

OpenACiM: An Open-Source SRAM-Based Approximate CiM Compiler

Yiqi Zhou¹, JunHao Ma¹, Xingyang Li², Yule Sheng¹, Yue Yuan¹, Yikai Wang¹, Bochang Wang¹, Yiheng Wu¹, Shan Shen^{1,*}, Wei Xing^{3,*}, Daying Sun^{1,*}, Li Li¹ and Zhiqiang Xiao⁴

¹Nanjing University of Science and Technology, Nanjing, 210094, China

²Beihang University, Beijing, 100191, China

³University of Sheffield, Sheffield, S10 2TN, United Kingdom

⁴The 58th Research Institute of China Electronics Technology Group Corporation, Wuxi, 214035, China

Abstract—The rise of data-intensive AI workloads has exacerbated the “memory wall” bottleneck. Digital Compute-in-Memory (DCiM) using SRAM offers a scalable solution, but its vast design space makes manual design impractical, creating a need for automated compilers. A key opportunity lies in approximate computing, which leverages the error tolerance of AI applications for significant energy savings. However, existing DCiM compilers focus on exact arithmetic, failing to exploit this optimization. This paper introduces OpenACM, the first open-source, accuracy-aware compiler for SRAM-based approximate DCiM architectures. OpenACM bridges the gap between application error tolerance and hardware automation. Its key contribution is an integrated library of accuracy-configurable multipliers (exact, tunable approximate, and logarithmic), enabling designers to make fine-grained accuracy-energy trade-offs. The compiler automates the generation of the DCiM architecture, integrating a transistor-level customizable SRAM macro with variation-aware characterization into a complete, open-source physical design flow based on OpenROAD and the FreePDK45 library. This ensures full reproducibility and accessibility, removing dependencies on proprietary tools. Experimental results on representative convolutional neural networks (CNNs) demonstrate that OpenACM achieves energy savings of up to 64% with negligible loss in application accuracy. The framework is available on [OpenACM:URL](#).

I. INTRODUCTION

The growth of data-intensive Artificial Intelligence (AI) applications has intensified the “memory wall” bottleneck in traditional von Neumann architectures. The constant data transfer between memory and processing units severely limits both performance and energy efficiency [1]. Compute-in-Memory (CiM) directly addresses this challenge by performing computations within the memory array, thereby minimizing data movement.

Among CiM approaches, Digital Compute-in-Memory (DCiM), which integrates logic into standard SRAM arrays, has emerged as a scalable and practical solution [2], [3]. However, the architectural flexibility of DCiM creates a vast design space spanning choices in bit-cell design, arithmetic circuits, and array organization. Manually exploring this space is prohibitively slow, rendering comprehensive Design Space

Exploration (DSE) impractical for the rapid hardware iteration that AI applications demand. This challenge necessitates sophisticated compilers to automate the design flow from high-level specification to physical layout.

Concurrently, approximate computing provides a powerful technique for improving energy efficiency. Many AI applications possess an inherent tolerance to minor computational errors [4], [5]. By strategically introducing controlled approximations into the hardware, designers can trade negligible losses in application-level accuracy for significant reductions in power consumption and area.

Existing DCiM compilers primarily focus on *exact* arithmetic [6]–[10] and lack support for approximation, thereby failing to exploit the full energy-saving potential permitted by error-tolerant applications. Moreover, current tools do not provide precision-configurable DCiM architectures that can adapt to varying levels of application-specific error tolerance. The reliance of many compilers on proprietary EDA toolchains further limits reproducibility and inhibits broader community-driven development [11].

To address these limitations, this work presents **OpenACM**, an Open-source SRAM-based Approximate CiM compiler. OpenACM is the first open-source framework that systematically integrates accuracy-aware approximate computing into a fully automated DCiM design flow. Enabled by flexible bit-width and precision configurability, OpenACM provides a unified platform for exploring energy-efficient CiM circuits under different application-specific accuracy constraints. Our key contributions are:

- **Approximate DCiM compiler with accuracy-configurable multipliers:** OpenACM integrates a multiplier library offering three families selectable under application accuracy constraints: (i) an exact 4-2 compressor-based multiplier; (ii) an approximate 4-2 compressor-based multiplier with tunable accuracy; and (iii) a logarithmic approximate multiplier. This enables fine-grained accuracy-energy trade-offs within a unified flow.
- **Complete open-source ecosystem integration:** The back-end (physical) design flow is implemented using the OpenROAD digital design flow [12] with FreePDK45 [13] to ensure full reproducibility without

This work is supported by the Fundamental Research Funds for the Central Universities under grant Nos. 30924012004 and 30925010605, and by the National Key Laboratory of Integrated Circuits and Microsystems under grant No. JCYQ2310803-1.

* Corresponding authors.

proprietary dependencies. The SRAM macro follows a FakeRAM2.0-style [14] template for seamless macro abstraction and Physical Design (PD) integration [14].

- **Variation-Aware SRAM Analysis:** Built upon OpenYield [15], OpenACM supports transistor-level customization of SRAM macros and integrates importance-sampling-based Monte Carlo (MC) simulation to accelerate library characterization under Process Voltage Temperature (PVT) variations, together with automated transistor sizing optimizations.

