

# A Few CNN Case Studies

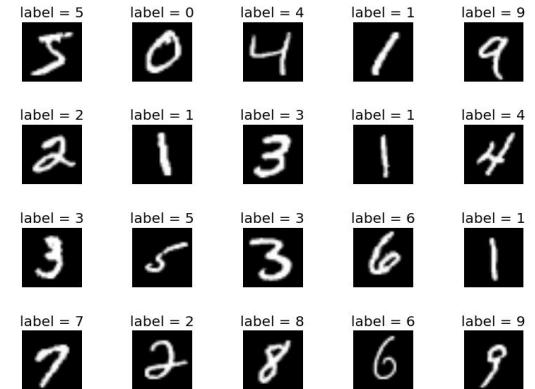
# Case Studies

## 1. Hand Written Digit Classification (LeNet - 1998)

**input:** a small single channel image

**output:** 10 outputs corresponding to the 10 digits 0-9.

60,000 training images, 10,000 test images



## 2. Image Net Classification – Annual world cup for CV

**input:** colored image

**output:** 1000 outputs corresponding to the 1000 object classes in the dataset

1.2 M training images and 100,000 test images

IMAGENET

- 1,000 object classes (categories).
- Images:
  - 1.2 M train
  - 100k test.



## CNNs on MNIST

### 1. LeNet (1998)

- 10 way neural network classifier
- Handwritten digits as an input
- Tolerant of various transformations like rotation and scale
- Was used by banks to recognize handwritten numbers on digitized checks
- 4 weight layers

## CNNs on ImageNet

1. AlexNet (2012)
  - First CNN to successfully be able classify ImageNet images
  - Improved benchmark performance (top-5) on this image dataset from 26% to 15%
  - 7 layers deep
2. ZF Net (2013)
  - Reduced the top-5 error rate to 11.2%
  - No major contributions
  - Also 7 layers deep
3. VGGNet (2014)
  - Simple and elegant
  - Reduced the top-5 error rate 7.2%
  - Did not win the competition, GoogleNet did!
  - 6 layers deep

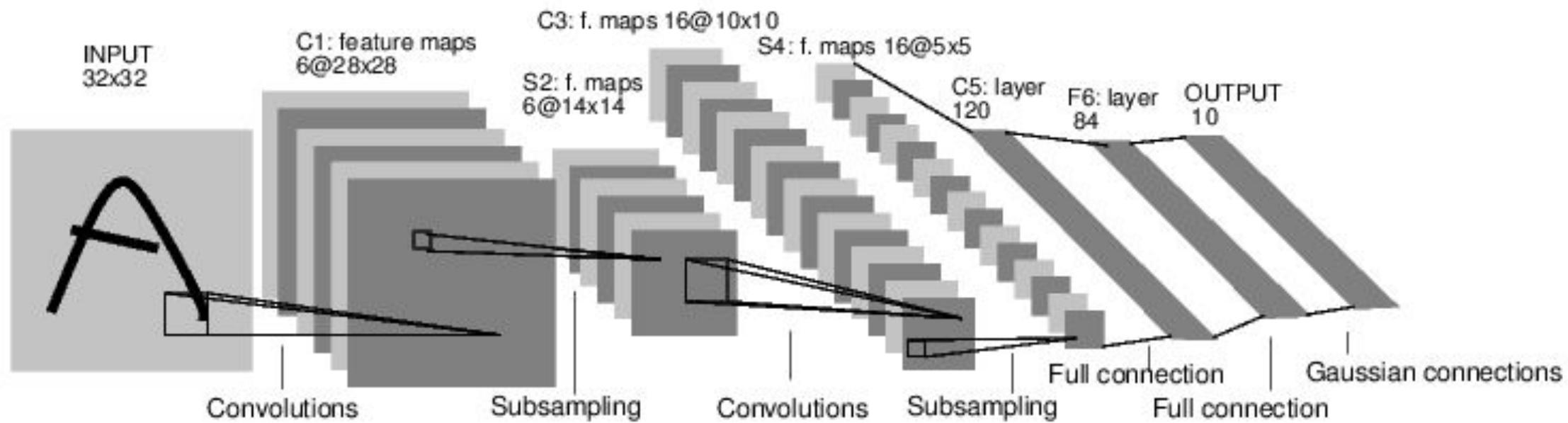
# Case Studies

## CNNs on ImageNet

4. GoogleNet (2014)
  - 2014 imagenet winner with top-5 error rate of 6.7%
  - Used inception modules
  - 22 layers deep and used side cost functions
5. ResNet (2015)
  - 2015 imagenet winner with top-5 error rate of 3.57
  - First truly deep network with 152 weight layers
6. CUIimage (2016)
  - 2016 imagenet winner with top-5 error rate of 2.99
  - Ensemble approach, not very interesting
7. SENet (2017)
  - 2016 ImageNet winner with top-5 error rate of 2.251
  - Work by Momenta
  - The last ImageNet challenge!

# Case Study: LeNet-5

[LeCun et al., 1998]

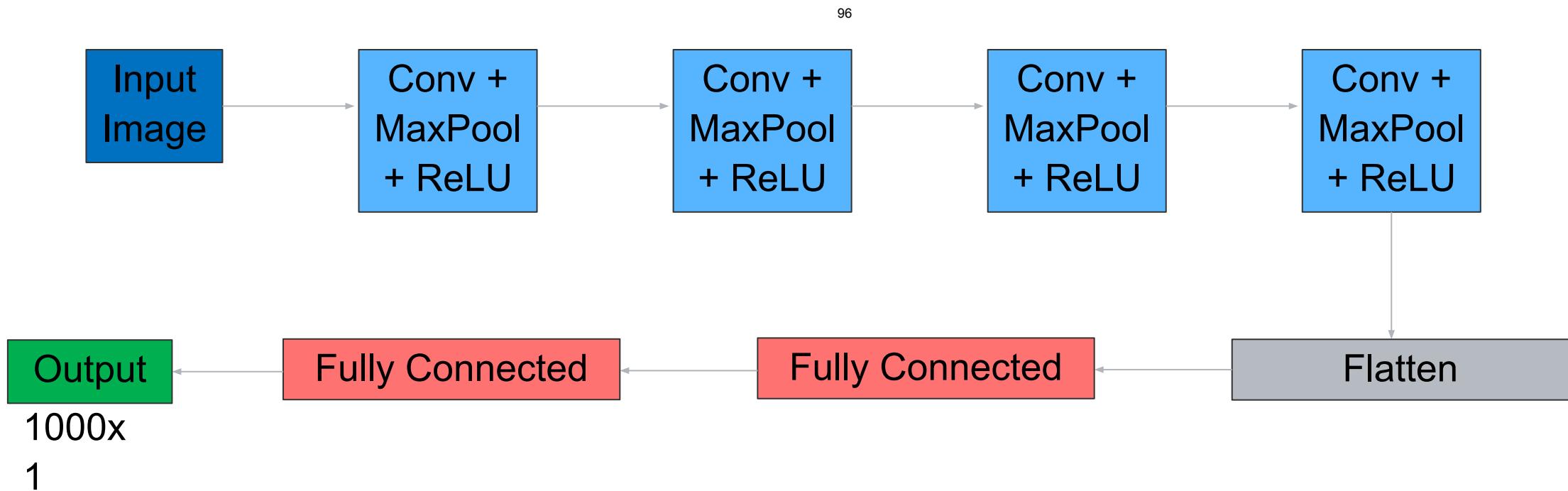


- Conv filters were 5x5, applied at stride 1
- Subsampling (Pooling) layers were 2x2 applied at stride 2
- Architecture is [CONV-POOL-CONV-POOL-CONV-FC]

Source: Gradient Based Learning Applied to Document Recognition, LeCun et al. (1998)

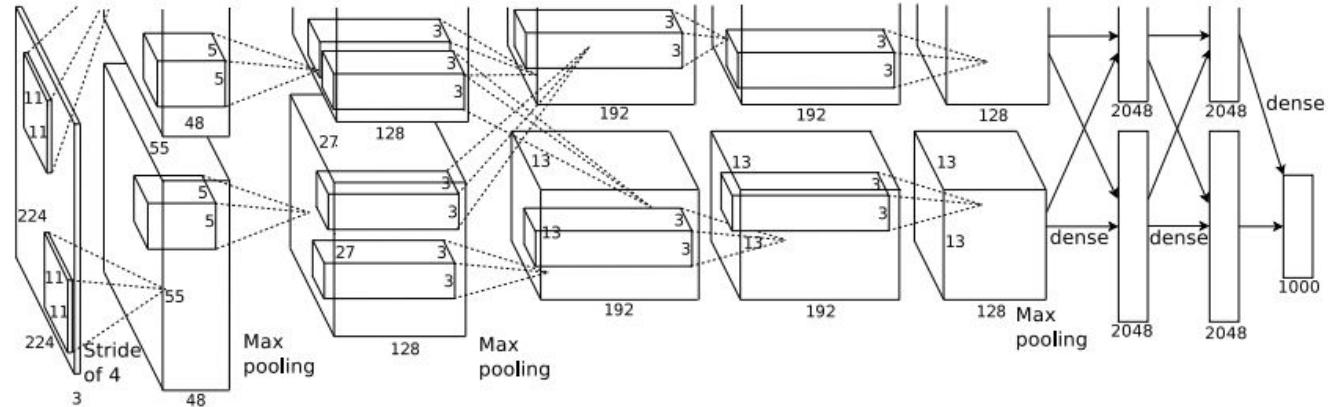
# Case Study: AlexNet

[Krizhevsky et al. 2012]



# Case Study: AlexNet

[Krizhevsky et al. 2012]

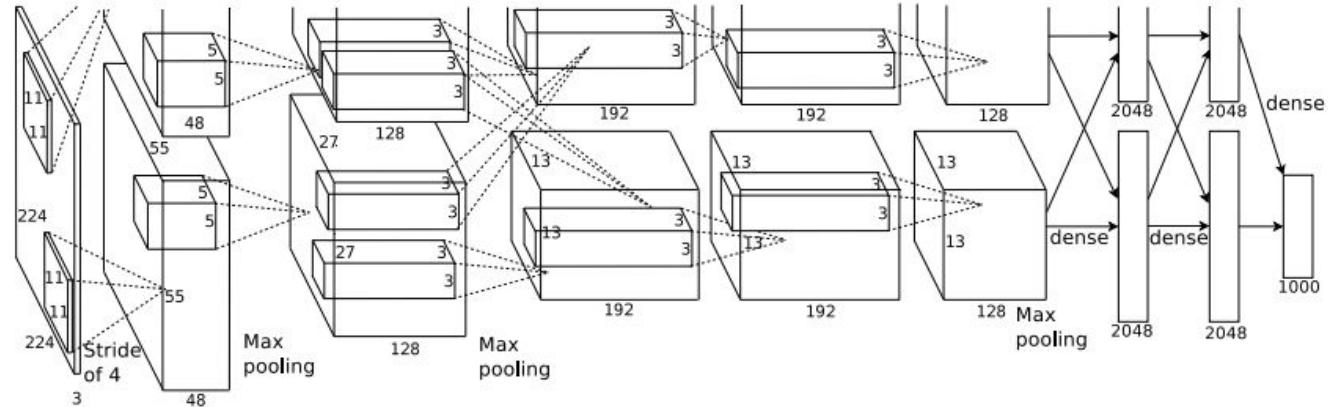


- **Input:** 227x227x3 images
- **First layer (CONV1):** 96 11x11 filters applied at stride 4
- Output volume size? (Hint:  $(227-11)/4+1$ )

Sourced with permission from: *ImageNet Classification with Deep Convolutional Neural Networks*, Krizhevsky et al. (2012)

# Case Study: AlexNet

[Krizhevsky et al. 2012]

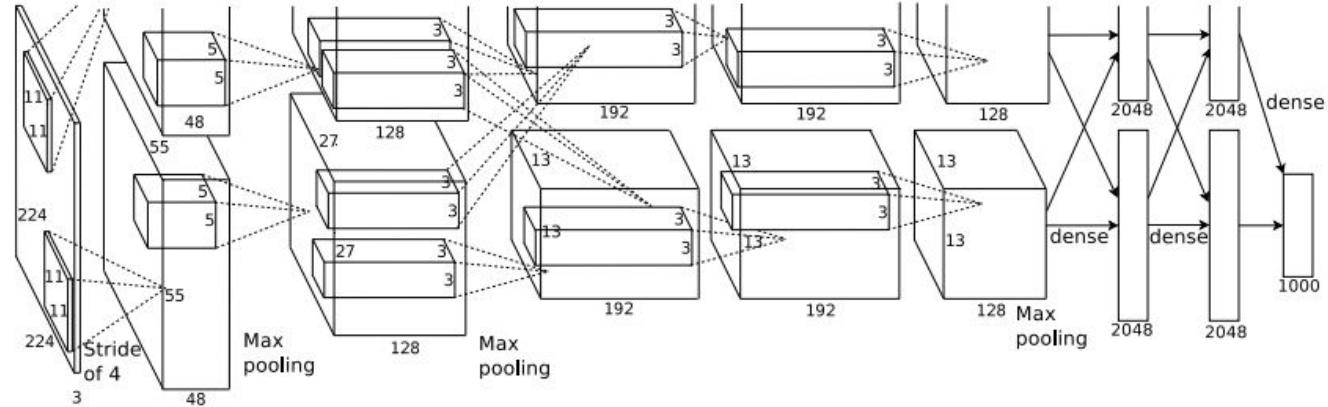


- **Input:** 227x227x3 images
- **First layer (CONV1):** 96 11x11 filters applied at stride 4
- **Output volume size:**  $(227-11)/4+1 = 55$  for each H and W, so 55x55x96

Sourced with permission from: *ImageNet Classification with Deep Convolutional Neural Networks*, Krizhevsky et al. (2012)

# Case Study: AlexNet

[Krizhevsky et al. 2012]

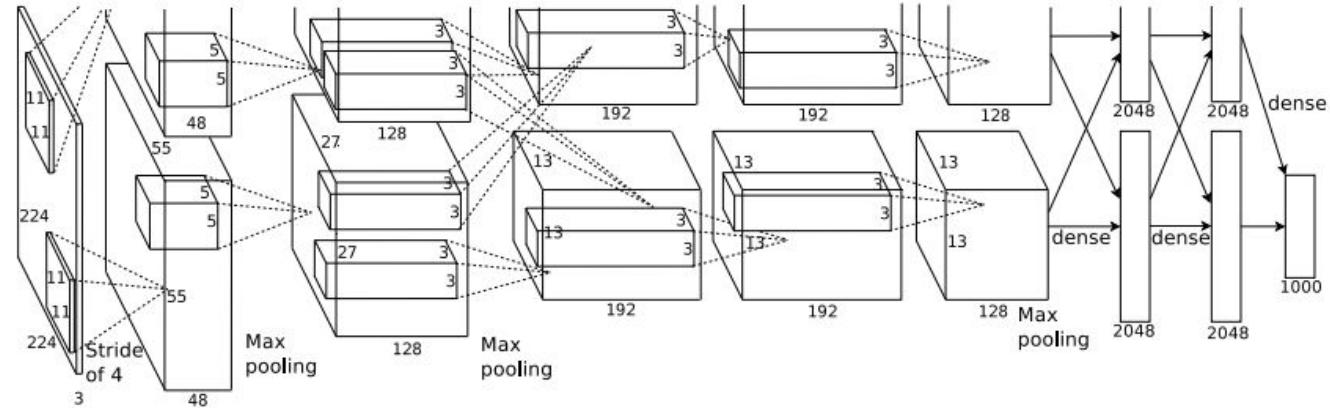


- **Input:** 227x227x3 images
- **First layer (CONV1):** 96 11x11 filters applied at stride 4
- **Output volume size:** 55x55x96
- Total number of parameters?

