# Highly Efficient Neuromorphic Computing Systems with Emerging Nonvolatile Memories
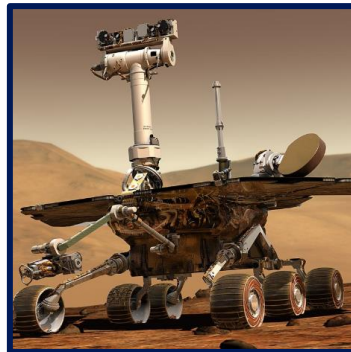
## Bonan Yan

Dept. Electrical & Computer Engineering

Duke University
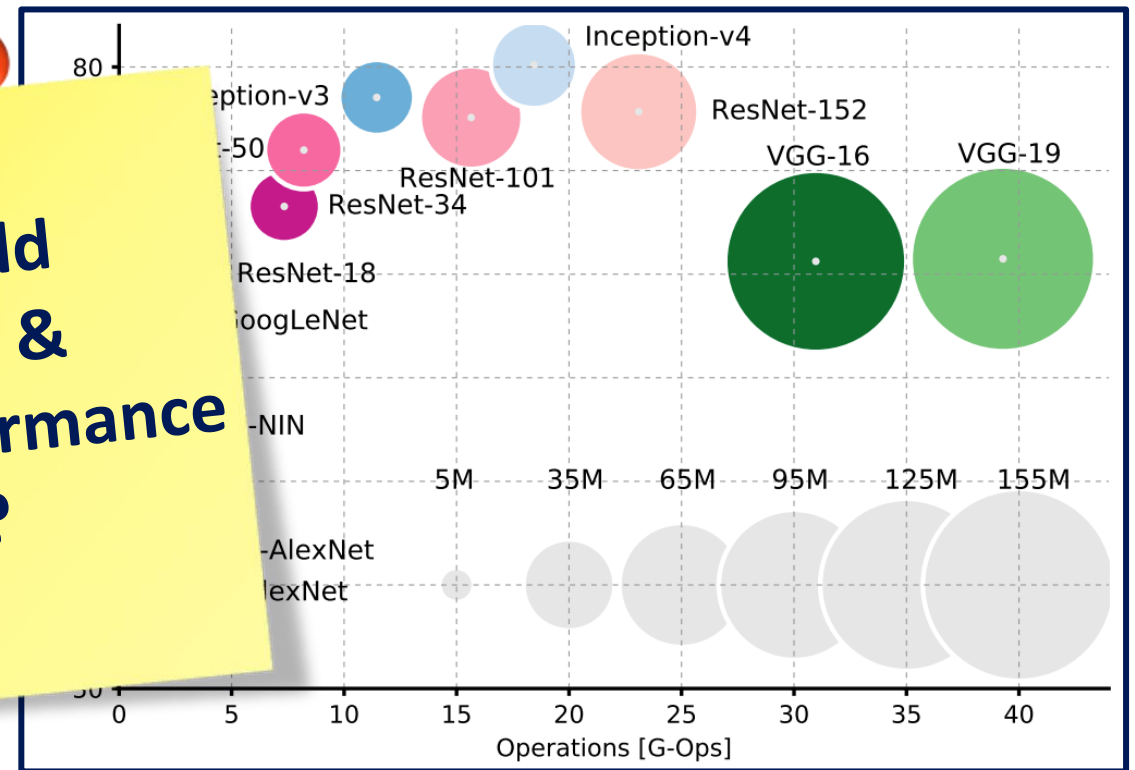
Slides available at: https://bonanyan.github.io/bn/

# Efficiency Is The Key to Ubiquitous AI

## Limited Power/Energy

## Better Accuracy Comes From Larger Models



How to build low-power & high-performance hardware?

# Overhead Dominated by Memories

## *Power*

On-Chip Fabric

Computation

12%

33%

IO

30%

22%

Clocks: 3%

CPU Caches

## *Area*



Interposer

Chiplet 1    Chiplet 2

## *Energy*

| Operations (32bit) | Energy (pJ) |
|---|---|
| Int ADD | 0.1 |
| Float ADD | 0.9 |
| Register File | 1 |
| Int Multiply | 3.1 |
| Float Multiply | 3.7 |
| SRAM Cache | 5 |
| DRAM Memory | 640 |

Source: AMD, Intel, [Mu-Shan Lin, et al, JSSC 2019]
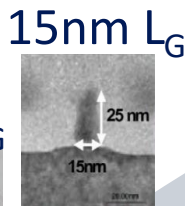
Duke

3

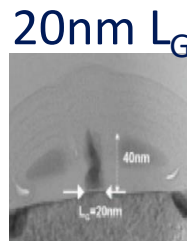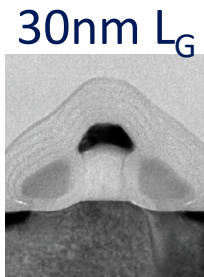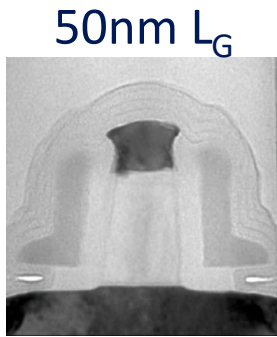# Specialized Hardware Enhances Efficiency

$$P = \alpha C V_{DD}^2 f$$

$P$: Power Dissipation
$\alpha$ : Activity Factor
$C$: Load Capacitance
$V_{DD}$: Power Supply
$f$ : Clock Frequency

15nm L$_G$

20nm L$_G$

30nm L$_G$

50nm L$_G$

MOSFET Scaling

Transistor Size &
Clock Frequency

Multicore to
Manycore

Domain-Specific

Google TPU v3
~200Watts
(4 TPU to run AlphaGo,
300x less power
consumption)

Duke

4

# Uniqueness of Neural Network Execution

**Multiply-Accumulate (MAC):**

inputs

weights

outputs

streaming

stationary

$$[x_1\ x_2\ \dots\ x_m]\begin{bmatrix} w_{11} & \cdots & w_{1n} \\ \vdots & & \vdots \\ w_{m1} & \cdots & w_{mn} \end{bmatrix} = [y_1\ y_2\ \dots\ y_n] \quad \Longrightarrow \quad y_j = \sum_{i=1}^{m} x_i \cdot w_{ji}$$

**Execution:**

MAC Unit

DRAM

Weight

Input Data

Partial Sum

× +

Up
Part

Inference:
Inputs change;
weights stay.
How to fix weights
to where MAC Units
without moving?

