DeepShift: Towards Multiplication-Less Neural Networks

Mostafa Elhoushi¹, Zihao Chen¹, Farhan Shafiq¹, Ye Henry Tian¹, and Joey Yiwei Li¹

¹Huawei Technologies , Toronto Hetrogeneous Compilers Lab , 19 Allstate Parkway, Markham, ON Canada

Abstract

Deployment of convolutional neural networks (CNNs) in mobile environments, their high computation and power budgets proves to be a major bottleneck. Convolution layers and fully connected layers, because of their intense use of multiplications, are the dominant contributer to this computation budget. This paper proposes to tackle this problem by introducing two new operations: convolutional shifts and fully-connected shifts, that replace multiplications all together with bitwise shift and sign flipping instead. For inference, both approaches may require only 6 bits to represent the weights. This family of neural network architectures (that use convolutional shifts and fully-connected shifts) are referred to as DeepShift models. We propose two methods to train DeepShift models: DeepShift-Q that trains regular weights constrained to powers of 2, and DeepShift-PS that trains the values of the shifts and sign flips directly. Training the DeepShift versions of ResNet18 architecture from scratch, we obtained accuracies of 92.33% on CIFAR10 dataset, and Top-1/Top-5 accuracies of 65.63%/86.33% on Imagenet dataset. Training the DeepShift version of VGG16 on ImageNet from scratch, resulted in a drop of less than 0.3% in Top-5 accuracy. Converting the pre-trained 32bit floating point baseline model of GoogleNet to DeepShift and training it for 3 epochs, resulted in a Top-1/Top-5 accuracies of 69.87%/89.62% that are actually higher than that of the original model. Further testing is made on various well-known CNN architectures. Last but not least, we implemented the convolutional shifts and fully-connected shift GPU kernels and showed a reduction in latency time of 25% when inferring ResNet18 compared to an unoptimized multiplication-based GPU kernels. The code is available online at https://github.com/mostafaelhoushi/DeepShift.

1. Introduction

Deep Neural Networks are increasingly being targeted for mobile and IoT applications. Devices at the edge have a lower power and price budget as well as constrained memory size. Moreover, the amount of communication between memory and compute also has a major role in the power requirements of a CNN. Moreover, if communication between device and cloud become necessary (e.g. in case of model updates etc), model size could affect the connectivity costs. Therefore, for mobile / IoT inference applications, model optimization, size reduction, faster inference and lower power consumption are key areas of research.

Several approaches are being considered to address this need and as such these efforts can be divided into a few categories. One approach is to build efficient models from the ground up resulting in novel network architectures, however that proves to be a task requiring a lot of training resource to try multiple variants of an architecture to find the best fit. Another approach is to start with a big model initially and prune it. Since among the many parameters in the network, some are redundant and dont contribute a lot to the output, hence a ranking is assigned to each parameters based on contribution to the output. Low ranking parameters can be done away with (pruned), without effecting the accuracy too much. Another technique is to start with a big model and reduce the model size by applying quantization to smaller bit-width floating or fixed-point numbers. In some cases the quantized models are retrained to regain some of the accuracy. Key attractions of these technique are that they can be easily applied to various kinds of networks and they not only reduce model size but also require less complex compute units on the underlying hardware. This results in smaller model footprint, less working memory (and cache), faster computation on supporting platforms and lower power consumption. Moreover, some optimization techniques replace multiplication with cheaper operations such as binary XNOR operations or bitwise shifts. Binary XNOR techniques may have high accuracy on small datasets such as MNIST or CIFAR10, but suffer high degradation on complex datasets such as Imagenet. On the other hand, bit-wise shift techniques have proven to show minimal drop in accuracy on Imagenet.

This paper presents an approach to reduce computation and power budget of CNNs by replacing regular multiplication-based convolution and linear operations (a.k.a fully-connected layer or matrix multiplication) with bitwise-shift-based convolution and linear operations respectively. Applying bitwise shift operation on an element is mathematically equivalent to multiplying it by a power of 2. We compare our work with 4 other similar approaches that replace multiplications with bitwise shift and sign flips: INQ [18], LogNN [13], and LogQuant [2]. We propose two approaches: DeepShift-Q and DeepShift-PS. Both our DeepShift approaches are based on training the shift values (i.e. the powers of 2) and sign flips either indirectly or directly, while INQ, LogNN, and LogQuant are based on rounding 32-bit floating-point weights to powers of 2. To the best of our knowledge, this paper is the first to propose training the shift values directly, which is the main novelty of the paper.

