

1 We sincerely thank the anonymous reviewers for their support and constructive comments.

2 **To Reviewer #1. Q1: Presentation.** A: Treatment group denotes the features to be researched while control group is the
 3 features for controlling irrelevant factors, which are instantiated as real and knockoff features respectively (Line 81-90
 4 in the main paper). We will refine the presentation and add a dedicated section for revisiting related-works.

5 **Q2: Different results in the baseline.** A: We collect the baseline results from the published papers and their test errors of
 6 the ‘original’ networks are slightly different. We re-implement the competing methods under the same baseline, and the
 results of ResNet-56 on CIFAR-10 with mean/std are shown below. More results will be included in the final version.

Method	Original Error (%)	Pruned Error (%)	Gap (%)	Params. ↓ (%)	FLOPs ↓ (%)
GAL (2019) [17]	6.30 ± 0.24	6.91 ± 0.14	0.61	42.4	50.0
SCP (Ours)	6.30 ± 0.24	6.36 ± 0.09	0.06	56.3	56.0

7 **Q3: Effectiveness of knockoffs.** A: Results of ResNet-50 on ImageNet
 8 using noise or knockoffs are shown below. Since ImageNet is more
 9 complex, e.g., 1,000 categories and images of high resolution (224 ×
 10 224), utilizing noise data is hard to obtain good results. More results
 11 will be included in the final version.

Method	Error (%)	Gap (%)	FLOPs ↓ (%)
No control	26.22	2.37	54.6
Noise	25.86	2.01	54.6
Ours with bias	24.74	0.89	54.6

12 **Q4: The term ‘filter’.** A: There are M filters in Eq. (1), and each of them is an $N \times k \times k$ tensor, where $k \times k$ is the
 13 kernel size (e.g., 3×3) and N is the number of input channels. We will refine the presentation around the definitions.

14 **Q5: Efficiency.** A: Learning based approaches (e.g., CP [8], GAL [17]) are compared in Table 1 and 2 in the main paper.
 15 These methods re-train the original network to learn the importance of filters. Compared with them, our method fixes
 16 the network weights and only tunes the control scales when discovering redundant filters, which is more efficient. The
 17 practical consuming time of pruning ResNet-56 (no fine-tuning) on CIFAR-10 is shown below (A V100 GPU).

Method	CP	GAP	Ours
Time (min)	46	25	16

18 **To Reviewer #2. Q1: Clarity issues.** A: 1) The term ‘scales’ denotes
 19 the magnitudes of β_l and $\tilde{\beta}_l$, which will be replaced with ‘scaling
 20 factors’ and defined at the beginning of the paper. 2) We adopt similar loss function for the discriminator as Knockoff-
 21 gan [9] and propose to aggregate multiple elements for efficiently generating knockoff data, which will be fully described
 22 in the final version. The novel contributions will also be discussed at the end of introduction as your suggestion.

23 **Q2: Pruning procedure in practice.** A: Unimportant filters in each layer will be pruned. As discussed in the main body
 24 (Line 184-186), for an arbitrary convolutional layer, filters with small ($\beta_j^l - \tilde{\beta}_j^l$) will be recognized as redundancy.

25 **Q3: Concept of knockoff data.** A: Specifically, knockoff data do not contain real objects of any category (e.g., goldfish,
 26 snail) in the real dataset, as shown in Figure 3 of the main body. We will add more explanations in the final version.
 27 **Q4: Limitation.** A: Thanks for this nice concern. 1) Control scales are distributed in range [0,1] as shown in Figure R1.
 28 2) It does not matter when scales are close to 0.5/0.5, since we focus on sorting scales of different filters, rather than β^l
 29 and $\tilde{\beta}^l$ of the same filter.

30 **To Reviewer #3. Q1: Figure 1.** A: We will replace the art picture in Fig. 1 with the actually generated knockoff data.
 31 **Q2: Example for swapping operation.** A: Suppose that $\mathcal{A}=[0.1,0.18,-$
 32 $0.1], \tilde{\mathcal{A}}=[0.13,0.16,-0.15]$, and then $[\mathcal{A}, \tilde{\mathcal{A}}]=[0.1,0.18,-0.1,0.13,0.16,-0.15]$. If
 33 $\tilde{S} = \{2\}$, the swapped feature $[\mathcal{A}, \tilde{\mathcal{A}}]_{\text{swap}(\tilde{S})}=[0.1,0.16,-0.1,0.13,0.18,-0.15]$.

34 **To Reviewer #4. Q1: Optimization procedure.** A: Eq. (6) is the general formu-
 35 lation for feature selection. In practice, we use Eq. (8) to avoid ℓ_0 -norm and then
 36 Adam optimizer can be applied.
 37 **Q2: $\beta^l, \tilde{\beta}^l$ in the feature selection layer.** A: Thanks for this constructive comment.
 38 To have an explicit understanding, we illustrate the change of β^l (i.e., the first
 39 conv layer in last stage of ResNet-56 on CIFAR-10) during the optimization in
 40 Figure R1. Wherein, each curve denotes the control scale ($\beta_j^l \in \beta^l$) of a specific
 41 convolution filter. Similarly, $\tilde{\beta}^l = 1 - \beta^l$ has the opposite phenomenon. Since most of existing deep neural networks are
 42 of heavy design for the accuracy reason, they will not collapse to $\beta^l = 1$. Visualizations on $\beta^l, \tilde{\beta}^l$ and corresponding
 43 discussions will be included in the final version.

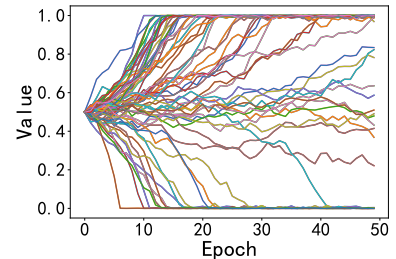


Figure R1: β^l w.r.t. epoch.

44 **Q3: Notions.** A: 1) We need to calculate knockoff feature $\tilde{\mathcal{A}}^{l+1}$, which is the $l+1$ -th layer’s feature map of the network
 45 with knockoff data as input, i.e., $\tilde{\mathcal{A}}^{l+1} = f^{l+1}(\tilde{X}, \mathcal{W}^{1:l+1})$. 2) Biases b^l and \tilde{b}^l are the modified terms from theoretical
 46 derivation, which ensure that the knockoff condition satisfies in the forward propagation of neural network from a
 47 theoretical perspective. These notions will be clarified more detailedly in the final version.
 48 **Q4: Different methods and backbone models.** A: We apply the proposed method on different backbone models to verify
 49 its generalization ability. The results of competing methods on different backbones are collected from their original
 50 papers for fair comparison. These methods have their own experimental settings and lack results on some backbones.
 51 For example, Hrank [16] did not report results on ResNet20, ResNet32 and MobileNetV2. Thus we do not include it.
 52 **Q5: Clarity.** A: Thanks, all of these typos and minor comments will be carefully fixed in the final version, and more
 53 discussions on the broader impact will be included.

54