# NEUROMORPHIC COMPUTING

## Towards Brain-like Energy Efficiency

**Sandia National Laboratories**

*Exceptional service in the national interest*

Suma George Cardwell, PhD

*Principal Member of Technical Staff*

*Sandia National Laboratories, USA*

**Parallel Processing and Applied Mathematics (PPAM) 2024**
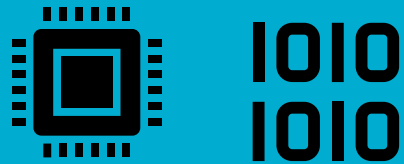
Ostrava, Czech Republic

September 9th, 2024

# MODERN COMPUTING LANDSCAPE

**DIGITAL COMPUTING**

CPU , GPU, FPGA

- Discrete-valued
- Deterministic
- High-Precision
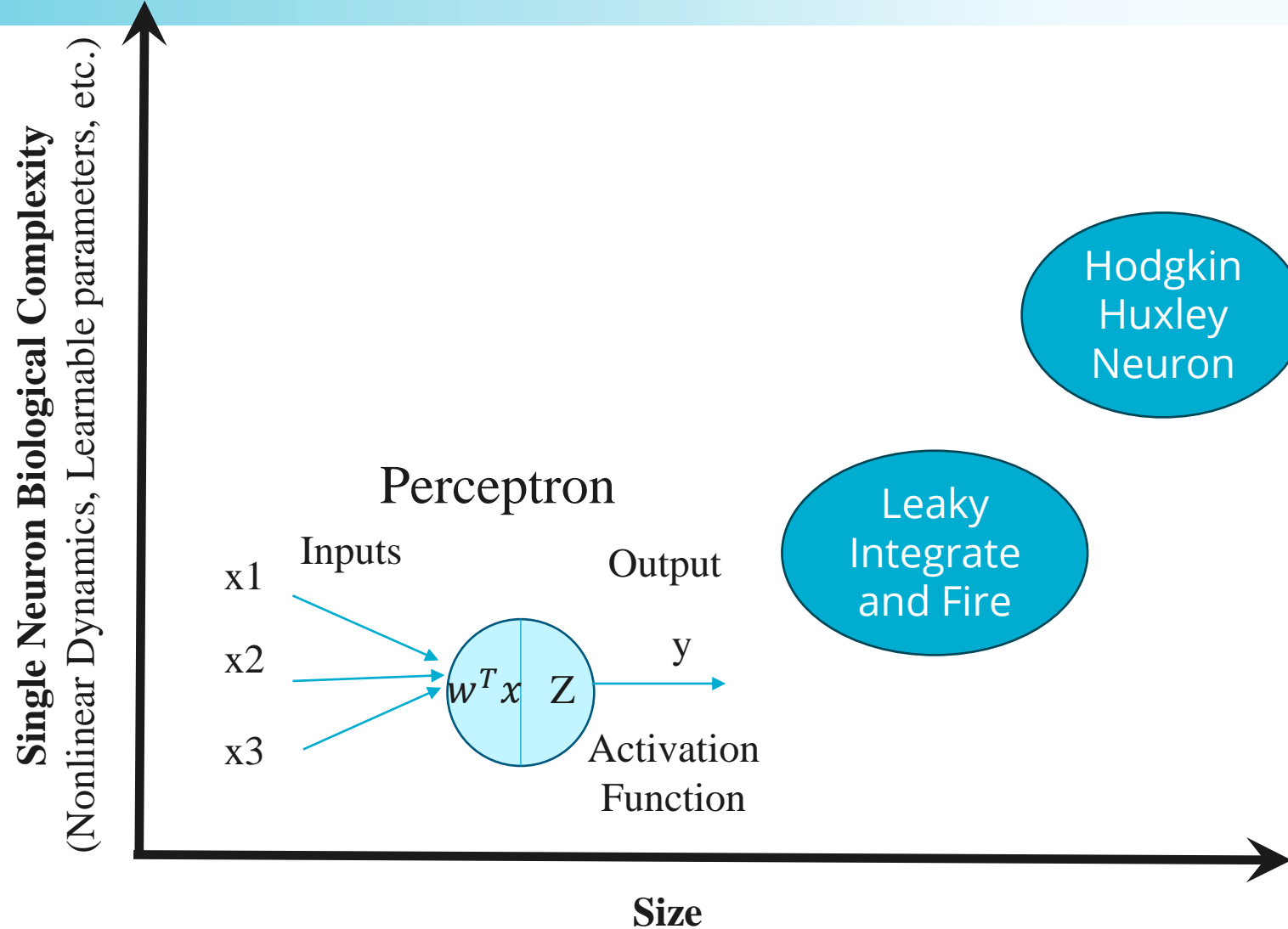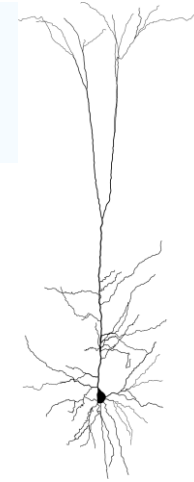- Simple(r) building blocks

**ANALOG COMPUTING**

Brain , FPAA

- Continuous-valued
- Stochastic
- Lower-Precision
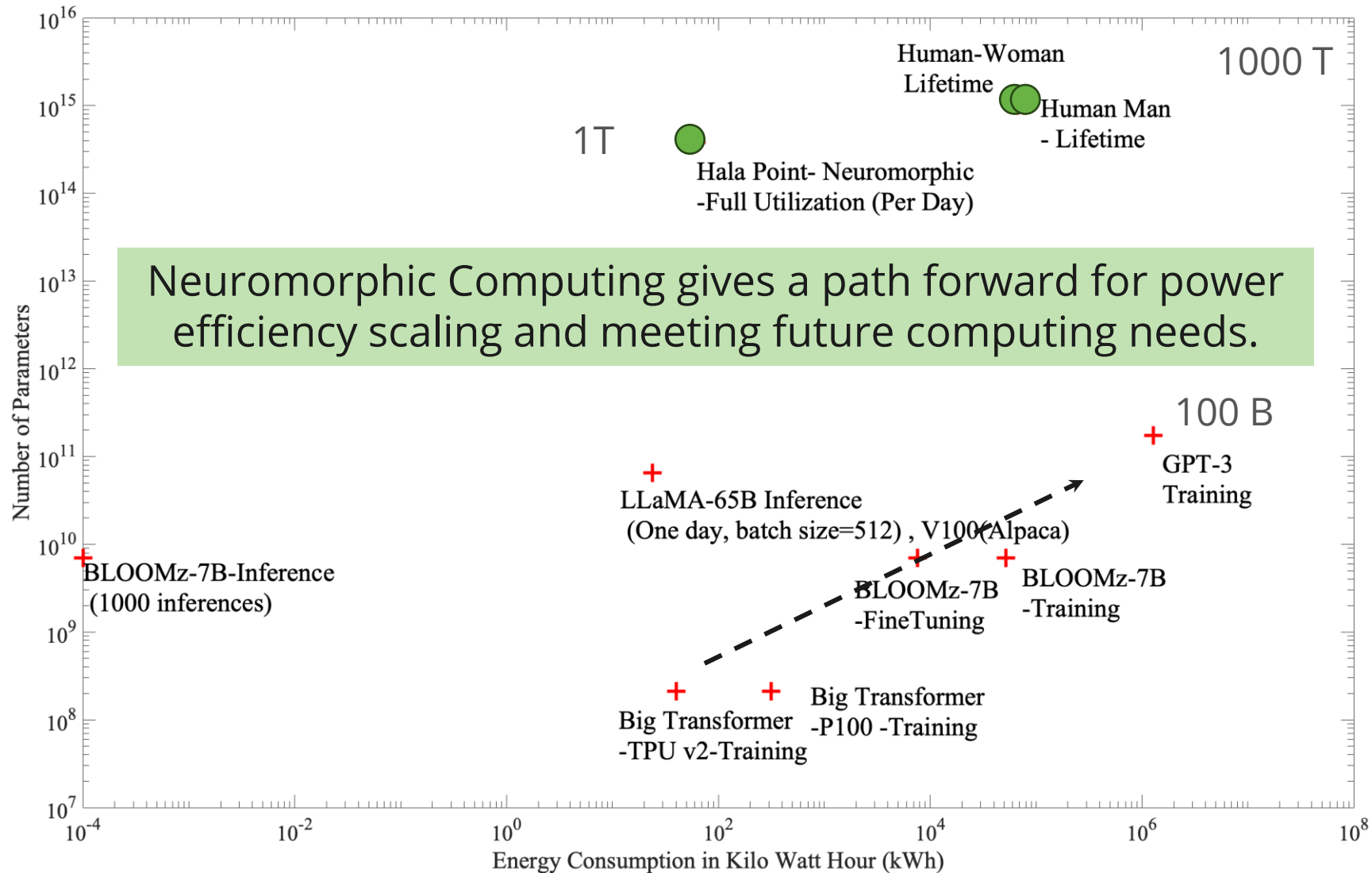- Complex building blocks (neurons)

Biological Neurons have rich dynamics

**Single Neuron Biological Complexity**
(Nonlinear Dynamics, Learnable parameters, etc.)

Hodgkin Huxley Neuron

Leaky Integrate and Fire

Perceptron

Inputs

Output

x1

x2

$w^T x$   Z   y

x3

Activation Function

**Size**

Neuromorphic Computing gives a path forward for power efficiency scaling and meeting future computing needs.
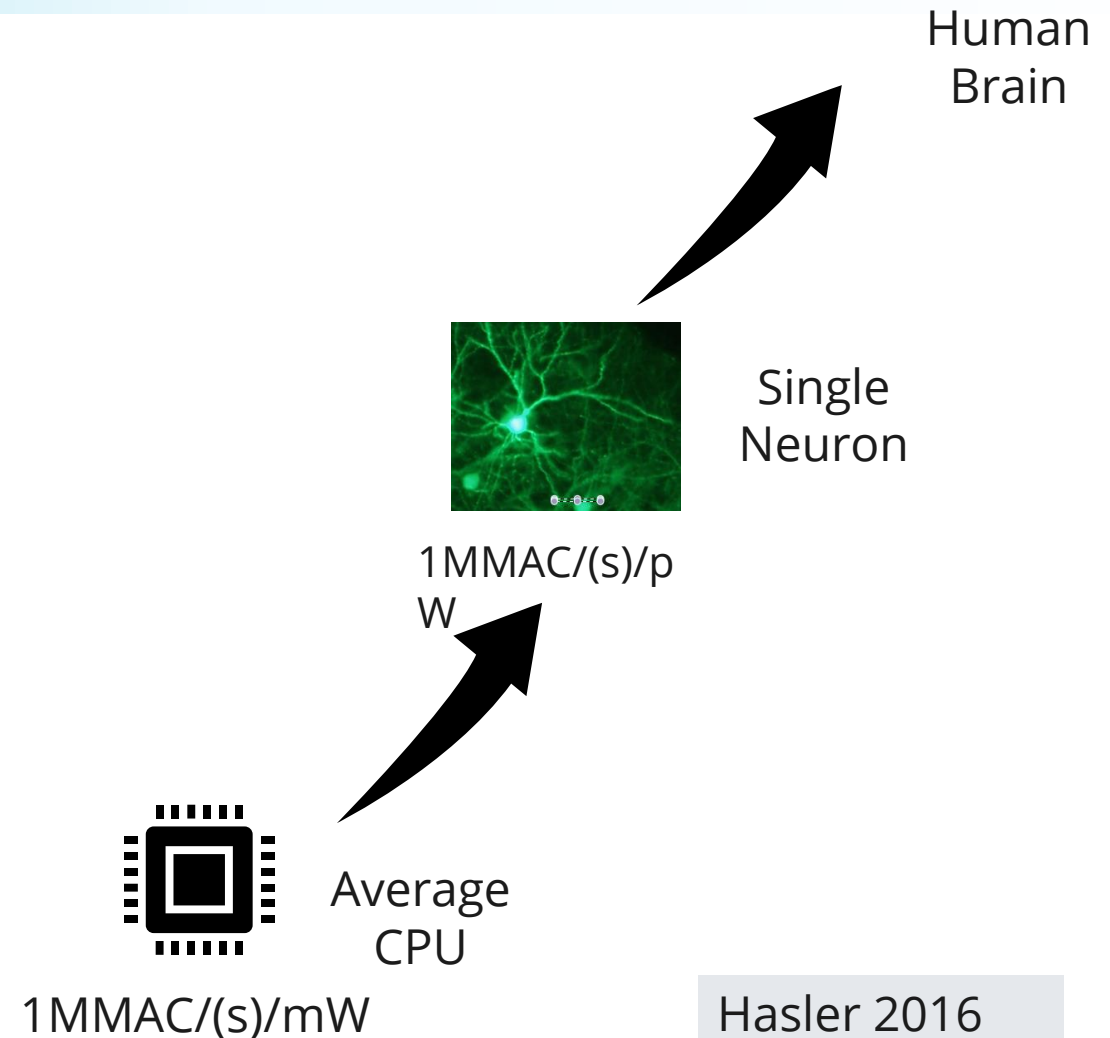
# NEUROMORPHIC COMPUTING