Duke

# Make Memory Access Less Expensive

| | CMOS Accelerator (Planar Integration) | Near-Memory Computing | In-Memory Computing |
|---|---|---|---|
| **Design Approach** | Off-Chip DRAM — Global Buffer 108KB — PE Array (168 PEs) — 3.5mm | Logic Die — 3D-stacked Memory — SoC — Interposer — High Bandwidth Memory | $\begin{bmatrix} x_1 \\ x_2 \\ \dots \\ x_m \end{bmatrix}$ Row Circuitry — Col Circuitry — $[y_1 \quad y_2 \quad \dots \quad y_n]$ |
| **Energy to Access Weights** | 200×~6× of MAC Energy | 55×~6× of MAC Energy | **0** |

DRAM        On-Chip Buffer

Source: AMD, Y.-H. Chen, JSSCC,2017

# Key Idea of In-Memory Computing

inputs        weights

outputs

$$[V_1 \quad V_2 \quad V_3] \begin{bmatrix} G_{11} & G_{12} & G_{13} \\ G_{21} & G_{22} & G_{23} \\ G_{31} & G_{32} & G_{33} \end{bmatrix} = [I_1 \quad I_2 \quad I_3]$$

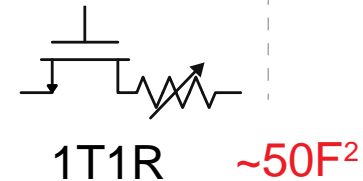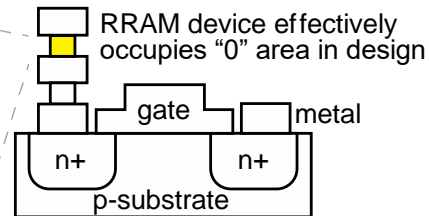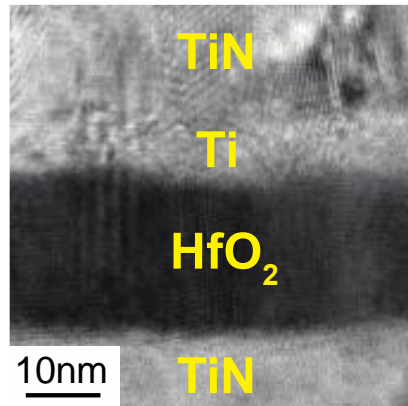- Weight Matrix Stored as Conductance $\boldsymbol{G}$

- Rely on Analog Computation (Kirchhoff's Current Law) for "almost free"
  - Multiplication: $I = V \cdot G$
  - Addition: $I^{column} = I_1^{row} + I_2^{row} + I_2^{row}$

- Ideal Nanoscale Devices for $\boldsymbol{G}$:
  - Programmable Conductance
  - Multi-Level Cell
  - Small Footprint/High Density
  - Compatible with Existing CMOS Process



## Duke

# Memristors for In-Memory Computing

## Also Called **R**esistive **R**andom **A**ccess **M**emory, RRAM or ReRAM

Cross-section TEM



TiN
Ti
HfO₂
TiN
10nm

RRAM device effectively occupies "0" area in design

gate    metal

n+    n+

p-substrate

1T1R    ~50F²



Programmable resistor w/ analog states

ISSCC: Intel adds embedded ReRAM to 22nm portfolio

January 03, 2019

TSMC to start embedded RRAM production in 2019

According to reports, Taiwan Semiconductor Manufacturing Company (TSMC) is aiming to start producing embedded RRAM chips in 2019 using a 22 nm process. This will be initial "risk production" to gauge market reception.

|  | Multi-Level Cell | Cell Area | R/W Speed |
|---|---|---|---|
| SRAM | × | large | Fast |
| DRAM | × | medium | Medium |
| 1T1R | √ | medium | Medium Fast |
| Flash | √ | small | Slow |

**Duke**

# My work: Emerging Memory-Centric Design

- **Circuits & Systems Implementation**

  - Spike-based Interface [DAC'15, DAC'18, DAC'20]

  - Implementation of Neural Networks [VLSI'19, DAC'20]

- **Tolerate/Exploit Non-ideal Behavior of Memristors**

  - Device Nonlinearity [ISCAS'16, IEDM'17, IEDM'19]

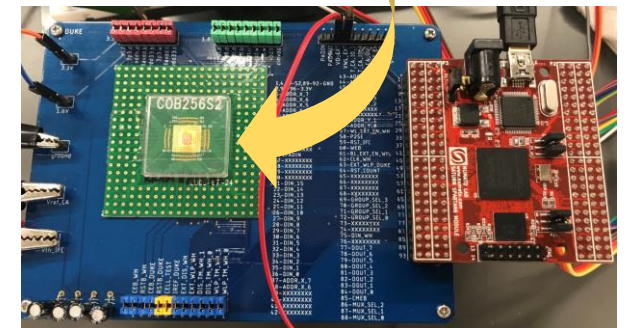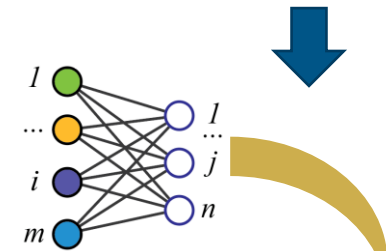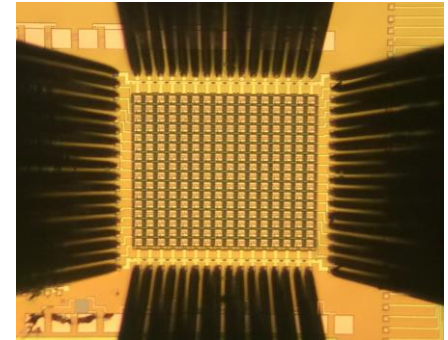  - Read Disturbance [ICCAD'17], Hard Fault [ITC'19]
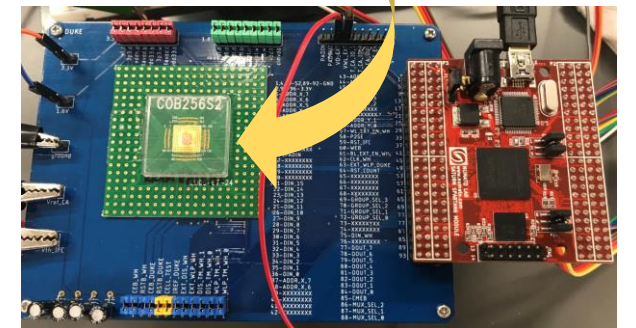
**Duke**

# My work: Emerging Memory-Centric Design



- **Circuits & Systems Implementation**

  - Spike-based Interface [DAC'15, DAC'18, DAC'20]

  - Implementation of Neural Networks [VLSI'19, DAC'20]

- **Tolerate/Exploit Non-ideal Behavior of Memristors**

  - Device Nonlinearity [ISCAS'16, IEDM'17, IEDM'19]

  - Read Disturbance [ICCAD'17], Hard Fault [ITC'19]

# Duke

# Conventional ADC is Too Large

**The Level-based Design**

- Compatible to existing signal processing

- High speed computation



256x256 1T1R Array

Actual Layout:



**Unable to Compute Parallel!**

ADC          ADC

Based on 1.66MF$^2$ 8bit ADC
by K. Ohhata (JSSC 2019)

Duke

11

# My Approach: Spiking Interface Circuit

## What Better Designs Look Like

- Compute Parallelly (Massive)
- Need Light-Weight Interface Circuitry



## The Spike-based Design

- Closer to biological system
- Extremely high power efficiency

# Spike Conversion

**Current Amplifier (CA)**

**Integrate & Fire Circuit (IFC)**

Column Current

BL

Vth

Vcap

Cm

- Comparator +

Reset Transistor

CLK

$\overline{\text{CME}}$ — Computing mode is enabled

ISNA_EN

cycle [i] with input [i]

cycle [i+1] with input [i+1]

Buffered spikes

Spike freq

Linear Region

Nonlinear Region

Column Current

Spike Conversion Circuit & Controller 3152x3152 µm²

*B. Yan, et al., ISCAS, 2016*

13

# Spike Conversion



Positive feedback to accelerate startup

Phase compensation

Scale column current

Stabilize BL

CA drives BL

SL sink current

Tradeoff between large input current range and response speed

Enlarge phase margin tolerating capacitor positive feedback

- 5.12 GOPS @ 8-bit precision
- 200ns latency @ 8-bit precision
- 43x area reduction vs. ADC by K. Ohhata (JSSC 2019)

