# Braindrop: A Mixed-Signal Neuromorphic System that Presents Clean Abstractions

Kwabena Boahen*
Bioengineering & Electrical Engineering
Stanford University

*Cofounder & Chief Scientific Adviser, Femtosense Inc.*

*26 March 2019*

*Chris Eliasmith*

# Deep learning is huge —in the cloud

* Backprop learning is powerful

  * Networks <u>deep in space or time</u>

    * Space is discretized into <u>layers</u>

    * Time is discretized into <u>steps</u>

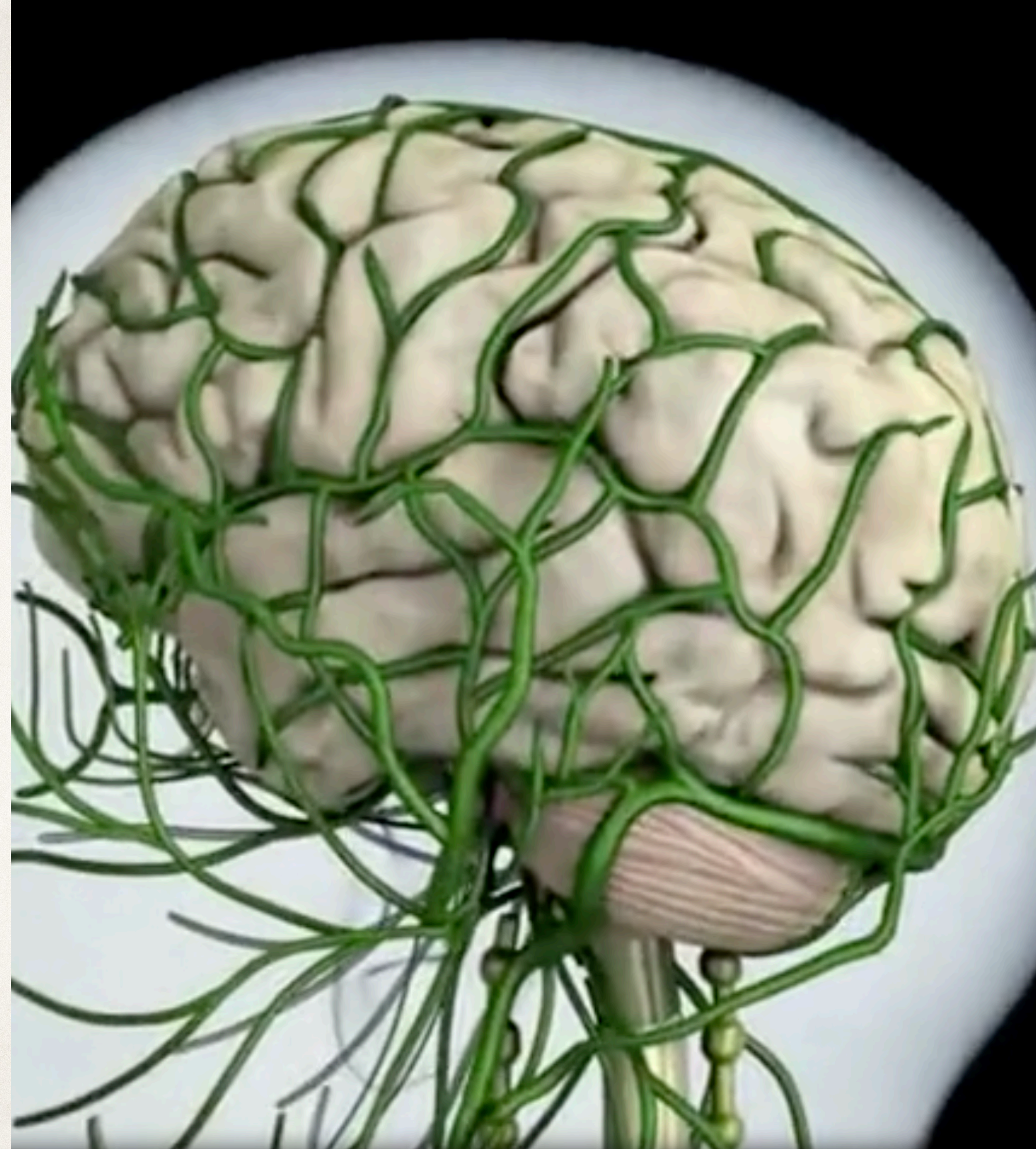  * Unit's output must be <u>differentiable</u> (with respect to outputs of units feeding it)

# Backprop's constraints limit design-space

✤ Cannot take advantage of:

  ✤ Physical <u>space</u> (its continuous)

  ✤ Real <u>time</u> (its also continuous)

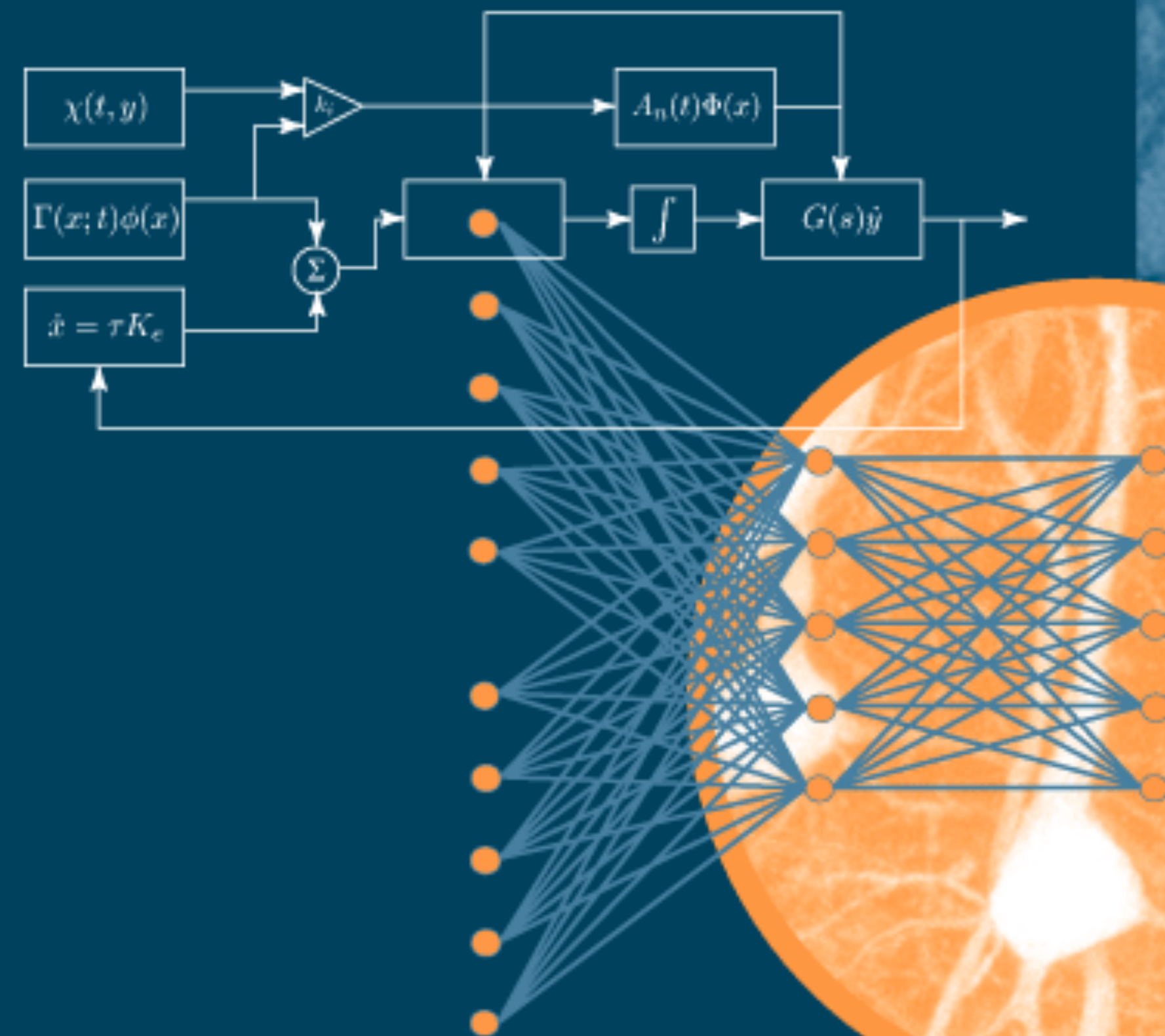  ✤ Non-differentiable signals (e.g., <u>spikes</u>)

