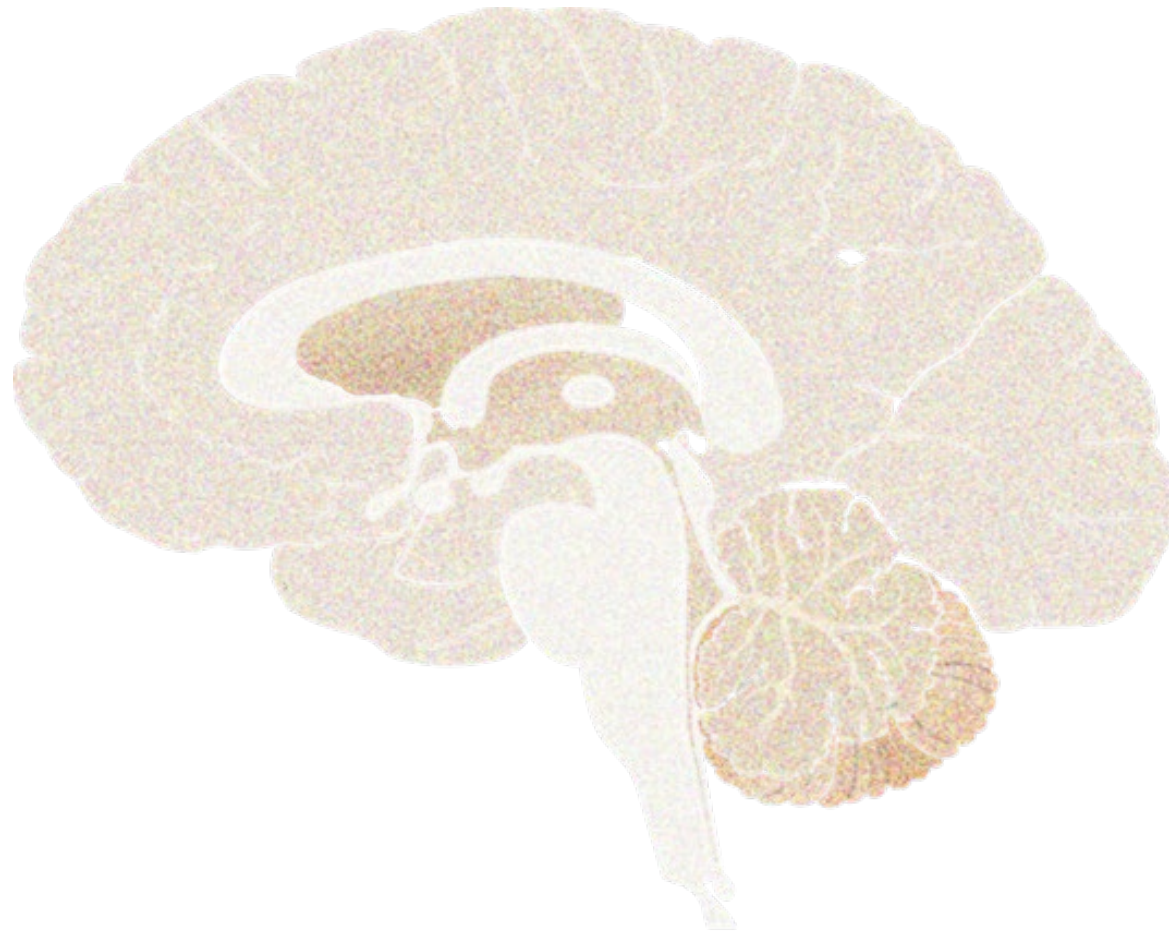




# A THEORETICAL FRAMEWORK FOR TIME, SPACE, AND ENERGY SCALING OF NEUROMORPHIC ALGORITHMS

Brad Aimone

Center for Computing Research  
jbaimon@sandia.gov



We don't know how this works. Or for that matter, we don't even know what it is doing....

Hypothesis: Anything the brain does should be consistent with what the brain is.

**The algorithm must match the architecture.**



This is our unknown



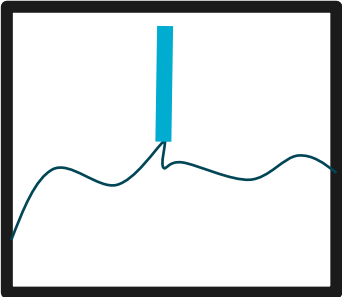
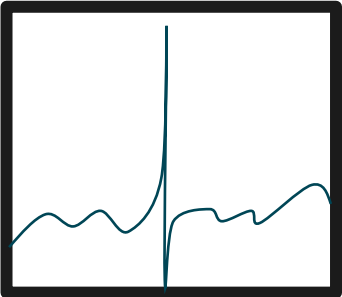
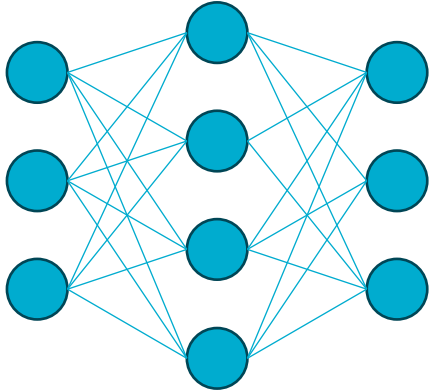
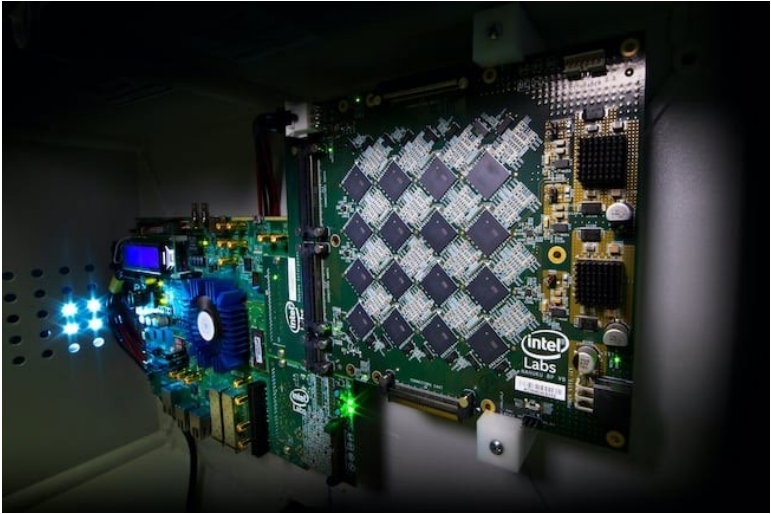
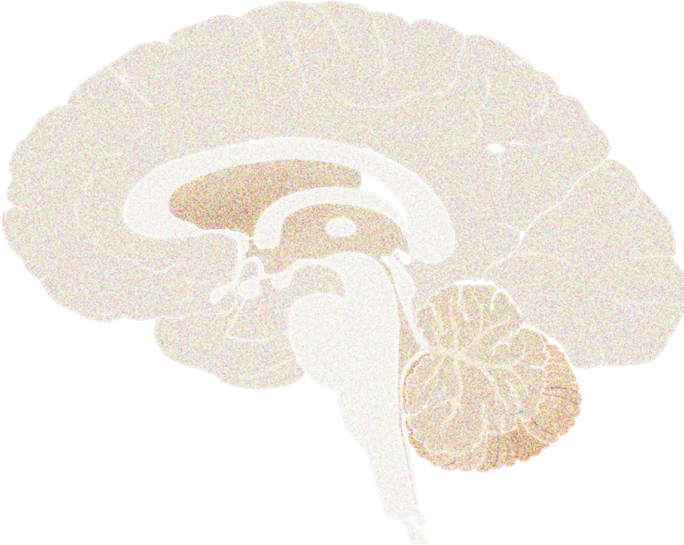
The **algorithm** must match the **architecture**.



Can we approximate this?

Neuromorphic computing is a potential approximation.  
We need to better define what neuromorphic computing is.

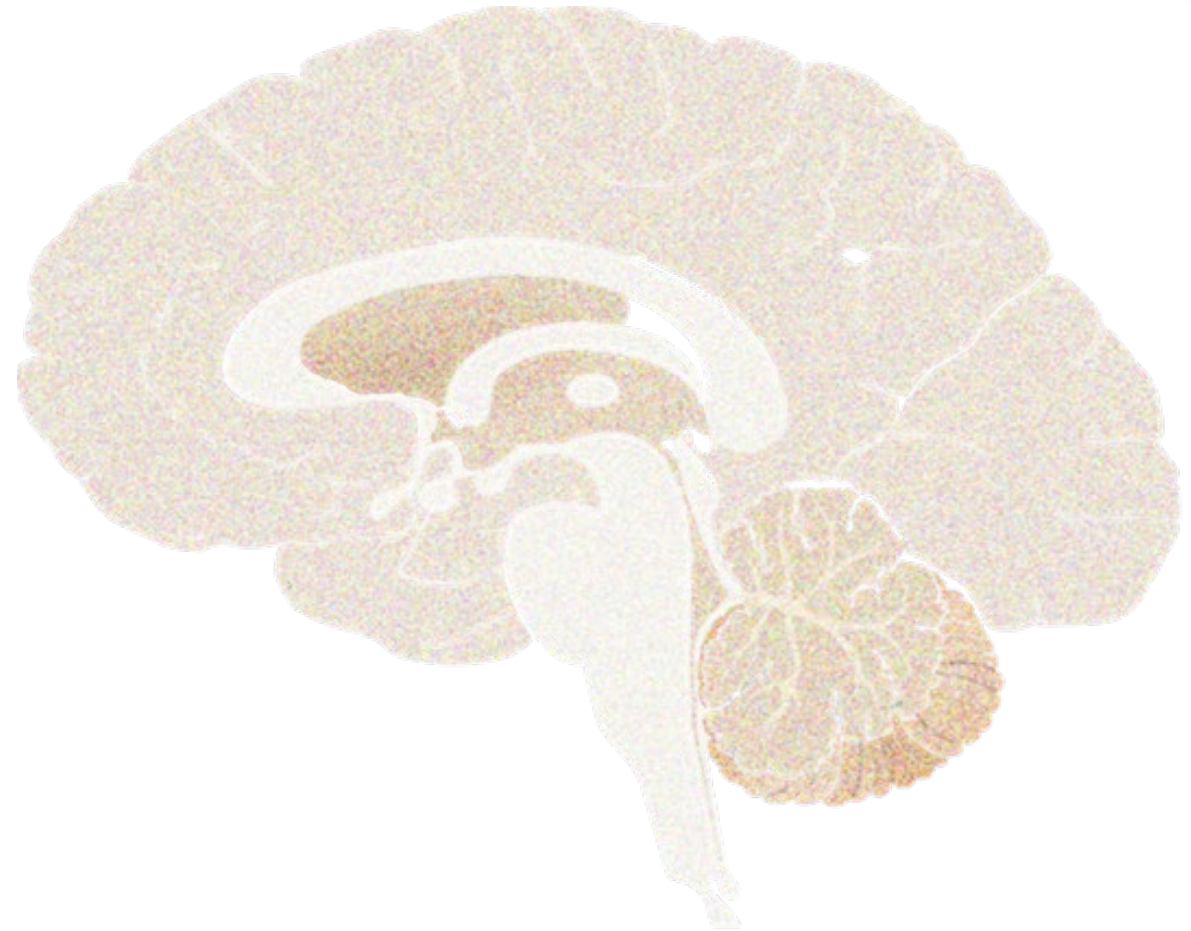
# ARE WE READY FOR A COMMON FRAMEWORK FOR ANY NEURAL COMPUTATION?



# ARE WE READY FOR A COMMON FRAMEWORK FOR ANY NEURAL COMPUTATION?



- ✓ The brain is an existence proof



# ARE WE READY FOR A COMMON FRAMEWORK FOR ANY NEURAL COMPUTATION?



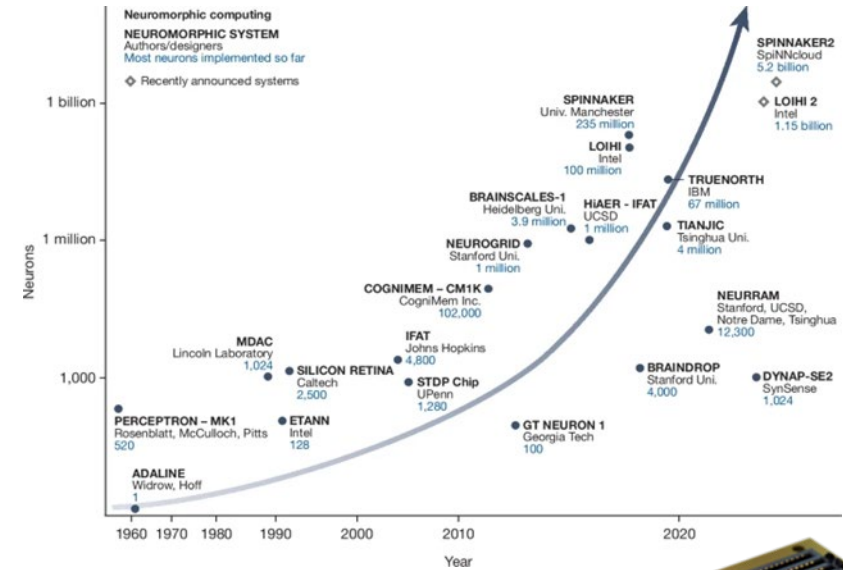
- ✓ The brain is an existence proof
- ✗ The brain is hard to understand



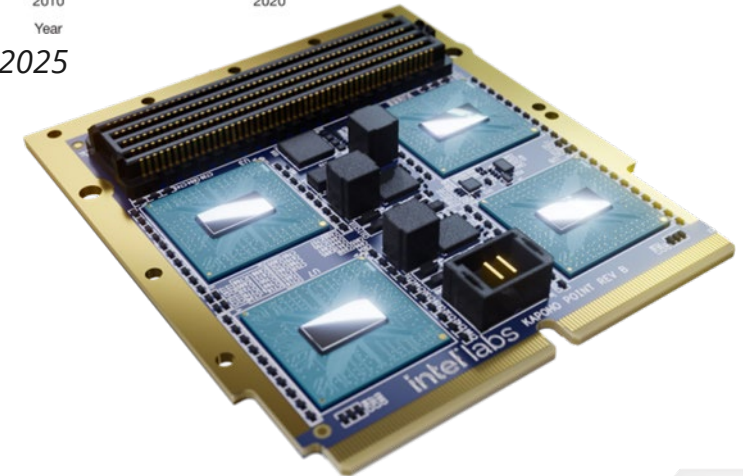
# ARE WE READY FOR A COMMON FRAMEWORK FOR ANY NEURAL COMPUTATION?



- ✓ The brain is an existence proof
- ✗ The brain is hard to understand
- ✓ Scalable hardware exists



Kudithipudi et al, Nature 2025



# ARE WE READY FOR A COMMON FRAMEWORK FOR ANY NEURAL COMPUTATION?



- ✓ The brain is an existence proof
- ✗ The brain is hard to understand
- ✓ Scalable hardware exists
- ✗ People don't really know how to program it

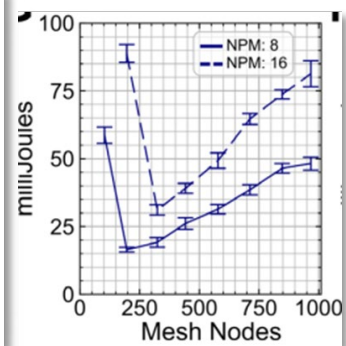
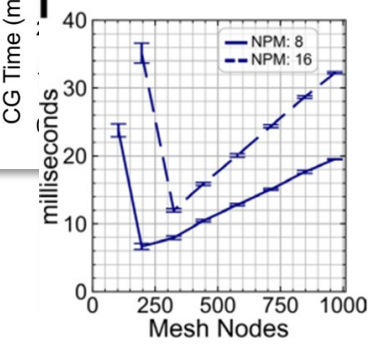
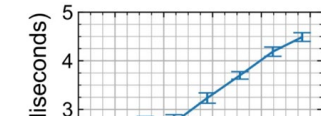
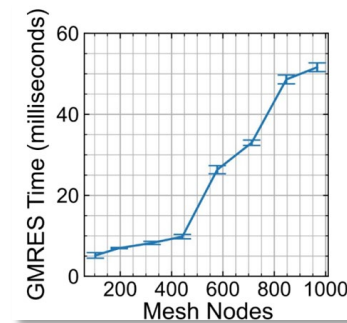
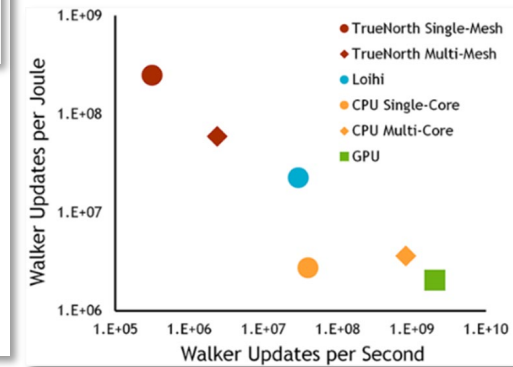
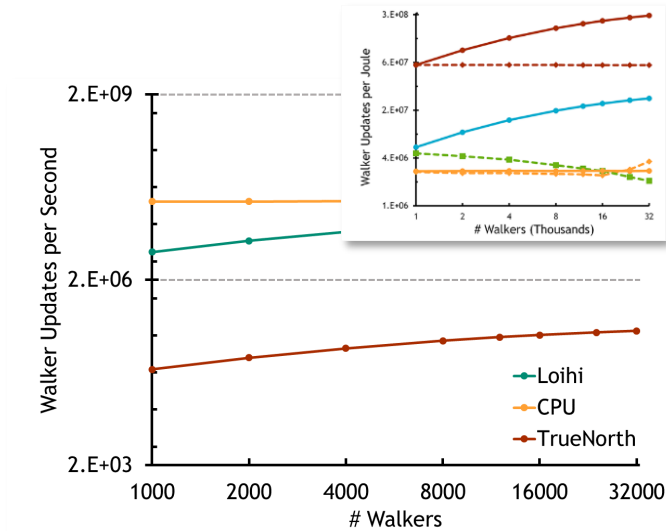




# ARE WE READY FOR A COMMON FRAMEWORK FOR ANY NEURAL COMPUTATION?



- ✓ The brain is an existence proof
- ✗ The brain is hard to understand
- ✓ Scalable hardware exists
- ✗ People don't really know how to program it
- ✓ There are algorithm successes
- ✗ Results are largely heuristic and empirical



# A START TOWARDS A THEORETICAL FRAMEWORK FOR NEUROMORPHIC COMPUTING



- *Challenge:* Neuromorphic computing has been hardware motivated. As such, the field lacks the theoretical foundation that conventional and quantum computing have.
  - There is no widely agreed upon “model of neural computation”
  - There is no standard method for formal analysis of neuromorphic algorithms, architectures, etc.

