# IMPROVING LOW-BIT POST TRAINING QUAN TIZATION: A DATA-FREE APPROACH

Anonymous authors

Paper under double-blind review

#### ABSTRACT

Post-training quantization (PTO) without access to real data is enabling efficient model optimization and deployment in scenarios where privacy or proprietary constraints restrict the use of original datasets. Traditional data free quantization methods rely on Batch Normalization (BN) statistics from the trained fullprecision model to generate calibration dataset for quantization. However, this reliance on BN statistics limits their applicability to deep neural networks (DNNs) without BN layer such as AlexNet. In this paper, we propose a calibration dataset generation algorithm that is agnostic to BN statistics, leveraging just the backpropagation to create synthetic images for PTQ. We also demonstrate that it is not necessary to include samples from every target category in the calibration dataset to get the representative activation ranges for quantization. Extensive experiments with both large and lightweight models on large-scale image classification tasks demonstrate that our method consistently improves quantization performance across various DNN architectures, especially in low-bit settings. Notably, in 4-bit quantization, we achieve an improvement of 3.42% in top-1 accuracy for the ResNet18 model and 3.14% for the InceptionV3 model compared to the stateof-the-art (SOTA) DSG method. Importantly, we use very few synthetic samples for quantization compared to other methods.

027 028 029

030

025

026

004

010 011

012

013

014

015

016

017

018

019

021

#### 1 INTRODUCTION

031 DNNs have achieved remarkable success across a wide range of applications, driven by advance-032 ments in computational power, large-scale datasets, and innovative algorithms (LeCun et al., 2015). 033 These networks have revolutionized the field of computer vision. They power image classifica-034 tion (Krizhevsky et al., 2012), object detection (Girshick et al., 2014), and semantic segmentation (Long et al., 2015) systems with unprecedented accuracy. With these domains, DNNs have 035 made significant strides in robotics, enhancing the ability of robots to perceive and interact with their environments autonomously (Mnih et al., 2015). In autonomous driving, deep learning mod-037 els play a crucial role in enabling vehicles to understand and navigate complex road conditions safely (Bojarski et al., 2016). The medical field has also seen transformative impacts, with DNNs aiding in disease diagnosis, medical imaging analysis, and personalized treatment planning (Esteva 040 et al., 2017). However, deploying these neural networks on resource-constrained devices remains 041 a considerable challenge due to their substantial memory requirements and intensive computational 042 demands (Howard, 2017). The advancement of hardware capable of low-precision computations 043 has made quantization a favored method for addressing these challenges (Jacob et al., 2018). Quan-044 tizing the floating-point values of weights and/or activations in a neural network to integers can significantly reduce the model size without altering the architecture.

Quantization methods are generally categorized into two types: quantization-aware training (QAT) and PTQ. QAT incorporates quantization into the training process, allowing the model to learn and adapt to the quantized weights and activations. This approach typically results in higher accuracy compared to PTQ, as the model is optimized to perform well under the constraints of quantization (Choi et al., 2018). While the performance degradation is minimal with QAT, the process can be computationally intensive and time-consuming (Nahshan et al., 2021). PTQ involves applying quantization techniques to a pre-trained model without training process (Cai et al., 2020; Qin et al., 2023). One of the challenges in quantization is determining the range of activation values in DNN. Many quantization methods are designed for data-driven scenarios, requiring access to either training or

054 validation data, or relying on a small set of unlabeled calibration data to accurately represent the activation values (Banner et al., 2019). However, real data may not always be available due to some 056 constraints, particularly with proprietary data. To mitigate this, data-free quantization techniques 057 have been developed to allow for the quantization without needing access to real data (Cai et al., 058 2020). Among existing data-free quantization methods, generative approaches create calibration data by aligning the distribution of synthetic data with the BN statistics of the trained full-precision model (Xu et al., 2020). Generative approaches can greatly reduce the accuracy gap between quan-060 tized models that use synthetic data and those that use real data (Qin et al., 2023; Cai et al., 2020). 061 However, relying on BN statistics limits the applicability of these methods to neural networks which 062 don't have BN layer such as AlexNet (Krizhevsky et al., 2012). In our work, we demonstrate that it 063 is not necessary to depend on the BN statistics of the full-precision model to generate a calibration 064 dataset for achieving optimal performance from the quantized model. Additionally, the common 065 practice when curating a calibration dataset is to include at least one image from each target class 066 to calibrate the image classification model (Zhang et al., 2023). We show that optimal performance 067 can be achieved without having a sample from each category and that fewer samples can suffice. 068

In this paper, we primarily explore the method of generating synthetic calibration dataset for PTQ of DNN models. We demonstrate that the BN statistics of the original floating-point model are not necessary to create effective synthetic data. The main contributions of this paper are:

- We propose a method to generate synthetic data using the trained full-precision model agnostic to BN statistics, making our approach applicable to any model architecture.
- We experimentally demonstrate that it is not necessary to include a sample from each target category in the calibration dataset; selecting only a few target classes is sufficient to create an effective calibration dataset which in turn reduces PTQ time.
- Through extensive PTQ comparisons on several standard network architectures such as ResNet18/20/50 (He et al., 2016), SqueezeNext (Gholami et al., 2018), InceptionV3 (Szegedy et al., 2016), and lightweight architectures like ShuffleNet (Zhang et al., 2018), and MobileNetV2/V3 (Sandler et al., 2018; Howard et al., 2019), we show that our method significantly outperforms existing generative quantization methods with calibration. Specifically, in 4-bit precision setting on the ResNet18 and InceptionV3 models, our method improves the top-1 accuracy by over 3.42% and 3.14% compared to the SOTA DSG method.

