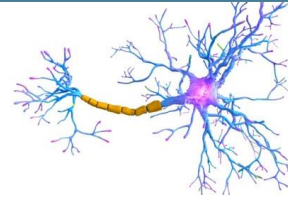# Modeling and Simulation Challenges of Neuromorphic Architectures

**Suma George Cardwell**

Cognitive and Emerging Computing

Sandia National Laboratories

**ModSim 2023**

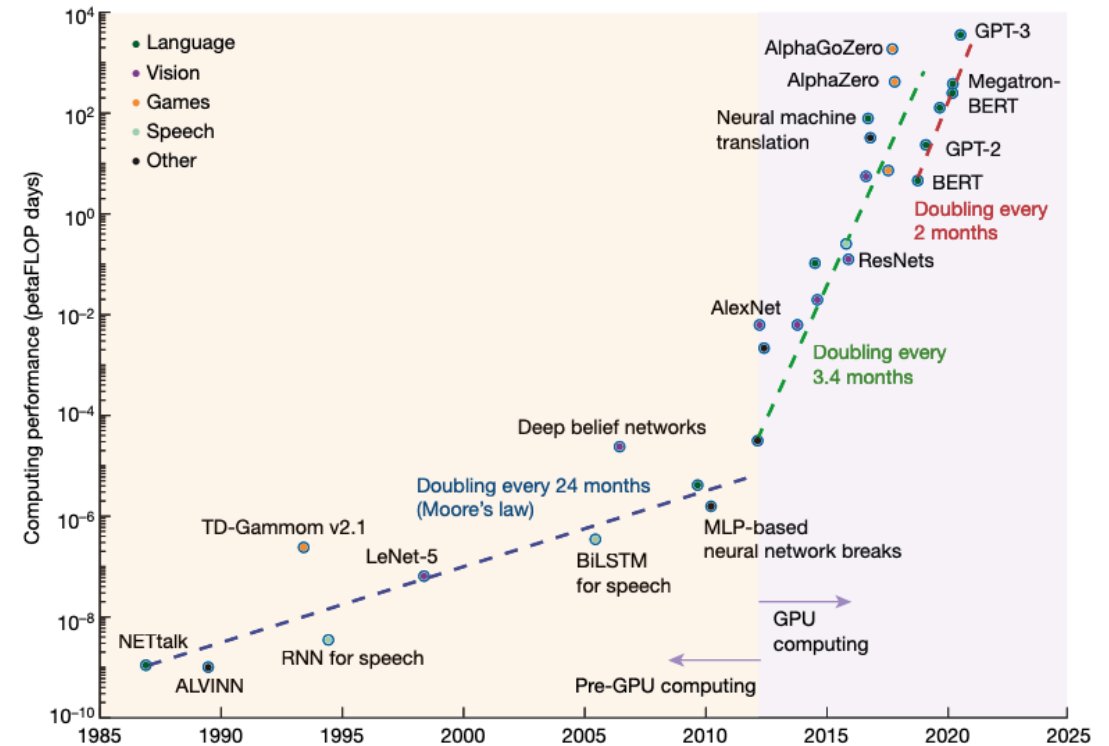Workshop on Modeling & Simulation of Systems and Applications

August 9th, 2023

# FUNDAMENTAL CHALLENGES IN COMPUTING

- Limits of scaling have ushered in the "Golden Age of Computer Architecture" Hennessy & Patterson 2019

- Inefficiency of generality

- Performance saturation

- Growing demands for HPC and SWaP-constrained edge computing



AI compute demands are increasing

Mehonic & Kenyon 2022, Open AI Research Blog

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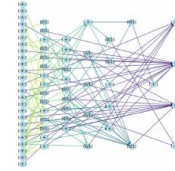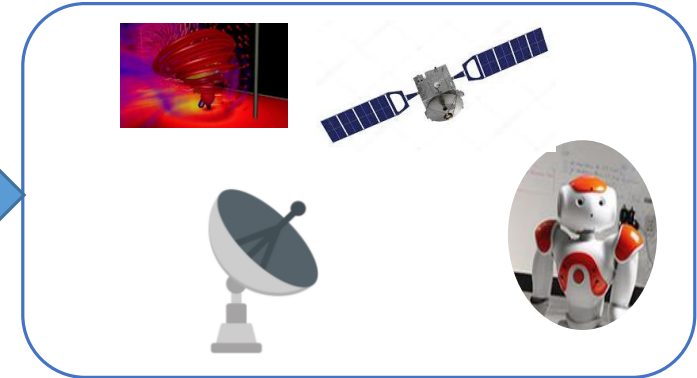
# COMPUTING LANDSCAPE

**Sensors**

**Algorithms**

- Scientific Computation
- Machine Learning
- Brain-derived algorithms
- Signal Processing

**Applications**

**50 billion IoT devices by 2030**

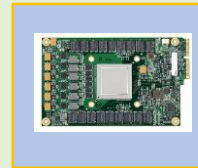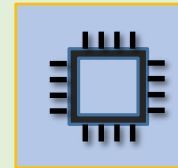**Modern Computing**

