

ADDITIVE POWERS-OF-TWO QUANTIZATION: A NON-UNIFORM DISCRETIZATION FOR NEURAL NETWORKS

Anonymous authors

Paper under double-blind review

ABSTRACT

We proposed Additive Powers-of-Two (APoT) quantization, an efficient non-uniform quantization scheme that attends to the bell-shaped and long-tailed distribution of weights in neural networks. By constraining all quantization levels as a sum of several Powers-of-Two terms, APoT quantization enjoys overwhelming efficiency of computation and a good match with weights' distribution. A simple reparameterization on clipping function is applied to generate better-defined gradient for updating of optimal clipping threshold. Moreover, weight normalization is presented to refine the input distribution of weights to be more stable and consistent. Experimental results show that our proposed method outperforms state-of-the-art methods, and is even competitive with the full-precision models demonstrating the effectiveness of our proposed APoT quantization. For example, our 3-bit quantized ResNet-34 on ImageNet only drops 0.3% Top-1 and 0.2% Top-5 accuracy without bells and whistles, while the computation of our model is approximately $2\times$ less than uniformly quantized neural networks.

1 INTRODUCTION

Deep Neural Networks (DNNs) have made a significant improvement for various real-world applications. However, the huge memory and computational cost impede the mass deployment of DNNs, e.g., on resource-constrained devices. To the reduce memory footprint and computational burden, several model compression methods such as quantization (Zhou et al., 2016), pruning (Han et al., 2015) and low-rank decomposition (Denil et al., 2013) have been widely explored.

In this paper, we focus on the neural network quantization for efficient inference. Two operations are involved in the quantization process, namely clipping and projection. The clipping operation sets a full precision number to the range boundary if it is outside of the range; the projection operation maps each number (after clipping) into a predefined quantization level (a fixed number). We can see that both operations incur information loss. A good quantization method should resolve the two following questions/challenges, which correspond to two contradictions respectively.

How to determine the optimal clipping threshold to balance clipping range and projection resolution? The resolution indicates the interval between two quantization levels; the smaller the interval, the higher the resolution. The first contradiction is that given a fixed number of bits to represent weights, the range and resolution are inversely proportional. For example, a larger range can clip fewer weights; however, the resolution becomes smaller and thus damage the projection. Note that slipshod clipping of outliers can jeopardize the network a lot (Zhao et al., 2019) although they may only take 1-2% of all weights in one layer. Previous works have tried either pre-defined (Cai et al., 2017; Zhou et al., 2016) or trainable (Choi et al., 2018b) clipping thresholds, but how to find the optimal threshold during training automatically is still not resolved.

How to design quantization levels with consideration for both the computational efficiency and the distribution of weights? Most of the existing quantization approaches (Cai et al., 2017; Gong et al., 2019), use uniform quantization although non-uniform quantization can usually achieve better accuracy (Zhu et al., 2016). The reason is that projection against uniform quantization levels are much more hardware-friendly (Zhou et al., 2016). However, empirical study (Han et al., 2015) has shown that weights in a layer of DNN follow a bell-shaped and long-tailed distribution instead of a uniform distribution. In other words, a fair percentage of weights concentrate around the mean; and a few weights are of relatively high magnitude and out of the quantization range (called outliers).

The second contradiction is: considering the bell-shaped distribution of weight, it is well-motivated to assign higher resolution (i.e. smaller quantization interval) around the mean; however, such non-uniform quantization levels will introduce high computational overhead. Powers-of-Two quantization levels (Miyashita et al., 2016; Zhou et al., 2017) are then proposed because of its cheap multiplication implemented by *shift* operations on hardware, and super high resolution around the mean. However, the vanilla powers-of-two quantization method only increases the resolution near the mean and ignores other regions at all when the bit-width is increased. Consequently, it assigns inordinate quantization levels for a tiny range around the mean.

To resolve the two challenges above, we propose multiple novel quantization techniques as follows:

1. We introduce the Additive Powers-of-Two (APoT) quantization scheme for the weights of DNNs. APoT is a non-uniform quantization scheme, in which the quantization levels is a sum of several PoT terms and can adapt well to the bell-shaped distribution of weights. APoT quantization enjoys a $2\times$ speed-up compared with uniform quantization on both general and specific hardware.
2. We propose a Reparameterized Clipping Function (RCF) that can compute a more accurate gradient for the clipping threshold and thus facilitate the optimization of the clipping threshold.
3. We introduce weight normalization for neural network quantization. Normalized weights in the forward pass are more stable and consistent for clipping and projection.

Experimental results show that our proposed method outperforms state-of-the-art methods, and is even competitive with the full-precision implementation with higher computational efficiency. Specifically, our 3-bit quantized ResNet-34 on ImageNet only drops 0.3% Top-1 and 0.2% Top-5 accuracy.

2 METHODOLOGY

In this section, we first give some preliminaries in uniform quantization and PoT quantization. Then we introduce our APoT quantization and RCF. Lastly, we present a weights normalization technique for quantized neural networks.

2.1 PRELIMINARIES

Suppose kernels in a convolutional layer are represented by a 4D tensor $\mathbf{W} \in \mathbb{R}^{C_{out} \times C_{in} \times K \times K}$, where C_{out} and C_{in} are the number of output and input channels respectively, and K is the kernel size. We denote the quantization of the weights as

$$\hat{\mathbf{W}} = \Pi_{\mathcal{Q}(\alpha, b)}[\mathbf{W}, \alpha], \quad (1)$$

where α is clipping threshold and the clipping function $[\cdot, \alpha]$ clips weights into $[-\alpha, \alpha]$. After clipping, each element of \mathbf{W} is projected by $\Pi(\cdot)$ onto the quantization levels. We denote $\mathcal{Q}(\alpha, b)$ for a set of quantization levels, where b is the bit-width. For uniform quantization, the quantization levels are defined as

$$\mathcal{Q}^u(\alpha, b) = \alpha \times \left\{0, \frac{\pm 1}{2^{b-1} - 1}, \frac{\pm 2}{2^{b-1} - 1}, \dots, \pm 1\right\}. \quad (2)$$

For every floating-point number, uniform quantization maps it to a b -bit fixed-point representation (quantization levels) in $\mathcal{Q}^u(\alpha, b)$. Note that α is stored separately as a full-precision floating-point number for each whole \mathbf{W} . Convolution is done against the quantization levels first and the results are then multiplied by α . Arithmetical computation, e.g., convolution, can be implemented using low-precision fixed point operations on hardware, which are substantially cheaper than their floating-point contradictory (Goldberg, 1991). Nevertheless, uniform quantization does not match the distribution of weights, which is typically bell-shaped. A straightforward solution is to assign more quantization levels (higher resolution) for the peak of the distribution and fewer levels (lower resolution) for the tails. However, it is difficult to implement the arithmetical operations for the non-uniform quantization levels efficiently.