#### 084 085

087

073

075

076

077

078

079

081

082

# 2 RELATED WORK

In this section, we first review and organize existing research on quantization into two primary methodologies: QAT and PTQ. Additionally, we categorize these methods based on whether they require any form of training or validation data for the quantization process.

091 092

## 2.1 QUANTIZATION AWARE TRAINING

QAT is an advanced technique in which the quantization process is integrated into the training phase 094 of the neural network. (Jacob et al., 2018) integrate low-precision computations typically used dur-095 ing inference into the forward pass of training, while maintaining standard backpropagation with 096 floating-point weights and biases to allow precise adjustments. This approach allows the model 097 to learn and adapt to the quantization noise, resulting in higher accuracy for the quantized model 098 compared to PTQ. (Courbariaux et al., 2016) introduce the concept of binarized neural networks (BNNs), where both weights and activations are constrained to +1 or -1. They replace most arith-100 metic operations with bit-wise operations to reduce memory size. (Krishnamoorthi, 2018) evaluated 101 various quantization methods and bit-widths, revealing that even with per-channel quantization, 102 lightweight networks do not reach baseline accuracy with INT8 PTQ and require QAT. Using an 103 annealing learning rate schedule and a very small final learning rate, (McKinstry et al., 2018) show 104 that many networks can be fine-tuned for just one epoch after quantizing to INT8 and still reach 105 baseline accuracy. (Zhou et al., 2016) introduce DoReFa-Net, a framework for training neural networks with low bit-width weights, activations, and gradients. They also propose QAT techniques to 106 ensure effective learning with low precision. (Rastegari et al., 2016) presented XNOR-Net, a QAT-107 trained binary CNN architecture, achieving impressive accuracy on ImageNet and demonstrating

the feasibility of highly efficient binary networks for real-world tasks. PACT (Choi et al., 2018)
 dynamically adjusts the clipping values to minimize quantization error, leading to higher accuracy,
 especially for low-bit width quantization.

QAT maintains quantized model accuracy close to the full precision model while significantly reducing the model size and computational complexity. However, it requires more computational resources and time due to additional quantization operations, making it more complex and resourceintensive than standard training (Nahshan et al., 2021). Additionally, QAT implementation is challenging, needing modifications to the training pipeline.

116 117 118

#### 2.2 POST TRAINING QUANTIZATION

119 PTQ is a widely used technique aimed at reducing the memory and computational requirements of 120 neural networks without the need for retraining. Unlike QAT, which incorporates quantization into 121 the training process, PTQ is applied after the model has been fully trained using a small calibration 122 dataset created from the original training or validation dataset. To quantize activations in DNNs, it is essential to determine the activation ranges. (Banner et al., 2019) address this by analytically find-123 ing the activation clipping ranges. They demonstrate that the activations typically follow bell-curve 124 statistics, fitting either Laplace or Gaussian distributions, and they formulate the clipping process as 125 an optimization problem. (Zhao et al., 2019) introduce Outlier Channel Splitting (OCS), a technique 126 designed to handle outliers in the distribution of weights and activations that can negatively impact 127 quantization. OCS works by identifying a small number of channels containing outliers, duplicating 128 them, and then halving the values in those channels, which moves the affected outliers toward the 129 center of the distribution. (Li et al., 2021) introduce a novel method called Block Reconstruction 130 Quantization (BRECQ), which focuses on optimizing the quantization of neural network weights in 131 a block-wise manner. BRECQ achieves this through local optimization and layer-wise error correc-132 tion, significantly reducing quantization-induced errors without requiring retraining. (Nagel et al., 133 2020) present a novel adaptive rounding technique, AdaRound, that dynamically determines whether to round weights and activations up or down during quantization. This optimization-based approach 134 minimizes the quantization error and preserves model accuracy without requiring retraining. Unlike 135 AdaRound, which confines quantized weights to be within  $\pm 1$  of their rounded values, AdaQuant 136 (Hubara et al., 2021) grants more freedom. It achieves this by independently optimizing each layer's 137 weights and quantization parameters, using the calibration set to minimize the mean-squared-error 138 between the original and quantized layer outputs. (Cai et al., 2020) introduces a zero-shot quanti-139 zation approach that eliminates the need for real calibration data by generating synthetic data. This 140 synthetic data is used to calibrate the quantization parameters, allowing for effective quantization 141 without access to real data. (Qin et al., 2023) addresses the homogenization of synthetic data for 142 quantization and improves quantized model performance.

