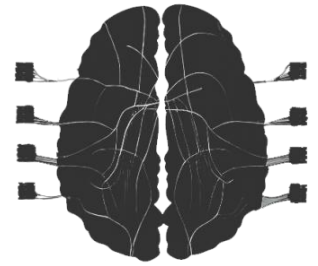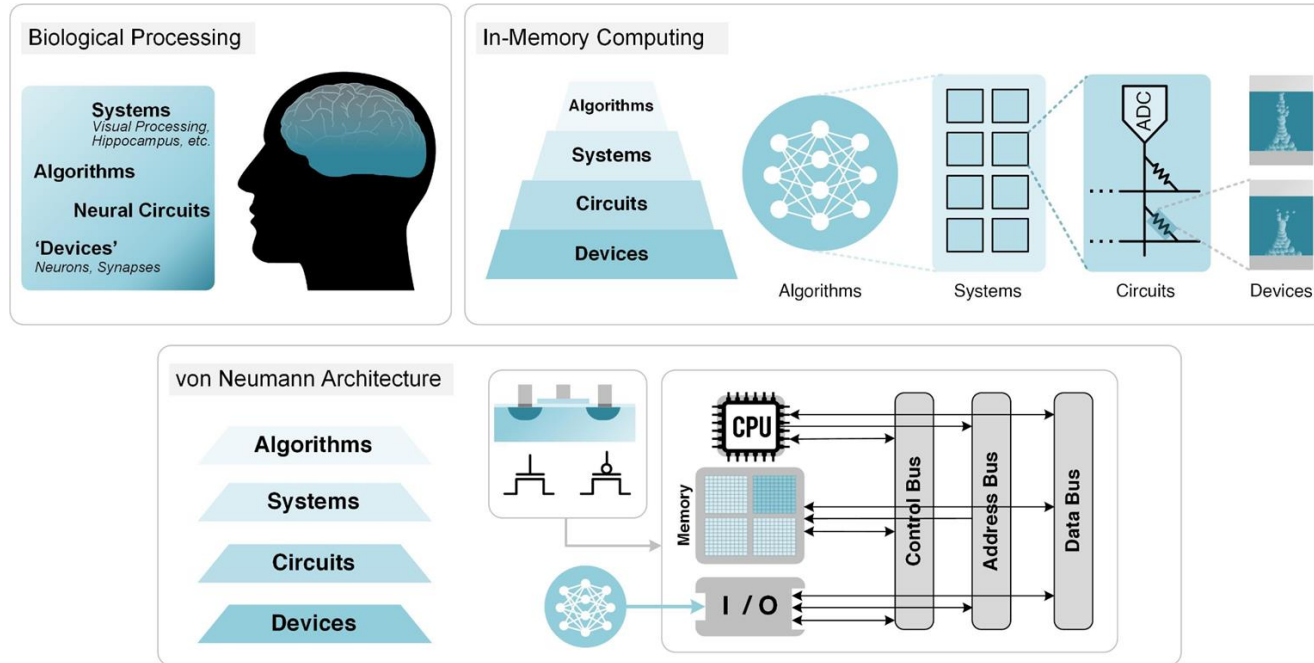# Open-Source Neuromorphic Computing

Jason K. Eshraghian
Assistant Professor, ECE, UC Santa Cruz

13 December 2023

# UCSC Neuromorphic Computing Group
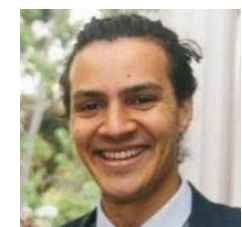


Lab Logo Generated by Stable Diffusion

**Biological Processing**

Systems
*Visual Processing, Hippocampus, etc.*

Algorithms

Neural Circuits

'Devices'
*Neurons, Synapses*

**In-Memory Computing**

Algorithms
Systems
Circuits
Devices

Algorithms · Systems · Circuits · Devices

**von Neumann Architecture**

Algorithms
Systems
Circuits
Devices

CPU · Control Bus · Address Bus · Data Bus · Memory · I / O

## Ph.D. Students
Binh Nguyen
Ruijie Zhu
Juan Lu
Coen Arrow
Assel Kembay

Power bill of training AI in data centers: **>$1 million**

Power budget of brains: **~12 watts**

At the UCSC NCG, our goal is to bridge this gap by applying principles from neuroscience to build better technology.
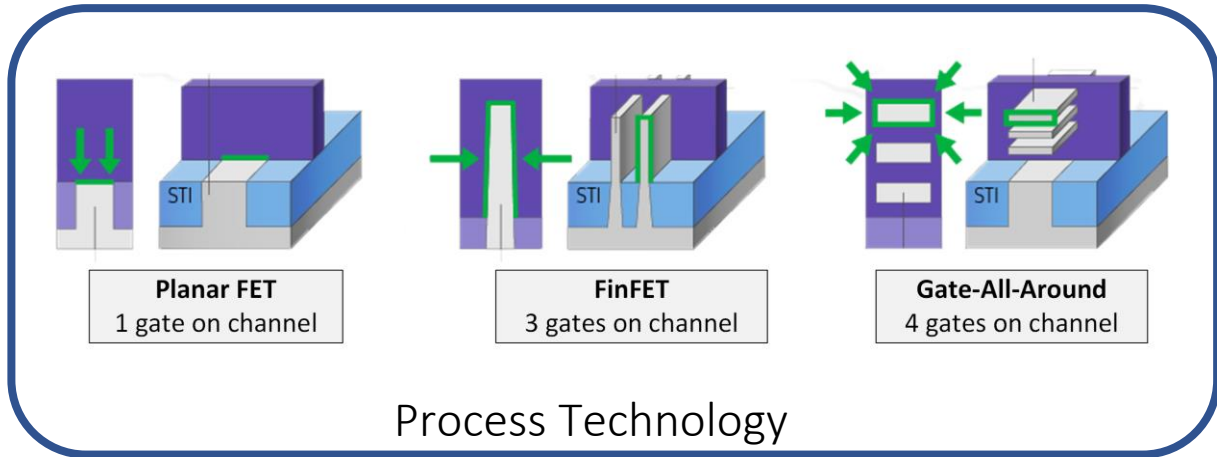
## Undergraduate Students
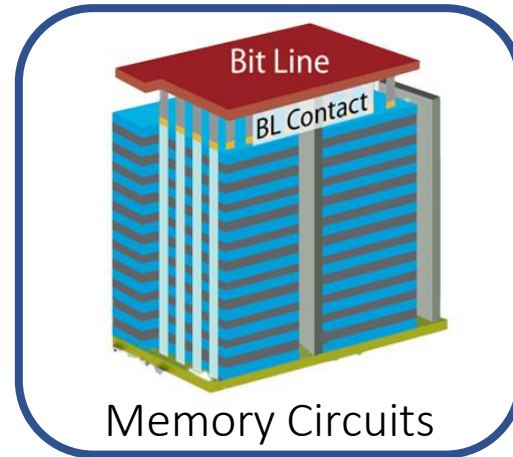Farhad Modaresi
Sreyes Venkatesh
Skye Gunasekaran
Hannah Cohen-Sandler

Ruhai Lin
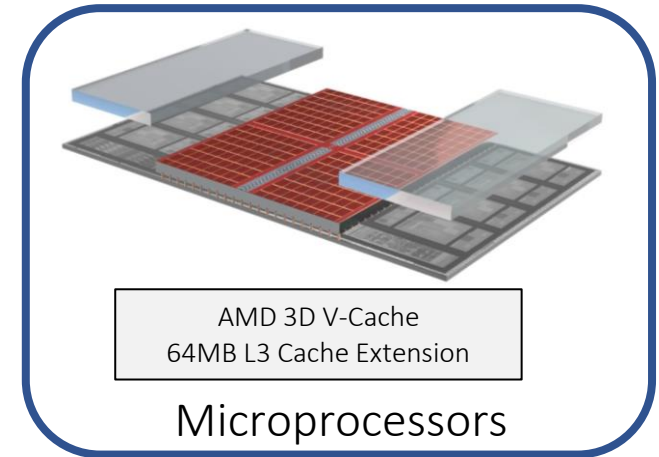Dylan Louie
Sahil Konjarla

# Circuits from 1D→ 2D → 3D → ...4D?



Process Technology

Samsung Newsroom, 2019.

Planar FET
1 gate on channel

FinFET
3 gates on channel

Gate-All-Around
4 gates on channel



Memory Circuits

Western Digital, 2018.

