

# Exploring Complexity Reduction in Deep Learning

Sourya Dey, Peter Beerel, Keith Chugg

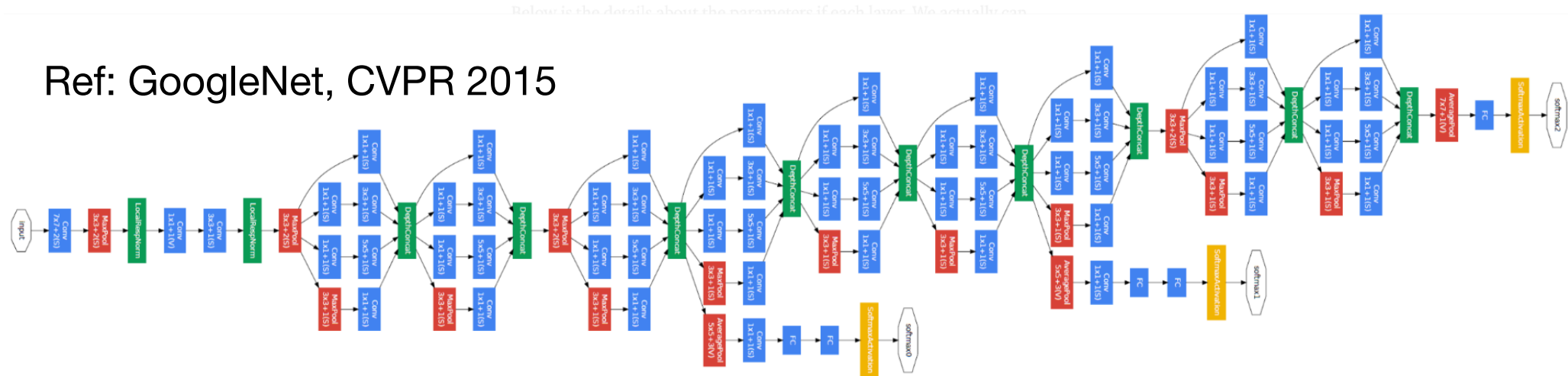
## (1) Problem Statement

Neural networks *need a lot of manual tuning*

- Architecture, layers (discrete)
- Hyperparameter values (continuous)

Neural networks have *massive complexity*

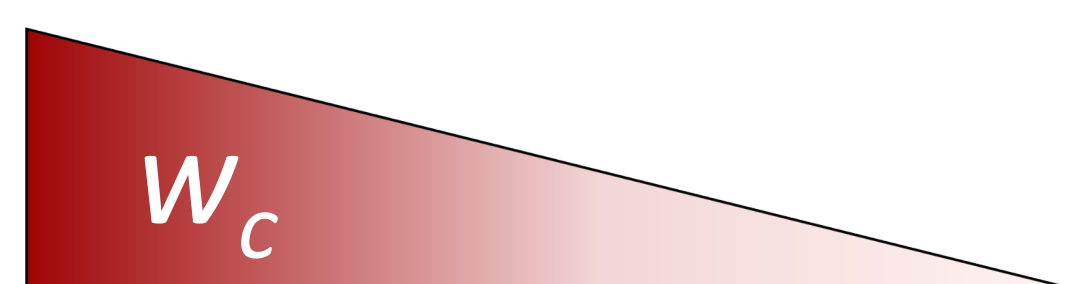
Ref: GoogleNet, CVPR 2015



Our research goal: **Automate the search for low complexity networks which give good performance**

Optimization objective:

$$f = f_p(\text{Performance}) + w_c * f_c(\text{Complexity})$$



Current focus: CNNs

Quick to train  
Bad performance

Good performance  
Too long to train

## (2) Approaches

Search space is both continuous and discrete

Each point  $x$  is a neural network to be trained

*Evaluating  $f$  is expensive and noisy!*

Potential approaches

- Simulated annealing
- **Bayesian optimization**
- Evolutionary / genetic algorithms

