

Improving Quantization-aware Training of Low-Precision Network via Block Replacement on Full-Precision Counterpart

Chengting Yu^{1,2}, Shu Yang², Fengzhao Zhang², Hanzhi Ma^{1,2}, Aili Wang^{1,2,*}, Er-Ping Li^{1,2}

¹ College of Information Science and Electronic Engineering, Zhejiang University

² ZJU-UIUC Institute, Zhejiang University

chengting.21@intl.zju.edu.cn, ailiwang@intl.zju.edu.cn

arXiv:2412.15846v1 [cs.LG] 20 Dec 2024

Abstract—Quantization-aware training (QAT) is a common paradigm for network quantization, in which the training phase incorporates the simulation of the low-precision computation to optimize the quantization parameters in alignment with the task goals. However, direct training of low-precision networks generally faces two obstacles: 1. The low-precision model exhibits limited representation capabilities and cannot directly replicate full-precision calculations, which constitutes a deficiency compared to full-precision alternatives; 2. Non-ideal deviations during gradient propagation are a common consequence of employing pseudo-gradients as approximations in derived quantized functions. In this paper, we propose a general QAT framework for alleviating the aforementioned concerns by permitting the forward and backward processes of the low-precision network to be guided by the full-precision partner during training. In conjunction with the direct training of the quantization model, intermediate mixed-precision models are generated through the block-by-block replacement on the full-precision model and working simultaneously with the low-precision backbone, which enables the integration of quantized low-precision blocks into full-precision networks throughout the training phase. Consequently, each quantized block is capable of: 1. simulating full-precision representation during forward passes; 2. obtaining gradients with improved estimation during backward passes. We demonstrate that the proposed method achieves state-of-the-art results for 4-, 3-, and 2-bit quantization on ImageNet and CIFAR-10. The proposed framework provides a compatible extension for most QAT methods and only requires a concise wrapper for existing codes.

I. INTRODUCTION

Convolutional neural networks (CNNs) have made substantial advancements in image classification [15, 29], semantic segmentation [11, 17], object detection [43, 44], and image restoration [13]. With the large computational complexity of CNNs, there is an increasing demand for strategies to successfully deploy networks on resource-constrained devices with limited processing capability. Several strategies have been developed to train fast and compact neural networks that are both efficient in terms of speed and size, such as manually-designed architectures [7, 20, 47, 55], pruning [10, 32, 35], quantization [1, 8, 14, 22, 33], knowledge distillation [19, 45, 49], and neural architecture search. In this paper, we focus on network quantization, which attempts to transform the weights and/or activations of a full-precision (FP) network into low-precision (LP) approximations so that fixed-point

computations can be conducted without compromising too much accuracy.

Quantization-aware training (QAT) has become one of the most effective methods for achieving low-bit quantization while maintaining relatively high accuracy by simulating low-precision operators during training and optimizing quantizer parameters with task objectives [8, 27, 39]. The main difficulty in directly training a quantized network arises from the discretization stage, in which a discrete quantizer (i.e., scaled round function) converts a normalized value to one of the discrete levels. First, since a discrete quantizer is non-differentiable, the straight-through estimator (STE) [2] is often utilized during backpropagation [5, 8, 25, 31, 41, 53, 59]. Despite recent STE-based methods demonstrating reasonable performance [3, 8, 25, 31, 41], STE may produce an unexpected gradient mismatch issue [31, 52] with distinct forward and backward passes, resulting in an undesirable drop in accuracy [61]. Second, due to the limited representation capability of the discrete quantizer [26], the accuracy of extremely low-bit quantized networks (i.e., 4-, 3-, or 2-bit) is inevitably diminished. In order to enhance the low-bit representation, certain methods are suggested, known as non-uniform quantization [33, 34, 37, 51, 58]. This technique entails modifying the quantization resolution in accordance with the density of the real-valued distribution in order to find suitable correspondences with the activation and weight distributions [33, 51]. However, the increased inference burden presents them with major obstacles.

By consulting full-precision models for guidance as auxiliary supervision, optimal schemes are proposed to partially compensate for the inherent issues of QAT [26, 60, 61] through the incorporation of additional losses into intermediate layers to address the gradient approximation problem and provide regularization. In this regard, the simplest approach commonly employed in QAT is to initialize the low-precision model directly with the weights obtained from the pre-trained full-precision model [3, 8, 33]. Besides, training strategies such as knowledge distillation have been combined to learn a low-precision student network by distilling knowledge from a full-precision teacher [26, 42, 60]; additionally, some strategies use full-precision auxiliary routines for dealing with quantizer noise [4, 61].

In this paper, we propose a novel, neat, and efficient framework for QAT methods to utilize the guiding potential of full-precision counterparts properly. As shown in Fig. 1, instead of the straight and rough initialization, we build the mixed-precision models by converting the end FP blocks into LP blocks and then integrating those intermediate models into the entire training framework for implicit guidance. The weights of FP blocks are fixed (not changed during training), while the weights of LP blocks learn from both FP- and LP-backward routines. By grafting the LP model into its FP counterpart, the proposed framework allows the front LP blocks to pass the end FP blocks to ensure the reliability of backward gradients, as well as making the LP representation mirror the FP representation to accommodate the following FP feature extractor. The aforementioned concerns of both forward and backward passes are addressed in this case. Empirical experiments support our hypothesis: 1. Each intermediate model outperforms the vanilla LP model, and the performance tendency decreases from the top (pure FP model) to the bottom (pure LP model); 2. The guidance from intermediate models is positive and valid, leading FP blocks to converge toward training targets; 3. During training, the exchanging LP blocks gradually mimic their FP counterparts, executing an implicit regularization for FP representation. We validate the proposed framework based on uniform quantization [8], and demonstrate its improved performance over the state-of-the-art methods for low-bit network quantization.

Our contributions can be summarized as follows:

- We propose a neat framework for QAT that makes complete use of a pre-trained, full-precision model. Our approach is capable of producing high-performance, low-bit quantization models, without increasing the model complexity at the inference phase.
- We evaluate the proposed framework through empirical observations, and uncover the framework’s implicit insights in conjunction with our concerns and hypotheses.
- We demonstrate the effectiveness of our method under low-precision of 4-, 3-, and 4-bit widths, achieving state-of-the-art results on ImageNet.

II. RELATED WORK

Network Quantization. [14] presents a detailed description of CNN quantization and divides quantizer designs into two categories: post-training quantization (PTQ) and quantization-aware training (QAT). PTQ techniques often quantize a network without complete training of full data [1, 6, 9, 21, 38, 57]. QAT techniques commonly outperform PTQ for low-precision networks with proper training on original data [3, 5, 8, 12, 23, 25, 59]. However, non-idealities appear in both forward and backward owing to the discrete quantizers used by the QAT during training. Previous QAT approaches have been investigated to optimize the quantization parameters (e.g., clipping value, quantization step size) [3, 8] and introduce non-uniform quantization [33, 34, 37, 51, 58] in order to improve the representation of low-precision and effectively balance quantization error. Another line of QAT studies has attempted to tackle the gradient mismatch issue induced by

the non-differentiable quantizer during QAT [31, 61].