## TAKING INSPIRATION FROM BRAINS

**Functionality**

- Solves ill-structured problems with little training
- Online learning
- Transfer learning
- Continual learning

**Attributes**

- Computational Memory/In-memory computing
- More complexity and computation/ single unit
- Self-organizing/ Reconfigurable
- Spiking/event-driven communication
- Sub-threshold computation
- Stochasticity as a feature
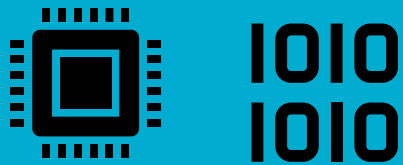- Dense local connectivity
- Massively parallel computation

Human
Brain

Single
Neuron

1MMAC/(s)/pW

Average
CPU

1MMAC/(s)/mW

Hasler 2016

# NEUROMORPHIC COMPUTING

Varied solutions proposed spanning digital, mixed-signal and beyond-CMOS

**DIGITAL CMOS**



CPU, GPU, FPGA

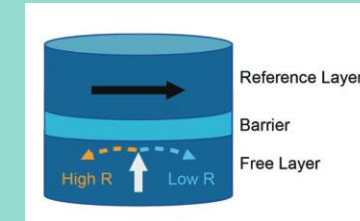- Scaled to 1.15 B/ 2B neurons
- Deterministic
- High-Precision

**ANALOG/ MIXED-SIGNAL**



NVM, Analog, Sub-threshold

- Scaled to 1 M
- Analog/ Stochastic
- Low-Precision

**BEYOND-CMOS**
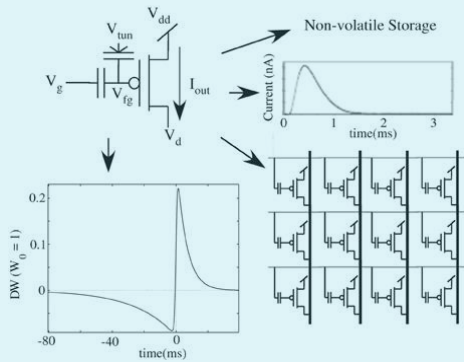


Memristors, FeFETs, MTJs etc.

- Large focus on in-memory computing
- Analog/ Stochastic
- Low-precision
- Integration with CMOS
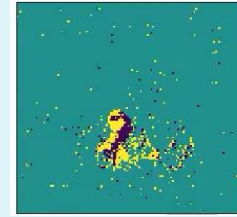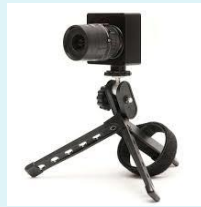
# NEUROMORPHIC BUILDING BLOCKS

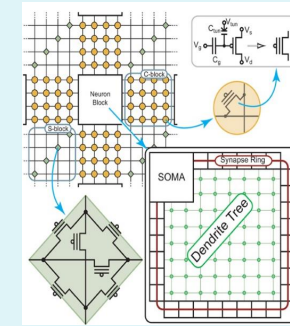Neuromorphic offers computational richness

## In-memory Computing With Synapses
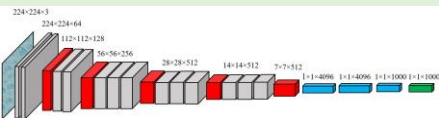


## Sensory Processing



Silicon Retina/ Event Sensor
Silicon Cochlea etc.

## Neural Processing



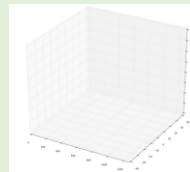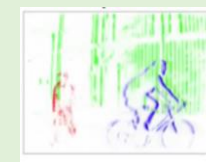Dendrites, Learning,
Multi-modal

## APPLICATIONS



AI/ML
(ANN, SNN)

Brain-Inspired
Algorithms
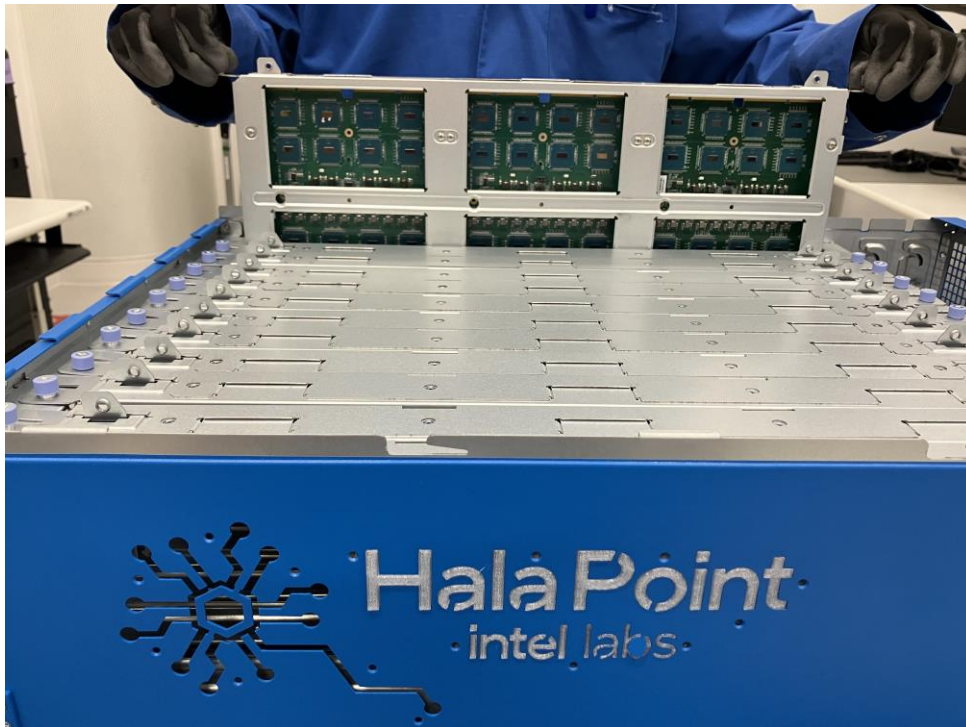
Scientific
Computing

Edge
Computing

High Performance
Computing

# NEUROMORPHIC SUPERCOMPUTER

Neuromorphic systems with billions of neurons



- SNL hosts Intel's Hala Point, which utilizes the Loihi 2 chips to realize one of the largest neuromorphic supercomputers in the world.

- **1.15 Billion neurons** and **trillions of synapses** with a total power consumption of only **2600 W**

- The Hala Point system incorporates 1,152 Loihi 2 processors, each of which can simulate a million neurons.

- Capable of **15 TOPS-per-watt at 8-bit precision** and does not require extensive data-processing or batching in advance.
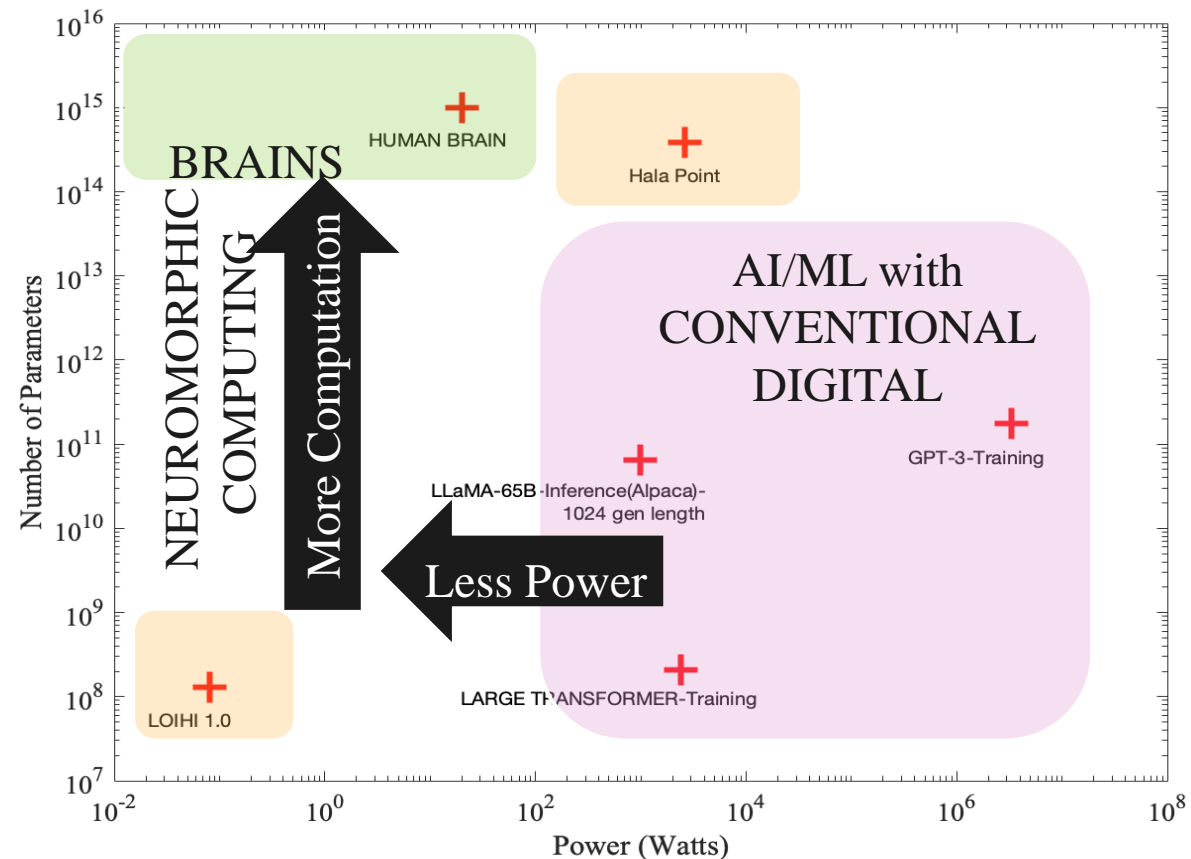
## Next-generation Neuromorphic Architectures for HPC and Edge

Key attributes we need:

- Increase complexity to get more computational/unit

- Leverage stochasticity as a feature

- Scalability and Complexity

- Codesign across scales

- Software tools for Codesign

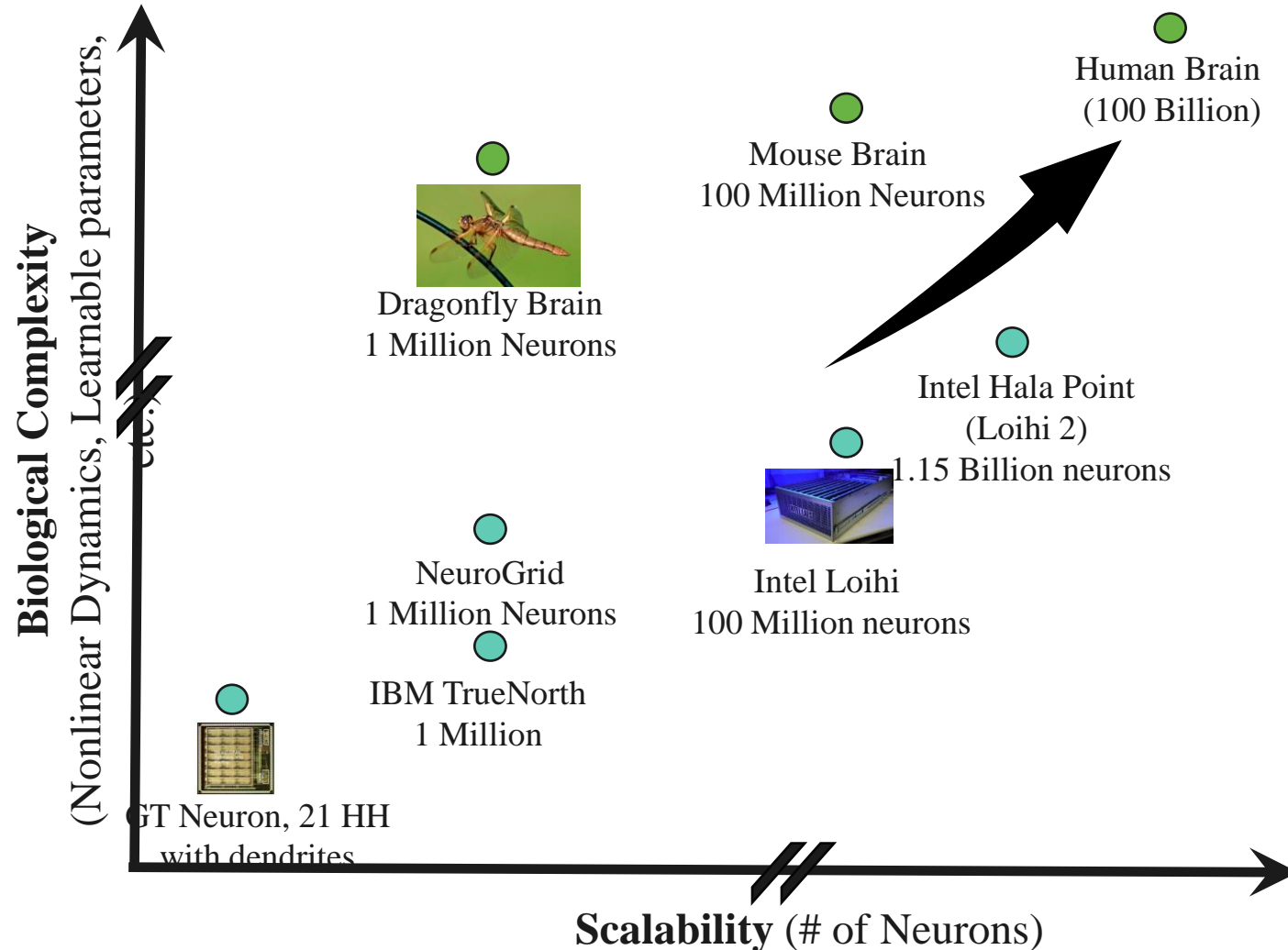- Application specific Solutions

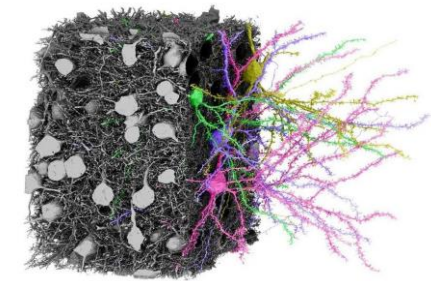- Design of Heterogeneous Architectures



Cardwell et al., Neuromorphic IOP 2024
*Under Review*

# OPPORTUNITIES IN NEUROMORPHIC COMPUTING

## NEUROMORPHIC COMPUTING NEEDS BOTH!



**Biological Complexity**
(Nonlinear Dynamics, Learnable parameters, etc.)

**Scalability** (# of Neurons)

Human Brain
(100 Billion)

Mouse Brain
100 Million Neurons

Dragonfly Brain
1 Million Neurons

Intel Hala Point
(Loihi 2)
1.15 Billion neurons

Intel Loihi
100 Million neurons

NeuroGrid
1 Million Neurons

IBM TrueNorth
1 Million

GT Neuron, 21 HH
with dendrites

- We need to increase computational efficiency as well as computational density for neuromorphic systems.

- We can improve the complexity of a single neuron.

- Novel devices and materials can help bridge this gap but codesign is a challenge.

Biological neurons have rich dynamics and a lot more computational power.

# DENDRITES IMPROVE COMPUTATION WITHIN A SINGLE NEURON
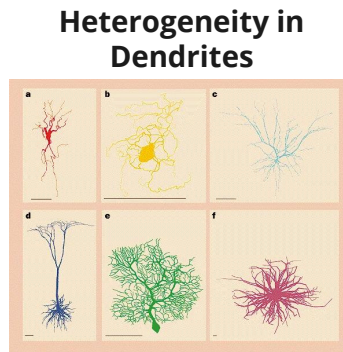
## SIMPLE NEURONS

- **Simpler neuron models requires larger neural networks.**

  - Power hungry ANNs

  - Inefficient scaling

  - Inefficient hardware implementations

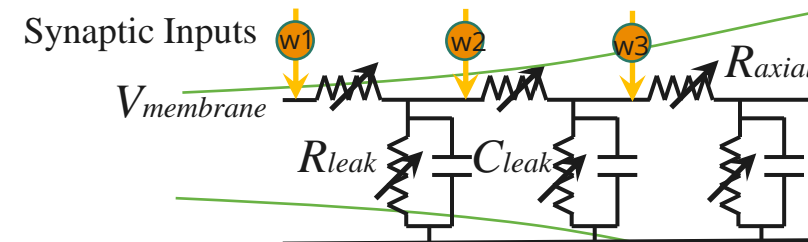  - Focus only on synapses and

### Our Current Research

- Dendrites supporting shunting Inhibition
  - Pattern Recognition (NICE 2023, ICONS 2023)
  - Direction-Selective (ICONS 2024)
- Software library for analog hardware-based dendrites (NeurIPS workshop, NICE 2023)
- DEND-NET: SNN with dendrites (Neuro IOP 2024 - *In Review*)
- Neuromorphic Design Space Exploration with SANA-FE tool in collaboration with UT Austin.

## COMPLEX NEURONS

- **Active Dendrites improves performance of Artificial Neural Networks**

  - More energy efficiency ("Neural Network within a neuron"), non-linear filtering

  - Heterogeneity

  - More computations/unit

  - Better Connectivity and fan-in (3D architectures)

  - Scalable with CMOS+X approaches.

- **Plastic Synapses**

- **Learning Mechanisms**

**Heterogeneity in Dendrites**

Segev, Nature 1998

Synaptic Inputs
$w1$ $w2$ $w3$
$V_{membrane}$ $R_{axial}$
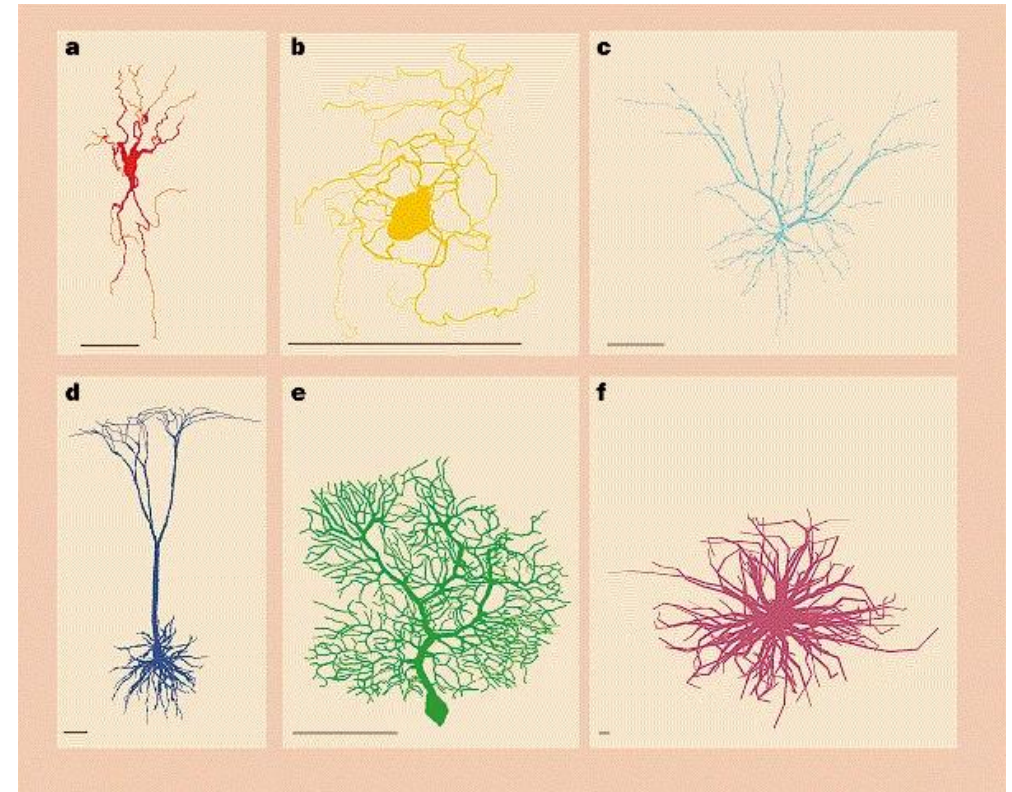$R_{leak}$ $C_{leak}$

**DENDRITE CABLE**

Analog devices and circuits well suited to model dendrites efficiently in hardware.

# COMPUTATION USING THE DENDRITES

Dendrites are not just wires!

- Over 70% of the neuron's volume (Stuart 2016)

- Great diversity of dendrites within a single brain and across animal species.

- Insights from neuroscience are foundational to the pursuits of neuromorphic computing.

- Biological dendrites are known for their complex physical structures that incorporate significant fan-in  (~10,000 inputs).
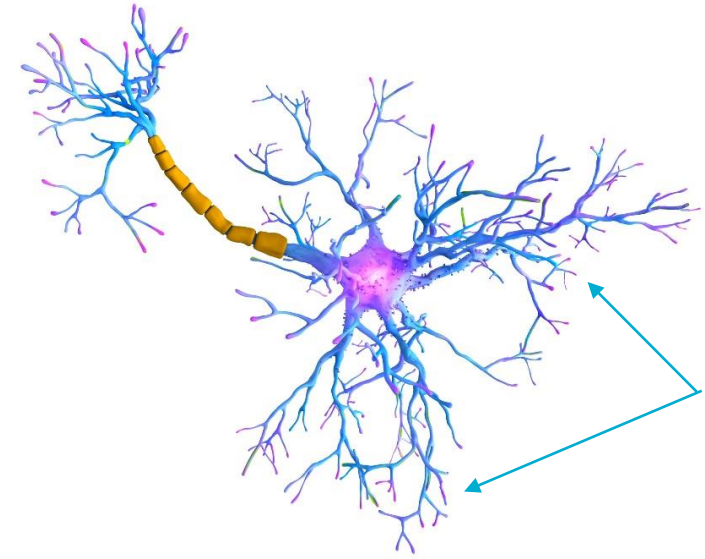


Segev, Nature 1998

# DENDRITIC TOOLKIT FOR COMPUTATION

Dendrites are tree-like structures that connect neur
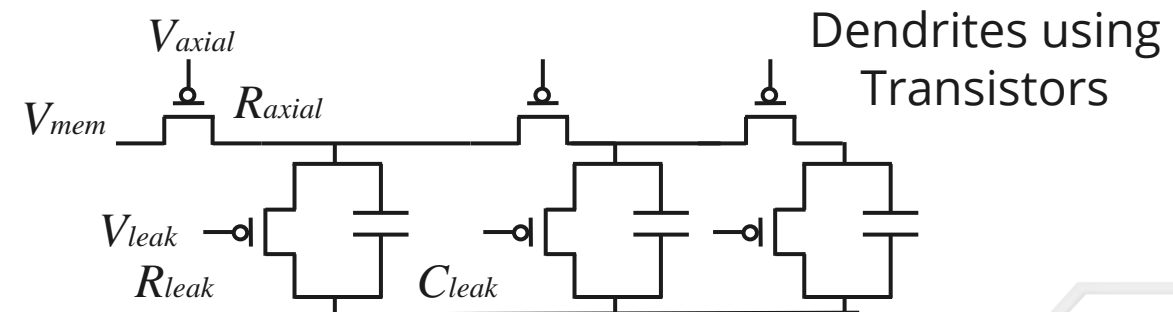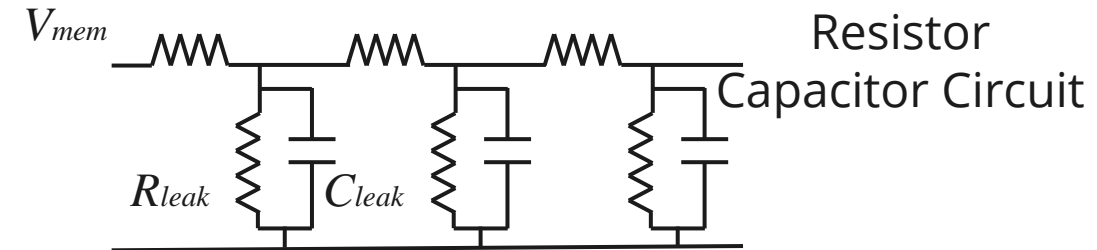


Dendrites

- **Dendrites are not *just* wires!**

- They can perform interesting computation like:

  - Coincidence Detection

  - Current Summation

  - Directional selectivity

  - Non-linear filtering

  - Amplification of Synaptic inputs

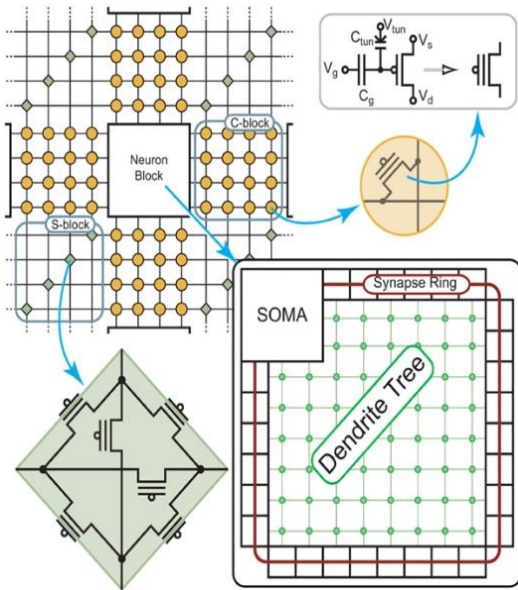    London 2005, Poirazi 2020

Increased Connectivity and Computation



Resistor Capacitor Circuit

$V_{mem}$

$R_{leak}$    $C_{leak}$

Dendrites using Transistors

$V_{axial}$

$V_{mem}$    $R_{axial}$

$V_{leak}$

$R_{leak}$    $C_{leak}$

Dendrites have been modeled to different degrees in neuromorphic hardware.



