



HAL
open science

Role of synaptic variability in resistive memory-based spiking neural networks with unsupervised learning

Denys R. B. Ly, Alessandro Grossi, Claire Fenouillet-Beranger, Etienne Nowak, Damien Querlioz, Elisa Vianello

► **To cite this version:**

Denys R. B. Ly, Alessandro Grossi, Claire Fenouillet-Beranger, Etienne Nowak, Damien Querlioz, et al.. Role of synaptic variability in resistive memory-based spiking neural networks with unsupervised learning. *Journal of Physics D: Applied Physics*, 2018, 51 (44), pp.444002. 10.1088/1361-6463/aad954 . hal-02939324

HAL Id: hal-02939324

<https://hal.science/hal-02939324>

Submitted on 21 Sep 2020

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

1
2
3
4
5
6
7
8
9
10 **Role of synaptic variability in Resistive**
11 **Memory-based Spiking Neural Networks with**
12 **unsupervised learning**
13

14
15 **Denys R. B. Ly¹, Alessandro Grossi¹, Claire**
16 **Fenouillet-Beranger¹, Etienne Nowak¹, Damien Querlioz²**
17 **and Elisa Vianello¹**

18 ¹ Univ. Grenoble Alpes, CEA, LETI, 38000 Grenoble, France

19 ² Centre de Nanosciences et de Nanotechnologies, CNRS, Univ. Paris-Sud,
20 Université Paris-Saclay, C2N - Orsay, 91405 Orsay Cedex, France

21 E-mail: denys.ly@cea.fr ; elisa.vianello@cea.fr
22
23

24 Submitted to: *J. Phys. D: Appl. Phys.*
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Abstract. Resistive switching memories (RRAMs) have attracted wide interest as adaptive synaptic elements in artificial bio-inspired Spiking Neural Networks (SNNs). These devices suffer from high cycle-to-cycle and cell-to-cell conductance variability, which is usually considered as a big challenge. However, biological synapses are noisy devices and the brain seems in some situations to benefit from the noise. It has been predicted that RRAM-based SNNs are intrinsically robust to synaptic variability. Here, we investigate this robustness based on extensive characterization data: we analyze the role of noise during unsupervised learning by Spike-Timing Dependent Plasticity (STDP) for detection in dynamic input data and classification of static input data. Extensive characterizations of multi-kilobits HfO₂-based Oxide-based RAM (OxRAM) arrays under different programming conditions are presented. We identify the trade-offs between programming conditions, power consumption, conductance variability and endurance features. Finally, the experimental results are used to perform system-level simulations fully calibrated on the experimental data. The results demonstrate that, similarly to biology, SNNs are not only robust to noise but a certain amount of noise can even improve the network performance. OxRAM conductance variability increases the range of synaptic values explored during the learning process. Moreover, the reduction of constraints on the OxRAM conductance variability allows the system to operate at low power programming conditions.

Keywords— Resistive switching memory (RRAM), artificial synapse, Spiking Neural Network (SNN), variability, unsupervised learning

1. Introduction

Noise is ubiquitous in any computational system, and brains are no exception. Neurons and synapses - the fundamental units of the brain - are noisy devices [1,2], due to effects such as the stochastic nature of ion channels [2], transmitter release [3], and background synaptic activity [1, 4, 5]. In the brain, up to 70% of presynaptic signals do not elicit postsynaptic signals (*synaptic failures*) [6, 7]. While it is clear that noise is present in brains, its implications are not entirely understood. Several studies have suggested that the brain may benefit from noise [8–12]. For instance, the predominance of synaptic failures could provide an energy-saving mechanism [13]. Noise may also help brains to explore possible solutions to a specific problem [14–16], preventing them from being stuck in suboptimal solutions [2, 17].

These considerations can have important implications in nanoelectronics, as today, multiple bio-inspired hardware architectures are being developed incorporating nanodevices. Many of these architectures en-

code neuron values as *spikes* [18–20] - in so-called Spiking Neural Networks (SNNs), which can lead to high energy-efficiencies. These architectures also incorporate the brain-inspired principle of learning, in the way the synaptic connections among neurons are created, modified and preserved. For this purpose, multiple works implement artificial synapses with filamentary-based Resistive Memory (RRAM) such as Oxide-based Memory (OxRAM). RRAMs consist of a capacitor-like Metal-Insulator-Metal structure, in series with a selector device. Memory cells can be integrated in the Back End Of Line (BEOL) with advanced Complementary Metal-Oxide-Semiconductor (CMOS) technology nodes (Figure 1). Upon the application of voltage pulses, RRAMs exhibit a reversible conductance switching behavior, and RRAMs can be integrated in dense arrays to connect many silicon neurons and used to implement spike-based learning mechanisms that change their conductance. Various approaches have been proposed to implement learning, such as the bio-inspired Spike-Timing Dependent Plasticity (STDP) [21–26]. STDP allows the SNN to learn the synaptic weights in an unsupervised way (*i.e.* training examples do not need to be labeled).

Despite the advantages listed above, RRAMs pose challenges. One drawback is the high conductance variability - both across cycles and cells - inducing non repeatable behavior [27, 28]. It has been demonstrated that RRAM-based SNNs are intrinsically robust to synaptic variability [22, 29–34], but a clear study explaining the origin of this robustness is still to be provided. In particular, it is still to be understood if artificial SNNs are simply robust to synaptic variability or if synaptic variability could be beneficial, in the same manner that noise might be beneficial to the brain [35].

In this work, we provide a comprehensive insight on RRAM electrical requirements for artificial SNN systems with unsupervised learning by stochastic STDP. The impact of RRAM characteristics (memory window, conductance variability, aging) on artificial SNN performance is investigated. A fully connected feed-forward neural network topology with leaky integrate and fire neurons and RRAM-based synapses is adopted. We focus on two different applications: a detection task and a classification task. In contrast to memory applications, we show that RRAM conductance variability is not necessarily detrimental for neuromorphic applications. On the contrary, it can

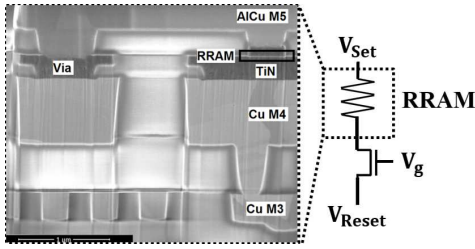


Figure 1. (Left) SEM cross-section of the TiN/HfO₂/Ti/TiN OxRAM cell integrated on the top of the fourth Cu metal layer. Both HfO₂ and Ti layers are 10 nm thick. (Right) Schematic view of the 1T1R cell configuration. The NMOS transistor is used as a selector device.

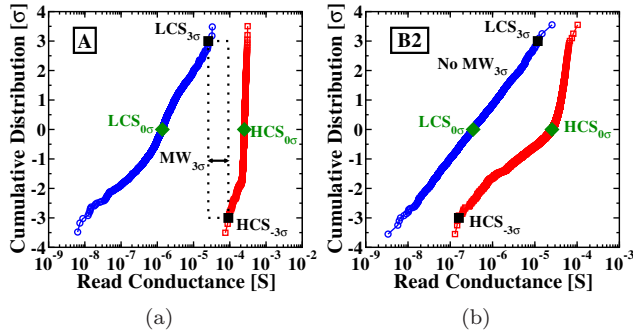


Figure 2. Cumulative distribution of the LCS and HCS measured on the 4 kbit array (a) after 1,000 switching cycles with condition A of Table 1 and (b) with condition B2 of Table 1. These distributions represent the cell-to-cell variability.

be exploited to improve the network performance while saving energy for programming.

