

Revised course work weights

Midterm 1 30% → 40%
 Midterm 2 30% → removed
 Project 40% → 60%

Proposal PPT presentation Final Report (written)

Recent topics

- power comparison of different architectures (parallel, pipeline, ...)
- Adiabatic circuits & energy recycling
- Multiple- V_{DD} designs & voltage level shifting
- Memristors
- Neuromorphic computing

Fig. 7.8 Prewitt edge detector. (a) Horizontal kernel (f_x). (b) Vertical kernel (f_y).

(a)	(b)
$\begin{bmatrix} -1 & 0 & 1 \\ -1 & 0 & 1 \\ -1 & 0 & 1 \end{bmatrix}$	$\begin{bmatrix} -1 & -1 & -1 \\ 0 & 0 & 0 \\ 1 & 1 & 1 \end{bmatrix}$
Horizontal (f_x)	Vertical (f_y)

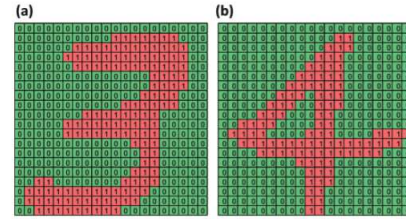


Fig. 7.9 Original 28×28 pixels of black-and-white MNIST handwritten digits: (a) "3" and (b) "4".

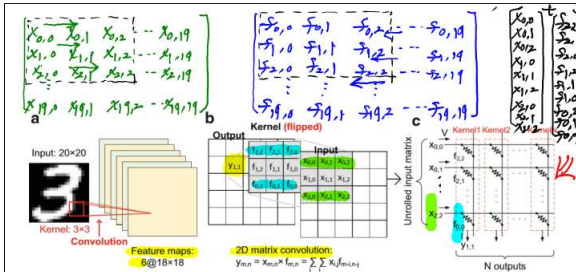


Fig. 7.7 (a) Input stage for CNN architecture. (b) 2D matrix convolution. (c) Implementation of convolution for multiple feature maps into a cross-point architecture by reduction of 2D kernel matrix into 1D column vector

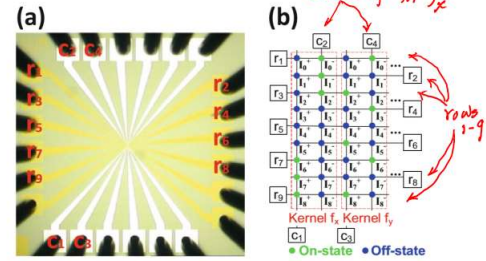


Fig. 7.11 (a) The microscopic top-view image of the fabricated 42×42 cross-point array. The probe card tips are touched on the pads. (b) The implementation of the Prewitt horizontal kernel (f_x) by programming the cells of C_1 and C_2 columns into on-state and off-state and the Prewitt vertical kernel (f_y) by programming the cells of C_3 and C_4 columns into on-state and off-state

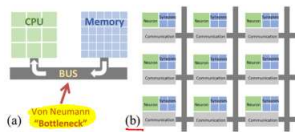


Fig. 11.3 In the Von Neumann architecture (a), data (both operations and operands) must move to and from the dedicated central processing unit (CPU) along a bus. In contrast, in a new Von Neumann architecture (b), distributed computation takes place at the location of the data, reducing the time and energy spent moving data around (Adapted from Burr et al. [11])

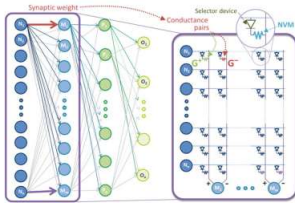


Fig. 11.2 Neuron-inspired non-Von Neumann computing (11.4) in which neurons activate each other through dense networks of programmable synaptic weights, can be implemented using dense crossbar arrays of nonvolatile memory (NVM) and selector device pairs (Adapted from Burr et al. [11])

