

Interconnect Parasitics and Partitioning in Fully-Analog In-Memory Computing Architectures

Md Hasibul Amin, Mohammed Elbtity, Ramtin Zand

Department of Computer Science and Engineering, University of South Carolina, Columbia, SC 29208, USA

Abstract—Fully-analog in-memory computing (IMC) architectures that implement both matrix-vector multiplication and non-linear vector operations within the same memory array have shown promising performance benefits over conventional IMC systems due to the removal of energy-hungry signal conversion units. However, maintaining the computation in the analog domain for the entire deep neural network (DNN) comes with potential sensitivity to interconnect parasitics. Thus, in this paper, we investigate the effect of wire parasitic resistance and capacitance on the accuracy of DNN models deployed on fully-analog IMC architectures. Moreover, we propose a partitioning mechanism to alleviate the impact of the parasitic while keeping the computation in the analog domain through dividing large arrays into multiple partitions. The SPICE circuit simulation results for a $400 \times 120 \times 84 \times 10$ DNN model deployed on a fully-analog IMC circuit show that a 94.84% accuracy could be achieved for MNIST classification application with 16, 8, and 8 horizontal partitions, as well as 8, 8, and 1 vertical partitions for first, second, and third layers of the DNN, respectively, which is comparable to the $\sim 97\%$ accuracy realized by digital implementation on CPU. It is shown that accuracy benefits are achieved at the cost of higher power consumption due to the extra circuitry required for handling partitioning.

Index Terms—analog computing, in-memory computing, interconnect parasitics, memristive technology, partitioning.

I. INTRODUCTION

With the increased computational demands of machine learning (ML) workloads, in-memory computing (IMC) [1] architectures have attracted considerable attention to address the processor-memory bottleneck in conventional von Neumann architectures through executing the logic functions directly on memory via changing the internal memory circuitry. Resistive random access memory (RRAM) [2], phase-change memory (PCM) [3], magnetoresistive random-access memory (MRAM) [4], and conductive bridging random access memory (CBRAM) [5] are some of the promising technologies that have been utilized in IMC architectures to realize matrix-vector multiplication (MVM) operation in DNNs. While memristive technologies are also leveraged in digital IMC architectures [6], [7] to realize logic functions such as XNOR, here we focus on analog IMC architectures due to their great potential for achieving outstanding energy efficiency. For instance, Imec, a world-leading research and innovation hub in nanoelectronics, has recently provided a blueprint towards 10,000 tera operations per second per watt (TOPS/W) DNN inference in [8], which is based on the memristive-based analog IMC architectures.

Despite the potential benefits of the analog IMC architectures, one of the major factors limiting their wide use in

practical ML applications is the large and energy-hungry signal conversion units required to change the computation domain from analog-to-digital (and vice versa) to compute the non-linear vector operations, *e.g.* activation functions in DNNs [9]. Recently, fully-analog IMC architectures are introduced that use memristive technologies to realize both MVM operations and activation functions within the same memory array [10]. These architectures remove the need for the signal conversion unit through maintaining the computation in the analog domain across various layers of DNNs and achieve significant energy and performance improvements. However, due to the fully-analog characteristic of these circuits, interconnect parasitics can have a major impact on the reliability and accuracy of the results obtained by these architectures. Thus, in this paper, we investigate the effect of interconnect parasitics on accuracy of fully-analog IMC architectures and propose an analog partitioning approach to resolve the parasitics effects while keeping the computation in the analog domain.

II. FULLY-ANALOG IMC ARCHITECTURES

Figure 1 shows a schematic of the fully-analog IMC architectures, which includes a network of tightly coupled subarrays interconnected through programmable switch blocks. Each IMAC subarray consists of memristive crossbar, differential amplifiers, and neuron circuits, as shown in Fig. 1 (b). The memristor crossbars compute the MVM operation in DNNs in the analog domain through various physical mechanisms such as Ohm's law and Kirchhoff's law in electrical circuits [11]. In particular, the multiplications are performed according to the Ohm's law ($I = GV$), and the accumulation operation is based on the conservation of charge described by the Kirchhoff's current law as expressed in the following equation, $I_j = \sum_j G_{ij}V_i$, where G_{ij} is the conductance of the resistive devices between neurons i and j , I_j is the input current of post-synaptic neuron j , and V_i is the output voltage of pre-synaptic neuron i .

