

# Genetic Algorithm-Based Energy-Aware CNN Quantization for Processing-In-Memory Architecture

Beomseok Kang<sup>✉</sup>, Anni Lu<sup>✉</sup>, Yun Long<sup>✉</sup>, *Student Member, IEEE*,

Daehyun Kim<sup>✉</sup>, *Graduate Student Member, IEEE*, Shimeng Yu<sup>✉</sup>, *Senior Member, IEEE*,  
and Saibal Mukhopadhyay, *Fellow, IEEE*

**Abstract**—We present a genetic algorithm based energy-aware convolutional neural network (CNN) quantization framework (EGQ) for processing-in-memory (PIM) architectures. EGQ predicts layer-wise dynamic energy consumption based on the number of ADC access. Also, EGQ automatically optimizes layer-wise weight/activation bitwidth that can reduce total dynamic energy with negligible accuracy loss. As EGQ requires basic CNN model information such as weight/activation dimensions to predict the dynamic energy, various models can be compressed by EGQ. We analyse the effectiveness of EGQ on the area, dynamic energy, and energy efficiency of PIM architectures for VGG-19, ResNet-18, and ResNet-50 using NeuroSim. We observe EGQ is an effective approach for the CNN models to reduce the dynamic energy in various PIM designs with SRAM, RRAM, and FeFET technologies. EGQ achieves 6.1 bit of average weight bitwidth and 6.3 bit of average activation bitwidth in ResNet-18, that improves energy efficiency by 6.5 $\times$  than the 16-bit model. For ResNet-18 with CIFAR-10, 2.5 bit and 3.9 bit of average weight and activation bitwidth are achieved. Both results show the negligible accuracy loss of 2%.

**Index Terms**—Quantization, convolutional neural network, genetic algorithms, processing-in-memory.

## I. INTRODUCTION

**D**EEP neural network (DNN) inference is computationally expensive for resource-constrained mobile devices. Processing-in-memory (PIM) is an attractive hardware architecture for energy-efficient DNN inference as processing units are integrated into memory [1]–[3]. Weight and activation quantization [4]–[8] are popular algorithms to reduce the computational cost of DNN. Recently, hardware-aware quantization algorithms have shown that the energy efficiency of DNN inference can be further improved by considering hardware architectures and algorithms together [9], [10].

Manuscript received March 27, 2021; revised August 8, 2021 and October 16, 2021; accepted November 1, 2021. Date of publication November 10, 2021; date of current version December 13, 2021. This work was supported in part by the Energy-Efficient Computing: from Devices to Architectures (E2CDA) Program of the National Science Foundation under Grant CCF1740197 and in part by the Semiconductor Research Corporation under Grant NCORE2762.003. This article was recommended by Guest Editor J.-S. Seo. (Corresponding author: Beomseok Kang.)

Beomseok Kang, Anni Lu, Daehyun Kim, Shimeng Yu, and Saibal Mukhopadhyay are with the Department of Electrical and Computer Engineering, Georgia Institute of Technology, Atlanta, GA 30332 USA (e-mail: beomseok@gatech.edu).

Yun Long is with Google, Mountain View, CA 94043 USA.

Color versions of one or more figures in this article are available at <https://doi.org/10.1109/JETCAS.2021.3127129>.

Digital Object Identifier 10.1109/JETCAS.2021.3127129

It is important to consider data movement such as memory access in CMOS-based architectures to effectively improve the energy efficiency [11]. However, PIM architectures avoid the large number of memory access, and rather analog-to-digital converters (ADC) access is the main energy overhead. This motivates energy-aware quantization for PIM to differentiate from the quantization for CMOS-based architectures.

Uniform bitwidth quantization is widely used and simple to adapt for various DNN models. However, significant accuracy loss is observed at very low bitwidth [12]. Sun *et al.* [13] shows uniform bitwidth quantization with negligible accuracy loss, but it needs to re-train models after the quantization. Flexible bitwidth quantization is a promising technique to improve compression ratio and/or reduce accuracy loss compared to uniform quantization [6], [8]. It is important to note that the bitwidth of DNN can be differently determined depending on quantization methods while showing the same accuracy level. In other words, optimizing the bitwidth only with the constraint of accuracy does not guarantee the optimal energy efficiency. However, most of the flexible quantization algorithms for PIM do not explicitly consider energy consumption [14], [15]. There are few works on energy-aware quantization algorithms for PIM with simplified energy cost [13], [16]. As the distribution of ADC energy in each layer is non-uniform [17], assigning different bitwidth with the consideration of the energy distribution can be effective in terms of energy efficiency. Considering ADC access in flexible bitwidth quantization algorithms for PIM is still required.

In this paper, we present a genetic algorithm based energy-aware CNN quantization framework for PIM architectures (EGQ) that automatically optimizes weight and activation bitwidth for each layer. EGQ is a flexible bitwidth quantization algorithm that can predict layer-wise dynamic energy consumption in PIM platforms. The predicted dynamic energy is directly utilized to evaluate and search bitwidth candidates that achieve high energy efficiency with low accuracy loss. As the large amount of dynamic energy is consumed by ADC in PIM [1], we mainly focus on the number of ADC access to predict the layer-wise dynamic energy consumption. EGQ only requires the basic CNN model information such as activation sizes (i.e. dimensions) and kernel sizes to predict dynamic energy. Also, EGQ does not need re-training, which means any pre-trained CNN models can be simply compressed by EGQ.

EGQ uses genetic algorithm to automatically search appropriate layer-wise weight and activation bitwidth. The concept of using genetic algorithm for flexible bitwidth quantization was recently proposed [7]. However, it does

not optimize activation bitwidth, and more importantly, it is a hardware-agnostic algorithm as ADC cost related to the PIM performance is not considered. The main contributions of this paper is to develop an automatic framework for searching optimal layer-wise weight and activation bitwidth considering the number of ADC access and the associated dynamic energy. We mainly focus on designing the cost function to estimate the number of ADC access, and genetic algorithm is used as an optimization method to prove that overall energy efficiency in PIM can be improved by our cost function. It is important to note that EGQ is not limited to genetic algorithm as the cost function can be adapted to other problem solvers such as reinforcement learning [16].

We evaluate the quantization performance of EGQ with three CNN models, VGG-19 [18], ResNet-18, and ResNet-50 [19]. Also, we analyse the effect of EGQ on the area, dynamic energy, and energy efficiency of PIM architectures for different CNN models using NeuroSim [17]. As PIM platforms can be combined with various memory devices such as SRAM and embedded non-volatile memories (eNVMs), we study how the effectiveness of EGQ changes for PIM designs with different memory technologies. Finally, we compare EGQ with previous works on PIM-aware quantization and genetic algorithm based quantization [7], [12], [13].

We observe EGQ is an effective approach to reduce dynamic energy in various PIM designs with SRAM, RRAM, and FeFET technologies. In particular, eNVMs show the higher improvement of energy efficiency with EGQ due to low leakage energy. EGQ improves the energy efficiency of ResNet-18 by  $6.5\times$  compared to the 16-bit model. Also, EGQ shows 15% - 103% higher energy efficiency than the existing quantization algorithms for PIM with 2% accuracy loss.

## II. RELATED WORKS

In this section, we introduce details of the related works on PIM-aware quantization [13], [15], [16], [20], [21], and compare these works with EGQ to state the differences between EGQ and these works.

### A. Energy-Unaware Quantization Algorithms for PIM

The paper by Zhu *et al.* [15] showed flexible weight and activation quantization for a RRAM-based PIM architecture with the consideration of storage and latency. This RRAM-aware quantization quantifies the storage based on the number of crossbars and the latency by the number of digital-to-analog converters (DAC) access as optimization objectives. Each quantities are related to weight and activation bitwidth, respectively. The storage and latency can indirectly reduce energy consumption, but energy is not reflected into a loss function in the research. Also, re-training is required after every iteration of the weight and activation quantization.

C. Zhang *et al.* [21] proposed a robust RRAM-based quantization framework by developing the quantization algorithm to minimize RRAM variation error. They calculate the numerical error between full-precision weight and weight with the noises introduced by quantization error and device variation. The error is used as a loss function to utilize gradient-descent and back-propagation algorithms. As the loss function is only related to weight parameters, activation bitwidth is not considered. Also, energy is not related to the loss function. Thus, the algorithm is only suitable for weight quantization.

Reinforcement learning is a widely used algorithm in AutoML. Qu *et al.* [20] used a reinforcement learning

based deep deterministic policy agent to search flexible weight bitwidth to maximize memory utilization. The memory utilization is estimated based on the weight mapping in PIM subarrays. Bitwidth candidates that have the high memory utilization and accuracy give the agent high rewards, so the agent automatically searches an optimal bitwidth candidate. However, the memory utilization does not consider activation bitwidth. As the large amount of energy consumption is resulted from activation [22], flexible weight and activation bitwidth are required for energy-aware quantization.

### B. Energy-Aware Quantization Algorithms for PIM

Sun *et al.* [13] proposed uniform bitwidth activation quantization for an energy-efficient PIM accelerator. In particular, they show a non-linear quantization scheme to reduce required ADC bitwidth. As ADC results in critical energy overhead in analog PIM, decreasing ADC bitwidth is an effective method to enhance the energy efficiency [23]. However, changing ADC bitwidth depending on CNN models requires variable resolution ADC macros that are difficult to design. Also, energy-aware regularization in a loss function was introduced in the paper. The energy consumption is estimated with the simple multiplication of the voltage and current at RRAM devices in PIM subarrays. However, the improvement due to the energy-aware regularization is marginal since the dynamic energy consumed by eNVMs is not critical.

Huang *et al.* [16] proposed reinforcement learning based flexible weight and activation quantization for PIM. A cost function is defined with weight, activation, and ADC bitwidth to incorporate hardware performance into the quantization algorithm. Two fractions are mainly considered in the cost function. One is the fraction of the number of parameters over total parameters, and the other one is the fraction of bitwidth over full precision. ADC bitwidth is multiplied to the fractions as an exponential term to reflect the large energy overhead of high-precision ADC. However, as the fractions only represent the compression ratio of weight and activation parameters, the cost function cannot accurately estimate the layer-wise energy consumption in PIM architectures.

### C. Genetic Algorithm-Based Quantization

Genetic algorithm based quantization was introduced by Long *et al.* [7]. The paper shows a flexible weight quantization algorithm named Q-PIM. Q-PIM shows automated quantization without re-training and the large amount of training dataset, so arbitrary DNN models can be efficiently quantized by Q-PIM. A cost function in Q-PIM is related to the compression ratio of weight and accuracy loss, but hardware cost is not considered. Also, activation bitwidth is not considered in the cost function. Hence, Q-PIM is suitable for memory-efficient quantization, but do not provide an automated, flexible, and energy-aware quantization algorithm for weight and activation in PIM designs.

