

An accurate, error-tolerant and energy-efficient neural network inference engine based on SONOS analog memory

T. Patrick Xiao, *Member, IEEE*, Ben Feinberg, *Member, IEEE*, Christopher H. Bennett, *Member, IEEE*, Vineet Agrawal, Prashant Saxena, Venkatraman Prabhakar, Krishnaswamy Ramkumar, Harsha Medu, Vijay Raghavan, *Member, IEEE*, Ramesh Chettuvetty, Sapan Agarwal, *Member, IEEE*, Matthew J. Marinella, *Senior Member, IEEE*

Abstract—We demonstrate SONOS (silicon-oxide-nitride-oxide-silicon) analog memory arrays that are optimized for neural network inference. The devices are fabricated in a 40nm process and operated in the subthreshold regime for in-memory matrix multiplication. Subthreshold operation enables low conductances to be implemented with low error, which matches the typical weight distribution of neural networks, which is heavily skewed toward near-zero values. This leads to high accuracy in the presence of programming errors and process variations. We simulate the end-to-end neural network inference accuracy, accounting for the measured programming error, read noise, and retention loss in a fabricated SONOS array. Evaluated on the ImageNet dataset using ResNet50, the accuracy using a SONOS system is within 2.16% of floating-point accuracy without any retraining. The unique error properties and high On/Off ratio of the SONOS device allows scaling to large arrays without bit slicing, and enables an inference architecture that achieves 20 TOPS/W on ResNet50, a $>10\times$ gain in energy efficiency over state-of-the-art digital and analog inference accelerators.

Index Terms—SONOS, charge trap memory, neuromorphic, neural network, analog, in-memory computing, inference accelerator

I. INTRODUCTION

NEURAL networks are increasingly being used to solve important problems in many application domains. As these workloads become larger and more complex, a number of domain-specific architectures have been developed to address the challenges of scalable performance and energy efficiency. Digital processors that are specialized for matrix operations, such as graphics processing units (GPUs) [1], field-programmable gate arrays (FPGAs) [2], and application-specific chips like Google’s Tensor Processing Unit (TPU) [3] are already being deployed in data centers and edge devices to process neural networks.

Accelerators based on analog computation inside memory arrays potentially yield further order-of-magnitude improvements in energy efficiency. Performing matrix computations inside a memory array exploits both analog parallelism to reduce processing energy and data locality to reduce data movement energy [4], [5], [6]. Neural network inference

is a particularly attractive application for these accelerators, as it involves only a single matrix computation primitive (matrix-vector multiplication, or MVM) and does not require frequent programming of the memory devices. Many non-volatile memory technologies have been explored as candidate devices for analog inference. However, despite significant recent progress in both devices and architectures [7], [8], [9], analog systems have yet to match the precision of digital inference accelerators. Analog precision is directly affected by noise, process variations, and parasitic effects. Although many analog accelerators have demonstrated high accuracy on simple datasets like MNIST and CIFAR-10 in the presence of these errors, their viability for a more realistic benchmark like the ImageNet dataset [10] is rarely tested. Various works have proposed device-aware training as a solution for analog errors [11], [12], [13], but these methods may not always integrate easily into existing training workflows and do not always recover the full digital accuracy.

This work investigates SONOS (silicon-oxide-nitride-oxide-silicon) charge trap memory as the weight storage element for neural network inference. The SONOS device, which has been commercialized for data storage applications [14], behaves functionally as a memristor. Charge that is stored persistently in the SONOS gate stack serves as a state variable that electrostatically modulates the channel conductance. The embedded SONOS devices in this work were fabricated in a 40nm foundry process, and statistics from thousands of devices were used to characterize programming errors, retention loss, and read noise. This experimental data is used to simulate the accuracy of a SONOS accelerator on ImageNet with the ResNet50 convolutional neural network (CNN) [15], [16].

Our novel finding is that by operating the SONOS transistor in the subthreshold regime, we obtain device properties that are uniquely well matched to a universal property of neural networks. The exponential current-voltage characteristic under subthreshold operation implies that the SONOS device can be programmed to have a very high current On/Off ratio ($>10^7$) and furthermore, a conductance near zero can be programmed with nearly zero absolute error. At the same time, weights in a neural network tend to be exponentially skewed toward values with low magnitude. This combination allows the SONOS device to represent the most abundant weights in a network with the highest precision, enabling high accuracy

T. P. Xiao, B. Feinberg, C. H. Bennett, S. Agarwal, and M. J. Marinella are with Sandia National Laboratories, Albuquerque, NM, 87111 USA. e-mail: txiao@sandia.gov, m@asu.edu.

V. Agrawal, P. Saxena, V. Prabhakar, K. Ramukar, H. Medu, V. Raghavan, and R. Chettuvetty are with Infineon Technologies LLC, San Jose, CA.

and robustness to programming errors. The high On/Off ratio further provides robustness to parasitic voltage drops even in a large array, which enhances energy efficiency.

These unique properties allow the SONOS accelerator to achieve within 2.16% of floating-point accuracy on ImageNet without any network retraining – the most accurate result to date for an analog accelerator that is based on characterized devices. We then design the circuits and architecture for a SONOS inference accelerator with a projected peak energy efficiency that exceeds 50 TOPS/W.

The remainder of this paper is organized as follows. Section II provides a brief background on analog in-memory inference accelerators. Section III describes the SONOS device and its measured properties. Section IV discusses how the SONOS device is used for MVM, and highlights the key properties that are uniquely suited for the inference application. Section V uses the device data to simulate the end-to-end accuracy of a SONOS inference accelerator on ImageNet. Sections VI describes the design of the peripheral circuitry for the SONOS array, and Section VII describes the system architecture of the inference accelerator, focusing on CNNs that use rectified linear (ReLU) activations. Section VIII evaluates the projected performance, area, and energy efficiency of the accelerator on different workloads.

II. ANALOG IN-MEMORY ACCELERATORS FOR NEURAL NETWORK INFERENCE

Inference operations in neural networks rely heavily on MVM computations [17]. Digital inference accelerators can efficiently execute MVMs using arrays of multiply-accumulate (MAC) units, but their system-level efficiency is limited by the energy to move operands between memory and the processor [3], [18]. Analog in-memory computing reduces inference energy in two important ways. First, MVMs are computed in the same arrays where the network’s weights are stored, greatly reducing data movement energy. Second, individual MACs can be executed more efficiently in the analog domain, and the full MVM operation can be processed in parallel. As shown in Fig. 1, each memory cell multiplies an input element (voltage) and its stored weight (conductance) to produce a current. The voltage can be proportional to the input element or can be binary, representing one bit of the input at a time. These currents are summed along an array column using Kirchoff’s current law to produce a dot product. The analog dot products are then digitized using an analog-to-digital converter (ADC) and sent to the next layer’s array.

An important question in designing analog accelerators is the precision of the weight values stored in the memory cells, which ultimately determines the inference accuracy. Many memory technologies have been proposed for in-memory MVM, with varying degrees of precision: these include resistive random access memory (ReRAM) or memristors [19], [20], [21], phase change memory (PCM) [11], [22], flash memory [23], [24], [25], [26], [27], and ferroelectric memory [28]. The weight precision is also affected by architecture-level design choices. The previously cited works map one weight to one device, or to a differential pair of devices for

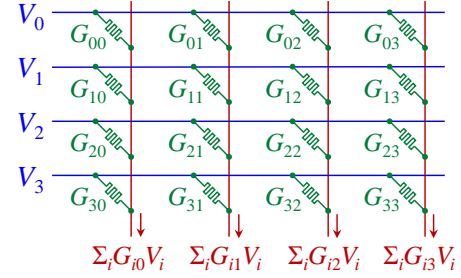


Fig. 1: Execution of an analog matrix-vector multiplication inside a resistive memory array.

signed weights. Many other works use bit slicing, where a B -bit weight is mapped to S devices that are each programmed to B/S bits of resolution [4], [5], [29], [30]. Bit slicing can represent arbitrary precisely weights with low-resolution devices, but it does not guarantee digital precision in the MVM computation [31]. Bit slices and their aggregation also incur significant area and energy overheads.

The amount of analog error that can be tolerated depends on the neural network application. Accuracy on simple image recognition tasks such as MNIST [32] and CIFAR-10 [33] can tolerate relatively large errors. A more realistic machine learning benchmark is the much more difficult ImageNet dataset [10]. ImageNet networks generally require about 8 bits of precision in the weights and activations [34]. Recently, Joshi *et al.* [11] showed (via simulation) that with a PCM device, the analog errors lead to a 7.8% drop in top-1 accuracy using ResNet34 [15]. After retraining the network to regularize for the expected magnitude of device errors, the accuracy drop was reduced to 1.6% relative to the original floating-point network. While device-aware retraining is a potentially viable option for analog inference, it may be complicated to integrate with state-of-the-art training workflows. Ideally, analog inference accelerators would provide near-digital accuracy with off-the-shelf neural networks without any retraining.

III. SONOS ANALOG MEMORY

Fig. 2 shows the embedded SONOS analog memory cell used in this work, which was fabricated in a 40nm foundry process. Each cell consists of a select transistor in series with a SONOS memory device. When used as an analog memory, the internal state variable of the SONOS device is the quantity of charge Q that resides in electronic trapping sites inside the silicon nitride layer. Charge is injected or removed from the nitride via Fowler-Nordheim tunneling of electrons or holes from the silicon channel [14], and is thereafter confined by the potential barriers of both the traps and the nitride-oxide interfaces. The amount of stored charge modulates the threshold voltage V_T of the SONOS transistor channel, which in turn controls the drain current I_D . While Q represents the state of the device, this state is electrically read out using the drain current I_D at a fixed bias. During an MVM, the bit line (BL), source line (SL), and control gate (CG) voltages are fixed. A binary voltage on the select gate (SG) switches the select transistor between its fully on and off states, and controls whether or not current flows between the SL and BL through the SONOS device.

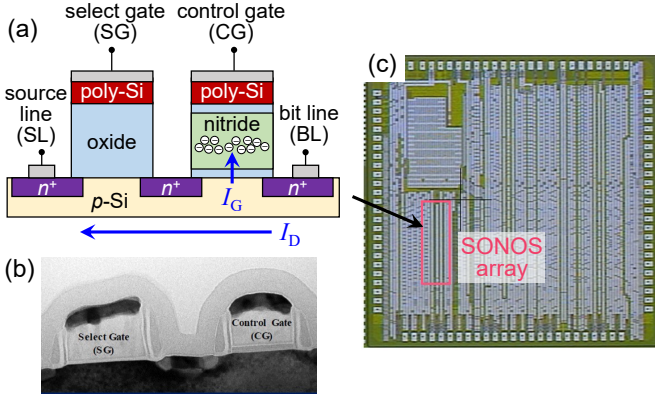


Fig. 2: (a) 40nm analog memory cell including a select transistor (left) and a SONOS transistor (right). (b) TEM cross section of the cell. (c) Die photo of the fabricated test chip with a 1024×1024 SONOS array.

