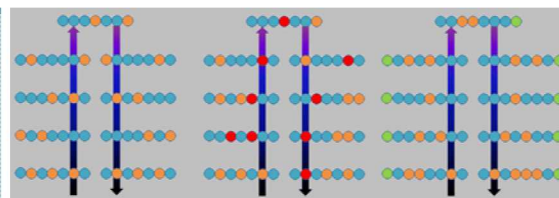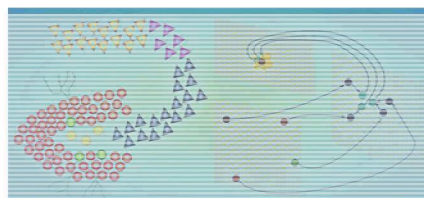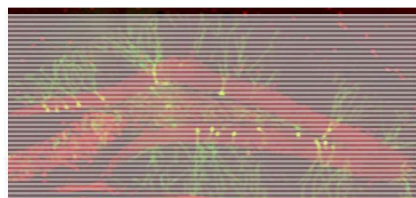# Computing with Spikes

*Everything from the deep learning to numerical applications*
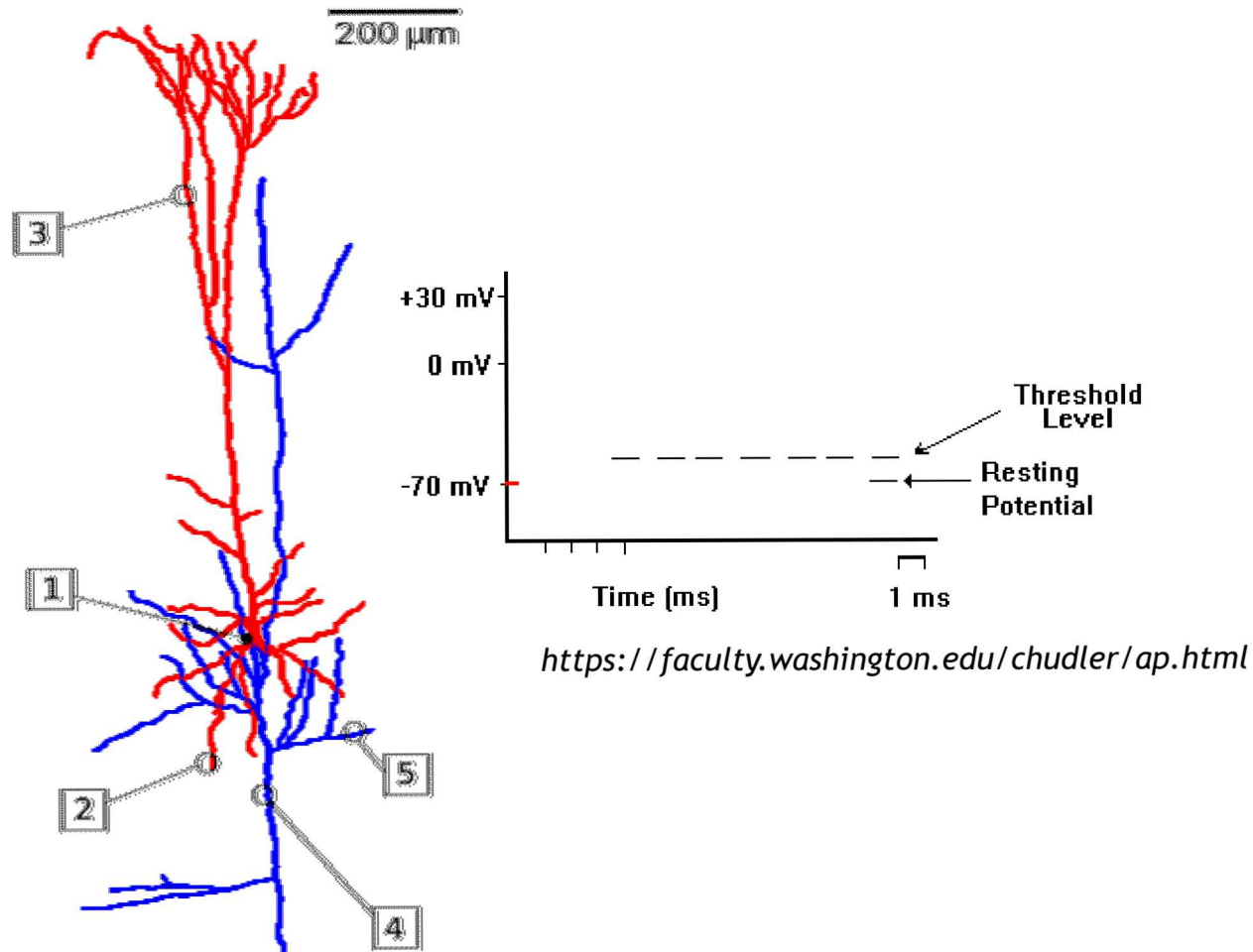
PRESENTED BY