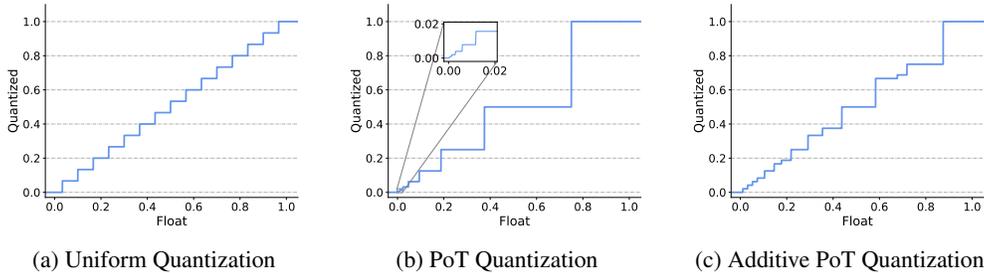


Figure 1: 4-bit quantization of unsigned data using three different quantization levels. APoT quantization has a more reasonable resolution assignment than both uniform and vanilla PoT quantization.

2.2 ADDITIVE POWERS-OF-TWO QUANTIZATION

To solve the contradiction between non-uniform resolution and hardware efficiency, Powers-of-Two (PoT) quantization (Miyashita et al., 2016; Zhou et al., 2017) is proposed by constraining quantization levels to be powers-of-two values or zero, i.e.,

$$\mathcal{Q}^p(\alpha, b) = \alpha \times \{0, \pm 2^{-2^{b-1}+1}, \pm 2^{-2^{b-1}+2}, \dots, \pm 2^{-1}, \pm 1\}. \quad (3)$$

Apparently, as a non-uniform quantizer, PoT has a higher resolution for the value range with denser weights because of its exponential property. Furthermore, multiplication between a Powers-of-two number 2^k and the other operand x can be implemented by bit-wise shift instead of bulky digital multipliers, i.e.,

$$2^k x = \begin{cases} x & \text{if } k = 0 \\ x \ll k & \text{if } k > 0, \\ x \gg k & \text{if } k < 0 \end{cases} \quad (4)$$

where \gg denotes the right *shift* operation and is computationally cheap, which only takes 1 clock cycle in modern CPU architectures.

However, we find that in PoT quantization does not benefit from more bits. Assume α is 1, as shown in Equation (3), when we increase the bit-width from b to $b + 1$, the interval $[0, \pm 2^{-2^{b-1}+1}]$ will be split into 2^b sub-intervals, whereas all other intervals remain unchanged. In other words, by increasing the bit-width, the resolution will increase only for $[-2^{-2^{b-1}+1}, 2^{-2^{b-1}+1}]$. We refer this phenomenon as the *rigid resolution* of PoT quantization. Take $\mathcal{Q}^p(1, 5)$ as an example, the two smallest positive levels are 2^{-15} and 2^{-14} , which is excessive fine-grained. In contrast, the two largest levels are 2^{-1} and 2^0 , whose interval is large enough to incur high projection error for weights between $[2^{-1}, 2^0]$, e.g., 0.75.

To tackle the *rigid resolution* problem, we propose Additive Powers-of-Two (APoT) quantization. Without loss of generality, in this section, we only consider unsigned numbers for simplicity¹. In APoT quantization, levels are viewed as a sum of several PoT terms as shown below,

$$\mathcal{Q}^a(\alpha, kn) = \gamma \times \left\{ \sum_{i=0}^{n-1} p_i \right\} \text{ where } p_i \in \left\{ 0, \frac{1}{2^i}, \frac{1}{2^{i+n}}, \dots, \frac{1}{2^{i+(2^k-1)n}} \right\}, \quad (5)$$

where γ is a scaling coefficient to make sure the maximum level in \mathcal{Q}^a is α . k is called the base bit-width, which is the bit-width for each additive term, and n is the number of additive terms. When the bit-width b and the base bit-width k is set, n can be calculated by $n = \frac{b}{k}$. The number of additive terms in APoT quantization can increase with bit-width b , which provides flexible resolution for the non-uniform levels.

We use $b = 4$ and $k = 2$ as an example to illustrate how APoT resolves the *rigid resolution* problem. For this example, we have $p_1 \in \{0, 2^0, 2^{-2}, 2^{-4}\}$, $p_2 \in \{0, 2^{-1}, 2^{-3}, 2^{-5}\}$ and $\gamma = 2\alpha/3$. Firstly, we can see the smallest positive quantization level in $\mathcal{Q}^a(1, 4)$ is $2^{-4}/3$. Compared with the original PoT levels, APoT allocates quantization levels prudently for the central area. Secondly,

¹To extend the solution for the signed number, we only need to add 1 more bit for the sign.

APoT generates 3 new quantization levels between 2^0 and 2^{-1} , to properly increase the resolution. Figure 1 compares the 3 quantization methods using 4 bits for range $[0, 1]$. APoT quantization has a reasonable distribution of quantization levels, with more levels in the peak area (near 0) and relatively higher resolution than the vanilla PoT quantization at the tail (near 1).