# How do we relax its constraints? (Part I)

✤ <u>Map</u> functional abstractions onto physical ones

✤ Two existing examples:

  ✤ <u>N</u>eural <u>E</u>ngineering <u>F</u>ramework (Eliasmith & Anderson 2003)

  ✤ <u>P</u>redictive <u>C</u>oding <u>F</u>ramework (Deneve et al. 2014)

## Neural Engineering

COMPUTATION, REPRESENTATION, AND DYNAMICS IN NEUROBIOLOGICAL SYSTEMS

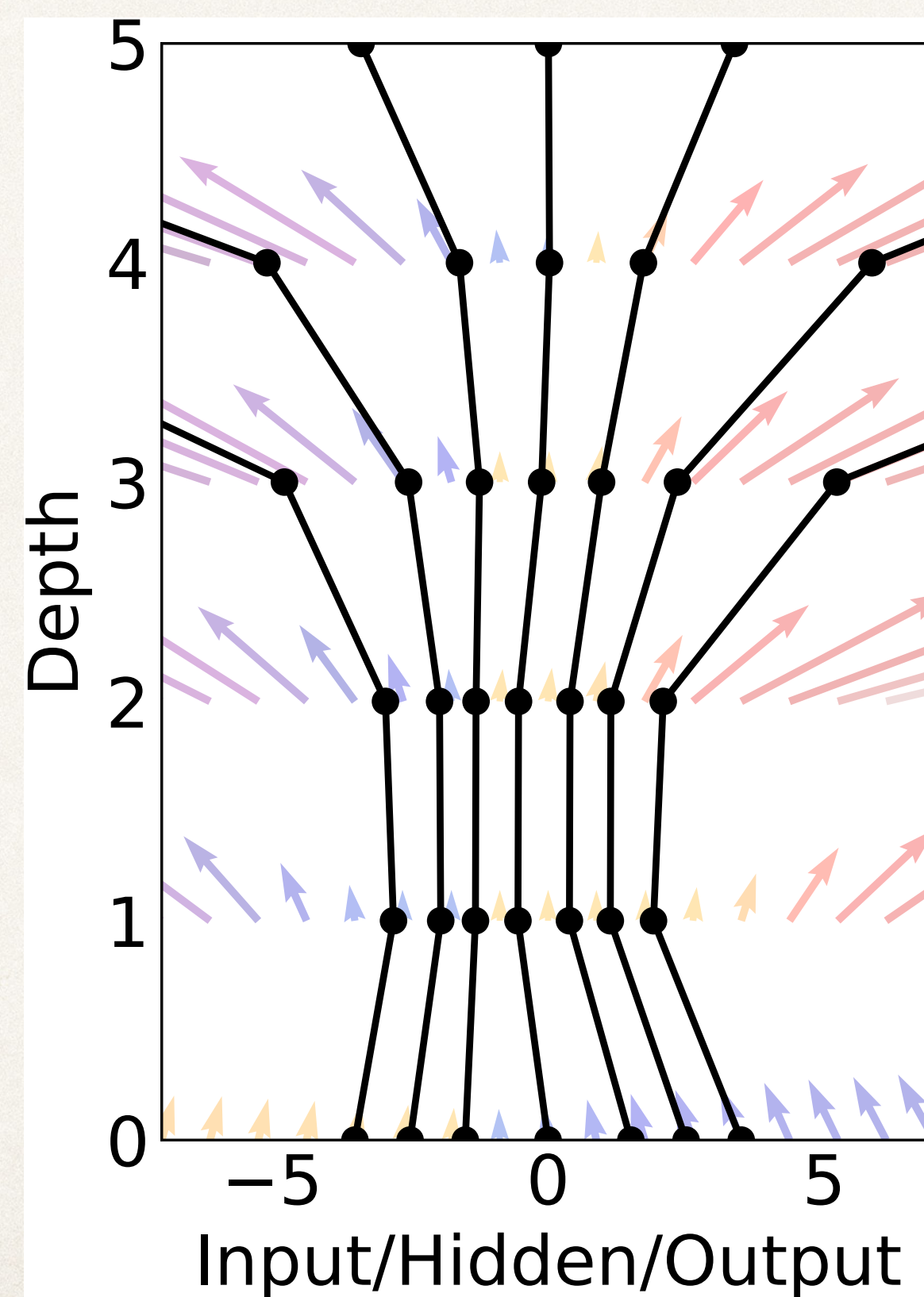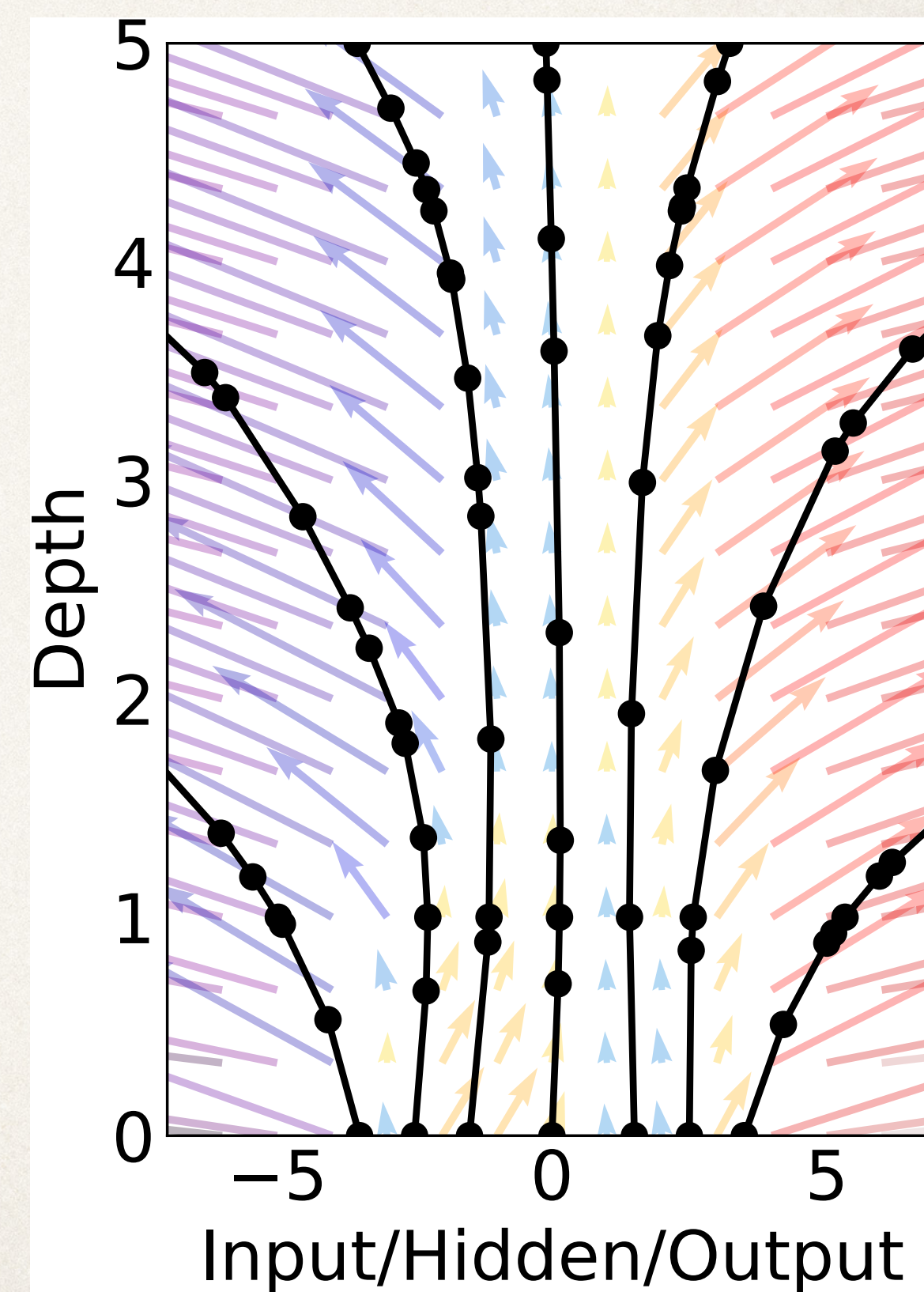Chris Eliasmith and Charles H. Anderson

