

MONTÉ CARLO DEEP NEURAL NETWORK ARITHMETIC

Anonymous authors

Paper under double-blind review

ABSTRACT

Quantization is a crucial technique for achieving low-power, low latency and high throughput hardware implementations of Deep Neural Networks. Quantized floating point representations have received recent interest due to their hardware efficiency benefits and ability to represent a higher dynamic range than fixed point representations, leading to improvements in accuracy. We present a novel technique, Monte Carlo Deep Neural Network Arithmetic (MCA), for determining the sensitivity of Deep Neural Networks to quantization in floating point arithmetic. We do this by applying Monte Carlo Arithmetic to the inference computation and analyzing the relative standard deviation of the neural network loss. The method makes no assumptions regarding the underlying parameter distributions. We evaluate our method on pre-trained image classification models on the CIFAR10 and ImageNet datasets. For the same network topology and dataset, we demonstrate the ability to gain the equivalent of bits of precision by simply choosing weight parameter sets which demonstrate a lower loss of significance from the Monte Carlo trials. Additionally, we can apply MCA to compare the sensitivity of different network topologies to quantization effects.¹

1 INTRODUCTION

Deep Neural Networks have achieved state-of-the-art performances in many machine learning tasks such as speech recognition (Collobert et al. (2011)), machine translation (Bahdanau et al. (2014)), object detection (Ren et al. (2015)) and image classification (Krizhevsky et al. (2012)). However, excellent performance comes at the cost of significantly high computational and memory complexity, typically requiring TeraOps of computation during inference and Gigabytes of storage. To overcome these complexities, compression methods have been utilized, aiming to exploit the inherent resilience of DNNs to noise. These engender representations which maintain algorithm performance but significantly improve latency, throughput and power consumption of hardware implementations. In particular, exploiting reduced numerical precision for data representations through quantization has been emphatically promising, whereby on customizable hardware, efficiency scales quadratically with each bit of precision.

Quantization of fixed-point arithmetic (Q-FX) for DNN inference has been extensively studied, and more recently there has been increasing interest in quantized floating point (Q-FP) arithmetic for both DNN inference and training (Wang et al. (2018)). Q-FP has the advantage of higher dynamic range compared to equivalent Q-FX representations and reduced hardware cost over single-precision floating point (FP). This has influenced application specific integrated circuits (ASICs) such as Google’s tensor processing unit (TPU), which supports 16-bit FP and soft processors such as Microsoft’s Project Brainwave which utilizes 8-bit FP.

To illustrate these hardware benefits, we synthesized arithmetic logic units (ALUs) in different formats and different precisions on an FPGA and present performance estimates in operations per second (OPs) and area estimates in Look-up Tables (LUTs) per operation (LUTs/Op) in Figure 1. As shown, 7-bit FP has comparable performance and area to 2/8 bit Q-FX and improves significantly over both 12-bit Q-FX and 8/9-bit FP, outlining substantial performance and area benefits from reducing precision by only 1/2-bits. Thus, if we can design networks which not only achieve

¹Source code will be available if the paper is accepted

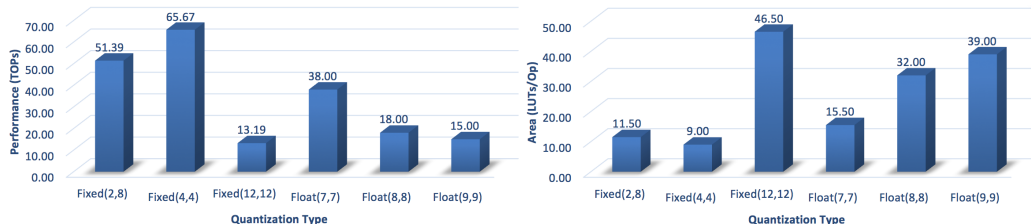


Figure 1: The estimated performance benefits in TeraOps (Left) and Area (Right) of different ALUs for computing a multiply-addition operation with (weight, activation) precision on a Xilinx VU9P FPGA operating at 500MHz.

high accuracy but are robust to quantization, architecting higher performing hardware solutions is possible.

Since in Q-FP, we are trying to represent the infinite set of real numbers using a finite number of bits, quantization and rounding artefacts will be introduced, with inaccuracies being cascaded along the computation graph (IEEE (1985); Goldberg (1991); Higham (2002)). This paper proposes a novel Monte Carlo Arithmetic (MCA) technique (Parker et al. (2000)) for determining the sensitivity of Deep Neural Networks to Q-FP representations. It allows hardware-software designers to quantify the impact of quantization, enabling more efficient systems to be discovered. We do this by exploiting Monte Carlo simulations under which rounding effects are randomized. This, in turn, allows one to infer the sensitivity of executing a computation graph to quantization effects.