The remainder of this paper is organized as follows: [Sec. II](#) reviews relevant background and related DCiM compilers; [Sec. III](#) presents the OpenACM framework and its core components; [Sec. IV](#) details the OpenROAD-based physical design flow; [Sec. V](#) reports post-layout results; and [Sec. VI](#) concludes with discussions and future directions. OpenACM is under continuous development, with several key future enhancements discussed in [Sec. VI](#).

II. BACKGROUND

A. Error-Tolerant Applications

The fundamental premise of this work rests on a powerful opportunity: the inherent error resilience of neural networks. Unlike traditional high-performance computing, DNNs are statistical models that learn to recognize patterns in noisy, real-world data. Their robustness is analogous to the human brain’s ability to identify a familiar face in a poorly lit photograph; perfect, high-fidelity input is not required for a correct outcome. This resilience translates directly to the hardware level. The millions of multiply-accumulate operations underpinning a network’s inference do not all require perfect mathematical accuracy. A small degree of computational error in circuits can often be introduced with little to no degradation in the final application accuracy. This presents a golden opportunity to trade this unneeded accuracy for dramatic savings in energy consumption and silicon area [16], [17].

B. Digital CiM Compilers

Automated DCiM compilers [2], [6]–[10] have significantly improved design productivity. However, none currently support configurable approximation or comprehensive yield-aware SRAM analysis. [Tab. I](#) highlights the key distinctions in precision, openness, and analysis depth. As detailed in the table, AutoDCIM [6] set the precedent for automation but is closed-source. Subsequent works like ARCTIC [7], SynDCIM [8], SEGA-DCIM [9], and DAMIL-DCIM [10] advanced specific aspects such as multi-precision support and physical layout, yet they remain proprietary and restricted to exact arithmetic. OpenC² [2] broke the open-source barrier but offers only basic SRAM analysis and lacks approximation features.

OpenACM distinguishes itself as the first open-source compiler to systematically integrate accuracy-configurable approximation with advanced variation-aware SRAM analysis. By bridging the gap between exact-computing tools and the error tolerance of AI applications, it provides a comprehensive solution for next-generation energy-efficient designs.

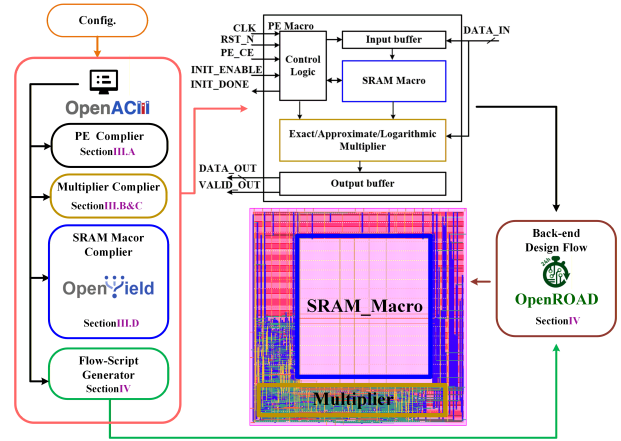


Fig. 1: The end-to-end open-source workflow enabled by OpenACM, comprised of its core hardware components and the EDA toolchain. The flow integrates a library of accuracy-configurable multipliers, a generator for a customizable SRAM macro, and a script generator for the steps in the PD.

III. OpenACM: APPROXIMATE DCiM COMPILER

This section overviews OpenACM’s flow, beginning with the overall architecture and followed by compiler components.

A. Overall Architecture

As illustrated in [Fig. 1](#), OpenACM is an end-to-end framework that generates a DCiM macro from architecture specifications and multiplier configurations, surpassing traditional DCiM compilers that rely solely on bit-width scaling for accuracy control.

The OpenACM compiler consists of four main components:

- 1) **Processing Element (PE) compiler** generates the control logic for the SRAM and multiplier, together with the associated input/output buffers. The PE first initializes the SRAM with the required data, and then performs multiplication between incoming data and the stored values, producing the final results.
- 2) **Multiplier compiler** provides (i) exact multipliers of arbitrary bit widths; (ii) approximate multipliers with configurable precision using approximate 4-2 compressors to optimize the partial-product reduction tree; and

TABLE I: The Landscape of Automated DCiM Compilers: A Missing Capability.

Compiler	Precision	Open?	SRAM Anal.	Multiplier
AutoDCIM [6]	Fixed-point	✗	Basic	Adder-tree
ARCTIC [7]	Int/FP with BIST	✗	Limited	Adder-tree
SynDCIM [8]	Multiple precision	✗	Basic	Adder-tree with 4-2 compressor
SEGA-DCIM [9]	Int/FP precision	✗	Limited	Adder-tree
DAMIL-DCIM [10]	Digital precision	✗	Basic	Adder-tree
OpenC ² [2]	Digital precision	✓	Basic	Adder-tree
OpenACM (Ours)	Configurable approximation	✓	Advanced	Exact / Approx 4-2 compressor / Log