Relation to other quantization methods. On the one hand, the fixed-point number representations used in the uniform quantization is a special case of APoT. When $k = 1$ in Equation (5), the quantization levels is a sum of b PoT terms or 0. In the fixed-point representations, each bit indicates one specific choice of the additive terms. On the other hand, when $k = b$, there is only one PoT term and $\mathcal{Q}^a(\alpha, b)$ becomes $\mathcal{Q}^p(\alpha, b)$, i.e., PoT quantization. We can conclude that when k decreases, APoT levels are decomposed into more PoT terms, and the distribution of levels becomes more uniform. Our experiments use $k = 2$, which is an intermediate choice between the uniform case ($k = 1$) and the vanilla PoT case ($k = b$).

Computation. Multiplication for fixed-point numbers can be implemented by shifting the multiplicand (i.e., the activations) and adding the partial product. The n in Equation (5) denotes the number of additive PoT terms in the multiplier (weights), and control the speed of computation. Since $n = \frac{b}{k}$, either decreasing b or increasing k can accelerate the multiplication. Compared with uniform quantization ($k = 1$), our method ($k = 2$) is approximately $2\times$ faster. As for the full precision α , it is a coefficient for all weights in a layer and can be multiplied only once after the multiply-accumulate operation is finished.

2.3 REPARAMETERIZED CLIPPING FUNCTION

Besides the projection operation, the clipping operation $[\mathbf{W}, \alpha]$ is also important for quantization. α is a threshold that determines the value range of weights in a quantized layer. Tuning the clipping threshold α is a key challenge because of the long-tail distribution of the weights. Particularly, if α is too large (e.g., the maximum absolute value of \mathbf{W}), $\mathcal{Q}(\alpha, b)$ would have a wide range and then the projection will lead to large error as a result of insufficient resolution for the weights in the central area; if α is too small, more outliers will be clipped slipshodly. Considering the distribution of weights can be complex and differs across layers and training steps, a static clipping threshold α for all layers is not optimal.

To jointly optimize the clipping threshold α and weights via SGD during training, Choi et al. (2018b) apply the Straight-Through Estimator (STE) (Bengio et al., 2013) to do the backward propagation for the projection operation. According to STE, we have

$$\frac{\partial \hat{\mathbf{W}}}{\partial \alpha} \approx \frac{\partial [\mathbf{W}, \alpha]}{\partial \alpha} = \text{sign}(\mathbf{W}) \text{ if } |\mathbf{W}| > \alpha \text{ else } 0. \quad (6)$$

where the weights outside of the range cannot contribute to the gradients, which results in inaccurate gradient approximation. To provide a refined gradient for the clipping threshold, we design a Reparameterized Clipping Function (RCF) as

$$\hat{\mathbf{W}} = \alpha \Pi_{\mathcal{Q}(1,b)} \left[\frac{\mathbf{W}}{\alpha}, 1 \right]. \quad (7)$$

Instead of directly clipping them to $[-\alpha, \alpha]$, RCF outputs a constant clipping range and rescales weights back after the projection, which is mathematically equivalent to Equation (1) during forward. In backpropagation, STE is adopted for the projection operation and the gradients of α are calculated by

$$\frac{\partial \hat{\mathbf{W}}}{\partial \alpha} = \begin{cases} \text{sign}(\mathbf{W}) & \text{if } |\mathbf{W}| > \alpha \\ \Pi_{\mathcal{Q}(1,b)} \frac{\mathbf{W}}{\alpha} - \frac{\mathbf{W}}{\alpha} & \text{if } |\mathbf{W}| \leq \alpha \end{cases} \quad (8)$$

Compared with the normal clipping function, RCF provides more accurate gradient signals for the optimization because both weights inside ($|\mathbf{W}| \leq \alpha$) and out of ($|\mathbf{W}| > \alpha$) the range can contribute to the gradient for the clipping threshold. Particularly, the outliers are responsible for the clipping and the weights in $[-\alpha, \alpha]$ are for projection. Therefore, the update of α considers both clipping and projection and finds a sweet spot between them.

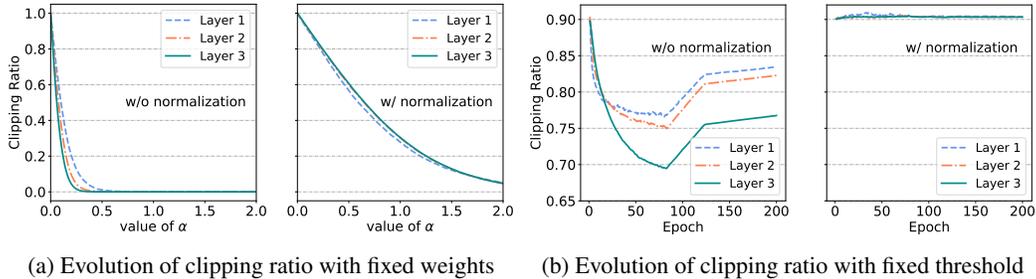


Figure 2: The evolution of clipping ratio of the first three layers in ResNet-20. (a) demonstrates clipping ratio is too sensitive to threshold to hurt its optimization without weights normalization. (b) shows that weights distribution after normalization is relatively more stable during training.

In experiments, we observe that the clipping threshold will become universally smaller when bit-width is reduced to guarantee sufficient resolution, which further validates the effectiveness of the gradient in RCF.

2.4 WEIGHT NORMALIZATION

In practice, we find that learning α for weights is quite an arduous process because the distribution of weights is pretty steep and changes frequently during training. As a result, jointly training the clipping threshold and weights parameters is hard to converge. Inspired by the crucial role of batch normalization (BN) (Ioffe & Szegedy, 2015) in activation quantization (Cai et al., 2017), we propose weight normalization (WN) to refine the distribution of weights with zero mean and unit variance,

$$\tilde{W} = \frac{W - \mu}{\sigma + \epsilon}, \text{ where } \mu = \frac{1}{I} \sum_{i=1}^I W_i, \sigma = \sqrt{\frac{1}{I} \sum_{i=1}^I (W_i - \mu)^2}, \quad (9)$$

where ϵ is a small number (typically 10^{-5}) for numerical stability, and I denotes the number of weights in one layer. Note that quantization of weights is applied right after this normalization.

