

# Training High-Performance and Large-Scale Deep Neural Networks with Full 8-bit Integers

Yukuan Yang<sup>a</sup>, Shuang Wu<sup>a</sup>, Lei Deng<sup>b\*</sup>, Tianyi Yan<sup>c</sup>, Yuan Xie<sup>b</sup>, Guoqi Li<sup>a\*</sup>

<sup>a</sup>Department of Precision Instrument, Center for Brain Inspired Computing Research, Tsinghua University, Beijing 100084, China.

<sup>b</sup>Department of Electrical and Computer Engineering, University of California, Santa Barbara, CA 93106, USA.

<sup>c</sup>School of Life Science, Beijing Institute of Technology, Beijing 100084, China.

**Abstract**—Deep neural network (DNN) quantization converting floating-point (FP) data in the network to integers (INT) is an effective way to shrink the model size for memory saving and simplify the operations for compute acceleration. Recently, researches on DNN quantization develop from inference to training, laying a foundation for the online training on accelerators. However, existing schemes leaving batch normalization (BN) untouched during training are mostly incomplete quantization that still adopts high precision FP in some parts of the data paths. Currently, there is no solution that can use only low bit-width INT data during the whole training process of large-scale DNNs with acceptable accuracy. In this work, through decomposing all the computation steps in DNNs and fusing three special quantization functions to satisfy the different precision requirements, we propose a unified complete quantization framework termed as “WAGEUBN” to quantize DNNs involving all data paths including W (Weights), A (Activation), G (Gradient), E (Error), U (Update), and BN. Moreover, the Momentum optimizer is also quantized to realize a completely quantized framework. Experiments on ResNet18/34/50 models demonstrate that WAGEUBN can achieve competitive accuracy on the ImageNet dataset. For the first time, the study of quantization in large-scale DNNs is advanced to the full 8-bit INT level. In this way, all the operations in the training and inference can be bit-wise operations, pushing towards faster processing speed, decreased memory cost, and higher energy efficiency. Our throughout quantization framework has great potential for future efficient portable devices with online learning ability.

**Keywords:** Neural Network Quantization, 8-bit Training, Full Quantization, Online Learning Device

## I. INTRODUCTION

Deep neural networks [1] have achieved state-of-art results in many fields like image processing [2], object detection [3], natural language processing [4], and robotics [5] through learning high-level features from a large amount of input data. However, due to the existence of a huge number of floating-point (FP) values and complex FP multiply-accumulate operations (MACs) in the process of network training and inference, the intensive memory overhead, large computational complexity, and high energy consumption impede the wide deployment of deep learning models. DNN quantization [6] which converts FP MACs to bit-wise operations is an effective way to reduce the memory and computation costs and improve the speed of deep learning accelerators.

With the deepening of research, DNN quantization gradually transfers from inference quantization (BWN [7], XNOR-

Net [8], ADMM [9]) to training quantization (DoReFa [10], GXNOR-Net [11], FP8 [12]). Usually, the inference quantization focuses on the forward pass; while the training quantization further quantizes the backward pass and weight updates. Recently, training quantization becomes a hot topic in the network compression community. Whereas, there are still two major issues in existing schemes. The first issue lies in the incomplete quantization, including two aspects: partial quantization and FP dependency. Partial quantization means that only parts of dataflows, not all of them, are quantized (e.g. DoReFa [10], GXNOR-Net [11] and QBP2 [13]); FP dependency still remains FP values during the training process (e.g. MP [14] and FP8 [12]). The second issue is that the quantization of batch normalization (BN) [15] is ignored by most schemes (e.g. MP-INT [16] and FX Training [17]). BN is an essential layer for the training of DNNs by addressing the problem of the internal covariate shift of each layer’s inputs, especially as the network deepens, allowing a much higher learning rate and less careful weight initialization.

Compared with all the studies above, WAGE [18] is the most thorough work of DNNs quantization, which quantizes the data including W (Weights), A (Activation), G (Gradient), E (Error), U (Update) and replacing each BN layer with a constant scaling factor. WAGE has achieved competitive results on LeNet [19], VGG [20], and AlexNet [21], providing a good inspiration for this work. However, we find that WAGE is difficult to be applied in large-scale DNNs due to the absence of BN layers. Besides, it is known that the gradient descent optimizer such as Momentum [22] or Adam [23] increases the stability and even helps get rid of the local optimum, thus the speed and final performance are significantly improved. A complete quantization should cover the entire training process, including W, A, G, E, U, BN, and the optimizer. Regretfully, up to now, there is still no such solution that can achieve this complete quantization, especially on large-scale DNNs.

To address the issues of incomplete quantization and ignored BN quantization mentioned above and extend quantization framework to large-scale datasets and networks with high performance, we propose a unified complete quantization framework termed “WAGEUBN” to constrain W, A, G, E, U, BN, and the optimizer in the low-bit integer (INT) space. To the best of our knowledge, WAGEUBN is the first complete quantization framework achieving high performance in large-scale datasets, where all computation steps and operands in DNNs are decomposed and quantized.

We mainly make the following efforts to create the complete quantization framework. Firstly, according to the various data

\*Corresponding authors. Email addresses: yyk17@mails.tsinghua.edu.cn (Y. Yang), wus15@mails.tsinghua.edu.cn (S. Wu), leideng@ucsb.edu (L. Deng), yantianyi@bit.edu.cn (T. Yan), Y. Xie (yuanxie@ece.ucsb.edu) and liguqi@mail.tsinghua.edu.cn (G. Li).

distributions and the role the quantized data plays in DNN training, we fuse three quantization functions to satisfy the different precision requirements. Furthermore, we propose a new storage and computing method by introducing a flag bit to expand the data coverage and solve the non-convergence problem caused by insufficient data representation. Last but not least, we quantize BN and Momentum optimizer for the first time, converting all FP operations in DNNs to bit-wise operations. Compared with the full precision DNNs, DNNs under the full 8-bit WAGEUBN framework can achieve about  $4\times$  memory saving. More importantly, the multiplication and accumulation operations of WAGEUBN, which are the main operations in DNNs, can perform  $>3\times$  and  $9\times$  faster in speed,  $10\times$  and  $>30\times$  lower in power,  $9\times$  and  $>30\times$  smaller in circuit area, respectively. Besides, the efficient INT8 multiplication and accumulation operations also make WAGEUBN a big step ahead of most existing quantization schemes in computational costs, whether it is FP8 [12], INT16 [16], FP16 [14] or INT32. In addition to the huge advantages in memory cost, computing speed, energy consumption and circuit area, WAGEUBN also shows competitive accuracy on large-scale networks (ResNet18/34/50) and dataset (ImageNet [24]). What's more, the hardware design of WAGEUBN is much simpler and more efficient because of the complete quantization. Due to the improvement of computing speed and saving in hardware resources and energy consumption, WAGEUBN provides a feasible idea for the architecture design of future efficient online learning chips used in portable devices with limited computational resources. The contributions of this work are twofold, which are summarized as follows:

- We address two main issues existing in most quantization schemes via fully quantizing all the data paths, including W, A, G, E, U, BN, and the optimizer, greatly reducing the memory and compute costs. What's more, we constrain the data to INT8 for the first time, pushing the training quantization to a new bit level compared with the existing FP16, INT16, and FP8 solutions.
- Our quantization framework is validated in large-scale DNN models (ResNet18/34/50) over ImageNet dataset and achieves competitive accuracy with much fewer overheads, indicating great potential for future portable devices with online learning ability.

The organization of this paper is as follows: Section II introduces the related work of DNN quantization; Section III details the WAGEUBN framework; Section IV presents the experiment results of WAGEUBN and the corresponding analyses; Section V summarizes this work and delivers the conclusion.

## II. RELATED WORK

With the wide applications of DNNs, the related compression technologies have been proposed rapidly, among which the quantization plays an important role. The development of DNN quantization can be divided into two stages, inference quantization and training quantization, according to the different quantization objects.

**Inference quantization:** Inference quantization starts from constraining W into  $\{-1, 1\}$  (BWN [7]), replacing complex

FP MACs with simple accumulations. BNN [25] and XNOR-Net [8] further quantize both W and A, making the inference computation dominated by bit-wise operations. However, extremely low bit-width quantization usually leads to significant accuracy loss. For example, when the bit width comes to  $<4$  bits, the accuracy degradation becomes obvious, especially for large-scale DNNs. Instead, the bit width of W and A for inference quantization can be reduced to 8 bits with little accuracy degradation. The study of inference quantization is sufficient for the deep learning inference accelerators. Whereas, this is not enough for efficient online learning accelerators because only the data in the forward pass are considered.

**Training quantization:** To further extend the quantization towards the training stage, DoReFa [10] trains DNNs with low bit-width W, A, and G, while leaving E and BN unprocessed. MP [14] and MP-INT [16] use FP16 and INT16 values, respectively, to constrain W, A, and G. Recently, FP8 [12] further pushes W, A, G, E, and U to 8, 8, 8, 8, and 16-bit FP values, respectively, still leaving BN untouched. QBP2 [13] replaces the conventional BN with range BN and constrains W, A, and E to INT8 values while calculating G with FP MACs. Recently, WAGE [18] adopts a layer-wise scaling factor instead of using the BN layer and quantizes W, A, G, E, and U to 2, 8, 8, 8, and 8 bits, respectively. Despite its thorough quantization, WAGE is difficult to be applied to large-scale DNNs due to the absence of powerful BN layers. In summary, there still lacks a complete INT8 quantization framework for the training of large-scale DNNs with high accuracy.

## III. WAGEUBN FRAMEWORK

The main idea of WAGEUBN is to quantize all the data in DNN training to INT8 values. In this section, we detail the WAGEUBN framework implemented in large-scale DNN models. The organization of this section is as follows: Subsection III-A introduces the straight-through estimator (STE) [26, 27] method which is accepted and used by most researchers to solve the non-differentiable problem of quantization; Subsection III-B and Subsection III-C describe the notations and quantization functions, respectively; Subsection III-D explains the specific quantization schemes for W, A, G, E, U, BN, and the Momentum optimizer, respectively; Subsection III-E goes through the overall implementation of WAGEUBN, including in both forward and backward passes; Subsection III-F summarizes the whole process and shows the pseudo codes.

