# Using magnetic tunnel junctions to compute like the brain

Mark Stiles with help from a cast of dozens



To understand how the brain works

To develop machine consciousness

To do cognitive computing more efficiently



· PML · NDCD · Alternative Computing Group

#### Hardware for Artificial Intelligence – NIST Gaithersburg



### Computers are designed to solve numerical problems, the brain excels at categorical problems



Satellite trajectories



Natural language processing

Digital Computer	Neuromorphic System
High precision numerical	Categorical

## Computers separate long term memory from processing, while in the brain, they are intimately connected



#### Small scale structure of the brain:



Digital Computer	Neuromorphic System
Memory-Processor Bottleneck	Collocation of processor and memory
High precision numerical	Categorical



#### Computers use a rigid encoding scheme, the brain prioritizes resiliency with more flexible schemes



Digital Computer	Neuromorphic System
Discrete, Deterministic, & Synchronous	Analog, Stochastic, & Asynchronous
Memory-Processor Bottleneck	Collocation of processor and memory
High precision numerical	Categorical



Neural networks provide the simplest implementation of brain-like cognitive computing (rate coding)



- Neurons output a rate representing a spike train.
- The rate from each neuron is multiplied by synaptic weights for each neuron pair.
- Neurons sum the incoming weights and output a non-linear function of the sum (activation function).
- Network non-linearly transforms space to bring similar inputs together

A whole range of approaches to more efficient computing



CMOS has developed well-defined abstraction layers, but with novel devices, designing across the stack is necessary



Details depend on what the device can do.



What is important depends on what the architecture needs



### Magnetic tunnel junctions are multifunctional devices

• Non-volatile embedded memory

• Spin-torque nano-oscillators



Superparamagnetic tunnel junctions



NIST



**Intel**: MRAM integrated into 22nm FinFET CMOS

Neuronal spiking



### Bigger energy barriers mean longer data retention, but more energy needed to write.



#### Magnetic tunnel junctions can be controlled by current.

Positive current Electron flow through tunnel junction Energy  $\rightarrow$  Spin-transfer torque low-resistance high-resistance **I**<sub>MTJ</sub> Parallel Antiparallel Negative current Energy high-resistance low-resistance Electron flow through heavy metal underlayer  $\rightarrow$  Spin-orbit torque Parallel Antiparallel

#### Superparamagnetic magnetic tunnel junctions: Control rate with spin-transfer torque (or spin-orbit torque)



· PML · NDCD · Alternative Computing Group

NIST

Superparamagnetic magnetic tunnel junctions: Control rate with spin-transfer torque (or spin-orbit torque)



**NIST** · PML · NDCD · Alternative Computing Group

#### Computing with superparamagnetic tunnel junctions



· PML · NDCD · Alternative Computing Group

NIST

#### Many reasons to incorporate Complementary Metal Oxide Semiconductor (CMOS)



Author: Cephidan, https://commons.wikimedia .org/wiki/File:LDD-MOS\_transistor\_-\_CMOS\_with\_STI.svg

#### Many reasons to incorporate Complementary Metal Oxide Semiconductor (CMOS)



#### Energies vary over many orders of magnitude





**A. Mizrahi, T. Hirtzlin**, A. Fukushima, H. Kubota, S. Yuasa, J. Grollier, and D. Querlioz, Nature Comm. 9, 1533 (2018);

**A. Mizrahi**, J. Grollier, D. Querlioz, and M. D. Stiles, J. Appl. Phys. 124, 152111 (2018).

### Bigger energy barriers mean longer data retention, but more energy needed to write.



NIST · PML · NDCD · Alternative Computing Group

# Population coding – represent continuous degrees of freedom with multiple spike rates



NIST

#### Simple sense, respond, measure, and learn test



Key advantage = stochastic analog to digital conversion

Continuous learning: the key to a robust system that can adapt to changes

Usually system is only trained once because it requires a lot of energy BUT unable to adapt to changes!



