

OPTIMAL BINARY QUANTIZATION FOR DEEP NEURAL NETWORKS

Anonymous authors

Paper under double-blind review

ABSTRACT

Quantizing weights and activations of deep neural networks results in significant improvement in inference efficiency at the cost of lower accuracy. A source of the accuracy gap between full precision and quantized models is the quantization error. In this work, we focus on the binary quantization, in which values are mapped to -1 and 1. We introduce several novel quantization algorithms: optimal 2-bits, optimal ternary, and greedy. Our quantization algorithms can be implemented efficiently on the hardware using bitwise operations. We present proofs to show that our proposed methods are optimal, and also provide empirical error analysis. We conduct experiments on the ImageNet dataset and show a reduced accuracy gap when using the proposed optimal quantization algorithms.

1 INTRODUCTION

A major challenge in the deployment of Deep Neural Networks (DNNs) is their high computational cost. Finding effective methods to improve run-time efficiency is still an area of research. We can group various approaches taken by researchers into the following three categories.

Hardware optimization: Specifically designed hardware are deployed to efficiently perform computations in ML tasks. **Compiler optimization:** Compression and fusion techniques coupled with efficient hardware-aware implementations, such as dense and sparse matrix-vector multiplication, are used. **Model optimization:** Run-time performance can also be gained by modifying the model structure and the underlying arithmetic operations. While hardware and compiler optimization is typically lossless (i.e. incurs no loss in model accuracy), model optimization trades-off computational cost (memory, runtime, or power) for model accuracy. For example, by scaling the width of the network (Zagoruyko & Komodakis, 2016). The goal of model optimization is to improve the trade-off between computational cost and model accuracy. The methods we describe in this work fall into this category.

1.1 RELATED WORKS

In this section, we summarize different model optimization techniques.

1.1.1 ARCHITECTURE OPTIMIZATION

One strategy to construct efficient DNNs is to define a template from which efficient computational blocks can be generated. Multiple instantiations of these blocks are then chained together to form a DNN. SqueezeNet (Iandola et al., 2016), MobileNets (Howard et al., 2017; Sandler et al., 2018), ShuffleNets (Zhang et al., 2018b; Ma et al., 2018), and ESPNets (Mehta et al., 2018; 2019) fall into this category. Complementary to these methods, NASNet (Zoph et al., 2018) and EfficientNet (Tan & Le, 2019) search for an optimal composition of blocks restricted to a computational budget (e.g., FLOPS) by changing the resolution, depth, width, or other parameters of each layer.

1.1.2 PRUNING AND COMPRESSION

Several methods have been proposed to improve runtime performance by detecting and removing computational redundancy. Methods in this category include low-rank acceleration (Jaderberg et al., 2014), the use of depth-wise convolution in Inception (Szegedy et al., 2015), sparsification of kernels

in deep compression (Han et al., 2015), re-training redundant neurons in DSD (Han et al., 2016b), depth-wise separable convolution in Xception (Chollet, 2017), pruning redundant filters in PFA (Suau et al., 2018), finding an optimal sub-network in lottery ticket hypothesis (Frankle & Carbin, 2018), and separating channels based on the features resolution in octave convolution (Chen et al., 2019). While some of these compression methods can be applied to a trained network, most add training-time constraints to create a computationally efficient model.

1.1.3 LOW-PRECISION ARITHMETIC AND QUANTIZATION

Another avenue to improve runtime performance (and the focus of this work) is the use of low-precision arithmetic. The idea is to use fewer bits to represent weights and activations by quantizing their representations. Some instances of these strategies already exist in AI compilers, where it is common to cast weights of a trained model from 32 bits to 16 or 8 bits. However, in general, post-training quantization reduces the model accuracy. This can be addressed by incorporating lower-precision arithmetic into the training process (during-training quantization), allowing the resulting model to better adapt to the lower precision. For example, in Gupta et al. (2015); Jacob et al. (2018) the authors use 16 and 8 bits fixed-point representation to train DNNs.

Using fewer bits results in dramatic memory savings. This has motivated research into methods that use a single bit to represent a scalar weight: In Courbariaux et al. (2015) the authors train models with weights quantized to the values in $\{-1, 1\}$. While this results in a high level of compression, model accuracy can drop significantly. Li et al. (2016) and Zhu et al. (2016) reduce the accuracy gap between full precision and quantized models by considering ternary quantization (using the values in $\{-1, 0, 1\}$), at the cost of slightly less compression.

To further improve the computational efficiency, the intermediate activation tensors (feature maps) can also be quantized. When this is the case, an implementation can use high-performance operators that act on quantized inputs, for example a convolutional block depicted in Figure 1(left). This idea has been explored in (Courbariaux et al., 2016; Rastegari et al., 2016; Zhou et al., 2016; Hubara et al., 2017; Mishra et al., 2017; Lin et al., 2017; Cai et al., 2017; Ghasemzadeh et al., 2018; Zhang et al., 2018a; Choi et al., 2018), and many other works.

We call a mapping from a tensor with full precision entries to a tensor with the same shape but with values in $\{-1, 1\}$ a **binary quantization**. When both weights and activations of a DNN are quantized using binary quantization, fast and power-efficient kernels which use bitwise operations can be implemented. Observe that the inner-product between two vectors with entries in $\{-1, 1\}$ can be written as bitwise XNOR operations followed by bit-counting (Courbariaux et al., 2016). However, the quantization of both weights and activations further reduces the model accuracy.

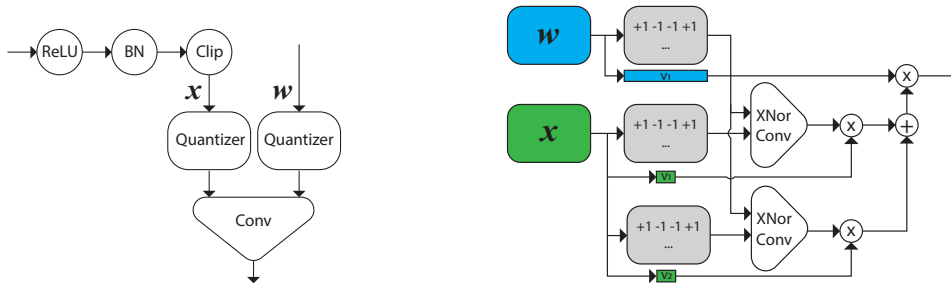


Figure 1: **Left:** The convolutional block used in this paper. **Right:** When both weights and activations are quantized using binary quantization, the convolution can be implemented efficiently using bitwise XNOR and bit-counting operations. See Section 3.2 for more details.

1.2 MAIN CONTRIBUTIONS

In this work, we analyze the accuracy of binary quantization when applied to both weights and activations of a DNN, and propose methods to improve the quantization accuracy:

- We present an analysis of the quantization error and show that scaled binary quantization is a good approximation (Section 2).
- We derive the optimal 1-bit (Section 3.1.1), 2-bits (Section 3.2.2), and ternary (Section 3.2.3) scaled binary quantization algorithms.
- We propose a greedy k -bits quantization algorithm (Section 3.2.4).
- Experiments on the ImageNet dataset show that the optimal algorithms have reduced quantization error, and lead to improved classification accuracy (Section 5).