### A. Straight-through estimator

In the early stage of the study, the non-differentiable problem in mathematical sense caused by quantization has been hindering the development of quantization research. However, since the straight-through estimator (STE) method was used to estimate the gradient of quantized data in BNN [25], almost all the works in the field of quantization have adopted this method to avoid the mathematical non-derivative problem [10, 28, 29, 30].

The STE method used in WAGEUBN can be illustrated as the following

$$\begin{aligned} \text{Forward : } x_q &= Q(x) \\ \text{Backward : } \frac{\partial L}{\partial x} &= \frac{\partial L}{\partial x_q} \end{aligned} \quad (1)$$

where  $x$  is the data to be quantized,  $Q(\cdot)$  is the quantization function that may be non-differentiable, and  $L$  denotes the objective function.

### B. Notations

Before introducing the WAGEUBN quantization framework formally, we need to define some notations. Considering the  $l$ -th layer of DNNs, we divide the forward pass of DNNs into four steps as described in Figure 1 (BN is divided into two steps: Normalization &  $Q_{BN}$  and Scale & Offset).

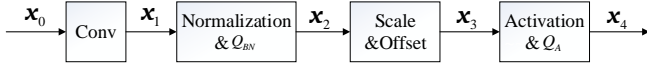


Fig. 1: Forward quantization of the  $l$ -th layer in DNNs.

Specifically, we have

$$\begin{aligned} \text{Input : } x_0^l &= x_4^{l-1} \\ \text{Conv : } x_1^l &= W_q^l x_0^l \\ \text{Normalization\&Q}_{BN} : x_2^l &= Q_{BN}\left(\frac{x_1^l - \mu_q^l}{\sigma_q^l}\right) \quad , \quad (2) \\ \text{Scale\&Offset : } x_3^l &= \gamma_q^l x_2^l + \beta_q^l \\ \text{Activation\&Q}_A : x_4^l &= Q_A[\text{relu}(x_3^l)] \end{aligned}$$

where  $x_4^{l-1}$  is the output of the  $(l-1)$ -th layer and  $x_0^l$  is the input of the  $l$ -th layer;  $W_q^l$  is the quantized weight for convolution;  $Q_{BN}$  is the quantization function used for constraining  $x_2^l$  to low-bit INT values, which will be given in Equation (13);  $\mu_q^l$  and  $\sigma_q^l$  are the quantized mean and standard deviation value of one mini-batch;  $\gamma_q^l$  and  $\beta_q^l$  are the quantized scale and bias used in the BN layer of DNNs;  $Q_A$  is the quantization function for activation which will be detailed in Equation (14);  $\text{relu}$  is the activation function which is commonly used in DNNs. Noting that every step in the forward pass is quantized,  $x_0^l, x_1^l, x_2^l, x_3^l, x_4^l$  are all integers. In order to make the notations used in the paper consistent, we make the following rules: subscript  $q$  denotes the data which have been quantized to INT values and superscript  $l$  denote the layer index.



Fig. 2: Backward quantization of the  $l$ -th layer in DNNs.

Different from most existing schemes, we define  $e$  and  $g$  respectively, where  $e$  represents the gradient of A (activation) which is used in the error backpropagation and  $g$  represents the gradient of W (weights) which is used in the weight update. Moreover, we quantize the BN layers, including both

the forward and backward passes, which is not well touched in most prior work. Similar to the forward pass, we divide the backward pass of the  $l$ -th layer into five steps as shown in Figure 2. According to the derivative chain rules, we have

$$\begin{aligned} Q_{E_1} : e_0^l &= Q_{E_1}\left(\frac{\partial L}{\partial x_4^l}\right) = Q_{E_1}(e_4^{l+1}) \\ \text{Activation : } e_1^l &= \frac{\partial L}{\partial x_3^l} = \frac{\partial L}{\partial x_4^l} \odot \frac{\partial x_4^l}{\partial x_3^l} = e_0^l \odot \frac{\partial x_4^l}{\partial x_3^l} \\ \text{Scale\&Offset : } e_2^l &= \frac{\partial L}{\partial x_2^l} = \frac{\partial L}{\partial x_3^l} \odot \frac{\partial x_3^l}{\partial x_2^l} = e_1^l \odot \gamma_q^l \\ \text{Norm\&Q}_{E_2} : e_3^l &= Q_{E_2}\left(\frac{\partial L}{\partial x_1^l}\right) = Q_{E_2}\left(\frac{\partial L}{\partial x_2^l} \odot \frac{\partial x_2^l}{\partial x_1^l}\right) \\ &= Q_{E_2}(e_2^l \odot \frac{\partial x_2^l}{\partial x_1^l}) \\ \text{Conv : } e_4^l &= \frac{\partial L}{\partial x_0^l} = \frac{\partial x_1^l}{\partial x_0^l} \frac{\partial L}{\partial x_1^l} = W_q^{lT} e_1^l \end{aligned} \quad (3)$$

where  $L$  is the loss function,  $e_4^{l+1}$  represents the error from the  $(l+1)$ -th layer, and  $\odot$  represents the Hadamard product. For vectors with the same dimension, such as  $\mathbf{a} = (a_1, a_2, \dots, a_n)$  and  $\mathbf{b} = (b_1, b_2, \dots, b_n)$ , we have:  $\mathbf{a} \odot \mathbf{b} = (a_1 b_1, a_2 b_2, \dots, a_n b_n)$ . Two quantization functions are used here:  $Q_{E_1}$  is the quantization function detailed as Equation (15) that converts high bit-width integers to low bit-width integers;  $Q_{E_2}$  detailed as Equation (17) is trying to convert FP values to low bit-width integers.  $W_q^{lT}$  is the transposed matrix of  $W_q^l$ , and  $\partial x_4^l / \partial x_3^l$  represents the gradient of activation. When  $\text{relu}$  is used as the activation function,  $\partial x_4^l / \partial x_3^l$  is a tensor containing only 0 and 1 elements.

According to the definitions given above, the gradients of  $W$ ,  $\gamma$ , and  $\beta$  can be summarized as follows

$$\begin{aligned} g_W^l &= \frac{\partial L}{\partial W^l} = \frac{\partial L}{\partial x_1^l} \frac{\partial x_1^l}{\partial W^l} = e_3^l x_0^{lT} \\ g_\gamma^l &= \frac{\partial L}{\partial \gamma^l} = \frac{\partial L}{\partial x_3^l} \odot \frac{\partial x_3^l}{\partial \gamma^l} = e_1^l \odot x_2^l \\ g_\beta^l &= \frac{\partial L}{\partial \beta^l} = \frac{\partial L}{\partial x_3^l} \odot \frac{\partial x_3^l}{\partial \beta^l} = e_1^l. \end{aligned} \quad (4)$$

To further reduce the bit width of  $G$  that will increase greatly after the multiplication, we have

$$\begin{aligned} g_{W_q}^l &= Q_{G_W}(g_W^l) \\ g_{\gamma_q}^l &= Q_{G_\gamma}(g_\gamma^l) \\ g_{\beta_q}^l &= Q_{G_\beta}(g_\beta^l) \end{aligned} \quad (5)$$

where  $Q_{G_W}$ ,  $Q_{G_\gamma}$ , and  $Q_{G_\beta}$  are quantization functions for the gradient of  $W$ ,  $\gamma$ , and  $\beta$ , respectively, which will be shown in Equation (18).

Some notations to be used below are also explained here.  $k_W$ ,  $k_A$ ,  $k_{G_W}$ ,  $k_E$  ( $k_{E_1}$  and  $k_{E_2}$ ), and  $k_{BN}$  are the bit width of  $W$ ,  $A$ ,  $G$ ,  $E$ , and  $BN$ , respectively.  $k_{WU}$ ,  $k_{\gamma U}$ , and  $k_{\beta U}$  are the bit width of  $W$ ,  $\gamma$ , and  $\beta$  update, which are also the bit width of data stored in memory.  $k_\gamma$ ,  $k_\beta$ ,  $k_\mu$ , and  $k_\sigma$  are the bit width of  $\gamma$ ,  $\beta$ ,  $\mu$ , and  $\sigma$ , respectively, used in the BN layer.  $k_{G_\gamma}$  and  $k_{G_\beta}$  are the bit width of  $\gamma$  and  $\beta$  gradient, respectively.  $k_{Mom}$  and  $k_{Acc}$  are the bit width of momentum coefficient ( $Mom$ ) and accumulation ( $Acc$ ), respectively, used in the Momentum optimizer. At last,  $k_{lr}$  is the bit width of

the learning rate.

### C. Quantization Functions

There are three quantization functions used in WAGEUBN. The direct-quantization function uses the nearest fixed-point values to represent the continuous values of W, A, and BN. The constant-quantization function for G is used to keep the bit width of U (update) fixed since G is directly related to U. Because U and the weights stored in memory have the same bit width, the bit width of weights stored in memory can be fixed, which is more hardware-friendly. The magnitude of E is very small, so the shift-quantization function reduces the bit width of E greatly compared with the direct-quantization function under the same precision.

#### (1) Direct-quantization function

The direct-quantization function simply approximates a continuous value to its nearest discrete state and is defined as

$$Q(\mathbf{x}, k) = \frac{\text{round}(\mathbf{x} \cdot 2^{k-1})}{2^{k-1}} \quad (6)$$

where  $k$  is the bit width, and  $\text{round}(\cdot)$  rounds a number to its nearest INT value.