A. Electrical Properties

The electrical properties of the SONOS device can be considered in analogy with a memristor [35]. Where in a metal-oxide or conductive-bridge memristor the state variable is typically the arrangement of defects or ions in the conducting channel [36], [37], the state variable in a SONOS device is the amount of trapped charge in the nitride layer Q , which lies outside the channel but electrostatically influences the channel's conductance G . This can be captured by the effect of the charge on the channel's threshold voltage:

$$I_D = G[V_T(Q)]V_{BL} \quad (1)$$

The dynamical equation for the state variable Q is:

$$\frac{dQ}{dt} = I_G(Q, V_{CG}, V_{SG}, V_{BL}, T) \quad (2)$$

where T is the temperature, and V_{BL} , V_{SG} , and V_{CG} are the voltages on the BL, SG, and CG terminals, respectively, relative to the SL terminal. The rate of change of the state variable Q is controlled by the net tunneling current I_G into the gate stack that fills the nitride traps. I_G is a nonlinear function of both the state Q and the terminal voltages, all of which control the electronic band alignment between the channel and the storage layer. During device programming, large terminal voltages are applied to induce a significant Fowler-Nordheim tunneling current and a relatively large change in state $|\Delta Q|$ over a timescale of microseconds to milliseconds [14]. When the device is under a smaller voltage bias, as during an MVM or in standby mode, dQ/dt is smaller by orders of magnitude. However, the very small tunneling current can still cause charge leakage over a much longer timescale (days, weeks, or more).

Though its behavior can be expressed in terms of the memristor equations, the three-terminal SONOS transistor differs from a two-terminal memristor in that the electrical readout occurs through a separate conduction path (across the Si channel) from the change in the state variable (through the SONOS gate stack). This allows the two processes, described by Equations (1) and (2), to be more readily decoupled. The charge Q affects the channel conductance electrostatically

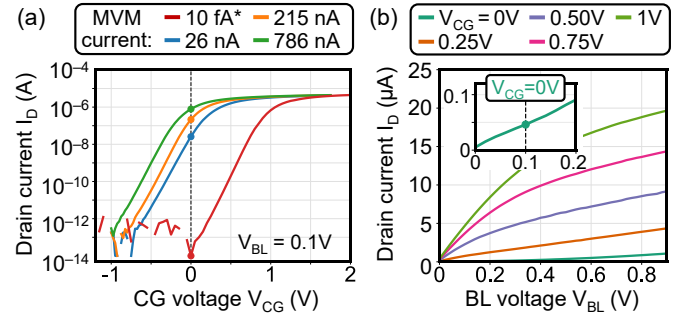


Fig. 3: (a) I_D - V_{CG} properties of the SONOS cell for four programmed states. The weight value is proportional to the value of I_D at the MVM operating point: $V_{CG} = 0V$, $V_{BL} = 0.1V$, $V_{SG} = 2.5V$. (b) I_D - V_{BL} properties for one programmed state at several V_{CG} . Inset zooms into the MVM operating point. *The current is below the measurement resolution.

through V_T , but the tunneling currents during the write process modify Q directly, with no explicit dependence on V_T .

Fig. 3(a) shows the measured cell current as a function of the gate voltage V_{CG} on the SONOS device, with the select transistor biased on. The I_D - V_{CG} characteristics are shown for four different levels of programmed charge Q in the SONOS device, where the red curve is the fully programmed (high V_T) state of the device. Fig. 3(b) shows the I_D - V_{BL} characteristic of the cell. For an MVM, we choose to operate the device in the subthreshold regime with $V_{CG} = 0V$ to reduce current consumption and to minimize the effect of programming errors (see Section III-B). We also use a small BL voltage of $V_{BL} = 0.1V$, where the cell is approximately linear (see inset) and thus G can be considered independent of V_{BL} .

When the SONOS device operates in the subthreshold regime with a fixed gate bias $V_{CG} = 0V$, its channel conductance is given by [38]:

$$G(V_T) = G_0 \exp\left[-\eta \frac{qV_T}{kT}\right] \quad (3)$$

where G_0 is a constant for a given device geometry and terminal bias, $\eta \approx 0.5$ is the gate efficiency, q is the electron charge, and k is the Boltzmann constant. The threshold voltage V_T changes approximately linearly with the stored charge [38]:

$$V_T(Q) = V_{T0} - \frac{Q}{C_N} \quad (4)$$

where V_{T0} is the value of V_T with $Q = 0$, and C_N is the effective capacitance between the control gate and the charge centroid of the nitride storage layer.

B. Programming Precision

A test chip was fabricated that contains a 1024×1024 SONOS array designed for analog MVM operation, shown in Fig. 2(c). A combination of process and write algorithm optimizations enable weight storage with high precision and long retention in the SONOS device. The material profile of the ONO stack was designed to minimize random dopant fluctuations, and a high-quality tunneling oxide reduced the density of interface states. A programming algorithm was developed that both programs the drain current I_D (measured at

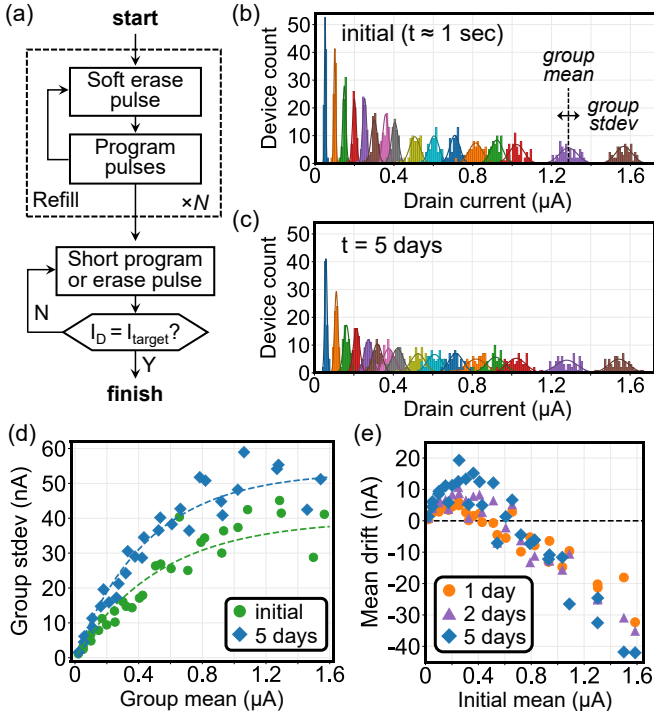


Fig. 4: (a) Algorithm for SONOS programming. (b) Distribution of cell currents just after programming. Each color is a group of 64 cells with the same target current. Each histogram bin is 10 nA wide. (c) Distribution of the same devices after five days. (d) Programming error in the SONOS current vs target current. The data are fit to saturating exponentials. (e) Drift in the distribution means of SONOS currents.

the MVM bias) precisely around a target current and optimizes data retention. Retention depends on the energy spectrum of the populated electronic traps in the nitride: shallow traps near the band edges have less retention than deep traps near the midgap, which provide greater charge confinement [39].

Fig. 4(a) schematically shows the programming procedure, which begins with several refill cycles that are performed for retention: in each cycle, the device is erased, then programmed via short pulses to a verified level. The soft erase pulse is designed to empty electrons from shallow traps, which are de-trapped more rapidly, but not deep traps. The subsequent program pulses fill both types of traps, but over multiple refill cycles a progressively larger portion of electrons are stored in the deep traps [40]. After the refill cycles, a write-verify scheme is used to approach the target current via short program and erase pulses. Programming at 85°C further promoted the ionization of shallow traps but not the deep traps [41]. Cells were programmed one row at a time. After a cell’s target current is reached, it is biased to prevent write disturb.

SONOS cells in the array were programmed to various target levels, and the current distributions are shown in Fig. 4(b) just after programming, and in Fig. 4(c) five days after programming. The currents are measured at the MVM operating bias of $V_{CG} = 0V$ and $V_{BL} = 0.1V$. Each color corresponds to a group of cells programmed to the same target. The histogram for each group can be fit to a normal distribution whose standard deviation is the expected programming error.

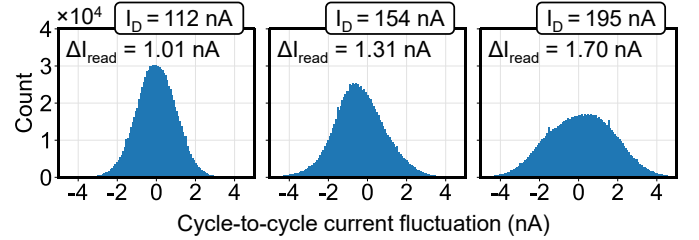


Fig. 5: The current of three SONOS cells is sampled 960,000 times at 2.5 MHz, and the statistics of the fluctuations between samples is shown. ΔI_{read} is one standard deviation.

The relationship between the mean and error currents is shown in Fig. 4(d). An important feature of the error characteristics is that at low currents, the error increases proportionally with current. This is a consequence of the subthreshold operation of the SONOS device, expressed in Equation (3). Inserting Equation (4) into Equation (3), then differentiating with respect to the state variable, we obtain:

$$dG = \frac{q\eta}{kTC_N} dQ \times G \quad (5)$$

Equation (5) shows that for a given amount of programming error dQ on the stored electronic charge, the error in the conductance dG is proportional to G . This is an intuitive result: since the conductance is exponential with Q , a large error dQ at low conductance results in a small *absolute* error in the conductance. Fig. 4(d) also shows that at drain currents above about 1 μA , and thus lower V_T , the error increases sublinearly because the device is transitioning out of the subthreshold regime. The proportionality between dG and G at low conductance has significant implications for neural network error tolerance. This will be discussed in Section IV.

C. Charge Leakage and Read Noise

The weights that are stored in the SONOS devices change over time due to leakage of charge from the nitride layer, caused by thermal de-trapping and trap-assisted tunneling through the oxide into the silicon substrate [42]. Fig. 4(e) shows that leakage causes the means of the distributions in Fig. 4(d) to drift. We find that compared to drift, the inference accuracy is more strongly affected by the increase in the widths of the distributions over time, shown in Fig. 4(c) and (d). The increase in error is due to the random, thermal nature of the leakage process. The variable rates of drift from device to device widen the current distributions over time.

Noise in the cell current also perturbs the weight value after programming. This noise is characterized by rapidly sampling the currents of individual devices at 2.5 MHz. The statistics of cycle-to-cycle fluctuations in the current of cells programmed to three different states is shown in Fig. 5. The noise has a $1/f$ power spectral density, and at the measured levels, the noise standard deviation is a small, fixed proportion ($\sim 0.87\%$) of the drain current I_D . Both of these properties suggest that the dominant noise component is flicker noise, which originates from random trapping/de-trapping of channel carriers at the Si/oxide interface or from random mobility fluctuations [43]. Fig. 5 also shows that $\sim 10^6$ read cycles caused no persistent change in the weight stored by the device.