This paper makes the following contributions:

- We introduce 2 different methods to train DeepShift networks, neural networks with powers of 2 weights: each method trains the powers of 2 at run-time, and when computing the parameters gradients at train-time (see Section 2).
- We conduct experiments for each training method, which show that it is possible to train DeepShift networks on MNIST, CIFAR-10 and Imagenet, from scratch or from a pre-trained full-precision weights, and achieve nearly state-of-the-art results (see Section 4).
- We developed a bitwise-shift-based matrix multiplication and convolution GPU kernel with which it is possible to run the DeepShift ResNet18 architecture with 25% faster inference times than the baseline ResNet18 model with an unoptimized GPU kernel (see Section 5.1).

The code for training and running our DeepShift networks using both methods is available online at https://github.com/mostafaelhoushi/DeepShift.

2. DeepShift Networks

As shown in Figure 1, the main concept of this paper is to replace multiplication with bitwise shift and sign flip. If the underlying binary representation of an input number, x

is in integer or fixed-point format, a bit-wise shift of \tilde{p} bits, where \tilde{p} is an integer, to the left (or right) is mathematically equivalent to multiplying by a positive (or negative power) of 2:

$$2^{\tilde{p}}x = \begin{cases} x << \tilde{p} & \text{if } \tilde{p} > 0\\ x >> \tilde{p} & \text{if } \tilde{p} < 0 \quad s.t.\tilde{p} \in \mathbb{Z}\\ x & \text{if } \tilde{p} = 0 \end{cases}$$
 (1)

For simplicity, we will use $<<\tilde{p}$ to implicitly denote a left shift if \tilde{p} is positive and right shift if \tilde{p} is negative.

Bitwise shift can only be equivalent to multiplying by a positive number, because $2^{\pm \tilde{p}} > 0$ for any real value of \tilde{p} . However, in neural networks, it is necessary for the training to have the equivalent of multiplying by negative numbers in its search space, especially in convolutional neural networks where filters with both positive and negative values contribute to detecting edges. Therefore, we also need to use the sign flip (a.k.a negation) operation. We use a ternary sign operator, i.e., an operator that can either leave its operand unchanged, replace it with 0, or changes the sign of its operand. The sign flip operation can be represented mathematically as:

$$\tilde{s}x = \begin{cases} -x & \text{if } \tilde{s} = -1\\ 0 & \text{if } \tilde{s} = 0\\ x & \text{if } \tilde{s} = +1 \end{cases} \quad s.t.\tilde{s} \in \{-1, 0, +1\} \quad (2)$$

Similar to bitwise shift, sign flip is a computationally cheap operation too as it involves returning the 2's complement of a number.

We introduce novel operators, LinearShift and ConvShift, that, in the forward pass, replace multiplication with bitwise shift and sign flip. Hence, the weight matrix, W, whether it is used for linear operation (i.e., Y=WX+b) or convolution operation (i.e., $Y=W\circledast X+b$), will be replaced by elementwise product of the shift matrix \tilde{P} and sign matrix \tilde{S} .

$$W \to \tilde{S} \cdot 2^{\tilde{P}}$$
 (3)

In the following sub-sections we shall present two different methods to train the LinearShift and ConvShift operators of our DeepShift models: DeepShift-Q and DeepShift-PS. Both approaches use Equation 3 for the forward pass but differ in the backward pass and differ in how \tilde{P} and \tilde{S} are derived. The two approaches are summarized in Figure 2.

2.1. DeepShift-Q

The DeepShift-Q approach, training is similar to the regular training approach of linear or convolution operations with weight W, but in the forward pass and backward pass the weight matrix is quantized to \tilde{W}_q by rounding it to the

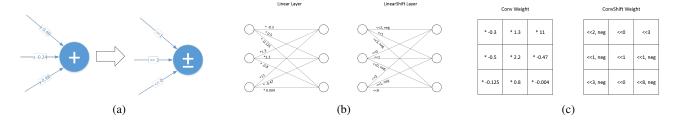


Figure 1: (a) Multiply and Accumulate (MAC), which is the main operation in most neural networks, is replaced in DeepShift with bitwise shift and addition/subtraction. (b) Weights of original linear operator vs. proposed shift linear operator. (c) Weights of original convolution operator vs. proposed shift convolution operator

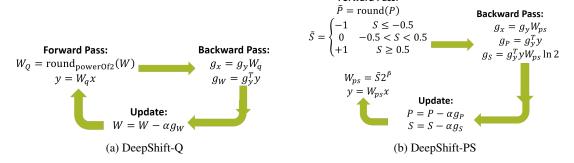


Figure 2: This paper presents 2 training algorithms to train the bitwise shift and negation parameters of a model: (a) DeepShift-Q, and (b) DeepShift-PS.

nearest power of 2:

$$S = \operatorname{sign}(W)$$

$$\tilde{S} = S \tag{4}$$

$$P = \log_2(\operatorname{abs}(W))$$

$$\tilde{P} = \operatorname{round}(P)$$
(5)

$$\tilde{W}_{a} = \tilde{S} \cdot 2^{\tilde{P}} \tag{6}$$

Hence, the forward pass for our LinearShift operator will be $Y = \tilde{W}_q X = (\tilde{S} \cdot 2^P) X$ and for our ConvShift operator will be $Y = \tilde{W}_q \circledast X + b = (\tilde{S} \cdot 2^{\tilde{P}}) \circledast X + b$.