Bit Line

BL Contact



Microprocessors
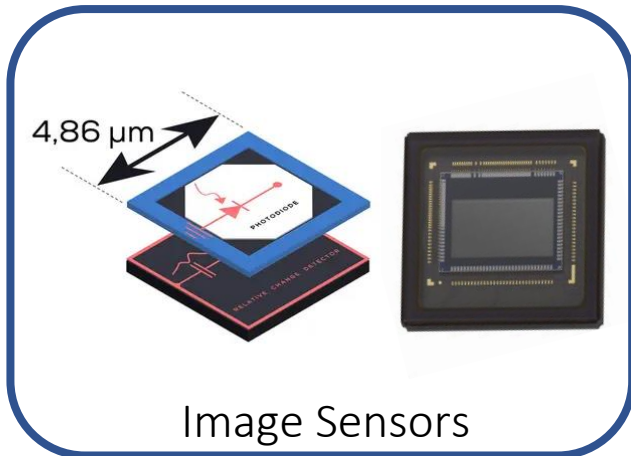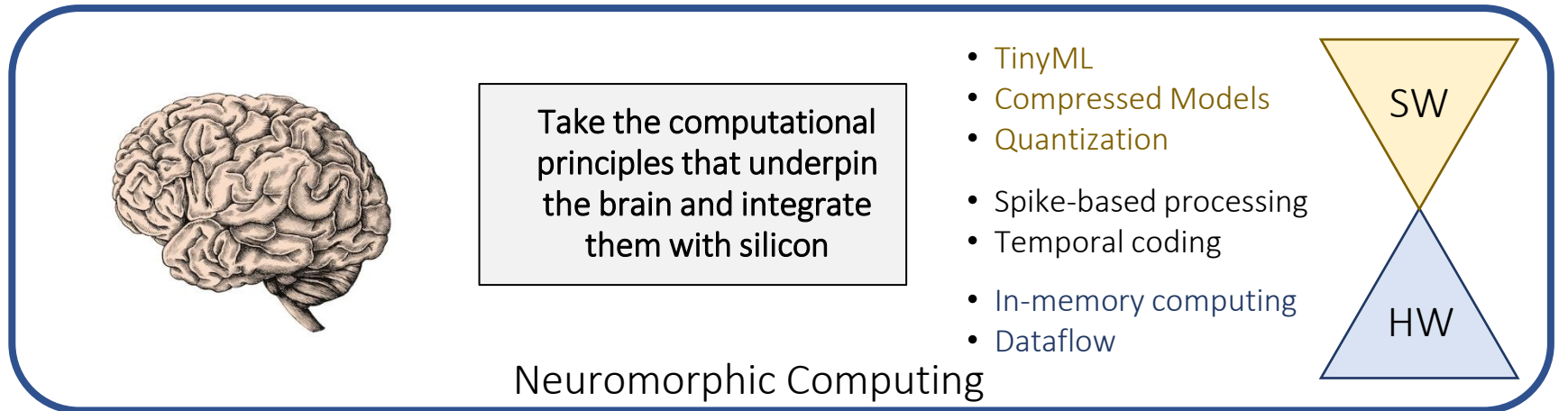
AMD, ISSCC 2022.

AMD 3D V-Cache
64MB L3 Cache Extension



Image Sensors

Sony & Prophesee, 2021.

4,86 μm



Neuromorphic Computing

Take the computational principles that underpin the brain and integrate them with silicon

- TinyML
- Compressed Models
- Quantization

- Spike-based processing
- Temporal coding

- In-memory computing
- Dataflow

SW
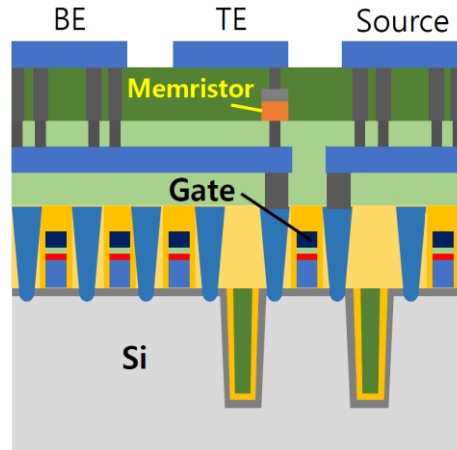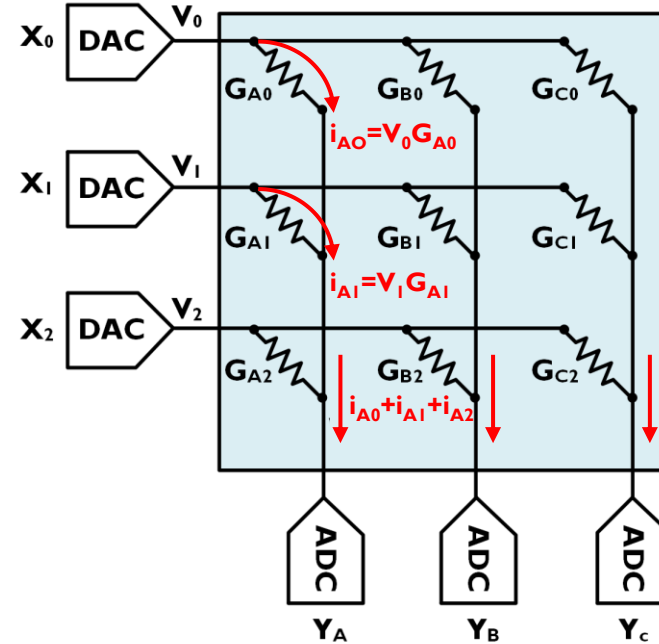
HW

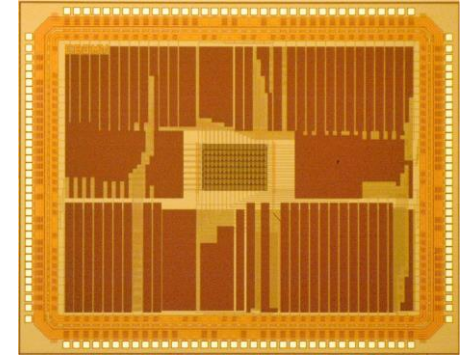# In-Memory Processing Using RRAM Crossbars



Valence Change Mechanism

BEOL metal-insulator-metal structure

Bit-line current summation

Chip micrograph: 65-nm mixed-signal 16x16 RRAM Crossbar

# In-Memory Processing Using RRAM Crossbars

Electrode

Oxide layer #2

Oxide layer #1

Electrode

### Input Data

### Neuron Weights

### Output Equations

$$[x_0 \quad x_1 \quad \cdots \quad x_n] * \begin{bmatrix} a_0 & b_0 & c_0 \\ a_1 & b_1 & c_1 \\ \vdots & \vdots & \vdots \\ a_n & b_n & c_n \end{bmatrix} = \begin{bmatrix} y_a = x_0 a_0 + x_1 a_1 + x_2 a_2 \\ y_b = x_0 b_0 + x_1 b_1 + x_2 b_2 \\ y_c = x_0 c_0 + x_1 c_1 + x_2 c_2 \end{bmatrix}$$

$G_{A2}$ $G_{B2}$ $G_{C2}$
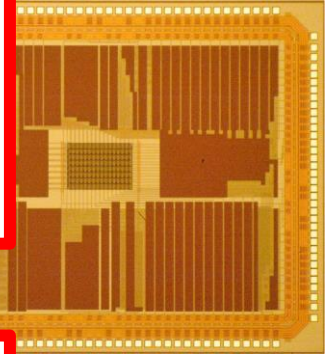
$i_{A0} + i_{A1} + i_{A2}$

### RRAM Conductance
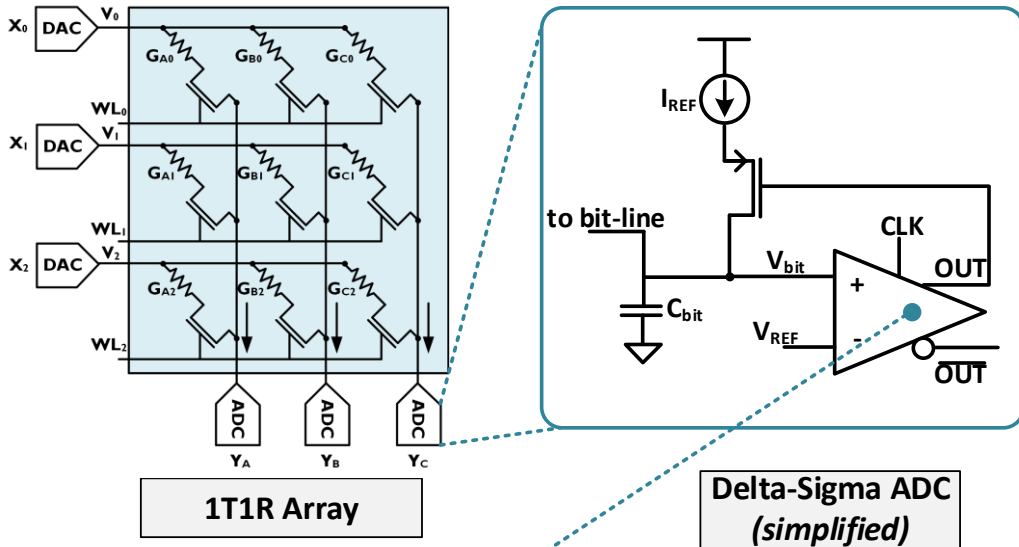
### Column Current

### Voltage Input

$$[V_0 \quad V_1 \quad \cdots \quad V_n] * \begin{bmatrix} G_{A0} & G_{B0} & G_{C0} \\ G_{A1} & G_{B1} & G_{C1} \\ \vdots & \vdots & \vdots \\ G_{An} & G_{Bn} & G_{Cn} \end{bmatrix} = \begin{bmatrix} I_A = V_0 G_{A0} + V_1 G_{A1} + V_2 G_{A2} \\ I_B = V_0 G_{B0} + V_1 G_{B1} + V_2 G_{B2} \\ I_C = V_0 G_{C0} + V_1 G_{C1} + V_2 G_{C2} \end{bmatrix}$$
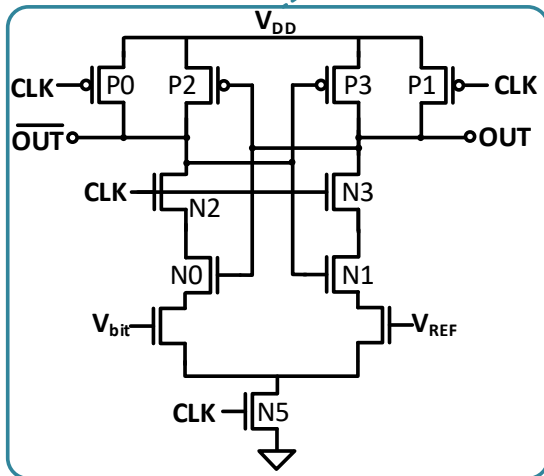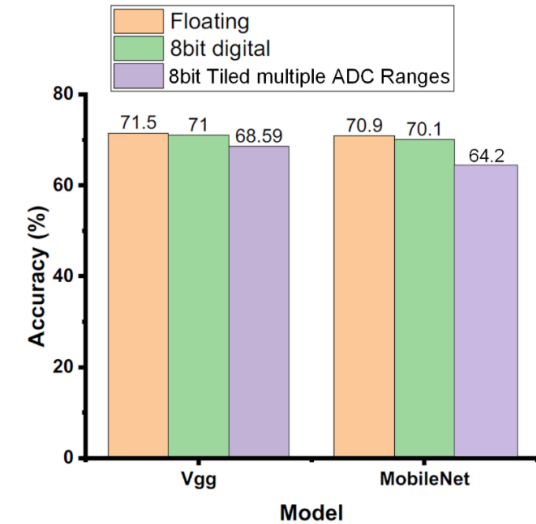
Valen
Me