Auxiliary supervision. Auxiliary supervision seeks to improve network training by integrating auxiliary objectives and extra losses into intermediate layers, hence combating the gradient issue and providing regularization [36, 40, 48, 56]. [61] was the first to suggest using full-precision branches to address the gradient problem in QAT, demonstrating the practicality of auxiliary supervision in QAT. Knowledge distillation is a common approach for network training [19], in which smaller networks are trained by transferring features from stronger teacher models, and it may also be thought of as auxiliary supervision. [42] were the first to propose incorporating distillation into QAT, with the full-precision model serving as the teacher model and distillation losses established on the low-precision model’s output layer for auxiliary supervision. Following research, such as QKD [26] and QFD [60], has enhanced the design of distillation branches for auxiliary supervision in QAT. We develop full-precision branches for auxiliary supervision in the proposed quantization framework by grafting the full-precision model to better address the quantizer’s gradient mismatch. Notably, our framework employs KD loss as one of the auxiliary learning sources based only on pre-trained FP counterparts to aid optimization, removing the need to train an extra teacher network. On this premise, we continue to outperform KD approaches without the assistance of large teachers.

III. APPROACH

As shown in Fig. 1 (b), we propose a general technique for enhancing quantized-aware training (QAT) in this section: the block-wise replacement framework (BWRP). Initially, in accordance with depth and resolution, the identical network structure of the low-precision (LP) model and its full-precision (FP) counterpart is uniformly partitioned into multiple block-wise sections. For instance, the block sections of the ResNets family correspond to the code implementation of res-blocks. Grafted mixed-precision (MP) models are obtained by progressively replacing the last low-precision blocks with pre-trained full-precision blocks. Throughout the training phase, only the low-precision component of the mixed-precision model is modified; the full-precision component remains constant. The parameters contained within individual LP blocks are synchronized and shared. Mixed-precision models provide full-precision branches that serve as auxiliary supervision for the low-precision backbone, which are neglected during the inference process. The primary objective of developing full-precision branches is to enhance the utilization of full-precision counterparts for guidance and to mitigate issues associated with limited representation in the forward pass and gradient estimation in the backward pass.

A. Quantization-aware Training

The fundamental concept underlying Quantization-aware Training (QAT) is to incorporate quantization simulation during network training, thus enabling the direct determination of the model’s optimal solution while subject to quantization constraints. Floating-point activation and weights are converted

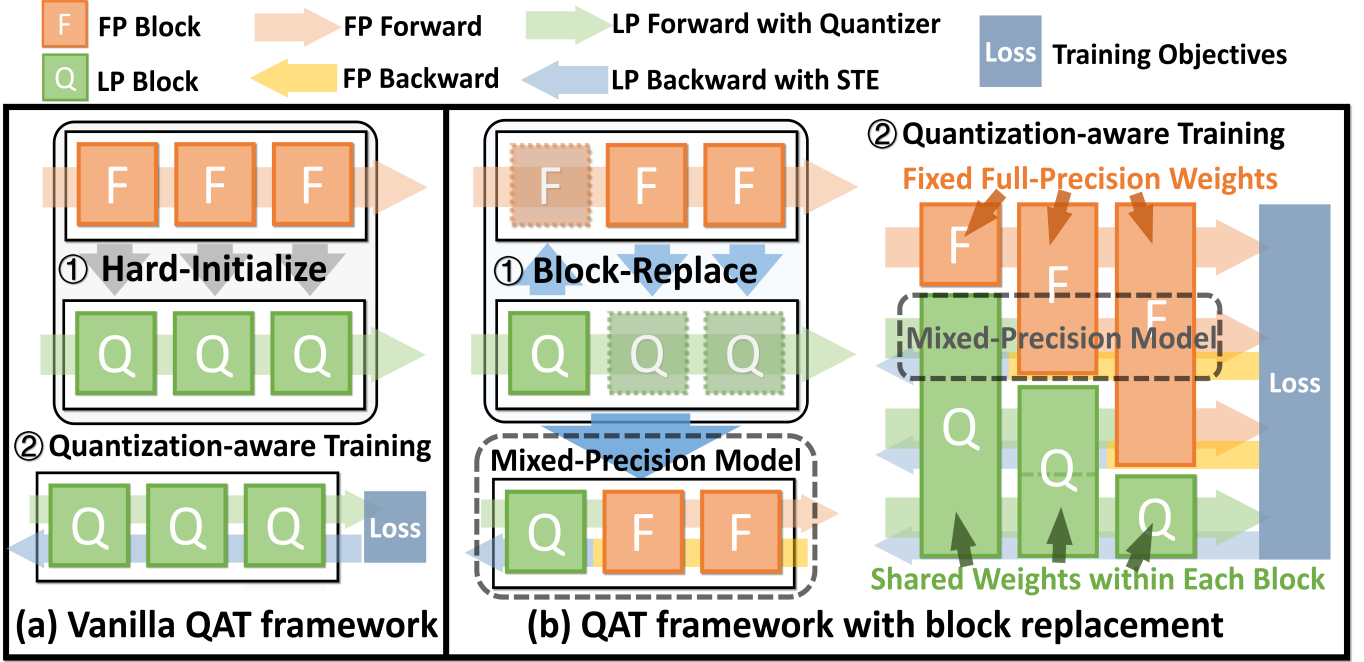


Fig. 1. Framework Overview. (a) The fundamental implementation of Quantization-aware Training (QAT) in which weight initialization is performed using full-precision counterparts. (b) The proposed block-wise replacement framework (BWRF) generates mixed-precision models during the training phase, employing full-precision counterparts for auxiliary supervision.

to a restricted low-bits representation using quantizers during the forward pass. The quantizer, denoted as $q(\cdot)$, accepts an input vector v and generates a pseudo-quantized output \hat{v} , as follows:

$$\hat{v} = q(v) = s \cdot \lfloor \text{clip}(\frac{v}{s}, N, P) \rfloor \quad (1)$$

where notation $\lfloor \cdot \rfloor$ denotes the round-to-nearest operator, s is the scaling factor, $\text{clip}(\cdot, N, P)$ represents the clipping function whose lower and upper quantization thresholds N and P , respectively.

