




Symmetry Regularization and Saturating Nonlinearity for Robust Quantization

Sein Park^{1*} , Yeongsang Jang^{2*} , and Eunhyeok Park^{1,2} 

* equal contribution

¹ Graduate School of Artificial Intelligence,

² Department of Computer Science and Engineering

POSTECH, Pohang, Korea

{seinpark, jangys, eh.park}@postech.ac.kr

Abstract. Robust quantization improves the tolerance of networks for various implementations, allowing reliable output in different bit-widths or fragmented low-precision arithmetic. In this work, we perform extensive analyses to identify the sources of quantization error and present three insights to robustify a network against quantization: reduction of error propagation, range clamping for error minimization, and inherited robustness against quantization. Based on these insights, we propose two novel methods called symmetry regularization (SymReg) and saturating nonlinearity (SatNL). Applying the proposed methods during training can enhance the robustness of arbitrary neural networks against quantization on existing post-training quantization (PTQ) and quantization-aware training (QAT) algorithms and enables us to obtain a single weight flexible enough to maintain the output quality under various conditions. We conduct extensive studies on CIFAR and ImageNet datasets and validate the effectiveness of the proposed methods.

Keywords: Robust Quantization, Post-training Quantization (PTQ), Quantization-aware Training (QAT)

1 Introduction

Deep learning algorithms have shown excellence in diverse applications, but the increasing memory footprint and computation overhead have become obstacles to utilizing them. To exploit the excellence of deep neural networks (DNNs) in practice, neural network optimization is becoming more and more important. Neural network quantization is a representative optimization technique beneficial to footprint reduction and performance improvement. Due to its practical advantages, advanced hardware is already equipped with low-precision support, such as the well-known float16, bfloat16, and int8-based operations [42,28,19,40,36,37], even with 4-bit or lower-precision acceleration [29,33,39,27,18]. With the aid of a judiciously designed quantization algorithm, we could enjoy the benefit of low-precision computation in reality.

However, quantization has the substantial limitation of accuracy degradation due to the limited representation capability. Many studies have been actively proposed to address this problem, and Quantization-aware training (QAT) is a representative approach where end-to-end training is applied to refine the simulated error of pseudo (or fake-) quantization [9,6,30,20,22,47]. QAT is advantageous in the sense of minimal accuracy degradation in the given bit-width. Recently, post-training quantization (PTQ) has emerged as an alternative approach that quantizes the pre-trained network without fine-tuning [2,3,46,10,22,26,25,5,16,41]. PTQ allows us to exploit the benefit of low-precision computation with a minimal number of training datasets, thereby having many more practical use cases compared to QAT.

Nonetheless, both QAT and PTQ have severe drawbacks: QAT requires access to the entire dataset and the expenses of an additional training stage. In addition, the model with QAT is specialized for the target precision and quantization scheme, thereby lacking robustness in different bit-widths. On the other hand, PTQ suffers from notable accuracy degradation to QAT due to the lower degree of freedom and insufficient information to compensate for the errors. Many studies have been proposed to overcome this limitation, but a bit-width of 8- or more is still required for the advanced networks [5,16,41].

Recently, alternative approaches have been proposed to enhance the robustness of networks against quantization [45,12,14,34]. Improving the robustness of networks has diverse advantages, including allowing the quantized model to maintain the accuracy in bit-widths other than in which it trained and preserving the quality of output with various quantization algorithms. Practically, these properties help to utilize the pareto-front optimal points of energy consumption and computation, such as exploiting a high-precision model when the resource (e.g., battery) is sufficient and reducing precision dynamically when the resource is scarce. In addition, numerous companies are now designing their own accelerators having divergent and fragmented low-precision implementations. When we need to support multiple accelerators, preparing a single low-precision model robust enough to endure the minor modification of different implementations could be an attractive option for rapid deployment. Robust quantization enables diverse appealing applications, having strong importance in practice.

In this work, we propose two novel methods to increase the robustness of neural networks based on three insights about the error component of quantization. The paper is organized as follows; first, we perform an extensive analyses to identify the source of errors from quantization and indicate three motivations to robustify the network: reduction of error propagation, range clamping for error minimization, and inherited robustness against quantization (Sec. 3). To address those motivations, we introduce two novel ideas: symmetry regularization (Sym-Reg) for the reduction of error propagation and saturating nonlinearity (SatNL) for the others (Sec. 4). According to our extensive experiments, the proposed methods are beneficial for maximizing the robustness of networks after QAT or PTQ, showing state-of-the-art results. (Sec. 5). We then clarify the limitation of this study in Sec. 6 and conclude the paper in Sec. 7.

2 Related Work

2.1 Quantization-aware Training and Post-training Quantization

QAT [9,6,30,20,22,47] shows the potential of low-precision computation, where the milestone networks (e.g., VGG[35], GoogleNet[38], and ResNet[13]) could be quantized into sub-4-bit without accuracy loss [9,6], and advanced light-weight networks (e.g., MobileNet-V2) could be quantized into 4-bit with negligible accuracy loss [30]. Meanwhile, PTQ [2,3,46,10,22,26,25,5,16] applies a conservative bit-width to maintain the quality of output, even though it offers performance benefits with relaxed constraints. In this work, we aim to maximize the benefits of QAT and PTQ through the advantages of robust quantization.

2.2 Robustness of Neural Networks

A line of work closely related to ours is the analysis of the robustness of neural networks, which has attracted attention recently. Currently, the Hessian-aware metric is often used to identify the robustness of neural networks. Previous studies pointed out that the second derivative of loss is a good approximation of network sensitivity [8,1], proposed a way to estimate the approximate Hessian metric efficiently, and showed the potential of sensitivity-aware quantization [8,7,44,23,41]. On the other hand, few studies have focused on easing the sensitivity of networks during the training phase to improve their generalization performance [11,21]. The proposed (adaptive) sharpness-aware minimization, (A)SAM, makes the network have lower Hessian spectra than the networks trained without it. It is expected to be beneficial for enhancing the endurance of the network for quantization, but we observe that the benefit of (A)SAM is degraded in the quantization domain. In this paper, we propose a novel idea, SatNL, to maximize the robustness of quantization with (A)SAM training.