#### (2) Constant-quantization function

The intention of constant-quantization function is to normalize a tensor firstly, then limit it to INT, and finally maintain its magnitude. It is governed by

$$\begin{aligned} R(\mathbf{x}) &= 2^{\text{round}(\log_2(\max(|\mathbf{x}|)))} \\ Sr(\mathbf{x}) &= \begin{cases} \lfloor \mathbf{x} \rfloor, & P_x = \lceil \mathbf{x} \rceil - \mathbf{x} \\ \lceil \mathbf{x} \rceil, & P_x = \mathbf{x} - \lfloor \mathbf{x} \rfloor \end{cases} \\ Norm(\mathbf{x}) &= \frac{\mathbf{x}}{R(\mathbf{x})} \\ Sd(\mathbf{x}) &= \text{clip}\{Sr[dr \cdot Norm(\mathbf{x})], -dr + 1, dr - 1\} \\ CQ(\mathbf{x}) &= \frac{Sd(\mathbf{x})}{2^{k_{GC}-1}}. \end{aligned} \quad (7)$$

The illustration of constant-quantization function is described in Figure 3. Here,  $R(\cdot)$  is used to project the maximum value of  $\mathbf{x}$  to its nearest fixed-point value, which is prepared for normalization;  $Sr(\cdot)$  is a stochastic rounding function used for converting a continuous float value to its nearby INT value in a probabilistic manner and  $P_x$  is the rounding probability;  $Norm(\cdot)$  denotes normalization and  $clip(\cdot)$  is a saturation function limiting the data range;  $Sd(\cdot)$  is to shift the distribution of  $\mathbf{x}$  and limit  $\mathbf{x}$  to INT values between  $-dr + 1$  and  $dr - 1$ . Here  $dr \in [2^{k-1}, 2^{k-2}, \dots, 1]$  limits the data range after mapping and decreases as the training goes on, presenting the same effect as reducing the learning rate. For example,  $Sd(\cdot)$  maps G to  $\{-127, -126, \dots, 126, 127\}$  and  $\{-63, -62, \dots, 62, 63\}$  in the early training stage ( $k = 8$ ,  $dr = 128$ , epoch in  $[0, 30]$ ) and later training stage ( $k = 7$ ,  $dr = 64$ , epoch in  $[30, 60]$ ), respectively.  $CQ(\cdot)$  is utilized to maintain the magnitude order of data, where  $2^{k_{GC}-1}$  is a constant scaling factor and  $k_{GC}$  is its bit width.

#### (3) Shift-quantization function

The shift-quantization function serves for the quantization of E and is defined as

$$d(k) = \frac{1}{2^{k-1}}$$

$$SQ(\mathbf{x}, k) = R(\mathbf{x}) \cdot \text{clip}\{Q[Norm(\mathbf{x}), k], -1 + d(k), 1 - d(k)\} \quad (8)$$

where  $d(\cdot)$  is the minimum interval for a  $k$ -bit INT and  $Q(\cdot)$  is the direct-quantization function defined in Equation (6).

The shift-quantization function normalizes E first, then converts E to fixed-point values, and finally uses a layer-wise scaling factor ( $R(\cdot)$  defined in Equation (7)) to maintain the magnitude. The differences between the constant-quantization function and the shift-quantization function mainly exist in two points: First, the constant-quantization uses a constant to keep the magnitude for hardware friendliness while the shift-quantization uses a layer-wise scaling factor; Second, the constant-quantization contains a stochastic rounding process while the shift-quantization function does not.

### D. Quantization Schemes in WAGEUBN

After introducing the quantization functions used in our WAGEUBN framework, we provide detailed quantization schemes.

#### (1) Weight Quantization

Since weights are stored and used as fixed-point values, weights should be also initialized discretely. An initialization method proposed by MSRA [31] has been evidenced helpful for faster training. The initialization of weights can be formulated as follows

$$\begin{aligned} W' &\sim N(0, \frac{1}{\sqrt{n_{in}}}) \\ W &= \text{clip}[Q(W', k_{WU}), -1 + d(k_{WU}), 1 - d(k_{WU})] \end{aligned} \quad (9)$$

where  $n_{in}$  is the layer's fan-in number, and  $k_{WU}$  is the bit width of weight update and the memory storage.

Because of the different bit width for weight storage and computation, it should be quantized from  $k_{WU}$  bits to  $k_W$  (the bit width of weights used for convolution) bits for convolution. In addition, we also limit the data range of W. Finally, the quantization function for W is

$$Q_W(\mathbf{x}) = \text{clip}[Q(\mathbf{x}, k_W), -1 + d(k_W), 1 - d(k_W)]. \quad (10)$$

#### (2) Batch Normalization Quantization

As aforementioned, BN plays an important role in training large-scale DNNs. WAGE [18] has proved that simple scaling layers are not enough to replace BN layers. Conventional BN layer can be divided into two steps as

$$\begin{aligned} \hat{\mathbf{x}} &= \frac{\mathbf{x} - \mu^l}{\sqrt{\sigma^{l2} + \epsilon}} \\ \mathbf{y} &= \gamma^l \hat{\mathbf{x}} + \beta^l \end{aligned} \quad (11)$$

where  $\mu^l$  and  $\sigma^l$  are the mean and standard, respectively, deviation of  $\mathbf{x}$  over one mini-batch in the  $l$ -th layer;  $\epsilon$  is a small positive value added to  $\sigma$  to avoid the case of dividing by zero;  $\gamma^l$  and  $\beta^l$  are the scale and offset parameters, respectively.



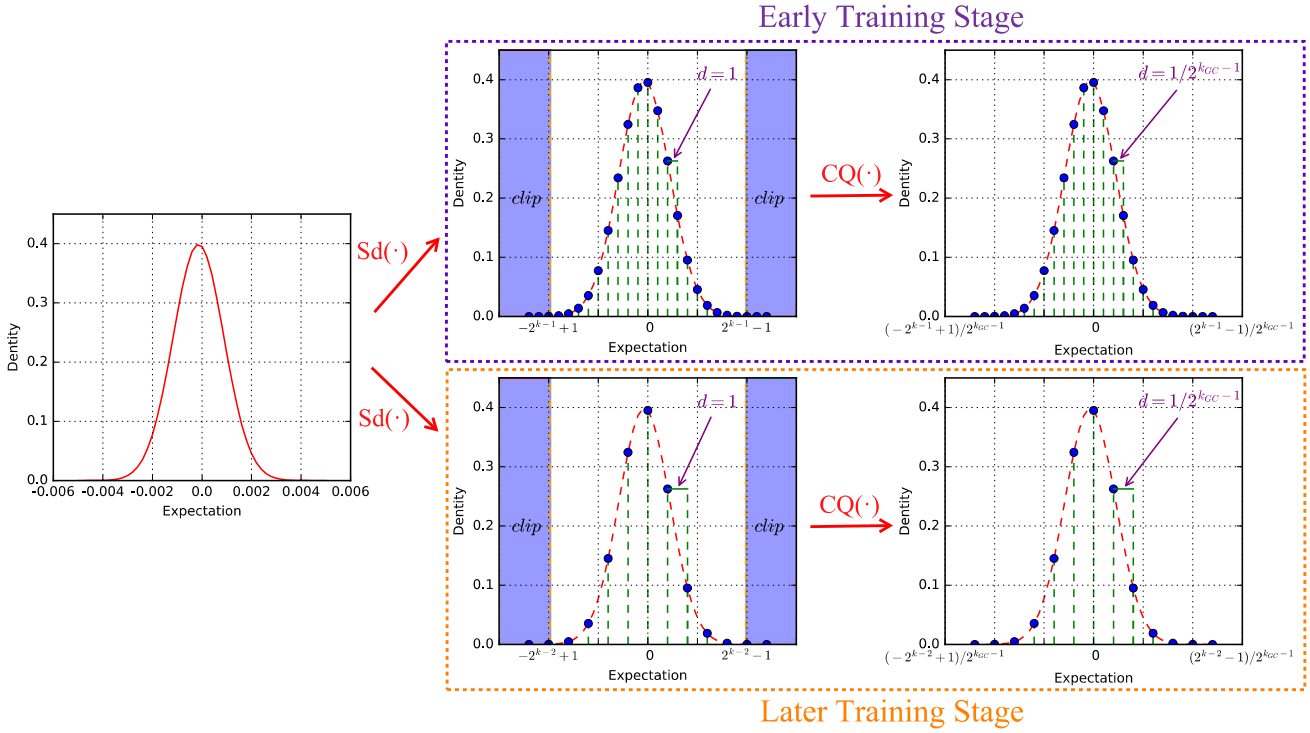


Fig. 3: Illustration of the constant-quantization function during training.

Under the WAGEUBN framework, the BN layer is also quantized. Through the operations described in Equation (12), all operands are quantized and all operations are bit-wise. Specifically, the quantization follows

$$\begin{aligned}\mu_q^l &= Q_\mu(\mu^l), \quad \sigma_q^l = Q_\sigma(\sigma^l) \\ \hat{x} &= Q_{BN}\left(\frac{x - \mu_q^l}{\sigma_q^l + \epsilon_q}\right) \\ \gamma_q^l &= Q_\gamma(\gamma^l), \quad \beta_q^l = Q_\beta(\beta^l) \\ y &= \gamma_q^l \hat{x} + \beta_q^l\end{aligned}\quad (12)$$

where  $Q_\mu, Q_\sigma, Q_\gamma, Q_\beta, Q_{BN}$  are the quantization functions converting the operands to fixed-point values defined as

$$\begin{aligned}Q_\mu(x) &= Q(x, k_\mu), \quad Q_\sigma(x) = Q(x, k_\sigma) \\ Q_\gamma(x) &= Q(x, k_\gamma), \quad Q_\beta(x) = Q(x, k_\beta) \\ Q_{BN}(x) &= Q(x, k_{BN}).\end{aligned}\quad (13)$$

And  $\epsilon_q$  is a small fixed-point value, playing the same role as  $\epsilon$  in Equation (11);  $k_\mu, k_\sigma, k_\gamma, k_\beta, k_{BN}$  are the bit width of  $\mu, \sigma, \gamma, \beta$  and  $\hat{x}$ , respectively.

### (3) Activation Quantization

After the convolution and BN layers in the forward pass, the bit width of operands increases due to the multiplication operation. To reduce the bit width and keep the input bit width of each layer consistent, we need to quantize the activations. Here, the quantization function for activations can be described as

$$Q_A(x) = Q(x, k_A) \quad (14)$$

where  $k_A$  is the bit width of activations.

### (4) Error Quantization

In Equation (3), we have given the definition of E and quantized E. Through investigating the importance of error propagation in DNN training, we find that the quantization of E is very essential for the model convergence. If E is naively quantized using the direct-quantization function, it will require a large bit width of operands to realize the convergence of DNNs. Instead, we use the following shift-quantization function

$$Q_{E_1}(x) = SQ(x, k_{E_1}) \quad (15)$$

where  $SQ(\cdot)$  is the shift-quantization function defined in Equation (8), and  $k_{E_1}$  is the bit width of  $e_0^l$  defined in Equation (3).

