

# Diverse Sample Generation: Pushing the Limit of Generative Data-free Quantization

Haotong Qin, Yifu Ding, Xiangguo Zhang, Jiakai Wang, Xianglong Liu\*, *Member, IEEE*, and Jiwen Lu, *Senior Member, IEEE*

**Abstract**—Generative data-free quantization emerges as a practical compression approach that quantizes deep neural networks to low bit-width without accessing the real data. This approach generates data utilizing batch normalization (BN) statistics of the full-precision networks to quantize the networks. However, it always faces the serious challenges of accuracy degradation in practice. We first give a theoretical analysis that the diversity of synthetic samples is crucial for the data-free quantization, while in existing approaches, the synthetic data completely constrained by BN statistics experimentally exhibit severe homogenization at distribution and sample levels. This paper presents a generic **Diverse Sample Generation (DSG)** scheme for the generative data-free quantization, to mitigate detrimental homogenization. We first slack the statistics alignment for features in the BN layer to relax the distribution constraint. Then, we strengthen the loss impact of the specific BN layers for different samples and inhibit the correlation among samples in the generation process, to diversify samples from the statistical and spatial perspectives, respectively. Comprehensive experiments show that for large-scale image classification tasks, our DSG can consistently quantization performance on different neural architectures, especially under ultra-low bit-width. And data diversification caused by our DSG brings a general gain to various quantization-aware training and post-training quantization approaches, demonstrating its generality and effectiveness.

**Index Terms**—data-free quantization, quantized neural networks, model compression, deep learning.

arXiv:2109.00212v3 [cs.CV] 20 Oct 2022

## 1 INTRODUCTION

WITH the advent of deep learning, the deep neural network has achieved a great success in a variety of fields, such as image classification [27], [34], [48], object detection [17], [18], [40], [45], semantic segmentation [14], [59], etc. Nevertheless, it is still a significant challenge to apply advanced neural networks on resource-limited devices for their high memory usage and expensive computation. With more and more hardware supporting low bit-width computations, network quantization emerges as an efficient method to compress and accelerate models [4], [12], [43], [44], [51], [52], [53], [54], [60]. Many quantization methods, called quantization-aware training (QAT), apply the following pipeline: considering the quantization function in the training process on the original dataset, and minimizing the loss caused by the quantization through backward propagation. Since QAT methods require the finetuning steps, it is considered to be time-consuming and computationally intensive [20], [26], [42]. Thus, quantization without training or finetuning process is also demanded in the industry, which is called post-training quantization (PTQ) in recent studies [2], [10], [32], [37], [58].

Most quantization approaches are designed for the data-driven scenario, and thus require real data in their quantization process. However, real data is not always accessible for privacy or security concerns (e.g., medical and user data). Fortunately, data-free quantization is proposed to

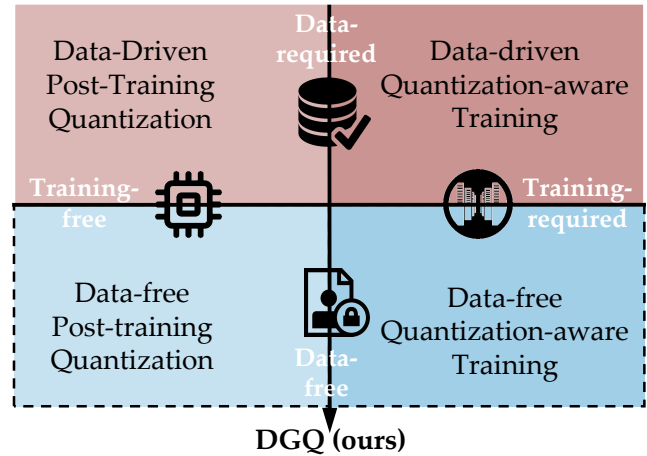


Fig. 1: The existing quantization approaches can be quartered by data-required or not and train/finetune-required or not. Our generic Diverse Sample Generation (DSG) method is proposed for data-free quantization approaches, including data-free post-training quantization and quantization-aware training.

quantize deep neural networks without accessing real data. As shown in Fig. 1, the data-free quantization approaches can also be classified as data-free PTQ and data-free QAT according to whether the model training is required. Among existing data-free quantization methods, generative methods calibrate or train networks using the "optimal synthetic data", which with the distribution best matches the Batch Normalization (BN) statistics of the original full-precision neural network. Generative data-free PTQ methods apply BN statistics loss to directly update the generated data and then progress a computation-saving calibration process [5],

- H. Qin, Y. Ding, X. Zhang, and X. Liu (corresponding author, E-mail: xlliu@buaa.edu.cn) are with the State Key Laboratory of Software Development Environment, Beihang University, China.
- J. Wang is with Zhongguancun Laboratory, China.
- J. Lu is with the Department of Automation, Tsinghua University, China.

Our code is released at <https://github.com/htqin/DSG>.

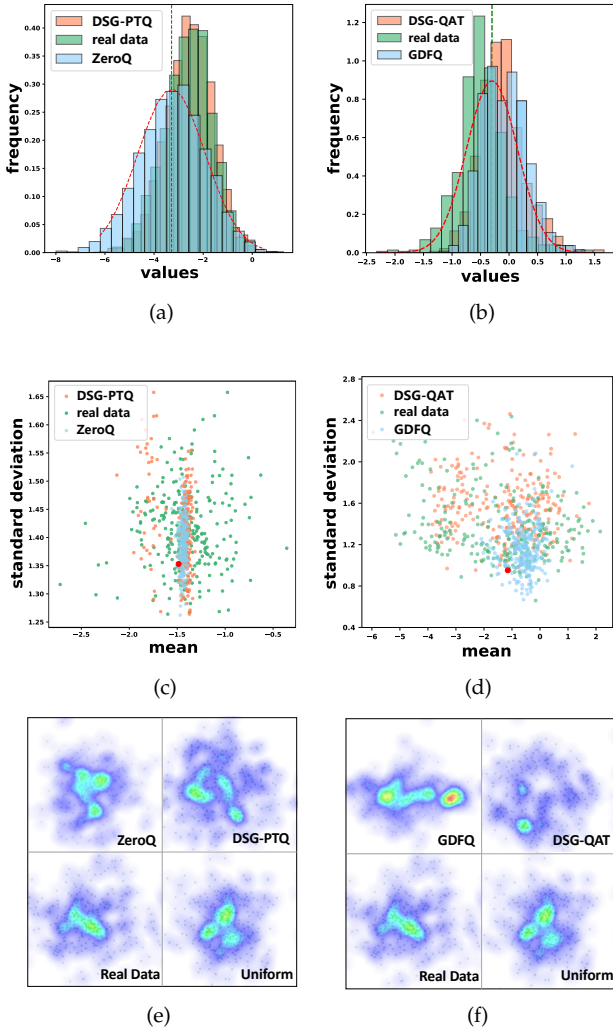


Fig. 2: Comparison between real data (random sampling from ImageNet) and data synthesized by generative data-free quantization methods (DSG and ZeroQ for PTQ in the first column, DSG and GDFQ for QAT in the second column) with 256 samples of each. (a) and (b) show the activation distribution of one channel in ResNet18 for PTQ and QAT, respectively. (c) and (d) are charts of the mean and standard deviation for one channel of synthetic data. (e) and (f) are the visualization of the synthetic data after dimensionality reduction with PCA [15].

[38], [56], while QAT ones generally introduce a separate generator for synthesizing data towards higher accuracy, and train it and the quantized network jointly [8], [25], [35], [47].

However, among existing quantization approaches, the data-driven methods always achieve significantly higher accuracy compared with data-free ones, whether for PTQ or QAT. One intuitive reason is that the difference between real data and synthetic data leads to the gap in the performance between data-driven and data-free quantization approaches. As a result, many approaches attempt to improve on synthetic data so that the performance of generative data-free quantization methods continues to approach data-driven quantization. In this work, we first give a theoretical analysis from a diversity view for generative data-free quantization. Our analysis shows that synthetic samples should be distributed as dispersed in the input domain while

satisfying given constraints. This is often described as the diversification of synthetic samples and would facilitate the optimization of the quantization to cover its potential input domain as possible. Compared with real training and testing datasets that usually approximately satisfy the Independent and Identically Distributed (IID) assumption, synthetic data generated according to certain constraints in generative methods are difficult to meet the diversity requirement.

Corresponding to our theoretical analysis, we also experimentally reveals that severe homogenization exists in the data generation processes of existing typical data-free generative quantization methods (ZeroQ and GDFQ methods for PTQ and QAT, respectively).

**First**, the synthetic data is constrained to match the BN statistics, and thereby the feature distribution might overfit the BN statistics in each layer when the data forward propagate in models. As shown in Fig. 2(a) and Fig. 2(b), the distribution of samples synthesized by existing PTQ and QAT methods almost strictly follows the Gaussian distribution with corresponding BN mean and variance, while the distribution of real data has an obvious offset and enjoys a more diverse distribution.

We consider the first phenomenon as the *distribution-level homogenization*.

**Second**, all samples of synthetic data are updated by one specific objective function in typical generative data-free quantization methods, all synthetic samples are applied to the same constraint and directly sum the loss term of each layer. In Fig. 2(c) and Fig. 2(d), the distribution statistics of real data are dispersed while the data generated by existing approaches are centralized.

**Third**, for data-free QAT methods, synthetic samples are synthesized by a generator network. In the process of learning the distributions of synthetic data for quantization, both the generator and the quantized network may converge to a trivial solution where the former learns to produce few modes exclusively, which is referred to by mode collapse and causes the synthetic data to clustered in sample space [13]. As Fig. 2(e) and Fig. 2(f) show, the synthetic data is aggregated while the real data is scattered.

The second and third phenomena are considered to be *sample-level homogenization* from the statistical and spatial perspectives.

In a word, the distribution-level homogenization means that the overall distribution of synthetic data is strictly restricted to the specific distribution with BN statistics, while the sample-level homogenization means the little differences among samples from the statistical and spatial perspectives. Only mitigating homogenization at one certain level cannot ensure the diversity of data at the other level.

To alleviate the accuracy degeneration of the quantized neural network caused by the homogenization of synthetic data, this paper presents a generic data generation scheme, **Diverse Sample Generation (DSG)**, for generative data-free quantization to enhance the diversity of the synthetic data. The proposed DSG scheme mainly relies on three novel techniques: *Slack Distribution Alignment (SDA)*: relax the distribution constraint of synthetic data by slacking the feature statistics alignment in each BN layer. *Layerwise Sample Enhancement (LSE)*: strengthen the impact of statistics loss of the specific BN layer for its corresponding synthetic sample

by applying a layerwise enhancement. *Sample Correlation Inhibition* (SCI): weaken the correlation among the synthetic samples by applying determinantal point processes loss to the intermediate features in the generation process. Among these techniques, SDA diversifies the synthetic data at the distribution level, while LSE and SCI diversify it at the sample level from the statistical and spatial perspective, respectively. Considering the generality, DSG focuses on the improvement of the generation process while could be almost decoupled from specific quantization methods. Therefore, the proposed techniques of our DSG can be flexibly applied to different PTQ and QAT quantization approaches effectively, improving the accuracy of quantized neural networks by diversifying synthetic data in the generation process.

Our scheme presents a novel perspective of data diversity for generative data-free quantization, and extensive experiments show that DSG significantly improves both data-free PTQ and QAT. The DSG performs remarkably well across several mainstream neural architectures, including VGG16bn [48] (the VGG16 with BN for the dense layers), ResNet18/20/50 [23], SqueezeNext [16], InceptionV3 [50], ShuffleNet [57], and MobileNetV2 [46], and surpasses the existing data-free methods in a wide margin on the large-scale image classification task and achieves the state-of-the-art (SOTA) results. The quantized networks quantized by our DSG even outperform networks quantized by the real data under various settings. The performance of the quantized neural networks trained by our DSG scheme is pushed to that of their full-precision counterparts, meanwhile, the quantization process gets rid of the data dependence.

We summarize the main contributions of this paper as:

- We first give a theoretical analysis for the utility of data diversity in generative data-free quantization. Our analysis shows that synthetic samples should be distributed as dispersed while satisfying given constraints, which would facilitate the optimization of quantization to cover its potential input domain as possible.
- We experimentally show and analyze the homogenization of synthetic data in existing generative data-free quantization. Our study presents the homogenization at the distribution and sample levels and reveals that the sample-level homogenization is not limited to the statistical perspective, but also the spatial perspective.
- We propose a novel generic DSG scheme for generative data-free quantization from a comprehensive perspective of data diversity, which effectively improves the data-free PTQ and QAT. DSG presents a novel Sample Correlation Inhibition (SCI), in conjunction with the Slack Distribution Alignment (SDA) and Layerwise Sample Enhancement (LSE) techniques to diversify synthetic data at distribution and sample level.
- We conduct a detailed ablation study on the proposed DSG, which presents the effectiveness of the proposed techniques (SDA, LSE, and SCI) in generative data-free PTQ and QAT. And the comprehensive evaluation of the DSG scheme shows that our DSG surpasses the existing SOTA methods by a wide margin on various neural architectures and bit-widths, which demonstrates that the diversity is an important property of high-quality synthetic data.

Note that our paper extends the preliminary conference paper [56]. This manuscript first gives a theoretical analysis that diversifying the synthetic samples is a crucial element for improving the data-free quantization process, and experimentally reveals that the sample-level homogenization is not limited to the statistical perspective, but also the spatial perspective. This work proposed a generic DSG scheme for generative data-free quantization, including both PTQ and QAT, while the previous scheme proposed in the conference paper can be regarded as a special case for PTQ. The DSG in this work further presents a novel SCI technique to tackle the sample-level homogenization from the spatial perspective, which can both significantly improve the performance of various generative data-free quantization approaches. We evaluate our DSG scheme on various neural architectures and compared it with more generative data-free quantization methods. Besides, we further add more analysis and discussion of synthetic data and evaluate the data with various methods. The results present that diversity is an important property of high-quality synthetic data in the generative data-free quantization.

## 2 RELATED WORK

**Data-driven Quantization.** As a widely used compression technique, quantization earns lots of attention in recent years. To mitigate the severe accuracy drop resulted from quantization, many studies, such as [20], [26], [54], [60], utilize quantization-aware training methods to regain the accuracy. Generally, they can always provide performance improvements. However, the training process is quite cumbersome for huge computational and time costs. Another critical deficiency is that the datasets are not always ready to use, especially for privacy and security concerns. Without any finetuning or training process, research of post-training quantization has been conducted to improve the trade-off between accuracy and computational cost. [2] introduces two closed-form analytical solutions targeting clipping approximation and per-channel bit allocation combined with bias-correction to push the limit of weights and activations quantization to 4-bit. [10] analyzes the mean squared error of quantization and minimizing it via the OMSE method. Unlike previous work which aims to find an appropriate clipping range to deal with outliers in the bell-shaped distribution, [58] proposes outlier channel splitting while remains a functionally identical network. Not restricted to the round-to-nearest method, [37] exploits a more flexible weight-rounding mechanism called AdaRound, which optimizes task loss due to quantization w.r.t. the preactivation. However, these methods still require access to the limited data for a more resilient and better-quantized model.

