

Lecture 8: CNN Architectures

Reminder: A2 due today!

Due at 11:59pm

Remember to run the validation script!

Soon: Assignment 3!

Modular API for backpropagation

Fully-connected networks

Dropout

Update rules: SGD+Momentum, RMSprop, Adam

Convolutional networks

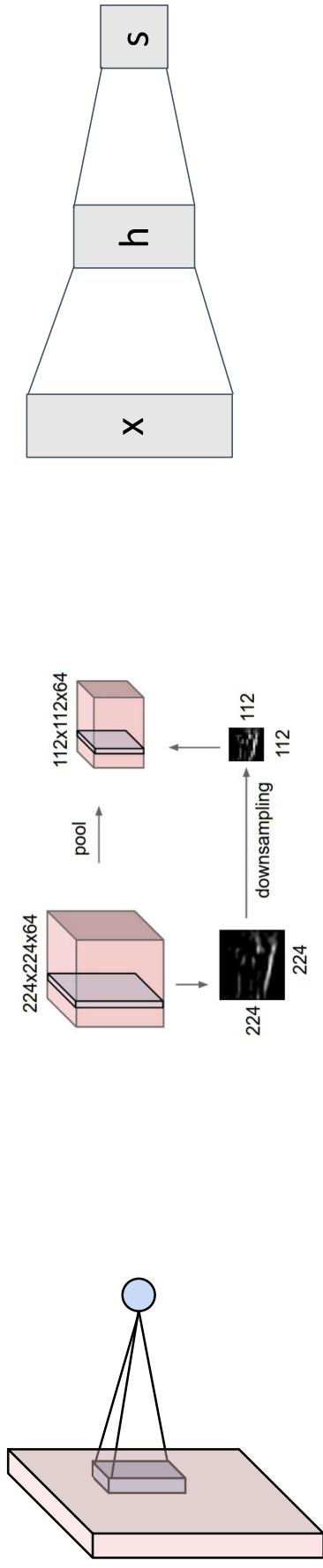
Batch normalization

Will be released today or tomorrow

Will be due two weeks from the day it is released

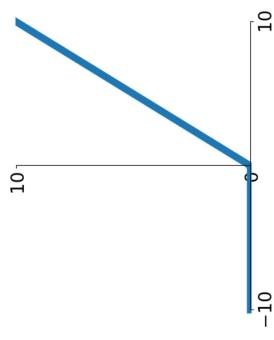
Last Time: Components of Convolutional Networks

Convolution Layers Pooling Layers Fully-Connected Layers

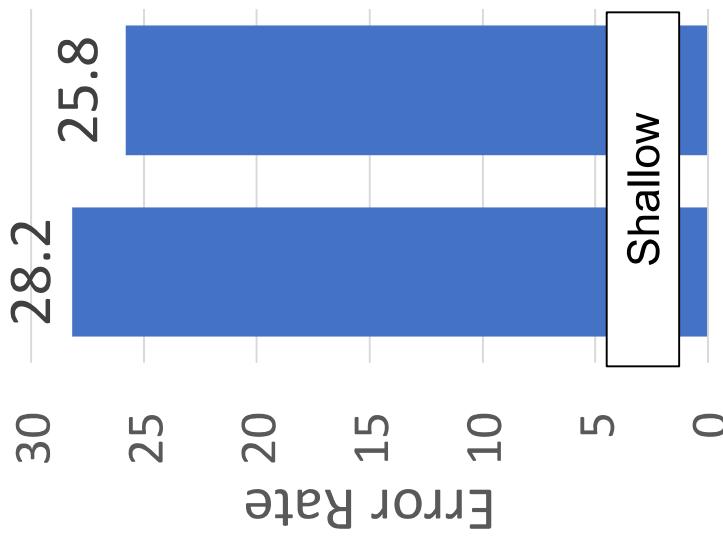


Activation Function Normalization

$$\hat{x}_{i,j} = \frac{x_{i,j} - \mu_j}{\sqrt{\sigma_j^2 + \epsilon}}$$



ImageNet Classification Challenge



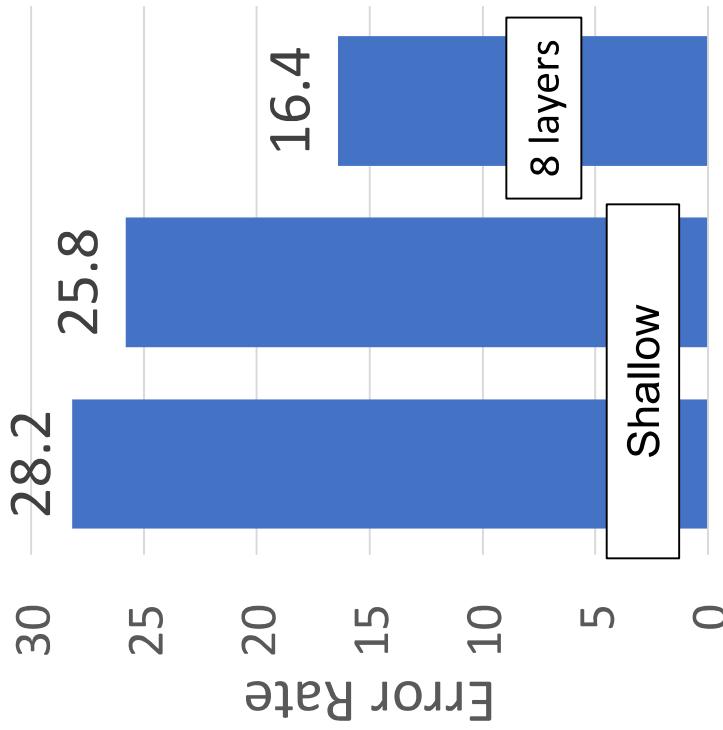
Lin et al
Sanchez &
Perronnin

Justin Johnson

Lecture 8 - 5

September 30, 2019

ImageNet Classification Challenge

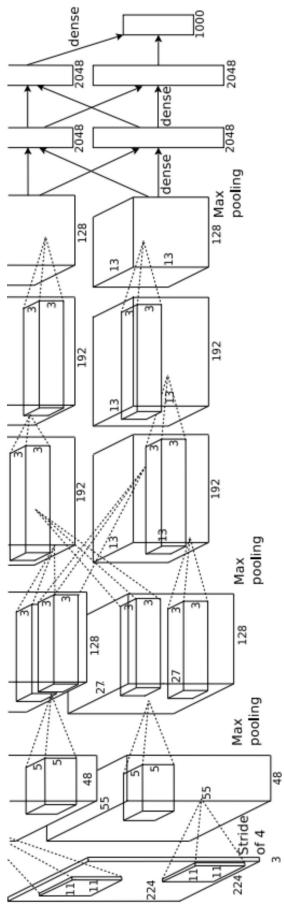


Lin et al Sanchez & Perronnin Krizhevsky et al
(AlexNet) (AlexNet)

Justin Johnson

Lecture 8 - 6

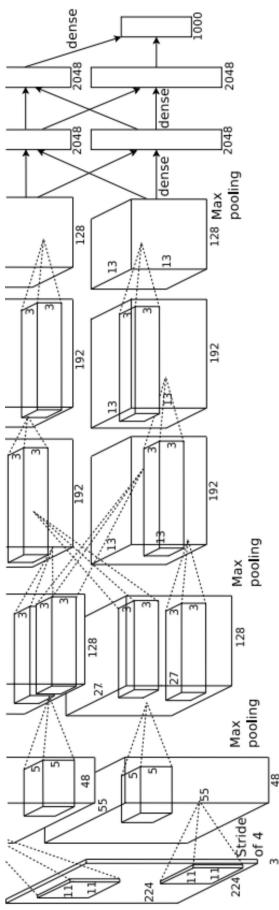
September 30, 2019



AlexNet

- 227 × 227 inputs
- 5 Convolutional layers
- Max pooling
- 3 fully-connected layers
- ReLU nonlinearities

AlexNet



Used “Local response normalization”;
Not used anymore

227 x 227 inputs
5 Convolutional layers
Max pooling
3 fully-connected layers
ReLU nonlinearities

Trained on two GTX 580 GPUs – only
3GB of memory each! Model split
over two GPUs

AlexNet

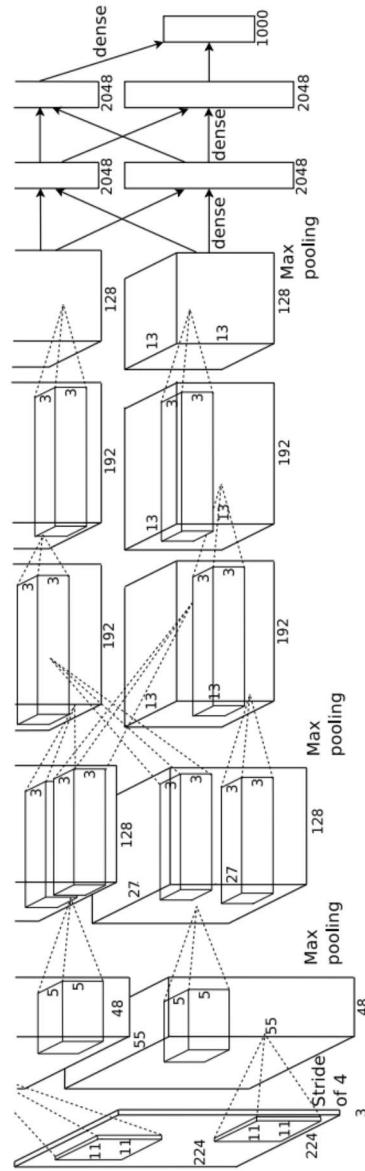
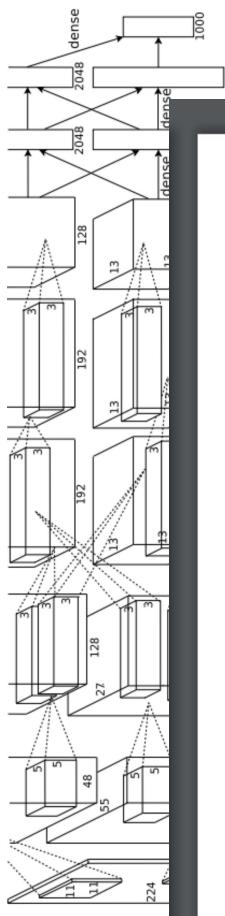
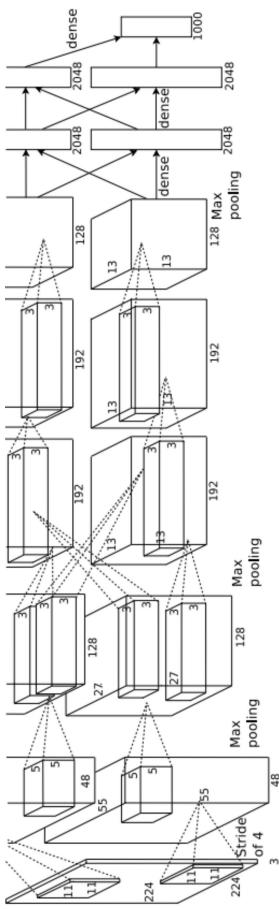


Figure 2: An illustration of the architecture of our CNN, explicitly showing the delineation of responsibilities between the two GPUs. One GPU runs the layer-parts at the top of the figure while the other runs the layer-parts at the bottom. The GPUs communicate only at certain layers. The network's input is 150,528-dimensional, and the number of neurons in the network's remaining layers is given by 253,440–186,624–64,896–64,896–43,264–4096–4096–1000.

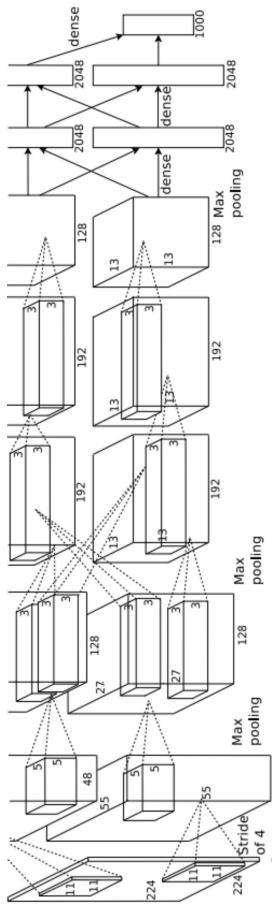
AlexNet



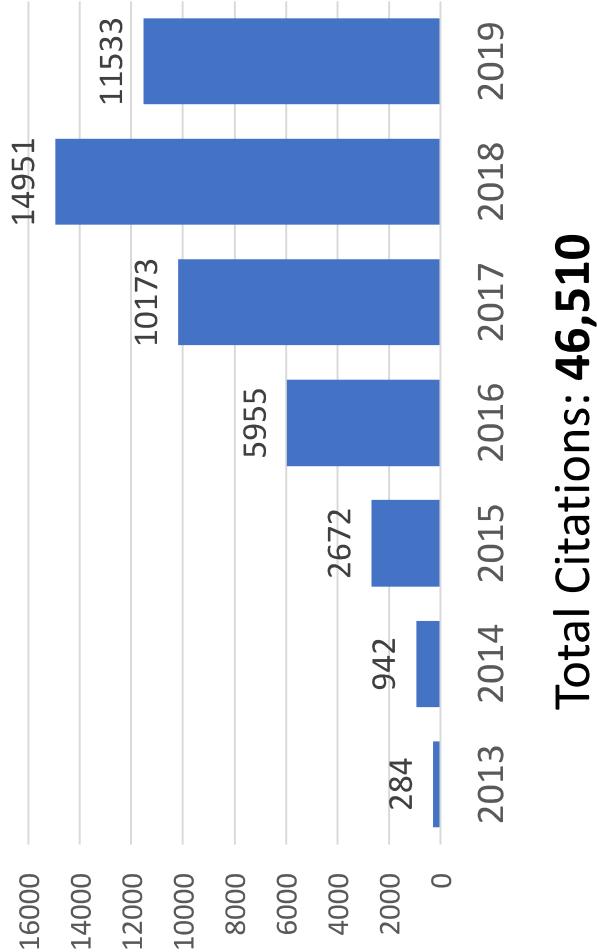
AlexNet Citations per year
(As of 9/30/2019)



AlexNet



AlexNet Citations per year
(As of 9/30/2019)



Citation Counts

Darwin, "On the origin of species", 1859: **50,007**

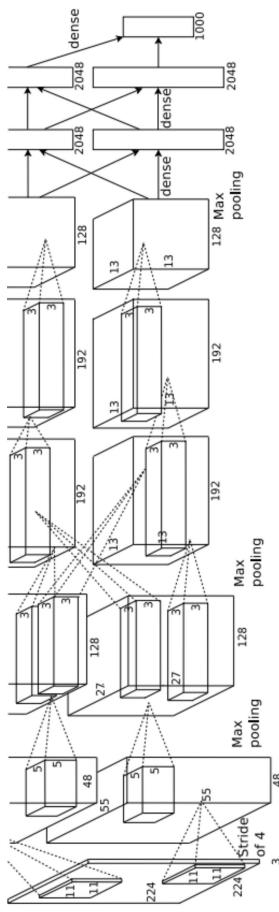
Shannon, "A mathematical theory of communication", 1948: **69,351**

Watson and Crick, "Molecular Structure of Nucleic Acids", 1953: **13,111**

ATLAS Collaboration, "Observation of a new particle in the search for the Standard Model Higgs boson with the ATLAS detector at the LHC", 2012: **14,424**

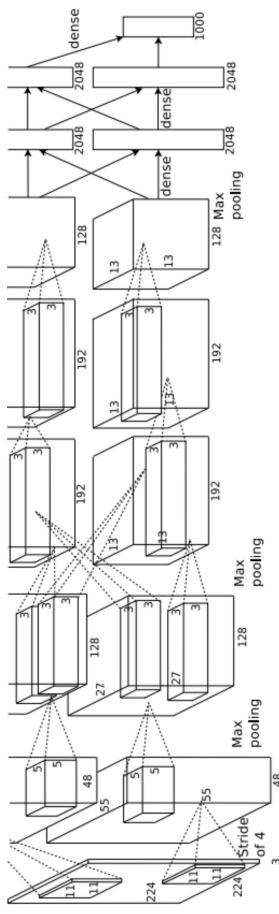
Figure copyright Alex Krizhevsky, Ilya Sutskever, and Geoffrey Hinton, 2012. Reproduced with permission.

AlexNet



Layer	Input size	Layer	Output size					
Layer	C	H / W	filters	kernel	stride	pad	C	H / W
conv1	3	227	64	11	4	2	?	

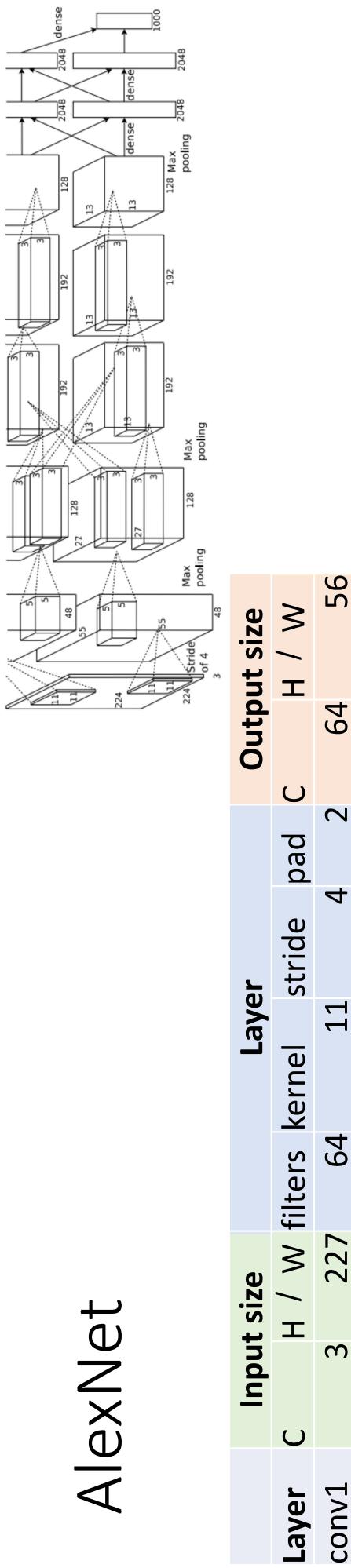
AlexNet



Layer	Input size	Layer	Output size					
Layer	C	H / W	filters	kernel	stride	pad	C	H / W
conv1	3	227	64	11	4	2	64	?

Recall: Output channels = number of filters

AlexNet



$$\begin{aligned}
 \text{Recall: } W' &= (W - K + 2P) / S + 1 \\
 &= 227 - 11 + 2 * 2) / 4 + 1 \\
 &= 220/4 + 1 = 56
 \end{aligned}$$

AlexNet

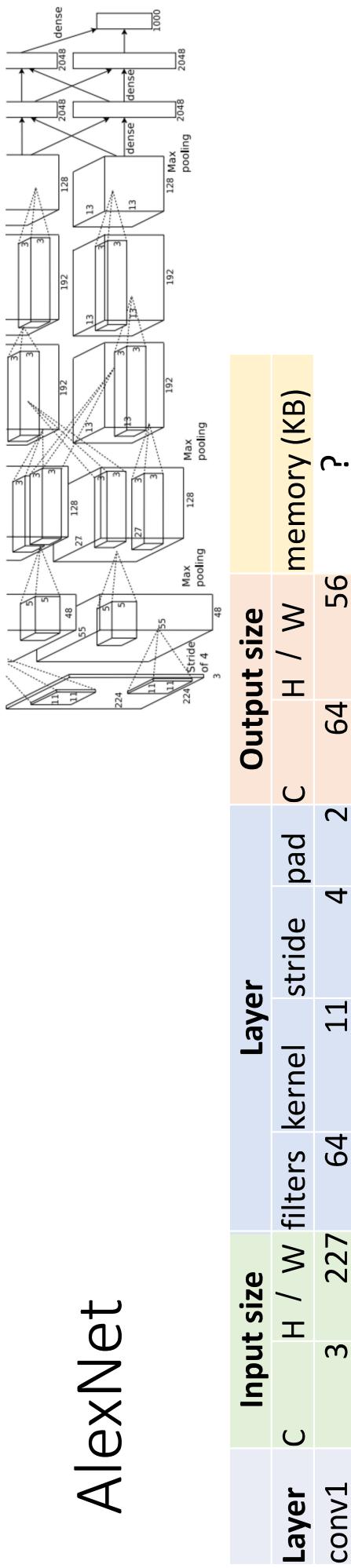
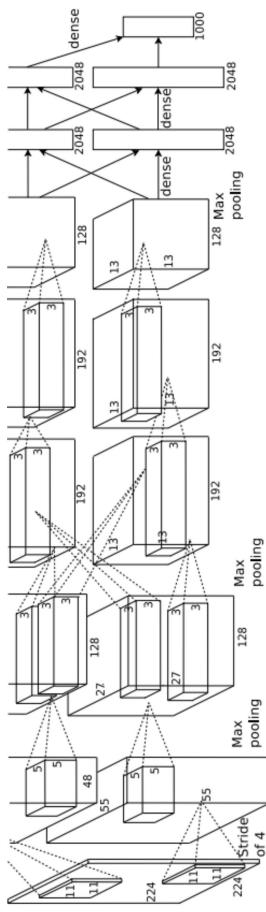


Figure copyright Alex Krizhevsky, Ilya Sutskever, and Geoffrey Hinton, 2012. Reproduced with permission.

AlexNet



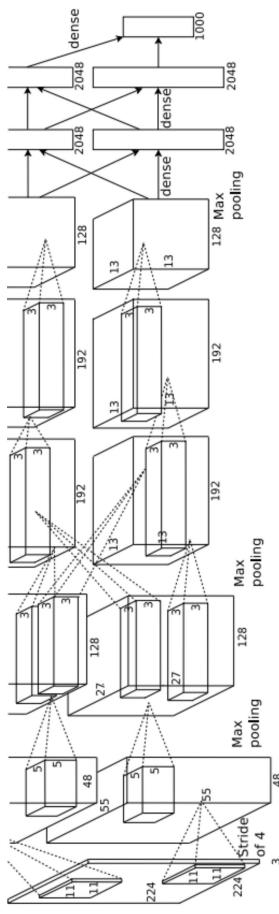
Number of output elements = $C * H' * W'$

$$= 64 * 56 * 56 = 200,704$$

Bytes per element = 4 (for 32-bit floating point)

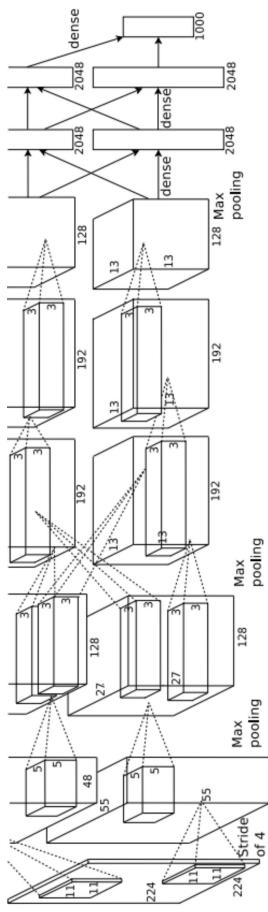
$$\begin{aligned} \text{KB} &= (\text{number of elements}) * (\text{bytes per elem}) / 1024 \\ &= 200704 * 4 / 1024 \\ &= \mathbf{784} \end{aligned}$$

AlexNet



	Input size			Layer			Output size			
Layer	C	H / W	filters	kernel	stride	pad	C	H / W	memory (kB)	params (k)
conv1	3	227	64	11	4	2	64	56	784	?

AlexNet

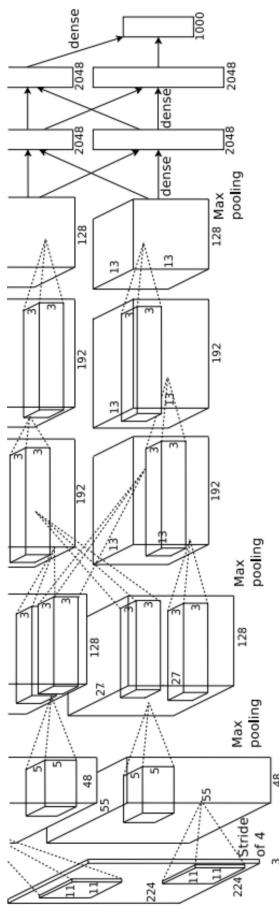


Layer	Input size	Layer	Output size				memory (kB)	params (k)
	C H / W		filters	kernel	stride	pad	C	H / W
conv1	3 227	conv1	64	11	4	2	64	56

$$\begin{aligned}
 \text{Weight shape} &= C_{\text{out}} \times C_{\text{in}} \times K \times K \\
 &= 64 \times 3 \times 11 \times 11 \\
 \text{Bias shape} &= C_{\text{out}} = 64 \\
 \text{Number of weights} &= 64 * 3 * 11 * 11 + 64 \\
 &= \mathbf{23,296}
 \end{aligned}$$

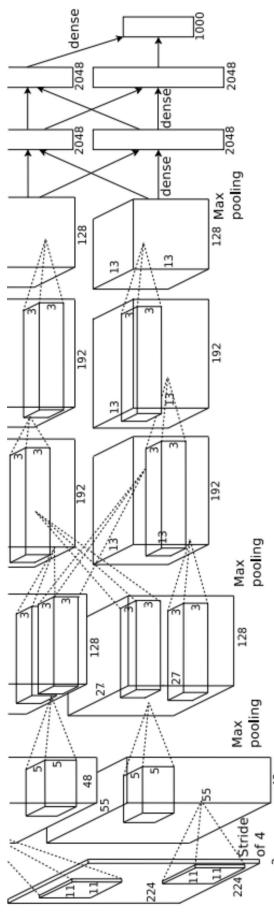
Figure copyright Alex Krizhevsky, Ilya Sutskever, and Geoffrey Hinton, 2012. Reproduced with permission.

AlexNet



	Input size			Layer			Output size					
Layer	C	H / W	filters	kernel	stride	pad	C	H / W	memory (kB)	params (k)	flop (M)	
conv1	3	227	64	11	4	2	64	56	784	23	?	

