

Less is More: Towards Compact CNNs

Hao Zhou¹, Jose M. Alvarez² and Fatih Porikli^{2,3}

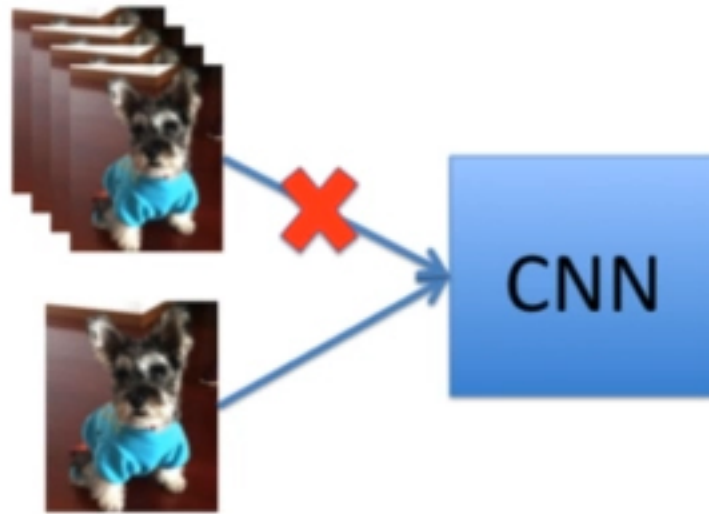
¹University of Maryland, College Park, USA

²Data61/CSIRO, Canberra, Australia

³Australian National University, Canberra, Australia

Motivation

1. CNNs are very large (Millions of parameters)
2. Large memory footprint



Motivation

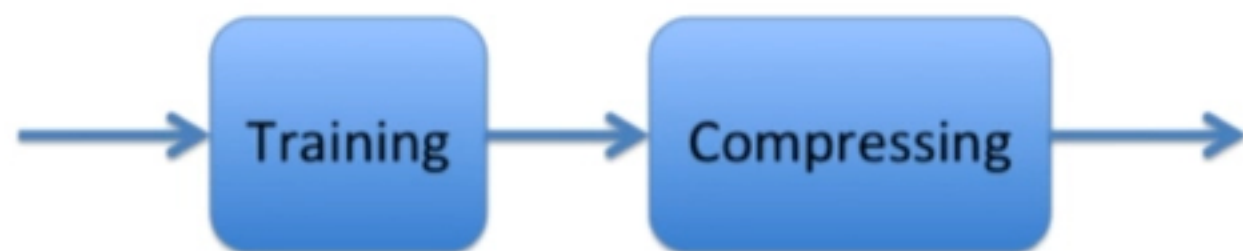
What we did

AlexNet: 60M  14M

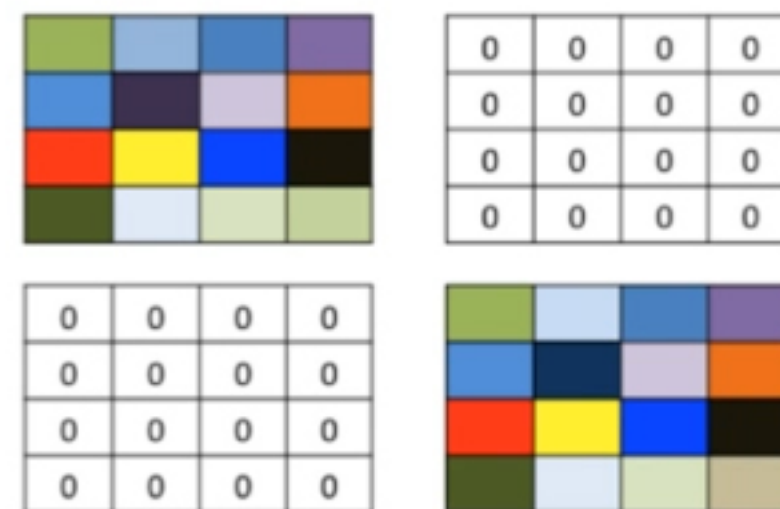
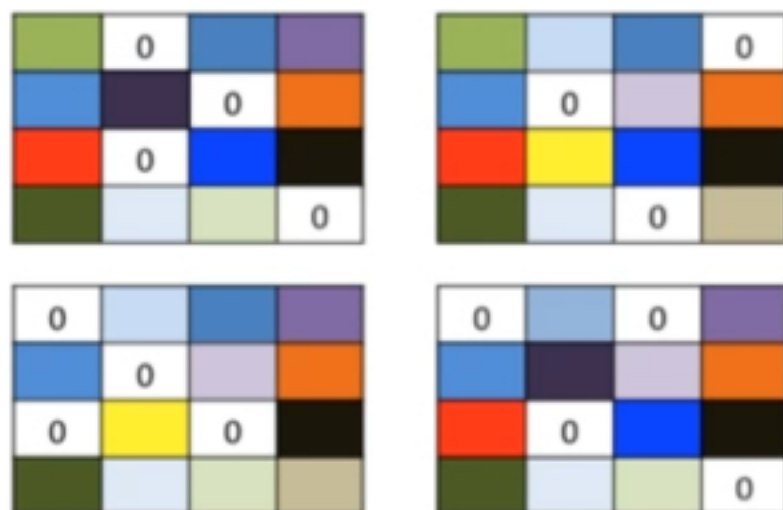
VGG: 133M  74M

