

ARTICLE

DOI: 10.1038/s41467-018-04482-4

OPEN

Implementation of multilayer perceptron network with highly uniform passive memristive crossbar circuits

F. Merrikh Bayat¹, M. Prezioso¹, B. Chakrabarti¹, H. Nili¹, I. Kataeva² & D. Strukov¹

The progress in the field of neural computation hinges on the use of hardware more efficient than the conventional microprocessors. Recent works have shown that mixed-signal integrated memristive circuits, especially their passive (OT1R) variety, may increase the neuro-morphic network performance dramatically, leaving far behind their digital counterparts. The major obstacle, however, is immature memristor technology so that only limited functionality has been reported. Here we demonstrate operation of one-hidden layer perceptron classifier entirely in the mixed-signal integrated hardware, comprised of two passive 20×20 metal-oxide memristive crossbar arrays, board-integrated with discrete conventional components. The demonstrated network, whose hardware complexity is almost $10\times$ higher as compared to previously reported functional classifier circuits based on passive memristive crossbars, achieves classification fidelity within 3% of that obtained in simulations, when using ex-situ training. The successful demonstration was facilitated by improvements in fabrication technology of memristors, specifically by lowering variations in their I - V characteristics.

¹Electrical and Computer Engineering Department, University of California, Santa Barbara, CA 93117, USA. ²DENSO CORP, 500-1 Minamiyama, Komenoki-cho, Nisshin 470-0111, Japan. These authors contributed equally: F. Merrikh Bayat, M. Prezioso. Correspondence and requests for materials should be addressed to I.K. (email: IRINA_KATAEVA@denso.co.jp) or to D.S. (email: strukov@ece.ucsb.edu)

Started more than half a century ago, the field of neural computation has known its ups and downs, but since 2012, it exhibits an unprecedented boom triggered by the dramatic breakthrough in the development of deep convolutional neuromorphic networks^{1,2}. The breakthrough³ was enabled not by any significant algorithm advance, but rather by the use of high performance graphics processors⁴, and the further progress is being fueled now by the development of even more powerful graphics processors and custom integrated circuits^{5–7}. Nevertheless, the energy efficiency of these implementations of convolutional networks (and other neuromorphic systems^{8–11}) remains well below that of their biological prototypes^{12,13}, even when the most advanced CMOS technology is used. The main reason for this efficiency gap is that the use of digital operations for mimicking biological neural networks, with their high redundancy and intrinsic noise, is inherently unnatural. On the other hand, recent works have shown^{11–16} that analog and mixed-signal integrated circuits, especially using nanoscale devices, may increase the neuromorphic network performance dramatically, leaving far behind both their digital counterparts and biological prototypes and approaching the energy efficiency of the brain. The background for these advantages is that in such circuits the key operation performed by any neuromorphic network, the vector-by-matrix multiplication, is implemented on the physical level by utilization of the fundamental Ohm and Kirchhoff laws. The key component of this circuit is a nanodevice with adjustable conductance G —essentially an analog nonvolatile memory cell—used at each crosspoint of a crossbar array, and mimicking the biological synapse.

Though potential advantages of specialized hardware for neuromorphic computing had been recognized several decades ago^{17,18}, up until recently, adjustable conductance devices were mostly implemented using the standard CMOS technology¹³. This approach was used to implement several sophisticated, efficient systems—see, e.g., refs. 14,15. However, these devices have relatively large areas leading to higher interconnect capacitance and hence larger time delays. Fortunately, in the last decade, another revolution has taken place in the field of nanoelectronic memory devices. Various types of emerging nonvolatile memories are now being actively investigated for their use in fast and energy-efficient neuromorphic networks^{19–41}. Of particular importance, is the development of the technology for programmable, nonvolatile two-terminal devices called ReRAM or memristors^{42,43}. The low-voltage conductance G of these devices may be continuously adjusted by the application of short voltage pulses of higher, typically >1 V amplitude⁴². These devices were used to demonstrate first neuromorphic network providing pattern classification^{21,26,28,30,32,40}. The memristors can have a very

low chip footprint, which is determined only by the overlap area of the metallic electrodes, and may be scaled down below 10 nm without sacrificing their endurance, retention, and tuning accuracy, with some of the properties (such as the ON/OFF conductance ratio) being actually improved⁴⁴.

Much of the previous very impressive demonstrations of neuromorphic networks based on resistive switching memory devices, including pioneering work by IBM^{25,34}, were based on the so-called 1T1R technology, in which every memory cell is coupled to a select transistor^{22,27–31}. The reports of neuromorphic functionality based on passive 0T1R or 1D1R circuits (in which acronyms stand for 0 Transistor or 1 Diode +1 Resistive switching device per memory cell, respectively) have been so far very limited^{26,39}, in part due to much stricter requirement for memristors' I - V uniformity for successful operation. The main result of this paper is the experimental demonstration of a fully functional, board-integrated, mixed-signal neuromorphic network based on passively integrated metal-oxide memristive devices. Our focus on 0T1R memristive crossbar circuits is specifically due to their better performance and energy-efficiency prospects, which can be further improved by three-dimensional monolithical integration^{45–47}. Due to the extremely high effective integration density, three-dimensional memristive circuits will be instrumental in keeping all the synaptic weights of a large-scale artificial neural networks locally, thus cutting dramatically the energy and latency overheads of the off-chip communications. The demonstrated network is comprised of almost an order of magnitude higher number of devices as compared to the previously reported neuromorphic classifiers based on passive crossbar circuits²⁶. The inference, the most common operation in applications of deep learning, is performed directly in a hardware, which is different from many previous works that relied on post-processing the experimental data with external computer to emulate the functionality of the whole system^{25–27,34,39,40}.

