# Quantized Neural Networks: Characterization and Holistic Optimization

Yoonho Boo, Sungho Shin, and Wonyong Sung

Department of Electrical and Computer Engineering
Seoul National University
dnsgh337@snu.ac.kr; sungho.develop@gmail.com; wysung@snu.ac.kr

**Abstract.** Quantized deep neural networks (QDNNs) are necessary for low-power, high throughput, and embedded applications. Previous studies mostly focused on developing optimization methods for the quantization of given models. However, quantization sensitivity depends on the model architecture. Therefore, the model selection needs to be a part of the QDNN design process. Also, the characteristics of weight and activation quantization are quite different. This study proposes a holistic approach for the optimization of QDNNs, which contains QDNN training methods as well as quantization-friendly architecture design. Synthesized data is used to visualize the effects of weight and activation quantization. The results indicate that deeper models are more prone to activation quantization, while wider models improve the resiliency to both weight and activation quantization. This study can provide insight into better optimization of QDNNs.

**Keywords:** Quantized Deep Neural Network, Activation Quantization, Weight Quantization, Holistic approach

## 1 Introduction

Deep neural network (DNN) applications frequently demand extremely large models for an improved performance, which consumes a large amount of computation power not only for training but also for inference [45,8]. Thus, it is necessary to reduce their complexity for implementation on embedded devices. Various DNN compression methods have recently been devised to reduce the computational cost, power consumption, and storage space. Network quantization is one well-known method for substituting 32-bit floating-point weights with low bit-width numbers that usually employ one to four bits. Specifically, the performance of a quantized DNN (QDNN) is mostly maintained when retraining is applied after weight quantization [18,7]. Meanwhile, activation quantization has also been studied to reduce the computational cost and working memory footprint [31,43]. Therefore, activation quantization is particularly effective for DNN models with a large hidden-state dimension, such as convolutional neural networks (CNNs). Most previous studies on QDNN optimization have focused on the quantization number formats and training methods. The goal of these

previous studies has been reducing the performance gap between the floating-point and quantized models. However, not all networks can be quantized in the same manner. Some networks are more robust to weight quantization, whereas some others are not [36]. Optimizing a QDNN requires understanding the characteristics of such quantization errors.

In this paper, we visualize the characteristics of the quantization errors and their effects on the performance of QDNNs when the model architecture and sizes are different. We use synthetic data and DNN models for error characteristic visualization. Based on the analysis results, we adopt two simple training methods to compensate weight and activation quantization errors; fine-tuning with cyclic learning rate scheduling for improved generalization and applying regularization term that reduces noise amplification through propagation. Experiments are conducted using CIFAR-10, ImageNet, and PASCAL VOC 2012 semantic image segmentation benchmark. The contributions of this paper are as follows:

- We visualize the errors from the weight and activation quantization. The results indicate that the effect of the weight quantization error reduces the generalization capability whereas activation quantization error induces noise.
- We show that increasing the width of a DNN model helps to mitigate the quantization effects of both the weight and activation whereas increasing the depth only decreases the weight quantization error.
- We reduce the weight and activation quantization errors by employing training methods that improve the generalization capability and a regularization term that increases the noise robustness, respectively.
- Our work is a holistic approach for the optimization of QDNN by examining the quantization effects of weights and activations, and also the architectural change.

## 2   Backgrounds

### 2.1   Related Works on Network Quantization

Most DNN models are trained using 32-bit floating-point numbers. Apparently, DNN models do not demand 32-bit precision. Many quantization methods have been developed, some of which use an extremely small bit-width for a weight representation, such as 1-bit binary [7,31] or 2-bit ternary [10,44]. The signal-to-quantization-noise ratio (SQNR) of several weight quantizers was also compared [25]. Quantization noise has been measured to find a better training scheme [17] or optimal quantization precision [33]. Activation quantization has also been developed to lower the computational costs [5]. An efficient QDNN implementation on embedded systems has also been studied [13,1]. The weight quantization effects usually depend on the model size; small DNN models tend to show considerable performance degradation after quantization [36]. In particular, increasing the number of parameters in CNNs reduces the quantization sensitivity [27].

However, considering the purpose of model compression, the number of parameters needs to be constrained. A recent study showed that weight quantization up to certain bits does not reduce the memorization capacity [2]. During the last several years, residual connections have been developed mainly for improved training of neural networks [15]. Architectural modifications of increasing the width or moving the location of activation and batch normalization have been studied [41,16]. These architectural changes also affect the quantization sensitivity.

The activation quantization has not been discussed as much as weight quantization, and most studies have not distinguished the effects of activation and weight quantization [31,20]. It has been observed that activation usually demands more bits than weights [43]. The different quantization approaches for the weight and activation are applied in [29] because the latter was not suitable for cluster-based quantization. Many studies have shown that a DNN can be vulnerable to noise. Even with an extremely small amount of noise, the inference of a DNN can easily be manipulated [37,22]. Previous studies have shown that quantizing the input makes it robust to adversarial attacks by reducing the amount of noise [40]. Several studies have shown that a QDNN can help defend from adversarial attacks [30,11]. However, QDNNs become more vulnerable to adversarial attacks than floating-point models when the noise exceeds a certain level [26]. The adversarial noise becomes larger at the deeper layers [24].

## 2.2   Revisit of QDNN Optimization

The process of uniform quantization for a DNN involves the following two steps, namely, clipping and quantization:

$$\hat{x} = Clip\ (x, \alpha, \beta), \quad Q(x) = \Delta \lfloor \frac{\hat{x}}{\Delta} + 0.5 \rfloor. \tag{1}$$

The parameters of the DNNs are signed values so that the clip value $\beta = -\alpha = \Delta(2^{n-1} - 1)$ where $n$ is the number of bits used to represent each parameter. Parameter quantization is mainly applied to the weights. For the sake of simple structure, all fixed-point weights in a layer share one scale factor $\Delta$. The activation quantization is used to lower the computational cost and the size of the working memory for inference. When using the ReLU activation, the hidden vectors are represented with unsigned values and $\alpha$ and $\beta$ becomes 0 and $\Delta(2^n - 1)$, respectively. Low-precision quantized networks require training in a fixed-point domain to improve the performance as follows:

$$W_t^q = Q(W_t) \tag{2}$$

$$E_t = f(x_t, y_t, W_t^q) \tag{3}$$

$$W_{t+1} = W_t - \alpha \frac{\partial E_t}{\partial W_t^q}, \tag{4}$$

where $E_t$ is the loss computed through the model, $f(\cdot)$, at $t$-th iteration. Forward and backward propagations are conducted using quantized weights and activation. However, the computed gradients need to be added to the floating-point

weights because those gradients are relatively small compared to the step size $\Delta$ [7]. It is known that QDNNs perform better when retrained from a pretrained model at floating-point than trained from the scratch [43,20].

Most QDNN studies quantize the weights or activation according to Equation (1), although the processes of obtaining $\Delta$, $\alpha$, and $\beta$ are different [18,5]. However, the effect of each quantization on the inference is quite different. The quantization errors are $\epsilon_W = W - Q(W)$ and $\epsilon_a = a - Q(a)$ where $\epsilon_W$ and $\epsilon_a$ are errors owing to the quantization of the weight and activation, respectively. Because the trained weights have fixed values during inferences, $\epsilon_W$ is a constant error. In other words, weight quantization can be modeled as the process of distorting the weights of the DNNs. This changes the direction of the input-prediction mapping function of the DNNs and thus causes distorted results at the inference. By contrast, $\epsilon_a$ is an error that depends on the input applied during the inference process. Depending on the remainder of the hidden vector divided by $\Delta$, the direction or magnitude of the error may change. That is, $\epsilon_a$ induces noise with a maximum magnitude of $\frac{\Delta}{2}$.