$\chi(t, y)$

$\Gamma(x; t)\phi(x)$

$\dot{x} = \tau K_e$

$k_i$

$\Sigma$

$A_n(t)\Phi(x)$

$\int$

$G(s)\hat{y}$

# How do we relax its constraints? (Part II)

- ✤ Train networks continuous in time and space

  - ✤ Known as <u>dynamical systems</u>

- ✤ An existing example:

  - ✤ <u>N</u>eural <u>O</u>rdinary <u>D</u>ifferential Equations (Duvenaud et al. 2018)
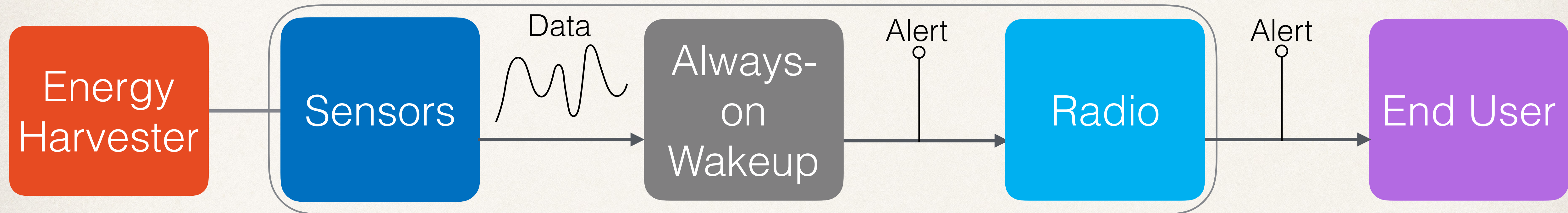


Residual Network

ODE Network

# What's the payoff? Learning at the edge



✤ Exploit <u>physical primitives</u> to implement physical abstractions

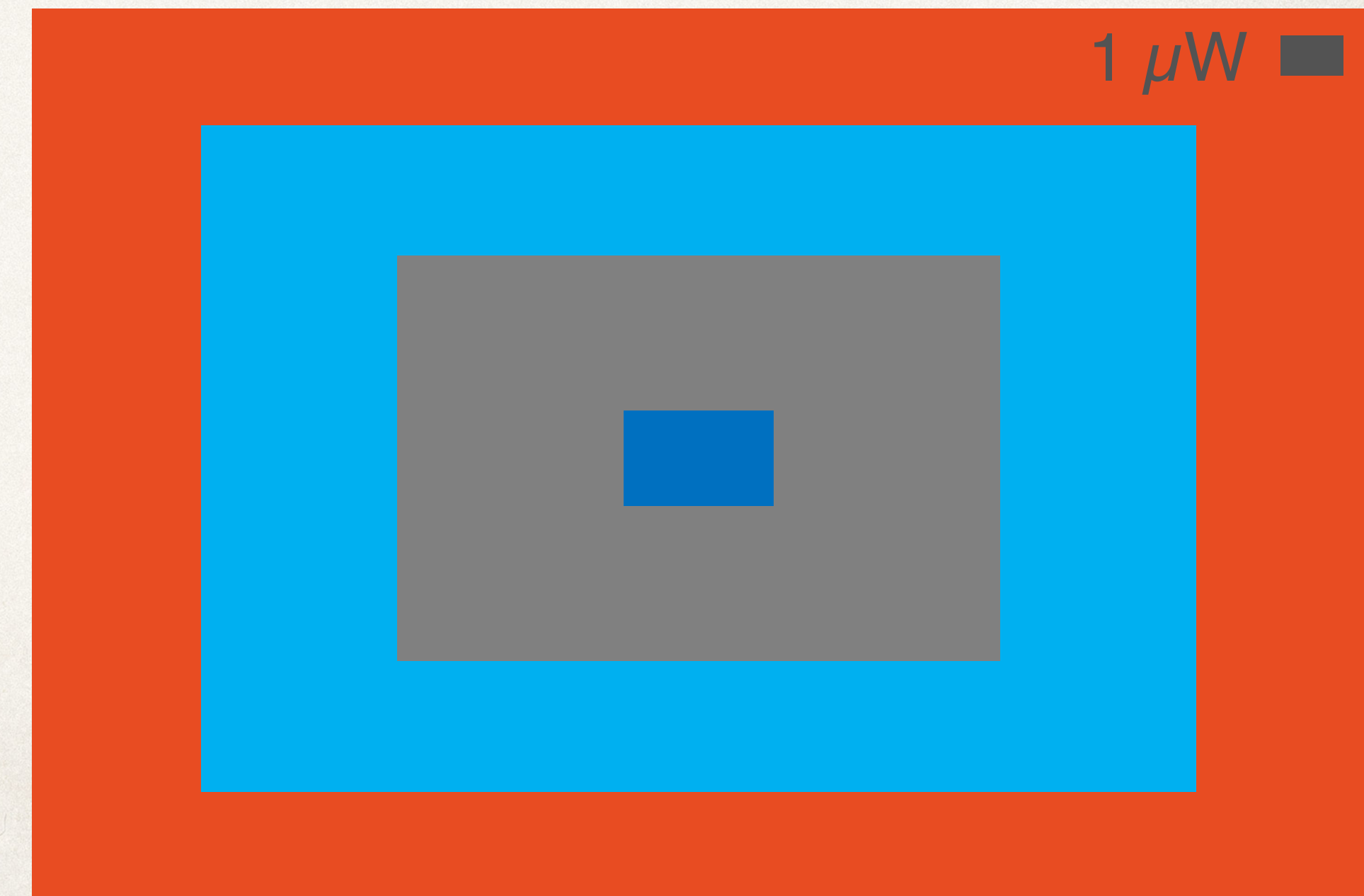✤ Reap dramatic gains in <u>energy-efficiency</u>