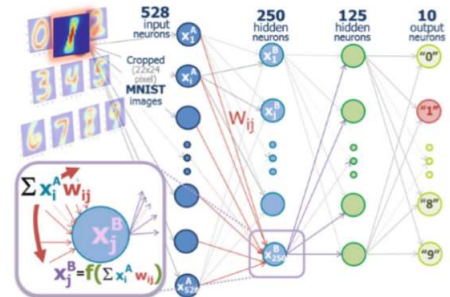
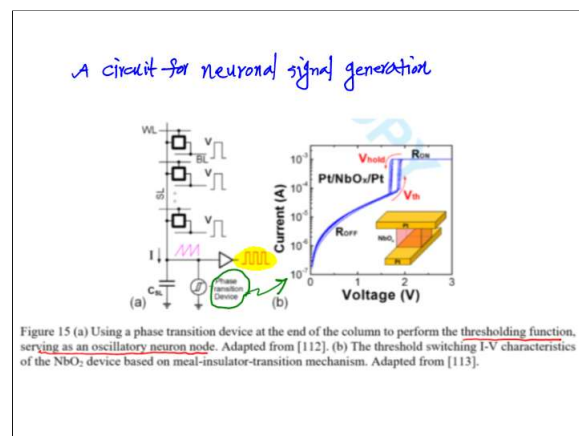
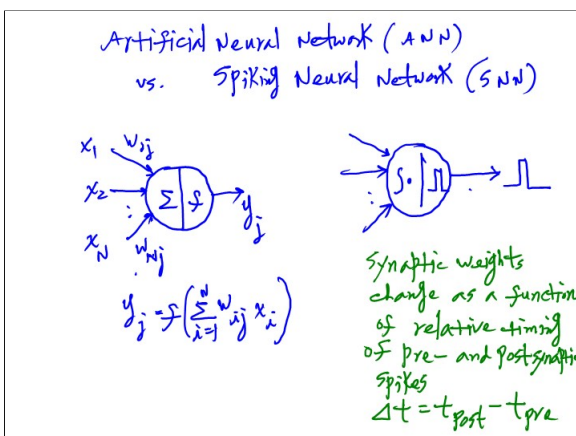
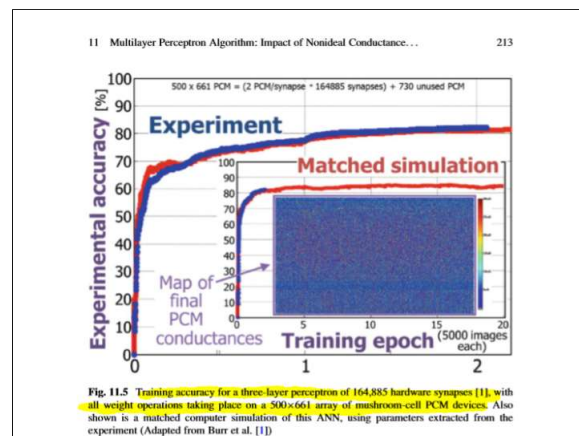
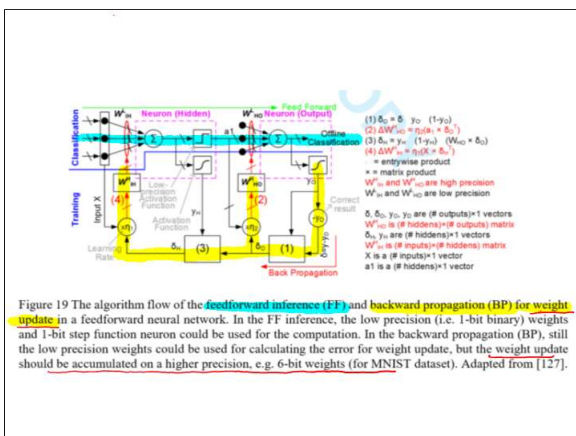
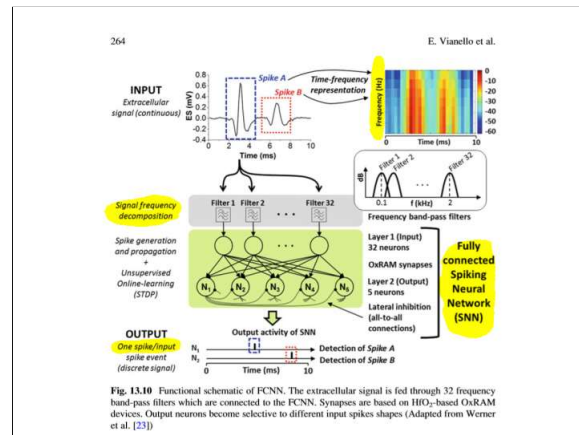
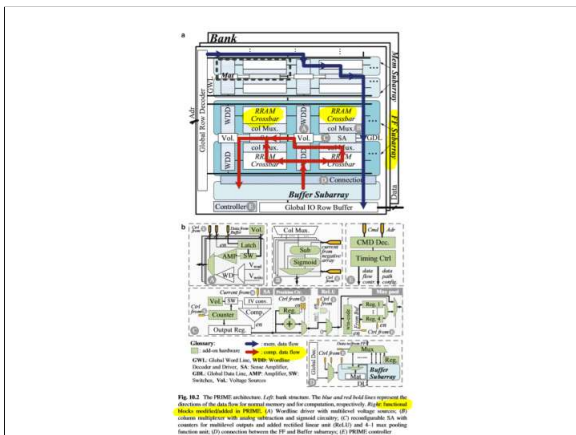


Fig. 11.4 In forward evaluation of a multilayer perceptron, each layer's neurons drive the next layer through weights w_{ij} and a nonlinearity $f(\cdot)$. Input neurons are driven by input (for instance, pixels from successive MNIST images (cropped to 28×28)). The ten output neurons classify which digit was presented (Adapted from Burr et al. [11])