## A Theoretical Framework for Time, Space, and Energy Scaling of Neuromorphic Algorithms

James B Aimone  
Neural Exploration & Research Laboratory  
Sandia National Laboratories  
Albuquerque, NM 87185  
jbaimon@sandia.gov

Revised March 2026

### Abstract

Neuromorphic computing (NMC) is increasingly viewed as a low-power alternative to conventional von Neumann architectures such as central processing units (CPUs) and graphics processing units (GPUs), however the computational value proposition has been difficult to define precisely.

Here, we propose a computational framework for analyzing NMC algorithms and architectures. Using this framework, we demonstrate that NMC can be analyzed as general-purpose and programmable even though it differs considerably from a conventional stored-program architecture. We show that the time and space scaling of idealized NMC has comparable time and footprint tradeoffs that align with that of a theoretically infinite processor conventional system. In contrast, energy scaling for NMC is significantly different than conventional systems, as NMC energy costs are event-driven. Using this framework, we show that while energy in conventional systems is largely determined by the scheduled operations determined by the structural algorithm graph, the energy of neuromorphic systems scales with the activity of the algorithm, that is the activity trace of the algorithm graph. Without making strong assumptions on NMC or conventional costs, we demonstrate which neuromorphic algorithm formulations can exhibit asymptotically improved energy scaling when activity is sparse and decaying over time. We further use these results to identify which broad algorithm families are more or less suitable for NMC approaches.

### 1 Introduction

Neuromorphic computing (NMC) approaches to computation have been proposed for many years [1–3], and today NMC hardware is increasingly available and can be implemented at scales approaching  $10^9$  neurons in silicon [4], although still these systems remain far from the brain’s complexity [5]. One implication of the success of implementing large-scale NMC systems is that it is increasingly evident that the biggest challenge facing the neuromorphic field, at least in the near-term, is the identification of which applications are well suited for its use. While scalable NMC platforms were originally motivated, in part, by large-scale brain simulations [6–9], the growing need for low-power solutions at both the edge and in high-performance computing systems has increased interest in NMC to address energy challenges in computation broadly.

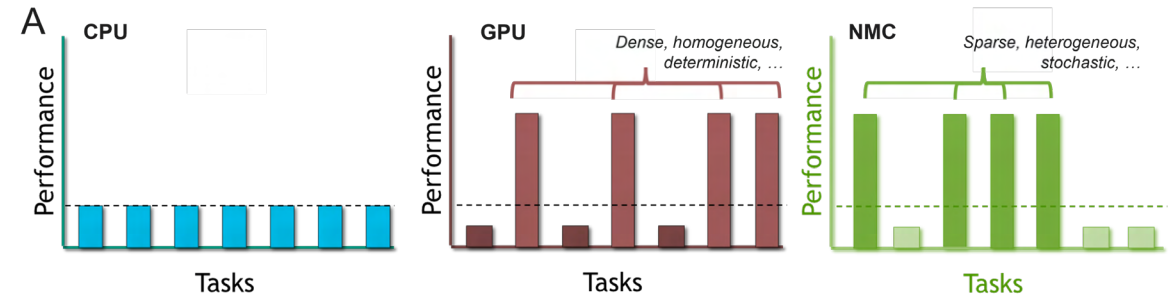
Key to their value proposition is that today’s digital NMC systems are generically programmable—aside from some constraints of fan-in/fan-out, these systems can implement arbitrary computational graphs. Conceptually, this compatibility includes any algorithm that can be formulated as a threshold gate (TG) circuit, artificial neural network (ANNs), or any other more complex ensemble of neurons. As the ability to implement TGs confers Turing completeness [10], and ANNs confer universal function approximation [11], we can deduce that NMC is both universal in terms of algorithm compatibility and its ability to approximate functions.

Given this universality, the relevant question of NMC is not “*can* neuromorphic solve this task?”, but rather “*should* NMC be used to solve this task?”. While NMC is by definition general purpose in potential, based on the ‘No Free Lunch’ theorem as applied to computer architectures, it should be expected that it will

# HOW DOES NMC RELATE TO GPUS AND OTHER VON NEUMANN ARCHITECTURES?



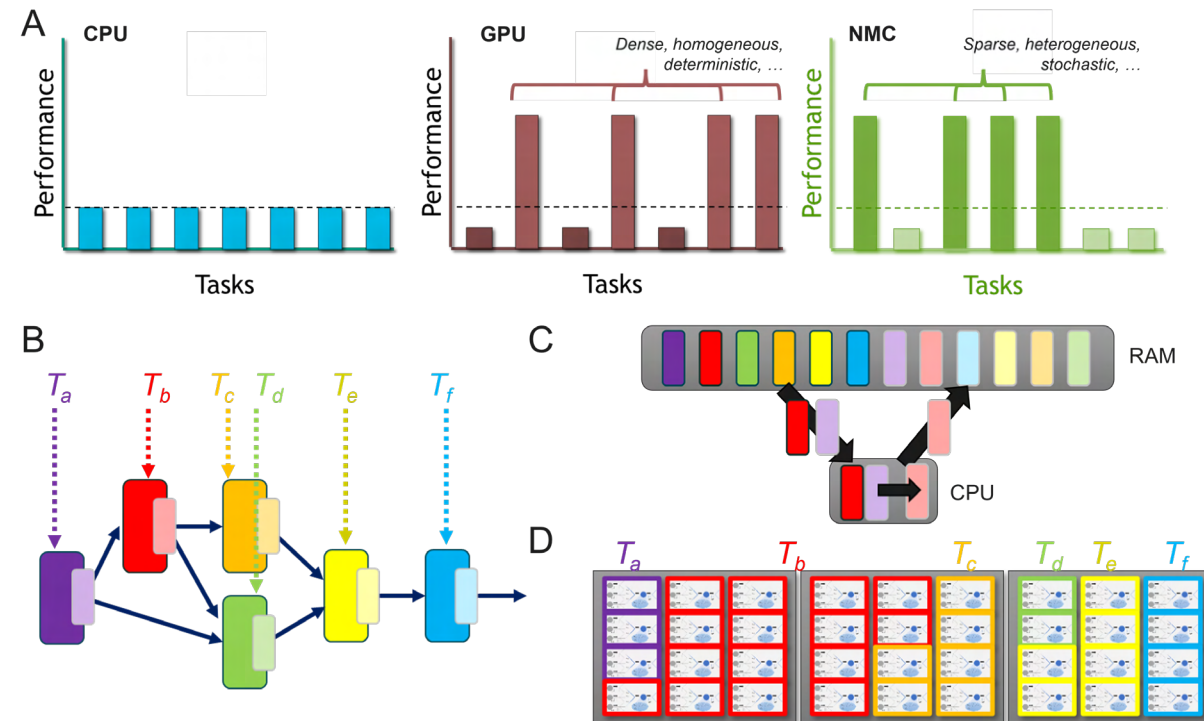
- Expectation that neuromorphic systems will excel at some tasks but be worse at others
  - “Specialized general purpose”, like GPUs but specialized to different tasks
  - “No Free Lunch”



# HOW DOES NMC RELATE TO GPUS AND OTHER VON NEUMANN ARCHITECTURES?



- Expectation that neuromorphic systems will excel at some tasks but be worse at others
  - “Specialized general purpose”, like GPUs but specialized to different tasks
  - “No Free Lunch”
- Neuromorphic is fully *in memory computing*
  - No stored program
  - Allows event-driven, asynchronous operation



# OUTLINE

1. Formal definitions
2. First crack at a complexity analysis
3. Algorithm successes



# OUTLINE

1. Formal definitions
2. First crack at a complexity analysis
3. Algorithm successes



# FORMAL DEFINITION OF NEURONS, SYNAPSES, NEURAL ARCHITECTURES AND NEURAL ALGORITHMS



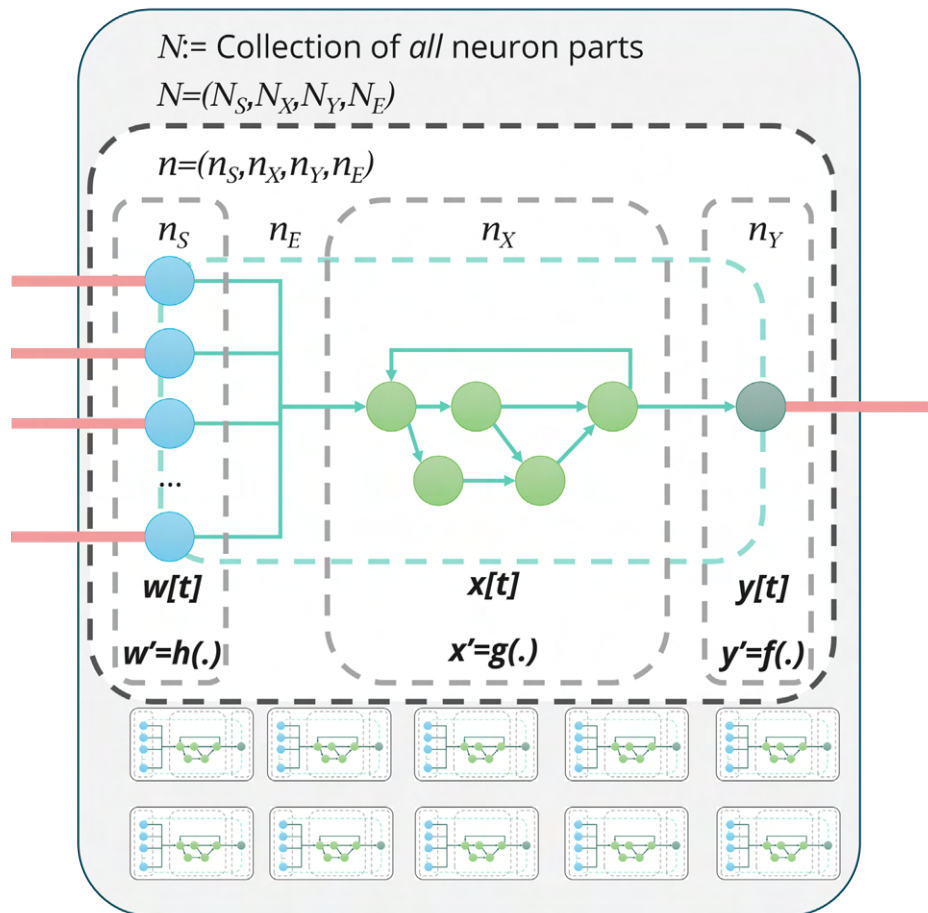
**Definition (Neuron).** A neuron is a compute element represented as a (potentially cyclic) directed subgraph  $n = (n_S, n_X, n_Y, n_E)$ , where:

- $n_S$ : Synaptic input nodes, representing incoming spiking signals from other neurons ( $s[t]$ ) and exhibiting synaptic dynamics related to the representation ( $w[t]$ ), and evolution ( $h(\cdot)$ ) of synaptic weight, synaptic delay, stochasticity, and synaptic plasticity.
- $n_X$ : Internal compute nodes describing the neuron's internal state variables ( $x[t]$ ) and the neuron's internal state dynamics ( $g(\cdot)$ ).
- $n_Y$ : Output nodes, representing the neuron's observable output ( $y[t]$ ) and the generation of that output ( $f(\cdot)$ ). For spiking neurons,  $y[t] \in \{0, 1\}$ .
- $n_E$ : Internal edges, representing directed connections between  $n_S$ ,  $n_X$ , and  $n_Y$ . The internal edges  $n_E$  may contain self-loops and directed cycles to represent intrinsic dynamics (e.g., decay, refractoriness, or adaptation).

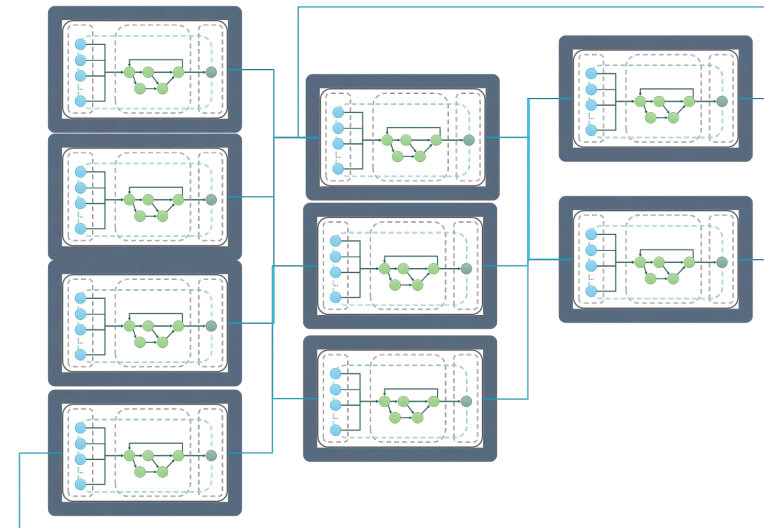
# FORMAL DEFINITION OF NEURONS, SYNAPSES, NEURAL ARCHITECTURES AND NEURAL ALGORITHMS



**Definition S1.1 (Neuron).** A neuron is a compute element represented as a (potentially cyclic) directed subgraph  $n = (n_S, n_X, n_Y, n_E)$ , where:



$\mathcal{N} :=$  Collection of *all* neurons



# FORMAL DEFINITION OF NEURONS, SYNAPSES, NEURAL ARCHITECTURES AND NEURAL ALGORITHMS



**Definition (Synapse).** Let  $N$  represent the set of isolated neuron subgraphs in a neuromorphic system, where each synapse is a directed synaptic communication channel from a source neuron  $n_i \in N$  to a target neuron  $n_j \in N$ . We take the set of synaptic communication channels to be a set of directed edges

$$S \subseteq N_Y \times N_S$$

Each synapse connects an output node  $n_{yi} \in n_{i,Y}$  of the source neuron  $n_i$  to a synaptic input node  $n_{sj} \in n_{j,S}$  of the target neuron  $n_j$  via a directed edge  $s_{ij} \in S$ , where:

- $n_{yi}$ : The output node of the source neuron  $n_i$ , representing the observable spiking output  $y_i[t] \in \{0, 1\}$ .
- $n_{sj}$ : The synaptic input node of the target neuron  $n_j$ , representing post-synaptic state associated with the connection, including synaptic weight  $w_{ij}[t]$  and any synaptic dynamics (e.g., delay, stochasticity, filtering, and plasticity).
- $s_{ij}$ : The directed edge connecting  $n_{yi}$  to  $n_{sj}$ , representing discrete spike-event transmission between neurons,  $s_{ij}[t] \in \{0, 1\}$ .

