

# NIPQ: Noise proxy-based Integrated Pseudo-Quantization

Juncheol Shin<sup>\*2</sup>, Junhyuk So<sup>\*1</sup>, Sein Park<sup>2</sup>, Seungyeop Kang<sup>3</sup>, Sungjoo Yoo<sup>3</sup>, and Eunhyeok Park<sup>1,2</sup>

<sup>1</sup>Department of Computer Science and Engineering, POSTECH

<sup>2</sup>Graduate School of Artificial Intelligence, POSTECH

<sup>3</sup>Department of Computer Science and Engineering, Seoul National University

{jchshin, junhyukso, seinpark, eh.park}@postech.ac.kr, {ksy8653, yeonbin}@snu.ac.kr

## Abstract

*Straight-through estimator (STE), which enables the gradient flow over the non-differentiable function via approximation, has been favored in studies related to quantization-aware training (QAT). However, STE incurs unstable convergence during QAT, resulting in notable quality degradation in low precision. Recently, pseudo-quantization training has been proposed as an alternative approach to updating the learnable parameters using the pseudo-quantization noise instead of STE. In this study, we propose a novel noise proxy-based integrated pseudo-quantization (NIPQ) that enables unified support of pseudo-quantization for both activation and weight by integrating the idea of truncation on the pseudo-quantization framework. NIPQ updates all of the quantization parameters (e.g., bit-width and truncation boundary) as well as the network parameters via gradient descent without STE instability. According to our extensive experiments, NIPQ outperforms existing quantization algorithms in various vision and language applications by a large margin.*

## 1. Introduction

Neural network quantization is a representative optimization technique that reduces the memory footprint by storing the activation and weight in a low-precision domain. In addition, when hardware acceleration is available (e.g., low-precision arithmetics [37, 48, 51, 57] or bit-serial operations [15, 19, 50]), it also brings a substantial performance boost. These advantages make network inference affordable in large-scale servers as well as embedded devices [2, 39], which has helped popularize it in various applications. However, quantization has a critical disadvantage, quality degradation due to limited degrees of freedom.

<sup>\*</sup>Equal contribution.

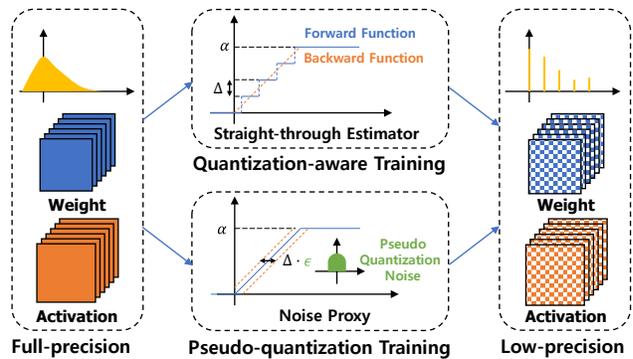


Figure 1. The proposed NIPQ as an alternative to QAT with STE.

The way to train the networks accurately within limited precision is critical and receiving much attention these days.

To mitigate the accuracy degradation, quantization-aware training (QAT) has emerged that trains a neural network with quantization operators to adapt to low precision. While the quantization operator is not differentiable, the straight-through estimator (STE) [5] allows the backpropagation of the quantized data based on linear approximation [20, 60]. This approximation works well in redundant networks with moderate precision ( $>4$ -bit). Thus, not only early studies [9, 41, 60] but also advanced ones [8, 17, 54] have proposed diverse QAT schemes based on STE and shown that popular neural networks (i.e., ResNet-18 [21]) can be quantized into 4-bit without accuracy loss.

However, STE-based QAT bypasses the approximated gradient, not the true gradient, and many studies have pointed out that it can incur instability and bias during training [20, 34, 36, 38, 42, 46]. For instance, PROFIT [38] points out that the instability is a major source of accuracy drop for the optimized networks (e.g., MobileNet-v2/v3 [24, 43]). More recently, an alternative training scheme, pseudo-quantization training (PQT) based on pseudo-quantization

noise (PQN), has been proposed [3, 10, 44] to address the instability of STE-based QAT. During PQT, the behavior of the quantization operator is simulated via PQN, and the learnable parameters are updated based on the proxy of quantization. While those studies are applied only to the weight, they can stabilize the training process significantly compared to QAT with STE and show the potential of PQT.

Nevertheless, the existing PQT algorithms have room for improvement. Various STE-based studies [8, 17, 60] have shown that truncation contributes significantly to reducing quantization errors, but even the advanced PQT studies [10, 44] use a naive min-max quantization. Integrating the truncation on the PQT framework can greatly reduce quantization error and enable us to exploit a PQT scheme for activation, which has an input-dependent distribution and requires a static quantization range. In addition, there is no theoretical support for whether PQT guarantees the optimal convergence of the quantization parameters. Intuitive interpretation exists, but proof of whether quantization parameters are optimized after PQT has yet to be provided.

In this paper, we propose a novel PQT method (Figure 1) called Noise proxy-based Integrated Pseudo-Quantization (NIPQ) that quantizes all activation and weight based on PQN. We present a novel idea, called Noise proxy, that shares the same quantization hyper-parameters (e.g., truncation boundary and bit-width) with the existing STE-based algorithm LSQ(+) [6, 17]. However, noise proxy allows to update the quantization parameters, as well as the network parameters, via gradient descent with PQN instead of STE. Subsequently, an arbitrary network can be optimized in a mixed-precision representation without STE instability. Our key contributions are summarized as follows:

- NIPQ is the first PQT that integrates truncation in addition to discretization. This extension not only further reduces the weight quantization error but also enables PQT for activation quantization.
- NIPQ optimizes an arbitrary network into the mixed-precision with awareness of the given resource constraint without human intervention.
- We provide theoretical analysis showing that NIPQ updates the quantization hyperparameters toward minimizing the quantization error.
- We provide extensive experimental results to validate the utility of NIPQ. It outperforms all existing mixed-precision quantization schemes by a large margin.

## 2. Related Work

Due to its practical usefulness, various studies have been proposed for multi-bit linear quantization [8, 17, 26, 60]. Most of these studies are based on STE, and the proposed

schemes have evolved by designing the quantization function to enable the optimization of the quantization parameters via gradient descent. However, in optimized networks such as MobileNet-v2, accuracy loss induced by STE instability has been reported for both activation [20] and weight [38]. They addressed the instability via a newly designed pipeline and non-linear approximation but suffered from the increased complexity and cost of QAT. In this work, NIPQ updates the quantization parameters without STE approximation, enabling stable convergence without additional cost or complexity and, most importantly, outperforming the existing methods by a large margin.

Mixed-precision studies focus on allocating layer-wise or group-wise bit-width in consideration of the precision sensitivity of each layer to minimize accuracy drop within the given resource constraints. Various methods have been proposed, e.g., RL-based [16, 33, 53], Hessian-based [13, 14, 58], and differentiable [52, 54] algorithms, but RL-based and Hessian-based methods are relatively complex, requiring a lot of parameter tuning, and differentiable algorithms still suffer from STE approximation error. In the present work, we reinterpret the differentiable bit-width tuning in terms of PQT framework, enabling the layer-wise bit-width tuning via gradient descent without STE instability.

Robust quantization [1, 7, 47] aims to guide the convergence of the network toward a smooth and flat loss surface based on additional regularization. The robustness of neural networks is highly beneficial for deploying noisy devices or low-precision ALUs. In the case of noise proxy, it inherently improves the robustness during PQT with PQN. As far as we know, we observe for the first time that the robustness of activation quantization can be enhanced (Section 7.2).

The benefit of the PQT has been demonstrated in diverse studies [10, 44] in various perspectives. However, previous studies utilize PQT with PQN in a limited domain only for weight and have yet to provide theoretical support for the convergence of PQT corresponding to the quantization. We extend the idea of PQT to both activation and weight by integrating the idea of truncation in addition to discretization and prove that the optimization of quantization parameters through noise proxy minimizes the actual quantization error. In addition, our work is the first to demonstrate that the PQT-based pipeline outperforms STE-based mixed-precision quantization algorithms by a large margin.

## 3. Simple Example Problem

Before introducing the details, we first propose a simple problem for straightforward explanation. The objective is to minimize the difference from the target data  $t \in \mathbb{R}^N$  and the quantized value of the learnable parameters  $x \in \mathbb{R}^N$ .

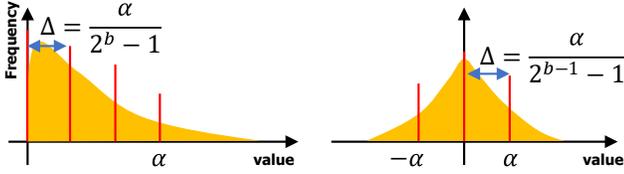


Figure 2. Visualization of the 2-bit quantization with hyper-parameters (e.g.,  $\alpha$  for truncation boundary and  $b$  for bit-width) for non-negative distribution (left) and symmetric distribution (right)

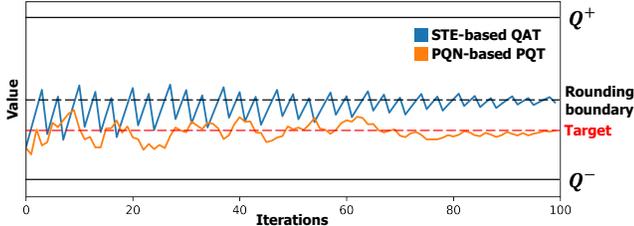


Figure 3. The trajectory of a scalar input during the minimization of  $\mathcal{L}^Q$  through SGD with STE-based QAT and PQN-based PQT.  $Q^+$  represents the positive quantization level and  $Q^-$  the negative.

The example loss function is defined as follows:

$$\mathcal{L}_{exp}^Q(x, \alpha, b) = \sum_{i=0}^{N-1} \left( t_i - Q(x_i | \alpha, b) \right)^2, \quad (1)$$

where subscript  $i$  means the  $i$ -th element of the data,  $t$  is arbitrary distribution, and  $Q(\cdot)$  is a given quantization function, parameterized  $\alpha$  (truncation boundary) and  $b$  (bit-width). To minimize the loss, we need to update  $x$ ,  $\alpha$ , and  $b$  judiciously. We utilize this example for the rest of the paper.