*Conventional Digital*

*Novel Computing Paradigms*

| CPUs | GPUs | FPGAs | TPUs | ASICs | Neuromorphic | Quantum |

**Digital Neuromorphic**

**Analog/Mixed-signal Neuromorphic**

**Beyond CMOS devices**

Intel Loihi
Davies 2018

SpiNNaker
Furber 2016

GT Neuron
Brink 2013

DAVIS 240C,
DYNAPSEL

NeuroGrid
Benjamin 2014

**RERAM**
Marinella et al., 2016

**Mott- Memristor**
Kumar et al., 2020

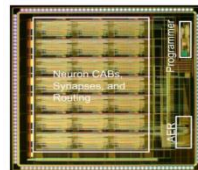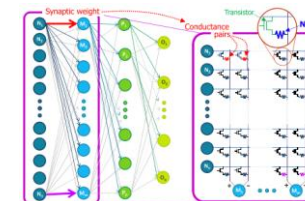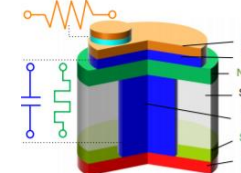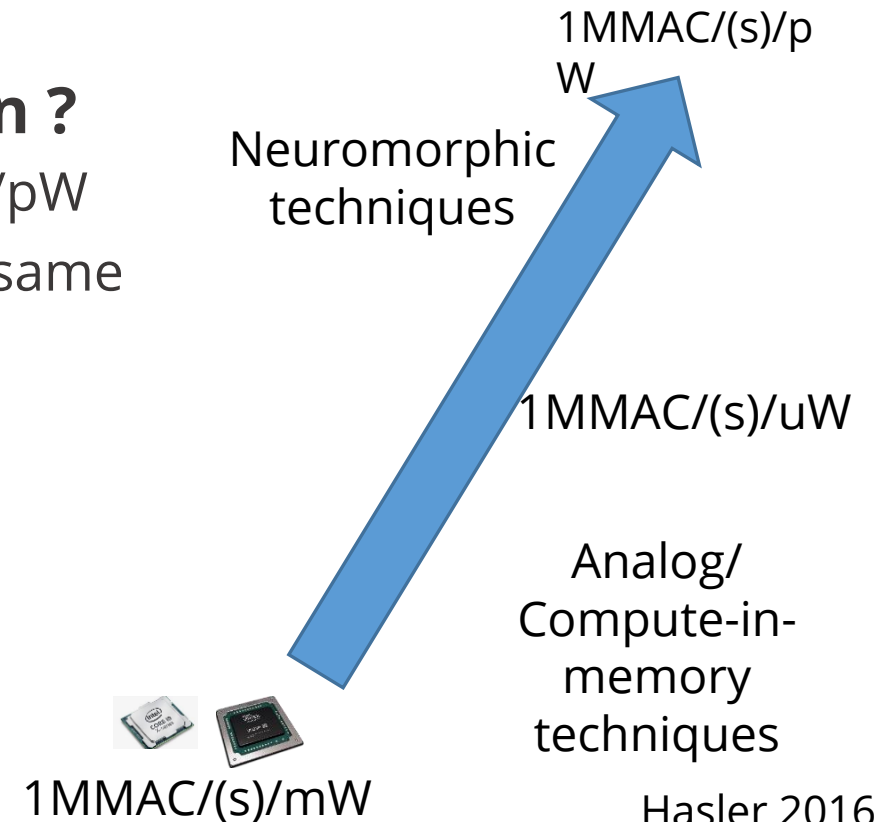# NEUROMORPHIC COMPUTING: INSPIRED BY THE BRAIN

## Brain and Computing: Why make the connection ?

- High computational efficiency, Single neuron ~1MMAC/pW
- Processing and memory operations performed by the same components
- Self-organizing system
- Online learning
- Solving ill-structured problems
- Transfer learning
- Spiking/event driven communication, subthreshold computation

1MMAC/(s)/pW

Neuromorphic techniques

1MMAC/(s)/uW

Analog/ Compute-in-memory techniques

1MMAC/(s)/mW

Hasler 2016

Neuromorphic techniques will be disruptive to how we develop our computing systems

MMAC: Million Multiply Accumulates

# NEUROMORPHIC COMPUTING: DIVERSE SOLUTIONS

## Digital Neuromorphic



Intel Loihi/ Loihi 2.0

SNL hosts Intel's 50 million neural supercomputer



Scaled to a billion neurons

SpiNNaker/ SpiNNaker 2



IBM TrueNorth

ODIN (Open-source)

## Analog/Mixed-Signal



GT Neuron

INI, ETH Zurich



Stanford Neurogrid

NeuRRAM UCSD/Tsinghua

## Beyond CMOS Devices



Mott- Memristor

ECRAM



RRAM Crossbar

MTJ



Brink et al., 2013

# NEUROMORPHIC BUILDING BLOCKS

**Analog Crossbars using NVMs**



Neuromorphic offers computational richness we can leverage, to move beyond today's computational limitations.

**Winner-Take-All**



Lazzaro et al. 1988

**Dendritic Processing**



George Cardwell et al. 2013

**Neural Path Planning**



Koziol et al. 2013

**Random Walks**

Smith et al. 2021



**APPLICATIONS**

AI/ML (ANN, SNN)

Brain-inspired algorithms

Scientific Computing

Edge Computing

**Learning Synapses**



Many different models for neurons, synapses, online learning and dendrites.

**Silicon Cochlea**



Liu et al. 2020

**Silicon Retina/ Event Sensor**



Delbruck et al. 2020

Posch et al. 2014

# CHALLENGE: SCALABILITY VS. COMPLEXITY



Dragonfly Brain
1 Million Neurons

Mouse Brain
(100 Million)

Human Brain
(100 Billion)

Analog
Neuromorphic

Stanford NeuroGrid
1 Million Neurons

Huge
Gap

IBM TrueNorth
1 Millions

Intel Loihi
100 Million
neurons

Biological Complexity

Scalability (# of Neurons, Synapses)

However, to achieve brain-like complexity we need both scaling and rich dynamics.

- Solving ill-structured problems
- Online learning
- Transfer learning

Understanding fundamental mechanisms in neuroscience, translated to algorithms and models will influence next-generation devices, architectures and intelligent computing systems

# INCREASING "BIOLOGICAL COMPLEXITY"

Increase computational efficiency
and
Increase computational density

Novel devices and materials
can help bridge this gap.

