Training CNNs with Selective Allocation of Channels

<u>Jongheon Jeong</u>¹ Jinwoo Shin^{1,2}

¹Korea Advanced Institute of Science and Technology (KAIST) ²AITRICS

ICML 2019

Recent Advances in CNN Architectures

• Larger model, better accuracy \rightarrow "efficient" architecture design



Recent Advances in CNN Architectures

- Larger model, better accuracy \rightarrow "efficient" architecture design
- Motivation: Current CNN design typically focus on static layers



"Static" Layers?

- Current convolutional layers allocate parameters uniformly across channels
- "Static": The allocation is fixed until the end of training



Channel Inefficiency in "Static" CNNs

- Current convolutional layers allocate parameters uniformly across channels
- **"Static"**: The allocation is **fixed** until the end of training
- What if the layer contains unnecessary channels to compute?
 - The parameters allocated for those channels can hardly utilized



Channel Inefficiency in "Static" CNNs

- Current convolutional layers allocate parameters uniformly across channels
- **"Static"**: The allocation is **fixed** until the end of training
- What if the layer contains unnecessary channels to compute?
 - The parameters allocated for those channels can hardly utilized
- Moreover, # channels are usually set without much care in recent CNNs
 - Can we utilize them in training for efficiency?



^{*} Source: He et al. "Deep Residual Learning for Image Recognition", CVPR 2016

Framework: Training with "Dynamic" Re-wiring

- Idea: Incorporating function-preserving operations on channels
- A re-wiring operation f can be called at any time in training



Framework: Training with "Dynamic" Re-wiring

- Idea: Incorporating function-preserving operations on channels
- A re-wiring operation f can be called at any time in training
 - Connectivity is updated without affecting the overall training



Framework: Training with "Dynamic" Re-wiring

- Idea: Incorporating function-preserving operations on channels
- A re-wiring operation f can be called at any time in training
 - Connectivity is updated without affecting the overall training
 - Manipulation on **channels** rather than **parameters** \rightarrow architecture-agnostic



- Idea: Incorporating function-preserving operations on channels
- Two re-wiring operations: dealloc & realloc
 - 1. dealloc: Release unimportant channels \rightarrow pruning parameters
 - 2. realloc: Replicate important channels \rightarrow re-using the pruned parameter



- Two operations during training: dealloc & realloc
 - 1. Channel de-allocation (dealloc): Release "unimportant" channels
 - Drops unimportant channels based on a measure $||\Delta_{-i}||_2$ we defined



- Two operations during training: dealloc & realloc
 - 1. Channel de-allocation (dealloc): Release "unimportant" channels
 - Drops unimportant channels based on a measure $||\Delta_{-i}||_2$ we defined
 - Equivalent to pruning parameters corresponding to the de-allocated channels



- Two operations during training: dealloc & realloc
 - 1. Channel de-allocation (dealloc): Release "unimportant" channels
 - Drops unimportant channels based on a measure $||\Delta_{-i}||_2$ we defined
 - Equivalent to pruning parameters corresponding to the de-allocated channels
 - Key point: dealloc has to be function-preserving:

 $\operatorname{Conv}(\texttt{dealloc}(\mathbf{X});\mathbf{W})\approx\operatorname{Conv}(\mathbf{X};\mathbf{W})$



- Two operations during training: dealloc & realloc
 - Key point: dealloc has to be function-preserving $\operatorname{Conv}(\operatorname{dealloc}(\mathbf{X}); \mathbf{W}) \approx \operatorname{Conv}(\mathbf{X}; \mathbf{W})$



• We measure **expected channel damage** for channel importance

$$\Delta_{-i} := \frac{1}{HW} \sum_{h,w} \mathbb{E}_{\mathbf{X}} [\underbrace{\operatorname{Conv}(\mathbf{X}; \mathbf{W})}_{\text{original}} - \underbrace{\operatorname{Conv}(\mathbf{X}; \mathbf{W}_{-i})}_{\text{after pruning channel } i}]_{:,h,w} \in \mathbb{R}^{O}$$