The gradients of the backward pass can be expressed as:

$$\frac{\partial C}{\partial X} = \frac{\partial C}{\partial Y} \frac{\partial Y}{\partial X} = \frac{\partial C}{\partial Y} \tilde{W}_{q}$$

$$\frac{\partial C}{\partial W} = \frac{\partial C}{\partial Y} \frac{\partial Y}{\partial \tilde{W}_{q}} \frac{\partial \tilde{W}_{q}}{\partial W}$$

$$\frac{\partial C}{\partial b} = \frac{\partial C}{\partial Y}$$
(7)

where $\frac{\partial C}{\partial Y}$ is the gradient input to the operator (derivative of the model loss (a.k.a cost function), C, with respect to the

operator output), $\frac{\partial C}{\partial X}$ is the gradient output to the operator (derivative of the model loss with respect to the operator input)¹, and $\frac{\partial C}{\partial W}$ is the derivative of the model loss with respect to the operator weights ².

Since W_q is a rounded value of W, the differentiation $\frac{\partial \tilde{W}_q}{\partial W}$ is zero at all points except at the points where W is a power of 2. This will result in W not being updated during backpropagation, and hence no learning or updating of W. We solve this issue by using a straight-through estimator (STE) [17]:

$$\frac{\partial \tilde{W}_q}{\partial W} \cong 1 \tag{8}$$

Hence, $\frac{\partial C}{\partial W}$ is set to:

$$\frac{\partial C}{\partial W} \cong \frac{\partial C}{\partial Y} \frac{\partial Y}{\partial \tilde{W}_a} \tag{9}$$

Figure 2a summarizes the forward and backward passes of DeepShift-Q.

The have simplified the expression of $\frac{\partial C}{\partial X}$ in Equation 7. To be exact, it is equivalent to $\frac{\partial C}{\partial y} \tilde{W}_q^T$ for LinearShift and $X \circledast \frac{\partial C}{\partial y}$ for ConvShift.

2For LinearShift, $\frac{\partial C}{\partial W} = x^T \frac{\partial C}{\partial y}$ and for ConvShift $\frac{\partial C}{\partial W} = x^T \frac{\partial C}{\partial y}$

 $W_{rot180} \circ \circledast \frac{\partial C}{\partial u}$

2.2. DeepShift-PS

DeepShift-PS on the other hand directly uses the shift, P, and sign flip, S, as the trainable parameters:

$$\begin{split} \tilde{P} &= \operatorname{round}(P) \\ \tilde{S} &= \operatorname{sign}(\operatorname{round}(S)) \\ \tilde{W}_{ps} &= \tilde{S} \cdot 2^{\tilde{P}} \end{split} \tag{10}$$

Note the definition of sign flip matrix, \tilde{S} , implies that it is a ternary value rather than a binary parameter, as each element in the matrix can take one of 3 values:

$$\tilde{s} = \begin{cases} -1 & \text{if } s \le -0.5\\ 0 & \text{if } 0 < s < 0.5\\ +1 & \text{if } s \ge 0.5 \end{cases}$$
(11)

where s is an element of matrix S and \tilde{s} is an element of \tilde{S} . The gradients of the backward pass can be expressed as:

$$\begin{split} \frac{\partial C}{\partial X} &= \frac{\partial C}{\partial Y} \frac{\partial Y}{\partial X} \\ &= \frac{\partial C}{\partial Y} W_{ps}^T \\ \frac{\partial C}{\partial P} &= \frac{\partial C}{\partial Y} \frac{\partial Y}{\partial \tilde{W}_{ps}} \frac{\partial \tilde{W}_{ps}}{\partial \tilde{P}} \frac{\partial \tilde{P}}{\partial P} \\ &= \frac{\partial C}{\partial Y} \frac{\partial Y}{\partial \tilde{W}_{ps}} W_{ps}^T \ln 2 \frac{\partial \tilde{P}}{\partial P} \\ \frac{\partial C}{\partial S} &= \frac{\partial C}{\partial Y} \frac{\partial Y}{\partial \tilde{W}_{ps}} \frac{\partial \tilde{W}_{ps}}{\partial \tilde{S}} \frac{\partial \tilde{S}}{\partial S} \\ \frac{\partial L}{\partial b} &= \frac{\partial L}{\partial Y} \end{split}$$

