

Control Variate Approximation for DNN Accelerators

Georgios Zervakis*, Ourania Spantidi[†], Iraklis Anagnostopoulos[†], Hussam Amrouch[‡], and Jörg Henkel*

*Chair for Embedded Systems (CES), Karlsruhe Institute of Technology, Karlsruhe, Germany

[†]Department of Electrical, Computer and Biomedical Engineering, Southern Illinois University, Carbondale, U.S.A.

[‡]Chair of Semiconductor Test and Reliability (STAR), University of Stuttgart, Stuttgart, Germany

*{georgios.zervakis,henkel}@kit.edu, [†]{ourania.spantidi,iraklis.anagno}@siu.edu, [‡]amrouch@iti.uni-stuttgart.de

Abstract—In this work, we introduce a control variate approximation technique for low error approximate Deep Neural Network (DNN) accelerators. The control variate technique is used in Monte Carlo methods to achieve variance reduction. Our approach significantly decreases the induced error due to approximate multiplications in DNN inference, without requiring time-exhaustive retraining compared to state-of-the-art. Leveraging our control variate method, we use highly approximated multipliers to generate power-optimized DNN accelerators. Our experimental evaluation on six DNNs, for Cifar-10 and Cifar-100 datasets, demonstrates that, compared to the accurate design, our control variate approximation achieves same performance and 24% power reduction for a merely 0.16% accuracy loss.

Index Terms—Approximate Computing, Arithmetic Circuits, Control Variate, Deep Neural Networks, Low Power, MAC Array

I. INTRODUCTION

Recent Deep Neural Networks (DNNs) have brought advancements in many fields over the past years. This wide variety of fields includes various applications, such as computer vision, speech recognition, scientific computing and many more [1]. However, the late DNN advancements have brought up the immense demands for computational power as well as for energy efficiency. This need is intensified especially in deployments on smart and IoT edge devices, since they have very restricted energy and computing resources. In addition, the increased DNN workload has led to the emerge of customized hardware DNN accelerators [1].

The main arithmetic computation during inference is the multiply-accumulate (MAC) operation. DNNs perform millions of MAC operations on their convolution and fully-connected layers. Therefore, DNN accelerators integrate thousands of MAC units. For example, Google TPU [1] comprises 64K MAC units while Google Edge TPU contains a 4K MAC units. This vast number of MAC units combined with high parallelization leads to very high power demands [2].

Recently, approximate computing emerged as a promising solution to develop energy/power-efficient circuits. Approximate computing leverages the intrinsic error resilience of a vast number of application domains to improve their energy profile at the cost of some accuracy loss [3]. Since DNNs are characterized by error-resilience, they inherently become appropriate candidates for approximate computations [3], [4]. Driven by this high potential for energy efficiency, significant research interest is shown in the design of approximate DNN accelerators [3]–[9]. State of the art mainly approximates the multipliers of the DNN accelerator, since the multiplier is the most power consuming component of the MAC unit. Modern DNNs are becoming gradually deeper, ending up with multiple layers that can differ significantly in error resilience [7]. In addition, the errors due to approximate circuits are not constant but they are highly input dependent [8]. Hence, different layers within the neural network or different neural networks require different approximation to satisfy an accuracy threshold [6]–[8]. Moreover, recent research showed that the deeper the neural network, the more sensitive it becomes to any approximation [6]. To overcome the aforementioned limitations,

existing works mainly apply retraining to mitigate the accuracy loss due to approximation [4], [5]. However, for DNNs, retraining targeting a specific approximate accelerator is very time consuming. In addition, in many cases the training set might not be available (e.g., proprietary models) and thus, retraining might be infeasible [7].

In this work, we introduce a *control variate approximation technique* that improves the accuracy of approximate DNN accelerators without the need to perform time-overwhelming retraining. Leveraging the accuracy improvement of our control variate approximation, we use [10] to replace the multipliers of the DNN accelerator with aggressive approximate ones (i.e., with high error but also with high power gain). The approximate multiplier of [10] is selected because i) it delivers very high power reduction, albeit its high error and ii) its error can be analytically modeled with simple equations. Our analysis demonstrates that our technique nullifies the mean value of the convolution error, and decreases its variance. Our experimental evaluation, over six DNNs trained on Cifar-10 and Cifar-100, shows that, for a merely 0.16% average accuracy loss, our approach delivers 24% power reduction and same performance compared with the accurate design. Our control variate approximation technique can be applied with any approximate multiplier as long as it is based on mathematical formulation instead of custom of logic simplification.

Our novel contributions within this paper are as follows:

- (1) We propose a *control variate approximation technique* for generating approximate DNN accelerators. Our technique does not require re-training and delivers high accuracy by controlling the error induced by the approximate multiplications.
- (2) We demonstrate that our technique improves the accuracy up to 21%, on average, compared to exactly the same approximate DNN accelerator, but without our proposed control variate approximation.
- (3) Leveraging our proposed control variate approximation, we can integrate highly approximate multipliers, thus significantly reduce the power consumption (more than 24%) of the whole DNN accelerator. In addition, our technique outperforms state-of-the-art works, achieving more than 3.8x higher energy reduction.