As mentioned above, we use  $Q_{E_1}$  and  $Q_{E_2}$  for the error quantization. However, the precision requirements of  $Q_{E_1}$  and  $Q_{E_2}$  vary a lot. Experiments show that  $k_{E_1} = 8$  affects little on accuracy while  $k_{E_2} \leq 8$  will cause the non-convergence of large-scale DNNs when using  $SQ(\cdot)$  as the quantization function.  $k_{E_2} = 16$  is a proper value for the training of DNNs with minimum accuracy degradation. More analyses will be given in Subsection IV-E. Here we will provide two versions of  $Q_{E_2}$ , the 16-bit and 8-bit versions. The 16-bit  $Q_{E_2}$  is defined as

$$Q_{E_2}(x) = SQ(x, k_{E_2}) \quad (16)$$

where  $k_{E_2}$  is the bit width of  $e_3^l$  defined in Equation (3).

Experiments have proved the data range covered by 8-bit  $Q_{E_2}$  ( $k_{E_2} = 8$ ) is not sufficient to train DNNs. In order

to expand the coverage of quantization function while still maintaining a low bit width, we introduce a layer-wise scaling factor  $Sc$  and a flag bit. Then, to distinguish it from  $Q_{E_2}$  defined in Equation (16), we name the quantization function  $Q_{E_2}$  and the quantization process is governed by

$$Sc = \frac{R(\mathbf{x})}{2^{k_{E_2}-1}}$$

$$Q_{E_2}(\mathbf{x}) = \begin{cases} Sc \cdot \text{clip} \left[ \text{round}(\frac{\mathbf{x}}{Sc}), \min, \max \right], & |\frac{\mathbf{x}}{Sc}| \geq 1 \\ Sc \cdot Q(\frac{\mathbf{x}}{Sc}, k_{E_2}), & |\frac{\mathbf{x}}{Sc}| < 1 \end{cases} \quad (17)$$

where  $k_{E_2} = 8$ ,  $\min = -2^{k_{E_2}-1}$ , and  $\max = 2^{k_{E_2}-1}$ .

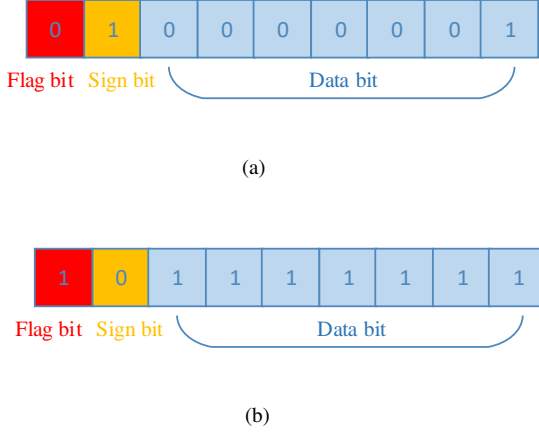


Fig. 4: Data format of 9-bit integers.

By introducing a layer-wise scaling factor and a flag bit, we can expand the data coverage greatly. Details can be found in Figure 4. The flag bit is used to indicate whether the absolute value of  $x$  stored in memory is less than the layer-wise scaling factor (e.g., 0 represents  $|x| < Sc$  and 1 represents  $|x| \geq Sc$ ). The sign bit is used to denote the positive or negative direction of the value. The data bit follows the conventional binary format. According to the definition, the values stored in Figure 4(a) and 4(b) are  $+Sc/128$  and  $-127 \times Sc$  when  $k_{E_2} = 8$ , respectively. Therefore, the 9-bit data format can cover almost the same data range as the direct 15-bit quantization described in Equation (16). Since the flag bit is just used for judgment, the effective value for computation is still INT8.

#### (5) Gradient Quantization

The gradient is another important part in DNN training because it is directly related to the weight update. The rules for calculating and quantizing the gradients of  $W$ ,  $\gamma$ , and  $\beta$  are described as Equation (4) and (5). Since  $e_1^l$ ,  $e_3^l$ ,  $x_0^l$ , and  $x_2^l$  are all fixed-point values, the conventional FP MACs operations can be replaced with bit-wise operations during the process of calculating  $g_W^l$ ,  $g_\gamma^l$ , and  $g_\beta^l$ . The quantization functions are defined to further reduce the bit width of gradients and prepare

for the next step of the optimizer. Specifically, we have

$$\begin{aligned} Q_{G_W}(\mathbf{x}) &= CQ(\mathbf{x}, k_{G_W}) \\ Q_{G_\gamma}(\mathbf{x}) &= Q(\mathbf{x}, k_{G_\gamma}) \\ Q_{G_\beta}(\mathbf{x}) &= Q(\mathbf{x}, k_{G_\beta}) \end{aligned} \quad (18)$$

where  $CQ(\cdot)$  is the constant-quantization function defined in Equation (7);  $k_{G_W}$ ,  $k_{G_\gamma}$ , and  $k_{G_\beta}$  are the bit width of the gradient of  $W$ ,  $\gamma$ , and  $\beta$ , respectively.

#### (6) Momentum Optimizer Quantization

Momentum optimizer is one of the most common optimizers used in DNN training, especially for classification tasks. For the  $i$ -th training step of the  $l$ -th layer, the conventional Momentum optimizer works as follows

$$Acc_i^l = Mom \cdot Acc_{i-1}^l + g_i^l \quad (19)$$

where  $Acc_i^l$  and  $Acc_{i-1}^l$  are the accumulation in the  $i$ -th and  $(i-1)$ -th training step, respectively;  $Mom$  is a constant value used as a coefficient;  $g_i^l$  is the gradient of  $W$ ,  $\gamma$ , or  $\beta$ .

Momentum optimizer under the WAGEUBN framework is trying to constrain all operands to fixed-point values. The process can be formulated as

$$\begin{aligned} Acc_i^l &= Mom \cdot Acc_{(i-1)q}^l + g_{iq}^l \\ Acc_{iq}^l &= Q_{Acc}(Acc_i^l) \end{aligned} \quad (20)$$

where  $Acc_{(i-1)q}^l$  is the quantized accumulation in the  $(i-1)$ -th training step;  $g_{iq}^l$  is the quantized gradient of  $W$ ,  $\gamma$ , or  $\beta$ ;  $Q_{Acc}(\cdot)$  is the quantization function defined as

$$Q_{Acc}(\mathbf{x}) = Q(\mathbf{x}, k_{Acc}). \quad (21)$$

To guarantee the consistency of bit width, we further set

$$k_{G_\gamma} = k_{G_\beta} = k_{GC} = k_{Mom} + k_{Acc} - 1. \quad (22)$$

#### (7) Update Quantization

The parameter update is the last step in the training of each mini-batch. Different from conventional DNNs where the learning rate can take any FP value, the learning rate under WAGEUBN must also be a fixed-point value and the bit width of update is directly related to the bit width of learning rate. The update under quantized Momentum optimizer can be described as

$$\begin{aligned} \Delta W &= lr \cdot Acc_i^l \\ W^l &= W^l - \Delta W \end{aligned} \quad (23)$$

where  $\Delta W$  is the update of  $W$  with  $k_{WU}$  bits, and  $lr$  is the fixed-point learning rate with  $k_{lr}$  bits. The updates of  $\gamma$  and  $\beta$  are the same as in Equation (23). According to Equation (20), (22), and (23), we have

$$\begin{aligned} k_{WU} &= k_{\gamma U} = k_{\beta U} = k_{Mom} + k_{Acc} + k_{lr} - 2 \\ &= k_{GC} + k_{lr} - 1 \\ &= k_{G_\gamma} + k_{lr} - 1 \\ &= k_{G_\beta} + k_{lr} - 1. \end{aligned} \quad (24)$$

Through our evaluations, the precision of the update has the

greatest impact on the accuracy of DNNs because it is the last step to constrain the parameters. Thus, we need to set a reasonable bit width for update to balance the model accuracy and memory cost.

#### E. Quantization Framework

Given the quantization details of W, A, G, E, U, BN, and the Momentum optimizer, the overall quantization framework is depicted in Figure 5. Under this framework, conventional FP MACs can be replaced with bit-wise operations. Here, the forward pass of the  $l$ -th layer in DNNs is divided into three parts: Conv (convolution), BN, and activation.  $x_0^l, x_1^l, x_2^l, x_3^l$ , and  $x_4^l$  are defined in Equation (2). The weights ( $W^l$ ) are stored as  $k_{WU}$ -bit integers and then  $Q_W$  maps  $W^l$  to  $k_W$ -bit INT values ( $W_q^l$ ) before convolution. After convolution,  $MEAN \& Q_\mu$  and  $STD \& Q_\sigma$  operations are used to calculate the mean and standard deviation of  $x_1^l$  in one mini-batch and then quantize them to  $k_\mu$  and  $k_\sigma$  bits, respectively. The  $BW \& Q_{BN}$  operation constrains  $x_2^l$  to  $k_{BN}$  bits in BN. Similar to W,  $\gamma$  and  $\beta$  are stored as  $k_{\gamma U}$  and  $k_{\beta U}$ -bit integers and used in  $k_\gamma$  and  $k_\beta$  bits ( $\gamma_q^l$  and  $\beta_q^l$ ) after the  $Q_\gamma$  and  $Q_\beta$  quantization, respectively. After the second step of BN, activation and quantization are implemented with the  $ACT \& Q_A$  operation, reducing the increased bit width to  $k_A$  bits again and preparing inputs for the next layer.

The backward pass of the  $l$ -th layer is much more complicated than the forward pass, including error propagation, gradient of weight, gradient of BN, Momentum optimizer, and weight update. In the process of error propagation,  $e_0^l, e_1^l, e_2^l, e_3^l$ , and  $e_4^l$  are defined in Equation (3) and there are two locations needing quantization using  $Q_{E_1}$  and  $Q_{E_2}$ .  $Q_{E_1}$  reduces the bit width of  $e_4^{l+1}$  from  $k_{E_2} + k_W - 1$  to  $k_{E_1}$ .  $ACT'$  is the derivative of activation function ( $relu$ ) and  $Q_{E_2}$  is used to constrain  $e_3^l$  to  $k_{E_2}$  bits. In the phase of calculating the gradients of weights and BN,  $Q_{G_W}$ ,  $Q_{G_\gamma}$ , and  $Q_{G_\beta}$  are leveraged to reduce the increased bit width caused by the multiplication operations.

All parameters of the Momentum optimizer in the  $i$ -th training step are quantized. Different from the conventional learning rate with FP value, WAGEUBN requires a discrete learning rate so that the bit width of weight updates can be controlled. The updates of  $\gamma$  and  $\beta$  in BN layers are also similar to the weight update, which are omitted in Figure 5 for simplicity.