AlexNet



	Input size		Layer		Output size			memory (kB)	params (k)	flop (M)
Layer	C	H / W	filters	kernel	stride	pad	C			
conv1	3	227	64		11	4	2	64	56	784

Number of floating point operations (multiply+add)

= (number of output elements) * (ops per output elem)

$$= (C_{\text{out}} \times H' \times W') * (C_{\text{in}} \times K \times K)$$

$$= (64 * 56 * 56) * (3 * 11 * 11)$$

$$= 200,704 * 363$$

$$= 72,855,552$$

AlexNet

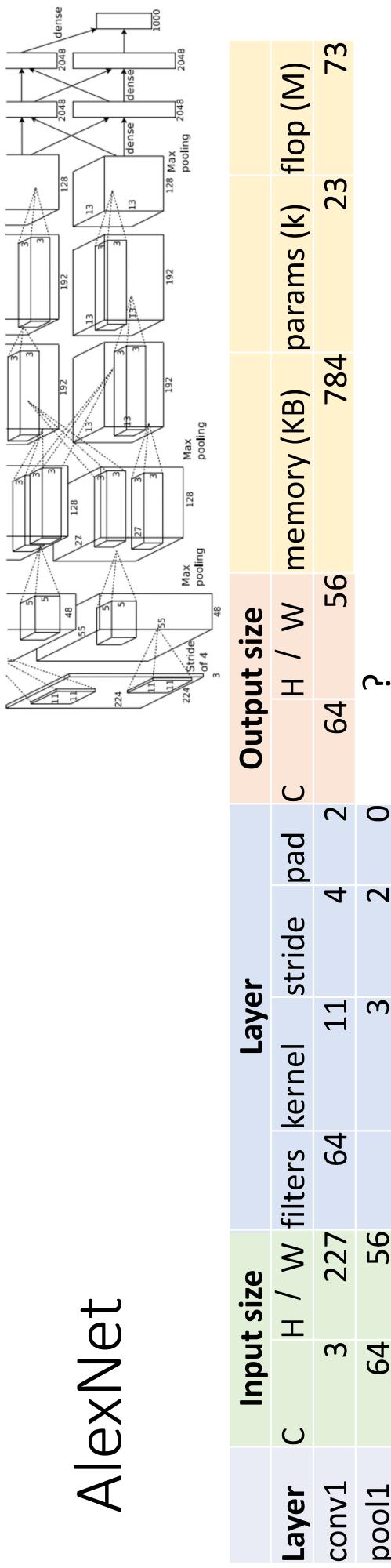
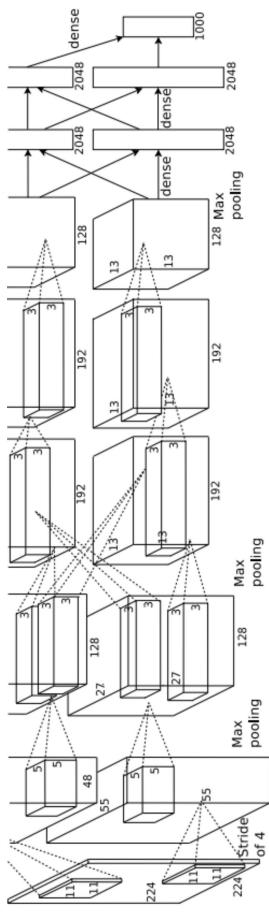


Figure copyright Alex Krizhevsky, Ilya Sutskever, and Geoffrey Hinton, 2012. Reproduced with permission.

AlexNet



	Input size			Layer			Output size				
Layer	C	H / W	filters	kernel	stride	pad	C	H / W	memory (KB)	params (k)	flop (M)
conv1	3	227	64		11	4	2	64	56	784	23
pool1	64	56			3	2	0	64	27		73

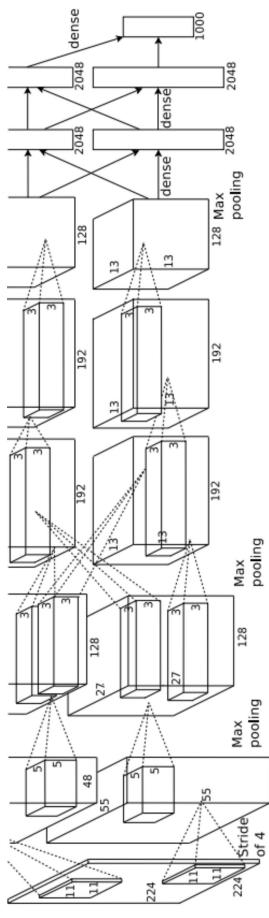
For pooling layer:

$$\# \text{output channels} = \# \text{input channels} = 64$$

$$\begin{aligned} W' &= \text{floor}((W - K) / S + 1) \\ &= \text{floor}(53 / 2 + 1) = \text{floor}(27.5) = 27 \end{aligned}$$

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

AlexNet



Layer	C	H / W	filters	kernel	stride	pad	C	H / W	memory (KB)	params (k)	flop (M)
conv1	3	227	64		11	4	2	64	56	784	23
pool1	64	56			3	2	0	64	27	182	?

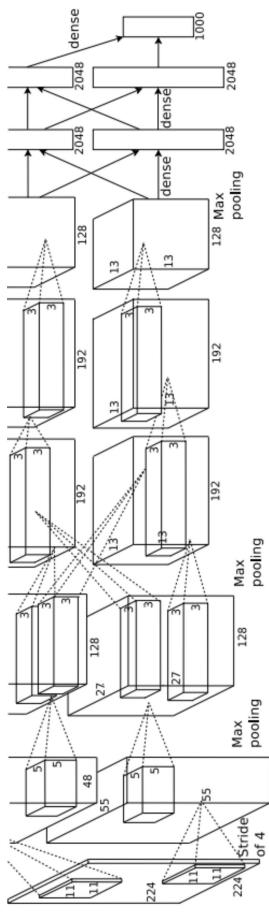
$$\# \text{output elems} = C_{\text{out}} \times H' \times W'$$

$$\text{Bytes per elem} = 4$$

$$\begin{aligned} \text{KB} &= C_{\text{out}} * H' * W' * 4 / 1024 \\ &= 64 * 27 * 4 / 1024 \\ &= \mathbf{182.25} \end{aligned}$$

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

AlexNet

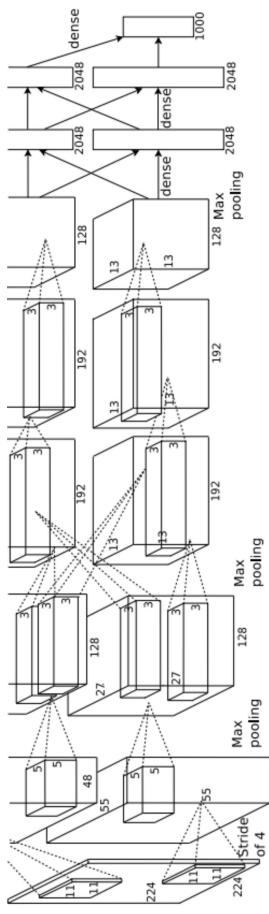


Input size

Layer	C	H / W	filters	kernel	stride	pad	C	H / W	memory (KB)	params (k)	flop (M)
conv1	3	227	64		11	4	2	64	56	784	23
pool1	64	56			3	2	0	64	27	182	0 ?

Pooling layers have no learnable parameters!

AlexNet



	Input size			Layer			Output size					
Layer	C	H / W	filters	kernel	stride	pad	C	H / W	memory (KB)	params (k)	flop (M)	
conv1	3	227	64	11	4	2	64	56	784	23	73	
pool1	64	56		3	2	0	64	27	182	0	0	

Floating-point ops for pooling layer

$$= (\text{number of output positions}) * (\text{flops per output position})$$

$$= (C_{\text{out}} * H' * W') * (K * K)$$

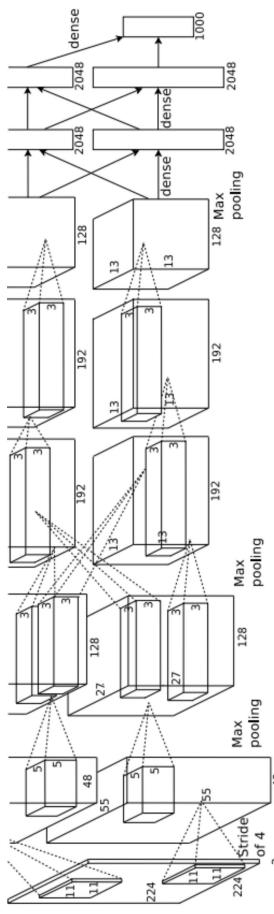
$$= (64 * 27 * 27) * (3 * 3)$$

$$= 419,904$$

$$= \mathbf{0.4 \text{ MFLOP}}$$

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

AlexNet



Layer	Input size			Layer			Output size			memory (KB)	params (k)	flop (M)
	C	H / W	filters	kernel	stride	pad	C	H / W				
conv1	3	227	64	11	4	2	64	56	784	23	73	
pool1	64	56		3	2	0	64	27	182	0	0	
conv2	64	27	192	5	1	2	192	27	547	307	224	
pool2	192	27		3	2	0	192	13	127	0	0	
conv3	192	13	384	3	1	1	384	13	254	664	112	
conv4	384	13	256	3	1	1	256	13	169	885	145	
conv5	256	13	256	3	1	1	256	13	169	590	100	
pool5	256	13		3	2	0	256	6	36	0	0	
flatten	256	6						9216	36	0	0	

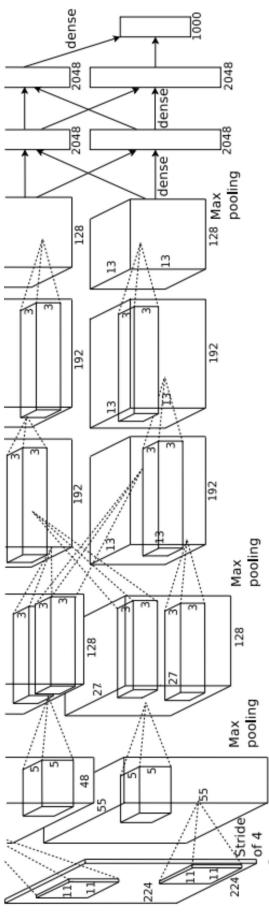
$$\text{Flatten output size} = C_{\text{in}} \times H \times W$$

$$= 256 * 6 * 6$$

$$= \mathbf{9216}$$

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

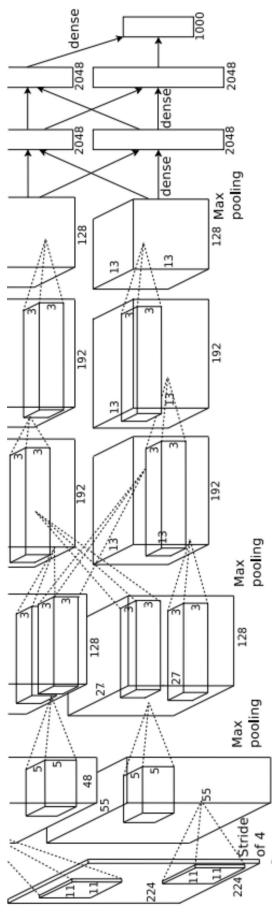
AlexNet



	Input size			Layer			Output size				
Layer	C	H / W	filters	kernel	stride	pad	C	H / W	memory (KB)	params (k)	flop (M)
conv1	3	227	64	11	4	2	64	56	784	23	73
pool1	64	56		3	2	0	64	27	182	0	0
conv2	64	27	192	5	1	2	192	27	547	307	224
pool2	192	27		3	2	0	192	13	127	0	0
conv3	192	13	384	3	1	1	384	13	254	664	112
conv4	384	13	256	3	1	1	256	13	169	885	145
conv5	256	13	256	3	1	1	256	13	169	590	100
pool5	256	13		3	2	0	256	6	36	0	0
flatten	256	6					9216		36	0	0
fc6	9216		4096							37,749	38

$$\begin{aligned}
 \text{FC params} &= C_{\text{in}} * C_{\text{out}} + C_{\text{out}} \\
 &= 9216 * 4096 + 4096 \\
 &= 37,725,832
 \end{aligned}
 \quad
 \begin{aligned}
 \text{FC flops} &= C_{\text{in}} * C_{\text{out}} \\
 &= 9216 * 4096 \\
 &= 37,748,736
 \end{aligned}$$

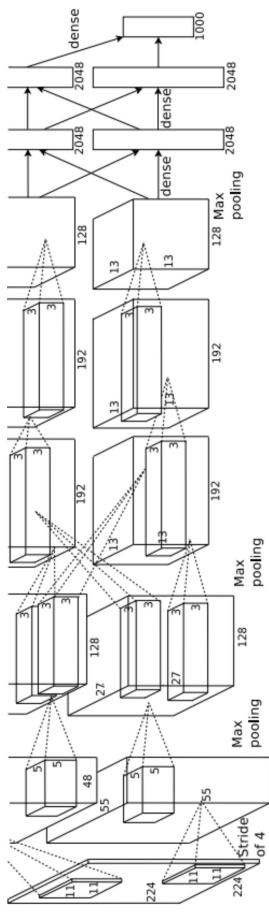
AlexNet



Layer	Input size			Layer			Output size			memory (KB)	params (k)	flop (M)
	C	H / W	filters	kernel	stride	pad	C	H / W				
conv1	3	227	64	11	4	2	64	56	784	23	73	
pool1	64	56		3	2	0	64	27	182	0	0	
conv2	64	27	192	5	1	2	192	27	547	307	224	
pool2	192	27		3	2	0	192	13	127	0	0	
conv3	192	13	384	3	1	1	384	13	254	664	112	
conv4	384	13	256	3	1	1	256	13	169	885	145	
conv5	256	13	256	3	1	1	256	13	169	590	100	
pool5	256	13		3	2	0	256	6	36	0	0	
flatten	256	6					9216		36	0	0	
fc6	9216		4096				4096		16	37,749	38	
fc7	4096		4096				4096		16	16,777	17	
fc8	4096		1000				1000		4	4,096	4	

AlexNet

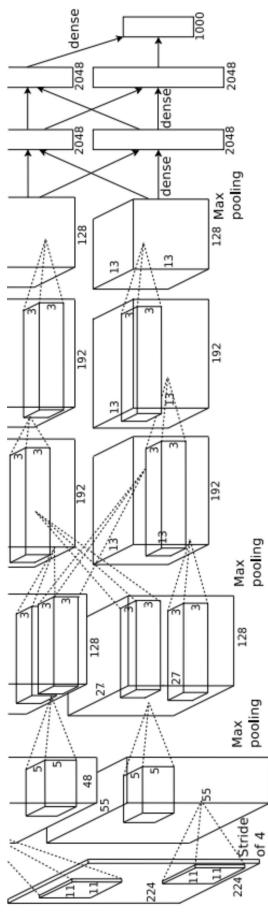
How to choose this?
Trial and error = (



Layer	Input size	Layer					Output size		memory (KB)	params (k)	flop (M)
	C	H / W	filters	kernel	stride	pad	C	H / W			
conv1	3	227	64	11	4	2	64	56	784	23	73
pool1	64	56		3	2	0	64	27	182	0	0
conv2	64	27	192	5	1	2	192	27	547	307	224
pool2	192	27		3	2	0	192	13	127	0	0
conv3	192	13	384	3	1	1	384	13	254	664	112
conv4	384	13	256	3	1	1	256	13	169	885	145
conv5	256	13	256	3	1	1	256	13	169	590	100
pool5	256	13		3	2	0	256	6	36	0	0
flatten	256	6					9216		36	0	0
fc6	9216		4096				4096		16	37,749	38
fc7	4096		4096				4096		16	16,777	17
fc8	4096		1000				1000		4	4,096	4

AlexNet

Interesting trends here!

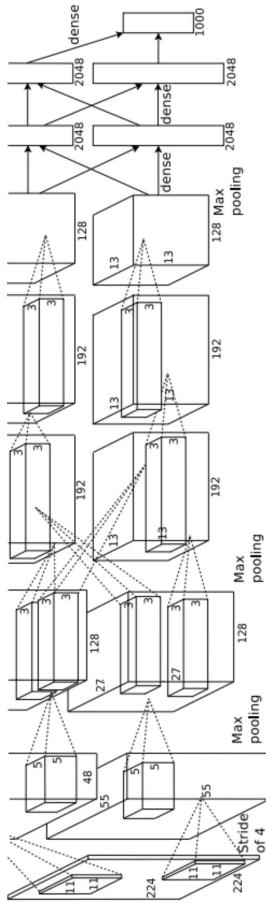


Layer	C	H / W	filters	kernel	stride	pad	C	H / W	Output size	memory (KB)	params (k)	flop (M)
conv1	3	227	64	11	4	2	64	56	784	23	73	
pool1	64	56		3	2	0	64	27	182	0	0	
conv2	64	27	192	5	1	2	192	27	547	307	224	
pool2	192	27		3	2	0	192	13	127	0	0	
conv3	192	13	384	3	1	1	384	13	254	664	112	
conv4	384	13	256	3	1	1	256	13	169	885	145	
conv5	256	13	256	3	1	1	256	13	169	590	100	
pool5	256	13		3	2	0	256	6	36	0	0	
flatten	256	6					9216		36	0	0	
fc6	9216		4096				4096		16	37,749	38	
fc7	4096		4096				4096		16	16,777	17	
fc8	4096		1000				1000		4	4,096	4	

AlexNet

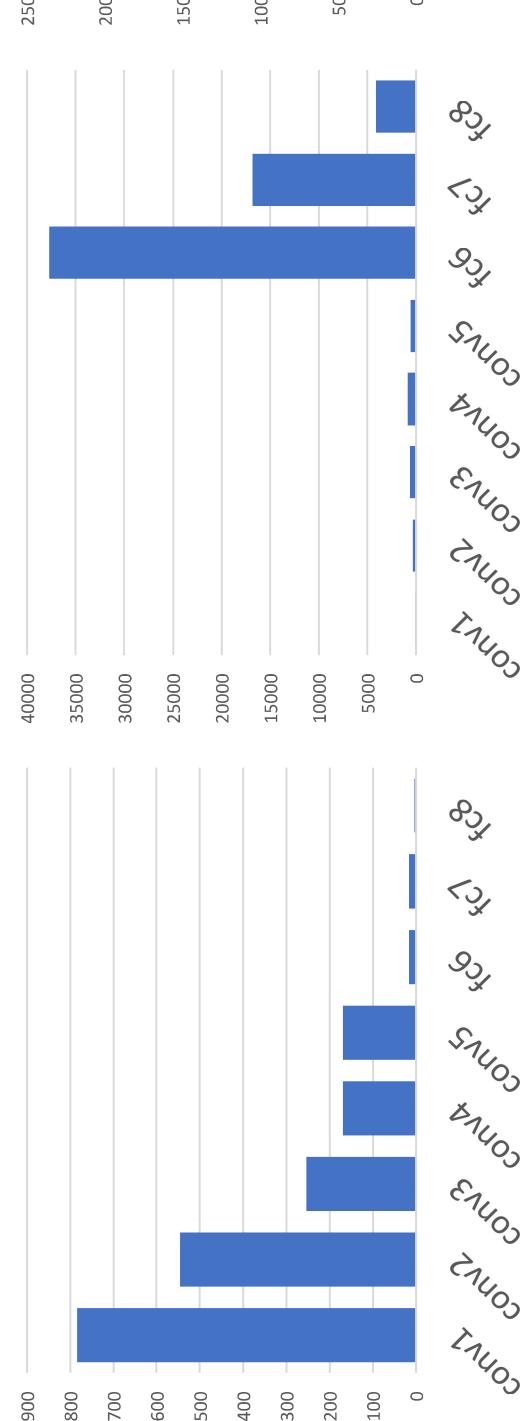
Most of the **memory usage** is in the early convolution layers

Nearly all **parameters** are in the fully-connected layers



Most floating-point ops occur in the convolution layers

MFLOP

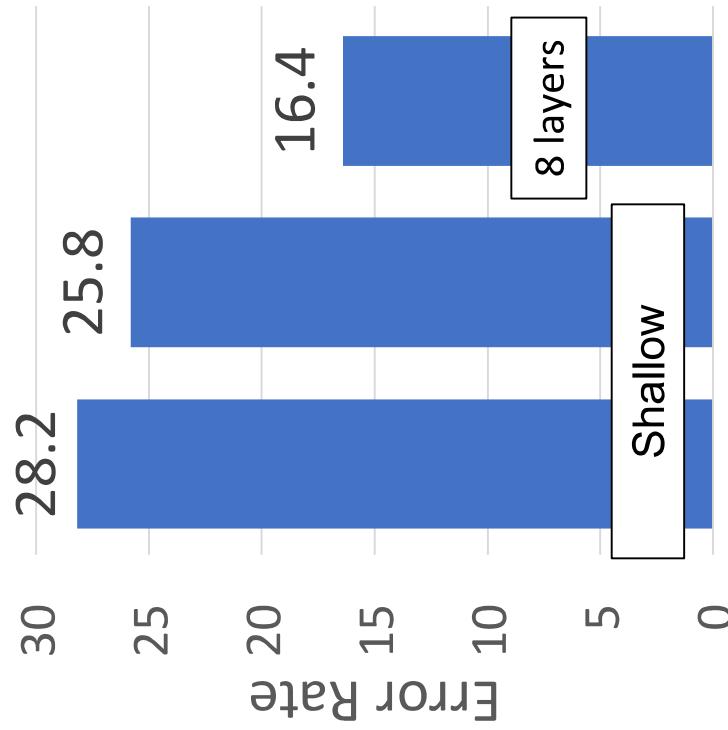


Justin Johnson

Lecture 8 - 31

September 30, 2019

ImageNet Classification Challenge



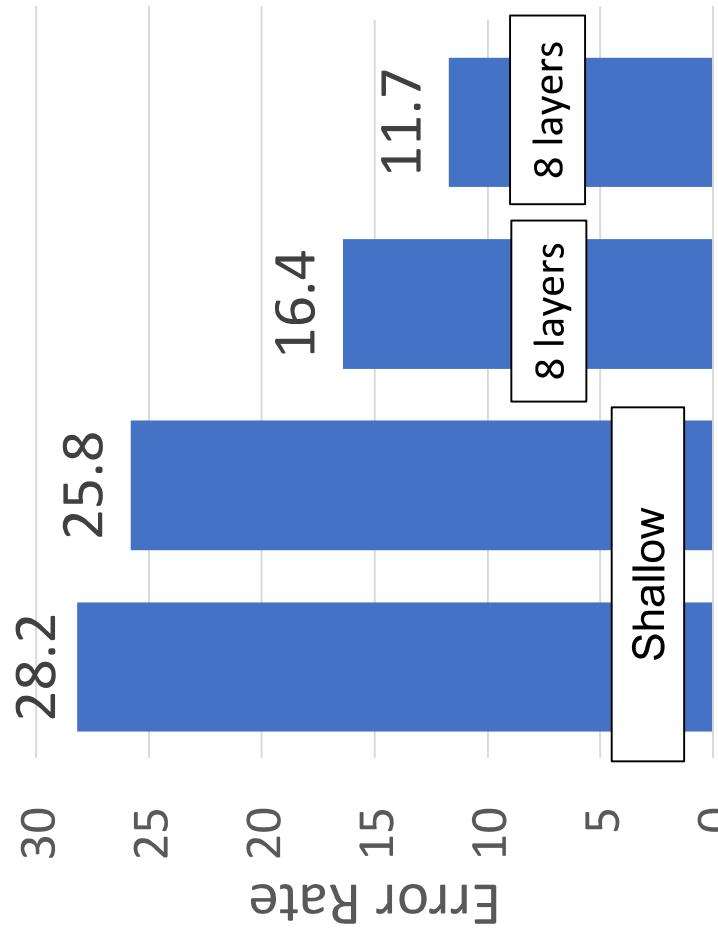
Lin et al Sanchez & Perronnin Krizhevsky et al
(AlexNet)

Justin Johnson

Lecture 8 - 32

September 30, 2019

ImageNet Classification Challenge



Lin et al Sanchez & Perronnin Krizhevsky et al (AlexNet) Zeiler & Fergus

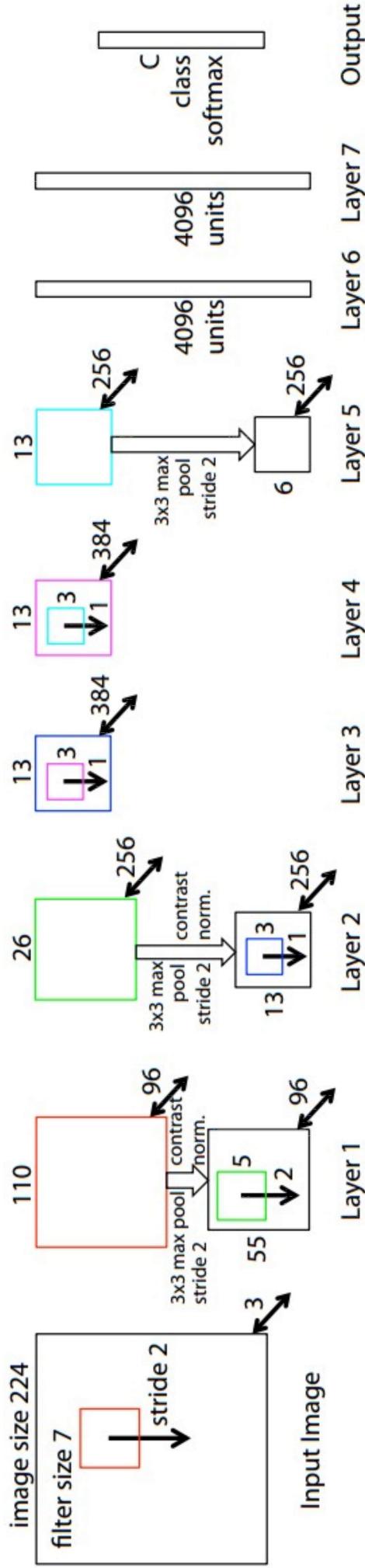
Justin Johnson

Lecture 8 - 33

September 30, 2019

ZFNet: A Bigger AlexNet

ImageNet top 5 error: 16.4% -> 11.7%



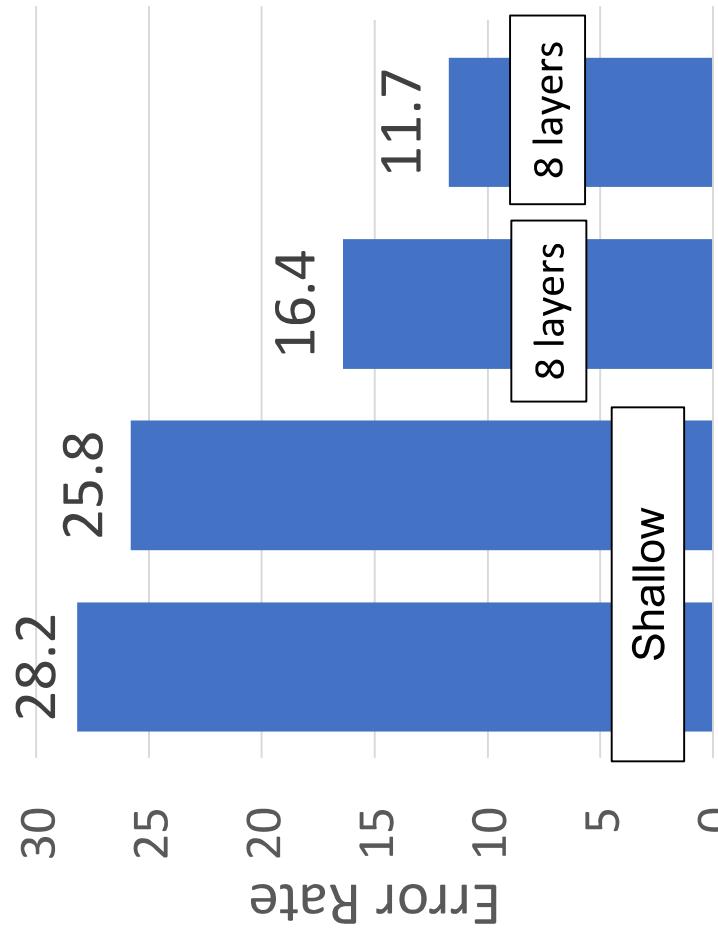
AlexNet but:

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

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

More trial and error = (

ImageNet Classification Challenge



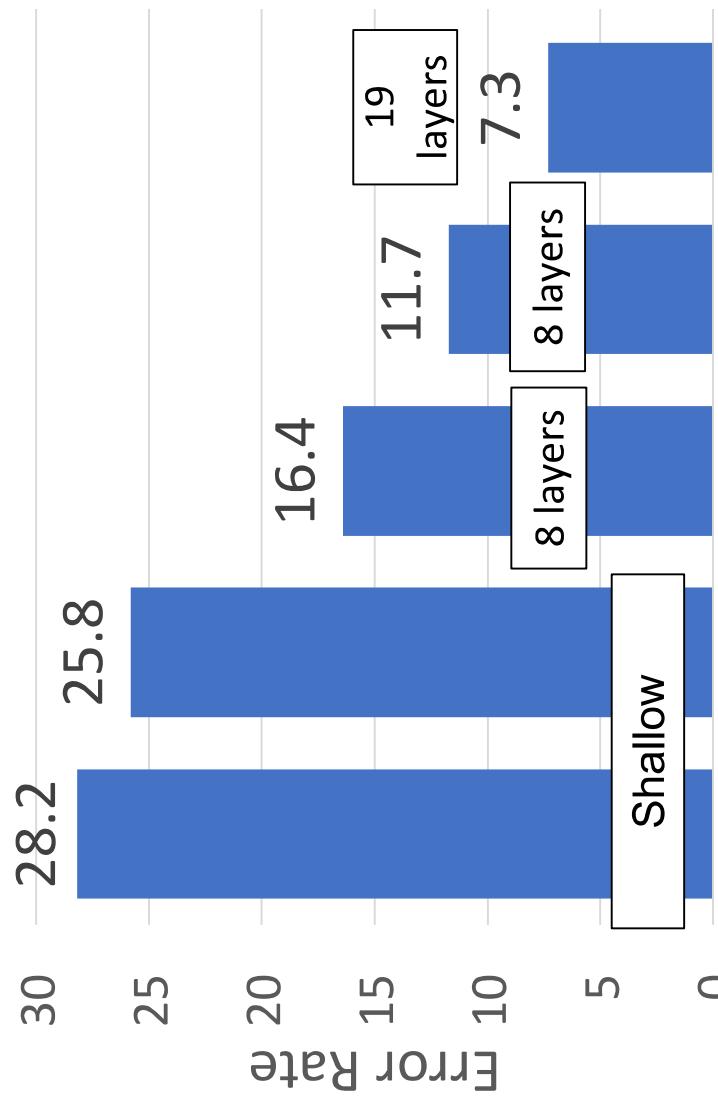
Lin et al Sanchez & Perronnin Krizhevsky et al (AlexNet) Zeiler & Fergus

Justin Johnson

Lecture 8 - 35

September 30, 2019

ImageNet Classification Challenge



Justin Johnson
Lecture 8 - 36

Lin et al Sanchez & Perronnin Krizhevsky et al (AlexNet) Zeiler & Fergus Simonyan & Zisserman (VGG)

September 30, 2019

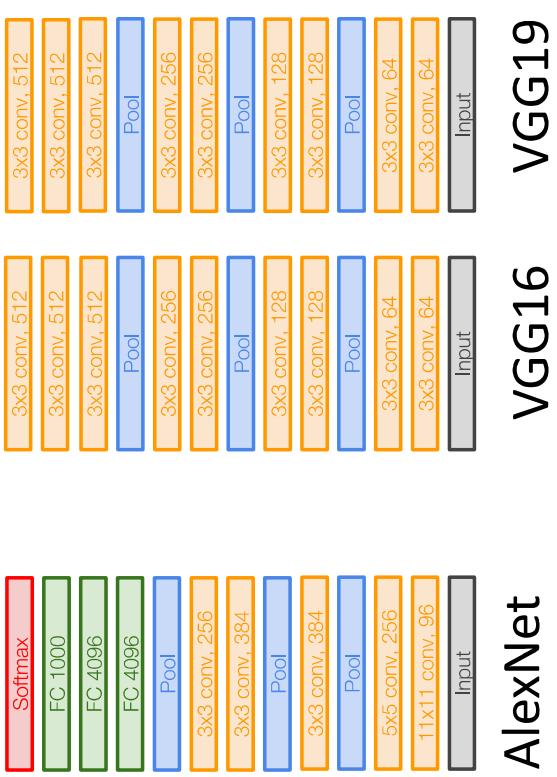
VGG: Deeper Networks, Regular Design

VGG Design rules:

All conv are 3×3 stride 1 pad 1

All max pool are 2×2 stride 2

After pool, double #channels



VGG: Deeper Networks, Regular Design

VGG Design rules:

All conv are 3×3 stride 1 pad 1

All max pool are 2×2 stride 2

After pool, double #channels

Network has 5 convolutional stages:

Stage 1: conv-conv-pool

Stage 2: conv-conv-pool

Stage 3: conv-conv-pool

Stage 4: conv-conv-conv-[conv]-pool

Stage 5: conv-conv-conv-[conv]-pool

(VGG-19 has 4 conv in stages 4 and 5)

AlexNet

VGG16

VGG19

Simonyan and Zissermann, "Very Deep Convolutional Networks for Large-Scale Image Recognition", ICLR 2015

VGG: Deeper Networks, Regular Design

VGG Design rules:

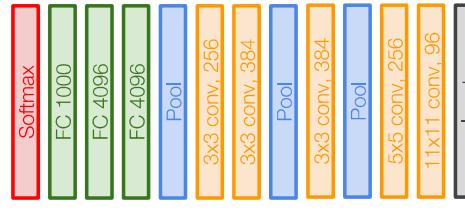
All conv are 3x3 stride 1 pad 1

All max pool are 2x2 stride 2

After pool, double #channels

Option 1:

Conv(5x5, C -> C)



FLOPs: $25C^2HW$

AlexNet

VGG16

VGG19

Simonyan and Zissermann, "Very Deep Convolutional Networks for Large-Scale Image Recognition", ICLR 2015

VGG: Deeper Networks, Regular Design

VGG Design rules:

All conv are 3x3 stride 1 pad 1

All max pool are 2x2 stride 2

After pool, double #channels

Option 1:

Conv(5x5, C -> C)
Conv(3x3, C -> C)

Option 2:

Conv(3x3, C -> C)
Conv(3x3, C -> C)

Params: $25C^2$
FLOPs: $25C^2HW$

Params: $18C^2$
FLOPs: $18C^2HW$

AlexNet

VGG16 VGG19

Simonyan and Zissermann, "Very Deep Convolutional Networks for Large-Scale Image Recognition", ICLR 2015

VGG: Deeper Networks, Regular Design

VGG Design rules:

All conv are 3x3 stride 1 pad 1

All max pool are 2x2 stride 2

After pool, double #channels

Two 3x3 conv has same receptive field as a single 5x5 conv, but has fewer parameters and takes less computation!

Option 1:

Conv(5x5, C -> C)

Conv(3x3, C -> C)

Params: $25C^2$

FLOPs: $25C^2\text{HW}$

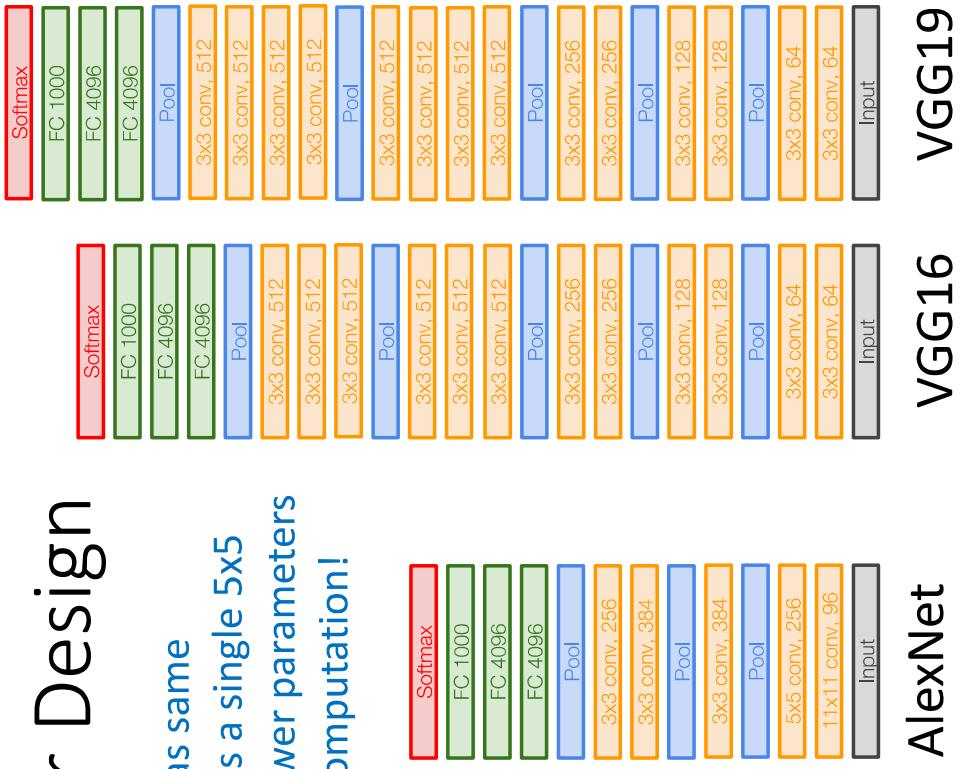
Option 2:

Conv(3x3, C -> C)

Conv(3x3, C -> C)

Params: $18C^2$

FLOPs: $18C^2\text{HW}$



Simonyan and Zissermann, "Very Deep Convolutional Networks for Large-Scale Image Recognition", ICLR 2015

VGG: Deeper Networks, Regular Design

VGG Design rules:

All conv are 3×3 stride 1 pad 1

All max pool are 2×2 stride 2

After pool, double #channels

Input: $C \times 2H \times 2W$

Layer: $\text{Conv}(3 \times 3, C \rightarrow C)$

Memory: $4HW\mathcal{C}$

Params: $9\mathcal{C}^2$

FLOPs: $36HW\mathcal{C}^2$

AlexNet

VGG16

VGG19

Simonyan and Zissermann, "Very Deep Convolutional Networks for Large-Scale Image Recognition", ICLR 2015

VGG: Deeper Networks, Regular Design

VGG Design rules:

All conv are 3×3 stride 1 pad 1

All max pool are 2×2 stride 2

After pool, double #channels

Input: $2C \times H \times W$
Layer: $\text{Conv}(3 \times 3, C \rightarrow 2C)$

Memory: $4HWC$
Params: $9C^2$
FLOPs: $36HWC^2$

Memory: $2HWC$
Params: $36C^2$
FLOPs: $36HWC^2$

AlexNet VGG16 VGG19

Simonyan and Zissermann, "Very Deep Convolutional Networks for Large-Scale Image Recognition", ICLR 2015

VGG: Deeper Networks, Regular Design

VGG Design rules:

All conv are 3x3 stride 1 pad 1

All max pool are 2×2 stride 2

After pool, double #channels

Input: $2C \times H \times W$
Layer: Conv(3x3, C->C)
Conv(3x3, 2C -> 2C)

Memory: 4HWC

Params: 9C²

FLOPs: $36HW\mathcal{C}^2$

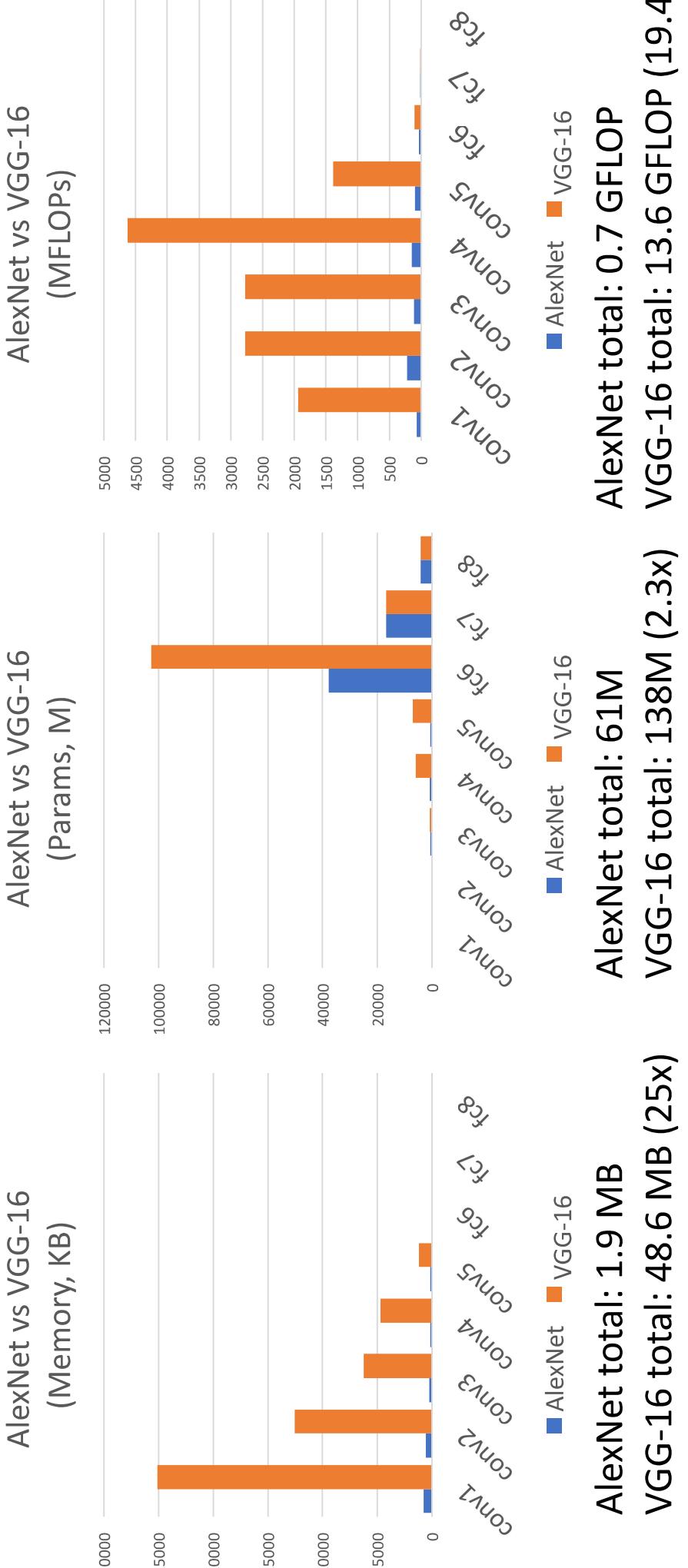
Justin Johnson

Lecture 8 - 44

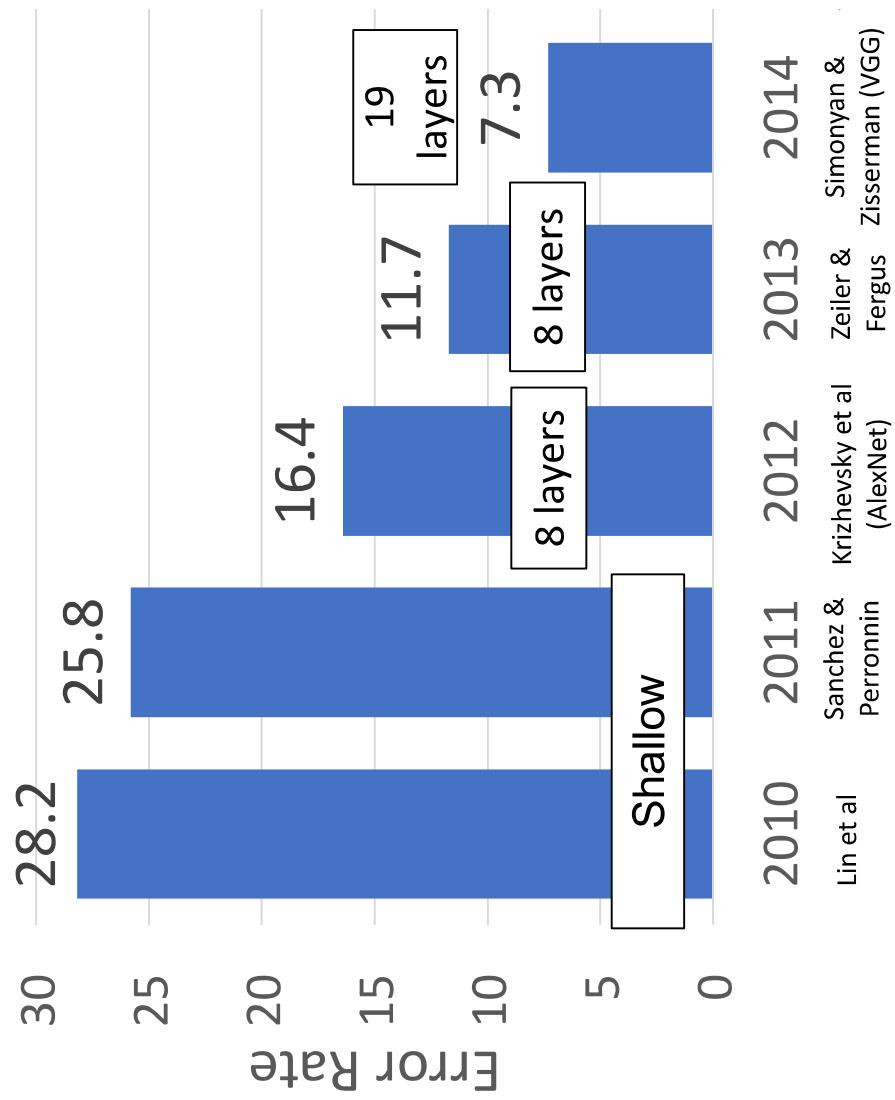
卷之三

September 30, 2019

AlexNet vs VGG-16: Much bigger network!



ImageNet Classification Challenge

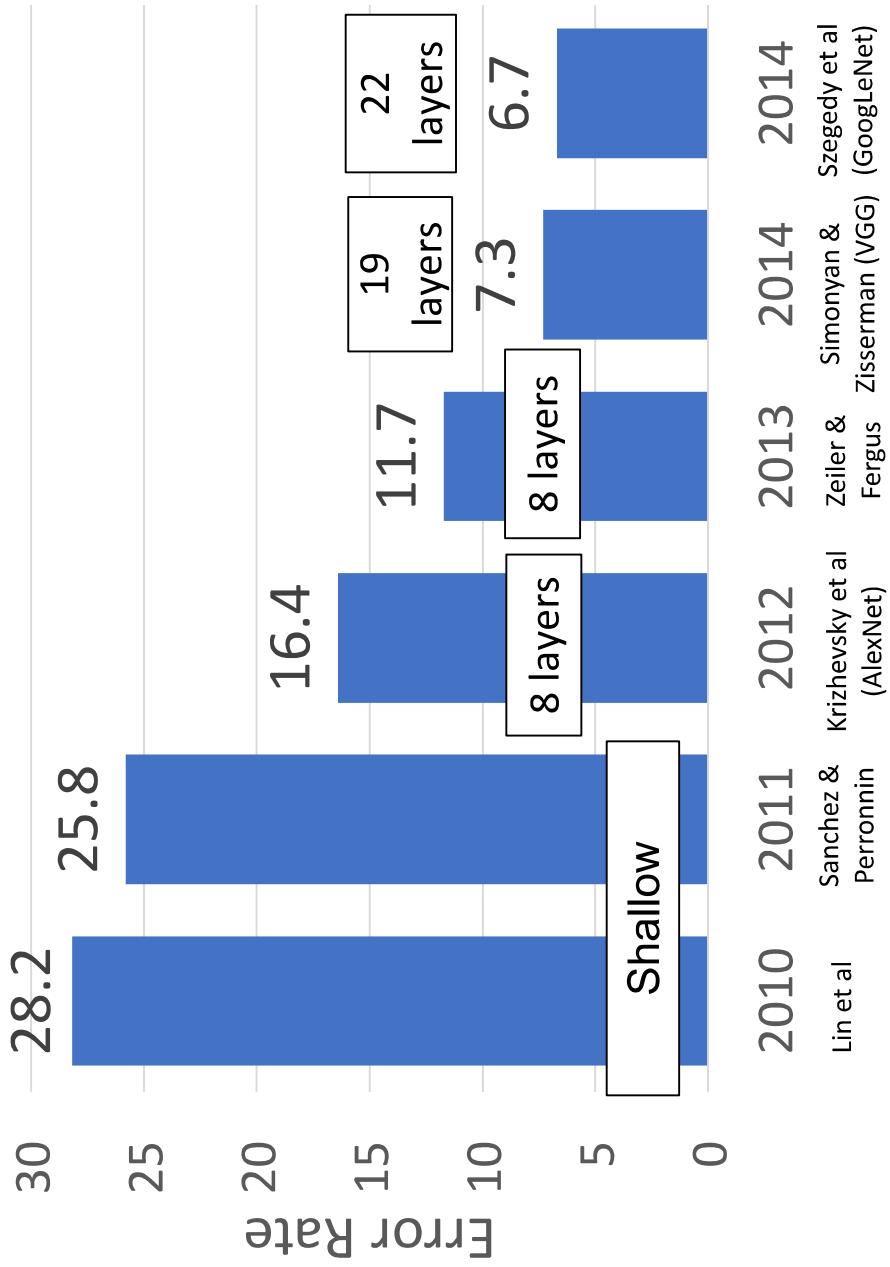


Justin Johnson

Lecture 8 - 46

September 30, 2019

ImageNet Classification Challenge



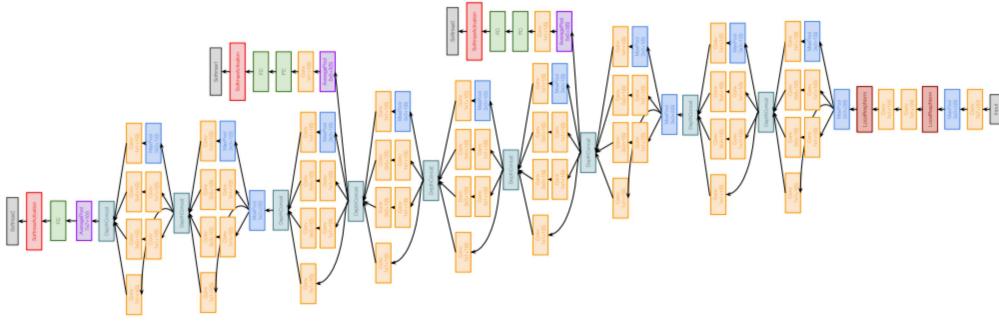
Justin Johnson

Lecture 8 - 47

September 30, 2019

GoOGLeNet: Focus On Efficiency

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



Szegedy et al, "Going deeper with convolutions", CVPR 2015

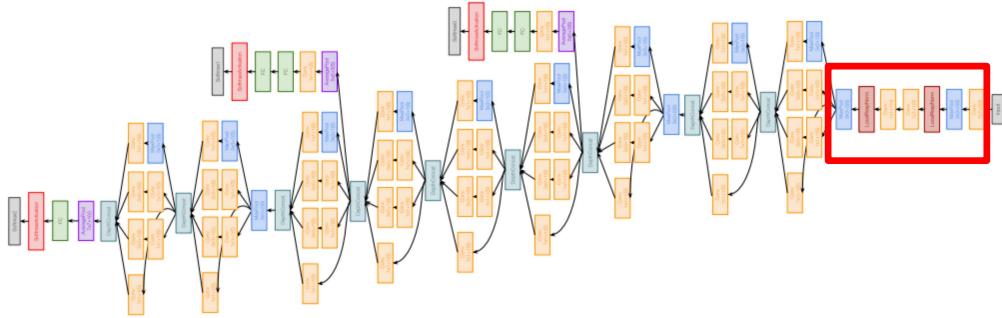
Justin Johnson

Lecture 8 - 48

September 30, 2019

GoOGLeNet: Aggressive Stem

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



Szegedy et al, "Going deeper with convolutions", CVPR 2015

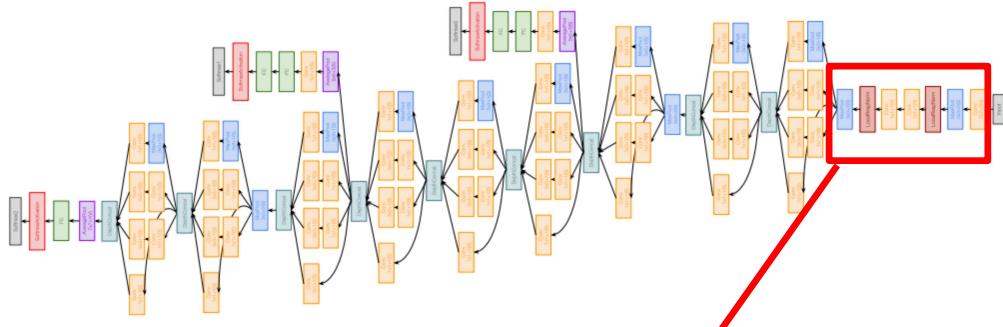
Justin Johnson

Lecture 8 - 49

September 30, 2019

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		
	C	H	W	filters	kernel	stride	pad	C	H/W
conv	3	224	224	64	7	2	3	64	112
max-pool	64	112	112	64	3	2	1	64	56
conv	64	56	56	64	1	1	0	64	56
conv	64	56	56	192	3	1	1	192	56
max-pool	192	56	56	3	2	1	1	192	28
								588	0
									1

Total from 224 to 28 spatial resolution:

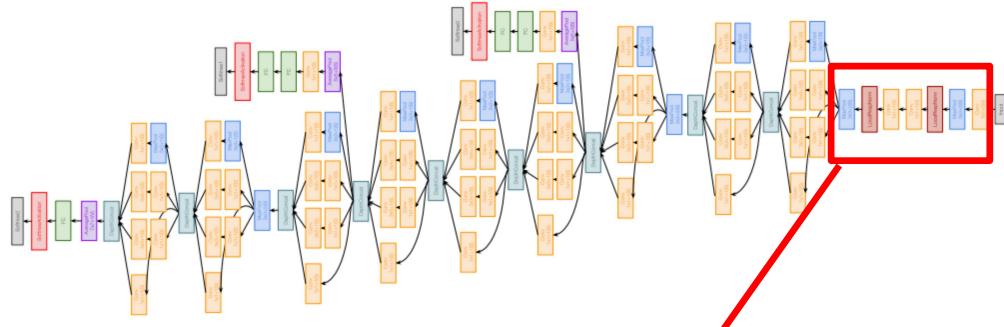
Memory: 7.5 MB

Params: 124K

MFLOP: 418

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		
	C	H	W	filters	kernel	stride	pad	C	H/W
conv	3	224	224	64	7	2	3	64	112
max-pool	64	112	112	64	3	2	1	64	56
conv	64	56	56	64	1	1	0	64	56
conv	64	56	56	192	3	1	1	192	56
max-pool	192	56	56	3	2	1	1	192	28
								588	0
									1