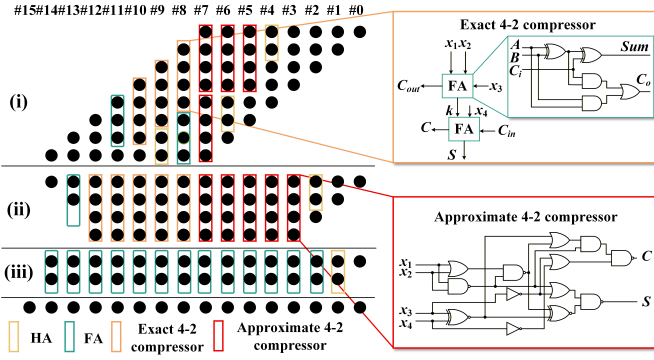


Fig. 2: Block diagram of the proposed 8-bit approximate multiplier. It consists of three stages: (i) partial-product generation, (ii) a configurable reduction tree employing exact or approximate 4-2 compressors on selected low-order columns, and (iii) a final carry-propagate adder.

(iii) logarithmic multipliers to further improve energy efficiency, especially for large bit-width designs.

- 3) **SRAM macro compiler** generates 6T SRAM arrays of arbitrary dimensions, along with the necessary control and read/write circuitry.
- 4) **Flow-script generator** produces the backend design scripts required by OpenROAD and leverages it to complete the PD flow. The compiler ultimately delivers a physically implemented CiM macro together with sign-off artifacts, enabling designers to meet application-specific accuracy and energy-efficiency requirements without relying on proprietary layout generation tools.

B. Approximate Multiplier

Existing DCiM compilers mainly use fixed-point representations with varying bit widths, offering only coarse-grained accuracy control. However, diverse application error tolerances demand finer tuning. To address this, we introduce approximate computing techniques that leverage approximate compressors under a given bit width for more precise accuracy adjustment. In SynDCiM [8], the use of exact 4-2 compressors was proposed to optimize adder trees, but their effectiveness remained constrained by the strict requirement of exact computation. In contrast, OpenACM integrates approximate multipliers based on approximate 4-2 compressors, enabling flexible trade-offs between accuracy and PPA, and significantly expanding the design space of DCiM beyond exact computing. For approximate multipliers based on 4-2 compressors, the circuit structure is typically depicted in Fig. 2, which presents an 8-bit multiplier. In the initial stage, black dots indicate Partial Products (PPs) generated by AND gates from the two input operands. These PPs are compressed through two intermediate stages, ultimately resulting in two rows of PPs in the final stage. These rows are then summed to produce the output of the multiplier. Throughout the three stages, various combinational logic circuits, including HAs, FAs, and 4-2 compressors, are employed to compress the PPs. To optimize resource consumption while maintaining acceptable accuracy,

approximate 4-2 compressors are commonly applied in the lower 8 bits of the PPs, specifically in columns #0 to #7, as highlighted by the red box in Fig. 2. The framework is fully scalable and supports approximate multipliers of arbitrary bit widths. Designers can either tailor approximate 4-2 compressors to meet specific accuracy requirements or adopt widely-used approximate 4-2 designs [18]–[23]. Moreover, OpenACM also supports the generation of exact multipliers of any bit width, enabling flexible selection of multipliers that best satisfy the accuracy requirements of different applications.

C. Logarithmic Multiplier

Although approximate 4-2 compressor multipliers improve area and power efficiency, applications with high error tolerance allow further accuracy relaxation. To extend the design space beyond compressor-based approximation, OpenACM integrates an energy-efficient Logarithmic Multiplier (LM), which naturally provides finer-grained accuracy-PPA trade-offs.

The conventional LM is derived from Mitchell’s Algorithm (MA) [24], which approximates multiplication in the logarithmic domain. For an operand N , it can be expressed as $N = 2^k(1 + x)$, where k denotes the position of the most significant “1” and x represents the fractional part. This expression can also be reformulated as $x \cdot 2^k = N - 2^k$. Consequently, the product of two operands A and B can be written as:

$$A \times B = 2^{k_1+k_2} + (A - 2^{k_1})2^{k_2} + (B - 2^{k_2})2^{k_1} + (A - 2^{k_1})(B - 2^{k_2}) \quad (1)$$

Where k_1 and k_2 represent the leading-one position of A and B , respectively, the $(A - 2^{k_1})(B - 2^{k_2})$ is taken as the Error Part (EP) and the remainder is taken as the Approximate Part (AP). After taking the logarithm on both sides, the AP can be computed using only shift and addition operations, while the EP is typically neglected or approximated. Neglecting the EP introduces significant error; therefore, we propose an adder-free dynamic compensation strategy to handle the EP term.

For an n -bit multiplier, the maximum Rounding Error (RE) associated with the leading-one position k can be expressed as:

$$RE_k = (2^{k+1} - 2^k)/2 = 2^{k-1}, \quad k \in \{1, 2, \dots, n-2\}. \quad (2)$$

The Worst-Case Error (WCE) depends on which operand is chosen for rounding in the EP, where $Q_1 = A - 2^{k_1}$ and $Q_2 = B - 2^{k_2}$. If the smaller operand is rounded, the WCE is $4^{n-2} - 2^{n-3}$; if the larger operand is rounded, the WCE reduces to $3 \cdot 4^{n-3}$. To minimize the WCE, the proposed algorithm dynamically selects and scales the larger operand in the EP to generate compensation values.