### D. Contributions in EGQ Over Prior Works

Most of the related works show the storage-aware, latency-aware, or utilization-aware algorithms. Although there are some energy-aware algorithms, they require complex multi-precision ADCs or used overly simplified energy cost without considering the details of data layout/flow for weight mapping and input data transfer. In contrast, EGQ is an energy-aware algorithm for PIM and uses a cost function that

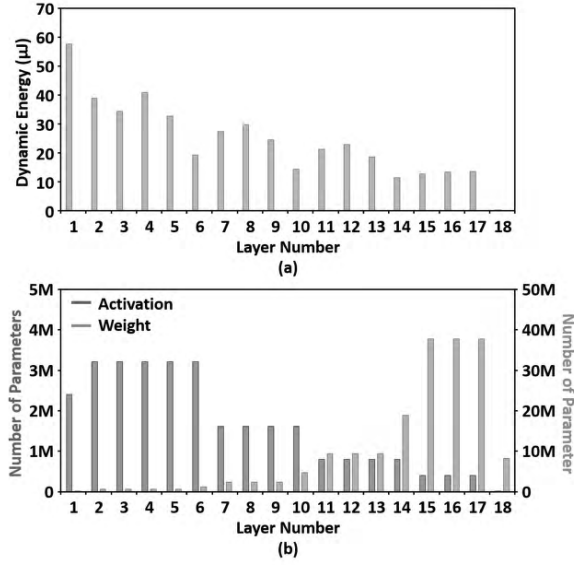


Fig. 1. The layer-wise dynamic energy distribution and the number of parameters of ResNet-18 with 16-bit precision. ImageNet dataset is used for simulation. (a) shows the dynamic energy simulation result from NeuroSim with 7nm SRAM technology node [17]. (b) shows the layer-wise number of weight parameters (orange) and activation parameters (blue).

includes weight and activation compression ratio, the number of ADC access, and accuracy loss. The cost function estimates energy using the number of ADC access which needs to consider the weight mapping and input data transfer in PIM. Accordingly, EGQ is a more suitable quantization framework for energy-efficient inference in PIM. As EGQ is based on genetic algorithm, the advantages of Q-PIM such as efficient search, automation, and no re-training are preserved. The overall complexity of EGQ can be slightly higher than Q-PIM due to activation quantization. However, the energy-awareness does not affect the complexity as the parameters required to estimate the number of ADC access is stored once at a initialization step. Thus, EGQ enables energy-aware flexible weight and activation quantization with the similar complexity of Q-PIM.

### III. ENERGY CONSUMPTION IN PIM ARCHITECTURE

#### A. Energy Model of NeuroSim

The methodology of the energy model in NeuroSim is shown in Fig. 2. Memory array dynamic energy of word-lines and bit-lines are computed by their corresponding  $CV^2$ . The number of open rows and columns are separately decided by input pattern and mux sharing. For peripheral circuits, the transistor sizes are predefined according to required drivability or experience data validated with real silicon. With transistor gate and parasitic capacitance parameters extracted from PTM model,  $CV^2$  at all nodes can be estimated and sum up. Except for ADC, whose energy estimation is from some polynomial functions fitted from SPICE simulation with various weight patterns and  $V_{ref}$  swept in a reasonable range. Finally, all the dynamic energy are sum up following the circuits and operations mapped from algorithm. For example, Table I shows how to compute the energy consumed in once reading of one whole memory array in parallel mode. The peripheral circuits and operations would vary a little with different modes. Buffer, interconnect and DRAM energy are also sum up to get the system-level energy.

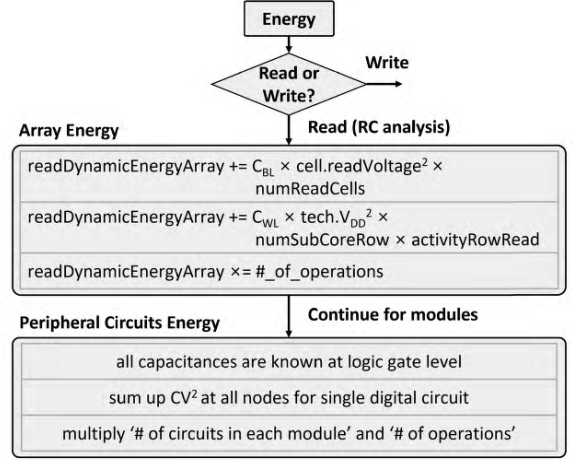


Fig. 2. Energy estimation flow in NeuroSim. numReadCells in Array Energy means the number of read columns per time. activityRowRead is determined by the number of ones in input vector.

TABLE I

ENERGY MODEL OF ARRAY AND PERIPHERAL MODULES IN NEUROSIM

Module	Energy Model
Array	$(C_{WL} \times V_{read}^2 \times \text{numOpenRow} + C_{BL} \times V_{DD}^2) \times \frac{\text{numCol}}{\text{numMuxShare}} \times \text{numMuxShare}$
Switch matrix	$\text{unitEnergy} \times \text{numOpenRow} \times \text{numMuxShare}$
Mux	$\text{unitEnergy} \times \text{numCol}$
Mux decoder	$\text{unitEnergy} \times \text{numMuxShare}$
ADC	Sum up of all unit energy referenced from fitting function
Shift-add	$\text{unitEnergy} \times \frac{\text{numCol}}{\text{numMuxShare}} \times \text{numMuxShare}$
Buffer	$\text{unitReadOrWriteEnergy} \times \frac{\text{numData}}{\text{interfaceWidth}}$
H-Tree	$(\text{unitLengthWireEnergy} + \text{unitLengthRepeaterEnergy}) \times \text{length} \times \text{numData}$
DRAM	$\text{energyPerBit} \times \text{numData}$

Though the energy model in NeuroSim considers the effect of different weight patterns on the read dynamic energy of ADC, EGQ assumes that the read dynamic energy of ADC access is constant regardless of the resistance and capacitance of reading columns. It simplifies the dynamic energy estimation to be proportional to the number ADC access, so we observe slight error of the estimation in Fig. 6.

#### B. Dynamic Energy Consumption of ResNet

Fig. 1 shows the layer-wise dynamic energy consumption of ResNet-18 in PIM platforms [17] and the layer-wise number of weight and activation parameters. Most of the dynamic energy is consumed in early stages as the early layers have relatively large input feature maps and small kernels than the later layers. The energy distribution is similar to the distribution in the number of activation parameters. The first layer consumes the largest dynamic energy, but the number of activation parameter at the first layer is not largest. This indicates assigning lower bits to the layers with many parameters does not guarantee optimal energy efficiency. Thus, it is important to predict the distribution of dynamic energy before determining weight and activation bitwidth. In this section, we discuss how to predict the dynamic energy based on weight mapping and input data transfer in PIM.

#### C. PIM Architecture

PIM-based DNN accelerators with various memory devices have attracted many interests due to high energy

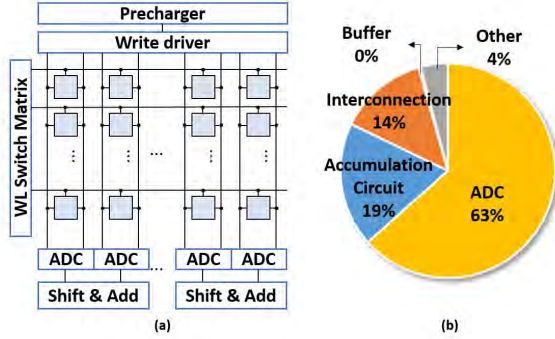


Fig. 3. (a) shows the schematic of a SRAM-based subarray in PIM architectures. Blue squares are SRAMs. (b) shows the dynamic energy consumption of circuit components in PIM. 16-bit ResNet-18 is used for dynamic energy simulation.

efficiency [1], [2], [24]. Most of the PIM architectures accelerate vector and matrix multiplication (VMM) operation based on parallel analog computations by bit-line current summation [1], [2], [24]. Fig. 3(a) shows the basic schematic of a SRAM subarray in PIM architectures [25]. We also assume the weight within the subarray is stored as multi-bit words, where each cell represents a single bit. We assume the input to a word-line is a serial bit-wise binary value which does not require DAC [26]. Once the input vector is applied to the word-lines, the current flowing through the bit-line is determined by weight values stored in the SRAM. The current is converted to a digital value by ADC without row-by-row access, and then the converted value is used as an input feature map (IFM) of the next layer. The parallel operation significantly increases computational efficiency in PIM architectures.

#### D. Energy Overhead of ADC

The majority of energy consumption in PIM comes from read peripheral circuits, mainly, ADC. Fig. 3(b) shows the dynamic energy consumption of circuit components in a PIM platform estimated by modified NeuroSim [17]. We use ResNet-18 with ImageNet dataset for the inference simulation.  $128 \times 128$  subarray size, 5-bit ADC, and 7nm SRAM technology are used for the simulation. More details about modified NeuroSim are described in the next section. We observe that ADC consumes about 63% of total dynamic energy. As ADC dynamic energy is proportional to the number of ADC access, the lower number of ADC access can effectively reduce the total dynamic energy. EGQ assumes the hardware design parameters such as the size of subarrays and ADC precision are fixed, and find optimal bitwidth for each layer to reduce the total number of ADC access.

#### E. Number of ADC Access

The number of ADC access in PIM architectures is determined by two factors, the length of input vectors and the number of subarrays to store weight matrices. All weight parameters are assumed to be stored on a single chip, and each layer uses different subarrays. Fig. 4 shows the schematic of IFMs and kernels for convolution operation. First, the number of subarrays for a weight matrix (i.e. kernel) is related to weight bitwidth and a weight dimension. Each kernel is unrolled along the columns in subarrays. The number of required rows are same with the multiplication of the width, height, and input channels of a kernel. (see gray rows in Fig. 4.) The number of required columns for each kernel

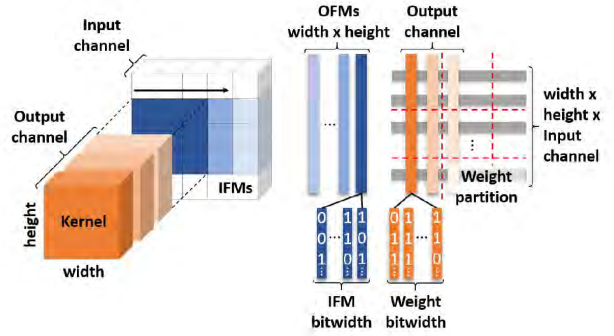


Fig. 4. The left figure shows the schematic of convolution operation with blue IFMs and orange kernels. The right figure shows the weight mapping on subarrays and the serial bit-wise input of IFMs to subarrays.

is same with weight bitwidth. Thus, the total memory size is determined by each kernel size, the number of kernels, and weight bitwidth. As we assume the size of subarrays is not flexible, the weight matrix should be partitioned by the subarray size. Red-dot lines in Fig. 4 indicate subarray boundaries. Next, the length of input vectors is related to the bitwidth of IFMs. IFMs with high precision need to be converted to long serial binary vectors that proportionally increase the number of ADC access. For convolution layers, the length of input vectors also relies on an output feature map (OFM) size. Kernels capture partial information from IFMs, and each capture constitutes one OFM. Thus, the number of captures is determined by the OFMs size, not the IFMs size. Blue bars in Fig. 4 show that the length of input vectors is proportional to IFM bitwidth and OFMs size.

### IV. PROPOSED APPROACH

#### A. Flow of Genetic Algorithm-Based Quantization

Genetic algorithm is a well-known algorithm for complex search and optimization. As flexible bitwidth quantization is also a search problem, EGQ uses genetic algorithm to determine appropriate bitwidth for each layer. Fig. 5 shows the overall flow and details of each stage in EGQ. The initialization step prepares weight and activation bitwidth candidates. The size of kernels and OFMs are stored at the initialization step as EGQ uses them repeatedly at the evaluation step. Fig. 5(b) and Fig. 5(c) are the lower and upper bound of ResNet-18 which constrains the range of initial candidates. The candidates are randomly sampled based on the two boundaries. Detail process to determine the upper and lower bound is explained in the previous study [7]. The candidates are evaluated based on fitness values at the evaluation step. Fig. 5(d) shows 15 candidates and the fitness values of each candidate. The red fitness values mean three good candidates, and they become parent candidates at the re-generation step. Other candidates are eliminated. The re-generation step includes crossover and mutation. Two parents from the three are randomly selected, and children candidates are again randomly sampled in the range of the two parents. EGQ repeatedly evaluates and reproduces candidates searching appropriate weight and activation bitwidth candidates. Final bitwidth can be further compressed by the greedy search based fine tuning [7].