Total from 224 to 28 spatial resolution:

Memory: 7.5 MB

Params: 124K

MFLOP: 418

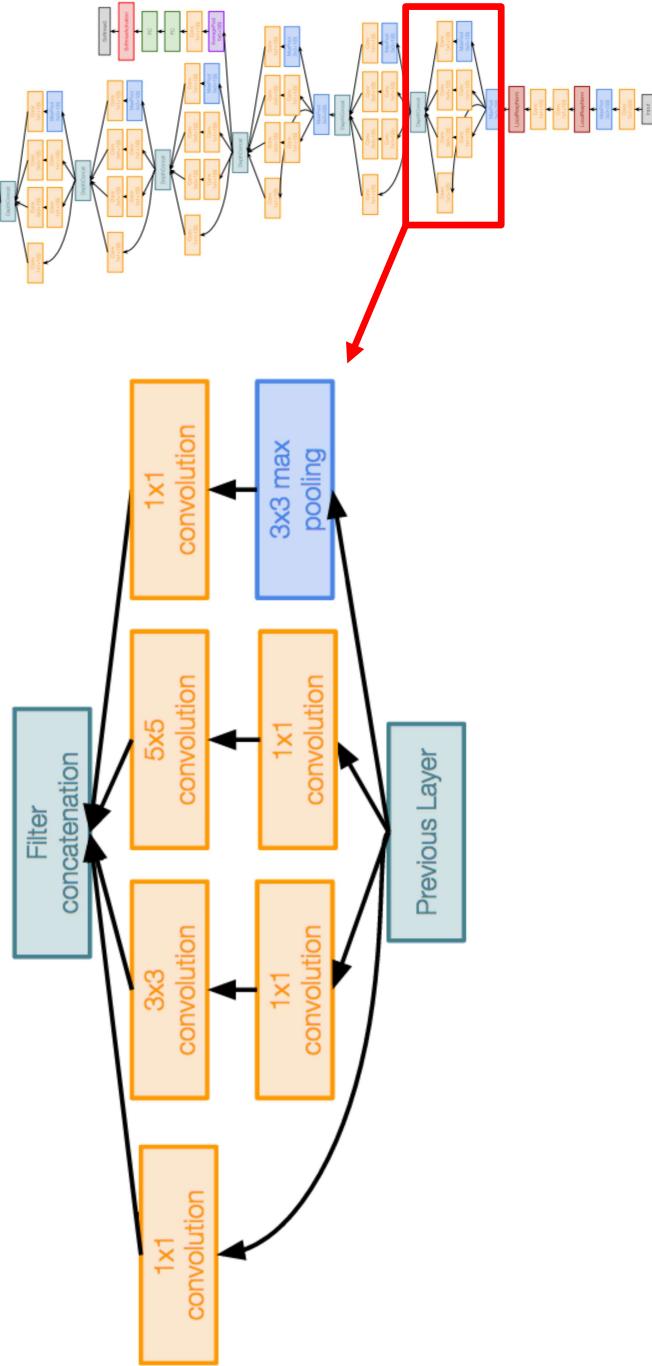
Compare VGG-16:
Memory: 42.9 MB (5.7x)
Params: 1.1M (8.9x)
MFLOP: 7485 (17.8x)

GoOGLeNet: Inception Module

Inception module

Local unit with parallel branches

Local structure repeated many times throughout the network



Szegedy et al., "Going deeper with convolutions", CVPR 2015

Justin Johnson

Lecture 8 - 52

September 30, 2019

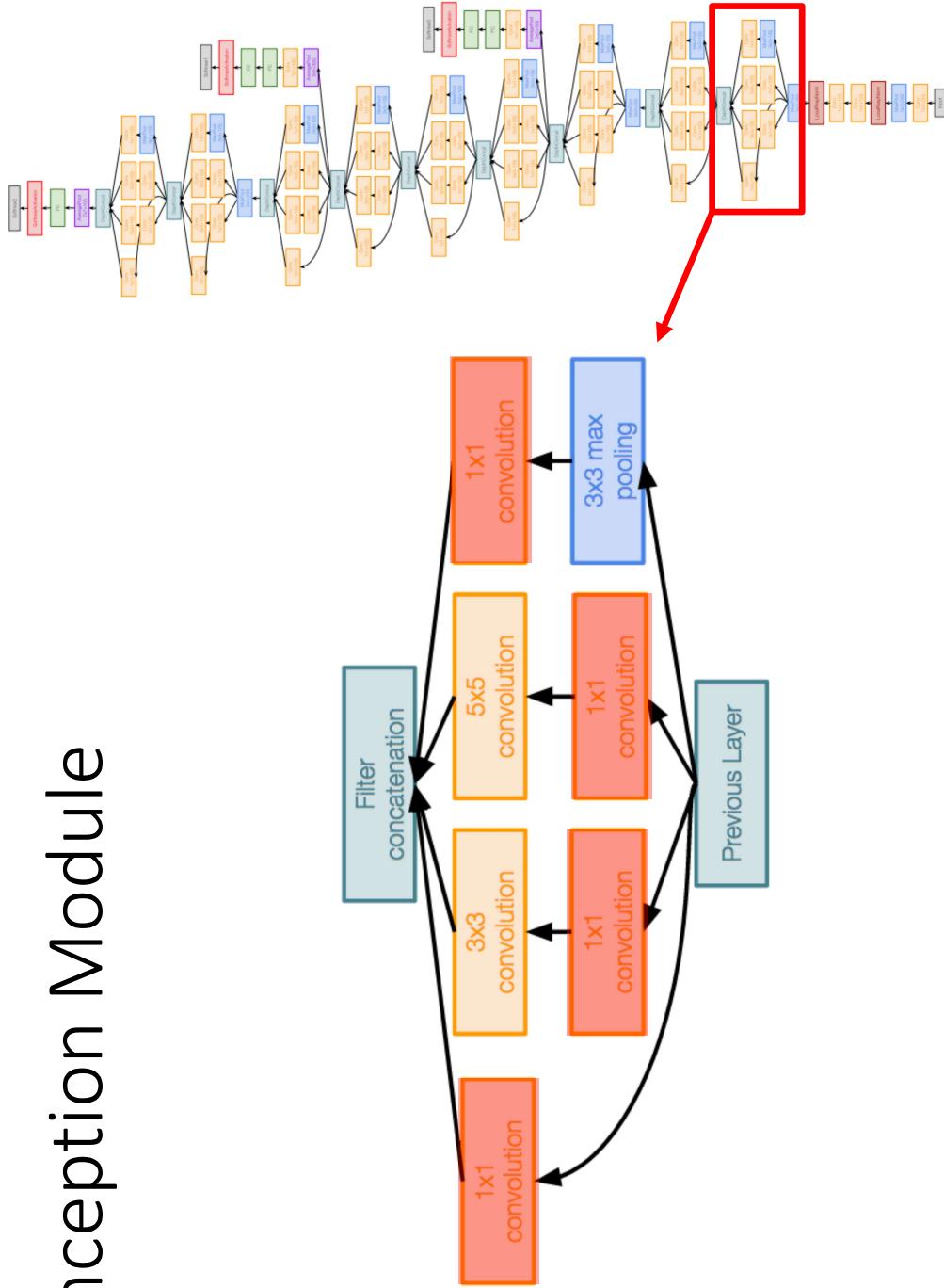
GoOGLeNet: Inception Module

Inception module

Local unit with parallel branches

Local structure repeated many times throughout the network

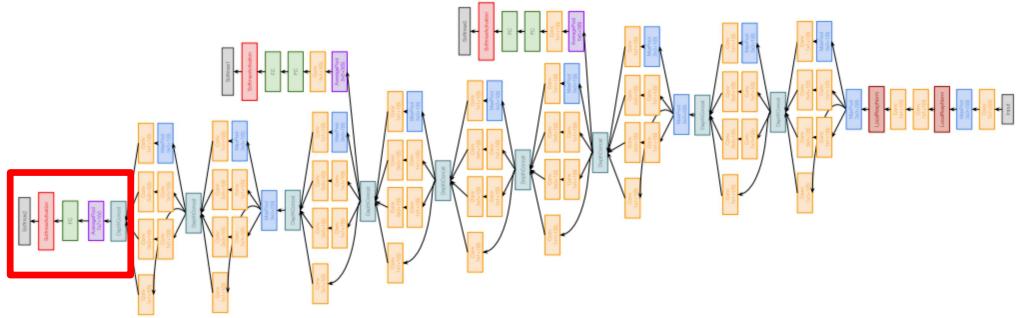
Uses 1x1 “Bottleneck” layers to reduce channel dimension before expensive conv (we will revisit this with ResNet!)



Szegedy et al., “Going deeper with convolutions”, CVPR 2015

GoOGLeNet: Global Average Pooling

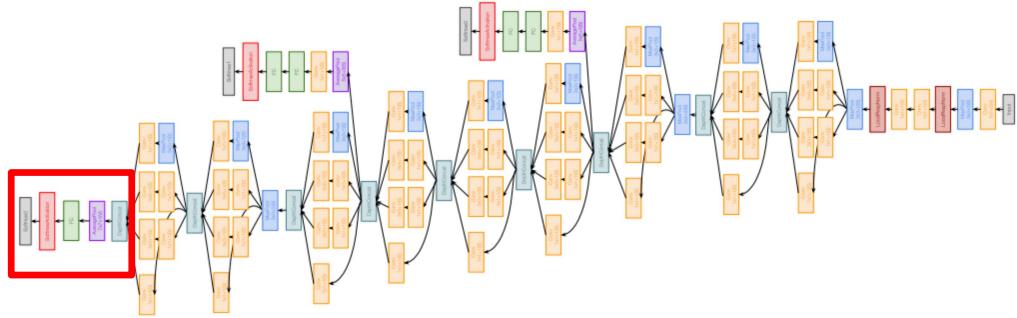
No large FC layers at the end! Instead uses **global average pooling** to collapse spatial dimensions, and one linear layer to produce class scores
(Recall VGG-16: Most parameters were in the FC layers!)



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

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!)

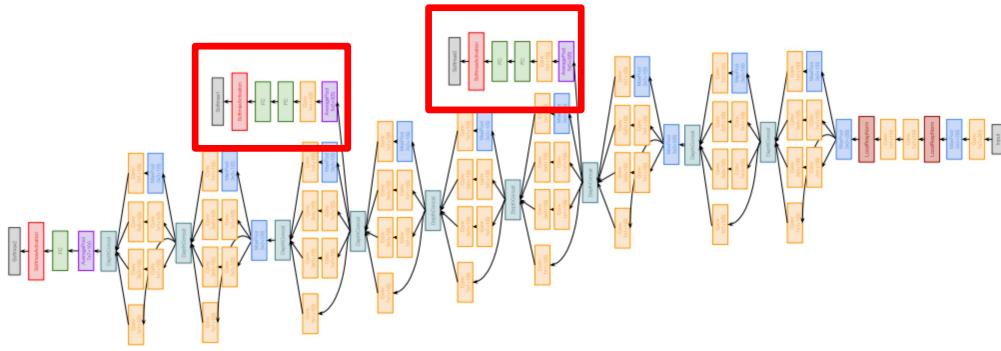


	Input size	Layer			Output size		
Layer	C	H/W	filters	kernel stride	pad	C	H/W
avg-pool	1024	7		7	1	0	1024
fc	1024		1000			1000	

Compare with VGG-16:

	Input size	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

GoOGLeNet: Auxiliary Classifiers

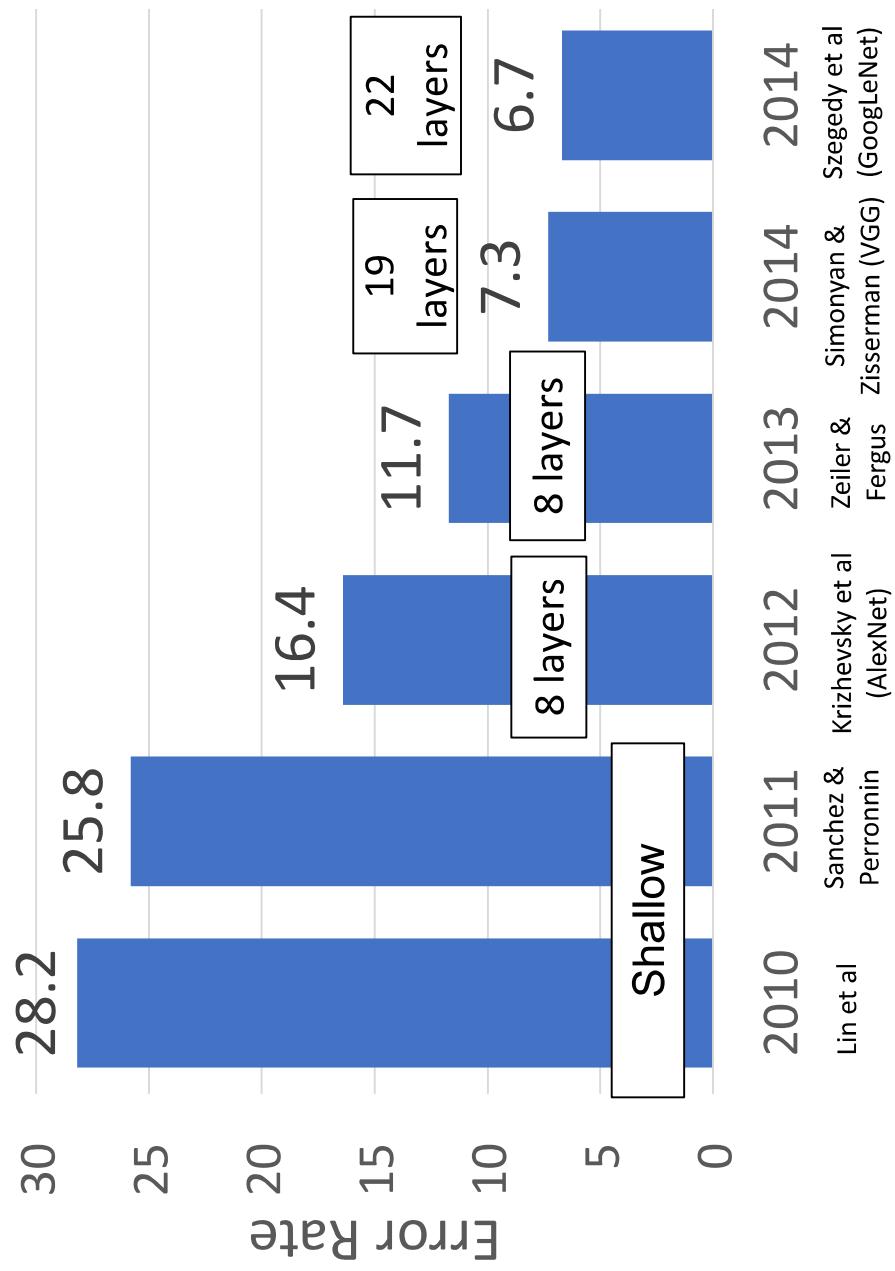


Training using loss at the end of the network didn't work well:
Network is too deep, gradients don't propagate cleanly

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

ImageNet Classification Challenge

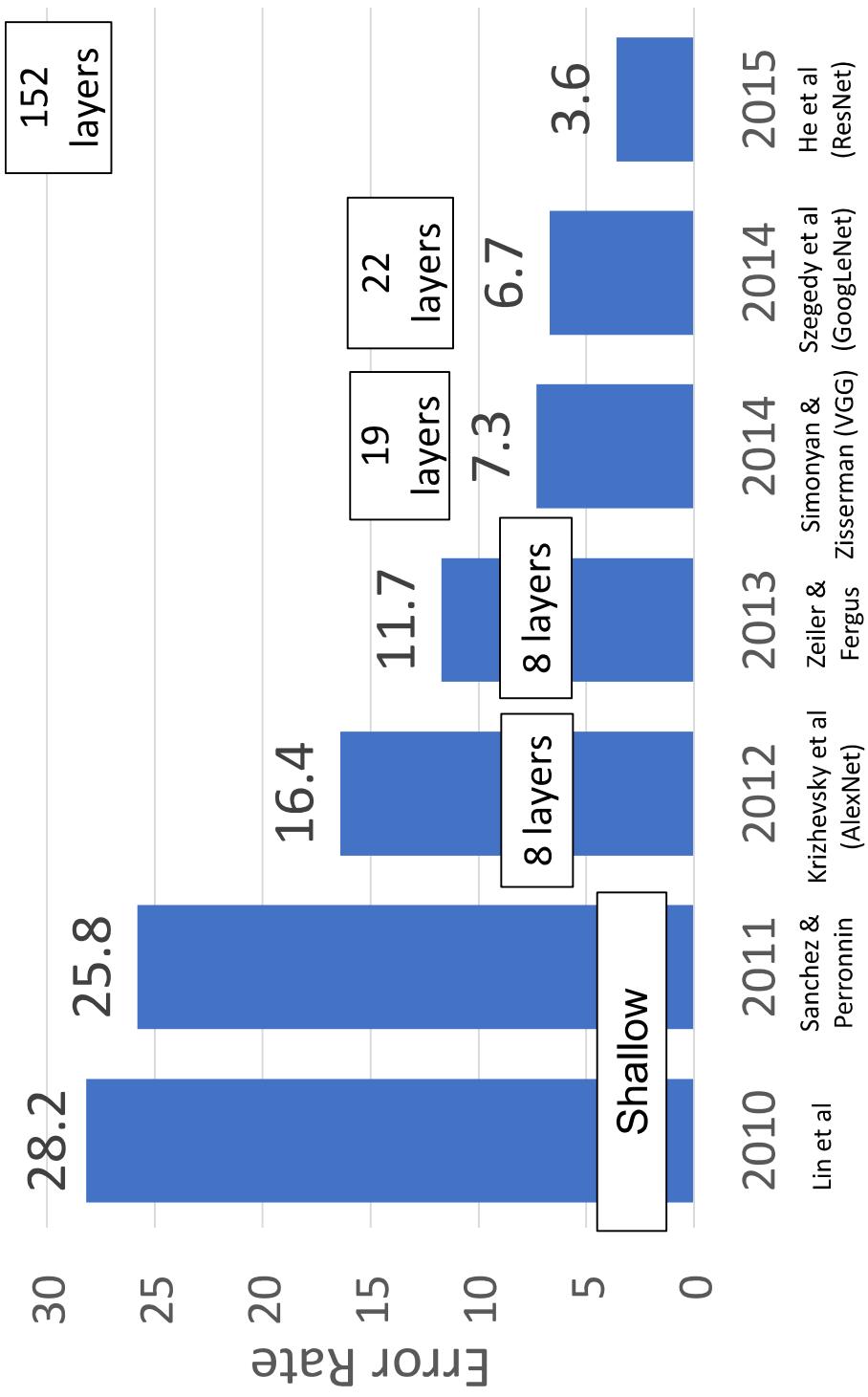


Justin Johnson

Lecture 8 - 57

September 30, 2019

ImageNet Classification Challenge



Justin Johnson

Lecture 8 - 58

September 30, 2019

Lin et al Sanchez & Perronnin Krizhevsky et al (AlexNet) Zeiler & Fergus Simonyan & Zisserman (VGG) Szegedy et al (GoogleNet) He et al (ResNet)

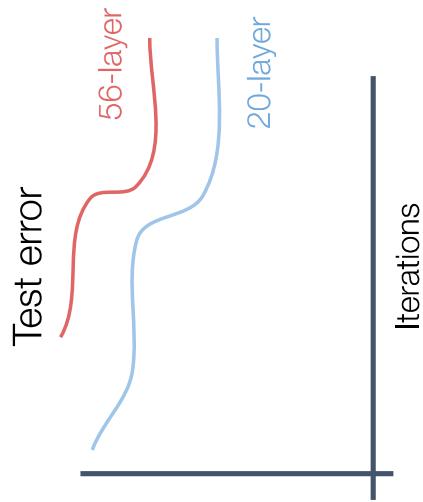
Residual Networks

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

Residual Networks

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

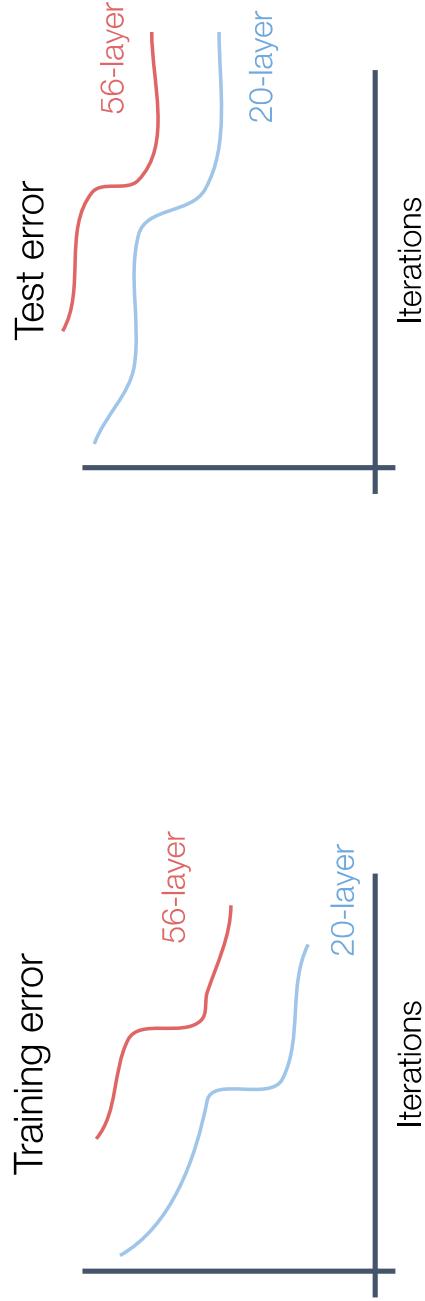
Deeper model does worse than
shallow model!



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

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**

Residual Networks

A deeper model can emulate a shallower model: copy layers from shallower model, set extra layers to identity

Thus deeper models should do at least as good as shallow models

Hypothesis: This is an optimization problem. Deeper models are harder to optimize, and in particular don't learn identity functions to emulate shallow models

Residual Networks

A deeper model can emulate a shallower model: copy layers from shallower model, set extra layers to identity

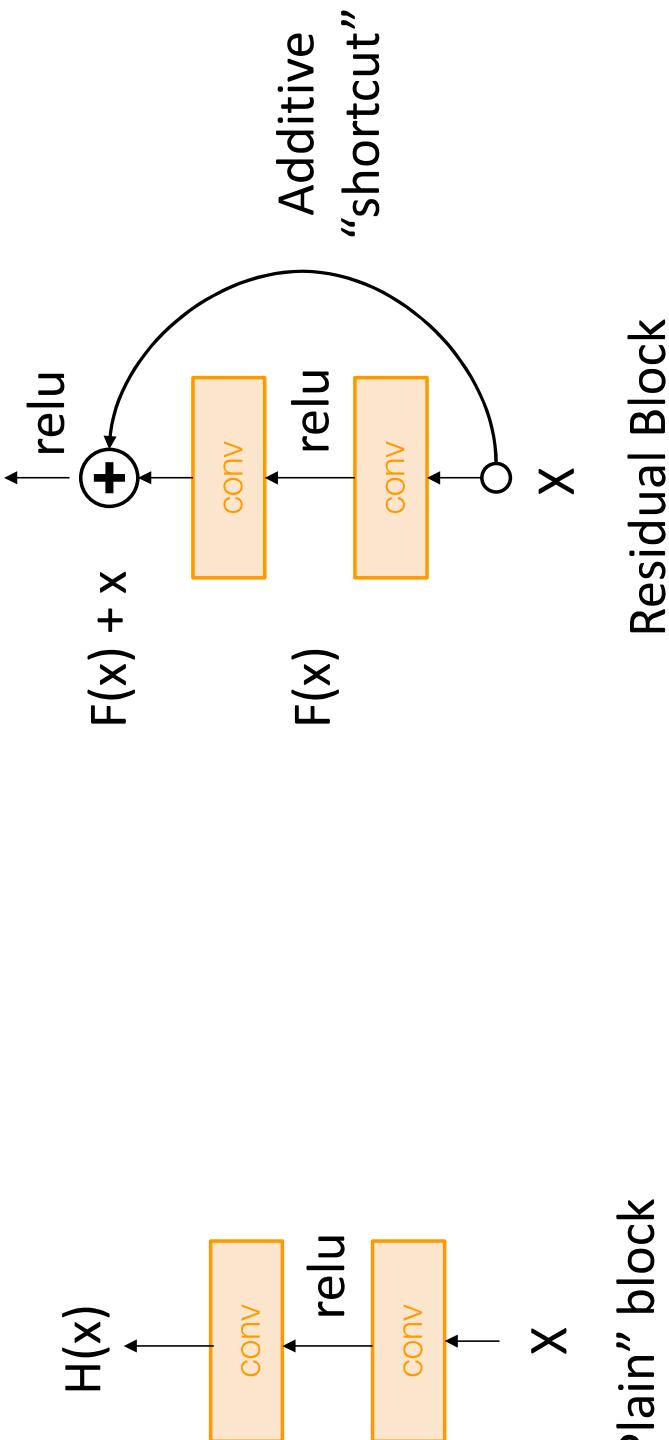
Thus deeper models should do at least as good as shallow models

Hypothesis: This is an optimization problem. Deeper models are harder to optimize, and in particular don't learn identity functions to emulate shallow models

