# **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>

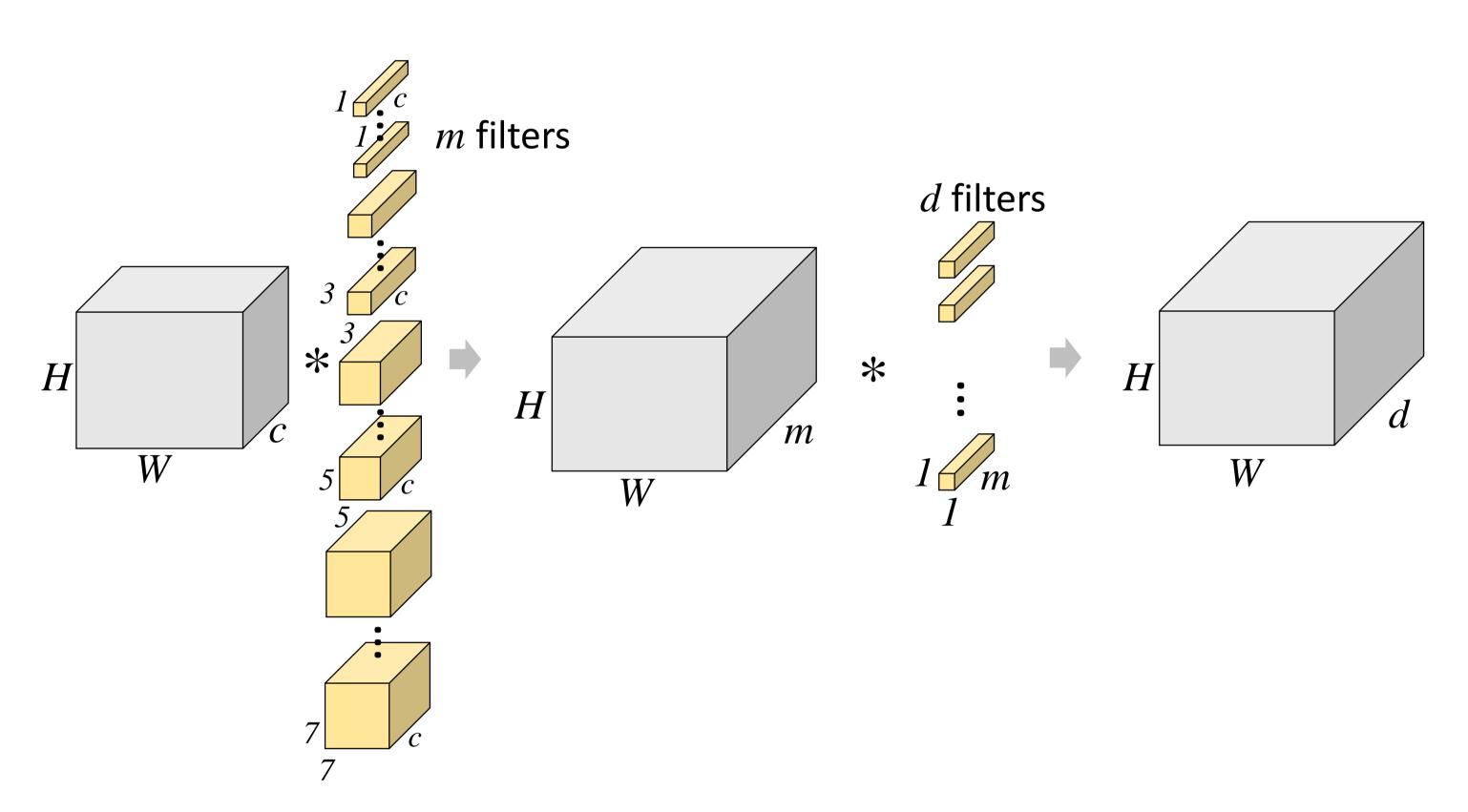
<sup>1</sup>University of Cambridge, <sup>2</sup>Microsoft Research, Cambridge

# Summary

- We train CNNs with 'root modules', a novel sparse connection structure based on filter groups.
- In effect our networks learn a *basis* for the channel extents of filters, based on compact filter groups
- s, while maintaining or aster and use f Our models are factoria increasing accuracy
- 6 fewer FLOPS, and ResNet 50, our model has 4 fewer parameters, 4 is 31% (12%) faster on a CPU (GPU)
- % fewer FLOPS ResNet 200, our model has 48% fewer parameters and 279
- ► GoogLeNet, our model has 7% fewer parameters and is 21% (16%) faster
- on a CPU (GPU)

