

Exploring the Feasibility of Using 3-D XPoint as an In-Memory Computing Accelerator

MASOUD ZABIHI¹ (Graduate Student Member, IEEE), SALONIK RESCH¹, HÜSREV CILASUN¹,
ZAMSHED I. CHOWDHURY¹ (Graduate Student Member, IEEE), ZHENGYANG ZHAO¹,
ULYA R. KARPUZCU¹ (Member, IEEE), JIAN-PING WANG¹ (Fellow, IEEE), and
SACHIN S. SAPATNEKAR¹ (Fellow, IEEE)

Department of Electrical and Computer Engineering, University of Minnesota, Minneapolis, MN 55455 USA
CORRESPONDING AUTHOR: M. ZABIHI (zabih003@umn.edu)

This work was supported in part by the National Science Foundation (NSF) under Award CCF-1725420 and Award CCF-1763761.

This article has supplementary downloadable material available at <https://doi.org/10.1109/JXCDC.2021.3112238>, provided by the authors.

ABSTRACT This article describes how 3-D XPoint memory arrays can be used as in-memory computing accelerators. We first show that thresholded matrix-vector multiplication (TMVM), the fundamental computational kernel in many applications including machine learning (ML), can be implemented within a 3-D XPoint array without requiring data to leave the array for processing. Using the implementation of TMVM, we then discuss the implementation of a binary neural inference engine. We discuss the application of the core concept to address issues such as system scalability, where we connect multiple 3-D XPoint arrays, and power integrity, where we analyze the parasitic effects of metal lines on noise margins. To assure power integrity within the 3-D XPoint array during this implementation, we carefully analyze the parasitic effects of metal lines on the accuracy of the implementations. We quantify the impact of parasitics on limiting the size and configuration of a 3-D XPoint array, and estimate the maximum acceptable size of a 3-D XPoint subarray.

INDEX TERMS 3-D XPoint, in-memory computing, matrix-vector multiplication, neural network, phase-change memory (PCM).

I. INTRODUCTION

WITH the rapidly increasing sizes of datasets and challenges in transistor scaling in recent years, the need for new computing paradigms is felt more than ever [1]. In today's computing systems, large portions of computation energy and time are wasted for transferring data back and forth between the processor and the memory [2]. One approach is to bring the processor and memory closer to each other and build a near-memory platform that places the computing engine adjacent to the memory, and hence reduce the energy and time overhead for data transfer. Another approach that even more significantly reduces the time overhead and energy is to use the memory device as the computational unit and built a true in-memory computing platform. We follow the latter approach.

The substrate that we work on is 3-D XPoint [3], a class of memory technology that fills a unique place within the memory hierarchy between solid state storage drive (SSD) and the system main memory. In comparison with the NAND-based SSD (which is the most ubiquitous storage device available

today [4]), it has the advantage of being faster, denser, and more scalable. Its nonvolatility differentiates it from competing technologies such as NAND-based SSDs and dynamic random access memories (DRAMs), although NAND-based SSDs are more cost-effective today and DRAMs are faster. Other emerging nonvolatile technologies face limitations: stand-alone PCM must deal with resistance drift, where the cell resistance increases over time [5]; FeFET is handicapped by its large operating voltage and limited endurance [6]; MRAMs require an access transistor (unlike 3-D XPoint), leading to a larger cell size than 3-D XPoint; ReRAM is not commercially viable to the level of 3-D XPoint and MRAM.

The operation and performance of 3-D XPoint as a memory unit are discussed in [7]–[11]. In our work, rather than focusing again on the memory aspects of 3-D XPoint, we explore the possibility of exploiting 3-D XPoint arrays to perform in-memory computation. This means not only that 3-D XPoint can function as a storage unit, but also that it can perform computation inside its array without the need for the data to leave the array. Therefore, unlike conventional computational

systems, the information can be processed locally rather than being sent to a processor through the memory hierarchy.

The analysis in this article considers wire non-idealities and physical design of 3-D XPoint subarray. We first show the implementation of thresholded matrix-vector multiplication (TMVM), which is a building block for neural networks (NNs) and deep-learning applications. Second, using this core operation, we discuss the implementation of a neural network inference engine. Finally, we discuss how to enable 3-D XPoint for more complex versions of these implementations (e.g., multibit operations and multilayer NNs).

For in-memory computing platforms, wire resistances are a substantial source of non-ideality that must be taken into account during the implementations [12]. Some works attempt to analyze the parasitic effects of wires but do not consider all contributing factors with realistic layout considerations [13], [14]. In [15], a framework is presented to incorporate the effects of nonidealities in 2-D resistive crossbar performance. In [16], an analytical approach is developed to study the effects of the parasitic of wires for the implementation on spintronics computational RAM.

We discuss the feasibility of using 3-D XPoint as an in-memory computing engine for neuromorphic applications, and evaluate its performance for MNIST digit recognition. We present novel methods for the implementation of MVM and NN on 3-D XPoint. We use 3-D stacked PCM memory layers in the 3-D XPoint subarray to compute and store the computation results entirely inside the array, without sending the data to the periphery of the array. Our method is scalable by using multiple arrays to handle a large computational workload. In addition, multibit operations are supported in our methodology. We evaluate a realistic size of the 3-D XPoint subarray and metal features (based on the ASAP7 7-nm technology [17], [18]) for accurate electrical correctness. We develop a comprehensive method to analyze the impact of wire parasitics of wires in the 3-D XPoint subarray and devise a methodology to determine the maximum size of a 3-D XPoint subarray that ensures electrically correct operation.

Next, we discuss the structure of 3-D XPoint in Section II. In Section III, we describe the implementation of TMVM, and NN. In Section IV, we explore the methods for more complex implementations. We develop the models for the effect of wire parasitics in Section V, evaluate the results of our analysis in Section VI, and then conclude the article in Section VII.

II. BACKGROUND

Some recent works study the implementation of logic operations using a 3-D crossbar array architecture. In [19], a double-layer Pt/HfO_{2-x}/TiN ReRAM crossbar array is used and it is experimentally shown that the array can implement MVM and CNN. Additional peripheral circuitry (e.g., AND gates) are required for obtaining the computation results. In [20], it is shown that stateful logic operation can be performed on a memristive TiO₂-based 3-D crossbar array. The adverse effects of wire non-idealities are not incorporated in the implementations. Similarly, in [21], the authors map logic operations on a memristive 3-D crossbar without considering

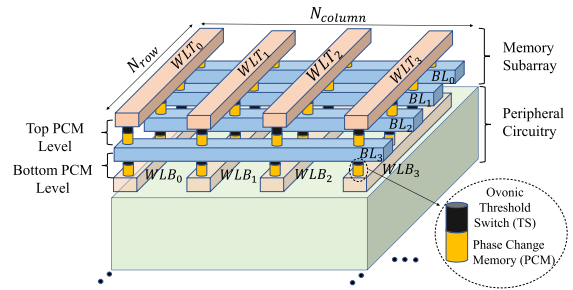


FIGURE 1. Structure of a 3-D XPoint subarray. The CMOS peripheral circuitry is located underneath the memory subarray.

wire parasitic or technology design rules. In [22], the authors use 3-D memristive crossbars for neuromorphic computation. For the implementation of a neuron, additional amplifier circuitry is required at the periphery of the crossbar.