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

Residual Networks

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



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

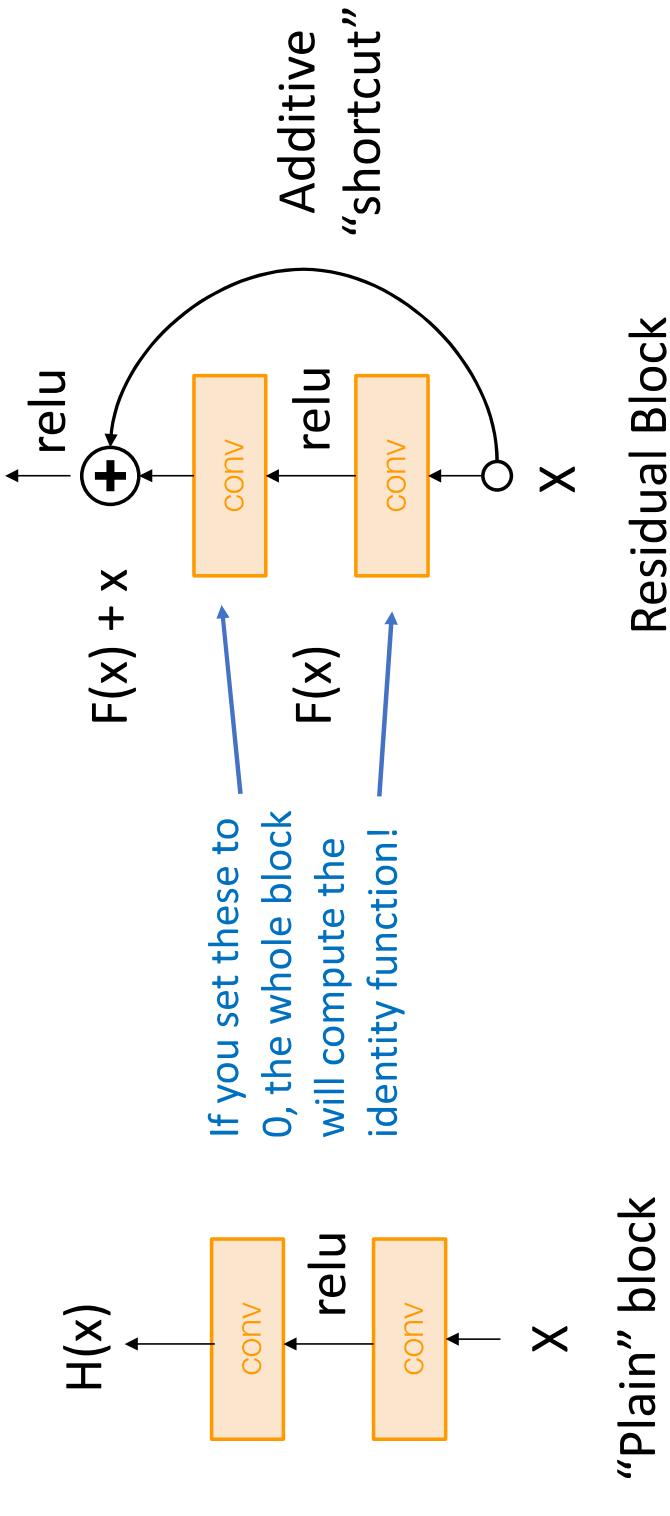
Justin Johnson

Lecture 8 - 64

September 30, 2019

Residual Networks

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

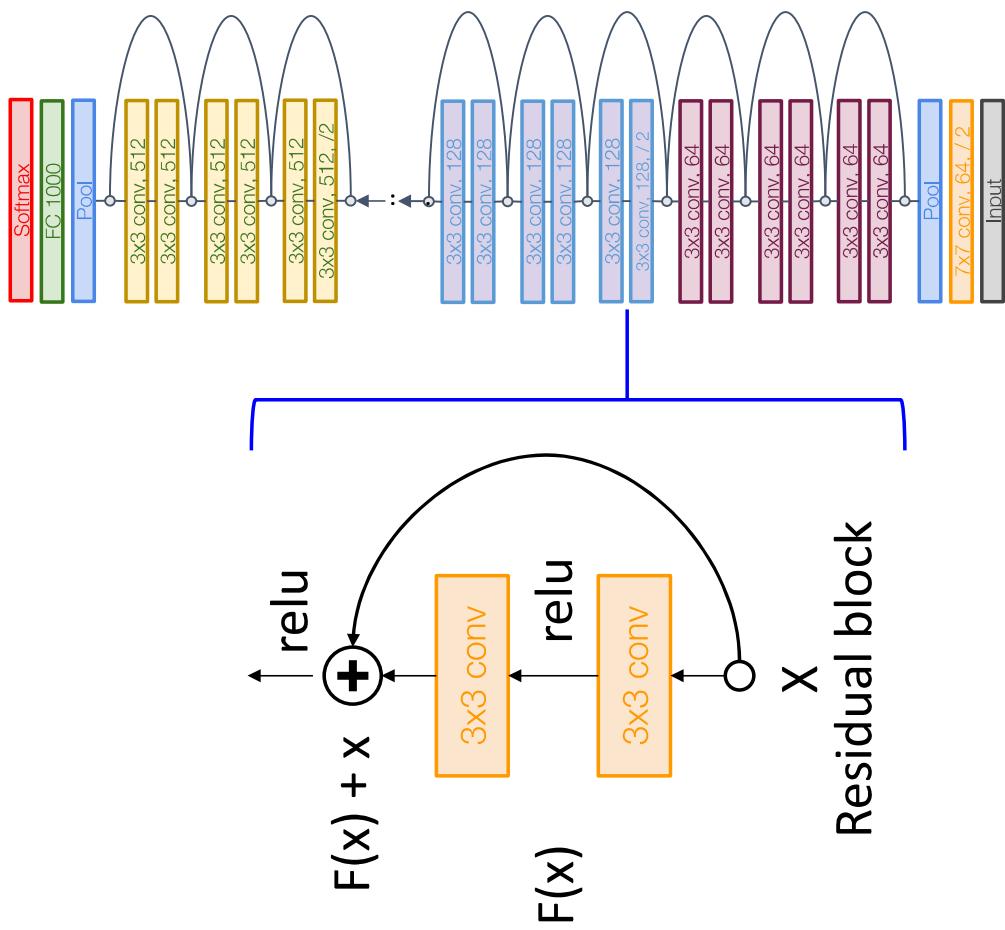


Residual Networks

A residual network is a stack of many residual blocks

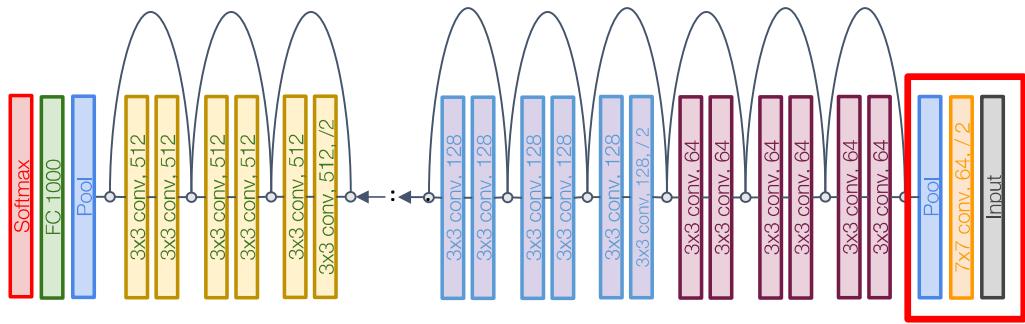
Regular design, like VGG: each residual block has two 3×3 conv

Network is divided into **stages**: the first block of each stage halves the resolution (with stride-2 conv) and doubles the number of channels



Residual Networks

Uses the same aggressive **stem** as GoogleNet to downsample the input 4x before applying residual blocks:

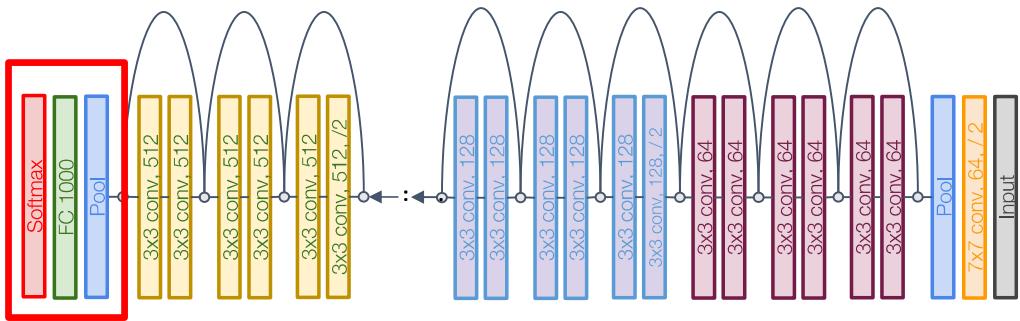


Layer	C	H/W filters	kernel	stride	pad	C	H/W	memory (KB)	Output size		params (k)	flop (M)
									Input size	Layer		
conv	3	224	64	7	2	3	64	112			3136	9 118
max-pool	64	112		3	2	1	64	56			784	0 2

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

Residual Networks

Like GoogLeNet, no big fully-connected-layers: instead use
global average pooling and a single linear layer at the end



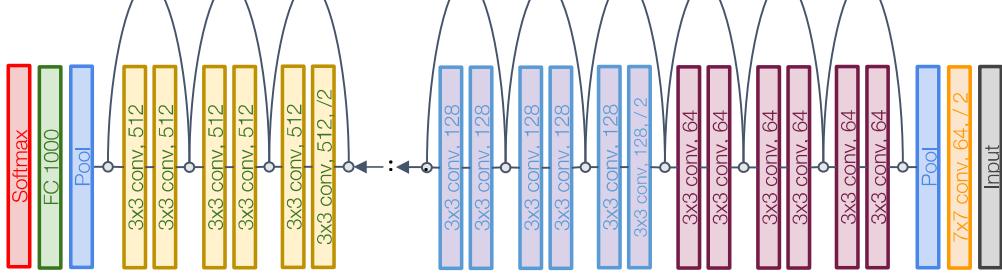
Residual Networks

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



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

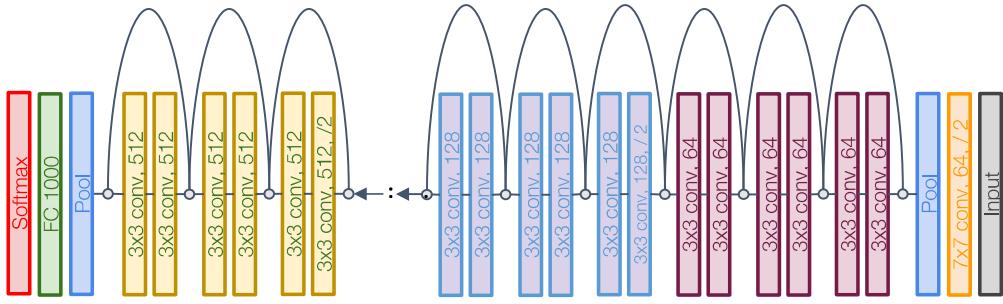
Residual Networks

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

ImageNet top-5 error: 8.58
GFLOP: 3.6



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

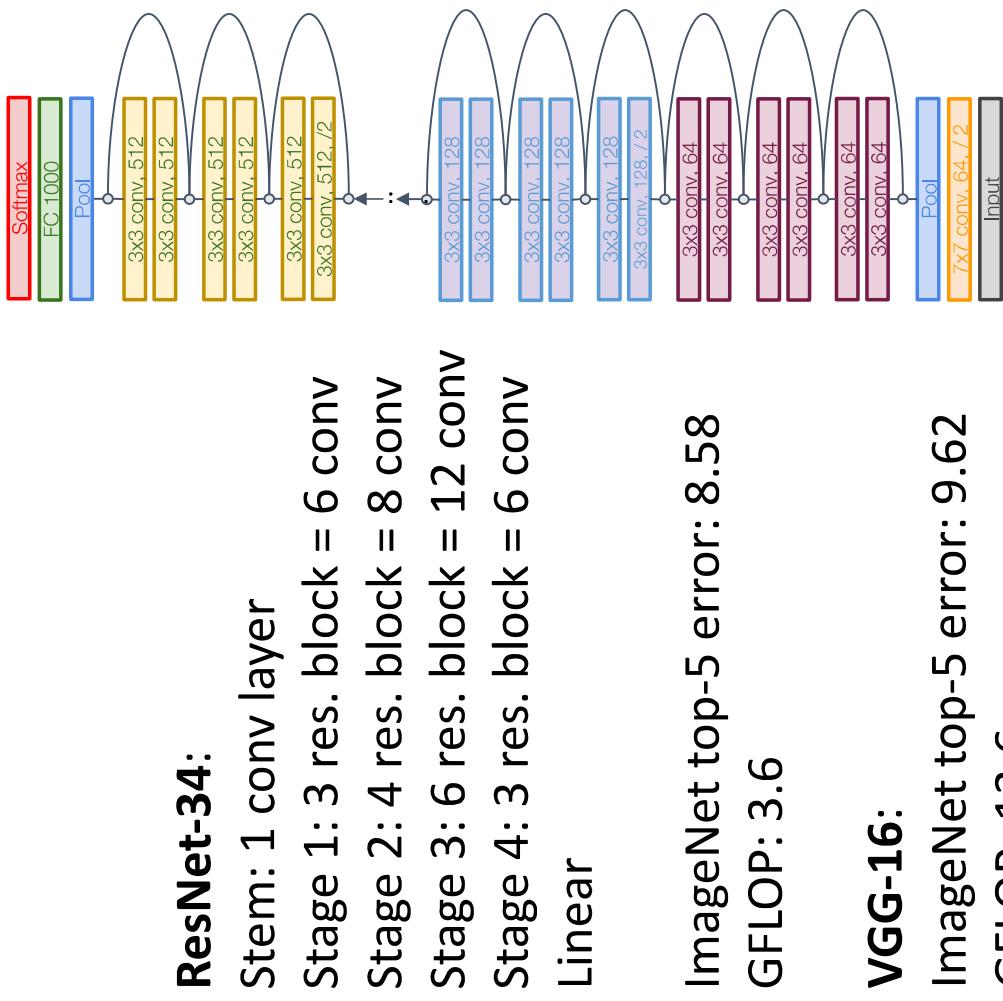
He et al., "Deep Residual Learning for Image Recognition", CVPR 2016
Error rates are 224x224 single-crop testing, reported by [torchvision](#)

Residual Networks

ResNet-18:

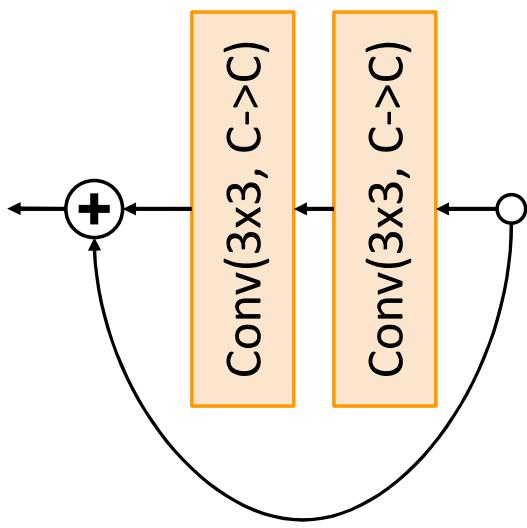
Stem: 1 conv layer
Stage 1 (C=64): 2 res. block = 4 conv
Stage 2 (C=128): 2 res. block = 4 conv
Stage 3 (C=256): 2 res. block = 4 conv
Stage 4 (C=512): 2 res. block = 4 conv
Linear

ImageNet top-5 error: 10.92
GFLOP: 1.8



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

Residual Networks: Basic Block



“Basic”
Residual block

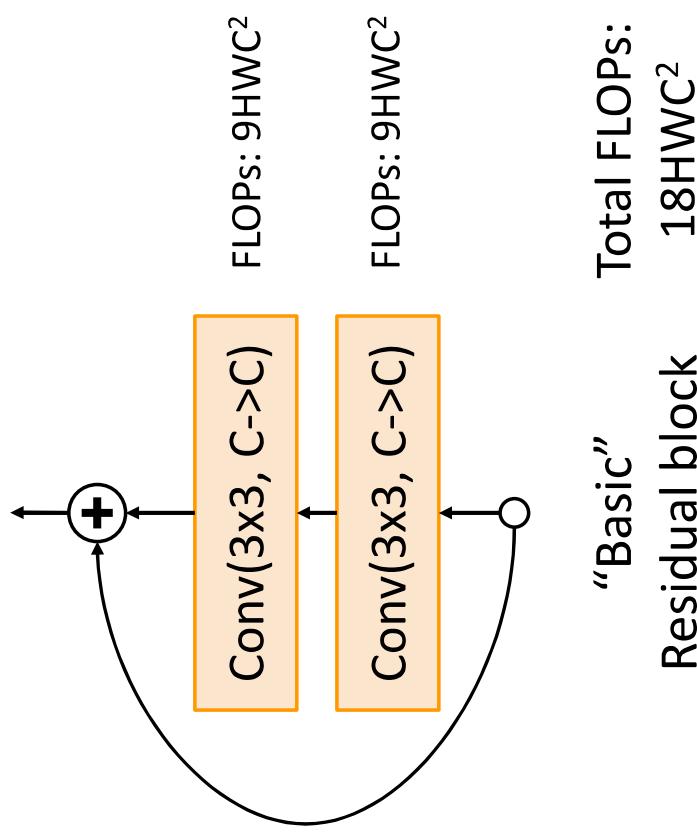
He et al, “Deep Residual Learning for Image Recognition”, CVPR 2016

Justin Johnson

Lecture 8 - 72

September 30, 2019

Residual Networks: Basic Block



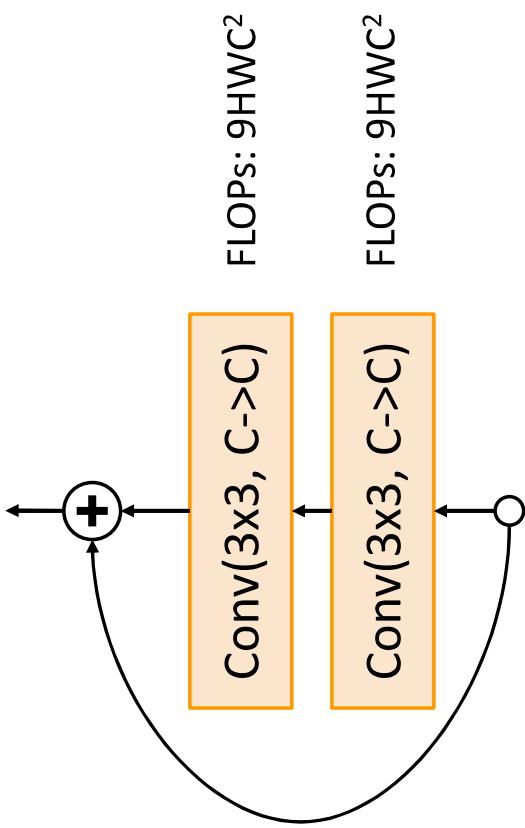
He et al, “Deep Residual Learning for Image Recognition”, CVPR 2016

Justin Johnson

Lecture 8 - 73

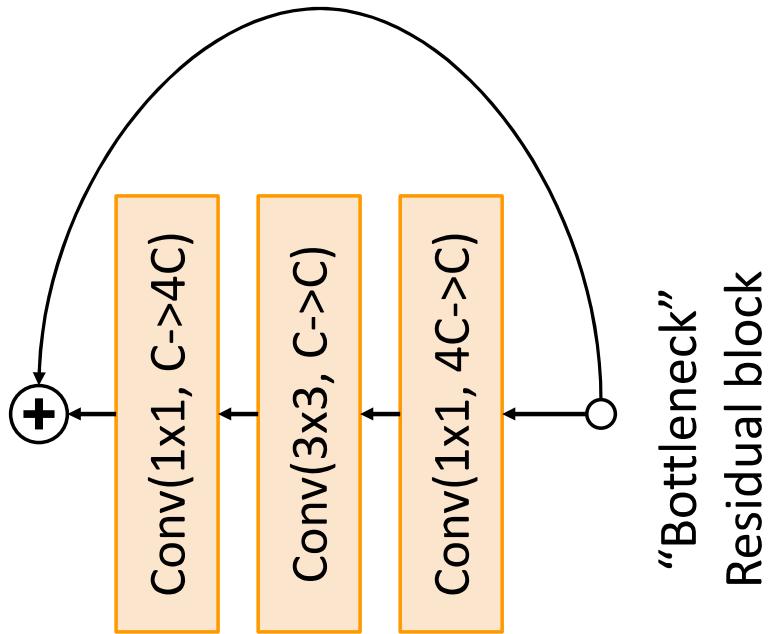
September 30, 2019

Residual Networks: Bottleneck Block



“Basic”
Residual block

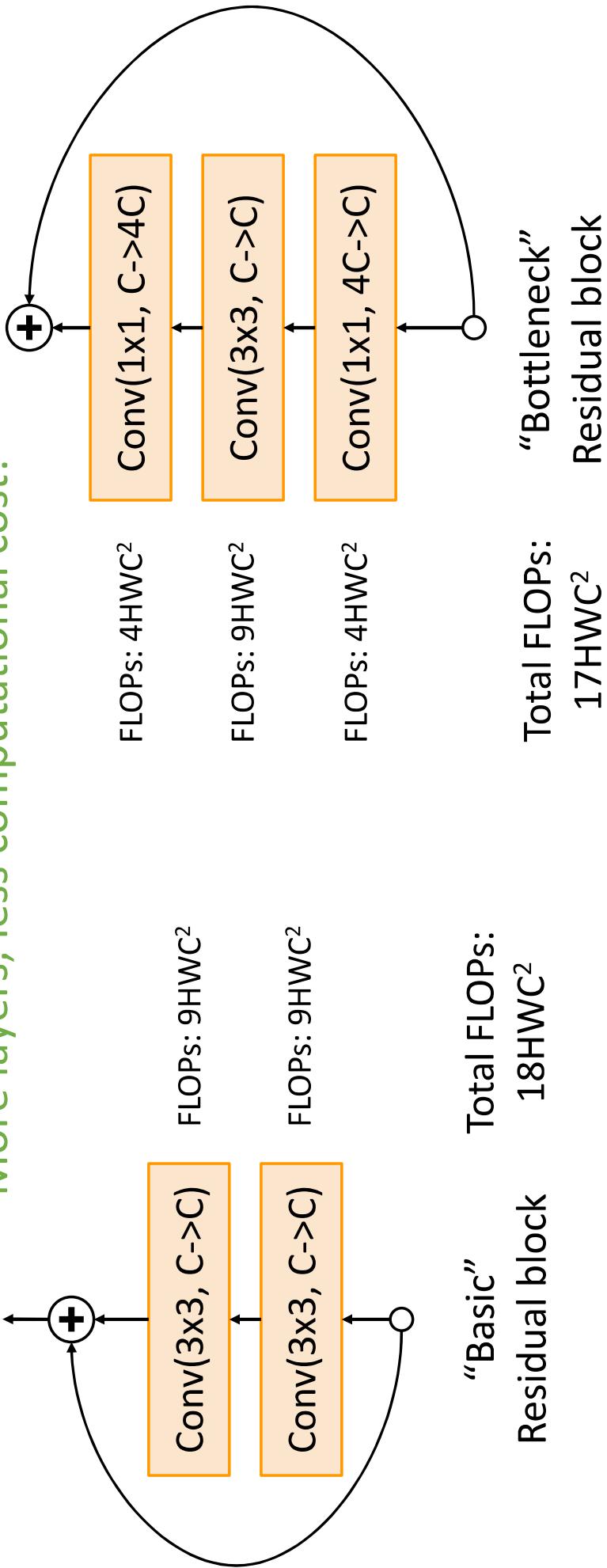
Total FLOPs:
 $18HWC^2$



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

Residual Networks: Bottleneck Block

More layers, less computational cost!



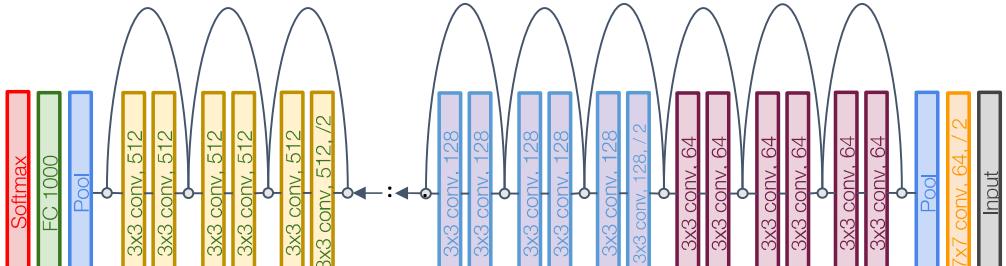
He et al, "Deep Residual Learning for Image Recognition", CVPR 2016

Justin Johnson

Lecture 8 - 75

September 30, 2019

Residual Networks



		Stage 1	Stage 2	Stage 3	Stage 4			
Block type	Stem layers	Blocks	Layers	Blocks	Layers	Blocks	Layers	FC
ResNet-18	Basic	1	2	4	2	4	2	4
ResNet-34	Basic	1	3	6	4	8	6	12

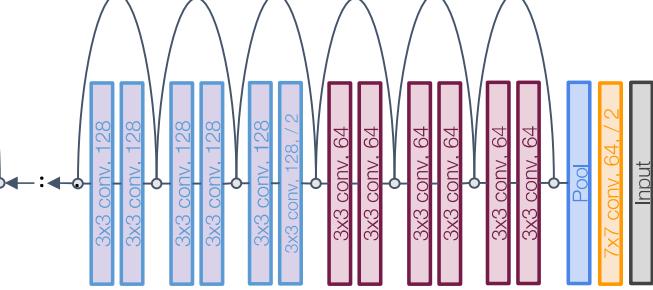
	ImageNet top-5 error	GFLOP
ResNet-18	10.92	1
ResNet-34	8.58	3.6

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

Residual Networks

ResNet-50 is the same as ResNet-34, but replaces Basic blocks with Bottleneck Blocks.
This is a great baseline architecture for many tasks even today!

		Stage 1	Stage 2	Stage 3	Stage 4	FC	
Block type	Stem layers	Blocks	Layers	Blocks	Layers	Blocks	Layers
ResNet-18	Basic	1	2	4	2	4	2
ResNet-34	Basic	1	3	6	4	8	6
ResNet-50	Bottle	1	3	9	4	12	6



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

Residual Networks

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

		Stage 1	Stage 2	Stage 3	Stage 4		
	Block type	Stem layers	Blocks	Layers	Blocks	Layers	FC
ResNet-18	Basic	1	2	4	2	4	2
ResNet-34	Basic	1	3	6	4	8	6
ResNet-50	Bottle	1	3	9	4	12	6
ResNet-101	Bottle	1	3	9	4	12	6
ResNet-152	Bottle	1	3	9	8	24	36

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

Residual Networks

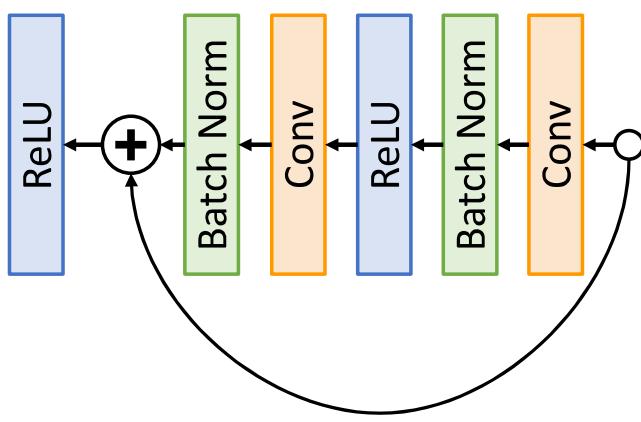
- Able to train very deep networks
- Deeper networks do better than shallow networks (as expected)
- Swept 1st place in all ILSVRC and COCO 2015 competitions
- Still widely used today!