Results

Integrated memristors. The passive 20×20 crossbar arrays with Pt/ Al_2O_3 / TiO_{2-x} /Ti/Pt memristor at each crosspoint were fabricated using a technique similar to that reported in ref. 26 (Fig. 1). Specifically, the bilayer binary oxide stack was deposited using low temperature reactive sputtering method. The crossbar electrodes were evaporated using oblique angle physical vapor deposition (PVD) and patterned by lift-off technique using lithographical masks with 200-nm lines separated by 400-nm gaps. Each crossbar electrode is contacted to a thicker (Ni/Cr/Au 400 nm) metal line/bonding pad, which are formed at the last step of the fabrication process. As evident from Fig. 1a, b, due to the

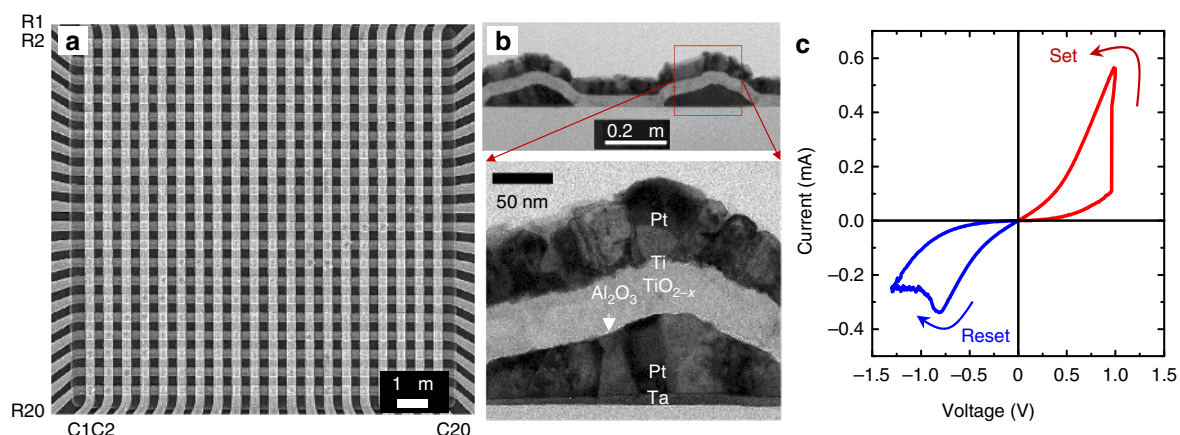


Fig. 1 Passive memristive crossbar circuit. **a** A top-view SEM and **b** cross-section TEM images of 20×20 Pt/ Al_2O_3 / TiO_{2-x} /Ti/Pt crossbar circuit; **c** A typical I - V switching curve

utilized undercut in the photoresist layer and tilted PVD sputtering in the lift-off process, the metal electrodes have roughly triangular shape with ~ 250 nm width. Such shape of the bottom electrodes ensured better step coverage for the following processing layers and, in particular, helped to reduce the top electrode resistance. The externally measured (pad-to-pad) crossbar line resistance for the bonded chip is around 800Ω . It is similar to that of smaller crossbar circuit reported in ref.²⁶ due to the dominant contribution of the contact between crossbar electrode and thicker bonding lines.

Majority of the devices required an electroforming step which consisted of one-time application of a high current (or voltage) ramp bias. We have used both increasing amplitude current and

voltage sweeps for forming but did not see much difference in the results of the forming procedure (Fig. 2). This could be explained by the dominant role of capacitive discharge from the crossbar line during forming, which cannot be controlled well by external current source or current compliance. The devices were formed one at a time, and to speed up the whole process, an automated setup has been developed—see Methods section for more details. The setup was used for early screening of defective samples and has allowed a successful forming and testing of numerous crossbar arrays (Fig. 2). Specially, about 1–2.5% of the devices in the crossbar arrays, i.e., 10 or less out of 400 total, could not be formed with the algorithm parameters that we used. (It might have been possible to form even these devices by applying larger stress but we have not tried it in this experiment to avoid permanently damaging the crossbar circuit.) Typically, the failed devices were stuck at some conductance state, comparable to the range of conductances utilized in the experiment, and as a result have negligible impact on the circuit functionality.

Memristor I - V characteristics are nonlinear (Fig. 1c) due to the alumina barrier between the bottom electrode and the switching layer. I - V 's nonlinearity provides sufficient selector functionality to limit leakage currents in the crossbar circuit, and hence reduce disturbance of half-selected devices during conductance tuning. It is worth mentioning that the demonstrated nonlinearity is weaker as compared to state-of-the-art selector devices that are developed in the context of memory applications. However, our analysis (Supplementary Note 1) shows that strengthening I - V nonlinearity would only reduce power consumption during very infrequent tuning operation but otherwise have no impact on the more common inference operation in the considered neuromorphic applications.

Most importantly, memristive devices in the fabricated 20×20 crossbar circuits have uniform characteristics with gradual (analog) switching. The distributions of the effective set and reset voltages are sufficiently narrow (Fig. 2) to allow precise tuning of devices' conductances to the desired values in the whole array (Fig. 3, Supplementary Fig. 12), which is especially challenging in the passive integrated circuits due to half-select disturbance. For example, an analog tuning was essential for other demonstrations based on passive memristive circuits, though was performed with much cruder precision^{19,39}. A comparable tuning accuracy was demonstrated in ref.⁴⁰, though for less dense but much more robust to variations 1T1R structures, in which each memory cell is coupled with a dedicated low-variation transistor. Furthermore, memristors can be retuned multiple times without noticeable aging—see Supplementary Note 2 for more details.

