# Quantization Meets OOD: Generalizable Quantization-aware Training from a Flatness Perspective

**Jiacheng Jiang**
SIGS, Tsinghua University
Shenzhen, China
jiangjc23@mails.tsinghua.edu.cn

**Yuan Meng**[†]
Key Laboratory of Pervasive
Computing, Ministry of Education
Department of Computer Science and
Technology, Tsinghua University
Beijing, China
yuanmeng@mail.tsinghua.edu.cn

**Chen Tang**
MMLab, The Chinese University of
Hong Kong
Hong Kong, China

**Han Yu**
Department of Computer Science and
Technology, Tsinghua University
Beijing, China

**Qun Li**
SIGS, Tsinghua University
Shenzhen, China

**Zhi Wang**[†]
SIGS, Tsinghua University
Shenzhen, China
wang_zhi@tsinghua.edu.cn

**Wenwu Zhu**[†]
Key Laboratory of Pervasive
Computing, Ministry of Education
Department of Computer Science and
Technology, Tsinghua University
Beijing, China
wwzhu@tsinghua.edu.cn

## Abstract

Current quantization-aware training (QAT) methods primarily focus on enhancing the performance of quantized models on in-distribution (I.D) data, while overlooking the potential performance degradation on out-of-distribution (OOD) data. In this paper, we first substantiate this problem through rigorous experiment, showing that *QAT can lead to a significant OOD generalization performance degradation*. Further, we find the contradiction between the perspective that flatness of loss landscape gives rise to superior OOD generalization and the phenomenon that QAT lead to a sharp loss landscape, can cause the above problem. Therefore, we propose a flatness-oriented QAT method, FQAT, to achieve generalizable QAT. Specifically, i) FQAT introduces a layer-wise freezing mechanism to mitigate the gradient conflict issue between dual optimization objectives (i.e., vanilla QAT and flatness). ii) FQAT proposes an disorder-guided adaptive freezing algorithm to dynamically determines which layers to freeze at each training step, effectively addressing the challenges caused by interference between layers. A gradient disorder metric is designed to help the algorithm identify unstable layers during training. Extensive experiments on influential OOD benchmark demonstrate the superiority of our method

over state-of-the-art baselines under both I.D and OOD image classification tasks. **Code:** https://github.com/JachinJiang/Quantization-Meets-OOD

## CCS Concepts

• **Computing methodologies → Neural networks**; **Transfer learning**; • **Computer systems organization → Embedded software**.

## Keywords

Quantization; OOD; SAM; Freeze

## 1 Introduction

Quantization is one of the most effective technique for compressing deep neural networks (DNNs) to enable efficient inference and on-device execution while preserving high accuracy of computer vision (CV) tasks [6, 8, 14, 34, 55]. By transforming pre-trained weights and activations from the high-bit (a.k.a., full precision) format, such as FP32, to more compact low-bit formats, such as INT8, quantization significantly reduces power consumption and speeds up inference, making it ideal for deploying neural networks on resource-limited edge devices [8, 18, 38]. Quantization-aware
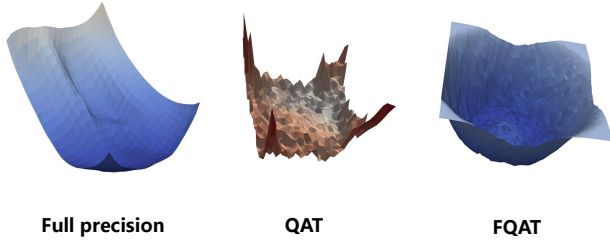
---

[†]Corresponding author.

**Figure 1: Visualization of loss landscapes of full precision model, quantized model (W4A4 with existing QAT method LSQ), quantized model (W4A4 with FQAT) on Sketch domain of PACS.**

training (QAT) simulates quantization during training, allowing the model to adapt to quantization effects and achieve better performace than post-training quantzation [14, 34].

Despite the notable success of QAT methods, the existing works all focus on retaining the performance of full-precision models on in-distribution (I.D) data [8, 21, 27, 39–41, 46, 55], which may result in the performance degradation on out-of-distribution (OOD) data. However, OOD data is common in edge applications, such as foggy pedestrian detection in autonomous driving [15, 22]. Unfortunately, to the best of our knowledge, no prior research has systematically explored this potential risk. In this work, we take the first step by addressing a critical question:

*Does quantization degrade the OOD generalization performance of a full-precision CV model?*

We perform systematic analysis on a representative class of QAT methods that with learned scale factors [8, 19, 21] and arrived at the following conclusion: **existing QAT methods lead to a decline in the model's OOD generalization performance** Furthermore, we identify a key inconsistency with a prior study [17], which we trace back to data leakage in their experimental setup.

Drawing inspiration from generalization techniques originally proposed for full-precision models—particularly sharpness-aware minimization (SAM) [9], we note that flatter loss landscapes are associated with improved OOD generalization [4, 45, 50]. This indicates, however, a sharp loss landscape induced by current QAT methods [27], which likely accounts for the observed performance degradation. Fig. 1 visualizes the loss landscape of the full-precision model and the quantified model. The loss landscape of the quantified model is sharp compared to the full-precision model.

Based on the aforementioned findings, we propose a flatness-aware quantization-aware training method, FQAT, which incorporated a flatness-oriented optimization objective into the vanilla quantization-aware training objective to achieve generalizable quantized DNNs. However, the optimization of FQAT is non-trivial due to the following challenges. First, it is challenging to optimize the key quantization parameter, *scale factor*, due to its involvement in two optimization functions (i.e., vanilla QAT and flatness) with conflicting gradients. We find that the flatness-oriented gradient exhibits greater instability. Therefore, we propose a layer-wise QAT-gradient freezing method to mitigate the instability introduced by

quantization. Second, it is challenging to formulate the freezing policy. The existing freezing methods [29, 36, 40] fail to consider the interference between scale factors of different layers under two optimization objectives, leading to poor generalization performances. To tackle these challenges, we design a gradient disorder metric to identify more stable layers during training. Based on this metric, we propose an adaptive freezing algorithm that leverages historical gradient information to automatically select layers for freezing quantization gradients in the future, thereby ensuring stable training for the flatness objective. The last figure of Fig. 1 also visualizes how our approach can effectively improve flatness of loss landscape.

In summary, we have made the following contributions:

- To the best of our knowledge, this is the first work to propose and substantiate that quantization damages the OOD performance of full-precision CV models.
- We propose FQAT, a generalizable QAT method designed from the perspective of flatness. FQAT introduces a layer-wise freezing mechanism to mitigate gradient conflicts between standard QAT and flatness objectives. Additionally, we develop a disorder-guided adaptive freezing algorithm, which dynamically adjusts the layers to freeze at each training step based on historical gradient disorder.
- Extensive experiments on PACS, OfficeHome and Domain-Net validate the effectiveness of FQAT. Compared to the baselines, FQAT outperforms across all experimental settings and even surpasses the full-precision model in 8-bit quantization. FQAT improves both I.D and OOD performance, with greater gains in OOD (e.g., Under the 3-bit setting of PACS, OOD gain (+5.24%) exceeds I.D gain (+2.61%)). This highlights effectiveness of FQAT in enhancing OOD generalization.

## 2 Preliminaries

In this paper, we adopt the symmetric uniform quantization function for both weight and activation: $\hat{v} = Q(v; s) = s \times \left\lfloor \text{clip}\left(\frac{v}{s}, l, u\right) \right\rceil$, where $v$ and $\hat{v}$ represent the full-precision value and the quantized value, respectively. The operator $\lfloor \cdot \rceil$ denotes rounding to the nearest integer, and $s$ is a learnable scale factor in QAT [8, 41]. The clip function ensures values stay within the bounds $[l, u]$. In $b$-bit quantization, for activation quantization, we set $l = 0$ and $u = 2^b - 1$; for weight quantization, we set $l = -2^{b-1}$ and $u = 2^{b-1} - 1$. Furthermore, to overcome the non-differentiability of the rounding operation, the Straight-Through Estimator (STE) [2] is employed to approximate the gradients: $\frac{\partial \mathcal{L}}{\partial v} \approx \frac{\partial \mathcal{L}}{\partial \hat{v}} \cdot 1_{l \le \frac{v}{s} \le u}$.

## 3 OOD Generalization Assessment of Quantized DNNs

This chapter will address the critical question posed in this paper and offer a detailed analysis.

**Experimental Settings.** In assessment pipeline, the evaluated model is obtained through two steps: 1) obtaining a full-precision OOD generalizable model, 2) performing quantization-oriented training. Following the recent mainstream OOD model evaluation protocol [4, 45, 48, 50], we conduct step one by fine-tuning the self-supervised pretrained model with an OOD generalization method

**Table 1: LSQ results on three datasets. Experiments were run on seeds 0 and 23, with results reported as mean±std. LSQ: representative QAT method; LSQ+SAGM: Directly introducing SAGM objective to QAT; Ours: proposed FQAT; Val: I.D validation accuracy; Test: OOD test accuracy. The best performance is highlighted in bold.**

| Dataset | Method | Bit-width (W/A) | Val | Test |
|---|---|---|---|---|
| | ERM | Full | 62.04 | 40.95 |
| | LSQ | 4/4 | 60.95±0.38 | 39.14±0.24 |
| | LSQ+SAGM | 4/4 | 61.79±0.26 | 40.07±0.23 |
| | Ours | 4/4 | **62.52±0.25** | **40.60±0.19** |
| DomainNet | LSQ | 5/5 | 60.43±0.37 | 38.46±0.16 |
| | LSQ+SAGM | 5/5 | 62.73±0.42 | 40.51±0.37 |
| | Ours | 5/5 | **63.19±0.36** | **40.93±0.25** |
| | LSQ | 8/8 | 60.97±0.31 | 39.10±0.16 |
| | LSQ+SAGM | 8/8 | 63.21±0.36 | 41.10±0.35 |
| | Ours | 8/8 | **63.39±0.21** | **41.27±0.18** |
| | ERM | Full | 96.42 | 85.29 |
| | LSQ | 3/3 | 77.39±0.68 | 54.35±2.23 |
| | LSQ+SAGM | 3/3 | 74.98±1.49 | 50.60±1.37 |
| | Ours | 3/3 | **77.59±1.21** | **55.84±1.48** |
| PACS | LSQ | 4/4 | 80.44±0.82 | 57.4±1.58 |
| | LSQ+SAGM | 4/4 | 78.98±1.46 | 55.66±1.67 |
| | Ours | 4/4 | **81.26±1.12** | **59.42±1.25** |
| | LSQ | 5/5 | 80.07±1.05 | 57.38±1.33 |
| | LSQ+SAGM | 5/5 | 80.43±2.21 | 58.69±2.53 |
| | Ours | 5/5 | **81.97±1.64** | **60.41±2.43** |
| | ERM | Full | 78.44 | 60.31 |
| OfficeHome | LSQ | 3/3 | 59.58±2.07 | 38.15±1.85 |
| | LSQ+SAGM | 3/3 | 60.69±1.75 | 39.98±1.72 |
| | Ours | 3/3 | **62.12±1.11** | **41.30±1.23** |