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

Improving Residual Networks: Block Design

Original ResNet block



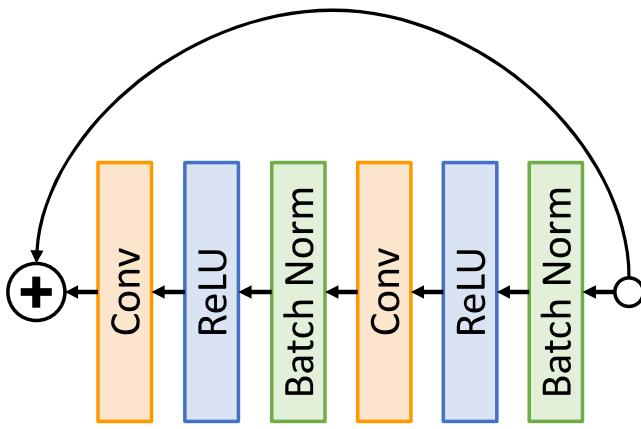
Note ReLU **after** residual:

Cannot actually learn
identity function since
outputs are nonnegative!

Note ReLU **inside** residual:

Note ReLU inside residual:
Can learn true identity
function by setting Conv
weights to zero!

“Pre-Activation” ResNet Block



Note ReLU **after** residual:

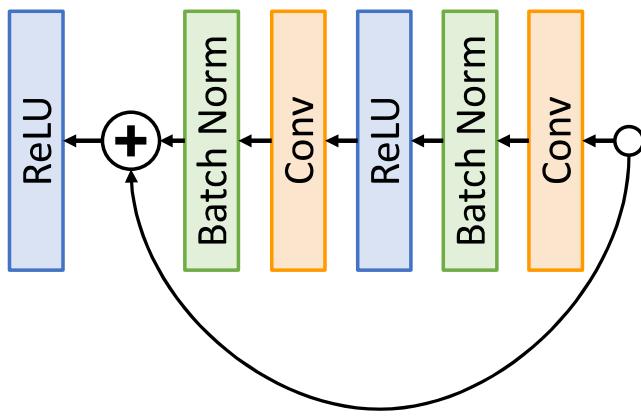
Cannot actually learn
identity function since
outputs are nonnegative!

Note ReLU **inside** residual:

Note ReLU inside residual:
Can learn true identity
function by setting Conv
weights to zero!

Improving Residual Networks: Block Design

Original ResNet block

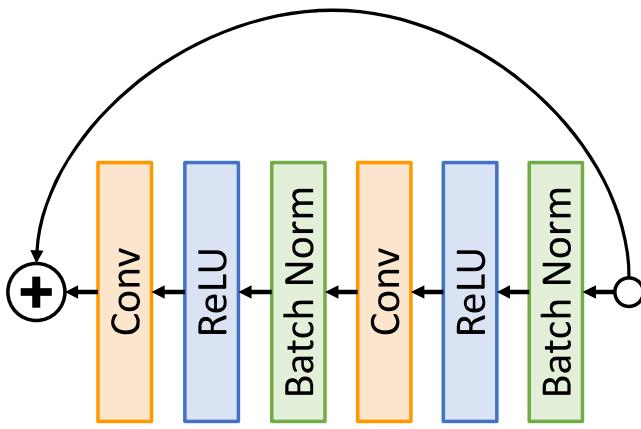


Slight improvement in accuracy
(ImageNet top-1 error)

ResNet-152: 21.3 vs **21.1**
ResNet-200: 21.8 vs **20.7**

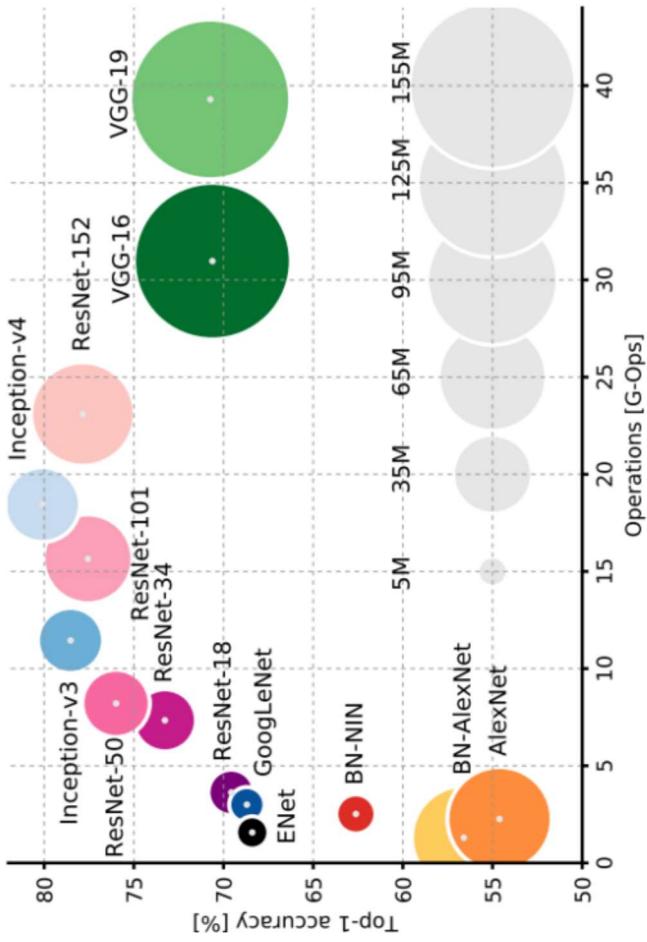
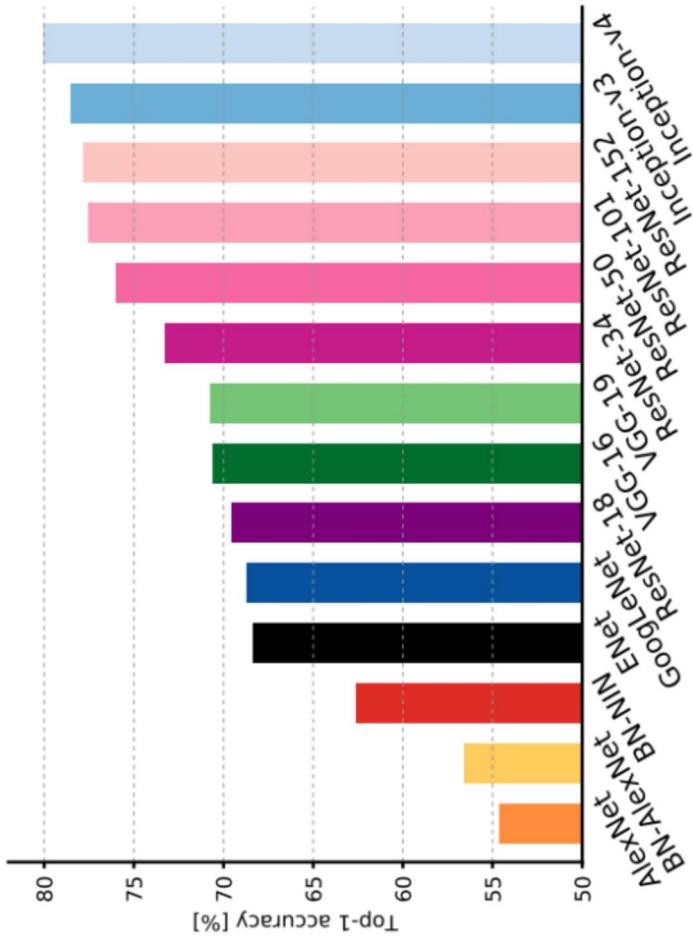
Not actually used that much in
practice

“Pre-Activation” ResNet Block



He et al., “Identity mappings in deep residual networks”, ECCV 2016

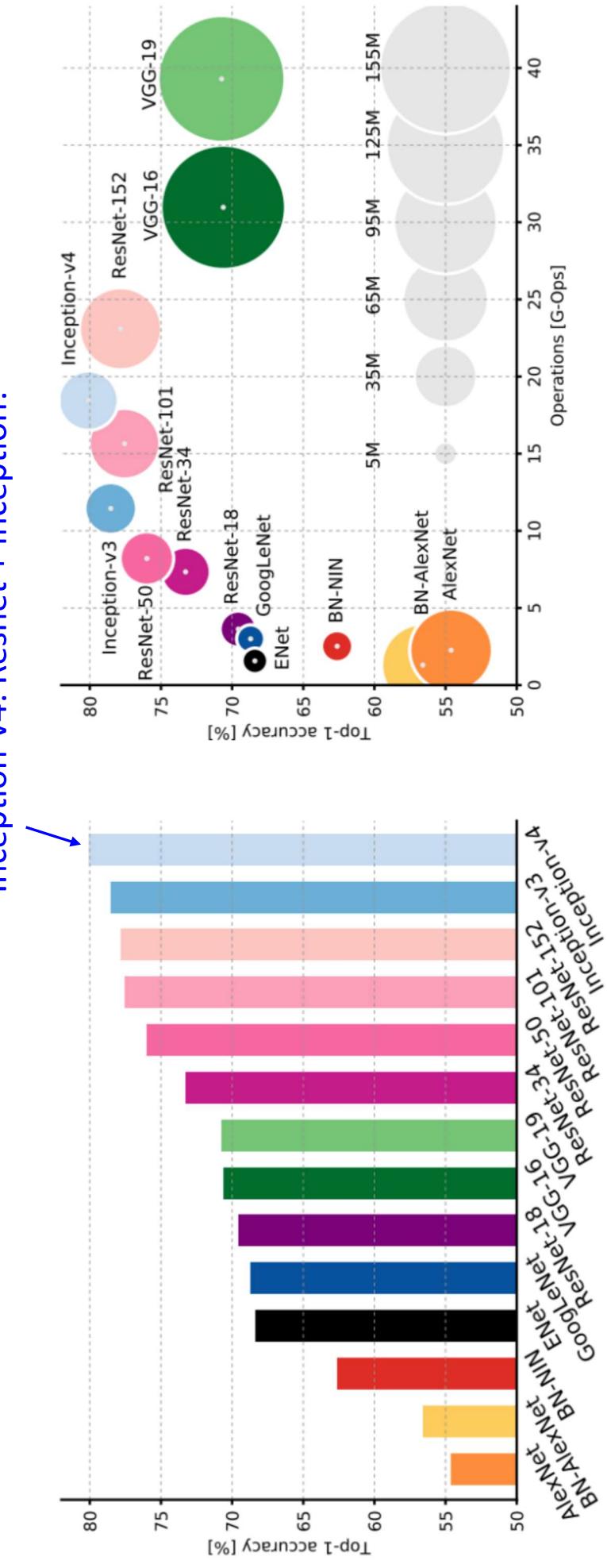
Comparing Complexity



Canziani et al, "An analysis of deep neural network models for practical applications", 2017

Comparing Complexity

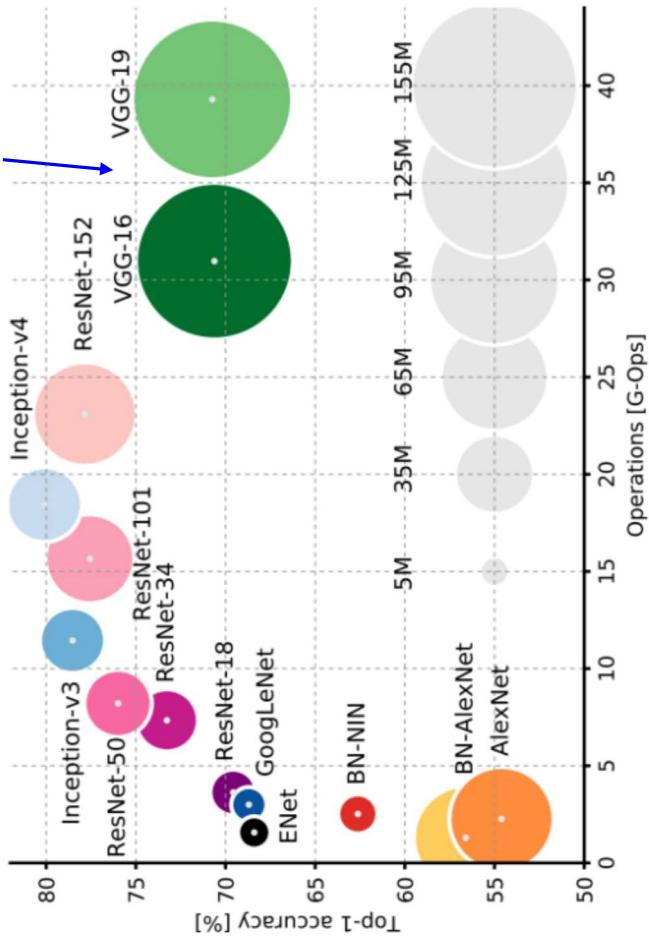
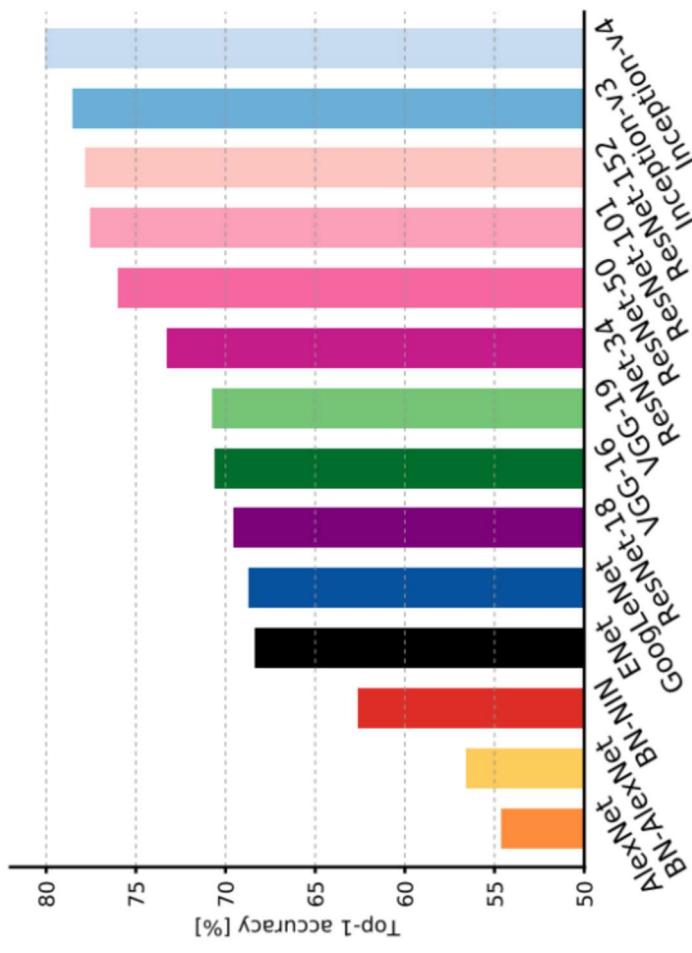
Inception-v4: Resnet + Inception!



Canziani et al, "An analysis of deep neural network models for practical applications", 2017

Comparing Complexity

VGG: Highest
memory, most
operations



Canziani et al, "An analysis of deep neural network models for practical applications", 2017

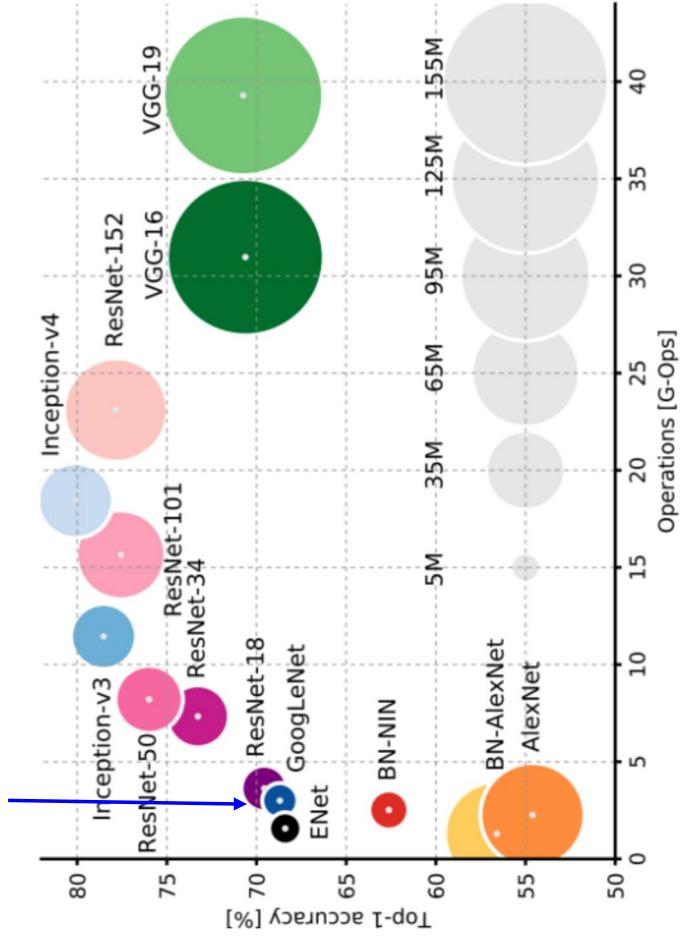
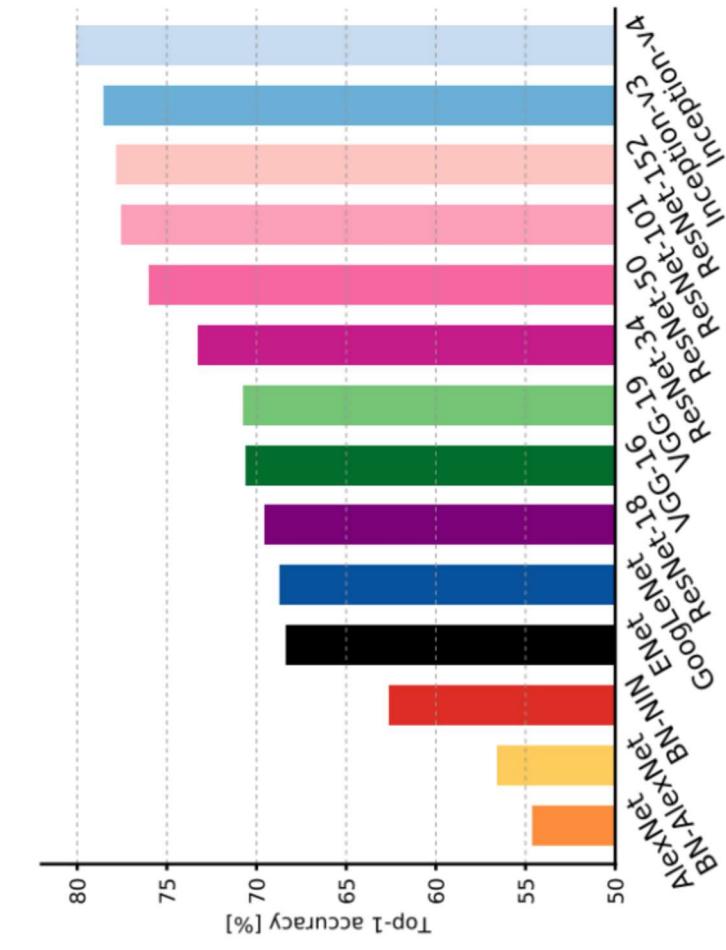
Justin Johnson

Lecture 8 - 84

September 30, 2019

Comparing Complexity

GoogLeNet:
Very efficient!



Canziani et al, "An analysis of deep neural network models for practical applications", 2017

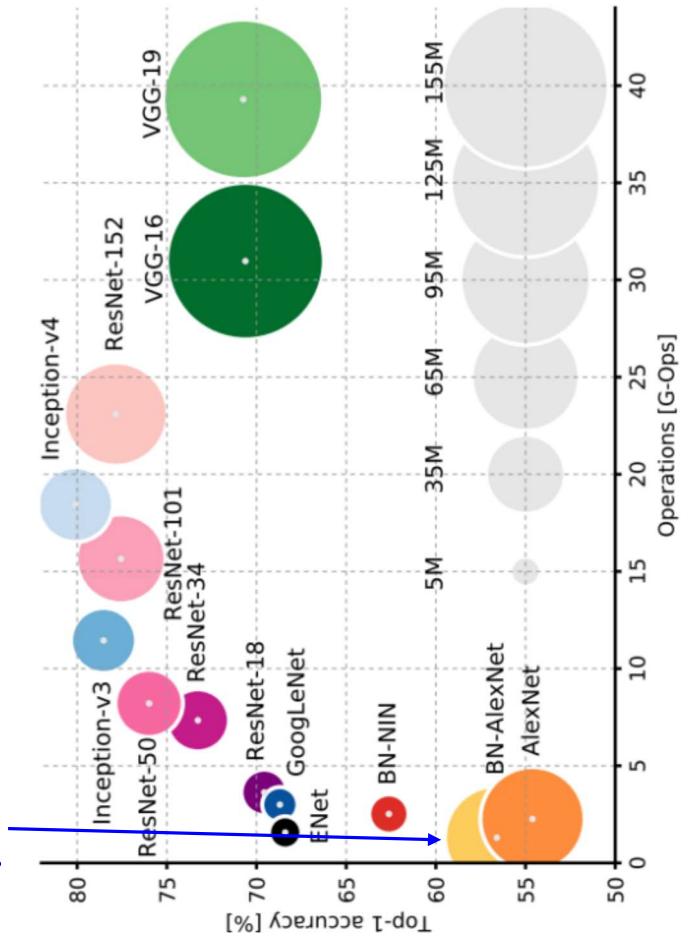
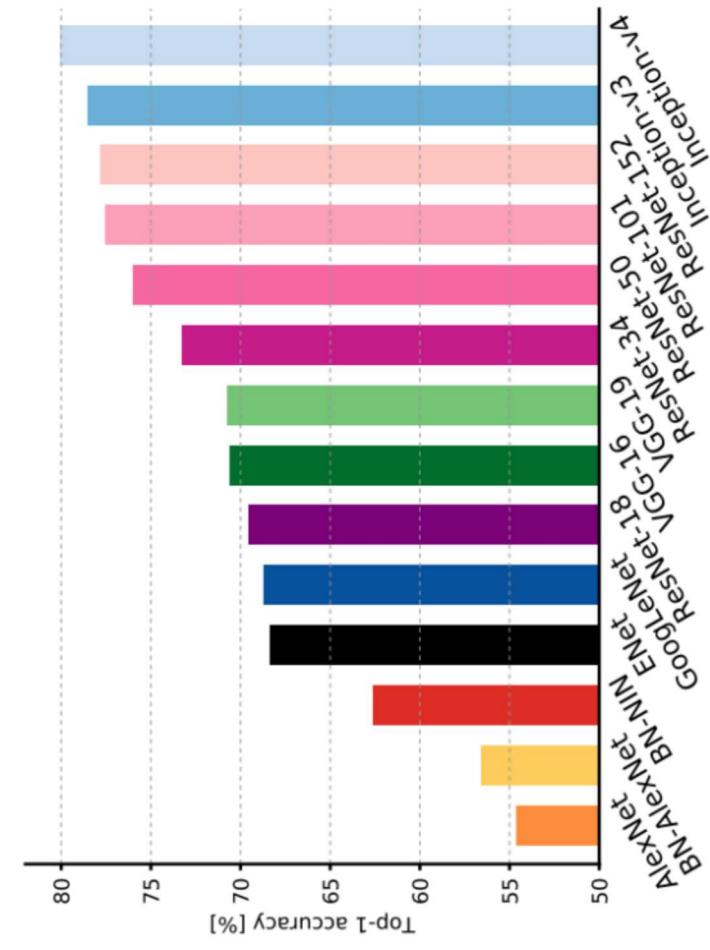
Justin Johnson

Lecture 8 - 85

September 30, 2019

Comparing Complexity

AlexNet: Low
compute, lots
of parameters



Canziani et al, "An analysis of deep neural network models for practical applications", 2017

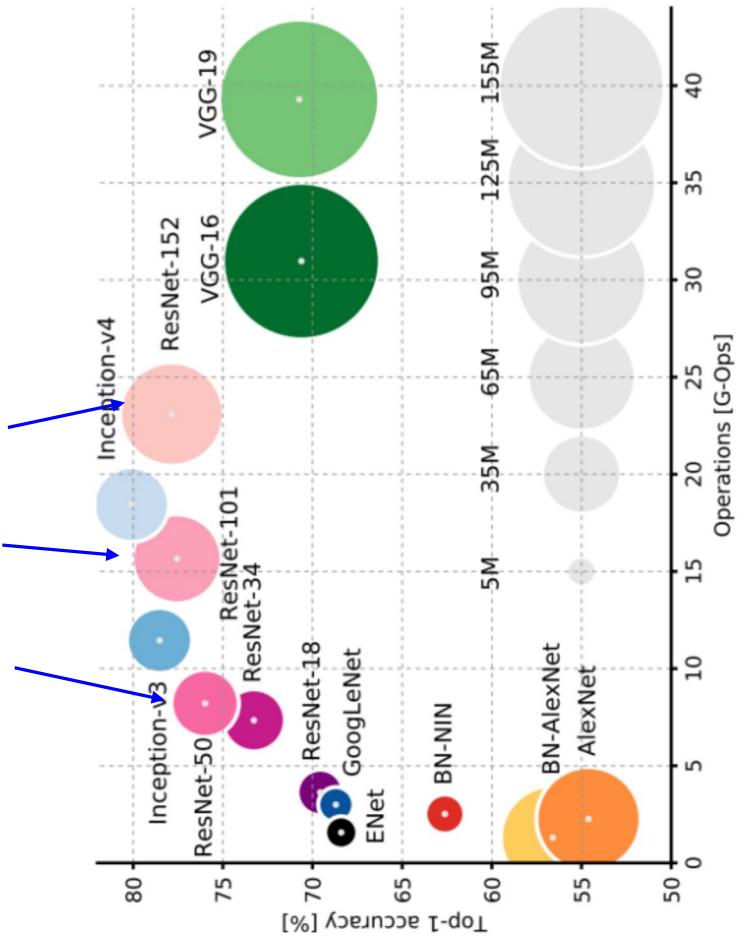
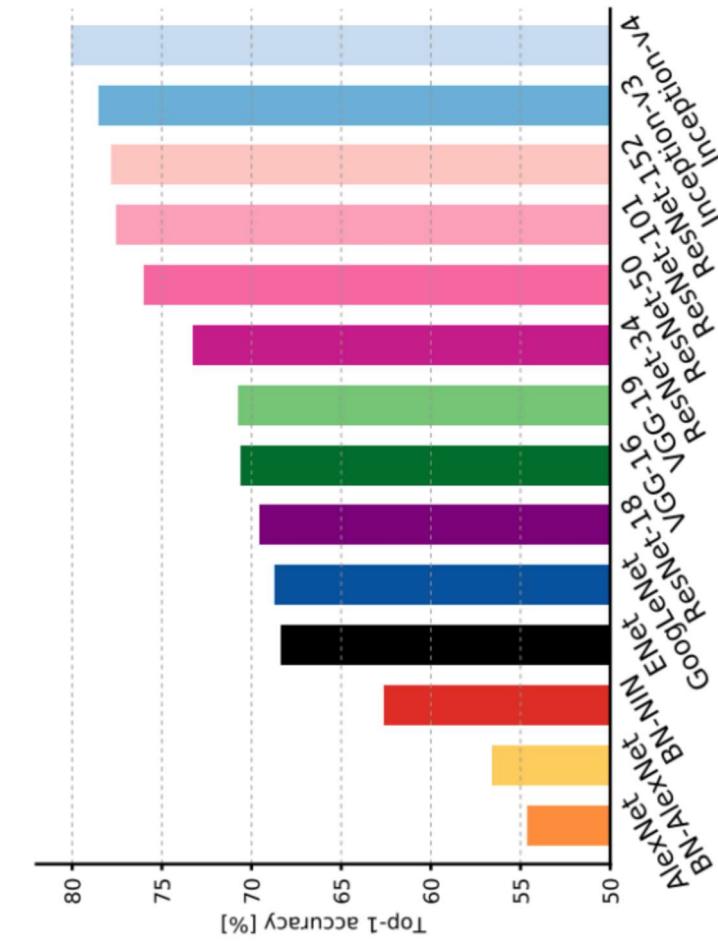
Justin Johnson

