

Restricted Connectivity in Deep Neural Networks

Yani Ioannou

University of Cambridge

April 17, 2017



Overview

Introduction

Research Overview

PhD 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.
- ▶ 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

Motivation



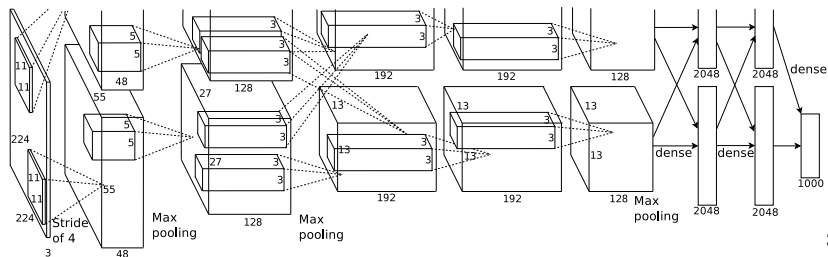
Imagenet Large-Scale Visual Recognition Challenge



- ▶ Imagenet Large-Scale Visual Recognition Challenge¹.
- ▶ 1.2 Million Training Images, 1000 classes.
- ▶ 50,000 image validation/test set.
 - ▶ In 2012 Alex Krizhevsky won challenge with CNN².
 - ▶ 'AlexNet' was 26.2% better than second best, 15.3%.
- ▶ State-of-the-art beats human error (5%).

¹Russakovsky et al., "ImageNet Large Scale Visual Recognition Challenge".

²Krizhevsky, Sutskever, and Hinton, "ImageNet Classification with Deep Convolutional Neural Networks".

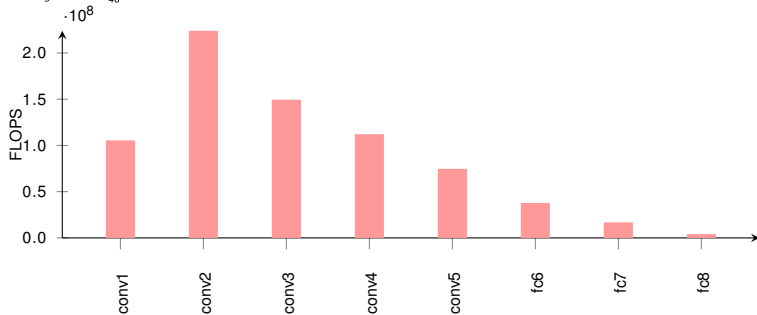
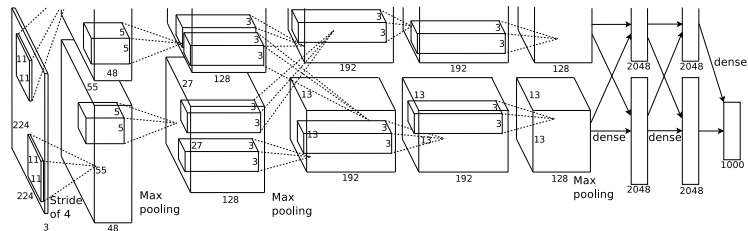


3

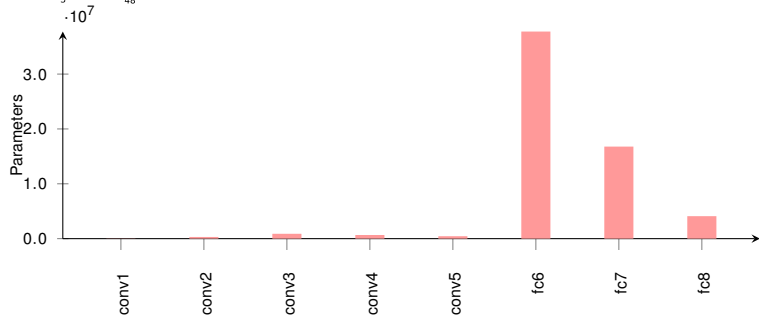
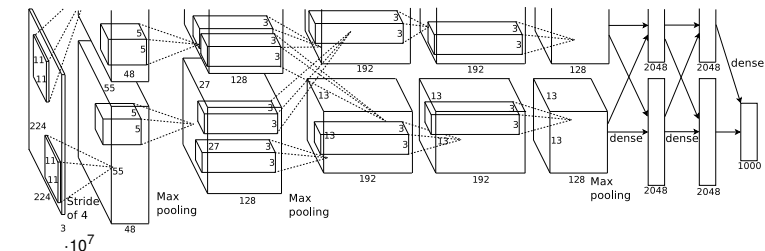
- ▶ ≈ 61 million parameters
- ▶ ≈ 724 million FLOPS (per-sample)
- ▶ Imagenet has 1.28 million training samples ($227 \times 227 \times 3$)
- ▶ Images of dimensions ($227 \times 227 \times 3$) ≈ 200 billion pixels

³Krizhevsky, Sutskever, and Hinton, "ImageNet Classification with Deep Convolutional Neural Networks".