2 LOW-RANK BINARY QUANTIZATION

Binary quantization (that maps entries of a tensor to $\{-1, 1\}$) of weights and activation tensors of a neural network can significantly reduce the model accuracy. A remedy to retrieve this accuracy loss is to scale the binarized tensors with few full precision values. For example, Hubara et al. (2017) learn a scaling for each channel from the parameters of batch-normalization, and Rastegari et al. (2016) scale the quantized activation tensors using the channel-wise average of pixel values.

In this section, using low-rank matrix analysis, we analyze different scaling strategies. We conclude that multiplying the quantized tensor by a single scalar, which is computationally the cheapest option, has approximately the same accuracy as the more expensive alternatives.

We first introduce the rank- k binary quantization— an approximation to a matrix $\mathbf{X} \in \mathbb{R}^{m \times n}$:

$$\mathbf{X} \simeq \mathbf{X}_k \odot \mathbf{S}, \quad (1)$$

where $\mathbf{X}_k \in \mathbb{R}^{m \times n}$ is a rank- k matrix, $\mathbf{S} \in \{-1, 1\}^{m \times n}$, and \odot is element-wise multiplication (Hadamard product). Note that this approximation is also defined for tensors, after appropriate reshaping. For example, for an image classification task, we can reshape the output of a layer of a DNN with shape $h \times w \times n$, where h , w , and n are height, width, and number of channels, respectively, into an $m \times n$ matrix with $m = hw$ rows and one column per channel.

We define the error of a rank- k binary quantization as $\|\mathbf{X} - \mathbf{X}_k \odot \mathbf{S}\|_F$, where $\|\cdot\|_F$ is the Frobenius norm. Entries of \mathbf{S} are in $\{-1, 1\}$, therefore, the quantization error is equal to $\|\mathbf{X} \odot \mathbf{S} - \mathbf{X}_k\|_F$. Note that $\|\mathbf{X} \odot \mathbf{S}\|_F^2$ (the total energy), which is equal to sum of the squared singular values, is the same for any \mathbf{S} . Different choices of \mathbf{S} change the distribution of the total energy among components of the Singular Value Decomposition (SVD) of $\mathbf{X} \odot \mathbf{S}$. The optimal rank- k binary quantization is achieved when most of the energy of $\mathbf{X} \odot \mathbf{S}$ is in its first k components.

In Rastegari et al. (2016), the authors proposed to quantize the activations by applying the sign function and scale them by their channel-wise average. We can formulate this scaling strategy as a special rank-1 binary quantization $\mathbf{X} \simeq \mathbf{a}\mathbf{1}^\top \odot \text{sign}(\mathbf{X})$, where

$$a_i = \frac{\sum_{j=1}^n |\mathbf{X}_{i,j}|}{n} \quad \text{for } 1 \leq i \leq m, \quad \text{sign}(x) = \begin{cases} -1 & \text{if } x < 0 \\ 1 & \text{if } x \geq 0 \end{cases}, \quad (2)$$

and $\mathbf{1}$ is an n -dimensional vector with all entries 1.

In this paper, we show that the optimal rank-1 binary quantization is given by $\mathbf{S} = \text{sign}(\mathbf{X})$ and $\mathbf{X}_1 = \text{truncated}_1\text{-SVD}(|\mathbf{X}|)$, where $\text{sign}(\mathbf{X})$ is the element-wise sign of \mathbf{X} , and $\text{truncated}_1\text{-SVD}(|\mathbf{X}|) = \sigma_1 \mathbf{u}_1 \mathbf{v}_1^\top$ is the first component of the SVD of $\mathbf{X} \odot \text{sign}(\mathbf{X}) = |\mathbf{X}|$ (for the proof see Appendix A).

We empirically analyze the accuracy of the optimal rank-1 binary quantization for a random matrix \mathbf{X} , where its entries are sampled independently from the standard normal distribution. This is a relevant example since after the application of Batch Normalization (BN) (Ioffe & Szegedy, 2015) activation tensors are expected to have a similar distribution (see Figure 1). In Figure 2(left), we show the distribution of energy $(\sigma_i^2(|\mathbf{X}|)/\|\mathbf{X}\|_F^2)$ for $|\mathbf{X}|$, where the first singular value captures most of the energy $\sigma_1^2(|\mathbf{X}|)/\|\mathbf{X}\|_F^2 \simeq 0.64$. Moreover, we plot the first left and right singular vectors of $|\mathbf{X}|$ in Figure 2(right) and observe that they are almost constant vectors. Therefore, a scalar multiple of $\text{sign}(\mathbf{X})$ approximates \mathbf{X} well: $\mathbf{X} \simeq \sigma_1 \mathbf{u}_1 \mathbf{v}_1^\top \odot \text{sign}(\mathbf{X}) \simeq v \mathbf{1}\mathbf{1}^\top \odot \text{sign}(\mathbf{X}) = v \text{sign}(\mathbf{X})$, where $v \in \mathbb{R}_{\geq 0}$. We call this **scaled binary quantization**.

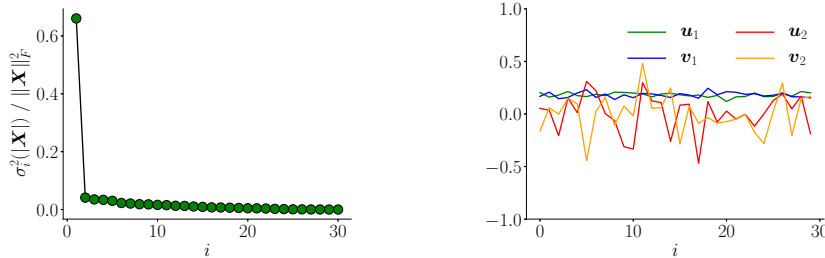


Figure 2: **Left:** Distribution of energy for $|\mathbf{X}|$, where $\mathbf{X} \in \mathbb{R}^{30 \times 30}$ is a standard normal random matrix. **Right:** Entries of the first left and right singular vectors of $|\mathbf{X}|$ (shown in green and blue) are almost constant.

3 SCALED BINARY QUANTIZATION

In Section 2 we showed that scaled binary quantization is a good approximation to activation and weight tensors of a DNN. Next we show how we can further improve the accuracy of scaled binary quantization using more bits. To simplify the presentation (1) we flatten matrix $\mathbf{X} \in \mathbb{R}^{m \times n}$ in to a vector $\mathbf{x} \in \mathbb{R}^N$ with $N = mn$, and (2) we assume the entries of \mathbf{x} are different realizations of a random variable x with underlying probability distribution p . In practice, we compute all statistics using their unbiased estimators from vector \mathbf{x} (e.g., $\sum_i x_i / N$ is an unbiased estimator of $\mathbb{E}_{x \sim p}[x]$). Furthermore, for $f : \mathbb{R} \rightarrow \mathbb{R}$, we denote element-wise application of f to entries of \mathbf{x} by $f(\mathbf{x})$. The quantized approximation of \mathbf{x} is denoted by \mathbf{x}^q , and the error of quantization is $\|\mathbf{x} - \mathbf{x}^q\|_2$.