Fig. 1 shows the overall structure of a 3-D XPoint subarray. A two-level PCM stack is integrated at the top of CMOS peripheral circuitry. The storage device is based on phase-change memory (PCM) technology, which is connected to a compatible ovonic threshold switch (OTS) made of AsTe-GeSiN [23]–[25]. Word lines at the top (WLTs), word lines at the bottom (WLBs), and bit lines (BLs) in the middle provide the current path to each individual memory cell [26]. The compatibility of the junction of PCM and OTS devices is a key factor in allowing access to individual cells without facing sneak path problems [27]. The total number of PCM cells in the 3-D XPoint subarray with N_{row} rows and N_{column} columns is $(2 N_{row} \times N_{column})$, with half in the top PCM level and half in the bottom PCM level, as shown in the figure.

PCM is a nonvolatile memory technology exploiting Ge–Sb–Te (GST) alloys (e.g., Ge₂Sb₂Te₅) as the storage medium [28]. PCM has two states: a crystalline phase with high conductance (G_C) and an amorphous phase with low conductance (G_A). The GST alloy transition between amorphous and crystalline states is triggered by changing the temperature level [29], [30]. In early explorations of PCM technology (1970s–early 2000s), the temperature level was changed using a laser source [31]. The state-of-the-art research on PCM is focused on using electrical impulses to change the temperature, and hence the state, of the PCM device by applying an electric current (or voltage) pulse across the PCM device [32].

Fig. 2(a) shows that applying a fast high-amplitude current pulse of amplitude I_{RESET} (called the RESET pulse) heats up the GST material to the melting temperature T_{melt} (~ 600 °C or higher [30]), erasing the previous periodic and ordered atomic arrangement. After quenching, the new disordered atomic structure will be frozen, making the transition from high conductance crystalline state to low conductance amorphous state possible. To change the state of the GST from amorphous to crystalline, a slow, relatively low amplitude current pulse of amplitude I_{SET} (called the SET pulse) must pass through the GST material. The SET current pulse causes the GST material to heat up to crystalline temperature T_{cryst} (~ 400 °C [30]). Over a long SET time of several tens of nanoseconds, this is a high enough temperature (still lower

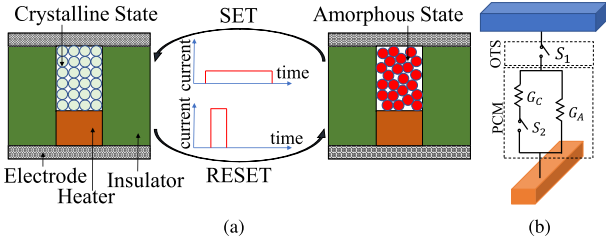


FIGURE 2. PCM model: (a) transition between amorphous and crystalline phases by applying SET and RESET pulses across a pillar-type PCM device, and (b) PCM cell can be modeled using a resistive circuit with two voltage control switches [9].

than T_{melt}) for the reconfiguration and crystallization of the previous amorphous atomic region to the crystalline state. The desirable PCM characteristics are a lower amplitude of the RESET current and a shorter SET time. A RESET current as low as $10 \mu\text{A}$ and a SET time as low as 25 ns for individual PCM devices is already demonstrated with sub-20 nm scalability, high endurance 10^{12} cycles, and a projected ten-year retention time at 210°C [5].

Fig. 2(b) shows the electrical model of PCM cell. The resistance across the PCM cell can be modeled by two voltage-controlled switches [9]. Depending on the status of switches S_1 and S_2 , different currents flow between two lines connected to the terminals of the PCM cell, determined by G_A and G_C . The ON/OFF states of the memory cell are determined by OTS: If the voltage level across the OTS of a cell is larger than a threshold, the cell is considered to be ON, and it is OFF otherwise. In today's technologies, the OTS conductance for the OFF state is up to $10^8 \times$ smaller than for the ON state.

The value stored in the PCM device can represent either logic 1 (crystalline phase) or logic 0 (amorphous phase). Three memory operations available in 3-D XPoint: write logic 1 (using the fast high-amplitude SET pulse), write logic 0 (using long low-amplitude RESET pulse), and read. For the memory read operation, since it is undesirable to change the state of the PCM cell, a pulse with relatively very small amplitude will be applied, increasing the temperature slightly above the ambient temperature but below T_{cryst} (and of course T_{melt}).

III. REALIZATION OF IN-MEMORY COMPUTING

A. IMPLEMENTATION OF TMVM

TMVM is a fundamental step in the implementation of many applications, and is a fundamental computational kernel in machine learning (ML) applications. Using 3-D XPoint as the TMVM computation engine can tremendously decrease the ML computational workload, as the data do not need to leave the 3-D XPoint array during the computation.

To show how the first step of a TMVM, let us multiply, without thresholding, matrix $G \in \mathbb{R}^{(N_x+1) \times (N_y+1)}$ and vector $V = [V_0 V_1 V_2, \dots, V_{N_x}]^T \in \mathbb{R}^{(N_x+1)}$, where G is given by

$$G = \begin{pmatrix} G_{0,0} & G_{0,1} & \cdots & G_{0,N_y} \\ G_{1,0} & G_{1,1} & \cdots & G_{1,N_y} \\ \vdots & \vdots & \cdots & \vdots \\ G_{N_x,0} & G_{N_x,1} & \cdots & G_{N_x,N_y} \end{pmatrix}. \quad (1)$$

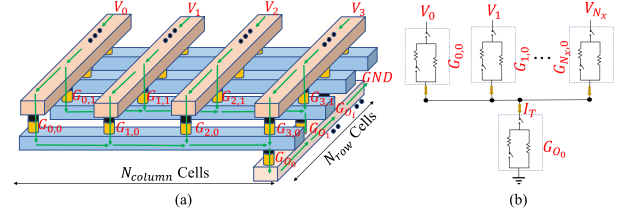


FIGURE 3. (a) Using 3-D XPoint as an in-memory computing engine for TMVM of GV. (b) Equivalent circuit model for the implementation of a dot product (to calculate O_0).

This computes $O = [O_0 O_1 O_2, \dots, O_{N_y}]^T \in \mathbb{R}^{(N_y+1)}$ where each element of vector O is a dot product. For example

$$O_0 = G_{0,0}V_0 + G_{1,0}V_1 + \cdots + G_{N_x,0}V_{N_x}. \quad (2)$$

This is computed in the 3-D XPoint array by applying voltages across a set of conductances to produce a current O_0 .

Today's PCM cells can only store binary values. Hence, we assume that elements of matrix G and vector V are binary. To implement a "neuron-like" operation using TMVM, the O_0 value and computed O_i values are followed by a thresholding operation. In (2), if the sum of products exceeds the current required to flip the output bit, then logic 1 is stored as the conductance, G_{O_0} , of the PCM cell O_0 ; otherwise, the stored logic value is 0. Similarly, for other O_i s, the values after thresholding are stored as the conductance states, G_{O_i} .

