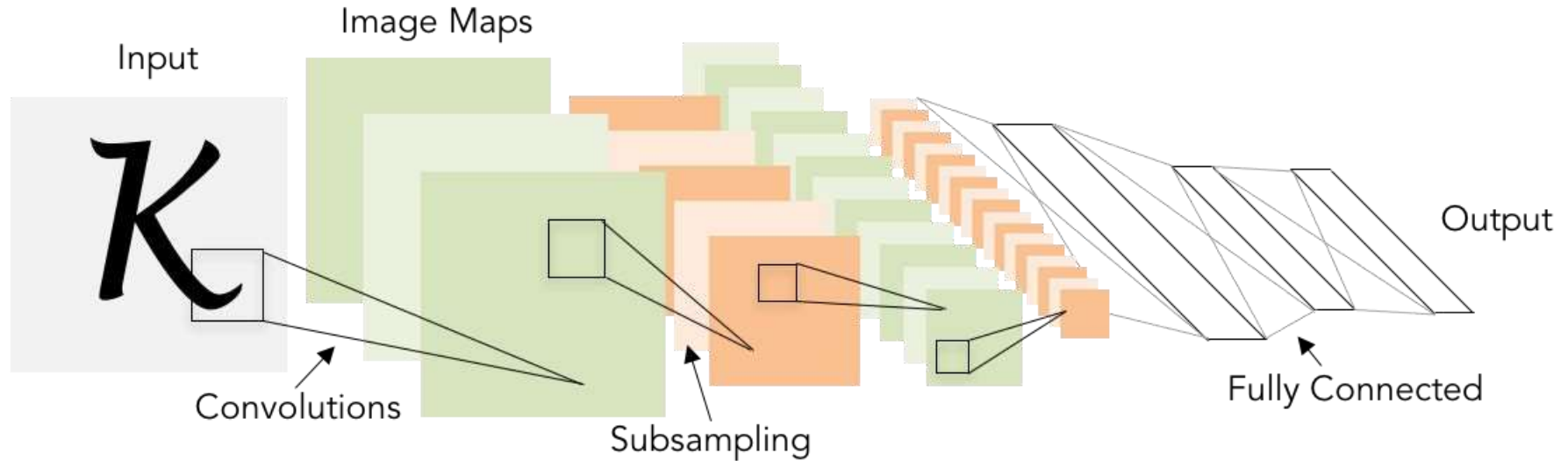# CNN Architectures

# Today: CNN Architectures

## Case Studies

- AlexNet
- VGG
- GoogLeNet
- ResNet

# Review: 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]

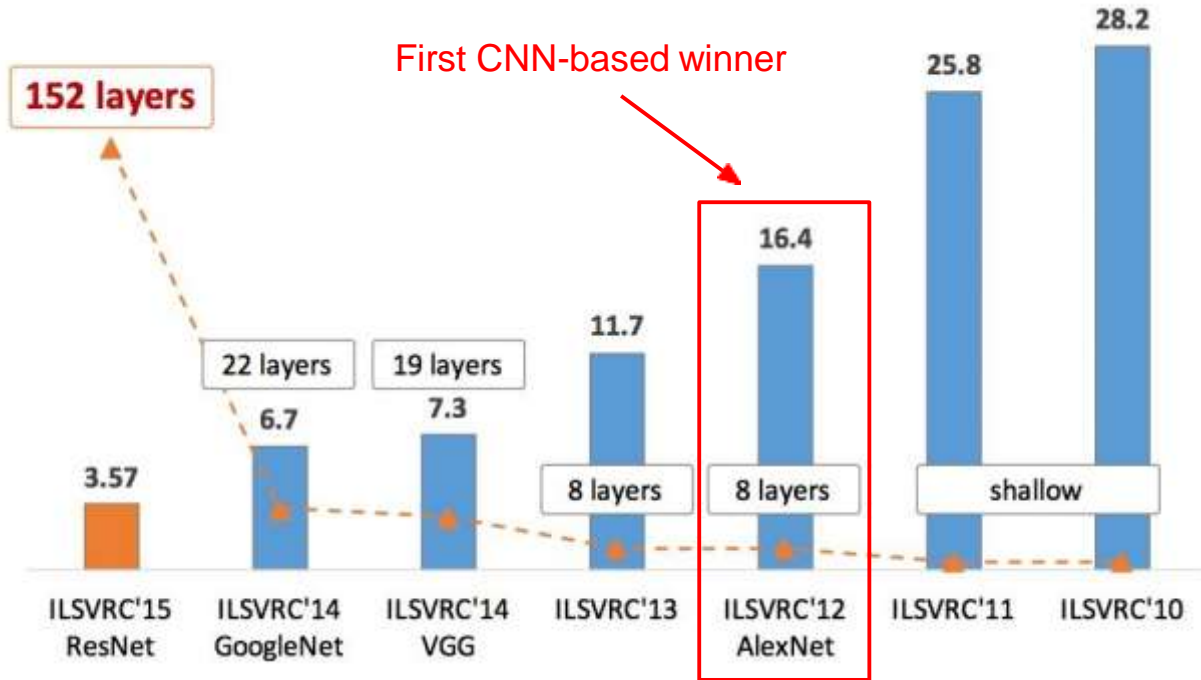# ImageNet Large Scale Visual Recognition Challenge (ILSVRC) winners



First CNN-based winner

Figure copyright Kaiming He, 2016. Reproduced with permission.

# Case Study: AlexNet

*[Krizhevsky et al. 2012]*



Full (simplified) AlexNet architecture:
[227x227x3] INPUT
[55x55x96] CONV1: 96 11x11 filters at stride 4, pad 0
[27x27x96] MAX POOL1: 3x3 filters at stride 2
[27x27x96] NORM1: Normalization layer
[27x27x256] CONV2: 256 5x5 filters at stride 1, pad 2
[13x13x256] MAX POOL2: 3x3 filters at stride 2
[13x13x256] NORM2: Normalization layer
[13x13x384] CONV3: 384 3x3 filters at stride 1, pad 1
[13x13x384] CONV4: 384 3x3 filters at stride 1, pad 1
[13x13x256] CONV5: 256 3x3 filters at stride 1, pad 1
[6x6x256] MAX POOL3: 3x3 filters at stride 2
[4096] FC6: 4096 neurons
[4096] FC7: 4096 neurons
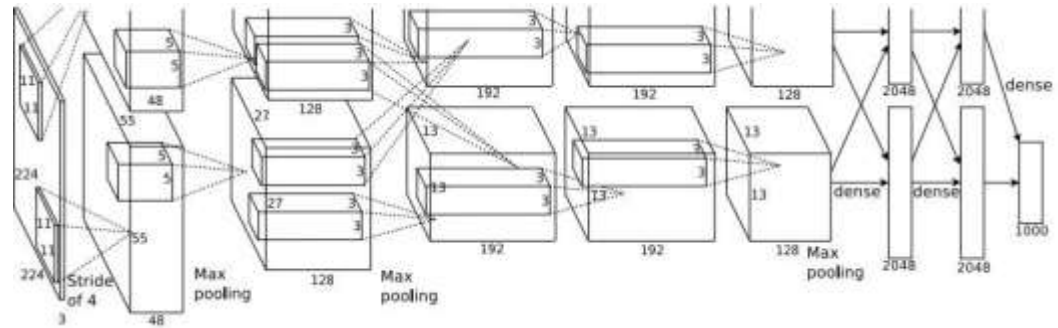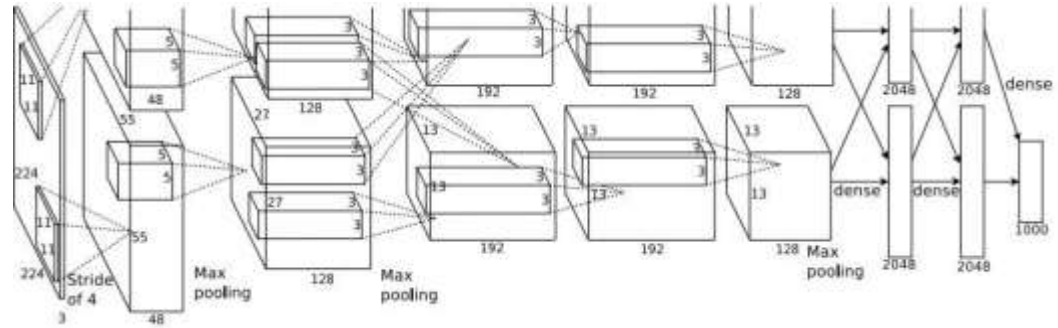[1000] FC8: 1000 neurons (class scores)

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

# Case Study: AlexNet

*[Krizhevsky et al. 2012]*



Full (simplified) AlexNet architecture:
[227x227x3] INPUT
[55x55x96] CONV1: 96 11x11 filters at stride 4, pad 0
[27x27x96] MAX POOL1: 3x3 filters at stride 2
[27x27x96] NORM1: Normalization layer
[27x27x256] CONV2: 256 5x5 filters at stride 1, pad 2
[13x13x256] MAX POOL2: 3x3 filters at stride 2
[13x13x256] NORM2: Normalization layer
[13x13x384] CONV3: 384 3x3 filters at stride 1, pad 1
[13x13x384] CONV4: 384 3x3 filters at stride 1, pad 1
[13x13x256] CONV5: 256 3x3 filters at stride 1, pad 1
[6x6x256] MAX POOL3: 3x3 filters at stride 2
[4096] FC6: 4096 neurons
[4096] FC7: 4096 neurons
[1000] FC8: 1000 neurons (class scores)