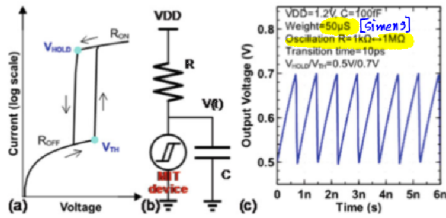


Fig. 9.12 (a) Hysteresis I-V characteristics of a 1T1R device. (b) Circuit configuration of an oscillation neuron node with 1T1R device and RRAM synaptic weight. (c) SPICE simulation waveform of the oscillation neuron

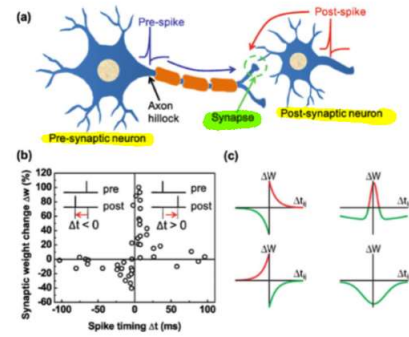


Fig. 2.9 Synaptic plasticity. (a) Relative timings of neuronal spikes from the presynaptic neuron and the postsynaptic neuron determine the weight change in synapse. (b) Synaptic weight change is plotted as a function of relative timing of pre- and post-spikes (Reprinted with permission from Bi and Poo [82]). (c) Diverse forms of STDP (Reprinted with permission from Shewell et al. [76]).

Spike Timing Dependent Plasticity (STDP)

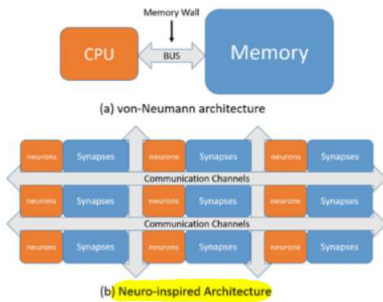


Fig. 1.1 A revolutionary shift of the computing paradigm from the computation-centric (von Neumann architecture) to the data-centric (neuro-inspired architecture)

Table 1.1 Categories of different design options for hardware implementation of neuro-inspired computing. Representative prototypes are shown

	Off-the-shelf technologies	CMOS ASIC	Emerging resistive synaptic devices
Digital representation	GPUs [9] FPGAs [10]	TPU [13] CNN accelerators [11, 12]	Analog synapses: UCSB's 12×12 crossbar array [18] Binary synapses: ASU/Tsinghua's 16 Mb RRAM macro [19]
Spike representation	SpiNNaker [14]	Analog neuron: HICANN [15] Digital neuron: TrueNorth [16]	IBM's 256×256 PCM array with STDP neuron circuits [20]

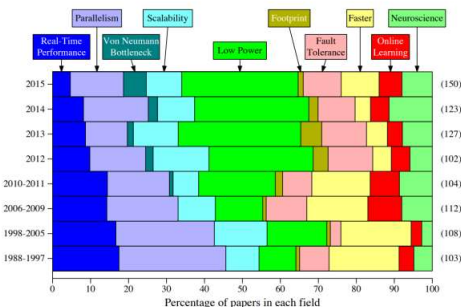


Fig. 3 Ten different motivations for developing neuromorphic systems and over time, the percentage of the papers in the literature that have indicated that motivation is one of the primary reasons they have pursued the development of neuromorphic systems

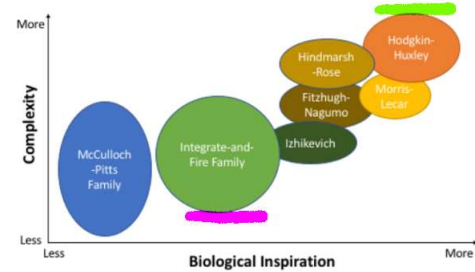


Fig. 5. A qualitative comparison of neuron models in terms of biological inspiration and complexity of the neuron model.

Ionic Channels in Hodgkin-Huxley Model

A Memristive System (Chua and Kang, Proc. of the IEEE, 1976)

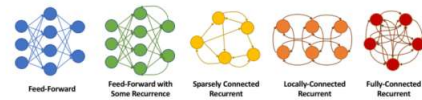
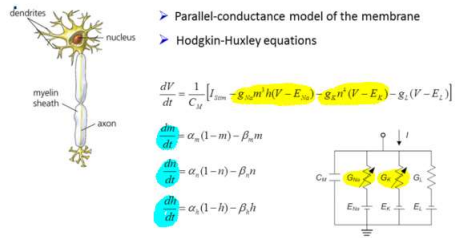


Fig. 8. Different network topologies that might be desired for neuromorphic systems. Determining the level of connectivity that is required for a neuromorphic implementation and then finding the appropriate hardware that can accommodate that level of connectivity is often a non-trivial exercise.