2. Experimental characterization

2.1. Resistive Memory device characteristics

In this work, we focus on HfO₂-based OxRAM cells, integrated in the BEOL of a 130 nm CMOS logic process [28]. The memory element integration starts on top of the fourth metal layer (Cu). The cross-section of a 300 nm diameter OxRAM device is shown in Figure 1 (Left). The OxRAM devices are composed of a TiN/HfO₂/Ti/TiN stack, where both HfO₂ and Ti layers are 10 nm thick. A NMOS transistor in series with the memory element is used as a selector device (1T1R), as depicted in Figure 1 (Right), which allows each memory device of the array to be read from and written to individually. Each 1T1R structure in the matrix is addressed using a Bit Line and a Source Line, which connect to the top electrode of the device and the source of the transistor respectively, and a Word Line which connects to the gate of the transistor, and which also regulates the compliance current during the Set and Reset operations (I_{cc}).

All measurements presented here have been performed on a 4 kbit 1T1R array. Figure 2 (a) shows

Condition		A	B1	B2	C
Voltage [V]	V_{Set}	2	2	2	2
	V_{Reset}	2.5	2.5	2.5	2.5
I_{cc} [μ A]		200	57	20	600
Energy [pJ/spike]	E_{Set}	40	11.4	4	120
	E_{Reset}	50	14.3	5	150
$\sigma_{G,HCS}$ [$\log_{10}(S)$]		0.03	0.3	0.5	0.02
$\sigma_{G,LCS}$ [$\log_{10}(S)$]		0.5	0.6	0.5	0.6
$MW_{3\sigma}$ [#]		3	1.3	0.014	370
Endurance [#]		10^6	10^4	10^7	10^2

Table 1. Programming conditions used in this work, with $t_{pulse} = 100ns$.

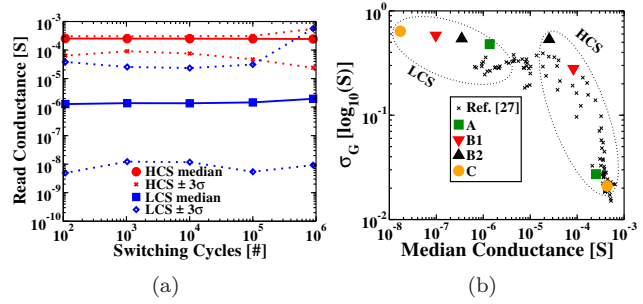


Figure 3. (a) Endurance with programming conditions A in Table 1. No smart algorithm is applied. (b) Conductance variability as a function of median conductance values for different programming conditions. Conductance variability is defined in Equation 2.

the cumulative distribution of High Conductance State (HCS) and Low Conductance State (LCS), measured with the programming conditions A of Table 1 after 1,000 Set/Reset cycles. Figure 2 (b) shows the same measurement with condition B2 in Table 1, which uses lower programming power consumption conditions than condition A. The cumulative distributions are measured on all cells on the 4 kbit array, after 1 cycle (1 Set and 1 Reset operations). In memory applications, OxRAMs are used to store one bit of information in their conductance value: OxRAMs in HCS are associated with a binary ‘1’ value, and OxRAMs in LCS are associated with a binary ‘0’ value. Therefore, it is fundamental that HCS and LCS distributions do not overlap, as in the situation of Figure 2 (a), and unlike the situation in Figure 2 (b). The appropriate separation of HCS and LCS distributions is characterized by the Memory Window at 3σ ($MW_{3\sigma}$), the ratio between the high conductance value at -3σ , $HCS_{-3\sigma}$, and the low conductance value at 3σ , $LCS_{+3\sigma}$, of the conductance distributions (Figure 2 (a)):

$$MW_{3\sigma} = \frac{HCS_{-3\sigma}}{LCS_{3\sigma}} \quad (1)$$

Figure 3 (a) shows the evolution of HCS and LCS during one million Set/Reset cycles with programming

conditions A. Solid lines represent the median values of HCS and LCS distributions ($HCS_{0\sigma}$ and $LCS_{0\sigma}$), which remain constant for the 10^6 switching cycles. However, the conductance values at $\pm 3\sigma$, represented by dotted lines, evidence an increase of conductance variability in both HCS and LCS due to OxRAM aging, which causes a reduction of the memory window. After 10^5 switching cycles, HCS and LCS distributions start overlapping and it is no longer possible to use the OxRAM for memory applications. At 10^6 cycles, oxide breakdowns can be observed in some cells, causing these *broken* cells to be stuck in HCS.

$MW_{3\sigma}$, endurance performance and variability of both HCS and LCS depend on the programming conditions (compliance current and the amplitude of Set/Reset voltage pulses) [36, 37]. Figure 3 (b) shows the conductance variability as a function of the median conductance value, measured on the 4 kbit array for different programming conditions. The conductance variability is defined as the standard deviation of the base-10 logarithm of the conductance distributions:

$$\begin{aligned}\sigma_{G,HCS} &= \text{std}[\log_{10}(G_{HCS})] \\ \sigma_{G,LCS} &= \text{std}[\log_{10}(G_{LCS})]\end{aligned}\quad (2)$$

Conductance variability is constant for conductance values lower than $77.5 \mu\text{S}$ and then decreases with the median conductance value. In order to increase the memory window, it is necessary to apply stronger Reset programming conditions to decrease the LCS median conductance value, and/or apply stronger Set programming conditions to decrease HCS variability and increase HCS median conductance value. However, this implies an increase in power consumption. In addition, it has been demonstrated that a trade-off exists between memory window and endurance performance: higher memory windows imply lower endurance [36, 38]. In this work, we focus on four representative programming conditions, which are reported on Figure 3 (b). Table 1 summarizes the parameters of each condition:

- A: compromise between endurance and $MW_{3\sigma}$ (suited condition for standard memory applications);
- B1 and B2: low power consumption, high variability in both HCS and LCS and low $MW_{3\sigma}$ (cannot be used for memory applications due to the low window margin);
- C: highest $MW_{3\sigma}$ among the four conditions, high power consumption, low HCS variability and low endurance.

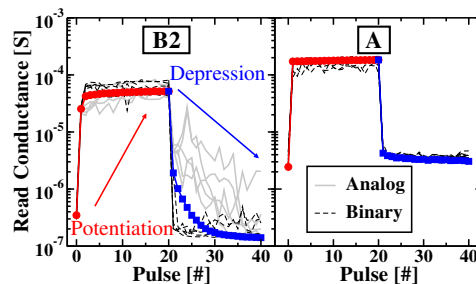


Figure 4. Conductance evolution during the application of a series of 20 identical Set pulses and Reset pulses ((Left) condition B2 and (Right) condition A of Table 1. Grey lines are representative of ten single cells. Blue and red lines corresponds to the median value calculated on 4 kbit cells. Pulse 0 is the conductance value before the first Set.