**Data-free Quantization.** Several recent studies have explored quantization methods without training and even calibration datasets. [38] uses weight equalization to acquire similar weight ranges among channels and bias correction for activations. However, quantizing models to run in 6-bit or lower bit-width is a non-trivial task. Works like [21], [5] and [25] leverage BN statistics to reconstruct "realistic" data. Both [21] and [47] utilizes knowledge distillation scheme to finetune the model, while it might cost a little more time. [21] proposes an inception scheme to generate data,

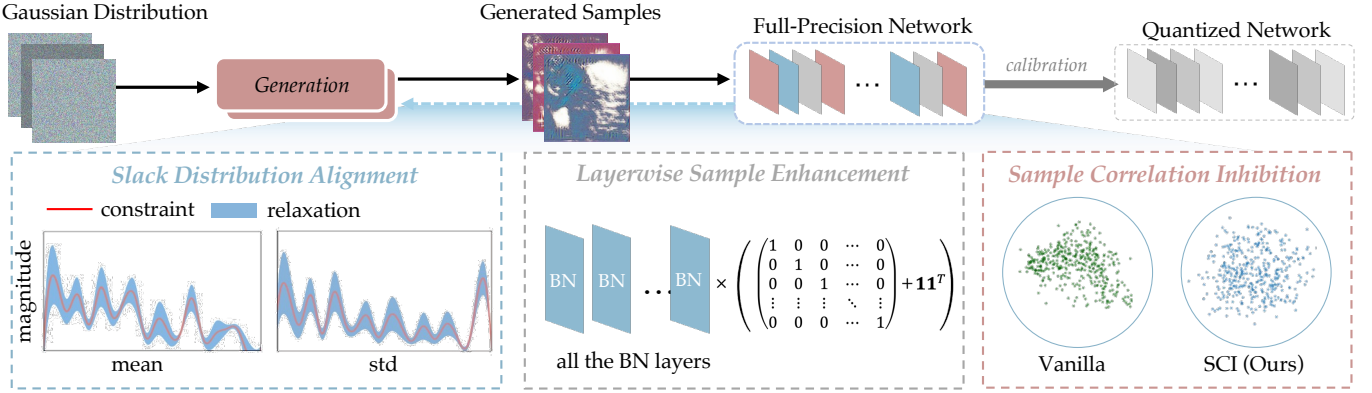


Fig. 3: The overview of the Diverse Sample Generation (DSG) scheme for data-free quantization. The proposed three techniques are applied to the generation process: Slack Distribution Alignment (SDA) relaxes BN statistics constraint in each layer; Layerwise Sample Enhancement (LSE) provides a specific loss term for each sample; Sample Correlation Inhibition (SCI) diversifies the synthetic data in sample space. The first one diversifies synthetic data at the distribution level, and the latter two at the sample level

while [25] utilizes a generative modeling to deal with the distribution of BN statistics. Besides, [5] can work on multi-bit quantization utilizing the Pareto frontier. [8] and [35] focuses on adversarial approaches. [8] proposes a data-free adversarial knowledge distillation method, which minimizes the maximum distance between the outputs of the full-precision teacher model and the quantized student model for any adversarial samples from a generator. [35] instead not only considers the output of model layer but also maximizes the discrepancy of feature map in intermediate inter-channel between teacher and student model by the proposed Channel Relation Map. Although the above methods use distinct tricks to synthesize data and improve model accuracy, they all have similar limitations and we take [5] as an example to illustrate in Section 3.3 specifically.

### 3 METHODS

In this section, we first revisit the previous works of generative data-free quantization. Then, we reveal the data homogenization phenomenon in existing methods and improve the generation process to diversify the synthetic samples.

#### 3.1 Preliminaries

Existing data-driven quantization methods require original training and testing datasets to calibrate or finetune the quantized network for higher performance. But the real data is not accessible in many practical applications for privacy or security concerns, and thereby general quantization schemes cannot be applied directly to these scenarios. The data-free quantization is proposed to quantize the network without access to real data.

Since the real data is inaccessible, the knowledge of the real data in the pre-trained network should be fully exploited in data-free quantization. The statistics (mean and standard deviation) of BN layers in full-precision models fit the original real dataset in the training process. Therefore, most data-free quantization schemes construct the synthetic data by using BN statistics loss to utilize the information in the BN layers. The following optimization objective enables

the distribution of synthetic data  $\mathbf{x}^s$  to best match the BN statistics, including the mean and standard deviation:

$$\min_{\mathbf{x}^s} \mathcal{L}_{\text{BN}} = \sum_{i=1}^N \|\tilde{\boldsymbol{\mu}}_i^s - \boldsymbol{\mu}_i\|_2^2 + \|\tilde{\boldsymbol{\sigma}}_i^s - \boldsymbol{\sigma}_i\|_2^2, \quad (1)$$

where  $\tilde{\boldsymbol{\mu}}_i^s$  and  $\tilde{\boldsymbol{\sigma}}_i^s$  are the mean and standard deviation of feature distribution of synthetic data  $\mathbf{x}^s$  at  $i$ -th BN layer, while  $\boldsymbol{\mu}_i$  and  $\boldsymbol{\sigma}_i$  are that of pre-trained full-precision model  $M$ . By minimizing the BN statistics loss  $\mathcal{L}_{\text{BN}}$ , the synthetic data or the generator learns a distribution of input data to best match BN statistics of each layer.

Among most of the mainstream generative data-free quantization methods, the PTQ methods directly update the data from random vectors during the generation process and then use the synthetic data to calibrate the quantized neural network. The typical optimization objective in generation process of these methods can be expressed as

$$\min_{\mathbf{x}^s} \mathcal{L}_{\text{PTQ}}(\mathbf{x}^s, M), \quad (2)$$

where  $M$  denotes the full-precision model, and the loss function  $\mathcal{L}_{\text{PTQ}}$  is usually the BN statistic loss  $\mathcal{L}_{\text{BN}}$ . The synthetic data  $\mathbf{x}^s$  is then used to calibrate the quantized model, a process that is usually completely separate from the previous generation process.

While the generative data-free QAT methods train a separate generator model to generate the synthetic data and use them to train the quantized neural network. The typical generation process of these methods is expressed as follows:

$$\min_G \mathcal{L}_{\text{QAT}}(G(\mathbf{x}^{s*}), M), \quad \mathbf{x}^s = G(\mathbf{x}^{s*}), \quad (3)$$

where  $\mathbf{x}^{s*}$  is an initial tensor drawn from standard Gaussian distribution  $\mathcal{N}(0, 1)$ , the loss function  $\mathcal{L}_{\text{QAT}}$  usually consists of the BN statistical loss and the cross-entropy function of the generator  $G$ , *etc.*, and  $\mathbf{x}^s$  and  $y$  are synthetic data and the corresponding label, respectively. And the quantized neural network  $Q$  is trained simultaneously using synthetic data:

$$\min_Q \mathbb{E}_{\mathbf{x}^s, y} [\mathcal{C}(Q, (\mathbf{x}^s; y))], \quad (4)$$

where  $\mathcal{C}(\cdot, \cdot)$  is a loss function for quantized neural network  $Q$ , such as cross-entropy loss and mean squared error, and  $y$  is the corresponding labels of synthetic data  $\mathbf{x}^s$ .

### 3.2 Diversity of Synthetic Data for Quantization

In generative data-free quantization and other related generative works, diverse synthetic data is always more welcome intuitively. But there is little relevant theoretical analysis to support this intuition, especially in the context of the data synthesis with strict constraints. It is important to understand the theoretical motivation behind this intuition. Therefore, we give a theoretical analysis to show that in various generative data-free quantization approaches, the diversity of synthetic data facilitates the quantization process.

**Definition 1.** For a well-trained full-precision model  $M$ , the optimization objective of the whole generative data-free quantization process can be abstracted as

$$\min_{\mathbf{x}^s, Q} \mathcal{L}_1(\mathbf{x}^s, M) + \mathcal{L}_2(Q, \mathbf{x}^s), \quad \mathbf{x}^s \subseteq \mathcal{X},$$

where  $\mathbf{x}^s$  is the synthetic data,  $M$  and  $Q$  are the full-precision and quantized networks, respectively, and  $\mathcal{X}$  denotes the input domain of these networks.  $\mathcal{L}_1$  and  $\mathcal{L}_2$  are constraints for the data generation and network quantization.

According to the discussion in Section 3.1, we can summarize that, in generative data-free quantization methods, the key elements are synthetic data, full-precision model, and corresponding quantized model. The full-precision model is well-trained and fixed, while the quantized model and synthetic data are optimizable. Therefore, we abstract and define the general generative data-free quantization process as Definition 1. In this summarized definition, a crucial implicit premise is that the synthetic data  $\mathbf{x}^s$  should reflect the input domain  $\mathcal{X}$  in the process. The optimization process of the quantized model  $Q$  will perform well only if this premise is well satisfied. However, the input domain  $\mathcal{X}$  is not explicit in the process, and the real data, which is usually regarded as the ideal set reflecting the input domain, even cannot be obtained.

**Lemma 1.** For any input domains  $\mathcal{X}$  that includes multiples classes (at least 2) of samples, it can be modeled as several independent high-density  $\{\mathcal{R}_{H1}, \dots, \mathcal{R}_{Hh}\}$  and low-density  $\{\mathcal{R}_{L1}, \dots, \mathcal{R}_{Ll}\}$  sub-regions divided by possible decision surfaces, where  $h \geq 1$  and  $l \geq 0$ .

Therefore, for generative data-free quantization, it is necessary to study how to allow synthetic samples to reflect the properties of an unknown input domain well. Fortunately, there are two classical assumptions about the input domain concluding multiple classes of data in deep learning, and they help reveal the properties of the ideal set of synthetic samples. (1) *Low-density assumption.* The assumption indicates that the boundary of class regions should be more likely to avoid the high-density centers and pass through the low-density areas. (2) *Smoothness assumption.* The assumption states that for two points in the input domain that are close by in the input space, the corresponding labels should be the same. As the Lemma 1, these assumptions allow us to model the potential input domain as several independent sub-regions, and we provide the proof of this modeling in Appendix A.1.

When we model a potential input domain  $\mathcal{X}$  in generative data-free quantization in this way, consider a discrete set of synthetic samples  $\mathbf{x}^s = \{x_0^s, x_1^s, \dots, x_N^s\}$ , we can intuitively get that a necessary condition for the sample set to well

reflect the input domain is to well reflect all sub-regions of the input domain, because each sub-regions are divided as independent and are not representative for others. And for the sample generation, the specific shape and characteristics of the potential input domain are uncertain, so the discussion of the input domain should face all possible input domains rather than a specific one.

**Theorem 1.** Given a set of all possible input domains  $\mathbf{X} = \{\mathcal{X}_0, \mathcal{X}_1, \dots\}$ , whose  $i$ -th element can be denoted as  $\mathcal{X}_i$  with scale  $V^i$  and consists of several sub-regions  $\{\mathcal{R}_1^i, \dots, \mathcal{R}_{K^i}^i\}$  with scales  $\{V_1^i, \dots, V_{K^i}^i\}$ , and the number  $K^i \geq 2$  is unknown yet limited. Consider a sample set  $\mathbf{x}^s = \{x_0^s, \dots, x_N^s\} \subset \mathcal{X}^*$ , where  $\mathcal{X}^* = \mathbb{E}(\mathbf{X})$  denotes the potential input domain and the differences inside each sub-region of  $\mathcal{X}$  is neglected. When the set  $\mathbf{x}^s$  satisfies that for  $\forall x_i^s \in \mathbf{x}^s$ ,  $p(x_i^s \in \mathcal{R}_j^*) = \frac{V_j^*}{V^*}$ , the information reflecting from all possible input domains  $\mathbf{X}$  by the sample set  $\mathbf{x}^s$  will be the maximized in mathematical expectation, where  $V^* = \sum_{k=0}^{K^*} V_k^*$ .

Theorem 1 shows that the optimized strategy to make synthetic samples adequately reflect the unknown input domain in expectation is to correlate the distribution of samples with the spatial scales of the sub-regions, *i.e.*, dispersing synthetic samples uniformly in the potential input domain instead of gathering collapse. We present the proof of the theorem in Appendix A.2. In practice, this means that synthetic samples should be distributed as dispersed in the domain as possible while satisfying given constraints, which is often described as the diversification of synthetic samples.

### 3.3 Homogenization of Synthetic Data

The generative data-free quantization well sidesteps the issue of lacking accessibility of real data, but we have to concede the huge performance drop compared with networks quantized with real data. As the theoretical analysis in Section 3.2, we show that the diversity of synthetic data is beneficial for generative data-free quantization, thus we attempt to find the performance bottleneck of existing methods from this aspect. By experimental analysis for the typical generative data-free PTQ and QAT methods with the almost independent generation process (ZeroQ and GDFQ, respectively), we discover the homogenization phenomena from the existing image synthesizing methods in data-free quantization, which degrade the fidelity and quality of synthetic data. We conclude the homogenization of synthetic data into the distribution level and the sample level to elaborate on them.

#### 3.3.1 Distribution-level Homogenization

Since the performance of existing data-driven quantization significantly surpasses that of data-free quantization in usual, the original training data (real data) is generally considered as the upper limit of the performance of synthetic data. So as Eq. (1) shows, methods constrained by  $\mathcal{L}_{BN}$  are dedicated to generating samples by confining the feature statistics of synthetic data around BN statistics of the pre-trained model. However, the data from the real world is always more diverse, while that generated by fitting BN statistics suffers constrained features as shown in Fig. 2(a) and Fig. 2(b). The activation distribution of data generated in ZeroQ and GDFQ well fit the normal distribution of BN statistics while that of

real data deviates from BN statistics. When we calculate the Wasserstein distance between the Gaussian distribution with BN statistics and the feature distribution of synthetic data, the index values of the existing methods are significantly smaller than the real data (real data 0.120 vs ZeroQ 0.040 in PTQ, and real data 0.134 vs GDFQ 0.116 in QAT). We consider it as distribution-level homogenization causing by the strict constraint applied to the synthesizing process makes the synthetic data overfit BN statistics in each layer.