ip micrograph:
mixed-signal 16x16
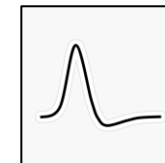RAM Crossbar

# In-Memory Processing Using RRAM Crossbars



**1T1R Array**



**Delta-Sigma ADC (simplified)**



**Pre-Charged Sense-Amplifier**

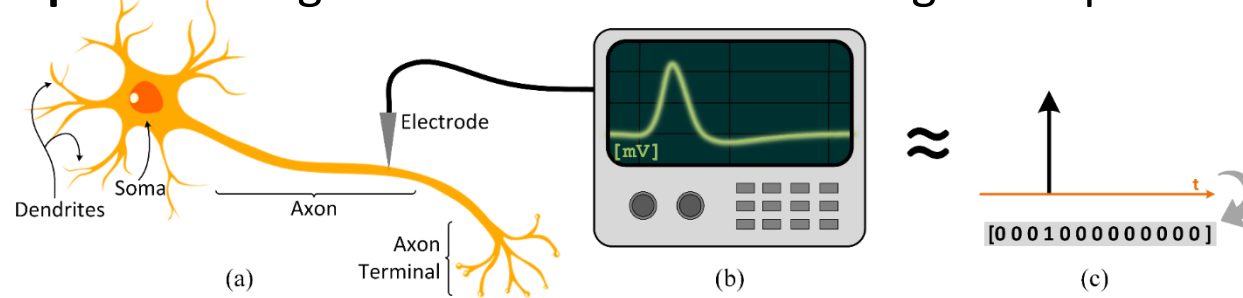| | Our Work | Mythic |
|---|---|---|
| **Technology** | 65 nm | 40 nm |
| **CLK Freq (MHz)** | 40 - ADC  100 - RRAM | 170 |
| **Power (W)** | 1.89 | 5 |
| **TOPS/W** | 5.90 | 4 |
| **Quantization** | Reconfigurable | 8-bit |



Top-5 ImageNet classification accuracy

## Challenges

- 1.5-3x improvements are insufficient to justify startup costs for a new process
- Scaling challenges: advanced processes are not optimized for analog operation: thermal noise
  - ADCs still <u>dominate latency</u>, and <u>~50% energy consumption</u>
    - Removal of ADCs requires algorithm-level optimization

spike-based processing?

# Toward Biological Networks: The Three S's

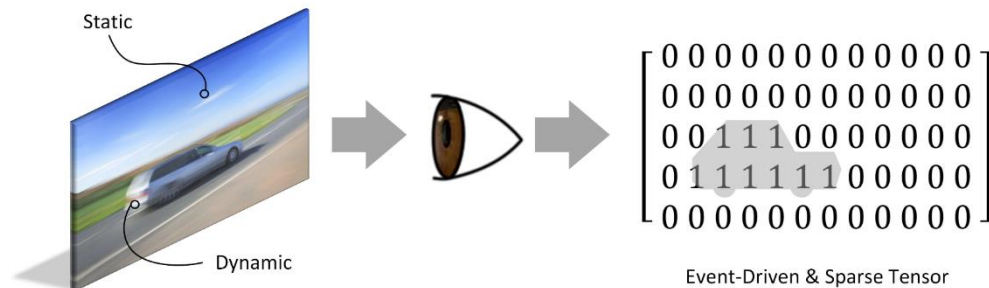**Spikes:** Biological neurons interact via single-bit spikes



**Sparsity:** Biological neurons spend most of their time at rest, setting most activations to *zero* at any given time

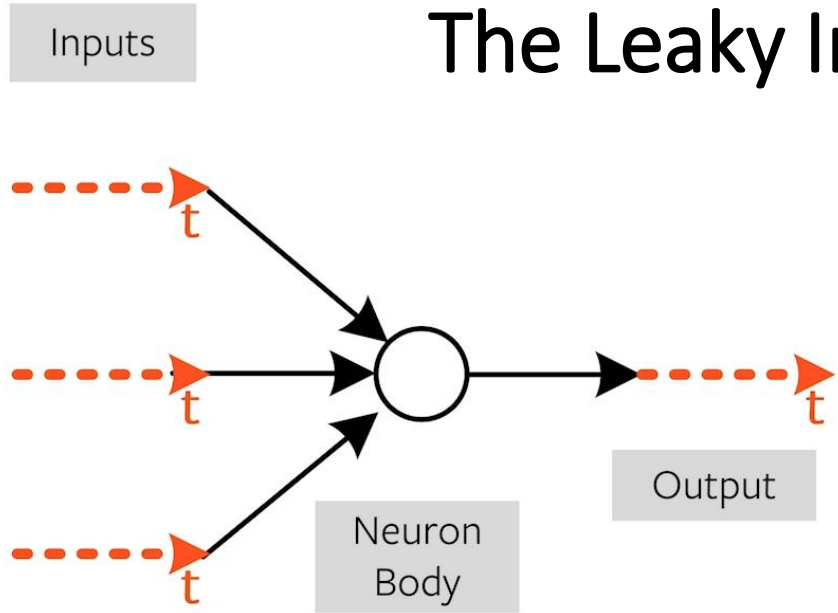$$[0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 7, 0, 0, 0, 0, 0, 0, 0, 0, 0, 5]$$

*"7 at position 10; 5 at position 20"*

**Static Suppression (aka Event-driven Processing):** The sensory periphery only processes information when there is *new* information to process
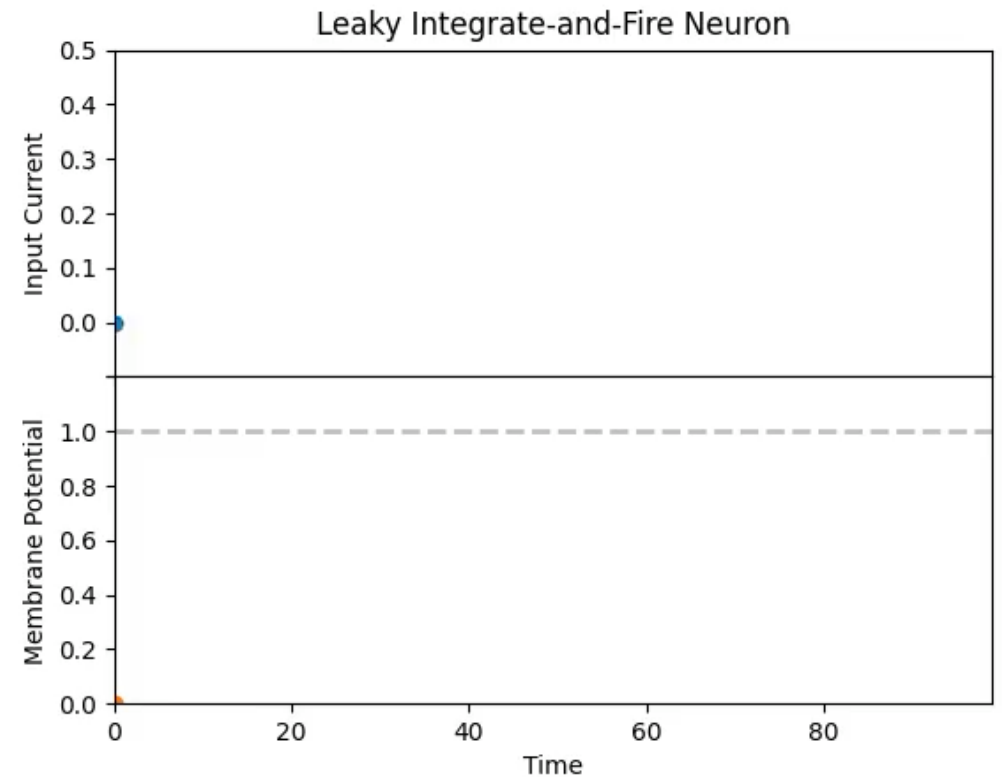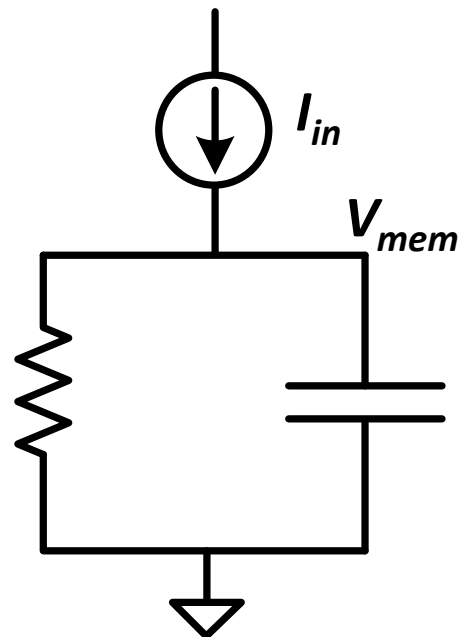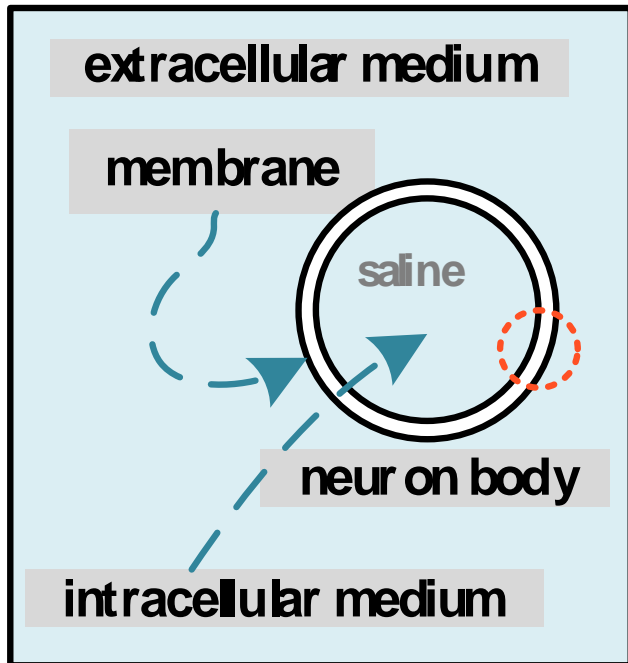