2.2. Implementation of synaptic elements and learning rule with Resistive Memories

Many implementations of OxRAM synapses seek an analog conductance modulation under identical pulses in both programming directions: when consecutive Set (Potentiation) or Reset (Depression) pulses are applied, the conductance gradually increases or decreases, respectively [22, 32–34, 39, 40]. Figure 4 reports the conductance response when a series of 20 identical Set and Reset pulses are applied on the 4 kbit OxRAM array. Light grey curves are the conductance response of single OxRAM cells with an analog behavior. Dotted black lines are the conductance response of single OxRAM cells with a binary behavior, *i.e.* an abrupt switching between LCS and HCS. Only ten single cells are plotted for the sake of clarity. Red and blue curves are the median conductance value extracted from 4 kbit cells after potentiation and depression respectively. Low and high power programming conditions (B2 and A in Table 1) are used. For low power programming conditions (Figure 4 (Left)), the evolution of the median conductance value shows an analog switching during depression. Unfortunately, this behavior is difficult to control across a large array due to the strong cell-to-cell conductance variability. Some cells present a binary behavior and only switch between two distinct states. Moreover, even in the cells presenting an analog-like switching behavior, the conductance increase (decrease) after a Set (Reset) pulse is random from cell to cell and from pulse to pulse. For high power programming conditions (Figure 4 (Right)), most of the OxRAM cells (more than 90%) present a binary behavior.

To overcome these limitations, we have proposed a synaptic compound of multiple (n) OxRAM cells connected in parallel plus a probabilistic programming scheme [24, 27]. The circuit implementation is depicted in Figure 5 (a). The number of synaptic levels is defined by the number of OxRAM cells operating

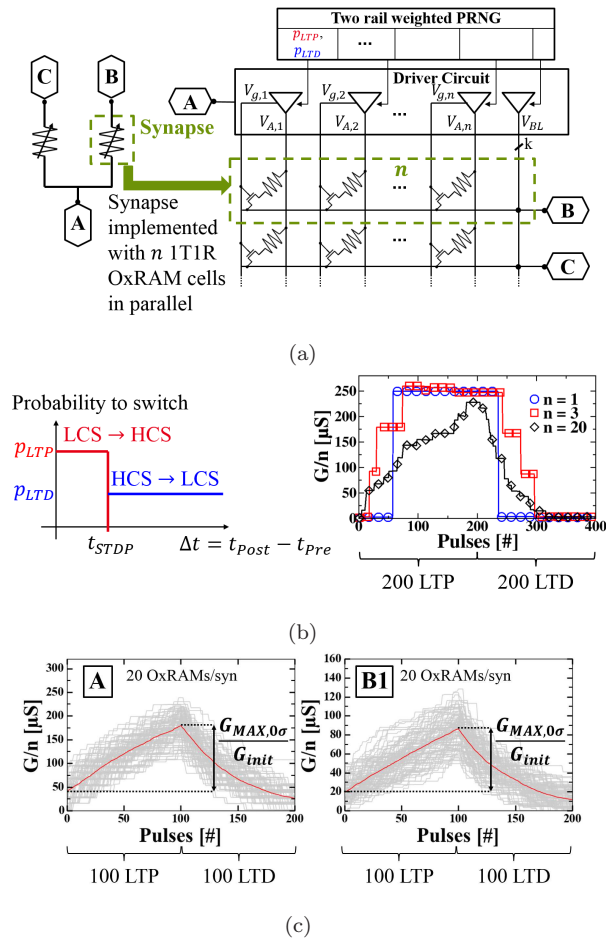


Figure 5. (a) OxRAM-based synapse implementation. The Pseudo Random Number Generator (PRNG) is used to tune the switching probabilities. (b) Stochastic STDP rule and conductance evolution of an OxRAM-based synapse composed by 1, 3 and 20 OxRAM cells in parallel when 200 potentiation pulses followed by 200 depression pulses are applied. Condition A of Table 1 is used. (c) Conductance evolution of an OxRAM-based synapses composed of 20 OxRAM cells when 100 potentiation pulses and 100 depression pulses with condition A (Left) and B1 (Right) are applied. The red line represents the mean conductance over 100 synapses; grey lines represent the evolution of each synapse.

in parallel. In order to define the conductance state of each OxRAM device (HCS or LCS), we associate this implementation with a stochastic Spike-Timing Dependent Plasticity (STDP) rule [41] - a simplified form of the bio-inspired STDP rule [42, 43]. When the presynaptic neuron spikes just before the postsynaptic neuron spikes, a Long-Term Potentiation event (LTP) occurs and each OxRAM cell of the synapse has a probability p_{LTP} to switch to the HCS. Otherwise, a Long Term Depression event (LTD) occurs and each OxRAM cell of the synapse has a probability p_{LTD} to switch to the LCS (Figure 5 (b)). The switching probabilities (p_{LTP} and p_{LTD}) can be defined using an external pseudo random number generator circuit block. Figure 5 (c) shows the impact of the programming conditions

((Left) condition A and (Right) condition B1) on the conductance evolution of a synapse composed of 20 OxRAM cells operating in parallel when 100 LTP pulses followed by 100 LTD pulses are applied. Grey lines show the conductance evolution of 100 different synapses, the red line represents the median conductance over the 100 synapses. As a probabilistic learning rule is used, the impact of device variability on the synaptic conductance response plays a secondary role with respect to the stochasticity introduced by the probabilistic STDP learning rule. In both cases we observe a gradual increase of the conductance as a function of the number of pulses. The ratio between the maximum conductance value (*i.e.* all the devices are in the HCS) and the initial conductance value (*i.e.* one device is in the HCS while the others are in the LCS), $G_{MAX,0\sigma}/G_{init}$, is about 4.5 for both programming conditions A and B1. $G_{MAX,0\sigma}/G_{init}$ is reduced to 4.3 for programming conditions B2.

3. Implications for a learning system

We now investigate the impact of OxRAM conductance variability for two different applications implemented with spiking neural networks: a detection task and a classification task. The network performance is assessed by means of system-level simulations with a special purpose neuromorphic hardware simulator [44]. The detailed OxRAM physical characterization presented in Section 2 has been implemented into physical models to understand how device properties translate in terms of learning. Variability effects due to peripheral circuits only (such as neuron variability [45]) are intentionally not taken into account. For each programming condition, the real conductance distributions measured on the 4 kbit array have been used to perform the simulations.

3.1. Network topology

Both applications are based on a one-layer fully connected feed-forward neural network topology: each neuron of the first layer is connected to each neuron of the second layer with a synaptic element. A detailed description of the simulated networks for detection and classification is provided in the Appendix. For the car detection application, the input layer corresponds to an image sensor composed of 128x128 spiking pixels, fully connected to an output layer of 60 neurons. The F1-score is used to assess network performance (see Appendix). For the digit classification application, the input layer corresponds to input images, composed of 28x28 pixels and converted into spikes with a spike frequency encoding. The input layer is fully connected to an output layer of 500 neurons. The Classification Rate (CR) is computed as the ratio

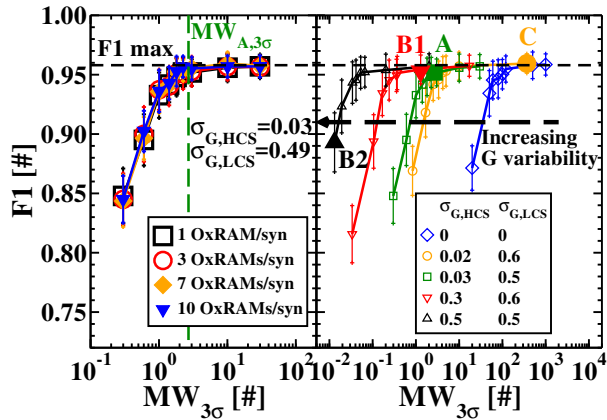


Figure 6. (Left) F1-score as a function of $MW_{3\sigma}$ defined in Equation 1 for different numbers of OxRAMs per synapse. The LCS and HCS distributions measured under condition A on the 4 kbit array were used (Figure 2 (a)). (Right) F1-score as a function of $MW_{3\sigma}$. One device per synapse is used. The LCS and HCS distributions measured on the 4 kbit array for the four conditions of Table 1 and an artificial case with zero variability were used. The $MW_{3\sigma}$ was varied by a translation of the LCS distributions to lower or higher conductance values.

between the number of successfully classified digits and the number of input digits presented. We implemented synaptic elements with the OxRAM compound presented in Figure 5 (a). Networks are trained with the unsupervised stochastic STDP rule and lateral inhibition.