Contributions

Others



Ours



Our method

Idea: adding sparse constraints to neurons.

$$\min_{\hat{\mathbf{W}}} \psi(\hat{\mathbf{W}}) + g(\hat{\mathbf{W}})$$

The diagram shows the optimization objective $\min_{\hat{\mathbf{W}}} \psi(\hat{\mathbf{W}}) + g(\hat{\mathbf{W}})$. Below the term $\psi(\hat{\mathbf{W}})$ is a box labeled "Loss for CNNs" with an arrow pointing to $\hat{\mathbf{W}}$ in the term. Below the term $g(\hat{\mathbf{W}})$ is a box labeled "Sparse Constraints" with an arrow pointing to $\hat{\mathbf{W}}$ in the term.

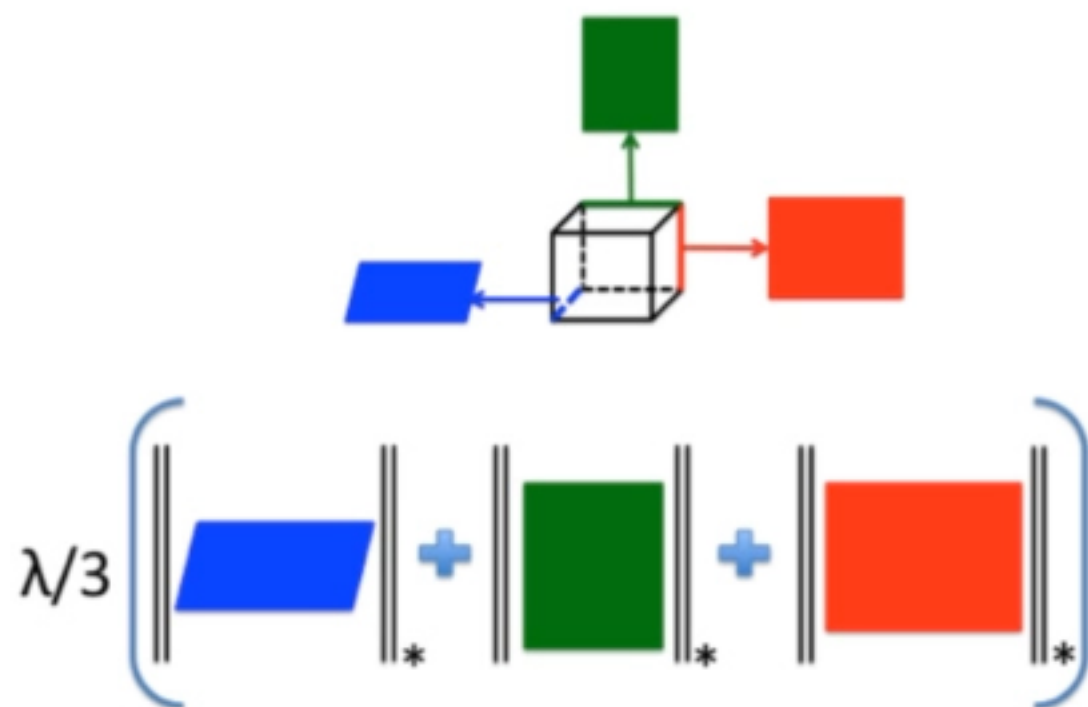
Forward-backward splitting:

→ Forward: Backprop $\hat{\mathbf{W}}^* \leftarrow \hat{\mathbf{W}} - \tau \frac{\psi(\hat{\mathbf{W}})}{\hat{\mathbf{W}}}$

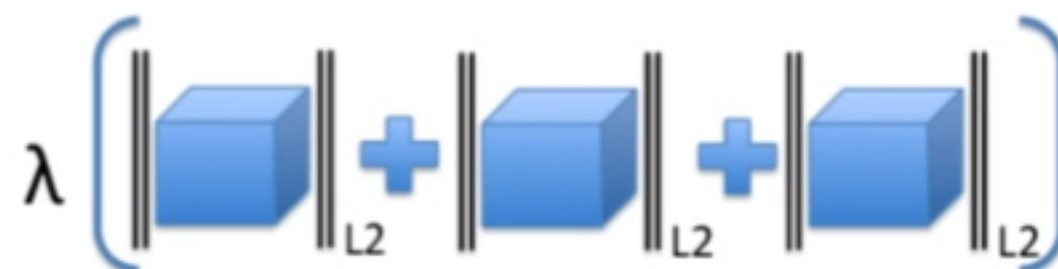
→ Backward: Sparsity $\hat{\mathbf{W}} \leftarrow \arg \min_{\hat{\mathbf{W}}} g(\hat{\mathbf{W}}) + \frac{1}{2\tau} \|\hat{\mathbf{W}} - \hat{\mathbf{W}}^*\|^2$