Normalization is important to provide a relatively consistent and stable input distribution to both clipping and projection functions for smoother optimization of α over different layers and iterations during training. Besides, making the mean of weights to be zero can help make full use of the symmetric design of the quantization levels.

Here, we conduct a case study of ResNet-20 on CIFAR10 to illustrate how normalization for weights can help quantization. For a certain layer (at a certain training step) in ResNet-18, We firstly fix the weights, let α go from 0 to $\max(|W|)$ and plot the curve of clipping ratio (i.e. the proportion of clipped weights). As shown in Figure 2a, the change of clipping ratio is much smoother after quantization. As a result, the optimization of α will be significantly smoother. In addition, normalization also makes the distribution of weights quite more consistent over training iterations. We fix the value of α , and visualize clipping ratio over training iterations in Figure 2b. After normalization, the same α will result in almost the same clipping ratio, which improves the consistency of optimization goal for α . More experimental analysis demonstrating the effectiveness of the normalization on weights can be found in Appendix A.

2.5 TRAINING AND DEPLOYING

For activation, we use uniform quantization as we observe APoT does not outperform uniform quantization significantly for activation. During backpropagation, we use STE when computing the gradients of weights, i.e. $\frac{\partial \tilde{W}}{\partial W} = 1$. The detailed training procedure is shown in Algorithm 1.

For efficient inference, we discard the full precision weights W and only store the quantized weights \tilde{W} . Which means our proposed normalization on weights does not occupy any memory or computation resources during inference. Compared with other uniform quantization methods, APoT quantization is more efficient and effective during inference.

Algorithm 1 Forward and backward procedure for an APoT quantized convolutional layer

Input: input activations \mathbf{X}_{in} , the full precision weight tensor \mathbf{W} , the clipping threshold for weights and activations α_W, α_X , the bit-width b of quantized tensor.

Output: the output activations \mathbf{X}_{out}

- 1: Normalize weights \mathbf{W} to $\tilde{\mathbf{W}}$
- 2: Apply RCF and APoT quantization to the normalized weights $\hat{\mathbf{W}} = \alpha_W \Pi_{Q^a(1,b)} \lfloor \frac{\tilde{\mathbf{W}}}{\alpha_W}, 1 \rfloor$
- 3: Apply RCF and uniform quantization to the activations $\hat{\mathbf{X}}_{in} = \alpha_X \Pi_{Q^u(1,b)} \lfloor \frac{\mathbf{X}_{in}}{\alpha_X}, 1 \rfloor$
- 4: Compute the output activations $\mathbf{X}_{out} = Conv(\hat{\mathbf{W}}, \hat{\mathbf{X}}_{in})$
- 5: Compute the loss \mathcal{L} and the gradients $\frac{\partial \mathcal{L}}{\partial \mathbf{X}_{out}}, \frac{\partial \mathcal{L}}{\partial \hat{\mathbf{X}}_{in}}, \frac{\partial \mathcal{L}}{\partial \hat{\mathbf{W}}}$
- 6: Compute the gradients of convolution $\frac{\partial \mathcal{L}}{\partial \mathbf{X}_{in}}, \frac{\partial \mathcal{L}}{\partial \mathbf{W}}$
- 7: Compute the gradients for clipping threshold $\frac{\partial \mathcal{L}}{\partial \alpha_W}, \frac{\partial \mathcal{L}}{\partial \alpha_X}$ based on Equation (8)
- 8: Compute the gradients to the full precision weights $\frac{\partial \mathcal{L}}{\partial \mathbf{W}} = \frac{\partial \mathcal{L}}{\partial \hat{\mathbf{W}}} \frac{\partial \hat{\mathbf{W}}}{\partial \mathbf{W}}$
- 9: Update $\mathbf{W}, \alpha_W, \alpha_X$ with learning rate $\eta_W, \eta_{\alpha_W}, \eta_{\alpha_X}$

3 RELATED WORKS

Non-Uniform Quantization. Several methods are proposed to adapt to the non-uniform distribution of weights. LQ-Nets(Zhang et al., 2018) learns quantization levels based on the quantization error minimization (QEM) algorithm. Distillation (Polino et al., 2018) optimizes the quantization levels directly to minimize the task loss which reflects the behavior of their teacher network. These methods use finite floating-point numbers to quantize weights (and activations), bringing extra computation overhead compared with linear quantization. Logarithmic quantizers (Zhou et al., 2017; Miyashita et al., 2016) leverage powers-of-2 values to accelerate the computation by *shift* operations; however, they suffer from the *rigid resolution* problem.

Jointly Training. Many works have been explored to optimize the quantization parameters (e.g., α) and the weights parameters simultaneously. Zhu et al. (2016) learns positive and negative scaling coefficients respectively. LQ-Nets jointly train these parameters to minimize the quantization error. QIL (Jung et al., 2019) introduces a learnable transformer to change the quantization intervals and optimize them based on the task loss. PACT (Choi et al., 2018b) parameterizes the clipping threshold in activations and optimize it through gradient descent. However, in PACT, the gradient of α is not accurate, which only includes the contribution from outliers and ignores the contribution from other weights.

Weight Normalization. Previous works on weights normalization mainly focus on addressing the limitations of BatchNorm (Ioffe & Szegedy, 2015). Salimans & Kingma (2016); Hoffer et al. (2018) decouple direction from magnitude to accelerate the training procedure. Weight Standardization (Qiao et al., 2019) normalizes weights to zero mean and unit variance during forward pass. However, there is limited literature that studies the normalization on weights for neural network quantization. (Zhu et al., 2016) uses feature-scaling to normalize weights by dividing the maximum absolute value.

4 EXPERIMENT

In this section, we validate our proposed method on ImageNet-ILSVRC2012 (Russakovsky et al., 2015) and CIFAR10 (Krizhevsky et al., 2009) datasets. We also conduct a number of ablation studies investigating different quantizers. The checkpoints of quantized models are released anonymously².