3.1 1-BIT QUANTIZATION

A 1-bit scaled binary quantization of \mathbf{x} is:

$$\mathbf{x} \simeq \mathbf{x}^q = v s(\mathbf{x}), \quad (3)$$

which is determined by a scalar $v \in \mathbb{R}_{\geq 0}$ and a function $s : \mathbb{R} \rightarrow \{-1, 1\}$. Finding the optimal 1-bit scaled binary quantization can be formulated as the following optimization problem:

$$\begin{aligned} \underset{v, s}{\text{minimize}} \quad & \int_{-\infty}^{+\infty} p(x) (v s(x) - x)^2 dx \\ \text{s.t.} \quad & s : \mathbb{R} \rightarrow \{-1, 1\} \\ & v \in \mathbb{R}_{\geq 0} \end{aligned} \quad (4)$$

3.1.1 OPTIMAL 1-BIT ALGORITHM

The solution of problem (4) is given by $v = \mathbb{E}_{x \sim p}[|x|]$ and $s(x) = \text{sign}(x)$ (for the proofs see Appendix B). Therefore, for a vector \mathbf{x} the optimal scaled binary quantization is given by

$$\mathbf{x} \simeq \mathbf{x}^q = \frac{\sum_i |x_i|}{N} \text{sign}(\mathbf{x}), \quad (5)$$

where $\frac{\sum_i |x_i|}{N}$ is an unbiased estimator of $\mathbb{E}_{x \sim p}[|x|]$.

3.2 k -BITS QUANTIZATION

We can further improve the accuracy of scaled binary quantization by adding more terms to the approximation (3). A k -bits scaled binary quantization of \mathbf{x} is

$$\mathbf{x} \simeq \mathbf{x}^q = \sum_{i=1}^k v_i s_i(\mathbf{x}), \quad (6)$$

which is determined by k scalars $v_1 \geq v_2 \geq \dots \geq v_k \geq 0$ and functions $s_i : \mathbb{R} \rightarrow \{-1, 1\}$.

When both weights, \mathbf{w} , and activations, \mathbf{x} , are quantized using scaled binary quantization (6), their inner-product can be written as:

$$\langle \mathbf{x}^q, \mathbf{w}^q \rangle = \sum_{i=1}^{k^a} \sum_{j=1}^{k^w} v_i^a v_j^w \langle \mathbf{s}_i^a, \mathbf{s}_j^w \rangle, \quad (7)$$

where $\mathbf{x}^q = \sum_{i=1}^{k^a} v_i^a \mathbf{s}_i^a$ and $\mathbf{w}^q = \sum_{j=1}^{k^w} v_j^w \mathbf{s}_j^w$ are quantized activations and weights with k^a and k^w bits, respectively, $\mathbf{s}_i^a = \mathbf{s}_i^a(\mathbf{x})$, and $\mathbf{s}_j^w = \mathbf{s}_j^w(\mathbf{w})$. This inner-product can be computed efficiently using bitwise XNors followed by bit-counting (see Figure 1(right) for an example with $k^a = 2$ and $k^w = 1$).

Finding the optimal k -bits scaled binary quantization can be formulated as the following optimization problem:

$$\begin{aligned} \text{minimize}_{s_i, v_i} \quad & \int_{-\infty}^{+\infty} p(x) \left(\left(\sum_{i=1}^k v_i s_i(x) \right) - x \right)^2 dx \\ \text{s.t.} \quad & \forall 1 \leq i \leq k \quad s_i : \mathbb{R} \rightarrow \{-1, 1\} \\ & v_1 \geq v_2 \geq \dots \geq v_k \geq 0 \end{aligned} \quad (8)$$

This is a non-convex optimization problem for which no general global optimal solution has been found to the best of our knowledge. In this paper we provide an efficient algorithm to optimally solve (8) for $k = 2$ in Section 3.2.2. We also propose a greedy algorithm for general k in Section 3.2.4.

Discussion: A general k -bits quantizer maps full precision values to an arbitrary set of 2^k numbers, not necessarily in the form of (6). The optimal quantization in this case can be computed using the Lloyd’s algorithm (Lloyd, 1982). While a general k -bits quantization has more representation power compared to k -bits scaled binary quantization, it does not allow an efficient implementation based on bitwise operations. Fixed-point representation (as opposed to floating point) is also in the form of (6) with an additional constant term. However, fixed-point quantization uniformly quantizes the space, therefore, it can be significantly inaccurate for small values of k .

3.2.1 FOLDABLE QUANTIZATION

In this section, we introduce a special family of k -bits scaled binary quantizations that allow fast computation of the quantized values. We name this family of quantizations **foldable**. A k -bits scaled binary quantization (6) is foldable if

$$s_i(x) = \text{sign}\left(x - \sum_{j=1}^{i-1} v_j s_j(x)\right) \quad \text{for } i = 1, 2, \dots, k \quad (9)$$

When the foldable condition is satisfied, given v_i ’s, we can compute the $s_i(\mathbf{x})$ ’s in (6) efficiently by applying the sign function. In this work, we show that the solution to (8) is foldable for $k = 1$ (see Appendix B) and $k = 2$ (see Appendix C).

3.2.2 OPTIMAL 2-BITS ALGORITHM

In this section, we present the optimal 2-bits binary quantization algorithm, the solution of (8) for $k = 2$.

The optimal 2-bits binary quantization is foldable and the scalars v_1 and v_2 should satisfy the following optimality conditions (for the proof see Appendix C):

$$v_1 = \frac{1}{2} (\mathbb{E}_{\mathbf{x} \sim p}[|\mathbf{x}| \mid |\mathbf{x}| > v_1] + \mathbb{E}_{\mathbf{x} \sim p}[|\mathbf{x}| \mid |\mathbf{x}| \leq v_1]) \quad (10)$$

$$v_2 = \frac{1}{2} (\mathbb{E}_{\mathbf{x} \sim p}[|\mathbf{x}| \mid |\mathbf{x}| > v_1] - \mathbb{E}_{\mathbf{x} \sim p}[|\mathbf{x}| \mid |\mathbf{x}| \leq v_1]) \quad (11)$$

In Figure 3 we visualize the conditional expectations that show up in (10) for a random variable x with standard normal distribution. The optimal v_1 lies on the intersection of the identity line and average of the conditional expectations in (10).

For a given vector $\mathbf{x} \in \mathbb{R}^N$ we can solve for v_1 in (10) efficiently. We substitute the conditional expectations in (10) by conditional average operators as their unbiased estimators. (10) implies that for the optimal v_1 , the average of the entries in $|\mathbf{x}|$ smaller than v_1 (an estimator of $\mathbb{E}_{\mathbf{x} \sim p}[|\mathbf{x}| \mid |\mathbf{x}| \leq v_1]$) and the average of the entries greater than v_1 (an estimator of $\mathbb{E}_{\mathbf{x} \sim p}[|\mathbf{x}| \mid |\mathbf{x}| > v_1]$) should be equidistant from v_1 . Note that (10) may have more than one solution, which are local minima of the objective function in (8). We find all the values that satisfy this condition in $\mathcal{O}(N \log N)$ time. We first sort entries of \mathbf{x} based on their absolute value and compute their cumulative sum. Then with one pass we can check whether (10) is satisfied for each element of \mathbf{x} . We evaluate the objective function in (8) for each local minima, and retain the best. After v_1 is calculated v_2 is simply computed from (11). As explained in Section 4, this process is only done during the training. In our experiments, finding the optimal 2-bits quantization increased the training time by 35% compared to the 2-bits greedy algorithm (see Section 3.2.4). Since the optimal 2-bits binary quantization is foldable, after recovering v_1 and v_2 , we have $s_1(\mathbf{x}) = \text{sign}(\mathbf{x})$ and $s_2(\mathbf{x}) = \text{sign}(\mathbf{x} - v_1 \text{sign}(\mathbf{x}))$.