During the backward pass, QAT encounters the non-differentiability of the round function in Eq. 1, which restricts gradient-based training. A commonly used approach to address the issue is to approximate the true gradient for quantizers using the straight-through estimator (STE) [2, 18], which simply approximates the gradient of round function as:

$$\frac{\partial \hat{v}}{\partial x} (\lfloor x \rfloor) = 1 \quad (2)$$

Then, the gradient passing quantizers can be obtained:

$$\frac{\partial \hat{v}}{\partial v} = \begin{cases} 1, & \text{if } N < v < P; \\ 0, & \text{otherwise.} \end{cases} \quad (3)$$

We follow LSQ [8] and APoT [33] in practice, which allow scale parameters of quantizers to be learned; the gradient of s can be derived from Eq. 1 and Eq. 3:

$$\frac{\partial \hat{v}}{\partial s} = \begin{cases} -v/s + \lfloor v/s \rfloor, & \text{if } N < v/s < P; \\ N, & \text{if } v/s \leq N; \\ P, & \text{if } v/s \geq P. \end{cases} \quad (4)$$

In accordance with the standard configuration of QAT, the quantizers for activations and weights are established prior to

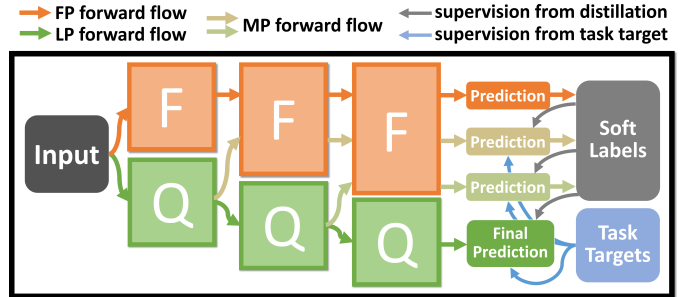


Fig. 2. Implementation of BWRF Training. Mixed-precision models are implemented implicitly through the utilization of overlapping LP forward flows. Both task targets and model predictions are regarded as loss sources for training.

the linear operators (e.g., convolutional and fully-connected layers) to ensure that matrix multiplication can occur in fixed-point domains [5, 8, 33]. During code implementation, the functionalities of quantizers are encapsulated within linear operators to facilitate their efficient integration into the structure of the corresponding FP model.

B. Framework Formulation

Consider the entire definition of the framework. Given the LP model Q and its corresponding FP model F ,

$$Q = \{Q_1, Q_2, \dots, Q_n\}, F = \{F_1, F_2, \dots, F_n\} \quad (5)$$

where n represents the quantity of blocks, Q_i and F_i denote the i -th block that is separated according to the image resolutions. During training, the weights of the FP model remain

Algorithm 1: BWRP Training w.r.t the low-precision backbone network Q and the full-precision counterpart network F .

Input: Mini-batch $\{x, y\}$; trainable weights $\{\theta_i\}_{i \leq n}$ within the low-precision blocks $\{Q_i\}_{i \leq n}$; the full-precision blocks $\{F_i\}_{i \leq n}$.

Output: Updated weights $\{\theta_i\}_{i \leq n}$.

- 1 Compute the FP output $y_F = F(x)$;
 - 2 Initialize $x_{Q_0} = x$;
 - 3 Continuously compute the intermediate features via the LP blocks $\{x_{Q_i} | x_{Q_i} = Q_i(x_{Q_{i-1}}; \theta_i)\}_{i \in [1, n]}$;
 - 4 Obtain the LP output $y_Q = x_{Q_n}$;
 - 5 **for** $k = 1$ to $n - 1$ **do**
 - 6 Initialize $x_{M^k} = x_{Q_k}$;
 - 7 Continuously compute the intermediate features via the FP blocks $\{x_{M_i^k} | x_{M_i^k} = F_i(x_{M_{i-1}^k})\}_{i \in [k+1, n]}$;
 - 8 Obtain the MP output $y_{M^k} = x_{M_n^k}$
 - 9 **end**
 - 10 Compute the loss L according to Eq. 8 - Eq. 11 ;
 - 11 Propagate the gradients for LP weights

$$\frac{\partial L}{\partial \theta_i} = \frac{\partial L}{\partial x_{Q_i}} \frac{\partial x_{Q_i}}{\partial \theta_i} = \left(\frac{\partial L}{\partial x_{Q_{i+1}}} \frac{\partial x_{Q_{i+1}}}{\partial x_{Q_i}} + \frac{\partial L}{\partial x_{M_{i+1}^k}} \frac{\partial x_{M_{i+1}^k}}{\partial x_{Q_i}} \right) \frac{\partial x_{Q_i}}{\partial \theta_i};$$
 - 12 Update $\{\theta_i\}_{i \leq n}$ based on gradients.
-

constant, whereas the weights of the LP model, represented by $\{\theta_i\}_{i \leq n}$, are trainable. We define the k -th mixed-precision (MP) model with implementing the block replacement starting from k -th blocks (between $\{Q_i\}_{i > k}$ and $\{F_i\}_{i > k}$), as:

$$M^k = \{Q_1, \dots, Q_k, F_{k+1}, \dots, F_n\} \quad (6)$$

where M^k comprises the first k blocks of LP models with trainable weights $\{\theta_i\}_{i \leq k}$. It is notable to remark that the implementation of MP within the framework is implicit. As depicted in Fig. 2, in the computing graph, the LP forward flow is reused for MP flows in consideration of the calculation consistency of the M^k and Q in the first k blocks.

C. Training Objectives

Given the input x with task target y , the prediction of each branch is computed as:

$$\begin{aligned} y_Q &= Q(x; \{\theta_i\}_{i \leq n}) \\ y_{M^k} &= M^k(x; \{\theta_i\}_{i \leq k}) \\ y_F &= F(x) \end{aligned} \quad (7)$$

The LP model's output denoted as y_Q , represents the final training result of our framework. y_F and y_{M^k} are employed as soft labels during QAT training to aid in directing the training of LP blocks. As depicted in Fig. 2, the training targets are divided into two components: task targets y , and soft labels y_F, y_{M^k} . The training objective designated as L_{target} , which

pertains to the classification target, is applied uniformly to both Q and M :

$$L_{target} = L_{ce}(y_Q, y) + \sum_{k=1}^{k < n} \alpha_k \cdot L_{ce}(y_{M^k}, y) \quad (8)$$

where L_{ce} represents the cross-entropy for the classification task, α_k are the hyper-parameters utilized to balance the losses of FP branches. MP models are expected to perform the same functions as the LP backbone, wherein the alignment between the end LP blocks and their FP branches (e.g., $\{Q_{k+1}, \dots, Q_n\}$ and $\{F_{k+1}, \dots, F_n\}$) is implicitly applied, and the first LP blocks (e.g., $\{Q_1, \dots, Q_k\}$) are necessary to extract functional features applicable to both the FP and LP branches. The knowledge distillation objectives, represented as $L_{distill}$, are designed to enhance the performance of the LP and MP models. This is achieved through the utilization of soft label objectives, which are determined by the KL distance, L_{kd} , in accordance with the vanilla logit-based KD method [19]:

$$\begin{aligned} L_{distill} &= L_{kd}(y_Q, y_F) + L_{kd}(y_Q, y_{M_{n-1}^{avg}}) \\ &+ \sum_{k=1}^{k < n} \alpha_k \left(L_{kd}(y_{M^k}, y_F) + L_{kd}(y_{M^k}, y_{M_{k-1}^{avg}}) \right) \end{aligned} \quad (9)$$

where the $y_{M^k}^{avg}$ reflects the ensemble of preceding MP models and is computed using the average of the previous prediction:

$$y_{M^k}^{avg} = \frac{1}{k+1} \left(y_F + \sum_{j=1}^{j \leq k} y_{M^j} \right) \quad (10)$$

Self-distillation [54] is the underlying principle of these designs, which facilitates the constructive mutual influence of the various self-components and is supported by empirical results.