## 4. Motivation

In this work, we focus on linear quantization, especially for a more specific case that has an edge on hardware acceleration (e.g., affine layer-wise quantization for activation and symmetry channel-wise quantization for weight) [51]. To get optimal quality within restricted resources, we need to tune the quantization parameters, shown in Figure 2. In this section, we analyze the limitations of existing STE-based and PQT-based algorithms and provide insight for next-step innovation.

### 4.1. Limitation of STE-based Quantization

STE [5] allows the gradient flow over the quantized data by approximating the discretization as linear identity mapping. STE-based QAT updates the quantization parameters via gradient descent, which enables the joint update of network and quantization parameters toward minimizing the target loss value. Many prior works show successful results relying on this approximation [8, 17, 26, 34, 38, 54, 60].

However, STE-based QAT has a critical limitation. As shown in Figure 3, the scalar value never converges to the

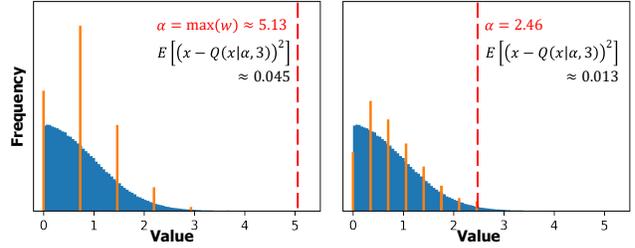


Figure 4. The comparison of 3-bit min-max quantization (left) and 3-bit quantization with truncation (right).

target value; instead, it keeps oscillating near the rounding boundary in the middle of the two nearest quantization levels. This is because when the scalar value is updated to cross the rounding boundary, it is mapped to a different quantization level. Then the sign of the gradient is reversed, updating the scalar value toward the opposite direction. The same process is repeated during training. Even at the later iterations, when the learning rate gradually decreases, the scalar value still oscillates near the rounding boundary.

The frequent flipping of the quantization value induces instability in network convergence. Moreover, the oscillation near the rounding boundary becomes the major source of large quantization errors. Even at the end of the training,  $Q(x | \alpha, b)$  can be mapped to either  $Q^+$  and  $Q^-$ , while  $Q^-$  has a minimal difference from the target value. Many studies have pointed out that this phenomenon degrades the quality of the quantized network, especially for the optimized networks [36, 38, 42].

### 4.2. Pros and Cons of Previous PQN-based PQT

To tackle the limitation of STE-based QAT, an alternative quantization pipeline, PQN-based PQT, has been proposed [10, 44]. During PQT, no quantization operator is used for forward and backward operation, but the behavior of the quantization operator is simulated via PQN. The learnable parameters are updated in a differentiable way without STE instability based on the proxy of quantization, becoming robust in the low-precision domain. As shown in Figure 3, the scalar value with PQT is updated stably toward the target value with small oscillation. Especially as the learning rate gets decreased, oscillation amplitude becomes reduced, eventually converging to the target value. After the PQT, the scalar value has been firmly mapped to  $Q^-$ , leading to minimal quantization error. Compared to STE-based QAT, PQN-based PQT is expected to have lower quantization error after the training.

However, existing PQT-related studies still have critical limitations. The advanced methods [10, 44] utilize min-max quantization, which uses the minimum/maximum value of data for the quantization range. Many STE-based studies point out that truncation is essential to minimize the quanti-

zation error in the low-precision domain [8, 17, 26, 60]. As shown in Figure 4, applying quantization based on the minimum/maximum value of the data wastes invaluable quantization levels, resulting in higher quantization error. As a long tail or outlier pushes the boundary by a large margin, truncation reduces overall quantization error significantly, especially for near-zero values. In addition, quantization with truncation is crucial for extending PQT for activation quantization; because the activation is input-dependent, the static boundary needs to be trained over a training dataset for stable and fast inference. Unfortunately, the existing studies do not provide the theoretical integration of truncation on top of a PQN-based PQT framework.

## 5. NIPQ: Noise Proxy-based Integrated Pseudo-Quantization

In this work, we propose a novel PQT pipeline to compensate for the quality degradation after quantization without STE instability. Moreover, the proposed scheme includes the support of updating truncation, resulting in accurate quantization not only for weight but also for activation. Besides, we adopt the gradient-based bit-width optimization via continuous relaxation [10, 54]. This extension allows us to update all of the hyper-parameters for linear quantization via gradient descent, which greatly simplifies the layer-wise optimization process.

In this section, we define the proposed idea and present the analytical support for the optimal convergence of quantization parameters based on PQT. Our method is composed of a pair of algorithms: STE-based quantization  $Q(x|\alpha, b)$  [6, 17] and the corresponding noise proxy  $\tilde{Q}(x|\alpha, b)$ . The algorithms are expressed as follows:

$$Q(x|\alpha, b) = \begin{cases} 0, & x \leq 0 \\ \lfloor x/\Delta \rfloor \cdot \Delta, & 0 < x < \alpha \\ \alpha, & x \geq \alpha, \end{cases} \quad (2)$$

$$\tilde{Q}(x|\alpha, b) = \begin{cases} 0, & x \leq 0 \\ x + \epsilon \cdot \Delta, & 0 < x < \alpha \\ \alpha, & x \geq \alpha, \end{cases} \quad (3)$$

where  $\Delta = \frac{\alpha}{2^b - 1}$  represents the quantization step size or the difference between quantization levels,  $\epsilon$  represents the PQN, and  $\lfloor \cdot \rfloor$  means the rounding operation<sup>1</sup>.

During PQT,  $\tilde{Q}(x|\alpha, b)$  is applied on top of the quantization targets, activation and weight, and the learnable parameters are updated via gradient descent. Because  $\alpha$  and  $b$  are included in the gradient path, the training with noise proxy updates all of the quantization parameters as well as the network's learnable parameters jointly. If we set the target loss

<sup>1</sup>Please note that we present the quantization algorithm for non-negative distribution data, e.g., activation after ReLU, for simplicity. However, this approach holds the same for weight, except that quantization index is used in  $[-2^{bit-1}, 2^{bit-1} - 1]$ .

function,  $\mathcal{L}_{target}$  as the mixture of the task loss  $\mathcal{L}_{task}$ , and the cost loss  $\mathcal{L}_{cost}$ , for memory footprint or computation as follows:

$$\mathcal{L}_{target} = \mathcal{L}_{task} + \lambda \cdot \mathcal{L}_{cost}, \quad (4)$$

where  $\lambda$  is a hyper-parameter of balancing two loss functions. The layer-wise bit-width  $b$  is assigned within the available resource budget, and the corresponding  $\alpha$  is tuned appropriately to achieve the best quality possible.

After PQT,  $Q(x|\alpha, b)$  is used for the true quantization during inference. To achieve the high-quality output with the true quantization operator, the tuned  $\alpha$  and  $b$  through PQT with noise proxy should also have the lowest loss with  $Q(x|\alpha, b)$ . In the following subsections, we provide theoretical support that the optimization of  $\alpha$  and  $b$  through PQT guarantees the optimal convergence  $Q(x|\alpha, b)$  when we sample PQN judiciously.

### 5.1. Input Activation Convergence

First of all, we evaluate the convergence property of the learnable input tensor  $x$ . In STE-based algorithm, the gradient of input data is bypassed through STE algorithm. Thereby the derivative is expressed as follows:

$$\frac{\partial Q(x|\alpha, b)}{\partial x} = \begin{cases} 0, & x \leq 0 \\ 1, & 0 < x < \alpha \\ 0, & x \geq \alpha. \end{cases} \quad (5)$$

In the case of the example in Eq. 1, the gradient of  $x_i$  within the quantization interval ( $x_i \in [0, \alpha]$ ) is calculated as follows:

$$\frac{\partial \mathcal{L}_{exp}^Q(x, \alpha, b)}{\partial x_i} = 2(Q(x_i|\alpha, b) - t_i). \quad (6)$$

Therefore,  $x_i$  converges to  $t_i$  theoretically as the training progresses. However, as shown in Figure 3,  $x_i$  in practice fluctuates near the rounding boundary, except  $Q(x_i|\alpha, b) = t_i$  exist, due to the instability of STE-based QAT.

On the other hand, the gradient of input data regarding noise proxy is expressed as follows:

$$\frac{\partial \tilde{Q}(x|\alpha, b)}{\partial x} = \begin{cases} 0, & x \leq 0 \\ 1, & 0 < x < \alpha \\ 0, & x \geq \alpha, \end{cases} \quad (7)$$

which is exactly the same as the  $Q(\cdot)$  condition. When we update the example loss function via the proposed algorithm, the gradient of  $x_i$  within the quantization interval is evaluated as follows:

$$\frac{\partial \mathcal{L}_{exp}^{\tilde{Q}}(x, \alpha, b)}{\partial x_i} = 2(\tilde{Q}(x_i|\alpha, b) - t_i) = 2(x_i + \epsilon \cdot \Delta - t_i). \quad (8)$$

The gradient of  $x_i$  points toward the target,  $t_i$ , but with drift and noise from PQN. However, the drift can be amortized when the average of PQN is zero-centered ( $E[\epsilon] = 0$ ).

Note that the distribution of actual quantization errors is empirically known to have zero-mean [55]. In addition, because the update step size is proportional to the learning rate, the oscillation amplitude becomes smaller when the learning rate becomes smaller as learning progresses; eventually,  $x_i$  converges to the target value, as shown in Figure 3. Unlike the STE-based QAT, noise proxy-based PQT precisely evaluates the effect of oscillation induced by PQN and discards the bias of PQN gradually during training.

Besides, the convergence property of oscillation near the target value is greatly beneficial to enhance the robustness of network, as pointed out in several previous studies [7]. When some noise,  $\delta$ , exists, the objective function can be approximated via Taylor expansion as follows:

$$E[\mathcal{L}(x + \delta)] \quad (9)$$

$$\approx E[\mathcal{L}(x) + \delta \cdot \nabla_x \mathcal{L}(x) + \frac{1}{2} \delta^T \cdot \nabla_x^2 \mathcal{L}(x) \cdot \delta] \quad (10)$$

$$\approx \mathcal{L}(x) + \frac{E[\delta^2]}{2} \text{Tr}\{\nabla_x^2 \mathcal{L}(x)\}, \quad (11)$$

where the term related to the first derivative is removed since  $E[\delta] = 0$ , and the off-diagonal elements of the second derivative term become 0 because it relates to the expectation of multiplication of two i.i.d. samples. When the loss is converged to the minima point, the sum of eigenvalues,  $\nabla_x^2 \mathcal{L}(x)$ , should have a non-negative value. To minimize the average loss after training, the loss surface is guided to be converged to the minima having a lower Hessian trace value [7], resulting in enhanced network robustness.