Fig. 3(a) shows the implementation of the TMVM on a 3-D XPoint subarray with $(2N_{\text{row}} \times N_{\text{column}})$ PCM cells ($(N_{\text{row}} \times N_{\text{column}})$ cells each at top PCM level and bottom PCM level) where $N_{\text{column}} = N_x + 1$ and $N_{\text{row}} = N_y + 1$. For clarity, as compared to Fig. 1, only the lines and PCM cells engaged in the computation are shown, and the rest of the lines and the PCM cells (at the bottom) are removed from the figure. All elements of O will be calculated simultaneously and are stored in the same column with N_{row} PCM cells. Considering that today's 3-D XPoint cannot store multiple values in a cell, we assume that elements of vector V and G are binary.

The conductances G are first programmed in the top PCM level by memory write operations or by previous computation.

- 1) Before the computation starts, cells that store G_{O_i} s at the bottom are preset to logic 0.
- 2) Then, voltages V_i , $0 \leq i \leq N_x$ are applied to the word lines WLTs connected to input cells located at top. If V_i represents logic 1, voltage V_{DD} is applied ($V_i \leftarrow V_{\text{DD}}$) to the WLT $_i$ and the current that flows through the corresponding input cell is proportional to $G_{0,i}V_{\text{DD}}$.
- 3) If V_i represents logic 0, WLT $_i$ is floated ($V_i \leftarrow \text{float}$) and no current passes through the corresponding PCM cell.
- 4) The summation of currents (I_T) from input cells flows through the G_{O_0} in a time interval t_{SET} . Based on the values of V_i and $G_{i,0}$, different currents pass through the input cells that store $G_{i,0}$. If $I_T > I_{\text{SET}}$, the state of G_{O_0} changes to logic 1. However, we require $I_T < I_{\text{RESET}}$ to avoid erroneous computation.

To calculate the minimum and maximum allowable applied voltage (V_{DD}) to the lines, we consider a simplified electrical model for the implementation of a dot product (e.g., for O_0) shown in Fig. 3(b). Current I_T can be written as follows:

$$I_T = G_{O_0} \frac{\sum_{i=0}^{N_x} V_i G_{i,0}}{\sum_{i=0}^{N_x} G_{i,0} + G_{O_0}}. \quad (3)$$

When the computation begins, $G_{O_0} \approx G_A$ since the preset is 0, and $I_T(t = 0)$ is small (of the order of few hundred nAs). However, by the passage of time, the amorphous region near the heater in the output PCM starts to turn crystalline, resulting in increasing G_{O_0} (and consequently I_T) and heat (generated by the flow of more electric current). If the applied voltage V_{DD} is large enough to provide a current larger than I_{SET} , crystallization repeats until a threshold point where the whole amorphous region in the output PCM turns into a crystalline region with high conductivity, representing logic 1. On the other hand, the V_{DD} must not be so large that the generated temperature exceeds T_{melt} , causing erroneous computation.

To calculate the V_{DD} range for the accurate implementation of the described dot product, we analyze the two cases corresponding to V_{min} (the minimum acceptable voltage) and V_{max} (the maximum acceptable voltage). For the V_{min} case, we assume that all V_i s, and all $G_{i,0}$ s represent logic 1, i.e., V_{DD} voltages are applied to all WLTs and the conductances of input cells are in the high conductance state corresponding to G_C . In this case, from (3), $I_T = (N_x + 1)/(N_x + 2)G_C V_{DD}$. Since $I_{SET} \leq I_T \leq I_{RESET}$, the V_{min} requirement implies a first constraint, requiring that V_{DD} to lie in the range

$$\mathcal{R}_1 = \left[\left(\frac{N_x + 2}{N_x + 1} \right) \left(\frac{I_{SET}}{G_C} \right), \left(\frac{N_x + 2}{N_x + 1} \right) \left(\frac{I_{RESET}}{G_C} \right) \right]. \quad (4)$$

For the V_{max} case, all V_i s are set to logic 1, while all $G_{i,0}$ s are set to logic 0. Since the result of the dot product should be at logic 0, we expect that the preset value of PCM stored O_0 remains intact. At the maximum voltage possible, we hypothetically assume that the applied voltage pulse should be below the level required to change the output state from logic 0 to logic 1, even if conductances of all input cells are G_A . In other words, $I_T = ((N_x + 1)G_A G_C)/((N_x + 1)G_A + G_C)V_{DD} < I_{SET}$, i.e., the output state cannot be altered. Therefore, the second set of constraints require V_{DD} to lie in the range

$$\mathcal{R}_2 = \left[0, \left(\frac{(N_x + 1)G_A + G_C}{(N_x + 1)G_A G_C} \right) I_{SET} \right]. \quad (5)$$

The acceptable range for V_{DD} is $\mathcal{R}_1 \cap \mathcal{R}_2$. Therefore, the minimum and maximum acceptable voltages are $V_{min} = \min(\mathcal{R}_1)$ and $V_{max} = \min(\max(\mathcal{R}_1), \max(\mathcal{R}_2))$, respectively.

B. IMPLEMENTATION OF NN

Using the TMVM implementation, we implement a neuro-morphic inference engine. Fig. 4(a) shows a single-layer NN with N inputs and P outputs. Fig. 4(b) shows the data layout for the NN implementation on a 3-D XPoint subarray. The top PCM cells are allocated for storing the weights ($W_{i,j}$ s), similar

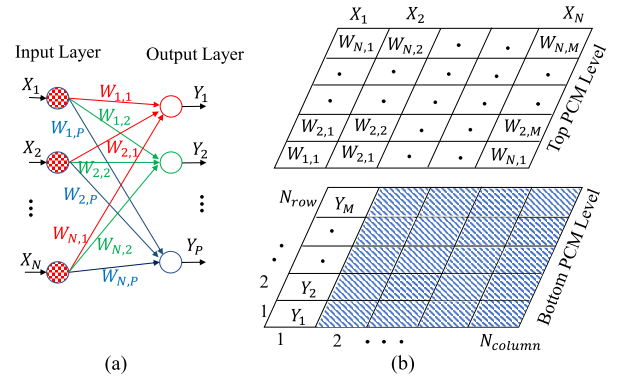


FIGURE 4. NN implementation: (a) Single-layer neuromorphic inference engine. (b) Data layout for the NN implementation.

to $G_{i,j}$ s in TMVM that were stored in the top PCM cells, and the bottom PCM cells are allocated for storing the outputs (Y_j s), similar to O_j s in TMVM. The inputs (X_i s) are applied to WLTs as voltage pulses (similar to V_i s in TMVM). If $N \leq N_{column}$ and $P \leq N_{row}$, then all Y_j s can be determined simultaneously in one step. The output elements of the NN can be stored in any column at the bottom (here, we choose column 1). In Fig. 4(b), all other cells at the bottom patterned by diagonal stripes are not engaged in the computation of Y_j s, i.e., WLBs connected to these cells are floated.

