

A spiking recurrent neural network with phase change memory neurons and synapses for the accelerated solution of constraint satisfaction problems

Giacomo Pedretti¹, Member, IEEE, Piergiulio Mannocci¹, Shahin Hashemkhani¹, Valerio Milo¹, Member, IEEE, Octavian Melnic¹, Elisabetta Chicca², Member, IEEE and Daniele Ielmini¹, Fellow, IEEE

¹Dipartimento di Elettronica, Informazione e Bioingegneria, Politecnico di Milano and IU.NET, Milano 20133, Italy

²Faculty of Technology and Cognitive Interaction Technology Center of Excellence (CITEC), Bielefeld University, Bielefeld, Germany

Data-intensive computing applications such as object recognition, time series prediction and optimization tasks are becoming increasingly important in several fields including smart mobility, health and industry. Because of the large amount of data involved in the computation, the conventional von Neumann architecture suffers from excessive latency and energy consumption due to the memory bottleneck. A more efficient approach consists of in-memory computing (IMC), where computational operations are directly carried out within the data. IMC can take advantage of the rich physics of memory devices, such as their ability to store analogue values to be used in matrix-vector multiplication (MVM) and their stochasticity which is highly valuable in the frame of optimization and constraint satisfaction problems (CSPs). This work presents a stochastic spiking neuron based on a phase change memory (PCM) device for the solution of CSPs within a Hopfield recurrent neural network (RNN). In the RNN, the PCM cell is used as the integrating element of a stochastic neuron, supporting the solution of a typical CSP, namely a Sudoku puzzle in hardware. Finally, the ability to solve Sudoku puzzles using RNNs with PCM-based neurons is studied for increasing size of Sudoku puzzle by a compact simulation model, thus supporting our PCM-based RNN for data-intensive computing

I. INTRODUCTION

Optimization problems are among the most intensive computing tasks for several application fields, such as industry, finance and transport. In general, optimization is carried out by several iterations to identify the global minimum of a certain cost function. In each iteration, a conventional digital system must access the memory to fetch input data and upload the temporary output, which is time and energy consuming. To enable a more efficient optimization, a non-von Neumann architecture can be adopted, to eliminate the latency and energy spent for shuttling the data between the memory and the central processing unit (CPU) [1]. An example of non-von Neumann computing architecture is the concept of in-memory computing (IMC) where the computation is executed directly

This article has received funding from the European Research Council (ERC) under the European Union's Horizon 2020 research and innovation programme (grant agreement No. 648635) and the Cluster of Excellence Cognitive Interaction Technology "CITEC" (EXC 277) at Bielefeld University, which is funded by the German Research Foundation (DFG). Corresponding author: D. Ielmini (email: daniele.ielmini@polimi.it).

within the memory array. For instance, IMC can efficiently accelerate the typical multiply-accumulate (MAC) operation, which is the foundation for modern digital accelerators for artificial intelligence (AI) and optimization [2]. Emerging memory devices such as phase change memory (PCM) [3], [4] and resistive random-access memory (RRAM) [5], [6], offer scalable, efficient and CMOS-compatible solutions to store analogue information as the conductance value. Several IMC demonstrators have thus been reported for accelerating neural network training [7], [8], inference [9], image processing [10] and the solution of algebraic problems [11]–[14].

In a constraint satisfaction problem (CSP), the objective is to find a set of states satisfying a collection of constraints. Typical CSPs include Max-SAT, Max-Cut, graph coloring and the Sudoku puzzle [15]. The latter is indicated as an NP-complete problem in its generic form where the time complexity for solving the problems rapidly increases with the size [16]. CSPs can be implemented by Hopfield recurrent neural networks (RNNs) where the constraints are mapped with synaptic weights, whereas the electrical stimulation allows to minimize the energy cost function to find the solution of the optimization problem [17], [18]. Note that the solution of a CSP in a Hopfield RNN becomes increasingly difficult when the number of the local minima increases, because the network state can be trapped within a local minimum [17]. To circumvent this limitation, the stochastic computational annealing is generally adopted, where the external stimulation is suitably mixed with random noise to help the system escape from local minima [19], [20]. Various solutions have been proposed for practical hardware implementation of computational annealing with CMOS circuits [22]–[26], FPGA [27], quantum computing [28], [29], photonic computing [30] and IMC [31]–[33]. The IMC implementation facilitates two key operations in the computational annealing, namely: (i) the matrix-vector multiplication (MVM) among neuron output signals and synaptic weights which is accelerated in the crosspoint memory array [2], and (ii) the random network stimulation, which can take advantage of the stochastic memory behaviors such as random telegraph noise (RTN) [34] and 1/f noise [35]. While the intrinsic memory noise has been shown to be fruitfully adopted as entropy source for hardware-based true random

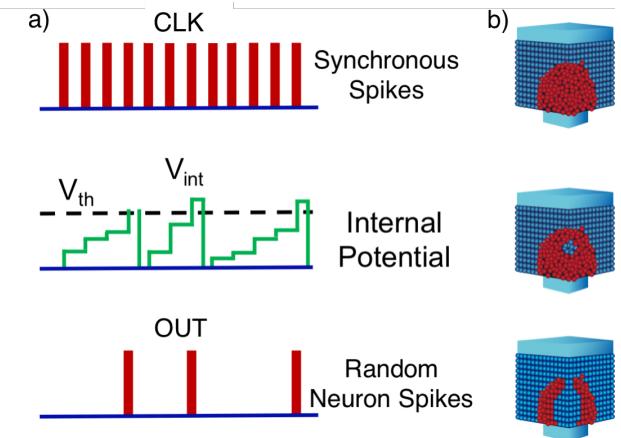


Fig. 1. (a) Operation principle of the proposed stochastic neuron. A train of synchronous spikes applied to the neuron leads to stochastic gradual increase of the internal potential V_{int} . When V_{int} reaches a threshold V_{th} , a spike is generated and V_{int} is restored to zero. (b) Sketch of the gradual crystallization process in PCM devices.

number generators (TRNGs) [36], [37], the same approach is not easily applicable to computational annealing, especially where a fine control of the annealing temperature is needed for dynamic cooling [38]. In fact resistive memory devices suffer from resistance broadening [35], namely the spread of read noise increases with time, which makes the control of stochasticity less controllable. The adoption of a non-physical, pseudo-random number generator (PRNG) such as the linear-feedback shift register (LFSR) was previously proposed for providing the entropy source in stochastic annealing [40]. However, a physics-based entropy source such the PCM can provide true, tunable stochastic input with higher quality of the random noise [41]. Tunable stochastic properties of the memory device, such as the stochastic switching [2], [39], [41], [43]–[45], may also be explored to solve CSPs with an IMC approach.

