

# Neuromorphic Computing: From Devices to Integrated Circuits

V. Saxena<sup>1, a)</sup>Electrical and Computer Engineering Department,  
University of Delaware, Newark, DE.

(Dated: 11 November 2020)

A variety of nonvolatile memory (NVM) devices including the resistive Random Access Memory (RRAM) are currently being investigated for implementing energy-efficient hardware for deep learning and artificial intelligence at the edge. RRAM devices are employed in the form of dense crosspoint or crossbar arrays. In order to exploit the high-density and low-power operation of these devices, circuit designers need to accommodate their non-ideal behavior and consider their impact on circuit design and algorithm performance. Hybrid integration of RRAMs with standard CMOS technology is spurring development of large-scale neuromorphic system-on-a-chip (NeuSoC). This review article provides an overview of neuromorphic integrated circuits using hybrid CMOS-RRAM integration with an emphasis on spiking neural networks (SNNs), device nonidealities, their associated circuit design challenges, and potential strategies for their mitigation. An overview of various SNN learning algorithms and their co-development with devices and circuits is discussed. Finally a comparison of NVM-based fully-integrated neuromorphic ICs is presented along with a discussion on their future evolution.

PACS numbers: 84.30.-r, 84.35.+i,

Keywords: Artificial Intelligence (AI), Back-end-of-the-line (BEOL) Processing, CMOS Neurons, Edge-AI, Non-Volatile Memory (NVM), Neuromorphic Computing, Resistive Random Access Memory (RRAM).

## I. INTRODUCTION

Over the last decade, *deep neural networks (DNNs)*, or deep learning, have emerged as the next wave of Artificial Intelligence (AI) which has been propelled by the advances in specialized hardware, open-source software and availability of datasets. Hardware platforms such as the graphics processing units (GPUs) and application-specific integrated circuits (ASICs) for AI acceleration enable parallel processing of a large amount of data to train DNN models. Training or learning in DNN models is performed using variants of the gradient-descent *back-propagation*, or *Backprop*, algorithm<sup>1-3</sup> which is both computationally and energy intensive due to the massive amounts of data continuously shuttled between memory and processing units. Recently, low-power GPUs and ASICs have appeared for deep learning inference for Edge-AI<sup>4,5</sup>. However, training is performed on a GPU-based server or Cloud infrastructure using software frameworks such as TensorFlow<sup>6</sup> and PyTorch<sup>7</sup>. In view of diminishing returns from such architectures with the near-end of Moore's scaling<sup>8</sup>, semiconductor industry's International Roadmap for Devices and Systems (IRDS) looks forward to *Beyond-Moore* or *post-CMOS* technologies to conceive radically new computing architectures for AI workloads<sup>9</sup>. This requires a cross-layer investigation of non von Neumann computing architectures across the entire devices, circuits and algorithms hierarchy.

Several classes of emerging non-volatile memory (NVM) devices are currently being investigated for their application in analog implementation of Edge-AI hard-

ware. These devices can be two- or three-terminal and employ a wide range of material systems and associated physical mechanisms to achieve multi-level non-volatile memory states. Moreover, these NVM devices need to be monolithically integrated with standard CMOS technology to enable hybrid integrated circuit design. Deep learning hardware realizations include *vector-by-matrix multipliers* (VMMs), and neural-inspired or *Neuromorphic* computing circuits. The NVM devices are employed in the form of crossbar, or cross-point, arrays with or without selectors along with CMOS circuits at the periphery of the array.

While these emerging in-memory computing architectures show promise, in order to exploit the high-density and operation of these devices, *integrated circuit* (IC) designers need to accommodate their realistic behavior and nonidealities. This is particularly important for optimizing hybrid transistor-NVM circuit design for performance, area and power consumption. NVM nonidealities include device variability, low resistances offered by the two-terminal devices, resolution and stability of multi-level states, nonlinearity and device endurance. Moreover, novel in-situ learning algorithms must be developed which can take advantage of the localized in-memory computing to minimize von Neumann bottlenecks.

Several recent review articles in the literature focus on emerging devices for in-memory computing<sup>10-16</sup>, neuromorphic learning algorithms<sup>17</sup>, pathways and survey of neuromorphic hardware architectures<sup>18,19</sup>. This article provides the motivation and overview for spike-based Neuromorphic hardware that can be realized using emerging NVMs, with a specific focus on the Resistive Random Access Memory (RRAM), *aka* the *memristors*, and *spiking neural network* (SNN) algorithms. The unique aspect of this review is that it focuses on bridging

<sup>a)</sup>Electronic mail: vsaxena@udel.edu

the large gap between the expectations produced by the experimental results from emerging NVM devices and the requirements set by the neuromorphic system architectures. This is attempted by considering the device non-idealities that the circuit designers need to accommodate and the hardware bottlenecks to be surmounted when implementing spike-based deep learning in neuromorphic hardware. Rest of the article is organized as follows: Section II introduces SNNs and their use in neuromorphic computing. Section III provides an overview of the RRAM devices and discusses the impact of device nonidealities on circuit design considerations and system-level performance. Section IV discusses CMOS integrated circuits needed to realize neuromorphic computing. Section V provides a brief overview of the learning algorithms for deep SNNs. Finally, Section VI benchmarks the performance of fabricated NVM-based neuromorphic computing ICs followed by a concluding discussion.

## II. SPIKING NEURAL NETWORKS (SNNs) AND NEUROMORPHIC COMPUTING

Energy-efficiency of DNNs realized on GPUs and ASICs based on von Neumann architectures is fundamentally limited by the energy and latency cost of the ‘distance’ between storage memory and processing units<sup>20</sup>. In a radical contrast, a biological brain stores memory and performs localized computing using similar neural motifs with extremely high energy-efficiency, thus making a compelling inspiration for in-memory computing for DNNs.

Fig. 1 illustrates a neuromorphic computing architecture based on 1T1R crosspoint arrays with neuron circuits at the array periphery. In a fully-connected neural network, each array realizes a network layer which are connected with each other using on-chip and/or off-chip interconnects. These interconnects employ asynchronous spike-based communication using a protocol such as the address-event representation (AER) protocol, which is widely used in neuromorphic sensors and processors<sup>21–23</sup>. The neural network weights are stored in the NVM memory array. The pre-neurons concurrently drive the rows (or the wordlines), as opposed to random access. The pre-neuron activations (or voltages) are weighted by the conductance of the synaptic weights and the resulting current is summed and integrated on the post-neurons connected to the columns (or the bitlines).

In the past decade, advances in spike-based models with localized plasticity mechanisms such as the *spike-timing-dependent-plasticity* (STDP) and its feedback-based modulation<sup>24–30</sup> have opened new avenues in neuromorphic computing research. For example, theoretical studies have suggested STDP-like plasticity mechanisms can be used to train two-layer SNNs in-situ without trading-off their parallelism<sup>31–34</sup>. SNNs essentially encode information using asynchronous spatiotemporal ‘spikes.’ A spikes sequence can represent the input sig-

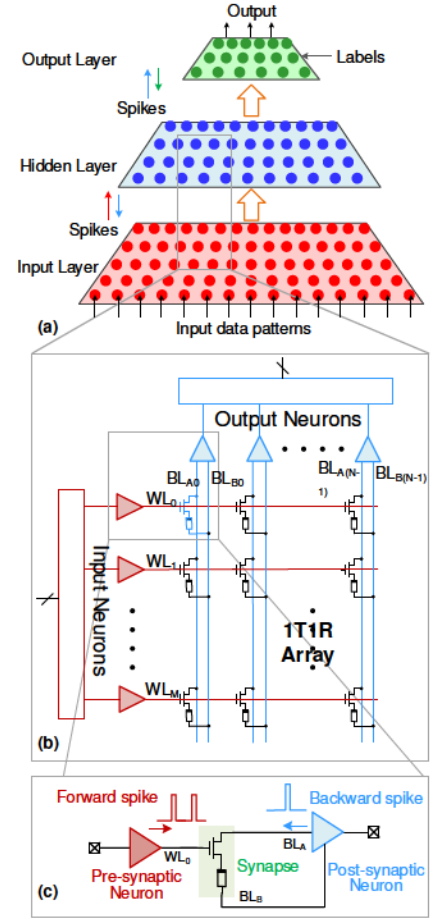


FIG. 1. A neuromorphic computing architecture built around 1T1R RRAM crosspoint arrays: (a) A fully-connected spiking neural network (SNN) showing input, hidden and output layers of spiking neurons, (b) a crosspoint RRAM array with input and output neurons, (c) a single synapse between the input and output neurons with localized weight updates.

nal either using *rate-coding* where the average spike rate represents a real valued signal, or temporal coding where spike delay (or latency) encodes the same information. In biology, sensory information is encoded as a combination of both rate and latency coding.

### A. Leaky Integrate and Fire Neurons

Biological neurons exhibit complex spike signal processing and spike generation behavior with some of the spike filtering occurring in the dendrites<sup>35</sup>. As a result, wide range of models have been developed in the literature to emulate their response<sup>21,36–40</sup>. However, in neuromorphic computing only the simple salient features useful for learning and inference are adapted into the SNN architecture. A single-compartment *leaky integrate-and-fire* (LIF) neuron is a simple and commonly-used neural motif. Network weights,  $w_{ij}$ , are stored in ‘synapses’ which



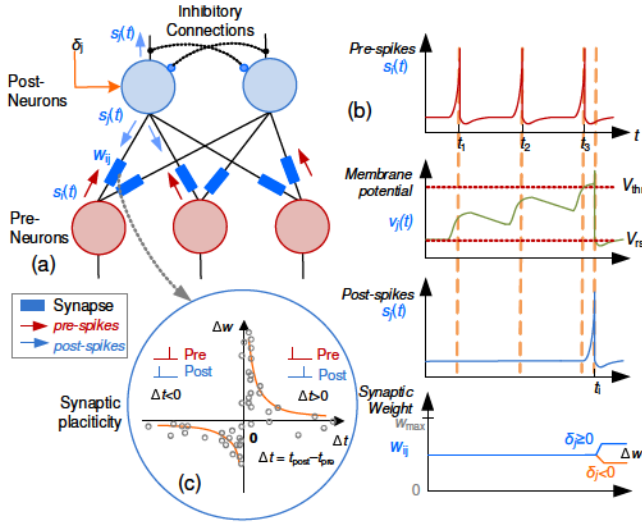


FIG. 2. (a) Section of a two-layer spiking neural network illustrating pre and post spike flow. (b) Transient waveforms for these spikes, membrane potential, and weight update modulated by an external feedback error,  $\delta$ . (c) Spike-timing dependent plasticity (STDP) learning window where the weight update depends upon the temporal delay,  $\Delta t$ , between the pre and post spikes at a synapse.

can be updated based on Hebbian learning, i.e. the correlation of the spiking activity of the pre-synaptic and post-synaptic neurons. During feedforward passes, a LIF neuron integrates its input spikes,  $s_i(t)$ , weighted by the synapses, into its membrane potential

$$v_{m,j}(t) = \sum_{i=1}^N w_{ij} \int_{-\infty}^t s_i(t) \otimes h(t) dt \quad (1)$$

where  $h(t) = e^{-t/\tau_m} u(t)$  is the impulse response that incorporates neuron's leaky behavior. When the membrane potential of a neuron crosses its firing threshold, i.e.  $v_{m,j}(t) > V_{thr,j}$ , the neuron produces a post-spike,  $s_j(t)$ . A firing event resets the membrane potential to a resting potential,  $V_{rst}$  and propagates the post-spike in the forward as well as the backward direction. This is described by Eq. 2 below and illustrated in Fig. 2.

$$\text{at } t = t_j^f \quad v_{m,j}^l(t) > V_{thr,j}^l : \begin{cases} v_{m,j}^l(t) \leftarrow V_{rst} \\ s_j(t) \leftarrow g(t - t_j^f) \end{cases} \quad (2)$$

After firing, a neuron enters a refractory, or silent, period before it can fire again. During the refractory period of duration  $\tau_{refr}$ , membrane potential rests at  $V_{rst}$ .  $\tau_{refr}$  sets the maximum rate at which a neuron can fire, i.e.  $\nu_{max} \geq \frac{1}{\tau_{refr}}$ , and thus the resolution for quantization in the time-domain.

Neurons can form inhibitory and/or excitatory connections with other neurons as seen in Fig. 2. Inhibitory connections prevent other neurons from firing if one of

the neurons has fired in a given time window. Excitatory connections make connected neurons fire simultaneously. Inhibition is used to implement competitive learning in an population of neurons to realize Winner-take-all (WTA) motifs<sup>41</sup>.

## B. Synapses and Plasticity

Synapses essentially implement weights in the SNN where the pre-spikes are weighted by the conductance of the synapse before they are integrated in the neuron. In case of online learning, synaptic weights are in-situ updated on the neuromorphic chip. Weight update in a synapse is governed by a plasticity rule based on Hebbian correlation which can be modulated by an error feedback, leading to an approximated three-factor learning rule<sup>42</sup>

$$\Delta w_{ij} = f(t_i, t_j, \delta_j) \quad (3)$$

Here,  $t_i$  and  $t_j$  are pre and post spike times, and  $\delta_j$  is the error feedback. Several plasticity rules have appeared in literature based on electrical probing of biological neurons and update rule derivations in computational neuroscience<sup>26,29,43</sup>. STDP rule is a two-factor rule based on the time difference between post and pre spikes  $\Delta t_{ij} = t_{post} - t_{pre} = t_j - t_i$  and expressed as

$$\Delta w_{ij} = \begin{cases} a^+ e^{-\frac{\Delta t_{ij}}{\tau^+}}, & \Delta t_{ij} = t_j - t_i \geq 0 \\ -a^- e^{-\frac{\Delta t_{ij}}{\tau^-}}, & \Delta t_{ij} = t_j - t_i < 0 \end{cases} \quad (4)$$

Here,  $a^+$  and  $\tau^+$  are the voltage and temporal parameters for the *long-term potentiation* (LTP) when the synaptic weight increases. Conversely, Here,  $a^-$  and  $\tau^-$  are the corresponding parameters for the *long-term depression* (LTD) when the synaptic weight decreases.

Fundamentally, spike-based computing simplifies the need for precise multiplication and replaces it by scaling of bi-level, or binary, spikes by synaptic conductances. This results in simpler digital neuromorphic hardware implementations and precludes the need for four-quadrant multiplication in analog neuromorphic realizations. Consequently, SNNs can perform computing with simpler hardware and consume very small amount of energy only when the spike events occur<sup>23,44,45</sup>. Due to their unique characteristics, SNNs are better realized on an event-driven neuromorphic platform instead of a von Neumann computer.

## C. Analog Mixed-Signal Neuromorphic Platforms

As discussed earlier, neural-inspiration provides the motivation for realizing dense and low-power neuromorphic hardware by implementing neuron and synapses using analog mixed-signal circuits. Advances in analog neuromorphic circuits include the Neurogrid hardware<sup>46</sup>,

where subthreshold biomimetic CMOS circuits were developed to reproduce dynamics occurring in biological neural networks. However, the fundamental limitation of such analog architectures is that the weights are dynamically stored and updated on capacitors, which leak away in a few seconds, thus limiting any long-term learning. Bistability of analog weights has been used as a stop-gap solution<sup>47,48</sup>, but recent studies on deep neural networks have determined that 4-bit or higher synaptic resolution is needed to realize SNNs with reasonable classification accuracy<sup>49,50</sup>. Other solutions include using *Floating Gate*, or *NOR Flash*, devices for realizing multilevel non-volatile synaptic weights<sup>51</sup>. However, despite of their excellent retention, floating-gate memory suffers from low endurance of  $< 10^5$  write cycles. This limits the number of times the neural network weights can be updated and thus significantly limit on-chip training capability.

#### D. Neuromorphic Computing using Emerging Nonvolatile Memory Devices

Several emerging *nonvolatile memory* (NVM) devices including RRAMs, *spin-torque transfer RAM* (STTRAM), *phase-change RAM* (PCRAM), also generally referred to as memristors, have been considered for their use in Edge-AI and neuromorphic computing<sup>20,55–59</sup>. Recently, ferroelectric-field-effect transistors (FeFETs) have come under focus for their low power performance<sup>54,60</sup>. In the last decade, these NVM devices have been extensively investigated as a high-density potential replacement for Flash memory<sup>20,52,61</sup> with their advantages summarized in Table I. As the research community gain a better understanding of the physical mechanisms for electrical switching in the RRAM devices, there is a trend towards developing novel applications by leveraging these devices for computing<sup>57,62–67</sup>.

