

Learning in School: Multi-teacher Knowledge Inversion for Data-Free Quantization

Yuhang Li¹², Feng Zhu², Ruihao Gong²³, Mingzhu Shen², Fengwei Yu², Shaoqing Lu², Shi Gu¹

¹University of Electronic Science and Technology of China

²Sensetime Research, ³Beihang University

Abstract

User data confidentiality protection is becoming a rising challenge in the present deep learning research. In that case, data-free quantization has emerged as a promising method to conduct model compression without the need for user data. With no access to data, model quantization naturally becomes less resilient and faces a higher risk of performance degradation. Prior works propose to distill fake images by matching the activation distribution given a specific pre-trained model. However, this fake data cannot be applied to other models easily and is optimized by an invariant objective, resulting in the lack of generalizability and diversity whereas these properties can be found in the natural image dataset. To address these problems, we propose Learning in School (LIS) algorithm, capable to generate the images suitable for all models by inverting the knowledge in multiple teachers. We further introduce a decentralized training strategy by sampling teachers from hierarchical courses to simultaneously maintain the diversity of generated images. LIS data is highly diverse, not model-specific and only requires one-time synthesis to generalize multiple models and applications. Extensive experiments prove that LIS images resemble natural images with high quality and high fidelity. On data-free quantization, our LIS method significantly surpasses the existing model-specific methods. In particular, LIS data is effective in both post-training quantization and quantization-aware training on the ImageNet dataset and achieves up to 33% top-1 accuracy uplift compared with existing methods.

1. Introduction

To enable powerful deep learning models on the embedded and mobile devices without sacrificing performance, various model compression techniques have been discovered. For example, neural network quantization [12, 19, 24, 41] converts 32-bit floating-point models into low-bit fixed point models and benefits from the acceleration of fixed-

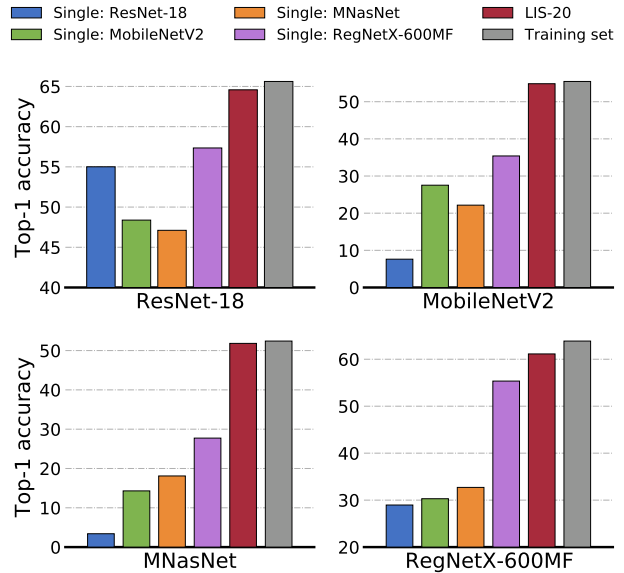


Figure 1. When data synthesized on one model, it has poor performance on other models’ quantization. Here different color bar indicates different data source, 4 different chart denotes the 4 different target model when performing W2A8 quantization.

point computation and less memory consumption. Network pruning [9, 14, 37] focuses on reducing the redundant neural connections and find a sparse network. Knowledge Distillation (KD) [18, 32] transfer the knowledge in the large teacher network to small student networks.

However, one cannot compress the neural networks aggressively without the help of data. As an example, most full precision models can be safely quantized to 8-bit by directly rounding the parameters to their nearest integers [23, 31]. However, when the bit-width goes down to 4, we have to perform quantization-aware training to compensate for the accuracy loss. Unfortunately, due to the privacy protection issue¹, one cannot get user data easily. Moreover, the whole ImageNet dataset contains 1.2M images (more

¹https://ec.europa.eu/info/law/law-topic/data-protection_en

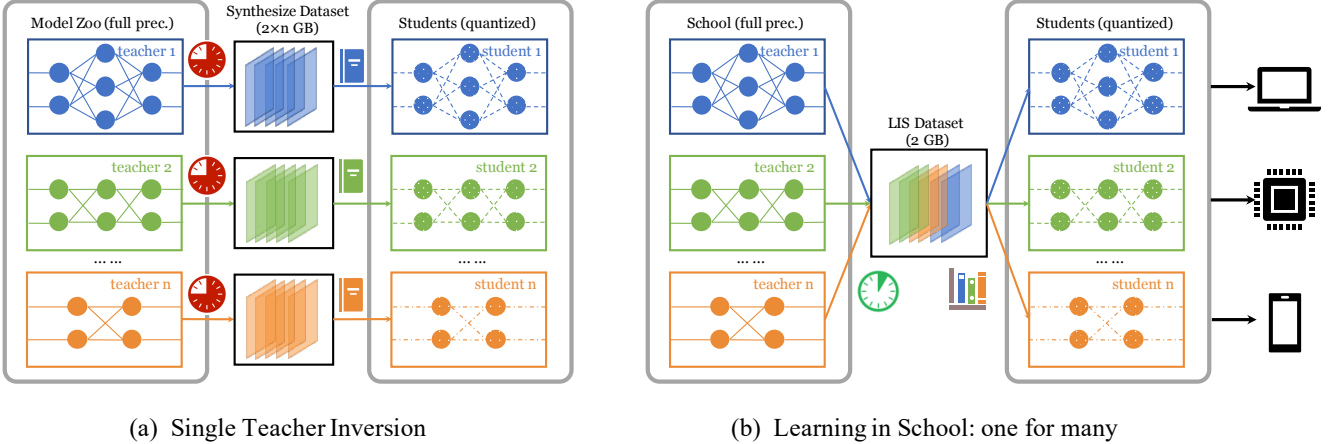


Figure 2. The overall pipeline of the proposed methods. (a) Conventional single teacher inversion needs to generate images one-by-one, which greatly delays the model production cycle. (b) LIS pipeline incorporates knowledge from every teacher and generates a one-for-many synthesized dataset, after which the compressed model can be rapidly deployed to users.

than 100 gigabytes), which consumes much more memory space than the model itself. Therefore, the data-free model quantization is more demanding now.

