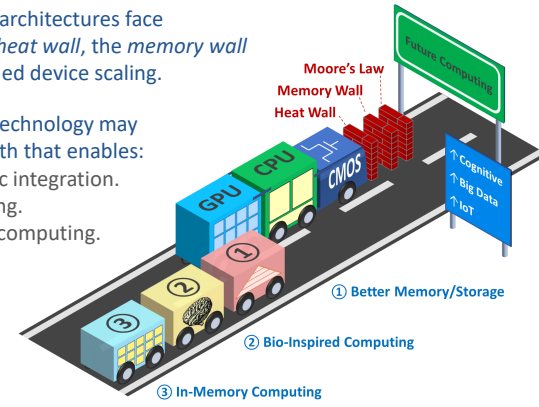# RRAM fabric for neuromorphic and in-memory computing applications

**Mohammed Zidan and Wei Lu**

**University of Michigan**
**Electrical Engineering and Computer Science**

## The Race Towards Future Computing Solutions

- Conventional computing architectures face challenges including the *heat wall*, the *memory wall* and difficulties in continued device scaling.

- Developments in RRAM technology may provide an alternative path that enables:
  - Hybrid memory–logic integration.
  - Bioinspired computing.
  - Efficient in-memory computing.

① Better Memory/Storage
② Bio-Inspired Computing
③ In-Memory Computing

*Lu Group*
*U. Michigan*

---

## Two-Terminal Memory Devices and Crossbar Arrays

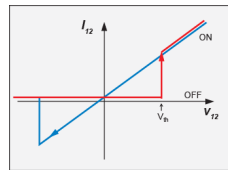**Single-cell structure**    **Crossbar Structure**

Switching Medium

Crossbar

CMOS

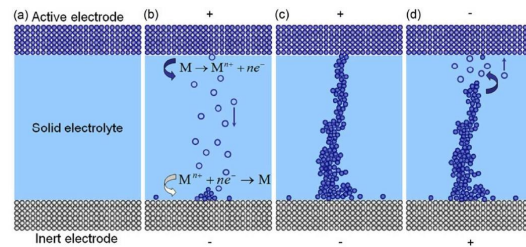**Resistive memory (RRAM)**, *memory + resistor* **(memristor)**

**Hysteretic resistive switches and crossbar structures**

- Simple structure
  - Formed by two-terminal devices
  - Not limited by transistor scaling
- Ultra-high density
  - NAND-like layout, cell size $4F^2$
  - Terabit potential
- Large connectivity
- Memory, logic/neuromorphic applications

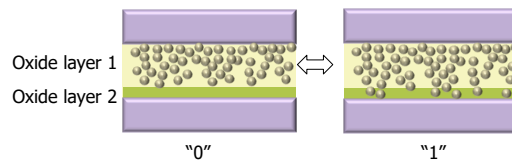$I_{12}$  ON  OFF  $V_{th}$  $V_{12}$

Lu group

*Lu Group*
*U. Michigan*

---

## Physically reconfigurable materials and devices: Resistive Memory

**ElectroChemical Metallization Cell (ECM, CBRAM)**

(a) Active electrode  (b)  (c)  (d)

$M \rightarrow M^{n+} + ne^-$

Solid electrolyte

$M^{n+} + ne^- \rightarrow M$

Inert electrode

- *Creating "new" materials on the fly*
- Active electrode material + inert dielectric
- "Filament" based on electrode material injection and redox at electrodes
- Switching layer facilitates ionic movement
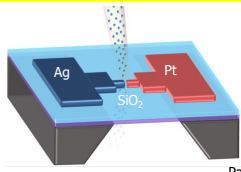
**Valency Change Cell (VCM)**

Oxide layer 1
Oxide layer 2

"0"   "1"

- *Modulating exiting material properties*
- Filament based on oxygen exchange between two oxide layers
- Electrode plays minor role
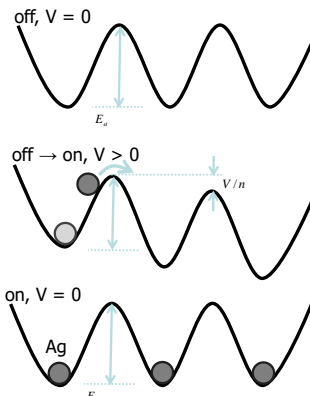
*Lu Group*
*U. Michigan*