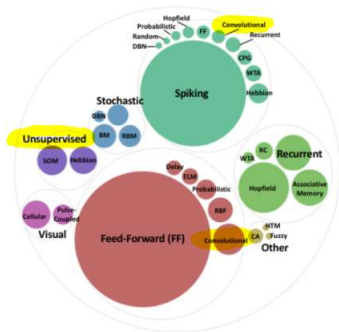


Fig. 7. A breakdown of network models in neuromorphic implementations, grouped by overall type and sized to reflect the number of associated papers.

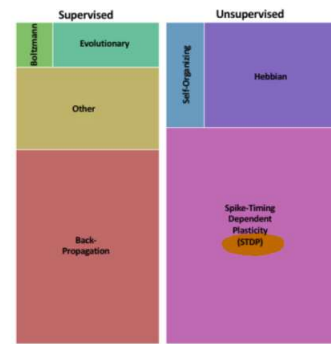


Fig. 9. An overview of on-chip training/learning algorithms. The size of the box corresponds to the number of papers in that category.

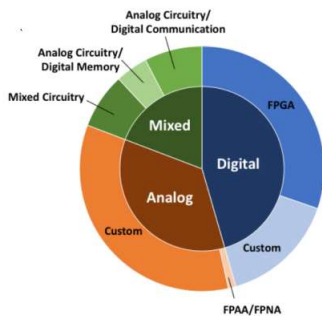


Fig. 10. An overview of hardware implementations in neuromorphic computing. These implementations are relatively basic hardware implementations and do not contain the more unusual device components discussed in Section V-B.

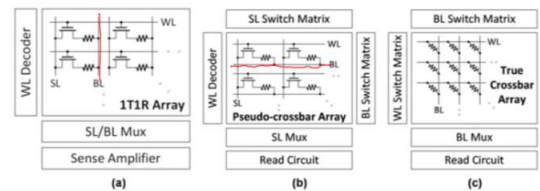


Fig. 9.1 Representative resistive NVM arrays with peripheral circuits are shown. (a) 1T1R array with row-by-row operation. (b) Pseudo-crossbar array implemented by rotating BL in 1T1R memory array. (c) True crossbar array without selector transistors for array-level parallelism

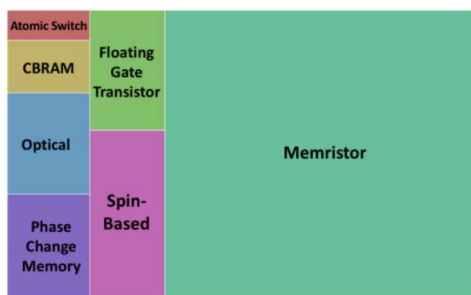


Fig. 11. Device-level components and their relative popularity in neuromorphic systems. The size of the boxes corresponds to the number of works referenced that have included those components.

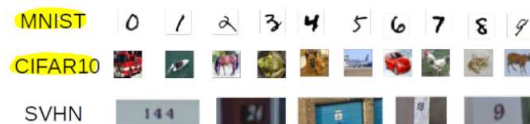


Fig. 14. Examples from different **image data sets** (MNIST [2618], CIFAR10 [2619], and SVHN [2620]) to which neuromorphic systems have been applied for classification purposes.

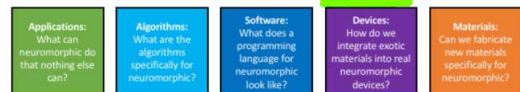


Fig. 15. Major neuromorphic computing research challenges in different fields.



Functional Circuitry on Commercial Fabric via Textile-Compatible Nanoscale Film Coating Process for Fibertronics

Hagyuul Baet, Byung Chul Jang¹§, Hongkeun Park¹, Soo-Ho Jung¹, Hye Moon Lee¹ §, Jun-Young Park¹, Seung-Bae Jeon¹, Gyeongho Son¹, Il-Woong Tchoi¹, Kyoungsik Yoo¹, Sung Gap Im¹, Sung-Yool Choi¹§, and Yang-Kyu Choi¹ §

¹School of Electrical Engineering, ²School of Chemical and Biomolecular Engineering, and ³Graphene/2D Materials Research Center, Korea Advanced Institute of Science and Technology (KAIST), 291 Daehak-ro, Yuseong-gu, Daejeon 34141, South Korea

¹ Powder and Ceramics Division, Korea Institute of Materials Science (KIMS), 797 Chanwondae-ro, Changwon, 51508, South Korea

Nano Lett. 2017, 17 (10), pp 6443–6452

DOI: 10.1021/acs.nanolett.7b03435

Publication Date (Web): September 11, 2017

Copyright © 2017 American Chemical Society

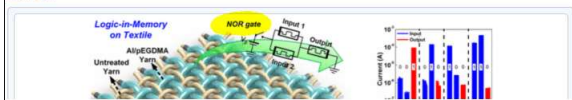
*E-mail: ykchoi@ee.kaist.ac.kr, *E-mail: sungyool.choi@kaist.ac.kr