These emerging NVM devices are employed in the form of *crossbar* or *cross-point* arrays with a diode or transistor selector—cell with *one transistor and one RRAM* (1T1R). These NVM arrays are being investigated for DNN computations in the analog domain. These include VMM for Edge-AI<sup>58,68</sup> and synaptic weights in a mixed-signal neuromorphic hardware<sup>23,69</sup>. Emerging NVM arrays are promising for neuromorphic computing as they provide: (1) a very high synaptic density with low leakage power, (2) localized in-memory learning similar to biological synaptic plasticity, (3) very low power consumption with event-driven updates, and (4) non-volatility of weights<sup>23</sup>. This review focuses on two-terminal RRAM devices for neuromorphic computation.

### III. RRAM DEVICE CHARACTERISTICS AND PROCESS INTEGRATION

#### A. RRAM Devices

Resistance switching in emerging nanoscale resistive memory devices has sustained interest with the goal of high-density and lower power replacement for NVM-based embedded memory and computing applications<sup>9,10,52,70</sup>. RRAM arrays are considered as a suitable alternative to the Flash-based solutions due to the following reasons:

- (i) The filamentary nature of resistance switching has the potential to scale well beyond the sub-10nm feature size<sup>61,71</sup>.
- (ii) Lower switching voltages allow low power operation and compatibility with scaled-CMOS<sup>952</sup>.
- (iii) Very simple planar two-terminal structures and fabrication-friendly materials facilitate integration with standard CMOS technology.
- (iv) These devices demonstrate biologically plausible plasticity (i.e. weight update) behavior in several experiments,<sup>57,62,72–74</sup> and therefore have emerged as an ideal candidate for realizing electrical equivalent of biological synapses.

Several categories of RRAMs have been intensely pursued by the device community: (1) Electrochemical memory (ECM) *aka* Conductive Bridge RAM (CBRAM), (2) Valence Change Memory (VCM) *aka* Oxygen Vacancy based RAM (OxRAM), (3) Thermochemical memory (TCM), and (4) Interfacial or 2D switching RRAMs<sup>75–77</sup>. This review primarily focuses on CBRAM and OxRAM devices for neuromorphic computing, which are graphically depicted in Fig. 3. All these RRAMs are essentially metal-insulator-metal (MIM) structures where the resistance between the two electrodes can be changed in a non-volatile manner by either filamentary or interfacial-type switching.

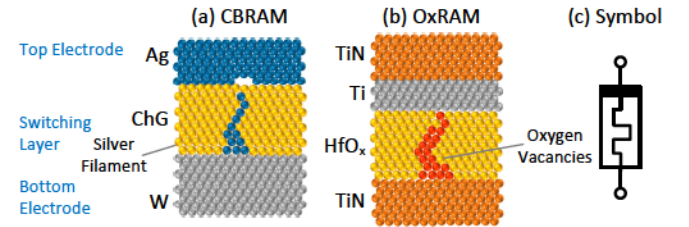


FIG. 3. Graphical illustration of the cross-section of CBRAM and OxRAM NVM devices and their circuit symbol. Device cross-sections are not to the scale.

**CBRAM:** CBRAM uses a dielectric layer sandwiched between two electrodes and the resistance is changed by forming an conductive electrochemical bridge. As shown in Fig. 3(a), the anode (top electrode) is electrochemically active and is made of metals such as Silver (Ag)<sup>78–81</sup> or Copper-based alloys (e.g. CuTe<sub>x</sub>)<sup>82</sup>.



TABLE I. Comparison of memory technologies for in-memory computing. Data reproduced from Refs.<sup>20,52–54</sup>.

Parameters	SRAM	DRAM	NOR Flash	PCRAM	STTRAM	OxRAM	CBRAM	FeFET
Cell size	100F <sup>2</sup>	7F <sup>2</sup>	5F <sup>2</sup>	4F <sup>2</sup>	12F <sup>2</sup>	4-6F <sup>2</sup>	4-6F <sup>2</sup>	24F <sup>2</sup>
Density	Low	High	High	V. High	High	V. High	V. High	High
Write Latency	1ns	5ns	10 $\mu$ -1ms	100ns	2-25ns	10ns	10ns	3ns
Read Latency	1ns	20-80ns	50ns	10ns	2-25ns	1-10ns	10ns	2ns
Write Energy (pJ/bit)	<1	<1	100	2-25	0.1-2.5	0.1-3	0.1-10	0.1
Leakage Power	High	Medium	Low	Low	Low	Low	Low	Low
Endurance (write cycles)	> 10 <sup>15</sup>	> 10 <sup>15</sup>	10 <sup>5</sup>	10 <sup>8</sup>	10 <sup>12</sup>	10 <sup>9</sup>	10 <sup>8</sup>	> 10 <sup>5</sup>
Switching Polarity	NA	NA	Bipolar	Unipolar	Bipolar	Bipolar	Bipolar	Bipolar
MLC Capability	✗	✗	4-8 bits	4-6 bits	2 bits	4-6 bits	2 bits	3 bits
MLC Retention	✗	✗	Years	R-drift	Tunneling	R-drift	R-drift	Tunneling
3D Stacking	✗	✗	✗	✓	✗	✓	✓	✓

The cathode is an inert electrode and usually realized using Platinum (Pt), Nickel (Ni), or Tungsten (W). The dielectric is a thin film layer which comprises of either a chalcogenide (such as Ge<sub>x</sub>Se<sub>1-x</sub>), amorphous silicon<sup>71,72</sup>, Al<sub>2</sub>O<sub>3</sub>, MO<sub>x</sub>, or HfO<sub>x</sub><sup>82</sup>. The reader is referred to<sup>81,83</sup> for the CBRAM fabrication details. An initial step called ‘forming’ step is required to introduce silver into the unformed or pristine switching layer. This is achieved by either applying a sufficiently large voltage pulse (*aka* electroforming) or by photo-diffusion using optical illumination<sup>84,85</sup>.

**OxRAM:** OxRAMs are fabricated as an MIM stack with transition metal oxides as the dielectric. The mechanism for resistance switching is the formation of a conductive filament due to the migration of oxygen vacancy defects<sup>86,87</sup>. Several insulators have been explored in the literature including HfO<sub>x</sub>, Ta<sub>2</sub>O<sub>x</sub>, TiO<sub>x</sub>, ZrO<sub>x</sub>, and NiO. Since forming an oxygen vacancy in these oxides requires higher energy, a metallic layer of Hf (or Ti, La, Zr) is used as a *scavenger electrode* that scavenges oxygen vacancies from HfO<sub>x</sub><sup>88</sup>. This stack is shown in Fig. 3(b). The Hf/HfO<sub>x</sub> (or Ti/HfO<sub>x</sub>, etc.) interface acts as the active electrode<sup>89</sup> which is capped with material such as TiN, Ni, TaN, ITO, or Al<sup>87</sup>. Incorporation of the scavenger electrode leads to better switching control, reduced variability and higher retention<sup>88</sup>. Doping of OxRAM active layer with dopants such as Ti or Ge increases the oxygen vacancies and has been explored to realize forming-free devices<sup>87</sup>.

**Other RRAM Variants:** Example of TCM devices include transition metal oxide cells such as NiO and HfO, where filamentary switching occurs due to current-based heating and the resulting stoichiometric changes due to temperature increase<sup>75,77,90</sup>. Interfacial switching RRAMs include Nb-doped SrTiO<sub>3</sub> (NbO) and doped perovskites such as Pr<sub>0.7</sub>Ca<sub>0.3</sub>MnO<sub>3</sub> (PCMO). These RRAMs are also referred to as 2D switching devices as the resistance is dependent on cross-section area, as opposed to the filamentary switching devices<sup>75,91,92</sup>.

## B. RRAM Electrical Characteristics

**CBRAM:** A typical hysteresis current–voltage (I–V) characteristics in a metal–insulator–metal (MIM) structure are shown in Fig. 4. These electrical switching characteristics were obtained from experimental characterization of Ag/Ge<sub>20</sub>Se<sub>80</sub>/W CBRAM devices which were fabricated by Mitkova group at Boise State University<sup>81,83</sup>. Here, triangular voltage sweeps were applied across the device and the current was measured using a Semiconductor Parameter Analyzer (SPA) such as Keysight B1500. The I–V sweeps were performed for several settings of the compliance current,  $I_{cc}$ , ranging from 50nA to 10 $\mu$ A. The device state during the I–V switching characteristics are described as follows:

(A) Here, initially the device is in the *High-resistance state* (HRS), erased, or Off state.

(B) When a sufficiently large positive voltage greater than the program threshold voltage ( $V_{th}^+$ ) is applied on the top electrode, silver is oxidized and the Ag<sup>+</sup> ions start moving towards the cathode and forming a bridge (or filament) in the process.

(C & D) The silver ion bridge eventually forms a high-conductivity path between the electrodes realizing the *Low-resistance state* (LRS), programmed, or On state. This is referred to as the Program or Set operation.

(E) Conversely, if a voltage more negative than the erase threshold voltage ( $V_{th}^-$ ) is applied across the CBRAM, the filament is dissolved and the device reverts to the HRS. This is called the Erase or Reset operation.

The electrochemical process of conductive filament formation is inherently stochastic and varies across devices and switching cycles, where silver filaments of varying geometry can be formed in the amorphous switching layer. The variability is reduced with scaling to sub-15nm cross-section area as only one dominant filament can be formed.

CBRAM switching thresholds depend upon the active electrode material and switching layer used in the device. In the CBRAM fabricated by author’s collaborators with their experimental switching characteristics shown in Fig. 4, the program threshold is  $V_{th}^+ = 0.7V$  and the

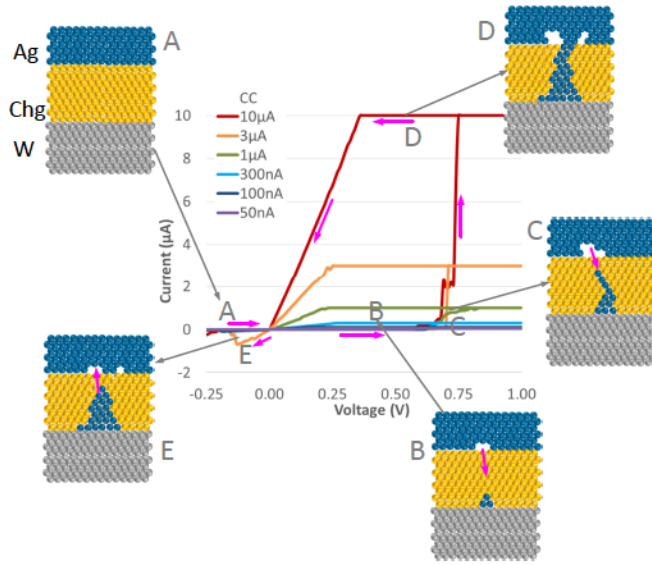


FIG. 4. Experimental I-V characteristics of a CBRAM device with graphical illustration of Program (A-D) and Erase (E) operations.

erase threshold is  $V_{th}^- = -0.1V$ . The CBRAM switching behavior is bipolar, i.e. the Set and Reset operations require different polarity. CBRAM' bipolar switching is asymmetric as a very small negative voltage is needed to break the filament (bridge) at the anode and revert the device to HRS<sup>81,93,94</sup>.

Compliance current is an important parameter that sets the maximum current in the device during the electroforming, and the subsequent operation of the device, and thus essentially determines the LRS resistance,  $R_{LRS}$ . This is illustrated in Fig. 5, where the resistance distribution for varied values of  $I_{cc}$  is plotted. For a large  $I_{cc} > 10\mu A$ , a thick filament is formed resulting in a narrower distribution with lower  $R_{LRS}$  values. This regime is suitable for digital applications with binary or bistable memory.

For lower  $I_{cc} \sim 1\mu A$ , a weak filament is formed with a wider distribution and higher resistance values. These characteristics give rise to the possibility of realizing analog-like behavior where the RRAM resistance can be programmed to one of the multilevel resistance states. This can be thought of as a *multi-level cell* (MLC) with 2-bit or higher resolution. The filament structure determines the switching speed, repeatability, LRS resistance range, multilevel behavior, and the retention of devices.

**OxRAM:** Electrical characteristics of OxRAM follow a I-V hysteresis loop similar to CBRAM with a few differences. Resistance switching in Ti/HfO<sub>x</sub> (Hf/HfO<sub>x</sub>, Ta/Ta<sub>2</sub>O<sub>x</sub> and similar stacks) OxRAMs occurs due to the formation of a conductive filament of oxygen vacancies. The metallic Hf, Ti or similar layers, facilitate scavenging of oxygen atoms from HfO<sub>x</sub> layer. The energy of breaking oxygen (O) is compensated by the exothermic Hf-O bond formation energy<sup>88,89</sup>.

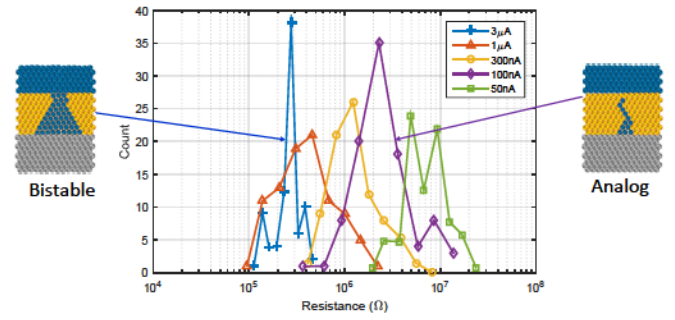


FIG. 5. CBRAM resistance distribution for several values of compliance current. Analog-like states are observed when a weak filament is formed. A thick filament leads to digital switching behavior.

When a positive voltage is applied to the active electrode, positive oxygen vacancies ( $V_O^+$ ) are created at the Ti/HfO<sub>x</sub> interface, which acts as a  $V_O^+$  reservoir. The positively charged vacancies diffuse towards the inert (TiN) electrode held at a lower potential, thus forming a conductive filament. As the filament grows in length, the electron conduction mechanism changes from trap-assisted tunneling to Poole-Frenkel hopping to eventually Ohmic conduction when the filament is fully formed and consigns the OxRAM to LRS<sup>89</sup>. A compliance current ensures that complete breakdown of HfO<sub>x</sub> dielectric is avoided. A 1T1R cell configuration is typically used to set this compliance current in the array. When the voltage polarity across the electrodes is reversed, O<sup>-</sup> diffuses back into HfO<sub>x</sub> and the conductive filament is dissolved leading to HRS.

OxRAMs devices based on transition metal switching layer have been extensively characterized using 1R as well as 1T1R arrays<sup>82,88,95–99</sup>. The program and erase thresholds vary widely depending upon specific device structure and material used. For example in devices from Leti<sup>95,100</sup>, the Program/Set thresholds in the range of  $V_{th}^+ \approx 1V$  and Erase/Reset is  $V_{th}^- \approx -1.5V$ . The bipolar switching in OxRAM is more symmetric than CBRAMs as a larger negative voltage is needed to diffuse the O<sup>-</sup> ions back and dissolve the  $V_O^+$  bridge. The LRS variability depends on  $I_{cc}$  and the pulse duration during the initial forming step. The HRS variability depends upon the target resistance range<sup>95</sup>, which is of interest for analog (or multilevel) applications. For HRS, cycle-to-cycle variability tends to higher than the cell-to-cell variability. Thus a smart program-and-verify algorithm is needed to place the multi-level OxRAM in a specific resistance state. The forming voltage,  $\approx 2V$ , is larger than the program threshold. A smaller  $I_{cc} \approx 5\mu A$  and shorter pulse width ( $\approx 100ns$ ) avoids stuck-at-LRS faults and places the OxRAM in analog region of operation (10kΩ–1MΩ). The program/Set operation can be abrupt in OxRAMs due to a positive feedback loop between electrical field/temperature and oxygen vacancies. The ana-



log behavior can be improved by introducing a *thermal enhanced layer* (TEL) to provide a better control over incremental resistance changes<sup>101?</sup>.

**Other RRAMs:** TCM devices exhibit apolar or unipolar switching, i.e. the Set and Reset operations have the same polarity<sup>75,90</sup>. The advantage of unipolar switching is that a simple diode-based selector can be employed forming a 1D1R cell. Disadvantages of TCM devices include the higher energy consumption in the thermal-activated switching process<sup>90</sup>. Interfacial switching devices such as NbO and PCMO exhibit bipolar switching by modulating the Schottky barrier barrier at the metal-oxide interface<sup>75,91,92</sup>. An advantage of the interfacial switching RRAMs is that they exhibit an asymmetric nonlinearity, which realizes a selector-free self-rectifying cell. Also, forming-free PCMO RRAM devices have been demonstrated which is advantageous for on-chip integration<sup>102</sup>.

