

Appendices for DSA: More Efficient Budgeted Pruning via Differentiable Sparsity Allocation

Xuefei Ning^{1*}, Tianchen Zhao^{2*}, Wenshuo Li¹, Peng Lei², Yu Wang¹, and Huazhong Yang¹

¹ Department of Electronic Engineering, Tsinghua University

² Department of Electronic Engineering, Beihang University

foxdoraame@gmail.com, ztc16@buaa.edu.cn, yu-wang@tsinghua.edu.cn

1 Topological Grouping and Budget Model \mathcal{F}

Alg. 1 describes the procedure of topological grouping to handle the topological constraints introduced by shortcuts, depthwise convolutions and so on. An example of grouping convolutions in two consecutive residual blocks is shown in Fig. 1. All the incoming connections of the normal convolutions are removed, and then the convolutions in each connected component belong to the same topological group (i.e., share the same keep ratio $\alpha^{(k)}$ and masks $m_{i=1, \dots, C}^{(k)}$).

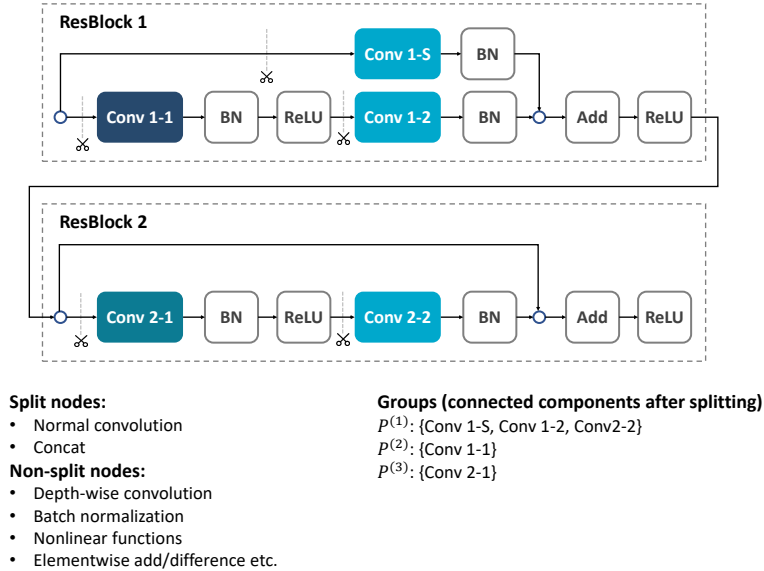


Fig. 1. An example of the topological grouping procedure

* Both authors contributed equally to this work.

Algorithm 1 Topological grouping procedure

- 1: Construct the computational directed acyclic graph G
 - 2: Removing all the incoming connections at split nodes (operation nodes that is not a channel-wise operation): normal convolution, concat operation.
NOTE: Non-split nodes include all channel-wise operations: depthwise convolution, element-wise add, ReLU, batch normalization, etc.
 - 3: Find the connected components $\{P^{(k)}\}_{k=1, \dots, K}$. All the convolution layers in each $P^{(k)}$ share the same keep ratio $\alpha^{(k)}$ and masks $m_{i=1, \dots, C}^{(k)}$
-

A budget model \mathcal{F} is needed for measuring the resource consumption $\mathcal{F}(\mathcal{A})$ corresponding to the sparsity allocation \mathcal{A} . Taking FLOPs as an example, $\mathcal{F}(\mathcal{A})$ could be represented as $\mathcal{F}(\mathcal{A}) = \mathcal{A}^T \mathcal{F}_A \mathcal{A} + \mathcal{F}_B^T \mathcal{A}$. We summarize the calculation procedure of \mathcal{F}_A and \mathcal{F}_B in Alg. 2.

Algorithm 2 Calculation of $\mathcal{F}_A, \mathcal{F}_B$ (\mathcal{F} for FLOPs resource)

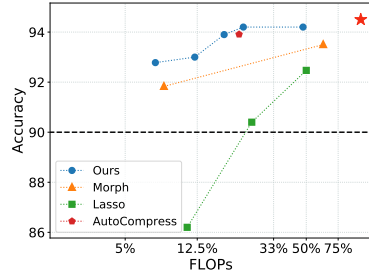
- 1: G : the directed acyclic graph of operations
 - 2: K : number of connected components
 - 3: $\mathcal{F}_A = \mathbf{0}_{K \times K}, \mathcal{F}_B = \mathbf{0}_K$
 - 4: Convolution node attributes: 1) C : output channel number; 2) k : the index of the connected components that the convolution node belongs to; 3) kss : kernel spatial size (e.g. $3 \times 3 = 9$); 4) oss : output spatial size (e.g. $16 \times 16 = 256$)
 - 5: **for all** convolution node M in G **do**
 - 6: **if** $n.type == DEPTHWISE.CONV$ **then**
 - 7: $\mathcal{F}_B[M.k] += 2 \times M.c \times M.kss \times M.oss$
 - 8: **else**
 - 9: $stack = [predecessor(M)]$
 - 10: **while** $stack$ **do**
 - 11: $n = stack.pop()$
 - 12: **if** $n.type == CONCAT$ **then**
 - 13: **for** pn in $predecessor(n)$ **do**
 - 14: $stack.push(pn)$
 - 15: **end for**
 - 16: **else if** $n.type$ in $ELEMENTWISE.OPs$ (e.g. ADD, ReLU) **then**
 - 17: $stack.push(predecessor(n)[0])$
 - 18: **else if** $n.type == NORMAL.CONV$ **then**
 - 19: $\mathcal{F}_A[M.k, n.k] += 2 \times n.C \times M.C \times M.kss \times M.oss$
 - 20: **else if** $n.type == INPUT$ **then**
 - 21: $\mathcal{F}_B[M.k] += 2 \times n.C \times M.C \times M.kss \times M.oss$
 - 22: **else**
 - 23: $stack.push(n)$
 - 24: **end if**
 - 25: **end while**
 - 26: **end if**
 - 27: **end for**
 - 28: **return** $\mathcal{F}_A, \mathcal{F}_B$
-

2 Additional Results

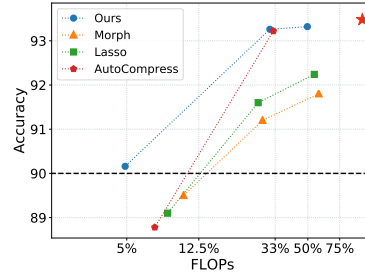
Table. 1 and Fig. 2 show the results of pruning ResNet-18 and VGG-16 on CIFAR-10.

Table 1. Pruning results of ResNet-18 and VGG-16 on CIFAR-10

Method	ResNet-18 (94.0%)		VGG-16 (93.48%)	
	Acc.	FLOPs ratio	Acc.	FLOPs ratio
AutoCompress	93.91%	4.7× (21.28%)	93.22%	3.1× (32.26%)
	-	-	88.78%	14.0× (7.14%)
DSA (Ours)	94.19%	4.5× (22.46%)	93.26%	3.2× (30.85%)
	93.90%	5.7× (17.54%)	90.16%	20.4× (4.90%)



(a) ResNet-18



(b) VGG-16

Fig. 2. Pruning results of ResNet-18, VGG-16 on CIFAR-10