

Structured Sparsification with Joint Optimization of Group Convolution and Channel Shuffle

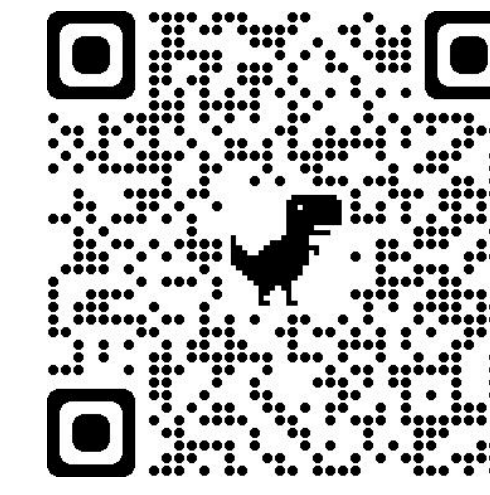
Xin-Yu Zhang¹, Kai Zhao¹, Taihong Xiao², Ming-Ming Cheng¹, Ming-Hsuan Yang²

¹Nankai University, ²University of California, Merced

arXiv: <https://arxiv.org/abs/2002.08127>

Github: <https://github.com/Sakura03/StrucSpars>

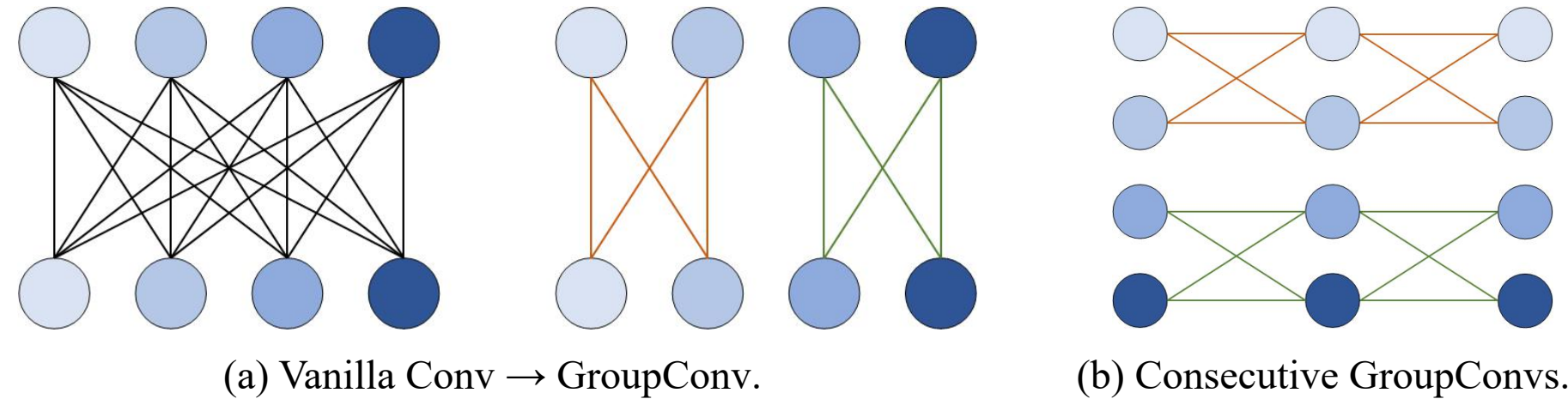
Email: xinyuzhang@mail.nankai.edu.cn



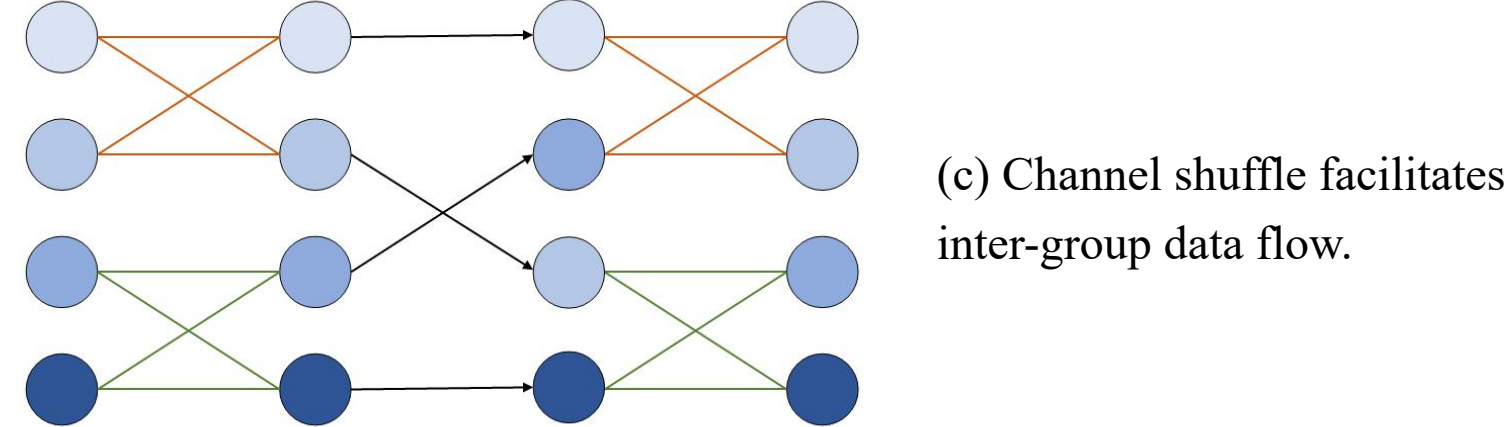
Background

Group convolution (GroupConv) for model compression.