### C. RRAM Retention and Endurance

Retention and endurance of both CBRAM and OxRAM have been extensively studied for binary storage applications<sup>52,53</sup>. Development of RRAMs for analog computing where each cell can be programmed with multilevel states is a current focus of research<sup>103–107</sup>. Ti/HfO<sub>x</sub> OxRAMs have demonstrated endurance of  $\approx 10^9$  cycles, ultra-fast writes ( $\sim 5$ ns), and  $R_{off}/R_{on}$  ratio exceeding 100<sup>52,88,98</sup>. Multilevel resistance states are programmed by modulating either the voltage amplitude or the pulse-width of the program/erase pulses, and setting a compliance current. A program-and-verify scheme is employed for iterative programming to achieving tight resistance distributions while minimizing the read/write disturb between the cells<sup>108</sup>. Consequently, the device is stressed multiple times per MLC write event and has a bearing on overall device lifetime. Recent work on MLC OxRAM characterization has shown that it can be programmed up to 64-levels, or 6-bit per cell resolution<sup>66,108–110</sup>. However, the programmed states have been observed to relax and their resistances drift over time on a timescale of hours to weeks<sup>66,108,111</sup>. Several retention failure modes have been observed depending upon the device material stack and structure<sup>112</sup>. Also, the resistance drift is accelerated as the chip temperature is increased due to the increased Brownian motion of conducting ions in the switching layer<sup>66</sup>.

CBRAMs have also shown multi-level resistance states with 2-bit/cell storage and  $\geq 10^5$ s retention by using appropriate compliance current settings<sup>53,113</sup>. Based on literature survey, OxRAMs have been more amenable to analog behavior when compared to CBRAMs, but the latter can be improved by using optimized bilayer materials<sup>53</sup>.

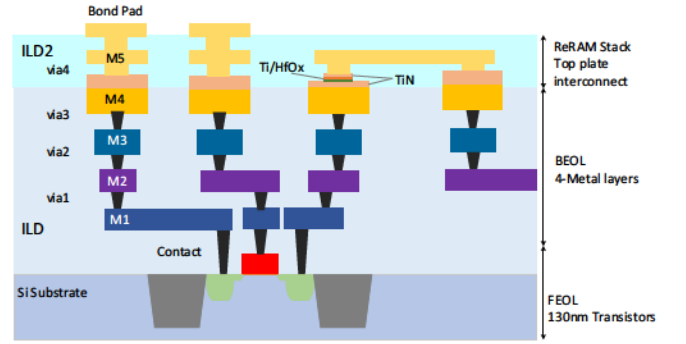


FIG. 6. Graphical illustration of integration of RRAM stack in the BEOL of a 4-Metal CMOS process as described in Ref.<sup>100</sup>

### D. Process Integration of RRAM with CMOS

Integration of emerging NVM devices with standard CMOS technology is essential for the development of large scale neuromorphic circuits. To allow flexibility when optimizing the NVM device stack, a preferred approach is to integrate RRAM, or any other emerging NVM devices, in the back-end-of-the-line (BEOL) of a CMOS process. OxRAMs allow straightforward integration with a CMOS process flow as the materials used are compatible with CMOS fabrication. This is due to the fact that HfO<sub>x</sub> and similar materials are used as high-k gate oxide in scaled CMOS process nodes. Recently, few foundries have integrated OxRAM and CBRAMs into their CMOS fabrication flows<sup>20,69,114,115</sup>. Fig. 6 shows a cross-section representation of the OxRAM integrated in the BEOL of a 4-metal CMOS process<sup>100</sup>. 1T1R RRAM arrays from 16Kbits to 2Mbits have been fabricated using this platform with memory size scaled down to 40nm<sup>67,96</sup>. CBRAM with Ag filament is challenging as silver is a contaminant for silicon-based devices. Thus, copper (Cu) based CBRAMs have been in focus as Cu is already used in the CMOS BEOL flow<sup>53</sup>.

### E. Challenges with RRAM Devices for Neuromorphic Computing

Hybrid CMOS-NVM circuits have been proposed to achieve dense integration of CMOS neurons and emerging memory for neuromorphic computing ICs<sup>116–121</sup>. However, RRAM synapses incur several limitations which are summarized as follows<sup>23</sup>:

1. **Variability:** Individual RRAM characteristics including the switching threshold voltages and intermediate resistances vary from cell-to-cell and cycle-to-cycle depending upon the resistance state and the compliance current<sup>53,66,82,95,112</sup>. Impact of variability on STDP updates was experimentally studied on a crossbar array<sup>122</sup>. The cell-to-cell

variation of switching thresholds is a major challenge that needs to be addressed by careful program/erase pulse waveform design<sup>122</sup>. In an 1T1R cell, transistor variations also contribute to variations in the overall synapse. The RRAM resistance also exhibits temperature-dependent switching with a decrease in  $R_{HRS}$  due to the increase in leakage current with temperature<sup>107,123</sup>.

2. **Low-resistance:** Typical RRAMs exhibit 1-10k $\Omega$  LRS resistance. A large synapse fan-out incurs significant static power consumption in the CMOS driver circuits. A 1T1R synapse can mitigate this by employing a source-degenerated transconductor<sup>124</sup> or cascode topology<sup>124,125</sup>.
3. **Resolution, Retention and State Drift:** Even though RRAM can realize multi-level cell (MLC) storage<sup>66</sup>, the resistance-drift over time presents challenges in their use as stable analog synapses over a long duration. Furthermore, the resistance drift has strong dependence on the chip temperature which can make it difficult to estimate degradation of the synaptic states beyond acceptable performance limits.
4. **Endurance:** Continual on-chip training is restricted by the maximum number of write cycles allowed before the devices wear out. Also, MLC devices with shorter retention time will lead to neural network classification performance degradation if states are not periodically restored.

#### IV. NEUROMORPHIC INTEGRATED CIRCUITS

As seen earlier in Fig. 1, RRAM based neural networks are built around 1R or 1T1R crosspoint arrays driven by pre-neurons. The post-neurons sum the weighted currents and produce spike patterns. The neurons are located on the periphery where pre-neurons drive the rows (or the wordlines) and post-neurons on the columns (or the bitlines) accumulate weighted synaptic currents and generate output spikes.

##### A. CMOS-RRAM Synapses

An 1R synapse is simply a crosspoint synapse without a select transistor which is programmed to a resistance state of  $R_M$ <sup>126,127</sup>. The LIF neuron provides a virtual ground where the input spike voltages are weighted by the conductance,  $G_M = \frac{1}{R_M}$ , of the synapse and then integrated on the neuron. On the other hand, a transconductor synapse is realized with an 1T1R array where the select transistor is biased appropriately to set the synapse transconductance<sup>124,128</sup>. Fig. 7 shows a 1T1R synapse circuit which employs the transistor as a source-degenerated transconductor ( $G_{syn}$ ) and converts the pre-synaptic spikes to synaptic currents ( $i_{spk}$ ). Here, the pre-neuron output drives the input capacitance of the tran-

sistor ( $\approx 5\text{-}10\text{fF}$  in 65nm CMOS). The bottom (inert) electrode is connected to  $V_{post}$  terminal, which stays at ground during integration and switches to  $V_{rst}$  when the post-neuron fires.

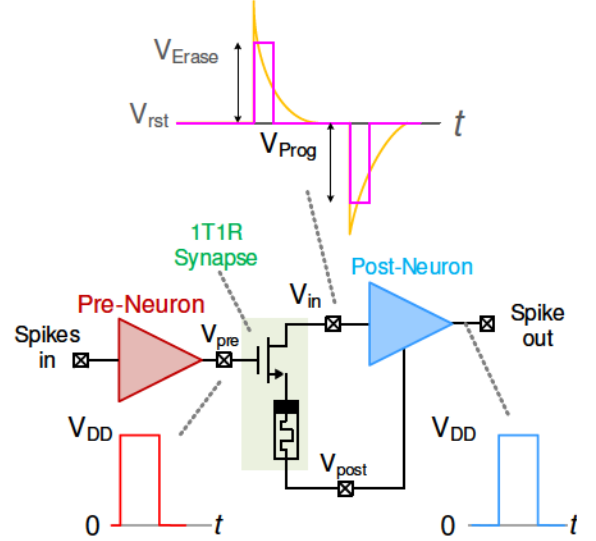


FIG. 7. A transconductance-type 1T1R synapse with in-situ weight updates as described in Ref.<sup>124</sup>.

The synaptic current in Fig. 7 is expressed as

$$i_{spk}(t) \approx G_{syn} \cdot v_{pre}(t) \quad (5)$$

where the transconductance of the overall synapse is given by

$$G_{syn} = \frac{g_m G_M}{g_m + G_M} \approx g_m || G_M \quad (6)$$

Here, the transistor is biased in saturation with small-signal transconductance,  $g_m$ , and output resistance,  $r_o$ .

The relationship between  $i_{spk}$  and  $v_{pre}$  for the 1T1R synapse is nonlinear. This nonlinearity shown in Fig. 8 and depends upon the RRAM state and also vary with process and temperature. The resulting synapse transconductance,  $G_{syn}$ , is also not constant and dependent upon  $v_{pre}$ . By employing binary spikes, the I-V nonlinearity is mitigated as a straight line can be fitted across the two points on the I-V characteristics with a constant slope. Also, the spacing between various synaptic states can be made linear by controlled programming of the MLC states, or by accommodating the resulting nonlinearity at the algorithmic level.

The synapse allows weight updates based on the correlation of pre-synaptic and post-synaptic spikes and error feedback ( $\delta$ ) as illustrated in Fig. 9. Here, voltage waveform engineering is performed to translate the time delay between the pre- and post-spikes ( $\Delta t$ ) into the application of program and erase pulse across the RRAM device to increase or decrease the weight, i.e. the



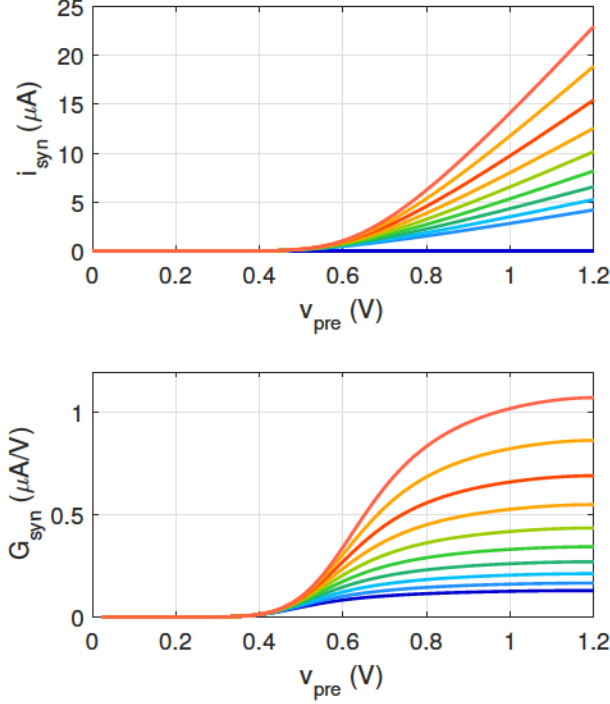


FIG. 8. (left) Simulated DC  $i_{syn}$  versus  $v_{pre}$  characteristics for a 1T1R synapse with 65nm CMOS and an RRAM model. The RRAM is initialized in the model in ten distinct states. The select transistor is sized with  $\frac{W}{L} = \frac{600nm}{60nm}$  with its drain at  $V_{rst} = 0.6V$  and source at  $0V$ . (right) Synapse transconductance,  $G_{syn}$ , as a function of input voltage.

synapse transconductance,  $\Delta w = \Delta G_{syn}$ <sup>128</sup>. The pre-neuron only observes the input capacitance of the transistor and can drive a large synaptic fanout using digital CMOS buffers (as opposed to driving several resistances in parallel as in the 1R array). The output resistance of the transconductor synapse is given by

$$R_{out} \approx g_m \tau_o \cdot R_M \quad (7)$$

which significantly reduces the loading at the input terminal,  $V_{in}$ , of the post-neuron in the firing phase. When pre and post pulses overlap, the  $V_{post}$  node is loaded by

$$R_{post} \approx R_M + \tau_o \quad (8)$$

The select transistor should be sized such that during inference, the voltage dropped across  $R_M$  shouldn't exceed the threshold,  $V_{th}^+$  to avoid disturbing the synaptic weight. Also, during the weight-update phase, a large fraction of the program/erase voltage should drop across  $R_M$  affect LTP/LTD. These constraints need to work in concert with the waveform design seen in Figure 9. Significant post-neuron loading and the associated energy consumption is avoided by ensuring sparse overlap of pre and post spikes at the algorithm level.

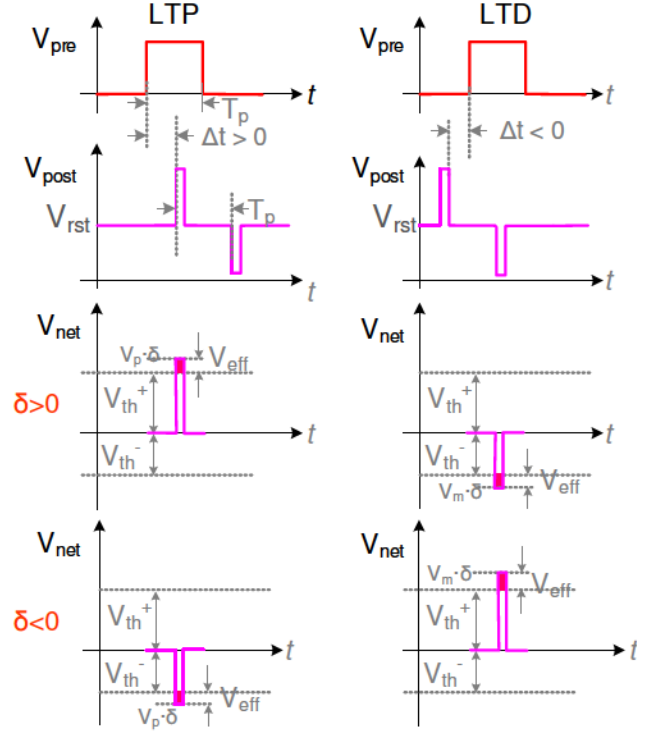


FIG. 9. An example of waveform engineering in the 1T1R synapse to translate the temporal correlation of spikes into synapse weight update as described in Refs.<sup>124, 128</sup>. Rectangular as well as exponentially decaying pulse can be used to customize the STDP learning function, although only the former is shown here.

The pre and post waveforms are designed such that during LTP, a positive voltage pulse greater than the  $V_{th}^+$  appears across the device. Conversely, during LTD, a voltage pulse more negative than  $V_{th}^-$  is applied across the device. An error feedback signal,  $\delta$ , can be employed to modify the post-waveform by changing its sign and amplitude in order to implement a desired three-factor learning rule as in Eq. 3. In addition to waveform engineering with RRAMs, novel devices have been engineered with bio-realistic internal dynamics to support STDP-type updates. These include Ag-in-oxide memristor (Ag/SiO<sub>x</sub>Ny/Pt) with diffusive dynamics<sup>129</sup> and second-order oxygen vacancy dynamics in a Pt/WO<sub>x</sub>/W device<sup>130</sup>.

## B. Array Programming Circuits

The immediate application of NVM-based NeuSoCs is in transfer learning for inference applications. Also, on-chip learning requires suitable initialization of synaptic weights. Moreover, it may be desirable to read out learned weights from a particular chip. Thus, RRAM array program and read circuits become a necessary feature (in addition to the neurons).

For SNNs with unsigned weights a 1T1R synapse is employed. However, for VMMs and higher accuracy DNNs, signed weights become necessary. These are realized using a differential 2T2R topology where the conductances of two devices are subtracted in the analog domain<sup>131–133</sup>

$$G_{syn,ij} = G_{ij}^+ - G_{ij}^- \quad (9)$$

Here,  $G_{ij}^\pm$  are the conductances of positive and negative branches.