**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%

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

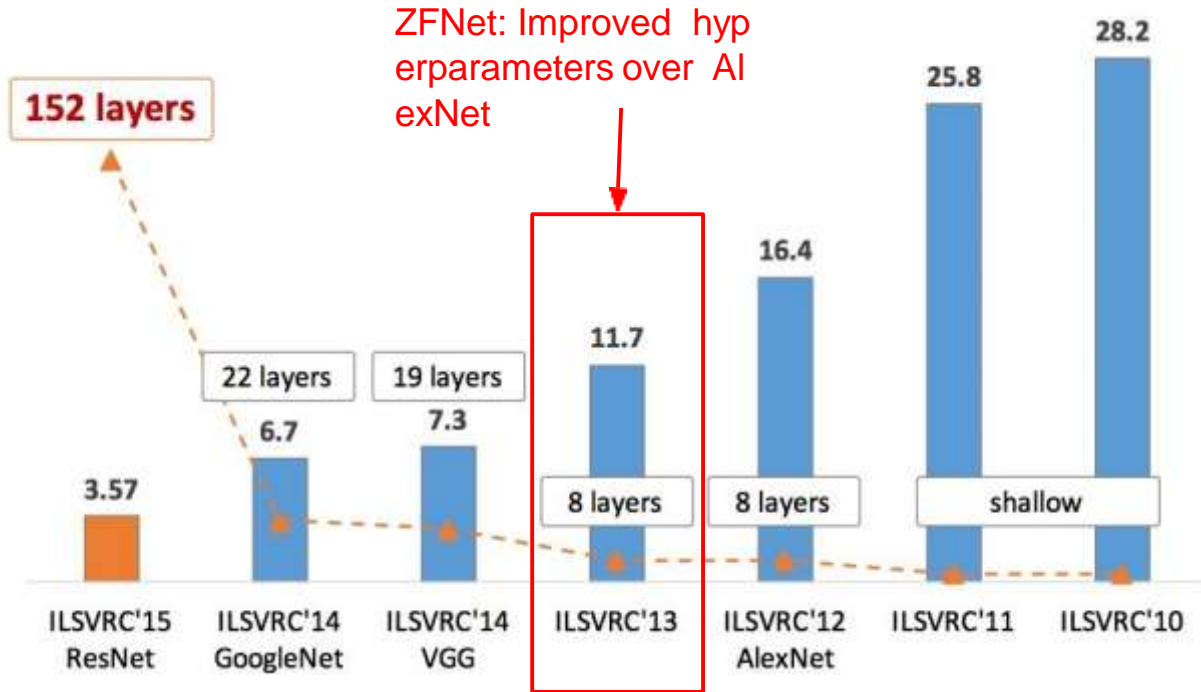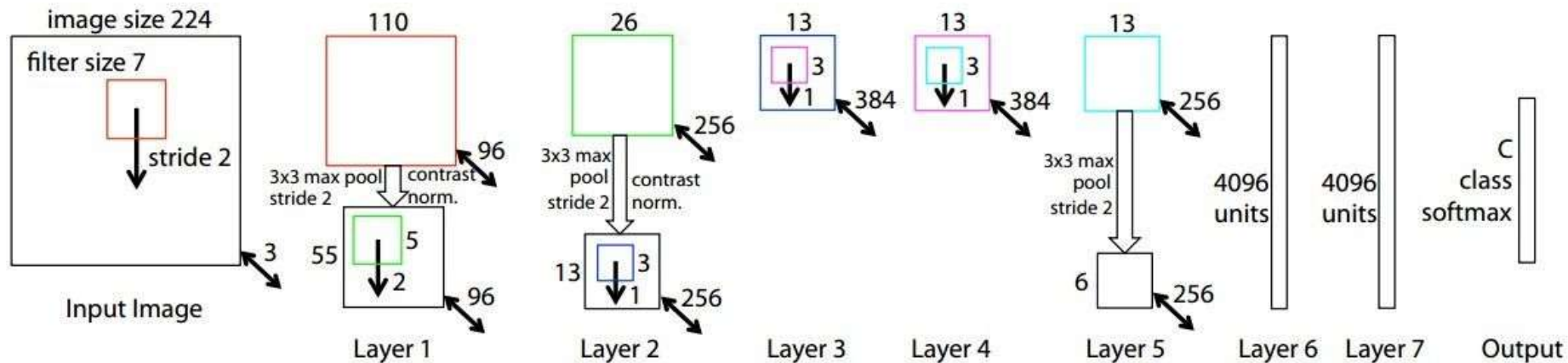# ImageNet Large Scale Visual Recognition Challenge (ILSVRC) winners



Figure copyright Kaiming He, 2016. Reproduced with permission.

# ZFNet

TODO: remake figure

AlexNet but:
CONV1: change from (11x11 stride 4) to (7x7 stride 2)
CONV3,4,5: instead of 384, 384, 256 filters use 512, 1024, 512

ImageNet top 5 error: 16.4% -> 11.7%

8

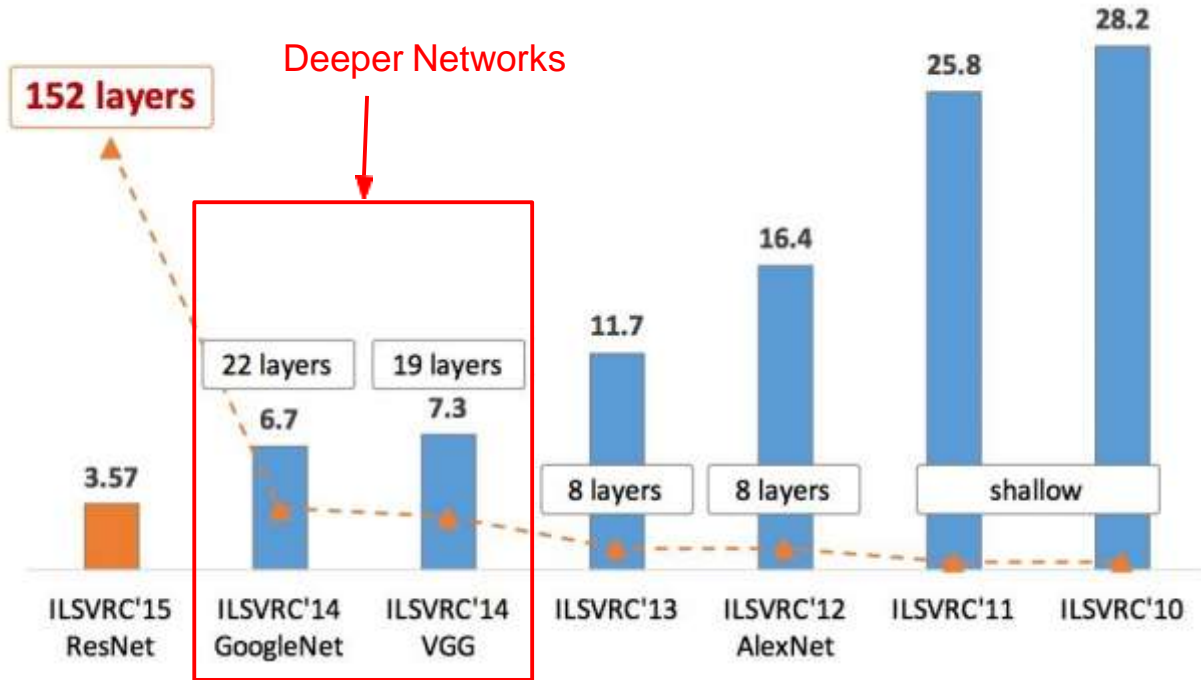# ImageNet Large Scale Visual Recognition Challenge (ILSVRC) winners



Figure copyright Kaiming He, 2016. Reproduced with permission.

# Case Study: VGGNet

*[Simonyan and Zisserman, 2014]*

Small filters, Deeper networks
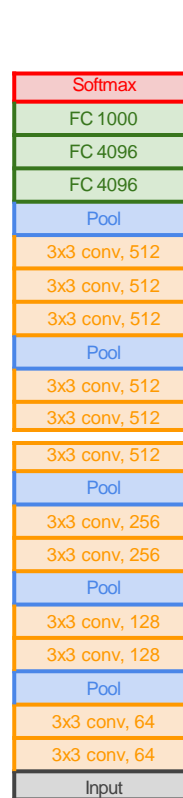
8 layers (AlexNet)
-> 16 - 19 layers (VGG16Net)

Only 3x3 CONV stride 1, pad 1
and 2x2 MAX POOL stride 2
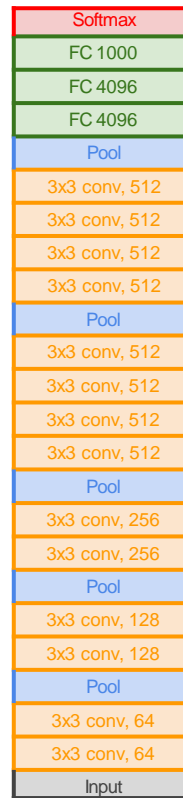
11.7% top 5 error in ILSVRC'13
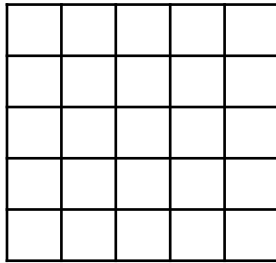(ZFNet)
-> 7.3% top 5 error in ILSVRC'14



AlexNet          VGG16          VGG19
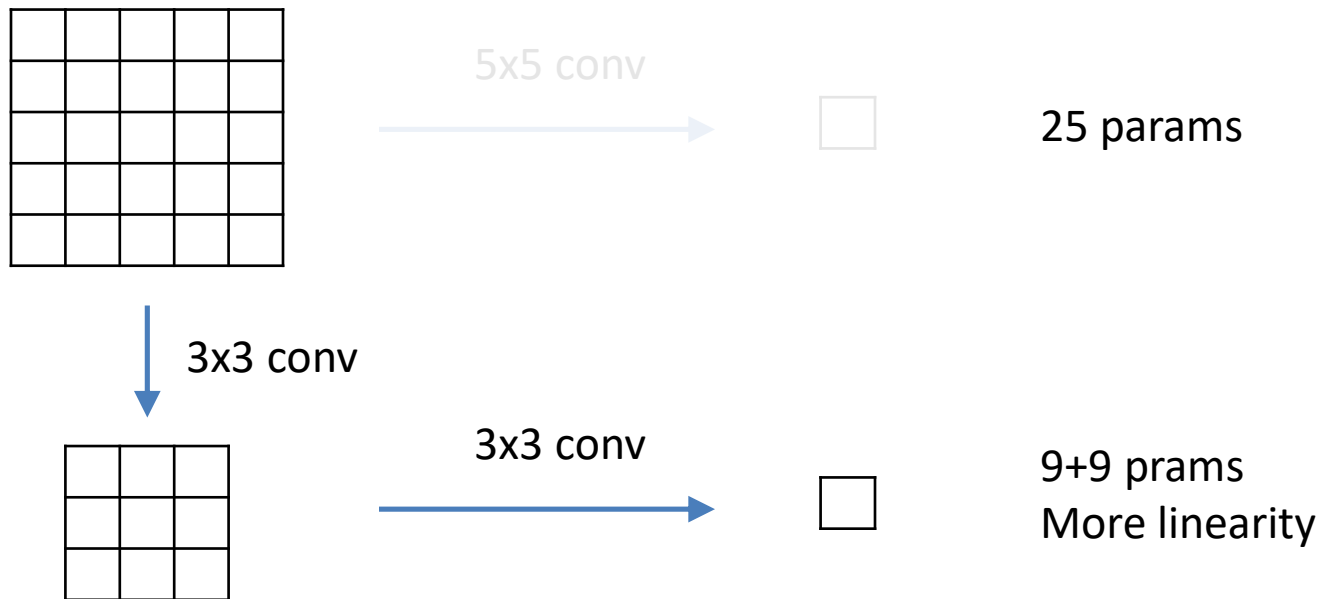
# Case Study: VGGNet

*[Simonyan and Zisserman, 2014]*

5x5 conv

25 params

# Case Study: VGGNet

*[Simonyan and Zisserman, 2014]*

5x5 conv → 25 params

3x3 conv

3x3 conv → 9+9 prams
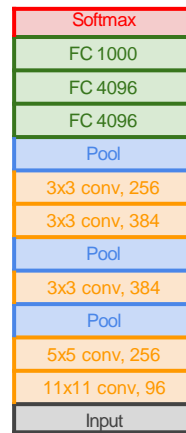More linearity

# Case Study: VGGNet

*[Simonyan and Zisserman, 2014]*
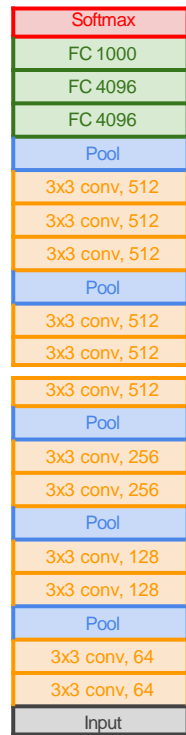
Q: Why use smaller filters? (3x3 conv)