**LIF neuron**
- Single passive compartment
- Spikes
- Limited dynamics
- Relatively easy to scale

**Biological neuron**
- Dendrites = intricate structure and dense connectivity
- Complex pattern of active conductances
- Rich dynamics , multiple patterns of spiking, subthreshold computation
- More computational power, not compact

LIF: leaky Integrate and Fire

# DENDRITIC TOOLKIT FOR COMPUTATION

Dendrites are tree-like structures that connect neurons synapses to its soma.

**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

# SINGLE NEURON MULTIPLICATION



Groschner et al., Nature 2022

Leveraging Inhibition

Shunting Inhibition/
Leveraging Leakage
Conductance



Schemmel, Johannes, et al. , IEEE IJCNN, 2017.



Dendrites with Active Channel, Ramakrishnan et al., IEEE TBIOCAS, 2013.

NMDA /Ca



Lobula giant movement detector (LGMD) of locusts
Gabbiani et al., Nature 2002

Multiplication based on dendritic subtraction of two converging inputs encoded logarithmically, followed by exponentiation through active membrane conductances.

# SINGLE NEURON MULTIPLICATION

*Algorithms*

*Devices & Circuits*

*Physics of Computing*

from fan-shaped body of *Drosophila* brain

$R = \mathbf{A}*f(x)$

$R = f(x) - \mathbf{A}$

Lu et al. 2022

Chance & Cardwell
NICE 2023

Shunting Inhibition in
Neuromorphic Dendrite

# NEUROMORPHIC CODESIGN

*Algorithms*  *Devices & Circuits*  *Physics of Computing*

## Coordinate transformations from Dragonflies to Neuromorphic Hardware

Lead PI: Frances Chance, SNL

*October 2021*

Gonzalez-Bellido, UMN

**DRAGONFLY EXPERIMENTS**

Fovea Layer    Sensory Layer    Motor Layer

Retinal Layer

Chance 2020

**IEEE Spectrum**

The Brain of A Tiny Hunter

*August 2021*

**COMPUTATIONAL MODEL**

GT FPAA        Intel's Loihi

George Cardwell 2016        Davies 2018

SNL, Baylor

**NEUROMORPHIC IMPLEMENTATION**

Increased collaboration between neuroscience and neuromorphic engineering will facilitate  development of novel neural-inspired architectures.

U.S. DEPARTMENT OF ENERGY | Office of Science

DOE ASCR (FY21-24)

Department of Energy
Advanced Scientific Computing Research

# DRAGONFLY WITH DENDRITES

Excitatory
Input

Shunt Input ⟶ ⟶ out

$V_{axial}$ | $R_{axial}$

$V_{mem}$

$V_{leakage}$
$R_{leakage}$

$C_{leakage}$

Sub-threshold region of
operation

Visual
Input

Sensorimotor
Dendrites

Motor
Commands

Fovea
Shunt



— Leakage Variation w/o SI Vv=0.5
— Leakage Variation w/ SI Vv=0.4
— Leakage Variation w/ SI Vv=0.3
— Leakage Variation w/ SI Vv=0.2
— Leakage Variation w/ SI Vv=0.1

out (V)

peak input (A)

Sensorimotor
Dendrites

Visual
Input

Motor
Commands

Fovea
Shunt

Multiplication in
Dendrites

Sensorimotor
Dendrites Response

$$S_{ij} = f_i(x)g_j(y)$$

Cardwell & Chance
ICONS 23

# DIRECTION-SELECTIVE DENDRITES



Inputs

Compartment 0    Compartment 1    Compartment Destination

Pattern 1: UP

Pattern 2: DOWN

Spike Generator

Event Sensor Output

Nahuku (Loihi Chips)
Davies 2018

Cardwell & Chance ICONS 23

A. Input Spikes

DOWN

UP

B. Compartment 0 Spikes

C. Compartment 1 Spikes

D. Destination Compartment Current

E. Destination Compartment Voltage

F. Destination Spikes

Loihi TimeStep

# DIRECTION-SELECTIVE CELLS FOR COMPLEX PATTERNS

MT (Middle Temporal)

MT+

dorsal medial superior temporal

V1

Optic Flow Estimation

Directional Selectivity

Speed Selectivity

Slow speeds → Fast Speeds

MSTd

Radial Cells

Steinmetz et al. 2022

# CHALLENGE: CODESIGN TOOLS

| Co-Design Tools for Novel Architectures | Next-generation Neuromorphic Architectures |

**ATHENA**
Analytical Tool for analog and neuromorphic ML accelerator

**SANA-FE**
Neuromorphic Architecture Exploration

**AI-Enhanced Codesign**
Reinforcement Learning/Evolutionary for Circuit and System design

**COINFLIPS**
Probabilistic Neural Computing, Leverage stochasticity in beyond-CMOS devices

**DRAGONFLY**
Dendritic processing, Coordinate transformation from Dragonflies to Neuromorphic hardware,

ASC-AML (FY20-22)

SNL LDRDs (FY21-24)

SNL LDRDs (FY21-23)

DOE ASCR/BES (FY21-24)

DOE ASCR (FY21-24)

External Collaborators: UT Austin, Intel, NCSU, Infineon Memory Solutions, Georgia Tech, UMN, Baylor University, UT Knoxville, Temple University, NYU, ORNL

# 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.

- In the process of making ATHENA open-source.