In order to program multi-level states in 1T1R (or 2T2R) arrays, an *Incremental Step Pulse Program* (ISPP) scheme is employed, which was adapted from earlier MLC Flash memory designs<sup>134</sup>. ISPP for RRAM arrays has been demonstrated in the recent literature using a semiconductor parameter analyzer (SPA). In these experiments, the voltage pulse profiles are set after careful characterization of cell-to-cell and cycle-to-cycle variability of the devices<sup>135–137</sup>.

In an ISPP scheme, the program (or erase) voltage is applied to the device, while setting a compliance current,  $I_{cc}$ , and its state is verified using a read operation. The program (erase) pulse magnitude is gradually increased and the state is verified till a target synapse conductance state is achieved<sup>135,136</sup>. Alternatively, pulse width can also be progressively increased (or decreased) to set the desired conductance state<sup>138</sup>, however, modification of the program pulse magnitude was found to be more effective in terms of device variability and program/erase time<sup>136</sup>.

On-chip ISPP implementations include 1D arrays of row DACs and column ADCs with addressing logic and an on-chip RISC processor to isolate a particular MLC cell and then program it into a desired state<sup>139</sup>. Here, the DACs drive the wordlines with analog inputs,  $x_i$ , and the each of the columns digitize the weighted sums,  $y_j = \sum_i w_{ij} x_i$ . Backend processing is handled in the digital domain using the processor.

### C. CMOS Neuron Circuit Design

There is a significant body of work over the past several decades on low-power spiking neurons and synapses using subthreshold CMOS circuits, which is documented in<sup>21</sup> and references therein. These circuits operated in the range of kHz spike rate, similar to biology, and were optimized for ultra low power emulation of neurophysiological behavior. However, with the advent of NVM synapses, a renewed focus is on in-memory neuromorphic computing for deep neural networks. Recently, several spiking neurons designs have been proposed for neuromorphic computing<sup>36,40,140,142,145</sup>, however the challenges associated with their circuit integration with practical RRAM devices weren't considered.

In the context of NVM-based neuromorphic computing, CMOS neuron design can vary widely depending upon the desired functionality. For inference-only NeuSoC architectures, the neuron needs to integrate synaptic

currents, make decisions when the membrane potential crosses the threshold (i.e. fire), and generate post spikes. On the other hand, for in-situ learning functionality, the neuron has to generate the waveform seen in Figs. 7 and 9 and drive the RRAMs (Eq. 7). Moreover, the neuron should also allow for a three-factor learning rule that incorporates the feedback error ( $\delta$ ).

Ultra-low power neurons with  $\approx 100$ nW static power consumption can be designed for inference-only neuromorphic ICs<sup>146</sup>. However, neurons designs that interface with RRAMs need sufficient drive capability and additional circuitry for enabling in-situ learning<sup>36,119,147</sup>.

Opamp-based neuron designs were introduced in<sup>118,141</sup> and designed to drive a large RRAM fan-out. Event-driven LIF neurons were demonstrated in 180nm CMOS which drove 1R RRAM crosspoint arrays with in situ STDP-based learning as illustrated in Fig. 10<sup>23,41,148,149</sup>. Here, the CBRAMs were integrated using wire bonding with the CMOS neurons and a neural circuit with associative learning (i.e. a Pavlov's Dog experiment) was demonstrated. In order to accommodate a wide variety of material stack used in the CBRAM, the neuron spike voltage and temporal profile parameters were designed to be digitally programmable. This neuron design had a bias current of  $13\mu A$  in the integration mode and dynamically-biased with  $56\mu A$  in the firing mode. The neuron's class-AB output stage could source up to  $1.4$ mA current to drive an equivalent RRAM load of  $10\Omega$  with a power efficiency of  $97\%$ <sup>41,148,149</sup>.

1T1R synapses not only significantly relax current sourcing requirement for the neurons, but also enable digital CMOS drive (i.e. with full-swing output between 0 and  $V_{DD}$ ) as seen in Fig. 7. Fig. 11 shows an event-driven neuron that is adapted for the 1T1R crosspoint arrays<sup>124</sup>. The neuron operates asynchronously in two event-driven modes with a shared opamp— *integration* and *firing*. The neuron nominally operates in the integration mode, where the opamp is biased with very low bias current and configured as a leaky integrator with the virtual ground at the rest potential  $V_{rst} = V_{DD}/2$ ,  $V_{DD}$  being the supply voltage. The incoming weighted current spikes,  $i_i = \sum_j w_{ij} \cdot v_{pre,j}(t)$ , are integrated on the capacitor  $C_m$  resulting in the membrane voltage

$$V_{m,j}(t) = V_{rst} + \sum_i \frac{G_{m,ij}}{C_{m,j}} \int_0^t v_{pre,i}(t) \otimes h(t) \cdot dt \quad (10)$$

Here,  $h(t) = e^{-\frac{t}{\tau_{lk}}} u(t)$  is the impulse response of the LIF circuit during integration mode where  $\tau_{lk} = R_{lk} C_{m,j}$  is the leak time-constant;  $R_{lk}$  is realized using the MOSFET,  $M_{lk}$ , biased in triode.

An asynchronous comparator is used to compare the membrane potential,  $V_{m,j}$ , with the threshold voltage,  $V_{thr}$ . When a positive crossing occurs, the neuron switched into the fire mode where its reconfigured as a voltage follower/buffer. Concurrently, an output spike with full digital levels is created and propagated forward.



TABLE II. Performance comparison of recent CMOS integrate and fire neurons.

Design	Architecture	Technology	Synapse Type	Spike Rate	In-situ Learning	Power	Energy-efficiency (/spike/synapse)	Area ( $\mu m^2$ )
indiveri et. al. 2006 <sup>7</sup>	Subthreshold	0.35 $\mu m$	CMOS, Bistable	200Hz	✓	-	90pJ	2,573
Basu and Hasler 2010 <sup>7</sup>	Subthreshold	0.35 $\mu m$	Floating Gate	100Hz	✗	1.74nW	17.4pJ	2,740
Cruz-Albrecht et. al. 2012 <sup>140</sup>	Capacitive	90nm	CMOS, Dynamic	100Hz	✓	40pW	400fJ	442
Wu et. al. 2015 <sup>141</sup>	Opamp	0.18 $\mu m$	1R array <sup>†</sup>	0.1-1MHz	✓	23.4 $\mu W$ 95.4 $\mu W$ <sup>‡</sup>	9.3pJ	12,100
Sahoo 2017 <sup>142</sup>	Ring VCO	65nm	None	0.4-1.5MHz	✗	-	-	-
Larras et. al. 2017 <sup>143</sup>	Current summing	65nm	Digital, Binary	-	✗	-	7fJ	41,820
Sourikopoulos et. al. 2017 <sup>144</sup>	Subthreshold	65nm	None	25kHz	✗	100pW	4fJ	35
Saxena 2020 <sup>124</sup>	Opamp	0.18 $\mu m$	1T1R array	0.1-1MHz	✓	16.2 $\mu W$ 64.8 $\mu W$ <sup>‡</sup>	8.1fJ $\Delta$ 40fJ	5,625

<sup>†</sup> An LRS resistance of  $R_{LRS} = 1k\Omega$  is assumed.

<sup>‡</sup> Opamp is dynamically biased with higher current during the firing mode.

$\Delta$  Energy-efficiency when used in inference mode only.

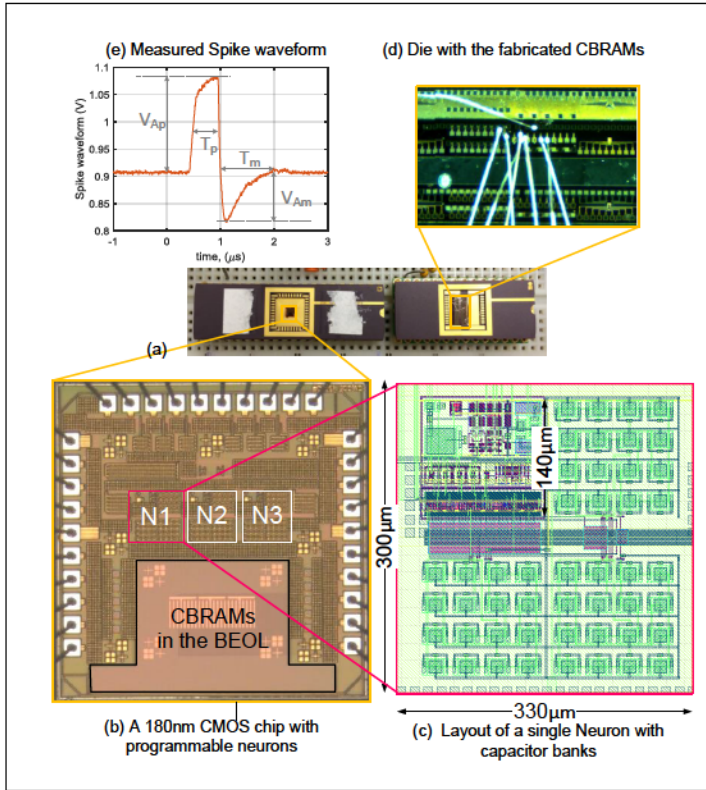


FIG. 10. (a) An opamp-based event-driven integrate and fire neuron to drive 1R CBRAM synapses as described in Refs.<sup>41,141,149</sup>, (b) a 180nm CMOS chip with LIF neurons, (c) Layout of a single neuron with capacitor banks, (d) CBRAM devices as described in Ref.<sup>93</sup>, (e) An output spike generated by the neuron.

If on-chip learning is enabled then the waveform seen in Fig. 9 is generated by a switched-capacitor circuit. The phase control circuit generates the strobes  $\phi_{int}$  and  $\phi_{fire}$  for the switched-capacitor circuits, and  $\phi_1$  and  $\phi_2$  for waveform generation. The waveform parameters can be digitally configured for a specific RRAM material stack and can be changed on the fly. The waveform drives the positive opamp input while the opamp is configured as a voltage follower. The voltage follower provides the suffi-

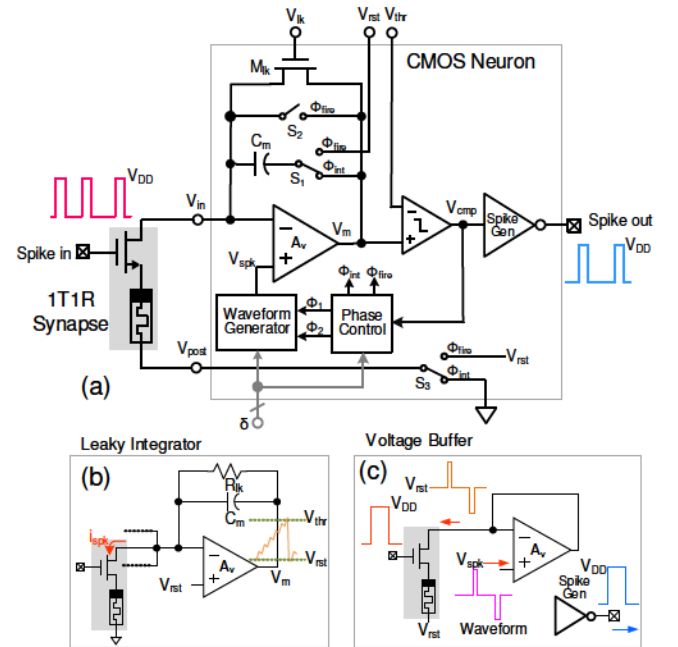


FIG. 11. (a) Schematic of the event-driven CMOS neuron for 1T1R array architecture as described in Ref.<sup>124</sup>, (b) Neuron in the integration mode, (c) Neuron in the firing mode where a weight update can take place.

cient drive current for the RRAMs that are in the STDP window where pre- and post-spikes overlap<sup>124</sup>. After the fire phase is concluded, the neuron enters a refractory period, where it is reconfigured back in the integration mode and  $V_{m,j}$  is reset to  $V_{rst}$ .

The CMOS neuron from Fig. 11 was designed 180nm CMOS process with a  $V_{DD} = 1.8V$ <sup>124</sup>. A compact model for the RRAM<sup>150</sup> was implemented in Verilog-A. The total bias current was  $9\mu A$  and and energy-efficiency of 40fJ/spike/synapse. Compared to prior neurons that were designed to drive 1R synapses<sup>141</sup>, this 1T1R design only drives a the small input capacitance of the access transistor during inference, and resistive load only dur-

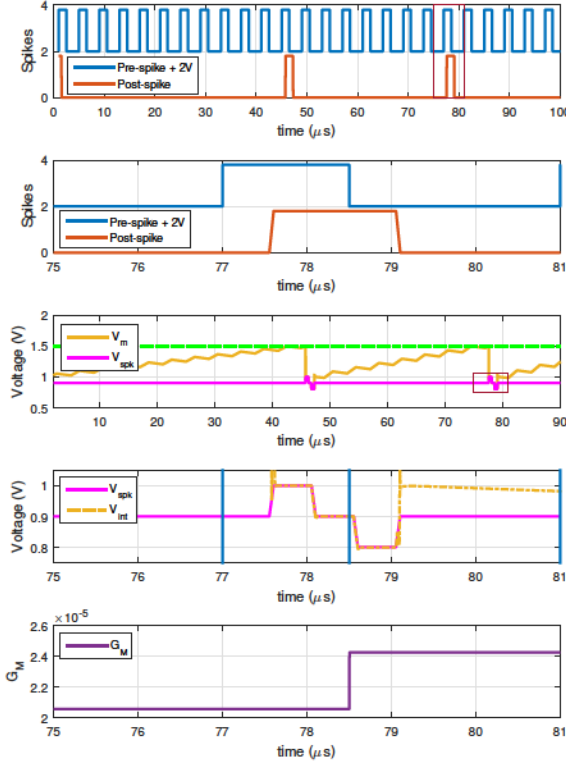


FIG. 12. Plots generated from a transient simulation of the 180nm CMOS neuron for driving 1T1R synapses seen in Figure 11. The asynchronous input spikes are integrated into the membrane potential,  $V_m$ . When  $V_m$  crosses  $V_{thr}$ , the neuron fires and a full-swing output spike is generated. Also, a spike waveform,  $V_{spk}$ , is generated and applied across the 1T1R synapse, and thus the RRAM device, by the opamp which is configured as a voltage follower. In this case, the conductance of the synapse is increased as the positive portion of the spike waveform overlaps with the pre-spike pulse.

ing the training/learning phases.

A performance comparison of recent integrate and fire neurons from the literature is provided in Table II. As evidenced by this comparison, very few architectures have addressed the circuit interfacing challenges with the RRAM devices. There is significant room for reduction in neuron energy consumption by employing scaled CMOS nodes, and techniques including adaptive biasing (with static bias currents in the nAs), relaxing the transition times on the program/erase pulses and power gating. However, the designs will continue to co-evolve with spike-based learning algorithms.

## V. NEUROMORPHIC DEEP LEARNING ALGORITHMS

Realization of DNNs in a low-power hardware has been of great interest in applications where neural network processing needs to occur close to the sensors with real-time inference and with minimal reliance on the Cloud infrastructure. This is driven by two divergent paradigms. The first view is driven by the need to demonstrate competitive performance on real-world applications compared to Backprop-based standard DNNs. The second view is of exploring the underlying mechanism behind distributed and continual learning in biological brains, which can learn from small amount of data in an unsupervised or semi-supervised manner. A brief overview of evolving neuromorphic learning algorithms is presented in this section and their performance is compared in Table III.

### A. Transfer Learning

A straightforward approach is to first train a DNN comprised of differentiable non-spiking neurons using Backprop, regularization methods such as Dropout, and optimization techniques including batch normalization<sup>3</sup>. Then the DNN is converted to its equivalent SNN by employing rate based coding, scaling the weights based on the spiking activity, and adjusting the thresholds to minimize absence of spikes or neuron saturation due to over-spiking<sup>156,171</sup>. There is a perceptible performance drop in classification accuracy in the DNN to SNN conversion<sup>156–158,171</sup>. Transfer learning has also been used to initialize a deep SNN which is then further trained using Backprop<sup>172</sup>. To facilitate transfer learning, a neuromorphic IC should have a mechanism to download a pre-trained model and program the on-chip synapses to the desired weights. Furthermore, these weights decay over time due to resistance drift and a *drift compensation* scheme is needed to maintain model accuracy over longer time scales<sup>173</sup>. Transfer learning has achieved highest accuracy of 99.44%<sup>158</sup> for the MNIST handwritten digits dataset compared to the best in-class DNN accuracy of 99.79%<sup>151</sup>. For CIFAR-10 dataset, this gap is within 2% of the Alexnet DNN.