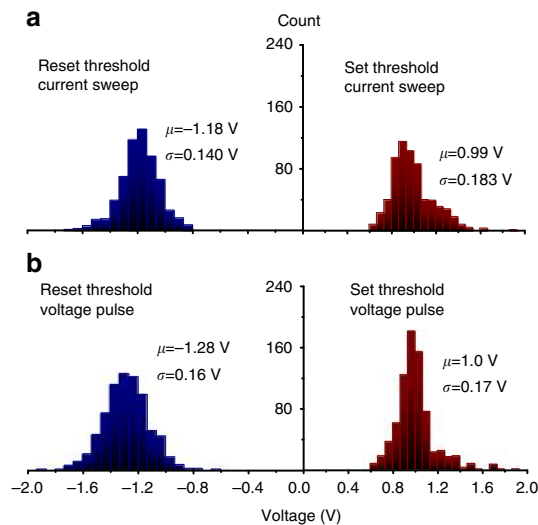


Fig. 2 Set and reset threshold statistics. The data are shown for seven 20×20 -device crossbar arrays at memristor switching with **a** current and **b** voltage pulses. The set/reset thresholds are defined as the smallest voltages at which the device resistance is increased/decreased by $>5\%$ at the application of a voltage or current pulse of the corresponding polarity. The legends show the corresponding averages and standard deviations for the switching threshold distributions. Note that the variations are naturally better when only considering devices within a single crossbar circuit, and in addition, excluding memristors at the edges of the circuit, which typically contribute to the long tails of the histograms. For example, excluding these devices, μ is $1.0 \text{ V}/-1.2 \text{ V}$ and σ is $0.13 \text{ V}/0.15 \text{ V}$ for voltage controlled set/reset for one of the crossbars used in the experiment

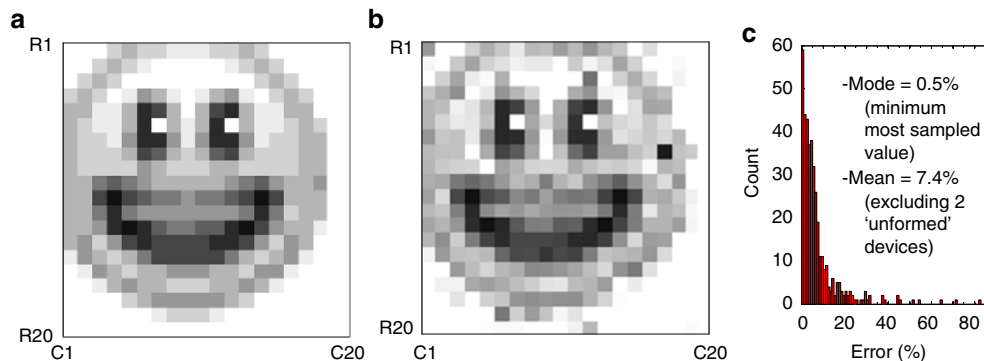


Fig. 3 High precision tuning. **a** The desired “smiley face” pattern, quantized to 10 gray levels. **b** The actual resistance values measured after tuning all devices in 20×20 memristive crossbar with the nominal 5% accuracy, using the automated tuning algorithm⁴⁸, and **c** the corresponding statistics of the tuning errors, which is defined as normalized absolute difference between the target and actual conductance values. On panel **a**, the white/black pixels correspond to $96.6 \text{ K}\Omega/7 \text{ K}\Omega$, measured at 0.2 V bias. The tuning was performed with $500\text{-}\mu\text{s}$ -long voltage pulses with amplitudes in a $[0.8 \text{ V}, 1.5 \text{ V}]/[-1.8 \text{ V}, -0.8 \text{ V}]$ range to increase/decrease device conductance. (Supplementary Fig. 3 shows absolute values of resistances and absolute error for the data on panels **b** and **c**, respectively)

