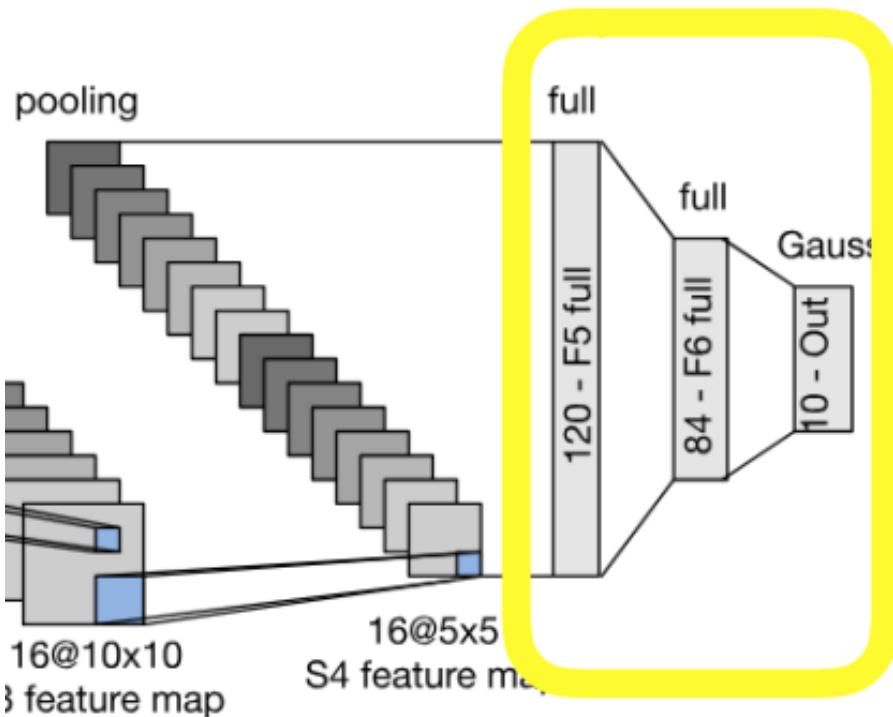


The Last Layer(s)



- Convolution layers need relatively few parameters

$$c_i \times c_o \times k^2$$

- Last layer needs many parameters for n classes

$$c \times m_w \times m_h \times n$$

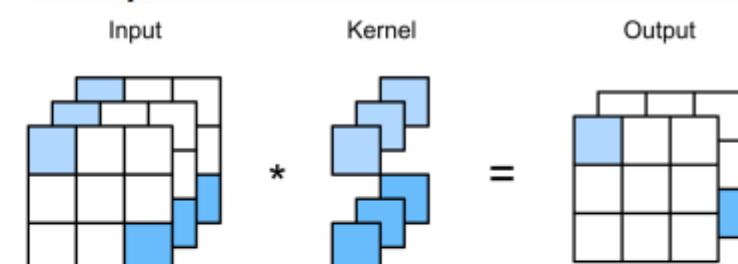
- LeNet $16 \times 5 \times 5 \times 120 = 48k$
- AlexNet $256 \times 5 \times 5 \times 4096 = 26M$
- VGG $512 \times 7 \times 7 \times 4096 = 102M$



gluon-cv.mxnet.io

Breaking the Curse of the Last Layer

- Key Idea
 - **Get rid of the fully connected last layer(s)**
 - Convolutions and pooling reduce resolution
(e.g. stride of 2 reduces resolution 4x)
- Implementation details
 - Reduce resolution progressively
 - Increase number of channels
 - Use **1x1 convolutions** (they only act per pixel)
- **Global average pooling in the end**

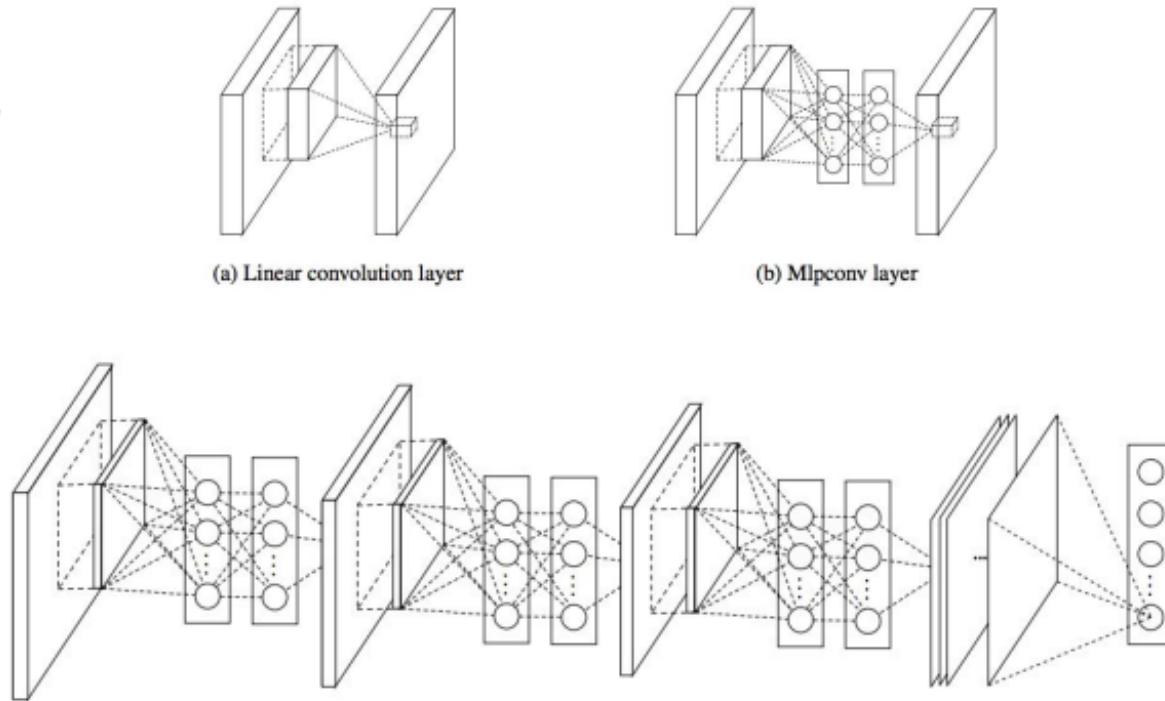


github.com/mynet-in

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



Figures copyright Lin et al., 2014. Reproduced with permission.

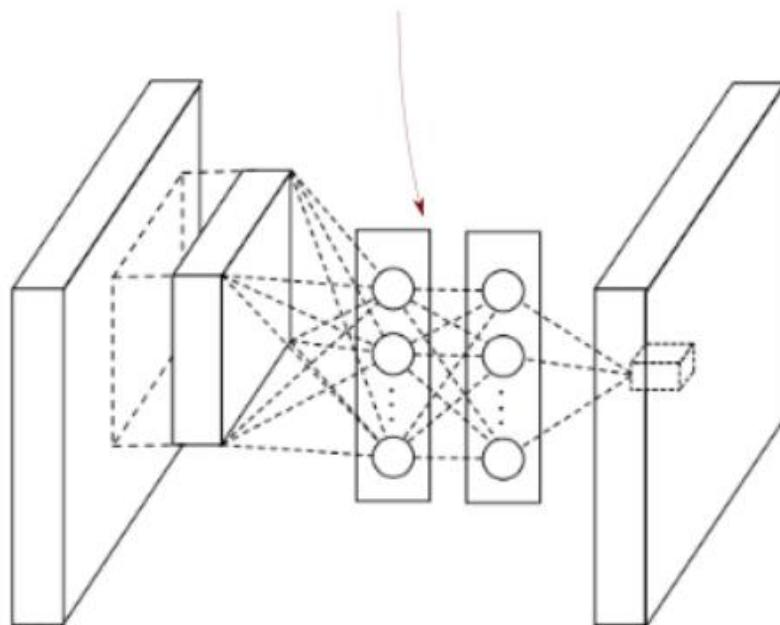
What's a 1x1 convolution anyway?

- Extreme case
1x1 image with n channels
- Equivalent to MLP
- Pooling allows for
translation invariance of
detection (e.g. 5x5)



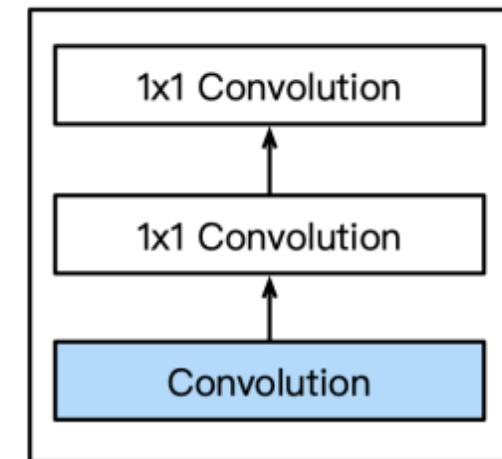
courses.d2l.ai/berkeley-stat-157

Non linear mapping introduced by mlpconv layer consisting of multiple fully connected layers with non linear activation function.

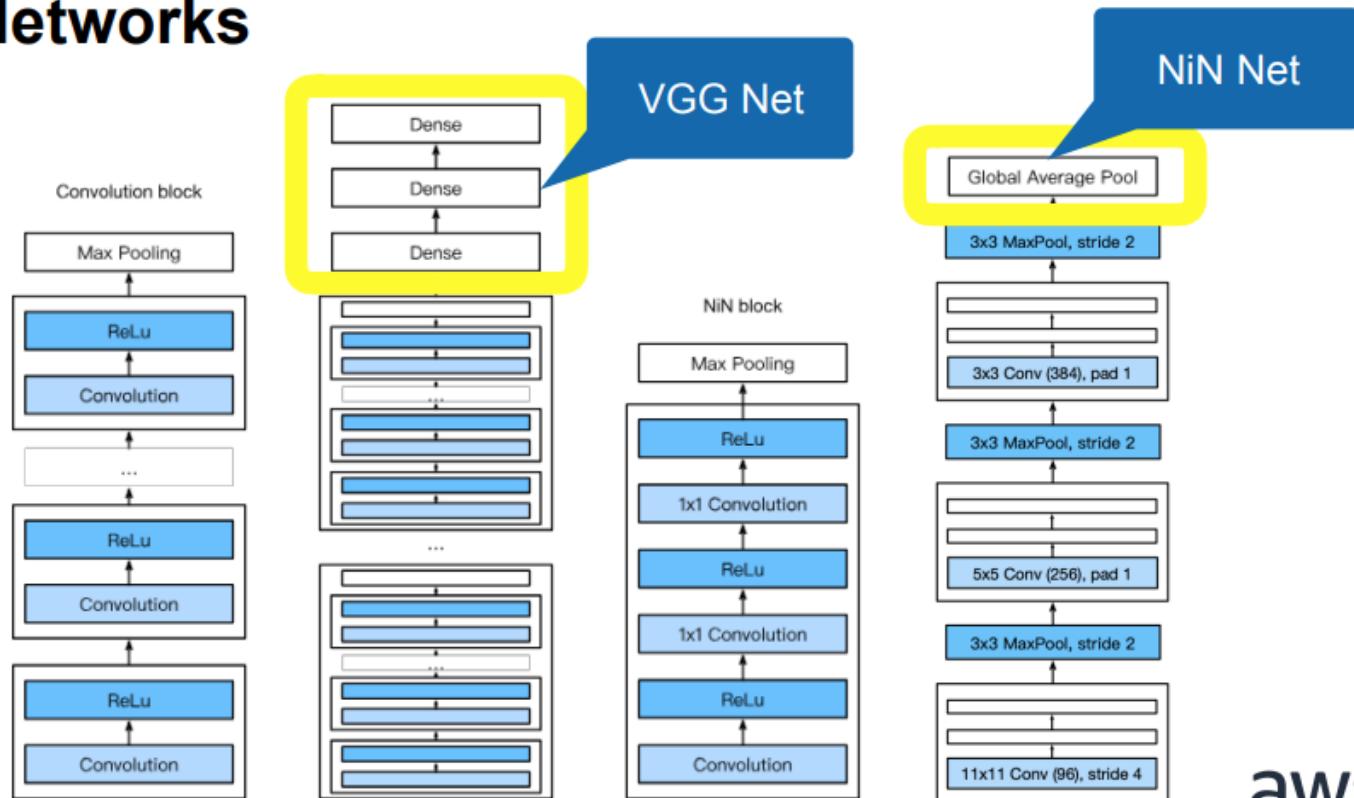


NiN Block

- A convolutional layer
 - kernel size, stride, and padding are hyper-parameters
- Following by two 1x1 convolutions
 - 1 stride and no padding, share the same output channels as first layer
 - Act as dense layers



NiN Networks



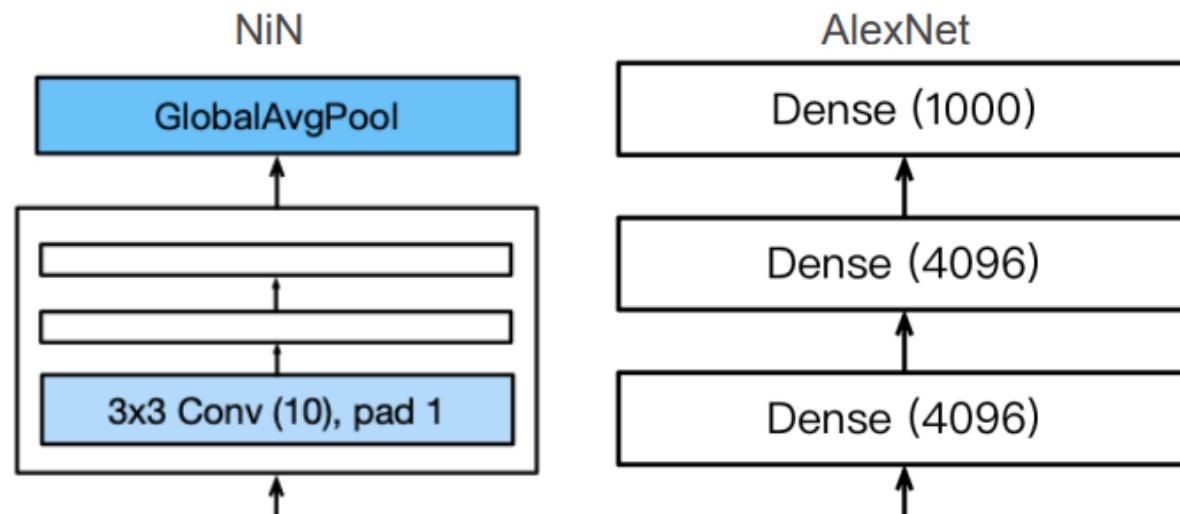
gluon-cv.mxnet.io

Zhang, A., Lipton, Z. C., Li, M., & Smola, A. J. (2021). Dive into deep learning. *arXiv preprint arXiv:2106.11342*.



NiN Last Layers

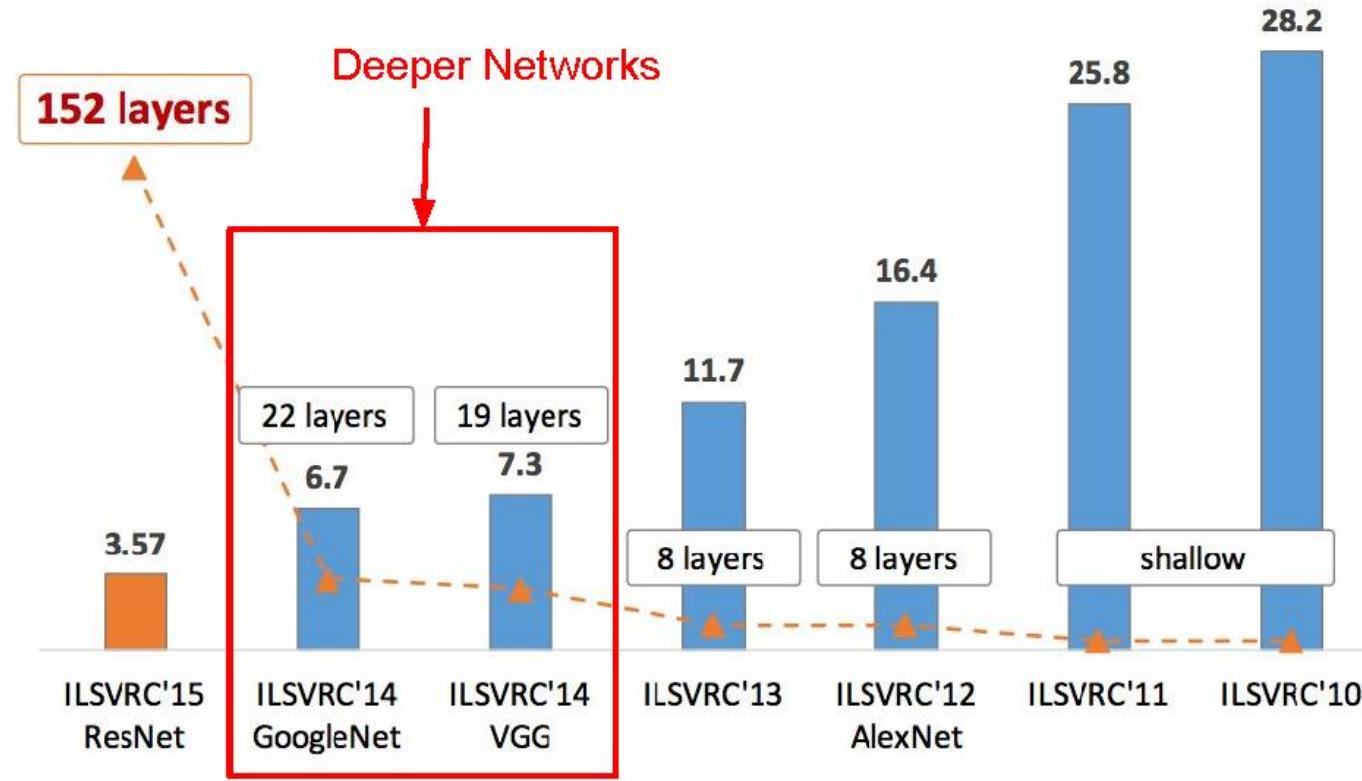
- Replaced AlexNet's dense layers with a NiN block
- Global average pooling layer to combine outputs



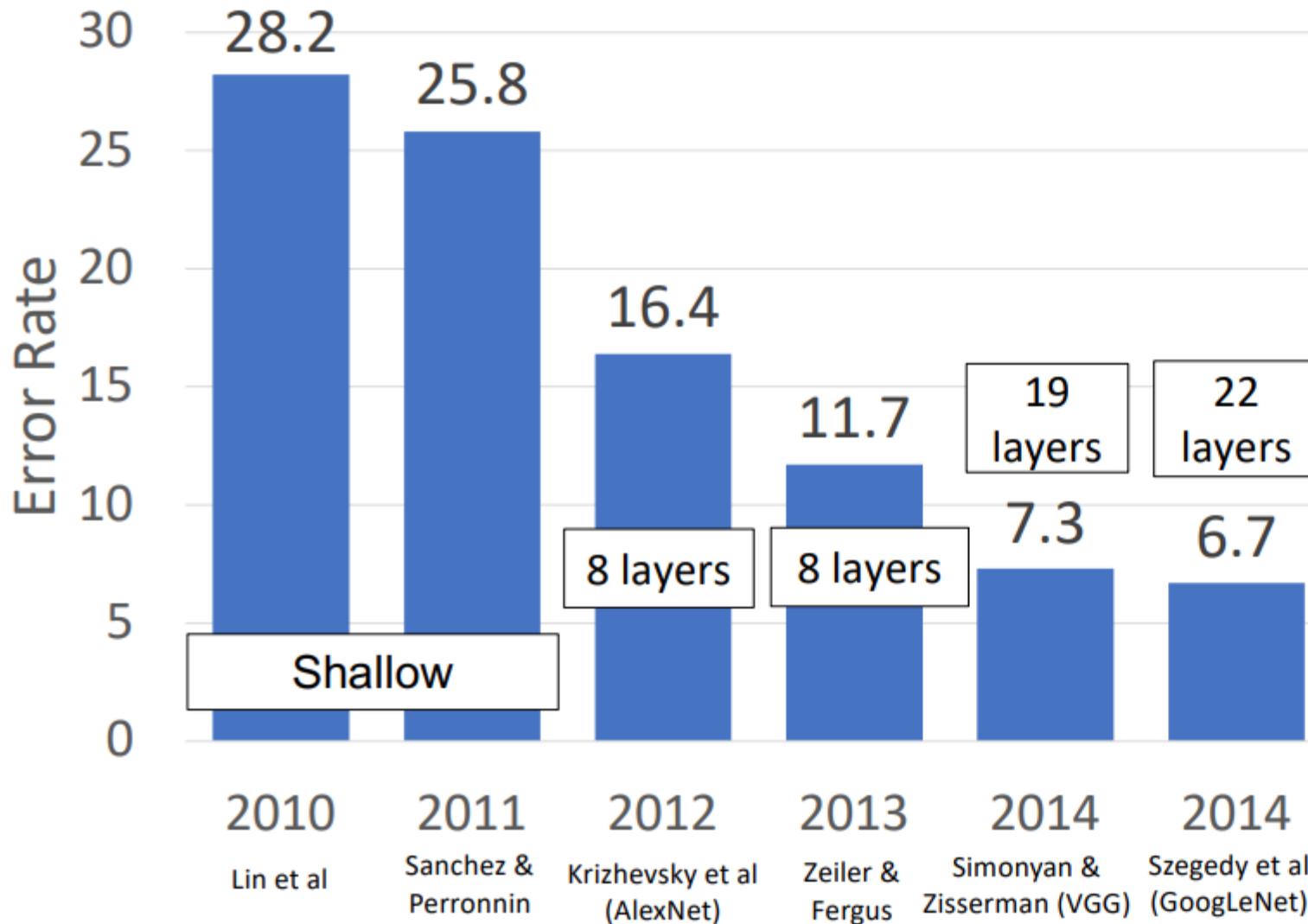
aws

Zhang, A., Lipton, Z. C., Li, M., & Smola, A. J. (2021). Dive into deep learning. *arXiv preprint arXiv:2106.11342*.