## 5.2. Alpha and Bit-width Convergence

In addition to the input data, the noise proxy should guarantee the convergence of the quantization parameters in the optimal point that minimizes the loss with  $Q(x|\alpha, b)$ . First, let's consider the gradient of  $\alpha$  regarding quantization and noise proxy function.

$$\frac{\partial Q(x|\alpha, b)}{\partial \alpha} = \sum_{x_i \geq \alpha} 1 + \sum_{0 < x_i < \alpha} \frac{1}{N_{lv}} \lfloor \frac{N_{lv}}{\alpha} x \rfloor - \frac{x}{\alpha}, \quad (12)$$

$$\frac{\partial \tilde{Q}(x|\alpha, b)}{\partial \alpha} = \sum_{x_i \geq \alpha} 1 + \sum_{0 < x_i < \alpha} \epsilon \cdot \frac{1}{N_{lv}}, \quad (13)$$

where  $N_{lv} = 2^b - 1$  denotes the number of levels.

Two gradients have a difference only for the case when  $x \in (0, \alpha)$ . However, note that Equations 12 and 13 become identical when PQN  $\epsilon$  is sampled from the quantization noise distribution  $\lfloor x/\Delta \rfloor - x/\Delta$ . The same conclusion is observable for the number of quantization levels, which

is expressed as:

$$\frac{\partial Q(x|\alpha, b)}{\partial N_{lv}} = \sum_{0 < x_i < \alpha} \frac{x}{N_{lv}} - \frac{\alpha}{N_{lv}^2} \lfloor \frac{N_{lv}}{\alpha} x \rfloor, \quad (14)$$

$$\frac{\partial \tilde{Q}(x|\alpha, b)}{\partial N_{lv}} = \sum_{0 < x_i < \alpha} -\epsilon \frac{\alpha}{N_{lv}^2}. \quad (15)$$

Therefore, when we use the sampled noise following the quantization error distribution, the gradient of quantization parameters becomes identical to the true quantization operator while still enjoying the benefit of stochasticity of PQN.

In short, when we apply NIPQ whose noise is sampled from the quantization noise distribution, the gradient direction of  $\alpha$  and  $b$  is identical in both  $Q(x|\alpha, b)$  and  $\tilde{Q}(x|\alpha, b)$ . Therefore, when  $\alpha$  converges to the optimal point via PQT where  $\frac{\partial \mathcal{L}^{\tilde{Q}}}{\partial \alpha} = \sum_i \frac{\partial \mathcal{L}}{\partial \tilde{Q}_i(x_i|\alpha, b)} \cdot \frac{\partial \tilde{Q}_i(\cdot)}{\partial \alpha} = 0$ , then we can easily validate that  $\frac{\partial \mathcal{L}^Q}{\partial \alpha} = \sum_i \frac{\partial \mathcal{L}}{\partial Q_i(x_i|\alpha, b)} \cdot \frac{\partial Q_i(\cdot)}{\partial \alpha} = 0$ . Likewise, when  $\frac{\partial \mathcal{L}^{\tilde{Q}}}{\partial b} = 0$ , then  $\frac{\partial \mathcal{L}^Q}{\partial b} = 0$  and that is the point where the loss value with  $Q(x|\alpha, b)$  is minimized. Therefore, optimizing the quantization parameters  $\alpha$  and  $b$  based on NIPQ also minimizes the target loss with true quantization.

## 6. Additional Details

As explained above, NIPQ enables the compensation of quality degradation of the quantized network without STE approximation. However, in practice, we need to consider several details to maximize obtainable accuracy.

### 6.1. Stochastic Rounding for Bit-width Parameter

First, during PQT phase, we use the bit-width with stochastic rounding as follows:

$$b \leftarrow 2 + 14 \cdot \text{Sigmoid}(\beta), \quad (16)$$

$$b \leftarrow \lfloor b + U(-0.5, 0.5) \rfloor, \quad (17)$$

where  $U(x, y)$  represents the uniform noise in  $[x, y]$ , and  $\beta$  is a trainable value in the continuous domain corresponding to the bit-width  $b$ . The bit-width is narrowed into a range of 2 to 16 bits, then mapped to the integer value via stochastic rounding. The idea of continuous relaxation of bit-width and updating it through gradient descent has been proposed in several previous studies [10, 54]. They have used the continuous value as it is, but we find that these settings yield suboptimal convergence because of the domain difference in that only discrete bit-width can be a candidate of precision during inference. To resolve the domain gap, we propose the bit-width update with stochastic rounding. It is an unbiased estimator but makes a decision in discrete space, whereby the domain difference could be mitigated. Our experiment shows that this update improves the final accuracy by a large margin, especially in the sub-4-bit domain. The experiments are included in the Supplementary Materials.

## 6.2. Stabilization in the Late Training Stage

In addition, we need to update the batch normalization layers after finishing NIPQ training (BN update). PQN in NIPQ induces instability in the running statistics of the normalization layers, so we update the batch normalization statistics with the true quantization operator after finishing training. We also observe additional performance improvement by switching from NIPQ to STE-based QAT at the later stage of training and updating the last few epochs with a very small learning rate (QAT finetune), while the amount of improvement is far smaller than BN update. We speculate that QAT finetune gives an opportunity to finetune the learnable parameters as well as normalization statistics of the network to mitigate the minor instability from noise. The details are included in the Supplementary Materials.

## 6.3. PQN Sampling

Finally, we provide the convergence condition in the last section in that PQN should be sampled following the quantization error distribution. However, we observe that NIPQ with a uniform noise  $U(-\frac{\Delta}{2}, \frac{\Delta}{2})$  also shows the comparable result to the noise with the desired property<sup>2</sup>. In this configuration, the relative frequency of the quantization error is ignored. This makes the quantization noise overestimated, and the low bit tends not to be selected, but a similar result is achievable by making  $\lambda$  of the cost function larger. In particular, the difference in final quality is negligible if the QAT finetune is applied. In practice, sampling of quantization error distribution slows down the training significantly. Uniform noise with QAT finetune can be considered as a practical approach of reducing PQN sampling overhead while exploiting the benefit of PQT.

## 7. Experimental Setup

In order to validate the performance of NIPQ, we conduct comprehensive studies to measure the output quality and observe diverse properties in various applications. The results in this paper are obtained based on PQN sampled from the quantization noise, unless otherwise specified explicitly. Likewise, QAT finetune is used at the last few epochs, otherwise specified explicitly. We implement our algorithm on PyTorch [40] library, and the code is available at our repository<sup>3</sup>. The details of target loss and training pipeline are available in the Supplementary Materials.

### 7.1. Sensitivity-aware Mixed-precision

Mixed-precision quantization aims at improving accuracy while minimizing the overall storage footprint and

<sup>2</sup>Unlike the previous studies [3, 10], PQN from normal distribution shows inferior performance in our algorithm.

<sup>3</sup><https://github.com/ECOLab-POSTECH/NIPQ>

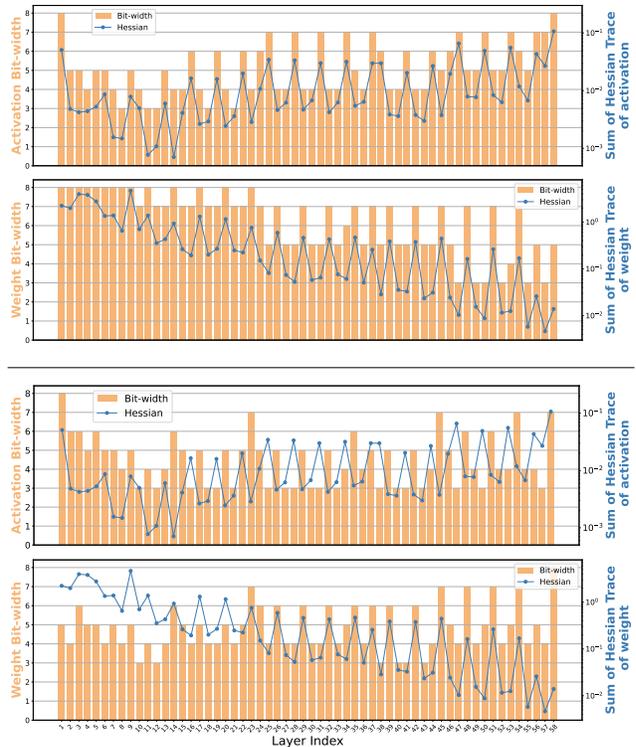


Figure 5. Layer-wise assigned bit-width of activation (top) and weight (bottom) and corresponding sensitivity measure (the sum of Hessian trace) of MobileNet-v2 on the CIFAR-100 dataset [28]. The average precision of activation weight is restricted to 4-bit (top) and the computation is restricted to 1.5 GBOPs (bottom).

computational cost. To maintain accuracy within the resource budget, more bit-width should be assigned to the sensitive layer to minimize overall quantization errors. Recently, the Hessian of the loss function has often been used as a metric of sensitivity against quantization [13, 14, 58]. The smaller the sum of the Hessian trace, the lower the sensitivity, and vice versa. Figure 5 (top) visualizes the assigned bit-width of activation and weight and the corresponding sum of the Hessian trace estimated via the Hutchinson algorithm [13] when applying NIPQ to MobileNet-v2 on the CIFAR-100 dataset. Each plot is obtained by optimizing the weights and activation independently. As shown in the figure, the more sensitive the layer is, the more bit-width is assigned. The redundant layers tend to have a fewer bit-width to match the target objective. Note that NIPQ does not have any additional stages that measure the sensitivity of the layer. Instead, we just train the network with an additional penalty term to restrict the average precision to 4-bit. The noise proxy algorithm allocates precision by itself, considering the sensitivity of the target layer, thereby quantizing the network with the highest accuracy as efficiently as possible within the constraints.

