

# Few Shot Network Compression via Cross Distillation

Haoli Bai<sup>1</sup>, Jiaxiang Wu<sup>2</sup>, Irwin King<sup>1</sup>, Michael Lyu<sup>1</sup>  
 {hlbai, king, lyu}@cse.cuhk.edu.hk, jonathanwu@tencent.com

<sup>1</sup>The Chinese University of Hong Kong, <sup>2</sup>Tencent AI Lab

Looking for  
internship!



## Introduction

Most prevalent network compression methods require fine-tuning with sufficient training data to ensure accuracy, which could be challenged by privacy and security issues. Network compression with few shot training instances is a new direction for research.

**Our Work:** We propose cross distillation, a novel network compression approach specialized for few shot training samples. The proposed method offers a general framework compatible with pruning or quantization.

## Methods

### Overall Framework

#### Task Definition

Given an over-parameterized teacher network  $\mathcal{F}^T$ , our goal is to learn a compact student network  $\mathcal{F}^S$ . We proceed in a layer-wise manner:

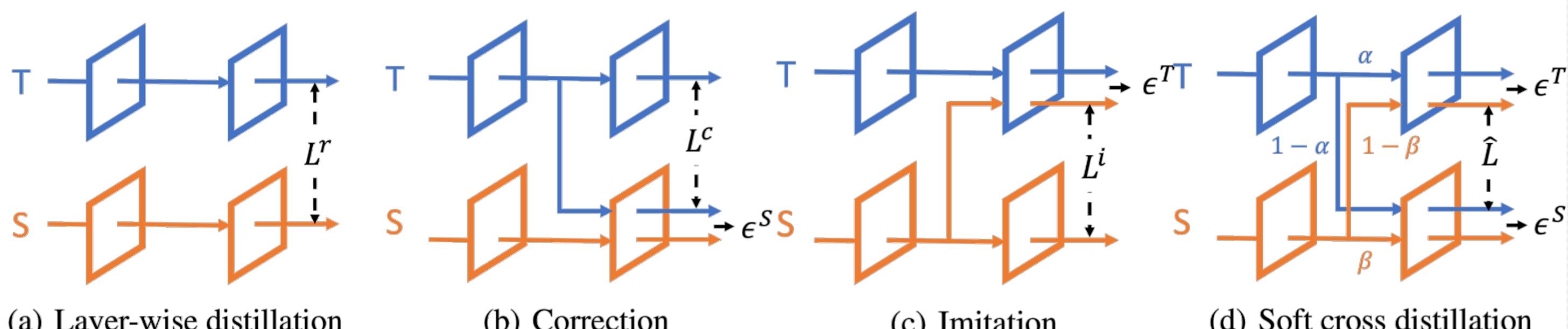
$$\mathbf{W}_*^S = \arg \min_{\mathbf{W}^S} \frac{1}{N} \mathcal{L}^r(\mathbf{W}^S) + \lambda \mathcal{R}(\mathbf{W}^S),$$

where  $\mathcal{L}^r(\mathbf{W}^S) = \|\sigma(\mathbf{W}^T * \mathbf{h}^T) - \sigma(\mathbf{W}^S * \mathbf{h}^S)\|_F^2$  is the **estimation error**, and  $\mathcal{R}(\mathbf{W}^S)$  is some regularization for sparsity/quantization.

⚠ Fewer training samples → larger estimation errors;

⚠ Estimation errors accumulate and propagate layer-wisely.

#### Cross Distillation



#### (b) Correction loss – reduce the historically accumulated errors

$$\mathcal{L}^c(\mathbf{W}^S) = \|\sigma(\mathbf{W}^T * \mathbf{h}^T) - \sigma(\mathbf{W}^S * \mathbf{h}^T)\|_F^2$$

#### (c) Imitation loss – teacher-aware accumulated errors on student

$$\mathcal{L}^i(\mathbf{W}^S) = \|\sigma(\mathbf{W}^T * \mathbf{h}^S) - \sigma(\mathbf{W}^S * \mathbf{h}^S)\|_F^2$$

#### (d) Soft cross distillation – a trade-off

$$\hat{\mathcal{L}}(\mathbf{W}^S) = \|\sigma(\mathbf{W}^T * \hat{\mathbf{h}}^T) - \sigma(\mathbf{W}^S * \hat{\mathbf{h}}^S)\|_F^2$$

#### (\* ) Convex combination between (b) and (c) – another trade-off

$$\tilde{\mathcal{L}} = \mu \mathcal{L}^c + (1 - \mu) \mathcal{L}^i, \quad \mu \in [0, 1]$$

### Combination with Pruning/Quantization

We adopt proximal mapping update for different  $\mathcal{R}(\mathbf{W}^S)$

$$\mathbf{W}_{t+1}^S = \text{Prox}_{\lambda \mathcal{R}}(\mathbf{W}_t^S - \eta \nabla \tilde{\mathcal{L}}(\mathbf{W}_t^S)) \in R^{c_o \times c_i \times k \times k}$$

