

# VecQ: Minimal Loss DNN Model Compression With Vectorized Weight Quantization

Cheng Gong<sup>ID</sup>, Yao Chen, Ye Lu<sup>ID</sup>, Tao Li<sup>ID</sup>, Cong Hao, and Deming Chen<sup>ID</sup>, *Fellow, IEEE*

**Abstract**—Quantization has been proven to be an effective method for reducing the computing and/or storage cost of DNNs. However, the trade-off between the quantization bitwidth and final accuracy is complex and non-convex, which makes it difficult to be optimized directly. Minimizing direct quantization loss (DQL) of the coefficient data is an effective local optimization method, but previous works often neglect the accurate control of the DQL, resulting in a higher loss of the final DNN model accuracy. In this paper, we propose a novel metric, called Vector Loss. Using this new metric, we decompose the minimization of the DQL to two independent optimization processes, which significantly outperform the traditional iterative L2 loss minimization process in terms of effectiveness, quantization loss as well as final DNN accuracy. We also develop a new DNN quantization solution called VecQ, which provides minimal direct quantization loss and achieve higher model accuracy. In order to speed up the proposed quantization process during model training, we accelerate the quantization process with a parameterized probability estimation method and template-based derivation calculation. We evaluate our proposed algorithm on MNIST, CIFAR, ImageNet, IMDB movie review and THUCNews text data sets with numerical DNN models. The results demonstrate that our proposed quantization solution is more accurate and effective than the state-of-the-art approaches yet with more flexible bitwidth support. Moreover, the evaluation of our quantized models on Salient Object Detection (SOD) tasks maintains comparable feature extraction quality with up to  $16\times$  weight size reduction.

**Index Terms**—DNN compression, DNN quantization, vectorized weight quantization, low bitwidth, vector loss

## 1 INTRODUCTION

DEEP Neural Networks (DNNs) have been widely adopted in machine learning based applications [1], [2]. However, besides DNN training, DNN inference is also a computation-intensive task which affects the effectiveness of DNN based solutions [3], [4], [5]. Neural network quantization employs low precision and low bitwidth data instead of high precision data for the model execution. Compared to the DNNs with floating point with 32-bit width (FP32), the quantized model can achieve up to  $32\times$  compression rate with an extremely low-bitwidth quantization [6]. The low-bitwidth processing, which reduces the cost of the inference by using less memory and reducing the complexity of the

multiply-accumulate operation, improves the efficiency of the execution of the model significantly [5], [7].

However, lowering the bitwidth of the data often brings accuracy degradation [4], [8], [9]. This requires the quantization solution to balance between computing efficiency and final model accuracy. However, the quantitative trade-off is non-convex and hard to optimize – *the impact of the quantization to the final accuracy of the DNN models is hard to formulate*.

Previous methods neglect the quantitative analysis of the Direct Quantization Loss (DQL) of the weight data and make the quantization decision empirically while directly evaluating the final model accuracy [6], [10], [11], [12], [13] thus only achieving unpredictable accuracy.

In order to achieve higher training accuracy, finding an optimal quantization solution with minimal loss during the training of the learning kernels is effective and practical. One way of finding a local optimal solution is to minimize the DQL of the weight data, which is widely used in the current quantization solutions [14], [15], [16], [17], [18].

As shown in Fig. 1,  $w_f$  denotes the full-precision weight and  $w_q$  is the value after quantization. Conventional quantization methods regard  $w_f$  as a point (set as origin in Fig. 1) in Euclidean Space, and  $w_q$  is a point which is close to  $w_f$  in a discrete data space. The discrete data space contains a certain number of data points that can be represented by the selected bitwidth. Therefore, the Square of Euclidean Distance (Square 2-norm or called L2 distance [19]), represented as  $\|w_f - w_q\|_2^2$ , between the original weight data and the quantized data is simply used as the loss of the quantization process, which is going to be reduced [15], [16], [17], [18].

Although the L2 based solutions are proven to be effective and provide good training results in terms of accuracy

- Cheng Gong is with the College of Computer Science, Nankai University, Tianjin 300071, China, and also with the Tianjin Key Laboratory of Network and Data Security Technology, Tianjin 300350, China. E-mail: cheng-gong@mail.nankai.edu.cn.
- Yao Chen is with the Advanced Digital Sciences Center, Singapore 138602. E-mail: yao.chen@adsc-create.edu.sg.
- Ye Lu is with the College of Cyber Science, Nankai University, Tianjin 300071, China. E-mail: luyue@nankai.edu.cn.
- Tao Li is with the College of Computer Science, Nankai University, Tianjin 300071, China, and also with the Tianjin Key Laboratory of Network and Data Security Technology, Tianjin 300350, China. E-mail: litao@nankai.edu.cn.
- Cong Hao is with the Electrical and Computer Engineering, the Grainger College of Engineering, University of Illinois at Urbana-Champaign, Champaign, IL 61820 USA. E-mail: congh@illinois.edu.
- Deming Chen is with the Electrical and Computer Engineering, the Grainger College of Engineering, University of Illinois at Urbana-Champaign, Champaign, IL 61820 USA, and also with the Advanced Digital Sciences Center, Singapore 138602. E-mail: dchen@illinois.edu.

Manuscript received 15 Oct. 2019; revised 8 Mar. 2020; accepted 1 Apr. 2020. Date of publication 18 May 2020; date of current version 7 Apr. 2021. (Corresponding authors: Ye Lu, Tao Li and Deming Chen.) Recommended for acceptance by Xuehai Qian and Yanzhi Wang. Digital Object Identifier no. 10.1109/TC.2020.2995593

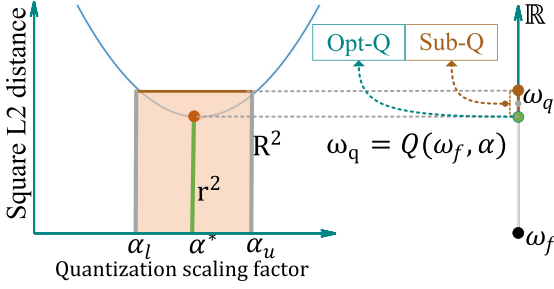


Fig. 1. Uncertainty of using L2 to evaluate the quantization loss.  $\alpha^*$  is the optimal scaling factor for the quantization in the range of  $[\alpha_l, \alpha_u]$ . The optimal distance and achievable distance are denoted as  $r^2$  and  $R^2$ .

of the model and bitwidth of the weight data, such solutions still have some major issues. (1) Traditional L2 based optimizations generally rely on iterative search methods and can get stuck at local optima. As shown in Fig. 1, the quantized results usually fall into the sub-optimal space (Sub-Q) instead of the optimal value (Opt-Q). Even with an additional quantization scaling factor  $\alpha$ , which could help to reduce the differences between the original and quantized data, traditional methods still can not avoid considerable accuracy loss during the quantization process. (2) The process of L2 based quantization focuses on each of the individual weight data and neglects the distribution and correlations among these data points in a kernel or a layer.

To address the issues above, instead of directly minimizing L2 loss, we propose a more accurate quantization loss evaluation metric called *Vector Loss*. Minimizing vector loss is much efficient than directly and iteratively minimizing L2, and results in higher final DNN accuracy; we also propose an algorithm to guide the quantization of the weight data effectively. We construct the weights into a vector  $w_f$  rather than scalar data, to take advantage of the characteristic that the loss between vectors can be decomposed into orientation loss and modulus loss, which are independent of each other. As a result, we are able to achieve the minimal loss of the weight quantization for DNN training.

In this paper, we will demonstrate that using vectorization loss as an optimization objective is better than directly optimizing the L2 distance of the weights before and after quantization. Based on our proposed vectorized quantization loss measurement, we further propose a **Vectorized Quantization** method (**VecQ**) to better explore the trade-off between computing efficiency and the accuracy loss of quantization.

In summary, our contributions are as follows:

- We propose a new metric, Vector Loss, as the loss function for DNN weight quantization, which can provide effective quantization solution compared to traditional methods.
- A new quantization training flow based on the vectorized quantization process is proposed, named VecQ, which achieves better model accuracy for different bitwidth quantization target.
- Parametric estimation and computing template are proposed to reduce the cost of probability density estimation and derivative calculation of VecQ to speed up the quantization process in model training.
- Extensive experiments show that VecQ achieves a lower accuracy degradation under the same training

settings when compared to the state-of-the-art quantization methods in the image classification task with the same DNN models. The evaluations on Saliency Object Detection (SOD) task also show that our VecQ maintains comparable feature extraction quality with up to  $16\times$  weight size reduction.

This paper is structured as follows. Section 2 introduces the related works. In Section 3, the theoretical analysis of the effectiveness of vector loss compared to L2 loss is presented. Section 4 presents the detailed approach of VecQ. Section 5 proposes the fast solution for our VecQ quantization as well as the integration of VecQ into the DNN training flow. Section 6 presents the experimental evaluations and Section 7 concludes the paper.

## 2 RELATED WORKS AND MOTIVATION

As an effective way to compress DNNs, many quantization methods have been explored [6], [8], [9], [10], [11], [12], [13], [14], [15], [16], [17], [18], [19], [20], [21], [22], [23], [24], [25], [26], [27], [28]. These quantization methods can be roughly categorized into 3 different types based on their objective functions for the quantization process:

- Methods based on heuristic guidance of the quantization, e.g., directly minimizing the final accuracy loss;
- Methods based on minimizing Euclidean Distance of weight data before and after quantization;
- Other methods such as training with discrete weights and teacher-student network.

In this section, we first introduce the existing related works based on their different categories and then present our motivation for vectorized quantization.

### 2.1 Heuristic Guidance

The heuristic methods usually directly evaluate the impact of the quantization on the final output accuracy. They often empirically iterate the training process to improve the final accuracy. For example, the BNNs [6] proposed a binary network for fast network inference. It quantizes all the weights and activations in a network to 2 values,  $\{-1, +1\}$ , based on the sign of the data. Although it provides a DNN with 1-bit weights and activations, it is hard to converge without Batch Normalization layers [29] and leads to a significant accuracy degradation when compared to full-precision networks. The Binary Connect [10] and Ternary Connect [11] sample the original weights into binary or ternary according to a sampling probability defined by the value of the weights (after scaling to  $[0,1]$ ). All these works do not quantify the loss during the quantization, so that only the final accuracy is the guideline of the quantization.

Quantization methods in [9], [12] convert the full-precision weights to fixed-point representation by dropping the least significant bits without quantifying the impact.

INQ [23] iteratively processes weight partition, quantization and re-training method until all the weights are quantized into powers-of-two or zeros.

STC [13] introduces a ternary quantization which first scales the weights into the range of  $[-1, 1]$ , and then quantizes all scaled weights into ternary by uniformly partitioning them. Thus, the values located in  $[-1, -1/3]$  and  $[1/3, 1]$  are quantized to -1, 1 and the rest of them are set to 0.

TTQ [19] introduces a ternary quantization which quantizes full-precision weights to ternary by a heuristic threshold but with two different scaling factors for positive and negative values, respectively. The scaling factors are optimized during the back propagation.