Recently, many works [5, 38, 15] succeed to *invert* images from a pretrained model (a.k.a. teacher-net). They try to match the activations’ distribution by comparing the recorded running mean and running variance in the Batch Norm [20] layer. However, these works put their focus on developing a better criterion of single teacher inversion. As a result, the synthesized images distilled by one teacher are over-fitted and cannot be extended to another. For example, in DeepInversion [38], synthesizing 215k 224×224 resolution images from ResNet-50v1.5 [1] requires 2.8k V100 GPU hours. These images cannot be used for another teacher directly. As we envisioned in Fig. 1, data inverted from ResNet-18 only achieves 3.4% accuracy on MNasNet quantization. Therefore, single teacher inversion requires additional thousands of GPU hours to adapt to another model’s quantization. Distilling only one model may also suffer from the diversity issue. Over 100k images are trained to match the same activation distribution, therefore these images could be overlapped in the loss landscape. Yin *et al.* [38] realize this problem and propose Adaptive Deep-Inversion. However, this method is teacher-student-specific since a unique teacher-student combination will generate a unique dataset, and therefore, lack of flexibility.

In this work, we propose **Learning in School (LIS)** knowledge inversion to generate images that can generalize different models and enjoy high diversity, which contributes to an overall improvement for data-free quantization. LIS conducts *one-time joint inversion of multiple teachers* based on the simple rule that the natural ImageNet dataset can be used to do quantization on any models. Therefore, the ideal synthesized images should generalize any target model’s

quantization. Notwithstanding this observation, distilling data from multiple teachers is a non-trivial task because the loss scale differs in each model and the different models extract the features differently. To guarantee the generalizability of the synthesized data, we first build a ‘school’ that teaches different types of knowledge, and then distill the images from miscellaneous teachers. To guarantee the diversity and stability of data, we arrange hierarchical courses and randomly sample two teachers each time from different courses in the school. Fig. 1 shows our LIS data consistently reaches the highest performance on each model’s quantization. The contributions in this paper are threefold:

1. **Generalizability:** We build a *school* that contains different aspect of knowledge and distill images with the help of multi-teacher. As a consequence, the data can generalize well to many models.
2. **Diversity:** We propose decentralized training and aggregation to optimize LIS data, therefore the data shares overwhelming diversity and behaves more like natural images with high-fidelity.
3. **Efficiency:** We show only a one-time synthesis of LIS data can be applied to various models. Experimental results indicate that LIS data not only performs better than model-specific data but also reduces the synthesis time on multi-model quantization.

2. Related Works

Data-Driven Model Compression Data is an essential requirement in model compression. For example, automated exploring compact neural architectures [42, 25] requires data to continuously train and evaluate sub-networks. Besides the neural architecture search, quantization is also a prevalent method to compress the full precision networks. For the 8-bit case where the weight quantization nearly does

not affect accuracy, only a small subset of calibration images is needed to determine the activation range in each layer, which is called Post-training Quantization [2, 23]. AdaRound [30] learns the rounding mechanism of weights and improve post-training quantization by reconstructing each layer outputs. Quantization-aware finetuning [10] can achieve near-to-original accuracy even when weights and activations are quantized to INT3. But this method requires a full training dataset as we mentioned. ZeroQ [5] and the *Knowledge Within* [15] use distilled dataset to perform data-free quantization, but their methods are model-specific, i.e., one generated dataset can only be used for one model’s quantization. Apart from quantization, knowledge distillation [18] is also widely explored. However, most literature does not focus on data-free scenarios.

Image Synthesis Prior to *inverting* images from a pre-trained model, there are many works studying generating high-fidelity and high-resolution images. For example, Generative Adversarial Networks (GAN) [13, 28] jointly train a generator and a discriminator to synthesize images. BigGAN [3] sets the state-of-the-art GAN by applying large-scale training and improves stability. Some works like [6, 7] apply GAN to generate training images and learn the student network. However, their work cannot extend to large-scale experiments (e.g. ImageNet) easily and requires a complicated process for training generators, which is not practical to conduct fast quantization. A parallel ax in image synthesis is *model inversion* [26, 27]. Mordvintsev *et al.* proposed DeepDream [29] to ‘dream’ objects’ features onto images from a single pre-trained model. Recently, DeepInversion [38] uses the BN statistics variable as an optimization metric to distill the data and obtain high-fidelity images. BN scheme has also achieved improvements in other tasks [5, 15]. Despite their notable progress, they ignore the true property of real images, i.e. diversity and generalizability on every model.

3. Motivation

In this section, we will briefly discuss the background of how to distill images from a single pretrained model, then we will raise two concerns about this method.

3.1. Single Teacher Inversion

Suppose a trainable image X with the size of $[w, h, c]$ (in ImageNet dataset [8], the size is $224 \times 224 \times 3$) and a pretrained network \mathcal{A} , the Inceptionism [29] can invert the knowledge by assigning the image with a random label Y , since the networks have already captured the class information. Using the cross-entropy loss, the image can be optimized by

$$\min_X L_{CE}(\mathcal{A}(X), Y). \quad (1)$$

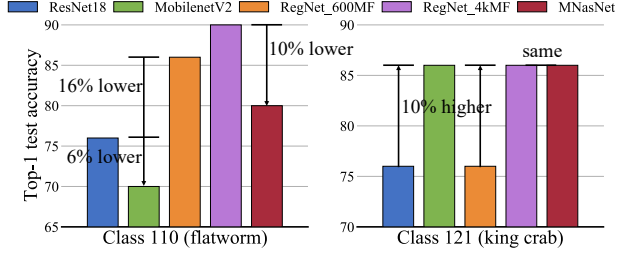


Figure 3. Different teacher has different expertise (or feature extractor). The test accuracy differs at two different class.