**Table 2: EWGS results on two datasets. Experiments were run on seeds 0 and 23, with results reported as mean±std. EWGS: representative QAT method; EWGS+SAGM: Directly introducing SAGM objective to QAT; Ours: proposed FQAT; Val: I.D validation accuracy; Test: OOD test accuracy. The best performance is highlighted in bold.**

| Dataset | Method | Bit-width (W/A) | Val | Test |
|---|---|---|---|---|
| | ERM | Full | 62.04 | 40.95 |
| | EWGS | 4/4 | 60.34±0.78 | 38.74±0.55 |
| | EWGS+SAGM | 4/4 | 61.12±0.30 | 39.98±0.27 |
| | Ours | 4/4 | **61.41±0.50** | **40.23±0.34** |
| DomainNet | EWGS | 5/5 | 61.21±0.16 | 39.02±0.24 |
| | EWGS+SAGM | 5/5 | 62.23±0.47 | 40.44±0.34 |
| | Ours | 5/5 | **62.81±0.46** | **40.88±0.26** |
| | EWGS | 8/8 | 61.04±0.17 | 39.24±0.10 |
| | EWGS+SAGM | 8/8 | 63.02±0.30 | 40.92±0.21 |
| | Ours | 8/8 | **63.11±0.18** | **41.01±0.16** |
| | ERM | Full | 78.44 | 60.31 |
| OfficeHome | EWGS | 3/3 | 60.32±2.41 | 39.09±2.93 |
| | EWGS+SAGM | 3/3 | 60.12±1.21 | 39.12±1.51 |
| | Ours | 3/3 | **62.13±1.40** | **41.99±1.96** |

based on training dataset of the OOD benchmark. The fine-tuning dataset is also used for QAT in step two.

There are three critical choices in the experiments: *pre-trained model, full-precision OOD generalization method, and QAT method*. Previous OOD generalization research often employs ImageNet-pre-trained ResNet50 as the pre-trained model [13, 45, 50]. However, a recent study [48] highlights that ImageNet pre-training introduces severe data leakage issues, thereby compromising the validity of evaluations. Specifically, it becomes unclear whether performance improvements stem from retaining I.D performance or from genuine enhancements in OOD performance. To address this concern, MoCoV2-pre-trained ResNet50 is recommended as the pre-trained model [48]. Furthermore, based on the sound pre-trained model, many OOD generalization methods underperform Empirical Risk Minimization (ERM) [42, 48]. Therefore, we adopt ERM as our full-precision OOD generalization method. Finally, we select two representative QAT methods LSQ [8] and EWGS [21] to handle the quantization process.

**Conclusion:** *Quantization can lead to OOD generalization performance degradation.*

In Table 1 and 2, we present the experimental results on three OOD benchmarks: DomainNet, PACS and OfficeHome. By comparing the experimental results of ERM and QAT methods across
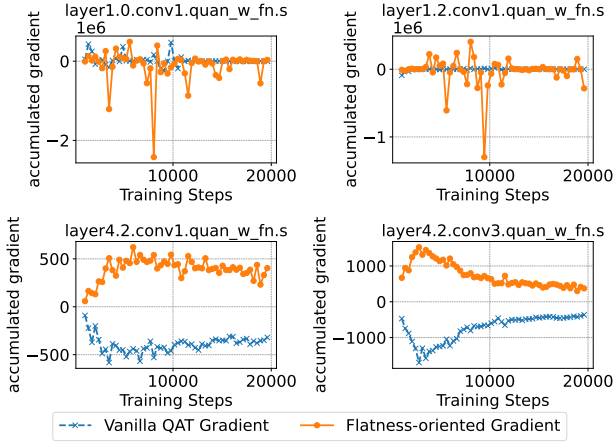
four different bit-widths, we find that quantization leads to significant performance degradation. It is worth mentioning that a related work [17] has reached a conclusion inconsistent with ours. We notice that they use ImageNet-pre-trained ResNet50 [11] as pretrained model, which has been criticized for potential data leakage issues in OOD generalization evaluation [48]. Since the related work has not been open-sourced [17], we conducted ablation experiment within our own workflow by replacing the pre-trained model. The 4-bit LSQ performance of Art domain in PACS improved from 51.07 (MoCoV2 pre-trained) to 76.51 (ImageNet pre-trained, without hyperparameter tuning), similar with the related work [17]. This result validates the soundness of our workflow and clarifies the inconsistency in conclusions.

SAM [4, 45, 50] is an important research direction in the field of OOD generalization. It suggests that enhancing the flatness of the model's loss landscape can effectively improve its OOD generalization ability. Specifically, flatness means that perturbations to the model's weights result in minimal fluctuations in its performance, thereby contributing to OOD generalization. However, recent studies have revealed that quantization tends to sharpen the loss landscape [27]. Based on these facts, we propose the following proposition and, in the next section, present an optimization method from the perspective of flatness.

**Proposition:** *The sharp loss landscape caused by QAT may contribute to the degradation of OOD generalization ability.*

## 4 Method

In this paper, we propose a flatness-oriented QAT (FQAT) method to maintain the OOD generalization ability of the full-precision model during quantization. In this chapter, we first present the optimization objective of FQAT. Next, we analyze the gradient conflict issue of quantization parameters caused by dual optimization

**Figure 2: Results of cumulative gradients every 350 steps in the 4-bit test on the PACS ART domain, revealing conflicts of V-QAT gradient and flatness-oriented gradient.**

objectives and propose a layer-wise freezing mechanism. Finally, we introduce the disorder-guided adaptive freezing algorithm of FQAT.

## 4.1 Optimization Objective of FQAT

Following SAGM [45], we adopt three optimization objectives for flatness-oriented minimization over the training distributions $\mathcal{D}$: (a) empirical risk loss $\mathcal{L}(\theta; \mathcal{D})$, (b) perturbed loss $\mathcal{L}_p(\theta; \mathcal{D})$, and (c) the surrogate gap $h(\theta) \triangleq \mathcal{L}_p(\theta; \mathcal{D}) - \mathcal{L}(\theta; \mathcal{D})$, where $\theta$ is the parameters of DNNs. Minimizing $\mathcal{L}(\theta; \mathcal{D})$ and $\mathcal{L}_p(\theta; \mathcal{D})$ finds low-loss regions, while minimizing $h(\theta)$ ensures a flat minimum. This combination improves both training performance and OOD generalization. Hence, the overall optimization objective is:

$$\min_{\theta} \big( \mathcal{L}(\theta; \mathcal{D}), \mathcal{L}_p(\theta; \mathcal{D}), h(\theta) \big) \quad (1)$$

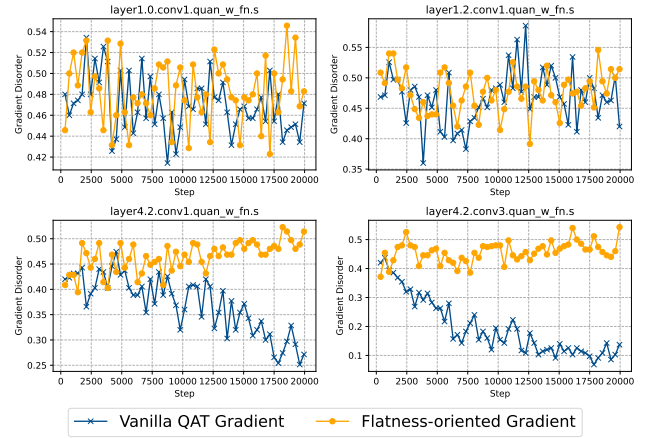SAGM [45] proposed the overall objective can be achieved by the following formulation :

$$\min_{\theta} \big( \mathcal{L}(\theta; \mathcal{D}) + \mathcal{L}_p(\theta - \alpha \nabla \mathcal{L}(\theta; \mathcal{D}); \mathcal{D}) \big), \quad (2)$$

where $\alpha$ is the hyperparameter. Then we incorporate QAT quantizer $Q(\theta; \mathbf{s})$ to the full-precison objective of flatness:

$$\min_{\theta, \mathbf{s}} \mathcal{L}\left(Q(\theta; \mathbf{s}); \mathcal{D}\right) + \mathcal{L}_p\Big( Q\big(\theta - \alpha \nabla \mathcal{L}\left(Q(\theta; \mathbf{s}); \mathcal{D}\right); \mathbf{s}\big); \mathcal{D} \Big), \quad (3)$$

where $\mathbf{s}$ denote the trainable parameters (scaling factors) of quantizers.

Specifically, during the optimization process, a gradient is obtained through the original QAT process. Based on this gradient, the original parameters are perturbed to obtain perturbed parameters. Then, a flat-related gradient is obtained by performing the QAT process with these perturbed parameters. By jointly update original parameters by the original quantization gradient and the flat-related gradient, a flatter loss landscape with low loss value is achieved. More details about the optimization objective are provided in supplementary materials.