- GroupConv has less parameters and lower complexity (Fig. (a));
- However, if multiple GroupConv are stacked sequentially, the inter-group communication will be eliminated (Fig. (b));

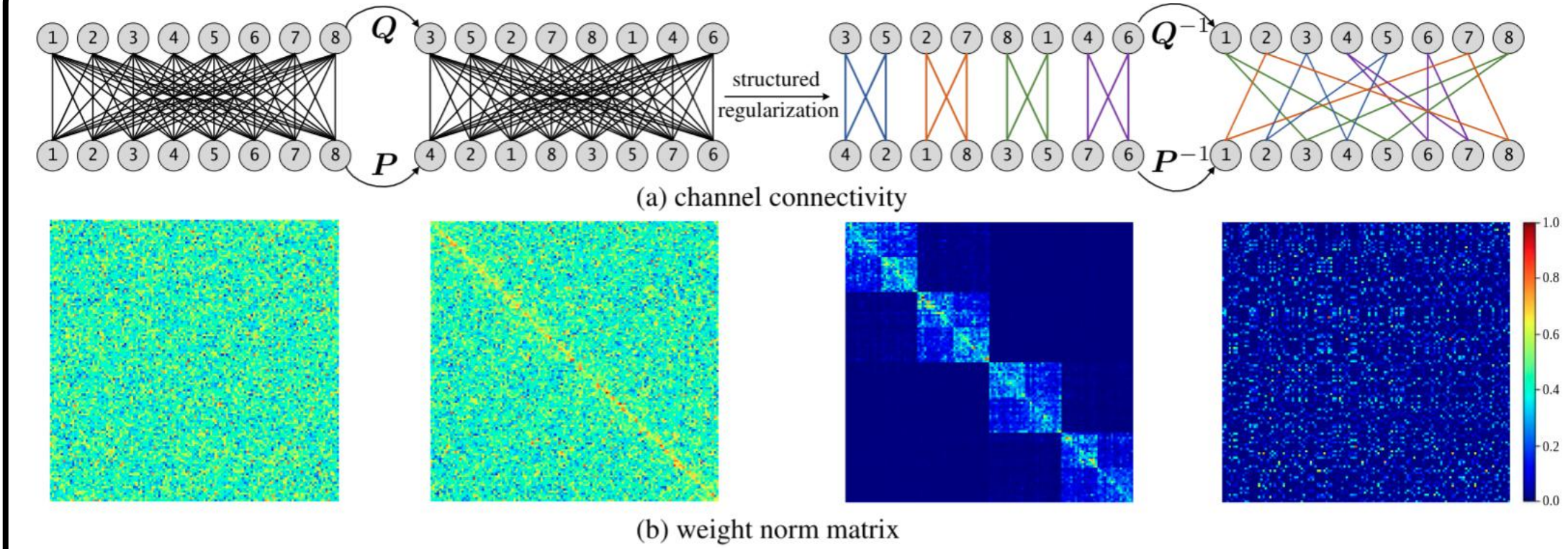


- ShuffleNet introduces a *channel shuffle* operation to re-arrange channels for each group, but the channel permutation still follows a pre-defined scheme (Fig. (c));



- This work proposes a *learnable channel shuffle* mechanism which unifies the norm-based criteria and the learning of channel permutation.