Hodgkin Huxley
Neurons with
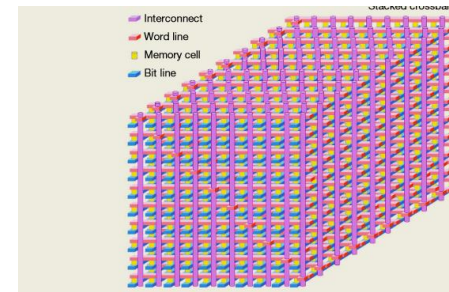Active Dendrites
Ramakrishnan, 2013

- Active Dendrites with Calcium and NMDA channels: BrainScales

- Floating-gate based active dendrites: Georgia Tech Neuron Chips

- Dendrocentric Learning with mult-gate FeFETS (Boahen 2020)
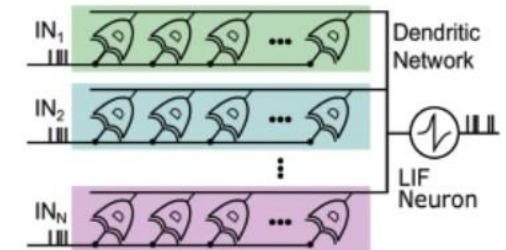
- DenRAM: Dendrites with Resistive RAM (Payvand 2024)
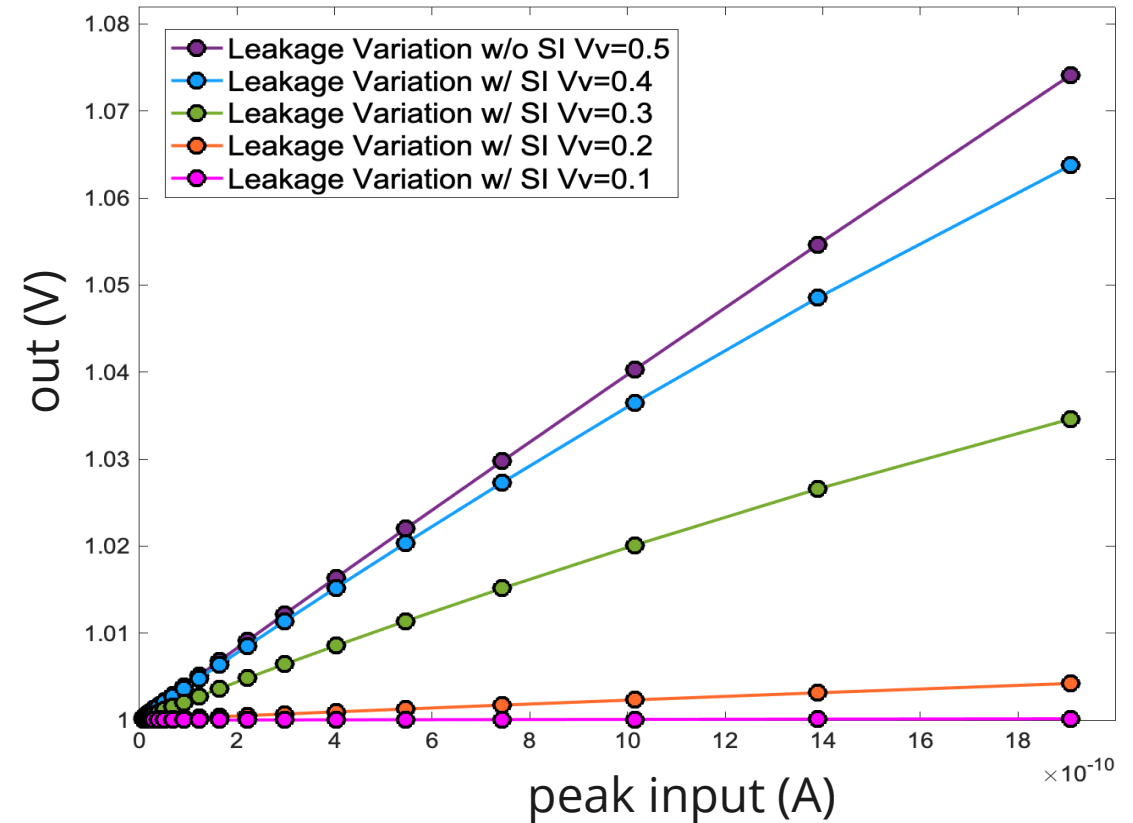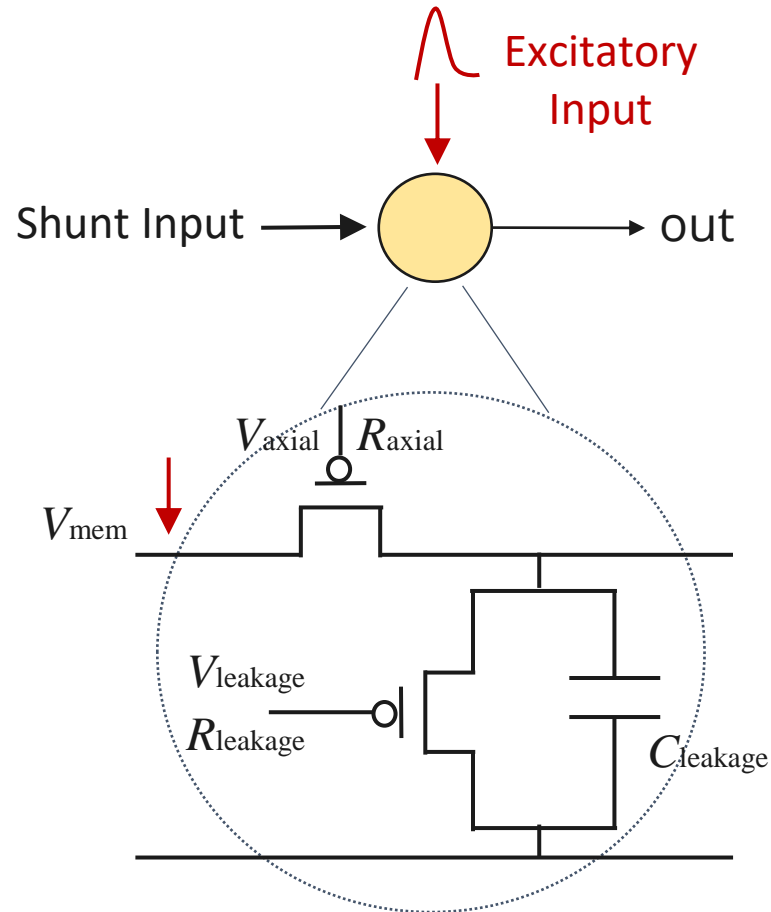


NanoDendrite
Multi-gate FeFET



Proposed Dendrite
3D Architecture



DenRAM with RRAM
dendrites

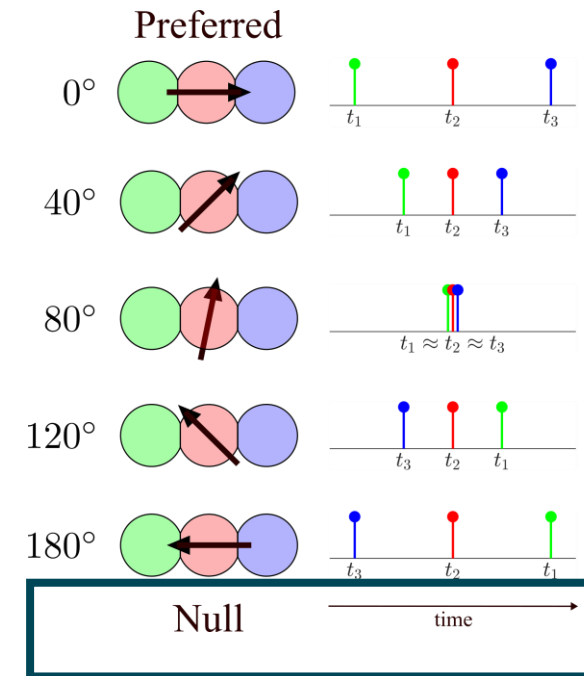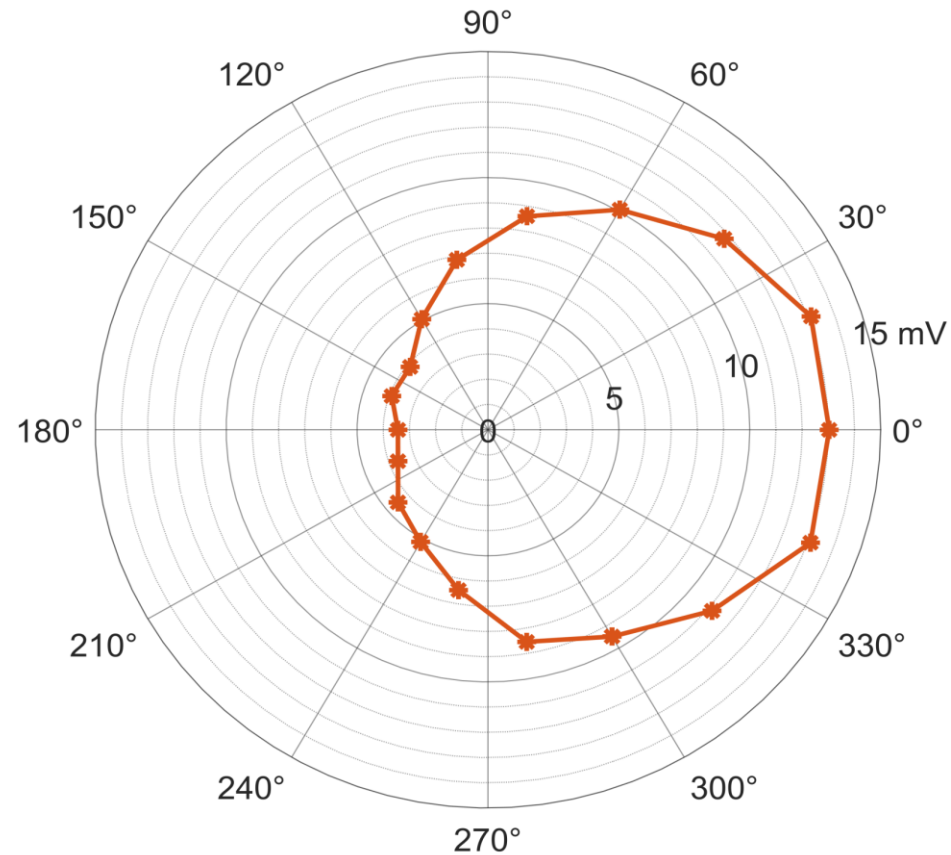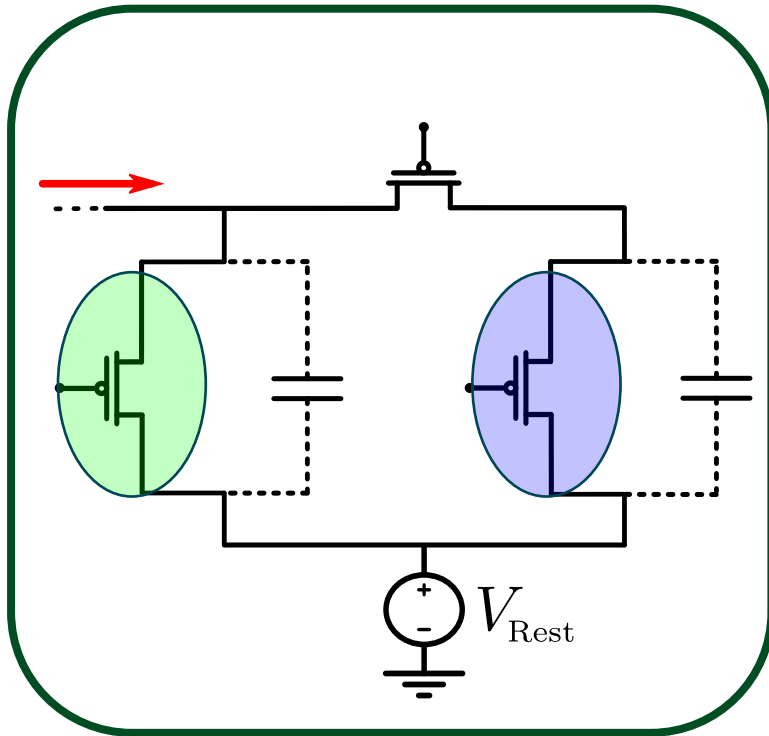Using shunting inhibition for gain modulation