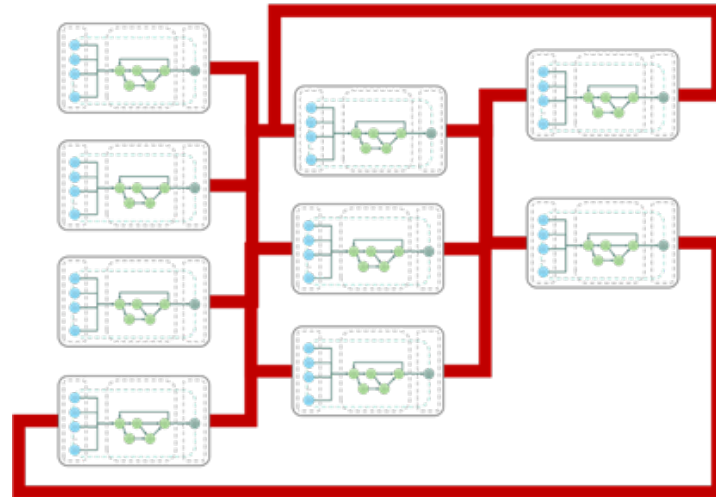
# FORMAL DEFINITION OF NEURONS, SYNAPSES, NEURAL ARCHITECTURES AND NEURAL ALGORITHMS



**Definition S1.2 (Synapse).** Let  $N$  represent the set of isolated neuron subgraphs in a neuromorphic system, where each synapse is a directed synaptic communication channel from a source neuron  $n_i \in N$  to a target neuron  $n_j \in N$ . We take the set of synaptic communication channels to be a set of directed edges

$$S \subseteq N_Y \times N_S$$

$S$  := Collection of *all* connections



# FORMAL DEFINITION OF NEURONS, SYNAPSES, NEURAL ARCHITECTURES AND NEURAL ALGORITHMS



**Definition (Neuromorphic Computing Architecture).** A neuromorphic computing architecture is a distributed computing system capable of representing a set of neurons  $N$  and a set of synaptic communication channels  $S$  as a directed graph  $G_{N,arch} = (N_{arch}, S_{arch})$ , where:

- $N_{arch}$ : The union of neuron components of the set of isolated neuron subgraphs  $N$ , where each neuron  $n \in N$  is represented as  $n = (n_S, n_X, n_Y, n_E)$ , as defined in Definition S1.1.
- $S_{arch}$ : The set of directed synaptic communication channels between neuron subgraphs, where each  $s_{ij} \in S_{arch}$  connects a source output node  $n_{yi}$  to a target synaptic input node  $n_{sij}$ , as defined in Definition S1.2

# FORMAL DEFINITION OF NEURONS, SYNAPSES, NEURAL ARCHITECTURES AND NEURAL ALGORITHMS



**Definition (Neuromorphic Algorithm).** A neuromorphic algorithm is an algorithm implemented on a neuromorphic computing architecture  $G_{N,arch} = (N_{arch}, S_{arch})$ . The algorithm is specified by (i) a selection of an algorithm subgraph  $G_{N,alg} = (N_{alg}, S_{alg})$  and (ii) a parameterization of the selected neuron-local state and update rules, where:

- $N_{alg}$ : The combined graph of selected neuron subgraphs, representing the computational nodes of the algorithm.
- $S_{alg}$ : The selected set of synaptic communication channels, representing the directed edges between distinct neurons. The algorithm graph  $G_{N,alg} = (N_{alg}, S_{alg})$  represents the ideal computational graph for the algorithm, with  $N_{alg} \subseteq N_{arch}$ ,  $S_{alg} \subseteq S_{arch}$ .

# UPSHOT OF FRAMEWORK DEFINITIONS



- Clarifies neuromorphic algorithm [ $G_N = (N, S)$ ] resource count
  - $N$  (neurons) captures all computation
    - Neurons can have arbitrary complexity, but this is likely a constant cost
    - *Opens up neuron complexity vs algorithm scaling direction*
  - $S$  (synapses) captures the non-local communication
    - If fan-out/fan-in is bounded, then synapse count scales no worse than neuron count
    - *Allows us to isolate distinct energy costs*

# OUTLINE

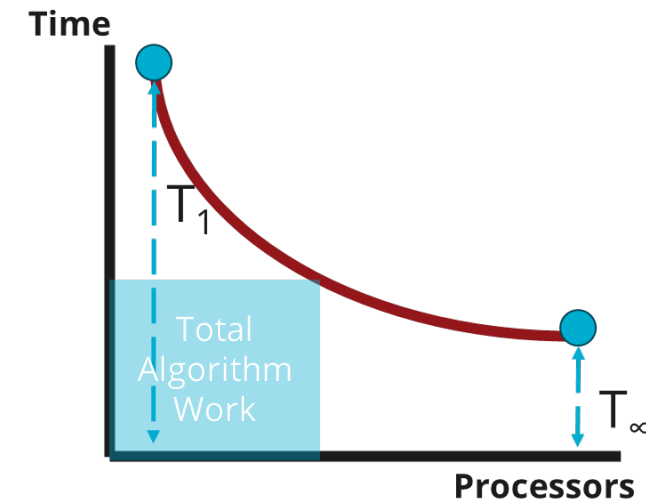
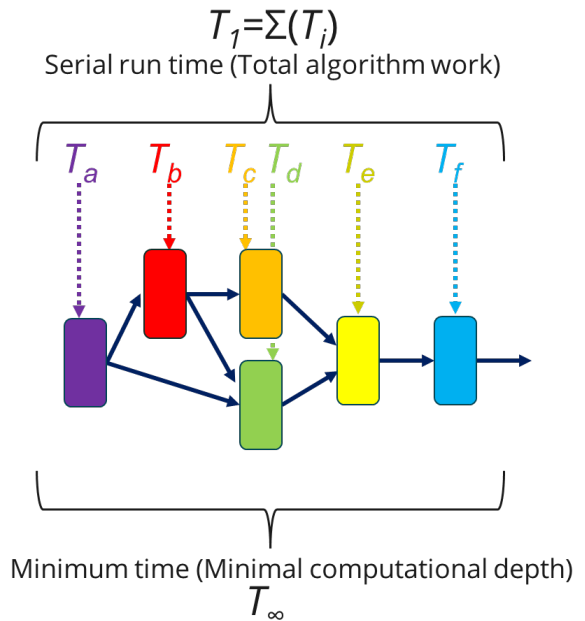
1. Formal definitions
2. First crack at a complexity analysis
3. Algorithm successes



# NEUROMORPHIC ALGORITHM TIME & SPACE COSTS



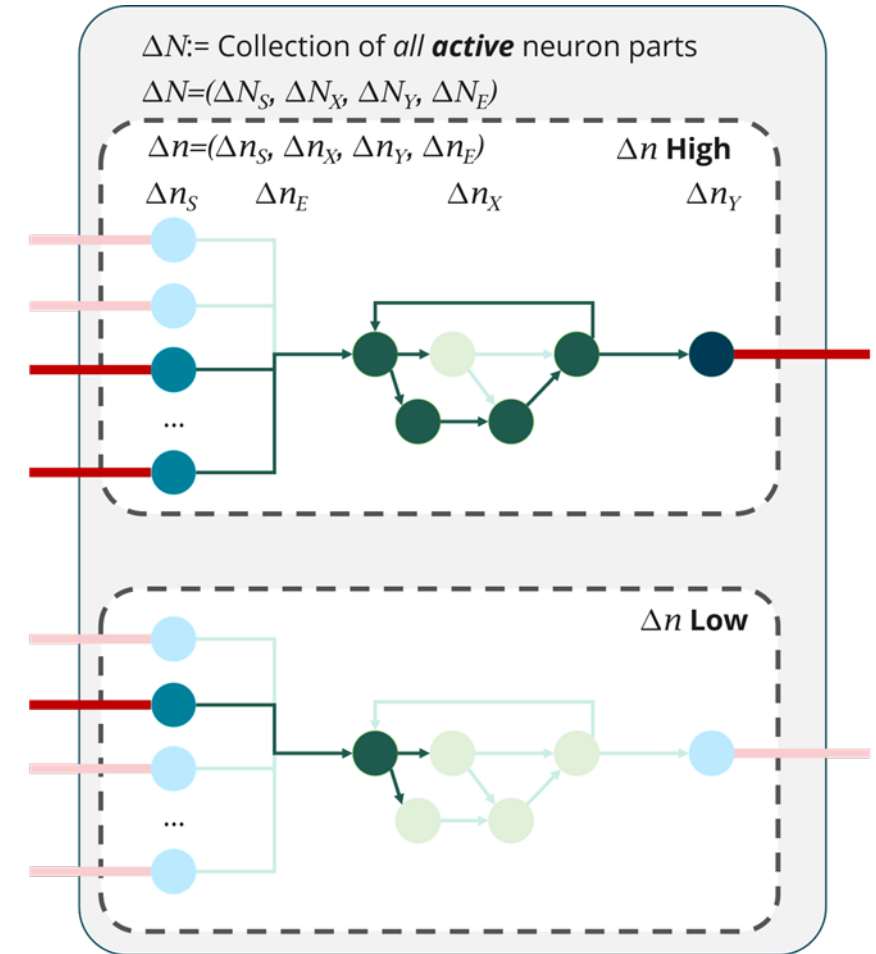
- Ideal NMC executes at parallel depth limits
  - Runtime is governed by the critical path of the time-unrolled graph, i.e.  $C_T \sim T_\infty (G_N^{(T)})$ .
  - *Brent's Theorem applies to neuromorphic algorithms*
- No generic asymptotic time/space advantage
  - NMC aligns with classical parallel depth/work tradeoffs rather than changing complexity class.
- The tradeoff shifts to instantiated footprint
  - Because the “program is the circuit,” footprint scales with deployed graph size.



# NEUROMORPHIC ENERGY COSTS REQUIRE A TRACE-BASED VIEW



- Structure is not enough for energy
  - The instantiated graph  $G_N = (N, S)$  tells us what can happen, not what actually happens.
- Event-driven systems pay *only* for realized activity
  - Rnergy depends on which neuron updates and synaptic events occur, i.e. on  $\Delta G_N(t) = (\Delta N, \Delta S)$ .
  - Energy use is directly related to physical changes in in-memory hardware (flipping a bit, charging a wire)
- This is less convenient a priori, but unavoidable
  - For neuromorphic (and biological) computation, energy cannot in general be inferred from structure alone.



# NEUROMORPHIC ALGORITHM ENERGY COSTS



- Conventional energy tracks scheduled work
  - Operations and memory accesses give energy proportional to executed work, up to hierarchy-dependent constants.
- Neuromorphic energy tracks cumulative activity
  - $C_E(G_N; T) \propto \sum_t |\Delta N(t)| + \sum_t |\Delta S(t)|$
  - $C_E(G_N; T) = O(\sum_{t=0}^{T-1} |\Delta S(t)|)$
- Advantage is conditional, not universal
  - If activity is dense, no gain; if traces are sparse/decaying, energy can scale with realized activity rather than worst-case structure.

# OUTLINE

1. Formal definitions
2. First crack at a complexity analysis
3. Algorithm successes





## Structural Metrics

- Graph homogeneity
  - $H(t) := \max_{op} |v \in V_t: \text{type}(v) = op|$
- Structural reuse factor
  - $R_T(G_N) = \frac{T_1(G_N^{(T)})}{|G_N|}$
- Fan-out
  - $K(G_N) := \max_{n_y \in N_Y} \text{deg}_S^+(n_y)$
  - $\bar{K}(G_N) := \frac{|S|}{|N_Y|}$

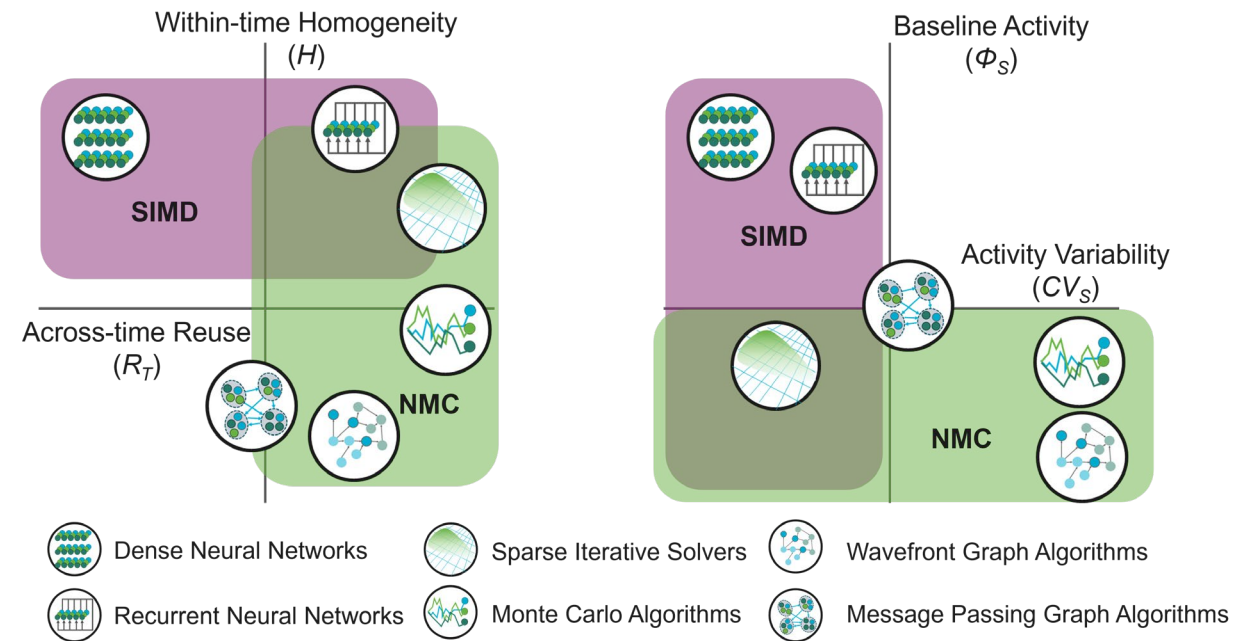
## Trace Metrics

- Synaptic activity / sparsity
  - $\Phi_S(G_N; T) = \frac{1}{|S|T} \sum_{t=0}^{T-1} |\Delta S(t)|$
- Activity decay / growth
  - $\Lambda_S(G_N; T) = \frac{1}{T-1} \sum_{t=0}^{T-2} (a_S(t+1) - a_S(t))$
- Activity variability
  - $CV_S(G_N; T) = \frac{\sqrt{\text{Var}(A_S)}}{E[A_S]}$

$$A_S := \sum_{t=0}^{T-1} |\Delta S(t)| \quad a_S(t) := |\Delta S(t)|$$

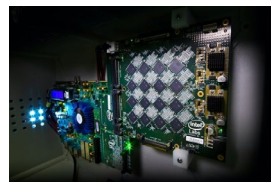
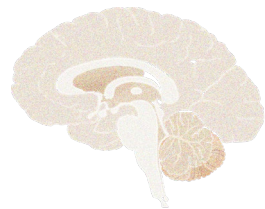
# UPSHOT OF THIS FRAMEWORK

- Neuromorphic advantage *in energy*
  - Formally simply another extreme parallel architecture, but perhaps can realize near-ideal parallel scaling
  - Energy advantage is conditional on sparse activity ( $\Delta G$  low)
  - Empirical demonstration on DTMC random walks and finite element simulations.
- Advantage relative to GPUs
  - High use of temporal recurrence (allows amortization of spatial cost)
  - High heterogeneity, thread divergence, stochasticity, etc.

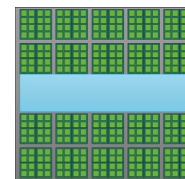




## The Brain / Neuromorphic



## GPUs / SIMD Architectures



### Structural Metrics

Structural Heterogeneity

Recurrence / Structural Reuse

Local Connectivity

Structural Homogeneity

Feed-forward / Acyclic

Global Connectivity

### Trace Metrics

Sparse Activity

Time-dependent Activity

- Thread Divergence
- Stochasticity okay

Dense Activity

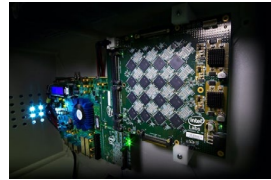
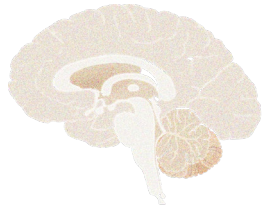
Time-homogeneous Activity

- Thread Convergence
- Stochasticity challenging

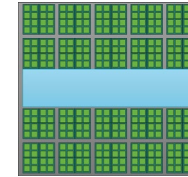
# ARTIFICIAL FEED FORWARD NEURAL NETWORKS



## The Brain / Neuromorphic



## GPUs / SIMD Architectures



### Structural Metrics

Structural Heterogeneity



Structural Homogeneity

Recurrence / Structural Reuse



Feed-forward / Acyclic

Local Connectivity



Global Connectivity

### Trace Metrics

Sparse Activity



Dense Activity

Time-dependent Activity



Time-homogeneous Activity

Thread Divergence



Thread Convergence

Stochasticity okay

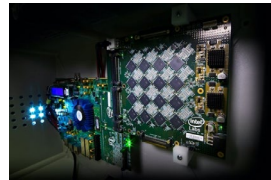
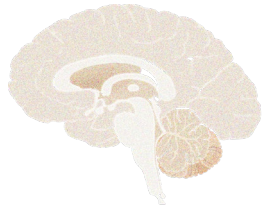


Stochasticity challenging

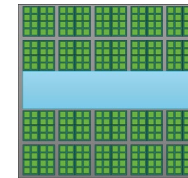
# GRAPH ANALYTICS (MESSAGE PASSING, ETC)



## The Brain / Neuromorphic



## GPUs / SIMD Architectures



### Structural Metrics

Structural Heterogeneity	
Recurrence / Structural Reuse	
Local Connectivity	

Structural Homogeneity
Feed-forward / Acyclic
Global Connectivity

### Trace Metrics

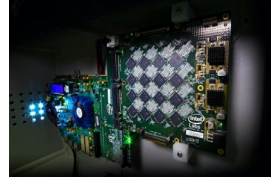
Sparse Activity	
Time-dependent Activity	
Thread Divergence	
Stochasticity okay	