**Harvest**
Vibration: 500 $\mu$W

**Sense**
Accelerometer: 6 $\mu$W

**Compute**
Neuroprocessor: 100 $\mu$W

**Communicate**
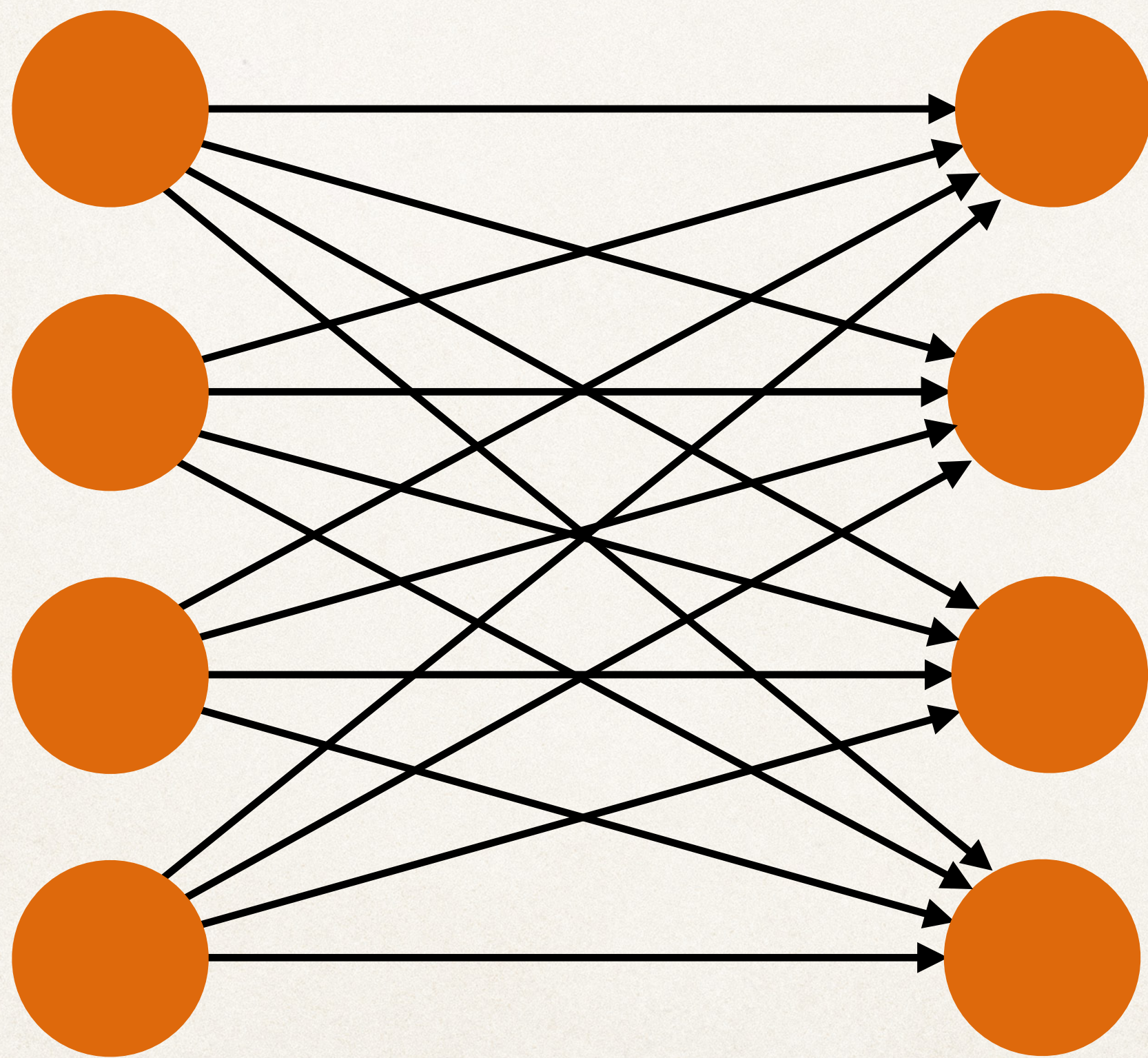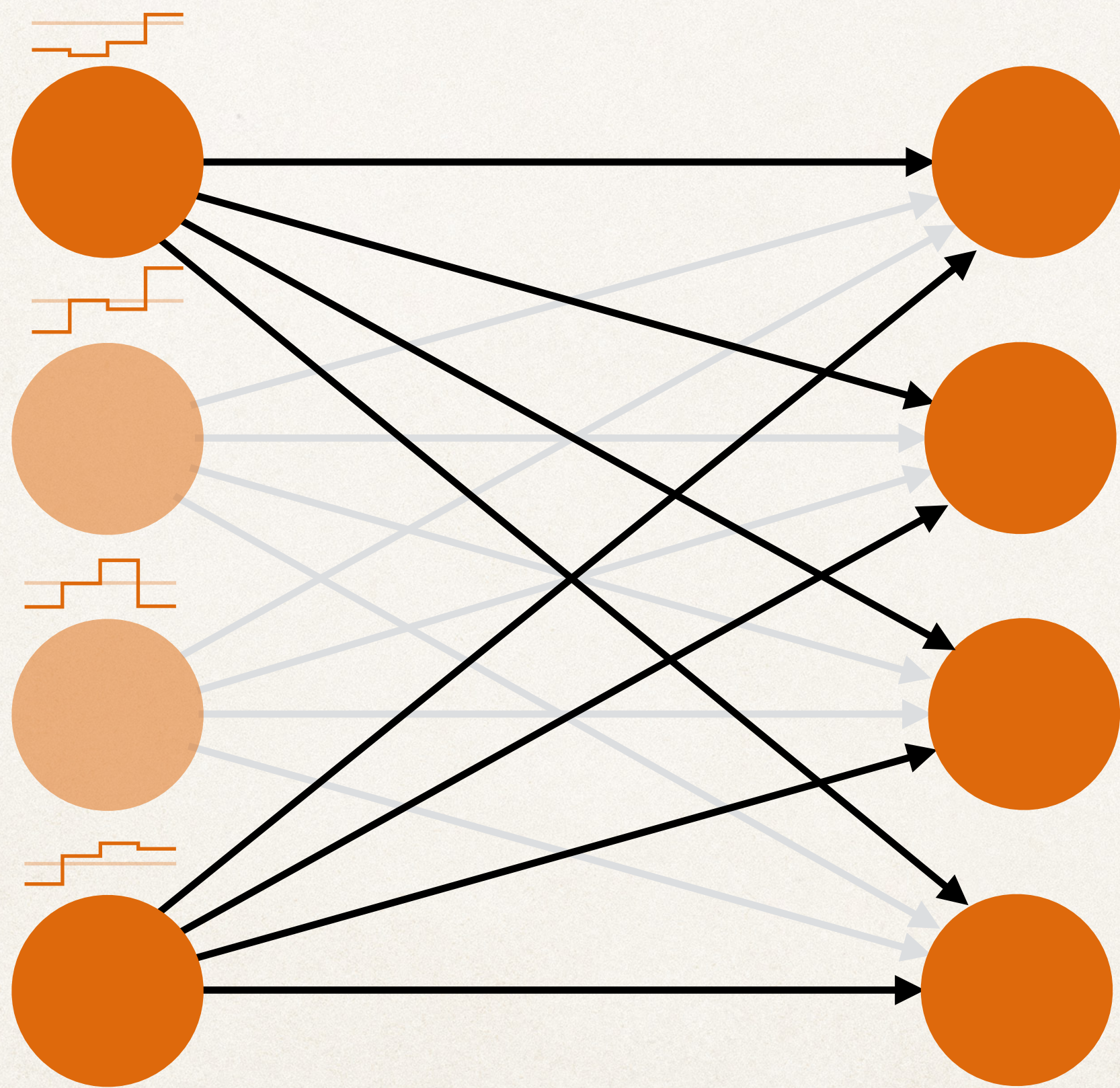BTLE (1% on-time): 280 $\mu$W

