

1 Supplementary Material

2 1 Additional Studies

3 Besides the evaluation of various models on different datasets, we also perform additional studies to
4 obtain deep understandings of our proposed channel independence-based filter pruning approach.

5 1.1 Relationship between Channel Independence and Importance of Feature Map

6 We use a numerical example to demonstrate the relationship between Channel Independence (CI)
7 and importance of feature maps. Here for the following example 3×4 matrix, each of its rows
8 denotes one vectorized feature map of one channel. Our goal is to identify the least important row
9 that can be represented by other rows. Intuitively, Row-1 or Row-2 should be removed due to their
10 linear dependence. Furthermore, because the l_2 -norm of Row-2 is less than that of Row-1, Row-2 is
11 expected to be the least important one.

$$\begin{pmatrix} 0.9 & 0.8 & 1.1 & 1.2 \\ 0.81 & 0.72 & 0.99 & 1.08 \\ 0.8 & 0.9 & 1.2 & 1.1 \end{pmatrix} \quad (1)$$

Now according to Equation 3, we can obtain the CI of each row as shown in Table 1:

Table 1: CI of each row.

CI of Row-1	0.696
CI of Row-2	0.549
CI of Row-3	0.827

12

13 And it is seen that Row-2 is assigned as the smallest CI, which is consistent with our expectation.

14 1.2 Balance between Pruning and Task Performance

15 In the context of model compression, high pruning rate and high accuracy cannot be always achieved at
16 the same time – an efficient compression approach should provide good balance between compression
17 performance and task performance. Fig. 1 shows the change of accuracy of the pruned ResNet-50
18 on ImageNet dataset via using our approach with respect to different pruning ratios. It can be seen
19 that our approach can effectively reduce the number of model parameters and FLOPs with good
20 performance on test accuracy.

21 1.3 Accuracy-Pruning Rate Trade-off Curves of Different Pruning Methods

22 We study the accuracy-pruning rate trade-off curves of different pruning methods (CHIP, SCOP,
23 HRank) for ResNet-50 on ImageNet. The results are shown in Fig. 2.

24 1.4 Quantified Sensitiveness of Channel Independence to Input Data

25 To analyze the potential sensitiveness of channel independence to input data (as indicated in **Question**
26 **#3**), Fig. 4 in the main paper visualizes the average channel independence with different batches of
27 input images to show that the channel independence is not sensitive to the change of inputs. In this
28 supplementary material, we further quantify the sensitiveness. To be specific, for each batch of input
29 data (batch size = 128), we form a length-64 vector consisting of the average channel independence
30 for all the 64 feature maps of one layer (ResNet-56_55) in ResNet-56 model on CIFAR-10 dataset,
31 and then we calculate the **Pearson correlation coefficient** among different channel independence
32 vectors that correspond to different batches. As shown in Table 2, those vectors are highly correlated
33 with each other though they are generated from different input batches, thereby demonstrating the
34 low sensitiveness of channel independence metric to the input data.

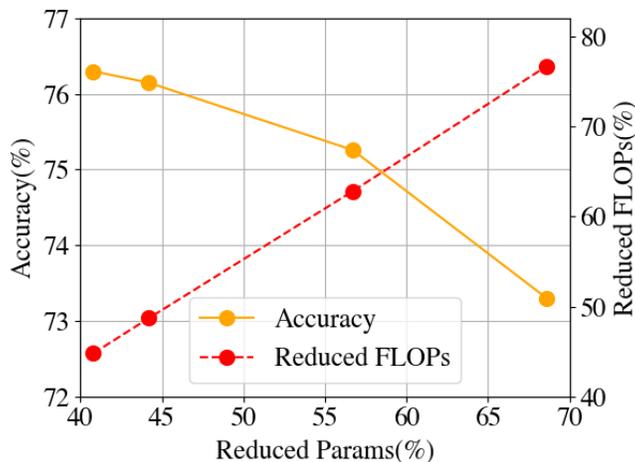


Figure 1: The accuracies and computational costs of our pruned ResNet-50 model with respect to different pruning ratios (on ImageNet dataset).