Our MCA technique is highly sensitive, allowing very small differences in quantization behaviour to be detected. The technique makes no assumptions regarding data distributions and directly measures the effects of quantization on the problem under study. This allows us to provide insights into the precision requirements of any inference network for any given dataset. Additionally we can use the technique to select weight parameters which are more robust to floating point rounding. The theoretical and practical contributions of this work can be summarized as follows:

- We introduce a novel and rigorous analysis technique, Monte Carlo Deep Neural Network Arithmetic (MCA), which measures the sensitivity of Deep Neural Network inference computation to floating point rounding error.
- When applied to neural network inference, we show MCA can determine the precision requirements of different networks and is an exquisitely sensitive test which can detect small differences between different neural network topologies and weight sets.
- We demonstrate that while a network with the same topology but different weights may have the same loss and cross-validation error, their sensitivity to quantization can be vastly different. Using the CIFAR10 and ImageNet datasets, we introduce a method to choose weights which are more robust to rounding error, resulting in a greatly improved accuracy-area tradeoff over state-of-the-art methods.

It is worth noting that although we consider convolutional neural networks for image classification, this method could be applied for any neural network model architectures and applications. Moreover, while the experiments in this paper are limited to inference, it may be possible to apply the same idea to analyze training algorithms.

2 RELATED WORK

Low-precision representations of deep learning have been extensively studied. Many training methods have been developed for fixed point inference (Jacob et al. (2018); Faraone et al. (2018); Zhou et al. (2016)) and training (Wu et al. (2018b); Yang et al. (2019); Gupta et al. (2015); Sakr & Shanbhag (2019)). Other methods have also utilized Q-FP arithmetic for inference and training whilst maintaining full-precision accuracy. (Micikevicius et al. (2018)) implemented 16-bit FP arithmetic training whilst storing a 32-bit FP master copy for the weight updates. Additionally, (Wang et al. (2018); Mellempudi et al. (2019)) with 8-bit FP arithmetic whilst using a 16-bit FP copy of the

weights and 16/32-bits for the accumulator. Aside from these applications of low-precision arithmetic in deep learning, techniques for determining per-layer sensitivity to quantization have also been studied (Choi et al. (2016); Sakr & Shanbhag (2018)). Further, other studies have successfully determined the minimum fixed point precision requirements for a given DNN accuracy threshold (Sakr et al. (2017)). In FP arithmetic, rounding of inexact values to their nearest FP approximation has been studied in several publications (Higham (2002); Wilkinson (1994)). This has led to techniques for tracking information lost from finite precision arithmetic in various computational graphs using random perturbation such as Monte Carlo Arithmetic (Parker et al. (2000); Frechtling & Leong (2015)). To the best of our knowledge, our work is the first to present a technique for determining the sensitivity of DNNs to floating point rounding. These ideas can be very usefully applied to extending the limits of low-precision representations in deep learning applications.

3 BACKGROUND

In this section, we describe background theory upon which our technique for determining the sensitivity of DNNs to floating point rounding is based.

3.1 FLOATING POINT ARITHMETIC

The IEEE-754 binary floating point format (IEEE (1985)) represents real numbers x by a subset of the form:

$$x = (-1)^{s_x}(1 + m_x)2^{e_x} \quad (1)$$

where s is the sign bit, e is the base-2 exponent of x in binary floating point arithmetic and m_x is the mantissa of x . Such number formats can be described as a (s_x, e_x, m_x) tuple. In binary form the representation is $(b_s, b_{e_1}, b_{e_2}, \dots, b_{e_x}, b_{m_1}, b_{m_2}, \dots, b_{m_x})$. The infinite set of real numbers \mathbb{R} is represented in a computer with $B_x = s_x + e_x + m_x$ bits, and we define the finite set of real numbers representable in floating point format as *exact values*, $\mathbb{F} \subset \mathbb{R}$. Real numbers which aren't representable are rounded to their nearest exact value and we call this set of numbers *inexact values*, \mathbb{I} , where $\mathbb{I} \cup \mathbb{F} = \mathbb{R}$.

The approximation $\hat{x} = \mathbb{F}(x) = x(1 + \theta)$, given, $x \in \mathbb{I}$ introduces rounding error into the computation. The value of $\delta = \left\| \frac{x - \hat{x}}{x} \right\|$, represents the relative error which is a function of the machine hardware precision, p , as $\delta \leq \epsilon$, where $\epsilon = 2^{-p}$ (IEEE (1985); Goldberg (1991); Higham (2002)).

In general, inexactness can be caused by finite representations or errors propagating from earlier parts of the computation. Often the primary cause of error in floating point arithmetic is *catastrophic cancellation* which is typically the cause of horrific numerical inaccuracy from numerical analysis literature. Catastrophic cancellation occurs when for example, two near equal FP numbers, sharing k significant digits, are subtracted from one another as shown in (2). (Higham (2002)).

$$\begin{array}{l} 0. f_1 f_2 \dots f_k f_{(k+1)} \dots f_t \\ - 0. f_1 f_2 \dots f_k g_{(k+1)} \dots g_t \\ \hline = 0. 0 \ 0 \dots 0 \ h_{(k+1)} \dots h_t \end{array} \quad (2) \quad \begin{array}{l} 0. f_1 f_2 \dots f_k f_{(k+1)} \dots f_t r_{(t+1)} \dots r_p \\ - 0. f_1 f_2 \dots f_k g_{(k+1)} \dots g_t \hat{r}_{(t+1)} \dots \hat{r}_p \\ \hline = 0. 0 \ 0 \dots 0 \ h_{(k+1)} \dots h_t i_{(t+1)} \dots i_p \end{array} \quad (3)$$

In normalized form, the leading zeros are removed by shifting the result to the left and adjusting the exponent accordingly. The result is $0.h_{k+1} \dots h_t i_1 \dots i_k$ which has only $(t - k)$ accurate digits and digits i which are unknown. Additionally, the remaining accurate digits h are most likely affected by rounding error in previous computations. This can significantly magnify errors, especially in computing large computational graphs such as that of state-of-the-art DNNs.