In the rest of this paper, we analyze how such differences in quantization errors affect the performance of QDNNs under various model architectures. The weight quantization method in [18] and the PACT activation quantization [5] are adopted for our experiments. QDNNs are retrained from pretrained floating-point models.

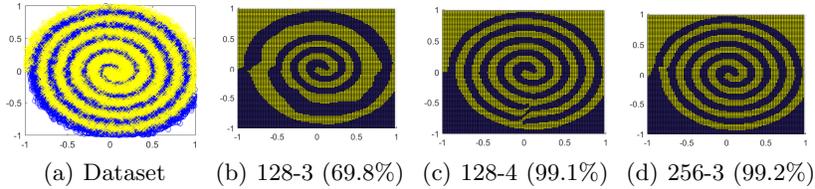## 3    Visualization of Quantization Errors using Synthetic Dataset

### 3.1    Synthetic Dataset Generation

Most DNNs and their training samples used in real tasks have extremely high dimensions, and it is therefore very difficult to discern the effects of quantization. For a visualization analysis of QDNNs, we synthesized 2D inputs whose elements consist of $x$ and $y$ axes. The training dataset is composed of inputs $S \in \mathbb{R}^2$ and $C \in (0, 1)$, which are used to train FCDNNs for binary classification. The training dataset is synthesized through two steps. First, the core samples, $s_c$, mapped to a label, $c$, are generated using the following equation:

$$s_0 \in \left\{ (x, y) \middle| \begin{cases} x^2 + y^2 = (2i+1)^2 r^2, & y > 0 \\ (x-r)^2 + y^2 = (2i+2)^2 r^2, & y \leq 0 \end{cases} \right\}, \tag{5}$$

$$s_1 \in \left\{ (x, y) \middle| \begin{cases} x^2 + y^2 = (2i+2)^2 r^2, & y > 0 \\ (x-r)^2 + y^2 = (2i+1)^2 r^2, & y \leq 0 \end{cases} \right\}. \tag{6}$$

In our experiments, $i$ is within $\{0, 1, 2, 3, 4\}$ and $r$ is 0.1. For each $i$, two semicircles correspond to one label. In each semicircle, we sample 100 points by increasing the angle linearly. Therefore, the total number of core samples is 2,000. Next, we generate subsamples by adding Gaussian noise to each core sample as follows:

(a) Dataset     (b) 128-3 (69.8%)   (c) 128-4 (99.1%)   (d) 256-3 (99.2%)

**Fig. 1.** Illustrations of the dataset and prediction results from FCDNNs. (a) The synthetic dataset used to train FCDNNs. (b) An example of the evaluation results when the model is too small to learn to the data distribution. (c, d) Examples of large models. The values in parentheses are the accuracy for the correct answer.

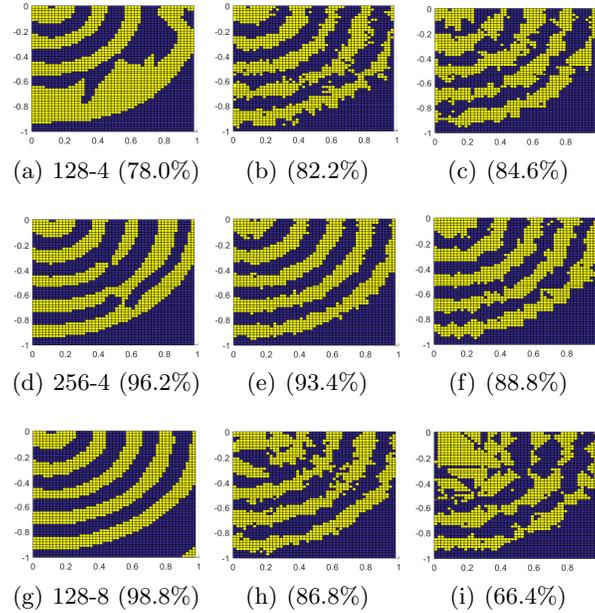$$s' = s + \epsilon, \quad \epsilon \sim \mathcal{N}(0, \frac{1}{3}r\mathbf{I}). \tag{7}$$

Nine subsamples for each core sample are used, and the samples are mapped to the same labels. As a result, the total number of datasets used for training is 10,000 for each label. The distribution of the generated dataset is shown in Fig. 1 (a).

We quantize the FCDNN trained using the generated dataset and analyze the difference between the weight and activation quantization. In particular, we visualize the errors when varying the model depth or width. Two datasets are adopted for the evaluation of the trained model. The first is a test dataset. The correct answer dataset is constructed by dividing the area corresponding to each label by the radius of the semicircle. This dataset is used for quantitative analysis of the trained QDNNs by measuring accuracy. The other is a grid dataset that consists of $(x, y) \in \{(x, y)|0 < x < 1, 0 < y < 1\}$. By visualizing the prediction on the $x$-$y$ plane, we can analyze whether the input-prediction mapping of a DNN is distorted or as added noise.

### 3.2   Results on Synthetic Dataset

We devise artificial DNN models for testing with the synthetic dataset. Fully-connected DNN (FCDNN) models are employed with varying depth and width. In addition, models with residual connections are also considered. We indicate the experimental models as "width" - "depth" of FCDNNs. The prediction results of floating-point models using the evaluation dataset are shown in Fig. 1. The 128-3 FCDNN shows quite a different prediction result from the actual data distribution. Fig. 1 (c,d) show that increasing the depth to 4 or the width to 256 is sufficient to learn the synthesized dataset quite faithfully. For the remaining experiments, we represent only the bottom-right quarter circle for detailed visualization. All QDNN results are reported after retraining.
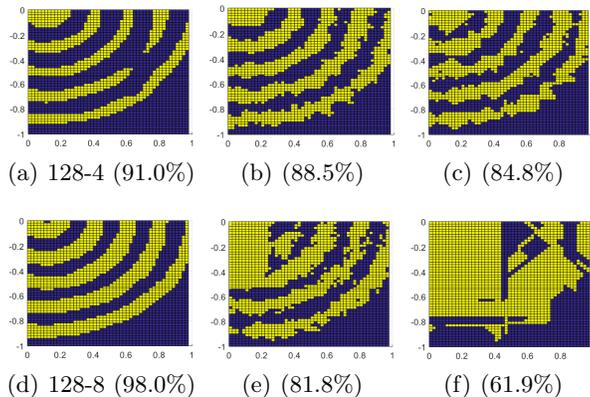
Fig. 2 (**Left**) shows the effects of the weight quantization according to the width and depth of FCDNNs. We can see that weight quantization distorts

(a) 128-4 (78.0%)          (b) (82.2%)          (c) (84.6%)

(d) 256-4 (96.2%)          (e) (93.4%)          (f) (88.8%)

(g) 128-8 (98.8%)          (h) (86.8%)          (i) (66.4%)

**Fig. 2.** Prediction results of QDNN as the width and depth increase. (**Left**) 2-bit weights. (**Middle**) 2-bit activations. (**Right**) 2-bit weights and activations.

the input-prediction mapping of the DNN. The evaluation results of the weight quantized models resemble that of a small floating-point model, such as the 128-3 FCDNN shown in Fig. 1 (b). The experiment results show that the decrease in learning ability occurs similarly when the model size is reduced or the weights are quantized. As studied in [28], increasing the model size helps generalization. Our experiments show that the generalization capability decreases as the precision of the parameters is lowered. Thus, the effect of reduced generalization capability due to the weight quantization is not noticeable when the model size is large enough. The distortion with 2-bit weight quantization is barely found when the layer width or the number of layers is increased.