In this work, we propose a Hopfield RNN for computational annealing based on stochastic spiking neurons, in analogy with the biological brain [46]. The PCM device acts as the source of noise for generating random spikes [35]. First, we show an experimental demonstration of a PCM-based stochastic neuron with tunable output frequency of the generated spikes. After characterizing the integrating neuron element, we implement a Hopfield RNN with PCM synapses [48]. The stochastic RNN is demonstrated for the solution of a 2×2 Sudoku puzzle in hardware. Finally, the convergence of the solution for various annealing algorithms and puzzle sizes up to 16×16 is studied by simulations to allow for the comparison with other types of hardware Sudoku solvers.

II. STOCHASTIC NEURON

A stochastic spiking neuron can act as computational primitive for solving complex CSP problems. Figure 1(a) illustrates the operation concept of the proposed stochastic spiking neuron. A deterministic train of spikes of frequency f_{clk} (top) is accumulated by the neuron. The membrane potential V_{int} , representing the input integral, is stored as a suitable state

variable of the neuron device, such as the device conductance in the case of the PCM (center). As a threshold potential V_{th} is reached, the neuron releases a spike (bottom) while the membrane potential is reset to zero to reinitialize the integration process. Thanks to the stochastic integration of the memory device [43], where the state variable update is affected by variations, the output spikes are randomly generated in time, thus providing the fundamental basis of the stochastic neuron.

A. Stochastic PCM crystallization

The neuron integration function was implemented by a PCM device, where each applied pulse causes a partial crystallization in the amorphous volume. Among the 2-terminal nonvolatile memories, the PCM is one of the most promising concepts thanks to many ideal properties, including high switching speed, low current operation and tunable analogue resistance [3], [4]. PCM devices exhibit two resistance states, associated to the crystalline and amorphous phases of the chalcogenide active material, e.g., $\text{Ge}_2\text{Sb}_2\text{Te}_5$, or GST. To amorphize the GST, a reset voltage pulse is applied above the melting voltage V_m , thus leading to a transition to the liquid phase, followed by a rapid freezing into the amorphous phase, corresponding to the high resistance state (HRS) or reset state. To crystallize the GST, a set voltage pulse is applied usually below V_m and above the threshold voltage V_t for threshold switching [50]. The set pulse causes Joule heating and consequent crystallization within the amorphous volume, thus leading to the low resistance state (LRS).

The crystallization process can also be executed gradually by applying a train of voltage pulses, each inducing partial crystallization within the amorphous volume. Figure 1(b) shows a schematic illustration of the gradual crystallization of a PCM device, starting from the HRS (top) corresponding to a complete amorphous phase, to an intermediate state (center) with some material already in the crystalline phase, until the LRS with fully crystalline phase is reached (bottom). Applications of the PCM as analog weight in neural network accelerators has been widely demonstrated [7], [11].

A key issue for analogue conductance update is the statistical variability of the PCM device, where the pulse-induced increase in conductance changes from cycle to cycle due to the stochastic nature of the crystallization process. As a result, a PCM can also be used as stochastic entropy source for a PCM neuron [43]. To study the variation of gradual crystallization dynamics, we characterized a PCM device with one-transistor-one-resistor (1T1R) structure. Figure 2(a) shows the measured conductance of a PCM device as a function of the number of pulses of voltage V_A , normalize with the initial conductance G_0 . The measurement was repeated 100 times, and each time the device was reinitialized in the HRS by a reset pulse to reach the initial conductance G_0 . After an initial incubation phase where the applied pulse causes no change of conductance, G/G_0 steeply increases as a result of the cumulative crystallization within the amorphous volume and eventually saturates to a value G_{sat}/G_0 . The onset of crystallization shows statistical variation in the same device,

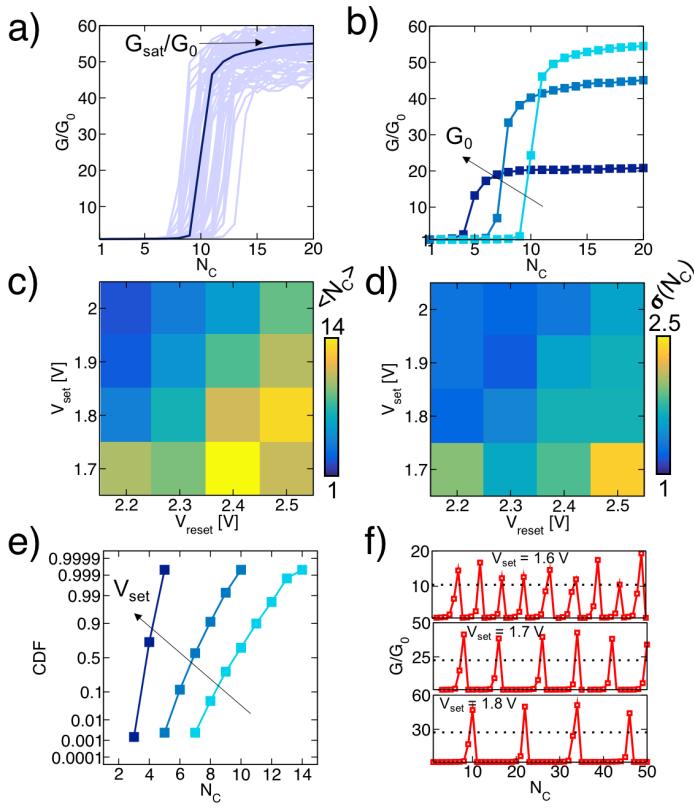


Fig. 2. (a) Relative PCM conductance G/G_0 as a function of N_c for a given set and reset voltage. Both the individual 100 measurements and their average G/G_0 are shown. (b) Average G/G_0 for increasing G_0 , i.e., decreasing initial amorphous volume. As G_0 decreases, the number of pulses in the incubation phase increases (c) Average number of cycles N_c to reach a given conductance threshold $G_{sat}/2$ and (d) its standard deviation $\sigma(N_c)$. (e) Cumulative distribution function of N_c to reach the conductance threshold at increasing V_{set} . (f) Measured conductance change G/G_0 as a function of N_c for increasing V_{set} values namely 1.6 V (top), 1.7 V (center) and 1.8 V (bottom). The conductance is initialized to G_0 every time the threshold (dashed line) is reached, thus resulting in a stochastic train of spikes.

which can be attributed to the stochastic nucleation and growth processes in the amorphous volume [51]. Fig. 2(b) shows the average conductance change G/G_0 for increasing conductance G_0 of the initial HRS as a function of the number of programming pulses. The initial conductance G_0 impacts on all the parameters of the update characteristics, including the incubation number of pulses, the slope of the G increase and the G_{sat} value.