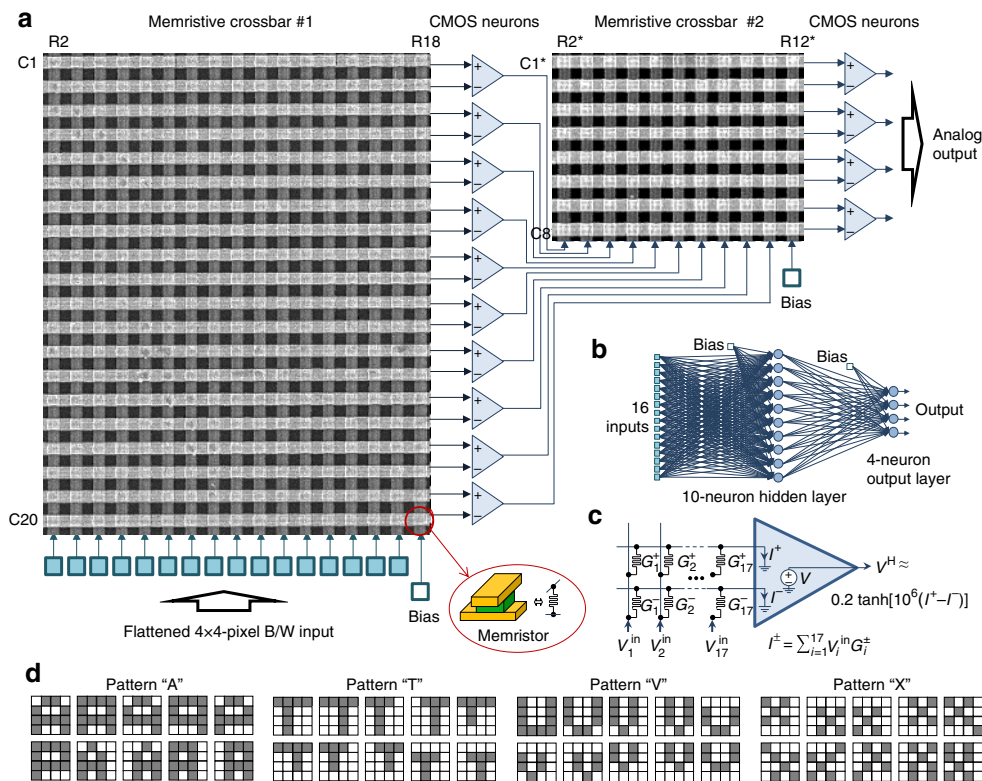


Fig. 4 Multilayer perceptron classifier. **a** A perceptron diagram showing portions of the crossbar circuits involved in the experiment. **b** Graph representation of the implemented network; **c** Equivalent circuit for the first layer of the perceptron. For clarity, only one hidden layer neuron is shown; **d** A complete set of training patterns for the 4-class experiment, stylistically representing letters “A”, “T”, “V” and “X”

Multilayer perceptron implementation. Two 20×20 crossbar circuits were packaged and integrated with discrete CMOS components on two printed circuit boards (Supplementary Fig. 2b) to implement the multilayer perceptron (MLP) (Fig. 4). The MLP network features 16 inputs, 10 hidden-layer neurons, and 4-outputs, which is sufficient to perform classification of 4×4 -pixel black-and-white patterns (Fig. 4d) into 4 classes. With account of bias inputs, the implemented neural network has 170 and 44 synaptic weights in the first and second layers, respectively.

The integrated memristors implement synaptic weights, while discrete CMOS circuitry implements switching matrix and neurons. Each synaptic weight is implemented with a pair of memristors, so that 17×20 and 11×8 contiguous subarrays were involved in the experiment (Fig. 4a), i.e., almost all of the available memristors in the first crossbar and about a quarter of the devices in the second one. The switching matrix was implemented with analog discrete component multiplexers and designed to operate in two different modes. The first one is utilized for on-board forming of memristors as well as their conductance tuning during weight import. In this operation mode, the switching matrix allows the access to any selected row and column and, simultaneously, the application of a common voltage to all remaining (half-selected) crossbar lines, including an option of floating them. The voltages are generated by an external parameter analyzer. In the second, inference mode the switching matrix connects the crossbar circuits to the neurons as shown in Fig. 4a and enables the application of ± 0.2 V inputs, corresponding to white and black pixels of the input patterns. Concurrently, the measurement of output voltages of the perceptron network is carried out. The whole setup is controlled by a general-purpose computer (Supplementary Fig. 2c).

The neuron circuitry is comprised of three distinct stages (Supplementary Fig. 2a). The first stage consists of inverting

operational amplifier, which maintains a virtual ground on the crossbar row electrodes. Its voltage output is a weighted sum between the input voltages, applied to crossbar columns (Fig. 4a), and the conductances of the corresponding crosspoint devices. The second stage op-amp computes the difference between two weighted sums calculated for the adjacent line of the crossbar. The operational amplifier’s output in this stage is allowed to saturate for large input currents, thus effectively implementing tanh-like activation function. In the third and final stage of the neuron circuit, the output voltage is scaled down to be within -0.2 V to $+0.2$ V range before applying it to the next layer. The voltage scaling is only implemented for the hidden layer neurons to ensure negligible disturbance of the state of memristors in the second crossbar array.

With such implementation, perceptron operation for the first and second layers is described by the following equations:

$$V_j^H \approx 0.2 \tanh \left[10^6 (I_j^+ - I_j^-) \right], I_j^\pm = \sum_{i=1}^{17} V_i^{\text{in}} G_{ij}^{(1)\pm} \quad (1)$$

$$V_k^{\text{out}} \approx 10^6 (I_k^+ - I_k^-), I_k^\pm = \sum_{j=1}^{11} V_j^H G_{jk}^{(2)\pm} \quad (2)$$

Here V^{in} , V^H , V^{out} are, respectively, perceptron input, hidden layer output, and perceptron output voltages. $G^{(1)\pm}$ and $G^{(2)\pm}$ are the device conductances in the first and second crossbar circuits, with \pm superscripts denoting a specific device of a differential pair, while I^\pm are the currents flowing into the corresponding neurons. j and k are hidden and output neuron indexes, while i is the pixel index of an input pattern. The additional bias inputs V_{17}^{in} and V_{11}^H are always set to $+0.2$ V.