Each positive and negative weight can be realized through adjusting the relative conductance of two memristive devices that are connected to a differential amplifier. As shown in Fig. 1 (b), differential amplifiers are connected to two consecutive rows in the crossbar, *i.e.* representing positive and negative weights, and generate an output current of $I_{O,i} = \sum_{k=1}^n (I_{k,i}^+ - I_{k,i}^-)$ for the i th row, where n is the total number of input nodes in the layer. Whereas $I_{k,i}^+ \propto V_k G_{k,i}^+$ and $I_{k,i}^- \propto V_k G_{k,i}^-$, thus, $I_{O,i} \propto \sum_{k=1}^n V_k (G_{k,i}^+ - G_{k,i}^-)$, in which $G_{k,i}^+$ and $G_{k,i}^-$ are the conductance of resistive devices that are shown in

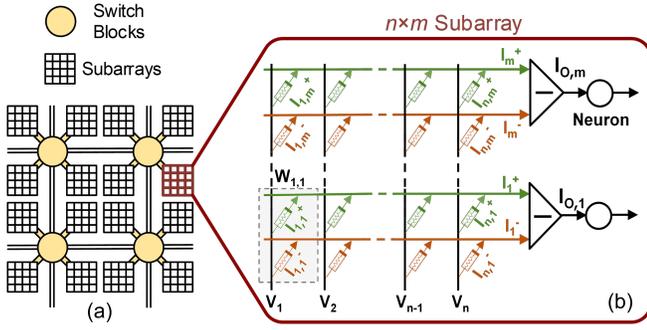


Fig. 1. (a) Fully-analog IMC architecture, (b) An $n \times m$ subarray.

green and red color in Fig. 1 (b), respectively. Finally, the outputs of the differential amplifiers are connected to the analog neurons to compute the activation functions. In this architecture, each subarray computes both MVM operations and neurons' activation functions in a single DNN layer and passes the result to its downstream neighbor IMC subarrays that can compute the next layer.

III. INTERCONNECT PARASITICS CALCULATION

Parasitic interconnect resistance (R_W) and capacitance (C_W) are a function of wire geometry and material properties of interconnections in analog IMC subarrays. Scaling up the size of arrays increases R_W and C_W leading to an increase in the latency of IMC circuits, thus limiting their operating clock frequency. Furthermore, increasing R_W reduces the read margin that can impact the accuracy of analog IMC circuits [12]. Figure 2(c) shows the parasitic model for one bitcell in the IMC array. We use the most common equation to find the parasitic resistances for interconnects:

$$R_W = \rho \frac{L}{W.T}, \quad (1)$$

where ρ , L , W and T are the resistivity, length, width, and thickness of the metal wire, respectively. Resistivity is commonly a fixed parameter for a specific metal. However, for the sub-micron technology nodes, the resistivity increases due to the surface and grain boundary scattering as the metal width gets comparable to the mean free path of electrons [13]. These two well-known scattering effects are quantified using Fuchs-Sondheimer (FS) [14] and Mayadas-Shatzkes (MS) [15] model respectively.

$$\frac{\rho_{FS}}{\rho_{Cu}} = 1 + (1-p) \frac{l_0}{W} \quad (2)$$

$$\frac{\rho_{MS}}{\rho_{Cu}} = \left[1 - \frac{3\alpha}{2} + 3\alpha^2 - 3\alpha^3 \ln \left(1 + \frac{1}{\alpha} \right) \right]^{-1} \quad (3)$$

where $\alpha = \frac{l_0}{d} \frac{R}{1-R}$, ρ_{Cu} is the resistivity of bulk Cu ($1.9 \times 10^{-9} \Omega m$), l_0 is the mean free path of electrons in Cu (39 nm), W is the width of the metal wire, p is the specular scattering fraction, d is the average grain size and R is the probability for electrons to reflect at the grain boundary. R and p are assumed

to be 0.3 and 0.25, respectively, and d is assumed to be equal to the wire width as mentioned in the literature [16], [17]. The two scattering effects are combined using Matthiessen's rule which results in the following equation [18]:

$$\begin{aligned} \frac{\rho}{\rho_{Cu}} &= 1 + \left(\frac{\rho_{FS}}{\rho_{Cu}} - 1 \right) + \left(\frac{\rho_{MS}}{\rho_{Cu}} - 1 \right) \\ &= (1-p) \frac{l_0}{W} + \left[1 - \frac{3\alpha}{2} + 3\alpha^2 - 3\alpha^3 \ln \left(1 + \frac{1}{\alpha} \right) \right]^{-1} \end{aligned} \quad (4)$$