II. RELATED WORKS

There has been great interest around approximate computing for neural network inference. [5] employed approximate multipliers to different convolution layers and [4] proposed a compact and energy-efficient multiplier-less artificial neuron. However, [4], [5] are based on retraining to recover accuracy loss caused by the usage of approximation. Similar to [5], [9] utilized approximate multipliers and introduce an error compensation module. However, [5], [9] are evaluated on the LeNet network, a very shallow architecture which cannot provide the amount of operations recent DNNs do. Therefore, both of these methods can be deemed inapplicable in modern scenarios which require deeper network architectures. The authors in [7] propose a non-uniform architecture that utilized approximate multipliers from [11]. Their work tunes the weights accordingly and avoids

retraining. [7] applies layer-wise approximation, while power-gating the unused approximate multipliers. In [8], approximate multipliers with reconfigurable accuracy at run-time are generated. Similar to [7], they also apply layer-wise approximation, which is narrowing down potential benefits. In [6], approximate reconfigurable multipliers with low variance are generated using [8]. Then, [6] proposed a mapping algorithm to implement weight-oriented approximate inference in which the accuracy level of the approximate multiplier is determined by the weight's value. Nevertheless, to enable runtime reconfiguration and accuracy control, [6], [8] deliver limited energy savings due to the additional hardware.

Distinguish from Existing State of the Art: Our technique controls the error of the approximate multiplications and improves the accuracy without requiring re-training. Consequently, this high accuracy improvement enables the exploitation of highly approximated multipliers, maximizing the power gains.

III. CONTROL VARIATE APPROXIMATION

This section describes our control variate approximation technique and presents an error analysis of the approximate convolution. The core operation of a convolution is given by:

$$G = B + \sum_{j=1}^k W_j \cdot A_j, \quad (1)$$

where B is the bias of the neuron, W_j are the weights, and A_j are the input activations. In our approximate architecture and error analysis, approximate multipliers are used to replace the accurate multipliers of the DNN accelerator. We denote ϵ_j the multiplication error of the product $W_j \cdot A_j$. By error, we refer to the difference between the accurate and the approximate results. Thus, $\epsilon_j = W_j \cdot A_j - W_j \cdot A_j|_{approx}$. Given (1), the convolution error, namely ϵ_G , is equal to:

$$\epsilon_G = B + \sum_{j=1}^k W_j \cdot A_j - B - \sum_{j=1}^k W_j \cdot A_j|_{approx} = \sum_{j=1}^k \epsilon_j. \quad (2)$$

The error value of an approximate multiplier can be considered as a random variable, and is therefore defined by its mean value and variance [12]. Denoting the mean and variance of the approximate multiplier by μ_{AM} and σ_{AM}^2 , the mean and variance of the convolution operation are given by:

$$\begin{aligned} E[\epsilon_G] &= E\left[\sum_{j=1}^k \epsilon_j\right] = k\mu_{AM} \\ \text{Var}(\epsilon_G) &= \text{Var}\left(\sum_{j=1}^k \epsilon_j\right) = k\sigma_{AM}^2. \end{aligned} \quad (3)$$

Note that, the error values ϵ_j are independent variables and thus their covariance is zero [6], [12].

Hence, even if the approximate multiplier features small error (small μ_{AM} and σ_{AM}^2), the convolution error is significantly higher since it is proportional to the filter's size as (3) demonstrates. In [6], approximate multipliers with systematic error are employed and a constant correction term is used to compensate for the mean error (i.e., $E[\epsilon_G]$). However, even in this case, *the error of the convolution is still high, since it is defined by its high variance ($\text{Var}(\epsilon_G)$)*.

In our work, we propose the utilization of a control variate technique to reduce the convolution error variance. A control variate is an easily evaluated random variable, with known mean, that is highly correlated with our variable of interest. To implement our control variate approximation technique, instead of (1), we compute (4).

$$G^* = B + \sum_{j=1}^k W_j \cdot A_j|_{approx} + V \quad (4)$$

For the approximate multiplication ($W_j \cdot A_j|_{approx}$), we use the state-of-the-art power efficient perforated multipliers [10]. In [10], the m least partial products are perforated, i.e., they are not generated, and thus, they are removed from the accumulation tree. The latter results to smaller tree architectures that exhibit high delay, area, and power gains [10]. Moreover, a significant feature of [10] is that it applies functional approximation. Hence, its error does not depend on any carry propagation and it can be modeled in a mathematically rigorous manner. Nevertheless, our control variate approximation and our analysis hereafter can be employed with any approximate multiplier that exhibits a similar behavior with [10], i.e., the induced error can be described by an analytical model.