Brad Aimone

# What is spiking?



Pyramidal Cell -- Wikipedia



https://faculty.washington.edu/chudler/ap.html
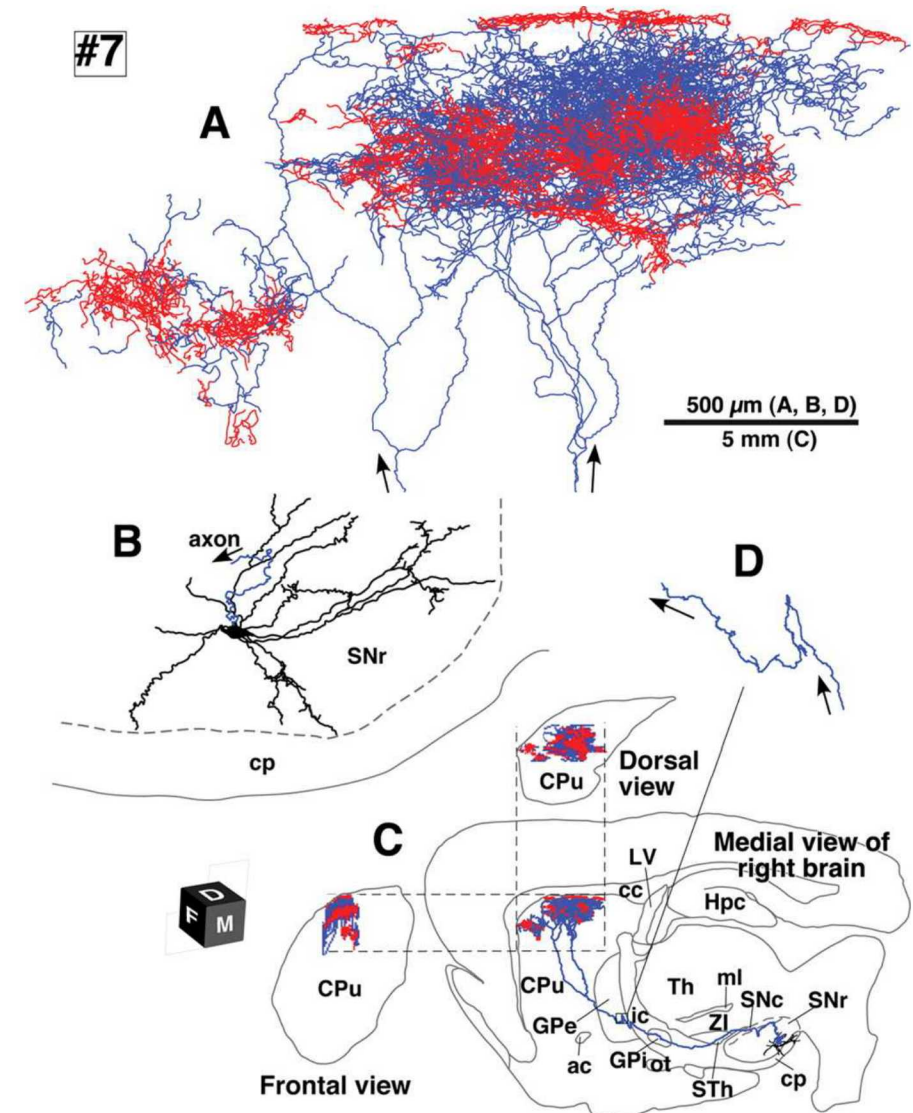


