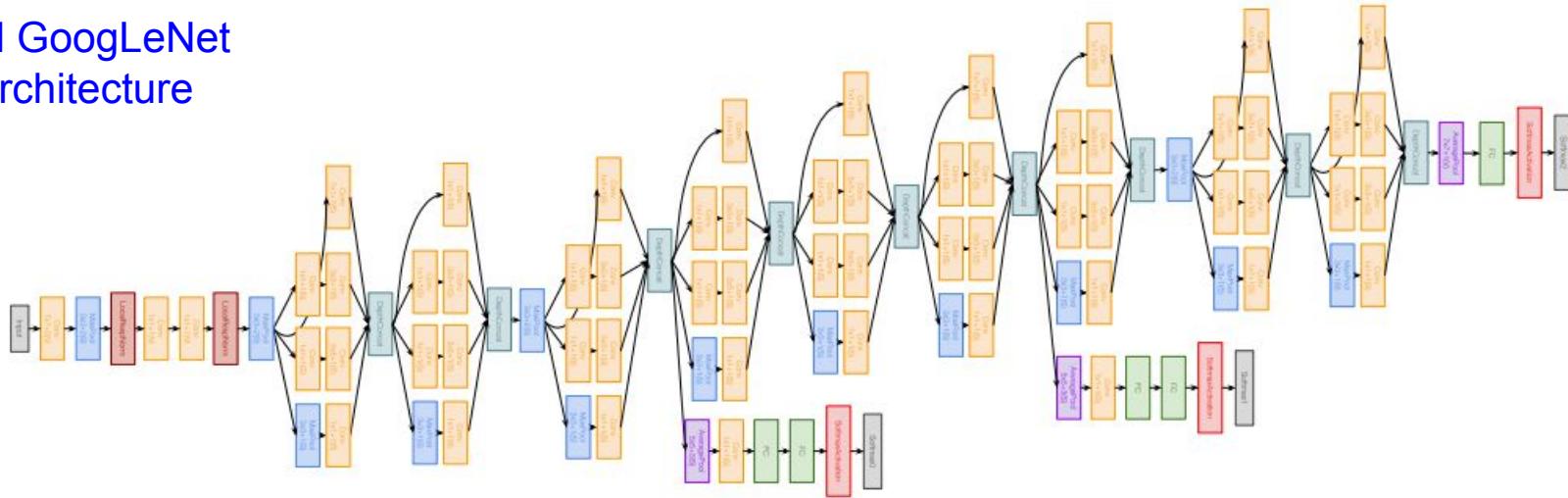
Full GoogLeNet  
architecture



Auxiliary classification outputs to inject additional gradient at lower layers  
(AvgPool-1x1Conv-FC-FC-Softmax)

# GoogLeNet [Szegedy et al., 2014]

Full GoogLeNet  
architecture

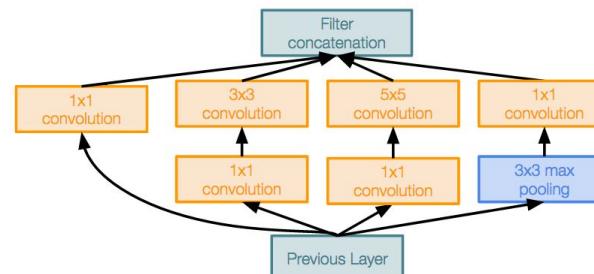


22 total layers with weights (including each parallel layer in an Inception module)

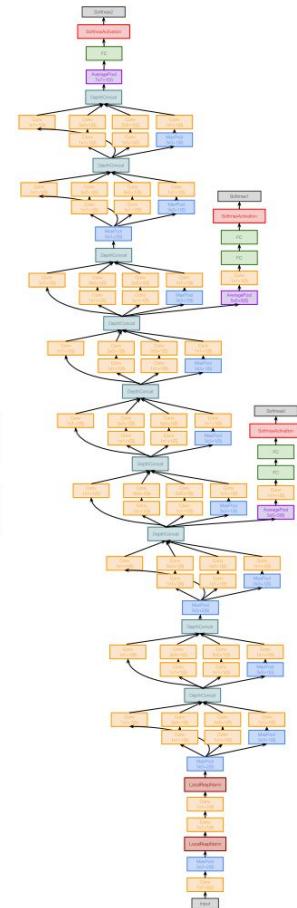
# GoogLeNet [Szegedy et al., 2014]

Deeper networks, with computational efficiency

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

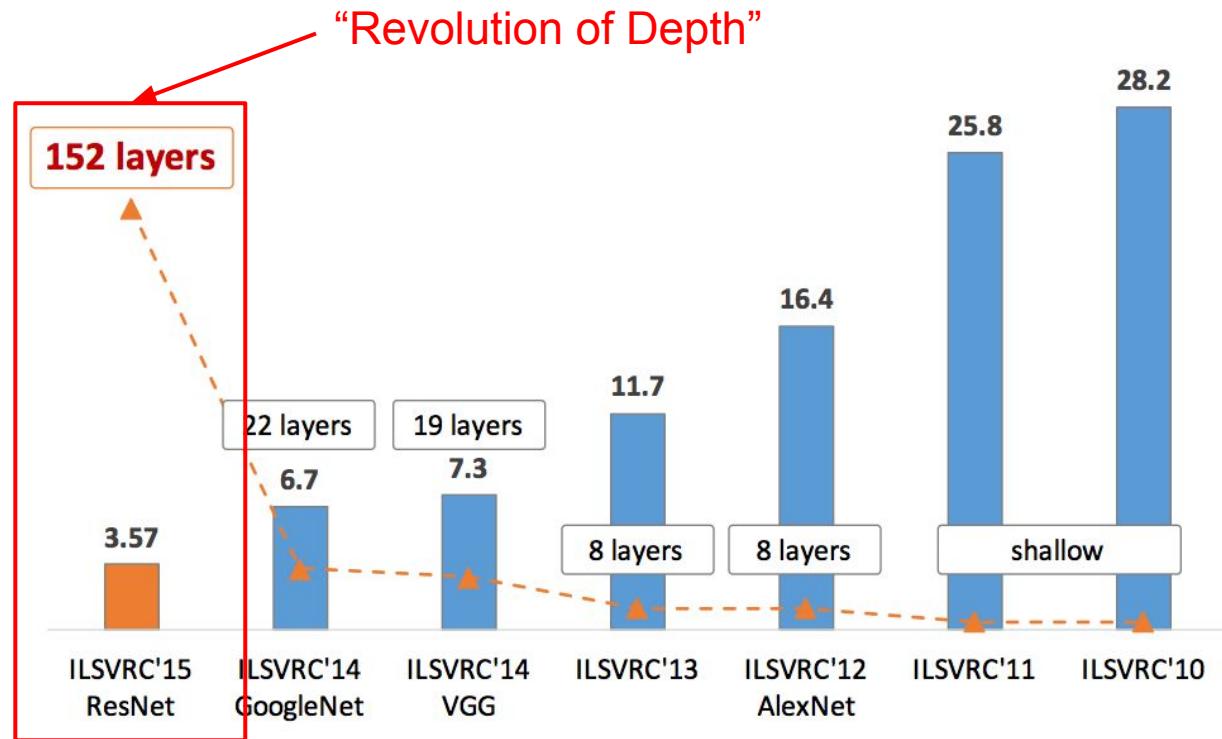


Inception module



# ResNet [He et al., 2015]

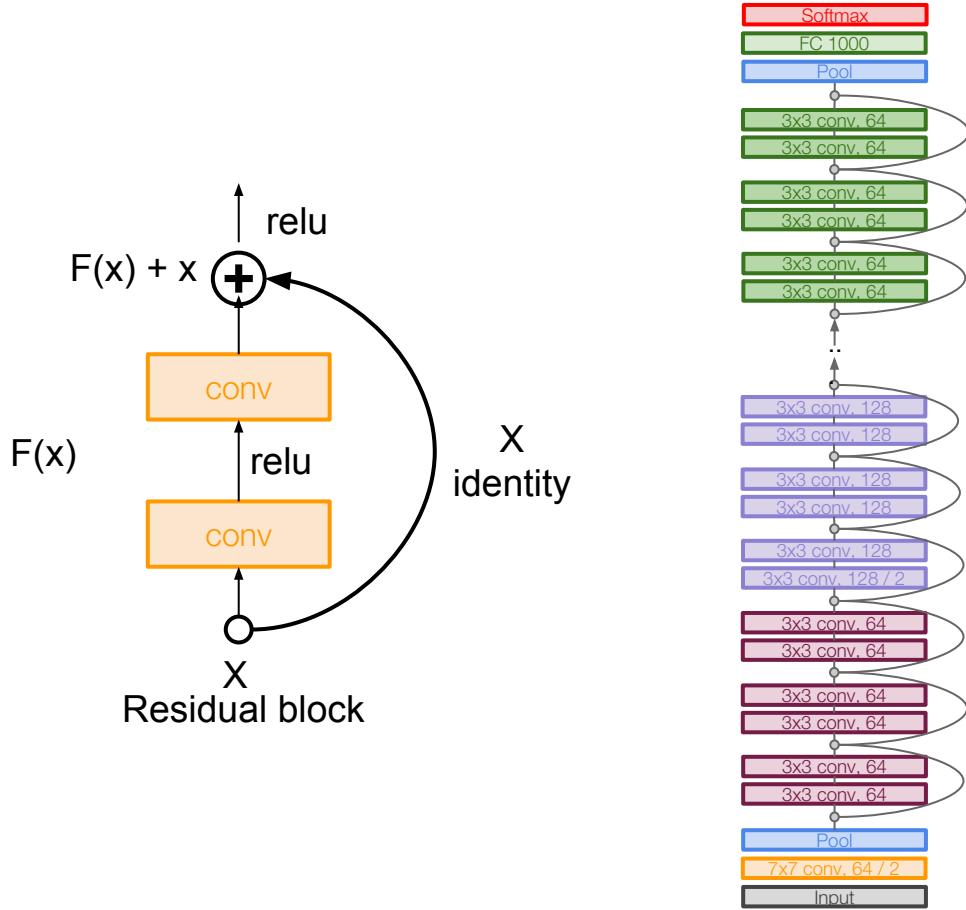
ImageNet Large Scale Visual Recognition Challenge (ILSVRC) winners



# ResNet [He et al., 2015]

Very deep networks using residual connections

- 152-layer model for ImageNet
- ILSVRC'15 classification winner (3.57% top 5 error)
- Swept all classification and detection competitions in ILSVRC'15 and COCO'15!

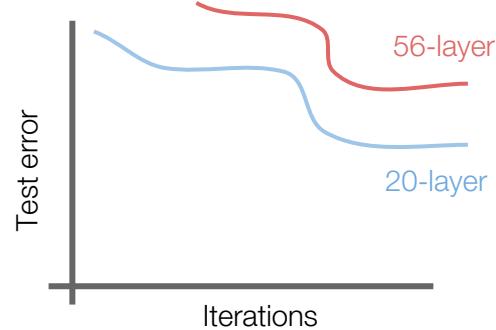
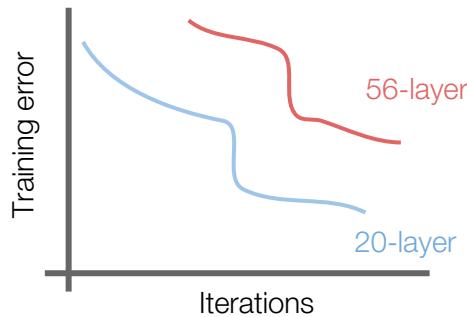


# ResNet [He et al., 2015]

What happens when we continue stacking deeper layers on a “plain” convolutional neural network?

# ResNet [He et al., 2015]

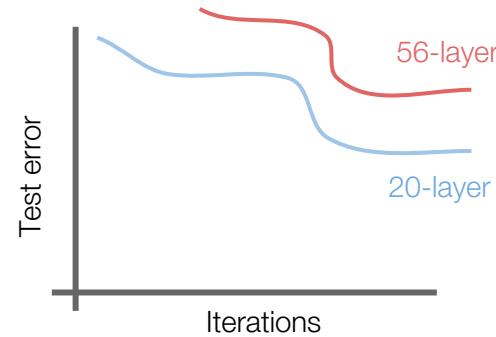
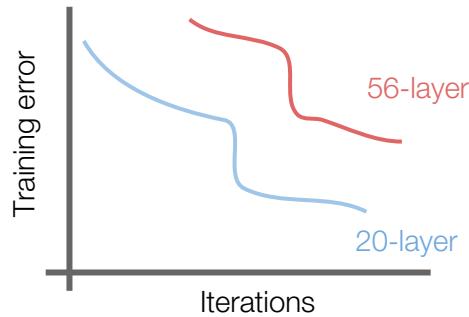
What happens when we continue stacking deeper layers on a “plain” convolutional neural network?



Q: What's strange about these training and test curves?  
[Hint: look at the order of the curves]

# ResNet [He et al., 2015]

What happens when we continue stacking deeper layers on a “plain” convolutional neural network?



56-layer model performs worse on both training and test error  
-> The deeper model performs worse, but it's not caused by overfitting!

# ResNet [He et al., 2015]

Hypothesis: the problem is an *optimization* problem, deeper models are harder to optimize

# ResNet [He et al., 2015]

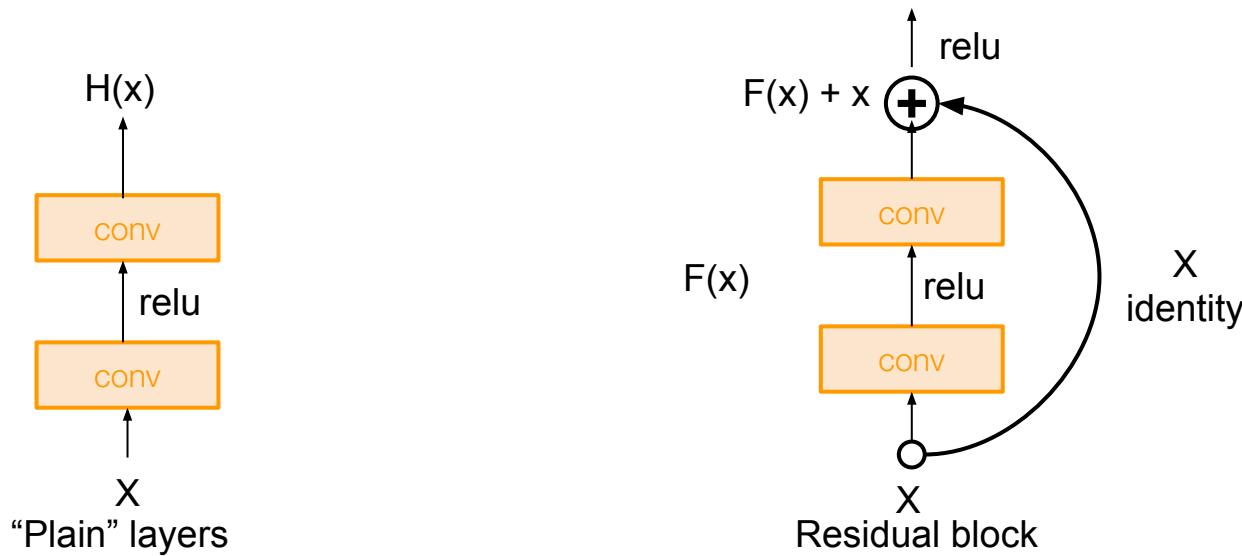
Hypothesis: the problem is an *optimization* problem, deeper models are harder to optimize

The deeper model should be able to perform at least as well as the shallower model.

