



**Hewlett Packard
Enterprise**

Inspired by the Brain: Computing from Architecture to Devices

R. Stanley Williams
HPE Senior Fellow

**Neuromorphic Computing Symposium
July 18, 2017**



THE END OF MOORE'S LAW

IEEE

AIP
aip.org

What's Next?

The end of Moore's law could be the best thing that has happened in computing since the beginning of Moore's law.

The Chua Lectures: A 12-Part Series at HPE Labs

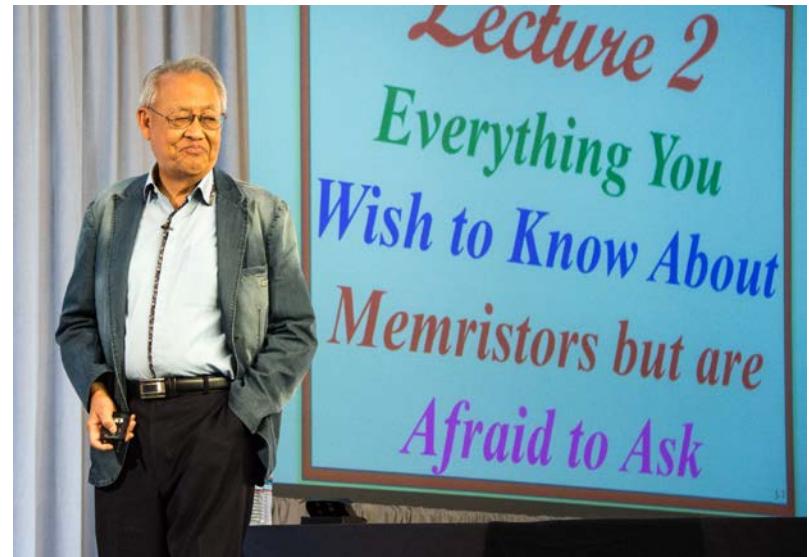
From Memristors and Cellular Nonlinear Networks to the Edge of Chaos

<https://www.youtube.com/playlist?list=PLtS6YX0YOX4eAQ6lrOZSta3xjRXzpcXyi>

or enter “The Chua Lectures” into your favorite browser

What's missing?

‘Linearize then analyze’ is not valid for understanding nanodevices or neurons – a new nonlinear dynamical theory of ‘electronic’ circuits is needed, and was developed 50 years ago by Leon Chua, father of nonlinear circuit theory and Cellular NNs.



Structure of a US Neuromorphic Science Computing Program

1. Connect Theory of Computation with Neuroscience and Nonlinear Dynamics
 - e.g. Boolean logic, CNN, Bayesian Inference, Energy-Based Models, Markov Chain
2. Architecture of the Brain and Relation to Computing and Learning
 - Theories of Mind: Albus, Eliasmith, Grossberg, Mead, many others
3. Simulation of Computational Models and Systems – need supercomputers!
4. System Software, Algorithms & Apps – Make it Programmable/Adaptable
5. Chip Design – System-on-Chip: Accelerators, Learning and Controllers
 - Compatible with standard processors, memory and network fabric (Gen-Z)
6. Chip Processing and Integration – Full Service Back End of Line on CMOS
 - DoE Nanoscale Science Research Centers (NSRCs) – e.g. CINT
7. Devices and Materials – *in situ* and *in operando* test and measurement
 - Most likely materials will be adopted from Non-Volatile Memory

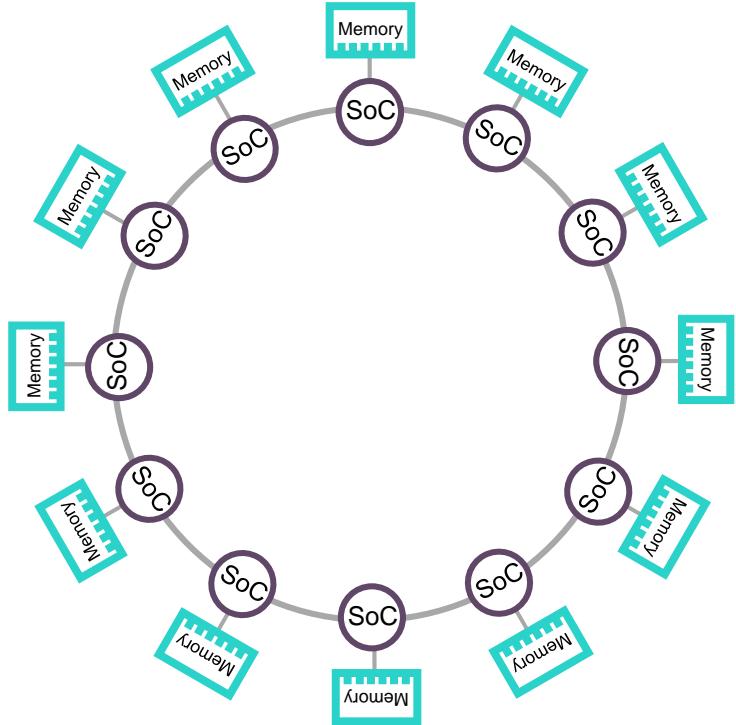
Future exponential increases in computer performance and efficiency will require multiple advances:

- Memory-centric computing – no von Neumann bottleneck
- Gen-Z – high performance open fabric to democratize computing
(<http://genzconsortium.org/>)
- Dot-product engine: memristor-based vector-matrix multiplication accelerator for neural nets and signal processing
- Mimicking Synapse and Neuron Dynamics with Memristors
- Chaos as a computing resource for constrained optimization problem solving (Hopfield network)

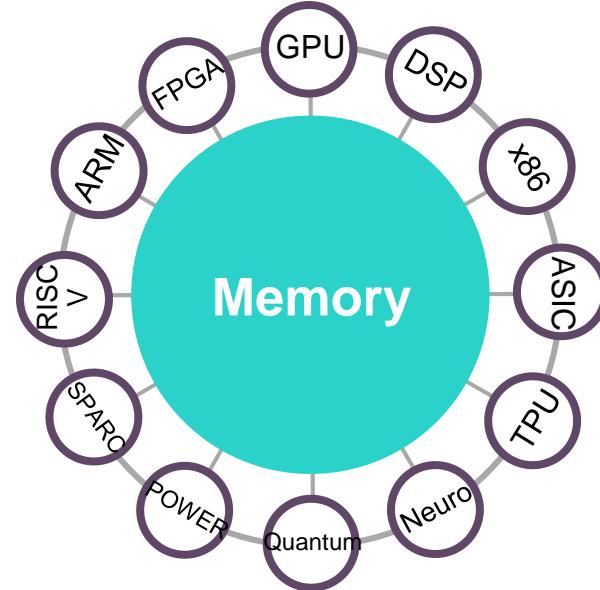
Systems and Architectures



Devices



From processor-centric computing...
the traditional von Neumann architecture



...to Memory-Driven Computing
**with Gen-Z open (nonproprietary) fabric,
 computing is plug and play**

Memory-Driven Computing (MDC) is a reality: The Machine

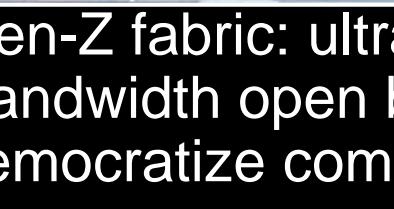
Fast, persistent memory

Combining memory and storage in a stable environment to increase processing speed and improve energy efficiency



Fast memory fabric

Using photonics where necessary to eliminate distance and create otherwise impossible topologies



Task-specific processing

Optimizing processing from general to specific tasks



New and Adapted software

Radically simplifying programming and enabling new applications that we can't even begin to build today