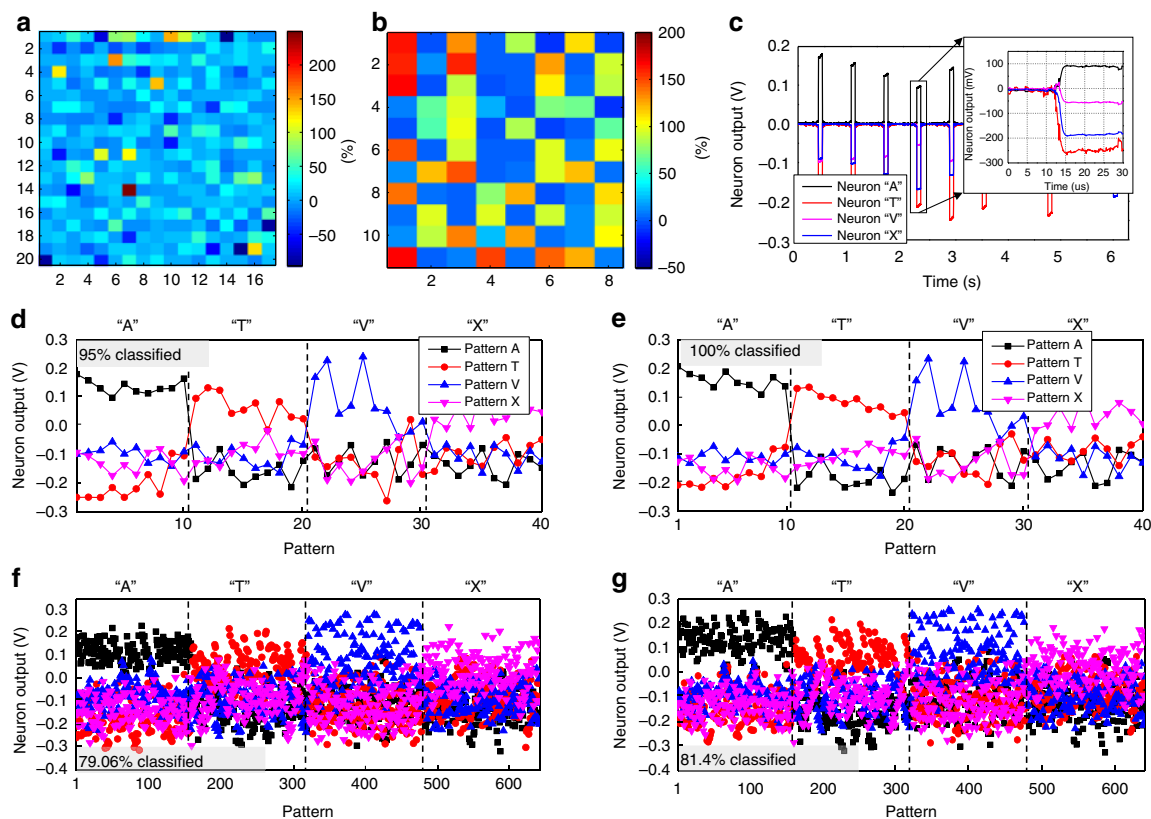


Fig. 5 Ex-situ training experimental results. **a, b** The normalized difference between the target and the actual conductances after tuning in **a** the first and **b** the second layer of the network for the hardware-oblivious training approach; **c** Time response of the trained network for 6 different input patterns, in particular showing less than 5 μ s propagation delay. Perceptron output voltage for **d, f** hardware-oblivious and **e, g** hardware-aware ex-situ training approaches, with **d-g** panels showing measured results for training/test patterns

Pattern classification. In our first set of experiments, the multi-layer perceptron was trained ex-situ by first finding the synaptic weights in the software-implemented network, and then importing the weights into the hardware. Because of limited size of the classifier, we have used custom 4-class benchmark, which is comprised of a total of 40 training (Fig. 4d) and 640 test (Supplementary Fig. 4) 4×4 -pixel black and white patterns representing stylized letters “A”, “T”, “V”, and “X”. As Supplementary Fig. 5 shows, the classes of the patterns in the benchmark are not linearly separable and the use of multi-bit (analog) weights significantly improve performance for the implemented training algorithm.

In particular, the software-based perceptron was trained using conventional batch-mode backpropagation algorithm with mean-square error cost function. The neuron activation function was approximated with tangent hyperbolic with a slope specific to the hardware implementation. We assumed a linear I - V characteristics for the memristors, which is a good approximation for the considered range of voltages used for inference operation (Fig. 1c). During the training the weights were clipped within (10 μ s, 100 μ s) conductance range, which is an optimal range for the considered memristors.

In addition, two different approaches for modeling weights were considered in the software network. In the simplest, hardware-oblivious approach, all memristors were assumed to be perfectly functional, while in a more advanced, hardware-aware approach, the software model utilized additional information about the defective memristors. These were the devices whose conductances were experimentally found to be stuck at some values, and hence could not be changed during tuning.

The calculated synaptic weights were imported into the hardware by tuning memristors’ conductances to the desired values using an automated write-and-verify algorithm⁴⁸. The

stuck devices were excluded from tuning for the hardware-aware training approach. To speed up weight import, the maximum tuning error was set to 30% of the target conductance (Fig. 5a, b), which is adequate import precision for the considered benchmark according to the simulation results (Supplementary Fig. 5). Even though tuning accuracy was often worse than 30%, the weight errors were much smaller and, e.g., within 30% for 42 weights (out of 44 total) in the second layer of the network (Supplementary Fig. 6). This is due to our differential synapses implementation, in which one of the conductances was always selected to have the smallest (i.e., 10 μ s) value and the cruder accuracy was used for tuning these devices because of their insignificant contribution to the actual weight.

After weight import had been completed, the inference was performed by applying ± 0.2 V inputs specific to the pattern pixels and measuring four analog voltage outputs. Figure 5c shows typical transient response. Though the developed system was not optimized for speed, the experimentally measured classification rate was quite high—about 300,000 patterns per second and was mainly limited by the chip-to-chip propagation delay of analog signals on the printed circuit board.