1 $\mu$W

# Minimizing energy



$$E_{op} = N_{active} N_{conn}$$
$$= (\rho_{active} N)(\rho_{conn} N)$$

# Minimizing energy: Temporal sparsity



$$E_{op} = N_{active} N_{conn}$$
$$= (\rho_{active} N)(\rho_{conn} N)$$

# Temporal sparsity: Spikes



$$E_{op} = N_{active} N_{conn}$$
$$= (\rho_{active} N)(\rho_{conn} N)$$

# Minimizing energy: Spatial sparsity



$$E_{op} = N_{active}N_{conn}$$
$$= (\rho_{active}N)(\rho_{conn}N)$$

# Spatial sparsity: Analog convolving

$$E_{op} = N_{active} N_{conn}$$
$$= (\rho_{active} N)(\rho_{conn} N)$$

COMPUTATION

ANALOG | DIGITAL

COMMUNICATION

ANALOG | DIGITAL

Analog computer
1900-1950

Digital computer
1950-????

Make point about top-left corner being lowest power

# Digital versus Analog: 1 day versus 1000 yrs



**4.3B-transistor processor**
**10.35Wh battery**

*28 nm FDSOI thick-oxide transistor*

2.5 hours
1 day
10 days
100 days
2.74 years
27.4 years
274 years
2747 years

# Analog Challenge I: Heterogeneity

# Silicon neurons' tuning-curves (Braindrop)



484 Braindrop
neurons at 26°C

# Analog Challenge I: Thermal Sensitivity

$$I_\mu(T) = I_{0_{\mathrm{nom}}} e^{\langle \gamma_1 \rangle \left(1 - \frac{T_{\mathrm{nom}}}{T}\right)} e^{\frac{(1-\kappa)V_{\mathrm{BS}}}{U_T}}$$

$$\times e^{\frac{\kappa V_{\mathrm{GS}}}{U_T}} e^{\left(\lambda_1 \frac{T_{\mathrm{nom}}}{T} + \lambda_2\right)\Delta V_{\mathrm{DS}}} \left(1 - e^{-\frac{V_{\mathrm{DS}}}{U_T}}\right)$$

* A subthreshold transistor's current ($I_\mu$) is exponentially sensitive to temperature

  * $T$ is the absolute temperature

  * $U_T = kT/q$ is the thermal voltage

* Across a 50°C range, the current changes by 1.5 to 3 decades

# Tuning-curves' thermal sensitivity (Braindrop)



*Reid, Montoya, & Boahen 2019*

4 Braindrop neurons for 0 to 38°C

# Approximating functions

$$\mathbf{f} = \mathbf{A}\mathbf{d} \Rightarrow \mathbf{d} = (\mathbf{A}^{\mathrm{T}}\mathbf{A})^{-1}\mathbf{A}^{\mathrm{T}}\mathbf{f}$$

✤ The desired function $f(x)$ is expressed as a weighted sum of the neural *tuning curves* $a_i(x)$

✤ The weights—called decoders—are labeled $d_i$

$$
\begin{bmatrix} f(x_1) \\ f(x_2) \\ \vdots \\ f(x_Q) \end{bmatrix}
=
\begin{bmatrix} a_1(x_1) & a_2(x_1) & \cdots & a_N(x_1) \\ a_1(x_2) & a_2(x_2) & \cdots & a_N(x_2) \\ \vdots & \vdots & \ddots & \vdots \\ a_1(x_Q) & a_2(x_Q) & \cdots & a_N(x_Q) \end{bmatrix}
\begin{bmatrix} d_1 \\ d_2 \\ \vdots \\ d_N \end{bmatrix}
$$

3 to 30 Braindrop neurons at 26°C

# Thermally robust computation (Braindrop)



256 Braindrop neurons from 0 to 38°C

*Reid, Montoya, & Boahen 2019*

# NEF: Decode-Transform-Encode

*Eliasmith & Anderson 2003*



**❶ Spike**

$$\langle \delta_{x_i} \rangle_t = \lceil 0, \alpha_i(I_i + \beta_i) \rceil$$

Somas emit unit-area deltas $\delta_{x_i}$ at rates $\langle \delta_{x_i} \rangle_t$ dictated by their input current $I_i$.