A solution by construction is copying the learned layers from the shallower model and setting additional layers to identity mapping.

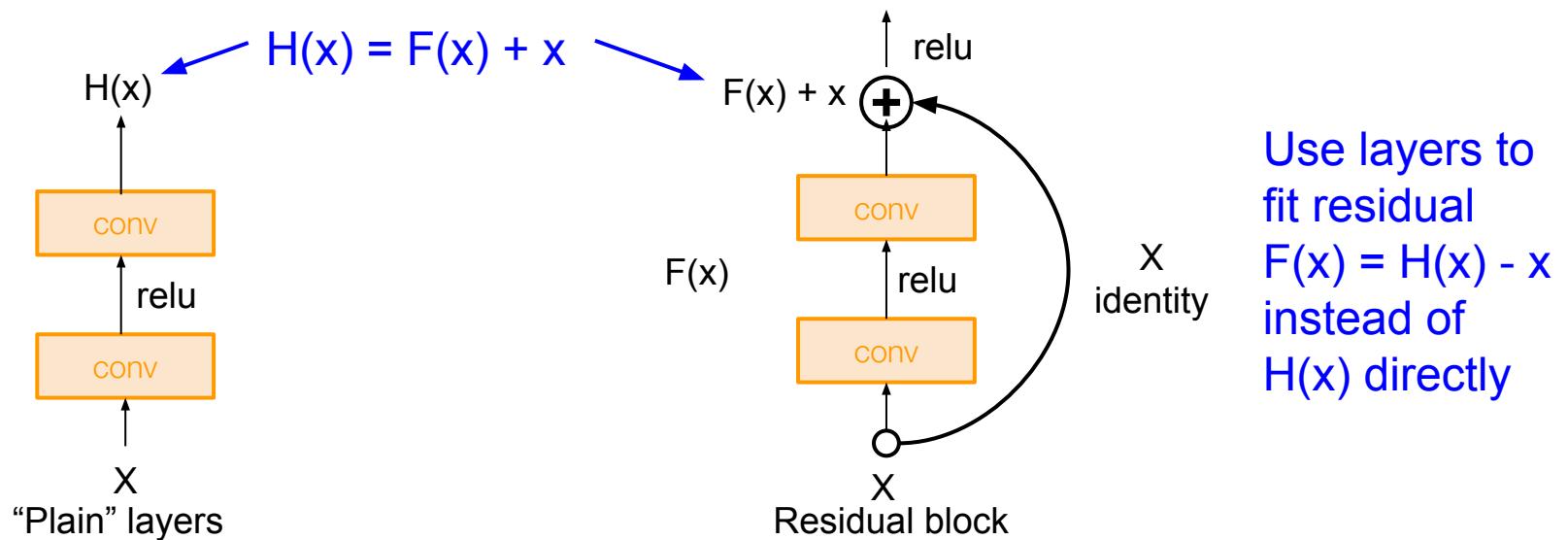
# ResNet [He et al., 2015]

Solution: Use network layers to fit a residual mapping instead of directly trying to fit a desired underlying mapping



# ResNet [He et al., 2015]

Solution: Use network layers to fit a residual mapping instead of directly trying to fit a desired underlying mapping

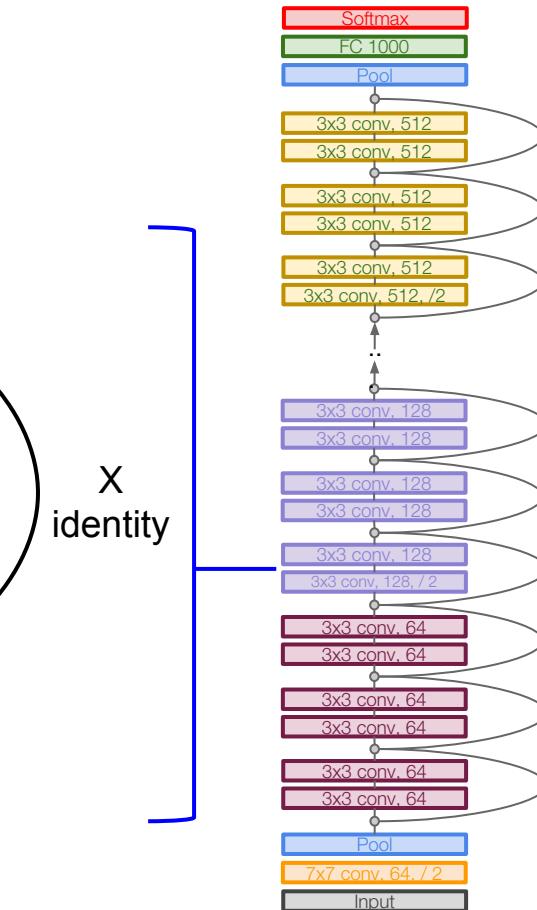
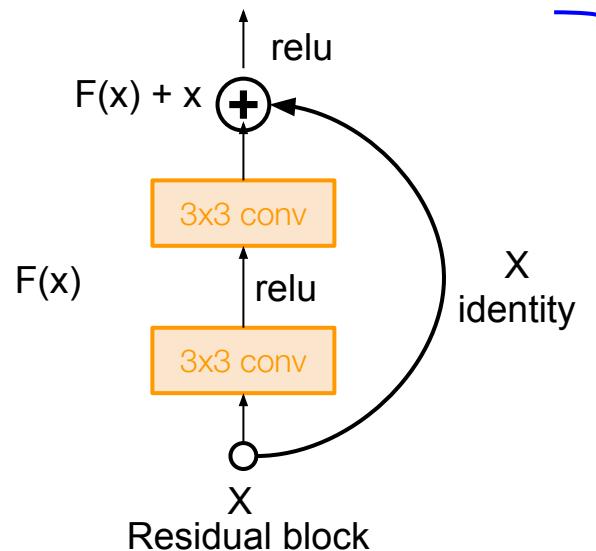


He et al. hypothesize that it is easier to optimize the residual mapping than to optimize the original, unreferenced mapping

# ResNet [He et al., 2015]

Full ResNet architecture:

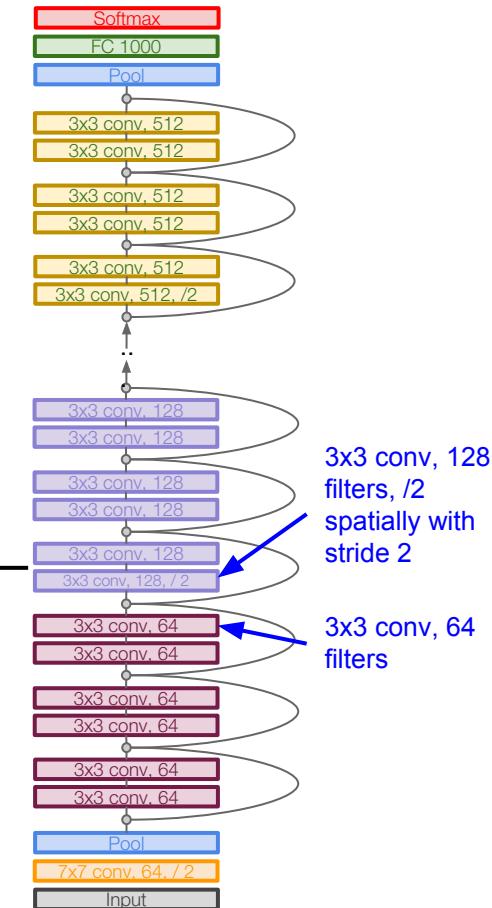
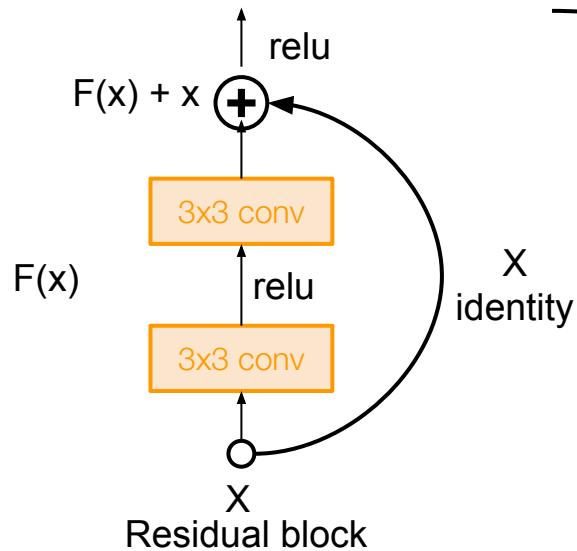
- Stack residual blocks
- Every residual block has two 3x3 conv layers



# ResNet [He et al., 2015]

Full ResNet architecture:

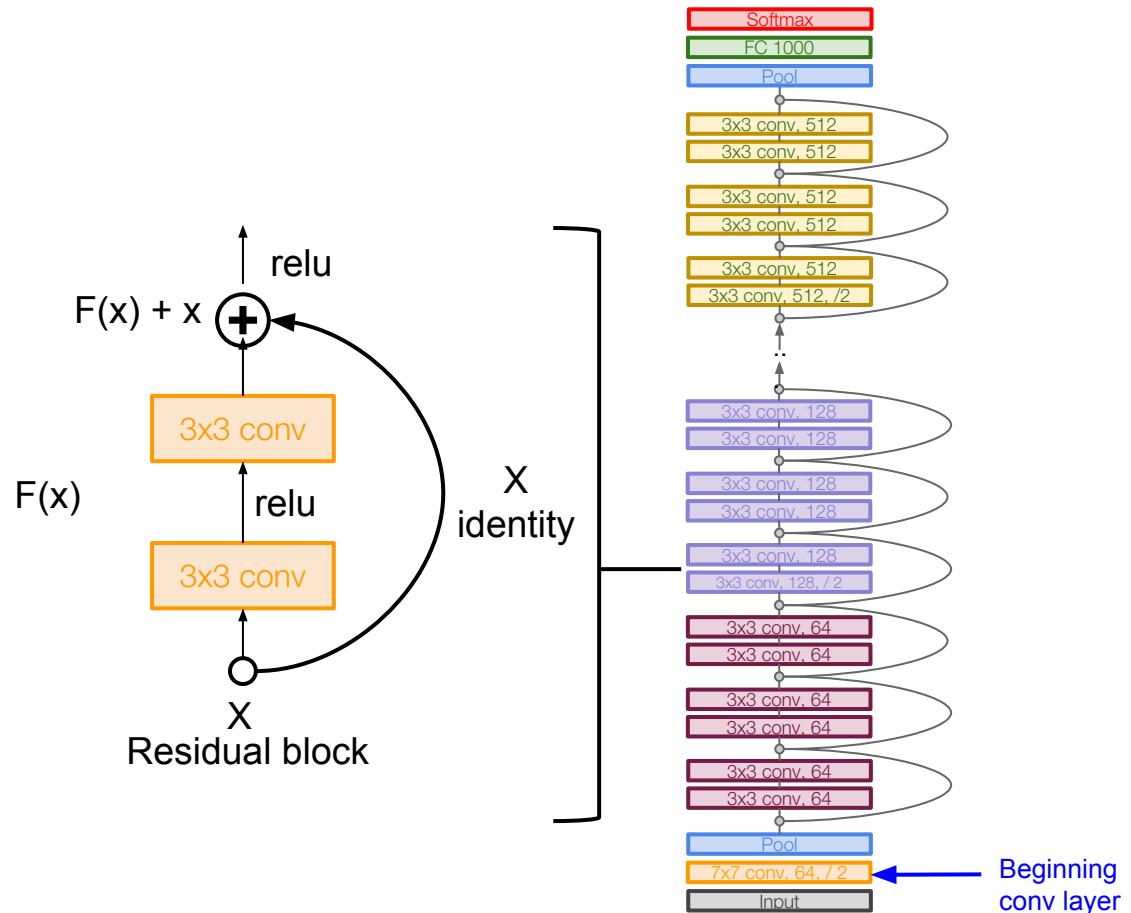
- Stack residual blocks
- Every residual block has two 3x3 conv layers
- Periodically, double # of filters and downsample spatially using stride 2 (/2 in each dimension)



# ResNet [He et al., 2015]

Full ResNet architecture:

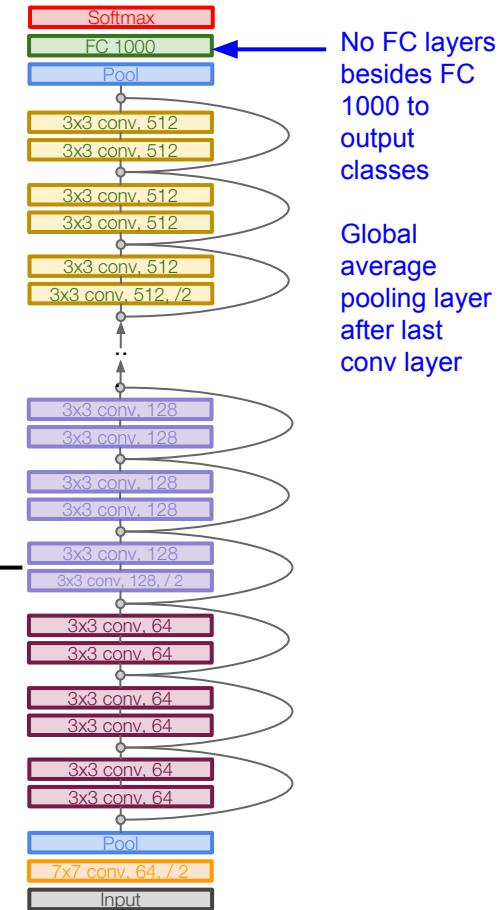
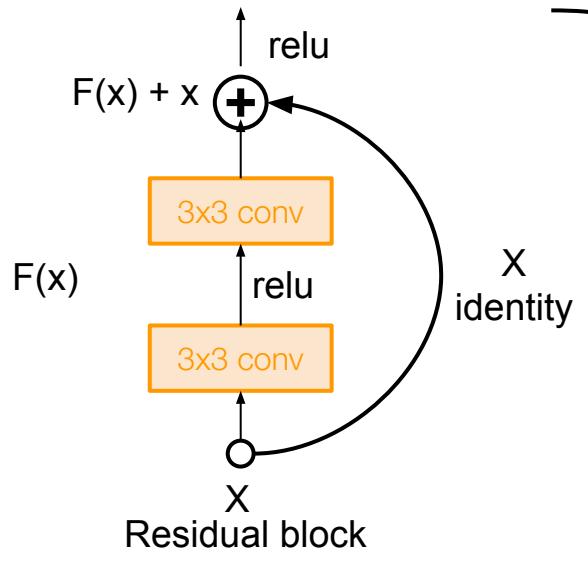
- Stack residual blocks
- Every residual block has two 3x3 conv layers
- Periodically, double # of filters and downsample spatially using stride 2 (/2 in each dimension)
- Additional conv layer at the beginning



