

## What is the State of Neural Network Pruning? Davis Blalock\*, Jose J. González\*, Jonathan Frankle, John Guttag

- made it impossible to tell what works best





- efficiency
- pretained model



- Pruning works!
- Measurable efficiency gains
- Better than random
- But not as helpful as using a better architecture



This research was funded by two Qualcomm Innovation Fellowships, a "La Caixa" Fellowship, Wistron Corporation, and Quanta Computer

## Lack of Experimental Standardization



### Reality

- Incomplete curves, or single points
- No clear progress over time



### **But Why?**

Number of Papers a Given Paper Compares T

Butuset		
ImageNet	VGG-16	22
CIFAR-10	ResNet-56	14
ImageNet	ResNet-50	14
ImageNet	CaffeNet 11	
ImageNet	AlexNet	9
CIFAR-10	CIFAR-VGG	8
ImageNet	ResNet-34 6	
ImageNet	ResNet-18	6
CIFAR-10	ResNet-110	5

# pruning algorithm

Model (+				
-2.1	4.6	0.8		
0.2	1.5	-4.9		
-2.5	2.7	4.2		
-0.3	5.0	3.1		







## ShrinkBench

Standardized datasets, architectures, preprocessing, finetuning, evaluation, etc

Flexible mask-based API to allow arbitrary sparsity and



## **Additional Pitfalls Uncovered**

Initial weights are an important confounder