If either operand in (2) is inexact, then the digits h are no more significant than any other sequence of digits. Yet, FP arithmetic has no mechanism of recording this loss of significance. By padding both our operands with random digits r and \hat{r} in (3), the resulting digits i are randomized. If k digits are lost in the result, then k random digits will be in the normalized result and when computed over many random trials, the results will disagree on the trailing k digits. In this case, we are able to detect catastrophic cancellation because the randomization over many trials provides a statistical simulation of round-off errors. We can use techniques from numerical analysis such as Monte Carlo methods to appropriately insert precision-dependant randomization in this way.

3.2 MONTE CARLO ARITHMETIC

Monte Carlo methods can be used to analyze rounding by representing inexact values as random variables (Parker et al. (2000); Frechtling & Leong (2015)). The real value x , as represented in (1), can be modelled to t digits, using:

$$\text{inexact}(x, t, \delta) = x + 2^{e_x - t} \delta = (-1)^{s_x} (m_x + 2^{-t} \delta) 2^{e_x} \quad (4)$$

where $\delta \in U(-\frac{1}{2}, \frac{1}{2})$ is a uniformly distributed random variable and t is a positive integer representing the *virtual precision* of concern. For the same input x in (4), we can run many Monte Carlo trials which will yield different values on each trial, where $0 < t < p$ so that the MCA can be run accurately on a computer with machine precision p . The ability to vary t is useful because it allows us to then evaluate the hardware precision requirements of a given system or computational graph for a given DNN.

MCA is a method to model the effect of rounding on a computational graph by randomizing all arithmetic operations. The randomization is applied for both generating inexact operands and also in rounding. In each operation using MCA, ideally both catastrophic cancellation and rounding error can be detected. An operation using MCA is defined as:

$$x \circ y = \text{round}(\text{inexact}(\text{inexact}(x) \circ \text{inexact}(y))) \quad (5)$$

where $\circ \in (+, -, \times, \div)$. By applying the inexact function to both operators we make it possible to detect catastrophic cancellation. Furthermore, applying the inexact function to the operation output and then rounding this value implements random rounding and hence is used to detect rounding error (Parker & Langley (1997)). Hence, for the same input into the system, each trial will yield different operands and output.

After applying MCA, we use random sampling to simulate Monte Carlo trials. The number of trials is an important consideration because, if insufficient, it can produce adverse effects on results. For each trial, we collect data on the resulting output of the system. This output data then constitutes a distribution on which we can apply statistical methods to understand its behaviour (Parker et al. (2000)). With a sufficient number of Monte Carlo trials and virtual precision t , the expected value of the output from these trials will equal the value from using real arithmetic. We can use this data to determine the total number of digits lost to rounding error and the minimum precision required to avoid a total loss of significance.

3.3 ANALYSIS

It has been noted in (Wilkinson (1994); Goldberg (1991)) that the relative error, δ is limited by $\delta \leq 2^{-p}$ for binary FP systems. With this definition of relative error, we can determine the expected number of significant binary digits available from a p -digit FP system as $p \geq -\log_2(\delta)$. These definitions can be adapted for MCA by replacing the precision of the FP system, p by the virtual precision t of an MCA operation. Thus, the relative error of an MCA operation is, for virtual precision t , is $\delta \leq 2^{-t}$ and the expected number of significant binary digits in a t -digit MCA operations is at least t . Using this definition and the proof provided in (Parker & Langley (1997)), the total significant binary digits in a set of MCA results is:

$$s' = \log_2 \frac{\mu}{\sigma} \quad (6)$$

where μ is the mean and σ the standard deviation of the MCA results. The output of the system should be some scalar value so that we can perform such analysis. The total number of base-2 significant digits lost in an MCA result set, K , is defined as:

$$K = t - s' = t - \log_2 \left(\frac{\mu}{\sigma} \right) = \log_2 \Theta + t, \quad (7)$$

where $\Theta = \frac{\sigma}{\mu}$ ($\mu \neq 0$) is the relative standard deviation (RSD) of the MCA results. When applied to DNN inference, K is a measure of the sensitivity of the network to floating point rounding. The method for implementing this is discussed in the section.

The virtual precision t controls the perturbation strength applied by the inexact function. For a given K , as we reduce t , the RSD should increase according to equation 7. At some point, an unexpected loss of significance (Frechtling & Leong (2015)) is encountered due to the nonlinear effects of quantization. The value at which this occurs is defined as t_{min} .