Figure 5d, e shows classification results for the considered benchmark using the two different approaches. (In both software simulations and hardware experiments, the winning class was determined by the neuron with maximum output voltage.) The generalization functionality was tested on a 640 noisy test patterns (Supplementary Fig. 4), obtained by flipping one of the pixels in the training images (Fig. 4d). The experimentally measured fidelity on a training and test set patterns for the hardware-oblivious approach were 95% and 79.06%, respectively (Fig. 5d, f), as compared to 100% and 82.34% achieved in the

software (Supplementary Fig. 5). As expected, the experimental results were much better for hardware-aware approach, i.e., 100% for the training patterns and 81.4% for the test ones (Fig. 5e, g).

It should be noted that the achieved classification fidelity on test patterns is far from ideal 100% value due to rather challenging benchmark. In our demonstration, the input images are small and addition of noise, by flipping one pixel, resulted in many test patterns being very similar to each other. In fact, many of them are very difficult to classify even for a human, especially distinguishing between test patterns ‘V’ and ‘X’.

In our second set of experiments, we have trained the network in-situ, i.e., directly in a hardware²¹. (Similar to our previous work²⁶, only inference stage was performed in a hardware during such in-situ training, while other operations, such as computing and storing the necessary weight updates, were assisted by an external computer.) Because of limitations of our current experimental setup, we implemented in-situ training using fixed-amplitude training pulses, which is similar to Manhattan rule algorithm. The classification performance for this method was always worse as compared to that of both hardware-aware and hardware-oblivious ex-situ approaches. For example, the experimentally measured fidelity for 3-pattern classification task was 70%, as compared to 100% classification performance achieved on training set using both ex-situ approaches. This is expected because in ex-situ training the feedback from read measurements of the tuning algorithm allows to effectively cope with switching threshold variations by uniquely adjusting write pulse amplitude for each memristor, which is not the case for the fixed-amplitude weight update (Supplementary Fig. 7). We expect that fidelity of in-situ trained network can be further improved using variable-amplitude implementation⁴⁹.

Discussion

We believe that the presented work is an important milestone towards implementation of extremely energy efficient and fast mixed-signal neuromorphic hardware. Though demonstrated network has rather low complexity to be useful for practical applications, it has all major features of more practical large-scale deep learning hardware—a nonlinear neuromorphic circuit based on metal-oxide memristive synapses integrated with silicon neurons. The successful board-level demonstration was mainly possible due to the advances in memristive circuit fabrication technology, in particular much improved uniformity and reliability of memristors.

Practical neuromorphic hardware should be able to operate correctly under wide temperature ranges. In the proposed circuits, the change in memristor conductance with ambient temperature (Supplementary Fig. 9) is already partially compensated by differential synapse implementation. Furthermore, the temperature dependence of I - V characteristics is weaker for higher conductive states (Supplementary Fig. 9). This can be exploited to improve robustness with respect to variations in ambient temperature, for example, by setting the device conductances within a pair to $G_{\text{BIAS}} \pm G/2$, where G_{BIAS} is some large value. An additional approach is to utilize memristor, with conductance G_M , in the feedback of the second operational amplifier stage of the original neuron circuit (Supplementary Fig. 2a). In this case, the output of the second stage is proportional to $\Sigma_i V_i^{\text{in}}(G_i^+ - G_i^-)/G_M$ with temperate drift further compensated assuming similar temperature dependence for the feedback memristor.

Perhaps the only practically useful way to scale up the neuromorphic network complexity further is via monolithical integration of memristors with CMOS circuits. Such work has already been started by several groups^{19,30}, including ours⁴⁷. We envision that the most promising implementations will be based on passive memristor technology, i.e., similar to the one demonstrated in this paper, because it is suitable for monolithical back-end-of-line integration of multiple crossbar layers⁴⁶. The three dimensional nature of such

circuits⁵⁰ will enable neuromorphic networks with extremely high synaptic density, e.g., potentially reaching 10^{13} synapses in one square centimeter for 100-layer 10-nm memristive crossbar circuits, which is only hundred times less compared to the total number of synapses in a human brain. (Reaching such extremely high integration density of synapses would also require increasing crossbar dimensions—see discussion of this point in Supplementary Note 1.)

Storing all network weights locally would eliminate overhead of the off-chip communication and lead to unprecedented system-level energy efficiency and speed for large-scale networks. For example, the crude estimates showed that energy-delay product for the inference operation of a large-scale deep learning neural networks implemented with mixed-signal circuits based on the 200-nm memristor technology similar to the one discussed in this paper could be six orders of magnitude smaller as compared to that of the advanced digital circuits, while more than eight orders of magnitude smaller when utilizing three-dimensional 10-nm memristor circuits⁵¹.

Methods

Automated forming procedure. To speed up the memristor forming, an algorithm for its automation was developed (Supplementary Fig. 1a). In general, the algorithm follows a typical manual process of applying an increasing amplitude current sweep to form a memristor. To avoid overheating during voltage controlled forming, the maximum current was limited by the current compliance implemented with external transistor connected in series with biased electrode.

In the first step of the algorithm, the user specifies a list of crossbar devices to be formed, the number of attempts, and the algorithm parameters specific to the device technology, including the initial (I_{start}) and the final minimum (I_{min}) and maximum (I_{max}) values, and step size (I_{step}) for the current sweep, the minimum current ratio (A_{min}), measured at 0.1 V, which user requires to register successful forming, reset voltage V_{reset} , and the threshold resistance of pristine devices (R_{TH}), measured at 0.1 V. The specified devices are then formed, one at a time, by first checking the pristine state of the device.