Moreover, few studies have tried to minimize quantization error in the view of robustness. GDRQ [45] and BR [12] attempted to improve the robustness of networks via regularization in QAT tasks. Gradient l_1 -regularization [14] lowered the sensitivity of networks for quantization via regularizing l_1 -norm of gradient, and KURE [34] regularized the weights in a uniform distribution to have minimal accuracy drop after the quantization. The two studies [14,34] are the most relevant to ours in that sharing the same objective of robust quantization. However, according to our observation, the former becomes unstable in advanced networks (i.e., MobileNet-V2/V3) and the latter is orthogonal to ours. Our methods show results comparable to those of KURE, and we can maximize the endurance of the network by applying ours with KURE jointly, as will be shown in Sec. 5.

3 Motivation

In this section, we provide the analysis of the error sources induced by quantization and explain the motivations to enhance the robustness of the networks for each error source.

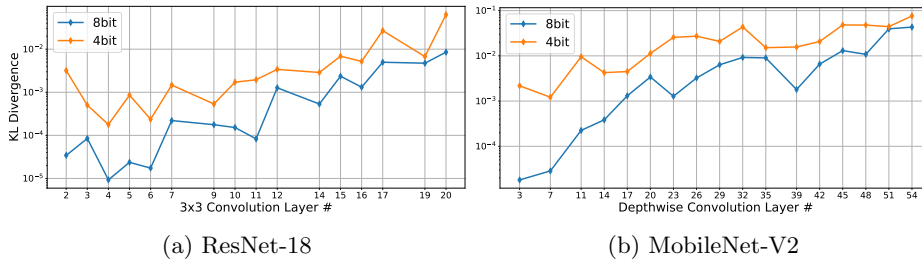


Fig. 1: KL divergence of 3x3 convolution and depthwise convolution output between original model and model quantized with PTQ [2] (a) ResNet-18 (b) MobileNet-V2.

3.1 Reduction of Error Propagation

Several previous studies have indicated that quantization introduces the distortion of statistics compared to the original distribution [24,10]. Moreover, when we apply quantization over the entire network, each layer introduces additional distortion to the output. As a result, the error continues to propagate and accumulate over the networks, as shown in Fig. 1, resulting in a large amount of accuracy degradation. Many studies related to PTQ have attempted to mitigate this problem by explicitly minimizing the difference in statistics before and after the quantization [10,4]. However, in this paper, we propose an alternative approach to minimize the difference in statistics on any quantization algorithms.

To achieve this goal, we focus on minimizing the biased quantization error problem [25,4]. Consider the linear or convolution operation $y = W \cdot x$, where W is an arbitrary fixed weight and x is an activation assumed i.i.d. variable with $E_i[x_i] = \mu_x$. In this condition, we can estimate the expected value of a single output unit $y_j = \sum_i^N W_{j,i} \cdot x_i$:

$$E_j[y_j] = N\mu_x E_i[W_{j,i}], \quad (1)$$

where N is the number of elements and i is the index of input x . When we apply the quantization to the weight, the expected value of output drifts due to the distortion of the weight as follows:

$$E_j[\tilde{y}_j - y_j] = N\mu_x \cdot \left(E_i[Q(W_{j,i})] - E_i[W_{j,i}] \right), \quad (2)$$

where y_j and \tilde{y}_j are the j -th output with the original weight and quantized weight, respectively. To minimize error propagation, we should minimize the difference of averaged weight in the output-channel dimension. However, satisfying Eq. (2) strictly for any quantization algorithm is highly challenging. To ease the difficulty of the objective, we adopt the additional condition as given by:

$$E_i[Q(W_{j,i})] = E_i[W_{j,i}] = 0. \quad (3)$$

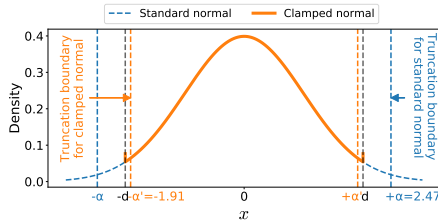


Fig. 2: Histogram of standard normal vs clamped normal ($d = \pm 2$), and their truncation boundaries.

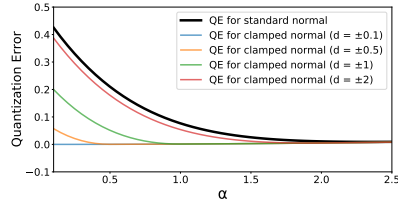


Fig. 3: Analysis of quantization error for standard normal and clamped normal.

By forcing the mean of the weights before and after the quantization toward 0, Eq. (2) is satisfied as a sufficient condition. When the full-precision weight is symmetric, the mean of the full-precision weight is zero. In addition, when we apply a symmetric quantization, which is commonly used for the weight quantization due to hardware compatibility [17,42,40], the drift in the positive values could be amortized by the drift in the negative values; thereby, the mean of the quantized weight is also zero.

In summary, if we can force the full-precision weight in a symmetric distribution, the statistics distortion and error propagation after the quantization could be minimized. Unlike the explicit bias correction process [25], weight symmetry inherits the robustness against bias drift, enabling us effortless transition to different quantization policies.

To guide the convergence of weight toward the symmetric distribution, we propose a novel regularization in Sec. 4.1. When we apply this regularization during pre-train the model, the difference in statistics before and after PTQ is reduced. Furthermore, the statistics distortion is also minimized when we utilize the fine-tuned weight after QAT in different bit-widths without an additional fine-tuning stage, helping to maintain the quality of output.

3.2 Range Clamping for Error Minimization

In linear quantization, the quantization levels are evenly distributed in between the truncation boundaries. Applying quantization to the full-precision tensor induces quantization error, which can be decomposed into the truncation error and the rounding error. The truncation and rounding errors are inevitable because of the limited number of quantization levels, but the difference of domain (i.e., the unbounded full-precision and the truncated quantization) enlarges the quantization error. A previous study [31] pointed out that the data with infrequent but large values require a widened truncation boundary, which increases the rounding error significantly. To mitigate the quantization error, the domain of full-precision data should be narrowed to a bounded range.

