

NEUROMORPHIC COMPUTATION

Sparse codes from memristor grids

The adjustable resistive state of memristors makes it possible to implement sparse coding algorithms naturally and efficiently.

Bruno A. Olshausen and Christopher J. Rozell

Each waking moment, the human brain is inundated with nearly a gigabit per second of image data from the eyes. Massive cortical circuits within our occipital lobes must efficiently process this data stream to sift out the relevant pieces of information. Although neuroscience is still only beginning to understand the detail of how these cortical circuits work, some general principles underlying the process were proposed decades ago. One such principle is sparse coding, an idea initially proposed in the early 1970s by the neuroscientist Horace Barlow¹.

Barlow hypothesized — based both on experimental observation and on theoretical grounds — that the cortex reformats sensory data to form a complete representation with the smallest number of active neurons. Importantly, this reformatting is adapted to the statistics of the input stream so that neurons become selective to commonly occurring patterns. In this way, a sparse code may facilitate learning and the forming of associations at higher levels of processing by making explicit the structure and features occurring in the input signal, with a high efficiency in terms of the energy consumption.

In the mid-1990s, it was shown through computer simulation of a neural network trained on a large database of natural image

patches that sparse coding could account for the shapes of receptive fields of neurons in primary visual cortex². Since then, the idea has not only gained acceptance in neuroscience but has also become a widely adopted tool in modern signal processing³ and computer vision⁴, and it forms the basis of a number of deep unsupervised-learning algorithms⁵. Also in view of its energy efficiency, sparse coding is an extremely promising approach to deal with the data streams we face in modern technology.

A key challenge in employing sparse coding is that finding optimal codes for each input is a highly nonlinear, neural computation that also changes as the system learns by exposure to more data. In particular, the fundamental computations — when carried out on an ordinary central processing unit — involve a time-consuming, iterative procedure consisting of repeated application of inner-products between neural weight vectors and input error signals. Although a central processing unit is reprogrammable, allowing adaptation with learning, it is particularly inefficient for the required iterative computations. This is because the signal flows occurring in a network of highly interconnected, asynchronous neurons need to be digitally simulated. In contrast, asynchronous analog circuits may be more efficient at performing

the computation, but are generally not able to adapt the circuit structure to learn with exposure to new inputs.

Now, writing in *Nature Nanotechnology*, Sheridan *et al.* report on a promising approach to this problem by exploiting memristors, two-terminal devices whose resistance value can be dynamically adjusted analogously to a synapse between neurons⁶. The researchers fabricate a network of artificial neurons fully interconnected via a 32×32 crossbar array of WO_x -based analog memristors. When a set of voltages is applied as input, the network allows for a natural and immediate computation of inner products by weighting each voltage by the corresponding synaptic conductance value via Ohm's law and by subsequently summing each resulting current via Kirchhoff's law. Sheridan *et al.* program a set of learned image features in the memristors' conductance values, which allows the circuit as a whole to compute the sparse code of an image when applied to the input.

The hardware discussed by Sheridan *et al.* implements a particular sparse coding algorithm called a locally competitive algorithm⁷, which has a straightforward mapping onto a recurrent neural network. Each neuron computes an inner product between its weight vector and the input, the result is passed through a threshold, and the output is fed back through the weights to reconstruct the input. The resulting error is treated as the next input signal and the process is repeated until the network converges to the optimal sparse code. In particular, after convergence, only a few output neurons corresponding to features contained in the input image are left active.

While ideally the network dynamics evolve in continuous-time, Sheridan *et al.* implement the feedforward and feedback phases by explicitly clocking each phase. In that sense, this constitutes something of a hybrid approach in which each neuron's inner product is computed in an analog manner, but the network as a whole settles on a solution in a step by step, clocked fashion.

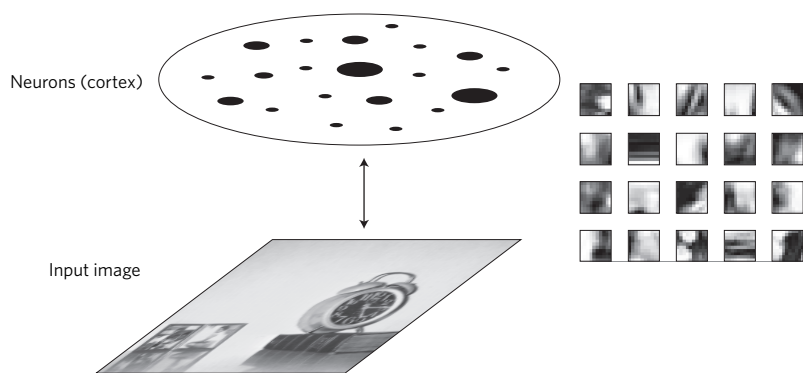


Figure 1 | Sparse coding of an image. Sparse coding transforms data (in this case, an image) into a representation in which only a small fraction of units (neurons) are active — the size of each dot denotes that unit's activity. The data is thus compactly represented in terms of a small number of elementary patterns (shown on the right). Because sparse codes use few active units, they also consume less energy.