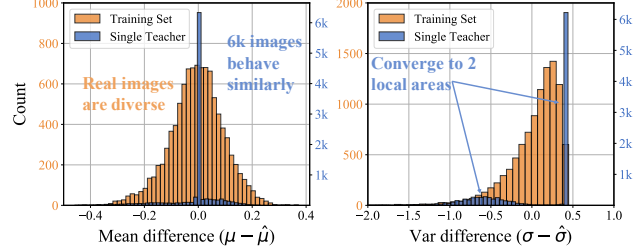


Figure 4. Difference between the distribution and the BN stored distribution on 10k images. Single teacher inversion data shares similar activation distribution and thus is not diverse.

Recently, [5, 15, 38] observe that the pretrained networks have stored the activation statistics in the BN layers (i.e. running mean and running variance for inference). Consequently, it is reasonable for synthesized images to mimic the activation distribution of the natural images in the network. Therefore, assuming the activation (regardless of the batch) in each layer is Gaussian distributed, the BN statistics loss can be defined as

$$L_{BN} = \sum_{i=1}^{\ell} (\|\mu_i(X) - \hat{\mu}_i\|_2 + \|\sigma_i^2(X) - \hat{\sigma}_i^2\|_2), \quad (2)$$

where $\mu_i(X)$ ($\sigma_i^2(X)$) is the mean (variance) of the synthesized images activation in (i) -th layer while $\hat{\mu}_i$ ($\hat{\sigma}_i^2$) is the stored running mean (variance) in the BN layer. Note that we can replace the MSE loss to Kullback-Leibler divergence loss as did in [15].

Last, a image prior loss is imposed on X to ensure the images are generally smooth. In [15], the prior loss is defined as the MSE between X and its Gaussian blurred version $\varepsilon(X)$. In this work, we use the prior loss defined in [38]: $L_{prior}(X) = \mathcal{R}_{TV}(X) + \lambda_{\ell_2} \mathcal{R}_{\ell_2}(X)$, which is the sum of variance and norm regularization. Combining these three losses, the final minimization objective of knowledge inversion for a single teacher can be formulated as:

$$\min_X \lambda_1 L_{CE}(X) + \lambda_2 L_{BN}(X) + \lambda_3 L_{prior}(X) \quad (3)$$

3.2. Insufficient Generalizability

For image synthesis tasks, the real ImageNet dataset could be viewed as the global minimum, which can be uti-

lized to perform model quantization on any neural architectures. However, we find that the data synthesized from one teacher cannot be directly applied to another different architecture. Example results demonstrated in Fig. 1 show that data synthesized on ResNet-18 gets bad quantization results (only 3.4%) on MNasNet.

We conjecture there might be two potential reasons. On the one hand, different teachers perceive different expertise (i.e. unique feature extractor) [4]. As Fig. 3 illustrated, ResNet-18 performs 6% better than MobileNetV2 in flatworm classification, but 10% worse in the king crab classification. If the synthesis only utilizes its own teacher’s expertise, the image quality cannot be guaranteed and it naturally cannot be easily applied to another. On the other hand, each teacher encodes the distribution characteristics of the original input images under its own feature space (i.e. unique BN statistics). Just inverting images from one model tends to overfit to the partial distribution information stored in this model. As a consequence, single teacher inversion generalizes poor in other networks and we have to repeatedly distill the data for each unique teacher network.

3.3. Low Dataset Diversity

Single teacher inversion also suffers from a low diversity issue. Denote the synthesized dataset as $\mathbf{X} = [X_1, X_2, \dots, X_N]$, we can find that all N images are synthesized from the same optimization objective, i.e. Eq. (3). If we are synthesizing more than 10k images, they inevitably converge to a concentrated area in the loss surface. In Fig. 4, we evaluate the difference between the activation mean (variance) and the running mean (variance) in BN of 10k images in ResNet-18 first layer’s first channel. It can be seen that the more than 60% synthesized images (by single ResNet-18) have almost the same mean value, which significantly restricts the diversity of the images. In comparison, the training set has a rather smooth and broad distribution. We think the inherent reason is all the N images share a similar convergence in single teacher inversion since they are under the same supervision. Therefore, we must increase image diversity to generate different kinds of images.

4. Methodology

In this section, we introduce the proposed *Learning in School* algorithm, which can improve both the generalizability and the diversity of the dataset.

4.1. Multi-teacher Inversion

Before introducing the proposed algorithm, we would like to discuss *how a deep learning model can generalize well with test data?* Consider the learnable parameters W in a deep network \mathcal{A} , we want to minimize the discrepancy

Table 1. Learning paradigm for model optimization and data optimization (knowledge inversion).

Method	Supervision	Training	Evaluation
<i>Model optimization</i>	Data	Optimize → Model	Generalize → Test data
<i>Data optimization</i>	Model	Optimize → Data	Generalize → Test models

between the model output $\mathcal{A}(X)$ and the ground truth label Y over the data distribution $P(X, Y)$. However, the real data distribution is unavailable in most practical cases, therefore we can only optimize the *empirical version* of the loss function [39], denoted by:

$$\min_W \frac{1}{N} \sum_{i=1}^N L_{CE}(\mathcal{A}(X_i), Y_i). \quad (4)$$

By optimizing the above objective, the network can learn the mapping from images to labels in a distribution spanned by these N training images. If these training images are sufficient to represent the data distribution and the network is well-optimized, the model can correctly predict the test data. To prevent over-fitting on the training data, data augmentation are widely used to expand the distribution of training data. We call this learning process as *model optimization* shown in Table 1.

