

# Structural Priors in Deep Networks

Yani Ioannou

University of Cambridge

August 29, 2017



# Overview

## Introduction

## Research Overview

PhD Research

Collaborative Research

## Motivation

## Structural Priors

Spatial Structural Priors

Filter-wise Structural Priors

## Summary/Future Work

Collaborative Research

# Introduction



- ▶ Ph.D. student in the Department of Engineering at the University of Cambridge.
- ▶ Funded by a Microsoft Research PhD Scholarship
- ▶ Supervised by Professor Roberto Cipolla, head of the Computer Vision and Robotics group in the Machine Intelligence Lab, and Dr. Antonio Criminisi, a principal researcher at Microsoft Research.



# Research Overview



- ▶ Decision Forests, Convolutional Networks and the Models in-Between.  
Y. Ioannou, D. Robertson, D. Zikic, P. Kotschieder, J. Shotton, M. Brown, A. Criminisi.  
MSR Technical Report 2015
- ▶ \*Training CNNs with Low-Rank Filters for Efficient Image Classification.  
Y. Ioannou, D. Robertson, J. Shotton, R. Cipolla, A. Criminisi.  
ICLR 2016
- ▶ \*Deep roots: Improving CNN efficiency with hierarchical filter groups.  
Y. Ioannou, D. Robertson, R. Cipolla, A. Criminisi.  
CVPR 2017

---

\*To be presented in this talk

## ▶ **Medical Computer Vision**

- ▶ Segmentation of brain tumor tissues with convolutional neural networks.

D. Zikic, Y. Ioannou, M. Brown, A. Criminisi. *MICCAI-BRATS 2014*

- ▶ Using CNNs for Malaria Diagnosis.  
Intellectual Ventures/Gates Foundation

## ▶ **Adversarial Examples**

Measuring Neural Net Robustness with Constraints.

O. Bastani, Y. Ioannou, L. Lampropoulos, D. Vytiniotis, A. Nori, A. Criminisi. *NIPS 2016*

## ▶ **Neural Network Design**

Refining Architectures of Deep Convolutional Neural Networks.

S. Shankar, D. Robertson, Y. Ioannou, A. Criminisi, R. Cipolla. *CVPR 2016*

# Motivation



## Imagenet Large-Scale Visual Recognition Challenge



- ▶ Imagenet Large-Scale Visual Recognition Challenge<sup>2</sup>.
- ▶ 1.2 Million Training Images, 1000 classes.
- ▶ 50,000 image validation/test set.
  - ▶ In 2012 Alex Krizhevsky won challenge with CNN<sup>3</sup>.
  - ▶ 'AlexNet' was 26.2% better than second best, 15.3%.
- ▶ State-of-the-art beats human error (5%).

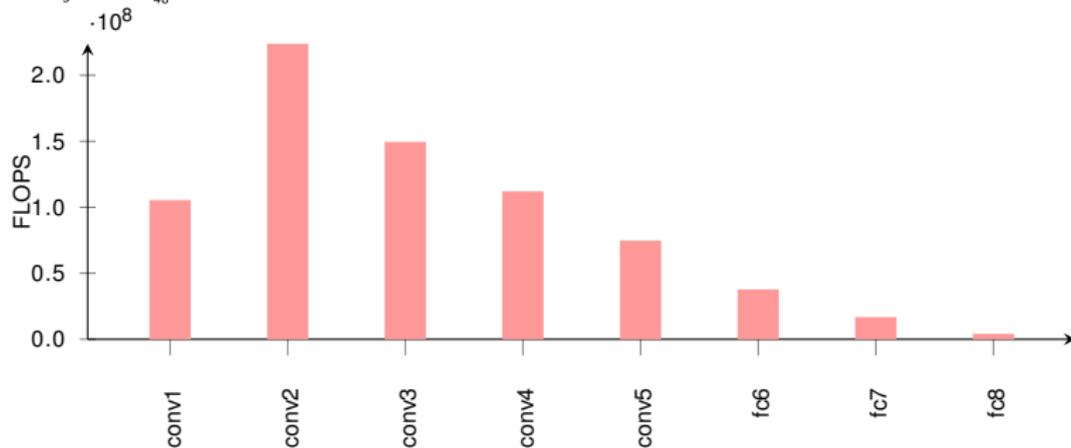
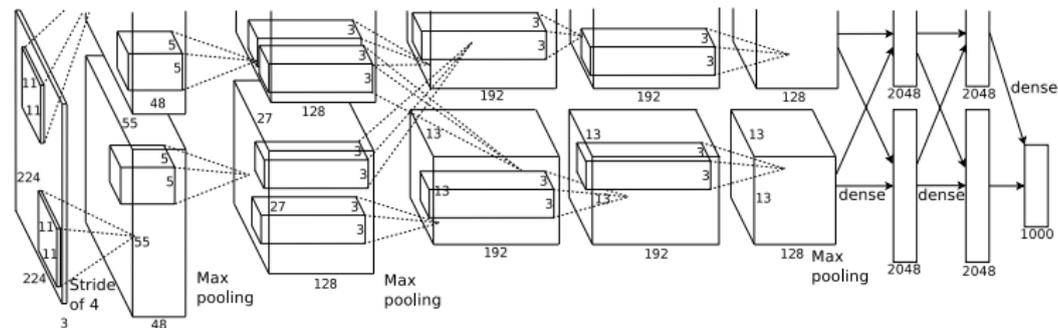
---

<sup>2</sup>Russakovsky et al., "ImageNet Large Scale Visual Recognition Challenge".

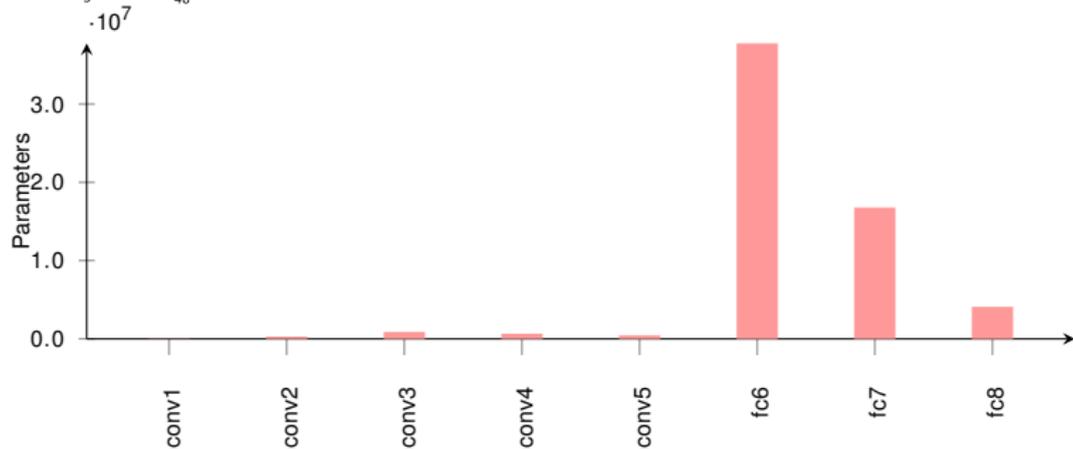
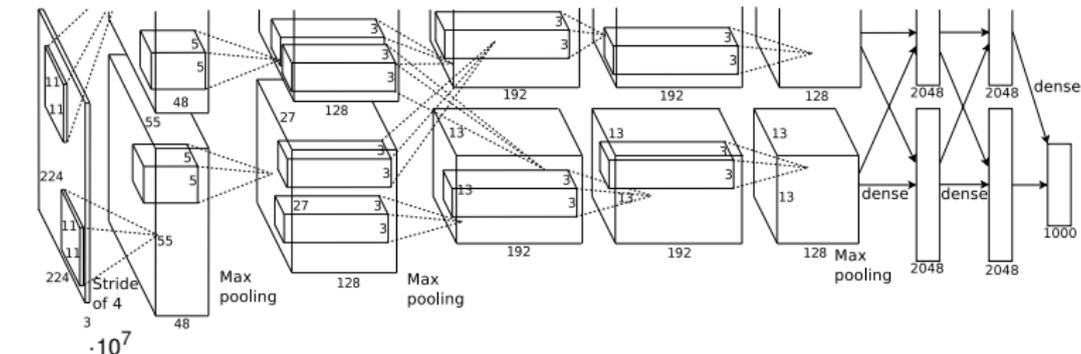
<sup>3</sup>Krizhevsky, Sutskever, and Hinton, "ImageNet Classification with Deep Convolutional Neural Networks".