Memristor technology,
ongoing transfer from lab to
commercial product



Developing new
accelerators from novel
device behavior



High Bandwidth Low Latency

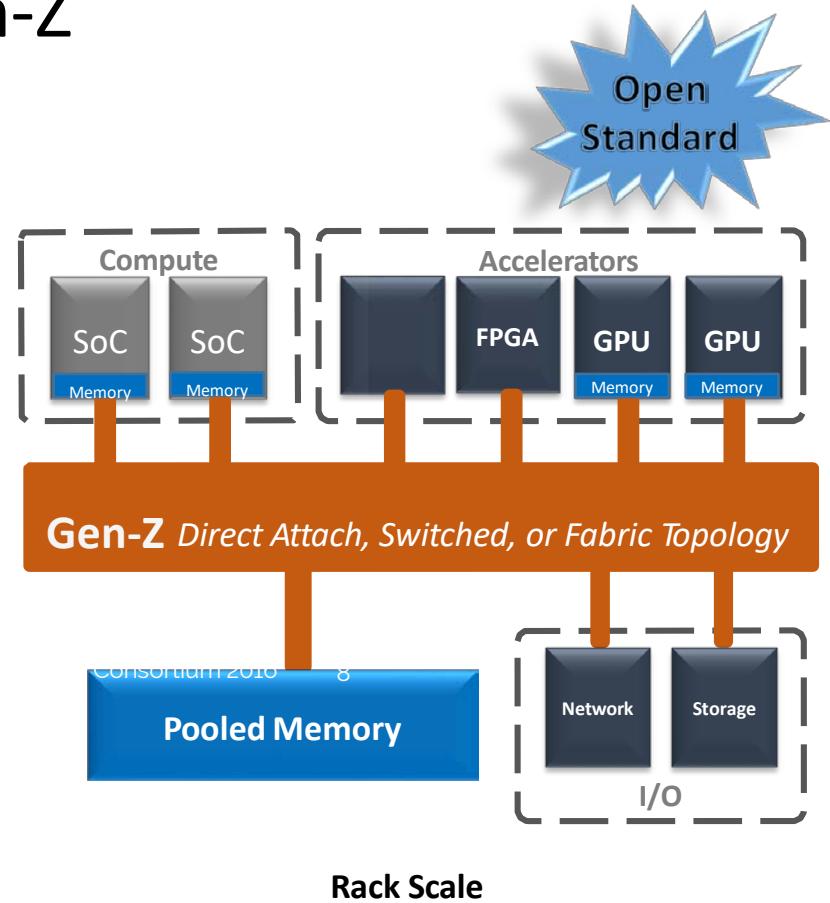
- Memory Semantics – simple Reads and Writes
- From tens to several hundred GB/s of bandwidth
- Sub-100 ns load-to-use memory latency

Advanced Workloads & Technologies

- Real time analytics
- Enables data centric and hybrid computing
- Scalable memory pools for in memory applications
- Abstracts media interface from SoC to unlock new media innovation

Secure Compatible Economical

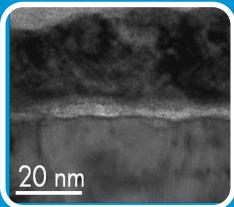
- Provides end-to-end secure connectivity from node level to rack scale
- Supports unmodified OS for SW compatibility
- Graduated implementation from simple, low cost to highly capable and robust
- Leverages high-volume IEEE physical layers and broad, deep industry ecosystem



Founding Members of the Gen-Z Consortium (38 today)

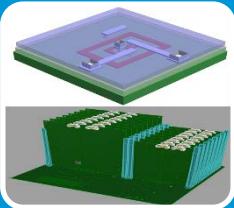


Brain Inspired Accelerators: devices → circuits → applications



Materials engineering – device properties

- High resistance
- Large OFF/ON ratio
- Low switching current
- Linear electronic transport

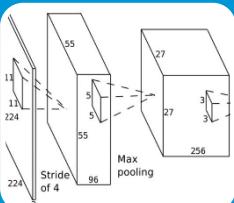


Cell integration, arrays, and circuits

- Integration/fab of memristors with selectors/transistors
- Optimizing reading and writing circuits
- Construct platform to write/read device arrays
- Compact modeling of memristor characteristics

Architectures, Algorithms, Applications

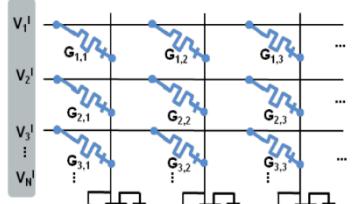
- Matching capability to favorable applications
- Benchmarking and optimizing performance; GOPS/Watts
- Quantifying sources of error and bottlenecks



Application-aware research
into materials/cell/circuit

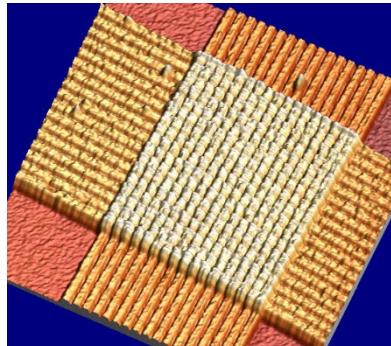
Dot Product Engine: memristor arrays accelerate vector-matrix multiplication (C. elegans – 800 papers/yr)

Input
Voltage
vector



Output
current

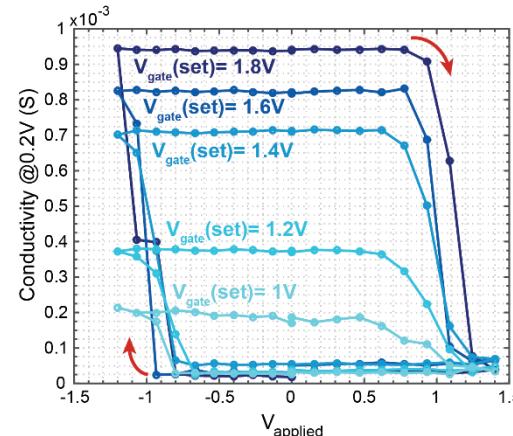
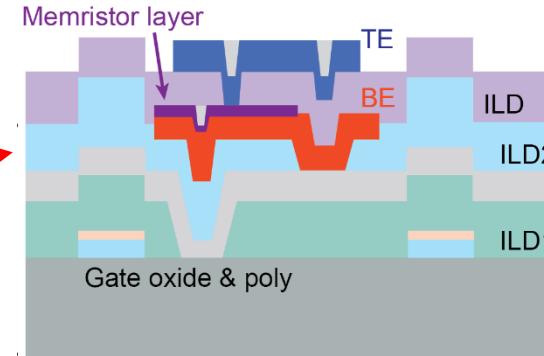
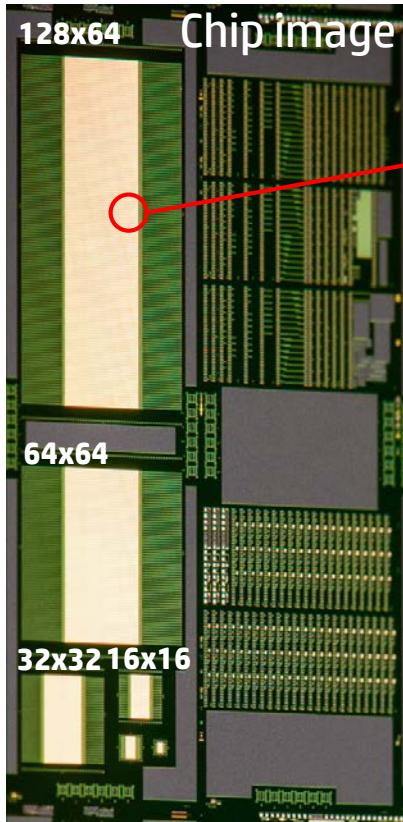
$$I_j^0 = \sum_i G_{ij} \cdot V_i^I$$



- Parallel multiply & add through Kirchoff's and Ohm's laws
1961, K. Steinbuch "Die Lernmatrix" – suggests using "ferromagnetic toroids"
- Memristors as highly scalable, tunable analog resistors
High ON/OFF ratio (~10⁵), supporting multiple levels
- Well suited for streaming workloads like neural nets
- Many ways to scale up
Memristor levels, array size, wire pitch, 3D layer, DAC/ADC speed & width etc.
- Performance (execution time) improvements >1000x and energy efficiency >100x over GPUs for particular applications

The DPE: memristor-based analog computing platform

Flexible system for programming and computing on arrays of integrated Transistor-Memristors (1T1M)