Lecture 8 - 86

September 30, 2019

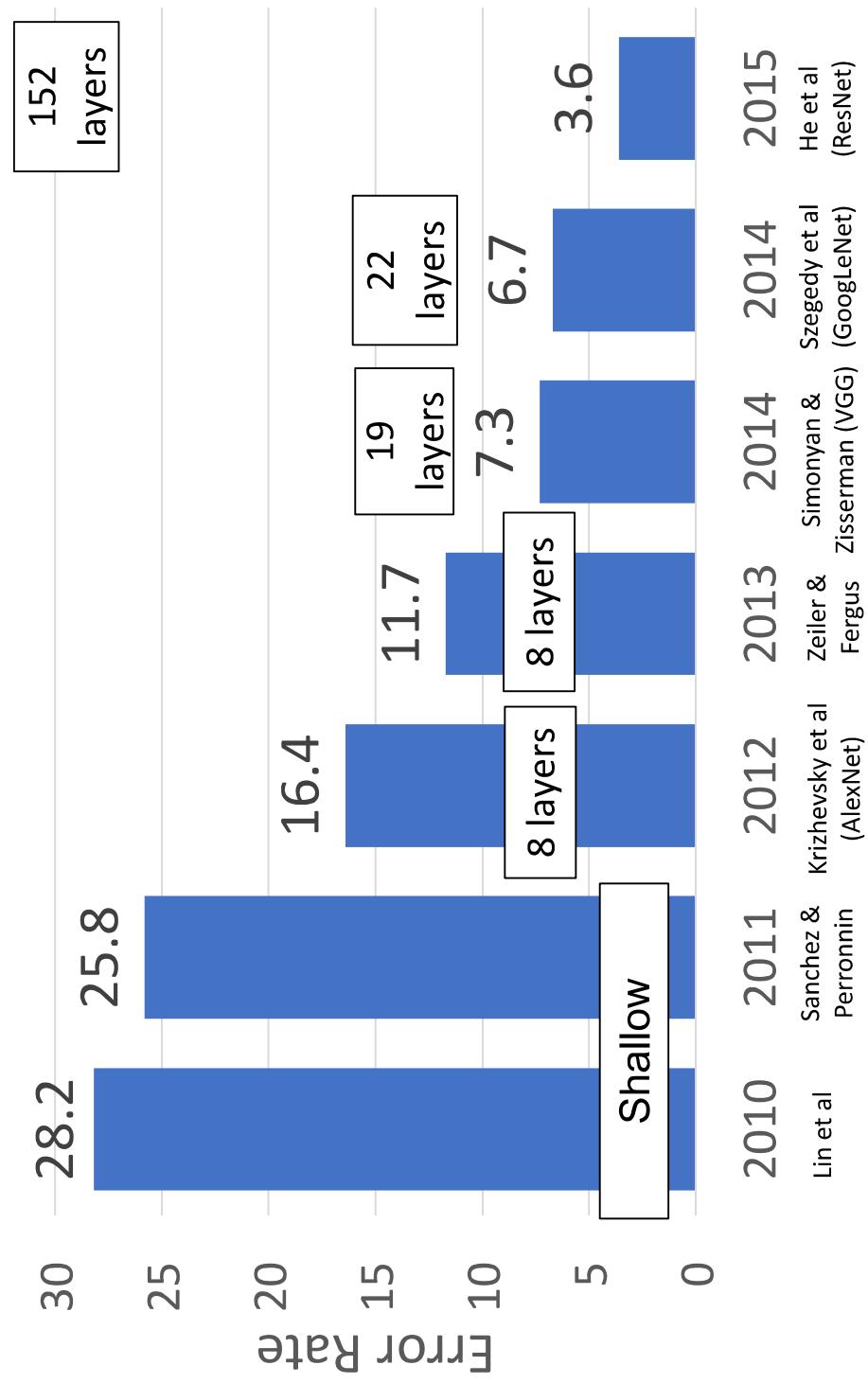
Comparing Complexity

ResNet: Simple design,
moderate efficiency,
high accuracy



Canziani et al, "An analysis of deep neural network models for practical applications", 2017

ImageNet Classification Challenge

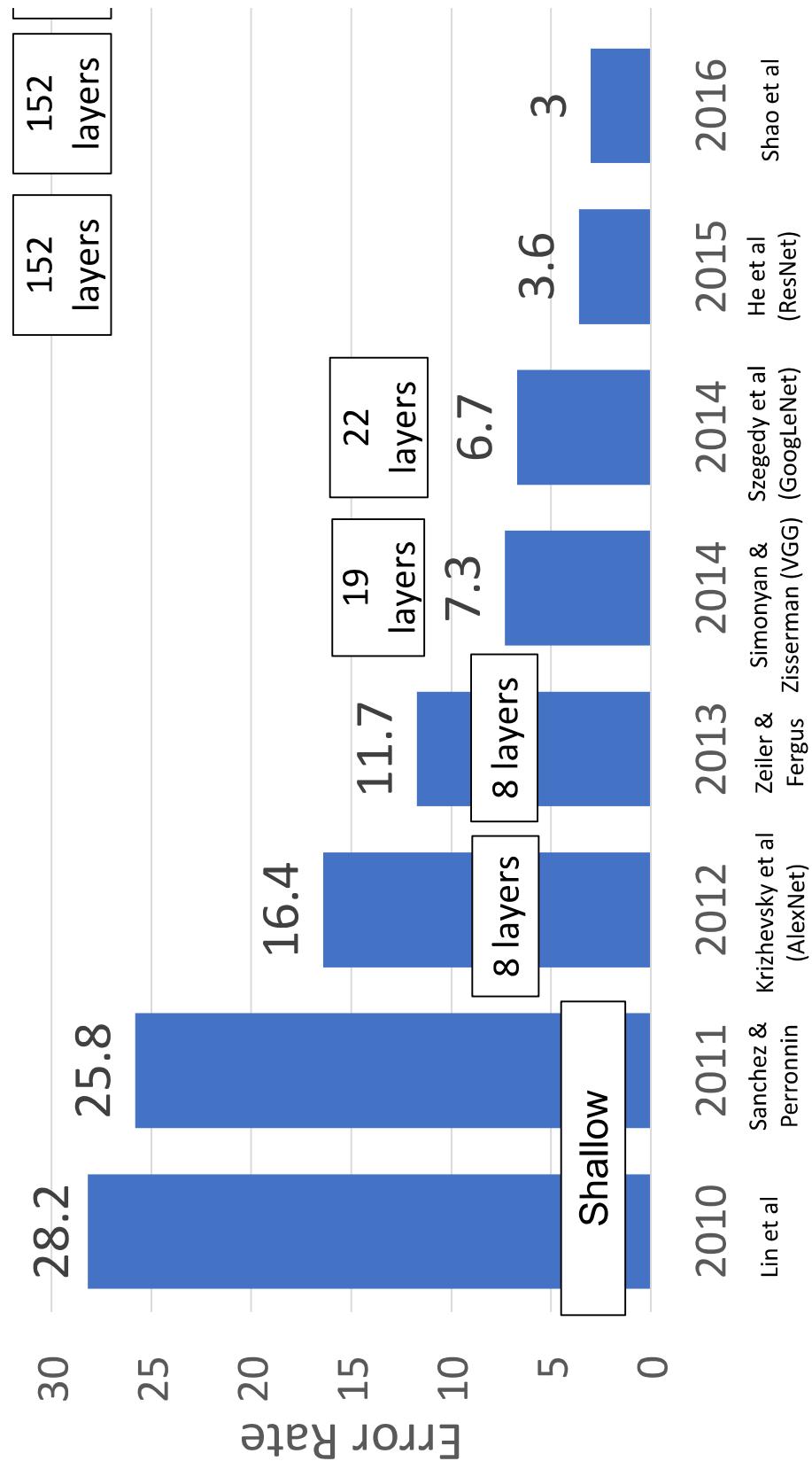


Justin Johnson

Lecture 8 - 88

September 30, 2019

ImageNet Classification Challenge



Justin Johnson

Lecture 8 - 89

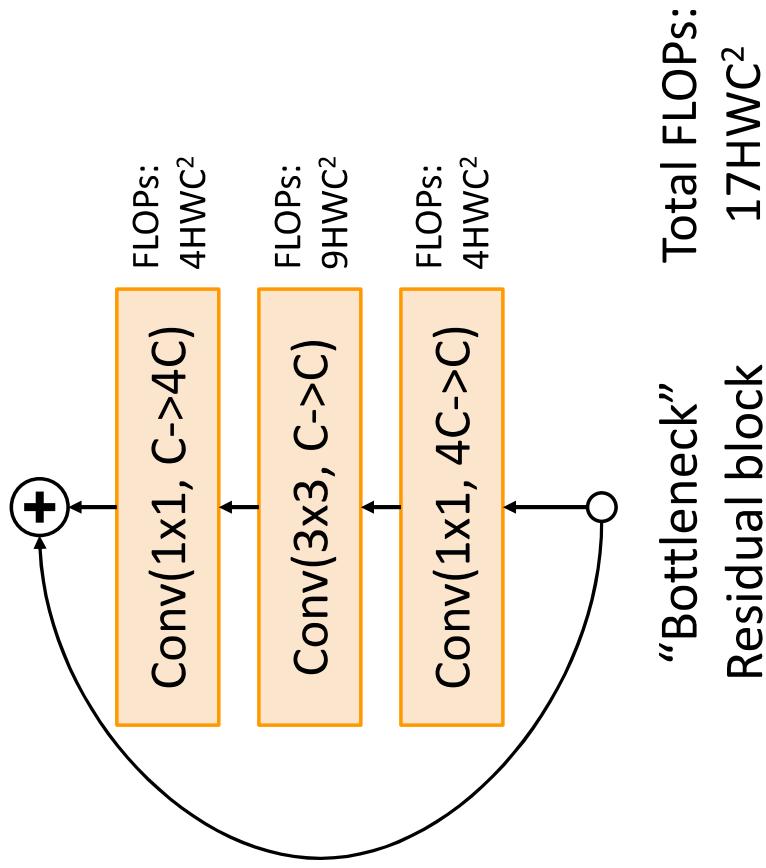
September 30, 2019

ImageNet 2016 winner: Model Ensembles

Multi-scale ensemble of Inception, Inception-Resnet,
Resnet, Wide Resnet models

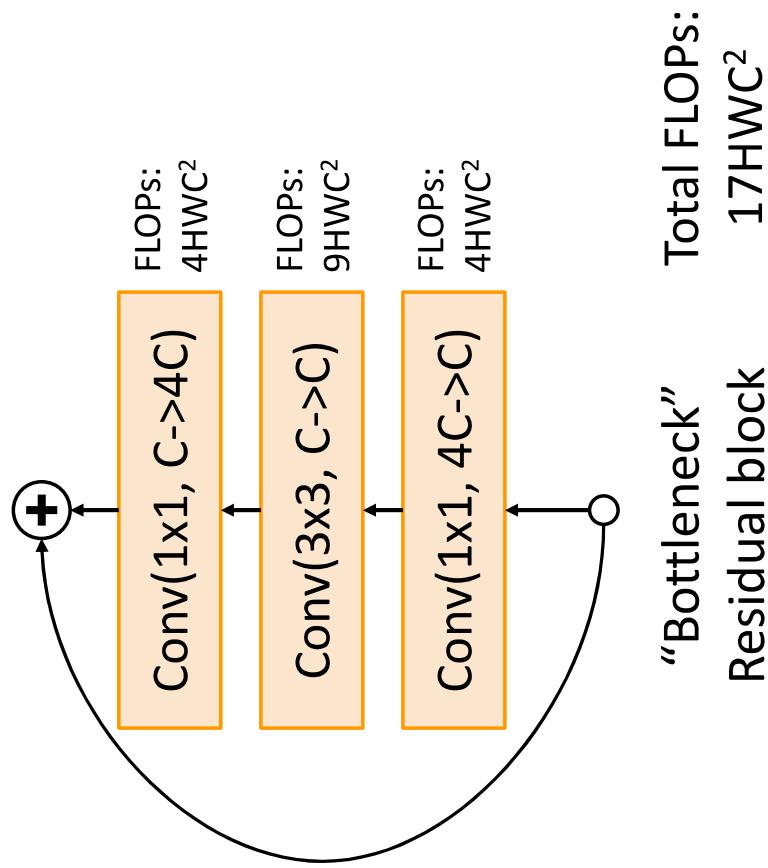
	Inception-v3	Inception-v4	Inception-Resnet-v2	Resnet-200	Wrn-68-3	Fusion (Val.)	Fusion (Test)
Err. (%)	4.20	4.01	3.52	4.26	4.65	2.92 (-0.6)	2.99

Improving ResNets



Improving ResNets: ResNeXt

G parallel pathways

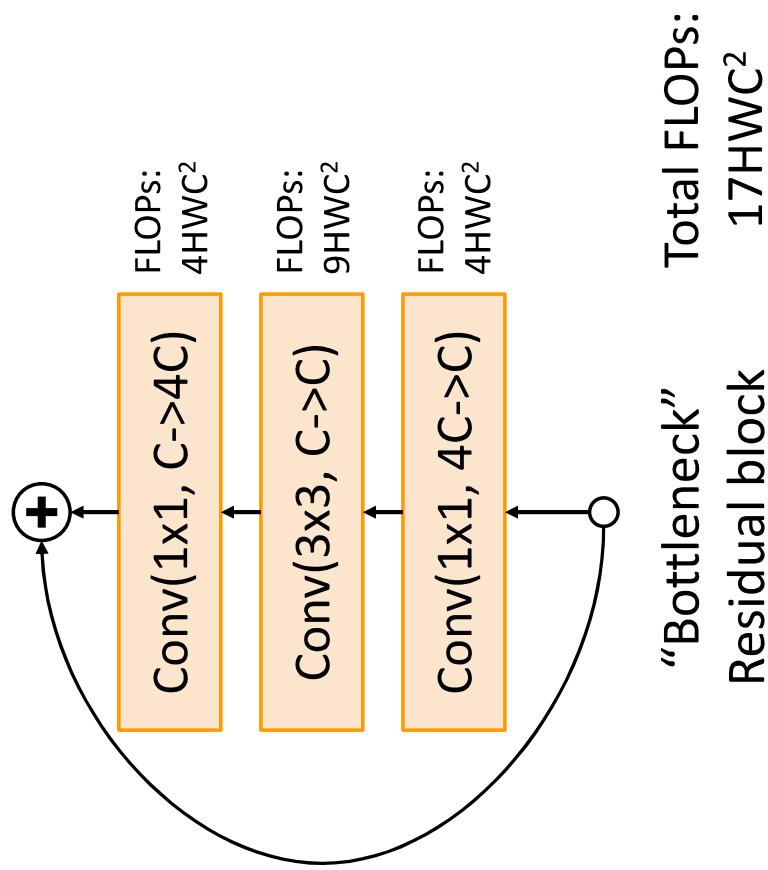


“Bottleneck”
Residual block
Total FLOPs:
 $17\text{HW}C^2$

Xie et al, “Aggregated residual transformations for deep neural networks”, CVPR 2017

Improving ResNets: ResNeXt

G parallel pathways



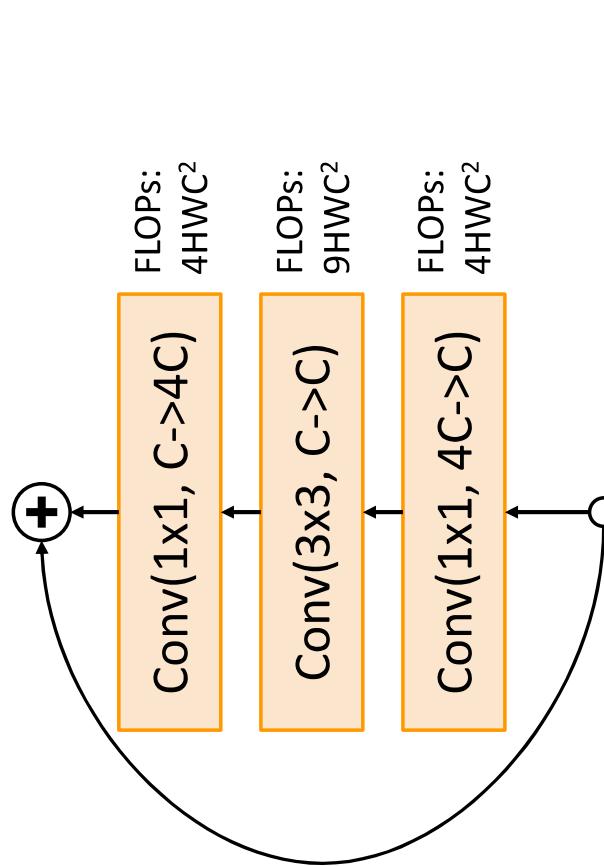
“Bottleneck”
Residual block

Total FLOPs:
 $17HWc^2$

Total FLOPs:
 $(8Cc + 9c^2)*HWG$

Improving ResNets: ResNeXt

$\underbrace{\quad\quad\quad}_{G \text{ parallel pathways}}$



FLOPs:
4HW C^2

FLOPs:
9HW C^2

FLOPs:
4HW C^2

Total FLOPs:
17HW C^2

“Bottleneck”
Residual block

Total FLOPs:
 $(8CC + 9C^2) * HWG$

Xie et al, “Aggregated residual transformations for deep neural networks”, CVPR 2017

Example: C=64, G=4, c=24; C=64, G=32, c=4

Justin Johnson

Lecture 8 - 94

September 30, 2019

Grouped Convolution

Convolution with groups=1:
Normal convolution

Input: $C_{in} \times H \times W$

Weight: $C_{out} \times C_{in} \times K \times K$

Output: $C_{out} \times H' \times W'$

FLOPs: $C_{out} C_{in} K^2 H W$

All convolutional kernels touch
all C_{in} channels of the input

Grouped Convolution

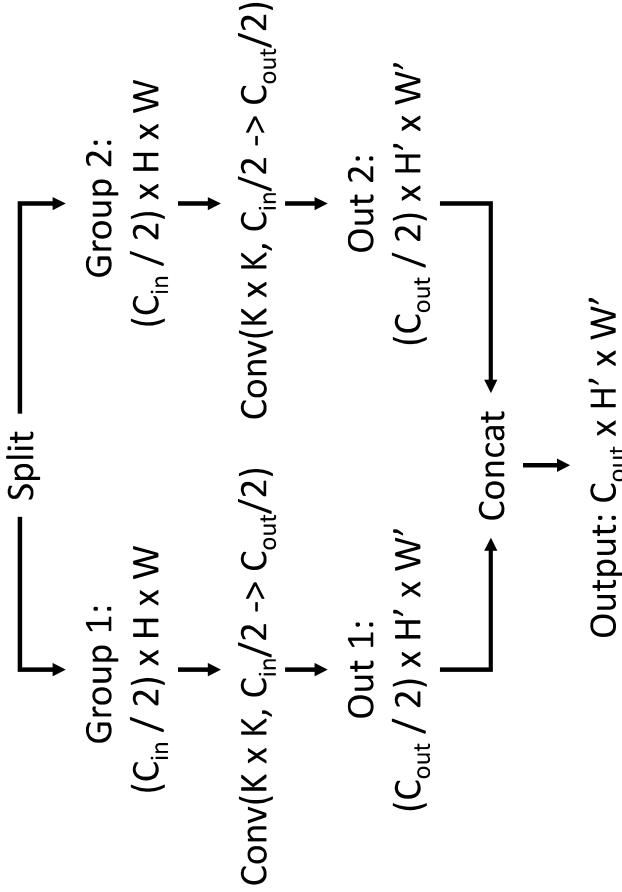
Convolution with groups=1:
Normal convolution

Input: $C_{in} \times H \times W$
Weight: $C_{out} \times C_{in} \times K \times K$
Output: $C_{out} \times H' \times W'$
FLOPs: $C_{out} C_{in} K^2 H W$

All convolutional kernels touch
all C_{in} channels of the input

Convolution with groups=2:
Two parallel convolution layers that
work on half the channels

Input: $C_{in} \times H \times W$



Grouped Convolution

Convolution with groups=1:
Normal convolution

Input: $C_{in} \times H \times W$
Weight: $C_{out} \times C_{in} \times K \times K$
Output: $C_{out} \times H' \times W'$
FLOPs: $C_{out} C_{in} K^2 HW$

All convolutional kernels touch
all C_{in} channels of the input

Convolution with groups=G:
 G parallel conv layers; each “sees”
 C_{in}/G input channels and produces
 C_{out}/G output channels

Input: $C_{in} \times H \times W$
Split to $G \times [(C_{in}/G) \times H \times W]$
Weight: $G \times (C_{out}/G) \times (C_{in}/G) \times K \times K$
 G parallel convolutions
Output: $G \times [(C_{out}/G) \times H' \times W']$
Concat to $C_{out} \times H' \times W'$
FLOPs: $C_{out} C_{in} K^2 HW/G$

Grouped Convolution

Convolution with groups=1:
Normal convolution

Input: $C_{in} \times H \times W$
Weight: $C_{out} \times C_{in} \times K \times K$
Output: $C_{out} \times H' \times W'$
FLOPs: $C_{out} C_{in} K^2 HW$

All convolutional kernels touch
all C_{in} channels of the input

Depthwise Convolution

Special case: $G=C_{in}$, $C_{out} = nC_{in}$
Each input channel is convolved
with n different $K \times K$ filters to
produce n output channels

Convolution with groups=G:
 G parallel conv layers; each “sees”
 C_{in}/G input channels and produces
 C_{out}/G output channels

Input: $C_{in} \times H \times W$
Split to $G \times [(C_{in}/G) \times H \times W]$
Weight: $G \times (C_{out}/G) \times (C_{in}/G) \times K \times K$
 G parallel convolutions
Output: $G \times [(C_{out}/G) \times H' \times W']$
Concat to $C_{out} \times H' \times W'$
FLOPs: $C_{out} C_{in} K^2 HW/G$

Grouped Convolution in PyTorch

PyTorch convolution gives an option for groups!

Conv2d

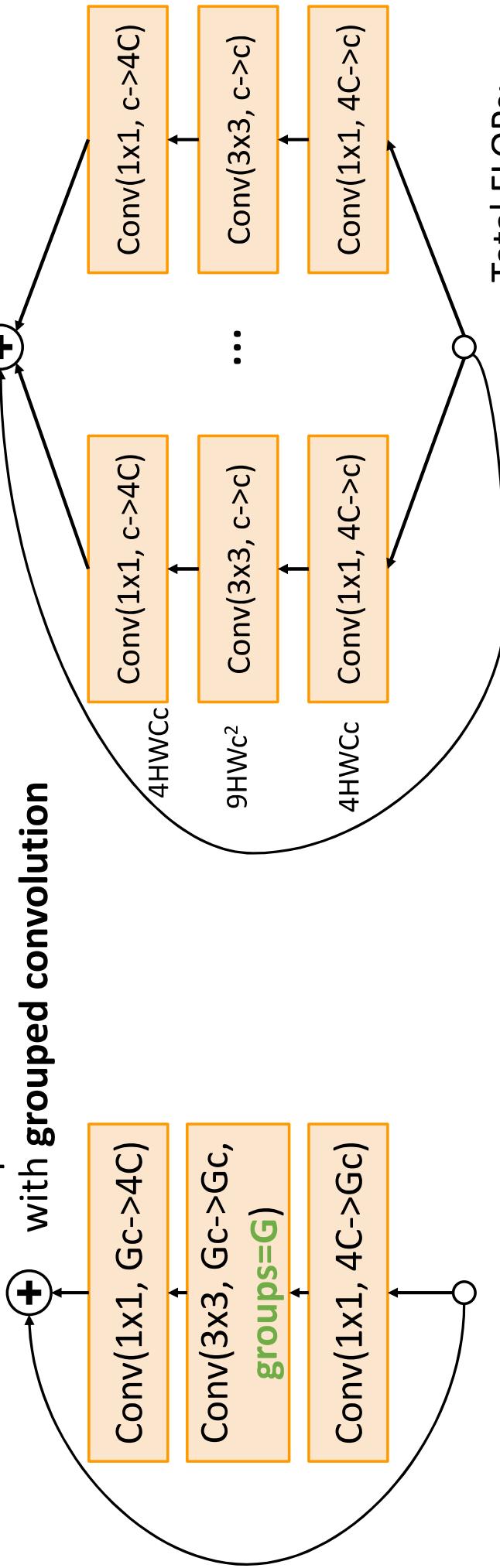
```
CLASS torch.nn.Conv2d(in_channels, out_channels, kernel_size,
stride=1, padding=0, dilation=1, groups=1, bias=True,
padding_mode='zeros')
```

[SOURCE]

Improving ResNets: ResNeXt

Equivalent formulation
with **grouped convolution**

G parallel pathways



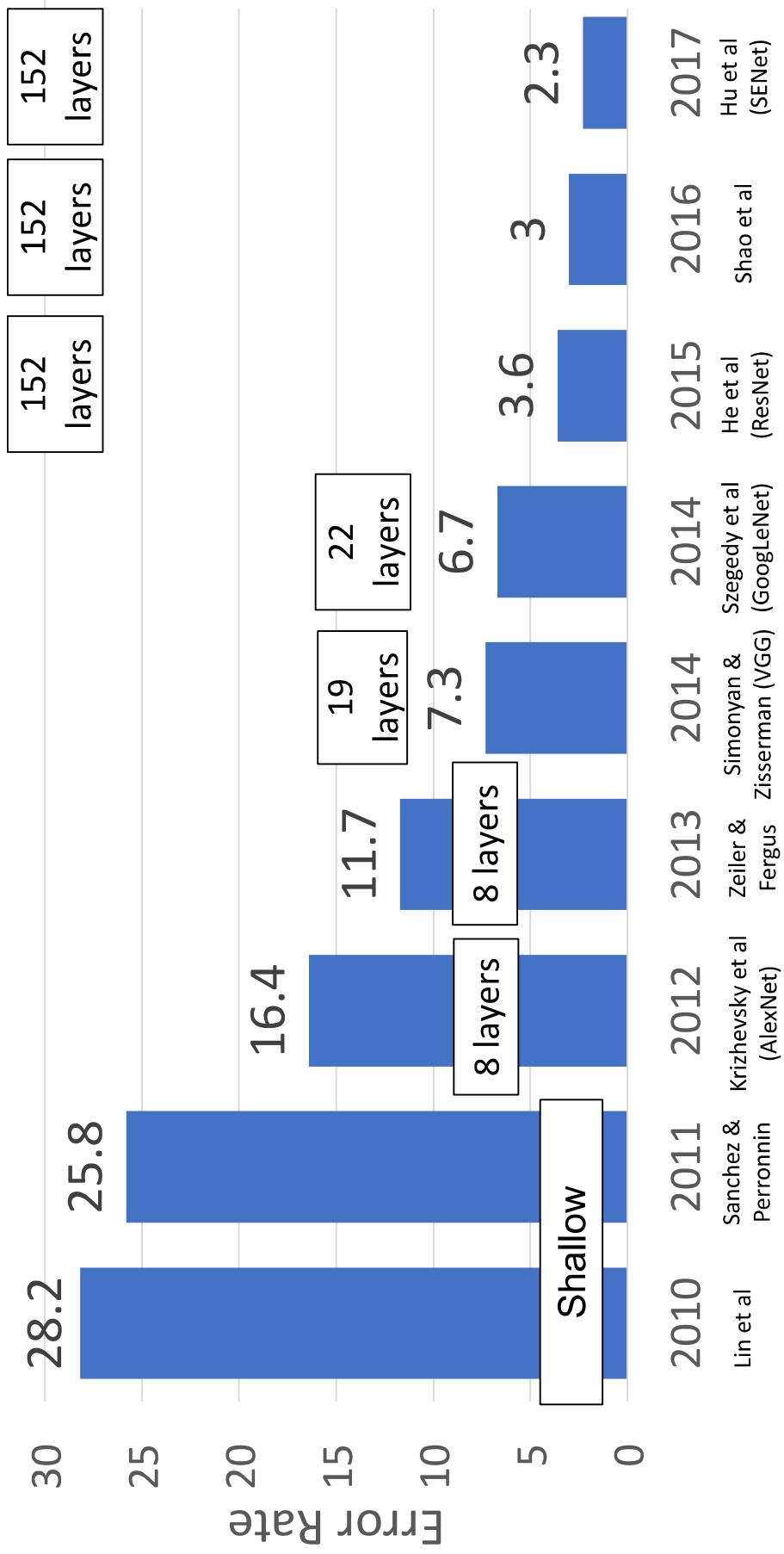
Xie et al., "Aggregated residual transformations for deep neural networks", CVPR 2017

ResNeXt: Maintain computation by adding groups!

