

## **Efficient Neural Network Compression**

Namhoon Lee

University of Oxford

3 May 2019

### A Challenge in Deep Learning: *Overparameterization*

Large neural networks require:



memory & computations



power consumption

### A Challenge in Deep Learning: *Overparameterization*

Large neural networks require:

Critical to resource constrained environments



memory & computations



power consumption



embedded systems e.g., mobile devices



real-time tasks e.g., autonomous car

#### Network compression

The goal is to reduce the *size* of neural network <u>without compromising accuracy</u>.



- Network pruning
  - : reduce the number of parameters

- Network pruning
  - : reduce the number of parameters
- Network quantization
  - : reduce the precision of parameters

- Network pruning
  - : reduce the number of parameters
- Network quantization
  - : reduce the precision of parameters

Others: knowledge distillation, conditional computation, etc.

- Network pruning
  - : reduce the number of parameters
- Network quantization
  - : reduce the precision of parameters

Others: knowledge distillation, conditional computation, etc.

### Network pruning

**Different forms** 

- Parameters (weights, biases)
- Activations (neurons)

can be done structured way (e.g., channel, filter, layer)

## Network pruning

**Different forms** 

- Parameters (weights, biases)
- Activations (neurons)

can be done structured way (e.g., channel, filter, layer)

#### **Different principles**

- Magnitude based
- Hessian based
- Bayesian

## Network pruning

**Different forms** 

- Parameters (weights, biases)
- Activations (neurons)

can be done structured way (e.g., channel, filter, layer)

#### **Different principles**

- Magnitude based
- Hessian based
- Bayesian

 $\Rightarrow$  remove > 90% parameters

• Hyperparameters with weakly grounded heuristics

(e.g., layer-wise threshold [5], stochastic pruning rule [2])

- [1] Learning both weights and connections for efficient neural network, Han et al. NIPS'15
- [2] Dynamic network surgery for efficient dnns, Guo et al. NIPS'16.
- [3] Learning-compression algorithms for neural net pruning, Carreira-Perpinan & Idelbayev. CVPR'18.
- [4] Variational dropout sparsifies deep neural networks, Molchanov et al. ICML'17.
- [5] Learning to prune deep neural networks via layer-wise optimal brain surgeon, Dong et al. NIPS'17.
- [6] Learning Sparse Neural Networks through L0 Regularization, Louizos et al. ICLR'18

- Hyperparameters with weakly grounded heuristics (*e.g.,* layer-wise threshold [5], stochastic pruning rule [2])
- Architecture specific requirements

(e.g., conv/fc separate prune in [1])

- [1] Learning both weights and connections for efficient neural network, Han et al. NIPS'15
- [2] Dynamic network surgery for efficient dnns, Guo et al. NIPS'16.
- [3] Learning-compression algorithms for neural net pruning, Carreira-Perpinan & Idelbayev. CVPR'18.
- [4] Variational dropout sparsifies deep neural networks, Molchanov et al. ICML'17.
- [5] Learning to prune deep neural networks via layer-wise optimal brain surgeon, Dong et al. NIPS'17.
- [6] Learning Sparse Neural Networks through L0 Regularization, Louizos et al. ICLR'18

- Hyperparameters with weakly grounded heuristics (*e.g.*, layer-wise threshold [5], stochastic pruning rule [2])
- Architecture specific requirements

(e.g., conv/fc separate prune in [1])

• Optimization difficulty

(e.g., convergence in [3, 6])

- [1] Learning both weights and connections for efficient neural network, Han et al. NIPS'15
- [2] Dynamic network surgery for efficient dnns, Guo et al. NIPS'16.
- [3] Learning-compression algorithms for neural net pruning, Carreira-Perpinan & Idelbayev. CVPR'18.
- [4] Variational dropout sparsifies deep neural networks, Molchanov et al. ICML'17.
- [5] Learning to prune deep neural networks via layer-wise optimal brain surgeon, Dong et al. NIPS'17.
- [6] Learning Sparse Neural Networks through L0 Regularization, Louizos et al. ICLR'18

- Hyperparameters with weakly grounded heuristics (*e.g.*, layer-wise threshold [5], stochastic pruning rule [2])
- Architecture specific requirements

(e.g., conv/fc separate prune in [1])