When perforating the m least partial products, the multiplication error of [10] is obtained by:

$$\begin{aligned} \epsilon_j &= W_j \cdot A_j - W_j \cdot A_j|_{approx} \\ &= W_j \cdot A_j - W_j \cdot (A_j - A_j \bmod 2^m) \\ &= W_j \cdot x_j, \text{ with } x_j = A_j \bmod 2^m = A_j \& (2^m - 1) \end{aligned} \quad (5)$$

Note that, for each filter in the DNN, the values of the weights (W_j) are fixed and are known a priori (after training and quantization). Thus, only the x_j change at run-time. In addition, $x_j \in [0, 2^m - 1]$ and thus requires only m -bits.

Considering the multiplication error ϵ_j in (5), we set our control variate V equal to:

$$V = \sum_{j=1}^k v_j = \sum_{j=1}^k x_j \cdot C_j, \quad v_j = x_j \cdot C_j \quad (6)$$

where C_j is a constant. Hence, ϵ_j and V_j feature a perfect linear correlation. Obviously, selecting $C_j = W_j$ will deliver accurate results. However, calculating $x_j \cdot W_j$ is computationally expensive and neglects the gains (area, power) of the perforated multiplier, since calculating V requires k multiplications and $k-1$ additions. Note that a control variate is easily evaluated. For this reason, we set $C_j = C$, $\forall C_j$. As a result, V is given by:

$$V = C \cdot \sum_{j=1}^k x_j \text{ and } v_j = x_j \cdot C. \quad (7)$$

Note that to calculate V in (7), only $k-1$ additions and 1 multiplication are required.

Given (7), the approximate convolution (4) is then written as:

$$\begin{aligned} G^* &= B + \sum_{j=1}^k (W_j \cdot A_j - \epsilon_j) + \sum_{j=1}^k v_j \\ &= G - \sum_{j=1}^k (\epsilon_j - v_j) \end{aligned} \quad (8)$$

and thus, the approximate convolution error equals:

$$\epsilon_{G^*} = \sum_{j=1}^k (\epsilon_j - v_j) = \sum_{j=1}^k (x_j \cdot (W_j - C)). \quad (9)$$

Hence, the variance $\text{Var}(\epsilon_{G^*})$ of the error of the approximate convolution, i.e., ϵ_{G^*} , is calculated from:

$$\begin{aligned} \text{Var}(\epsilon_{G^*}) &= \sum_{j=1}^k \text{Var}(\epsilon_j - v_j) \\ &= \sum_{j=1}^k ((W_j - C)^2 \cdot \text{Var}(x_j)) \\ &= \underbrace{\frac{(2^m - 1)(2^m + 1)}{12}}_{\text{Var}(x_j)} \sum_{j=1}^k (W_j - C)^2. \end{aligned} \quad (10)$$

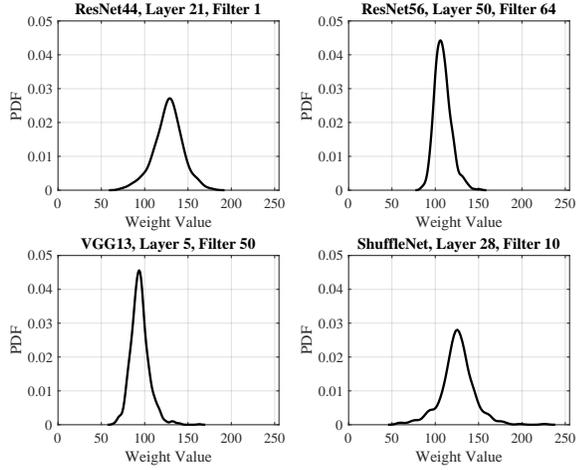


Fig. 1. Weight distribution of randomly selected filters of various NNs. Four examples are depicted.

As a result, $\text{Var}(\epsilon_{G^*})$ is minimized when:

$$\begin{aligned} \frac{d}{dC} \text{Var}(\epsilon_{G^*}) = 0 &\Rightarrow \\ C = \mathbb{E}[W_j] &= \frac{1}{k} \sum_{j=1}^k W_j. \end{aligned} \quad (11)$$

Proofs of (10) and (11) are simple and they are omitted due to space limitation. Note that $C \neq 0$, i.e., variance without our control variate (as in (3)). In addition, note that the more squeezed the weights' distributions is (i.e., concentrated close to $\mathbb{E}[W_j]$) the closer $\text{Var}(\epsilon_{G^*})$ is to zero. Fig. 1, shows the weights' distribution for four different examples. In Fig. 1, the neural networks and the respective filters and layers, were randomly selected out of the neural networks we consider in Section V. Similar results are obtained for the rest filters and neural networks. As shown in Fig. 1, for all the examined filters, the majority of the weights is well concentrated in a closed region (squeezed dispersion in Fig. 1). Hence, this feature boosts the efficiency of our variance reduction method, as explained above.