Stack of three 3x3 conv (stride 1) layers has same **effective receptive field** as one 7x7 conv layer
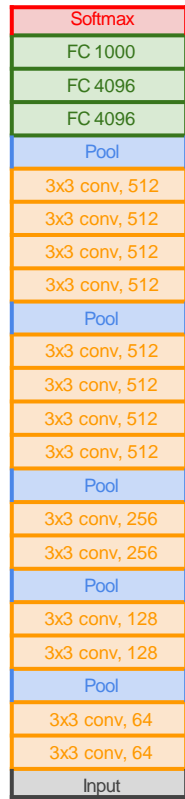
But deeper, more non-linearities

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



AlexNet      VGG16      VGG19

INPUT: [224x224x3]        memory:  224*224*3=150K   params: 0    (not counting biases)
CONV3-64: [224x224x64]  memory:  224*224*64=3.2M    params: (3*3*3)*64 = 1,728
CONV3-64: [224x224x64]  memory:  224*224*64=3.2M    params: (3*3*64)*64 = 36,864
POOL2: [112x112x64]  memory:  112*112*64=800K     params: 0
CONV3-128: [112x112x128]  memory:  112*112*128=1.6M      params: (3*3*64)*128 = 73,728
CONV3-128: [112x112x128]  memory:  112*112*128=1.6M      params: (3*3*128)*128 = 147,456
POOL2: [56x56x128]  memory:  56*56*128=400K    params: 0
CONV3-256: [56x56x256]  memory: 56*56*256=800K  params: (3*3*128)*256 = 294,912   C
ONV3-256: [56x56x256]  memory: 56*56*256=800K  params: (3*3*256)*256 = 589,824   CO
NV3-256: [56x56x256]  memory: 56*56*256=800K  params: (3*3*256)*256 = 589,824
POOL2: [28x28x256]  memory:  28*28*256=200K    params: 0
CONV3-512: [28x28x512]  memory: 28*28*512=400K  params: (3*3*256)*512 = 1,179,648   C
ONV3-512: [28x28x512]  memory: 28*28*512=400K  params: (3*3*512)*512 = 2,359,296   CO
NV3-512: [28x28x512]  memory: 28*28*512=400K  params: (3*3*512)*512 = 2,359,296
POOL2: [14x14x512]  memory:  14*14*512=100K    params: 0
CONV3-512: [14x14x512]  memory:  14*14*512=100K  params: (3*3*512)*512 = 2,359,296   C
ONV3-512: [14x14x512]  memory:  14*14*512=100K  params: (3*3*512)*512 = 2,359,296   CO
NV3-512: [14x14x512]  memory:  14*14*512=100K  params: (3*3*512)*512 = 2,359,296
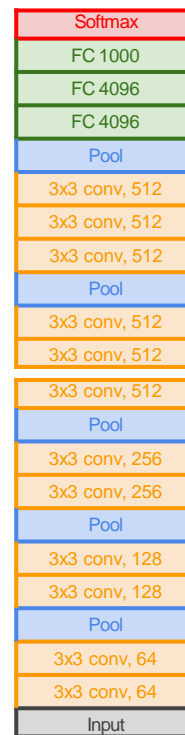POOL2: [7x7x512]  memory:  7*7*512=25K  params: 0
FC: [1x1x4096]  memory:  4096  params: 7*7*512*4096 = 102,760,448   F
C: [1x1x4096]  memory: 4096  params: 4096*4096 = 16,777,216
FC: [1x1x1000]  memory: 1000  params: 4096*1000 = 4,096,000

TOTAL memory: 24M * 4 bytes ~= 96MB / image (only forward! ~*2 for bwd)
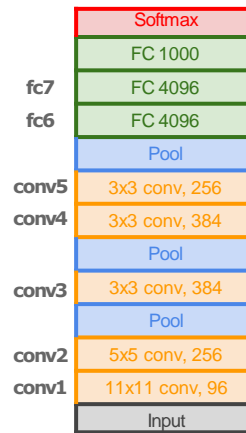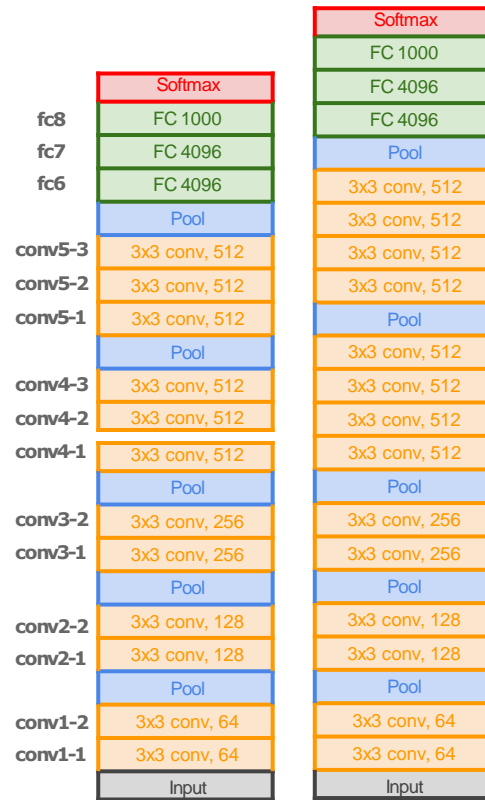TOTAL params: 138M parameters



VGG16

14

# Case Study: 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



AlexNet    VGG16    VGG19

15

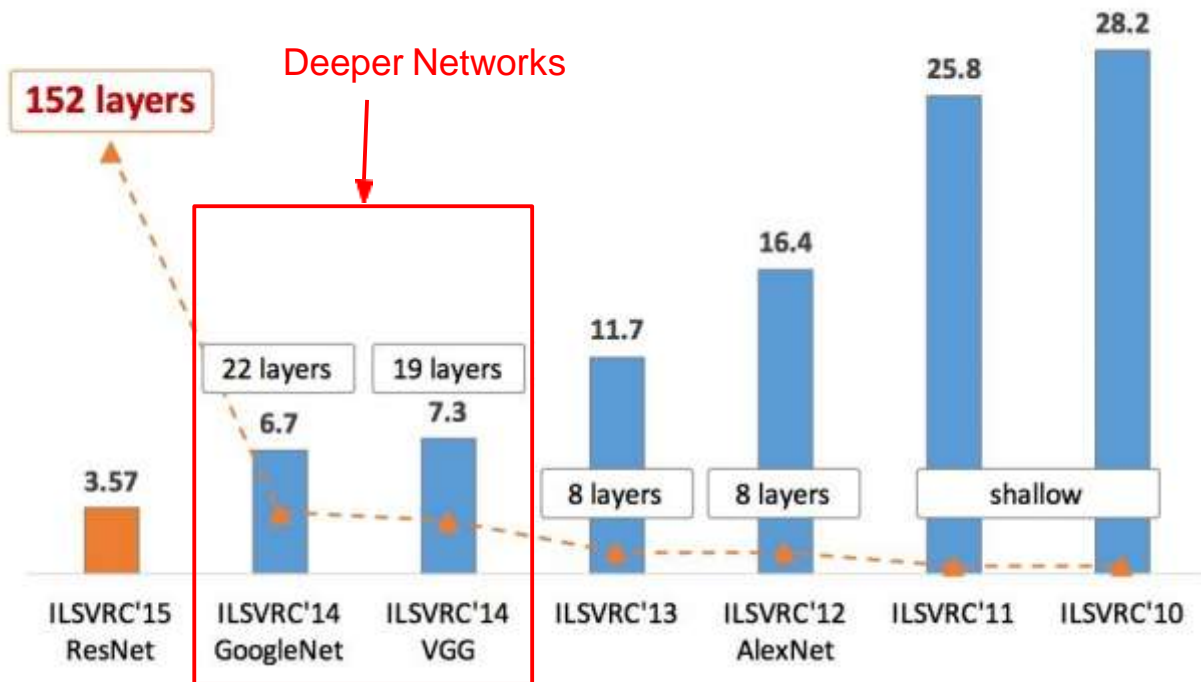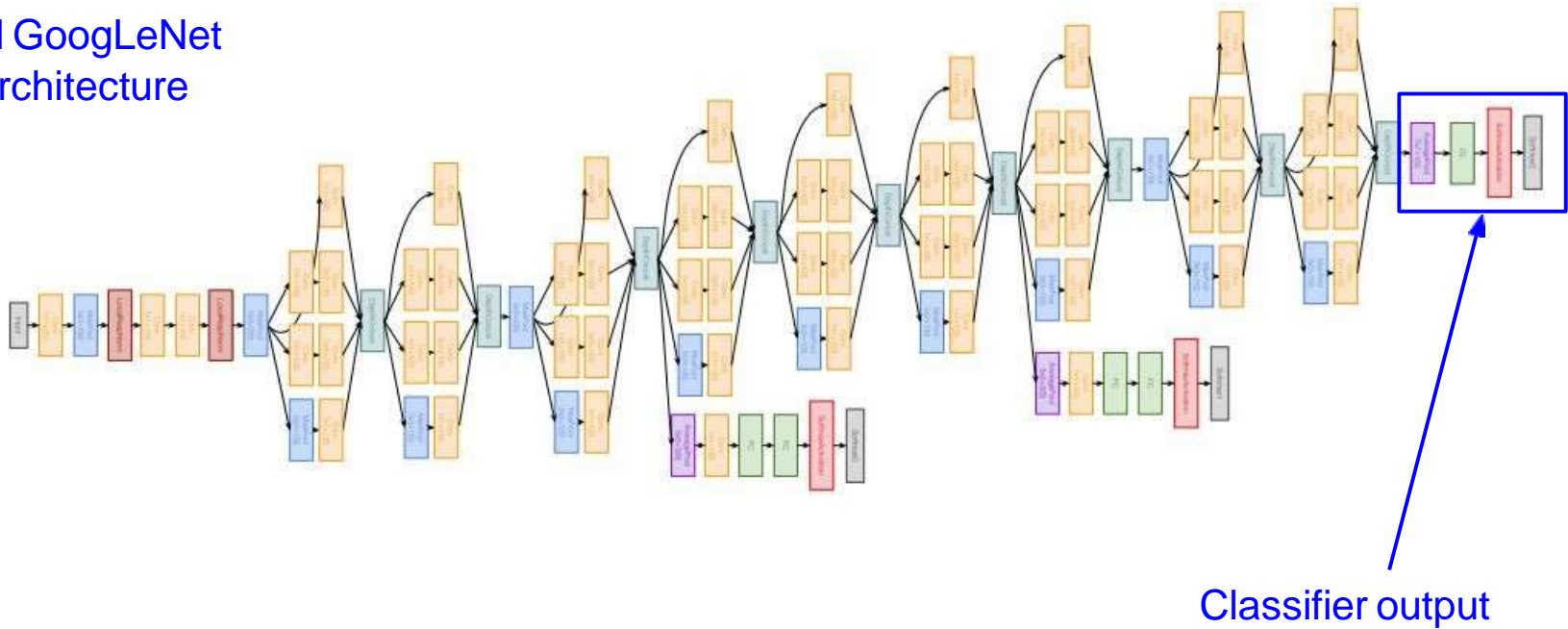# ImageNet Large Scale Visual Recognition Challenge (ILSVRC) winners



Figure copyright Kaiming He, 2016. Reproduced with permission.