Now, let us consider the inversion of the role between data and model. Instead of optimizing a model that can predict the classification of test images, we are optimizing an image that can have correct activation distribution and output in test models. Formally, we are interested in *how a single image can generalize well with different models?* The solution could be similar in Eq. (4), i.e., we can minimize the BN loss and CE loss over the model distribution. Model distribution, in most cases, is also unavailable like data distribution. As a result, we could leverage the empirical loss function by providing M ($M > 1$) pre-trained models as a school and then minimize the average loss:

$$\min_X \frac{1}{M} \sum_{i=1}^M (\lambda_1 L_{CE}(\mathcal{A}_i(X), Y) + \lambda_2 L_{BN}(X, \hat{\mu}_i, \hat{\sigma}_i)) + \lambda_3 L_{prior}(X). \quad (5)$$

As shown in Table 1, if the training models are sufficient to represent the model distribution, the optimized image can behave like natural image since all the pre-trained models are optimized by natural images. We call this process as *data optimization* (or knowledge inversion). In this paper, we build a training model set called *school*, which contains multiple pre-trained teachers, and the images can distill the knowledge from various teachers.

4.2. Decentralized Training and Aggregation

As shown in Fig. 5, multi-teacher inversion improves the generalizability across multiple models. Unfortunately,

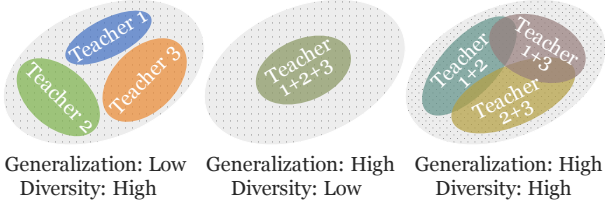


Figure 5. **Left:** Single teacher inversion produced insufficient generalized data, but aggregate these data can increase diversity. **Middle:** jointly optimizing all teachers will lead to one central area set which may degenerate the diversity. **Right:** decentralized training can ensure images won’t converge closely together and aggregation can span the local data across a wide range.

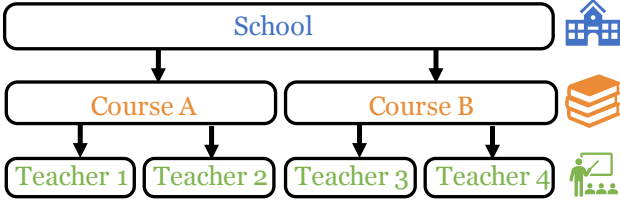


Figure 6. Hierarchical sampling of the teacher. We divide the teachers into groups (Course in the figure) based their architectures, then the teachers are sampled from different groups.

the synthesized images still have the diversity issue as we mentioned in Sec. 3.3. The multi-teacher inversion results in local area convergence since the supervision knowledge from teacher models (Eq. (5)) is invariant. Thus, the N synthesized images face similar convergence in optimization, leading to indistinguishable behaviors and a reduction in diversity.

To achieve both good generalizability and diversity of synthesized images, we propose decentralized training and aggregation operations to vary the supervision in knowledge inversion, e.g. use two different models to train X_1 and X_2 separately. First, for each batch of data, we randomly sample m different teachers in the school (M teachers in total) and use them to distill the data. This step is referred to as decentralized training (DT), because images are under variant supervision in the school. After DT, we can simply perform aggregation to collect all images. These images contain rich information and are highly diverse because they absorb the knowledge from miscellaneous teachers. The resulted images are shown in Fig. 5 right side, which converges to different areas since the supervision is variant. We find 2 teachers per-batch in DT is sufficient to achieve a good balance between generalizability and diversity. Increasing m brings little improvements and will greatly slow down the synthesis process. In section 5.2, we provide quantitative results of the impact of m .

Hierarchical Sampling. One potential drawback of DT is it contains partial knowledge in the school, which might not generalize well to all test models. To solve this, we adopt

Table 2. School with 20 teachers.

School					
Course-A Normal Conv. Model	Acc.-1	Course-B Group Conv. Model*	Acc.-1	Course-C Depthwise Conv. Model	Acc.-1
ResNet-18	71.04	RegNet-200MF	68.65	MobileNetV2-0.75	70.46
ResNet-34	74.03	RegNet-400MF	72.28	MobileNetV2-1.0	73.15
ResNet-50	77.01	RegNet-600MF	74.03	MobileNetV2-1.4	75.82
ResNet-50ad [†]	79.03	RegNet-800MF	75.23	MobileNetV2-2.0	77.81
ResNet-101ad	80.23	RegNet-1600MF	77.29	MNASNet-0.75	72.47
ResNet-152ad	80.87	RegNet-3200MF	78.72	MNASNet-1.0	74.02
		RegNet-4000MF	79.37	MNASNet-1.3	75.55

*We choose RegNetX without Squeeze-and-Extract module.

[†] Average down layer and deep stem layer as described in bag of tricks [17].

hierarchical sampling for teacher selection. Specifically, the teachers are first grouped to form a *course* based on their architecture family. For example, the ResNet-18 and ResNet-50 will be grouped in the ResNet course. Then, only one teacher is randomly sampled from one course to ensure the data can learn from multiple sources. This step is to prevent data from learning knowledge of the same course because we assume the knowledge in the same course could be similar. The hierarchical sampling process is demonstrated in Fig. 6.

In this work, we recruit 20 teachers in our school to generate the images learned in 3 courses: (1) *ResNet* [16] family with normal convolution, (2) *RegNet* [33] family with group convolution, (3) *MobileNet* and *MNASNet* [35, 36] family with depthwise-separable convolution. The detailed architectures information as well as the test Top-1 accuracy are reported in Table 2.

Adaptive Loss Weight. In decentralized training, a different batch of data will face different teachers, and we find the BN statistic loss will vary along with models. This is because the depth, width in each model are not identical. And it is not practical to manually set the loss weight (λ in Eq. 5) for each batch of data. To address this problem, we adopt a similar strategy in [21] to learn the loss weight by backpropagation. In particular, we first normalized each loss term to 1 after the first calculation of the loss function. Denote the normalized loss term as \hat{L} , the adaptive loss weight is formulated by

$$\min_{X, \alpha} = \sum_i \left(\frac{1}{\alpha_i^2} \hat{L}_i(X) + \alpha_i^2 \right), \quad (6)$$

where $\hat{L}_i(X)$ is the normalized loss term, e.g. BN statistic loss and cross-entropy loss, and α_i is the learnable loss weight to balance the loss function. For α , based on inequality of arithmetic and geometric means we know that $\frac{1}{\alpha^2} L + \alpha^2$ has a minimum $2\sqrt{L}$ when $\alpha^2 = \sqrt{L}$, therefore Eq. (6) can tune the loss weight based on the magnitude of each normalized loss term and prevent gradient domination. Finally, we can add a fixed constant λ to decide the overall importance between loss terms since the loss terms are normalized to $[0, 1]$. Together with DT, we formalize the LIS procedure in algorithm 1.