Additionally, to incorporate error compensation without hardware overhead, the proposed compensation algorithm generates compensation values within a unique range. Since operands Q_1 and Q_2 are derived from shifted values with their leading ones removed, the range of the compensation value $\text{round}(Q_1)Q_2$ is strictly less than $2^{k_1+k_2}$. This guarantees

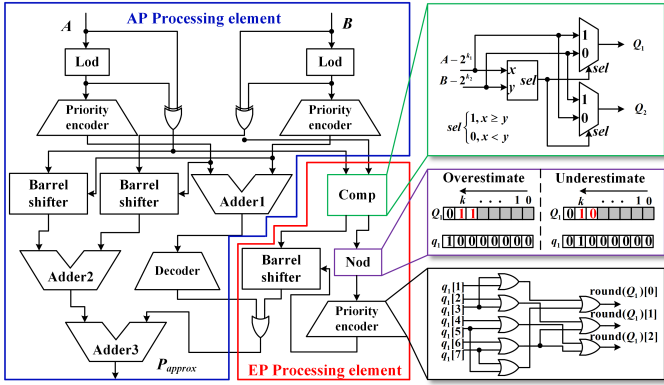


Fig. 3: Block diagram of the proposed 8-bit logarithmic multiplier. It approximates multiplication by (i) base-2 logarithmic encoding using a leading-one detector and a small LUT, (ii) addition in the log domain, and (iii) antilogarithmic decoding via a barrel shifter and LUT, optionally with error-compensation.

that when the compensation is added to $2^{k_1+k_2}$, no carry is generated. Thus, a bitwise OR gate can replace an FA without affecting correctness. The final approximate product can be expressed as:

$$P_{approx} = 2^{k_1+k_2} | \text{round}(Q_1)Q_2 + (A-2^{k_1})2^{k_2} + (B-2^{k_2})2^{k_1} \quad (3)$$

Finally, the proposed LM is shown in Fig. 3. In the AP processing element, two Leading-one Detectors (LoDs) and two priority encoders are implemented to achieve the leading-one position of two operands, k_1 and k_2 represent the leading-one position of two operands, respectively. XOR gates are used to remove the leading-one from two operands, $A - 2^{k_1}$ and $B - 2^{k_2}$ can be obtained, k_1 and k_2 are added by the Adder1, the result of the Adder1 is decoded to the $2^{k_1+k_2}$. $A - 2^{k_1}$ and $B - 2^{k_2}$ are shifted k_2 and k_1 bits by two barrel shifters. The partial sum of $(A - 2^{k_1})2^{k_2}$ and $(B - 2^{k_2})2^{k_1}$ are added by the Adder2. In the EP processing element, the product of $A - 2^{k_1}$ and $B - 2^{k_2}$ are estimated. Two operands are compared in the COMP, and the larger is overestimated to 2^{k+1} or underestimated to 2^k . Then the compensated result can be obtained by shifting the smaller operand. The partial sum of the $2^{k_1+k_2}$ and the result of the EP processing element can be achieved by an OR gate, and intermediate values are combined by Adder3 to produce the approximate product. This architecture easily scales to different bit widths by proportionally widening each module.

D. SRAM Macro

The SRAM macro provides the memory substrate for DCiM operations. As shown in Fig. 4, OpenACM adopts a compact, banked, and subarrayed 6T design with hierarchical Word-Line (WL) decoders/drivers, PREcharge (PRE), write drivers, optional column multiplexers, and differential Sense Amplifiers (SAs). Reads precharge BL/BLB, assert WL, and latch via SA; writes drive BL/BLB while WL is asserted. This structure is intentionally minimal to ease timing closure and tiling across

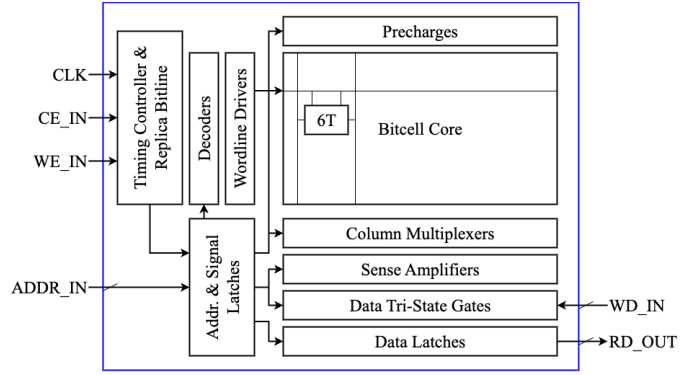


Fig. 4: OpenACM SRAM macro: banked, subarrayed 6T array with hierarchical WL, PRE, write drivers, column MUX, and differential SAs.

banks/subarrays. We highlight the following features in the implementation of the SRAM macro:

- 1) MC- and yield-aware SRAM characterization (from OpenYield [15]): OpenACM runs Monte Carlo SPICE with process variations (local mismatch) and built-in yield analysis algorithms to extract distributions of access latency, SNM (read/write/hold), and dynamic/leakage power; to improve efficiency, it employs the Importance Sampling (IS) method to bias rare-event regions, substantially reducing the number of simulations required and accelerating LIB characterization while preserving accuracy at the target yield.
- 2) This macro supports fine-grained, compiler-visible tunable knobs for co-optimization: rows/cols, word width, bank/subarray count, column-mux ratio, and timing controls (e.g., SAE, precharge), which map directly to design-space exploration with approximate multipliers for accuracy-PPA co-design.
- 3) Flow readiness and reproducibility: abstract views and footprints align with open-source flows (OpenROAD) and FakeRAM-style [14] abstract macros to enable black-box integration during place-and-route; see Sec. IV for flow details.