To study the stochastic variations of crystallization, Figures 2(c) and 2(d) show the mean value $\langle N_c \rangle$ and the deviation $\sigma(N_c)$, respectively, of the number of incoming set pulses N_c to reach a threshold conductance $G_{th} = G_{sat}/2$ as function of the applied V_{set} and the pre-programming V_{reset} pulse to reach the different desired values of G_0 . The average number of pulses to crystallization decreases with V_{set} , since a higher set voltage induces a more abrupt crystallization. Conversely, the number of pulses to crystallization increases with increasing V_{reset} , because of the larger initial amorphous volume that needs to be crystallized. Similarly, the deviation of the number of pulses increases with increasing V_{reset} and decreasing V_{set} , which correctly tracks the behavior of the average number of

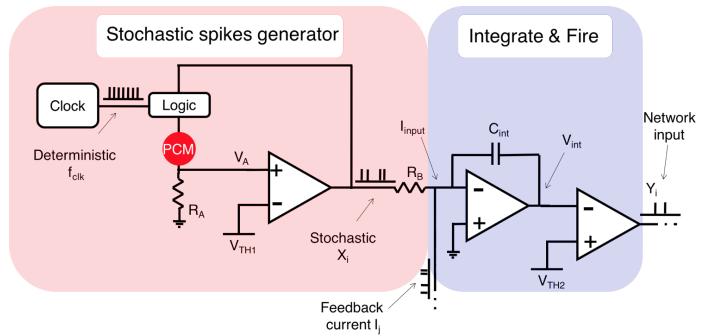


Fig. 3. Schematic of the neuron circuit with the stochastic spike generator (left) and the I&F unit (right). Input deterministic spikes with clock frequency f_{clk} stimulate the PCM device which is connected in series with a fixed resistance R_A . As the voltage across R_A exceeds the comparator threshold V_{TH1} , a stochastic spike is emitted. The output spikes X_i stimulate a current across resistance R_B , which is summed with the feedback current from the RNN and integrated by the I&F unit. As the internal potential V_{int} exceeds the second comparator threshold V_{TH2} , the I&F unit generated a spike Y_i , which is then propagated within the synaptic network of the RNN.

pulses.

The statistical analysis is summarized in the cumulative distributions of the number of crystallizing pulses at increasing V_{set} for a fixed $V_{reset} = 2.4$ V in Fig. 2(e). All distributions show a Gaussian-like behavior. Note that the standard deviation of N_c in Fig. 2(e) controls the statistics of the output spikes of the stochastic neuron in Figure 1(a), hence the annealing dynamics in the RNN. Therefore, the ability to tune the average number and spread of N_c in Figure 2(e) is deeply beneficial for the hardware solution of CSPs. Figure 2(f) shows the measured G/G_0 as function of the number of pulses for $V_{set} = 1.6$ V (top), $V_{set} = 1.7$ V (center) and 1.8 V (bottom). To reproduce the response of the integrate and fire (I&F) neuron, the PCM device was reset every time G/G_0 exceeded the threshold value, which is indicated as a dashed line. It is possible to observe the variation of the number of pulses between every fire event, which supports the stochastic behavior of the PCM neuron. Note that PCM device offers the unique physical property of stochastic integration of Fig. 2, which would not be equally feasible in other types of memory device, such as resistive switching memory or magnetic spin torque memory.

B. Stochastic neuron circuit

Figure 3 shows a schematic illustration of the stochastic neuron circuit which includes two stages, namely (i) a stochastic spike generator, based on a PCM stochastic seed, and (ii) an I&F output stage. A deterministic train of spikes with clock frequency f_{clk} is applied across the PCM in series with a load resistance R_A , thus acting as a voltage divider. As the PCM conductance increases from the initial value G_0 because of gradual crystallization, the voltage divider output V_A evaluated at the PCM bottom electrode (BE) between every input spike increases. The voltage V_A is compared with the threshold V_{TH1} of a comparator. As V_A reaches V_{TH1} , an output spike is generated and the PCM conductance is restored to G_0 by a feedback signal. The correct voltage applied to

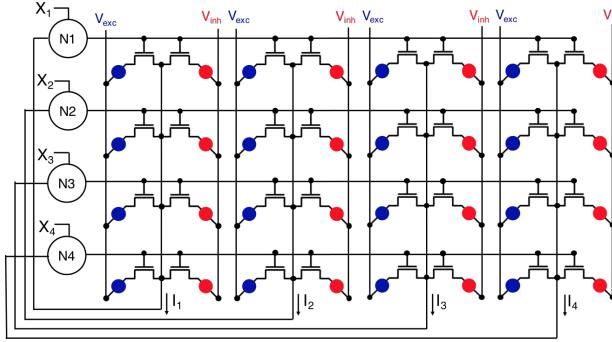


Fig. 4. Schematic illustration of the Hopfield RNN. Neurons N_i represent the I&F units of the stochastic neuron of Figure 3. The input X_i is a stochastic signal generated by the stochastic spike generator of Figure 3. Synapses consist of 1T1R PCM devices organized in excitatory (blue) and inhibitory (red) paths. The neuron output is applied to the gates of the 1T1R in the same row, whereas the synaptic source currents are collected along the same columns and fed back to the neurons. The TEs of the excitatory and inhibitory synapses are biased at V_{exc} and V_{inh} , respectively.

the PCM device is synchronized by a control logic circuit. The stochastic voltage spikes X_i of average frequency f_{input} , are then converted into a spiking current by the resistor R_B and summed with the feedback column current spikes I_j collected from the RNN. The total current is then integrated on a capacitor C_{int} , thus causing the membrane potential V_{int} to increase spike after spike. As V_{int} reaches the threshold V_{TH2} of the second comparator, a spike is generated and applied to the i -th row of the RNN. The neuron response was implemented by a Monte Carlo (MC) model of the stochastic PCM device using the parameters extracted from Figure 2, and implemented as a stochastic primitive for solving CSP problems. The MC simulations were carried on a Matlab environment.

III. HARDWARE RNN

Figure 4 shows the hardware implementation of a Hopfield RNN with PCM-based synapses and neurons [48], [49]. Each neuron N_i represents the I&F unit of Figure 3, while the stochastic unit to generate stimulating signal X_i is not shown. The input stimulation to every neuron N_i is thus the stochastic signal X_i generated by the stochastic spike generator block of Figure 3. Each synaptic unit consists of two PCM devices with 1T1R structure, acting as the excitatory synapse and the inhibitory synapse, respectively. In each synaptic element, the gate terminals of the excitatory and inhibitory synapses are tied together, as well as shared with all other synaptic elements along the same row in the RNN. The source terminals are also connected and shared among all the synaptic elements in the same column of the RNN. Finally, the top electrodes (TEs) of the excitatory and inhibitory synapses are all biased to a positive read voltage V_{exc} and a negative read voltage V_{inh} , respectively, to induce the corresponding column currents. As a result, the overall synaptic weight G_{ij} can be obtained from the difference between the excitatory conductance G_{ij}^+ and the inhibitory conductance G_{ij}^- , according to $G_{ij} = G_{ij}^+ - G_{ij}^-$. The neuron output signal Y_i controls the synaptic gates along the row, whereas the source currents along each column are