•Ag/SiO$_2$/Pt structure, sputtered SiO$_2$ film
•The filament grows from the IE backwards toward the AE
•Branched structures were observed with wider branches pointing to the AE

Completed filament

Partially formed filaments

Ag

Pt

200nm

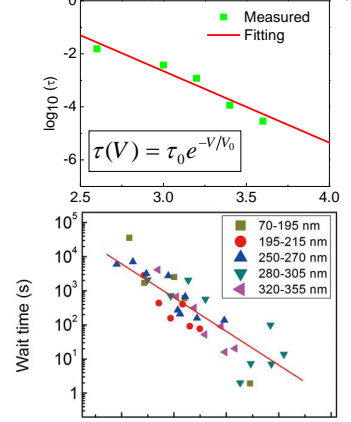**Y. Yang, Gao, Chang, Gaba, Pan, and W. Lu, Nature Communications, 3, 732, 2012.** 5

Lu group
U. Michigan

off, V = 0

$E_a$

off → on, V > 0

$V/n$

on, V = 0

Ag

$E_a$

$$v = 2d\lambda e^{-E_a/k_BT}\left(e^{qEd/2k_BT} - e^{-qEd/2k_BT}\right)$$

Measured
Fitting

$$\tau(V) = \tau_0 e^{-V/V_0}$$

$\log_{10}(\tau)$

Wait time (s)

70-195 nm
195-215 nm
250-270 nm
280-305 nm
320-355 nm

$$\Gamma = 1/\tau = \upsilon e^{-E_a^{'}(V)/k_BT}$$

$$E_a^{'}(V) = E_a - V/2n$$
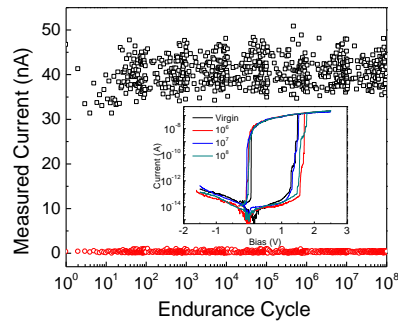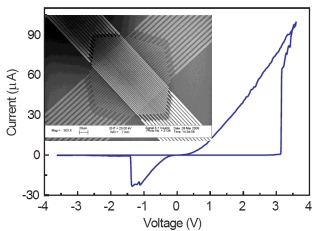
**Jo, Kim, Lu Nano Lett. 9, 496-500 (2009).**

•Filament formation is a thermally activated process.
•Activation energy reduced by applied bias.
•Speed is a ca. *exponential* function of voltage.

Current (μA)

Voltage (V)

Output (V)

Time (uS)

Endurance Cycles —— 10$^4$, —— 10$^5$
—— 10$^6$, —— 10$^7$

Lu

Measured Current (nA)

Endurance Cycle

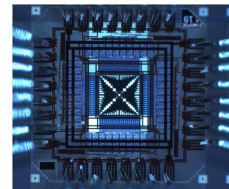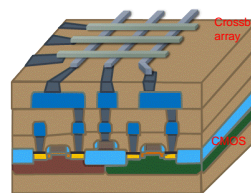Current (A)

Bias (V)

Virgin
10$^6$
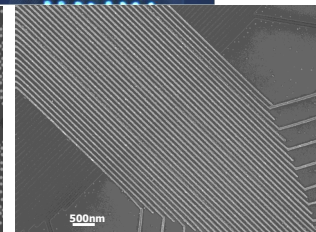10$^7$
10$^8$

➢ 1e6 on/off
➢ 1e8 W/E endurance
➢ Switching speed ~10ns

**Jo, Kim, W. Lu, Nano Lett., 8, 392 (2008)**
**Kim, Jo, W. Lu, Appl. Phys. Lett. 96, 053106 (2010)**

Lu Group
U. Michigan

Crossbar array

CMOS

•Low-temperature process, RRAM array fabricated on top of CMOS
•CMOS provides address mux/demux
•RRAM array: 100nm pitch, 50nm linewidth with density of 10Gbits/cm$^2$
•CMOS units – larger but fewer units needed. 2n CMOS cells control n$^2$ memory cells