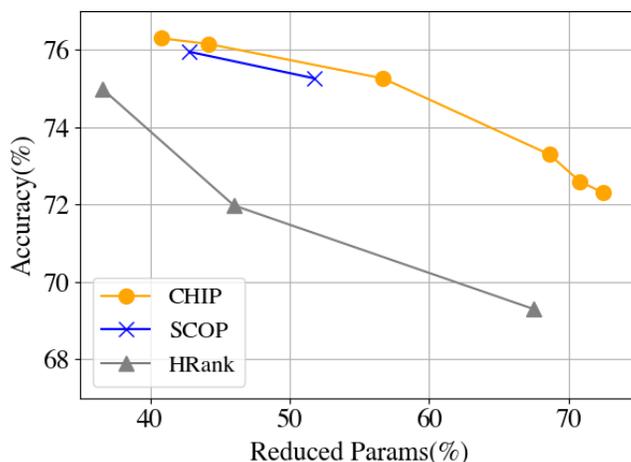


Figure 2: The accuracies of ResNet-50 model from different methods (CHIP, SCOP, HRank) with respect to different pruning ratios (on ImageNet dataset).

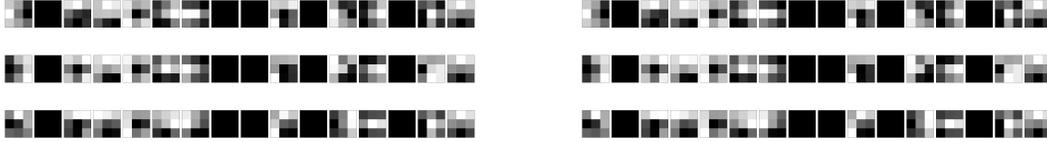
35 1.5 Is Additional Adjustment of Importance Ranking Needed?

36 As analyzed in **Question #4**, a potential extension of our approach is to introduce an additional
 37 phase to further adjust the importance ranking from the training data, once our one-shot channel
 38 independence-based pruning is finished. To be specific, an even better channel-wise pruning mask
 39 strategy could be further learned built upon the mask determined by our approach as the initialization.
 40 Intuitively, this data-driven strategy might potentially provide an extra performance improvement.

41 To explore this potential opportunity, we conduct experiments for different models on different
 42 datasets. Our empirical observation is that an additional learning phase for the pruning mask does
 43 not bring an extra accuracy increase. Fig. 3 visualizes the same part of filters in Conv1 layer of
 44 VGG-16 without and with additional pruning mask training. It is seen that there is nothing change for
 45 the selected filters to be pruned before and after using the trained mask. Our experiments for other
 46 models on other datasets also show the same phenomenon. Therefore we conclude that additional
 47 adjustment on the pruning mask is not required in the context of our channel independence-based
 48 filter pruning.

Table 2: Pearson correlation coefficient among 5 length-64 different channel independence vectors of ResNet-56_55 layer (containing 64 output feature maps) with 5 different input batches (CIFAR-10 dataset).

-	Vector-1	Vector-2	Vector-3	Vector-4	Vector-5
Vector-1	1	0.907	0.850	0.911	0.821
Vector-2	0.907	1	0.880	0.899	0.901
Vector-3	0.850	0.880	1	0.913	0.913
Vector-4	0.911	0.899	0.913	1	0.881
Vector-5	0.821	0.901	0.913	0.881	1



(a) Visualization of filters without further pruning mask adjustment.

(b) Visualization of filters with further pruning mask adjustment.

Figure 3: Visualization of filters in Conv1 layer of VGG-16 model on CIFAR-10 dataset. Here we only show the first 16 out of 64 filters of this layer due to the space limitation. **Left:** the pruned filters using our approach. **Right:** the pruned filters after further pruning mask adjustment with the mask determined by our approach as initialization. x-axis represents different filters and y-axis represents different input channels. The kernel size is 3×3 . Black kernels are the pruned ones.