# ResNet [He et al., 2015]

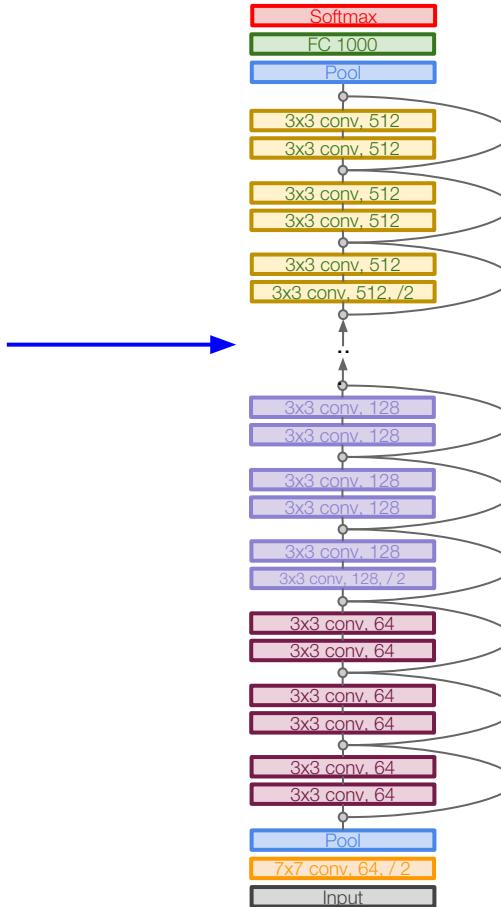
Full ResNet architecture:

- Stack residual blocks
- Every residual block has two 3x3 conv layers
- Periodically, double # of filters and downsample spatially using stride 2 (/2 in each dimension)
- Additional conv layer at the beginning
- No FC layers at the end (only FC 1000 to output classes)



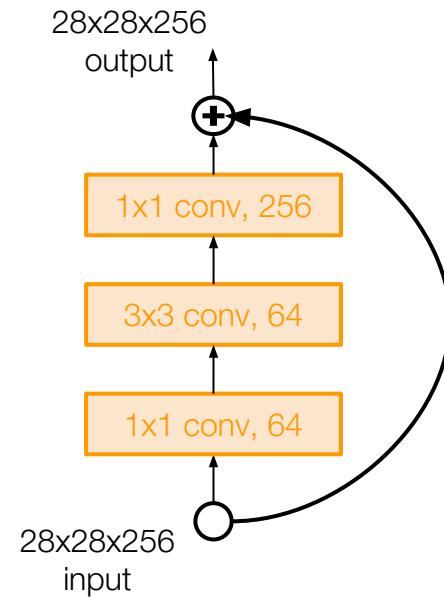
# ResNet [He et al., 2015]

Total depths of 34, 50, 101, or  
152 layers for ImageNet



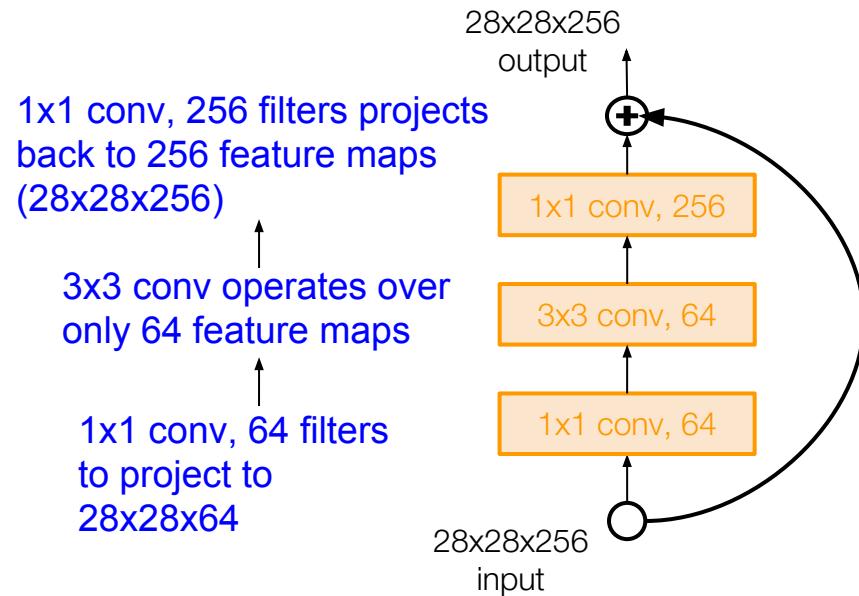
# ResNet [He et al., 2015]

For deeper networks  
(ResNet-50+), use “bottleneck”  
layer to improve efficiency  
(similar to GoogLeNet)



# ResNet [He et al., 2015]

For deeper networks  
(ResNet-50+), use “bottleneck”  
layer to improve efficiency  
(similar to GoogLeNet)



# ResNet [He et al., 2015]

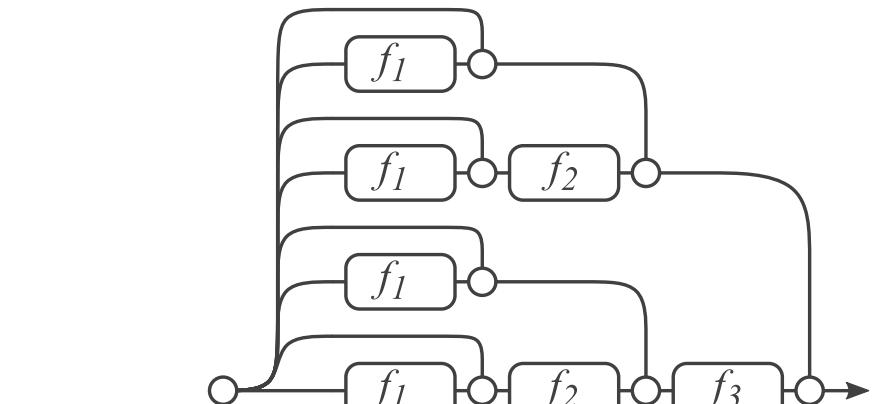
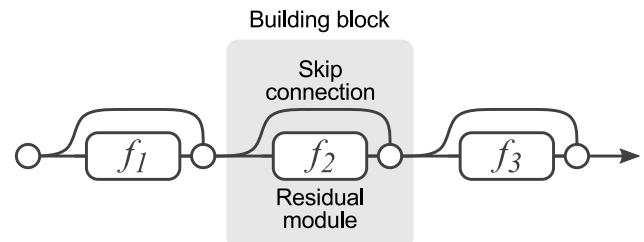
Training ResNet in practice:

- Batch Normalization after every CONV layer
- Xavier/2 initialization from He et al.
- SGD + Momentum (0.9)
- Learning rate: 0.1, divided by 10 when validation error plateaus
- Mini-batch size 256
- Weight decay of 1e-5
- No dropout used

why

# ResNet [He et al., 2015]

Residual Networks Behave Like Ensembles of Relatively Shallow Networks  
[A. Veit et al., NIPS, 2016]



(a) Conventional 3-block residual network

(b) Unraveled view of (a)

N residual blocks



$2^N$  paths

# ResNet [He et al., 2015]

Residual Networks Behave Like Ensembles of Relatively Shallow Networks  
[A. Veit et al., NIPS, 2016]

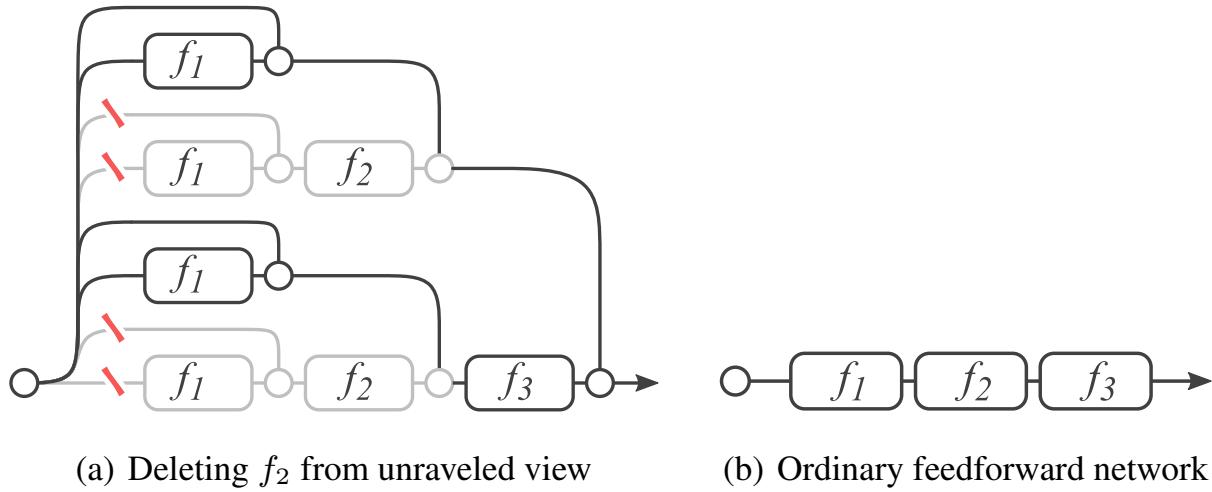


Figure 2: Deleting a layer in residual networks at test time (a) is equivalent to zeroing half of the paths. In ordinary feed-forward networks (b) such as VGG or AlexNet, deleting individual layers alters the only viable path from input to output.

# ResNet [He et al., 2015]

Residual Networks Behave Like Ensembles of Relatively Shallow Networks  
[A. Veit et al., NIPS, 2016]

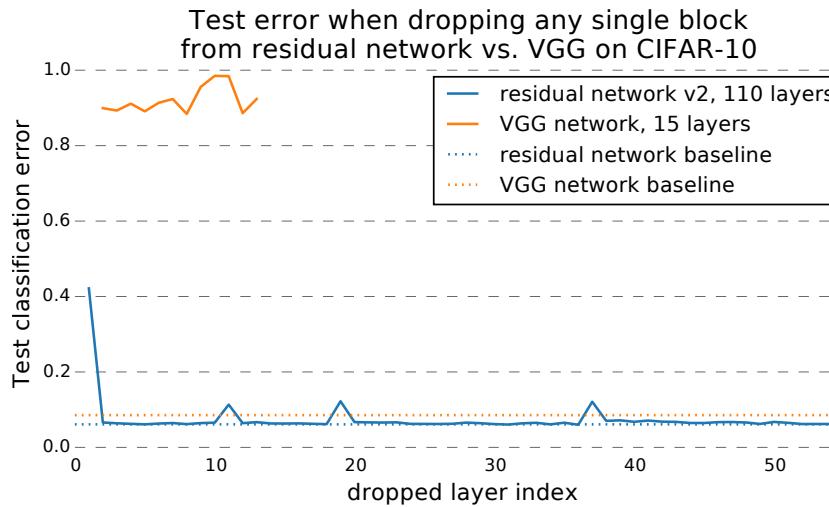


Figure 3: Deleting individual layers from VGG and a residual network on CIFAR-10. VGG performance drops to random chance when any one of its layers is deleted, but deleting individual modules from residual networks has a minimal impact on performance. Removing downsampling modules has a slightly higher impact.

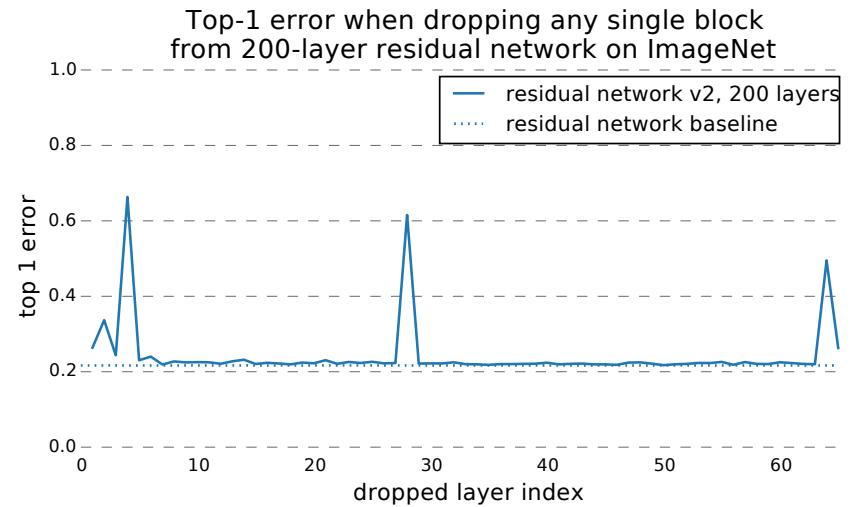


Figure 4: Results when dropping individual blocks from residual networks trained on ImageNet are similar to CIFAR results. However, downsampling layers tend to have more impact on ImageNet.

