



# Compressed Sensing With Approximate Message Passing Using In-Memory Computing

Manuel Le Gallo<sup>1</sup>, Abu Sebastian<sup>1</sup>, *Senior Member, IEEE*, Giovanni Cherubini<sup>1</sup>, *Fellow, IEEE*, Heiner Giefers, *Senior Member, IEEE*, and Evangelos Eleftheriou, *Fellow, IEEE*

**Abstract**—In-memory computing is a promising non-von Neumann approach where certain computational tasks are performed within resistive memory units by exploiting their physical attributes. In this paper, we propose a new method for fast and robust compressed sensing (CS) of sparse signals with approximate message passing recovery using in-memory computing. The measurement matrix for CS is encoded in the conductance states of resistive memory devices organized in a crossbar array. In this way, the matrix-vector multiplications associated with both the compression and recovery tasks can be performed by the same crossbar array without intermediate data movements at potential  $O(1)$  time complexity. For a signal of size  $N$ , the proposed method achieves a potential  $O(N)$ -fold recovery complexity reduction compared with a standard software approach. We show the array-level robustness of the scheme through large-scale experimental demonstrations using more than 256k phase-change memory devices.

**Index Terms**—Approximate message passing (AMP), compressed sensing (CS), in-memory computing, phase-change memory (PCM).

## I. INTRODUCTION

IN-MEMORY computing is an attractive approach for performing computationally expensive tasks of a high-level algorithm in an energy-efficient manner. For instance, crossbar arrays of resistive memory (memristive) devices can be used to store a matrix and perform analog matrix-vector multiplications at constant  $O(1)$  time complexity without intermediate movements of data. This capability can be exploited in a wide range of applications from neural network inference to solving systems of linear equations [1]–[3].

Another well-suited application domain is that of complex optimization problems such as compressed sensing (CS)

Manuscript received June 15, 2018; revised August 6, 2018; accepted August 8, 2018. Date of publication August 29, 2018; date of current version September 20, 2018. This work was supported in part by the European Research Council through the European Union's Horizon 2020 Research and Innovation Program under Grant 682675 and in part by the European Union's Horizon 2020 Research and Innovation Program through the project MNEMOSENE under Grant 780215. The review of this paper was arranged by Editor C. Monzio Compagnoni. (Corresponding authors: Manuel Le Gallo; Abu Sebastian.)

The authors are with Cloud and Computing Infrastructure, IBM Research Zurich, 8803 Rüschlikon, Switzerland (e-mail: anu@zurich.ibm.com; ase@zurich.ibm.com; cbi@zurich.ibm.com; hgiefers@gmail.com; ele@zurich.ibm.com).

Color versions of one or more of the figures in this paper are available online at <http://ieeexplore.ieee.org>.

Digital Object Identifier 10.1109/TED.2018.2865352

recovery. CS is an active research field in signal processing, which attempts to perform sampling and compression simultaneously via a measurement matrix and allows the recovery of a high-dimensional signal from low-dimensional noisy measurements. CS is used in various applications, such as MRI, facial recognition, holography, audio restoration, and in mobile phone camera sensors. In a camera sensor, the approach allows to significantly reduce the acquisition energy per image or equivalently increase the image frame rate, by capturing only few measurements, e.g., 10%, instead of the whole image. However, CS recovery algorithms are usually complex, and conventional implementations are confronted with limited scalability owing to the large number of operations involved and high memory requirements. In-memory computing promises to significantly reduce the memory and computing resources needed to solve the problem as well as its computational complexity, at the cost of potentially reducing solution accuracy.

In Internet of Things systems, it may be desirable to design implementations of CS with reconstruction on the same device, e.g., a sensor, using very low power, in order to have energy-efficient signal acquisition while at the same time not having to send the compressed signal to the cloud for reconstruction. Moreover, implementations of CS that can deal with very large measurement matrices may be desirable in applications where signals are received by large sensor arrays, as, for example, envisaged for the Square Kilometre Array [4], where the signal size may be on the order of  $10^8$ .

In this paper, we propose an implementation of a CS recovery algorithm, namely approximate message passing (AMP), based on memristive crossbar arrays, of which we presented a preliminary version in [5]. We experimentally investigate the impact of this memristive implementation on the performance of AMP, in particular on the reconstruction accuracy. The benefits and limitations of the memristive implementation are discussed for three use cases of the AMP algorithm, namely linear estimation, CS with soft thresholding, and compressive imaging with image denoising.

## II. OVERVIEW OF COMPRESSED SENSING

### A. Problem Setting

The basic idea of CS is to acquire few sampling measurements from a high-dimensional signal and subsequently to

recover that signal accurately. The compressive measurements can be thought of as a linear mapping of a signal  $x_0$  of length  $N$  to a measurement vector  $y$  of length  $M < N$ . If this process is linear, it can be modeled by an  $M \times N$  measurement matrix  $A$ . The CS reconstruction problem is to determine the signal  $x_0$  from the measurements  $y$  when sampled as

$$y = Ax_0 + w \quad (1)$$

where  $w$  represents the measurement noise. CS asserts that signals can be recovered from fewer samples than dictated by the Shannon–Nyquist theorem if they are sparse, that is, if their information rate is lower than the Nyquist rate. If the signal  $x_0$  is sparse in some transform domains, we can represent it as  $x_0 = \Psi\zeta$ , where  $\zeta$  contains only a few ( $k$ ) nonnegligible elements. It can be shown that if  $\Psi$  is incoherent with  $A$ ,  $\zeta$  can be recovered from  $y$  when  $M < N$ , as long as  $k$  is sufficiently small.  $\Psi$  represents the inverse transform matrix, for example, an inverse wavelet transform. CS is fundamentally different from transform coding, which is used, for example, in JPEG or MPEG compression. In the latter, the signal  $x_0$  needs to be fully acquired, then the transform  $\zeta$  is computed, and the largest  $k$  transform coefficients and their locations are kept so that the signal can be reconstructed. In CS, however, only  $M < N$  measurements of  $x_0$  are acquired while still being able to reconstruct the signal accurately. The downside is the cost of complex CS reconstruction algorithms.

In the case of a sparse signal  $x_0$  and  $w = 0$ , a reconstruction of  $x_0$  from  $y$  is obtained by solving the basis pursuit (BP)  $L_1$  minimization problem. An alternative formulation known as BP denoising (BPDN) extends BP to the more realistic noisy measurement case with  $w \neq 0$ . The solution of both BP and BPDN can be obtained by convex optimization using linear programming (LP) algorithms. However, the high computational complexity of LP represents an obstacle for the large problem sizes that occur very often in applications.

An appealing alternative to LP algorithms is offered by iterative thresholding algorithms because of their low computational complexity. One particular iterative thresholding scheme to recover  $x_0$  from  $y$  is of the form

$$\begin{aligned} x^{t+1} &= \eta_t(A^*z^t + x^t) \\ z^t &= y - Ax^t. \end{aligned} \quad (2)$$

Here,  $x^t \in \mathbb{R}^N$  is the current estimate of  $x_0$  at iteration  $t$ ,  $z^t \in \mathbb{R}^M$  is the current residual,  $\eta_t(\cdot)$  is a (typically nonlinear) function,  $A^*$  denotes the transpose of  $A$ , and  $x^0 = 0$ . However, while offering low complexity, the sparsity-undersampling tradeoff achieved by algorithm (2), that is, the smallest value that  $M$  can take given a certain sparsity of  $x_0$  to successfully recover the signal, is usually less favorable than for LP-based reconstruction.