Dense Activity
Time-homogeneous Activity
Thread Convergence
Stochasticity challenging

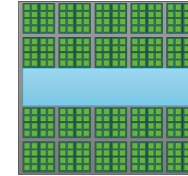
# SPARSE LINEAR SOLVERS (CONJUGATE GRADIENT, ETC)



## The Brain / Neuromorphic

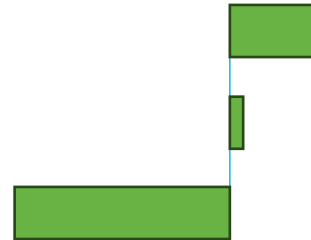


## GPUs / SIMD Architectures



### Structural Metrics

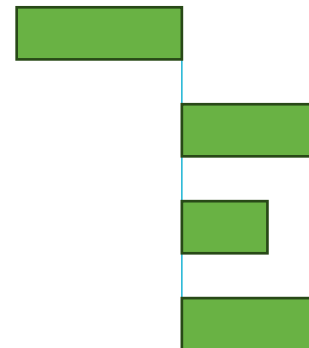
Structural Heterogeneity  
Recurrence / Structural Reuse  
Local Connectivity



Structural Homogeneity  
Feed-forward / Acyclic  
Global Connectivity

### Trace Metrics

Sparse Activity  
Time-dependent Activity  
Thread Divergence  
Stochasticity okay



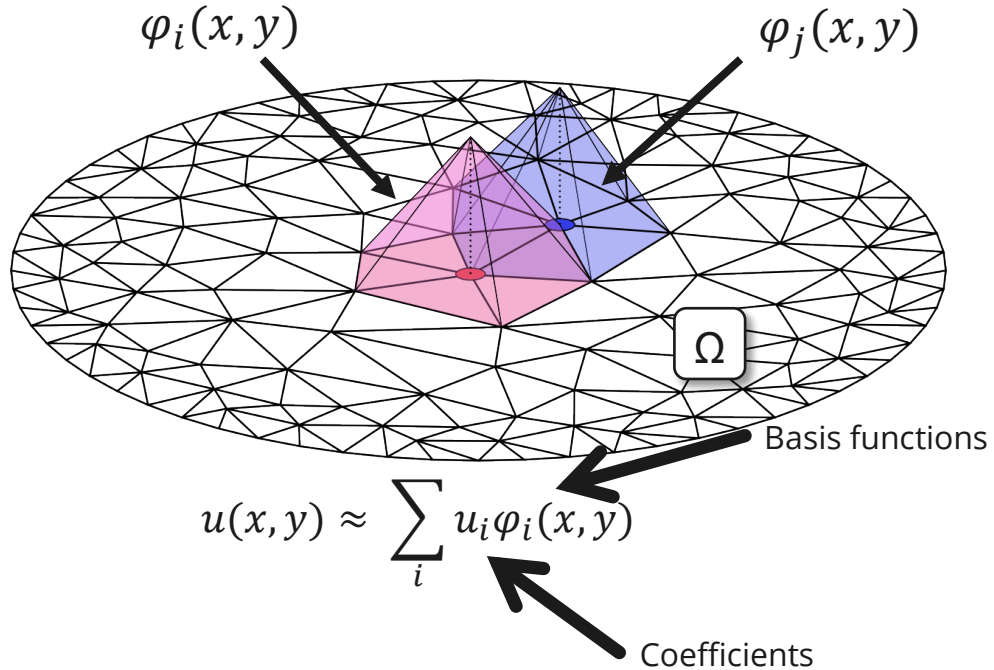
Dense Activity  
Time-homogeneous Activity  
Thread Convergence  
Stochasticity challenging

# CAN WE TACKLE FEM WITH PROBABILISTIC NEURAL HARDWARE?



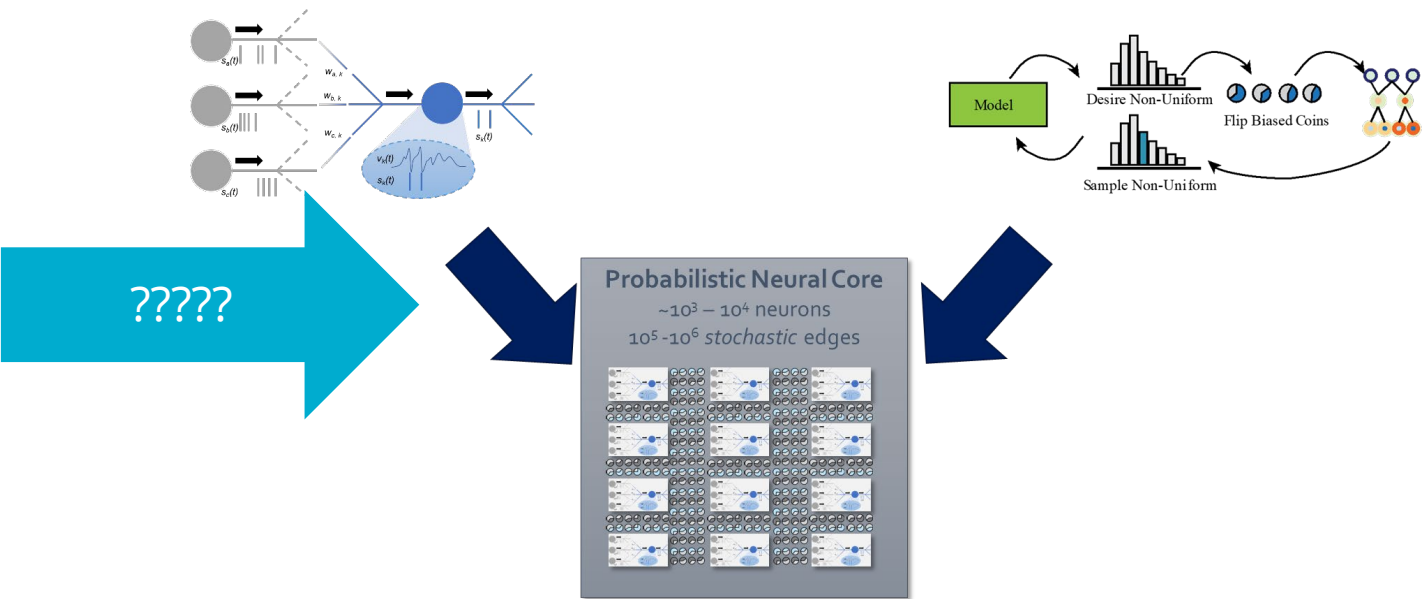
$$\nabla^2 u = f$$

Finite Element Methods - Gold Standard but expensive



Sparse, linear system

$$A\vec{u} = \vec{b}$$



**NEUROMORPHIC FINITE ELEMENT METHODS?**



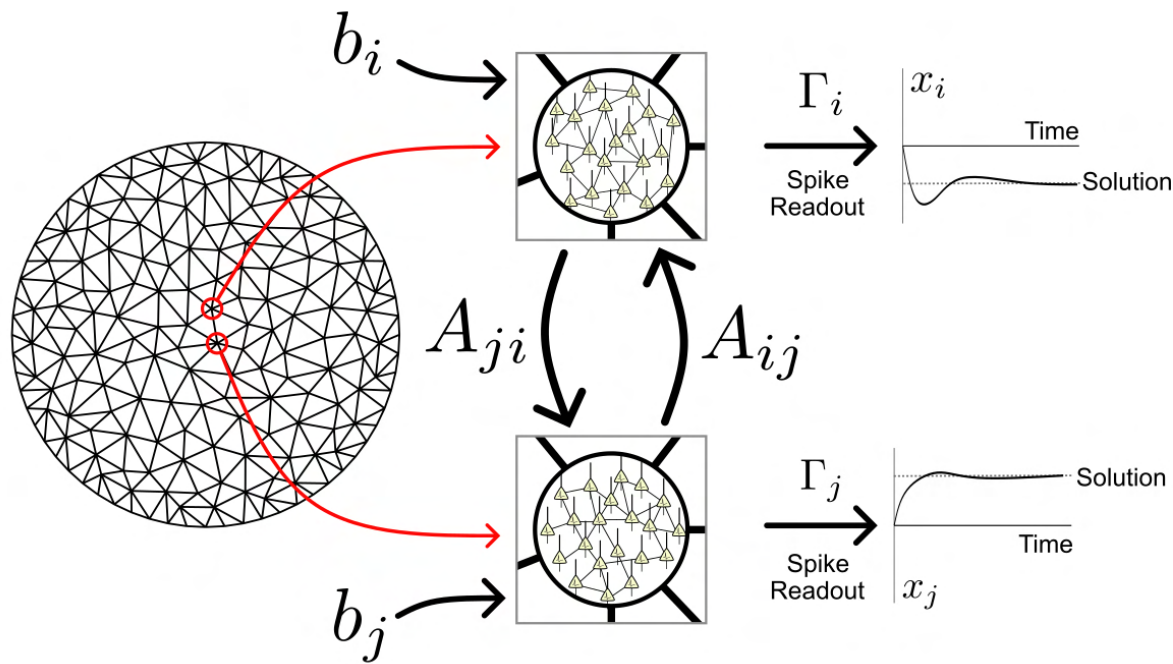
Poisson Equation

$$\nabla^2 u = f$$



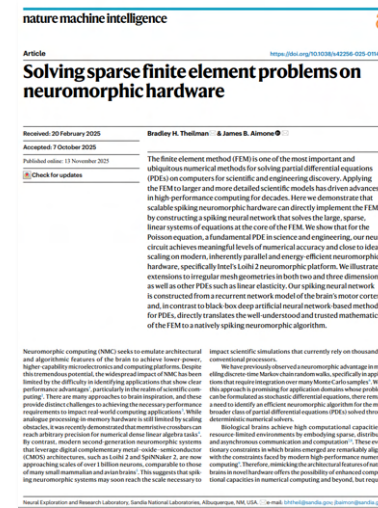
Sparse Linear System

$$A\vec{x} = \vec{b}$$



## Neuromorphic virtues:

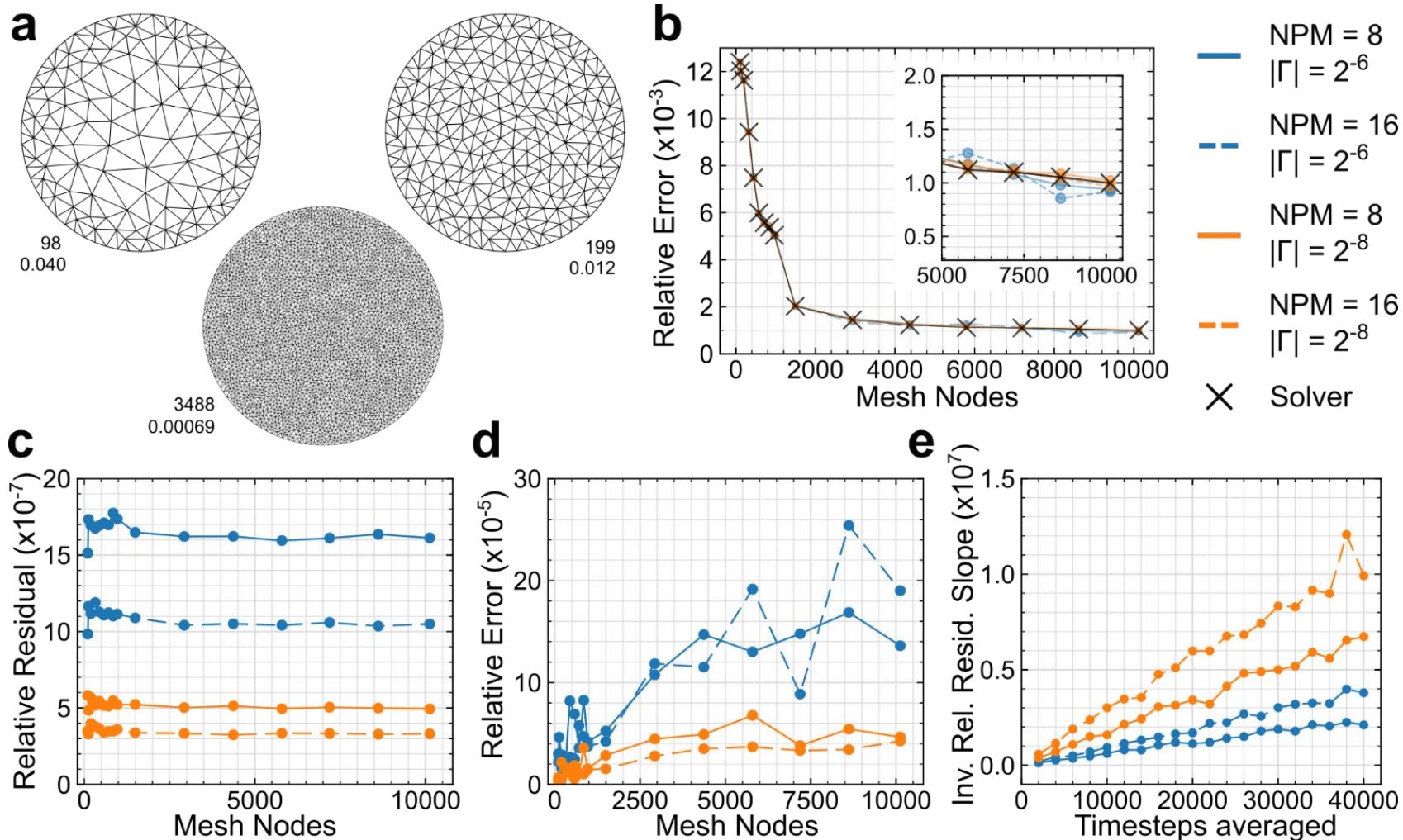
- Locally dense, globally sparse connectivity
- Scalable: neighbors  $\sim O(1)$
- Sparse spiking activity



# NEUROFEM IN ACTION



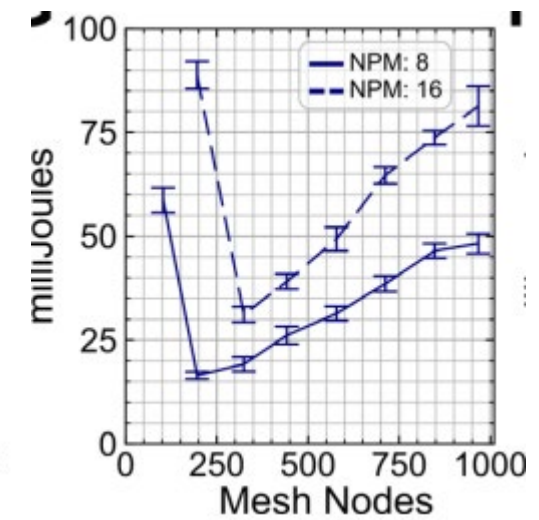
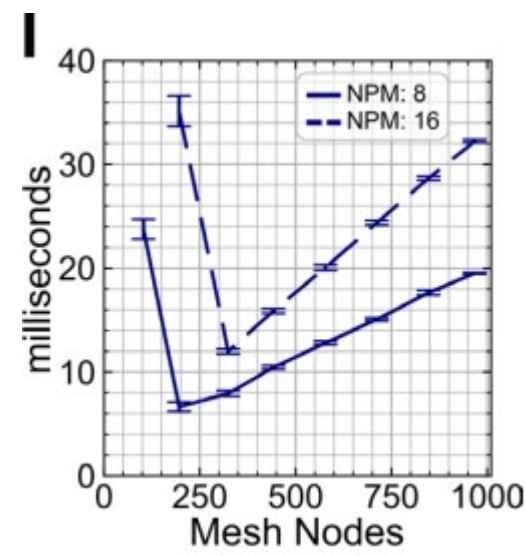
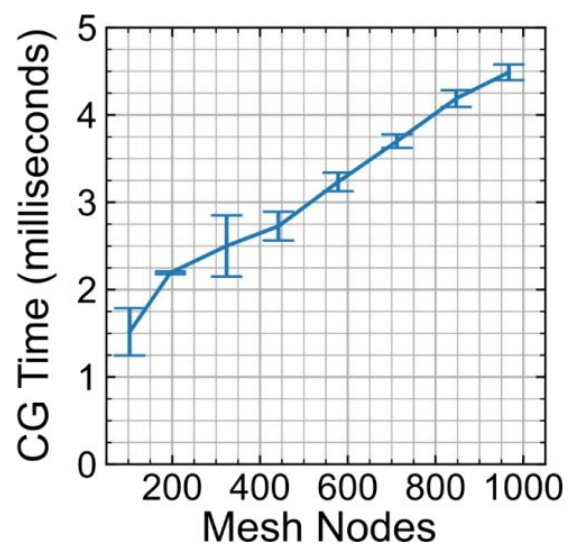
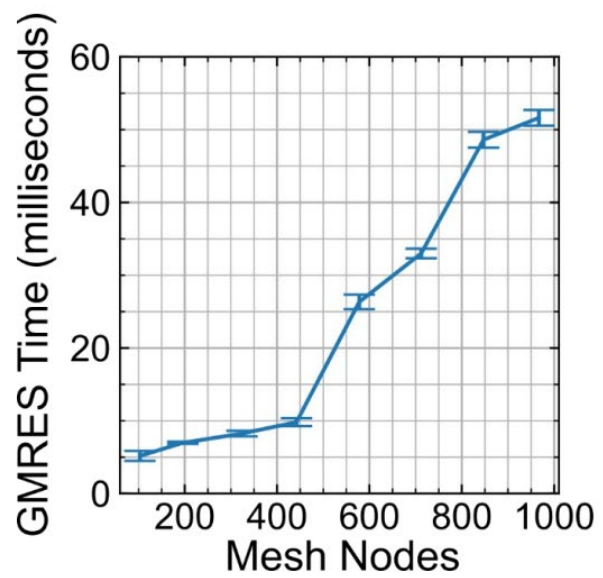
# NEUROFEM PROVIDES SOLVER-QUALITY SOLUTIONS



