

Exploring Complexity Reduction in Deep Learning

Sourya Dey, Peter Beerel, Keith Chugg



USC
Viterbi
School of Engineering
Ming Hsieh Department
of Electrical and
Computer Engineering

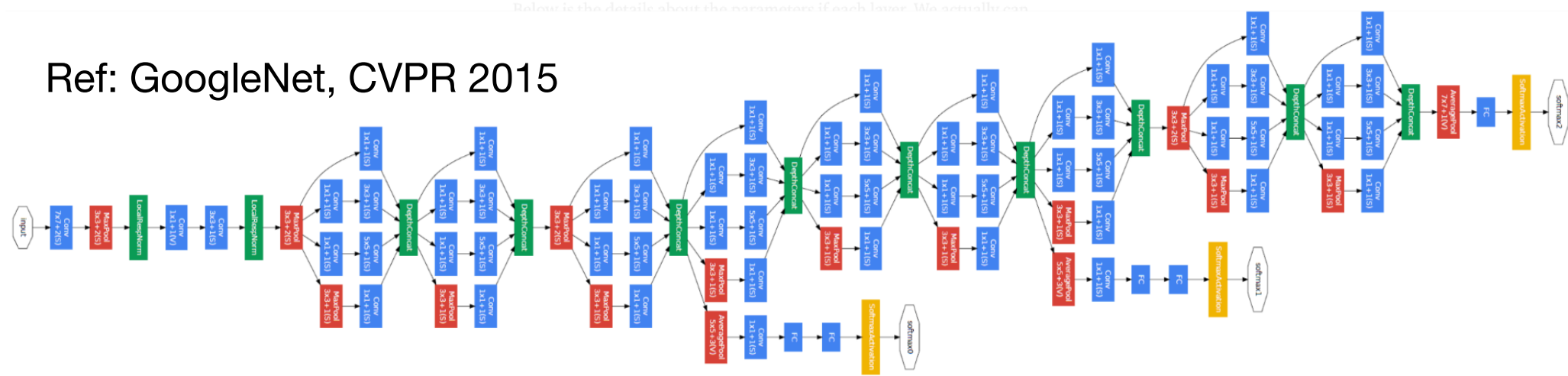
(1) Problem Statement

Neural networks *need a lot of manual tuning*

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

Neural networks have *massive complexity*

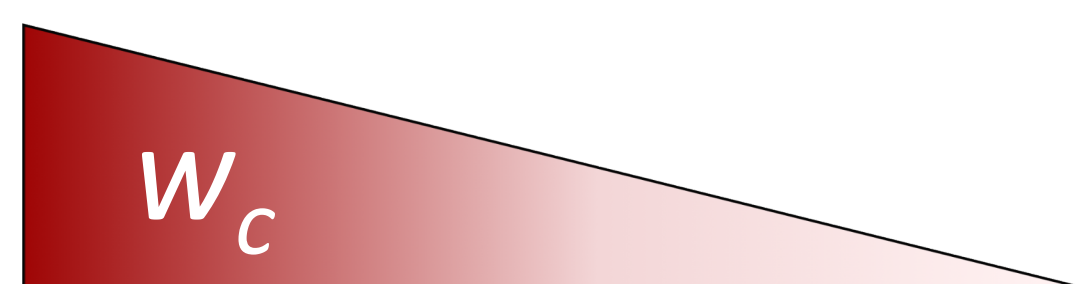
Ref: GoogleNet, CVPR 2015



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

Optimization objective:

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



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**

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

Get potential new networks via **expected improvement**

- Expensive f evaluations are minimized
- Kernel can model noise

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

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*:

- | | |
|----------------------------------|--|
| - # convolutional layers | } Level 1:
Basic structure |
| - # channels in each | |
| - Downsampling (strides/pooling) | |
| - Kernel sizes | } Level 2:
Parameter adjustments |
| - Batch normalization (yes/no) | |
| - Grouped convolutions | |
| - # classification layers | } Level 3:
Classifier |
| - Densities | |
| - Weight decay coefficients | |
| - Learning rate | } Level 4:
Training hyperparameters |
| - Learning rate decay | |
| - Batch size | |

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.