**❷ Decode**

$$\delta_{y_j} = \sum_i D_{ji}\delta_{x_i}, D \in \mathbb{R}^{D \times N}$$

Deltas are then scaled by their decode weight and merged together.

**❸ Transform**

$$\delta_{z_k} = \sum_j T_{kj}\delta_y, T \in \mathbb{R}^{D \times D}$$

Transform works the same way as Decode.

**❹ Encode**

$$\tau \dot{I}_l = -I_l + \sum_k E_{lk}\delta_{z_k}, E \in \mathbb{R}^{N \times D}$$

Synaptic filters superpose and low-pass filter weighted deltas to produce output currents ($I_l$) that feed the next soma layer.

# Digital thinning and analog convolving



**❶ Spike**

$$\langle \delta_{x_i} \rangle_t =$$
$$\lceil 0, \alpha_i (I_i + \beta_i) \rceil$$

Somas emit
delta trains
(as in Figure 2).

**❷ Decode**

$$\langle \delta_{y_j} \rangle_t = \sum_i D_{ji} \langle \delta_{x_i} \rangle_t$$
$$D \in [-1, 1]^{D \times N}$$

Weighted deltas are
accumulated to produce a
stream of unit-area deltas.

**❸ Transform**

$$\langle \delta_{z_k} \rangle_t = \sum_j T_{kj} \langle \delta_{y_j} \rangle_t$$
$$T \in [-1, 1]^{D \times D}$$

Transform still works the
same as Decode ($T_{kj}=1$
in this example).

**❹ Sparse Encode**

$$\tau I_l = -I_l + \sum_k S_{lk} \delta_{z_k}$$
$$S \in \{-1, 0, 1\}^{N \times D}$$

Each accumulator's deltas
are sent to a subset of
synaptic filters (tap-points).

**❺ Convolve**

$$I_m = \sum_l \hat{d}_\gamma (m - l) I_l$$

Filter outputs ($I_l$) are
convolved ($I_m$) and sent
to the next soma layer.

**1** Somas' spikes are muxed into address-events.

**2** Address-events are translated into base-addresses.

**3** Decoding vector and buckets' states are retrieved and the latter are updated. A tag is emitted if threshold is exceeded.

**4** Tags are buffered.

**5** Tags are translated into base-addresses.

**6** Corresponding transform's column and associated buckets' states are retrieved, states updated, and tags emitted.

**7** Tags are translated into address-events with associated polarities.

**8** Address-events are demuxed to excite or inhibit synaptic filters.

**9** Filters' output currents are distributed to somas.

IO/Router

WM   AM

ACCUMULATOR

FIFO   TAT

AER:RX   Synaptic Filters   Diffusor   Somas   AER:TX   PAT

# Braindrop: 4096 neurons in 28nm FDSOI CMOS



2 mm

# Programming Environment (Nengo)



$$a = f(x)$$
$$b = g(y)$$
$$\dot{z} = h(z) + a + b$$

```
1   import nengo
2   import numpy as np
3   model = nengo.Network()
4 ▾ with model:
5       x = nengo.Node(lambda t: np.cos(2*np.pi*t))
6       y = nengo.Node(lambda t: np.cos(4*np.pi*t))
7       a = nengo.Ensemble(n_neurons=256, dimensions=1)
8       b = nengo.Ensemble(n_neurons=256, dimensions=1)
9       nengo.Connection(x, a)
10      nengo.Connection(y, b)
```

**Encode** ①
Synaptic filters' states, $x(t)$, are projected onto encoding vectors, $E_i$, to drive somas.

**Somas** ②
The $i^{\text{th}}$ soma's input current, $J_i = E_i x(t)$, drives it to spike at rate $\langle \delta_{x_i} \rangle_t$.

$x_1(t)$
$c_1(t) \rightarrow$

$c_2(t) \rightarrow$
$x_2(t)$

$E_{i1}$   $D_{1i}$
$D_{1i}\delta_{x_i}$
$E_{i2}$   $D_{2i}$

**Decode** ③
Deltas with area equal to components of the neuron's decoding vector, $D_i^{\text{T}}$, replace spikes.

④ **Synaptic Filters**
Merged delta trains are lowpass filtered.
As a result, $\tau_{\text{syn}}\dot{x}(t) + x(t) = \sum_i D_i^{\text{T}} \delta_{x_i} + c(t)$,
where $c(t)$ is a vector of injected currents.

✤ To emulate the dynamical system

$$\tau_{\text{dyn}}\dot{x}(t) = f(x) + u(t)$$

✤ Choose decoding weights such that, after synaptic filtering,

$$\sum_i D_i^{\text{T}} \delta_{x_i} \approx \tau_{\text{syn}}/\tau_{\text{dyn}} f(x) + x$$

✤ And set

$$c(t) = \tau_{\text{syn}}/\tau_{\text{dyn}} u(t)$$

*Eliasmith & Anderson 2003*

# Tapped delay-line (Braindrop)



Delay Network    Delay Length (s)    State Vector

384 Braindrop neurons

$$\theta\ddot{\mathbf{x}}(t) = A\mathbf{x}(t) + Bc(t)$$

$$c(t - \theta') \approx C_{\theta'/\theta}\mathbf{x}(t), \quad 0 \leq \theta' \leq \theta$$