# NEUROFEM IS SIMILAR IN SPEED TO GMRES AND CONJUGATE GRADIENT

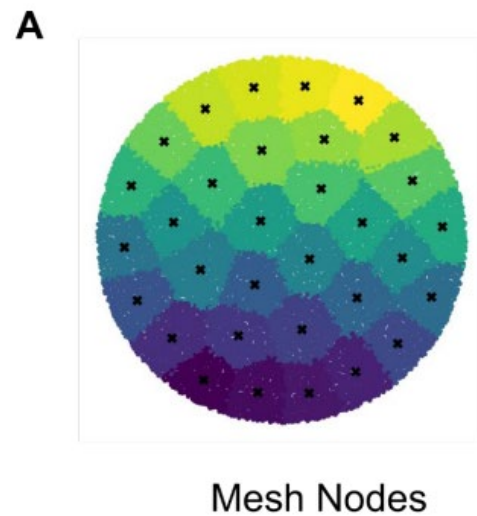


Neuromorphic solution is general

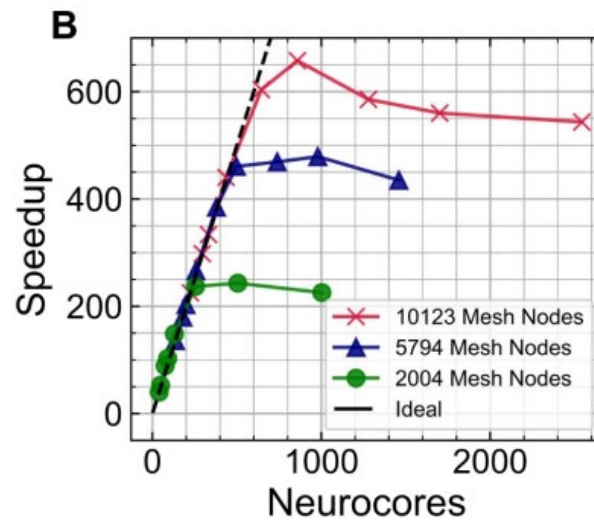


Conjugate gradient is faster at small scales, but can only be used for symmetric matrices

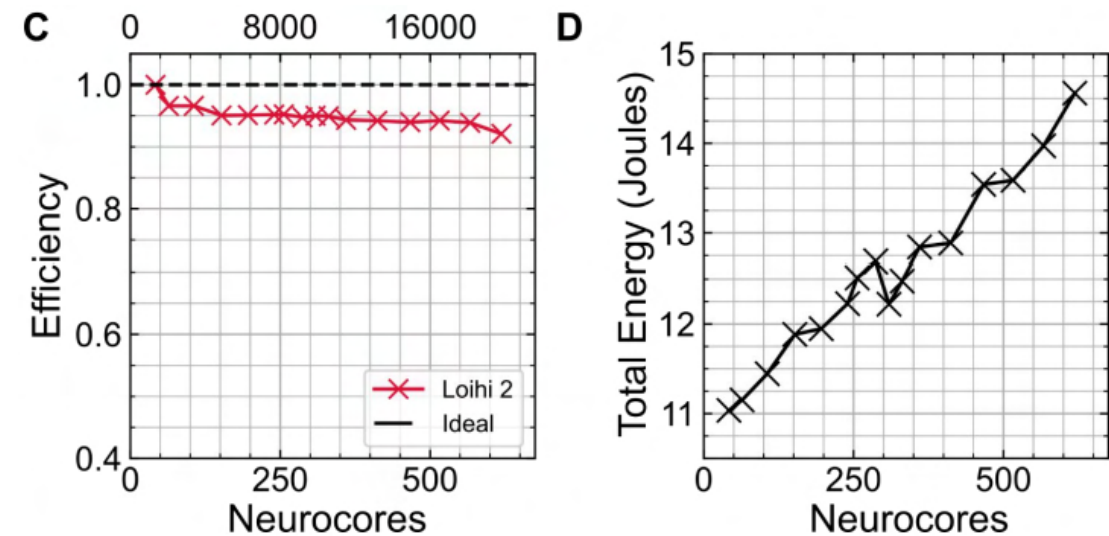
# NEUROFEM SHOWS COMPELLING STRONG AND WEAK SCALING



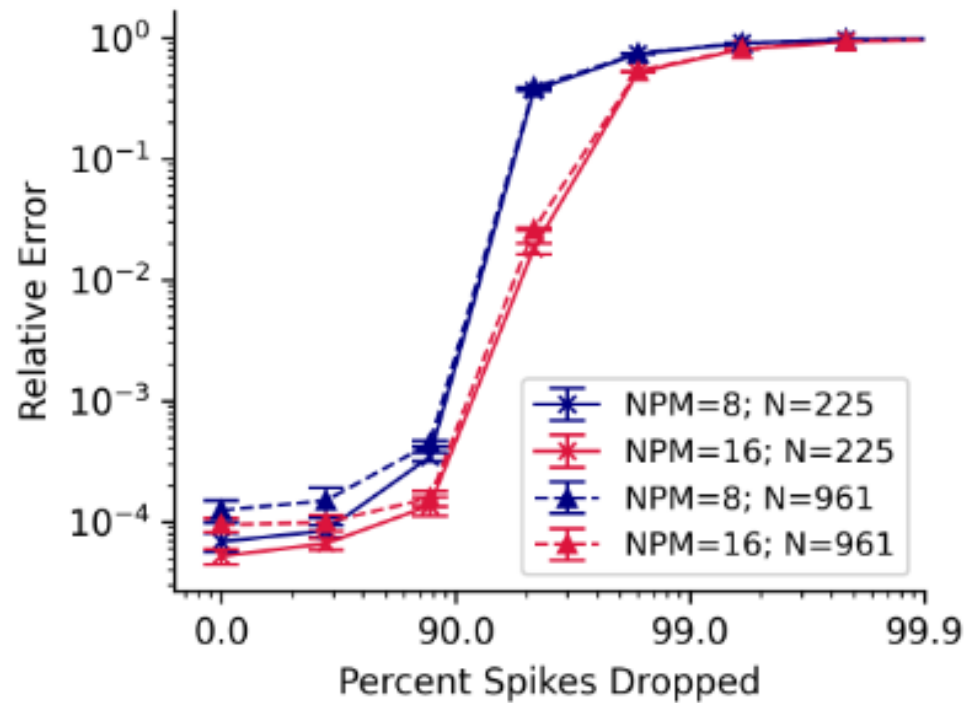
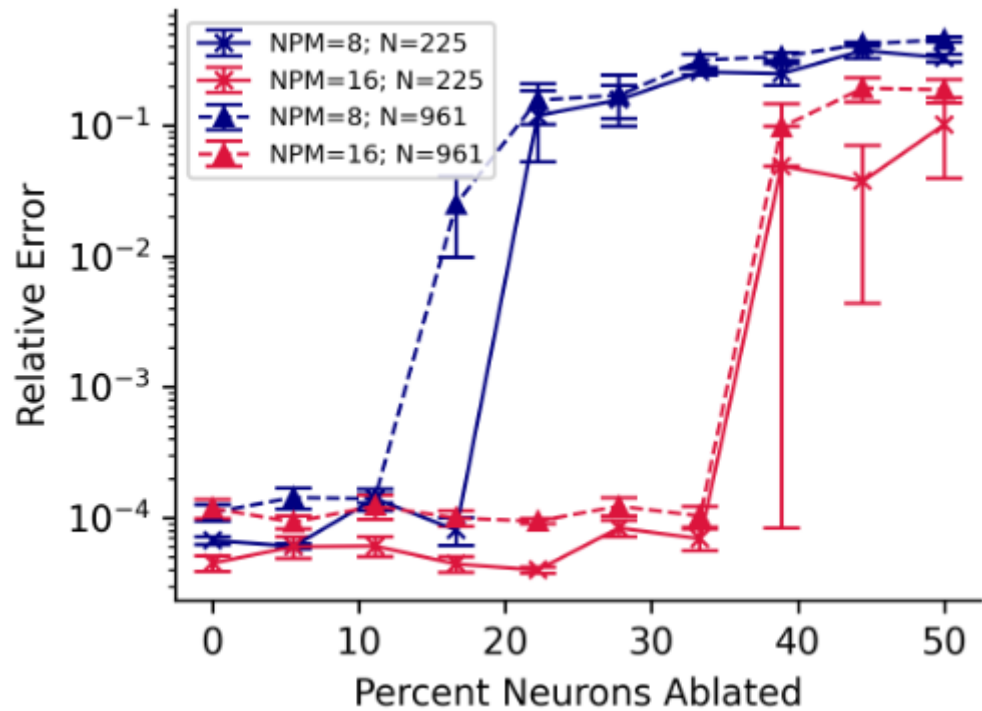
## Strong Scaling



## Weak Scaling



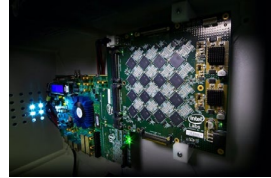
# NEUROFEM IS INTRINSICALLY FAULT AND ERROR TOLERANT



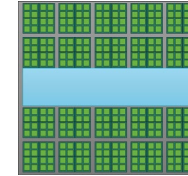
# SPARSE LINEAR SOLVERS (NEUROFEM)



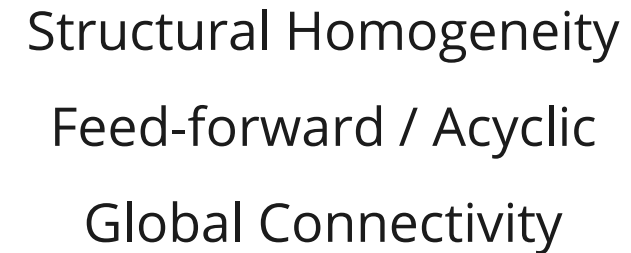
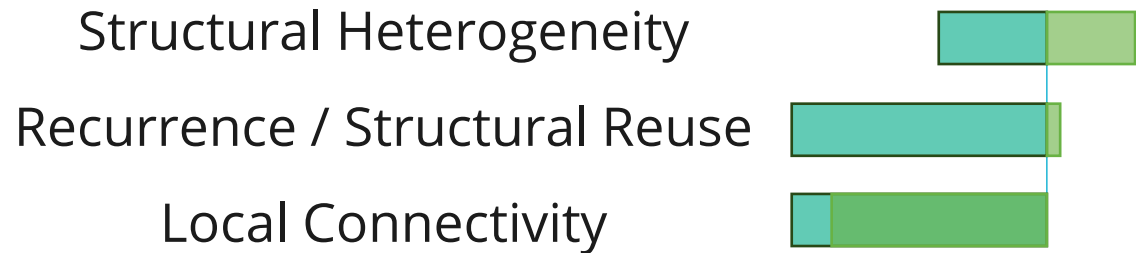
## The Brain / Neuromorphic



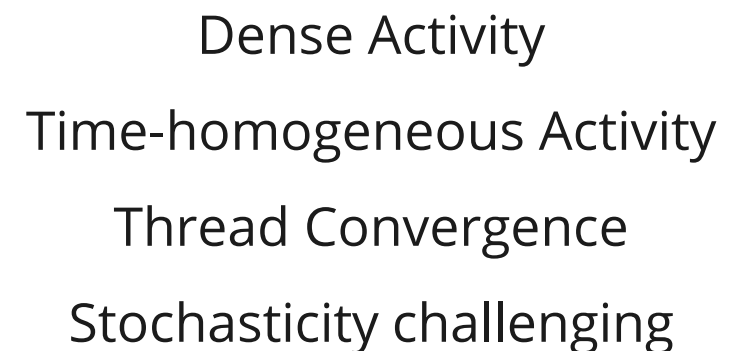
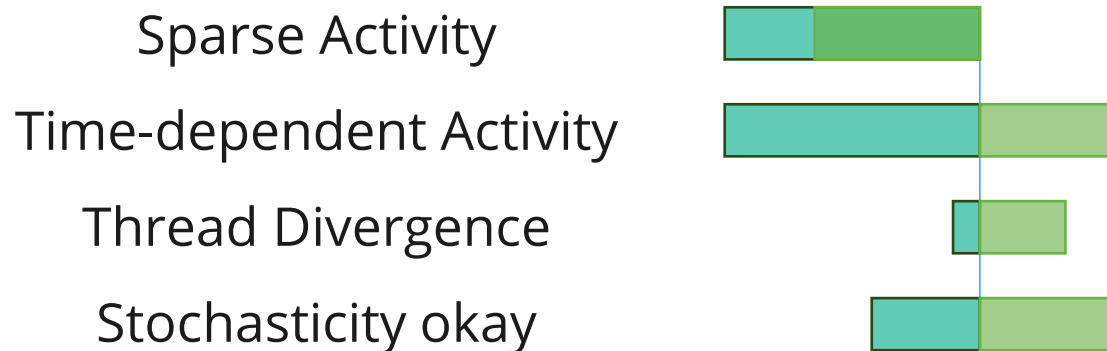
## GPUs / SIMD Architectures



### Structural Metrics



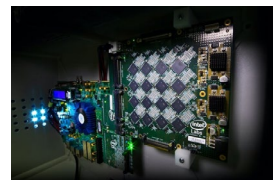
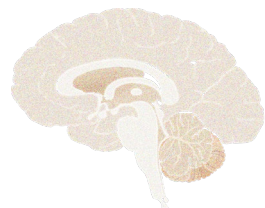
### Trace Metrics



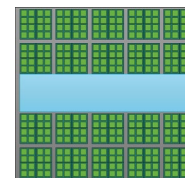
# MONTE CARLO



## The Brain / Neuromorphic



## GPUs / SIMD Architectures



### Structural Metrics

Structural Heterogeneity		█	Structural Homogeneity
Recurrence / Structural Reuse	█		Feed-forward / Acyclic
Local Connectivity		█	Global Connectivity

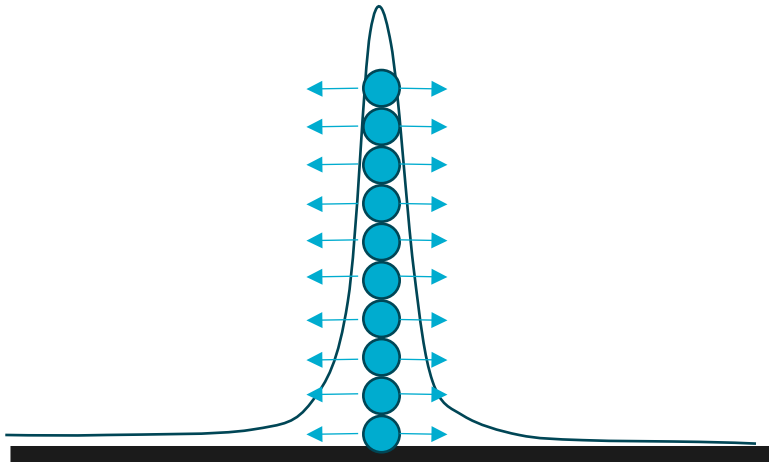
### Trace Metrics

Sparse Activity	█	Dense Activity
Time-dependent Activity	█	Time-homogeneous Activity
Thread Divergence	█	Thread Convergence
Stochasticity okay	█	Stochasticity challenging

# CAN WE REFORMULATE MONTE CARLO FOR NEUROMORPHIC?



Initial State



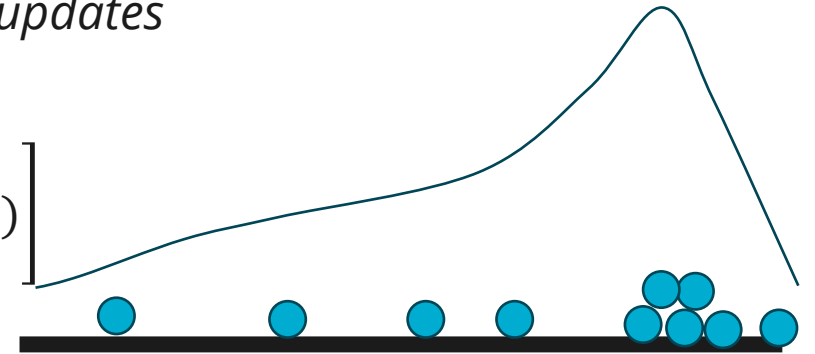
$$\frac{dx_i}{dt} = f(x_i, X, \dots)$$

*K particles*

*K location-dependent updates*

$$\begin{bmatrix} \uparrow \\ x_i(t+1) \\ \downarrow \end{bmatrix} = f \left( \begin{bmatrix} \uparrow \\ x_i(t) \\ \downarrow \end{bmatrix}, X, \dots \right) + \begin{bmatrix} \uparrow \\ x_i(t) \\ \downarrow \end{bmatrix}$$