In particular, if the measured resistance of as-fabricated memristor is lower than the defined threshold value, then the device is already effectively pre-formed by annealing. In this case, the forming procedure is not required, and the device is switched into the low conducting state to reduce leakage currents in the crossbar during the forming of the subsequent devices from the list.

Alternatively, a current sweep (or voltage) is applied to the device to form the device. If forming is failed, the amplitude of the maximum current in a sweep is increased and the process is repeated. (The adjustment of the maximum sweep current is performed manually in this work but could be easily automated as well.) If the device could not be formed within allowed number of attempts, the same forming procedure is performed again after resetting all devices in the crossbar to the low conductive states. The second try could still result in successful forming, if the failure to form in the first try was because of large leakages via on-state memristors that were already formed. Even though all formed devices are reset immediately after forming, some of them may be accidentally turned on during forming of other devices. Finally, if a device could not be formed within allowed number of attempts for the second time, it is recorded as defective.

Experimental setup. Supplementary Fig. 2 shows additional details of the MLP implementation and the measurement setup. We have used AD8034 discrete operational amplifiers for the CMOS-based neurons and ADG1438 discrete analog multiplexers to implement on-board switch matrix.

Data availability. The data that support the plots within this paper and other findings of this study are available from the corresponding author upon reasonable request.

Received: 28 November 2017 Accepted: 2 May 2018

Published online: 13 June 2018

References

1. LeCun, Y., Bengio, Y. & Hinton, G. Deep learning. *Nature* **521**, 436–444 (2015).
2. Schmidhuber, J. Deep learning in neural networks: an overview. *Neural Netw.* **61**, 85–117 (2015).
3. Krizhevsky, A., Sutskever, I. & Hinton, G. Imagenet classification with deep convolutional neural networks. *Proc. Adv. Neural Inf. Process. Syst.* **12**, 1097–1105 (2012).