An application for the proposed NN implementation is handwritten digit recognition of MNIST dataset with 10-K test images [33]. Analyzing each MNIST test image can be performed using a similar NN shown in Fig. 4(a). Here $P = 10$, as in MNIST each image represents a digit (from 0 to 9). In each computational step, $\lfloor (N_{row}/P) \rfloor$ images can be processed and stored in a column.

IV. ENABLING MORE COMPLEX IMPLEMENTATIONS

In this section, we discuss three concepts that enable more complex computations. Then, we provide the implementation of a multilayer NN as an example. Our approach can also be extended to perform multibit computation directly using the principles in [34].

A. 3-D XPOINT WITH FOUR STACKED LEVEL OF PCM CELLS

Industry projections show that the next generation of 3-D XPoint will have four-level stacked PCM cells [35]. If the number of PCM levels increases, then the volume of stored information per footprint area increases, and more complex implementations are possible. Although a two-level subarray of PCM cells is sufficient to implement any NN (see Section IVD), we will illustrate how we can use a four-level subarray of PCM cells to implement a multilayer NN by exploiting the extra PCM levels. The NN in Fig. 5 has three layers: an input layer, a hidden layer, and an output layer. At the top PCM level, the first set of weights are stored. In the next PCM level, the hidden layer data are calculated, and by applying the second set of weights as voltage pulses, we obtain the outputs Y_i of the NN at the third PCM level.

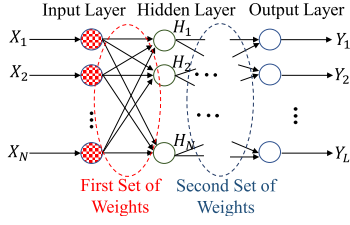


FIGURE 5. Multilayer NN with an input, hidden, and output layer.

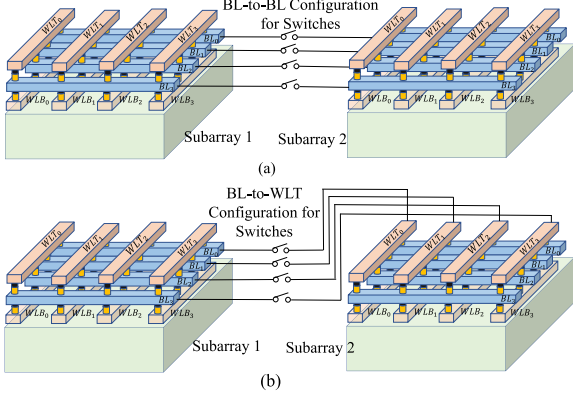


FIGURE 6. Two configurations for communication between 3-D XPoint subarrays: (a) switches connect BLs of subarray 1 to BLs of Subarray 2, and (b) switches connect BLs of the subarray 1 to WLTs of subarray 2.

B. SCALABILITY OF 3-D XPoint TO LARGE COMPUTATIONS

We can connect multiple 3-D XPoint subarrays to create a larger array to handle computations with higher matrix dimensions. In Fig. 6(a), switches connect BLs of subarray 1 to those of subarray 2, enabling current flow from the BLs of subarray 1 to those of subarray 2. The WLB of subarray 2 that is scheduled to store the computation results will be connected to ground, while all other WLBs not engaging in the computation (in both subarrays 1 and 2) are floated. Hence, the computation results in subarray 1, are being calculated and stored at the bottom PCM level of subarray 2 (BL-to-WLT). In Fig. 6(b), switches connect BLs of the subarray 1 to WLT of subarray 2. In this configuration, the results are being calculated at the top PCM level of subarray 2. The status of lines during the computation for these two configurations are listed in the Supplementary Materials.

C. MULTILAYER NN IMPLEMENTATION IN A TWO-LEVEL PCM STACK

We now illustrate how a multilayer NN can be implemented using a three-layer PCM stack. As an example, we discuss the implementation of the three-layer NN (shown in Fig. 5) using two two-level 3-D XPoint subarrays. Let us assume that the NN is required to analyze 10-K images of the MNIST dataset. The data layout of this implementation is illustrated in Fig. 7 using two subarrays connected with BL-to-WLT configuration [see Fig. 6(b)]. The first set of weights is stored at the top PCM cells of subarray 1. The inputs (X_0, X_1, \dots, X_N) are applied as the voltages to the WLTs of

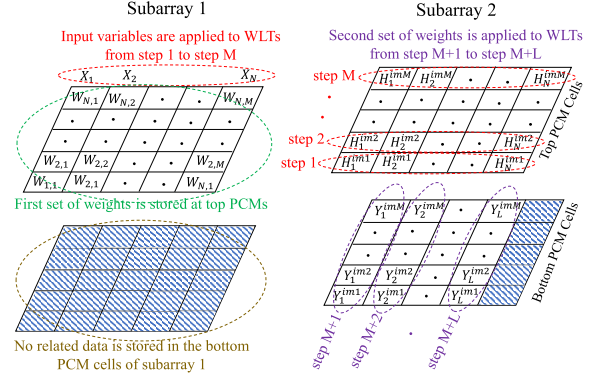


FIGURE 7. Data layout for the implementation of three-layer NN.

subarray 1. We assume that at each time step, the hidden layer values (H_1, H_2, \dots, H_N) for a specific image from MNIST are being processed. For example, the hidden layer values of the second image ($H_1^{im2}, H_2^{im2}, \dots, H_N^{im2}$) is calculated in step 2. Assuming that we calculate the hidden layer values of $M (=N_{\text{row}})$ images in each set of computation, we require M steps to calculate and store the hidden layer values of M images at the top PCM cells of subarray 2. In each of these steps, the corresponding BL in subarray 2 is connected to GND, and the remaining BLs in subarray 2 are floated. After all hidden layer values are stored at the top PCM cells of subarray 2, we apply the second set of the weights (as voltage pulses) to the WLTs of subarray 2. At each column at the bottom PCM cells of subarray 2, the outputs (Y_i s) of M images are calculated and stored.