**Figure 3: Results of V-QAT gradient and flatness-oriented gradient disorder of scale factors over 350 steps in the 4-bit test on the PACS ART domain, revealing in certain layers, the gradient disorder of V-QAT gradient decreases significantly as training progresses.**

**Table 3: Results of 4-bit quantization tests with perturbed scale factors on the Clipart and Infograph domains (Domain-Net). A / B indicates OOD test accuracy of Clipart and Infograph. Original OOD Performance is 60.21 / 15.81. The notation x% denotes that the scale factor is multiplied by x%. ↓: degradation; ↑: improvement; —: no change. Proved that the scale factors has not fully converged, and the sensitivity of different layers varies**

| Layer | 80% | 90% | 110% | 120% |
|---|---|---|---|---|
| L3.0.c1.w.s | 60.30↑ / 15.93↑ | 60.15↓ / 15.94↑ | 59.96↓ / 15.62↓ | 59.82↓ / 15.38↓ |
| L3.0.c1.a.s | 60.47↑ / 16.12↑ | 60.31↑ / 15.90↑ | 60.10↓ / 15.72↓ | 59.93↓ / 15.65↓ |
| L1.0.c1.w.s | 60.25↑ / 15.60↓ | 60.14↑ / 15.61↓ | 60.32↑ / 15.48↓ | 60.18↓ / 15.27↓ |
| L1.0.c1.a.s | 60.23↑ / 15.81– | 60.22↑ / 15.85↑ | 60.26↑ / 15.78↓ | 60.24↑ / 15.67↓ |

## 4.2 Layer-wise Freezing Mechanism for Gradient Conflict of Quantizer

As demonstrated by the results of *LSQ+SAGM* and *EWGS+SAGM* in Table 1, the OOD performance improve in certain cases (e.g., 5bit on PACS, 4, 5, 8bit on DomainNet), though directly adopting the objective. This phenomenon suggests the potential for enhancing quantized model performance from a flatness perspective. However, in some cases (for example, 3, 4bit on PACS), performance degrades significantly. We hypothesize that the improper optimization of the scale factor may be the cause of this performance degradation. The scale factor, used to portray the characteristic of weight and activation distribution [41], is highly sensitive to the perturbations [8, 29]. Optimization of quantization parameters (i.e., the scale factor in the quantizer) is of importance for OOD generalization.

Compared to full-precision training, Eq. (3) introduces several scale factors $\mathbf{s}$ in two optimization objective functions, thus generating two sets of gradients. One set originates from the optimization of vanilla QAT, denoted as $\mathbf{g}_{\text{va}}$ and derived from $\mathcal{L}(\cdot)$. The other

set is the newly introduced flatness-oriented gradient, aimed at enhancing generalization ability and denoted as $\mathbf{g}_{\text{flat}}$, derived from $\mathcal{L}_p(\cdot)$.

To investigate the interaction between these two gradients, we visualized their sum during the training process. As shown at the top of Figure 2, the flatness-oriented gradient (yellow) exhibits significantly higher volatility and more outliers compared to the vanilla QAT gradient (blue). Moreover, for certain layers, $\mathbf{g}_{\text{va}}$ and $\mathbf{g}_{\text{flat}}$ are in opposite directions and tend to cancel each other out (bottom of Figure 2). Here we define this phenomenon as the *sub-optimal equilibrium state*. This conflict between $\mathbf{g}_{\text{va}}$ and $\mathbf{g}_{\text{flat}}$ prevents the scale factors from fully converging, which can substantially degrade the performance of QAT in OOD scenarios.

To analyze gradient behavior from the perspective of directional fluctuation, we define the gradient disorder as a metric to quantify the degree of directional variability during training.

**Definition 4.1. Gradient Disorder:** For every $K$ steps of training, we define two gradient sequences: $G_1 = \{\mathbf{g}_1, \mathbf{g}_2, \ldots, \mathbf{g}_{K-1}\}$ and $G_2 = \{\mathbf{g}_2, \mathbf{g}_3, \ldots, \mathbf{g}_K\}$, where $\mathbf{g}_j$ denotes the gradient at step $j$. Let $\text{sgn}(\cdot)$ denote the element-wise sign function. The *gradient disorder* is defined as:

$$\delta = \frac{1}{K-1} \sum_{i=1}^{K-1} \mathbb{1}\left(\text{sgn}(G_1^{(i)}) \neq \text{sgn}(G_2^{(i)})\right), \quad (4)$$
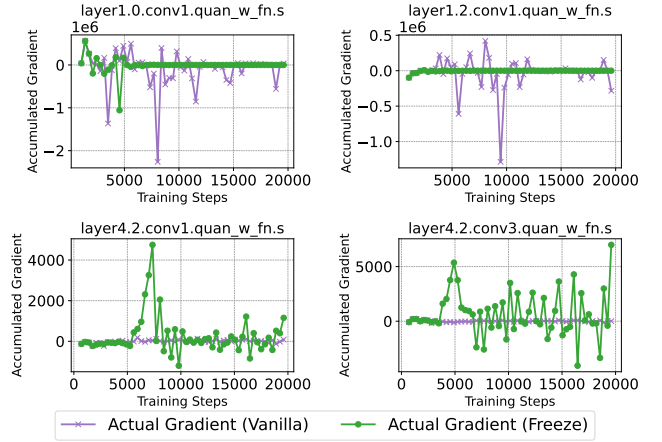
where $\mathbb{1}(\cdot)$ is the indicator function. $\delta$ measures the proportion of adjacent gradients with opposite directions in the gradient sequence, indicating the degree of directional fluctuation.

Figure 3 shows that in some layers, the gradient disorder of $\mathbf{g}_{\text{va}}$ decreases significantly during training, implying increasingly consistent gradient directions—a somewhat counterintuitive result. In contrast, $\mathbf{g}_{\text{flat}}$ gradient disorder remains high across layers. Layers with lower $\mathbf{g}_{\text{va}}$ disorder (bottom of Figure 3) display opposite and similar-magnitude gradients in Figure 2, suggesting a tendency to settle into a *sub-optimal equilibrium state*.

**Assumption 4.2. Impact of Scale Factors with Low $\mathbf{g}_{\text{va}}$ Gradient Disorder on Training:** Scale factors with overly low $\mathbf{g}_{\text{va}}$ gradient disorder hinder their $\mathbf{g}_{\text{flat}}$ training, causing insufficient convergence and indirectly affecting other scale factors convergence.

To verify this, we draw inspiration from weight freezing strategies [29, 36, 40] and conduct an experiment using the gradient disorder of $\mathbf{g}_{\text{va}}$ as a freezing indicator. Specifically, we freeze the $\mathbf{g}_{\text{va}}$ of scale factors once they begin to exhibit low gradient disorder, allowing only $\mathbf{g}_{\text{flat}}$ to update the parameters until the end of training.

Figure 2 shows the gradients of $\mathbf{g}_{\text{va}}$ and $\mathbf{g}_{\text{flat}}$ without freezing. In the top layers, $\mathbf{g}_{\text{flat}}$ exhibits large fluctuations and outliers, while in the bottom layers, $\mathbf{g}_{\text{flat}}$ and $\mathbf{g}_{\text{va}}$ reach *sub-optimal equilibrium state*. Figure 3 demonstrates that after some training, the gradient disorder of $\mathbf{g}_{\text{va}}$ decreases in the bottom layers. Thus, our freezing strategy freezes $\mathbf{g}_{\text{va}}$ in these layers after a certain number of steps, using only $\mathbf{g}_{\text{flat}}$ for updates. Note that the original update actual gradient is $\mathbf{g}_{\text{va}} + \mathbf{g}_{\text{flat}}$, but after freezing, the update gradient consists solely of $\mathbf{g}_{\text{flat}}$. Figure 4 compares the actual gradients with and without freezing, showing that the gradients of the upper layers
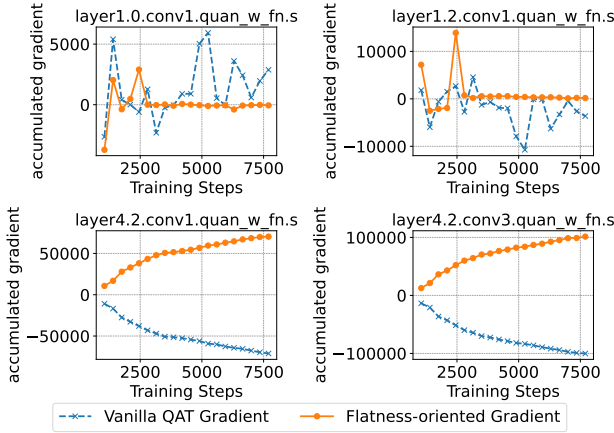


**Figure 4: Results of the 4-bit quantization test on the PACS ART domain demonstrate the impact of layer freezing over 350 training steps, with V-QAT gradient disorder serving as an indicator. The results show that freezing specific layers (the bottom layers) strengthens the training of their actual gradients (only flatness-oriented). Meanwhile, the actual gradients (V-QAT and flatness-oriented) of the unfrozen layers become more stable, with fewer outliers. This demonstrates that V-QAT gradient disorder effectively regulates gradient fluctuations across different layers, providing a reliable metric for optimizing layer-wise training dynamics.**

(high disorder, unfreeze) stabilize significantly after other layers freezing. In contrast, in the bottom layers (low disorder, freeze), after a certain number of steps, the actual gradient only involves $\mathbf{g}_{\text{flat}}$ and continues to evolve, indicating continued training rather than insufficient convergence due to the cancellation of $\mathbf{g}_{\text{va}}$. A more detailed explanation is as follows: regarding the dynamics of $g_{\text{flat}}$, once freezing is applied (e.g., to layer4.2.conv1 at 5k training steps), $g_{\text{flat}}$ exhibits two distinct phases: (1) an initial phase of unidirectional updates between approximately 5k and 7.5k steps, confirming the suppressive influence of $g_{\text{va}}$; and (2) a subsequent phase of balanced oscillations from 7.5k to 20k steps, indicating that the model is undergoing stable convergence, as shown in Fig 4. This confirms the hypothesis and validates that gradient disorder can serve as an effective indicator for guiding which layers to freeze, helping alleviate training instability.