4 MONTE CARLO DEEP NEURAL NETWORK ARITHMETIC

We now describe a methodology for applying MCA techniques to DNN computation, allowing us to understand the sensitivity of a given network computational graph and/or its weight parameter representation.

4.1 NETWORK MODEL

We consider a generalized non-linear neural network with an output vector \mathbf{y} , input data vector \mathbf{x} and learnable weight parameter vector \mathbf{w} , whereby $\mathbf{y} = f(\mathbf{x}; \mathbf{w})$. To compute \mathbf{y} , several layers consisting of general matrix multiplication (GEMM) operations (such as convolutions and fully-connected layers) between the layer input \mathbf{x}_1 and weight parameters \mathbf{w}_1 ($\mathbf{x}_1 \otimes \mathbf{w}_1$), followed by a non-linear activation function to compute intermediate layer outputs \mathbf{y}_1 . The output of a given layer becomes the input to the subsequent layer. A loss function is defined as the objective function we want to minimize throughout training by updating \mathbf{w} at each mini-batch via an optimizer such as stochastic gradient descent. For a given set of input data X , the total network loss during inference is calculated by applying a loss function $L(f(\mathbf{x}; \mathbf{w}), \hat{\mathbf{y}}(\mathbf{x}))$ where $\hat{\mathbf{y}}(\mathbf{x})$ is the target ground truth output for \mathbf{x} . The total loss for X is then a scalar output, such that:

$$L(X; \mathbf{w}) = \frac{1}{|X|} \sum_{\mathbf{x} \in X} L(f(\mathbf{x}; \mathbf{w}), \hat{\mathbf{y}}(\mathbf{x})) \quad (8)$$

Naively applying MCA to each operation (fine-grained MCA) as described in (4) poses significant computational difficulties for DNN models. We observe two primary issues with employing fine-grained MCA to a DNN computational graph:

- Firstly, the number of required trials for Monte Carlo experiments to generate robust results can typically be in the hundreds or thousands. As DNN inference of state-of-the-art networks typically consist of billions of operations, the computational requirements of applying MCA after each operation will be very large, making the technique impractical.
- Using the accuracy as the system output for MCA experiments involves averaging classification results over many images. Averaging significantly reduces the standard deviation of the results, meaning results across different virtual precisions become indistinguishable and the standard deviation is potentially zero for high values of t .

4.2 MONTE CARLO NETWORK INFERENCE

To reduce the computational cost of Monte Carlo experiments, we employ a coarse-grained approach to MCA for GEMM operations. Conveniently, these can be naturally implemented in modern machine learning frameworks such as PyTorch. Furthermore, to minimize the amount of averaging in our results, the output of the loss function can be used as the system scalar output, rather than the accuracy. In this case, infinitesimal perturbations in layer operands are more likely to produce observable changes in the output. We can then apply (4) to the DNN inference computational problem in (8). The inexact equation for a given network layer operation with n inputs becomes:

$$\mathbf{y}_1 = \text{round}(\text{inexact}(\text{inexact}(\text{input}_1) \circ \text{inexact}(\text{input}_2) \circ \dots \text{inexact}(\text{input}_n))) \quad (9)$$

as opposed to (5), our operands in this case are vectors and \circ represents a neural network layer operation. For example, the output from performing a GEMM operation can thus be represented by:

$$\mathbf{y}_1 = \text{round}(\text{inexact}(\text{inexact}(\mathbf{x}_1) \otimes \text{inexact}(\mathbf{w}_1))) \quad (10)$$

This applies the inexact function to each operand and then performs the full GEMM operation (multiplies and additions). The inexact function is then applied to all output values of the operation and then rounded. The form in (9) is applied to each edge of neural network computational graph where multiply (division) and/or add (subtract) operations are performed. Hence, it is not applied to operations such as MaxPool and ReLU. As an example, in Figure 2 we show where the inexact function is applied for a residual block with folded batch normalization, which is a repeating sequence of layers found in ResNet models (He et al. (2015)). At the final output of the network, the loss is computed with (8) and the inexact function is applied to the outputs \mathbf{y} and also the loss output scalar value

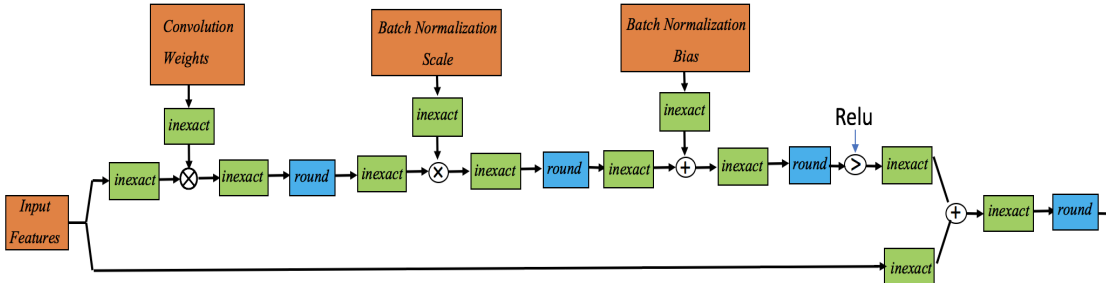


Figure 2: Applying MCA to residual blocks found in ResNet. For each operation, the inexact function is applied to each input operand and both the inexact function and baes-2 rounding is applied to the resulting output.