Experimental Demonstration using analog circuits

FPAA 2.9V

## Dendrite Circuit



$V_{\mathrm{Rest}}$



Preferred

0°

40°

80°

120°

180°

Null

time

Parker et. al IEEE ICONS 2024

Dendrites for pattern and velocity estimation

Super-Pixel

w11  w12  w13

w21  w22  w23

...

wn1  wn2  wn3
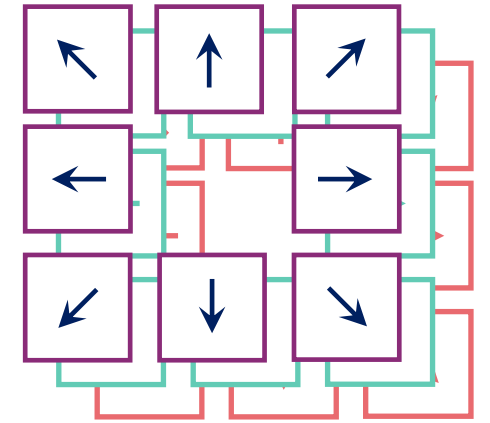
Image Sensor
(NXM pixels)

Feature
Extraction

Complex
Features

Using shunting inhibition for gain modulation

Cardwell et al. *In Review*



Original Dragonfly
Interception Circuit
(Chance 2020)

Dragonfly NN with
Dendrites

Comparison with original
interception circuit

# CO-DESIGN IS CHALLENGING

Software tools are critical for design and co-optimization



Application/ Algorithm Design

System Architecture Design

Codesign

Circuits/ System Design

Materials/Device Design

## Comparing networks with LIF and LIF +Dendrites



256 / 16 Signals to LIF Neuron

16 Signals to LIF Neuron

MSE Of LIF and Dendritic Networks

Leveraging inherent properties of dendrites

- Implemented a "Dendrite Pooling Layer" for use in ANN

- Trained ResNet18 on CIFAR-10 for 300 epochs

  - ResNet + Dendrite layer took significantly longer to train

  - Simplified ODE layer adds state and loops
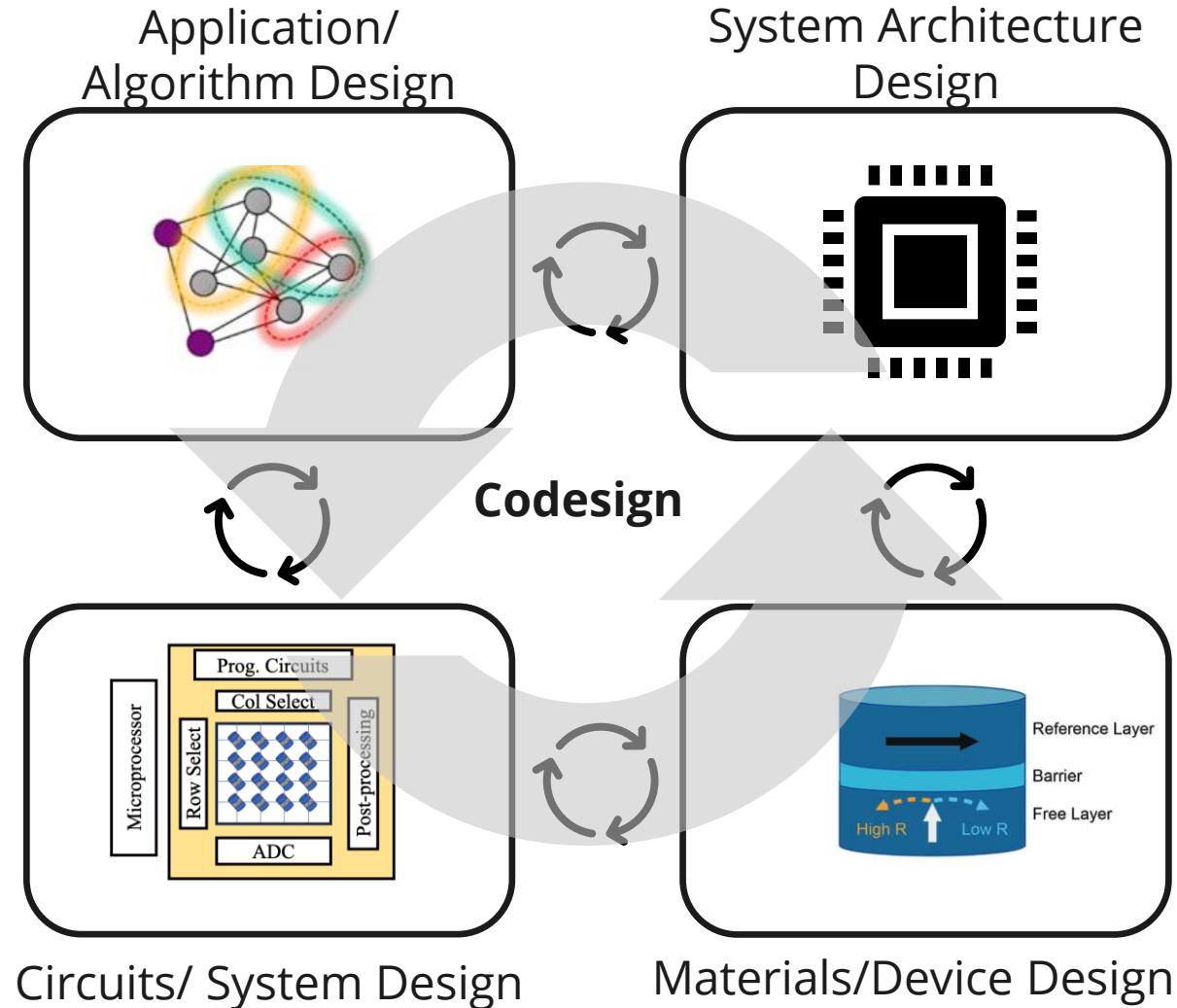
- Found accuracy to be comparable

- Dendritic pooling has potential in ANNs



| Pooling Layer Hardware | Energy |
|---|---|
| Digital Nvidia Jetson | 504.41 uJ |
| Dendritic pooling | 0.265 uJ |

Plagge et al., IEEE ICRC 2024

# DENDRITE ENABLED SPIKING NEURAL NETWORKS

- Implemented Torch library with a dendritic chain

  - Simplified version of the complex ODE dendrite solution.

  - Wrapped dynamics into a set of constants and parameters.

- Dendrites support SNNTorch & Non-Spiking Torch.

- Provides a trainable 1-D chain of dendrites

| Value | Type |
|---|---|
| **Lambda** | "**Spatial**" constant: Represents Distance |
| **Tau** | "**Temporal**" constant: Capacitance and Resistance |
| **Leak** | Signal loss for each tap |
| **Input Weight** | Increases or Decreases signal strength |



Hardware Constraints

Learned Parameters

Plagge et al., IEEE NICE 2024

SANA-FE: Specialized tool to explore novel neuromorphic architectures

- Rapidly estimate performance of neuromorphic architectures for design-space exploration

- General & extensible spiking H/W simulator

- Model functional behavior & track performance

- Schedule messages & intra-core interactions

- Calibrate simulator to real-world systems

- Accurately predicts latency & energy of gesture categorization spiking neural network (SNN)

**SANA-FE: Simulating Advanced Neuromorphic Architectures for Fast Exploration**

Configuration & Input Spikes

```
Simulator Kernel

build architecture
initialize network

for all timesteps:
    get external inputs
    for all tiles:
        for all cores:
            process neuron
    send messages

    write results
```

Architecture Description

Mapped Spiking Neural Network

Performance Estimates

J. Boyle et al. ICONS 2023

# DVS GESTURE RECOGNITION APPLICATION

- Predict energy & performance for larger real-world neuromorphic applications

  - SNN trained on DVS gesture data-set [Massa'20]

