



Paper



Code

# Practical Network Acceleration with Tiny Sets

Guo-Hua Wang Jianxin Wu

State Key Laboratory for Novel Software Technology, Nanjing University



## 1. Introduction

In this paper, we propose an algorithm named PRACTISE to accelerate deep networks with tiny training sets.

● PRACTISE: practical network acceleration with tiny sets

● Deep networks: ResNet, MobileNet, etc

● Tiny training sets: 50~1000 images

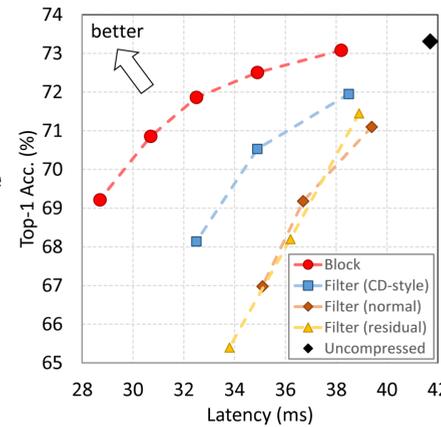
➤ Key property

● High latency-accuracy performance

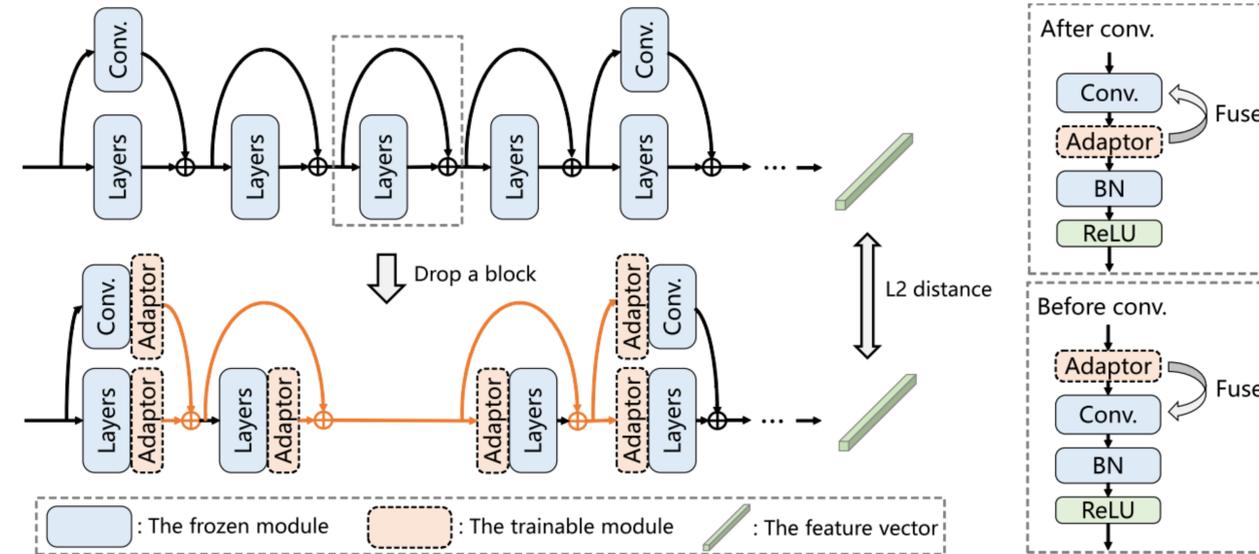
● Small training sets

● Robust to out-of-domain images

● Fast training (1.5 hours)



## 3. Method



### Algorithm 1: PRACTISE

**Input:** The original model  $\mathcal{M}_O$ , the number of dropped blocks  $k$ , the tiny training data  $\mathcal{D}_T$

Test the latency of  $\mathcal{M}_O$ ;

**for each block  $\mathcal{B}_i$  do**

Drop  $\mathcal{B}_i$  to obtain the pruned model  $\mathcal{M}_{P(\mathcal{B}_i)}$ ;

Test latency of  $\mathcal{M}_{P(\mathcal{B}_i)}$  and find  $\tau(\mathcal{B}_i)$  (Eq. 2);

Insert adaptors;

Compute  $\mathcal{R}(\mathcal{B}_i)$  with  $\mathcal{D}_T$  (Eq. 1);

Compute the score  $s(\mathcal{B}_i)$  (Eq. 3);

Add  $\mathcal{B}_i$  back and remove all adaptors;

Choose the top  $k$  blocks with the minimum scores;

Drop these  $k$  blocks to obtain  $\mathcal{M}_P$ ;

Finetune  $\mathcal{M}_P$  with  $\mathcal{D}_T$  by minimizing  $\mathcal{L}$  (Eq. 4);

**return** The pruned model  $\mathcal{M}_P$

$$\mathcal{R}(\mathcal{B}_i) = \min_{\alpha} \mathbb{E}_{x \sim p(x)} \|\mathcal{M}_O(x; \theta) - \mathcal{M}_{P(\mathcal{B}_i)}(x; \theta \setminus \mathcal{B}_i, \alpha)\|_F^2 \quad (1)$$

$$\tau(\mathcal{B}_i) = \frac{lat_{\mathcal{M}_O} - lat_{\mathcal{M}_{P(\mathcal{B}_i)}}}{lat_{\mathcal{M}_O}} \quad (2)$$