Final State

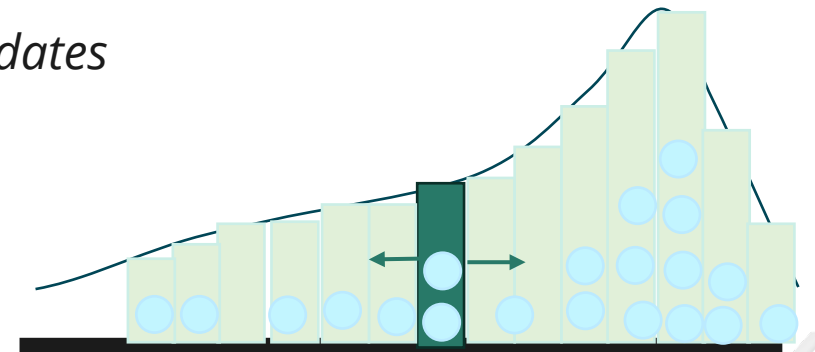
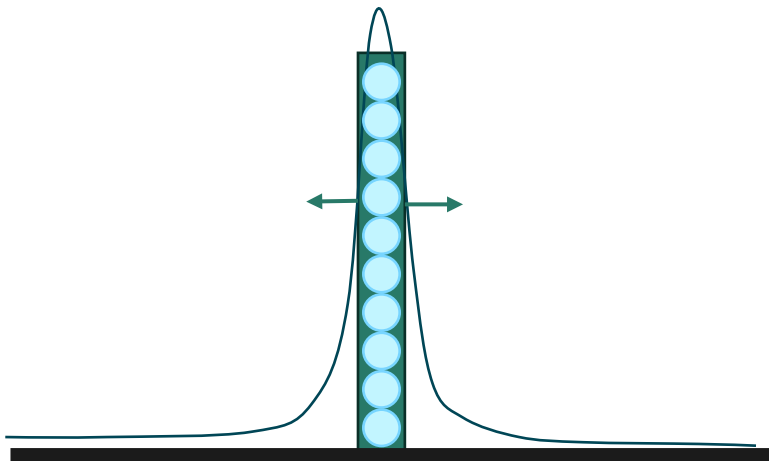


$$\frac{dm_i}{dt} = g_i(X, \dots)$$

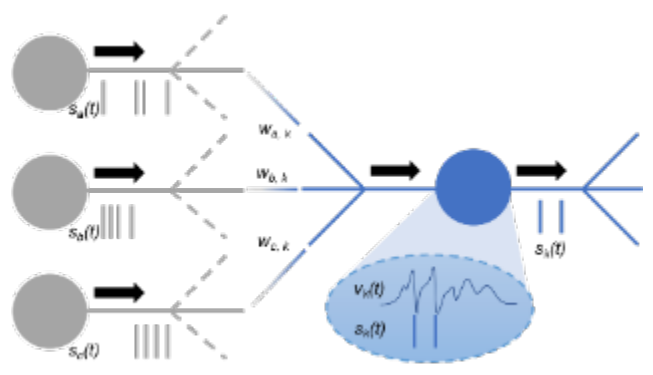
*M locations*

*M location-specific updates*

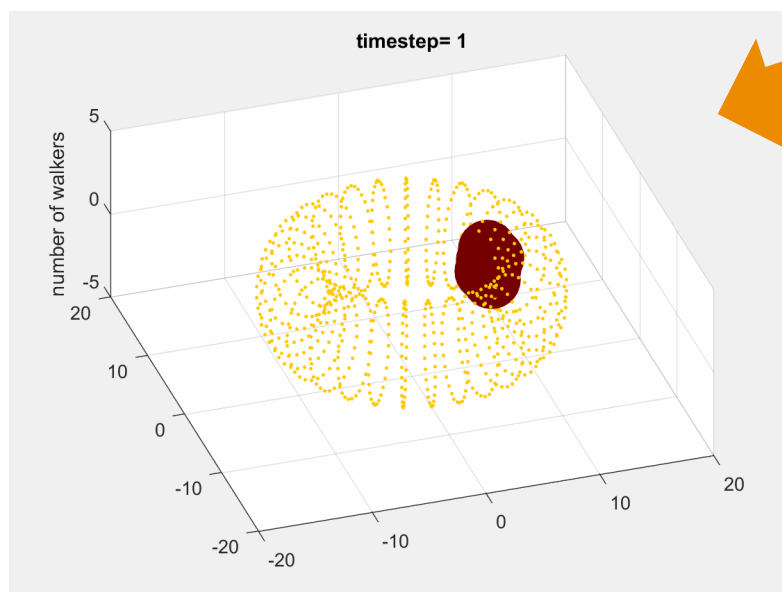
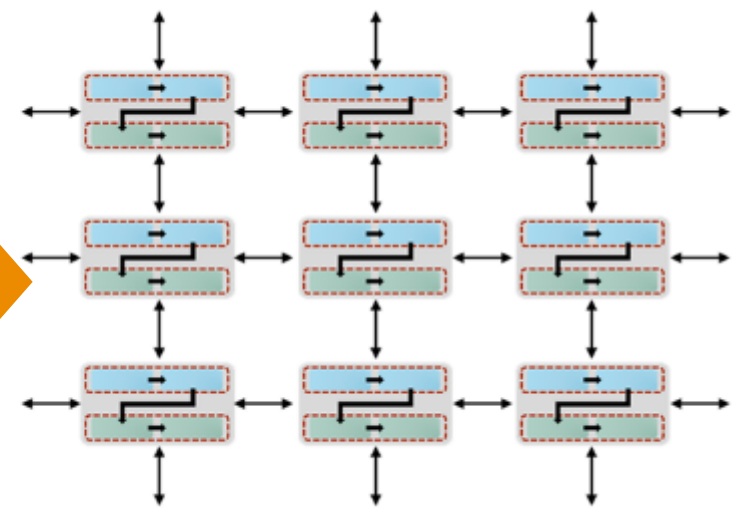
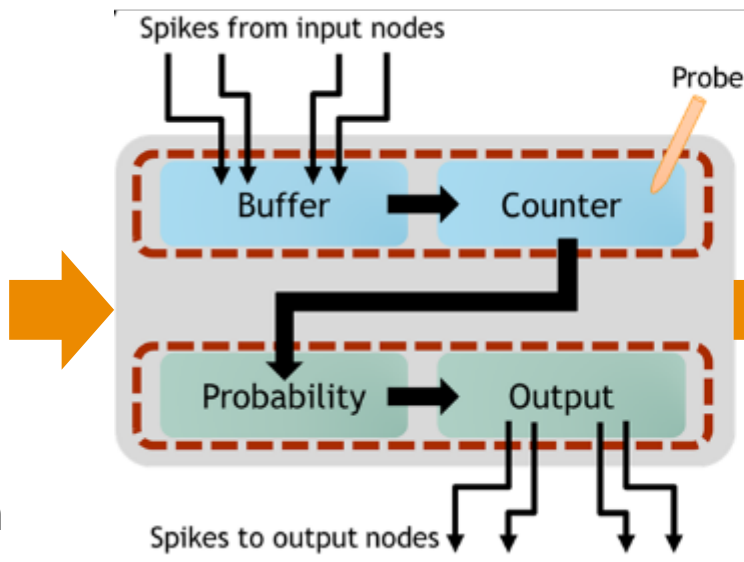
$$\begin{bmatrix} \uparrow \\ m_i(t+1) \\ \downarrow \end{bmatrix} = g_i(X, \dots) + \begin{bmatrix} \uparrow \\ m_i(t) \\ \downarrow \end{bmatrix}$$



# USE NEURONS TO REPRESENT STATE SPACE OF MONTE CARLO AND USE SPIKES TO REPRESENT PARTICLES



Leaky integrate-and-fire neuron



nature electronics ARTICLES  
<https://doi.org/10.1038/s41928-021-00705-7>

### Neuromorphic scaling advantages for energy-efficient random walk computations

J. Darby Smith, Aaron J. Hill, Leah E. Reeder, Brian C. Frank, Richard B. Lehoucq, Ojas Parekh, William Severa and James B. Aumane

Neuromorphic computing, which aims to replicate the computational structure and architecture of the brain in synthetic hardware, has typically focused on artificial intelligence applications. What is less explored is whether such brain-inspired hardware can provide value beyond cognitive tasks. Here we show that the high degree of parallelism and configurability of spiking neuromorphic architectures makes them well suited to implement random walks via discrete-time Markov chains. These random walks are useful in Monte Carlo methods, which represent a fundamental computational tool for solving a wide range of numerical computing tasks. Using IBM's TrueNorth and Intel's Loihi neuromorphic computing platforms, we show that our neuromorphic computing algorithm for generating random walk approximations of diffusion offers advantages in energy-efficient computation compared with conventional approaches. We also show that our neuromorphic computing algorithm can be extended to more sophisticated jump-diffusion processes that are useful in a range of applications, including financial economics, particle physics and machine learning.

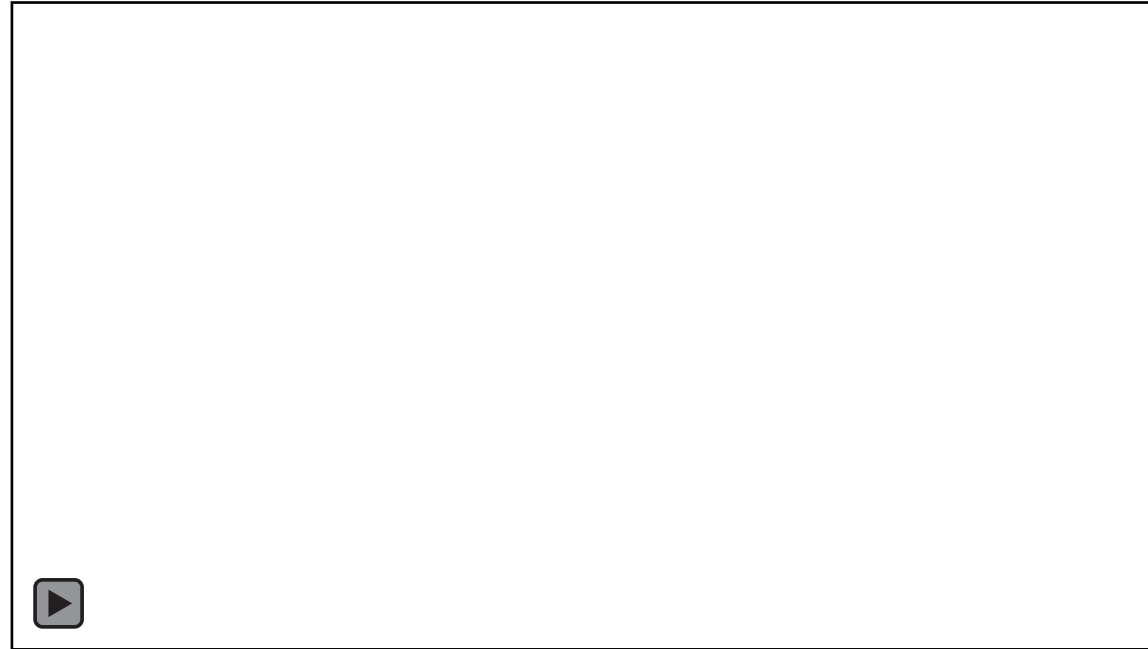
Despite the increasing ability to develop large-scale neural hardware, the theoretical value of neuromorphic hardware remains unclear—unlike quantum computing that offers clear fundamental advantages at scale. Nevertheless, there are several architectural features of most nervous systems that could yield advantages including the high degree of connectivity between neurons, the collocation of processing and memory, and the use of action potentials (inferred as spikes) to communicate. Algorithm research for spiking neuromorphic hardware has primarily focused on its suitability for deep learning and other emerging artificial intelligence (AI) algorithms. Such applications are straightforward, given the alignment of neural architectures with neural networks, and it can be expected that the value of neuromorphic computing will grow as AI algorithms derive further inspiration from the brain. However, the impact of neuromorphic computing beyond cognitive applications is less certain.

Quantum computing has shown how emerging hardware can have an impact beyond its original inspiration: it was conceived as a means for efficient chemistry simulations, but is now recognized as useful in a much broader range of applications. Unlike quantum computing, which faces technical challenges in scaling up, time scaling compared with the von Neumann architecture and still requiring less total energy to perform the same computation.

Observing a neuromorphic advantage for non-cognitive applications should not be taken as a given since the specialization of computer architectures to improve performance on a subset of tasks will likely result in degraded performance in other tasks. Therefore, observing a neuromorphic advantage on non-cognitive applications would demonstrate that neuromorphic computing can have a broader impact than previously assumed and provide a concrete framework by which to develop the technology. Although a definitive neuromorphic advantage (as defined here) has not yet been demonstrated for non-cognitive applications, there are three categories of such computing tasks that appear well suited for neuromorphic computing: linear algebra, in which the high fan-in of neurons can be used to realize known theoretical advantages of threshold gate (TG) logic; graph analytical tasks that can leverage the configurability and parallelization of neural circuits; and sampling steady-state distributions for a wide range of potential applications using stochastic neural circuits.

In this Article, we show that large-scale neuromorphic hardware can offer a neuromorphic advantage on a fundamental

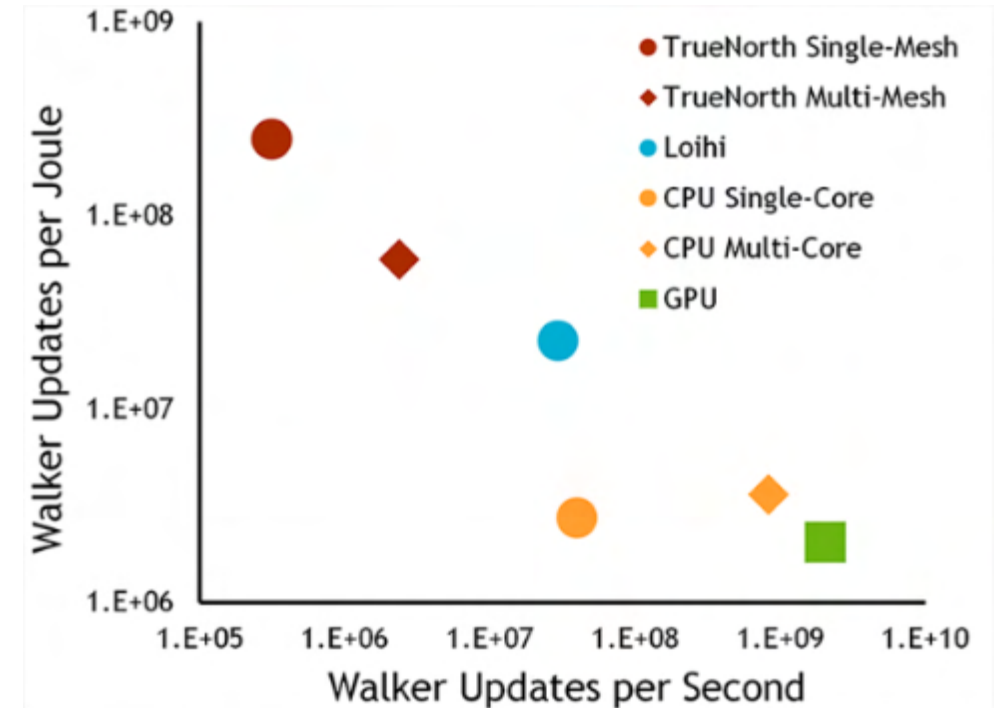
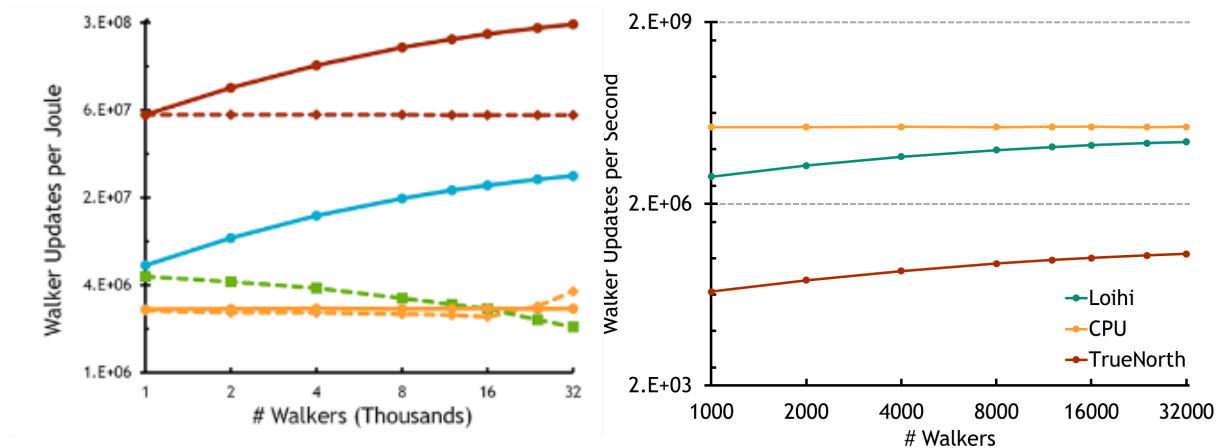
NEUROMORPHIC COMPUTING ADVANTAGE APPEARS TO BE WHEN AN ALGORITHM CAN SPLIT THE TASK ACROSS COMPUTATIONAL GRAPH WITH SPARSE COMMUNICATION



# WE CAN IDENTIFY A NEUROMORPHIC ADVANTAGE FOR SIMULATING RANDOM WALKS



We define a *neuromorphic advantage* as an algorithm that shows a demonstrable **advantage** in terms of one resource (e.g., energy) while exhibiting comparable **scaling** in other resources (e.g., time).