Fig. 2 (**Middle**) shows the activation quantization results. The effect of the activation quantization is very different from that of the model capacity reduction in a DNN. Activation quantization appears to add noise to the prediction results. Although both the weight and activation quantization errors degrade the performance of DNNs, they behave in a completely different manner. When the activation is quantized to 2 bits, increasing the depth does not mitigate the noise added to the prediction results. Rather, the noise tends to worsen with weight quantization when the depth increases. Activation quantization is related to the dimension of each layer rather than the capacity of the model. Activation quantization in wide FCDNNs (Fig. 2 (e)) is more robust than in deep FCDNNs (Fig. 2 (h)). The effect of noise from the activation quantization is reduced because the
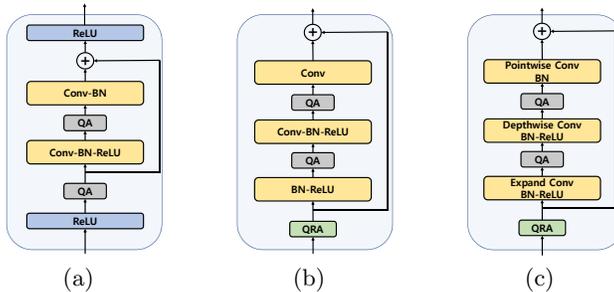
(a) 128-4 (91.0%)        (b) (88.5%)        (c) (84.8%)

(d) 128-8 (98.0%)        (e) (81.8%)        (f) (61.9%)

**Fig. 3.** Prediction results of QDNNs with residual connection. (**Left**) 2-bit weights. (**Middle**) 2-bit activations. (**Right**) 2-bit weights and activations.

number of dimensions of the hidden vector received as the input in each layer is increased. When quantizing both the weight and activation, the two errors are combined, and result in a noisy and distorted prediction, as shown in Fig. 2 (**Right**). The more results of different QDNNs are reported in Appendix **A**.

   We also analyze the effect of residual connections on activation quantization. The residual connection helps train DNNs with an extremely large number of layers [15]. In our experiment, the residual connections are implemented by adding each hidden output to the activation of the next layer. The result is multiplied by 0.5 to preserve the scale of the intermediate results. Fig. 3 shows the results of the quantizing activation when the residual connection is applied. Residual connections help alleviate the distortion through weight quantization. However, FCDNNs with residual connections are more sensitive to activation quantization than the original models. With residual connections, the outputs of the quantized hidden layer are summed so that the activation quantization noise is also added. As a result, applying residual connections shows more noisy prediction in deep models, such as 128-8 FCDNNs.

## 4   QDNN Optimization with Architectural Transformation and Improved Training

The visualization results with the synthesized data show that weight quantization decreases the generalization capability of DNNs, while activation quantization induces noised inference. Also, the quantization effects depend on the architecture very much. Based on this observation, we employ three approaches for QDNN optimization. The first one is modifying the architecture quantization-friendly. The second one is the training method for improved generalization. This technique is intended to reduce the effects of weight quantization. The third one is applying the regularization term that limits the amount of activation noise.

**Fig. 4.** Types of residual blocks. (a) basic block, (b) pre-activation block, and (c) depthwise block. QA and QRA denote the activation quantization operations.

### 4.1  Architecture Transformation for Improved Robustness to Quantization

Deep CNN models are hard to train because of the gradient vanishing problem. The residual architecture was developed to solve this problem [15,39,38]. In CNN with residual connections, increasing the depth, often over 100, usually helps to improve the performance. Of course, widening the networks also increases the performance [41]. When the number of parameters is limited, increasing the depth is usually preferred because the model complexity rises in proportional to the depth, but squarely proportional to the width. However, our work in in Section 3 shows that deep CNN models are very prone to activation quantization. The conventional approach for QDNN design is developing the best performing floating-point model, and then quantizing it in the best way possible. In this case, the best performing floating-point model prefers deeper ones, which are, however, prone to activation quantization. Thus, we need to consider the effects of weight and activation quantization even for the initial floating-point model design. Wide CNN models, which are considered parameter inefficient, often show better performance than deep ones when the activation is severely quantized.

Recent CNN models employ various residual blocks for improved performance or parameter efficiency. The most well-known residual blocks for CNNs are shown in Fig. 4. The depthwise block employed to MobileNetV2 [34] helps to reduce the number of parameters and computations. However, the quantization performances of these blocks are not well studied.

### 4.2  Cyclical Learning Rate Scheduling for Improved Generalization

We adopt the cyclic learning rate scheduling (CLR) as a way to reduce the effects of weight quantization. CLR increases and decreases the learning rate periodically, while conventional training usually reduces the learning rate in one direction. This method is known to increase the generalization capability of the model by leading to a flat loss surface [35]. Among a few different cyclical learning rate scheduling algorithms, we choose the one that alters the learning

rate discretely, which is known to be more effective for generalization [19]. The maximum and minimum boundaries of the CLR are determined between the 100 and 0.1 times of the last learning rate of the retraining procedure, respectively. The learning rate changes 8 times in one cycle and exponentially decreases or increases. The CLR scheduling for each task is illustrated in Appendix **B**.

### 4.3 Regularization for Limiting the Activation Noise Amplification

Training methods to increase noise robustness of DNNs have been studied in the field of adversarial training. Parseval networks reduce the Lipschitz constant so that the noise of the input is not amplified as the layer increases [6]. In particular, [26] shows that activation quantized DNNs exacerbate performance degradation due to adversarial noise, and add the regularization term to the loss to keep the Lipschitz constant of each layer small. The regularization term is as follows:

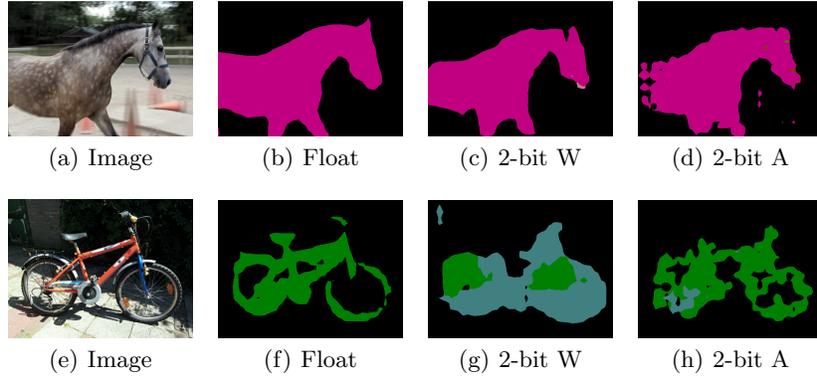$$L_{Lip} = \frac{1}{2} \sum_{W_l} ||W_l^T W_l - I||^2.$$ (8)

Note that convolution kernels are reshaped to $(k \times k \times c_{in}, c_{out})$ where $k$, $c_{in}$, and $c_{out}$ are the kernel size, input channels, and output channels, respectively. $L_{Lip}$ was applied to enhance the adversarial attack robustness of activation quantized DNNs [26]. We show that $L_{Lip}$ can reduce the noise due to activation quantization itself. Also, we compare the effect of the regularization term on the performance of weight quantized DNNs.
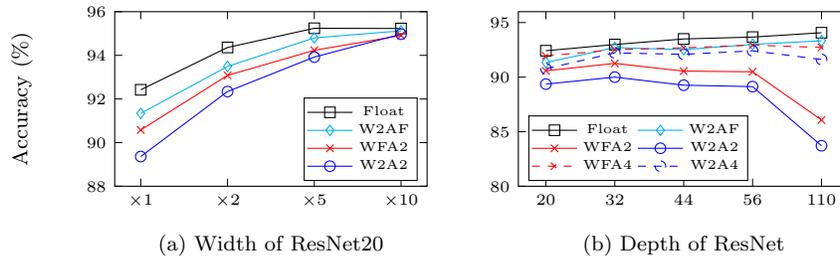
## 5 Experimental Results