Recently, Donoho *et al.* [6] proposed an AMP algorithm, which adds a simple modification to (2) that substantially improves the sparsity-undersampling tradeoff without significantly increasing the computational complexity. The AMP algorithm is formulated as [7]

$$\begin{aligned} x^{t+1} &= \eta_t(A^*z^t + x^t) \\ z^t &= y - Ax^t + \frac{N}{M}z^{t-1}\langle \eta'_{t-1}(A^*z^{t-1} + x^{t-1}) \rangle \end{aligned} \quad (3)$$

where  $\langle v \rangle \equiv N^{-1} \sum_{n=1}^N v_n$  denotes the average of a vector  $v$ ,  $\eta'_t$  represents the derivative of  $\eta_t$ ,  $x^t \in \mathbb{R}^N$  is the current estimate of  $x_0$  at iteration  $t$ ,  $z^t \in \mathbb{R}^M$  is the current residual,  $A^*$  denotes the transpose of  $A$ , and  $x^0 = 0$ . With respect to iterative thresholding (2), AMP includes the additional term  $(N/M)z^{t-1}\langle \eta'_{t-1}(A^*z^{t-1} + x^{t-1}) \rangle$  in the computation of the residual, which is shown to substantially improve the sparsity-undersampling tradeoff [6]. AMP has the remarkable property that its solutions are governed by a state evolution whose fixed points (when unique) yield the true posterior means, in the limit  $M, N \rightarrow \infty$ , with the ratio  $M/N$  fixed, and assuming that the elements of  $A$  are independent identically distributed (i.i.d.) Gaussian random variables  $A_{mn} \sim N(0, 1/M)$  [7].

### B. Compressed Sensing Hardware Implementations

Many works have focused on efficient hardware implementations for the acquisition of compressed measurements, such as in a camera sensor [8]–[10]. In an image sensor, the measurement matrix is typically binary, and the measurement acquisition can be implemented either in the optical domain [8] or on-chip [9], [10]. Efficient implementations of single-shot imaging have been demonstrated with scalability up to  $256 \times 256$  pixels consuming less than 100 mW of power and showing no loss in signal-to-noise ratio (SNR) compared with normal (not compressed) capture [10]. In such implementations, the reconstruction algorithm is typically not implemented on-chip, and therefore, reconstruction has to be done off-line.

For CS reconstruction, a number of implementations have been reported on field-programmable gate arrays (FPGAs) and application-specific integrated circuit (ASIC) designs. ASIC implementations of the orthogonal matching pursuit algorithm [11] and the AMP algorithm [12] have been presented, as well as the FPGA implementations of both [13]. Very recently, an implementation of the second-order cone program recovery algorithm for CS based on memristive crossbar arrays has been proposed [14], however, without experimental validation.

In this paper, we propose an implementation of the AMP algorithm based on memristive crossbar arrays, whereby the memristive arrays are used to perform the required matrix-vector multiplications. We aim to provide a robust set of experimental results of this implementation using phase-change memory (PCM) arrays. In comparison with typical high-precision implementations on GPUs or FPGAs, reconstruction with a memristive implementation will exhibit lower accuracy. The expectation is that the energy efficiency and scalability of a memristive implementation will allow to deal with much larger signals than in a typical high-precision implementation and will yield faster and low-power solutions, at the cost of a reduced reconstruction accuracy, which may, however, be considered acceptable in many applications.

## III. REALIZATION USING IN-MEMORY COMPUTING

### A. Implementation of Compressed Sensing With AMP Recovery Using Resistive Memory Arrays

The key idea of realizing CS using in-memory computing relies on the encoding of the elements of  $A$  as conductance

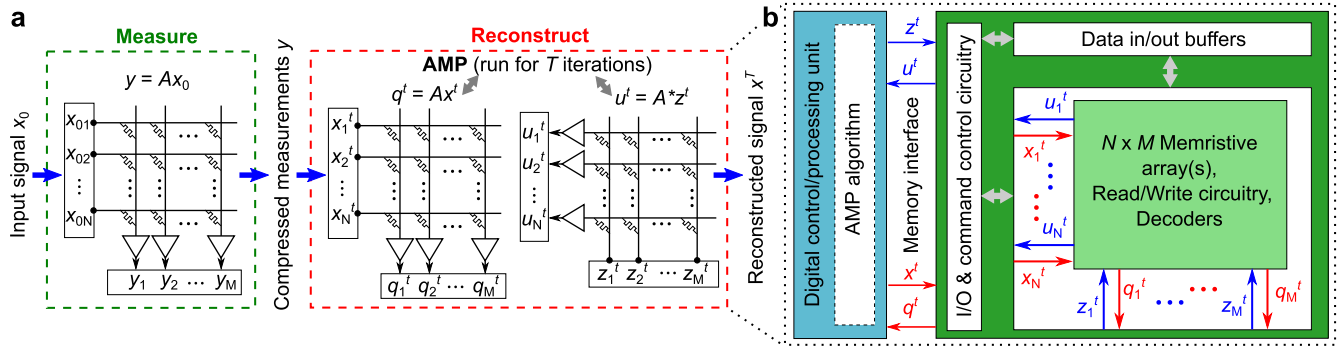


Fig. 1. (a)  $N \times M$  memristive crossbar encoding the measurement matrix  $A$  used to acquire the CS measurements and to realize the matrix-vector computations of the AMP recovery algorithm. (b) Architecture of the memristive implementation of AMP.

values of memristive devices organized in a crossbar array, as shown in Fig. 1(a). One possible method to program the conductance values is by an iterative program-and-verify procedure. The compressed measurements (1) are acquired by applying  $x_0$  as voltages to the crossbar rows via digital-to-analog conversion and obtaining  $y$  through analog-to-digital conversion of the resulting output currents at columns. The positive and negative elements of  $A$  can be coded on separate devices together with a subtraction circuit, whereas negative vector elements can be applied as negative voltages.

Once the matrix  $A$  has been programmed in the crossbar array and the measurements  $y$  have been obtained, the AMP algorithm can be implemented as shown in Fig. 1(b). The AMP algorithm is run in a dedicated processing unit, whereas the computation of  $q^t = Ax^t$  and  $u^t = A^*z^t$  is performed using the (same) crossbar array. The vector  $q^t$  is computed by applying  $x^t$  as voltages to the rows and reading back the resulting currents on the columns, and  $u^t$  by applying  $z^t$  as voltages to the columns and reading back the resulting currents on the rows. In a memristive crossbar, it has been argued that the matrix-vector multiplications can be performed with constant time complexity  $O(\gamma)$ , where  $\gamma$  is independent of the crossbar size [3]. The reason is that the computation is performed in parallel through Kirchhoff's circuit law locally at the same place where the matrix data are stored. Therefore, the complexity of (3) is potentially reduced from  $O(MN)$  to  $O(N)$  if  $A$  is dense, as it is the case for  $A$  with i.i.d. Gaussian elements. The precise value of  $\gamma$  will depend on the read current settling time and the time required to digitize the current by the peripheral circuitry. Consequently, larger crossbars may eventually lead to higher  $\gamma$  if some of the readout circuitry must be shared across columns/rows and multiplexed.