The framework is then trained under the integrated objectives, as shown in Fig. 2:

$$L = L_{target} + L_{distill} \quad (11)$$

IV. EXPERIMENTS

In this section, we evaluate the effectiveness of the proposed framework compared to alternative QAT methods. We focus on quantization performance under the low-bit widths of 4/3/2. In this case, the issue of performance degradation resulting from quantizers becomes apparent and deserves particular consideration. The primary outcomes of the proposed method, along with a comparison to alternative QAT methods, will be plainly outlined in Section 4.2. A comprehensive description of the ablation studies for the method will be provided in Section 4.3. In Section 4.4, the additional influence of the proposed method on QAT training will be analyzed by incorporating the visualization results.

A. Experimental Setup

Datasets and Networks. We perform experiments on two standard image classification datasets: ImageNet [46] and CIFAR-10 [28]. ImageNet contains about 1.2 million training and 50K validation images of 1,000 object categories.

TABLE I
 COMPARISON OF VALIDATION ACCURACY ON IMAGENET USING THE RESNET-18 ARCHITECTURE. THE WEIGHTS AND ACTIVATIONS ARE QUANTIZED TO THE SAME PRECISION THROUGHOUT 4-, 3-, AND 2-BITS. QAT METHODS ARE CATEGORIZED BASED ON UNIFORM AND NON-UNIFORM QUANTIZATION MANNERS. LSQ[†] DENOTES REPLICATED BENCHMARKS IN ACCORDANCE WITH OUR EXPERIMENTAL SETTING FOR UNBIASED COMPARISONS. THE PERFORMANCE GAIN ACHIEVED WITH BWRF IS DENOTED BY Δ .

Network	Manner	Method	4 bits		3 bits		2 bits	
			Top1	Top5	Top1	Top5	Top1	Top5
ResNet18 FP(72.26)	Non-uniform	LQ-Nets	69.3	88.8	68.2	87.9	64.9	85.9
		QIL	70.1	-	69.2	-	65.7	-
		DAQ	70.5	-	69.6	-	66.9	-
		APOT	70.7	89.6	69.9	89.2	67.3	87.5
		LCQ	71.5	-	70.6	-	68.9	-
	Uniform	PACT	69.2	89	68.1	88.2	64.4	85.6
		DoReFa-Net	68.1	88.1	67.5	87.6	64.7	84.4
		DSQ	69.6	88.9	68.7	-	65.2	-
		LSQ	71.1	90.0	70.2	89.4	67.6	87.6
		LSQ+	70.8	-	69.3	-	66.8	-
		EWGS	70.6	-	69.7	-	67.0	-
		BR	70.8	89.6	69.9	89.1	67.2	87.3
		QKD	71.4	90.3	70.2	89.9	67.4	87.5
		QFD	71.1	89.8	70.3	89.4	67.6	87.8
		LSQ [†]	70.4	89.4	69.3	88.8	65.3	86.2
		BWRF (ours)	71.9	90.5	70.8	89.9	67.7	87.9
Δ	+1.5	+1.1	+1.5	+1.1	+2.4	+1.7		

CIFAR10 contains 60,000 32x32 color images in 10 classes, with 6,000 images per class. In our experiments, images from the ImageNet training set were randomly cropped to 224x224 pixels and randomly flipped for data augmentation. For the validation set, images were first resized to 256x256 pixels, then a 224x224 center crop was taken to evaluate consistent image sizes. All images in the dataset underwent normalization. Similar preprocessing was done for CIFAR-10 as on ImageNet. We evaluated our method in the Reset series [16], including ResNet18, ResNet20, ResNet34, ResNet50 and ResNet56. Among them, ResNet18, ResNet34 and ResNet50 were used for ImageNet, while ResNet20 and ResNet56 were used for CIFAR10. We inserted quantizers for weights or activations before convolution operators in each layer to perform quantization. 8-bit quantization was applied in the first and last layers.

Quantization Setting. In accordance with the configurations established in prior QAT methods[14, 21, 31, 33], we implemented quantizers for activations and weights prior to convolution operators in each layer to simulate quantization. While the final fully-connected layer and the initial convolutional layer were both quantized to 8 bits, all other convolutional layers were quantized to a consistent low-bit width. The weights of the low-precision network are consistently initialized using the pre-trained, full-precision counterparts. Training details are appended in the supplement.

Training Details. We use SGD with a momentum of 0.9 for all cases. On imageNet, the batch size was 1024 for all bit

widths for ResNet18. For 4-bit quantization of ResNet18, the initial learning rate was set to 1e-2 with a weight decay of 1e-4. For 3-bit ResNet18, the settings were the same. For 2-bit ResNet18, the learning rate was set to 4e-2 and weight decay to 2e-5. For ResNet34, the batch size was 1024 for all except 2-bit. The 4-bit quantization used a learning rate of 2e-2 and weight decay of 1e-4. The 3-bit settings used a learning rate of 1e-2 and weight decay of 1e-4. The 2-bit ResNet34 used a learning rate of 2e-2, batch size of 512, and weight decay of 2e-5. For ResNet50, the batch size was 512 for all bit widths. The 4-bit quantization used a learning rate of 2e-2 and weight decay of 1e-4. The 3-bit settings used a learning rate of 1e-2 and weight decay of 1e-4. The 2-bit ResNet50 used a learning rate of 2e-2 and weight decay of 2e-5. Models were all trained for 120 epochs with the learning rate divided by 10 every 30 epochs. For FP counterparts, we trained ResNet18 with knowledge distillation and applied pre-trained ResNet-34 and ResNet-50 from TIMM [50]. On CIFAR-10, the initial learning rate was set consistently to 4e-2 for 2, 3, and 4-bit levels with a weight decay of 1e-4 for 3/4-bit and 3e-5 for 2-bit. Models were trained for 300 epochs with the multistep schedule, dividing the learning rate by 10 at the 150th and 225th epochs.

B. Main Results

Results on ImageNet. The performance evaluation of the proposed methods is presented on the ImageNet dataset, utilizing the ResNet-18, ResNet-34, and ResNet-50 architectures,

TABLE II
COMPARISON OF VALIDATION ACCURACY ON IMAGENET USING THE RESNET-34 AND RESNET-50 ARCHITECTURES. THE WEIGHTS AND ACTIVATIONS ARE QUANTIZED TO THE SAME PRECISION THROUGHOUT 4-, 3-, AND 2-BITS. QAT METHODS ARE CATEGORIZED BASED ON UNIFORM AND NON-UNIFORM QUANTIZATION MANNERS.