#### Structured Pruning

$$\mathcal{R}(\mathbf{W}^S) = \|\mathbf{W}^S\|_{2,1} = \sum_i \|\mathbf{W}_i^S\|_2$$

$$\text{Prox}_{\lambda \|\cdot\|_2}(\mathbf{W}_i^S) = \max(1 - \frac{\lambda}{\|\mathbf{W}_i^S\|_2}, 0) \cdot \mathbf{W}_i^S$$

#### Unstructured Pruning

$$\mathcal{R}(\mathbf{W}^S) = \|\mathbf{W}^S\|_1 = \sum_{i,j,h,w} |W_{ijhw}^S|$$

$$\text{Prox}_{\lambda \|\cdot\|_1}(W_{ijhw}^S) = \begin{cases} W_{ijhw}^S - \lambda & W_{ijhw}^S > \lambda \\ 0 & |W_{ijhw}^S| \leq \lambda \\ W_{ijhw}^S + \lambda & W_{ijhw}^S < -\lambda \end{cases}$$

#### Quantization

$$\mathcal{R}(\mathbf{W}^S) = \Pi_Q(g(\mathbf{W}^S)) \quad Q = \{0, \frac{\pm 1}{2^{B-1}-1}, \frac{\pm 2}{2^{B-1}-1}, \dots, \pm 1\}$$

$$\text{Prox}_{\Pi_Q}(W_{ijhw}^S) = 0 \text{ if } W_{ijhw}^S \in Q \text{ else } \infty$$

### Theoretical Analysis

**Theorem 1.** Suppose both  $\mathcal{F}^T$  and  $\mathcal{F}^S$  are  $L$ -layer convolutional neural networks followed by the un-pruned softmax fully-connected layer. If the activation functions  $\sigma(\cdot)$  are Lipschitz-continuous such as ReLU(), the gap of softmax cross entropy  $\mathcal{L}^{ce}$  between the network logits  $\mathbf{o}^T = \mathcal{F}^T(\mathbf{x})$  and  $\mathbf{o}^S = \mathcal{F}^S(\mathbf{x})$  can be bounded by

$$|\mathcal{L}^{ce}(\mathbf{o}^T; \mathbf{y}) - \mathcal{L}^{ce}(\mathbf{o}^S; \mathbf{y})| \leq C \tilde{\mathcal{L}}_L + \sum_{l=1}^{L-1} \prod_{k=l}^L C'_k(\mu) \tilde{\mathcal{L}}_l, \quad (6)$$

where  $C$  and  $C'(\mu)$  are constants and  $C'(\mu)$  is linear in  $\mu$ .

## Experiments

### Main Results

#### Structured Pruning with ResNet-34 on ImageNet

Methods	50	100	500	1	2	3
L1-norm	72.94±0.00	72.94±0.00	72.94±0.00	72.94±0.00	72.94±0.00	72.94±0.00
BP	83.18±1.86	84.32±1.29	85.34±0.89	85.76±0.73	86.05±0.51	86.29±0.56
FSKD	82.53±1.52	84.58±1.13	86.67±0.78	87.08±0.76	87.23±0.52	87.20±0.43
FitNet	86.86±1.81	87.12±1.63	87.73±0.96	87.66±0.84	88.61±0.76	<b>89.32±0.78</b>
ThiNet	85.67±1.57	85.54±1.39	86.97±0.89	87.42±0.76	87.52±0.68	87.53±0.50
CP	86.34±1.24	86.38±1.37	87.41±0.80	88.03±0.66	87.98±0.49	88.21±0.37
Ours-NC	86.51±1.71	86.61±1.20	87.92±0.75	87.98±0.60	88.63±0.49	88.82±0.38
Ours	86.95±1.59	87.60±1.13	88.34±0.69	88.17±0.73	88.57±0.40	88.59±0.41
<b>Ours-S</b>	<b>87.42±1.69</b>	<b>87.73±1.17</b>	<b>88.60±0.82</b>	<b>88.40±0.61</b>	<b>88.84±0.48</b>	88.87±0.35