Using the C value, obtained in (11), that minimizes the variance, we compute the the mean convolution error $\mathbb{E}[\epsilon_{G^*}]$:

$$\begin{aligned} \mathbb{E}[\epsilon_{G^*}] &= \sum_{j=1}^k \mathbb{E}[\epsilon_j - v_j] \\ &= \sum_{j=1}^k \mathbb{E}[x_j] \cdot W_j - \sum_{j=1}^k \mathbb{E}[x_j] \cdot \mathbb{E}[W_j] \\ &= \underbrace{\frac{(2^m - 1)}{2}}_{\mathbb{E}[x_j]} \left(\sum_{j=1}^k W_j - k \cdot \mathbb{E}[W_j] \right) \\ &= 0 \end{aligned} \quad (12)$$

As a result, the proposed control variate approximation method with $V = \mathbb{E}[W_j] \sum_{j=1}^k x_j$, effectively nullifies the mean error of the approximate convolution and also manages to decrease its variance. In other words, the error distribution is constrained in a squeezed region around zero. Hence, high convolution accuracy is expected. However, as (10) shows, the larger the m is, the larger the error variance will be and thus the accuracy loss.

IV. APPROXIMATE DNN ACCELERATOR IMPLEMENTATION

In this work, as our DNN accelerator use case, we consider a micro-architecture similar to the Google TPU [1]. TPU comprises

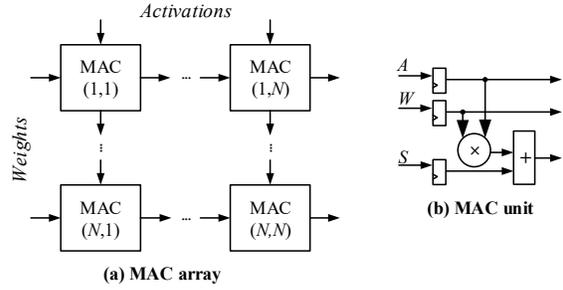


Fig. 2. The a) accurate systolic MAC array and b) MAC unit.

a large $N \times N$ systolic MAC array. Fig. 2 illustrates the accurate systolic MAC array. Fig. 3 depicts how the accurate MAC array (Fig. 2) is modified to implement our proposed control variate approximation. As shown in Fig. 3, our approximate MAC array features N rows but $N + 1$ columns. In the first N columns, MAC* units are used while the $N + 1$ column uses MAC+ units. The first N columns with the MAC* units calculate the convolution result $B + \sum_{j=1}^k W_j \cdot A_j|_{approx}$. The MAC+ unit in the last column is responsible for adding the control variable $V = \mathbb{E}[W_j] \sum_{j=1}^k x_j$ to the partial sum of the first N columns.

In the accurate MAC array, each MAC unit comprises an 8-bit multiplier and a $\lceil \log_2(N \times (2^{16} - 1)) \rceil$ -bit adder to avoid accumulation overflow [6]. For a 64×64 MAC array, the size of the adder is 22-bit.

In the approximate MAC array, the MAC* unit uses the perforated multipliers [10] to perform the approximate multiplication $W_j \cdot A_j|_{approx}$. The approximate product requires $16 - m$ bits, where m is the number of perforated partial products. Hence, we can also decrease the size of the adder by m bits. To achieve this, we have to shift the final partial sum m places left, as we will explain later. Moreover, for the partial sum input of the first column MAC*, instead of B we use $B[7 : m]$, in order to align the two inputs of the adder. In addition, each MAC* has to compute the partial sum required to calculate V (i.e., $\sum_{j=1}^k x_j$). As discussed in the previous section, the size of x_j is m -bit and thus, a $\lceil \log_2(N \times (2^m - 1)) \rceil$ -bit adder is required. Note, however, that $m < 8$. For a 64×64 MAC array and $m = 2$, the size of the adder is only 8 bits.

Overall, each MAC* computes the following:

$$\begin{aligned} P_j^* &= W_j \cdot A_j[7 : m] \\ sum_j &= sum_{j-1} + P_j^*, \quad sum_0 = B[7 : m] \\ sumX_j &= sumX_{j-1} + A_j[m - 1 : 0], \quad sumX_0 = 0 \end{aligned} \quad (13)$$

In terms of hardware requirements, the multiplication requires the accumulation of m fewer partial products and thus, $8 \cdot m$ fewer full adders are required [13]. In addition, the adder that computes sum_j requires m fewer full adders, while the adder that computes $sumX_j$ requires $\lceil \log_2(N \times (2^m - 1)) \rceil - 1$ full adders and 1 half adder [13]. Hence, the full adders reduction is about $(9 \cdot m - \lceil \log_2(N \times (2^m - 1)) \rceil + 0.5)$.

Moreover, $sumX_j$ and sum_j are independent and are computed in parallel. Hence, the adder $sumX_j$ is not on the critical path of MAC* and thus, a slow ripple-carry adder can be used to save power. Furthermore, MAC* can operate at lower delay than the accurate MAC, since the delay of both the multiplier and the sum_j adder is decreased. The multiplier has to accumulate m fewer partial products (i.e., shorter accumulation tree [10]), while the adder has to sum m

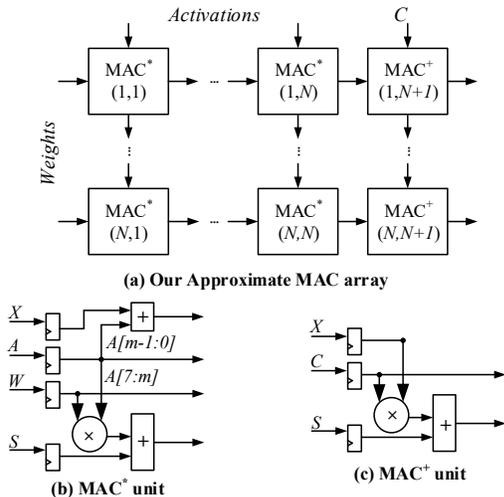


Fig. 3. a) Our approximate systolic MAC array, b) the MAC* unit, and c) the MAC+ unit.

fewer bits. Therefore, this delay slack enables downsizing of gates of the critical paths and boosts further the area and power savings [14].