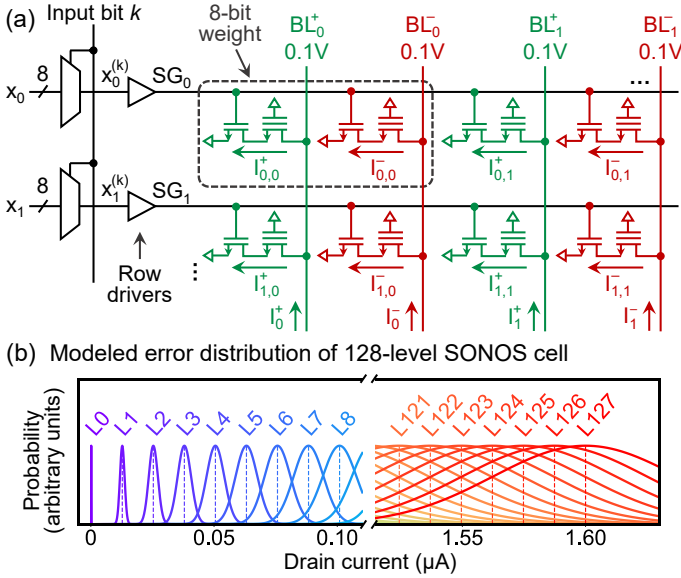


Fig. 6: (a) Analog MVM in a SONOS array, with 8-bit weights and 8-bit inputs. (b) Modeled error distribution of a SONOS cell that is used to approximately represent a 7-bit value.

IV. ERROR-TOLERANT MVM IN A SONOS ARRAY

In this section, we argue that the SONOS device properties are uniquely well suited to maximize the accuracy and energy efficiency of analog neural inference. The accuracy benefit will be quantified in the next section.

The SONOS array used for in-memory MVM is shown in Fig. 6(a). An input value x_i is applied one bit at a time to its corresponding row in the array. A row driver circuit raises the digital high logic level to the 2.5V required to turn on the select transistor. The SLs are held at 0V. Each cell is either connected to the SL when activated (high input bit) or left floating (low input bit). The BLs are held at 0.1V by the current readout circuits (not shown) that provide a virtual ground node. The input bits are applied sequentially and integrated as described in Section VI. Since there is only one value of V_{BL} for an activated cell, the current and conductance are equivalent representations of the stored value.

A signed weight W_{ij} is mapped onto the difference in current of two SONOS cells: $I_{ij}^+ - I_{ij}^-$. Currents from the positive and negative cells are accumulated on separate BLs, and these currents are subtracted by analog peripheral circuitry as described in Section VI. The weight's magnitude is stored in one of the two cells, while the other cell is left in the minimum current state, depending on the sign. To handle ImageNet networks, we choose to map 8-bit weights to each differential pair so that each device encodes a 7-bit magnitude. Note that based on the error properties in Fig. 4(d), the device does not have enough precision to have 128 equally-spaced and well-separated current levels across its entire dynamic range. Thus, we use the SONOS device as an *approximate* 7-bit memory: the programming circuitry targets 128 digital levels, but due to programming errors, the error distributions for nearby levels can overlap. These levels are statistically resolvable over a large number of devices, but there is no guarantee that a single device will be programmed to the correct level.

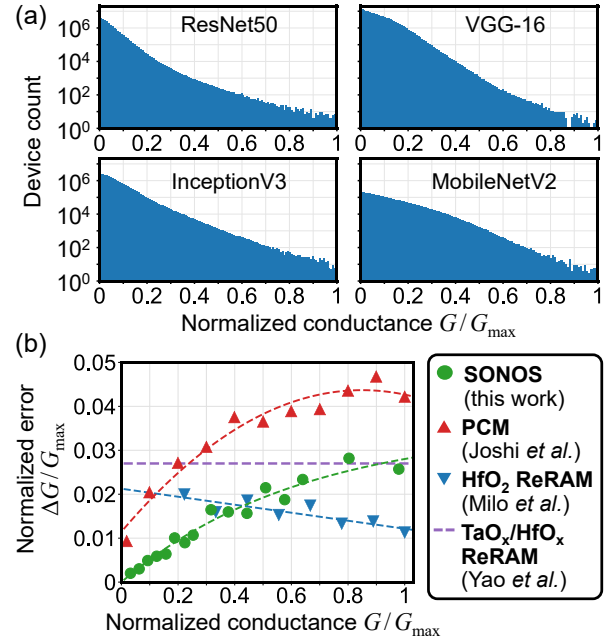


Fig. 7: (a) Distribution of SONOS conductances in systems that implement four popular ImageNet CNNs. (b) Comparison of the error normalized to the maximum utilized conductance of SONOS devices and published PCM [11], HfO₂ memristors [46], and TaO_x/HfO_x memristors [21]. The distributions just after programming are used if available. A parabolic fit is shown for PCM and a linear fit is shown for HfO₂ devices.

Fig. 6(b) shows the modeled error distribution of the 128 levels, equally spaced between I_{\min} and $I_{\max} = 1.6 \mu\text{A}$ ($G_{\max} = 16 \mu\text{S}$). The I_{\min} level is implemented by a state with a much higher threshold than the others, shown by the red curve in Fig. 3(a). The value of I_{\min} is below the measurement resolution, but the On/Off ratio I_{\max}/I_{\min} is at least 10^7 . The state-dependent widths of the distributions in Fig. 6(b) are obtained from a saturating exponential fit to the data in Fig. 4(d) to interpolate between the measured points. At the lowest current levels, the distributions are well-separated, indicating high precision. The highest current levels have significant overlap and thus high uncertainty. The I_{\max} level maps to the weight with the largest absolute value in a given layer.

Fig. 7(a) shows the global distribution of conductances, normalized to G_{\max} , for systems that map four popular ImageNet CNNs. In every case, the conductances are heavily skewed toward zero, which results from two factors. First, most of the network weights have low magnitude; this is a very common feature of neural networks [44], [45]. Second, weights are mapped to conductances in proportion to their magnitude; this is a consequence of differential conductance mapping and the very high On/Off ratio of the SONOS device, which allows zero-valued weights to map to near-zero conductances.

The skewed conductance distributions in Fig. 7(a) combine well with another property of the SONOS cell: the proportionality of the conductance error dG to the conductance G when the conductance is low. This was shown in Fig. 4(d) and is reprised in Fig. 7(b) in normalized units (green). Most importantly, the lowest conductance levels have nearly zero

error. Thus, the most frequently used weights in the neural network are matched with the most precise conductance levels in the cell. Subthreshold operation of the SONOS device is the key: it enables both a high On/Off ratio and state-proportional error, as shown by Equation (5). The vast majority of neural network weights have nearly zero error, providing intrinsic error tolerance for inference.

Fig. 7(b) also contrasts the SONOS device with three recent works with published programming error distributions, which were all used for neural network inference: the PCM device in Joshi *et al.* [11], the HfO₂ memristors in Milo *et al.* [46], and the TaO_x/HfO_x memristors in Yao *et al.* [21]. The programming precision depends both on the device and the programming scheme. The key distinctive property of the SONOS device is that the programming error approaches zero as the weight approaches zero. As mentioned above, this is important since values near zero are the most abundant weights in a neural network. The PCM devices are the most similar among the compared works, but its error in the lowest state remains at $\sim 1\%$ of the utilized conductance range.

The devices in Fig. 7(b) are fairly similar in their maximum conductance: G_{\max} is 16 μS for SONOS, 25 μS for PCM, 20 μS for the TaO_x/HfO_x memristor, and 225 μS for the HfO₂ memristor. However, because it operates in the subthreshold regime, the SONOS device has a much higher On/Off ratio ($>10^7$) compared to the other devices; the On/Off ratio is between 10 and 100 for the two memristors and ~ 100 for PCM based on the error. The On/Off ratio is important because with differential cells and a heavily skewed weight distribution, the majority of cells on a BL will use the minimum or near-minimum conductance level. If the On/Off ratio is insufficient, these small cell currents can accumulate to become a large BL current. For a given array size, larger BL currents require more power-hungry peripheral circuits and also increase the effect of parasitic IR voltage drops in the array, which can degrade accuracy as we show in Section V.

The large On/Off ratio of the SONOS cell ensures that zero-valued weights conduct negligible current, reducing accumulated BL currents. Similarly, the thick gate oxide and large threshold voltage ($>2\text{V}$) of the select transistor ensure that non-activated cells also draw negligible current during an MVM. This allows an increase in the array size (number of rows) while keeping the BL currents relatively small, and without incurring MVM accuracy losses due to BL parasitic resistance. Since the energy consumption of an analog MVM tends to be dominated by peripheral circuits, a larger array directly increases energy efficiency by amortizing these energy costs over more MAC operations.

V. INFERENCE ACCURACY EVALUATION

This section uses the characterized device properties in Section III and array design in Section IV to simulate the accuracy of a SONOS-based analog inference accelerator. To emulate a realistic machine learning application, accuracy is evaluated using the ResNet50-v1.5 network on the ImageNet dataset [10]. We use the reference implementation of this network from the MLPerf Inference Benchmark [15], [16],

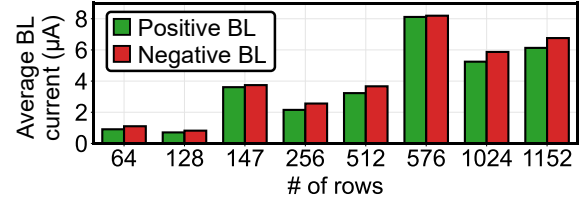


Fig. 8: BL current vs SONOS array size for a system that implements ResNet50, based on inference simulation data on ImageNet. Each value is an average over all arrays of the same size, and over all MVMs in 30 inference operations.

which has also been used to benchmark digital inference accelerators. We refer to this network henceforth as ResNet50. All weights are quantized to 8 bits before being mapped to hardware. Weights are not retrained or modified (besides quantization) before being transferred to the arrays.

A. Accuracy Simulator

The accuracy of a SONOS inference accelerator is predicted using CrossSim [47], an end-to-end neural network simulation tool that models the various analog errors in in-memory MVM. For every SONOS cell, a random programming error ΔI_{prog} is added to the target current. This error is sampled from a normal distribution, whose target-dependent standard deviation is derived from the data in Fig. 4(d). Leakage is modeled using the time-dependent standard deviation and drift (Fig. 4(e)) in the current distributions. Read noise is modeled by sampling a new error ΔI_{read} based on Fig. 5 that adds to ΔI_{prog} on every MVM. Array parasitic resistance and ADC quantization are modeled as described below. Digital operations such as ReLU activations are assumed to be error-free.

Convolutions are unrolled into a sequence of sliding window MVMs, executed on arrays of size $K_x K_y N_{\text{ic}} \times N_{\text{oc}}$ as described by Shafiee *et al.* [5] ($K_x \times K_y$ is the 2D filter size, N_{ic} and N_{oc} are input and output channel dimensions). During inference, batch normalization layers are folded into the matrix of a preceding convolution to reduce the digital processing overhead between layers [34]. Bias weights are stored and added digitally, as they can lie in a different numerical range from the weights that are stored in the array.