Event-Driven & Sparse Tensor

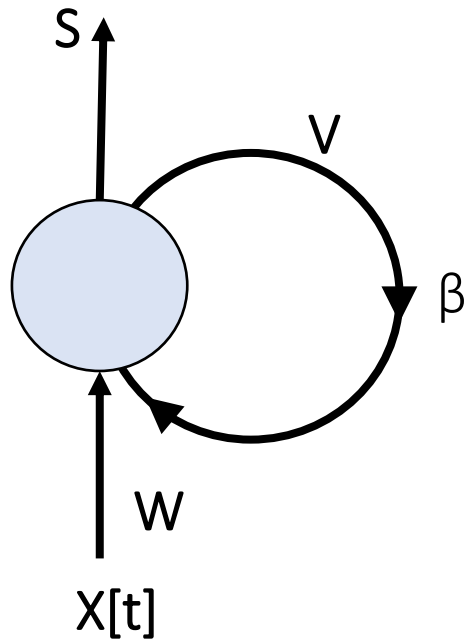# The Leaky Integrate-and-Fire Neuron

Inputs

Output

Neuron Body

- Bilipid thin-film membrane surrounded by ions: **capacitive**
- Ion-leakage/transfer: **resistive**
- The leaky integrate-and-fire neuron is just a 1st-order low-pass filter, i.e., an **RC circuit**

**The neuroscientists stole from electrical engineers so it's time to steal back from them**
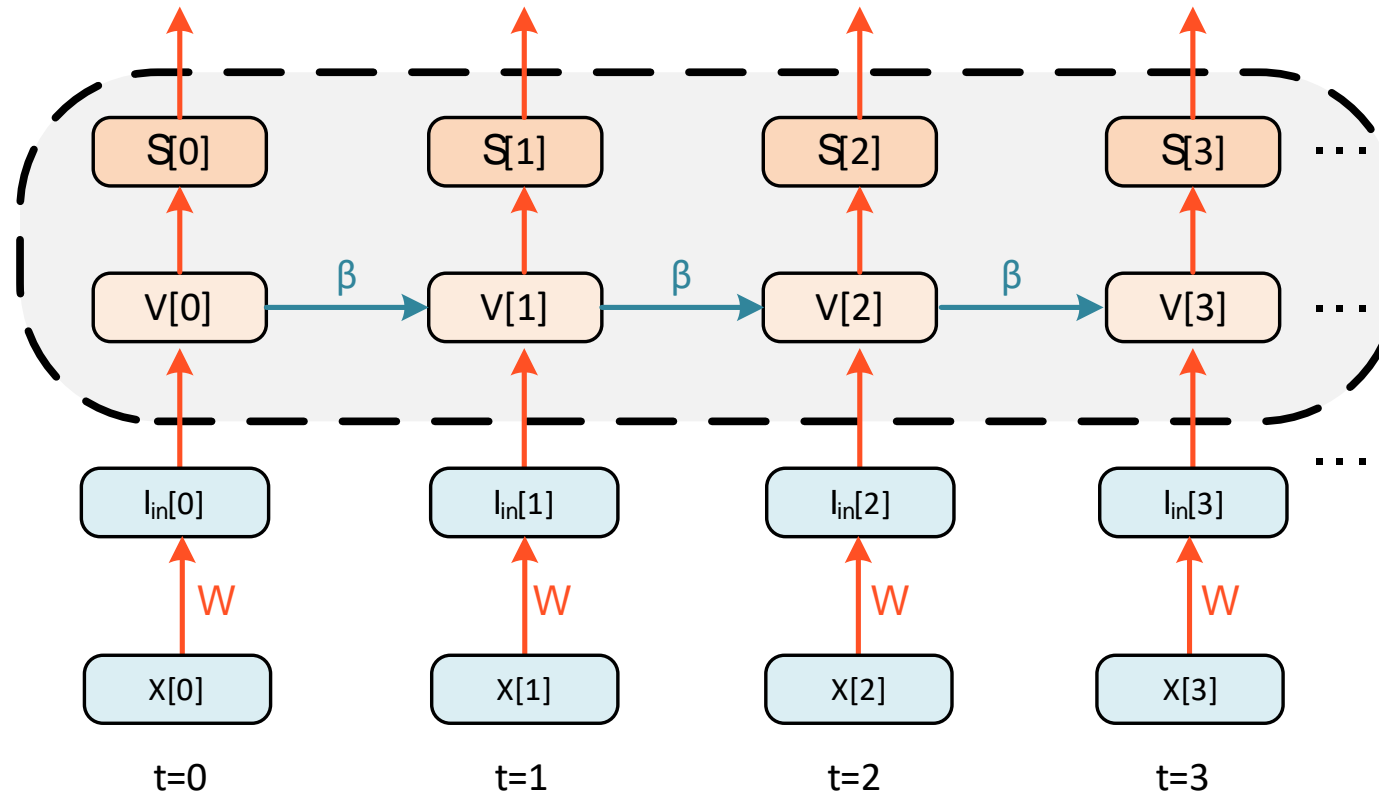
extracellular medium