# AlexNet Complexity - FLOPS

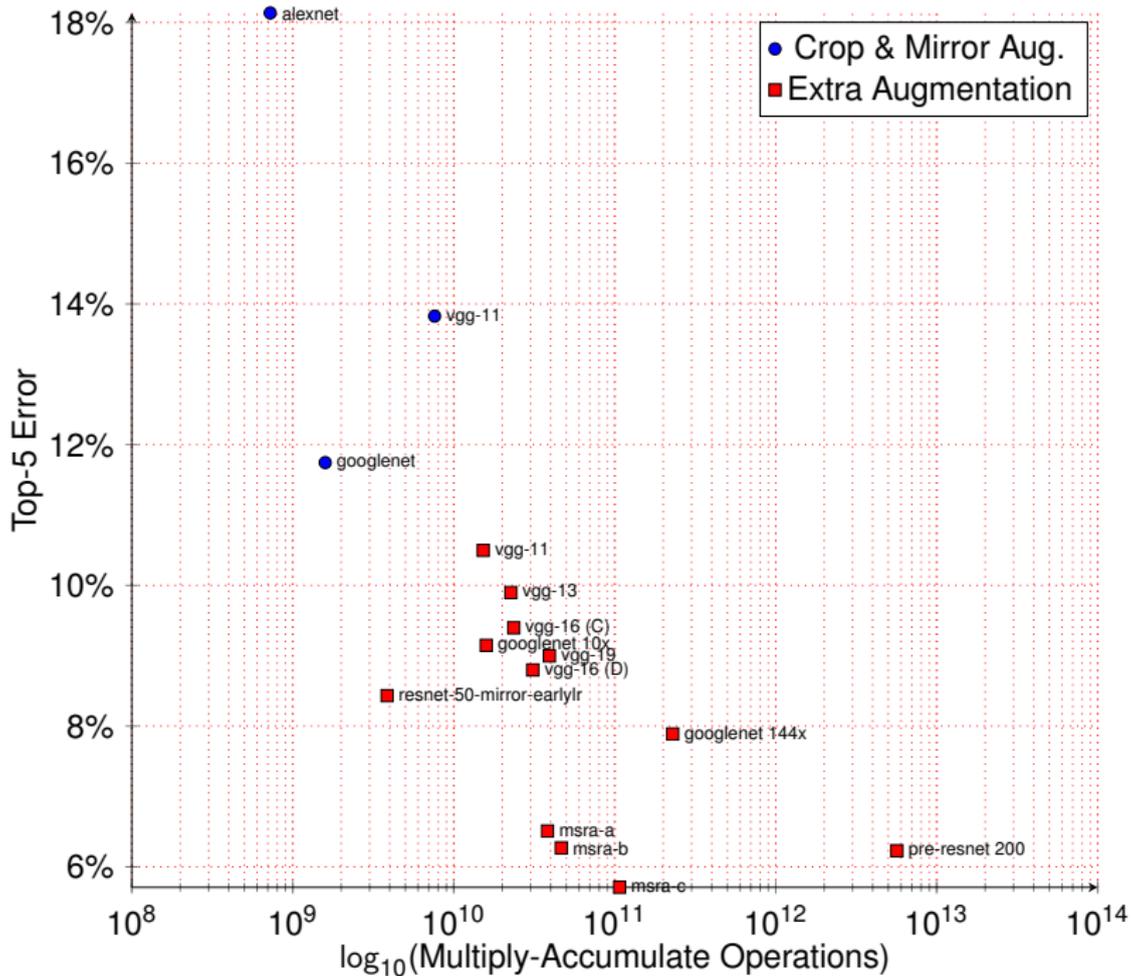


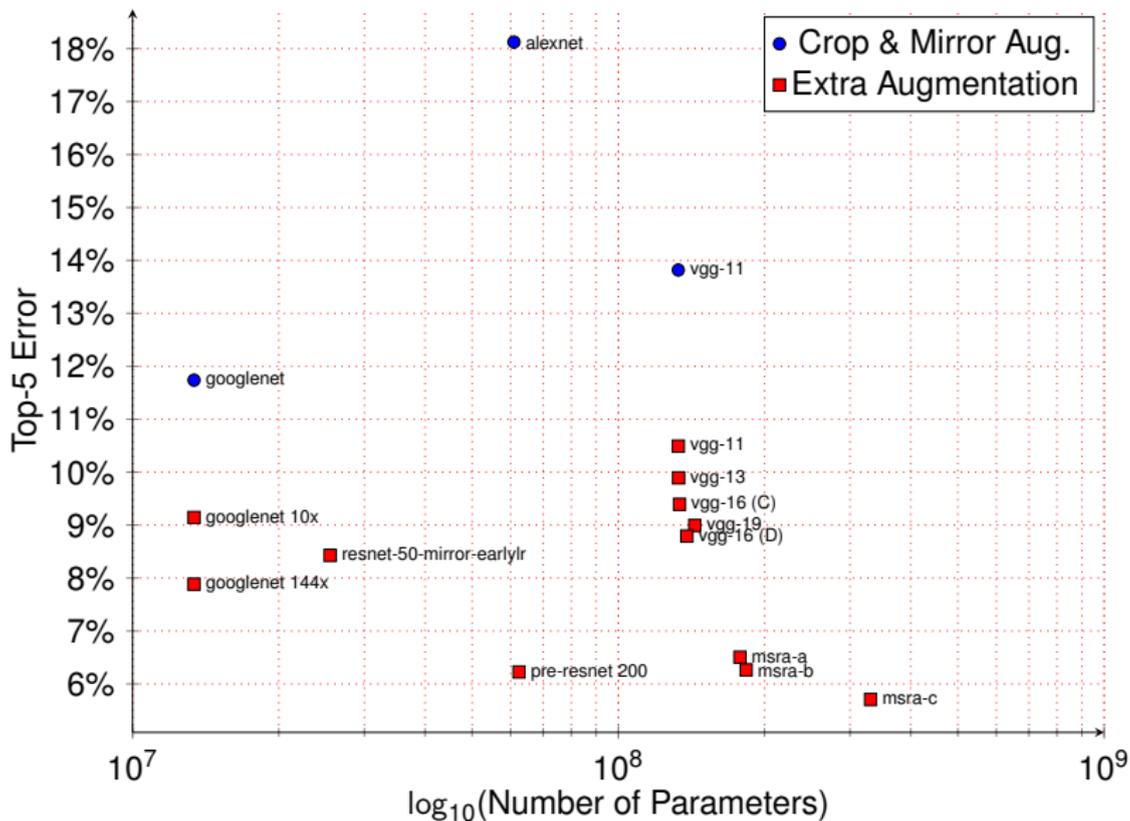
# AlexNet Complexity - Parameters



96% in fully connected layers

≈





# The Problem

- ▶ Creating a massively over-parameterized network, has consequences
- ▶ Training time: Translates into 2-3 weeks of training on 8 GPUs! (ResNet 200)
- ▶ Forward pass (ResNet 50): 12 ms GPU, 621 ms CPU
- ▶ Forward pass (GoogLeNet): 4.4 ms GPU, 300 ms CPU

But what about the practicalities of using deep learning:

- ▶ on embedded devices
- ▶ realtime applications
- ▶ backed by distributed/cloud computing

# Compression/Representation

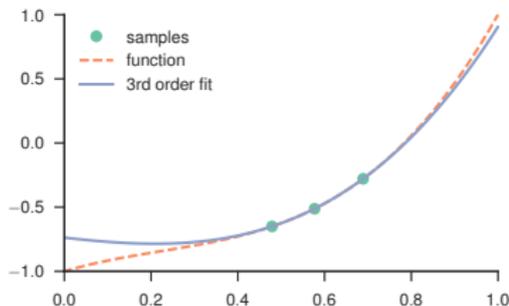
Isn't that already being addressed?

- ▶ Approximation (compression/pruning) of neural networks
- ▶ Reduced representation (8-bit floats/binary!)

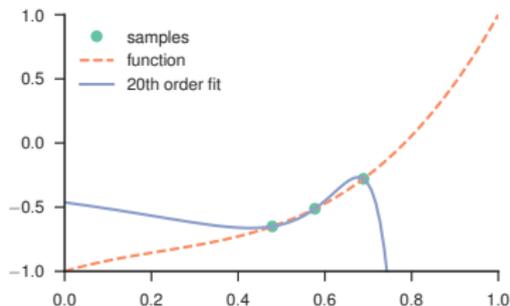
Allow us to have a trade off in compute v.s. accuracy.

*These methods will still apply to any network.* Instead, let's try to address the fundamental problem of over-parameterization.