# Previous Work: Learning a Basis for Filter Size

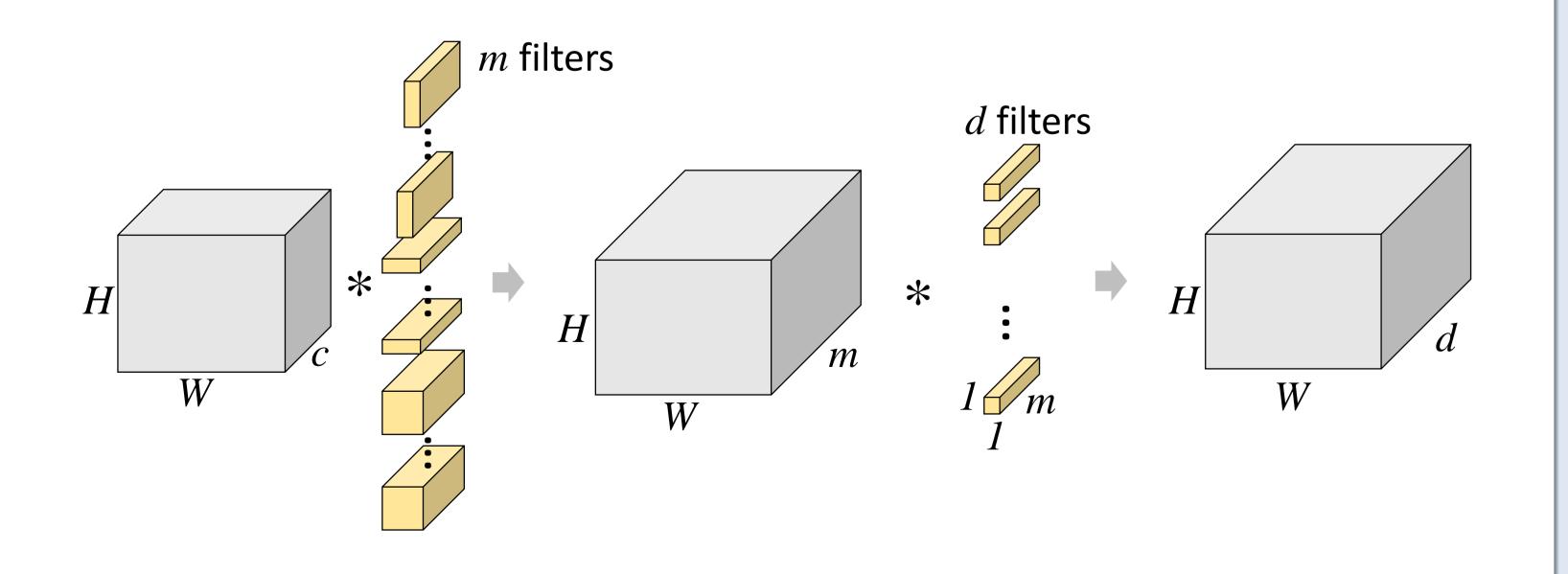
▶ In [1], linear combination of different sized filters is learned, *i.e.* a basis space for filters:



- Learns a limited number of large filters (7 $\times$ 7, 5 $\times$ 5), and a great number of small filters  $(3 \times 3, 1 \times 1)$
- Motivation: expect most image correlations to be highly localized, *i.e.* many small filters. However, a few may require larger, more complex filters
- We can learn an effective filter of full size, but limited parameterization, by learning a basis – the  $1 \times 1$  filters can learn a linear combination of the basis

# Previous Work: Learning a Low-Rank Basis for Filters

▶ In [2], linear combination of low-rank filters is learned, *i.e.* a basis space for