### 5.1 Visualizing the Effects of Quantization on the Segmentation Task

We first visualize the weight and activation quantization effects using a segmentation task. The PASCAL VOC 2012 dataset [9] is used. The dataset contains 1,464 training images and 1,449 validation images. Each image is labeled at pixel-level with 20 object classes and a background class. The MobileNetV2 [34] is adopted, which is trained according to DeepLabV3 [4] with 10,582 augmented training images [14] [1]. The model is used as a floating-point pretrained model after fine-tuned using the original 1,464 training images for 30K iterations. The output stride is 16 and Atrous Spatial Pyramid Pooling (ASPP) [3] is not applied. The performance is measured in terms of mean intersection over union (mIOU) without multi-scaling and flipping input images. Only the original training images are used for retraining and fine-tuning. The retraining is conducted for 30K iterations with a batch size of 16. The initial learning rate is 1e-3 and the learning rate policy is the same as [4].

---

[1] We obtained the pretrained model from https://github.com/tensorflow/models/tree /master/research/deeplab

(a) Image      (b) Float      (c) 2-bit W      (d) 2-bit A



(e) Image      (f) Float      (g) 2-bit W      (h) 2-bit A

**Fig. 5.** Visualization of quantization errors on the PASCAL VOC segmentation benchmark. W and A are abbreviations for the weight and activation, respectively. Activation outputs are retained in floating-point precision on weight quantized model, and vice-versa.



(a) Width of ResNet20      (b) Depth of ResNet

**Fig. 6.** Performance of quantized ResNet on the CIFAR-10 testset according to the (a) width of the ResNet20 and (b) depth of the ResNet. Legends represent the precision of 'weights' (W) and 'activation' (A). 'F' denotes the floating-point precision.

The segmentation results of the retrained QDNNs are visualized in Fig. 5. Either weights or activations are quantized to 2 bits. When the object is simple, as shown in Fig. 5 (a), the weight quantized model seems to perform the segmentation fairly well. However, the results with the activation quantized model contains some noise on the section where the background and the object colors are similar. In the segmentation of a complex one, the weight quantized model fails to find the characteristics of the object, as shown in Fig. 5 (e). We can even consider that the activation noise corrupted model segments the bicycle more faithfully than the weight quantized model. The visualization results imply that weight quantization degrades the generalization ability, and activation quantization induces noise. The experiment with the segmentation task confirms the observation with the synthetic dataset in Section 3. More segmentation results of QDNNs are compared in Appendix **C**.
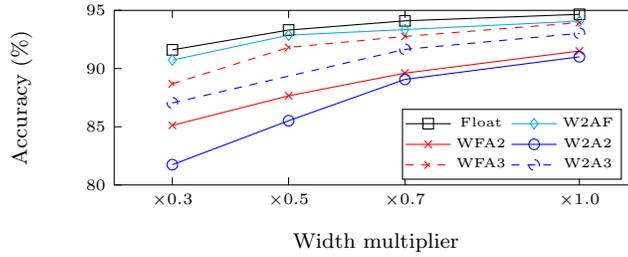
**Fig. 7.** CIFAR-10 test accuracy (%) of MobileNetV2 according to the width multiplier.



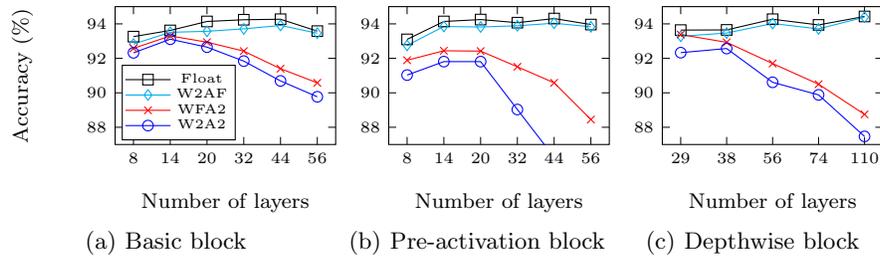(a) Basic block          (b) Pre-activation block     (c) Depthwise block

**Fig. 8.** CIFAR-10 test accuracy (%) of quantized ResNet according to the block types with varying depth and width. Note that the number of parameters of all experimented models are about 1M.

## 5.2   The Width and Depth Effects on QDNNs

We analyze the weight and activation quantization sensitivities when the depth and width of ResNet vary using the CIFAR-10 dataset [21]. The depth refers the number of layers and the width is that of channels in a CNN. Simple data augmentation techniques, cropping and flipping, are applied as suggested in [23]. The batch size is 128 and the number of epochs for pretraining is 200. The SGD optimizer with a momentum of 0.9 is used. The learning rate starts at 0.1 and decays by 0.1 times at 100 and 150 epochs. The L2 regularization is applied with a scale of 5e-4. Quantized models are retrained for 100 epochs with the initial learning rate of 0.01, and the learning rate decreases by a factor of 0.1 at 50 and 80 epochs. We do not employ the L2 regularization when retraining of the quantized model. The number of layers is either 20, 32, 56, or 110 and the width multiplier ranges from x1 to x10 times of the original ResNet. The experimental model denoted as ResNet×$I$ employs $I$ times larger number of channels. The width-expanded ResNets are compared in Fig. 6 (a). As the width of the ResNet20 increases, both weight and activation quantization errors decrease. The performance of ResNets when the depth of layers increases is shown in Fig. 6 (b). When the depth increases, the 2-bit weight quantized models show improved performance, but those with 2-bit quantized activations

exhibit degraded performance. At least, 4-bit activation quantization is needed for deep ResNet.

We also evaluate the quantization sensitivity of the MobileNetV2 [34]. The width multiplier is used to control the number of parameters of MobileNetV2. The MobileNetV2 performance on CIFAR-10 with various width-multiplication factors is shown in Fig. 7. Note that the MobileNetV2×1.0 is the same with the original model except that the first stride of 2 is replaced to 1. As the width decreases, the performances of full-precision and weight quantized models degrade gradually. However, the activation quantized models exhibit severe performance loss as the width decreases. The results indicate that DNNs with small widths are more vulnerable to the activation quantization.

### 5.3   QDNN Architecture Selection under the Parameter Constraint

We compare the quantization sensitivity of CNNs according to the model structure under the constraint on the number of parameters. CNNs with three different block types: basic, pre-activation, depthwise blocks are evaluated using CIFAR-10 dataset. The number of parameters of these models is approximately one million. The performance of the quantized CNNs is shown in Fig. 8. The experimented model configurations are summarized in Appendix **D**. When the number of parameters is comparable, increasing the depth to a certain range helps to improve the performance of full precision and weight quantized models. However, the models with 2-bit activation quantization show poor performances when the depth increases beyond certain numbers. We can also find that depthwise blocks are more robust to the activation quantization. The 2-bit CNNs with the basic or pre-activation blocks show severe performance degradation when the depth is over 32 while the performance with depthwise blocks is improves until the depth of 38. For the CIFAR-10 dataset, the best performing quantized ResNet can be designed by choosing the depth of 14 with the basic blocks or 38 with the depthwise ones.