- Difference after pruning channel $i \rightarrow$ function-preserving property
- The value is averaged over the spatial dimension *h*, *w*
- $\Delta \in \mathbb{R}^{I \times O}$: Expected Channel Damage Matrix (ECDM)

- Two operations during training: dealloc & realloc
 - Key point: dealloc has to be function-preserving $\operatorname{Conv}(\operatorname{dealloc}(\mathbf{X}); \mathbf{W}) \approx \operatorname{Conv}(\mathbf{X}; \mathbf{W})$



• We measure **expected channel damage** for channel importance

$$\Delta_{-i} := \frac{1}{HW} \sum_{h,w} \mathbb{E}_{\mathbf{X}} [\underline{\operatorname{Conv}(\mathbf{X}; \mathbf{W})}_{\text{original}} - \underline{\operatorname{Conv}(\mathbf{X}; \mathbf{W}_{-i})}_{\text{after pruning channel } i}]_{:,h,w} \in \mathbb{R}^{O}$$

- Difference after pruning channel $i \rightarrow$ function-preserving property
- The value is averaged over the spatial dimension h, w
- $\Delta \in \mathbb{R}^{I \times O}$: Expected Channel Damage Matrix (ECDM)
- Challenge: Computing Δ_{-i} requires a marginalization over X 🤔



Selective Convolutional Layer

- Two operations during training: dealloc & realloc
 - We measure **expected channel damage** for channel importance

$$\Delta_{-i} := \frac{1}{HW} \sum_{h,w} \mathbb{E}_{\mathbf{X}} [\underbrace{\operatorname{Conv}(\mathbf{X}; \mathbf{W})}_{\text{original}} - \underbrace{\operatorname{Conv}(\mathbf{X}; \mathbf{W}_{-i})}_{\text{after pruning channel } i}]_{:,h,w} \in \mathbb{R}^{O}$$

- Challenge: Computing Δ_{-i} requires a marginalization over ${f X}$
- Idea: Use a BatchNorm layer to approximate Δ_{-i}
 - 1. BatchNorm $(x) \sim \mathcal{N}(\beta, \gamma^2)$: β, γ are BatchNorm parameters
 - 2. $\mathbf{X} = \operatorname{ReLU}(BN(\mathbf{Y}))$: A common design in modern CNNs

$$\Delta_{-i} \approx \underbrace{\left(|\gamma_i| \phi_{\mathcal{N}} \left(\frac{\beta_i}{|\gamma_i|} \right) + \beta_i \Phi_{\mathcal{N}} \left(\frac{\beta_i}{|\gamma_i|} \right) \right)}_{\text{p.d.f.} \text{ activity level } \textbf{c.d.f.}} \cdot \underbrace{\sum_{k=1}^{K^2} \mathbf{W}_{i,:,k}}_{\text{magnitude}}$$



Χ

- Two operations during training: dealloc & realloc
 - **Challenge**: Computing Δ_{-i} requires a marginalization over **X**
 - Idea: Use a BatchNorm layer to approximate Δ_{-i}
 - BatchNorm(x) ~ $\mathcal{N}(\beta, \gamma^2)$: β, γ are BatchNorm parameters
 - $\mathbf{X} = \operatorname{ReLU}(\operatorname{BN}(\mathbf{Y}))$: A common design in modern CNNs
 - Are those assumptions realistic in practice?