## 4.3 Disorder-guided Adaptive Freezing Algorithm

To prevent suboptimal convergence resulting from the persistent freezing of $\mathbf{g}_{\text{va}}$ in certain layers without timely unfreezing, we design an adaptive freezing strategy: for every $K$ steps, we assess the gradient disorder $\delta_{t,\mathbf{s}_i}$ of $\mathbf{g}_{\text{va}}$ for each scale factor $\mathbf{s}_i$, calculated from the $\mathbf{g}_{\text{va}}$ sequence stored over the previous $K$ steps. If $\delta_{t,\mathbf{s}_i}$ falls below the threshold $r$, we freeze $\mathbf{g}_{\text{va}}$ of $\mathbf{s}_i$ for the next $K$ steps, updating only with $\mathbf{g}_{\text{flat}}$. Otherwise, we unfreeze $\mathbf{g}_{\text{va}}$, updating with $\mathbf{g}_{\text{va}} + \mathbf{g}_{\text{flat}}$. The complete procedure is detailed in supplementary materials.

**Figure 5: Results of cumulative gradients every 2111 steps in the 3-bit test on the DomainNet Clipart and Infograph domainsrevealing few anomalous gradients, with V-QAT gradient dominating.**

This adaptive selective freezing strategy effectively enhances model performance in OOD scenarios by adaptively managing gradient updates, while avoiding convergence issues due to excessive freezing.

## 5 Experiment

### 5.1 Experimental Setup and Implementation Details

**Datasets and evaluation protocol.** We conduct a comprehensive evaluation on three widely used OOD datasets: PACS [24], containing 9,991 images across 7 categories and 4 domains; OfficeHome [43] with 15,588 images spanning 65 categories and 4 domains; DomainNet [37], the largest benchmark comprising 586,575 images across 345 categories and 6 domains. We basically follow the evaluation protocol of DomainBed [10], including the optimizer, data split, and model selection, where we adopt train-domain validation as our model selection strategy for all algorithms in our experiments. For PACS and OfficeHome, for each time we treat one domain as the test domain and other domains as training domains, which is the leave-one-domain-out protocol commonly adopted in OOD evaluation. For DomainNet, following Yu et al. [48], we divide the domains into three groups: (1) *Clipart* and *Infograph*, (2) *Painting* and *Quickdraw*, and (3) *Real* and *Sketch*. Then we employ the leave-one-group-out protocol, where we treat one group of two domains as test domains and other two groups as training domains each time. For the number of training steps, for full-precision models we set it as 5,000 for PACS and OfficeHome, 15,000 for DomainNet following Cha et al. [4], while for quantization training we use 20,000 for PACS and OfficeHome, 50,000 for DomainNet. To reduce time cost, for quantization training we conduct validation and testing for DomainNet only after 45,000 steps. Each performance is reported as the average of two runs with seeds 0 and 23, while the ablation study is conducted with a single seed (seed 0) to reduce computational costs.

**Quantization.** We follow established practices in QAT literature by employing the QAT method LSQ [8] and EWGS [21] to quantize both weights and activations. The quantization scale factors are learned with a fixed learning rate of $1 \times 10^{-5}$. We use Mean Squared Error (MSE) range estimation [35] to determine the quantization parameters for weights and activations. Due to the risk of test data information leakage of supervised pretrained weights revealed by Yu et al. [48], we employ MoCo-v2 [5] pretrained ResNet-50 as initialization as recommended. Then we fine-tune the model using Empirical Risk Minimization (ERM) to obtain a full-precision model with generalization capabilities, which serves as the baseline for quantization. The weights and activations are fully quantized, except for the first convolutional layer, which quantizes only the activations, and the final linear layer, which remains unquantized, striking a balance between efficiency and model capacity. For DomainNet, we apply LSQ and EWGS to quantize model to 4, 5, and 8bit precision, as 3bit quantization yields minimal conflicts (Figure 5). For PACS, we use LSQ with 3, 4, and 5bit quantization. For OfficeHome, both LSQ and EWGS are applied at 3bit precision.

### 5.2 Hyperparameter settings

Given the substantial computational resources required by the original DomainBed setup, we adjust the hyperparameter search space and conduct grid search to reduce computational cost following SAGM [45]. The search space of learning rate is {1e-5, 3e-5, 5e-5}, and the dropout rate is fixed as zero. The batch size of each training domain is set as 32 for PACS and OfficeHome, 24 for DomainNet. Following SAM [9], we fix the hyperparameter $\rho = 0.05$. Following SAGM [45], we set $\alpha$ in Equation (3) as 0.001 for PACS and Office-Home, 0.0005 for DomainNet, and set weight decay as 1e-4 for PACS and OfficeHome, 1e-6 for DomainNet. For PACS and OfficeHome, the gradient disorder threshold $r$ is selected from {0.28, 0.30, 0.32} for 3-bit, 4-bit and 5-bit quantization. The number of freeze steps is selected from {300, 350, 400} for both 4-bit and 5-bit quantization, and from {100, 150, 200} for 3-bit quantization. For DomainNet, $r$ is selected from {0.20, 0.25} for both 4-bit and 5-bit quantization, and from {0.3, 0.35} for 8-bit quantization. The number of freeze steps is chosen from {3000, 4000} for both 4-bit and 5-bit quantization, and from {1500, 2000} for 8-bit quantization. We determine the hyperparameter space based on single-domain observation and apply it to other test domains, validating the robustness of our method. To reduce the high computational cost, we first select the shared hyperparameters, i.e. learning rate, weight decay, through grid search, which serve as the base hyperparameter configuration. Then we fix the base configuration and conduct further grid search on our specific hyperparameters, i.e. freeze steps, freeze threshold.

### 5.3 Main Results

We conducted a comprehensive evaluation of our method on the PACS, OfficeHome and DomainNet datasets across various quantization bit-widths (see Table 1), highlighting three key advantages: **(i) significant improvements in I.D and OOD performance**, **(ii) greater enhancement in OOD performance compared to I.D performance**, and **(iii) improved training stability**. Across all quantization bit-widths, our method achieved the best performance in both validation and test sets on the three datasets under the

quantization schemes. On the DomainNet dataset, the I.D performance surpassed the full-precision performance at most bits, with our method achieving an OOD test accuracy of 41.27% (LSQ) and 41.01% (EWGS) at 8-bit quantization, exceeding the full-precision accuracy of 40.95% and setting a new state-of-the-art result. At 4-bit and 5-bit quantization, the OOD performance also approached full-precision levels. In terms of OOD performance preservation, despite the baseline accuracy of the validation set being significantly higher than that of the test set, the test accuracy gains across various bit-widths on the DomainNet dataset (compared to directly introducing SAGM) remained close to the corresponding validation set gains. On the PACS dataset, the test accuracy gains consistently exceeded the validation set gains compare to directly introducing SAGM (3-bit: +5.24% test vs +2.61% validation ; 4-bit: +3.07% test vs +2.37% validation; 5-bit: +1.72% test vs +1.54% validation), further confirming the significant enhancement of our method's generalization capability beyond I.D optimization. For OfficeHome dataset, both LSQ and EWGS exhibit significant OOD performance degradation at low-bit quantization, while our method achieves substantial gains. In terms of training stability, our method focuses on resolving gradient conflicts. Experimental results show that, in most cases, the standard deviations of our method's I.D and OOD performance are smaller than those of directly introducing SAGM, validating the effectiveness of our gradient coordination mechanism. More experimental results are provided in supplementary materials.

## 5.4 Ablation Study

In our analysis, we validated the effectiveness of freezing $g_{va}$ when gradient disorder falls below a specific threshold, coupled with periodically reselecting the freeze set to stabilize quantization training in the OOD scenario. This naturally leads to the question: what would occur if these strategies were modified? For instance, what if scale factors with gradient disorder above the threshold were frozen instead, or both gradients ($g_{va}$ and $g_{flat}$) were frozen simultaneously, or unfreezing were avoided after the initial freeze? Furthermore, what would happen if a decoupled alternating update approach were adopted, where one step updates $g_{va}$ and the next updates $g_{flat}$? Exploring these variations would offer deeper insights into the effectiveness of each component of our method.

Table 4 highlights the importance of our proposed modules. The *w/o Unfreeze* strategy shows a significant performance drop on PACS, emphasizing the necessity of our phased freezing strategy. The *Reverse Freeze* strategy also results in notable accuracy degradation, validating the effectiveness of our freezing metric. Other strategies further demonstrate the superiority of our gradient selection methodology, confirming the robustness and effectiveness of our approach.

## 5.5 Hyperparameter Sensitivity Analysis

Figure 6 demonstrates the robustness of our method within a reasonable hyperparameter search space. We evaluated the performance of our approach on PACS (4-bit) using fixed hyperparameters across all domains (note that grid-searching for optimal hyperparameters per domain would yield better results). The results show that our method outperforms the baseline in both validation and test accuracy in most cases, indicating its insensitivity to hyperparameter
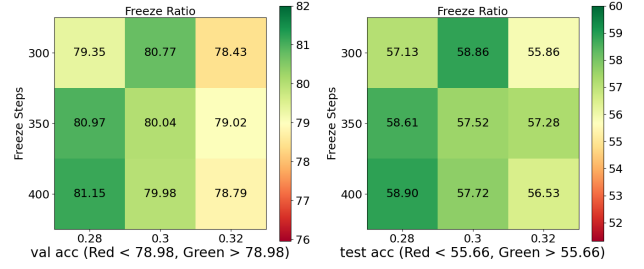


Figure 6: 4-bit LSQ PACS (seed 0, 23): Val acc (I.D mean) and test acc (OOD mean) were evaluated via grid search on freeze steps and freeze ratio. Each cell shows mean results for all domains with the corresponding hyperparameters. Green: improvement over directly introducing flat objective; red: degradation. Baseline: Val 78.98 / Test 55.66. Most configurations show improvements, proving our method's robustness to hyperparameters within a reasonable range.