We change the depth and width of the ResNets and evaluate the quantization performance on the ImageNet dataset [32]. Data augmentation methods and hyperparameters are the same as [20]. Pre-activation blocks are employed and the shortcut signals are quantized to 8 bits. 4-level 2-bit weight quantization is applied. The experimented ResNet structures are compared in Table 1. Since ResNets for ImageNet classification have four groups of residual blocks, which are separated by the convolution layers with the stride of 2, the number of blocks is represented as a list: [first, second, third, fourth] groups. For example, ResNet18×1.4 means that the depth is 18 with 8 number of blocks and the initial channel width is 90. The number of channels increases by a factor of 2 at every convolution layer with the stride of 2.

The performance degradation by the quantization is much lower when the model is shallow under a comparable number of parameter. However, the top-1 accuracy of 2-bit ResNet18×1.4 is 0.3% lower than the 2-bit ResNet34 because the floating-point performance is too low. When the depth, $L$, is 26 and the initial channel width, $D_{init}$ is 70, the top-1 accuracy of the floating-point model is 0.6%

**Table 1.** Performance comparison of 2-bit ResNet on ImageNet validation set. '$L$' and '$D_{init}$' represent the number of layers and initial channel dimension, respectively. '$B$' is the stack of the residual blocks and represented separately according to the stride of 2.

| Model | | | | | | Top-1 Acc (%) | | |
|---|---|---|---|---|---|---|---|---|
| Type | $L$ | $D_{init}$ | $B$ | Params. $(10^6)$ | Ops. $(10^9)$ | Float | W2A2 | Diff. |
| ResNet34 [42] | 34 | 64 | [3, 4, 6, 3] | 21.8 | 3.7 | 73.8 | 69.8 | 4.0 |
| ResNet34 [12] | 34 | 64 | [3, 4, 6, 3] | 21.8 | 3.7 | 73.8 | 70.0 | 3.8 |
| ResNet34 | 34 | 64 | [3, 4, 6, 3] | 21.8 | 3.7 | 73.6 | 70.5 | 3.1 |
| ResNet18×1.4 | 18 | 90 | [2, 2, 2, 2] | 22.8 | 3.5 | 72.6 | 70.2 | **2.4** |
| ResNet26×1.1 | 26 | 70 | [3, 3, 3, 3] | 21.4 | 3.3 | 73.0 | **70.6** | **2.4** |
| ResNet50 [42] | 50 | 64 | [3, 4, 6, 3] | 25.6 | 4.1 | 76.4 | 71.5 | 4.9 |
| ResNet50 | 50 | 64 | [3, 4, 6, 3] | 25.6 | 4.1 | 76.3 | 72.7 | 3.6 |
| ResNet44×1.1 | 44 | 72 | [3, 4, 5, 2] | 25.0 | 4.6 | 75.5 | **73.4** | **2.1** |
| ResNet101 | 101 | 64 | [3, 4, 23, 3] | 44.6 | 7.9 | 77.5 | 20.1 | 57.4 |
| ResNet50×1.3 | 50 | 85 | [3, 4, 6, 3] | 44.2 | 7.2 | 77.2 | **73.7** | **3.5** |

lower but the 2-bit quantized model achieves 0.1% higher top-1 accuracy when compared to the ResNet34. The performance improvements is more significant when the model is deeper. The top-1 accuracy degradation of 2-bit ResNet50 is 4.9 [42] and 3.6(Ours). However, 2-bit ResNet44×1.1 shows 2.1% top-1 accuracy drop. As a result, 2-bit ResNet44×1.1 achieves the 0.7% top-1 accuracy improvement compared to the ResNet50 with the similar number of parameters and operations. The floating-point ResNet101 shows the top-1 accuracy of 77.5%, which outperforms shallow and wider models with a comparable number of parameters. However, the 2-bit quantization of ResNet101 degrades the performance significantly, showing only about 20.1% top-1 accuracy. As observed in Section 5.2, the 2-bit activation quantized model degrades the performance dramatically when the model is very deep. On the other hand, the ResNet with the depth of 50 achieves 73.5% top-1 accuracy even after 2-bit quantization. These results indicate that the performance of QDNNs can be improved simply by designing proper depth and width of the model.

## 5.4   Results of Training Methods on QDNNs

We assess the effects of the training methods on QDNN optimization using the segmentation task. The effects of applying the CLR for improved generalization and adding the regularization term, $L_{Lip}$, for noise robustness are summarized in Table 2. The mIOU of the floating-point pretrained model is 76.97, that of the retrained 4-bit weight is 73.63, and that of retrained 4-bit activation is 74.79. CLR is applied with a period of 3K iterations and fine-tuning was performed

**Table 2.** Performance improvements of quantized MobileNetV2 on the PASCAL VOC 2012 validation set. $n_W$ and $n_A$ represent the precision of the weights and activations, respectively.

| Method | $n_W$ | $n_A$ | mIOU |
|---|---|---|---|
| Pretrained model | Float | Float | 76.97 |
| Retrain (baseline) | 4-bit | Float | 73.63 |
| Fine-tune with CLR | 4-bit | Float | **74.15** |
| Retrain with $L_{Lip}$ | 4-bit | Float | 73.22 |
| Retrain (baseline) | Float | 4-bit | 74.79 |
| Fine-tune with CLR | Float | 4-bit | 74.71 |
| Retrain with $L_{Lip}$ | Float | 4-bit | **74.99** |

**Table 3.** Performance in terms of mIOU on the PASCAL VOC 2012 validation set when both weights and activations are quantized.

| Method | $n_W$ / $n_A$ (bits) | | | |
|---|---|---|---|---|
| | 8/8 | 6/6 | 4/4 | 3/3 |
| Retrain (baseline) | **76.56** | 75.79 | 71.16 | 59.64 |
| Retrain ($L_{Lip}$) | 76.34 | 76.23 | 71.93 | 59.85 |
| $L_{Lip}$ + CLR | 76.46 | **76.40** | **72.56** | **60.74** |

for 15K iterations (i.e. 5 cycles). $L_{Lip}$ is added to the loss after multiplying the scaling factor of 1e-4. Applying the CLR increases the mIOU of the 4-bit weight quantized model by 0.5 but it is not effective for the activation quantized model. Retraining with the regularization term that reduces the Lipschitz constant improves the mIOU of the 4-bit activation quantized model. However, the performance of the weight quantized model is decreased when $L_{Lip}$ is applied. The performances of QDNNs when both weights and activations are quantized are shown in Table 3. The results show that the proposed approach for reducing the quantization effects of weights by CLR and activations by adding the Lipschitz loss works well when the precision is equal to or lower than 6-bit. But these techniques are not effective when the precision of quantization is 8-bit or larger. Even, the $L_{Lip}$ regularization degrades the performance of the 8-bit quantized model.

We also evaluate the effects of the training methods for the classification task. The performance degradation by severe quantization can be alleviated with $L_{Lip}$ and CLR as shown in Table 4. $L_{Lip}$ is added to the loss with a factor of 1e-4 and CLR is applied for 40 epochs with the learning rates between 1e-3 and 1e-5. Although we can improve the performance of QDNN considerably for deep networks by applying CLR and $L_{Lip}$ constraint, the best performing QDNN can be found when the depth is 32.