AlexNet Complexity - FLOPS



AlexNet Complexity - Parameters



96% in fully connected layers



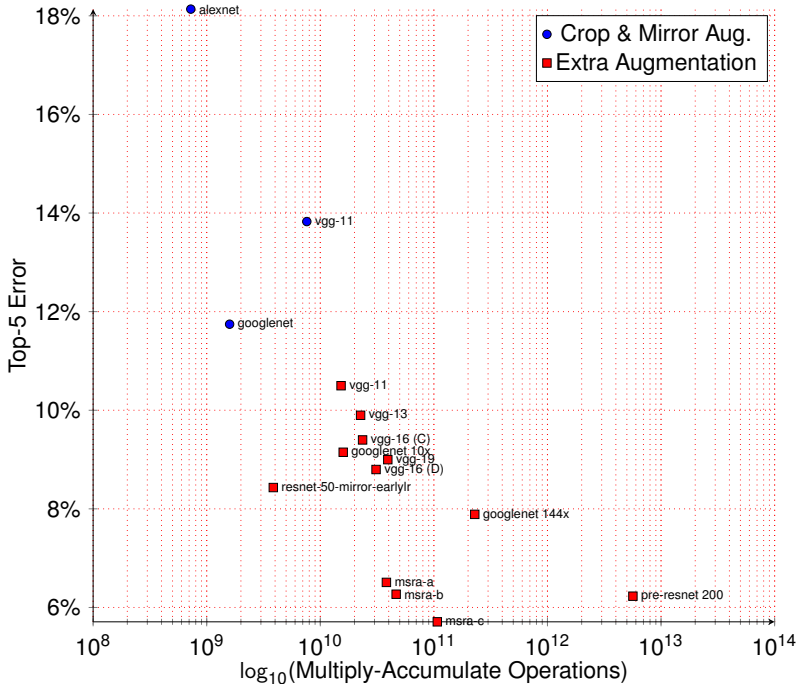
*“If you give me any sized dataset, what you ought to do **if you want good generalization** is get yourself into the small data regime. That is, however big your dataset, you ought to **make a much bigger model** so that that’s small data.”*

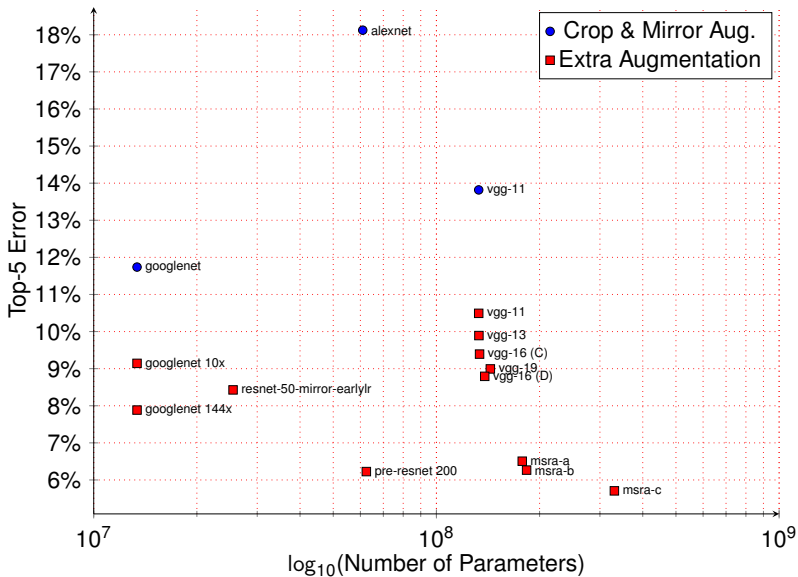
*“So, I think what the brain is doing is . . . **regularizing the hell** out of it, and that’s a better thing to do than what statisticians used to think you should do, which is have a small model.”*