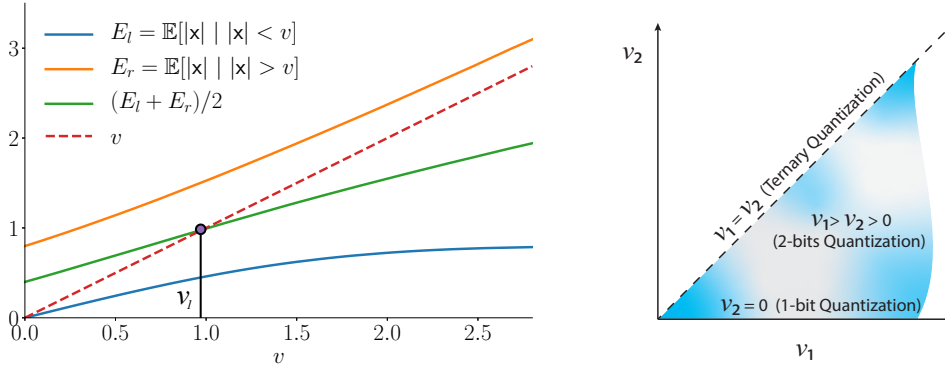


Figure 3: **Left:** The conditional expectations in (10) for a random variable x with standard normal distribution. The optimal value for 2-bits quantization is shown with a solid dot. **Right:** Optimization domain of (8) for $k=2$. The boundaries correspond to 1-bit and ternary quantizations.

3.2.3 OPTIMAL TERNARY ALGORITHM

The optimization domain of (8) for $k = 2$ is illustrated in Figure 3(right). The boundaries of the domain, $v_2 = 0$ and $v_1 = v_2 = v$, correspond to 1-bit binary and ternary (Li et al., 2016) quantizations, respectively. The scaled ternary quantization maps each full precision value x to $\{-2v, 0, 2v\}$. Ternary quantization needs 2-bits for representation. However, when a hardware with sparse calculation support is available, for example as in EIE (Han et al., 2016a), using ternary quantization can be more efficient compared to general 2-bits quantization. The optimal scaled ternary quantization is foldable and the scalar v should satisfy the following optimality condition (for the proof see Appendix D):

$$v = \frac{1}{2} \mathbb{E}_{\mathbf{x} \sim p}[|\mathbf{x}| \mid |\mathbf{x}| > v] \quad (12)$$

The process of solving for v in (12) is similar to that of solving for v_1 in (10) as described above.

3.2.4 k -BITS GREEDY ALGORITHM

In this section, we propose a greedy algorithm to compute k -bits scaled binary quantization, which we call Greedy Foldable (GF). It is given in Algorithm 1.

Algorithm 1: k -bits Greedy Foldable (GF) binary quantization: compute \mathbf{x}^q given \mathbf{x}

```

 $\mathbf{r} \leftarrow \mathbf{x}$ 
for  $i \leftarrow 1$  to  $k$  do
     $v_i \leftarrow \text{mean}(\text{abs}(\mathbf{r}))$ 
     $\mathbf{s}_i \leftarrow \text{sign}(\mathbf{r})$  // element-wise sign. For gradient of sign use STE.
     $\mathbf{r} \leftarrow \mathbf{r} - v_i \mathbf{s}_i$  // compute new residual.
end
return  $\mathbf{x} - \mathbf{r}$ 

```

In GF algorithm we compute a sequence of residuals. At each step, we greedily find the best s_i and v_i for the current residual using the optimal 1-bit binary quantization (5). Note that for $k = 1$ the GF is the same as the optimal 1-bit binary quantization.

Few of the other papers that have tackled the k -bits binary quantization to train quantized DNNs are as follows. In ReBNet (Ghasemzadeh et al., 2018), the authors proposed an algorithm similar to Algorithm 1, but considered v_i 's as trainable parameters to be learned by back-propagation. Lin et al. (2017) and Zhang et al. (2018a) find k -bits binary quantization via alternating optimization for s_i 's and v_i 's. Note that, all these methods produce sub-optimal solutions.

4 TRAINING BINARY NETWORKS

The loss functions in our quantized neural networks are non-differentiable due to the sign function in the quantizers. To address this challenge we use the training algorithm proposed in Courbariaux et al. (2015). To compute the gradient of the sign function we use the Straight Through Estimator (STE) (Bengio et al., 2013): $d/dx \text{sign}(x) = \mathbf{1}_{|x| \leq 1}$. During the training we keep the full precision weights and use Stochastic Gradient Descent (SGD) to gradually update them in back-propagation. In the forward-pass, only the quantized weights are used.

During the training we compute quantizers (for both weights and activations) using the online statistics, i.e., the scalars in a k -bits scaled binary quantization (6) are computed based on the observed values. During the training we also store the running average of these scalars. During inference we use the stored quantized scalars to improve the efficiency. This procedure is similar to the update of the batch normalization parameters in a standard DNN training (Ioffe & Szegedy, 2015).

5 EXPERIMENTS

We conduct experiments on the ImageNet dataset (Deng et al., 2009) using the ResNet-18 architecture (He et al., 2016). The details of the architecture and optimization parameters are provided in Appendix E.

We conduct three sets of experiments: (1) evaluate quantization error of activations of a pre-trained DNN, (2) evaluate the quantization error based on the classification accuracy of a post-training quantized network, and (3) evaluate the classification accuracy of during-training quantized networks. We report the quantization errors of the proposed binary quantization algorithms (optimal 1-bit, optimal 2-bits, optimal ternary, and the greedy foldable), and also compare their classification accuracies with the state-of-the-art algorithms BWN-Net (Rastegari et al., 2016), XNor-Net (Rastegari et al., 2016), TWN-Net (Li et al., 2016), DoReFa-Net (Zhou et al., 2016), ABC-Net (Lin et al., 2017), and LQ-Net (Zhang et al., 2018a).

5.1 QUANTIZATION ERROR OF ACTIVATIONS

To quantify the errors of the introduced binary quantization algorithms we adopt the analysis performed by Anderson & Berg (2017). They show that the angle between the full precision vector \mathbf{x} and its quantized version \mathbf{x}^q can be used as a measure of the accuracy of a quantization scheme. They prove that when $\mathbf{x}^q = \text{sign}(\mathbf{x})$, and the underlying distribution of elements of $\mathbf{x} \in \mathbb{R}^N$ is $p(x) = \mathcal{N}(0, 1)$, the angle between \mathbf{x} and \mathbf{x}^q converges to ~ 37 degrees for sufficiently large N .