The parasitic capacitances play major roles in determining the latency of analog IMC circuits. To obtain good accuracy, we fixed the sampling time at 1ns considering the overall latency due to the addition of parasitic capacitances. We use the Sakurai-Tamaru model [19] for calculating the parasitic capacitance per length:

$$\begin{aligned} C_W &= \epsilon \times \frac{1}{2} \left[1.15 \left(\frac{W}{H} \right) + 2.8 \left(\frac{W}{H} \right)^{0.222} \right] \\ &+ \epsilon \times 2 \left[0.03 \left(\frac{W}{H} \right) + 0.83 \left(\frac{T}{H} \right) - 0.07 \left(\frac{T}{H} \right)^{0.222} \right] \\ &\times \left(\frac{S}{H} \right)^{-1.34} \end{aligned} \quad (5)$$

where $\epsilon = 20\epsilon_0$ is the dielectric permittivity of the inter-metal space, W and T are the width and thickness of the metal line, $H = 20nm$ is the inter-metal layer spacing and S is the inter-wire spacing [12]. Here, we leverage equations (1) to (5) to model the interconnect parasitics in the SPICE circuit simulations of analog IMC architectures.

IV. ANALOG HORIZONTAL AND VERTICAL PARTITIONING

As the interconnect parasitic impacts can severely degrade the accuracy of the fully-analog IMC circuits, we propose an analog horizontal and vertical partitioning technique to decrease R_W and C_W without requiring to convert the signals from analog domain to digital. Figures 2 (a) and 2 (b) provide a schematic of the horizontal and vertical partitioning circuitry, respectively. For the horizontal partitioning, a layer of demultiplexers (DEMUX) is added to the output of the crossbars, which distributes the output currents corresponding to the matrix-vector multiplication results to either neurons in the same subarray for normal non-partitioned operation, or to the next subarray as partial products of that particular partition. Moreover, we locate switches on the output of the crossbars before DEMUX circuits to identify whether the generated output currents should be accumulated with the currents arriving from other subarrays (i.e. partitions) or not. Using these peripheral circuitry and necessary signaling to control the switches and DEMUX circuits, IMC circuit can handle the horizontal partitioning in the analog domain. Figure 2 (a) shows an example of horizontal partitioning with two partitions $H_P = 2$. For vertical partitioning, an $n \times m$ array is divided into multiple $n \times k_i$ subarrays, in which $m = \sum_{i=0}^{V_P} k_i$,