49 How to find the best combination of the largest $CI(\{A_{b_i}^l\}_{i=1}^m)$? Given one image randomly sampled
50 from total images, metricized feature maps $\{A_{b_i}^l\}_{i=1}^m$ of l -th layer are generated after the inference.
51 Firstly, we calculate the CI upon our Algorithm 1. Then, we initial the score of M_{b_1, \dots, b_m}^l based
52 on the normalized CI . That is, if we are not going to train the M_{b_1, \dots, b_m}^l to change b_1, \dots, b_m , the
53 nuclear norm and index of pruned filters from $\{A_{b_i}^l\}_{i=1}^m$ equals to the result from our Algorithm
54 1. Therefore, this initialization can be viewed as a baseline for $CI(A_{b_i}^l)$. Secondly, we train the
55 M_{b_1, \dots, b_m}^l using the MSE loss functions to minimize the gap between the Upper Bound and current
56 nuclear norm under sparsity 83.3% from VGG-16. With optimizer of ADAM and SGD, the learning
57 rate is set from 0.1 to 0.001 and the weight decay is set from 0.05 to 5. Among each possible pair of
58 above hyperparameters, we get the pruned filters of maximal $CI(\{A_{b_i}^l\}_{i=1}^m)$ from what we desire.
59 To sum up, although there has not been proven theoretically, we find that index of pruned filter
60 generated from our method almost equals to index of pruned filters from global optimal methods in
61 experiments.

62 2 Detailed Setting of κ^l and Pruning Ratios

63 In this section, we provide the details of κ^l (number of preserved filters) and pruning ratios of all
64 layers. On CIFAR-10, we report the κ^l and pruning ratios for ResNet-56, ResNet-110 and VGG-16.
65 On ImageNet, κ^l and pruning ratios are reported for ResNet-50.

66 2.1 κ^l (Number of Preserved Filters of All Layers)

67 2.1.1 ResNet-56

68 **For overall sparsity 42.8%, layer-wise κ^l are :** [16, 9, 13, 9, 13, 9, 13, 9, 13, 9, 13, 9, 13,
69 9, 13, 9, 13, 19, 27, 19, 27, 19, 27, 19, 27, 19, 27, 19, 27, 19, 27, 19, 27, 19, 27, 19, 27, 38, 64, 38, 64, 38, 64,
70 38, 64, 38, 64, 38, 64, 38, 64, 38, 64, 38, 64]

115 **2.2.2 ResNet-110**

116 **For overall sparsity 48.3%, layer-wise pruning ratios are :** [0.0, 0.35, 0.22, 0.35, 0.22, 0.35,
117 0.22, 0.35, 0.22, 0.35, 0.22, 0.35, 0.22, 0.35, 0.22, 0.35, 0.22, 0.35, 0.22, 0.35, 0.22, 0.35,
118 0.22, 0.35, 0.22, 0.35, 0.22, 0.35, 0.22, 0.35, 0.22, 0.35, 0.22, 0.35, 0.22, 0.45, 0.22, 0.45, 0.22, 0.45,
119 0.22, 0.45, 0.22, 0.45, 0.22, 0.45, 0.22, 0.45, 0.22, 0.45, 0.22, 0.45, 0.22, 0.45, 0.22, 0.45, 0.22, 0.45,
120 0.22, 0.45, 0.22, 0.45, 0.22, 0.45, 0.22, 0.45, 0.22, 0.45, 0.22, 0.45, 0.22, 0.45, 0.22, 0.45, 0.0, 0.45, 0.0, 0.45,
121 0.0, 0.45, 0.0, 0.45, 0.0, 0.45, 0.0, 0.45, 0.0, 0.45, 0.0, 0.45, 0.0, 0.45, 0.0, 0.45, 0.0, 0.45, 0.0, 0.45,
122 0.0, 0.45, 0.0, 0.45, 0.0, 0.45, 0.0, 0.45, 0.0, 0.45, 0.00]