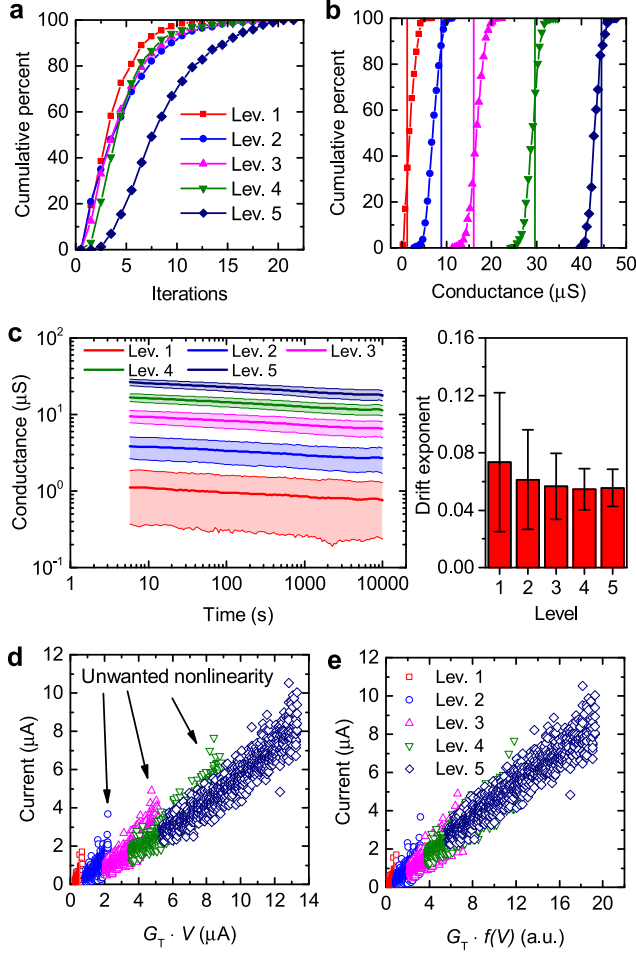
### B. Physical Implementation on Prototype PCM Chip

We implemented CS with AMP recovery using a prototype multilevel PCM chip that contains 1 million usable PCM cells. PCM is a resistive memory technology that is based on the rapid and reversible transition between the crystalline and amorphous phases of certain materials by the application of suitable electrical pulses. Each PCM cell consists of a PCM

device in series with an access transistor. The PCM devices are based on doped- $\text{Ge}_2\text{Sb}_2\text{Te}_2$  and are integrated into the prototype chip in a 90-nm CMOS baseline technology [15]. In addition to the PCM cells, the prototype chip integrates the circuitry for cell addressing, on-chip analog-to-digital converter (ADC) for cell readout, and voltage- or current-mode cell programming. The PCM chip is interfaced to a hardware platform comprising two FPGA boards and an analog front-end board. The layout, picture, and specifications of the experimental PCM chip with integrated read/write circuitry can be found in [5].

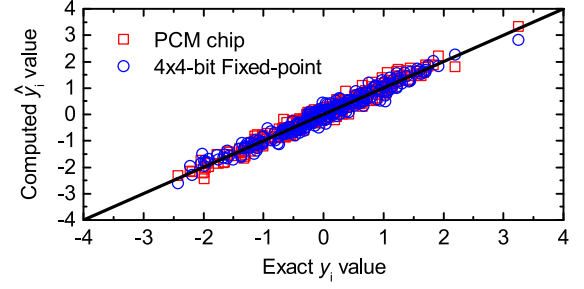
The selection of one PCM device is done by serially addressing a word line and a bitline (BL). For reading a PCM device, the selected BL is biased to a constant voltage (typically 0–300 mV) by a voltage regulator via a voltage generated off-chip. The sensed current is integrated by a capacitor, and the resulting voltage is then digitized by the on-chip 8-bit cyclic ADC. The total time of one read is 1  $\mu\text{s}$ . The readout characteristic is calibrated via on-chip reference polysilicon resistors. For programming a PCM device, a voltage generated off-chip is converted on-chip into a programming current. This current is then mirrored into the selected BL for the desired duration of the programming pulse. Each programming pulse is a box-type rectangular pulse ( $\sim 1$  ns rise/fall times) with a duration of 400 ns and an amplitude varying between 0 and 500  $\mu\text{A}$ . Iterative programming involving a sequence of program-and-verify steps is used to program the PCM devices to the desired conductance values [16]. After each programming pulse, a verify step is performed, and the value of the device conductance programmed in the previous iteration is read at a voltage of 0.2 V. The programming current applied to the PCM device in the subsequent iteration is adapted according to the sign of the value of the error between the target level and the read value of the device conductance. The total time of the one program-and-verify step is approximately 2.5  $\mu\text{s}$ . The array can be erased (RESET) using the maximum amplitude pulse of 500  $\mu\text{A}$  and reprogrammed at will, and each cell can sustain approximately  $10^9$  programming pulses.

In our implementation of CS with AMP recovery, the element-by-element multiplications of the matrix-vector products were realized in the PCM chip, and the remaining operations were implemented in software. The elements of  $A$

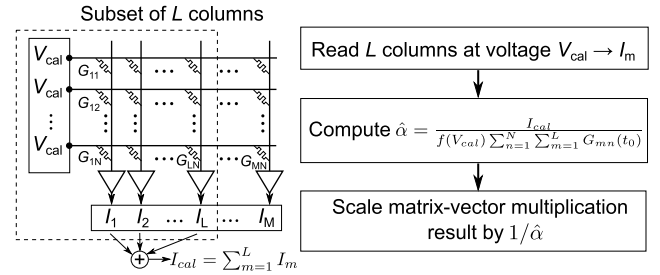


**Fig. 2.** Iterative programming of five representative conductance levels [vertical lines in (b)] on 5000 devices of the PCM chip. (a) Number of iterations needed for the convergence of the iterative programming algorithm. (b) Conductance distributions at approximately  $50 \mu\text{s}$  after programming. (c) Evolution of the mean conductance values of the five programmed levels versus time; filled areas represent the standard deviation for each level, and the plot on the right shows the calculated drift exponent  $\nu$  of the five levels computed from  $G(t) = G(t_0)(t/t_0)^{-\nu}$ . (d) Readout current of the 5000 programmed PCM devices for a voltage range  $0\text{--}0.3 \text{ V}$ , plotted versus  $G_T \cdot V$ , where  $V$  is the applied voltage and  $G_T$  is the target conductance of the different levels. (e) Readout current plotted versus  $G_T \cdot f(V)$ , where  $f(V) = V + 5V^3$ .