Motivated by this limitation, we propose a straightforward idea of introducing the range clamping to the full-precision weight¹, as shown in Fig. 2. Assume that the original weight follows a normal distribution whose PDF is $f(x) \sim N(0, 1)$, just as the convention of previous studies [2,34]. When we apply the b-bit symmetric quantization toward minimizing the L2 norm, the quantization error can be estimated as follows [2]:

$$\begin{aligned} \text{Quantization Error} &= E[(W - Q(W))^2] \\ &\approx \underbrace{2 \cdot \int_{\alpha}^{\infty} f(x) \cdot (x - \alpha)^2 dx}_{\text{truncation error}} + \underbrace{\frac{\alpha^2}{3 \cdot 2^{2b}}}_{\text{rounding error}}, \end{aligned} \quad (4)$$

where α is the truncation boundary that minimizes $\|W - Q(W)\|_2$. On the other hand, the clamping modifies the distribution of weight, whose cumulative distribution function $G(x)$ is expressed as

$$G(x; d) = \begin{cases} 0, & x \leq -d \\ F(x) + F(-|d|), & -d < x < d \\ 1, & d \leq x, \end{cases} \quad (5)$$

where d is the newly introduced clamping target and $F(x)$ is the cumulative distribution function of $f(x)$. Then, the quantization error of the clamped distribution is expressed as

$$\begin{aligned} \text{Quantization Error}' &= E[(W' - Q(W'))^2] \\ &\approx \underbrace{2 \cdot \left(F(-|d|) \cdot (d - \alpha')^2 + \int_{\alpha'}^d f(x) \cdot (x - \alpha')^2 dx \right)}_{\text{truncation error}} + \underbrace{\frac{\alpha'^2}{3 \cdot 2^{2b}}}_{\text{rounding error}}, \end{aligned} \quad (6)$$

where α' is the truncation boundary that minimizes $\|W' - Q(W')\|_2$.

When we compare the errors of Eq. (4) and Eq. (6), Eq. (6) always has a smaller error than Eq. (4), as shown in Fig. 3. The proofs of Eq. (4) to Eq. (6) are provided in the supplementary material.

This analysis indicates that the range clamping of full-precision data is beneficial for minimizing quantization error. Thereby, if we train a network with the range clamping nonlinearity, the network could have a strong endurance for quantization. In addition, the quantization error could be minimized in the different precision, resulting in low accuracy loss other than the bit-width we quantized. Indeed, a similar idea was addressed in MobileNet-V2 [32] in terms of ReLU6 for activation. The range clamping can be seen as an extension of the idea of ReLU6 for weight.

One possible drawback of this idea is the accuracy degradation due to the limited degree of freedom. For instance, when we set the clamping target close

¹ Please note that we intentionally use different expressions to distinguish quantization's truncation and the clamping of full-precision data.

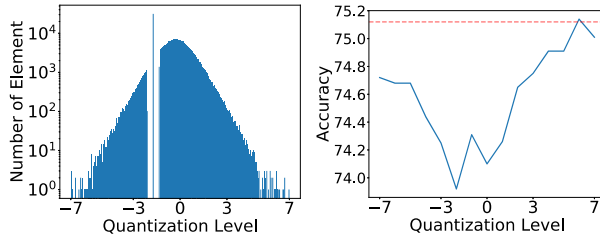


Fig. 4: Example of single-level quantization and corresponding result. The dashed line represents the baseline full-precision accuracy, and the solid line shows that of the quantized network.

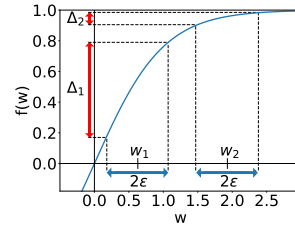


Fig. 5: Saturating non-linearity for weight and the adversarial boundary of ASAM.

to zero, the quantization error becomes negligible, but the training from scratch may fail or show poor accuracy. Therefore, it is necessary to find a sweet spot that reduces the quantization error by a large margin while maintaining the quality of output. Empirically, we determine the practical implementation of the truncation that has negligible quality degradation while minimizing quantization loss significantly and provide it in Sec. 4.2.

3.3 Inherited Robustness against Quantization

In addition to the range clamping for error minimization, we inspect another approach that enhances the inherited robustness of networks against quantization based on Hessian-aware loss sharpness minimization. Recently, several studies have focused on enhancing the generalization of neural networks based on the training pipeline aware of the sharpness of loss surface [11,21]. Those studies have aimed to guide the convergence of networks toward the flat minimum having smaller Hessian values, where the minor distortion of weight could be ignored without affecting the output. From a similar perspective, the Hessian of weight is utilized as an important metric for measuring the sensitivity of networks regarding quantization [8,7,23]. Motivated by these examples, we try to adopt the Hessian-aware training to enhance the robustness of networks for quantization by guiding the convergence of networks into a smooth loss surface. We utilize (A)SAM [11,21], and empirically observe that training with (A)SAM is beneficial for minimizing quantization error.

However, we also observe that there is room for improvement in terms of enhancing robustness for quantization with (A)SAM, because the quantization sensitivity differs depending on the value of quantization levels. Fig. 4 shows the accuracy degradation after applying single-level PTQ to the weights of MobileNet-V2 on CIFAR-100, where the weights corresponding to the specific quantization level are quantized while the rest remain in full-precision. As shown in the figure, the accuracy degradation increases when the smaller weights are quantized. This indicates that the weights near zero are more vulnerable to quantization than the weights having large values. We speculate that because the majority of weights

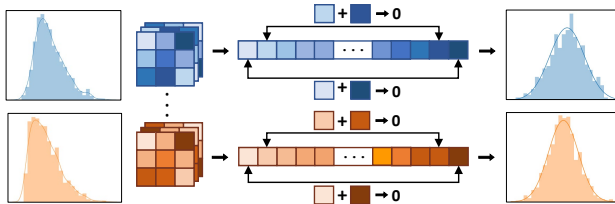


Fig. 6: Visualization of proposed SymReg.

are concentrated near zero, the accumulated error of quantization is inversely proportional to the magnitude of the value. Thereby, to maintain the quality of output with quantization, the sharpness minimization should be applied in different strengths depending on the magnitude of weight.

To realize the aforementioned objective, we propose a novel idea that introduces the saturating nonlinearity to the weight, as visualized in Fig. 5. When we apply (A)SAM with the nonlinearity, the robustness boundary of (A)SAM, which is equally assigned in the weight before the nonlinearity, covers a different range in output depending on the slope of the nonlinearity. SatNL has the largest slope near zero; as a result, the effect of the (A)SAM algorithm turns friendly to the quantization.

4 Implementation

Based on the motivations in the previous section, we propose the practical implementations, called symmetry regularization (SymReg) and saturating nonlinearity (SatNL).

4.1 Symmetry Regularization

In Sec. 3.1, we claimed that the symmetric weight could minimize the error propagation. To realize this, we propose an additional regularization called symmetry regularization (SymReg).