# WHAT PDES CAN NEURAL RANDOM WALKS ADDRESS?



## ***Class of Partial Integro-Differential Equations:***

$$\begin{aligned} \frac{\partial}{\partial t} u(t, \mathbf{x}) = & \frac{1}{2} \sum_{i,j} (\mathbf{a}\mathbf{a}^\top)_{i,j}(t, \mathbf{x}) \frac{\partial^2}{\partial x_i \partial x_j} u(t, \mathbf{x}) + \sum_i b_i(t, \mathbf{x}) \frac{\partial}{\partial x_i} u(t, \mathbf{x}) \\ & + \lambda(t, \mathbf{x}) \int \left( u(t, \mathbf{x} + \mathbf{h}(t, \mathbf{x}, q)) - u(t, \mathbf{x}) \right) \phi_Q(q; t, \mathbf{x}) dq \\ & + c(t, \mathbf{x})u(t, \mathbf{x}) + f(t, \mathbf{x}), \quad \mathbf{x} \in \mathbb{R}^d, t \in [0, \infty). \end{aligned}$$

## ***Stochastic Process:***

$$d\mathbf{X}(t) = \mathbf{b}(t, \mathbf{X}(t))dt + \mathbf{a}(t, \mathbf{X}(t))d\mathbf{W}(t) + \mathbf{h}(t, \mathbf{X}(t), q)dP(t; Q, \mathbf{X}(t)).$$

NMC Hardware Simulates This Stochastic Process



Monte Carlo Approximates This Expectation

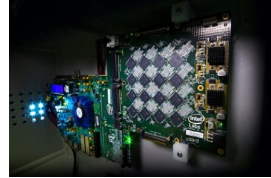
## ***Solution to initial value problem ( $u(0, \mathbf{x}) = g(\mathbf{x})$ ):***

$$u(t, \mathbf{x}) = \mathbb{E} \left[ g(\mathbf{X}(t)) \exp \left( \int_0^t c(s, \mathbf{X}(s)) ds \right) + \int_0^t f(s, \mathbf{X}(s)) \exp \left( \int_0^s c(\ell, \mathbf{X}(\ell)) d\ell \right) ds \middle| \mathbf{X}(0) = \mathbf{x} \right].$$

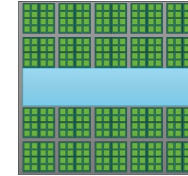
# MONTE CARLO



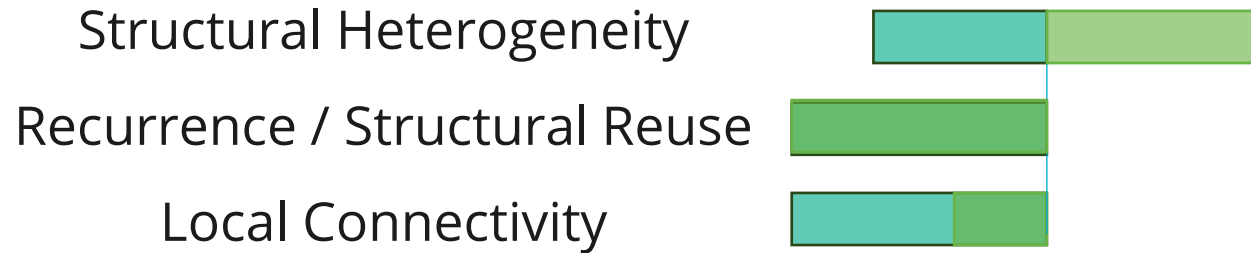
## The Brain / Neuromorphic



## GPUs / SIMD Architectures

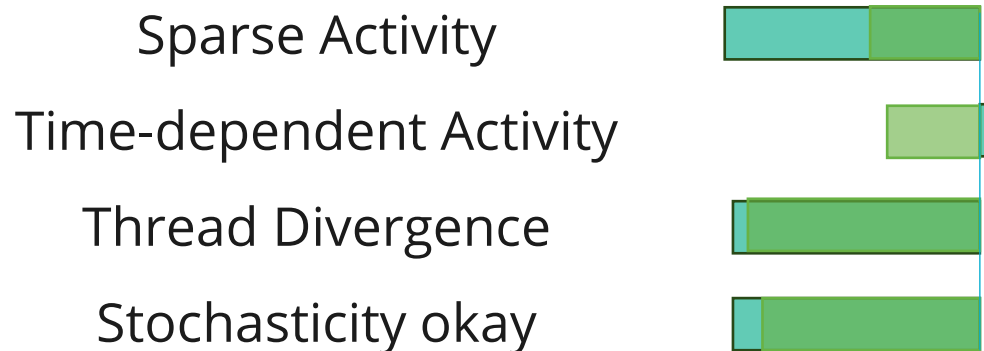


### Structural Metrics



Structural Homogeneity  
Feed-forward / Acyclic  
Global Connectivity

### Trace Metrics

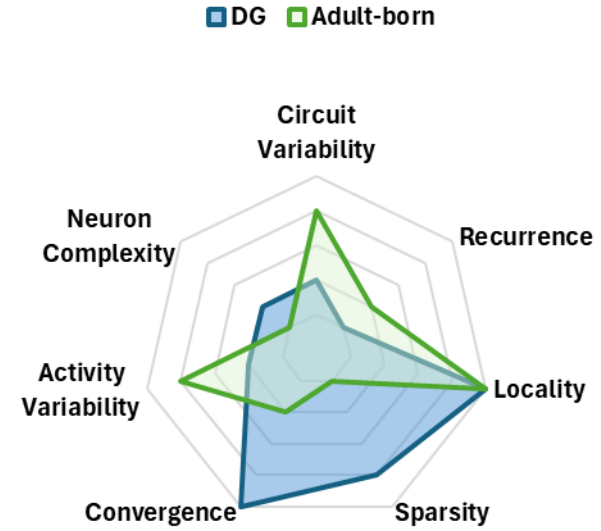
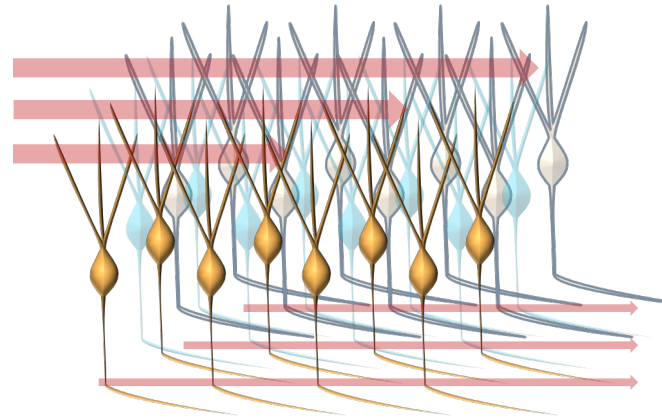


Dense Activity  
Time-homogeneous Activity  
Thread Convergence  
Stochasticity challenging

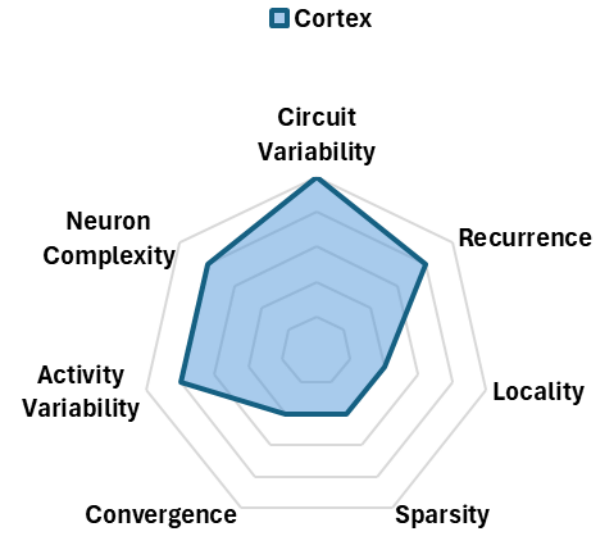
# WHAT ABOUT THE BRAIN?



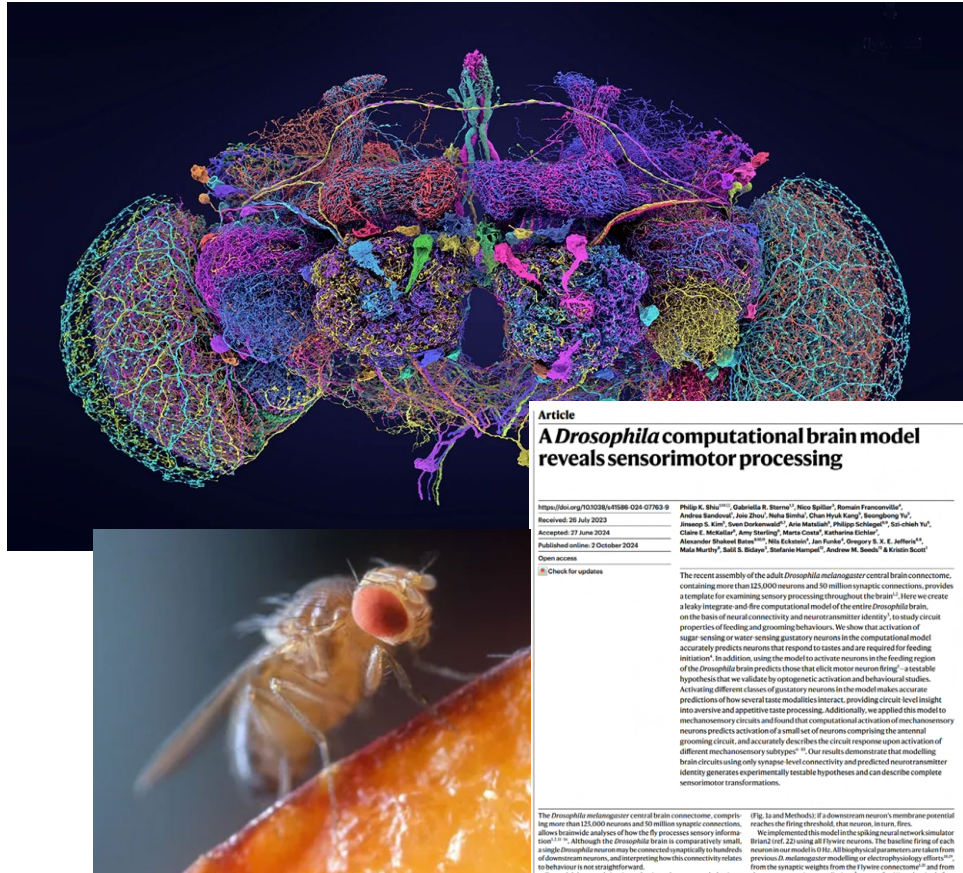
## Dentate Gyrus



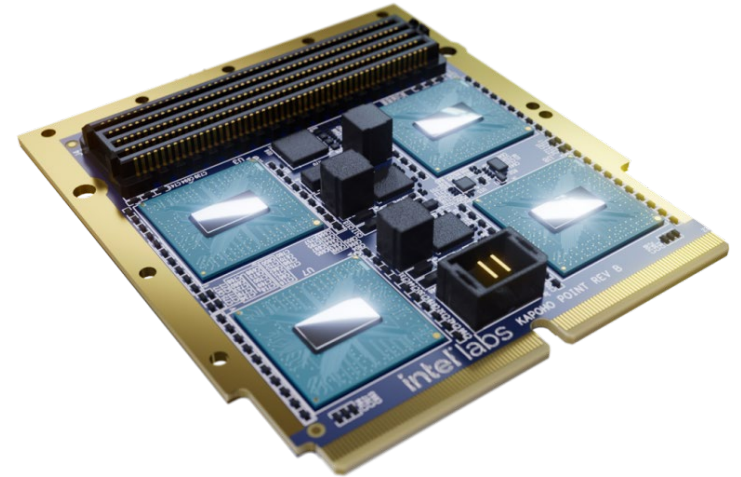
## Cortex



# REAL BRAINS HAVE NON-TRIVIAL STRUCTURE



125,000 neurons  
50 million synapses



## Article **A *Drosophila* computational brain model reveals sensorimotor processing**

<https://doi.org/10.1038/s41586-024-07762-9>  
Received: 20 July 2024  
Accepted: 27 June 2024  
Published online: 2 October 2024  
Open access

**Philip K. Shiu<sup>1,2,3\*</sup>, Gabriella E. Burns<sup>4</sup>, Nico Spillier<sup>5</sup>, Roman Francomilla<sup>6</sup>, Andrea Bardeoni<sup>7</sup>, Joo Zhou<sup>8</sup>, Naha Brink<sup>9</sup>, Chan Hyuk Kang<sup>10</sup>, Beonghwi Yu<sup>11</sup>, Jinsong S. Kim<sup>12</sup>, Soren Dockenwald<sup>13</sup>, Arne Marziani<sup>14</sup>, Philipp Schöberl<sup>15</sup>, So-Heek Yu<sup>16</sup>, Clara E. Motulsky<sup>17</sup>, Amy Swearing<sup>18</sup>, Maria Costa<sup>19</sup>, Katharina Eichler<sup>20</sup>, Alexander Bialas<sup>21</sup>, Jia-E. Kim<sup>22</sup>, Jan Farkas<sup>23</sup>, Gergely S. E. Jellison<sup>24</sup>, Mala Murthy<sup>25</sup>, Sallie S. Bidaye<sup>26</sup>, Stefania Harapel<sup>27</sup>, Andrew M. Sved<sup>28</sup> & Kristin Scott<sup>29</sup>**

The recent assembly of the adult *Drosophila melanogaster* central brain connectome, containing more than 125,000 neurons and 50 million synaptic connections, provides a template for examining sensory processing throughout the brain<sup>1</sup>. Here we create a fully integrate-and-fire computational model of the entire *Drosophila* brain, on the basis of neural connectivity and neurotransmitter identity<sup>2</sup>, to study circuit properties of feeding and grooming behaviours. We show that activation of sugar-sensing or water-sensing gustatory neurons in the computational model accurately predicts neurons that respond to tastes and are required for feeding initiation<sup>3</sup>. In addition, using the model to activate neurons in the feeding region of the *Drosophila* brain predicts those that elicit motor neuron firing<sup>4</sup>—a testable hypothesis that we validate by optogenetic activation and behavioural studies. Activating different classes of gustatory neurons in the model makes accurate predictions of how several taste modalities interact, providing circuit-level insight into overdrive and appetitive taste processing. Additionally, we applied this model to mechanosensory circuits and found that computational activation of mechanosensory neurons predicts activation of a small set of neurons comprising the antennal grooming circuit, and accurately describes the circuit response upon activation of different mechanosensory subtypes<sup>5</sup>. Our results demonstrate that modelling brain circuits using only synapse-level connectivity and predicted neurotransmitter identity generates experimentally testable hypotheses and can describe complete sensorimotor transformations.

