## Slimming ResNet by Slimming Shortcut

Donggyu Joo, Doyeon Kim, and Junmo Kim

SIIT Lab. KAIST





Statistical Inference and Information Theory Laboratory

## **Introduction & Objectives**



## **Existing Pruning**

- $\checkmark$  Pruning is applied to only inside of residual block
- ✓ Unstable bottleneck-like structure is obtained



### **Proposed** SSPruning

- $\checkmark$  Pruning is applied to both shortcut and block.
- $\checkmark$  We obtain the optimal structure without restriction

#### **Method** 1. Shortcut Separation



(a) Original ResNet architecture. All shortcuts are connected together.

(b) Shortcut Separation: Additional  $1 \times 1$  convolutional is added at each shortcut to separate them.

(c, d) To remove the channels related to the addition layer (marked in purple), we only need to consider the channels in layers marked in blue.



- ✓ The importance learning gate (ILG) layer is applied after each addition layer.
- ✓ ILG is composed of *n* parameters  $\alpha_1, \alpha_2, ..., \alpha_n \in R$  where *n* is the number of channels in the corresponding location. These *n* parameters are initialized as 1, and independently multiplied to the output of each corresponding channel.

#### **1. Comparative Experiments**

- ✓ Datasets: CIFAR-10, Network: ResNet-32 and ResNet-110
- ✓ SSPruning- $\alpha$ : Result of the proposed method ( $\alpha$ % of FLOPs are pruned from the original network).
- ✓ Our SSPruning achieves the state-of-the-art performance in several aspects.

Network	Method	Baseline Acc (%)	Pruned Acc (%)	$\begin{array}{c} \mathrm{Acc} \downarrow \\ (\%) \end{array}$	Flops $\downarrow$ (%)
ResNet-32	MIL [32] SFP [33] FPGM [12] SSPruning-50 SSPruning-60	92.33 92.63 92.63 <b>92.54</b> <b>92.54</b>	90.74 92.08 91.93 <b>92.80</b> <b>92.36</b>	1.59 0.55 0.70 <b>-0.26</b> <b>0.18</b>	31.2 41.5 53.2 <b>50.1</b> <b>60.1</b>
ResNet-110	PFEC [34] SFP [33] NISP [29] FPGM [12] SSPruning-50 SSPruning-60	93.53 93.68 93.68 93.65 93.65	93.30 93.38 93.85 93.99 93.76	0.23 0.30 0.18 -0.17 <b>-0.34</b> <b>-0.11</b>	38.6 40.8 43.8 52.3 <b>50.1</b> <b>60.1</b>

#### **2. Ablation on Shortcut Pruning**

- ✓ Datasets: CIFAR-10, Network: ResNet-32
- ✓ Pruning of both shortcut and inblock outperforms traditional inblock pruning in various settings.

Experiment for SSPruning							
	Pruned Acc (%)						
Flops $\downarrow$ (%)	InBlock Only	Shortcut Only	Both (proposed)				
50%	92.58	92.16	92.80				
60%	92.08	91.57	92.36				

Experiment for Random Selection						
	Pruned Acc(%)					
Flops $\downarrow$ (%)	InBlock Only	Shortcut Only	Both			
50%	92.32	92.29	92.51			
60%	91.84	91.85	92.03			

(a)

(b)

# THANK YOU

## Slimming ResNet by Slimming Shortcut

ICPR 2020





Statistical Inference and Information Theory Laboratory