4.1 EVALUATION ON IMAGENET

We compare our methods with several strong state-of-the-art methods on ResNet-18 and ResNet-34 (He et al., 2016), including DoReFa-Net (Zhou et al., 2016), PACT (Choi et al., 2018b), LQ-

²<https://github.com/codes4paper/ICLR2020APoT>

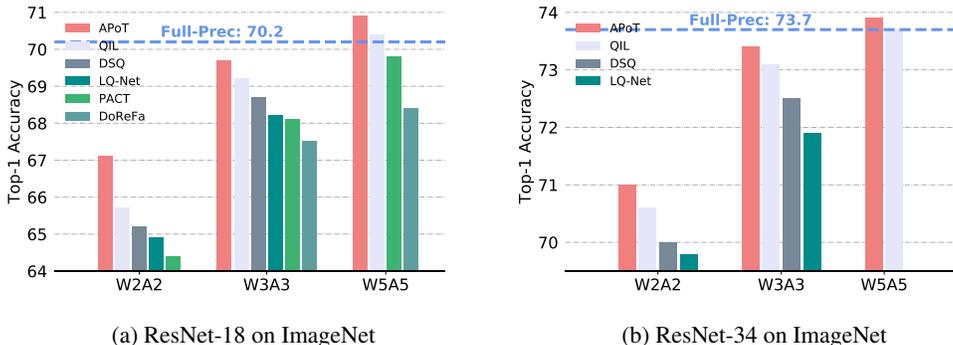


Figure 3: Top-1 accuracy comparison on ImageNet. Our proposed method is approaching and even exceeding full-precision accuracy.

Net (Zhang et al., 2018), DSQ (Gong et al., 2019), QIL (Jung et al., 2019). Both weights and activations of the networks are quantized for comparisons.

For our proposed APoT quantization algorithm, two configurations of the bit-width, i.e., 3 and 5 ($k = 2$ and $b = 1$ or 2 in Equation (5))³ are tested. We also test the 2-bit quantized network. For the 2-bit symmetric quantization levels for weights, $Q(\alpha, 2) = \{\pm\alpha, 0\}$, therefore only RCF and WN are used in this setting. To obtain a reasonable initialization, we follow Lin et al. (2017); Jung et al. (2019) to use a progressive way to initialize our model. For instance, the 5-bit quantized model is initialized from the pre-trained full precision one⁴, while the 3-bit networks is initialized from the trained 5-bit model. More details of the implementation are in the Appendix.

Overall results are shown in Figure 3. The results of DoReFa-Net are taken from Choi et al. (2018b), and the other results are quoted from the original papers. It can be observed that our 5-bit quantized network achieves even higher accuracy than the full precision model (0.7% Top-1 improvement on ResNet-18 and 0.2% Top-1 improvement on ResNet-34). Our 3-bit quantized networks are also approaching full-precision accuracy and only drop 0.5% and 0.3% accuracy on ResNet-18 and ResNet-34. When b is further reduced to 2, our model still outperforms baselines, which demonstrates the effectiveness of RCF and WN.

Table 1: Accuracy comparison of ResNet architectures on CIFAR10

Models	Methods	Accuracy(%)		
		2 bit	3 bit	5 bit
ResNet-20 (FP: 91.6)	DoReFa-Net (Zhou et al., 2016)	88.2	89.9	90.5
	PACT (Choi et al., 2018b)	89.7	91.1	91.7
	LQ-Net (Zhang et al., 2018)	90.2	91.6	-
	PACT+SAWB+fpzc (Choi et al., 2018a)	90.5	-	-
	APoT quantization (Ours)	91.0	92.2	92.3
ResNet-56 (FP: 93.2)	PACT+SAWB+fpzc (Choi et al., 2018a)	92.5	-	-
	APoT quantization (Ours)	92.9	93.9	94.0

4.2 EVALUATION ON CIFAR10

We quantize ResNet-20 and ResNet-56 (He et al., 2016) on CIFAR10 for evaluation. Again, we adopt progressive initialization and choose the quantization bit as 2, 3 and 5. More implementations and the test error curves can be found in the Appendix.

³Note that one bit is used for the sign of the weights.

⁴https://pytorch.org/docs/stable/_modules/torchvision/models/resnet.html

Table 1 summarizes the accuracy of our APoT in comparison with baselines. For 3 bit and 5 bit models, APoT quantization surpasses the full precision accuracy by 0.6 to 0.7%. It is worthwhile to note that all state-of-the-arts methods in the table use 4 levels to quantize weights into 2 bits. Our model only employs ternary weights for 2-bit representation and still achieves the highest accuracy.

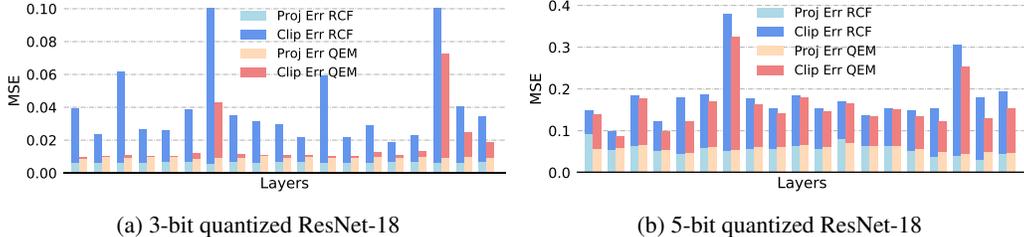


Figure 4: A summary of projection error and clipping error in different layers.