#### B. Genetic Algorithm-Based Energy-Aware Quantization

In this subsection, we present the fitness function of EGQ and compare that with the fitness function of Q-PIM.

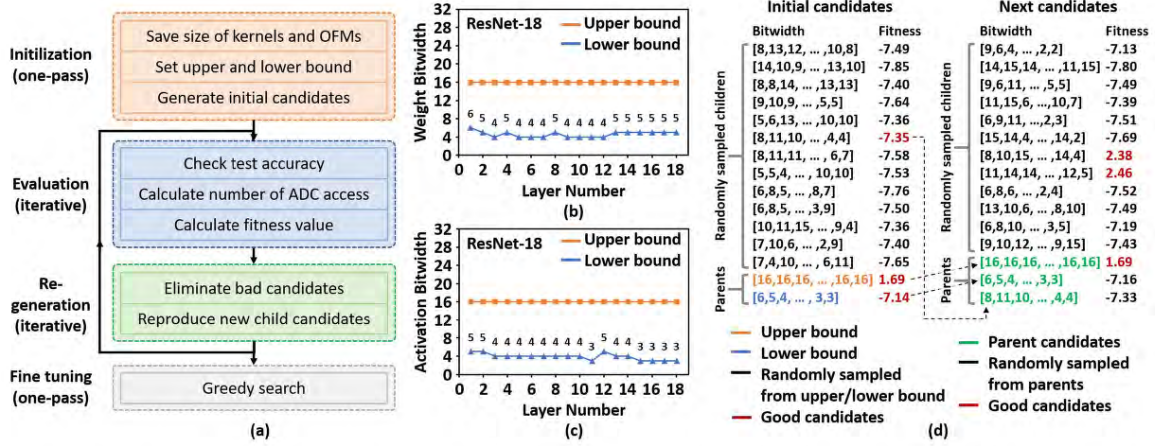


Fig. 5. (a) is the overall flow of EGQ. (b) and (c) show the upper and lower bound of ResNet-18. The upper bound is heuristically determined. (d) describes the first iteration of EGQ. The left half of (d) is the result from the initialization step, and the right area means next candidates after the first evaluation and re-generation step.

The fitness function of EGQ is given by

$$F(C) = \alpha \cdot C_W + \beta \cdot C_A + \gamma \cdot C_{ADC} + \delta \cdot \text{Accuracy} \quad (1)$$

where  $C$  represents a bitwidth candidate;  $C_W$ ,  $C_A$ , and  $C_{ADC}$  are the compression ratio of weight, activation (i.e. IFMs), and the number of ADC access, respectively.  $\alpha$ ,  $\beta$ ,  $\gamma$ , and  $\delta$  are weighting factors. EGQ optimizes the fitness function to find a bitwidth candidate that has high compression ratio and energy efficiency with low accuracy loss. The fitness function of Q-PIM is given by [7]

$$F(C) = -\alpha \cdot \sum_{i=1}^N W_i \cdot P_{W,i} - \beta \cdot \text{Error} \quad (2)$$

where  $W_i$  is the weight bitwidth of  $i^{th}$  layer;  $P_{W,i}$  is the number of weight parameters; Error represents accuracy loss. There are two main differences between EGQ and Q-PIM. First, the fitness function of Q-PIM does not consider the activation compression ratio. As equation (2) has only two terms related to the weight compression ratio and accuracy loss, it cannot be used for flexible activation quantization. Next, equation (2) does not have the connection between the quantization algorithm and hardware platforms, which is the main contribution of EGQ. The fitness function of EGQ considers the number of ADC access in PIM platforms. Thus, EGQ is more suitable for energy efficient quantization for PIM.

$C_{ADC}$  represents the compression ratio of the number of ADC access compared to the 32-bit model. It is assumed that the dynamic energy consumption is proportional to the number of ADC access as discussed above, so  $C_{ADC}$  reflects the expected total dynamic energy to the fitness function.  $C_{ADC}$  is defined by

$$C_{ADC} = 1 - \frac{\sum_{i=1}^N ADC_i(W_i, A_i)}{\sum_{i=1}^N ADC_i(32, 32)} \quad (3)$$

$$ADC_i = S_i \times OFM_{w,i} \times OFM_{h,i} \times A_i \quad (4)$$

$$S_i = \left\lceil \frac{K_{IC,i} \times K_{w,i} \times K_{h,i}}{s} \right\rceil \times \left\lceil \frac{K_{OC,i} \times W_i}{s} \right\rceil \quad (5)$$

where  $ADC_i$  is the number of ADC access;  $W_i$  and  $A_i$  are weight and activation bitwidth;  $S_i$  is the number of subarrays.  $OFM_{w,i}$  and  $OFM_{h,i}$  are the width and height of a OFM.

$K_{IC,i}$ ,  $K_{w,i}$ ,  $K_{h,i}$ ,  $K_{OC,i}$  are the input channel, width, height, and output channel of a kernel, respectively.  $s$  is a subarray size such as 128. As EGQ reduces weight and activation bitwidth, the ratio in equation (3) also becomes smaller. The ratio is subtracted from 1 to make good candidates have high fitness values.  $OFM_{w,i}$  and  $OFM_{h,i}$  are varying depending on layers. Fully-connected layers do not have width and height dimensions, so we regard both  $OFM_{w,i}$  and  $OFM_{h,i}$  as 1.  $S_i$  is the multiplication of the row and column number of subarrays. The first operand in  $S_i$  is the row number of subarrays, and the second operand is the column number of subarrays. Each operand is in a ceil function in equation (5). As  $ADC_i$  changes depending on candidates and iterations, the size of OFMs and kernels should be saved in EGQ. For the accuracy in the fitness function, we set an accuracy threshold that can give the penalty value of  $-10$ , if a candidate shows lower accuracy than the threshold. All quantization experiments here use 2% accuracy threshold. The weighting factors would change the behavior of EGQ, but they are set to 1 as default. We compare the impact of the weighting factors in the next section.

$C_W$  and  $C_A$  are determined by the ratio of the number of weight and activation parameters. Each of them are given by

$$C_W = 1 - \frac{\sum_{i=1}^N W_i \times P_{W,i}}{\sum_{i=1}^N 32 \times P_{W,i}} \quad (6)$$

$$C_A = 1 - \frac{\sum_{i=1}^N A_i \times P_{A,i}}{\sum_{i=1}^N 32 \times P_{A,i}} \quad (7)$$

where  $P_{W,i}$  and  $P_{A,i}$  mean the number of weight and activation parameters. The total number of the binary bits for weight parameters is obtained by the product of  $P_{W,i}$  and weight bitwidth,  $W_i$ . We assume weights in batch normalization layers are floating point, so weights in convolution layers and fully-connected layers are mainly compressed. As  $W_i$  at the denominator is 32, the ratio in equation (6) represents how much the weight of  $i^{th}$  layer is compressed compared to the 32-bit model.  $C_A$  can be explained in the same way as  $C_W$ . EGQ calculates the compression ratio of each candidate every iteration. Hence, the number of layer-wise weight and activation parameters should be stored in the initialization step.

EGQ is chosen to be layer-wise granular to limit the search space and hence, improve convergence time. EGQ can be

modified to support channel-wise optimization, but search space significantly increases. For example, ResNet-50 has 50 layers and 22,720 channels. The granularity of EGQ can be switched if equation (4), (7), and (8) are modified.  $N$  in the each equation represents the number of layers. EGQ estimates the compression ratio of weight, activation, and ADC access with total  $N$  number of precision assigned to each layer. If  $N$  is modified to represent the number of channel, EGQ can be adapted to be channel-wise granular.

### C. Linear Quantization

Linear quantization is a widely used scheme for DNN quantization. EGQ also uses linear quantizer for both weight and activation. The quantization of weight with  $n$ -bit precision is defined by

$$Q(x) = \text{round}\left(\frac{x(2^{n-1} - 1)}{\max(|x|)}\right) \times \frac{\max(|x|)}{2^{n-1} - 1} \quad (8)$$

Weight range is first normalized and expanded to  $2^n$  range. As weight generally includes negative values, the range is from  $-2^{n-1}$  to  $2^{n-1} - 1$ . Round function makes the expanded float values to be integers, thus introduces quantization error. The other multiplication term is for dequantization from the integer range to the float range again. For activation quantization, range can be different depending on the location of a quantizing layer. If activation quantization is directly connected to the OFMs, activation quantization uses same formula with weight quantization. However, if activation quantization is after ReLU function, there are no negative values. Thus, the range becomes from 0 to  $2^n - 1$ . EGQ quantizes activation after ReLU function.

## V. SIMULATION RESULTS

EGQ is implemented in Pytorch for accuracy analysis. VGG-19, ResNet-18, ResNet-50 with ImageNet dataset are mainly used for the evaluation. MobileNet-V2 and SqueezeNet are also explored to compare with prior works in the later section. Dynamic energy simulation is performed by modified NeuroSim mainly with 7nm SRAM technology node.

### A. Modified NeuroSim

DNN+NeuroSim is an integrated framework to benchmark processing-in-memory (PIM) accelerators for deep neural networks, with hierarchical design options from device-level, to circuit-level and up to algorithm-level [17]. It supports flexible PIM array design options with different device technologies (from SRAM to eNVMs) with various peripheral circuitry modules, which have been validated with SPICE simulations and actual device data. The framework supports automatic algorithm to hardware mapping for diverse models, and evaluates chip-level performance like area, latency and energy consumption. The automatic floorplan and dataflow design of NeuroSim is inherently compatible with various algorithm and device precisions. NeuroSim is originally designed to simulate the uniform bit-width model. To accommodate the flexible bit-width, layer-wise quantization is employed to extract the real trace data to feed the hardware performance estimation, and the different activation and weight precisions are assigned to corresponding layers. The digital resources like shift-add, accumulator and activation are designed to accommodate the maximum precision. We set a subarray size to  $128 \times 128$ , and ADC precision to 5 bit for all simulation.

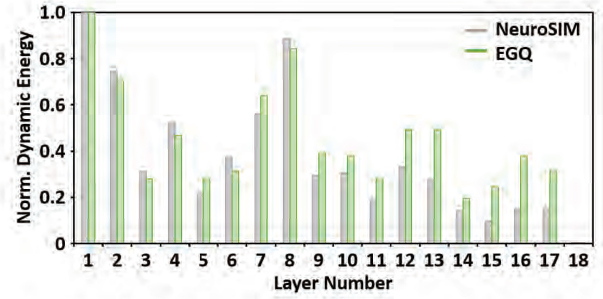


Fig. 6. Normalized dynamic energy comparison between EGQ and NeuroSim. The prediction from EGQ is based on the number of ADC access. Both results are obtained from ResNet-18 with same precision.