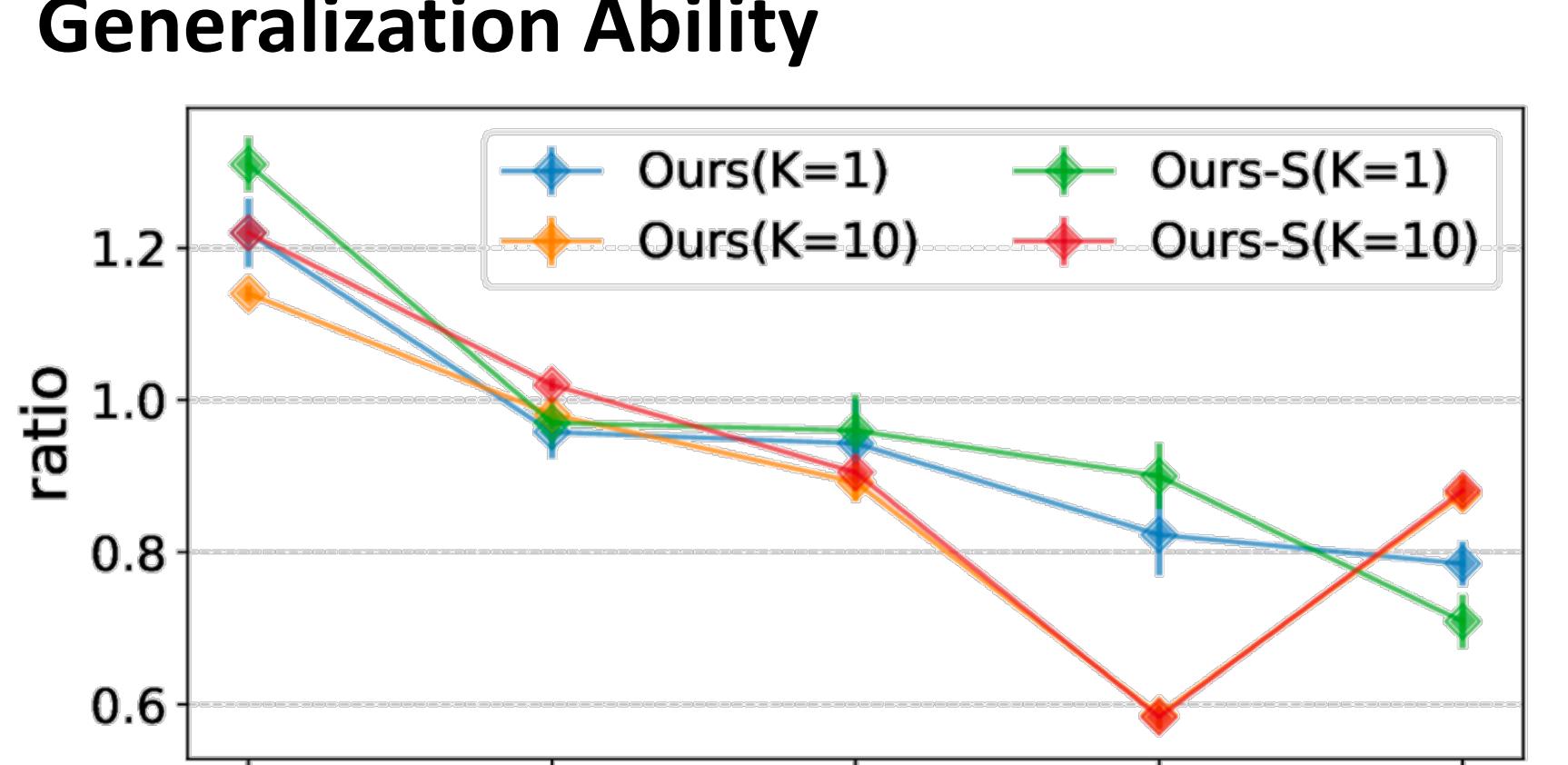
#### Unstructured Pruning with VGG-16 on ImageNet

Methods	50	100	500	1	2	3
L1-norm	0.5±0.00	0.5±0.00	0.5±0.00	0.5±0.00	0.5±0.00	0.5±0.00
BP	42.87±2.07	48.78±1.43	65.47±1.15	71.25±0.97	74.85±0.71	76.04±0.48
FitNet	52.66±2.93	57.09±2.14	76.59±1.45	80.14±1.23	82.27±0.70	83.14±0.51
Ours-NC	78.73±1.78	83.29±1.12	85.04±0.93	85.36±0.61	85.21±0.41	85.49±0.46
Ours	<b>83.81±1.49</b>	86.21±1.09	87.19±0.96	87.61±0.82	87.78±0.45	87.86±0.39
Ours-S	83.67±1.52	<b>86.72±1.23</b>	<b>87.82±1.04</b>	<b>88.14±0.74</b>	<b>88.23±0.61</b>	<b>88.38±0.43</b>

#### Weight Quantization with ResNet-56 on Cifar-10

	K=1		K=5	
	W2A32	W4A32	W2A32	W4A32
Ours-NC	72.48±1.94	85.75±0.96	84.67±1.89	91.09±0.37
Ours	<b>80.92±2.23</b>	<b>90.42±0.53</b>	<b>86.11±1.97</b>	<b>91.23±0.45</b>

#### Generalization Ability



#### Sensitivity Analysis