L. From analyzing the behavior of the loss, we infer the sensitivity of the accuracy of the system to floating point rounding. By using coarse-grained MCA for the GEMM operations, we won’t be able to detect every single instance of catastrophic cancellation and random rounding. However, we significantly reduce execution time over fine-grained MCA and show in the next section that we can still retrieve valuable information about our system. In fact, for one trial with one batch of 32 images on ImageNet running on a Nvidia Titan Xp GPU, the speed up of regular inference without MCA is only $1.05\times$. We also note that fine-grained MCA could be applied with a simple modification and would be possible given appropriate customized hardware support for parallel Monte Carlo computations such as in (Yeung et al. (2011)).

5 EXPERIMENTAL RESULTS

In this section, we present experimental results for applying MCA to exemplary convolutional neural networks. We use the CIFAR10 and ImageNet image classification datasets to compare MobileNet-v2 (Sandler et al. (2018)), EfficientNet (Tan & Le (2019)), AlexNet, ResNet (He et al. (2015)), SqueezeNet(Iandola et al. (2016)) and MnasNet (Tan et al. (2018)). For CIFAR-10, we use a batch size of 128, whilst we use a batch size of 32 for ImageNet experiments. Cross-entropy is used as the loss function for both datasets. For a given network, dataset and weight representation, we run 1000 Monte Carlo trials and then compute the relative standard deviation of the network loss/accuracy for different values of t , using methods discussed in Section 3.3. The network loss output from each trial is computed with a single batch of images from the training dataset, showing the ability to achieve these results without requiring validation/test data. When reporting Q-FP results, we use the quantization function from (Wang et al. (2018)) with stochastic rounding².

5.1 DISTINGUISHING WEIGHT PARAMETERS REPRESENTATIONS

As discussed in Section 3.1, the inexactness in FP arithmetic largely depends on the numerical value of operands. Thus, two instances of the same network and dataset which have the same accuracy, but vastly different weight representations, will likely produce differing sensitivities to FP rounding. We can measure this using MCA, by calculating their loss of significance values, K . We first train 8 instances of EfficientNet-b0 and MobileNet-V2 on the CIFAR10 dataset from scratch with random initialization from (Glorot & Bengio (2010)), all achieving within 1% accuracy of one another.

We then test their percentage accuracy decrease from using post-training quantization (i.e. no fine-tuning) with varying Q-FP precisions. Evidently in Figure 3, we see that the models with higher K values typically experience a larger drop in Q-FP accuracy, indicating they are more sensitive to floating point rounding error. Notably, the model with lowest K for 7-bit MobileNet-v2 experiences a lower percentage accuracy drop than three of the 8-bit models. In this case, MCA model selection enables the saving of a bit of precision while achieving smaller accuracy decrease than some of the trained 8-bit Q-FP models.

²<https://github.com/Tiiiger/QPyTorch>

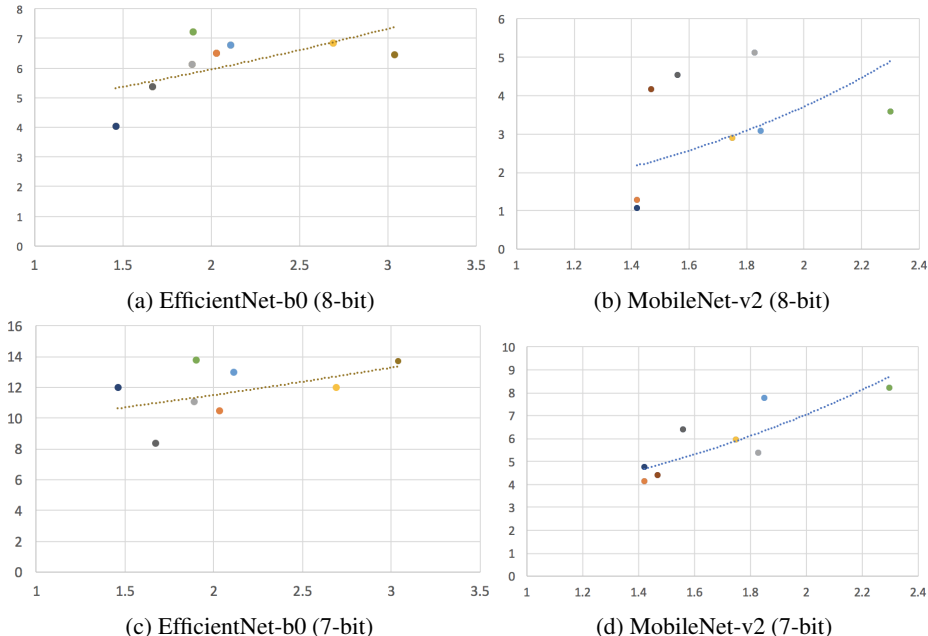


Figure 3: Accuracy decrease as a percentage of single-precision accuracies (y-axis) vs K (x-axis) for 8 different trained models when using post-training Q-FP representations on the CIFAR-10 dataset.