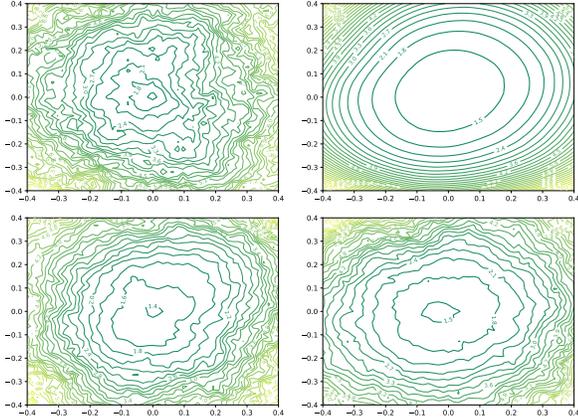


Figure 6. Loss landscapes [29] of 3-bit quantized MobileNet-V2 fine-tuned by STE-based LSQ [17] (top left), 3-bit NIPQ, BN update (top right), mixed-precision NIPQ (avg. 3-bit), BN update (bottom left), mixed-precision NIPQ (avg. 3-bit), QAT finetune (bottom right) on the CIFAR-100 dataset.

On the other hand, the automated bit-width allocation enables tuning of network considering complex metric that is practically improbable with the conventional Hessian-based methods [13, 14, 58]. Figure 5 (bottom) shows the assigned bit-width of activation and weight when restricting the computation cost (bit-operations, BOPs [4, 38]) as a 1.5G BOPs, which is equal to the computation cost of the quantized 4-bit model using PACT [8] or LSQ [17]. Unlike Figure 5 (top), the bit-width of activation is slightly misaligned with the sum of the hessian trace. In Figure 5 (bottom), we aim to optimize the average bit-width of activation and weight independently. Those two configurations are optimized via disjoint target losses thereby, each precision is assigned proportionally to the sensitivity of activation and weight separately. However, in the case of BOPs, the computation cost is proportional to the product of bit-width of activation and weight. Thereby, when we restrict the overall computation cost, the activation bit-width is allocated with awareness of the activation sensitivity as well as the corresponding weight sensitivity, and vice versa. This experimental result indicates that assigning the bit-width based on layer-wise sensitivity naively might not be an optimal policy for computation-aware quantization.

## 7.2. Robust Quantization for Practical Deployment

Another strength of NIPQ is the enhanced robustness of the network against unexpected noise. This brings diverse advantages to deploying the network in practice. For instance, the analog sum-based device [45] has significant energy efficiency but has inevitable noise induced by process variation or temperature drift, which causes instability of output. Even in the case of digital NPU, the low-precision implementation is fragmented because there are hundreds

	MNLI	CoLA	MRPC	QNLI	QQP	RTE	SST-2	STS-B
FP	84.46	56.53	90.62	90.66	87.76	71.19	92.32	88.64
DiffQ-MP-6 [10]	81.66	0.04	85.85	50.54	86.65	47.29	89.91	74.86
NIPQ-MP-6	<b>83.87</b>	<b>59.60</b>	<b>89.68</b>	<b>90.92</b>	<b>87.48</b>	<b>68.95</b>	<b>91.97</b>	<b>88.11</b>
DiffQ-MP-3 [10]	31.82	0.02	81.22	49.46	0	47.29	50.92	11.27
NIPQ-MP-3	<b>81.88</b>	<b>53.03</b>	<b>88.62</b>	<b>90.23</b>	<b>86.68</b>	<b>69.31</b>	<b>91.51</b>	<b>88.42</b>

Table 1. Results of weight quantization for BERT-base on GLUE

of hardware manufacturers [47]. When we deploy a neural network to the target device, the implementation difference can introduce unexpected noise on the data, resulting in accuracy degradation. The robustness of the network allows accuracy to be maintained in this environment, so securing this property is a significant advantage.

As explained in Section 5.1, NIPQ regularizes the loss surface. Figure 6 visualizes the sharpness of the loss surface by measuring the change of loss values after adding noise on top of the trained weight along the two random vectors [29]. As shown in the figure, NIPQ converges to a flat and smooth loss surface (top left vs. right). When layer-wise mixed precision is enabled, the loss surface becomes complex (top right vs. bottom left). We speculate that mixed-precision optimization exploits the sensitivity to trade robustness for better output quality. NIPQ enables such a trade-off in an automatic manner by forcing the network to adapt to the pseudo-noise that generalizes the quantization noise.

Besides, NIPQ enhances the robustness of the quantization parameters as well as the network parameters. Check the Supplementary Materials for the detailed results.

## 7.3. Importance of Truncation for Quantization

The last line of the related study is DiffQ [10], which is based on PQN-based PQT for weight quantization. The key difference from NIPQ is that DiffQ uses linear quantization based on the min-max value of weight instead of truncation. However, to minimize quantization errors with the limited bit-width, the presence of truncation is extremely helpful. In general, the data distribution of a natural network follows a bell-shaped distribution, where the majority of data is concentrated near zero. Because min-max quantization increases the quantization interval near zero, the quantization error greatly increases. Table 1 compares the accuracy after the quantization between DiffQ and NIPQ. Note that only the weight is quantized, and bit-width is allocated per layer. As shown in the table, NIPQ outperforms DiffQ by a large margin in the same bit-width. Integration of truncation on quantization minimizes the overall quantization error significantly, especially when the precision is limited.

## 7.4. Quantization of Large-scale Vision Tasks

To demonstrate the outstanding performance of the proposed method, we apply NIPQ for large-scale vision applications, including ImageNet [11] classification, multi-scale super-resolution, and VOC [18] object detection.

Model	Method	Bit(W/A)	BOPs(G)	Top-1
ResNet-18 [21]	FP	32/32	1857.6	70.54
	FP+KD	32/32	1857.6	72.17
	PACT [8]	4/4	34.7	69.2
	LSQ [17]	4/4	34.7	69.39
	DJPQ [54]	MP-BOPs	35.0	69.3
	HAQ [53]	MP-BOPs	34.4	69.2
	HAWQ [14]	MP-BOPs	34.0	68.5
	HAWQ-V3 [58]	MP-BOPs	34.0	68.5
	HAWQ-V3 [58]	MP-BOPs	72.0	70.2
	NIPQ-MP-4	3.89/3.98	45.5	<b>69.84</b>
	NIPQ-BOPs	4.53/3.67	34.1	<b>69.47</b>
	NIPQ-MP-4+KD	4.09/4.17	50.73	<b>71.83</b>
NIPQ-BOPs+KD	4.47/3.65	34.2	<b>71.24</b>	
MobileNet-v2 [43]	FP	32/32	306.8	72.6
	FP+KD	32/32	306.8	73.41
	DSQ [20]	4/4	15.8	64.8
	LSQ [17]	4/4	15.8	70.46
	DJPQ [54]	MP-BOPs	7.9	69.3
	HAQ [53]	MP-BOPs	8.3	69.5
	DuQ+KD [38]	4/4	5.3	69.86
	NIPQ-MP-4	3.79/4.00	8.67	<b>71.52</b>
	NIPQ-BOPs	6.19/5.31	8.25	<b>72.26</b>
	NIPQ-BOPs	5.19/4.24	5.37	<b>70.92</b>
	NIPQ-MP-4+KD	3.80/4.01	8.73	<b>72.16</b>
	NIPQ-BOPs+KD	6.21/5.35	8.31	<b>72.94</b>
NIPQ-BOPs+KD	5.24/4.24	5.34	<b>71.58</b>	
MobileNet-v3 [24]	FP	32/32	218.7	74.52
	FP+KD	32/32	218.7	75.8
	PACT [8]	4/4	3.46	67.98
	PACT+KD [8]	4/4	3.46	70.16
	DuQ [38]	4/4	3.46	69.50
	DuQ+KD [38]	4/4	3.46	71.01
	NIPQ-MP-4	3.97/3.99	7.77	<b>72.41</b>
	NIPQ-BOPs	5.02/3.71	3.26	<b>70.96</b>
	NIPQ-MP-4+KD	4.05/4.00	7.96	<b>74.49</b>
	NIPQ-BOPs+KD	4.89/3.69	3.29	<b>72.41</b>

Table 2. Top-1 accuracy (%) of quantized networks on ImageNet dataset. MP-BOPs represents the mixed-quantization with the bit-operations (BOPs) constraint while MP-N is that with N-bit average bit-width. ‘\*’ denotes the first and last layers remaining 8-bit, and KD denotes the knowledge distillation [22].

	Target Bit-width (Weight / Activation)						
	FP/FP	5/8	4/8	3/8	5/5	4/4	3/3
DoReFa [60]	0.857	0.593	0.541	0.359	0.588	0.498	0.288
PACT [8]	0.857	0.845	0.835	0.806	0.838	0.811	0.708
LSQ [17]	0.857	0.85	0.843	0.815	0.843	0.823	0.761
NIPQ-MP	0.857	<b>0.851</b>	<b>0.848</b>	<b>0.833</b>	<b>0.848</b>	<b>0.836</b>	<b>0.801</b>

Table 3. mAP comparison of Yolov5-S [25] on the PASCAL VOC [18].

First, we apply quantization to diverse networks on the ImageNet [11] classification task, and performs comparisons with various existing studies. Existing studies report mixed results with/without using well-known teacher-student-based knowledge distortion, so we report the accuracy for both cases, with EfficientNet-B0 as a teacher when necessary. Only the input image has 8-bit precision, and every layer in the network, including the first and last layers, is quantized via noise proxy.

Table 2 shows top-1 accuracy of the quantized networks.

NIPQ shows outstanding results for the optimized but hard to quantize networks such as MobileNet-v2/v3, as well as redundant networks such as ResNet-18. As shown in the table, existing methods are inferior to the proposed method, which achieves high accuracy in the same bit-width or bit-operations. These outstanding results come from two factors: first, NIPQ allows the network to converge to a more robust space without STE-induced instability, and second, within the resource budget, the quantization hyperparameters could be automatically tuned without the intervention of any hand-crafted manipulations. The benefit of these properties is maximized in optimized networks such as MobileNet-v2/v3. Note that when the average precision is constrained, NIPQ tries to increase accuracy with additional operations and vice versa. The automated tuning allows us to quantize the network considering our target goal.

