

A. Implementation Details

In this section, we describe the implementation details of the compared methods.

As mentioned in Sec. 4, we implemented FSKD for the pruned ResNet models according to its official codes¹, and we re-ran the official CD codes² to get its results. We only replaced the pruned models while keeping the hyperparameters unchanged. As for BP, we also used SGD as the optimizer, and the initial learning rate, weight decay, and momentum were $1e-3$, $1e-4$, and 0.9 , respectively. For KD, we also used SGD and used the same learning rate, weight decay, and momentum as those in BP, in which we set the temperature $\tau = 2.0$ and the loss balancing factor was $\alpha = 0.7$.

B. Extra Results

In this part, we show results with fewer FLOPs pruned, results with connections pruned, and results on the CIFAR-10 dataset.

Pruning less FLOPs. As we reported in the main paper, we used ‘Prune-C Normal’ and ‘Prune-C Residual’ to prune the ResNet-34 model (*cf.* Tables 1 and 4, respectively). Here, we illustrated the ResNet-34 results of ‘Prune-B Normal’ (*cf.* Table 9) and ‘Prune-B Residual’ (*cf.* Table 10) on ILSVRC-2012. As shown in Table 9 and Table 10, MiR can outperform other methods regardless of pruning a large amount or a small amount of FLOPs, and the gap between MiR_{after} and MiR_{before} decreased when pruning less FLOPs.

Unstructured pruning results. We implemented our MiR to fine-tune the models which were pruned by an unstructured manner (*i.e.*, connection pruning). We used the ℓ_1 -norm weight pruning method to prune the less important weights and compared our MiR with the CD method under 90% weights pruned. Following the same settings in CD, we only pruned the weights in conv. layers, which means the weights in batch normalization layers and fully-connected layers were kept (but we can still update the weights in these layers). And 90% weights here means pruning 90% weights every conv. layer.

The representation ability was damaged because so much information was lost when pruning 90% weights with smaller ℓ_1 -norm. So we tried a progressive way to prune and fine-tune weights. We pruned 20% weights and then fine-tuned with 400 iterations until we got a 90% pruned sparse model. At last, we fine-tuned the 90% pruned model with 4000 iterations. As for the results of the CD method, we re-ran the official codes and reported the mean and std. of top-1 accuracy.

CIFAR-10 results. We also conducted ResNet-56 on the CIFAR-10 dataset following the same pruning setting

in CD [1]. ResNet-56 is a customized model for small datasets like CIFAR-10, and our pre-trained ResNet-56 model has 93.39%/99.87% of top-1/top-5 accuracy. We directly compared our MiR_{before} with FSKD [18], FitNet [24], ThiNet [23], CP [14] and CD [1]. As for our MiR, no extra augmentation was used except for random horizontal flip, and training settings were the same as mentioned in Sec. 4.

In Table 12, we used 1/2/3/5/10/50 samples per class to fine-tune the pruned models and reported the mean and std. of top-1 accuracy under five independent trials. Results marked with * were copied from CD. As shown in Table 12, our MiR worked well in all cases.

¹<https://github.com/LTH14/FSKD>

²<https://github.com/haolibai/Cross-Distillation>

Table 9. Mean and standard deviation of top-1/top-5 accuracy (%) on ILSVRC-2012. We used ‘Prune-B Normal’ to prune ResNet-34 and compared different methods with different training sizes. We used 50, 100, 500 random samples, and N -way- K -shot (N/K in the top row) settings. All the results were reported with five trials. **Bold** denotes the best results.