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



**SANA-FE**
UT Austin Collaboration

Configuration & Input Spikes

**Hardware Simulator**

**build architecture**
**initialize network**

**for all timesteps:**
   **get external inputs**
   **for all neurons**
      **process neuron**
      **receive messages**

**estimate energy, latency**

Performance & Energy Estimates

Architecture Description

Mapped Spiking Neural Network

Boyle et al., ICONS 2023

# SPIKING ARCHITECTURE TEMPLATE

## Tile-based architecture
- Network-on-chip connecting neural cores

## Many cores per tile
- Cores simulate group of mapped neurons
- Local shared memory

## Core pipeline
- Axon stage
- Synapse stage
- Dendrite stage
- Soma stage

# CHALLENGE: ARCHITECTURE DESCRIPTION

Describes different H/W architectures

- Represents different existing & future spiking designs based on common features
- Defines compute elements of chip
- YAML-based, flexible & extensible



```
architecture:
 name: demo
 tile:
  - name: demo_tile[0..7]
    attributes:
      energy_east_west: 1e-12
      latency_east_west: 2e-9
      ...
    core:
      - name: demo_core[0..3]
        soma:
          - name: core_lif
            attributes:
              energy_spiking: 68e-12
              latency_spiking: 30e-9
...
```

Boyle et al., ICONS 2023

# PERFORMANCE MODELING RESULTS

For randomized spiking inputs on the application SNN



- Detailed breakdown of on-chip activity on Loihi
- Captures dynamic energy and latency trends

➢ Detailed insight into H/W behavior

Boyle et al., ICONS 2023

# RESULTS FOR OTHER NEUROMORPHIC BENCHMARKS

Predict performance & energy for larger real-world neuromorphic applications

- SNN trained on DVS gesture data-set
- 18,678 neurons across 6 layers
- Mapped to 45 Loihi cores out of 128



**Image reproduced from [Massa,'20]**

# SIMULATOR SPEED RESULTS

- Compared to existing spiking simulator (NeMo)



- Simulating IBM TrueNorth architecture
- Randomized SNN with 80% of spikes intra-core, 20% spikes between cores
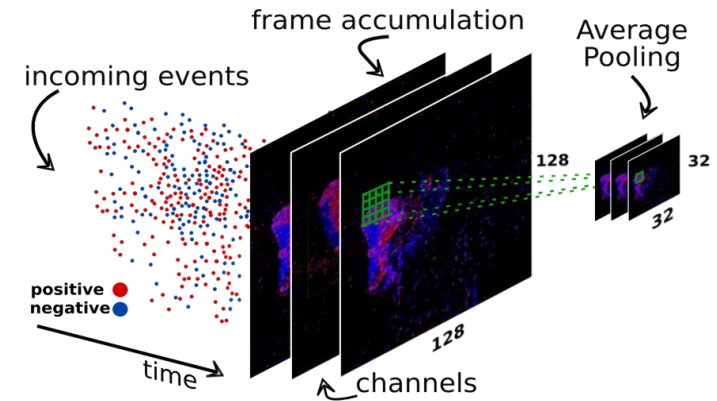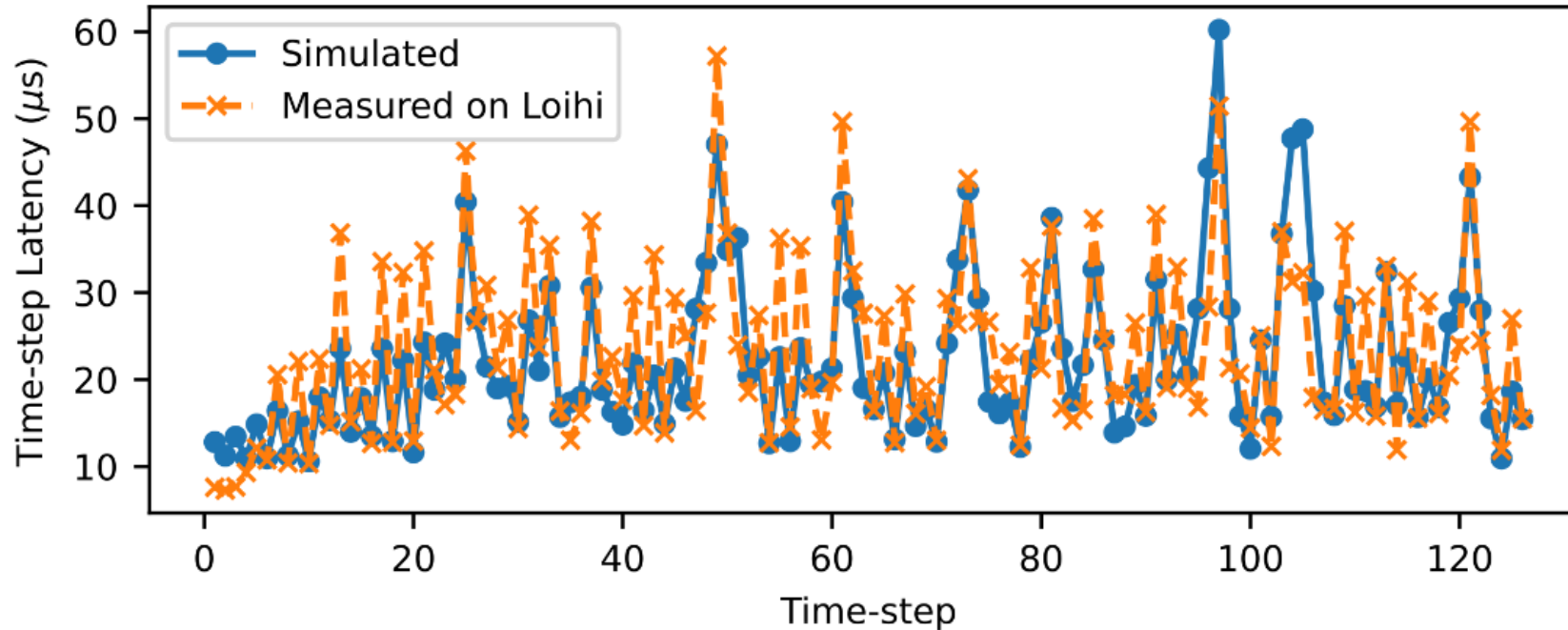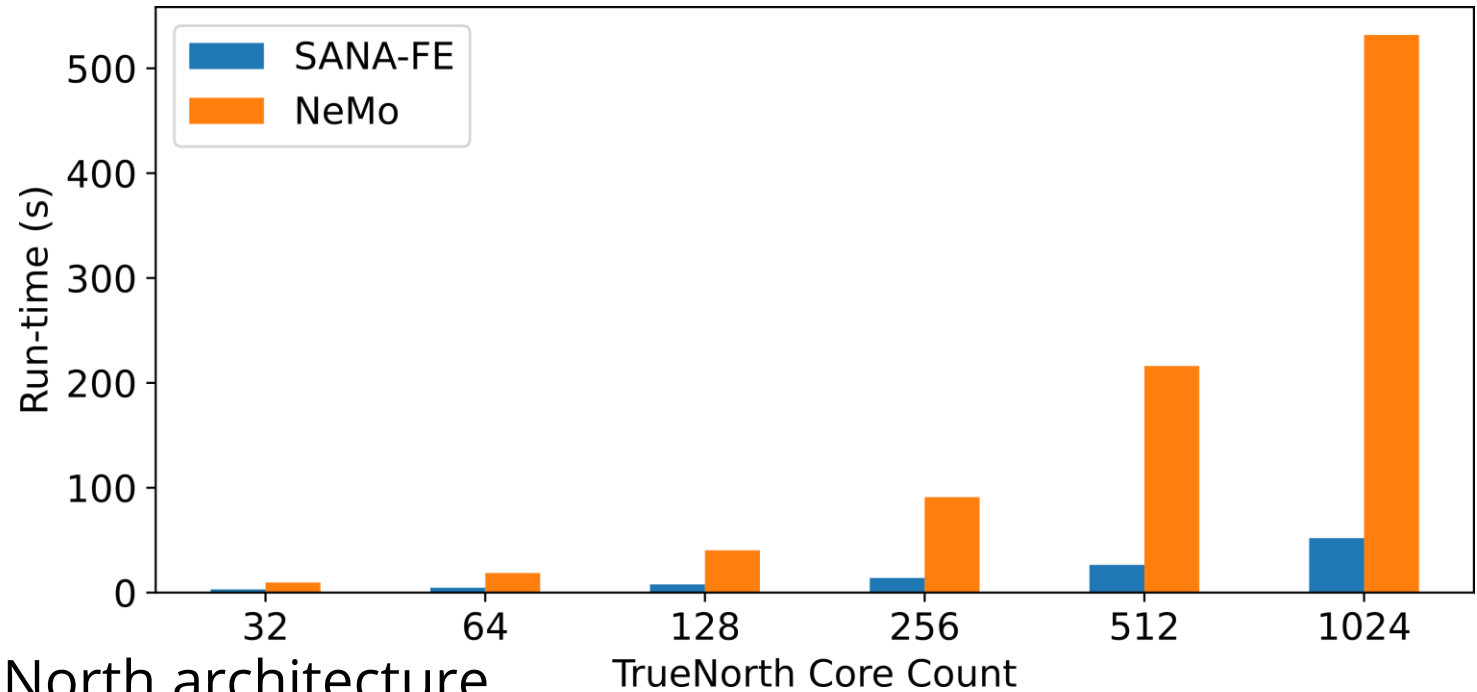- **Over 10x faster than NeMo for 1024 cores**