Clockwise from top: IBM TrueNorth, Intel Loihi, SpiNNaker, BrainScales

# Why spiking?

- ➤ Event-driven
  - ➤ Only expend energy when neuron crosses threshold

- ➤ Reliable and efficient over long distances
  - ➤ Neurons often project across brain or whole body…

- ➤ Robust to noise
  - ➤ Away from threshold, biophysical noise should not accidently cause spikes

# Correcting some common misconceptions about spiking

➤ *Spiking is necessary for brain-like computation*
  ➤ Reality: The advantage of spiking is efficiency and reliability over distance, not computability.

    It changes the tradeoffs between time, power, and space

➤ *Spiking does not offer anything for algorithms*
  ➤ Reality: Spiking facilitates developing algorithms that more directly leverage time in computing
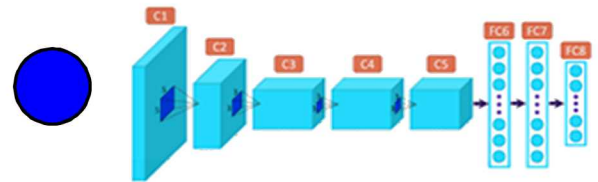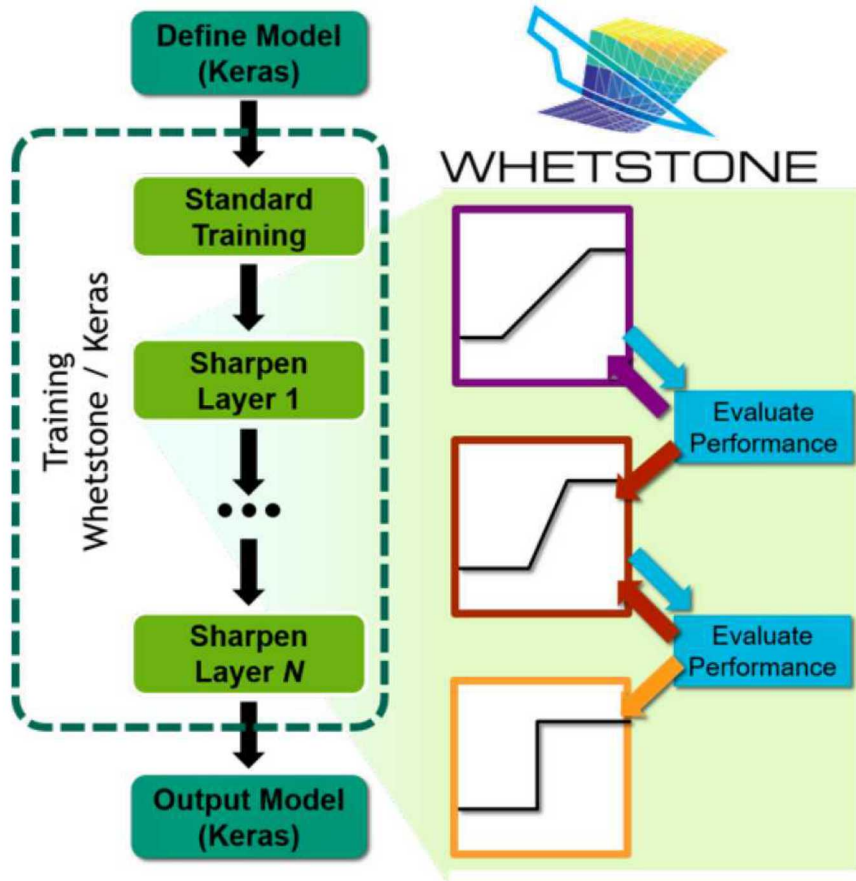
➤ *Spiking reduces the accuracy of algorithms*
  ➤ Reality: Not necessarily! Spiking does lower the precision of communication, but often this precision is unnecessary or can be compensated for in other ways.