Furthermore, by adhering to the FakeRAM2.0-style abstract memory interface, our SRAM macro can be seamlessly integrated into other open-source projects that already use FakeRAM macros (e.g., OpenROAD's tinyRocket tutorial [25] that uses fakeram45_256x16/fakeram7_256x32), enabling drop-in replacement in broader SoC flows.

IV. DESIGN FLOW: OPENROAD'S TOOLCHAIN

As shown in Fig. 5, OpenACM accomplishes circuit implementation through a standard digital design flow. OpenACM enables automated circuit generation via Python scripts, supporting configurable parameters such as the number of SRAM rows/columns, multiplier bit-width, approximate multipliers based on 4-2 compressors, and logarithmic multipliers. For compressor-based approximate multipliers, users can further specify the compressor type and the combination strategy of different approximate compressors. Once the configuration

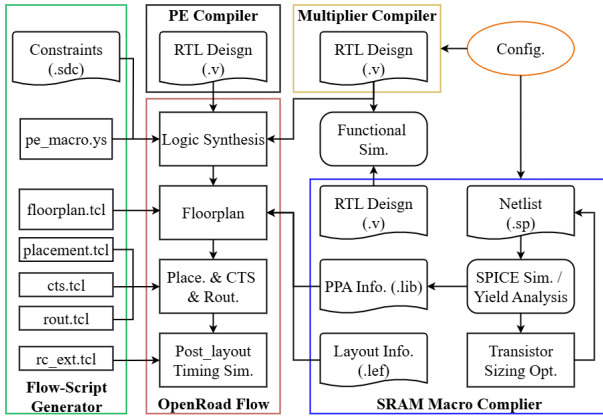


Fig. 5: Overall design flow of OpenACM, which is built on top of OpenROAD and OpenYield.

parameters are defined, the toolchain generates: (i) an SRAM behavioral model, (ii) an LEF abstract for physical integration, (iii) LIB timing/power/area views derived from characterization, (iv) RTL for the PE control logic and multipliers, and (v) the scripts required for OpenROAD to complete the backend physical design. Note that the current release does not generate the exact SRAM layout; the SRAM is integrated as a black-box hard macro during place-and-route. As we want to support transistor sizing customization, the GDSII generation of the SRAM macro is still under development and will be left for future work. OpenROAD implements the standard-cell logic and performs top-level integration around the SRAM abstract, with STA conducted using OpenSTA. After detailed routing, parasitic extraction (SPEF) and SDF generation are applied to the synthesized logic and interconnect (excluding SRAM internals) for post-layout timing simulation, ensuring the design meets the specified functionality and performance requirements.

V. EXPERIMENTS

We evaluate the effectiveness of OpenACM through a set of experiments targeting three key research questions. All CiM macros are generated using the FreePDK45 [13] technology and implemented with the OpenROAD flow [12], from synthesis to layout. Power and energy are derived from post-layout switching activity simulations, while SRAM behavior is characterized using Xyce [26] circuit simulations.

A. PPA Comparison

For a fair comparison, it is worth noting that TOPS/W is not reported since it mainly reflects architecture-level throughput and is heavily affected by interconnect and scheduling, making it unsuitable for PE macro-level evaluation where SRAM access and control dominate performance. Instead, we present area, delay, and power, which directly reveal the circuit efficiency of our PE and provide a more reliable basis for system-level scaling. Tab. II presents the post-layout results of SRAM-multiplier systems generated by OpenACM under a 100 MHz clock and a 0.5 pF load. The compiler automatically assembles the SRAM arrays, control units, and three arithmetic cores:

TABLE II: Post-layout performance of OpenACM-generated SRAM-multiplier systems at 100 MHz and a 0.5 pF output load.

SRAM	Multiplier Type	Delay (ns)	Area (μm^2)			Power (W)
			Logic	SRAM	P&R	
16×8 (8-bit width)	OpenC ²	5.22	1431	7052	8483	2.82E-04
	Exact	5.22	1079		8131	2.45E-04
	Log-our	5.22	1173		8225	2.82E-04
	Appro4-2	5.22	939		7991	2.11E-04
32×16 (16-bit width)	OpenC ²	5.24	4842	16910	21752	1.15E-03
	Exact	5.24	3568		20478	1.08E-03
	Log-our	5.24	2402		19312	6.15E-04
	Appro4-2	5.24	2633		19543	7.58E-04
64×32 (32-bit width)	OpenC ²	5.24	19734	48642	68376	7.00E-03
	Exact	5.24	10132		58774	4.03E-03
	Log-our	5.24	4960		53602	1.45E-03
	Appro4-2	5.24	9331		57973	3.36E-03

an exact multiplier using exact 4-2 compressors, an approximate design based on approximate 4-2 compressors (Appro4-2), and the proposed logarithmic multiplier (Log-our). Since OpenACM supports arbitrary combinations of approximate 4-2 compressors, we adopt the widely used and highly cited Yang1 [22] as a representative example to demonstrate the energy-efficiency advantages of approximate CiM compiler. For completeness, we also include a representative adder-tree-based implementation from OpenC² as a baseline.