Cite this: Nano Lett. 17, 10, 6443–6452

RIS Citation

GO

Abstract



volatile power-hungry electronic components, and modest battery storage. Here, we report a novel poly(ethylene glycol dimethacrylate) (pEGDMA)-textile memristive nonvolatile logic-in-memory circuit, enabling normally off computing, that can overcome those challenges. To form the metal electrode and resistive switching layer, strands of cotton yarn were coated with aluminum (Al) using a solution dip coating method, and the pEGDMA was conformally applied using an initiated chemical vapor deposition process. The intersection of two Al/pEGDMA coated yarns becomes a unit memristor in the lattice structure. The pEGDMA-Textile Memristor (ETM), a form of crossbar array, was interwoven using a grid of Al/pEGDMA coated yarns and untreated yarns. The former were employed in the active memristor and the latter suppressed cell-to-cell disturbance. We experimentally demonstrated for the first time that the basic Boolean functions, including a half adder as well as NOT, NOR, OR, AND, and NAND logic gates, are successfully implemented with the ETM crossbar array on a fabric substrate. This research may represent a breakthrough development for practical wearable and smart fibertronics.

The research team also demonstrated that the basic Boolean functions, including NOT, NOR, OR, AND, and NAND logic gates, were reliably implemented for data processing and storage within the ETM-based circuit array. Furthermore, they experimentally demonstrated a self-adder, which is a kind of logic functional block. It was just comprised of only 5 ETMs. The results of the study show the feasibility of the ETM circuit for energy-efficient wearable electronics.

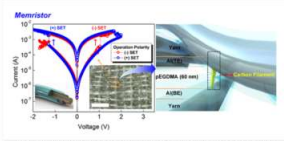


Figure 2. Current-voltage characteristics of the fabricated ETM array on the cotton substrate with a pEGDMA film thickness of 80 nm. The inset shows the optical microscope image of the fabricated ETM array and the schematic image of the Al/pEGDMA-coated yarn. The memristor device operates under low operating voltage, ranging both from -1.5 to $+1.5$ V and from -1.5 to $+1.5$ V in negative-bias SET and positive-bias SET operation, respectively. The right figure shows the formation of a conductive carbon filament bridging between top and bottom electrodes via the pEGDMA.

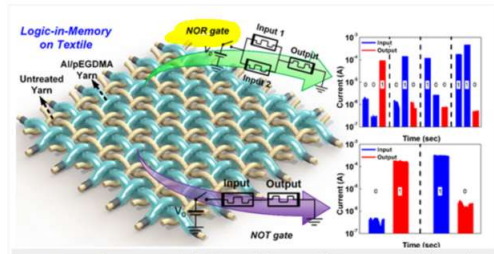


Figure 3. Schematic view of the fabricated device with logic circuits and electrical measured data of the NOR and NOT gate on fabric.

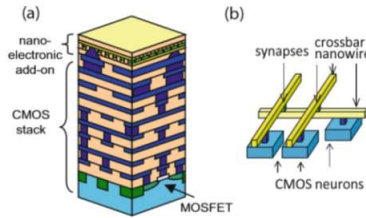


Fig. 6.1 CMOL circuits. (a) A cartoon of a hybrid CMOS/memristor integrated circuit. (b) The example of three CMOS cells (neurons) interconnected via corresponding crossbar nanowires (dendrites and axons) and cross-point memristive devices (synapses), which are located above CMOS layer

Fig. 1.2 An analogy between a biologic synapse and the resistive synaptic device

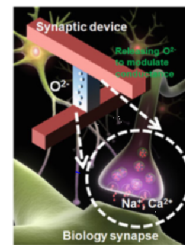
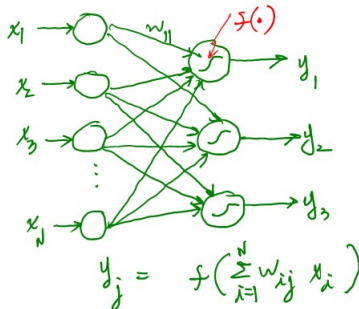


Table 1.2 Summary of the desirable performance metrics for synaptic devices

Performance metrics	Desired targets
Device dimension	<10 nm
Multilevel states' number	$>10^3$
Energy consumption	<10 fJ/programming pulse
Dynamic range	$>10^3$
Retention	>10 years ^a
Endurance	$>10^5$ updates ^a

Note: ^aThese numbers are application dependent



Linearity in Weight Update The linearity in weight update refers to the linearity of the curve between the device conductance and the number of identical programming pulses. Ideally, this should be a linear relationship for the direct mapping of the weights in the algorithms to the conductance in the devices. However, the resistive synaptic devices generally have the nonlinearity in weight update (see Fig. 1.3). The trajectory of the long-term-potential (LTP) process that increases the conductance differs from that of the long-term-depression (LTD) process that decreases the conductance. The weight tends to saturate at the end of LTP or LTD processes. This nonlinearity is undesired because the change of the weight (ΔW) depends on the current weight (W), or in other words, the weight update has a history dependence. Recent results have shown that this nonlinearity has caused the learning accuracy loss in the neural networks [41, 42].