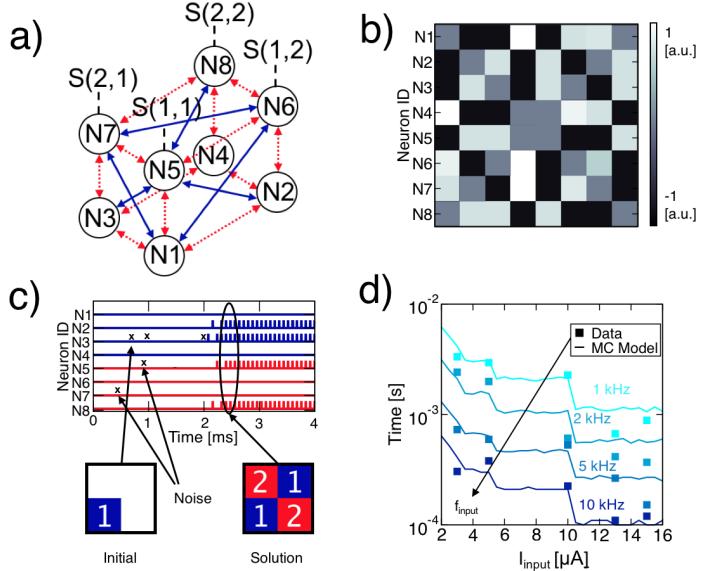


Fig. 5. Schematic illustration of the synaptic connectivity to solve a 2×2 Sudoku puzzle, with inhibitory connections (dashed red arrows) and excitatory connections (solid blue arrows). A 2×2 Sudoku solver needs $N^3 = 8$ neurons and $N^6 - N^3 = 56$ synapses connecting every neuron to each other (without self-connection) (b) Map of synaptic conductance for a 2×2 Sudoku programmed in a PCM array. (c) Experimental solution of a 2×2 Sudoku, including spikes X_j of the external stimulation, corresponding to the initial condition of number '2' in position (2,1), and spikes Y_j , reflecting the spiking activity of each neuron in the RNN. The solution is achieved after about 2 ms of stimulation. (d) Experiments (squares) and MC simulations (lines) of the solution time for a 2×2 Sudoku puzzle as a function of the stimulation amplitude I_{input} and average frequency f_{input} .

collected and applied back to the neurons for integration. For instance, N_1 controls all gates in the first row of the RNN while the synaptic currents in the first column are all collected by Kirchhoff's law, forming the internal signal I_1 , which is applied back to N_1 . The synaptic current I_j induced by the output spike voltage V_i of neuron N_i is given by $I_j = \sum_i G_{ij} V_i$, thus accelerating the physical MVM within the RNN by IMC. Note that, according to the Hopfield topology of the RNN, synapses in all diagonal positions are omitted to prevent self-excitation/inhibition in any neuron.

IV. HARDWARE SOLUTION OF A SUDOKU PUZZLE

To implement a certain problem with the RNN, the constraints need to be correctly mapped in the synaptic weights, i.e., the conductance values G_{ij} . Considering a Sudoku puzzle of size N , the constraints can be mapped in N layers of $N \times N$ neurons, where each neuron corresponds to a certain number in a certain position of the puzzle (e.g. number '1' in position $S(1,1)$). Each layer corresponds to a possible number in the puzzle, e.g. '1' or '2' for the 2×2 Sudoku. The neurons can be thus rearranged in a $N \times N \times N$ matrix [23], where the entry corresponding to '1' indicates a firing neuron and the entry corresponding to '0' indicates a silent neuron. Figure 5(a) shows the constraints for a 2×2 Sudoku problem, represented by 2 layers of 4 neurons each. Every neuron represents the neuron circuit of Figure 3. Solid lines indicate excitatory connections, where a number in a certain position

is exciting the same number in a different row/column, or the other number in the same row/column. Dashed lines instead indicate inhibitory connections, where a number inhibits the same number in the same row/column, or the other number in the same position. In larger Sudoku puzzles, there are also excitatory connections from any neuron to any possible other neuron that does not violate the constraints. For example, in a regular size Sudoku (with $N = 9$), the number '1' will excite any neuron on the same row/column with the numbers '2'-'9'. Figure 5(b) shows the conductance map for a 2×2 Sudoku, which was implemented in two 8×8 PCM arrays of Figure 4, one for excitatory and one for inhibitory synapses. In this RNN, the neuron spikes are applied to its corresponding row, while the currents are collected on the columns and fed back according to the schematic in Figure 4.

A. Experimental solution of a Sudoku puzzle

We carried out experiments for the solution of a 2×2 Sudoku puzzle with the stochastic spiking RNN. The 2×2 Sudoku solution can only contain number 1 and 2, each appearing only once in each row/column, i.e., only solutions (1,2;2,1) and (2,1;1,2) are possible. Figure 5(c) shows the measured train of spikes X_j and Y_j , namely the stochastic stimulation and the neuron spiking output, respectively, in Figure 4. An initial guess $S(2,1) = 1$ is given as external stimulation, corresponding to neuron N_3 being externally stimulated by a stochastic spiking train X_3 at relatively-high average frequency, whereas all other neurons are only subject to random spikes at lower average frequency. Random spiking in non-stimulated neurons is necessary to prevent trapping in a local minimum of the cost function in the RNN. Note that a stochastic implementation is not necessary to solve a 2×2 Sudoku, however the simplicity of the Sudoku puzzle allows to clearly illustrate the solution algorithm and the spiking signals in Fig. 5(c). The stochastic spikes were generated with the MC model by assuming a stimulation of the PCM with a deterministic train of spikes, where f_{clk} and $f_{clk}/10$ were used to generate the stochastic stimulation of average frequency f_{input} and $f_{input}/10$, and then uploaded on a microcontroller (μ C). The latter was programmed to serve as I&F neuron, with the output connected to the gates of 1T1R PCM synapses. The synaptic currents were collected, converted into voltage signals by a transimpedance amplifier (TIA), digitalized with an analog to digital converter (ADC) and fed back into the μ C. The system was temporized by a clock with frequency $f_{clk} = 10 \text{ kHz}$, which also limits the maximum f_{input} . The experimental results show that after about 2 ms , neurons N_2 , N_3 , N_5 and N_8 start to regularly fire at high frequency whereas all other neurons remain silent. This corresponds to a stable attractor [48] of the RNN and to the minimum of the cost function, thus yielding the hardware solution of the Sudoku puzzle. The configuration of spiking neurons was sustained even after the external stimuli have been removed, thus indicating the stability of the attractor state.