The quantization method in [26] (denoted as QAT) employs the *affine mapping* of integers to real values with two constant parameters: Scale and Zero-point. It first subtracts the Zero-point parameter from data (weights/activation), then divides the data by a scaling factor and obtains the quantized results with rounding operation and affine mapping. The approach of TQT [27] follows QAT but with the improvement of constraining the scale-factors into power-of-2 and relates them to trainable thresholds.

## 2.2 Optimizing Euclidean Distance

In order to provide better accuracy control, reducing the Euclidean Distance of the data before and after quantization becomes a popular solution.

Xnor-Net [20] adds a scaling factor on the basis of BNs [6] and calculates the optimal scaling factor to minimize the distance of the weights before and after quantization. The scaling factor boosts the convergence of the model and improves the final accuracy. The following residual quantization method in [24] adopts Xnor-Net [20] to further compensate the errors produced by single binary quantization to improve the accuracy of the quantized model.

TWN [15] proposes an additional threshold factor together with the scaling factor for ternary quantization. The optimal parameters (scaling factor and threshold factor) are still based on the optimization of the Euclidean Distance of weights before and after quantization. TWN achieves better final accuracy than Xnor-Net and BNs.

Extremely low bit method (ENN) proposed in [16] quantizes the weights into the exponential values of 2 by iteratively optimizing the L2 distance of the weights before and after quantization.

TSQ [17] presents a two-step quantization method, which first quantizes the activation to low-bit values, and then fixes it and quantizes the weights into ternary. TSQ employs scaling factor for each of the kernel, resulting in a limited model size reduction.

$\mu$ L2Q [18] first shifts the weights of a layer to a standard normal distribution with a shifting parameter and then employs a linear quantization for the data. The uniform parameter considers the distribution of the weight data, which provides better loss control compared to simply optimizing the Euclidean Distance during the quantization.

Several other works [14], [28] adopt k-means with irregular non-linear quantization. Although the values are clustered before quantization, the final results are still obtained with the optimization of the Euclidean Distance between the original values and the quantized ones.

## 2.3 Other Works

Besides the heuristic and Euclidean Distance approaches, there are still many other works focusing on low-precision DNN training.

GXNOR-Net [25] utilizes the discrete weights during training instead of the full-precision weights. It regards the

discrete values as states and projects the gradients in backward propagation as the transition of the probabilities to update the weights directly, hence, providing a network with ternary weights.

T-DLA [8] quantizes the scaling factor of ternary weights and full-precision activation into fixed-point numbers and constrains the quantization loss of activation values by adopting infinite norms. Compared with [9], [12], it shifts the available bitwidth to the most effective data portion to make full use of the targeted bitwidth.

In TNN [21], the authors design a method using ternary student network, which has the same network architecture as the full-precision teacher network, aiming to predict the output of the teacher network without training on the original datasets.

In HAQ [22], the authors proposed a range parameter — all weights out of the range are truncated and the weights within the range are linearly mapped to discrete values. The optimal range parameter was obtained by solving the KL-divergence of the weights during the quantization.

However, comparing to the heuristic guidance and Euclidean Distance based methods, the approaches above either focus on a specific bitwidth or perform worse in terms of the accuracy of the trained DNNs.

## 2.4 Motivation of the VecQ Method

For the sake of simplicity, in this paper, we will use L2 distance to represent the squared L2 distance. We have witnessed the effectiveness of the L2 distance-based methods among all the existing approaches. However, as explained in the introduction, there are still two defects that lead to inaccurate DQL measurement.

The first defect is that the traditional way of using L2 distance as the loss function for optimization usually cannot be solved accurately and efficiently [14], [15], [17], [18], [28], even with an additional scaling factor  $\alpha$  to scale the data into proper range. As shown in Fig. 1, the quantization function with the additional scaling factor  $\alpha$  to improve the accuracy [15], [18] is denoted as  $w_q = Q(w_f, \alpha)$ ; the L2 distance curve between  $w_f$  and  $w_q$  with the change of  $\alpha$  is drawn in blue. It has a theoretical optimal solution when  $\alpha = \alpha^*$  with a L2 distance of  $r^2$ , shown as the green dot. However, only the solutions  $\alpha \in [\alpha_l, \alpha_u]$  with the L2 distance ranging in  $[r^2, R^2]$  could be obtained due to the lack of solvable expressions [15], [18], leading to an inefficient and inaccurate quantization result. Additionally, even the methods involving k-means for clustering of the weights still fall into the sub-optimal solution space [14], [28]. Their corresponding approximated quantized weights are located in the **Sub-Q** space colored with orange.

The second defect of traditional methods is they neglect the correlation of the weights within the same kernel or layer, but only focus on the difference between single values. Even with the k-means based solutions, the distribution of the weights in the same layer is ignored in the quantization process. However, the consideration of the distribution of the weight data is proven to be effective for the accuracy control in the existing approaches [8], [18].

We discover that when we represent quantization loss of the weight for a kernel or a layer using vector distance



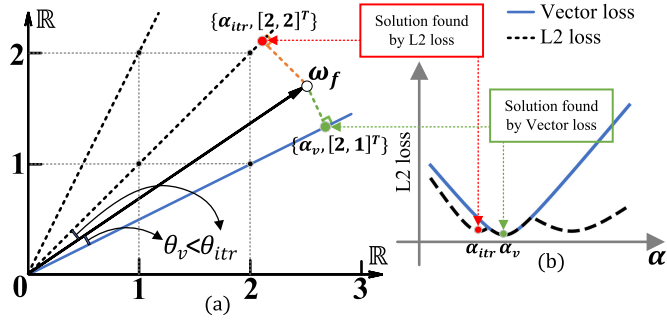


Fig. 2. The solutions based on L2 loss and Vector loss. (a) Solutions of different methods; (b) L2 loss with different methods and  $\alpha$  values.

instead of L2 distance, it will intrinsically solve the two problems mentioned above. We focus on two attributes of a vector, orientation and modulus. Meanwhile, we define a Quantization Angle that represents the intersection angle between the original weight vector and the quantized vector. As a result, the vector distance between the two is naturally determined by the Quantization Angle and their modulus. Therefore, in this work, we evaluate DQL with Vector Loss by leveraging the vector distance, which involves both quantization angle and vector modulus. Note the orientation of the vector is related to quantization angle. When the quantization angle is 0, the orientations of the two intersecting vectors are the same. When the angle is not 0, the orientations of the two vectors are different. For the sake of simplicity, in the rest of the paper, we will use orientation loss and modulus loss to represent the Vector Loss. To the best of our knowledge, there is no previous work that leverages the vector loss for DNN weight quantization.

In this work, we demonstrate that the vector loss can provide effective quantization solution and hence achieve a smaller DQL for the weight data quantization during the model training, which helps for achieving a higher model accuracy. Based on this, we propose VecQ, which carries out the quantization process based on vector loss. We also propose a fast parameter estimation method and a computation template to speed up our vectorized quantization process for easier deployment of our solution.

### 3 VECTORIZED DISTANCE VERSUS L2 DISTANCE

Before introducing our vectorized quantization method, we first explain the effectiveness of loss control with vector loss using two data points as an example for simplicity. Assume a DNN layer with only two weights, denoted as  $\{w_{f1}, w_{f2}\}$  whose values are  $\{2.5, 1.75\}$ . The weights will be quantized into  $k$  bits. The quantization loss based on L2 distance is denoted as  $J_{l2}$  and the quantization solution set is expressed in the format of  $\{\alpha, (v_1, v_2)\}$ , where  $\alpha$  is the floating point scaling factor and  $v_1, v_2$  are the discrete values in the quantization set  $\{v_1, v_2\} \in Q = \{-2^{k-1}, -2^{k-1} + 1, \dots, -1, 0, 1, \dots, 2^{k-1} - 1\}$ , then we get

$$J_{l2} = (w_{f1} - \alpha v_1)^2 + (w_{f2} - \alpha v_2)^2. \quad (1)$$

Let  $w_f = [w_{f1}, w_{f2}]^T$  be a vector from the origin point  $(0,0)$ , and its quantized value is  $w_q = \alpha v$  and  $v = [v_1, v_2]^T$ .  $J_{l2}$

could also be represented as the squared modulus of the distance between vector  $w_f$  and  $w_q$ . The  $J_{l2}$  is calculated as:

$$J_{l2} = \|w_f - \alpha \cdot v\|_2^2, \quad (v_i \in Q). \quad (2)$$

As shown in Fig. 2, there are only two dimensions in the solution space, each representing one weight. Each dimension contains  $2^k$  possible values. The values of the possible solutions are located on the black dotted line in the Fig. 2 due to the full precision scaling factor  $\alpha$ . The quantization angle between  $w_f$  and the quantized version is denoted as  $\theta$ .

However, due to the non-convex characteristic of optimizing L2 loss under the  $k$ -bit constraint [16], the result based on the iterative method may be found as the red point in the figure with the solution of  $\{\alpha_{itr}, v_{itr}\} = \{1.0625, [2, 2]^T\}$  and an angle of  $\theta_{itr}$ , which is the first sub-optimal solution point on the curve of the L2 loss vs  $\alpha$  values in Fig. 2 (b), ignoring the solution with lower loss which is the second extreme value on the curve.

Therefore, instead of directly using L2 distance as the quantization loss function, we use Vector Loss, denoted as  $J_v$ , to measure the difference between vectors: the original weight vector  $w_f$  and the quantized weight vector  $w_q$ , to obtain the quantization solution. We define the vector loss  $J_v = J(w_f, w_q)$  as a composition of two independent losses, orientation loss  $J_o$  and modulus loss  $J_m$ , and

$$J_v = J_o + J_m. \quad (3)$$

$J_o$  and  $J_m$  are computed as:

$$\begin{aligned} J_o &= 1 - \cos \theta, \quad \left( \cos \theta = \frac{\alpha v \cdot w_f}{|\alpha v| |w_f|} \right) \\ &= 1 - e_v e_{w_f} \\ &= 1 - \sum_{i=1}^d (e_{v_i} e_{w_{fi}}) \\ J_m &= \|w_f - \alpha v\|_2^2, \end{aligned} \quad (4)$$

where  $e_v$  and  $e_{w_f}$  represent the unit vector for  $v$  and  $w_f$ .  $w_f$  is a weight vector of a layer of a DNN containing  $d$  weights.

Given the definition of  $J_v$ , we now discuss the effectiveness of minimizing  $J_v$ .

First, according to Eq. 4, the optimization of the orientation loss is only defined by the unit vectors  $e_v$  and  $e_{w_f}$ , where  $\alpha$  value is not affected. Second, when the  $e_v$  is determined, the optimization of the modulus loss is to find the optimal value of  $\alpha$ . Therefore, the optimization of  $J_o$  and  $J_m$  can be achieved independently. If both  $\theta$  and  $\alpha$  are continuous, both orientation and modulus loss minimizations are convex and can be solved with the optimal solutions. However, due to the integer constraints coming from the quantized bitwidth, we need to take the floating-point results and find the corresponding integer values. As a result, the actual achievable  $\theta_v$  of our solution can be different from the optimal value for it. Thus, our solution represents an approximated solution of the optimal solution.