As shown in Fig. 3, the last column of each row contains a MAC+ unit. MAC+ computes the following:

$$V = C \cdot \text{sum}X_N, \text{ where } C = E[W_j] \quad (14)$$

$$G^* = \{\text{sum}_N, B[m-1:0]\} + V \quad (15)$$

MAC+ requires an accurate $\lceil \log_2(N \times (2^m - 1)) \rceil \times 8$ -bit multiplier to compute the product of (14), i.e., $V = E[W_j] \sum_{j=1}^k x_j$. Furthermore, a $\lceil \log_2(N \times (2^{16} - 1)) \rceil$ -bit adder is required to produce the final output. The length of this adder and the length of the adder required by the accurate MAC unit are the same. Note that in (15), the partial sum sum_N , i.e., $\sum_{j=1}^k W_j \cdot A_j|_{\text{approx}}$, is shifted left m places and in its m -LSBs the missing m -LSBs of B , i.e., $B[m-1:0]$, are added. As a result, we achieve to both i) shift left to the required position the partial sum and ii) add the missing bits of the bias.

If the delay of MAC+ is higher than the delay of the accurate MAC, we pipeline MAC+ in order to sustain the same operating frequency. Hence, if the MAC+ unit requires l cycles, then the latency overhead of our proposed approximate MAC array is l cycles per convolution layer. However, considering that the inference phase requires thousands of cycles for each convolution layer [15], this overhead is negligible. Moreover, the MAC+ requires the value $C = E[W_j]$ to calculate V . This value can be transferred to the DNN accelerator among with the filters' weights. The memory overhead of the latter, is only 8 bits per filter.

As aforementioned, our approach replaces the accurate MAC units with the MAC* ones. Hence, the power and area of $N \times N$ MAC units are reduced. However, N additional MAC+ units are required to add V . The impact of MAC+ becomes negligible as the size of the MAC array increases. The gains due to MAC* increase quadratically w.r.t. N while the overhead due to MAC+ increases linearly w.r.t. N . Table I presents a theoretical estimation of the full adders reduction for perforation values $m = 1$ and $m = 2$. In addition, Table I shows both the reduction due to MAC* and the overhead due to the additional additional MAC+. The number of required full adders is calculated based on [13]. As shown in Table I, as N (MAC array size) or m (the perforation parameter) increase, the number of full adders required, decreases significantly. In addition, as N increases, the reduction due to MAC* becomes significantly larger

TABLE I
THEORETICAL EVALUATION OF FULL ADDERS (FA) REDUCTION

N	FA decrease due to MAC*	FA Increase due to MAC+	Total FA Decrease
$m = 1$			
16	1408	760	648
32	4608	1776	2832
48	8064	3048	5016
64	14336	4064	10272
$m = 2$			
16	3200	984	2216
32	11776	2224	9552
48	24192	3720	20472
64	43008	4960	38048

than the increase due to MAC+. As a result, the impact of MAC+ is insignificant compared to the gains due MAC*. Even for small MAC arrays, e.g., $N = 16$, when $m = 1$, the reduction due to MAC* is 2.59x higher than the overhead due to MAC+. Similarly, when $m = 2$, the respective value is 3.25x. Therefore, Table I validates our hypothesis in Section III that the hardware cost of our proposed control variate method is negligible.

Table I provides a simplistic theoretical evaluation of the proposed architecture in order to provide an initial insight on the impact of the MAC* and MAC+ units. For this reason, we use the gains in full adders as a representative metric. Components such as registers and partial product generation circuits are not considered. In addition, the impact of logic downsizing is not evaluated either. Section V provides a comprehensive experimental evaluation of the hardware gains.

V. EXPERIMENTAL RESULTS

In this section, we examine the efficiency our proposed control variate approximation in terms of area, power, and accuracy. For our analysis we employ industry-strength tools. All the MAC arrays evaluated in this section, are synthesized using Synopsys Design Compiler and are mapped to a commercial 14nm technology library. Synthesis is performed using the "compile_ultra" command and targeting performance optimization. Each accurate MAC array is synthesized at its critical path delay, while the approximate ones are synthesized at the critical path delay of the respective accurate one. Hence, the accurate and the respective approximate MAC arrays feature the same clock frequency. Post-synthesis timing simulations are performed using Mentor Questasim to obtain precise switching activity. Then, Synopsys PrimeTime is used to calculate the power consumption. For power analysis, we simulate the examined MAC arrays for 10,000 inference cycles to obtain accurate switching activity estimation. Running post-synthesis timing simulations for the entire inference phase is infeasible due to the vast time required [6]. The inference accuracy is captured by extending the approximate TensorFlow implementation of [16] with our control variate approximation. Six NNs of varying size, depth, and architecture are considered. The NNs are trained on the Cifar-10 and Cifar-100 datasets.