# ResNet [He et al., 2015]

Residual Networks Behave Like Ensembles of Relatively Shallow Networks  
[A. Veit et al., NIPS, 2016]

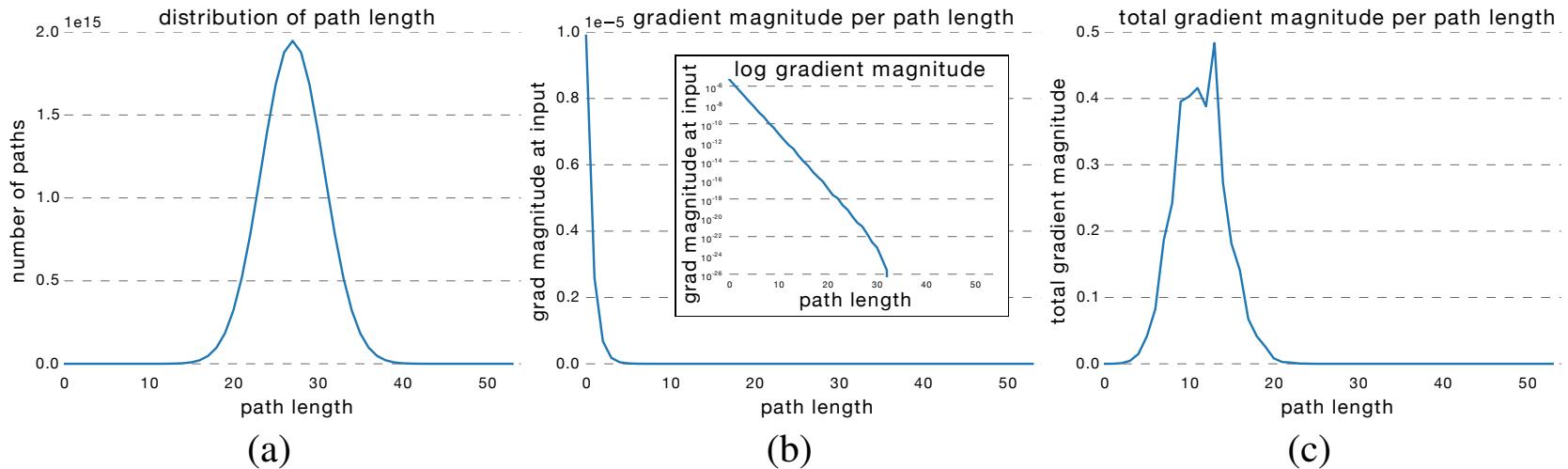


Figure 6: How much gradient do the paths of different lengths contribute in a residual network? To find out, we first show the distribution of all possible path lengths (a). This follows a Binomial distribution. Second, we record how much gradient is induced on the first layer of the network through paths of varying length (b), which appears to decay roughly exponentially with the number of modules the gradient passes through. Finally, we can multiply these two functions (c) to show how much gradient comes from all paths of a certain length. Though there are many paths of medium length, paths longer than  $\sim 20$  modules are generally too long to contribute noticeable gradient during training. This suggests that the effective paths in residual networks are relatively shallow.

# ResNet [He et al., 2015]

Residual Networks Behave Like Ensembles of Relatively Shallow Networks  
[A. Veit et al., NIPS, 2016]

He's response in “Aggregated Residual Transformations for Deep Neural Networks”

**Ensembling.** Averaging a set of independently trained networks is an effective solution to improving accuracy [24], widely adopted in recognition competitions [33]. Veit *et al.* [40] interpret a single ResNet as an ensemble of shallower networks, which results from ResNet’s *additive* behaviors [15]. Our method harnesses additions to aggregate a set of transformations. But we argue that it is imprecise to view our method as ensembling, because the members to be aggregated are trained jointly, not independently.

# ResNet [He et al., 2015]

## Experimental Results

- Able to train very deep networks without degrading (152 layers on ImageNet, 1202 on Cifar)
- Deeper networks now achieve lower training error as expected
- Swept 1st place in all ILSVRC and COCO 2015 competitions

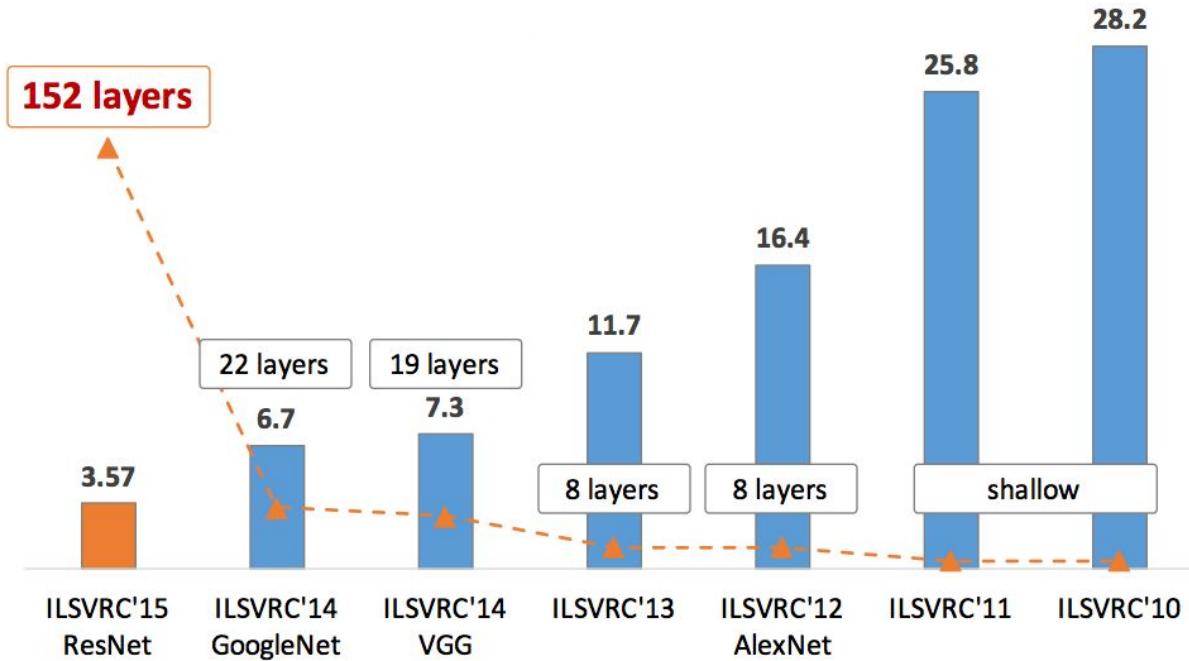
### MSRA @ ILSVRC & COCO 2015 Competitions

- **1st places** in all five main tracks
  - ImageNet Classification: “Ultra-deep” (quote Yann) **152-layer nets**
  - ImageNet Detection: **16%** better than 2nd
  - ImageNet Localization: **27%** better than 2nd
  - COCO Detection: **11%** better than 2nd
  - COCO Segmentation: **12%** better than 2nd

ILSVRC 2015 classification winner (3.6% top 5 error) -- better than “human performance”! (Russakovsky 2014)

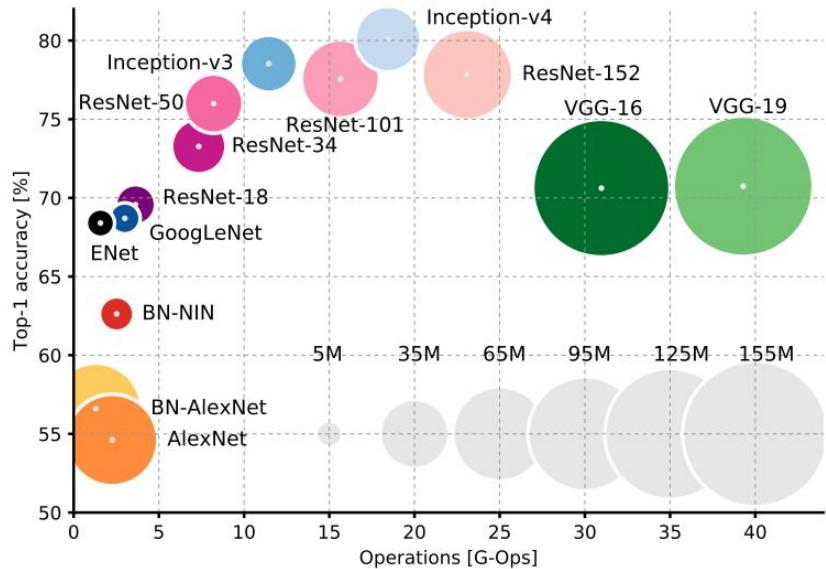
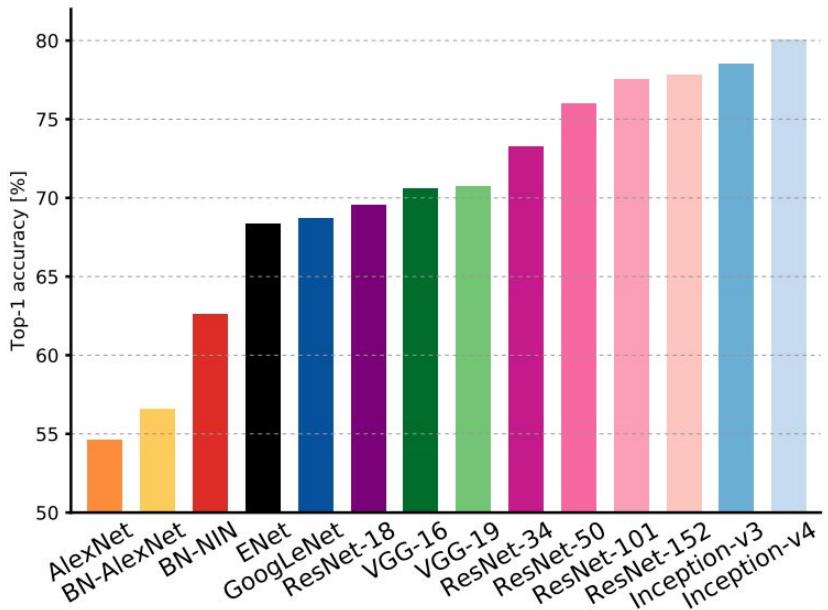
# ResNet [He et al., 2015]

ImageNet Large Scale Visual Recognition Challenge (ILSVRC) winners



# Comparison

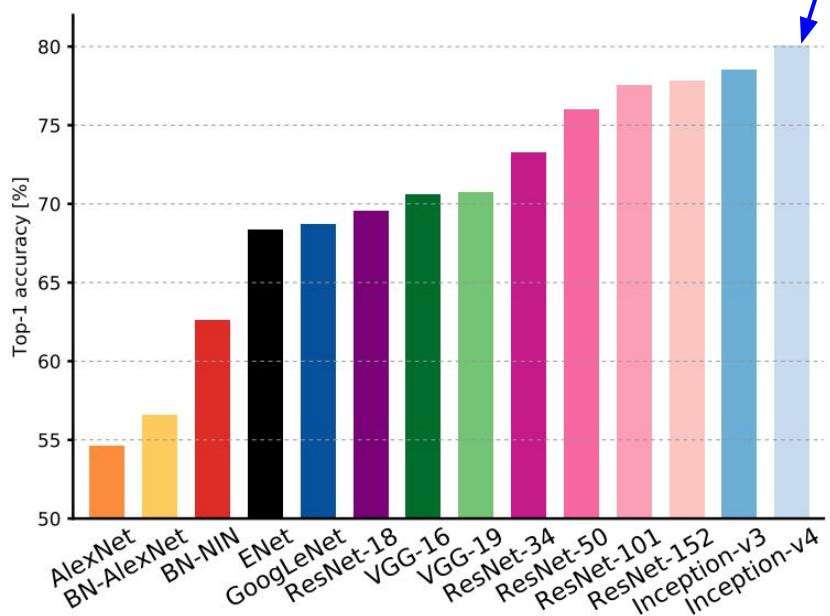
Comparing complexity...



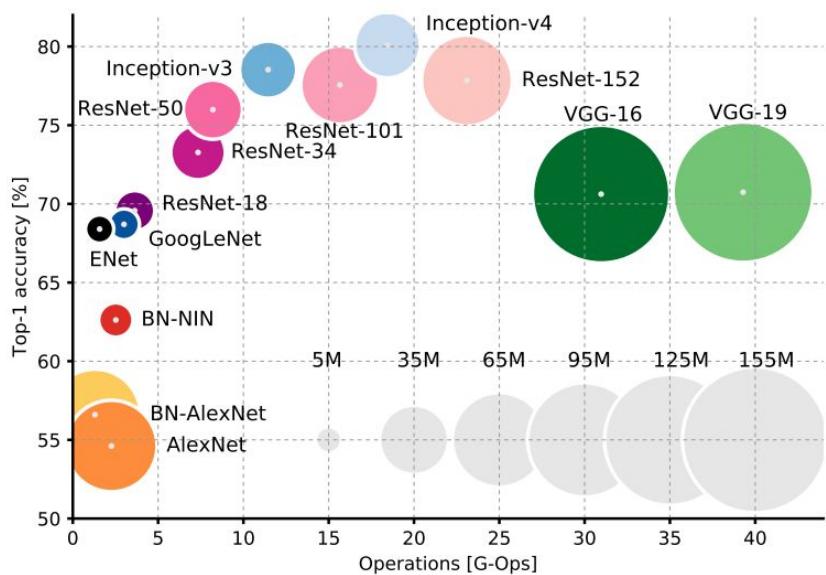
An Analysis of Deep Neural Network Models for Practical Applications, 2017.

# Comparison

Comparing complexity...



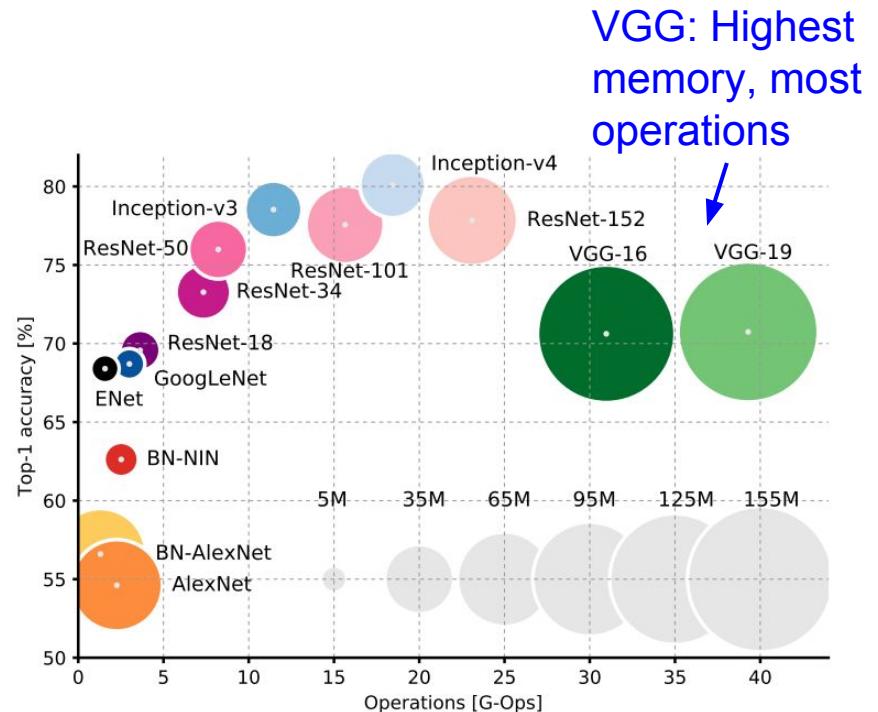
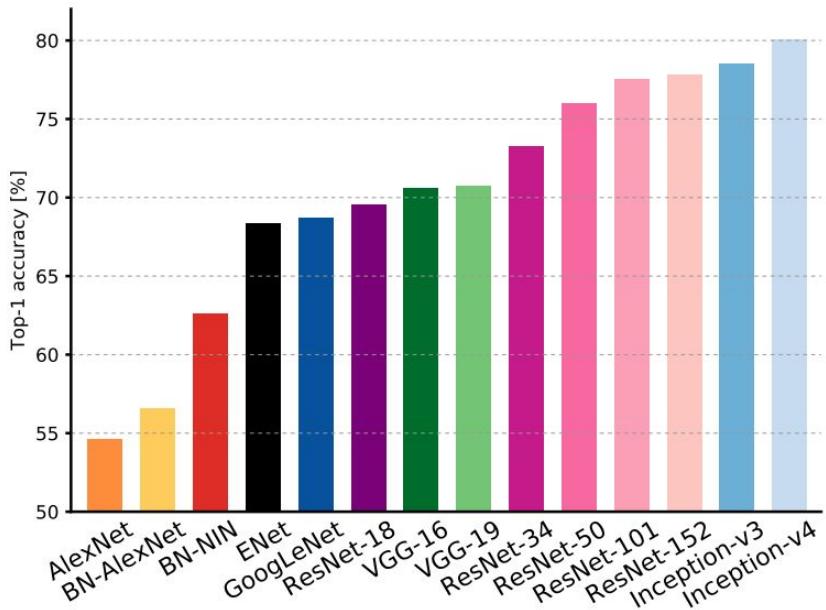
Inception-v4: Resnet + Inception!