### 3.3.2 Sample-level Homogenization

Besides homogenization at the distribution level, synthetic data in existing generative data-free quantization also suffer the homogenization at the sample level.

**From the statistical perspective**, in existing generative data-free methods, such as ZeroQ and GDFQ, all samples are constrained by the identical form of BN statistics loss. So there might be little difference among statistics of synthetic samples in one batch, while samples from the real world are diverse. As shown in Fig. 2(c) and Fig. 2(d), the points of the mean and standard deviation of ZeroQ and GDFQ data are significantly overlapping and all points gather together closely near BN statistics. However, statistics of real data and our DSG data have larger variance and are more dispersed. Taking PTQ as an example, the variance of the statistics of the real sample is 0.029, that much larger than the variance of data generated by ZeroQ, which is 0.009.

**From the spatial perspective**, the synthetic samples generated in different ways have different behaviors. For existing data-free methods, all samples are synthesized by the same generation strategy. The samples generated by GDFQ are concentrated in spatial perspective, as shown in Fig. 2(f), while samples of our DSG and real data are more scattered. The heatmap shows the degree of aggregation in the synthetic data sample space. The densest cluster of GDFQ samples appears red in the heatmap and the density index is up to 2052 near the cluster center, while the highest density index of the real sample is only 1681. For PTQ, the index is 1902, and its homogenization phenomenon is not significant as QAT but is more severe than the real data. We regard these as sample-level homogenization from the statistical and spatial perspectives, which leads to a significant performance drop for the networks quantized by synthetic data.

In a nutshell, for the generation process in data-free quantization, the homogenization exists in two levels, *i.e.*, distribution level and sample level. Thus the models quantized by these homogenized data can hardly achieve high accuracy. In this work, we devote diversify the synthetic data to improve the generative data-free quantization. We propose a generic data-free quantization scheme to tackle the homogenization problems by diversifying synthetic data from different perspectives, which we call Diverse Sample Generation (DSG) scheme. The synthetic data generated by our scheme enables the quantized neural network to achieve high accuracy.

## 3.4 Diverse Sample Generation

We present a generic Diverse Sample Generation (DSG) scheme for both generative data-free PTQ and QAT approaches (Fig. 3), including three novel techniques to alleviate homogenization: Slack Distribution Alignment in

Section 3.4.1 is to relax the statistical constraint and thus mitigate distribution-level homogenization, while Layerwise Sample Enhancement in Section 3.4.2 and Sample Correlation Inhibition in Section 3.4.3 diversify the synthetic data at the sample level from distribution and spatial perspectives, respectively.

### 3.4.1 Slack Distribution Alignment

To deal with the distribution-level homogenization, we propose Slack Distribution Alignment (SDA). Specifically, we apply the relaxation constants during the generation process for matching original BN statistics, which can be intuitively regarded as margins for means and standard deviation. With the relaxation constants, the distribution statistics of synthetic samples do not need to fit the BN statistics strictly, and thus the data can be more diverse in distribution. The loss term of SDA is  $\mathbf{l}_{\text{SDA}} = [l_{\text{SDA}_1}, l_{\text{SDA}_2}, \dots, l_{\text{SDA}_N}]^T$ , where  $N$  is the number of BN layers in the quantized neural network, and the  $l_{\text{SDA}_i}$  for  $i$ -th BN layer is expressed as follow:

$$l_{\text{SDA}_i} = \|\max(|\tilde{\boldsymbol{\mu}}_i^s - \boldsymbol{\mu}_i| - \delta_i, 0)\|_2^2 + \|\max(|\tilde{\boldsymbol{\sigma}}_i^s - \boldsymbol{\sigma}_i| - \gamma_i, 0)\|_2^2, \quad (5)$$

where  $\delta_i$  and  $\gamma_i$  denote the relaxation constants for the mean and standard deviation statistics of features at the  $i$ -th BN layers, respectively. After introducing relaxation to the constraint between statistics of synthetic data and BN layers, the mean and standard deviation of generated samples can fluctuate within a certain range. Benefiting from the relaxed constraints, the synthetic samples are more diverse and their feature distribution behaves closer to that of real data.

As the ideal data to quantize networks, our experimental evidence in Fig. 2 shows that the feature statistics of real data still exist gaps compared with BN statistics. Thus, we attempt to approximate the gap between the statistics of real data and BN statistics of the original model as a reference to determine the value of  $\delta_i$  and  $\gamma_i$ . For real-world data, when the number of samples achieves a huge amount, the components and features of inputs approximate the Gaussian assumption. Based on the central limit theorem, we suppose the whole potential input set is generally consistent with a Gaussian distribution.

Thus, we set the initial value of  $\delta_i$  and  $\gamma_i$  by a batch of samples that were randomly initialized from the Gaussian distribution: First, we input 1024 samples initialized by Gaussian distribution with  $\mu = 0$  and  $\sigma = 1$  into the full-precision network and then save the feature statistics at each BN layer, including the mean and standard deviation of the feature distribution. Then, we measure the margin between saved statistics and the corresponding BN statistics, then take percentiles of the margin to initialize  $\delta_i$  and  $\gamma_i$ :

$$\delta_i = |\tilde{\boldsymbol{\mu}}_i^0 - \boldsymbol{\mu}_i|_\epsilon, \quad \gamma_i = |\tilde{\boldsymbol{\sigma}}_i^0 - \boldsymbol{\sigma}_i|_\epsilon, \quad (6)$$

where  $\epsilon \in (0, 1]$  is a hyper-parameter determining the degree of relaxation, and  $\tilde{\boldsymbol{\mu}}_i^0/\tilde{\boldsymbol{\sigma}}_i^0$  are mean/standard deviation of the feature of Gaussian initialized data  $\mathbf{x}^0$  at the  $i$ -th BN layer.  $|\tilde{\boldsymbol{\mu}}_i^0 - \boldsymbol{\mu}_i|_\epsilon$  and  $|\tilde{\boldsymbol{\sigma}}_i^0 - \boldsymbol{\sigma}_i|_\epsilon$  denote the  $\epsilon$  percentile of  $|\tilde{\boldsymbol{\mu}}_i^0 - \boldsymbol{\mu}_i|$  and  $|\tilde{\boldsymbol{\sigma}}_i^0 - \boldsymbol{\sigma}_i|$ , respectively. Intuitively, the larger the value of  $\epsilon$ , the more relaxing the constraints are in Eq. (5). We empirically set the default value of  $\epsilon$  as 0.9 to deal with outliers.

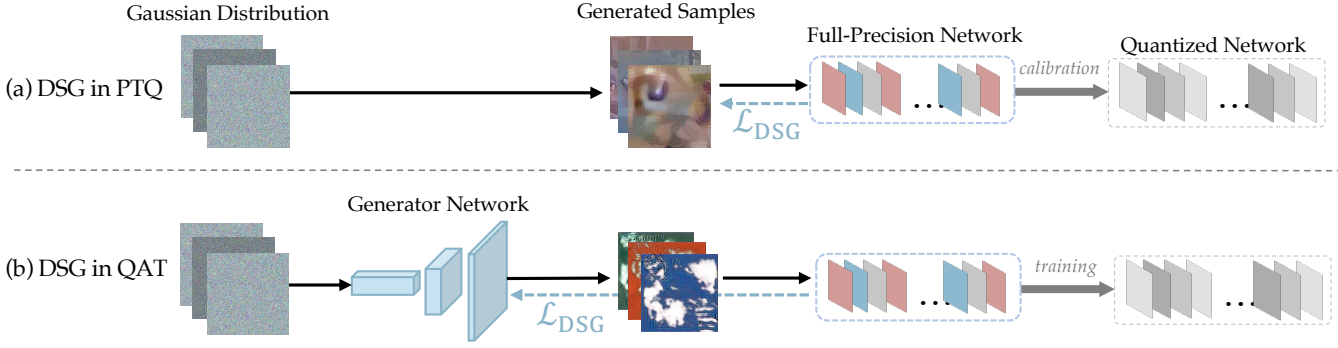


Fig. 4: The procedures of our DSG for specific quantization approaches. The case (a) is the DSG in PTQ and the case (b) is in QAT.

### 3.4.2 Layerwise Sample Enhancement

For the generation processes in existing generative data-free quantization approaches, all synthetic samples are always generated by the same initialization and optimization strategy. Specifically, the samples share the same objective function that directly sums the loss of each layer, *i.e.*, all the BN layers receive the same degree of attention from each sample, and thus results in homogenization among the synthetic samples. We propose Layerwise Sample Enhancement (LSE) to tackle sample-level homogenization. LSE enables each sample to focus on the statistics of different layers by enhancing the specific loss term of the sample in the optimization process so that the synthetic samples of our DSG in one batch could be diverse.

Given a well-trained full-precision network containing the  $N$  number of BN layers, earlier generative data-free quantization methods treat each layer the same way without particular predilection or bias. But in LSE, we engineer  $N$  different loss terms and apply each of them to the specific data sample. Here, we suppose to generate  $N$  samples in one batch, equaling the number of BN layers of the quantized network. The loss term of LSE for the  $i$ -th sample is as follows:

$$l_{\text{LSE}_i} = \frac{1}{N} (\mathbf{1} + \mathbf{h}_i)^T \mathbf{l}, \quad (7)$$

where  $\mathbf{1}$  denotes a column vector of all ones,  $\mathbf{h}_i$  is an  $N$ -dim one-hot column vector where the 1 is in the  $i$ -th position,  $\mathbf{l} = [l_1, l_2, \dots, l_N]^T$  is the vector of the original loss terms for all BN layers. The set of loss functions for all samples can be expressed as  $\mathbf{X}_{\text{LSE}} \mathbf{l} = [l_{\text{LSE}_1}, l_{\text{LSE}_2}, \dots, l_{\text{LSE}_N}]^T$ , which is an  $N$ -dim column vector and its  $i$ -th element represents the loss function for  $i$ -th sample in the batch. Thus, for the whole set of samples  $\mathbf{x}^s$ , the effect exerted by LSE can be formulated as

$$\mathbf{X}_{\text{LSE}} = \frac{1}{N} (\mathbf{1}\mathbf{1}^T + [\mathbf{h}_1, \mathbf{h}_2, \dots, \mathbf{h}_N]), \quad (8)$$

where  $\mathbf{X}_{\text{LSE}}$  can be seen as the enhancement matrix for loss terms. In this way, we focus on every BN layer of the original model respectively and generate corresponding samples for each layer. By such layerwise enhancement, we optimize the whole batch of synthetic samples simultaneously.

### 3.4.3 Sample Correlation Inhibition

The above two proposed methods are based on the relaxation of generation constraints, however, the state of samples at the spatial level still cannot be perceived and adjusted

directly during the generation process. Therefore, in addition to improving the loss function constraining the statistics of synthetic data, we construct the unsupervised Sample Correlation Inhibition (SCI) loss for the samples to sample-level diversify the synthetic data from the spatial perspective. SCI applies the determinant point process loss to spatially disperse the intermediate features in the generator, thereby inhibiting the correlation among these synthetic samples.

The determinant point process  $\mathcal{P}$  can express negative interactions among samples by the similarity kernel  $\mathbf{K}_{\mathbf{S}}$  of the elements in a set  $\mathbf{S}$ , where a large value of  $\mathbf{K}_{ij}; i, j \in \mathbf{S}$  reduces the likelihood of both elements to appear together in a diverse subset [3], [29]. The related works are widely used to generate or select diverse samples in some existing generation tasks for diversifying samples closer to the real data [6], [13], [30], [33]. However, real data is inaccessible in generative data-free quantization, and thereby it cannot be used to construct the constraint of the diversity of synthetic data.

Therefore, we utilize random samples in our SCI instead of the real data, and the random data is applied as a lower bound of diversity to alleviate the homogenization of synthetic samples. Inspired by Theorem 1, each dimension of these vectors is initialized by the uniform distribution  $U[0, 1)$  independently and randomly. The vectors initialized in this way are for uniforming approximately on spatial, which are fixed during the quantization process and expected to compose a base set that covers the potential input domain uniformly and is de-homogenized. In this way, the homogenization of synthetic data is limited to a level lower than its initial state.

So far we are able to build the complete flow of our SCI. The intermediate features correspond to the data  $\mathbf{x}^s$  in the generation process is expressed as  $\mathbf{f} = \{f_1, f_2, \dots, f_N\}$ , where  $f_i$  is the feature of the corresponding synthetic sample. And for the set of features  $\mathbf{f}$ , we construct a set of noise vectors  $\mathbf{r} = \{r_1, r_2, \dots, r_N\}$  of the same size. Then we construct the similarity kernels  $\mathbf{K}_{\mathbf{f}}$  and  $\mathbf{K}_{\mathbf{r}}$  of the features  $\mathbf{f}$  and noise vector set  $\mathbf{r}$  as the following decomposition proposed in [13]:  $\mathbf{K}_{\mathbf{f}} = \phi(\mathbf{f})^T \phi(\mathbf{f})$  and  $\mathbf{K}_{\mathbf{r}} = \phi(\mathbf{r})^T \phi(\mathbf{r})$ , where all  $\phi(f_i) \in \phi(\mathbf{f})$  and  $\phi(r_i) \in \phi(\mathbf{r})$  are the  $\ell_2$  normalized vectors that guarantees the similarity kernels  $\mathbf{K}_{\mathbf{f}}$  and  $\mathbf{K}_{\mathbf{r}}$  to be real positive semidefinite matrices. The main idea of our SCI is to learn a similarity kernel  $\mathbf{K}_{\mathbf{f}}$  of feature  $\mathbf{f}$  and make it less than the similarity kernel  $\mathbf{K}_{\mathbf{r}}$  of random data  $\mathbf{r}$ . In this way, we keep the correlation among the features of

---

**Algorithm 1:** The process of DSG scheme in PTQ

---

**Input:** pretrained model  $M$  with  $N$  BN layers, training iterations  $T$ .  
**Output:** synthetic data:  $\mathbf{x}^s$ , quantized network  $Q$ .  
Initialize  $\mathbf{x}^s$  from Gaussian distribution  $\mathcal{N}(0, 1)$ ;  
Initialize  $\mathbf{x}^0$  from Gaussian distribution  $\mathcal{N}(0, 1)$ ;  
Get  $\mu_i$  and  $\sigma_i$  from BN layers of  $M$ ,  $i = 1, 2, \dots, N$ ;  
Forward propagate  $M(\mathbf{x}^0)$  and gather activations;  
Compute  $\delta_i$  and  $\gamma_i$  using Eq. (6);  
**for all**  $t = 1, 2, \dots, T$  **do**  
    Forward propagate  $M(\mathbf{x}^s)$  and gather activations;  
    Get  $\hat{\mu}_i^s$  and  $\hat{\sigma}_i^s$  from activations;  
    Compute all  $l_{\text{SDA}i}$  using Eq. (5);  
    Descend  $\mathcal{L}_{\text{DSG-PTQ}}$  using Eq. (10) and update data  $\mathbf{x}^s$ ;  
Get synthetic data  $\mathbf{x}^s$ ;  
Calibrate quantized network  $Q$  using data  $\mathbf{x}^s$ .