Network	Manner	Method	4 bits		3 bits		2 bits	
			Top1	Top5	Top1	Top5	Top1	Top5
ResNet34 FP(76.32)	Non-uniform	LQ-Nets	-	-	71.9	90.2	69.8	89.1
		QIL	73.7	-	73.1	-	70.6	-
		APOT	73.8	91.6	73.4	91.1	70.9	89.7
		LCQ	74.3	-	74.0	-	72.7	-
	Uniform	DSQ	72.8	-	72.5	-	70.0	-
		LSQ	74.1	91.7	73.4	91.4	71.6	90.3
		EWGS	73.9	-	73.3	-	71.4	-
		QKD	74.6	92.1	73.9	91.4	71.6	90.3
		QFD	74.7	92.3	73.9	91.7	71.7	90.4
	BWRf (ours)		75.9	92.6	74.3	91.7	71.7	90.4
ResNet50 FP(80.11)	Non-uniform	LQ-Nets	75.1	92.4	74.2	91.6	71.5	90.3
		APOT	76.6	93.1	75.8	92.7	73.4	91.4
		LCQ	76.6	-	76.3	-	75.1	-
	Uniform	PACT	76.5	93.2	75.3	92.6	72.2	90.5
		DeReFa-Nets	71.4	89.8	69.9	89.2	67.1	87.3
		LSQ	76.7	93.2	75.8	92.7	73.7	91.5
		QKD	77.3	93.6	76.4	93.2	73.9	91.6
	BWRf (ours)		79.0	94.2	77.8	93.5	74.2	91.6

respectively. The obtained results are summarized in Tab. I and Tab. II, which show the performance of our BWRf in the presence of low precision bit widths. "LSQ[†]" denotes the outcomes achieved via the preliminary implementation outlined in section 3.1, while combined with the configurations specified in section 4.1. These results formed the baseline for the subsequent developments incorporated into our framework. Our proposed framework consistently outperforms other QAT methods with respect to R18, R34, and R50. Significantly, our framework employing uniform quantization surpasses the state-of-the-art results of non-uniform quantization techniques in almost all cases. This indicates that the framework offers a dependable approach to achieve uniform quantization, which is particularly well-suited for inferring fixed-points computation.

Results on CIFAR-10. Tab. III presents the performance on CIFAR-10, showcasing a consistent enhancement in comparison to alternative QAT methods. The experimental findings suggest that BWRf can effectively convert the FP model to 3/4 bit with minimal performance degradation. It is worth mentioning that the BWRf model exhibits superior performance when operating at 4 bits compared to the FP model. This can be attributed to the utilization of mechanisms resembling self-supervision and auxiliary supervision, which will be elaborated upon in Section 4.3 following ablation experiments.

TABLE III
COMPARISON OF VALIDATION ACCURACY ON CIFAR-10 USING RESNET-20 AND RESNET-56. LSQ[†] DENOTES REPLICATED BENCHMARKS IN ACCORDANCE WITH OUR EXPERIMENTAL SETTING FOR UNBIASED COMPARISONS.

Network	Method	Accuracy (%)		
		4 bits	3 bits	2 bits
ResNet20 FP(92.96)	DoReFa-Net	90.5	89.9	88.2
	PACT	91.7	91.1	89.7
	LQ-Net	-	91.6	90.2
	PACT+SAWB+fp _{sc}	-	-	90.5
	QKD	93.1	92.7	90.5
	APoT	92.3	92.2	91.0
	LSQ [†]	92.67	92.34	90.74
BWRf (ours)		93.13	92.76	90.84
ResNet56 FP(94.46)	PACT+SAWB+fp _{sc}	-	-	92.5
	APoT	94.0	93.9	92.9
	LSQ [†]	93.74	93.69	92.84
	BWRf (ours)	94.69	94.31	93.03

C. Ablation Study

Ablation of training objectives. The results of the ablation experiment regarding the training objectives of BWRf are presented in Tab. IV. The foundational implementation of

TABLE IV
ABLATION STUDY OF TRAINING OBJECTIVES. THE 4-BIT RESULTS ARE OBTAINED USING RESNET-20 AND RESNET-56 ON CIAFR-10.

y_Q	y_{M_k}	y_Q	y_{M_k}	w/ y^{avg}	Accuracy (%)	
w/ target y		w/ FP label y_F			R20	R56
✓					92.67	93.74
✓	✓				92.89	94.02
✓		✓			92.91	94.23
✓	✓	✓			93.02	94.35
✓	✓	✓	✓		93.09	94.54
✓	✓	✓	✓	✓	93.13	94.69

TABLE V
ABLATION STUDY OF PRUNING MP MODELS. THE 4-BIT RESULTS ARE OBTAINED USING RESNET-20 AND RESNET-56 ON CIAFR-10.

Network	base	w/ M^1	w/ M^2	w/ M^1 & M^2
ResNet20	92.67	92.97	93.02	93.13
ResNet56	93.74	94.42	94.56	94.69

QAT is training exclusively with $L_{ce}(y_Q, y)$; training with both $L_{kd}(y_Q, y_F)$ and $L_{ce}(y_Q, y)$ is analogous to distilling LP using the FP model. We gradually integrate training objectives posterior to section 3.2 as auxiliary supervision. The overall performance continues to improve with the incorporation of additional loss sources, indicating that these loss sources are effective for the auxiliary supervision of QAT. It is important to highlight that the structure implemented utilizing $L_{ce}(y_{M_k}, y)$ is similar to the methodology employed to implement multi-exit for classification in DSN [30]. Nevertheless, the branch that BWRf employs remains consistent, and the FP model verifies that the branch utilized as the classifier is rational. Furthermore, the implementation carried out utilizing MP outputs may be regarded as a unique form of self-supervised methodology [54] that integrates untrained branches for the purpose of extracting intermediate features; the effectiveness of this approach is assessed via empirical experiments. Overall, the integration of these auxiliary training objectives improves the quality of the QAT training framework and can be further integrated to yield better outcomes.

Ablation of MP branches. The outcomes of ablation for MP branches are presented in Tab. V, which illustrates the beneficial effects of setting MP branches within the framework. As the branches are pruned, the performance returns to its initial state. It is noteworthy to mention that M^2 provides enhanced auxiliary supervision when compared to M^1 , which implies

TABLE VI
EXTENSION TO NON-UNIFORM QUANTIZATION. THE RESULTS ARE OBTAINED USING RESNET-20 AND RESNET-56 ON CIAFR-10.