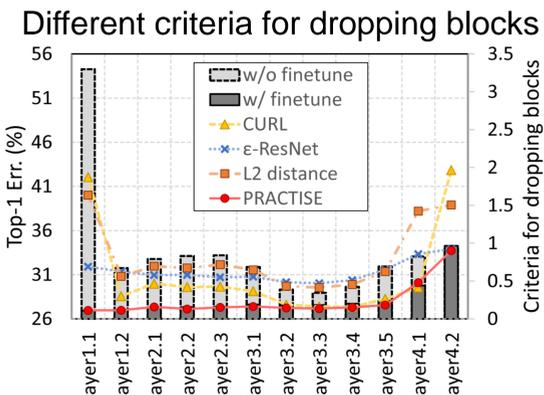
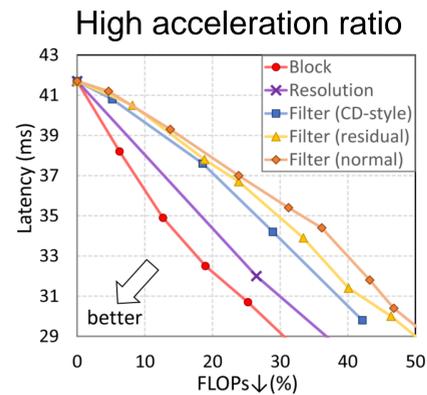
$$s(\mathcal{B}_i) = \frac{\mathcal{R}(\mathcal{B}_i)}{\tau(\mathcal{B}_i)} \quad (3) \quad \mathcal{L} = \|\mathcal{M}_O(x; \theta_O) - \mathcal{M}_P(x; \theta_P)\|_F^2 \quad (4)$$

## 2. Motivation

➤ **Latency-accuracy** vs. FLOPs-accuracy: more practical

➤ **Drop blocks**: higher acceleration ratio, more convex optimization

➤ **Recoverability**: predict the finetuned acc. to prune blocks



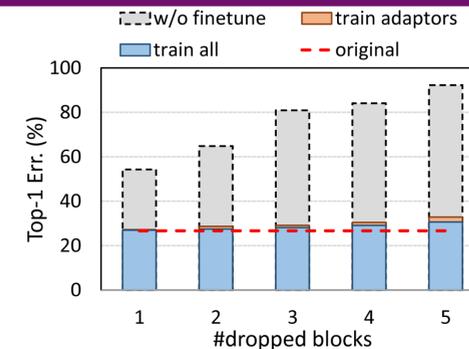
## 4. Experiments

### Different criteria for dropping blocks

#Dropped blocks	1	2	3	4	5
Random	71.12	71.27	67.89	64.27	63.47
CURL [21]	72.33	71.08	69.14	65.48	64.97
$\epsilon$ -ResNet [34]	72.51	71.20	69.00	67.90	64.75
L2 distance	72.51	71.20	69.00	68.61	64.93
PRACTISE	<b>73.02</b>	<b>72.52</b>	<b>71.92</b>	<b>70.86</b>	<b>69.35</b>

### ImageNet results

Method	Latency (ms)	50	100	500	1000
BP (filter)	33.8 (18.9% ↓)	24.2±0.92/52.7±1.36	27.6±0.41/56.7±0.62	42.9±0.28/70.5±0.27	51.2±0.32/76.5±0.16
BP (block)	<b>32.5 (22.1% ↓)</b>	<b>60.6±0.62/83.5±0.42</b>	<b>61.6±0.31/84.3±0.36</b>	<b>65.0±0.19/86.5±0.20</b>	<b>66.8±0.18/87.5±0.13</b>
KD [10]	33.8 (18.9% ↓)	30.1±0.69/57.7±1.10	33.1±0.43/61.0±0.53	45.7±0.26/72.2±0.25	50.5±0.29/75.9±0.23
FSKD [12]	33.8 (18.9% ↓)	31.1±0.90/56.5±1.10	36.6±0.44/63.1±0.46	42.8±0.49/69.1±0.58	44.9±0.20/70.5±0.29
MiR [30]	33.8 (18.9% ↓)	59.9±0.30/83.2±0.31	62.1±0.22/84.8±0.18	65.4±0.07/87.0±0.03	66.6±0.05/87.7±0.04
PRACTISE	<b>32.5 (22.1% ↓)</b>	<b>68.0±1.36/88.2±0.77</b>	<b>70.4±0.42/89.7±0.23</b>	<b>71.8±0.07/90.5±0.02</b>	<b>71.9±0.05/90.6±0.04</b>



## 5. Contributions & Conclusions

➤ Argue that the FLOPs-accuracy tradeoff is a misleading metric for few-shot compression, and advocate that the **latency-accuracy** tradeoff is more crucial in practice.

➤ The first to reveal **dropping blocks** great potential in few-shot compression.

➤ Propose a new concept **recoverability** to measure the difficulty of recovering each block, and in determining the priority to drop blocks.

➤ Propose **PRACTISE**, an algorithm for accelerating networks with tiny sets of images.

➤ The extraordinary performance: For 22.1% latency reduction, PRACTISE surpasses the previous state-of-the-art (SOTA) method on average by 7.0%.