500nm

**Kim, Gaba, Wheeler, Cruz-Albrecht, Srivinara, W. Lu Nano Lett., 12, 389–395 (2012).**

8

Lu Group
U. Michigan

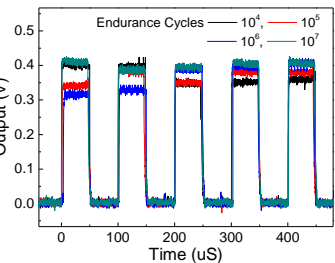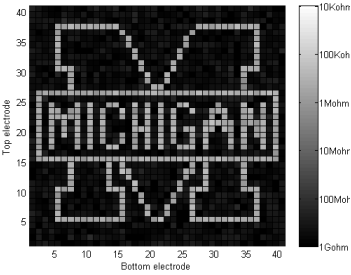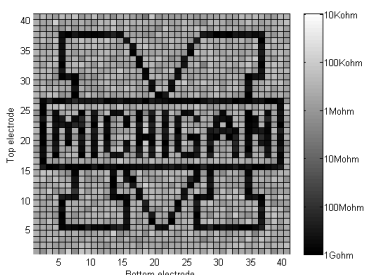- **Crossbar array operation, array written followed by read**
- **Programming and reading through integrated CMOS address decoders**
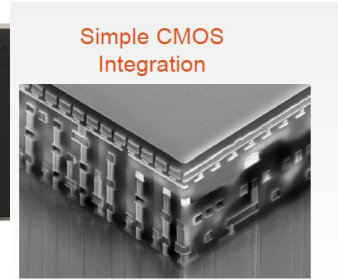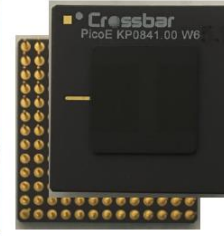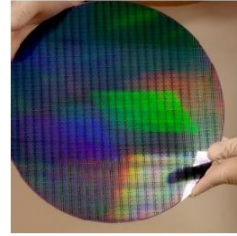- **Each bit written with a single pulse**

Stored/retrieved array 1

Stored/retrieved array 2



**Results from a 40x40 crossbar array integrated on CMOS**

Kim, Gaba, Wheeler, Cruz-Albrecht, Srivinara, W. Lu *Nano Lett.*, 12, 389–395 (2012).

*Lu Group*
*U. Michigan*

- **CMOS** Compatible
- **3D** Stackable, Scalable Architecture – Low thermal budget process
- **Architectures** proven include multiple Via schemes and Subtractive etching
- **Crossbar Inc** founded in 2010, $100M VC funding to date
- **Commercial Products** offered in 2016 based on 40nm CMOS

Simple CMOS Integration



Crossbar

*Lu Group*
*U. Michigan*

---

| Embedded memory with **1T1R** | SCM with 1T1R or **1TnR** | Mass storage with **1TnR** | FPGC configurable logic, CPU with **1T1R** |
|---|---|---|---|



- Monolithic logic/memory integration
- Different memory components integrated on the same chip
- Flexibility of speed/density/cost

**Different approaches for improving computing efficiency (depending on the application):**

- **Bring memory as close to logic as possible, still largely based on conventional architecture**

- **Neuromorphic computing in artificial neural networks**

- More bio-inspired, taking advantage of the internal ionic dynamics at different time scales

- Other compute applications based on vector-matrix multiplications

Towards a general in-memory computing fabric based on a common physical substrate