were mapped to conductance values between  $0$  and  $50 \mu\text{S}$  and programmed on four PCM devices averaged per element using iterative programming, with a conductance margin of  $1.74 \mu\text{S}$  per device, that is, the iterative algorithm converges when the programmed conductance reaches a value within at most  $1.74 \mu\text{S}$  from the target value. The matrix is programmed only once before CS is performed. Fig. 2(a) shows the number of programming cycles required, and Fig. 2(b) and (c) show the conductance distributions for five representative levels. Here, only five levels are shown for clarity, but in our experiments, the conductance may assume any value in the range  $0\text{--}50 \mu\text{S}$ . We mapped the vector elements to voltage values in the range  $0\text{--}0.3 \text{ V}$  using a nonlinear mapping  $f(V)$  to account for the slight nonlinearity of the current–voltage ( $I\text{--}V$ ) characteristics of the PCM devices [17]. The effect of this mapping is shown in Fig. 2(d) and (e), where each point corresponds to the



**Fig. 3.** Comparison of the precision in the computation of  $y = Ax_0$  by the experimental PCM chip and  $4 \times 4$ -bit multiplications.  $A$  is a  $256 \times 256$  Gaussian matrix coded in the PCM chip,  $x_0$  is a  $256$ -long Gaussian vector applied as voltages, and  $y_i$  is the  $i$ th element of  $y$ .



**Fig. 4.** Calibration procedure to prevent errors due to conductance drift.

current of one PCM device measured at the applied voltage. The accuracy of the matrix-vector computation with our PCM chip for a  $256 \times 256$  matrix with i.i.d. Gaussian elements is comparable to that of a fixed-point implementation where the matrix and vector elements are quantized to 4 bits, as shown in Fig. 3.

To prevent errors in the multiplication results due to conductance drift of the PCM devices, we developed a drift calibration procedure which consists in periodically reading the summed current of  $L$  columns in the array during an experiment. Those  $L$  columns contain devices programmed to known conductance values  $G_{mn}(t_0)$ , and therefore, by reading them periodically at a constant voltage  $V_{\text{cal}}$ , we can compensate for a global conductance shift, as shown in Fig. 4. This procedure is especially simple because  $L$  can be chosen to be small, enough to get sufficient statistics, and the sum  $\sum_{n=1}^N \sum_{m=1}^L G_{mn}(t_0)$  needs to be computed only once. The additional operations for drift calibration can be efficiently implemented and are not expected to incur significant time/power overhead. Reading the subset of  $L$  columns of the crossbar can be done while the PCM array is idle, i.e., when the digital unit performs the additional computations of the recovery algorithm, and additional means are needed to perform the  $L$  current summations as well as computing and storing  $\hat{\alpha}$ . They could be implemented either with on-chip digital circuitry or in the control/processing unit. In our experiments, the calibration procedure was performed in the control unit on  $L = 40$  columns after every five matrix-vector multiplications.

## IV. EXPERIMENTAL RESULTS

### A. Linear Estimation

First, we study the simple use case of linear estimation, where the vector  $x_0$  is not sparse and its entries are i.i.d.

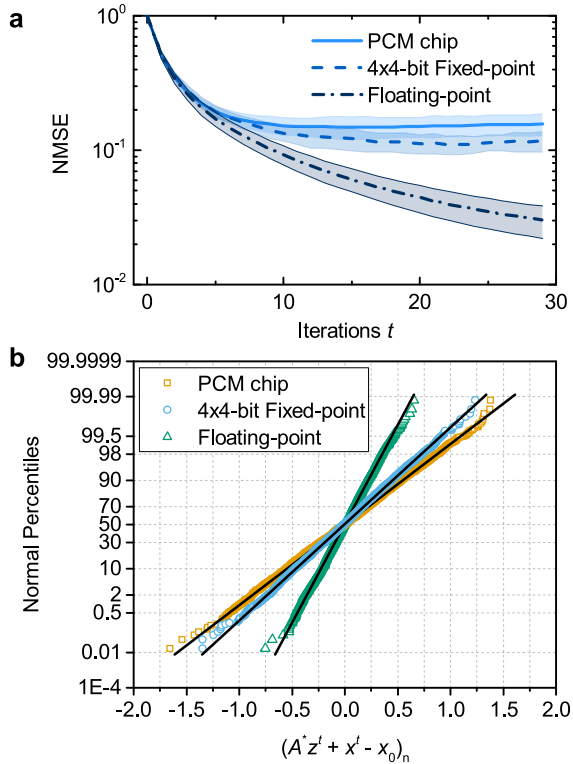


Fig. 5. (a) Normalized mean square error as a function of the number of AMP iterations for linear estimation with  $N = M = 256$ . The filled areas represent the standard deviation over 16 different realizations of  $A$  and  $x_0$ . (b) Empirical distribution of the effective noise  $A^*z^t + x^t - x_0$  at the last AMP iteration  $t = 29$  for the three implementations. All the 16 experiments were used to build the empirical distributions.

Gaussian  $N(0, 1)$ . In this case, the optimal AMP algorithm uses  $\eta_t(x) = \lambda_t x$  with  $\lambda_t = (1/1 + \tau_t^2)$ , where  $\tau_t^2$  is the variance of the empirical distribution of  $A^*z^t + x^t - x_0$ , which can be seen as the effective noise of the algorithm at iteration  $t$  [7].  $\tau_t^2$  can be estimated by  $\hat{\tau}_t^2 = \|z^t\|_2^2/M$ , which is shown to be a good approximation of the variance of  $A^*z^t + x^t - x_0$  in the large system limit [18].

We implemented this algorithm on the PCM chip for a random signal  $x_0$  of size  $N = 256$  and  $M = N$  measurements. The  $M \times N$  measurement matrix  $A$  was programmed in the PCM chip with i.i.d. Gaussian elements normalized, such that the norm of its columns is approximately 1 [7]. The measurements  $y$  were obtained by applying  $x_0$  as voltages on the PCM chip after matrix  $A$  had been programmed, thus realizing  $Ax_0$  in hardware. Subsequently,  $x_0$  was reconstructed with AMP using the PCM chip to compute the matrix-vector operations  $Ax^t$  and  $A^*z^t$ , as shown in Fig. 1(a). We performed the experiment 16 times for 16 different realizations of randomly generated  $A$  and  $x_0$  and reported the mean and standard deviation of the normalized mean square error (NMSE)  $\|x^t - x_0\|_2^2/\|x_0\|_2^2$  over those 16 experiments. The different realizations of  $A$  and  $x_0$  were chosen, such that proper convergence of the AMP algorithm was obtained.<sup>1</sup>

<sup>1</sup>Due to the small system size ( $N = 256$ ), AMP does not converge properly for all combinations of randomly generated  $A$  and  $x_0$ . In the experiments, we ensured that for all realizations of  $A$  and  $x_0$  chosen, the NMSE neither floors nor starts monotonically increasing in the floating-point implementation within the number of AMP iterations performed, in the case 29.

The evolution of the NMSE between the original and reconstructed signals is shown in Fig. 5(a). The NMSE decreases as  $1/(1+t)$  for the floating-point implementation as dictated by state evolution [7]. For the PCM chip and an implementation where the multiplications in  $Ax^t$  and  $A^*z^t$  are done in  $4 \times 4$ -bit fixed-point arithmetic, the NMSE floors at values of approximately 0.15 and 0.12, respectively. However, the initial convergence rate of AMP is not affected by the inexact implementations. This finding will be further confirmed in the next experiments of Sections IV-B and IV-C.