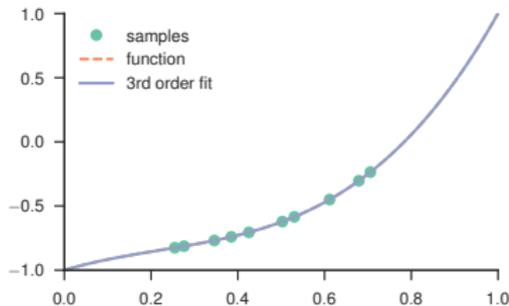
# Generalization and Num. Parameters



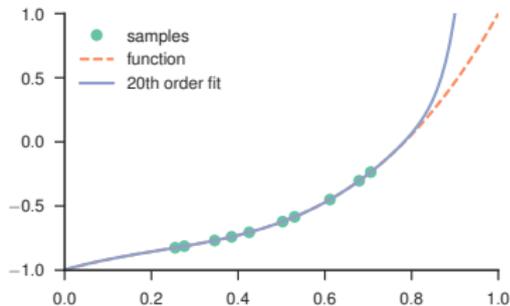
(a) 3<sup>rd</sup>-order poly., 3 points



(b) 20<sup>th</sup>-order poly., 3 points



(c) 3<sup>rd</sup>-order poly., 10 points



(d) 20<sup>th</sup>-order poly., 10 points

- ▶ When fitting a curve, we often have little idea of what order polynomial would best fit the data!
- ▶ Weak Prior - Regularization.
  - ▶ Prior is knowing only that our model is over-parameterized
  - ▶ This restricts the model to effectively use only a small number of the parameters
- ▶ Strong Priors - Structural.
  - ▶ With more prior information on the task, *e.g.* from the convexity of the polynomial, we may imply that a certain order polynomial is more appropriate, and restrict learning to some particular orders.

- ▶ Deep networks need many more parameters than data points because they aren't just learning to model data, but also learning what *not* to learn.
- ▶ Idea: Why don't we help the network, through structural priors, not to learn things it doesn't need to?

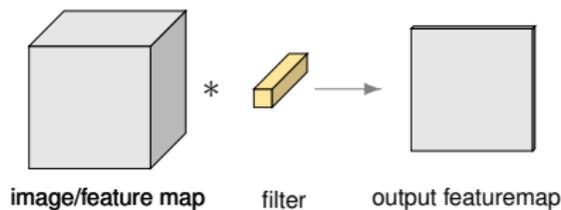
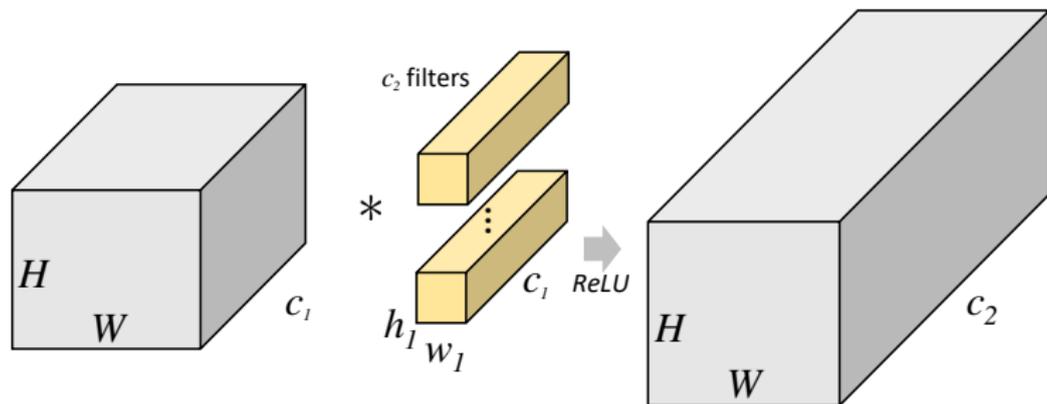


The Atlas Slave  
(Accademia, Florence)

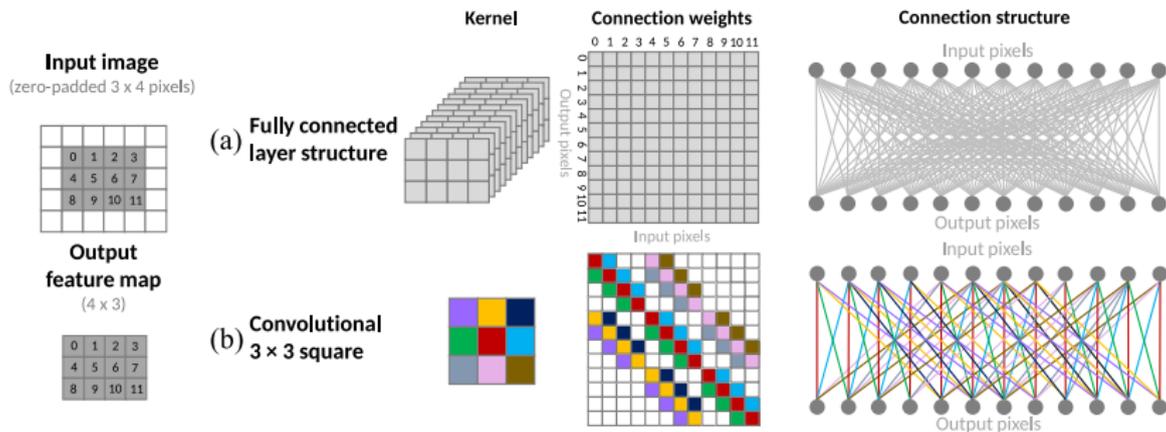
# Structural Priors



# Typical Convolutional Layer



# Sparsity of Convolution



- ▶ Convolutional Neural Networks (CNNs) are structural priors for natural images
- ▶ Local connectivity - local correlations are important in natural images, *e.g.* edges
- ▶ Shared parameters - we know we don't need to re-learn filters for every pixel

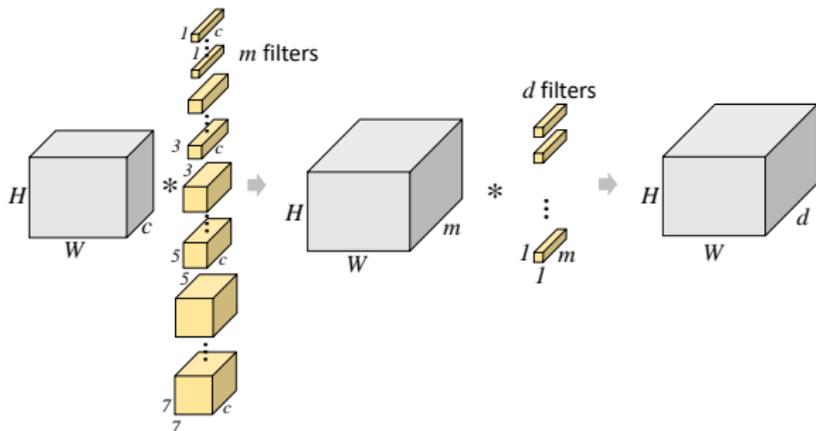
# Why are CNNs uniformly structured?

*“The marvelous powers of the brain emerge not from any single, uniformly structured connectionist network but from highly evolved arrangements of smaller, specialized networks which are interconnected in very specific ways.”*

Marvin Minsky  
Perceptrons (1988 edition)

- ▶ Deep networks are largely monolithic (uniformly connected), with few exceptions
- ▶ Why don't we try to structure our networks closer to the specialized components required for learning images?

- ▶ In<sup>5</sup>, linear combination of different sized filters is learned, *i.e.* a basis space for filters:



- ▶ Motivation: expect most image correlations to be highly localized, *i.e.* many small filters. However, a few may require larger, more complex filters

<sup>5</sup>Szegedy et al., "Going Deeper with Convolutions".

# Spatial Structural Priors

Published as a conference paper at ICLR 2016

---

## TRAINING CNNs WITH LOW-RANK FILTERS FOR EFFICIENT IMAGE CLASSIFICATION

Yani Ioannou<sup>1</sup>, Duncan Robertson<sup>2</sup>, Jamie Shotton<sup>2</sup>, Roberto Cipolla<sup>1</sup> & Antonio Criminisi<sup>2</sup>

<sup>1</sup>University of Cambridge, <sup>2</sup>Microsoft Research

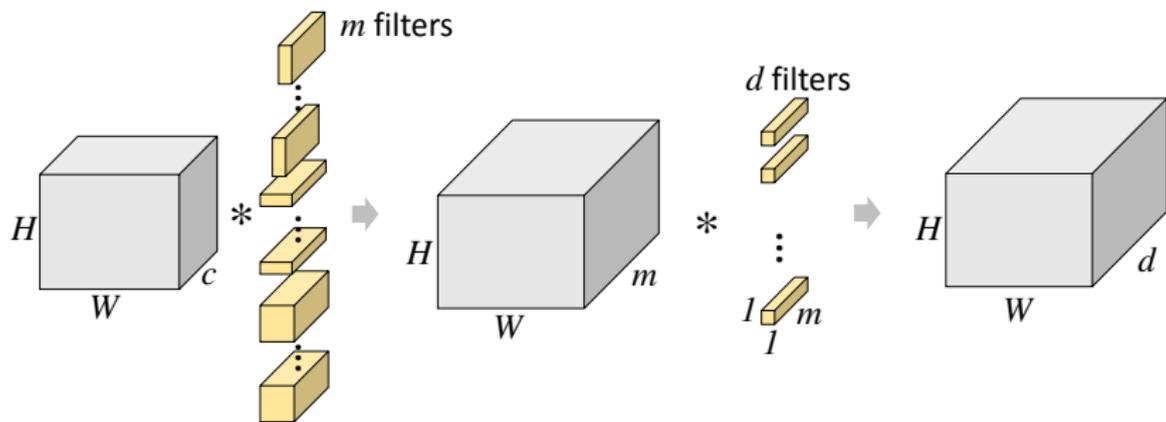
{yai20, rc10001}@cam.ac.uk, {a-durobe, jamiesho, antcrim}@microsoft.com