A. Hardware Evaluation

First, we evaluate the power and area gains of our proposed approximate MAC array compared to the accurate one. Both MAC arrays are implemented using the optimized arithmetic components of the Synopsys DesignWare Library, as typically done in commercial flows. Note that, unlike logic approximation approaches [6]–[8],

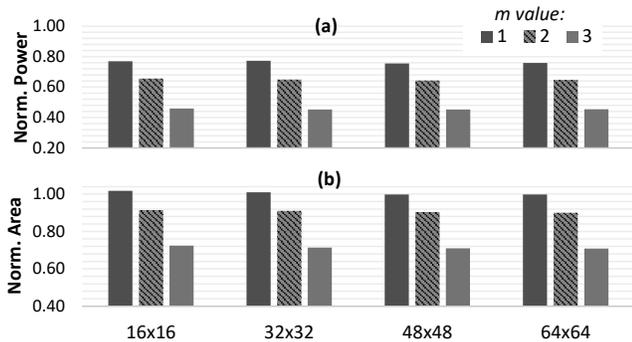


Fig. 4. The a) power and b) area of our control variate approximation for varying m values and MAC array sizes. The area and power values are normalized over the corresponding values of the respective accurate design.

the perforated multiplier, we use, applies functional approximation and thus, it can be implemented with the optimized DesignWare components [10], [17]. Finally, note that for the examined sizes, i.e., 16×16 to 64×64 , the MAC⁺ unit required 1 cycle and thus, the overall overhead is only 1 cycle per convolution layer.

Fig. 4 presents the area and power savings delivered by our technique for varying MAC array sizes and perforation values (i.e., m). As shown in Fig. 4a, our approximate MAC arrays achieve large power reduction, ranging from 22.8% up to 54.7%. In addition, it is observed that the power gain depends mainly on the perforation value m . For $m = 1$, the power reduction ranges from 22.8% to 24.2%. For $m = 2$, the respective values are from 34.5% to 35.7%, while for $m = 3$, the power decreases from 54.1% to 54.8%. This behavior verifies that the impact of the MAC⁺ units is negligible and thus, the power reduction of the MAC array is mainly defined by the power reduction delivered by the MAC* units. In Fig. 4b, the area reduction follows the same trend, i.e., it is defined by the m value and is slightly affected by the array size. The area gain goes up to 29% for $m=3$. For $m = 1$, the area gains due to the reduction in combinational logic is compensated by the increase in the number of registers required in the MAC* units (compared to the accurate MAC unit). Hence, for $m = 1$ the area remains almost the same.

In Table II, we evaluate the impact of the MAC⁺ units that are required to add V . Table II presents the area and power of the MAC⁺ units w.r.t. the area and power of the entire approximate MAC array. For all the examined scenarios the MAC⁺ units occupy at most the 1.49% of the entire array. Similarly, the power consumed by the MAC⁺ units is up to 1.87% of the total power consumption. It is noteworthy, that these values correspond to the smallest MAC array (only 16×16) and the highest approximation ($m = 3$). As shown in Table II, as the MAC array size increases, the impact of the MAC⁺ units becomes insignificant, validating our claim in Section IV. For example, the respective values for the 64×64 MAC array are only 0.40% and 0.49%. Finally, Fig. 4 and Table II demonstrate the scalability of our technique/architecture, i.e., the achieved power savings are sustained independently of the MAC array size.

B. Accuracy Evaluation

Next, we evaluate the accuracy delivered by our control variate approximation. For our analysis, we consider six NNs and in Table III we report the accuracy loss for varying perforation values m . Note that the accuracy depends only on m and not on the size of the MAC array. Moreover, in order to analyze the efficiency of our control variate approximation, we report the accuracy when using the same

TABLE II
EVALUATION OF THE AREA AND POWER OVERHEADS OF MAC⁺

Percentage of Total Area (%)				
m	16×16	32×32	48×48	64×64
1	1.06	0.55	0.38	0.28
2	1.18	0.61	0.41	0.31
3	1.49	0.77	0.53	0.40
Percentage of Total Power (%)				
m	16×16	32×32	48×48	64×64
1	1.15	0.58	0.40	0.30
2	1.32	0.68	0.46	0.35
3	1.87	0.97	0.64	0.49