Geoffery Hinton

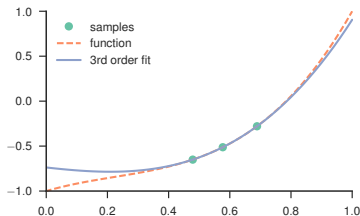
⁴Cambridge, June 2015

⁴<http://sms.cam.ac.uk/media/2017973>

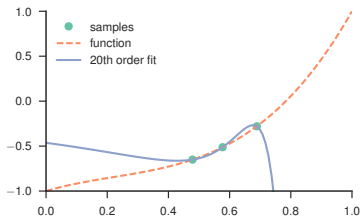




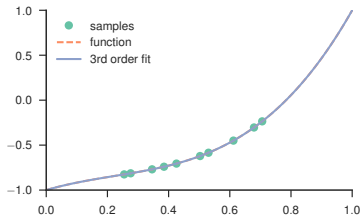
Generalization and Num. Parameters



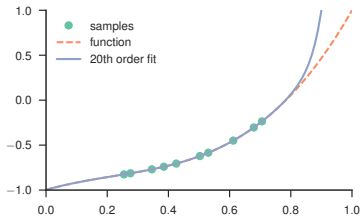
(a) 3rd-order poly., 3 points



(b) 20th-order poly., 3 points



(c) 3rd-order poly., 10 points



(d) 20th-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. Knowing only that our model is over-parameterized is a relatively weak prior, however we can encode this into the fit by using a regularization penalty. This restricts the model to effectively use only a small number of the parameters, by adding a penalty for example on the L2 norm of the model weights.
- ▶ 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.

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

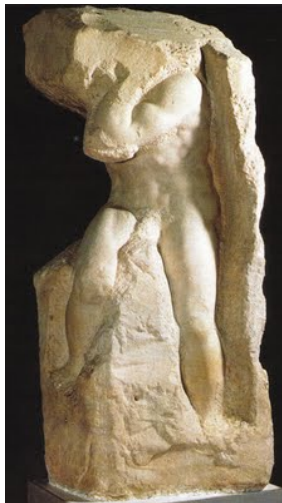
Isn't that already being addressed?

- ▶ 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.

- ▶ 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?

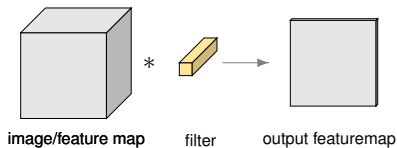
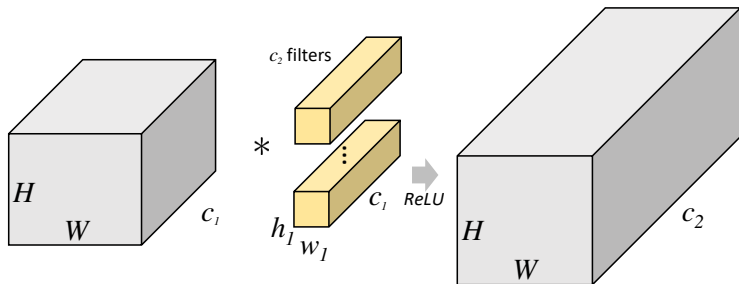


The Atlas Slave
(Accademia, Florence)

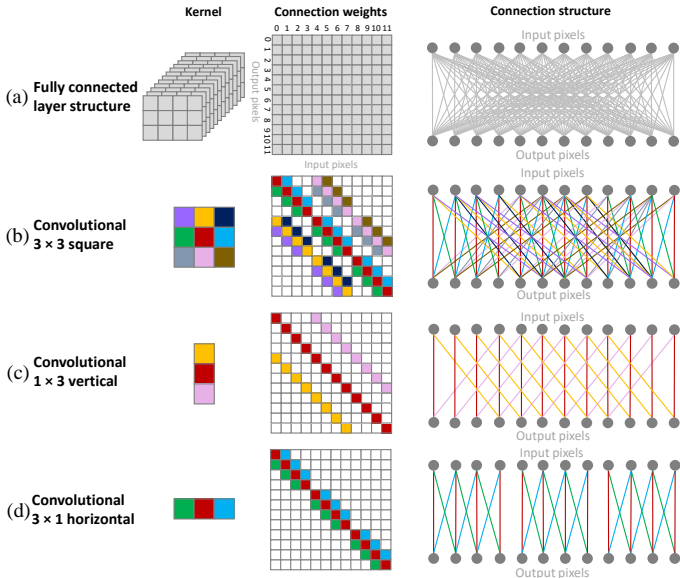
Structural Priors



Typical Convolutional Layer



Sparsity of Convolution



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?