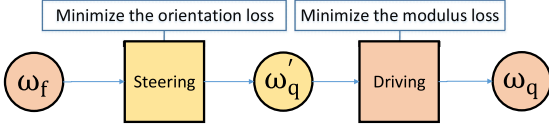


Fig. 3. The overall flow of quantization process, including both steering and driving stage.

Furthermore, denote the quantization result of optimizing  $J_v$  as  $w_q^v$ , we have:

$$\begin{aligned} |w_q^v| &= |w_f| \cos \theta_v \\ |w_q^{tr}| &= |w_f| \cos \theta_{itr}, \end{aligned} \quad (5)$$

Typically, our approximated result can satisfy  $\theta_v \leq \theta_{itr}$ , then we could easily achieve:

$$\begin{aligned} J_{l2} &= |w_f - w_q^v|^2 \\ &= |w_f|^2 + |w_q^v|^2 - 2|w_f||w_q^v| \cos \theta_v \\ &= |w_f|^2 - |w_f|^2 \cos^2 \theta_v \\ &\leq |w_f|^2 - |w_f|^2 \cos^2 \theta_{itr} \\ &= |w_f - w_q^{itr}|^2 \end{aligned} \quad (6)$$

which shows that with vector loss, as long as we optimize the orientation loss to find a better quantization angle, we could achieve an effective solution for weight quantization. Guided by the observation above, we have the final optimal solution for the example points located on the blue solid line in the Fig. 2 as  $\{\alpha_v, v_v\} = \{1.35, [2, 1]^T\}$ , which provides a smaller DQL.

Base on our proposed vector loss metric, we will discuss our algorithms to obtain the vector loss based quantization solution as well as the methods to speed up the algorithms for practical purpose.

## 4 VECTORIZED QUANTIZATION

VecQ is designed to follow the theoretical guideline we developed in Section 3. First of all, the vectorization of the weights is introduced. Then the adoption of the vector loss in VecQ is explained in detail. Finally, the process of VecQ quantization with two critical stages is presented.

### 4.1 Vectorization of Weights

For the weight set  $W_f(l)$  of layer  $l$ , we flatten and reshape them as a  $N \times M \times K^2$ -dimension vector  $w_f(l)$ .  $N$ ,  $M$  indicate the number of input channel and output channel for this layer, and  $K$  indicates the size of the kernel for this layer. For simplicity, we use  $w_f$  to represent the weight vector of a certain layer before quantization.

### 4.2 Loss Function Definition

We use the vectorized loss instead of Euclidean Distance during the quantization process. Since solving the orientation and modulus loss independently could achieve the optimal solution for each of them, we illustrate the quantization loss as is defined in Eq. 3:

$$J(w_f, w_q) = J_o(w_f, w_q) + J_m(w_f, w_q), \quad (7)$$

to provide the constraint during the quantization process.

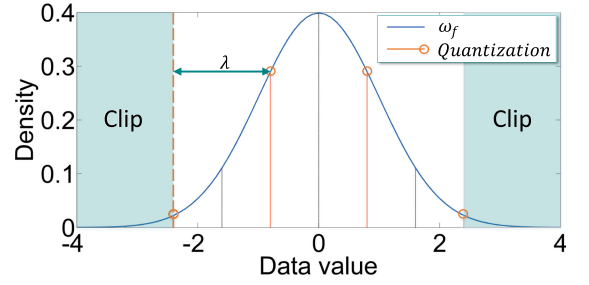


Fig. 4. Linear quantization with interval  $\lambda$ .

### 4.3 Overall Process

According to our analysis in Section 3, the orientation loss  $J_o$  indicates the optimized quantization angle and the modulus loss  $J_m$  indicates the optimized scale at this angle. Therefore, our quantization takes two stages to minimize the two losses independently, which are defined as *steering stage* and *driving stage* as shown in Fig. 3. In the steering stage, we adjust the orientation of the weight vector to minimize the orientation loss. Then, we fix the orientation and only scale the modulus of the vector at the driving stage to minimize the modulus loss.

Let  $w_f \in \mathbb{R}^{N \times M \times K^2}$  be the weight vector of the layer of a DNN in the real space and the  $w_q \in \mathbb{Q}^{N \times M \times K^2}$  be the quantized weight vector in the uniformly discrete subspace. First, steering  $w_f$  to  $w_q'$ :

$$w_q' = \text{steer}(w_f). \quad (8)$$

Where  $w_q'$  is an orientation vector that disregards the modulus of the vector and only focuses on the orientation. Second, along with the determined orientation vector  $w_q'$ , we search the position of the modulus and “drive” to the optimal position with minimum modulus loss. The quantized vector  $w_q$  is achieved by driving the  $w_q'$ .

$$w_q = \text{drive}(w_q'). \quad (9)$$

The complete quantization process is the combination of the two stages. The final target is reducing the loss between the original weight  $w_f$  and the quantized results  $w_q$ . The entire quantization process is represented as

$$w_q = Q(w_f) = \text{drive}(\text{steer}(w_f)). \quad (10)$$

### 4.4 Steering Stage

The purpose of the steering stage is to search for an optimized orientation vector, which has the least orientation loss with  $w_f$  to minimize the  $J_o$ .

As shown in Fig. 4,  $w_f$  is the weight in floating point representation and it would be quantized into  $k$ -bit representation. It means, there are total  $2^k$  values that can be used to represent the values in  $w_f$ , where each of them is denoted as  $q_i, i \in [1, 2^k]$ . We adopt linear quantization method, where an interval  $\lambda = q_i - q_{i-1}$  is used to represent the distance between two quantized values.

$$\text{steer}(w_f, \lambda, k) = \text{Clip}\left\{\left\lfloor \frac{w_f}{\lambda} - 0.5 \right\rfloor, -2^{k-1}, 2^{k-1} - 1\right\} + 0.5. \quad (11)$$

The vector with floating data  $w_f$  is quantized to a vector with discrete data by an rounding ( $\lfloor \cdot \rfloor$ ) operation for each of the values in the vector. The data are limited to the range of  $[-2^{k-1}, 2^{k-1} - 1]$  by extended clip ( $Clip()$ ) operation. The subtraction of 0.5 is used to avoid aggregation at 0 position and guarantees the maximum number of rounding values to be  $2^k$ .

Given a  $k$  for the number of bits as the quantization target, the intermediate quantized weight is

$$w'_q = \text{steer}(w_f, \lambda, k), \quad (12)$$

which has the minimized orientation loss with the  $\lambda$  as an interval parameter. When  $k$  is fixed,  $\lambda$  decides the orientation loss between  $w_f$  and  $w'_q$ . In order to minimize the loss, we only need to find the optimal  $\lambda$ :

$$\lambda^* = \arg \min_{\lambda} (J_o(w_f, w'_q)). \quad (13)$$

Finding the optimal  $\lambda$  requires several processes with high computational complexities; the detailed processes and the corresponding fast solution is presented in Section 5.

#### 4.5 Driving Stage

In the driving stage, we minimize the modulus loss  $J_m$  between the orientation vector  $w'_q$  obtained from the steering stage and the original weight vector  $w_f$ . Since we focus on the modulus in this stage, only the scaling factor  $\alpha$  is involved.

$$\text{drive}(w'_q) = \alpha w'_q. \quad (14)$$

Here we only need to find the optimal  $\alpha$  to minimize the modulus loss.

$$\alpha^* = \arg \min_{\alpha} (J_m(w_f, w_q)), \quad (15)$$

where  $\alpha^*$  can be easily obtained by finding the extreme of  $J_m$  with

$$\frac{dJ_m}{d\alpha} = -2w_q^T (w_f - \alpha w'_q) = 0, \quad (16)$$

and the value of  $\alpha^* = \frac{w_q^T w_f}{w_q^T w'_q}$ .

Finally, with the two stages above, the quantized value of the  $w_f$  is determined by  $\lambda$  and  $\alpha$ :

$$w_q = Q(w_f), \left( \lambda \rightarrow \lambda^*, \alpha^* = \frac{w_q^T w_f}{w_q^T w'_q} \right). \quad (17)$$

### 5 FAST QUANTIZATION

With the proposed quantization method in Section 4, the minimum quantization loss is achieved when the optimal  $\lambda$  in Eq. (13) is found. However, as one of the most critical processes in the model training, the computational overhead of quantization leads to inefficient training of the model. In order to address this issue, in this section, we first analyze the computational complexity to calculate the value of  $\lambda$ . Then, we propose a fast solution based on our proposed fast probability estimation and computation template. In the end, the detailed implementation of our quantization solution together with the fast solver is integrated into our training flow.

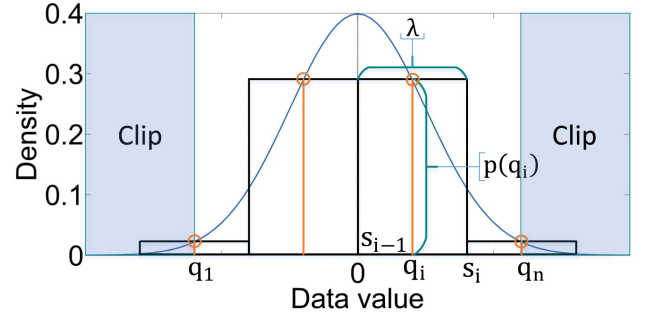


Fig. 5. The illustration of quantization regions and symbols.

#### 5.1 Analysis of the Optimal $\lambda$

The most computational intensive process in our quantization is the steering stage, specifically, the process to solve Eq. (13). However, Eq. (13) can not be solved directly due to the clipping and rounding operations. Instead of directly using the values in  $w_f$ , we involve the distribution of the values in  $w_f$  to support a more general computation of the  $J_o$ . The probability density of the value  $t$  in  $w_f$  is noted as  $p(t)$ .

According to the steering method in Eq. (11), each value  $t \in w_f$  is projected to a value  $q \in w'_q$ ; the  $q$  values are linearly distributed into  $n$  ( $n = 2^k$ ) with a uniform distance defined by interval  $\lambda$ . As shown in Fig. 5, the light blue curve is the distribution of values in  $w_f$  and orange dots are the values after quantization, represented as  $q \in w'_q$ .  $p(q_i)$  indicates the probability of  $q_i$  in  $w_f$ .

Specifically, the data within range  $(s_{i-1}, s_i]$  ( $s_i - s_{i-1} = \lambda, i = 1, \dots, n-1$ ) is replaced by the single value  $q_i$  and the data out of the range  $[q_1, q_n]$  are forced to be truncated and set to the nearest  $q_i$ . As given in Eq. (11),  $q_1 = -2^{k-1}$  and  $q_n = 2^{k-1} - 1$ .