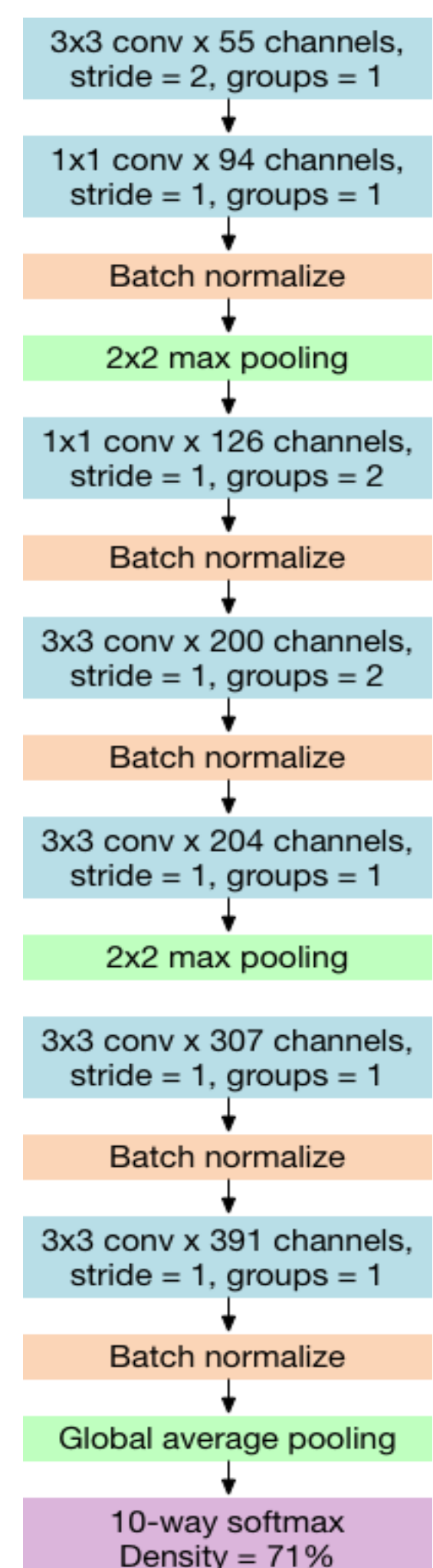
$f_p(\text{Performance}) = 1$ - Best validation acc

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

Can set w_c according to desired tradeoff

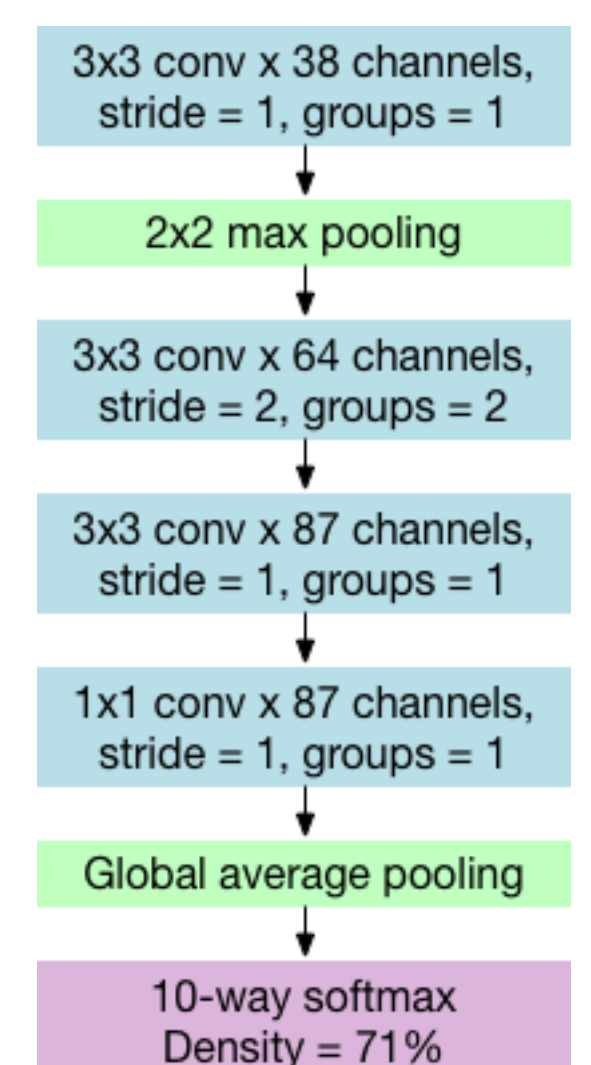
(4) Results so far

$w_c = 0.1$ (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)