Sample  $f(\cdot)$  and model via a **Gaussian process**

$$f(\mathbf{X}_{1:n}) \sim \mathcal{N}\left(\boldsymbol{\mu}, \boldsymbol{\Sigma}\right) \quad \boldsymbol{\Sigma} = \begin{bmatrix} k(x_1, x_1) & \cdots & k(x_1, x_n) \\ \vdots & \text{Covariance} & \vdots \\ k(x_n, x_1) & \cdots & k(x_n, x_n) \end{bmatrix}$$

kernel  $k$

Get potential new networks via **expected improvement**

- Expensive  $f$  evaluations are minimized
- Kernel can model noise

$$EI(\mathbf{x}) = (f^* - \mu)P\left(\frac{f^* - \mu}{\sigma}\right) + \sigma p\left(\frac{f^* - \mu}{\sigma}\right)$$

$f^*$  = Current optimal value

## (3) Research Methodology

Bayesian optimization can fail if search space is too big

Given a problem, **divide search space into levels:**

- |  |   |   |
|--|---|---|
| <ul style="list-style-type: none"> <li>- # convolutional layers</li> <li>- # channels in each</li> <li>- Downsampling (strides/pooling)</li> </ul> | } | <p><b>Level 1:</b><br/>Basic structure</p>          |
| <ul style="list-style-type: none"> <li>- Kernel sizes</li> <li>- Batch normalization (yes/no)</li> <li>- Grouped convolutions</li> </ul>           | } | <p><b>Level 2:</b><br/>Parameter adjustments</p>    |
| <ul style="list-style-type: none"> <li>- # classification layers</li> <li>- Densities</li> <li>- Weight decay coefficients</li> </ul>              | } | <p><b>Level 3:</b><br/>Classifier</p>               |
| <ul style="list-style-type: none"> <li>- Learning rate</li> <li>- Learning rate decay</li> <li>- Batch size</li> </ul>                             | } | <p><b>Level 4:</b><br/>Training hyperparameters</p> |

Our prior work on pre-defined sparsity

S. Dey, K. Huang, P. A. Beerel and K. M. Chugg, "Pre-Defined Sparse Neural Networks with Hardware Acceleration," in *IEEE JETCAS*, vol. 9, no. 2, pp. 332-345, June 2019.

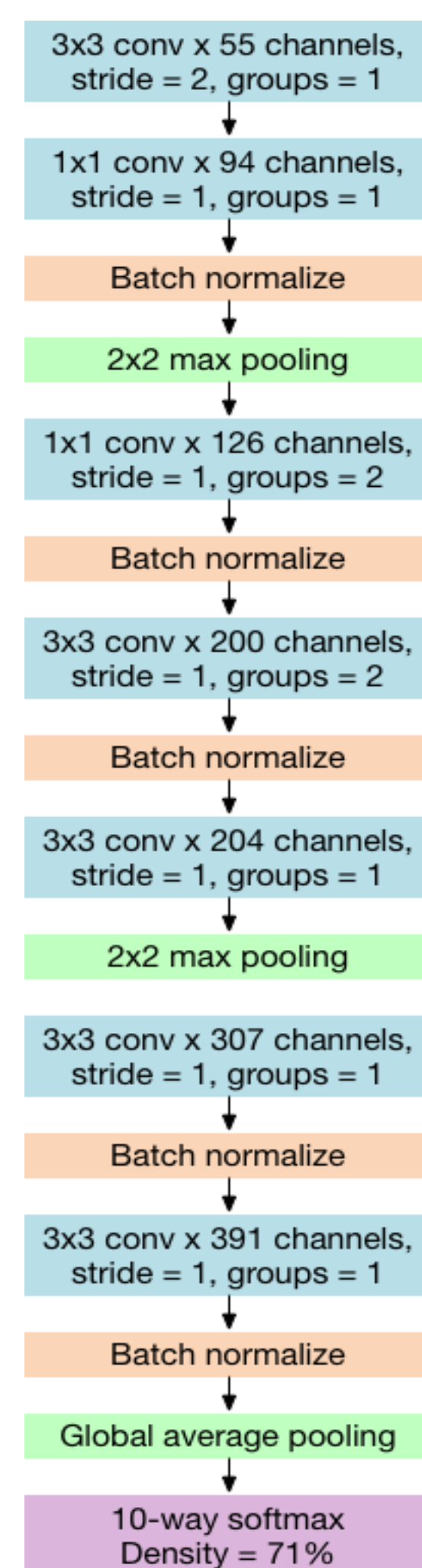
$f_p(\text{Performance}) = 1$  - Best validation acc