The weight symmetry is achievable when every weight has a corresponding mate with an identical magnitude but a different sign. After sorting the weights and assigning the index in ascending order, we can make pairs where one element is selected in ascending order \tilde{w}_i while the other is in descending order \tilde{w}_{N-i} . When we calculate the $L1$ norm of the sum of each pair, the expectation should become zero as follows:

$$\frac{2}{N} \sum_{n=1}^{\frac{N}{2}} |\tilde{w}_n + \tilde{w}_{N-n}| = 0. \quad (7)$$

Based on this intuition, we design a layer-wise SymReg that guides the convergence of weight into the symmetric distribution as shown in Fig. 6. The

SymReg is defined as:

$$\mathcal{L}_{sym1} = \frac{2}{C \cdot N} \sum_{c=1}^C \sum_{i=1}^{\frac{N}{2}} |w_i^c + w_{N-i}^c|, \quad (8)$$

where w_i^c represents the i -th smallest element in the c -th channel.

\mathcal{L}_{sym1} is applied when we train the full-precision network in addition to the conventional loss functions. However, if the restriction of \mathcal{L}_{sym1} is too severe, it could lead to a minor accuracy degradation. To minimize it, we propose the relaxed SymReg, which measures the symmetry in a 2:2 relation with more degree of freedom instead of a 1:1 relation. We empirically observe that more than 3:3 relation degrades the benefit of SymReg.

$$\mathcal{L}_{sym2} = \frac{4}{C \cdot N} \sum_{c=1}^C \sum_{i=1}^{\frac{N}{4}} |w_{2i}^c + w_{2i+1}^c + w_{N-2i}^c + w_{N-2i-1}^c|. \quad (9)$$

In the experiment, we combine \mathcal{L}_{sym1} and \mathcal{L}_{sym2} adequately to minimize accuracy degradation while exploiting the benefit of propagation error minimization. The overall loss is expressed as follows:

$$\mathcal{L} = \mathcal{L}_{CE} + \lambda_1 \cdot \mathcal{L}_{sym1} + \lambda_2 \cdot \mathcal{L}_{sym2}. \quad (10)$$

4.2 Saturating Nonlinearity (SatNL)

In Sec. 3.2, we showed that the truncation of full-precision weight could be highly beneficial for minimizing quantization error. In addition, in Sec. 3.3, we showed the empirical analysis indicating that the quantization sensitivity differs depending on the magnitude of the weight. To resolve these two problems simultaneously, we propose applying a specialized nonlinearity function f on top of the weight as follows: $conv(W, X) \rightarrow conv(f(W), X)$.

The nonlinearity should satisfy three properties. 1. It needs to be an odd function. Because we assume that the weight is quantized by the symmetric quantization, it would be better to have an identical range of the negative and positive regions to maximize the quantization level efficiency. 2. The range of output needs to be bounded. To satisfy the criteria of Sec. 3.2, the weight after the nonlinearity should be narrowed to a certain range. 3. The slope is gradually decreased as the input value is increased. To maximize the benefit of (A)SAM, the nonlinearity should be saturated as the value is increased. Empirically, we choose the hyper-tangent (tanh) function as nonlinearity which satisfies all three conditions. Note that the normalized tanh was used in QAT studies [47], but there is a fundamental difference in terms of the intention of using it. In this study, we exploit the tanh function intentionally to maximize the robustness against quantization and turn (A)SAM algorithm friendly to quantization. Other nonlinearity functions could be applicable when the three conditions are met. According to our experiments, the impact on the final accuracy is negligible regardless of the nonlinearity functions, while the desired properties for robustness are valid. Additional analysis is in the supplementary material.

5 Experiments

To show the superiority of the proposed methods, we conduct extensive studies on CIFAR-100 and ImageNet datasets with representative networks (i.e., ResNet-18 [13], MobileNet-V2 [32], and MobileNet-V3 [15]). We use layer-wise asymmetric quantization for activation and output-channel-wise symmetric quantization for weight to enable acceleration on existing hardwares [17,43,40,42].

For QAT, we apply LSQ [9], which is the advanced differentiable quantization scheme that allows the quantization of ResNet-18 into 3-bit without accuracy loss in the ImageNet. For PTQ, we apply the ACIQ [2], AdaQuant [16], and QDrop [41] algorithms. ACIQ is a well-known PTQ that analytically finds the optimal quantization boundary, and QDrop is a state-of-the-art PTQ algorithm with small calibration sets. By utilizing multiple quantization algorithms with different properties, we aim to validate the universal robustness.

Moreover, we use a common practice that fixes the bit-width of the first and the last layers into 8-bit. All of the other layers are quantized into the given bit-width identically unless explicitly specified otherwise. SymReg is not applied to the depthwise convolution and the linear layer in MobileNet-V2/V3. In 3×3 depthwise convolution, SymReg degrades the expression capability significantly by forcing one out of nine elements to become zero. In ASAM, ρ is fixed as 1 except 0.2 for MobileNet-V3 because ASAM with high ρ becomes unstable. The hyper-parameters of SymReg λ_1/λ_2 are empirically set as 0.1/0.1 for ImageNet and CIFAR-100. The details of the training parameters (i.e., learning rate, decay, epochs, etc.), are provided in the supplementary material.

5.1 Robustness of Bit-Precision for PTQ

Table 1 shows the accuracy degradation after PTQ of ResNet-18, and MobileNet-V2/V3 on the ImageNet dataset. Our/baseline models are trained from scratch with/without the proposed ideas and then quantized into low-precision with PTQ algorithms. Our experiments show that the proposed methods are beneficial for minimizing the accuracy degradation after PTQ in every point regardless of the PTQ details. When we combine ours with the advanced PTQ (i.e., QDrop), we can quantize ResNet-18 into 4-bit with accuracy loss of less than 1 %. Compared to the baseline, we reduce 1.02 % of the top-1 accuracy degradation. In addition, the synergy of the advanced PTQ algorithm and our methods enables sub-8-bit quantization for the advanced networks with minimal accuracy degradation. QDrop with ours achieves 69.87 % in 4-bit MobileNet-V2, which is the highest accuracy after 4-bit PTQ to the best of our knowledge. Notably, when we select the layer-wise bit-width depending on the sensitivity of layers, where the depthwise and squeeze-excitation layers are quantized into 8-bit and the rest of the layers are quantized into 4-bit with AdaQuant, MobileNet-V2/V3 shows 2.96 %/4.34 % of accuracy degradation respectively². In this configuration, we can enjoy the benefit of 4-bit computation in 85.3 % of computation in the case

² Note that the column of mixed-precision results is omitted in Table 1 for brevity.

