

A. Appendix

A.1. Implementation Details

For all experiments, unless stated otherwise, we use the following hyperparameters: for classification tasks, 1024 images are generated, and for object detection, 32 images are used. Each result represents an average of five experiments. We perform 1000 optimization iterations, using the RAdam optimizer [36] with an initial learning rate of 16 and a reduced learning rate on the plateau scheduler.

For ϕ_{prep} , we begin by applying a 3×3 Gaussian smoothing filter. Then, we randomly apply horizontal flipping to the images with a probability of 0.5. Next, we perform random cropping by starting with images that are 32 pixels larger in both height and width than the final output size and then cropping them to the desired output shape. The final images are center-cropped.

A.2. Derivation of Equation (5)

Here, we provide the derivation of Eq. (5). We begin with the data mean decomposition, which is followed by the decomposition of variance.

$$\begin{aligned}
 [\bar{\mu}_l(\mathcal{D})]_i &\triangleq \frac{1}{M} \sum_{m=1}^M \frac{1}{d_l} \sum_{j=1}^{d_l} [\mathbf{y}_l^{(m)}]_{i,j} \\
 &= \frac{1}{N} \sum_{n=1}^N \frac{1}{K} \sum_{k \in \mathcal{S}_n} \frac{1}{d_l} \sum_{j=1}^{d_l} [\mathbf{y}_l^{(k)}]_{i,j} \\
 &= \frac{1}{N} \sum_{n=1}^N [\bar{\mu}_l(\mathcal{B}_n)]_i, \tag{9}
 \end{aligned}$$

where \mathcal{S}_n is the set of index correspond to the n^{th} batch. The first step follows directly from the definition in (3a). In the second step, we apply the linearity of summation. Finally, we use the notion of an empirical mean over batch. Next, we present the relationship between the variance and the second moment:

$$\begin{aligned}
 [\bar{\sigma}_l(\mathcal{D})]_i &\triangleq \frac{1}{M} \sum_{m=1}^M \frac{1}{d_l} \sum_{j=1}^{d_l} \left([\mathbf{y}_l^{(m)}]_{i,j} - [\bar{\mu}_l(\mathcal{D})]_i \right)^2 \\
 &= [\bar{\mu}_l(\mathcal{D})]_i^2 + \frac{1}{M} \sum_{m=1}^M \frac{1}{d_l} \sum_{j=1}^{d_l} [\mathbf{y}_l^{(m)}]_{i,j}^2 - 2 [\bar{\mu}_l(\mathcal{D})]_i \frac{1}{M} \sum_{m=1}^M \frac{1}{d_l} \sum_{j=1}^{d_l} [\mathbf{y}_l^{(m)}]_{i,j} \\
 &= \frac{1}{N} \sum_{n=1}^N \frac{1}{K} \sum_{k \in \mathcal{S}_n} \frac{1}{d_l} \sum_{j=1}^{d_l} [\mathbf{y}_l^{(k)}]_{i,j}^2 - [\bar{\mu}_l(\mathcal{D})]_i^2 \\
 &= \frac{1}{N} \sum_{n=1}^N \left[\bar{M}_\ell^{(2)}(\mathcal{B}_n) \right]_i - [\bar{\mu}_l(\mathcal{D})]_i^2, \tag{10}
 \end{aligned}$$

The first step follows from the definition in (3b). In the second step, we expand the parentheses. In the third step,

Table 5. Top-1 accuracy on ImageNet-1k validation set with models quantized using BRECQ [30] in an academic quantization setting. The models are quantized using four different data generation algorithms and real data

Method	ResNet-18	ResNet-50	MBV2
	71.06	77.0	72.49
W4A4 Academic Quantized			
ZeroQ [4]	68.67 \pm 0.09	74.26 \pm 0.12	69.27 \pm 0.22
IntraQ [60]	67.78 \pm 0.27	66.08 \pm 0.67	69.06 \pm 0.13
GENIE [25]	69.67 \pm 0.05	74.90 \pm 0.06	69.13 \pm 0.05
DGH (Ours)	69.56 \pm 0.04	74.72 \pm 0.12	69.23 \pm 0.08
Real Data	69.74 \pm 0.07	74.90 \pm 0.06	69.30 \pm 0.10
W2A4 Academic Quantized			
ZeroQ [4]	59.73 \pm 0.17	63.53 \pm 0.22	27.42 \pm 1.04
IntraQ [60]	49.17 \pm 0.56	44.03 \pm 1.76	22.8 \pm 3.17
GENIE [25]	64.37 \pm 0.17	69.53 \pm 0.04	48.23 \pm 3.24
DGH (Ours)	64.15 \pm 0.10	68.91 \pm 0.06	45.83 \pm 1.44
Real Data	65.86 \pm 0.10	70.28 \pm 0.03	54.46 \pm 1.44

we apply the linearity of summation and the notion of empirical mean. Finally, we use the notion of an empirical moment over the batch.