# Case Study: GoogLeNet

*[Szegedy et al., 2014]*

Full GoogLeNet
architecture



Classifier output

# Case Study: GoogLeNet

*[Szegedy et al., 2014]*

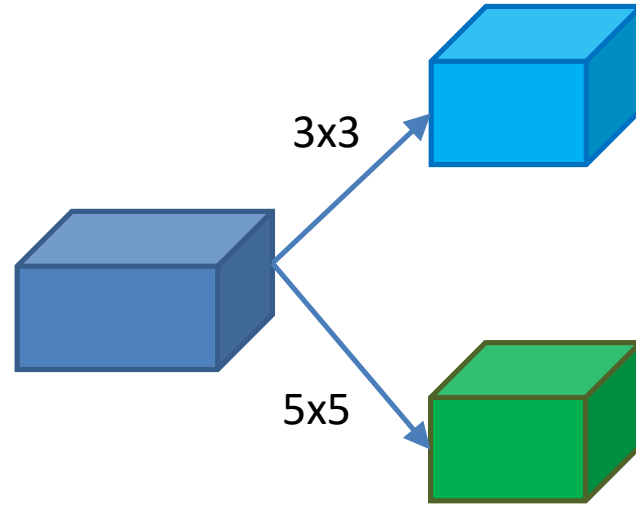"Inception module": design a good local network topology (network within a network) and then stack these modules on top of each other



Inception module

# Case Study: GoogLeNet

*[Szegedy et al., 2014]*



Naive Inception module

Inception module with dimension reduction
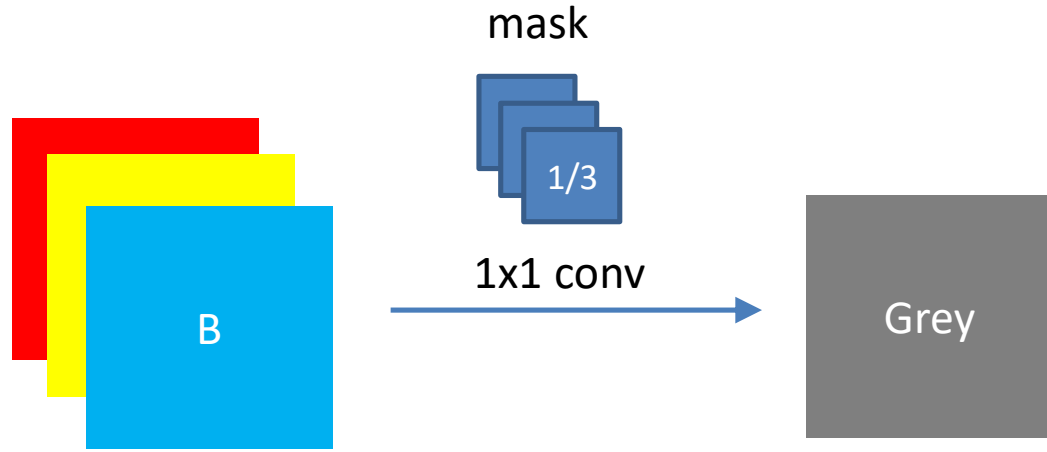
# Case Study: GoogLeNet

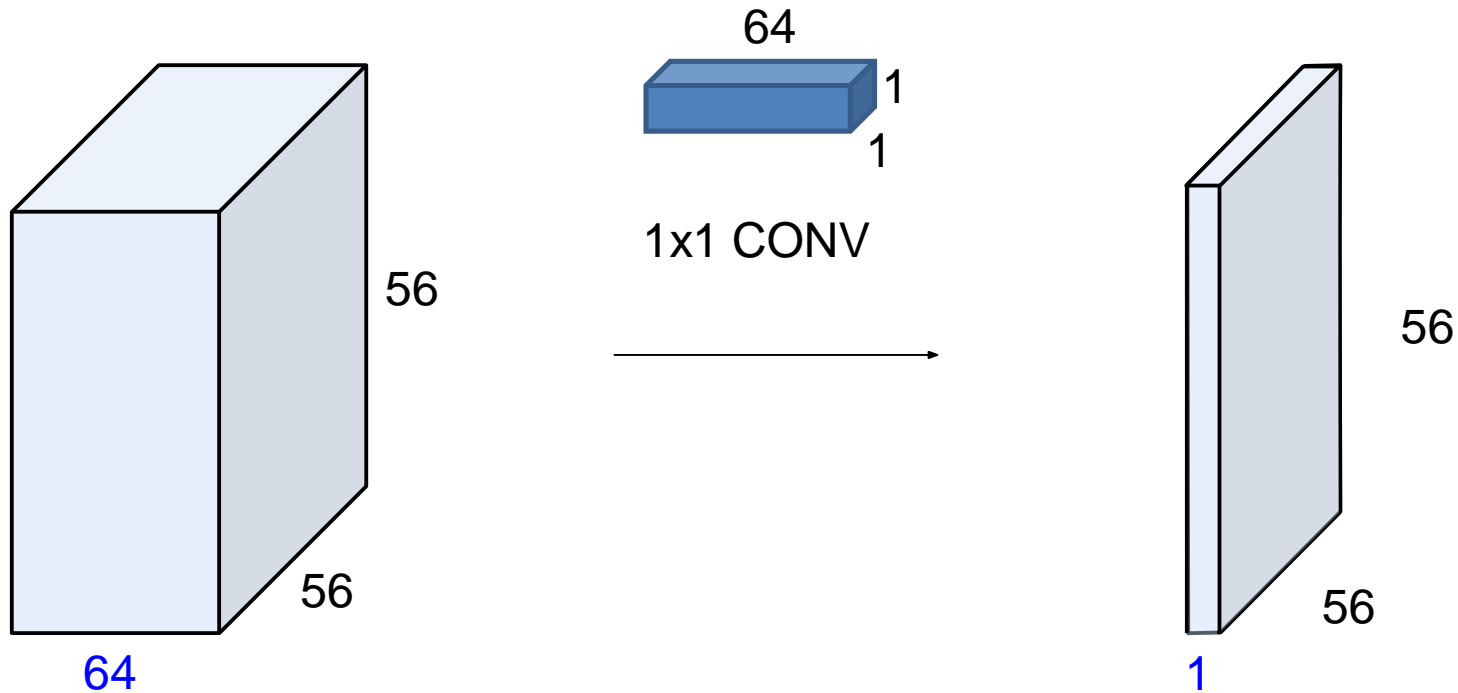*[Szegedy et al., 2014]*



Naive Inception module
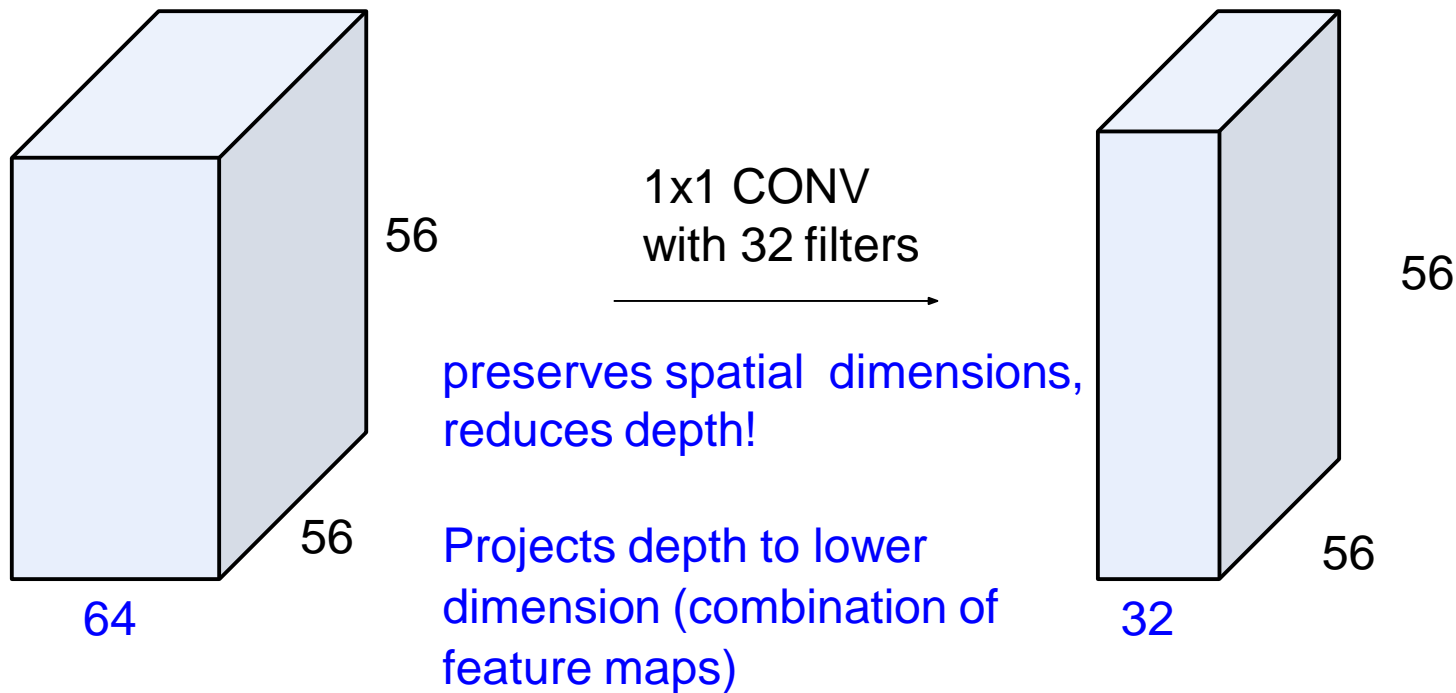
3x3

5x5

# Reminder: 1x1 convolutions

# Reminder: 1x1 convolutions

# Reminder: 1x1 convolutions

64

1

1

1x1 CONV

56

56

64

56

56

1

# Reminder: 1x1 convolutions

1x1 CONV
with 32 filters

56

56

64

→

56

56

32

preserves spatial dimensions, reduces depth!