-1.5

-1.0

-0.5

Mean

- Two operations during training: dealloc & realloc
 - Challenge: Computing Δ_{-i} requires a marginalization over ${f X}$
- 5

- Are those assumptions realistic in practice?
 - 1. For a fixed channel, most of the pixel distributions are unimodal (a)
 - 2. The means and variances are clustered for each channel (d)
 - 3. The trend still exists even the model re-initialized (c, d)





- Two operations during training: dealloc & realloc
 - Challenge: Computing Δ_{-i} requires a marginalization over ${f X}$

X

- Are those assumptions realistic in practice?
 - Two properties of CNN can be responsible:
 - 1. CLT from the linear, weighted summing nature of convolution
 - 2. Translation-equivariance of convolution on spatial dimensions



- **Two operations** during training: **dealloc** & **realloc**
 - 1. Channel de-allocation (dealloc): Release "unimportant" channels
 - Drops unimportant channels based on a measure $||\Delta_{-i}||_2$ we defined
 - Equivalent to pruning parameters corresponding to the de-allocated channels
 - Key point: dealloc has to be function-preserving



- Two operations during training: dealloc & realloc
 - 2. Channel re-allocation (realloc): Replicate "important" channels into the released area



- Channels with high $||\Delta_{-i}||_2$ are copied into the de-allocated channels
- Challenge: Naïvely copying a channel in does NOT give any benefit 🤔
 - Due to the **linearity** of convolutional layers
 - A convolutional filter will see the exact same patch of the two images

- Two operations during training: dealloc & realloc
 - 2. Channel re-allocation (realloc): Replicate "important" channels into the released area



• Channels with high $||\Delta_{-i}||_2$ are copied, but with spatial shifting bias $b = (b^h, b^w)$

- Two operations during training: dealloc & realloc
 - 2. Channel re-allocation (realloc): Replicate "important" channels into the released area



- Two operations during training: dealloc & realloc
 - 2. Channel re-allocation (realloc): Replicate "important" channels into the released area



- Framework: Incorporating re-wiring operations in training
- Two operations during training: dealloc & realloc
- Dynamic re-wiring of parameters \rightarrow selective kernel expansion



- Framework: Incorporating re-wiring operations in training
- Two operations during training: dealloc & realloc
- Dynamic re-wiring of parameters \rightarrow selective kernel expansion



- Framework: Incorporating re-wiring operations in training
- Two operations during training: dealloc & realloc
- Flexible training: model reduction \leftrightarrow accuracy improvement



Experiments: Improving Modern CNNs

- Selective convolution can be readily applied to various existing CNNs
- Training each convolutional layer with dealloc + realloc

			Error rates (%)							
Model	Params	Method	CIFAR-10	CIFAR-10	0	Fashie	on-MNIST		Tiny-In	nageNet
DenseNet-40 (bottleneck, $k = 12$)	0.21M	Baseline SelectConv	6.62±0.15 6.09±0.10 (-8.01%)	29.9±0.1 28.8±0.1 (-3.42%)		5.03±0.07 4.73±0.06 (-5.96%)		%)	45.8±0.2 44.4±0.2 (-3.03%)	
DenseNet-100 (bottleneck, $k = 12$)	1.00M	Baseline SelectConv	4.51±0.04 4.29±0.08 (-4.88%)	22.8±0.3 22.2±0.1 (-	2.8±0.3 2.2±0.1 (-2.64%)		4.70±0.06 4.58±0.05 (-2.55%)		41.0±0.1 39.9 ±0.3 (-2.78%)	
ResNet-164 (bottleneck, pre-act)	1.66M	Baseline SelectConv	4.23±0.15 3.92±0.14 (-7.33%)	21.3±0.2 20.9±0.2 (-	1.97%)	4.53± 4.37 ±	0.04 0.03 (-3.53	%)	37.7±0. 37.5 ±0.	4 2 (-0.56 %)
ResNeXt-29 ($8 \times 64d$)	33.8M	Baseline SelectConv	3.62±0.12 3.39±0.14 (-6.36%)	18.1±0.1 17.6 ±0.1 (-	(-2.92%) 4.40±0 4.27±0		±0.07 ± 0.06 (- 2.95 %)		31.7±0.3 31.4±0.3 (-0.88%)	
					Model		Params	М	lethod	Error (%)
(top) Results on CIFAR-10/100, FMNIST and Tiny-ImageNet (right) Results on ImageNet dataset					DenseNet-121 $(k = 32)$ 7.98M		7.98M	Ba	aseline	24.7
							SelectConv		24.4	