• Optimization difficulty

(e.g., convergence in [3, 6])

• Pretraining step

([1,2,3,4,5,6]; almost all)

- [1] Learning both weights and connections for efficient neural network, Han et al. NIPS'15
- [2] Dynamic network surgery for efficient dnns, Guo et al. NIPS'16.
- [3] Learning-compression algorithms for neural net pruning, Carreira-Perpinan & Idelbayev. CVPR'18.
- [4] Variational dropout sparsifies deep neural networks, Molchanov et al. ICML'17.
- [5] Learning to prune deep neural networks via layer-wise optimal brain surgeon, Dong et al. NIPS'17.
- [6] Learning Sparse Neural Networks through L0 Regularization, Louizos et al. ICLR'18

- Hyperparameters with weakly grounded heuristics (*e.g.*, layer-wise threshold [5], stochastic pruning rule [2])
- Architecture specific requirements (e.g., conv/fc separate prune in [1])
- Optimization difficulty (*e.g.*, convergence in [3, 6])
- Pretraining step

([1,2,3,4,5,6]; almost all)



- [1] Learning both weights and connections for efficient neural network, Han et al. NIPS'15
- [2] Dynamic network surgery for efficient dnns, Guo et al. NIPS'16.
- [3] Learning-compression algorithms for neural net pruning, Carreira-Perpinan & Idelbayev. CVPR'18.
- [4] Variational dropout sparsifies deep neural networks, Molchanov et al. ICML'17.
- [5] Learning to prune deep neural networks via layer-wise optimal brain surgeon, Dong et al. NIPS'17.
- [6] Learning Sparse Neural Networks through L0 Regularization, Louizos et al. ICLR'18

## We want ...

No hyperparameters

No iterative prune -- retrain cycle

No pretraining

No large data

## We want ..

No hyperparameters

No iterative prune -- retrain cycle

No pretraining

No large data

Single-shot pruning prior to training



# SNIP: Single-shot Network Pruning based on Connection Sensitivity

N. Lee, T. Ajanthan, P. Torr

International Conference on Learning Representations (ICLR) 2019

### Objective

• Identify important parameters in the network and remove unimportant ones

$$\min_{\mathbf{w}} L(\mathbf{w}; \mathcal{D}) = \min_{\mathbf{w}} \frac{1}{n} \sum_{i=1}^{n} \ell(\mathbf{w}; (\mathbf{x}_i, \mathbf{y}_i)) ,$$
  
s.t.  $\mathbf{w} \in \mathbb{R}^m, \quad \|\mathbf{w}\|_0 \le \kappa .$ 

### Objective

• Identify important parameters in the network and remove unimportant ones

$$\min_{\mathbf{w}} L(\mathbf{w}; \mathcal{D}) = \min_{\mathbf{w}} \frac{1}{n} \sum_{i=1}^{n} \ell(\mathbf{w}; (\mathbf{x}_i, \mathbf{y}_i)) ,$$
  
s.t.  $\mathbf{w} \in \mathbb{R}^m, \quad \|\mathbf{w}\|_0 \le \kappa .$ 

• Measure the effect of removing each parameter on the loss

$$\Delta L_j(\mathbf{w}; \mathcal{D}) = L(\mathbf{1} \odot \mathbf{w}; \mathcal{D}) - L((\mathbf{1} - \mathbf{e}_j) \odot \mathbf{w}; \mathcal{D}) ,$$

• Measure the effect of removing each parameter on the loss

$$\Delta L_j(\mathbf{w}; \mathcal{D}) = L(\mathbf{1} \odot \mathbf{w}; \mathcal{D}) - L((\mathbf{1} - \mathbf{e}_j) \odot \mathbf{w}; \mathcal{D}) ,$$

• Measure the effect of removing each parameter on the loss

$$\Delta L_j(\mathbf{w}; \mathcal{D}) = L(\mathbf{1} \odot \mathbf{w}; \mathcal{D}) - L((\mathbf{1} - \mathbf{e}_j) \odot \mathbf{w}; \mathcal{D}),$$

• The greedy way is prohibitively expensive to perform: O(m!)