143 144

145

## 2.3 DATA-DEPENDENT QUANTIZATION

146 Data-dependent quantizations are the techniques that leverage the characteristics and distribution of 147 the training data to enhance the efficiency and accuracy of quantized neural networks. Different fixed precision data-dependent QAT schemes have been proposed in the literature. (Gupta et al., 2015) 148 use 16 bits word length for weights, biases, and other intermediate variables during training with 149 stochastic rounding. (Wang et al., 2019) mine channel-wise interaction to learn interacted bit count 150 to prevent performance degradation in binary neural networks. Unlike traditional methods (Zhuang 151 et al., 2021) first optimized the network with quantized weights and then with quantized activations. 152 To address the increased training time and computational costs of data-driven QAT methods, PTQ 153 methods using real data have been proposed. (Banner et al., 2018) introduce Analytical Clipping for 154 Integer Quantization (ACIQ) to optimize activation clipping, a per-channel bit allocation policy to 155 minimize mean-squared-error, and a bias-correction method to address inherent biases in quantized 156 weights. This technique enhances the accuracy and efficiency of 4-bit quantized convolutional neural 157 networks. A low-bit precision linear quantization framework is proposed by (Choukroun et al., 158 2019) in which the optimal solution to the quantization problem is found via mean-squared-error 159 optimization with fixed precision constraints. (Zhao et al., 2019) sampled activation distribution using 512 training images to determine quantization grid points, then use this grid during testing. 160 Generally, when quantizing the DNNs, assigning each floating-point weight to its nearest fixed-161 point value is the common approach but (Nagel et al., 2020) show that this isn't the best approach. In AdaRound they use a better weight-rounding mechanism for PTQ that adapts to the data and the task loss. Even though these methods provide good quantization techniques, they still need access to training or validation data.

165

166 167 168

#### 2.4 DATA-INDEPENDENT QUANTIZATION

170 Due to privacy or proprietary constraints, access to training or validation datasets is not always possible. Consequently, several works have been proposed to perform quantization using synthetic 171 data to determine activation ranges. DFQ by (Nagel et al., 2019) leverages the scale-equivariance 172 property of activation functions to equalize weight ranges within the network and also corrects 173 quantization-induced biases in the error. DFQ works fine on 8-bit quantization, but the accuracy 174 drops significantly for lower bits. (Xu et al., 2020) use a generator to generate the synthetic data and 175 fine-tune the quantized model. Unlike other generative data-free quantization approaches, (Qian 176 et al., 2023) determine if the knowledge carried out by generated samples is informative or not to 177 the quantized model. They generate samples with adaptive ability to boost the quantized model 178 performance. (Choi et al., 2020) propose an adversarial knowledge distillation framework that min-179 imizes the possible loss for a worst case via adversarial learning by constraining a generator in the 180 adversarial learning framework. The Inception scheme is proposed by (Haroush et al., 2020) to gen-181 erate the data; they proposed knowledge distillation-based methods to improve accuracy. Methods like (Haroush et al., 2020; Cai et al., 2020; Qin et al., 2023) purely use BN statistics of the floating 182 point model to create synthetic data. (Qin et al., 2023) show that the synthetic data generated by 183 constraining BN statistics suffers serious homogenization at both the distribution level and sample 184 level and improves it by slacking distribution alignment and layerwise sample enhancement. While 185 the aforementioned methods employ various techniques to generate synthetic data for quantization and enhance model accuracy, they typically require fine-tuning or rely on BN statistics from the 187 full-precision model, which we demonstrate is not necessary for achieving good performance. 188

189

190 191

## 3 MOTIVATION

192 193

194 BN is a commonly utilized technique during the training of neural networks to stabilize and ac-195 celerate the convergence process by normalizing activations within each layer. However, when it 196 comes to generating synthetic data for PTO, the primary objective shifts to synthesizing inputs that produce a wide and representative range of activations within the network. This objective can be 197 effectively met through direct optimization techniques such as backpropagation, which iteratively 198 refines synthetic inputs using the model's gradients. By doing so, we can ensure that these synthetic 199 inputs accurately represent the target classes without the need for BN statistics. Furthermore, Qin 200 et al. (2023) demonstrate that samples generated solely from BN statistics exhibit minimal statistical 201 variation compared to real-world samples, which are inherently diverse. To address this, we intro-202 duce Gaussian noise during the optimization process to enhance the variability and diversity of the 203 generated synthetic data. 204

Additionally, while generating representative datasets for quantization, the prevailing practice has 205 been to select at least one sample from each class to ensure comprehensive coverage of the activation 206 space (Zhang et al., 2023). This approach is based on the assumption that including every class 207 is necessary to capture the full range of activations required for effective quantization. However, 208 this method can be computationally intensive and time-consuming, especially for large datasets 209 with numerous classes. (Banner et al., 2018) show that the statistical distributions of weights and 210 activations follow a bell curve, indicating that large values occur very rarely compared to small 211 values. Motivated by this for a more efficient process, we investigate whether a smaller subset of 212 classes could provide an effective activation range. By focusing on a reduced number of target 213 classes, our goal was to streamline the calibration data generation process while maintaining the efficacy of quantization. This exploration was driven by the hypothesis that a subset of classes could 214 still encompass the diverse activation patterns needed for accurate model quantization as described 215 in subsection 5.4.

