

Deeper Deep Networks

presented by:
Spencer Cappallo

Overview

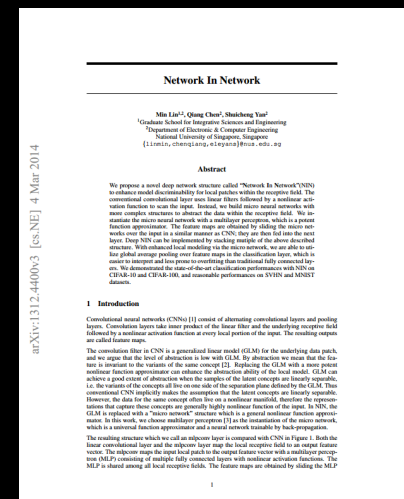
- Three recent papers discussing “deeper” deep networks
- All achieve state-of-the-art results
- High-level overview of new ideas in these networks

Network in Network

Lin, Min, Qiang Chen, and Shuicheng Yan. "Network In Network." *arXiv preprint arXiv:1312.4400* (2013).

Basic Idea:

Uses small “micro networks” as a function approximator to replace conventional convolution operation



Network in Network: Intuition

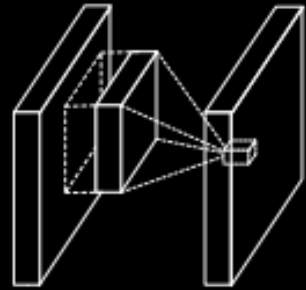
- Activations correspond to latent concepts
- Convolutional filters act as linear binary classifiers for these latent concepts in local patches
- These filters work well when the local latent concepts are linearly separable, but instead often very nonlinear
- NIN instead opts for nonlinear network structure to replace convolutional filters

Network in Network: mlpconv layers

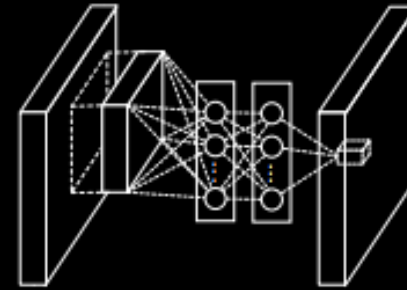
Use multilayer perceptrons as universal function approximator in place of standard convolutions

- Easily incorporated into backpropagated network
- Uses ReLUs
- Shared hidden units

Network in Network: mlpconv layers



(a) Linear convolution layer



(b) Mlpconv layer

Hidden layers are shared between output feature maps

- Cross channel information

Network in Network: Structure

Replace fully connected layers with global average pooling

- More interpretable: direct connection between categories and feature maps
 - Enforces a correspondence between last feature map and category
- Prevents overfitting

Network in Network: Structure

The network used in the paper:

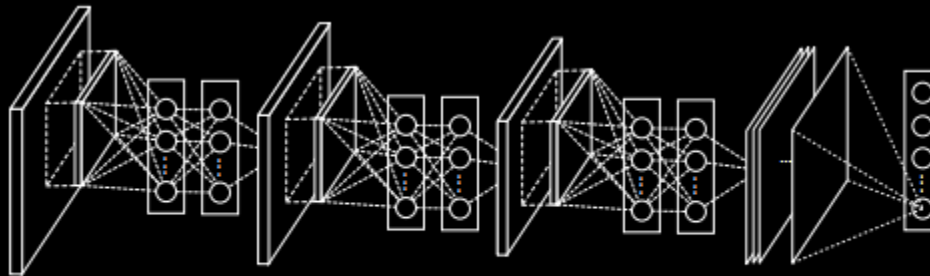


Figure 2: The overall structure of Network In Network. In this paper the NINs include the stacking of three mlpconv layers and one global average pooling layer.

- Three mlpconv layers
- Global average pooling layer
- No FC layers

Network in Network: Results

Comparison of global average pooling to fully connected layers on CIFAR-10:

Method	Testing Error
mlpconv + Fully Connected	11.59%
mlpconv + Fully Connected + Dropout	10.88%
mlpconv + Global Average Pooling	10.41%

Improvement over fully connected layers. This result will also be repeated in the next paper.