TABLE II  
ENERGY-AWARE AND ENERGY-UNAWARE  
QUANTIZATION RESULTS OF RESNET-18

Model Parameter	ResNet-18		ResNet-50	
	$\gamma = 0$	$\gamma = 1$	$\gamma = 0$	$\gamma = 1$
Normalized Number of ADC Access	0.30	0.26	0.36	0.28
Weight Bitwidth	8.1	7.1	8.1	8.0
Activation Bitwidth	7.6	7.6	8.0	8.3
Total Dynamic Energy ( $\mu$ J)	136.2	125.9	414.5	333.2
ADC Dynamic Energy ( $\mu$ J)	79.3	72.3	207.4	154.4
Total Leakage Energy ( $\mu$ J)	54.9	40.0	724.1	673.5
Energy Efficiency (TOPS/W)	27.9	32.2	9.0	10.2
Area ( $\text{mm}^2$ )	18.7	17.5	49.1	47.7
Baseline Accuracy (%)	69.8	69.8	76.2	76.2
Compressed Accuracy (%)	67.7	68.0	74.3	74.2

### B. ResNet Quantization Evaluation

Fig. 6 shows the layer-wise normalized dynamic energy of quantized ResNet-18 with flexible bitwidth. The prediction from EGQ based on the normalized number of ADC access closely matches the results from NeuroSim.

We compare ResNet-18 and ResNet-50 with energy-unaware (i.e.  $\gamma = 0$ ) and energy-aware (i.e.  $\gamma = 1$ ) cases in Fig. 7. Also, compression-unaware (i.e.  $\alpha, \beta = 0$ ) and compression-aware (i.e.  $\alpha, \beta = 1$ ) cases are compared in the same figure. As  $C_{ADC}$  is the ratio of the number of ADC access, the weighting factor of  $C_{ADC}$  ( $\gamma$ ) in the fitness function controls the energy awareness of EGQ. Similarly,  $\alpha$  and  $\beta$  control the compression awareness for weight and activation, respectively. Other weighting factors in both experiments are set to 1. Normalized ADC access in Fig. 7(a), (c), (e), and (g) means how much the number of ADC access is compressed compared to the 16-bit model.

In both ResNet-18 and ResNet-50, the energy-aware cases effectively reduce the number of ADC access than the energy-unaware cases showing the effectiveness of including  $C_{ADC}$  in the fitness function of EGQ. Fig. 7(b) and (f) show the dynamic energy of the circled candidates in Fig. 7(a) and (e). Both figures show the improved total dynamic energy in the energy-aware cases. As the difference in the normalized number of ADC access at the 100<sup>th</sup> iteration is larger in ResNet-50, we observe that the total dynamic energy of ResNet-50 is more reduced than ResNet-18. Table II summarizes the total dynamic energy, total leakage energy, and energy efficiency of the ResNet models. As we assume all weight parameters in each model are mapped to on-chip memory in PIM platforms, large leakage energy is observed in both models during the layer-wise dynamic energy simulation. In particular, Table II shows the total leakage energy is even

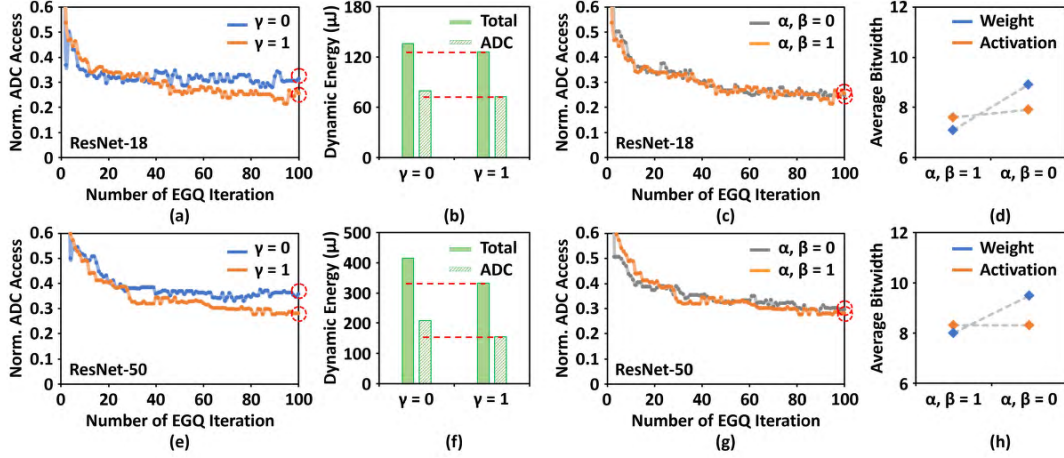


Fig. 7. (a) and (e) show the effectiveness of  $\gamma$  in EQG for ResNet-18 and ResNet-50, respectively. y-axis means the normalized number of ADC access. (b) and (f) show the dynamic energy comparison of the 100<sup>th</sup> EQG iteration results for ResNet-18 and ResNet-50, respectively. (c) and (g) compare the effectiveness of  $\alpha, \beta$  to the normalized number of ADC access. (d) and (h) show the average weight and activation bitwidth of the last candidates in (c) and (g). The weighting factors not mentioned in each figure are set to 1.

TABLE III  
LAYER-WISE BITWIDTH OF RESNET-18 COMPRESSED BY  
ENERGY-AWARE AND ENERGY-UNAWARE QUANTIZATION

ResNet-18	$\gamma = 0$	Weight	12,12,10,7,12,11,15,13,14,12,8,7,11,10,9,5,6,12
		Activation	8,9,7,9,11,5,3,7,8,6,5,5,10,10,6,5,9,10
	$\gamma = 1$	Weight	10,9,6,10,11,10,7,10,8,12,10,7,7,7,6,5,13
		Activation	8,9,6,6,3,8,13,12,7,9,4,10,10,8,5,9,8
ResNet-50	$\gamma = 0$	Weight	13,7,11,8,6,10,15,4,15,10,13,7,6,14,7,9,8,11,6,11,10,14,11,9,12,8,6,11,13,10,12,12,14,6,11,6,7,12,14,10,5,7,9,10,5,7,8,8,7,7
		Activation	11,11,11,7,6,10,8,7,11,11,7,4,6,8,13,7,5,11,15,5,13,8,6,6,11,11,12,8,14,11,7,8,10,9,5,12,12,8,11,9,10,8,13,12,10,13,5,13,10,8
	$\gamma = 1$	Weight	8,14,9,9,8,7,4,12,8,10,5,12,9,13,6,12,11,12,11,10,13,6,14,6,7,5,8,10,10,7,14,8,10,10,10,11,10,13,5,14,9,6,6,12,6,12,13,4,5,6
		Activation	9,9,11,10,7,8,4,14,11,15,8,5,8,12,8,8,4,7,13,7,8,11,6,8,11,5,7,6,11,6,8,4,9,4,6,10,6,8,4,5,5,12,10,4,11,5,8,8,5,9

higher than the total dynamic energy in ResNet-50. The energy efficiency is improved by 15.4% for ResNet-18, and 13.3% for ResNet-50, but the two cases show almost same accuracy (Table II).

We observe that compression awareness does not play a significant role in the number of ADC access in Fig. 7(c) and (g). As both cases include  $\gamma = 1$ , the normalized ADC access is reduce to the similar level of the energy-aware case. Fig. 7(d) and (h) compare the average weight and activation bitwidth at the circled candidates in Fig. 7(c) and (g). The compression-aware cases show lower weight and activation bitwidth. In particular, the weight bitwidth of the compression-aware cases is 1.5 bit lower in ResNet-18 and 1.8 bit lower in ResNet-50. Activation bitwidth is 0.3 bit difference in ResNet-18, and same in ResNet-50. We expect similar dynamic energy consumption in both compression-unaware and compression-aware cases, but the larger weight bitwidth will increase required memory footprint and leakage energy. Thus,  $\alpha$  and  $\beta$  should be considered together with  $\gamma$  in the fitness function to achieve compact and energy efficient models.

Table III summarizes the layer-wise weight and activation bitwidth for the energy awareness experiment. Fig. 8 visualizes the layer-wise weight and activation bitwidth in the

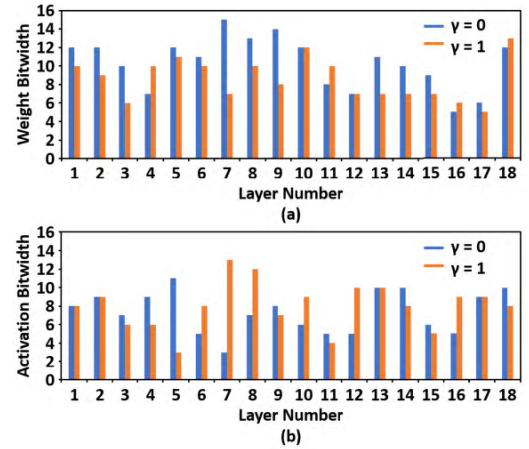


Fig. 8. The distribution of layer-wise weight (a) and activation bitwidth (b) in ResNet-18. Energy-aware case and energy-unaware case are compared.

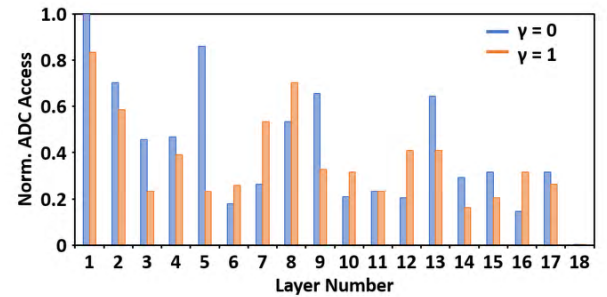


Fig. 9. The distribution of the normalized number of ADC access in ResNet-18. Energy-aware case and energy-unaware case are compared.

energy-aware and energy-unaware cases of ResNet-18. Fig. 9 shows the normalized number of ADC access for both cases. The 5<sup>th</sup> layer shows huge difference in the normalized ADC access. This is because the 5<sup>th</sup> layer has the much higher activation precision in the energy-unaware case while the weight bitwidth is similar. The 7<sup>th</sup> layer has also higher activation precision in the energy-unaware case, but the smaller weight precision mitigates the effect of the activation precision to the normalized ADC access. We observe that most of the

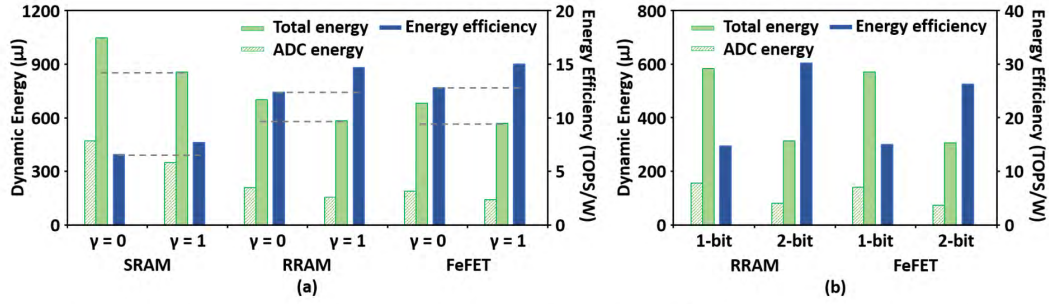


Fig. 10. (a) shows the dynamic energy and energy efficiency of ResNet-50 with single-bit SRAM, RRAM, and FeFET. Energy-unaware and energy-aware cases are compared with the memory devices. (b) compares the dynamic energy and energy efficiency of ResNet-50 between single-bit and multi-bit eNVMs. Energy-aware case is used for the comparison in (b).