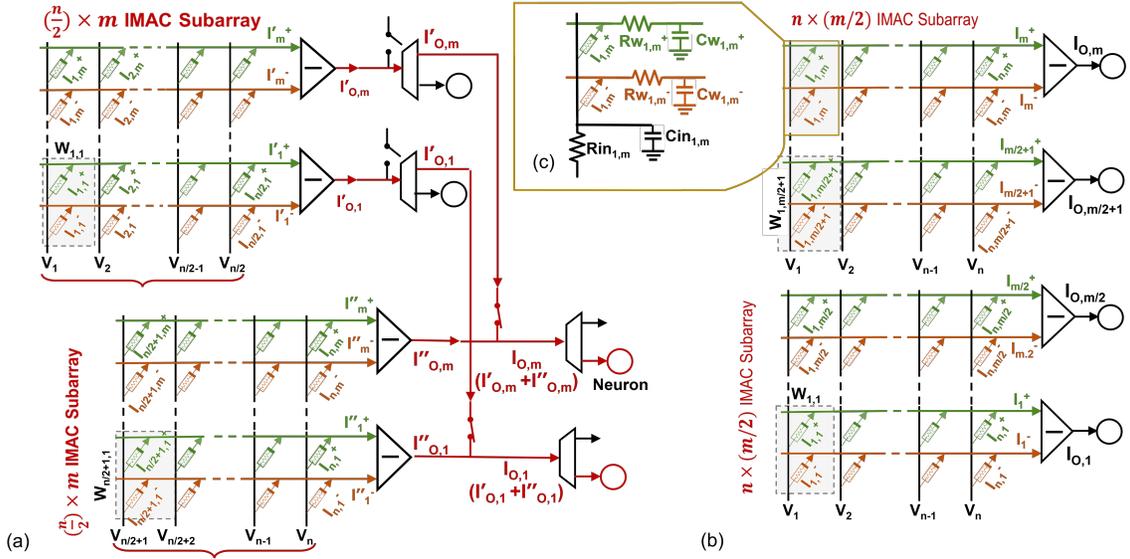


Fig. 2. (a) Horizontal partitioning ($H_P = 2$), and (b) vertical partitioning ($V_P = 2$) in an analog IMC array. (c) Parasitic capacitance and resistance model.

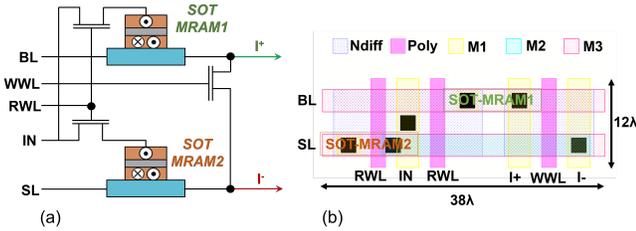


Fig. 3. (a) The SOT-MRAM based synapse bitcell. (b) Layout design.

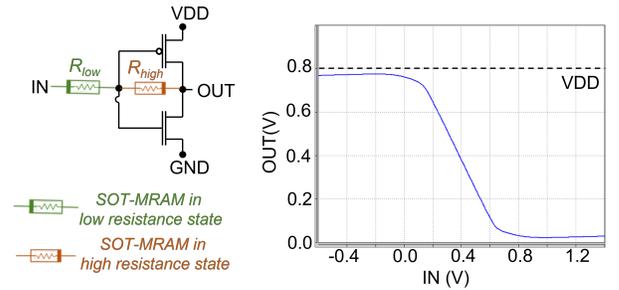


Fig. 4. The memristive sigmoid neuron circuit and SPICE simulation.

where V_P is the total number of vertical partitions. Figure 2 (b) shows a sample of vertical partitioning with $V_P = 2$.

V. SIMULATION RESULTS AND DISCUSSION

In this section, we implement a $400 \times 120 \times 84 \times 10$ DNN in SPICE circuit simulator for MNIST [20] handwritten image classification with 20×20 pixels. We use the 14nm High-Performance PTM-MG FinFET model [21] along with the V_{DD} and V_{SS} voltages of 0.8V and -0.8V, respectively. Based on the 18nm gate length and the 22nm Fin height of the PTM 14nm FinFET model [22], the layout design parameter λ and the metal thickness are fixed to 9nm and 22nm, respectively.

Here, spin orbit torque MRAM (SOT-MRAM) device model [23] is utilized to implement both synaptic structures and activation functions, as shown in Fig. 3 and Fig. 4, respectively. We use an analog sigmoidal neuron that includes two resistive devices and a CMOS-based inverter [10]. The resistive devices in the neuron's circuit create a voltage divider that reduces the slope of the linear operating region in the inverter leading to a smooth high-to-low output voltage transition, which enables the realization of a *sigmoid* activation function.

First, we study the effect of partitioning on accuracy and power consumption of the $400 \times 120 \times 84 \times 10$ DNN implemented on a fully-analog IMC circuit. We select the number

TABLE I
EFFECT OF PARTITIONING ON THE ACCURACY AND POWER CONSUMPTION OF FULLY-ANALOG IMC CIRCUITS. THE DIMENSIONS OF THE LAYERS ARE $L_1=400 \times 120$, $L_2=120 \times 84$, AND $L_3=84 \times 10$.