The solution can be easily accelerated by increasing the clock frequency and the stimulating currents. This is shown in Figure 5(d), which reports the computing time to solve Sudoku as a

function of the input current I_{input} of the stimulating stochastic spikes X_i for increasing spiking frequency f_{input} . Data points represent the average over 5 experiments conducted on our RNN. The solution becomes faster as I_{input} and f_{input} increase, since the activated neurons generate random spikes at higher average frequency. It should also be noted that the solution can be further accelerated by optimizing the neuron threshold, which was set to $V_{th} = 1 \text{ V}$ in the experiment of Fig. 5c, considering an equivalent capacitor of $C_{int} = 100 \text{ pF}$. Simulation results from a MC model of the network, including PCM variability in neurons and synapses, are also shown in Fig. 5(d). The simulation results clearly show a staircase behavior of the computing time, where the step change corresponds to I_{input} being a submultiple of $I_{th} = V_{th} \cdot C_{int} \cdot f_{clk}$, steps are in fact clearly visible in correspondence of $I_{input} = 10 \mu\text{A}$ or $I_{input} = 5 \mu\text{A}$.

V. TEMPERATURE OPTIMIZATION

With the developed MC model, we simulated various Sudoku problems with increasing size from 4 to 16, to study the RNN performance and the correct tuning of the random spikes, which can be viewed as an equivalent temperature in the simulated annealing process to reach the global minimum of the cost function. To properly solve the puzzles, we assumed an RNN with N^3 neurons and $N^6 - N^3$ PCM synapses encoding all constraints of the Sudoku. We then ran MC simulations to evaluate the success probability P_{sol} , namely the probability of reaching the right solution, for a fixed iteration number and variable input and noise frequency to study the optimal tuning of the stochastic neurons. Figure 6 shows the calculated P_{sol} for increasing size, namely (a) 6×6 , (b) 9×9 , (c) 12×12 and (d) 15×15 . P_{sol} is shown as a function of the probability P_{input} of generating an input spike, namely the ratio between the number of stimulating stochastic spikes X_j and the number of deterministic spikes of frequency f_{clk} in Fig. 3, and the probability P_{noise} of generating a noise spike, which is defined similarly to P_{input} but referred to random noise spikes. Each simulation was run for 1000 cycles and was repeated for 100 times. The position of the maximum P_{sol} moves to lower P_{input} and higher P_{noise} for increasing N , thus indicating an increasing need for stochasticity for increasing Sudoku size. This can be explained by the number of local minima increasing with N , thus resulting in a higher noise contribution, hence temperature, to prevent trapping within a local minimum. On the other hand, an excessive temperature may instead lead to an unstable result, where the RNN can escape also from the global minimum. Note that this method combines the stochastic spike timing and the stochastic PCM conductance variations as entropy sources of the stochastic annealing, whereas previous works only considered stochastic conductance variations [31]–[33]. The simulation results suggest that a tunable source of stochastic spikes is essential for an efficient solution of CSPs in hardware.

VI. EFFICIENCY AND SCALING

To study the impact of Sudoku size on CSP complexity, we ran MC simulations to evaluate the error probability

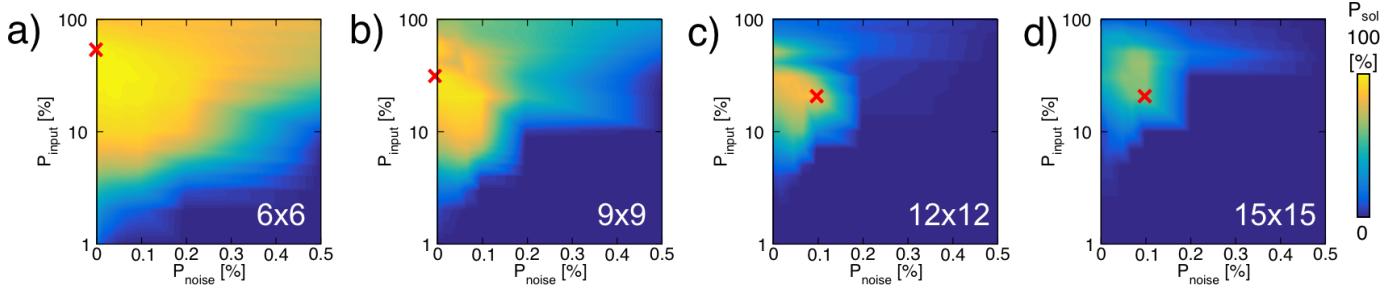


Fig. 6. Probability P_{sol} to solve a Sudoku puzzles for increasing size, namely (a) 6×6 , (b) 9×9 , (c) 12×12 and (d) 15×15 , as a function of the probability P_{input} of generating an input spike and the probability P_{noise} of generating a noise spike. The maximum P_{sol} is marked, indicating that more stochasticity is needed to solve the problem at increasing N .

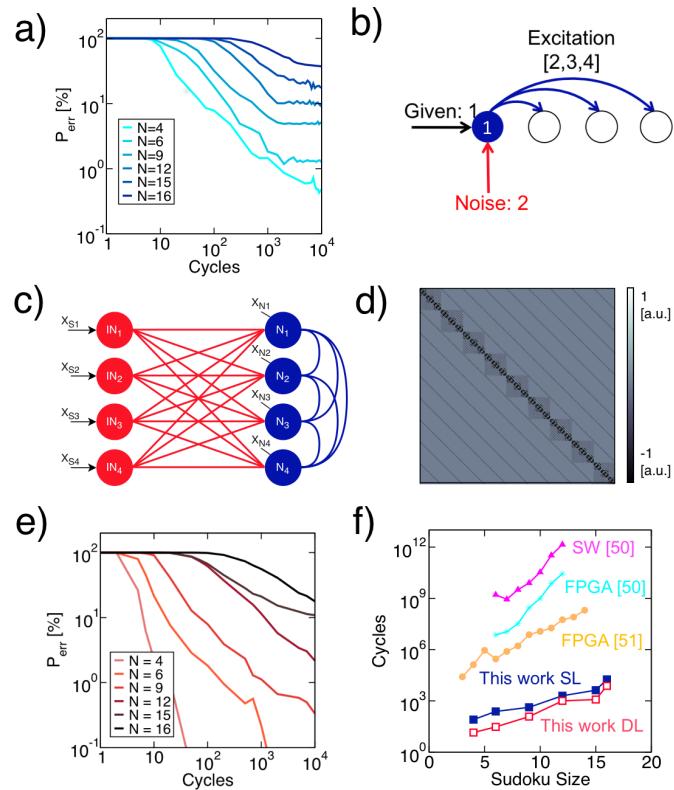


Fig. 7. (a) Error probability P_{err} as a function of the number of cycles for increasing Sudoku size. (b) Schematic illustration of the possible conflict while solving a Sudoku puzzle, where neurons can be excited and inhibited at the same time by spiking neurons. (c) Schematic of the double-layer network to prevent conflicting excitation/inhibition and improve the convergence. (d) Matrix of the synaptic weights of the feedforward input layer. (e) P_{err} as a function of the number of cycles for increasing Sudoku size for the double-layer network. (f) Performance of the Sudoku solver, namely number of iteration cycles for the solution as a function of size, for the single-layer RNN and the double-layer RNN, compared with other solvers from the literature.