216	Algorithm 1 Generate Calibration Dataset
217	<b>Require:</b> Full precision model $F$ , total number of classes $N$ , number of target classes $M$ , iterations
218	$T$ , learning rate $\alpha$ , threshold $\epsilon$
219	<b>Ensure:</b> Synthetic sample x for each selected class
220	1: Randomly select $M$ unique target classes from N total classes
221	2: Initialize an empty array $\mathbf{X}$ to store synthetic samples
222	3: Define loss function $\mathcal{L}$
223	4: for each target class $c$ in $M$ do
224	5: Initialize $\mathbf{x} \sim \mathcal{N}(0, 1)$
225	6: Set target label $y = c$
226	7: <b>for</b> $t = 1$ to $T$ <b>do</b>
227	8: $\mathbf{x} \leftarrow \mathbf{x} - \alpha \nabla_{\mathbf{x}} \mathcal{L}(F(\mathbf{x}), y)$
228	9: $\mathbf{x} \leftarrow \mathbf{x} + \mathcal{N}(0, 0.1)$
220	10: <b>if</b> $F(\mathbf{x}) = y$ and $\mathcal{L}(F(\mathbf{x}), y) < \epsilon$ then
229	11: break
230	12: end if
231	13: end for
232	14: Append $\mathbf{x}$ to $\mathbf{X}$
233	15: end for
234	16: return X
225	

#### 4 Method

237 238

254

260 261

262

263 264

265

269

239 The algorithm 1 for generating synthetic images begins by randomly selecting M unique target classes from a total of N available classes, ensuring diversity in the synthetic dataset. A loss func-240 tion,  $\mathcal{L}$ , is defined to measure the discrepancy between the model's prediction and the target label. 241 For each selected target class, the process initiates by creating a synthetic image,  $\mathbf{x}$ , with values 242 drawn from a normal distribution,  $\mathcal{N}(0,1)$ . The target label, y, is set to the current class label c. 243 Inspired by the Fast Gradient Sign Method (FGSM) introduced by (Goodfellow et al., 2014), the 244 algorithm then enters an iterative loop for a specified number of iterations, T. In each iteration, the 245 synthetic image x is updated by performing gradient descent, where the gradient of the loss function 246 with respect to x is subtracted, scaled by a learning rate  $\alpha$  as shown in equation 1. This step is anal-247 ogous to the FGSM approach, which perturbs the input in the direction of the gradient to maximize 248 the loss. To introduce variability in the samples, Gaussian noise,  $\mathcal{N}(0,0.1)$ , is added to the image. 249 The algorithm checks if the model's prediction for x matches the target label, and if the loss falls 250 below a predefined threshold,  $\epsilon$ . If both conditions are satisfied, the loop breaks early, indicating that a satisfactory synthetic image has been generated. This process is repeated for each of the M target 251 classes, and the resulting synthetic images are stored in an array X. Finally, the array of synthetic 252 images is returned as the output. 253

$$\mathbf{x} = \mathbf{x} - \alpha \nabla_{\mathbf{x}} \mathcal{L}(F(\mathbf{x}), y) \tag{1}$$

## 5 EXPERIMENTS

In this section, we perform an in-depth assessment of the performance of our approach on image classification tasks using a series of comprehensive experiments.

5.1 IMPLEMENTATION DETAILS

We used PyTorch (Paszke et al., 2019) to implement the proposed approach due to its robust automatic differentiation capabilities. All the experiments are performed on NVIDIA A100 GPU, and the pre-trained full precision models are obtained from pytorchcv<sup>1</sup>. We generate the synthetic data

<sup>&</sup>lt;sup>1</sup>pytorchv:https://pypi.org/project/pytorchcv/

271	Table 1: MobileNetV2/V3 on ImageNet						Table 2: ResNet20 and VGG16bn on C					
272	Model	Method	W-bit	A-bit	Top-1		FARIU					
273		Baseline	32	32	72.97%		Model	Method	W-bit	A-bit	Top-1	
274		SelectQ	4	4	10.88%			Baseline	32	32	94.08%	
275		Ours	4	4	11.98%			Real Data	4	4	87.38%	
276	MobileNetV2	SelectQ	6	6	70.25%		ZeroQ	4	4	85.39%		
277		Ours	6	6	70.12%		ResNet20	DSG	4	4	87.79%	
278		SelectQ	8	8	72.84%			Ours	4	4	89.26%	
270		Ours	8	8	72.87%			Real Data	6	6	93.80%	
279		Baseline	32	32	75.34%		10031 (0120	ZeroQ	6	6	93.33%	
280		SelectQ	4	4	0.36%			DSG	6	6	93.55%	
281		Ours	4	4	19.29%			Ours	6	6	93.64%	
282	MobileNetV3	SelectQ	6	6	60.04%			Real Data	8	8	93.95%	
283	WIODIICI VCI V J	Ours	6	6	<b>72.84</b> %			ZeroQ	8	8	93.94%	
200		SelectQ	8	8	75.04%		DSG	8	8	93.97%		
204		Ours	8	8	75.05%			Ours	8	8	94.00%	
285								Baseline	32	32	93.86%	
286								Real Data	4	4	92.50%	
287								ZeroQ	4	4	91.79%	
288								DSG	4	4	92.89%	
200								Ours	4	4	93.17%	
209							VGG16bn	Real Data	6	6	93.48%	
290								ZeroQ	6	6	93.45%	
291							DSG	6	6	93.68%		
292							Ours	6	6	93.85%		
293							Real Data	8	8	93.59%		
204								ZeroQ	8	8	93.53%	
234								DSG	8	8	93.61%	
295								Ours	8	8	93.79%	
296												

using algorithm 1 and use the independent calibration process from (Cai et al., 2020). In our experiments, we apply quantization to all layers and clip the activations on a per-layer basis. We utilize a Gaussian distribution to initialize the synthetic image as well as to add noise to it. For the hyperparameters, we empirically find the optimum value of total number of target labels M for each model and bit setting but never go beyond 35, iterations T to 100, learning rate  $\alpha$  in the range [0.1, 0.2], and threshold  $\epsilon$  to 0.001.