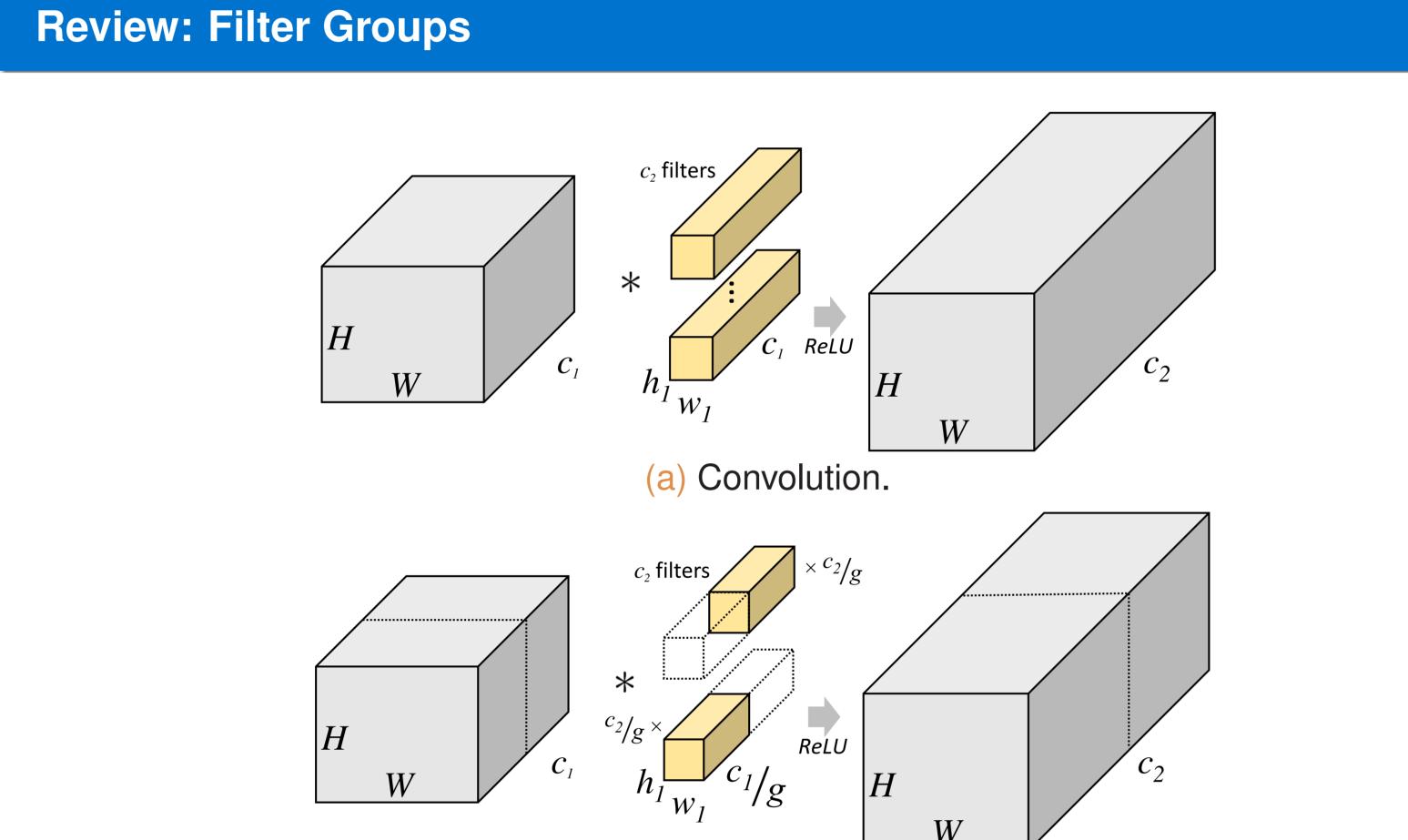
- A set of filters of different shape (similar to 'Inception', but low-rank and of different orientation [1])
- On the following layer, use  $d \times [1 \times 1 \times m]$  filters to linearly combine
- Only learning a low-rank filter in the spatial dimensions filters maintain full channel depth