Sourced with permission from: *ImageNet Classification with Deep Convolutional Neural Networks*, Krizhevsky et al. (2012)

# Case Study: AlexNet

[Krizhevsky et al. 2012]

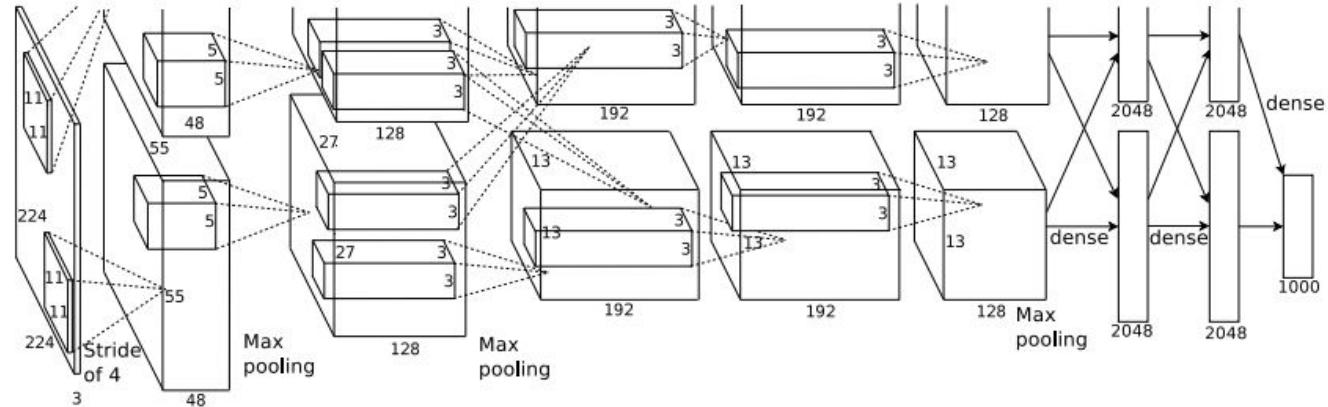


- **Input:** 227x227x3 images
- **First layer (CONV1):** 96 11x11 filters applied at stride 4
- **Output volume size:** 55x55x96
- **Total number of parameters:**  $(11*11*3)*96 = 35K$

Sourced with permission from: *ImageNet Classification with Deep Convolutional Neural Networks*, Krizhevsky et al. (2012)

# Case Study: AlexNet

[Krizhevsky et al. 2012]

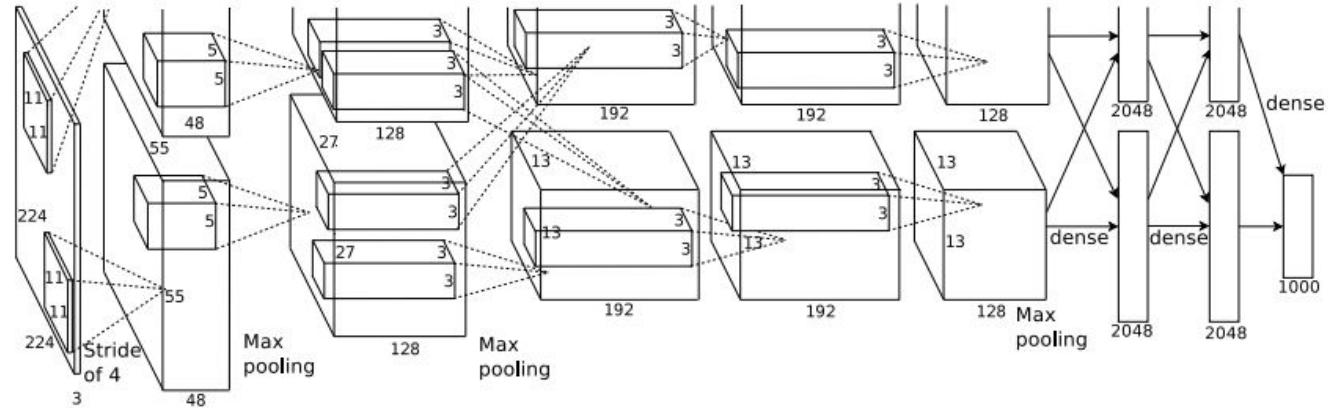


- **Input:** 227x227x3 images
- **After CONV1:** 55x55x96
- **Second layer (POOL1):** 3x3 filters applied
- What is the output volume size? (Hint:  $(55-3)/2+1 = 27$ )

Sourced with permission from: *ImageNet Classification with Deep Convolutional Neural Networks*, Krizhevsky et al. (2012)

# Case Study: AlexNet

[Krizhevsky et al. 2012]

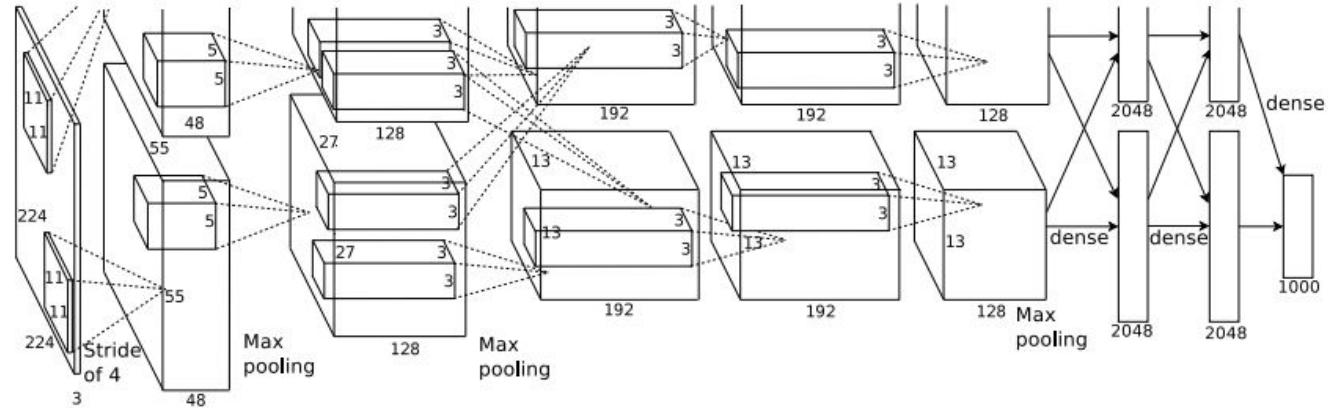


- **Input:** 227x227x3 images
- **After CONV1:** 55x55x96
- **Second layer (POOL1):** 3x3 filters applied at stride 2
- Output volume: 27x27x96
- What is the number of parameters?

Sourced with permission from: *ImageNet Classification with Deep Convolutional Neural Networks*, Krizhevsky et al. (2012)

# Case Study: AlexNet

[Krizhevsky et al. 2012]

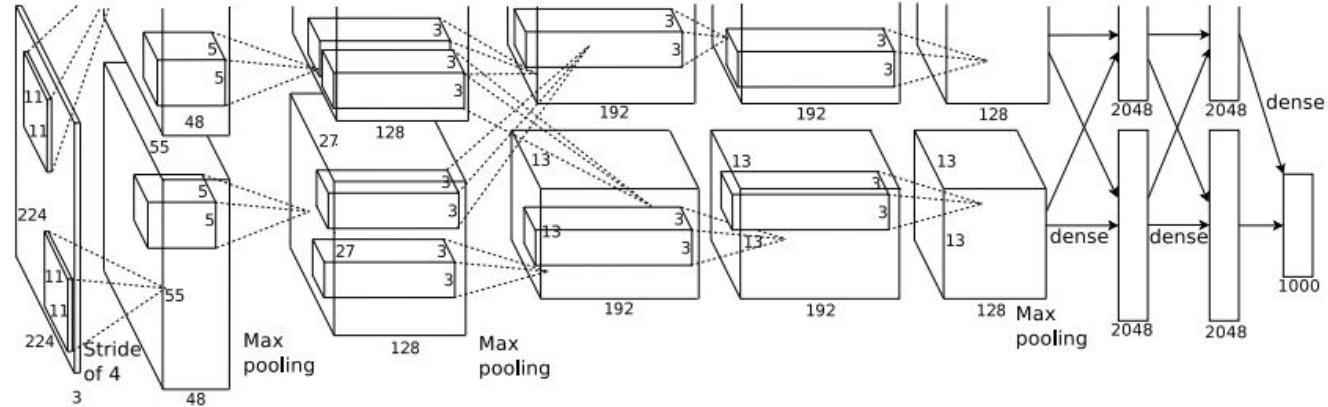


- **Input:** 227x227x3 images
- **After CONV1:** 55x55x96
- **Second layer (POOL1):** 3x3 filters applied at stride 2
- Output volume: 27x27x96
- What is the number of parameters: **0!**

Sourced with permission from: *ImageNet Classification with Deep Convolutional Neural Networks*, Krizhevsky et al. (2012)

# Case Study: AlexNet

[Krizhevsky et al. 2012]

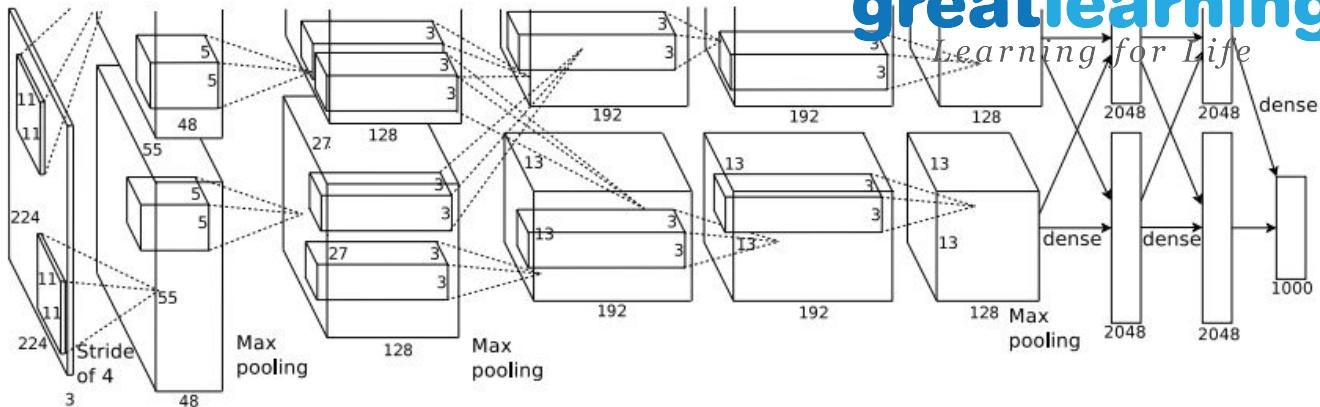


- **Input:** 227x227x3 images
- **After CONV1:** 55x55x96
- **After POOL1:** 27x27x96

Sourced with permission from: *ImageNet Classification with Deep Convolutional Neural Networks*, Krizhevsky et al. (2012)

# Case Study: AlexNet

[Krizhevsky et al. 2012]



## 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

## Finishing with:

[4096] FC6: 4096 neurons

[4096] FC7: 4096 neurons

[1000] FC8: 1000 neurons (class scores)

Sourced with permission from: ImageNet Classification with Deep Convolutional Neural Networks, Krizhevsky et al. (2012)

# Case Study: AlexNet

[Krizhevsky et al. 2012]

greatlearning

## 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

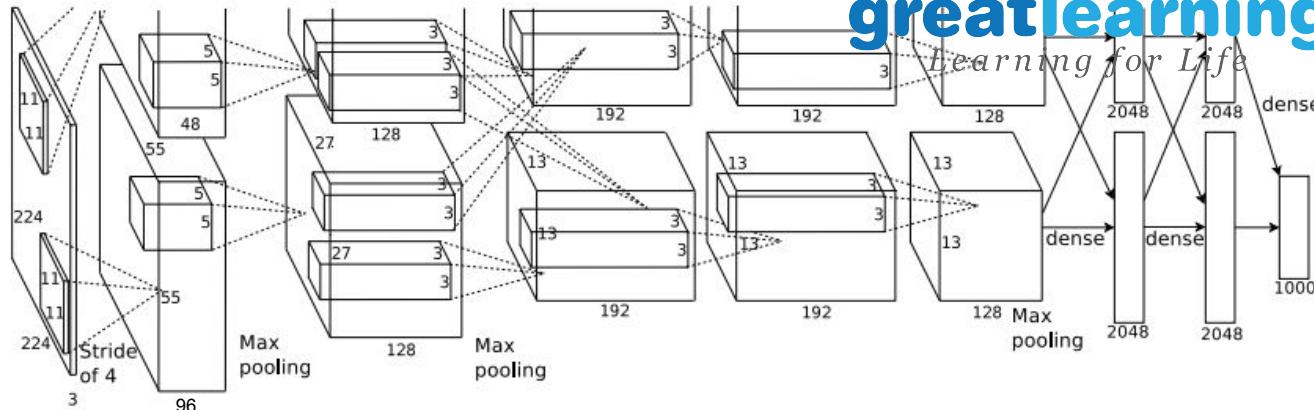
[13x13x256] CONV5: 256 3x3 filters at stride 1, pad 1

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

[4096] FC6: 4096 neurons

[4096] FC7: 4096 neurons

[1000] FC8: 1000 neurons (class scores)



## Salient points:

- Popularized use of ReLU in Vision
- Used Norm layers (not common anymore)
- Heavy data augmentation
- Dropout 0.5 in only last few fully-connected
- 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% improved to 15.4%**

Sourced with permission from: *ImageNet Classification with Deep Convolutional Neural Networks*, Krizhevsky et al. (2012)

# Case Study: AlexNet

[Krizhevsky et al. 2012]

## Architecture:

[227x227x3] INPUT

[55x55x96] CONV1: 96 11x11 filters at stride 4, pad 0 [55x55x48] x 2

[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

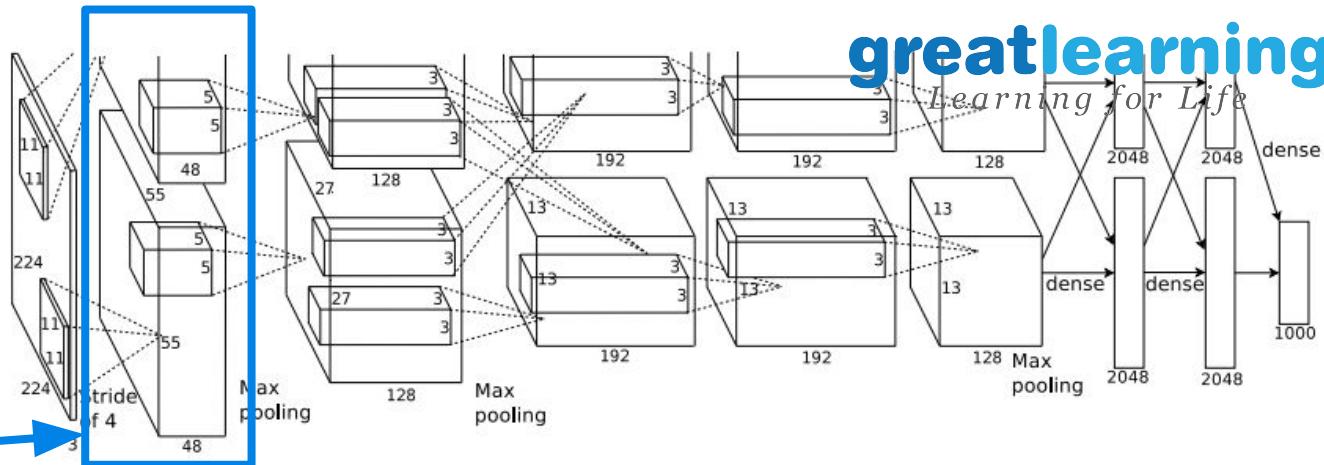
[13x13x256] CONV5: 256 3x3 filters at stride 1, pad 1

[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 GTX580 GPU with only 3 GB of memory. Network spread across 2 GPUs, half the feature maps on each GPU.