Programming memristor arrays for computation

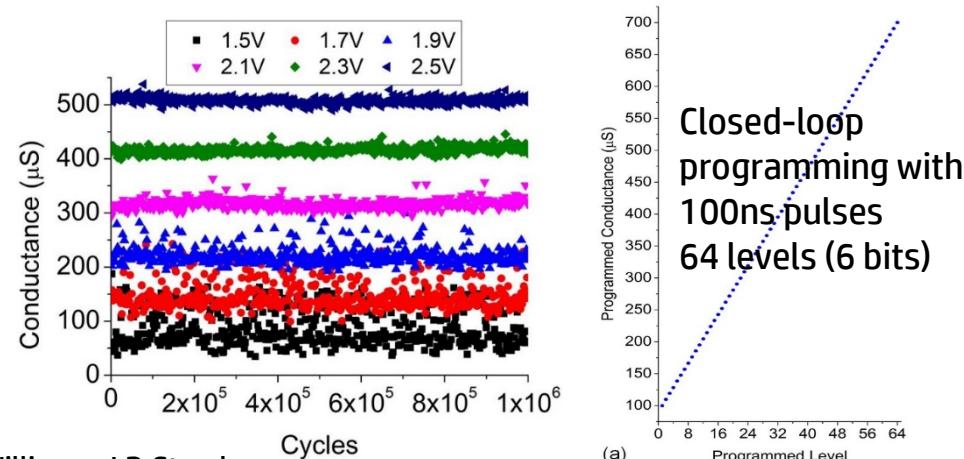


64x128 array = 8192 memristors

Each “pixel” is a continuously tunable memristor cell

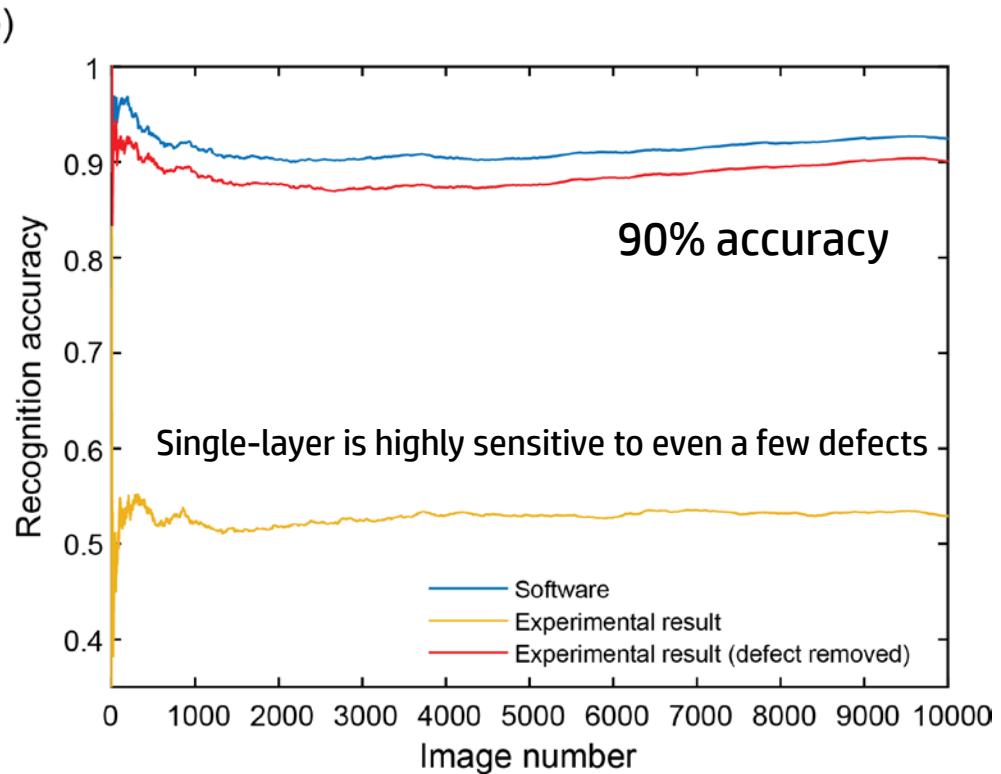
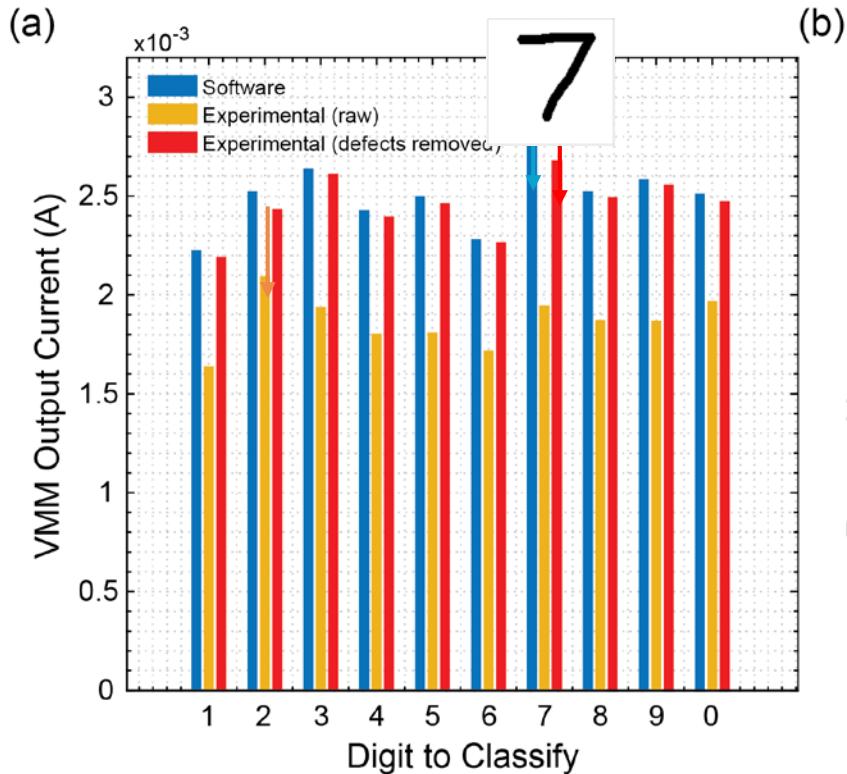
Grayscale - 182 distinct conductance levels (~7-8 bits)

Each memristor can be reprogrammed >1e6 times



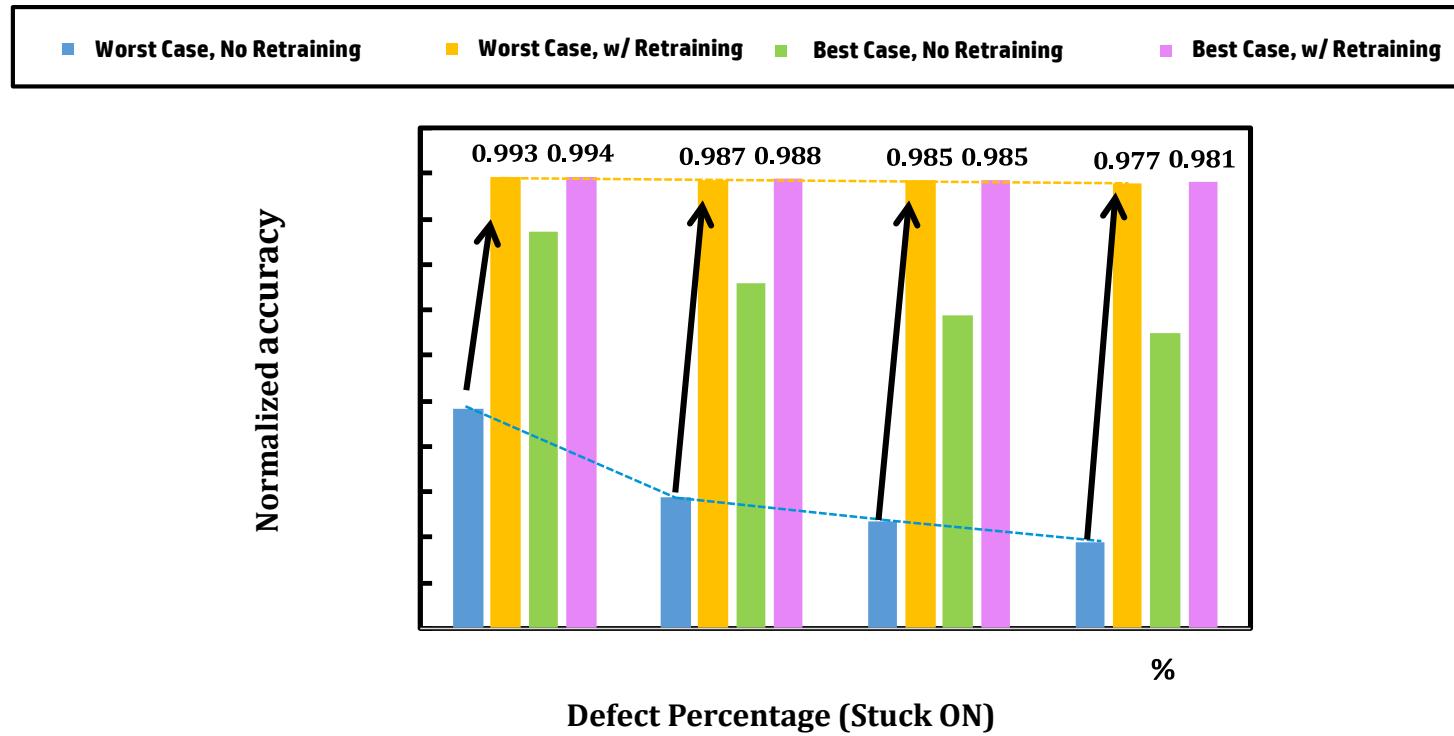
E.J. Merced-Grafals, N. Dávila, N. Ge, R.S. Williams, J.P. Strachan
Nanotechnology 27, 365202 (2016).