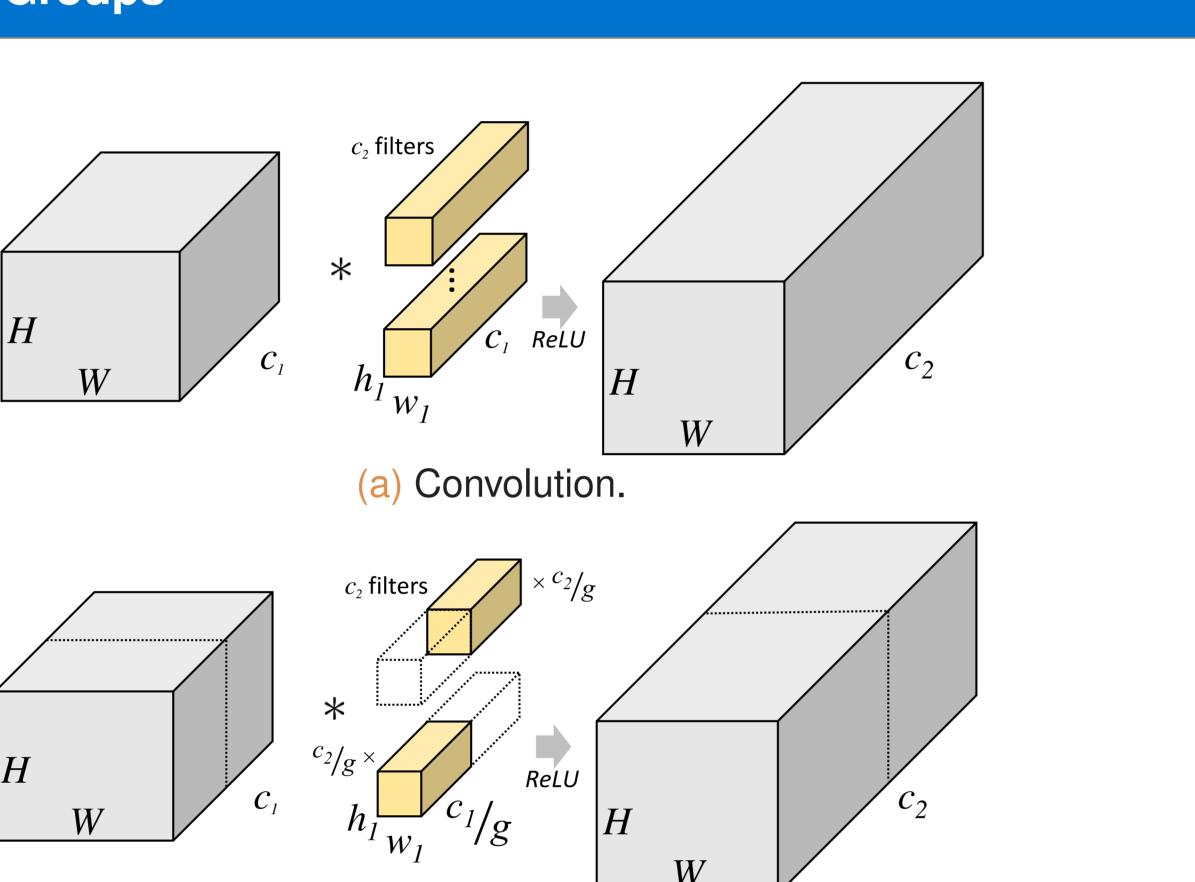


- These previous works have learned a more compact representation by learning a basis, informed by our prior knowledge of the task However, they have only looked at the spatial extents of a convolutional filter –
- what about the 'channel' extents?

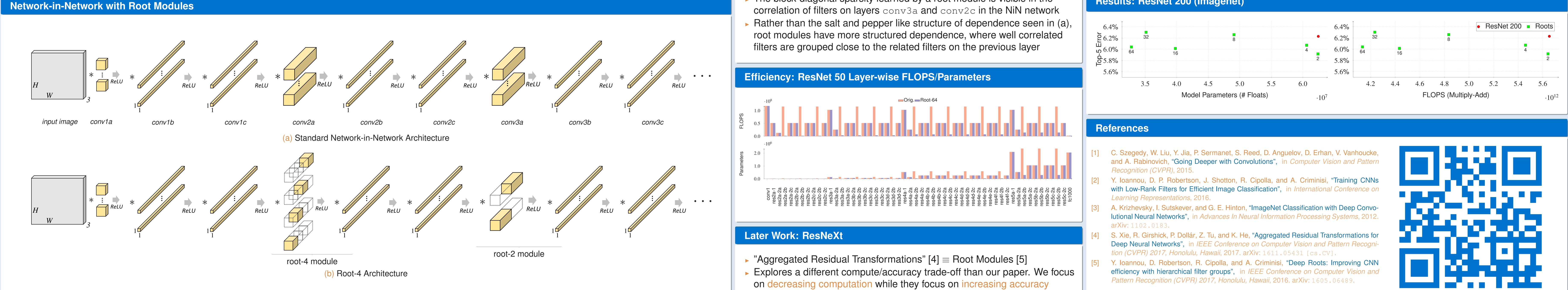
Idea

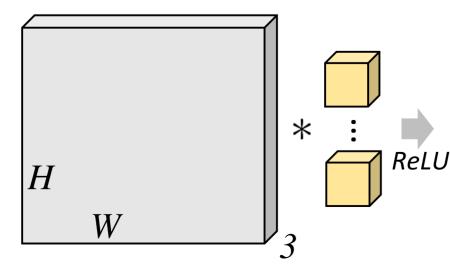
- In a CNN, each filter learns a connection to e