Sourced with permission from: ImageNet Classification with Deep Convolutional Neural Networks, Krizhevsky et al. (2012)

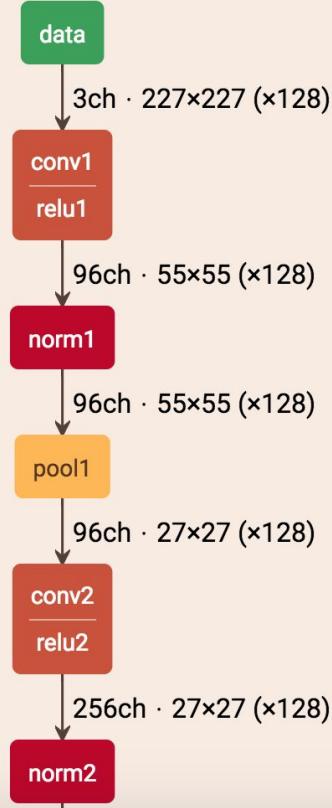
# A tool to analyze deep networks

<http://dgschwend.github.io/netscope/#/editor>

```

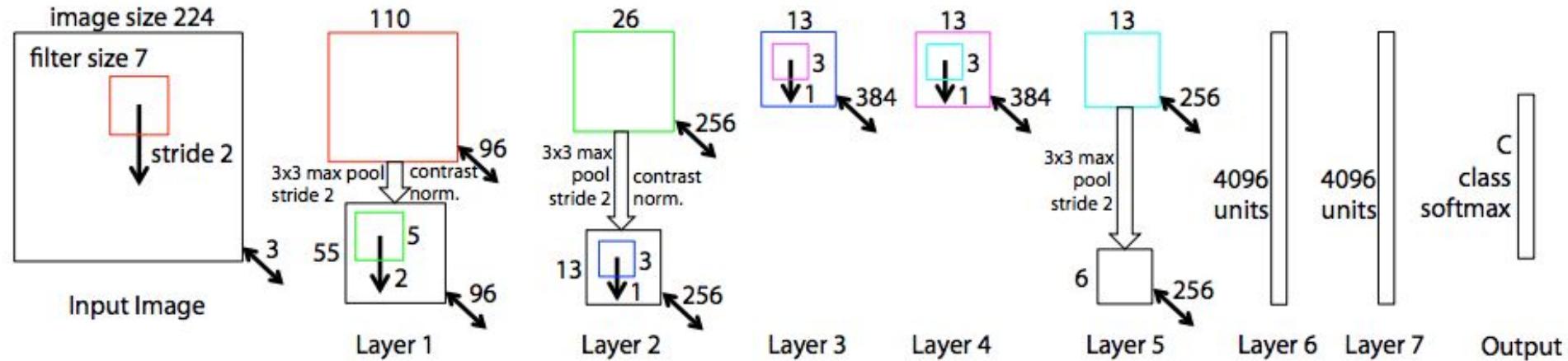
1 name: "AlexNet"
2 layer {
3   name: "data"
4   type: "Data"
5   top: "data"
6   input_param {
7     shape: {
8       dim: 128
9       dim: 3
10      dim: 227
11      dim: 227
12    }
13  }
14 }
15 layer {
16   name: "conv1"
17   type: "Convolution"
18   bottom: "data"
19   top: "conv1"
20   param {
21     lr_mult: 1
22     decay_mult: 1
23   }
24   param {
25     lr_mult: 2
26     decay_mult: 0
27   }
28   convolution_param {
29     num_output: 96
30     kernel_size: 11
31     stride: 4
32     weight_filler {
33       type: "gaussian"
34       std: 0.01
35     }
36     bias_filler {
37       type: "constant"
38       value: 0
39     }
40   }
41 }
42 
```

AlexNet (edit)



# Case Study - ZFNet

[Zeiler and Fergus, 2013]



Similar to AlexNet with the following differences:

**CONV1:** (7x7 stride 2) instead of (11x11 stride 4)

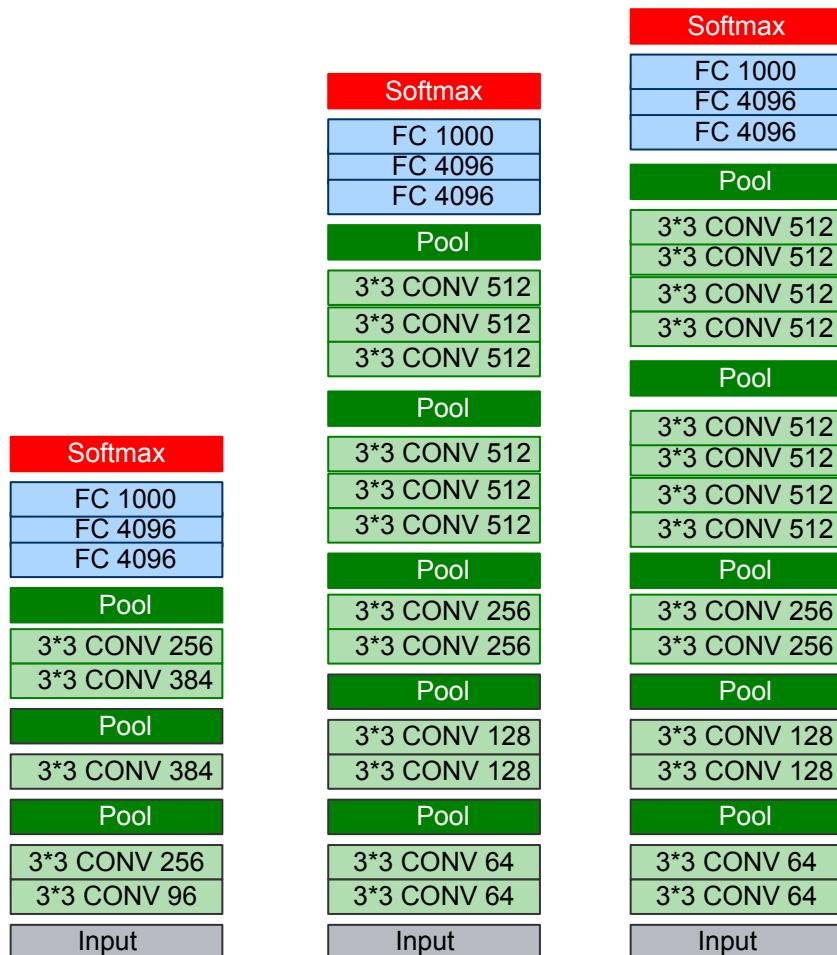
**CONV3,4,5:** 512, 1024, 512 filters instead of 384, 384, 256 respectively

Reduced top 5 error on ImageNet  
From **15.4%** To **14.8%**  
*Later brought down to 11.2%*

Sourced with permission from 'Visualizing and Understanding Convolutional Networks' by Zeiler, Fergus (2013)

# Case Study: VGGNet

[Simonyan and Zisserman, 2014]



**AlexNet**

**VGG 16**

**VGG 19**

**This model used:**

- Smaller filters**  
But
- Deeper networks**

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

**Why use smaller filters? (3x3 conv)**

Answer: Stack of three 3x3 conv (stride 1) layers has same effective receptive field as one 7x7 but deeper, more non-linearities and fewer parameters.

# Case Study: VGGNet

[Simonyan and Zisserman, 2014]

This model used:

- Smaller filters
- But
- Deeper networks

3x3 CONV stride 1, pad 1

2x2 MAX POOL stride 2

Improved from 11.2% top 5 error  
in ILSVRC 2013

To 7.3% top 5 error

And yet, this model did not win!

Sourced with permission from: 'Very deep  
large-scale image recognition, Simonyan & Zisserman (2015)

Convolutional networks for visual recognition. All Rights Reserved. Unauthorized use or distribution prohibited

ConvNet Configuration					
A	A-LRN	B	C	D	E
11 weight layers	11 weight layers	13 weight layers	16 weight layers	16 weight layers	19 weight layers
input (224 × 224 RGB image)					
conv3-64	conv3-64 LRN	conv3-64 <b>conv3-64</b>	conv3-64 conv3-64	conv3-64 conv3-64	conv3-64 conv3-64
maxpool					
conv3-128	conv3-128	conv3-128 <b>conv3-128</b>	conv3-128 conv3-128	conv3-128 conv3-128	conv3-128 conv3-128
maxpool					
conv3-256 conv3-256	conv3-256 conv3-256	conv3-256 conv3-256	conv3-256 conv3-256 <b>conv1-256</b>	conv3-256 conv3-256 <b>conv3-256</b>	conv3-256 conv3-256 conv3-256 <b>conv3-256</b>
maxpool					
conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512 <b>conv1-512</b>	conv3-512 conv3-512 <b>conv3-512</b>	conv3-512 conv3-512 conv3-512 <b>conv3-512</b>
maxpool					
conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512 <b>conv1-512</b>	conv3-512 conv3-512 <b>conv3-512</b>	conv3-512 conv3-512 conv3-512 <b>conv3-512</b>
maxpool					
FC-4096					
FC-4096					
FC-1000					
soft-max					

Table 2: Number of parameters (in millions).

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

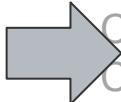
Best performing model

ConvNet Configuration			
B	C	D	
13 weight layers	16 weight layers	16 weight layers	19
put (224 × 224 RGB image)			
conv3-64	conv3-64	conv3-64	cc
<b>conv3-64</b>	conv3-64	conv3-64	cc
maxpool			
conv3-128	conv3-128	conv3-128	co
<b>conv3-128</b>	conv3-128	conv3-128	co
maxpool			
conv3-256	conv3-256	conv3-256	co
conv3-256	conv3-256	conv3-256	co
<b>conv1-256</b>		<b>conv3-256</b>	co
maxpool			
conv3-512	conv3-512	conv3-512	co
conv3-512	conv3-512	conv3-512	co
<b>conv1-512</b>		<b>conv3-512</b>	co
maxpool			
conv3-512	conv3-512	conv3-512	co
conv3-512	conv3-512	conv3-512	co
<b>conv1-512</b>		<b>conv3-512</b>	co
maxpool			
FC-4096			
FC-4096			
FC-1000			
soft-max			



INPUT: [224x224x3]  
 CONV3-64: [224x224x64]  
 CONV3-64: [224x224x64]  
 POOL2: [112x112x64]  
 CONV3-128: [112x112x128]  
 CONV3-128: [112x112x128]  
 POOL2: [56x56x128]  
 CONV3-256: [56x56x256]  
 CONV3-256: [56x56x256]  
 CONV3-256: [56x56x256]  
 POOL2: [28x28x256]  
 CONV3-512: [28x28x512]  
 CONV3-512: [28x28x512]  
 CONV3-512: [28x28x512]  
 POOL2: [14x14x512]  
 CONV3-512: [14x14x512]  
 CONV3-512: [14x14x512]  
 CONV3-512: [14x14x512]  
 POOL2: [7x7x512]  
 FC: [1x1x4096]  
 FC: [1x1x4096]  
 FC: [1x1x1000]

ConvNet Configuration			
B	C	D	
13 weight layers	16 weight layers	16 weight layers	19
put (224 x 224 RGB image)			
conv3-64	conv3-64	conv3-64	cc
<b>conv3-64</b>	conv3-64	conv3-64	cc
maxpool			
conv3-128	conv3-128	conv3-128	co
<b>conv3-128</b>	conv3-128	conv3-128	co
maxpool			
conv3-256	conv3-256	conv3-256	co
conv3-256	conv3-256	conv3-256	co
<b>conv1-256</b>		<b>conv3-256</b>	co
maxpool			
conv3-512	conv3-512	conv3-512	co
conv3-512	conv3-512	conv3-512	co
<b>conv1-512</b>		<b>conv3-512</b>	co
maxpool			
conv3-512	conv3-512	conv3-512	co
conv3-512	conv3-512	conv3-512	co
<b>conv1-512</b>		<b>conv3-512</b>	co
maxpool			
FC-4096			
FC-4096			
FC-1000			
soft-max			

INPUT: [224x224x3]  
 CONV3-64: [224x224x64]  
 CONV3-64: [224x224x64]  
 POOL2: [112x112x64]  
 CONV3-128: [112x112x128]  
 CONV3-128: [112x112x128]  
 POOL2: [56x56x128]  
 CONV3-256: [56x56x256]  
  
 CONV3-256: [56x56x256]  
 CONV3-256: [56x56x256]  
 POOL2: [28x28x256]  
 CONV3-512: [28x28x512]  
 CONV3-512: [28x28x512]  
 CONV3-512: [28x28x512]  
 POOL2: [14x14x512]  
 CONV3-512: [14x14x512]  
 CONV3-512: [14x14x512]  
 CONV3-512: [14x14x512]  
 POOL2: [7x7x512]  
 FC: [1x1x4096]  
 FC: [1x1x4096]  
 FC: [1x1x1000]

## MEMORY

$$224 \times 224 \times 3 = 150K$$

$$224 \times 224 \times 64 = 3.2M$$

$$224 \times 224 \times 64 = 3.2M$$

$$112 \times 112 \times 64 = 800K$$

$$112 \times 112 \times 128 = 1.6M$$

$$112 \times 112 \times 128 = 1.6M$$

$$56 \times 56 \times 128 = 400K$$

$$56 \times 56 \times 256 = 800K$$

$$56 \times 56 \times 256 = 800K$$

$$28 \times 28 \times 256 = 200K$$

$$28 \times 28 \times 512 = 400K$$

$$28 \times 28 \times 512 = 400K$$

$$28 \times 28 \times 512 = 400K$$

$$14 \times 14 \times 512 = 100K$$

$$7 \times 7 \times 512 = 25K$$

$$4096$$

$$4096$$

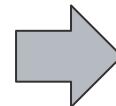
$$1000$$

Total memory:

$$24M * 4 \text{ bytes} \approx 93MB/\text{image}$$

Only for forward. What if we include backward?