All designs are evaluated using the same multiplication workloads to ensure fair power comparison. The results show that the critical delay is nearly constant (5.2 ns) across all multipliers, indicating SRAM dominates system timing. Approximate designs achieve significant area reduction: the logarithmic multiplier cuts logic area by 33% for 32×16 and 51% for 64×32 . For energy, the 4-2 compressor design saves up to 14% power at 16×8 , favoring small-scale cases, while the logarithmic design excels at larger scales, reducing power by nearly 64% compared to the exact multiplier in 64×32 and also outperforming Appro4-2. By contrast, the adder-tree-based architecture (OpenC² [2]) exhibits consistently higher area and power consumption across all configurations. It is worth noting that for small bit-widths (e.g., 8-bit), the logarithmic multiplier exhibits a slightly higher overhead than the 4-2-based structure. This is because the logarithm and anti-logarithm estimation modules constitute a relatively large portion of the short datapath, introducing additional logic complexity and switching activity. In contrast, the 4-2 compressor network remains structurally compact and therefore offers advantages in both area and dynamic power at low bit-widths. As bit-width increases, the relative overhead of the logarithmic modules diminishes, enabling the logarithmic multiplier to deliver substantial area and energy benefits in medium- and large-scale configurations.

B. Accuracy-Constrained Applications

Image blending [27] and edge detection [23] are used to evaluate the practical performance of the approximate multipliers generated by OpenACM. In image blending, an 8-bit unsigned multiplier processes two grayscale images pixel by

TABLE III: PSNR Comparison of Different Approximate Multipliers for Various Image Processing Tasks.

Processing Task	Test Image	Multiplier Type		
		Appro4-2	Log-our	LM [24]
Image Blending	Lake & Mandril	67.19 dB	32.01 dB	26.08 dB
	Jetplane & Boat	70.93 dB	37.17 dB	22.10 dB
	Cameraman & Lake	69.81 dB	43.22 dB	24.82 dB
Edge Detection	Boat	66.21 dB	46.43 dB	38.77 dB
	Cameraman	67.55 dB	45.61 dB	38.37 dB
	Jetplane	66.20 dB	44.13 dB	39.07 dB

TABLE IV: Influence of approximate multiplier on Top-1 and Top-5 scores for the ResNet-18 network.

Multiplier Type	Top-1	Top-5	NMED	MRED
Exact	0.677	0.873	-	-
Appro4-2	0.668	0.880	1.70E-09	1.27E-10
Log-our	0.680	0.870	4.40E-03	1.55E-02
LM [24]	0.610	0.842	2.79E-02	9.44E-02

pixel, with results scaled back to 8 bits. For edge detection, convolution and squaring employ a 16-bit signed approximate multiplier, while the square root is computed exactly. Image quality is measured by peak signal-to-noise ratio (PSNR), using the exact multiplier as the baseline. Typically, PSNR below 30 dB indicates visible degradation, whereas values above 40 dB imply near-identical quality.

As shown in Tab. III, Appro4-2 achieves high accuracy and can effectively replace exact multipliers in image processing. Log-our further improves accuracy over the LM [24]. In image blending, LM produces unsatisfactory results with PSNR generally below 30 dB, while Log-our consistently surpasses this threshold, making it viable for energy-efficient applications. In edge detection, LM reaches about 30 dB, whereas Log-our achieves 45 dB, delivering quality comparable to the exact design. Overall, the proposed OpenACM architecture supports both accuracy-critical and energy-constrained scenarios, enabling flexible circuit selection tailored to diverse application needs.

Neural networks are widely applied in computer vision and image classification, where approximate multipliers can often replace exact ones during forward inference. We utilized the pre-trained ResNet-18 [28] model on the ILSVRC2012 dataset as a baseline and tested various approximate multipliers for the Top-1 and Top-5 accuracies of this neural network. The floating-point weights and inputs are quantized to 32-bit fixed-point for computations. In addition, the normalized mean error distance (NMED) and mean relative error distance (MRED) are calculated as error metrics of the approximate multipliers for reference.

As shown in Tab. IV, both Appro4-2 and Log-our reduce power consumption significantly without degrading classification accuracy. Specifically, Appro4-2 achieves 17% power savings, while Log-our achieves 64% compared to the exact multiplier. Thanks to the inherent error resilience of deep neural networks, such reductions can be realized with no

TABLE V: Comparison of MC and MNIS Yield Analysis Methods on Various SRAM Sizes.