We set  $s_0 = -\infty$  and  $s_n = \infty$  to ease the formulation. Based on the distribution of data in  $w_f$ , the expanding terms of  $J_o(w_f, w'_q)$  in Eq. (13) can be obtained as follows,

$$\begin{cases} |w_f| = \sqrt{\int_{-\infty}^{\infty} t^2 p(t) dt} \\ |w'_q| = \sqrt{\sum_{i=1}^n (\int_{s_{i-1}}^{s_i} q_i^2 p(t) dt)} \\ w_f w'_q = \sum_{i=1}^n (\int_{s_{i-1}}^{s_i} t q_i p(t) dt) \end{cases} \quad (18)$$

The  $J_o(w_f, w'_q)$  is represented as

$$J_o(w_f, w'_q) = 1 - \frac{\sum_{i=1}^n (\int_{s_{i-1}}^{s_i} t q_i p(t) dt)}{\sqrt{\int_{-\infty}^{\infty} t^2 p(t) dt} \cdot \sqrt{\sum_{i=1}^n (\int_{s_{i-1}}^{s_i} q_i^2 p(t) dt)}}. \quad (19)$$

Since the linear quantization is adopted with the fixed interval of  $\lambda$ ,  $q_i$  and  $s_i$  can be easily derived by the following equations,

$$q_i = (i - \frac{n}{2} - 0.5)\lambda, (i = 1, \dots, n) \quad (20)$$

$$s_i = \begin{cases} -\infty, & i = 0 \\ (i - \frac{n}{2})\lambda, & i = 1, \dots, n-1 \\ \infty, & i = n \end{cases} \quad (21)$$

So  $J_o(w_f, w'_q)$  is only related to  $\lambda$  and Eq. (13) can be solved by solving

$$\frac{dJ_o(w_f, w'_q)}{d\lambda} = 0. \quad (22)$$

Concluding from the discussion above, for each  $w_f$ , three steps are necessary:

- Estimating probability density of the values in  $w_f$ .
- Solving Eq. (22) and getting the optimal  $\lambda$ .
- Using the optimal  $\lambda$  to obtain the final quantization results.

However, the first two steps (probability density estimation and derivative calculation) are complex and costly in terms of CPU/GPU time and operations, which limit the training speed.

## 5.2 Fast Solver

For the computation-intensive probability density estimation and derivative calculation, we propose two methods to speed up the processes, which are fast parametric estimation and computing template, respectively.

### 5.2.1 Fast Probability Estimation

There are two methods for probability density estimation: parametric estimation and non-parametric estimation.

Non-parametric estimation is usually used for fitting the distribution of data without prior knowledge. It requires all the density and probability of the data to be estimated individually, which will lead to a huge computational overhead.

We take the widely adopted non-parametric estimation method, Kernel Density Estimation (KED) as an example, to illustrate the complexity of non-parametric estimation.

$$p(t) = \frac{1}{nh} \sum_{i=1}^n K\left(\frac{t_i - t}{h}\right). \quad (23)$$

Here  $p(t)$  is the probability density of  $t$ .  $n$  is the number of the samples and  $K(x)$  is the non-negative kernel function that satisfies  $K(x) \geq 0$  and  $\int K(x)dx = 1$ .  $h$  is the smoothing parameter. The time complexity of computing all probability densities  $p(x)$  is  $O(n^2)$  and the space complexity is  $O(n)$  because all the probability densities need to be computed and stored.

Parametric estimation is used for the data which have a known distribution and only computes some parameters of the distributions instead. It could be processed fast with the proper assumption of the distribution. Thus, we adopt a parametric estimation method in our solution.

There is prior knowledge of the weights of the layers of DNNs, which assumes that they are obeying normal distribution with the mean value of 0 so that the training could be conducted easily and the model could provide better generalization ability [17], [18]:

$$w_f \sim \mathcal{N}(0, \sigma^2), \quad (24)$$

$\sigma$  is the standard derivation of  $w_f$ . Based on this prior knowledge, we can use parametric estimation to estimate the probability density of  $w_f$  and the only parameter that is needed to be computed during the training is  $\sigma$ . The

effectiveness of this prior distribution is also proven by the final accuracy we obtain in the evaluations.

With Eq. (24), we have

$$p(t) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left(-\frac{t^2}{2\sigma^2}\right). \quad (25)$$

Therefore, parametric estimation only requires the standard deviation  $\sigma$ , which could be calculated with

$$\sigma^2 = E(w_f^2) - E(w_f)^2. \quad (26)$$

Here  $E(\cdot)$  computes the expectation. Hence, the time complexity of computing  $\sigma$  is reduced to  $O(n)$  and the space complexity is reduced to  $O(1)$ .

### 5.2.2 Computing Template

After reducing the complexity of computing the probability, finding optimal  $\lambda$  is still complex and time-consuming. In order to improve the computing efficiency of this step, we propose a computing template-based method.

Since the weights of a layer  $w_f$  obey normal distribution,  $\mathcal{N}(0, \sigma^2)$ , they can be transformed to the standard normal distribution:

$$\varphi = \frac{w_f}{\sigma} \sim \mathcal{N}(0, 1), \quad (27)$$

Then, we could use  $\varphi$  to compute  $J_o(\varphi, w'_q)$  instead of using  $w_f$ , because of:

$$\begin{aligned} J_o(\varphi, w'_q) &= 1 - \frac{\varphi w'_q}{|\varphi||w'_q|} \\ &= 1 - \frac{(w_f/\sigma)w'_q}{|w_f/\sigma||w'_q|} \\ &= J_o(w_f, w'_q). \end{aligned} \quad (28)$$

Here,  $\varphi$  is the computing template for  $w_f$ , because it has the same orientation loss with  $w'_q$  as  $w_f$ . By choosing this computing template, solving Eq. (22) is equivalent to solve the substitute equation Eq. (29).

$$\lambda^* = \arg \min_{\lambda} (J_o(\varphi, w'_q)). \quad (29)$$

Since  $\varphi \sim \mathcal{N}(0, 1)$ , the probability of value  $t$  is:

$$p(t) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{t^2}{2}\right). \quad (30)$$

After the probability  $p(t)$  is obtained, the orientation loss function  $J_o(\varphi, w'_q)$  can be expressed as a function only relating to  $\lambda$  and the targeting bitwidth  $k$  for the quantization.

$$J_o(\varphi, w'_q) \sim J_o(\lambda, k), \quad (31)$$

$J_o(\varphi, w'_q)$  is a convex function with the condition of  $k > 1$ . However, it is constant when  $k = 1$  because the angle between the weight vector and the vector constructed with the signs of values in the weights is constant. Due to its independence at  $k = 1$ , we set  $\lambda$  to 1 for the convenience of the following process. We plot the curve of  $J_o(\lambda, k)$  in Fig. 6 with the change of  $\lambda$  under different  $k$  bits.



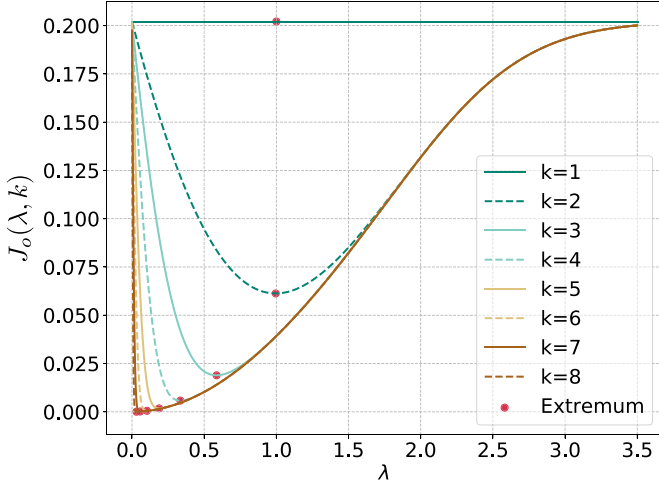


Fig. 6. The orientation loss ( $J_o(\varphi, w'_q) \sim J_o(\lambda, k)$ ) with different  $k$  and  $\lambda$  values. The extrema are marked.

TABLE 1  
Optimal Value of  $\lambda$  for  $J_o(\varphi, w'_q)$  With Bitwidth  $k$

$k$	1	2	3	4	5	6	7	8	$> 8$
$\lambda$	(0, $\infty$ )	0.9957	0.5860	0.3352	0.1881	0.1041	0.0569	0.0308	$6/2^k$

The optimal  $\lambda$  values for all bitwidth settings obtained by solving Eq. (31) is shown in Table 1. The loss is small enough when the targeted bitwidth is greater than 8, so we omit the results for them. With the template above, we only need to solve  $J_o(w_f, w'_q)$  once to find the optimal  $\lambda$ , and then apply it to all quantization without repetitively calculating it. In other words, simply looking up the corresponding value in this table can obtain the optimal parameter thus reducing the complexity and intensity of the computation, which significantly speeds up the training process.

### 5.3 DNN Training Integration

We integrate our VecQ quantization into the DNN training flow for both the weight data and the activation data, as shown in Fig. 7.

**Weight Quantization.** For layer  $l$ , during the forward propagation, we first quantize the weights with full precision ( $w_f(l)$ ) into the quantized values ( $w_q(l)$ ), then use the quantized weights to compute the output ( $z(l)$ ). During the backward propagation, the gradient is calculated with  $w_q(l)$  instead of  $w_f(l)$  and propagated. In the final update process, the gradient  $g(l)$  of  $w_q(l)$  is used to update  $w_f(l)$  [12].

**Activation Quantization.** Inspired by the Batch Normalization (BN) technique, instead of using pre-defined

TABLE 2  
The Image Classification Datasets Attributes

Datasets	MNIST	CIFAR10	ImageNet
Image size	28×28×1	32×32×3	224×224×3
# of Classes	10	10	1000
# of Images	60000	50000	1281167
# of Pixels (log10)	7.67	8.19	11.29

TABLE 3  
The Sentiment and Text Classification Datasets

Datasets	IMDB	THUCNews
Objectives	Movie reviews	Text classification
# of Classes	2	10
# of samples	50000	65000

distribution, we compute the distribution parameter of the activation outputs  $p(t)$  and update it with Exponential Moving Average. During the inference, the distribution parameter is employed as a linear factor to the activation function [29]. The  $A(l)$  is the activation output of layer  $l$ , and  $Activation(\cdot)$  is the non-linear activation function following the convolution or fully-connected layers, such as Sigmoid, Tanh, ReLU.

## 6 EVALUATIONS

We choose Keras v2.2.4 as the baseline DNN training framework [30]. The layers in Keras are rewritten to support our proposed quantization mechanism as presented in Section 5.3. Specifically, all the weight data in the DNNs are quantized to the same bitwidth in the evaluation of VecQ, including the first and last layers. Our evaluations are conducted on two classic tasks: (1) image classification and (2) salient object detection (SOD). The evaluation results for image classification are compared to state-of-the-art results with the same bitwidth configuration and the SOD results are compared to the state-of-the-art solutions that are conducted with the FP32 data type.

### 6.1 Classification

Image classification is the basis of many computer vision tasks, so the classification accuracy of the quantized model is representative for the effectiveness of our proposed solution.