---



---

**Algorithm 2:** The process of DSG scheme in QAT

---

**Input:** pretrained model  $M$  with  $N$  BN layers, training iterations  $T$ .  
**Output:** generator network  $G$ , quantized network  $Q$ .  
Initialize  $\mathbf{x}^{s*}$  from Gaussian distribution  $\mathcal{N}(0, 1)$  and specify label  $y$ ;  
Get  $\mu_i$  and  $\sigma_i$  from BN layers of  $M$ ,  $i = 1, 2, \dots, N$ ;  
Forward propagate  $M(\mathbf{x}^{s*})$  and gather activations;  
Compute  $\delta_i$  and  $\gamma_i$  using Eq. (6);  
**for all**  $t = 1, 2, \dots, T$  **do**  
    Generate data  $(\mathbf{x}^s, y)$  and gather activations;  
    Forward propagate  $M(\mathbf{x}^s)$  and gather activations;  
    Descend  $\mathcal{L}_{\text{DSG-QAT}}$  using Eq. (11) and update  $G$ ;  
    Update the quantized network  $Q$  using data  $(\mathbf{x}^s, y)$ ;  
Train generator network  $G$  and generate data  $(\mathbf{x}^s, y)$ ;  
Train quantized network  $Q$  using data  $(\mathbf{x}^s, y)$ .

---

synthetic samples lower than that among random vectors, thus alleviating the homogenization of synthetic data. Since matching two similarity kernels directly is an unconstrained optimization problem [31], we construct the loss using the major characteristics: eigenvalues and eigenvectors of kernels. Hence, our SCI loss term is defined as follows:

$$l_{\text{SCI}} = \max \left( \sum_i \left( \|\lambda_{\mathbf{f}}^i - \lambda_{\mathbf{r}}^i\|_2 - \hat{\lambda}_{\mathbf{f}}^i \cos(v_{\mathbf{f}}^i, v_{\mathbf{r}}^i) \right), 0 \right), \quad (9)$$

where  $\lambda_{\mathbf{f}}^i$  and  $\lambda_{\mathbf{r}}^i$  are the  $i^{\text{th}}$  eigenvalues of  $\mathbf{K}_{\mathbf{f}}$  and  $\mathbf{K}_{\mathbf{r}}$ , respectively.  $v_{\mathbf{f}}^i$  and  $v_{\mathbf{r}}^i$  are the eigenvectors, and  $\hat{\lambda}_{\mathbf{f}}^i$  is the min-max normalized version of the eigenvalues applied to alleviate the effect of noisy structures.

### 3.5 Specialization for different quantization approaches

The motivation of our DSG is to improve the performance of various generative data-free quantization methods by proposing a novel and generic data generation scheme. Therefore, here we derive and present the most typical specializations of DSG for generative data-free PTQ and QAT approaches (Fig. 4).

**For the data-free PTQ approach**, DSG applies SDA to each layer for relaxed BN statistics loss terms, integrates LSE to introduce diverse degrees of attention among different samples, and uses SCI to diversify the distribution of samples in the potential input domain. Therefore, the overall loss  $\mathcal{L}_{\text{DSG-PTQ}}$  can be defined as follows involving the above techniques:

$$\mathcal{L}_{\text{DSG-PTQ}} = (\mathbf{X}_{\text{LSE}} \mathbf{I}_{\text{SDA}})^T \mathbf{1} + l_{\text{SCI}}. \quad (10)$$

This optimization loss  $\mathcal{L}_{\text{DSG-PTQ}}$  is used to optimize synthetic data directly. After the generation process is complete, the optimized synthetic data is applied to calibrate the quantized neural network. The calibration process is usually completely separate from the generation process, and the typical calibration methods include [5] and [19]. The DSG for PTQ is elaborated in Algorithm 1, and it is summarized in Fig. 4 as the case (a).

**For the data-free QAT approach**, DSG applies an alternating optimization strategy to update the generator  $G$  and the quantized neural network  $Q$ . As shown in the Eq. (3), we

initialize a set of random noise samples and specified labels, feeding it in generator  $G$  to generate the synthetic data and then training the generator  $G$ . The loss function  $\mathcal{L}_{\text{DSG-QAT}}$  for the generator can be expressed as

$$\mathcal{L}_{\text{DSG-QAT}} = (\mathbf{X}_{\text{LSE}} \mathbf{I}_{\text{SDA}})^T \mathbf{1} + \mathbf{I}_{\text{SCI}}^T \mathbf{1} + \mathcal{L}_{\text{CE}}^G(G), \quad (11)$$

where  $\mathbf{I}_{\text{SCI}} = \frac{1}{B} [l_{\text{SCI}1}, l_{\text{SCI}2}, \dots, l_{\text{SCI}B}]^T$  in Eq. (11) is the SCI loss terms for intermediate features of each blocks ( $B$  blocks in total) in generator  $G$ .  $\mathcal{L}_{\text{CE}}^G(G)$  is the cross-entropy loss.

The data synthesized by generator  $G$  is applied to train the quantized neural network according to the usual practice [47], aiming to close the prediction probability distribution of the quantized neural network to that of the full-precision model. In the quantization process, DSG train the generator  $G$  and the quantized neural network alternately in every epoch, and the pre-trained full-precision network is fixed. The DSG for QAT is elaborated in Algorithm 2 and is further illustrated in Fig. 4 as the case (b).

### 3.6 Analysis and Discussion

The techniques in our DSG scheme, SDA, LSE, and SCI, alleviate the homogenization of synthetic data at the distribution level and sample level. In Fig. 2, we statistically observe the behavior of the three types of data, *i.e.*, data generated by our DSG scheme, data generated by other SOTA generative data-free quantization methods and real data. Here we analyze and discuss the effect of our DSG scheme in diversifying synthetic data in generative data-free quantization methods.

We first analyze the synthetic data generated by our DSG scheme in PTQ. Compared with the data generated by the existing generative data-free PTQ method (ZeroQ), the behavior of samples generated with our DSQ is more similar to that of real data statistically. As shown in Fig. 2(a), the distribution of our DSG data does not strictly fall in close vicinity to BN mean and standard deviation, instead, they are more various in value and thus more fluctuate in frequency. As introduced in Section 3.3, the Wasserstein distance between BN statistics and our DSG synthetic data is 0.124 (compared with 0.040 of ZeroQ). And in Fig. 2(c), the statistics of the data generated by ZeroQ are centralized, while that of the DSG samples are dispersed. The phenomenon means that our synthetic data significantly alleviates the sample-level homogenization of existing data-free generative PTQ



methods from the statistical perspective. As mentioned in Section 3.3, we calculate the variance of the statistics of our DSG data and compare it with ZeroQ in PTQ, which is 0.029 vs. 0.009. Fig. 2(e) presents the distribution of different types of synthetic data in the sample space. The density index near the cluster center of the synthesized sample generated by DSG is only 1571, which is even lower than that of real data 1681 and far lower than ZeroQ 1902.

For the DSG scheme in QAT, due to the utility of SDA and LSE, the distribution of our DSG data is not constrained strictly and the synthetic samples enjoy more diverse statistics, as Fig. 2(b) and Fig. 2(d) show. The SCI further weakens the correlation among samples by constructing unsupervised loss for intermediate features in the generation process to alleviate the homogenization of the generated samples in sample space. From Fig. 2(f), we observe that the density index of the synthesized sample generated by DSG is far lower than the existing data-free QAT method, which means that the sample-level homogenization is greatly alleviated from the spatial perspective.

Therefore, the data generated by our DSG could be more diverse from various perspectives comprehensively, which may have the potential to be an alternative for different quantization approaches that real data is not accessible.

## 4 EXPERIMENTS

In this section, we conduct comprehensive experiments to validate the performance of our DSQ scheme on image classification tasks. We first conduct an ablation study to show the effects of different components, including SDA, LSE, and SCI in Section 4.1. Then in Section 4.2, we compare DSQ with SOTA data-free PTQ and QAT methods respectively across various network architectures. We evaluate the DSG on CIFAR10 [28] and ImageNet (ILSVRC12) [11] datasets in PTQ while on ImageNet dataset in QAT. Finally, in Section 4.3, we conduct a further study on synthetic data. Specifically, we analyze and discuss the data diversity, and integrate and evaluate our synthetic data with various calibration methods and data-driven quantization methods. The results show that good diversity is an important property of high-quality synthetic data.

**DSG scheme.** In the data-free PTQ, the proposed DSG scheme is used for generating synthetic data while the independent calibration processes are as [5] and [19], and the effectiveness is evaluated by measuring the accuracy of quantized models. Unless otherwise specified, the calibration process for DSG in PTQ is the same as [5]. For the DSG scheme in QAT, we train the generator network to synthesize the data and use it to finetune the quantized network. The generator in the DSG scheme in QAT is constructed following ACGAN [39] as in [47].

**Network architectures.** We evaluate our DSG scheme in a wide range of network architectures with various bit-width to prove the versatility of our method. We employ VGG16bn [48], ResNet20/18/50 [23], SqueezeNext [16], InceptionV3 [50], ShuffleNet [57], and MobileNetV2 [46] with various bit-widths, including W4A4 (means 4-bit weight and 4-bit activation), W6A6, W8A8, etc.

**Implementation details.** The proposed scheme is implemented by PyTorch for the sake of the powerful automatic

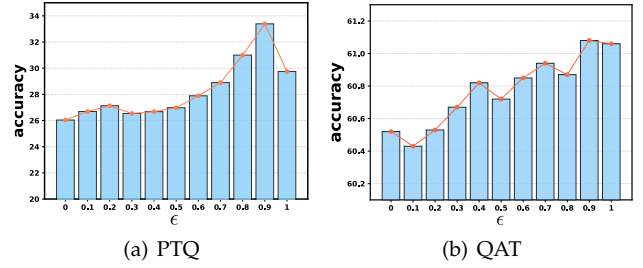


Fig. 5: The accuracy comparison of different  $\epsilon$  values in Eq. (6) on ResNet18 in PTQ and QAT. As  $\epsilon$  varies from 0 to 0.9, the final accuracy is mainly on the rise. But a significant drop is caused by the outliers when  $\epsilon = 1$ .

Method	No FT	W-bit	A-bit	Top-1
Baseline	–	32	32	71.47
Vanilla (ZeroQ)	✓	4	4	26.04
Sample Correlation Inhibition	✓	4	4	39.53
Layerwise Sample Enhancement	✓	4	4	27.12
Slack Distribution Alignment	✓	4	4	33.39
DSG (Ours)	✓	4	4	<b>39.90</b>
Vanilla (GDFQ)	✗	4	4	60.52
Layerwise Sample Enhancement	✗	4	4	60.71
Sample Correlation Inhibition	✗	4	4	60.92
Slack Distribution Alignment	✗	4	4	61.08
DSG (Ours)	✗	4	4	<b>62.18</b>

TABLE 1: Ablation study for DSG scheme on ResNet18. We abbreviate quantization bits used for weights as "W-bit" (for activations as "A-bit"), top-1 test accuracy as "Top-1".

differentiation mechanism. We adopt Gaussian distribution as initialization for the data generation process in our DSG scheme. In our experiments, we quantize all the layers. And the activation is clipped in a layerwise manner. As for hyper-parameter (e.g., the number of iterations to provide synthetic data), we mostly follow the released official implementations or the models and settings clarified in their original paper for a fair comparison [5], [19], [47], [56]. For the training and finetuning process of DSG in QAT, we use the Adam and SGD optimizer in the experiments, where momentum is 0.9 and weight decay is 1e-4. For CIFAR10, we train quantized networks and generators for 400 epochs. The learning rates are initialized to 1e-4 and 1e-3, respectively, and both of them are decayed by 0.1 for every 100 epochs. For ImageNet, we set the initial learning rate of the quantized model as 1e-6, and other training settings are the same as those on CIFAR10.

### 4.1 Ablation Study

We investigate the effect of the proposed LSE, SDA, and SCI techniques for our DSG scheme in data-free PTQ and QAT by ablation experiments. We use the ResNet18 architecture with the ImageNet dataset to evaluate our method under the W4A4 bit-width setting, which can show the effect of each part more obviously.

#### 4.1.1 Effect of SDA

We first analyze the effectiveness of SDA. As discussed earlier, hyper-parameter  $\epsilon$  determines the degree of relaxation in

SDA. Therefore, in Fig. 5, we verify the impact of different values of  $\epsilon$  in the interval  $(0, 1]$  with a moderate step size 0.1, and further add  $\epsilon = 0$  to serve as the vanilla case that the fitting of BN statistics is constrained without relaxation. In PTQ (Fig. 5(a)), when the  $\epsilon$  varies from 0 to 0.9, the final performance increases gradually from 26.04% to 33.39%. In QAT (Fig. 5(b)), the accuracy also presents an overall upward trend (from 60.52% to 61.08%) as the  $\epsilon$  varies from 0 to 0.9. However, when  $\epsilon = 1$ , the accuracy drops more or less in both scenarios, which is 3.61% in PTQ and 0.02% in QAT. These phenomena forcefully prove that by way of relaxing the constraints, the SDA method brings generated data a certain offset when fitting the BN statistic distribution, which contributes to feature diversity and consequently improves the performance. And when  $\epsilon$  is set to 1, the degree of slack reaches the limit as all the outliers in  $\tilde{\mu}_i^0$  and  $\tilde{\sigma}_i^0$  are taken into account, which means the feature distribution of generated data might be out of the reasonable range. Therefore, it might result in better feature diversity in synthetic data but would harm the final accuracy of the quantized network. Therefore, we set the default value of  $\epsilon$  as 0.9 empirically to balance divergence and the impact of outliers.

As the results shown in TABLE 1, the vanilla data-free PTQ method without SDA suffers a severe accuracy degradation by 7.35%, and the SDA also provides 0.56% accuracy promotion in data-free QAT (vanilla 60.52% vs. SDA 61.08%). The results reflect that SDA is essential. Compared to the other two parts of our scheme, SDA may provides a major contribution to the final performance.

#### 4.1.2 Effect of LSE

The TABLE 1 also shows that LSE can improve the performance in both data-free QAT and PTQ. As the ablation results show, in PTQ, LSE method gives a non-negligible increment compared with ZeroQ by 1.08%. As for in QAT, it also helps the DSG scheme to have a slight improvement compared with the vanilla method by 0.19%.