An Analysis of Deep Neural Network Models for Practical Applications, 2017.

# Comparison

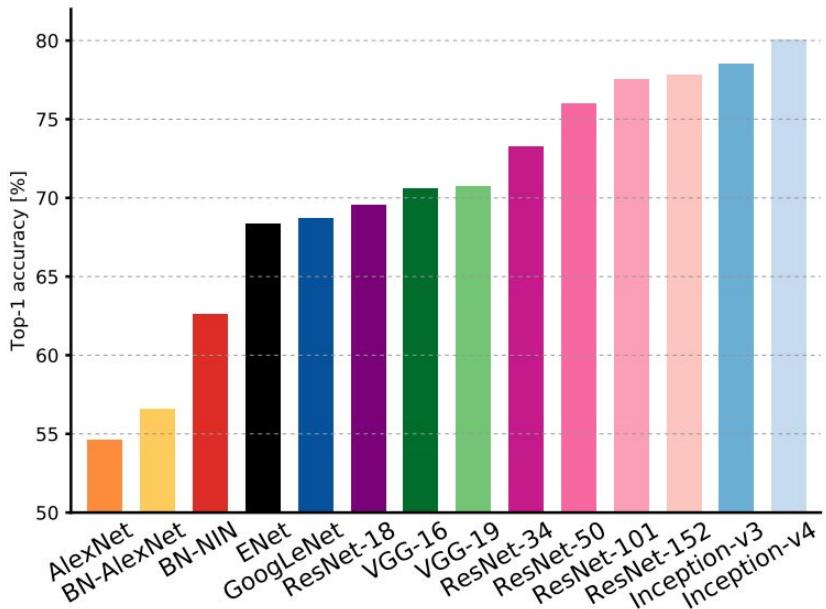
Comparing complexity...



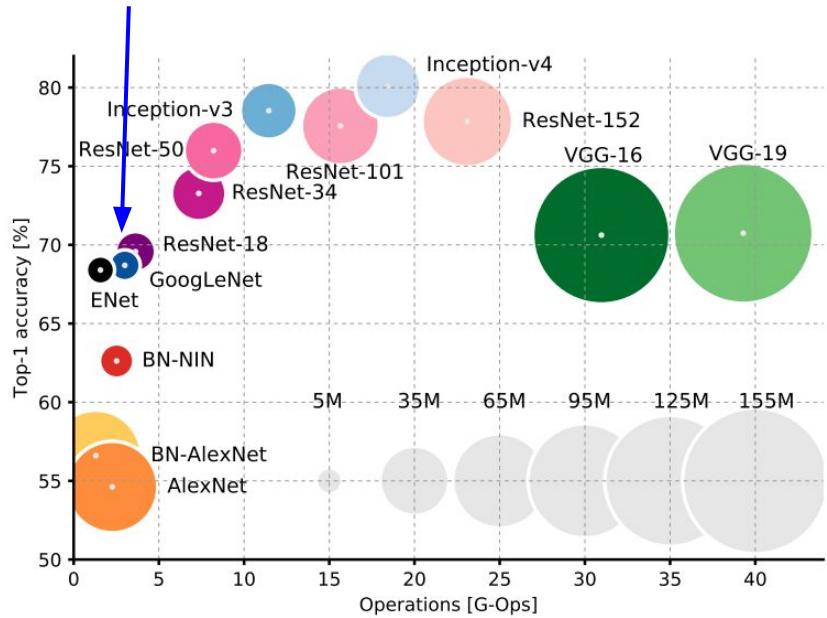
An Analysis of Deep Neural Network Models for Practical Applications, 2017.

# Comparison

Comparing complexity...



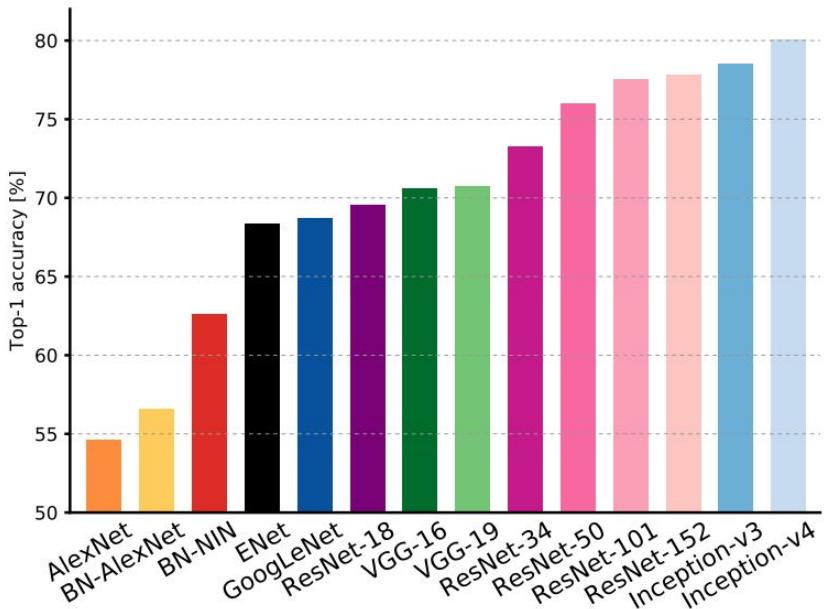
GoogLeNet:  
most efficient



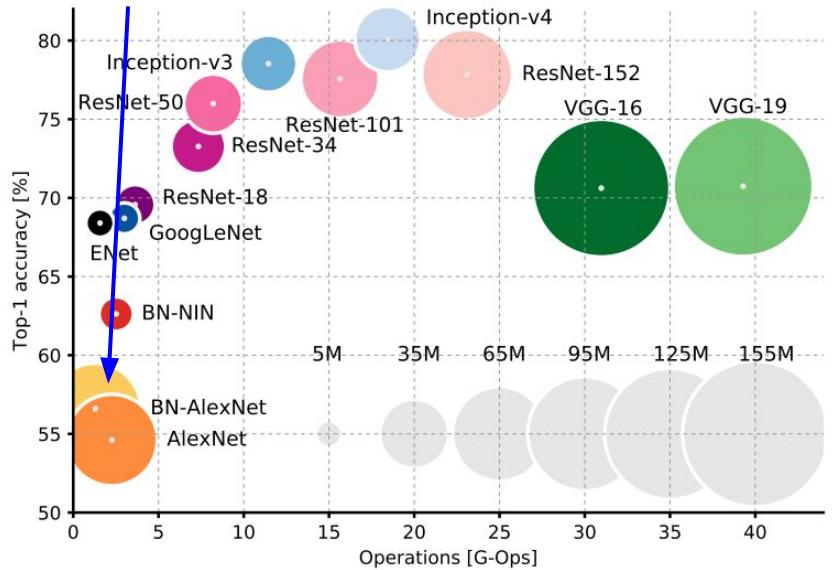
An Analysis of Deep Neural Network Models for Practical Applications, 2017.

# Comparison

## Comparing complexity...



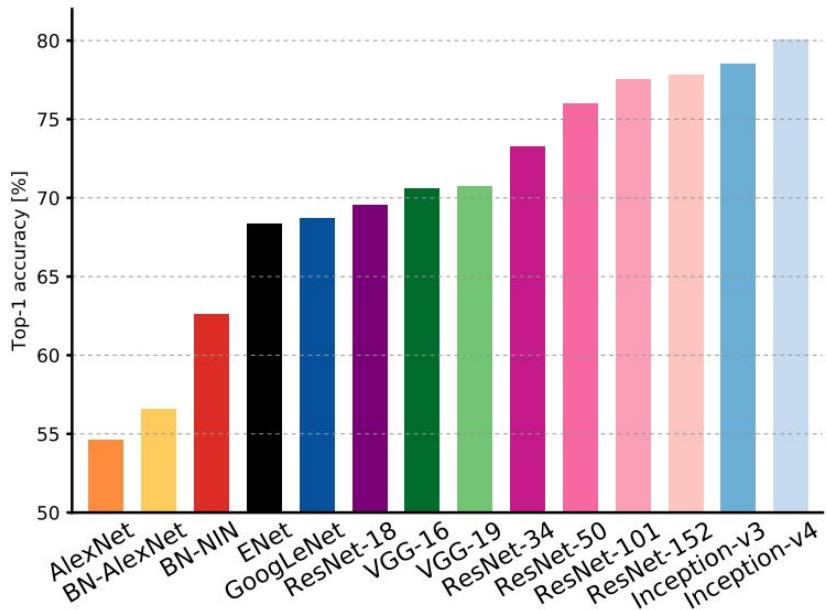
AlexNet:  
Smaller compute, still memory heavy, lower accuracy



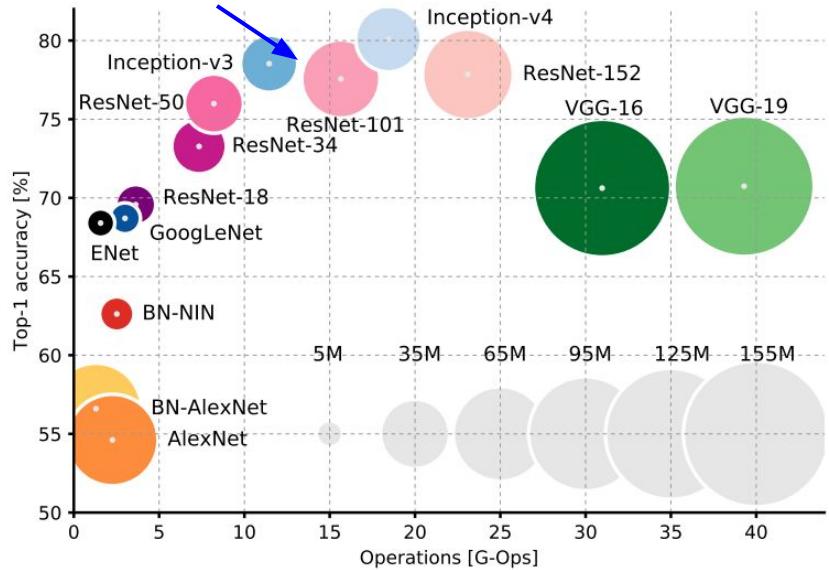
An Analysis of Deep Neural Network Models for Practical Applications, 2017.

# Comparison

## Comparing complexity...



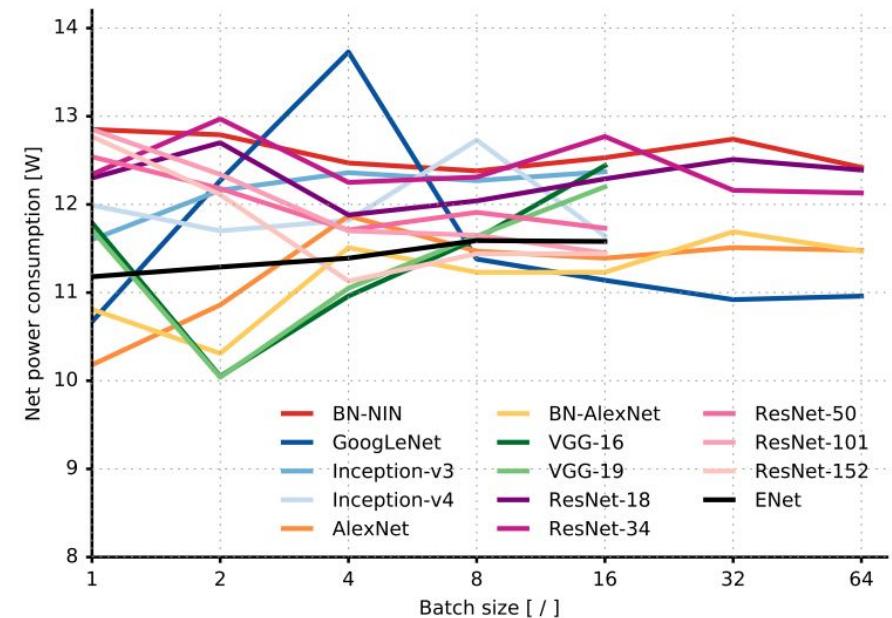
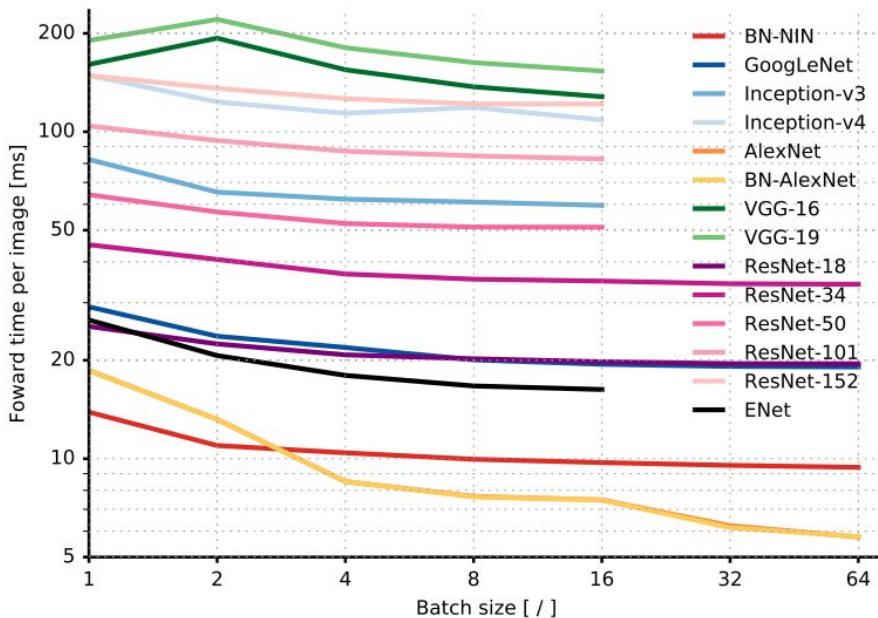
ResNet:  
Moderate efficiency depending on  
model, highest accuracy



An Analysis of Deep Neural Network Models for Practical Applications, 2017.