➤ *Spiking requires paying a time penalty*
  ➤ Reality: Not always! Some coding schemes are actually time advantageous – e.g., you can implement very fast threshold gate circuit algorithms on spiking hardware
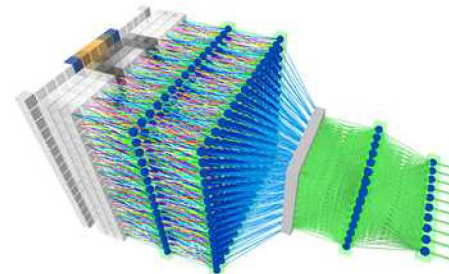
# What can you do with spiking neurons?



**Spiking deep neural networks**
- Whetstone allows us to use spiking communication with *no time penalty* and minimal accuracy reduction

*Severa et al., Nature Machine Intelligence, Feb 2019*
*Vineyard et al., NICE Proceedings, 2019*

# What can you do with spiking neurons?

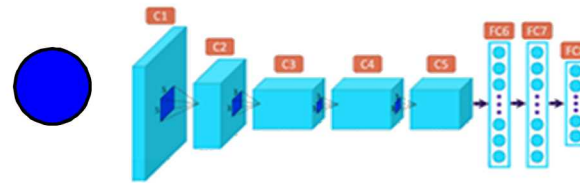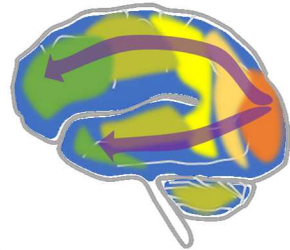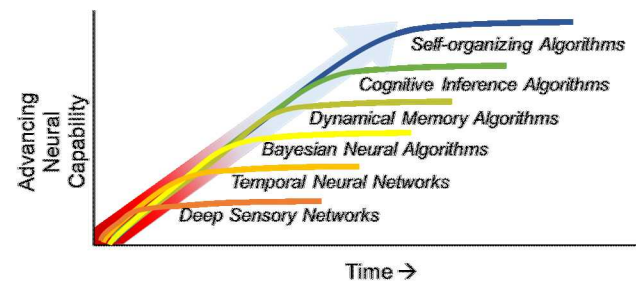| Algorithm Class | Current Algorithms | Inspiration | Application |
|---|---|---|---|
| Deep Vision Processing | Deep Convolutional Networks (VGG, AlexNet, GoogleNet, etc.), HMax, Neocognitron | Hierarchy of sensory nuclei and early sensory cortices | Static feature extraction (e.g., images) & pattern classification |
| Temporal Neural Networks | Deep Recurrent Networks (e.g., long short-term memory), Hopfield Networks | Local recurrence of most biological neural circuits, especially higher sensory cortices | Dynamic feature extraction (e.g., videos, audio) & classification |
| Bayesian Neural Algorithms | Predictive Coding, Hierarchical Temporal Memory, Recursive Cortical Networks | Substantial reciprocal feedback between "higher" and "lower" sensory cortices | Inference across spatial and temporal scales |
| Dynamical Memory and Control Algorithms | Liquid State Machines, Echo State Networks, Neural Engineering Framework | Continual dynamics of hippocampus, cerebellum, and prefrontal and motor cortices | Online learning content-addressable memory & adaptive motor control |
| Cognitive Inference Algorithms | Reinforcement learning (e.g., Deep Q-learning) Neural Turing Machines | Integration of multiple modalities and memory into prefrontal cortex, which provides top-down influence on sensory processing | Context and experience dependent information processing and decision making |
| Self-organizing Algorithms | Neurogenesis Deep Learning | Initial development and continuous refinement of neural circuits to specific input and outputs | Automated neural algorithm development for unknown input and output transformations |