Single Layer NN for MNIST classification

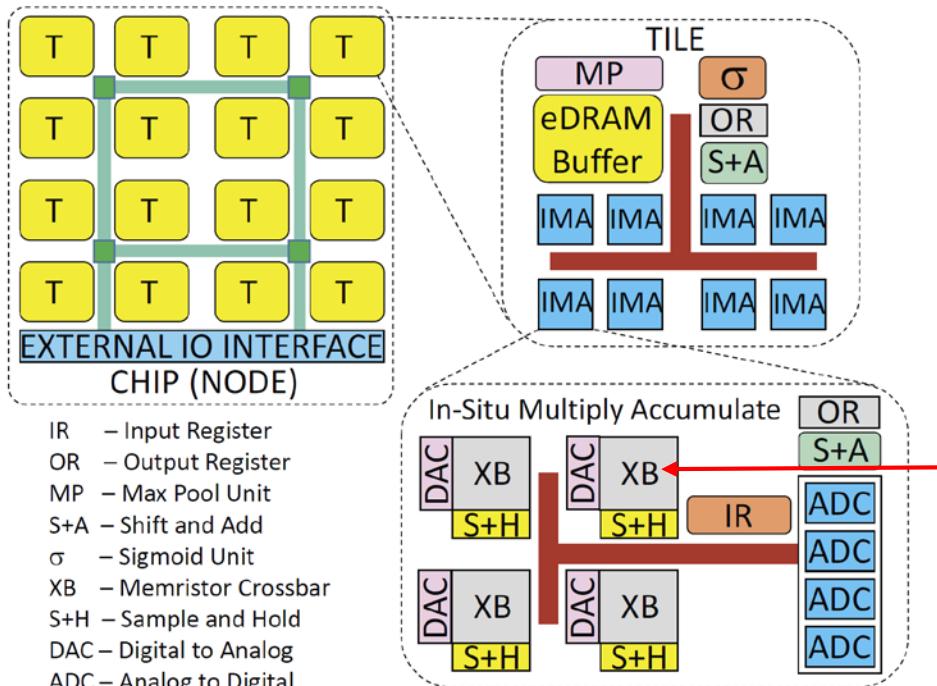


Dealing with “defects” - retraining

Retrain the network - surrounding weights can compensate for bad devices



“ISAAC” - architecture acceleration using DPEs



Store and re-use Kernels in non-volatile memristors – reduce data fetching

Heavy pipelining

Speedup of >5,000x over GPUs and 800x less energy

Speedup of 15x over Digital ASIC and 5.5x less energy

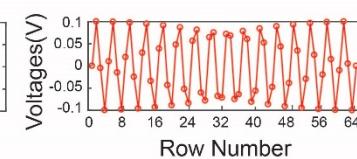
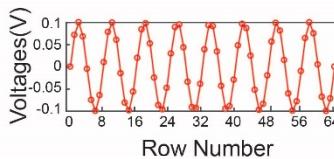
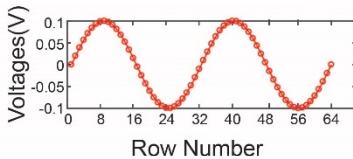
128x128 memristor Xbar arrays, 2 bits per cell, operating at 10 MHz

A. Shafiee, et al., “ISAAC: A Convolutional Neural Network Accelerator with In-Situ Analog Arithmetic in Crossbars”, *International Symposium on Computer Architecture (ISCA) 2016*.

Signal processing on a DPE system

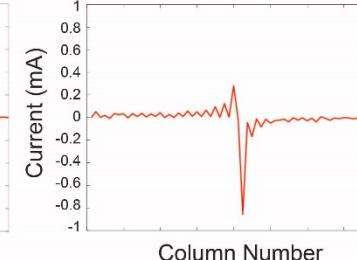
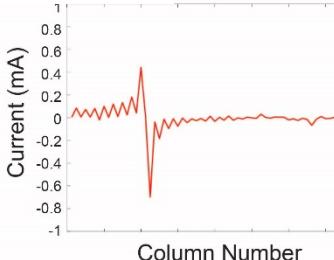
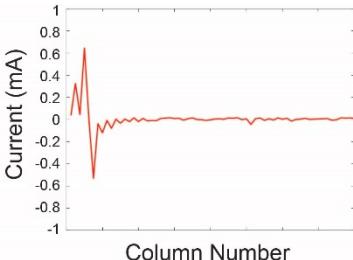
(UMass collaboration) – instantaneous cosine transform

Time-domain inputs (applied to rows)

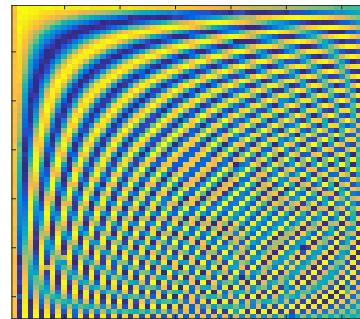
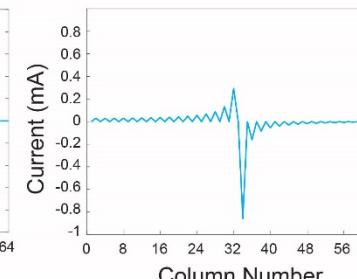
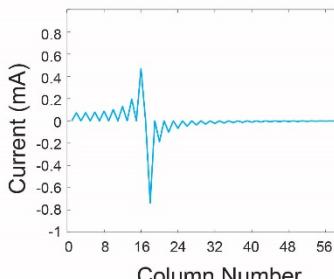
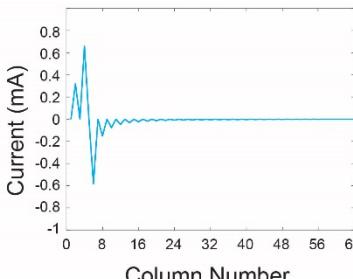


Freq-domain (DCT) outputs (read from columns)

Real-time
experimental
data from DPE

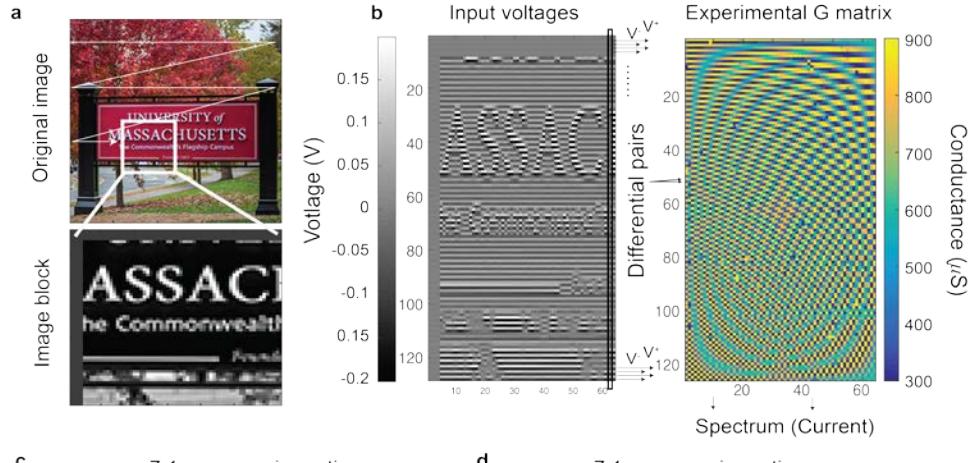


Software
cosine
transform



C. Li, et al.,
submitted
Umass Amherst

Image Compression with cosine transform on the DPE



Software encoder

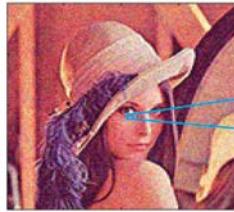


Experimental crossbar encoder

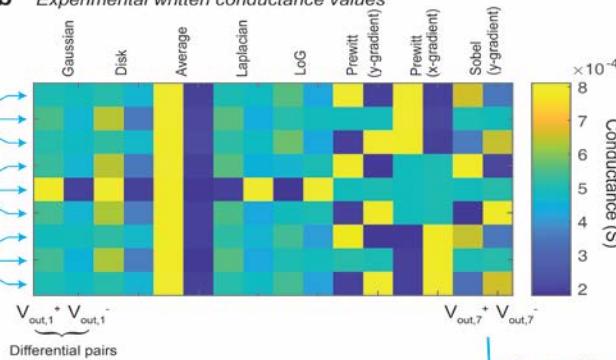
C. Li, et al., *submitted*
UMass Amherst

Convolutional Filters with the DPE

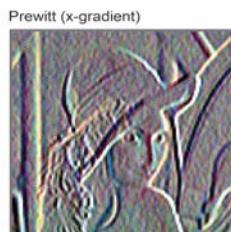
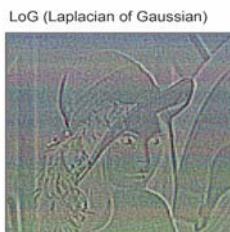
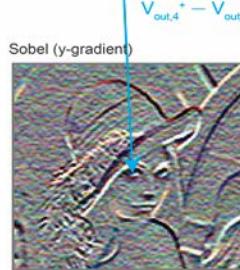
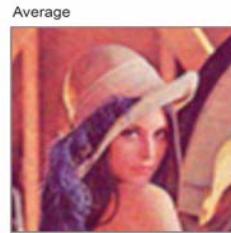
a Original Image with noises