V. ANALYZING INTERCONNECT PARASITIC EFFECTS

To ensure the electrical correctness of the implementations in in-memory computing platforms, we must consider non-idealities due to wire parasitic effects [12], [16]. As an example, we consider the implementation of a TMVM illustrated in Fig. 3(a). In the equivalent circuit model shown in Fig. 8, the WLTs, BLs, and WLBs have nonzero parasitics that cause a voltage drop in the current path across the 3-D XPoint subarray, that may potentially lead to errors in the results of TMVM. Let G_x and G_y be the conductances of the segments of BLs and WLTs, respectively. The conductances for WLT and WLB are considered equal (both G_y) due to the symmetry and equal allocation of metal resources to WLTs and WLBs. We use G_{ij} to denote the conductance of PCM cell (i, j) at the top level, and G_{Oj} s to denote conductances of a column of PCM cells at the bottom level. In the worst case, each row performs an identical operation, and carries an equal current I_{row} . The total voltage drop to the last row is

$$\frac{I_{\text{row}}}{G_y} + \frac{2I_{\text{row}}}{G_y} + \dots + \frac{N_{\text{row}}I_{\text{row}}}{G_y} = \frac{N_{\text{row}}(N_{\text{row}} + 1)I_{\text{row}}}{2G_y} \quad (6)$$

where the first, second, and last terms on the left side of the equation are for voltage drops of Segment $_{N_{\text{row}}}$, Segment $_{N_{\text{row}}-1}$, and Segment $_1$, respectively. The voltage drop of the last row increases quadratically with the number of rows, and this causes a significant limit on the accuracy of the implementations [12], [16]. Hence, it is important to find

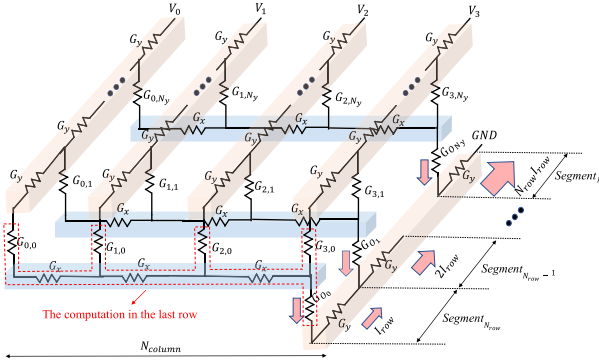


FIGURE 8. Equivalent circuit model for the TMVM implementation with considering wire parasitics.

the maximum allowable subarray size in which the voltage drop does not impair the electrical of implementations.

During the computation, the resistive network shown in Fig. 8 can have different configurations based on the applied voltage to WLTs. For example, if $V_2 \leftarrow \text{float}$, then all $G_{2,j}$ s and their connected parasitics must be removed from the equivalent circuit model of Fig. 8. To analyze parasitic effects, we consider the corner case for voltage drop, where only $V_0 \leftarrow V_{DD}$ and the rest of the V_i s are floated, resulting in minimum equivalent conductance for inputs and wire parasitics. Moreover, for the corner case, we assume that inputs and outputs are located N_{column} columns away from each other (the farthest possible distance). The value of all inputs assumed to be 1, and therefore, the current flows from the inputs of the TMVM computation must be sufficient enough to change the state of the output of the TMVM computation. An excessive voltage drop across the input and output cells causes a failure in the TMVM implementation discussed earlier.

The rows far from the drivers have larger parasitics between them and the driver. In particular, for the last row (farthest from the driver, see Fig. 8), the voltage drop is the worst. If the electrical correctness for the last row does not hold up, the implementations would be unacceptable. We observe the rest of the circuit from the last row and calculate (for the worst case) the Thevenin resistance (R_{th}) and Thevenin voltage (V_{th}) [see Fig. 9(a)]. We define the Thevenin coefficient, $\alpha_{\text{th}} = (V_{\text{th}}/V_{DD})$, and its value is between 0 and 1. Both R_{th} and α_{th} can be obtained analytically using a recursive approach explained in Appendix. Both are functions of parameters such as N_{row} , N_{column} , PCM cell width (W_{cell}) and length (L_{cell}) as well as other parameters of PCM and wire devices. Fig. 9(b) and (c) shows R_{th} and α_{th} for different N_{row} values. For the smaller N_{row} , PCM resistances are the dominant resistances, and the parasitic effect of wires is minimum. When the N_{row} increases, the values of the collective parasitic resistances become comparable to those of PCM devices and hence can degrade the electric correctness of the subarray. The configuration of lines is based on configuration 1 that will be discussed in Table 1 in the next section.

There are negligible parasitics between the first row and the driver, and the voltage range that ensures accuracy of

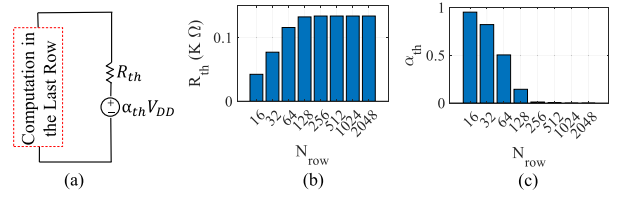


FIGURE 9. (a) Thevenin equivalents can be observed from the last row, (b) effects of N_{row} on R_{th} , (c) and on α_{th} .

TABLE 1. Different configurations of metal lines in the 3-D XPoint subarray based on ASAP7 design rules.

Config	WLT	WLB	BL	$W_{\text{min}} \times L_{\text{min}}$
1	M3	M1	M2	36nm×36nm
2	M3, M6, M8	M1, M7, M9	M2, M4, M5	48nm×80nm
3	M3, M5, M6, M8	M1, M4, M7, M9	M2	36nm×80nm

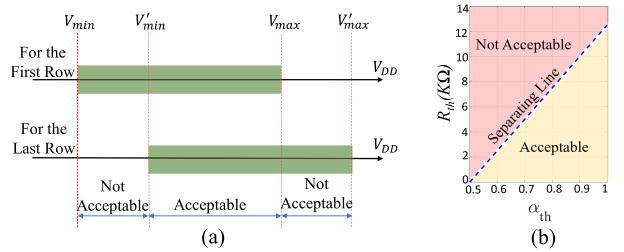


FIGURE 10. (a) Calculated voltage ranges for the first and last rows. (b) Acceptable and unacceptable regions in the $(\alpha_{\text{th}}, R_{\text{th}})$ plane.

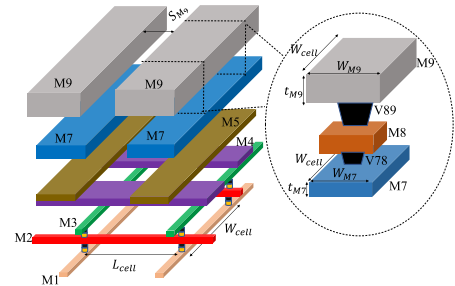


FIGURE 11. Multimetal layer configuration can be utilized for the design of WLTs, BLs, and WLBs of 3-D XPoint subarray.

computing in the first row is closer to $[V_{\text{min}}, V_{\text{max}}]$ (discussed in Section III) than that of the last row. For the last row, values of α_{th} and R_{th} are significantly affected by parasitics. Let us assume that the new voltage range ensures electrical correctness of the last row is $[V'_{\text{min}}, V'_{\text{max}}]$. The voltage ranges for the first row and last row are shown in Fig. 10(a). We use the voltage ranges of first row and last row as two corner cases (with least and most voltage drops, respectively), and we find a voltage range that satisfy both corner cases; the obtained voltage range guarantees the electrical correctness for intermediate rows as well. The final acceptable voltage range is the overlap between two voltage ranges shown in Fig. 10(a), $[V'_{\text{min}}, V_{\text{max}}]$, ensuring all of the rows from first row to the last row receiving the proper voltage.