Projects depth to lower dimension (combination of feature maps)

# Case Study: GoogLeNet

256

128

3x3

256 X 3X3 X 128
= 294,912

# Case Study: GoogLeNet

256

128

3x3

256 X 3X3 X 128
= 294,912

1x1

128

128

3x3

256 X 1X1 X 128
+ 128 X 3X3 X 128
= 189,224

# Case Study: GoogLeNet

*[Szegedy et al., 2014]*



Inception module with dimension reduction

3x3 max pooling, stride=1



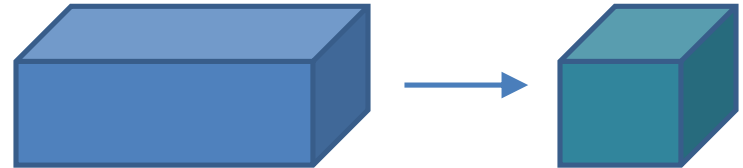Feature map          Enhanced feature map

1x1 Convolution

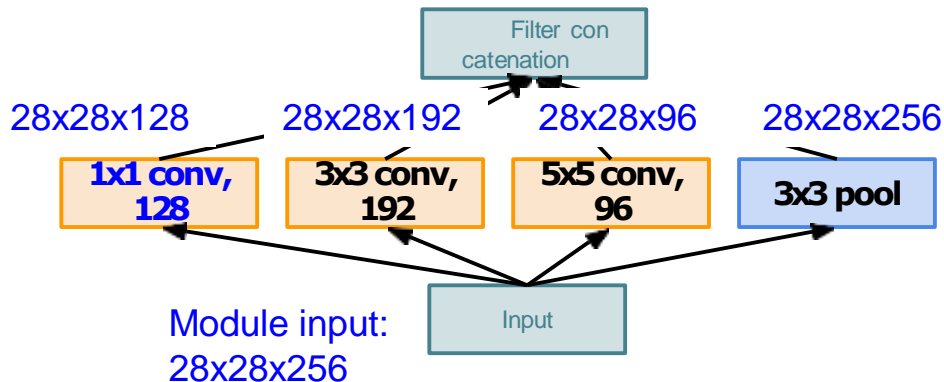# Case Study: GoogLeNet

*[Szegedy et al., 2014]*

Example:

Q3:What is output size after filter concatenation?

28x28x(128+192+96+256) = 28x28x672



Filter concatenation

28x28x128    28x28x192    28x28x96    28x28x256

1x1 conv, 128    3x3 conv, 192    5x5 conv, 96    3x3 pool

Module input: 28x28x256

Input

Naive Inception module

**Conv Ops:**
[1x1 conv, 128] 28x28x128x1x1x256
[3x3 conv, 192] 28x28x192x3x3x256
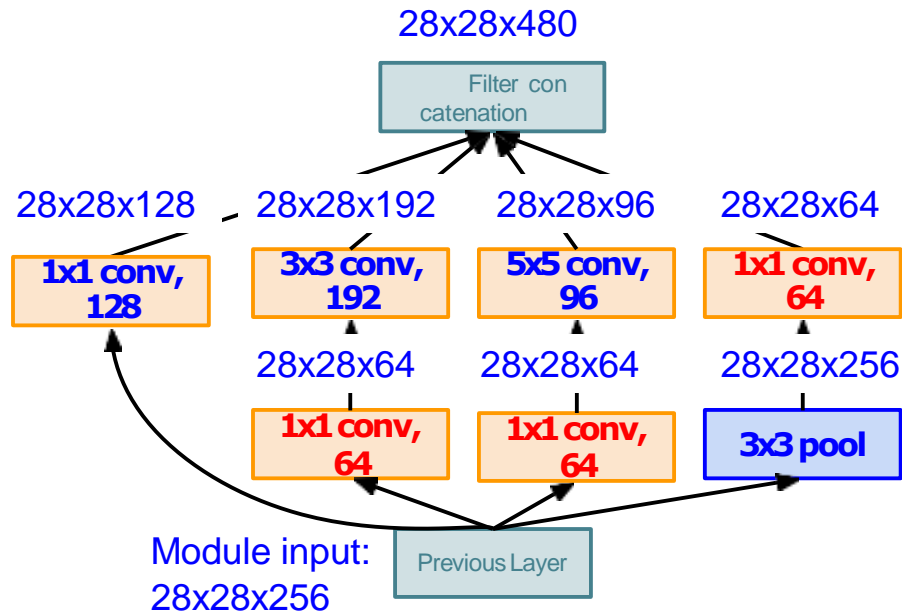[5x5 conv, 96] 28x28x96x5x5x256
**Total: 854M ops**

Very expensive compute

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

# Case Study: GoogLeNet

*[Szegedy et al., 2014]*



Inception module with dimension reduction

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

**Conv Ops:**
[1x1 conv, 64] 28x28x64x1x1x256
[1x1 conv, 64] 28x28x64x1x1x256
[1x1 conv, 128] 28x28x128x1x1x256
[3x3 conv, 192] 28x28x192x3x3x64
[5x5 conv, 96] 28x28x96x5x5x64
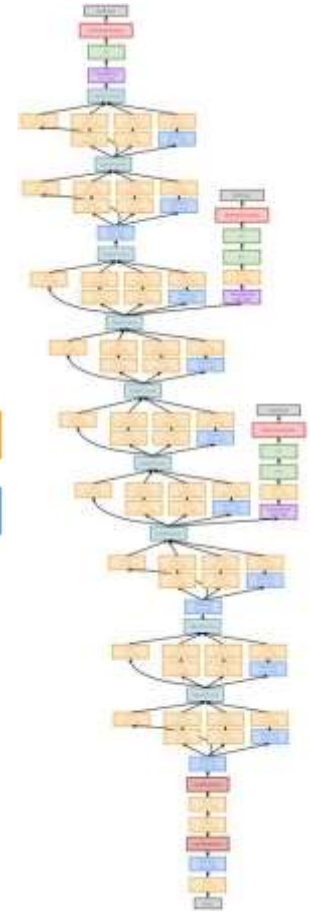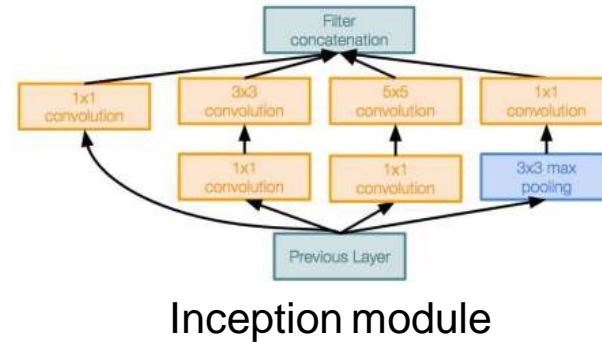[1x1 conv, 64] 28x28x64x1x1x256
**Total: 358M ops**

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

# Case Study: GoogLeNet

*[Szegedy et al., 2014]*

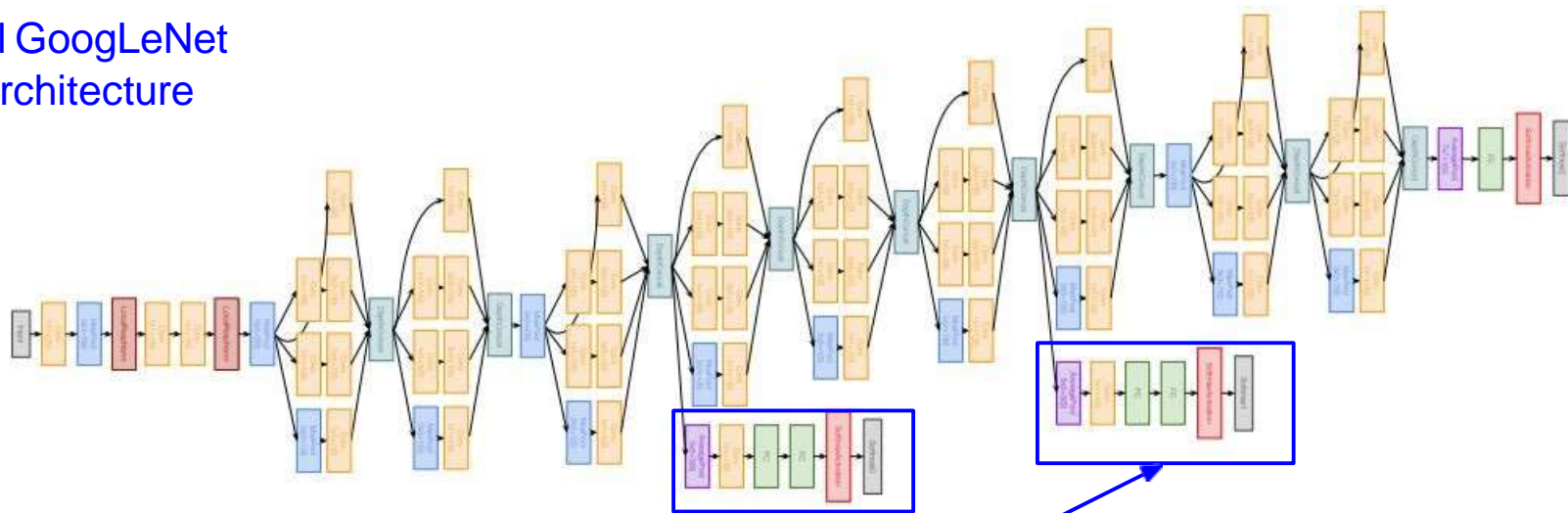Deeper networks, with computational efficiency

- 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)
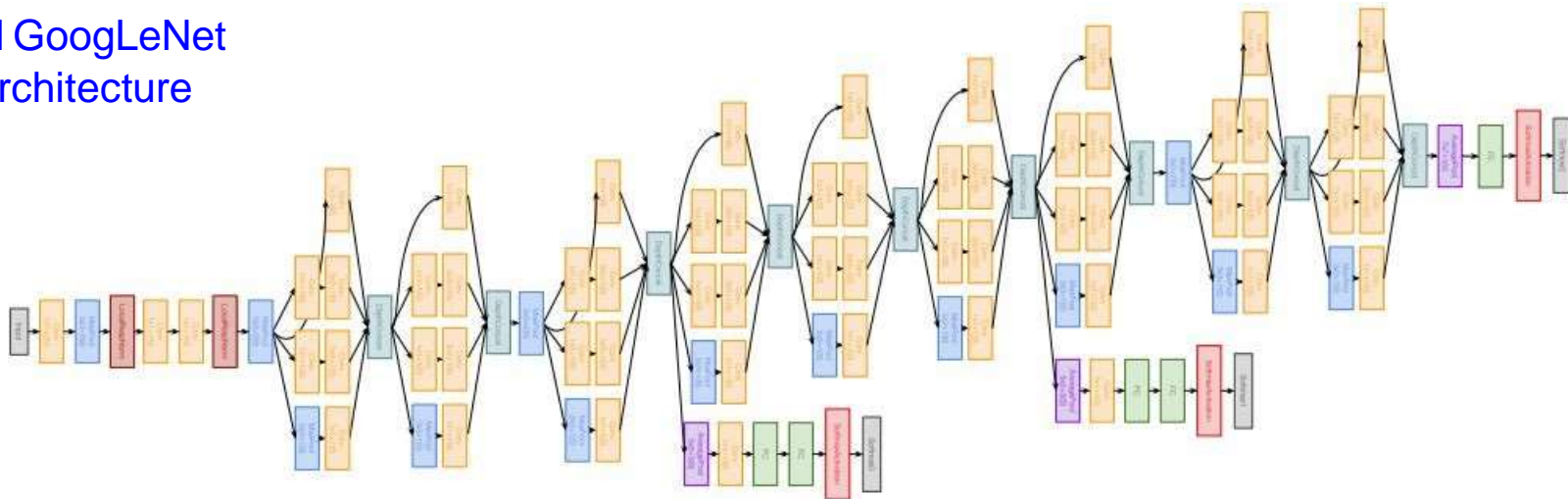
Inception module