4.3 REVISITING QUANTIZATION ERROR

Typically, quantization error is defined as the mean squared error between weights $\tilde{\mathbf{W}}$ and $\hat{\mathbf{W}}$ before and after quantization respectively. The quantization error can be viewed as a sum of clipping error and projection error, defined as

$$E_{clip} = \frac{1}{I} \sum_{|\tilde{\mathbf{W}}_i| > \alpha} \left(|\tilde{\mathbf{W}}_i| - \alpha \right)^2, E_{proj} = \frac{1}{I} \sum_{|\tilde{\mathbf{W}}_i| \leq \alpha} (\tilde{\mathbf{W}}_i - \hat{\mathbf{W}}_i)^2. \quad (10)$$

Previous methods (Zhang et al., 2018; Cai et al., 2017) seek to minimize the quantization error (i.e. $\min(E_{clip} + E_{proj})$) to obtain the optimal clipping threshold, while RCF is directly optimized by the final training loss to balance projection error and clipping error. We compare the Quantization Error Minimization (QEM) method with our RCF on the quantized ResNet-18 model. Figure 4 gives an overview of clipping error and projection error using RCF or QEM. For the 5-bit quantized model, RCF has a much higher quantization error. The projection error obtained by RCF is lower than QEM and QEM significantly reduces the clipping error. Based on the truth that RCF has better final accuracy, we can infer that projection error has a higher priority in RCF. When quantizing to 3-bit, the clipping error in RCF still exceeds QEM except for the first quantized layer. This means RCF can identify whether the projection is more important than the clipping over different layers and bit-width. Generally, the insight behind is that simply minimizing the quantization error may not be the best choice and it is more direct to optimize threshold with respect to training loss.

4.4 QUANTIZER COMPARISON

Since other quantization methods do not use RCF or normalization on weights, the effectiveness of our APoT quantization is not well evaluated. Hence we implement our algorithm with different quantization levels $\mathcal{Q}(a, b)$ while keeping other settings unchanged. For 5 bit quantized models, we evaluate uniform quantization, APoT quantization, and the vanilla PoT quantization. As for 3 bit, APoT resembles PoT, therefore we evaluate APoT quantization and uniform quantization. Table 2 shows the Top-1 accuracy of ResNet-18 on ImageNet with the above quantization methods. For both 3-bit and 5-bit setting, APoT achieves higher accuracy than uniform quantization as well as higher hardware efficiency.

Table 2: Quantizer comparison

Bit-width	5 / 5	3 / 3
APoT	70.9	69.7
Uniform	70.7	69.4
PoT	70.3	-

REFERENCES

Yoshua Bengio, Nicholas Léonard, and Aaron Courville. Estimating or propagating gradients through stochastic neurons for conditional computation. *arXiv preprint arXiv:1308.3432*, 2013.

- Zhaowei Cai, Xiaodong He, Jian Sun, and Nuno Vasconcelos. Deep learning with low precision by half-wave gaussian quantization. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 5918–5926, 2017.
- Jungwook Choi, Pierce I-Jen Chuang, Zhuo Wang, Swagath Venkataramani, Vijayalakshmi Srinivasan, and Kailash Gopalakrishnan. Bridging the accuracy gap for 2-bit quantized neural networks (qnn). *arXiv preprint arXiv:1807.06964*, 2018a.
- Jungwook Choi, Zhuo Wang, Swagath Venkataramani, Pierce I-Jen Chuang, Vijayalakshmi Srinivasan, and Kailash Gopalakrishnan. Pact: Parameterized clipping activation for quantized neural networks. *arXiv preprint arXiv:1805.06085*, 2018b.
- Misha Denil, Babak Shakibi, Laurent Dinh, Marc’Aurelio Ranzato, and Nando De Freitas. Predicting parameters in deep learning. In *Advances in neural information processing systems*, pp. 2148–2156, 2013.
- David Goldberg. What every computer scientist should know about floating-point arithmetic. *ACM Computing Surveys (CSUR)*, 23(1):5–48, 1991.
- Ruihao Gong, Xianglong Liu, Shenghu Jiang, Tianxiang Li, Peng Hu, Jiazhen Lin, Fengwei Yu, and Junjie Yan. Differentiable soft quantization: Bridging full-precision and low-bit neural networks. *arXiv preprint arXiv:1908.05033*, 2019.
- Song Han, Huizi Mao, and William J Dally. Deep compression: Compressing deep neural networks with pruning, trained quantization and huffman coding. *arXiv preprint arXiv:1510.00149*, 2015.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 770–778, 2016.
- Elad Hoffer, Ron Banner, Itay Golan, and Daniel Soudry. Norm matters: efficient and accurate normalization schemes in deep networks. In *Advances in Neural Information Processing Systems*, pp. 2160–2170, 2018.
- Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. *arXiv preprint arXiv:1502.03167*, 2015.
- Sangil Jung, Changyong Son, Seohyung Lee, Jinwoo Son, Jae-Joon Han, Youngjun Kwak, Sung Ju Hwang, and Changkyu Choi. Learning to quantize deep networks by optimizing quantization intervals with task loss. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 4350–4359, 2019.
- Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. Technical report, Citeseer, 2009.
- Xiaofan Lin, Cong Zhao, and Wei Pan. Towards accurate binary convolutional neural network. In *Advances in Neural Information Processing Systems*, pp. 345–353, 2017.
- Daisuke Miyashita, Edward H Lee, and Boris Murmann. Convolutional neural networks using logarithmic data representation. *arXiv preprint arXiv:1603.01025*, 2016.
- Antonio Polino, Razvan Pascanu, and Dan Alistarh. Model compression via distillation and quantization. *arXiv preprint arXiv:1802.05668*, 2018.
- Siyuan Qiao, Huiyu Wang, Chenxi Liu, Wei Shen, and Alan Yuille. Weight standardization. *arXiv preprint arXiv:1903.10520*, 2019.
- Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, et al. Imagenet large scale visual recognition challenge. *International journal of computer vision*, 115(3):211–252, 2015.
- Tim Salimans and Durk P Kingma. Weight normalization: A simple reparameterization to accelerate training of deep neural networks. In *Advances in Neural Information Processing Systems*, pp. 901–909, 2016.

- Dongqing Zhang, Jiaolong Yang, Dongqiangzi Ye, and Gang Hua. Lq-nets: Learned quantization for highly accurate and compact deep neural networks. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pp. 365–382, 2018.
- Ritchie Zhao, Yuwei Hu, Jordan Dotzel, Chris De Sa, and Zhiru Zhang. Improving neural network quantization without retraining using outlier channel splitting. In *International Conference on Machine Learning*, pp. 7543–7552, 2019.
- Aojun Zhou, Anbang Yao, Yiwen Guo, Lin Xu, and Yurong Chen. Incremental network quantization: Towards lossless cnns with low-precision weights. *arXiv preprint arXiv:1702.03044*, 2017.
- 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.
- Chenzhuo Zhu, Song Han, Huizi Mao, and William J Dally. Trained ternary quantization. *arXiv preprint arXiv:1612.01064*, 2016.