#### F. Overall Algorithm

Given the framework of WAGEUBN, we summarize the entire quantization process and present the pseudo codes for both the forward and backward passes as shown in Algorithm 1 and 2, respectively.

### IV. RESULTS

#### A. Experimental Setup

To verify the effectiveness of the proposed quantization framework, we apply WAGEUBN on ResNet18/34/50 on ImageNet dataset. We provide two versions of WAGEUBN,

---

#### Algorithm 1 Forward pass of $l$ -th layer

---

Convolution:

$$x_0^l \leftarrow x_4^{l-1}, W_q^l \leftarrow Q_W(W^l, k_W)$$

$$x_1^l \leftarrow W_q^l x_0^l$$

BN:

$$\mu_q^l \leftarrow Q_\mu(\mu^l), \sigma_q^l \leftarrow Q_\sigma(\sigma^l)$$

$$x_2^l \leftarrow Q_{BN}\left(\frac{x_1^l - \mu_q^l}{\sigma_q^l}\right)$$

$$\gamma_q^l \leftarrow Q_\gamma(\gamma^l), \beta_q^l \leftarrow Q_\beta(\beta^l)$$

$$x_3^l \leftarrow \gamma_q^l x_2^l + \beta_q^l$$

Activation:

$$x_4^l \leftarrow Q_A(relu(x_3^l))$$


---

---

#### Algorithm 2 Backward pass of $l$ -th layer

---

Error propagation:

$$e_0^l \leftarrow Q_{E_1}(e_4^{l+1})$$

$$e_1^l \leftarrow e_0^l \odot \frac{\partial x_4^l}{\partial x_3^l}$$

$$e_2^l \leftarrow e_1^l \odot \gamma_q^l$$

$$e_3^l \leftarrow Q_{E_2}(e_2^l \odot \frac{\partial x_2^l}{\partial x_1^l})$$

$$e_4^l \leftarrow W_q^{lT} e_3^l$$

Gradient computation:

$$g_W^l \leftarrow e_3^l x_0^{lT}$$

$$g_{W_q}^l \leftarrow Q_{G_W}(g_W^l)$$

$$g_\gamma^l \leftarrow e_1^l \odot x_2^l$$

$$g_{\gamma_q}^l \leftarrow Q_{G_\gamma}(g_\gamma^l)$$

$$g_\beta^l \leftarrow e_1^l$$

$$g_{\beta_q}^l \leftarrow Q_{G_\beta}(g_\beta^l)$$

Momentum optimizer( $i$ -th step):

$$Acc_i^l = Mom \cdot Acc_{(i-1)q}^l + g_{iq}^l(g_{W_q}^l, g_{\gamma_q}^l, \text{or } g_{\beta_q}^l)$$

$$Acc_{iq}^l = Q_{Acc}(Acc_i^l)$$

Weight updates( $i$ -th step):

$$\Delta W = lr \cdot Acc_i^l$$

$$W^l = W^l - \Delta W$$


---

one with full 8-bit INT where  $k_W, k_A, k_{G_W}, k_{E_1}, k_{E_2}, k_\gamma$ , and  $k_\beta$  are equal to 8. The other version has 16-bit  $k_{E_2}$ . The only difference between the 16-bit  $E_2$  version and the full 8-bit version exists in the quantization function  $Q_{E_2}$  (see Equation (16) and (17), respectively).  $k_{G_\gamma}, k_{G_\beta}$ , and  $k_{GC}$  are 15 and we set  $k_{mom}, k_{Acc}, k_{lr}$ , and  $k_{WU}$  to 3, 13, 10, and 24 respectively to satisfy Equation (22) and (24). In addition, we set  $k_{BN}, k_\mu$ , and  $k_\sigma$  to 16. Since W, A, G, and E occupy the majority of memory and compute costs, their bit-width values are reduced as much as possible. Other parameters occupying much less

## Training Deep Neural Networks with 8-bit integers

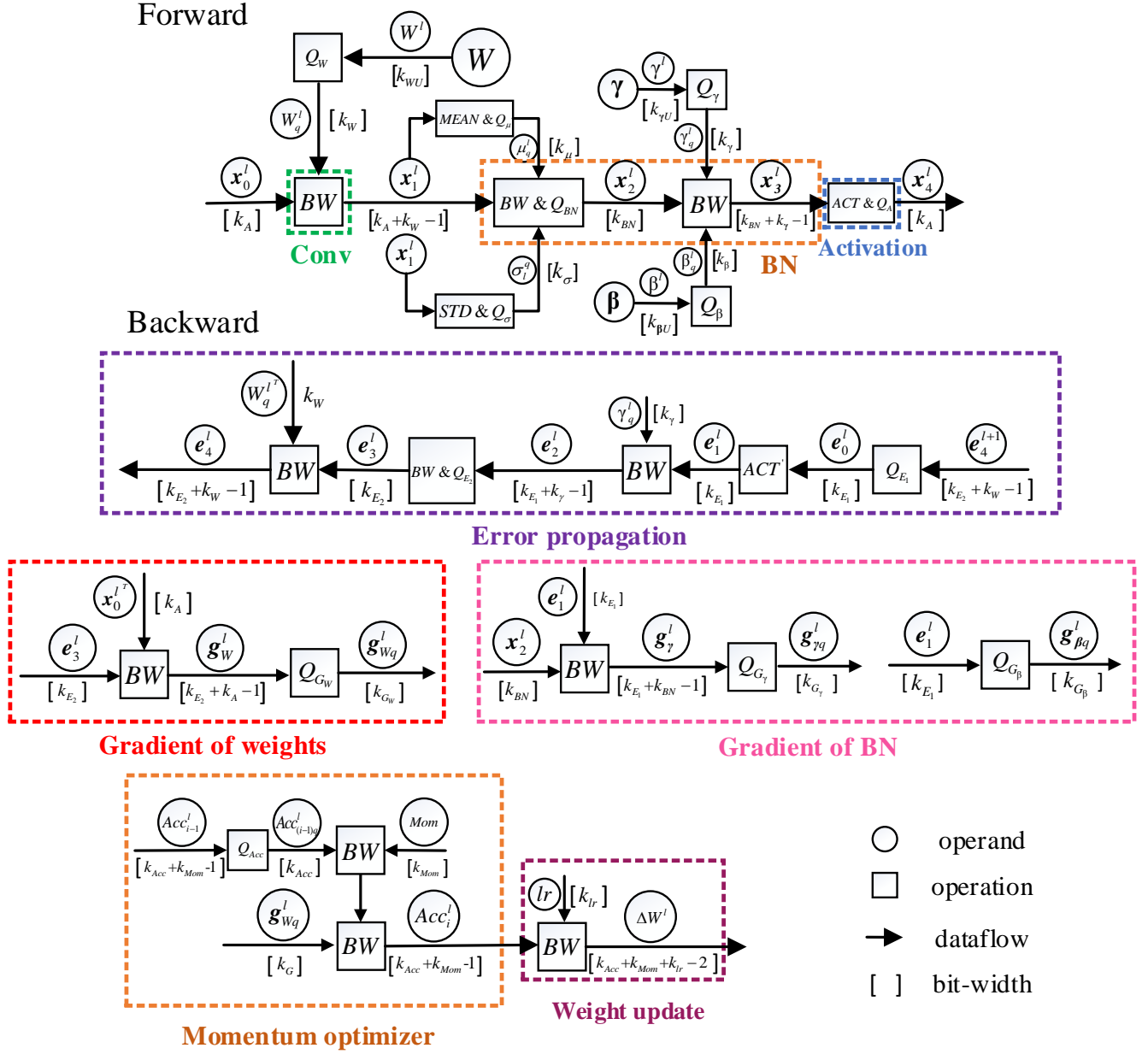


Fig. 5: Overview of the WAGEUBN quantization framework. “BW” denotes bit-wise operations.

resources can increase the bit width to maintain the accuracy, e.g.,  $\mu_q^l$ ,  $\sigma_q^l$ ,  $g_{\gamma q}^l$ , and  $g_{\beta q}^l$  in BN layers.

The first and last layers are believed to differ from the rest because of their interface with network inputs and outputs. The quantization of these two layers will cause significant accuracy degradation compared to hidden layers and they just consume few overheads due to the small number of neurons. Therefore, we do not quantize the first and last layers, as previous work did [32, 33].

For the quantization frameworks, it’s a common issue of saturated learning and vanishing gradients because the data in DNNs may be clipped very small after quantization and the model may be stuck on a local optimum because of the few

updates. And this situation can be more serious when the quantization framework is used in large-scale complex datasets. WAGEUBN has made many efforts to alleviate this problem. Firstly, the newly designed shift-quantization function ensures that the quantized error in WAGEUBN framework is on the same magnitude as that of a traditional FP network, making the errors (gradients of activations) not vanishing in the process of back propagation. Secondly, based on the observation that it is the orientation rather than the magnitude of gradients (gradients of weights) that guides deep neural networks (DNNs) to converge, it’s more important to keep the orientation rather than the magnitude. What’s more, the updates ( $\Delta W$ ) jointly determined by the learning rate and gradients are the data that



ultimately determines the final changes of the network. The much bigger learning rate setting in WAGEUBN (minimum learning rate: WAGEUBN  $1.95 \times 10^{-3}$  ( $2^{-9}$ ) VS floating-point network  $5 \times 10^{-6}$ ) can make up for the problem of update magnitude caused by the smaller quantized gradients. Through the technics above, the problems of saturated learning and dead gradients can be solved. In addition to the above measures which can help the model get rid of local optimum, the proposed quantized Momentum optimizer, fixed-point updates which is precise enough and a proper batch size selection can also reduce the probability of falling into a local optimum greatly.