Spatial Structural Priors

Published as a conference paper at ICLR 2016

TRAINING CNNs WITH LOW-RANK FILTERS FOR EFFICIENT IMAGE CLASSIFICATION

Yani Ioannou¹, Duncan Robertson², Jamie Shotton², Roberto Cipolla¹ & Antonio Criminisi²

¹University of Cambridge, ²Microsoft Research

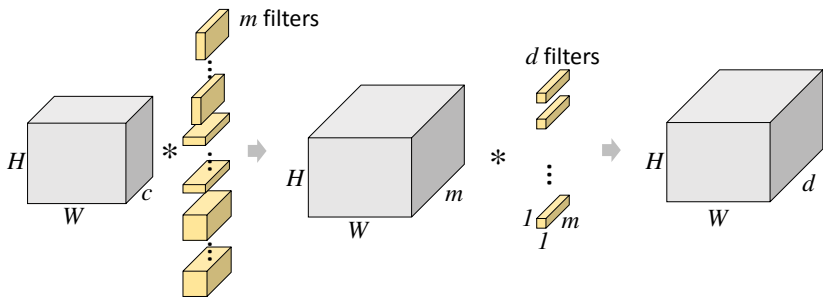
{yai20, rc10001}@cam.ac.uk, {a-durobe, jamiesho, antcrim}@microsoft.com

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

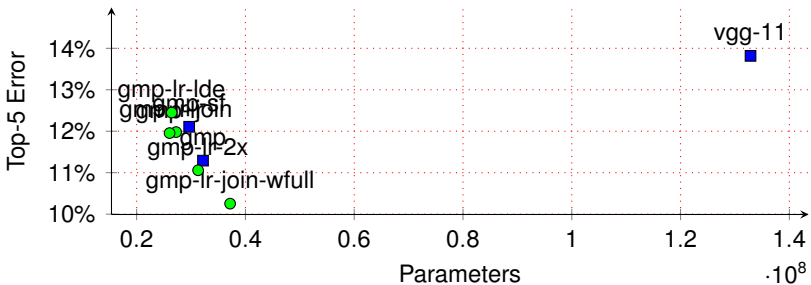
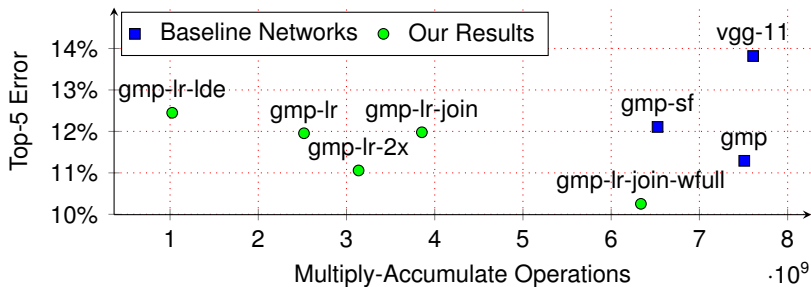
Proposed Method

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

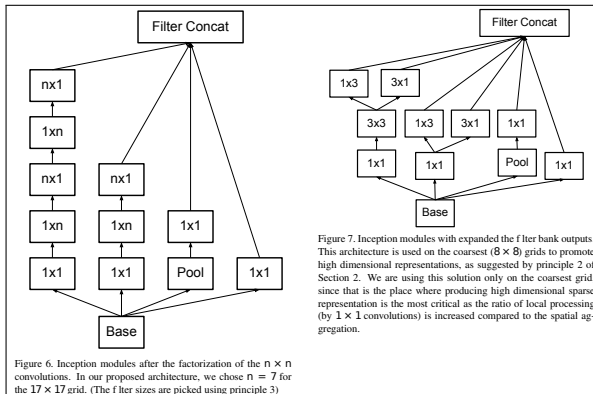


Figure 6. Inception modules after the factorization of the $n \times n$ convolutions. In our proposed architecture, we chose $n = 7$ for the 17×17 grid. (The filter sizes are picked using principle 3)

Figure 7. Inception modules with expanded filter bank outputs. This architecture is used on the coarsest (8×8) 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 1×1 convolutions) is increased compared to the spatial aggregation.