An important feature of AMP is that the effective noise  $A^*z^t + x^t - x_0$  is approximately Gaussian [18]. This allows the asymptotically exact analysis of AMP whereby the variance of this noise can be computed exactly from state evolution for any  $t$  when  $N \rightarrow \infty$  [7]. Moreover, the variance can be used as an input to the function  $\eta_t$  in order to optimally denoise this Gaussian noise [7]. For iterative thresholding (2), the effective noise is generally not Gaussian, and state evolution does not hold [6], [7]. Hence, it is important to verify whether the Gaussianity of this noise is affected by the PCM implementation. We obtained the effective noise  $A^*z^t + x^t - x_0$  at the last AMP iteration for the three implementations. We found no clear departure from a Gaussian distribution for both the PCM and fixed-point implementations [see Fig. 5(b)]. The tails which deviate from an exact Gaussian distribution close to percentiles 0.01 and 99.99 observed in all three implementations are likely a consequence of the small system size ( $N = 256$ ).

## B. Compressed Sensing With Soft-Thresholding

In this use case, the vector  $x_0$  is  $k$ -sparse, i.e., it contains  $k$  nonzero elements, and its nonzero elements are i.i.d. Gaussian  $N(0, 1)$ . In order to reconstruct  $x_0$  from the measurements  $y$ , we use the AMP algorithm (3) with a sequence of soft-threshold functions  $\eta_t(x)$  defined as [6]

$$\eta_t(x) = \begin{cases} x - \tau_t, & \text{if } x > \tau_t \\ 0, & \text{if } -\tau_t \leq x \leq \tau_t \\ x + \tau_t, & \text{if } x < -\tau_t \end{cases} \quad (4)$$

with thresholds  $\tau_t = \|z^t\|_2/\sqrt{M}$ . For the soft-threshold function (4), the term  $(N/M)z^{t-1}\langle\eta'_{t-1}(A^*z^{t-1} + x^{t-1})\rangle$  in the AMP algorithm can be calculated explicitly and yields  $(N/M)z^{t-1}\langle\eta'_{t-1}(A^*z^{t-1} + x^{t-1})\rangle = (1/M)z^{t-1}\|\eta_{t-1}(A^*z^{t-1} + x^{t-1})\|_0$ , where  $\|x\|_0$  denotes the number of nonzero elements of  $x$ .

We performed the experiments for a random signal  $x_0$  of size  $N = 256$  and  $k = 64$  randomly distributed nonzero elements. We tested cases for sampling rates of  $M/N = 1$  (no compression) and  $M/N = 0.75$ , each with 16 different realizations of randomly generated  $A$  and  $x_0$ . The evolution of the NMSE between the original and reconstructed signals is shown in Fig. 6(a). As in the previous use case, the initial convergence rate of AMP is unaffected by the approximate multiplications done in the PCM chip, and the magnitude of the NMSE floor obtained with the PCM chip is comparable to the  $4 \times 4$ -bit fixed-point implementation. When using a lower sampling rate  $M/N = 0.75$ , the convergence rate of

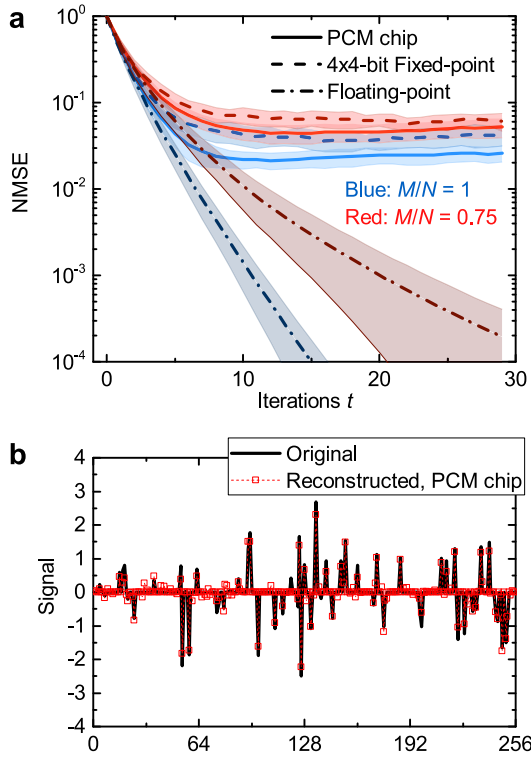


Fig. 6. (a) Normalized mean square error versus the number of AMP iterations for CS with soft thresholding; filled areas represent the standard deviation over 16 different realizations of  $A$  and  $x_0$ . (b) Example of the original and reconstructed signals for the PCM implementation with  $M/N = 0.75$ .

AMP decreases, and the NMSE floor increases for the inexact implementations compared with  $M/N = 1$ .

In certain applications, it is sufficient to recover only the sparsity pattern of  $x_0$ , without being concerned with the exact values of the nonzero elements. We show in Fig. 6(b) the original and reconstructed signals for one of the experiments performed with 0.75 sampling rate. We see that the general shape and the sparsity pattern of the signal are well recovered in the PCM implementation. Thus, in applications where the reconstruction accuracy is not of paramount importance, the accuracy obtained with our current prototype PCM chip may already be sufficient.

### C. Compressive Imaging With Image Denoising

Compressive imaging refers to performing CS on image signals. The elements of  $x_0$  thus represent the pixel intensities of an image. The goal is to acquire the image with  $M \ll N$  measurements and to reconstruct it accurately. A general methodology for compressive imaging with AMP was recently introduced by Metzler *et al.* [18]. They developed an extension of the AMP algorithm that uses a denoiser within its iterations. The proposed algorithm is given by

$$\begin{aligned} x^{t+1} &= D_{\tau_t}(A^*z^t + x^t) \\ z^t &= y - Ax^t + \frac{1}{M}z^{t-1}\text{div}D_{\tau_{t-1}}(A^*z^{t-1} + x^{t-1}) \\ \tau_t^2 &= \|z^t\|_2^2/M \end{aligned} \quad (5)$$

where  $D_{\tau}$  denotes a denoiser, which takes as input a signal plus Gaussian noise and an estimate of the standard deviation

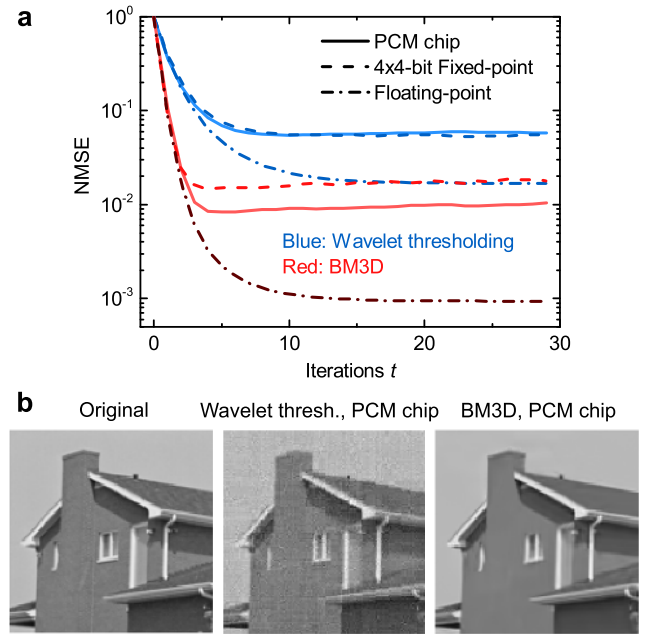


Fig. 7. (a) Evolution of the NMSE in image reconstruction for wavelet thresholding and BM3D denoisers with  $M/N = 1/2$ . (b) Original and reconstructed images with the PCM implementation.

of that noise  $\tau$ , and  $\text{div}D_{\tau}(x) = \sum_{n=1}^N (\partial D_{\tau}(x)_n / \partial x_n)$  denotes the divergence of the denoiser, where  $D_{\tau}(x)_n$  is the  $n$ th element of  $D_{\tau}(x)$  and  $x_n$  is the  $n$ th element of  $x$ .