ConvNet Configuration			
B	C	D	
13 weight layers	16 weight layers	16 weight layers	19
put (224 × 224 RGB image)			
conv3-64	conv3-64	conv3-64	cc
<b>conv3-64</b>	conv3-64	conv3-64	cc
maxpool			
conv3-128	conv3-128	conv3-128	co
<b>conv3-128</b>	conv3-128	conv3-128	co
maxpool			
conv3-256	conv3-256	conv3-256	co
conv3-256	conv3-256	conv3-256	co
<b>conv1-256</b>		<b>conv3-256</b>	co
maxpool			
conv3-512	conv3-512	conv3-512	co
conv3-512	conv3-512	conv3-512	co
<b>conv1-512</b>		<b>conv3-512</b>	co
maxpool			
conv3-512	conv3-512	conv3-512	co
conv3-512	conv3-512	conv3-512	co
<b>conv1-512</b>		<b>conv3-512</b>	co
maxpool			
FC-4096			
FC-4096			
FC-1000			
soft-max			



INPUT: [224x224x3]  
 CONV3-64: [224x224x64]  
 CONV3-64: [224x224x64]  
 POOL2: [112x112x64]  
 CONV3-128: [112x112x128]  
 CONV3-128: [112x112x128]  
 POOL2: [56x56x128]  
 CONV3-256: [56x56x256]  
 CONV3-256: [56x56x256]  
 CONV3-256: [56x56x256]  
 POOL2: [28x28x256]  
 CONV3-512: [28x28x512]  
 CONV3-512: [28x28x512]  
 CONV3-512: [28x28x512]  
 POOL2: [14x14x512]  
 CONV3-512: [14x14x512]  
 CONV3-512: [14x14x512]  
 CONV3-512: [14x14x512]  
 POOL2: [7x7x512]  
 FC: [1x1x4096]  
 FC: [1x1x4096]  
 FC: [1x1x1000]

138 Million  
total parameters!

224*224*3=150K	0
(3*3*3)*64 = 1,728	
(3*3*64)*64 = 36,864	
0	
(3*3*64)*128 = 73,728	
(3*3*128)*128 = 147,456	
0	
(3*3*128)*256 = 294,912	
(3*3*256)*256 = 589,824	
(3*3*256)*256 = 589,824	
0	
(3*3*256)*512 = 1,179,648	
(3*3*512)*512 = 2,359,296	
(3*3*512)*512 = 2,359,296	
0	
(3*3*512)*512 = 2,359,296	
(3*3*512)*512 = 2,359,296	
(3*3*512)*512 = 2,359,296	
0	
7*7*512*4096 = 102,760,448	
4096*4096 = 16,777,216	
4096*1000 = 4,096,000	

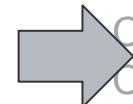
ConvNet Configuration			
B	C	D	
13 weight layers	16 weight layers	16 weight layers	19
put (224 x 224 RGB image)			
conv3-64	conv3-64	conv3-64	cc
<b>conv3-64</b>	conv3-64	conv3-64	cc
maxpool			
conv3-128	conv3-128	conv3-128	co
<b>conv3-128</b>	conv3-128	conv3-128	co
maxpool			
conv3-256	conv3-256	conv3-256	co
conv3-256	conv3-256	conv3-256	co
<b>conv1-256</b>		<b>conv3-256</b>	co
maxpool			
conv3-512	conv3-512	conv3-512	co
conv3-512	conv3-512	conv3-512	co
<b>conv1-512</b>		<b>conv3-512</b>	co
maxpool			
conv3-512	conv3-512	conv3-512	co
conv3-512	conv3-512	conv3-512	co
<b>conv1-512</b>		<b>conv3-512</b>	co
maxpool			
FC-4096			
FC-4096			
FC-1000			
soft-max			

INPUT: [224x224x3]  
**224\*224\*64=3.2M**  
**224\*224\*64=3.2M**  
 POOL2: [112x112x64]  
 CONV3-128: [112x112x128]  
 CONV3-128: [112x112x128]  
 POOL2: [56x56x128]  
 CONV3-256: [56x56x256]  
**CONV3-256: [56x56x256]**  
**CONV3-256: [56x56x256]**  
 POOL2: [28x28x256]  
 CONV3-512: [28x28x512]  
 CONV3-512: [28x28x512]  
 CONV3-512: [28x28x512]  
 POOL2: [14x14x512]  
 CONV3-512: [14x14x512]  
 CONV3-512: [14x14x512]  
 CONV3-512: [14x14x512]  
 POOL2: [7x7x512]  
 FC: [1x1x4096]  
 FC: [1x1x4096]  
 FC: [1x1x1000]

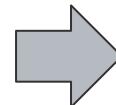
## MEMORY

224\*224\*3=150K  
**224\*224\*64=3.2M**  
**224\*224\*64=3.2M**  
 112\*112\*64=800K  
 112\*112\*128=1.6M  
 112\*112\*128=1.6M  
 56\*56\*128=400K  
 56\*56\*256=800K  
**56\*56\*256=800K**  
**56\*56\*256=800K**  
 28\*28\*256=200K  
 28\*28\*512=400K  
 28\*28\*512=400K  
 28\*28\*512=400K  
 14\*14\*512=100K  
 14\*14\*512=100K  
 14\*14\*512=100K  
 14\*14\*512=100K  
 7\*7\*512=25K  
 4096  
 4096  
 1000

Most memory in early CONV layers



ConvNet Configuration			
B	C	D	
13 weight layers	16 weight layers	16 weight layers	19
put (224 × 224 RGB image)			
conv3-64	conv3-64	conv3-64	cc
<b>conv3-64</b>	conv3-64	conv3-64	cc
maxpool			
conv3-128	conv3-128	conv3-128	co
<b>conv3-128</b>	conv3-128	conv3-128	co
maxpool			
conv3-256	conv3-256	conv3-256	co
conv3-256	conv3-256	conv3-256	co
<b>conv1-256</b>		<b>conv3-256</b>	co
maxpool			
conv3-512	conv3-512	conv3-512	co
conv3-512	conv3-512	conv3-512	co
<b>conv1-512</b>		<b>conv3-512</b>	co
maxpool			
conv3-512	conv3-512	conv3-512	co
conv3-512	conv3-512	conv3-512	co
<b>conv1-512</b>		<b>conv3-512</b>	co
maxpool			
FC-4096			
FC-4096			
FC-1000			
soft-max			



INPUT: [224x224x3]  
 CONV3-64: [224x224x64]  
 CONV3-64: [224x224x64]  
 POOL2: [112x112x64]  
 CONV3-128: [112x112x128]  
 CONV3-128: [112x112x128]  
 POOL2: [56x56x128]  
 CONV3-256: [56x56x256]  
 CONV3-256: [56x56x256]  
 CONV3-256: [56x56x256]  
 POOL2: [28x28x256]  
 CONV3-512: [28x28x512]  
 CONV3-512: [28x28x512]  
 CONV3-512: [28x28x512]  
 POOL2: [14x14x512]  
 CONV3-512: [14x14x512]  
 CONV3-512: [14x14x512]  
 CONV3-512: [14x14x512]  
 POOL2: [7x7x512]  
 FC: [1x1x4096]  
 FC: [1x1x4096]  
 FC: [1x1x1000]

## MEMORY

224*224*3=150K	0
(3*3*3)*64 = 1,728	
(3*3*64)*64 = 36,864	
112*112*64=800K	0
112*112*128=1.6M	(3*3*64)*128 = 73,728
112*112*128=1.6M	(3*3*128)*128 = 147,456
56*56*128=400K	0
56*56*256=800K	(3*3*128)*256 = 294,912
56*56*256=800K	(3*3*256)*256 = 589,824
56*56*256=800K	(3*3*256)*256 = 589,824
28*28*256=200K	0
28*28*512=400K	(3*3*256)*512 = 1,179,648
28*28*512=400K	(3*3*512)*512 = 2,359,296
28*28*512=400K	(3*3*512)*512 = 2,359,296
14*14*512=100K	0
14*14*512=100K	(3*3*512)*512 = 2,359,296
14*14*512=100K	(3*3*512)*512 = 2,359,296
14*14*512=100K	(3*3*512)*512 = 2,359,296
7*7*512=25K	0
7*7*512*4096 = 102,760,448	
4096*4096 = 16,777,216	
4096*1000 = 4,096,000	

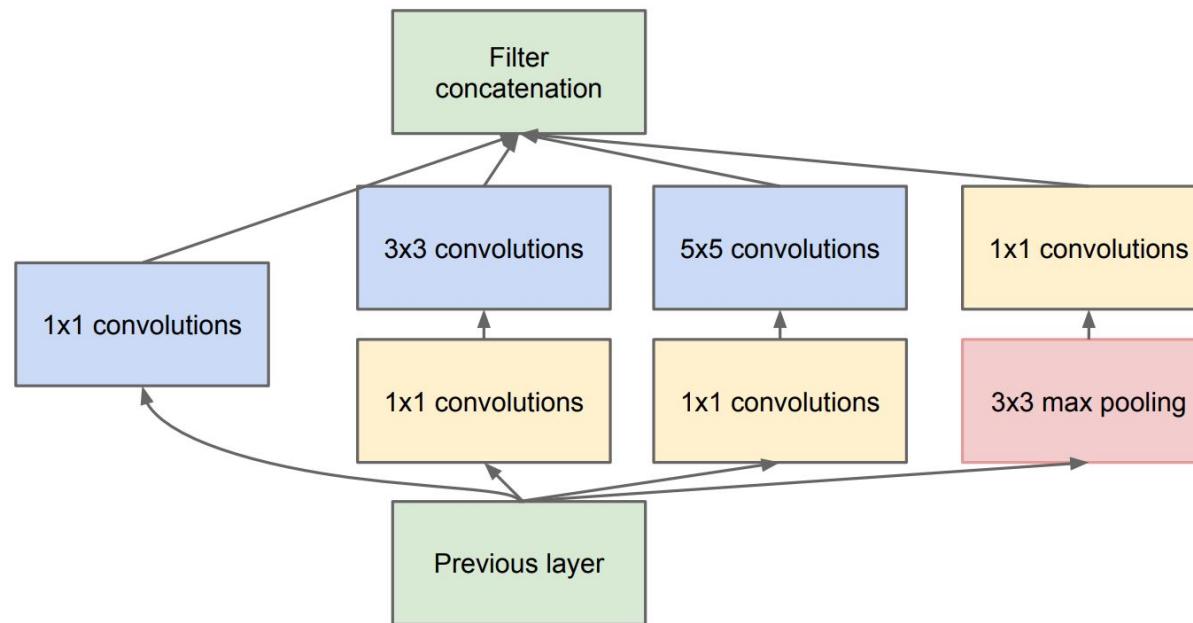
Most parameters  
in late FC

1000

Parameters not including biases

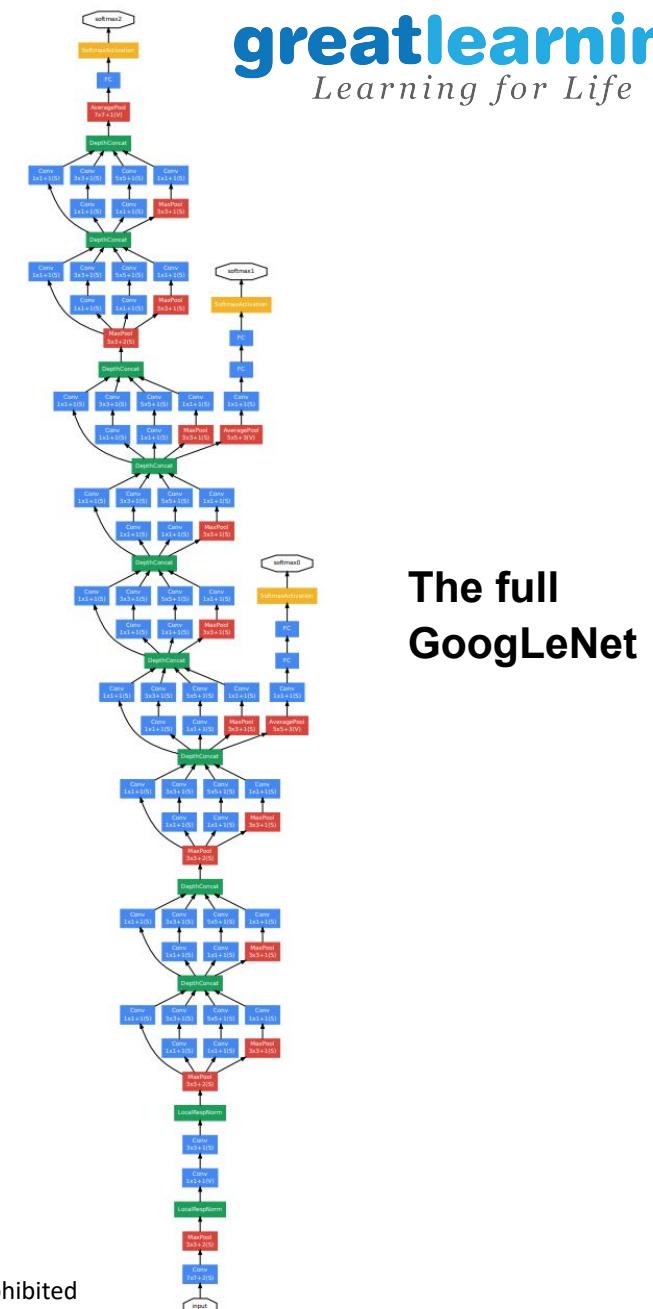
# Case Study: GoogLeNet

[Szegedy et al., 2014]



### Inception module – with dimension reductions

Winner of ILSVRC 2014 with **6.7%** top 5 error



# The full GoogLeNet

Sourced with permission from: 'Going Deeper with Convolutions', Szegedy et al. (2014)

# Case Study: GoogLeNet

[Szegedy et al., 2014]

type	patch size/ stride	output size	depth	#1×1	#3×3 reduce	#3×3	#5×5 reduce	#5×5	pool proj	params	ops
convolution	7×7/2	112×112×64	1							2.7K	34M
max pool	3×3/2	56×56×64	0								
convolution	3×3/1	56×56×192	2		64	192				112K	360M
max pool	3×3/2	28×28×192	0								
inception (3a)		28×28×256	2	64	96	128	16	32	32	159K	128M
inception (3b)		28×28×480	2	128	128	192	32	96	64	380K	304M
max pool	3×3/2	14×14×480	0								
inception (4a)		14×14×512	2	192	96	208	16	48	64	364K	73M
inception (4b)		14×14×512	2	160	112	224	24	64	64	437K	88M
inception (4c)		14×14×512	2	128	128	256	24	64	64	463K	100M
inception (4d)		14×14×528	2	112	144	288	32	64	64	580K	119M
inception (4e)		14×14×832	2	256	160	320	32	128	128	840K	170M
max pool	3×3/2	7×7×832	0								
inception (5a)		7×7×832	2	256	160	320	32	128	128	1072K	54M
inception (5b)		7×7×1024	2	384	192	384	48	128	128	1388K	71M
avg pool	7×7/1	1×1×1024	0								
dropout (40%)		1×1×1024	0								
linear		1×1×1000	1							1000K	1M
softmax		1×1×1000	0								

This model has only 5 million parameters! (Removes FC layers completely)

**Compared to AlexNet, this model has:** 12X less params | 2x more compute | 6.67% top-5 error rate vs. 16.4%

# Case Study: ResNet

[He et al., 2015]

greatlearning  
Learning for Life

Winner of ILSVRC 2015  
3.6% top-5 error!