$f_c(\text{Complexity}) = \text{Normalized training time per epoch}$

Can set  $w_c$  according to desired tradeoff

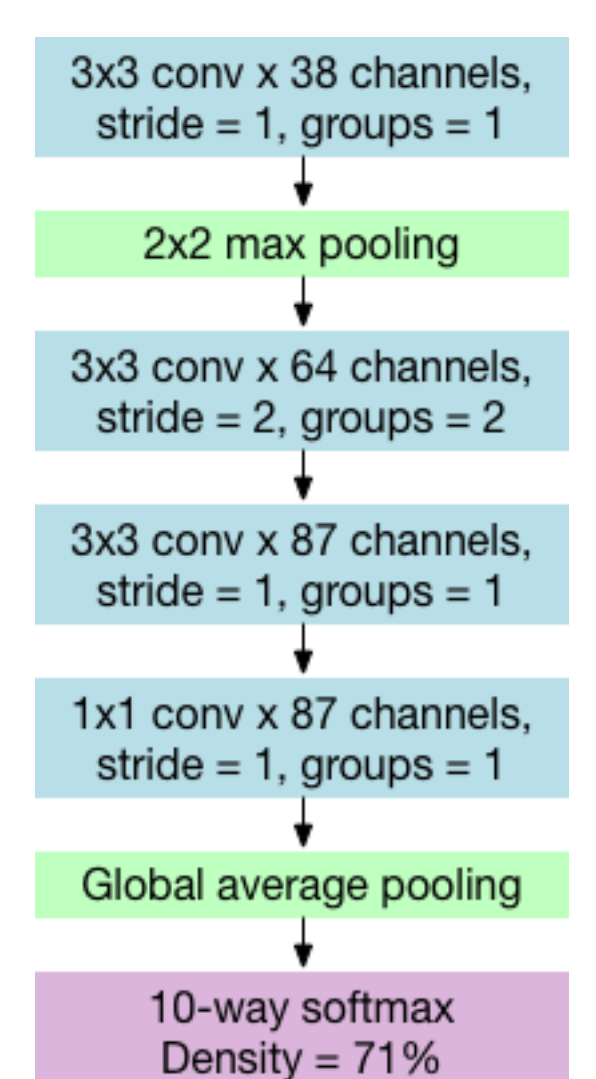
## (4) Results so far

$w_c = 0.1$  (Balanced case)



Learning rate = 4.4e-3, Decay = 0.992  
Weight decay = 2.3e-3, Batch size = 338  
Best val acc = 82% in 30 eps

$w_c = 1$   
(More focus on low complexity)



Learning rate = 4.3e-3, Decay = 0.999  
Weight decay = 4e-4, Batch size = 501  
Best val acc = 74% in 30 eps

Dataset used:  
CIFAR-10 images  
of size 32x32 x 3  
channels (no aug)

Contact Information: [souryade@usc.edu](mailto:souryade@usc.edu), [pabeerel@usc.edu](mailto:pabeerel@usc.edu), [chugg@usc.edu](mailto:chugg@usc.edu)

This work is partly supported by National Science Foundation, USA, Grant #1763747, and by Defense Threat Reduction Agency, USA.