### B. Semi-supervised Learning

Spike-based unsupervised or semi-supervised learning algorithms are based on the neural-inspired WTA motifs of LIF neurons with inhibition and competition. Their ability to learn spatiotemporal patterns using STDP was first studied using software simulations<sup>174</sup> and was then applied to vision tasks<sup>31,175</sup>. Subsequently, analytical modeling showed that WTA neurons with STDP learning realize a spiking version of the Expectation Maximization algorithm<sup>34</sup>. This two layer motif learns from a small



TABLE III. Comparison of deep SNN Algorithms and their benchmarking against the highest performing non-spiking DNNs.

Reference	Algorithm/Neuron Type	Network	Dataset	Accuracy
Standard Deep Neural Networks				
Simard et. al. 2003 <sup>151</sup>	Backprop with ReLU	784-20c5-2s-100-10o	MNIST	99.6%
Wan et. al. 2013 <sup>152</sup>	Backprop with ReLU and DropConnect	-	MNIST	99.79%
Krizhevsky et. al. 2012 <sup>153</sup>	Alexnet: Backprop with ReLU	5-layer ConvNet	CIFAR-10	88.91%
He et. al. 2016 <sup>154</sup>	ResNet: Backprop with ReLU	ResNet-1001	CIFAR-10	95.08%
VGGNet 2017 <sup>155</sup>	Backprop with ReLU	16-layer ConvNet	ImageNet <sup>∇</sup>	93.2%
SNNs with Transfer Learning				
Diehl et. al. 2015 <sup>156</sup>	Fully connected with Rate-based IFN	784-1200-1200-10o	MNIST	98.68%
	ConvNet with Rate-based IFN	784-12c5-2s-64c5-2s-10o		99.14%
			784-16c5-2s-64c5-2s-10o	MNIST
Hunsberger et. al. 2016 <sup>157</sup>	ConvNet with Rate-based LIFN	Alexnet	CIFAR-10	82.95%
		-	ImageNet	76.2%
		LeNet	MNIST	99.44%
Rueckauer 2017 <sup>158</sup>	ConvNet with Rate-based LIFN	Alexnet	CIFAR-10	88.82%
		VGG-16	ImageNet <sup>∇</sup>	84.86%
		Inception-V3	ImageNet <sup>∇</sup>	92.04%
Sengupta et. al. 2019 <sup>159</sup>	ResNet with Rate-based LIFN	ResNet-20	CIFAR-10	87.46%
		ResNet-34	ImageNet	86.43%
Unsupervised Learning SNNs				
Diehl & Cook 2015 <sup>160</sup>	STDP-WTA, LIFN with inhibition	784-1600-10o	MNIST	95%
Kheradpisheh et. al. 2016 <sup>161</sup>	STDP-WTA, IFN with latency coding	784-30c5-2s-100c5-2s-10o	MNIST	98.4%
Vaila et. al. 2019 <sup>162</sup>	Binarized STDP-WTA, Surrogate gradients	784c5-2s-30c5-s7-500-70-10o	MNIST	98.49%
			EMNIST	85.3%
Spike-based Backpropagation				
Neftci et. al. 2016 <sup>163</sup>	Event-driven RBP with Rate-based IFN	784-500-10o	MNIST	97.04%
Lee et. al. 2016 <sup>164</sup>	Backprop with Rate-based IFN	784-800-10o	MNIST	99.31%
O'Connor et. al. 2017 <sup>165</sup>	Backprop with Delta coding	784-200-200-10o	MNIST	98.36%
Kulkarni et. al. 2018 <sup>166</sup>	NormAD with IFN	784-12c3-10o	MNIST	98.17%
Mostafa 2018 <sup>167</sup>	Backprop with Temporal coding	784-400-400-10o	MNIST	97.55%
Shrestha et. al. 2018 <sup>168</sup>	SLAYER with Rate-based IFN	784-12c5-2s-64c5-2s-10o	MNIST	99.36%
Mostafa et. al. 2018 <sup>169</sup>	DNN with Synthetic gradients, ReLU <sup>†</sup>	784-1000-1000-1000-10o	MNIST	98.7%
		ConvNet3 <sup>△</sup>	CIFAR-10	89.1%

<sup>†</sup> This work uses non-spiking ReLUs but can be adapted to SNNs using rate-based or temporal-coding<sup>170</sup>.

<sup>△</sup> Network for CIFAR-10 is 32x32x3-96c5-s3-128-c5-s3-256c5-s3-2048-2048-10o.

<sup>∇</sup> Top-5 accuracy on ImageNet dataset.

number of samples and the SNN weights converge to the log probability of the patterns in the training dataset<sup>34</sup>.

The above mentioned two-layer SNN, without a hidden layer, was applied to the MNIST handwritten digits recognition task<sup>41,160</sup>. Using competitive learning, each neuron learned to fire on distinct inputs. The output labels are either assigned after the learning is completed<sup>160</sup> or only the intended output neurons are allowed to fire<sup>41</sup>. This semi-supervised SNN achieved a classification accuracy of 94% for four digits and 83% on all the ten digits (with 10 output neurons) with around 1000 training samples for each image label<sup>41</sup>. Higher classification accuracy was demonstrated by using a large number of competing neurons ( $\approx 5000$ ) leading to a maximum classification accuracy of 95%<sup>160</sup>.

The WTA motifs can be organized in a ConvNet architecture as shown in Fig. 13 where only one neuron per kernel is allowed to fire across all the feature maps<sup>161,177</sup>. Stacking of these spiking ConvNet motifs to improve classification performance was explored next<sup>161,162,177,178</sup>. While the ConvNet layers learn in an unsupervised manner, a fully-connected read-out layer which is trained using Backprop is employed as seen in Fig. 13. It was found that greedy stacking of more than two unsupervised learning ConvNet layers didn't improve the classification accuracy. In essence, each WTA with STDP layer can be thought of as if its performing unsupervised clustering over the input feature space. This

learning occurs with very few samples. However, in the absence of a mechanism to assign credit across layers based on the output classification error, the resulting accuracy doesn't improve by increasing the SNN depth. These deep semi-supervised SNNs have demonstrated a maximum accuracy of 98.5% for the MNIST handwritten digit dataset<sup>161,162,177-179</sup>, and have been shown to be suitable for incremental learning of tasks<sup>180</sup>.

### C. Backpropagation-based Learning

In spite of the desirable features of the semi-supervised learning SNNs, there is a classification accuracy gap between semi-supervised SNNs and Backprop-trained DNNs. As a result, there has been sustained interest in adapting Backprop to SNNs<sup>163</sup>. Backprop, along with the ConvNet layers, is the workhorse for deep learning and minimizes the output classification error by propagating error gradients backward from the output layer,  $L$ , to the lower network layers,  $1 \leq l < L$ . This is analytically described by the four Backprop equations below<sup>3</sup>

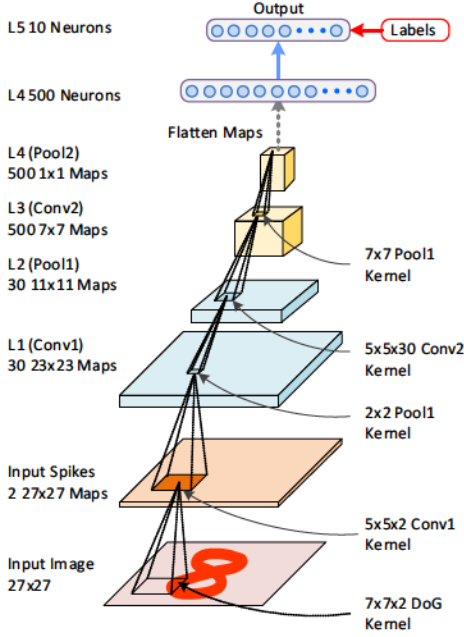


FIG. 13. A Spike-based ConvNet for MNIST handwritten digits dataset:  $27 \times 27 \times 2$ -30c5-2s-500c7-7s-10o as described in Ref.<sup>176</sup>. Edge detection is performed using On-center and Off-center Difference of Gaussian (DoG) kernels to result in a  $27 \times 27 \times 2$  image.

$$\delta^L = \nabla_a \mathcal{C} \odot \sigma'(z^L) \quad (11a)$$

$$\delta^l = ((W^{l+1})^T \delta^{l+1}) \odot \sigma'(z^l) \quad (11b)$$

$$\frac{\partial \mathcal{C}}{\partial w_{jk}^l} = a_k^{l-1} \delta_j^l \quad (11c)$$

$$\frac{\partial \mathcal{C}}{\partial \theta_j^l} = \delta_j^l \quad (11d)$$

Here,  $\mathcal{C}$  is a formulation of the output loss, or cost, function and given by<sup>3</sup>

$$\mathcal{C} = \frac{1}{2N} \sum_x \|y(x) - a^L(x)\|^2 \quad (12)$$

where  $N$  is the batch size, and  $x$  and  $y$  represent a sample and the corresponding label in the dataset respectively. Furthermore,  $a_j^l$  represents the activation for the  $i^{th}$  neuron in the  $l^{th}$  layer ( $1 \leq l \leq L$ ),  $\sigma(\cdot)$  is the equivalent activation function,  $\delta_j^l$  is the backpropagated error for neuron  $j$  in layer  $l$ , and  $\theta_j^l \equiv V_{rst}$  is the spike threshold for the neuron. Also,  $z_j^l = \sum_k w_{jk}^l \nu_k^{l-1}$  is analogous to the neuron membrane potential,  $V_{m,j}$ . SNNs can either be formulated using rate-based or spike latency coding<sup>167,181</sup>.

For every training batch, each of the network weights

are updated using the gradients as<sup>163</sup>

$$w_k^{(t+1)} = w_k^{(t)} - \eta \frac{\partial \mathcal{C}}{\partial w_k} \quad (13a)$$

$$\theta_j^{(t+1)} = \theta_j^{(t)} - \eta \frac{\partial \mathcal{C}}{\partial \theta_j} \quad (13b)$$

where  $\eta$  is the learning rate. Since the activation of spiking neurons is discontinuous, direct computation of gradients is not feasible. As a result, equivalent linearized or stochastic differentiable neuron models have been derived to obtain a continuous activation function,  $\sigma(\cdot)$ , with a well defined derivative,  $\sigma'(\cdot)$ <sup>163,168</sup>.

## 1. Challenges with Spike-based Backprop

Since Backprop has evolved using von Neumann architectures, it assumes that the network-wide weights, neuron activations, and their derivatives are always accessible from a high-density memory. This memory has a latency and energy cost associated with data access which percolates through the von Neumann bottlenecks<sup>163,170</sup>. Neuromorphic computers, on the other hand, aim to minimize this back and forth shuttling of data by performing localized computing inside the memory itself; however, the nature of data flow associated with Backprop presents the following architectural challenges<sup>23,163</sup>:

- 1. Weight Transport Problem:** In order to compute the weight updates at layer  $l$  using Eq. 11b, the transpose of the weight matrix for layer  $(l+1)$ ,  $(W^{l+1})^T$ , must be available while evaluating weights connecting layers  $(l-1)$  and  $l$ , which poses challenges for hardware design.
- 2. Non-concurrence:** Data flow must alternate between forward and backward passes during each minibatch which limits learning on real-time streaming data.
- 3. Differentiability and Precision:** Derivatives need to be computed with high-precision or approximated to simpler functions.
- 4. Temporal Credit Assignment Problem:** During training of a DNN as shown in Fig. 13, the network layers undergo a forward pass and then wait for the gradients to be propagated in the reverse direction. This poses a temporal credit assignment problem where future errors are needed to update weights based on current spike correlations. Consequently, the lower layers in the network are frozen till the weight update in the backward pass takes place<sup>182</sup>.

## 2. Recent Advances in Spike-based Backprop

Several techniques have appeared in spike-based Backprop adaptations where the constraints of standard Back-



prop are relaxed to simplify neuromorphic hardware design. Random Backprop (RBP) or feedback alignment (FA) algorithm eliminates the symmetry constraint<sup>183</sup>. Here, instead of  $(W^{l+1})^T$  in the weight update computation in Eq. 11c, a fixed random matrix,  $B^{l+1}$  is used. This replaces Eq. 11c by<sup>183</sup>

$$\delta_{FA}^l = (B^{l+1} \delta^{l+1}) \odot \sigma'(z^l) \quad (14)$$

Switching to the fixed matrix simplifies the symmetry constraint in the architecture without incurring significant decrease in classification accuracy<sup>183</sup>. Random Backprop was adapted to event-driven neurons in an early demonstration of deep learning in SNNs<sup>163</sup>.

In the Backprop Eqs. 11a to 11d seen earlier, a continuous, and thus differentiable, model for the neuron activation function,  $\sigma(\cdot)$ , was assumed. Recent results in deep SNNs have developed continuous models for the LIF neurons by considering signal noise to soften the hard decision thresholds. In one of these differential neuron models, the derivative of spiking activation was approximated as:

$$\sigma'(z_j^l) = \alpha e^{-\beta|z_j^l - \theta|} \quad (15)$$

where,  $\alpha$  and  $\beta$  are model constants, and  $\theta$  is the neuron threshold<sup>168</sup>.

By combining RBP with event-driven spikes and approximate rate-based gradients<sup>163</sup>, an accuracy of 97.4% for the MNIST dataset was obtained with a three-layer fully-connected network. In another rate-coding SNN, the WTA structure was preserved in the weight update rules and resulted in an MNIST accuracy of 98.71%<sup>164</sup>. In latency or time-coded SNNs, SpikeProp<sup>181</sup> was an early work where the post neuron spike times were linearized to compute gradients. Another learning rule called *normalized approximate gradient descent (NormAD)* was proposed for temporally-coded SNNs, where an approximate gradient was defined by linearizing the membrane potential as a function of input spikes<sup>166,184</sup>. In a more recent work, the explicit spike delays were employed to compute gradients without any linearization<sup>167</sup>. This work resulted in a reduced MNIST classification accuracy at 97.55%.

A recent SNN training algorithm, called SLAYER, adapted spike-based Backprop where learning occurs in both the weights as well as the axonal delays using the stochastic exponential neuron model seen in Eq. 15<sup>168</sup>. In this work, a temporal credit assignment scheme is employed where the error is backpropagated through time. The algorithm was implemented in a GPU and achieved the highest SNN classification accuracy for MNIST dataset at 99.44%<sup>168</sup>. The requirement of propagating error back in time was solved by the SuperSpike<sup>185</sup> algorithm. SuperSpike employs a three-factor learning rule with synaptic eligibility traces to solve the temporal credit assignment problem<sup>170</sup>.

From the comparison in Table III, it can be seen that transfer learning achieves the highest accuracy for SNNs

which comes very close to the classification performance for the AlexNet-size DNNs. Steady progress has been made in adapting Backprop to SNNs with the recently reported algorithms demonstrating competitive performance on the MNIST and CIFAR-10 datasets while relaxing the hardware bottlenecks associated with the standard Backprop.

## VI. NVM-BASED NEUROMORPHIC ICS

Early analog neuromorphic ICs comprised of low-power subthreshold analog neuron and synapse circuits and were primarily intended for emulating ion channel kinetics in biological neural circuits<sup>21,46</sup>. Mixed-signal ICs employed SRAM with DACs to realize synapses along with analog LIF neurons<sup>49,191</sup>. These were scaled to wafer-level SNN implementations and demonstrated several neurobiological as well as neuromorphic computing tasks<sup>192,193</sup>.

Progress in digital neuromorphic hardware platforms has led to the realization asynchronous event-driven (as opposed to clock driven) computing ICs that communicate information on and across the chips using digital spikes. The most pertinent examples of digital neuromorphic chips are IBM's TrueNorth chip<sup>44</sup>, the recent Loihi chip from Intel<sup>45</sup>, and the two versions of SpiN-Naker systems from the European Brain Project<sup>194-196</sup>.

Table IV presents NVM-based Neuromorphic ICs in the literature along with their performance benchmarking. Development of in-memory computing neuromorphic ICs initially focused on small-scale NVM device arrays for characterization of device switching, multi-level or analog states, variability, retention and endurance<sup>41,117,137,197-199</sup>. In several of these works, device results were extrapolated to DNNs or SNNs which were entirely simulated in software.