#### 4.1.3 Effect of SCI

Then we evaluate SCI in data-free quantization. The motivation of our SCI is to alleviate the sample-level homogenization from the spatial perspective. TABLE 1 shows the effects of using SCI. In QAT, compared with the vanilla method, integrating with SCI helps to acquire better accuracy, which is 0.40% higher. Also, in PTQ, the SCI method gains 13.49% accuracy compared with using pure random inputs.

Intuitively, these three techniques are motivated by different observations, and also the processes are carefully engineered that they hardly interfere with each other. In short, SDA and LSE improve the loss related to BN statistics, aiming to prevent the generated samples from overfitting to BN statistics and makes the samples focus on the statistics of different layers. While the SCI for data-free quantization constructs an unsupervised loss to constrain the features in the generator to holding distances among them, so that mitigate the sample homogenization. The results show that, in data-free QAT, DSG scheme equipped with these techniques obtains 1.66% improvement in total with ResNet18 under W4A4 setting, which is up to 62.18%. As for data-free PTQ, DSG can also boost the performance to 13.86% with all the techniques.

Arch	Method	No D	No FT	W-bit	A-bit	Top-1	
ResNet20	Baseline	-	-	32	32	94.08	
	Real Data	✗	✓	4	4	87.38	
	ZeroQ	✓	✓	4	4	85.39	
	DSG (Ours)	✓	✓	4	4	<b>87.79</b>	
	Real Data	✗	✓	6	6	93.80	
	ZeroQ	✓	✓	6	6	93.33	
	DSG (Ours)	✓	✓	6	6	<b>93.55</b>	
	Real Data	✗	✓	8	8	93.95	
	ZeroQ	✓	✓	8	8	93.94	
	DSG (Ours)	✓	✓	8	8	<b>93.97</b>	
	VGG16bn	Baseline	-	-	32	32	93.86
		Real Data	✗	✓	4	4	92.50
ZeroQ		✓	✓	4	4	91.79	
DSG (Ours)		✓	✓	4	4	<b>92.89</b>	
Real Data		✗	✓	6	6	93.48	
ZeroQ		✓	✓	6	6	93.45	
DSG (Ours)		✓	✓	6	6	<b>93.68</b>	
Real Data		✗	✓	8	8	93.59	
ZeroQ		✓	✓	8	8	93.53	
DSG (Ours)		✓	✓	8	8	<b>93.61</b>	

TABLE 2: Results of data-free PTQ methods with ResNet20 and VGG16bn on CIFAR10.

## 4.2 Comparison with SOTA Methods

We extensively evaluate our DSG scheme on a wide range of architectures on CIFAR10 and ImageNet datasets for the image classification tasks. We denote the bit-width setting of the quantized network as  $WwAa$  where  $w$  is the bit-width for weight and  $a$  is that for activation, like W8A8, W6A6, W4A4, etc.

### 4.2.1 Comparison with Data-free PTQ Methods

To evaluate the advantage of the proposed scheme in PTQ, we first compare our DSG against other data-free PTQ methods (ZeroQ [5], DFQ [38], ACIQ [1], MSE [7], KL [49], SQuant [19], and OCS [58]) on CIFAR10 and ImageNet datasets. Among these methods, ZeroQ and SQuant are typical generative data-free PTQ methods, which reconstruct synthetic data and calibrate the quantized network in different ways. Thus, we evaluate the generation performance of our DSG in PTQ with calibration methods from ZeroQ and SQuant, denoted as DSG<sup>1</sup> and DSG<sup>2</sup>. Other quantization methods use weight equalization or analytical clip range to improve the network performance. Besides, we additionally compare our method with OCS, which is also a PTQ method but requires real data for calibration. We evaluate these methods on various bit-width settings, and the results on CIFAR10 and ImageNet dataset are shown as TABLE 2 and TABLE 3, respectively.

Specifically, on CIFAR10 [28] dataset, we evaluate our DSG with ResNet20 [24] and VGG16bn [48]. The results are shown in TABLE 2. Under all settings on the CIFAR10 dataset, our DSG far exceeds existing SOTA methods. And

(a) ResNet18 and ResNet50							(b) SqueezeNext, InceptionV3, and ShuffleNet							
Arch	Method	No D	No FT	W-bit	A-bit	Top-1	Arch	Method	No D	No FT	W-bit	A-bit	Top-1	
ResNet18	Baseline	-	-	32	32	71.47	SqueezeNext	Baseline	-	-	32	32	69.38	
	Real Data	✗	✓	4	4	65.22		Real Data	✗	✓	6	6	66.51	
	DFQ	✓	✓	4	4	0.10		ZeroQ	✓	✓	6	6	39.83	
	ACIQ	✓	✓	4	4	7.19		SQuant	✓	✓	6	6	67.34	
	MSE	✓	✓	4	4	15.08		DSG <sup>1</sup> (Ours)	✓	✓	6	6	<b>66.23</b>	
	KL	✓	✓	4	4	16.27		DSG <sup>2</sup> (Ours)	✓	✓	6	6	<b>68.05</b>	
	ZeroQ	✓	✓	4	4	26.04		Real Data	✗	✓	8	8	69.23	
	SQuant	✓	✓	4	4	66.14		ZeroQ	✓	✓	8	8	68.01	
	DSG <sup>1</sup> (Ours)	✓	✓	4	4	<b>39.90</b>		SQuant	✓	✓	8	8	69.22	
	DSG <sup>2</sup> (Ours)	✓	✓	4	4	<b>66.67</b>		DSG <sup>1</sup> (Ours)	✓	✓	8	8	<b>69.27</b>	
	Real Data	✗	✓	6	6	71.18		DSG <sup>2</sup> (Ours)	✓	✓	8	8	<b>69.37</b>	
	ACIQ	✓	✓	6	6	61.15		Baseline	-	-	32	32	78.80	
	KL	✓	✓	6	6	61.34		Real Data	✗	✓	4	4	73.50	
	MSE	✓	✓	6	6	66.96		ZeroQ	✓	✓	4	4	12.00	
	DFQ	✓	✓	6	6	67.30		SQuant	✓	✓	4	4	73.26	
	ZeroQ	✓	✓	6	6	69.74		DSG <sup>1</sup> (Ours)	✓	✓	4	4	<b>57.17</b>	
	DSG <sup>1</sup> (Ours)	✓	✓	6	6	<b>70.46</b>		DSG <sup>2</sup> (Ours)	✓	✓	4	4	<b>74.02</b>	
	DSG <sup>2</sup> (Ours)	✓	✓	6	6	<b>71.18</b>		Real Data	✗	✓	6	6	78.59	
	Real Data	✗	✓	8	8	71.48		ZeroQ	✓	✓	6	6	75.14	
	ACIQ	✓	✓	8	8	68.78		SQuant	✓	✓	6	6	78.30	
DFQ	✓	✓	8	8	69.70	DSG <sup>1</sup> (Ours)	✓	✓	6	6	<b>78.12</b>			
KL	✓	✓	8	8	70.69	DSG <sup>2</sup> (Ours)	✓	✓	6	6	<b>78.59</b>			
MSE	✓	✓	8	8	71.01	Real Data	✗	✓	8	8	78.79			
ZeroQ	✓	✓	8	8	71.43	ZeroQ	✓	✓	8	8	78.70			
SQuant	✓	✓	8	8	71.43	SQuant	✓	✓	8	8	78.79			
DSG <sup>1</sup> (Ours)	✓	✓	8	8	<b>71.49</b>	DSG <sup>1</sup> (Ours)	✓	✓	8	8	<b>78.81</b>			
DSG <sup>2</sup> (Ours)	✓	✓	8	8	<b>71.46</b>	DSG <sup>2</sup> (Ours)	✓	✓	8	8	<b>78.85</b>			
ResNet50	Baseline	-	-	32	32	77.72	InceptionV3	Baseline	-	-	32	32	65.07	
	Real Data	✗	✓	4	4	68.13		Real Data	✗	✓	6	6	56.25	
	ZeroQ	✓	✓	4	4	8.20		ZeroQ	✓	✓	6	6	39.92	
	DFQ	✓	✓	4	4	10.32		SQuant	✓	✓	6	6	60.25	
	SQuant	✓	✓	4	4	70.80		DSG <sup>1</sup> (Ours)	✓	✓	6	6	<b>60.71</b>	
	DSG <sup>1</sup> (Ours)	✓	✓	4	4	<b>56.12</b>		DSG <sup>2</sup> (Ours)	✓	✓	6	6	<b>61.94</b>	
	DSG <sup>2</sup> (Ours)	✓	✓	4	4	<b>68.30</b>		Real Data	✗	✓	8	8	64.52	
	Real Data	✗	✓	6	6	76.84		ZeroQ	✓	✓	8	8	64.46	
	OCS	✗	✓	6	6	74.80		SQuant	✓	✓	8	8	64.68	
	ZeroQ	✓	✓	6	6	75.56		DSG <sup>1</sup> (Ours)	✓	✓	8	8	<b>64.87</b>	
	SQuant	✓	✓	6	6	77.05		DSG <sup>2</sup> (Ours)	✓	✓	8	8	<b>64.97</b>	
	DSG <sup>1</sup> (Ours)	✓	✓	6	6	<b>76.90</b>		Real Data	✗	✓	8	8	64.52	
	DSG <sup>2</sup> (Ours)	✓	✓	6	6	<b>77.22</b>		ZeroQ	✓	✓	8	8	64.46	
	Real Data	✗	✓	8	8	77.70		SQuant	✓	✓	8	8	64.68	
	ZeroQ	✗	✓	8	8	77.67		DSG <sup>1</sup> (Ours)	✓	✓	8	8	<b>64.87</b>	
	DSG <sup>1</sup> (Ours)	✓	✓	8	8	<b>77.72</b>		DSG <sup>2</sup> (Ours)	✓	✓	8	8	<b>64.97</b>	
	DSG <sup>2</sup> (Ours)	✓	✓	8	8	<b>77.83</b>								

TABLE 3: Results of data-free PTQ methods with (a) ResNet18, ResNet50, (b) SqueezeNext, InceptionV3, and ShuffleNet on ImageNet. Here, "Arch" means the network architectures, "No D" means that none of the data is used to assist quantization, "No FT" stands for no finetuning (retraining). "Real Data" represents using real training data and quantization methods in ZeroQ (without any finetuning).

we highlight that the quantized network calibrated with DSG data completely surpasses the network calibrated with real data on various settings. For example, under the W4A4

settings of ResNet20 and VGG6bn, our method exceeds the calibration with real data by 0.41% and 0.39%, respectively.

As listed in TABLE 3, the results over various network architectures, including ResNet18/50 [23], SqueezeNext [16],

(a) ResNet18 and ResNet50							(b) ShuffleNet, MobileNetV2, and InceptionV3						
Arch	Method	No D	No FT	W-bit	A-bit	Top-1	Arch	Method	No D	No FT	W-bit	A-bit	Top-1
ResNet18	Baseline	-	-	32	32	71.74	ShuffleNet	Baseline	-	-	32	32	65.07
	Real Data	✗	✗	4	4	63.87		Real Data	✗	✗	4	4	29.18
	GDFQ	✓	✗	4	4	60.52		GDFQ	✓	✗	4	4	21.78
	DSG (Ours)	✓	✗	4	4	<b>62.18</b>		DSG (Ours)	✓	✗	4	4	<b>29.71</b>
	Real Data	✗	✗	6	6	71.42		Real Data	✗	✗	6	6	62.89
	Integer-Only	✗	✗	6	6	67.30		GDFQ	✓	✗	6	6	60.12
	GDFQ	✓	✗	6	6	70.43		DSG (Ours)	✓	✗	6	6	<b>61.37</b>
	DSG (Ours)	✓	✗	6	6	<b>71.12</b>		Real Data	✗	✗	8	8	62.95
	Real Data	✗	✗	8	8	71.44		GDFQ	✓	✗	8	8	64.03
	DFC	✓	✗	8	8	69.57		DSG (Ours)	✓	✗	8	8	<b>64.76</b>
	RVQuant	✗	✗	8	8	70.01		Baseline	-	-	32	32	71.88
	GDFQ	✓	✗	8	8	71.43		Real Data	✗	✗	4	4	66.39
	DSG (Ours)	✓	✗	8	8	<b>71.54</b>		GDFQ	✓	✗	4	4	51.30
	ResNet50	Baseline	-	-	32	32		77.72	DSG (Ours)	✓	✗	4	4
Real Data		✗	✗	4	4	70.27	Real Data	✗	✗	6	6	72.11	
GDFQ		✓	✗	4	4	55.65	Integer-Only	✗	✗	6	6	70.90	
RVQuant		✗	✗	4	4	64.90	GDFQ	✓	✗	6	6	70.98	
ZAQ		✓	✗	4	4	70.06	GZNQ	✓	✗	6	6	71.12	
DSG (Ours)		✓	✗	4	4	<b>71.96</b>	DSG (Ours)	✓	✗	6	6	<b>71.48</b>	
Real Data		✗	✗	6	6	77.56	Real Data	✗	✗	8	8	72.92	
GDFQ		✓	✗	6	6	76.59	GDFQ	✓	✗	8	8	70.17	
ZS-CGAN		✓	✗	6	6	76.82	ZAQ	✓	✗	8	8	71.43	
DSG (Ours)		✓	✗	6	6	<b>77.25</b>	DSG (Ours)	✓	✗	8	8	<b>72.90</b>	
Real Data		✗	✗	8	8	77.66	Baseline	-	-	32	32	78.80	
GDFQ		✓	✗	8	8	77.51	Real Data	✗	✗	4	4	73.51	
DSG (Ours)		✓	✗	8	8	<b>77.64</b>	GDFQ	✓	✗	4	4	70.39	
ResNet50		Baseline	-	-	32	32	77.72	DSG (Ours)	✓	✗	4	4	<b>72.01</b>
	Real Data	✗	✗	4	4	70.27	Real Data	✗	✗	6	6	78.81	
	GDFQ	✓	✗	4	4	55.65	GDFQ	✓	✗	6	6	77.20	
	RVQuant	✗	✗	4	4	64.90	DSG (Ours)	✓	✗	6	6	<b>78.60</b>	
	ZAQ	✓	✗	4	4	70.06	Real Data	✗	✗	8	8	79.00	
	DSG (Ours)	✓	✗	4	4	<b>71.96</b>	GDFQ	✓	✗	8	8	78.62	
	Real Data	✗	✗	6	6	77.56	DSG (Ours)	✓	✗	8	8	<b>78.94</b>	
	GDFQ	✓	✗	6	6	76.59	Baseline	-	-	32	32	78.80	
	ZS-CGAN	✓	✗	6	6	76.82	Real Data	✗	✗	4	4	73.51	
	DSG (Ours)	✓	✗	6	6	<b>77.25</b>	GDFQ	✓	✗	4	4	70.39	
	Real Data	✗	✗	8	8	77.66	DSG (Ours)	✓	✗	4	4	<b>72.01</b>	
	GDFQ	✓	✗	8	8	77.51	Real Data	✗	✗	6	6	78.81	
	DSG (Ours)	✓	✗	8	8	<b>77.64</b>	GDFQ	✓	✗	6	6	77.20	
	ResNet50	Baseline	-	-	32	32	77.72	DSG (Ours)	✓	✗	6	6	<b>78.60</b>
Real Data		✗	✗	4	4	70.27	Real Data	✗	✗	8	8	79.00	
GDFQ		✓	✗	4	4	55.65	GDFQ	✓	✗	8	8	78.62	
RVQuant		✗	✗	4	4	64.90	DSG (Ours)	✓	✗	8	8	<b>78.94</b>	
ZAQ		✓	✗	4	4	70.06	Baseline	-	-	32	32	78.80	
DSG (Ours)		✓	✗	4	4	<b>71.96</b>	Real Data	✗	✗	4	4	73.51	
Real Data		✗	✗	6	6	77.56	GDFQ	✓	✗	4	4	70.39	
GDFQ		✓	✗	6	6	76.59	DSG (Ours)	✓	✗	4	4	<b>72.01</b>	
ZS-CGAN		✓	✗	6	6	76.82	Real Data	✗	✗	6	6	78.81	
DSG (Ours)		✓	✗	6	6	<b>77.25</b>	GDFQ	✓	✗	6	6	77.20	
Real Data		✗	✗	8	8	77.66	DSG (Ours)	✓	✗	6	6	<b>78.60</b>	
GDFQ		✓	✗	8	8	77.51	Real Data	✗	✗	8	8	79.00	
DSG (Ours)		✓	✗	8	8	<b>77.64</b>	GDFQ	✓	✗	8	8	78.62	
ResNet50		Baseline	-	-	32	32	77.72	DSG (Ours)	✓	✗	8	8	<b>78.94</b>
	Real Data	✗	✗	4	4	70.27	Baseline	-	-	32	32	78.80	
	GDFQ	✓	✗	4	4	55.65	Real Data	✗	✗	4	4	73.51	
	RVQuant	✗	✗	4	4	64.90	GDFQ	✓	✗	4	4	70.39	
	ZAQ	✓	✗	4	4	70.06	DSG (Ours)	✓	✗	4	4	<b>72.01</b>	
	DSG (Ours)	✓	✗	4	4	<b>71.96</b>	Real Data	✗	✗	6	6	78.81	
	Real Data	✗	✗	6	6	77.56	GDFQ	✓	✗	6	6	77.20	
	GDFQ	✓	✗	6	6	76.59	DSG (Ours)	✓	✗	6	6	<b>78.60</b>	
	ZS-CGAN	✓	✗	6	6	76.82	Real Data	✗	✗	8	8	79.00	
	DSG (Ours)	✓	✗	6	6	<b>77.25</b>	GDFQ	✓	✗	8	8	78.62	
	Real Data	✗	✗	8	8	77.66	DSG (Ours)	✓	✗	8	8	<b>78.94</b>	
	GDFQ	✓	✗	8	8	77.51	Baseline	-	-	32	32	78.80	
	DSG (Ours)	✓	✗	8	8	<b>77.64</b>	Real Data	✗	✗	4	4	73.51	