Programming Energy Consumption The estimated energy consumption per synaptic event is around $1 \sim 10$ fJ in biological synapses. Most RRAM/CBRAM devices show a programming energy around 100 fJ ~ 10 pJ, while most PCM devices may have even higher programming energy $10 \sim 100$ pJ. The fundamental challenge is that it is much more difficult (thus paying more energy) to move the ions/defects in solid-state devices than moving calcium ions in the liquid environment in biological synapses. A back-of-envelope calculation is given as follows. In biological synapses, the spike voltage is ~ 10 mV, the ionic current ~ 1 nA, and the spike period ~ 1 ms; therefore, the energy is about 10 fJ. In resistive synaptic devices, the typical programming voltage is ~ 1 V, and the programming current is typically $> \mu$ A; although the programming speed can be accelerated less than the real time to be $< \mu$ s, still the energy is on the order of pJ. Further device engineering is thus needed to reduce the energy consumption.

Retention and Endurance During the online training, the weights are frequently updated, and the data retention requirement can be relaxed. When the training is complete, the resistive synaptic should behave as a long-term memory with a data retention in the order of 10 years at elevated temperature similarly as the

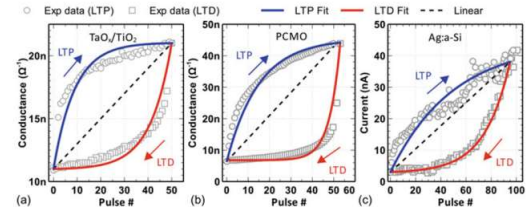


Fig. 1.3 The measured nonlinearity in the weight update reported from the literature: (a) TaO_x/TiO₂ device [39], (b) PCMO device [36], and (c) Ag-a-Si device [33]

requirement of NVM. The number of endurance is much application dependent, relying on how many weight updates are required in the training processes. For a relatively simple task (i.e., the MNIST handwritten digit recognition [43]), 60,000 training images with 50 training epochs (to repeated) give a maximum weight update possibility to be 3×10^6 updates. Actually not every synapse is updated in the training; thus, an endurance $\sim 10^4$ is sufficient for training MNIST dataset [19]. However, considering more challenging tasks (i.e., ImageNet challenge [44]), much more endurance may be required.

Uniformity and Variability Poor uniformity or significant variability in emerging NVMs is a major barrier for digital memory applications. In contrast, the neural networks promise robustness against device variations. The device variations could partially be tolerated by two mechanisms: the massive (thus maybe redundant) connections between neuron nodes by synaptic arrays and the iterative weight update process during the training. The degree of variations that can be tolerated at the system level strongly depends on the network architecture and the accuracy required by the target application. The device-algorithm co-simulations have shown the reasonable robustness against device variations in different neural networks [42, 45].

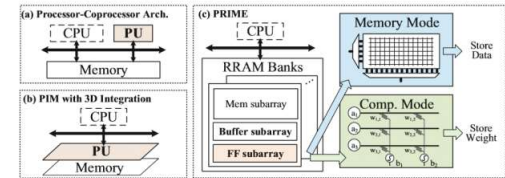
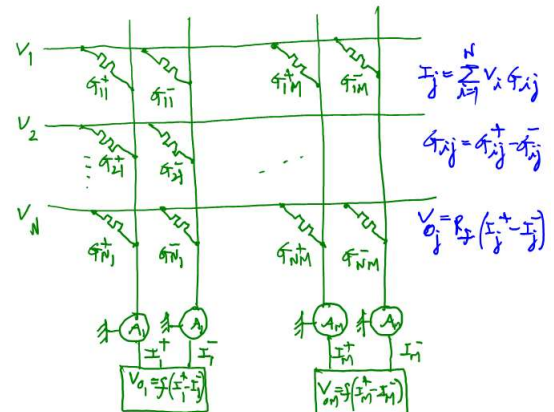
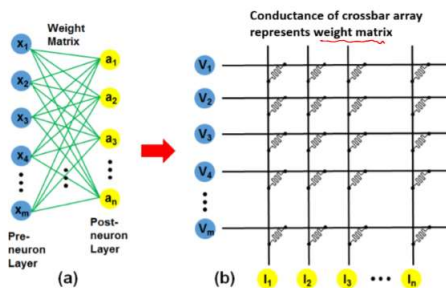


Fig. 10.1 (a) Traditional processor-coprocessor architecture with shared memory; (b) PIM architecture using 3D integration technologies; (c) PRIME design



Forward Propagation

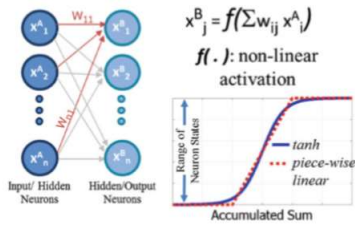


Fig. 11.21 Forward propagation operation in a deep neural network. The multiply-accumulate operation occurs on the crossbar array. Neuron circuitry must handle the nonlinear squashing function (Adapted from Fumarola et al. [9])