$$(12)$$

We use the straight through estimator to set $\frac{\partial \tilde{P}}{\partial P} \approxeq 1$ and $\frac{\partial \tilde{S}}{\partial S} \approxeq 1$. To evaluate, $\frac{\partial \tilde{W}_{ps}}{\partial S}$ we use the following derivation, treating the gradient of the sign operation as a straight-through estimator $(\partial \text{sign}(x) = \partial x)^3$:

$$\tilde{S} = \operatorname{sign}(\tilde{W_{ps}})
\partial \tilde{S} = \partial \operatorname{sign}(\tilde{W_{ps}})
\partial \tilde{S} \approx \partial \tilde{W_{ps}}
\frac{\partial \tilde{W}_{ps}}{\partial \tilde{S}} \approx 1$$
(13)

Hence, the gradients of the weights are set to:

$$\begin{split} \frac{\partial C}{\partial P} &\approxeq \frac{\partial C}{\partial Y} \frac{\partial Y}{\partial \tilde{W}_{ps}} \tilde{W}_{ps} \ln 2 \\ \frac{\partial C}{\partial S} &\approxeq \frac{\partial C}{\partial Y} \frac{\partial Y}{\partial \tilde{W}_{ps}} \end{split} \tag{14}$$

Figure 2b summarizes the forward and backward passes of DeepShift-PS.

3. Implementation

To implement the forward and backward passes, we follow an approach similar to that of [6]. We used PyTorch to define the forward pass and backward pass for the 2 custom ops: LinearShift and for ConvShift. In order to emulate the precision of an actual bitwise shift hardware implementation, the input data to the LinearShift and ConvShift operators and their bias parameters are rounded to fixed-point format precision before applying the forward pass.

To initialize the weights, we use Kamming initialization for W_q in DeepShift-Q, and uniform random distribution to initialize P and S in DeepShift-PS.

In DeepShift-PS, L2 normalization has been slightly modified to occur on W_{PS} rather than on P on S, i.e. the regularization term added is $\sum W_{PS}^2 = \sum 2^P S^2$ rather than $\sum P^2 + \sum S^2$.

4. Benchmark Results

We have tested the training and inference results on 3 datasets: MNIST [8], CIFAR10 [7], and Imagenet [3]. For each dataset, we have tested a group of architectures. We have explored various evaluations for the models:

 Original Version: evaluating the original architecture with standard convolution and linear operators,

2. DeepShift Versions:

- (a) Train from Scratch: start with randomly initialized weights, convert all the convolution and linear operators to their shift counterparts, and train from scratch using either training method Deepshift-Q or DeepShift-PS,
- (b) Train from Pre-trained Baseline: start with baseline model that is pretrained with 32-bits floating points weights, convert all the convolution and linear operators to their shift counterparts, convert the weights (using Equations 4 and 4) and train using either training method Deepshift-Q or DeepShift-PS.

We have noticed - in terms of validation accuracy results - that the best optimizer for DeepShift-PS is Rectified Adam (RAdam) [12] while the best optimizer for DeepShift-Q and regular 32-bit floating point training is stochastic gradient descent (SGD).

4.1. MNIST Data Set

Two simple models were trained and tested on the MNIST dataset:

 $^{^3}$ We have tried an alternative method to express $\frac{\partial \tilde{W}_{ps}}{\partial \tilde{S}}$ as $abs(W_{ps})$ but the training accuracy was very low.

Table 1: DeepShift accuracy results on MNIST dataset

Method	Train from Scratch	Train from Pre-trained			
Simple FC					
Original	96.92%	-			
DeepShift-Q	97.03%	94.91%			
DeepShift-PS	98.26%	98.26%			
Simple CNN					
Original	98.75%	-			
DeepShift-Q	98.81%	99.15%			
DeepShift-PS	99.12%	99.16%			

- Simple FC: a simple fully-connected model consisting of 3 linear layers with feature output sizes 512, 512, and 10 respectively. Dropout layers with probability of 0.2 were inserted in between the layers. All intermediate layers had a ReLu activation following it.
- Simple CNN: a model consisting of 2 convolutional layers and 2 linear layers. The 2 convolutional layers had output channels sizes 20 and 50 respectively, and both had kernel sizes of 5x5 and strides of 1. Max pooling layers of window size 2x2 followed by ReLu activation were inserted after each convolution layer. The linear layers had output feature sizes of 500 and 10 respectively.