Network in Network: Results

Table 1: Test set error rates for CIFAR-10 of various methods.

Method	Test Error
Stochastic Pooling [11]	15.13%
CNN + Spearmint [14]	14.98%
Conv. maxout + Dropout [8]	11.68%
NIN + Dropout	10.41 %
CNN + Spearmint + Data Augmentation [14]	9.50%
Conv. maxout + Dropout + Data Augmentation [8]	9.38%
DropConnect + 12 networks + Data Augmentation [15]	9.32%
NIN + Dropout + Data Augmentation	8.81 %

Table 3: Test set error rates for SVHN of various methods.

Method	Test Error
Stochastic Pooling [11]	2.80%
Rectifier + Dropout [18]	2.78%
Rectifier + Dropout + Synthetic Translation [18]	2.68%
Conv. maxout + Dropout [8]	2.47%
NIN + Dropout	2.35%
Multi-digit Number Recognition [19]	2.16%
DropConnect [15]	1.94%

Table 2: Test set error rates for CIFAR-100 of various methods.

Method	Test Error
Learned Pooling [16]	43.71%
Stochastic Pooling [11]	42.51%
Conv. maxout + Dropout [8]	38.57%
Tree based priors [17]	36.85%
NIN + Dropout	35.68 %

Table 4: Test set error rates for MNIST of various methods.

Method	Test Error
2-Layer CNN + 2-Layer NN [11]	0.53%
Stochastic Pooling [11]	0.47%
NIN + Dropout	0.47%
Conv. maxout + Dropout [8]	0.45 %

* Note that NIN still needs dropout

Network in Network: Take-aways

- Additional non-linearity may significantly improve discriminative abilities of layers
- Replacing fully connected layers with the global average pooling seems to improve performance
- Impressive performance on datasets tested

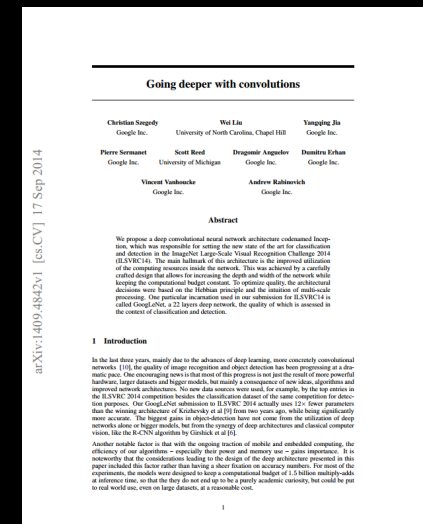
...but how well do these ideas scale up?

GoogLeNet

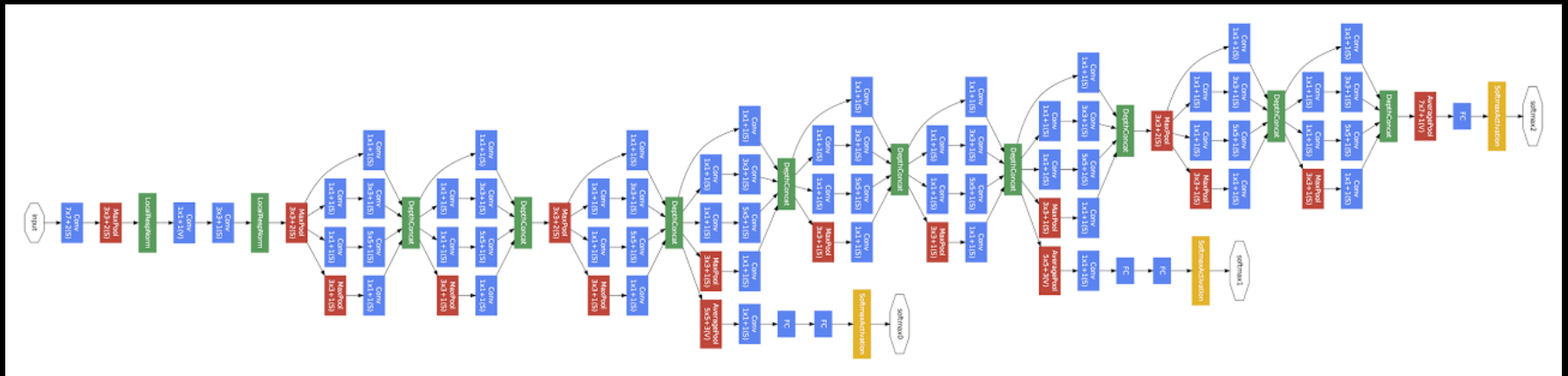
Szegedy, Christian, et al. "Going deeper with convolutions." *arXiv preprint arXiv:1409.4842* (2014).