**Table 4.** CIFAR-10 test accuracy (%) improvements when retrained with $L_{Lip}$ and fine-tuned using CLR.

| $n_W$, $n_A$ | ResNet depth | | | |
|---|---|---|---|---|
| | 20 | 32 | 56 | 110 |
| Float | 92.42 | 92.99 | 93.67 | 94.07 |
| 2-bit | 89.36 | 90.00 | 89.13 | 83.71 |
| 2-bit ($L_{Lip}$ + CLR) | 89.68 | 90.60 | 90.24 | 87.77 |

## 6   Concluding remarks

We have presented a holistic approach for the optimization of quantized deep neural networks (QDNNs). We first visualize the effects of the weight and activation quantization error using a synthetic dataset. The result clearly shows that the effects of weight and activation quantization are different. Especially, activation quantization severely degrades the performance of deep models. Through additional experiments with real tasks, we confirm that the optimal model structure under a parameter constraint is different for the full-precision and quantized DNNs because floating-point models usually prefer the deep networks but QDNNs tend to show improved performances on wide networks. We also show the effects of the DNN training schemes for improved generalization and noise reduction to optimize QDNNs. The proposed holistic approach can yield much better QDNN when compared to the conventional design approaches that start from the best performing floating-point models and optimize them using elaborate quantization and training methods.
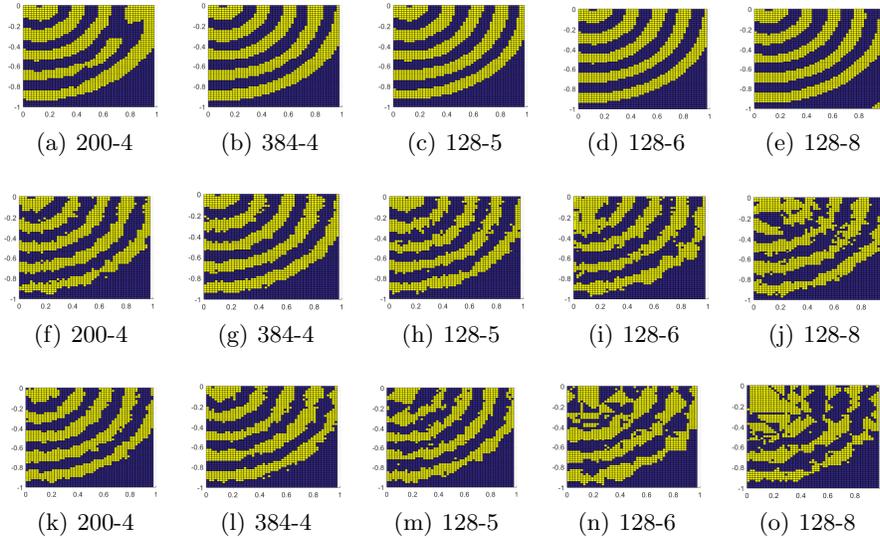
# References

1. Andri, R., Cavigelli, L., Rossi, D., Benini, L.: Yodann: An ultra-low power convolutional neural network accelerator based on binary weights. In: 2016 IEEE Computer Society Annual Symposium on VLSI (ISVLSI). pp. 236–241. IEEE (2016)
2. Boo, Y., Shin, S., Sung, W.: Memorization capacity of deep neural networks under parameter quantization. In: ICASSP 2019-2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP). pp. 1383–1387. IEEE (2019)
3. Chen, L.C., Papandreou, G., Kokkinos, I., Murphy, K., Yuille, A.L.: Deeplab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected crfs. IEEE transactions on pattern analysis and machine intelligence **40**(4), 834–848 (2017)
4. Chen, L.C., Papandreou, G., Schroff, F., Adam, H.: Rethinking atrous convolution for semantic image segmentation. arXiv preprint arXiv:1706.05587 (2017)
5. Choi, J., Wang, Z., Venkataramani, S., Chuang, P.I.J., Srinivasan, V., Gopalakrishnan, K.: PACT: Parameterized clipping activation for quantized neural networks. arXiv preprint arXiv:1805.06085 (2018)
6. Cisse, M., Bojanowski, P., Grave, E., Dauphin, Y., Usunier, N.: Parseval networks: Improving robustness to adversarial examples. In: Proceedings of the 34th International Conference on Machine Learning-Volume 70. pp. 854–863. JMLR. org (2017)
7. Courbariaux, M., Bengio, Y., David, J.P.: Binaryconnect: Training deep neural networks with binary weights during propagations. In: Advances in Neural Information Processing Systems (NIPS). pp. 3123–3131 (2015)
8. Devlin, J., Chang, M.W., Lee, K., Toutanova, K.: Bert: Pre–training of deep bidirectional transformers for language understanding. arXiv preprint arXiv:1810.04805 (2018)
9. Everingham, M., Eslami, S.A., Van Gool, L., Williams, C.K., Winn, J., Zisserman, A.: The pascal visual object classes challenge: A retrospective. International journal of computer vision **111**(1), 98–136 (2015)
10. Fengfu, L., Bo, Z., Bin, L.: Ternary weight networks. In: NIPS Workshop on EMDNN. vol. 118, p. 119 (2016)
11. Galloway, A., Taylor, G.W., Moussa, M.: Attacking binarized neural networks. In: International Conference on Learning Representations (ICLR) (2018)
12. Gong, R., Liu, X., Jiang, S., Li, T., Hu, P., Lin, J., Yu, F., Yan, J.: Differentiable soft quantization: Bridging full-precision and low-bit neural networks. In: Proceedings of the IEEE International Conference on Computer Vision. pp. 4852–4861 (2019)
13. Gupta, S., Agrawal, A., Gopalakrishnan, K., Narayanan, P.: Deep learning with limited numerical precision. In: International Conference on Machine Learning (ICML). pp. 1737–1746 (2015)
14. Hariharan, B., Arbeláez, P., Bourdev, L., Maji, S., Malik, J.: Semantic contours from inverse detectors. In: 2011 International Conference on Computer Vision. pp. 991–998. IEEE (2011)
15. He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. In: 2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR). pp. 770–778. IEEE (2016)
16. He, K., Zhang, X., Ren, S., Sun, J.: Identity mappings in deep residual networks. In: European conference on computer vision. pp. 630–645. Springer (2016)