**Duke**

*B. Yan, et al., IEEE VLSI Symp. on Tech. 2019*

14

# How to Use Spiking-Based Design to Execute Neural Networks?

# In Situ Nonlinear Activation (ISNA) Function

**Fully-Connected (FC) Layer**



**Convolutional (Conv) Layer**



Single-layer Inference Operation:

**D** • Step 1: Load data from buffer

**A** • Step 2: Vector-matrix multiplication

**D** • Step 3: Nonlinear activation function

**D** • Step 4: Pooling

**D** • Step 5: Store results to buffer

**D** : digital domain     **A** : analog domain

*B. Yan, et al., IEEE VLSI Symp. on Tech. 2019*

# In Situ Nonlinear Activation (ISNA) Function

**Fully-Connected (FC) Layer**



**Convolutional (Conv) Layer**



Single-layer Inference Operation:

**D** • Step 1: Load data from buffer

**A** • Step 2: Vector-matrix multiplication
   • Step 3: Nonlinear activation function

**D** • Step 4: Pooling

**D** • Step 5: Store results to buffer

**D** : digital domain   **A** : analog domain

Combine Step 2 & Step 3 to simplify PE operation:
Use linear + nonlinear regions

**_B. Yan_**, _et al., IEEE VLSI Symp. on Tech. 2019_

17

# Adjust Activation Function



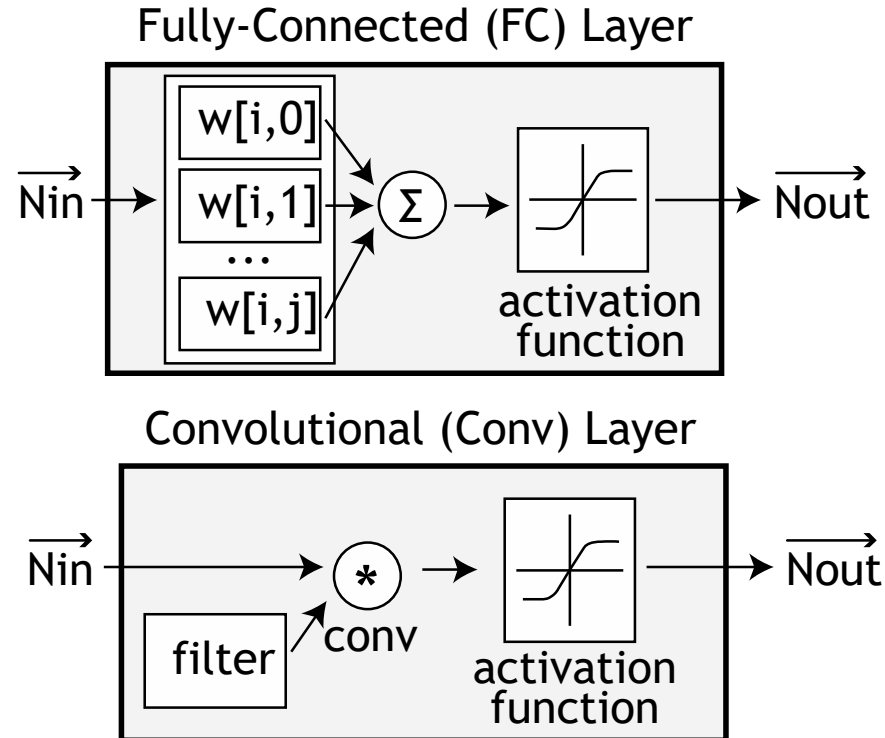Vref: Read voltage of BL
- Vref ↑, BL current ↑
- Proportional tuning

Vth: Threshold of capacitor charging/discharging
- Vth ↓, Charging/discharging ↑
- Distorted tuning



Measured ISNA behavior

$(V_{ref}, V_{th})=$
- (0.30V, 0.35V)
- (0.30V, 0.40V)
- (0.35V, 0.35V)
- (0.35V, 0.40V)

Function clipped-relu

≈

tanh (x>0)

*B. Yan*, et al., IEEE VLSI Symp. on Tech. 2019

# Mixed-Signal Chip Architecture

Load Well-Trained Weights →

Read/Write Controller
[FIFO interface included]

Read/Write Circuits

WL Decoders

WL Drivers

Column Decoder

Memristor Crossbar Array

Column Decoder

Input Vectors →

Compute Controller

Spike Conversion Circuits
(In Situ Nonlinear Activation)

↓

Output Vectors

*B. Yan*, et al., IEEE VLSI Symp. on Tech. 2019

Duke

19

# Mixed-Signal Design Flow

Load Well-Trained Weights →

Synthesis Design Flow

Using Verilog HDL

Custom Design Flow

Transistor-level Design

Input Vectors →

| Read/Write Controller [FIFO interface included] | | Read/Write Circuits |
|---|---|---|
| WL Decoders | WL Drivers | Column Decoder |
| | | Memristor Crossbar Array |
| | | Column Decoder |
| Compute Controller | | Spike Conversion Circuits (In Situ Nonlinear Activation) |

↓

Output Vectors

Duke

# Chip Summary



RRAM Nonvolatile CIM PE

| | |
|---|---|
| Computing Controller | Counter |
| *ISNA* | CA+IFC |
| | BL/SL Decoder |
| WL Driver | RRAM Array |
| *RRAM Macro* | |
| Memory Controller | RRAM RD/WR Logic |

| Technology | 150nm CMOS +HfO$_x$ RRAM |
|---|---|
| Macro Capacity | 64K (256×256) |
| Clock Frequency | 50MHz |
| Energy Efficiency | 0.257pJ/Mac |
| Average Power | 1.52 mW |
| Layer-wise Latency | 200ns |
| Real-time Benchmarks | 3-layer perceptrons, LeNet-4, LetNet-5 |