Table 4: Ablation Study: LSQ experiments on seed 0 evaluate Val (I.D validation mean accuracy) and Test (OOD test mean accuracy). Compared variants include: Ours (proposed FQAT), Alter Update (Alternate V-QAT and flatness-oriented gradients updates), Freeze Both (freeze QAT and flatness-oriented gradients), w/o Unfreeze (no unfreezing after initial freeze), and Reverse Freeze (freeze above the threshold, unlike the original). Best performance is in bold, validating the effectiveness of our design.

| Dataset | Method | Bit-width (W/A) | Val | Test |
|---|---|---|---|---|
| DomainNet | Ours | 4/4 | 62.45 | **40.60** |
| | Alter Update | 4/4 | 57.41 | 37.15 |
| | Freeze Both | 4/4 | 61.03 | 39.36 |
| | w/o Unfreeze | 4/4 | **62.50** | 40.53 |
| | Reverse Freeze | 4/4 | 59.32 | 38.24 |
| PACS | Ours | 4/4 | **81.20** | **59.33** |
| | Alter Update | 4/4 | 80.66 | 58.36 |
| | Freeze Both | 4/4 | 79.74 | 56.98 |
| | w/o Unfreeze | 4/4 | 79.50 | 57.20 |
| | Reverse Freeze | 4/4 | 80.16 | 56.89 |

variations. This suggests that simply defining a reasonable search range is sufficient for practical use.

During the experiments, we randomly selected one domain to observe the gradient disorder of $g_{va}$ under different freeze steps. Based on empirical observations, we determined the search space for freeze steps and freeze ratio, which was then directly applied to other domains without further tuning.

## 5.6 Loss Surface Visualization

Following the approach in [25], Figure 7 compares the loss surface visualizations across the four domains of PACS (4-bit LSQ) when using SAGM directly versus our proposed method. The results demonstrate that our method consistently achieves flatter loss surfaces with lower loss values across all domains, highlighting its effectiveness in optimizing the loss landscape. We also visualize

**Figure 7: Visualization of the loss landscape across different domains (PACS 4bit), where the categories from left to right are Art, Cartoon, Photo, and Sketch. Blue indicates low loss values, while red represents high loss values. The top row displays the results of LSQ + SAGM, and the bottom row shows the results using our proposed method. Our approach consistently achieves flatter loss surfaces with lower loss values across all four domains in PACS, demonstrating its effectiveness in optimizing the loss landscape.**

loss surface of DomainNet (Clipart & Infograph domains) in supplementary materials. LSQ + SAGM shows a smoother loss surface compared to PACS, which further explains why our approach is more significantly effective on PACS than on DomainNet.

## 6 Related Work

### 6.1 Quantizaion-aware Training on CV

Quantizaion-aware training (QAT) involves inserting simulated quantization nodes and retraining the model, which achieves a better balance between accuracy and compression ratio [14, 34]. DoReFa [55] and PACT [6] use low-precision weights and activations during the forward pass and utilize STE techniques [2] during backpropagation to estimate gradients of the piece-wise quantization functions. LSQ [8] adjusts the quantization function by introducing learnable step size scaling factors. EWGS [21] scales the gradients to precisely adjust the quantization error component in the gradients, reducing gradient conflicts caused by quantization errors. SAQ [27] demonstrates that quantized models make the loss surface sharper. Recently, some works have explored the possibility of improving quantization performance by freezing unstable weights to further enhance results [29, 36, 40]; however, these methods have only considered the identically distributed (I.D) data. Due to distribution shifts in test data—which often occur in practical applications—the quality and reliability of quantized models cannot be guaranteed [12].

### 6.2 OOD Generalization on CV

In practical applications, when deploying machine learning models, test data distribution may differ from the training distribution, a common phenomenon known as distribution shift [20, 28, 47]. OOD generalization methods aim to enhance a model's ability to generalize to OOD distributions [44, 53]. Common strategies include domain alignment [26, 33, 52], meta learning [1, 7, 23], data augmentation [3, 54], disentangled representation learning [49] and

utilization of causal relations [31, 32]. Inspired by previous studies of flat minima [9, 16, 30, 51, 56], flatness-aware methods start to gain attention and exhibit remarkable performance in OOD generalization [4, 45, 50], such as SAGM [45], improving generalization ability by optimizing the angle between weight gradients to reduce gradient conflicts. As far as we know, existing works all focus on improving the OOD generalization ability of full-precision models, while neglecting the generalization issues of quantized models. Our work is the first to fill this gap by identifying the issue and offering a solution.

## 7 Conclusion and Future Work

Our work validates a crucial yet overlooked problem: quantization harms the OOD gereralization performance of full-precision CV models. Further, we proposes FQAT, a flatness-oriented QAT method, incorporating a layer-wise freezing mechanism to mitigate gradient conflicts and a disorder-guided adaptive freezing algorithm for dynamic layer adjustment. The extensive experiments on several OOD benchmarks demonstrate the superiority of the proposed method over state-of-the-art baselines.

**Limitations and future work.** Due to limitations in experimental resources, we did not explore alternative flatness objectives to further investigate their effects on scale factor gradients. Additionally, we recommend evaluating the robustness of our approach in conjunction with various full-precision OOD generalization techniques. Our experiments also revealed that different domains exhibit varying sensitivities to the scale factor, suggesting that a deeper investigation into the relationship between domain characteristics and scale factor dynamics could offer promising directions for future optimization. Finally, our current analysis is primarily conducted under the ReLU activation setting, where all QAT zero points are aligned at zero. Extending the analysis to other activation functions and quantization settings remains an important avenue for future work.

# Acknowledgments

# References

[1] Yogesh Balaji, Swami Sankaranarayanan, and Rama Chellappa. 2018. Metareg: Towards domain generalization using meta-regularization. *Advances in neural information processing systems* 31 (2018).

[2] Yoshua Bengio, Nicholas Léonard, and Aaron Courville. 2013. Estimating or propagating gradients through stochastic neurons for conditional computation. *arXiv preprint arXiv:1308.3432* (2013).

[3] Fabio M Carlucci, Antonio D'Innocente, Silvia Bucci, Barbara Caputo, and Tatiana Tommasi. 2019. Domain generalization by solving jigsaw puzzles. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 2229–2238.

[4] Junbum Cha, Sanghyuk Chun, Kyungjae Lee, Han-Cheol Cho, Seunghyun Park, Yunsung Lee, and Sungrae Park. 2021. Swad: Domain generalization by seeking flat minima. *Advances in Neural Information Processing Systems* 34 (2021), 22405–22418.

[5] Xinlei Chen, Haoqi Fan, Ross Girshick, and Kaiming He. 2020. Improved Baselines with Momentum Contrastive Learning. arXiv:2003.04297 [cs.CV] https://arxiv.org/abs/2003.04297

[6] Jungwook Choi, Zhuo Wang, Swagath Venkataramani, Pierce I-Jen Chuang, Vijayalakshmi Srinivasan, and Kailash Gopalakrishnan. 2018. Pact: Parameterized clipping activation for quantized neural networks. *arXiv preprint arXiv:1805.06085* (2018).

[7] Qi Dou, Daniel Coelho de Castro, Konstantinos Kamnitsas, and Ben Glocker. 2019. Domain generalization via model-agnostic learning of semantic features. *Advances in neural information processing systems* 32 (2019).

[8] Steven K Esser, Jeffrey L McKinstry, Deepika Bablani, Rathinakumar Appuswamy, and Dharmendra S Modha. 2019. Learned step size quantization. *arXiv preprint arXiv:1902.08153* (2019).

[9] Pierre Foret, Ariel Kleiner, Hossein Mobahi, and Behnam Neyshabur. 2020. Sharpness-aware minimization for efficiently improving generalization. *arXiv preprint arXiv:2010.01412* (2020).

[10] Ishaan Gulrajani and David Lopez-Paz. 2020. In search of lost domain generalization. *arXiv preprint arXiv:2007.01434* (2020).

[11] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2016. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 770–778.

[12] Qiang Hu, Yuejun Guo, Maxime Cordy, Xiaofei Xie, Wei Ma, Mike Papadakis, and Yves Le Traon. 2022. Characterizing and understanding the behavior of quantized models for reliable deployment. *arXiv preprint arXiv:2204.04220* (2022).

[13] Zeyi Huang, Haohan Wang, Eric P Xing, and Dong Huang. 2020. Self-challenging improves cross-domain generalization. In *Computer vision–ECCV 2020: 16th European conference, Glasgow, UK, August 23–28, 2020, proceedings, part II 16*. Springer, 124–140.

[14] Itay Hubara, Yury Nahshan, Yair Hanani, Ron Banner, and Daniel Soudry. 2021. Accurate post training quantization with small calibration sets. In *International Conference on Machine Learning*. PMLR, 4466–4475.

[15] Brody Huval, Tao Wang, Sameep Tandon, Jeff Kiske, Will Song, Joel Pazhayampallil, Mykhaylo Andriluka, Pranav Rajpurkar, Toki Migimatsu, Royce Cheng-Yue, et al. 2015. An empirical evaluation of deep learning on highway driving. *arXiv preprint arXiv:1504.01716* (2015).

[16] Pavel Izmailov, Dmitrii Podoprikhin, Timur Garipov, Dmitry Vetrov, and Andrew Gordon Wilson. 2018. Averaging weights leads to wider optima and better generalization. *arXiv preprint arXiv:1803.05407* (2018).

[17] Saqib Javed, Hieu Le, and Mathieu Salzmann. 2024. QT-DoG: Quantization-aware Training for Domain Generalization. *arXiv preprint arXiv:2410.06320* (2024).

[18] Norman P Jouppi, Cliff Young, Nishant Patil, David Patterson, Gaurav Agrawal, Raminder Bajwa, Sarah Bates, Suresh Bhatia, Nan Boden, Al Borchers, et al. 2017. In-datacenter performance analysis of a tensor processing unit. In *Proceedings of the 44th annual international symposium on computer architecture*. 1–12.