### ABSTRACT

We propose a new method for creating computationally efficient convolutional neural networks (CNNs) by using low-rank representations of convolutional filters. Rather than approximating filters in previously-trained networks with more efficient versions, we learn a set of small basis filters from scratch; during training, the network learns to combine these basis filters into more complex filters that

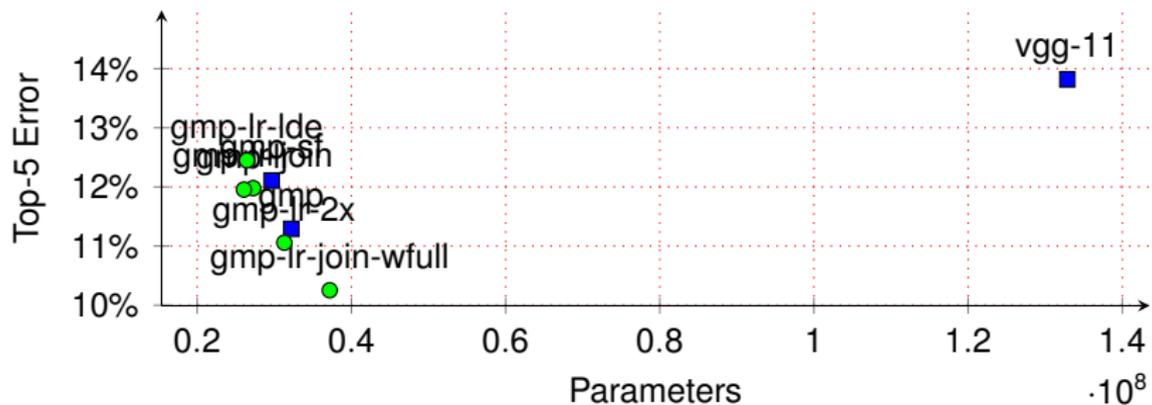
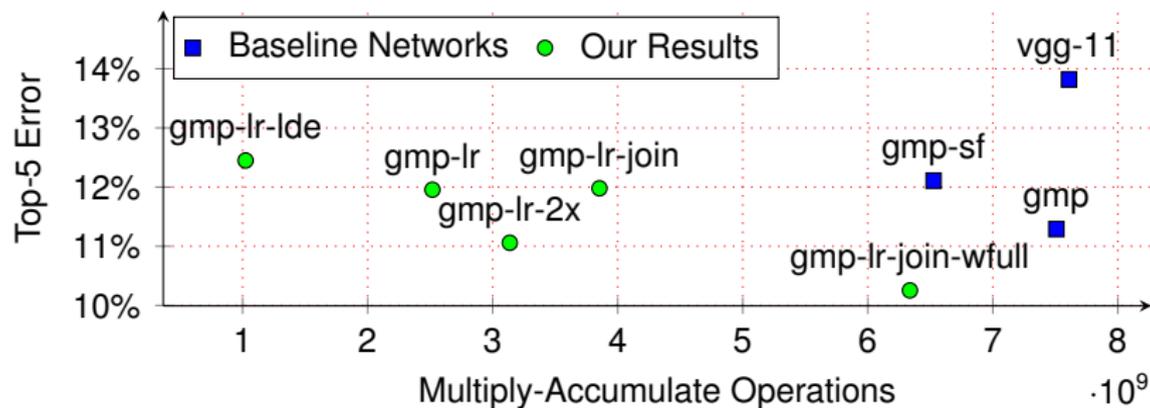
# Proposed Method

Learning a basis for filters.



- ▶ A learned basis of vertical/horizontal rectangular filters and square filters!
- ▶ Shape of learned filters is a full  $w \times h \times c$ .
- ▶ But what can be effectively learned is limited by the number and complexity of the components.

# VGG/Imagenet Results



# Imagenet Results

- ▶ VGG-11 (low-rank): **24%** smaller, **41%** fewer FLOPS
- ▶ VGG-11 (low-rank/full-rank mix): **16%** fewer FLOPS with **1% lower error** on ILSRVC val, but 16% larger.
- ▶ GoogLeNet: **41%** smaller, **26%** fewer FLOPS

Or better results if you tune it on GoogLeNet more. . .

# Rethinking the Inception Architecture for Computer Vision

Christian Szegedy  
Google Inc.

szegedy@google.com

Vincent Vanhoucke  
vanhoucke@google.com

Sergey Ioffe  
sioffe@google.com

Jonathon Shlens  
shlens@google.com

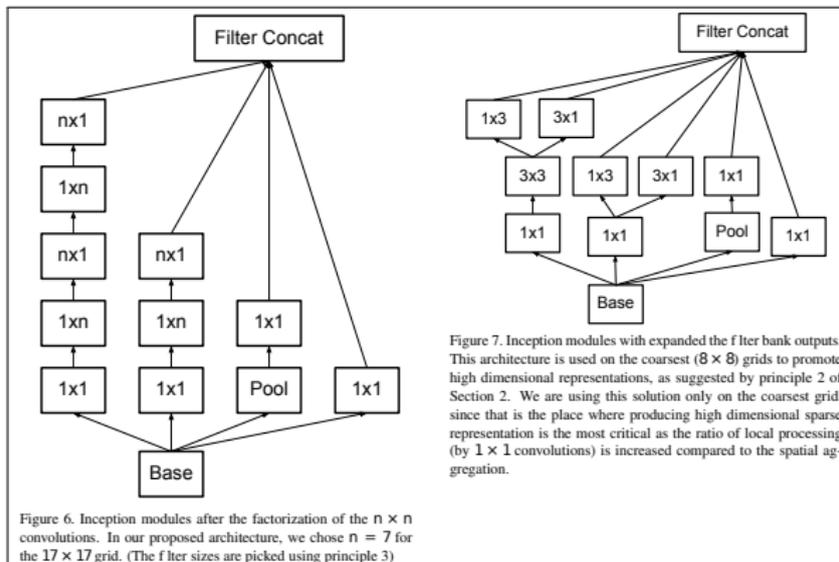


Figure 6. Inception modules after the factorization of the  $n \times n$  convolutions. In our proposed architecture, we chose  $n = 7$  for the  $17 \times 17$  grid. (The filter sizes are picked using principle 3)

Figure 7. Inception modules with expanded filter bank outputs. This architecture is used on the coarsest ( $8 \times 8$ ) grids to promote high dimensional representations, as suggested by principle 2 of Section 2. We are using this solution only on the coarsest grid, since that is the place where producing high dimensional sparse representation is the most critical as the ratio of local processing (by  $1 \times 1$  convolutions) is increased compared to the spatial aggregation.

Convo  
of-the-art  
tasks. Sin  
to become  
ous bench  
computatio  
for most  
for traini  
count are  
mobile vi  
ing ways  
the addre  
factorize  
benchmark  
challenge  
the state

single frame evaluation using a network with a computational cost of 5 billion multiply-adds per inference and with using less than 25 million parameters. With an ensemble of 4 models and multi-crop evaluation, we report 3.5% top-5 error and 17.3% top-1 error.

## 1. Introduction

Since the 2012 ImageNet competition [16] winning en-

signifi-  
gains  
signifi-  
cantly.  
per con-  
perform-  
increas-  
Also,  
appli-  
where  
needed,  
n[4].

ure of  
evalu-  
On the  
at [20]  
at con-  
exam-

ple, GoogleNet employed only 5 million parameters, which represented a  $12 \times$  reduction with respect to its predecessor AlexNet, which used 60 million parameters. Furthermore, VGGNet employed about  $3 \times$  more parameters than AlexNet.