123 **For overall sparsity 68.3%, layer-wise pruning ratios are :** [0.0, 0.5, 0.4, 0.5, 0.4, 0.5, 0.4, 0.5,
124 0.4, 0.5, 0.4, 0.5, 0.4, 0.5, 0.4, 0.5, 0.4, 0.5, 0.4, 0.5, 0.4, 0.5, 0.4, 0.5, 0.4, 0.5, 0.4, 0.5, 0.4,
125 0.5, 0.4, 0.5, 0.4, 0.5, 0.4, 0.65, 0.4, 0.65, 0.4, 0.65, 0.4, 0.65, 0.4, 0.65, 0.4, 0.65, 0.4, 0.65,
126 0.4, 0.65, 0.4, 0.65, 0.4, 0.65, 0.4, 0.65, 0.4, 0.65, 0.4, 0.65, 0.4, 0.65, 0.4, 0.65, 0.4, 0.65,
127 0.4, 0.65, 0.0, 0.65, 0.0, 0.65, 0.0, 0.65, 0.0, 0.65, 0.0, 0.65, 0.0, 0.65, 0.0, 0.65, 0.0, 0.65, 0.0, 0.65,
128 0.0, 0.65, 0.0, 0.65, 0.0, 0.65, 0.0, 0.65, 0.0, 0.65, 0.0, 0.65, 0.0, 0.65, 0.0, 0.65, 0.0, 0.65, 0.0]

129 **2.2.3 VGG-16**

130 **For overall sparsity 81.6%, layer-wise pruning ratios are :** [0.21, 0.21, 0.21, 0.21, 0.21, 0.21,
131 0.21, 0.75, 0.75, 0.75, 0.75, 0.75, 0]

132 **For overall sparsity 83.3%, layer-wise pruning ratios are :** [0.3, 0.3, 0.3, 0.3, 0.3, 0.3, 0.3, 0.75,
133 0.75, 0.75, 0.75, 0.75, 0]

134 **For overall sparsity 87.3%, layer-wise pruning ratios are :** [0.45, 0.45, 0.45, 0.45, 0.45, 0.45,
135 0.45, 0.78, 0.78, 0.78, 0.78, 0.78, 0]

136 **2.2.4 ResNet-50**

137 **For overall sparsity 40.8%, layer-wise pruning ratios are :** [0.0, 0.35, 0.35, 0.1, 0.35, 0.35, 0.1,
138 0.35, 0.35, 0.1, 0.35, 0.35, 0.1, 0.35, 0.35, 0.1, 0.35, 0.35, 0.1, 0.35, 0.35, 0.1, 0.35, 0.35, 0.1, 0.35,
139 0.35, 0.1, 0.35, 0.35, 0.1, 0.35, 0.35, 0.1, 0.35, 0.35, 0.1, 0.35, 0.35, 0.1, 0.35, 0.35, 0.0, 0.35, 0.35,
140 0.0, 0.35, 0.35, 0.0]

141 **For overall sparsity 44.2%, layer-wise pruning ratios are :** [0.0, 0.38, 0.38, 0.12, 0.38, 0.38,
142 0.12, 0.38, 0.38, 0.12, 0.38, 0.38, 0.12, 0.38, 0.38, 0.12, 0.38, 0.38, 0.12, 0.38, 0.38, 0.12, 0.38, 0.38,
143 0.12, 0.38, 0.38, 0.12, 0.38, 0.38, 0.12, 0.38, 0.38, 0.12, 0.38, 0.38, 0.12, 0.38, 0.38, 0.12, 0.38, 0.38,
144 0.0, 0.38, 0.38, 0.0, 0.38, 0.38, 0.0]

145 **For overall sparsity 56.7%, layer-wise pruning ratios are :** [0.0, 0.5, 0.5, 0.25, 0.5, 0.5, 0.25,
146 0.5, 0.5, 0.25, 0.5, 0.5, 0.25, 0.5, 0.5, 0.25, 0.5, 0.5, 0.25, 0.5, 0.5, 0.25, 0.5, 0.5, 0.25, 0.5, 0.5, 0.25,
147 0.5, 0.5, 0.25, 0.5, 0.5, 0.25, 0.5, 0.5, 0.25, 0.5, 0.5, 0.25, 0.5, 0.5, 0.25, 0.5, 0.5, 0.0, 0.5, 0.5, 0.0, 0.5, 0.5, 0.0]