### B. Training Curve

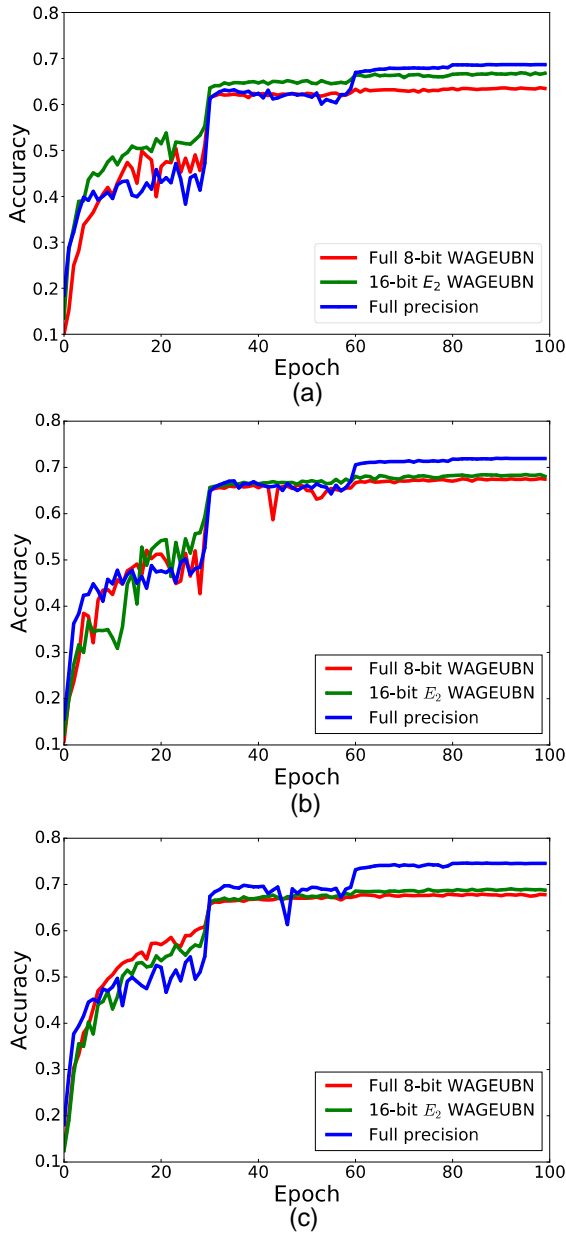


Fig. 6: Training curves under the WAGEUBN framework: (a) ResNet18; (b) ResNet34; (c) ResNet50.

Figure 6 illustrates the accuracy comparison between vanilla DNNs (FP32), DNNs with the full 8-bit version of WAGEUBN, and DNNs with the 16-bit  $E_2$  version of WAGEUBN. The initial learning rate and momentum coefficient under WAGEUBN and full precision are slightly different. To ensure the best recognition accuracy, we use the official parameter settings of TensorFlow [34], where the initial learning rate and momentum coefficient are set to 0.05 (batch size is 128) and 0.9 while those of WAGEUBN framework are set to 0.05078125 ( $26 \times 2^{-9}$ , 10-bit integer) and 0.75 ( $3 \times 2^{-2}$ , 3-bit integer). And this may cause the phenomenon that WAGEUBN converges faster at the beginning of training. During epoch 30 and 60, we have reduced the learning rate, which is a general practice in training process [35, 36]. The training curves show that there is little difference between vanilla DNNs and the ones under the WAGEUBN framework when the training epoch is less than 60, which reflects the effectiveness of our approach. As the epoch evolves, the accuracy gap begins to grow because the learning rate in vanilla DNNs is much lower than that in WAGEUBN, such as  $5 \times 10^{-6}$  v.s.  $1.95 \times 10^{-3}$  ( $k_{lr} = 10$ ), thus the update of vanilla DNNs is more precise than that under the WAGEUBN framework. We can further improve the accuracy by reducing the learning rate, while the bit width values of learning rate and update need to increase accordingly at the expense of more overheads.

Table I quantitatively presents the accuracy comparison between vanilla DNNs and WAGEUBN DNNs on the ImageNet dataset. We have achieved the state-of-the-art accuracy on large-scale DNNs with full 8-bit INT quantization. The 16-bit  $E_2$  WAGEUBN only loses 3.46% mean accuracy compared with the vanilla DNNs. Because the bit width of most data keeps the same between the full 8-bit WAGEUBN and the 16-bit  $E_2$  WAGEUBN, the overhead difference between them is negligible. The DNNs under WAGEUBN framework have achieved an accuracy that is comparable to FP8 [12] and QBP2 [13]. And although MP [14] and MP-INT [16] can achieve the accuracy close to the full precision networks, the computational cost is much higher than WAGEUBN because of the floating-point data type and higher bit width, which will be detailed in Section IV-F. Compared with the vanilla DNNs, about  $4\times$  memory size shrink, much faster processing speed, and much less energy and circuit area can be achieved under the proposed WAGEUBN framework.

### C. Quantization Strategies for W, A, G, E, and BN

In our WAGEUBN framework, we use different quantization strategies for W, A, G, E, and BN, i.e.  $Q(\cdot)$  for W, A, and BN;  $CQ(\cdot)$  for G, and  $SQ(\cdot)$  for E. Different quantization strategies are based on the data distribution, data sensitivity, and hardware friendliness. Figure 7 shows the distribution comparison between W, BN ( $x_2^l$  defined in Equation (2)), A, G (weight gradient), and E ( $e_0^l$ ,  $e_3^l$  defined in Equation (3)) before and after quantization.

According to the definition, the resolution of the direct-quantization function is  $2^{-7}$  when the bit width equals 8 and there is no limitation on the data range. Because W, BN, and

TABLE I: Accuracy of vanilla DNNs and WAGEUBN DNNs on ImageNet dataset.

Network	$k_W$	$k_A$	$k_{G_W}$	$k_{E_1}$	$k_{E_2}$	$k_{WU}$	Accuracy Top-1/Top-5(%)
ResNet18	32	32	32	32	32	32	68.70/88.37
	8	8	8	8	16	24	67.40/87.63
	8	8	8	8	8	24	64.79/85.81
ResNet34	32	32	32	32	32	32	71.99/90.56
	8	8	8	8	16	24	68.50/87.96
	8	8	8	8	8	24	67.63/87.70
ResNet50	32	32	32	32	32	32	74.66/92.13
	8	8	8	8	16	24	69.07/88.45
	8	8	8	8	8	24	67.95/88.01

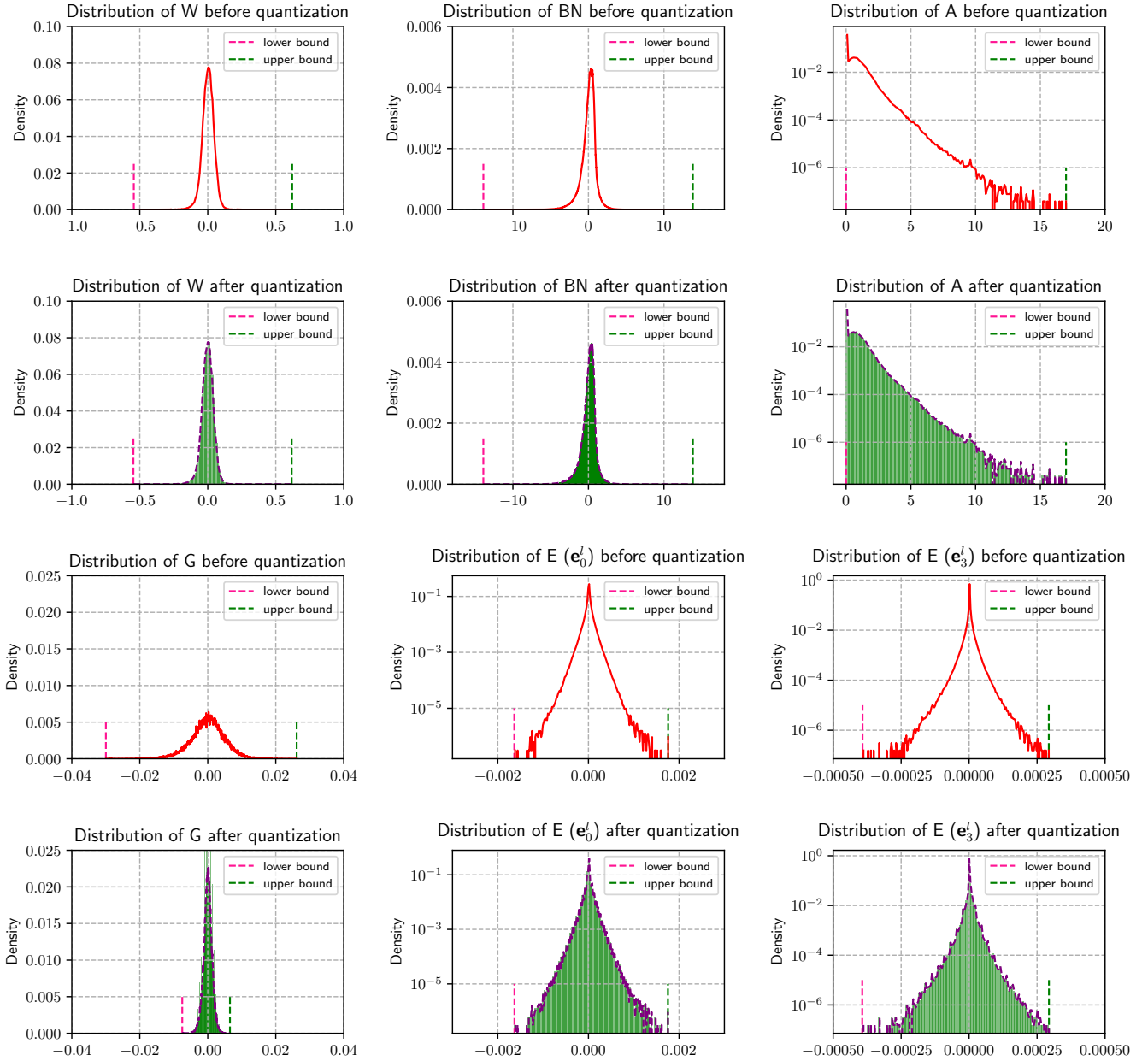


Fig. 7: Data distribution comparison between W, BN, A, G, and E before and after quantization.

A in the inference stage directly affect the loss function and further influence the backpropagation, the quantization of W, BN, and A should be as precise as possible to avoid the loss fluctuation. This is guaranteed for the reason that the resolution of the direct-quantization function is enough for W, BN, and A, which indicates that the direct-quantization function barely changes their data distributions.

The constant-quantization function has a resolution of  $2^{-14}$  and the data range after quantization is about  $[-2^{-7}, 2^{-7}]$  in the case of  $k_{GC} = 15$  and  $k = 8$ .  $k$  will decrease as the training epoch goes on, causing the data range reduction. Figure 7 reveals that the constant-quantization function changes the data distribution of G greatly while the network accuracy has not declined much as a result. The reason behind this phenomenon is that it is the orientation rather than the magnitude of gradients that guides DNNs to converge. In the meantime, it is easy to ensure that the bit width of updates can be fixed when  $k_{GC}$  is fixed, which is more hardware-friendly since the bit width of weights stored in memory can also be fixed during training.