Reverse Propagation

Fig. 11.22 Reverse propagation operation in a deep neural network. Multiply-accumulate operation on δ occurs on the crossbar array. Neuron circuitry must handle generation and multiplication of the derivative of the squashing function (Adapted from Fumarola et al. [9])

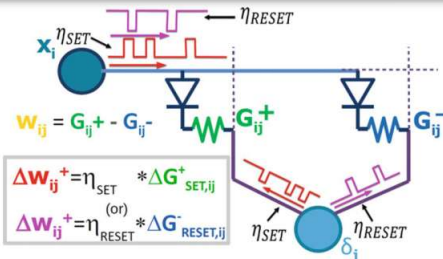
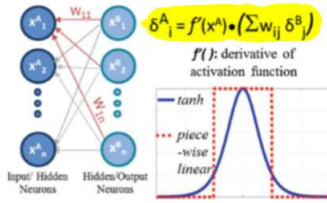


Fig. 11.18 Schematic showing crossbar-compatible [1] weight-update rule for analog bidirectional NVMs. Weight increases (decreases) can be implemented either as a SET operation on G^+ (G^-) or a RESET operation on G^- (G^+) devices. Asymmetry in the partial SET and RESET operation is compensated by applying a different learning rate parameter (η_{SET} , η_{RESET}) that modulates the number of pulses fired from the neurons into the array (Adapted from Fumarola et al. [9])

Bottom electrode (Pt)	Top electrode (W)									
	1	2	3	4	5	6	7	8	9	10
1	5.5E+3	1.9E+3 247.1E+0	6.1E+3	1.6E+3 160.3E+0	273.2E+0	4.3E+3	3.6E+3	5.3E+3		
2	3.0E+3 252.4E+0 315.9E+0	2.5E+3	1.2E+3	6.0E+3	4.2E+3	2.6E+3	3.8E+3	5.1E+3		
3	2.3E+3 1.7E+3 450.5E+0	32.2E+3	2.4E+3	5.5E+3	5.6E+3 119.3E+0 160.2E+0	4.9E+3				
4	2.3E+3	42.2E+3	1.4E+3	2.2E+3 37.2E+0	1.7E+3	5.8E+3 77.1E+0	1.70E+3	6.6E+3		
5	2.4E+3	12.2E+3	1.2E+3	4.6E+3 145.5E+0	1.1E+3	3.3E+3 134.0E+0	9.2E+3	5.5E+3		
6	3.0E+3	11.6E+3 15.7E+0	4.0E+3	1.4E+3	1.2E+3	6.6E+3	5.2E+3	2.9E+3	5.2E+3	
7	9.8E+3	5.5E+3 872.9E+0	12.2E+3	1.5E+3	1.9E+3	10.6E+3	2.0E+3	9.2E+3	4.8E+3	
8	9.4E+3	4.2E+3	11.2E+3 294.3E+0 455.9E+0	6.0E+3	5.6E+3 284.6E+0 780.4E+0	10.3E+3				
9	3.1E+3	3.5E+3	5.2E+3	2.1E+3 3978.0E+0	12.1E+3	3.8E+3	2.8E+3	4.7E+3	7.3E+3	
10	389.4E+0	16.2E+3	4.8E+3 618.9E+0 305.9E+0	7.7E+3	3.7E+3	5.0E+3 251.5E+3 313.9E+0				

On/Off ratio:

$10^1 \sim 10^2$

$10^2 \sim 10^3$

$10^3 \sim 10^4$

$10^4 \sim 10^5$

Fig. 3.6 Similar on/off ratio of 100 bits in the 1 kb PCMO-synapse array

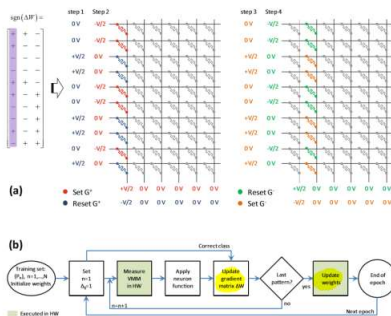
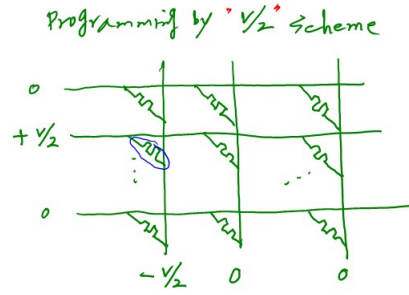


Fig. 6.9 (a) The first four steps of crossbar update. The sign of the gradient matrix (on the left), which is obtained after one epoch of training, specifies the direction of the state update for each device in crossbar circuit, i.e., whether to incrementally set or reset the device. The update is performed using the V/2 scheme with appropriate chosen voltages (on the right). The voltage shown in red/green and blue/orange are for the first and second steps, respectively. (b) Flow chart of the training algorithm. Gray boxes show the steps implemented in hardware, while all remaining steps were emulated in software