Array Size	Number of Partitions						Accuracy	Power (W)
	Horizontal (H_P)			Vertical (V_P)				
	L1	L2	L3	L1	L2	L3		
32×32	13	4	3	4	3	1	91.71%	2.640
64×64	7	2	2	2	2	1	84.16%	1.592
128×128	4	1	1	1	1	1	15.43%	0.826
256×256	2	1	1	1	1	1	13.17%	0.829
512×512	1	1	1	1	1	1	10.42%	0.927
32×32	16	8	8	8	8	1	94.84%	3.375

of partitions for each layer based on various dimensions of IMC subarrays, as listed in Table I. For instance, every layer in the targeted DNN can be deployed on an IMC architecture with 512×512 subarrays without any partitioning, while if we use 256×256 subarrays, the first layer that includes 400 inputs must be divided into two horizontal partitions to fit into two 256×256 subarrays. The results listed in Table I show that without partitioning the deployed model fails to provide reliable classification. It can be seen that as the number of horizontal and vertical partitions increases, the accuracy improves due to the decrease in the length of the interconnects,

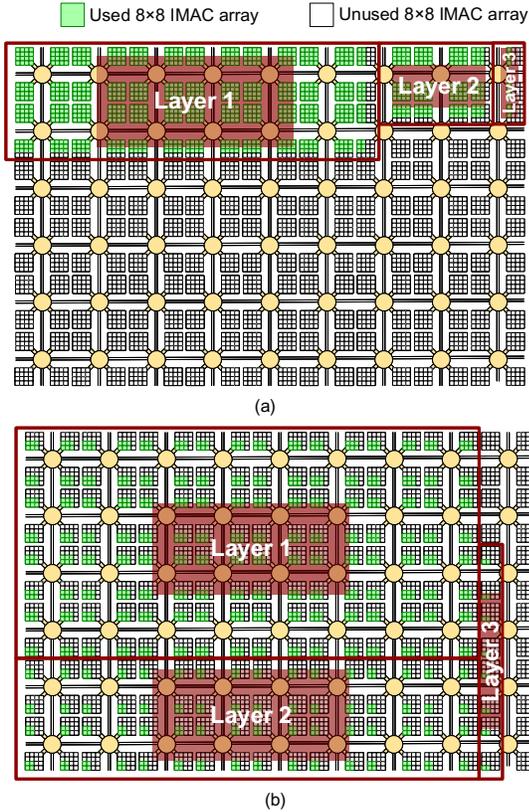


Fig. 5. The deployment of a $400 \times 120 \times 84 \times 10$ DNN on a fully-analog IMC architecture with 32×32 subarrays. (a) Maximum subarray utilization with $H_P = [13, 4, 3]$ and $V_P = [4, 3, 1]$. (b) Highly-partitioned deployment using $H_P = [16, 8, 8]$ and $V_P = [8, 8, 1]$.

and consequently their parasitic resistances. However, this is achieved at the cost of increased power consumption due to the extra circuitry added to handle partitioning.

The last row of Table I shows the results for a highly-partitioned case with $H_P = [16, 8, 8]$ and $V_P = [8, 8, 1]$ horizontal and vertical partitions for each layer, respectively. This means that assuming an analog IMC architecture with 32×32 subarrays, the deployed model does not use the entire capacity of the subarrays, as shown in Fig. 5 (b). This deployment scenario results in a high accuracy of 94.84%, which is close to the $\sim 97\%$ accuracy realized by the full-precision digital implementations on CPU. However, it is achieved at the cost of higher power consumption and more distributed deployment of DNN model on the architecture that leads to a higher area utilization.

Finally, we investigate the impact of bitcell size on the accuracy and power consumption of the analog IMC circuits. The distance between the wires and the length of the metal lines in an analog IMC subarray depends on the size of the synapse bitcell, which affects the interconnect parasitic resistances and capacitances as described in Section III. Figure 6 shows a non-ideal layout design for the SOT-MRAM based synapse, which leads to a larger bitcell area compared to what is realized in Figure 3. Table II provides the accuracy and

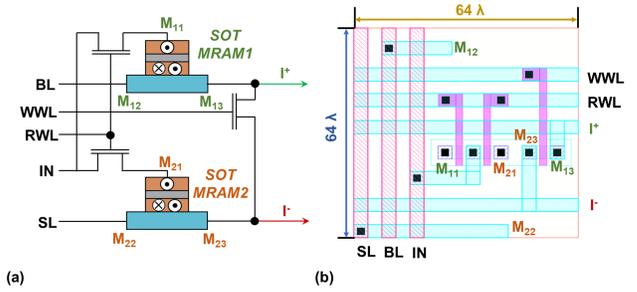


Fig. 6. (a) The SOT-MRAM based synapse bitcell, (b) non-ideal layout design.

TABLE II
EFFECTS OF PARTITIONING ON THE DEPLOYMENT OF A $400 \times 120 \times 84 \times 10$ DNN ON AN ANALOG IMC CIRCUIT WITH NON-IDEAL SYNAPSE LAYOUT DESIGN.