$P_{err} = 1 - P_{sol}$ for increasing N . Figure 7(a) shows the calculated error probability P_{err} as a function of the number of computing cycles for increasing N between 4 and 16. The error probability increases with N for a given computational cycle. Conversely, the computing speed to solve the Sudoku problem with sufficiently low P_{err} decreases for increasing N . While the system becomes more unreliable for bigger problems, computation can be parallelized on more than one

memory array to enhance the probability of reaching a correct solution [33].

The algorithm can also be improved to reduce P_{err} and accelerate the annealing process by properly taking into account conflicts among constraints at the hardware level. For instance, Figure 7(b) shows possible conflicts within the first row of a 4×4 Sudoku. An initial condition, namely $S(1,1) = 1$ promotes with excitatory synapses all the other numbers on the same row, namely $S(1,j) = [2, 3, 4]$, for $j \neq 1$. At the same time a random noise spike could activate the neuron coding for digit 2 on the same cell, thus conflicting with initial condition and leading to escape from the correct global minimum [18]. To avoid this, the neurons coding for digits directly conflicting with initial condition should be inhibited from firing. This can be achieved by properly transforming the initial condition such that an excitatory stimulation is provided to neuron coding for givens, whereas an inhibitory stimulation is provided for neurons coding for digit conflicting with the givens following the Sudoku constraints.

To this purpose, we propose a double-layer network (DL) as shown in Figure 7(c), where a feedforward layer is added as input layer to the RNN for filtering the input conditions and inhibit wrong stimulations. Noise is only injected in the second layer, while the inhibitions given by the input layer also control the annealing temperature by acting as a cooling effect when the correct solution is reached. In fact, if noise stimulates a neuron N_i which is in contrast with the input condition, the first layer will inject a negative current to N_i to prevent its activation. Figure 7(d) shows the synaptic weights of the input layer for a 9×9 Sudoku indicating all the inhibitions to the recurrent layer.

Figure 7(e) shows the error probability P_{err} as a function of the number of cycles for increasing size of the Sudoku size by adopting the DL network of Figure 7(c). Figure 7(f) summarizes the computing speed, evaluated as the number of cycles to reach $P_{err} = 1\%$, for the single-layer (SL) RNN and the DL network, indicating that the computing speed is clearly improved by the DL network. The performance of the RNN is compared with state-of-the-art systems for solving Sudoku with FPGA implementations [52], [53] and software approaches [52]. The results indicate that the IMC approach allows to accelerate the solution of CSP by about 4 orders of magnitude compared with FPGA and 7 orders of magnitude

compared with software-based solvers. This is due to (i) the compact implementation of stochastic spikes generator and (ii) the low latency MAC operation of in-memory computing compared with other techniques [54]. Moreover, the novel DL network further improves the performance of the CSP solver and takes full advantage of the compact MAC core. Our implementation shows reduced number of cycles to solution also compared with state-of-the-art analog neuromorphic processor [26] which takes about 50 cycles to solve a 4×4 Sudoku, compared with just 14 cycles of our DL network. These results support the use of the PCM technology for computational annealing techniques.

VII. CONCLUSION

We have developed a brain-inspired spiking RNN for solving CSP problems with stochastic PCM neurons and PCM synapses. First, the stochasticity behavior of the PCM device was experimentally studied during gradual crystallization. Then, the PCM-based stochastic neuron was implemented in a RNN for solving Sudoku puzzles, which were experimentally validated on a small scale (2×2). A MC model was developed to describe the network stochastic behavior and predict multiple experimental results as a function of stimulating current and frequency. The model was then used to scale the system and study the error probability as a function of the problem size. Finally, a DL spiking neural network was designed to further reduce the error probability. The results demonstrate the superior performance of our system compared with the state-of-art implementations, confirming IMC as a promising approach to accelerate the hardware solution of CSPs.