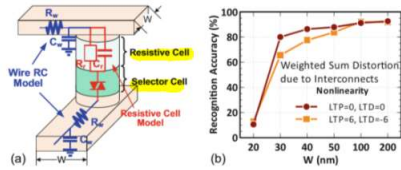
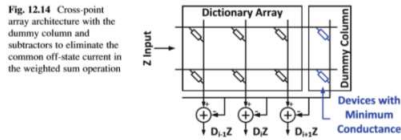


Fig. 12.15 (a) Sub-circuit module of a synaptic device cell (W is wire width). The cell consists of a resistive synaptic device and a selector. The resistive cell has capacitor (C_s) in parallel with the cell resistor (R_c). There are also wire resistors (R_w) and capacitors (C_w) for top and bottom interconnect. Sub-circuit is duplicated for the entire array to perform SPICE simulation. (b) Learning accuracy with different wire widths. Smaller wire width will degrade the learning accuracy due to the IR drop along the interconnects

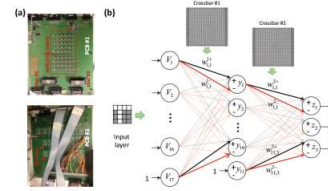


Fig. 6.11 (a) The memristive MLP fabricated on two printed circuit boards, one including two 20×20 crossbars with peripheral circuitries while the other one implements neurons. (b) High-level diagram of the implemented MLP. Each set of weights is implemented with one crossbar

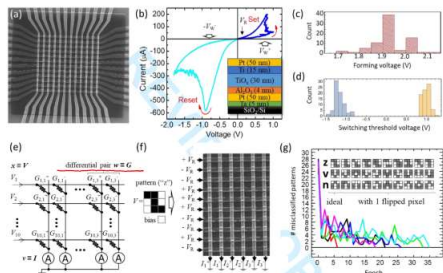
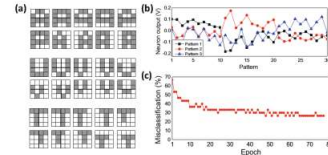


Figure 10 Perceptron classifier demonstration: (a) integrated 12×12 crossbar with an $\text{AlO}_x/\text{TiO}_x$ memristor at each cross-point; (b) typical I - V curve of a formed memristor; histograms of forming voltages (c) and effective switching threshold voltages (d) for set and reset transitions; (e) perceptron implementation using a 10×6 fragment of the memristive crossbar; (f) example of the classification operation for a specific input pattern; and (g) the convergence of network outputs, in the process of training, to the perfect (zero-error) set, for 6 different initial states. The classification was considered successful when the output signal corresponding to the correct class of the applied pattern was larger than all other outputs. The insets in panels (b) and (g) show device's cross-section and the used input pattern set, correspondingly. On panel (d), the positive / negative switching threshold voltages were defined as the smallest amplitudes of 500- μs voltage pulses that caused resistance change by more than $2 \text{ k}\Omega$ in memristors pre-set to their high / low resistive states. Adapted from [37, 101].

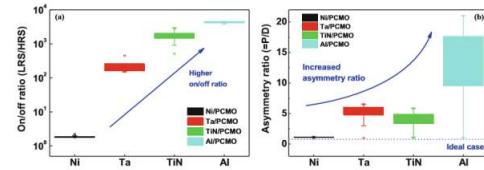


Fig. 3.13 Each sample exhibited different (a) on/off ratio and (b) asymmetry ratio

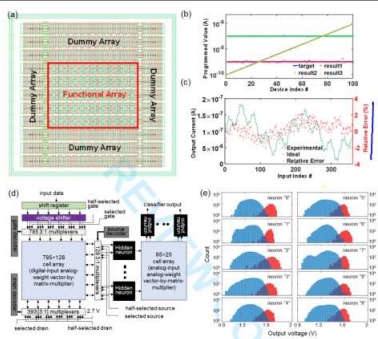


Figure 12 NOR flash memory circuits redesigned for neuro-inspired computing: (a) Layout of a 55-nm vector-matrix multiplication circuit with a 16×16 cell array and auxiliary pass-gates and (b, c) its experimental test results, for (b) cell tuning (measured vs. target weights) and (c) 4-input vector-by-vector multiplication. The four inputs are quasi-DC currents sampled from sine functions with different frequencies. 2-layer MLP based on 160-nm industrial-grade floating-gate devices: (a) high-level architecture (with the weight tuning circuitry for the 2nd array not shown for clarity), and (b) histograms of output voltages for all 10,000 MNIST test patterns. The classification of one pattern takes time below 1 μs and energy below 20 nJ. Adapted from [107, 108].

B. T. Murphy
Yield Model

$$Y = Y_0 e^{-D_0 A} > P_{\text{single transistor}}^N$$

D_0 = defect density
 A = chip area containing the N transistor circuit