System equipped with continuous learning





Here,  $\Delta E = 15 k_B T$  magnetic tunnel junctions can be used for memory with no precision loss



#### Using unreliable synapses lowers the energy consumption



#### Using unreliable synapses lowers the energy consumption

100 neurons Normalized power consumption 46 neurons 0.1 Write current 20 neurons  $I \propto \Delta E_w$ 0.01 10 neurons (Sato et al., APL 2014) 0.001 2 3 Write energy 5 10 11 12 9 6 Error (%)  $Energy \propto \Delta E_w^2$ 

Write current  $I \propto \Delta E_w$ (Sato et al., APL 2014) Write energy  $Energy \propto \Delta E_w^2$ Example: For 3% precision  $\rightarrow$  54 neurons in each population (2916 weights)

 $\rightarrow \Delta E_{w} = 12 k_{B}T$ 





**M. W. Daniels**, **A. Madhavan**, P. Tatatchian, A. Mizrahi, and M. D. Stiles, Physical Review Applied, 13, 034016 (2020).

Stochastic computing represents numbers as probabilistic bit streams



- Resilient to white noise
  - Incorrect bit  $\Rightarrow$  error  $\approx O(1)$
  - Compare to binary: error  $\approx O(2^N)$
  - Naturally coherent with bilevel devices (magnetic tunnel junctions).
- Interpreting these "spike trains" as **rates** or **probabilities** restricts us to 0 ≤ p ≤ 1.

Unbiased superparamagnetic tunnel junction plus CMOS for low energy random bitstreams



**NIST** • PML • NDCD • Alternative Computing Group

AND gates for multiplication – Programmable bitstream generator: 10 fJ per bit



Typical stochastic computing uses Linear Feedback Shift Registers  $\rightarrow$  Short period pseudorandom numbers

## Synapses from bitstream generators (weights) and AND gates (multiplication)



Programmable Bit stream generator

Neurons from OR gates for nonlinear summation





 $p_{\rm OR} = p_a + p_b - p_a p_b$ 

OR-gate neuron does not work with correlated bit streams

#### Low area, high efficiency neural network with stochastic information throughout the calculation



Efficient random bitstream generation and avoiding domain translation – energy efficiency on a standard problem

- Energy/performance tradeoff by tuning circuit parameters.
- Running for longer or shorter total time ( $\sim N$ ) gives energy/performance tradeoff.



Future: implement primitives to check feasibility with realistic MTJs

· PML · NDCD · Alternative Computing Group NIST

### Collaborators



Brian Hoskins, **Matthew Daniels**, Guru Khalsa, Mark Anders, Jonathan Goodwill, Jabez McClelland, Nikolai Zhitenev, William Rippard, Matthew Pufall, Emilie Jué, <u>Mark Stiles</u>



#### Alice Mizrahi, Advait Madhavan, Philippe Talatchian



Siyuan Huang, Gina Adam

UC SANTA BARBARA

George Tzimpragos, Tim Sherwood



Jacob Torrejon, Mathieu Riou, Flavio Abreu Araujo, Paolo Bortolotti, Vincent Cros, Julie Grollier



Tifenn Hirtzlin, Damien Querlioz

Sumito Tsunegi, Kay Yakushiji, Akio Fukushima, Hitoshi Kubota, Shinji Yuasa



Energy Frontier Research Center Department of Energy

Vivek Amin, Jonathan Gibbons, Paul Haney, Axel Hoffmann, Julie Grollier



NIST

### Using magnetic tunnel junctions to compute like the brain

- Emulating features of the brain can increase efficiency for cognitive computing.
- Applications require designing across the computational stack
  - Architecture
  - Encoding
  - Circuitry
  - Devices



- CMOS is likely to play an important role in any room-temperature computing scheme.
- Magnetic tunnel junctions
  - Multifunctional
  - Already in CMOS Fabs

Architecture	Encoding	Devices
Population coding	Fluctuation rates	Superparamagnetic tunnel junctions
Synaptic weights	Binary	Unstable magnetic tunnel junctions
Neural network with stochastic computing	Random bitstreams	Superparamagnetic tunnel junctions, Repurposed CMOS Logic gates