# SANA-FE

**Generic & extensible**
- User-defined architecture & SNN
- Supports range of spiking architectures

**Fast & accurate**
- Time-step based approach
- Detailed hardware activity for each time-step
- Accurately estimates performance & energy

**Future work**
- Support other existing architectures & scale to larger designs
- Adapt other neuromorphic benchmark applications
- Model analog architectures & novel devices
- Integrate with other frameworks e.g., SST, Fugu & Lava

Access at: https://github.com/SLAM-Lab/sana-fe

Prof. Andreas Gerstlauer's SLAM Lab @ UT Austin



Boyle et al., ICONS 2023

# AI-ENHANCED CODESIGN: COINFLIPS

Lead PI: Brad Aimone



**Microelectronics Codesign
Award DOE ASCR/BES (FY21-24)**
Department of Energy
Advanced Scientific Computing Research
Basic Energy Sciences

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

https://coinflipscomputing.org/

**AI-Guided Approach**

Unfair coins can be combined with AI-designed neural circuits to allow sampling of application desired probability distributions, avoiding accept/reject steps.

We leveraged evolutionary algorithms for circuit design and optimization

- **LEAP** (Library of Evolutionary Algorithms in Python)
- **EONS** (Evolutionary Optimization for Neuromorphic Systems)- Schuman et al. , 2020

We used abstracted device models for TD and MTJ to capture functionality and energy usage.

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

# AI-GUIDED CODESIGN OF PROBABILSITIC CIRCUITS



$$\mathbb{P}[\text{Coin } 1 = H \text{ and Coin } 2 = H] = \frac{1}{2}$$
$$\mathbb{P}[\text{Coin } 1 = H \text{ and Coin } 2 = T] = \frac{1}{6}$$
$$\mathbb{P}[\text{Coin } 1 = T \text{ and Coin } 2 = H] = \frac{1}{6}$$
$$\mathbb{P}[\text{Coin } 1 = T \text{ and Coin } 2 = H] = \frac{1}{6}$$

Probabilistic Mixing Algorithm

Optimized weight values for each device

Multi-objective optimization of weights of fitness function for optimal KL divergence, biased weight and energy usage.