b Experimental written conductance values



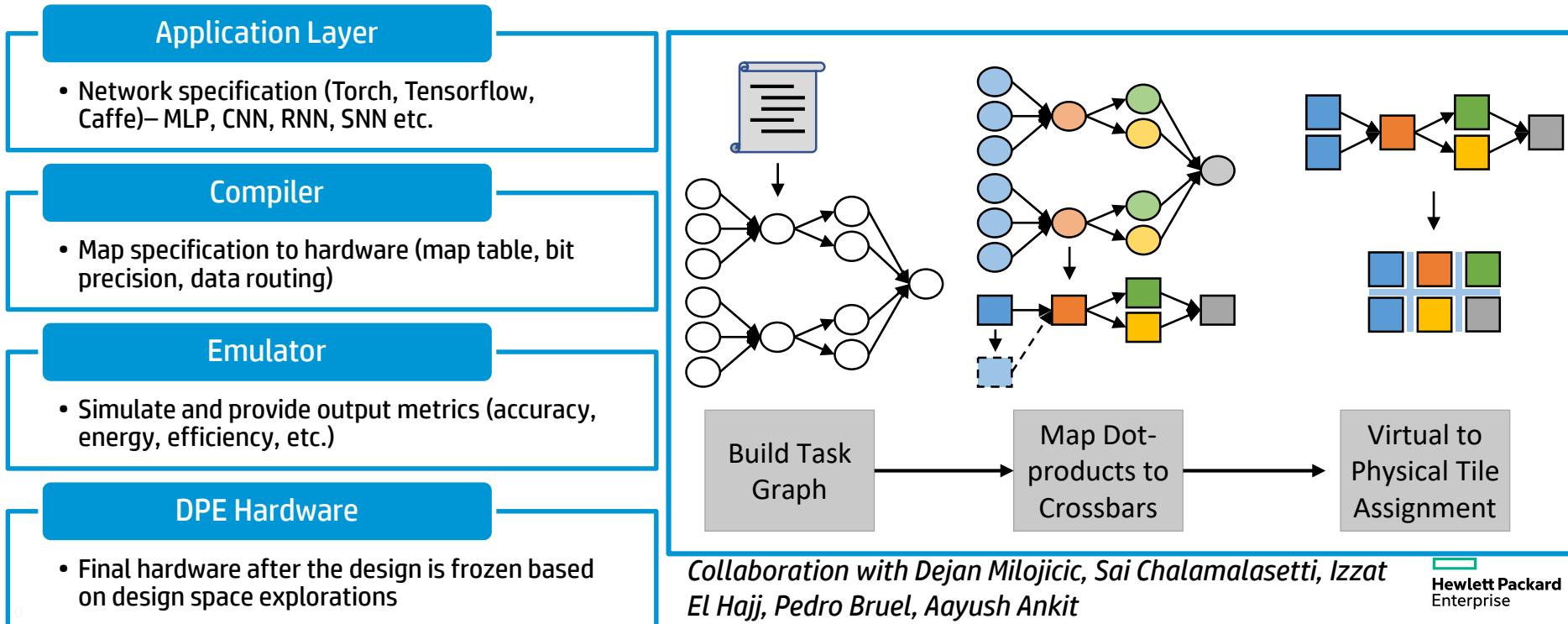
c Experimental crossbar outputs:



C. Li, et al., submitted
Umass Amherst

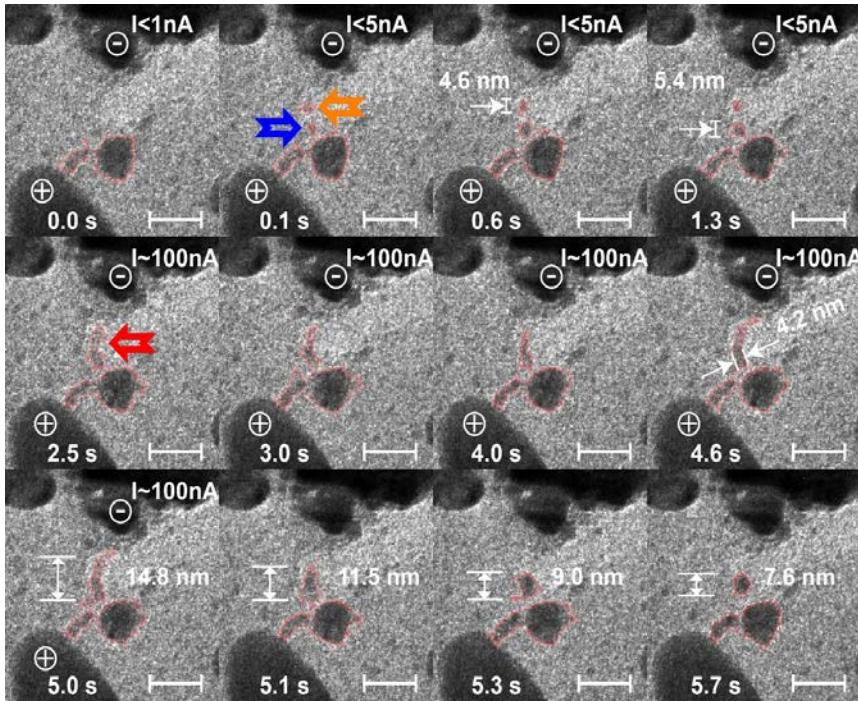
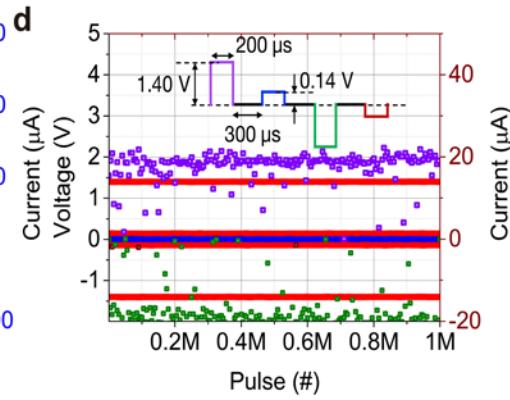
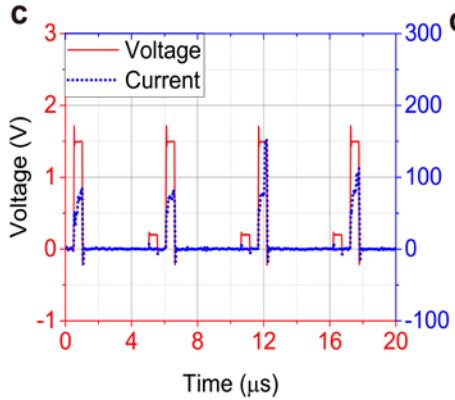
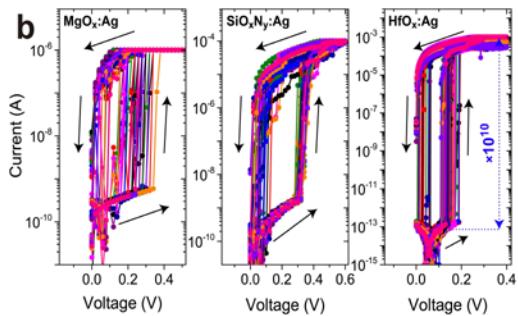
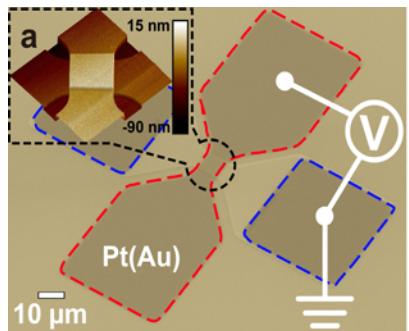
Writing software for a memristor accelerator

- Initial development of OS/Compiler/HW emulator
- Co-designed with hardware efforts
- Generalize/Broaden to other matrix-heavy applications (medical imaging, sci computing, network management)