# Case Study: GoogLeNet

*[Szegedy et al., 2014]*

Full GoogLeNet
architecture



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

# Case Study: GoogLeNet

*[Szegedy et al., 2014]*

Full GoogLeNet
architecture



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

# ImageNet Large Scale Visual Recognition Challenge (ILSVRC) winners
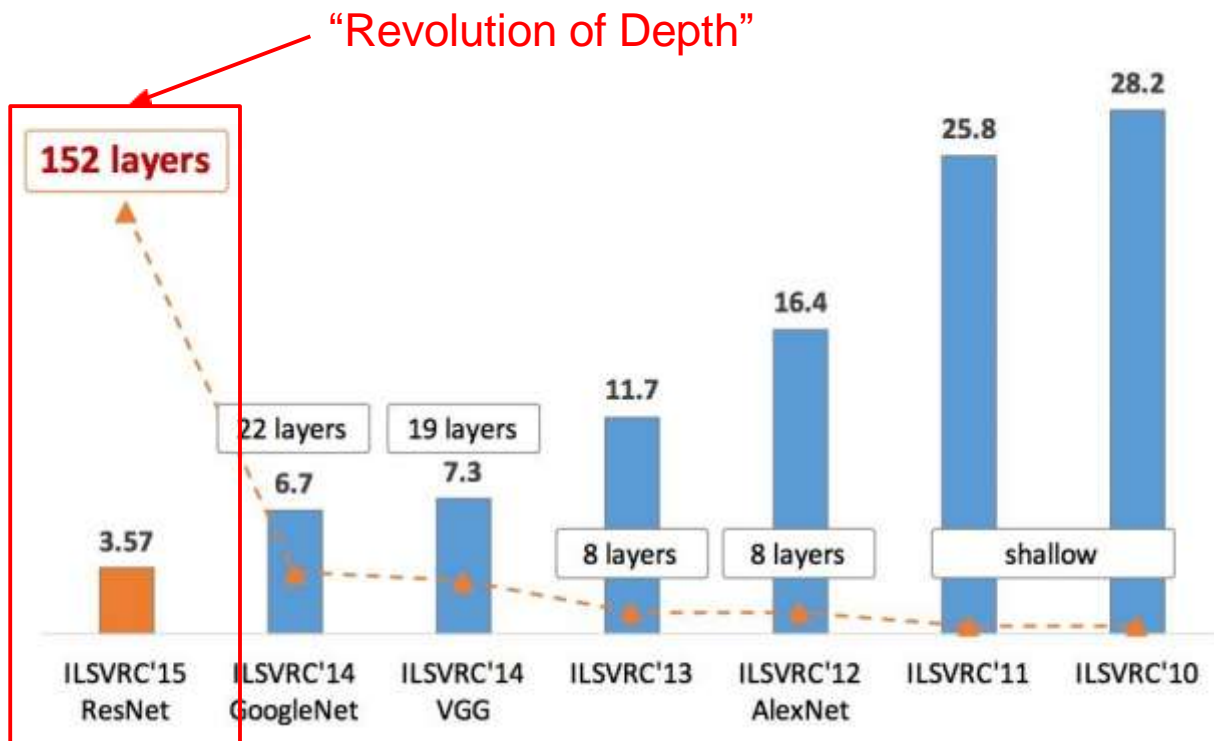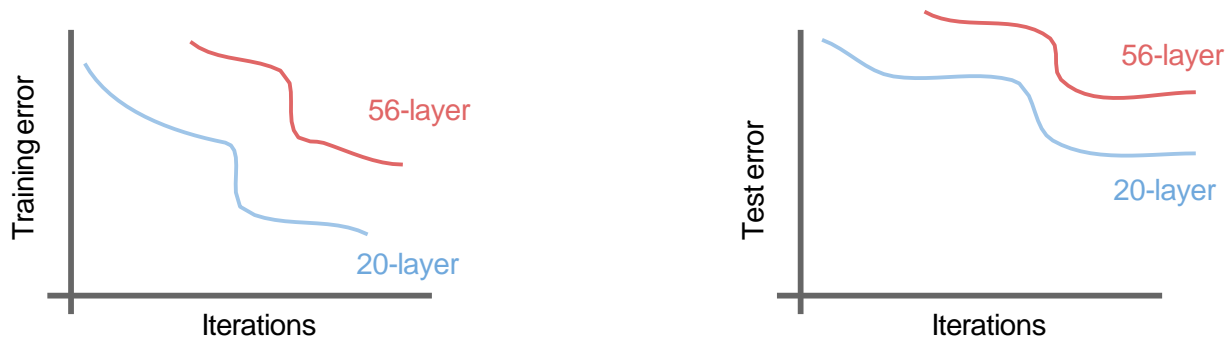


"Revolution of Depth"

Figure copyright Kaiming He, 2016. Reproduced with permission.

# Case Study: 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
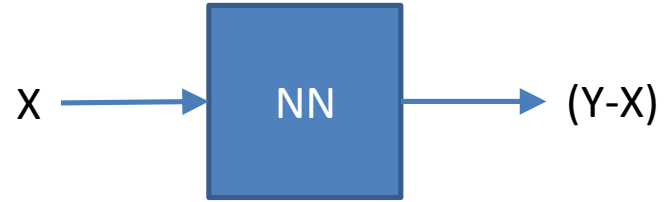-> The deeper model performs worse, but it's not caused by overfitting!
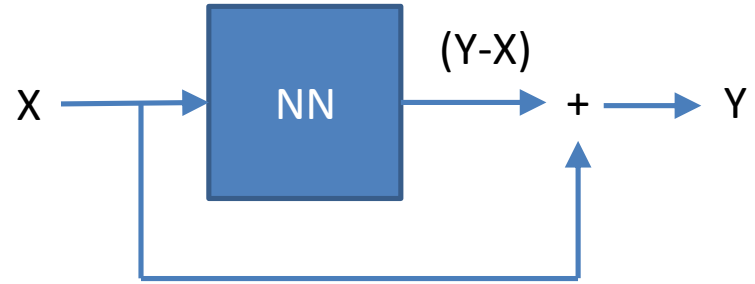
# Case Study: ResNet

| X | Y |
|---|---|
| 1 | 0.9 |
| 2 | 2.1 |
| 3 | 3.0 |
| 4 | 4.2 |

X → NN → Y

# Case Study: ResNet

| X | Y | Y-X |
|---|---|-----|
| 1 | 0.9 | -0.1 |
| 2 | 2.1 | 0.1 |
| 3 | 3.0 | 0.0 |
| 4 | 4.2 | 0.2 |

X → NN → (Y-X)

# Case Study: ResNet

| X | Y | Y-X |
|---|---|-----|
| 1 | 0.9 | -0.1 |
| 2 | 2.1 | 0.1 |
| 3 | 3.0 | 0.0 |
| 4 | 4.2 | 0.2 |

X → NN → (Y-X) → + → Y

# Case Study: ResNet

*[He et al., 2015]*

| X | Y |
|---|---|
| 1 | 0.9 |
| 2 | 2.1 |
| 3 | 3.0 |
| 4 | 4.2 |

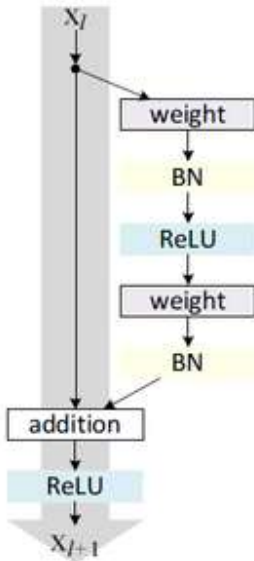# Case Study: ResNet

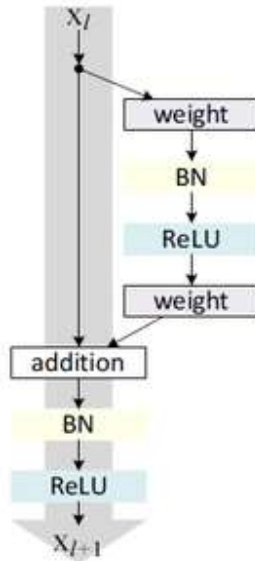*[He et al., 2015]*

# Case Study: ResNet
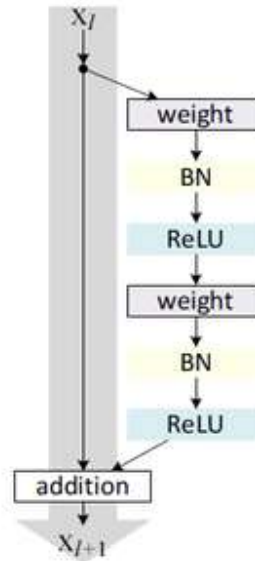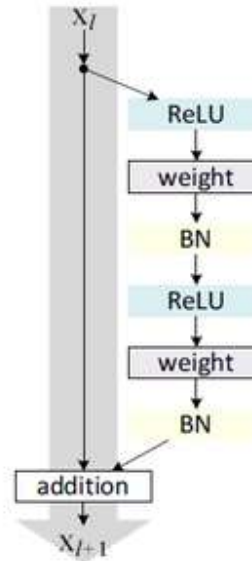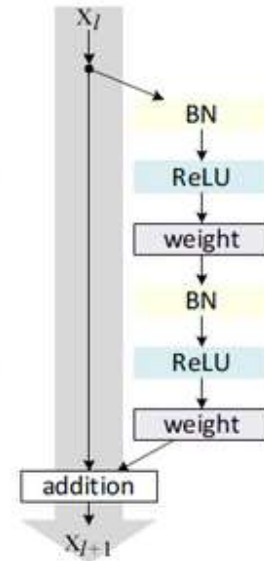
*[He et al., 2015]*