## Spiking deep neural networks
- Whetstone allows us to use spiking communication with *no time penalty* and minimal accuracy reduction

## Neuroscience-constrained algorithms
- Computation incorporates broad range of neural plasticity and dynamics
- *Generally still unexplored from algorithms perspective*

*Aimone JB, Communications of ACM, April 2019*

# What can you do with spiking neurons?

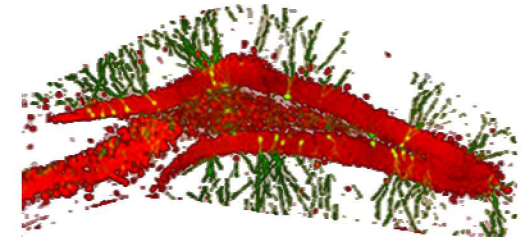*Treat neurons as powerful logic gates*

*Algorithms are circuits...*
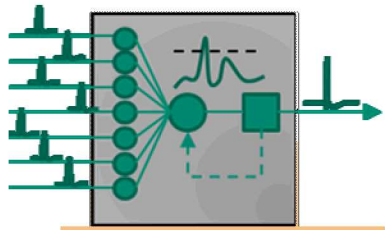
**Spiking deep neural networks**
- Whetstone allows us to use spiking communication with *no time penalty* and minimal accuracy reduction

**Neuroscience-constrained algorithms**
- Computation incorporates broad range of neural plasticity and dynamics
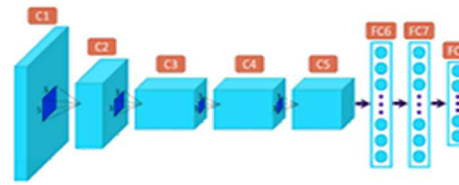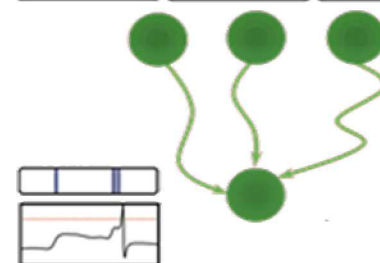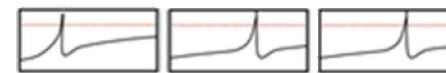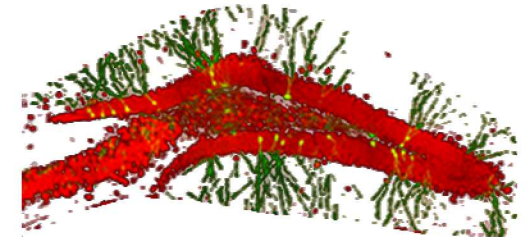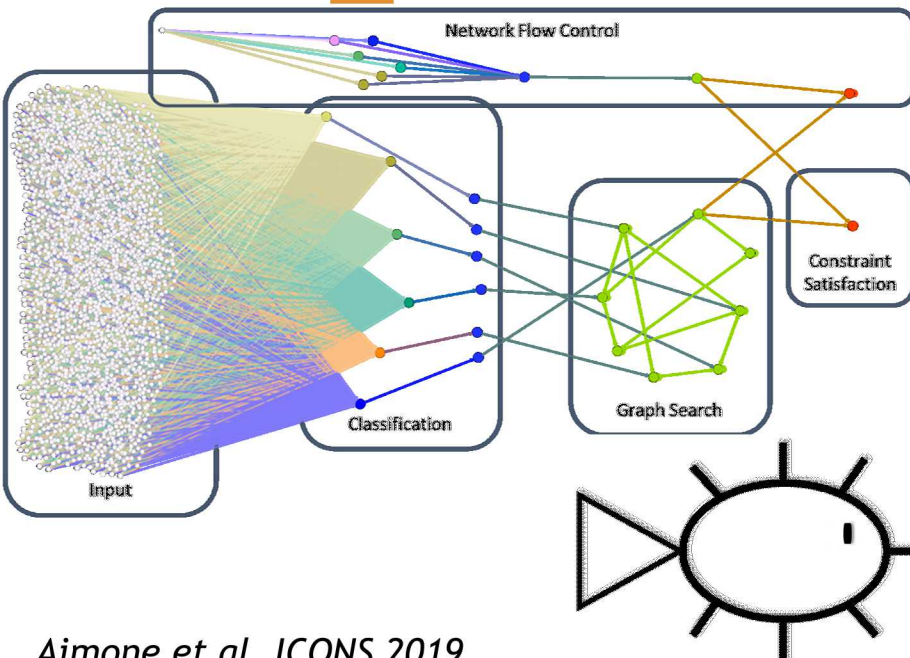- *Generally still unexplored from algorithms perspective*