Overview



- Learn the permutation matrices (\mathbf{P} , \mathbf{Q}) and use them to perform channel shuffle;
- Structurally regularize the convolutional weights to induce a group structure;
- Remove the sparsified weights and form the GroupConv;
- Shuffle back to the original ordering and obtain a *structurally sparse representation of the convolutional weights*.

Challenges

- Under what criteria to learn the channel shuffle (\Rightarrow linear programming);
- How to induce a group structure on the convolutional weights (\Rightarrow structured regularization).

Learning of Connectivity

Observations.

- Weight norm matrix of GroupConv \rightarrow *block-diagonal* matrix;
- Channel shuffle \rightarrow *row/column permutation* of weight norm matrix.

Formulation.

- Since weight norm matrix cannot be permuted to be an exact block-diagonal one, the aim of channel shuffle is to make it “*as block-diagonal as possible*”:

$$\begin{aligned} \min_{\mathbf{P}, \mathbf{Q}} \mathbf{P} \mathbf{S} \mathbf{Q} \otimes \mathbf{R} \\ \text{s.t. } \mathbf{P} \in \mathcal{P}^{\text{C}_{\text{out}}} \text{ and } \mathbf{Q} \in \mathcal{P}^{\text{C}_{\text{in}}}, \end{aligned} \quad (1)$$

where \mathbf{S} is the weight norm matrix, i.e., $S_{ji} = \|W_{j,i,:}\|$, \mathbf{R} is a cost matrix, \mathcal{P}^N denotes the set of $N \times N$ permutation matrices, and \otimes denotes the operator of element-wise multiplication and summation over all entries.

Solution.

- Since problem (1) is NP-hard, we adopt two techniques to make it solvable:
 - the set \mathcal{P}^N of permutation matrices is relaxed to its convex hull --- the set of *doubly-stochastic matrices*, i.e., the *Birkhoff polytope*:

$$\mathcal{B}^N = \{\mathbf{X} \in \mathbb{R}_+^{N \times N} : \mathbf{X} \mathbf{1}_N = \mathbf{1}_N, \mathbf{X}^T \mathbf{1}_N = \mathbf{1}_N\};$$
 - the variables \mathbf{P} and \mathbf{Q} in (1) are updated alternatively, namely, the *coordinate descent* algorithm is used;

- With the two tricks, problem (1) is adapted as (when updating \mathbf{P}):

$$\begin{aligned} \min_{\mathbf{P}} \mathbf{P} \otimes \mathbf{R} \mathbf{Q}^T \mathbf{S}^T \\ \text{s.t. } \mathbf{P} \in \mathcal{B}^{\text{C}_{\text{out}}}; \end{aligned} \quad (2)$$

- In (2), the objective function is linear in \mathbf{P} , and \mathcal{B}^N is a *simplex*, so (2) is essentially a *linear program*, which can be solved by the *network simplex* method.