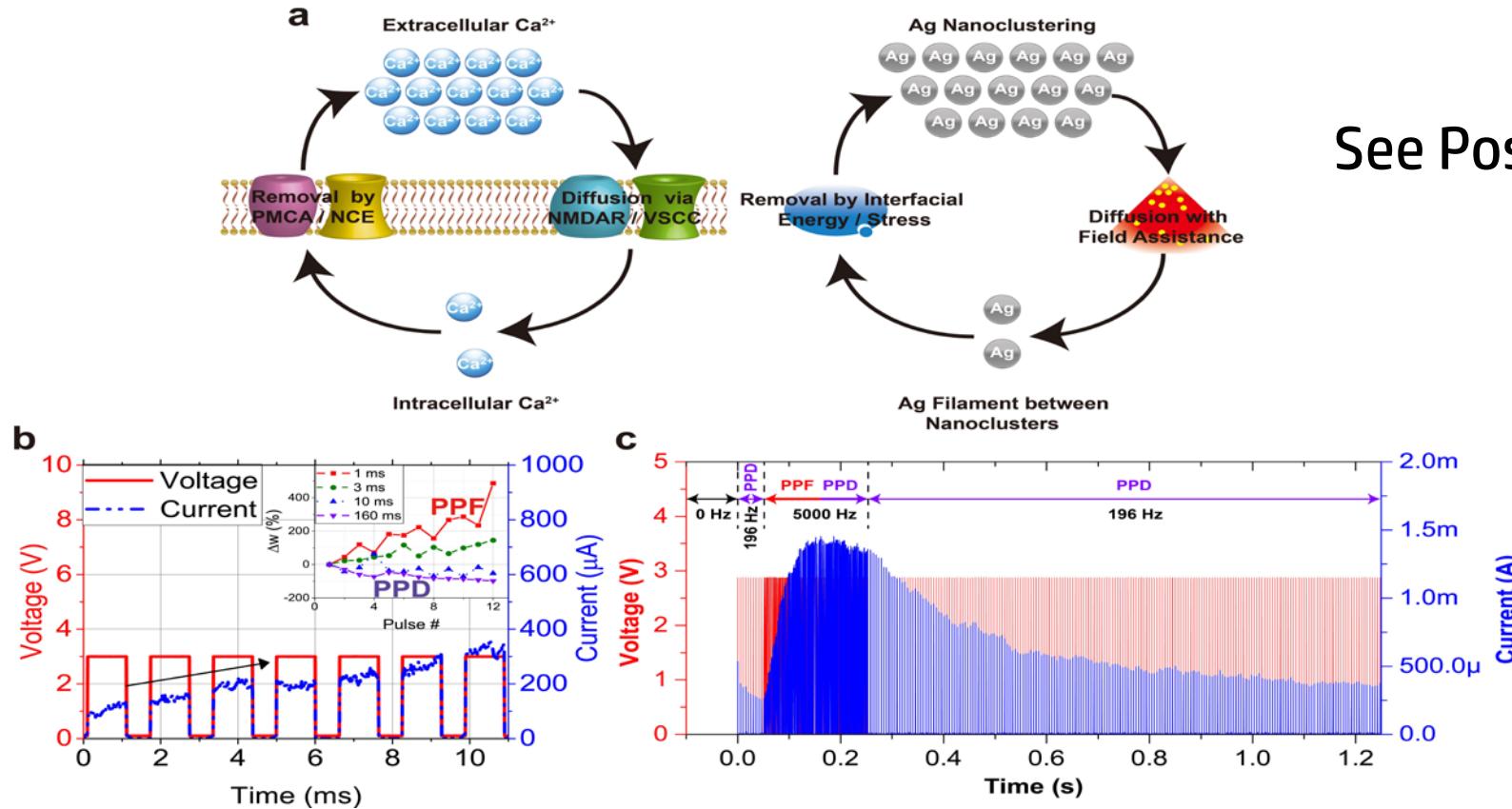


Collaboration with Dejan Milojicic, Sai Chalamalasetti, Izzat El Hajj, Pedro Bruel, Aayush Ankit

Ag in Metal Oxide Diffusive Memristor – Emulating Dynamics

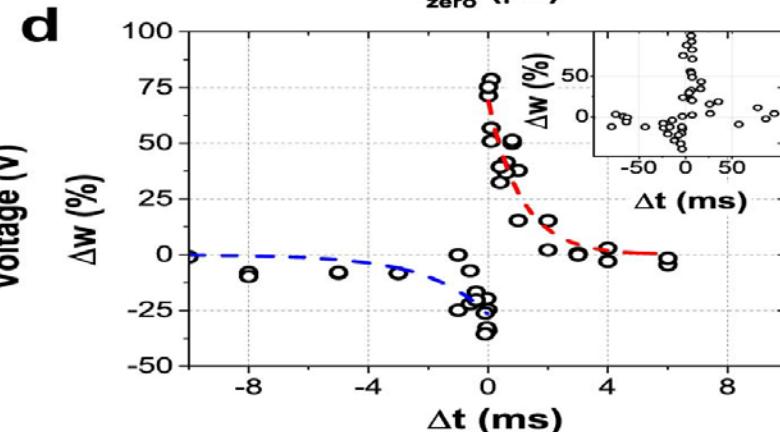
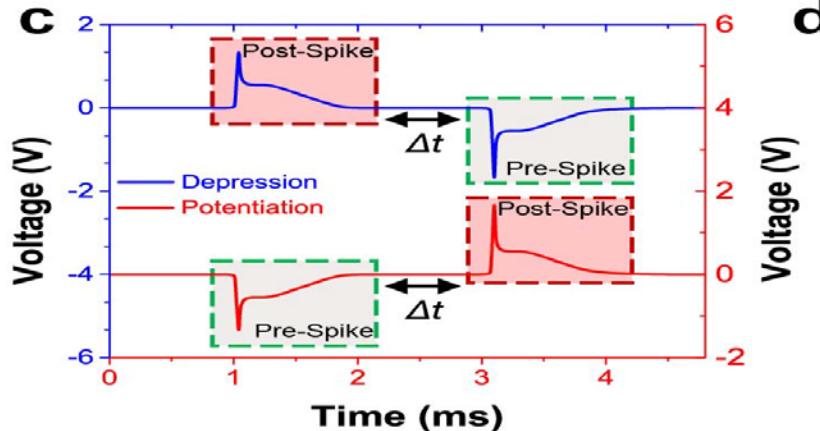
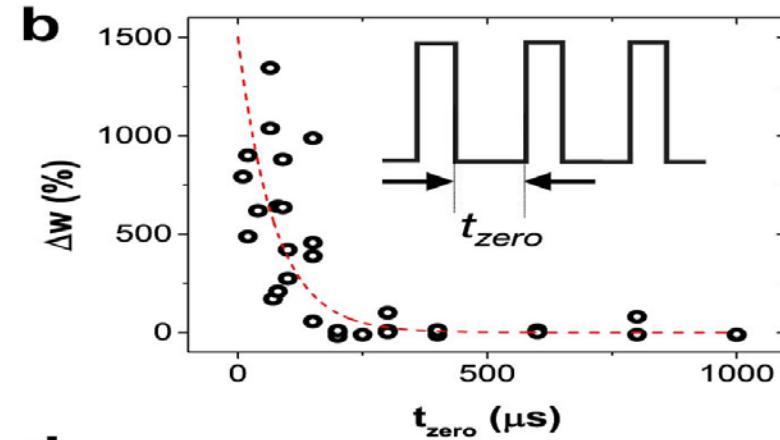
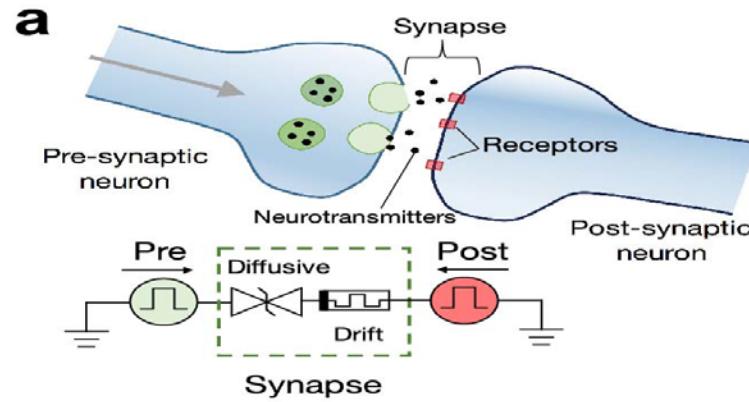


Ag Diffusive Memristor Mimics Ca^{2+} Ion Dynamics in Synaptic Gap



See Poster!

Diffusive Decay Provides an Extended Window for STDP



Modified Hopfield network for optimization problems

Turing believed that intelligence required randomness and would necessarily lead to mistakes

Chua has shown that neurons are ‘poised on the edge of chaos’, and that chaos leads to complexity and emergent behavior