The contribution by Sheridan *et al.* represents a critical missing piece required for a system that combines the computational efficiency previously seen in implementations of asynchronous analog circuits for sparse coding⁸ with the ability to learn the statistics of new data as it streams into the system. This combination could open a new avenue for computing approaches that merge the efficiency of neuromorphic computing with the data-driven learning that is critical to modern machine learning. Furthermore, beyond learning the statistics of data, advances toward such adaptive systems may allow the use of on-chip learning⁹ to achieve

highly effective computations in analog systems by allowing chips to compensate for the significant device variations that are inherently present in the manufacturing process. To this end, sparse coding models that utilize fully local learning rules¹⁰ appear particularly promising. □

Bruno A. Olshausen is at Helen Wills Neuroscience Institute and School of Optometry, University of California, Berkeley, California 94720, USA. Christopher J. Rozell is in the School of Electrical and Computer Engineering, Georgia Institute of Technology, Georgia 30332, USA. e-mail: baolshausen@berkeley.edu; crozell@gatech.edu

References

1. Barlow, H. B. *Perception* **1**, 371–394 (1972).
2. Olshausen, B. A. & Field, D. J. *Nature* **381**, 607–609 (1996).
3. Elad, M., Figueiredo, M. A. T. & Ma, Y. *Proc. IEEE* **98**, 972–982 (2010).
4. Wright, J. *et al. Proc. IEEE* **98**, 1031–1044 (2010).
5. Zeiler, M. D., Taylor, G. W. & Fergus, R. in *2011 IEEE Int. Conf. Computer Vision* <http://doi.org/fzcf3f> (2011).
6. Sheridan, P. M. *et al. Nat. Nanotech.* **12**, 784–789 (2017).
7. Rozell, C. J., Johnson, D. H., Baraniuk, R. G. & Olshausen, B. A. *Neural Comput.* **20**, 2526–2563 (2008).
8. Shapero, S., Charles, A. S., Rozell, C. & Hasler, P. *IEEE J. Em. Sel. Top. C* **2**, 530–541 (2012).
9. Prezioso, M. *et al. Nature* **521**, 61–64 (2015).
10. Zylberberg, J., Murphy, J. T. & DeWeese, M. R. *PLoS Comput. Biol.* **7**, e1002250 (2011).

Published online: 22 May 2017

PHOTODETECTORS

A heated junction

Resonant photonic structures made of thermoelectric materials can convert light into electricity without wavelength limitations.

Ming Zhou and Zongfu Yu

Photodetection is an essential process in cameras, light sensors, solar cells and optical communications, covering a wide spectral range from the visible to the far infrared. Most conventional photodetectors use either semiconductors or resistive bolometers. Semiconductor photodetectors are fast and sensitive but are unresponsive to energies below the semiconductor bandgap. Below this cutoff wavelength, resistive bolometry is used, providing efficient photodetection, but with slow speed. Writing in *Nature Nanotechnology*, Mauser *et al.* now report a new type of photodetector that combines the thermoelectric (TE) effect and plasmonic resonance, potentially offering fast photodetection without a cutoff wavelength¹.

Figure 1 illustrates the structure of the resonant TE photodetector. The central nanowire made of bismuth telluride and antimony telluride is the photoactive region. This is connected to two metallic pads of p- and n-doped TE materials, acting as positive and negative electrodes, respectively. The nanowire is specifically designed to realize a unique thermophotonic function: it exploits the doped electrodes to form a TE junction in the middle; and, at the same time, it is shaped to support a guided optical resonance. When in operation, the optical resonance converts incident light

into localized heat, which then drives the TE junction to produce an electrical signal.

This occurs because a temperature gradient creates a heat flow that carries electrons and holes. In p-doped materials, the heat flow carries holes from the hot to cold region, whereas in n-doped materials, electrons move from the hot to the cold region. By using both n- and p-doped materials, the TE junction generates a voltage when the temperature in the middle is higher than that at the two ends (Fig. 1).

However, a TE junction is generally insensitive to incident light. To use a TE junction for photodetection, light must be first converted to localized heat. And this is where the optical properties of the nanowire come into play. The nanowire is a nanoscale resonator and can collect incident light from an area much larger than its geometrical cross section. Such a concentration effect, already exploited² in semiconductors for photodetectors, single molecule imaging and solar cells, is used by Mauser *et al.* in a metallic nanowire to convert light to heat through ohmic losses. The pads at the two ends of the device are highly reflective, and hence remain cool, creating a temperature gradient that drives the voltage across the TE junction.

This device offers a few advantages compared to conventional photodetectors. It can work over an extremely broad spectral

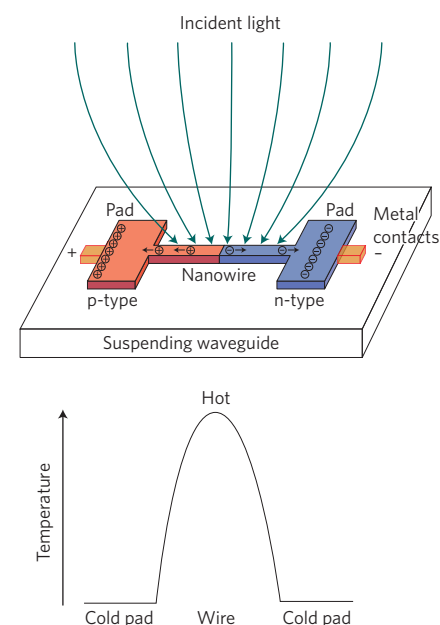


Figure 1 | Schematic of the resonant thermoelectric (TE) photodetector and its light-to-heat conversion mechanism. A nanowire made of TE p- and n-doped materials converts light into localized heat. The pads reflect most of the incident light and stay cool. The heat flow carries electrons to the n-doped pad and holes to the p-doped pad, generating an electrical voltage difference across the TE junction.