Microsoft  
Research

## 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

\*improvements are relative numbers

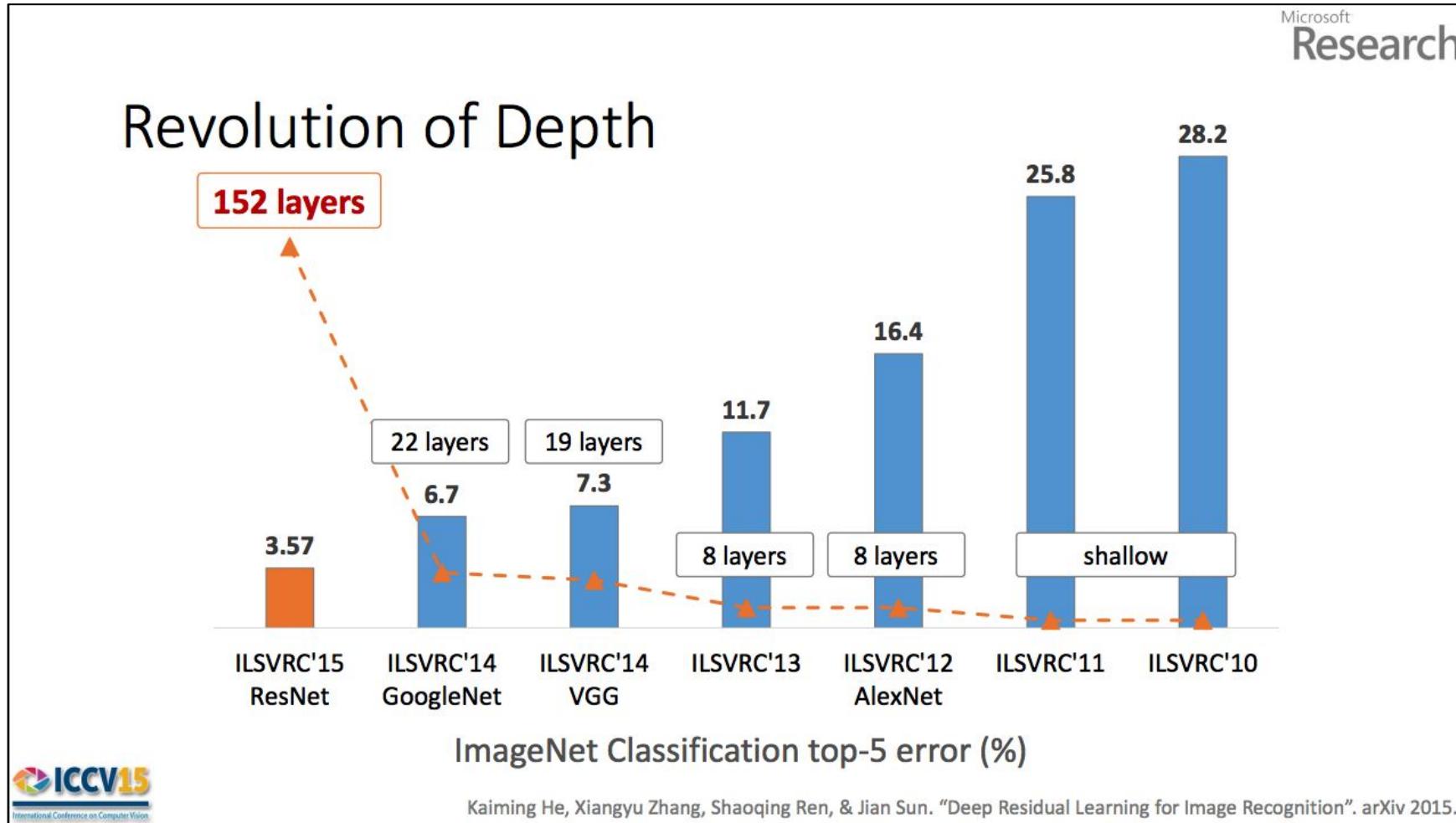


Kaiming He, Xiangyu Zhang, Shaoqing Ren, & Jian Sun. “Deep Residual Learning for Image Recognition”. arXiv 2015.

Sourced with permission from: Deep Residual Learning for Image Recognition, Kaiming He (2015)

# Case Study: ResNet

[He et al., 2015]

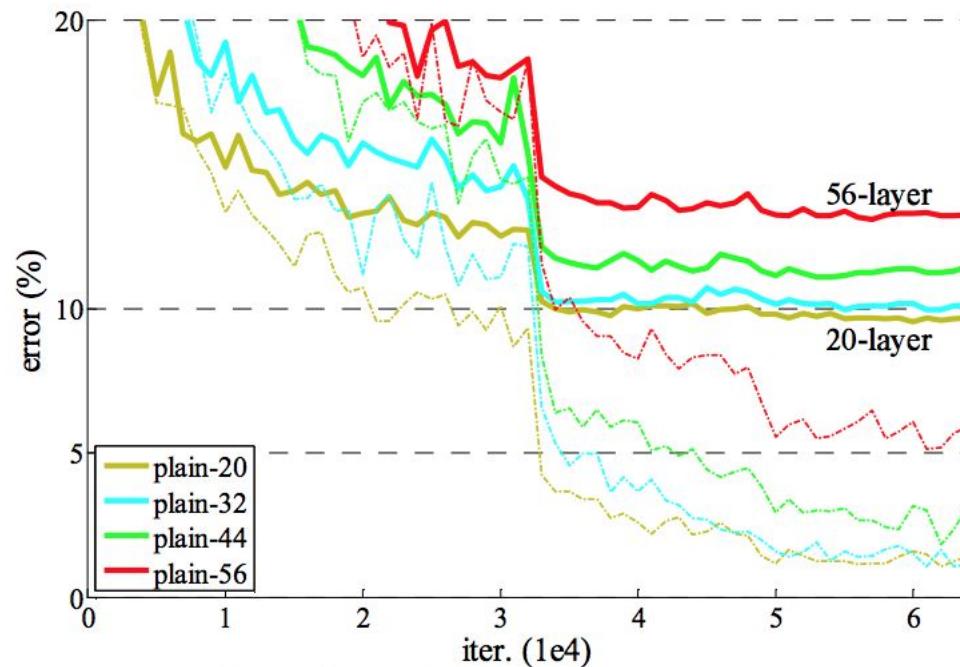


Sourced with permission from: Deep Residual Learning for Image Recognition, Kaiming He (2015)

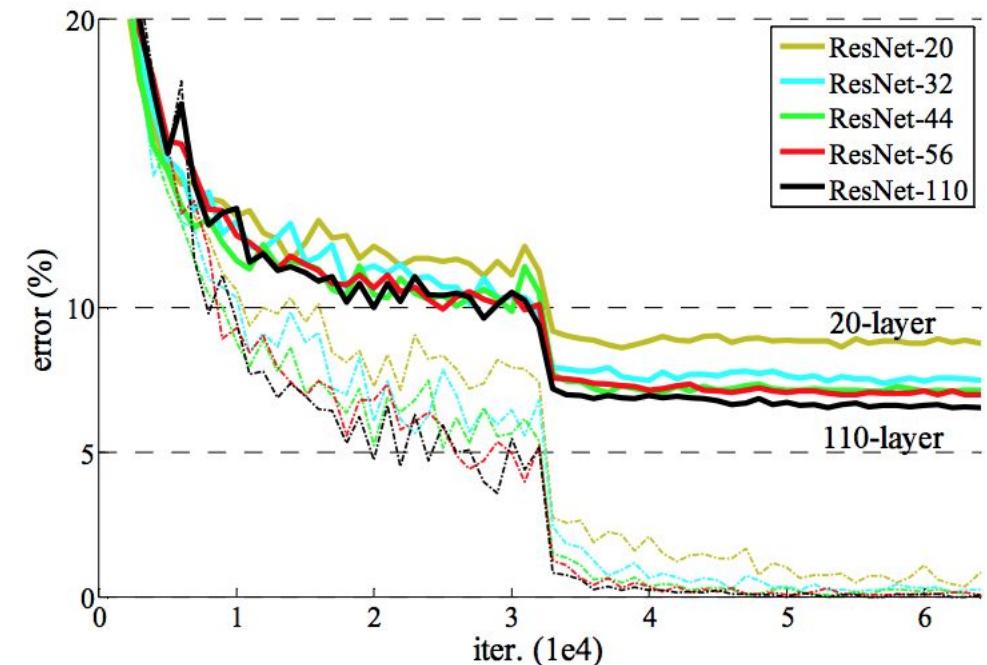
# Case Study: ResNet

[He et al., 2015]

## Experiments on CIFAR-10



**Plain Networks**



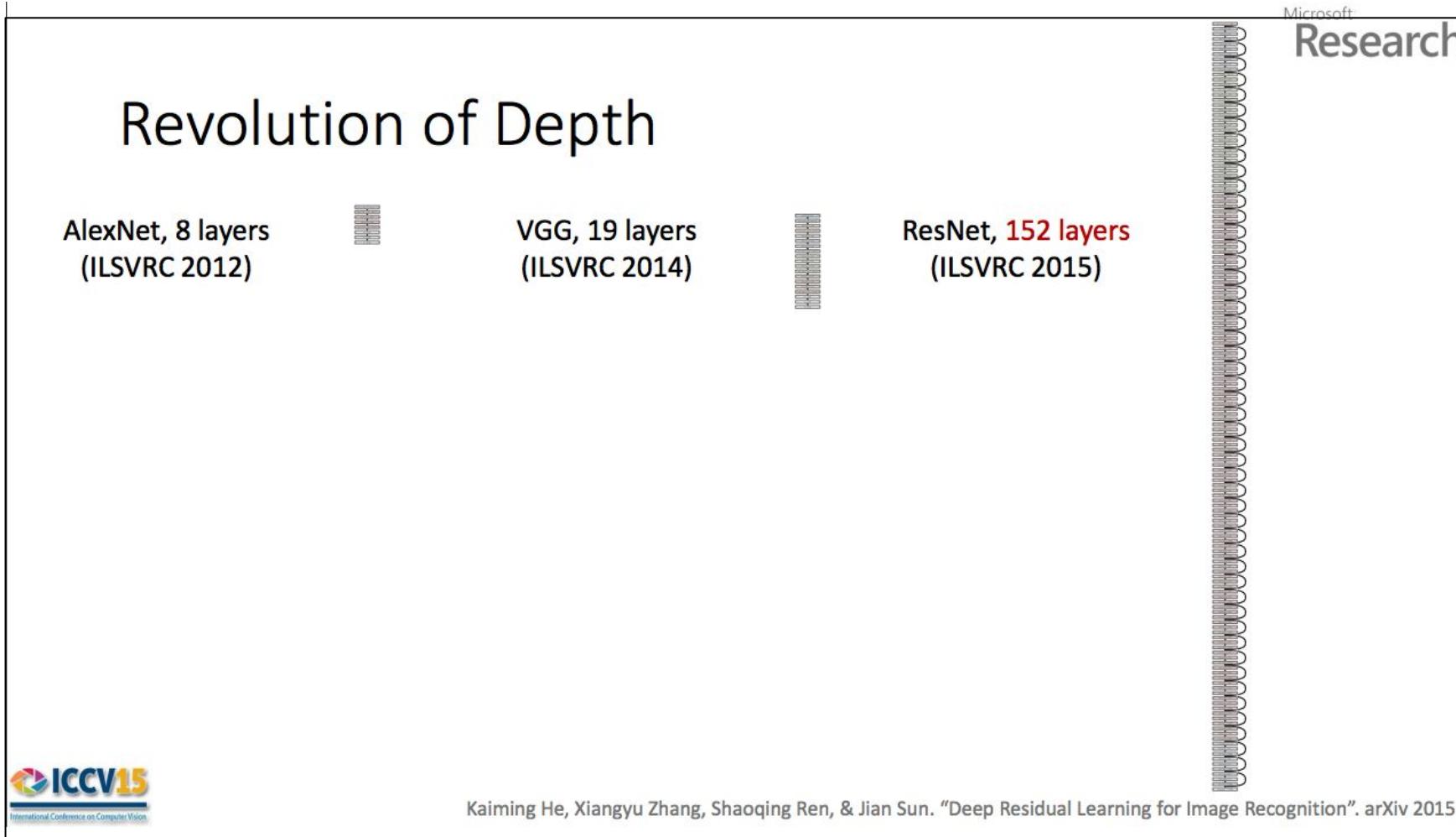
**ResNets**

Sourced with permission from: Deep Residual Learning for Image Recognition, Kaiming He (2015)

# Case Study: ResNet

[He et al., 2015]

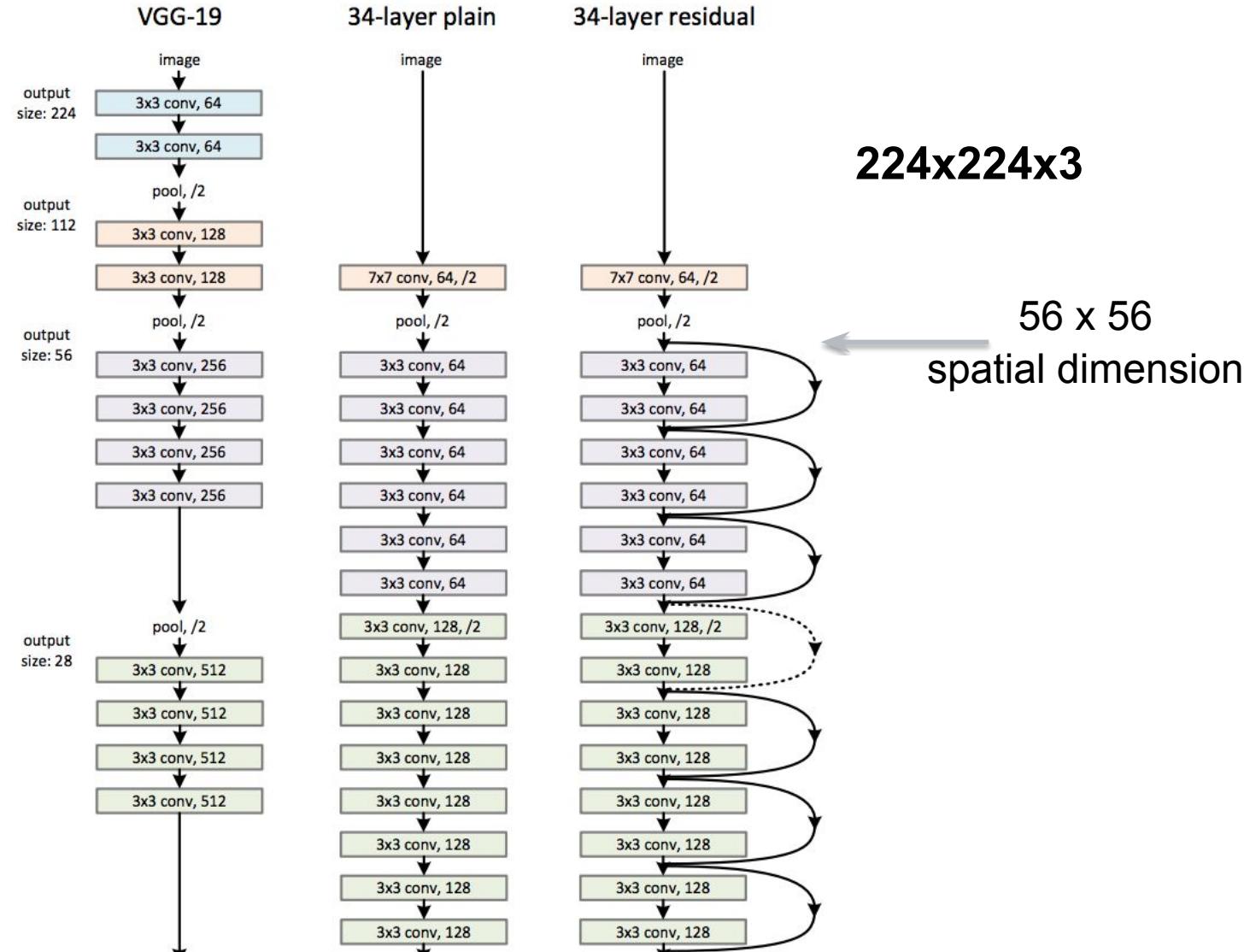
- 2-3 weeks of training on 8 GPU machine
- At runtime: Faster than VGGNet
- (*even with 8x more layers*)