A HOW DOES NORMALIZATION HELP QUANTIZATION

In this section, we show some experimental results to illustrate the effect of our weights normalization in quantization neural networks.

A.1 WEIGHTS DISTRIBUTION

We visualize the density distribution of weights before normalization \mathbf{W} and after normalization $\tilde{\mathbf{W}}$ during training to demonstrate its effectiveness.

Figure 5a demonstrates the density distribution of the fifth layer of the 5-bit quantized ResNet-18, from which we can see that the density of the unnormalized weights could be extensive high (> 8) in the centered area. Such distribution indicates that even a tiny change of clipping threshold would bring a significant effect on clipping when α is small, as shown in Figure 2a. , which means a small learning rate for α is needed. However, if the learning rate is too small, the change of α cannot follow the change of weights distribution because weights are also updated according to Figure 5a. Thus it is unfavourable to train the clipping threshold for unnormalized weights, while Figure 5b shows that the normalized weights can have a stable distribution. Furthermore, the dashed line in the figure indicates \mathbf{W} usually do not have zero mean, which may not utilize the symmetric design of quantization levels.

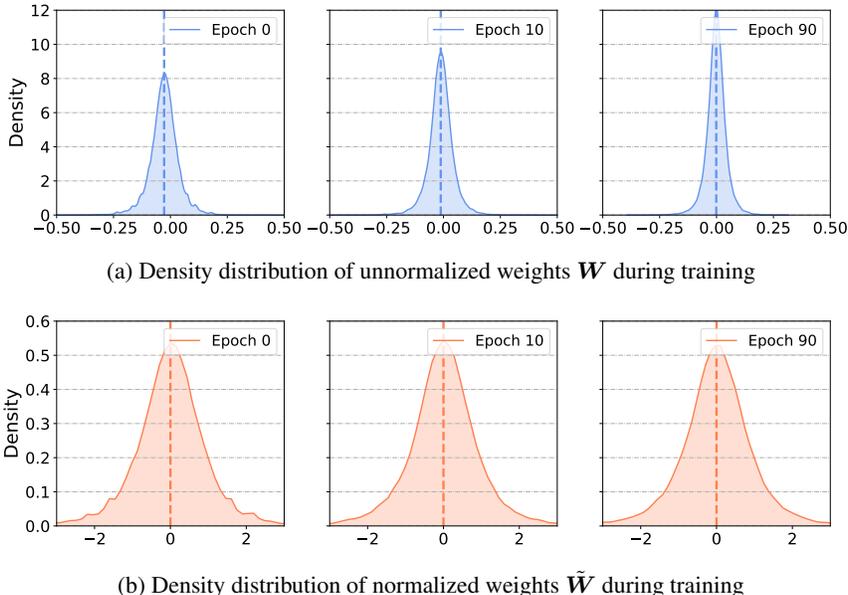


Figure 5: When weights are normalized the distribution of weights are more stable. The dashed line shows the mean value of weights.

A.2 TRAINING BEHAVIOR

The above experiments use normalization during training to compare the distribution of weights. In this section, we compare the training of quantization neural networks with and without normalization to investigate the real effect of WN. Here, we train a 3-bit quantized (full precision for activations) ResNet-20 from scratch, and compare the results under different learning rate for α . The results is shown in Table 3, from which we can find that if weights are normalized during training, the network can converge to descent performances and is robust to the learning rate of clipping threshold. However, if the weights are not normalized, the network would diverge if the learning rate for α is too high. Even if the learning rate is set to a lower value, the network do not outperform the normalized one. Based on the training behaviors, we can conclude that the learning for clipping threshold without WN in QNNs need a careful choice.

Table 3: Accuracy comparison of 3-bit quantized ResNet-20 on CIFAR10.

Learning Rate	0.1	0.01	0.001	0.0001
w/ Normalization	91.6	91.7	91.6	91.8
w/o Normalization	0.2	0.2	62.8	84.7

B EXPERIMENTS DETAILS

Table 4: Performance in terms of Top-1 and Top-5 accuracy comparison of ResNet-18 and ResNet-34 on ImageNet with several methods.

Methods	Bit-width (W / A)	ResNet-18		ResNet-34	
		Top-1	Top-5	Top-1	Top-5
Full-Precision	32 / 32	70.2	89.4	73.7	91.3
ABC-Nets (Lin et al., 2017)	5 / 5	65.0	85.9	68.4	88.2
DoReFa-Net (Zhou et al., 2016)	5 / 5	68.4	88.3	-	-
PACT (Choi et al., 2018b)	5 / 5	69.8	89.3	-	-
QIL (Jung et al., 2019)	5 / 5	70.4	-	73.7	-
APoT quantization (Ours)	5 / 5	70.9	89.7	73.9	91.6
ABC-Nets (Lin et al., 2017)	3 / 3	61.0	83.2	66.4	87.4
DoReFa-Net (Zhou et al., 2016)	3 / 3	67.5	87.6	-	-
PACT (Choi et al., 2018b)	3 / 3	68.1	88.2	-	-
LQ-Net (Zhang et al., 2018)	3 / 3	68.2	87.9	71.9	90.2
DSQ (Gong et al., 2019)	3 / 3	68.7	-	72.5	-
QIL (Jung et al., 2019)	3 / 3	69.2	-	73.1	-
APoT quantization (Ours)	3 / 3	69.7	88.9	73.4	91.1
DoReFa-Net (Zhou et al., 2016)	2 / 2	62.6	84.6	-	-
PACT (Choi et al., 2018b)	2 / 2	64.4	85.6	-	-
LQ-Net (Zhang et al., 2018)	2 / 2	64.9	85.9	69.8	89.1
DSQ (Gong et al., 2019)	2 / 2	65.2	-	70.0	-
QIL (Jung et al., 2019)	2 / 2	65.7	-	70.6	-
PACT+SAWB+fp32 (Choi et al., 2018a)	2 / 2	67.0	-	-	-
APoT quantization (Ours)	2 / 2	67.1	87.2	71.0	89.9