Demo video online: https://bit.ly/AICHIP

## Duke

Die Photo of
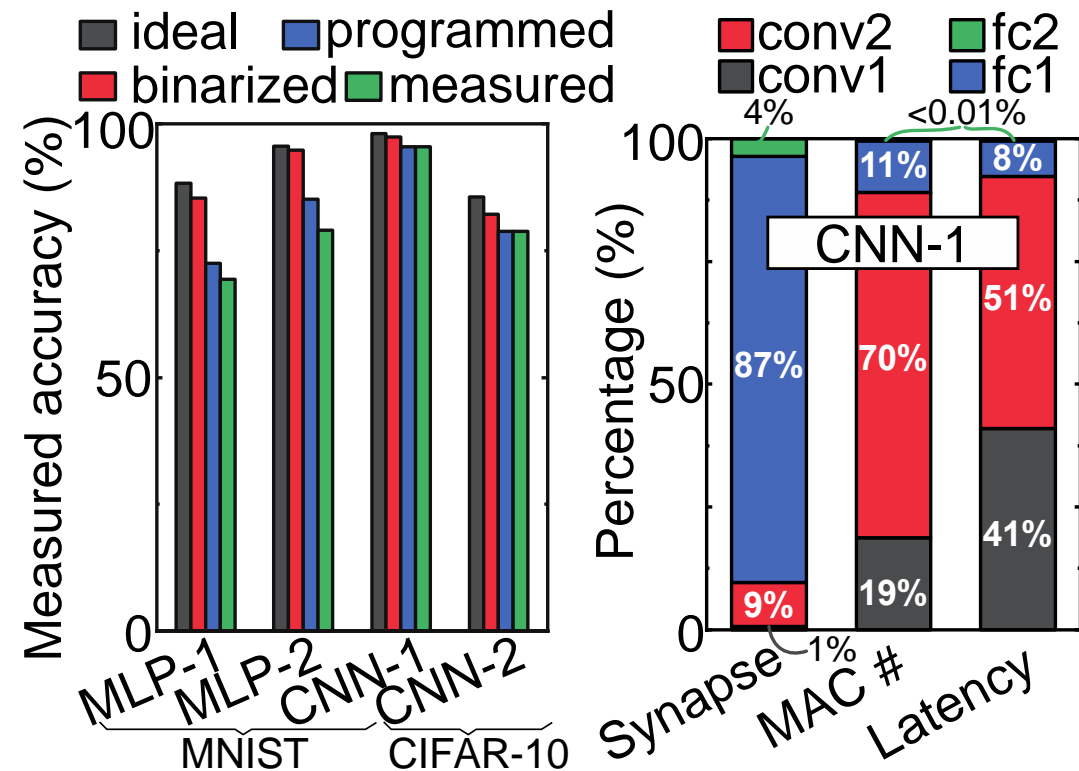RRAM Neuro-computing Engine

GUI

Camera

# Evaluation: Measured Neural Network Results

**MNIST:**

- MLP-1: Single-layer perceptron

- MLP-2: 2-layer perceptron
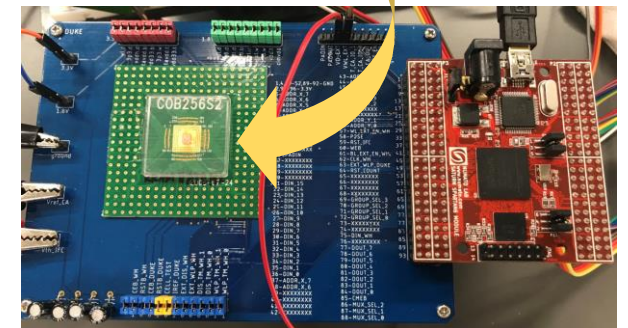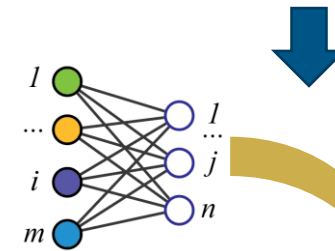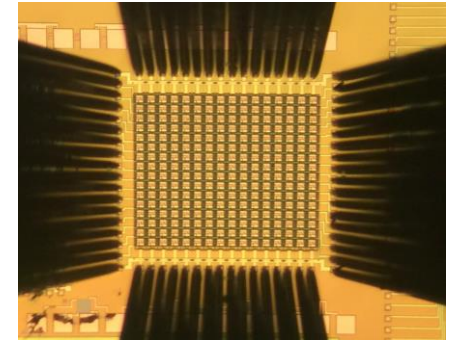
- CNN-1: 4-Layer LeNet

**CIFAR-10:**

- CNN-2: 5-Layer LeNet

*B. Yan, et al., IEEE VLSI Symp. on Tech. 2019*

Duke

# My work: Emerging Memory-Centric Design
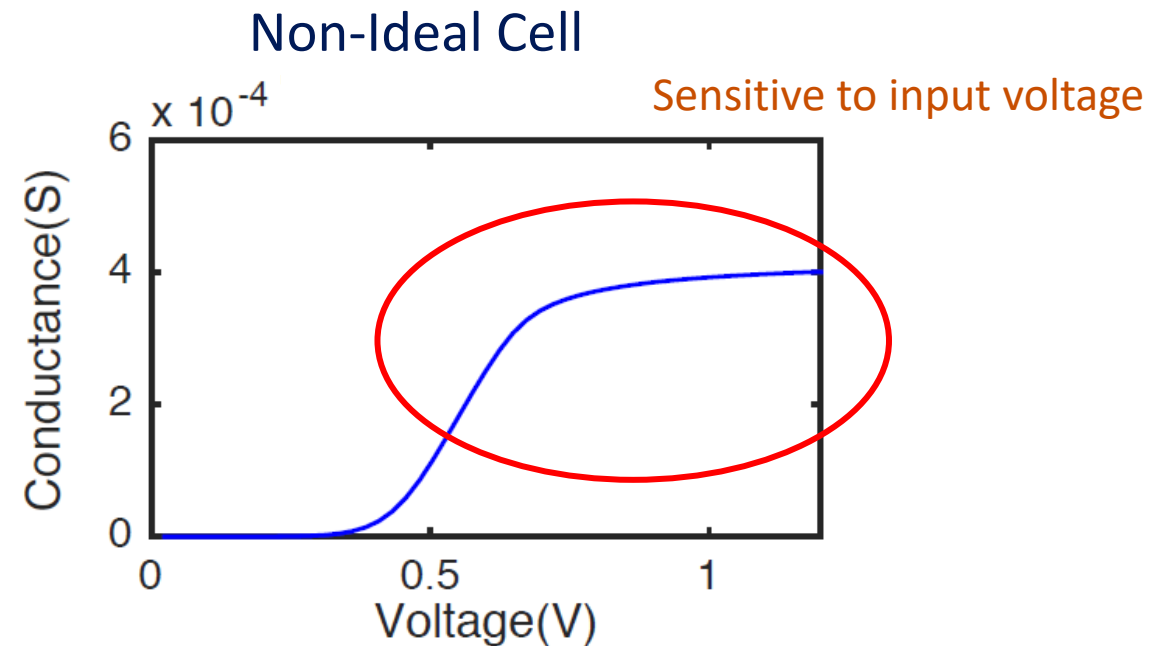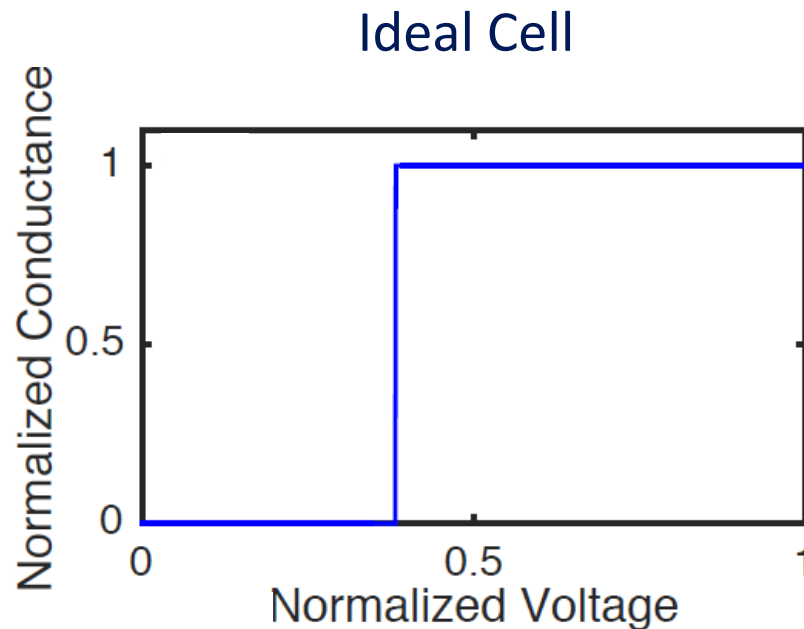
- **Circuits & Systems Implementation**