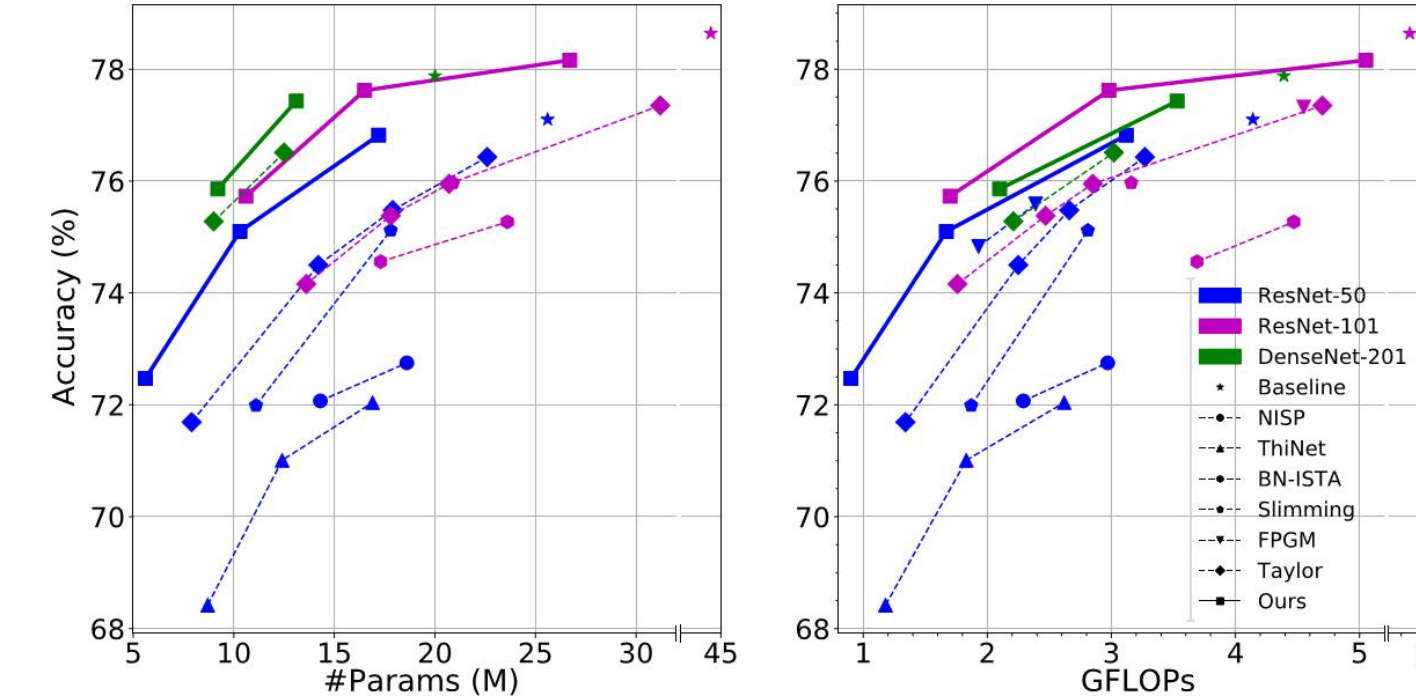
Discussion.

- By linear programming theory and Birkhoff’s theorem, at least one solution of (2) is exactly a permutation matrix;
- Therefore, feasible region \mathcal{B}^N is naturally reduced to \mathcal{P}^N .

Experiments

- Comparison with SOTA pruning methods on ImageNet;
- Trade-off between accuracy and complexity;

Methods	#Params.(10 ⁶) ↓	GFLOPs ↓	Acc.(%) ↑
ResNet-50			
Baseline	25.6	4.14	77.10
Slimming-20%	17.8	2.81	75.12
Taylor-19%	17.9	2.66	75.48
StrucSpars-35%	17.2	3.12	76.82
Taylor-28%	14.2	2.25	74.50
StrucSpars-65%	10.3	1.67	75.10
Taylor-44%	7.9	1.34	71.69
Slimming-50%	11.1	1.87	71.99
StrucSpars-85%	5.6	0.90	72.47



- Impact of channel shuffle mechanism.

Config.	ResNet-50-65%		ResNet-101-65%	
Acc.	Top-1	Top-5	Top-1	Top-5
FINETUNE	75.10	92.52	77.62	93.72
FROMSCRATCH	75.02	92.46	77.14	93.53
SHUFFLENET	74.97	92.41	76.91	93.38
RANDOM	69.45	89.45	73.16	91.44
NOshuffle	73.30	91.39	75.31	92.64

Structured Regularization

- Despite channel shuffle, the group structure cannot be formed spontaneously.