  - 18,678 neurons across 6 layers

  - Mapped to 45 Loihi cores out of 128

# DVS GESTURE DESIGN SPACE EXPLORATION

- Design-space exploration using DVS gesture application
  - Loihi-based designs, traded-off core count (c) vs neurons per core (n)
  - Optimum design had 170 cores, 21% faster than **Loihi (128 cores)**
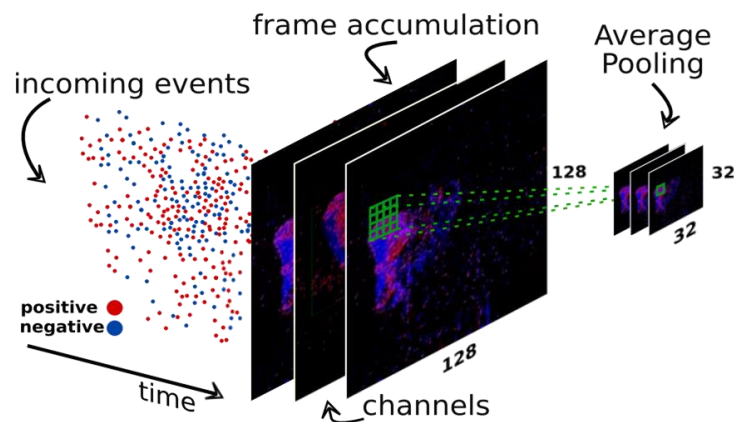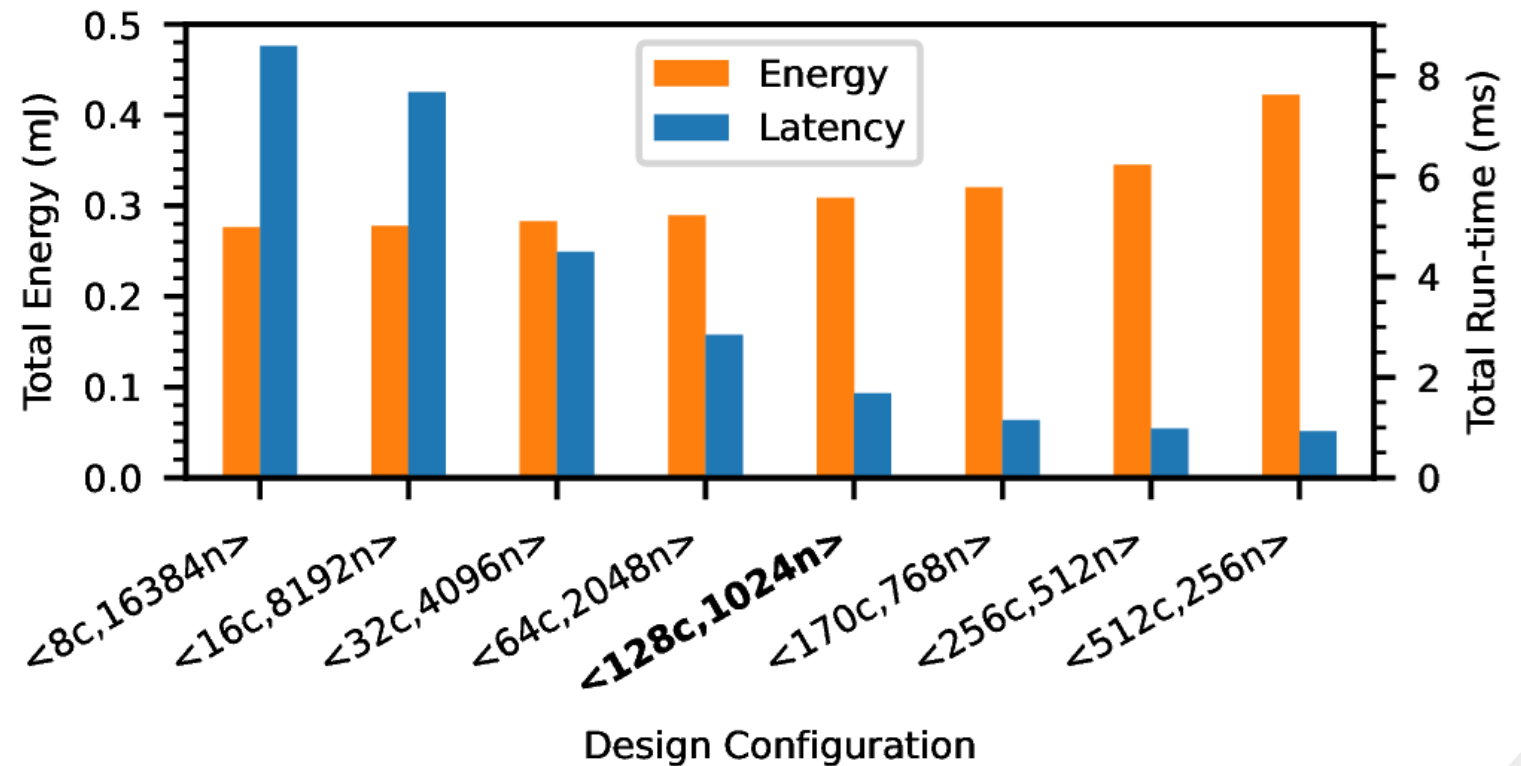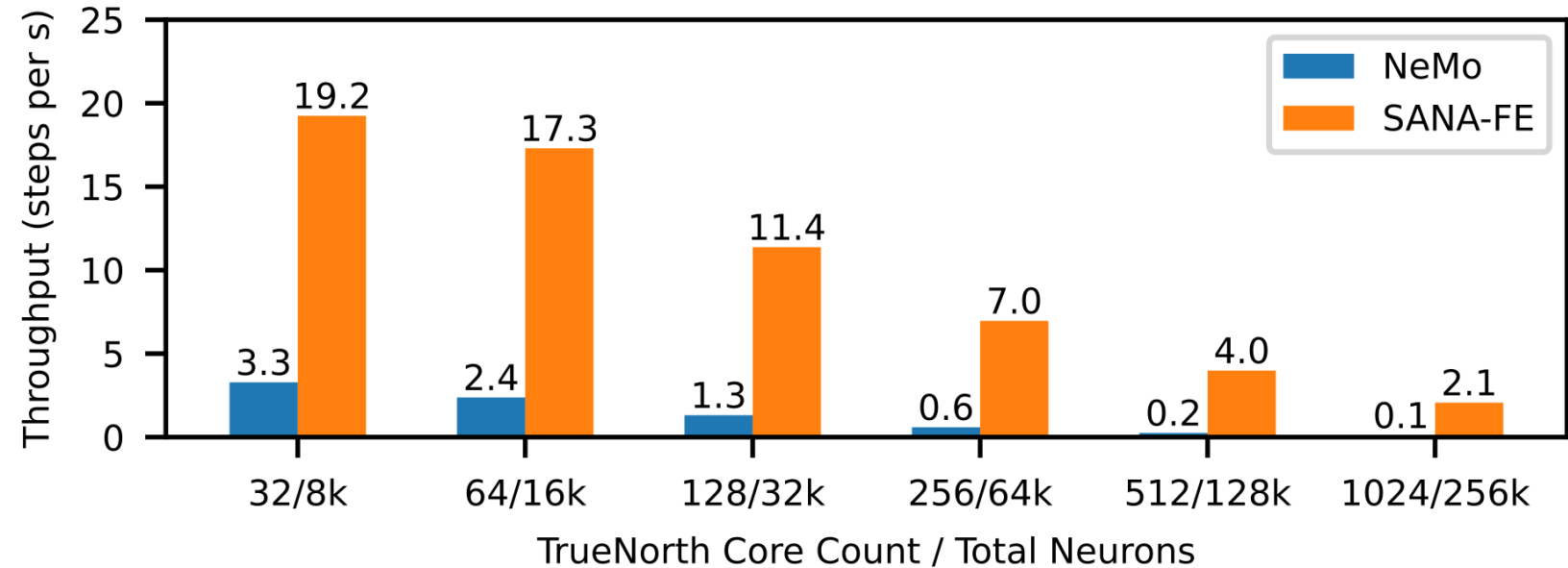  - Design-space sweep took 29 s



Image reproduced from [Massa,'20]

# SIMULATOR SPEED RESULTS

- **Compared to existing discrete-event based spiking simulator (NeMo)**



⓾ Simulating IBM TrueNorth architecture
⓾ Randomized SNN with 80% of spikes intra-core, 20% spikes between cores

➤ **Over 20x faster than NeMo for 1024 cores**

# DENDRITES IN NEUROMORPHIC ARCHITECTURE

- Further develop links with Dendrite-SNN hardware simulations – SanaFe

- Work on a spiking self-attention network with dendrites:

  - Dendritic attention layer (Temporal coherence and context)

  - Dendritic pooling layers (More efficient summary layer)

- Other compelling network designs

- Release as stand-alone library or as SNNTorch add-on

SanaFe – Hardware Simulator



An in-progress tool to estimate timing and energy of neuromorphic systems.
Currently supports Loihi.
Analog components are WIP

## Embracing Heterogeneity

- Move past general-purpose solutions and only use them for prototyping

- Heterogeneous and reconfigurable systems

- Get the best performance based on system needs.

- System designed with application in mind for optimization

- Near-sensor processing



'Truly Heterogeneous Computing', Cardwell et al., SMC 2020

We are in the "golden-age of computer architectures" - Patterson



1. Cui et al. arXiv 2024

2. Karki et al. Journal of
   Materials Research 2024

**OBJECTIVE:** Leverage stochasticity in computing by exploiting the underlying physics of emerging random number generator (RNG) devices to build probabilistic neural architectures.

Collaborators: NYU, ORNL, Temple University, UT-Austin and UT-Knoxville, USA

# TRUE RANDOM NUMBER GENERATION (TRNG)

## APPLICATION: HIGH ENERGY PHYSICS

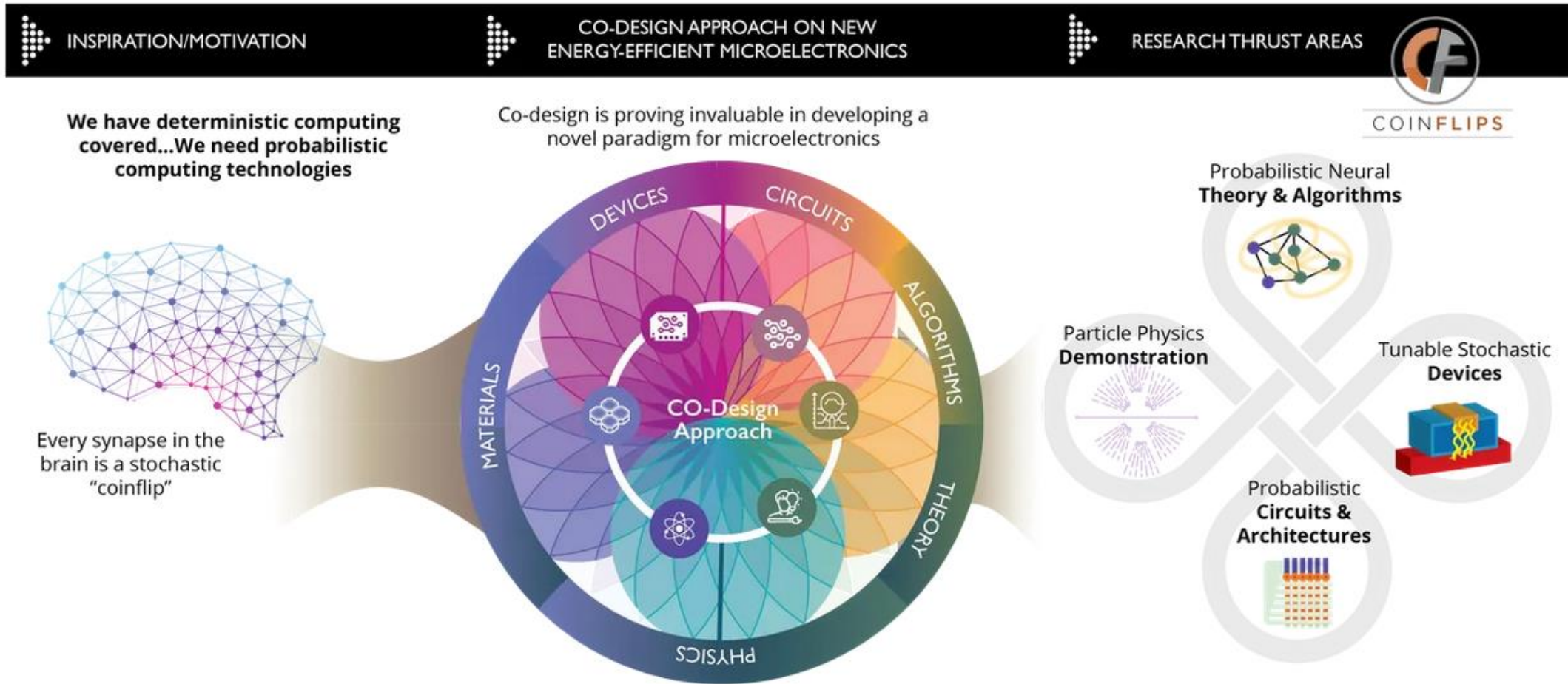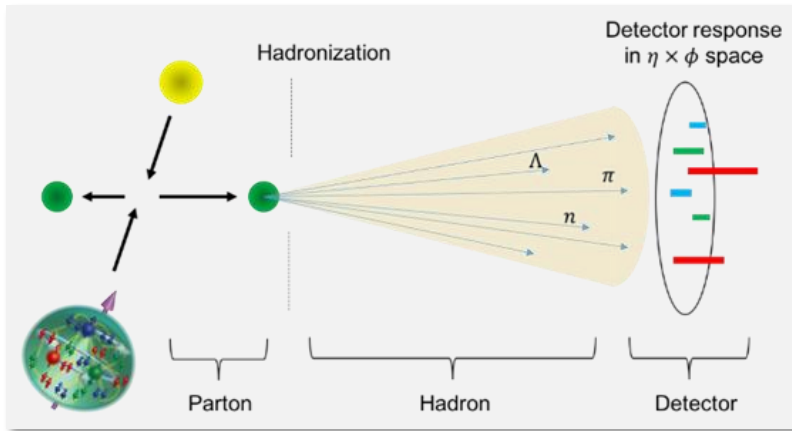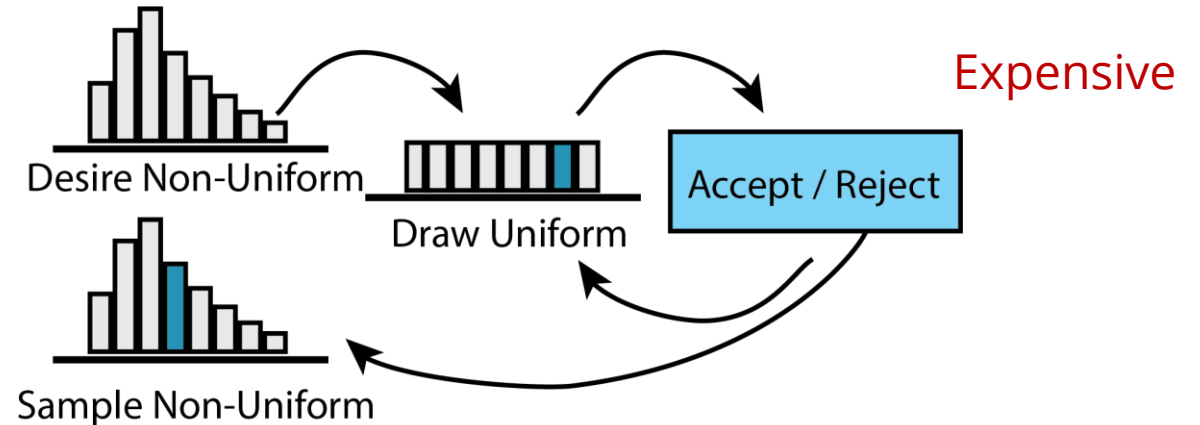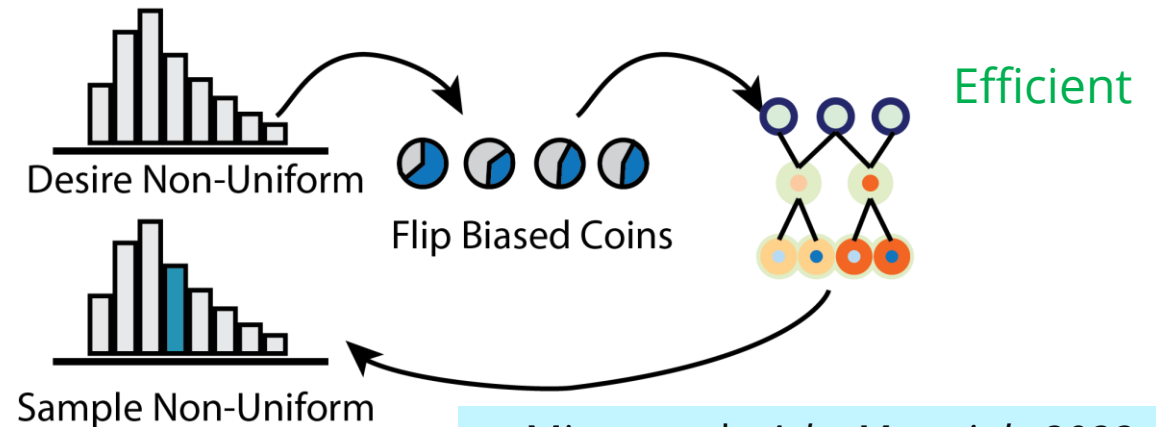- Current PRNG Methods take 40—50% of CPU compute time.