Network	Manner	4 bits	3 bits	2 bits
ResNet20	uniform	93.13	92.76	90.84
	Non-uniform	93.16	92.82	91.06
ResNet56	uniform	94.69	94.31	93.03
	Non-uniform	94.74	94.38	93.14

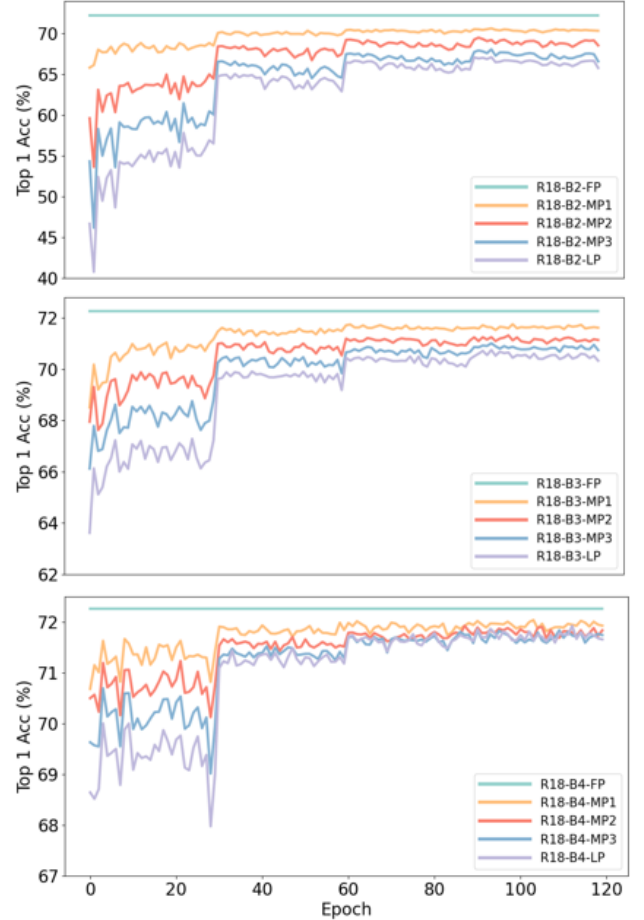


Fig. 3. Validation accuracy of mixed-precision models and low-precision backbone during training. The results are obtained by ResNet-18 on ImageNet.

that the replacement between intermediate blocks provides more reliable support.

D. Analysis and Discussion

Extension to non-uniform quantization. The adaptability of BWRf to the non-uniform design was assessed, as illustrated in Tab. VI. It is not unexpected that performance can be marginally enhanced by combining BWRf with a non-uniform quantizer, considering that non-uniform quantization offers a more reliable representation at lower levels of bit widths. Nevertheless, upon comparison, the distinction between uniform and non-uniform is not particularly significant. For increased adaptability, we suggest implementing uniform quant in accordance with BWRf training.

Analysis of MP performance. An illustration of the performance trajectory of intermediate mixed-precision models throughout the training procedure can be found in Fig. 3. MP generally exhibits performance that lies between that of FP and LP, with a gradual deterioration noted as the number of LP blocks increases. This observation aligns with the expectations that were held. LP inevitably encounters a decrease in performance due to the representation capacity

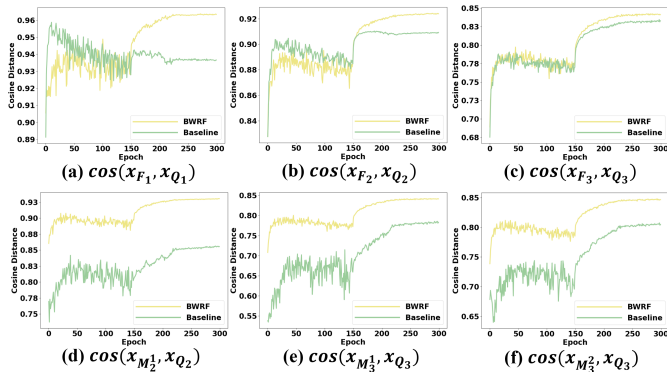


Fig. 4. **Measures of feature similarity.** The results of cosine distances are obtained by ResNet-18 on CIFAR-10 under 4-bit quantization.

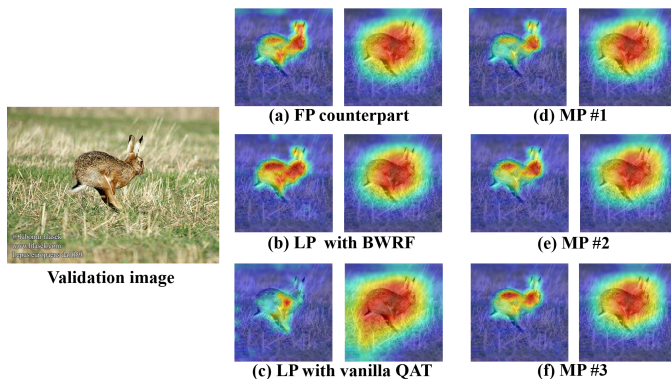


Fig. 5. **Visualization results of class activation mapping (CAM).** Two visualized targets are established on the third and fourth blocks’ output layers. (a-b) The results of the full-precision model F and low-precision model Q trained with BWRF. (c) The results of the vanilla model with baseline implementation. (d-f) The results of mixed-precision models M^1 , M^2 , M^3 , respectively.

limitation that occurs during the phased block replacements from the FP block to the LP block. This phenomenon is at its most conspicuous at two and three bits. It is worth mentioning that the differentiation between the branches of the 4-bit model does not become immediately evident during the late training phase. The discrepancy between the FP and the first MP appears to be the most significant obstacle. Therefore, to enhance the accuracy of ancillary supervision, an intuitive optimization technique involves further dividing the first LP block into several subblocks.

Analysis on distance between intermediate features. As shown in Fig. 4, to obtain insights regarding the comparison between quantized and full-precision features in the QAT process, the cosine distance between intermediate features of LP and MP/FP is computed throughout the training phase. A comparative analysis of the situation is performed utilizing the initial LSQ method. The degree of alignment between the features of FP and LP is illustrated in Fig. 4 (a)-(c), which depicts the direct correspondence between the two designs. The cosine distance of the LP model undergoing BWRF training approaches 1 during the late training phases, indicating a more optimal fit with FP. This finding provides empirical validation for the hypothesis that the auxiliary supervision

implemented in the BWRF configuration can assist the LP model in incorporating the previous FP representation into the LP features. The distance relationship between LP and MP features is illustrated in Fig. 4 (d)-(f). This relationship can be used to determine the degree of block-wise similarity between FP and LP, thereby offering a more comprehensive view of the fitting situation between FP and LP. As an example, Fig. 4 (d) illustrates the value of $\cos(x_{Q_2}, x_{M_2^1})$, which quantifies the separation between $Q_2(x_{Q_1})$ and $F_2(x_{Q_1})$ in order to represent the matching distance between Q_2 and F_2 with respect to a given Q_1 block. In block-wise distance detection, BWRF is invariably greater than the baseline method; therefore, as expected, fitting LP and FP via block replacement is valid. Additionally, it is critical to note that the model prioritizes the fitting of block-wise representations Fig. 4 (d)-(f) when the block replacement setting is enabled. Consequently, the initial depiction of the overall features of FP may be marginally delayed Fig. 4 (a)-(b); however, it possesses the capability to converge towards higher values gradually.