  - Spike-based Interface [DAC'15, DAC'18, DAC'20]

  - Implementation of Neural Networks [VLSI'19, DAC'20]

- **Tolerate/Exploit Non-ideal Behavior of Memristors**

  - Device Nonlinearity [ISCAS'16, IEDM'17, IEDM'19]

  - Read Disturbance [ICCAD'17], Hard Faults [ITC'19]



**Duke**

# Non-Ideal Memristor - I

- Cell Nonlinearity: conductance varies when applied with different voltages

Ideal Cell

Non-Ideal Cell

Sensitive to input voltage

- Solution: Current Amplifier to Clamp Cell Voltage (shown in previous circuit design part)

**Duke**

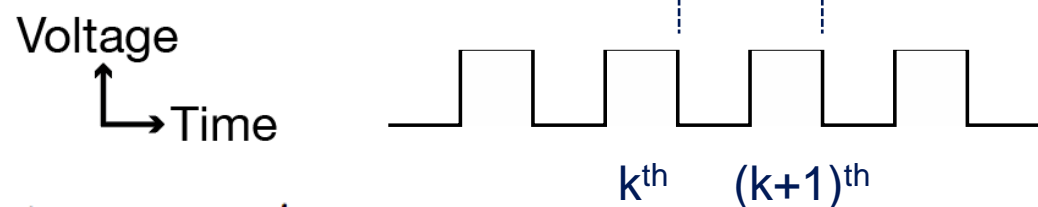*B. Yan*, *et al., ISCAS, 2016*

# Non-Ideal Memristor - II

- Memristance Drift: conductance/memristance gradually deviates from original values under read voltage (read disturbance)

- Characteristic:
  - Very slow to observe



*B. Yan*, et al., ICCAD, 2017

# How to Mitigate Memristance Drift for Inference?

$$E_{mse} = \sum_{j}^{n} \left( t_j - \sum_{i}^{m} w_{ij} x_i \right)^2 \qquad E'_{mse} = \sum_{j}^{n} \left( t_j - \sum_{i}^{m} w_{ij} x_i - \sum_{i}^{m} \Delta w_{ij} x_i \right)^2$$
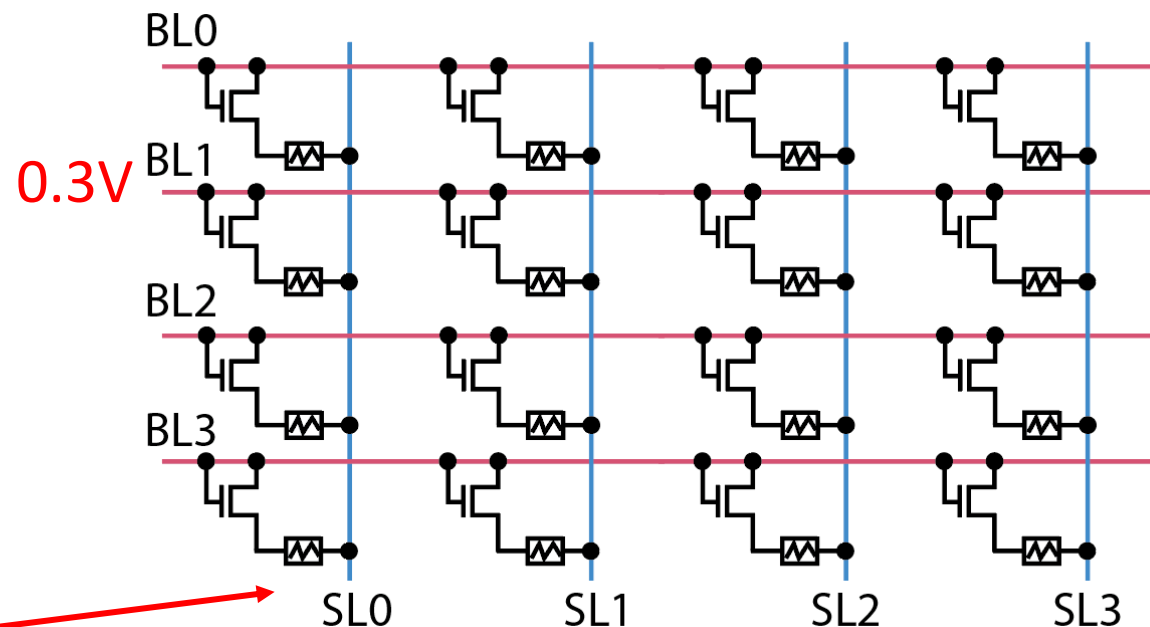
$E_{MSE}$: MSE Error Function
$t_j$: label for each neuron

Voltage

Time

$k^{th}$     $(k+1)^{th}$

$$\Delta E_{mse} = E'_{mse} - E_{mse}$$

$$= -2 \sum_{j}^{n} \left[ \left( t_j - \sum_{i}^{m} w_{ij} x_i \right) \left( \sum_{i}^{m} \Delta w_{ij} x_i \right) \right]$$

Minimize Degraded Error Function:

$$\frac{\partial \Delta E_{mse}}{\partial \Delta w_{ij}} = -2 \left( t_j - \sum_{i}^{m} w_{ij} x_i \right) x_i \leq 0$$