Table 1: Post-training Q-FP quantization accuracy when using K for model selection rather than cross-validation accuracy on the full CIFAR10 validation set

Precision (Weight, Act.)	MobileNet-v2		EfficientNet	
	MCA ($K = 1.42$)	(Wang et al. (2018)) ($K = 1.56$)	MCA ($K = 1.46$)	(Wang et al. (2018)) ($K = 1.89$)
(32,32)	89.6	90.3	93.7	94.6
(8,8)	88.5	86.7	90.0	89.0
(6,6)	79.0	76.6	64.6	56.0
(5,5)	48.4	40.8	22.0	18.6

5.2 COMPARISON TO PREVIOUS WORK

One practical use case from the insights gained by MCA is model selection for quantization. Typically when quantizing a given model trained on a given dataset, cross validation is used to choose the model with the highest accuracy and the sensitivity to quantization is assumed to be the same across models. As discussed, for Q-FP representations this is not necessarily the case. We can use K from MCA to predict which models will be more robust to quantization. To demonstrate this, in Table 1 we compare post-training quantization results for model selection based on K from MCA vs CIFAR10 accuracy. Evidently, even though the full-precision accuracy is as much as 0.9% higher initially, after quantizing the network to 8-5 bits, the accuracy of the network chosen by smallest K , always has significantly higher accuracy.

5.3 NETWORK COMPARISON

Modern DNNs consist of convolutional blocks with highly varying computational graphs (Wu et al. (2018a); Howard et al. (2017)). Using MCA we can also compute and compare their sensitivities to floating point rounding error to determine which networks will be robust to Q-FP representations. For a given virtual precision t , we calculate the RSD from the MCA results of pre-trained models

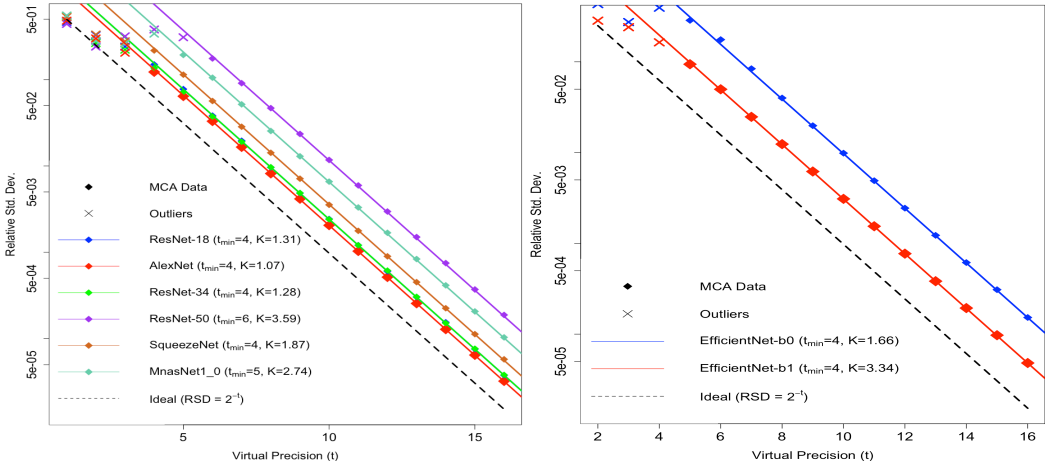


Figure 4: Comparison of the RSD of the loss for various networks (Right) and EfficientNet variants (Left) at different virtual precisions on the ImageNet dataset.

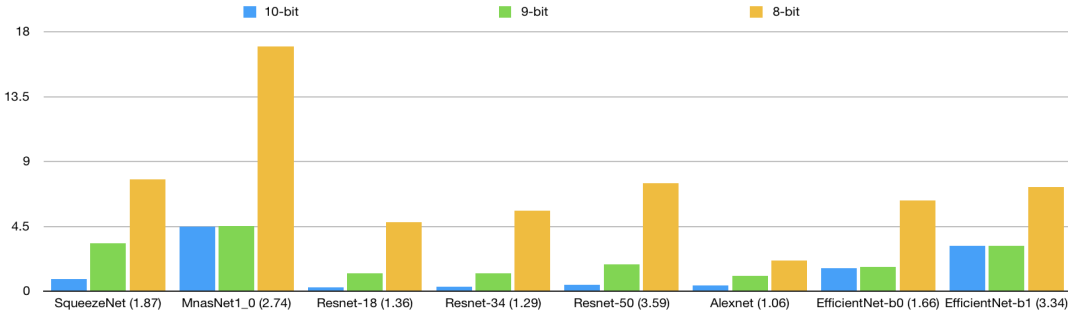


Figure 5: Q-FP Percentage accuracy decrease for different ImageNet models

trained on the ImageNet dataset from PyTorch³ ⁴. From here we can then assign a K value to each network and compare their loss of significance. In Figure 4 we show the RSD of each network differing values of t and run linear regression analysis over our data points. The distance of the regression lines to the ideal line represents the values of K , as described in (7). From the MCA results, AlexNet is the least sensitive to rounding and ResNet-50 is the most, with various models in between these two. Additionally we compare EfficientNet at two different model scales and evidently the larger model has much larger sensitivity. We then also compare the accuracy percentage decrease of all models at 10, 9 and 8-bit post-training Q-FP in Figure 5. At 8-bit Q-FP, besides MnasNet which experiences a large accuracy drop, K is able to predict accuracy degradation amongst networks. This implies MCA provides very valuable information about Q-FP model design.