- TRNGs leveraging stochastic devices can lead to significant energy and latency savings.



Pierog et al., Phy Rev. 2022

## RNG TODAY



## COINFLIPS APPROACH



Misra et al., *Adv. Materials 2022*

# AI TO ACCELERATE CODESIGN FOR EMERGING COMPUTING

AI-Guided Tools a force multiplier for large design space

- Computing demands are constantly increasing.

- Emerging computing techniques can alleviate these challenges.

- However, the design space is huge and optimization is needed across the stack.

- AI-guided techniques can alleviate these challenges.

Misra et al., *Adv. Materials 2022*

Application/
Algorithm Design

System Architecture
Design

**AI-Guided
Codesign**

Prog. Circuits
Col Select
Microprocessor
Row Select
Post-processing
ADC

Reference Layer
Barrier
Free Layer
High R      Low R

Circuits/ System Design

Materials/Device Design

# AI-GUIDED METHODS FOR EMERGING COMPUTING

## EVOLUTIONARY OPTIMIZATION

- Evolutionary algorithms
  - LEAP: Library of Evolutionary Algorithms in Python
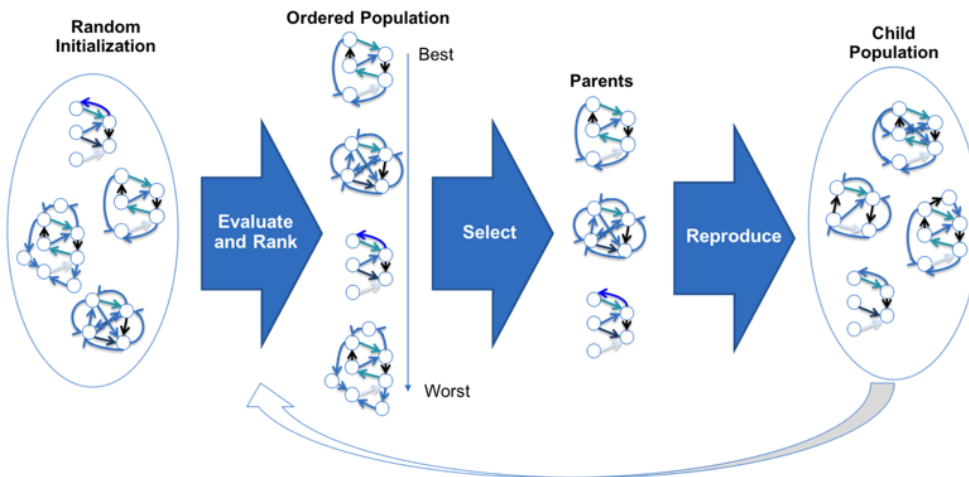  - EONS: Evolutionary Optimization for Neuromorphic Systems

## REINFORCEMENT LEARNING

- Trains agent to make optimal decisions in an environment to maximize rewards.
  - Agent trained on PPO policy
  - Environment: Physics-based Device model





**Other AI Approaches: Physics-aware machine learning, Generative AI, Neural ODEs**

**FITNESS FUNCTION**

$$f(w, p1, p2, q1, q2) = \boldsymbol{\omega_1} KL(p_1, p_2, q_1, q_2) + \boldsymbol{\omega_2}(\sum_{i=1}^{2} |p_i - 0.5| + \sum_{i=1}^{2} |q_i - 0.5|) + \boldsymbol{\omega_3} EN(p_1, p_2, q_1, q_2)$$

| Kullback-Leibler Divergence | Difference of weight from a fair coin | Energy of a coinflip |
|---|---|---|



SHE: Spin Hall Effect

Multi-objective optimization of weights ω1, ω2, ω3 for optimal KL divergence and energy usage of MTJ-SHE devices

Cardwell et al., IEEE ICRC 2022

Discover best device and material characteristics for TRNG

Device Check



**ACTIONS**

Update Device and Materials Parameters

**AGENT/POLICY PPO**

**ENVIRONMENT MTJ Models**

Rewards Based on Metrics of Interest

**REWARDS**

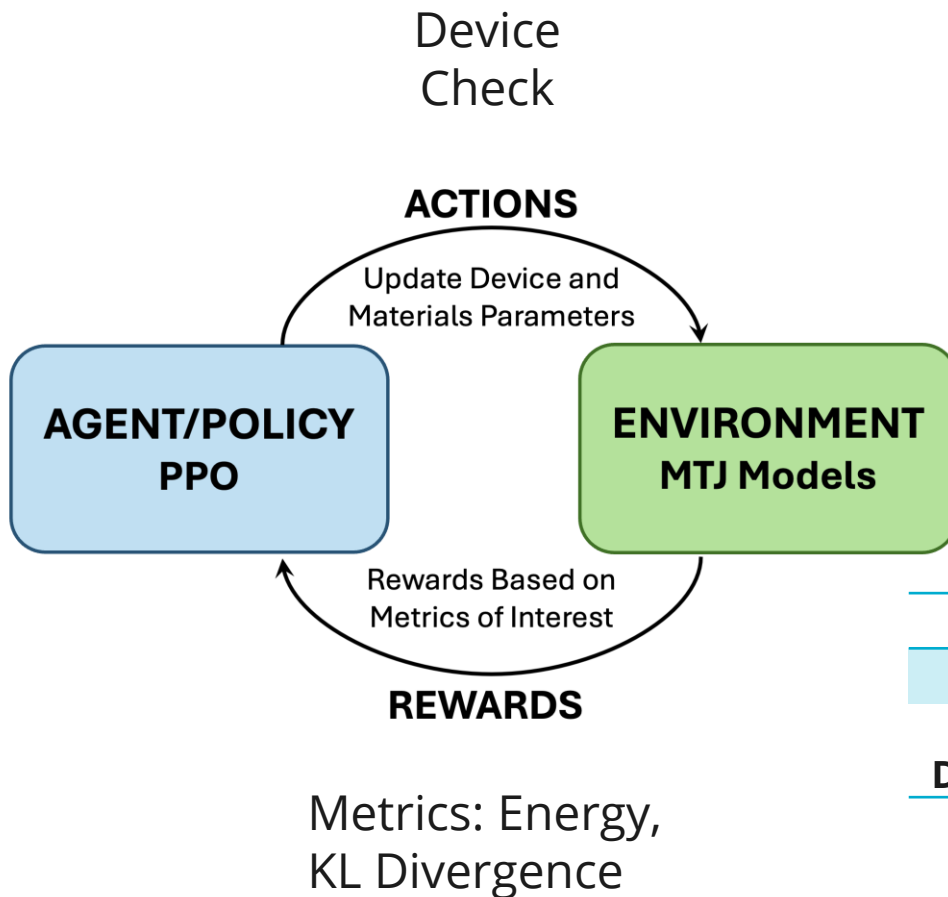Metrics: Energy, KL Divergence

- RL agent trained using PPO policy to find best device and material configuration.

- Device validity checked.

- Best configuration found normalizing for energy and KL-divergence.

| Metric | Best Config |
|---|---|
| **Energy (J)** | $2.568 \times 10^{-14}$ |
| **KL-Divergence** | 0.0532 |



PDF Comparison (w/ Normalization)

Cardwell et al., *IEEE ISCAS 2024*
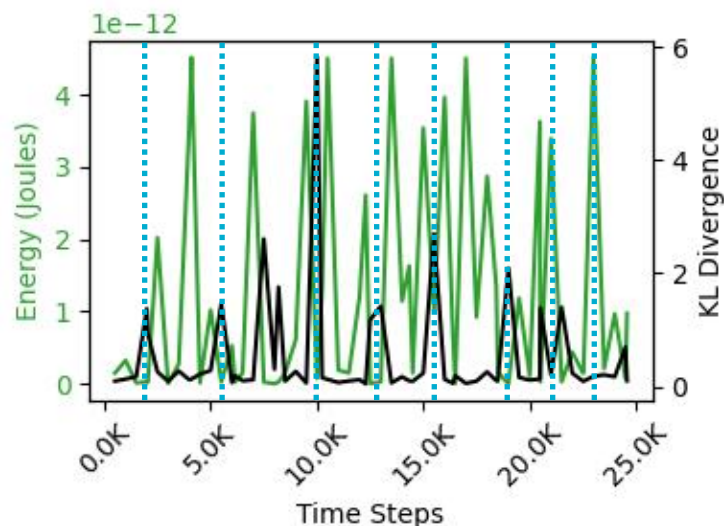
# DEVICE PERFORMANCE TRADEOFFS

## TRAINING RL AGENT

- Agent had to balance both energy and KL-divergence for optimization which seemed to have an inverse relationship.
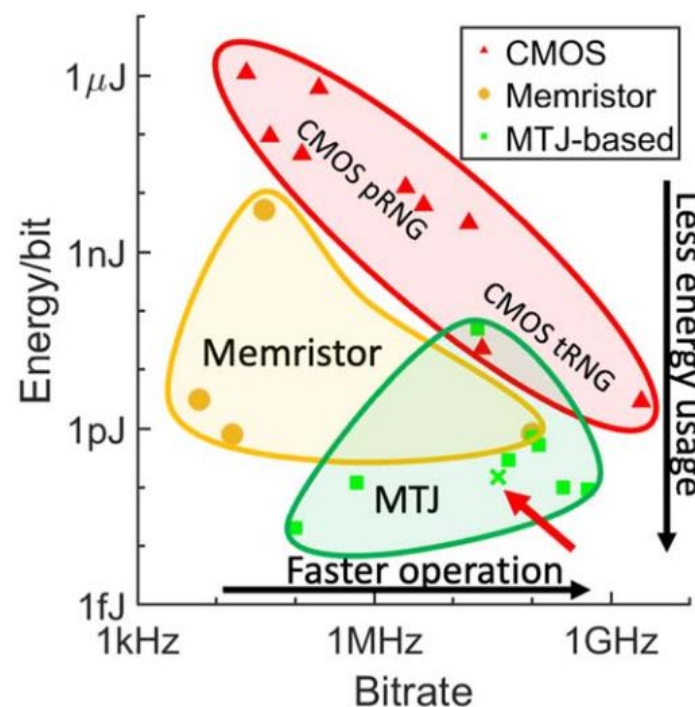
- Reward schema is extremely important.