We tested this algorithm using the  $128 \times 128$  pixel “house” image shown in Fig. 7(b) as signal  $x_0$ . We implemented the two following denoisers.

- 1) *Wavelet Thresholding*: It transforms the signal into a wavelet basis, thresholds the coefficients, and then inverts the transform. If  $W$  denotes the wavelet transform, this denoiser is defined as  $D_{\tau_t}(x) = W^{-1}\eta_t(Wx)$ . We used the soft-threshold function (4) as  $\eta_t$  and 2-D Haar wavelet transform. The divergence of this denoiser can be calculated explicitly and yields  $\text{div}D_{\tau_{t-1}}(A^*z^{t-1} + x^{t-1}) = \|\eta_{t-1}(W(A^*z^{t-1} + x^{t-1}))\|_0$ , which is the number of nonzero elements of the thresholded sparsified estimate.
- 2) *Block Matching 3-D Collaborative Filtering (BM3D)*: It can be considered a combination of nonlocal means (averaging weighted neighboring pixels) and wavelet thresholding. The term  $\text{div}D_{\tau_{t-1}}(A^*z^{t-1} + x^{t-1})$  cannot be calculated explicitly and thus is estimated using the Monte Carlo procedure described in [18]. The divergence is estimated with  $\text{div}D_{\tau}(x) \simeq (b^*/\epsilon)(D_{\tau}(x + \epsilon b) - D_{\tau}(x))$  for small  $\epsilon$  and vector  $b$  with elements i.i.d.  $N(0, 1)$ . BM3D performs much better on images than wavelet thresholding because images are not exactly sparse in the wavelet domain.

The length of  $x_0$  in this experiment is  $N = 16384$ . For such a large value of  $N$ , it is not possible to code all elements of an  $M \times N$  Gaussian matrix in our PCM hardware, which has only 1 million usable devices. To overcome this difficulty, we use a block-based compression approach, whereby a small measurement matrix  $H$  of size  $M_s \times N_s$  is used, with

TABLE I  
PSNR (IN dB) OF THE  $128 \times 128$  "HOUSE" IMAGE RECONSTRUCTIONS

	PCM chip	Fixed-point	Floating-point
Wavelet thresh.	27.15	27.39	32.50
BM3D	34.58	32.27	45.06

$N_s = 256$ . We perform measurements on consecutive  $16 \times 16$  pixel blocks using the same measurement matrix  $H$ . In order to obtain uncorrelated measurements and ensure the convergence of AMP, we perform a (fixed) random permutation  $P$  of the pixel intensities before doing the measurements. The matrix  $A$  can thus be written as  $A = \text{blkdiag}(H)P$ , where  $\text{blkdiag}(H)$  is an  $M \times N$  matrix with  $N/N_s$  main diagonal blocks matrices  $H$ , where it is assumed that  $N$  is a multiple of  $N_s$  and  $M_s/N_s = M/N$ . The elements of  $H$  are i.i.d.  $\sim N(0, 1/M_s)$ .

We programmed a  $128 \times 256$  Gaussian measurement matrix  $H$  in the PCM chip (sampling rate  $M/N = 1/2$ ), divided the image into  $16 \times 16$  pixel blocks, and compressed each block individually with the PCM chip. Subsequently, the image was reconstructed with algorithm (5) using the PCM chip to compute the matrix-vector operations  $Ax^t$  and  $A^*z^t$ . In Fig. 7(a), we show the NMSE evolution for the PCM, fixed-point, and floating-point implementations for wavelet thresholding and BM3D denoisers. The peak SNR<sup>2</sup> (PSNR) at the last AMP iteration is reported in Table I. It can be seen that using a better denoiser (e.g., BM3D) results in a lower final NMSE in the PCM and fixed-point implementations. It indicates that denoisers can be used effectively to improve the reconstruction accuracy by mitigating the errors from the PCM chip. Moreover, the convergence rate of AMP is only affected by the choice of the denoiser but not by the approximate implementations.

## V. DISCUSSION

There are several reasons why AMP is well suited for a memristive implementation. First, matrix  $A$  does not change over iterations, and thus, only read operations are performed during AMP reconstruction. Therefore, matrix  $A$  needs to be programmed only once and will be retained in the array thanks to the nonvolatility of the PCM devices. The read operations that are performed during reconstruction require significantly less power than programming and thus can be heavily parallelized. With the 90-nm PCM technology used in this paper, we estimate the read energy to be between 1 and 100 fJ per device depending on the programmed resistance state, compared with approximately 100 pJ for programming (assuming five program-and-verify iterations). Moreover, unlike programming endurance, the read endurance (at least in PCM) is essentially unlimited; hence, this implementation is favorable with respect to device reliability issues and will not lead to device degradation due to excessive reprogramming at every iteration.

The effect of device imperfections and failures on the final reconstruction NMSE is discussed in [5]. We found that the AMP recovery can tolerate conductance variations due to programming errors (up to 20%) and up to 20% stuck-SET and

stuck-RESET device failures. Device imperfections that have a detrimental effect on the reconstruction accuracy include the device conductance noise (most dominant effect) and the  $I$ - $V$  nonlinearity. Finally, the achievable reconstruction NMSE is ultimately limited by the resolution of the digital-to-analog/analog-to-digital converters used at the input/output of the crossbar array.

To quantify the potential energy gains of the memristive implementation over a digital design, based on the figures currently achieved with our prototype PCM chip, we made an FPGA design that operates at the same speed and the same precision at which we expect a PCM-based crossbar to perform [5]. In (3), the matrix-vector multiplications are the most expensive operations, so we compared the memristive crossbar analog multiplier with a 4-bit FPGA multiplier design. The 4-bit matrix elements are stored in the FPGA block-RAM, and 32 dot-product units operate in parallel to compute a  $256 \times 256$  matrix-vector product in  $1.2 \mu\text{s}$ . The dynamic power consumption achieved with this design is 800 mW [5]. In a  $256 \times 256$  PCM-based crossbar, the dynamic power dissipation in the devices for one read operation would be in the order of 13.1 mW (read current of  $1 \mu\text{A}$  per device at 0.2 V). Thus, a  $256 \times 256$  PCM-based crossbar in the 90-nm technology operating at  $1 \mu\text{s}$  cycle time plus two 8-bit ADCs operating at 125 MS/s to convert the current (12-mW/GS/s power consumption) is expected to consume 16.2 mW, which is 50 times less than the FPGA design. The power advantage arises because only read operations, which consume little energy, are performed in the memristive crossbar for multiplications.