The effect on the loss can be approximated by

- 1. auxiliary variables representing the connectivity of parameters
- 2. derivative of the loss w.r.t. these indicator variables

1. Introduce c

$$\min_{\mathbf{c},\mathbf{w}} L(\mathbf{c} \odot \mathbf{w}; \mathcal{D}) = \min_{\mathbf{c},\mathbf{w}} \frac{1}{n} \sum_{i=1}^{n} \ell(\mathbf{c} \odot \mathbf{w}; (\mathbf{x}_i, \mathbf{y}_i)) ,$$
  
s.t.  $\mathbf{w} \in \mathbb{R}^m ,$   
 $\mathbf{c} \in \{0, 1\}^m, \quad \|\mathbf{c}\|_0 \le \kappa ,$ 

1. Introduce c

$$\min_{\mathbf{c},\mathbf{w}} L(\underline{\mathbf{c} \odot \mathbf{w}}; \mathcal{D}) = \min_{\mathbf{c},\mathbf{w}} \frac{1}{n} \sum_{i=1}^{n} \ell(\mathbf{c} \odot \mathbf{w}; (\mathbf{x}_i, \mathbf{y}_i)) ,$$
  
s.t.  $\mathbf{w} \in \mathbb{R}^m ,$   
 $\mathbf{c} \in \{0, 1\}^m, \quad \|\mathbf{c}\|_0 \le \kappa ,$ 

1. Introduce c

$$\min_{\mathbf{c},\mathbf{w}} L(\mathbf{c} \odot \mathbf{w}; \mathcal{D}) = \min_{\mathbf{c},\mathbf{w}} \frac{1}{n} \sum_{i=1}^{n} \ell(\mathbf{c} \odot \mathbf{w}; (\mathbf{x}_i, \mathbf{y}_i)) ,$$
s.t.  $\mathbf{w} \in \mathbb{R}^m ,$   
 $\mathbf{c} \in \{0, 1\}^m, \quad \|\mathbf{c}\|_0 \le \kappa ,$ 

2. Derivative w.r.t. c

57.8 (J D)

$$\Delta L_j(\mathbf{w}; \mathcal{D}) \approx g_j(\mathbf{w}; \mathcal{D}) = \left. \frac{\partial L(\mathbf{c} \odot \mathbf{w}; \mathcal{D})}{\partial c_j} \right|_{\mathbf{c} = \mathbf{1}} = \left. \lim_{\delta \to 0} \frac{L(\mathbf{c} \odot \mathbf{w}; \mathcal{D}) - L((\mathbf{c} - \delta \mathbf{e}_j) \odot \mathbf{w}; \mathcal{D})}{\delta} \right|_{\mathbf{c} = \mathbf{1}}$$

1. Introduce c

$$\min_{\mathbf{c},\mathbf{w}} L(\mathbf{c} \odot \mathbf{w}; \mathcal{D}) = \min_{\mathbf{c},\mathbf{w}} \frac{1}{n} \sum_{i=1}^{n} \ell(\mathbf{c} \odot \mathbf{w}; (\mathbf{x}_i, \mathbf{y}_i)) ,$$
s.t.  $\mathbf{w} \in \mathbb{R}^m ,$   
 $\mathbf{c} \in \{0, 1\}^m, \quad \|\mathbf{c}\|_0 \le \kappa ,$ 

2. Derivative w.r.t. c

$$\Delta L_j(\mathbf{w}; \mathcal{D}) \approx g_j(\mathbf{w}; \mathcal{D}) = \left. \frac{\partial L(\mathbf{c} \odot \mathbf{w}; \mathcal{D})}{\partial c_j} \right|_{\mathbf{c} = \mathbf{1}} = \lim_{\delta \to 0} \left. \frac{L(\mathbf{c} \odot \mathbf{w}; \mathcal{D}) - L(\mathbf{c} - \delta \mathbf{e}_j) \odot \mathbf{w}; \mathcal{D}}{\delta} \right|_{\mathbf{c} = \mathbf{1}}$$

- $\partial L/\partial cj$  is an infinitesimal version of  $\Delta Lj$
- measures the rate of change of L w.r.t. infinitesimal change in cj from 1  $\rightarrow$  1  $\delta$
- computed efficiently in one forward-backward pass using auto differentiation, for all j at once

Reference: Understanding black-box predictions via influence functions, Koh & Liang. ICML'17