The computational cost of Inception is also much lower than VGGNet or its higher performing successors [6]. This has made it feasible to utilize Inception networks in big-data scenarios [17], [13], where huge amount of data needed to be processed at reasonable cost or scenarios where memory

# Filter-wise Structural Priors

## Deep Roots: Improving CNN Efficiency with Hierarchical Filter Groups

Yani Ioannou<sup>1</sup>    Duncan Robertson<sup>2</sup>    Roberto Cipolla<sup>1</sup>  
Antonio Criminisi<sup>2</sup>

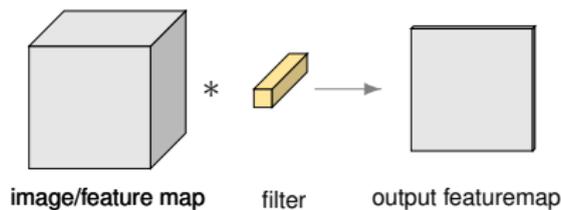
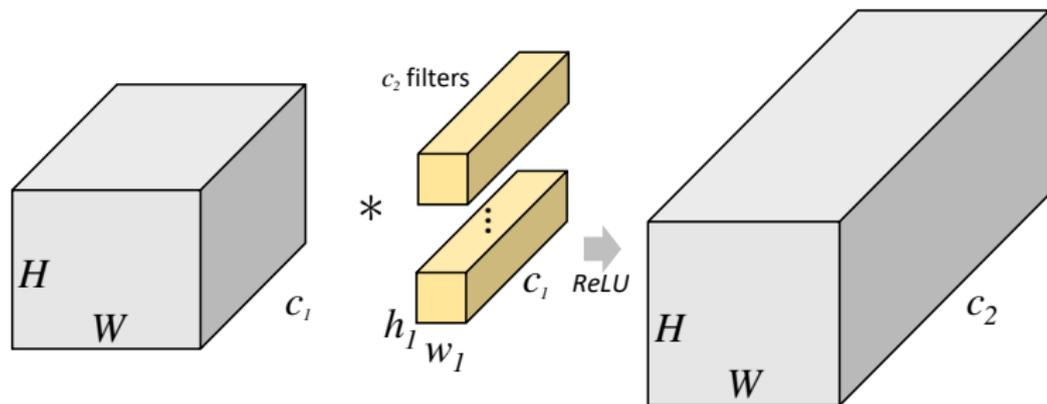
<sup>1</sup>University of Cambridge, <sup>2</sup>Microsoft Research

### Abstract

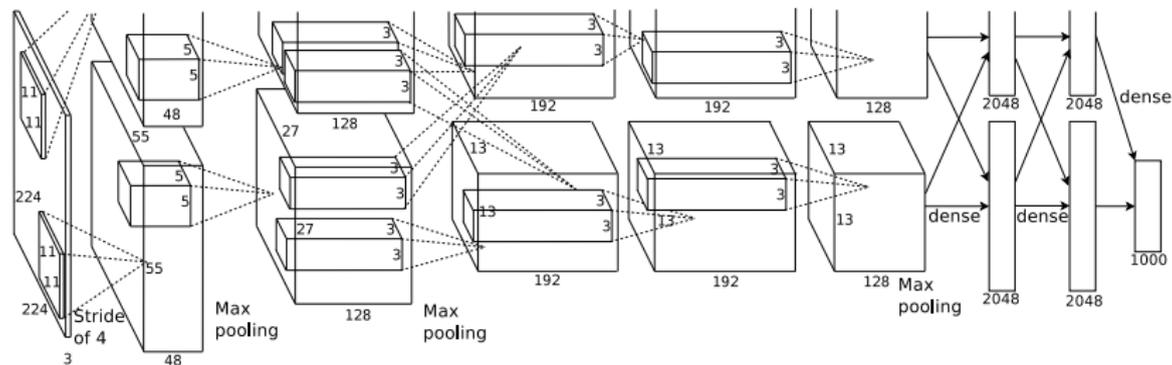
*We propose a new method for creating computation-ally efficient and compact convolutional neural networks*

be achieved by weight decay or dropout [5]. Furthermore, a carefully designed sparse network connection structure can also have a regularizing effect. Convolutional Neural Networks (CNNs) [6, 7] embody this idea, using a sparse

# Typical Convolutional Layer

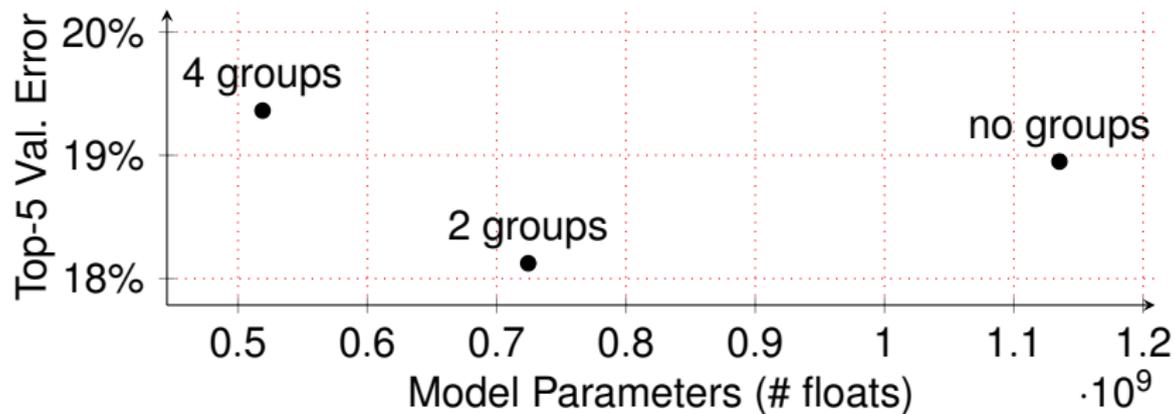
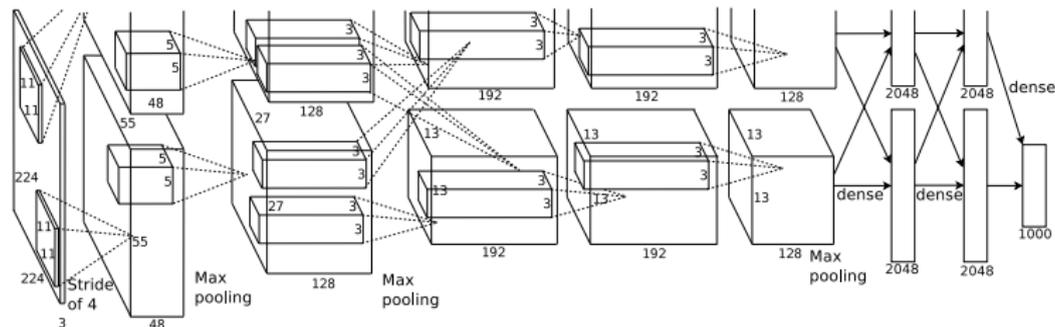


# AlexNet Filter Grouping

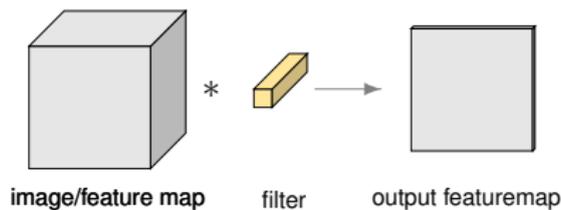
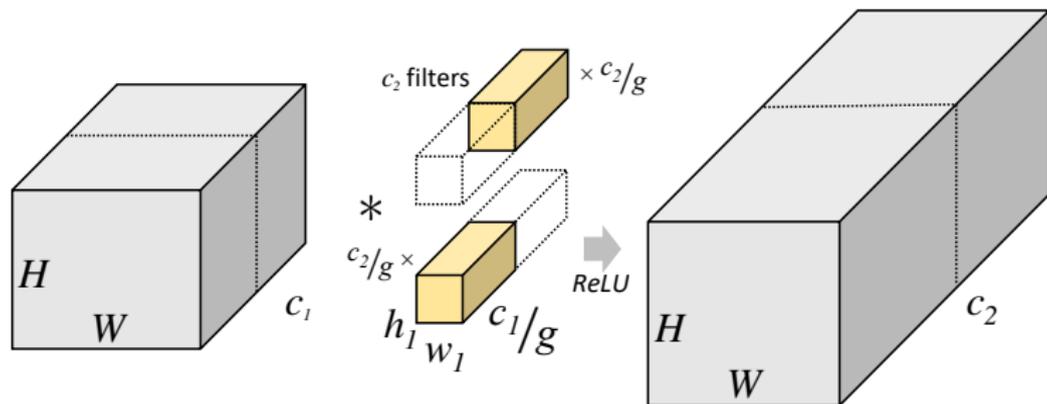


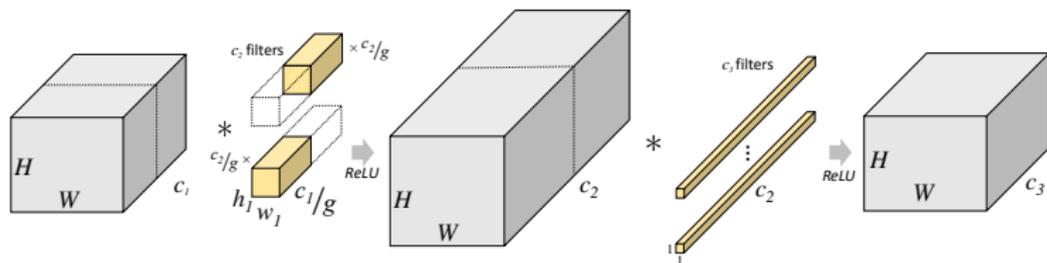
- ▶ Uses 2 filter groups in most of the convolutional layers
- ▶ Allowed training across two GPUs (model parallelism)