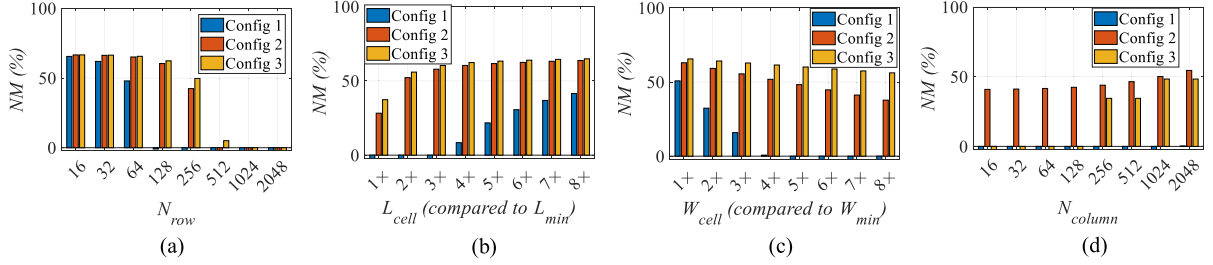


FIGURE 12. NMs of the three metal line configurations: (a) changing N_{row} (while $N_{column} = 128$, $L_{cell} = 4L_{min}$, and $W_{cell} = W_{min}$), (b) changing L_{cell} (while $N_{column} = 128$, $N_{row} = 128$, and $W_{cell} = W_{min}$), (c) changing W_{cell} (while $N_{column} = 128$, $N_{row} = 64$, and $L_{cell} = 4L_{min}$), and (d) changing N_{column} (while $N_{row} = 256$, $L_{cell} = 4L_{min}$, and $W_{cell} = W_{min}$).

The noise margin (NM) in implementations is defined by

$$NM = \frac{V_{max} - V'_{min}}{V_{mid}} \quad (7)$$

where $V_{mid} = (V_{max} + V'_{min})/2$. Clearly, we desire $NM \geq 0$. In Fig. 10(b), the acceptable and unacceptable regions in the (α_{th}, R_{th}) plane is shown. The NM on the separating line is 0, above it NM is negative (unacceptable), and below it NM is positive (acceptable). Our goal is to choose wire configurations so that the corresponding (α_{th}, R_{th}) of the design falls into the acceptable region with maximum NM possible.

VI. RESULTS AND DISCUSSION

A. NM EVALUATION

To realistically analyze the effect of the parasitics, we assumed that metal layers in 3-D XPoint are constructed based on ASAP7 design rules [17], [18] (see Fig. 11). We can create different configurations for allocating metal lines to WLTs, WLBS, and BLs. Table 1 lists three possible configurations. In configuration 1, only M1, M2, and M3 (the first three metal lines) in ASAP7 are exploited for 3-D XPoint, and they are allocated to WLB, BL, and WLT, respectively. For configurations 2 and 3, we assume that other than M1, M2, and M3, the other metal layers (M4 to M9) can also be allocated to the 3-D XPoint lines. In configuration 2, we allocate M4 and M5 to the BLs, and M6 to M9 are allocated equally between WLTs and WLBS. In configuration 3, we assume that all metals from M4 to M9 are allocated equally between WLTs and WLBS; no extra top metal lines are allocated to BLs. We report the minimum cell width (W_{min}) and length (L_{min}) for each configuration based on the minimum required width of a line and space between adjacent lines in each layer. The values of parameters for metal lines and PCM devices are available in the Supplementary Material.

1) NM IMPROVES WITH INCREASING N_{row}

Fig. 12(a) shows NMs of different N_{row} values. NM is significantly sensitive to N_{row} . For N_{row} as large as 2048, the implementations are not valid due to excessive voltage drop, and hence negative NM. Configuration 3 provides the best NM, because more metal resources dedicated to WLT and WLB causes smaller parasitics in the current path across rows.

2) NM IMPROVES WITH INCREASING L_{cell}

Fig. 12(b) shows the NMs for different L_{cells} (for each configuration, values are normalized to L_{min} listed in Table 1). By increasing L_{cell} , the width of the WLTs and WLBS increases, decreasing the parasitic resistances related to WLTs and WLBS.

3) NM DECREASES WITH INCREASING W_{cell}

Fig. 12(c) shows the NMs for different W_{cell} (for each configuration, values are normalized to W_{min} listed in Table 1). By increasing W_{cell} , the length of the WLTs and WLBS increases, and consequently, parasitics related to WLTs and WLBS considerably increase. Therefore, for all cases, smaller W_{cell} causes larger NM.

4) NM REMAINS UNCHANGED WITH INCREASING N_{column}

Fig. 12(d) shows that the increase in N_{column} does not affect NM significantly. By increasing N_{column} , parasitics of BLs increase. However, since the BL resistances are in series with those of PCM devices with orders of magnitude larger resistance, the increase in BL resistance does not affect NM.

B. IMPLEMENTING NNs ON THE 3-D XPoint SUBSTRATE

We list the performance of various 3-D XPoint subarrays of various sizes for the digit recognition of MNIST dataset in Table 2. Each MNIST test image is scaled to 11×11 as in [36], a transformation that maintains 91% recognition accuracy and reduces computation. We use configuration 3 in all cases, as it provides the best NM among all alternatives. For the smallest subarray with size 64×128 , NM is the maximum among all cases. For the largest subarray with size 1024×1024 , we increase L_{cell} by $2.6\times$ (compare to that of 64×128 subarray) to decrease the parasitics of lines. Consequently, we achieve acceptable NM of 34.5%. With this relatively large subarray, we have more parallelism that allow to process a larger number of MNIST images in each computational step, reducing the total execution time ($17\times$ faster than that of 64×128 subarray). The energy per image is similar for all cases because the subarray sizes listed in Table 2 are large enough to allow fully processing an 11×11 MNIST image locally without extra data movement between subarrays or peripheral circuitry.

The table also shows the impact of 10% process variation on NM. In the second to the last column, interconnect resistances are changed by 10%. For small arrays, the change in

TABLE 2. Evaluation of different subarray sizes for digit recognition application.

Subarray Size	Cell Size (nm \times nm)	#Image per Step	Energy per Image	Subarray Area (μm^2)	Execution Time (μs)	V_{max} (V)	V'_{min} (mV)	NM	NM w/ 10% Var.
64 \times 128	36 \times 240	6	21.5pJ	62.9	133.3	1.25	636.2	65.1%	64.9%
128 \times 256	36 \times 320	12	21.5pJ	335.5	66.7	1.25	650.6	63.1%	62.7%
256 \times 512	36 \times 400	25	20.7pJ	1677.7	32.0	1.25	681.0	58.9%	58.1%
512 \times 1024	36 \times 480	51	20.3pJ	8053.0	15.7	1.25	732.5	52.2%	50.8%
1024 \times 2048	36 \times 640	102	20.3pJ	42949.6	7.8	1.25	882.2	34.5%	31.5%

NM is negligible because the resistance of the PCM element and OTS dominate wire resistance. For larger arrays, this effect becomes more noticeable (but is still not significant). In the last column, PCM device parameters as well as interconnect resistances are simultaneously varied by 10% to find the worst NM due to variations of all parameters. NMs for all arrays are still positive and acceptable (between 12.4% for the largest subarray and 46.6% for the smallest subarray).