TABLE 4: Results of data-free QAT methods with ResNet18, ResNet50, InceptionV3, ShuffleNet, and MobileNetV2 on ImageNet.

InceptionV3 [50], and ShuffleNet [57], demonstrate that our proposed DSG significantly outperforms previous sample generation methods. As the results shown, the accuracy of quantized networks calibrated with our DSG data under the W8A8 setting is almost not decreased and even surpasses the full-precision (W32A32) baseline networks on the ResNet18 and InceptionV3 architectures. The higher accuracy might be attributed to making full use of the potential of the quantized network. Quantizing the network to W8A8 maintains greater representation ability of the network compared with the lower bit-width quantization (such as W4A4) and brings less quantization error. Thus, under this setting, our diversified data can alleviate the performance degradation of the network quantization while even leading to a better

solution compared with the full-precision network. And the advantage of our DSG gets more evident when the bit-width becomes lower. For example, our DSG outperforms ZeroQ on SqueezeNext on the W6A6 setting by more than 26%, and even outperforms real data by 1.54%. And it is noteworthy that in W4A4 cases with ResNet18 architecture, our DSG surpasses ZeroQ by 13.86%, and even surpasses real data by a notable 1.45% with the compared SQuant calibration, which is up to 66.67%. Moreover, with InceptionV3 under W4A4 settings, our DSG helps to acquire 74.02% accuracy, which is higher than SQuant and real data by eminent 0.76% and 0.52%, respectively.

### 4.2.2 Comparison with Data-free QAT Methods

Furthermore, to demonstrate the applicability of our DSG scheme in data-free QAT, we compare it with existing QAT methods, such as DFC [21], RVQuant [41], Integer-Only [26], ZAQ [35], GZLNQ [25], ZS-CGAN [9], and GDFQ [47]. These methods are engineered to finetune and update the network parameters. DFC is a finetuning method to recover accuracy for ultra-low bit-width cases, which uses Inceptionism [36] to facilitate the generation of data with random labels. The other mentioned methods finetune the quantized network to improve the accuracy of the network. Among them, ZAQ and GDFQ introduce a generator to synthesize data and use the generated data to finetune the network. We evaluate these methods on various bit-width settings in the more challenging large-scale image classification task. The results on ImageNet are shown as TABLE 4.

We test these methods on ResNet18/50 [23], InceptionV3 [50], ShuffleNet [57], and MobileNetV2 [46]. The results in TABLE 4 show that our DSG enjoys the best performance. Concretely, as can be seen, regardless of network architecture, our DSG method surpasses many previous methods, including DFC, RVQuant, Integer-Only, and GDFQ in various bit-width settings. And it is noteworthy that, after finetuning with our DSG data, quantized ResNet18 and InceptionV3 even surpass their full-precision counterparts under the W8A8 setting. Comparing with other data-free QAT training methods under W4A4 setting, our DSG outperforms GDFQ by 1.66% and 16.31% with ResNet18 and ResNet50 respectively, also higher than RVQuant with real data by 7.06% with ResNet50. Besides the mainstream neural architectures, we also evaluate our DSG over existing well-designed lightweight architectures. The proposed DSG shows an overwhelming advantage over GDFQ with a similar training pipeline, 9.16% higher with MobileNetV2 under W4A4, 7.93% higher with ShuffleNet under W4a4, and 1.62% higher with InceptionV3 under W4A4.

In short, our DSG scheme outperforms other competitors over a wide range of experiments in data-free PTQ and QAT, including various bit-width, different network architectures, and two datasets. All the results forcefully demonstrate that the diversity of generated data is significant to calibrate the model for higher accuracy, especially in ultra-low bit-width conditions. If the synthetic data are homogeneous or even identical at the distribution and sample level, it would be ineffectual when used to quantize the network. Our DSG effectively diversifies the synthetic data and thus improves the performance of the quantized neural network.

### 4.3 Further Study on Synthetic Data

To further demonstrate the data diversifying caused by our DSG scheme has a general gain to the quantized network, we provide more experiments as corroborations to support our viewpoint. In our paper, we have conducted a bunch of experiments on our synthetic data to evaluate its effectiveness in improving the network performance, including analyzing the data diversity and integrating the synthetic data generated by our method to different calibration methods and data-driven quantization methods. The experiments show that benefits by the increase in the diversity of synthetic data, the network performance in various methods is significantly improved.

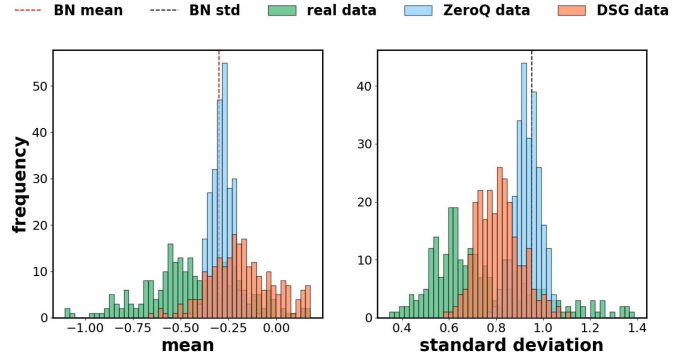


Fig. 6: Mean and standard deviation of the activations in one channel of ResNet18 when feeding different types of data (with 256 samples), including the synthetic data generated by ZeroQ and DSG and the real data. Each sample generated by ZeroQ behaves similarly overfitting BN statistics compared with real data, which shows the homogenization at both distribution and sample level. Our DSG data enjoys the diversity close to real data to obtain the accurate quantized network.

It shows that good diversity is an important property of high-quality synthetic data.

#### 4.3.1 Analysis and Discussion for Data Diversity

In this section, we further discuss the visualization results of our DSG scheme. First, for data-free PTQ, we show the distribution of statistics of the real data, vanilla data (ZeroQ), and DSG data (ours) in Fig. 6, which explains that SDA and LSE alleviate the homogenization from the aforementioned two levels. Then, for data-free QAT, we visualized samples of vanilla data (GDFQ), DSG data (ours), real data, and Gaussian data in Fig. 7 to show the additional effect of SCI that diversifies samples by inhibiting feature correlation. We also visualize the synthetic samples of our method and other generative data-free quantization methods (ZeroQ, GDFQ) in Fig. 8 to visually show the effect of data diversifying.

In Fig. 6, it can be obviously investigated that DSG samples behave more like real data than vanilla data on the offset of mean and standard deviation statistics, which corresponds to the diversity at the distribution level of our generated data scheme. Especially, the SDA plays an important role to make the distribution diverse by slacking the constraint of statistics during the generation process. This phenomenon proves that our DSG scheme diversifies the synthetic data at the distribution level. Moreover, the DSG scheme also generates data with a larger variance compared to the vanilla scheme, which implies that our data samples are widely dispersed and more in line with the real situation. Especially, both SDA and LSE jointly promote diversity at the sample level, which might be useful in providing more content information.

The motivation of the SCI method is to inhibit the correlation among features in the generator, and thus the synthetic samples are separated, as shown in Fig. 7. We visualize the data generated by the vanilla method (GDFQ) and our SCI method, as well as the real data and Gaussian random data. At the right-bottom corner of each subplot, we measure the diversity of each set of features by summarizing the elements in the similarity kernel  $\mathbf{K}_f$  as an index  $s = \sum_{i,j \in f_k} [\mathbf{K}_{i,j}]$ , which indicates the feature distance between samples in

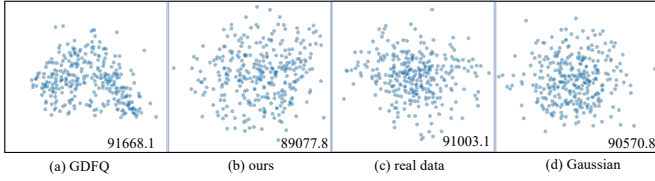


Fig. 7: Visualization of samples generated by (a) GDFQ method, (b) our SCI method, and (c) real data, (d) Gaussian random data. We collect 512 samples for each type of data randomly and reduce dimensions with PCA. Samples generated by our SCI method are more scattered than GDFQ samples and close to that of real data.

spatial perspective. The larger the  $s$ , the more similar the samples. As can be seen from Fig. 7, the Gaussian random data is dispersive and stochastic. And it is also quite obvious that the real data is naturally distinct. However, samples generated by the vanilla GDFQ method, as shown in the subplot (a), are generally more concentrated in comparison, some of which are even overlapped. So it has the biggest  $s$ . Instead, our SCI utilizes the random Gaussian data as initialization and inhibits the correlation between samples to make the data more dispersed than Gaussian sampling. The samples generated by our method, which has the smallest  $s$ , are more dispersed than those generated by the vanilla method and seem close to real data. It demonstrates that the SCI method helps to avoid homogenization from the spatial perspective, and thus diversifies the synthetic samples from each other.

We visualize some synthetic samples of our method and other generative data-free quantization methods (ZeroQ, GDFQ). Take a closer look at Fig. 8, comparing pictures of DSG-PTQ and ZeroQ above, both of which generate images directly and update them iteratively. It is obvious that pictures generated by our DSG-PTQ method seem to have diverse colors with fine-grained textures and coarse-grained figures, while pictures of ZeroQ seem identical and have little difference in between. Meanwhile, below two sets of pictures are generated by DSG-QAT and GDFQ respectively. Owing to the generator, there are much more possibilities in the generation process and finally exhibit in both sets of images. It can be seen with naked eyes that samples of our DSG-QAT have distinct colors with higher saturation degrees, and white images inside the pictures have uncertain patterns. However, although samples generated by GDFQ have different colors and details inside the pictures as well, these samples have a uniform style in texture, which is tanglesome and poor in diversity.

### 4.3.2 Evaluation with Calibration Methods

Additionally, to verify the versatility and robustness of the performance gains of our DSG method, we evaluate our DSG scheme with different calibration methods and compare it with other data generation methods under the same setting. The calibration methods include Percentile, EMA, and MSE. As shown in TABLE 5, no matter which calibration methods is integrated, our DSG scheme substantially outstrips ZeroQ by 6.52%, 9.36%, and 0.61%, respectively. The results strongly suggest that DSG can achieve the leading performance in

Method	No D	W-bit	A-bit	Quant	Top-1
Baseline	-	-	32	32	71.47
Real Data	✗	4	4	Vanilla	31.86
ZeroQ	✓	4	4	Vanilla	26.04
DSG (Ours)	✓	4	4	Vanilla	<b>39.90</b>
Real Data	✗	4	4	Percentile	42.83
ZeroQ	✓	4	4	Percentile	32.24
DSG (Ours)	✓	4	4	Percentile	<b>38.76</b>
Real Data	✗	4	4	EMA	42.67
ZeroQ	✓	4	4	EMA	32.31
DSG (Ours)	✓	4	4	EMA	<b>41.67</b>
Real Data	✗	4	4	MSE	41.45
ZeroQ	✓	4	4	MSE	39.39
DSG (Ours)	✓	4	4	MSE	<b>40.00</b>

TABLE 5: Evaluation with calibration on ImageNet. We evaluate our DSG scheme with various calibration methods (Percentile, EMA, MSE) on ResNet18, and Vanilla means the calibration method adopted by ZeroQ, which simply obtain the quantizer by the maximum and minimum of the weight and activation.