Submission to ILSVRC2014 challenge

- 1st place Classification
- 1st place Object Detection with additional Training Data



GoogLeNet



Convolution/FC

Max Pooling Softmax

Concatenation

Looks like a big, ugly mess.

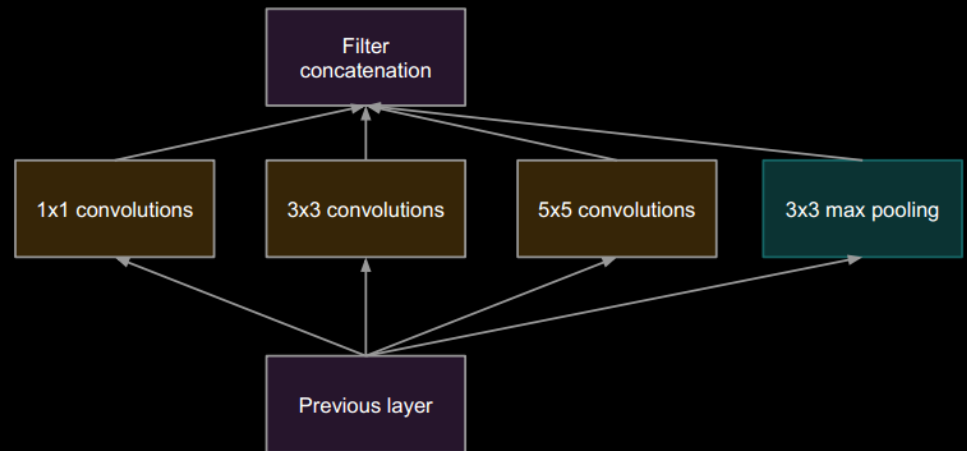
Fortunately, if we break it down a bit it's not *too* bad.

GoogLeNet: Inception Module

Idea 1:

Use 1x1, 3x3, and 5x5 convolutions in parallel to capture a variety of structures

Also add a parallel max pooling path

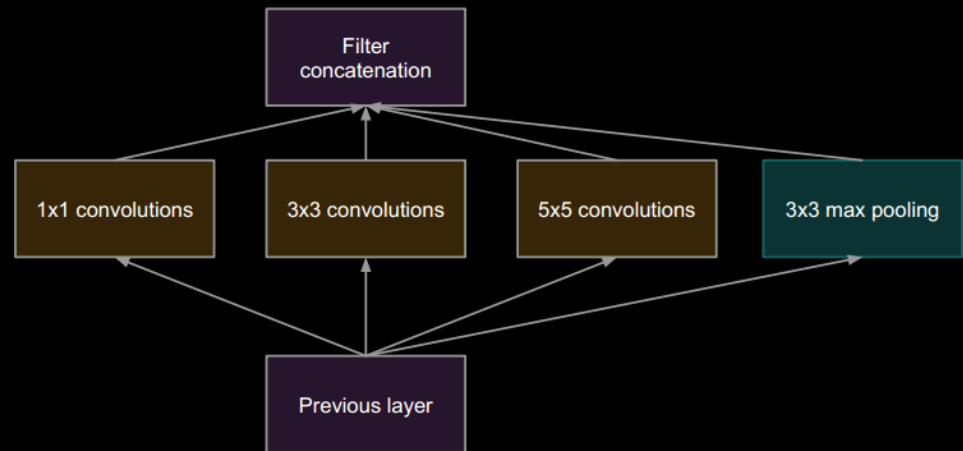


GoogLeNet: Inception Module

Idea 1:

Use 1x1, 3x3, and 5x5 convolutions in parallel to capture a variety of structures

Also add a parallel max pooling path



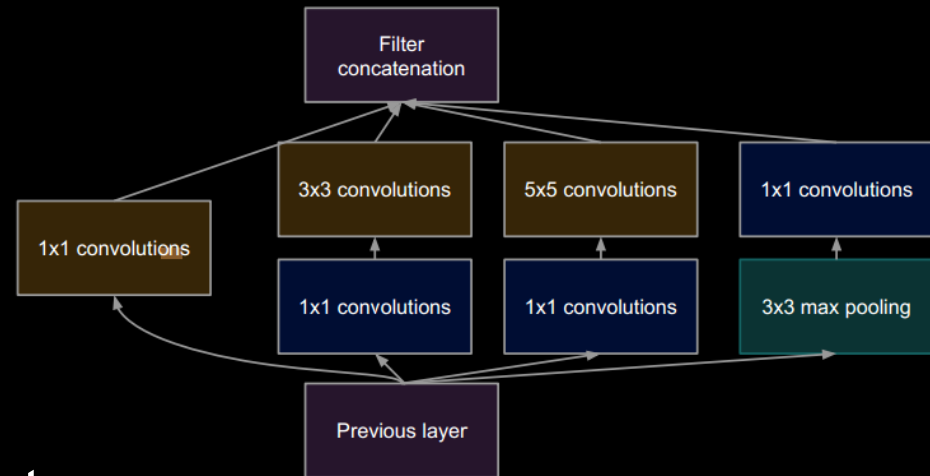
The problem: Computational Expense quickly balloons

GoogLeNet: Inception Module

Idea 2:

Use 1x1 convolutional layers for dimensional reduction.

- Limits computational blow up from increasing parameters
- The 1x1 convolutions also use ReLUs, so provide an added element of non-linearity