B. Parasitic Resistance Effects

As described in Section IV, the accumulated currents on the rows and columns during an MVM constrain the maximum array size. The larger these currents, the larger the parasitic voltage drops across the array interconnects. This effect is negligible on the SG lines of the SONOS array since the transistor gates draw almost no current. However, the effect can be significant on the BLs. Parasitic voltage drops cause the BL terminal of the cells to depart from 0.1V in a nonuniform manner across the array, distorting the weight matrix.

Fig. 8 shows the average simulated BL currents for differently sized SONOS arrays that map the layers of ResNet50. A maximum of 1152 rows is assumed; smaller matrices use fewer rows, while a matrix with more than 1152 rows is partitioned over multiple arrays. Even with 1152 cells per

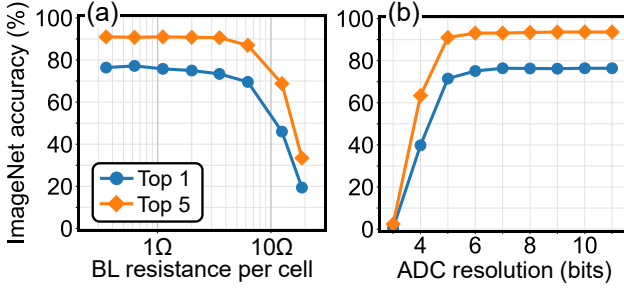


Fig. 9: (a) ImageNet accuracy (500 images) with ResNet50 vs BL parasitic resistance, (b) ImageNet accuracy (1000 images) vs ADC resolution. All results assume 1152 maximum rows. SONOS errors and noise are not included.

column, the average BL current is only several times larger than the maximum cell current of $1.6 \mu\text{A}$. This results from several factors: the skewed distribution of weight values in the neural network, the high On/Off ratio of the SONOS device and the select transistor, and the fact that the most significant bits of the input vector are sparse. All of these factors help reduce the impact of parasitic resistance on accuracy [48].

For computational tractability, parasitic resistance effects are modeled by treating the SONOS device as a linear resistor around $V_{\text{BL}} = 0.1\text{V}$, which is a good approximation based on Fig. 3(b). A circuit simulation is performed for each input bit; eight per MVM. Fig. 9(a) shows the accuracy on a subset of ImageNet versus the BL resistance between two adjacent cells. The accuracy begins to fall for a parasitic resistance larger than 1Ω , which is $\sim 60,000$ times smaller than the minimum resistance of a SONOS cell. A unit cell resistance smaller than 1Ω can readily be achieved by metal interconnects in our 40 nm process, indicating that an array with 1152 rows is not limited by parasitic resistance effects.

C. ADC Resolution

Like the array size, the ADC resolution is a critical system parameter that affects both accuracy and energy efficiency. To obtain high accuracy at the minimum feasible resolution, the ADC’s input range must be calibrated to minimize quantization and clipping errors [34]. As described in Section VI, the final analog output after the integration of all input bits is a voltage V_{out} , which is then digitized. Values outside the ADC’s range are clipped to the lowest or highest ADC level.

To calibrate the range of each layer’s ADC, we first collect statistics (via simulation) on each layer’s output voltages using 500 ImageNet images. For every layer, the range (V_{min} , V_{max}) is then found that minimizes the error $\epsilon = \|\vec{V}_{\text{out}} - \vec{V}_{\text{Q}}\|_1$ over the collected output voltages \vec{V}_{out} . \vec{V}_{Q} is obtained by clipping and quantizing \vec{V}_{out} to N bits in this range, where $N = 12$ was found to be optimal. If a matrix is partitioned across multiple arrays, they share the same ADC limits.

Fig. 9(b) shows the accuracy on a subset of ImageNet versus ADC resolution, after every layer’s ADC range is calibrated. Based on this result, we use an ADC resolution of 8 bits across the full accelerator, which incurs very little loss in accuracy. Afterwards, the digital dot products from different partitions are added (if necessary), the biases are added, and activation

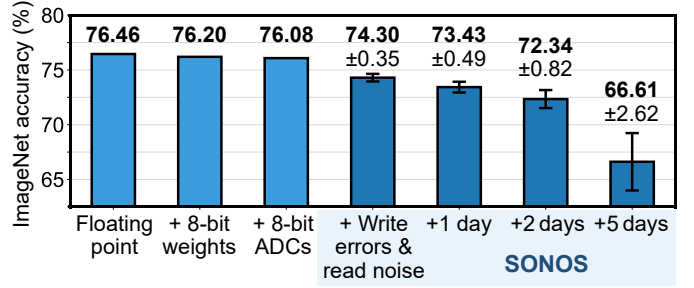


Fig. 10: ImageNet top-1 accuracy (50,000 images) using ResNet50. SONOS arrays have at most 1152 rows.

TABLE I: Inference accuracy comparison of technologies

	PCM [11]	SONOS
Process node	90 nm	40 nm
Expected weight error in ResNet50 (% of G_{max})	1.42% (0.35 μS / 25 μS)	0.18% (0.028 μS / 16 μS)
Evaluated network	ResNet34 (ImageNet)	ResNet50 (ImageNet)
Floating-point accuracy	73.20%	76.46%
In-accelerator accuracy (no re-training)	65.41% \pm 3.6% (-7.79%)	74.30% \pm 0.35% (-2.16%)
In-accelerator accuracy (with re-training)	71.62% \pm 0.4% (-1.58%)	-

and/or pooling functions are applied. The results are quantized again to 8 bits to produce the inputs to the next layer. The ranges of these 8-bit inputs are optimized using the same method as above. Section VI-B and VII will describe how the calibrated ADC and activation ranges are encoded in hardware.

D. Accuracy of the SONOS Accelerator

Fig. 10 shows the accuracy of ResNet50 on the full ImageNet test set when deployed on the SONOS accelerator. The various circuit- and device-related non-idealities are added cumulatively. Cell programming errors are re-sampled ten times to obtain the variance in accuracy. The programming errors (Fig. 4) have a much larger effect on accuracy than the read noise (Fig. 5). For computational tractability, parasitic resistance is not included; based on Fig. 9(a), its effects on accuracy are negligible using our process.

Just after programming the weights, the SONOS accelerator attains an average top-1 ImageNet accuracy of 74.30%, which is a 2.16% loss relative to floating-point. Table I shows that this is a significant improvement over a recent simulation result based on characterized PCM devices, the only other device that has been benchmarked on ImageNet. The very low error of the SONOS cell at low currents leads to a much lower expected error for a given neural network, and this is responsible for its higher inference accuracy. The cell errors are low enough on average that their effect remains small even when they accumulate over as many as 1152 rows. The higher precision and intrinsic error tolerance of the SONOS device reduce the need for specialized network re-training prior to deployment in the proposed analog hardware, enabling greater compatibility with state-of-the-art training workflows.

The accelerator gradually loses accuracy due to charge leakage from the SONOS cells, as shown in Fig. 10. To restore high accuracy, the weights need to be refreshed every one to five days from a reference, depending on application needs.

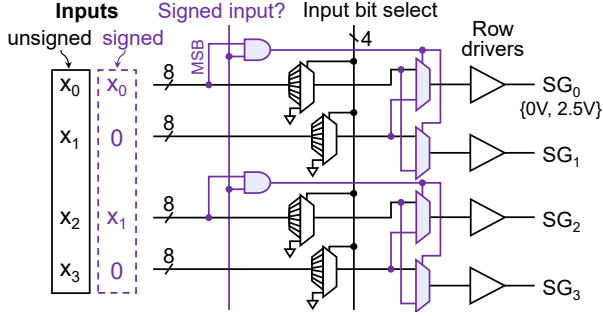


Fig. 11: Row peripheral circuitry. For unsigned inputs, the elements in violet can be ignored. For signed inputs, a control unit ensures that they are applied to every other input port.

Device robustness to 100,000 write cycles is sufficient for ten-year operation at this refresh rate. The high-voltage pulses used for SONOS programming and erase pulses can introduce defects in the tunnel oxide over time, which has been observed to shift the fully programmed/erased V_T states and gradually reduce data retention [49]. The effect of 100,000 write cycles has been shown to be minimal on 40nm digital SONOS memory technology [14], but analog states are more sensitive, and it is possible that gradually more frequent refreshes are needed over the system’s lifetime to maintain high accuracy.

VI. AN ENERGY-EFFICIENT SONOS MVM CORE

This section describes the design of the analog MVM core, which consists of the SONOS array and its peripheral circuitry. The core takes 8-bit digital inputs, executes an analog MVM with its stored 8-bit weights, and produces 8-bit digital outputs. Circuit components are designed and simulated in the 40nm foundry process used for the embedded SONOS arrays.

A. Row Peripheral Circuits

Fig. 11 shows the circuit that supplies the 8-bit inputs of an MVM. One bit at a time of the input is selected, and a row driver converts the digital high logic level to the 2.5V required to turn on the select transistor. To handle both signed inputs and signed weights, the array in Fig. 6 must use four cells per weight, rather than two. In this case, two cells implement the weight W , and the other two cells on a second row implement its inverse $-W$. The signed input logic, shown in violet in Fig. 11, drives one of the two rows based on the input sign bit. In CNNs, the use of ReLU activations means that typically only the first layer has signed inputs, and the other layers can be configured to use two cells per weight to save area.

B. Bit Line Integrators

Each pair of BLs in Fig. 6 is connected to an integrator, shown in Fig. 12, that subtracts the positive and negative BL currents, integrates the difference on a capacitor, and accumulates the dot products for different input bits. The positive BL and negative BL are each connected to a current conveyor (CCII) with a current gain of -1 and $+1$, respectively, from port X to Z. The two currents are subtracted at the output node to implement the differential conductance representation

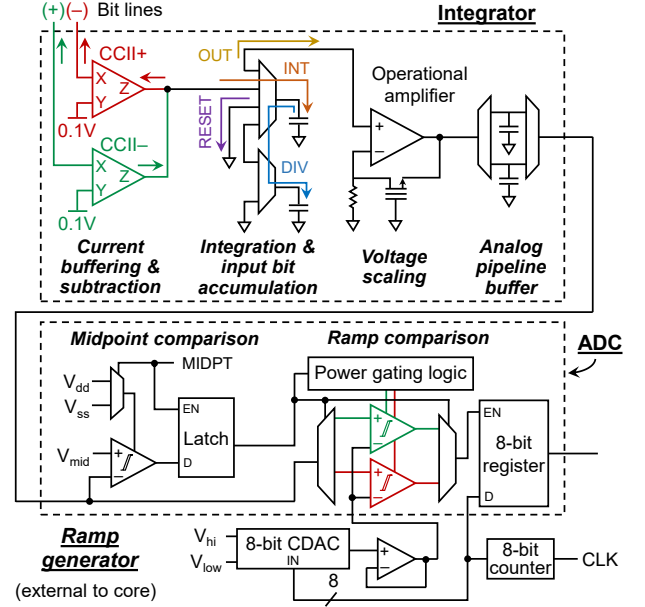


Fig. 12: Peripheral circuitry connected to each pair of bit lines. The lower three blocks implement the ADC.

of weights. The CCII also acts as voltage follower between X and Y. This is used to maintain a virtual ground on each BL, ensuring a proper voltage bias on the SONOS cells. A CCII design with low input resistance is used to keep the BL voltage as close as possible to 0.1V [50], [51].