### A. NOR Flash Architectures

Flash-based neuromorphic architectures have been studied for the past several years<sup>18</sup>. However, the recent interest in NVM-based VMMs led to their first hardware demonstration using established NOR Flash arrays integrated in a standard CMOS technology, along with the array programming and read circuitry<sup>187</sup>. In this work, a three-layer analog neural network was demonstrated using 180nm NOR Flash array<sup>187</sup>. The chip was programmed using a transfer learning approach with 6-bit analog precision and demonstrated an MNIST accuracy of 94.7%. An important observation was that the cell conductances decayed by around 13% over the 7 months storage period, however, the classification accuracy had minimal degradation and it remained above 94%.

TABLE IV. Comparison of NVM-based Neuromorphic ICs for Deep Learning.

Design	Technology	Architecture	Synapse Resolution	On-chip learning	Network Size	Dataset, Classf. Acc.	Energy-efficiency
Yu 2016 <sup>186</sup>	130nm CMOS + TaO <sub>x</sub> /HfO <sub>2</sub> OxRAM	Binary VMM	1-bit	✗	400-200-10o	MNIST, 96.5%	-
Guo 2017 <sup>187</sup>	180nm NOR Flash	Analog VMM	6-bit Analog	✗	784-256c5-10o	MNIST, 94.7% CIFAR-10, 84.8% <sup>†</sup>	20nJ/inference
Li 2018 <sup>188</sup>	2μm CMOS + Ta/HfO <sub>x</sub> OxRAM	R array	Analog	✓	64-10o	8x8 UCI, 91.7%	-
Wang 2018 <sup>189</sup>	Ta/HfO <sub>x</sub> OxRAM	1T1R array	Analog	✓	8-8o	Custom patterns	-
Cai 2019 <sup>69</sup> & Correll 2019 <sup>139</sup>	180nm CMOS + WOx RRAM	1R array, Mixed-signal	Analog	✓	54x108 array	5x5 images	8.5GOPS/W
CEA-Leti 2019 <sup>115</sup>	130nm CMOS + Ti/HfO <sub>x</sub> OxRAM	1T1R, SNN	1-bit	✗	784-10o	MNIST, 84%	180pJ/SynOp
Xue 2019 <sup>190</sup>	55nm CMOS + RRAM	VMM, 1T1R	3-bit Digital	✗	1Mb Macro	CIFAR-10, 85.52% <sup>‡</sup>	53.17TOPS/W
Hirtzlin 2020 <sup>131</sup>	130nm CMOS + Ti/HfO <sub>x</sub> OxRAM	2T2R	1-bit	✗	Off-chip	MNIST, 98.3% <sup>†</sup> CIFAR-10, 87.5% <sup>†</sup> ImageNet, 69.7% <sup>†</sup>	20-30pJ/SynOp
Nandakumar 2020 <sup>173</sup>	90nm CMOS + GST PCRAM	2T1R, LIFN	7-bit Analog	✓	784-10o	MNIST, 70%	-
Liu 2020 <sup>133</sup>	130nm CMOS + RRAM	VMM, 2T2R	1 to 8-bit Digital	✗	784-100-10o	MNIST, 94.4%	78.4TOPS/W
Wan 2020 <sup>146</sup>	130nm CMOS + TaO <sub>x</sub> RRAM	IFN, 1T1R	1-bit	✗	225-60	MNIST RBM	74TOPS/W
Xue 2020 <sup>114</sup>	22nm CMOS + RRAM	VMM, 1T1R	4-bit Digital	✗	1Mb Macro	CIFAR-10, 90.19% <sup>‡</sup> CIFAR-100, 64.15% <sup>‡</sup>	121.38TOPS/W

<sup>†</sup> Simulated result.<sup>‡</sup> Result is using off-chip computation.

## B. CMOS-RRAM Architectures

Initial demonstrations focused on very simple pattern learning tasks using small-scale RRAM memory arrays. For example, synapses were interfaced with discrete electronic circuits to demonstrate a small-scale network<sup>200</sup>. In an early work on event-driven RRAM-compatible neuron design, a 3-neuron associative SNN was experimentally demonstrated<sup>148</sup>, and then adapted to  $8 \times 8$  UCI handwritten digit dataset<sup>41</sup>.

These were followed by in-silicon demonstration of shallow two-layer neural networks. For example, online learning of binarized neural networks was demonstrated using TaO<sub>x</sub>/HfO<sub>2</sub> RRAMs in 130nm CMOS and resulted in an MNIST accuracy of 96.5%<sup>186</sup>. In another work, in-situ learning in a two-layer network was demonstrated using Ta/HfO<sub>x</sub> memristor array integrated with 2μm CMOS<sup>188</sup>. However, in this work, the neuron activations were simulated in software<sup>188</sup>. Furthermore,  $8 \times 8$  Ta/HfO<sub>x</sub>/Pd 1T1R arrays were demonstrated to learn basic patterns under unsupervised training<sup>189</sup>.

With the recent integration of RRAMs with CMOS transistors in a foundry process, neuromorphic ICs with medium to large-scale integration of circuits with RRAMs have begun to appear. A majority of these works target inference-only applications by leveraging the high density of RRAM arrays with analog-domain multiply and accumulate (MAC) operations<sup>114,115,131,133,146,186,190</sup>. Among the inference-only demonstration chips, either binarized weights were employed, or parallel RRAM cells were used to emulate a multibit synapse. In the latter case, unit-weighted binary 1T1R cells were used as a current-DAC to realize multi-bit synapses<sup>114,131,133,146,190</sup>.

As an example of large-scale integration of OxRAMs with 130nm CMOS, a 2Kb differential binary RRAM (2T2R) array with integrated column sense-amps was fabricated, and on-chip inference was exhibited for MNIST, CIFAR and ImageNet datasets with competitive classification performance<sup>131</sup>. In another recent work<sup>133</sup> using a similar technology, binarized 2T2R synapses were used for signed weights in order to demonstrate a multi-bit VMM. Again, in this work multi-bit weights were realized by combining several 2T2R cells in parallel. The weighted currents were integrated and digitized using a successive approximation register (SAR) ADC. This chip-scale demonstration fully-integrated a 784-100-10 network with an FPGA-based back-end, and demonstrated 94.4% MNIST accuracy<sup>133</sup>. Overall, a full-chip integration that demonstrates true multilevel-cell 1T1R synapses for on-chip inference is yet to be seen (other than the attempts in the works<sup>69,139</sup>).

So far, only a few designs have attempted fully-integrated on-chip learning where the challenges associated with analog synapses need to be addressed<sup>69,139</sup>. This CMOS-RRAM IC prototype incorporates  $54 \times 108$  WOx RRAM crossbar array integrated with 180nm CMOS<sup>69,139</sup>. The chip also includes arrays of 6-bit DAC and 13-bit column ADCs and a RISC processor for digital backend. The row DAC produces voltage pulses of fixed width proportional to the input value. These DAC pulses are weighted by the RRAM-based VMM and then integrated in charge-domain on the integrating-type column ADCs, thus realizing a mixed-signal VMM with digital input and output vectors.



## VII. DISCUSSION AND CONCLUSION

From the discussion in Section VI, it is evident that fully CMOS-based digital and mixed-signal neuromorphic ICs feature very high-level of system integration and functionality. This is due to the maturity of design, modeling and verification infrastructure for CMOS technology. However, CMOS realizations either exhibit lower neurosynaptic density, employ binarized volatile synapses, and can be limited by a von Neumann bottleneck.

While digital neuromorphic ICs have made significant progress in low-power realizations of deep neural networks, NVMs have a role to play in another order of magnitude improvement in neurosynaptic density and energy-efficiency. Among the NVMs, NOR Flash-based architectures have demonstrated multi-layer neural networks for inference, all due to the maturity of floating-gate devices, proven multilevel cell capability, and their longer retention times.

As seen in Table IV, RRAM and PCRAM-based neuromorphic inference ICs have shown steady progress, while on-chip training architectures are still in their infancy. Hybrid CMOS-RRAM inference SoCs using binarized RRAMs have demonstrated higher neurosynaptic density, competitive on-chip classification accuracies, and higher energy-efficiency approaching 100 TOPS/W. RRAMs also promise in-situ training capability due to their significantly higher reported endurance approaching  $10^9$  write cycles. However, several device-level challenges such as the controllability of multilevel states and state drift need to be addressed. Resistance state drift in RRAM synapses degrades the classification performance of the neural network models, that rely on multi-bit synapses, only in a few hours and thus restoration of the states will have to be addressed at the circuit as well as algorithmic level.

As far as algorithms are concerned, a direct adaptation of Backprop to SNNs may not be the actual algorithm responsible for cognitive ‘computation’ occurring in the biological brains. Nevertheless, it provides an intermittent solution to embedded AI applications desired by the computing community. Needless to say, development of learning algorithms for SNN is a promising area of research and together with developments in the field of computational neuroscience, it may lead to better understanding of brain computation. However, going forward with the development of large-scale neuromorphic computing architectures, these algorithms will synergistically evolve by accommodating the realistic behavior of synaptic devices and by alleviating the hardware bottlenecks that arise when deep learning algorithms are mapped to in-memory computing hardware.

## ACKNOWLEDGMENTS

The author would like to thank Dr. Tomas Wu for CBRAM characterization, Prof. Maria Mitkova for providing CBRAM samples fabricated by Dr. Muhammad Rizwan Latif, Prof. John Chiasson and Ruthvik Vaila for the results on unsupervised learning in SNNs.

## DATA AVAILABILITY

Data sharing is not applicable to this article as no new data were created or analyzed in this study.

- <sup>1</sup>Y. LeCun, Y. Bengio, and G. Hinton, *Nature* **521**, 436 (2015).
- <sup>2</sup>I. Goodfellow, Y. Bengio, and A. Courville, *Deep Learning* (MIT press, 2016).
- <sup>3</sup>M. Nielsen, *Neural Networks and Deep Learning*, 1st ed. (2017).
- <sup>4</sup>M. H. Ionica and D. Gregg, *IEEE Micro* **35**, 6 (2015).
- <sup>5</sup>P. Blouw, X. Choo, E. Hunsberger, and C. Eliasmith, arXiv preprint arXiv:1812.01739 (2018).
- <sup>6</sup>M. Abadi, P. Barham, J. Chen, Z. Chen, A. Davis, J. Dean, M. Devin, S. Ghemawat, G. Irving, M. Isard, M. Kudlur, J. Levenberg, R. Monga, S. Moore, D. G. Murray, B. Steiner, P. Tucker, V. Vasudevan, P. Warden, M. Wicke, Y. Yu, and X. Zheng, in *OSDI*, Vol. 16 (2016) pp. 265–283.
- <sup>7</sup>A. Paszke, S. Gross, S. Chintala, G. Chanan, E. Yang, Z. DeVito, Z. Lin, A. Desmaison, L. Antiga, and A. Lerer, (2017).
- <sup>8</sup>A. Fuchs and D. Wentzlaff, in *2019 IEEE International Symposium on High Performance Computer Architecture (HPCA)* (IEEE, 2019) pp. 1–14.
- <sup>9</sup>IRDS, “International Roadmap for Devices and Systems (IRDS): More Moore Technologies,” (2017).
- <sup>10</sup>R. Waser, R. Dittmann, G. Staikov, and K. Szot, *Advanced materials* **21**, 2632 (2009).
- <sup>11</sup>D. Ielmini and H.-S. P. Wong, *Nature Electronics* **1**, 333 (2018).
- <sup>12</sup>A. Sebastian, M. Le Gallo, R. Khaddam-Aljameh, and E. Eleftheriou, *Nature Nanotechnology*, **1** (2020).
- <sup>13</sup>Z. Wang, H. Wu, G. W. Burr, C. S. Hwang, K. L. Wang, Q. Xia, and J. J. Yang, *Nature Reviews Materials*, **1** (2020).
- <sup>14</sup>W. Zhang, R. Mazzarello, M. Wuttig, and E. Ma, *Nature Reviews Materials* **4**, 150 (2019).
- <sup>15</sup>A. Sebastian, M. Le Gallo, and E. Eleftheriou, *Journal of Physics D: Applied Physics* **52**, 443002 (2019).
- <sup>16</sup>Q. Xia and J. J. Yang, *Nature materials* **18**, 309 (2019).
- <sup>17</sup>M. Pfeiffer and T. Pfeil, *Frontiers in neuroscience* **12**, 774 (2018).
- <sup>18</sup>J. Hasler and H. B. Marr, *Frontiers in neuroscience* **7**, 118 (2013).
- <sup>19</sup>C. D. Schuman, T. E. Potok, R. M. Patton, J. D. Birdwell, M. E. Dean, G. S. Rose, and J. S. Plank, arXiv preprint arXiv:1705.06963 (2017).
- <sup>20</sup>Z. You and S. Wei, “White Paper on AI Chip Technologies,” (2018).
- <sup>21</sup>S.-C. Liu, *Event-Based Neuromorphic Systems* (John Wiley & Sons, 2015).
- <sup>22</sup>S. Moradi, N. Qiao, F. Stefanini, and G. Indiveri, *IEEE transactions on biomedical circuits and systems* **12**, 106 (2017).
- <sup>23</sup>V. Saxena, X. Wu, I. Srivastava, and K. Zhu, *Journal of Low Power Electronics and Applications* **8**, 34 (2018).
- <sup>24</sup>W. Gerstner, R. Kempter, J. L. van Hemmen, and H. Wagner, *Nature* **383**, 76 (1996).
- <sup>25</sup>G.-q. Bi and M.-m. Poo, *The Journal of neuroscience* **18**, 10464 (1998).
- <sup>26</sup>G.-q. Bi and M.-m. Poo, *Annual review of neuroscience* **24**, 139 (2001).
- <sup>27</sup>P. J. Sjöström, G. G. Turrigiano, and S. B. Nelson, *Neuron* **32**, 1149 (2001).