Our method — sparse constraints

Sparse Constraints

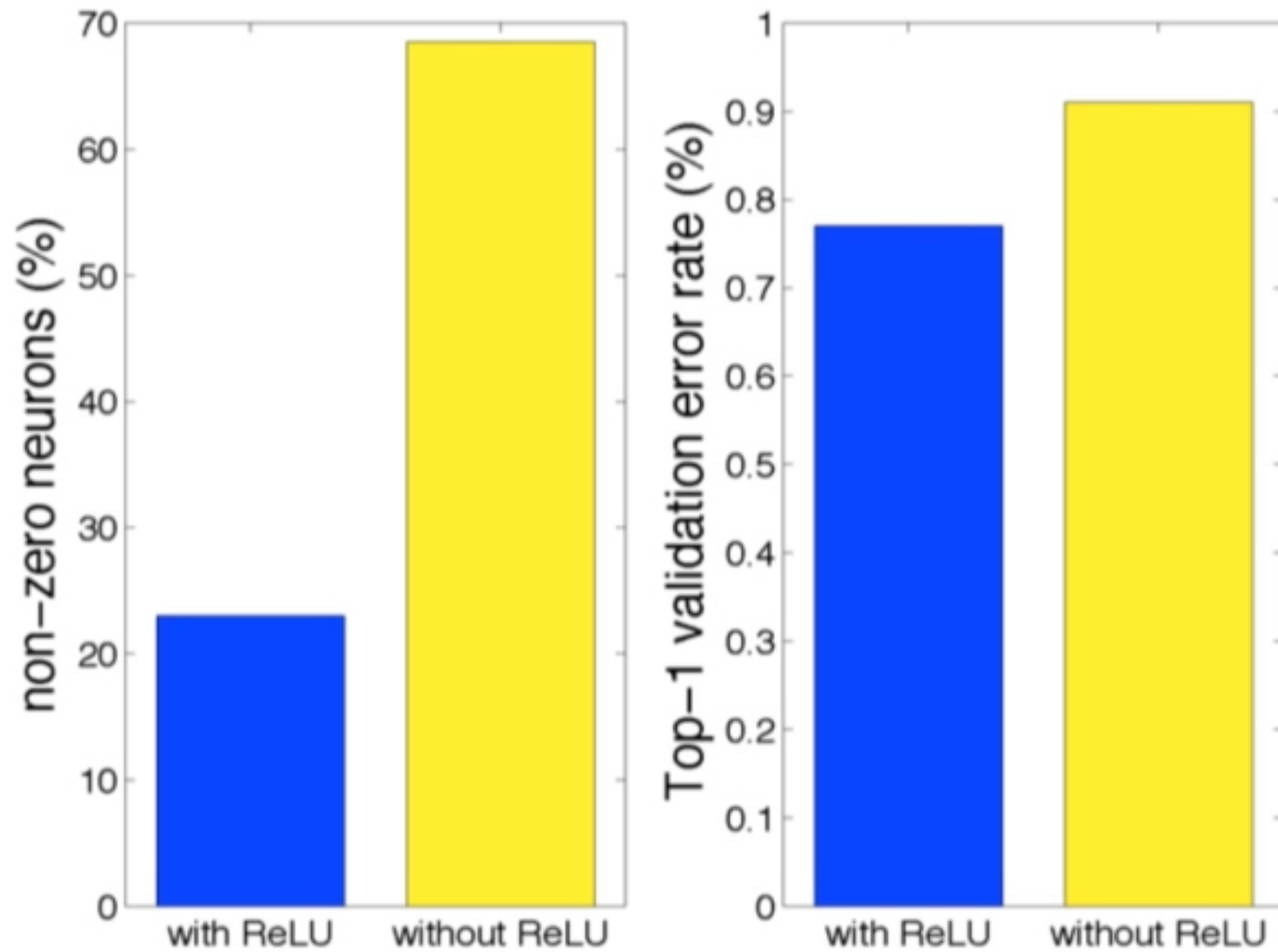


Tensor Low Rank



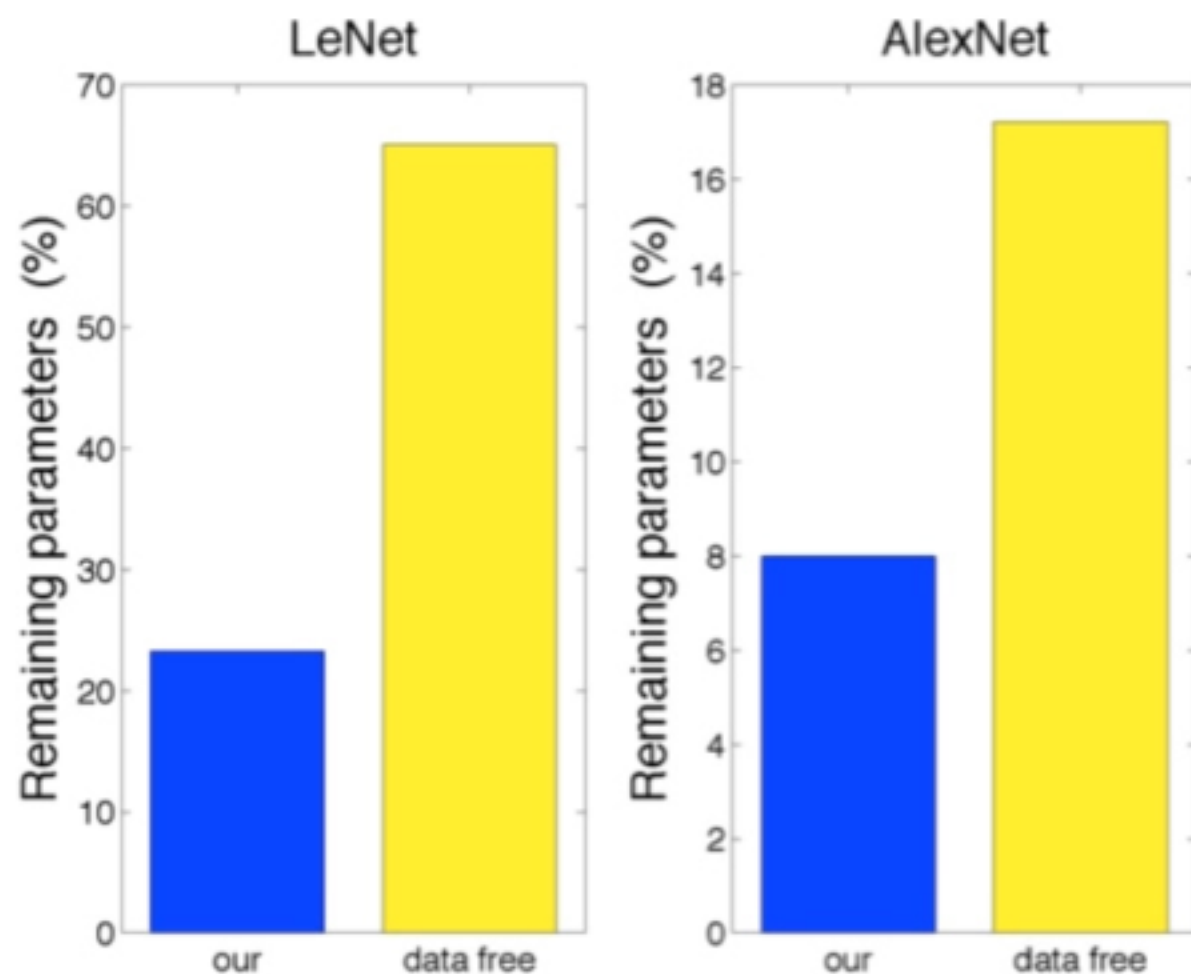
Group Sparsity