Metric	Monte Carlo			MNIS [29]			Speedup
	P_f	FoM	#Sim.	P_f	FoM	#Sim.	
16×2	1.6E-4	0.1	55,600	3.2E-4	0.05	2,985	18×
32×2	6.4E-2	0.17	22,900	1.7E-2	0.15	2,260	10×
64×2	3.9E-3	0.05	41,500	1.5E-3	0.03	4,260	9.7×

loss of inference quality. Notably, although Appro4-2 introduces smaller absolute errors, its one-sided distribution leads to systematic deviations in accumulated results. In contrast, Log-our produces bidirectional errors resembling zero-mean perturbations, which act as noise regularization and enhance generalization. Consequently, Log-our attains slightly higher Top-1 accuracy than the exact multiplier, while Appro4-2 shows only minor Top-5 improvement. The conventional LM suffers from large errors, causing a severe accuracy drop, further validating the effectiveness of the proposed Log-our.

C. IS-Based Yield Analysis

Tab. V presents a comparison of yield estimation between Monte Carlo (MC) and Mean-shifted Importance Sampling (MNIS) [29]. To further accelerate circuit-level simulations, we use trimmed SRAM arrays with only two columns ($N \times 2$), while retaining the full WL parasitics of the original arrays. This preserves the WL RC loading seen by the drivers but reduces the number of simulated bitline columns, yielding substantial runtime savings. The figure of merit (FoM) shown in the table is defined as $\text{std}(P_f)/P_f$, and $\text{std}(P_f)$ represents the standard deviation of the estimated failure rate P_f . As shown, MC requires a substantially larger number of simulations (e.g., 41,500 for the 64×2 case), whereas MNIS achieves comparable accuracy with significantly fewer runs (only 4,260), yielding a speedup of approximately 9.7×. Similar trends are observed for smaller circuits, where MNIS achieves about 18× speedup for the 16×2 case and 10× for the 32×2 case. Furthermore, as the circuit size increases, MNIS continues to exhibit superior scalability and robustness, thereby reinforcing its effectiveness. Overall, these results highlight the advantages of importance sampling-based approaches in yield estimation and demonstrate their strong potential to replace conventional MC methods in large-scale circuit characterization.

VI. CONCLUSION

In this work, we presented OpenACM, the first open-source, accuracy-aware compiler for SRAM-based approximate DCiM architectures. While OpenACM provides a solid foundation, several key extensions remain. Our near-term priorities include completing the automated layout generator for custom SRAM macros to enable a fully physical-aware flow, and extending the front-end and multiplier library to support native floating-point operations. We also plan to develop an automated DSE engine to jointly optimize multiplier choice, precision, and array organization for application-specific efficiency. We believe OpenACM will continue to support and accelerate research in energy-efficient computing.