Method	No D	No FT	W-bit	A-bit	Top-1
Baseline	-	-	32	32	69.76
Real Data	✗	✓	6	6	59.16
ZeroQ	✓	✓	6	6	58.12
DSG (Ours)	✓	✓	6	6	<b>58.69</b>
Real Data	✗	✓	8	8	69.22
ZeroQ	✓	✓	8	8	65.75
DSG (Ours)	✓	✓	8	8	<b>68.88</b>

TABLE 6: Evaluation with DFQ on ImageNet. We use cross-layer equalization and bias correction proposed by DFQ to perform per-layer quantization on ResNet18.

various experimental settings, and thus the improvement is versatile and robust for various calibration methods.

We also evaluate the synthetic data with DFQ [38], which is a calibration method for data-free quantization. Specifically, DFQ has proposed cross-layer range equalization to equalize the different channel ranges of weight in per-layer quantization and bias correction which is to eliminate the biased quantization error. Both of the two techniques rely on the statistics of BN layers following the convolution layer. Therefore, BN layers are needed to calibrate the corresponding activations, so they have to proceed behind each convolution layer, which results in DFQ only working on specific network architectures and cannot be commonly practiced. Fortunately, generative methods, such as ZeroQ and our DSG, can work on arbitrary architectures, and the statistics of activations can take the place of BN statistics and thus be used in the DFQ method. TABLE 6 shows the closeups of two generative data-free quantization methods, *i.e.*, ZeroQ and DSG, in conjunction with the DFQ method. Results show that our DSG outperforms ZeroQ by 0.57% and 3.13% in W6A6 and W8A8 cases.

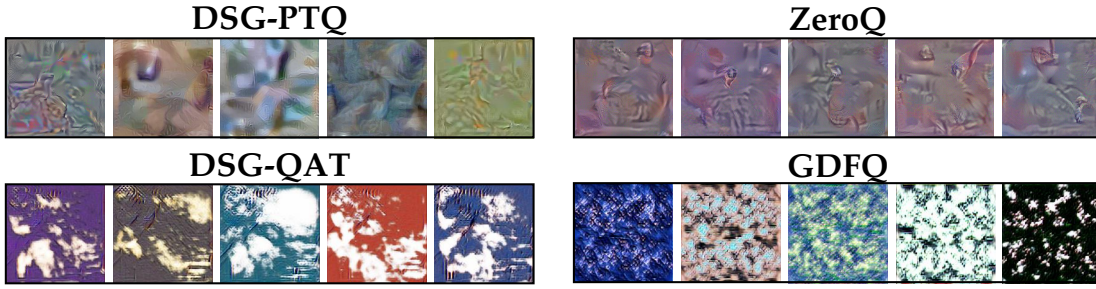


Fig. 8: Visualization of synthetic data for ZeroQ, GDFQ, and our DSG method in PTQ and QAT scenario. Each showcases 5 examples.

Arch	Method	No D	Label	Prior	W-bit	A-bit	Top-1
ResNet18	Real Data	✗	✗	✗	3	32	64.16
	ZeroQ	✓	✗	✗	3	32	49.86
	DSG (Ours)	✓	✗	✗	3	32	<b>56.09</b>
	DSG (Ours)	✓	✓	✗	3	32	<b>58.27</b>
	DSG (Ours)	✓	✓	✓	3	32	<b>61.32</b>
	Real Data	✗	✗	✗	4	32	68.42
	ZeroQ	✓	✗	✗	4	32	63.86
	DSG (Ours)	✓	✗	✗	4	32	<b>66.87</b>
	DSG (Ours)	✓	✓	✗	4	32	<b>67.09</b>
	DSG (Ours)	✓	✓	✓	4	32	<b>67.78</b>
	Real Data	✗	✗	✗	5	32	69.21
	ZeroQ	✓	✗	✗	5	32	68.39
	DSG (Ours)	✓	✗	✗	5	32	<b>68.97</b>
	DSG (Ours)	✓	✓	✗	5	32	<b>69.02</b>
	DSG (Ours)	✓	✓	✓	5	32	<b>69.16</b>
	MobileNetV2	Real Data	✗	✗	✗	3	32
ZeroQ		✓	✗	✗	3	32	11.07
DSG (Ours)		✓	✓	✓	3	32	<b>45.40</b>
Real Data		✗	✗	✗	4	32	68.37
ZeroQ		✓	✗	✗	4	32	56.16
DSG (Ours)		✓	✓	✓	4	32	<b>58.13</b>

TABLE 7: AdaRound on ImageNet with ResNet18 and MobileNetV2. We evaluate the DSG scheme on AdaRound, one of the SOTA methods of data-driven post-training quantization, which learns how to quantize weights using several batches of unlabeled samples. We adopt "Label" [22] and image prior ("Prior") [55] techniques to evaluating our DSG scheme further.

### 4.3.3 Evaluation with Data-driven Quantization Methods

Experiments above are conducted on calibration methods that optimize the clipping value for activations. We further evaluate our DSG data with a novel data-driven PTQ method named Adaround [37], which utilizes a rounding approach to quantize weights. Meanwhile, we introduce two other data generation tricks into our DSG scheme, *e.g.*, generating data with labels [22] provides class information from parameters, image prior [55] avoids generating unpractical scenes or unrecognizable patterns. We have conducted different bit-

width for both weights and activations on different network architectures including ResNet18 and MobileNetV2 (TABLE 7). And we generate 1024 samples in every single experiment. TABLE 7 shows that our DSG scheme outperforms ZeroQ by a wide margin when solely quantizing weights and preserving full-precision for activations. It is notable that in ultra-low bit-width settings (*i.e.*, 3-bit), our DSG surpasses ZeroQ by 11.46% on ResNet18 and a surprising 34.33% on MobileNetV2. We also tried quantizing activations to 8-bit, and the results prove that our DSG is more robust to quantization of parameters, which outstrips ZeroQ by 6.06% with ResNet18 on ImageNet. Besides, we have the observation that the model gains further improvements when these tricks are applied together, even close the performance applying real data. Because our DSG is orthogonal with labels and image prior methods, so these methods can jointly boost the accuracy performance without any inconsistency.

## 5 CONCLUSION

In this paper, we first revisit the sample generation process in generative data-free PTQ and QAT quantization and then give a theoretical analysis that the diversity of synthetic samples is crucial for the data-free quantization and reveal the homogenization of synthetic data in the distribution and sample levels. In this paper, we propose a novel Diverse Sample Generation (DSG) scheme for generative data-free quantization, to address the deficiencies of previous methods which severely debases the quality of the synthetic data and further harms the performance of the quantized network. Our scheme has been evaluated on a variety of bit widths and neural architectures, and the results forcefully demonstrate the effectiveness and versatility of the DSG scheme. It shows notable accuracy improvements in ultra-low bit-width cases (*e.g.* W4A4). Moreover, benefiting from the enhanced diversity, the performance of the network is significantly improved in various methods integrating with synthetic data, which demonstrates that diversity is an important property of high-quality synthetic data. We hope our work can provide directions for future research on data-free quantization.

## ACKNOWLEDGMENTS

This work was supported by National Natural Science Foundation of China (62022009, 61872021), Beijing Nova Program of Science and Technology (Z191100001119050), State Key Lab of Software Development Environment (SKLSDE-2020ZX-06).

## REFERENCES

- [1] R Banner, Y Nahshan, E Hoffer, and D Soudry. Post training 4-bit quantization of convolution networks for rapid-deployment. *CoRR, abs/1810.05723*, 2018.
- [2] Ron Banner, Yury Nahshan, Elad Hoffer, and Daniel Soudry. Post-training 4-bit quantization of convolution networks for rapid-deployment, 2019.
- [3] Alexei Borodin. Determinantal point processes. *arXiv preprint arXiv:0911.1153*, 2009.
- [4] Adrian Bulat and Georgios Tzimiropoulos. Hierarchical binary cnns for landmark localization with limited resources. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 42(2):343–356, 2020.
- [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.
- [6] Laming Chen, Guoxin Zhang, and Hanning Zhou. Fast greedy map inference for determinantal point process to improve recommendation diversity. In *Conference on Neural Information Processing Systems (NeurIPS)*, 2018.
- [7] Tianqi Chen, Mu Li, Yutian Li, Min Lin, Naiyan Wang, Minjie Wang, Tianjun Xiao, Bing Xu, Chiyuan Zhang, and Zheng Zhang. Mxnet: A flexible and efficient machine learning library for heterogeneous distributed systems. *arXiv preprint arXiv:1512.01274*, 2015.
- [8] Yoojin Choi, Jihwan Choi, Mostafa El-Khamy, and Jungwon Lee. Data-free network quantization with adversarial knowledge distillation, 2020.
- [9] Yoojin Choi, Mostafa El-Khamy, and Jungwon Lee. Zero-shot learning of a conditional generative adversarial network for data-free network quantization. In *2021 IEEE International Conference on Image Processing (ICIP)*, pages 3552–3556. IEEE, 2021.
- [10] Yoni Choukroun, Eli Kravchik, Fan Yang, and Pavel Kisilev. Low-bit quantization of neural networks for efficient inference. In *2019 IEEE/CVF International Conference on Computer Vision Workshop (ICCVW)*, pages 3009–3018. IEEE, 2019.
- [11] Jia Deng, Wei Dong, Richard Socher, Li Jia Li, Kai Li, and Fei Fei Li. Imagenet: a large-scale hierarchical image database. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2009.
- [12] Yueqi Duan, Jiwen Lu, Ziwei Wang, Jianjiang Feng, and Jie Zhou. Learning deep binary descriptor with multi-quantization. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 41(8):1924–1938, 2019.
- [13] Mohamed Elfeki, Camille Couprie, Morgane Riviere, and Mohamed Elhoseiny. Gdpp: Learning diverse generations using determinantal point processes. In *International Conference on Machine Learning (ICML)*, 2019.
- [14] Mark Everingham, Luc Van Gool, Christopher K. I. Williams, John Winn, and Andrew Zisserman. The pascal visual object classes challenge. *International Journal of Computer Vision (IJCV)*, 2010.
- [15] DF Frey and RA Pimentel. Principal component analysis and factor analysis. 1978.
- [16] Amir Gholami, Kiseok Kwon, Bichen Wu, Zizheng Tai, Xiangyu Yue, Peter Jin, Sicheng Zhao, and Kurt Keutzer. Squeezenext: Hardware-aware neural network design, 2018.
- [17] Ross Girshick. Fast r-cnn. In *International Conference on Computer Vision (ICCV)*, 2015.
- [18] Ross Girshick, Jeff Donahue, Trevor Darrell, and Jitendra Malik. Rich feature hierarchies for accurate object detection and semantic segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2014.
- [19] Cong Guo, Yuxian Qiu, Jingwen Leng, Xiaotian Gao, Chen Zhang, Yunxin Liu, Fan Yang, Yuhao Zhu, and Minyi Guo. SQuant: On-the-fly data-free quantization via diagonal hessian approximation. In *International Conference on Learning Representations*, 2022.
- [20] Suyog Gupta, Ankur Agrawal, Kailash Gopalakrishnan, and Pritish Narayanan. Deep learning with limited numerical precision. In *International Conference on Machine Learning (ICML)*, pages 1737–1746, 2015.
- [21] 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.
- [22] M. Haroush, I. Hubara, E. Hoffer, and D. Soudry. The knowledge within: Methods for data-free model compression. In *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 8491–8499, 2020.
- [23] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016.
- [24] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 770–778, 2016.
- [25] Xiangyu He, Jiahao Lu, Weixiang Xu, Qinghao Hu, Peisong Wang, and Jian Cheng. Generative zero-shot network quantization. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3000–3011, 2021.
- [26] Benoit Jacob, Skirmantas Kligys, Bo Chen, Menglong Zhu, Matthew Tang, Andrew Howard, Hartwig Adam, and Dmitry Kalenichenko. Quantization and training of neural networks for efficient integer-arithmetic-only inference. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2704–2713, 2018.
- [27] Krizhevsky, Alex, Sutskever, Ilya, Hinton, and E. Geoffrey. Imagenet classification with deep convolutional neural networks. *Communications of the ACM*, 2017.
- [28] Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. 2009.
- [29] Alex Kulesza and Ben Taskar. Determinantal point processes for machine learning. *arXiv preprint arXiv:1207.6083*, 2012.
- [30] Chun-Liang Li, Wei-Cheng Chang, Yu Cheng, Yiming Yang, and Barnabás Póczos. Mmd gan: Towards deeper understanding of moment matching network. *Advances in neural information processing systems*, 30, 2017.
- [31] Fuxin Li, Yunshan Fu, Yu-Hong Dai, Cristian Sminchisescu, and Jue Wang. Kernel learning by unconstrained optimization. In *Artificial Intelligence and Statistics*. PMLR, 2009.
- [32] Yuhang Li, Ruihao Gong, Xu Tan, Yang Yang, Peng Hu, Qi Zhang, Fengwei Yu, Wei Wang, and Shi Gu. Brcq: Pushing the limit of post-training quantization by block reconstruction. In *International Conference on Learning Representations*, 2021.
- [33] Fanghui Liu, Xiaolin Huang, Yudong Chen, and Johan AK Suykens. Random features for kernel approximation: A survey on algorithms, theory, and beyond. *IEEE Transactions on Pattern Analysis & Machine Intelligence*, (01):1–1, 2021.
- [34] Jing Liu, Bohan Zhuang, Zhuangwei Zhuang, Yong Guo, Junzhou Huang, Jinhui Zhu, and Mingkui Tan. Discrimination-aware network pruning for deep model compression. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, pages 1–1, 2021.
- [35] Yuang Liu, Wei Zhang, and Jun Wang. Zero-shot adversarial quantization. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2021.
- [36] A. Mordvintsev, Christopher Olah, and M. Tyka. Inceptionism: Going deeper into neural networks. 2015.
- [37] Markus Nagel, Rana Ali Amjad, Mart van Baalen, Christos Louizos, and Tijmen Blankevoort. Up or down? adaptive rounding for post-training quantization, 2020.
- [38] Markus Nagel, Mart van Baalen, Tijmen Blankevoort, and Max Welling. Data-free quantization through weight equalization and bias correction. In *International Conference on Computer Vision (ICCV)*, 2019.
- [39] Augustus Odena, Christopher Olah, and Jonathon Shlens. Conditional image synthesis with auxiliary classifier gans. In *International Conference on Machine Learning (ICML)*, 2017.
- [40] Jiangmiao Pang, Kai Chen, Jianping Shi, Huajun Feng, Wanli Ouyang, and Dahua Lin. Libra r-cnn: Towards balanced learning for object detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019.
- [41] Eunhyeok Park, Sungjoo Yoo, and Peter Vajda. Value-aware quantization for training and inference of neural networks, 2018.
- [42] Haotong Qin, Zhongang Cai, Mingyuan Zhang, Yifu Ding, Haiyu Zhao, Shuai Yi, Xianglong Liu, and Hao Su. Bipointnet: Binary neural network for point clouds, 2020.
- [43] Haotong Qin, Ruihao Gong, Xianglong Liu, Xiao Bai, Jingkuan Song, and Nicu Sebe. Binary neural networks: A survey. *Pattern Recognition*, 2020.
- [44] Haotong Qin, Ruihao Gong, Xianglong Liu, Mingzhu Shen, Ziran Wei, Fengwei Yu, and Jingkuan Song. Forward and backward information retention for accurate binary neural networks. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2020.
- [45] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster r-cnn: Towards real-time object detection with region proposal networks, 2016.
- [46] Mark Sandler, Andrew G. Howard, Menglong Zhu, Andrey Zhmoginov, and Liang-Chieh Chen. Mobilenetv2: Inverted residuals and linear bottlenecks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018.
- [47] Xu Shoukai, Li Haokun, Zhuang Bohan, Liu Jing, Cao Jiezhong,