1. Introduce c

$$\min_{\mathbf{c},\mathbf{w}} L(\mathbf{c} \odot \mathbf{w}; \mathcal{D}) = \min_{\mathbf{c},\mathbf{w}} \frac{1}{n} \sum_{i=1}^{n} \ell(\mathbf{c} \odot \mathbf{w}; (\mathbf{x}_i, \mathbf{y}_i)) ,$$
  
s.t.  $\mathbf{w} \in \mathbb{R}^m ,$   
 $\mathbf{c} \in \{0, 1\}^m, \quad \|\mathbf{c}\|_0 \le \kappa ,$ 

2. Derivative w.r.t. c

$$\Delta L_j(\mathbf{w}; \mathcal{D}) \approx g_j(\mathbf{w}; \mathcal{D}) = \left. \frac{\partial L(\mathbf{c} \odot \mathbf{w}; \mathcal{D})}{\partial c_j} \right|_{\mathbf{c} = \mathbf{1}} = \lim_{\delta \to 0} \left. \frac{L(\mathbf{c} \odot \mathbf{w}; \mathcal{D}) - L((\mathbf{c} - \delta \mathbf{e}_j) \odot \mathbf{w}; \mathcal{D})}{\delta} \right|_{\mathbf{c} = \mathbf{1}}$$

3. Connection sensitivity

**6**2 (2)

$$s_j = rac{|g_j(\mathbf{w}; \mathcal{D})|}{\sum_{k=1}^m |g_k(\mathbf{w}; \mathcal{D})|} \; .$$

#### Prune at initialization

- Measure CS on untrained networks prior to training
  - $\rightarrow$  Or zero gradients at pretrained
- Sample weights from a dist. with architecture aware variance
   → Ensure the variance of weights to remain throughout the network ([1])
- Alleviate the dependency on the weights in computing CS

 $\rightarrow$  Remove the pretraining requirement, architecture dep. hyperparameters

#### LeNets



#### LeNets



Method	Criterion	LeNet- $\bar{\kappa}$ (%)	-300-100 err. (%)	LeNet $\bar{\kappa}$ (%)	-5-Caffe err. (%)	Pretrain	# Prune	Additional hyperparam.	Augment objective	Arch. constraints
Ref.	-	_	1.7		0.9	-	_	_	-	_
LWC	Magnitude	91.7	1.6	91.7	0.8	$\checkmark$	many	$\checkmark$	×	$\checkmark$
DNS	Magnitude	98.2	2.0	99.1	0.9	$\checkmark$	many	$\checkmark$	X	$\checkmark$
LC	Magnitude	99.0	3.2	99.0	1.1	$\checkmark$	many	$\checkmark$	$\checkmark$	×
SWS	Bayesian	95.6	1.9	99.5	1.0	$\checkmark$	soft	$\checkmark$	$\checkmark$	×
SVD	Bayesian	98.5	1.9	99.6	0.8	$\checkmark$	soft	$\checkmark$	$\checkmark$	×
OBD	Hessian	92.0	2.0	92.0	2.7	$\checkmark$	many	$\checkmark$	×	×
L-OBS	Hessian	98.5	2.0	99.0	2.1	$\checkmark$	many	$\checkmark$	×	$\checkmark$
SNIP (ours)	Connection sensitivity	95.0 98.0	<b>1.6</b> 2.4	98.0 99.0	<b>0.8</b> 1.1	×	1	×	×	×

Method	Criterion	LeNet- $\bar{\kappa}$ (%)	-300-100 err. (%)	LeNet $\bar{\kappa}$ (%)	-5-Caffe err. (%)	Pretrain	# Prune	Additional hyperparam.	Augment objective	Arch. constraints
								VI 1	5	
Ref.		0 <del></del>	1.7		0.9	_			83 <del></del> 5	-
LWC	Magnitude	91.7	1.6	91.7	0.8	$\checkmark$	many	$\checkmark$	×	$\checkmark$
DNS	Magnitude	98.2	2.0	99.1	0.9	$\checkmark$	many	$\checkmark$	×	$\checkmark$
LC	Magnitude	99.0	3.2	99.0	1.1	$\checkmark$	many	$\checkmark$	$\checkmark$	×
SWS	Bayesian	95.6	1.9	99.5	1.0	$\checkmark$	soft	$\checkmark$	$\checkmark$	×
SVD	Bayesian	98.5	1.9	99.6	0.8	$\checkmark$	soft	$\checkmark$	$\checkmark$	×
OBD	Hessian	92.0	2.0	92.0	2.7	$\checkmark$	many	$\checkmark$	×	×
L-OBS	Hessian	98.5	2.0	99.0	2.1	$\checkmark$	many	$\checkmark$	×	$\checkmark$
SNIP (ours)	Connection sensitivity	$95.0 \\ 98.0$	<b>1.6</b> 2.4	$98.0 \\ 99.0$	0.8 1.1	×	1	×	×	×