membrane

saline

neuron body

intracellular medium

$I_{in}$

$V_{mem}$

Leaky Integrate-and-Fire Neuron

# A Recurrent Representation

## Unrolled Computational Graph



**Spiking Dynamics**

$$S\,[t+1] = H(V[t+1] - V_{thr})$$

**Membrane Potential Dynamics**
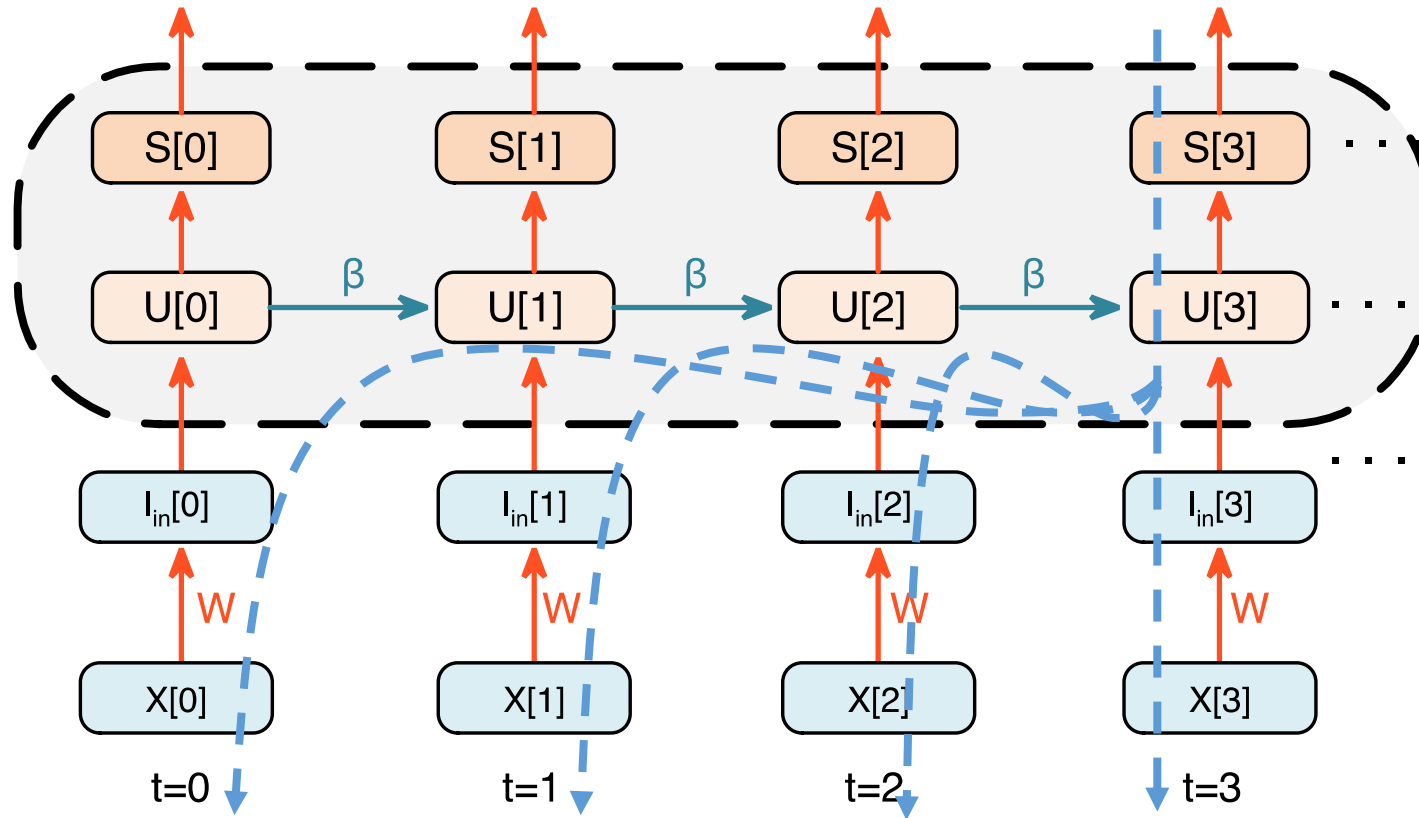
$$V\,[t+1] = \beta V[t] + (1-\beta)I_{in}[t]$$

The leaky integrate-and-fire neuron is now compatible with all the tricks and hacks that go with training deep learning models.

*Note:*
$(1-\beta)$ is removed, and
$I_{in}[t] = WX[t]$

# How do we train models that change over time?

Error Backpropagation ...Through Time

Unrolled Computational Graph



**Spatial Credit Assignment**
How does a loss assign 'blame' to a weight that is spatially far?

**Temporal Credit Assignment**
All states throughout history must be stored to run the BPTT algorithm
i.e., the entire graph must be stored
Memory complexity: **O(nT)**

# snnTorch

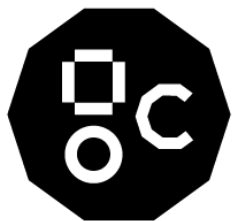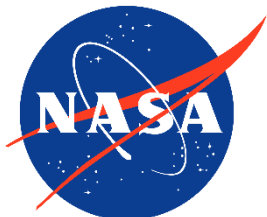Gradient-based Learning with Spiking Neural Networks

Python package for gradient-based optimization of SNNs

real-time online learning

seamless integration with PyTorch

CUDA + IPU accelerated

neuromorphic HW compatible

github.com/jeshraghian/snntorch

| Tutorial | Title | Colab Link |
|---|---|---|
| Tutorial 1 | Spike Encoding with snnTorch | CO Open in Colab |
| Tutorial 2 | The Leaky Integrate and Fire Neuron | CO Open in Colab |
| Tutorial 3 | A Feedforward Spiking Neural Network | CO Open in Colab |
| Tutorial 4 | 2nd Order Spiking Neuron Models (Optional) | CO Open in Colab |
| Tutorial 5 | Training Spiking Neural Networks with snnTorch | CO Open in Colab |
| Tutorial 6 | Surrogate Gradient Descent in a Convolutional SNN | CO Open in Colab |
| Tutorial 7 | Neuromorphic Datasets with Tonic + snnTorch | CO Open in Colab |

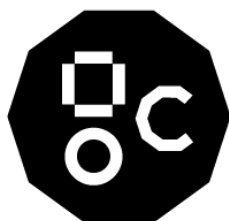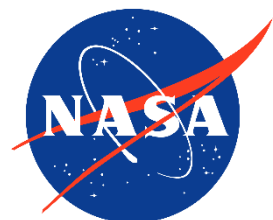| Advanced Tutorials | Colab Link |
|---|---|
| Population Coding | CO Open in Colab |
| Regression: Part I - Membrane Potential Learning with LIF Neurons | CO Open in Colab |
| Regression: Part II - Regression-based Classification with Recurrent LIF Neurons | CO Open in Colab |
| Accelerating snnTorch on IPUs | — |

Python package for gradient-based optimization of SNNs

real-time online learning

seamless integration with PyTorch

CUDA + IPU accelerated

neuromorphic HW compatible

github.com/jeshraghian/snntorch

# RRAM-based SNN Processing



**ANN via Current Summation**

(a)

Communicating noise-prone analog signals down long lengths of wire: noise-prone

**SNN via Integration**

*Active Synapse (RRAM Cell)*

Communicating noise-robust digital signals down long lengths of wire: noise-tolerant

An RRAM Approach to Spiking Neuron Dynamics

- Charge-based integration on the bit-line capacitance
- State decay using heterogeneous RRAM time constants

*J.K. Eshraghian et al., "Memristor-based Binarized SNNs" IEEE Nanotech. Mag. 2022*

# RRAM-based SNN Processing



right hand counter clockwise

Communicating noise-prone analog signals down long lengths of wire: noise-prone



Decay line

**SNN via Integration**

$X_1(t)$  I

0

$X_2(t)$

Active Synapse (RRAM Cell)

$Y_A(t)$  $Y_B(t)$  $Y_C(t)$

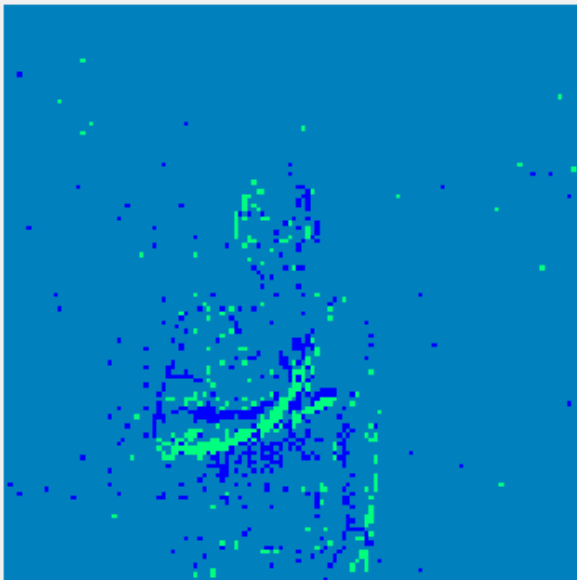Communicating noise-robust digital signals down long lengths of wire: noise-tolerant

An RRAM Approach to Spiking Neuron Dynamics

- Charge-based integration on the bit-line capacitance
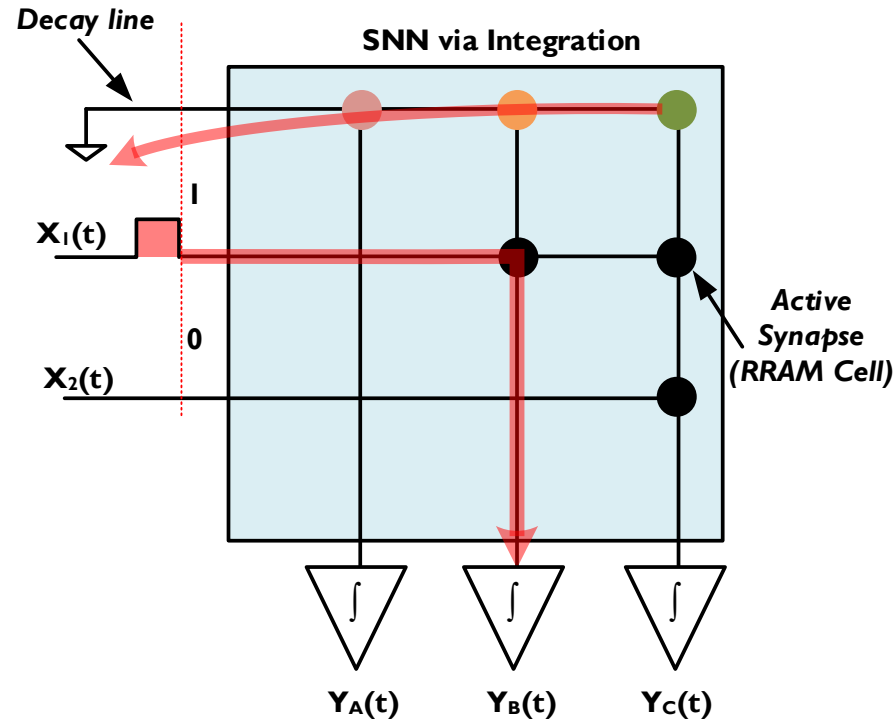- State decay using heterogeneous RRAM time constants
- Pre-charged sense amplifier for spiking dynamics (replace ADC)



$V_{bit}$  $V_{ref}$  $V_{in}$

t

$V_{out}$  spike!

t

*J.K. Eshraghian et al., "Memristor-based Binarized SNNs" IEEE Nanotech. Mag. 2022*