Algorithm 1: LIS Data Synthesis

Input: School (pretrained model zoo), number of batch N , loss weight scale λ , learnable loss weight α , subset size m for DT, training iterations T .

Initialize school, group teachers into different courses;

for all $i = 1, 2, \dots, N$ -th batch of images X_i **do**

 Initialize random label Y_i ;

 Initialize adaptive loss weight $\alpha = 1$;

 Randomly select m courses;

 Randomly sample one teacher for each course;

for all $t = 1, 2, \dots, T$ -iteration **do**

for all $j = 1, 2, \dots, m$ -th teacher model \mathcal{A}_j **do**

 Compute BN statistic loss $L_{\text{BN}}(X_i, \hat{\mu}_j, \hat{\sigma}_j)$;

 Compute CE loss $L_{\text{CE}}(\mathcal{A}_j(X_i), Y_i)$;

 Compute image prior loss $L_{\text{prior}}(X_i)$;

 Normalize each loss term and combine them;

 Descend final loss objective and update X_i, α

Aggregate all synthesized images (X_1 to X_N);

return LIS dataset (one-for-many)

5. Experiments

We conduct our experiments on the ImageNet dataset because generating high-fidelity 224×224 images is still a challenging task. We set the subset size of decentralized training $m = 2$ except we mention it. We use Adam [22] optimizer to optimize the images. Most hyper-parameters and implementation are aligned up with [38], such as multi-resolution training pipeline and image clip after the update. We optimize the images for 5k iteration and use a learning rate of 0.25 followed by a cosine decay schedule. α are optimized by SGD with a constant learning rate of $1e-3$. λ for $L_{\text{BN}}, L_{\text{CE}}, \mathcal{R}_{\text{TV}}, \mathcal{R}_{\ell_2}$ is set to $\{3, 2, 2, 0.1\}$. Training 100k LIS dataset images requires approximately one day on 16 NVIDIA TESLA V100 GPUs.

5.1. Analysis of the Synthesized Images

We have presented some qualitative evaluation in Fig. 7. To test the generalizability and diversity of the generated images, we report the average classification accuracy on three different courses. We also report the Inception Score (IS) [34] to evaluate the image quality and diversity.

Table 3. Classification accuracy evaluated at 20 different teachers and the Inception Score (IS) metric of the synthesized images.

Methods	Size	Acc.-A*	Acc.-B	Acc.-C	IS
DeepDream-R50 [29]	224	35.0	27.0	12.7	6.2
DeepInversion-R50 [38]	224	89.1	83.8	82.2	60.6
BigGAN [3]	256	-	-	-	178.0
SAGAN [40]	128	-	-	-	52.5
LIS-A	224	98.9	96.7	96.6	132.9
LIS-20	224	97.9	97.7	97.9	154.0

*Average Top-1 accuracy in Course-A.



Figure 7. Example images synthesized by 20 different teachers (Labels: oystercatcher, alp, otterhound, monastery, daisy, goldfish, cartoon, traffic light, bell cote, piano, syringe, red wine, mountain tent, padlock, park bench, quill pen).

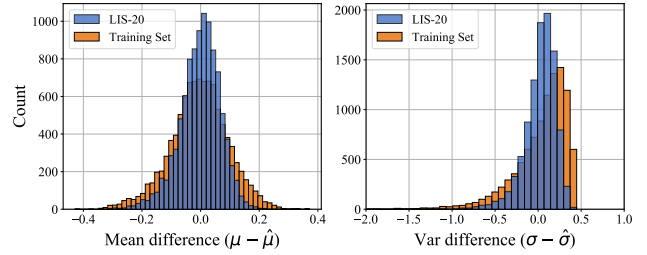


Figure 8. Difference between the activation mean (variance) and the BN running mean (variance). LIS-20 dataset is highly overlapped with the training dataset.

It can be noted from Table 3 that single teacher inversion (from ResNet-50) has lower accuracy in other courses. The proposed LIS-20 achieves near 98% accuracy in three courses. This means that our generated images have obvious class characteristics that all models can recognize. We also compare with some GAN-based image synthesis method where LIS can achieve comparable Inception Score. To see the effect of hierarchical sampling, we train LIS-A, a dataset learned only from teachers in Course-A. LIS-A can outperform LIS-20 in Course-A but generalize poor in other courses. Moreover, LIS-A has slightly lower IS than LIS-20 which verifies our assumption that multi-teacher can increase image diversity.

Last, we conduct the same experiments as did in Fig. 4. Fig. 8 shows that our LIS dataset has a similar distribution with the original real images, which confirms the positive effect of improving diversity in the proposed algorithm.

Table 4. ImageNet top-1 accuracy comparison on bias correction and adaptive rounding.

Model	Bias Correction + Calibration					Adaptive Rounding + Calibration				
	Precision	Training Set	DI [38]	LIS-A	LIS-20	Precision	Training Set	DI [38]	LIS-A	LIS-20
ResNet-18	6/6	70.25	70.24	70.24	70.22	4/8	70.57	70.30	70.48	70.63
FP:71.04	4/4	46.25	45.35	46.32	47.24	2/8	65.62	55.01	64.68	64.58
ResNet-50	6/6	76.02	75.82	75.85	76.18	4/8	76.37	75.90	76.25	76.38
FP:77.00	4/4	61.85	61.59	62.56	63.52	2/8	70.51	56.44	69.01	68.99
MobileNetV2	6/6	71.51	71.24	71.50	71.51	4/8	71.76	70.88	71.65	71.57
FP:72.49	4/4	23.56	22.58	21.37	34.00	2/8	55.45	27.53	54.03	54.84
MNASNet	6/6	67.57	69.12	65.28	68.11	4/8	72.57	69.97	72.36	72.54
FP:74.02	4/4	57.13	58.33	54.24	59.05	2/8	52.42	18.10	50.83	51.73
RegNetX-600MF	6/6	72.68	72.28	72.45	72.68	4/8	72.90	72.62	72.76	72.91
FP:73.71	4/4	40.95	28.28	34.42	45.97	2/8	63.87	55.35	60.36	61.13

Table 5. Classification accuracy of calibration.