#### 6.1.1 Evaluation Settings

**Datasets and DNN Models.** The MNIST, CIFAR10 [31] and ImageNet [32] datasets are selected for image classification

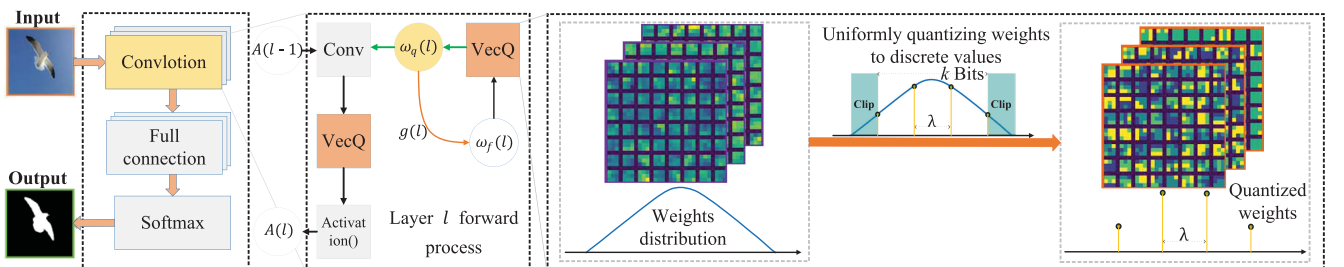


Fig. 7. Integrated quantization process in DNN training.



TABLE 4  
The Models for ImageNet

Models	AlexNet	ResNet-18	MobileNetV2
Convs	5	21	35
DepConvs	-	-	17
BNs	7	19	52
FCs	3	1	1
Parameters (M)	50.88	11.7	3.54

<sup>1</sup>Convs indicate the vanilla convolution layers, DepConvs are the depthwise convolution layers [35]. BNs stand for the Batch Normalization layers [29] and FCs are the full-connection layers.

evaluations; the IMDB movie reviews [33] and THUC-News [34] for Chinese text datasets are selected for the sentiment and text classification evaluations. The detailed information of the datasets are listed in Tables 2 and 3.

For MNIST dataset, Lenet5 with 32C5-BN-MP2-64C5-BN-MP2-512FC-10Softmax is used, where C stands for the Convolutional layer and the number in front denotes the output feature channel number and the number behind is the kernel size; BN stands for the Batch Normalization layer; FC represents the Fully-connected layer and the output channel number is listed in front of it; MP indicates the max pooling layer followed with the size of the pooling kernel. The mini-batch size is 200 samples and the initial learning rate is 0.01 and it is divided by 10 at epoch 35 and epoch 50 for a total of 55 training epochs.

For CIFAR10 dataset, a VGG-like network [1] with the architectural configuration as 64C3-BN-64C3-BN-MP2-128C3-BN-128C3-BN-MP2-256C3-BN-256C3-BN-MP2-1024FC-10Softmax is selected. A simple data augmentation which pads 4 pixels on each side and randomly crops the  $32 \times 32$  patches from the padded image or its horizontal flip is adopted during the training. Only the original  $32 \times 32$  images are evaluated in the test phase. The network is trained with mini-batch size 128 for a total of 300 epoch. The initial learning rate is 0.01 and decays 10 times at epoch 250 and 290.

For the ImageNet dataset, we select 3 famous DNN models, which are AlexNet [36], ResNet-18 [2] and MobileNetV2

[35]. ImageNet dataset contains 1000 categories and the size of the image is relatively bigger [1], [2], [32], [36]. We use the Batch Normalization (BN) layer instead of the original Local Response Normalization (LRN) layer in AlexNet for a fair comparison with [16], [17], [19]. The numbers of the different layers and the parameter sizes are listed in Table 4.

The architecture of the model for IMDB movie review [33] sentiment classification and the THUCNews [34] text classification are 128E-64C5-MP4-70LSTM-1Sigmoid and 128E-128LSTM-128FC-10Softmax, where E denotes Embedding layer and the number in front of it represents its dimension; LSTM is the LSTM layer and the number in front is the number of the hidden units. In addition, we quantize all layers including Embedding layer, Convolutional layer, Fully-connected layer and LSTM layer for these two models. Specifically, for the LSTM layer, we quantize the input features, outputs and the weights, but left the intermediate state and activation of each timestamp untouched since the quantization of them will not help to reduce the size of the model.

*Evaluation Metrics.* The final classification accuracy results on the corresponding datasets are taken as the evaluation metrics in image classification tasks as used in many other works [9], [15], [16], [17]. Moreover, the Top1 and Top5 classification accuracy results are presented simultaneously on all the models for ImageNet dataset for comprehensive evaluation as used in [1], [2].

### 6.1.2 Bitwidth Flexibility

VecQ is designed to support a wide range of targeted bitwidths. We conduct a series of experiments to verify the impact of bitwidth on model accuracy and model size reduction. The bitwidth in the following evaluations ranges from 1 to 8.

The accuracy results and model sizes for the image classification models are shown in Table 5. We first discuss weight quantization only. There is a relatively higher accuracy degradation when the bitwidth is set to 1. But starting from 2 bits and up, the accuracy of the models recovers to less than 1.37 percent drop when compared to the FP32 version except MobileNetV2. With the increase of bitwidth, the accuracy

TABLE 5  
The Accuracy and Model Size With Different Bitwidth Targets

W/A <sup>1</sup>	LeNet5		VGG-like [1]		Alexnet [36]		ResNet-18 [2]		MobileNetV2 [35]	
	Size(M)	Acc	Size(M)	Acc	Size(M)	Top1/Top5	Size(M)	Top1/Top5	Size(M)	Top1/Top5
32/32	6.35	99.40	20.44	93.49	194.10	60.01/81.90 <sup>2</sup>	44.63	69.60/89.24 <sup>3</sup>	13.50	71.30/90.10 <sup>4</sup>
1/32	0.21	99.34	0.67	90.39	6.21	55.06/77.78	1.45	65.58/86.24	0.67	53.78/77.07
2/32	<b>0.41</b>	<b>99.53</b>	1.31	92.94	12.27	59.31/81.01	2.68	68.23/88.10	1.09	64.67/85.24
3/32	0.60	99.48	1.94	93.02	18.33	60.36/82.40	4.24	68.79/88.45	1.50	69.13/88.35
4/32	0.80	99.47	2.58	93.27	24.39	61.21/82.94	5.63	69.80/89.11	1.92	71.89/90.38
5/32	1.00	99.47	3.22	93.37	30.45	61.65/83.19	7.02	69.98/89.15	2.33	71.47/90.15
6/32	1.20	99.49	3.86	93.51	36.51	62.01/83.32	8.42	69.81/88.97	2.74	72.23/90.61
7/32	1.40	99.48	<b>4.49</b>	<b>93.52</b>	42.57	62.09/83.44	9.81	70.17/89.09	<b>3.16</b>	<b>72.33/90.62</b>
8/32	1.60	99.48	5.13	93.50	<b>48.63</b>	<b>62.22/83.54</b>	<b>11.20</b>	<b>70.36/89.20</b>	3.57	72.24/90.66
2/8	0.41	99.43	1.31	92.46	12.27	58.04/80.21	2.68	67.91/88.30	1.09	63.34/84.42
4/8	0.80	99.53	2.58	93.37	24.39	61.22/83.24	5.63	68.41/88.76	1.92	71.40/90.41
8/8	1.60	99.44	5.13	93.55	48.63	61.60/83.66	11.20	69.86/88.90	3.57	72.11/90.68

<sup>1</sup>W/A denotes the quantizing bits of weights and activation, respectively.

<sup>2</sup>Results of AlexNet with Batch Normalization layers are cited from [37].

<sup>3</sup>Results of ResNet18 are cited from [38].

<sup>4</sup>Results are cited from the document of Keras [30].

TABLE 6  
Evaluation Results for LSTM-Based Models

W/A	IMDB		THUCNews	
	Size (M)	Acc.	Size (M)	Acc.
32/32	10.07	84.98 <sup>1</sup>	3.01	94.74
2/32	0.63	85.54	0.19	93.99
4/32	1.26	<b>86.24</b>	0.38	94.47
8/32	2.52	85.53	0.75	94.53
2/8	0.63	85.40	0.19	94.00
4/8	1.26	84.67	0.38	94.09
8/8	2.52	85.72	0.75	94.43

<sup>1</sup>The results of full precision model is from [30].

of the quantized model is improved. The highest accuracy of LeNet-5 and VGG-like are 99.53 and 93.52 percent at 2 bits and 7 bits, respectively, which outperform the accuracy results with FP32. The highest accuracy of AlexNet, ResNet-18 and MobileNetV2 are obtained at 8-bit with 62.22 percent (Top1), 8-bit with 70.36 percent (Top1) and 7-bit with 72.33 percent (Top1), respectively, and all of them outperform the values obtained with FP32. Overall, the best accuracy improvement of the models when compared to the full precision versions for the five models are 0.13 0.03, 2.21, 0.76 and 1.03 percent, at 2, 7, 8, 8 and 7 bits, respectively when activation is maintained as 32 bits. The table also contains the accuracy of the models with 8-bit activation quantization. Although the activation quantization leads to a degradation of the accuracy for most of the models (except VGG-like), the results are in general comparable with the models with FP32 data type.

The accuracy and model size of the sentiment classification and text classification models are shown in Table 6. Our solution easily outperforms the models trained with FP32. Even with the activation quantization, the accuracy results are still well maintained. The results also indicate that adopting appropriate quantization of the weight data improves the ability of generalization of the DNN models. In another word, right quantization achieves higher classification accuracy for the tests.

### 6.1.3 Comparison With State-of-the-Art Results

We collected the state-of-the-art accuracy results of DNNs quantized to different bitwidth with different quantization methods, and compared them to the results from VecQ with the same bitwidth target. The detailed bitwidth support of the comparison methods are listed in Table 7. Note here, when the quantization of the activations are not enabled, the SFB and SFN are not applicable to VecQ.

The comparisons are shown in Table 8. The final accuracy based on VecQ increased by up to 0.62%, 3.80% for LeNet-5 and VGG-like when compared with other quantization methods, respectively. There is also up to 3.26%/2.73% (top1/top5) improvement of the accuracy of AlexNet, 4.84%/2.84% of ResNet-18 and 6.60%/4.00% of MobileNetV2 when compared to the state-of-the-art methods. For all the 3 datasets and 5 models, the quantized models with VecQ achieve higher accuracy than almost all state-of-the-art methods with the same bitwidth. However, when bitwidth target is 1, the quantized models of AlexNet and Resnet-18 based on VecQ

TABLE 7  
Detailed Settings of the Quantization Methods Collected From the Literature

Methods	Weights			Activation			FConv	IFC	LFC
	Bits	SFB	SFN	Bits	SFB	SFN			
ReBNet [24]	1	32	1	3	32	3	-	Y	N
BC [10]	1	-	0	32	-	-	-	-	-
BWN [20]	1	32	1	32	-	-	-	-	-
BPWN [15]	1	32	1	32	-	-	-	N	N
TWN [15]	2	32	1	32	-	-	-	N	N
TC [11]	2	-	0	32	-	-	-	-	-
TNN [21]	2	-	0	2	-	0	-	-	-
TTQ [19]	2	32	2	32	-	-	N	Y	N
INQ [23]	2,3,4,5	-	0	32	-	-	-	-	-
FP [9]	2,4,8	-	0	32	-	-	-	N	N
uL2Q [18]	1,2,4,8	32	1	32	-	-	Y	Y	8
ENN [16]	1,2,3	32	1	32	-	-	-	-	-
TSQ [17]	2	32	c <sup>4</sup>	2	32	2	N	Y	N
DC [14] <sup>2</sup>	2,3,4	32	4,8,16	32	-	-	8	-	1
HAQ [22] <sup>3</sup>	flexible	32	1	32	-	-	8	-	1
QAT [26]	8	32	1	8	32	1	N	-	N
TQT [27]	8	-	1	8	-	1	8	-	8
VecQ	1-8	32	1	32,8	-32	-1	Y	Y	Y

<sup>1</sup>**Weights** and **Activation** denote the quantized data of the model. **Bits** refer to quantized bitwidth of methods; **SFB** is the bitwidth for the scaling-factor; **SFN** is the number of the scaling-factors. **FConv**, **IFC**, and **LFC** represent whether the First Convolution layer, the Internal Fully-Connected layers and the Last Fully-Connected layer are quantized.