*Classical pain points of IMC RRAM are largely addressed by spikes*

right hand counter clockwise

SNN via Integration

An RRAM Approach to Spiking Neuron Dynamics

$X_0(t)$

$X_1(t)$

$X_2(t)$

0

I

0

*Active Synapse (RRAM Cell)*

Variability

ADC overhead

Bandlimited at scale

Endurance

$Y_A(t)$  $Y_B(t)$  $Y_C(t)$

- Charge-based integration on the bit-line capacitance
- State decay using heterogeneous RRAM time constants
- Pre-charged sense amplifier for spiking dynamics (replace ADC)
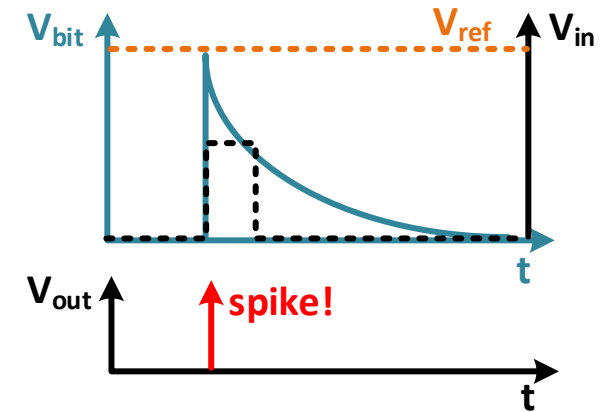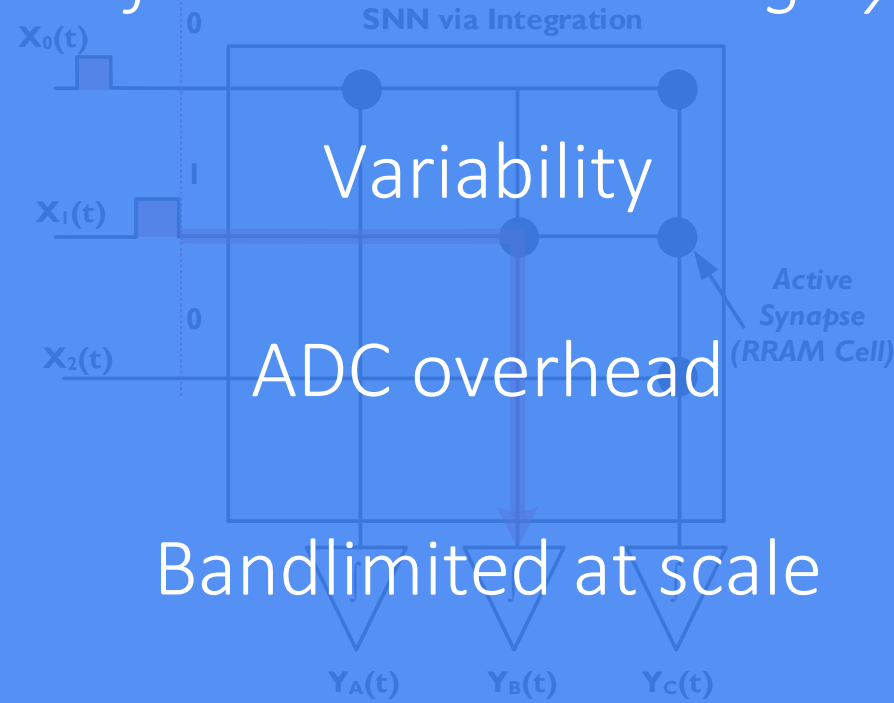
$V_{bit}$  $V_{ref}$  $V_{in}$

$V_{out}$  spike!  t

t

Communicating noise-prone analog signals down long lengths of wire: noise-prone

Communicating more robust digital signals down long lengths of wire: noise-tolerant

*J.K. Eshraghian et al., "Memristor-based Binarized SNNs" IEEE Nanotech. Mag. 2022*

# RRAM-based SNN Processing

*Classical pain points of IMC RRAM are largely addressed by spikes*

~~Variability~~ Heterogeneity

~~ADC overhead~~ Single-bit Spikes

~~Bandlimited at scale~~ Sparse data movement

~~Endurance~~ Weight updates only occur at spike times

J.K. Eshraghian et al., "Memristor-based Binarized SNNs" IEEE Nanotech. Mag. 2022
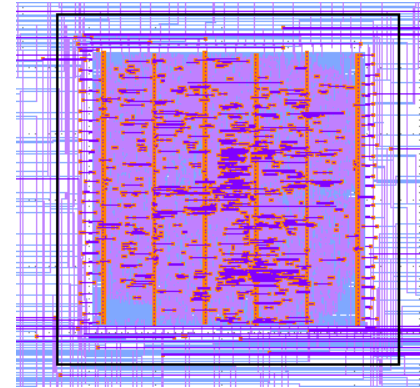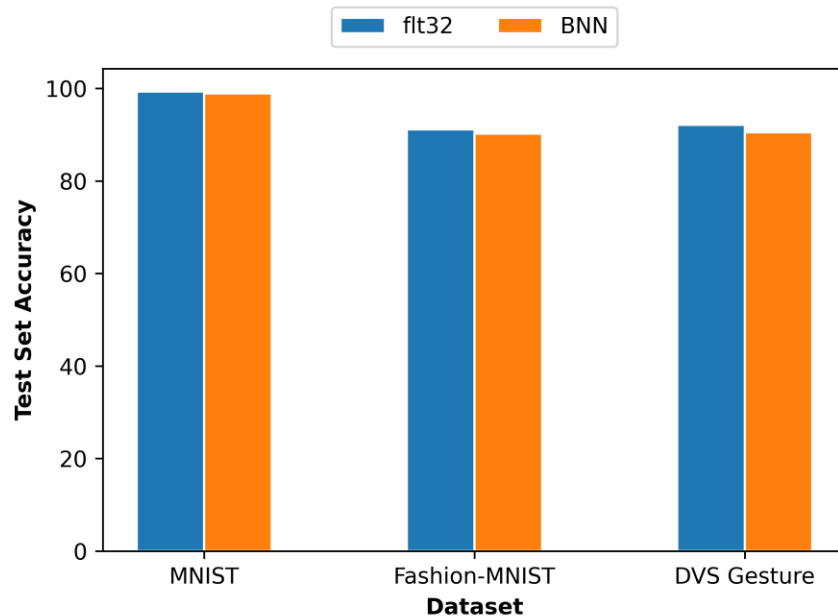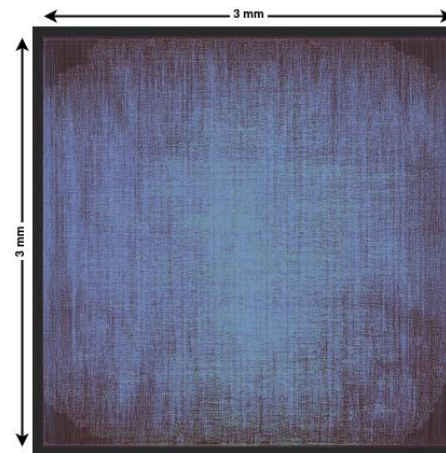
# SkyWater 130 Neuromorphic Accelerator

## Open-Source Neuromorphic IP

- Google-sponsored tape-outs using SkyWater 130nm process
- Deep learning success was in part due to open-source; let's port this hype train over to silicon
- 2x successful tape-outs + 1x tape-in…



MPW6 BSNN Streaming Accelerator

- Online event streaming accelerator
- HD Event Cameras: ~10 MHz spike rate
- Die area: 0.09 mm$^2$
- Average Power: 11 nW
- GPIO: 50 MHz peak, Core: 210 MHz



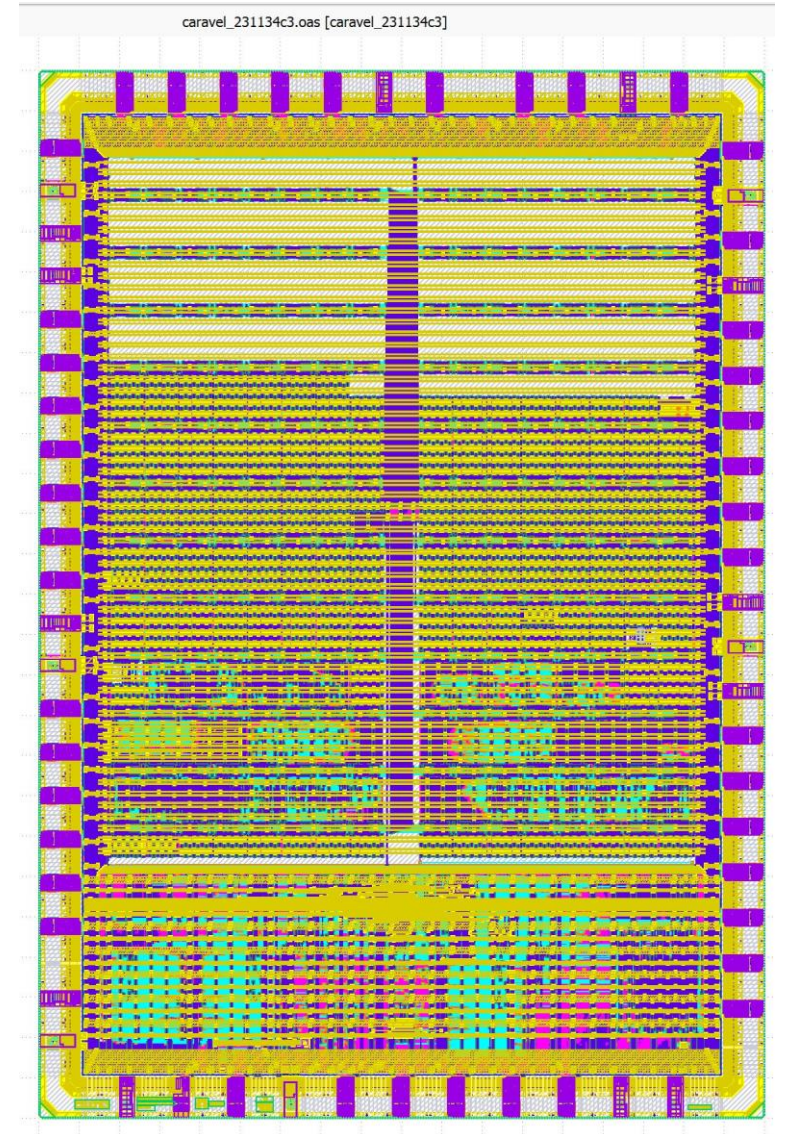MPW8 SNN Accelerator
Lead Designer: Farhad Modaresi