Visualization with class activation mapping. Fig. 5 shows the Class Activation Mapping (CAM) generated by the Layer-CAM method [24]. The discovery that the LP, when applied to BWRF, yields outcomes similar to those of the FP counterpart, indicates that the LP’s feature extraction and gradient backpropagation computation align with those of the FP, thus validating the LP’s ability to effectively convey the FP representation. The progressive step results from FP to LP are displayed across MPs’ results, illustrating that continuous and logical intermediate step models can be derived from FP to LP through the sequential block-wise replacement process.

V. CONCLUSION

In this paper, we propose the block-wise replacement framework (BWRF) for quantization-aware training (QAT) by applying auxiliary supervision for low-bit models via the FP counterpart in order to mitigate the representation limitation and gradient mismatch issues that arise from the incorporation of discrete quantizers. We analyze and discuss the latent insights of the proposed framework through empirical experiments. To be emphasized, we obtain state-of-the-art results with uniform quantization settings for low-bit widths of 4/3/2, and we even outperform non-uniform methods in most instances. The framework itself is neat and flexible, necessitating only a concise wrapper for code implementations. We believe that BWRF can facilitate network quantization by providing extensions to QAT methods.

REFERENCES

- [1] Ron Banner, Yury Nahshan, and Daniel Soudry. Post training 4-bit quantization of convolutional networks for rapid-deployment. *Advances in Neural Information Processing Systems*, 32, 2019.
- [2] Yoshua Bengio, Nicholas Léonard, and Aaron Courville. Estimating or propagating gradients through stochastic neurons for conditional computation. *arXiv preprint arXiv:1308.3432*, 2013.

- [3] Yash Bhalgat, Jinwon Lee, Markus Nagel, Tijmen Blankevoort, and Nojun Kwak. Lsq+: Improving low-bit quantization through learnable offsets and better initialization. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*, pages 696–697, 2020.
- [4] Yoonho Boo, Sungho Shin, Jungwook Choi, and Wonyong Sung. Stochastic precision ensemble: self-knowledge distillation for quantized deep neural networks. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pages 6794–6802, 2021.
- [5] Jungwook Choi, Zhuo Wang, Swagath Venkataramani, Pierce I-Jen Chuang, Vijayalakshmi Srinivasan, and Kailash Gopalakrishnan. Pact: Parameterized clipping activation for quantized neural networks. *arXiv preprint arXiv:1805.06085*, 2018.
- [6] Yoni Choukroun, Eli Kravchik, Fan Yang, and Pavel Kisilev. Low-bit quantization of neural networks for efficient inference. In *2019 IEEE/CVF International Conference on Computer Vision Workshop (ICCVW)*, pages 3009–3018. IEEE, 2019.
- [7] Jiequan Cui, Pengguang Chen, Ruiyu Li, Shu Liu, Xiaoyong Shen, and Jiaya Jia. Fast and practical neural architecture search. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 6509–6518, 2019.
- [8] Steven K Esser, Jeffrey L McKinstry, Deepika Bablani, Rathinakumar Appuswamy, and Dharmendra S Modha. Learned step size quantization. In *International Conference on Learning Representations*, 2020.
- [9] Alexander Finkelstein, Uri Almog, and Mark Grobman. Fighting quantization bias with bias. *arXiv preprint arXiv:1906.03193*, 2019.
- [10] Jonathan Frankle and Michael Carbin. The lottery ticket hypothesis: Finding sparse, trainable neural networks. *arXiv preprint arXiv:1803.03635*, 2018.
- [11] Ross Girshick, Jeff Donahue, Trevor Darrell, and Jitendra Malik. Rich feature hierarchies for accurate object detection and semantic segmentation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 580–587, 2014.
- [12] 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. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 4852–4861, 2019.
- [13] Yiwen Guo, Anbang Yao, Hao Zhao, and Yurong Chen. Network sketching: Exploiting binary structure in deep cnns. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 5955–5963, 2017.
- [14] Tiantian Han, Dong Li, Ji Liu, Lu Tian, and Yi Shan. Improving low-precision network quantization via bin regularization. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 5261–5270, 2021.
- [15] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.
- [16] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.
- [17] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. Mask r-cnn. In *Proceedings of the IEEE international conference on computer vision*, pages 2961–2969, 2017.
- [18] Geoffrey Hinton, Nitish Srivastava, and Kevin Swersky. Neural networks for machine learning lecture 6a overview of mini-batch gradient descent. *Cited on*, 14 (8):2, 2012.
- [19] Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. Distilling the knowledge in a neural network. *arXiv preprint arXiv:1503.02531*, 2015.
- [20] Andrew Howard, Mark Sandler, Grace Chu, Liang-Chieh Chen, Bo Chen, Mingxing Tan, Weijun Wang, Yukun Zhu, Ruoming Pang, Vijay Vasudevan, et al. Searching for mobilenetv3. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 1314–1324, 2019.
- [21] Itay Hubara, Yury Nahshan, Yair Hanani, Ron Banner, and Daniel Soudry. Improving post training neural quantization: Layer-wise calibration and integer programming. *arXiv preprint arXiv:2006.10518*, 2020.
- [22] Benoit Jacob, Skirmantas Kligys, Bo Chen, Menglong Zhu, Matthew Tang, Andrew Howard, Hartwig Adam, and Dmitry Kalenichenko. Quantization and training of neural networks for efficient integer-arithmetic-only inference. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2704–2713, 2018.
- [23] Sambhav Jain, Albert Gural, Michael Wu, and Chris Dick. Trained quantization thresholds for accurate and efficient fixed-point inference of deep neural networks. *Proceedings of Machine Learning and Systems*, 2:112–128, 2020.
- [24] Peng-Tao Jiang, Chang-Bin Zhang, Qibin Hou, Ming-Ming Cheng, and Yunchao Wei. Layercam: Exploring hierarchical class activation maps for localization. *IEEE Transactions on Image Processing*, 30:5875–5888, 2021.
- [25] 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. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4350–4359, 2019.
- [26] Jangho Kim, Yash Bhalgat, Jinwon Lee, Chirag Patel, and Nojun Kwak. Qkd: Quantization-aware knowledge distillation. *arXiv preprint arXiv:1911.12491*, 2019.
- [27] Raghuraman Krishnamoorthi. Quantizing deep convolutional networks for efficient inference: A whitepaper. *arXiv preprint arXiv:1806.08342*, 2018.
- [28] Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. 2009.