# Comparison

## Forward pass time and power consumption



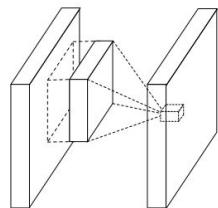
An Analysis of Deep Neural Network Models for Practical Applications, 2017.

# Other Architectures to know

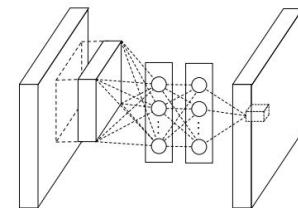
## Network in Network (NiN)

[Lin et al. 2014]

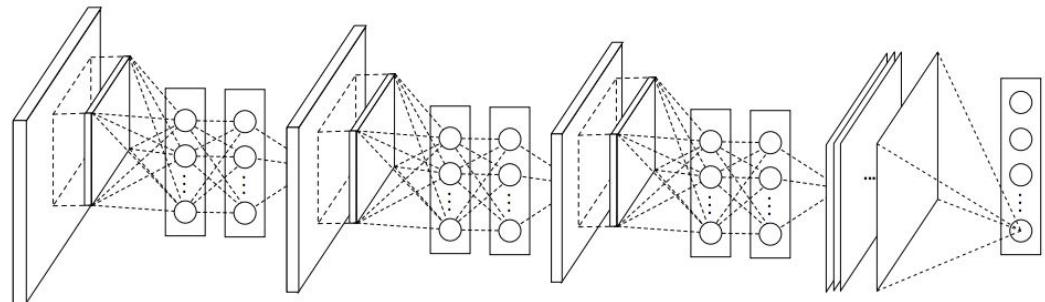
- Mlpconv layer with “micronetwork” within each conv layer to compute more abstract features for local patches
- Micronetwork uses multilayer perceptron (FC, i.e. 1x1 conv layers)
- Precursor to GoogLeNet and ResNet “bottleneck” layers
- Philosophical inspiration for GoogLeNet



(a) Linear convolution layer



(b) Mlpconv layer



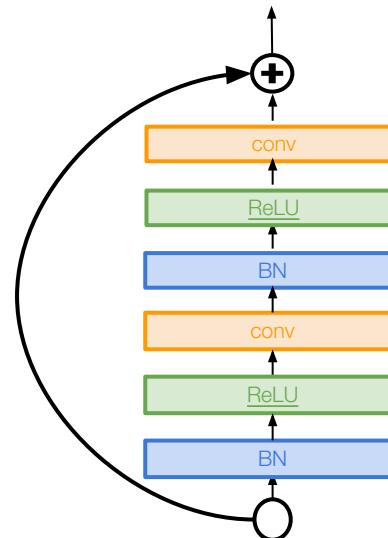
# Other Architectures to know

Improving ResNets...

## Identity Mappings in Deep Residual Networks

[He et al. 2016]

- Improved ResNet block design from creators of ResNet
- Creates a more direct path for propagating information throughout network (moves activation to residual mapping pathway)
- Gives better performance



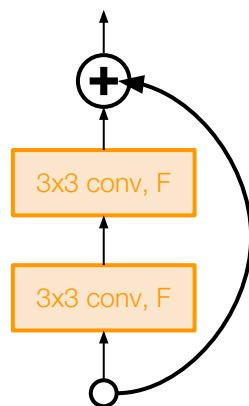
# Other Architectures to know

Improving ResNets...

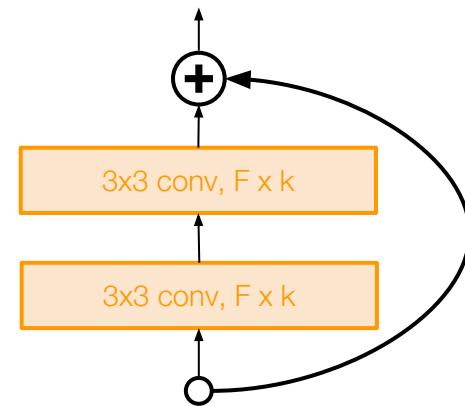
## Wide Residual Networks

[Zagoruyko et al. 2016]

- Argues that residuals are the important factor, not depth
- Use wider residual blocks ( $F \times k$  filters instead of  $F$  filters in each layer)
- 50-layer wide ResNet outperforms 152-layer original ResNet
- Increasing width instead of depth more computationally efficient (parallelizable)



Basic residual block



Wide residual block

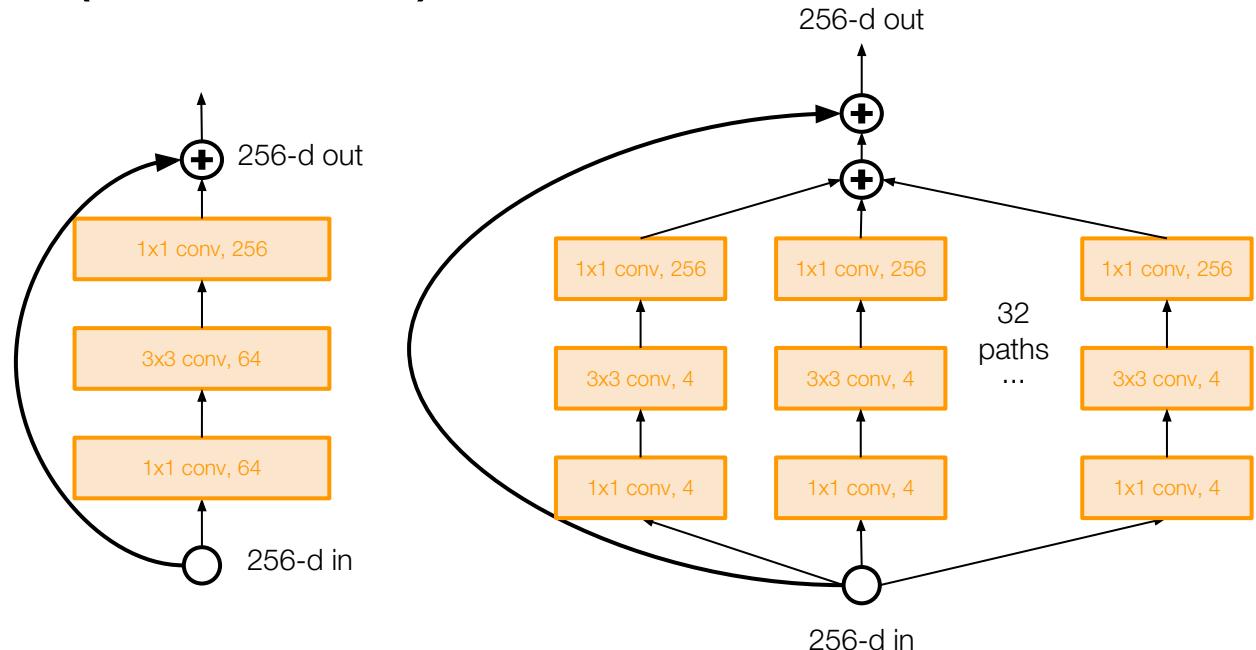
# Other Architectures to know

Improving ResNets...

## Aggregated Residual Transformations for Deep Neural Networks (ResNeXt)

[Xie et al. 2016]

- Also from creators of ResNet
- Increases width of residual block through multiple parallel pathways (“cardinality”)
- Parallel pathways similar in spirit to Inception module



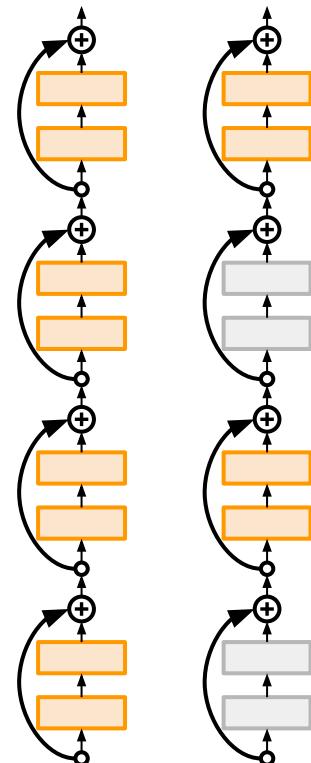
# Other Architectures to know

Improving ResNets...

## Deep Networks with Stochastic Depth

[Huang et al. 2016]

- Motivation: reduce vanishing gradients and training time through short networks during training
- Randomly drop a subset of layers during each training pass
- Bypass with identity function
- Use full deep network at test time



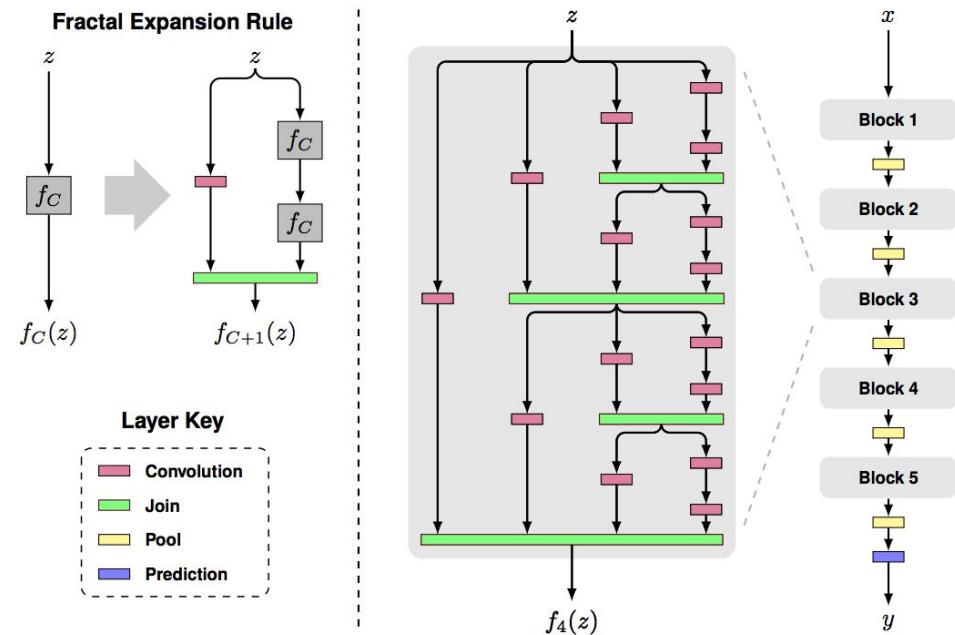
# Other Architectures to know

Beyond ResNets...

## FractalNet: Ultra-Deep Neural Networks without Residuals

[Larsson et al. 2017]

- Argues that key is transitioning effectively from shallow to deep and residual representations are not necessary
- Fractal architecture with both shallow and deep paths to output
- Trained with dropping out sub-paths
- Full network at test time



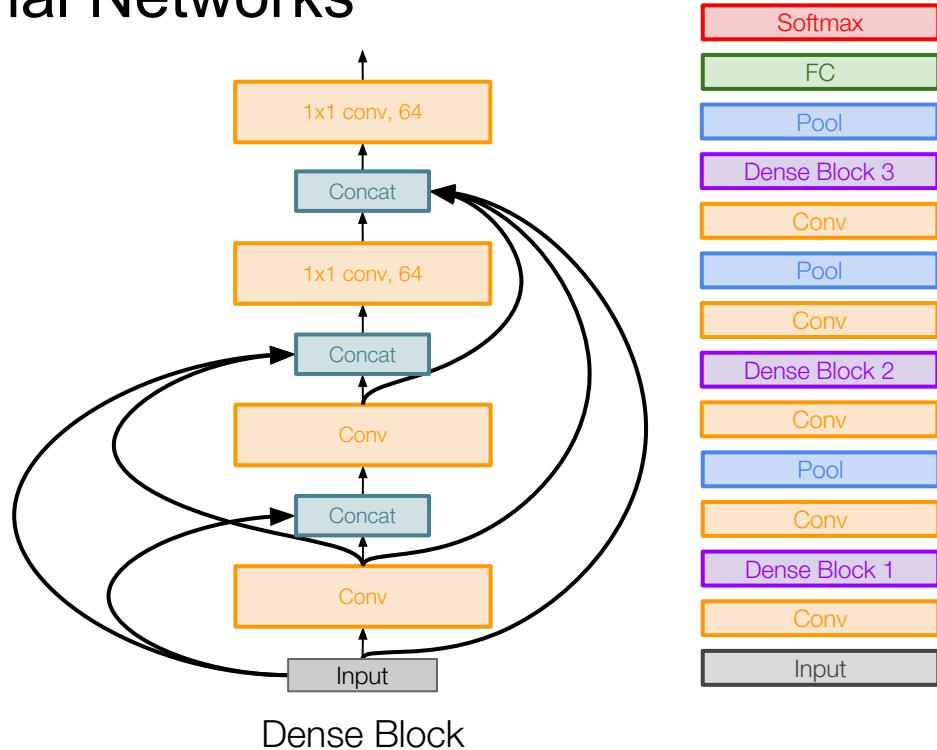
# Other Architectures to know

Beyond ResNets...