ImageNet Large Scale Visual Recognition Challenge (ILSVRC) winners

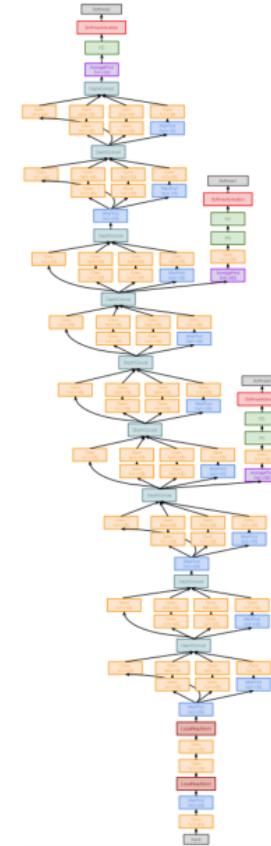


ImageNet Classification Challenge



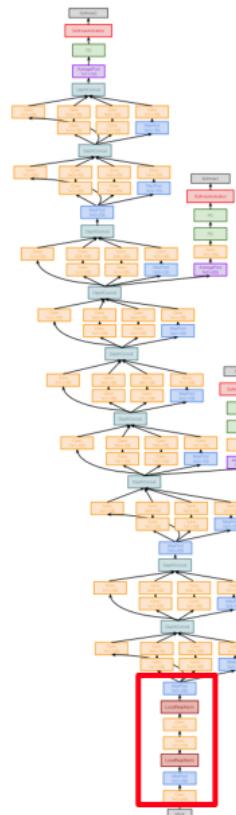
GoogLeNet: Focus on Efficiency

Many innovations for efficiency: reduce parameter count, memory usage, and computation



GoogLeNet: Aggressive Stem

Stem network at the start aggressively downsamples input
(Recall in VGG-16: Most of the compute was at the start)



GoogLeNet: Aggressive Stem

Stem network at the start aggressively downsamples input
(Recall in VGG-16: Most of the compute was at the start)

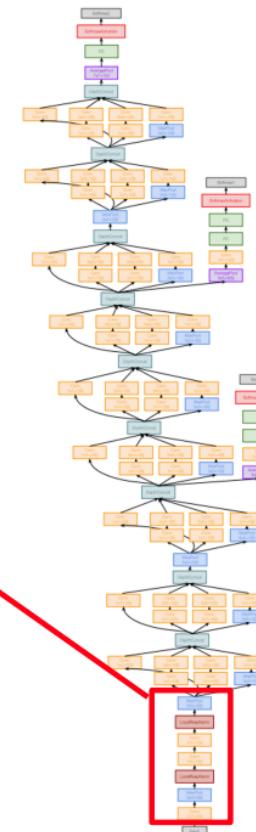
Layer	Input size			Layer			Output size			memory (KB)	params (K)	flop (M)
	C	H / W	filters	kernel	stride	pad	C	H/W				
conv	3	224	64	7	2	3	64	112	3136	9	118	
max-pool	64	112		3	2	1	64	56	784	0	2	
conv	64	56	64	1	1	0	64	56	784	4	13	
conv	64	56	192	3	1	1	192	56	2352	111	347	
max-pool	192	56		3	2	1	192	28	588	0	1	

Total from 224 to 28 spatial resolution:

Memory: 7.5 MB

Params: 124K

MFLOP: 418

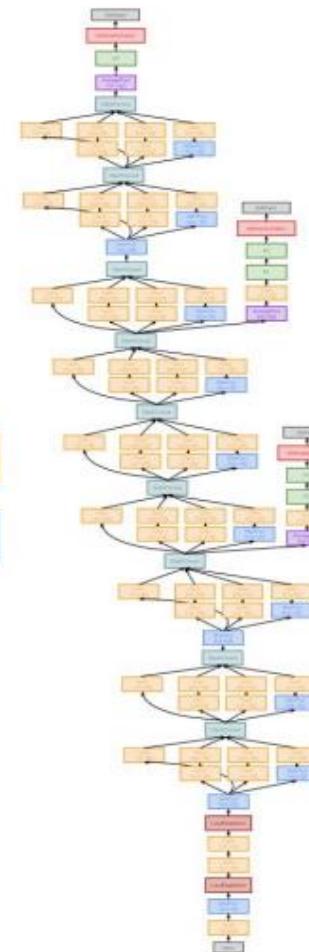
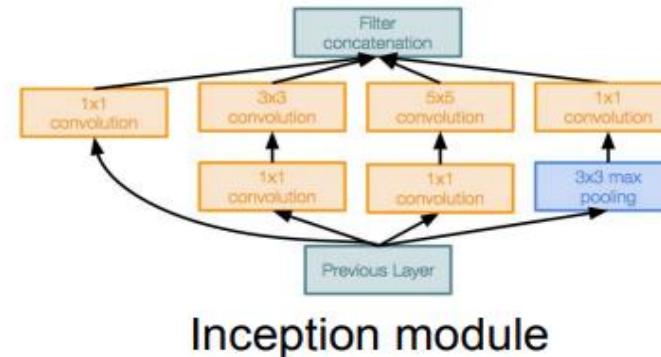


Case Study: GoogLeNet

[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)



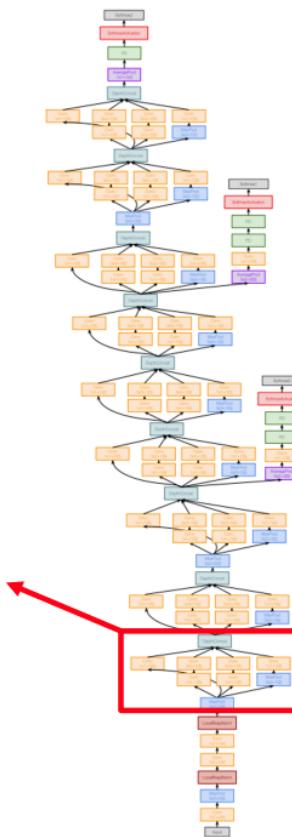
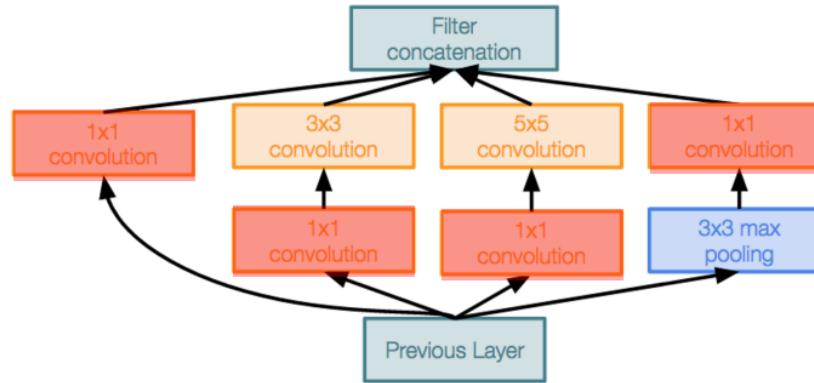
GoogLeNet: Inception Module

Inception module

Local unit with parallel branches

Local structure repeated many times throughout the network

Uses 1x1 “Bottleneck” layers to reduce channel dimension before expensive conv (we will revisit this with ResNet!)



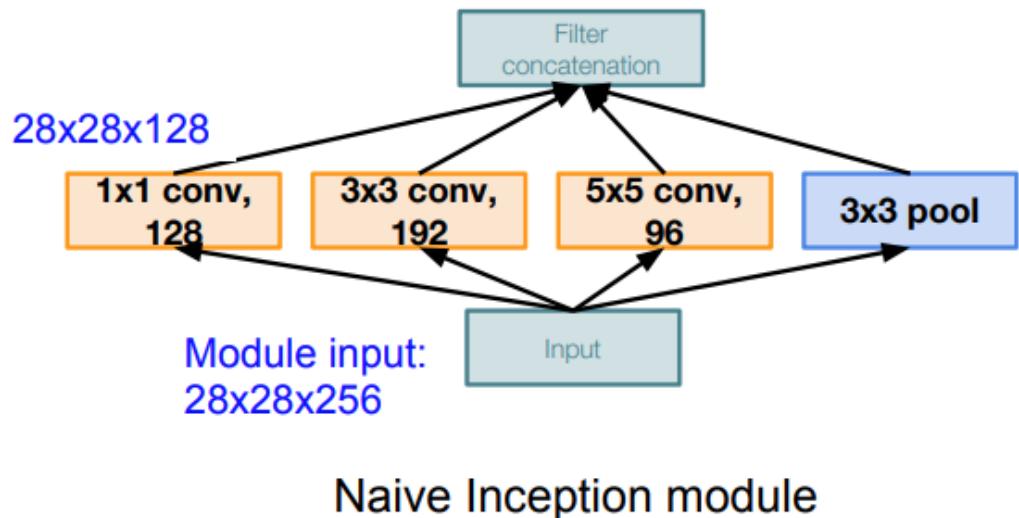
Case Study: GoogLeNet

[Szegedy et al., 2014]

Example:

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

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



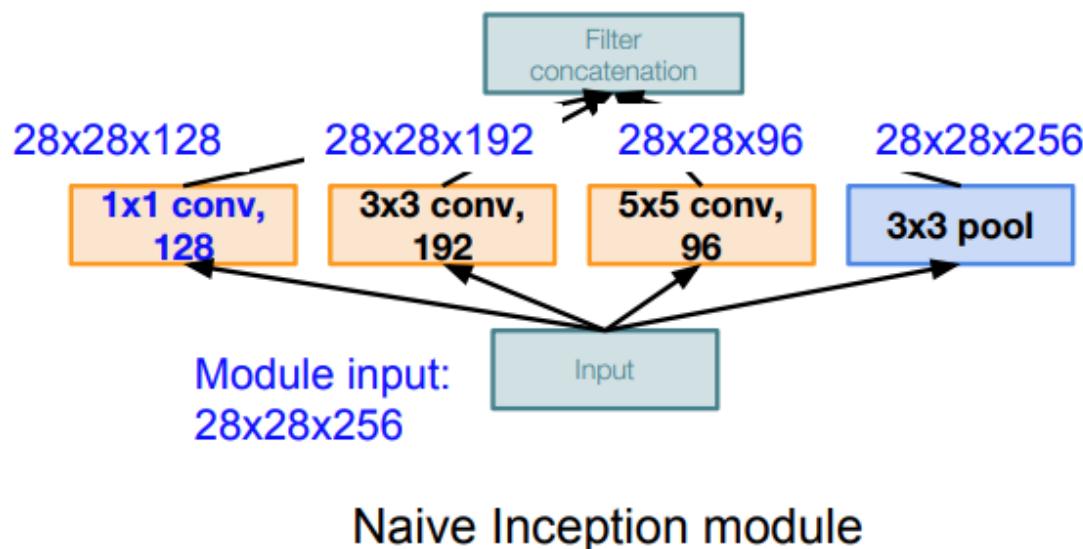
Case Study: GoogLeNet

[Szegedy et al., 2014]

Example:

Q3: What is output size after
filter concatenation?

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



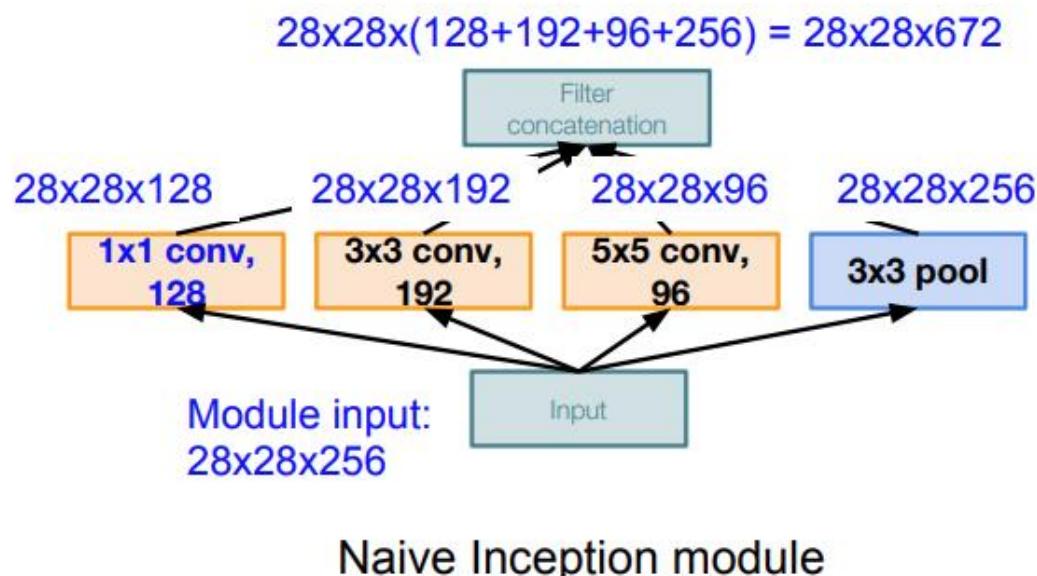
Q: What is the problem with this?
[Hint: Computational complexity]

Case Study: GoogLeNet

[Szegedy et al., 2014]

Example:

Q3: What is output size after
filter concatenation?



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

Conv Ops:

[1x1 conv, 128] 28x28x128x1x1x256
[3x3 conv, 192] 28x28x192x3x3x256
[5x5 conv, 96] 28x28x96x5x5x256

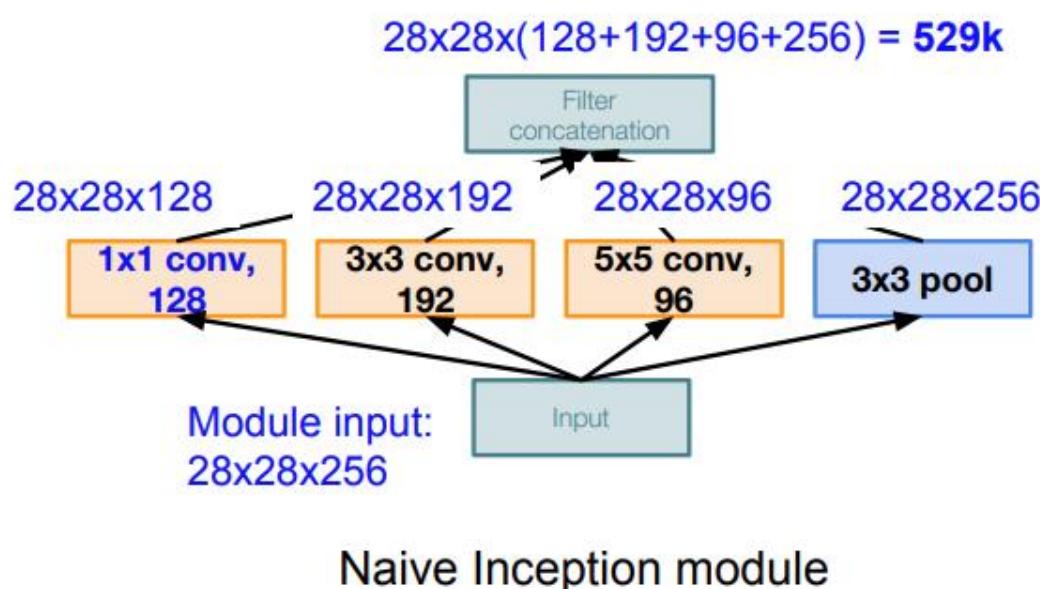
Total: 854M ops

Case Study: GoogLeNet

[Szegedy et al., 2014]

Example:

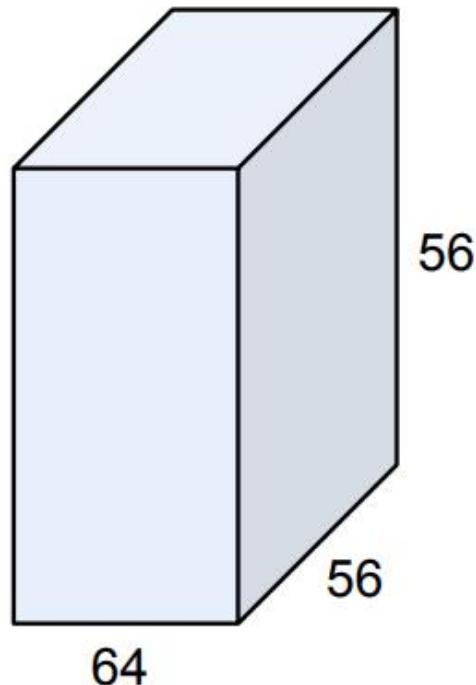
Q3: What is output size after filter concatenation?



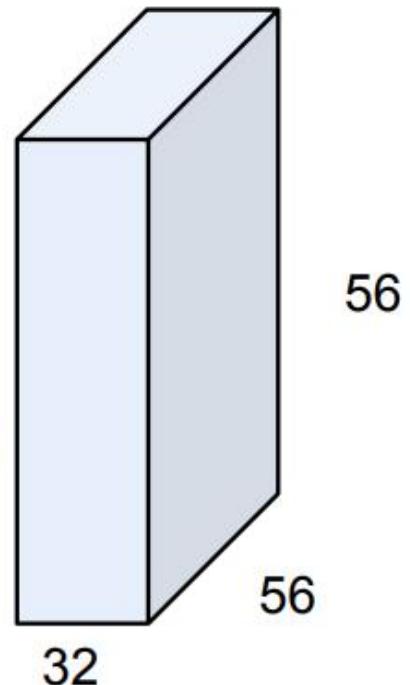
Q: What is the problem with this?
[Hint: Computational complexity]

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

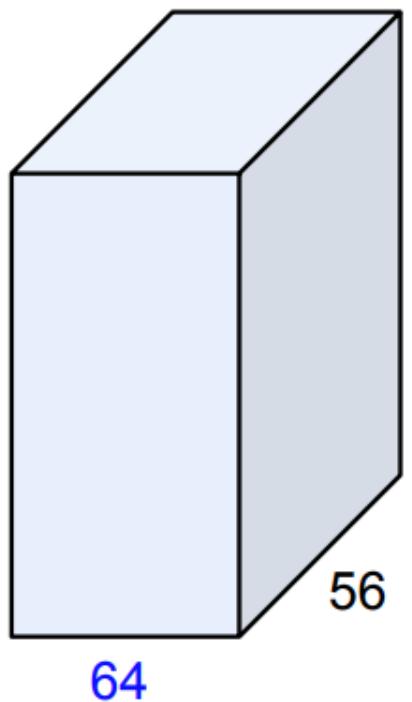
Reminder: 1x1 convolutions



1x1 CONV
with 32 filters
 $\xrightarrow{\hspace{1cm}}$
(each filter has size
1x1x64, and performs a
64-dimensional dot
product)



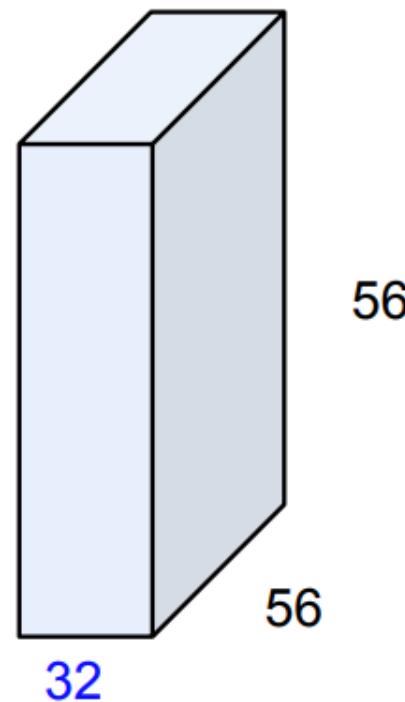
Reminder: 1x1 convolutions



1x1 CONV
with 32 filters

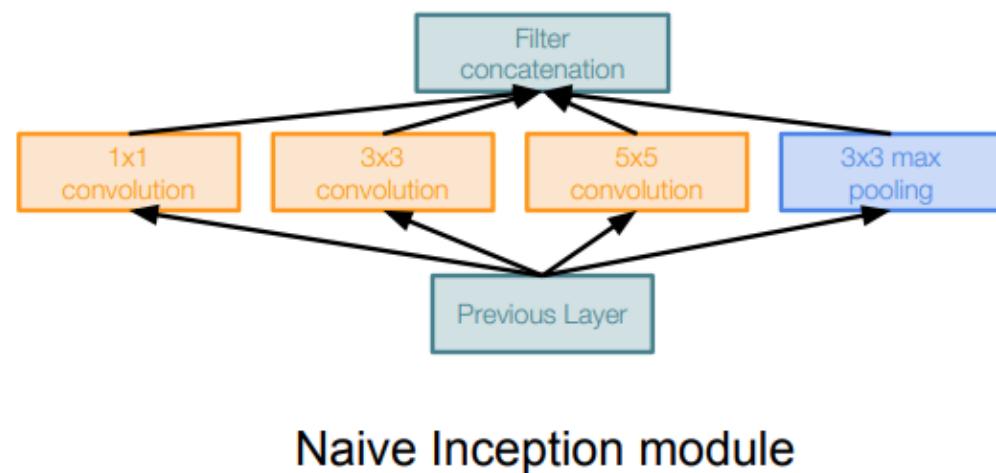
preserves spatial
dimensions, reduces depth!