3.2. Impact of OxRAM-based synapses characteristics on the network performance

The impact of OxRAM-based synapses characteristics (number of levels, OxRAM conductance variability, memory window, and aging) on the learning process of the networks designed for detection and classification is investigated.

Detection task: The first step was to study the impact of the number of synaptic levels and the OxRAM memory window on the network performance. To vary the number of synaptic levels, the number n of OxRAMs per synapse is modified. Figure 6 (Left) shows F1-score as a function of the memory window at 3σ (Equation 1) for different numbers of OxRAMs per synapse. Each point has been averaged over 20 simulations. We used LCS and HCS distributions measured under the condition A (Figure 2 (a)). $MW_{3\sigma}$ is modified by a translation of the LCS distribution to higher (decrease of $MW_{3\sigma}$) or lower (increase of $MW_{3\sigma}$) conductance values with respect to the value obtained under condition A (green dashed line in Figure 6 (Left)). This allows to decouple the impact of $MW_{3\sigma}$ from the impact of conductance variability. Surprisingly, the SNN performance is independent of the number of devices per synapse: a binary synapse is

sufficient for this type of application. We obtained the same result with the other LCS and HCS distributions from Table 1 (not shown). By contrast, the essential parameter to improve SNN performance is the $MW_{3\sigma}$: F1 increases with the $MW_{3\sigma}$ and it saturates at F1-score of about 0.96 for a memory window at 3σ larger than 3.

Second, we studied the impact of the conductance variability. We simulated the proposed application with the four LCS and HCS distributions measured under the four programming conditions presented in Table 1, plus the artificial case of a synapse with zero variability ($\sigma_{G,LCS}=0$ and $\sigma_{G,HCS}=0$) (Figure 6 (Right)). The different $MW_{3\sigma}$ values were obtained by translating the LCS distribution to lower or higher conductance values. This allows to decouple the impact of $MW_{3\sigma}$ from the conductance variability. The $MW_{3\sigma}$ corresponding to the experimental results are highlighted by a filled symbol. The higher the conductance variability, the lower the $MW_{3\sigma}$ required to reach the maximum score F1=0.96. For $\sigma_{G,LCS}=\sigma_{G,HCS}=0.5$ (condition B2, black line in Figure 6 (Right)), a $MW_{3\sigma}$ larger than 0.05 is required to reach the maximum score whereas with no variability (synapse with no conductance variability, blue line in Figure 6 (Right)), a $MW_{3\sigma}$ of at least 200 is necessary. Increasing the OxRAM synaptic variability is therefore a way to relax the constraints on the minimal $MW_{3\sigma}$ required. This can be explained by the increased dynamic range with the higher conductance variability, *i.e.* the increased range of synaptic values that can be reached during the training phase. After the training phase with the STDP learning rule, potentiated synapses (OxRAMs in HCS) represent relevant inputs, *i.e.* synapses transmitting spikes generated by a car passing on the motorway, and depressed synapses (OxRAMs in LCS) represent noisy inputs. This is well illustrated on the 2D conductance mapping after learning in the top left of Figure A1 (a) (see Appendix). As the size of a car is relatively small compared to the size of the video, the majority of the synaptic weights has to be weak (OxRAM in LCS) with a tail of stronger connections (OxRAM in HCS) in order to achieve high performance after the training phase. In our simulations, high performance after the training phase was reached (F1 \approx 0.96) if the number of OxRAM cells with conductance lower than 20 μ S was at least 50 times more numerous than the number of synapses with conductance higher than 100 μ S (see Figure 11 (Left)). The increase of both conductance variability and memory window allows for an increase of the ratio between the conductance values of potentiated and depressed synapses. It is worth noting that even with low energy programming conditions (B1, F1=0.96) we

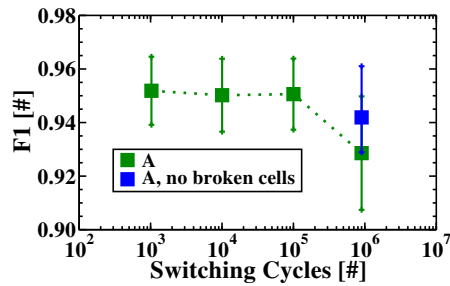


Figure 7. Impact of the OxRAM aging on F1-score. Simulations have been calibrated using the data of Figure 3 (a). Both cell-to-cell and cycle-to-cycle variability are taken into account.

have a score as good as with high energy ones (C, $F1=0.96$). The experimental condition B1 works well for neuromorphic applications whereas it cannot be used in a memory application due to its high HCS variability (no memory window). However, for the experimental condition B2, OxRAM works neither for memory nor neuromorphic applications. A decrease of F1 is observed with the experimental condition A (optimized for standard memory applications) but is still acceptable ($F1=0.95$) if we can tolerate a loss of performance with an increase in endurance.

Finally, we studied the impact of the OxRAM aging with endurance on network performance. Both cell-to-cell and cycle-to-cycle variability are taken into account. We extracted the conductance distribution during cycling on the 4 kbit array, up to one million cycles for the condition A (*cf* Figure 3 (a)) and we used these data to evaluate the impact of OxRAM aging on F1-score. The results are shown in Figure 7. We can maintain a constant F1-score of 0.95 until 10^5 cycles after which F1 plummets ($F1=0.93$). The degradation of F1 at 10^6 cycles is not due to the increase in conductance variability and decrease of $MW_{3\sigma}$ but to the broken cells (cells stuck in HCS). Upon removal of the broken cells (1%) from the distributions (blue square), it is possible to move back up to a score of 0.95.

Classification task: A similar study on the impact of OxRAM-based synapses characteristics on the network for digit classification (*cf* Figure A2 (a)) is performed. First, we investigated the impact of the number of synaptic levels and the OxRAM memory window, then the conductance variability on the SNN performance. Figure 8 reports the Classification Rate (CR) as a function of the Memory Window at 3σ ($MW_{3\sigma}$). The LCS and HCS distributions measured under the condition A and B2, plus the synapse with zero variability ($\sigma_{G,LCS}=0$ and $\sigma_{G,HCS}=0$) are used. The different curves correspond to a different number of OxRAM devices per synapse. In contrast to detection task, *the CR is independent of the $MW_{3\sigma}$ for all*

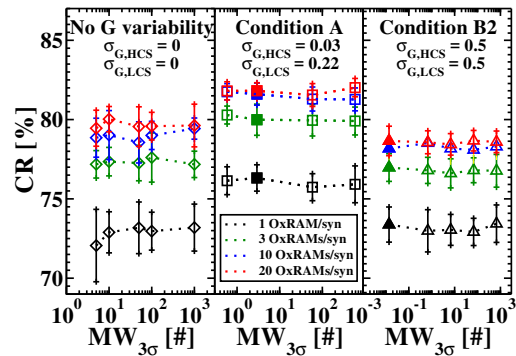


Figure 8. Classification Rate (CR) as a function of $MW_{3\sigma}$ defined in Figure 2 (a), for different numbers of OxRAMs per synapse. The LCS and HCS distributions measured on the 4 kbit array for the conditions A and B2 of Table 1 and an artificial case with zero variability were used. The $MW_{3\sigma}$ was varied by a translation of the LCS distribution to lower or higher conductance values.