Sourced with permission from: Deep Residual Learning for Image Recognition, Kaiming He (2015)

# Case Study: ResNet

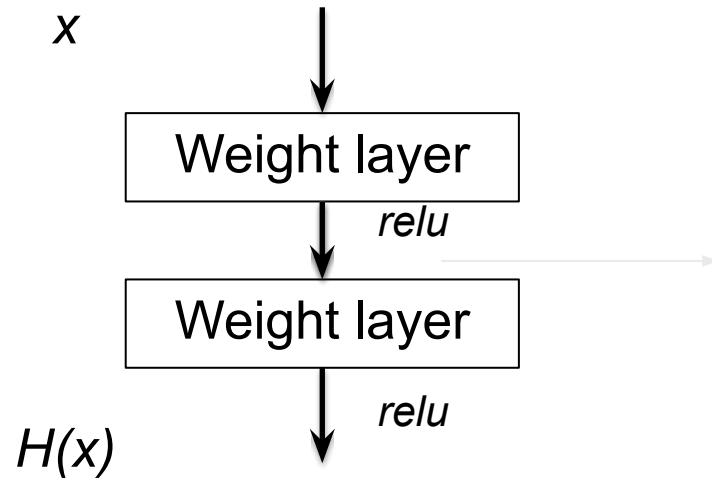
[He et al., 2015]



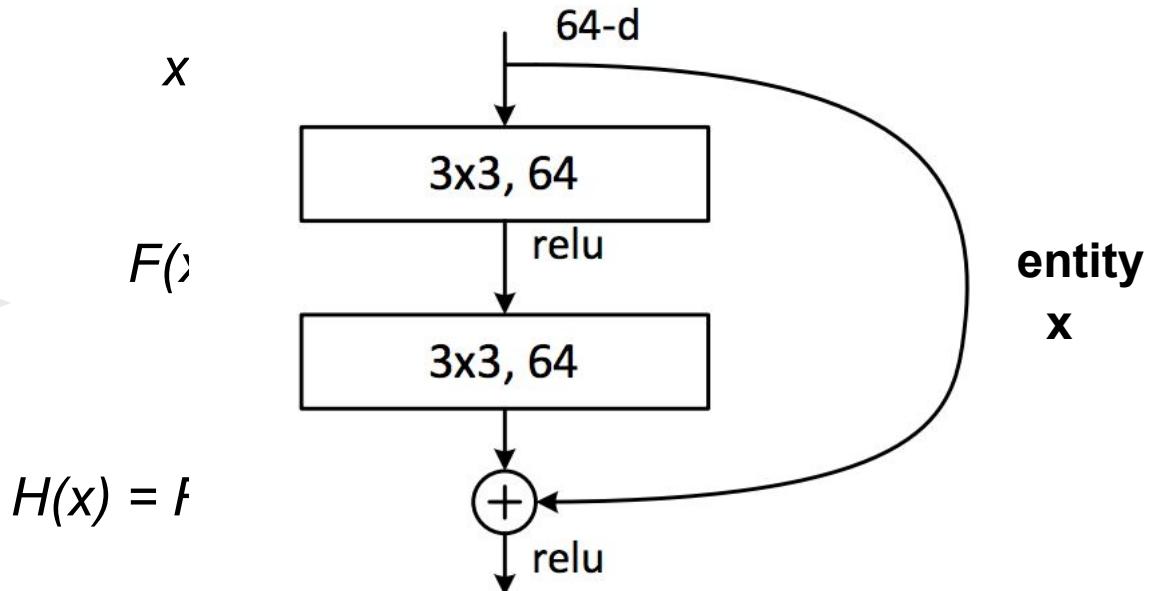
Sourced with permission from: Deep Residual Learning for Image Recognition, Kaiming He (2015)

# Case Study: ResNet

[He et al., 2015]



**Plain Network**



**ResNet**

Sourced with permission from: Deep Residual Learning for Image Recognition, Kaiming He (2015)

# Case Study: ResNet

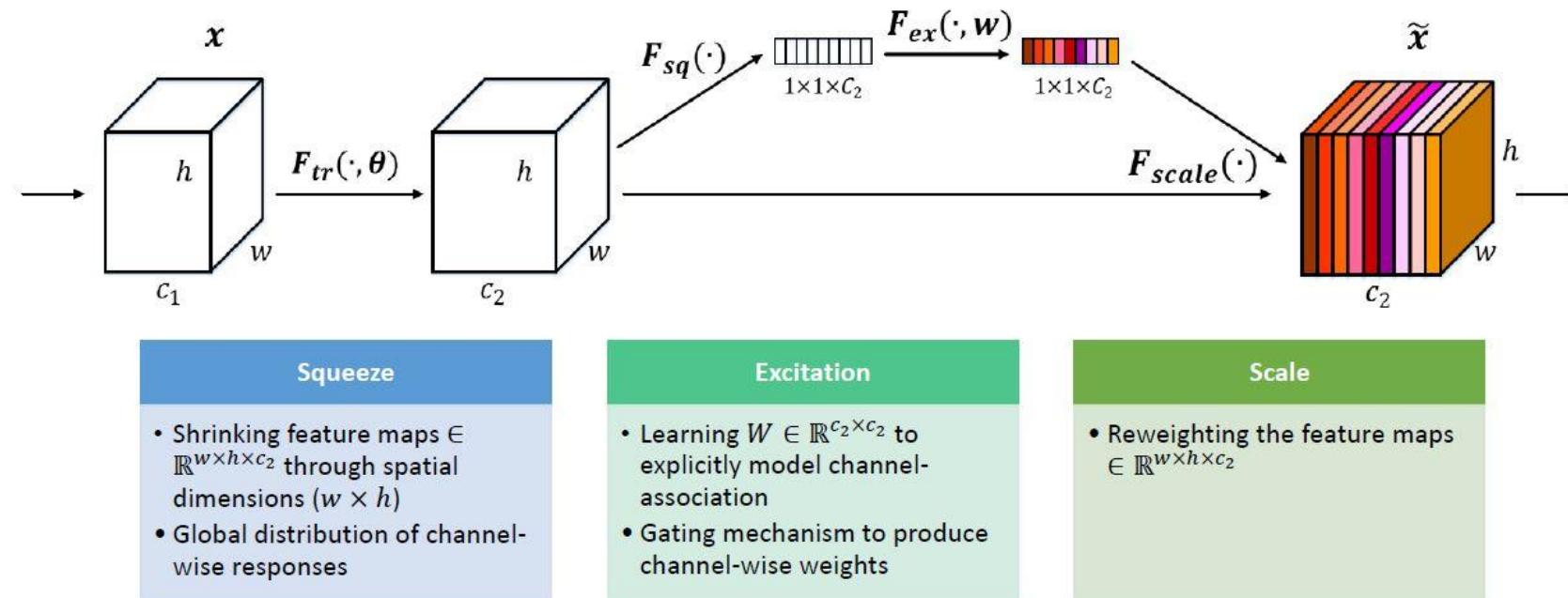
[He et al., 2015]

- 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

- CUIImage was the winner with the ensemble approach.
- Classification error is down to 3.0% from 3.6% last year.
- Pretty boring, best model is just an ensemble
- [https://www.reddit.com/r/MachineLearning/comments/54jiyy/large\\_scale\\_visual\\_recognition\\_challenge\\_2016/](https://www.reddit.com/r/MachineLearning/comments/54jiyy/large_scale_visual_recognition_challenge_2016/)
- <http://image-net.org/challenges/LSVRC/2016/results#loc>

# ILVRC 2017, Squeeze & Excitation Network

- Squeeze and Excitation block that can be added to a Conv Layer
- Add parameters to each channel of a convolutional block so that the network can adaptively adjust the weighting of each feature map.

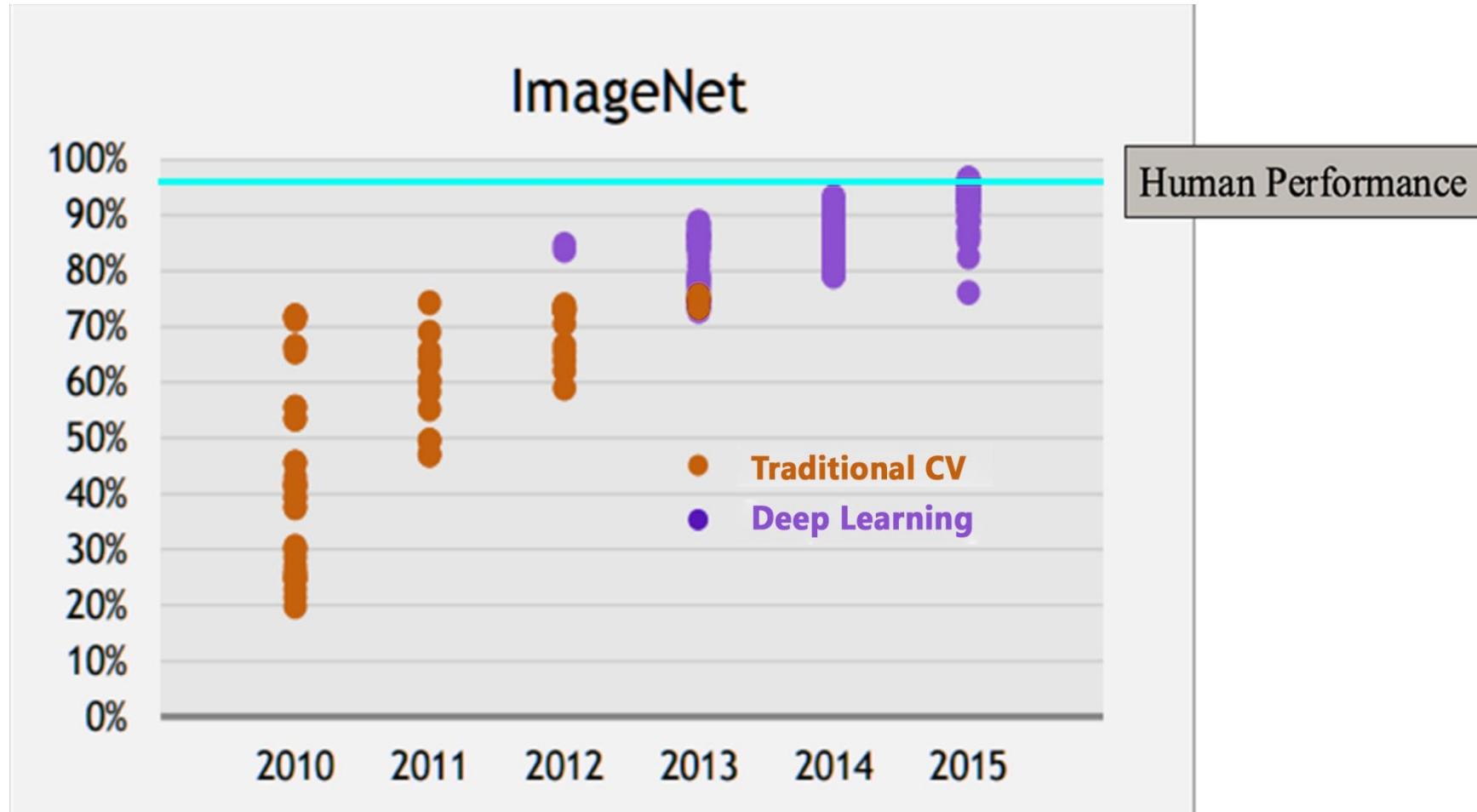


Sourced with permission from: Squeeze and Excitation Networks, ILSVRC 2017 presentation, Jie Hu et al. (2017)

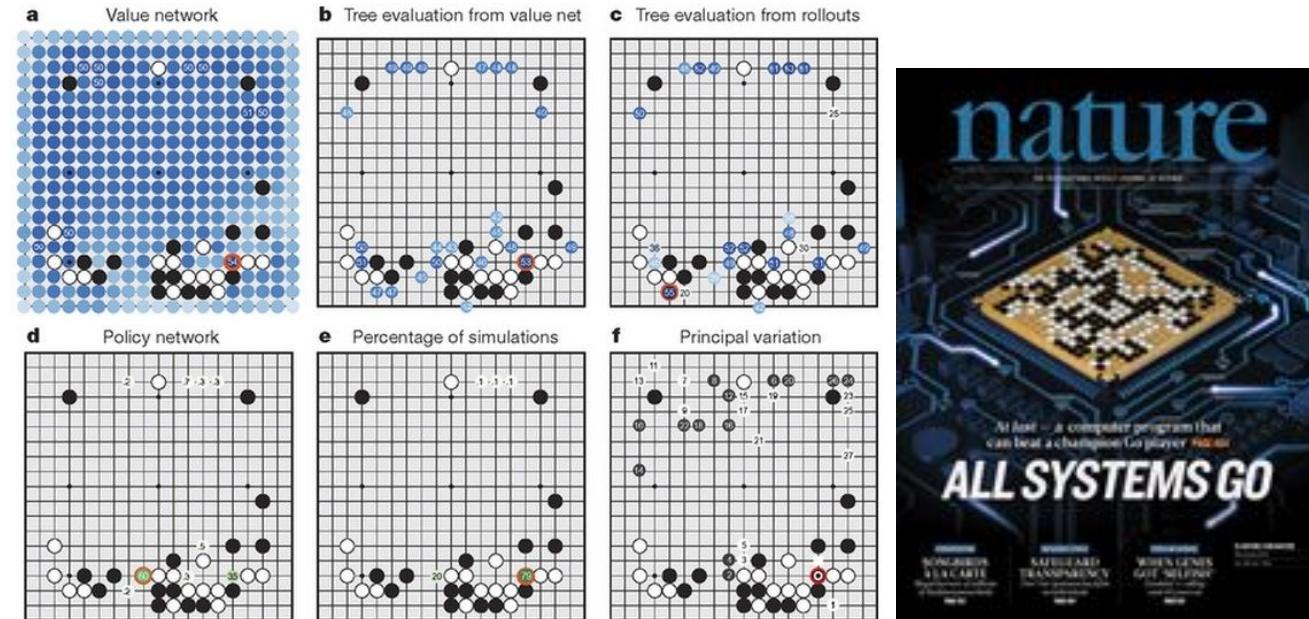
# ILVRC 2017, Squeeze & Excitation Network

- Winning entry comprised a small ensemble of SENets that employed a standard multi-scale and multi-crop fusion strategy
- **2.251%** top-5 error on the test set
- Nearly 25% improvement on the winning entry of 2016 (2.99% top-5 error)
- One of the high-performing networks is constructed by integrating SE blocks with a modified ResNeXt

# Why ConvNets?



# Case Study: DeepMind's AlphaGo



Images Source: 'Mastering the game of Go without human knowledge', Nature, David Silver et al. (2017)

# Case Study: DeepMind's AlphaGo

The input to the policy network is a  $19 \times 19 \times 48$  image stack consisting of 48 feature planes. The first hidden layer zero pads the input into a  $23 \times 23$  image, then convolves  $k$  filters of kernel size  $5 \times 5$  with stride 1 with the input image and applies a rectifier nonlinearity. Each of the subsequent hidden layers 2 to 12 zero pads the respective previous hidden layer into a  $21 \times 21$  image, then convolves  $k$  filters of kernel size  $3 \times 3$  with stride 1, again followed by a rectifier nonlinearity. The final layer convolves 1 filter of kernel size  $1 \times 1$  with stride 1, with a different bias for each position, and applies a softmax function. The match version of AlphaGo used  $k = 192$  filters; Fig. 2b and Extended Data Table 3 additionally show the results of training with  $k = 128, 256$  and  $384$  filters.