Array Size	Partitioning						Accuracy	Power (W)
	Horizontal (H_P)			Vertical (V_P)				
	L1	L2	L3	L1	L2	L3		
32×32	13	4	3	4	3	1	73.64%	1.747
64×64	7	2	2	2	2	1	28.44%	0.926
128×128	4	1	1	1	1	1	11.35%	0.476
256×256	2	1	1	1	1	1	11.35%	0.478
512×512	1	1	1	1	1	1	11.35%	0.479
32×32	16	8	8	8	8	1	94.04%	2.774

power consumption results for various partitioning scenarios for the $400 \times 120 \times 84 \times 10$ DNN workload deployment on an analog IMC architecture with non-ideal synapse layout design. Accuracy comparisons show that a $\sim 55\%$ accuracy drop for the non-ideal IMC architecture with $H_P = [7, 2, 2]$ and $V_P = [2, 2, 1]$ partitions can reduce to less than 1% accuracy drop for the highly-partitioned scenario with $H_P = [16, 8, 8]$ and $V_P = [8, 8, 1]$ partitions. This shows that increasing the number of partitions can potentially compensate for the imperfections in the layout design at the cost of higher power and area consumption.

VI. CONCLUSION

Herein, we focused on the impacts of interconnect parasitics on the accuracy of DNN models deployed on the fully-analog IMC architectures. The initial simulation results show that without any mechanisms to resolve the parasitic effects, a $400 \times 120 \times 84 \times 10$ DNN model can barely achieve 15% accuracy for MNIST classification. Thus, we proposed a horizontal and vertical partitioning mechanism to alleviate the parasitic impacts, while maintaining the computation in the analog domain. This is particularly important in fully-analog IMC architectures which are designed to remove the need for signal conversion units through implementing nonlinear activation functions as well as matrix-vector multiplications in the analog domain. Our proposed partitioning mechanism has shown to be effective to diminish the parasitic impacts such that more than 94% accuracy could be realized for two different ideal and non-ideal layout design for the IMC circuit.