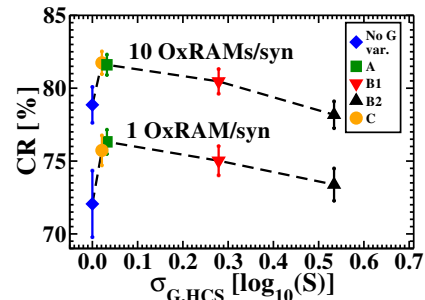


Figure 9. Classification Rate (CR) as a function of variability in the HCS, for 1 and 10 OxRAMs per synapse. The LCS and HCS distributions measured on the 4 kbit array for the conditions of Table 1 and an artificial case with zero variability were used.

the studied distributions. The network performance depends on:

- **Number of synaptic levels:** The CR increases with the number n of OxRAMs per synapse and saturates after 10 OxRAMs/syn. This is in agreement with the studies performed in [22] and [40].
- **Synaptic variability:** the LCS and HCS distributions measured under the programming conditions A allows to improve the performance with respect to the synapse with zero variability for a given number of OxRAM per synapse. However, similar performances are achieved with condition B2 (high conductance variability) and the synapse with zero variability.

To quantify and understand the impact of conductance variability, we calculated the SNN performance for the four programming conditions of Table 1 plus the synapse with zero variability. Figure 9 plots the CR as a function of HCS variability for a synapse composed by one and ten OxRAM cells. As $MW_{3\sigma}$ has no impact (Figure 8), we simulated the experimental $MW_{3\sigma}$ for

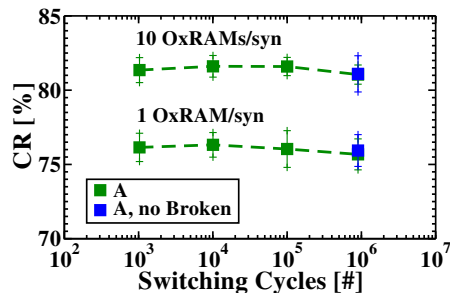


Figure 10. Impact of the OxRAM aging on the Classification Rate (CR). Simulations have been calibrated using the data of Figure 3 (a). Both cell-to-cell and cycle-to-cycle variability are taken into account.

each programming condition and $MW_{3\sigma}=5$ for the case with no variability. Network performance is maximal for $\sigma_{G,HCS}\approx 0.03$ (conditions A and C, $CR\approx 82\%$ for 10 OxRAMs/synapse) and is degraded when HCS variability is too high (condition B2, $CR=78.6\%$) or when there is no variability at all (synapse with zero variability, $CR=79.6\%$). Note that, as shown in Section 3.2, $G_{MAX,0\sigma}/G_{init}$ is reduced when there is no HCS variability or the HCS variability is too high. If we replace the cumulative OxRAM synapse by just one device, the maximum score is 76%. These results are far from the best one obtained for the same dataset, which used a supervised off-line learning approach and millions of adjustable parameters [46]. However, our results compare well to on-line supervised neural network with backpropagation and a similar number of adjustable parameters (82.9%) [32] and to previously published results on on-line unsupervised learning (82.9%) [47].

Figure 10 shows the impact of OxRAM aging on the Classification Rate. Both cell-to-cell and cycle-to-cycle variability are taken into account. Simulations have been calibrated using the data of Figure 3 (a) measured on the 4 kbit array. As HCS variability values during aging remain around 0.03, the CR varies little in that range, as shown in Figure 9. Therefore, with 10 OxRAMs/syn, we can sustain a constant score $CR\approx 81.5\%$ until 10^6 cycles.

3.3. Comparison between detection and classification tasks

We now explain the surprising difference in how detection and classification tasks are affected by device characteristics. Figure 11 shows the synaptic weight distribution after the training phase for the car detection task (Left) and the digit classification task (Right), obtained with the synapse with no conductance variability and $MW_{3\sigma}=200$ (blue), and programming conditions B1 (red) and A (green) respectively. These conditions allow maximizing the network performance. For the detection task, the maximal F1-score (0.96) is

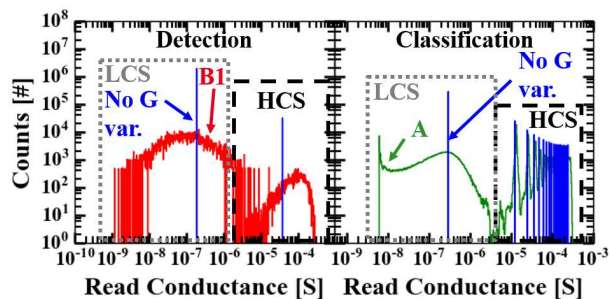


Figure 11. (Left) Synaptic weight distribution after the training phase for the detection task, for condition B1 (red, $MW_{3\sigma}=1.3$) and the ideal synapse with no conductance variability (blue, $MW_{3\sigma}=200$), with 1 OxRAM per synapse. Both conditions allow to achieve a score $F1=0.96$. (Right) Synaptic weight distribution during the inference phase for the classification task obtained with condition A (green, $MW_{3\sigma}=3$) and the synapse with no conductance variability (blue, $MW_{3\sigma}=200$), with 20 OxRAMs/syn. CR of 81.8% and 79.5% is achieved respectively.

reached if, after training, there are two synaptic populations: (1) potentiated synapses (OxRAMs in HCS) and (2) depressed synapses (OxRAMs in LCS). The fundamental requirement is that a ratio higher than 200 exists between the peaks of the potentiated (HCS) and depressed (LCS) synaptic distributions. Therefore, both memory window and conductance variability are beneficial as they increase the dynamic range of synaptic weight values available during training, facilitating the separation of the LCS and HCS peaks after training. For the classification task, multilevel conductance synapses are necessary to achieve the best performance. The number of OxRAM cells per synapse defines the number of levels. As parallel conductances sum up, the value of the equivalent synaptic weight is approximately $nHCS_{0\sigma}$, where n is the number of OxRAMs in HCS. Unlike the detection task, where the network exploits both HCS and LCS distributions, only the HCS distribution defines the synaptic value. Consequently, the LCS distribution and the memory window do not affect the network performance as shown in Figure 8. On the other hand, the classification task is sensitive to the HCS distribution as shown in Figure 9. If the HCS variability is too high (condition B2), the synaptic weight distribution after the training phase has only 6 synaptic levels (not shown) instead of the 20 levels achieved with condition A (Figure 11 (Right)), thus explaining the reduced performance for the condition B2 (Figure 9).