While PCM devices were used for the experiments presented in this paper, other memory devices could be considered to perform the analog matrix-vector multiplications in the proposed CS implementation. Potential candidates include metal-oxide resistive random-access memory [3], NOR Flash [19], and static random-access memory [20]. The main advantages of PCM for this application are its multilevel capability along with fast read/write latency and nonvolatility; however, the PCM programming current is generally higher than the other technologies, and resistance drift poses additional challenges that need to be addressed. Assessing different technologies for in-memory computing should account for array-level variability, device noise, and accuracy/ease of device programming in addition to latency and power consumption.

In the ASIC implementation of AMP reported in [12], the multiply-accumulate (MAC) units and the matrix generating unit take most of the chip area and are responsible for most of the power consumption, which amounts to  $> 90\%$  in the proposed AMP-M design for arbitrary matrices. In such an implementation, matrix  $A$  would have to be explicitly stored [in off-chip dynamic random-access memory (DRAM)], or its coefficient would have to be generated on the fly at every AMP iteration. In a memristive implementation, matrix  $A$  is stored in the memristive array(s) in a nonvolatile manner, thus avoiding the need of a unit to generate its coefficients or using an off-chip DRAM, while still being able to reprogram it without redesigning the entire circuit. Moreover, by computing

<sup>2</sup>PSNR =  $10 \log_{10}(255^2 / (\|\hat{x} - x_0\|_2^2 / N))$ , where  $\hat{x}$  is the estimate of  $x_0$ .

the matrix-vector multiplications inside the memristive array, the use of MAC units, which are expensive in both power and area when implemented in CMOS, is completely avoided.

Furthermore, a remarkable property of AMP is that its convergence rate is independent of the precision of the matrix-vector multiplications. This is a highly desirable property for this type of implementation, as the number of AMP iterations needed for reconstruction will not be larger than in a floating-point implementation. We also found that the NMSE floor due to computational errors can be lowered by using appropriate denoisers within AMP. Obviously, using a complex denoiser, such as BM3D, might not be efficient from an implementation point of view, because the speedup obtained by performing the matrix-vector multiplications in the memristive array may be overcompensated by the time required to apply the denoiser. However, an interesting avenue would be to design a denoiser that is specifically aimed at removing the computational errors from the memristive array.

Regarding the limitations of the memristive implementation, the computational errors from the memristive array are currently the biggest drawback. Very accurate reconstruction cannot be currently achieved with our prototype PCM chip, which performs with a precision similar to that of a matrix-vector product in the  $4 \times 4$ -bit fixed-point implementation. However, the precision of analog in-memory computation is expected to improve as the technology matures, e.g., with concepts such as projected memory to reduce the noise and drift [21]. The precision could be further increased by mapping a single column of the matrix across multiple physical columns of an array encoding different bits and applying the input vector to the array one or several bits at a time, still performing in-memory computing, at the expense of area and energy penalty, and additional support required by the peripheral circuitry.

Another limitation is that, for CS applications, it might be hard to justify the memristive implementation versus a digital implementation with a 1-bit measurement matrix, as the latter shows no loss in SNR for the compressed measurement acquisition and no multipliers are needed for a binary matrix [10]. However, this type of implementation is limited to one specific application only, i.e., only a binary measurement matrix is supported, whereas a memristive implementation can be used for any arbitrary measurement matrix. Moreover, such efficient implementations currently only acquire the compressed measurements and do not support reconstruction, which has to be done off-chip. The attractiveness of the memristive implementation is that both compression and reconstruction could be done on the same platform.

## VI. CONCLUSION

We propose an implementation of CS with AMP recovery based on the memristive crossbar arrays. The measurement matrix elements are programmed as conductance values of memristive devices in crossbar arrays, which are used to perform the matrix-vector multiplications in both the compression and the recovery algorithm. In this way, the computational complexity of AMP recovery is potentially reduced from

$O(MN)$  to  $O(N)$ . We tested this implementation experimentally for three use cases of AMP using more than 256k PCM devices in a prototype multilevel PCM chip to perform the matrix-vector multiplications. We found that the convergence rate of AMP is not affected by performing the matrix-vector multiplications in the PCM array. The accuracy achieved with our prototype PCM chip is comparable to that of a fixed-point implementation where the matrix and vector elements are quantized to 4 bits. In applications where the reconstruction accuracy is not of paramount importance, the memristive implementation could represent a viable solution to provide more efficient AMP reconstruction than a full von Neumann implementation.

## ACKNOWLEDGMENT

The authors would like to thank N. Papandreou and U. Egger for experimental help, L. Kull and T. Toifl for discussions, and M. Brightsky for providing the PCM devices used in this paper.