6 CONCLUSION

We present a novel, extremely sensitive, technique to quantify rounding error in DNNs. This is the first method to successfully compare the sensitivity of networks to floating point rounding error. Ultimately, this technique provides a tool for enabling the design of networks which perform higher when quantized. We do this by applying concepts from Monte Carlo Arithmetic theory to DNN computation. Furthermore, we show that by calculating the loss of significance metric K from MCA, on the CIFAR10 and ImageNet datasets, we can compare network sensitivities to floating point rounding error and gain valuable insights to potentially design better neural networks. This is an important contribution due to the increasing interest in low-precision floating point arithmetic for efficient DNN hardware systems. The theoretical and practical contributions of this paper will likely translate well to analyzing floating point rounding in backpropagation in future work.

³<https://github.com/pytorch/vision/tree/master/torchvision>

⁴<https://github.com/rwightman/pytorch-image-models>

REFERENCES

- Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. Neural machine translation by jointly learning to align and translate, 2014. URL <http://arxiv.org/abs/1409.0473>. cite arxiv:1409.0473Comment: Accepted at ICLR 2015 as oral presentation.
- Yoojin Choi, Mostafa El-Khamy, and Jungwon Lee. Towards the limit of network quantization. *ArXiv*, abs/1612.01543, 2016.
- Ronan Collobert, Jason Weston, Léon Bottou, Michael Karlen, Koray Kavukcuoglu, and Pavel Kuksa. Natural language processing (almost) from scratch. *J. Mach. Learn. Res.*, 12:2493–2537, November 2011. ISSN 1532-4435. URL <http://dl.acm.org/citation.cfm?id=1953048.2078186>.
- Julian Faraone, Nicholas J. Fraser, Michaela Blott, and Philip H. W. Leong. SYQ: learning symmetric quantization for efficient deep neural networks. In *2018 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2018, Salt Lake City, UT, USA, June 18-22, 2018*, pp. 4300–4309, 2018. doi: 10.1109/CVPR.2018.00452. URL http://openaccess.thecvf.com/content_cvpr_2018/html/Faraone_SYQ_Learning_Symmetric_CVPR_2018_paper.html.
- Michael Frechtling and Philip H. W. Leong. Mcalib: Measuring sensitivity to rounding error with monte carlo programming. *ACM Trans. Program. Lang. Syst.*, 37(2):5:1–5:25, April 2015. ISSN 0164-0925. doi: 10.1145/2665073. URL <http://doi.acm.org/10.1145/2665073>.
- Xavier Glorot and Yoshua Bengio. Understanding the difficulty of training deep feedforward neural networks. In *In Proceedings of the International Conference on Artificial Intelligence and Statistics (AISTATS10)*. Society for Artificial Intelligence and Statistics, 2010.
- David Goldberg. What every computer scientist should know about floating-point arithmetic. *ACM Comput. Surv.*, 23(1):5–48, March 1991. ISSN 0360-0300. doi: 10.1145/103162.103163. URL <http://doi.acm.org/10.1145/103162.103163>.
- Suyog Gupta, Ankur Agrawal, Kailash Gopalakrishnan, and Pritish Narayanan. Deep learning with limited numerical precision. In *Proceedings of the 32Nd International Conference on International Conference on Machine Learning - Volume 37, ICML’15*, pp. 1737–1746. JMLR.org, 2015. URL <http://dl.acm.org/citation.cfm?id=3045118.3045303>.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 770–778, 2015.
- Nicholas J. Higham. *Accuracy and Stability of Numerical Algorithms*. Society for Industrial and Applied Mathematics, Philadelphia, PA, USA, 2nd edition, 2002. ISBN 0898715210.
- 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*, abs/1704.04861, 2017.
- Forrest N. Iandola, Matthew W. Moskewicz, Khalid Ashraf, Song Han, William J. Dally, and Kurt Keutzer. Squeezenet: Alexnet-level accuracy with 50x fewer parameters and <1mb model size. *CoRR*, abs/1602.07360, 2016. URL <http://arxiv.org/abs/1602.07360>.
- IEEE. *IEEE standard for binary floating-point arithmetic*. Institute of Electrical and Electronics Engineers, New York, 1985. Note: Standard 754–1985.
- 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 *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2018.

Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. In F. Pereira, C. J. C. Burges, L. Bottou, and K. Q. Weinberger (eds.), *Advances in Neural Information Processing Systems 25*, pp. 1097–1105. Curran Associates, Inc., 2012. URL <http://papers.nips.cc/paper/4824-imagenet-classification-with-deep-convolutional-neural-networks.pdf>.