4. Conclusion

An extensive study of conductance variability, power consumption and aging of multi-kilobits OxRAM arrays over the full operation range has been presented. The experimental results were used to perform system-

level simulations of SNNs designed for (1) detection in dynamic patterns and (2) classification of static patterns applications. In comparison with previous studies on artificial SNNs [22, 29–32], we demonstrated that SNNs are not only robust to synaptic variability, but can also draw benefit from it. Noise can be beneficial as it increases the range of synaptic weight values available during learning. For detection applications, OxRAM technology is well-suited to implement synaptic elements as only one OxRAM device per synapse is needed, and their electrical characteristics enables to achieve maximal performance at low power consumption (less than 20 pJ/spike). On the other hand, for classification applications, multilevel conductance synapses are necessary to achieve the best performance; a synaptic compound composed of at least 10 OxRAMs per synapse is required. The maximal performance was achieved with a conductance variability in HCS of roughly 0.03 that can be achieved with programming energy of about 40 pJ/spike. This study provides guidelines to optimize the programming conditions for OxRAM-based synapses in SNNs capable of unsupervised learning by STDP. It also highlights that memory devices for neuromorphic applications may be more optimally used in different physical regimes than for conventional memory applications.

Acknowledgments

This work has been partially supported by the h2020 NeuRAM³ project. The authors would like to thank T. Werner, T. Dalgaty and T. Hirtzlin for fruitful discussion.

Appendix: Network topology

Both applications are based on a one-layer fully connected feed-forward neural network topology: each neuron of the first layer is connected to each neuron of the second layer with a synaptic element.

Car tracking

Figure A1 (a) presents the network simulated for the detection task. A video of cars passing on a six-lane wide motorway is recorded using Address Event Representation format by a Dynamic Vision Sensor [48] and it represents the input data [49]. An input pixel generates a spike each time there is a change of luminosity at its location in the visual field. Each input pixel is connected with two synapses to every output neuron to denote an increase (ON synapse) or decrease (OFF synapse) in illumination respectively. A similar network has been implemented in [39] and [41] exploiting multi-level Phase-Change Memory and

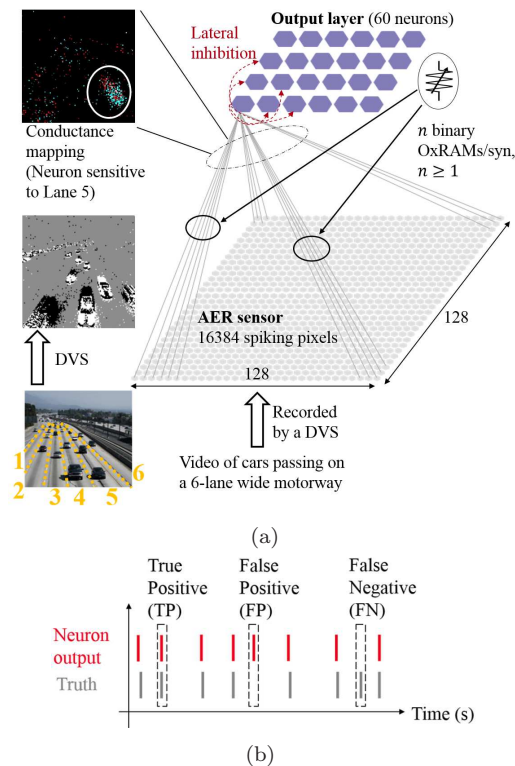


Figure A1. (a) Simulated neural network for the car tracking application, trained with a stochastic STDP learning rule and lateral inhibition. (b) Example of spiking activity of one output neuron (red) and the actual traffic (a grey spike corresponds to a car passing on the lane). True Positive (TP) events, False Positive (FP) events and False Negative (FN) events are put in evidence.

binary Conductive Bridge RRAM synapses respectively. In this work, we adopted the OxRAM technology presented in section 2. The total number of OxRAM devices is $128 \times 128 \times 2 \times 60 \times n = 1,966,080n$, where n is the number of OxRAM cells per synapse. Output neurons are implemented with the Leaky Integrate and Fire model [21]. Many implementations of such neurons have been proposed with CMOS technology [50]. Note that after an output neuron fires a spike, it cannot integrate any incoming spikes for a refractory period $T_{refrac} = 218$ ms. It also prevents all the other neurons of the layer from integrating incoming spikes for a period $T_{inhibit} = 29.9$ ms, referred to as *lateral inhibition*.

The network is trained with the unsupervised stochastic STDP rule presented in Figure 5 (b), with $p_{LTP} = 0.13$ and $p_{LTD} = 0.2$. After a training phase, every output neuron becomes sensitive to a specific lane. An example of the 2D conductance mapping of one output neuron after training is shown in the top left of Figure A1 (a). A potentiated ON synapse (resp. OFF synapse) of an input pixel is represented with a red (resp. blue) dot. When both ON and OFF synapses are potentiated, the resulting color is grey. When both ON and OFF synapses are depressed, the

resulting color is black. As a result of the training phase we can observe a pool of potentiated synapses (circled in white) denoting the sensitivity of this neuron to a car at this specific position on the motorway. When a car passes at that position, the neuron spikes. In this example, the output neuron is sensitive to the lane 5; the neuron spikes whenever a car passes on that lane. Figure A1 (b) sketches the spiking activity of one output neuron (red) and the actual traffic (a grey spike corresponds to a car passing on the lane). If the neuron detects a car, we have a True Positive (TP) event. If it spikes with no car passing, we have a False Positive (FP) event. If it misses a car, we have a False Negative (FN) event. We use the F1-score as a metric to assess network performance:

$$F1 = \frac{2TP}{2TP + FN + FP} \quad (\text{A.1})$$

F1 ranges from 0 to 1, with F1=1 being the best performance. Each output neuron becomes sensitive to one lane. Since there are 60 output neurons and only 6 lanes, several neurons become sensitive to the same lane. As more cars pass on the lanes 4 and 5, more neurons are sensitive to these lanes than to the lane 6, the least active lane. To assess network performance, only the most sensitive neuron for each lane is considered.

Digit classification

Figure A2 (a) presents the network simulated for the classification task. The Mixed National Institute of Standards and Technology (MNIST) dataset is used for the training and testing [51]. The input layer converts the input digit with a spike frequency encoding: each input neuron generates a spike train with a spiking rate f_{input} proportional to the grey level of the pixel. Synaptic elements are implemented with n OxRAMs in parallel as in Figure 5 (a). The network is trained with the unsupervised stochastic STDP rule (Figure 5 (b)), with $p_{LTP}=0.01$, $p_{LTD}=0.02$, $T_{\text{refrac}}=1$ ns and lateral inhibition with $T_{\text{inhibit}}=10$ μs .