Model	Groups	Group width	Top-1 Error
ResNet-50	1	64	23.9
ResNeXt-50	2	40	23
ResNeXt-50	4	24	22.6
ResNeXt-50	8	14	22.3
ResNeXt-50	32	4	22.2
ResNet-101		1	64
ResNeXt-101		2	40
ResNeXt-101		4	24
ResNeXt-101		8	14
ResNeXt-101		32	4

Adding groups improves performance **with same computational complexity!**

ImageNet Classification Challenge

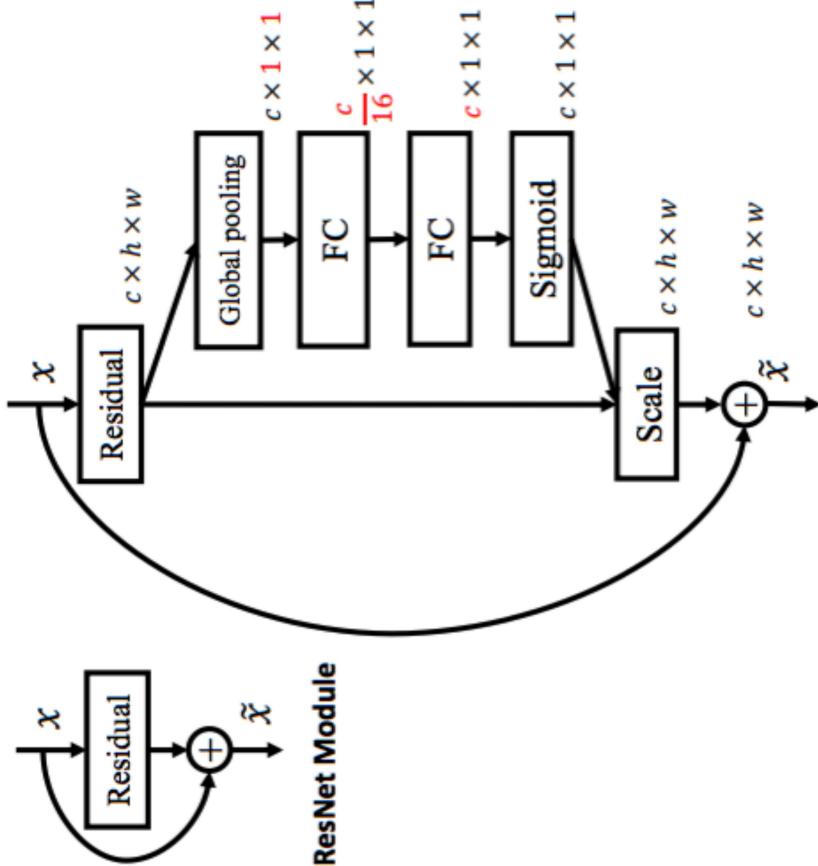


Justin Johnson

Lecture 8 - 102

September 30, 2019

Squeeze-and-Excitation Networks

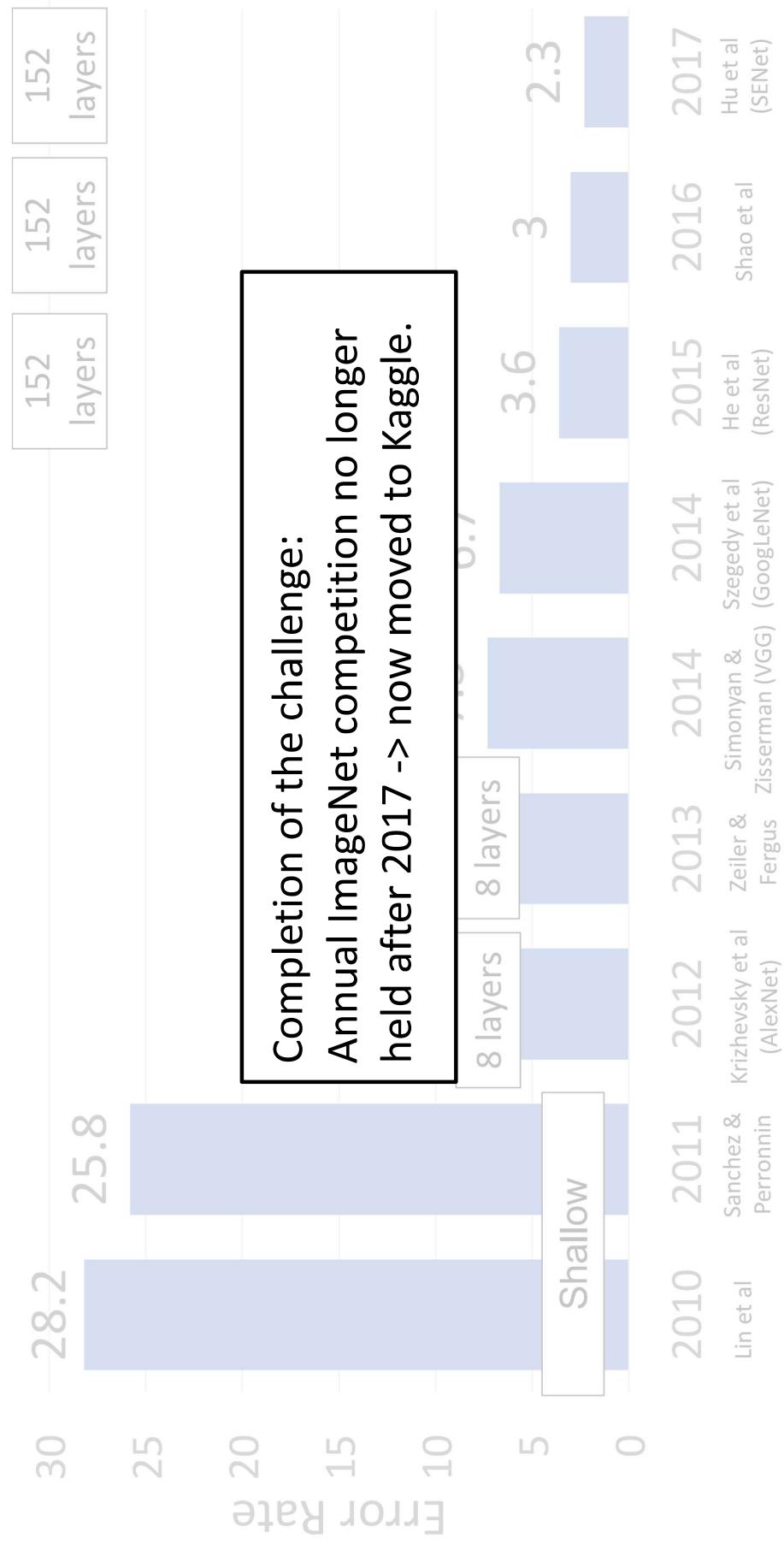


Adds a "Squeeze-and-excite" branch to each residual block that performs global pooling, full-connected layers, and multiplies back onto feature map

Adds **global context** to each residual block!

Won ILSVRC 2017 with ResNeXt-152-SE

ImageNet Classification Challenge

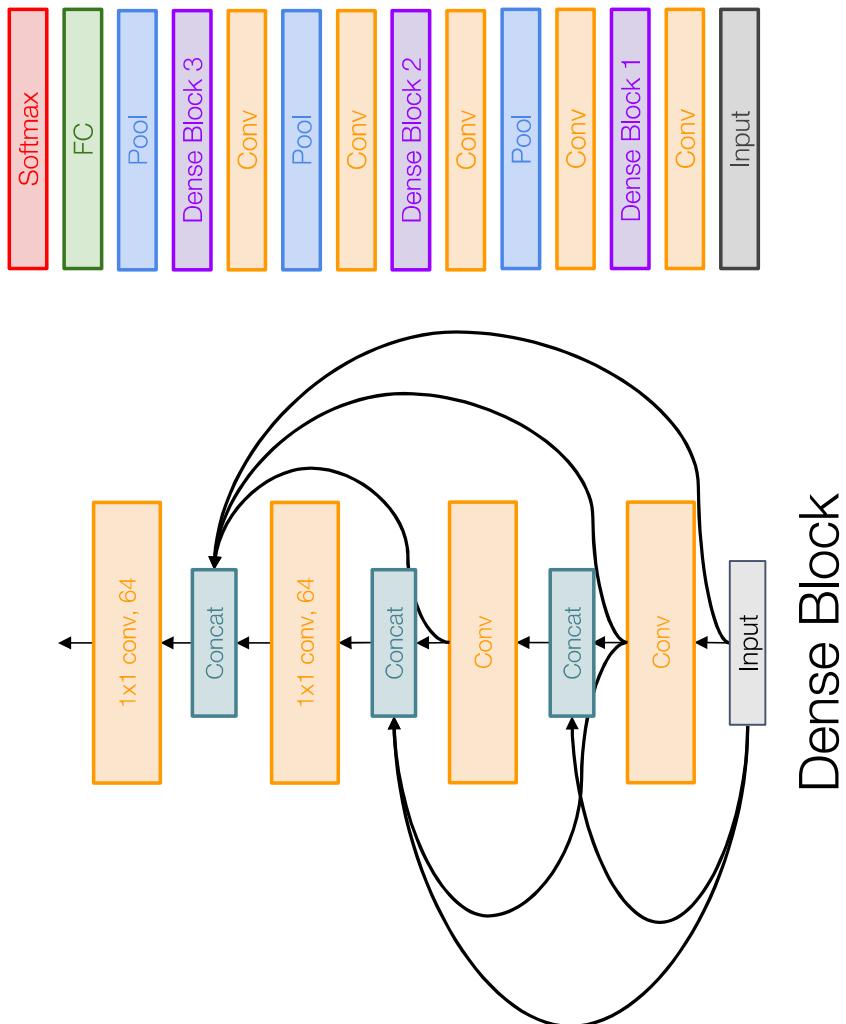


Justin Johnson

Lecture 8 - 104

September 30, 2019

Densely Connected Neural Networks



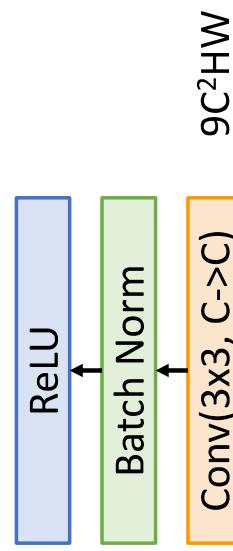
Dense blocks where each layer is connected to every other layer in feedforward fashion

Alleviates vanishing gradient,
strengthens feature propagation,
encourages feature reuse

MobileNets: Tiny Networks (For Mobile Devices)

Standard Convolution Block

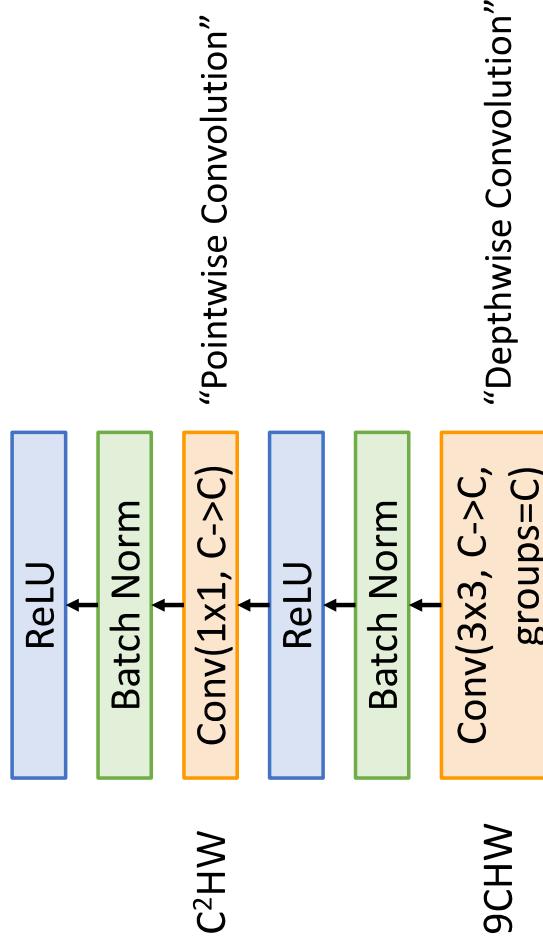
Total cost: $9C^2HW$



$$\begin{aligned} \text{Speedup} &= 9C^2/(9C+C^2) \\ &= 9C/(9+C) \\ &\Rightarrow 9 \text{ (as } C \rightarrow \infty \text{)} \end{aligned}$$

Depthwise Separable Convolution

Total cost: $(9C + C^2)HW$



Howard et al, "MobileNets: Efficient Convolutional Neural Networks for Mobile Vision Applications", 2017

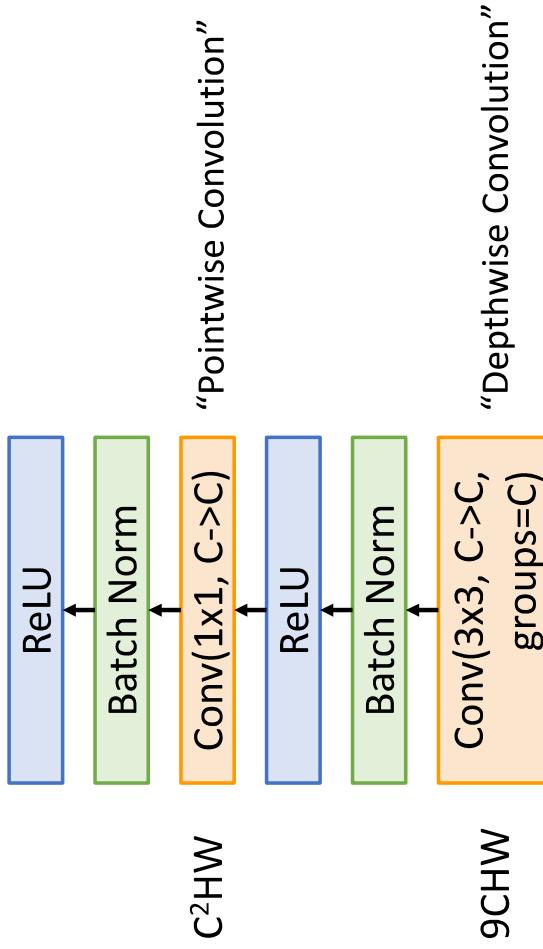
MobileNets: Tiny Networks (For Mobile Devices)

Depthwise Separable Convolution

Total cost: $(9C + C^2)HW$

Also related:

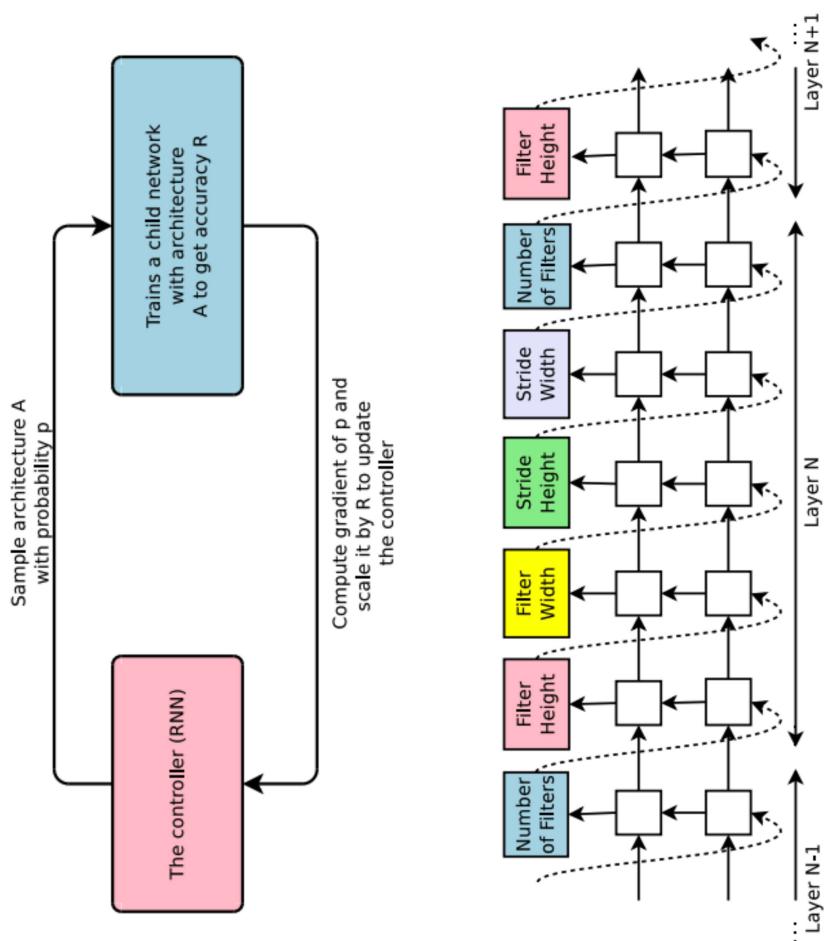
ShuffleNet: Zhang et al, CVPR 2018
MobileNetV2: Sandler et al, CVPR 2018
ShuffleNetV2: Ma et al, ECCV 2018



Neural Architecture Search

Designing neural network architectures is hard – let's automate it!

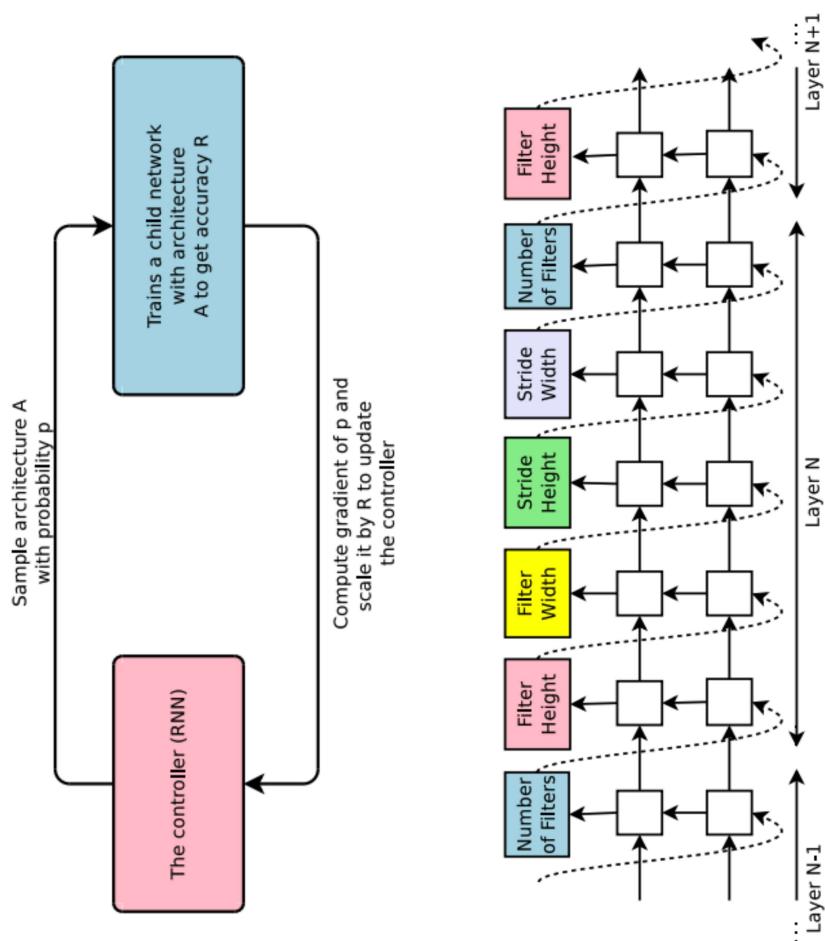
- One network (**controller**) outputs network architectures
- Sample **child networks** from controller and train them
- After training a batch of child networks, make a gradient step on controller network (Using **policy gradient**)
- Over time, controller learns to output good architectures!



Neural Architecture Search

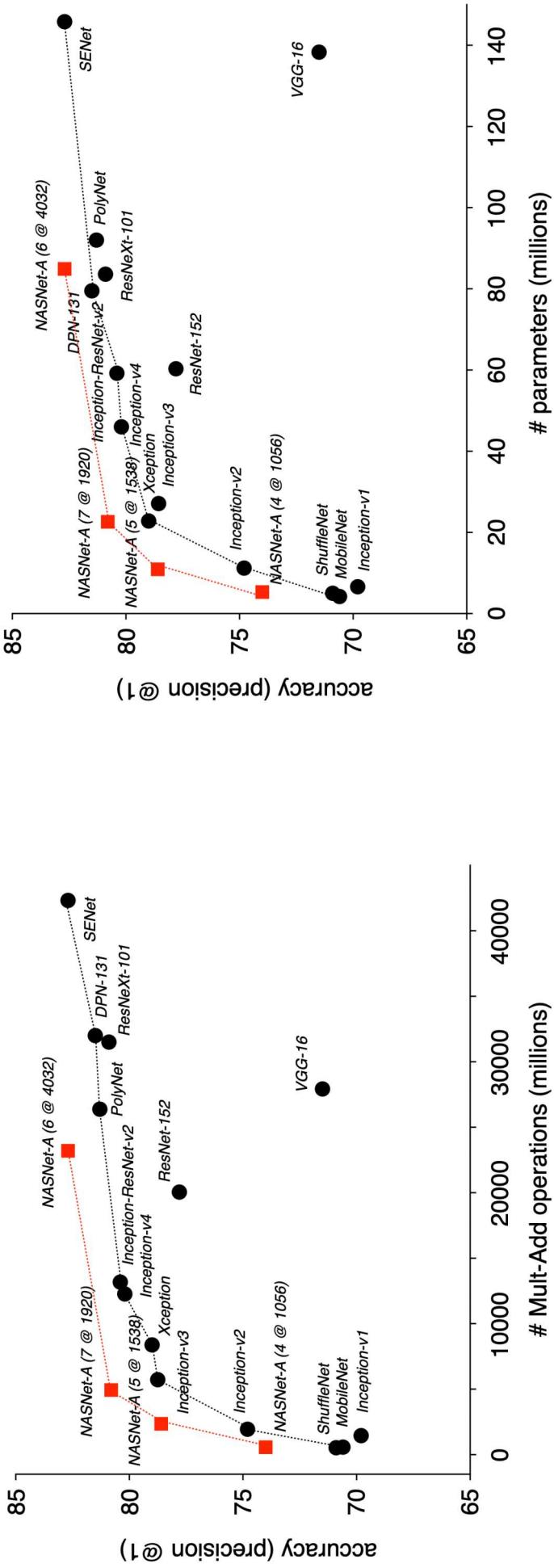
Designing neural network architectures is hard – let's automate it!

- One network (**controller**) outputs network architectures
- Sample **child networks** from controller and train them
- After training a batch of child networks, make a gradient step on controller network (Using **policy gradient**)
- Over time, controller learns to output good architectures!
- **VERY EXPENSIVE!! Each gradient step on controller requires training a batch of child models!**
- **Original paper trained on 800 GPUs for 28 days!**
- **Followup work has focused on efficient search**



Neural Architecture Search

Neural architecture search can be used to find efficient CNN architectures!



CNN Architectures Summary

Early work (AlexNet -> ZFNet -> VGG) shows that **bigger networks work better**

GoogLeNet one of the first to focus on **efficiency** (aggressive stem, 1x1 bottleneck convolutions, global avg pool instead of FC layers)

ResNet showed us how to train extremely deep networks – limited only by GPU memory! Started to show diminishing returns as networks got bigger

After ResNet: **Efficient networks** became central: how can we improve the accuracy without increasing the complexity?

Lots of **tiny networks** aimed at mobile devices: MobileNet, ShuffleNet, etc

Neural Architecture Search promises to automate architecture design

Which Architecture should I use?

Don't be a hero. For most problems you should use an off-the-shelf architecture; don't try to design your own!

If you just care about accuracy, **ResNet-50** or **ResNet-101** are great choices

If you want an efficient network (real-time, run on mobile, etc) try
MobileNets and **ShuffleNets**

Next Time:
Deep Learning Hardware and Software