The shift-quantization retains the magnitude order and omits the general values whose absolute value is less than  $2^{-7}$  when  $k = 8$ . The 8-bit shift-quantization works well for the quantization of error after activation ( $e_1^l$  defined in Equation (3)). However, we find the shift-quantization is not enough for the quantization of errors between Conv and BN ( $e_3^l$  defined in Equation (3)). Therefore, the newly designed quantization function in Equation (17) named 8-bit Flag  $Q_{E_2}$  is utilized. Figure 7 shows that the distribution of E ( $e_3^l$ ) is almost the same before and after quantization, revealing the validity of the 8-bit Flag  $Q_{E_2}$  quantization function.

#### D. Accuracy Sensitivity Analysis

To compare the influences of W, A, G, E, and BN quantization individually, we quantize them to 8-bit INT separately with the FP32 update. Taking  $k_W = 8$  as an example, we quantize only W to 8-bit INT and leaving others (A, BN, G, E, and U) still kept in FP32. The quantization function for single data used here is the same as what Section III-D describes (Equation (17) is used for the error quantization when  $k_{E_2} = 8$ ).

The results of ResNet18 under the WAGEUBN framework with single data quantization is shown in Table II. The accuracy of single data quantization reflects the difficulty degree when quantizing W, A, G, E, and BN, separately. From the table, we can see that the quantization of E, especially  $e_3^l$  defined in Equation (3), makes the most impacts on accuracy. In addition, we find that the accuracy heavily fluctuates during training when  $e_3^l$  is constrained to 8-bit INT, which does not appear in the quantization of other data. To sum up, the E data, especially  $e_3^l$ , demands the highest precision and is the most sensitive component under our WAGEUBN framework.

Because the BN process in the forward pass and the average gradients in the backward pass all involve the batch size, the sensibility of batch size to network performance under WAGEUBN is also explored. We have done comparative experiments and the results are illustrated in Figure 8. When

the batch size changes between 128 and 16, the accuracy of the full precision DNNs won't drop significantly, while that of the full 8-bit WAGEUBN DNNs has a relatively large decline in accuracy when the batch size equals 16. One reason for this lies in that there is a momentum used to update the moving means and standard deviations in a batch for full precision DNNs while WAGEUBN abandons this considering the computational cost. Besides, we think it is normal because quantization will inevitably result in the loss of some precise information. And the loss of information will bring a slightly higher batch size sensitivity. However, the robustness of WAGEUBN is still quite good, because there is a significant decline in accuracy only when the batch size is less than a very small number 16.

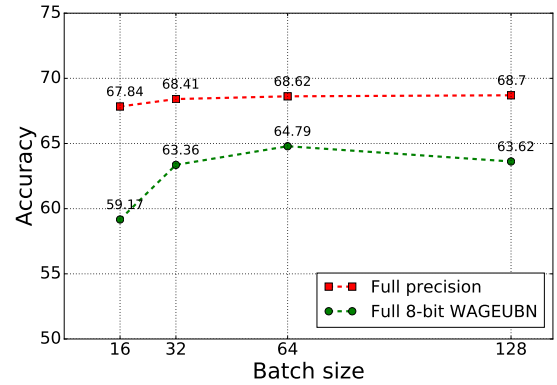


Fig. 8: Accuracy sensibility analysis of batch size under WAGEUBN.

#### E. Analysis of the Error Quantization between Conv and BN

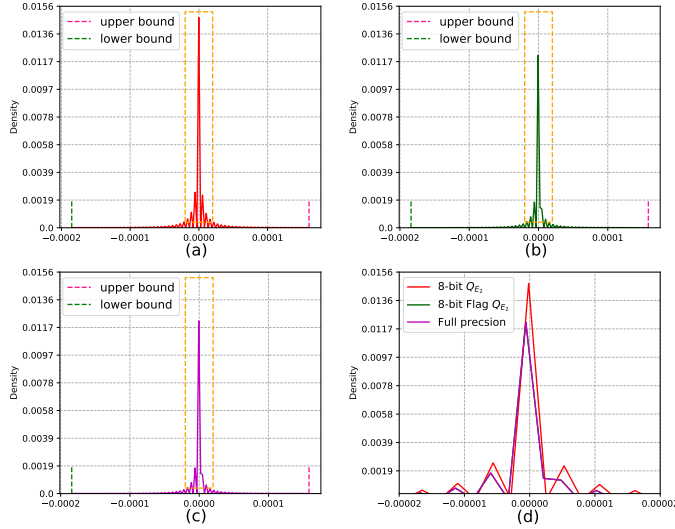
Error backpropagation is the foundation of DNN training. If the error of E caused by quantization is too large, the convergence of DNNs will be degraded. Especially, because the error quantization between convolution and BN ( $e_3^l$ ) is directly related to the weight update of the  $l$ -th layer, the impact of  $e_3^l$  quantization on the model accuracy is critical.

To further analyze the reason why 8-bit  $Q_{E_2}$  (defined in Equation (16) where  $k_{E_2} = 8$ ) causes the non-convergence of DNNs and compare the distributions of  $e_3^l$  under 8-bit  $Q_{E_2}$ , 8-bit Flag  $Q_{E_2}$  (defined in Equation (17) where  $k_{E_2} = 8$ ), and full precision, the data distribution of  $e_3^l$  of the first quantized layer on ResNet18 is shown in Figure 9. From the figure, we can see that the  $e_3^l$  distributions under 8-bit  $Q_{E_2}$  quantization and full precision differs a lot and those of 8-bit Flag  $Q_{E_2}$  quantization and full precision are almost the same. The major difference between the 8-bit  $Q_{E_2}$  quantization and full precision lies in the interval of  $[-2^{-8}R(e_3^1), 2^{-8}R(e_3^1)]$  ( $R(\cdot)$  is defined in Equation (7)), where the 8-bit  $Q_{E_2}$  quantization forces the data in this range to zero.

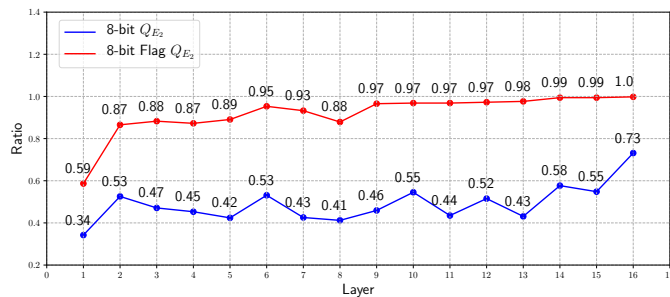
The only difference between 8-bit  $Q_{E_2}$  and 8-bit Flag  $Q_{E_2}$  quantization functions lies in the data range. Theoretically, the covered data range of 8-bit  $Q_{E_2}$  and 8-bit Flag  $Q_{E_2}$  are about  $[-R(e_3^l), -2^{-8}R(e_3^l)] \cup \{0\} \cup [2^{-8}R(e_3^l), R(e_3^l)]$  and  $[-R(e_3^l), -2^{-15}R(e_3^l)] \cup \{0\} \cup [2^{-15}R(e_3^l), R(e_3^l)]$ ,

TABLE II: Accuracy sensitivity under WAGEUBN with single data quantization on ResNet18.

Bit-width	$k_W = 8$	$k_{BN} = 8$	$k_A = 8$	$k_{GW} = 8$	$k_{E_1} = 8$	$k_{E_2} = 8$
Accuracy Top-1/Top-5 (%)	67.98/88.02	68.01/87.96	67.74/87.89	67.88/87.89	67.88/87.92	67.08/87.44


 Fig. 9: Data distributions of  $e_3^l$ : (a) 8-bit  $Q_{E_2}$ , (b) 8-bit Flag  $Q_{E_2}$ , (c) full precision; (d) Data distribution comparison.

respectively. Because the distribution of  $e_3^l$  is not uniform, the range covered by different quantization methods varies a lot. The data ratios (Proportion of non-zero values after quantization.) of 8-bit  $Q_{E_2}$  and 8-bit Flag  $Q_{E_2}$  quantization functions covered by each layer of ResNet18 are illustrated as Figure 10. Although the larger values take greater impacts on the model accuracy in the process of error propagation, the smaller values also contain useful information and occupy the majority. Compared with 8-bit Flag  $Q_{E_2}$ , the data ratio 8-bit  $Q_{E_2}$  covers is too little because of the smaller data range. That is to say, although the most important information contained by the larger values is retained, the information contained by the smaller values is ignored, resulting in the non-convergence of DNNs. In addition, there is also a rough trend that the data ratio decreases as the network becomes shallower, either in the 8-bit  $Q_{E_2}$  or 8-bit Flag  $Q_{E_2}$  quantization.


 Fig. 10: Data ratios of 8-bit  $Q_{E_2}$  and 8-bit Flag  $Q_{E_2}$  quantization methods covered by each layer of ResNet18.

### F. Cost Discussion

Although it is recognized that DNN quantization can greatly reduce memory and compute costs, resulting in lower energy consumption, quantitative analysis is rarely seen in recent research. In order to compare the full INT8 quantization with other precision solutions (FP32, INT32, FP16, INT16, and FP8) more clearly, we have simulated the processing speed, power consumption, and circuit area for single multiplication and accumulation operation on FPGA platform. Figure 11 shows the results. With FP32 as the baseline, taking the multiplication operation as an example, INT8 can perform  $>3\times$  faster in speed,  $10\times$  lower in power, and  $9\times$  smaller in circuit area. Similarly, compared with FP32, the speed of INT8 accumulation is about  $9\times$  faster, and the energy consumption and circuit area are reduced by  $>30\times$ . In addition, the INT8 multiplication and accumulation operations are more advantageous than other data type operations, whether it is FP8[12], INT16[16], FP16[14] or INT32. In conclusion, the proposed full INT8 quantization has great advantages in hardware overheads, whether in terms of memory cost, processing speed, power consumption, and circuit area. Given the huge advantages in hardware resources and computing speed of quantization, the idea of low-precision computing and the quantization functions in WAGEUBN can also extend to other research fields which involve large-scale matrix operations, such as large-scale control systems [37], weather forecasting models [38], etc.