[19] Sangil Jung, Changyong Son, Seohyung Lee, Jinwoo Son, Jae-Joon Han, Youngjun Kwak, Sung Ju Hwang, and Changkyu Choi. 2019. Learning to quantize deep networks by optimizing quantization intervals with task loss. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 4350–4359.

[20] Pang Wei Koh, Shiori Sagawa, Henrik Marklund, Sang Michael Xie, Marvin Zhang, Akshay Balsubramani, Weihua Hu, Michihiro Yasunaga, Richard Lanas Phillips, Irena Gao, et al. 2021. Wilds: A benchmark of in-the-wild distribution shifts. In *International conference on machine learning*. PMLR, 5637–5664.

[21] Junghyup Lee, Dohyung Kim, and Bumsub Ham. 2021. Network quantization with element-wise gradient scaling. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 6448–6457.

[22] Jesse Levinson, Jake Askeland, Jan Becker, Jennifer Dolson, David Held, Soeren Kammel, J Zico Kolter, Dirk Langer, Oliver Pink, Vaughan Pratt, et al. 2011. Towards fully autonomous driving: Systems and algorithms. In *2011 IEEE intelligent vehicles symposium (IV)*. IEEE, 163–168.

[23] Da Li, Yongxin Yang, Yi-Zhe Song, and Timothy Hospedales. 2018. Learning to generalize: Meta-learning for domain generalization. In *Proceedings of the AAAI conference on artificial intelligence*, Vol. 32.

[24] Da Li, Yongxin Yang, Yi-Zhe Song, and Timothy M. Hospedales. 2017. Deeper, Broader and Artier Domain Generalization. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*.

[25] Hao Li, Zheng Xu, Gavin Taylor, Christoph Studer, and Tom Goldstein. 2018. Visualizing the loss landscape of neural nets. *Advances in neural information processing systems* 31 (2018).

[26] Ya Li, Mingming Gong, Xinmei Tian, Tongliang Liu, and Dacheng Tao. 2018. Domain generalization via conditional invariant representations. In *Proceedings of the AAAI conference on artificial intelligence*, Vol. 32.

[27] Jing Liu, Jianfei Cai, and Bohan Zhuang. 2021. Sharpness-aware quantization for deep neural networks. *arXiv preprint arXiv:2111.12273* (2021).

[28] Jiashuo Liu, Zheyan Shen, Yue He, Xingxuan Zhang, Renzhe Xu, Han Yu, and Peng Cui. 2021. Towards out-of-distribution generalization: A survey. *arXiv preprint arXiv:2108.13624* (2021).

[29] Shih-Yang Liu, Zechun Liu, and Kwang-Ting Cheng. 2023. Oscillation-free quantization for low-bit vision transformers. In *International Conference on Machine Learning*. PMLR, 21813–21824.

[30] Yong Liu, Siqi Mai, Xiangning Chen, Cho-Jui Hsieh, and Yang You. 2022. Towards efficient and scalable sharpness-aware minimization. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 12360–12370.

[31] Fangrui Lv, Jian Liang, Shuang Li, Bin Zang, Chi Harold Liu, Ziteng Wang, and Di Liu. 2022. Causality inspired representation learning for domain generalization. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 8046–8056.

[32] Divyat Mahajan, Shruti Tople, and Amit Sharma. 2021. Domain generalization using causal matching. In *International conference on machine learning*. PMLR, 7313–7324.

[33] Krikamol Muandet, David Balduzzi, and Bernhard Schölkopf. 2013. Domain generalization via invariant feature representation. In *International conference on machine learning*. PMLR, 10–18.

[34] Markus Nagel, Rana Ali Amjad, Mart Van Baalen, Christos Louizos, and Tijmen Blankevoort. 2020. Up or Down? Adaptive Rounding for Post-Training Quantization. In *Proceedings of the 37th International Conference on Machine Learning (Proceedings of Machine Learning Research, Vol. 119)*, Hal Daumé III and Aarti Singh (Eds.). PMLR, 7197–7206. https://proceedings.mlr.press/v119/nagel20a.html

[35] Markus Nagel, Marios Fournarakis, Rana Ali Amjad, Yelysei Bondarenko, Mart van Baalen, and Tijmen Blankevoort. 2021. A White Paper on Neural Network Quantization. arXiv:2106.08295 [cs.LG] https://arxiv.org/abs/2106.08295

[36] Markus Nagel, Marios Fournarakis, Yelysei Bondarenko, and Tijmen Blankevoort. 2022. Overcoming oscillations in quantization-aware training. In *International Conference on Machine Learning*. PMLR, 16318–16330.

[37] Xingchao Peng, Zijun Huang, Ximeng Sun, and Kate Saenko. 2019. Domain agnostic learning with disentangled representations. In *International conference on machine learning*. PMLR, 5102–5112.

[38] Jiantao Qiu, Jie Wang, Song Yao, Kaiyuan Guo, Boxun Li, Erjin Zhou, Jincheng Yu, Tianqi Tang, Ningyi Xu, Sen Song, et al. 2016. Going deeper with embedded FPGA platform for convolutional neural network. In *Proceedings of the 2016 ACM/SIGDA international symposium on field-programmable gate arrays*. 26–35.

[39] Juncheol Shin, Junhyuk So, Sein Park, Seungyeop Kang, Sungjoo Yoo, and Eunhyeok Park. 2023. NIPQ: Noise proxy-based integrated pseudo-quantization. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 3852–3861.

[40] Chen Tang, Yuan Meng, Jiacheng Jiang, Shuzhao Xie, Rongwei Lu, Xinzhu Ma, Zhi Wang, and Wenwu Zhu. 2024. Retraining-free model quantization via oneshot weight-coupling learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 15855–15865.

[41] Chen Tang, Kai Ouyang, Zhi Wang, Yifei Zhu, Wen Ji, Yaowei Wang, and Wenwu Zhu. 2022. Mixed-precision neural network quantization via learned layer-wise importance. In *European Conference on Computer Vision*. Springer, 259–275.

[42] Vladimir Vapnik. 1991. Principles of risk minimization for learning theory. *Advances in neural information processing systems* 4 (1991).

[43] Hemanth Venkateswara, Jose Eusebio, Shayok Chakraborty, and Sethuraman Panchanathan. 2017. Deep hashing network for unsupervised domain adaptation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 5018–5027.

[44] Jindong Wang, Cuiling Lan, Chang Liu, Yidong Ouyang, Tao Qin, Wang Lu, Yiqiang Chen, Wenjun Zeng, and S Yu Philip. 2022. Generalizing to unseen domains: A survey on domain generalization. *IEEE transactions on knowledge and data engineering* 35, 8 (2022), 8052–8072.

[45] Pengfei Wang, Zhaoxiang Zhang, Zhen Lei, and Lei Zhang. 2023. Sharpness-aware gradient matching for domain generalization. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 3769–3778.

[46] Jiaming Yang, Chenwei Tang, Caiyang Yu, and Jiancheng Lv. 2024. GWQ: Group-Wise Quantization Framework for Neural Networks. In *Asian Conference on Machine Learning*. PMLR, 1526–1541.

[47] Han Yu, Jiashuo Liu, Xingxuan Zhang, Jiayun Wu, and Peng Cui. 2024. A survey on evaluation of out-of-distribution generalization. *arXiv preprint arXiv:2403.01874* (2024).

[48] Han Yu, Xingxuan Zhang, Renzhe Xu, Jiashuo Liu, Yue He, and Peng Cui. 2024. Rethinking the evaluation protocol of domain generalization. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 21897–21908.

[49] Hanlin Zhang, Yi-Fan Zhang, Weiyang Liu, Adrian Weller, Bernhard Schölkopf, and Eric P Xing. 2022. Towards principled disentanglement for domain generalization. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 8024–8034.

[50] Xingxuan Zhang, Renzhe Xu, Han Yu, Yancheng Dong, Pengfei Tian, and Peng Cui. 2023. Flatness-aware minimization for domain generalization. In *Proceedings*

of the IEEE/CVF International Conference on Computer Vision. 5189–5202.

[51] Xingxuan Zhang, Renzhe Xu, Han Yu, Hao Zou, and Peng Cui. 2023. Gradient norm aware minimization seeks first-order flatness and improves generalization. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 20247–20257.

[52] Shanshan Zhao, Mingming Gong, Tongliang Liu, Huan Fu, and Dacheng Tao. 2020. Domain generalization via entropy regularization. *Advances in neural information processing systems* 33 (2020), 16096–16107.

[53] Kaiyang Zhou, Ziwei Liu, Yu Qiao, Tao Xiang, and Chen Change Loy. 2022. Domain generalization: A survey. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 45, 4 (2022), 4396–4415.

[54] Kaiyang Zhou, Yongxin Yang, Yu Qiao, and Tao Xiang. 2021. Domain generalization with mixstyle. *arXiv preprint arXiv:2104.02008* (2021).

[55] Shuchang Zhou, Yuxin Wu, Zekun Ni, Xinyu Zhou, He Wen, and Yuheng Zou. 2016. Dorefa-net: Training low bitwidth convolutional neural networks with low bitwidth gradients. *arXiv preprint arXiv:1606.06160* (2016).

[56] Juntang Zhuang, Boqing Gong, Liangzhe Yuan, Yin Cui, Hartwig Adam, Nicha C Dvornek, James s Duncan, Ting Liu, et al. 2022. Surrogate Gap Minimization Improves Sharpness-Aware Training. In *International Conference on Learning Representations*.