Projects depth to lower
dimension (combination of
feature maps)

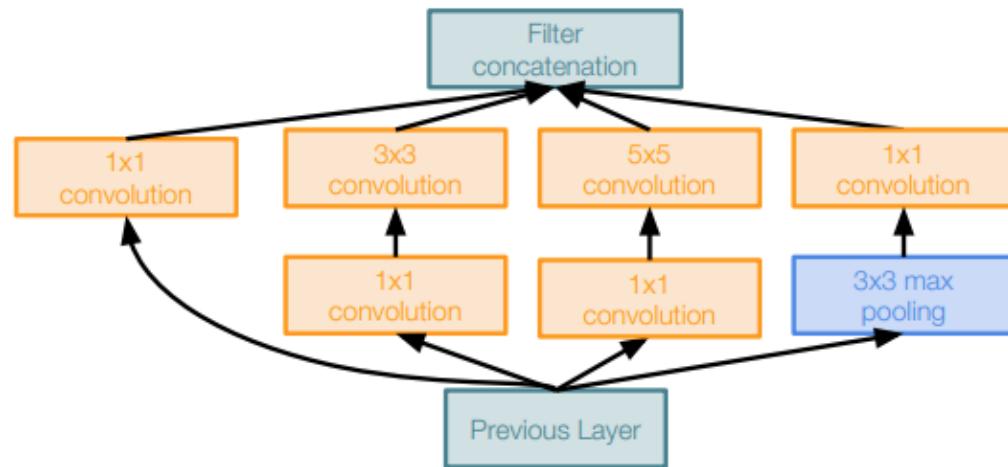


Case Study: GoogLeNet

[Szegedy et al., 2014]



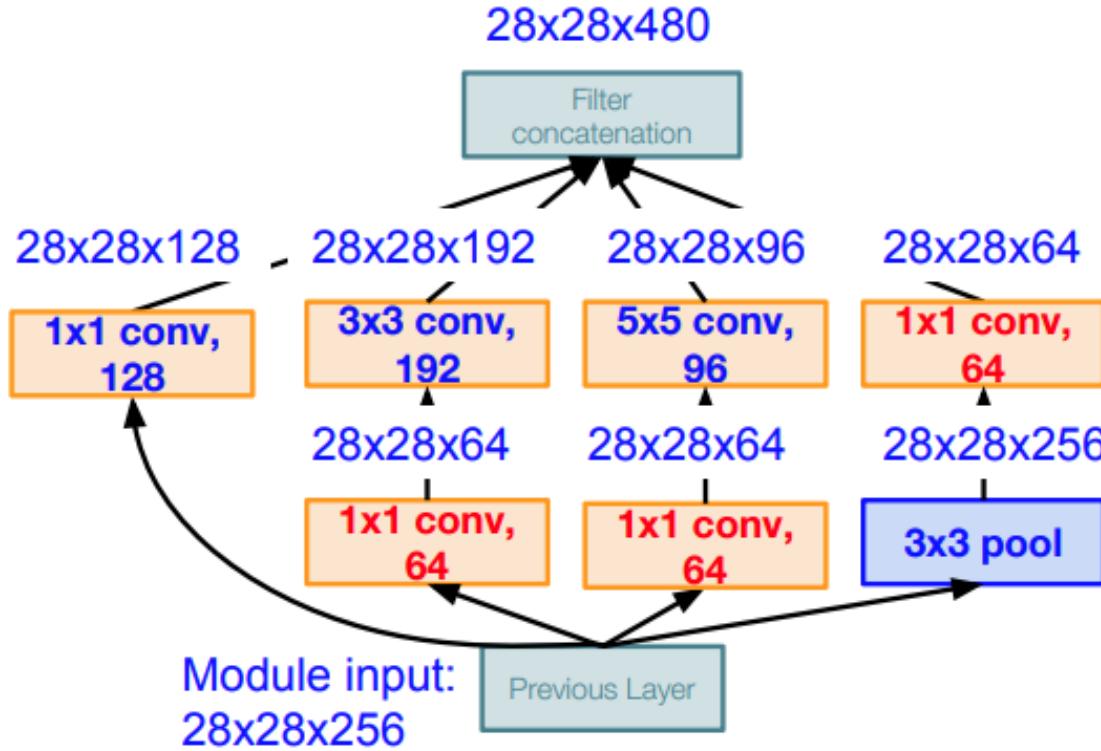
Naive Inception module



Inception module with dimension reduction

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

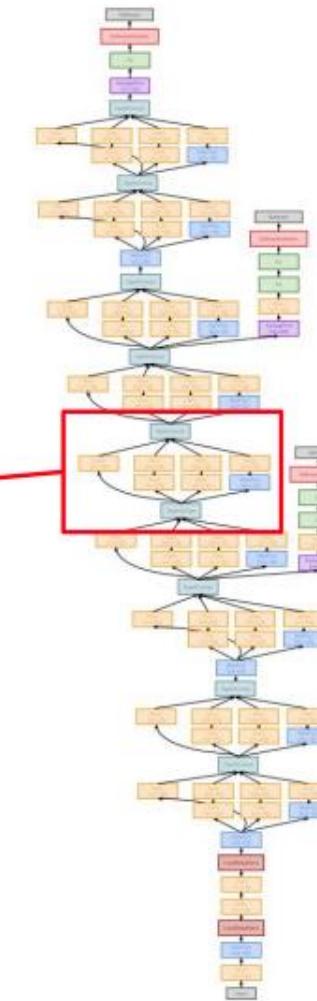
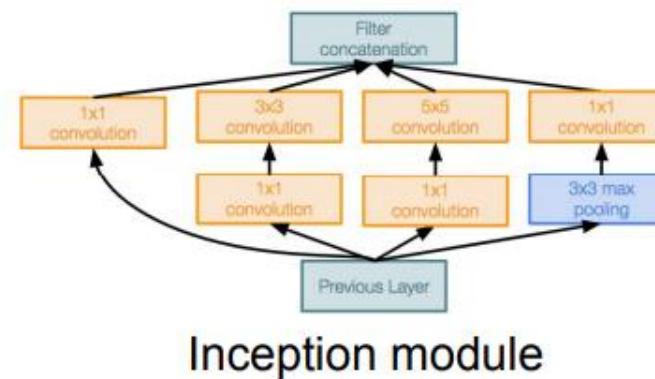
Total: 358M ops

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

Case Study: GoogLeNet

[Szegedy et al., 2014]

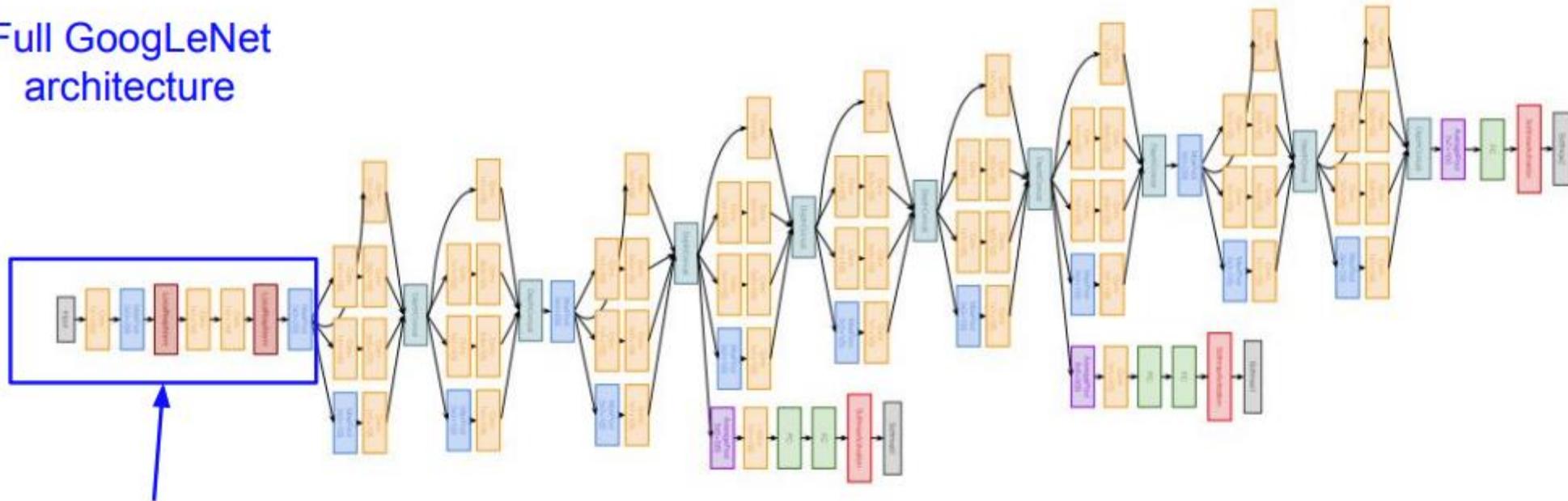
Stack Inception modules
with dimension reduction
on top of each other



Case Study: GoogLeNet

[Szegedy et al., 2014]

Full GoogLeNet
architecture

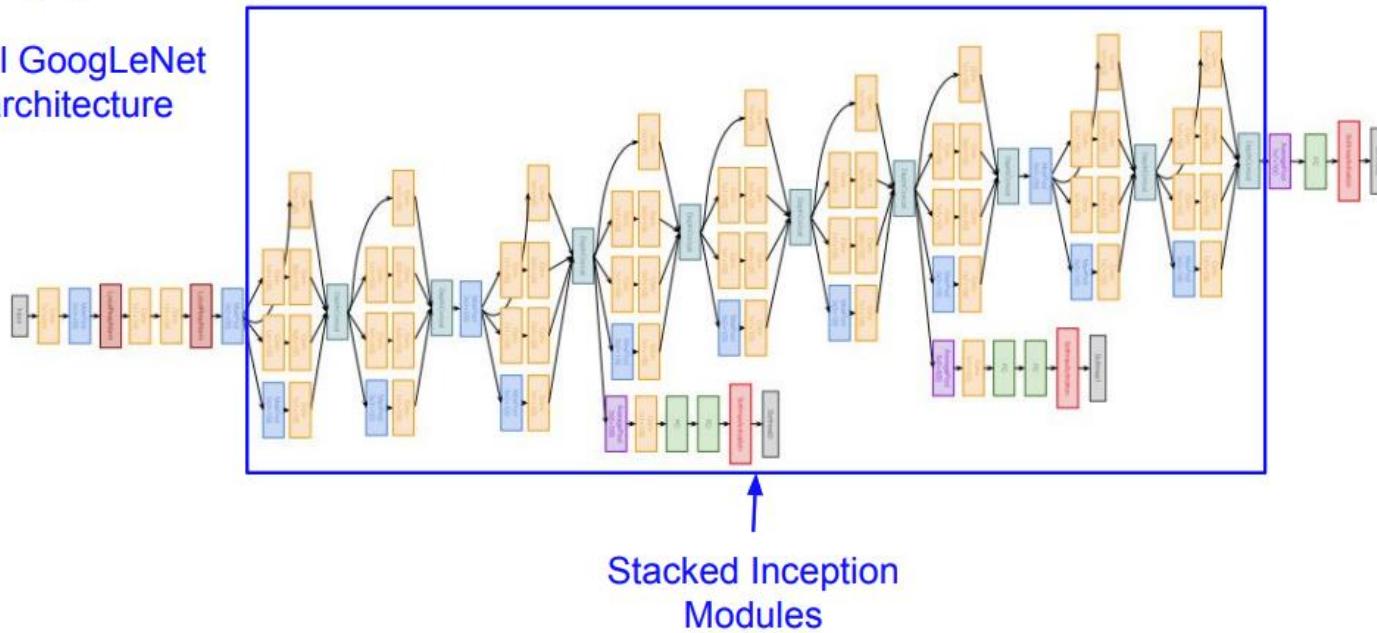


Stem Network:
Conv-Pool-
2x Conv-Pool

Case Study: GoogLeNet

[Szegedy et al., 2014]

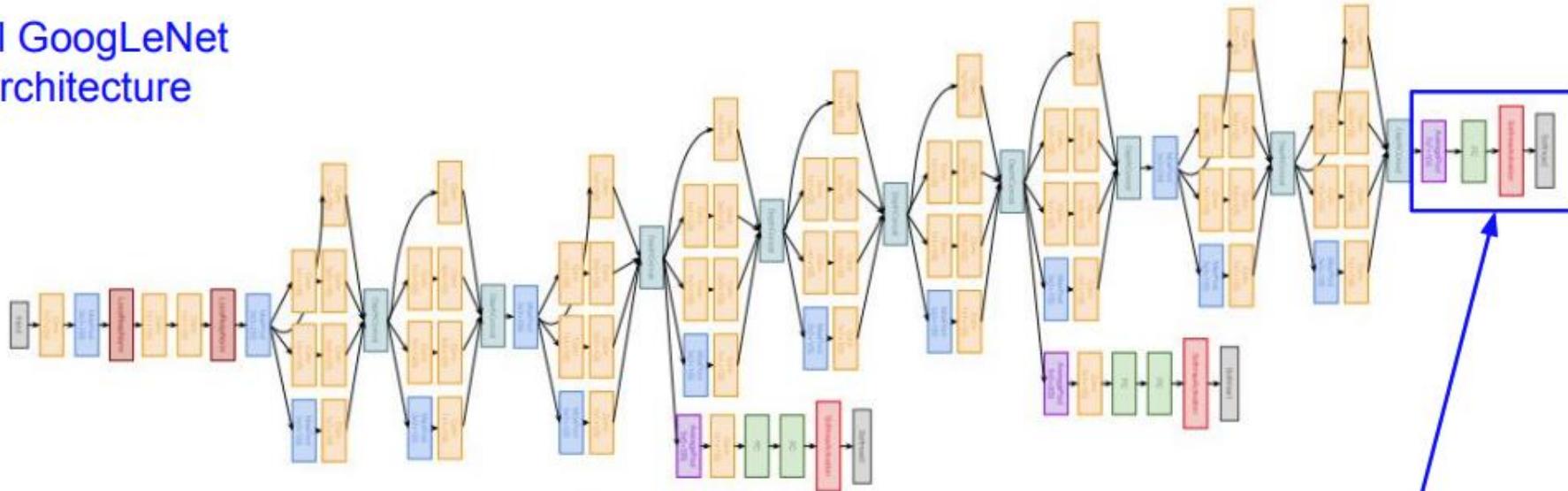
Full GoogLeNet architecture



Case Study: GoogLeNet

[Szegedy et al., 2014]

Full GoogLeNet
architecture

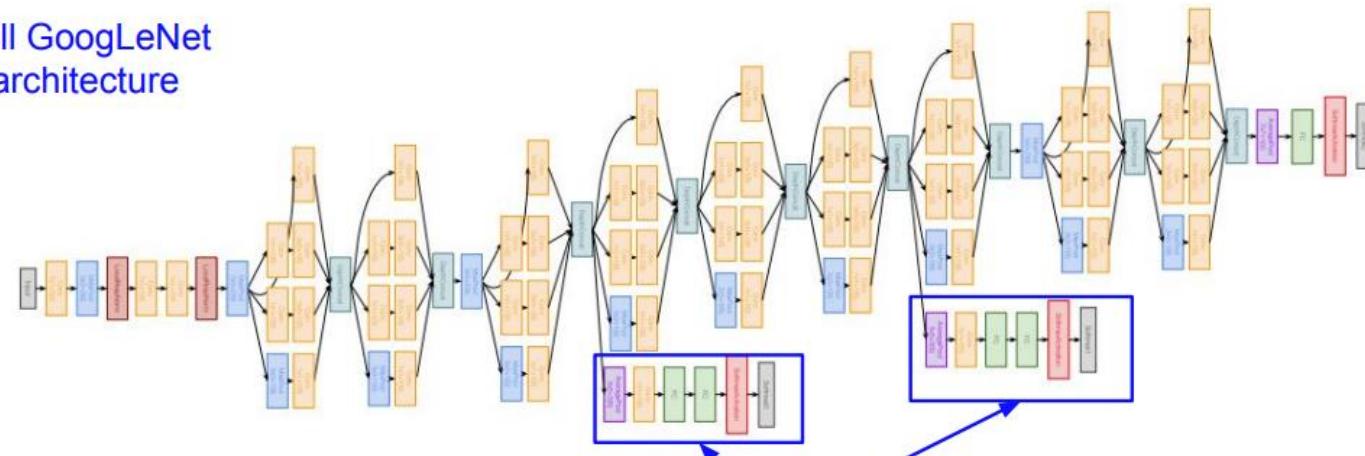


Classifier output
(removed expensive FC layers!)

Case Study: GoogLeNet

[Szegedy et al., 2014]

Full GoogLeNet
architecture



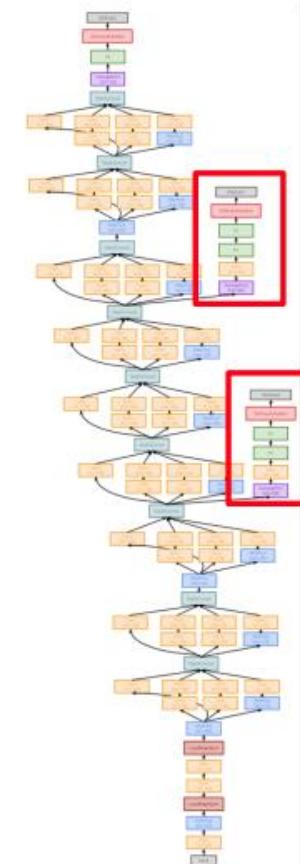
Auxiliary classification outputs to inject additional gradient at lower layers
(AvgPool-1x1Conv-FC-FC-Softmax)

GoogLeNet: Auxiliary Classifiers

Training using loss at the end of the network didn't work well:
Network is too deep, gradients don't propagate cleanly

As a hack, attach “auxiliary classifiers” at several intermediate points
in the network that also try to classify the image and receive loss

GoogLeNet was before batch normalization! With BatchNorm no
longer need to use this trick

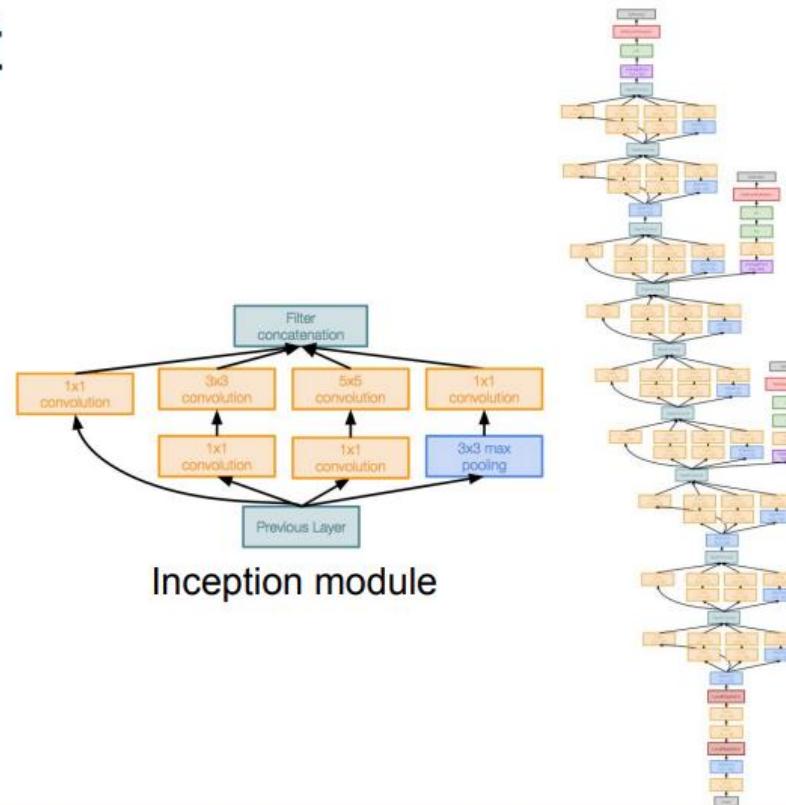


Case Study: GoogLeNet

[Szegedy et al., 2014]

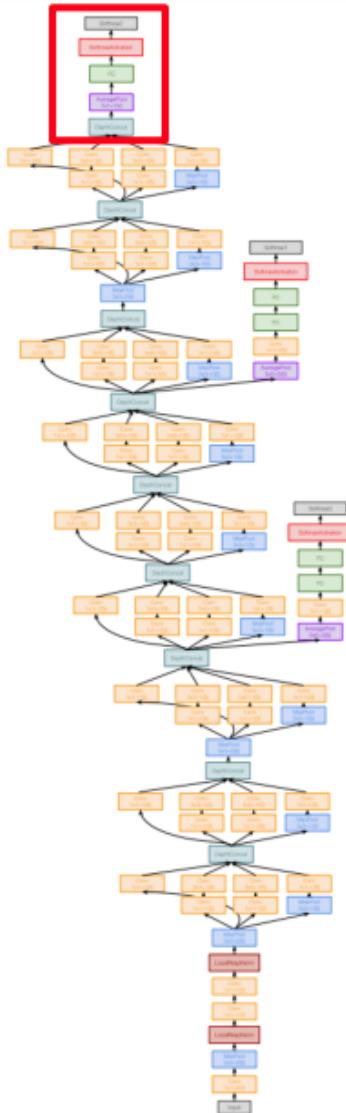
Deeper networks, with computational efficiency

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



GoogLeNet: Global Average Pooling

No large FC layers at the end! Instead uses “global average pooling” to collapse spatial dimensions, and one linear layer to produce class scores
(Recall VGG-16: Most parameters were in the FC layers!)



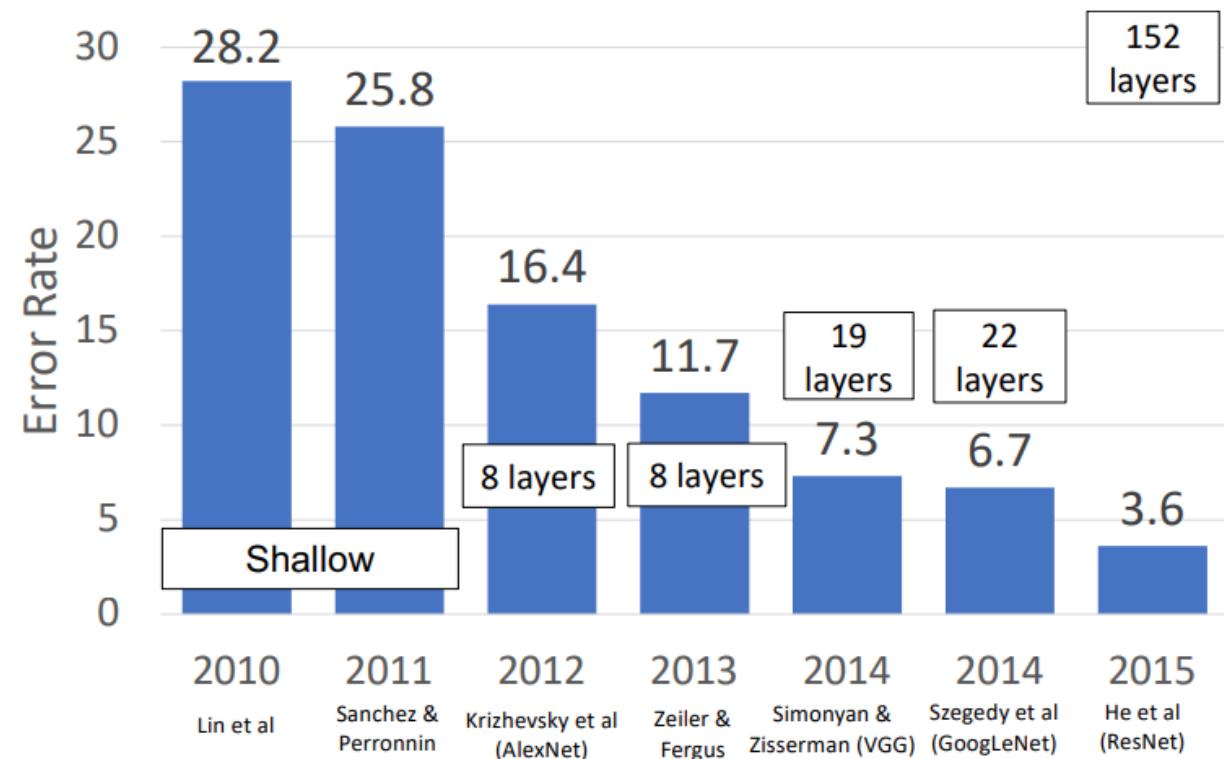
	Input size		Layer				Output size				
Layer	C	H/W	filters	kernel	stride	pad	C	H/W	memory (KB)	params (k)	flop (M)
avg-pool	1024	7		7	1	0	1024	1	4	0	0
fc	1024		1000				1000		0	1025	1

Compare with VGG-16:

Layer	C	H/W	filters	kernel	stride	pad	C	H/W	memory (KB)	params (K)	flop (M)
flatten	512	7					25088		98		
fc6	25088			4096			4096		16	102760	103
fc7	4096			4096			4096		16	16777	17
fc8	4096			1000			1000		4	4096	4

- **Global Average Pooling** is a pooling operation designed to replace fully connected layers in classical CNNs. The idea is to generate one feature map for each corresponding category of the classification task in the last mlpconv layer. Instead of adding fully connected layers on top of the feature maps, we take the average of each feature map, and the resulting vector is fed directly into the softmax layer.
- One advantage of global average pooling over the fully connected layers is that it is more native to the convolution structure by enforcing correspondences between feature maps and categories. Thus the feature maps can be easily interpreted as categories confidence maps. Another advantage is that there is no parameter to optimize in the global average pooling thus overfitting is avoided at this layer. Furthermore, global average pooling sums out the spatial information, thus it is more robust to spatial translations of the input.

ImageNet Classification Challenge

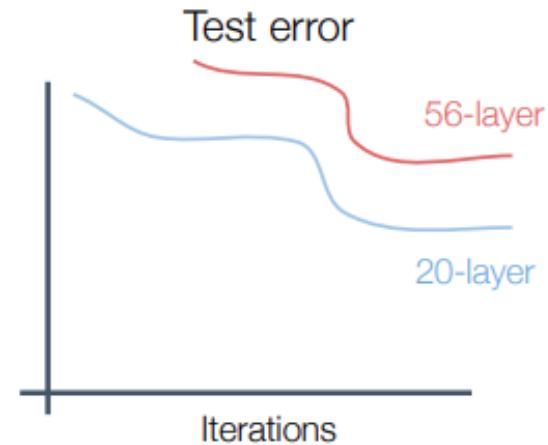


Residual Networks

Once we have Batch Normalization, we can train networks with 10+ layers.
What happens as we go deeper?