A learning rate of 0.01, a momentum of 0.0, as well as batch size of 64 was used to train. The training plot is shown in Figure 3. The accuracy on the validation set is shown in Table 1.

4.2. CIFAR10 Data Set

ResNet18 and MobileNetv2 models were trained on scratch on CIFAR10 usnig DeepShit-PS method. The models were trained with momentum of 0.9 and weight decay of 1×10^{-4} . The loss criterion used was categorical cross entropy. The learning rate used to train was 0.1 that decays by a factor of 0.1 at epochs 80, 120, and 160, the number of epochs for training was 200, and the batch size was 128. For ResNet18, we have added an additional test of training the baseline (multiplication-based) 32-bit floating point weights using RAdam, to compare it with DeepShift-PS that uses RAdam and ensure the comparisons are fair. The best accuracy reached by the baseline ResNet18 model with SGD was 94.86% and with RAdam was 94.06%.

For training starting with a pre-trained model, the learning rate was 1×10^{-4} , and the number of epochs was 15.

The training plots are shown in 4 while the numerical results for evaluating on the validation set for various scenarios are shown in Table 2.

Table 2: DeepShift accuracy results on CIFAR10 dataset

Method	Train from Scratch	Train from Pre-trained		
ResNet18				
Original DeepShift-Q DeepShift-PS	94.86% 94.03% 92.34%	- 94.52% 94.52%		
MobileNet				
Original DeepShift-Q DeepShift-PS	92.42% 93.53% 91.97%	92.75% 92.48%		

4.3. Imagenet Data Set

For Imagenet, ResNet18 and VGG16 were trained from scratch with momentum of 0.9 and weight decay of 1×10^{-4} . The loss criterion used was categorical cross entropy, the number of epochs was 90, the learning rate used to train from scratch was 0.1 that decays by a factor of 0.1 every 30 epochs. As can be seen from the training plot in Figure 5, and the results in Table 3, the best accuracy reached when training ResNet18 from scratch was less than the full-precision model by around 4% in Top-5 accuracy and less than the 3% in Top-1 accuracy, but for VGG16 the drop was less than 0.3%. For ResNet18, DeepShift-PS was slighlty better than DeepShift-Q in Top-1 accuracy, but vice versa for DeepShift-PS. Training VGG16 using DeepShift-PS is under progress. We have shown comparison to other work that trained models from scratch multiplication replaced with XNOR operations.

For training from a pre-trained model, we show the results for ResNet18, ResNet50, AlexNet, VGG16, and GoogleNet. We see that for GoogleNet, the accuracy obtained was higher than the original model. We show results to compare with other work that converted 32-bit weights to powers of 2.

5. Efficiency Analysis

Binary multiplication for integer or fixed-point format consists of multiple bit-wise shifts, ANDs, and additions. Floating point (FP) multiplication consist of more steps such as adding the exponents, and normalizing. In terms of CPU cycles, taking 32-bit Intel Atom instruction processor as an example, integer and FP multiplication instruction take 5 to 6 clock cycles, while a bit-wise shift instruction takes only 1 clock cycle ([4]). ([1]) shows for 16-bit architectures, average power and area of multipliers are at least $9.7\times$, and $1.45\times$, respectively, of bitwise shifter: 22.05 nW to 54.83 nW vs 1.32 nW to 2.27 nW, and 8.7 K to 29 K transistor count vs 2 K to 6 K only. Hence, replacing multiplica-

Table 3: DeepShift accuracy results on Imagenet dataset and comparison with other methods. Columns "A" and "W" show the number of bits used for each model to represent actiations and weights in each method. For fair comparison, we also reported a separate column for change from baseline that each paper had, as different papers reported different baselines.