## Densely Connected Convolutional Networks

[Huang et al. 2017]

- Dense blocks where each layer is connected to every other layer in feedforward fashion
- Alleviates vanishing gradient, strengthens feature propagation, encourages feature reuse



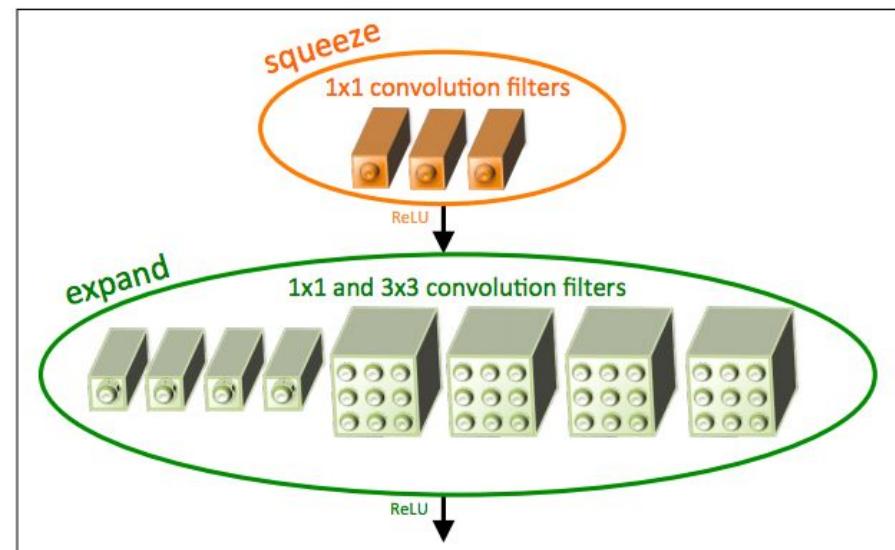
# Other Architectures to know

Efficient networks...

SqueezeNet: AlexNet-level Accuracy With 50x Fewer Parameters and <0.5Mb Model Size

[Iandola et al. 2017]

- Fire modules consisting of a 'squeeze' layer with 1x1 filters feeding an 'expand' layer with 1x1 and 3x3 filters
- AlexNet level accuracy on ImageNet with 50x fewer parameters
- Can compress to 510x smaller than AlexNet (0.5Mb)



# Summary

## Case Studies

- AlexNet
- VGG
- GoogLeNet
- ResNet

## Also....

- NiN (Network in Network)
- Wide ResNet
- ResNeXT
- Stochastic Depth
- DenseNet
- FractalNet
- SqueezeNet

# Summary

- VGG, GoogLeNet, ResNet all in wide use, available in model zoos
- ResNet current best default
- Trend towards extremely deep networks
- Significant research centers around design of layer / skip connections and improving gradient flow
- Even more recent trend towards examining necessity of depth vs. width and residual connections

# Deep networks, why and how

# Why deep networks

Note that a shallow network with a single, huge, hidden layer is already a universal approximator

# Why deep networks

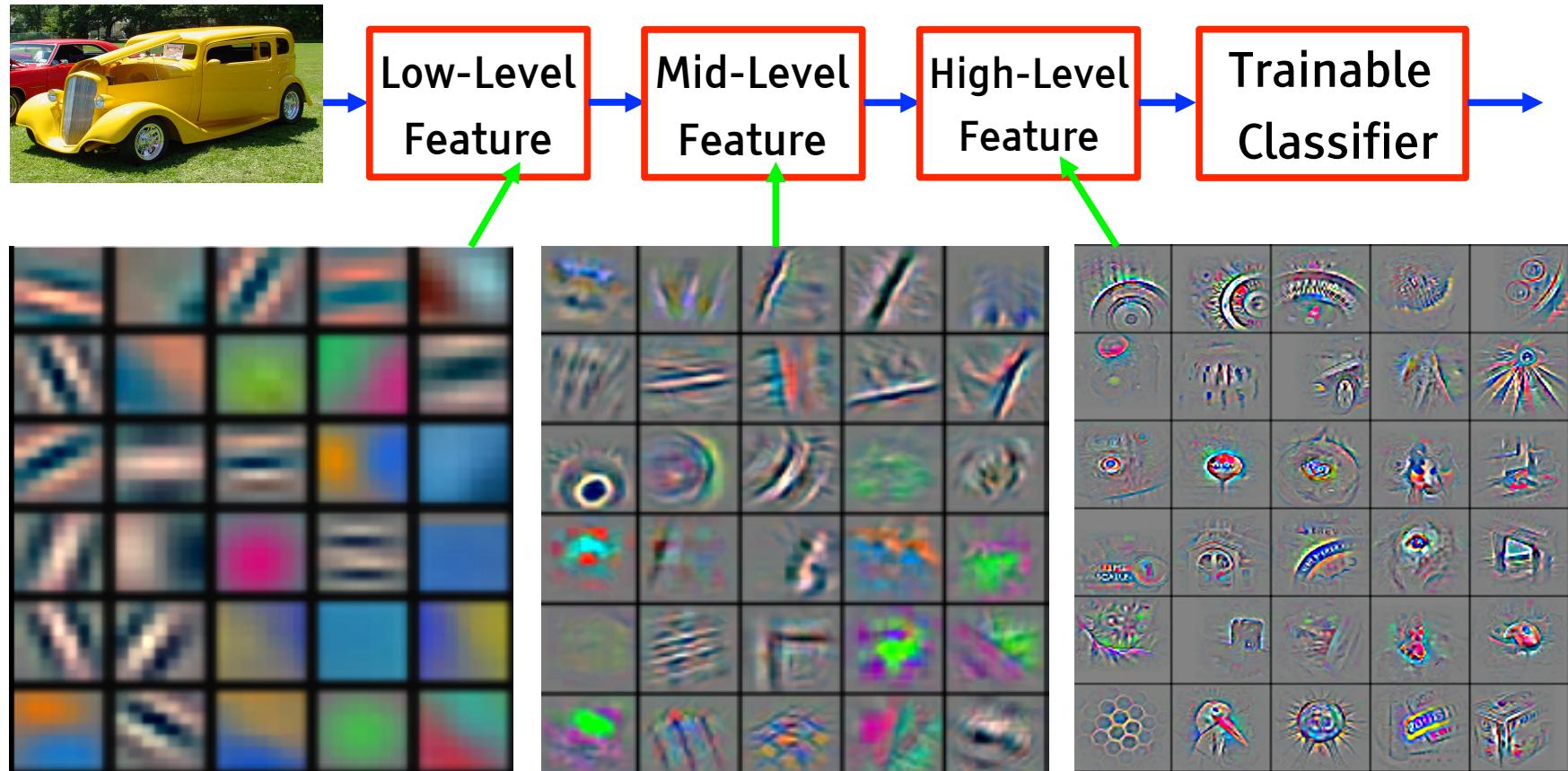
Note that a shallow network with a single, huge, hidden layer is already a universal approximator

The advantages of deeper networks:

1. Additional prior: compositionality is useful to describe the world around us efficiently
2. Exponential expressiveness of deep networks

# The advantage of deep networks

1. Additional prior: compositionality is useful to describe the world around us efficiently



# The advantage of deep networks

## 2. Exponential expressiveness of deep networks

Some functions compactly represented with  $k$  layers may require exponentially many more parameters with 2 layers

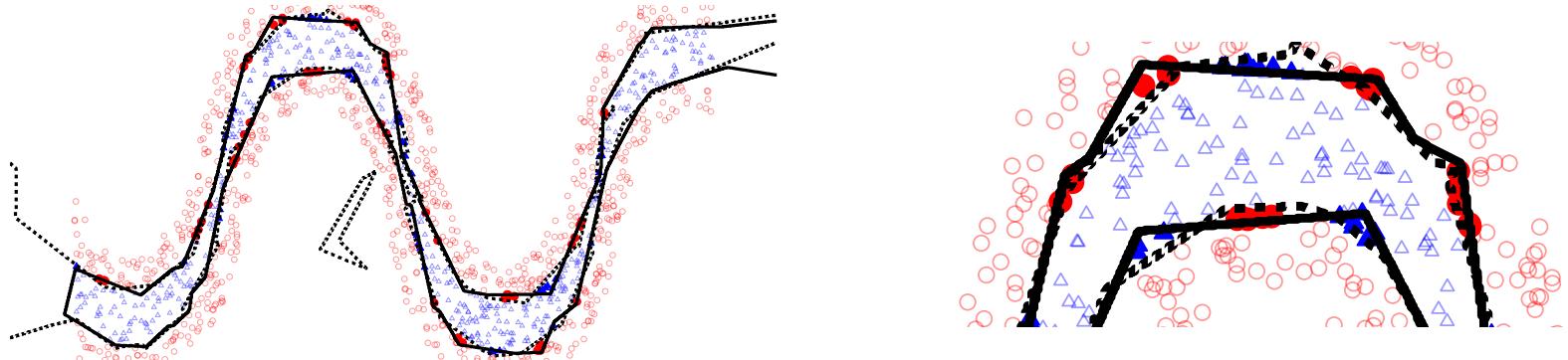


Figure 1: Binary classification using a shallow model with 20 hidden units (solid line) and a deep model with two layers of 10 units each (dashed line). The right panel shows a close-up of the left panel. Filled markers indicate errors made by the shallow model.

# But deep networks are hard to train

## Problem1: Vanishing/exploding gradients

- **Phenomenon:** the gradient of early layers become extremely small or large
- **When:** hamper convergence from the beginning
- **Reason:** multiplication of several parameters
- **How to avoid:**
  - Recurrent networks use the same parameter at each time step, hence gives rise to especially pronounced difficulties. Let's discuss next time.
  - Feed-forward networks do not use the same parameters at each layer, and this problem has been largely addressed by normalized initialization and intermediate normalization layers (e.g., batch normalization).

# But deep networks are hard to train

## Problem2: Degradation problem

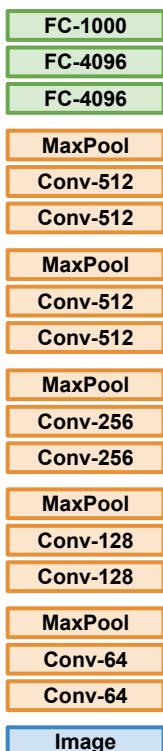
- **Phenomenon:** with the network depth increasing, accuracy gets saturated and then degrades rapidly.
- **When:** after deeper networks are able to start converging
- **Reason:** long-range dependency
- **How to avoid:**
  - adding extra copies of the output that are attached to the intermediate hidden layers of the network, as in GoogLeNet
  - adding skip connections between layers reduce the length of the shortest path from the lower layer's parameters to the output, as in ResNet

# Transfer Learning with CNNs

# Transfer learning with CNNs

## Transfer Learning with CNNs

### 1. Train on Imagenet

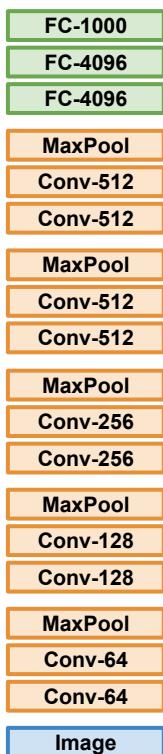


Donahue et al, "DeCAF: A Deep Convolutional Activation Feature for Generic Visual Recognition", ICML 2014  
Razavian et al, "CNN Features Off-the-Shelf: An Astounding Baseline for Recognition", CVPR Workshops 2014

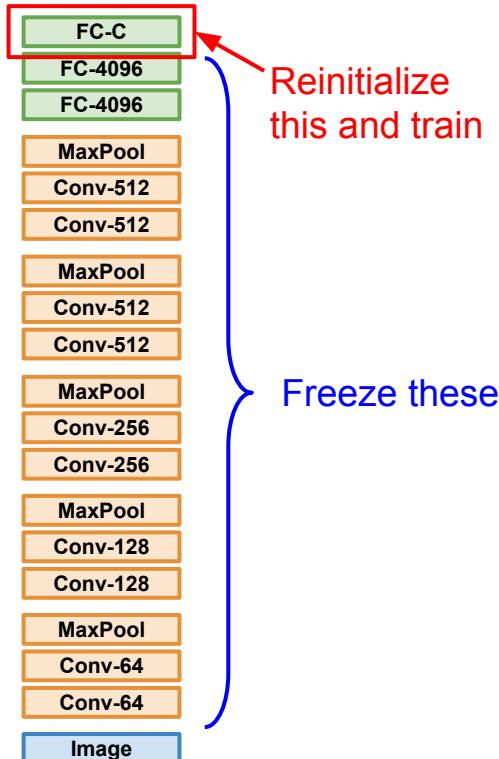
# Transfer learning with CNNs

## Transfer Learning with CNNs

1. Train on Imagenet



2. Small Dataset (C classes)

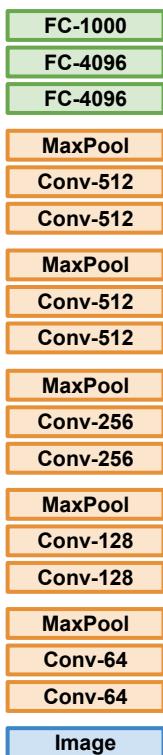


Donahue et al, "DeCAF: A Deep Convolutional Activation Feature for Generic Visual Recognition", ICML 2014  
Razavian et al, "CNN Features Off-the-Shelf: An Astounding Baseline for Recognition", CVPR Workshops 2014

# Transfer learning with CNNs

## Transfer Learning with CNNs

### 1. Train on Imagenet



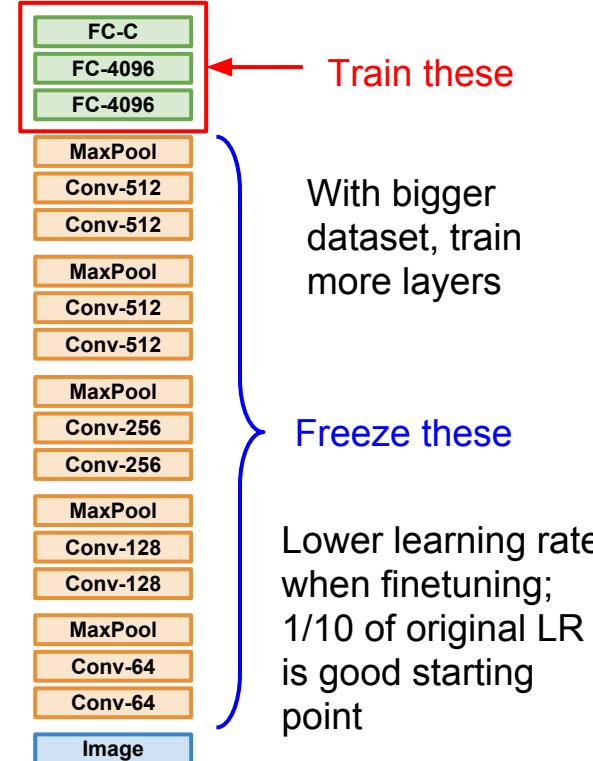
### 2. Small Dataset (C classes)