TABLE III
ACCURACY EVALUATION OVER SIX NEURAL NETWORKS TRAINED ON CIFAR-10 AND CIFAR-100

Accuracy Loss (%)						
NN on Cifar-10	$m = 1$		$m = 2$		$m = 3$	
	Ours ⁺	w/o V [*]	Ours	w/o V	Ours	w/o V
googlenet	-0.16	0.35	0.00	4.13	1.95	31.78
ResNet44	0.03	0.73	0.83	4.49	5.75	32.94
ResNet56	0.49	0.60	1.25	5.36	5.94	39.01
shufflenet	0.12	3.40	-0.48	7.6	6.53	31.74
vgg13	-0.01	0.23	-0.30	0.76	0.69	6.76
vgg16	-0.13	0.19	0.38	1.66	3.86	8.71
Average	0.06	0.92	0.28	4.00	4.12	25.16
NN on Cifar-100	$m = 1$		$m = 2$		$m = 3$	
	Ours	w/o V	Ours	w/o V	Ours	w/o V
googlenet	0.05	2.43	1.47	13.19	7.02	44.52
ResNet44	0.77	1.02	1.82	14.17	11.27	43.4
ResNet56	-0.34	3.02	2.25	15.83	14.34	44.59
shufflenet	0.20	9.47	1.09	5.92	6.57	15.84
vgg13	0.89	3.01	1.91	5.89	3.98	14.7
vgg16	-0.03	3.96	0.03	2.41	2.8	11.54
Average	0.26	3.82	1.43	9.57	7.92	29.10

⁺ Accuracy achieved when using the approximate multiplier of [10] with our proposed control-variate approximation.

^{*} Accuracy achieved when using only the approximate multiplier of [10].

perforated multipliers [10] (same m) without our proposed control variate technique (i.e., [10] w/o V). Negative values in Table III refer to accuracy improvement due to the use of approximation [7]. As shown in Table III, the average accuracy loss of our method for Cifar-10 is 0.06%, 0.28%, and 4.12% for $m = 1$, $m = 2$, and $m = 3$, respectively. The corresponding values for the more challenging Cifar-100 dataset are 0.26%, 1.43%, and 7.92%. As a result, our technique achieves $\sim 24\%$ power reduction for negligible accuracy loss, i.e., 0.16% on average on both datasets for $m = 1$. The power gains rise to $\sim 35\%$ ($m = 2$) for an average accuracy loss of only 0.85%. Finally, for 6.02% average accuracy loss ($m = 3$), the power savings jump to $\sim 55\%$.

Table III highlights also the efficiency of our control variate approximation in decreasing the convolution error. Compared with [10] w/o V , i.e., same approximation without our control variate, our technique achieves 2%, 6%, and 21% higher accuracy, on average, for $m = 1$, $m = 2$, and $m = 3$, respectively. As a result, considering Tables II-III, our proposed control variate approximation delivers significant accuracy improvement for minimal hardware cost.

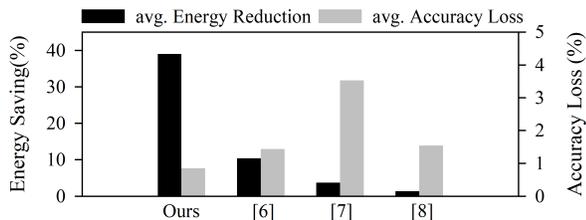


Fig. 5. Comparison of our technique against the state of the art [6]–[8]. The average energy reduction and average accuracy loss w.r.t. the NNs examined in Table III is reported.

C. Comparison with State of the Art

Finally, we compare our control variate approximation against three recent state-of-the-art works [6]–[8] that similar to our approach do not require retraining. In [7] a non-uniform architecture is used, that contains several approximate multiplier types from [11]. Each convolution layer uses only one multiplier type and the rest ones are power-gated. However, this leads to high area overhead and throughput loss due to the power-gated multipliers. For the fairness of the evaluation, we consider a uniform architecture for [7], with only one approximate multiplier type. In addition, note that the non-uniform approach of [7] is orthogonal with our work, since our technique also supports instantiating several MAC arrays with varying m values. For our comparison, we consider a 64×64 MAC array.

The architectures of [6]–[8] are based on the multipliers of EvoApprox library [11]. On the other hand, as typically done in commercial design flows, we employed the industry-level DesignWare library for our implementation. We implemented 64×64 MAC arrays using [6]–[8] and we evaluated their hardware characteristics. However, since [6]–[8] are based on suboptimal multipliers, the obtained MAC arrays feature significantly higher power, area, and delay compared with our *accurate* MAC array implemented with the DesignWare components. For this reason, we use also EvoApprox library [11] as the base of the MAC array implemented with our technique. For our technique, we select the perforation value $m = 2$ since, as shown above, it features high power reduction for moderate accuracy loss.