#### 5.2 EVALUATION

307 To demonstrate the effectiveness of our approach, we evaluate it on various network architectures with different bit settings. Our experiments include AlexNet (Krizhevsky et al., 2012), 308 VGG16bn (Simonyan & Zisserman, 2014), ResNet18/20/50 (He et al., 2016), SqueezeNext (Gho-309 lami et al., 2018), InceptionV3 (Szegedy et al., 2016), ShuffleNet (Zhang et al., 2018), and Mo-310 bileNetV2/V3 (Sandler et al., 2018; Howard et al., 2019). We assess these models using various 311 bit-width configurations, such as W4A4 (4-bit weights and 4-bit activations), W6A6, and W8A8. 312 We use validation datasets from ImageNet (Deng et al., 2009) and CIFAR10 (Krizhevsky et al., 313 2009) to evaluate our approach, measuring the effectiveness by assessing the top-1 accuracy of the 314 quantized models. 315

316 317

270

297 298

299

300

301

302

303

304 305

306

#### 5.3 COMPARISON WITH SOTA METHODS

318 To evaluate the benefits of our proposed PTQ scheme, we compare our method with other data-free 319 PTQ approaches, such as DSG (Qin et al., 2023), ZeroQ (Cai et al., 2020), DFQ (Nagel et al., 2019), 320 ACIQ (Banner et al., 2018), MSE (Chen et al., 2015), KL (Sung et al., 2015), and OCS (Zhao et al., 321 2019), on CIFAR10 and ImageNet datasets. Notably, DSG and ZeroQ are representative generative data-free PTQ methods that reconstruct synthetic data and calibrate the quantized network. We 322 assess these methods under various bit-width configurations, with the results presented in table 2 323 for the CIFAR10 dataset and table 1, 3, and 4 for the ImageNet dataset. For MobilNetV2/V3 on ImageNet dataset we compare our method with SelectQ<sup>2</sup> (Zhang et al., 2023) which uses training data for quantization.

On the CIFAR10 dataset, we evaluate our method with ResNet20 and VGG16bn, as shown in ta-327 ble 2. Our method consistently outperforms other methods across all bit-widths. Specifically, for 328 ResNet20, our method improved accuracy by approximately 1.47% over DSG in the 4-bit setting. 329 For the ImageNet dataset, we conducted experiments on MobileNetV2, MobileNetV3, ResNet18/50, 330 SqueezeNext, InceptionV3, ShuffleNet, and AlexNet models. Our method consistently outperforms 331 other quantization methods across different bit-widths. Notably, for MobileNetV3, our method 332 achieved a significant improvement of 18.93% over SelectQ in the 4-bit setting, reaching 19.29%. 333 In the case of ResNet18, our method achieved the highest accuracy in the 4-bit setting, with a 3.42% 334 improvement over DSG, reaching 43.32%. For ResNet50, our method improved accuracy by 1.50% over DSG in the 4-bit setting, achieving 57.62%. InceptionV3 showed a significant improvement 335 with our method, achieving 60.31% in the 4-bit setting, which is 3.14% higher than DSG, while 336 slightly improving the accuracies in 6 and 8-bit settings. In table 3 we also demonstrate the quanti-337 zation results of AlexNet model using our method which is not shown by other methods due to the 338 unavailability of BN layer in AlexNet. 339

Furthermore similar to methods such as ZeroQ and DSG we require a small number of synthetic samples to achieve effective PTQ. Empirically, we have determined the effective sample size for each model and bit setting. For example, we use 25 samples for ResNet18 in a 4-bit setting and 35 samples for ResNet50 in a 4-bit setting. Importantly, we never exceed 35 samples for effective quantization across all models and bit settings. This substantial reduction in sample size does not compromise quantization performance, making our method highly efficient for various classification models.

Overall, our quantization method demonstrated superior performance across various models and bit widths with fewer samples, particularly in the 4-bit setting. It consistently achieved high accuracy,
 making it suitable for resource-constrained environments without significant loss in performance.



364

366

367 368

369

362

350 351 352

353

354

355

356

357

359

360 361

Figure 1: Quantized model performance with different calibration set size

5.4 EFFECT OF SAMPLE SIZE ON QUANTIZATION

To evaluate the necessity of having a large sample size for calibration datasets in model quantization, we conducted two experiments using the ResNet18, InceptionV3, and ShuffleNet models.