4. NVIDIA. GP100 Pascal Whitepaper. *NVIDIA.com* <https://images.nvidia.com/content/pdf/tesla/whitepaper/pascal-architecture-whitepaper.pdf> (2016).
5. Chen, Y. H., Krishna, T., Emer, J. S. & Sze, V. Eyeriss: an energy-efficient reconfigurable accelerator for deep convolutional neural networks. *IEEE J. Solid-State Circuits* **52**, 127–138 (2017).
6. Moons, B., Uytterhoeven, R., Dehaene, W. & Verhelst, M. in *IEEE International Solid-State Circuits Conference (ISSCC)* 246–257 (IEEE, 2017).
7. Jouppi, N. P. et al. in *Proc. of the 44th Annual International Symposium on Computer Architecture* 1–12 (ACM, 2017).
8. Merolla, P. A. et al. A million spiking-neuron integrated circuit with a scalable communication network and interface. *Science* **345**, 668–673 (2014).
9. Benjamin, B. V. et al. Neurogrid: a mixed-analog-digital multichip system for large-scale neural simulations. *Proc. IEEE* **102**, 699–716 (2014).
10. Furber, S. B., Galluppi, F., Temple, S. & Plana, S. The SpiNNaker project. *Proc. IEEE* **102**, 652–665 (2014).
11. Indiveri, G. et al. Neuromorphic silicon neuron circuits. *Front. Neurosci.* **5**, 73 (2011).
12. Likharev, K. K. CrossNets: neuromorphic hybrid CMOS/nanoelectronic networks. *Sci. Adv. Mat.* **3**, 322–331 (2011).
13. Hasler, J. & Marr, H. B. Finding a roadmap to achieve large neuromorphic hardware systems. *Front. Neurosci.* **7**, 118 (2013).
14. Chakrabartty, S. & Cauwenberghs, G. Sub-microwatt analog VLSI trainable pattern classifier. *IEEE J. Solid-State Circuits* **42**, 1169–1179 (2007).
15. George, S. et al. A programmable and configurable mixed-mode FPAA SoC. *IEEE Trans Very Large Scale Integr. Syst.* **24**, 2253–2261 (2016).
16. Merrikh Bayat, F. et al. High-performance mixed-signal neurocomputing with nanoscale floating-gate memory cell arrays. *IEEE Trans. Neural Netw. Learn. Syst.* <https://doi.org/10.1109/TNNLS.2017.2778940> (2018).
17. Mead, C. *Analog VLSI and Neural Systems* (Addison-Wesley Longman Publishing Co. Inc., Boston, MA, USA, 1989).
18. Sarpeshkar, R. Analog versus digital: extrapolating from electronics to neurobiology. *Neural Comput.* **10**, 1601–1638 (1998).
19. Kim, K. H. et al. A functional hybrid memristor crossbar-array/CMOS system for data storage and neuromorphic applications. *Nano. Lett.* **12**, 389–395 (2011).
20. Suri, M. et al. in *2012 International Electron Devices Meeting* 235–238 (IEEE, 2012).
21. Alibart, F., Zamanidoost, E. & Strukov, D. B. Pattern classification by memristive crossbar circuits using ex situ and in situ training. *Nat. Commun.* **4**, 2072 (2013).
22. Eryilmaz, S. B. et al. Brain-like associative learning using a nanoscale non-volatile phase change synaptic device array. *Front. Neurosci.* **8**, 205 (2014).
23. Kaneko, Y., Nishitani, Y. & Ueda, M. Ferroelectric artificial synapses for recognition of a multishaded image. *IEEE Trans. Electron Devices* **61**, 2827–2833 (2014).
24. Piccolboni, G. et al. in *2015 International Electron Devices Meeting (IEDM)* 447–450 (IEEE, 2015).
25. Kim, S. et al. in *2015 International Electron Devices Meeting (IEDM)* 443–446 (IEEE, 2015).
26. Prezioso, M. et al. in *2015 International Electron Devices Meeting (IEDM)* 455–458 (IEEE, 2015).
27. Li, C. et al. Analogous signal and image processing with large memristor crossbars. *Nat. Electron.* **1**, 52–59 (2017).
28. Chu, M. et al. Neuromorphic hardware system for visual pattern recognition with memristor array and CMOS neuron. *IEEE Trans. Ind. Electron.* **62**, 2410–2419 (2015).
29. Hu, S. G. et al. Associative memory realized by a reconfigurable memristive Hopfield neural network. *Nat. Commun.* **6**, 7522 (2015).
30. Yu, S. et al. in *2016 International Electron Devices Meeting (IEDM)* 416–419 (IEEE, 2016).
31. Hu, M., Strachan, J. P., Li, Z. & Williams, R. S. in *2016 17th International Symposium on Quality Electronic Design (ISQED)* 374–379 (ISQED, 2016).
32. Emelyanov, A. V. et al. First steps towards the realization of a double layer perceptron based on organic memristive devices. *AIP Adv.* **6**, 111301 (2016).
33. Serb, A. et al. Unsupervised learning in probabilistic neural networks with multi-state metal-oxide memristive synapses. *Nat. Commun.* **7**, 12611 (2016).
34. Burr, G. W. et al. Experimental demonstration and tolerancing of a large-scale neural network (165000 synapses) using phase-change memory as the synaptic weight element. *IEEE Trans. Electron Devices* **62**, 3498–3507 (2015).
35. Ambrogio, S. et al. Neuromorphic learning and recognition with one-transistor-one-resistor synapses and bistable metal oxide RRAM. *IEEE Trans. Electron Devices* **63**, 1508–1515 (2016).
36. Choi, S., Shin, J. H., Lee, J., Sheridan, P. & Lu, W. D. Experimental demonstration of feature extraction and dimensionality reduction using memristor networks. *Nano. Lett.* **17**, 3113–3118 (2017).
37. Wang, Z. et al. Memristors with diffusive dynamics as synaptic emulators for neuromorphic computing. *Nat. Mater.* **16**, 101–108 (2017).
38. van de Burgt, Y. et al. A non-volatile organic electrochemical device as a low-voltage artificial synapse for neuromorphic computing. *Nat. Mater.* **16**, 414–418 (2017).
39. Sheridan, P. M. et al. Sparse coding with memristor networks. *Nat. Nanotechnol.* **12**, 784–789 (2017).
40. Yao, P. et al. Face classification using electronic synapses. *Nat. Commun.* **8**, 15199 (2017).
41. Boyn, S. et al. Learning through ferroelectric domain dynamics in solid-state synapses. *Nat. Commun.* **8**, 14736 (2017).
42. Yang, J. J., Strukov, D. B. & Stewart, D. R. Memristive devices for computing. *Nat. Nanotechnol.* **8**, 13–24 (2013).
43. Wong, P. H.-S. et al. Metal-oxide RRAM. *Proc. IEEE* **100**, 1951–1970 (2012).
44. Govoreanu, B. et al. in *2011 International Electron Devices Meeting* 729–732 (IEEE, 2011).
45. Gao, B. et al. Ultra-low-energy three-dimensional oxide-based electronic synapses for implementation of robust high-accuracy neuromorphic computation systems. *ACS Nano* **8**, 6998–7004 (2014).
46. Adam, G. C. et al. 3D memristor crossbars for analog and neuromorphic computing applications. *IEEE Trans. Electron Devices* **64**, 312–318 (2017).
47. Chakrabarti, B. et al. A multiply-add engine with monolithically integrated 3D memristor crossbar/CMOS hybrid circuit. *Nat. Sci. Rep.* **7**, 42429 (2017).
48. Alibart, F., Gao, L., Hoskins, B. D. & Strukov, D. B. High precision tuning of state for memristive devices by adaptable variation-tolerant algorithm. *Nanotechnology* **23**, 075201 (2012).
49. Kataeva, I. et al. in *The International Joint Conference on Neural Networks* 1–8 (IEEE, 2015).
50. Strukov, D. B. & Williams, R. S. Four-dimensional address topology for circuits with stacked multilayer crossbar arrays. *Proc. Natl Acad. Sci. USA* **106**, 20155–20158 (2009).
51. Ceze, L. et al. in *2016 74th Annual Device Research Conference (DRC)* 1–2 (IEEE, 2016).

Acknowledgements

This work was supported by DARPA under contract HR0011-13-C-0051UPSIDE via BAE Systems, Inc., by NSF grant CCF-1528205, and by the DENSO CORP., Japan. Useful discussions with G.Adam, B.Hoskins, X.Guo and K.K. Likharev are gratefully acknowledged.

Author contributions

F.M.B., M.P., I.K. and D.S. conceived the original concept and initiated the work. M.P. and B.C. fabricated devices. F.M.B., M.P., B.C. and I.K. developed the characterization setup. F.M.B., M.P., B.C. and H.N. performed measurements. F.M.B., I.K. and D.S. performed simulations and estimated performance. D.S. wrote the manuscript. All discussed results.

Additional information

Supplementary Information accompanies this paper at <https://doi.org/10.1038/s41467-018-04482-4>.

Competing interests: The authors declare no competing interests.

Reprints and permission information is available online at <http://npg.nature.com/reprintsandpermissions/>

Publisher's note: Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.



Open Access This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons license, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons license, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons license and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this license, visit <http://creativecommons.org/licenses/by/4.0/>.

© The Author(s) 2018