Deeper model does worse than shallow model!

Initial guess: Deep model is **overfitting** since it is much bigger than the other model



Deep Residual Learning for Image Recognition

Kaiming He Xiangyu Zhang Shaoqing Ren Jian Sun
Microsoft Research
`{kahe, v-xiangz, v-shren, jiansun}@microsoft.com`

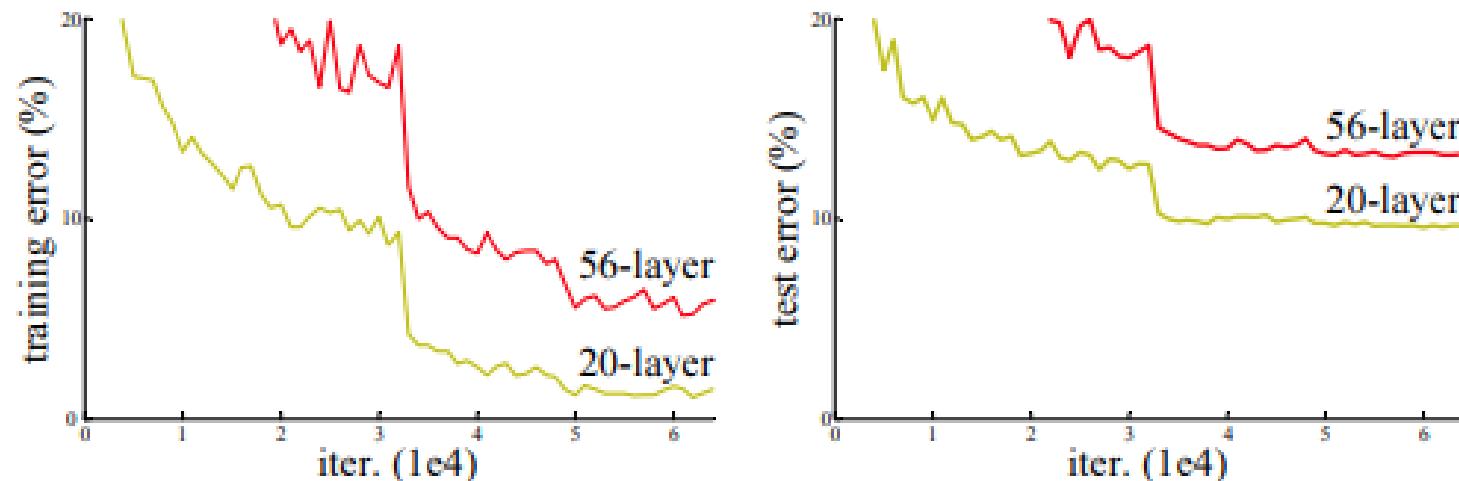
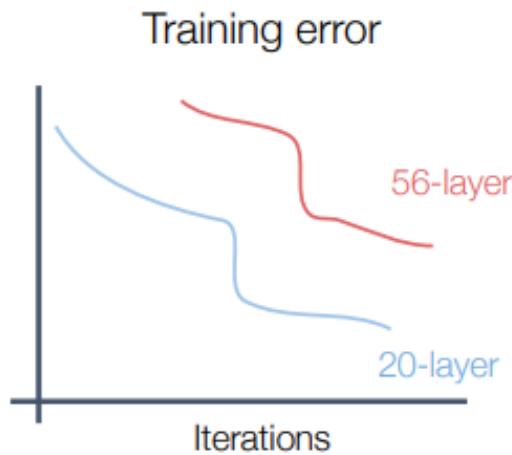


Figure 1. Training error (left) and test error (right) on CIFAR-10 with 20-layer and 56-layer “plain” networks. The deeper network has higher training error, and thus test error. Similar phenomena on ImageNet is presented in Fig. 4.

He, K., Zhang, X., Ren, S., & Sun, J. (2016). Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition* (pp. 770-778).

Residual Networks

Once we have Batch Normalization, we can train networks with 10+ layers.
What happens as we go deeper?



In fact the deep model seems to be **underfitting** since it also performs worse than the shallow model on the training set! It is actually **underfitting**

Case Study: ResNet

[He et al., 2015]

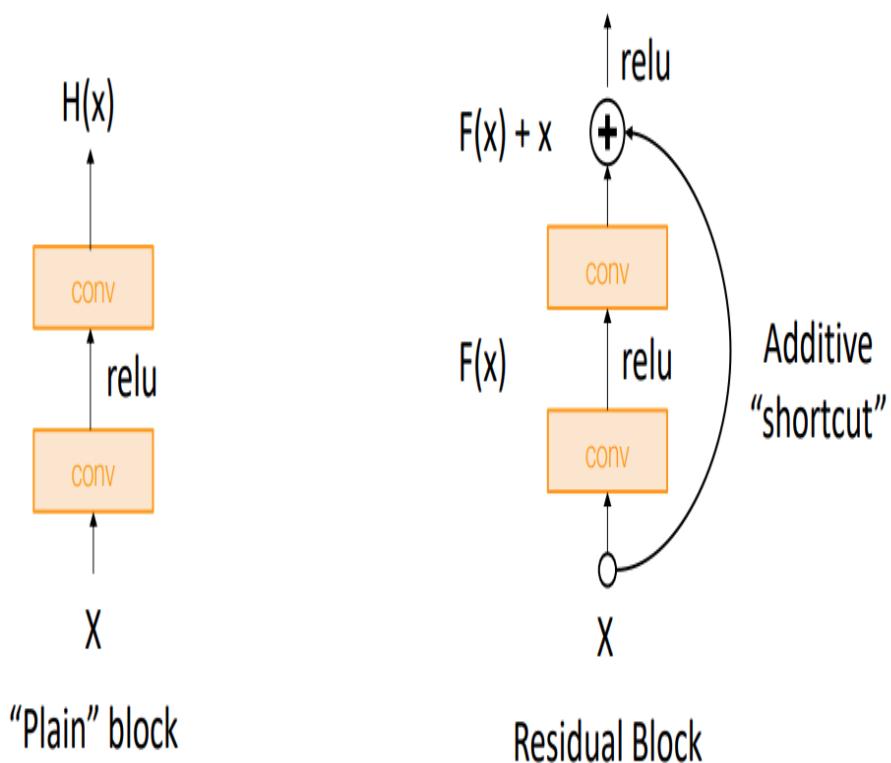
Hypothesis: the problem is an *optimization* problem, deeper models are harder to optimize

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.

Residual Networks

Solution: Change the network so learning identity functions with extra layers is easy!



Deep residual networks consist of many stacked "Residual Units". Each unit can be expressed in a general form:

$$y_l = h(x_l) + F(x_l, W_l)$$
$$x_{l+1} = f(y_l)$$

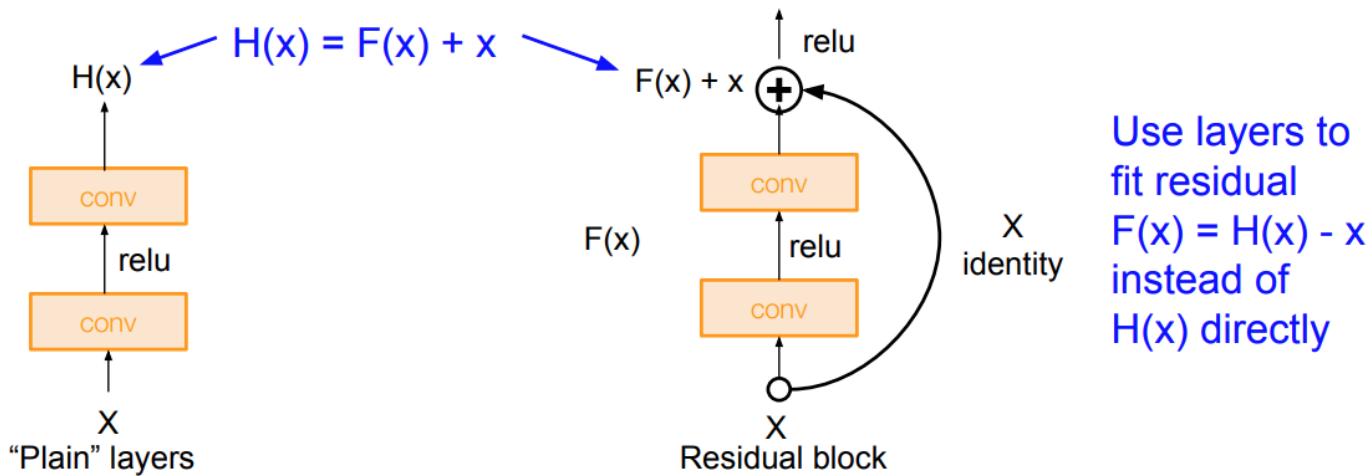
where x_l and x_{l+1} are input and output of the $l - th$ unit, and F is a residual function. In the last paper, $h(x_l) = x_l$ is an identity mapping and f is a ReLU function.

The central idea of ResNets is to learn the additive residual function F with respect to $h(x_l)$, with a key choice of using an identity mapping $h(x_l) = x_l$. This is realized by attaching an identity skip connection ("shortcut").

Case Study: ResNet

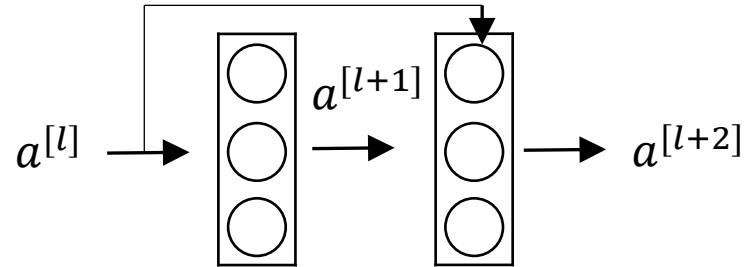
[He et al., 2015]

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



If the identity mapping $f(x)=x$ is the desired underlying mapping, the residual mapping amounts to $g(x)=0$ and it is thus easier to learn: we only need to push the weights and biases of the upper weight layer (e.g., fully connected layer and convolutional layer) within the dotted-line box to zero.

Residual Blocks



$$z^{[l+1]} = W^{[l+1]} a^{[l]} + b^{[l+1]}$$

“linear”

$$a^{[l+1]} = g(z^{[l+1]})$$

“relu”

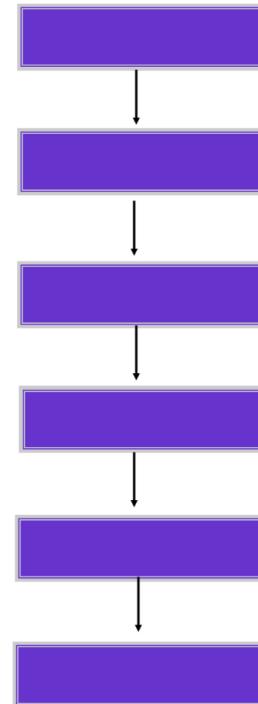
$$z^{[l+2]} = W^{[l+2]} a^{[l+1]} + b^{[l+2]}$$

“output”

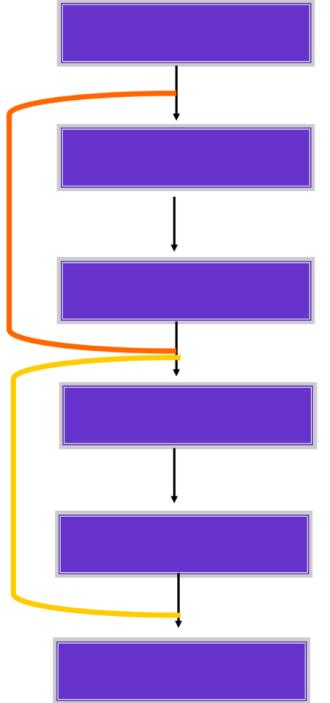
$$a^{[l+2]} = g(z^{[l+2]} + a^{[l]})$$

“relu on (output **plus input**)”

PLAIN CONNECTIONS



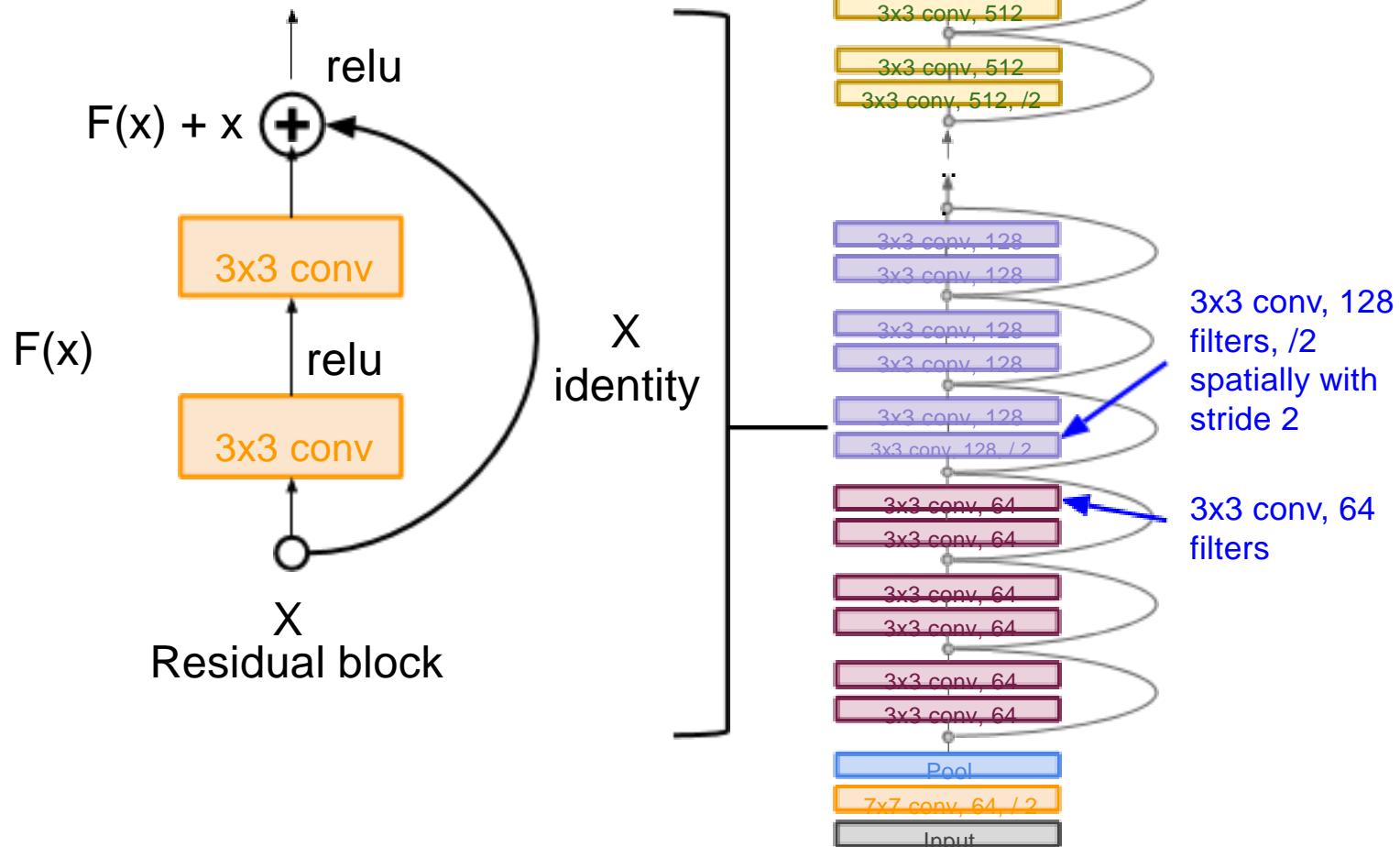
RESIDUAL CONNECTIONS



ResNet Architecture

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)



Intro

ResNet

Technical details

Results

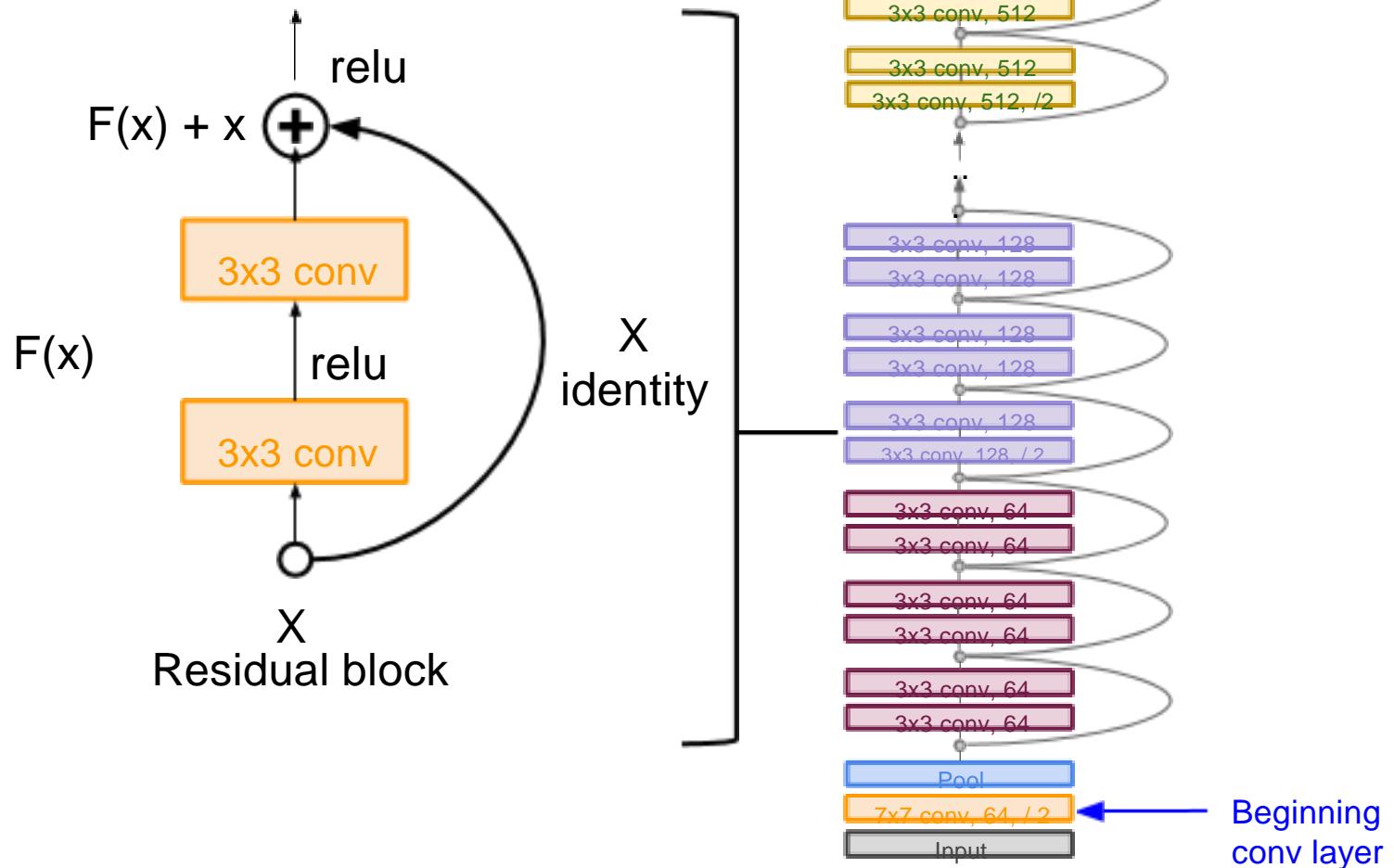
ResNet 1000

Comparison

ResNet Architecture

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



Intro

Results

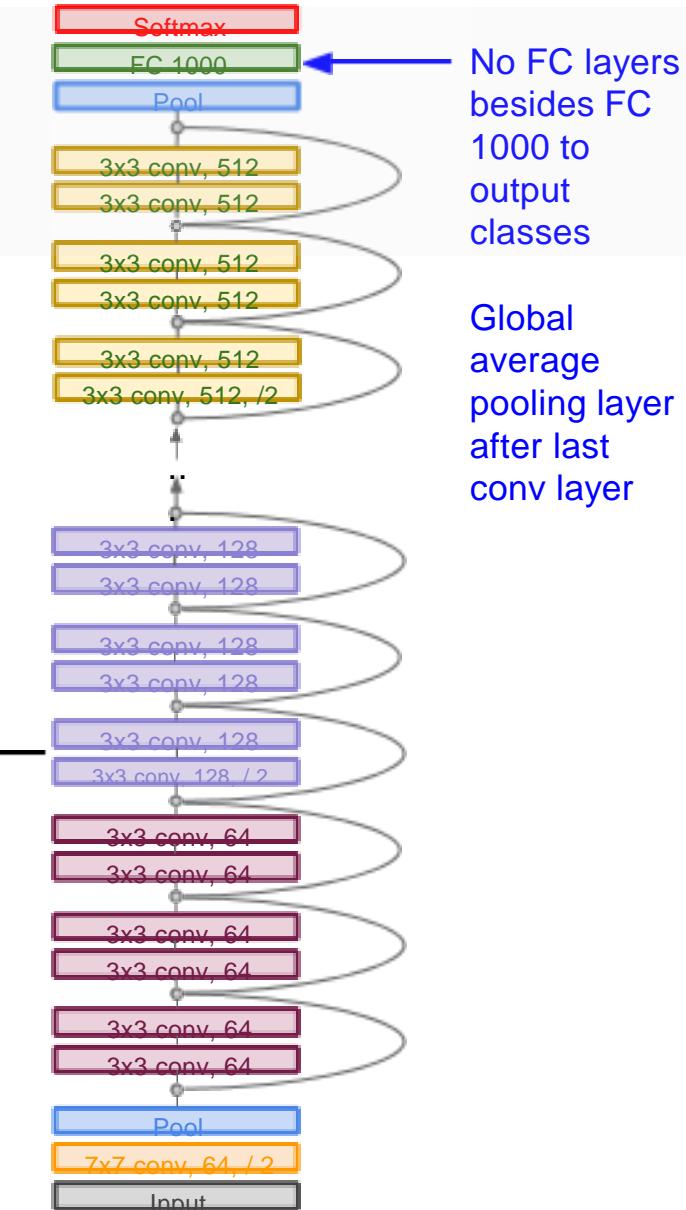
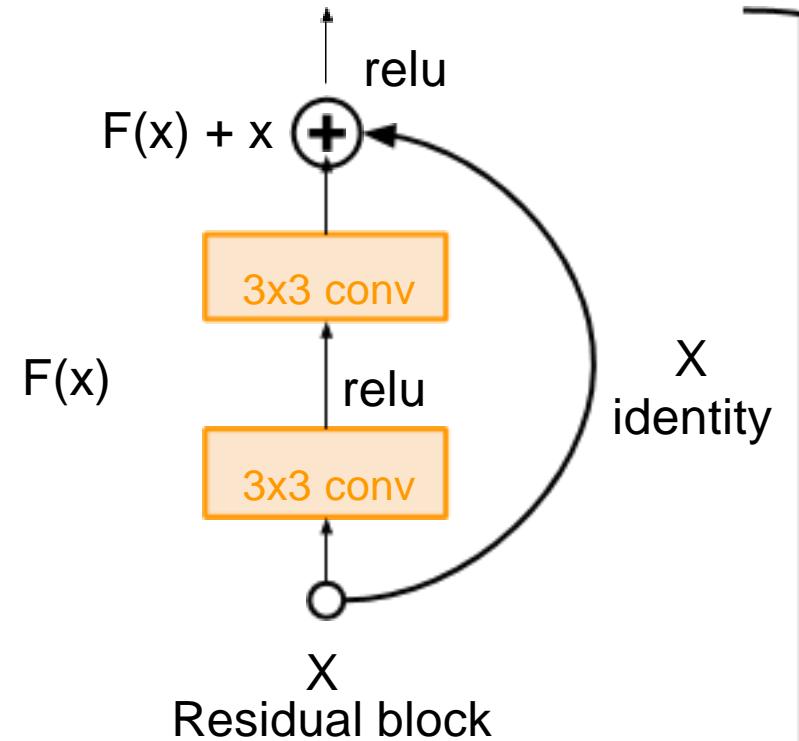
ResNet 1000