In the first experiment, we tested different unique sample sizes ranging from 15 to 1000 and measured the top-1 accuracy of the quantized models using our synthetic dataset generation method and reported in figure 1. Contrary to the common belief that a larger calibration dataset size leads to better quantized model performance, our findings reveal that the best results are achieved with a smaller sample size. For ResNet18, the highest accuracy of 43.32% was obtained with 25 samples, while for InceptionV3, the highest accuracy of 60.32% was also achieved with 25 samples. Interestingly, increasing the sample size beyond this point resulted in a decline in accuracy. For instance,

<sup>&</sup>lt;sup>2</sup>Importantly note that SelectQ requires access to the training data

					Model	Method	W-bit	A-bit	To
Model	Method	W-bit	A-bit	Top-1		Baseline	32	32	71.
-	Baseline	32	32	69.38%		Real Data	4	4	65.
	Real Data	6	6	66.51%		DFQ	4	4	0.
	ZeroQ	6	6	39.83%		ACIQ	4	4	7.
	DSG	6	6	66.23%		MSE	4	4	15
SqueezeNext	Ours	6	6	66.30%		KL	4	4	16
	Real Data	8	8	69.23%		ZeroQ	4	4	26
	ZeroQ	8	8	68.01%		DSG	4	4	39
	DSG	8	8	69.27%		Ours	4	4	43
	Ours	8	8	69.21%		Real Data	6	6	71
	Baseline	32	32	78.80%		ACIQ	6	6	61
	Real Data	4	4	73.50%	PacNat18	KL	6	6	61
	ZeroQ	4	4	12.00%	ixesiver10	MSE	6	6	66
	DSG	4	4	57.17%		DFQ	6	6	67
	Ours	4	4	60.31%		ZeroQ	6	6	69
-	Real Data	6	6	78.59%		DSG	6	6	70
	ZeroQ	6	6	75.14%		Ours	6	6	70
InceptionV3	DSG	6	6	78.12%		Real Data	8	8	71
1	Ours	6	6	78.35%		ACIQ	8	8	68
	Real Data	8	8	78.79%		DFQ	8	8	69
	ZeroO	8	8	78.70%		KL	8	8	70
	DSG	8	8	78.81%		MSE	8	8	71
	Ours	8	8	78.83%		ZeroQ	8	8	71
	Baseline	32	32	65.07%		DSG	8	8	71
	Real Data	6	6	56.25%		Ours	8	8	71
	ZeroO	6	6	39.92%		Baseline	32	32	77
	DSG	6	6	60.71%		Real Data	4	4	68
C1 C2 NT -	Ours	6	6	60.60%		ACIQ	4	4	61
ShuffleNet	Real Data	8	8	64.52%		ZeroQ	4	4	8.
	ZeroO	8	8	64.46%		DFQ	4	4	10
	DSG	8	8	64.87%		DSG	4	4	56
	Ours	8	8	64.93%		Ours	4	4	57.
	Baseline	32	32	59.04%		Real Data	6	6	76
	Ours	4	4	45.97%	DecNet50	ZeroQ	6	6	75
AlexNet	Ours	6	6	57.22%	Resilet50	DSG	6	6	76
	Ours	8	8	57.47%		Ours	6	6	77.
	1	-	-			Real Data	8	8	77
						ZeroQ	8	8	77
						DSG	8	8	77.
						Ours	8	8	77

Table 3: SqueezeNext, InceptionV3, and 379 ShuffleNet on ImageNet

#### Table 4: ResNet18/50 on ImageNet

40 40 41

378

- 41
- 41

414 415

416

417

421

423

the accuracy for ResNet18 dropped to 42.26% with 1000 samples, and for InceptionV3, it decreased 418 to 58.41%. This counterintuitive outcome suggests that using a smaller dataset can be more effective 419 for PTQ. Our results indicate that an optimal sample size exists, beyond which additional samples 420 may introduce noise or redundancy, negatively impacting the quantized model's performance. These findings underscore the efficiency and effectiveness of our method, which achieves good quantiza-422 tion results with significantly fewer samples.

In the second experiment, we plotted the minimum and maximum activation ranges for each con-424 volution layer in the ResNet18 and ShuffleNet models, using datasets created with 25 samples and 425 1000 samples, respectively. The results, illustrated in figures 2 and 3, show that the activation ranges 426 for 25 samples are remarkably close to those obtained with 1000 samples. This observation indi-427 cates that even with a significantly smaller number of samples, the activation ranges remain stable 428 and representative of the model's behavior. Consequently, this suggests that it is not necessary to 429 have one sample from each class to create a representative dataset for post-training quantization. 430

Overall, our findings demonstrate that a smaller, more manageable dataset can be effectively used 431 for calibration, simplifying the process and reducing the need for extensive calibration datasets.

8







## 6 CONCLUSION

In this paper, we have introduced a method for post-training data-free quantization that generates synthetic data using the trained full-precision model independent of BN statistics. This makes our approach versatile and applicable to any model architecture. Our experimental results demonstrate that it is not necessary to include samples from each target category in the calibration dataset; select-ing only a few target classes is sufficient to create an effective calibration dataset. Our method con-sistently outperforms existing generative data-free quantization methods with calibration, as demon-strated through extensive comparisons across various standard network architectures. These include ResNet18/50, SqueezeNext, InceptionV3, as well as lightweight architectures like ShuffleNet and MobileNetV2/V3. Notably, our approach shows significant improvements in 4-bit precision set-tings. For instance, on the ResNet18 model, our method increases the top-1 accuracy by over 3.42% compared to SOTA DSG method. These findings underscore the effectiveness and generalizabil-ity of our approach, highlighting its potential to achieve high accuracy with lower bit-widths and fewer calibration samples. This makes it a promising solution for efficient model deployment in resource-constrained environments.

References

Ron Banner, Yury Nahshan, Elad Hoffer, and Daniel Soudry. Aciq: Analytical clipping for integer quantization of neural networks. 2018.