Method	Precision	Acc.-1	Precision	Acc.-1
ResNet-18, ImageNet FP accuracy: 71.04				
Training set	W8A8	70.82	W6A6	69.37
DFQ [31]	W8A8	69.7	W6A6	66.3
ZeroQ [5]	W8A8	71.43/70.85*	W6A6	69.16*
KW [15]	W8A8	69.93/70.72*	W6A6	69.47*
LIS-20	W8A8	70.81	W6A6	69.60
MobileNetV2, ImageNet FP accuracy: 72.49				
Training set	W8A8	72.43	W6A6	70.02
DFQ [31]	W8A8	71.2	W6A6	-
ZeroQ [5]	W8A8	72.91/71.90*	W6A6	66.57*
KW [15]	W8A8	71.26/69.55*	W6A6	65.63*
LIS-20	W8A8	72.38	W6A6	69.93

*Self-implemented results.

5.2. Post-training Quantization

In this section, we utilize the images generated by LIS to conduct Post-Training Quantization (PTQ) on ImageNet, which requires a small amount of data. **Here we validate three types of PTQ methods, including Naïve Calibration, Bias Correction, and AdaRound.**

I. Naïve Calibration In this experiment, we directly determine the min-max range of activation with **the calibration data**. We sample one batch with 64 images to measure the activation range and fix it.

We compare ZeroQ [5] and *The Knowledge Within* [15] (abbreviated as KW in the following) on ResNet-18 and MobileNetV2. **To fairly compare results,** we implement ZeroQ and KW to the same pre-trained models. Results are demonstrated in Table 5. It can be seen that 8-bit quantization is easy to obtain FP-level performance. However, 6-bit quantization requires a higher quality of the synthesized images, and our LIS-20 data can improve up to **4.3%** accuracy compared to single teacher inversion methods.

II. Bias Correction Weight quantization will inevitably cause shifts in activation distribution. To address the bias

in the outputs, bias correction [11, 31] is proposed to absorb the quantization bias into the convolution layer bias. The absorption requires sufficient input images to correctly estimate the expectation of the output bias. Thus it can effectively evaluate the quality of generated images.

We evaluate 6/4-bit quantization with weights bias correction and activation calibration. The results in summarized in left part of Table 4. We also implement the single teacher inversion baseline DeepInversion [38] (**denoted by DI**) and align all training hyper-parameters. It can be seen that single teacher inversion generally performs well when bit-width is high, but the accuracy will drop in W4A4 situation. In most situations, the proposed *Learning in School* algorithm can achieve almost the same performances as real images and higher results than single teacher inversion (up to **17.7%** accuracy), demonstrating the high-quality and high-fidelity of the synthesized images. Interestingly, we observe that LIS-20 generally performs better than LIS-A, which demonstrates the diversity and the generalizability of multiple courses is higher than a single course.

III. Adaptive Rounding AdaRound [30] adopts the layer-by-layer feature reconstruction to optimize the rounding mechanism of weight quantization and thus also requires an accurate input distribution. We use 1024 images with a batch size of 32 to optimize the rounding of the weights. Each layer is optimized by 20k iterations. Training a quantized ResNet-18 only takes 20 minutes on a single GPU.

We use [38] as single teacher baseline. Table 6 shows that LIS obtains the same level as real images in W4A8 quantization. In W2A8, the difference between LIS and single teacher is much more significant. For example, LIS-20 outperforms single teacher by **12.55%** accuracy in ResNet-50 and **33.63%** accuracy in MNasNet. Compared to natural images, our LIS algorithm only degrades **1%~2%** accuracy.

Cross Validation To verify the generalizability of LIS, we conduct AdaRound experiments (2-bit weight quantization)

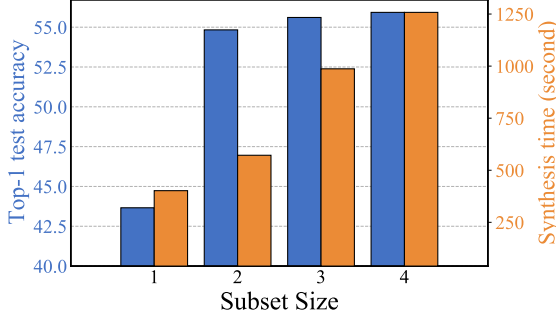


Figure 9. Ablation study on the effect of subset size in DT. We show that the 2 teachers per-batch is sufficient for generating high-quality images for data-free quantization.

on 4 different models. The results have been visualized in Fig. 1, we can find single teacher inversion performs poor on other models. For example, data synthesized on ResNet-18 can only achieve 7.46% accuracy on MobileNetV2 and 3.4% accuracy on MNasNet. However, we also observe that data generated by RegNet can *improve* the results of a single teacher, e.g. RegNet data improve 9.63% accuracy on MNasNet compared with its original teacher, demonstrating some single teacher can generate better images and therefore it is necessary to combine multi-teacher knowledge.

Ablation Study To evaluate the effect of sampling size in decentralized training, we conduct an ablation study on MobileNetV2 W2A8 quantization. The results are envisioned in Fig. 9, where we find that the synthesized images inverted from one teacher only obtain 43% accuracy. Two teachers per-batch can boost the performance to near 55% and is already close to the real images. Increasing sampling size will have little improvement in accuracy but lead to more and more synthesis time.