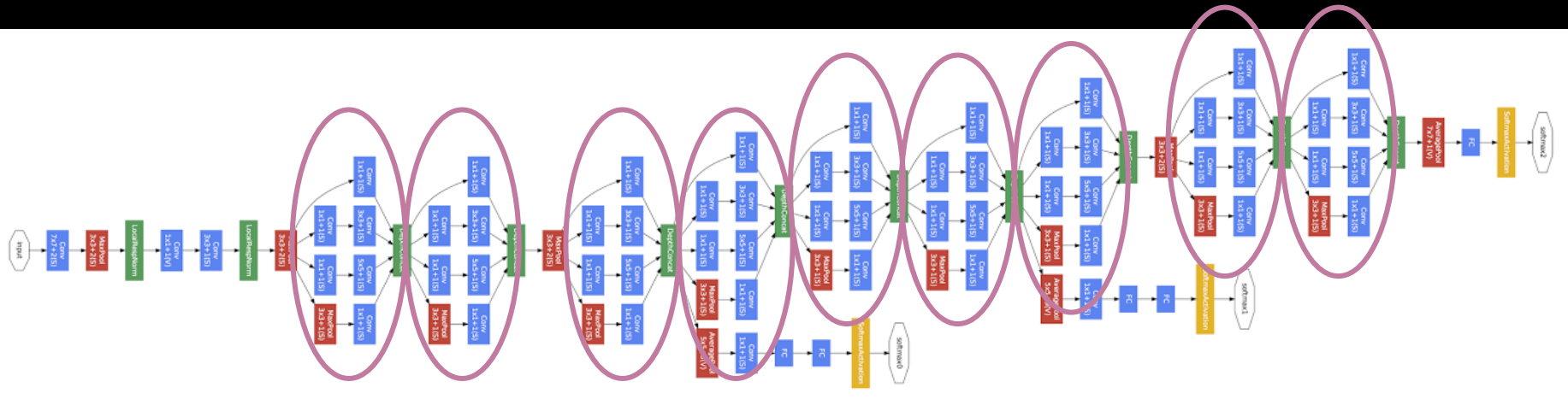


GoogLeNet

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

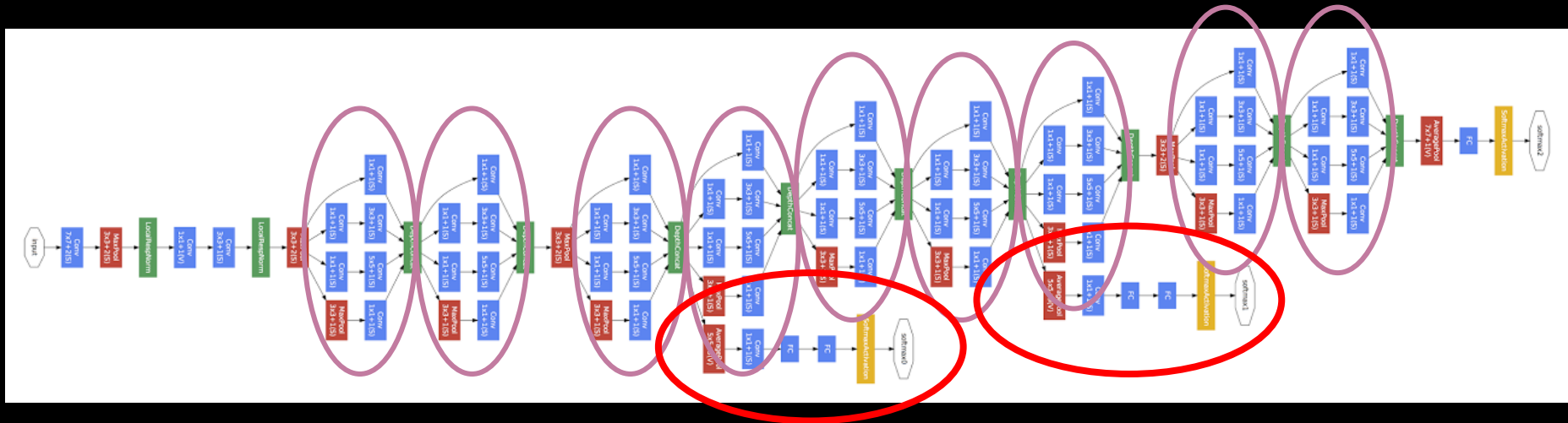
Note, they also replaced FC with avg pool

GoogLeNet



Not so scary after all, just 9 of these “inception modules” stacked on top of each other

GoogLeNet



Not so scary after all, just 9 of these “inception modules” stacked on top of each other

...but wait, what are *these things*?

GoogLeNet: Auxiliary Classifiers

Problem:

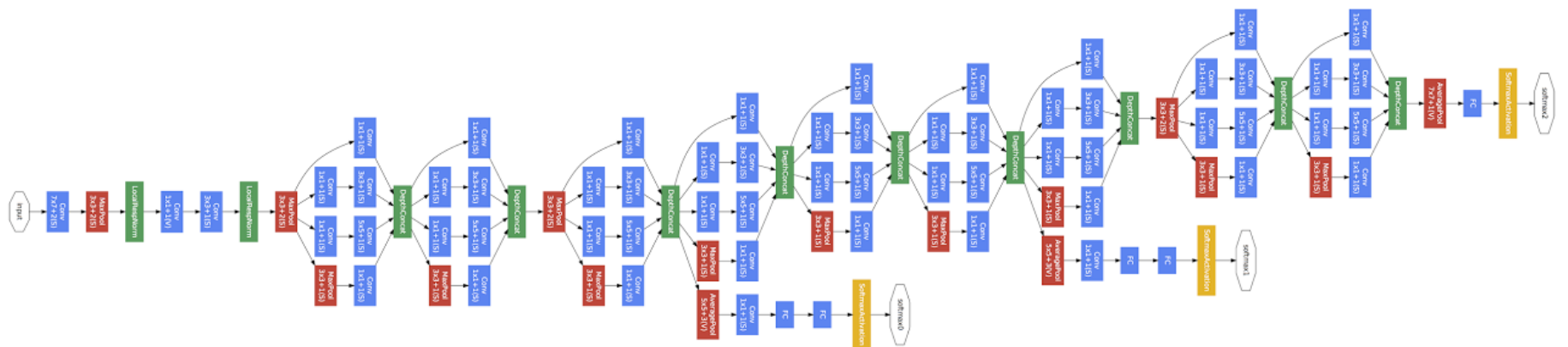
The depth of the network raises concerns about the effectiveness of the backpropagating gradient

Their solution?