## REFERENCES

- [1] M. Le Gallo *et al.*, "Mixed-precision in-memory computing," *Nature Electron.*, vol. 1, no. 4, pp. 246–253, Apr. 2018, doi: [10.1038/s41928-018-0054-8](https://doi.org/10.1038/s41928-018-0054-8).
- [2] P. M. Sheridan, F. Cai, C. Du, W. Ma, Z. Zhang, and W. D. Lu, "Sparse coding with memristor networks," *Nature Nanotechnol.*, vol. 12, pp. 784–789, May 2017, doi: [10.1038/nnano.2017.83](https://doi.org/10.1038/nnano.2017.83).
- [3] M. Hu *et al.*, "Memristor-based analog computation and neural network classification with a dot product engine," *Adv. Mater.*, vol. 30, no. 9, p. 1705914, 2018, doi: [10.1002/adma.201705914](https://doi.org/10.1002/adma.201705914).
- [4] G. Cherubini, P. Hurley, M. Simeoni, and S. Kazemi, "Imaging in radio interferometry by iterative subset scanning using a modified AMP algorithm," in *Proc. ICASSP*, Mar. 2016, pp. 3326–3330, doi: [10.1109/ICASSP.2016.7472293](https://doi.org/10.1109/ICASSP.2016.7472293).
- [5] M. Le Gallo, A. Sebastian, G. Cherubini, H. Giefers, and E. Eleftheriou, "Compressed sensing recovery using computational memory," in *IEDM Tech. Dig.*, Dec. 2017, pp. 28.3.1–28.3.4, doi: [10.1109/IEDM.2017.8268469](https://doi.org/10.1109/IEDM.2017.8268469).
- [6] D. L. Donoho, A. Maleki, and A. Montanari, "Message-passing algorithms for compressed sensing," *Proc. Nat. Acad. Sci. USA*, vol. 106, no. 45, pp. 18914–18919, Nov. 2009, doi: [10.1073/pnas.0909892106](https://doi.org/10.1073/pnas.0909892106).
- [7] M. Bayati and A. Montanari, "The dynamics of message passing on dense graphs, with applications to compressed sensing," *IEEE Trans. Inf. Theory*, vol. 57, no. 2, pp. 764–785, Feb. 2011, doi: [10.1109/TIT.2010.2094817](https://doi.org/10.1109/TIT.2010.2094817).
- [8] M. F. Duarte *et al.*, "Single-pixel imaging via compressive sampling," *IEEE Signal Process. Mag.*, vol. 25, no. 2, pp. 83–91, Mar. 2008, doi: [10.1109/MSP.2007.914730](https://doi.org/10.1109/MSP.2007.914730).
- [9] L. Jacques, P. Vanderghyest, A. Bibet, V. Majidzadeh, A. Schmid, and Y. Leblebici, "CMOS compressed imaging by random convolution," in *Proc. ICASSP*, Apr. 2009, pp. 1113–1116, doi: [10.1109/ICASSP.2009.4959783](https://doi.org/10.1109/ICASSP.2009.4959783).
- [10] Y. Oike and A. El Gamal, "CMOS image sensor with per-column  $\Sigma\Delta$  ADC and programmable compressed sensing," *IEEE J. Solid-State Circuits*, vol. 48, no. 1, pp. 318–328, Jan. 2013, doi: [10.1109/JSSC.2012.2214851](https://doi.org/10.1109/JSSC.2012.2214851).
- [11] P. Maechler, P. Greisen, B. Sporrer, S. Steiner, N. Felber, and A. Burg, "Implementation of greedy algorithms for LTE sparse channel estimation," in *Proc. ASIOMAR*, Nov. 2010, pp. 400–405, doi: [10.1109/ACSSC.2010.5757587](https://doi.org/10.1109/ACSSC.2010.5757587).
- [12] P. Maechler *et al.*, "VLSI design of approximate message passing for signal restoration and compressive sensing," *IEEE J. Emerg. Sel. Topics Circuits Syst.*, vol. 2, no. 3, pp. 579–590, Sep. 2012, doi: [10.1109/JETCAS.2012.2214636](https://doi.org/10.1109/JETCAS.2012.2214636).
- [13] L. Bai, P. Maechler, M. Muehlberghuber, and H. Kaeslin, "High-speed compressed sensing reconstruction on FPGA using OMP and AMP," in *Proc. 19th IEEE Int. Conf. Electron., Circuits Syst. (ICECS)*, Dec. 2012, pp. 53–56, doi: [10.1109/ICECS.2012.6463559](https://doi.org/10.1109/ICECS.2012.6463559).
- [14] S. Liu, A. Ren, Y. Wang, and P. K. Varshney, "Ultra-fast robust compressive sensing based on memristor crossbars," in *Proc. ICASSP*, Mar. 2017, pp. 1133–1137, doi: [10.1109/ICASSP.2017.7952333](https://doi.org/10.1109/ICASSP.2017.7952333).



- [15] M. Breitwisch *et al.*, "Novel lithography-independent pore phase change memory," in *VLSI Symp. Technol. Dig.*, Jun. 2007, pp. 100–101, doi: [10.1109/VLSIT.2007.4339743](https://doi.org/10.1109/VLSIT.2007.4339743).
- [16] N. Papandreou *et al.*, "Programming algorithms for multilevel phase-change memory," in *Proc. ISCAS*, May 2011, pp. 329–332, doi: [10.1109/ISCAS.2011.5937569](https://doi.org/10.1109/ISCAS.2011.5937569).
- [17] M. Le Gallo, M. Kaes, A. Sebastian, and D. Krebs, "Subthreshold electrical transport in amorphous phase-change materials," *New J. Phys.*, vol. 17, no. 9, p. 093035, Sep. 2015, doi: [10.1088/1367-2630/17/9/093035](https://doi.org/10.1088/1367-2630/17/9/093035).
- [18] C. A. Metzler, A. Maleki, and R. G. Baraniuk, "From denoising to compressed sensing," *IEEE Trans. Inf. Theory*, vol. 62, no. 9, pp. 5117–5144, Sep. 2016, doi: [10.1109/TIT.2016.2556683](https://doi.org/10.1109/TIT.2016.2556683).
- [19] X. Guo *et al.*, "Fast, energy-efficient, robust, and reproducible mixed-signal neuromorphic classifier based on embedded NOR flash memory technology," in *IEDM Tech. Dig.*, Dec. 2017, pp. 6.5.1–6.5.4, doi: [10.1109/IEDM.2017.8268341](https://doi.org/10.1109/IEDM.2017.8268341).
- [20] A. Biswas and A. P. Chandrakasan, "Conv-RAM: An energy-efficient SRAM with embedded convolution computation for low-power CNN-based machine learning applications," in *IEEE ISSCC Dig. Tech. Papers*, Feb. 2018, pp. 488–490, doi: [10.1109/ISSCC.2018.8310397](https://doi.org/10.1109/ISSCC.2018.8310397).
- [21] W. W. Koelmans, A. Sebastian, V. P. Jonnalagadda, D. Krebs, L. Dellmann, and E. Eleftheriou, "Projected phase-change memory devices," *Nat. Commun.*, vol. 6, Sep. 2015, Art. no. 8181, doi: [10.1038/ncomms9181](https://doi.org/10.1038/ncomms9181).



**Manuel Le Gallo** received the bachelor's degrees from the Ecole Polytechnique de Montréal, Montreal, QC, Canada, and Ecole Polytechnique (X), Palaiseau, France, and the M.Sc. and Ph.D. degrees from ETH Zürich, Zürich, Switzerland, in 2014 and 2017, respectively.

He joined IBM Research Zurich, Rüschlikon, Switzerland, in 2013, where he is currently a Post-Doctoral Researcher. His current research interests include using phase-change memory devices for non-von Neumann computing.



**Abu Sebastian** (M'03–SM'11) is currently a Principal Research Staff Member and a Master Inventor at IBM Research Zurich, Rüschlikon, Switzerland. He was a contributor to several key projects in the field of storage and memory technologies. His current research interests include non-von Neumann computing for applications, such as artificial intelligence and machine learning.



**Giovanni Cherubini** (S'80–M'82–SM'94–F'06) received the Laurea degree (*summa cum laude*) from the University of Padova, Padua, Italy, and the M.S. and Ph.D. degrees from the University of California at San Diego, La Jolla, CA, USA, all in electrical engineering.

He joined IBM Research–Zurich, Rüschlikon, Switzerland, in 1987. He holds over 200 patents in the areas of signal processing, control, data storage, and communication systems.



**Heiner Giefers** (M'13–SM'17) received the Diploma and Ph.D. degrees from the University of Paderborn, Paderborn, Germany, in 2006 and 2012, respectively.

In 2013, he joined IBM Research Zurich, Rüschlikon, Switzerland, where he was involved in energy-efficient computing, reconfigurable architectures, and hardware–software codesign. Since 2018, he has been a Professor of cloud computing with the South Westphalia University of Applied Sciences, Iserlohn, Germany.



**Evangelos Eleftheriou** (F'02) received the Ph.D. degree in electrical engineering from Carleton University, Ottawa, ON, Canada, in 1985.

In 1986, he joined the IBM Research Zurich, Rüschlikon, Switzerland, as a Research Staff Member, where he has been holding various management positions since 1998 and is currently responsible for the neuromorphic computing activities.