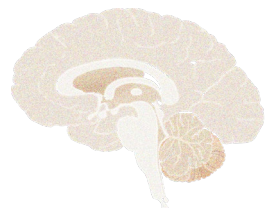
The *Drosophila melanogaster* central brain connectome, comprising more than 125,000 neurons and 50 million synaptic connections, allows brainwide analyses of how the fly processes sensory information<sup>1,2</sup>. Although the *Drosophila* brain is comparatively small, a single *Drosophila* neuron may be connected synaptically to hundreds of downstream neurons, and interpreting how this connectivity relates to behaviour is not straightforward. To model the neural circuit mechanisms that generate behaviour, we implement a simple leaky integrate-and-fire model using the connectome of reconstructed electron microscopy neurons<sup>1,2</sup> as well as neurotransmitter predictions for each neuron<sup>3</sup>. In this model, spiking of a neuron alters the membrane potential of downstream neurons in proportion to the connectivity from the upstream neuron<sup>4</sup>. We implemented the model's ability to predict circuit activity in two systems: feeding initiation and antennal grooming. We regularly examined

<sup>1</sup>Department of Molecular and Cell Biology and Helen Wills Neuroscience Institute, University of California, Berkeley, CA, USA. <sup>2</sup>University of Rochester Medical Center, Department of Neurobiology, New York, NY, USA. <sup>3</sup>Max Planck Florida Institute for Neuroscience, Jupiter, FL, USA. <sup>4</sup>10000 Institute for Genome Sciences and Policy, Johns Hopkins University, Baltimore, MD, USA. <sup>5</sup>Department of Biological Sciences, Singapore National University, Singapore. <sup>6</sup>Department of Neurobiology, Princeton University, Princeton, NJ, USA. <sup>7</sup>Department of Zoology, University of Cambridge, Cambridge, UK. <sup>8</sup>Neurobiology Division, MRC Laboratory of Molecular Biology, Cambridge, UK. <sup>9</sup>Center for Neural Circuits and Behavior, The University of Texas at Dallas, Dallas, TX, USA. <sup>10</sup>Department of Neurobiology and Behavioral and Cognitive Neuroscience, Howard Hughes Medical Institute, Chevy Chase, MD, USA. <sup>11</sup>Department of Neurobiology, University of Pennsylvania, Philadelphia, PA, USA. <sup>12</sup>Department of Psychology, University of California, San Diego, San Diego, CA, USA. <sup>13</sup>Department of Psychology, University of California, San Diego, San Diego, CA, USA. <sup>14</sup>Department of Psychology, University of California, San Diego, San Diego, CA, USA. <sup>15</sup>Department of Psychology, University of California, San Diego, San Diego, CA, USA. <sup>16</sup>Department of Psychology, University of California, San Diego, San Diego, CA, USA. <sup>17</sup>Department of Psychology, University of California, San Diego, San Diego, CA, USA. <sup>18</sup>Department of Psychology, University of California, San Diego, San Diego, CA, USA. <sup>19</sup>Department of Psychology, University of California, San Diego, San Diego, CA, USA. <sup>20</sup>Department of Psychology, University of California, San Diego, San Diego, CA, USA. <sup>21</sup>Department of Psychology, University of California, San Diego, San Diego, CA, USA. <sup>22</sup>Department of Psychology, University of California, San Diego, San Diego, CA, USA. <sup>23</sup>Department of Psychology, University of California, San Diego, San Diego, CA, USA. <sup>24</sup>Department of Psychology, University of California, San Diego, San Diego, CA, USA. <sup>25</sup>Department of Psychology, University of California, San Diego, San Diego, CA, USA. <sup>26</sup>Department of Psychology, University of California, San Diego, San Diego, CA, USA. <sup>27</sup>Department of Psychology, University of California, San Diego, San Diego, CA, USA. <sup>28</sup>Department of Psychology, University of California, San Diego, San Diego, CA, USA. <sup>29</sup>Department of Psychology, University of California, San Diego, San Diego, CA, USA.

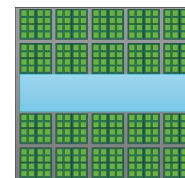
# BRAIN SIMULATIONS ARE THE IDEAL APPLICATION



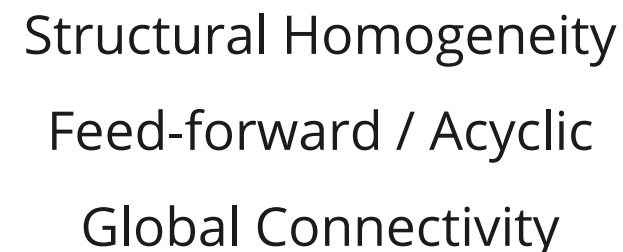
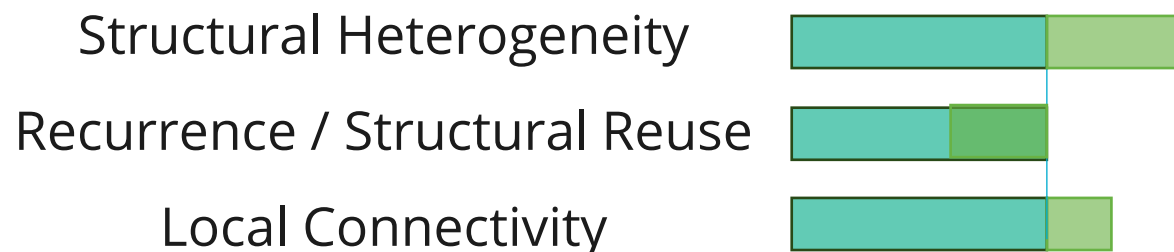
## The Brain / Neuromorphic



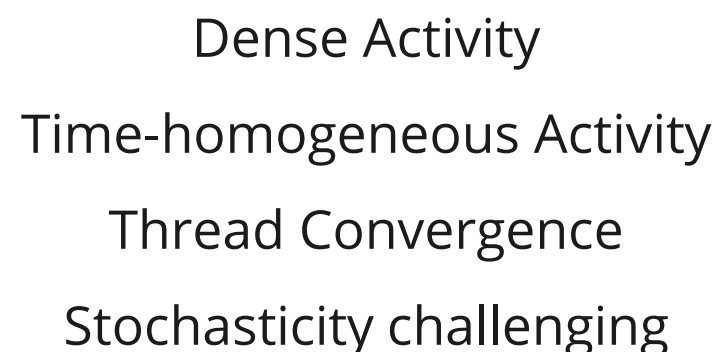
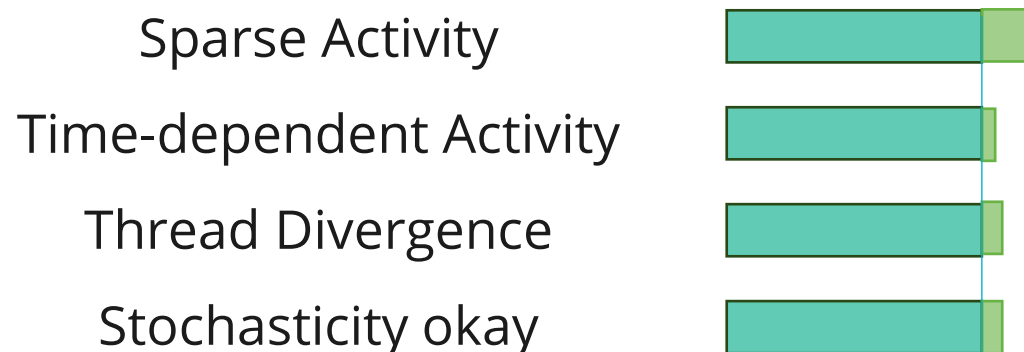
## GPUs / SIMD Architectures



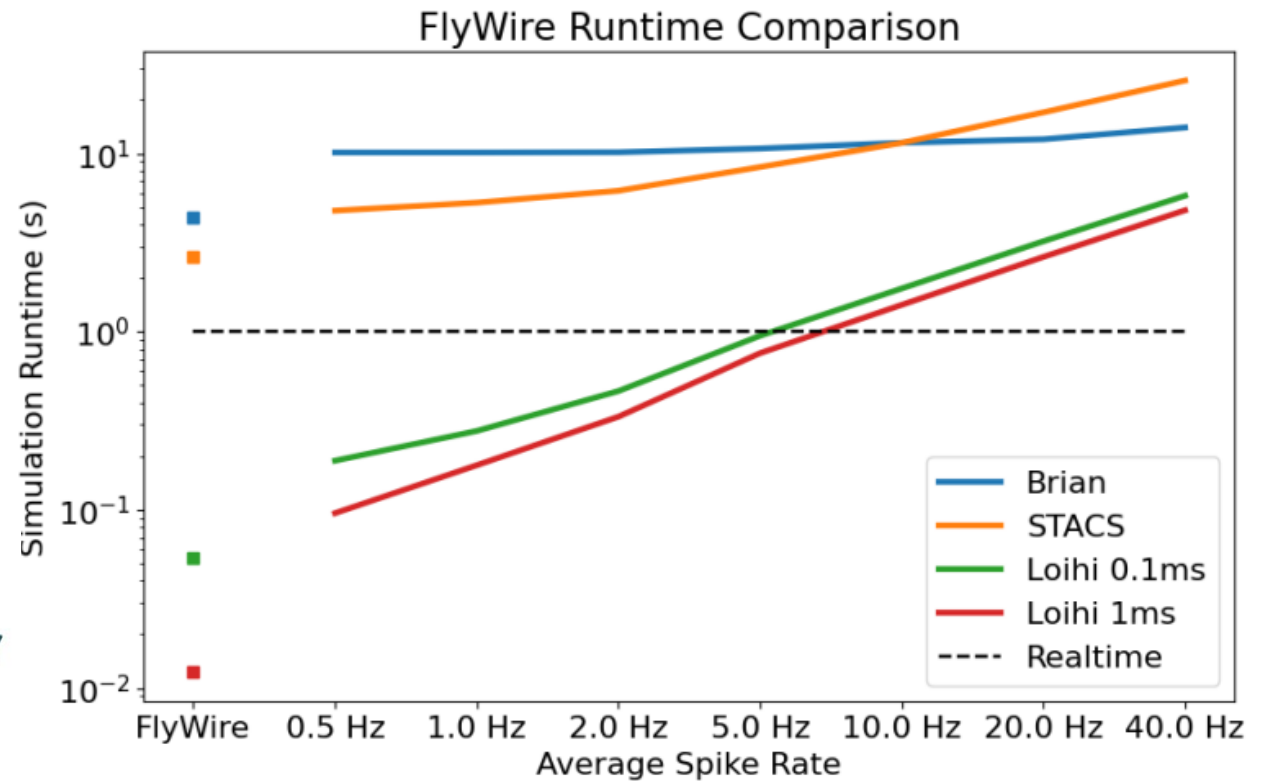
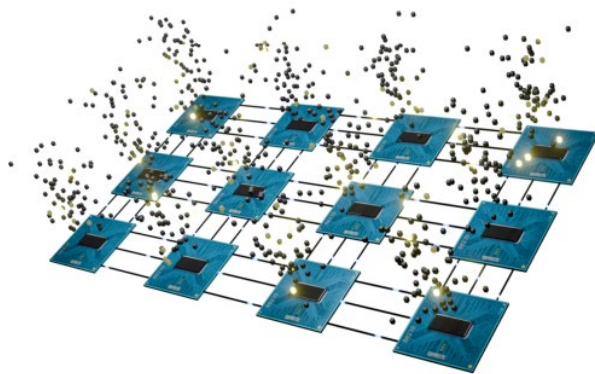
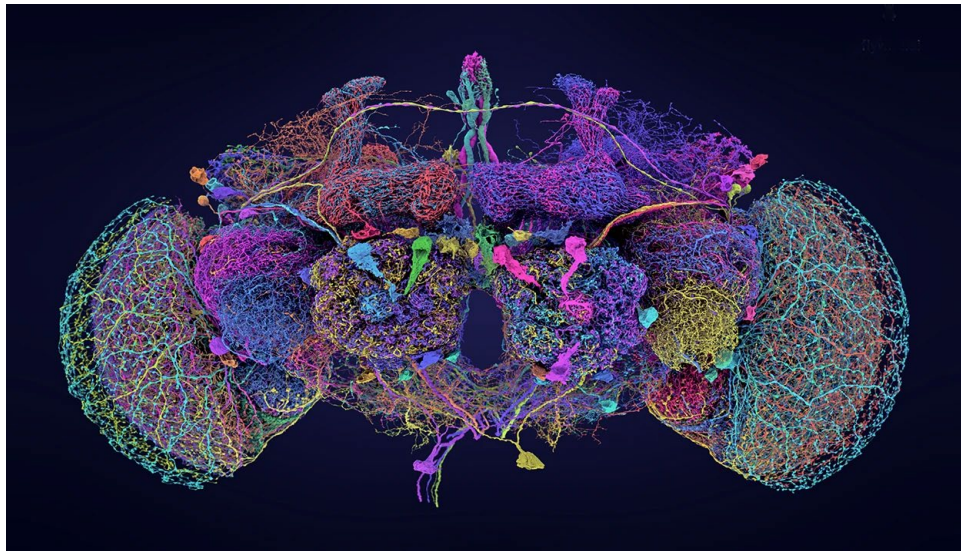
### Structural Metrics



### Trace Metrics



# ALGORITHM MATCHES THE ARCHITECTURE: NEUROMORPHIC HARDWARE CAN ACCELERATE BIOLOGICAL SIMULATIONS BY >100X

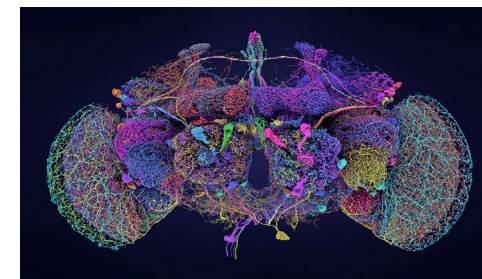
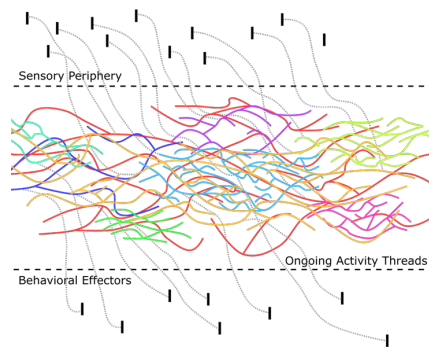
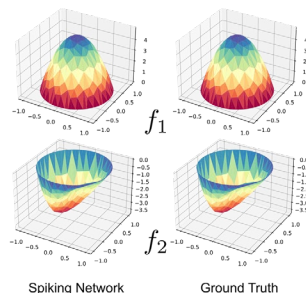
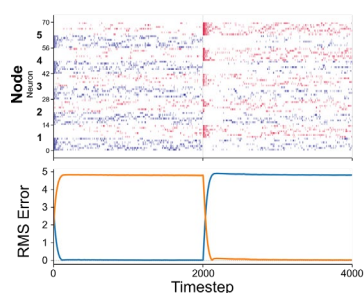


Wang et al., Arxiv 2025  
<https://arxiv.org/pdf/2508.16792>



# The **algorithm** must match the **architecture**.

- We must break out of representation-centric connectionist mindset
- Algorithms based on trajectories of spikes over time and space?
- Is the goal to understand uncertainties and distributions?
- *Time to revisit what learning means?*
- Processing in memory
- Event-driven
- Local & sparse activity
- *Heterogeneity, stochasticity, divergence okay!*



# THANK YOU!



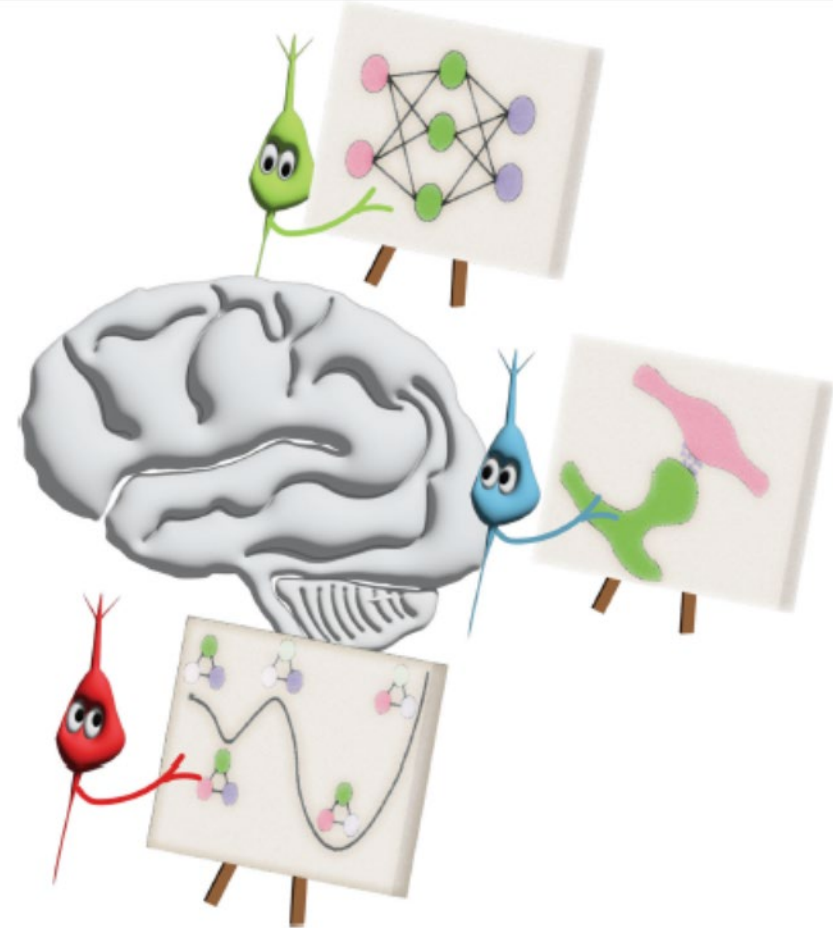
Any questions to [jbaimon@sandia.gov](mailto:jbaimon@sandia.gov)



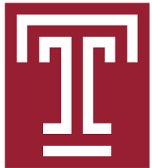
**Neural Exploration & Research Lab**  
COGNITIVE & EMERGING COMPUTING

Craig M Vineyard, Suma G Cardwell, Corinne M Teeter, William Severa, J Darby Smith, Felix Wang, Fred Rothganger, Michael Krygier, Cale Crowder, Bradley H Theilman, William Chapman, Ryan Dellana, Mark Plagge, Efrain Gonzalez, Srideep Musuvathy

Aaron Hill, Shashank Misra, Yang Ho, Brady Taylor, Chris Allemang, Rich Lehoucq, Brian Franke, Ojas Parekh



UC San Diego



UTSA



U.S. DEPARTMENT OF  
**ENERGY**

Office of  
Science



LABORATORY DIRECTED  
RESEARCH & DEVELOPMENT