17. Hou, L., Kwok, J.T.: Loss-aware weight quantization of deep networks. International Conference on Learning Representations (ICLR) (2018)
18. Hwang, K., Sung, W.: Fixed-point feedforward deep neural network design using weights +1, 0, and -1. In: Signal Processing Systems (SiPS), 2014 IEEE Workshop on. pp. 1–6. IEEE (2014)
19. Jastrzkebski, S., Kenton, Z., Arpit, D., Ballas, N., Fischer, A., Bengio, Y., Storkey, A.: Three factors influencing minima in SGD. arXiv preprint arXiv:1711.04623 (2017)
20. Jung, S., Son, C., Lee, S., Son, J., Han, J.J., Kwak, Y., Hwang, S.J., Choi, C.: Learning to quantize deep networks by optimizing quantization intervals with task loss. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 4350–4359 (2019)
21. Krizhevsky, A., Hinton, G.: Learning multiple layers of features from tiny images. Tech. rep., Citeseer (2009)
22. Kurakin, A., Goodfellow, I., Bengio, S., Dong, Y., Liao, F., Liang, M., Pang, T., Zhu, J., Hu, X., Xie, C., et al.: Adversarial attacks and defences competition. In: The NIPS'17 Competition: Building Intelligent Systems, pp. 195–231. Springer (2018)
23. Lee, C.Y., Xie, S., Gallagher, P., Zhang, Z., Tu, Z.: Deeply-supervised nets. In: Artificial Intelligence and Statistics. pp. 562–570 (2015)
24. Liao, F., Liang, M., Dong, Y., Pang, T., Hu, X., Zhu, J.: Defense against adversarial attacks using high-level representation guided denoiser. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 1778–1787 (2018)
25. Lin, D., Talathi, S., Annapureddy, S.: Fixed point quantization of deep convolutional networks. In: International Conference on Machine Learning. pp. 2849–2858 (2016)
26. Lin, J., Gan, C., Han, S.: Defensive quantization: When efficiency meets robustness. International Conference on Learning Representations (ICLR) (2019)
27. Mishra, A., Nurvitadhi, E., Cook, J.J., Marr, D.: WRPN: wide reduced-precision networks. arXiv preprint arXiv:1709.01134 (2017)
28. Neyshabur, B., Li, Z., Bhojanapalli, S., LeCun, Y., Srebro, N.: Towards understanding the role of over-parametrization in generalization of neural networks. arXiv preprint arXiv:1805.12076 (2018)
29. Park, E., Ahn, J., Yoo, S.: Weighted-entropy-based quantization for deep neural networks. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 5456–5464 (2017)
30. Rakin, A.S., Yi, J., Gong, B., Fan, D.: Defend deep neural networks against adversarial examples via fixed anddynamic quantized activation functions. arXiv preprint arXiv:1807.06714 (2018)
31. Rastegari, M., Ordonez, V., Redmon, J., Farhadi, A.: Xnor-net: Imagenet classification using binary convolutional neural networks. In: Proceedings of the European Conference on Computer Vision (ECCV). pp. 525–542. Springer (2016)
32. Russakovsky, O., Deng, J., Su, H., Krause, J., Satheesh, S., Ma, S., Huang, Z., Karpathy, A., Khosla, A., Bernstein, M., et al.: Imagenet large scale visual recognition challenge. International Journal of Computer Vision $115$(3), 211–252 (2015)
33. Sakr, C., Shanbhag, N.: Minimum precision requirements for deep learning with biomedical datasets. In: 2018 IEEE Biomedical Circuits and Systems Conference (BioCAS). pp. 1–4. IEEE (2018)

34. Sandler, M., Howard, A., Zhu, M., Zhmoginov, A., Chen, L.C.: MobileNetV2: Inverted residuals and linear bottlenecks. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 4510–4520 (2018)
35. Smith, L.N.: Cyclical learning rates for training neural networks. In: 2017 IEEE Winter Conference on Applications of Computer Vision (WACV). pp. 464–472. IEEE (2017)
36. Sung, W., Shin, S., Hwang, K.: Resiliency of deep neural networks under quantization. arXiv preprint arXiv:1511.06488 (2015)
37. Szegedy, C., Zaremba, W., Sutskever, I., Bruna, J., Erhan, D., Goodfellow, I., Fergus, R.: Intriguing properties of neural networks. arXiv preprint arXiv:1312.6199 (2013)
38. Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A.N., Kaiser, Ł., Polosukhin, I.: Attention is all you need. In: Advances in neural information processing systems. pp. 5998–6008 (2017)
39. Wu, Y., Schuster, M., Chen, Z., Le, Q.V., Norouzi, M., Macherey, W., Krikun, M., Cao, Y., Gao, Q., Macherey, K., et al.: Google's neural machine translation system: Bridging the gap between human and machine translation. arXiv preprint arXiv:1609.08144 (2016)
40. Xu, W., Evans, D., Qi, Y.: Feature squeezing: Detecting adversarial examples in deep neural networks. arXiv preprint arXiv:1704.01155 (2017)
41. Zagoruyko, S., Komodakis, N.: Wide residual networks. arXiv preprint arXiv:1605.07146 (2016)
42. Zhang, D., Yang, J., Ye, D., Hua, G.: LQ-Nets: Learned quantization for highly accurate and compact deep neural networks. In: Proceedings of the European Conference on Computer Vision (ECCV). pp. 365–382 (2018)
43. Zhou, S., Wu, Y., Ni, Z., Zhou, X., Wen, H., Zou, Y.: Dorefa-net: Training low bitwidth convolutional neural networks with low bitwidth gradients. arXiv preprint arXiv:1606.06160 (2016)
44. Zhu, C., Han, S., Mao, H., Dally, W.J.: Trained ternary quantization. International Conference on Learning Representations (ICLR) (2017)
45. Zoph, B., Vasudevan, V., Shlens, J., Le, Q.V.: Learning transferable architectures for scalable image recognition. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 8697–8710 (2018)
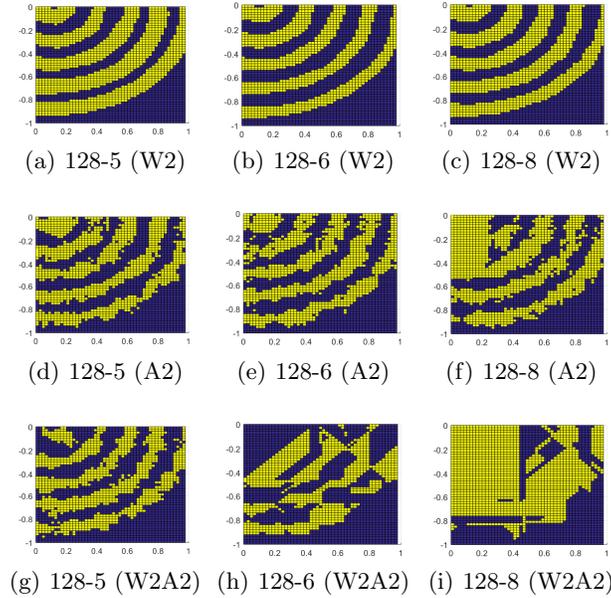
# Appendix

## A. Additional Visualization Examples on the Synthetized Dataset



|     |     |     |     |     |
|-----|-----|-----|-----|-----|
| (a) 200-4 | (b) 384-4 | (c) 128-5 | (d) 128-6 | (e) 128-8 |
| (f) 200-4 | (g) 384-4 | (h) 128-5 | (i) 128-6 | (j) 128-8 |
| (k) 200-4 | (l) 384-4 | (m) 128-5 | (n) 128-6 | (o) 128-8 |

**Fig. 9.** Synthetic data visualization results on FCDNNs as the width and depth increase. The experimented FCDNNs are denoted as '$D$'-'$H$', where $D$ is the dimension of hidden layers and $H$ is the number of layers. The models are either 2-bit wegiht quantized (**first row**), 2-bit activation quantized (**second row**), or 2-bit weight and activation quantized (**third row**).
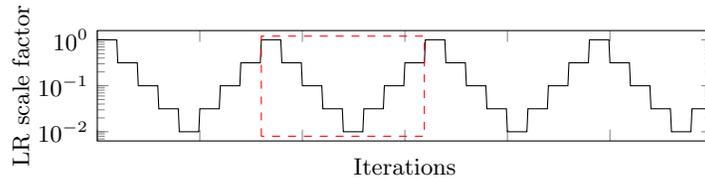
Additional examples of visualization with synthetic dataset are summarized in Fig. 9. The 200-4 FCDNN and 128-6 FCDNN contain 81K and 66K parameters, respectively. The deeper networks, 128-6 and 128-8 FCDNNs, are more robust to weight quantization. However, the wider one, 200-4 FCDNN is more resilient to activation quantization than the deeper networks. As discussed in Section 3 of the main contents, increasing the width alleviates the effects of weight and activation quantization errors, whereas increasing the depth only reduces the effects of weight quantization error.

The visualization results with the residual connections are shown in Fig. 10. As discussed in the main contents, increasing the depth with residual connections rather amplifies the activation quantization error when the shortcut outputs are quantized.