The difference current is integrated on a 0.5 pF capacitor, producing a voltage (INT phase in Fig. 12). The current is integrated for 10 ns per input bit to overcome noise and circuit transients. Analog results from different input bits are accumulated using the successive integration and rescaling (SIR) technique [52]. After integration, the capacitor is disconnected from the BLs and is discharged to a second identical capacitor to halve its voltage (DIV). The capacitor is then reconnected to the BLs and the current from the next input bit is integrated on the accumulated charge. After all bits have been integrated, the capacitor is connected to the next stage (OUT). SIR effectively performs the shift-and-add aggregation of input bits in the analog domain, and reduces the ADC overhead by $8\times$ as it requires only one ADC conversion for all eight bits.

Before being passed to the ADC, the integrated voltage is scaled by a tunable-gain amplifier to properly place the MVM result within the range of the ADC. This step is necessary to obtain high accuracy, as explained in Section V-C. The tunable gain is implemented using an operational amplifier in a non-inverting configuration, with a programmable SONOS device acting as a feedback resistor. The required gain can be found offline and does not change during inference.

The final voltage output is stored on a capacitor. Analog capacitive double buffering is used to concurrently process two MVMs within the core: one in the analog MVM and integration step, and one in the ADC step. While the ADC reads one capacitor to digitize one MVM, the analog result of the next MVM is written to the other capacitor.

MVM accuracy can be degraded by variation and offset in the integrator components. In particular, errors in the inte-

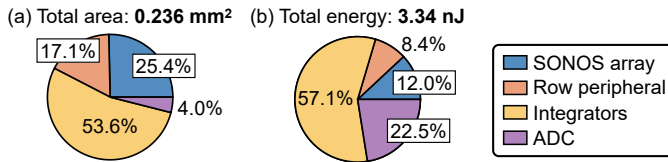


Fig. 13: Breakdown of (a) area and (b) energy consumption for a SONOS core implementing an 1152×256 matrix.

grating capacitors and operational amplifiers can cause MVM outputs on different BLs to be scaled by slightly different factors, distorting the results. These effects can be partially corrected. The programmable feedback resistor can be separately calibrated on each BL pair to adjust the amplifier gain to compensate for capacitance variations. Offsets in the current conveyors, amplifiers, or comparators can be compensated with an additional row of calibrated SONOS devices that inject fixed currents onto the BLs [50].

C. Analog-to-Digital Converters

To achieve high speed, low area, and low energy, the core uses an 8-bit parallel ramp ADC, shown in Fig. 12 after the integrator. The ADC uses a 256-level staircase ramp between two reference voltages, which is generated by passing an 8-bit, 1 GHz digital counter into a capacitive digital-to-analog-converter (CDAC). The ramp generator is shared by multiple cores. Every integrator’s output voltage is compared to the shared ramp. On the clock cycle when the two signals are equal, the comparator switches and the 8-bit counter value is latched in a register. The comparator, whose design is based on Marinella *et al.* [50], uses internal positive feedback to obtain a 1 ns response time. By the end of the 256 ns ramp, the 8-bit digital dot products are available for the full MVM.

To maximize the ADC’s voltage resolution, the ramp starts from the negative voltage rail and/or ends at the positive rail. Two comparators are combined that together provide rail-to-rail input dynamic range: one with an NMOS input stage (green) and one with a PMOS input stage (red). To reduce power consumption, one of the two comparators is gated off during the ramp depending on the input value. This is decided by a third comparator (black), which compares the input to the ramp midpoint and stores the result in a latch. This third comparator is gated on for 2 ns (MIDPT) prior to the ramp.

D. Core Energy Efficiency and Area

Fig. 13 shows the breakdown of the core’s energy and area among its components for an 1152×512 SONOS array, which supports up to an 1152×256 MVM (576×256 if inputs are signed). These dimensions were chosen to maximize the number of rows while also balancing utilization by small and large matrices in a typical CNN.

Energy consumption is computed based on SPICE simulations and metal interconnect properties of the 40nm process. The energy of the array and row drivers is adjusted for the average SONOS currents and the average values of the input bits when running ResNet50 on ImageNet, obtained by simulation. The ramp generator uses the CDAC and reference voltage buffer designs in Kull *et al.* [53] and the digital counter

design in Baert *et al.* [54], with properties scaled to the technology node, supply voltage, and 8-bit DAC resolution. The core area is a sum of circuit block areas rather than a physical layout; other than the SONOS array, 25% of the area is assumed to be wiring, which is typical for this process.

The integrators make up the largest share of both the area and energy of the core, since SIR greatly reduces the energy cost of the ADCs. The SONOS core achieves a peak energy efficiency of 177 TOPS/W (11.3 fJ/MAC) for a 1152×256 MVM, which fully utilizes the array and peripheral circuits. Since the energy is dominated by the peripherals, the efficiency is reduced for smaller matrices that use fewer rows, where the BL circuit costs are amortized over fewer operations.

VII. SONOS INFERENCE ARCHITECTURE

The SONOS inference architecture is hierarchically organized to support large-scale CNN processing. Several MVM cores implementing the same layer are combined with activation memory and digital processing elements to form a *tile*. To implement the full neural network, tiles are interconnected by a concentrated mesh network similar to prior work [5], [29]. Weight data is kept stationary inside the non-volatile SONOS arrays during inference, while activation data is moved within and between tiles. Computation is pipelined across the system to allow different tiles to work concurrently on different computations within a layer, different layers within the same inference operation, or on separate inference operations.

The accelerator uses a uniform and modular tile design, shown in Fig. 14(a), which flexibly implements different layer types. MVM layers, including convolutions and fully-connected layers, are executed mostly inside the analog cores. Non-MVM layers such as pooling and element-wise addition are processed using the tile’s digital arithmetic logic unit (ALU). MVM layers and non-MVM layers share the same tile architecture but follow different datapaths within the tile. Computation is pipelined through the tile’s hardware components to maximize throughput. The colored hardware blocks in Fig. 14(a) correspond to the colored pipeline stages in Fig. 14(b) and (c) for MVM and non-MVM layers, respectively. We refer to the length of one pipeline stage as a machine cycle, which is set to 295 clock cycles at 1 GHz.

To enable concurrent execution of different pipeline stages within a tile, double buffering is used in the pipeline registers (yellow). In a given machine cycle, one buffer is populated by outputs from the preceding stage while the other holds inputs that are read by the next stage. Fig. 14(d) shows the tile area breakdown, where the control unit is not designed but conservatively estimated to take up $\sim 4\%$ of the area.

The dataflow pipeline consists of the following stages:

1) Data in: 8-bit inputs are received through the interconnection network and placed in 32 parallel FIFO queues. Each FIFO is connected to the write port of one of 32 banks in the local SRAM memory.

2) Memory write/read: In the first half of the machine cycle, activations are written from the FIFOs to their corresponding SRAM banks. At most 4608 values (4.5 kB) can be written to SRAM in this stage, equal to the maximum number

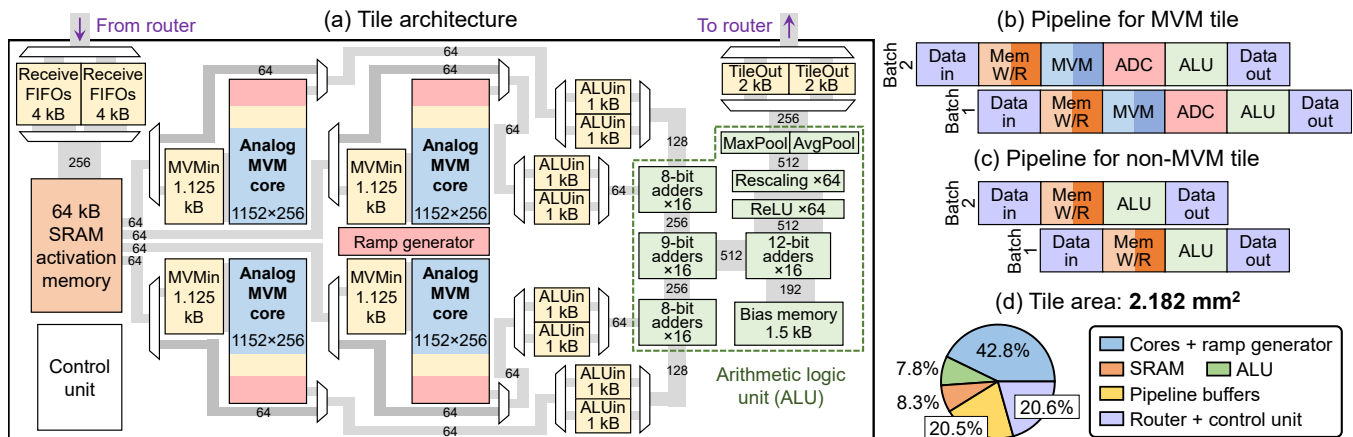


Fig. 14: (a) Tile architecture, (b-c) dataflow pipeline for MVM/non-MVM tiles, (d) Breakdown of tile area.

of inputs consumed by the four cores in one machine cycle. In the second half-cycle, activations are read from SRAM and loaded into the MVMMin buffers of each core. If the same activation is needed at multiple row positions, it is copied during the Data In stage to multiple banks. In a non-MVM tile, activations are read from memory to the ALUin buffers, bypassing the MVM and ADC stages.

3) MVM: In the first half-cycle, each core reads the inputs from the MVMMin buffers and executes an 8-bit analog MVM (all input bits integrated). In the second half-cycle, the MVMMin buffers are disconnected from the cores and new values are loaded from SRAM. During this time, the core performs voltage scaling of the analog dot products and stores the results in its internal analog pipeline buffer.

4) ADC: The analog MVM results are converted to 8-bit digital values, which requires 258 clock cycles. In the remainder of the machine cycle, the values are transferred from each core's 256 output registers to the ALUin buffers.

5) ALU: Non-MVM computations are performed in this stage using the ALU on operands fetched from the ALUin buffers. The ALU receives up to 64 total inputs simultaneously from four ALUin buffers and has a worst-case delay that allows a new set of inputs to be processed every four clock cycles. ALU components can be bypassed if not used. The results are stored in the TileOut buffer.

In MVM tiles, the ALU sums the results from multiple cores (if needed), adds the bias values, and applies the ReLU function. The same adders can be used by element-wise addition layers. The MaxPool and AvgPool units are implemented with trees of adders and digital magnitude comparators. Where possible, pooling is done in the ALU stage of an MVM tile rather than in a dedicated pooling tile.

Though the ALU takes 8-bit operands, intermediate results within the ALU can have up to 13 bits of resolution. These must be reduced to 8-bit activations that can be used by the next layer, with calibrated numerical ranges as described in Section V-C. Scaling to these ranges is accomplished via integer multiplication by a layer-specific scaling factor, followed by a bit shift to implement a power-of-two division. The relevant eight bits are then extracted from the result.