VII. CONCLUSION

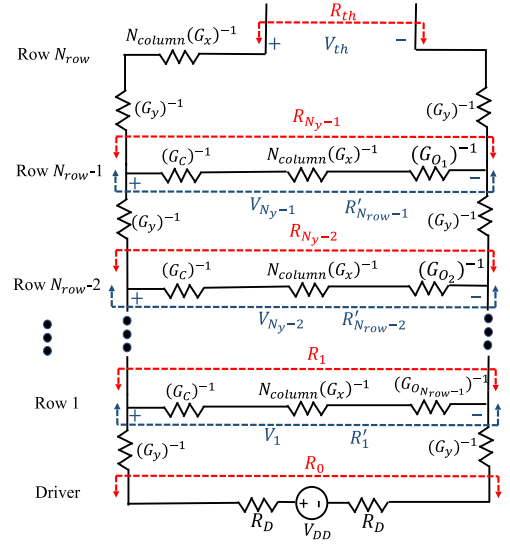
We have presented methods for the implementations of TMVM, NN, and 2-D convolution on 3-D XPoint. To ensure the accuracy of the implementations, we considered wire parasitics in our implementations. We have demonstrated that interconnect parasitics have a significant effect on the implementations performance and have developed a comprehensive model for analyzing this impact. Using this methodology, we have developed guidelines for the 3-D XPoint Subarray size and configurations based on ASAP7 technology design rules. We used different size 3-D XPoint subarrays for digit recognition of MNIST dataset. Using the our methodology methodology, we design a relatively large subarray of 2 Mb with acceptable NM of 34.5%, providing the opportunity for processing more images per step without any energy overhead.

APPENDIX

We derive recursive expressions for calculating R_{th} and V_{th} of a ($N_{row} \times N_{column}$) subarray of 3-D XPoint.

Within the footprint area of a cell ($W_{cell} \times L_{cell}$), we define G_y (representing the conductance for WLT and WLB segments) and G_x (representing the conductance of BL segment). Fig. 13 shows the equivalent simplified circuit model for the implementation TMVM in the corner case. Row i is separated from its predecessor by conductance G_y at each end. The input cell is connected to the output cell, N_{column} columns away. The conductances in the last row are rearranged to create a two-port structure consisting of the PCM conductances so that the rest of the network can be modeled using Thevenin Equivalents (R_{th} and V_{th}).

For configuration 1 (listed in Table 1), $G_y = G_{M_1} = G_{M_3}$ (assuming similar wire conductance for WLTs and WLBs) and $G_x = G_{M_2}$ in which the conductance, G_{M_k} , is given by: $G_{M_k}^{-1} = (\rho_{M_k} L_{M_k}) / (t_{M_k} W_{M_k})$, where ρ_{M_k} , L_{M_k} , t_{M_k} , and W_{M_k} are, respectively, the resistivity, length, thickness, and width in metal layer k (see the Supplementary Material). For configurations 2 and 3, the equivalent conductance of the wire segment must be calculated based on the multimetal layer


FIGURE 13. Notations used for calculating Thevenin resistance (R_{th}) and Thevenin voltage (V_{th}) shown on the circuit model for the implementation of TMVM in the worst case scenario.

configuration of a given segment. For example, in configuration 2, G_y (representing a segment conductance of WLT) is obtained by $G_y = G_{M_3} + G_{M_6} + G_{M_8}$.

To calculate R_{th} and V_{th} , we derive recursive expressions. For conciseness, we define the resistance, R_{row_i} , of row i as

$$R_{row_i} = N_{column} (G_x)^{-1} + (G_c)^{-1} + (G_{O_{N_{row}-i}})^{-1}. \quad (8)$$

We can obtain R_{th} , using the notations in Fig. 13, as

$$R_{th} = 2 (G_y)^{-1} + N_{column} (G_x)^{-1} + R_{N_{row}-1} \quad (9)$$

where $R_{N_{row}-1}$ is calculated using the recursive expression

$$R_i = (R_{row_i}) \left\| \left(R_{i-1} + 2 (G_y)^{-1} \right) \right. \quad (10)$$

The base case corresponds to the driver row that precedes the first row, and is $R_0 = 2 R_D$, as seen in Fig. 13.

To compute V_{th} , as illustrated in Fig. 13, we first compute the intermediate variable R'_j , which corresponds to the effective downstream resistance (away from the source) seen from node j . The computation proceeds in a recursive fashion from the last row toward the first as

$$R'_{j-1} = (R_{row_{j-1}}) \left\| \left(R'_j + 2 (G_y)^{-1} \right) \right. \quad (11)$$

with the base case $R'_{N_{row}-1} = R_{row_{N_{row}-1}}$. Having computed R'_j , we may now compute $V_{th} = V_{N_{row}}$, using a recursive computation on V_i