Here we use the real data distribution. We trained a full precision network. We compute the activation tensors at each layer for a set of 128 images. In Figure 4 we show the angle between the full precision and quantized activations for different layers. When the optimal quantization is used, a significant reduction in the angle is observed compared to the greedy algorithm. The optimal 2-bits quantization is even better than the greedy 4-bits quantization for later layers of the network. Furthermore, the accuracy of the optimal quantization has less variance with respect to different input images and different layers of the network.

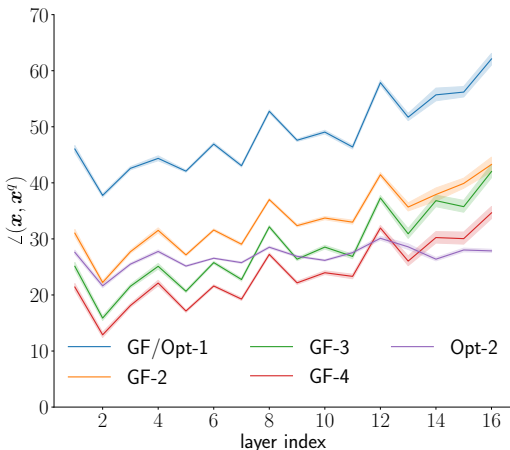


Figure 4: The angle between the full precision and the quantized activations for different layers of a trained full precision ResNet-18 architecture on ImageNet. The 95% confidence interval over different input images is shown.

Method	k^a	k^w	Top-1	Top-5
Post-GF	32	1	0.1	0.5
Post-GF	32	2	0.3	1.1
Post-GF	32	3	1.4	4.6
Post-GF	32	4	5.3	14.1
Post-Opt	32	2	5.3	13.9
Opt	1	1	54.3	77.3
GF	2	1	59.7	81.7
GF	3	1	61.1	82.7
GF	4	1	61.3	82.8
Opt	T	1	58.3	80.6
Opt	2	1	60.4	82.2
FP	32	32	69.6	89.2

Table 1: Validation accuracy of a quantized ResNet-18 trained on ImageNet. k^a and k^w are number of bits to quantize activations and weights, respectively. Post refers to post-training quantization. The training setup is the same for all quantized experiments. T, Opt, GF, and FP refer to ternary, optimal, Greedy Foldable, and full precision, respectively.

5.2 POST-TRAINING QUANTIZATION

In this section we apply post-training quantization to the weights of a pre-trained full precision network. We then use the quantized network for inference and report the classification accuracy. This procedure can result in an acceptable accuracy for a moderate number of bits (e.g., 16 or 8). However, the error significantly grows with a lower number of bits, which is the case in this experiment. Therefore, we only care about the relative differences between different quantization strategies. This experiment demonstrates the effect of quantization errors on the accuracy of the quantized DNNs.

The results are shown in the top half of Table 1. When the optimal 2-bits quantization is used, significant accuracy improvement (more than one order of magnitude) is observed compared to the greedy 2-bits quantization, which illustrate the effectiveness of the optimal quantization.

5.3 DURING-TRAINING QUANTIZATION

To achieve higher accuracy we apply quantization during the training, so that the model can adapt to the quantized weights and activations. In the bottom half of Table 1, we report the accuracies of the during-training quantized DNNs. We use 1-bit binary quantization for weights, and use different quantization algorithms for activations. When quantization is applied during-training, significantly higher accuracies are achieved. Similar to the previous experiments the optimal quantization algorithm achieves a better accuracy compared to the greedy.

In Table 2 we report results from the related works in which DNNs with quantized weights and/or activations are trained. All results correspond to training ResNet-18 architecture on the ImageNet dataset for the classification task. We report the mean and standard deviation of the model accuracy over 5 runs when our algorithms are used. Note that for 1-bit quantization the Greedy Foldable (GF) algorithm is the same with the optimal 1-bit binary quantization. In Opt* we used $2\times$ larger batch-size compared to Opt but with the same number of optimization steps. As shown in the Table 2 the proposed quantization algorithms match or improve the accuracies of the state-of-the-art binary neural networks.

Method	k^a	k^w	Val. top-1	Val. top-5
XNor-Net (Rastegari et al., 2016)	1	1	51.2	73.2
Opt	1	1	54.3 ± 0.2	77.4 ± 0.2
Opt	T	1	58.3 ± 0.1	80.6 ± 0.0
DoReFa-Net (Zhou et al., 2016) ^a	2	1	53.4	-
LQ-Net (Zhang et al., 2018a)	2	1	62.6	84.3
HWGQ-Net (Cai et al., 2017)	2	1	59.6	82.2
GF	2	1	59.7 ± 0.2	81.7 ± 0.1
Opt	2	1	60.4 ± 0.2	82.2 ± 0.1
Opt*	2	1	62.4 ± 0.1	83.6 ± 0.1
GF	3	1	61.1 ± 0.1	82.7 ± 0.1
ABC-Net (Lin et al., 2017)	3	3	61.0	83.2
DoReFa-Net (Zhou et al., 2016)	4	1	59.2	81.5
GF	4	1	61.3 ± 0.2	82.8 ± 0.1
BWN-Net (Rastegari et al., 2016)	32	1	60.8	83.0
Opt	32	1	64.2 ± 0.4	85.2 ± 0.1
TWN-Net (Li et al., 2016)	32	T	61.8	84.2
Opt	32	T	64.4 ± 0.1	85.4 ± 0.1
FP	32	32	69.6	89.2

Table 2: Comparison of state-of-the-art quantization methods for ResNet-18 architecture trained on the ImageNet dataset. Opt and GF are the proposed optimal and greedy foldable quantization algorithms, respectively. T and FP refer to ternary and full precision network, respectively.

^aThis result is taken from (Zhang et al., 2018a).

6 CONCLUSION

In this work, we analyze the accuracy of binary quantization to train DNNs with quantized weights and activations. We discuss methods to improve the accuracy of quantization, namely scaling and using more bits.

We introduce the rank- k binary quantization, as a general scaling scheme. Based on a singular value analysis we motivate using the scaled binary quantization, a computationally efficient scaling strategy. We also define a general k -bits scaled binary quantization. We provide provably optimal 2-bits and ternary quantizations. In addition, we propose a greedy k -bits quantization algorithm. We show results for post and during-training quantization, and demonstrate significant improvement in accuracy when optimal quantization is used. We compare the proposed quantization algorithms with state-of-the-art binary neural networks on the ImageNet dataset and show improved classification accuracies.

REFERENCES

- Alexander G Anderson and Cory P Berg. The high-dimensional geometry of binary neural networks. *arXiv preprint arXiv:1705.07199*, 2017.
- 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.
- Yunpeng Chen, Haoqi Fang, Bing Xu, Zhicheng Yan, Yannis Kalantidis, Marcus Rohrbach, Shuicheng Yan, and Jiashi Feng. Drop an octave: Reducing spatial redundancy in convolutional neural networks with octave convolution. *arXiv preprint arXiv:1904.05049*, 2019.