<sup>2</sup>The results of DC are from [22].

<sup>3</sup>HAQ is a mix-precision method; results here are from the experiments that only quantize the weight data.

<sup>4</sup>TSQ introduces a floating-point scaling factor for each convolutional kernel, so the SFN equals to the number of kernels.

perform worse due to the reason that we have quantized all the weights into 1 bit, including first and last layers, which are different from the counterparts that are using higher bitwidth for the first or last layers. This also leads to more accuracy degradation at low bitwidth on the lightweight network MobileNetV2. This, however, allows us to provide an even smaller model size. Besides, the solution called BWN [20] and ENN [16] for AlexNet has 61M parameters [20], while ours is 50.88M because eliminating the layer paddings in the intermediate layers leads to less weights for fully-connected layers. When compared to TWN,  $\mu$ L2Q and TSQ, VecQ achieves significantly higher accuracy at the same targeted bitwidth, which also indicates that our vectorized quantization is superior to the L2 loss based solutions.

### 6.1.4 Analysis of $\lambda$ Values

In order to evaluate the accuracy of our theoretical  $\lambda$  in Table 1, we choose the last convolutional layers from different models to calculate the actual  $\lambda$  values. Since  $\lambda$  is the quantization interval, the range of (0,3] of it covers more than 99.74 percent of the layer data, so the actual  $\lambda$  is obtained by exhaustively searching the values in the range of (0,3] with the precision at 0.001. The comparison is shown in Table 9. As we can learn from Table 9, there are differences between the theoretical  $\lambda$ s and the actual values. However, the final results in terms of accuracy in the previous subsection is maintained and not effected by the small bias of  $\lambda$ , which proves the effectiveness of our solution.

## 6.2 Salient Object Detection

Salient object detection aims at standing out the region of the salient object in an image. It is an important evaluation that

TABLE 8  
The Comparison With Other State-of-the-Art Quantization Solutions

Bitwidth	Methods	Datasets&Models				
		MNIST	Cifar10 [31]	ImageNet [32]		
		LeNet5	VGG-like [1]	AlexNet [36]	ResNet18 [2]	MobileNetV2 [35]
32	FP32	99.40	93.49	60.01/81.90	69.60/89.24	71.30/90.10
1	ReBNet [24]	98.25	86.98	41.43/-	-	-
	BC [10]	98.82	-	-	-	-
	BWN [20]	-	-	56.80/79.40	60.80/83.00	-
	BPWN [15]	99.05	-	-	57.50/81.20	-
	ENN [16]	-	-	57.00/79.70	64.80/86.20	-
	$\mu$ L2Q [18]	99.06	89.02	-	66.24/86.00	-
	VecQ	<b>99.34</b>	<b>90.39</b>	55.06/77.78	65.58/86.24	53.78/77.07
	Mean-Imp <sup>3</sup>	0.55	2.39	3.32/-1.77	3.24/2.14	-/-
2	FP [9]	98.90	-	-	-	-
	TWN [15]	99.35	-	57.50/79.80 <sup>1</sup>	61.80/84.20	-
	TC [11]	98.85	-	-	-	-
	TNN [21]	98.33	-	-	-	-
	TTQ [19]	-	-	57.50/79.70	66.60/87.20	-
	INQ [23]	-	-	-	66.02/87.13	-
	ENN	-	-	58.20/80.60	67.00/87.50	-
	TSQ [17]	-	-	58.00/80.50	-	-
	$\mu$ L2Q	99.12	89.50	-	65.60/86.12	-
	DC <sup>4</sup>	-	-	-	-	58.07/81.24
	VecQ	<b>99.53</b>	<b>92.94</b>	<b>59.31/81.01</b>	<b>68.23/88.10</b>	<b>64.67/85.24</b>
	Mean-Imp	0.62	3.44	1.51/0.86	2.83/1.67	6.60/4.00
3	INQ	-	-	-	68.08/88.36	-
	ENN <sup>2</sup>	-	-	60.00/82.20	68.00/88.30	-
	DC	-	-	-	-	68.00/87.96
	VecQ	99.48	93.02	<b>60.36/82.40</b>	<b>68.79/88.45</b>	<b>69.13/88.35</b>
	Mean-Imp	-	-	0.36/0.20	0.75/0.12	1.13/0.39
4	FP	99.10	-	-	-	-
	INQ	-	-	-	68.89/89.01	-
	$\mu$ L2Q	99.12	89.80	-	65.92/86.72	-
	DC	-	-	-	-	71.24/89.93
	VecQ	<b>99.47</b>	<b>93.27</b>	61.21/82.94	<b>69.80/89.11</b>	<b>71.89/90.38</b>
	Mean-Imp	0.36	3.47	-/-	2.39/1.24	0.65/0.45
5	INQ	-	-	57.39/80.46	68.98/89.10	-
	VecQ	99.47	93.37	<b>61.65/83.19</b>	<b>69.98/89.15</b>	71.47/90.15
	Mean-Imp	-	-	3.26/2.73	1.00/0.05	-/-
6	HAQ	-	-	-	-	66.75/87.32
	VecQ	99.49	93.51	62.01/83.32	69.81/88.87	<b>72.23/90.61</b>
	Mean-Imp	-	-	-/-	-/-	5.48/3.29
7	VecQ	99.48	93.52	62.09/83.44	70.17/89.09	72.33/90.62
8	FP	99.10	-	-	-	-
	$\mu$ L2Q	99.16	89.70	-	65.52/86.36	-
	QAT-c <sup>5</sup>	-	-	-	-	71.10/-
	TQT-wt-th	-	-	-	-	71.80/90.60
	VecQ	<b>99.48</b>	<b>93.50</b>	62.22/83.54	<b>70.36/89.20</b>	<b>72.24/90.66</b>
	Mean-Imp	0.35	3.80	-/-	4.84/2.84	0.79/0.06

<sup>1</sup>The results of TWN on AlexNet are from [16].

<sup>2</sup>ENN adopts 2 bits shifting results noted as  $\{-4, +4\}$  [16].

<sup>3</sup>Mean-Imp denotes **Mean** of accuracy **Improvement** results compared to the state-of-the-art methods.

<sup>4</sup>The results are from HAQ [22].

<sup>5</sup>The results are from [39].

provides good visible results. Previous experiments show that 2 bits can achieve a good trade-off between accuracy and bitwidth reduction. In this section, only 2 bit quantization for the weights with VecQ is used as the quantization method for the DNN models.

### 6.2.1 Evaluation Settings

**Datasets.** All models are trained with the training data in the MSRA10K dataset (80 percent of the entire dataset) [40]. After training, the evaluation is conducted on multiple

datasets, including the MSRA10K (20 percent of the entire dataset), ECSSD [40], HKU-IS [40], DUTs [40], DUT-OMRON [40] and the images containing target objects and existing ground truth maps in the THUR15K [41]. The details of the selected datasets are shown in Table 10. All images are resized to  $224 \times 224$  for the training and test.

**Models.** The famous end-to-end semantic segmentation models i.e., U-Net [40], FPN [42], LinkNet [43] and UNet++ [44] are selected for the comprehensive comparison. Their detailed information is shown in Table 11. The models are

TABLE 9  
Optimal Value of  $\lambda$  for  $J_o(\phi, w'_q)$  With Bitwidth  $k$

$k$	2	3	4	5	6	7	8
AlexNet	0.9700	0.5700	0.3300	0.1900	0.1000	0.0600	0.0300
ResNet18	1.0000	0.6100	0.3600	0.2100	0.1200	0.0600	0.0300
MobileNetV2	1.0000	0.5900	0.3400	0.1900	0.1100	0.0600	0.0300
MeanError	0.0171	-0.0019	-0.0024	-0.0256	-0.0178	-0.0094	0.0023

TABLE 10  
The Datasets for Salient Object Detection

Datasets	Images	Contrast
ECSSD	1000	High
HKU-IS	4000	Low
DUTs	15572	Low
DUT-OMRON	5168	Low
MSRA10K	10000 (80/20%)	High
THUR15K	6233	Low

based on the same ResNet-50 backbone as encoder and initialized with the weights trained on the ImageNet dataset.

*Evaluation Metrics.* We choose 4 widely used metrics for a comprehensive evaluation: (1) Mean Absolute Error (MAE) [40], (2) Maximal F-measure (MaxF) [40], (3) Structural measure (S-measure) [40], and (4) Enhanced-alignment measure (E-measure) [45].

The MAE measures the average pixel-wise absolute error between the output map  $F$  and the ground-truth mask  $G$ .

$$MAE = \frac{1}{m} \sum_{k=1}^m \left( \frac{1}{H_k \times W_k} \sum_{i=1}^{H_k} \sum_{j=1}^{W_k} |G_k(i, j) - F_k(i, j)| \right). \quad (33)$$

Here  $m$  is the number of the samples,  $H_k$  and  $W_k$  are the height and width of  $G_k$ .

F-measure comprehensively evaluates the *Precision* and *Recall* with a weight parameter  $\beta$ .

$$F_\beta = \frac{(1 + \beta^2) Precision + Recall}{\beta^2 Precision + Recall}, \quad (34)$$

$\beta^2$  is empirically set to 0.3 [40]. The MaxF is the maximal  $F_\beta$  value.

S-measure considers the object-aware and region-aware structure similarities.

TABLE 11  
The Models for Salient Object Detection

Models	U-Net	FPN	LinkNet	UNet++
<b>Backbone</b>	ResNet-50	ResNet-50	ResNet-50	ResNet-50
<b>Convs</b>	64	67	69	76
<b>Parameters (M)</b>	36.54	28.67	28.78	37.7
<b>Model size (M)</b>	139.37	109.38	109.8	143.81
<b>Q. size (M)<sup>1</sup></b>	9.05	7.20	7.24	9.35

<sup>1</sup>Q. size (M) stands for the size of the VecQ quantized model.