- 154 kB of dual-power 2 kB SRAM macros
- 9 mm$^2$ core, 21.89 mm$^2$ memory
- SRAM: 40 MHz, Core: 20 MHz
- Custom precision (up to 8-bits)

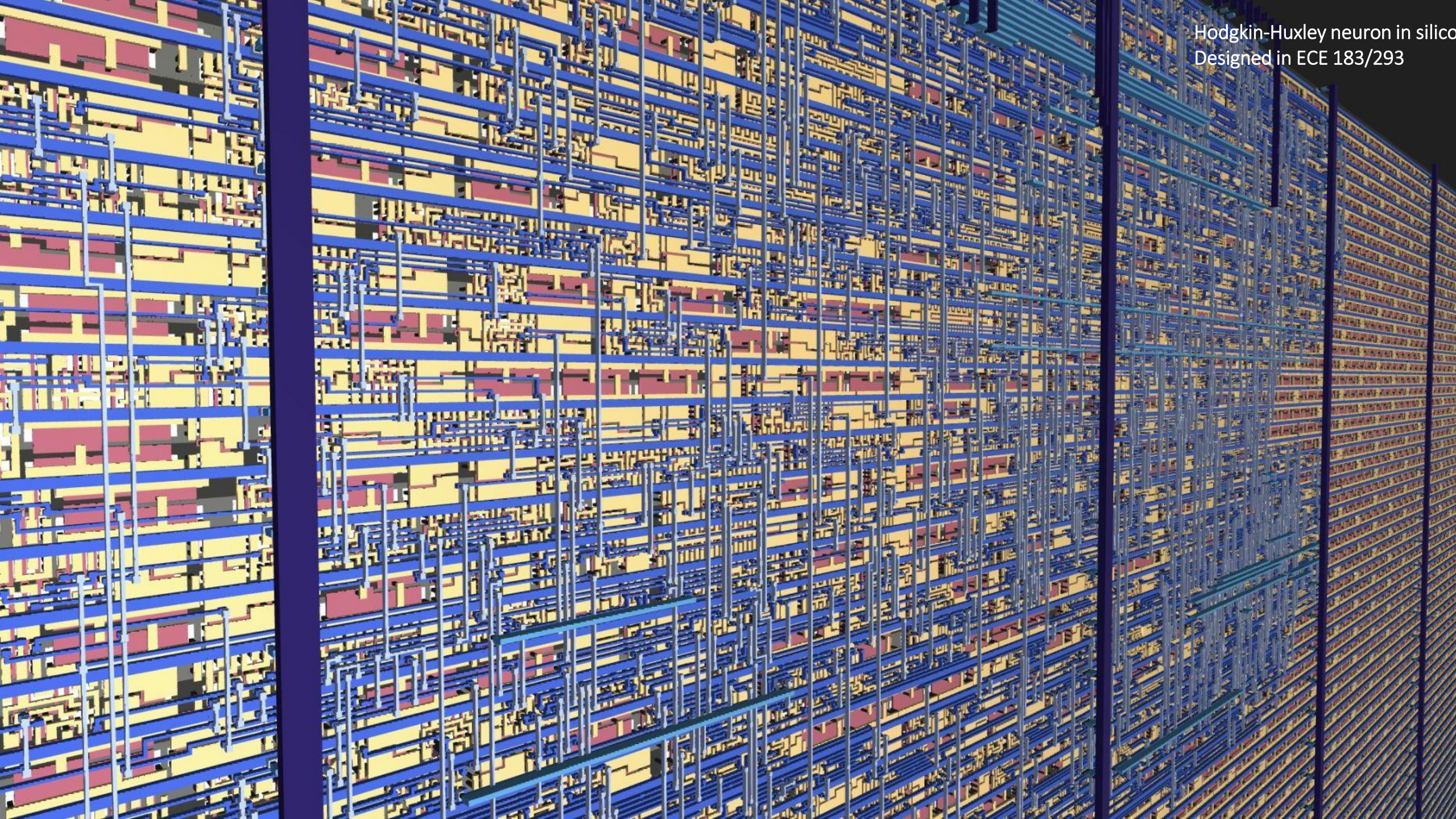# The next generation of neuromorphic engineers

31 first-time chip designers at UCSC taking the
"Brain-Inspired Machine Learning" class
defined their own design and submitted a chip.

It is currently being manufactured and will ship
back in April 2023.

www.tinytapeout.com


caravel_231134c3.oas [caravel_231134c3]

Hodgkin-Huxley neuron in silico
Designed in ECE 183/293

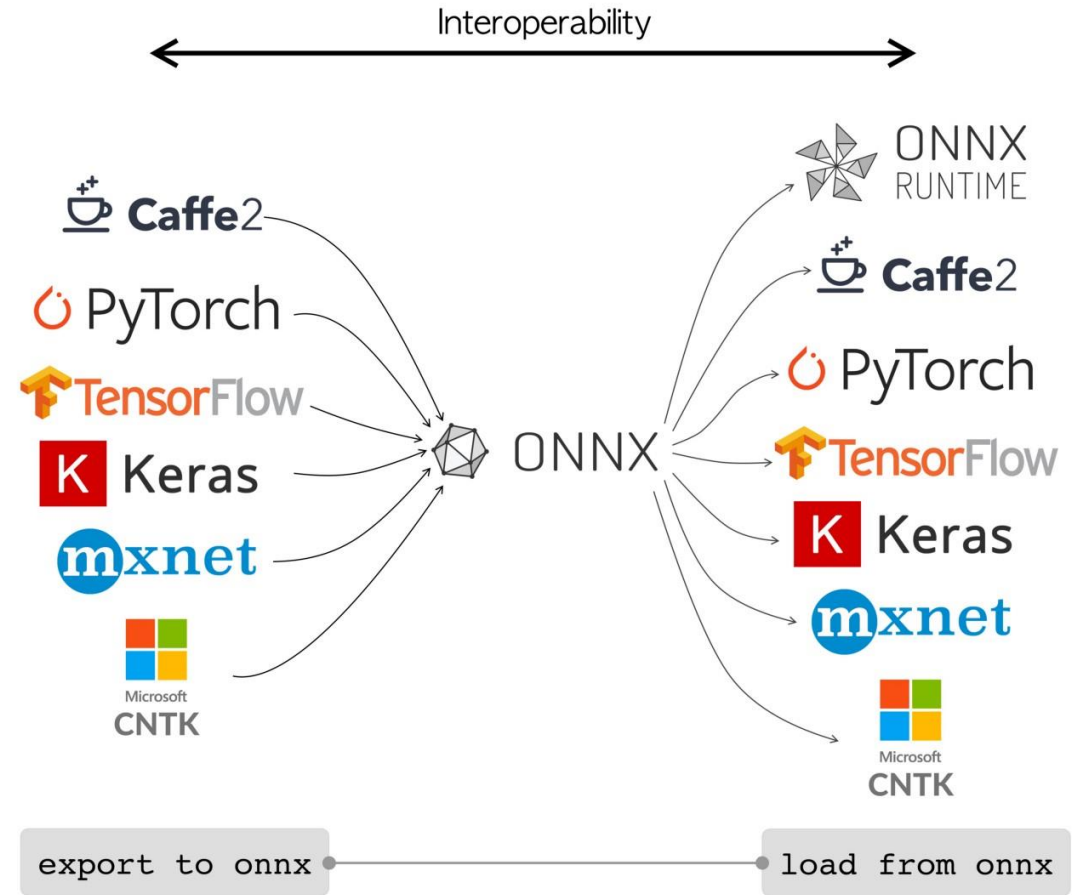How many of you have converted your garage to a semiconductor manufacturing plant?
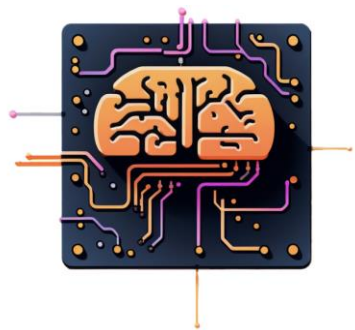
I didn't think so.

# How can we port dynamical, time-varying neural networks across different systems?

Solutions already exist for vanilla deep learning.