485 Ron Banner, Yury Nahshan, and Daniel Soudry. Post training 4-bit quantization of convolutional networks for rapid-deployment. *Advances in Neural Information Processing Systems*, 32, 2019.

486 Mariusz Bojarski, Davide Del Testa, Daniel Dworakowski, Bernhard Firner, Beat Flepp, Prasoon 487 Goyal, Lawrence D Jackel, Mathew Monfort, Urs Muller, Jiakai Zhang, et al. End to end learning 488 for self-driving cars. arXiv preprint arXiv:1604.07316, 2016. 489 Yaohui Cai, Zhewei Yao, Zhen Dong, Amir Gholami, Michael W Mahoney, and Kurt Keutzer. 490 Zeroq: A novel zero shot quantization framework. In Proceedings of the IEEE/CVF conference 491 on computer vision and pattern recognition, pp. 13169–13178, 2020. 492 493 Tianqi Chen, Mu Li, Yutian Li, Min Lin, Naiyan Wang, Minjie Wang, Tianjun Xiao, Bing Xu, 494 Chiyuan Zhang, and Zheng Zhang. Mxnet: A flexible and efficient machine learning library for 495 heterogeneous distributed systems. arXiv preprint arXiv:1512.01274, 2015. 496 Jungwook Choi, Zhuo Wang, Swagath Venkataramani, Pierce I-Jen Chuang, Vijayalakshmi Srini-497 vasan, and Kailash Gopalakrishnan. Pact: Parameterized clipping activation for quantized neural 498 networks. arXiv preprint arXiv:1805.06085, 2018. 499 500 Yoojin Choi, Jihwan Choi, Mostafa El-Khamy, and Jungwon Lee. Data-free network quantization 501 with adversarial knowledge distillation. In Proceedings of the IEEE/CVF Conference on Com-502 puter Vision and Pattern Recognition Workshops, pp. 710–711, 2020. Yoni Choukroun, Eli Kravchik, Fan Yang, and Pavel Kisilev. Low-bit quantization of neural net-504 works for efficient inference. In 2019 IEEE/CVF International Conference on Computer Vision 505 Workshop (ICCVW), pp. 3009-3018. IEEE, 2019. 506 507 Matthieu Courbariaux, Itay Hubara, Daniel Soudry, Ran El-Yaniv, and Yoshua Bengio. Binarized 508 neural networks: Training deep neural networks with weights and activations constrained to+ 1 509 or-1. arXiv preprint arXiv:1602.02830, 2016. 510 Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hi-511 erarchical image database. In 2009 IEEE conference on computer vision and pattern recognition, 512 pp. 248–255. Ieee, 2009. 513 514 Andre Esteva, Brett Kuprel, Roberto A Novoa, Justin Ko, Susan M Swetter, Helen M Blau, and 515 Sebastian Thrun. Dermatologist-level classification of skin cancer with deep neural networks. 516 nature, 542(7639):115-118, 2017. 517 Amir Gholami, Kiseok Kwon, Bichen Wu, Zizheng Tai, Xiangyu Yue, Peter Jin, Sicheng Zhao, and 518 Kurt Keutzer. Squeezenext: Hardware-aware neural network design. In Proceedings of the IEEE 519 conference on computer vision and pattern recognition workshops, pp. 1638–1647, 2018. 520 521 Ross Girshick, Jeff Donahue, Trevor Darrell, and Jitendra Malik. Rich feature hierarchies for ac-522 curate object detection and semantic segmentation. In Proceedings of the IEEE conference on 523 computer vision and pattern recognition, pp. 580-587, 2014. 524 Ian J Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing adversarial 525 examples. arXiv preprint arXiv:1412.6572, 2014. 526 527 Suyog Gupta, Ankur Agrawal, Kailash Gopalakrishnan, and Pritish Narayanan. Deep learning with 528 limited numerical precision. In International conference on machine learning, pp. 1737–1746. 529 PMLR, 2015. 530 Matan Haroush, Itay Hubara, Elad Hoffer, and Daniel Soudry. The knowledge within: Methods for 531 data-free model compression. In Proceedings of the IEEE/CVF Conference on Computer Vision 532 and Pattern Recognition, pp. 8494-8502, 2020. 533 534 Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recog-535 nition. In Proceedings of the IEEE conference on computer vision and pattern recognition, pp. 536 770–778, 2016. 537 Andrew Howard, Mark Sandler, Grace Chu, Liang-Chieh Chen, Bo Chen, Mingxing Tan, Weijun 538 Wang, Yukun Zhu, Ruoming Pang, Vijay Vasudevan, et al. Searching for mobilenetv3. In Pro-539 ceedings of the IEEE/CVF international conference on computer vision, pp. 1314–1324, 2019.