Convo
of-the-ar
tasks. Sin
to becom
ous benc
putations
for most
for traini
count are
mobile vi
ing ways
the addc
factorize
benchma.
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-

sifica-
! gains
ignifi-
mains.
p con-
perform-
increas-
Also,
appli-
where
eered,
n[4].

ure of
evalu-
On the
t [20]
t 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 Ioannou¹ Duncan Robertson² Roberto Cipolla¹
Antonio Criminisi²

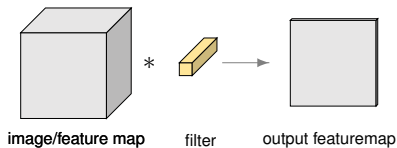
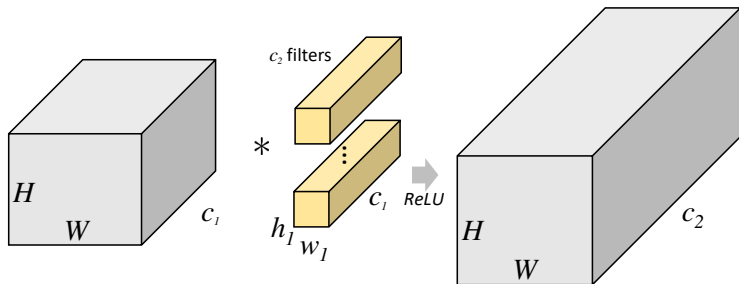
¹University of Cambridge, ²Microsoft Research

Abstract

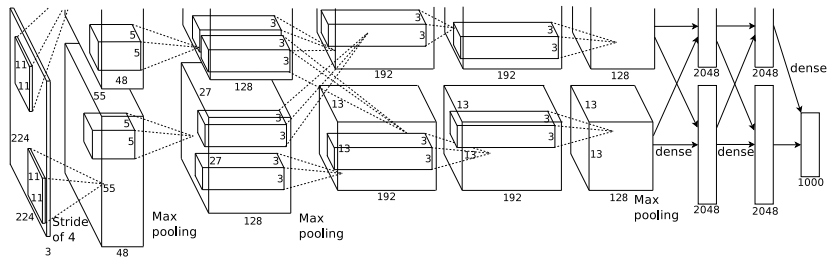
We propose a new method for creating computation-ally efficient and compact convolutional neural networks

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

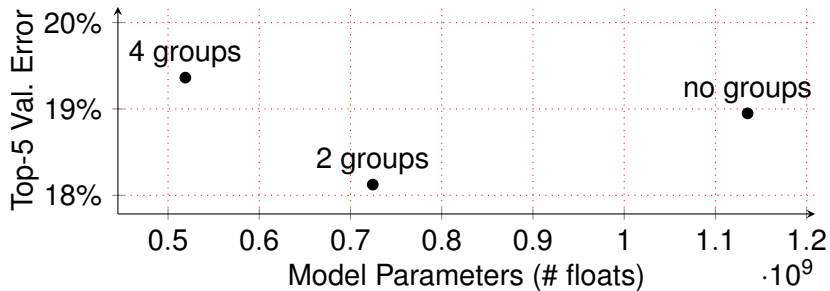
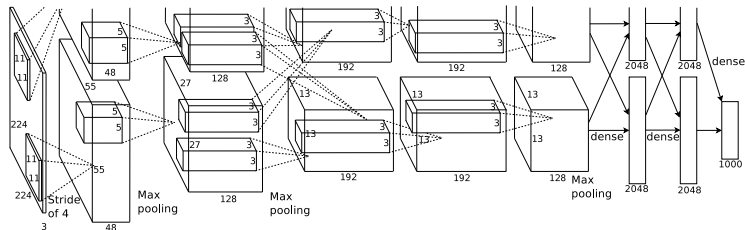


AlexNet Filter Grouping

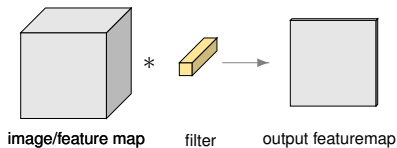
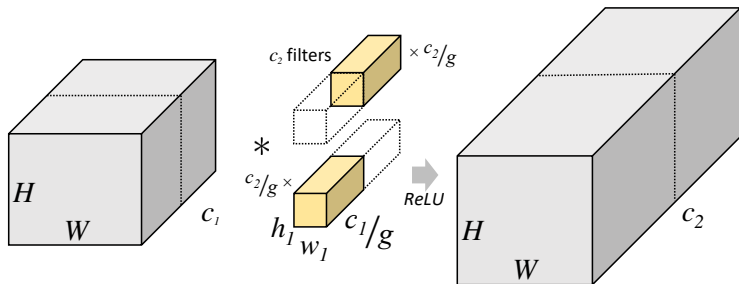