During the training phase, each output neuron becomes sensitive to a specific class of digit, for example the output neuron 94 becomes sensitive to the class of digit '8' as illustrated in the 2D conductance mapping of Figure A2 (a). After training, each output neuron is associated with the digit it is the most sensitive to - this represents the class of the neuron. To assess network performance during the testing phase, the Classification Rate (CR) is computed as shown on Figure A2 (b). Each input digit is presented to the network for 350 μs and the output neuron that spikes the most within this time window corresponds to the network response. If the class of digit of this most active neuron is the input digit, the digit is

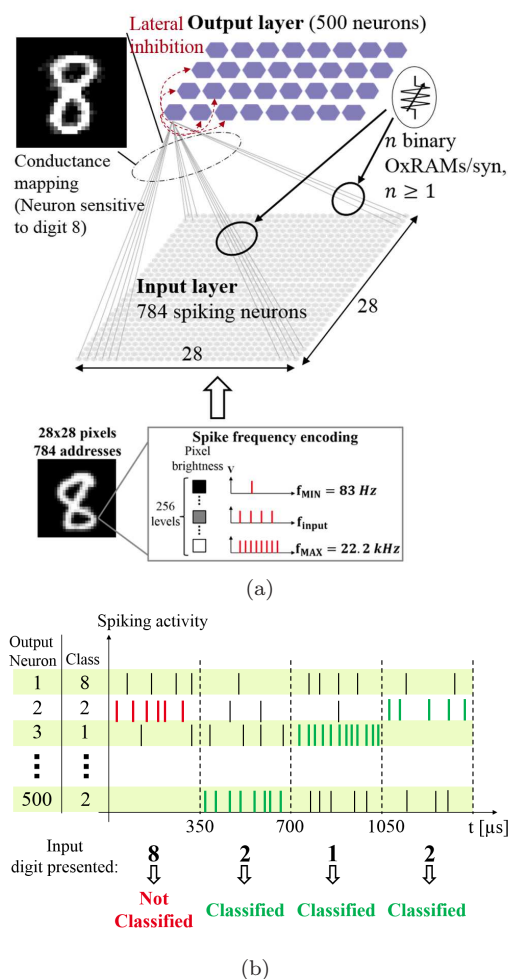


Figure A2. (a) Simulated neural network for the digit classification application, trained with a stochastic STDP learning rule and lateral inhibition. (b) Example of spiking activity of four output neurons when four different input digits are presented. If the class of digit of the most active neuron is the input digit, the digit is successfully classified (green), otherwise the digit is not classified (red).

successfully classified (green spikes in Figure 7 (b)). If its class is different from the input digit, the digit is not classified (red spikes in Figure A2 (b)). The Classification Rate is calculated as the ratio between the number of successfully classified digits ($n_{\text{classified}}$) and the number of input digits (n_{input}):

$$CR = \frac{n_{\text{classified}}}{n_{\text{input}}} \quad (\text{A.2})$$

As there are multiple ways to handwrite the same digit, increasing the number of output neurons allows for an improvement of network performance as demonstrated in [22]. Indeed, this enables the network to have at its disposal several neurons specialized to the same digit, and more precisely to have neurons specialized in different handwritings of the same digit. As shown in [22], the increase of CR with the number of output neurons saturates after 500 output neurons.

- [1] Yosef Yarom and Jorn Hounsgaard. Voltage fluctuations in neurons: Signal or noise? *Physiological Reviews*, 91(3):917–929, 2011.
- [2] A. Aldo Faisal, Luc P. J. Selen, and Daniel M. Wolpert. Noise in the nervous system. *Nat Rev Neurosci*, 9(4):292–303, Apr 2008.
- [3] C. Allen and C. F. Stevens. An evaluation of causes for unreliability of synaptic transmission. *Proc Natl Acad Sci U S A*, 91(22):10380–10383, 1994.
- [4] Rasmus M. Birn. The behavioral significance of spontaneous fluctuations in brain activity. *Neuron*, 56(1):8–9, 2007.
- [5] Michael D. Fox, Abraham Z. Snyder, Justin L. Vincent, and Marcus E. Raichle. Intrinsic fluctuations within cortical systems account for intertrial variability in human behavior. *Neuron*, 56(1):171–184, 2007.
- [6] Neal A. Hessler, Aneil M. Shirke, and Roberto Malinow. The probability of transmitter release at a mammalian central synapse. *Nature*, 366:569, Dec 1993.
- [7] J. Gerard G. Borst. The low synaptic release probability in vivo. *Trends in Neurosciences*, 33(6):259–266, 2010.
- [8] Mark D. McDonnell and Derek Abbott. What is stochastic resonance? definitions, misconceptions, debates, and its relevance to biology. *PLOS Computational Biology*, 5(5):1–9, 05 2009.
- [9] G. Bard Ermentrout, Roberto F. Galán, and Nathaniel N. Urban. Reliability, synchrony and noise. *Trends in Neurosciences*, 31(8):428–434, 2008.
- [10] David C. Knill and Alexandre Pouget. The bayesian brain: the role of uncertainty in neural coding and computation. *Trends in Neurosciences*, 27(12):712–719, 2004.
- [11] Edwin R Lewis, Kenneth R Henry, and Walter M Yamada. Essential roles of noise in neural coding and in studies of neural coding. *Biosystems*, 58(1):109–115, 2000.
- [12] Anthony Randal McIntosh, Natasa Kovacevic, and Roxane J. Itier. Increased brain signal variability accompanies lower behavioral variability in development. *PLOS Computational Biology*, 4(7):1–9, 07 2008.
- [13] Mark S. Goldman. Enhancement of information transmission efficiency by synaptic failures. *Neural Computation*, 16(6):1137–1162, 2004.
- [14] Stefan Habenschuss, Zeno Jonke, and Wolfgang Maass. Stochastic computations in cortical microcircuit models. *PLOS Computational Biology*, 9(11):1–28, 11 2013.
- [15] Peter G. H. Clarke. The limits of brain determinacy. *Proceedings of the Royal Society of London B: Biological Sciences*, 279(1734):1665–1674, 2012.
- [16] Uri Rokni, Andrew G. Richardson, Emilio Bizzi, and H. Sebastian Seung. Motor learning with unstable neural representations. *Neuron*, 54(4):653–666, 2007.
- [17] Astrid A. Prinz, Dirk Bucher, and Eve Marder. Similar network activity from disparate circuit parameters. *Nature Neuroscience*, 7:1345, Nov 2004.
- [18] Carver Mead. *Analog VLSI and Neural Systems*. Addison-Wesley Longman Publishing Co., Inc., Boston, MA, USA, 1989.
- [19] G. Indiveri, F. Corradi, and N. Qiao. Neuromorphic architectures for spiking deep neural networks. In *2015 IEEE International Electron Devices Meeting (IEDM)*, pages 4.2.1–4.2.4, Dec 2015.
- [20] Alice Mizrahi, Tifenn Hirtzlin, Akio Fukushima, Hitoshi Kubota, Shinji Yuasa, Julie Grollier, and Damien Querlioz. Neural-like computing with populations of superparamagnetic basis functions. *Nature Communications*, 9(1):1533, 2018.
- [21] O. Bichler, D. Querlioz, S. J. Thorpe, J. P. Bourgoin, and C. Gamrat. Unsupervised features extraction from asynchronous silicon retina through spike-timing-dependent plasticity. In *The 2011 International Joint Conference on Neural Networks*, pages 859–866, 2011.
- [22] D. Querlioz, O. Bichler, A. F. Vincent, and C. Gamrat. Bioinspired programming of memory devices for implementing an inference engine. *Proceedings of the IEEE*, 103(8):1398–1416, Aug 2015.
- [23] T. Werner, E. Vianello, O. Bichler, A. Grossi, E. Nowak, J. F. Nodin, B. Yvert, B. DeSalvo, and L. Perniola. Experimental demonstration of short and long term synaptic plasticity using oxram multi k-bit arrays for reliable detection in highly noisy input data. In *2016 IEEE International Electron Devices Meeting*, pages 16.6.1–16.6.4, 2016.
- [24] E. Vianello, T. Werner, O. Bichler, A. Valentian, G. Molas, B. Yvert, B. De Salvo, and L. Perniola. Resistive memories for spike-based neuromorphic circuits. In *2017 IEEE International Memory Workshop*, pages 1–6, May 2017.
- [25] M. R. Azghadi, B. Linares-Barranco, D. Abbott, and P. H. W. Leong. A hybrid cmos-memristor neuromorphic synapse. *IEEE Transactions on Biomedical Circuits and Systems*, pages 434–445, 2017.
- [26] Konstantin Zarudnyi, Adnan Mehonic, Luca Montesi, Mark Buckwell, Stephen Hudziak, and Anthony J. Kenyon. Spike-timing dependent plasticity in unipolar silicon oxide rram devices. *Frontiers in Neuroscience*, 12:57, 2018.
- [27] D. Garbin, E. Vianello, O. Bichler, Q. Raffhay, C. Gamrat, G. Ghibaudo, B. DeSalvo, and L. Perniola. Hfo₂-based oxram devices as synapses for convolutional neural networks. *IEEE Transactions on Electron Devices*, 62(8):2494–2501, 2015.
- [28] A. Grossi et al. Fundamental variability limits of filament-based rram. In *2016 IEEE International Electron Devices Meeting*, pages 4.7.1–4.7.4, 2016.
- [29] D. Garbin, E. Vianello, O. Bichler, M. Azzaz, Q. Raffhay, P. Candelier, C. Gamrat, G. Ghibaudo, B. DeSalvo, and L. Perniola. On the impact of oxram-based synapses variability on convolutional neural networks performance. In *Proceedings of the 2015 IEEE/ACM International Symposium on Nanoscale Architectures (NANOARCH15)*, pages 193–198, 2015.
- [30] T. Werner, D. Garbin, E. Vianello, O. Bichler, D. Cattaert, B. Yvert, B. De Salvo, and L. Perniola. Real-time decoding of brain activity by embedded spiking neural networks using oxram synapses. In *2016 IEEE International Symposium on Circuits and Systems (ISCAS)*, pages 2318–2321, May 2016.
- [31] M. Suri, V. Parmar, A. Kumar, D. Querlioz, and F. Alibart. Neuromorphic hybrid rram-cmos rbm architecture. In *2015 15th Non-Volatile Memory Technology Symposium (NVM-TS)*, pages 1–6, Oct 2015.
- [32] G. W. Burr, R. M. Shelby, C. di Nolfo, J. W. Jang, R. S. Shenoy, P. Narayanan, K. Virwani, E. U. Giacometti, B. Kurdi, and H. Hwang. Experimental demonstration and tolerancing of a large-scale neural network (165,000 synapses), using phase-change memory as the synaptic weight element. In *2014 IEEE International Electron Devices Meeting*, pages 29.5.1–29.5.4, Dec 2014.
- [33] Alexander Serb, Johannes Bill, Ali Khiat, Radu Berdan, Robert Legenstein, and Themis Prodromakis. Unsupervised learning in probabilistic neural networks with multi-state metal-oxide memristive synapses. *Nature Communications*, 7:12611, 2016.
- [34] Erika Covi, Stefano Brivio, Alexander Serb, Themis Prodromakis, Marco Fanciulli, and Sabina Spiga. Analog memristive synapse in spiking networks implementing unsupervised learning. *Frontiers in Neuroscience*, 10:482, 2016.
- [35] G. Pedretti, V. Milo, S. Ambrogio, R. Carboni, S. Bianchi, A. Calderoni, N. Ramaswamy, A. S. Spinelli, and D. Ielmini. Stochastic learning in neuromorphic