Comparison

ResNet Architecture

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)



Residual Networks

ResNet-18:

Stem: 1 conv layer

Stage 1 (C=64): 2 res. block = 4 conv

Stage 2 (C=128): 2 res. block = 4 conv

Stage 3 (C=256): 2 res. block = 4 conv

Stage 4 (C=512): 2 res. block = 4 conv

Linear

ImageNet top-5 error: 10.92

GFLOP: 1.8

ResNet-34:

Stem: 1 conv layer

Stage 1: 3 res. block = 6 conv

Stage 2: 4 res. block = 8 conv

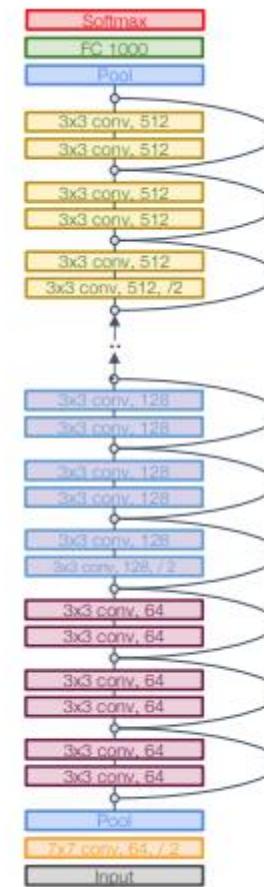
Stage 3: 6 res. block = 12 conv

Stage 4: 3 res. block = 6 conv

Linear

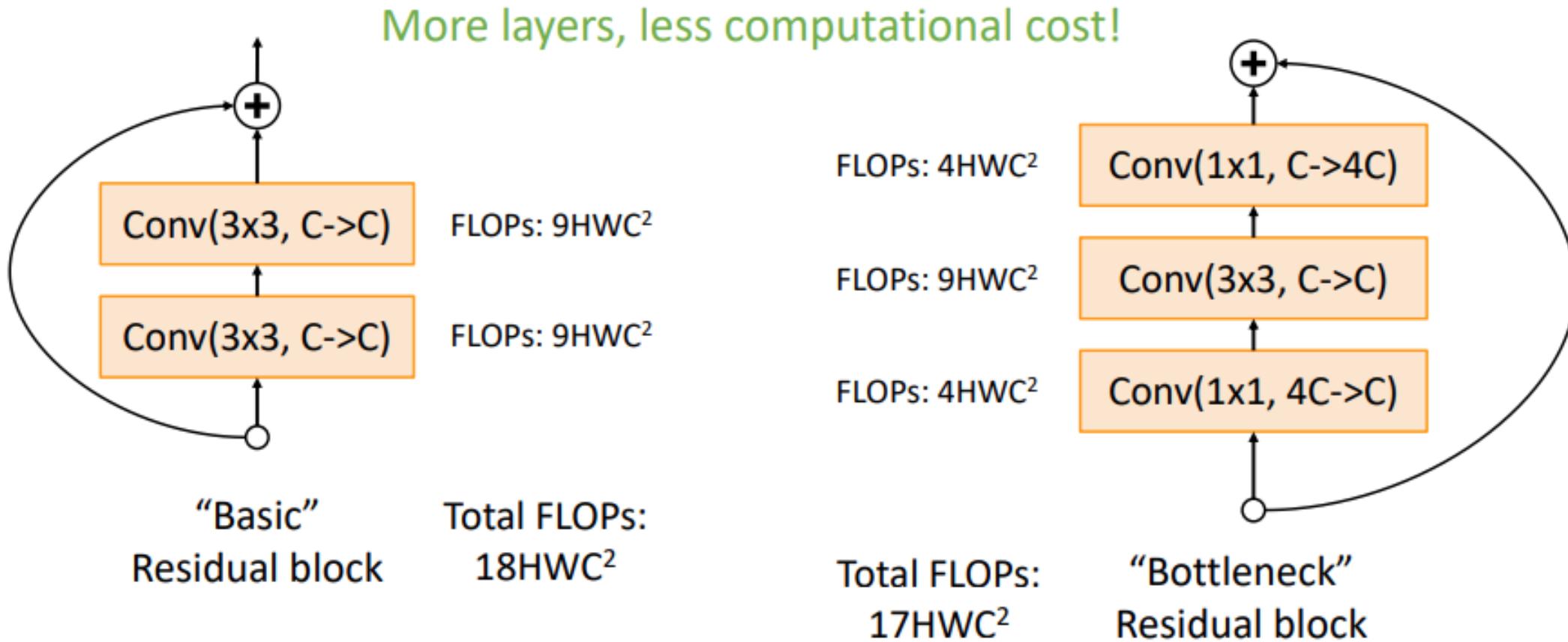
ImageNet top-5 error: 8.58

GFLOP: 3.6



He et al, "Deep Residual Learning for Image Recognition", CVPR 2016
Error rates are 224x224 single-crop testing, reported by [torchvision](#)

Residual Networks: Bottleneck Block



Residual Networks

Deeper ResNet-101 and ResNet-152 models are more accurate, but also more computationally heavy

	Block type	Stem layers	Stage 1			Stage 2			Stage 3			Stage 4			FC layers	ImageNet top-5 error
			Blocks	Layers	Blocks	Layers	Blocks	Layers	Blocks	Layers	Blocks	Layers	Blocks	Layers		
ResNet-18	Basic	1	2	4	2	4	2	4	2	4	1	1.8	10.92			
ResNet-34	Basic	1	3	6	4	8	6	12	3	6	1	3.6	8.58			
ResNet-50	Bottle	1	3	9	4	12	6	18	3	9	1	3.8	7.13			
ResNet-101	Bottle	1	3	9	4	12	23	69	3	9	1	7.6	6.44			
ResNet-152	Bottle	1	3	9	8	24	36	108	3	9	1	11.3	5.94			

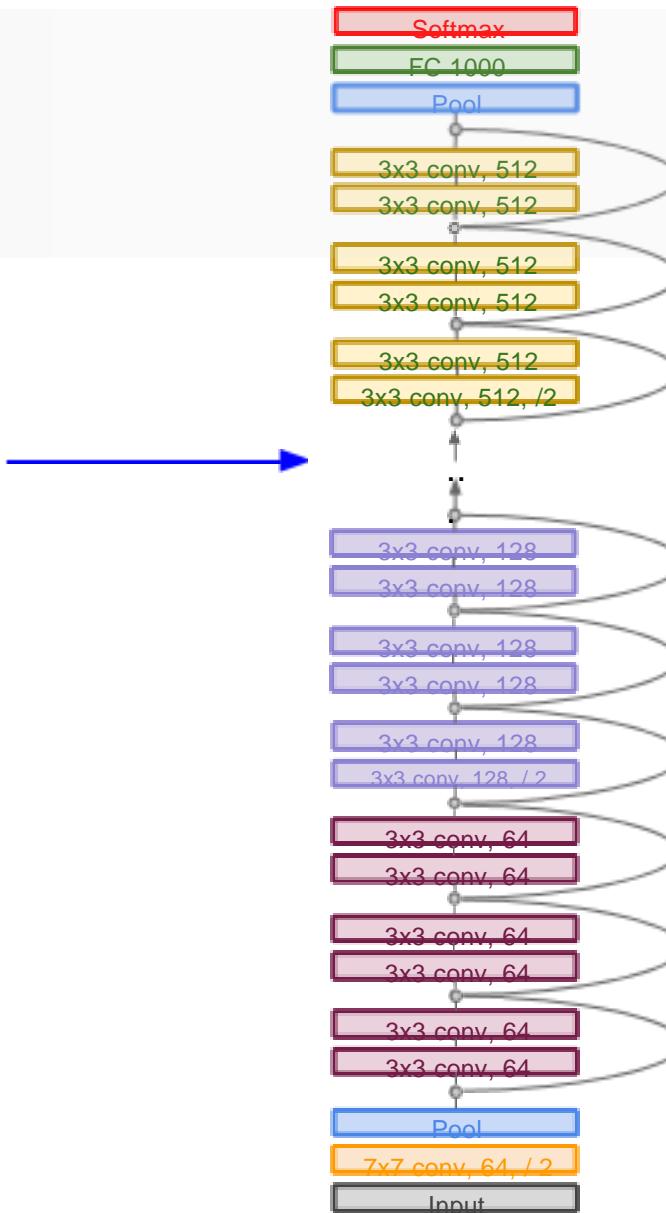


He et al, "Deep Residual Learning for Image Recognition", CVPR 2016

Error rates are 224x224 single-crop testing, reported by [torchvision](#)

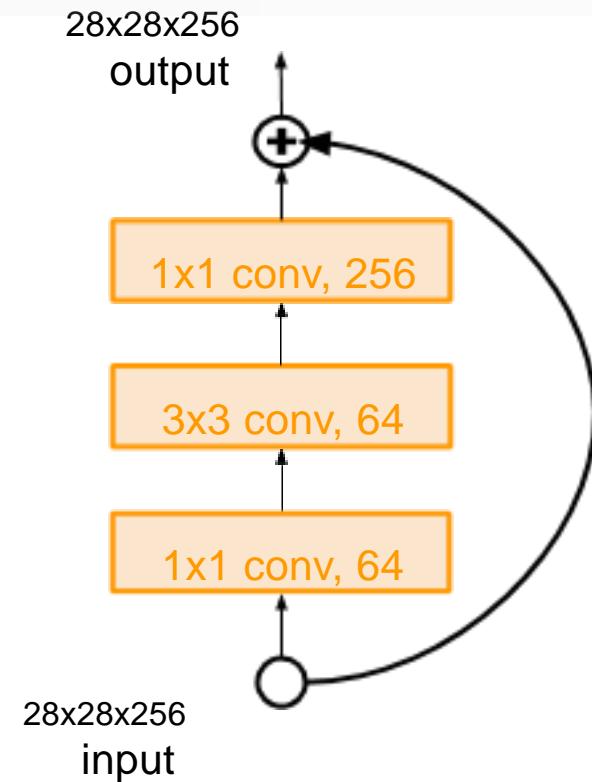
ResNet Architecture

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



ResNet Architecture

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



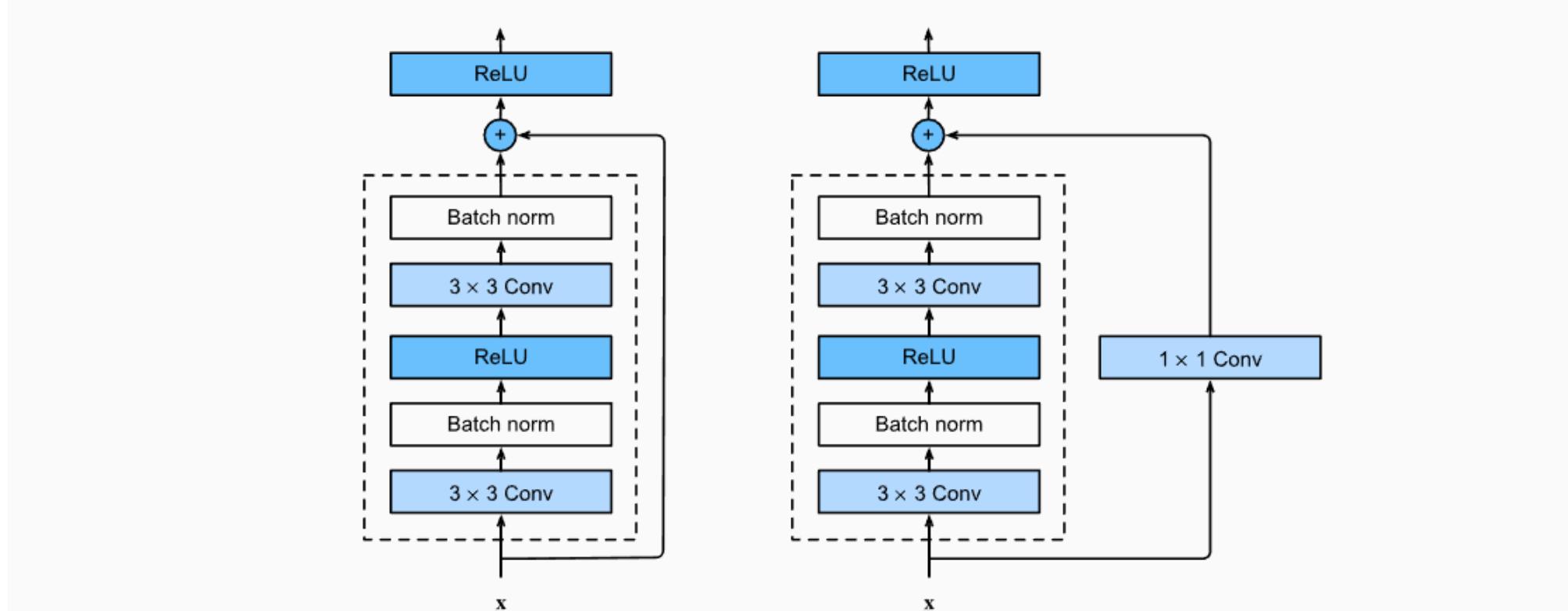


Fig. 8.6.3 ResNet block with and without 1×1 convolution, which transforms the input into the desired shape for the addition operation.

Zhang, A., Lipton, Z. C., Li, M., & Smola, A. J. (2021). Dive into deep learning. *arXiv preprint arXiv:2106.11342*.

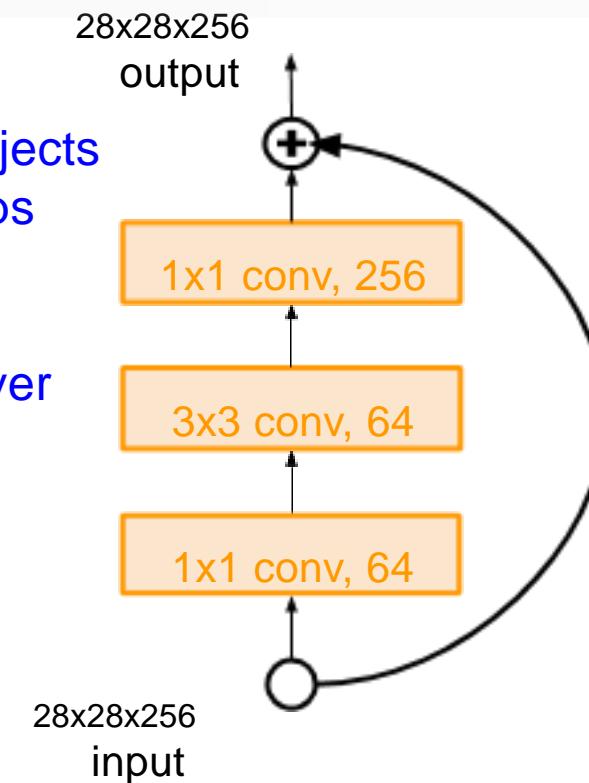
ResNet Architecture

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

1x1 conv, 256 filters projects
back to 256 feature maps
(28x28x256)

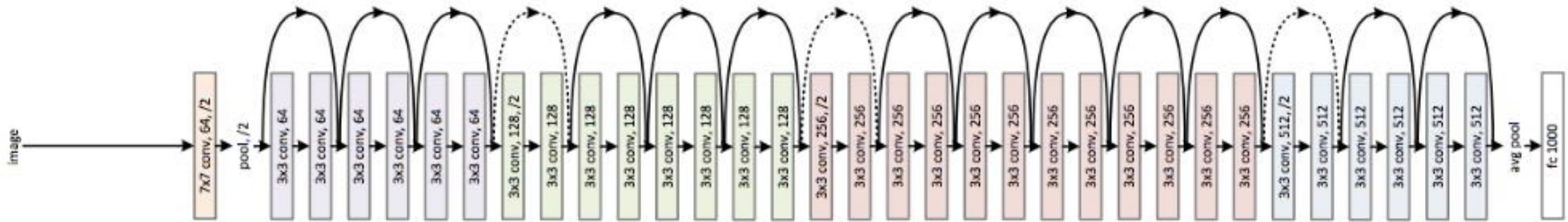
3x3 conv operates over
only 64 feature maps

1x1 conv, 64 filters
to project to
28x28x64



Residual Blocks (skip connections)

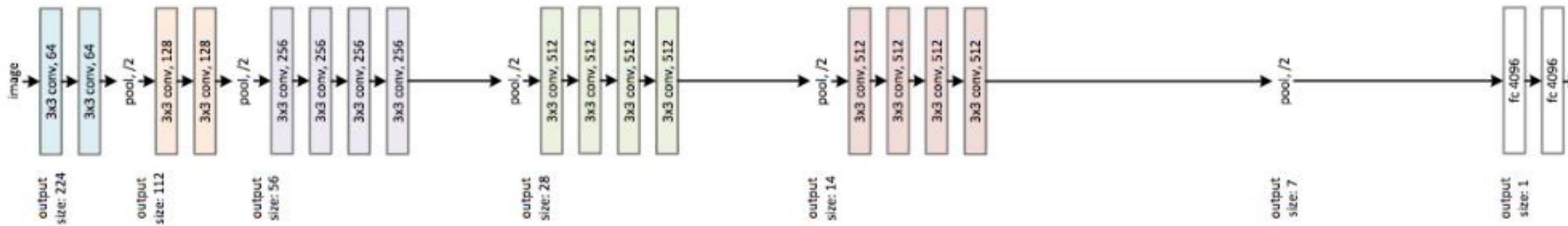
34-layer residual



34-layer plain



VGG-19



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.

Case Study: ResNet

[He et al., 2015]

ILSVRC 2015 winner (3.6% top 5 error)



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

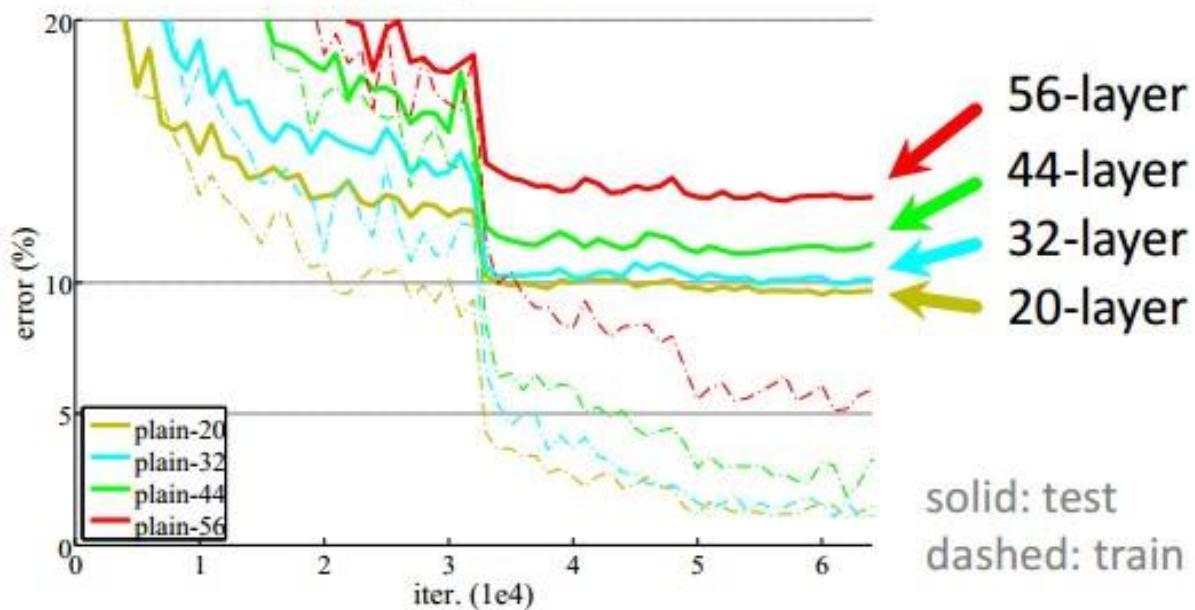


Kaiming He, Xiangyu Zhang, Shaoqing Ren, & Jian Sun. “Deep Residual Learning for Image Recognition”. arXiv 2015.