- 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*, 2018.
- François Chollet. Xception: Deep learning with depthwise separable convolutions. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 1251–1258, 2017.
- Matthieu Courbariaux, Yoshua Bengio, and Jean-Pierre David. Binaryconnect: Training deep neural networks with binary weights during propagations. In *Advances in neural information processing systems*, pp. 3123–3131, 2015.
- Matthieu Courbariaux, Itay Hubara, Daniel Soudry, Ran El-Yaniv, and Yoshua Bengio. Binarized neural networks: Training deep neural networks with weights and activations constrained to ± 1 or -1 . *arXiv preprint arXiv:1602.02830*, 2016.
- Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pp. 248–255. Ieee, 2009.
- Jonathan Frankle and Michael Carbin. The lottery ticket hypothesis: Finding sparse, trainable neural networks. *arXiv preprint arXiv:1803.03635*, 2018.
- Mohammad Ghasemzadeh, Mohammad Samragh, and Farinaz Koushanfar. Rebnet: Residual binarized neural network. In *2018 IEEE 26th Annual International Symposium on Field-Programmable Custom Computing Machines (FCCM)*, pp. 57–64. IEEE, 2018.
- Suyog Gupta, Ankur Agrawal, Kailash Gopalakrishnan, and Pritish Narayanan. Deep learning with limited numerical precision. In *International Conference on Machine Learning*, pp. 1737–1746, 2015.
- 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.
- Song Han, Xingyu Liu, Huizi Mao, Jing Pu, Ardavan Pedram, Mark A Horowitz, and William J Dally. Eie: efficient inference engine on compressed deep neural network. In *2016 ACM/IEEE 43rd Annual International Symposium on Computer Architecture (ISCA)*, pp. 243–254. IEEE, 2016a.
- Song Han, Jeff Pool, Sharan Narang, Huizi Mao, Enhao Gong, Shijian Tang, Erich Elsen, Peter Vajda, Manohar Paluri, John Tran, et al. Dsd: Dense-sparse-dense training for deep neural networks. *arXiv preprint arXiv:1607.04381*, 2016b.
- 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.
- Andrew G Howard, Menglong Zhu, Bo Chen, Dmitry Kalenichenko, Weijun Wang, Tobias Weyand, Marco Andreetto, and Hartwig Adam. Mobilenets: Efficient convolutional neural networks for mobile vision applications. *arXiv preprint arXiv:1704.04861*, 2017.
- Itay Hubara, Matthieu Courbariaux, Daniel Soudry, Ran El-Yaniv, and Yoshua Bengio. Quantized neural networks: Training neural networks with low precision weights and activations. *The Journal of Machine Learning Research*, 18(1):6869–6898, 2017.
- Forrest N Iandola, Song Han, Matthew W Moskewicz, Khalid Ashraf, William J Dally, and Kurt Keutzer. Squeezenet: Alexnet-level accuracy with 50x fewer parameters and 0.5 mb model size. *arXiv preprint arXiv:1602.07360*, 2016.
- Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. *arXiv preprint arXiv:1502.03167*, 2015.
- Benoit Jacob, Skirmantas Kligys, Bo Chen, Menglong Zhu, Matthew Tang, Andrew Howard, Hartwig Adam, and Dmitry Kalenichenko. Quantization and training of neural networks for efficient integer-arithmetic-only inference. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 2704–2713, 2018.

- Max Jaderberg, Andrea Vedaldi, and Andrew Zisserman. Speeding up convolutional neural networks with low rank expansions. *arXiv preprint arXiv:1405.3866*, 2014.
- Fengfu Li, Bo Zhang, and Bin Liu. Ternary weight networks. *arXiv preprint arXiv:1605.04711*, 2016.
- Xiaofan Lin, Cong Zhao, and Wei Pan. Towards accurate binary convolutional neural network. In *Advances in Neural Information Processing Systems*, pp. 345–353, 2017.
- Stuart Lloyd. Least squares quantization in pcm. *IEEE transactions on information theory*, 28(2): 129–137, 1982.
- Ningning Ma, Xiangyu Zhang, Hai-Tao Zheng, and Jian Sun. Shufflenet v2: Practical guidelines for efficient cnn architecture design. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pp. 116–131, 2018.
- Sachin Mehta, Mohammad Rastegari, Anat Caspi, Linda Shapiro, and Hannaneh Hajishirzi. Espnet: Efficient spatial pyramid of dilated convolutions for semantic segmentation. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pp. 552–568, 2018.
- Sachin Mehta, Mohammad Rastegari, Linda Shapiro, and Hannaneh Hajishirzi. Espnetv2: A light-weight, power efficient, and general purpose convolutional neural network. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 9190–9200, 2019.
- Asit Mishra, Eriko Nurvitadhi, Jeffrey J Cook, and Debbie Marr. Wprn: wide reduced-precision networks. *arXiv preprint arXiv:1709.01134*, 2017.
- Mohammad Rastegari, Vicente Ordonez, Joseph Redmon, and Ali Farhadi. Xnor-net: Imagenet classification using binary convolutional neural networks. In *European Conference on Computer Vision*, pp. 525–542. Springer, 2016.
- Mark Sandler, Andrew Howard, Menglong Zhu, Andrey Zhmoginov, and Liang-Chieh Chen. Mobilenetv2: Inverted residuals and linear bottlenecks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 4510–4520, 2018.
- Xavier Suau, Luca Zappella, and Nicholas Apostoloff. Network compression using correlation analysis of layer responses. 2018.
- Christian Szegedy, Wei Liu, Yangqing Jia, Pierre Sermanet, Scott Reed, Dragomir Anguelov, Dumitru Erhan, Vincent Vanhoucke, and Andrew Rabinovich. Going deeper with convolutions. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 1–9, 2015.
- Mingxing Tan and Quoc V Le. Efficientnet: Rethinking model scaling for convolutional neural networks. *arXiv preprint arXiv:1905.11946*, 2019.
- Sergey Zagoruyko and Nikos Komodakis. Wide residual networks. *arXiv preprint arXiv:1605.07146*, 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, 2018a.
- Xiangyu Zhang, Xinyu Zhou, Mengxiao Lin, and Jian Sun. Shufflenet: An extremely efficient convolutional neural network for mobile devices. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 6848–6856, 2018b.
- 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.
- Barret Zoph, Vijay Vasudevan, Jonathon Shlens, and Quoc V Le. Learning transferable architectures for scalable image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 8697–8710, 2018.