Method	A	W	Top-1		Top-5	
	Accuracy			Delta	Accuracy	Delta
ResNet18						
Original	32	32	69.76%	=	89.08%	-
Train from Scratch						
DeepShift-Q from Scratch [Ours]	32	6	65.28%	-4.48%	86.45%	-2.63%
DeepShift-PS from Scratch [Ours]	32	6	65.63%	-4.13%	86.33%	-2.75%
BWN [14]	32	1	60.80%	-9.30%	83.00%	-6.20%
TWN [10]	32	2	61.80%	-3.60%	84.20%	-2.56%
Train from Pre-Trained Model						
DeepShift-Q from Pre-trained [Ours]	32	6	68.67%	-1.09%	88.61%	-0.47%
DeepShift-PS from Pre-trained [Ours]	32	6	68.32%	-1.44%	88.41%	-0.67%
INQ [18]	32	5	68.98%	+0.71%	89.10%	+0.41%
ABC [11]	32	5	68.3%	-1.00%	87.9%	-1.30%
ResNet50						
Original	32	32	76.13%	-	92.86%	-
Train from Pre-Trained Model						
DeepShift-Q from Pre-trained [Ours]	32	6	75.32%	-0.81%	92.65%	-0.21%
DeepShift-PS from Pre-trained [Ours]	32	6	75.29%	-0.84%	92.55%	-0.31%
INQ [18]	32	5	74.81%	+1.59%	92.45%	+1.21%
GoogleNet						
Original	32	32	69.78%	-	89.53%	-
Train from Pre-Trained Model						
DeepShift-Q from Pre-trained [Ours]	32	6	69.51%	-0.27%	89.24%	-0.29%
DeepShift-PS from Pre-trained [Ours]	32	6	69.87%	+0.09%	89.62%	+0.09%
INQ [18]	32	5	69.02%	+0.13%	89.28%	+0.25%
LogQuant [2]	32	6	69.14%	-0.36%	88.86%	-0.28%
VGG16						
Original	32	32	71.59%	-	90.38%	-
Train from Scratch						
DeepShift-Q from Scratch [Ours]	32	6	70.87%	-0.71%	90.09%	-0.29%
Train from Pre-Trained Model						
DeepShift-Q from Pre-trained [Ours]	32	6	71.30%	-0.29%	90.27%	-0.09%
DeepShift-PS from Pre-trained [Ours]	32	6	71.12%	-0.47%	90.08%	-0.30%
INQ [18]	32	5	70.82%	+2.28%	90.30%	+1.65%
LogNN [13]	4	5	N/A	N/A	89.00%	-0.50%
LogQuant [2]	32	6	N/A	N/A	89.62%	-0.31%
AlexNet						
Original	32	32	56.52%		79.07%	
Train from Pre-Trained Model						
DeepShift-Q from Pre-trained [Ours]	32	6	54.97%	-1.55%	78.26%	-0.81%
DeepShift-PS from Pre-trained [Ours]	32	6	55.33%	-1.19%	78.40%	-0.67%
INQ [18]	32	6	57.39%	+0.15%	80.46%	+0.23%
LogNN [13]	4	5	N/A	N/A	75.1%	-1.7%

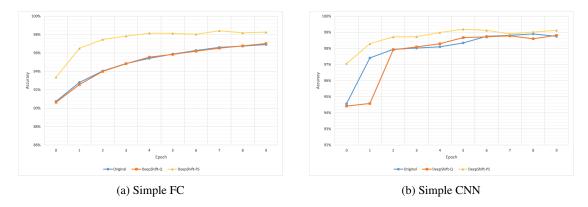


Figure 3: MNIST training from scratch

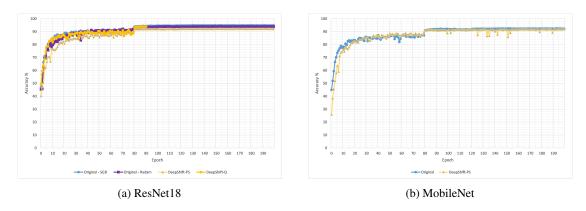


Figure 4: CIFAR10 training from scratch

tion with bitwise shift reduces energy. A common optimization in C++ compilers is to detect an integer multiplication with a (constant that is known at compilation time to be) a power of 2 and replace it with bitwise shift. Our contribution is enabling bitwise shift in training neural networks and achieving good accuracy.

DeepShift also achieves model storage compression by representing the weights with fewer bits. Since activations are represented us 32-bit fixed point activations, the range of shift values that we need to support are -32 to +32. Even if training results in shift values that are outside those ranges, they will be mathematically equivalent to multiplying by 0. Moreover, we noticed that shift values are rarely positive (i.e., most of the time $abs(W_Q) < 1$, hence P < 0 as $P = \log_2(abs(round(W_Q)))^4$). Therefore, we only need to support shift values of 0 to -32, and this requires 5 bits. Adding the bit required for the sign flip operation, we need 6 bits in total.

In addition to saving area and power for a potential custom hardware accelerator, our main advantages are to reduce the number of bits (thereby to reduce memory footprint of the model), and increase energy savings for transferring weights between different layers of memory of a GPU or CPU ([15]). That's needed in mobile and IoT applications, where memory footprint, even more so than FLOPs, are a primary constraint. Energy savings is crucial for applications like drones that fly on batteries, whether the computing platform is a low-end CPU, GPU, or FPGA. Also, bitwise shift operators are a good candidate for in-memory computing research [15], as they're cheaper in size and energy to place in or near memory than multiplication.