Table 1: Results of applying PTQ to baseline and network with proposed ideas on the ImageNet dataset. The values in the table represent the top-1 accuracy. The dashed cells represent the points where the PTQ fails to converge.

Model	PTQ	Method	Weight/activation bit-width configuration							
			FP	4/FP	3/FP	2/FP	6/6	5/5	4/4	3/3
ResNet-18	ACIQ[2]	Baseline	70.54	47.44	-	-	68.70	64.87	38.46	-
		Ours	70.92	69.22	49.06	-	70.02	68.99	66.65	42.95
	AdaQuant[16]	Baseline	70.54	69.29	66.18	3.23	70.17	69.55	67.67	57.57
		Ours	70.92	70.36	68.84	48.39	70.75	70.37	69.35	64.04
	QDrop[41]	Baseline	70.54	70.15	69.39	66.40	70.27	69.93	68.91	65.75
		Ours	70.92	70.69	70.06	66.95	70.81	70.57	69.93	67.45
MobileNet-V2	ACIQ[2]	Baseline	72.22	28.68	-	-	69.30	64.20	18.15	-
		Ours	72.87	70.07	40.79	-	71.07	68.66	58.30	6.25
	AdaQuant[16]	Baseline	72.22	70.67	59.80	-	71.52	70.72	63.70	-
		Ours	72.87	72.23	69.03	-	72.27	71.76	68.91	18.36
	QDrop[41]	Baseline	72.22	71.41	68.32	48.68	71.57	70.64	67.08	50.79
		Ours	72.87	72.44	71.18	61.68	72.61	72.05	69.87	62.55
MobileNet-V3	ACIQ[2]	Baseline	74.52	29.65	-	-	-	-	-	-
		Ours	74.43	61.95	1.04	-	-	-	-	-
	AdaQuant[16]	Baseline	74.52	72.92	64.17	-	72.73	68.95	43.88	-
		Ours	74.43	73.51	70.50	2.87	72.69	71.02	62.73	-

of MobileNet-V2. Without ours, the accuracy degradation is 7.25 % and 9.44 % respectively in mixed precision, which is too poor to be used in real applications. The combination of network robustness enhancement and sensitivity-aware quantization could be a good candidate for practical deployment. According to our experiment, SymReg and SatNL extended the training time by 2.53 % when we train MobileNet-V2 with ASAM on the ImageNet dataset. By spending this one-time overhead, a robust network that can minimize accuracy degradation regardless of the PTQ scheme can be achieved. Additional experiments compared with KURE are supported in supplementary materials.

5.2 Ablation Study

Table 2 shows the effect of the proposed methods based on the accuracy degradation with PTQ. When we add an additional component of the proposed methods progressively (+ SymReg, + SatNL, + All), the PTQ error is gradually reduced, showing that the proposed methods enhance the robustness of the network for PTQ in diverse aspects. SymReg and SatNL are beneficial for robustness but introduces a slight accuracy degradation in full-precision. However, the accuracy degradation could be mitigated with the assistance of ASAM. When we compare + ASAM and + SatNL + ASAM, the latter shows higher accuracy in full-precision and better robustness in lower precision, showing that the benefit of ASAM is boosted with SatNL.

When we compare the performance of the proposed methods with KURE, our methods lower accuracy degradation in the given bit-width than that of

Table 2: Results of ablation study of proposed methods on MobileNet-V2 at CIFAR-100 dataset. The weight is quantized into the given bit-width with ACIQ. "All" means every method (+ SymReg + SatNL + ASAM). The values in the table represent the mean and standard deviation of top-1 accuracy with eight trials (10 trials except min. and max. results). FP means full-precision.

	FP	4-bit	3-bit	2-bit
Baseline	74.70 \pm 0.16	73.32 \pm 0.32	66.81 \pm 0.86	6.92 \pm 2.68
+ SymReg	74.66 \pm 0.16	73.49 \pm 0.30	69.69 \pm 0.81	25.01 \pm 4.80
+ SatNL	74.80 \pm 0.11	73.43 \pm 0.36	68.65 \pm 1.23	14.45 \pm 2.93
+ SymReg + SatNL	74.41 \pm 0.10	73.34 \pm 0.31	69.66 \pm 0.62	33.68 \pm 5.87
+ ASAM	75.52 \pm 0.19	74.42 \pm 0.16	69.71 \pm 0.71	11.80 \pm 3.96
+ SatNL + ASAM	75.55 \pm 0.21	74.53 \pm 0.21	70.78 \pm 0.81	25.77 \pm 4.03
+ All	75.33 \pm 0.16	74.55 \pm 0.27	72.17 \pm 0.22	39.27 \pm 2.56
+ KURE	74.97 \pm 0.11	74.22 \pm 0.10	71.26 \pm 0.47	34.90 \pm 4.51
+ KURE + ASAM	75.57 \pm 0.18	74.96 \pm 0.11	72.48 \pm 0.46	42.73 \pm 4.71
+ KURE + All	75.41 \pm 0.3	74.91 \pm 0.29	73.34 \pm 0.25	37.58 \pm 2.99

KURE. Meanwhile, when we apply ASAM with KURE, the robustness becomes comparable to ours. Moreover, the implementation details of KURE and ours are orthogonal. Thus, we can maximize the robustness by combining KURE and ours (KURE + ALL). This could be a state-of-the-art method for enhancing the robustness of networks, as far as we know.

5.3 Robustness for Quantization Step Size

To validate the effect of the proposed methods, we measure the accuracy changes depending on the step size as shown in Fig. 7. In all cases, including PTQ and QAT, the proposed methods maintain the quality of output in the various step sizes. In the case of MobileNet-V2 for the ImageNet dataset, ours maintains 25.72 % higher accuracy compared to the baseline when the step size is changed by 8 %. In addition, ours shows comparable or better results to the previous best method, KURE, for the optimized network (i.e., MobileNet-V2). As indicated in the previous studies [34,17], the degree of freedom for the quantization step size could be restricted depending on the hardware implementation. For instance, some hardware supports step sizes having predefined values or limited resolution. In such a case, the robustness of the quantization step size is essential to maintain output quality after PTQ or QAT. Because our methods improve the robustness of the step size by a large margin, we expect that our methods could be helpful for the deployment of quantized networks in practice.