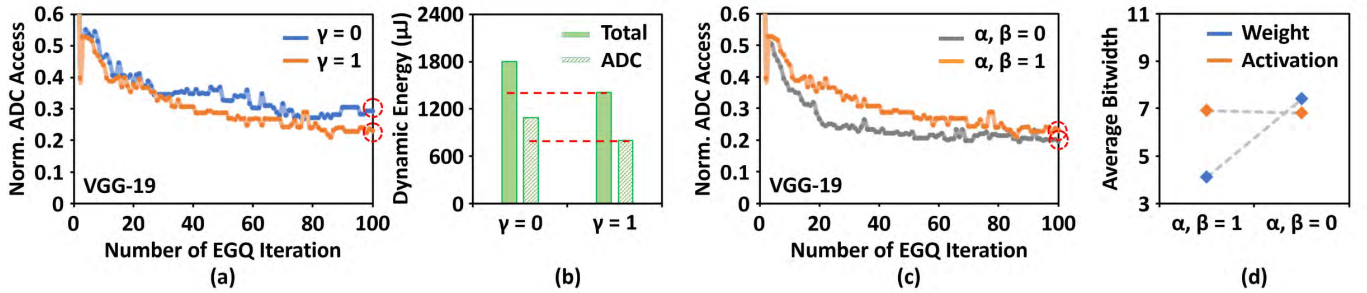


Fig. 11. (a) shows the effectiveness of  $\gamma$  in EGQ for VGG-19, respectively. (b) shows the dynamic energy comparison of 100<sup>th</sup> EGQ iteration results. (c) compare the effectiveness of  $\alpha, \beta$  to the normalized number of ADC access. (d) show the average weight and activation bitwidth of last candidates in (c). The weighting factors not mentioned in each figure are set to 1.

normalized ADC access is smaller in the energy-aware case due to the lower weight-bitwidth. Thus, the improved energy efficiency in the energy-aware case is reasonable.

Fig. 10 shows the performance comparison between the energy-unaware and energy-aware cases for ResNet-50 with various memory devices. We use single-bit SRAM, RRAM [27], and FeFET [28] in 22nm technology node in Fig. 10(a). We do not consider non-ideality and stochastic variations of the memory devices for NeuroSim simulation and accuracy analysis. The SRAM case shows higher total dynamic energy than the RRAM and FeFET cases. Also, ADC dynamic energy is smaller in the RRAM and FeFET cases. We observe that the amount of the total dynamic energy differences are similar with the one of ADC dynamic energy. This indicates that RRAM and FeFET energy efficiently access ADC than SRAM decreasing the total dynamic energy. As we use same the ResNet-50 model with same weight mapping and data transfer, the performance difference is mainly determined by memory device characteristics. The read dynamic energy of ADC is related to the resistance of the reading column in subarrays. Thus, the column resistance of RRAM and FeFET subarrays will be much smaller than SRAM subarrays.

The proportion of the ADC dynamic energy over the total dynamic energy in the RRAM and FeFET cases are small, but EGQ assumes ADC consumes most of the dynamic energy. However, the RRAM and FeFET cases also show better energy efficiency in the energy-aware cases. This is because the number of ADC access is related to the number of VMM operations in PIM. As read peripheral operations including ADC access are always followed by VMM operations, the energy-aware cases indirectly reduce the number of VMM operations and access to the other read peripheral circuits. The reduced total dynamic energy helps increasing the energy efficiency in the energy-aware cases, but total leakage energy

also plays a significant role in the energy efficiency for large models such as ResNet-50. eNVMs can significantly reduce the total leakage energy. For these reasons, Fig. 10(a) shows significant increases in the energy efficiency with RRAM and FeFET compared to the amount of the dynamic energy improvement.

Fig. 10(b) shows the dynamic energy and energy efficiency of the energy-aware cases for ResNet-50 with multi-bit memory devices. Total dynamic energy and ADC dynamic energy is significantly reduced in 2-bit RRAM and FeFET. This will be due to the lower number of required subarrays. As 2-bit memory devices use the half number of columns in subarrays, the number of required subarrays can be almost half of one for 1-bit memory devices. The number of ADC access is proportional to the number of subarrays. Hence, 2-bit memory devices achieve approximately two-times higher energy efficiency compared to 1-bit memory devices. Multi-bit eNVMs benefit not only the footprint of PIM platforms but also significantly the energy efficiency gain of EGQ. The key observations in the ResNet experiments are summarized:

- The number of ADC access is an effective approach to predict dynamic energy consumption in PIM.
- Including  $C_{ADC}$  in fitness function decreases the number of ADC access improving overall energy efficiency.
- Including  $C_W$  and  $C_A$  in fitness function reduces weight and activation bitwidth.
- EGQ is effective in SRAM, RRAM, and FeFET-based PIM

### C. VGGNet Quantization Evaluation

Fig. 11(a) and (b) show the comparison between the energy-unaware and energy-aware cases of VGG-19. EGQ shows the consideration of  $\gamma$  is also effective in VGG-19.

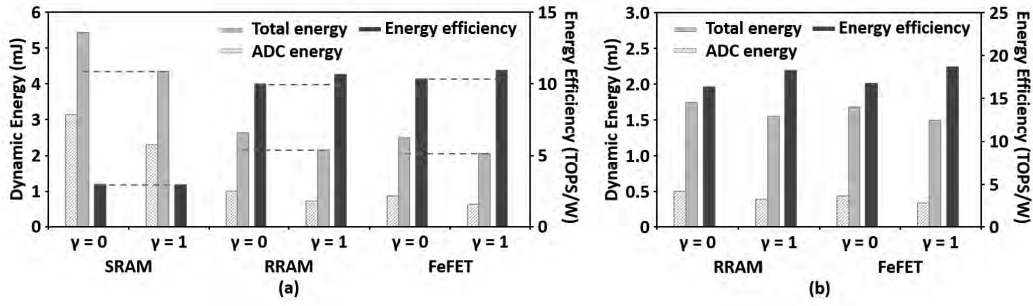


Fig. 12. (a) shows the dynamic energy and energy efficiency of VGG-19 with single-bit SRAM, RRAM, and FeFET. (b) shows the dynamic energy and energy efficiency of VGG-19 with 2-bit RRAM, and FeFET. Energy-unaware and energy-aware cases are compared with the memory devices.

TABLE IV  
ENERGY EFFICIENCY COMPARISON BETWEEN VGG-19 WITH  
AND WITHOUT FULLY-CONNECTED LAYERS

Model	VGG-19		VGG-19 without FC	
Parameter	$\gamma = 0$	$\gamma = 1$	$\gamma = 0$	$\gamma = 1$
Normalized Number of ADC Access	0.28	0.22	-	-
Weight Bitwidth	4.4	4.1	-	-
Activation Bitwidth	7.0	6.9	-	-
Total Dynamic Energy ( $\mu$ J)	1801.1	1403.4	1796.6	1258.1
ADC Dynamic Energy ( $\mu$ J)	1088.0	795.0	1086.0	792.0
Total Leakage Energy ( $\mu$ J)	8160.8	9510.0	1662.3	766.4
Energy Efficiency (TOPS/W)	3.9	3.6	11.3	19.3
Area ( $\text{mm}^2$ )	101.4	94.9	46.8	39.4
Baseline Accuracy (%)	72.4	72.4	-	-
Compressed Accuracy (%)	71.2	71.1	-	-

TABLE V  
LAYER-WISE BITWIDTH OF VGG-19 COMPRESSED BY ENERGY-AWARE  
AND ENERGY-UNAWARE QUANTIZATION

VGG-19	$\gamma = 0$	Weight	12,8,10,9,5,8,14,6,15,9,13,9,13,13,13,8,3,4,8
		Activation	12,5,13,8,7,6,4,6,12,9,6,7,9,12,14,7,9,9,11
	$\gamma = 1$	Weight	14,4,7,7,11,7,7,13,11,6,4,9,7,13,11,11,3,4,7
		Activation	13,7,8,7,5,6,5,7,15,6,5,6,12,5,12,5,14,6,14

Fig. 11(c) shows the compression-unaware case effectively reduces the number of ADC access at the early stage of iterations. However, the saturated level of the number of ADC access is similar with the compression-aware case at the 100<sup>th</sup> iteration. The average bitwidth of weight and activation is similar to the ResNet cases. Weight bitwidth is much smaller in the compression-aware case, but the difference in activation bitwidth is marginal.

Table IV summarizes the energy comparison between the two cases including total leakage energy. Layer-wise weight and activation bitwidth are summarized in Table V. We observe that VGG-19 consumes tremendous leakage energy. The total leakage energy of VGG-19 is much larger as fully-connected layers have way more weight parameters than the ResNet models. Specifically, the first fully-connected layer in VGG-19 has  $25088 \times 4096$  parameters, but one in ResNet-18 has  $512 \times 1000$ . Table IV shows the energy efficiency of the energy-aware case is slightly smaller than the energy-unaware case. As EGQ does not consider leakage energy, unexpected and too high leakage energy introduces the error to the energy efficiency.

We also compare the quantized VGG-19 and quantized VGG-19 without fully-connected layers in Table IV to clarify the impact of fully-connected layer to EGQ. VGG-19 without

FC in Table IV uses same models with the energy-aware and energy-unaware case, but the fully-connected layers are eliminated. Table IV shows that the energy efficiency is improved by 70.8% in the energy-aware case without the fully-connected layers. The energy efficiency tendency is flipped mainly because the total leakage energy is largely reduced. There is no meaningful change in the total dynamic energy as the length of input vectors is 1 in fully-connected layers indicating the number of ADC access is negligible compared to convolution layers.

Fig. 12(a) shows the comparison between the energy-unaware and energy-aware cases with various memory devices. We use the same 22nm technology node and hardware setting in Fig. 10. SRAM cases show energy-aware case has no longer lower energy efficiency than energy-unaware case. The energy efficiency in the SRAM cases is both 3.0 TOPS/W. Compared to the 7nm technology node in Table IV, we observe slightly smaller energy efficiency due to the larger dynamic energy of 22nm technology node. However, the SRAM cases with the 22nm technology node shows lower total leakage energy than the 7nm technology node. The RRAM and FeFET cases also decrease the total leakage energy significantly showing better the energy efficiency in energy-aware cases. By using eNVMs, we expect to reduce EGQ error to energy efficiency in VGG-like models.

Fig. 12(b) shows the comparison between the energy-unaware and energy-aware cases with multi-bit RRAM and FeFET. As multi-bit eNVMs can effectively reduce the number of subarrays, we observe much lower dynamic energy and leakage energy. Energy efficiency is improved by 12.3% in RRAM and 11.3% in FeFET. Large models including VGG-19 require the large number of subarrays to store them in on-chip memory, that increases leakage energy more than dynamic energy. Multi-bit eNVMs can be promising solution to effectively reduce the footprint of subarrays and the total leakage energy. As recent CNN models do not have really large fully-connected layers, we expect EGQ will not experience such error from the leakage energy in modern CNN models. The key observations in the VGGNet experiments are summarized:

- Considering ADC access and weight/activation bitwidth in the fitness function, reduces dynamic energy in PIM for VGG-19, similar to the ResNet models.
- Large fully-connected layers in VGG-19 like network and appreciable leakage energy of SRAM reduce overall energy-efficiency of EGQ.
- eNVMs mitigates the error and makes EGQ effective by decreasing leakage energy.