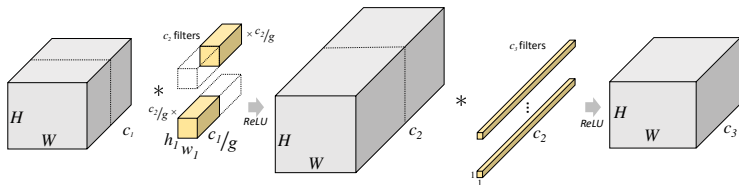
- ▶ Uses 2 filter groups in most of the convolutional layers
- ▶ Allowed training across two GPUs (model parallelism)

AlexNet Filter Grouping

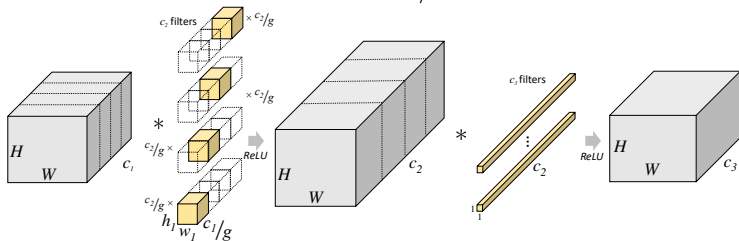


Grouped Convolutional Layer



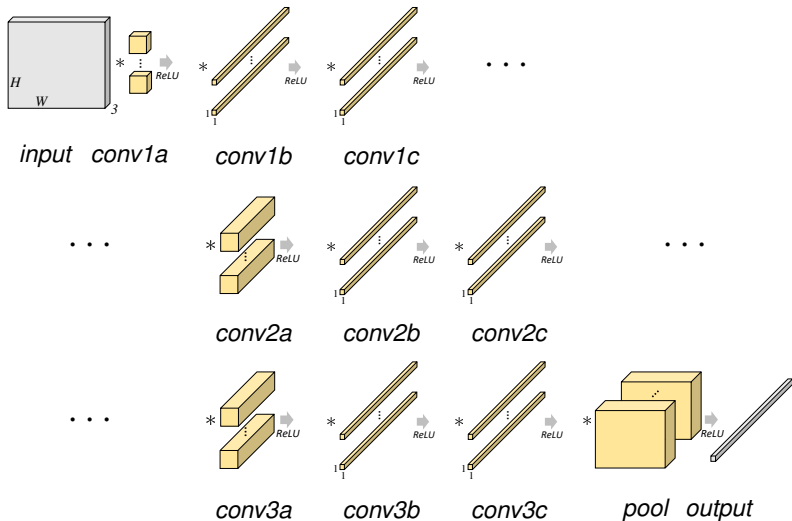


Root-2 Module: d filters in $g = 2$ filter groups, of shape $h \times w \times c/2$.

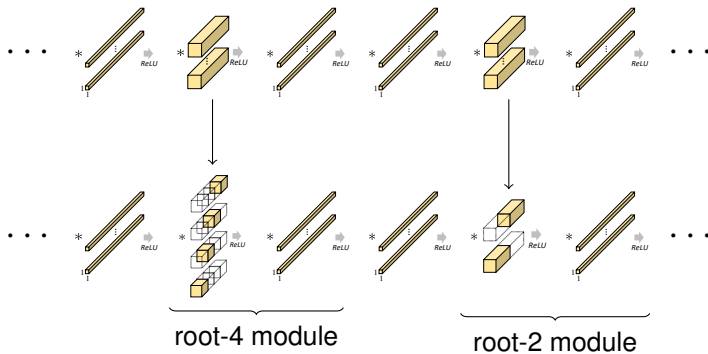


Root-4 Module: d filters in $g = 4$ filter groups, of shape $h \times w \times c/4$.