- Convolutional filters (yellow) typically have the same channel dimension  $(c_1)$  as the input feature maps (gray)
- colors [3], g independent groups of  $c_2/g$  filters operate With convolutional filter groups on a fraction  $c_1/g$  of the input feature map channels
- This change does not affect the dimensions of the input and output feature maps but significantly reduces computational complexity and the number of model parameters.



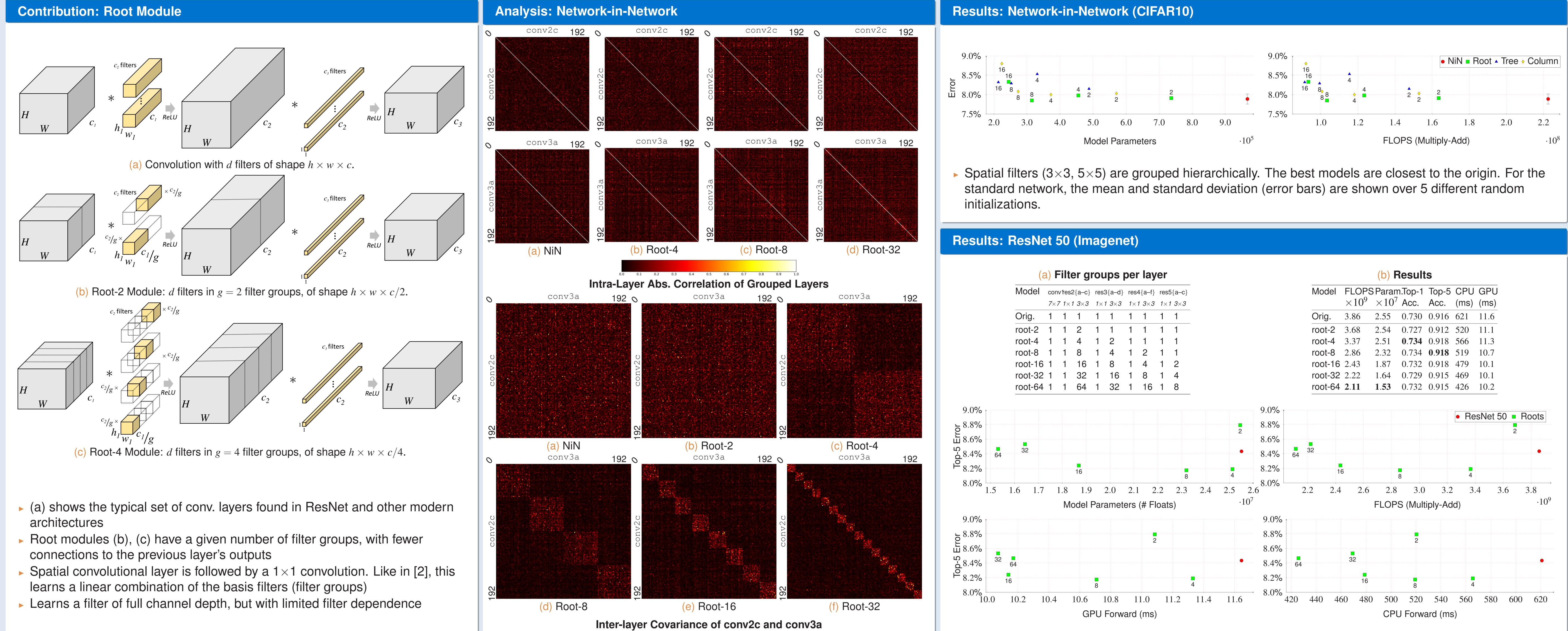


# Motivation: Learning a Basis for Inter-Layer Filter Dependence

r on the previous layer However, it has been shown that the learned weights in CNNs are sparse In a DPM/hierachical representation parts have sparse relationships Can we learn a more compact representation exploiting this sparsity?

) Convolution with filter groups.

Finite the term of term o



# Nicrosoft

- The block-diagonal sparsity learned by a root module is visible in the

INERS ( NOF

# Results: ResNet 200 (Imagenet)