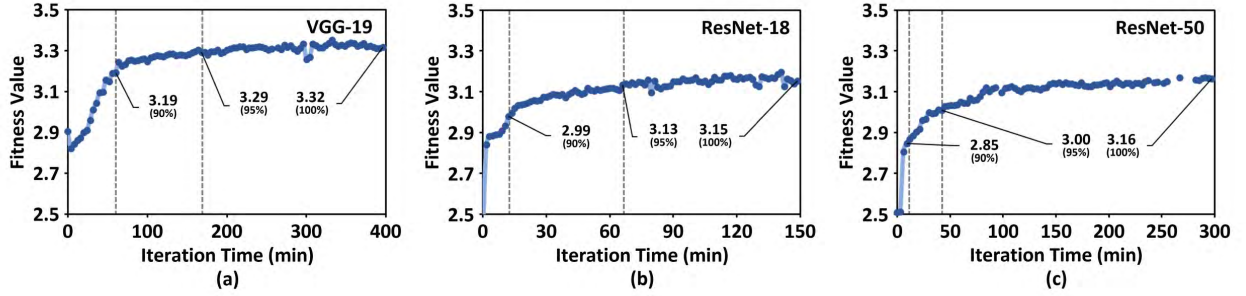


Fig. 13. Figures show the changes of fitness value during EGQ iterations. Energy-aware cases are used for VGG-19 (a), ResNet-18 (b), and ResNet-50 (c). The range of iteration time in three models is different, but the total number of EGQ iteration is same as 100. 90% and 95% of final fitness value are marked.

#### D. Run-Time of EGQ

We observe one EGQ iteration takes 1.5 minutes for ResNet-18, 3.0 minutes for ResNet-50, and 4.0 minutes for VGG-19 in GTX 1080Ti and i7-7700K environment. Fig. 13 shows saturation behavior of fitness values of three models. Fitness value abruptly decreases in VGG-19 and ResNet-50 sometimes because we set large penalty if the best candidate does not satisfy accuracy threshold. This time is mainly consumed at candidate evaluation step. Accuracy cost in the fitness function is estimated by average accuracy of 3,000 test images. The number of test data is heuristically determined, and it can be controlled depending on the required run-time. We can reduce the amount of the test data so that decrease the inference time proportionally. However, too small test dataset will not be desirable to correctly evaluate the average accuracy.

#### E. Balancing Weighting Factors in Fitness Function

Weighting factors ( $\alpha$ ,  $\beta$ , and  $\gamma$ ) reflect the relative importance of weight compression, activation compression, and ADC access. We determine the weighting factors as 0 or 1 in the previous section only to consider compression awareness and energy awareness. However, the overall energy efficiency can be affected by both weight, activation compression and ADC access. For example in Table IV, we observe large fully-connected layers in VGG models lead to tremendous leakage energy. In this case, higher  $\alpha$ ,  $\beta$  will be desirable as it can increase compression ratio and decrease leakage energy.

We explore other weighting factors for VGG-19, ResNet-18, and ResNet-50. They are determined based on the ratio of leakage energy and read dynamic energy of the 16-bit models. Specifically, we assume that the desirable ratio of the weighting factors follows leakage energy: read dynamic energy =  $\alpha$ :  $\gamma$  =  $\beta$ :  $\gamma$ . We constrain the sum of the weighting factors to 3 as the sum in initial energy-aware case is also 3. For example, 16-bit ResNet-50 shows 1.01mJ for leakage energy and 3.35mJ for read dynamic energy in 7nm SRAM based PIM. We set  $\alpha$ ,  $\beta$ , and  $\gamma$  for the model to 1.3, 1.3, and 0.4 as 3.35: 1.01 (mJ) = 1.3 ( $\alpha$ ): 0.4 ( $\gamma$ ) = 1.3 ( $\beta$ ): 0.4 ( $\gamma$ ). The weighting factors of VGG-19 and ResNet-18 can be similarly balanced. Table VI shows the energy efficiency of the three models with balanced weighting factors. We observe that VGG-19 shows slightly higher ADC dynamic energy while lower total leakage energy compared to energy-aware case in Table IV. ResNet-50 also shows similar results with VGG-19. However, ResNet-18 shows lower ADC dynamic energy with a bit higher total leakage energy. All three models prove that balancing weighting factors can further

TABLE VI  
ENERGY EFFICIENCY OF VGGNET, RESNET WITH  
BALANCED WEIGHTING FACTORS

Model	VGG-19	ResNet-18	ResNet-50
Parameter ( $\alpha$ , $\beta$ , $\gamma$ )	1.49, 1.49, 0.02	0.86, 0.86, 1.28	1.3, 1.3, 0.4
Weight Bitwidth	4.2	7.1	8.0
Activation Bitwidth	6.2	6.8	7.8
Total Dynamic Energy ( $\mu$ J)	1422.0	104.8	308.1
ADC Dynamic Energy ( $\mu$ J)	824.20	56.9	144.0
Total Leakage Energy ( $\mu$ J)	7263.1	43.0	526.9
Energy Efficiency (TOPS/W)	4.5	36.1	12.3
Baseline Accuracy (%)	72.4	69.8	76.2
Compressed Accuracy (%)	70.9	67.5	74.2

improve the overall energy efficiency compared with the energy-aware cases.

#### F. Comparison With Q-PIM

We compare EGQ and Q-PIM with hardware performance in this section [7]. Table VII shows the hardware performance evaluation of ResNet-18 with ImageNet dataset between EGQ and Q-PIM. We simulate the performance with various memory technologies such as SRAM, RRAM, and FeFET in the 22nm technology node. Table II shows 7.1 bit for the weight and 7.6 bit for the activation in energy-aware case for ResNet-18. However, we can further optimize the bitwidth by decreasing the upper bound in the initialization step and by the GS algorithm based fine tuning. EGQ achieves 6.1 bit for the weight and 6.3 bit for the activation for ResNet-18 to compare with previous works.

Fig. 14 shows the layer-wise weight and activation bitwidth of ResNet-18 from EGQ and Q-PIM. Fig. 15 shows the normalized number of ADC access of both which is calculated based on equation (4). The distribution of the weight and activation bitwidth is similar for both. However, the slight differences in the activation bitwidth at early layers result in the large differences in the number of ADC access. For example, the weight bitwidth is same, and the activation bitwidth is 2-bit difference in the first layer. As the highest number of ADC access is observed in the first layer, EGQ saves the total number of ADC access most from the first layer. Also the weight bitwidth of the 6<sup>th</sup> and 14<sup>th</sup> layers is higher in EGQ, but the difference of the number of ADC access at these layers are mitigated due to lower activation bitwidth.

TABLE VII  
RESNET-18 ENERGY EFFICIENCY COMPARISON BETWEEN EGQ  
AND Q-PIM WITH VARIOUS MEMORY DEVICES

Technology Node (LP)	22nm		22nm		22nm	
Memory Device	SRAM		RRAM		FeFET	
Cell Precision	1 bit		1 bit		1 bit	
Algorithm	EGQ	Q-PIM	EGQ	Q-PIM	EGQ	Q-PIM
Weight Bitwidth	6.1	5.2	6.1	5.2	6.1	5.2
Activation Bitwidth	6.3	8	6.3	8	6.3	8
Total Dynamic Energy ( $\mu$ J)	207.9	237.8	134.7	154.9	129.8	149.6
ADC Dynamic Energy ( $\mu$ J)	102.4	111.9	44.6	48.8	39.8	43.7
Total Leakage Energy ( $\mu$ J)	19.8	20.6	4.3	4.4	4.3	4.4
Energy Efficiency (TOPS/W)	23.4	20.7	38.4	33.5	39.8	34.7
Latency (ms/image)	23.0	27.9	19.5	23.8	19.5	23.8
Area ( $\text{mm}^2$ )	75.7	66.3	51.3	44.7	51.3	44.7
Baseline Accuracy (%)	69.8	69.8	69.8	69.8	69.8	69.8
Compressed Accuracy (%)	67.5	63.8	67.5	63.8	67.5	63.8

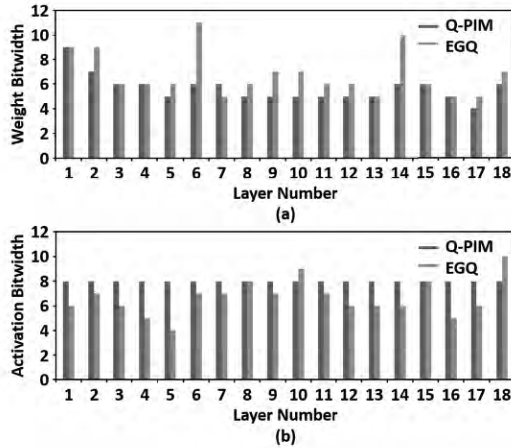


Fig. 14. The layer-wise weight and activation bitwidth for ResNet-18 from Q-PIM and EGQ. (a) shows the layer-wise weight bitwidth. (b) shows the layer-wise activation bitwidth.

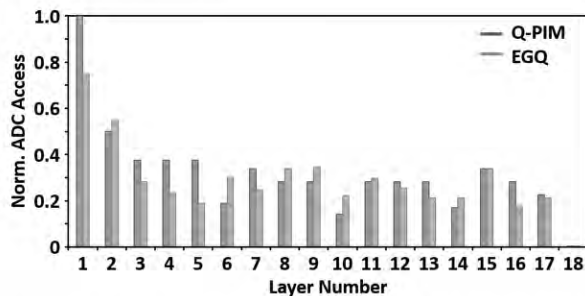


Fig. 15. The layer-wise number of ADC access for ResNet-18 from Q-PIM and EGQ.

Table VII shows the comparison between EGQ and Q-PIM with various memory devices. As EGQ has higher average weight bitwidth, we observe the chip area for EGQ is about 15% larger than Q-PIM. However, EGQ achieves 13.0%, 14.6%, and 14.7% higher energy efficiency than Q-PIM in

TABLE VIII  
COMPARISON AGAINST VARIOUS PIM-AWARE  
QUANTIZATION ALGORITHMS

Algorithm	EGQ	EGQ	[14]	[13]	[12]
Model	ResNet-18	ResNet-18	ResNet-18	ResNet-18	ResNet-18
Dataset	ImageNet	CIFAR-10	CIFAR-100	CIFAR-10	CIFAR-10
Weight Bitwidth	<b>6.1</b>	<b>2.5</b>	flexible	4	2
Activation Bitwidth	<b>6.3</b>	<b>3.9</b>	flexible	4	4
Energy Efficiency (TOPS/W)	<b>6.5×</b>	<b>18.6×</b>	3.2×	17.7×	29.8×
Baseline Accuracy (%)	<b>69.8</b>	<b>88.6</b>	70.9	88.9	89.1
Compressed Accuracy (%)	<b>67.5</b>	<b>86.4</b>	70.5	86.6	81.2

SRAM, RRAM, and FeFET, respectively. Total leakage energy is relatively negligible in the RRAM and FeFET cases consequently showing higher energy efficiency than SRAM. Latency per image is reduced by 17.6%, 18.1%, and 18.1% compared to Q-PIM in SRAM, RRAM, and FeFET, respectively. EGQ is not designed to consider the latency, but EGQ considers the number of ADC access that is related to the length of input vectors. Approximately 70% of the latency is due to interconnection indicating expensive input data transfer. As long input vectors require more data transfer, the number of data transfer is indirectly optimized by EGQ. Also, EGQ shows only 2.3% accuracy loss compared to the baseline accuracy from Pytorch pretrained model, but 6.0% accuracy loss is observed in Q-PIM.