**Spiking neural algorithms**
- Hand-crafted circuits of spiking neurons
- Model of parallel computation
- Energy efficiency through event-driven communication and high fan-in logic
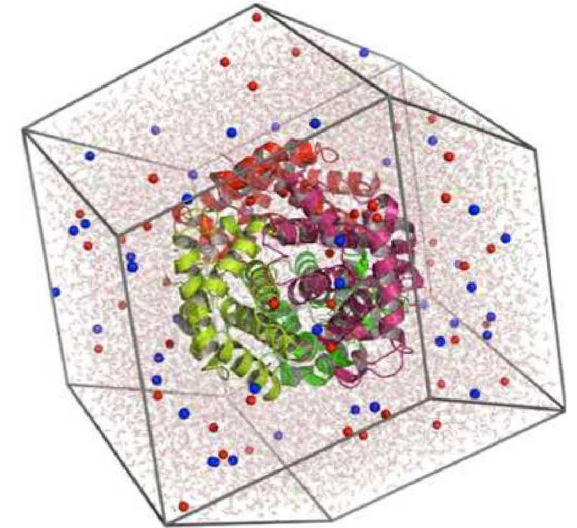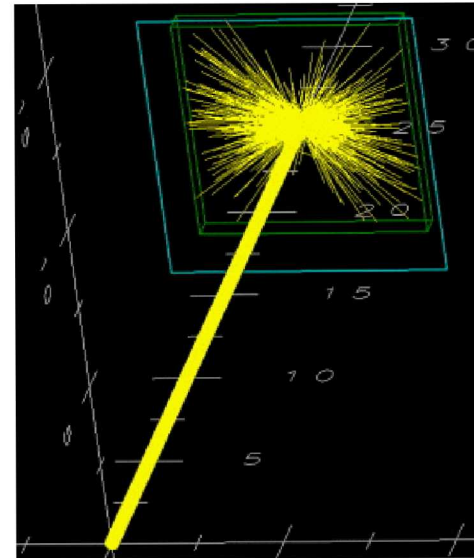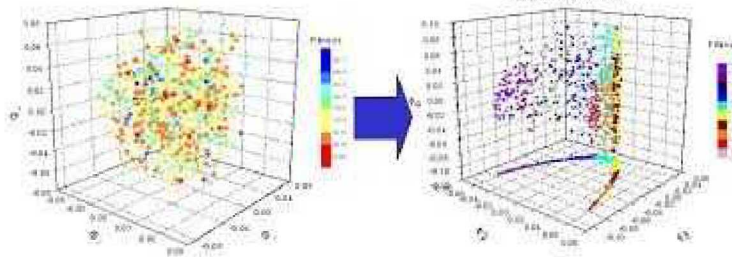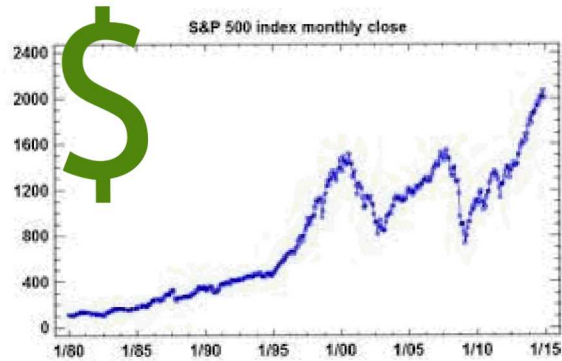
*Aimone et al, ICONS 2019*

Can spiking really be used to solve non-cognitive tasks efficiently?