# AlexNet Filter Grouping

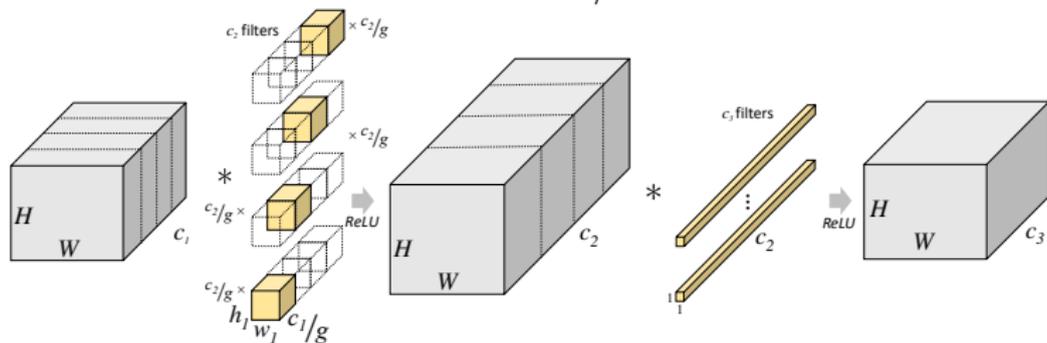


# Grouped Convolutional Layer



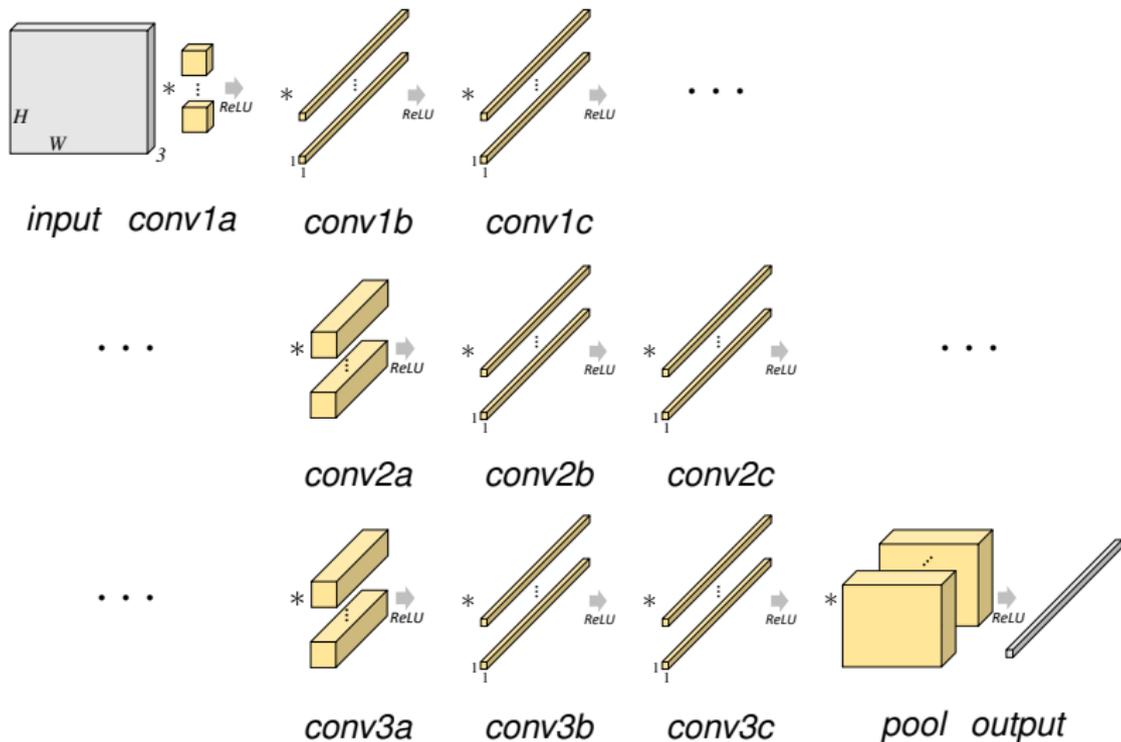


Root-2 Module:  $d$  filters in  $g = 2$  filter groups, of shape  $h \times w \times c/2$ .

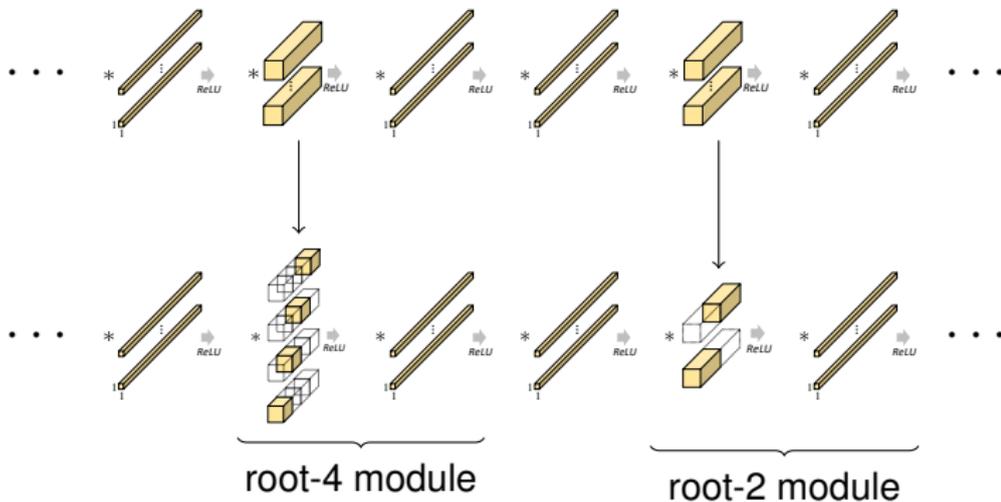


Root-4 Module:  $d$  filters in  $g = 4$  filter groups, of shape  $h \times w \times c/4$ .