#### G. Comparison With Prior Works

Table VIII compares EGQ with recent quantization algorithms for PIM [12]–[14]. A flexible bitwidth quantization algorithm is used in the paper by Vasquez *et al.* [14], and the others apply uniform bitwidth quantization algorithms [12], [13]. All algorithms are based on ResNet-18, but the dataset for test accuracy are different. Energy efficiency in the table means the improvement of the energy efficiency compared to ResNet-18 with 16-bit precision. EGQ achieves 6.5× higher energy efficiency compared to the 16-bit model with ImageNet dataset. Vasquez *et al.* shows small accuracy loss, but the improvement of energy efficiency is also small. The uniform quantization algorithms [12], [13] show very low bit quantization with CIFAR-10 dataset. Sun *et al.* [13] achieves 4 bit for both weight and activation. Cai *et al.* [12] shows more compressed model with 2-bit weight and 4-bit activation. However, huge accuracy loss of 7% is observed [12]. EGQ compresses weight and activation to 2.5 bit and 3.9 bit without huge loss of accuracy. Also, EGQ achieves lower energy consumption and memory size compared with the algorithm in [13].

We also compare EGQ with five different quantization methods in ResNet-50. Energy efficiency for each method is estimated in 7nm SRAM based PIM by NeuroSim. There are several well-known works on extremely low-precision quantization such as Dorefa [4] and PACT [5]. We observe that 2-bit weight and activation by Dorefa achieves the best energy efficiency, but it shows highest accuracy loss (9.8%). PACT shows 4-bit weight and activation with only 0.4% of accuracy loss with high energy efficiency. We expect that models compressed by these methods will show higher energy

TABLE IX  
COMPARISON AGAINST VARIOUS QUANTIZATION  
ALGORITHMS FOR RESNET-50

Algorithm	Re-training	Bitwidth (W, A)	Baseline Accuracy (%)	Compressed Accuracy (%)	Energy Efficiency (TOPS/W)
Dorefa [4]	One-pass	2, 2	76.9	67.1	257.1
PACT [5]	One-pass	4, 4	76.9	76.5	83.0
OMSE [29]	No need	4, 32	76.0	75.0	9.3
OCS [30]	No need	8, 7	76.1	74.5	14.5
EGQ (2%)	No need	6.8, 7.2	76.2	74.4	25.3
ACIQ [31]	No need	8, 4	76.1	71.5	25.9
EGQ (5%)	No need	6.2, 6.3	76.2	71.5	31.7

efficiency as compression ratio is generally much higher than EGQ. However, most of them require at least one-pass of re-training, large training dataset, or handcrafted tuning to reduce accuracy loss. For this reason, we mainly compare EGQ with re-training free quantization algorithms such as OMSE [29], OCS [30], and ACIQ [31].

OMSE achieves 4-bit weight and 32-bit activation. It successfully reduces the large leakage energy of ResNet-50, but 32-bit activation introduces high read dynamic energy. Energy efficiency of OMSE is 9.3TOPS/W which is the lowest among the algorithms. OCS achieves 8-bit weight and 7-bit activation, and it shows relatively balanced leakage energy and read dynamic energy. Hence, higher energy efficiency of 14.5TOPS/W is observed. EGQ shows 6.8-bit weight and 7.2-bit activation. Though average activation bitwidth is a bit higher than OCS, we observe that lower and more balanced leakage energy and read dynamic energy. EGQ achieves 25.3TOPS/W of higher energy efficiency with almost same accuracy compared to OMSE and OCS. We further investigate the performance of EGQ with ACIQ. ACIQ shows 8-bit weight and 4-bit activation with higher accuracy loss (4.6%). Energy efficiency of ACIQ is 25.9TOPS/W. We set the 5% of accuracy threshold for EGQ to achieve similar accuracy loss. We observe overall energy efficiency is improved by decreasing the accuracy threshold. EGQ shows 31.7TOPS/W energy efficiency and same accuracy with ACIQ. EGQ can automatically and efficiently find the optimal weight and activation bitwidth for energy efficient PIM.

#### H. Performance of EGQ on Low-Complexity Networks

We study the performance of EGQ with MobileNet-V2 [32] and SqueezeNet [33]. We also compare with the performance of several other quantization algorithms on these networks [29], [34], [35]. Table X summarizes the comparison for the two models. Hardware architecture is assumed to be 7nm SRAM based PIM same as preceding experiments. Normalized number of ADC access represents the number of ADC access of the compressed model divided by the one of 16-bit model. As NeuroSim does not support these models, we approximately calculate the total leakage energy and total read dynamic energy for the energy efficiency based on the energy model of NeuroSim.

For MobileNet-V2, Integer-only compresses weight and activation more than EGQ and Per-channel, and energy efficiency is also higher. However, it requires retraining while Per-channel and EGQ are retraining free quantization algorithms. EGQ shows lower weight and activation bitwidth than Per-channel with same accuracy. Also, normalized number of ADC access is lower than Per-channel, hence we observe

TABLE X  
COMPARISON AGAINST VARIOUS QUANTIZATION ALGORITHMS  
FOR MOBILENET-V2 AND SQUEEZENET

Algorithm	Integer-only [34]	Per-channel [35]	EGQ	OMSE [29]	EGQ
Model	MobileNet-V2		SqueezeNet		
Weight Bitwidth	6	8	7.8	4	6.7
Activation Bitwidth	6	8	7.5	32	7.9
Baseline Accuracy (%)	-	71.9	71.8	58.0	58.0
Compressed Accuracy (%)	70.9	69.7	69.7	55.4	55.5
Baseline Size (MB)	13.4	13.4	13.4	4.8	4.8
Compressed Size (MB)	2.5	3.3	3.3	0.6	1.0
Normalized Number of ADC Access	0.16	0.25	0.23	0.47	0.26
Energy Efficiency (TOPS/W)	39.7	23.1	24.5	11.9	32.2

higher energy efficiency in EGQ. EGQ shows relatively slight improvement of energy efficiency compared to ResNet-50. We observe that early layers in MobileNet-V2 is relatively more sensitive to quantization errors compared to VGGNet and ResNet models. Though early layers consumes more number of ADC access, EGQ does not assign lower bitwidth to the layers to avoid the huge accuracy loss, and it makes relatively small improvement of energy efficiency in MobileNet-V2. As EGQ is based on linear quantization scheme with maximum clipping threshold to simplify the algorithm, we expect that more advanced clipping method or non-linear quantization scheme will enhance the performance of EGQ in MobileNet like models. For SqueezeNet, OMSE shows 4-bit weight quantization with floating-point activation. EGQ achieves 6.7-bit weight and 7.9-bit activation with similar accuracy. Similar to the comparison between OMSE and EGQ in Table IX, we observe that EGQ shows more balanced leakage energy and read dynamic energy. Also, EGQ achieves much lower number of ADC access with higher energy efficiency. Effectiveness of EGQ over the energy efficiency improvement can be variable depending on the models, but EGQ is still effective in MobileNet-V2 and SqueezeNet.

#### I. Energy Distribution of Different Components

Energy distribution of hardware components with different quantization methods in Fig. 16 further explains how EGQ reduces the energy consumption in PIM. We choose three different methods, Q-PIM for ResNet-18, OCS for ResNet-50, and OMSE for SqueezeNet, where EGQ benefits the overall energy efficiency most in previous sections. EGQ reduces ADC energy consumption in all three methods as expected in Fig. 16. However, it is important to note that the rate of energy consumption decrease depends on the hardware components. For example, the decrease of leakage energy or accumulation circuits energy are more distinct compared to ADC energy in Fig. 16(b), on the other hand, EGQ achieves higher energy efficiency mainly due to ADC and IC energy in Fig. 16(c). Q-PIM and EGQ show almost similar energy distribution. This implies that energy consumption of different hardware components has different sensitivity to weight and activation precision depending the models. As energy-unaware methods

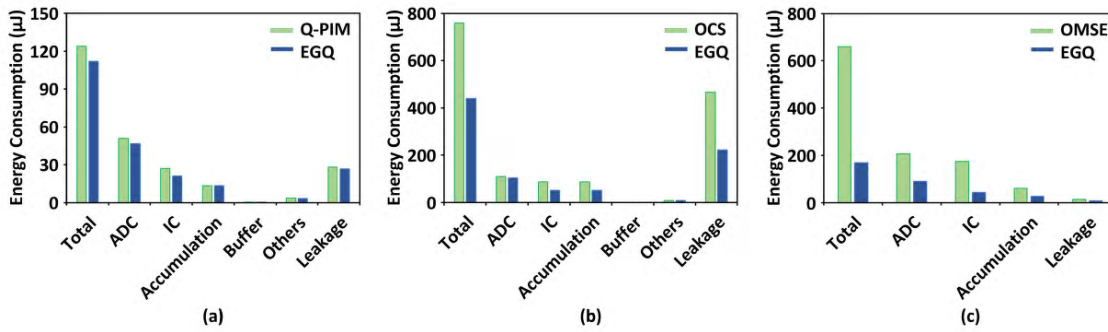


Fig. 16. Figures show the distribution of energy consumption in different hardware components and leakage energy. (a) compares Q-PIM and EGQ in ResNet-18, (b) compares OCS and EGQ in ResNet-50, and (c) compares OMSE and EGQ in SqueezeNet. 7nm SRAM based PIM design is used to estimate the energy consumption.

does not consider the different sensitivity, it can quantize weight more and eventually give us higher and suboptimal ADC energy. It explains how 6.7-bit weight and 7.9-bit in EGQ achieves higher energy efficiency compared to 4-bit weight and 32-bit activation in OMSE. Only lowering the activation bitwidth does not give us optimal energy efficiency as well, for example, OCS has lower activation bitwidth with lower energy efficiency. These results support that optimizing bit precision only with the constraint of accuracy does not guarantee the optimal energy efficiency.

## VI. CONCLUSION

This paper presents a genetic algorithm based energy-aware quantization method (EGQ) to run CNN models energy efficiently on PIM platforms. We discuss how to combine energy efficiency to a genetic algorithm based quantization method. EGQ predicts the dynamic energy consumption for inference based on the number of ADC access. ResNet-18, ResNet-50, and VGG-19 are used to evaluate the effectiveness of EGQ. We observe the ResNet models are more suitable for EGQ, and VGG-19 shows unexpected behavior due to too high leakage energy from large fully-connected layers. However, the error in EGQ introduced by the leakage energy is mitigated by eNVMs, consequently showing better energy efficiency. As most of modern CNN models do not have such large fully-connected layers, EGQ can effectively search quantization bitwidth of various CNN models for energy-efficient PIM. For future works, the consideration for more hardware details such as memory device types and ADC precision will further improve the energy efficiency of EGQ. We use a genetic algorithm due to the flexibility of the fitness function design. Other problem solvers can be adapted for future studies. Also, understanding the effect of device variation or non-ideality on mixed-precision design will be explored in the future.