Methods	50	100	500	1000/1	1000/2	1000/3
BP	48.3±1.55/76.4±0.94	49.3±0.44/77.4±0.41	57.9±0.19/82.5±0.09	62.0±0.27/84.5±0.20	63.7±0.23/85.5±0.13	64.6±0.14/86.0±0.10
KD	52.7±1.43/78.8±0.99	54.0±0.53/80.1±0.47	60.3±0.10/83.8±0.09	62.5±0.08/84.9±0.08	63.4±0.19/85.4±0.08	63.7±0.10/85.6±0.07
FSKD	55.8±0.38/80.2±0.26	59.6±0.35/83.1±0.17	63.7±0.13/85.8±0.04	64.8±0.07/86.4±0.08	65.3±0.06/86.7±0.07	65.5±0.11/86.8±0.05
CD	62.7±0.28/85.1±0.19	62.8±0.25/85.2±0.15	67.1±0.06/88.0±0.05	67.5±0.10/88.2±0.04	67.8±0.10/88.4±0.06	68.1±0.10/88.5±0.05
MiR _{after}	65.7±0.09/87.3±0.07	66.6±0.07/87.8±0.11	68.3±0.07/88.7±0.08	68.9±0.03/89.1±0.04	69.3±0.09/89.2±0.06	69.5±0.05/89.4±0.05
MiR _{before}	67.5 ±0.13/ 88.3 ±0.06	68.1 ±0.13/ 88.7 ±0.07	69.2 ±0.05/ 89.3 ±0.08	69.7 ±0.06/ 89.5 ±0.03	69.9 ±0.04/ 89.7 ±0.07	70.0 ±0.07/ 89.8 ±0.03

Table 10. Mean and standard deviation of top-1/top-5 accuracy (%) on ILSVRC-2012. We pruned ResNet-34 using ‘Prune-B Residual’ (cf. Table 2). **Bold** denotes the best results.

Methods	50	100	500	1000/1	1000/2	1000/3
BP	36.9±0.94/66.3±1.04	39.7±0.36/69.0±0.25	51.3±0.21/77.5±0.18	57.2±0.22/81.0±0.13	59.6±0.19/82.8±0.14	60.8±0.15/83.6±0.06
KD	42.4±0.48/70.1±0.68	44.9±0.40/72.2±0.35	54.0±0.18/79.0±0.16	57.2±0.15/81.0±0.11	58.6±0.10/82.0±0.09	58.9±0.10/82.2±0.04
FSKD	46.0±0.49/72.2±0.46	50.6±0.19/76.3±0.14	55.9±0.25/80.2±0.15	57.2±0.09/81.2±0.14	57.8±0.11/81.6±0.13	58.0±0.05/81.7±0.08
MiR _{after}	61.1±0.20/84.2±0.24	62.9±0.18/85.5±0.16	66.2±0.12/87.5±0.07	67.3±0.06/88.2±0.07	68.0±0.07/88.6±0.05	68.3±0.03/88.7±0.06
MiR _{before}	64.9 ±0.25/ 86.6 ±0.21	66.2 ±0.10/ 87.5 ±0.12	68.2 ±0.12/ 88.7 ±0.05	68.8 ±0.05/ 89.1 ±0.06	69.3 ±0.06/ 89.3 ±0.03	69.5 ±0.06/ 89.5 ±0.03

Table 11. The top-1 accuracy of non-structured pruning. We pruned every layer except the first conv., and the sparsity is 0.9. **Bold** denotes the best results.

	Methods	50	100	500	1000/1	1000/2	1000/3
0.9	CD	48.7±0.48	53.9±0.07	59.1±0.19	60.3±0.11	61.3±0.03	61.7±0.08
	MiR _{before}	52.7 ±0.53	55.5 ±0.56	60.9 ±0.29	63.2 ±0.11	64.3 ±0.13	64.7 ±0.02

Table 12. Mean and standard deviation of top-1 accuracy (%) on CIFAR-10. We pruned ResNet-56 with ‘Res-50%’ used in CD [1] and used 1/2/3/5/10/50 samples per class for tuning. Results were reported with five trials. Methods marked with * were copied from CD. **Bold** denotes the best results.

Methods	1	2	3	5	10	50
FSKD*	84.26±1.42	85.79±1.31	85.99±1.29	87.53±1.06	88.15±0.71	88.70±0.55
FitNet*	86.85±1.91	87.95±2.13	88.94±1.85	89.43±1.60	91.03±1.14	91.89±0.87
ThiNet*	88.40±1.26	88.76±1.18	88.95±1.19	89.54±0.84	90.36±0.76	90.89±0.49
CP*	88.53±1.37	88.69±1.09	88.79±0.94	89.39±0.80	89.91±0.69	90.45±0.43
CD*	89.00±1.59	89.45±1.43	89.56±1.32	90.14±1.19	90.82±0.79	91.24±0.33
MiR _{before}	89.27 ±0.21	90.43 ±0.12	90.70 ±0.15	91.14 ±0.23	91.57 ±0.14	92.16 ±0.11