In addition, we conduct an experiment to quantize the object detection task, which is known to be difficult to quantize. The difficulty is rapidly increased because we apply quantization to the advanced optimized network, Yolov5-S [25]. Table S6 shows the comparison of existing quantization studies in the same average bit-width. Due to large accuracy loss, existing 4-bit solutions are difficult to use in reality, but NIPQ shows practical, reliable quality in 4-bit precision. The results validate the stability of the NIPQ regardless of the difficulty of the target task.

Finally, to validate the superiority of NIPQ on the regression application, we apply quantization to the super-resolution task. In the case of super-resolution task, that conducts image restoration, the advantage of dynamic quantization, whose quantization parameters are updated regarding input data, is well demonstrated. Even in this case, NIPQ outmatches the best dynamic method, DDTB [59]. More results are available in the Supplementary Materials.

## 8. Conclusion

To achieve reliable output in low-precision, we propose NIPQ, which enhances the benefit of PQT pipeline by integrating the idea of truncation and providing theoretical support for it. NIPQ automates the layer-wise low-precision optimization for an arbitrary network by updating the network parameters and quantization parameters jointly via gradient descent toward minimizing the target loss. Our extensive results show that the proposed scheme outperforms all of the existing studies by a large margin for various vision and language applications.

## 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). We appreciate valuable comments from Myeonghwan Ahn at SNU.

## Supplementary Materials

### S-I. Overview

In this supplementary material, we present the details of our implementation and additional experimental results for various tasks and different datasets. We provide the following items:

- The detailed implementation of the cost loss function in Section [S-II](#).
- Detailed Configurations of our experiments in Section [S-III](#)
- The results of quantization for super resolution task in Section [S-IV](#).
- Additional experimental results of object detection on MS-COCO dataset in Section [S-V](#).
- An ablation study on the effect of stochastic rounding in Section [S-VI](#).
- An ablation study on the effect of late training stage in Section [S-VII](#).
- Experimental results on quantization parameter robustness in Section [S-VIII](#).
- The visualization of the quantization noise distribution in Section [S-IX](#)

### S-II. Cost Loss Function

In order to restrict the utilization of memory and computation resources, we introduce an additional cost loss function in addition to the target loss, as explained in Equation (4) of the main paper. The cost functions for the memory consumption  $\mathcal{L}_{cost-MP}$  and the computation cost  $\mathcal{L}_{cost-BOP}$  are defined as follows:

$$\mathcal{L}_{cost-MP} = \lambda_w h\left(\frac{\sum_i \lfloor b_i^w \rfloor \cdot e_i^w}{\sum_i e_i^w} - b_t\right) + \lambda_a h\left(\frac{\sum_i \lfloor b_i^a \rfloor \cdot e_i^a}{\sum_i e_i^a} - b_t\right), \quad (S1)$$

$$\mathcal{L}_{cost-BOP} = \lambda_b h(\sum_i \lfloor b_i^w \rfloor \cdot \lfloor b_i^a \rfloor \cdot OPS_i - b_t), \quad (S2)$$

where  $h(\cdot)$  denotes Huber loss,  $b_i^w/b_i^a$  denote the bit-width of  $i$ -th layer’s weight/activation,  $e_i^w/e_i^a$  are the number of elements in the  $i$ -th layer’s weight/activation,  $OPS_i$  is FLOPS of the  $i$ -th layer and  $b_t$  denotes the target bit-width.  $\mathcal{L}_{cost-MP}$  regularizes the average bit-width of activation/weight to the target bit, and  $\mathcal{L}_{cost-BOP}$  regularizes the sum of overall bit-operation (BOPs) to the target BOPs. Note that we utilize the bit-operations (BOPs) as a representative metric to measure the computation cost of a

neural network, which is commonly used in many previous studies [13, 14, 58]. However, any arbitrary differentiable function can be used as a drop-in replacement for the cost function, and NIPQ automatically optimizes the layer-wise bit-width to the sweet spot.

On the other hand, while the per-layer (or per-tensor) bit-width also requires rounding operation during forward operation, NIPQ is not applicable for the bit-width because it relies on the statistics of quantization error, but it is improvable to achieve the statistics for the scalar value. To overcome this limitation, we propose to update the bit-width via stochastic rounding with STE (Section [S-VI](#)).

### S-III. Experimental Configuration

In this paper, all experiments are conducted using GPU servers having 8 x NVIDIA GTX3090 with 24 GB VRAM with 2 x AMD 7313 (16 Core 32 T). The number of GPUs is selected to satisfy the minimum requirement of GPU memory for the target task. All of the experiments are implemented based on the PyTorch [40] framework (v1.12.1) [40]. Our source code is also provided. The additional details of training configuration, e.g., optimizer type, initial learning rate, decay policy, etc., are determined depending on the characteristics of applications and provided in the following paragraphs.

Table [S1](#) shows the configurations of ImageNet training for NIPQ results. In this experiment, we apply quantization to every convolution and linear layer, including the first and last layers. One exception is that the input of the first convolution layer is fixed as 8-bit. We use SGD with momentum optimizer and cosine annealing with warmup scheduling for learning rate adjustment [35].  $\eta_{min}$  is the final LR multiplier of cosine annealing, and  $\lambda_w$ ,  $\lambda_a$ , and  $\lambda_b$  are the hyper-parameter of resource constraints for the bit-width of weight, bit-width of activation, and BOPs, respectively. When knowledge distillation is triggered, we use EfficientNet-B0 [49] as a teacher network. We use the conventional dark-knowledge-based distillation [22].

Tables [S2](#) and [S3](#) show the detailed configurations of super-resolution task and object detection task, respectively. In both experiments, we keep the precision of the first and last layers as full-precision and apply low-precision quantization to the rest of the layers. In the super-resolution task, we use Adam optimizer [27] and cosine annealing scheduling for learning rate adjustment. In the object detection task, we use SGD with momentum optimizer and cosine annealing with warmup scheduling for learning rate adjustment. Like the image classification task,  $\eta_{min}$  is the final LR multiplier of cosine annealing, and  $\lambda_w$ ,  $\lambda_a$ , and  $\lambda_b$  represent the hyper-parameters of resource constraints for the bit-width of weight, bit-width of activation, and BOPs, respectively.

Table [S4](#) shows the detailed configurations of BERT-base [12] on the GLUE Task dataset. In this experiment,

Table S1. Fine-tuning configurations of ImageNet classification task.

Configuration		Epoch		SGD		Cosine annealing with warmup		$\lambda$		
		Stage-1	Stage-2	LR	Weight decay	Warmup len	$\eta_{min}$	$\lambda_w$	$\lambda_a$	$\lambda_b$
ResNet-18	ImageNet	40	3	0.04	$1 \times 10^{-5}$	3	$1 \times 10^{-3}$	1	1	1
MobileNet-v2	Cifar100	30	3	0.04	$5 \times 10^{-5}$	5	$1 \times 10^{-3}$	1	1	1
	ImageNet	40	3	0.04	$1 \times 10^{-5}$	3	$1 \times 10^{-3}$	1	1	3
MobileNet-v3	ImageNet	40	3	0.04	$1 \times 10^{-5}$	3	$1 \times 10^{-3}$	1	1	3

Table S2. Fine-tuning configurations of super-resolution task with EDSR.

Configuration		Epoch		Adam		Cosine annealing	$\lambda$		
		Stage-1	Stage-2	LR	Weight decay	$\eta_{min}$	$\lambda_w$	$\lambda_a$	$\lambda_b$
EDSR 4bit	DIV2K	30	10	$1 \times 10^{-4}$	0	$1 \times 10^{-3}$	15	15	-
EDSR 3bit	DIV2K	40	10	$1 \times 10^{-4}$	0	$1 \times 10^{-3}$	15	15	-

Table S3. Fine-tuning configurations of object detection task with YoloV5-S.

Configuration		Epoch		SGD		Cosine annealing with warmup		$\lambda$		
		Stage-1	Stage-2	LR	Weight decay	Warmup len	$\eta_{min}$	$\lambda_w$	$\lambda_a$	$\lambda_b$
YoloV5-S	Pascal VOC	30	5	0.0032	$3.6 \times 10^{-4}$	5	$1 \times 10^{-1}$	1	1	-
	COCO	35	5	0.0032	$3.6 \times 10^{-4}$	5	$1 \times 10^{-1}$	1	1	-

Table S4. Fine-tuning configurations of GLUE Dataset with BERT-base.

Configuration		Epoch		AdamW		Cosine annealing with warmup		$\lambda$		
		Stage-1	Stage-2	LR	Weight decay	Warmup len	$\eta_{min}$	$\lambda_w$	$\lambda_a$	$\lambda_b$
BERT-base	GLUE	25	5	1e-5	$1 \times 10^{-1}$	5	?	1	1	1

we modified the code from the huggingface-transformer [56] library. We apply weight quantization to every linear layer except the last classification head. Note that we do not quantize activation or word embedding. We use the AdamW optimizer and CosineLR scheduler for fine-tuning BERT except for the bit parameters because we find that AdamW can induce instability during training when the magnitude of the cost loss is too large. We use SGD with momentum optimizer for the bit parameters. Besides, we also find that  $\alpha$  and  $b$  parameters are not well trained when a single global learning rate is utilized ( $1e-5$ ). For fast and reliable convergence, we use the learning rate of  $1e-2$  for bit parameters. In addition, the gradient of  $\alpha$  is multiplied  $2^b - 1$  times over the global learning rate.

#### S-IV. Super Resolution Experiments

Table S5 shows the quantitative analysis of NIPQ on super resolution task, and Figure S1 visualizes the quality of the generated figures. We report PSNR as a quantitative measure, one of the well-known metrics in the area of super resolution. NIPQ outmatches the specialized quantization algorithm for super resolution, DDTB [59], which applies dynamic quantization that adjusts the quantization step size depending on the input data. These experimental results indicate that NIPQ works well in the regression task as well.