BL0

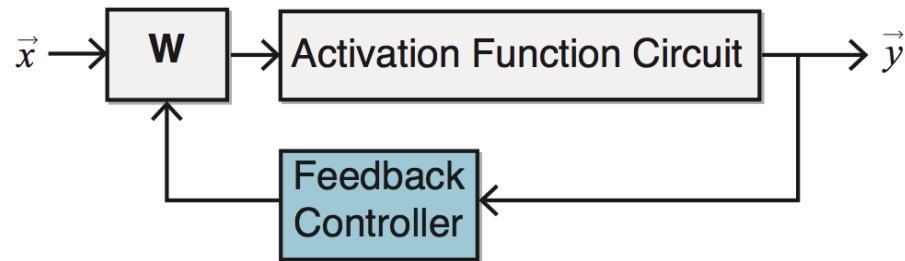0.3V   BL1

BL2

BL3

SL0     SL1     SL2     SL3
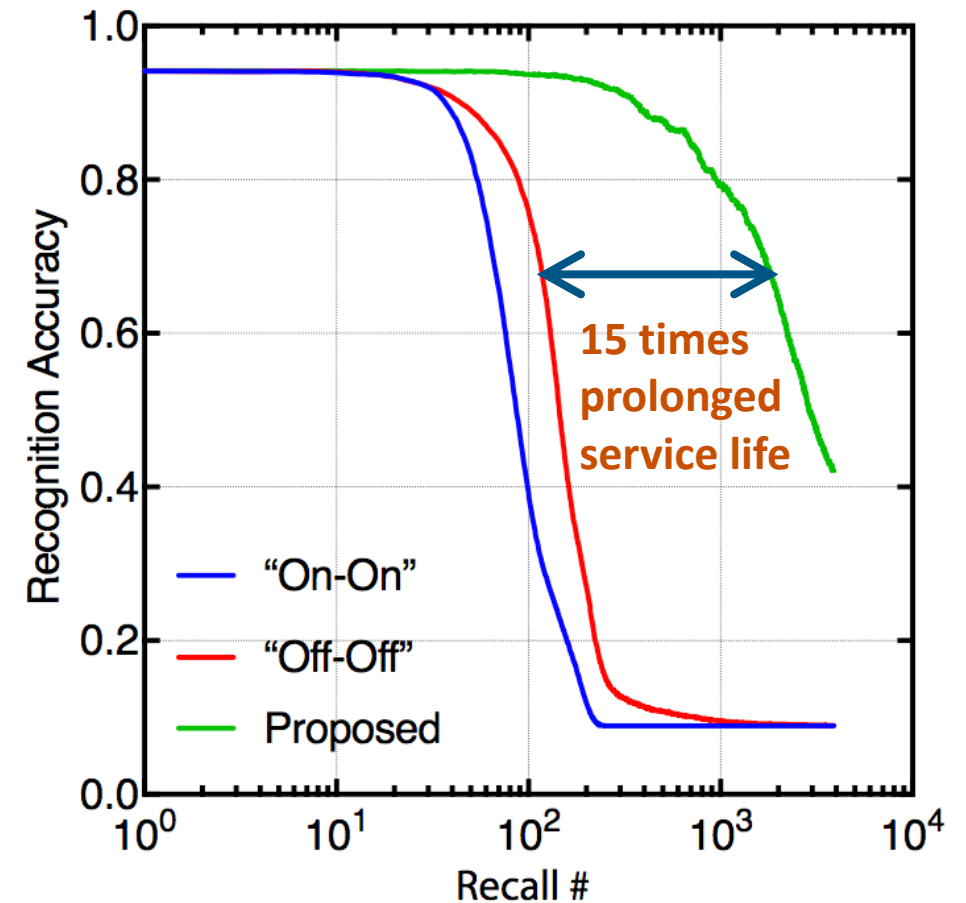
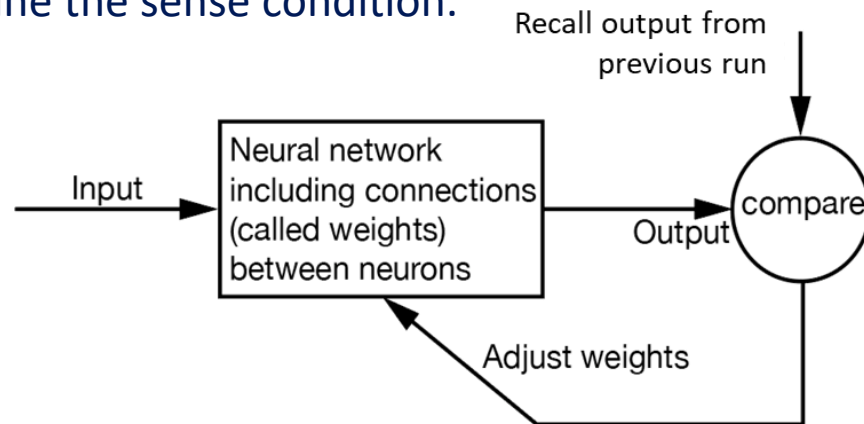Change Column Current Direction    0V/0.6V

Duke

# Closed-loop Design to Enhance Weight Stability

**Feedback controller:** Adjust the voltage condition to compensate the memristance drift.
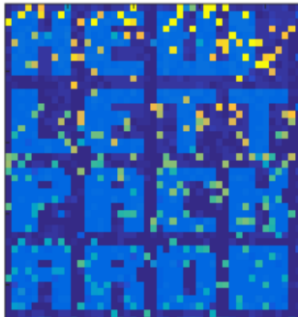


**"Arrogant principle":** last output is used as the label to determine the sense condition.



**15 times prolonged service life**

- "On-On"
- "Off-Off"
- Proposed

*B. Yan, et al, ICCAD, 2017*
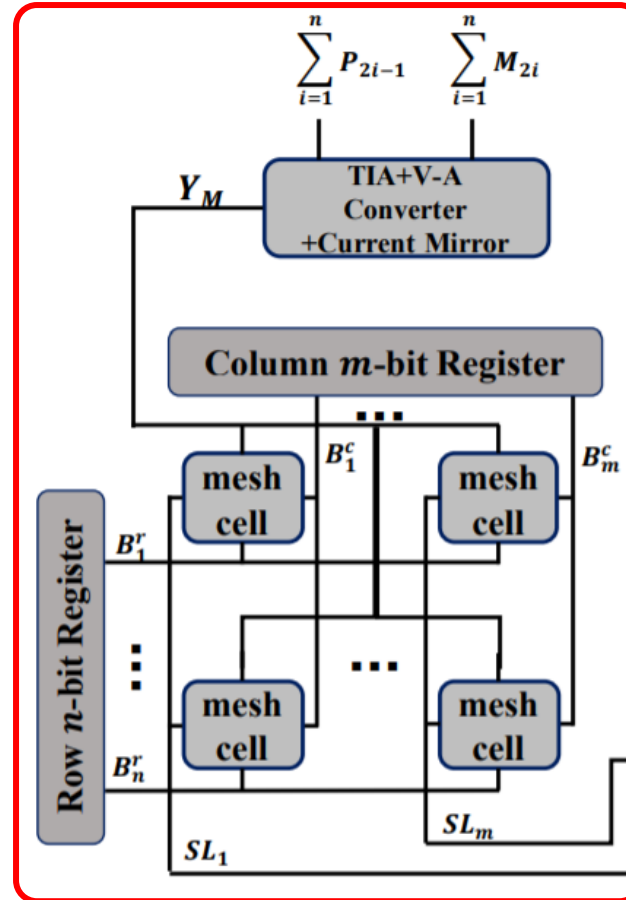
Duke

28

# Non-ideal Memristor-III

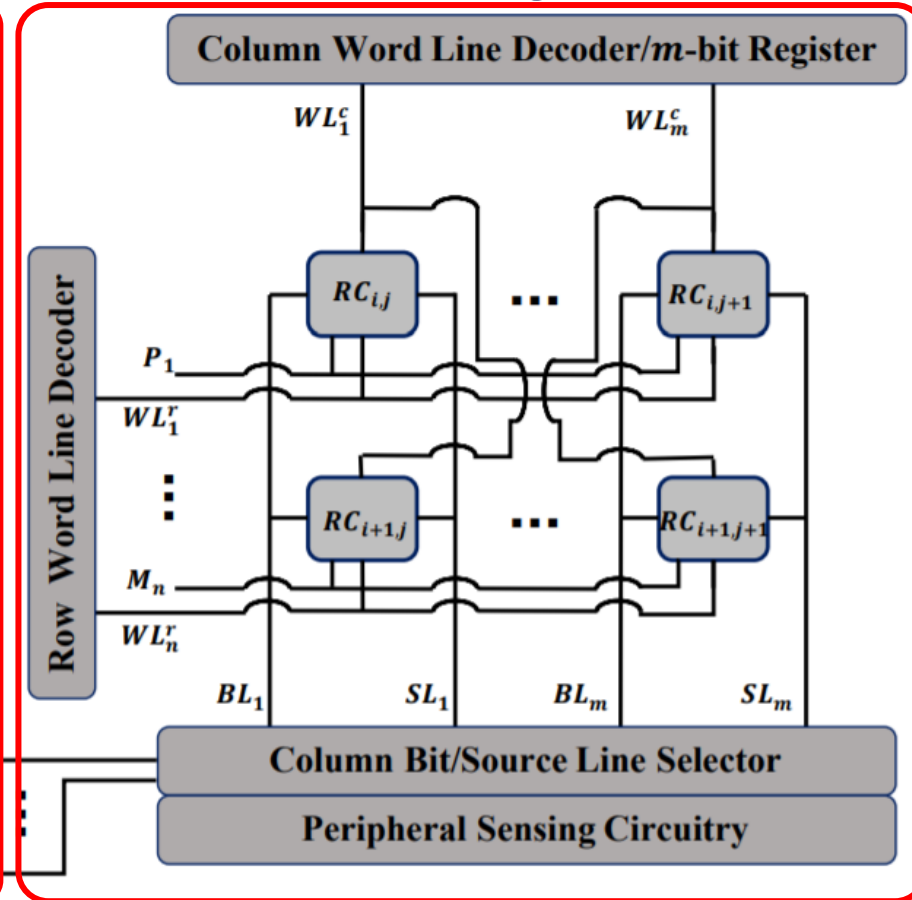- Hard Faults: Dead Cells that are not programmable