- Liang Chuangrun, and Tan Mingkui. Generative low-bitwidth data free quantization. In *European Conference on Computer Vision (ECCV)*, 2020.
- [48] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014.
- [49] Wonyong Sung, Sungho Shin, and Kyuyeon Hwang. Resiliency of deep neural networks under quantization. *arXiv preprint arXiv:1511.06488*, 2015.
- [50] Christian Szegedy, Vincent Vanhoucke, Sergey Ioffe, Jon Shlens, and Zbigniew Wojna. Rethinking the inception architecture for computer vision. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2818–2826, 2016.
- [51] Zhiqiang Tang, Xi Peng, Kang Li, and Dimitris N. Metaxas. Towards efficient u-nets: A coupled and quantized approach. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 42(8):2038–2050, 2020.
- [52] Frederick Tung and Greg Mori. Deep neural network compression by in-parallel pruning-quantization. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 42(3):568–579, 2020.
- [53] Ziwei Wang, Jiwen Lu, Ziyi Wu, and Jie Zhou. Learning efficient binarized object detectors with information compression. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, pages 1–1, 2021.
- [54] Ziwei Wang, Jiwen Lu, and Jie Zhou. Learning channel-wise interactions for binary convolutional neural networks. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, pages 1–1, 2020.
- [55] Hongxu Yin, Pavlo Molchanov, Zhizhong Li, Jose M. Alvarez, Arun Mallya, Derek Hoiem, Niraj K. Jha, and Jan Kautz. Dreaming to distill: Data-free knowledge transfer via deepinversion, 2020.
- [56] Xiangguo Zhang, Haotong Qin, Yifu Ding, Ruihao Gong, Qinghua Yan, Renshuai Tao, Yuhang Li, Fengwei Yu, and Xianglong Liu. Diversifying sample generation for accurate data-free quantization. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2021.
- [57] Xiangyu Zhang, Xinyu Zhou, Mengxiao Lin, and Jian Sun. Shufflenet: An extremely efficient convolutional neural network for mobile devices. *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018.
- [58] Ritchie Zhao, Yuwei Hu, Jordan Dotzel, Christopher De Sa, and Zhiru Zhang. Improving neural network quantization without retraining using outlier channel splitting. *arXiv preprint arXiv:1901.09504*, 2019.
- [59] Bohan Zhuang, Chunhua Shen, Mingkui Tan, Lingqiao Liu, and Ian Reid. Structured binary neural networks for accurate image classification and semantic segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019.
- [60] Bohan Zhuang, Mingkui Tan, Jing Liu, Lingqiao Liu, Ian Reid, and Chunhua Shen. Effective training of convolutional neural networks with low-bitwidth weights and activations. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, pages 1–1, 2021.

## APPENDIX A

### MAIN PROOFS AND DISCUSSION

#### A.1 Proof of Lemma 1

**Lemma A.1.** *For any input domains  $\mathcal{X}$  that includes multiples classes (at least 2) of samples, it can be modeled as several independent high-density  $\{\mathcal{R}_{H1}, \dots, \mathcal{R}_{Hh}\}$  and low-density  $\{\mathcal{R}_{L1}, \dots, \mathcal{R}_{Ll}\}$  sub-regions divided by possible decision surfaces, where  $h \geq 1$  and  $l \geq 0$ .*

*Proof.* As discussed in Section 3.2, the given input domain  $\mathcal{X}$  follows *Low-density assumption* and *Smoothness assumption*. And we use the method of proof by contradiction. If Lemma 1 is false, then at least one of the following is true:

- Low-density sub-regions in the input domain do not exist;
- High-density sub-regions in the input domain do not exist or are only one.

And below we prove that none of the above is true.

(1) Let us assume that the first term holds.

When there is no low-density area, according to the low-density assumption, the possible decision surface can only pass through the low-density area, so there is no possible decision surface in the input domain, that is, the input domain contains and contains only one meaningful category. This violates the premise that the input domain contains at least two classes.

(2) Then we assume that the second term holds.

When high-density regions do not exist in the input domain, the entire input domain is a low-density region. According to the smoothness assumption, since there is no high-density path between any two samples in the input domain, they cannot be classified into the same category by the decision surface. Therefore, any two samples in the input domain belong to different classes, that is, there is no determinable class division. This in principle violates the premise that the input domain contains multiple classes, that is, the input domain should be classifiable rather than distinct everywhere.

When there is only one high-density region in the input domain, also according to the smoothness assumption, only two samples in this high-density region can form a high-density path, i.e. belong to the same class. In other words, the input domain contains one and only one meaningful class, which also violates the premise that the input domain contains at least two classes.

So, to sum up, Lemma 1 must be true.  $\square$

#### A.2 Proof of Theorem 1

**Theorem A.1.** *Given a set of all possible input domains  $\mathbf{X} = \{\mathcal{X}_0, \mathcal{X}_1, \dots\}$ , whose  $i$ -th element can be denoted as  $\mathcal{X}_i$  with scale  $V^i$  and consists of several sub-regions  $\{\mathcal{R}_1^i, \dots, \mathcal{R}_{K^i}^i\}$  with scales  $\{V_1^i, \dots, V_{K^i}^i\}$ , and the number  $K^i \geq 2$  is unknown yet limited. Consider a sample set  $\mathbf{x}^s = \{x_0^s, \dots, x_N^s\} \subset \mathcal{X}^*$ , where  $\mathcal{X}^* = \mathbb{E}(\mathbf{X})$  denotes the potential input domain and the differences inside each sub-region of  $\mathcal{X}$  is neglected. When the set  $\mathbf{x}^s$  satisfies that for  $\forall x_i^s \in \mathbf{x}^s$ ,  $p(x_i^s \in \mathcal{R}_j^*) = \frac{V_j^*}{V^*}$ , the information reflecting from all possible input domains  $\mathbf{X}$  by the sample set  $\mathbf{x}^s$  will be the maximized in mathematical expectation, where  $V^* = \sum_{k=0}^{K^*} V_k^*$ .*

*Proof.* First, we discuss the properties of the potential input field  $X^* = E(X)$ . For the set of all possible input domains  $\mathbf{X} = \{\mathcal{X}_0, \mathcal{X}_1, \dots\}$ , since the consistent modeling among all its possible elements (as shown in Lemma 1), so for the defined potential input domain  $X^*$ , the number of its sub-regions  $K^*$  is  $K^* = \max(K^0, K^1, \dots)$ .

For any  $i$ -th input domain  $X_i$ , the  $K^* > K^i, R_{K^i+1} \dots R_{K^*}$  can be regarded as empty sub-regions, and there is no specific ordering among all its sub-regions  $R_0, R_1 \dots R_{K^*}$ . For the potential input domain  $X^*$ , the mathematical expectation of the scale of any sub-region  $R_j^*$  is  $V_j^* = \mathbb{E}_i \mathbb{E}_j(V_j^i)$ , so  $\forall i, j, V_i^* = V_j^*$ , i.e.,  $\forall j, V_j^* = V^*/K^*$ . And for  $X^* = \mathbb{E}(\mathbf{X})$ , since the density properties of each possible input domain are completely random,  $X^*$  can be seen as uniform in density in expectation.

Therefore, consider the set  $\mathbf{x}^s = \{x_0^s, \dots, x_N^s\}$  sampled from  $\mathcal{X}^*$ , for any sample  $x_i^s$ , the probability that it belongs to the  $j$ -th sub-region  $\mathcal{R}_j^*$  is  $p(x_i^s \in \mathcal{R}_j^*)$ . Since differences within sub-regions are ignored, samples in the same sub-region can be considered to have the exact same class. Therefore, maximizing the information amount of the sample set  $\mathbf{x}^s$  to reflect the latent region  $\mathcal{X}^*$  can be expressed as:

$$\max \mathcal{H}(\mathbf{x}^s) = - \sum_i \sum_j p(x_i^s \in \mathcal{R}_j^*) \log p(x_i^s \in \mathcal{R}_j^*) \quad (12)$$

Since the sampling of all samples is completely independent, the above formula is equivalent to

$$\sum_i \max \mathcal{H}_i(\mathbf{x}^s) = \sum_i \left( \max - \sum_j p(x_i^s \in \mathcal{R}_j^*) \log p(x_i^s \in \mathcal{R}_j^*) \right). \quad (13)$$

Since the samples are non-specific, we just need to discuss  $\max \mathcal{H}_i(\mathbf{x}^s)$  here. We simplify  $p(x_i^s \in \mathcal{R}_j^*)$  as  $p_j$ , and the optimization problem can be defined as:

$$\max \mathcal{H}_i(\mathbf{x}^s) = \max - \sum_{j=1}^{K^*} p_j \log p_j. \quad (14)$$

We introduce Lagrangian multiplier  $\lambda$ , the constructed Lagrangian function is:

$$L(p_j, \lambda) = - \sum_{j=1}^{K^*} p_j \log p_j + \lambda \left( \sum_{j=1}^{K^*} p_j - 1 \right) \quad (15)$$

And then solve as:

$$\begin{cases} \frac{\partial L(p_j, \lambda)}{\partial p_j} = 0 \\ \sum_{j=1}^{K^*} p_j = 1 \end{cases} \quad (16)$$

Then it is available that:

$$\frac{\partial L(p_j, \lambda)}{\partial p_j} = 0 \quad (17)$$

$$\Rightarrow \frac{\partial \left[ - \left( p_j \log p_j + \sum_{i \neq j} p_i \log p_i \right) + \lambda \left( p_j + \sum_{i \neq j} p_i - 1 \right) \right]}{\partial p_i} = 0 \quad (18)$$

$$\Rightarrow -(\log p_j + 1) + \lambda = 0 \quad (19)$$

$$\Rightarrow p_j = \frac{1}{K^*}. \quad (20)$$

Put  $p_j = 2^{\lambda-1}$  into  $\sum_{j=1}^{K^*} p_j = 1$  to get:

$$\sum_{i=1}^n 2^{\lambda-1} = 1 \Rightarrow 2^{\lambda-1} = \frac{1}{n} \Rightarrow p_i = \frac{1}{n}. \quad (21)$$

Bring  $p_i = \frac{1}{n}$  into  $H(X) = -\sum_{i=1}^n p_i \log p_i$  to get  $H(X) = \log n$ , so  $H(X) \leq \log n$ ; Therefore, when  $p_j = \frac{1}{K^*}$ , the information entropy  $\mathcal{H}(x^s)$  is maximized.

Since according to the above proof,  $\forall j, V_j^* = \frac{V^*}{K^*}$ , so when  $p(x_i^s \in \mathcal{R}_j^*) = \frac{V_j^*}{V^*} = 1/K^*$ , the information entropy  $\mathcal{H}(x^s)$  is maximized, that is, the sample set  $x^s$  can maximize the amount of information that reflects the potential area  $\mathcal{X}^*$  at this time. And Theorem 1 is proved.  $\square$

## APPENDIX B VISUALIZATION

In Fig. 9 we show the visualization of more synthetic samples. Compared to the synthetic samples in existing generative data-free methods like ZeroQ and GDFQ, the synthetic samples produced by our DSG have more diverse colors and textures in visualization, which verifies from another aspect that DSG improves synthetic samples through diversification.

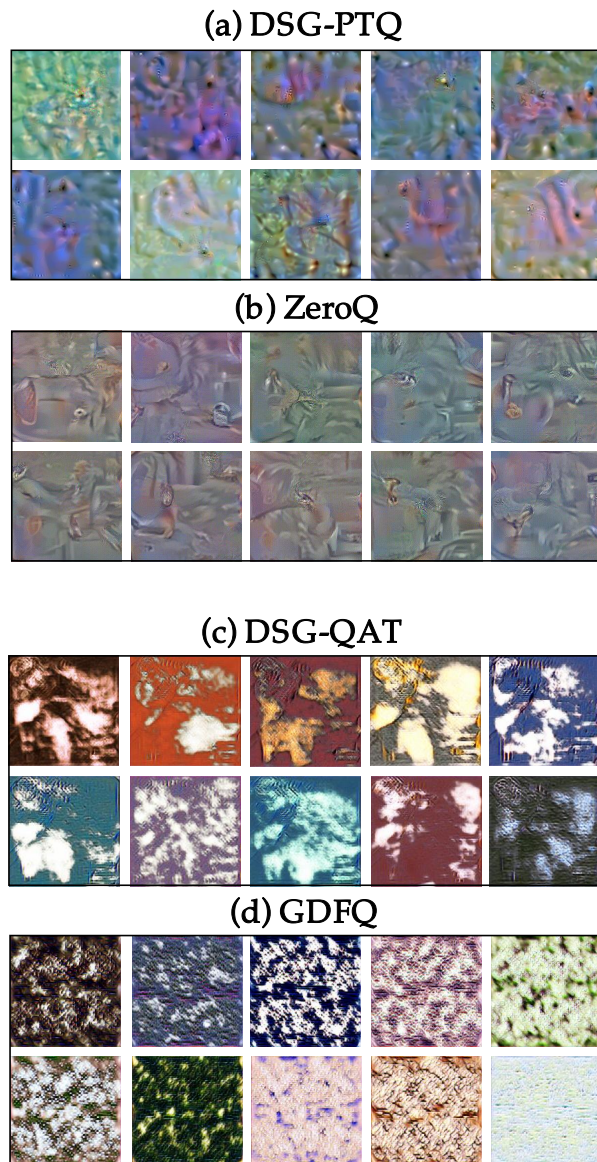


Fig. 9: Visualization of more synthetic data for ZeroQ, GDFQ, and our DSG method in PTQ and QAT approaches. Each showcases 10 examples.