Network	Method	Dataset							
		Set5		Set14		BSD100		Urban100	
		4bit	3bit	4bit	3bit	4bit	3bit	4bit	3bit
EDSRx2	DoReFa [60]	37.22	37.13	32.82	32.73	31.63	31.57	30.17	30
	TFLite [51]	37.64	37.33	33.24	32.98	31.94	31.76	31.11	30.48
	PACT [8]	37.57	37.36	33.2	32.99	31.93	31.77	31.09	30.57
	PAMS [30]	37.67	36.76	33.2	32.5	31.94	31.38	31.1	29.5
	DDTB [59]	37.72	37.51	<b>33.35</b>	33.17	<b>32.01</b>	31.89	<b>31.39</b>	31.01
	NIPQ	<b>37.74</b>	<b>37.66</b>	33.29	<b>33.20</b>	<b>32.01</b>	<b>31.95</b>	31.36	<b>31.13</b>
EDSRx4	DoReFa [60]	30.91	30.76	27.78	26.66	27.04	26.97	24.73	24.59
	TFLite [51]	31.54	31.05	28.2	27.92	27.31	27.12	25.28	24.85
	PACT [8]	31.32	30.98	28.07	27.87	27.21	27.09	25.05	24.82
	PAMS [30]	31.59	27.25	28.2	25.24	27.32	25.38	25.32	22.76
	DDTB [59]	<b>31.85</b>	31.52	<b>28.39</b>	28.18	<b>27.44</b>	27.3	<b>25.69</b>	25.33
	NIPQ	31.73	<b>31.62</b>	28.34	<b>28.25</b>	27.41	<b>27.36</b>	25.56	<b>25.39</b>

Table S5. PSNR comparison of quantized EDSR [31] of scale 4 and scale 2

#### S-V. Additional Experiments on Object Detection

We conduct an additional experiment on object detection task with the COCO dataset and report mAP on Table S6. NIPQ obtains the best results compared to existing quantization studies in the same average bit-width.

In addition, in Figure S2, we visualize the qualitative results of NIPQ on the VOC dataset. Bounding box regression and classification results of the quantized network are presented. As shown in the figure, NIPQ works surprisingly

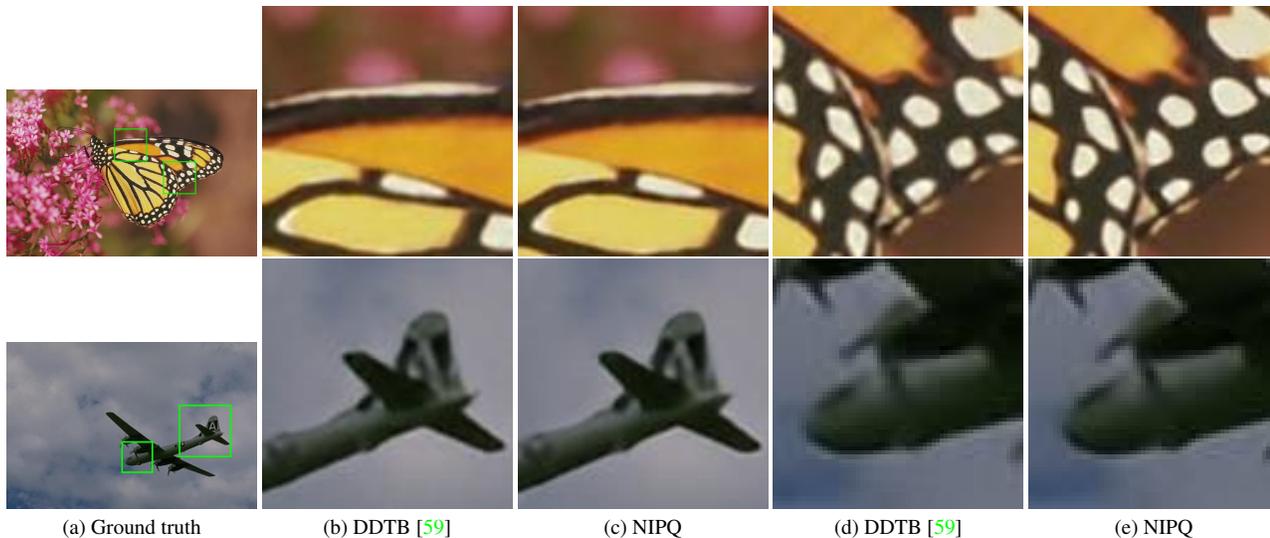


Figure S1. Qualitative results of super resolution on DIV2K dataset. EDSR $\times 4$  is quantized into 3-bit for both weights and activations.

Table S6. mAP comparison of YoloV5-S [25] on COCO dataset [32]

	Bit-width (Weight / Activation)			
	FP/FP	5/5	4/4	3/3
DoReFa [60]	0.354	0.266	0.24	0.191
PACT [8]	0.354	0.313	0.294	0.246
LSQ [17]	0.354	0.32	0.291	0.235
NIPQ	<b>0.354</b>	<b>0.33</b>	<b>0.317</b>	<b>0.284</b>

well in the 3-bit domain on the challenging object detection problem. YoloV5-S has a complicated structure, and the sensitivity of each layer is highly different. Because NIPQ has the ability to allocate the bit-width aware of the sensitivity automatically and enable stable convergence without STE instability, the quality of the quantized network outperforms all of the previous methods by a large margin.

## S-VI. Stochastic Rounding for Bit-width

While we propose an alternative training scheme for quantization instead of using STE, updating the bit-width is a remaining problem that is not addressed in the NIPQ pipeline. The proposed noise proxy is designed to update the learnable parameters by emulating the quantization operator based on PQN. However, the bit-width is assigned as a scalar value per the target tensor, and thereby it is impossible to aggregate the coarse-grained effect of the quantization operator. When we use rounding-based QAT with STE approximation, the bit-width also suffers from the instability of STE, resulting in highly unreliable result, as shown in Figure S3. Due to this limitation, many previous studies rely on the continuous approximation of bit-width during

training [10, 54] to avoid the instability problem. However, the representation mismatches to the domain of bit-width, resulting in suboptimal convergence in practice, especially when the target bit-width is in a sub-4-bit domain. In this paper, we propose an alternative idea to utilize the stochastic rounding of bit-width during training. Stochastic rounding is an unbiased estimator, so the learnable bit-width converges to the optimal point as the learning progresses. In addition, the bit-width is evaluated in the discrete domain during training, which mitigates the domain gap between training and inference. As shown in Figure S3, the stochastic rounding consistently draws the pareto-front line with small variance, which enables us to search for the best quantization configurations within the given resource budget.

## S-VII. Comparison for the Late Training Stage

Table S7. Comparison of accuracy regarding the late training stage. MobileNet-v2 is trained in 30 epochs and finetuned in 3 epochs on the CIFAR-100 dataset. The target computation overhead is 1.0 GBlops

	FP	Without Tuning	BN update	QAT finetune
Top-1	75.04	70.45	72.99	73.29

In Table S7, we show the results of NIPQ with different late training stage policies. As shown in the table, BN update offers a large performance benefit compared to the accuracy of the NIPQ training without the late stage tuning. Because PQN of NIPQ disturbs the statistics of normalization layers, the correction of the statistics is essential to maximize the accuracy in the inference phase. In addition, QAT finetune offers an additional performance improve-

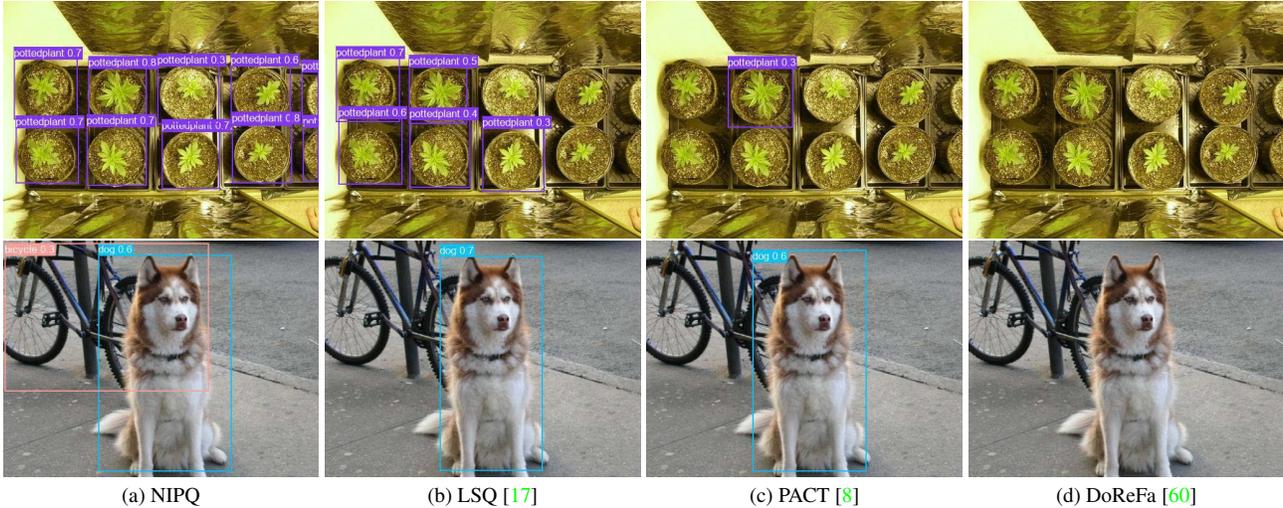


Figure S2. Qualitative results of object detection on the VOC dataset. Yolov5-S is quantized into 3-bit weights and activations according to each quantization method.

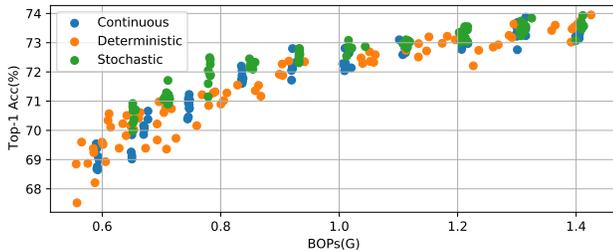


Figure S3. Accuracy comparison of MobileNet-v2 at CIFAR-100 dataset with different bit-width training strategies. We ran the experiments with 10 repetitions, increasing from 0.5 BOPs to 1.4 BOPs by 0.1 BOPs steps.

ment by giving an additional chance to adjust the learnable parameters of the entire network without the effect of PQN with a small learning rate, which enables the stabilization of network parameters near the optimal point with the tiniest effect of STE instability.