- [29] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. *Advances in neural information processing systems*, 25, 2012.
- [30] Chen-Yu Lee, Saining Xie, Patrick Gallagher, Zhengyou Zhang, and Zhuowen Tu. Deeply-supervised nets. In *Artificial intelligence and statistics*, pages 562–570. Pmlr, 2015.
- [31] Junghyup Lee, Dohyung Kim, and Bumsu Ham. Network quantization with element-wise gradient scaling. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 6448–6457, 2021.
- [32] Hao Li, Asim Kadav, Igor Durdanovic, Hanan Samet, and Hans Peter Graf. Pruning filters for efficient convnets. *arXiv preprint arXiv:1608.08710*, 2016.
- [33] Yuhang Li, Xin Dong, and Wei Wang. Additive powers-of-two quantization: An efficient non-uniform discretization for neural networks. *arXiv preprint arXiv:1909.13144*, 2019.
- [34] Zechun Liu, Kwang-Ting Cheng, Dong Huang, Eric P Xing, and Zhiqiang Shen. Nonuniform-to-uniform quantization: Towards accurate quantization via generalized straight-through estimation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4942–4952, 2022.
- [35] Zhuang Liu, Mingjie Sun, Tinghui Zhou, Gao Huang, and Trevor Darrell. Rethinking the value of network pruning. *arXiv preprint arXiv:1810.05270*, 2018.
- [36] Christos Louizos, Matthias Reisser, Tijmen Blankevoort, Efstratios Gavves, and Max Welling. Relaxed quantization for discretized neural networks. *arXiv preprint arXiv:1810.01875*, 2018.
- [37] Daisuke Miyashita, Edward H Lee, and Boris Murmann. Convolutional neural networks using logarithmic data representation. *arXiv preprint arXiv:1603.01025*, 2016.
- [38] Markus Nagel, Mart van Baalen, Tijmen Blankevoort, and Max Welling. Data-free quantization through weight equalization and bias correction. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 1325–1334, 2019.
- [39] Yury Nahshan, Brian Chmiel, Chaim Baskin, Evgenii Zheltonozhskii, Ron Banner, Alex M Bronstein, and Avi Mendelson. Loss aware post-training quantization. *Machine Learning*, 110(11-12):3245–3262, 2021.
- [40] Vladimir Nekrasov, Hao Chen, Chunhua Shen, and Ian Reid. Fast neural architecture search of compact semantic segmentation models via auxiliary cells. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 9126–9135, 2019.
- [41] Eunhyeok Park and Sungjoo Yoo. Profit: A novel training method for sub-4-bit mobilenet models. In *Computer Vision—ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part VI 16*, pages 430–446. Springer, 2020.
- [42] Antonio Polino, Razvan Pascanu, and Dan Alistarh. Model compression via distillation and quantization. *arXiv preprint arXiv:1802.05668*, 2018.
- [43] Joseph Redmon, Santosh Divvala, Ross Girshick, and Ali Farhadi. You only look once: Unified, real-time object detection. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 779–788, 2016.
- [44] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster r-cnn: Towards real-time object detection with region proposal networks. *Advances in neural information processing systems*, 28, 2015.
- [45] Adriana Romero, Nicolas Ballas, Samira Ebrahimi Kahou, Antoine Chassang, Carlo Gatta, and Yoshua Bengio. Fitnets: Hints for thin deep nets. *arXiv preprint arXiv:1412.6550*, 2014.
- [46] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, et al. Imagenet large scale visual recognition challenge. *International journal of computer vision*, 115:211–252, 2015.
- [47] Mark Sandler, Andrew Howard, Menglong Zhu, Andrey Zhmoginov, and Liang-Chieh Chen. Mobilenetv2: Inverted residuals and linear bottlenecks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4510–4520, 2018.
- [48] Christian Szegedy, Wei Liu, Yangqing Jia, Pierre Sermanet, Scott Reed, Dragomir Anguelov, Dumitru Erhan, Vincent Vanhoucke, and Andrew Rabinovich. Going deeper with convolutions. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1–9, 2015.
- [49] Lin Wang and Kuk-Jin Yoon. Knowledge distillation and student-teacher learning for visual intelligence: A review and new outlooks. *IEEE transactions on pattern analysis and machine intelligence*, 44(6):3048–3068, 2021.
- [50] R Wightman, H Touvron, and H Jégou. Resnet strikes back: An improved training procedure in timm. *arXiv 2021. arXiv preprint arXiv:2110.00476*.
- [51] Kohei Yamamoto. Learnable companding quantization for accurate low-bit neural networks. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 5029–5038, 2021.
- [52] Jiwei Yang, Xu Shen, Jun Xing, Xinmei Tian, Houqiang Li, Bing Deng, Jianqiang Huang, and Xian-sheng Hua. Quantization networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7308–7316, 2019.
- [53] Dongqing Zhang, Jiaolong Yang, Dongqiangzi Ye, and Gang Hua. Lq-nets: Learned quantization for highly accurate and compact deep neural networks. In *Proceedings of the European conference on computer vision (ECCV)*, pages 365–382, 2018.
- [54] Linfeng Zhang, Jiebo Song, Anni Gao, Jingwei Chen, Chenglong Bao, and Kaisheng Ma. Be your own teacher: Improve the performance of convolutional neural networks via self distillation. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 3713–3722, 2019.
- [55] Xiangyu Zhang, Xinyu Zhou, Mengxiao Lin, and Jian Sun. Shufflenet: An extremely efficient convolutional neural network for mobile devices. In *Proceedings of*

- the IEEE conference on computer vision and pattern recognition*, pages 6848–6856, 2018.
- [56] Hengshuang Zhao, Jianping Shi, Xiaojuan Qi, Xiaogang Wang, and Jiaya Jia. Pyramid scene parsing network. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2881–2890, 2017.
- [57] Ritchie Zhao, Yuwei Hu, Jordan Dotzel, Chris De Sa, and Zhiru Zhang. Improving neural network quantization without retraining using outlier channel splitting. In *International conference on machine learning*, pages 7543–7552. PMLR, 2019.
- [58] Aojun Zhou, Anbang Yao, Yiwen Guo, Lin Xu, and Yurong Chen. Incremental network quantization: Towards lossless cnns with low-precision weights. *arXiv preprint arXiv:1702.03044*, 2017.
- [59] Shuchang Zhou, Yuxin Wu, Zekun Ni, Xinyu Zhou, He Wen, and Yuheng Zou. Dorefa-net: Training low bitwidth convolutional neural networks with low bitwidth gradients. *arXiv preprint arXiv:1606.06160*, 2016.
- [60] Ke Zhu, Yin-Yin He, and Jianxin Wu. Quantized feature distillation for network quantization. 2023.
- [61] Bohan Zhuang, Lingqiao Liu, Mingkui Tan, Chunhua Shen, and Ian Reid. Training quantized neural networks with a full-precision auxiliary module. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 1488–1497, 2020.