5.1. GPU Implementation

While, various custom hardware designs have been proposed by [9], [16] and [5], we have implemented a GPU kernel that inferred a bit-wise-shift based ResNet18 in 1490 seconds on Imagenet compared to baseline of 2010 seconds. To the best of our knowledge, this is the first paper to implement a CUDA kernel convolution using bitwise shifts.

It is noteworthy that the baseline that we compared against is a custom multiplication-based convolution kernel which is slower than NVIDIA's cuDNN optimized kernel. Future work can be done to optimize our bitwise-shift CUDA kernels to compare them with NVIDIA's cuDNN multiplication-based kernels, such as fusing convolution

 $^{^4}$ This expression holds for both W_Q and W_PS

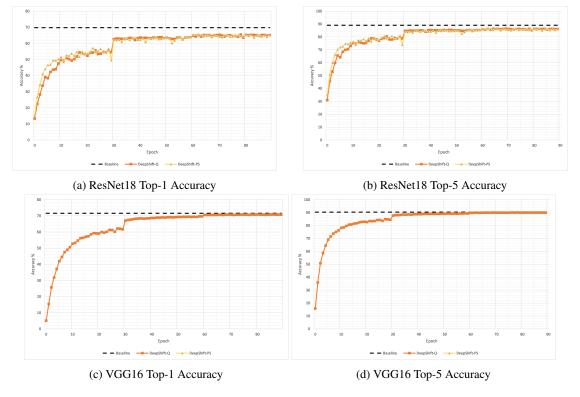


Figure 5: Imagenet training from scratch

with ReLu, tuning the tiling factors, and performing JIT (just-in-time) compilation of the CUDA kernels at runtime.

Some of the challenges that we faced when implementing the bitwise-shift convolution kernels, were that: - the GPU's instruction set architecture was not optimized for using smaller bit-width parameters that were needed to represent the shift and sign. Having to extract the 6-bit shift value and 1-bit sign value from 32-bit integers using masking and extra shifts incurred a lot of overhead, and - if-else statements on the values of the sign bits to determine whether to add or subtract during accumulation, or if-else conditions on whether the shift values are positive or negative to determine whether to shift left or right. If-else conditions stall instruction execution and are a bottleneck to pipelining instructions in hardware

These were bottlenecks that we had to mitigate in our code, and limited the speedup of the kernel compared to multiplication-based kernels. The implementation of the kernel can be found in our GitHub repo.

6. Related Work

[18] proposed Incremental Network Quantization (INQ) to quantize pre-trained full-precision DNNs with weights constrained to zeros and powers of two. This is accomplished by iteratively partitioning the weights into two sets, one of which is quantized while the other is retrained to

compensate for accuracy degradation.

LogNN [13] is another method that converts weights as well as activations to powers of 2. However, it requires manual empirical tuning rather than training. The authors of LogNN did present an algorithm to train weights and activations in powers of 2, but it was only implemented on pre-trained model on CIFAR10 dataset. Furthermore, the authors did not explain how sign flip of weights should be trained. The authors of LogNN also presented a hardware implementation in [9]. LogQuant [2] presented a method to convert the weights of a pre-trained model to powers of 2, such that no further training is needed and with minimal decrease in accuracy.

To the best of our knowledge, DeepShift is the first method to train a model from scratch using powers of 2.

7. Conclusion and Future Work

We have introduced DeepShift neural networks, which replace multiplications in the forward pass with bitwise operations - bitwise shift and negation, that can potentially lead to dramatic reduction in computation time, power consumption, and memory requirements during inference and training. We have proved that the accuracies of DeepShift networks on Imagenet and other datasets are near state-of-the-art.

While other methods that tackle neural network speed

ups such as BNNs perform well on small datasets (e.g., MNIST and CIFAR10) but suffer significant degradation on big datasets such as Imagenet, DeepShift networks proved that they are suitable for Imagenet.

Also, a GPU kernel of performing convolution using bitwise shift has been implemented as a proof-of-concept. Developing GPU instruction set architectures to make them more optimized for smaller bitwidth data types, and shift operations (e.g., an assembly instruction that can deduce bitwise shift should be done in the opposite direction if the shift value is negative) will help in making bitwise-shift-based training and inference more common. Moreover, existence of vector CPU instructions that can perform bitwise shifts given a vector of integers, and a vector of shift values, can make DeepShift more optimal on CPUs.

For future work, we suggest to look into training models from scratch with both activations and weights represented as powers of 2, similar to [13] that did that for inference only. Also, we would research into creating the bitwise-shift convolution kernels for CPU, as well as apply the well-known optimizations (fusing with activation, tiling tuning, and JIT compilation) for the GPU kernel.