Reinitialize  
this and train

Freeze these

### 3. Bigger dataset



Train these

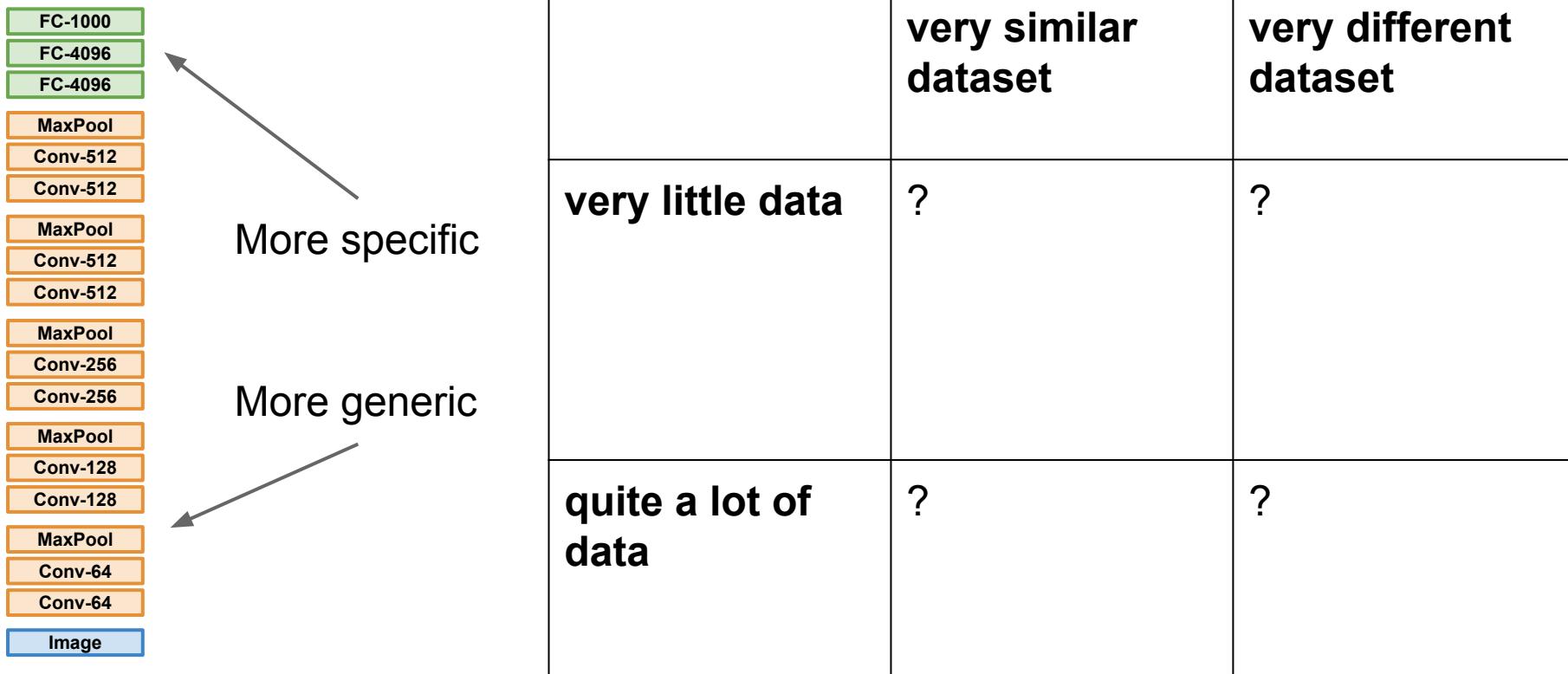
With bigger  
dataset, train  
more layers

Freeze these

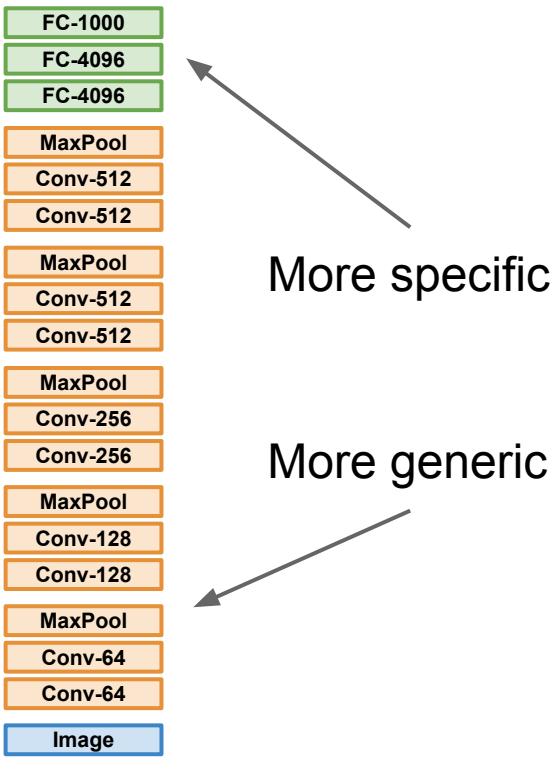
Lower learning rate  
when finetuning;  
1/10 of original LR  
is good starting  
point

Donahue et al, "DeCAF: A Deep Convolutional Activation Feature for Generic Visual Recognition", ICML 2014  
Razavian et al, "CNN Features Off-the-Shelf: An Astounding Baseline for Recognition", CVPR Workshops 2014

# Transfer learning with CNNs



# Transfer learning with CNNs

	<b>very similar dataset</b>	<b>very different dataset</b>
 <p>More specific</p> <p>More generic</p> <pre>graph TD; Image[Image] --&gt; Conv64_1[Conv-64]; Conv64_1 --&gt; MaxPool1[MaxPool]; MaxPool1 --&gt; Conv64_2[Conv-64]; Conv64_2 --&gt; MaxPool2[MaxPool]; MaxPool2 --&gt; Conv256[Conv-256]; Conv256 --&gt; MaxPool3[MaxPool]; MaxPool3 --&gt; Conv128[Conv-128]; Conv128 --&gt; MaxPool4[MaxPool]; MaxPool4 --&gt; Conv128_2[Conv-128]; Conv128_2 --&gt; MaxPool5[MaxPool]; MaxPool5 --&gt; Conv512[Conv-512]; Conv512 --&gt; MaxPool6[MaxPool]; MaxPool6 --&gt; Conv512_2[Conv-512]; Conv512_2 --&gt; MaxPool7[MaxPool]; MaxPool7 --&gt; Conv512_3[Conv-512]; Conv512_3 --&gt; MaxPool8[MaxPool]; MaxPool8 --&gt; Conv4096[Conv-4096]; Conv4096 --&gt; MaxPool9[MaxPool]; MaxPool9 --&gt; FC4096[FC-4096]; FC4096 --&gt; FC1000[FC-1000]</pre>	<b>very little data</b> Use Linear Classifier on top layer	?
	<b>quite a lot of data</b>	Finetune a few layers

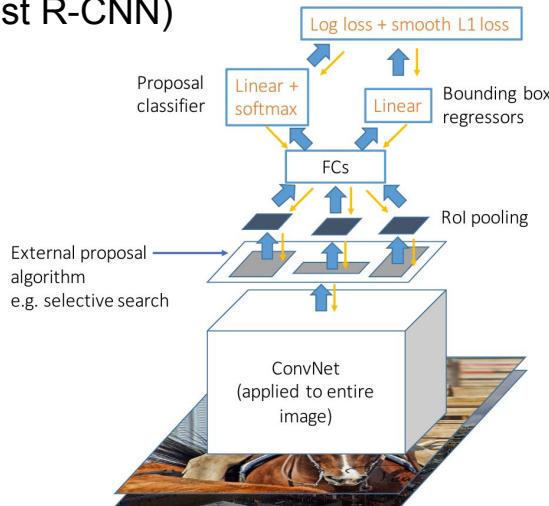
# Transfer learning with CNNs

	<b>very similar dataset</b>	<b>very different dataset</b>
<p>The diagram shows a vertical stack of layers. From bottom to top: a blue box labeled "Image", followed by four orange boxes labeled "Conv-64", "MaxPool", "Conv-64", and "MaxPool". Then there is a break in the sequence, indicated by a gap. Above this gap are five more layers: four orange boxes labeled "Conv-128", "MaxPool", "Conv-128", and "MaxPool", followed by a single green box labeled "FC-4096". At the very top is another green box labeled "FC-1000". A diagonal arrow points from the "Image" layer towards the "FC-1000" layer, with the text "More generic" written below it. Another diagonal arrow points from the "FC-4096" layer towards the "FC-1000" layer, with the text "More specific" written above it.</p> <pre>graph TD; Image[Image] --&gt; C1[Conv-64]; C1 --&gt; MP1[MaxPool]; MP1 --&gt; C2[Conv-64]; C2 --&gt; MP2[MaxPool]; MP2 --&gt; C3[Conv-128]; C3 --&gt; MP3[MaxPool]; MP3 --&gt; C4[Conv-128]; C4 --&gt; MP4[MaxPool]; MP4 --&gt; FC1[FC-4096]; FC1 --&gt; FC1000[FC-1000]</pre>	<b>very little data</b> Use Linear Classifier on top layer	You're in trouble... Try linear classifier from different stages
	<b>quite a lot of data</b>	Finetune a few layers Finetune a larger number of layers

# Transfer learning with CNNs

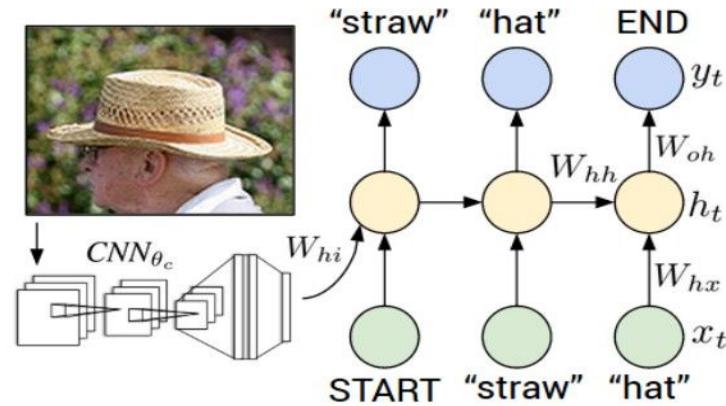
Transfer learning with CNNs is pervasive...  
(it's the norm, not an exception)

Object Detection  
(Fast R-CNN)



Girshick, "Fast R-CNN", ICCV 2015  
Figure copyright Ross Girshick, 2015. Reproduced with permission.

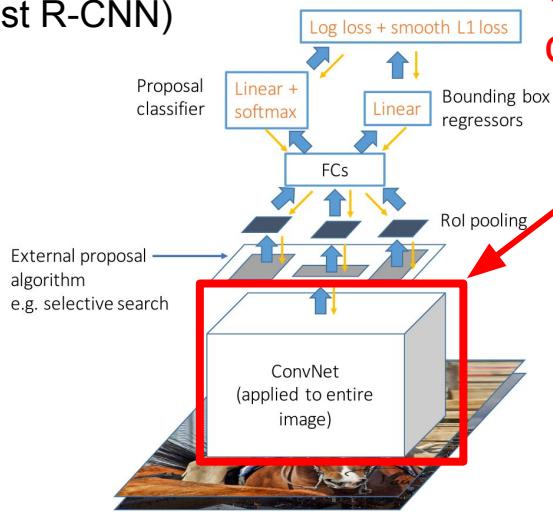
Image Captioning: CNN + RNN



Karpathy and Fei-Fei, "Deep Visual-Semantic Alignments for Generating Image Descriptions", CVPR 2015  
Figure copyright IEEE, 2015. Reproduced for educational purposes.

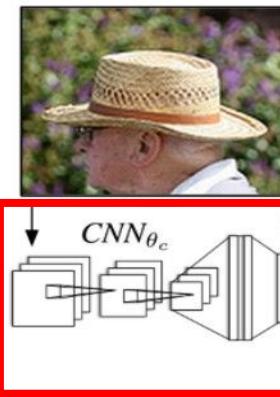
# Transfer learning is pervasive

Object Detection  
(Fast R-CNN)



CNN pretrained  
on ImageNet

Image Captioning: CNN + RNN

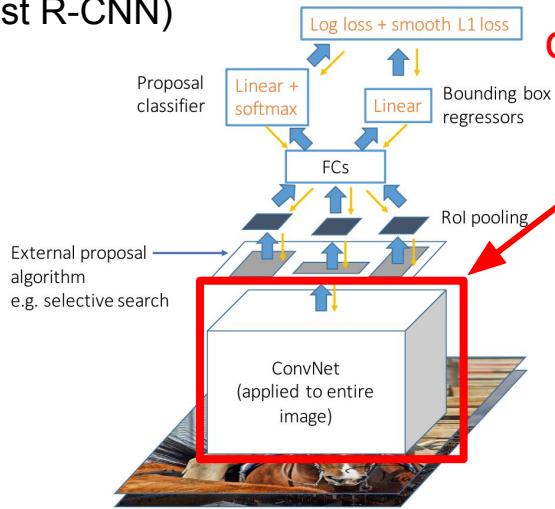


Girshick, "Fast R-CNN", ICCV 2015  
Figure copyright Ross Girshick, 2015. Reproduced with permission.

Karpathy and Fei-Fei, "Deep Visual-Semantic Alignments for Generating Image Descriptions", CVPR 2015  
Figure copyright IEEE, 2015. Reproduced for educational purposes.

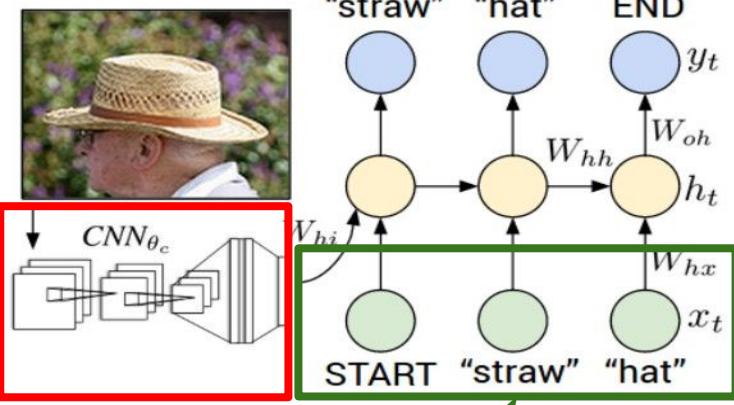
# Transfer learning is pervasive

Object Detection  
(Fast R-CNN)



CNN pretrained  
on ImageNet

Image Captioning: CNN + RNN



Word vectors pretrained  
with word2vec

Girshick, "Fast R-CNN", ICCV 2015  
Figure copyright Ross Girshick, 2015. Reproduced with permission.

Karpathy and Fei-Fei, "Deep Visual-Semantic Alignments for Generating Image Descriptions", CVPR 2015  
Figure copyright IEEE, 2015. Reproduced for educational purposes.