## REFERENCES

- [1] A. Shafiee *et al.*, "ISACC: A convolutional neural network accelerator with *in-situ*, analog arithmetic in crossbars," *ACM SIGARCH Comput. Archit. News*, vol. 44, no. 3, pp. 14–26, 2016.
- [2] P. Chi *et al.*, "PRIME: A novel processing-in-memory architecture for neural network computation in rram-based main memory," *ACM SIGARCH Comput. Archit. News*, vol. 44, no. 3, pp. 27–39, Jun. 2016.
- [3] C. Eckert *et al.*, "Neural cache: Bit-serial in-cache acceleration of deep neural networks," in *Proc. ACM/IEEE 45th Annu. Int. Symp. Comput. Archit. (ISCA)*, Jun. 2018, pp. 383–396.
- [4] S. Zhou, Y. Wu, Z. Ni, X. Zhou, H. Wen, and Y. Zou, "DoReFa-net: Training low bitwidth convolutional neural networks with low bitwidth gradients," 2016, *arXiv:1606.06160*.
- [5] J. Choi, S. Venkataramani, V. Srinivasan, K. Gopalakrishnan, Z. Wang, and P. Chuang, "Accurate and efficient 2-bit quantized neural networks," in *Proc. SysML*, 2019, pp. 1–5.
- [6] K. Wang, Z. Liu, Y. Lin, J. Lin, and S. Han, "HAQ: Hardware-aware automated quantization with mixed precision," in *Proc. Comput. Vis. Pattern Recognit. (CVPR)*, Jun. 2019, pp. 8612–8620.
- [7] Y. Long, E. Lee, D. Kim, and S. Mukhopadhyay, "Q-PIM: A genetic algorithm based flexible DNN quantization method and application to processing-in-memory platform," in *Proc. 57th ACM/IEEE Design Autom. Conf. (DAC)*, Jul. 2020, pp. 1–6.
- [8] Y. Zhou, S. M. Moosavi-Dezfooli, N. M. Cheung, and P. Frossard, "Adaptive quantization for deep neural network," in *Proc. AAAI Conf. Artif. Intell.*, 2018, pp. 4596–4604.
- [9] M. Rusci, A. Capotondi, and L. Benini, "Memory-driven mixed low precision quantization for enabling deep network inference on micro-controllers," 2019, *arXiv:1905.13082*.
- [10] C. Ding, S. Wang, N. Liu, K. Xu, Y. Wang, and Y. Liang, "REQ-YOLO: A resource-aware, efficient quantization framework for object detection on FPGAs," in *Proc. ACM/SIGDA Int. Symp. Field-Program. Gate Arrays*, Feb. 2019, pp. 33–42.
- [11] T. J. Yang, Y. H. Chen, and V. Sze, "Designing energy-efficient convolutional neural networks using energy-aware pruning," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, Jan. 2017, pp. 5687–5695.
- [12] Y. Cai, T. Tang, L. Xia, B. Li, Y. Wang, and H. Yang, "Low bit-width convolutional neural network on RRAM," *IEEE Trans. Comput.-Aided Design Integr. Circuits Syst.*, vol. 39, no. 7, pp. 1414–1427, Jul. 2020.
- [13] H. Sun, Z. Zhu, Y. Cai, X. Chen, Y. Wang, and H. Yang, "An energy-efficient quantized and regularized training framework for processing-in-memory accelerators," in *Proc. 25th Asia South Pacific Design Autom. Conf. (ASP-DAC)*, Jan. 2020, pp. 325–330.
- [14] K. Vasquez, Y. Venkatesha, A. Bhattacharjee, A. Moitra, and P. Panda, "Activation density based mixed-precision quantization for energy efficient neural networks," 2021, *arXiv:2101.04354*.
- [15] Z. Zhu *et al.*, "A configurable multi-precision CNN computing framework based on single bit RRAM," in *Proc. 56th Annu. Design Autom. Conf.*, Jun. 2019, pp. 1–6.
- [16] S. Huang *et al.*, "Mixed precision quantization for ReRAM-based DNN inference accelerators," in *Proc. 26th Asia South Pacific Design Autom. Conf.*, Jan. 2021, pp. 372–377.
- [17] X. Peng, S. Huang, Y. Luo, X. Sun, and S. Yu, "DNN+NeuroSim: An end-to-end benchmarking framework for compute-in-memory accelerators with versatile device technologies," in *IEDM Tech. Dig.*, Dec. 2019, p. 32.
- [18] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," 2014, *arXiv:1409.1556*.
- [19] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, Dec. 2016, pp. 770–778.
- [20] S. Qu, B. Li, Y. Wang, D. Xu, X. Zhao, and L. Zhang, "RaQu: An automatic high-utilization CNN quantization and mapping framework for general-purpose RRAM accelerator," in *Proc. 57th ACM/IEEE Design Autom. Conf. (DAC)*, Jul. 2020, pp. 1–6.
- [21] C. Zhang and P. Zhou, "A quantized training framework for robust and accurate ReRAM-based neural network accelerators," in *Proc. 26th Asia South Pacific Design Autom. Conf.*, Jan. 2021, pp. 43–48.

- [22] K. Samal, M. Wolf, and S. Mukhopadhyay, "Attention-based activation pruning to reduce data movement in real-time AI: A case-study on local motion planning in autonomous vehicles," *IEEE J. Emerg. Sel. Topics Circuits Syst.*, vol. 10, no. 3, pp. 306–319, Sep. 2020.
- [23] S. Yin, Z. Jiang, J.-S. Seo, and M. Seok, "XNOR-SRAM: In-memory computing SRAM macro for binary/ternary deep neural networks," *IEEE J. Solid-State Circuits*, vol. 55, no. 6, pp. 1733–1743, Jun. 2020.
- [24] Y. Long *et al.*, "A ferroelectric FET based power-efficient architecture for data-intensive computing," in *Proc. Int. Conf. Computer-Aided Design*, Nov. 2018, pp. 1–8.
- [25] X. Si *et al.*, "24.5 A twin-8T SRAM computation-in-memory macro for multiple-bit CNN-based machine learning," in *IEEE Int. Solid-State Circuits Conf. (ISSCC) Dig. Tech. Papers*, Feb. 2019, pp. 396–398.
- [26] X. Peng, R. Liu, and S. Yu, "Optimizing weight mapping and data flow for convolutional neural networks on processing-in-memory architectures," *IEEE Trans. Circuits Syst. I, Reg. Papers*, vol. 67, no. 4, pp. 1333–1343, Apr. 2020.
- [27] W. Wu *et al.*, "A methodology to improve linearity of analog RRAM for neuromorphic computing," in *Proc. IEEE Symp. VLSI Technol.*, Jun. 2018, pp. 103–104.
- [28] K. Ni *et al.*, "In-memory computing primitive for sensor data fusion in 28 nm HKMG FeFET technology," in *IEDM Tech. Dig.*, Dec. 2018, pp. 1–16.
- [29] R. Zhao, Y. Hu, J. Dotzel, C. D. Sa, and Z. Zhang, "Improving neural network quantization without retraining using outlier channel splitting," in *Proc. Int. Conf. Mach. Learn.*, 2019, pp. 7543–7552.
- [30] C. Y. E. Kravchik, F. Yang, and P. Kisilev, "Low-bit quantization of neural networks for efficient inference," in *Proc. ICCV Workshops*, 2019, pp. 3009–3018.
- [31] R. Banner, Y. Nahshan, E. Hoffer, and D. Soudry, "ACIQ: Analytical clipping for integer quantization of neural networks," 2018, *arXiv:1810.05723*.
- [32] M. Sandler, A. Howard, M. Zhu, A. Zhmoginov, and L.-C. Chen, "MobileNetV2: Inverted residuals and linear bottlenecks," in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit.*, Jun. 2018, pp. 4510–4520.
- [33] D. Degrijse and J. Souto, "Dimension invariants of outer automorphism groups," 2016, *arXiv:1602.04354*.
- [34] B. Jacob *et al.*, "Quantization and training of neural networks for efficient integer-arithmetic-only inference," in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit.*, Jun. 2018, pp. 2704–2713.
- [35] R. Krishnamoorthi, "Quantizing deep convolutional networks for efficient inference: A whitepaper," 2018, *arXiv:1806.08342*.



**Beomseok Kang** received the B.S. degree in electronic and electrical engineering from Sungkyunkwan University, Suwon, South Korea, in 2020. He is currently pursuing the Ph.D. degree in electrical and computer engineering with the Georgia Institute of Technology, Atlanta, GA, USA.

His current research interests include efficient deep learning algorithms and deep learning acceleration in processing-in-memory architectures.



**Anni Lu** received the B.S. degree in electronic information engineering from Tianjin University in 2019. She is currently pursuing the Ph.D. degree in electrical and computer engineering with the Georgia Institute of Technology, Atlanta, GA, USA.

Her current research interests include deep learning algorithms and hardware co-design for in-memory computing and hardware accelerator.



**Yun Long** (Student Member, IEEE) received the B.S. degree in microelectronics from Peking University, Beijing, China, in 2014, and the Ph.D. degree in electrical and computer engineering from the Georgia Institute of Technology, Atlanta, GA, USA, in 2019.

He is currently working with Google, focusing on TPU architecture and hardware-oriented ML models optimization.



**Daehyun Kim** (Graduate Student Member, IEEE) received the B.S. degree in semiconductor systems engineering from Sungkyunkwan University, Suwon, South Korea, in 2016. He is currently pursuing the Ph.D. degree in electrical and computer engineering with the Georgia Institute of Technology, Atlanta, GA, USA.

After graduation, he was with Samsung Electronics Company Ltd., Hwasung, South Korea, for two years. His current research interests include machine learning accelerator design and volatile/non-volatile memory designs with emerging technologies.



**Shimeng Yu** (Senior Member, IEEE) received the B.S. degree in microelectronics from Peking University in 2009, and the M.S. and Ph.D. degrees in electrical engineering from Stanford University in 2011 and 2013, respectively.

From 2013 to 2018, he was an Assistant Professor with Arizona State University. He is currently an Associate Professor of electrical and computer engineering with the Georgia Institute of Technology. His research expertise is on the emerging non-volatile memories for applications, such as deep

learning accelerator, in-memory computing, 3D integration, and hardware security. Among his honors, he was a recipient of the NSF Faculty Early CAREER Award in 2016, the IEEE Electron Devices Society (EDS) Early Career Award in 2017, the ACM Special Interests Group on Design Automation (SIGDA) Outstanding New Faculty Award in 2018, the Semiconductor Research Corporation (SRC) Young Faculty Award in 2019, the ACM/IEEE Design Automation Conference (DAC) Under-40 Innovators Award in 2020, and the IEEE Circuits and Systems Society (CASS) Distinguished Lecturer for the term 2021–2022.



**Saibal Mukhopadhyay** (Fellow, IEEE) received the B.E. degree in electronics and telecommunication engineering from Jadavpur University, Kolkata, India, in 2000, and the Ph.D. degree in electrical and computer engineering from Purdue University, West Lafayette, IN, USA, in 2006.

He is currently a Joseph M. Pettit Professor with the School of Electrical and Computer Engineering, Georgia Institute of Technology, Atlanta, GA, USA. He has authored or coauthored more than 200 papers in refereed journals and conferences, and holds five

U.S. patents. His research interests include design of energy-efficient, intelligent, and secure systems. His research explores a cross-cutting approach to design spanning algorithm, architecture, circuits, and emerging technologies. He was a recipient of the Office of Naval Research Young Investigator Award in 2012, the National Science Foundation CAREER Award in 2011, the IBM Faculty Partnership Award in 2009 and 2010, the SRC Inventor Recognition Award in 2008, the SRC Technical Excellence Award in 2005, and the IBM Ph.D. Fellowship Award for the period 2004–2005. He has received the IEEE TRANSACTIONS ON VERY LARGE SCALE INTEGRATION (VLSI) SYSTEMS Best Paper Award in 2014, the IEEE TRANSACTIONS ON COMPONENTS, PACKAGING AND MANUFACTURING TECHNOLOGY Best Paper Award in 2014, and multiple best paper awards in the International Symposium on Low power Electronics and Design in 2014, 2015, and 2016.