A OPTIMAL RANK-1 BINARY QUANTIZATION

In this section, we find the optimal rank-1 binary quantization of an m by n matrix \mathbf{X} that is discussed in Section 3:

$$\begin{aligned} & \underset{\mathbf{X}_1, \mathbf{S}}{\text{minimize}} && \|\mathbf{X} - \mathbf{X}_1 \odot \mathbf{S}\|_F \\ & \text{s.t.} && \mathbf{S} \in \{-1, 1\}^{m \times n} \\ & && \mathbf{X}_1 \in \mathbb{R}^{m \times n} \\ & && \text{rank}(\mathbf{X}_1) = 1 \end{aligned} \quad (13)$$

First, observe that the element-wise multiplication by -1 and $+1$ does not change the Frobenius norm. Therefore:

$$\min_{\mathbf{S}, \mathbf{X}_1} \|\mathbf{X} - \mathbf{X}_1 \odot \mathbf{S}\|_F = \min_{\mathbf{S}, \mathbf{X}_1} \|(\mathbf{X} - \mathbf{X}_1 \odot \mathbf{S}) \odot \mathbf{S}\|_F = \min_{\mathbf{S}, \mathbf{X}_1} \|\mathbf{X} \odot \mathbf{S} - \mathbf{X}_1\|_F \quad (14)$$

Furthermore, note that

$$\min_{\mathbf{S}, \mathbf{X}_1} \|\mathbf{X} \odot \mathbf{S} - \mathbf{X}_1\|_F^2 = \sigma_2^2(\mathbf{X} \odot \mathbf{S}) + \dots + \sigma_r^2(\mathbf{X} \odot \mathbf{S}) \quad (15)$$

Here $\sigma_i(\mathbf{X} \odot \mathbf{S})$ is the i 'th singular value of $\mathbf{X} \odot \mathbf{S}$ and r is its rank. In addition for any \mathbf{S} :

$$\sum_{i=1}^r \sigma_i^2(\mathbf{X} \odot \mathbf{S}) = \|\mathbf{X} \odot \mathbf{S}\|_F^2 = \|\mathbf{X}\|_F^2 \quad (16)$$

Hence, to minimize the sum in (15) we need to find an \mathbf{S} for which $\sigma_1^2(\mathbf{X} \odot \mathbf{S})$ is maximized:

$$\min_{\mathbf{S}, \mathbf{X}_1} \|\mathbf{X} \odot \mathbf{S} - \mathbf{X}_1\|_F^2 = \|\mathbf{X}\|_F^2 - \max_{\mathbf{S}} \sigma_1^2(\mathbf{X} \odot \mathbf{S}) \quad (17)$$

$\sigma_1(\mathbf{X} \odot \mathbf{S}) = \|\mathbf{X} \odot \mathbf{S}\|_2$ is the 2-norm of $\mathbf{X} \odot \mathbf{S}$. Therefore:

$$\max_{\mathbf{S}} \sigma_1^2(\mathbf{X} \odot \mathbf{S}) = \max_{\mathbf{S}} \max_{\|\mathbf{r}\|_2=1} \|(\mathbf{X} \odot \mathbf{S})\mathbf{r}\|_2^2 \quad (18)$$

For any \mathbf{S} and $\mathbf{r} \in \mathbb{R}^n$ we have $\|(\mathbf{X} \odot \mathbf{S})\mathbf{r}\|_2^2 \leq \|\mathbf{X}\|_2 \|\mathbf{r}\|_2^2$ since for $1 \leq i \leq m$ we have $|\sum_j S_{i,j} X_{i,j} r_j| \leq \sum_j |X_{i,j}| |r_j|$. Here $|\mathbf{X}| = \mathbf{X} \odot \text{sign}(\mathbf{X})$ is the element-wise absolute value of \mathbf{X} . Note that for $\mathbf{S} = \text{sign}(\mathbf{X})$ and \mathbf{r} with positive values the inequality becomes an equality. Therefore:

$$\max_{\mathbf{S}} \max_{\|\mathbf{r}\|_2=1} \|(\mathbf{X} \odot \mathbf{S})\mathbf{r}\|_2^2 = \max_{\|\mathbf{r}\|_2=1} \|\mathbf{X}\mathbf{r}\|_2^2 \quad (19)$$

Observe that the element-wise absolute value does not change the vector norm, i.e. $\|\mathbf{r}\|_2 = \|\mathbf{r}\|_2$, and hence $|\mathbf{r}|$ is a unit vector when \mathbf{r} is. Also for any \mathbf{r} we have $\|\mathbf{X}\mathbf{r}\|_2^2 \leq \|\mathbf{X}\|_2 \|\mathbf{r}\|_2^2$ since for $1 \leq i \leq m$ we have $|\sum_j |X_{i,j}| r_j| \leq \sum_j |X_{i,j}| |r_j|$. So we have

$$\max_{\|\mathbf{r}\|_2=1} \|\mathbf{X}\mathbf{r}\|_2^2 = \max_{\|\mathbf{r}\|_2=1} \|\mathbf{X}\mathbf{r}\|_2^2 = \sigma_1^2(|\mathbf{X}|) \quad (20)$$

Therefore, we showed that $\mathbf{S} = \text{sign}(\mathbf{X})$ and \mathbf{X}_1 equal to the best rank-1 approximation of $|\mathbf{X}|$ (i.e. the first term in its SVD) is a solution of (13).

B OPTIMAL SCALED BINARY QUANTIZATION

In this section, we solve the following optimization problem corresponding to the optimal 1-bit scaled binary quantization as discussed in Section 3.1:

$$\begin{aligned} & \underset{v, s}{\text{minimize}} && \int_{-\infty}^{+\infty} p(x)(vs(x) - x)^2 dx \\ & \text{s.t.} && s : \mathbb{R} \rightarrow \{-1, 1\} \\ & && v \in \mathbb{R}_{\geq 0} \end{aligned} \quad (21)$$

Here p is a probability distribution function. First, observe that:

$$\forall x \in \mathbb{R} : (v - x)^2 < (-v - x)^2 \quad \text{iff} \quad x > 0 \quad (22)$$

Therefore, the optimal choice for function s is $s(x) = \text{sign}(x)$. So we can rewrite (21) as follows:

$$\begin{aligned} & \underset{v}{\text{minimize}} && \int_{-\infty}^{+\infty} p(x)(v - |x|)^2 dx \\ & \text{s.t.} && v \in \mathbb{R}_{\geq 0} \end{aligned} \quad (23)$$

Setting the gradient of the objective function in (23) with respect to v to zero, we get:

$$v = \frac{\int_{-\infty}^{+\infty} p(x)|x|dx}{\int_{-\infty}^{+\infty} p(x)dx} = \mathbb{E}_{x \sim p}[|x|] \quad (24)$$

Hence, we showed that the optimal 1-bit scaled binary quantization maps x to $\mathbb{E}_{x \sim p}[|x|] \text{sign}(x)$.