Method	Criterion	LeNet- $\bar{\kappa}$ (%)	-300-100 err. (%)	LeNet $\bar{\kappa}$ (%)	-5-Caffe err. (%)	Pretrain	# Prune	Additional hyperparam.	Augment objective	Arch. constraints
								71 1	5	
Ref.		2. <del></del>	1.7		0.9	-	-		-	-
LWC	Magnitude	91.7	1.6	91.7	0.8	$\checkmark$	many	$\checkmark$	×	$\checkmark$
DNS	Magnitude	98.2	2.0	99.1	0.9	$\checkmark$	many	$\checkmark$	×	$\checkmark$
LC	Magnitude	99.0	3.2	99.0	1.1	$\checkmark$	many	$\checkmark$	$\checkmark$	×
SWS	Bayesian	95.6	1.9	99.5	1.0	$\checkmark$	soft	$\checkmark$	$\checkmark$	×
SVD	Bayesian	98.5	1.9	99.6	0.8	$\checkmark$	soft	$\checkmark$	$\checkmark$	×
OBD	Hessian	92.0	2.0	92.0	2.7	$\checkmark$	many	$\checkmark$	X	×
L-OBS	Hessian	98.5	2.0	99.0	2.1	$\checkmark$	many	$\checkmark$	×	$\checkmark$
SNIP (ours)	Connection sensitivity	95.0 98.0	<b>1.6</b> 2.4	98.0 99.0	<b>0.8</b> 1.1	×	1	×	×	×
	sensitivity	50.0	2.4	55.0	1.1					

Method	Criterion	LeNet- $\bar{\kappa}$ (%)	-300-100 err. (%)	LeNet $\bar{\kappa}$ (%)	-5-Caffe err. (%)	Pretrain	# Prune	Additional hyperparam.	Augment objective	Arch. constraints
Ref.		-	1.7	_	0.9	_			_	_
LWC	Magnitude	91.7	1.6	91.7	0.8	$\checkmark$	many	$\checkmark$	×	$\checkmark$
DNS	Magnitude	98.2	2.0	99.1	0.9	$\checkmark$	many	$\checkmark$	×	$\checkmark$
LC	Magnitude	99.0	3.2	99.0	1.1	$\checkmark$	many	$\checkmark$	$\checkmark$	×
SWS	Bayesian	95.6	1.9	99.5	1.0	$\checkmark$	soft	$\checkmark$	$\checkmark$	×
SVD	Bayesian	98.5	1.9	99.6	0.8	$\checkmark$	soft	$\checkmark$	$\checkmark$	×
OBD	Hessian	92.0	2.0	92.0	2.7	$\checkmark$	many	$\checkmark$	×	×
L-OBS	Hessian	98.5	2.0	99.0	2.1	$\checkmark$	many	$\checkmark$	×	$\checkmark$
SNIP (ours)	Connection sensitivity	$95.0 \\ 98.0$	<b>1.6</b> 2.4	98.0 99.0	<b>0.8</b> 1.1	×	1	×	×	×