REFERENCES

- [1] G. Armeniakos, G. Zervakis, D. Soudris, and J. Henkel, "Hardware approximate techniques for deep neural network accelerators: A survey," *ACM Computing Surveys*, vol. 55, no. 4, pp. 1–36, 2022.
- [2] T. Dong, S. Li, Y. Zuo, H. Jiang, Y. Ma, and S. Huang, "Openc 2: An open-source end-to-end hardware compiler development framework for digital compute-in-memory macro," in *2025 Design, Automation & Test in Europe Conference (DATE)*. IEEE, 2025, pp. 1–2.
- [3] K. Yoshioka, S. Ando, S. Miyagi, Y.-C. Chen, and W. Zhang, "A review of sram-based compute-in-memory circuits," *Japanese Journal of Applied Physics*, 2024.
- [4] C.-Y. Chen, J. Choi, K. Gopalakrishnan, V. Srinivasan, and S. Venkataramani, "Exploiting approximate computing for deep learning acceleration," in *2018 Design, Automation & Test in Europe Conference & Exhibition (DATE)*. IEEE, 2018, pp. 821–826.
- [5] Q. Zhang, T. Wang, Y. Tian, F. Yuan, and Q. Xu, "Approxann: An approximate computing framework for artificial neural network," in *2015 Design, Automation & Test in Europe Conference & Exhibition (DATE)*. IEEE, 2015, pp. 701–706.
- [6] J. Chen, F. Tu, K. Shao, F. Tian, X. Huo, C.-Y. Tsui, and K.-T. Cheng, "Autodcim: An automated digital cim compiler," in *2023 60th ACM/IEEE Design Automation Conference (DAC)*. IEEE, 2023, pp. 1–6.
- [7] H. Zhang, H. Zhu, S. He, M. Li, C. Wang, X. Xiong, H. Tian, X. Zeng, and C. Chen, "Arctic: Agile and robust compute-in-memory compiler with parameterized int/fp precision and built-in self test," in *2024 Design, Automation & Test in Europe Conference & Exhibition (DATE)*. IEEE, 2024, pp. 1–6.
- [8] K. Shao, F. Tian, X. Wang, J. Zheng, J. Chen, J. He, H. Wu, J. Chen, X. Guan, Y. Deng *et al.*, "Syndcim: A performance-aware digital computing-in-memory compiler with multi-spec-oriented subcircuit synthesis," in *2025 Design, Automation & Test in Europe Conference (DATE)*. IEEE, 2025, pp. 1–7.
- [9] H. Diao, H. Zhang, J. Song, H. Luo, Y. Lin, R. Wang, Y. Wang, and X. Tang, "Sega-dcim: Design space exploration-guided automatic digital cim compiler with multiple precision support," in *2025 Design, Automation & Test in Europe Conference (DATE)*. IEEE, 2025, pp. 1–7.
- [10] C. Wang, K. Hu, F. Yang, K. Zhu, and X. Zeng, "Damil-dcim: A digital cim layout synthesis framework with dataflow-aware floorplan and milp-based detailed placement," in *2025 Design, Automation & Test in Europe Conference (DATE)*. IEEE, 2025, pp. 1–7.
- [11] AnySilicon, "The ultimate guide to open source eda tools," 2024, available: <https://any silicon.com/the-ultimate-guide-to-open-source-eda-tools/>.
- [12] T. Ajayi, V. A. Chhabria, M. Fogaça, S. Hashemi, A. Hosny, A. B. Kahng, M. Kim, J. Lee, U. Mallappa, M. Neseem *et al.*, "Toward an open-source digital flow: First learnings from the openroad project," in *Proceedings of the 56th Annual Design Automation Conference 2019*, 2019, pp. 1–4, available: <https://openroad.readthedocs.io/>.
- [13] N. C. S. University, "Freepdk45: An open-source 45nm process design kit," 2024, available: <https://www.eda.ncsu.edu/freepdk/freepdk45/>.
- [14] ABKGroup, "Fakeram2.0: A fake memory compiler for physical design research," 2024, available: <https://github.com/ABKGroup/FakeRAM2.0>.
- [15] S. Shen, X. Li, Z. Liu, Y. Wang, Y. Wu, J. Ma, Y. Sun, and W. W. Xing, "Openyield: An open-source sram yield analysis and optimization benchmark suite," *arXiv preprint arXiv:2508.04106*, 2025.
- [16] S. Venkataramani, S. Chakradhar, K. Roy, and A. Raghunathan, "Approximate computing and the quest for computing efficiency," in *Proceedings of the 52nd Annual Design Automation Conference*, pp. 1–6, 2015.
- [17] V. Mrazek, Z. Vasicek, L. Sekanina, M. A. Hanif, and M. Shafique, "Hardware approximate techniques for deep neural network accelerators: A survey," *ACM Transactions on Design Automation of Electronic Systems (TODAES)*, vol. 27, no. 5, pp. 1–34, 2022.
- [18] O. Akbari, M. Kamal, A. Afzali-Kusha, and M. Pedram, "Dual-quality 4:2 compressors for utilizing in dynamic accuracy configurable multipliers," *IEEE Transactions on Very Large Scale Integration (VLSI) Systems*, vol. 25, no. 4, pp. 1352–1361, 2017.
- [19] M. Ha and S. Lee, "Multipliers with approximate 4–2 compressors and error recovery modules," *IEEE Embedded Systems Letters*, vol. 10, no. 1, pp. 6–9, 2018.
- [20] T. Kong and S. Li, "Design and analysis of approximate 4–2 compressors for high-accuracy multipliers," *IEEE Transactions on Very Large Scale Integration (VLSI) Systems*, vol. 29, no. 10, pp. 1771–1781, 2021.
- [21] A. Momeni, J. Han, P. Montuschi, and F. Lombardi, "Design and analysis of approximate compressors for multiplication," *IEEE Transactions on Computers*, vol. 64, no. 4, pp. 984–994, 2015.
- [22] Z. Yang, J. Han, and F. Lombardi, "Approximate compressors for error-resilient multiplier design," in *2015 IEEE International Symposium on Defect and Fault Tolerance in VLSI and Nanotechnology Systems (DFTS)*, 2015, pp. 183–186.
- [23] A. G. M. Strollo, E. Napoli, D. De Caro, N. Petra, and G. D. Meo, "Comparison and extension of approximate 4:2 compressors for low-power approximate multipliers," *IEEE Transactions on Circuits and Systems I: Regular Papers*, vol. 67, no. 9, pp. 3021–3034, 2020.
- [24] J. N. Mitchell, "Computer multiplication and division using binary logarithms," *IRE Transactions on Electronic Computers*, no. 4, pp. 512–517, 2009.
- [25] The OpenROAD Project, "Openroad: Open-source asic design for computer architects (micro 2022 tutorial repository)," GitHub repository, 2022, available: <https://github.com/The-OpenROAD-Project/micro2022tutorial>.
- [26] E. Keiter, J. Verley, T. Russo, R. Hoekstra, H. Thornquist, T. Takhtaganov, R. Schiek, T. Mei, R. Pawlowski, K. Santarelli *et al.*, "Xyce (™) parallel electronic simulator v. 7.6," Sandia National Laboratories (SNL), Albuquerque, NM, and Livermore, CA . . . , Tech. Rep., 2013.
- [27] F. Sabetzadeh, M. H. Moaiyeri, and M. Ahmadinejad, "A majority-based imprecise multiplier for ultra-efficient approximate image multiplication," *IEEE Transactions on Circuits and Systems I: Regular Papers*, vol. 66, no. 11, pp. 4200–4208, 2019.
- [28] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 770–778.
- [29] L. Dolecek, M. Qazi, D. Shah, and A. Chandrakasan, "Breaking the simulation barrier: Sram evaluation through norm minimization," in *Proceedings of the International Conference on Computer-Aided Design (ICCAD)*. IEEE, 2008, pp. 322–329.