REFERENCES

- [1] W. A. Wulf and S. A. McKee "Hitting the memory wall: implications of the obvious", *ACM SIGARCH Computer Architecture News* 23, pp. 20-24, 1995 doi: 10.1145/216585.216588.
- [2] D. Ielmini and H.-S. P. Wong, "In-memory computing with resistive switching devices, *Nat Electron* 1, pp. 333-343, 2018 doi:10.1038/s41928-018-0092-2.
- [3] S. Raoux, W. Welnic and D. Ielmini, "Phase change materials and their application to non-volatile memories", *Chem. Rev.* 110, pp. 240-267, 2010 doi: 10.1021/cr900040x.
- [4] G. W. Burr, et al., "Phase change memory technology," *Journal of Vacuum Science & Technology B* 28, 223, 2010 doi: 10.1116/1.3301579.
- [5] H. -S. P. Wong et al., "Metal-Oxide RRAM," *Proceedings of the IEEE*, 100 (6), pp. 1951-1970, 2012 doi: 10.1109/JPROC.2012.2190369.
- [6] D. Ielmini, "Resistive Switching Memories based on Metal Oxides: Mechanisms, Reliability and Scaling," *Semicond. Sci. Technol.* 31, 063002 2016 doi: 10.1088/0268-1242/31/6/063002.
- [7] S. Ambrogio, et al., "Equivalent-accuracy accelerated neural-network training using analogue memory", *Nature* 558, pp. 60-67, 2018 doi:10.1038/s41586-018-0180-5.
- [8] C. Li, et al., "Efficient and self-adaptive in-situ learning in multilayer memristor neural networks", *Nat Commun* 9, (2385), 2018 doi:10.1038/s41467-018-04484-2.
- [9] V. Milo, et al., "Multilevel HfO₂-based RRAM devices for low-power neuromorphic networks", *APL Materials* 7 (081120), 2019 doi:10.1063/1.5108650.
- [10] C. Li, et al., "Analogue signal and image processing with large memristor crossbars", *Nat Electron* 1, pp. 52-59, 2018 doi:10.1038/s41928-017-0002-z.
- [11] M. Le Gallo, et al., "Mixed-precision in-memory computing" *Nat Electron* 1, pp. 246-253, 2018 doi:10.1038/s41928-018-0054-8.
- [12] M. A. Zidan, et al., "A general memristor-based partial differential equation solver", *Nat Electron* 1, pp. 411-420, 2018 doi:10.1038/s41928-018-0100-6.
- [13] Z. Sun, et al., "Solving matrix equations in one step with cross-point resistive arrays", *PNAS* 116 (10), pp. 4123-4128, 2019 doi:10.1073/pnas.1815682116.
- [14] Z. Sun, G. Pedretti, A. Bricalli and D. Ielmini, "One-step regression and classification with cross-point resistive memory arrays", *Science Advances* 6: eaay2378, 2019 doi:10.1126/sciadv.aay2378.
- [15] M. R. Garey, "A guide to the theory of NP-completeness", *Computers and Intractability, W. H. Freeman and Company*, 1979 ISBN:978-0-7167-1045-5.
- [16] H. Simonis, "Sudoku as a constraint problem", *CP Workshop on Modeling and Reformulating Constraint Satisfaction Problems*, pp. 13-28, 2005.
- [17] J. J. Hopfield, "Searching for memories, sudoku, implicit check bits, and the iterative use of not-always-correct rapid neural computation", *Neural Computation*, vol. 20, no. 5, pp. 1119-1164, 2008 doi:10.1162/neco.2007.09-06-345.
- [18] S. Habenschuss, Z. Jonke, and W. Maass, "Stochastic computations in cortical microcircuit models", *PLOS Computational Biology*, vol. 9, pp. 1-28, 11 2013 doi:10.1371/journal.pcbi.1003311
- [19] E. H. L. Aarts and J. H. M Korst, "Simulated annealing and Boltzmann machines" *John Wiley*, 1988
- [20] V. Pavlovic, D. Schonfeld and G. Friedman , "Enhancement of Hopfield neural networks using stochastic noise processes", *Proc. Neural Networks for Signal Processing Soc. Workshop*, pp. 173-182, 2001 doi:10.1109/NNSP.2001.943122.
- [21] S. Kirkpatrick, C. D. Gelatt Jr. and M. P. Vecchi, "Optimization by Simulated Annealing", *Science* 220 (4598), pp. 671-680, 1983 doi: 10.1126/science.220.4598.671.
- [22] G. A. Fonseca Guerra and S. B. Furber, "Using stochastic spiking neural networks on SpiNNaker to solve constraint satisfaction problems", *Front. Neurosci.* 11:714, 2017 doi:10.3389/fnins.2017.00714.
- [23] J. Binas, G. Indiveri and M. Pfeiffer, "Spiking analog VLSI neuron assemblies as constraint satisfaction problem solvers", *IEEE International Symposium on Circuits and Systems (ISCAS)*, Montreal, QC, pp. 2094-2097, 2016 doi: 10.1109/ISCAS.2016.7538992.
- [24] H. Mostafa, L. Müller and G. Indiveri, "An event-based architecture for solving constraint satisfaction problems", *Nat Commun* 6, 8941, 2015 doi:10.1038/ncomms9941.
- [25] T. Takemoto, M. Hayashi, C. Yoshimura and M. Yamaoka, "2.6 A $2 \times 30k$ -Spin Multichip Scalable Annealing Processor Based on a Processing-In-Memory Approach for Solving Large-Scale Combinatorial Optimization Problems", *IEEE International Solid- State Circuits Conference - (ISSCC)*, San Francisco, CA, USA, pp. 52-54, 2019 doi: 10.1109/ISSCC.2019.8662517.
- [26] D. Liang and G. Indiveri, "A Neuromorphic Computational Primitive for Robust Context-Dependent Decision Making and Context-Dependent Stochastic Computation", *IEEE Trans. Circuits Syst. II*, vol. 66, no. 5, pp. 843-847, 2019 doi: 10.1109/TCSII.2019.2907848.
- [27] F. L. Traversa, C. Ramella, F. Bonani and M. Di Ventra, "Memcomputing NP-complete problems in polynomial time using polynomial resources and collective states", *Science Advances* 1 (6), e1500031, 2015 doi: 10.1126/sciadv.1500031.
- [28] S. Boixo, et al., "Evidence for quantum annealing with more than one hundred qubits", *Nature Phys* 10, 218-224, 2014 doi:10.1038/nphys2900.
- [29] V. S. Denchev, et. al., "What is the Computational Value of Finite-Range Tunneling?", *Phys. Rev. X* 6, 031015, 2016 doi: 10.1103/PhysRevX.6.031015.
- [30] R. Hamerly, et al., "Experimental investigation of performance differences between coherent Ising machines and a quantum annealer", *Science Advances* 5 (5), eaau0823, 2019 doi: 10.1126/sciadv.aau0823.
- [31] J. H. Shin, Y. J. Jeong, M. A. Zidan, Q. Wang and W. D. Lu, "Hardware Acceleration of Simulated Annealing of Spin Glass by RRAM Crossbar Array", *IEEE International Electron Devices Meeting (IEDM)*, San Francisco, CA, pp. 3.3.1-3.3.4, 2018 doi: 10.1109/IEDM.2018.8614698.
- [32] M. R. Mahmoodi, et al. , "An Analog Neuro-Optimizer with Adaptable Annealing Based on 64×64 0T1R Crossbar Circuit", *IEEE International Electron Devices Meeting (IEDM)*, San Francisco, CA, pp. 14.7.1-14.7.4, 2019 doi: 10.1109/IEDM19573.2019.8993442.
- [33] F. Cai, et al. , "Harnessing intrinsic noise in memristor Hopfield neural networks for combinatorial optimization". *ArXiv preprint arXiv:1903.11194*.
- [34] S. Ambrogio, et. al., "Statistical Fluctuations in HfO_x Resistive-Switching Memory: Part II - Random Telegraph Noise", *IEEE Transactions on Electron Devices*, vol. 61, no. 8, pp. 2920-2927, 2014 doi: 10.1109/TED.2014.2330202.
- [35] S. Ambrogio, S. Balatti, V. McCaffrey, D. Wang, and D. Ielmini, "Noise-induced resistance broadening in resistive switching memory (RRAM) -