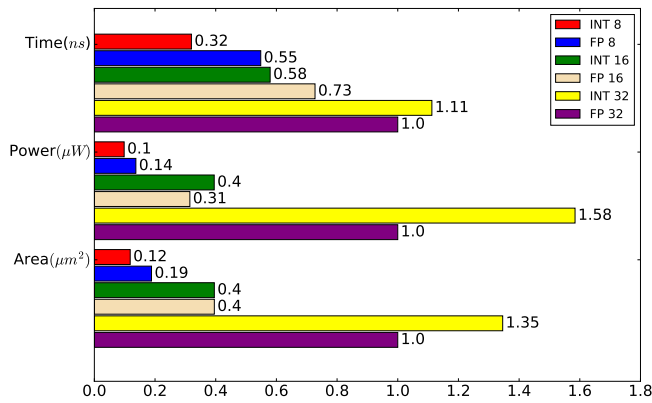
## V. CONCLUSIONS

We propose a unified framework termed as “WAGEUBN” to achieve a complete quantization of large-scale DNNs in both training and inference with competitive accuracy. We are the first to quantize DNNs over all data paths and promote DNN quantization to the full INT8 level. In this way, all the operations can be replaced with bit-wise operations, causing significant improvements in memory overhead, processing speed, circuit area, and energy consumption. Extensive experiments evidence the effectiveness and efficiency of WAGEUBN. This work provides a feasible solution for the online training acceleration of large-scale and high-performance DNNs and further shows the great potential for the applications in future efficient portable devices with online learning ability. Although great progress has been made, WAGEUBN still needs to build a special hardware architecture design and develop more applications. Future works could transfer to the design of computing architecture, memory hierarchy, interconnection infrastructure, and mapping tool to enable the specialized machine learning chips and more applications.

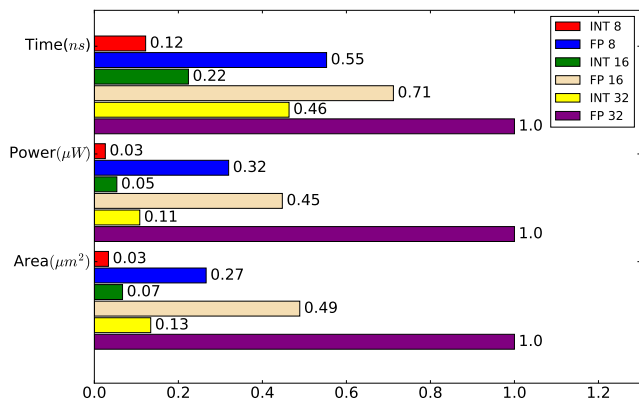
## ACKNOWLEDGMENT

The work was supported partially by the National Natural Science Foundation of China (No. 61603209, 61876215).





(a)



(b)

Fig. 11: Comparison of time, power, and area of single multiplication and accumulation operation under different quantization precision: (a) multiplication, (b) accumulation.

## REFERENCES

- [1] Yann LeCun, Yoshua Bengio, and Geoffrey Hinton. Deep learning. *nature*, 521(7553):436, 2015.
- [2] Junyuan Xie, Linli Xu, and Enhong Chen. Image denoising and inpainting with deep neural networks. In *International Conference on Neural Information Processing Systems*, 2012.
- [3] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster r-cnn: Towards real-time object detection with region proposal networks. In *Advances in neural information processing systems*, pages 91–99, 2015.
- [4] Ronan Collobert and Jason Weston. A unified architecture for natural language processing: Deep neural networks with multitask learning. In *Proceedings of the 25th international conference on Machine learning*, pages 160–167. ACM, 2008.
- [5] Bo Zhang, Luping Shi, and Sen Song. Creating more intelligent robots through brain-inspired computing. *special supplement: Brain-inspired intelligent robotics: The intersection of robotics and neuroscience sciences*, pages 4–9, Science Vol. 354, Issue 6318, pp.1445 (2016).
- [6] Wonyong Sung, Sungho Shin, and Kyuyeon Hwang. Resiliency of deep neural networks under quantization. *arXiv preprint arXiv:1511.06488*, 2015.
- [7] Matthieu Courbariaux, Yoshua Bengio, and Jean-Pierre David. Binaryconnect: Training deep neural networks with binary weights during propagations. In *Advances in neural information processing systems*, pages 3123–3131, 2015.
- [8] Mohammad Rastegari, Vicente Ordonez, Joseph Redmon, and Ali Farhadi. Xnor-net: Imagenet classification using binary convolutional neural networks. In *European Conference on Computer Vision*, pages 525–542. Springer, 2016.
- [9] Cong Leng, Zesheng Dou, Hao Li, Shenghuo Zhu, and Rong Jin. Extremely low bit neural network: Squeeze the last bit out with admm. 2018.
- [10] 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.
- [11] Lei Deng, Peng Jiao, Jing Pei, Zhenzhi Wu, and Guoqi Li. Gxnor-net: Training deep neural networks with ternary weights and activations without full-precision memory under a unified discretization framework. *Neural Networks*, 100:49–58, 2018.
- [12] Naigang Wang, Jungwook Choi, Daniel Brand, Chia-Yu Chen, and Kailash Gopalakrishnan. Training deep neural networks with 8-bit floating point numbers. In *Advances in neural information processing systems*, pages 7675–7684, 2018.
- [13] Ron Banner, Itay Hubara, Elad Hoffer, and Daniel Soudry. Scalable methods for 8-bit training of neural networks. In *Advances in Neural Information Processing Systems*, pages 5145–5153, 2018.
- [14] Paulius Micikevicius, Sharan Narang, Jonah Alben, Gregory Diamos, Erich Elsen, David Garcia, Boris Ginsburg, Michael Houston, Oleksii Kuchaiev, Ganesh Venkatesh, et al. Mixed precision training. *arXiv preprint arXiv:1710.03740*, 2017.
- [15] Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. *arXiv preprint arXiv:1502.03167*, 2015.
- [16] Dipankar Das, Naveen Mellempudi, Dheevatsa Mudigere, Dhiraj Kalamkar, Sasikanth Avancha, Kunal Banerjee, Srinivas Sridharan, Karthik Vaidyanathan, Bharat Kaul, Evangelos Georganas, et al. Mixed precision training of convolutional neural networks using integer operations. *arXiv preprint arXiv:1802.00930*, 2018.
- [17] Charbel Sakr and Naresh Shanbhag. Per-tensor fixed-point quantization of the back-propagation algorithm. *arXiv preprint arXiv:1812.11732*, 2018.
- [18] Shuang Wu, Guoqi Li, Feng Chen, and Luping Shi. Training and inference with integers in deep neural networks. *arXiv preprint arXiv:1802.04680*, 2018.
- [19] Yann LeCun, Léon Bottou, Yoshua Bengio, Patrick Haffner, et al. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.

- [20] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014.
- [21] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. In *Advances in neural information processing systems*, pages 1097–1105, 2012.
- [22] Ning Qian. On the momentum term in gradient descent learning algorithms. *Neural networks*, 12(1):145–151, 1999.
- [23] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- [24] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pages 248–255. Ieee, 2009.
- [25] Matthieu Courbariaux, Itay Hubara, Daniel Soudry, Ran El-Yaniv, and Yoshua Bengio. Binarized neural networks: Training deep neural networks with weights and activations constrained to+ 1 or-1. *arXiv preprint arXiv:1602.02830*, 2016.
- [26] Geoffrey Hinton, Nitsh Srivastava, and Kevin Swersky. Neural networks for machine learning. *Coursera, video lectures*, 264, 2012.
- [27] Yoshua Bengio, Nicholas Léonard, and Aaron Courville. Estimating or propagating gradients through stochastic neurons for conditional computation. *arXiv preprint arXiv:1308.3432*, 2013.
- [28] Itay Hubara, Matthieu Courbariaux, Daniel Soudry, Ran El-Yaniv, and Yoshua Bengio. Quantized neural networks: Training neural networks with low precision weights and activations. *The Journal of Machine Learning Research*, 18(1):6869–6898, 2017.
- [29] Shu-Chang Zhou, Yu-Zhi Wang, He Wen, Qin-Yao He, and Yu-Heng Zou. Balanced quantization: An effective and efficient approach to quantized neural networks. *Journal of Computer Science and Technology*, 32(4):667–682, 2017.
- [30] Qinyao He, He Wen, Shuchang Zhou, Yuxin Wu, Cong Yao, Xinyu Zhou, and Yuheng Zou. Effective quantization methods for recurrent neural networks. *arXiv preprint arXiv:1611.10176*, 2016.
- [31] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Delving deep into rectifiers: Surpassing human-level performance on imagenet classification. In *Proceedings of the IEEE international conference on computer vision*, pages 1026–1034, 2015.
- [32] Diwen Wan, Fumin Shen, Li Liu, Fan Zhu, Jie Qin, Ling Shao, and Heng Tao Shen. Tbn: Convolutional neural network with ternary inputs and binary weights. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 315–332, 2018.
- [33] Yoojin Choi, Mostafa El-Khamy, and Jungwon Lee. Learning low precision deep neural networks through regularization. *arXiv preprint arXiv:1809.00095*, 2018.
- [34] Martín Abadi, Ashish Agarwal, Paul Barham, Eugene Brevdo, Zhifeng Chen, Craig Citro, Greg S Corrado, Andy Davis, Jeffrey Dean, Matthieu Devin, et al. Tensorflow: Large-scale machine learning on heterogeneous distributed systems. *arXiv preprint arXiv:1603.04467*, 2016.
- [35] J Brownlee. Using learning rate schedules for deep learning models in python with keras, 2016.
- [36] Suki Lau. Learning rate schedules and adaptive learning rate methods for deep learning. *Towards Data Science*, 2017.
- [37] V Tharanidharan, R Sakthivel, B Kaviarasan, Faris Alzahrani, and S Marshal Anthoni. Finite-time boundedness of large-scale systems with actuator faults and gain fluctuations. *International Journal of Robust and Nonlinear Control*, 2019.
- [38] Filip Váňa, Peter Düben, Simon Lang, Tim Palmer, Martin Leutbecher, Deborah Salmond, and Glenn Carver. Single precision in weather forecasting models: An evaluation with the ifs. *Monthly Weather Review*, 145(2):495–502, 2017.