*Lu Group*
*U. Michigan*

*Lu Group*
*U. Michigan*

## Synapse – reconfigurable two-terminal resistive switches

Co-located memory-compute

High parallelism

post-neuron
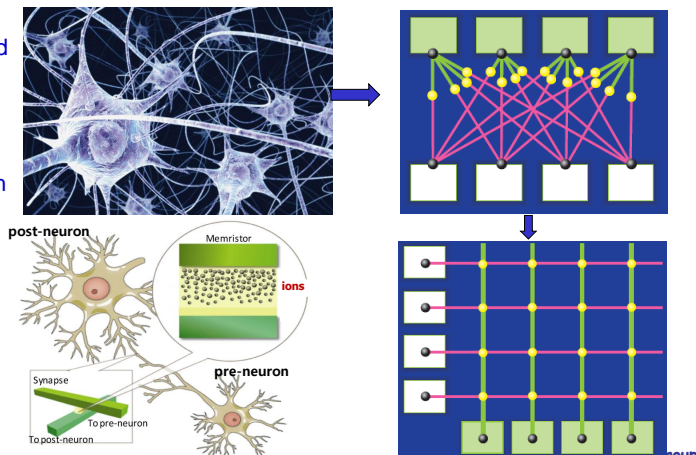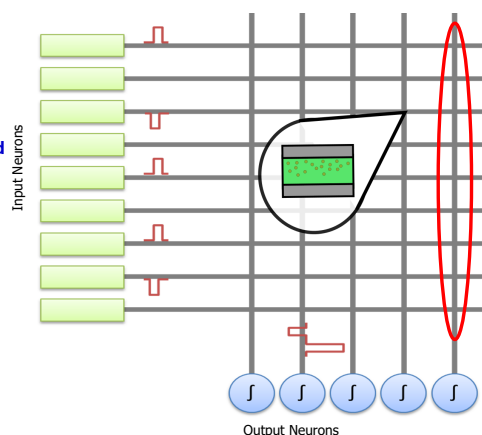
Memristor

ions

Synapse

To post-neuron

To pre-neuron

pre-neuron

S. H. Jo, T. Chang, I. Ebong, B. Bhavitavya, P. Mazumder, W. Lu, *Nano Lett.* **10**, 1297 (2010).

Lu Group
U. Michigan

**RRAM perform learning and inference functions**

**DARPA UPSIDE program**

- **RRAM weights form dictionary elements (features)**
- **Image input, Pixel intensity represented by widths of pulses**
- **Memristor array natively performs matrix operation**

$$\vec{I} = \vec{v} \cdot \vec{\vec{\Phi}}$$

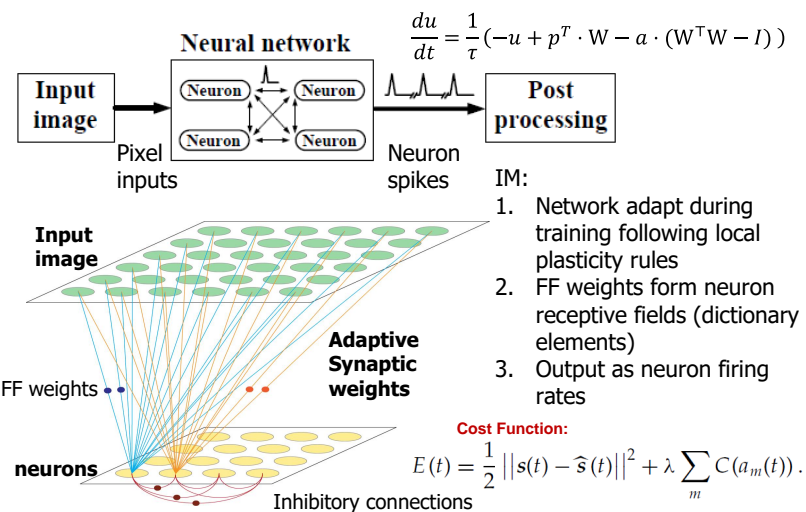- **Integrate and fire neurons**
- **Learning achieved by backpropagating spikes**

Input Neurons

Output Neurons

Lu Group
U. Michigan

---

# Neural Network for Image Processing based on Sparse Coding

**Neural network**

Input image → Neuron, Neuron, Neuron, Neuron → Post processing

Pixel inputs

Neuron spikes

$$\frac{du}{dt} = \frac{1}{\tau}(-u + p^T \cdot W - a \cdot (W^\mathsf{T}W - I))$$

Input image

Adaptive Synaptic weights

FF weights

neurons

Inhibitory connections

IM:
1. Network adapt during training following local plasticity rules
2. FF weights form neuron receptive fields (dictionary elements)
3. Output as neuron firing rates