- <sup>28</sup>N. L. Golding, N. P. Staff, and N. Spruston, *Nature* **418**, 326 (2002).
- <sup>29</sup>Y. Dan and M.-m. Poo, *Neuron* **44**, 23 (2004).
- <sup>30</sup>Y. Dan and M.-M. Poo, *Physiological reviews* **86**, 1033 (2006).
- <sup>31</sup>T. Masquelier and S. J. Thorpe, *PLoS computational biology* **3**, e31 (2007).
- <sup>32</sup>U. Weidenbacher and H. Neumann, in *Perception in Multimodal Dialogue Systems* (Springer, 2008) pp. 123–131.
- <sup>33</sup>H. Lee, *Unsupervised feature learning via sparse hierarchical representations* (Stanford University, 2010).
- <sup>34</sup>B. Nessler, M. Pfeiffer, L. Buesing, and W. Maass, *PLoS computational biology* **9**, e1003037 (2013).
- <sup>35</sup>D. Sterratt, B. Graham, A. Gillies, and D. Willshaw, *Principles of computational modelling in neuroscience* (Cambridge University Press, 2011).
- <sup>36</sup>R. Wang, T. J. Hamilton, J. Tapson, and A. van Schaik, in *Circuits and Systems (ISCAS), 2014 IEEE International Symposium on* (IEEE, 2014) pp. 1564–1567.
- <sup>37</sup>A. van Schaik, *Neural Networks* **14**, 617 (2001).
- <sup>38</sup>E. M. Izhikevich, *IEEE Transactions on neural networks* **14**, 1569 (2003).
- <sup>39</sup>A. van Schaik, C. Jin, A. McEwan, and T. J. Hamilton, in *Proceedings of 2010 IEEE International Symposium on Circuits and Systems* (IEEE, 2010) pp. 4253–4256.
- <sup>40</sup>G. Indiveri, B. Linares-Barranco, T. J. Hamilton, A. Van Schaik, R. Etienne-Cummings, T. Delbruck, S.-C. Liu, P. Dudek, P. Häfliger, S. Renaud, *et al.*, *Frontiers in neuroscience* **5**, 73 (2011).
- <sup>41</sup>X. Wu, V. Saxena, and K. Zhu, *IEEE Journal on Emerging and Selected Topics in Circuits and Systems (JETCAS)* **5**, 254 (2015).
- <sup>42</sup>Y. Bengio, T. Mesnard, A. Fischer, S. Zhang, and Y. Wu, *arXiv preprint arXiv:1509.05936* (2015).
- <sup>43</sup>J. Sjöström and W. Gerstner, *Scholarpedia* **5**, 1362 (2010).
- <sup>44</sup>P. A. Merolla, J. V. Arthur, R. Alvarez-Icaza, A. S. Cassidy, J. Sawada, F. Akopyan, B. L. Jackson, N. Imam, C. Guo, Y. Nakamura, B. Brezzo, I. Vo, S. K. Esser, R. Appuswamy, B. Taba, A. Amir, M. D. Flickner, W. P. Rish, R. Manohar, and D. S. Modha, *Science Magazine* **345**, 668 (2014).
- <sup>45</sup>M. Davies, N. Srinivasa, T.-H. Lin, G. Chinya, Y. Cao, S. H. Choday, G. Dimou, P. Joshi, N. Imam, S. Jain, Y. Liao, C.-K. Lin, A. Lines, R. Liu, D. Mathaikutty, S. McCoy, A. Paul, J. Tse, G. Venkataramanan, Y.-H. Weng, A. Wild, Y. Yang, and H. Wang, *IEEE Micro* **38**, 82 (2018).
- <sup>46</sup>K. Boahen, in *International Conference of the IEEE Engineering in Medicine and Biology Society* (2006).
- <sup>47</sup>G. Indiveri, E. Chicca, and R. Douglas, *IEEE Transactions On Neural Networks* **17**, 211 (2006).
- <sup>48</sup>V. Saxena, X. Wu, and K. Zhu, in *2018 IEEE International Symposium on Circuits and Systems (ISCAS)* (2018) pp. 1–5.
- <sup>49</sup>T. Pfeil, T. C. Potjans, S. Schrader, W. Potjans, J. Schemmel, M. Diesmann, and K. Meier, *arXiv preprint arXiv:1201.6255* (2012).
- <sup>50</sup>E. Neftci, S. Das, B. Pedroni, K. Kreutz-Delgado, and G. Cauwenberghs, *Frontiers in neuroscience* **7** (2013).
- <sup>51</sup>S. Brink, S. Nease, and P. Hasler, *Neural Networks* (2013).
- <sup>52</sup>A. Chen, *Solid-State Electronics* **125**, 25 (2016).
- <sup>53</sup>D. Jana, S. Roy, R. Panja, M. Dutta, S. Z. Rahaman, R. Mahapatra, and S. Maikap, *Nanoscale Research Letters* **10**, 188 (2015).
- <sup>54</sup>S. Dutta, C. Schafer, J. Gomez, K. Ni, S. Joshi, and S. Datta, *Frontiers in Neuroscience* **14** (2020).
- <sup>55</sup>D. B. Strukov, G. S. Snider, D. R. Stewart, and R. S. Williams, *Nature* **443**, 80 (2008).
- <sup>56</sup>A. Rothenbuhler, T. Tran, E. Smith, V. Saxena, and K. Campbell, *Journal of Low Power Electronics and Applications* **3**, 174 (2013).
- <sup>57</sup>D. Kuzum, R. G. Jeyasingh, B. Lee, and H.-S. P. Wong, *Nano letters* **12**, 2179 (2011).
- <sup>58</sup>G. W. Burr, R. M. Shelby, A. Sebastian, S. Kim, S. Kim, S. Sidler, K. Virwani, M. Ishii, P. Narayanan, A. Fumarola, L. L. Sanches, I. Boybat, M. Le Gallo, K. Moon, J. Woo, H. Hwang, and Y. Leblebici, *Advances in Physics: X* **2**, 89 (2017).
- <sup>59</sup>Z. Sun, X. Bi, H. Li, W.-F. Wong, and X. Zhu, *IEEE Transactions on Very Large Scale Integration (VLSI) Systems* **22**, 1281 (2013).
- <sup>60</sup>M. Jerry, P.-Y. Chen, J. Zhang, P. Sharma, K. Ni, S. Yu, and S. Datta, in *2017 IEEE International Electron Devices Meeting (IEDM)* (IEEE, 2017) pp. 6–2.
- <sup>61</sup>J. J. Yang, D. B. Strukov, and D. R. Stewart, *Nature nanotechnology* **8**, 13 (2013).
- <sup>62</sup>S. Yu, Y. Wu, R. Jeyasingh, D. Kuzum, and H.-S. P. Wong, *Electron Devices, IEEE Transactions on* **58**, 2729 (2011).
- <sup>63</sup>D. Kuzum, S. Yu, and H. P. Wong, *Nanotechnology* **24**, 382001 (2013).
- <sup>64</sup>S. Menzel, M. Waters, A. Marchewka, U. BÄřÄŕttger, R. Dittmann, and R. Waser, *Advanced Functional Materials* **21**, 4487 (2011).
- <sup>65</sup>R. Waser, R. Dittmann, G. Staikov, and K. Szot, *Advanced Materials* **21**, 2632 (2009).
- <sup>66</sup>M. Zhao, H. Wu, B. Gao, Q. Zhang, W. Wu, S. Wang, Y. Xi, D. Wu, N. Deng, S. Yu, H.-Y. Chen, and H. Qian, in *IEEE International Electron Devices Meeting (IEDM)* (IEEE, 2017) pp. 39–4.
- <sup>67</sup>A. Benoist, S. Blonkowski, S. Jeannot, S. Denorme, J. Damiens, J. Berger, P. Candelier, E. Vianello, H. Grampeix, J. Nodin, E. Jalaguier, and L. Perniola, in *2014 IEEE International Reliability Physics Symposium* (IEEE, 2014) pp. 2E–6.
- <sup>68</sup>X. Guo, F. M. Bayat, M. Prezioso, Y. Chen, B. Nguyen, N. Do, and D. B. Strukov, in *2017 IEEE Custom Integrated Circuits Conference (CICC)* (IEEE, 2017) pp. 1–4.
- <sup>69</sup>F. Cai, J. M. Correll, S. H. Lee, Y. Lim, V. Bothra, Z. Zhang, M. P. Flynn, and W. D. Lu, *Nature Electronics* **2**, 290 (2019).
- <sup>70</sup>D. Ielmini and R. Waser, *Resistive switching: from fundamentals of nanoionic redox processes to memristive device applications* (John Wiley & Sons, 2015).
- <sup>71</sup>S. H. Jo, K.-H. Kim, and W. Lu, *Nano letters* **9**, 870 (2009).
- <sup>72</sup>S. H. Jo, T. Chang, I. Ebong, B. B. Bhadviya, P. Mazumder, and W. Lu, *Nano letters* **10**, 1297 (2010).
- <sup>73</sup>K. Seo, I. Kim, S. Jung, M. Jo, S. Park, J. Park, J. Shin, K. P. Biju, J. Kong, K. Lee, B. Lee, and H. Hwang, *Nanotechnology* **22**, 254023 (2011).
- <sup>74</sup>Y. Li, Y. Zhong, L. Xu, J. Zhang, X. Xu, H. Sun, and X. Miao, *Scientific Reports* **3** (2013).
- <sup>75</sup>R. Waser, R. Dittmann, G. Staikov, and K. Szot, *Advanced materials* **21**, 2632 (2009).
- <sup>76</sup>R. Waser, D. Ielmini, H. Akinaga, H. Shima, H.-S. P. Wong, J. J. Yang, and S. Yu, *Resistive Switching: From Fundamentals of Nanoionic Redox Processes to Memristive Device Applications*, 1 (2016).
- <sup>77</sup>D. J. Wouters, in *Workshop on Memristive systems for Space Applications, Noordwijk, Netherlands* (2015).
- <sup>78</sup>M. Kozicki, M. Balakrishnan, C. Gopalan, C. Ratnakumar, and M. Mitkova, in *Symposium Non-Volatile Memory Technology 2005*. (Institute of Electrical & Electronics Engineers (IEEE)).
- <sup>79</sup>M. N. Kozicki and M. Mitkova, *Memory Devices Based on Mass Transport in Solid-state Electrolytes, Chapter 16 Nanotechnology. Volume 3*, edited by R. Waser (WILEY-VCHGmbH&Co KGaA and Weinheim, 2010).
- <sup>80</sup>M. N. Kozicki, M. Mitkova, and I. Valov, “Electrochemical metallization memories,” in *Resistive Switching* (Wiley-Blackwell, 2016) pp. 483–514.
- <sup>81</sup>M. R. Latif, M. Mitkova, G. Tompa, and E. Coleman, in *2013 IEEE Workshop on Microelectronics and Electron Devices (WMED)* (IEEE, 2013) pp. 1–4.
- <sup>82</sup>C. Nail, G. Molas, P. Blaise, G. Piccolboni, B. Sklenard, C. Cagli, M. Bernard, A. Roule, M. Azzaz, E. Vianello, *et al.*, in *2016 IEEE International Electron Devices Meeting (IEDM)* (IEEE, 2016) pp. 4–5.



- <sup>83</sup>M. R. Latif, *Nano-Ionic Redox Resistive RAM-Device Performance Enhancement Through Materials Engineering, Characterization and Electrical Testing*, Ph.D. thesis (2014).
- <sup>84</sup>C. Gopalan, Y. Ma, T. Gallo, J. Wang, E. Runnion, J. Saenz, F. Koushan, P. Blanchard, and S. Hollmer, *Solid-State Electronics* **58**, 54 (2011).
- <sup>85</sup>“Breakthrough resistive ram (reram) technology,”.
- <sup>86</sup>M. Lanza, *Materials* **7**, 2155 (2014).
- <sup>87</sup>X. Hong, D. J. Loy, P. A. Dananjaya, F. Tan, C. Ng, and W. Lew, *Journal of materials science* **53**, 8720 (2018).
- <sup>88</sup>L. Goux, A. Fantini, A. Redolfi, C. Chen, F. Shi, R. Degraeve, Y. Y. Chen, T. Witters, G. Groeseneken, and M. Jurczak, in *2014 Symposium on VLSI Technology (VLSI-Technology): Digest of Technical Papers* (IEEE, 2014) pp. 1–2.
- <sup>89</sup>S. Clima, K. Sankaran, Y. Y. Chen, A. Fantini, U. Celano, A. Belmonte, L. Zhang, L. Goux, B. Govoreanu, R. Degraeve, *et al.*, *physica status solidi (RRL)–Rapid Research Letters* **8**, 501 (2014).
- <sup>90</sup>D. J. Wouters, in *43rd IEEE Semiconductor Interface Specialists Conference, San Diego* (2012).
- <sup>91</sup>K. Moon, A. Fumarola, S. Sidler, J. Jang, P. Narayanan, R. M. Shelby, G. W. Burr, and H. Hwang, *IEEE Journal of the Electron Devices Society* **6**, 146 (2017).
- <sup>92</sup>A. Fumarola, S. Sidler, K. Moon, J. Jang, R. M. Shelby, P. Narayanan, Y. Leblebici, H. Hwang, and G. W. Burr, *IEEE Journal of the Electron Devices Society* **6**, 169 (2017).
- <sup>93</sup>M. Latif, T. Nichol, M. Mitkova, D. Tenne, I. Csarnovics, S. Kökenyesi, and A. Csik, in *2014 IEEE Workshop On Microelectronics And Electron Devices (WMED)* (IEEE, 2014) pp. 1–4.
- <sup>94</sup>M. Latif, I. Csarnovics, S. Kökenyesi, A. Csik, and M. Mitkova, *Canadian Journal of Physics* **92**, 623 (2013).
- <sup>95</sup>A. Grossi, E. Nowak, C. Zambelli, C. Pellissier, S. Bernasconi, G. Cibrario, K. El Hajjam, R. Crochemore, J. F. Nodin, P. Olivo, and L. Perniola, in *2016 IEEE International Electron Devices Meeting (IEDM)* (IEEE, 2016) pp. 4–7.
- <sup>96</sup>M. Azzaz, A. Benoist, E. Vianello, D. Garbin, E. Jalaguier, C. Cagli, C. Charpin, S. Bernasconi, S. Jeannot, T. Dewolf, G. Audoit, C. Guedj, S. Denorme1, P. Candelier1, C. Fenouillet-Beranger, and L. Perniola, in *2015 45th European Solid State Device Research Conference (ESSDERC)* (IEEE, 2015) pp. 266–269.
- <sup>97</sup>B. Govoreanu, G. Kar, Y. Chen, V. Paraschiv, S. Kubicek, A. Fantini, I. Radu, L. Goux, S. Clima, R. Degraeve, *et al.*, in *IEEE International Electron Devices Meeting (IEDM)* (IEEE, 2011) pp. 31–6.
- <sup>98</sup>S. Clima, Y. Chen, A. Fantini, L. Goux, R. Degraeve, B. Govoreanu, G. Pourtois, and M. Jurczak, *IEEE Electron Device Letters* **36**, 769 (2015).
- <sup>99</sup>Y. Y. Chen, R. Degraeve, S. Clima, B. Govoreanu, L. Goux, A. Fantini, G. S. Kar, G. Pourtois, G. Groeseneken, D. J. Wouters, *et al.*, in *IEEE International Electron Devices Meeting (IEDM)* (2012) pp. 20–3.
- <sup>100</sup>“Cmp-let memory advanced demonstrator 200mm (mad200),”.
- <sup>101</sup>P. Yao, H. Wu, B. Gao, S. B. Eryilmaz, X. Huang, W. Zhang, Q. Zhang, N. Deng, L. Shi, H.-S. P. Wong, *et al.*, *Nature communications* **8**, 1 (2017).
- <sup>102</sup>P. Kumbhare, I. Chakraborty, A. Khanna, and U. Ganguly, *IEEE Transactions on Electron Devices* **64**, 3967 (2017).
- <sup>103</sup>F. Alibart, L. Gao, B. D. Hoskins, and D. B. Strukov, *Nanotechnology* **23**, 075201 (2012).
- <sup>104</sup>W. He, H. Sun, Y. Zhou, K. Lu, K. Xue, and X. Miao, *Scientific Reports* **7**, 10070 (2017).
- <sup>105</sup>W. Wu, H. Wu, B. Gao, N. Deng, S. Yu, and H. Qian, *IEEE Electron Device Letters* **38**, 1019 (2017).
- <sup>106</sup>S. Yu, D. Kuzum, and H.-S. P. Wong, in *IEEE International Symposium on Circuits and Systems (ISCAS)* (IEEE, 2014) pp. 1062–1065.
- <sup>107</sup>F. Zahoor, T. Z. Azni Zulkifli, and F. A. Khanday, *Nanoscale Research Letters* **15**, 1 (2020).
- <sup>108</sup>M. Zhao, B. Gao, J. Tang, H. Qian, and H. Wu, *Applied Physics Reviews* **7**, 011301 (2020).
- <sup>109</sup>K. Beckmann, J. Holt, H. Manem, J. Van Nostrand, and N. C. Cady, *MRS Advances* **1**, 3355 (2016).
- <sup>110</sup>C. Li, M. Hu, Y. Li, H. Jiang, N. Ge, E. Montgomery, J. Zhang, W. Song, N. Dávila, C. E. Graves, *et al.*, *Nature Electronics* **1**, 52 (2018).
- <sup>111</sup>M. Zhao, B. Gao, Y. Xi, F. Xu, H. Wu, and H. Qian, *IEEE Journal of the Electron Devices Society* **7**, 1239 (2019).
- <sup>112</sup>P.-Y. Chen and S. Yu, in *2018 IEEE International Reliability Physics Symposium (IRPS)* (IEEE, 2018) pp. 5C–4.
- <sup>113</sup>U. Russo, D. Kamalanathan, D. Ielmini, A. L. Lacaita, and M. N. Kozicki, *IEEE transactions on electron devices* **56**, 1040 (2009).
- <sup>114</sup>C.-X. Xue, T.-Y. Huang, J.-S. Liu, T.-W. Chang, H.-Y. Kao, J.-H. Wang, T.-W. Liu, S.-Y. Wei, S.-P. Huang, W.-C. Wei, *et al.*, in *2020 IEEE International Solid-State Circuits Conference (ISSCC)* (IEEE, 2020) pp. 244–246.
- <sup>115</sup>A. Valentian, F. Rummens, E. Vianello, T. Mesquida, C. L.-M. de Boissac, O. Bichler, and C. Reita, in *2019 IEEE International Electron Devices Meeting (IEDM)* (IEEE, 2019) pp. 14–3.
- <sup>116</sup>F. Alibart, E. Zamanidoost, and D. B. Strukov, *Nature communications* **4** (2013).
- <sup>117</sup>T. Chang, Y. Yang, and W. Lu, *IEEE Circuits and Systems Magazine* **13**, 56 (2013).
- <sup>118</sup>T. Serrano-Gotarredona and B. Linares-Barranco, in *Electronics, Circuits and Systems (ICECS), 2012 19th IEEE International Conference on* (IEEE, 2012) pp. 949–952.
- <sup>119</sup>G. Indiveri, R. Legenstein, G. Deligeorgis, and T. Prodromakis, *Nanotechnology* **24**, 384010 (2013).
- <sup>120</sup>S. Saïghi, C. G. Mayr, T. Serrano-Gotarredona, H. Schmidt, G. Lecerf, J. Tomas, J. Grollier, S. Boyn, A. F. Vincent, D. Querlioz, S. La Barbera, F. Alibart, D. Vuillaume, O. Bichler, C. Gamrat, and B. Linares-Barranco, *Frontiers in neuroscience* **9** (2015).
- <sup>121</sup>V. Saxena, X. Wu, I. Srivastava, and K. Zhu, in *Proceedings of the 4th ACM International Conference on Nanoscale Computing and Communication* (ACM, 2017) p. 18.
- <sup>122</sup>M. Prezioso, M. Mahmoodi, F. M. Bayat, H. Nili, H. Kim, A. Vincent, and D. Strukov, *Nature communications* **9**, 1 (2018).
- <sup>123</sup>C. Walczyk, D. Walczyk, T. Schroeder, T. Bertaud, M. Sowinska, M. Lukosius, M. Frischke, D. Wolansky, B. Tillack, E. Miranda, *et al.*, *IEEE transactions on electron devices* **58**, 3124 (2011).
- <sup>124</sup>V. Saxena, in *IEEE International Symposium on Circuits & Systems (ISCAS)* (2020).
- <sup>125</sup>V. Saxena, in *(invited) IEEE Int. Midwest Symposium on Circuits and Systems (MWSCAS)* (2019).
- <sup>126</sup>X. Wu, V. Saxena, and K. A. Campbell, in *SPIE Sensing Technology+ Applications* (International Society for Optics and Photonics, 2014) pp. 911906–911906.
- <sup>127</sup>V. Saxena, in *IEEE International Midwest Symposium on Circuits and Systems (MWSCAS)* (2018) pp. 1–5.
- <sup>128</sup>S. Ambrogio, S. Balatti, V. Milo, R. Carboni, Z. Wang, A. Calderoni, N. Ramaswamy, and D. Ielmini, in *2016 IEEE Symposium on VLSI Technology* (IEEE, 2016) pp. 1–2.
- <sup>129</sup>Z. Wang, S. Joshi, S. E. Savelev, H. Jiang, R. Midya, P. Lin, M. Hu, N. Ge, J. P. Strachan, Z. Li, Z. Li, Q. Wu, M. Barnell, G.-L. Li, L. Xin, Huolin, R. S. Williams, Q. Xia, and J. J. Yang, *Nature materials* **16**, 101 (2017).
- <sup>130</sup>C. Du, W. Ma, T. Chang, P. Sheridan, and W. D. Lu, *Advanced Functional Materials* **25**, 4290 (2015).
- <sup>131</sup>T. Hirtzlin, M. Bocquet, B. Penkovsky, J.-O. Klein, E. Nowak, E. Vianello, J.-M. Portal, and D. Querlioz, *Frontiers in Neuroscience* **13** (2019).
- <sup>132</sup>V. Joshi, M. Le Gallo, S. Haefeli, I. Boybat, S. R. Nandakumar, C. Piveteau, M. Dazzi, B. Rajendran, A. Sebastian, and E. Eleftheriou, *Nature Communications* **11**, 1 (2020).