(a) 128-5 (W2)   (b) 128-6 (W2)   (c) 128-8 (W2)

(d) 128-5 (A2)   (e) 128-6 (A2)   (f) 128-8 (A2)

(g) 128-5 (W2A2)   (h) 128-6 (W2A2)   (i) 128-8 (W2A2)

**Fig. 10.** Synthetic data visualization results on FCDNNs with quantized residual connections.

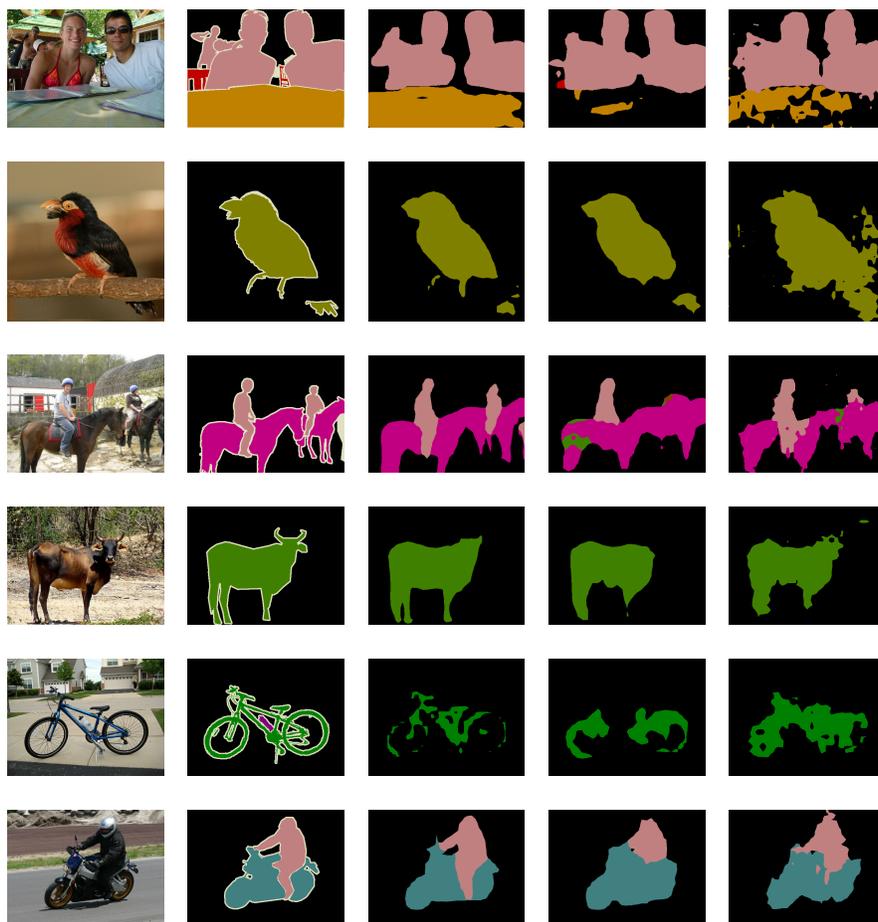## B. Details on Cyclic Learning Rate Scheduling



**Fig. 11.** Learning rate scale factor along training iterations. The red dashed box indicates a single cycle of CLR scheduling.

The cyclic learning rate (CLR) scheduling is shown in Fig. 11. The scale factor is multiplied to the base LR and changes 8 times in a single cycle; $\{1.0, \sqrt{0.1}, 0.1, 0.1\sqrt{0.1}, 0.01, 0.1\sqrt{0.1}, 0.1, \sqrt{0.1}\}$. The cycle period and base LR for each dataset are as follows:

**CIFAR-10:** The base LR is 1e-4, which is 10 times larger than the last LR of the retraining. The cycle period is 8 epochs, thus the LR scaling factor changes at every epoch.

**PASCAL VOC segmentation:** The last LR of the retraining becomes 0 with the polynomial policy for this task. We set the base LR to 1e-5, which is 10 times smaller than the initial LR of the retraining. The cycle period is 3K iterations.

## C. Additional Visualization Examples on the PASCAL VOC Segmentation Benchmark



**Fig. 12.** Visualization examples on the PASCAL VOC segmentation validation set. The first and second columns represent the input images and the ground truth labels. Third to fifth columns show the visualization results of the floating-point, 2-bit weight quantized, and 2-bit activation quantized models, respectively.

## D. Details of Experimented Model Structures

**Table 5.** Quantization performance of various model structures on CIFAR-10 testset. '$L$', '$D_{init}$', and '$B$' represent the number of layers, initial channel dimension, and number of blocks, respectively.

| Block type | Model | | | | Test Acc(%) | | | |
|---|---|---|---|---|---|---|---|---|
| | $L$ | $D_{init}$ | $B$ | Params. | Float | W2AF | WFA2 | W2A2 |
| Basic block | 8 | 59 | $1 \times 3$ | 1.05 M | 93.26 | 92.83 | 92.58 | 92.32 |
| | 14 | 40 | $2 \times 3$ | 1.10 M | 93.62 | 93.51 | **93.30** | **93.11** |
| | 20 | 31 | $3 \times 3$ | 1.03 M | 94.14 | 93.57 | 92.94 | 92.65 |
| | 32 | 24 | $5 \times 3$ | 1.06 M | 94.24 | 93.71 | 92.42 | 91.84 |
| | 44 | 20 | $7 \times 3$ | 1.04 M | **94.27** | **93.91** | 91.40 | 90.69 |
| | 56 | 18 | $9 \times 3$ | 1.09 M | 93.58 | 93.46 | 90.58 | 89.77 |
| Pre-activation block | 8 | 59 | $1 \times 3$ | 1.05 M | 93.10 | 92.75 | 91.89 | 91.03 |
| | 14 | 40 | $2 \times 3$ | 1.10 M | 94.14 | 93.85 | **92.44** | **91.81** |
| | 20 | 31 | $3 \times 3$ | 1.03 M | 94.25 | 93.82 | 92.42 | **91.81** |
| | 32 | 24 | $5 \times 3$ | 1.06 M | 94.07 | 93.89 | 91.51 | 89.03 |
| | 44 | 20 | $7 \times 3$ | 1.04 M | **94.31** | **94.04** | 90.58 | 86.26 |
| | 56 | 18 | $9 \times 3$ | 1.09 M | 93.96 | 93.83 | 88.45 | 85.45 |
| Depthwise block | 29 | 38 | $3 \times 3$ | 1.11 M | 93.64 | 93.28 | **93.38** | 92.33 |
| | 38 | 32 | $4 \times 3$ | 1.08 M | 93.65 | 93.47 | 92.93 | **92.57** |
| | 56 | 26 | $6 \times 3$ | 1.11 M | 94.27 | 94.02 | 91.70 | 90.61 |
| | 74 | 22 | $8 \times 3$ | 1.08 M | 93.94 | 93.71 | 90.50 | 89.88 |
| | 110 | 18 | $12 \times 3$ | 1.13 M | **94.43** | **94.39** | 88.75 | 87.47 |

The details of the model structure for Fig. 8 in the main contents are shown in Table 5. The number of parameters in all the experimented models is approximately one million. All CNNs consist of three block groups. '$B$' indicates the block structure; (number of blocks in each group) × (number of groups). The first block in second and third block groups have a stride of 2. $D$ represents the channel width of the first layer and the channel width increases by 2 times at the beginning of second and third block groups. Floating-point and weight quantized networks perform better on deeper models, while shallow and wider models shows better performance when activations are quantized.