REFERENCES

- [1] D. Ielmini and H.-S. P. Wong, "In-memory computing with resistive switching devices," *Nature Electronics*, vol. 1, no. 6, pp. 333–343, 2018.
- [2] S. Yin, X. Sun, S. Yu, and J. sun Seo, "High-throughput in-memory computing for binary deep neural networks with monolithically integrated rram and 90-nm cmos," *IEEE Transactions on Electron Devices*, vol. 67, pp. 4185–4192, 2020.
- [3] K. Spoon, S. Ambrogio, P. Narayanan, H. Tsai, C. Mackin, A. Chen, A. Fasoli, A. Friz, and G. W. Burr, "Accelerating deep neural networks with analog memory devices," in *2020 IEEE International Memory Workshop (IMW)*, 2020.
- [4] V. Ostwal, R. Zand, R. DeMara, and J. Appenzeller, "A novel compound synapse using probabilistic spin-orbit-torque switching for mtj-based deep neural networks," *IEEE Journal on Exploratory Solid-State Computational Devices and Circuits*, vol. 5, no. 2, pp. 182–187, 2019.
- [5] G. Molas, M. Harrand, C. Nail, and P. Blaise, "Advances in oxide-based conductive bridge memory (cbram) technology for computing systems," 2019.
- [6] S. Jain, A. Ranjan, K. Roy, and A. Raghunathan, "Computing in memory with spin-transfer torque magnetic ram," *IEEE Transactions on Very Large Scale Integration (VLSI) Systems*, vol. 26, no. 3, pp. 470–483, 2018.
- [7] S. Angizi, Z. He, A. Awad, and D. Fan, "Mrima: An mram-based in-memory accelerator," *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems*, vol. 39, no. 5, pp. 1123–1136, 2020.
- [8] S. Cosemans, B. Verhoef, J. Doevenspeck, I. A. Papistas, F. Catthoor, P. Debacker, A. Mallik, and D. Verkest, "Towards 10000tops/w dnn inference with analog in-memory computing – a circuit blueprint, device options and requirements," in *2019 IEEE International Electron Devices Meeting (IEDM)*, 2019, pp. 22.2.1–22.2.4.
- [9] A. Ankit, I. E. Hajj, S. R. Chalamalasetti, G. Ndu, M. Foltin, R. S. Williams, P. Faraboschi, W.-m. W. Hwu, J. P. Strachan, K. Roy, and D. S. Milojcic, "Puma: A programmable ultra-efficient memristor-based accelerator for machine learning inference," in *Proceedings of the Twenty-Fourth International Conference on Architectural Support for Programming Languages and Operating Systems*, ser. ASPLOS '19, New York, NY, USA: Association for Computing Machinery, 2019, p. 715–731. [Online]. Available: <https://doi.org/10.1145/3297858.3304049>
- [10] M. Elbitty, A. Singh, B. Reidy, X. Guo, and R. Zand, "An in-memory analog computing co-processor for energy-efficient cnn inference on mobile devices," *2021 IEEE Computer Society Annual Symposium on VLSI (ISVLSI)*, 2021, 2021.
- [11] M. Hu, J. P. Strachan, Z. Li, E. M. Grafals, N. Davila, C. Graves, S. Lam, N. Ge, J. J. Yang, and R. S. Williams, "Dot-product engine for neuromorphic computing: Programming 1t1m crossbar to accelerate matrix-vector multiplication," in *ACM/EDAC/IEEE Design Automation Conference (DAC)*, 2016, pp. 1–6.
- [12] F. L. Aguirre, S. M. Pazos, F. Palumbo, J. Suñé, and E. Miranda, "Application of the quasi-static memdiode model in cross-point arrays for large dataset pattern recognition," *IEEE Access*, vol. 8, pp. 202 174–202 193, 2020.
- [13] D. Josell, S. H. Brongersma, and Z. Tőkei, "Size-dependent resistivity in nanoscale interconnects," *Annual Review of Materials Research*, vol. 39, no. 1, pp. 231–254, 2009. [Online]. Available: <https://doi.org/10.1146/annurev-matsci-082908-145415>
- [14] K. Fuchs, "The conductivity of thin metallic films according to the electron theory of metals," *Mathematical Proceedings of the Cambridge Philosophical Society*, vol. 34, no. 1, p. 100–108, 1938.
- [15] A. F. Mayadas and M. Shatzkes, "Electrical-resistivity model for polycrystalline films: the case of arbitrary reflection at external surfaces," *Phys. Rev. B*, vol. 1, pp. 1382–1389, Feb 1970. [Online]. Available: <https://link.aps.org/doi/10.1103/PhysRevB.1.1382>
- [16] S. Rosnagel and T.-S. Kuan, "Alteration of cu conductivity in the size effect regime," *Journal of Vacuum Science & Technology B: Microelectronics and Nanometer Structures*, vol. 22, 01 2004.
- [17] W. Steinhögl, G. Schindler, G. Steinlesberger, M. Traving, and M. Engelhardt, "Comprehensive study of the resistivity of copper wires with lateral dimensions of 100 nm and smaller," *Journal of Applied Physics*, vol. 97, pp. 023 706–023 706, 12 2004.
- [18] T. Sun, B. Yao, A. P. Warren, K. Barnak, M. F. Toney, R. E. Peale, and K. R. Coffey, "Surface and grain-boundary scattering in nanometric cu films," *Phys. Rev. B*, vol. 81, p. 155454, Apr 2010. [Online]. Available: <https://link.aps.org/doi/10.1103/PhysRevB.81.155454>
- [19] T. Sakurai and K. Tamaru, "Simple formulas for two- and three-dimensional capacitances," *IEEE Transactions on Electron Devices*, vol. 30, no. 2, pp. 183–185, 1983.
- [20] Y. Lecun, L. Bottou, Y. Bengio, and P. Haffner, "Gradient-based learning applied to document recognition," *Proceedings of the IEEE*, vol. 86, no. 11, pp. 2278–2324, 1998.
- [21] "Predictive technology model," <http://ptm.asu.edu>.
- [22] S. Sinha, G. Yeric, V. Chandra, B. Cline, and Y. Cao, "Exploring sub-20nm finfet design with predictive technology models," in *DAC Design Automation Conference 2012*, 2012, pp. 283–288.
- [23] R. Zand, A. Roohi, and R. F. DeMara, "Fundamentals, modeling, and application of magnetic tunnel junctions," *Nanoscale Devices: Physics, Modeling, and Their Application*, p. 337, 2018.