5.3. Quantization-aware Training

Quantization-aware Training (QAT) aims to recover the performance of the quantization neural network in low-bit scenarios. In this work, we leverage the state-of-the-art QAT baseline: Learned Step Size Quantization [10]. In QAT, Straight-Through Estimator (STE) is adopted to compute the gradients of the latent weights. In the case of LSQ, STE is also used to estimate the gradients of quantization step size, yielding an optimal tradeoff between clipping error and rounding error. We synthesize 100k images and use a batch size of 128 to finetune the quantization neural network. During QAT, the full precision model serves as the teacher and we use KL loss with temperature $\tau = 3$ as the criterion. We also incorporate the intermediate feature loss proposed in [15]. We finetune the quantized model for 44000 steps and it only takes 2 hours on 8 GTX 1080TI to complete the finetuning.

We perform W4A4 quantization-aware training on 5 models. Additionally, W2A4 quantization is applied to

Table 6. ImageNet top-1 accuracy comparison on QAT.

Model	Precision	Training Set	DI [38] [†]	KW [15] [‡]	LIS-20
ResNet-18	4/4	68.69	66.79	67.95	68.21
FP:71.04	2/4	58.81	52.54	-	56.78
ResNet-50	4/4	74.09	71.89	-	73.39
FP:77.00	2/4	65.20	52.88	-	60.34
MobileNetV2	4/4	63.42	59.93	-	63.60
FP:72.49	4/4*	-	-	66.07	67.74
MNasNet	4/4	62.80	56.16	-	62.43
FP:74.02					
RegNet-600MF	4/4	67.77	64.03	-	67.90
FP:73.71					

[†]Baseline method with aligned hyper-parameters, [‡] results quoted from paper.
* 1×1 convolutions are 8-bit.

Table 7. Accuracy comparison on data-free mixed precision quantization.

Model	Method	Precision	Model Size	Acc.-1
ResNet-18	ZeroQ [5]	MP/4	5.57 MB	68.21
FP:71.04	LIS-20	MP/4	5.57 MB	70.53
ResNet-50	ZeroQ [5]	MP/8	12.17 MB	76.08
FP:77.00	LIS-20	MP/8	12.17 MB	76.29
MobileNetV2	ZeroQ [5]	MP/8	1.67 MB	69.44
FP:72.49	LIS-20	MP/8	1.67 MB	71.66

ResNet-18 and 50. Results are presented in Table 6. Note that the real training dataset only contains 100k images. In W4A4, we show the LIS-20 dataset merely does not drop accuracy compared to the natural images. In W2A4 quantization, the gap between synthesized images and natural images is much bigger, our LIS can recover **4.2%** accuracy in ResNet-18, and **7.5%** accuracy in ResNet-50.

5.4. Mixed Precision Quantization

In this section, we validate our method in the mixed-precision quantization as studied in [5]. To obtain an optimal bit-width configuration, we need to measure the KL divergence between the quantized model and the full precision model. Each layer’s sensitivity is stored in a look-up table, then we can use dynamic programming to determine the mixed-precision settings. We quantized the network using adaptive rounding and the results are shown in Table 7. Data-free mixed precision quantization using LIS method can increase **2.2%** accuracy on ResNet-18 and MobileNetV2.

6. Conclusion

In this work, we identify two major drawbacks in single teacher inversion method, generalizability and diversity. The proposed Learning in School algorithm improves the existing method by leveraging the knowledge of multi-teacher and decentralized training. LIS algorithm is efficient as it only requires one-time synthesis to generalize many models. Experimental results demonstrate that LIS establishes a new state-of-the-art for data-free quantization.