# Case Study: ResNet



(a) original

(b) BN after addition

(c) ReLU before addition
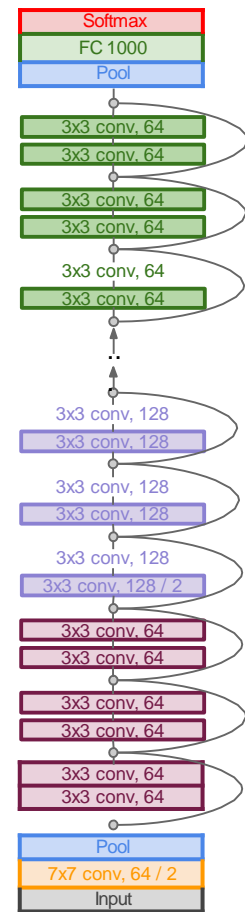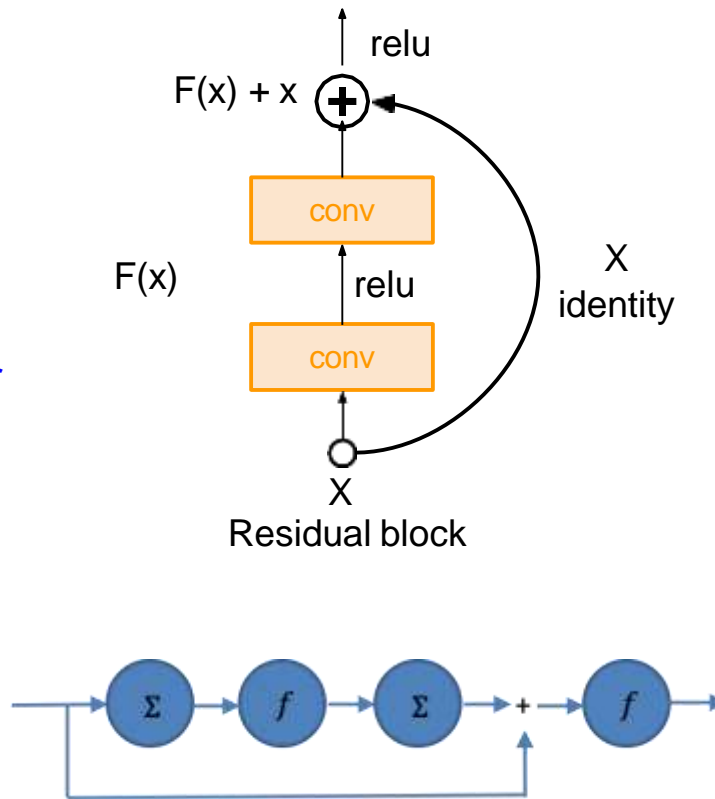
(d) ReLU-only pre-activation

(e) **full pre-activation**

# Case Study: 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)
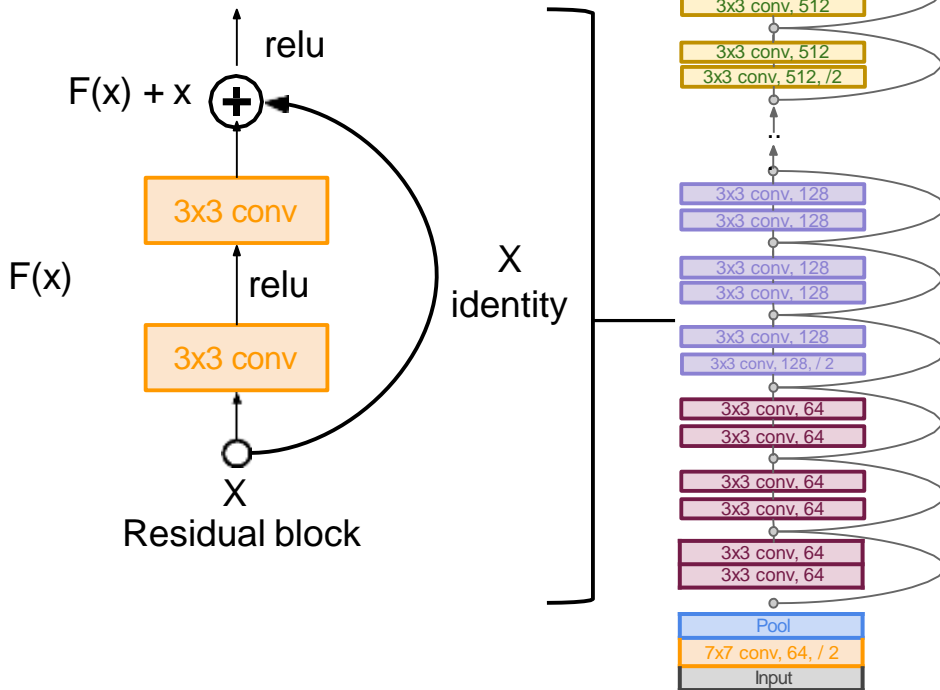- Swept all classification and detection competitions in I LSVRC'15 and COCO'15!



Residual block

# Case Study: 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)
- Global average pooling layer after last conv. layer



Residual block

# Case Study: 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

# ImageNet Large Scale Visual Recognition Challenge (ILSVRC) winners
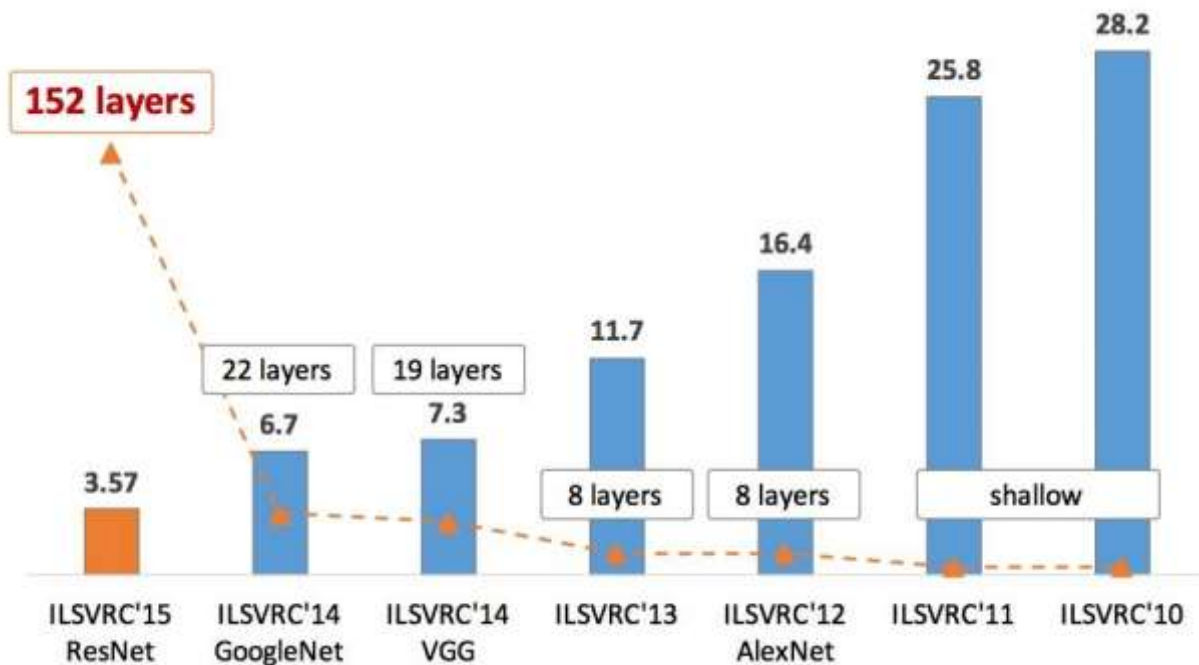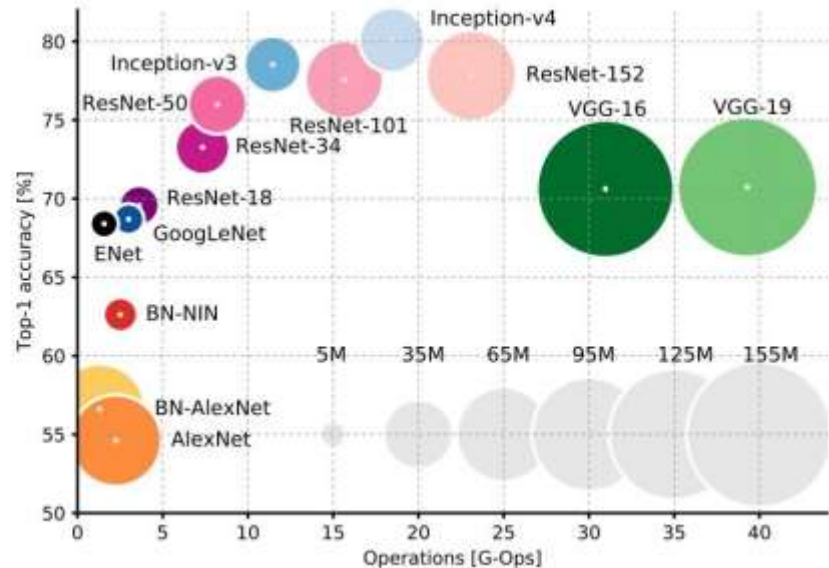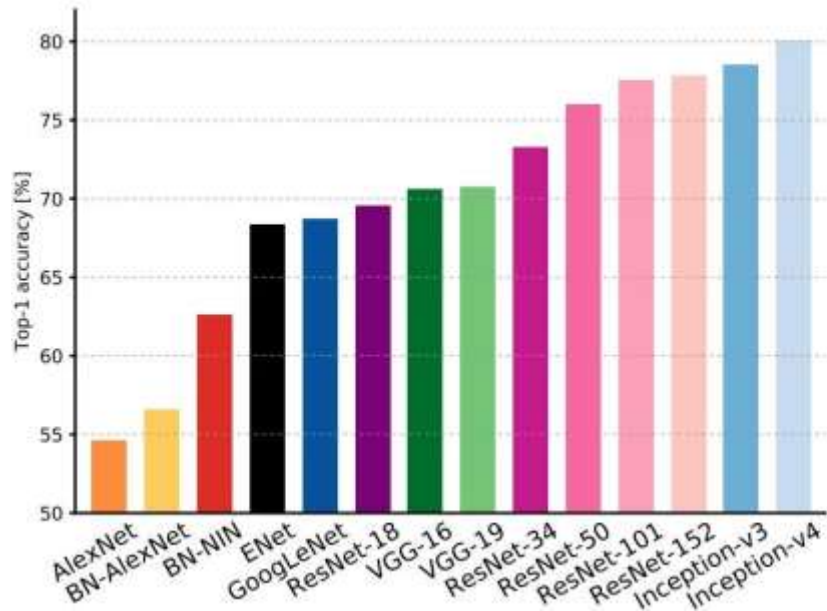