## A  Optimization Objective Explanation

The overall objective of SAGM can be achieved by the following formulation:

$$\min_{\theta}\big(\mathcal{L}(\theta;\mathcal{D}) + \mathcal{L}_p(\theta - \alpha\nabla_\theta\mathcal{L}(\theta;\mathcal{D});\mathcal{D})\big),$$

where $\alpha$ is a hyperparameter. The second term can be further rewritten as in SAGM:

$$\min_{\theta}\mathcal{L}\left(\theta + \hat{\epsilon} - \alpha\nabla_\theta\mathcal{L}(\theta;\mathcal{D});\mathcal{D}\right),$$

where

$$\hat{\epsilon} = \rho\frac{\nabla_\theta\mathcal{L}(\theta;\mathcal{D})}{\|\nabla_\theta\mathcal{L}(\theta;\mathcal{D})\|},$$

$\rho$ is a hyperparameter. The process begins with backpropagation based on the loss function $\mathcal{L}(\theta;\mathcal{D})$ to obtain the gradient $\nabla_\theta\mathcal{L}(\theta;\mathcal{D})$. The perturbation magnitude $\hat{\epsilon}$ is then calculated to perturb the original weights $\theta$ to obtain new weights $\theta_p = \theta + \hat{\epsilon}$. Next, the SAGM target perturbed weights are computed as:

$$\theta_{\text{SAGM}} = \theta_p - \alpha\nabla_\theta\mathcal{L}(\theta;\mathcal{D}).$$

Using these new weights, the loss is recalculated as $\mathcal{L}\left(\theta_{\text{SAGM}};\mathcal{D}\right)$, yielding a new gradient $\nabla_{\theta_{\text{SAGM}}}\mathcal{L}(\theta_{\text{SAGM}};\mathcal{D})$. Finally, the original weights $\theta$ are updated using both $\nabla_{\theta_{\text{SAGM}}}\mathcal{L}(\theta_{\text{SAGM}};\mathcal{D})$ and $\nabla_\theta\mathcal{L}(\theta;\mathcal{D})$, specifically by averaging the two gradients during implementation.

When directly integrating the SAGM objective into our quantization process, the key difference lies in the computation of gradients. The first gradient is obtained using quantized weights $Q(\theta,s)$ through the Straight-Through Estimator, i.e., $\nabla_\theta\mathcal{L}(Q(\theta,s);\mathcal{D})$. The perturbation magnitude is then calculated as:

$$\hat{\epsilon} = \rho\frac{\nabla_\theta\mathcal{L}(Q(\theta,\mathbf{s});\mathcal{D})}{\|\nabla_\theta\mathcal{L}(Q(\theta,\mathbf{s});\mathcal{D})\|},$$

resulting in $\theta' = \theta + \hat{\epsilon}$. For the second gradient, the actual computation uses the quantized perturbed weights $Q(\theta',\mathbf{s})$, yielding $\nabla_{\theta'}\mathcal{L}(Q(\theta',\mathbf{s});\mathcal{D})$. These two gradients are used to update $\theta$, while the scale parameter $\mathbf{s}$ remains unperturbed. During this process, $\mathbf{s}$ also receives two gradients, $\mathbf{g}_{\text{va}} = \nabla_s\mathcal{L}(Q(\theta,\mathbf{s});\mathcal{D})$ and $\mathbf{g}_{\text{flat}} = \nabla_s\mathcal{L}(Q(\theta',\mathbf{s});\mathcal{D})$, which are used to update $\mathbf{s}$.

## B  pseudocode

---

**Algorithm 1** Disorder-guided Adaptive Freezing Algorithm

---

**Require:** overall training steps $T$, step interval $K$, threshold $r$ and scale factors $\{\mathbf{s}_1, \mathbf{s}_2, \ldots, \mathbf{s}_n\}$

1: Init $t \leftarrow 0$, freeze$[\mathbf{s}_i] \leftarrow False$ for all $\mathbf{s}_i, i = 1, \cdots, n$
2: **while** $t < T$ **do**
3:     **for** $i = 1, \cdots, n$ **do**
4:         **if** freeze$[\mathbf{s}_i]$ **then**
5:             Update $\mathbf{s}_i$ with $\mathbf{g}_{\text{flat}}$
6:         **else**
7:             Update $\mathbf{s}_i$ with $\mathbf{g}_{\text{va}}, \mathbf{g}_{\text{flat}}$
8:         **end if**
9:     **end for**
10:     **if** $t \bmod K = 0$ **then**
11:         **for** $i = 1, \cdots, n$ **do**
12:             Compute $\delta_{t,\mathbf{s}_i}$
13:             **if** $\delta_{t,\mathbf{s}_i} < r$ **then**
14:                 freeze$[\mathbf{s}_i] \leftarrow True$
15:             **else**
16:                 freeze$[\mathbf{s}_i] \leftarrow False$
17:             **end if**
18:         **end for**
19:     **end if**
20:     $t \leftarrow t + 1$
21: **end while**
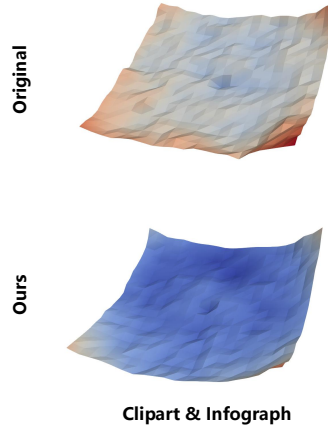
---

## C  More experimental results

**Figure 8: Visualization of the loss landscape of DomainNet Clipart & Infograph domains (4bit). The top row displays the results of LSQ + SAGM, and the bottom row shows the results using our proposed method.**

**Table 5: 4-bit LSQ Results on DomainNet Validation Set (val)**

| Method | Clipart | Infograph | Painting | Quickdraw | Real | Sketch | Avg |
|--------|---------|-----------|----------|-----------|------|--------|-----|
| LSQ | 65.35±0.98 | 65.35±0.98 | 59.71±0.15 | 59.71±0.15 | 57.80±0.02 | 57.80±0.02 | 60.95±0.38 |
| LSQ + SAGM | 66.13±0.36 | 66.13±0.36 | 61.15±0.05 | 61.15±0.05 | 58.09±0.37 | 58.09±0.37 | 61.79±0.26 |
| Ours | 67.19±0.01 | 67.19±0.01 | 61.83±0.68 | 61.83±0.68 | 58.54±0.05 | 58.54±0.05 | 62.52±0.25 |

**Table 6: 4-bit LSQ Results on DomainNet Test Set (test)**

| Method | Clipart | Infograph | Painting | Quickdraw | Real | Sketch | Avg |
|--------|---------|-----------|----------|-----------|------|--------|-----|
| LSQ | 59.86±0.59 | 15.21±0.44 | 44.81±0.13 | 14.74±0.01 | 52.51±0.19 | 47.72±0.10 | 39.14±0.24 |
| LSQ + SAGM | 60.95±0.22 | 15.73±0.09 | 46.75±0.08 | 16.09±0.20 | 52.23±0.39 | 48.69±0.38 | 40.07±0.23 |
| Ours | 61.13±0.13 | 16.13±0.01 | 47.02±0.79 | 16.43±0.01 | 53.47±0.06 | 49.43±0.14 | 40.6±0.19 |

**Table 7: 5-bit LSQ Results on DomainNet Validation Set (val)**

| Method | Clipart | Infograph | Painting | Quickdraw | Real | Sketch | Avg |
|--------|---------|-----------|----------|-----------|------|--------|-----|
| LSQ | 65.22±0.62 | 65.22±0.62 | 59.64±0.23 | 59.64±0.23 | 56.44±0.27 | 56.44±0.27 | 60.43±0.37 |
| LSQ + SAGM | 67.36±0.51 | 67.36±0.51 | 61.86±0.31 | 61.86±0.31 | 58.97±0.44 | 58.97±0.44 | 62.73±0.42 |
| Ours | 67.60±0.52 | 67.60±0.52 | 62.24±0.28 | 62.24±0.28 | 59.74±0.28 | 59.74±0.28 | 63.19±0.36 |

**Table 8: 5-bit LSQ Results on DomainNet Test Set (test)**

| Method | Clipart | Infograph | Painting | Quickdraw | Real | Sketch | Avg |
|--------|---------|-----------|----------|-----------|------|--------|-----|
| LSQ | 60.06±0.26 | 15.51±0.20 | 45.14±0.14 | 14.30±0.03 | 49.72±0.30 | 46.05±0.03 | 38.46±0.16 |
| LSQ + SAGM | 61.41±0.19 | 16.40±0.23 | 47.76±0.51 | 15.71±0.08 | 52.73±0.68 | 49.03±0.53 | 40.51±0.37 |
| Ours | 61.77±0.26 | 16.31±0.42 | 47.93±0.23 | 15.83±0.01 | 53.91±0.45 | 49.82±0.16 | 40.93±0.25 |

**Table 9: 8-bit LSQ Results on DomainNet Validation Set (val)**

| Method | Clipart | Infograph | Painting | Quickdraw | Real | Sketch | Avg |
|---|---|---|---|---|---|---|---|
| LSQ | 66.21±0.54 | 66.21±0.54 | 59.96±0.11 | 59.96±0.11 | 56.74±0.27 | 56.74±0.27 | 60.97±0.31 |
| LSQ + SAGM | 68.60±0.09 | 68.60±0.09 | 62.38±0.29 | 62.38±0.29 | 58.64±0.71 | 58.64±0.71 | 63.21±0.36 |
| Ours | 68.40±0.15 | 68.40±0.15 | 62.40±0.37 | 62.40±0.37 | 59.37±0.11 | 59.37±0.11 | 63.39±0.21 |

**Table 10: 8-bit LSQ Results on DomainNet Test Set (test)**

| Method | Clipart | Infograph | Painting | Quickdraw | Real | Sketch | Avg |
|---|---|---|---|---|---|---|---|
| LSQ | 60.18±0.11 | 16.02±0.30 | 46.19±0.06 | 13.96±0.01 | 51.22±0.25 | 47.05±0.21 | 39.1±0.16 |
| LSQ + SAGM | 62.49±0.13 | 17.45±0.07 | 48.22±0.34 | 15.26±0.04 | 53.80±1.02 | 49.37±0.51 | 41.1±0.35 |
| Ours | 62.23±0.17 | 17.24±0.10 | 48.33±0.30 | 15.40±0.18 | 54.54±0.21 | 49.90±0.13 | 41.27±0.18 |