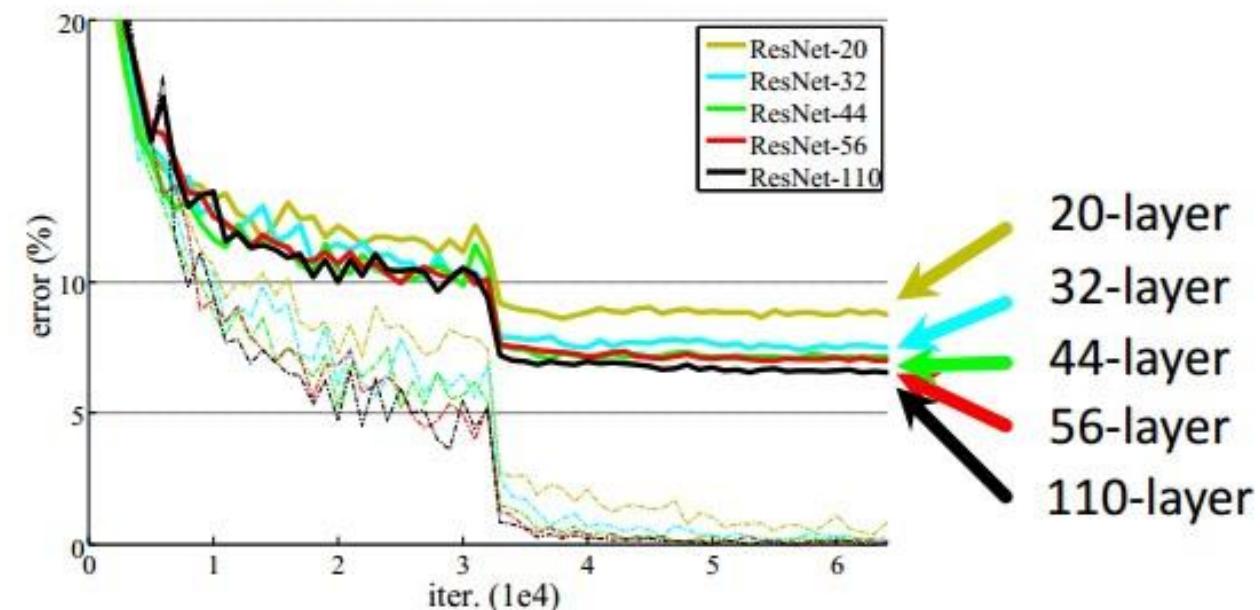
Slide from Kaiming He’s recent presentation <https://www.youtube.com/watch?v=1PGLj-uKT1w>

CIFAR-10 experiments

CIFAR-10 plain nets



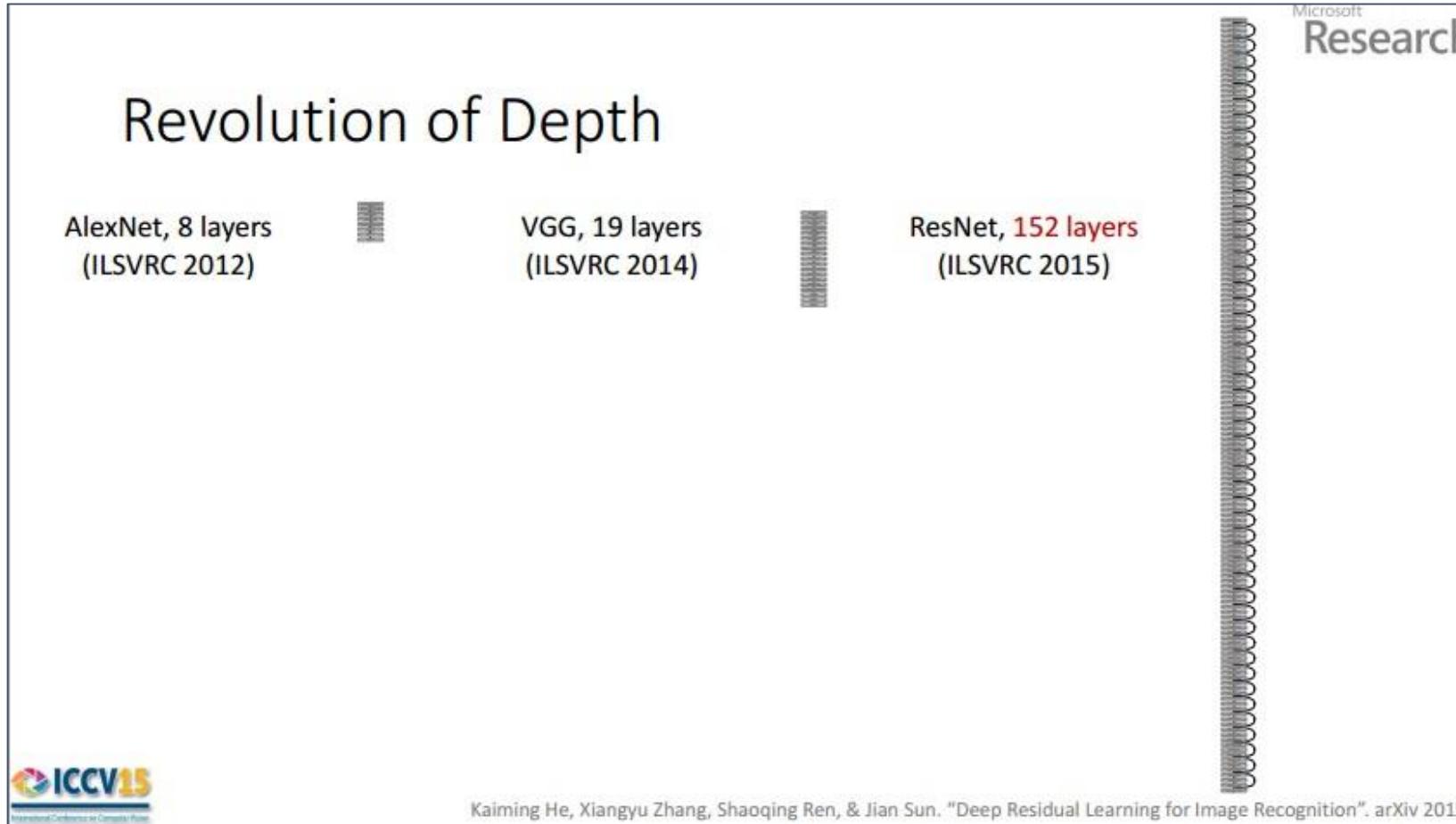
CIFAR-10 ResNets



Case Study: ResNet

[He et al., 2015]

ILSVRC 2015 winner (3.6% top 5 error)



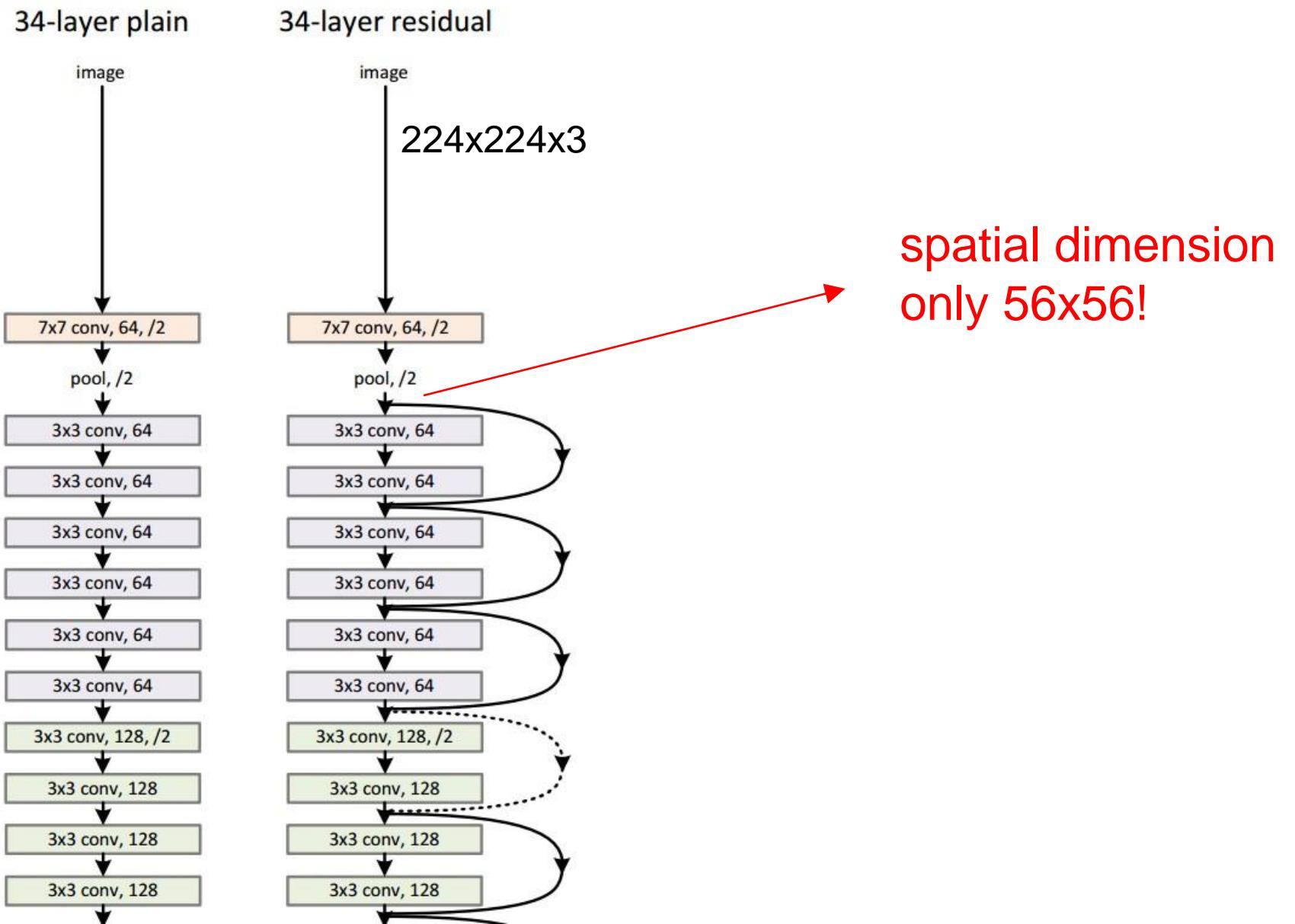
(slide from Kaiming He's recent presentation)

2-3 weeks of training
on 8 GPU machine

at runtime: faster
than a VGGNet!
(even though it has
8x more layers)

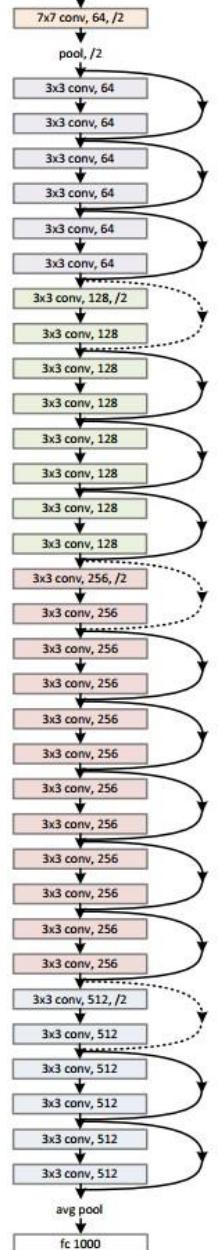
Case Study: ResNet

[He et al., 2015]



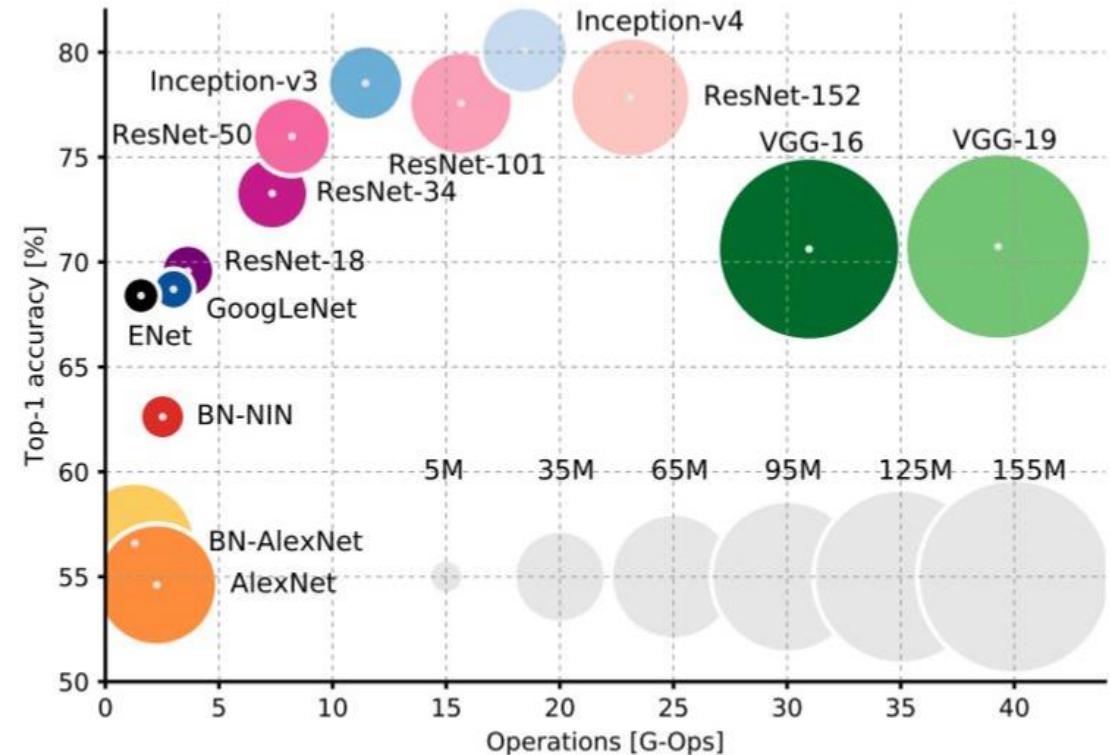
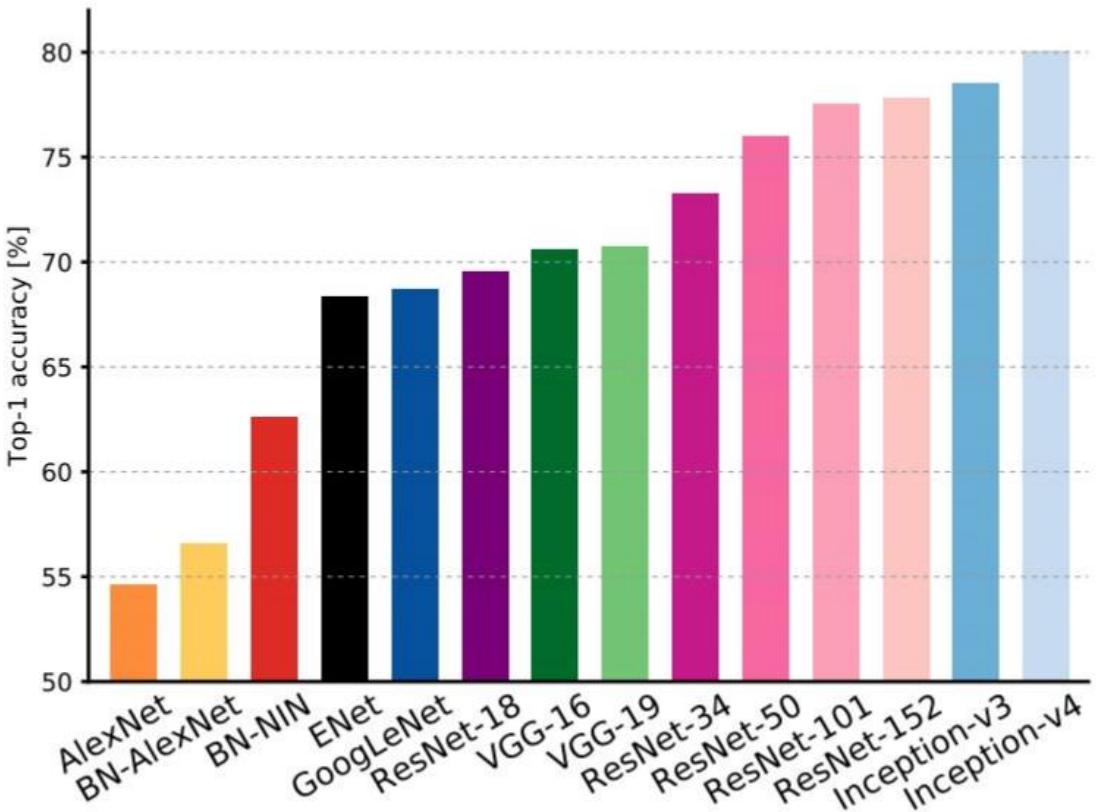
Layer Name	Output Size	ResNet-18
<code>conv1</code>	$112 \times 112 \times 64$	$7 \times 7, 64$, stride 2
<code>conv2_x</code>	$56 \times 56 \times 64$	3×3 max pool, stride 2 $\left[\begin{array}{l} 3 \times 3, 64 \\ 3 \times 3, 64 \end{array} \right] \times 2$
<code>conv3_x</code>	$28 \times 28 \times 128$	$\left[\begin{array}{l} 3 \times 3, 128 \\ 3 \times 3, 128 \end{array} \right] \times 2$
<code>conv4_x</code>	$14 \times 14 \times 256$	$\left[\begin{array}{l} 3 \times 3, 256 \\ 3 \times 3, 256 \end{array} \right] \times 2$
<code>conv5_x</code>	$7 \times 7 \times 512$	$\left[\begin{array}{l} 3 \times 3, 512 \\ 3 \times 3, 512 \end{array} \right] \times 2$
<code>average pool</code>	$1 \times 1 \times 512$	7×7 average pool
<code>fully connected</code>	1000	512×1000 fully connections
<code>softmax</code>	1000	

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



layer name	output size	18-layer	34-layer	50-layer	101-layer	152-layer
conv1	112×112			7×7, 64, stride 2		
conv2_x	56×56	$\begin{bmatrix} 3 \times 3, 64 \\ 3 \times 3, 64 \end{bmatrix} \times 2$	$\begin{bmatrix} 3 \times 3, 64 \\ 3 \times 3, 64 \end{bmatrix} \times 3$	$\begin{bmatrix} 1 \times 1, 64 \\ 3 \times 3, 64 \\ 1 \times 1, 256 \end{bmatrix} \times 3$	$\begin{bmatrix} 1 \times 1, 64 \\ 3 \times 3, 64 \\ 1 \times 1, 256 \end{bmatrix} \times 3$	$\begin{bmatrix} 1 \times 1, 64 \\ 3 \times 3, 64 \\ 1 \times 1, 256 \end{bmatrix} \times 3$
conv3_x	28×28	$\begin{bmatrix} 3 \times 3, 128 \\ 3 \times 3, 128 \end{bmatrix} \times 2$	$\begin{bmatrix} 3 \times 3, 128 \\ 3 \times 3, 128 \end{bmatrix} \times 4$	$\begin{bmatrix} 1 \times 1, 128 \\ 3 \times 3, 128 \\ 1 \times 1, 512 \end{bmatrix} \times 4$	$\begin{bmatrix} 1 \times 1, 128 \\ 3 \times 3, 128 \\ 1 \times 1, 512 \end{bmatrix} \times 4$	$\begin{bmatrix} 1 \times 1, 128 \\ 3 \times 3, 128 \\ 1 \times 1, 512 \end{bmatrix} \times 8$
conv4_x	14×14	$\begin{bmatrix} 3 \times 3, 256 \\ 3 \times 3, 256 \end{bmatrix} \times 2$	$\begin{bmatrix} 3 \times 3, 256 \\ 3 \times 3, 256 \end{bmatrix} \times 6$	$\begin{bmatrix} 1 \times 1, 256 \\ 3 \times 3, 256 \\ 1 \times 1, 1024 \end{bmatrix} \times 6$	$\begin{bmatrix} 1 \times 1, 256 \\ 3 \times 3, 256 \\ 1 \times 1, 1024 \end{bmatrix} \times 23$	$\begin{bmatrix} 1 \times 1, 256 \\ 3 \times 3, 256 \\ 1 \times 1, 1024 \end{bmatrix} \times 36$
conv5_x	7×7	$\begin{bmatrix} 3 \times 3, 512 \\ 3 \times 3, 512 \end{bmatrix} \times 2$	$\begin{bmatrix} 3 \times 3, 512 \\ 3 \times 3, 512 \end{bmatrix} \times 3$	$\begin{bmatrix} 1 \times 1, 512 \\ 3 \times 3, 512 \\ 1 \times 1, 2048 \end{bmatrix} \times 3$	$\begin{bmatrix} 1 \times 1, 512 \\ 3 \times 3, 512 \\ 1 \times 1, 2048 \end{bmatrix} \times 3$	$\begin{bmatrix} 1 \times 1, 512 \\ 3 \times 3, 512 \\ 1 \times 1, 2048 \end{bmatrix} \times 3$
	1×1	average pool, 1000-d fc, softmax				
FLOPs		1.8×10^9	3.6×10^9	3.8×10^9	7.6×10^9	11.3×10^9

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.

Intro

ResNet

Technical details

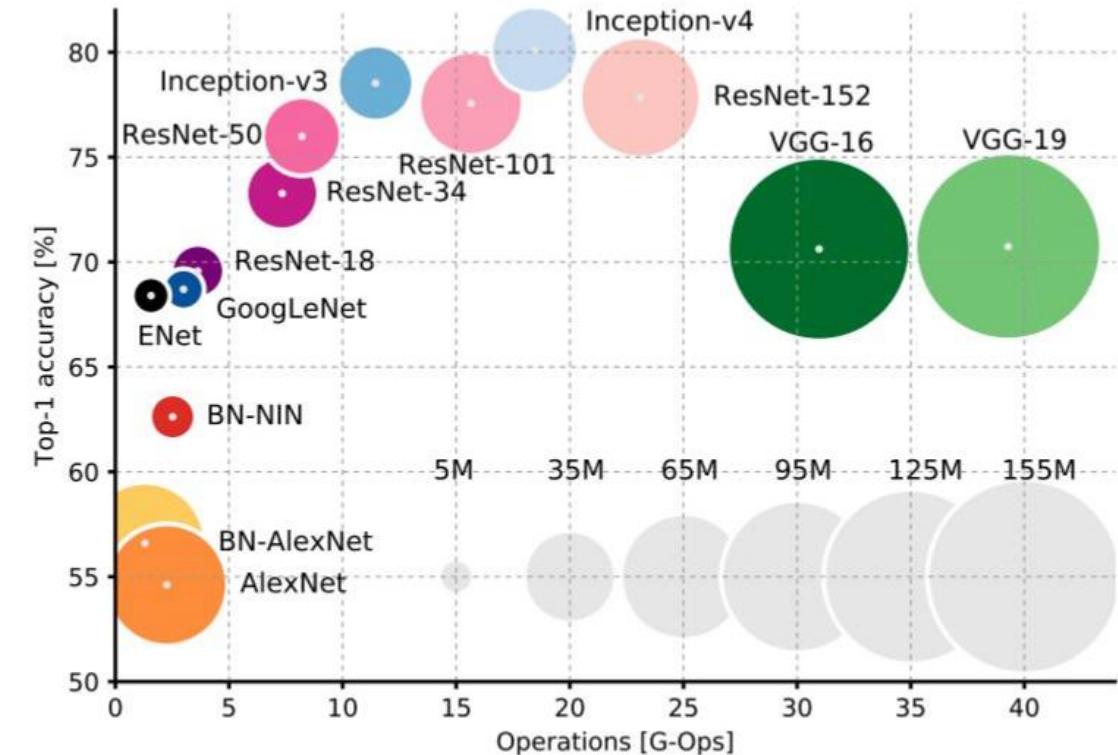
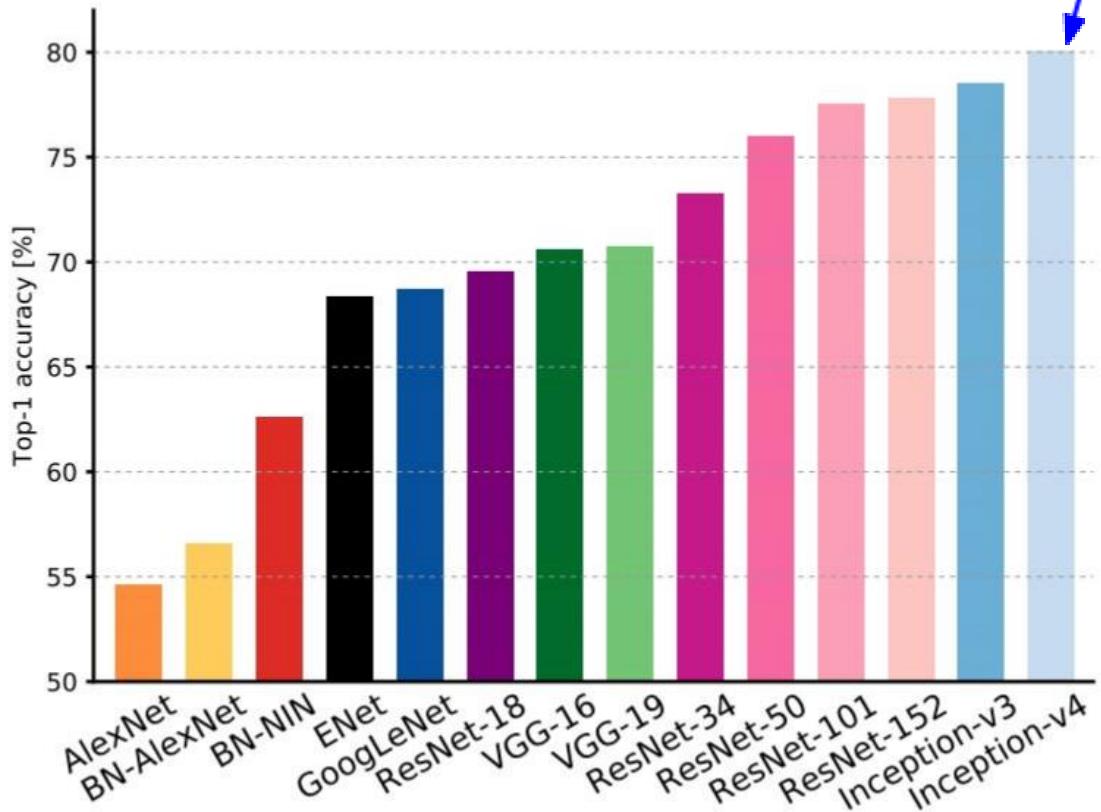
Results

ResNet 1000

Comparison

Comparing complexity...