Naveen Mellempudi, Sudarshan Srinivasan, Dipankar Das, and Bharat Kaul. Mixed precision training with 8-bit floating point. *CoRR*, abs/1905.12334, 2019. URL <http://arxiv.org/abs/1905.12334>.

Paulius Micikevicius, Sharan Narang, Jonah Alben, Gregory Diamos, Erich Elsen, David Garcia, Boris Ginsburg, Michael Houston, Oleksii Kuchaiev, Ganesh Venkatesh, and Hao Wu. Mixed precision training. In *International Conference on Learning Representations*, 2018. URL <https://openreview.net/forum?id=rlgs9JgRZ>.

Douglas Stott Parker and David Langley. Monte carlo arithmetic: exploiting randomness in floating-point arithmetic. Computer Science Department, University Of California, 1997.

Douglas Stott Parker, Brad Pierce, and Paul R. Eggert. Monte carlo arithmetic: how to gamble with floating point and win. *Computing in Science and Engineering*, 2:58–68, 2000.

Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster r-cnn: Towards real-time object detection with region proposal networks. In C. Cortes, N. D. Lawrence, D. D. Lee, M. Sugiyama, and R. Garnett (eds.), *Advances in Neural Information Processing Systems 28*, pp. 91–99. Curran Associates, Inc., 2015. URL <http://papers.nips.cc/paper/5638-faster-r-cnn-towards-real-time-object-detection-with-region-proposal-networks.pdf>.

Charbel Sakr and Naresh Shanbhag. An analytical method to determine minimum per-layer precision of deep neural networks. In *2018 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 1090–1094. IEEE, 2018.

Charbel Sakr and Naresh Shanbhag. Per-tensor fixed-point quantization of the back-propagation algorithm. In *International Conference on Learning Representations*, 2019. URL <https://openreview.net/forum?id=rkxanJ9Ym>.

Charbel Sakr, Yongjune Kim, and Naresh Shanbhag. Analytical guarantees on numerical precision of deep neural networks. In *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, pp. 3007–3016. JMLR. org, 2017.

Mark Sandler, Andrew G. Howard, Menglong Zhu, Andrey Zhmoginov, and Liang-Chieh Chen. Mobilenetv2: Inverted residuals and linear bottlenecks. *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 4510–4520, 2018.

Mingxing Tan and Quoc Le. EfficientNet: Rethinking model scaling for convolutional neural networks. In Kamalika Chaudhuri and Ruslan Salakhutdinov (eds.), *Proceedings of the 36th International Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning Research*, pp. 6105–6114, Long Beach, California, USA, 09–15 Jun 2019. PMLR. URL <http://proceedings.mlr.press/v97/tan19a.html>.

Mingxing Tan, Bo Chen, Ruoming Pang, Vijay Vasudevan, and Quoc V. Le. Mnasnet: Platform-aware neural architecture search for mobile. *ArXiv*, abs/1807.11626, 2018.

Naigang Wang, Jungwook Choi, Daniel Brand, Chia-Yu Chen, and Kailash Gopalakrishnan. Training deep neural networks with 8-bit floating point numbers. In *Proceedings of the 32Nd International Conference on Neural Information Processing Systems, NIPS’18*, pp. 7686–7695, USA, 2018. Curran Associates Inc. URL <http://dl.acm.org/citation.cfm?id=3327757.3327866>.

James H. Wilkinson. *Rounding Errors in Algebraic Processes*. Dover Publications, Inc., New York, NY, USA, 1994. ISBN 0486679993.

Bichen Wu, Alvin Wan, Xiangyu Yue, Peter Jin, Sicheng Zhao, Noah Golmant, Amir Gholaminejad, Joseph Gonzalez, and Kurt Keutzer. Shift: A zero flop, zero parameter alternative to spatial convolutions. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2018a.

Shuang Wu, Guoqi Li, Feng Chen, and Luping Shi. Training and inference with integers in deep neural networks. In *International Conference on Learning Representations*, 2018b. URL <https://openreview.net/forum?id=HJGXzmspb>.

Guandao Yang, Tianyi Zhang, Polina Kirichenko, Junwen Bai, Andrew Gordon Wilson, and Christopher De Sa. Swalp : Stochastic weight averaging in low precision training. In *ICML*, 2019.

Jackson H.C. Yeung, Evangeline F.Y. Young, and Philip H.W. Leong. A monte-carlo floating-point unit for self-validating arithmetic. In *Proceedings of the 19th ACM/SIGDA International Symposium on Field Programmable Gate Arrays, FPGA '11*, pp. 199–208, New York, NY, USA, 2011. ACM. ISBN 978-1-4503-0554-9. doi: 10.1145/1950413.1950453. URL <http://doi.acm.org/10.1145/1950413.1950453>.

Shuchang Zhou, Zekun Ni, Xinyu Zhou, He Wen, Yuxin Wu, and Yuheng Zou. Dorefa-net: Training low bitwidth convolutional neural networks with low bitwidth gradients. *ArXiv*, abs/1606.06160, 2016.

A APPENDIX

You may include other additional sections here.