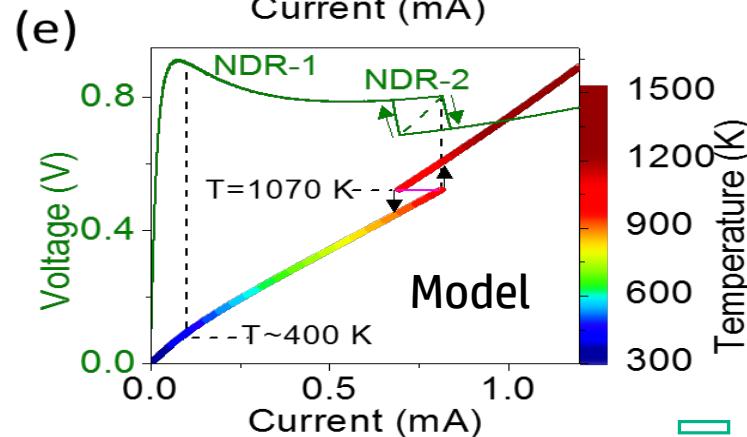
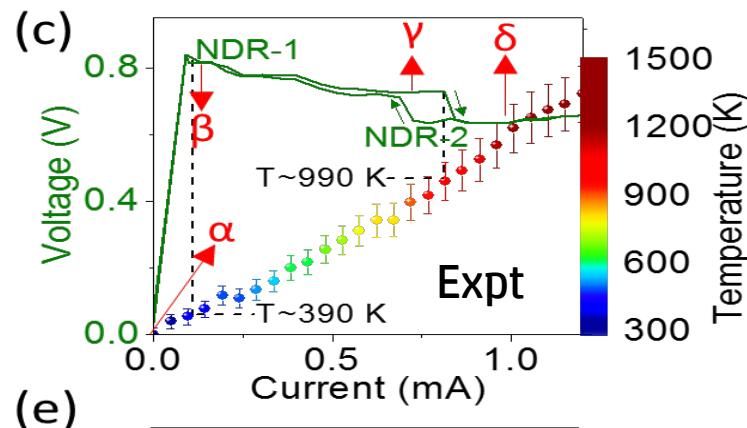
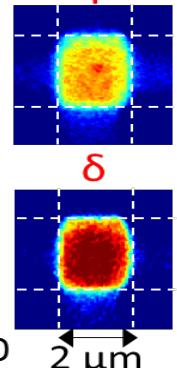
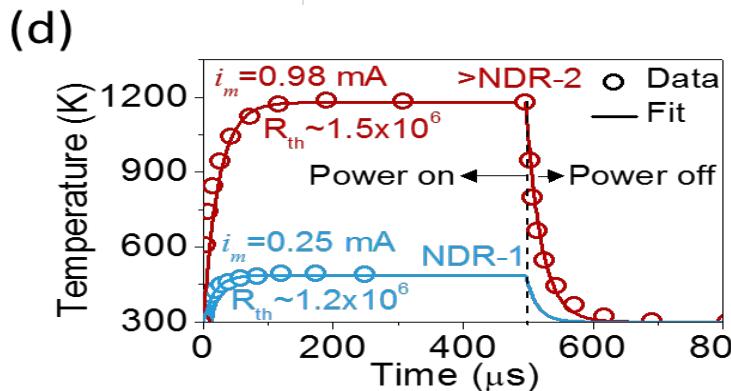
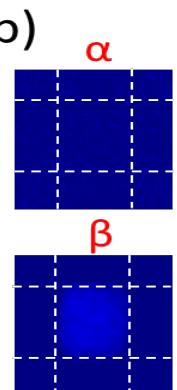
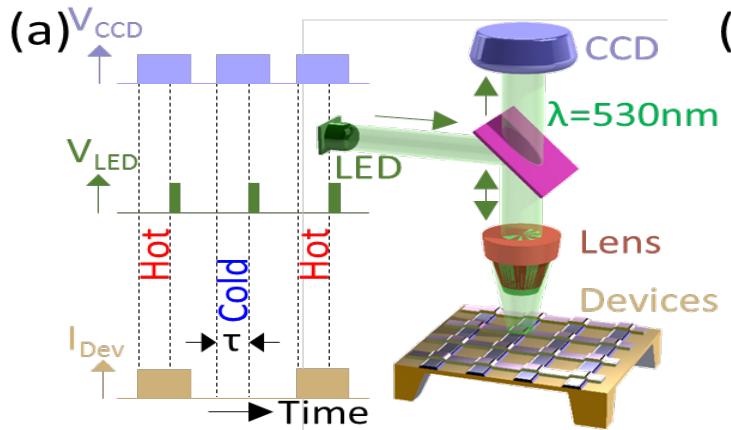
Epilepsy occurs when the neurons in the brain synchronize, i.e. lack of chaos

Chaos as a computing resource

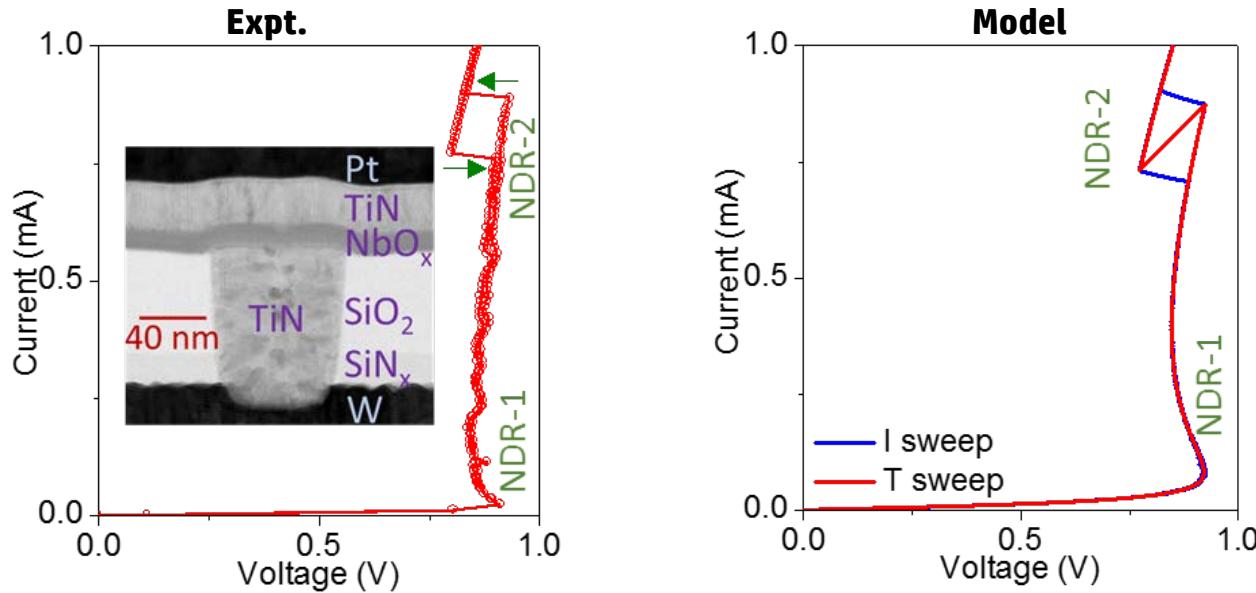
May allow construction of a highly scalable “annealing” machine

Benchmarking to compare to other Classical and Quantum Annealing machines

Thermoreflectance Imaging of NbO_2 current control behavior



Niobium oxide memristor: behavior and model



Model: Modified 3D Poole-Frenkel

$$i_m = \left[2A \times 10^4 e^{-\frac{0.301}{k_B T}} \left(\frac{k_B T}{\beta \sqrt{v_m}} \right)^2 \left\{ 1 + \left(\frac{\beta \sqrt{v_m/t}}{k_B T} - 1 \right) e^{\frac{\beta \sqrt{v_m/t}}{k_B T}} \right\} + \frac{2A \times 10^4 e^{-\frac{0.301}{k_B T}/t}}{2} \right] v_m$$

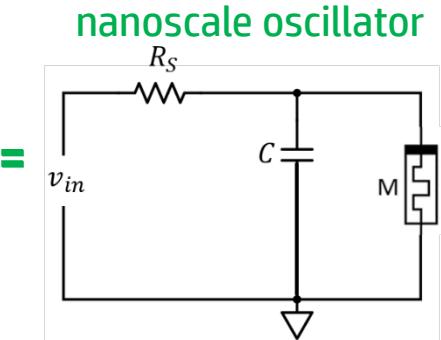
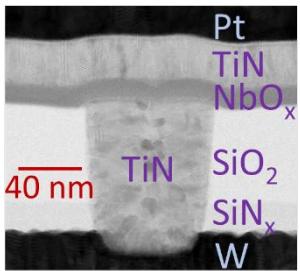
Temperature dynamics

$$\frac{dT}{dt} = \frac{i_m v_m}{C_{th}} - \frac{T - T_{amb}}{C_{th} R_{th}(T)}$$

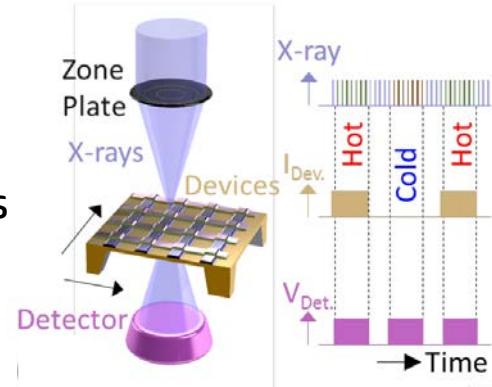
Mott transition in R_{th}

$$R_{th}(T) = 1.4 \times 10^6 \text{ (for } T \leq T_C\text{)} \text{ and } 2 \times 10^6 \text{ (for } T > T_C\text{)}$$

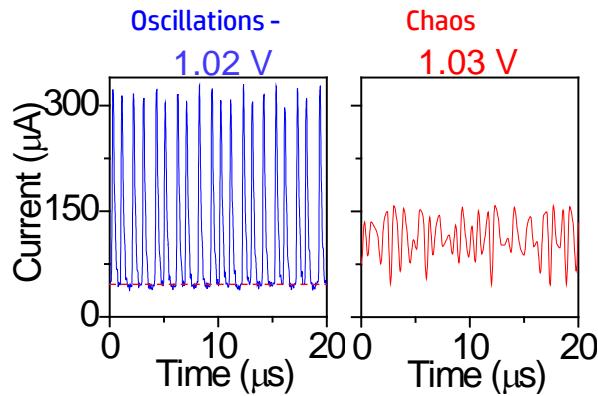
NbO₂ memristor: computing with chaos



in situ and *in operando* scanning transmission x-ray microscopy (STXM) at ALS critical for analysis of memristor operation