**Cost Function:**

$$E(t) = \frac{1}{2} ||s(t) - \widehat{s}(t)||^2 + \lambda \sum_m C(a_m(t)).$$

---

# Sparse Coding Implementation in RRAM Array

## Forward Pass
Update neurons/activities

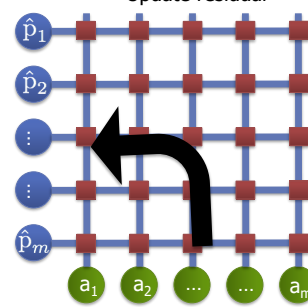$p_1$ $p_2$ ⋮ ⋮ $p_m$

$y_1$ $y_2$ ... ... $y_m$

$$y = p^\mathsf{T} W$$

$$\frac{du}{dt} = \frac{1}{\tau}(-u + p^T \cdot W - a \cdot (W^\mathsf{T}W - I))$$

$$\frac{du}{dt} = \frac{1}{\tau}(-u + (p - \hat{p})^\mathsf{T}W + a)$$
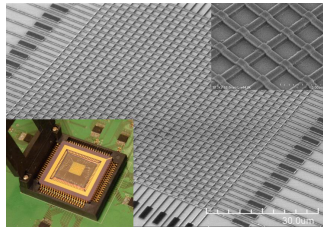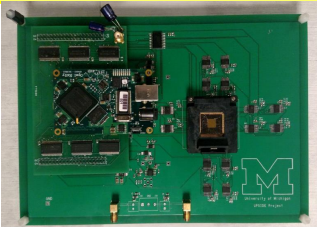
## Backward pass
Update residual

$\hat{p}_1$ $\hat{p}_2$ ⋮ ⋮ $\hat{p}_m$
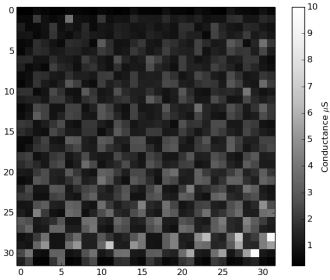
$a_1$ $a_2$ ... ... $a_m$

$$\hat{p} = aW^\mathsf{T}$$

Neuron membrane potential

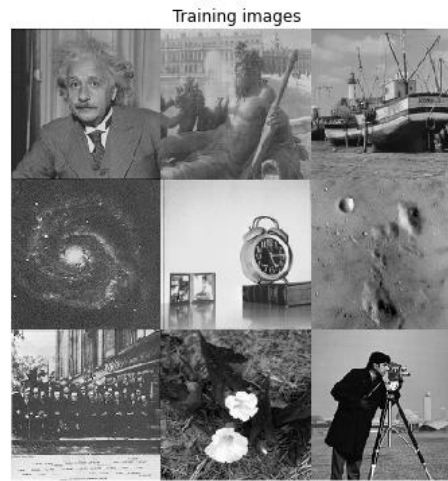Sheridan et al., *Nature Nanotechnology*, **12**, 784–789 (2017)

16

**32x32 memristor array**

- **Checkerboard pattern**
- **32 x 32 array**
- **Direct storage and read out**
- **No read-verify or re-programming**

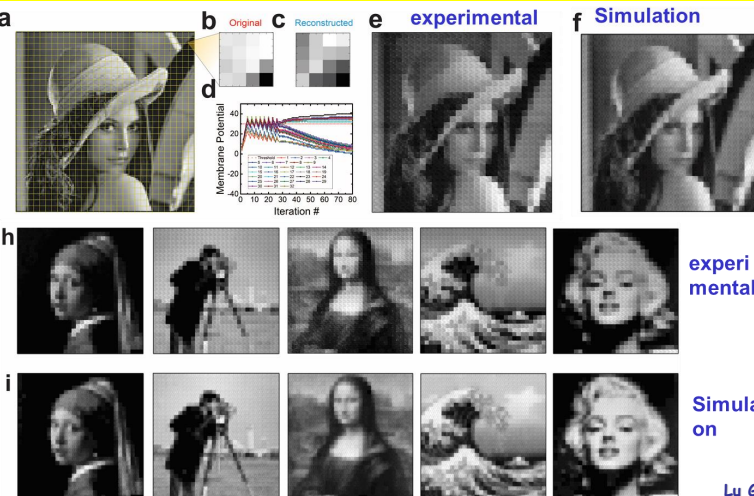Sheridan et al., Nature Nanotechnology 12, 784–789 (2017)