Energy and KL-Divergence

Cardwell et al., *IEEE ISCAS 2024*

## BENCHMARKING TRNGS

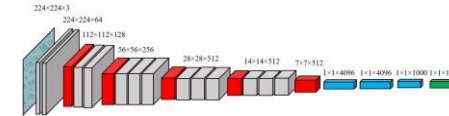- Comparing CMOS pRNG, tRNG, memristor tRNG and MTJ tRNG
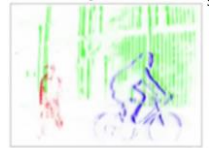
Maicke et al., *IOP Nano 2024*

**Re-think how we design computer architectures**

- We need more dynamics and complexity per computational unit.

- Leverage Stochasticity as a feature not a bug.

- We need systems that do not just process, but can learn, adapt and reconfigure.

- Novel integration: 3D architectures, wafer scale etc. for scaling and dense connectivity

- Many areas where neuromorphic can have impact from HPC to edge.

- Codesign tools can accelerate design space exploration and lead to creative solutions.
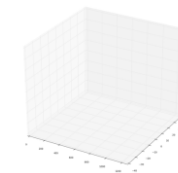
AI/ML
(ANN, SNN)

Edge
Computing

Brain-Inspired
Algorithms

High Performance
Computing

Scientific
Computing

Probabilistic
Computing

**Sandia National Laboratories NEUROMORPHIC**

# THANK YOU!

SGCARDW@SANDIA.GOV

# BACKUPS

# ATHENA
## (ANALYTICAL TOOL TO EVALUATE HETEROGENEOUS NEUROMORPHIC ARCHITECTURES)

- ATHENA will quickly evaluate performance metrics of analog architectures

- Developed as part of a larger ecosystem
  - Tools to enable next-generation hardware design prototyping

Plagge et al., International Conference on Rebooting Computing (ICRC) 2022

ATHENA – HARDWARE PERFORMANCE



Plagge et al., International Conference on Rebooting Computing (ICRC) 2022

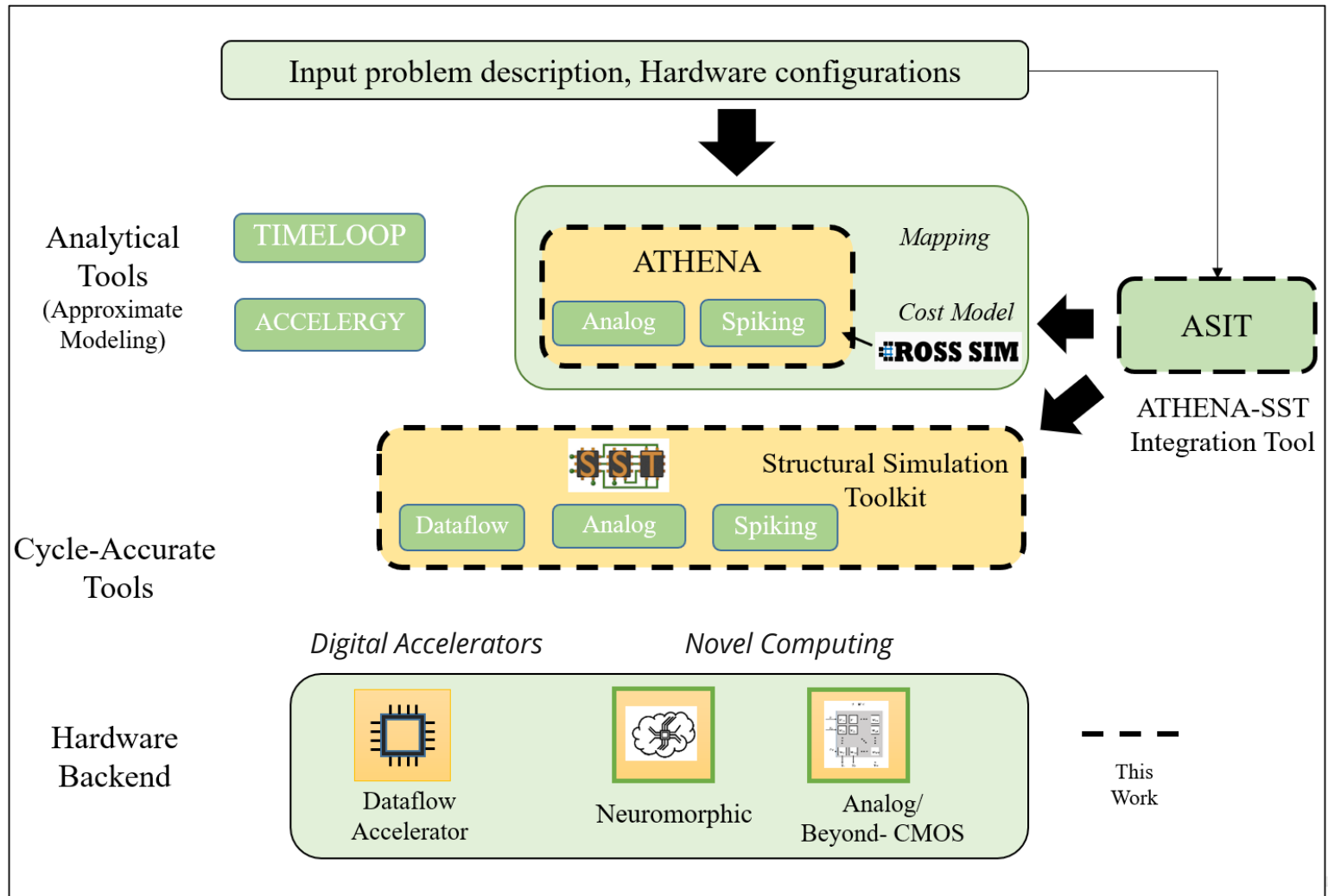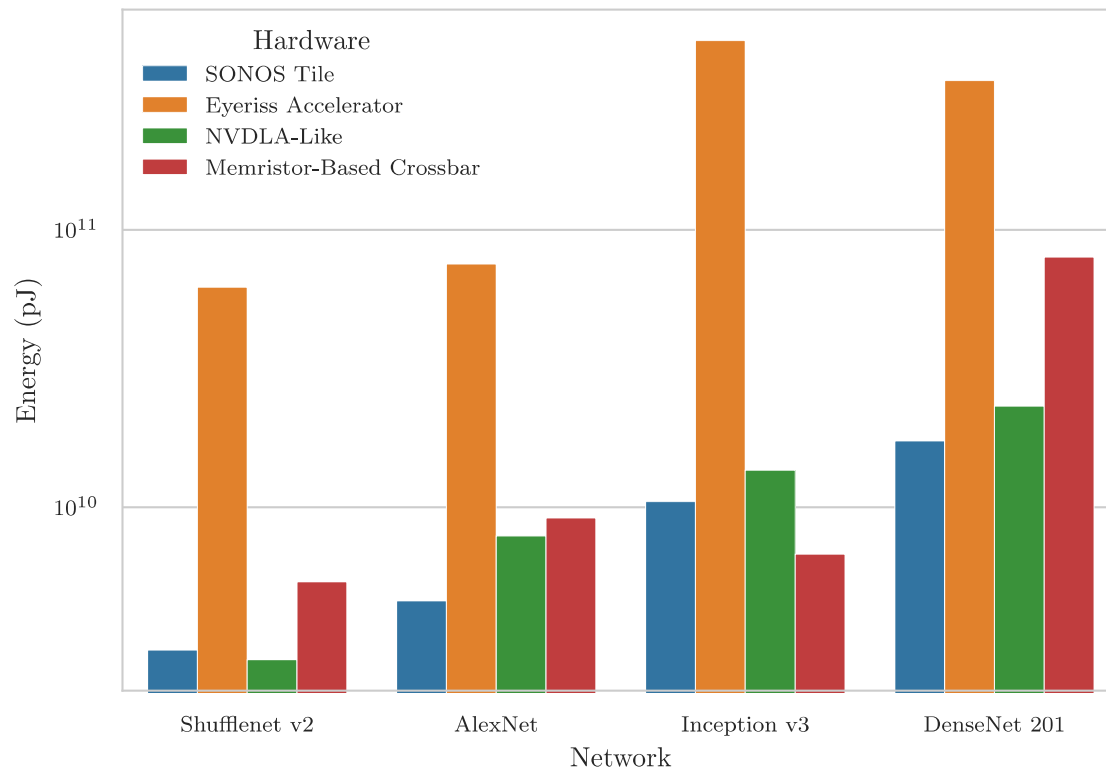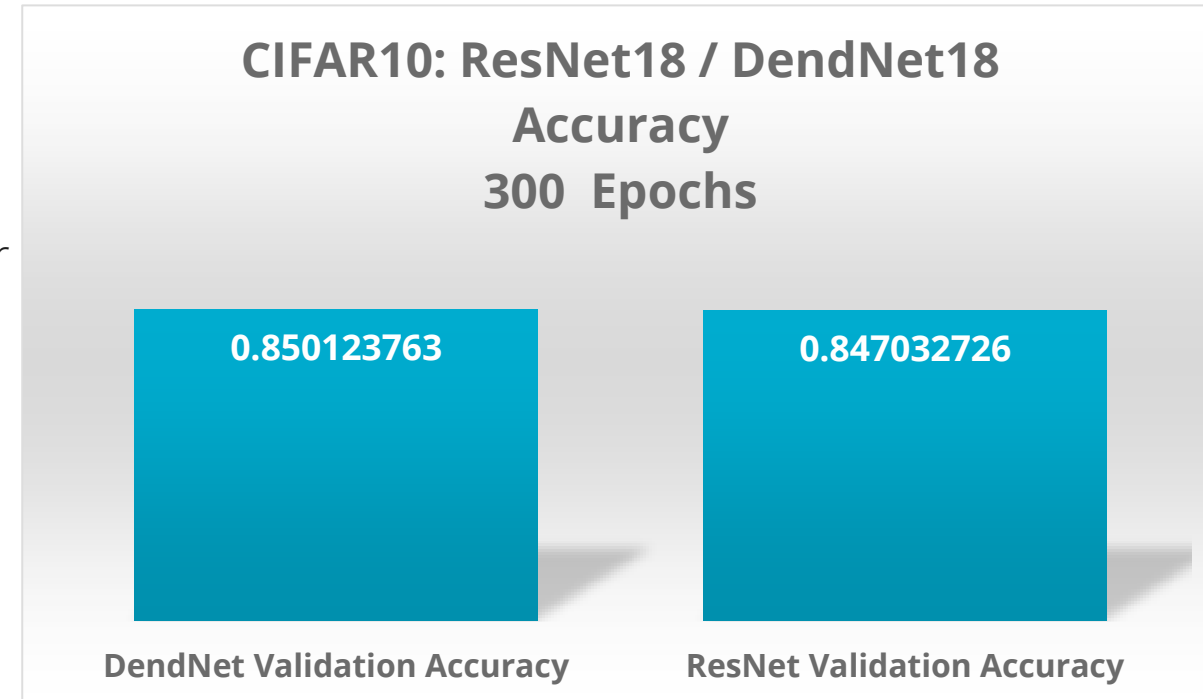- ATHENA was used to compare the performance of multiple hardware devices against various deep learning networks

- The SONOS tile-based architecture performed well across networks, with one notable exception: the Inception v3 network

- This performance difference could be explored – showing ATHENA's potential for codesign work

# ANN NETWORK APPLICATIONS – RESNET18

- Working with a graduate student intern at SNL

- Implemented a "Dendrite Pooling Layer" for use in AI ML

- Replaced traditional pooling layer with Dendrite Layer

- Trained ResNet18 on CIFAR-10 for 300 epochs
  - ResNet + Dendrite layer took significantly longer to train
  - Simplified ODE layer adds state and loops

- Found accuracy to be comparable
  - Dendritic pooling has potential in ANNs

**CIFAR10: ResNet18 / DendNet18 Accuracy 300 Epochs**

| 0.850123763 | 0.847032726 |
| --- | --- |
| **DendNet Validation Accuracy** | **ResNet Validation Accuracy** |

*Working with Priyam Mazumdar*

# ML NETWORK APPLICATIONS – RESNET18

- ResNet18 – Was slower to train with a dendritic layer

  – In hardware however, dendrites will be highly efficient

- Rough estimate of efficiency based on

  – Energy = $C(V_{mem} - Ek)V_{dd} = 500fJ$

    o C = 10pF

    o $V_{dd}$=2.5V

    o $V_{mem} - Ek$ = 100mV

  – Nvidia

| Pooling Layer on Digital Nvidia Jetson | 504.41 Micro Joules |
|---|---|
| Dendritic pooling | 0.265 Micro Joules |

*Rodrigues, Crefeda Faviola, Graham Riley, and Mikel Luján.*
*"Energy predictive models for convolutional neural networks on mobile platforms."*
*arXiv preprint arXiv:2004.05137 (2020).*