5.4 Robustness of Bit Precision for QAT

Because most of the existing QAT methods specialize the weight fine-tuning for the specific configuration, the accuracy of the quantized network in different bit-widths is reduced significantly. Meanwhile, ours enhance the robustness of

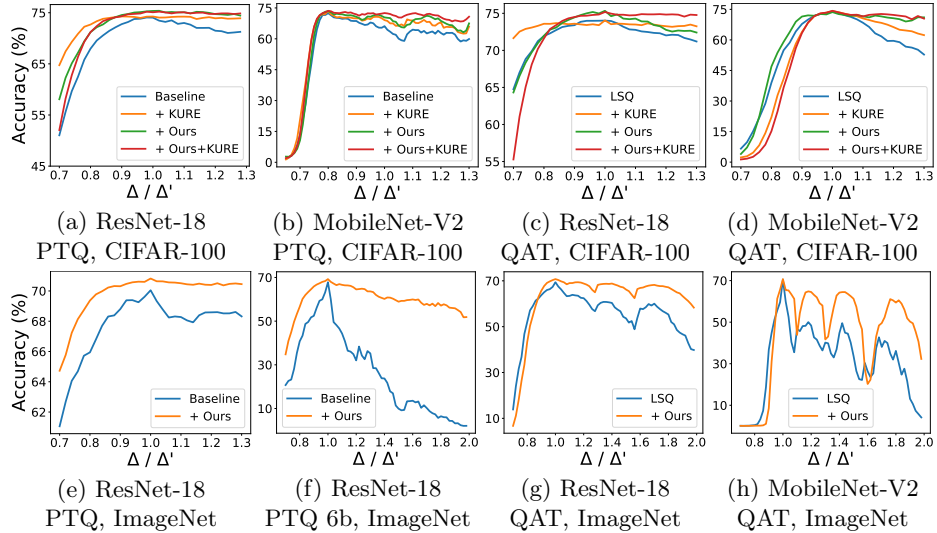


Fig. 7: Robustness of quantized network when we change step size of quantization operator for weight. The networks are optimized for the step size Δ' , and the accuracy is measured with the scaled step size Δ . All networks are quantized into 4-bit (except for (f) into 6-bit) with PTQ [16] and QAT [9], including activation and weight. Additional results are included in the supplementary material.

the quantized network, allowing stable accuracy in different bit-widths without fine-tuning. Fig. 8. shows the accuracy degradation depending on the operation bit-widths other than the one we trained via QAT. With the proposed methods, one can train a generic model that can produce reliable output in multiple bit-widths with existing QAT algorithms.

Table 3 shows the accuracy of QAT for ResNet-18/MobileNet-V2 in different bit-widths with/without the proposed methods. As shown in the table, the proposed fine-tuning scheme does not reduce the accuracy of the quantized network in all cases of ResNet-18 and 4-bit MobileNet-V2. In the case of 3-/2-bit MobileNet-V2, we speculated that the accuracy is slightly degraded due to the overlapped effect of strong regularization and the limited degree of freedom.

Table 3: Effect of our methods for top-1 accuracy of ResNet-18 and MobileNet-V2 quantized by LSQ [9] on the ImageNet dataset (90 epochs of fine-tuning).

Model	FP	Fine-tuning	4/4	3/3	2/2
ResNet-18	70.542	LSQ	69.39	68.80	66.26
		LSQ + Ours	70.74	69.73	66.58
MobileNet-V2	72.24	LSQ	70.46	67.51	44.87
		LSQ + Ours	71.16	66.93	43.41

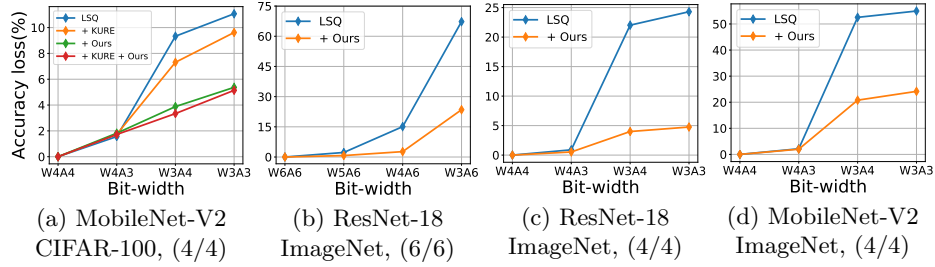


Fig. 8: Robustness of QAT model with and without proposed methods. (W/A) represents the initial bit-width of activation and weight, WXAY indicates the changed bit-width. The accuracy is measured without additional fine-tuning.

When applying quantization with 4-bit or higher precision, we can enjoy the benefit of robustness based on the proposed methods without losing accuracy.

6 Discussion

Our proposed methods are advantageous in minimizing quantization error after PTQ and QAT. However, one remaining important topic that we could not address in this paper is the robustness of the activation. In this paper, we rely on the PTQ/QAT algorithms for activation quantization. Unlike weight, activation is input-dependent, and the distribution is diverse depending on the behavior of nonlinear functions. This is a challenging problem and left as future work.

7 Conclusion

Enhancing the robustness of neural networks for quantization maximizes the benefit we can get from the low-precision operations. In this study, we reported three important motivations for minimizing the accuracy degradation after quantization: reduction of error propagation, range clamping for error minimization, and inherited robustness against quantization. Based on these insights, we proposed two novel ideas, symmetry regularization (SymReg) and saturating nonlinearity (SatNL). Our extensive experiments verified the advantages of the proposed methods, which significantly reduce the quantization error of diverse QAT and PTQ algorithms. Enhancing the robustness of quantization is achievable with negligible extra cost, but it enables us to exploit the benefit of low-precision computation with minimal accuracy degradation. We expect that the robustness of networks will minimize the deployment overhead for energy-efficient NPUs, thereby positively affecting the environment.

Acknowledgements. This work was supported by IITP grant funded by the Korea government (MSIT, No.2019-0-01906, No.2021-0-00105, and No.2021-0-00310), SK Hynix Inc. and Google Asia Pacific.