6) Data Out: 8-bit outputs in the TileOut buffer are routed to the destination tile. This stage coincides with the Data In

stage of the destination tile.

As described in Section V-A, a convolution is executed as a set of parallelizable sliding window (SW) MVMs. In a CNN, the early layers have small matrices that can be replicated several times within an array, allowing the parallel execution of multiple SWs in a core. Later layers have larger matrices that may need to be partitioned over multiple cores. To increase throughput at the expense of area, weights can be physically replicated across multiple cores and/or tiles that process different SWs in parallel [5].

VIII. PERFORMANCE, ENERGY, AND AREA EVALUATION

This section evaluates the performance, energy efficiency, and area of the proposed SONOS analog inference accelerator.

A. Architecture Simulator

To estimate the above metrics, a high-level simulator for the proposed accelerator was developed that models the inference pipeline in Fig. 14. The simulator uses a heuristic-based approach to assign tiles to neural network layers. The number of tiles allocated to a layer depends on the matrix shape and a target number of cycles to complete all SWs, which defines the replication factor for the weights. SWs are allocated among tiles in a manner that minimizes delay between layers and avoids overflow of the 64 kB SRAM tile memory.

After mapping a neural network, an inference dataflow simulation is conducted at the granularity of a machine cycle. In each cycle, every tile's SRAM is checked to determine whether the inputs are available for the cores to compute their next set of allocated SWs; if so, these are fetched to the MVMMin buffers. Inputs are removed from SRAM when they are not needed for any future computation.

The energy and area of the core components are estimated as described in Section VI-D. Peripheral circuits on rows or columns not utilized by a layer's matrix are gated off to save energy. The energy and delay of digital logic in the ALU are based on a standard cell library for the process. We use CACTI 7.0 [55] with a custom technology file for our 40nm process to obtain the area, access energy, access latency, and leakage of the SRAM memory. Network routers are modeled using Orion 3.0 [56]. For the concentrated mesh configuration, 8-port routers are used with 128-bit flits and two virtual channels.

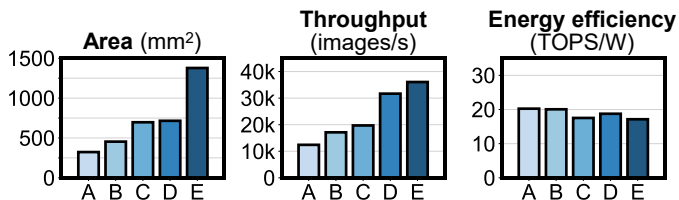


Fig. 15: Area, ImageNet inference throughput, and energy efficiency of five different system designs for ResNet50. Designs A-E use 148, 208, 296, 312, and 556 tiles, respectively.

B. System Scalability

Fig. 15 shows five accelerator designs for the same network (ResNet50) with different total tile counts. Adding more tiles generally increases the parallelism of SW computation, trading off inference throughput with area. Meanwhile, scaling the size of the system has a relatively small effect on the energy efficiency. For ResNet50, the efficiencies of all five designs are close to 20 TOPS/W, which is much lower than the peak core efficiency of 177 TOPS/W. This degradation results from energy contributions outside the core, and from the fact that not all layers efficiently utilize the SONOS arrays.

Fig. 16(a) shows the energy breakdown among the accelerator components for Design B in Fig. 15. For ResNet50, about 32% of the energy is consumed in the core and much of the remainder is spent on data movement, particularly in reads and writes to the tile SRAM buffer. This is in contrast to some earlier accelerators at the same ADC resolution [5], [29], where the ADC was the dominant energy consumer. The fact that the system’s energy is bottlenecked by the intra- and inter-tile communication of activations is a consequence of the high efficiency of the cores, despite their sub-optimal utilization by some layers. Significant further efficiency improvements would require a reduction in data movement energy.

C. Comparison of Neural Network Workloads

Fig. 16 compares the energy per inference of three ImageNet CNNs. There is a significant disparity in the total energy consumption of ResNet50 and ResNet34, which are nominally similar networks that differ only by 5% in the total number of operations. The key difference is that ResNet34 lacks the 1×1 convolutions that are common in ResNet50. The 1×1 convolution is a poor fit for in-memory MVM since it typically has too few inputs to efficiently amortize the energy cost of the BL peripheral circuits. This leads to a higher energy cost per MAC in these layers. Thus, the core operates closer to peak efficiency for ResNet34 and VGG-16.

Fig. 17(b) compares the accelerator’s performance and efficiency on the three networks above and the CNN in Fig. 17(a), whose topology was engineered to maximize the system’s energy efficiency. This network has 15.2M untrained weights and a structure that is similar to that of a typical CNN.

ResNet50 and ResNet34 have the longest end-to-end inference latency due to their depth, but ResNet34 has high throughput because multiple images can easily be pipelined through its simpler residual blocks. ResNet50 utilizes the available SONOS arrays poorly, largely due to its 1×1 convolutions. ResNet34 obtains superior utilization by eliminating

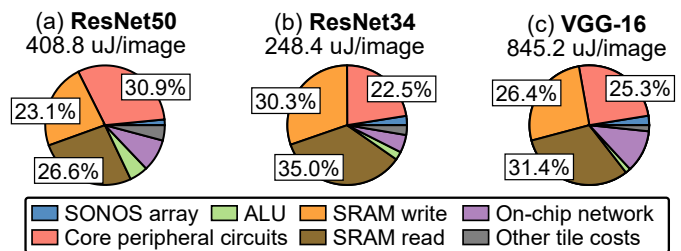


Fig. 16: Breakdown of energy consumption across hardware components for three ImageNet neural networks.

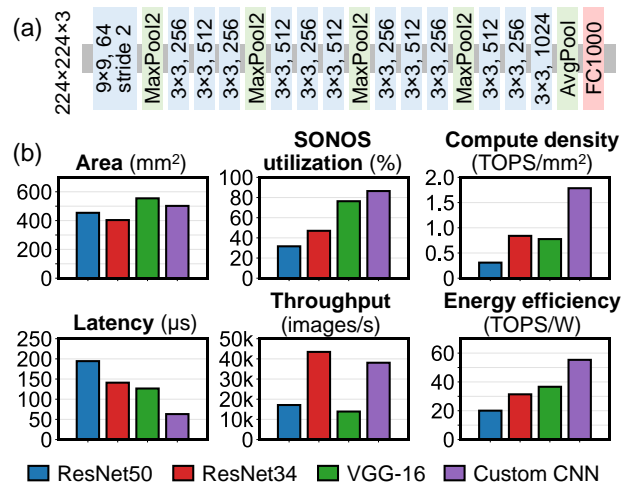


Fig. 17: (a) Custom CNN designed to maximize energy efficiency, (b) Comparison of accelerator performance on different ImageNet CNNs at approximately the same total area.

these convolutions, and VGG-16 improves upon this further by having very large matrices and wasting fewer arrays in non-MVM tiles (especially element-wise addition). However, the compute density of VGG-16 suffers from having very large fully-connected layers at the output of the network that consume the majority of the area but perform a minority of the operations. These tiles also consume significant static power while idle, reducing energy efficiency.

The observations above informed the design of the custom CNN. The network avoids large, infrequently used fully-connected layers; rapidly inflates the channel dimensions to more fully utilize every core; and rapidly compresses the spatial dimensions to reduce the number of SWs and improve throughput. The co-designed neural network and its hardware implementation reaches an efficiency over 55 TOPS/W, which represents close to the maximally efficient accelerator use case. This example shows that by tailoring the topology to the hardware, the neural network can more fully capitalize on the large energy and density benefits of the accelerator.

D. Comparison with other accelerators

Table II compares the proposed accelerator to several deployed digital inference accelerators and ISAAC, a memristor-based analog inference architecture [5]. The performance and efficiency of analog accelerators are simulated, while the digital accelerators are fully functional chips with measured performance.

TABLE II: Energy and performance comparison with other inference accelerators

	Google TPUv1 [3]	Nvidia A100 [1]	Alibaba HanGuang 800 [57]	Qualcomm Cloud AI 100 [58]	ISAAC [5]	This work
Process	28 nm	7 nm	12 nm	7 nm	32 nm	40 nm
Accelerator type	Digital	Digital	Digital	Digital	Analog, memristor	Analog, SONOS
Data resolution	8 bits	8 bits	8 bits	8 bits	16 bits	8 bits
Compute density (TOPS/mm²)	0.06 (average) 0.28 (peak)	0.36 (ResNet50)	0.80 (ResNet50)	–	0.48* (peak)	0.31* (ResNet50) 1.78* (custom)
Energy efficiency (TOPS/W)	2.3 (peak)	1.20 (ResNet50)	2.03 (ResNet50)	2.34 (ResNet50)	0.63* (peak)	20.1* (ResNet50) 55.3* (custom)

*The performance of ISAAC and this work are simulated, while others are measured hardware performance. ResNet50 performance values are based on MLPerf Inference v1.1 for the Nvidia A100 and Qualcomm Cloud AI 100, and v0.5 for the Hanguang 800 [16]. In TPUv1, ‘average’ refers to the average performance across six datacenter inference workloads [3]. For the digital accelerators, compute density and efficiency are based on vendor-published die area and thermal design power, if available. For the SONOS accelerator, ‘custom’ refers to the CNN in Fig. 17(a) which achieves near-peak efficiency.

Despite leveraging the efficiency of analog in-memory MVM, ISAAC provides very limited system-level efficiency gains over state-of-the-art digital accelerators. Our SONOS accelerator improves upon ISAAC and similar analog accelerators [4], [6], [29], [30] in two important ways, both enabled by the SONOS technology. First, the low SONOS programming error removes the need for bit slicing. ISAAC incurs a large overhead by encoding only two bits per cell, and requiring separate peripheral circuits for every 2-bit slice. As shown in Section V, the SONOS cells can obtain high ImageNet accuracy when operated as approximate memories. Second, the low errors and high On/Off ratio of the SONOS device enable comparatively large arrays to be used without degrading accuracy. Our SONOS arrays have $9\times$ more rows than the memristor arrays in ISAAC. Since analog operations inside the array are much cheaper in energy than peripheral circuit operations, large arrays directly increase energy efficiency.

IX. CONCLUSION

The intrinsic physical properties of SONOS charge trap memory are uniquely well matched to the statistical properties of modern neural networks, whose weights are heavily skewed toward zero. By operating the device in the subthreshold regime, where the current-voltage characteristics are exponential, two key properties can be exploited: a high conductance On/Off ratio ($>10^7$) and conductance programming errors that scale with conductance. Specifically, the error approaches zero in the limit of zero conductance, such that the most abundant weight values are represented with the highest precision. These properties are crucial in suppressing the accumulation of errors in a large array, whether the errors arise from programming or from parasitic voltage drops. The simulated accuracy of the SONOS accelerator (2.16% loss on ImageNet using ResNet50 without retraining) is the highest obtained so far in the literature, among works that are based on experimentally characterized devices. The precision and On/Off ratio of the device can be further leveraged to build a more energy-efficient analog accelerator that uses large arrays and avoids bit slicing. Without compromising accuracy or robustness, the proposed SONOS-based inference accelerator attains >20 TOPS/W efficiency for ResNet50 and a peak efficiency >55 TOPS/W with an optimally co-designed neural network.