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

Weights are customized for the device's behavior to target the best performance in terms of KL divergence and energy usage.

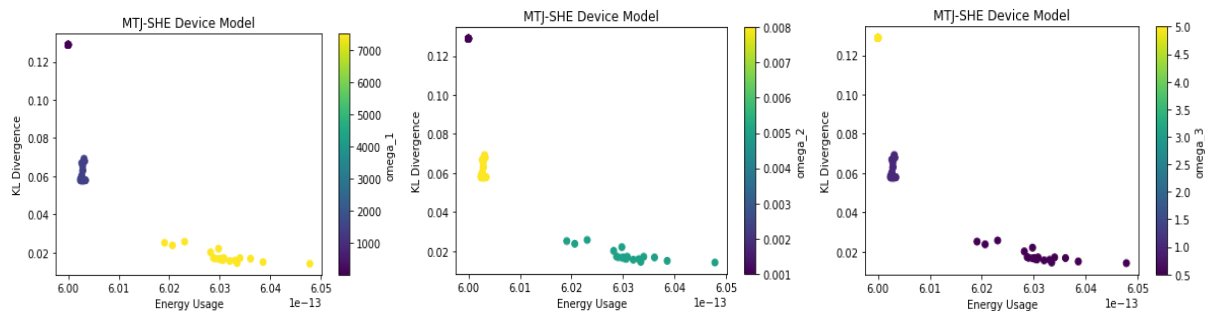One of the challenges in optimizing for both algorithms and devices was appropriately abstracting the device models and algorithmic constraints.

The functional models developed will also be evolved in time as new device data and research emerges.

Our framework can accommodate any emerging device type.

# AI-ENHANCED CODESIGN: NEURAL CIRCUITS

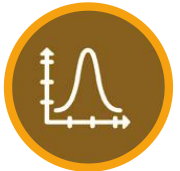We developed an RL algorithm approach which is capable of building very simple circuits.

**Simple delay line**

**Pattern Detection**

Training time to > 99% success

*Algorithms*

*Physics of Computing*

*Devices & Circuits*

Crowder et al. 2023 (MWCAS 2023)

# NEUROMORPHIC APPLICATIONS

**AI/ML Applications**
ANNs
SNNs

**Scientific Computing**
- Random Walks
High-fidelity Physics Simulations

**Edge Computing**
- Event sensors
- Spatio-temporal processing

**Brain-inspired Algorithms**
Dragonfly
Dendritic Processing

**Probabilistic Computing**
COINFLIPS

**Heterogeneous Computing Applications**

Novel Algorithm Complex Dynamics

Novel Architectures & Circuits

Software Stack and Tools

**CODESIGN**
Neuromorphic Computing

Leveraging Physics and Noise

AI-enhanced Codesign

Novel Devices and Materials

# CHALLENGES FOR NEXT-GENERATION OF NEUROMORPHIC SYSTEMS

- Algorithms are cognizant of architecture and device constraints.
- Leverage the complex dynamics of devices.
- Bio-inspired techniques, adoption in computing

- Heterogeneous architectures
- CoDesign to optimize communication and memory bottlenecks
- 3D architectures, Photonics

- Software tools to support design and development
- Integration with AI-enhanced techniques?
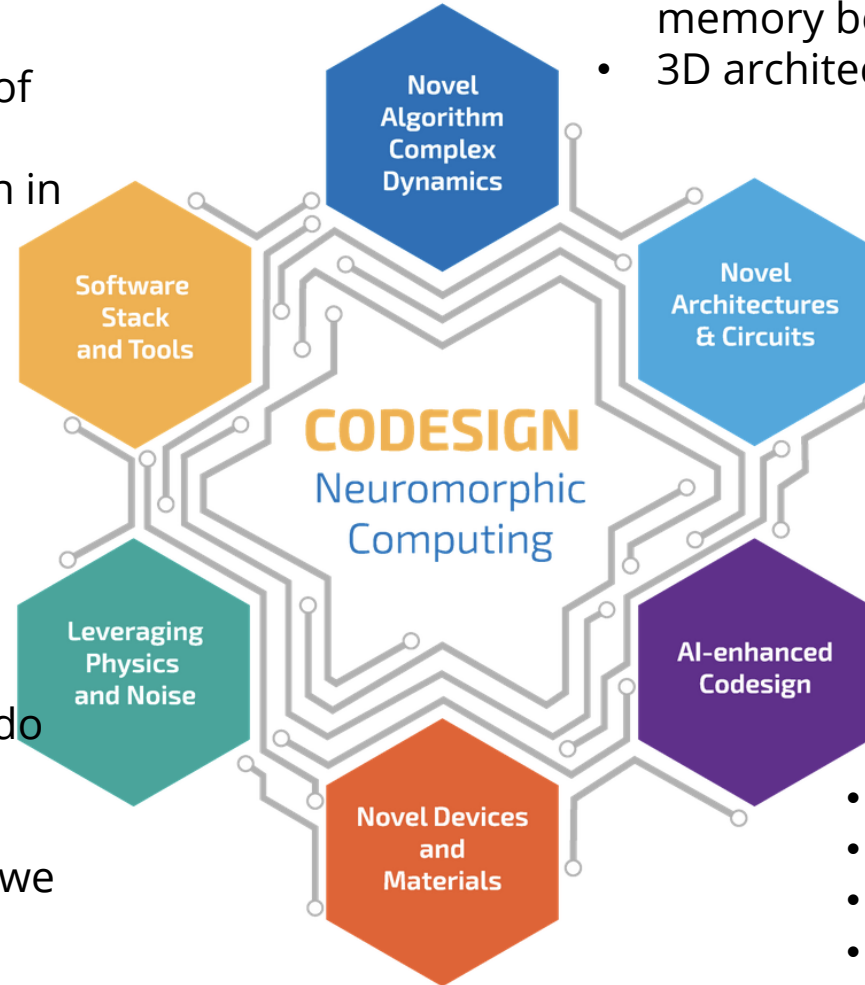
- How can AI-enhanced techniques accelerate scientific discovery?
- Different AI techniques at the device, circuit, system design and architecture level.
- Enable encoding of domain knowledge
- Enable concurrent contribution from researchers