Method	Criterion	LeNet- $\bar{\kappa}$ (%)	-300-100 err. (%)	LeNet $\bar{\kappa}$ (%)	-5-Caffe err. (%)	Pretrain	# Prune	Additional hyperparam.	Augment objective	Arch. constraints
Ref.	-	_	1.7		0.9	-	_		-	-
LWC	Magnitude	91.7	1.6	91.7	0.8	$\checkmark$	many	$\checkmark$	×	$\checkmark$
DNS	Magnitude	98.2	2.0	99.1	0.9	$\checkmark$	many	$\checkmark$	×	$\checkmark$
LC	Magnitude	99.0	3.2	99.0	1.1	$\checkmark$	many	$\checkmark$	$\checkmark$	×
SWS	Bayesian	95.6	1.9	99.5	1.0	$\checkmark$	soft	$\checkmark$	$\checkmark$	×
SVD	Bayesian	98.5	1.9	99.6	0.8	$\checkmark$	soft	$\checkmark$	$\checkmark$	×
OBD	Hessian	92.0	2.0	92.0	2.7	$\checkmark$	many	$\checkmark$	×	×
L-OBS	Hessian	98.5	2.0	99.0	2.1	$\checkmark$	many	$\checkmark$	×	$\checkmark$
SNIP (ours)	Connection sensitivity	95.0 98.0	<b>1.6</b> 2.4	98.0 99.0	<b>0.8</b> 1.1	×	1	×	×	×

Architecture	Model	Sparsity (%)	# Parameters	Error (%)	$\Delta$
Convolutional	AlexNet-s AlexNet-b VGG-C VGG-D VGG-like	90.0 90.0 95.0 95.0 97.0	$\begin{array}{rrrrrrrrrrrrrrrrrrrrrrrrrrrrrrrrrrrr$	$\begin{array}{cccccccccccccccccccccccccccccccccccc$	+0.87 +0.58 +0.45 +0.33 -0.26
Residual	WRN-16-8 WRN-16-10 WRN-22-8	$95.0 \\ 95.0 \\ 95.0 \\ 95.0$	$\begin{array}{rrrr} 10.0m \ \rightarrow \ 548k \\ 17.1m \ \rightarrow \ 856k \\ 17.2m \ \rightarrow \ 858k \end{array}$	$\begin{array}{rrrr} 6.21 \to & 6.63 \\ 5.91 \to & 6.43 \\ 6.14 \to & 5.85 \end{array}$	+0.42 +0.52 -0.29
Recurrent	LSTM-s LSTM-b GRU-s GRU-b	95.0 95.0 95.0 95.0	$\begin{array}{rrrr} 137k \rightarrow & 6.8k \\ 535k \rightarrow 26.8k \\ 104k \rightarrow & 5.2k \\ 404k \rightarrow & 20.2k \end{array}$	$\begin{array}{rrrrrrrrrrrrrrrrrrrrrrrrrrrrrrrrrrrr$	<b>−0.31</b> +0.20 +0.54 <b>−0.19</b>

Architecture	Model	Sparsity (%)	# Parameters	Error (%)	Δ
Convolutional	AlexNet-s AlexNet-b VGG-C VGG-D VGG-like	90.0 90.0 95.0 95.0 97.0	$\begin{array}{rrrrrrrrrrrrrrrrrrrrrrrrrrrrrrrrrrrr$	$\begin{array}{cccccccccccccccccccccccccccccccccccc$	+0.87 +0.58 +0.45 +0.33 -0.26
Residual	WRN-16-8 WRN-16-10 WRN-22-8	$95.0 \\ 95.0 \\ 95.0 \\ 95.0$	$\begin{array}{rrrrrrrrrrrrrrrrrrrrrrrrrrrrrrrrrrrr$	$\begin{array}{rrrr} 6.21 \ \to & 6.63 \\ 5.91 \ \to & 6.43 \\ 6.14 \ \to & 5.85 \end{array}$	+0.42 +0.52 -0.29
Recurrent	LSTM-s LSTM-b GRU-s GRU-b	95.0 95.0 95.0 95.0	$\begin{array}{rrrr} 137k \rightarrow & 6.8k \\ 535k \rightarrow & 26.8k \\ 104k \rightarrow & 5.2k \\ 404k \rightarrow & 20.2k \end{array}$	$ \begin{array}{rrrrrrrrrrrrrrrrrrrrrrrrrrrrrrrrrrrr$	-0.31 +0.20 +0.54 -0.19