Network-in-Network



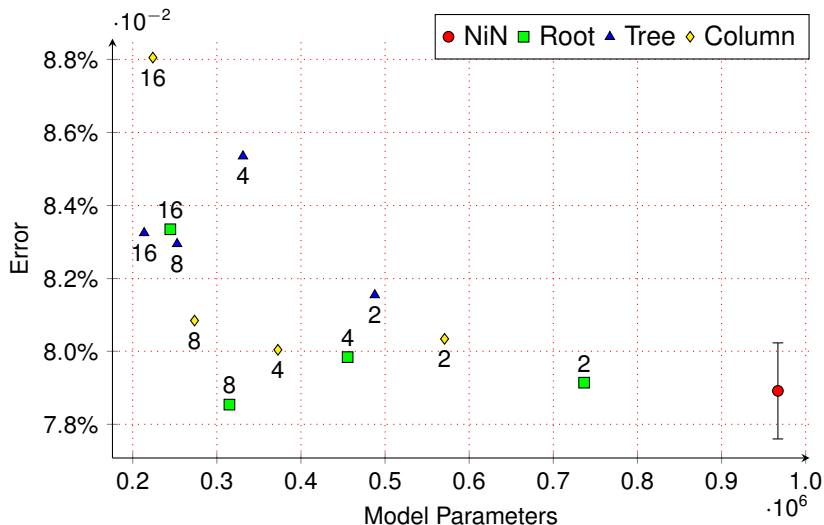
NiN Root Architectures



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

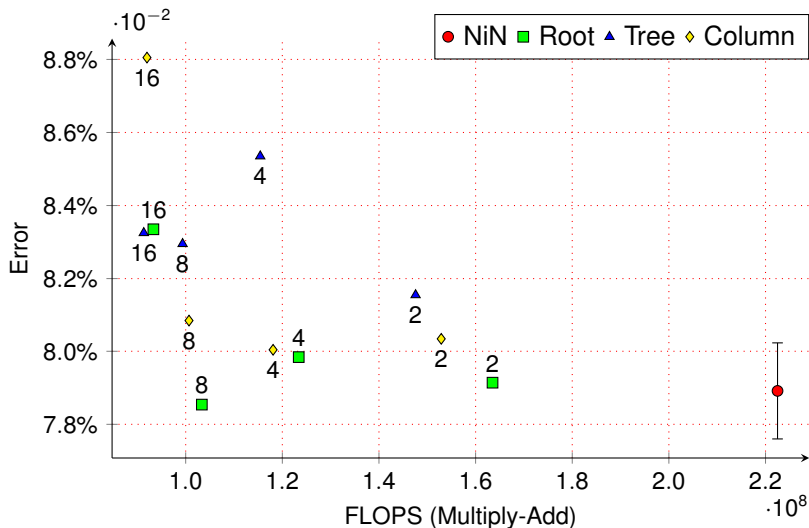
Model	conv1			conv2			conv3		
	<i>a</i>	<i>b</i>	<i>c</i>	<i>a</i>	<i>b</i>	<i>c</i>	<i>a</i>	<i>b</i>	<i>c</i>
	<i>5×5</i>	<i>1×1</i>	<i>1×1</i>	<i>5×5</i>	<i>1×1</i>	<i>1×1</i>	<i>3×3</i>	<i>1×1</i>	<i>1×1</i>
Orig.	1	1	1	1	1	1	1	1	1
root-2	1	1	1	2	1	1	1	1	1
root-4	1	1	1	4	1	1	2	1	1
root-8	1	1	1	8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