Figure copyright Kaiming He, 2016. Reproduced with permission.
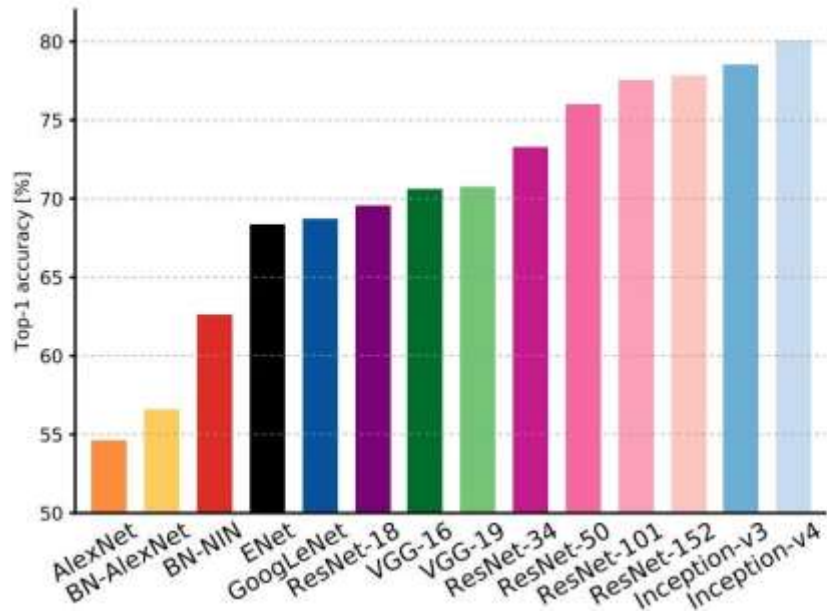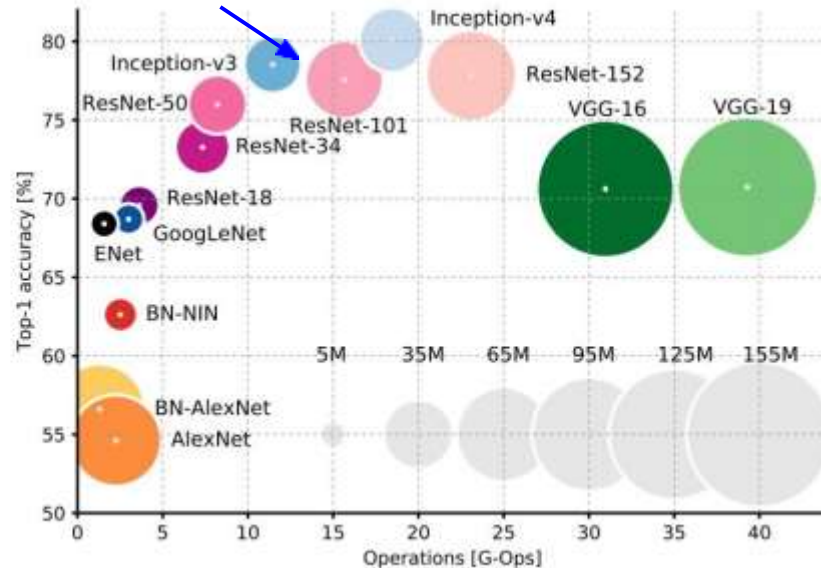
# Comparing complexity...



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

Figures copyright Alfredo Canziani, Adam Paszke, Eugenio Culurciello, 2017. Reproduced with permission.

# Comparing complexity...



ResNet:
Moderate efficiency depending on model, highest accuracy

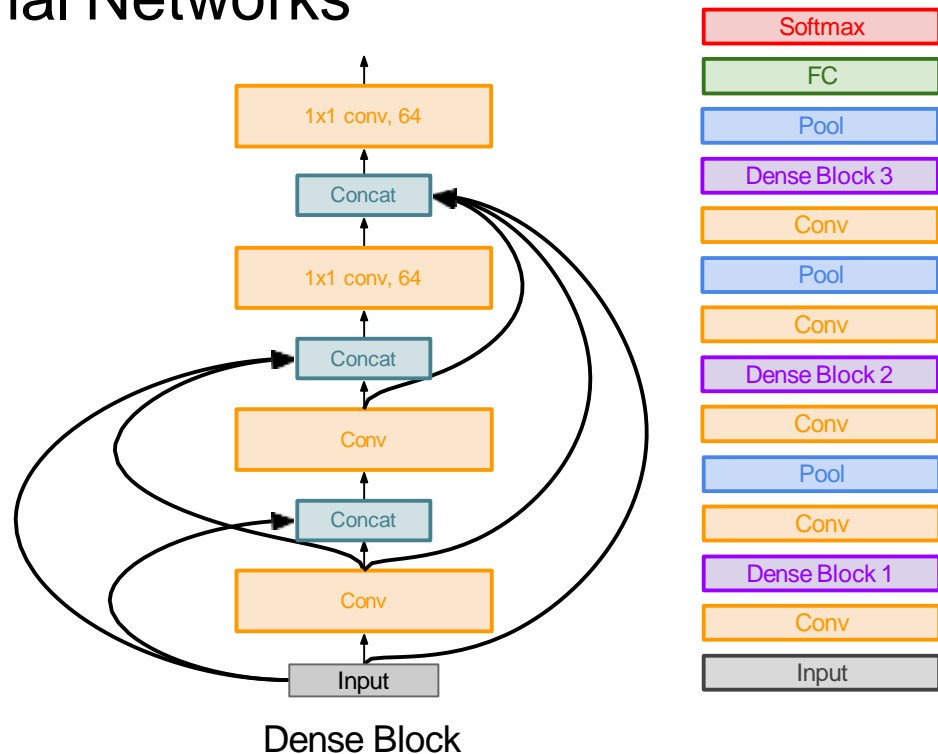An Analysis of Deep Neural Network Models for Practical Applications, 2017.

# 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, st rengthens feature propagation, encourages feature reuse
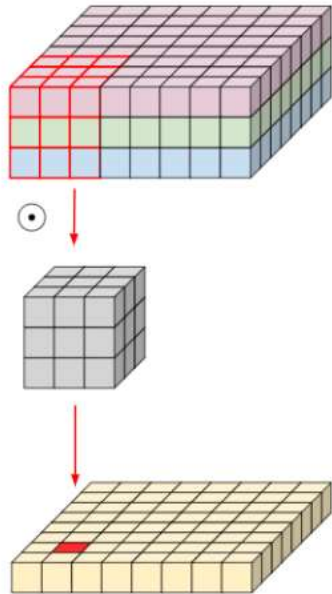


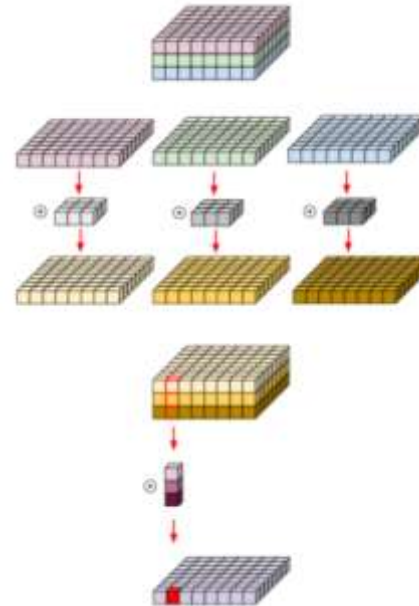Dense Block

# Summary: CNN Architectures

- 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
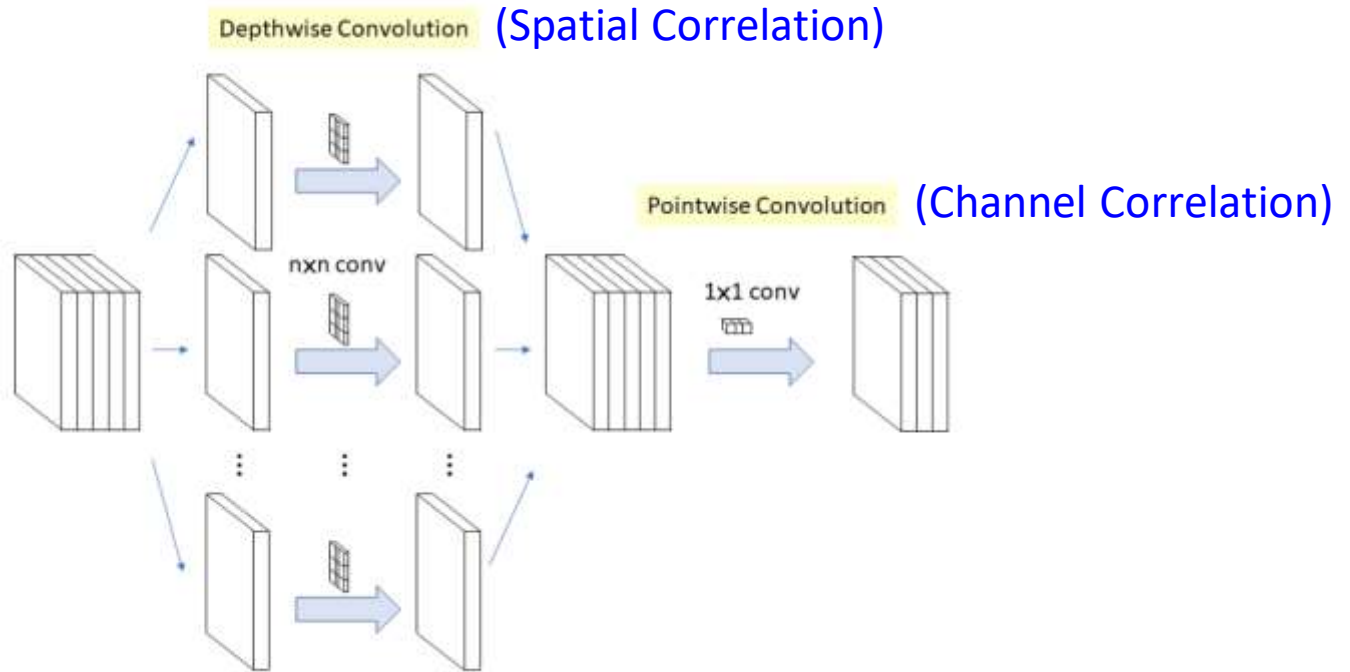
# Depthwise Separable Convolution
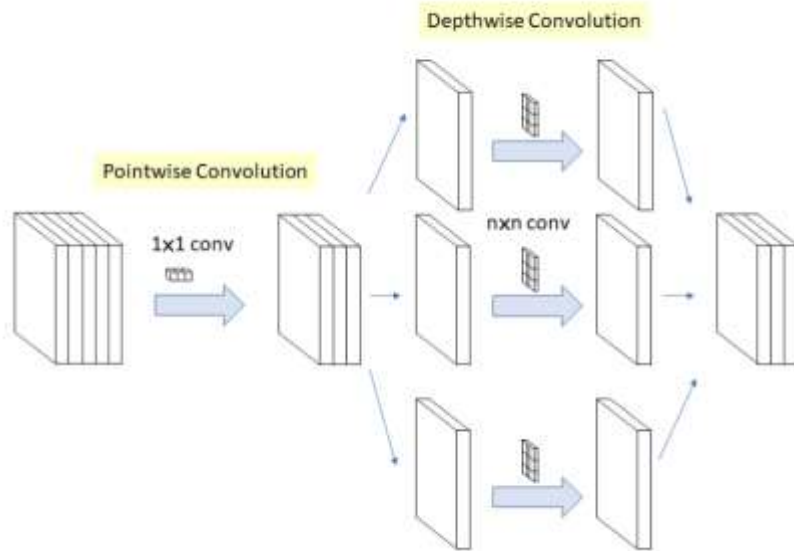
Regular Conv

Depthwise Separable Conv

# Depthwise Separable Convolution

# Depthwise Separable Convolution

**Xception**

**Depthwise Separable Convolution**