- Part I: Intrinsic cell behavior," *IEEE Trans. Electron Devices* 62, 3805-3811, 2015 doi: 10.1109/TED.2015.2475598.
- [36] Huang, C.-Y., Shen, W. C., Tseng, Y.-H., King, Y.-C. and Lin, C.-J., "A contact-resistive random-access-memory-based true random number generator," *IEEE Electron Device Lett.* 33, 1108-1110, 2012 doi: 10.1109/LED.2012.2199734.
- [37] Z. Wei, et. al., "True random number generator using current difference based on a fractional stochastic model in 40-nm embedded ReRAM," *2016 IEEE International Electron Devices Meeting (IEDM)*, San Francisco, CA, pp. 4.8.1-4.8.4, 2016 doi: 10.1109/IEDM.2016.7838349.
- [38] Y. Nourani and B. Andresen, "A comparison of simulated annealing cooling strategies" *J. Phys. A: Math. Gen.* 31 8373, 1998 doi:10.1088/0305-4470/31/41/011.
- [39] S. Ambrogio, et. al., "Statistical Fluctuations in HfO_x Resistive-Switching Memory: Part I - Set/Reset variability", *IEEE Transactions on Electron Devices*, vol. 61, no. 8, pp. 2920-2927, 2014 doi: 10.1109/TED.2014.2330202.
- [40] G. Cauwenberghs, "An analog VLSI recurrent neural network learning a continuous- time trajectory," *IEEE TNN* 7(20), 346-361, 1996 doi: 10.1109/72.485671
- [41] R. Carbone and D. Ielmini, "Stochastic Memory Devices for Security and Computing", *Adv. Electron. Mater.* 5 1900198, 2019 doi: 10.1002aelm.201900198.
- [42] S. Kumar, J.P. Strachan and R.S. Williams, "Chaotic dynamics in nanoscale NbO₂ Mott memristors for analogue computing", *Nature* 548, pp. 318-321, 2017 doi:10.1038/nature23307.
- [43] T. Tuma, et al., "Stochastic phase-change neurons", *Nature Nanotech* 11, pp. 693-699, 2016 doi:10.1038/nnano.2016.70.
- [44] A. Mizrahi, et al., "Neural-like computing with populations of superparamagnetic basis functions", *Nat Commun* 9, (1533), 2018 doi:10.1038/s41467-018-03963-w.
- [45] W.A. Borders, et al., "Integer factorization using stochastic magnetic tunnel junctions", *Nature* 573, pp. 390-393, 2019 doi:10.1038/s41586-019-1557-9.
- [46] W. Maass, "Noise as a resource for computation and learning in networks of spiking neurons", *Proceedings of IEEE* 102(5), pp. 860-880, 2014 doi: 10.1109/JPROC.2014.2310593.
- [47] A. Redaelli et al., "Impact of the current density increase on reliability in scaled BJT-selected PCM for high-density applications", *IEEE International Reliability Physics Symposium (IRPS)*, Anaheim, CA pp. 615-619, 2010 doi: 10.1109/IRPS.2010.5488760.
- [48] V. Milo, D. Ielmini and E. Chicca, "Attractor networks and associative memories with STDP learning in RRAM synapses," *IEEE International Electron Devices Meeting (IEDM)*, San Francisco, CA, pp. 11.2.1-11.2.4, 2017 doi: 10.1109/IEDM.2017.8268369.
- [49] G. Pedretti, V. Milo, S. Hashemkhani, et.al., "A spiking recurrent neural network with phase change memory synapses for decision making", *IEEE International Symposium on Circuits and Systems (ISCAS)*, Sevilla, Spain, 2020 (in press)
- [50] D. Ielmini, "Threshold switching mechanism by high-field energy gain in the hopping transport of chalcogenide glasses," *Phys. Rev. B* 78, 035308, 2008 doi: 10.1103/PhysRevB.78.035308.
- [51] U. Russo, D. Ielmini, A. Redaelli and A. L. Lacaita, "Intrinsic data retention in nanoscaled phase-change memories - Part I: Monte Carlo model for crystallization and percolation," *IEEE Trans. Electron Devices* 53, pp. 3032-3039, 2006 doi: 10.1109/TED.2006.885527.
- [52] P. Malakonakis, M. Smerdis, E. Sotiriades and A. Dollas, "An FPGA-based Sudoku Solver based on Simulated Annealing methods", *International Conference on Field-Programmable Technology*, Sydney, NSW, 2009, pp. 522-525, 2009 doi: 10.1109/FPT.2009.5377608.
- [53] M. Dittrich, T. B. Preufer and R. G. Spallek, "Solving Sudokus through an incidence matrix on an FPGA", *International Conference on Field-Programmable Technology*, Beijing, pp. 465-469, 2010 doi: 10.1109/FPT.2010.5681460.
- [54] X. Peng, et. al., "Inference engine benchmarking across technological platforms from CMOS to RRAM", in *Proceedings of the International Symposium on Memory Systems - MEMSYS 19*, Washington, District of Columbia, 2019, pp. 471-479 doi: 10.1145/3357526.3357566

research interests include the design of neuromorphic circuits for optimization and analog computing.

Piergiulio Mannocci Piergiulio Mannocci received his B.Sc. degree in electronics engineering from Politecnico di Milano in 2016 where is currently pursuing a M.Sc. in electronics engineering. His research interests include circuit implementation of neuromorphic and analog accelerators with emerging memories.

Shahin Hashemkhani received the B.S. degree in electronics at IAUCTB, Tehran, Iran and M.S. degrees in same major from the Politecnico di Milano, Milan, Italy, in 2014 and 2019, respectively, where he is currently pursuing the Ph.D. degree in electronics engineering. His main research interests are the design and characterization of neuromorphic network.

Valerio Milo (M'19) received the B.S., M.S. and Ph.D. (*cum laude*) in electronics engineering from Politecnico di Milano, Milan, Italy, in 2012, 2015 and 2019, respectively. He is currently a Post-Doctoral researcher at the Dipartimento di Elettronica, Informazione e Bioingegneria of Politecnico di Milano, Milano, Italy. His current research interests include design, modeling, and simulation of neuromorphic networks with resistive switching random access memory (RRAM) and phase change memory (PCM) for neuromorphic computing applications.

Octavian Melnic received the B.S. and M.S. degrees in electrical engineering from the Politecnico di Milano, Milano, Italy, in 2015 and 2018 respectively, where he is currently pursuing the Ph.D. degree in electrical engineering. His current research interests include characterization and modeling of phase change memories.

Elisabetta Chicca (M'06) received the Laurea degree (M.Sc.) in physics from the Università degli Studi di Roma 'La Sapienza', Italy, in 1999, the Ph.D. degree in natural science from the Department of Physics, Swiss Federal Institute of Technology Zurich, and the Ph.D. degree in neuroscience from the Neuroscience Center Zurich in 2006. Since 2017, she has been a Professor at CITEC and Faculty of Technology, Bielefeld University, leading the Neuromorphic Behaving Systems Research Group. Her research focuses on the development of neuromorphic full-custom VLSI models of neural circuits for brain-inspired computation, learning in spiking neural networks, learning in memristive devices and arrays, bio-inspired sensing, and motor control.

Daniele Ielmini (SM'09-F'19) is a Full Professor at the Dipartimento di Elettronica, Informazione, e Bioingegneria of Politecnico di Milano, Politecnico di Milano. He received the Ph.D. degree from Politecnico di Milano in 2000. He conducts research on emerging nanoelectronics devices, such as phase-change memory (PCM) and resistive switching memory (RRAM), and on novel computing with memory devices. Prof. Ielmini was a recipient of the Intel Outstanding Researcher Award in 2013, the ERC Consolidator Grant in 2014, and the IEEE EDS Rappaport Award in 2015.

Giacomo Pedretti (M'20) received the B.S., M.S. and Ph.D. (*cum laude*) in electronics engineering from Politecnico di Milano, Milan, Italy, in 2013, 2016 and 2020, respectively, where he is currently a Postdoc research associate. His