Experiments — ReLU

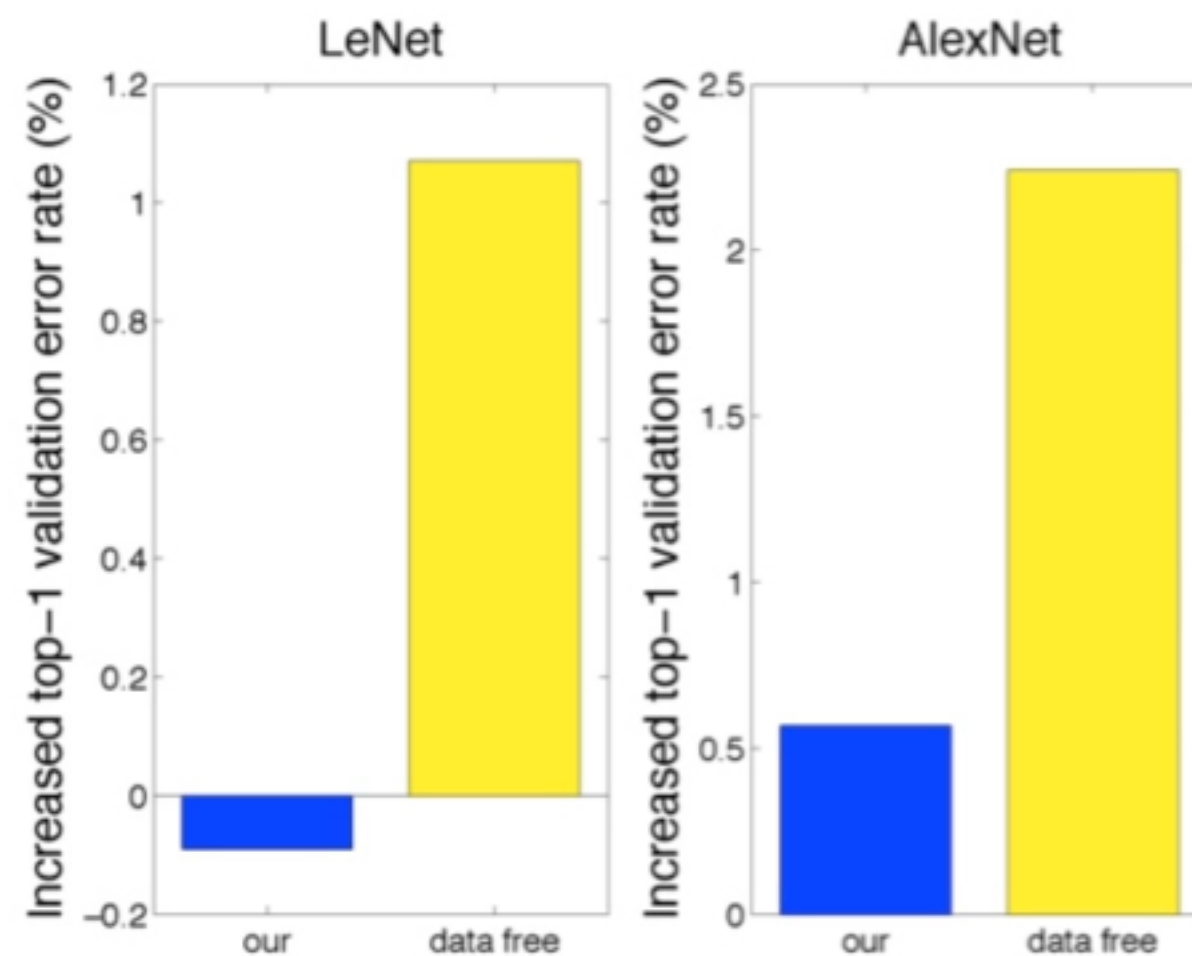


Conv2 on LeNet

Experiments



Non-zero parameters



Increased error rate

Comments?

Questions?

Welcome to poster

#09