References

1. Alizadeh, M., Behboodi, A., van Baalen, M., Louizos, C., Blankevoort, T., Welling, M.: Gradient l1 regularization for quantization robustness. In: 8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020. OpenReview.net (2020), <https://openreview.net/forum?id=ryxK0JBtPr>
2. Banner, R., Nahshan, Y., Hoffer, E., Soudry, D.: ACIQ: analytical clipping for integer quantization of neural networks. CoRR **abs/1810.05723** (2018), <http://arxiv.org/abs/1810.05723>
3. Banner, R., Nahshan, Y., Soudry, D.: Post training 4-bit quantization of convolutional networks for rapid-deployment. In: NeurIPS (2019)
4. Brock, A., De, S., Smith, S.L.: Characterizing signal propagation to close the performance gap in unnormalized resnets. In: International Conference on Learning Representations (2021), <https://openreview.net/forum?id=IX3Nnir2omJ>
5. Cai, Y., Yao, Z., Dong, Z., Gholami, A., Mahoney, M.W., Keutzer, K.: Zeroq: A novel zero shot quantization framework. 2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) pp. 13166–13175 (2020)
6. Choi, J., Wang, Z., Venkataramani, S., Chuang, P.I.J., Srinivasan, V., Gopalakrishnan, K.: PACT: Parameterized clipping activation for quantized neural networks (2018), <https://openreview.net/forum?id=By5ugjyCb>
7. Dong, Z., Yao, Z., Arfeen, D., Gholami, A., Mahoney, M.W., Keutzer, K.: Hawq-v2: Hessian aware trace-weighted quantization of neural networks. In: Larochelle, H., Ranzato, M., Hadsell, R., Balcan, M.F., Lin, H. (eds.) Advances in Neural Information Processing Systems. vol. 33, pp. 18518–18529. Curran Associates, Inc. (2020), <https://proceedings.neurips.cc/paper/2020/file/d77c703536718b95308130ff2e5cf9ee-Paper.pdf>
8. Dong, Z., Yao, Z., Gholami, A., Mahoney, M.W., Keutzer, K.: Hawq: Hessian aware quantization of neural networks with mixed-precision. 2019 IEEE/CVF International Conference on Computer Vision (ICCV) pp. 293–302 (2019)
9. Esser, S.K., McKinstry, J.L., Bablani, D., Appuswamy, R., Modha, D.S.: Learned step size quantization. In: International Conference on Learning Representations (2020), <https://openreview.net/forum?id=rkg066VKDS>
10. Finkelstein, A., Almog, U., Grobman, M.: Fighting quantization bias with bias. arXiv preprint arXiv:1906.03193 (2019)
11. Foret, P., Kleiner, A., Mobahi, H., Neyshabur, B.: Sharpness-aware minimization for efficiently improving generalization. In: International Conference on Learning Representations (2021), <https://openreview.net/forum?id=6Tm1mposlrM>
12. Han, T., Li, D., Liu, J., Tian, L., Shan, Y.: Improving low-precision network quantization via bin regularization. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 5261–5270 (2021)
13. He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. In: 2016 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2016, Las Vegas, NV, USA, June 27-30, 2016. pp. 770–778. IEEE Computer Society (2016). <https://doi.org/10.1109/CVPR.2016.90>, <https://doi.org/10.1109/CVPR.2016.90>
14. Hoffman, J., Roberts, D.A., Yaida, S.: Robust learning with jacobian regularization (2020), <https://openreview.net/forum?id=ryl-RTEYvB>
15. Howard, A., Sandler, M., Chu, G., Chen, L.C., Chen, B., Tan, M., Wang, W., Zhu, Y., Pang, R., Vasudevan, V., et al.: Searching for mobilenetv3. In: Proceedings

- of the IEEE/CVF International Conference on Computer Vision. pp. 1314–1324 (2019)
16. Hubara, I., Nahshan, Y., Hanani, Y., Banner, R., Soudry, D.: Accurate post training quantization with small calibration sets. In: Meila, M., Zhang, T. (eds.) Proceedings of the 38th International Conference on Machine Learning, ICML 2021, 18–24 July 2021, Virtual Event. Proceedings of Machine Learning Research, vol. 139, pp. 4466–4475. PMLR (2021), <http://proceedings.mlr.press/v139/hubara21a.html>
 17. Jacob, B., Kligys, S., Chen, B., Zhu, M., Tang, M., Howard, A., Adam, H., Kalenichenko, D.: 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)
 18. Jang, J., Lee, S., Kim, D., Park, H., Ardestani, A.S., Choi, Y., Kim, C., Kim, Y., Yu, H., Abdel-Aziz, H., Park, J., Lee, H., Lee, D., Kim, M.W., Jung, H., Nam, H., Lim, D., Lee, S., Song, J., Kwon, S., Hassoun, J., Lim, S., Choi, C.: Sparsity-aware and re-configurable NPU architecture for samsung flagship mobile soc. In: 48th ACM/IEEE Annual International Symposium on Computer Architecture, ISCA 2021, Valencia, Spain, June 14–18, 2021. pp. 15–28. IEEE (2021). <https://doi.org/10.1109/ISCA52012.2021.00011>, <https://doi.org/10.1109/ISCA52012.2021.00011>
 19. Jouppi, N.P., Yoon, D.H., Ashcraft, M., Gottscho, M., Jablin, T.B., Kurian, G., Laudon, J., Li, S., Ma, P.C., Ma, X., Norrie, T., Patil, N., Prasad, S., Young, C., Zhou, Z., Patterson, D.A.: Ten lessons from three generations shaped google’s tpuv4i : Industrial product. In: 48th ACM/IEEE Annual International Symposium on Computer Architecture, ISCA 2021, Valencia, Spain, June 14–18, 2021. pp. 1–14. IEEE (2021). <https://doi.org/10.1109/ISCA52012.2021.00010>, <https://doi.org/10.1109/ISCA52012.2021.00010>
 20. Jung, S.H., Son, C., Lee, S., Son, J., Han, J.J., Kwak, Y., Hwang, S.J., Choi, C.: Learning to quantize deep networks by optimizing quantization intervals with task loss. 2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) pp. 4345–4354 (2019)
 21. Kwon, J., Kim, J., Park, H., Choi, I.K.: Asam: Adaptive sharpness-aware minimization for scale-invariant learning of deep neural networks. In: Meila, M., Zhang, T. (eds.) Proceedings of the 38th International Conference on Machine Learning. Proceedings of Machine Learning Research, vol. 139, pp. 5905–5914. PMLR (18–24 Jul 2021), <https://proceedings.mlr.press/v139/kwon21b.html>
 22. Lee, J.H., Ha, S., Choi, S., Lee, W.J., Lee, S.: Quantization for rapid deployment of deep neural networks (2019), <https://openreview.net/forum?id=HkzZBi0cFQ>
 23. Li, Y., Gong, R., Tan, X., Yang, Y., Hu, P., Zhang, Q., Yu, F., Wang, W., Gu, S.: {BRECQ}: Pushing the limit of post-training quantization by block reconstruction. In: International Conference on Learning Representations (2021), <https://openreview.net/forum?id=P0Wv6hDd9XH>
 24. Lin, J., Gan, C., Han, S.: Defensive quantization: When efficiency meets robustness. In: International Conference on Learning Representations (2019), <https://openreview.net/forum?id=ryetZ20ctX>
 25. Nagel, M., Baalen, M.v., Blankevoort, T., Welling, M.: Data-free quantization through weight equalization and bias correction. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 1325–1334 (2019)
 26. Nahshan, Y., Chmiel, B., Baskin, C., Zheltonozhskii, E., Banner, R., Bronstein, A.M., Mendelson, A.: Loss aware post-training quantization. ArXiv [abs/1911.07190](https://arxiv.org/abs/1911.07190) (2021)