E-measure is designed for binary map evaluations. It combines the local pixel-wise difference and the global mean value of map for comprehensive evaluation. In our evaluation, the output map is first binarized to [0,1] by comparing with the threshold of twice of its mean value [40], then the binary map is evaluated with E-measure.

We also involve direct visual comparison of the full precision and quantized weights of the selected model, to provide a more visible comparison.

### 6.2.2 Results and Analysis

The quantitative comparison results are in Table 12. The \* indicates the quantized model based on VecQ besides the full precision model. Note here, the output sizes of the last layer of FPN and FPN\* are  $56 \times 56$  and then resized to  $224 \times 224$ . Overall, the performance degradation of the quantized models based on VecQ is less than 0.01 in most metric for all models, but the size of the quantized model is reduced by more than 93 percent when compared to the models using FP32.

As shown in Table 12, all the quantized models have a less than 0.01 degradation on MSRA10K with all evaluation metrics. This is because all the models are trained on the training subset of MSRA10K, so the features of the images in it are well extracted. The other 5 data sets are only used as testing datasets and there are more degradation with the evaluation metrics, especially in MaxF and S-measure, but the degradation is maintained within 0.02.

Comparing the quantized models with their corresponding full-precision models, FPN\* performs well on almost all the test sets with all the evaluation metrics (shown as *bold italics* numbers in Table 12), showing a

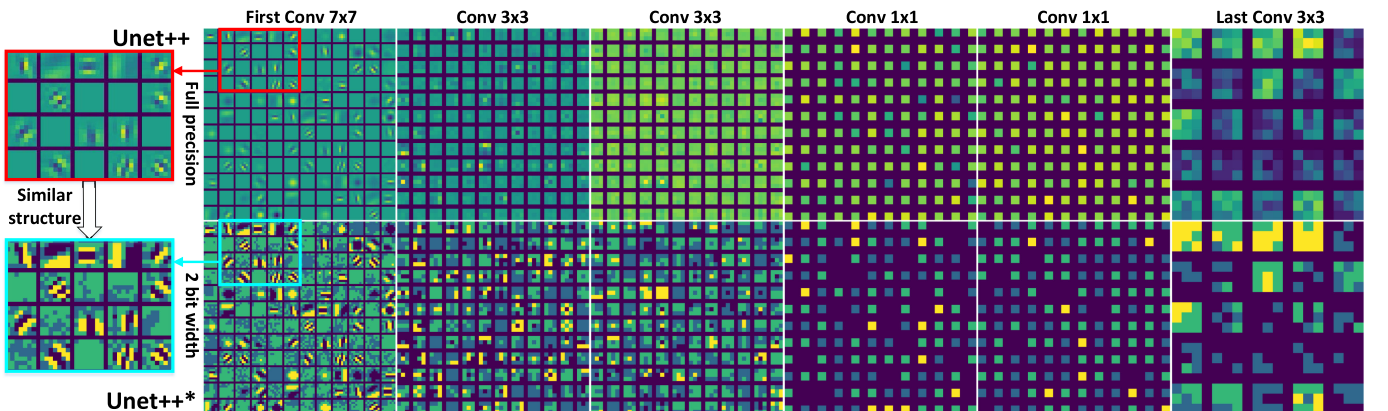


Fig. 8. The comparison of weights of convolutional layers in Unet++ and Unet++\*.



TABLE 12  
The Salient Object Detection Results

	Datasets	MSRA10K-test				ECSSD				HKU-IS				DUTs				DUT-OMRON				THUR15K			
Model	Size (M)	MAE ↓	MaxF ↑	S ↑	E ↑	MAE ↓	MaxF ↑	S ↑	E ↑	MAE ↓	MaxF ↑	S ↑	E ↑	MAE ↓	MaxF ↑	S ↑	E ↑	MAE ↓	MaxF ↑	S ↑	E ↑	MAE ↓	MaxF ↑	S ↑	E ↑
Unet	139.37	0.030	0.945	0.931	0.962	0.057	0.909	0.886	0.914	0.045	0.907	0.884	0.930	0.060	0.896	0.865	0.874	0.070	0.804	0.803	0.829	0.077	0.769	0.807	0.816
Unet*	9.05	0.032	0.940	0.926	0.959	0.064	0.896	0.871	0.906	0.050	0.893	0.870	0.923	0.065	0.885	0.852	0.871	0.071	0.793	0.795	0.835	0.081	0.749	0.797	0.815
Bias	93.51%	-0.002	0.004	0.005	0.003	-0.007	0.013	0.015	0.008	-0.005	0.014	0.015	0.007	-0.005	0.011	0.012	0.003	-0.001	0.011	0.008	-0.006	-0.005	0.020	0.010	0.001
FPN	109.38	0.043	0.920	0.899	0.949	0.070	0.882	0.854	0.901	0.059	0.875	0.848	0.911	0.072	0.875	0.835	0.861	0.081	0.777	0.772	0.812	0.087	0.750	0.778	0.804
FPN*	7.2	0.038	0.935	0.920	0.955	0.070	0.889	0.866	0.897	0.059	0.879	0.859	0.908	0.070	0.878	0.850	0.859	0.080	0.772	0.786	0.809	0.089	0.739	0.789	0.796
Bias	93.42%	0.005	-0.015	-0.022	-0.006	0.000	-0.007	-0.012	0.004	0.001	-0.005	-0.011	0.004	0.002	-0.003	-0.015	0.003	0.001	0.004	-0.014	0.003	-0.002	0.011	-0.011	0.007
Linknet	109.8	0.032	0.942	0.928	0.959	0.060	0.905	0.882	0.911	0.048	0.900	0.878	0.927	0.062	0.892	0.861	0.871	0.071	0.801	0.799	0.825	0.079	0.760	0.803	0.814
Linknet*	7.24	0.034	0.939	0.923	0.959	0.068	0.891	0.865	0.902	0.054	0.887	0.860	0.920	0.068	0.883	0.847	0.870	0.072	0.788	0.787	0.833	0.082	0.746	0.794	0.818
Bias	93.40%	-0.002	0.003	0.004	0.001	-0.008	0.014	0.017	0.008	-0.005	0.013	0.018	0.006	-0.005	0.010	0.014	0.001	-0.001	0.014	0.012	-0.008	-0.002	0.014	0.009	-0.004
UNet++	143.81	0.029	0.948	0.933	0.964	0.056	0.910	0.888	0.915	0.044	0.909	0.887	0.930	0.059	0.897	0.867	0.876	0.070	0.805	0.805	0.829	0.076	0.769	0.810	0.818
UNet++*	9.35	0.033	0.939	0.926	0.958	0.065	0.895	0.872	0.905	0.052	0.890	0.868	0.919	0.066	0.884	0.854	0.867	0.075	0.785	0.792	0.822	0.082	0.750	0.797	0.811
Bias	93.50%	-0.004	0.008	0.007	0.007	-0.009	0.015	0.016	0.010	-0.008	0.019	0.018	0.011	-0.007	0.013	0.013	0.009	-0.005	0.020	0.013	0.007	-0.006	0.019	0.013	0.007

<sup>†</sup>S and E stand for S-measure and E-measure, respectively. The ↑ indicates that the higher value shows better results and the ↓ is vice versa. The \* indicates the quantized models with VecQ-2 and the quantized model sizes are marked in italics. **Bias** row lists the difference between the full-precision model and the quantized model. A negative value is better in S↑, E↑, and MaxF↑ column but worse in the MAE↓ column.

better feature extraction capacity than the full precision version (FPN). Compared to other models, FPN outputs a relatively smaller prediction map ( $56 \times 56$ ). Although the backbone models are similar, the feature classification tasks in the FPN work on a smaller feature map with a similar size of coefficients. This provides a good chance for data quantization without effecting the detection ability because of the potential redundancy of the coefficients.

In order to further present the effectiveness of our proposed quantization solution, we print the weights of the convolution layers of the backbone model in Unet++ and Unet++\*, as shown in Fig. 8. There are three sizes of kernels involved in this model, which are  $7 \times 7$ ,  $3 \times 3$  and  $1 \times 1$ . In addition, we also compare the full precision and quantized weights of the last convolution layer for the selected model, which directly outputs the results for the detection.

In the first  $7 \times 7$  kernels, we notice a significant structural similarity between the full precision weights and the quantized weights. Since the very first layer of the backbone model extracts the overall features, the quantized weights provide a nice feature extraction ability. When the size of the kernels become smaller, we could still notice a good similarity between the full precision ones and quantized ones. Although the similarities are not significant in the Conv  $3 \times 3$  sets, they become obvious in the following Conv  $1 \times 1$  sets. The Last Conv group directly explains the comparable output results with the visible emphasized kernels and locations of the weights.

Overall, the quantized weights show a good similarity to the full precision ones in terms of the value and the location which ensures the high accuracy output when compared to the original full precision model.

## 7 CONCLUSION

In this paper, we propose a new quantization solution called VecQ. Different from the existing works, it utilizes the vector loss instead of adopting L2 loss to measure the loss of quantization. VecQ quantizes the full-precision

weight vector into a specific bitwidth with the least DQL and, hence, provides a better final model accuracy. We further introduce the fast quantization algorithm based on a reasonable prior knowledge of normally distributed weights and reduces the complexity of the quantization process in model training. The integration of VecQ into Keras [30] is also presented and used for our evaluations. The comprehensive evaluations have shown the effectiveness of VecQ on multiple datasets, models and tasks. The quantized low-bit models based on VecQ show comparable classification accuracy to models with FP32 datatype and outperform all the state-of-the-art quantization methods when the targeted bitwidth of the weights is higher than 2. Moreover, the experiments on salient object detection also show that VecQ can greatly reduce the size of the models while maintaining the performance of feature extraction tasks.

For future work, we will focus on the combination of non-linear quantization and illustrate an automated mixed-precision quantization with VecQ to achieve better performance improvement. The source code of the Keras built with VecQ could be found at <https://github.com/Gong-Cheng1919/VecQ>.

## ACKNOWLEDGMENTS

This work was supported in part by the National Natural Science Foundation (61872200), in part by the National Key Research and Development Program of China (2018YFB2100304, 2018YFB1003405), in part by the National Science Foundation of Tianjin (19JCZDJC31600, 19JCQNJC00600), and in part by the Open Project Fund of State Key Laboratory of Computer Architecture, Institute of Computing Technology, Chinese Academy of Sciences (CARCH201905). It is also supported in part by the National Research Foundation, Prime Minister's Office, Singapore under its Campus for Research Excellence and Technological Enterprise (CREATE) programme, and in part by the IBM-Illinois Center for Cognitive Computing System Research (C3SR)—a research collaboration as part of IBM AI Horizons Network.