*Neckar et al. 2019*

*Voelker & Eliasmith 2017*

# Measured Energy/op (pJ)

**1.1**

**3.8**

**15**

**28**

**6.4**

**①** Somas' spikes are muxed into address-events.

**②** Address-events are translated into base-addresses.

**③** Decoding vector and buckets' states are retrieved and the latter are updated. A tag is emitted if threshold is exceeded.
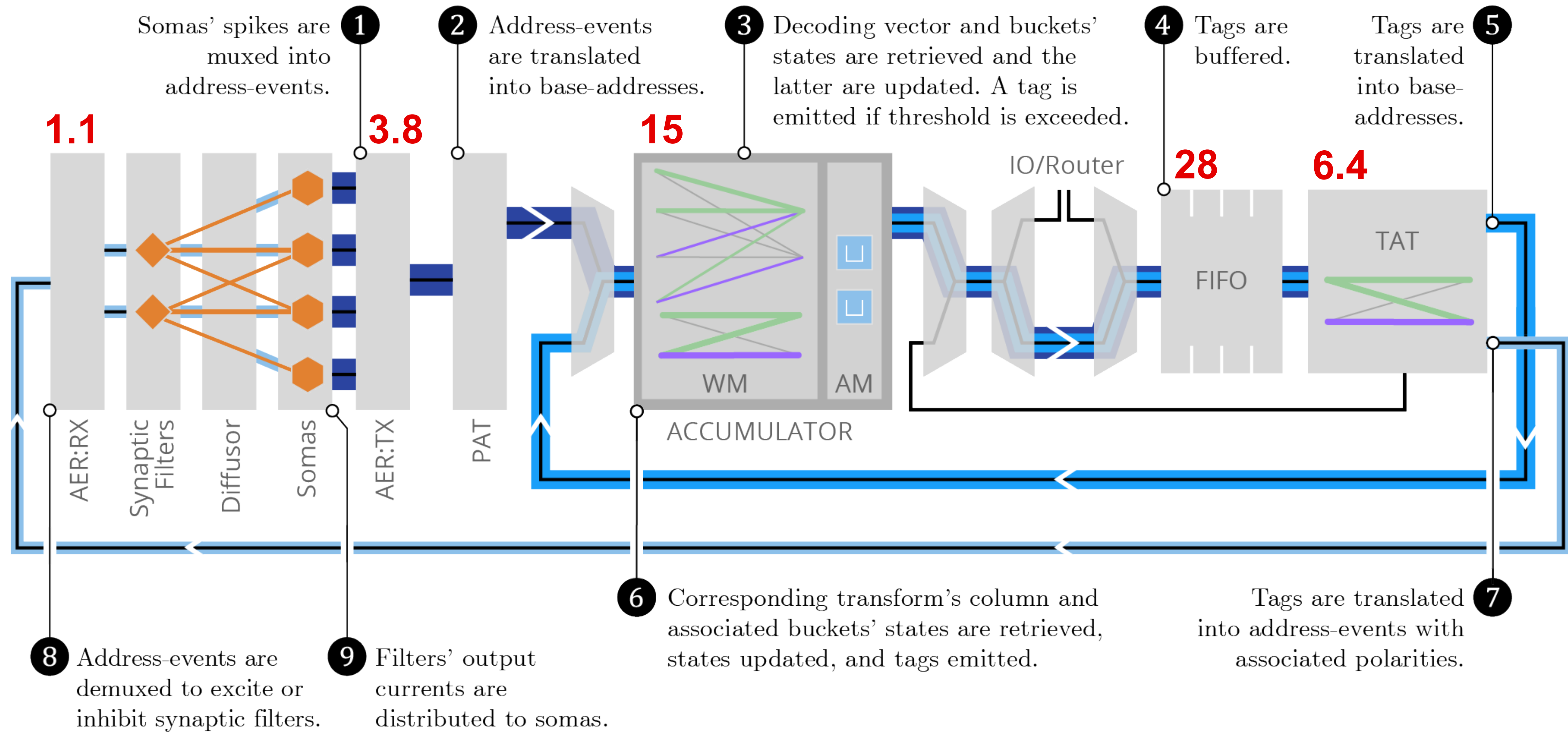
**④** Tags are buffered.

**⑤** Tags are translated into base-addresses.

IO/Router

AER:RX

Synaptic Filters

Diffusor

Somas

AER:TX

PAT

WM

AM

ACCUMULATOR

FIFO

TAT

**⑥** Corresponding transform's column and associated buckets' states are retrieved, states updated, and tags emitted.

Tags are translated into address-events with associated polarities.

**⑦**

**⑧** Address-events are demuxed to excite or inhibit synaptic filters.

**⑨** Filters' output currents are distributed to somas.
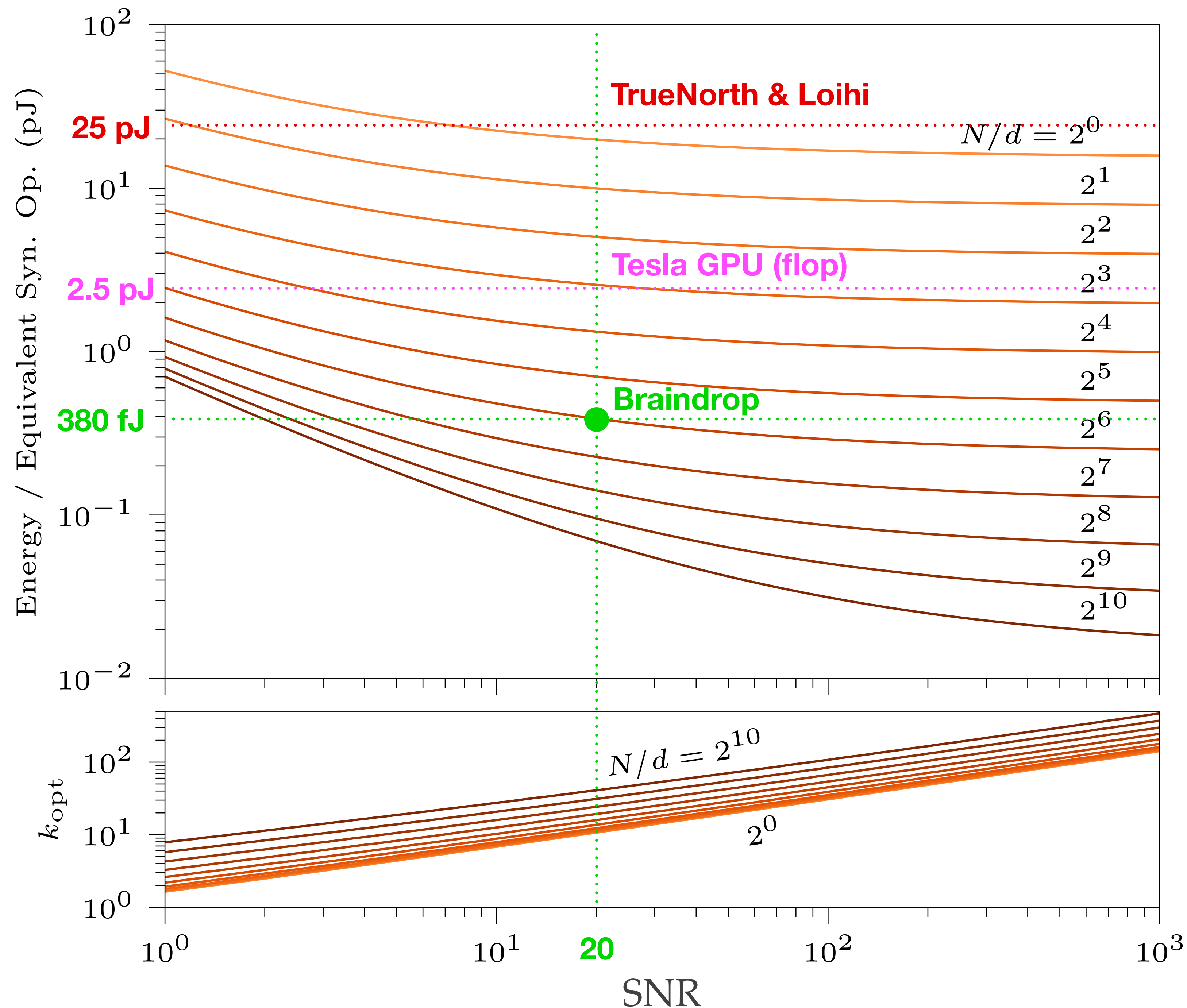
* Analog convolving fans out $d$ spike-trains to $N$ neurons; sparsifies spatially by $d/N$