27. Int4 precision for ai inference. <https://devblogs.nvidia.com/int4-for-ai-inference/> (2019), accessed: 2021-11-16
28. Nvidia a100 tensor core gpu architecture. <https://images.nvidia.com/aem-dam/en-zz/Solutions/data-center/nvidia-ampere-architecture-whitepaper.pdf> (2020), accessed: 2021-11-16
29. Park, E., Kim, D., Yoo, S.: Energy-efficient neural network accelerator based on outlier-aware low-precision computation. International Symposium on Computer Architecture (ISCA) (2018)
30. Park, E., Yoo, S.: Profit: A novel training method for sub-4-bit mobilenet models. In: European Conference on Computer Vision. pp. 430–446. Springer (2020)
31. Park, E., Yoo, S., Vajda, P.: Value-aware quantization for training and inference of neural networks. In: Ferrari, V., Hebert, M., Sminchisescu, C., Weiss, Y. (eds.) Computer Vision - ECCV 2018 - 15th European Conference, Munich, Germany, September 8-14, 2018, Proceedings, Part IV. Lecture Notes in Computer Science, vol. 11208, pp. 608–624. Springer (2018). https://doi.org/10.1007/978-3-030-01225-0_36, https://doi.org/10.1007/978-3-030-01225-0_36
32. Sandler, M., Howard, A., Zhu, M., Zhmoginov, A., Chen, L.C.: Mobilenetv2: Inverted residuals and linear bottlenecks. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 4510–4520 (2018)
33. Sharma, H., Park, J., Suda, N., Lai, L., Chau, B., Kim, J.K., Chandra, V., Esmaeilzadeh, H.: Bit fusion: Bit-level dynamically composable architecture for accelerating deep neural networks. International Symposium on Computer Architecture (ISCA) (2018)
34. Shkolnik, M., Chmiel, B., Banner, R., Shomron, G., Nahshan, Y., Bronstein, A.M., Weiser, U.C.: Robust quantization: One model to rule them all. In: Larochelle, H., Ranzato, M., Hadsell, R., Balcan, M.F., Lin, H.T. (eds.) Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual (2020), <https://proceedings.neurips.cc/paper/2020/hash/3948ead63a9f2944218de038d8934305-Abstract.html>
35. Simonyan, K., Zisserman, A.: Very deep convolutional networks for large-scale image recognition. In: International Conference on Learning Representations (2015)
36. Snapdragon neural processing engine sdk. <https://developer.qualcomm.com/docs/snpe/index.html> (2017), accessed: 2021-11-16
37. Song, J., Cho, Y., Park, J.S., Jang, J.W., Lee, S., Song, J.H., Lee, J.G., Kang, I.: 7.1 an 11.5 tops/w 1024-mac butterfly structure dual-core sparsity-aware neural processing unit in 8nm flagship mobile soc. International Solid-State Circuits Conference (ISSCC) (2019)
38. Szegedy, C., Liu, W., Jia, Y., Sermanet, P., Reed, S., Anguelov, D., Erhan, D., Vanhoucke, V., Rabinovich, A.: Going deeper with convolutions. In: Computer Vision and Pattern Recognition (CVPR) (2015), <http://arxiv.org/abs/1409.4842>
39. Tulloch, A., Jia, Y.: High performance ultra-low-precision convolutions on mobile devices. arXiv:1712.02427 (2017)
40. Tulloch, A., Jia, Y.: Quantization and training of neural networks for efficient integer-arithmetic-only inference. Conference on Computer Vision and Pattern Recognition (CVPR) (2018)
41. Wei, X., Gong, R., Li, Y., Liu, X., Yu, F.: Qdrop: Randomly dropping quantization for extremely low-bit post-training quantization. In: International Conference on Learning Representations (2022)

42. Wu, H.: NVIDIA Low Precision Inference on GPU. GPU Technology Conference (2019)
43. Wu, H., Judd, P., Zhang, X., Isaev, M., Micikevicius, P.: Integer quantization for deep learning inference: Principles and empirical evaluation. CoRR **abs/2004.09602** (2020), <https://arxiv.org/abs/2004.09602>
44. Yao, Z., Dong, Z., Zheng, Z., Gholami, A., Yu, J., Tan, E., Wang, L., Huang, Q., Wang, Y., Mahoney, M.W., Keutzer, K.: Hawqv3: Dyadic neural network quantization. In: ICML (2021)
45. Yu, H., Wen, T., Cheng, G., Sun, J., Han, Q., Shi, J.: Low-bit quantization needs good distribution. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops. pp. 680–681 (2020)
46. Zhao, R., Hu, Y., Dotzel, J., Sa, C.D., Zhang, Z.: Improving neural network quantization without retraining using outlier channel splitting. In: Chaudhuri, K., Salakhutdinov, R. (eds.) Proceedings of the 36th International Conference on Machine Learning, ICML 2019, 9–15 June 2019, Long Beach, California, USA. Proceedings of Machine Learning Research, vol. 97, pp. 7543–7552. PMLR (2019), <http://proceedings.mlr.press/v97/zhao19c.html>
47. Zhou, S., Ni, Z., Zhou, X., Wen, H., Wu, Y., Zou, Y.: Dorefa-net: Training low bitwidth convolutional neural networks with low bitwidth gradients. ArXiv **abs/1606.06160** (2016)