550

555

567

581

582

583

584

- Andrew G Howard. Mo-bilenets: Efficient convolutional neural networks for mo-bile vision applications. *arXiv preprint arXiv:1704.04861*, 2017.
- Itay Hubara, Yury Nahshan, Yair Hanani, Ron Banner, and Daniel Soudry. Accurate post training
   quantization with small calibration sets. In *International Conference on Machine Learning*, pp. 4466–4475. PMLR, 2021.
- 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*, pp. 2704–2713, 2018.
- Raghuraman Krishnamoorthi. Quantizing deep convolutional networks for efficient inference: A whitepaper. arXiv preprint arXiv:1806.08342, 2018.
- Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images.
   2009.
- 556 Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. *Advances in neural information processing systems*, 25, 2012.
- Yann LeCun, Yoshua Bengio, and Geoffrey Hinton. Deep learning. *nature*, 521(7553):436–444, 2015.
- Yuhang Li, Ruihao Gong, Xu Tan, Yang Yang, Peng Hu, Qi Zhang, Fengwei Yu, Wei Wang, and
   Shi Gu. Brecq: Pushing the limit of post-training quantization by block reconstruction. *arXiv* preprint arXiv:2102.05426, 2021.
- Jonathan Long, Evan Shelhamer, and Trevor Darrell. Fully convolutional networks for semantic segmentation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 3431–3440, 2015.
- Jeffrey L McKinstry, Steven K Esser, Rathinakumar Appuswamy, Deepika Bablani, John V Arthur,
   Izzet B Yildiz, and Dharmendra S Modha. Discovering low-precision networks close to full precision networks for efficient embedded inference. arXiv preprint arXiv:1809.04191, 2018.
- Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A Rusu, Joel Veness, Marc G Belle mare, Alex Graves, Martin Riedmiller, Andreas K Fidjeland, Georg Ostrovski, et al. Human-level
   control through deep reinforcement learning. *nature*, 518(7540):529–533, 2015.
- Markus Nagel, Mart van Baalen, Tijmen Blankevoort, and Max Welling. Data-free quantization through weight equalization and bias correction. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 1325–1334, 2019.
- Markus Nagel, Rana Ali Amjad, Mart Van Baalen, Christos Louizos, and Tijmen Blankevoort. Up or
   down? adaptive rounding for post-training quantization. In *International Conference on Machine Learning*, pp. 7197–7206. PMLR, 2020.
  - Yury Nahshan, Brian Chmiel, Chaim Baskin, Evgenii Zheltonozhskii, Ron Banner, Alex M Bronstein, and Avi Mendelson. Loss aware post-training quantization. *Machine Learning*, 110(11): 3245–3262, 2021.
- Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor
   Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: An imperative style, high performance deep learning library. *Advances in neural information processing systems*, 32, 2019.
- Biao Qian, Yang Wang, Richang Hong, and Meng Wang. Adaptive data-free quantization. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 7960–7968, 2023.
- Haotong Qin, Yifu Ding, Xiangguo Zhang, Jiakai Wang, Xianglong Liu, and Jiwen Lu. Diverse
   sample generation: Pushing the limit of generative data-free quantization. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 45(10):11689–11706, 2023.

- Mohammad Rastegari, Vicente Ordonez, Joseph Redmon, and Ali Farhadi. Xnor-net: Imagenet classification using binary convolutional neural networks. In *European conference on computer vision*, pp. 525–542. Springer, 2016.
- Mark Sandler, Andrew Howard, Menglong Zhu, Andrey Zhmoginov, and Liang-Chieh Chen. Mo bilenetv2: Inverted residuals and linear bottlenecks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 4510–4520, 2018.
- Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image
   recognition. *arXiv preprint arXiv:1409.1556*, 2014.
- Wonyong Sung, Sungho Shin, and Kyuyeon Hwang. Resiliency of deep neural networks under quantization. *arXiv preprint arXiv:1511.06488*, 2015.
- Christian Szegedy, Vincent Vanhoucke, Sergey Ioffe, Jon Shlens, and Zbigniew Wojna. Rethink ing the inception architecture for computer vision. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 2818–2826, 2016.
- Ziwei Wang, Jiwen Lu, Chenxin Tao, Jie Zhou, and Qi Tian. Learning channel-wise interactions for binary convolutional neural networks. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 568–577, 2019.
- Shoukai Xu, Haokun Li, Bohan Zhuang, Jing Liu, Jiezhang Cao, Chuangrun Liang, and Mingkui
  Tan. Generative low-bitwidth data free quantization. In *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XII 16*, pp. 1–17.
  Springer, 2020.
- Kiangyu Zhang, Xinyu Zhou, Mengxiao Lin, and Jian Sun. Shufflenet: An extremely efficient convolutional neural network for mobile devices. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 6848–6856, 2018.
- <sup>620</sup>
  <sup>621</sup> Zhao Zhang, Yangcheng Gao, Jicong Fan, Zhongqiu Zhao, Yi Yang, and Shuicheng Yan. Selectq:
  <sup>622</sup> Calibration data selection for post-training quantization. *Authorea Preprints*, 2023.
- Ritchie Zhao, Yuwei Hu, Jordan Dotzel, Chris De Sa, and Zhiru Zhang. Improving neural network
   quantization without retraining using outlier channel splitting. In *International conference on machine learning*, pp. 7543–7552. PMLR, 2019.
- Shuchang Zhou, Yuxin Wu, Zekun Ni, Xinyu Zhou, He Wen, and Yuheng Zou. Dorefa-net: Training low bitwidth convolutional neural networks with low bitwidth gradients. *arXiv preprint arXiv:1606.06160*, 2016.
- 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*, 44(10):6140–6152, 2021.
- 634 635

633

- 636 637
- 638
- 639 640

641

- 642
- 643
- 644

645

646 647