148 **For overall sparsity 68.6%, layer-wise pruning ratios are :** [0.0, 0.6, 0.6, 0.5, 0.6, 0.6, 0.5, 0.6,
149 0.6, 0.5, 0.6, 0.6, 0.5, 0.6, 0.6, 0.5, 0.6, 0.6, 0.5, 0.6, 0.6, 0.5, 0.6, 0.6, 0.5, 0.6, 0.6, 0.5, 0.6, 0.6, 0.5,
150 0.6, 0.6, 0.5, 0.6, 0.6, 0.5, 0.6, 0.6, 0.5, 0.6, 0.6, 0.0, 0.6, 0.6, 0.0, 0.6, 0.6, 0.0, 0.6, 0.6, 0.0]

151 **Checklist**

- 152 1. For all authors...
- 153 (a) Do the main claims made in the abstract and introduction accurately reflect the paper's
154 contributions and scope? [Yes] Please see "Technical Preview and Contribution" in
155 Section ??
- 156 (b) Did you describe the limitations of your work? [Yes] In Section ?? we indicate that
157 the calculation of channel independence of multiple feature maps is an approximated
158 method to save computational cost.
- 159 (c) Did you discuss any potential negative societal impacts of your work? [No] Model
160 compression can bring more energy-efficient deployment of deep learning, hence there
161 are not potential negative societal impacts.
- 162 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
163 them? [Yes] We have read the ethics review guidelines and follow them when preparing
164 the paper.
- 165 2. If you are including theoretical results...
- 166 (a) Did you state the full set of assumptions of all theoretical results? [N/A] This paper is
167 not a theoretical paper.
- 168 (b) Did you include complete proofs of all theoretical results? [N/A] This paper does not
169 contain theoretical proof.
- 170 3. If you ran experiments...
- 171 (a) Did you include the code, data, and instructions needed to reproduce the main experi-
172 mental results (either in the supplemental material or as a URL)? [Yes] The code is in
173 the github link.
- 174 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
175 were chosen)? [Yes] Please see "Experimental Setting" in Section ?? and Supplemental
176 materials
- 177 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
178 ments multiple times)? [No] We indeed have measured the error bars after multiple
179 runs. However, because 1) our results are very stable with respect to different random
180 seeds; and 2) The compared state-of-the-art works in Section ?? do not report error
181 bars, we do not list ours to be consistent with their reporting.
- 182 (d) Did you include the total amount of compute and the type of resources used (e.g., type
183 of GPUs, internal cluster, or cloud provider)? [Yes] Please see "Experimental Setting"
184 in Section ??
- 185 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 186 (a) If your work uses existing assets, did you cite the creators? [Yes] We cite the creators
187 of ImageNet and CIFAR-10 datasets in the Reference.
- 188 (b) Did you mention the license of the assets? [N/A] ImageNet and CIFAR-10 are public
189 datasets.
- 190 (c) Did you include any new assets either in the supplemental material or as a URL? [No]
191 We do not curate or release new assets.
- 192 (d) Did you discuss whether and how consent was obtained from people whose data you're
193 using/curating? [N/A] ImageNet and CIFAR-10 are public datasets and they are free to
194 download.
- 195 (e) Did you discuss whether the data you are using/curating contains personally identifiable
196 information or offensive content? [No] No personal identifiable or offensive content is
197 included in ImageNet or CIFAR-10 datasets.
- 198 5. If you used crowdsourcing or conducted research with human subjects...
- 199 (a) Did you include the full text of instructions given to participants and screenshots, if
200 applicable? [N/A] This paper is not related to research with human subjects.
- 201 (b) Did you describe any potential participant risks, with links to Institutional Review
202 Board (IRB) approvals, if applicable? [N/A] This paper is not related to research with
203 human subjects.

204
205
206

(c) Did you include the estimated hourly wage paid to participants and the total amount spent on participant compensation? [N/A] This paper is not related to research with human subjects.