Architecture	Model	Sparsity (%)	# Parameters	Error (%)	Δ
Convolutional	AlexNet-s AlexNet-b VGG-C VGG-D VGG-like	90.0 90.0 95.0 95.0 97.0	$\begin{array}{rrrrrrrrrrrrrrrrrrrrrrrrrrrrrrrrrrrr$	$\begin{array}{cccccccccccccccccccccccccccccccccccc$	+0.87 +0.58 +0.45 +0.33 -0.26
Residual	WRN-16-8 WRN-16-10 WRN-22-8	$95.0 \\ 95.0 \\ 95.0 \\ 95.0$	$\begin{array}{rrrrrrrrrrrrrrrrrrrrrrrrrrrrrrrrrrrr$	$\begin{array}{rrrr} 6.21 \ \to & 6.63 \\ 5.91 \ \to & 6.43 \\ 6.14 \ \to & 5.85 \end{array}$	+0.42 +0.52 -0.29
Recurrent	LSTM-s LSTM-b GRU-s GRU-b	$95.0 \\ 95.0 \\ 95.0 \\ 95.0 \\ 95.0$	$\begin{array}{rrrr} 137\mathrm{k} &\rightarrow & 6.8\mathrm{k} \\ 535\mathrm{k} &\rightarrow & 26.8\mathrm{k} \\ 104\mathrm{k} &\rightarrow & 5.2\mathrm{k} \\ 404\mathrm{k} &\rightarrow & 20.2\mathrm{k} \end{array}$	$\begin{array}{rrrrrrrrrrrrrrrrrrrrrrrrrrrrrrrrrrrr$	<b>−0.31</b> +0.20 +0.54 <b>−0.19</b>

Architecture	Model	Sparsity (%)	# Parameters	Error (%)	Δ
Convolutional	AlexNet-s AlexNet-b VGG-C VGG-D VGG-like	90.0 90.0 95.0 95.0 97.0	$\begin{array}{rrrrrrrrrrrrrrrrrrrrrrrrrrrrrrrrrrrr$	$\begin{array}{cccccccccccccccccccccccccccccccccccc$	+0.87 +0.58 +0.45 +0.33 -0.26
Residual	WRN-16-8 WRN-16-10 WRN-22-8	$95.0 \\ 95.0 \\ 95.0$	$\begin{array}{rrrr} 10.0m \ \rightarrow \ 548k \\ 17.1m \ \rightarrow \ 856k \\ 17.2m \ \rightarrow \ 858k \end{array}$	$\begin{array}{rrrr} 6.21 \to & 6.63 \\ 5.91 \to & 6.43 \\ 6.14 \to & 5.85 \end{array}$	+0.42 +0.52 -0.29
Recurrent	LSTM-s LSTM-b GRU-s GRU-b	95.0 95.0 95.0 95.0	$\begin{array}{rrrr} 137k \rightarrow & 6.8k \\ 535k \rightarrow 26.8k \\ 104k \rightarrow & 5.2k \\ 404k \rightarrow & 20.2k \end{array}$	$\begin{array}{rrrrrrrrrrrrrrrrrrrrrrrrrrrrrrrrrrrr$	<b>−0.31</b> +0.20 +0.54 <b>−0.19</b>

Visualize c in the first fc layer for varying data

- 1. curate a mini-batch
- 2. compute the connection sensitivity
- 3. create the pruning mask
- 4. visualize the first layer (fully connected)

Visualize c in the first fc layer for varying data

- 1. curate a mini-batch
- 2. compute the connection sensitivity
- 3. create the pruning mask
- 4. visualize the first layer (fully connected)



The input was digit 8.

Visualize c in the first fc layer for varying data

- 1. curate a mini-batch
- 2. compute the connection sensitivity
- 3. create the pruning mask
- 4. visualize the first layer (fully connected)



The input was digit 8.

Carrying out such inspection is not straightforward with other methods.



#### Carrying out such inspection is not straightforward with other methods.



The parameters connected to the discriminative part of image are retained.

#### **Prevent memorization**



[Fitting random labels] Understanding deep learning requires rethinking generalization, Zhang et al. ICLR'17

<u>The pruned network does not have sufficient capacity to fit the random labels.</u> <u>but is capable of performing the task.</u>

Simple Versatile Interpretable Paper: https://arxiv.org/abs/1810.02340

Code: https://github.com/namhoonlee/snip-public

Contact: http://www.robots.ox.ac.uk/~namhoon/