Inception-v4: Resnet + Inception!



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

Intro

ResNet

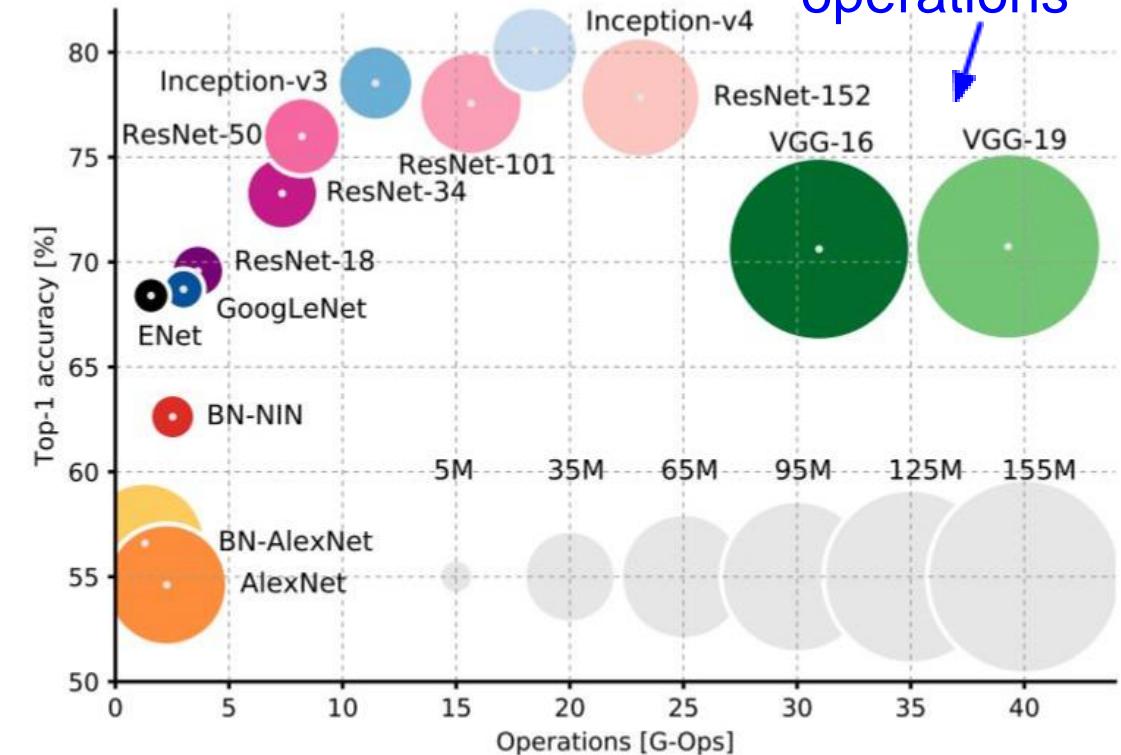
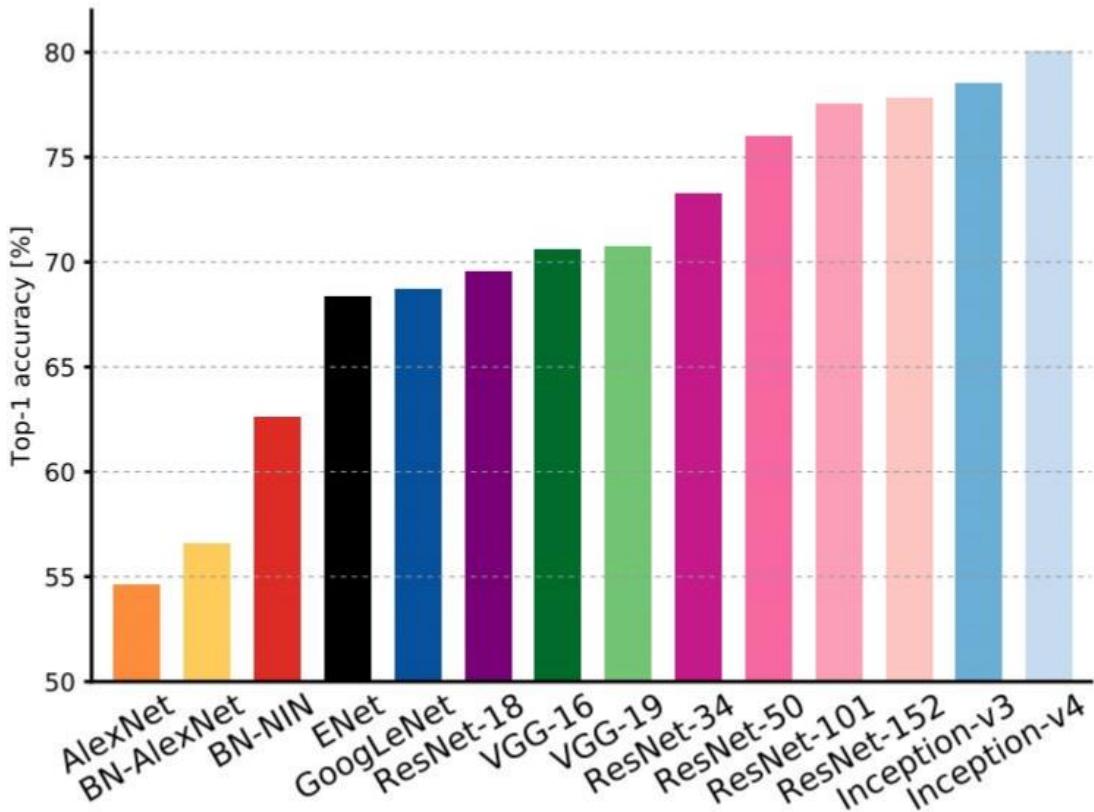
Technical details

Results

ResNet 1000

Comparison

Comparing complexity...



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

Intro

ResNet

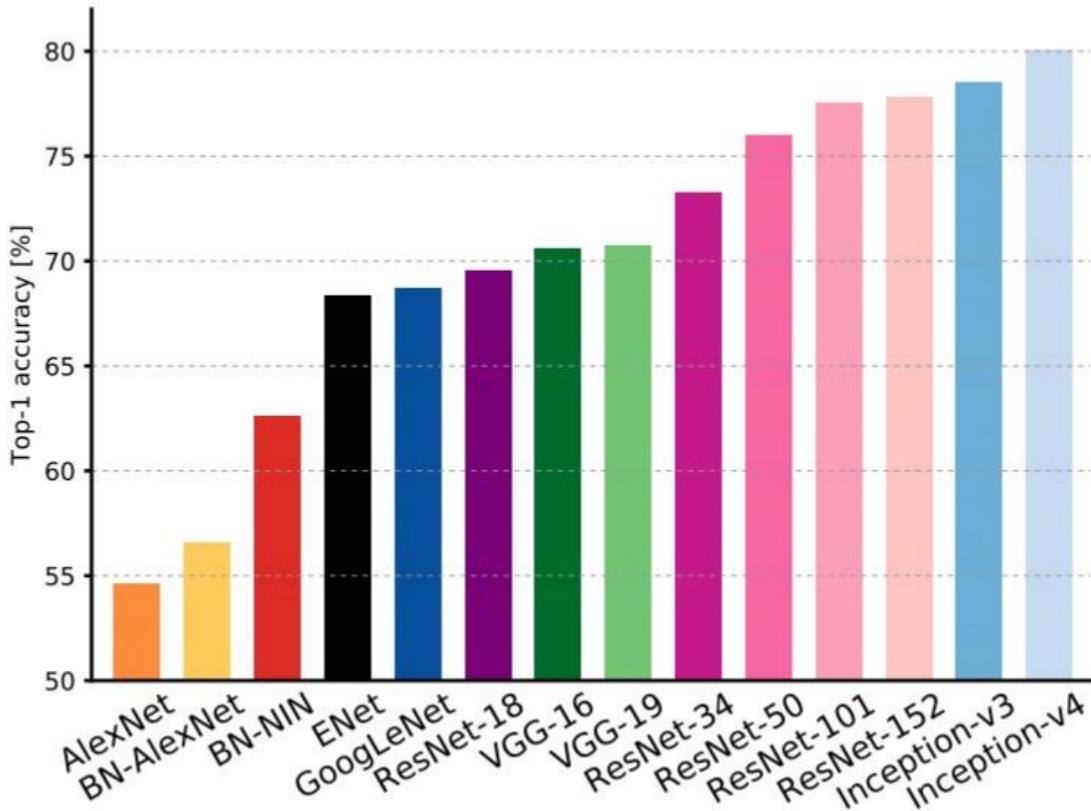
Technical details

Results

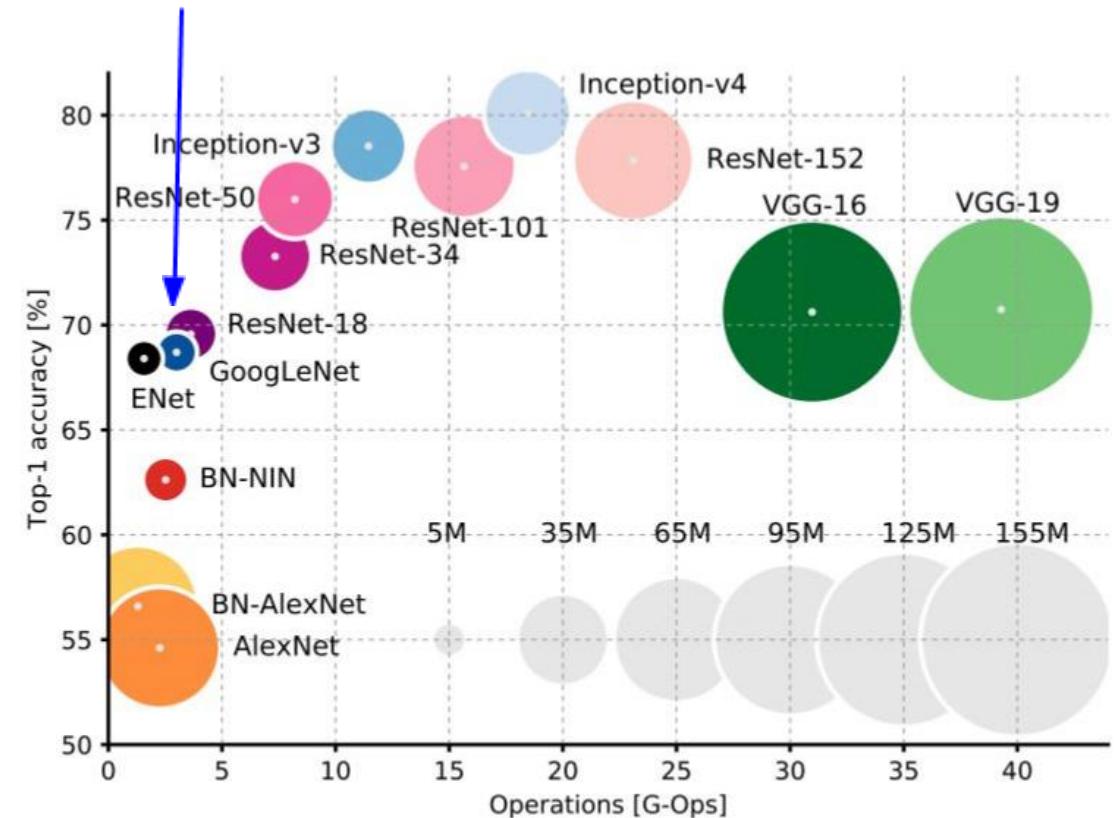
ResNet 1000

Comparison

Comparing complexity...



GoogLeNet:
most efficient



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

Intro

ResNet

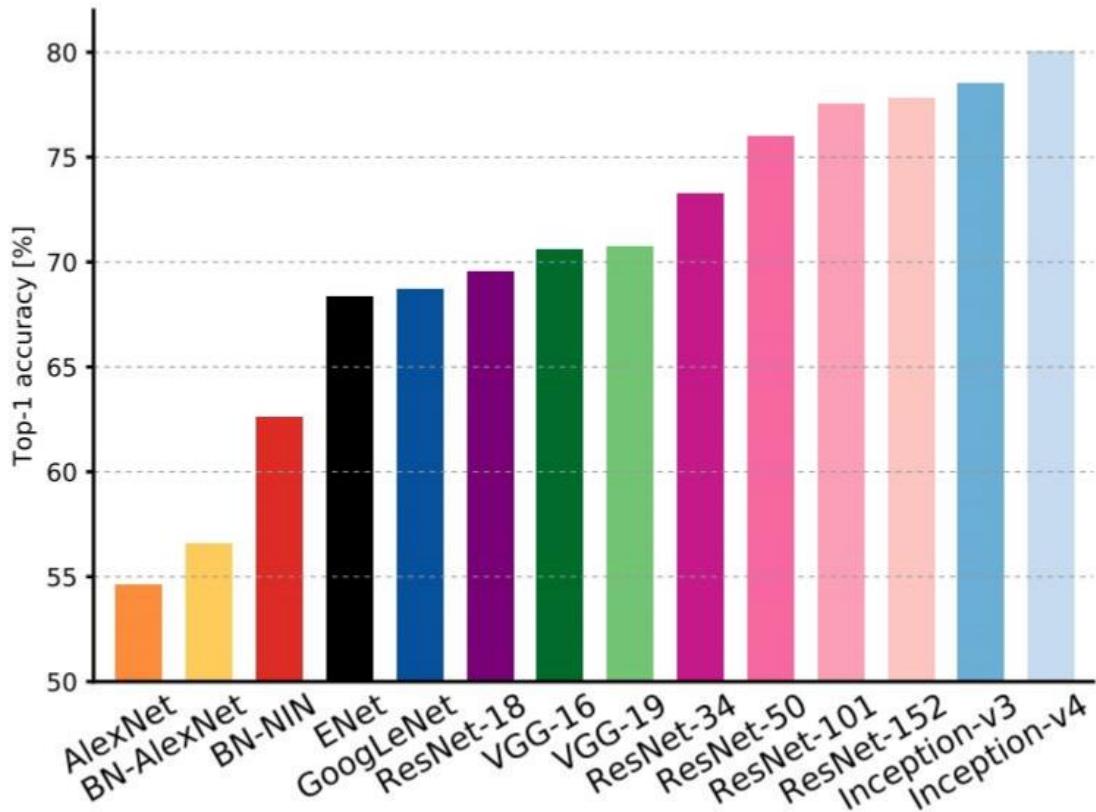
Technical
details

Results

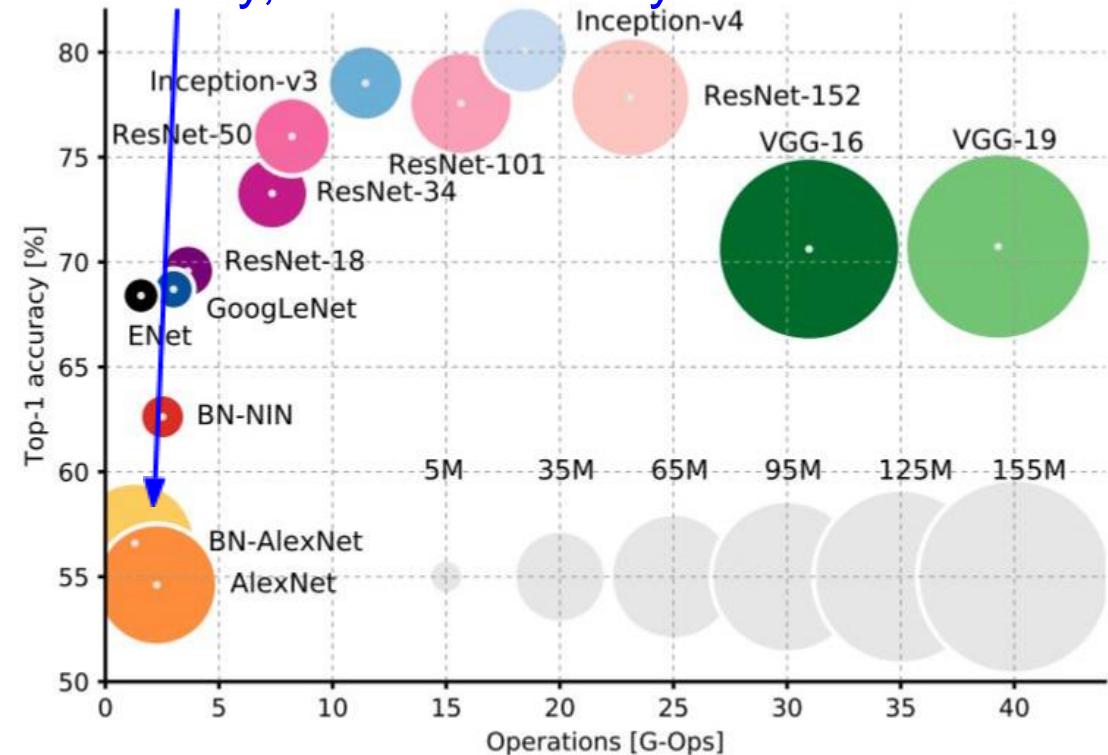
ResNet
1000

Comparison

Comparing complexity...



AlexNet:
Smaller compute, still memory heavy, lower accuracy



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

Intro

ResNet

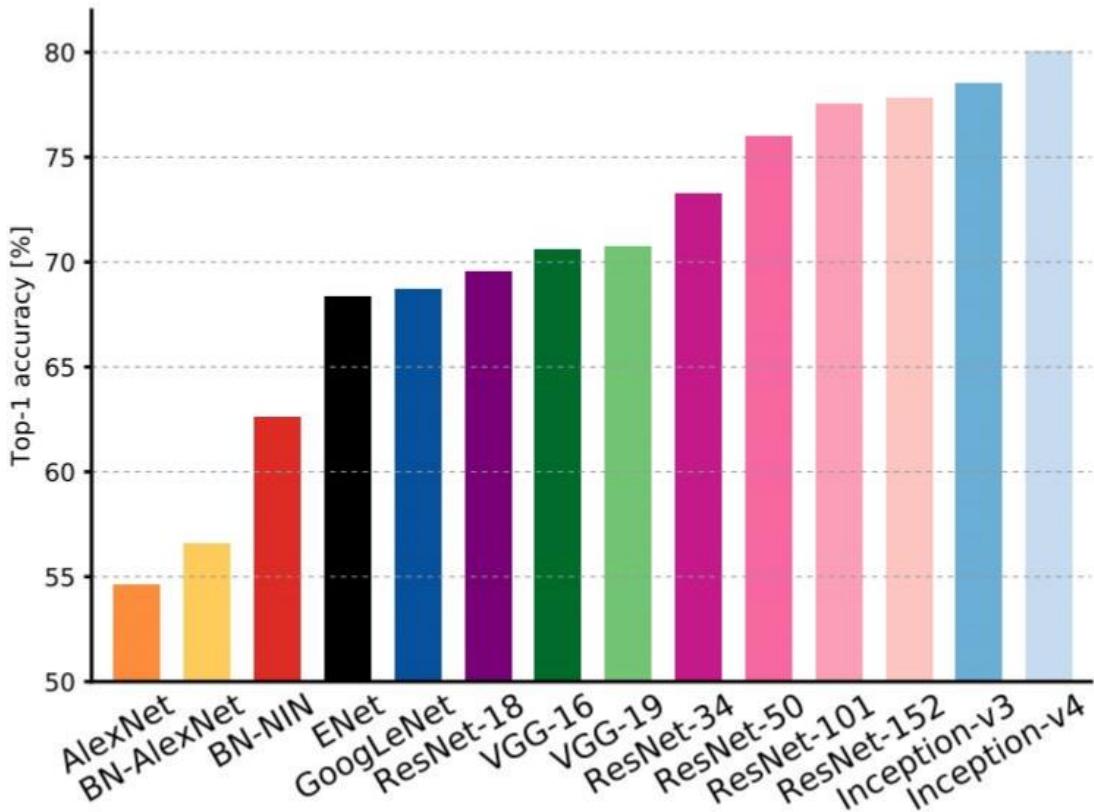
Technical details

Results

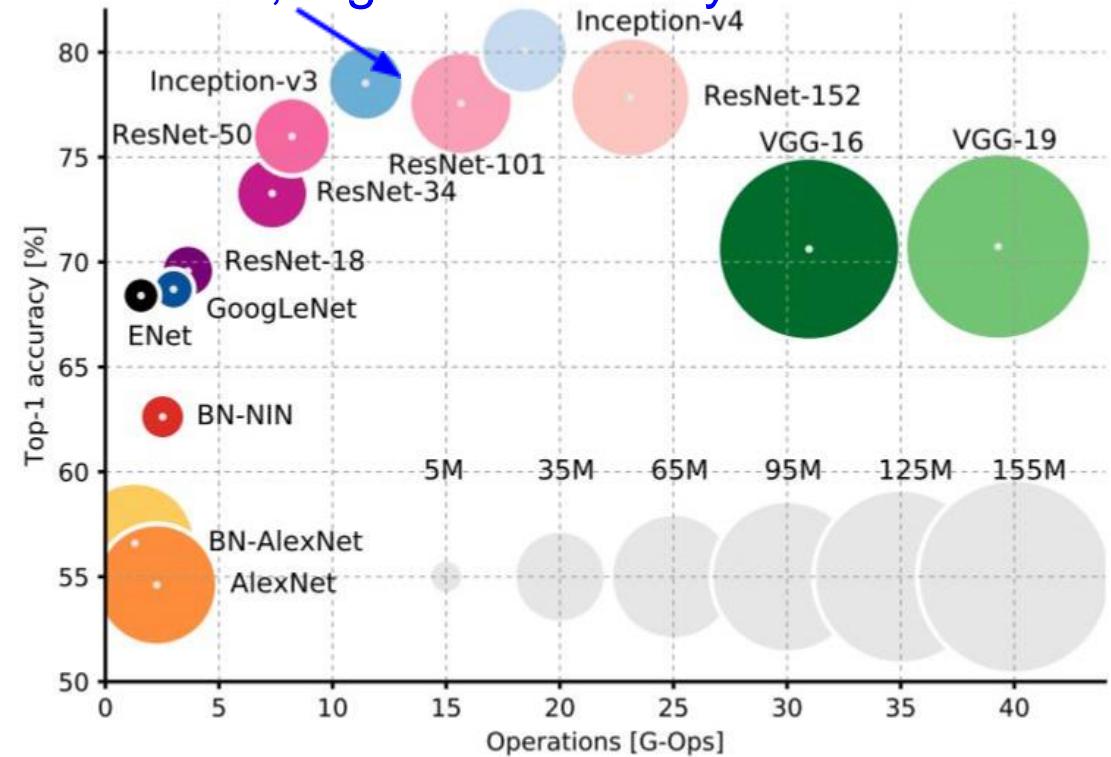
ResNet 1000

Comparison

Comparing complexity...