To compare our work against the state-of-the-art works [6]–[8], we evaluate, for each technique, the energy consumption and accuracy loss for all the NNs in Table III (for both Cifar-10 and Cifar-100). Energy consumption is calculated by $\text{cycles} \times \text{power} \times \text{delay}$. The cycles are obtain using the CNN cycle accurate simulator of ARM [15]. We evaluate the energy consumption since i) our technique requires 1 additional cycle per convolution layer and ii) [6], [8] use approximate reconfigurable multipliers with different configuration for each convolution layer. Fig. 5 reports, for each technique, the average energy reduction and accuracy loss. The energy reduction and accuracy loss are reported with respect to the corresponding values of the accurate design. The accurate MAC array is also implemented using the accurate multiplier (1JFF) of [11]. For [7], the results for the 125K multiplier are presented, since it delivered the the best energy-accuracy tradeof among all the approximate multipliers of [11]. As shown in Fig. 5, our technique delivers the highest energy savings. Compared with [6], [7], and [8] our technique achieves 3.8x, 10.5x, and 29.2x higher average energy reduction, respectively. In addition, as Fig. 5 shows, all the techniques achieve comparable accuracy, with our technique being the most accurate. However, in order to sustain high accuracy, [6]–[8] apply conservative approximation. As a result, [6]–[8] achieve very limited energy gains. On the other hand, our control variate approximation reduces significantly the error and

enables, thus, using highly approximate multipliers, such as [10]. Therefore, our control variate approximation boosts the obtained gains for minimal loss in accuracy.

VI. CONCLUSION

In this work, we introduced control variate approximation to increase the accuracy of approximate DNN accelerators without requiring any DNN retraining. Our mathematical analysis demonstrates that our technique mitigates the error induced by the approximate multipliers, by effectively nullifying the mean convolution error and reducing its variance. Our control variate approximation enables using aggressive approximate multipliers to design approximate DNN accelerators that boost the power savings.

ACKNOWLEDGMENT

This work is partially supported by the German Research Foundation (DFG) through the project “ACCROSS: Approximate Computing aCROSS the System Stack”.

REFERENCES

- [1] N. P. Jouppi *et al.*, “In-datacenter performance analysis of a tensor processing unit,” in *International Symposium on Computer Architecture*, 2017, pp. 1–12.
- [2] H. Amrouch *et al.*, “Npu thermal management,” *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems*, 2020.
- [3] H. Saadat *et al.*, “Minimally biased multipliers for approximate integer and floating-point multiplication,” *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems*, vol. 37, no. 11, pp. 2623–2635, 2018.
- [4] S. S. Sarwar *et al.*, “Energy-efficient neural computing with approximate multipliers,” *ACM Journal on Emerging Technologies in Computing Systems (JETC)*, vol. 14, no. 2, pp. 1–23, 2018.
- [5] V. Mrazek *et al.*, “Design of power-efficient approximate multipliers for approximate artificial neural networks,” in *Proceedings of the 35th International Conference on Computer-Aided Design*, 2016, pp. 1–7.
- [6] Z.-G. Tasoulas *et al.*, “Weight-oriented approximation for energy-efficient neural network inference accelerators,” *IEEE Transactions on Circuits and Systems I: Regular Papers*, 2020.
- [7] V. Mrazek *et al.*, “Alwann: Automatic layer-wise approximation of deep neural network accelerators without retraining,” in *International Conference on Computer-Aided Design (ICCAD)*, 2019, pp. 1–8.
- [8] G. Zervakis *et al.*, “Design automation of approximate circuits with runtime reconfigurable accuracy,” *IEEE Access*, vol. 8, pp. 53 522–53 538, 2020.
- [9] M. A. Hanif *et al.*, “Cann: Curable approximations for high-performance deep neural network accelerators,” in *2019 56th ACM/IEEE Design Automation Conference (DAC)*. IEEE, 2019, pp. 1–6.
- [10] G. Zervakis *et al.*, “Design-efficient approximate multiplication circuits through partial product perforation,” *IEEE Transactions on Very Large Scale Integration Systems*, vol. 24, no. 10, pp. 3105–3117, Oct 2016.
- [11] V. Mrazek *et al.*, “Evoapproxsb: Library of approximate adders and multipliers for circuit design and benchmarking of approximation methods,” in *Design, Automation & Test in Europe Conference & Exhibition*, 2017, pp. 258–261.
- [12] Chaofan Li *et al.*, “Joint precision optimization and high level synthesis for approximate computing,” in *Design Automation Conference (DAC)*, 2015, pp. 1–6.
- [13] V. Leon *et al.*, “Approximate hybrid high radix encoding for energy-efficient inexact multipliers,” *IEEE Transactions on Very Large Scale Integration (VLSI) Systems*, vol. 26, no. 3, pp. 421–430, 2018.
- [14] S. Venkataramani *et al.*, “Substitute-and-simplify: A unified design paradigm for approximate and quality configurable circuits,” in *Design, Automation and Test in Europe*, 2013, p. 1367–1372.
- [15] A. Samajdar *et al.*, “Scale-sim: Systolic cnn accelerator simulator,” *arXiv preprint arXiv:1811.02883*, 2018.
- [16] F. Vaverka *et al.*, “Tfapprox: Towards a fast emulation of dnn approximate hardware accelerators on gpu,” in *Design, Automation and Test in Europe Conference (DATE)*, 2020, p. 4.
- [17] H. Saadat *et al.*, “Realm: Reduced-error approximate log-based integer multiplier,” in *2020 Design, Automation Test in Europe Conference Exhibition (DATE)*, 2020, pp. 1366–1371.