Throw on auxiliary classifiers part way through

- Small convnets with a pooling layer, a 1x1 convolution layer, fully connected layers, and softmax loss layer on 1000 classes
- Combined with backprop'd loss with relative weight of 0.3
- (Removed at test time)

GoogLeNet



Now the structure should be less intimidating

GoogLeNet: Stats

12x fewer parameters than AlexNet

- The move away from fully connected layers near the top of the network helps with this

22 Layers deep

~2x more operations than AlexNet

GoogLeNet: Results

Team	Year	Place	Error (top-5)	Uses external data
SuperVision	2012	1st	16.4%	no
SuperVision	2012	1st	15.3%	Imagenet 22k
Clarifai	2013	1st	11.7%	no
Clarifai	2013	1st	11.2%	Imagenet 22k
MSRA	2014	3rd	7.35%	no
VGG	2014	2nd	7.32%	no
GoogLeNet	2014	1st	6.67%	no

Table 2: Classification performance

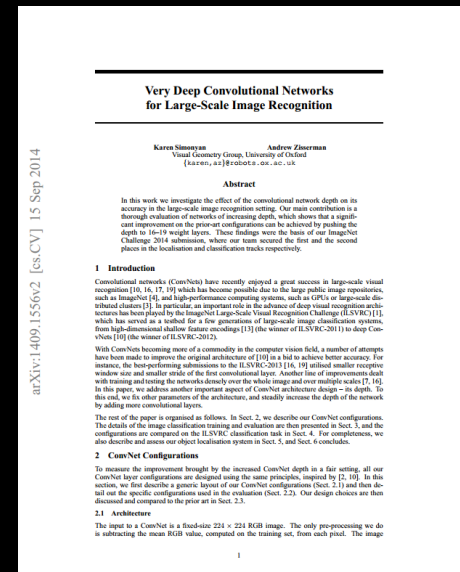
GoogLeNet: Take-aways

- Once more we see the replacement of a fully-connected layer with global average pooling
- 1x1 Convolutional filters similar to multilayer perceptrons in Network in Network paper
- Concatenation of different size convolutional filters
- Mid-network classification to improve backpropagation signal and increase mid-network discriminant abilities

Very Deep ConvNets

Simonyan, Karen, and Andrew Zisserman. "Very Deep Convolutional Networks for Large-Scale Image Recognition." *arXiv preprint arXiv:1409.1556* (2014).

How much effect does extra depth add?



Very Deep ConvNets

Basic idea:

- Stack a bunch of convolutional layers on top of each other, with occasional max-pooling
- All convolutional layers either 3x3 or 1x1
 - Stacks of 3x3 layers have equivalent receptive fields to larger convolutional filters
 - 1x1 convolutions being used here, again, to introduce extra non-linearity
(input/output channels' dimensions are equal here)

Very Deep ConvNets

Why stacks of 3x3 Convolutions?

- Added discriminative ability from more ReLU layers
- Effective receptive field equivalent to larger convolutions
- Fewer parameters

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

Table 2: **Number of parameters** (in millions).

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

Very Deep ConvNets: Results

Table 3: ConvNet performance at a single test scale.

ConvNet config. (Table 1)	smallest image side		top-1 val. error (%)	top-5 val. error (%)
	train (S)	test (Q)		
A	256	256	29.6	10.4
A-LRN	256	256	29.7	10.5
B	256	256	28.7	9.9
C	256	256	28.1	9.4
	384	384	28.1	9.3
	[256;512]	384	27.3	8.8
D	256	256	27.0	8.8
	384	384	26.8	8.7
	[256;512]	384	25.6	8.1
E	256	256	27.3	9.0
	384	384	26.9	8.7
	[256;512]	384	25.5	8.0

→ Deeper is better

Very Deep ConvNets: Take-aways

- Again we see deeper nets pushing the state of the art
- Once more, greater non-linearity improving network ability
 - Both through 1x1 convolutions and stacks of 3x3