- Leverage the physics of devices to do computation (analog)
- Embrace stochasticity of devices
- Analog devices are noisy. How can we incorporate this into algorithms?

**CODESIGN**
Neuromorphic Computing

Novel Algorithm Complex Dynamics

Novel Architectures & Circuits

Software Stack and Tools

Leveraging Physics and Noise

Novel Devices and Materials

AI-enhanced Codesign

- Novel devices with complex dynamics
- Radiation-hardened devices
- Reconfigurable devices
- Computational efficiency and computational density

# THANK YOU!!



**Neural Exploration and Research Lab**

## WE ARE HIRING!



**Careers**

**careers.sandia.gov**

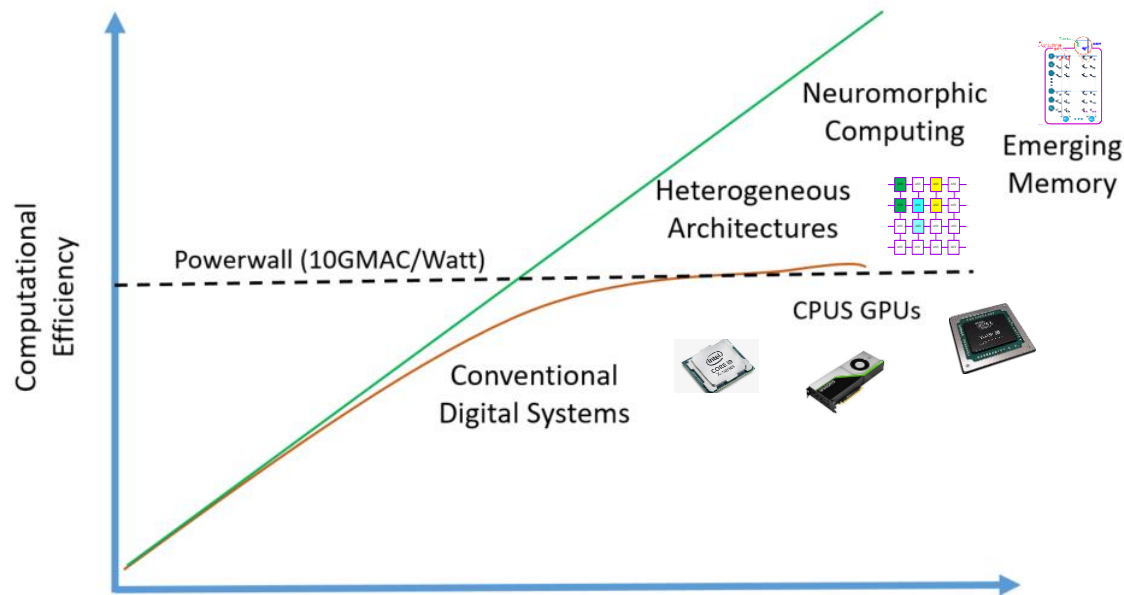**https://neuroscience.sandia.gov/**

Suma G. Cardwell sgcardw@sandia.gov

# FUTURE OF COMPUTING: HETEROGENEOUS ARCHITECTURES



Co-Design is critical to build the next-generation heterogeneous systems

Limits of scaling have ushered in the 'Golden Age of Computer Architecture'