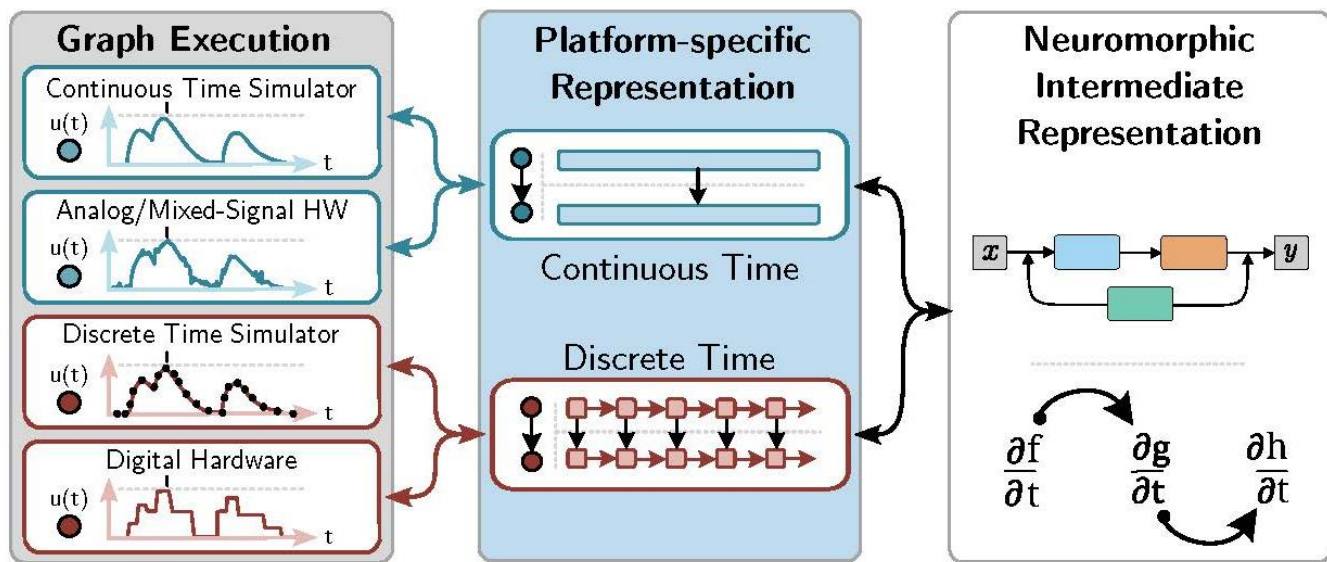Nothing but pain exists for spiking neural networks.

… until now.

# Neuromorphic Intermediate Representation

## A Unified Instruction Set for Interoperable Brain-Inspired Computing



- Provides a common standard to enable continuous-time, dynamical neural networks to interface with each other

- Allows the broader community to tap into commercial & exotic hardware with their software of choice – e.g., Intel Loihi

```
>> model.to_nir()
```

We wrote 1000s of lines of code so you only have to write 1.

Pederesen, Abreu, Eshraghian et al., arXiv, Nov. 2023.

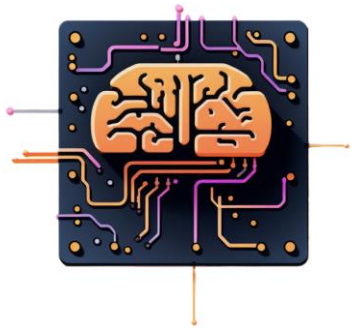# Neuromorphic Intermediate Representation

## A Unified Instruction Set for Interoperable Brain-Inspired Computing

- Provides a common standard to enable continuous-time, dynamical neural networks to interface with each other

- Allows the broader community to tap into commercial & exotic hardware with their software of choice – e.g., Intel Loihi

```
>> model.to_nir()
```

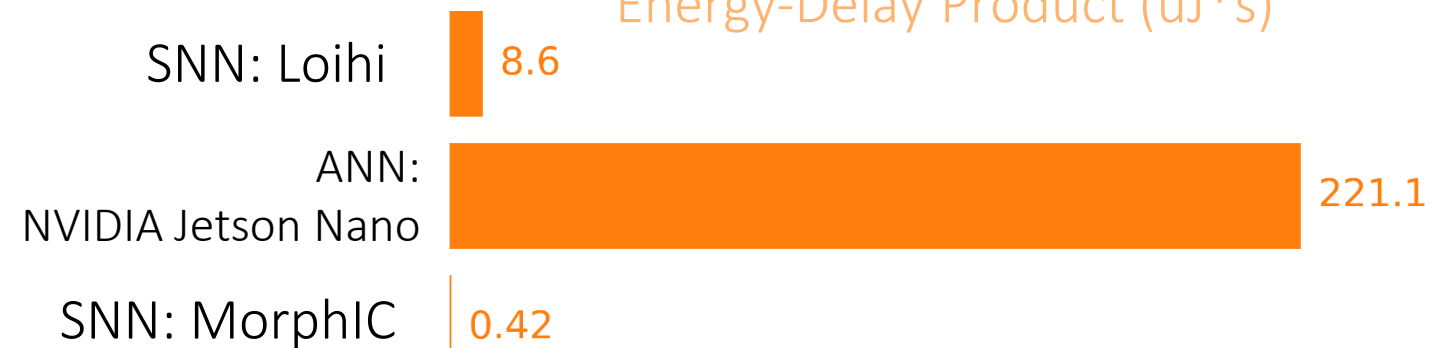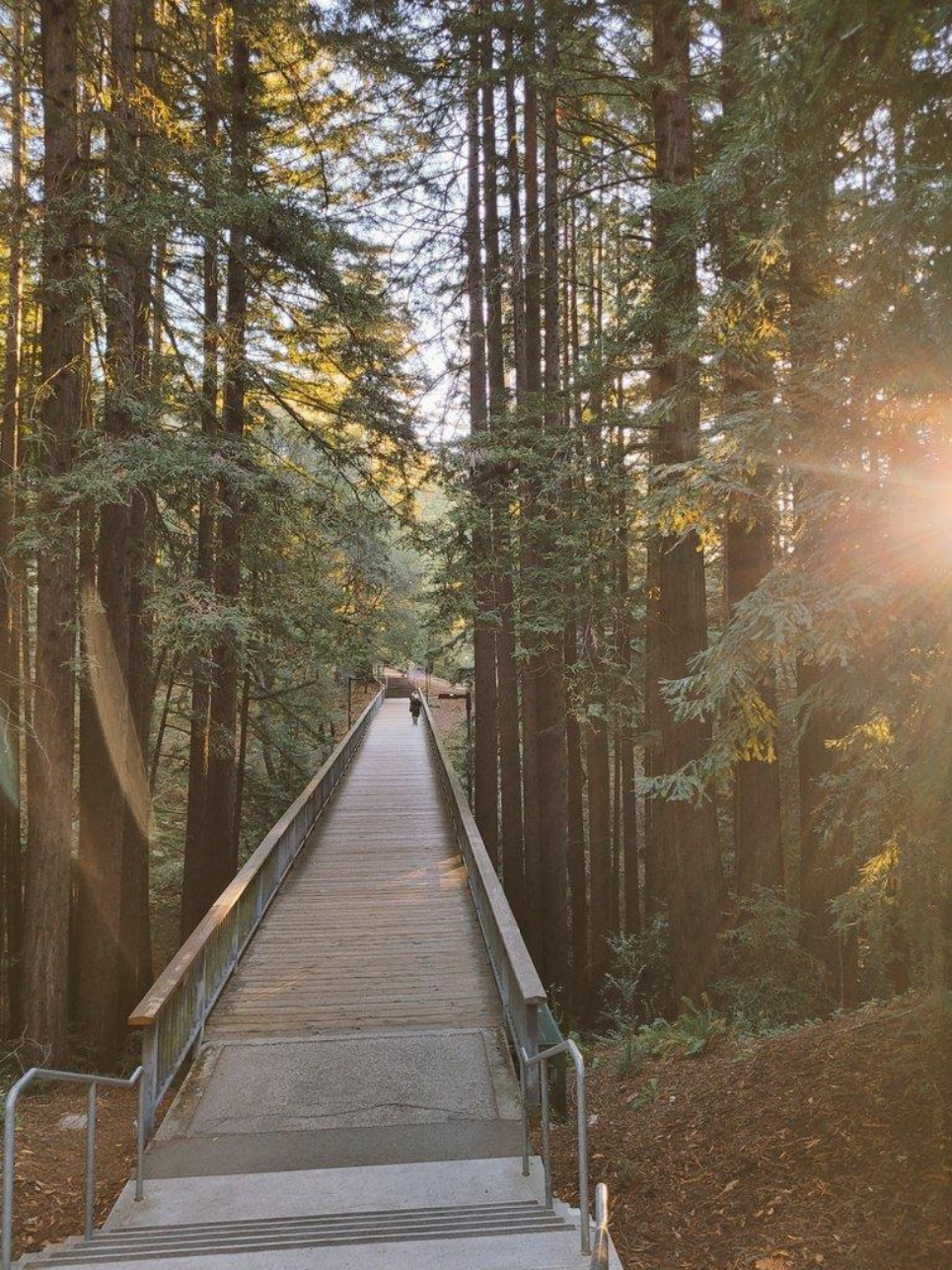We wrote 1000s of lines of code so you only have to write 1.

Server-scale neuromorphic HW: Loihi, Intel Labs

Edge neuromorphic hardware: SynSense, BrainChip

Pedersen, Abreu, Eshraghian et al., arXiv, Nov. 2023.

# SNN Evaluation – Multimodal Data



Accuracy

| | |
|---|---|
| SNN: Loihi | 96.0 |
| ANN: NVIDIA Jetson Nano | 95.4 |
| SNN: MorphIC | 89.4 |

Energy-Delay Product (uJ*s)

| | |
|---|---|
| SNN: Loihi | 8.6 |
| ANN: NVIDIA Jetson Nano | 221.1 |
| SNN: MorphIC | 0.42 |

*M. Rahimi-Azghadi, J.K. Eshraghian et al., IEEE Trans. Biomedical CAS, Dec. 2020*

# Open-Source Neuromorphic Computing

Jason K. Eshraghian
Assistant Professor, ECE, UC Santa Cruz

13 December 2023

UC SANTA CRUZ | Baskin Engineering