- hardware via spike timing dependent plasticity with rram synapses. *IEEE Journal on Emerging and Selected Topics in Circuits and Systems*, 8(1):77–85, March 2018.
- [36] E. Vianello et al. Resistive memories for ultra-low-power embedded computing design. In *2014 IEEE International Electron Devices Meeting*, pages 6.3.1–6.3.4, 2014.
- [37] A. Grossi, E. Vianello, C. Zambelli, P. Royer, J. P. Noel, B. Giraud, L. Perniola, P. Olivo, and E. Nowak. Experimental investigation of 4-kb rram arrays programming conditions suitable for tcam. *IEEE Transactions on Very Large Scale Integration (VLSI) Systems*, pages 1–9, 2018.
- [38] L. Perniola, G. Molas, G. Navarro, E. Nowak, V. Sousa, E. Vianello, and B. De Salvo. Universal signatures from non-universal memories: Clues for the future... In *2016 IEEE 8th International Memory Workshop (IMW)*, pages 1–3, May 2016.
- [39] M. Suri, O. Bichler, D. Querlioz, O. Cueto, L. Perniola, V. Sousa, D. Vuillaume, C. Gamrat, and B. DeSalvo. Phase change memory as synapse for ultra-dense neuromorphic systems: Application to complex visual pattern extraction. In *2011 International Electron Devices Meeting*, pages 4.4.1–4.4.4, Dec 2011.
- [40] Johannes Bill and Robert Legenstein. A compound memristive synapse model for statistical learning through stdp in spiking neural networks. *Front Neurosci*, 8:412, 2014.
- [41] M. Suri, O. Bichler, D. Querlioz, G. Palma, E. Vianello, D. Vuillaume, C. Gamrat, and B. DeSalvo. Cbram devices as binary synapses for low-power stochastic neuromorphic systems: Auditory (cochlea) and visual (retina) cognitive processing applications. In *2012 International Electron Devices Meeting*, pages 10.3.1–10.3.4, 2012.
- [42] Henry Markram, Joachim Lübke, Michael Frotscher, and Bert Sakmann. Regulation of synaptic efficacy by coincidence of postsynaptic aps and epsps. *Science*, 275(5297):213–215, 1997.
- [43] Guo qiang Bi and Mu ming Poo. Synaptic modification by correlated activity: Hebb’s postulate revisited. *Annual Review of Neuroscience*, 24(1):139–166, 2001.
- [44] O. Bichler, D. Roclin, C. Gamrat, and D. Querlioz. Design exploration methodology for memristor-based spiking neuromorphic architectures with the xnet event-driven simulator. In *2013 IEEE/ACM International Symposium on Nanoscale Architectures (NANOARCH)*, pages 7–12, 2013.
- [45] D. Querlioz, O. Bichler, P. Dollfus, and C. Gamrat. Immunity to device variations in a spiking neural network with memristive nanodevices. *IEEE Transactions on Nanotechnology*, 12(3):288–295, May 2013.
- [46] Dan Claudiu Ciresan, Ueli Meier, Luca Maria Gambardella, and Jürgen Schmidhuber. Deep big simple neural nets excel on handwritten digit recognition. *CoRR*, 2010.
- [47] Peter Diehl and Matthew Cook. Unsupervised learning of digit recognition using spike-timing-dependent plasticity. *Frontiers in Computational Neuroscience*, 9:99, 2015.
- [48] P. Lichtsteiner, C. Posch, and T. Delbruck. A 128 x 128 120 db 15 us latency asynchronous temporal contrast vision sensor. *IEEE Journal of Solid-State Circuits*, 43(2):566–576, Feb 2008.
- [49] Tobi Delbruck. Frame-free dynamic digital vision. *Proceedings of Intl. Symp. on Secure-Life Electronics (Tokyo, Japan: University of Tokyo)*, page 2126, 03 2008.
- [50] Giacomo et al. Indiveri. Neuromorphic silicon neuron circuits. *Frontiers in Neuroscience*, 5:73, 2011.
- [51] Y. Lecun, L. Bottou, Y. Bengio, and P. Haffner. Gradient-based learning applied to document recognition. *Pro-*