*Lu Group*
*U. Michigan*

Training images

9 Training Images
128x128px
4x4 patches
Trained in random order

Sheridan et al., Nature Nanotechnology 12, 784–789 (2017)

18

*Lu Group*
*U. Michigan*

---

# Image Reconstruction with RRAM Crossbar

# Improving Computing Efficiency using RRAM Arrays

a   b Original   c Reconstructed   e **experimental**   f **Simulation**
d Membrane Potential / Iteration #

h   **experimental**

i   **Simulation**

Sheridan et al., Nature Nanotechnology 12, 784–789 (2017)

*Lu Group*
*U. Michigan*

**Different approaches for improving computing efficiency (depending on the application):**

- Bring memory as close to logic as possible, still largely based on conventional architecture

- Neuromorphic computing in artificial neural networks

- More bio-inspired, taking advantage of the internal ionic dynamics at different time scales

- **Other compute applications based on vector-matrix multiplications**

Towards a general in-memory computing fabric based on a common physical substrate

*Lu Group*
*U. Michigan*

**Solving partial-differential equations (PDEs)**



a    b    c

$I_j = \sum V_i \cdot G_{i,j}$

**Solving an A·x=b problem in matrix form**
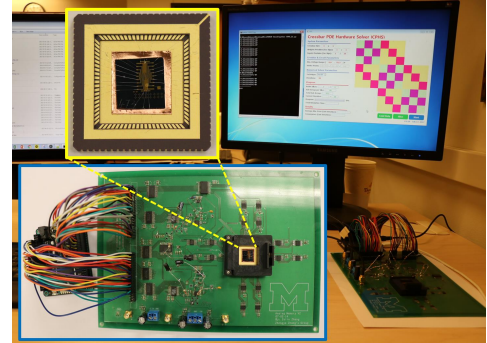
- A second order Poisson equation as a toy example,

$$\nabla^2 u = -2 \cdot \sin(x) \cdot \cos(y)$$

- The problem is solved using finite difference (FD), where matrix can be sliced into a set of few similar slices.



*Lu Group*
*U. Michigan*

## » Hardware Test bench

- The test board consists of: **(i)** RRAM crossbar, **(ii)** DACs to control the input signals, **(iii)** sense amplifiers and ADCs to sample the output current, **(iv)** MUXs to route the signals, and **(v)** FPGA to enables the software interface and control.



M. A. Zidan, Y.J. Jeong, J. Lee, B. Chen, S. Huang, M. J. Kushner, & W. D. Lu, Nature Electronics, 1, 411–420 (2018)
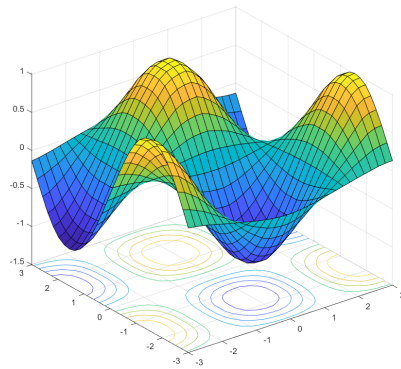
*Lu Group*
*U. Michigan*

---

# Hardware Prototyping

## » Solving a toy example

**Measured Results for the toy Example**



M. A. Zidan, Y.J. Jeong, J. Lee, B. Chen, S. Huang, M. J. Kushner, & W. D. Lu, Nature Electronics, **1,** 411–420 (2018)
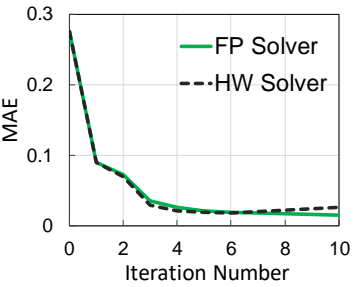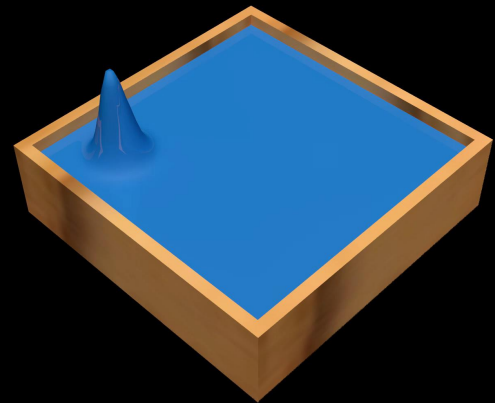
*Lu Group*
*U. Michigan*