### S-VIII. Robustness of the quantization parameters

NIPQ also enhances the robustness of the quantization parameters as well as the network parameter. Figure S4 visualizes the results of measuring the accuracy while changing the quantization step size or the truncation interval. The more robust the network, the more it can endure the change of the quantization configuration. As shown in the figure, NIPQ shows comparable or superior results to the previous best algorithm for robustness, KURE [47]. It is especially worthy that existing studies have focused on improving the robustness of weight only [1, 47], but NIPQ also improves

the robustness of activation as well by a large margin. To the best of our knowledge, this is for the first time that activation robustness can be improved, which is a crucial benefit of deploying networks in a noisy environment.

### S-IX. Quantization Noise Distribution

In Figure S5, we visualize the quantization noise distribution of ResNet-18’s 12-th convolution weight in different bit-widths. When applying quantization, not only rounding but also truncation is applied. The previous study argues that the quantization noise distribution follows a uniform distribution regardless of input distribution when the number of bits is sufficiently large [55]. According to our observation, the statement is held in practice when the bit-width is larger than 4-bit. However, in sub-4-bit precision, the distribution of noise seems to follow a bell-shaped curve instead of a uniform distribution. Due to these characteristics, conventionally uniform or gaussian distributions are often used to approximate PQN [10, 54]. However, as presented in this paper, the precise sampling of PQN following the quantization error distribution is essential to guarantee the convergence on the optimal point, while the uniform distribution shows comparable results in practice empirically. In this work, we realize the sampling process of quantization error distribution on GPU with practical performance as follows: first, the probability density function (PDF) of the quantization error distribution is estimated based on the histogram with 256 bins. Then, the distribution is sampled from the estimated PDF of the histogram. As shown in Figure S5, the sampled distribution precisely follows the quantization error distribution.

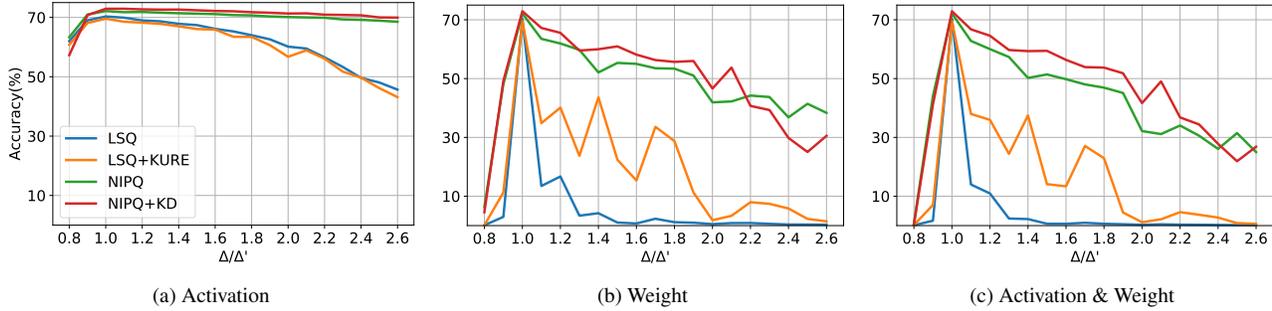


Figure S4. Robustness of quantized MobileNet-V2 on ImageNet against the change of  $\alpha$  of the quantization operator for weight and activation.  $\Delta'$  is the trained  $\alpha$  and  $\Delta$  is the scaled one. NIPQ+KD represents the quantized network with knowledge distillation [23] using EfficientNet-B0 [49] as a teacher.

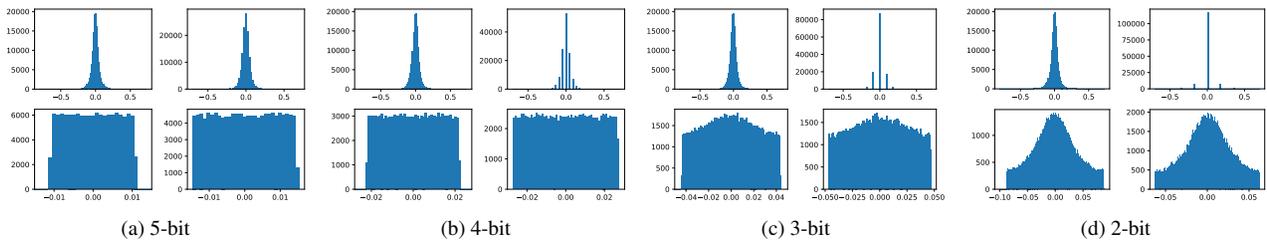


Figure S5. Quantization noise distribution of ResNet-18's 12-th convolution weight. (left top) 32 bit (FP) distribution. (right top) N-bit uniform quantized distribution. (left bottom) Real quantization noise distribution. (right bottom) Sampled quantization noise distribution.

## References

- [1] Milad Alizadeh, Arash Behboodi, Mart van Baalen, Christos Louizos, Tijmen Blankevoort, and Max Welling. Gradient  $\ell_1$  regularization for quantization robustness. *International Conference on Learning Representations (ICLR)*, 2020. 2, 4
- [2] Mário Almeida, Stefanos Laskaridis, Abhinav Mehrotra, Lukasz Dudziak, Ilias Leontiadis, and Nicholas D. Lane. Smart at what cost?: characterising mobile deep neural networks in the wild. *Internet Measurement Conference (IMC)*, 2021. 1
- [3] Chaim Baskin, Natan Liss, Yoav Chai, Evgenii Zheltonozhskii, Eli Schwartz, Raja Giryes, Avi Mendelson, and Alexander M. Bronstein. NICE: noise injection and clamping estimation for neural network quantization. *CoRR*, abs/1810.00162, 2018. 2, 6
- [4] Chaim Baskin, Eli Schwartz, Evgenii Zheltonozhskii, Natan Liss, Raja Giryes, Alexander M. Bronstein, and Avi Mendelson. UNIQU: uniform noise injection for the quantization of neural networks. *arXiv:1804.10969*, 2018. 7
- [5] Yoshua Bengio, Nicholas Léonard, and Aaron C. Courville. Estimating or propagating gradients through stochastic neurons for conditional computation. *CoRR*, abs/1308.3432, 2013. 1, 3
- [6] Yash Bhalgat, Jinwon Lee, Markus Nagel, Tijmen Blankevoort, and Nojun Kwak. LSQ+: improving low-bit quantization through learnable offsets and better initialization. *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020. 2, 4
- [7] Xiangning Chen and Cho-Jui Hsieh. Stabilizing differentiable architecture search via perturbation-based regularization. *International Conference on Machine Learning (ICML)*, 2020. 2, 5
- [8] Jungwook Choi, Zhuo Wang, Swagath Venkataramani, Pierce I-Jen Chuang, Vijayalakshmi Srinivasan, and Kailash Gopalakrishnan. PACT: parameterized clipping activation for quantized neural networks. *CoRR*, abs/1805.06085, 2018. 1, 2, 3, 4, 7, 8
- [9] Matthieu Courbariaux and Yoshua Bengio. Binarynet: Training deep neural networks with weights and activations constrained to +1 or -1. *CoRR*, abs/1602.02830, 2016. 1
- [10] Alexandre Défossez, Yossi Adi, and Gabriel Synnaeve. Differentiable model compression via pseudo quantization noise. *CoRR*, abs/2104.09987, 2021. 2, 3, 4, 5, 6, 7
- [11] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. *Computer Society Conference on Computer Vision and Pattern Recognition (CVPR)*, 2009. 7, 8
- [12] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018. 1
- [13] Zhen Dong, Zhewei Yao, Daiyaan Arfeen, Amir Gholami, Michael W. Mahoney, and Kurt Keutzer. HAWQ-V2: hessian aware trace-weighted quantization of neural networks. *Neural Information Processing Systems (NeurIPS)*, 2020. 2, 6, 7, 1