Acknowledgments

We thank Ahmed Eltantawy for providing us valuable feedback to implement the CUDA kernels. We thank Sara Elkerdawy for providing valuable insights during our discussions with her. We thank also the developers of PyTorch for providing example scripts and pre-trained model binary files to reproduce the accuracies of various models on the Imagenet dataset.

References

- [1] Abhijit Rameshwar Asati. A Comparative Study of High Performance CMOS Multipliers, Barrel Shifters and Modeling of NBTI Degradation in Nanometer Scale Digital VLSI Circuits. PhD thesis, Birla Institute of Technology and Science, Pilani, Rajasthan, India, 2009. 5
- [2] Jingyong Cai, Masashi Takemoto, and Hironori Nakajo. A deep look into logarithmic quantization of model parameters in neural networks. In *Proceedings of the 10th International Conference on Advances in Information Technology*, IAIT 2018, pages 6:1–6:8, New York, NY, USA, 2018. ACM. 2, 6, 8
- [3] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei. ImageNet: A Large-Scale Hierarchical Image Database. In CVPR09, 2009. 4
- [4] Agner Fog. Instruction tables: Lists of instruction latencies, throughputs and micro-operation breakdowns for intel, amd and via cpus. https://www.agner.org/optimize/instruction_tables.pdf, 2018. 5
- [5] Denis Gudovskiy and Luca Rigazio. ShiftCNN: Generalized low-precision architecture for inference of convolu-

- tional neural networks. arXiv preprint arXiv:1706.02393, 2017.
- [6] Itay Hubara, Matthieu Courbariaux, Daniel Soudry, Ran El-Yaniv, and Yoshua Bengio. Binarized neural networks. In D. D. Lee, M. Sugiyama, U. V. Luxburg, I. Guyon, and R. Garnett, editors, Advances in Neural Information Processing Systems 29, pages 4107–4115. Curran Associates, Inc., 2016. 4
- [7] A. Krizhevsky. Learning Multiple Layers of Features from Tiny Images. Technical report, University of Toronto, Department of Computer Science, 2009. 4
- [8] Yann LeCun and Corinna Cortes. MNIST handwritten digit database. 2010. 4
- [9] E. H. Lee, D. Miyashita, E. Chai, B. Murmann, and S. S. Wong. Lognet: Energy-efficient neural networks using logarithmic computation. In 2017 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), pages 5900–5904, March 2017. 7, 8
- [10] Fengfu Li and Bin Liu. Ternary weight networks. CoRR, abs/1605.04711, 2016. 6
- [11] Xiaofan Lin, Cong Zhao, and Wei Pan. Towards accurate binary convolutional neural network. In I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, *Advances in Neural Information Processing Systems 30*, pages 345–353. Curran Associates, Inc., 2017. 6
- [12] Liyuan Liu, Haoming Jiang, Pengcheng He, Weizhu Chen, Xiaodong Liu, Jianfeng Gao, and Jiawei Han. On the variance of the adaptive learning rate and beyond. arXiv preprint arXiv:1908.03265, 2019. 4
- [13] Daisuke Miyashita, Edward H. Lee, and Boris Murmann. Convolutional neural networks using logarithmic data representation. *ArXiv*, abs/1603.01025, 2016. 2, 6, 8, 9
- [14] Mohammad Rastegari, Vicente Ordonez, Joseph Redmon, and Ali Farhadi. Xnor-net: Imagenet classification using binary convolutional neural networks. *CoRR*, abs/1603.05279, 2016. 6
- [15] Vivienne Sze, Yu-Hsin Chen, Tien-Ju Yang, and Joel S. Emer. Efficient processing of deep neural networks: A tutorial and survey. *Proceedings of the IEEE*, 105(12):22952329, Dec 2017. 7
- [16] Sebastian Vogel, Mengyu Liang, Andre Guntoro, Walter Stechele, and Gerd Ascheid. Efficient hardware acceleration of cnns using logarithmic data representation with arbitrary log-base. In *Proceedings of the International Conference on Computer-Aided Design*, ICCAD '18, pages 9:1–9:8, New York, NY, USA, 2018. ACM. 7
- [17] Penghang Yin, Jiancheng Lyu, Shuai Zhang, Stanley J. Osher, Yingyong Qi, and Jack Xin. Understanding straight-through estimator in training activation quantized neural nets. In *International Conference on Learning Representations*, 2019. 3
- [18] Aojun Zhou, Anbang Yao, Yiwen Guo, Lin Xu, and Yurong Chen. Incremental network quantization: Towards lossless cnns with low-precision weights. *CoRR*, abs/1702.03044, 2017. 2, 6, 8