ACKNOWLEDGMENTS

This work was supported by the Laboratory Directed Research and Development program at Sandia National Labora-

tories, and was also funded in part by DTRA under grant no. HDTRA1-17-1-0038. This paper describes objective technical results and analysis. Any subjective views or opinions that might be expressed in the paper do not necessarily represent the views of the U.S. Department of Energy or the United States Government. Sandia National Laboratories is a multimission laboratory managed and operated by National Technology & Engineering Solutions of Sandia, LLC, a wholly owned subsidiary of Honeywell International Inc., for the U.S. Department of Energy’s National Nuclear Security Administration under contract DE-NA0003525.

REFERENCES

- [1] (2020, Nov.) NVIDIA data center deep learning product performance. NVIDIA. [Online]. Available: <https://developer.nvidia.com/deep-learning-performance-training-inference>
- [2] B. Gaide, D. Gaitonde, C. Ravishankar, and T. Bauer, “Xilinx Adaptive Compute Acceleration Platform: Versal architecture,” in *Intl. Symp. on Field-Programmable Gate Arrays*, 2019, pp. 84–93.
- [3] N. P. Jouppi *et al.*, “In-datacenter performance analysis of a Tensor Processing Unit,” in *Intl. Symp. on Computer Architecture (ISCA)*, June 2017, p. 1–12.
- [4] M. N. Bojnordi and E. Ipek, “Memristive Boltzmann machine: A hardware accelerator for combinatorial optimization and deep learning,” in *Intl. Symp. on High Performance Computer Architecture (HPCA)*, March 2016.
- [5] A. Shafiee, A. Nag, N. Muralimanohar, R. Balasubramanian, J. P. Strachan, M. Hu, R. S. Williams, and V. Srikumar, “ISAAC: A convolutional neural network accelerator with in-situ analog arithmetic in crossbars,” in *Intl. Symp. on Computer Architecture (ISCA)*, June 2016.
- [6] P. Chi, S. Li, S. Li, T. Zhang, J. Zhao, Y. Liu, Y. Wang, and Y. Xie, “PRIME: A novel processing-in-memory architecture for neural network computation in ReRAM-based main memory,” in *Intl. Symp. on Computer Architecture (ISCA)*, June 2016.
- [7] A. Sebastian, M. Le Gallo, R. Khaddam-Aljameh, and E. Eleftheriou, “Memory devices and applications for in-memory computing,” *Nature Nanotechnology*, vol. 15, no. 7, pp. 529–544, 2020.
- [8] T. P. Xiao, C. H. Bennett, B. Feinberg, S. Agarwal, and M. J. Marinella, “Analog architectures for neural network acceleration based on non-volatile memory,” *Applied Physics Reviews*, vol. 7, no. 3, p. 031301, 2020.
- [9] S. Yu, H. Jiang, S. Huang, X. Peng, and A. Lu, “Compute-in-memory chips for deep learning: Recent trends and prospects,” *Circuits and Systems Magazine*, vol. 21, no. 3, pp. 31–56, 2021.
- [10] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei, “ImageNet: a large-scale hierarchical image database,” in *Conf. on Computer Vision and Pattern Recognition (CVPR)*, June 2009.
- [11] V. Joshi, M. L. Gallo, S. Haefeli, I. Boybat, S. Nandakumar, C. Piveteau, M. Dazzi, B. Rajendran, A. Sebastian, and E. Eleftheriou, “Accurate deep neural network inference using computational phase-change memory,” *Nature Communications*, vol. 11, 2020.
- [12] Z. He, J. Lin, R. Ewetz, J.-S. Yuan, and D. Fan, “Noise injection adaption: End-to-end ReRAM crossbar non-ideal effect adaption for neural network mapping,” in *Design Automation Conf. (DAC)*, June 2019, pp. 57:1–57:6.

- [13] S. Jain and A. Raghunathan, "CxENN: Hardware-software compensation methods for deep neural networks on resistive crossbar systems," *ACM Trans. Embed. Comput. Syst.*, vol. 18, no. 6, Nov. 2019.
- [14] I. Kouznetsov *et al.*, "40 nm ultralow-power charge-trap embedded NVM technology for IoT applications," in *Intl. Memory Workshop (IMW)*, 2018, pp. 187–190.
- [15] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Conf. on Computer Vision and Pattern Recognition (CVPR)*, June 2016, pp. 770–778.
- [16] V. J. Reddi *et al.*, "MLPerf Inference Benchmark," in *Intl. Symp. on Computer Architecture (ISCA)*, June 2020.
- [17] A. Coates, B. Huval, T. Wang, D. J. Wu, A. Y. Ng, and B. Catanzaro, "Deep learning with COTS HPC systems," in *Intl. Conference on Machine Learning (ICML)*, 2013, pp. III–1337–III–1345.
- [18] V. Sze, Y. Chen, T. Yang, and J. S. Emer, "Efficient processing of deep neural networks: A tutorial and survey," *Proceedings of the IEEE*, vol. 105, no. 12, pp. 2295–2329, 2017.
- [19] 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 *Design Automation Conf. (DAC)*, June 2016, pp. 1–6.
- [20] C. Li *et al.*, "Analogue signal and image processing with large memristor crossbars," *Nature Electronics*, vol. 1, no. 1, pp. 52–59, 2018.
- [21] P. Yao, H. Wu, B. Gao, J. Tang, Q. Zhang, W. Zhang, J. J. Yang, and H. Qian, "Fully hardware-implemented memristor convolutional neural network," *Nature*, vol. 577, no. 7792, pp. 641–646, 2020.
- [22] S. Ambrogio *et al.*, "Reducing the impact of phase-change memory conductance drift on the inference of large-scale hardware neural networks," in *Intl. Electron Devices Meeting (IEDM)*, 2019, pp. 6.1.1–6.1.4.
- [23] X. Guo, F. M. Bayat, M. Bavandpour, M. Klachko, M. R. Mahmoodi, M. Prezioso, K. K. Likharev, and D. B. Strukov, "Fast, energy-efficient, robust, and reproducible mixed-signal neuromorphic classifier based on embedded NOR flash memory technology," in *Intl. Electron Devices Meeting (IEDM)*, Dec. 2017, pp. 6.5.1–6.5.4.
- [24] P. Wang, F. Xu, B. Wang, B. Gao, H. Wu, H. Qian, and S. Yu, "Three-dimensional NAND flash for vector-matrix multiplication," *IEEE Trans. on Very Large Scale Integration (VLSI) Systems*, vol. 27, no. 4, pp. 988–991, 2019.
- [25] L. Fick, D. Blaauw, D. Sylvester, S. Skrzyniarz, M. Parikh, and D. Fick, "Analog in-memory subthreshold deep neural network accelerator," in *Custom Integrated Circuits Conf. (CICC)*, May 2017, pp. 1–4.
- [26] Y. Du, L. Du, X. Gu, J. Du, X. S. Wang, B. Hu, M. Jiang, X. Chen, S. S. Iyer, and M. F. Chang, "An analog neural network computing engine using CMOS-compatible charge-trap-transistor (CTT)," *IEEE Trans. on Computer-Aided Design of Integrated Circuits and Systems*, vol. 38, no. 10, pp. 1811–1819, 2019.
- [27] C. H. Bennett *et al.*, "Device-aware inference operations in SONOS nonvolatile memory arrays," in *Intl. Reliability Physics Symp. (IRPS)*, 2020, pp. 3C.2.1–3C.2.6.
- [28] Y. Long, D. Kim, E. Lee, P. Saha, B. A. Mudassar, X. She, A. I. Khan, and S. Mukhopadhyay, "A ferroelectric FET-based processing-in-memory architecture for DNN acceleration," *IEEE Journal on Exploratory Solid-State Computational Devices and Circuits*, vol. 5, no. 2, pp. 113–122, 2019.
- [29] A. Ankit *et al.*, "PUMA: A programmable ultra-efficient memristor-based accelerator for machine learning inference," in *Intl. Conf. on Architectural Support for Programming Languages and Operating Systems (ASPLOS)*, April 2019, p. 715–731.
- [30] T.-H. Yang, H.-Y. Cheng, C.-L. Yang, I.-C. Tseng, H.-W. Hu, H.-S. Chang, and H.-P. Li, "Sparse ReRAM engine: Joint exploration of activation and weight sparsity in compressed neural networks," in *Intl. Symp. on Computer Architecture (ISCA)*, 2019, p. 236–249.
- [31] T. P. Xiao, B. Feinberg, C. H. Bennett, V. Prabhakar, P. Saxena, V. Agrawal, S. Agarwal, and M. J. Marinella, "On the accuracy of analog neural network inference accelerators," *arXiv preprint arXiv:2109.01262*, 2021.
- [32] Y. LeCun, "The MNIST database of handwritten digits," <http://yann.lecun.com/exdb/mnist/>.
- [33] A. Krizhevsky, "Learning multiple layers of features from tiny images," 2009.
- [34] B. Jacob, S. Kligys, B. Chen, M. Zhu, M. Tang, A. Howard, H. Adam, and D. Kalenichenko, "Quantization and training of neural networks for efficient integer-arithmetic-only inference," in *Conf. on Computer Vision and Pattern Recognition (CVPR)*, June 2018, pp. 2704–2713.
- [35] L. Chua and S. M. Kang, "Memristive devices and systems," *Proceedings of the IEEE*, vol. 64, no. 2, pp. 209–223, 1976.
- [36] D. B. Strukov, G. S. Snider, D. R. Stewart, and R. S. Williams, "The missing memristor found," *Nature*, vol. 453, no. 7191, pp. 80–83, 2008.
- [37] D. Ielmini, "Resistive switching memories based on metal oxides: mechanisms, reliability and scaling," *Semiconductor Science and Technology*, vol. 31, no. 6, p. 063002, may 2016.
- [38] S. M. Sze and K. K. Ng, *Physics of Semiconductor Devices*. John Wiley & Sons, 2006.
- [39] Y.-H. Hsiao, H.-T. Lue, M. Y. Lee, S.-C. Huang, T.-Y. Chou, S.-Y. Wang, K.-Y. Hsieh, R. Liu, and C.-Y. Lu, "A study of SONOS charge loss mechanism after hot-hole stressing using trap-layer engineering and electrical re-fill methods," in *Intl. Reliability Physics Symp. (IRPS)*, 2008, pp. 695–696.
- [40] H.-T. Lue, Y.-H. Hsiao, Y.-H. Shih, E.-K. Lai, K.-Y. Hsieh, R. Liu, and C.-Y. Lu, "Study of charge loss mechanism of SONOS-type devices using hot-hole erase and methods to improve the charge retention," in *Intl. Reliability Physics Symp. (IRPS)*, 2006, pp. 523–529.
- [41] V. Agrawal, V. Prabhakar, K. Ramkumar, L. Hinh, S. Saha, S. Samanta, and R. Kapre, "In-memory computing array using 40nm multibit SONOS achieving 100 TOPS/W energy efficiency for deep neural network edge inference accelerators," in *Intl. Memory Workshop (IMW)*, May 2020.
- [42] Y. Cai, Y. Luo, E. F. Haratsch, K. Mai, and O. Mutlu, "Data retention in MLC NAND flash memory: Characterization, optimization, and recovery," in *Intl. Symp. on High Performance Computer Architecture (HPCA)*, 2015, pp. 551–563.
- [43] K. Hung, P. Ko, C. Hu, and Y. Cheng, "A unified model for the flicker noise in metal-oxide-semiconductor field-effect transistors," *IEEE Trans. on Electron Devices*, vol. 37, no. 3, pp. 654–665, 1990.
- [44] S. Han, H. Mao, and W. J. Dally, "Deep Compression: Compressing deep neural networks with pruning, trained quantization and Huffman coding," *arXiv preprint arXiv:1510.00149*, 2015.
- [45] D. Miyashita, E. H. Lee, and B. Murmann, "Convolutional neural networks using logarithmic data representation," *arXiv preprint arXiv:1603.01025*, 2016.
- [46] V. Milo, F. Anzalone, C. Zambelli, E. Pérez, M. K. Mahadevaiah, s. G. Ossorio, P. Olivo, C. Wenger, and D. Ielmini, "Optimized programming algorithms for multilevel RRAM in hardware neural networks," in *Intl. Reliability Physics Symp. (IRPS)*, 2021, pp. 1–6.
- [47] S. Agarwal, S. J. Plimpton, R. K. Schiek, I. Richter, A. H. Hsia, D. R. Hughtart, R. B. Jacobs-Gedrim, C. D. James, and M. J. Marinella, (2017) CrossSim. [Online]. Available: <https://cross-sim.sandia.gov/>
- [48] T. P. Xiao, B. Feinberg, J. Rohan, C. Bennett, S. Agarwal, and M. Marinella, "Analysis and mitigation of parasitic resistance effects for analog in-memory neural network acceleration," *Semiconductor Science and Technology*, vol. 36, no. 11, p. 114004, 2021.
- [49] B. D. Salvo *et al.*, "Performance and reliability features of advanced nonvolatile memories based on discrete traps (silicon nanocrystals, SONOS)," *Trans. on Device and Materials Reliability*, vol. 4, no. 3, pp. 377–389, 2004.
- [50] M. J. Marinella, S. Agarwal, A. Hsia, I. Richter, R. Jacobs-Gedrim, J. Niroula, S. J. Plimpton, E. Ipek, and C. D. James, "Multiscale co-design analysis of energy, latency, area, and accuracy of a ReRAM analog neural training accelerator," *IEEE Journal on Emerging and Selected Topics in Circuits and Systems (JETCAS)*, vol. 8, no. 1, pp. 86–101, 2018.
- [51] F. Sequin and A. Fabre, "New second generation current conveyor with reduced parasitic resistance and bandpass filter application," *IEEE Trans. on Circuits and Systems I: Fundamental Theory and Applications*, vol. 48, no. 6, pp. 781–785, 2001.
- [52] M. Bavandpour, S. Sahay, M. R. Mahmoodi, and D. Strukov, "Efficient mixed-signal neurocomputing via successive integration and rescaling," *IEEE Trans. on Very Large Scale Integration (VLSI) Systems*, vol. 28, no. 3, pp. 823–827, 2020.
- [53] L. Kull, T. Toifl, M. Schmatz, P. A. Francese, C. Menolfi, M. Brändli, M. Kossel, T. Morf, T. M. Andersen, and Y. Leblebici, "A 3.1 mW 8b 1.2 GS/s single-channel asynchronous SAR ADC with alternate comparators for enhanced speed in 32 nm digital SOI CMOS," *IEEE Journal of Solid-State Circuits (JSSC)*, vol. 48, no. 12, pp. 3049–3058, 2013.
- [54] M. Baert and W. Dehaene, "20.1 a 5GS/s 7.2 ENOB time-interleaved VCO-based ADC achieving 30.5 fJ/conv-step," in *Intl. Solid-State Circuits Conf. (ISSCC)*, Feb. 2019, pp. 328–330.
- [55] R. Balasubramonian, A. B. Kahng, N. Muralimanohar, A. Shafiee, and V. Srinivas, "CACTI 7: New tools for interconnect exploration in innovative off-chip memories," *ACM Trans. on Architecture and Code Optimization (TACO)*, vol. 14, no. 2, pp. 14:1–14:25, June 2017.