$$V_j = \frac{R'_j}{2 (G_y)^{-1} + R'_j} V_{j-1} \quad (12)$$

in which $2 \leq j \leq N_{row} - 1$ and the base case is

$$V_1 = \frac{R'_1}{R'_1 + 2 (G_y)^{-1} + 2 R_D} V_b. \quad (13)$$

REFERENCES

- [1] A. McAfee, E. Brynjolfsson, T. H. Davenport, D. J. Patil, and D. Barton, "Big data: The management revolution," *Harvard Bus. Rev.*, vol. 90, pp. 60–68, Oct. 2012.
- [2] S. W. Keckler, W. J. Dally, B. Khailany, M. Garland, and D. Glasco, "GPUs and the future of parallel computing," *IEEE Micro*, vol. 31, no. 5, pp. 7–17, Sep. 2011.
- [3] *3D XPoint Technology*. Accessed: Apr. 2017. [Online]. Available: <https://www.intel.com/content/www/us/en/architecture-and-technology/intel-micron-3d-xpoint-webcast.html>
- [4] N. Agrawal, V. Prabhakaran, T. Wobber, J. D. Davis, M. Manasse, and R. Panigrahy, "Design tradeoffs for SSD performance," in *Proc. USENIX Annu. Tech. Conf.*, Jun. 2008, pp. 57–70.
- [5] M. Le Gallo and A. Sebastian, "An overview of phase-change memory device physics," *J. Phys. D, Appl. Phys.*, vol. 53, no. 21, May 2020, Art. no. 213002.
- [6] S. Beyer *et al.*, "FeFET: A versatile CMOS compatible device with game-changing potential," in *Proc. IEEE Int. Memory Workshop (IMW)*, May 2020, pp. 1–4.
- [7] K. Son *et al.*, "Signal integrity design and analysis of 3-D X-Point memory considering crosstalk and IR drop for higher performance computing," *IEEE Trans. Compon., Packag., Manuf. Technol.*, vol. 10, no. 5, pp. 858–869, May 2020.
- [8] K. Son *et al.*, "Modeling and verification of 3-dimensional resistive storage class memory with high speed circuits for core operation," in *Proc. IEEE Asia-Pacific Microw. Conf. (APMC)*, Dec. 2019, pp. 694–696.
- [9] K. Son, K. Cho, S. Kim, G. Park, K. Son, and J. Kim, "Modeling and signal integrity analysis of 3D XPoint memory cells and interconnections with memory size variations during read operation," in *Proc. IEEE Symp. Electromagn. Compat., Signal Integr. Power Integr.*, Jul. 2018, pp. 223–227.
- [10] Q. Lou, W. Wen, and L. Jiang, "3DICT: A reliable and QoS capable mobile process-in-memory architecture for lookup-based CNNs in 3D XPoint ReRAMs," in *Proc. IEEE/ACM Int. Conf. Comput.-Aided Design*, Nov. 2018, pp. 1–8.
- [11] J. Yang, B. Li, and D. J. Lilja, "Exploring performance characteristics of the optane 3D XPoint storage technology," *ACM Trans. Model. Perform. Eval. Comput. Syst.*, vol. 5, no. 1, pp. 1–28, Feb. 2020.
- [12] D. Ielmini and G. Pedretti, "Device and circuit architectures for in-memory computing," *Adv. Intell. Syst.*, vol. 2, pp. 2000049–1–2000049–19, Jul. 2020.
- [13] Y. Jeong, M. A. Zidan, and W. D. Lu, "Parasitic effect analysis in memristor-array-based neuromorphic systems," *IEEE Trans. Nanotechnol.*, vol. 17, no. 1, pp. 184–193, Jan. 2018.
- [14] N. M. G. F. Lalchhandama, K. Datta, and I. Sengupta, "Modelling and simulation of non-ideal MAGIC NOR gates on memristor crossbar," in *Proc. 8th Int. Symp. Embedded Comput. Syst. Design (ISED)*, Dec. 2018, pp. 124–128.
- [15] S. Jain, A. Sengupta, K. Roy, and A. Raghunathan, "RxNN: A framework for evaluating deep neural networks on resistive crossbars," *IEEE Trans. Comput.-Aided Design Integr. Circuits Syst.*, vol. 40, no. 2, pp. 326–338, Feb. 2021.
- [16] M. Zabihi *et al.*, "Analyzing the effects of interconnect parasitics in the STT CRAM in-memory computational platform," *IEEE J. Explor. Solid-State Comput. Devices Circuits*, vol. 6, no. 1, pp. 71–79, Jun. 2020.
- [17] L. T. Clark, V. Vashishtha, D. M. Harris, S. Dietrich, and Z. Wang, "Design flows and collateral for the ASAP7 7 nm FinFET predictive process design kit," in *Proc. IEEE Int. Conf. Microelectron. Syst. Educ. (MSE)*, May 2017, pp. 1–4.
- [18] L. T. Clark *et al.*, "ASAP7: A 7-nm FinFET predictive process design kit," *Microelectron. J.*, vol. 53, pp. 105–115, Jul. 2016.
- [19] Y. Kim *et al.*, "Kernel application of the stacked crossbar array composed of self-rectifying resistive switching memory for convolutional neural networks," *Adv. Intell. Syst.*, vol. 2, no. 2, Feb. 2020, Art. no. 1900116.
- [20] N. Xu *et al.*, "A stateful logic family based on a new logic primitive circuit composed of two antiparallel bipolar memristors," *Adv. Intell. Syst.*, vol. 2, no. 1, Jan. 2020, Art. no. 1900082.
- [21] N. G. Murali, P. S. Vardhan, F. Lalchhandama, K. Datta, and I. Sengupta, "Mapping of Boolean logic functions onto 3D memristor crossbar," in *Proc. 32nd Int. Conf. VLSI Design, 18th Int. Conf. Embedded Syst. (VLSID)*, Jan. 2019, pp. 500–501.
- [22] B. R. Fernando, Y. Qi, C. Yakopcic, and T. M. Taha, "3D memristor crossbar architecture for a multicore neuromorphic system," in *Proc. Int. Joint Conf. Neural Netw. (IJCNN)*, Jul. 2020, pp. 1–8.
- [23] M.-J. Lee *et al.*, "Highly-scalable threshold switching select device based on chalcogenide glasses for 3D nanoscaled memory arrays," in *IEDM Tech. Dig.*, Dec. 2012, pp. 2.6.1–2.6.3.
- [24] G. W. Burr *et al.*, "Access devices for 3D crosspoint memory," *J. Vac. Sci. Technol. B*, vol. 32, pp. 040802-1–040802-23, Jul. 2014.
- [25] D. Kau *et al.*, "A stackable cross point phase change memory," in *IEDM Tech. Dig.*, Dec. 2009, pp. 1–4.
- [26] W. Chien *et al.*, "A study on OTS-PCM pillar cell for 3-D stackable memory," *IEEE Trans. Electron Devices*, vol. 65, no. 11, pp. 5172–5179, Nov. 2018.
- [27] L. Shi, G. Zheng, B. Tian, B. Dkhil, and C. Duan, "Research progress on solutions to the sneak path issue in memristor crossbar arrays," *Nanosci. Adv.*, vol. 2, no. 5, pp. 1811–1827, May 2020.
- [28] J. Y. Wu *et al.*, "A low power phase change memory using thermally confined TaN/TiN bottom electrode," in *IEDM Tech. Dig.*, Dec. 2011, pp. 3.2.1–3.2.4.
- [29] H.-S. P. Wong *et al.*, "Phase change memory," *Proc. IEEE*, vol. 98, no. 12, pp. 2201–2227, Dec. 2010.
- [30] G. W. Burr *et al.*, "Recent progress in phase-change memory technology," *IEEE J. Emerg. Sel. Topics Circuits Syst.*, vol. 6, no. 2, pp. 146–162, Jun. 2016.
- [31] M. Wuttig and N. Yamada, "Phase-change materials for rewriteable data storage," *Nature Mater.*, vol. 6, pp. 824–832, Nov. 2007.
- [32] Q. Zheng, Y. Wang, and J. Zhu, "Nanoscale phase-change materials and devices," *J. Phys. D, Appl. Phys.*, vol. 50, no. 24, May 2017, Art. no. 243002.
- [33] Y. LeCun. *The MNIST Database of Handwritten Digits*. Accessed: Nov. 1998. [Online]. Available: <http://yann.lecun.com/exdb/mnist/>
- [34] P. Chi *et al.*, "PRIME: A novel processing-in-memory architecture for neural network computation in ReRAM-based main memory," in *Proc. 43rd Int. Symp. Comput. Archit. (ISCA)*, Jun. 2016, pp. 27–39.
- [35] *3D XPoint With Four-Level PCM Cells*. Accessed: Aug. 2020. [Online]. Available: <https://www.anandtech.com/show/15972/intel-previews-4layer-3d-xpoint-memory-for-secondgeneration-optane-ssds>
- [36] M. Liu, L. R. Everson, and C. H. Kim, "A scalable time-based integrate-and-fire neuromorphic core with brain-inspired leak and local lateral inhibition capabilities," in *Proc. IEEE Custom Integr. Circuits Conf. (CICC)*, Apr. 2017, pp. 1–4.

• • •