References

- [1] NVIDIA resnet50v1.5 training. <https://github.com/NVIDIA/DeepLearningExamples/tree/master/PyTorch/Classification/ConvNets/resnet50v1.5>. Accessed: Nov-9-2020. **2**
- [2] Ron Banner, Yury Nahshan, and Daniel Soudry. Post training 4-bit quantization of convolutional networks for rapid-deployment. In *Advances in Neural Information Processing Systems*, 2019. **3**
- [3] Andrew Brock, Jeff Donahue, and Karen Simonyan. Large scale gan training for high fidelity natural image synthesis. *arXiv preprint arXiv:1809.11096*, 2018. **3, 6**
- [4] Han Cai, Chuang Gan, Ligeng Zhu, and Song Han. Tiny transfer learning: Towards memory-efficient on-device learning. *arXiv preprint arXiv:2007.11622*, 2020. **4**
- [5] Yaohui Cai, Zhewei Yao, Zhen Dong, Amir Gholami, Michael W Mahoney, and Kurt Keutzer. Zeroq: A novel zero shot quantization framework. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13169–13178, 2020. **2, 3, 7, 8**
- [6] Hanting Chen, Yunhe Wang, Chang Xu, Zhaohui Yang, Chuanjian Liu, Boxin Shi, Chunjing Xu, Chao Xu, and Qi Tian. Data-free learning of student networks. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 3514–3522, 2019. **3**
- [7] Yoojin Choi, Jihwan Choi, Mostafa El-Khamy, and Jungwon Lee. Data-free network quantization with adversarial knowledge distillation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*, pages 710–711, 2020. **3**
- [8] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pages 248–255. Ieee, 2009. **3**
- [9] Xin Dong, Shangyu Chen, and Sinno Pan. Learning to prune deep neural networks via layer-wise optimal brain surgeon. In *Advances in Neural Information Processing Systems*, 2017. **1**
- [10] Steven K. Esser, Jeffrey L. McKinstry, Deepika Bablani, Rathinakumar Appuswamy, and Dharmendra S. Modha. Learned step size quantization. In *International Conference on Learning Representations*, 2020. **3, 8**
- [11] Alexander Finkelstein, Uri Almog, and Mark Grobman. Fighting quantization bias with bias. *arXiv preprint arXiv:1906.03193*, 2019. **7**
- [12] Ruihao Gong, Xianglong Liu, Shenghu Jiang, Tianxiang Li, Peng Hu, Jiazhen Lin, Fengwei Yu, and Junjie Yan. Differentiable soft quantization: Bridging full-precision and low-bit neural networks. *arXiv preprint arXiv:1908.05033*, 2019. **1**
- [13] Ishaan Gulrajani, Faruk Ahmed, Martin Arjovsky, Vincent Dumoulin, and Aaron C Courville. Improved training of wasserstein gans. In *Advances in neural information processing systems*, pages 5767–5777, 2017. **3**
- [14] Song Han, Huizi Mao, and William J Dally. Deep compression: Compressing deep neural networks with pruning, trained quantization and huffman coding. *arXiv preprint arXiv:1510.00149*, 2015. **1**
- [15] Matan Haroush, Itay Hubara, Elad Hoffer, and Daniel Soudry. The knowledge within: Methods for data-free model compression. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8494–8502, 2020. **2, 3, 7, 8**
- [16] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016. **5**
- [17] Tong He, Zhi Zhang, Hang Zhang, Zhongyue Zhang, Junyuan Xie, and Mu Li. Bag of tricks for image classification with convolutional neural networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 558–567, 2019. **5**
- [18] Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. Distilling the knowledge in a neural network. *arXiv preprint arXiv:1503.02531*, 2015. **1, 3**
- [19] Itay Hubara, Matthieu Courbariaux, Daniel Soudry, Ran El-Yaniv, and Yoshua Bengio. Quantized neural networks: Training neural networks with low precision weights and activations. *The Journal of Machine Learning Research*, 18(1):6869–6898, 2017. **1**
- [20] Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. *arXiv preprint arXiv:1502.03167*, 2015. **2**
- [21] Alex Kendall, Yarin Gal, and Roberto Cipolla. Multi-task learning using uncertainty to weigh losses for scene geometry and semantics. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 7482–7491, 2018. **5**
- [22] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014. **6**
- [23] Raghuraman Krishnamoorthi. Quantizing deep convolutional networks for efficient inference: A whitepaper. *arXiv preprint arXiv:1806.08342*, 2018. **1, 3**
- [24] Yuhang Li, Xin Dong, and Wei Wang. Additive powers-of-two quantization: An efficient non-uniform discretization for neural networks. In *International Conference on Learning Representations*, 2019. **1**
- [25] Hanxiao Liu, Karen Simonyan, and Yiming Yang. Darts: Differentiable architecture search. *arXiv preprint arXiv:1806.09055*, 2018. **2**
- [26] Aravindh Mahendran and Andrea Vedaldi. Understanding deep image representations by inverting them. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 5188–5196, 2015. **3**
- [27] Aravindh Mahendran and Andrea Vedaldi. Visualizing deep convolutional neural networks using natural pre-images. *International Journal of Computer Vision*, 120(3):233–255, 2016. **3**
- [28] Takeru Miyato, Toshiki Kataoka, Masanori Koyama, and Yuichi Yoshida. Spectral normalization for generative adversarial networks. *arXiv preprint arXiv:1802.05957*, 2018. **3**
- [29] Alexander Mordvintsev, Christopher Olah, and Mike Tyka. Inceptionism: Going deeper into neural networks. 2015. **3, 6**

- [30] Markus Nagel, Rana Ali Amjad, Mart van Baalen, Christos Louizos, and Tijmen Blankevoort. Up or down? adaptive rounding for post-training quantization. *arXiv preprint arXiv:2004.10568*, 2020. 3, 7
- [31] Markus Nagel, Mart van Baalen, Tijmen Blankevoort, and Max Welling. Data-free quantization through weight equalization and bias correction. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 1325–1334, 2019. 1, 7
- [32] Antonio Polino, Razvan Pascanu, and Dan Alistarh. Model compression via distillation and quantization. *arXiv preprint arXiv:1802.05668*, 2018. 1
- [33] Ilija Radosavovic, Raj Prateek Kosaraju, Ross Girshick, Kaiming He, and Piotr Dollár. Designing network design spaces. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10428–10436, 2020. 5
- [34] Tim Salimans, Ian Goodfellow, Wojciech Zaremba, Vicki Cheung, Alec Radford, Xi Chen, and Xi Chen. Improved techniques for training gans. In D. Lee, M. Sugiyama, U. Luxburg, I. Guyon, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 29, pages 2234–2242. Curran Associates, Inc., 2016. 6
- [35] Mark Sandler, Andrew Howard, Menglong Zhu, Andrey Zhmoginov, and Liang-Chieh Chen. Mobilenetv2: Inverted residuals and linear bottlenecks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4510–4520, 2018. 5
- [36] Mingxing Tan, Bo Chen, Ruoming Pang, Vijay Vasudevan, Mark Sandler, Andrew Howard, and Quoc V Le. Mnasnet: Platform-aware neural architecture search for mobile. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2820–2828, 2019. 5
- [37] Lucas Theis, Iryna Korshunova, Alykhan Tejani, and Ferenc Huszár. Faster gaze prediction with dense networks and fisher pruning. *arXiv preprint arXiv:1801.05787*, 2018. 1
- [38] Hongxu Yin, Pavlo Molchanov, Jose M Alvarez, Zhizhong Li, Arun Mallya, Derek Hoiem, Niraj K Jha, and Jan Kautz. Dreaming to distill: Data-free knowledge transfer via deep-inversion. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8715–8724, 2020. 2, 3, 6, 7, 8
- [39] Hongyi Zhang, Moustapha Cisse, Yann N Dauphin, and David Lopez-Paz. mixup: Beyond empirical risk minimization. *arXiv preprint arXiv:1710.09412*, 2017. 4
- [40] Han Zhang, Ian Goodfellow, Dimitris Metaxas, and Augustus Odena. Self-attention generative adversarial networks. In *International Conference on Machine Learning*, pages 7354–7363. PMLR, 2019. 6
- [41] Chenzhuo Zhu, Song Han, Huizi Mao, and William J Dally. Trained ternary quantization. *arXiv preprint arXiv:1612.01064*, 2016. 1
- [42] Barret Zoph and Quoc V Le. Neural architecture search with reinforcement learning. *arXiv preprint arXiv:1611.01578*, 2016. 2