* Digital thinning lets one per SNR spikes through; sparsifies temporally by $1/\text{SNR}$

# Task performance

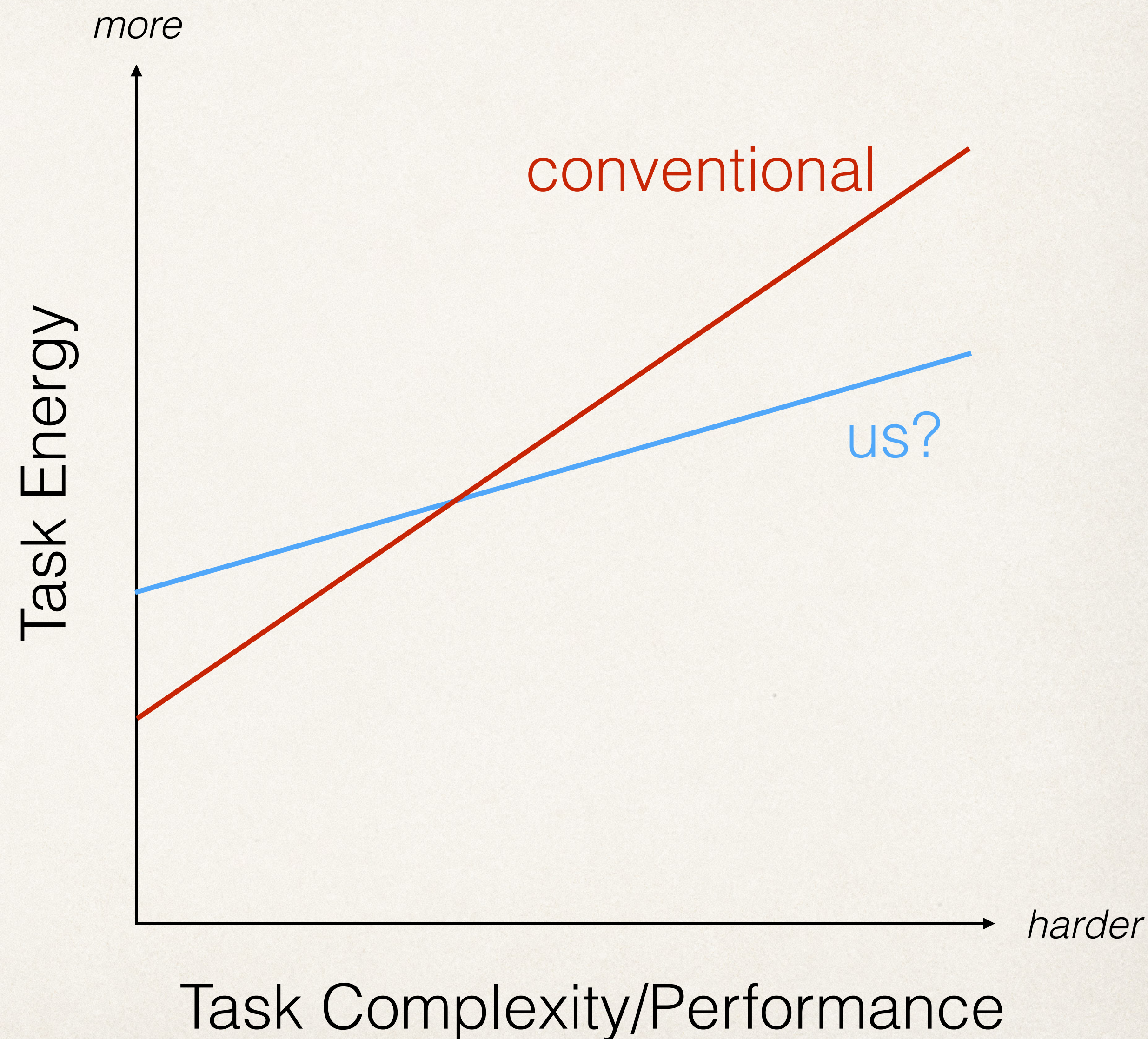- ✤ Two components:

  - ✤ Network design

  - ✤ Hardware design

$$E_{sys}(task) = E_{HW}(R_{net}(task))$$

$$E_{op} = N_{active} N_{conn}$$
$$= (\rho_{active} N)(\rho_{conn} N)$$

# Acknowledgments

# To learn more …

J Dethier, P Nuyujukian, C Eliasmith, T Stewart, S A Elassaad, K V Shenoy, and K Boahen, **A Brain-Machine Interface Operating with a Real-Time Spiking Neural Network Control Algorithm**, *Advances in Neural Information Processing Systems 24*, Curran Associates, Inc., pp 2213-21, 2011.

S Choudhary, S Sloan, S Fok, A Necker, E Trautmann, P Gao, T Stewart, C Eliasmith, and K Boahen, **Silicon Neurons that Compute**, *International Conference on Artificial Neural Networks, LNCS* vol VV, pp 121-128, Springer, Heidelberg, 2012.

S Menon, S Fok, A Neckar, O Khatib, and K Boahen, **Controlling Articulated Robots in Task-Space with Spiking Silicon Neurons**, *IEEE International Conference on Biomedical Robotics and Biomechatronics (BioRob)*, IEEE Press, pp 181-186, 2014.

K Boahen, **A Neuromorph's Prospectus**, *Computing in Science & Engineering*, vol 19, no 2, pp 14-28, IEEE Computer Society, Los Alamitos CA, USA, 2017.

E Kauderer-Abrams, A Gilbert, A Voelker, B Benjamin, and T C Stewart, and K Boahen, **A Population-Level Approach to Temperature Robustness in Neuromorphic Systems**, *IEEE International Symposium on Circuits and Systems (ISCAS)*, Baltimore MD, 2017.

A R Voelker, B V Benjamin, T C Stewart, K Boahen, and C Eliasmith, **Extending the Neural Engineering Framework for Nonideal Silicon Synapses**, *IEEE International Symposium on Circuits and Systems (ISCAS)*, Baltimore MD, 2017.

E Kauderer-Abrams and K Boahen, **Calibrating Silicon-Synapse Dynamics using Time-Encoding and Decoding Machines**, *IEEE International Symposium on Circuits and Systems (ISCAS)*, Baltimore MD, 2017.

*Proceedings of the IEEE, Jan 2019*

INVITED PAPER

# Braindrop: A Mixed-Signal Neuromorphic Architecture With a Dynamical Systems-Based Programming Model

*This paper provides an overview of a current approach for the construction of a programmable computing machine inspired by the human brain.*

By Alexander Neckar, Sam Fok, Ben V. Benjamin, Terrence C. Stewart, Nick N. Oza, Aaron R. Voelker, Chris Eliasmith, Rajit Manohar, *Senior Member IEEE*, and Kwabena Boahen, *Fellow IEEE*