- <sup>133</sup>Q. Liu, B. Gao, P. Yao, D. Wu, J. Chen, Y. Pang, W. Zhang, Y. Liao, C.-X. Xue, W.-H. Chen, *et al.*, in *2020 IEEE International Solid-State Circuits Conference-(ISSCC)* (IEEE, 2020) pp. 500–502.
- <sup>134</sup>K. Kanda, N. Shibata, T. Hisada, K. Isobe, M. Sato, Y. Shimizu, T. Shimizu, T. Sugimoto, T. Kobayashi, N. Kanagawa, *et al.*, *IEEE Journal of solid-state circuits* **48**, 159 (2012).
- <sup>135</sup>G. H. Kim, H. Ju, M. K. Yang, D. K. Lee, J. W. Choi, J. H. Jang, S. G. Lee, I. S. Cha, B. K. Park, J. H. Han, *et al.*, *Small* **13**, 1701781 (2017).
- <sup>136</sup>J.-C. Liu, T.-Y. Wu, and T.-H. Hou, *IEEE Transactions on Circuits and Systems II: Express Briefs* **65**, 617 (2018).
- <sup>137</sup>W. Zhang, X. Peng, H. Wu, B. Gao, H. He, Y. Zhang, S. Yu, and H. Qian, in *2019 56th ACM/IEEE Design Automation Conference (DAC)* (IEEE, 2019) pp. 1–6.
- <sup>138</sup>Z. Alamgir, K. Beckmann, J. Holt, and N. C. Cady, *Applied Physics Letters* **111**, 063111 (2017).
- <sup>139</sup>J. M. Correll, V. Bothra, F. Cai, Y. Lim, S. H. Lee, S. Lee, W. D. Lu, Z. Zhang, and M. P. Flynn, **6** (2020), 10.1109/JXCDC.2020.2992228.
- <sup>140</sup>J. M. Cruz-Albrecht, M. W. Yung, and N. Srinivasa, *IEEE transactions on biomedical circuits and systems* **6**, 246 (2012).
- <sup>141</sup>X. Wu, V. Saxena, K. Zhu, and S. Balagopal, *IEEE Transactions on Circuits and Systems II: Express Briefs* **62**, 1088 (2015).
- <sup>142</sup>B. D. Sahoo, in *Circuits and Systems (ISCAS), 2017 IEEE International Symposium on* (IEEE, 2017) pp. 1–4.
- <sup>143</sup>B. Larras, P. Chollet, C. Lahuec, F. Seguin, and M. Arzel, in *2017 IEEE International Symposium on Circuits and Systems (ISCAS)* (IEEE, 2017) pp. 1–4.
- <sup>144</sup>I. Sourikopoulos, S. Hedayat, C. Loyez, F. Danneville, V. Hoel, E. Mercier, and A. Cappy, *Frontiers in Neuroscience* **11**, 123 (2017).
- <sup>145</sup>A. Joubert, B. Belhadj, and R. Hélot, in *New Circuits and Systems Conference (NEWCAS), 2011 IEEE 9th International* (IEEE, 2011) pp. 9–12.
- <sup>146</sup>W. Wan, R. Kubendran, S. B. Eryilmaz, W. Zhang, Y. Liao, D. Wu, S. Deiss, B. Gao, P. Raina, S. Joshi, *et al.*, in *2020 IEEE International Solid-State Circuits Conference-(ISSCC)* (IEEE, 2020) pp. 498–500.
- <sup>147</sup>I. E. Ebong and P. Mazumder, *Proceedings of the IEEE* **100**, 2050 (2012).
- <sup>148</sup>X. Wu, V. Saxena, K. Zhu, and S. Balagopal, accepted in the *IEEE Transactions on Circuits and Systems II: Express Briefs* (2015), 10.1109/TCSII.2015.2456372.
- <sup>149</sup>X. Wu, *Analog Spiking Neuromorphic Circuits and Systems for Brain- and Nanotechnology-Inspired Cognitive Computing*, Ph.D. thesis (2016).
- <sup>150</sup>C. Yakopcic, T. M. Taha, G. Subramanyam, and R. E. Pino, *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems* **32**, 1201 (2013).
- <sup>151</sup>P. Y. Simard, D. Steinkraus, and J. C. Platt, in *null* (IEEE, 2003) p. 958.
- <sup>152</sup>L. Wan, M. Zeiler, S. Zhang, Y. Le Cun, and R. Fergus, in *International Conference on Machine Learning* (2013) pp. 1058–1066.
- <sup>153</sup>A. Krizhevsky, I. Sutskever, and G. E. Hinton, in *Advances In Neural Information Processing Systems* (2012) pp. 1097–1105.
- <sup>154</sup>K. He, X. Zhang, S. Ren, and J. Sun, in *European conference on computer vision* (Springer, 2016) pp. 630–645.
- <sup>155</sup>K. Simonyan and A. Zisserman, arXiv preprint arXiv:1409.1556 (2014).
- <sup>156</sup>P. U. Diehl, D. Neil, J. Binas, M. Cook, S.-C. Liu, and M. Pfeiffer, in *International Joint Conference on Neural Networks (IJCNN)* (2015) pp. 1–8.
- <sup>157</sup>E. Hunsberger and C. Eliasmith, arXiv preprint arXiv:1611.05141 (2016).
- <sup>158</sup>B. Rueckauer, I.-A. Lungu, Y. Hu, M. Pfeiffer, and S.-C. Liu, *Frontiers in neuroscience* **11**, 682 (2017).
- <sup>159</sup>A. Sengupta, Y. Ye, R. Wang, C. Liu, and K. Roy, *Frontiers in neuroscience* **13**, 95 (2019).
- <sup>160</sup>P. U. Diehl and M. Cook, *Frontiers in Computational Neuroscience* **9** (2015).
- <sup>161</sup>S. R. Kheradpisheh, M. Ganjtabesh, S. J. Thorpe, and T. Masquelier, *Neural Networks* **99**, 56 (2018).
- <sup>162</sup>R. Vaila, J. Chiasson, and V. Saxena, arXiv:1903.12272.
- <sup>163</sup>E. O. Neftci, C. Augustine, S. Paul, and G. Detorakis, *Frontiers in neuroscience* **11**, 324 (2017).
- <sup>164</sup>J. H. Lee, T. Delbruck, and M. Pfeiffer, *Frontiers in Neuroscience* **10** (2016).
- <sup>165</sup>P. O'Connor, E. Gavves, and M. Welling, arXiv preprint arXiv:1706.04159 (2017).
- <sup>166</sup>S. R. Kulkarni and B. Rajendran, *Neural Networks* **103**, 118 (2018).
- <sup>167</sup>H. Mostafa, *IEEE transactions on neural networks and learning systems* **29**, 3227 (2017).
- <sup>168</sup>S. B. Shrestha and G. Orchard, in *Advances in Neural Information Processing Systems* (2018) pp. 1412–1421.
- <sup>169</sup>H. Mostafa, V. Ramesh, and G. Cauwenberghs, *Frontiers in neuroscience* **12**, 608 (2018).
- <sup>170</sup>E. O. Neftci, H. Mostafa, and F. Zenke, *IEEE Signal Processing Magazine* **36**, 61 (2019).
- <sup>171</sup>E. Hunsberger and C. Eliasmith, arXiv preprint arXiv:1510.08829 (2015).
- <sup>172</sup>C. Lee, P. Panda, G. Srinivasan, and K. Roy, *Frontiers in neuroscience* **12**, 435 (2018).
- <sup>173</sup>S. Nandakumar, I. Boybat, M. Le Gallo, E. Eleftheriou, A. Sebastian, and B. Rajendran, *Scientific reports* **10**, 1 (2020).
- <sup>174</sup>T. Masquelier, R. Guyonneau, and S. J. Thorpe, *PloS One* **3**, e1377 (2008).
- <sup>175</sup>T. Masquelier, R. Guyonneau, and S. J. Thorpe, *Neural computation* **21**, 1259 (2009).
- <sup>176</sup>R. Vaila, J. Chiasson, and V. Saxena, arXiv preprint arXiv:1903.12272 (2019).
- <sup>177</sup>A. Tavanaei and A. S. Maida, arXiv preprint arXiv:1611.03000 (2016).
- <sup>178</sup>R. Vaila, J. Chiasson, and V. Saxena, arXiv preprint arXiv:2002.11843 (2020).
- <sup>179</sup>R. Vaila, J. Chiasson, and V. Saxena, in *International Conference on Neuromorphic Systems (ICONS)* (2019).
- <sup>180</sup>R. Vaila, J. Chiasson, and V. Saxena, arXiv preprint arXiv:2005.04167 (2020).
- <sup>181</sup>S. M. Bohte, J. N. Kok, and H. La Poutre, *Neurocomputing* **48**, 17 (2002).
- <sup>182</sup>M. Jaderberg, W. M. Czarnecki, S. Osindero, O. Vinyals, A. Graves, D. Silver, and K. Kavukcuoglu, in *International Conference on Machine Learning* (2017) pp. 1627–1635.
- <sup>183</sup>T. P. Lillicrap, D. Cownden, D. B. Tweed, and C. J. Akerman, *Nature communications* **7**, 13276 (2016).
- <sup>184</sup>N. Anwani and B. Rajendran, in *2015 international joint conference on neural networks (IJCNN)* (IEEE, 2015) pp. 1–8.
- <sup>185</sup>F. Zenke and S. Ganguli, *Neural computation* **30**, 1514 (2018).
- <sup>186</sup>S. Yu, Z. Li, P.-Y. Chen, H. Wu, B. Gao, D. Wang, W. Wu, and H. Qian, in *2016 IEEE International Electron Devices Meeting (IEDM)* (IEEE, 2016) pp. 16–2.
- <sup>187</sup>X. Guo, F. M. Bayat, M. Bavandpour, M. Klachko, M. Mahmoodi, M. Prezioso, K. Likharev, and D. Strukov, in *2017 IEEE International Electron Devices Meeting (IEDM)* (IEEE, 2017) pp. 6–5.
- <sup>188</sup>C. Li, D. Belkin, Y. Li, P. Yan, M. Hu, N. Ge, H. Jiang, E. Montgomery, P. Lin, Z. Wang, *et al.*, *Nature communications* **9**, 1 (2018).
- <sup>189</sup>Z. Wang, S. Joshi, S. Savelev, W. Song, R. Midya, Y. Li, M. Rao, P. Yan, S. Asapu, and Y. Zhuo, *Nature Electronics* **1**, 137 (2018).
- <sup>190</sup>C.-X. Xue, W.-H. Chen, J.-S. Liu, J.-F. Li, W.-Y. Lin, W.-E. Lin, J.-H. Wang, W.-C. Wei, T.-W. Chang, T.-C. Chang, *et al.*, in *2019 IEEE International Solid-State Circuits Conference-(ISSCC)* (IEEE, 2019) pp. 388–390.
- <sup>191</sup>J. Schemmel, J. Fieres, and K. Meier, in *Neural Networks, 2008. IJCNN 2008.(IEEE World Congress on Computational*



- Intelligence*). *IEEE International Joint Conference on* (IEEE, 2008) pp. 431–438.
- <sup>192</sup>J. Schemmel, D. Brüderle, A. Gribbl, M. Hock, K. Meier, and S. Millner, in *Circuits and systems (ISCAS), proceedings of 2010 IEEE international symposium on* (IEEE, 2010) pp. 1947–1950.
- <sup>193</sup>T. Pfeil, A. Grübl, S. Jeltsch, E. Müller, P. Müller, M. A. Petrovici, M. Schmuker, D. Brüderle, J. Schemmel, and K. Meier, *Frontiers in neuroscience* **7** (2013).
- <sup>194</sup>E. Painkras, L. Plana, J. Garside, S. Temple, S. Davidson, J. Pepper, D. Clark, C. Patterson, and S. Furber, in *Custom Integrated Circuits Conference (CICC), 2012 IEEE* (IEEE, 2012) pp. 1–4.
- <sup>195</sup>C. Mayr, S. Hoepfner, and S. Furber, arXiv preprint arXiv:1911.02385 (2019).
- <sup>196</sup>C. Liu, G. Bellec, B. Vogginger, D. Kappel, J. Partzsch, F. Neumärker, S. Höppner, W. Maass, S. B. Furber, R. Legenstein, *et al.*, *Frontiers in neuroscience* **12**, 840 (2018).
- <sup>197</sup>D. Garbin, E. Vianello, O. Bichler, Q. Rafhay, C. Gamrat, G. Ghibaudo, B. DeSalvo, and L. Perniola, *IEEE Transactions on Electron Devices* **62**, 2494 (2015).
- <sup>198</sup>K.-H. Kim, S. Gaba, D. Wheeler, J. M. Cruz-Albrecht, T. Husain, N. Srinivasa, and W. Lu, *Nano letters* **12**, 389 (2011).
- <sup>199</sup>M. Prezioso, F. Merrih-Bayat, B. Hoskins, G. C. Adam, K. K. Likharev, and D. B. Strukov, *Nature* **521**, 61 (2015).
- <sup>200</sup>S. Ambrogio, S. Balatti, V. Milo, R. Carboni, Z.-Q. Wang, A. Calderoni, N. Ramaswamy, and D. Ielmini, *IEEE Transactions on Electron Devices* **63**, 1508 (2016).