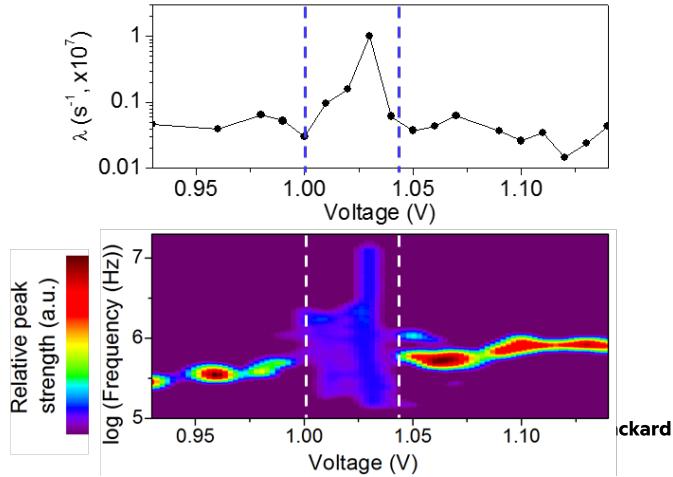
Dynamical Behavior:



Lyapunov exponent

$$\lambda = \lim_{t \rightarrow \infty} \lim_{\delta Z_0 \rightarrow 0} \frac{1}{t} \ln \frac{|\delta Z(t)|}{|\delta Z_0|}$$

Frequency analysis



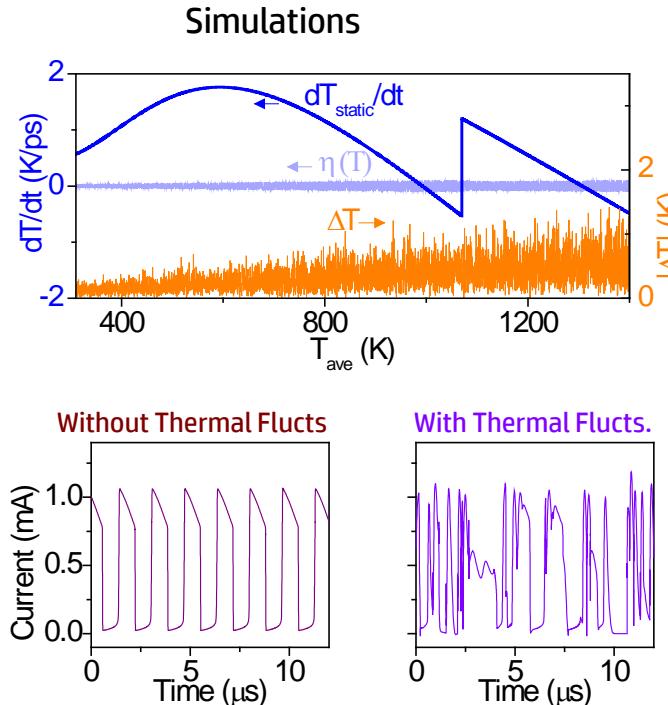
Able to describe by incorporating thermal fluctuations

$$\eta(T) = T \left(\frac{k_B}{C_{th}} \right)^{\frac{1}{2}} \frac{4\pi}{R_{th} C_{th}} \text{rand}(-1 \leftrightarrow 1)$$

Negative Differential Resistance in the system leads to positive feedback in presence of fluctuations

Thermal capacitance C_{th} shrinks with device size

→ increasing importance of thermal fluctuations

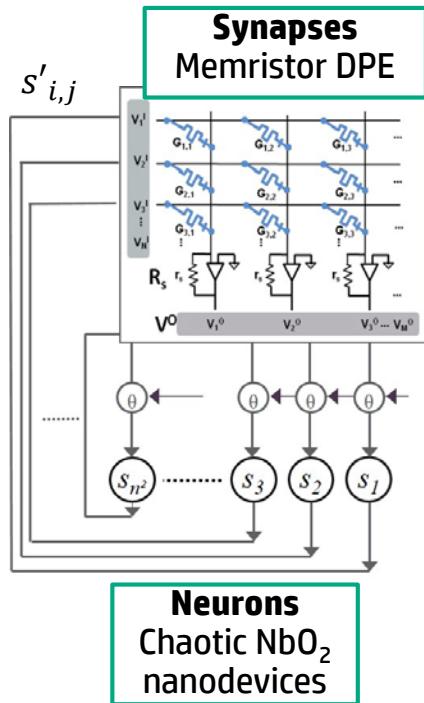


Modified Hopfield network for optimization problems



Traveling Salesman
problem

NP hard: $O(n^3 2^n)$

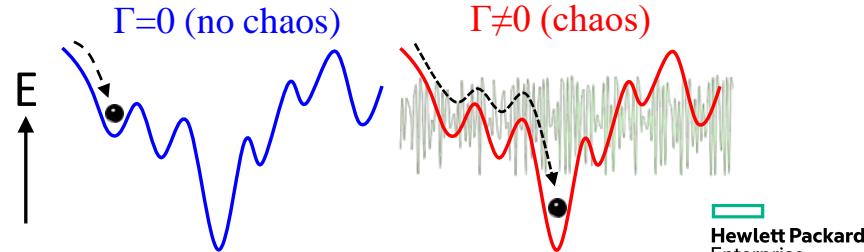


Encode any TSP instance in the DPE weight matrix

Defines an “energy” of the system to be minimized

$$E = -\frac{1}{2} \sum_i \sum_j s_{i,j} \sum_k \sum_l s_{k,l} w_{(i,j),(k,l)} + \sum_i \sum_j s_{i,j} \theta$$

Follow update rule: $s_{i,j} = \begin{cases} 1 & \text{if } Ws'_{i,j} > \theta \\ -1 & \text{if } Ws'_{i,j} < \theta \end{cases}$



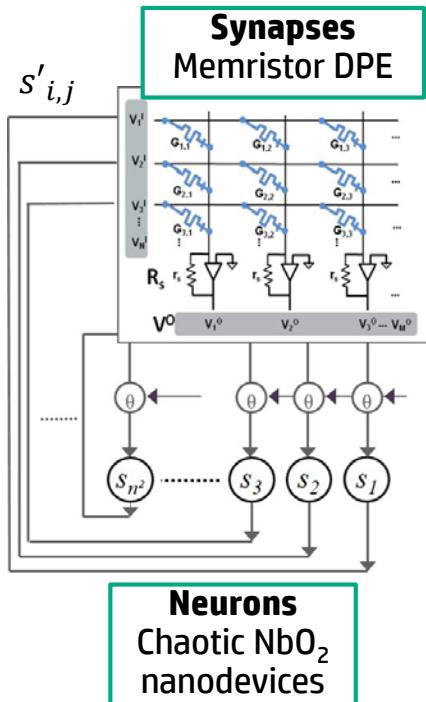
Modified Hopfield network for optimization problems



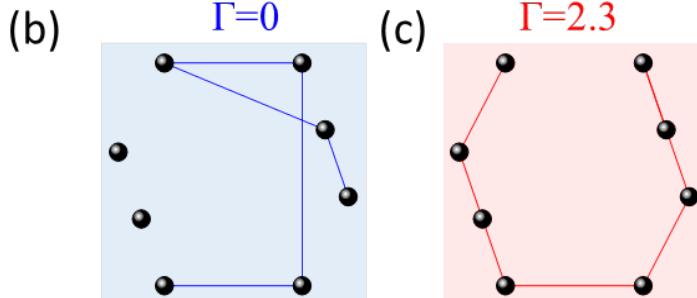
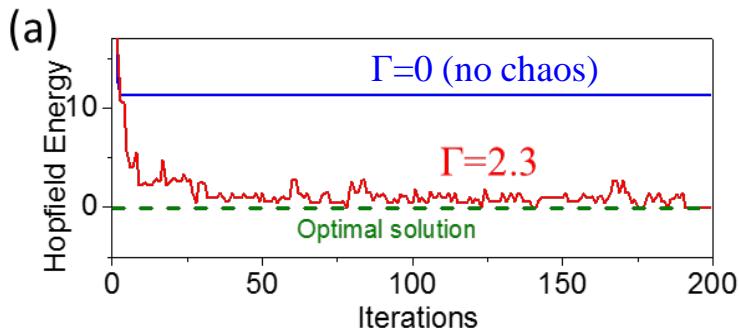
Traveling Salesman
problem

$O(n^2^n)$?

still NP hard



Example solutions w/ and w/o chaos



Acknowledgments



John Paul Strachan
Miao Hu
Suhas Kumar
Catherine Graves
Emmanuelle Merced
Dejan Milojicic
Ali Shafiee
Naveen Muralimanohar
Gary Gibson
Max Zhang



J. Joshua Yang
Qiangfei Xia
Can Li
Yunning Li
Saumil Joshi
Zhongrui Wang

Izzat El Hajj (Univ Illinois Urbana-Champagne)
Pedro Bruel (IME, Universidade de São Paulo)
Sergey Saveliev (Loughborough U, UK)

Funding assistance