**Table 11: 3-bit Quantization Results on PACS Validation Set (val)**

| Method | Art | Cartoon | Photo | Sketch | Avg |
|---|---|---|---|---|---|
| LSQ | 81.64±0.52 | 73.47±1.51 | 79.72±0.52 | 74.72±0.16 | 77.39±0.68 |
| LSQ + SAGM | 82.44±1.03 | 70.07±2.68 | 76.01±1.79 | 71.40±0.47 | 74.98±1.49 |
| Ours | 82.81±0.20 | 74.23±1.54 | 78.62±1.78 | 74.68±1.30 | 77.59±1.21 |

**Table 12: 3-bit Quantization Results on PACS Test Set (test)**

| Method | Art | Cartoon | Photo | Sketch | Avg |
|---|---|---|---|---|---|
| LSQ | 39.26±0.82 | 54.74±3.94 | 61.23±1.95 | 62.18±2.23 | 54.35±2.23 |
| LSQ + SAGM | 40.88±2.68 | 47.28±1.44 | 59.09±0.94 | 55.14±0.43 | 50.6±1.37 |
| Ours | 43.14±0.43 | 55.46±3.60 | 63.14±1.38 | 61.63±0.49 | 55.84±1.48 |

**Table 13: 4-bit Quantization Results on PACS Validation Set (val)**

| Method | Art | Cartoon | Photo | Sketch | Avg |
|---|---|---|---|---|---|
| LSQ | 86.57±1.71 | 78.30±0.43 | 81.61±0.80 | 75.28±0.32 | 80.44±0.82 |
| LSQ + SAGM | 85.18±1.02 | 77.05±3.11 | 80.19±1.60 | 73.50±0.11 | 78.98±1.46 |
| Ours | 85.93±0.82 | 80.08±1.65 | 81.98±1.81 | 77.04±0.21 | 81.26±1.12 |

**Table 14: 4-bit Quantization Results on PACS Test Set (test)**

| Method | Art | Cartoon | Photo | Sketch | Avg |
|---|---|---|---|---|---|
| LSQ | 46.25±4.82 | 57.22±0.88 | 63.51±0.26 | 62.63±0.35 | 57.4±1.58 |
| LSQ + SAGM | 45.97±0.52 | 53.94±2.72 | 62.57±2.10 | 60.16±1.35 | 55.66±1.67 |
| Ours | 47.41±1.83 | 58.02±1.36 | 65.98±0.71 | 66.28±1.11 | 59.42±1.25 |

**Table 15: 5-bit Quantization Results on PACS Validation Set (val)**

| Method | Art | Cartoon | Photo | Sketch | Avg |
|---|---|---|---|---|---|
| LSQ | 84.53±0.04 | 80.89±2.93 | 83.81±1.03 | 71.05±0.21 | 80.07±1.05 |
| LSQ + SAGM | 86.12±2.60 | 80.45±3.78 | 81.42±0.55 | 73.72±1.91 | 80.43±2.21 |
| Ours | 87.49±1.76 | 81.03±3.39 | 83.30±1.14 | 76.04±0.26 | 81.97±1.64 |

**Table 16: 5-bit Quantization Results on PACS Test Set (test)**

| Method | Art | Cartoon | Photo | Sketch | Avg |
|---|---|---|---|---|---|
| LSQ | 43.29±1.01 | 58.50±3.33 | 67.37±0.97 | 60.37±0.00 | 57.38±1.33 |
| LSQ + SAGM | 47.44±4.42 | 59.86±4.10 | 64.60±1.50 | 62.88±0.10 | 58.69±2.53 |
| Ours | 49.60±3.36 | 61.51±4.05 | 66.65±1.83 | 63.87±0.48 | 60.41±2.43 |

**Table 17: 8-bit EWGS Results on DomainNet Validation Set (val)**

| Method | Clipart | Infograph | Painting | Quickdraw | Real | Sketch |
|---|---|---|---|---|---|---|
| EWGS | 66.20±0.28 | 66.20±0.28 | 59.79±0.16 | 59.79±0.16 | 57.12±0.06 | 57.12±0.06 |
| EWGS + SAGM | 67.87±0.21 | 67.87±0.21 | 62.17±0.29 | 62.17±0.29 | 59.01±0.40 | 59.01±0.40 |
| Ours | 67.82±0.21 | 67.82±0.21 | 62.31±0.28 | 62.31±0.28 | 59.20±0.04 | 59.20±0.04 |

**Table 18: 8-bit EWGS Results on DomainNet Test Set (test)**

| Method | Clipart | Infograph | Painting | Quickdraw | Real | Sketch |
|---|---|---|---|---|---|---|
| EWGS | 60.22±0.07 | 16.17±0.21 | 45.84±0.18 | 13.94±0.01 | 51.84±0.08 | 47.40±0.07 |
| EWGS + SAGM | 62.09±0.07 | 16.76±0.12 | 48.16±0.31 | 15.51±0.18 | 53.89±0.27 | 49.13±0.33 |
| Ours | 61.99±0.09 | 16.75±0.01 | 48.35±0.17 | 15.51±0.12 | 54.12±0.31 | 49.32±0.24 |

**Table 19: 5-bit EWGS Results on DomainNet Validation Set (val)**

| Method | Clipart | Infograph | Painting | Quickdraw | Real | Sketch |
|---|---|---|---|---|---|---|
| EWGS | 66.32±0.03 | 66.32±0.03 | 60.20±0.12 | 60.20±0.12 | 57.11±0.32 | 57.11±0.32 |
| EWGS + SAGM | 66.98±0.95 | 66.98±0.95 | 60.42±0.23 | 60.42±0.23 | 59.30±0.22 | 59.30±0.22 |
| Ours | 67.32±1.04 | 67.32±1.04 | 61.18±0.16 | 61.18±0.16 | 59.93±0.17 | 59.93±0.17 |

**Table 20: 5-bit EWGS Results on DomainNet Test Set (test)**

| Method | Clipart | Infograph | Painting | Quickdraw | Real | Sketch |
|---|---|---|---|---|---|---|
| EWGS | 60.20±0.12 | 15.73±0.01 | 46.15±0.32 | 14.14±0.22 | 51.02±0.54 | 46.88±0.25 |
| EWGS + SAGM | 61.15±0.51 | 16.10±0.39 | 46.10±0.32 | 15.61±0.36 | 54.12±0.23 | 49.55±0.23 |
| Ours | 61.50±0.49 | 16.13±0.47 | 46.75±0.26 | 15.85±0.08 | 55.04±0.14 | 49.98±0.14 |

**Table 21: 4-bit EWGS Results on DomainNet Validation Set (val)**

| Method | Clipart | Infograph | Painting | Quickdraw | Real | Sketch |
|---|---|---|---|---|---|---|
| EWGS | 63.74±1.45 | 63.74±1.45 | 60.00±0.27 | 60.00±0.27 | 57.29±0.62 | 57.29±0.62 |
| EWGS + SAGM | 64.44±0.19 | 64.44±0.19 | 60.12±0.22 | 60.12±0.22 | 58.79±0.49 | 58.79±0.49 |
| Ours | 63.96±0.35 | 63.96±0.35 | 61.01±0.80 | 61.01±0.80 | 59.27±0.36 | 59.27±0.36 |

**Table 22: 4-bit EWGS Results on DomainNet Test Set (test)**

| Method | Clipart | Infograph | Painting | Quickdraw | Real | Sketch |
|---|---|---|---|---|---|---|
| EWGS | 58.91±0.75 | 14.63±0.49 | 45.62±0.20 | 14.46±0.20 | 51.64±0.88 | 47.15±0.79 |
| EWGS + SAGM | 59.82±0.11 | 15.71±0.21 | 45.63±0.40 | 16.04±0.01 | 53.52±0.50 | 49.18±0.40 |
| Ours | 59.04±0.21 | 15.35±0.15 | 46.18±0.77 | 16.32±0.12 | 54.48±0.41 | 50.02±0.37 |

**Table 23: 3-bit EWGS Results on OfficeHome Validation Set (val)**

| Method | Art | Clipart | Product | Real-World |
|---|---|---|---|---|
| EWGS | 66.48±0.97 | 58.33±3.74 | 56.10±1.88 | 60.36±3.03 |
| EWGS + SAGM | 62.83±3.25 | 58.81±0.05 | 54.57±1.18 | 64.27±0.37 |
| Ours | 68.17±1.79 | 60.63±0.31 | 56.09±2.34 | 63.64±1.15 |

**Table 24: 3-bit EWGS Results on OfficeHome Test Set (test)**

| Method | Art | Clipart | Product | Real-World |
|---|---|---|---|---|
| EWGS | 26.98±0.51 | 36.88±3.95 | 48.00±3.38 | 44.48±3.86 |
| EWGS + SAGM | 24.33±3.84 | 40.08±0.01 | 44.00±2.17 | 48.08±0.03 |
| Ours | 31.57±2.52 | 41.60±0.84 | 47.02±3.27 | 47.78±1.19 |

**Table 25: 3-bit LSQ Results on OfficeHome Validation Set (val)**

| Method | Art | Clipart | Product | Real-World |
|---|---|---|---|---|
| LSQ | 66.84±3.48 | 53.22±0.11 | 56.37±0.78 | 61.90±3.03 |
| LSQ + SAGM | 68.72±3.98 | 58.73±1.89 | 51.05±0.75 | 64.27±0.37 |
| Ours | 67.36±2.53 | 61.88±0.49 | 53.71±1.30 | 65.53±0.13 |

**Table 26: 3-bit LSQ Results on OfficeHome Test Set (test)**

| Method | Art | Clipart | Product | Real-World |
|---|---|---|---|---|
| LSQ | 27.16±3.48 | 33.71±0.20 | 46.44±1.31 | 45.28±2.42 |
| LSQ + SAGM | 30.74±4.74 | 40.48±1.39 | 40.61±0.72 | 48.08±0.03 |
| Ours | 29.43±2.60 | 42.37±0.56 | 44.16±0.49 | 49.24±1.25 |