- [56] A. B. Kahng, B. Lin, and S. Nath, "ORION3.0: A comprehensive NoC router estimation tool," *IEEE Embedded Systems Letters*, vol. 7, no. 2, pp. 41–45, 2015.
- [57] Y. Jiao, L. Han, and X. Long, "Hanguang 800 NPU - the ultimate AI inference solution for data centers," in *Hot Chips Symp.*, 2020, pp. 1–29.
- [58] (2021, Sept.) Qualcomm Cloud AI 100 Announcement. Qualcomm. [Online]. Available: <https://www.qualcomm.com/media/documents/files/qualcomm-cloud-ai-100-announcement.pdf>



T. Patrick Xiao (M'17) received the B.A. degree in physics and the Ph.D. degree in electrical engineering and computer sciences from the University of California, Berkeley in 2014 and 2019, respectively. He is currently a Senior Member of Technical Staff at Sandia National Laboratories. His research interests include physics-based or analog computing, low-power machine learning (ML) accelerators, non-volatile memory (NVM), radiation-hard electronics, and multi-scale modeling of devices through algorithms.

Ben Feinberg (S'15–M'19) received the B.S. degree in electrical and computer engineering and the M.S. and Ph.D. degrees in electrical engineering from the University of Rochester in 2012, 2014, and 2019, respectively. He is currently a Senior Member of Technical Staff at Sandia National Laboratories. His research is in computer architecture with an emphasis on memory-centric accelerators, heterogeneous architectures, and energy-efficient data representations.



Christopher H. Bennett (M'14) received the B.Sc. and M.Sc. degrees from Stanford University in 2011, the EMM-Nano M.Sc. degree from KU Leuven and the Chalmers University of Technology in 2014, and the Ph.D. degree from Université Paris-Saclay in 2018. He is currently a Senior Member of Technical Staff at Sandia National Laboratories. He is an emerging devices expert with varied background on polymer, spintronic, filamentary, and charge-based devices used for on-chip learning and edge inference scenarios. At Sandia, he explores these interests to

yield reliable, interpretable, and scalable future accelerator designs.



Vineet Agrawal received the B. Tech and M. Tech degrees from the Indian Institute of Technology (IIT) Madras in 2006. He is currently a Senior Design Director at Infineon Technologies, where he works on developing ultra-low power ML accelerators and differentiated memory products. Vineet has over 18 years of experience in advanced memory design and technology development and holds over 15 US patents. His research interests include the design and architecture of NVM based compute-in-memory accelerators for edge AI applications.

Prashant Saxena, photograph and biography not available at time of publication.



Venkatraman Prabhakar received the Ph.D. degree from the University of California at Los Angeles in 1996 and the B. Tech degree from IIT Madras, India. He is in charge of SONOS technology development at Infineon. He has extensive experience in integrating SONOS flash memories into advanced technology nodes from 130nm down to 22nm.



Krishnaswamy Ramukmar received the Ph.D. degree from IIT Bangalore, India in the field of Semiconductor Devices. He has been a faculty member at the Indian Institute for several years and visiting scientist at Rensselaer Polytechnic Institute in Troy, NY before joining Cypress Semiconductor in San Jose, CA in 1993. He is currently a VP Fellow at Infineon Technologies (after acquisition of Cypress) where he is involved in the development of SONOS technologies.



Harsha Medu received the B. Eng degree in Electronics and Communication from Visvesvaraya Technological University (VTU) in 2003 and a Master of Business Administration from the Indian Institute of Management - Bengaluru (IIMB) in 2014. He is currently a principal applications engineer at Infineon Technologies and has worked on the design and application of various NVM products.



Vijay Raghavan (M'98) received the B. Tech degree in electrical and electronics engineering from the Indian Institute of Technology, Chennai, India, in 1996, and the M.S. degree in electrical engineering from the Georgia Institute of Technology in 1997. He was with Cypress Semiconductor until 2020 and served as the Senior Director of Design Engineering. He is currently a Distinguished Engineer and the Center of Excellence Lead of the embedded NVM group at Infineon Technologies in Austin, TX. His areas of interest and expertise include memory (non-volatile and volatile) designs, analog designs, in-memory-compute hardware accelerator architecture and designs and radiation hardened designs. He has more than 45 patents in the field of memory and analog designs.

volatile and volatile) designs, analog designs, in-memory-compute hardware accelerator architecture and designs and radiation hardened designs. He has more than 45 patents in the field of memory and analog designs.



Ramesh Chettuvetty received his B. E. degree from the University of Calicut, Kerala in 2000. He spent more than a decade at Texas Instruments specializing in SAR, sigma delta and pipeline ADCs, DACs, power management, high-speed interfaces and high-performance timing solution circuits. He joined Cypress Semiconductor in 2011 and led various design, applications and marketing organizations and running business lines. Since 2019, he is with Infineon Technologies heading the marketing and applications organization of memory products. He is the lead

inventor of 3 patents on analog circuits, systems and AI architectures and has co-authored several publications on semiconductor circuits and systems.



Sapan Agarwal (M'06) received the B.S. degree in electrical engineering from the University of Illinois at Urbana-Champaign in 2007, and the Ph.D. degree in electrical engineering from the University of California, Berkeley in 2012. He is currently a Principal Member of Technical Staff at Sandia National Laboratories. He is interested in everything from explainable machine learning algorithms and neuromorphic systems to novel semiconductor devices.



Matthew J. Marinella (SM'17) received the Ph.D. degree in electrical engineering from Arizona State University in 2008. He is a Distinguished Member of the Technical Staff in Sandia National Laboratories' Microsystems Center, where he leads the Nonvolatile Memory Technology R&D Program, and several internal and externally funded research projects involving neuromorphic and low-power computing with emerging electronic devices. Dr. Marinella has authored or co-authored over 100 peer reviewed publications, given numerous invited and contributed talks, and presented several short courses on these topics. He has served in technical and leadership roles in Lab- and DOE-level initiatives on emerging computing and is a member of the SRC Decadal Plan Executive Committee, chairs Emerging Memory Devices for the IRDS Roadmap, and serves on various technical program committees.

talks, and presented several short courses on these topics. He has served in technical and leadership roles in Lab- and DOE-level initiatives on emerging computing and is a member of the SRC Decadal Plan Executive Committee, chairs Emerging Memory Devices for the IRDS Roadmap, and serves on various technical program committees.