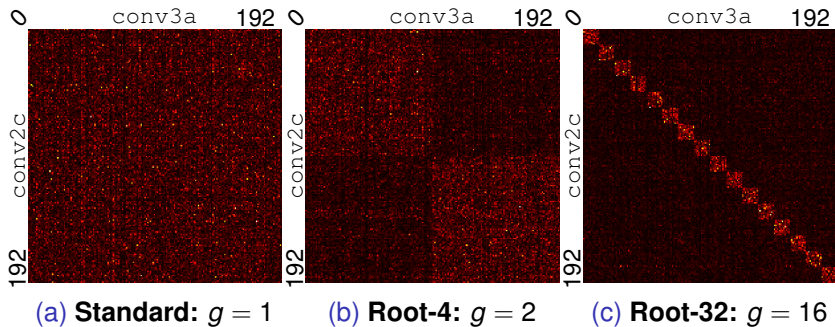
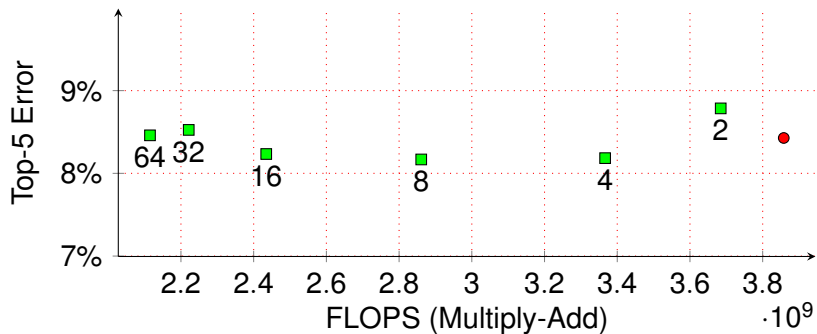
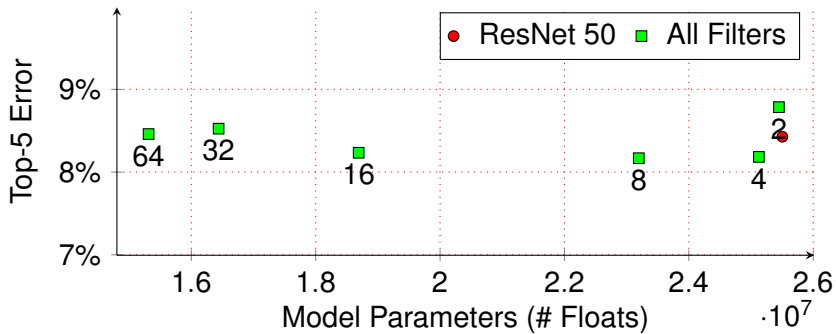
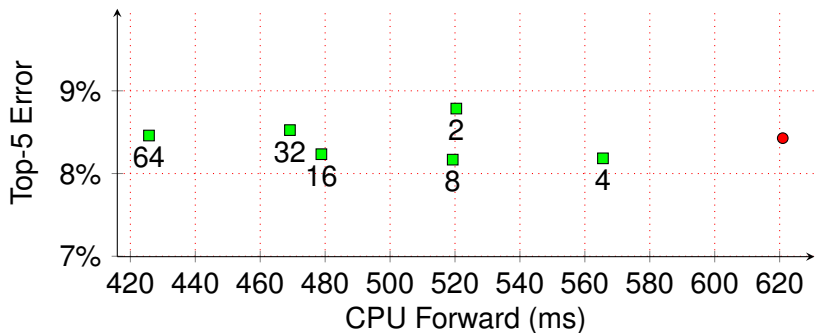
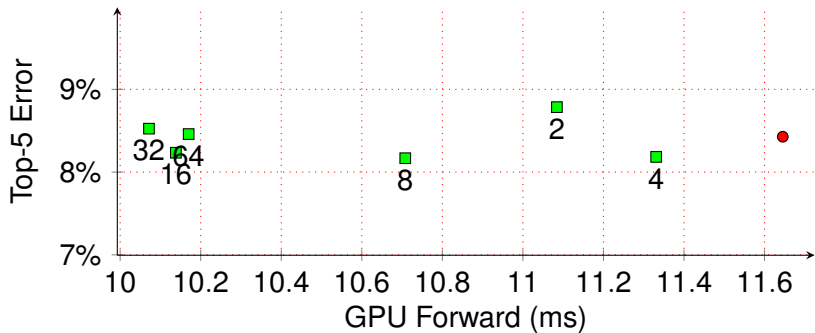


Figure: The block-diagonal sparsity learned by a root-module is visible in the correlation of filters on layers `conv3a` and `conv2c` in the NiN network.





Imagenet Results

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

- ▶ ResNet 50⁵: **40%** smaller, **45%** fewer FLOPS
- ▶ ResNet 200⁶: **44%** smaller, **25%** fewer FLOPS
- ▶ GoogLeNet: **7%** smaller, **44%** fewer FLOPS

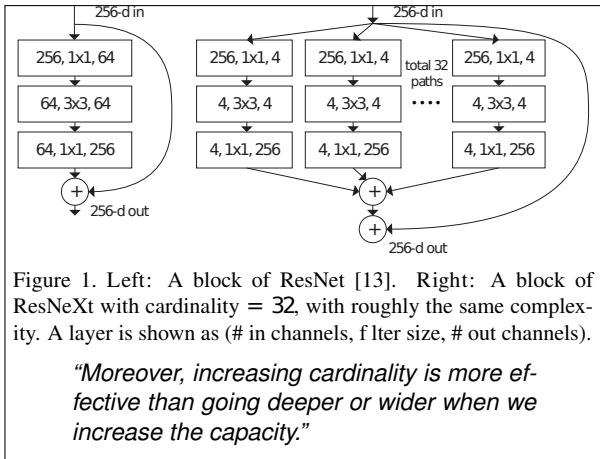
But when you also **increase the number of filters** . . .

⁵Caffe Re-implementation

⁶Based on Facebook Torch Model

Aggregated Residual Transformations for Deep Neural Networks

Saining Xie¹ Ross Girshick² Piotr Dollár² Zhuowen Tu¹ Kaiming He²
¹UC San Diego ²Facebook AI Research



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 1×1 convolutions), transformed by a set of specialized filters (3×3 , 5×5 , 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

- ▶ 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)