ResNet-50

(bottleneck)

23.9

23.4

Baseline

SelectConv

22.8M

Experiments: Improving Modern CNNs

- Selective convolution can be readily applied to various existing CNNs
- Reduction in error rates across all the tested architectures

			Error rates (%)				
Model	Params	Method	CIFAR-10	CIFAR-10)	Fashion-MNIST	Tiny-ImageNet
DenseNet-40 (bottleneck, $k = 12$)	0.21M	Baseline SelectConv	6.62±0.15 6.09±0.10 (-8.01%)	29.9±0.1 28.8±0.1 (-	3.42%)	5.03±0.07 4.73 ±0.06 (-5.96%	$\begin{array}{c} 45.8 \pm 0.2 \\ 44.4 \pm 0.2 \ \textbf{(-3.03\%)} \end{array}$
DenseNet-100 (bottleneck, $k = 12$)	1.00M	Baseline SelectConv	4.51±0.04 4.29±0.08 (-4.88%)	22.8±0.3 22.2±0.1 (-	2.64%)	4.70±0.06 4.58±0.05 (-2.55%	$\begin{array}{c} 41.0 \pm 0.1 \\ \textbf{39.9} \pm 0.3 \ \textbf{(-2.78\%)} \end{array}$
ResNet-164 (bottleneck, pre-act)	1.66M	Baseline SelectConv	4.23±0.15 3.92±0.14 (-7.33%)	21.3±0.2 20.9±0.2 (-	1.97%)	4.53±0.04 4.37±0.03 (-3.53%	37.7±0.4 6) 37.5±0.2 (-0.56%)
ResNeXt-29 ($8 \times 64d$)	33.8M	Baseline SelectConv	3.62±0.12 3.39±0.14 (-6.36%)	18.1±0.1 17.6 ±0.1 (-	2.92%)	4.40±0.07 4.27±0.06 (-2.95%	31.7±0.3 6) 31.4±0.3 (-0.88%)
					Model	Params	Method Error (%)
(top) Results on CIFAR-10/100, FMNIST and Tinv-ImageNet					DenseN	1 aranis	Baseline 24.7

• (right) Results on ImageNet dataset

• **Remark**: $23.6\% \rightarrow 23.0\% \approx 51$ more layers (official ResNet)

24.4

23.9

23.4

7.98M

22.8M

(k = 32)

ResNet-50

(bottleneck)

SelectConv

Baseline

SelectConv

Experiments: Mobile-targeted Architectures

- Selective convolution can further improve the "already-efficient" CondenseNet-182
- Training with dealloc \rightarrow model compression
- ECDM-based de-allocation outperforms existing dynamic channel pruning scheme
 - Replacing LGC to SConv → further compression



Model (CIFAR-10)	Params	FLOPs	Error (%)
ResNet-1001	16.1M	2,357M	4.62
WideResNet-28-10	36.5M	5,248M	4.17
ResNeXt-29 $(16 \times 64d)$	68.1M	$10,704 \mathrm{M}$	3.58
VGGNet-Slim [2]	2.30M	391M	6.20
ResNet-164-Slim [2]	1.10M	$275 \mathrm{M}$	5.27
CondenseNet-LGC-182 [1]	4.20M	513M	3.76
CondenseNet-SConv-182	$3.24\mathrm{M}$	422M	3.50

[1] Gao Huang et al. Condensenet: An efficient densenet using learned group convolutions. CVPR 2018.[2] Zhuang Liu et al. Learning efficient convolutional networks through network slimming. ICCV 2017.

Summary

- A new functionality: channel-selectivity
 - Dynamically choosing channels toward selective expansion of kernels
- We propose selective convolution = convolution + channel-selectivity
 - 1. Generic, easy to use: applicable to any kind of CNN
 - 2. Single-pass: no post-processing/re-training
 - 3. **On-demand**: accuracy improvement \leftrightarrow model compression
- We define a new metric of channel importance: expected channel damage