## Policy network:

INPUT:	[19x19x48]
CONV1: 192 5x5 filters , stride 1, pad 2	[19x19x192]
CONV2..12: 192 3x3 filters, stride 1, pad 1	[19x19x192]
CONV: 1 1x1 filter, stride 1, pad 0	[19x19] ( <i>probability map of promising moves</i> )

Excerpt Source: 'Mastering the game of Go without human knowledge', Nature, David Silver et al. (2017)

# Summary

- 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:  
 **$[(CONV-RELU)^*N-POOL?]^*M-(FC-RELU)^*K-SOFTMAX$**   
where N is usually up to ~5, M is large,  $0 \leq K \leq 2$ .
- But recent advances such as ResNet/GoogLeNet challenge this paradigm

**“ConvNets need a lot of data to train”?**



This is a myth

## Finetuning

ConvNets usually not trained from scratch

# Data needs for ConvNets

**Train once on massive data like  
ImageNet**



**Fine tune the network using your  
own (much smaller) data**



# Transfer Learning with CNNs



1.  
Train on  
ImageNet



2.  
If you have small  
dataset: fix all weights  
(treat CNN as fixed  
feature extractor),  
retrain only the  
classifier



3.  
If dataset is medium  
sized, “finetune”.  
Use the old weights  
as initialization, train  
the full network or  
only some of the  
higher layers

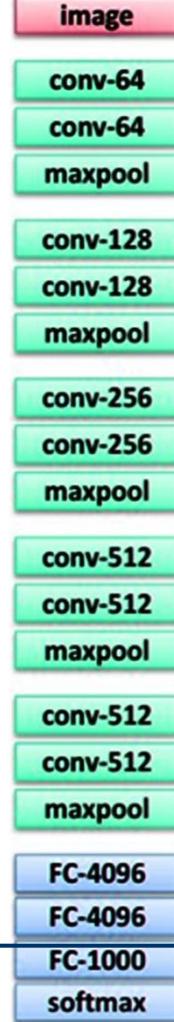
Swap softmax  
layer at end

Retrain bigger  
portion of network

# Transfer Learning with CNNs



**1.**  
Train on  
ImageNet



**2.**  
Small dataset:  
**Feature Extractor**



**3.**  
Medium dataset:  
**Finetune**  
Freeze

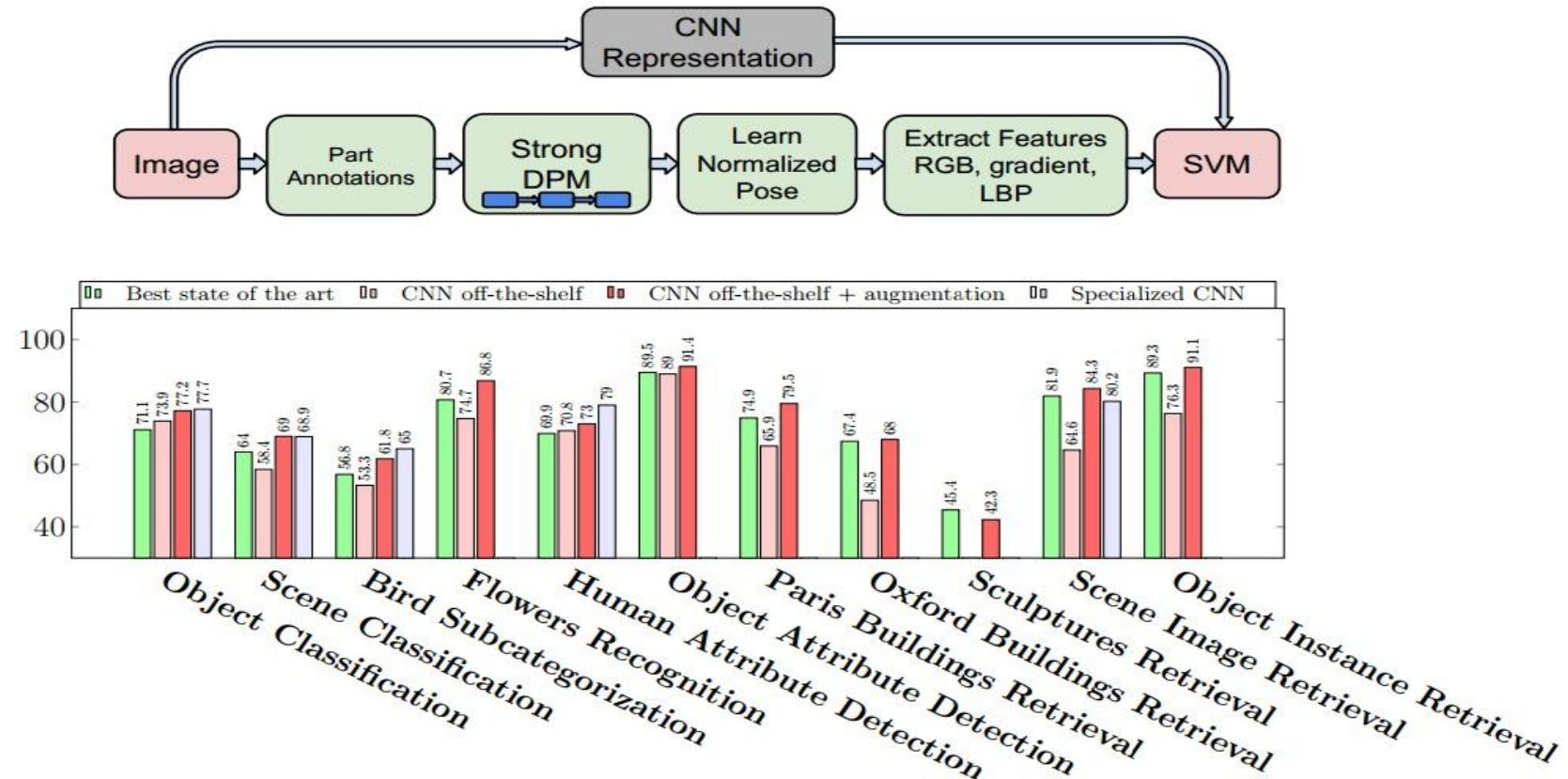
Train

## Rule of thumb:

- Use only ~1/10th of the original learning rate in finetuning top layer
- And ~1/100th in intermediate layers

# CNN Features off-the-shelf

[Razavian et al, 2014]



**“Recent results indicate that the generic descriptors extracted from the convolutional neural networks are very powerful.”**

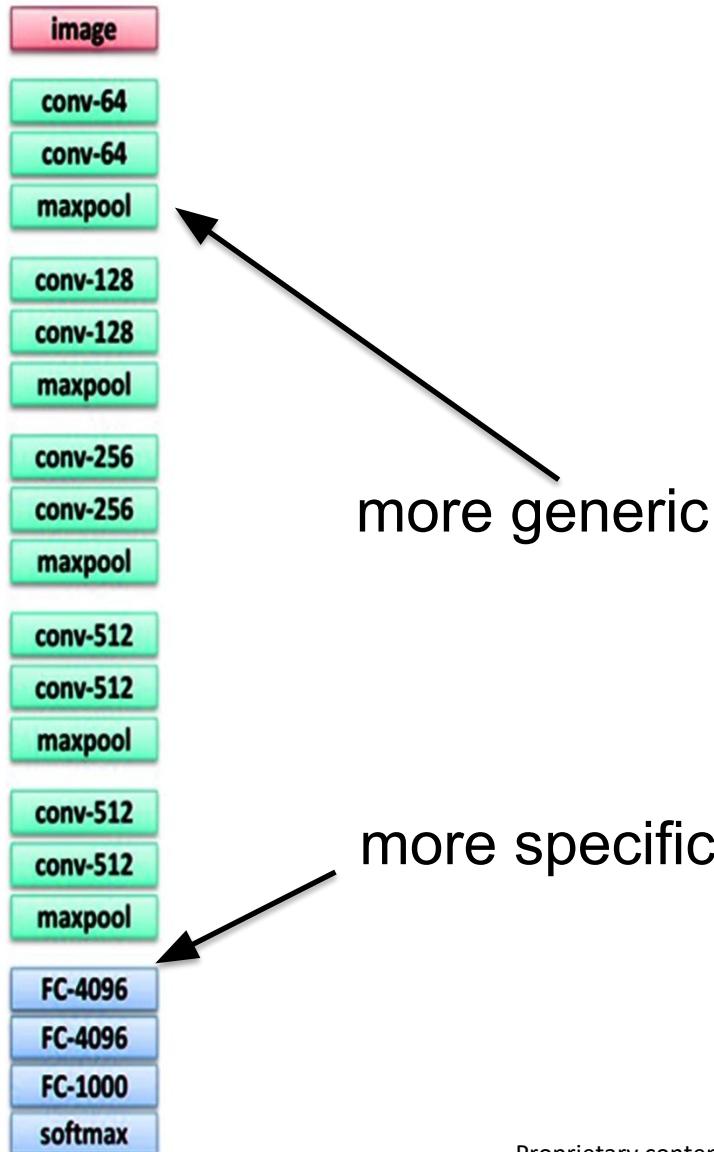
Source: ‘CNN Features off-the-shelf: An Astounding Baseline for Recognition’, Razavian et al. (2014)

# Deep Convolutional Activation for Generic Visual Recognition

[Donahue, Jia et al., 2013]

	DeCAF <sub>6</sub>	DeCAF <sub>7</sub>
LogReg	<b>40.94 ± 0.3</b>	40.84 ± 0.3
SVM	39.36 ± 0.3	40.66 ± 0.3
Xiao et al. (2010)		38.0

Source: 'DeCAF: A Deep Convolutional Activation Feature for Generic Visual Recognition, Donahue, Jia, et al., (2013)



**very little data**

**very similar dataset**

?

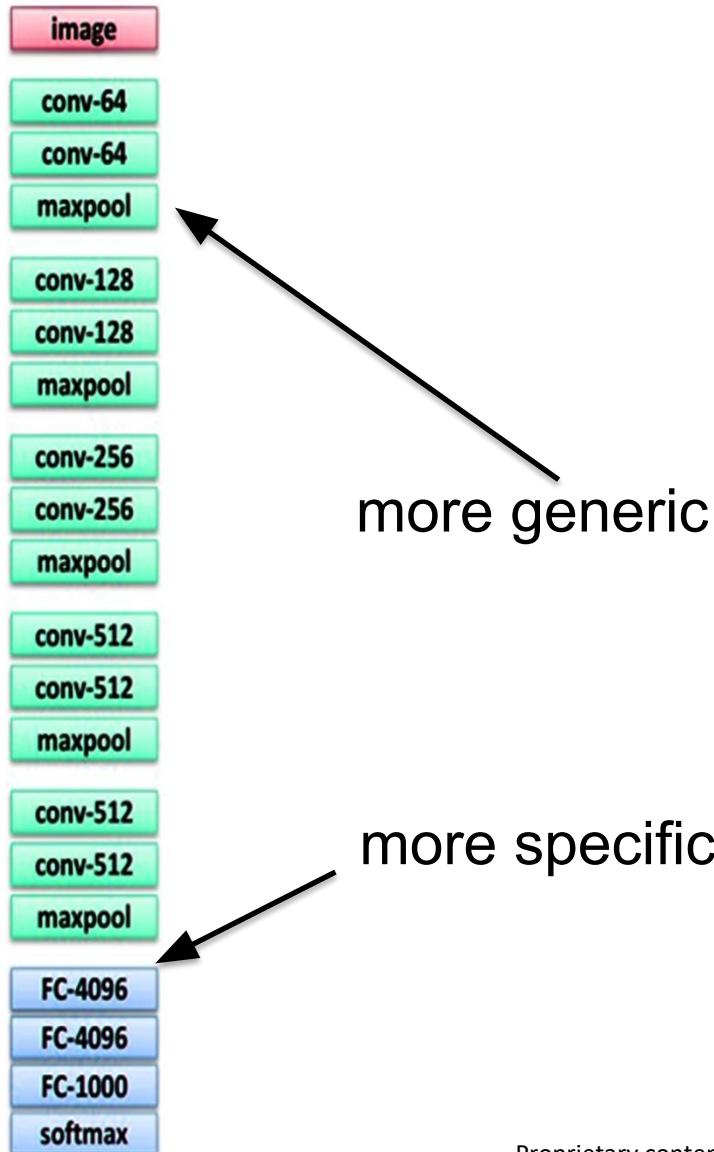
**very different dataset**

?

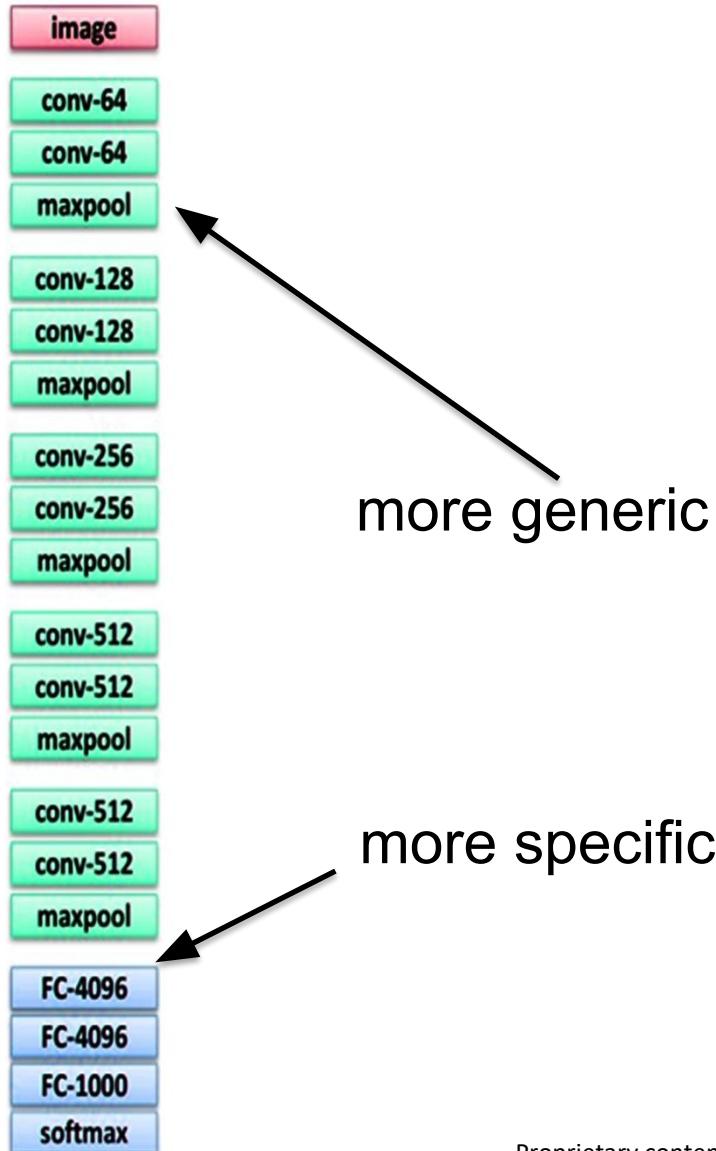
**quite a lot of data**

?

?



	<b>very similar dataset</b>	<b>very different dataset</b>
<b>very little data</b>	On the Top layer use Linear Classifier	?
<b>quite a lot of data</b>	Finetuning of few layers	?