# Spiking circuits can efficiently solve stochastic differential equations
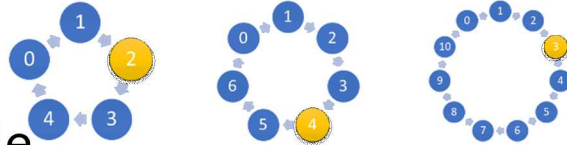
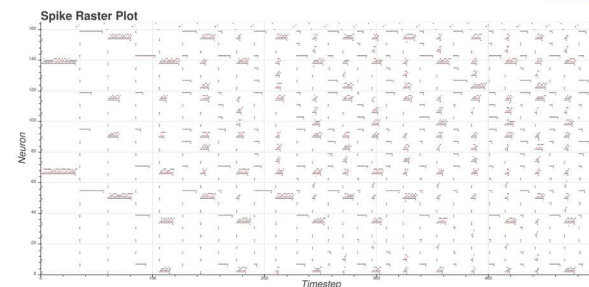Diffusion: $\dfrac{\partial C(x,t)}{\partial t} = D \dfrac{\partial^2 C(x,t)}{\partial x^2}$
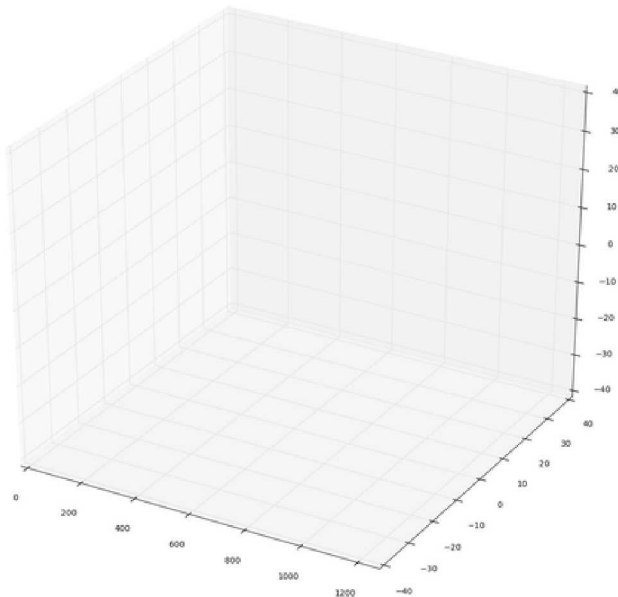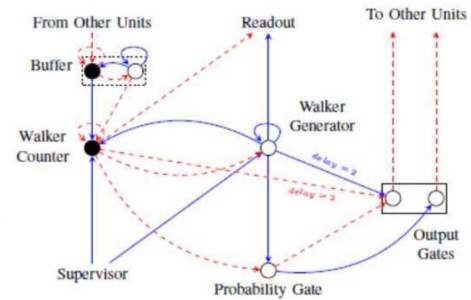
# Spiking circuits can efficiently solve stochastic differential equations

Diffusion:
$$\frac{\partial C(x,t)}{\partial t} = D\frac{\partial^2 C(x,t)}{\partial x^2}$$

Modular circuit of
spiking neurons
per random walk particle

RW counting circuit of
spiking neurons per
simulation mesh vertex

# There is an increasing diversity of spiking algorithms and applications

# A brief plug…



UNIVERSITÄT HEIDELBERG
ZUKUNFT SEIT 1386

## NICE 2020

March 17-20th, 2020

Neuro-Inspired
Computational Elements
Workshop

Abstracts due
November 1st, 2020

Im Neuenheimer Feld 227
D-69120 Heidelberg
Germany

Workshop: March 24-26th 2020
Tutorials: March 27th 2020

Heidelberg - Germany

Picture: fotolia.com / Sergey Borisov

Kirchhoff Institute for Physics