C OPTIMAL 2-BITS BINARY QUANTIZATION

In this section, we solve the following optimization problem corresponding to the optimal 2-bits binary quantization as discussed in Section 3.2.2:

$$\begin{aligned} & \underset{v_1, v_2, s_1, s_2}{\text{minimize}} && \int_{-\infty}^{+\infty} p(x)(v_1 s_1(x) + v_2 s_2(x) - x)^2 dx \\ & \text{s.t.} && s_1, s_2 : \mathbb{R} \rightarrow \{-1, 1\} \\ & && v_1 \geq v_2 \geq 0 \end{aligned} \quad (25)$$

First, we show that the optimal 2-bits binary quantization is foldable, i.e., $\forall x \in \mathbb{R} s_1(x) = \text{sign}(x)$ and $s_2(x) = \text{sign}(x - v_1 s_1(x))$. Observe that

$$\begin{aligned} f(x) = (v_1 s_1(x) + v_2 s_2(x) - x)^2 &= v_1^2 \left(1 + \frac{v_2}{v_1} s_1(x) s_2(x) - \frac{s_1(x)x}{v_1} \right)^2 \geq \\ & v_1^2 \left(1 + \frac{v_2}{v_1} s_1(x) s_2(x) - \frac{|s_1(x)x|}{v_1} \right)^2 = g(x) \end{aligned} \quad (26)$$

The inequality in (26) holds because $v_1 \geq v_2$, and therefore, $1 + \frac{v_2}{v_1} s_1(x) s_2(x) \geq 0$. The objective function in (25) is a weighted average of $f(x)$ with non-negative weights. For $x \in \mathbb{R}$ the inequality is strict if $s_1(x) \neq \text{sign}(x)$. In that case, flipping the value of both $s_1(x)$ and $s_2(x)$ reduces $f(x)$ to a strictly smaller value $g(x)$. Hence, the optimal solution of (25) should satisfy $s_1(x) = \text{sign}(x)$ for all $x \in \mathbb{R}$.

For any v_1 and s_1 if we consider $y = x - v_1 s_1(x)$, the problem reduces to the 1-bit binary quantization for y . Based on the result showed in Appendix B for the optimal solution we have $s_2(x) = \text{sign}(y) = \text{sign}(x - v_1 s_1(x))$. This completes the proof to show that the optimal 2-bits binary quantization is foldable.

Next, we find the optimal values for v_1 and v_2 . If we substitute $s_1(x) = \text{sign}(x)$ and $s_2(x) = \text{sign}(x - v_1 s_1(x))$ in (25) we can decompose \mathbb{R} into four segments and write:

$$\begin{aligned} e(v_1, v_2) &= \int_{-\infty}^{+\infty} p(x)(v_1 s_1(x) + v_2 s_2(x) - x)^2 dx = \\ & \int_{-\infty}^{-v_1} p(x)(x + v_1 + v_2)^2 dx + \int_{-v_1}^0 p(x)(x + v_1 - v_2)^2 dx + \\ & \int_0^{v_1} p(x)(x - v_1 + v_2)^2 dx + \int_{v_1}^{+\infty} p(x)(x - v_1 - v_2)^2 dx = \\ & \int_0^{v_1} q(x)(x - v_1 + v_2)^2 dx + \int_{v_1}^{+\infty} q(x)(x - v_1 - v_2)^2 dx \end{aligned} \quad (27)$$

Here $e(v_1, v_2)$ is the error as a function of v_1 and v_2 , and $q(x) = p(-x) + p(x)$ is the folded distribution function. Assuming the optimal point occurs in the interior of the domain, it should

satisfy the zero gradient condition: $\partial e/\partial v_1 = \partial e/\partial v_2 = 0$. Taking derivative from (27) with respect to v_1 and v_2 and set it to zero we get:

$$\begin{aligned} v_1 &= \int_0^{v_1} xq(x)dx + \int_{v_1}^{+\infty} xq(x)dx + v_2 \left(\int_0^{v_1} q(x)dx - \int_{v_1}^{+\infty} q(x)dx \right) \\ v_2 &= - \int_0^{v_1} xq(x)dx + \int_{v_1}^{+\infty} xq(x)dx + v_1 \left(\int_0^{v_1} q(x)dx - \int_{v_1}^{+\infty} q(x)dx \right) \end{aligned} \quad (28)$$

We can simplify (28) and rewrite v_1 and v_2 in terms of the following conditional expectations:

$$\begin{aligned} v_1 &= \frac{1}{2} (\mathbb{E}_{x \sim p}[|x| \mid |x| > v_1] + \mathbb{E}_{x \sim p}[|x| \mid |x| \leq v_1]) \\ v_2 &= \frac{1}{2} (\mathbb{E}_{x \sim p}[|x| \mid |x| > v_1] - \mathbb{E}_{x \sim p}[|x| \mid |x| \leq v_1]) \end{aligned} \quad (29)$$

Hence, the optimal values of v_1 and v_2 can be obtained by solving for v_1 and v_2 in (29). The optimization domain has two boundaries. One is when $v_2 = 0$. This reduces the problem to 1-bit binary quantization. The optimal solution in that case is discussed in Appendix B. The other boundary is when $v_1 = v_2$. This results in the ternary quantization. The optimal solution in this case is discussed in Appendix D.

D OPTIMAL TERNARY QUANTIZATION

In this section, we find the optimal symmetric ternary quantization, that is to map a full precision value x to a discrete set $\{-2v, 0, 2v\}$ as discussed in Section 3.2.3. Finding the optimal mapping can be formulated as the following optimization problem:

$$\begin{aligned} \underset{v, s_1, s_2}{\text{minimize}} \quad & \int_{-\infty}^{+\infty} p(x) (vs_1(x) + vs_2(x) - x)^2 dx \\ \text{s.t.} \quad & s_1, s_2 : \mathbb{R} \rightarrow \{-1, 1\} \\ & v \geq 0 \end{aligned} \quad (30)$$

First, the above form is a special case of (25), hence, using the same argument as in Appendix C we can show that there is a foldable optimal solution: $s_1(x) = \text{sign}(x)$ and $s_2(x) = \text{sign}(x - vs_1(x))$. Then the total error as a function of v can be written as:

$$e(v) = \int_0^v q(x)x^2 dx + \int_v^{+\infty} q(x)(x - 2v)^2 dx \quad (31)$$

where $q(x)$ is the folded probability distribution function. Taking derivative from (31) and setting it to zero, we get:

$$v = \frac{\int_v^{+\infty} xq(x)dx}{2 \int_v^{+\infty} q(x)dx} = \frac{1}{2} \mathbb{E}_{x \sim p}[|x| \mid |x| > v] \quad (32)$$

E DETAILS OF TRAINING RESNET ON IMAGENET

In this section, we explain the details of how the DNN results reported in this paper are produced. All results correspond to the ResNet-18 architecture trained on the ImageNet dataset for the classification task. We use the standard training and validation splits of the ImageNet data set. We followed a similar architecture as XNor-Net (Rastegari et al., 2016). The convolutional block that we use is depicted in Figure 1(left). We use ReLU non-linearity before the batch normalization as suggested by (Rastegari et al., 2016). Also, we find it important to use bounded dynamic range, and therefore clip the values to $[-d, d]$. We use $d = 2, 3, 5$, and 8 for $k = 1, 2, 3$, and 4 bits quantizations, respectively. Similar to the other binary neural networks for the first and last layers we use full precision. Also, as suggested by Choi et al. (2018) we use full precision short-cuts in ResNet architecture, which adds a small computational/memory overhead. We quantize weights per filter and activations per layer. As Cai et al. (2017) we use first-order polynomial learning-rate annealing schedule (from 10^{-1} to 10^{-4}) and train for 120 epochs. We do not use weight decay. For the data augmentation we use the standard methods used to train full precision ResNet architecture. For training we apply random resize and crop to 224×224 , followed by random horizontal flipping, color jittering, and lightening. For test we resize the images to 256×256 followed by a center cropping to 224×224 .