**very little data**

**very similar dataset**

**very different dataset**

On the Top layer  
use Linear  
Classifier

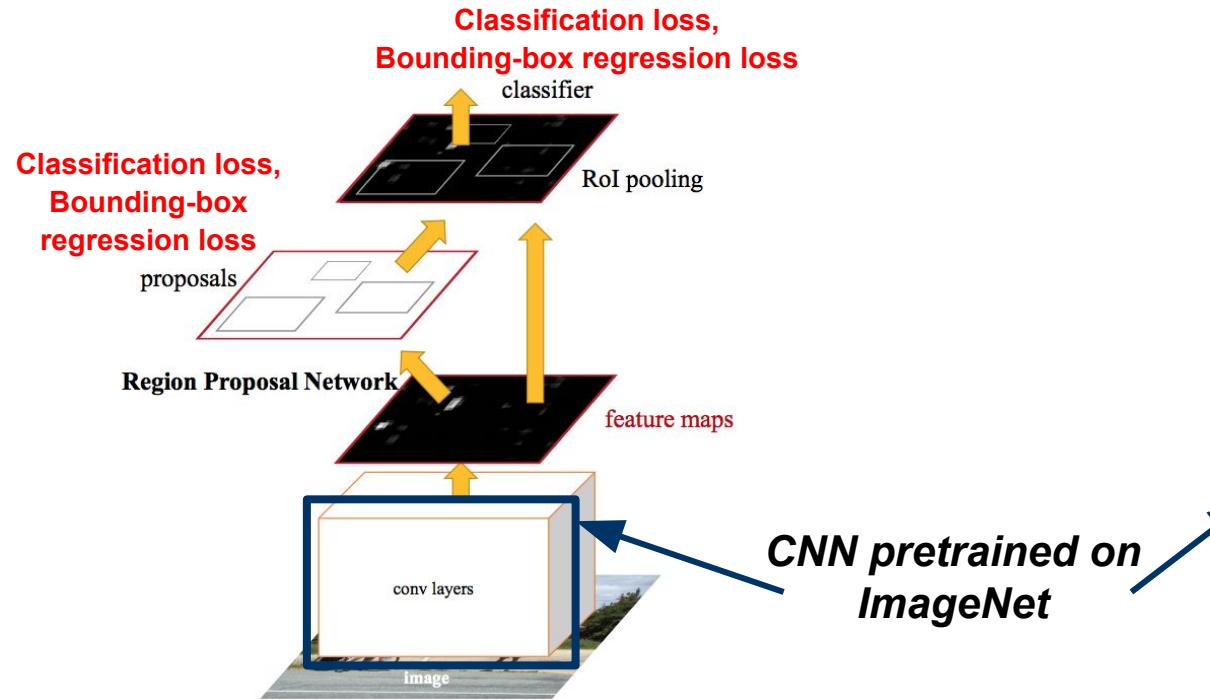
Try linear  
classifier from  
different stages

**quite a lot of  
data**

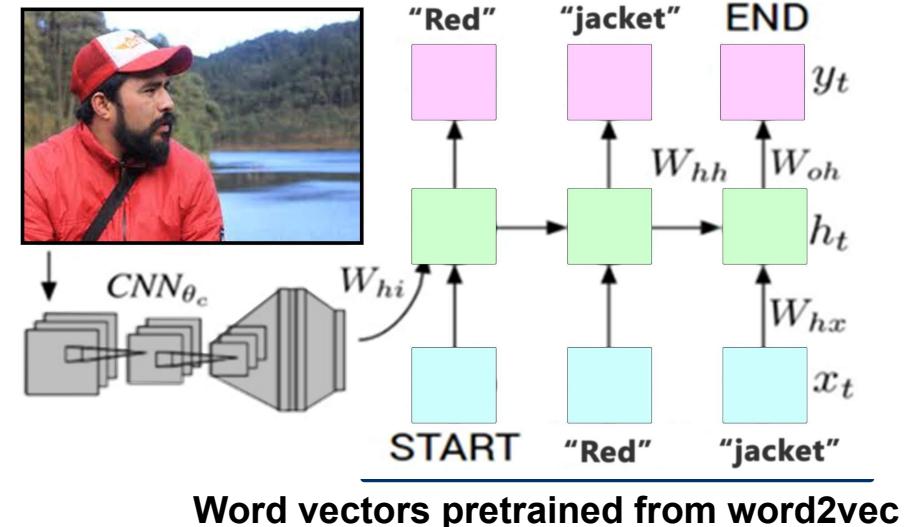
Finetuning of  
few layers

Finetune a  
larger number  
of layers

# Transfer learning with CNNs is common



**Object Detection**  
**Faster R-CNN**



**The Image Captioning problem**  
**CNN + RNN**

Sources: *Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks*, Ren, He et al. (2016)

# E.g. Caffe Model Zoo: Lots of pretrained ConvNets

<https://github.com/BVLC/caffe/wiki/Model-Zoo>

<https://github.com/szagoruyko/loadcaffe>

**Model Zoo**  
ELM edited this page 21 days ago - 56 revisions

Check out the model zoo documentation for details.

To acquire a model:

1. download the model gist by `./scripts/download_model_from_gist.sh <gist_id>`  
`<dirname>` to load the model metadata, architecture, solver configuration, and so on.  
(`<dirname>` is optional and defaults to `caffe/models`).
2. download the model weights by `./scripts/download_model_binary.py <model_dir>` where  
`<model_dir>` is the gist directory from the first step.

or visit the [model zoo documentation](#) for complete instructions.

**Berkeley-trained models**

- Finetuning on Flickr Style: same as provided in `models/`, but listed here as a Gist for an example.
- BVLC GoogLeNet: `models/bvlc_googlenet`.

**Network in Network model**

The Network In Network model is described in the following ICLR-2014 paper:

Network In Network  
M. Lin, Q. Chen, S. Yan  
International Conference on Learning Representations, 2014 (arXiv:1409.1556)

please cite the paper if you use the models.

Models:

- NIN-Imagenet: a small(29MB) model for Imagenet, yet performs slightly better than AlexNet, and fast to train. (Note: a more caffe-compatible version with correct convolutional weights shape: <https://drive.google.com/folderview?usp=driveweb>)
- NIN-CIFAR10: NIN model on CIFAR10, originally published in the paper Network In Network. The error rate of this model is 10.4% on CIFAR10.

**Models from the BMVC-2014 paper "Return of the Devil in the Details: Delving Deep into Convolutional Nets"**

The models are trained on the ILSVRC-2012 dataset. The details can be found on the [project page](#) or in the following BMVC-2014 paper:

Return of the Devil in the Details: Delving Deep into Convolutional Nets  
K. Chatfield, K. Simonyan, A. Vedaldi, A. Zisserman  
British Machine Vision Conference, 2014 (arXiv ref. cs.IA85.3531)

Please cite the paper if you use the models.

Models:

- VGG\_CNN\_S: 13.1% top-5 error on ILSVRC-2012-val
- VGG\_CNN\_M: 13.7% top-5 error on ILSVRC-2012-val
- VGG\_CNN\_M\_2048: 13.5% top-5 error on ILSVRC-2012-val
- VGG\_CNN\_M\_1024: 13.7% top-5 error on ILSVRC-2012-val
- VGG\_CNN\_M\_128: 15.6% top-5 error on ILSVRC-2012-val
- VGG\_CNN\_F: 16.7% top-5 error on ILSVRC-2012-val

**Models used by the VGG team in ILSVRC-2014**

**Places-CNN model from MIT.**  
Places CNN is described in the following NIPS 2014 paper:  
B. Zhou, A. Lapedriza, J. Xiao, A. Torralba, and A. Oliva  
Learning Deep Features for Scene Recognition using Places Database.  
Advances in Neural Information Processing Systems 27 (NIPS) spotlight, 2014.

The project page is [here](#)

Models:

- Places205-AlexNet: CNN trained on 205 scene categories of Places Database (used in NIPS14) with >2 million images. The architecture is the same as Caffe reference network.
- Hybrid-CNN: CNN trained on 1103 categories (205 scene categories from Places Database and 978 object categories from the train data of ILSVRC2012 (ImageNet) with >3.6 million images. The architecture is the same as Caffe reference network.
- Places205-GoogLeNet: GoogLeNet CNN trained on 205 scene categories of Places Database. It is used by Google in the deep dream visualization

**GoogLeNet GPU Implementation from Princeton.**  
We implemented GoogLeNet using a single GPU. Our main contribution is an effective way to initialize the network and a trick to overcome the GPU memory constraint by accumulating gradients over two training iterations.

- Please check <http://vision.princeton.edu/~li/GoogLeNet/> for more information. Pre-trained models on ImageNet and Places, and the training code are available for download.
- Make sure `cifc_to2` and `cifc3_fc` have `num_output = 1000` in the prototxt. Otherwise, the trained model would crash on test.

**Fully Convolutional Semantic Segmentation Models (FCN-Xs)**  
These models are described in the paper:

Fully Convolutional Models for Semantic Segmentation  
Jonathan Long, Evan Shelhamer, Trevor Darrell  
CVPR 2015  
arXiv:1411.4938

They are available under the same license as the Caffe-bundled models (i.e., for unrestricted use; see [http://caffe.berkeleyvision.org/model\\_zoo.html#v1-model-license](http://caffe.berkeleyvision.org/model_zoo.html#v1-model-license)).

These are pre-release models. They do not run in any current version of Caffe/caffe, as they require unmerged PRs. They should run in the preview branch provided at <https://github.com/longcw/caffe-fcn/tree/feature>. The FCN-32s PASCAL-Context model is the most complete example including network definitions, solver configuration, and Python scripts for solving and inference.

Models trained on PASCAL (using extra data from Harham et al. and finetuned from the ILSVRC-trained VGG-16 model):

- FCN-32s PASCAL: single stream, 32 pixel prediction stride version
- FCN-16s PASCAL: two stream, 16 pixel prediction stride version
- FCN-8s PASCAL: three stream, 8 pixel prediction stride version
- FCN-AlexNet PASCAL: AlexNet (CaffeNet) single stream, 32 pixel prediction stride version

To reproduce the validation scores, use the `seg11val` split defined by the paper in footnote 7. Since SBD and PASCAL VOC 11 segval intersect, we evaluate on the non-intersecting set for validation purposes.

Models trained on SIFT Flow (also finetuned from VGG-16):

- FCN-16s SIFT Flow: two stream, 16 pixel prediction stride version

Models trained on NYUDv2 (also finetuned from VGG-16, and using HHA features from Gupta et al. at <https://github.com/gupta/cnn-depth/>):

- FCN-32s NYUDv2: single stream, 32 pixel prediction stride version
- FCN-16s NYUDv2: two stream, 16 pixel prediction stride version

Models trained on PASCAL-Context including training model definition, solver configuration, and barebones solving script (finetuned from the ILSVRC-trained VGG-16 model):

- FCN-32s PASCAL-Context: single stream, 32 pixel prediction stride version
- FCN-16s PASCAL-Context: two stream, 16 pixel prediction stride version
- FCN-8s PASCAL-Context: three stream, 8 pixel prediction stride version

**CaffeNet fine-tuned for Oxford flowers dataset**  
<https://gist.github.com/limjoco/179e2305ca70bad01f>

This is the reference CaffeNet (modified AlexNet) fine-tuned for the Oxford 102 category flower dataset. The number of output in the inner product layer has been set to 102 to reflect the number of flower categories. Hyperparameter choices reflect those in Fine-tuning CaffeNet for Style Recognition on "Flickr Style". The global learning rate is reduced while the learning rate for the fully connected is increased relative to the other layers.

After 50,000 iterations, the top-1 error is 7% on the test set of 1,020 images.

**CNN Models for Salient Object Subtizing.**  
CNN models described in the following CVPR'15 paper "Salient Object Subtizing":

Salient Object Subtizing  
J. Zhang, S. Ma, M. Sameki, S. Sclaroff, M. Betke, Z. Lin, X. Shen, B. Price and R. CVPR, 2015.

Models:

- AlexNet: CNN model finetuned on the Salient Object Subtizing dataset (~5000 images). The architecture is the same as the Caffe reference network.
- VGG16: CNN model finetuned on the Salient Object Subtizing dataset (~5000 images). The architecture is the same as the VGG16 network. This model gives better performance than the AlexNet model, but is slower for training and testing.

**Deep Learning of Binary Hash Codes for Fast Image Retrieval**  
We present an effective deep learning framework to create the hash-like binary codes for fast image retrieval. The details can be found in the following "CVPRW'15 paper":

Deep Learning of Binary Hash Codes for Fast Image Retrieval  
K. Lin, H.-F. Yang, J.-W. Hsieh, C.-S. Chen  
CVPR 2015, DeepVision workshop

please cite the paper if you use the model:

- caffe-cvpr15: See our code release on Github, which allows you to train your own deep hashing model and create binary hash codes.
- CIFAR10-48bit: Proposed 48-bits CNN model trained on CIFAR10.

**Holistically-Nested Edge Detection**  
The model and code provided are described in the ICCV 2015 paper:

Holistically-Nested Edge Detection  
Seining Xie and Zhuowen Tu  
ICCV 2015

For details about training/evaluating HED, please take a look at <http://github.com/sxie/hed>.

Model trained on BSDS-500 Dataset (finetuned from the VGGNet):

- HED BSDS-500

**Translating Videos to Natural Language**  
These models are described in this NAACL-HLT 2015 paper:

Translating Videos to Natural Language Using Deep Recurrent Neural Networks  
S. Venugopalan, H. Xu, J. Donahue, M. Rohrbach, K. Mooney, R. Saenko  
NAACL-HLT 2015

More detail can be found on this project page.

Model:

Video2Text\_VGG\_mean\_pool: This model is an improved version of the mean pooled model described in the NAACL-HLT 2015 paper. It uses video frame features from the VGG-16 layer model. This is trained only on the YouTube video dataset.

Compatibility: These are pre-release models. They do not run in any current version of BVLC/caffe, as they require unmerged PRs. The models are currently supported by the `recurrent` branch of the Caffe fork provided at <https://github.com/jdonahue/caffe/tree/recurrent> and <https://github.com/vishashri/caffe/tree/recurrent>.

**VGG Face CNN descriptor**  
These models are described in this BMVC 2015 paper:

Deep Face Recognition  
Omkar M. Parkhi, Andrea Vedaldi, Andrew Zisserman  
BMVC 2015

More details can be found on this project page.

Model: VGG Face: This is the very deep architecture based model trained from scratch using 2.6 Million images of celebrities collected from the web. The model has been imported to work with Face from the original model trained using MatConvNet library.

If you find our models useful, please add suitable reference to our paper in your work.

**Yearbook Photo Dating**  
Model from the ICCV 2015 Extreme Imaging Workshop paper:

A Century of Portraits: Exploring the Visual Historical Record of American High School  
Shiry Ginosar, Kate Rakely, Brian Yin, Sarah Sachs, Alyseh Efros  
ICCV Workshop 2015

Model and protobuff files: Yearbook

**CCNN: Constrained Convolutional Neural Networks for Weakly Supervised Segmentation**  
These models are described in the ICCV 2015 paper:

Constrained Convolutional Neural Networks for Weakly Supervised Segmentation  
Deepend Pathak, Philipp Krähenbühl, Trevor Darrell  
ICCV 2015  
arXiv:1506.03048

These are pre-release models. They do not run in any current version of BVLC/caffe, as they require unmerged PRs. Full details, source code, models, prototxts are available here: CCNN.

Thank you!