ResNet:
Moderate efficiency depending on
model, highest accuracy



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

Intro

ResNet

Technical
details

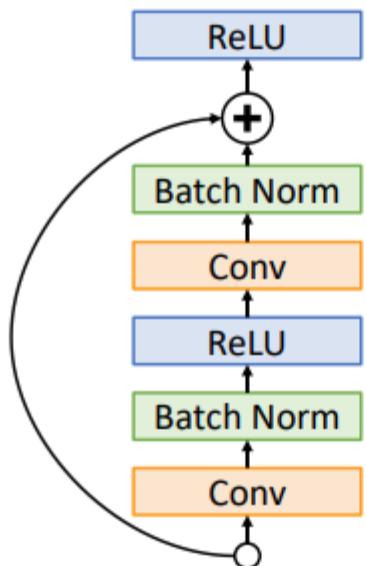
Results

ResNet
1000

Comparison

Improving Residual Networks: Block Design

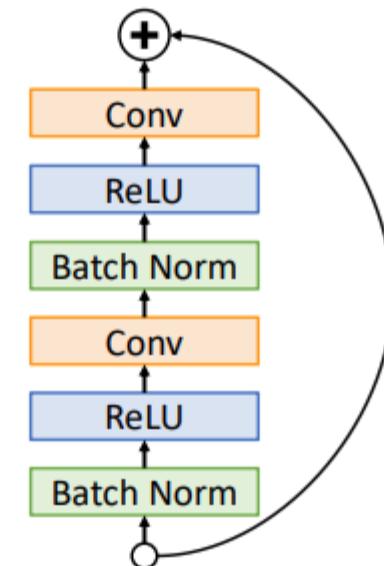
Original ResNet block



Note ReLU **after** residual:

Cannot actually learn identity function since outputs are nonnegative!

“Pre-Activation” ResNet Block

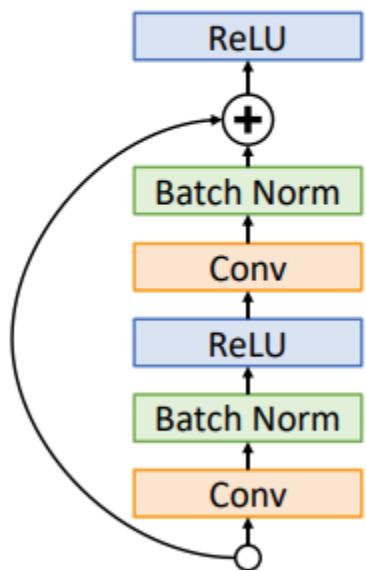


Note ReLU **inside** residual:

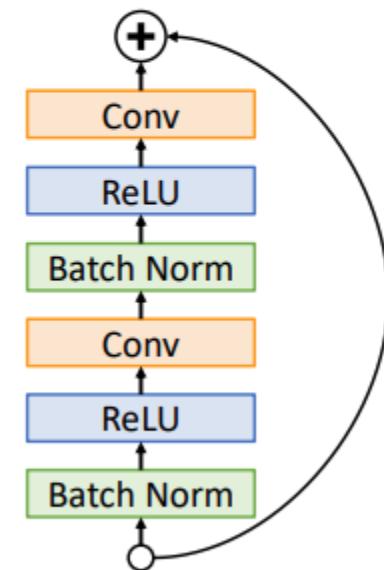
Can learn true identity function by setting Conv weights to zero!

Improving Residual Networks: Block Design

Original ResNet block



“Pre-Activation” ResNet Block



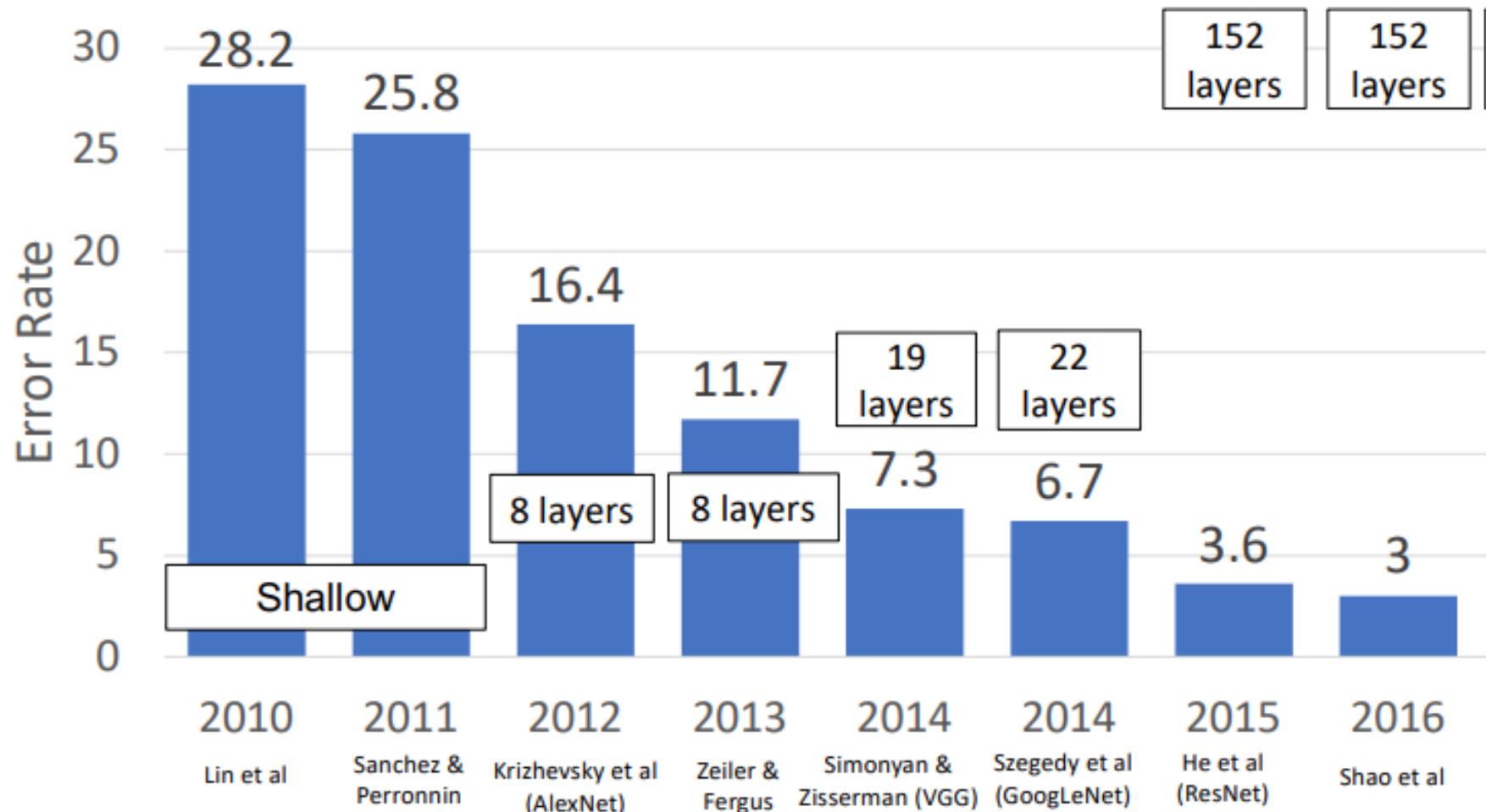
Slight improvement in accuracy
(ImageNet top-1 error)

ResNet-152: 21.3 vs **21.1**

ResNet-200: 21.8 vs **20.7**

Not actually used that much in practice

ImageNet Classification Challenge

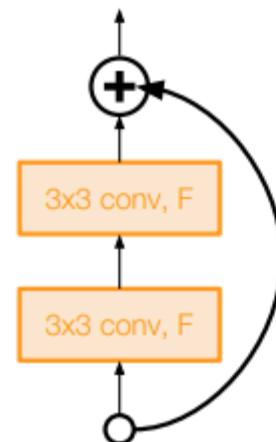


Improving ResNets...

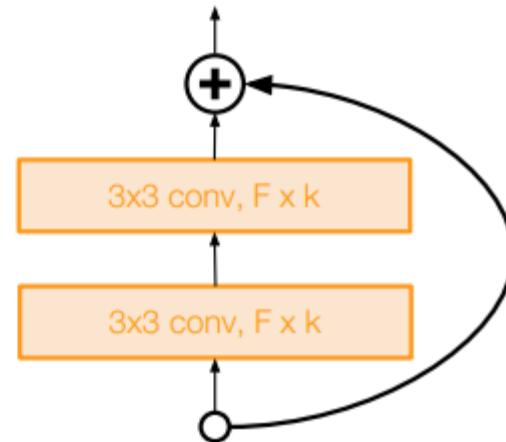
Wide Residual Networks

[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)

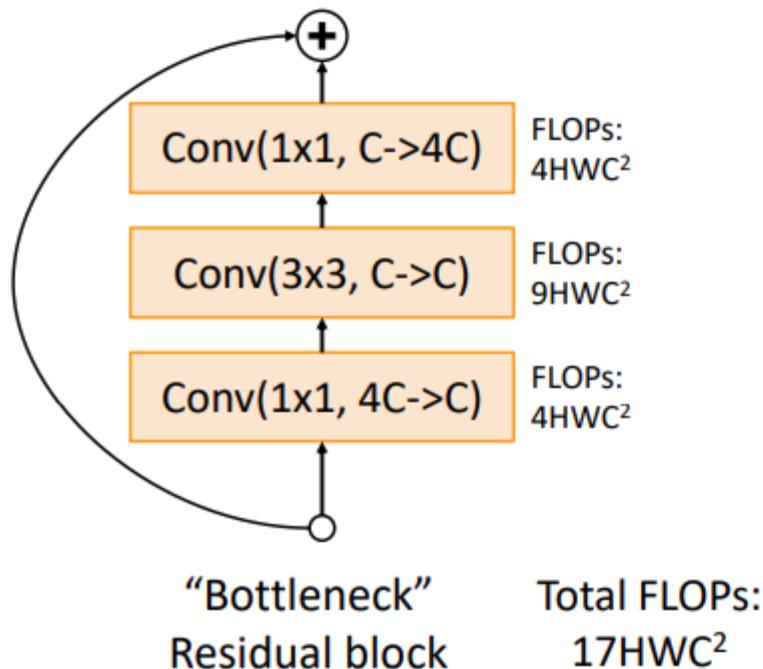


Basic residual block

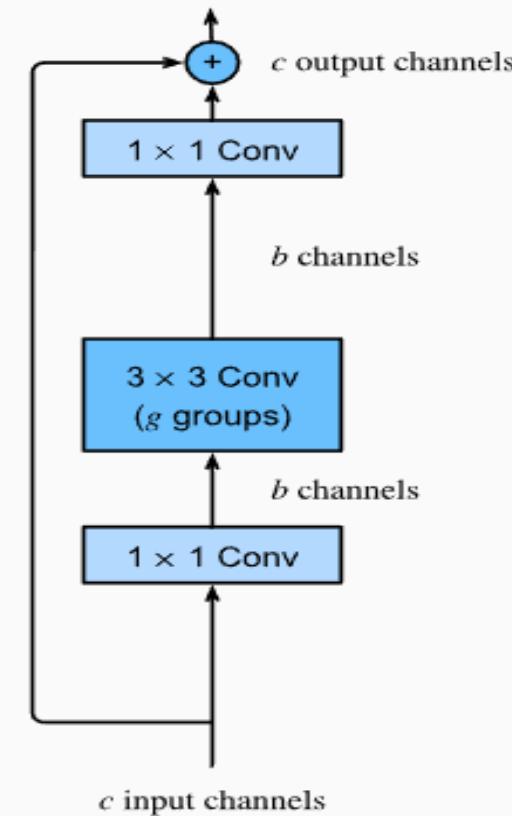
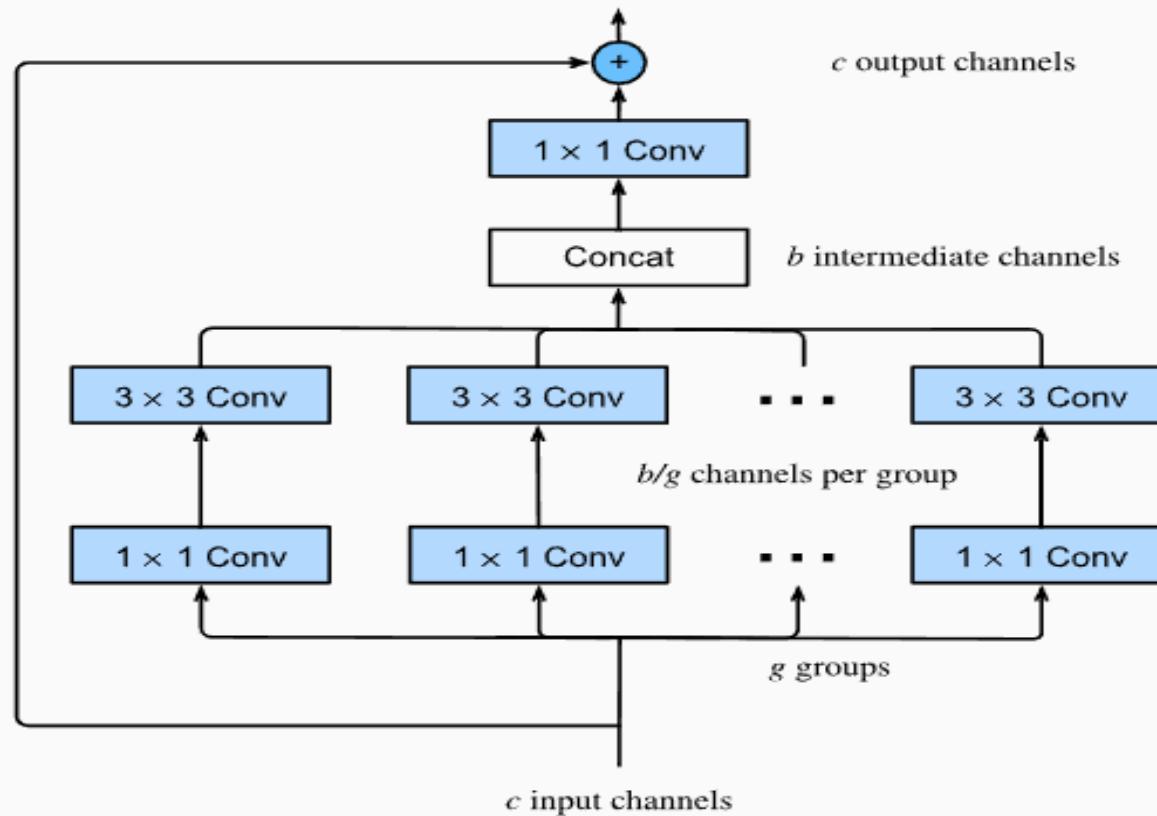


Wide residual block

Improving ResNets

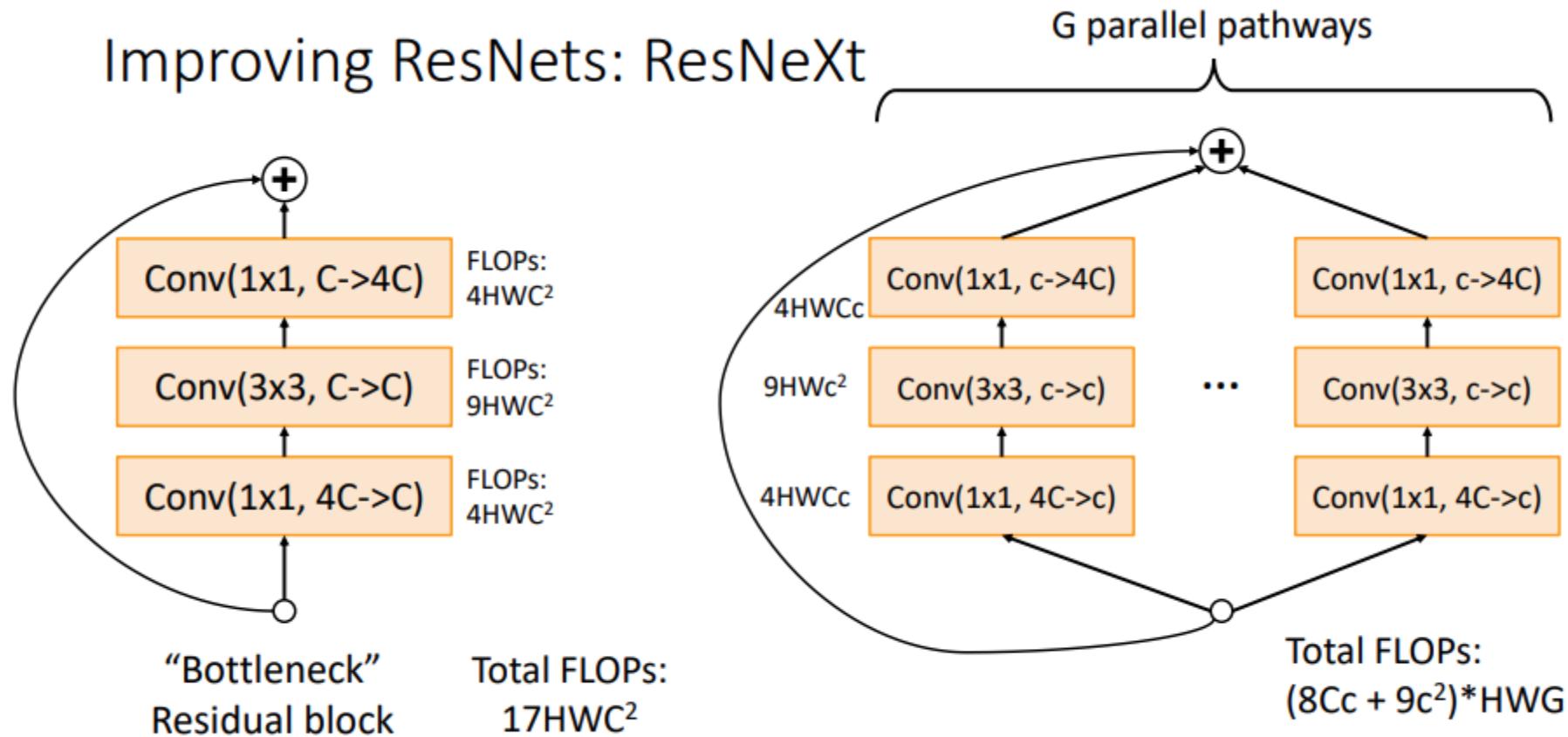


We can take some inspiration from the Inception block of [Fig. 8.4.1](#) which has information flowing through the block in separate groups. Applying the idea of multiple independent groups to the ResNet block of [Fig. 8.6.3](#) led to the design of ResNeXt ([Xie et al., 2017](#)). Different from the smorgasbord of transformations in Inception, ResNeXt adopts the same transformation in all branches, thus minimizing the need for manual tuning of each branch.



- Breaking up a convolution from c_i to c_o channels into one of g groups of size c_i/g generating g outputs of size c_o/g is called, quite fittingly, a *grouped convolution*. The computational cost (proportionally) is reduced from $O(c_i \cdot c_o)$ to $O(g \cdot (c_i/g) \cdot (c_o/g)) = O(c_i \cdot c_o/g)$, i.e., it is g times faster. Even better, the number of parameters needed to generate the output is also reduced from a $c_i \times c_o$ matrix to g smaller matrices of size $(c_i/g) \times (c_o/g)$, again a g times reduction. In what follows we assume that both c_i and c_o are divisible by g .
- The only challenge in this design is that no information is exchanged between the g groups. The ResNeXt block of Fig. amends this in two ways: the grouped convolution with a 3×3 kernel is sandwiched in between two 1×1 convolutions. The second one serves double duty in changing the number of channels back. The benefit is that we only pay the $O(c \cdot b)$ cost for 1×1 kernels and can make do with an $O(b^2/g)$ cost for 3×3 kernels. Similar to the residual block implementation in, the residual connection is replaced (thus generalized) by a 1×1 convolution.

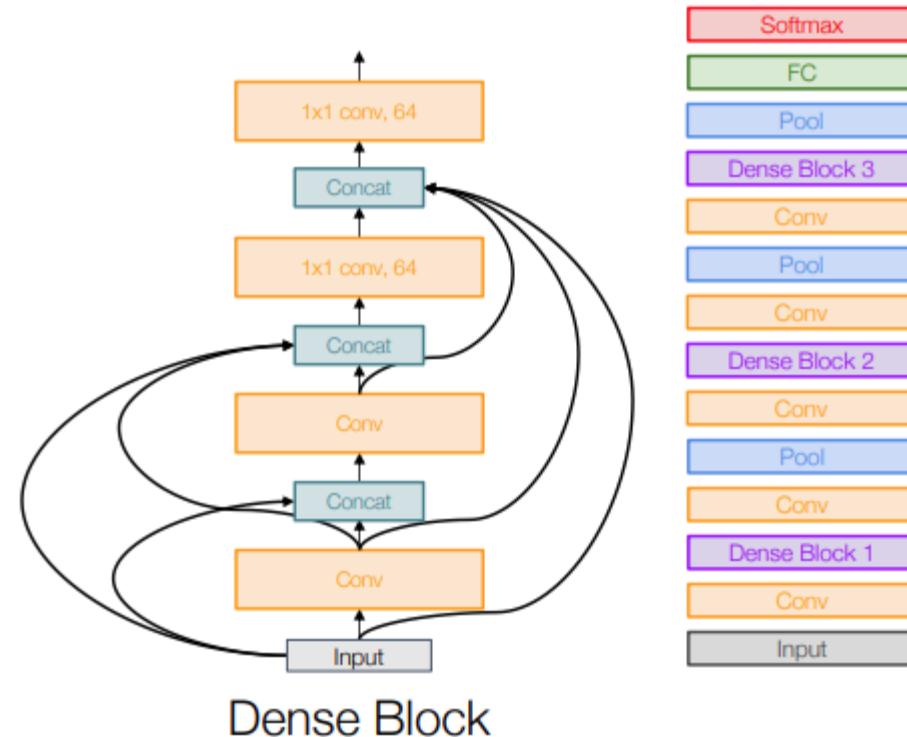
Improving ResNets: ResNeXt



Densely Connected Neural Networks

Dense blocks where each layer is connected to every other layer in feedforward fashion

Alleviates vanishing gradient, strengthens feature propagation, encourages feature reuse



“You need a lot of data if you want to
train/use CNNs”