## » Results Reconstructed as a 3D Animation
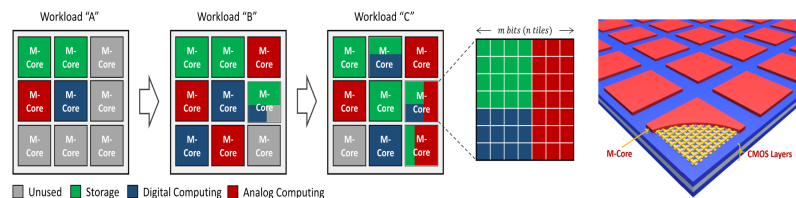


M. A. Zidan, Y.J. Jeong, J. Lee, B. Chen, S. Huang, M. J. Kushner, & W. D. Lu, Nature Electronics, , 1, 411–420 (2018)

Workload "A"  Workload "B"  Workload "C"  ← m bits (n tiles) →

☐ Unused  ☐ Storage  ☐ Digital Computing  ☐ Analog Computing

M-Core  CMOS Layers
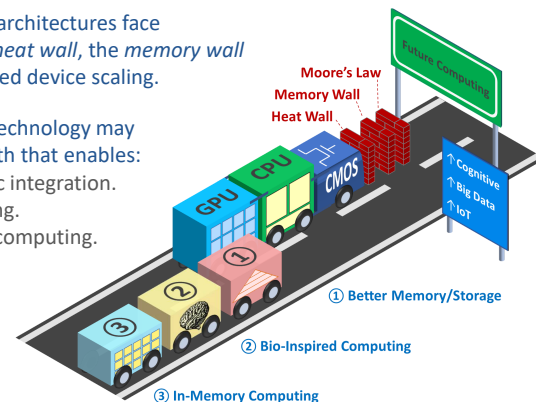
- **Memory-Computing Unit (MPU)**
- "General" purpose by design: the same hardware supports different tasks – low precision or high precision. Not just an neuromorphic accelerator
- Dense local connection, sparse global connection
- Run-time, dynamically reconfigurable. Function defined by software.

**M. Zidan, Y. Jeong, J. H. Shin, C. Du, Z. Zhang, and W. D. Lu, IEEE Trans Multi-Scale Comp Sys, DOI 10.1109/TMSCS.2017.2721160 (2017)**

**M. A. Zidan, J. P. Strachan, and W. D. Lu, Nature Electronics 1: 22–29 (2018)**

*Lu Group*
*U. Michigan*

- Conventional computing architectures face challenges including the *heat wall*, the *memory wall* and difficulties in continued device scaling.

- Developments in RRAM technology may provide an alternative path that enables:
  - Hybrid memory–logic integration.
  - Bioinspired computing.
  - Efficient in-memory computing.



Moore's Law
Memory Wall
Heat Wall

Future Computing
↑ Cognitive
↑ Big Data
↑ IoT

① Better Memory/Storage
② Bio-Inspired Computing
③ In-Memory Computing

**M. A. Zidan, J. P. Strachan, and W. D. Lu, Nature Electronics 1: 22–29 (2018)**

**M. Zidan, Y. Jeong, J. H. Shin, C. Du, Z. Zhang, and W. D. Lu, IEEE Trans Multi-Scale Comp Sys, DOI 10.1109/TMSCS.2017.2721160 (2017)**

*Lu Group*
*U. Michigan*

---

# Summary

**Different approaches for improving computing efficiency (depending on the application):**

- **Bring memory as close to logic as possible, still largely based on conventional architecture**

- **Neuromorphic computing in artificial neural networks**

- **More bio-inspired, taking advantage of the internal ionic dynamics at different time scales**

- **Other tasks based on vector-matrix multiplications**

⟹ **Towards a general in-memory computing fabric based on a common physical substrate**

*Lu Group*
*U. Michigan*