A.3. Additional Results

A.3.1 Academic Quantization Results

We validate our approach within the academic quantization scheme. In this setup, the first and last layers and the input to the second layer are quantized to 8 bits, while the output layer remains unquantized. Tab. 5 shows the Top-1 accuracy on the ImageNet-1k validation set using the BRECQ quantization algorithm under the academic quantization scheme, with activation bit-width set to 4 bits and weight bit-width set to 2 and 4 bits. The results demonstrate that the DGH achieves competitive results to GENIE [25], while greatly improving the image generation runtime since DGH does not use a generator. Specifically, GENIE [25] requires approximately two and a half hours on a V100 GPU to generate 1024 images for ResNet18, while DGH takes less than half an hour on an RTX 3090. Although ZeroQ [4] and IntraQ [60] may offer faster generation times than DGH, they deliver inferior performance in both academic and hardware-friendly quantization schemes.

A.3.2 The Effect of Varying Bit Widths

Here, we provide results offering further insights into the performance of DGH. We present the Top-1 accuracy on ImageNet-1k validation set using the BRECQ quantization algorithm in the hardware-friendly quantization setting, tested across various bit-widths. In this experiment, we keep the weight bit-width fixed at 8 bits and evaluate

IntraQ [60], GENIE [25], and DGH at different activation bit-widths. The results are shown in Fig. 6. Next, we fix the activation bit-width at 8 bits and evaluate IntraQ [60], GENIE [25], and DGH at varying weight bit-widths, with the results presented in Fig. 7. From both sets of results, we observe that the performance gap between DGH and IntraQ [60], GENIE [25] increases as the bit-width decreases. Note that we omit ZeroQ [4] as its accuracy is consistently near zero in all cases.

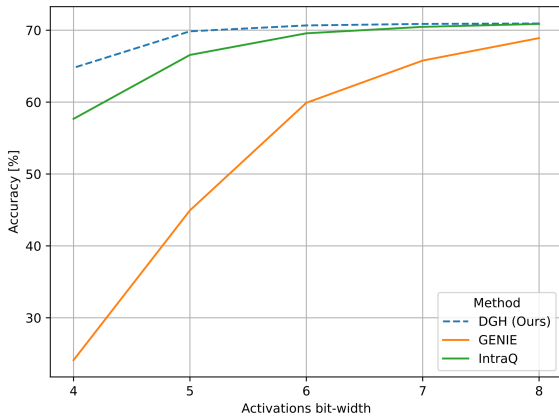


Figure 6. Top-1 accuracy on the ImageNet-1k validation set using the BRECQ quantization algorithm in a hardware-friendly quantization setting, with various activation bit-widths while the weight bit-width is fixed at 8 bits. The y-axis represents the accuracy, and the x-axis represents the activation bit-width.

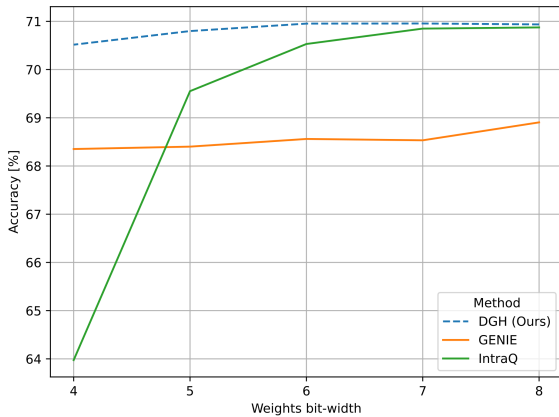


Figure 7. Top-1 accuracy on the ImageNet-1k validation set using the BRECQ quantization algorithm in a hardware-friendly quantization setting, with various weight bit-widths while the activation bit-width is fixed at 8 bits. The y-axis represents the accuracy, and the x-axis represents the weights bit-width.

A.3.3 Generated Images Using DGH

Next, we visualized the images generated by DGH from ResNet18 after 1k iteration, as shown in (Fig. 8), which were used in all experiments. Additionally, we present images generated after 40k iterations in (Fig. 9). In Fig. 8, we observe the initial emergence of shapes and patterns in the generated images. As shown in Fig. 9, with additional iterations, these patterns and class-specific features become increasingly distinct, reflecting improved image quality and clearer representation of different classes in the generated images.



Figure 8. Images generated from a ResNet18 using DGH with 1k Iterations.



Figure 9. Images generated from a ResNet18 using DGH with 40k Iterations.