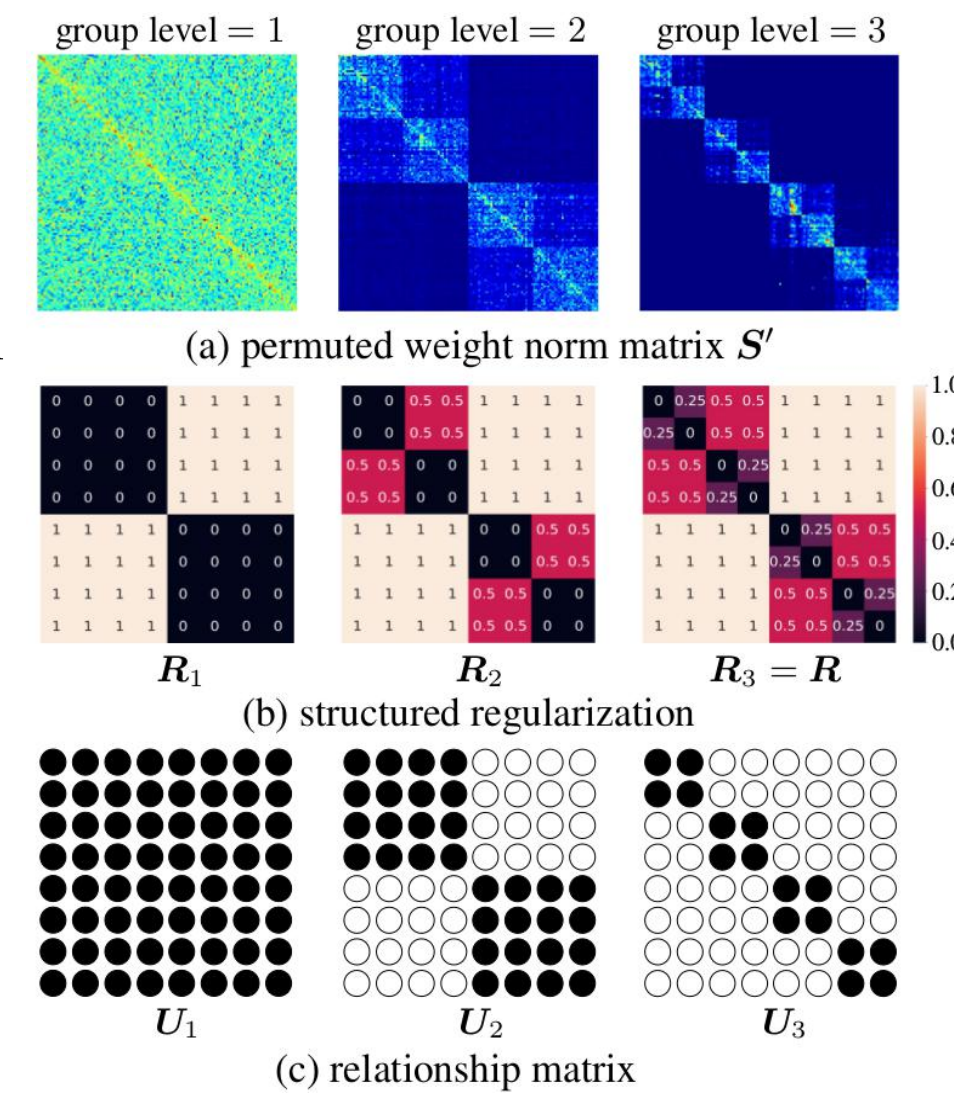
Structured Regularization.

- We impose *structured L_1 regularization* on the permuted weight norm matrix:

$$\mathcal{L}_{\text{reg}} = \mathbf{S}' \otimes \mathbf{R}_g,$$

where $\mathbf{S}' = \mathbf{P} \mathbf{S} \mathbf{Q}$ is the permuted weight norm matrix (Fig. (a)), \mathbf{R}_g denotes the structured regularization matrix (Fig. (b));

- Refer to our paper for details of \mathbf{R}_g and its relationship to the cost matrix \mathbf{R} .



Grouping Criteria.

- The grouping criteria is as follows:

$$g = \max \{ g : \mathbf{S}' \otimes \mathbf{U}_g \geq p \mathbf{S}' \otimes \mathbf{U}_1, g = 1, 2, \dots \},$$

where \mathbf{U}_g is the relationship matrix (Fig. (c)).

Future Perspective

- **Data-Driven Structured Sparsification.**

Can the weight norm fully represent weight importance to accuracy? How about making the sparsification process be guided by the data loss? We suggest *optimization-based meta-learning* techniques.

- **Progressive Sparsification.**

Spasified weights are removed progressively during training, leading to *finetune-free* training pipeline.

- **Combination with Filter Pruning.**

Filter pruning is beneficial within a strict memory budget. In order to jointly prune filters and learn a group structure, we suggest *group sparsity constraints*.