  *Noisy "Hewlett Packard", 2016*

- Sources of Hard Faults:
  - Fabrication
  - Beyond Endurance
- Our Solution: Error Correction Circuits

## Error Correction Circuit



## Core Engine



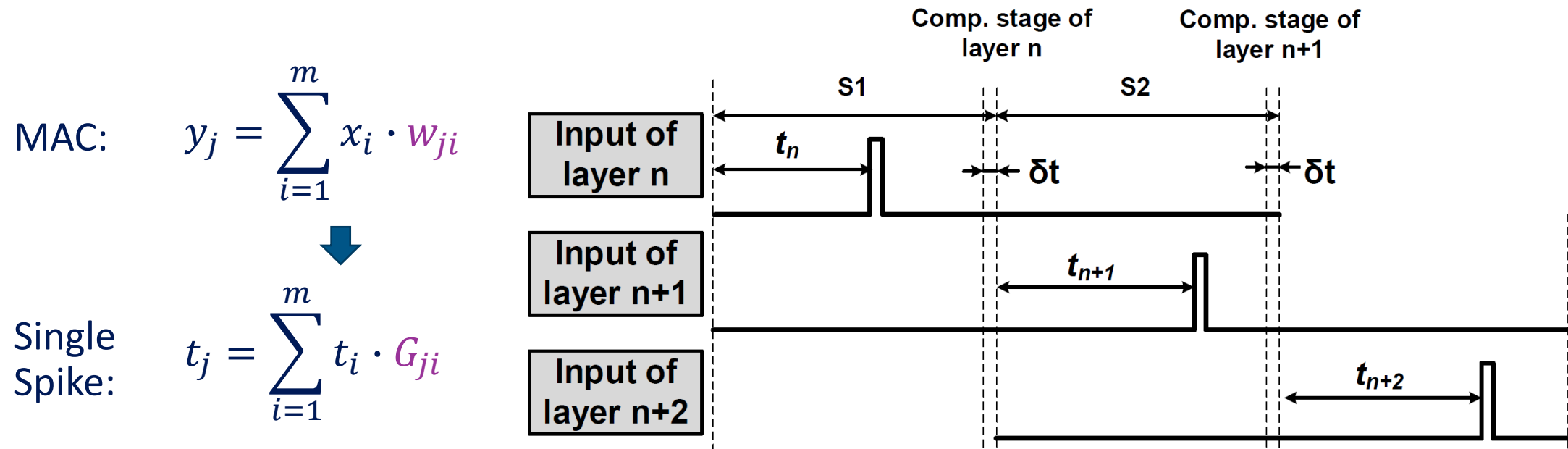*A. Chaudhuri, **B. Yan**, et al, ITC, 2019*

# Summary

- The Past and Now of AI Hardware

  - In-Memory Computing Eliminates Weights Movement

- My Work:

  - Spiking-based In-Memory Computing Engine Offer Very High Energy Efficiency & Good Performance

  - Clever Design Methodologies (e.g., Closed-Loop Sensing) is Effective to Tolerate Nonideal Features of Memristors

Duke

# Future Work I: Single Spike Processing Engine

Use Single Spike to Replace Multiple Spikes: ~10x Improvement of Energy Efficiency

MAC:
$$y_j = \sum_{i=1}^{m} x_i \cdot w_{ji}$$

Single Spike:
$$t_j = \sum_{i=1}^{m} t_i \cdot G_{ji}$$

Comp. stage of layer n     Comp. stage of layer n+1

S1     S2

Input of layer n    $t_n$    δt    δt

Input of layer n+1    $t_{n+1}$

Input of layer n+2    $t_{n+2}$

"ReSiPE: ReRAM-based Single-Spiking Processing-In-Memory Engine"
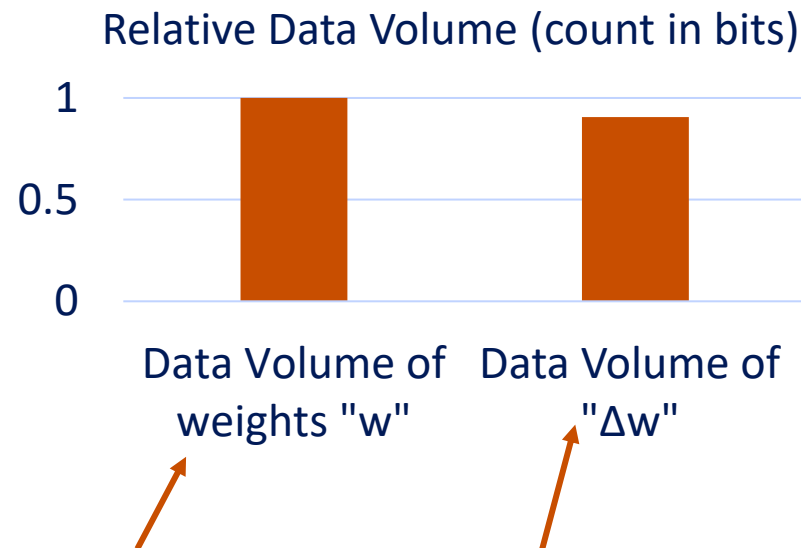Simulation Results Coming in July at Design Automation Conference (DAC) 2020

Duke

Expect to IC Tape-Out + Single-Spike Encoding Method [Potential Collaboration Opportunity]