Heterogeneous Architectures

Emerging memory

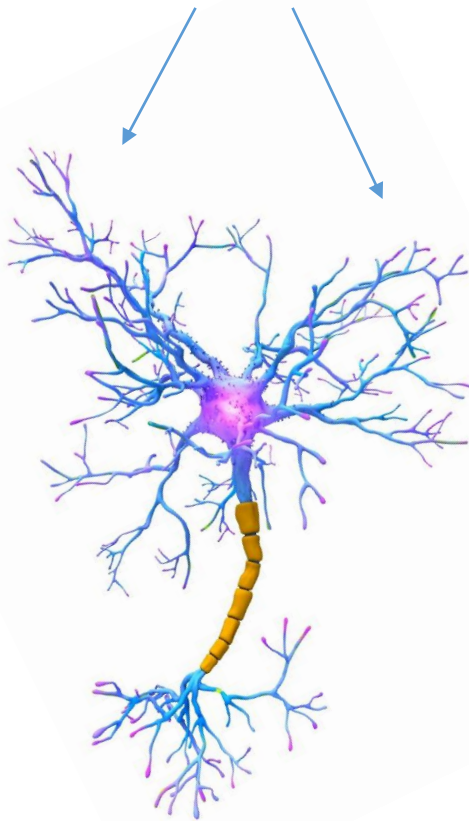Neuromorphic Accelerators

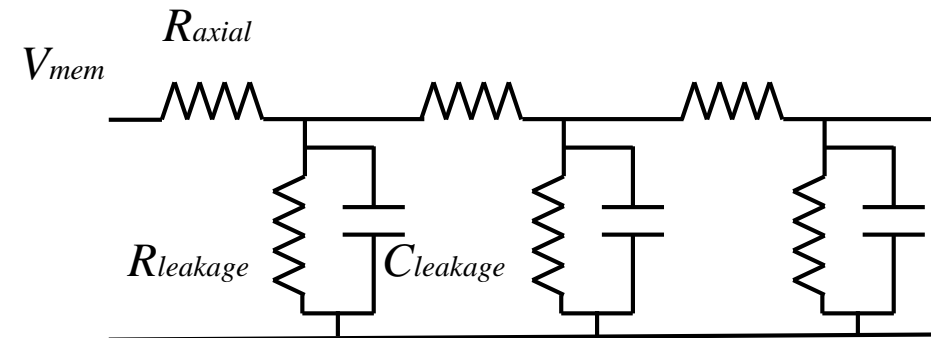Extremely Heterogeneous

Quantum Computing

5-10 years

15-20 years

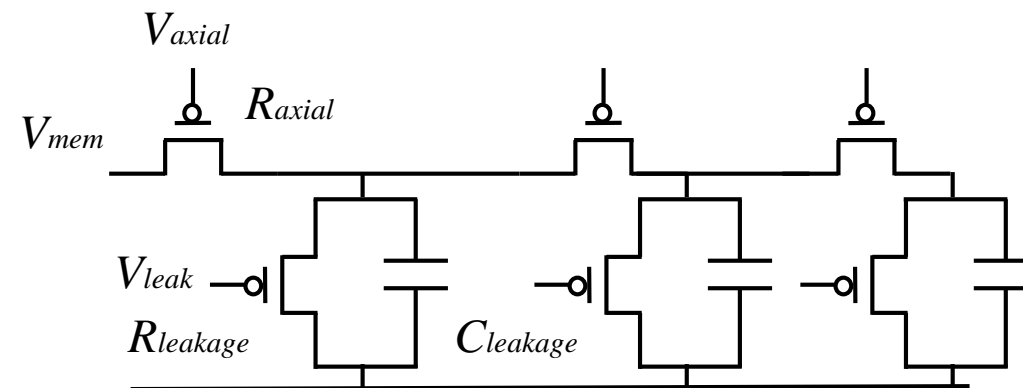'Truly Heterogeneous Computing', Cardwell et al., SMC 2020

# DENDRITE MODELING

## DENDRITES
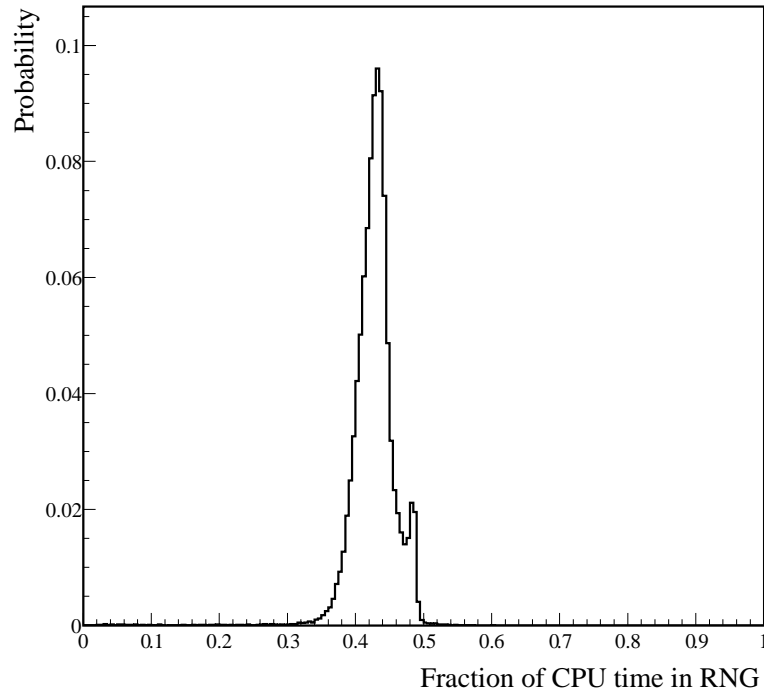
## Resistor-Capacitor Circuit



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

$R_{leakage}$    $C_{leakage}$

Rall's Cable Model

## CMOS-transistor based Dendrite

$V_{axial}$

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

$V_{leak}$

$R_{leakage}$    $C_{leakage}$

Nease et al. 2011

# COINFLIPS APPLICATION: NUCLEAR PHYSICS SIMULATIONS



Fraction of CPU time in RNG
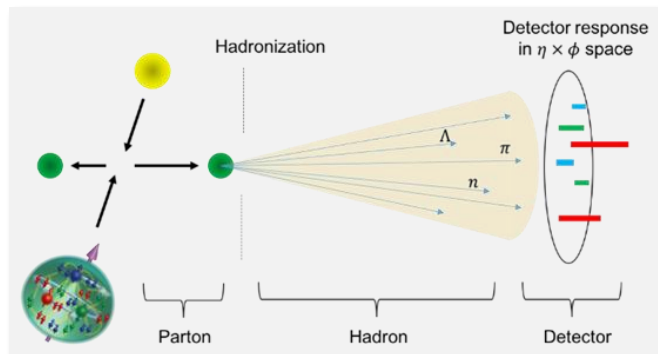


For a particular collider physics simulation [Pierog et al., Phy Rev. 2022], ~ 270K pseudo- random numbers needed for a single event, with billions of events needing to be simulated.

CPU time is ~ 30-50% of the total compute time

Direct random number generation leveraging stochastic devices can promise significant energy savings for such applications

Misra et al., *Advanced Materials 2022*

Random numbers are a limiting computational cost for some nuclear physics applications