## REFERENCES

- [1] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," in *Int. Conf. Learn. Representations (ICLR)*, 2015.
- [2] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2016, pp. 770–778.
- [3] Y. Chen, J. He, X. Zhang, C. Hao, and D. Chen, "Cloud-DNN: An open framework for mapping DNN models to cloud FPGAs," in *Proc. ACM/SIGDA Int. Symp. Field-Programmable Gate Arrays*, 2019, pp. 73–82.
- [4] C. Hao et al., "FPGA/DNN co-design: An efficient design methodology for IoT intelligence on the edge," in *Proc. 56th Annu. Des. Autom. Conf.*, 2019, pp. 1–6.
- [5] X. Zhang et al., "DNNbuilder: An automated tool for building high-performance DNN hardware accelerators for FPGAs," in *Proc. IEEE/ACM Int. Conf. Comput.-Aided Des.*, 2018, pp. 1–8.
- [6] I. Hubara, M. Courbariaux, D. Soudry, R. El-Yaniv, and Y. Bengio, "Binarized neural networks," in *Proc. Advances Neural Inf. Process. Syst.*, 2016, pp. 4107–4115.
- [7] J. Wang, Q. Lou, X. Zhang, C. Zhu, Y. Lin, and D. Chen, "Design flow of accelerating hybrid extremely low bit-width neural network in embedded FPGA," in *Proc. 28th Int. Conf. Field Programmable Logic Appl.*, 2018, pp. 163–1636.
- [8] Y. Chen et al., "T-DLA: An open-source deep learning accelerator for ternarized DNN models on embedded FPGA," in *Proc. IEEE Comput. Soc. Annu. Symp. (VLSI)*, 2019, pp. 13–18.
- [9] P. Gysel, M. Motamedi, and S. Ghiasi, "Hardware-oriented approximation of convolutional neural networks," *CoRR*, vol. abs/1604.03168, 2016.
- [10] M. Courbariaux, Y. Bengio, and J.-P. David, "Binaryconnect: Training deep neural networks with binary weights during propagations," in *Proc. 28th Int. Conf. Neural Inf. Process. Syst.*, 2015, pp. 3123–3131.
- [11] Z. Lin, M. Courbariaux, R. Memisevic, and Y. Bengio, "Neural networks with few multiplications," 2016.
- [12] S. Zhou, Y. Wu, Z. Ni, X. Zhou, H. Wen, and Y. Zou, "DoReFa-Net: Training low bitwidth convolutional neural networks with low bitwidth gradients," *CoRR*, vol. abs/1606.06160, 2016.
- [13] C. Jin, H. Sun, and S. Kimura, "Sparse ternary connect: Convolutional neural networks using ternarized weights with enhanced sparsity," in *Proc. 23rd Asia South Pacific Des. Autom. Conf.*, 2018, pp. 190–195.
- [14] S. Han, H. Mao, and W. J. Dally, "Deep compression: Compressing deep neural networks with pruning, trained quantization and Huffman coding," 2016.
- [15] F. Li, B. Zhang, and B. Liu, "Ternary weight networks," *CoRR*, vol. abs/1605.04711, 2016.
- [16] C. Leng, Z. Dou, H. Li, S. Zhu, and R. Jin, "Extremely low bit neural network: Squeeze the last bit out with ADMM," in *Proc. AAAI Conf. Artif. Intell.*, 2018, pp. 3466–3473.
- [17] P. Wang, Q. Hu, Y. Zhang, C. Zhang, Y. Liu, and J. Cheng, "Two-step quantization for low-bit neural networks," in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit.*, 2018, pp. 4376–4384.
- [18] C. Gong et al., "μ12q: An ultra-low loss quantization method for DNN," in *Proc. Int. Joint Conf. Neural Netw.*, 2019, pp. 1–8.
- [19] C. Zhu, S. Han, H. Mao, and W. J. Dally, "Trained ternary quantization," 2017.
- [20] M. Rastegari, V. Ordonez, J. Redmon, and A. Farhadi, "XNOR-Net: Imagenet classification using binary convolutional neural networks," in *Proc. Eur. Conf. Comput. Vis.*, 2016, pp. 525–542.
- [21] H. Alemdar, V. Leroy, A. Prost-Boucle, and F. Pétrot, "Ternary neural networks for resource-efficient AI applications," in *Proc. Int. Joint Conf. Neural Netw.*, 2017, pp. 2547–2554.
- [22] K. Wang, Z. Liu, Y. Lin, J. Lin, and S. Han, "HAQ: Hardware-aware automated quantization with mixed precision," in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit.*, 2019, pp. 8612–8620.
- [23] A. Zhou, A. Yao, Y. Guo, L. Xu, and Y. Chen, "Incremental network quantization: Towards lossless CNNs with low-precision weights," 2017.
- [24] M. Ghasemzadeh, M. Samragh, and F. Koushanfar, "ReBNet: Residual binarized neural network," in *Proc. IEEE 26th Annu. Int. Symp. Field-Programmable Custom Comput. Machines*, 2018, pp. 57–64.
- [25] L. Deng, P. Jiao, J. Pei, Z. Wu, and G. Li, "GXNOR-Net: Training deep neural networks with ternary weights and activations without full-precision memory under a unified discretization framework," *Neural Netw.*, vol. 100, pp. 49–58, 2018.
- [26] B. Jacob et al., "Quantization and training of neural networks for efficient integer-arithmetic-only inference," in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit.*, 2018, pp. 2704–2713.
- [27] S. R. Jain, A. Gural, M. Wu, and C. Dick, "Trained uniform quantization for accurate and efficient neural network inference on fixed-point hardware," *CoRR*, vol. abs/1903.08066, 2019.
- [28] M. Yu et al., "GradiVeQ: Vector quantization for bandwidth-efficient gradient aggregation in distributed CNN training," in *Proc. Conf. Neural Inf. Process. Syst.*, 2018, pp. 5129–5139.
- [29] S. Ioffe and C. Szegedy, "Batch normalization: Accelerating deep network training by reducing internal covariate shift," in *Proc. 32nd Int. Conf. Int. Conf. Mach. Learn.*, 2015, pp. 448–456.
- [30] F. Chollet et al., "Keras," 2015. [Online]. Available: <https://github.com/fchollet/keras>
- [31] A. Krizhevsky and G. Hinton, "Learning multiple layers of features from tiny images," *CiteSeer*, Tech. Rep., 2009.
- [32] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei, "Imagenet: A large-scale hierarchical image database," in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit.*, 2009, pp. 248–255.
- [33] A. L. Maas, R. E. Daly, P. T. Pham, D. Huang, A. Y. Ng, and C. Potts, "Learning word vectors for sentiment analysis," in *Proc. 49th Annu. Meeting Assoc. Comput. Linguistics: Human Lang. Technol.*, 2011, pp. 142–150.
- [34] M. Sum et al., "THUCTC: An efficient chinese text classifier," *GitHub Repository*, 2016. [Online]. Available: <https://github.com/thunlp/THUCTC>
- [35] M. Sandler, A. Howard, M. Zhu, A. Zhmoginov, and L.-C. Chen, "MobileNetV2: Inverted residuals and linear bottlenecks," in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit.*, 2018, pp. 4510–4520.
- [36] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "Imagenet classification with deep convolutional neural networks," in *Proc. Advances Neural Inf. Process. Syst.*, 2012, pp. 1097–1105.
- [37] M. Simon, E. Rodner, and J. Denzler, "Imagenet pre-trained models with batch normalization," *CoRR*, vol. abs/1612.01452, 2016.
- [38] S. Gross and M. Wilber, "Training and investigating residual nets," *Facebook AI Research*, vol. 6, 2016.
- [39] R. Krishnamoorthi, "Quantizing deep convolutional networks for efficient inference: A whitepaper," *CoRR*, vol. abs/1806.08342, 2018.
- [40] W. Wang, Q. Lai, H. Fu, J. Shen, and H. Ling, "Salient object detection in the deep learning era: An in-depth survey," *CoRR*, vol. abs/1904.09146, 2019.
- [41] M.-M. Cheng, N. Mitra, X. Huang, and S.-M. Hu, "Salient shape: Group saliency in image collections," *Vis. Comput.*, vol. 30, no. 4, pp. 443–453, 2014.
- [42] T.-Y. Lin, P. Dollár, R. Girshick, K. He, B. Hariharan, and S. Belongie, "Feature pyramid networks for object detection," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2017, pp. 2117–2125.
- [43] A. Chaurasia and E. Culurciello, "Linknet: Exploiting encoder representations for efficient semantic segmentation," in *Proc. IEEE Vis. Commun. Image Process.*, 2017, pp. 1–4.
- [44] Z. Zhou, M. M. R. Siddiquee, N. Tajbakhsh, and J. Liang, "UNet++: A nested U-Net architecture for medical image segmentation," in *Proc. Deep Learn. Med. Image Anal. Multimodal Learn. Clinical Decision Support*, 2018, pp. 3–11.
- [45] D.-P. Fan, C. Gong, Y. Cao, B. Ren, M.-M. Cheng, and A. Borji, "Enhanced-alignment measure for binary foreground map evaluation," in *Proc. 27th Int. Joint Conf. Artif. Intell.*, 2018, pp. 698–704.



**Cheng Gong** received the BEng degree in computer science from Nankai University, in 2016. He is currently working toward the PhD degree in the College of Computer Science, Nankai University. His main research interests include heterogeneous computing, machine learning and Internet of Things.



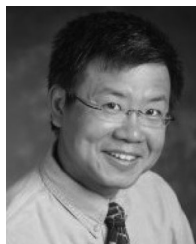
**Yao Chen** received the BS and PhD degrees from Nankai University, Tianjin, China, in 2010 and 2016, respectively. He is currently a research scientist with the Advanced Digital Sciences Center, Singapore, which is a research institute of University of Illinois at Urbana-Champaign. His research interests include reconfigurable computing, high level synthesis and high performance computing.



**Cong Hao** received the BS and MS degrees from Shanghai Jiao Tong University, in 2011 and 2014, respectively and the PhD degree in electronic engineering from Waseda University, Japan, in 2017. She is currently a postdoctoral researcher with the ECE Department, University of Illinois at Urbana-Champaign. Her current research interests include system-level and high-level synthesis, EDA techniques, and reconfigurable computing.



**Ye Lu** received the BS and PhD degrees from Nankai University, Tianjin, China, in 2010 and 2015, respectively. He is an associate professor with the College of Cyber Science, Nankai University now. His main research interests include embedded system, Internet of Things, and artificial intelligence.



**Deming Chen** (Fellow, IEEE) received the BS degree in computer science from the University of Pittsburgh, PA, in 1995, and the MS and PhD degrees in computer science from the University of California at Los Angeles, in 2001 and 2005, respectively. He is the Abel Bliss Endowed professor with the ECE Department, University of Illinois at Urbana-Champaign. He is an ACM Distinguished Speaker, and the editor-in-chief of the *ACM Transactions on Reconfigurable Technology and Systems* (TRETs). His current research interests include system-level and high-level synthesis, machine learning, GPU and reconfigurable computing, computational genomics, and hardware security.



**Tao Li** received the PhD degree in computer science from Nankai University, China, in 2007. He works at the College of Computer Science, Nankai University as a professor. He is the member of the IEEE Computer Society and the ACM, and the distinguished member of the CCF. His main research interests include heterogeneous computing, machine learning, and Internet of Things.

▷ **For more information on this or any other computing topic, please visit our Digital Library at [www.computer.org/csdl](http://www.computer.org/csdl).**