# Network-in-Network



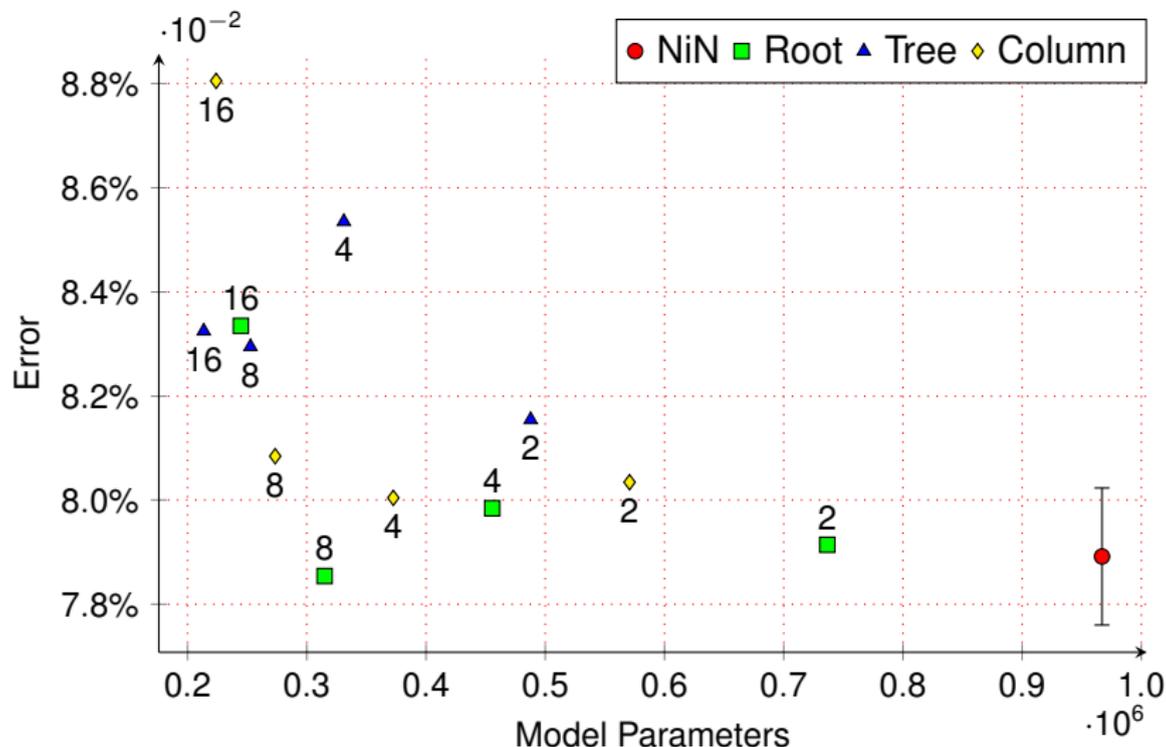
# NiN Root Architectures



**Network-in-Network.** Filter groups in each convolutional layer.

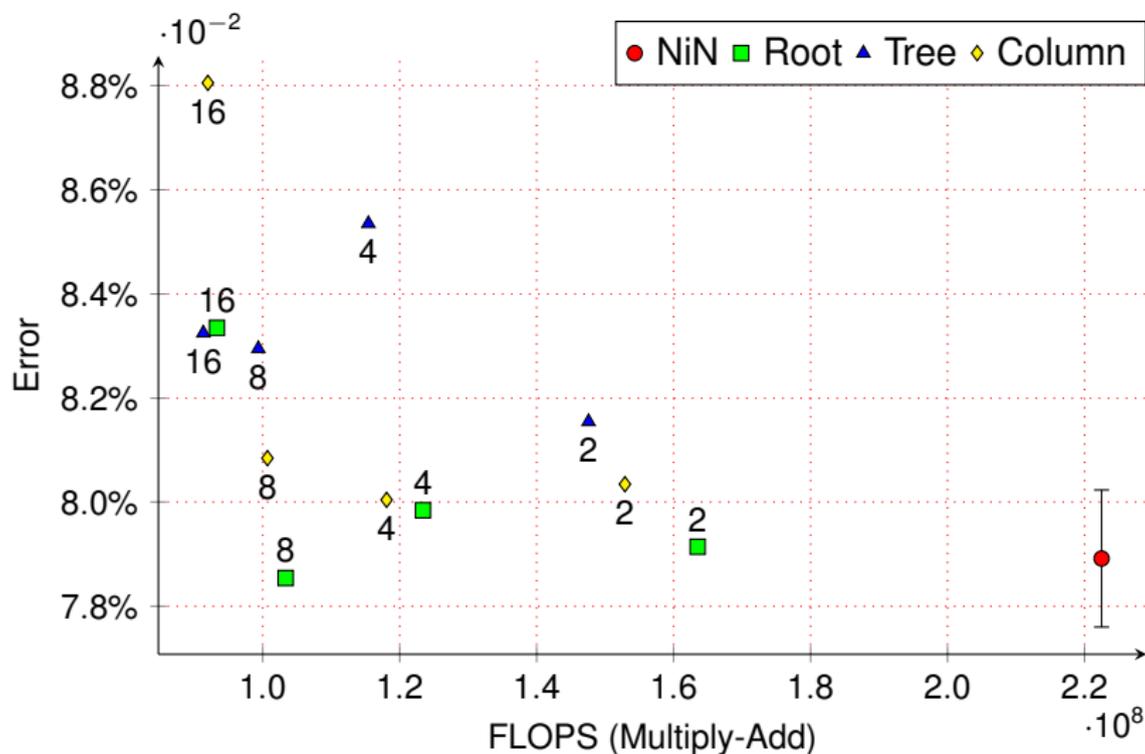
Model	conv1			conv2			conv3		
	<i>a</i>	<i>b</i>	<i>c</i>	<i>a</i>	<i>b</i>	<i>c</i>	<i>a</i>	<i>b</i>	<i>c</i>
	<i>5×5</i>	<i>1×1</i>	<i>1×1</i>	<i>5×5</i>	<i>1×1</i>	<i>1×1</i>	<i>3×3</i>	<i>1×1</i>	<i>1×1</i>
Orig.	1	1	1	1	1	1	1	1	1
root-2	1	1	1	2	1	1	1	1	1
root-4	1	1	1	4	1	1	2	1	1
root-8	1	1	1	8	1	1	4	1	1
root-16	1	1	1	16	1	1	8	1	1

# CIFAR10: Model Parameters v.s. Error



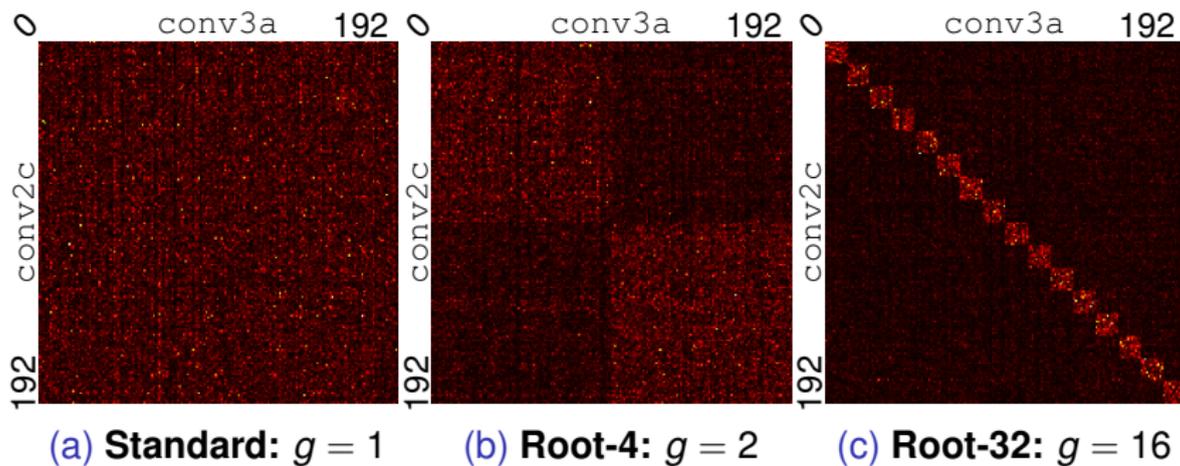
NiN: mean and standard deviation (error bars) are shown over 5 different random initializations.

# CIFAR10: FLOPS (Multiply-Add) v.s. Error.

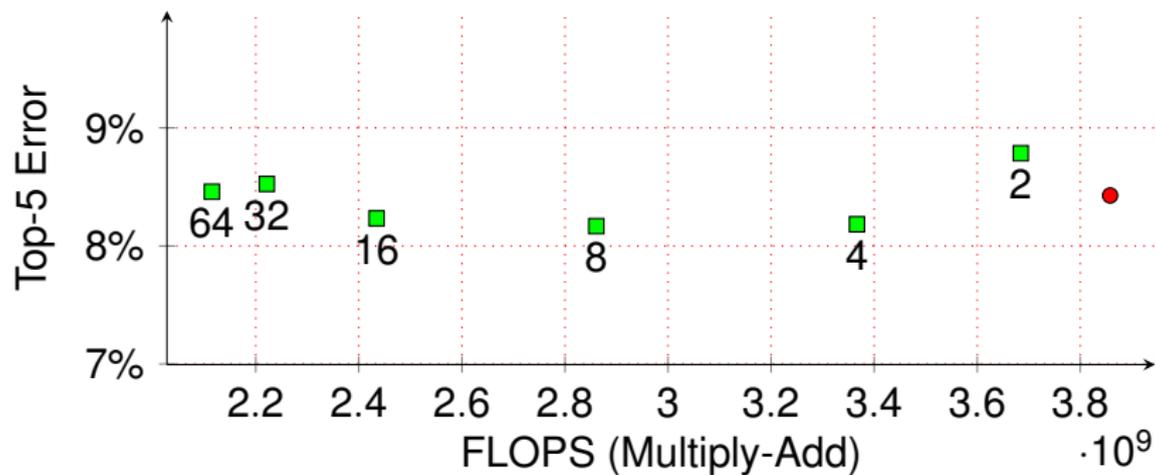
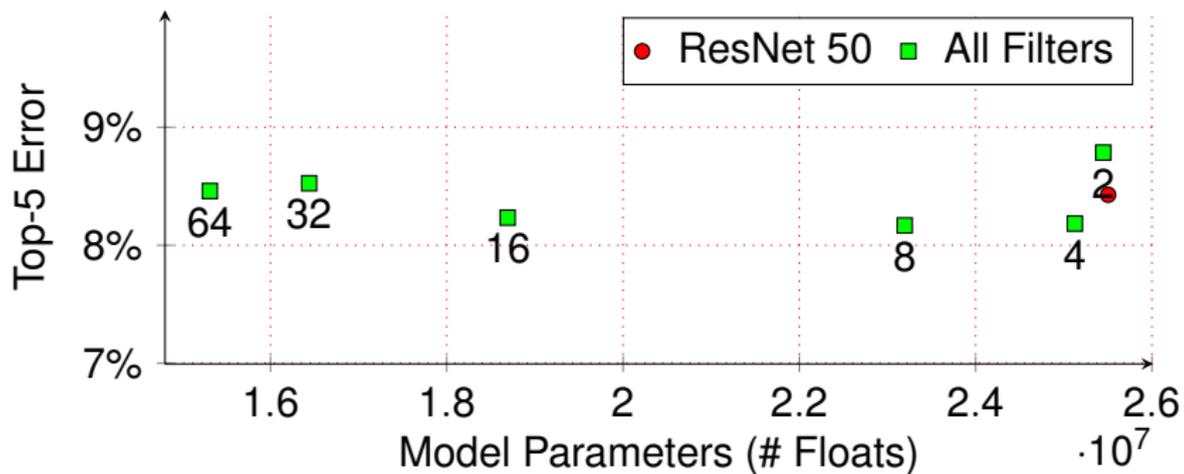