- [14] Zhen Dong, Zhewei Yao, Amir Gholami, Michael W. Mahoney, and Kurt Keutzer. HAWQ: hessian aware quantization of neural networks with mixed-precision. *International Conference on Computer Vision, (ICCV)*, 2019. 2, 6, 7, 8, 1
- [15] Charles Eckert, Xiaowei Wang, Jingcheng Wang, Arun Subramaniyan, Ravi R. Iyer, Dennis Sylvester, David T. Blaauw, and Reetuparna Das. Neural cache: Bit-serial in-cache acceleration of deep neural networks. *International Symposium on Computer Architecture (ISCA)*, 2018. 1
- [16] Ahmed T. Elthakeb, Prannoy Pilligundla, Fatemehsadat Mireshghallah, Amir Yazdanbakhsh, and Hadi Esmaeilzadeh. Releq : A reinforcement learning approach for automatic deep quantization of neural networks. *IEEE Micro*, 40(5):37–45, 2020. 2
- [17] Steven K. Esser, Jeffrey L. McKinstry, Deepika Bablani, Rathinakumar Appuswamy, and Dharmendra S. Modha. Learned step size quantization. *International Conference on Learning Representations (ICLR)*, 2020. 1, 2, 3, 4, 7, 8
- [18] Mark Everingham, Luc Van Gool, Christopher K. I. Williams, John M. Winn, and Andrew Zisserman. The pascal visual object classes (VOC) challenge. *Int. J. Comput. Vis.*, 88(2):303–338, 2010. 7, 8
- [19] Boyuan Feng, Yuke Wang, Tong Geng, Ang Li, and Yufei Ding. APNN-TC: accelerating arbitrary precision neural networks on ampere GPU tensor cores. *The International Conference for High Performance Computing, Networking, Storage and Analysis (SC)*, 2021. 1
- [20] Ruihao Gong, Xianglong Liu, Shenghu Jiang, Tianxiang Li, Peng Hu, Jiazhen Lin, Fengwei Yu, and Junjie Yan. Differentiable soft quantization: Bridging full-precision and low-bit neural networks. *International Conference on Computer Vision (ICCV)*, 2019. 1, 2, 8
- [21] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016. 1, 8
- [22] Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. Dark knowledge. *Presented as the keynote in BayLearn*, 2(2), 2014. 8, 1
- [23] Geoffrey Hinton, Oriol Vinyals, Jeff Dean, et al. Distilling the knowledge in a neural network. *arXiv preprint arXiv:1503.02531*, 2(7), 2015. 5
- [24] Andrew Howard, Ruoming Pang, Hartwig Adam, Quoc V. Le, Mark Sandler, Bo Chen, Weijun Wang, Liang-Chieh Chen, Mingxing Tan, Grace Chu, Vijay Vasudevan, and Yukun Zhu. Searching for mobilenetv3. *International Conference on Computer Vision (ICCV)*, 2019. 1, 8
- [25] Glenn Jocher, Alex Stoken, Jirka Borovec, Liu Changyu, Adam Hogan, L Diaconu, F Ingham, J Poznanski, J Fang, L Yu, et al. ultralytics/yolov5: v3. 1-bug fixes and performance improvements. *Zenodo*, 2020. 8, 3
- [26] Sangil Jung, Changyong Son, Seohyung Lee, JinWoo Son, Jae-Joon Han, Youngjun Kwak, Sung Ju Hwang, and Changkyu Choi. Learning to quantize deep networks by optimizing quantization intervals with task loss. *Conference on Computer Vision and Pattern Recognition, (CVPR)*, 2019. 2, 3, 4
- [27] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014. 1
- [28] Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. *CiteSeer*, 2009. 6
- [29] Hao Li, Zheng Xu, Gavin Taylor, Christoph Studer, and Tom Goldstein. Visualizing the loss landscape of neural nets. *Advances in Neural Information Processing Systems (NeurIPS)*, 2018. 7
- [30] Huixia Li, Chenqian Yan, Shaohui Lin, Xiawu Zheng, Baochang Zhang, Fan Yang, and Rongrong Ji. PAMS: quantized super-resolution via parameterized max scale. *European Conference on Computer Vision (ECCV)*, 2020. 2
- [31] Bee Lim, Sanghyun Son, Heewon Kim, Seungjun Nah, and Kyoung Mu Lee. Enhanced deep residual networks for single image super-resolution. In *2017 IEEE Conference on Computer Vision and Pattern Recognition Workshops, CVPR Workshops 2017, Honolulu, HI, USA, July 21-26, 2017*, pages 1132–1140. IEEE Computer Society, 2017. 2
- [32] Tsung-Yi Lin, Michael Maire, Serge J. Belongie, Lubomir D. Bourdev, Ross B. Girshick, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C. Lawrence Zitnick. Microsoft COCO: common objects in context. *CoRR*, abs/1405.0312, 2014. 3
- [33] Jing Liu, Jianfei Cai, and Bohan Zhuang. Sharpness-aware quantization for deep neural networks. *CoRR*, abs/2111.12273, 2021. 2
- [34] Zechun Liu, Baoyuan Wu, Wenhan Luo, Xin Yang, Wei Liu, and Kwang-Ting Cheng. Bi-real net: Enhancing the performance of 1-bit cnns with improved representational capability and advanced training algorithm. *European Conference on Computer Vision (ECCV)*, 2018. 1, 3
- [35] Ilya Loshchilov and Frank Hutter. Sgdr: Stochastic gradient descent with warm restarts. *arXiv preprint arXiv:1608.03983*, 2016. 1
- [36] Markus Nagel, Marios Fournarakis, Yelysei Bondarenko, and Tijmen Blankevoort. Overcoming oscillations in quantization-aware training. In *International Conference on Machine Learning*, pages 16318–16330. PMLR, 2022. 1, 3
- [37] Int4 precision for ai inference. <https://devblogs.nvidia.com/int4-for-ai-inference/>, 2019. Accessed: 2022-11-12. 1
- [38] Eunhyeok Park and Sungjoo Yoo. PROFIT: A novel training method for sub-4-bit mobilenet models. *European Conference on Computer Vision (ECCV)*, 2020. 1, 2, 3, 7, 8
- [39] Jongsoo Park, Maxim Naumov, Protonu Basu, Summer Deng, Aravind Kalaiah, Daya Shanker Khudia, James Law, Parth Malani, Andrey Malevich, Nadathur Satish, Juan Miguel Pino, Martin Schatz, Alexander Sidorov, Viswanath Sivakumar, Andrew Tulloch, Xiaodong Wang, Yiming Wu, Hector Yuen, Utku Diril, Dmytro Dzhulgakov, Kim M. Hazelwood, Bill Jia, Yangqing Jia, Lin Qiao, Vijay Rao, Nadav Rotem, Sungjoo Yoo, and Mikhail Smelyanskiy. Deep learning inference in facebook data centers: Characterization, performance optimizations and hardware implications. *CoRR*, abs/1811.09886, 2018. 1
- [40] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming

- Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: An imperative style, high-performance deep learning library. *Advances in neural information processing systems*, 32, 2019. [6](#), [1](#)
- [41] Mohammad Rastegari, Vicente Ordonez, Joseph Redmon, and Ali Farhadi. Xnor-net: Imagenet classification using binary convolutional neural networks. *European Conference on Computer Vision (ECCV)*, 2016. [1](#)
- [42] Charbel Sakr, Steve Dai, Rangha Venkatesan, Brian Zimmer, William Dally, and Brucek Khailany. Optimal clipping and magnitude-aware differentiation for improved quantization-aware training. *International Conference on Machine Learning (ICML)*, 2022. [1](#), [3](#)
- [43] Mark Sandler, Andrew G. Howard, Menglong Zhu, Andrey Zhmoginov, and Liang-Chieh Chen. Mobilenetv2: Inverted residuals and linear bottlenecks. *IEEE Conference on Computer Vision and Pattern Recognition, CVPR*, 2018. [1](#), [8](#)
- [44] Pedro Savarese, Xin Yuan, Yanjing Li, and Michael Maire. Not all bits have equal value: Heterogeneous precisions via trainable noise. *Advances in Neural Information Processing Systems (NeurIPS)*, 2022. [2](#), [3](#)
- [45] Ali Shafiee, Anirban Nag, Naveen Muralimanohar, Rajeev Balasubramonian, John Paul Strachan, Miao Hu, R. Stanley Williams, and Vivek Srikumar. ISAAC: A convolutional neural network accelerator with in-situ analog arithmetic in crossbars. *International Symposium on Computer Architecture (ISCA)*, 2016. [7](#)
- [46] Sungho Shin, Jinhwan Park, Yoonho Boo, and Wonyong Sung. Hhlp: Quantized neural networks training for reaching flat minima in loss surface. *The Thirty-Fourth AAAI Conference on Artificial Intelligence (AAAI)*, 2020. [1](#)
- [47] Moran Shkolnik, Brian Chmiel, Ron Banner, Gil Shomron, Yury Nahshan, Alex M. Bronstein, and Uri C. Weiser. Robust quantization: One model to rule them all. *Neural Information Processing Systems (NeurIPS)*, 2020. [2](#), [7](#), [4](#)
- [48] Jinook Song, Yunkyo Cho, Jun-Seok Park, Jun-Woo Jang, Sehwan Lee, Joon-Ho Song, Jae-Gon Lee, and Inyup Kang. 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. [1](#)
- [49] Mingxing Tan and Quoc V. Le. Efficientnet: Rethinking model scaling for convolutional neural networks. In Kamalika Chaudhuri and Ruslan Salakhutdinov, editors, *Proceedings of the 36th International Conference on Machine Learning, ICML 2019, 9-15 June 2019, Long Beach, California, USA*, volume 97 of *Proceedings of Machine Learning Research*, pages 6105–6114. PMLR, 2019. [1](#), [5](#)
- [50] Andrew Tulloch and Yangqing Jia. High performance ultra-low-precision convolutions on mobile devices. *arXiv:1712.02427*, 2017. [1](#)
- [51] Andrew Tulloch and Yangqing Jia. Quantization and training of neural networks for efficient integer-arithmetic-only inference. *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018. [1](#), [3](#), [2](#)
- [52] Stefan Uhlich, Lukas Mauch, Fabien Cardinaux, Kazuki Yoshiyama, Javier Alonso García, Stephen Tiedemann, Thomas Kemp, and Akira Nakamura. Mixed precision dnns: All you need is a good parametrization. *8th International Conference on Learning Representations, ICLR*, 2020. [2](#)
- [53] Kuan Wang, Zhijian Liu, Yujun Lin, Ji Lin, and Song Han. HAQ: hardware-aware automated quantization with mixed precision. *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019. [2](#), [8](#)
- [54] Ying Wang, Yadong Lu, and Tijmen Blankevoort. Differentiable joint pruning and quantization for hardware efficiency. *European Conference on Computer Vision (ECCV)*, 2020. [1](#), [2](#), [3](#), [4](#), [5](#), [8](#)
- [55] Bernard Widrow, Istvan Kollar, and Ming-Chang Liu. Statistical theory of quantization. *IEEE Transactions on instrumentation and measurement*, 45(2):353–361, 1996. [5](#), [4](#)
- [56] Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander M. Rush. Transformers: State-of-the-art natural language processing. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 38–45, Online, Oct. 2020. Association for Computational Linguistics. [2](#)
- [57] Hao Wu. NVIDIA Low Precision Inference on GPU. *GPU Technology Conference*, 2019. [1](#)
- [58] Zhewei Yao, Zhen Dong, Zhangcheng Zheng, Amir Gholami, Jiali Yu, Eric Tan, Leyuan Wang, Qijing Huang, Yida Wang, Michael W. Mahoney, and Kurt Keutzer. HAWQ-V3: dyadic neural network quantization. *Proceedings of the 38th International Conference on Machine Learning, ICML*, 2021. [2](#), [6](#), [7](#), [8](#), [1](#)
- [59] Yunshan Zhong, Mingbao Lin, Xunchao Li, Ke Li, Yunhang Shen, Fei Chao, Yongjian Wu, and Rongrong Ji. Dynamic dual trainable bounds for ultra-low precision super-resolution networks. *arXiv preprint arXiv:2203.03844*, 2022. [8](#), [2](#), [3](#)
- [60] Shuchang Zhou, Zekun Ni, Xinyu Zhou, He Wen, Yuxin Wu, and Yuheng Zou. Dorefa-net: Training low bitwidth convolutional neural networks with low bitwidth gradients. *CoRR*, abs/1606.06160, 2016. [1](#), [2](#), [3](#), [4](#), [8](#)