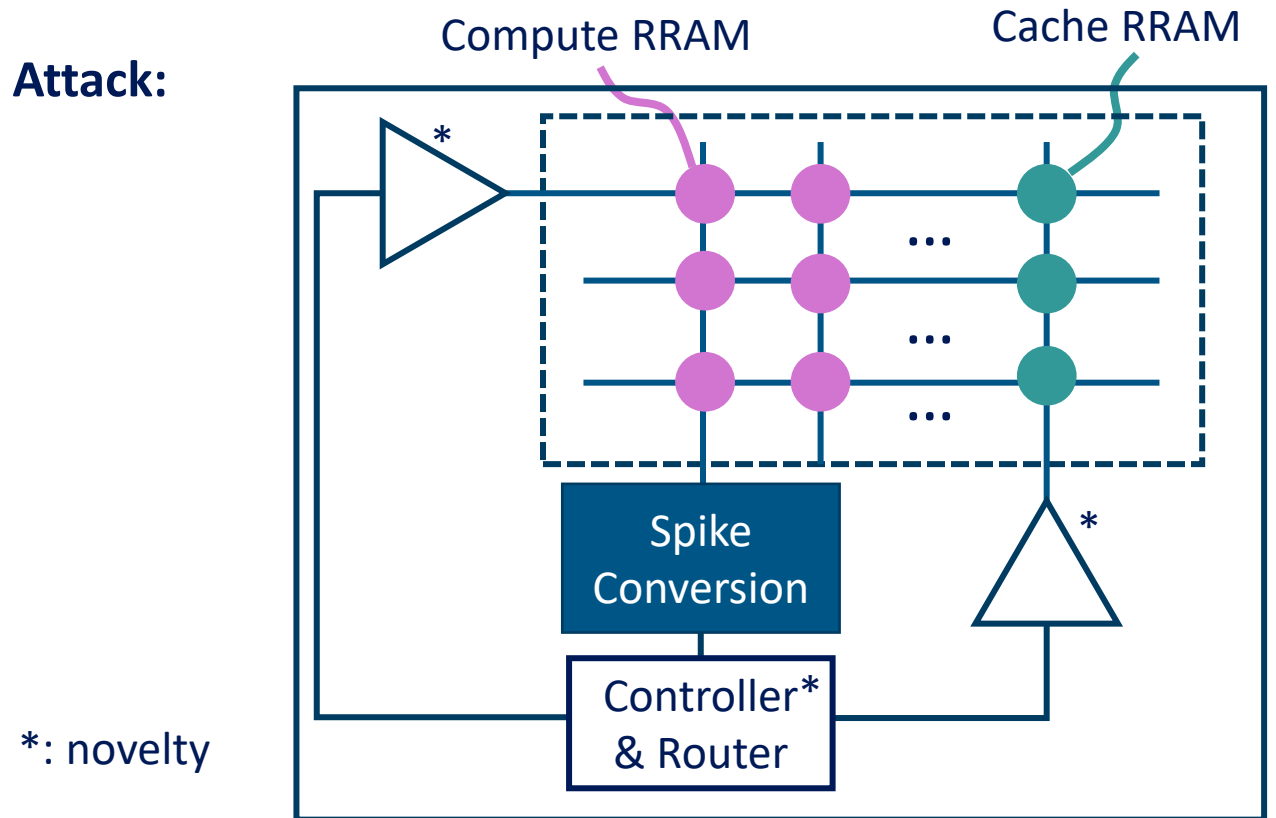# Future Work II: In-Memory Compute & Cache

## Challenge:

During Training:

Relative Data Volume (count in bits)



Data Volume of weights "w"

Data Volume of "Δw"

No Need to Transport Weight

**External** DDRAM Access (Expensive!)

## Attack:



Compute RRAM

Cache RRAM

Spike Conversion

Controller* & Router

*: novelty

**Prospect Applications:**

NN Training Acceleration
Self-Updatable In-Memory Computing Engine

# Long-Term Prospects

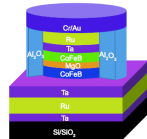| Compute | Control |
|---|---|

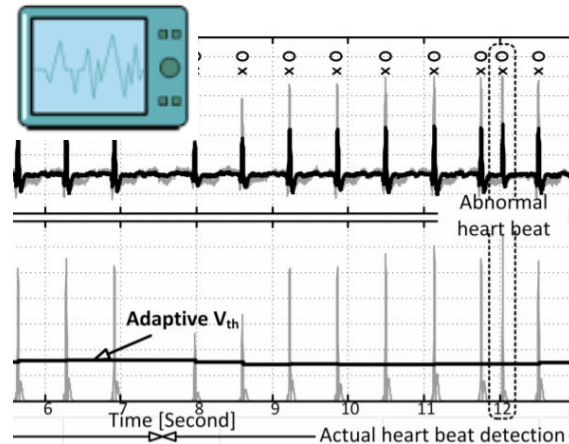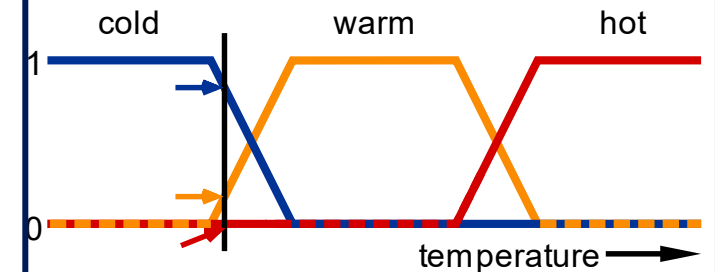**Rethink Compute Memory Hierarchies**



RRAM          PCM          MRAM

"Organ" for Compute & Cache w/ different emerging memory types

**In Situ Data Processing (Near-Sensor Computing)**



**Analog Content Addressable Memories**



Fuzzy Logic Hardware

# Many Thanks for the Support!

# Thanks for Listening & Qs!

slides available at:

https://bonanyan.github.io/bn/

Duke

# Backup Slides

# In-Memory Computing Candidates

| RRAM/Memristor | Phase-Change Memory | MRAM | Flash | SRAM |
|---|---|---|---|---|
|  |  |  |  |  |
| Nonvolatile | Nonvolatile | Nonvolatile | Nonvolatile | Volatile |
| Dense Cell Area | Dense Cell Area | Dense Cell Area | Dense Cell Area | Large Area |
| Bipolar Write | Unipolar Write | Bipolar Write | Erase-Needed Write | Bipolar Write |
| Fast Write | Fast Write | Fast Write | Slow Write | Ultra-Fast Write |
| Multi-level Cell | Multi-level Cell | Single-level Cell | Most Single-level Cell | Single-level Cell |

Duke

Zhang et al., JSSC 2017; Fick et al., CICC 2017; Zhang et al., T-ED 2015
Chen et al., ISSCC 2018; Boybat et al., Nature communications 2018

# Comparison

| Type | CIM Macro | | | | CIM PE | | Digital Processor |
|---|---|---|---|---|---|---|---|
| Work | VLSI'18 Panasonic | ISSCC'18 NTHU | ISSCC'18 NTHU | ISSCC'18 MIT | This work[†] | ISSCC'18 UIUC | TrueNorth [10] |
| Technology | 180nm | 65nm | 65nm | 65nm | 150nm | 65nm | 28nm |
| Synapse | 1T1R RRAM | 1T1R RRAM | 6T-SRAM | 10T-SRAM | 1T1R RRAM | 6T-SRAM | SRAM |
| Nonvolatility | Yes | Yes | No | No | Yes | No | No |
| Standby current | ~zero | ~zero | high | high | ~zero | high | high |
| Spiking NN | No | No | No | No | Yes | No | Yes |
| Capacity | 2M | 1M | 4K | 16K | 64K | 128K | 256M |
| Cell area [$F^2$] | — | 59 | 124 | 968 | 74 | ~256 | — |
| Normalized die area | 12× | — | — | 30× | 1× | 11× | ~17240× |
| Chip power [mW] | 15.8 | — | — | — | 1.52 | — | 204.4 |
| Activation precision | 1 bit | 3 bit | 1 bit | 7 bit | 1~8 bit | 8 bit | 1 bit |
| Power efficiency [TOPS/W] — MAC only | 20.7 | 16.95 | 55.8 | 28.1 | — | — | — |
| Power efficiency [TOPS/W] — MAC +Activation | — | — | — | — | 16.9 | 3.125 | 0.4 |
| On-chip Activation Function Integration | No | No | No | No | Yes (tanh) | Yes (relu) | Yes (relu) |
| FoM[*] | — | 0.86 | 0.45 | 0.20 | 1.83 | 0.098 | — |

FoM = energy efficiency maximum activation precision/cell area.

*B. Yan, et al., VLSI Symp. 2019*

**Duke**