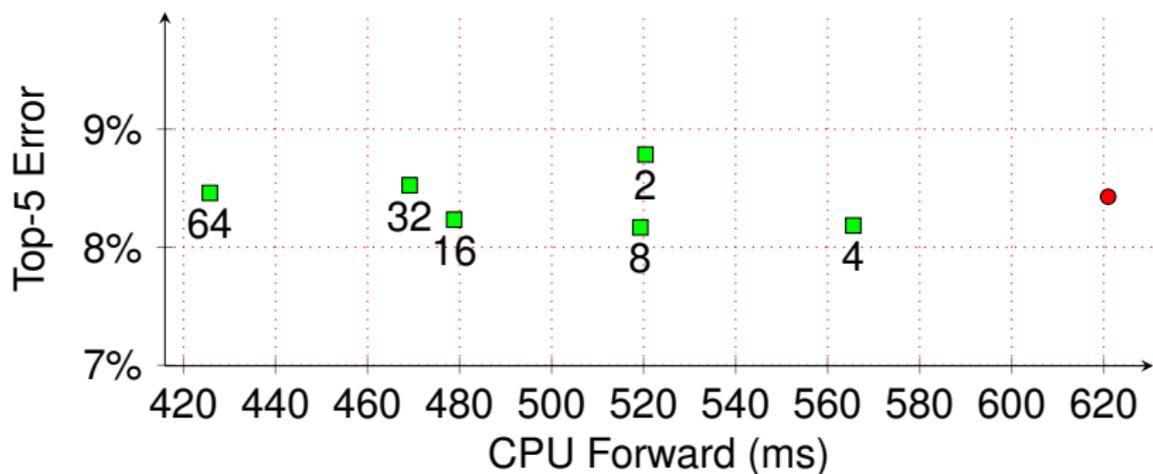
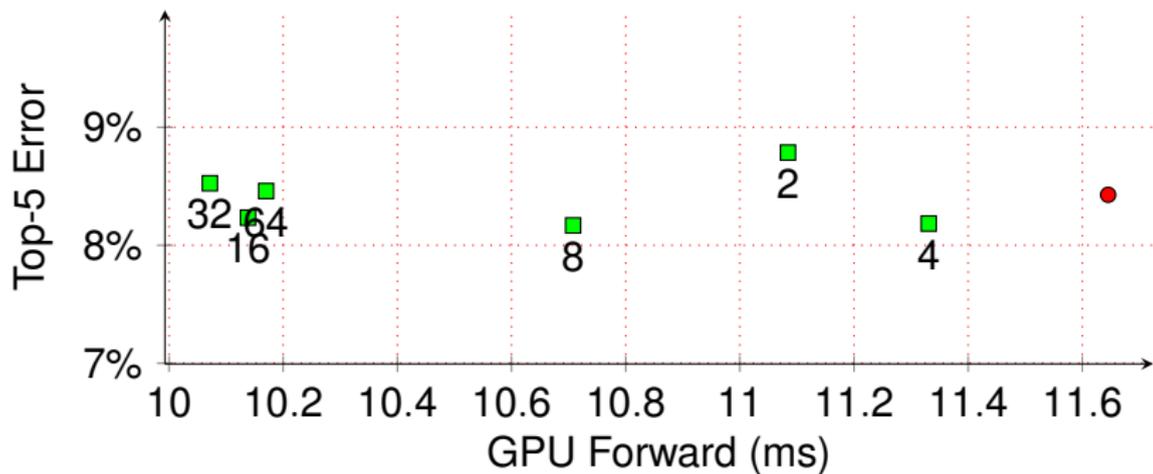
NiN: mean and standard deviation (error bars) are shown over 5 different random initializations.

# Inter-layer Filter Covariance



**Figure:** The block-diagonal sparsity learned by a root-module is visible in the correlation of filters on layers `conv3a` and `conv2c` in the NiN network.





# Imagenet Results

Networks with root modules have similar or higher accuracy than the baseline architectures with much less computation.

- ▶ ResNet 50<sup>6</sup>: **40%** smaller, **45%** fewer FLOPS
- ▶ ResNet 200<sup>7</sup>: **44%** smaller, **25%** fewer FLOPS
- ▶ GoogLeNet: **7%** smaller, **44%** fewer FLOPS

But when you also **increase the number of filters** . . .

---

<sup>6</sup>Caffe Re-implementation

<sup>7</sup>Based on Facebook Torch Model

# Aggregated Residual Transformations for Deep Neural Networks

Saining Xie<sup>1</sup>

Ross Girshick<sup>2</sup>

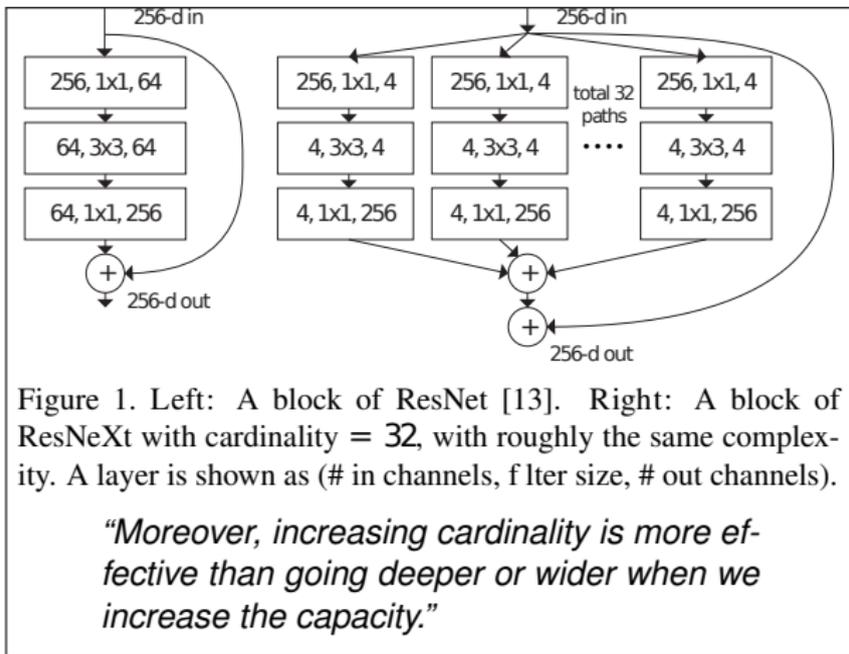
Piotr Dollár<sup>2</sup>

Zhuowen Tu<sup>1</sup>

Kaiming He<sup>2</sup>

<sup>1</sup>UC San Diego

<sup>2</sup>Facebook AI Research



We p  
ecture )  
by repe  
formati  
sults in  
only a f  
new din  
set of tr  
the dim  
dataset,  
conditi  
ity is al  
increasi  
or wide  
dename  
ILSVRC  
place. }  
set and  
than its

ck of  
plex-  
nels).

yper-  
nsion  
licity  
yper-  
GG-  
:cog-  
tasks

[37,  
gned

## 1. Intr

Research on visual recognition is undergoing a transition from “feature engineering” to “network engineering” [24, 23, 43, 33, 35, 37, 13]. In contrast to traditional hand-designed features (e.g., SIFT [28] and HOG [5]), features learned by neural networks from large-scale data [32] require minimal human involvement during training, and can be transferred to a variety of recognition tasks [7, 10, 27]. Nevertheless, human effort has been shifted to designing

topologies are able to achieve competing accuracy with low theoretical complexity. The Inception models have evolved over time [37, 38], but an important common property is a *split-transform-merge* strategy. In an Inception module, the input is split into a few lower-dimensional embeddings (by  $1 \times 1$  convolutions), transformed by a set of specialized filters ( $3 \times 3$ ,  $5 \times 5$ , etc.), and merged by concatenation. It can be shown that the solution space of this architecture is a strict subspace of the solution space of a single deep layer

# Summary/Future Work



- ▶ Using structural priors:
  - ▶ Models are **less computationally complex**
  - ▶ They also use **less parameters**
  - ▶ They significantly help generalization in **deeper networks**
  - ▶ They significantly help generalization with **larger datasets**
- ▶ Are amenable to **model parallelization** (as with original AlexNet), for better parallelism across gpus/nodes

## Future Work: Research

- ▶ We don't always have enough knowledge of the domain to propose good structural priors
- ▶ Our results (and follow up work) do show however that current methods of training/regularization seem to have limited effectiveness in DNNs learning such priors themselves
- ▶ How can we otherwise learn structural priors?

# Future Work: Applications

Both of these methods apply to most deep learning applications:

- ▶ Smaller model state – easier storage and synchronization
- ▶ Faster training and test of models behind ML cloud services
- ▶ Embedded devices/Tensor processing units

And more specific to each method

- ▶ Low-rank filters
  - ▶ Even larger impact for volumetric imagery (Microsoft Radiomics)
- ▶ Root Modules
  - ▶ Model parallelization (Azure/Amazon Cloud)