B.1 IMPLEMENTATIONS DETAILS

The ImageNet dataset consists of 1.2M training and 50K validation images. We use a standard data preprocess in the original paper (He et al., 2016). For training images, they are randomly cropped and resized to 224×224 . Validation images are center-cropped to the same size. We use the Pytorch official code⁵ to construct ResNet-18 and ResNet-34, and they are initialized from the released pre-trained model. We use stochastic gradient descent (SGD) with momentum of 0.9 to optimize both weight parameters and the clipping threshold simultaneously. Batch size is set to 1024 and learning rate starts from 0.1 with a decay factor of 0.1 at epoch 30,60,90. The network is trained up to 120 epochs and weight decay is set to 5×10^{-5} for ResNet-18 and 10^{-4} for ResNet-34. We do not quantize the first layer and the last layer as other methods did (Jung et al., 2019; Zhang et al., 2018; Choi et al., 2018b).

The CIFAR10 dataset contains 50K training and 10K test images with 32×32 pixel. The ResNet architectures for CIFAR10 (He et al., 2016) contains a convolutional layer followed by 3 residual blocks and a final FC layer. We train full precision ResNet-20 and ResNet-56 firstly and use them as initialization for quantized models. All networks were trained for 200 epochs with a mini-batch size of 128. SGD with momentum of 0.9 was adopted to optimize the parameters. Learning rate started from 0.1 and was scaled by 0.1 at epoch 80,120. Weight decay was set to 10^{-4} and we do not quantize the first and last layers.

⁵<https://github.com/pytorch/vision/blob/master/torchvision/models/resnet.py>

For clipping threshold α , we set 8.0 for activations and 3.0 for weights as initial value when training a 5-bit quantized model. The learning rate of α is set to 0.01 and 0.03 for weights and activations, respectively. During practice, we found that the learning rate of α merely do not influence the network performance. Different from PACT (Choi et al., 2018b), the update of α in our works already consider the projection error, so we do not require a relatively large L2-regularization for it. In practice, the network works fine when the weight decay for α is set to $1e-5$.

We report the top-1 and top-5 accuracy in Table 4 and the test error (and training error) of our quantized model in Figure 6. Our method achieves state-of-the-art accuracy. Though Choi et al. (2018a) only has 0.1% top-1 accuracy difference in 2-bit quantized ResNet-18, their method uses full precision shortcut to increase the network performance. Our method is more effective and efficient compared with the above uniform or non-uniform counterparts.

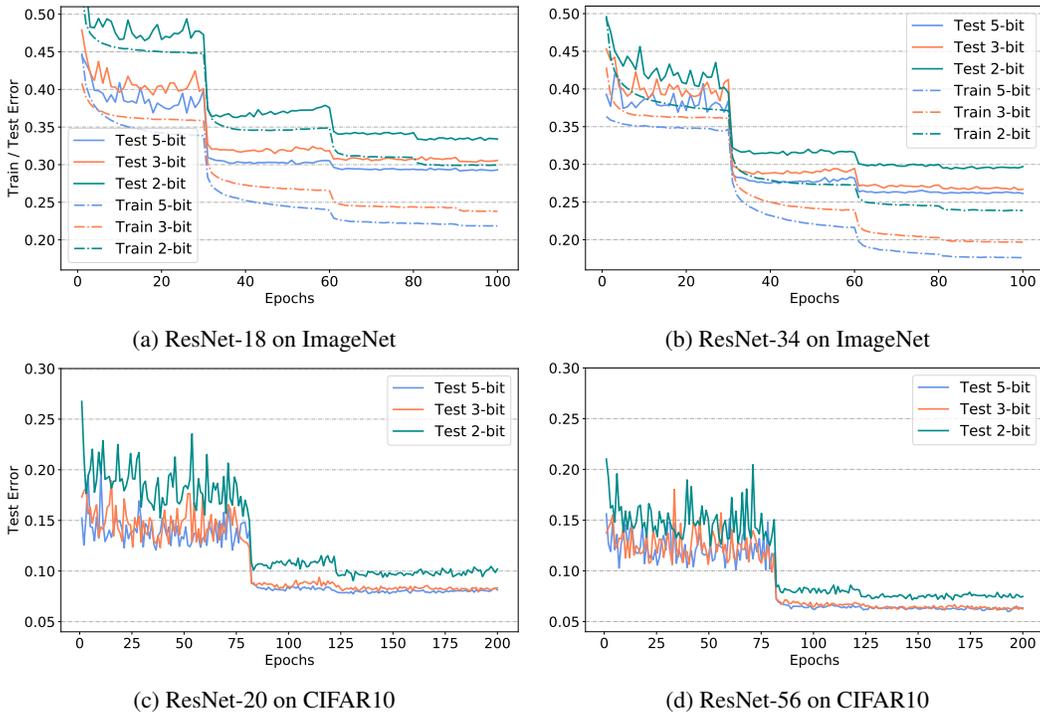


Figure 6: Training and test error with various bit-width quantized models.

B.2 APoT QUANTIZATION FOR ACTIVATIONS

In this section, we apply APoT quantization to activations. We also use two configurations ($k = 2, n = 1$ or $n = 2$). Note that activations after the ReLU function is positive, therefore we do not need additional sign bit for activations. For weights quantization, we use uniform quantization. Other settings are kept the same.

Table 5 summarizes the top-1 accuracy comparison of APoT and uniform quantization for activations. For 4-bit quantization, APoT achieves the same accuracy with uniform quantization, while the uniform quantization can surpass APoT in the 2-bit scenario. Our proposed APoT quantization do not have significant advantages when quantizing activations.

Table 5: Quantizer comparison

Bit-width	4 / 4	2 / 2
APoT	70.6	67.0
Uniform	70.6	67.1