# Swiss Physical Society Annual Meeting 2024

## Neuromorphic Intelligence: spiking neural network and on-line learning circuits for brain-inspired technologies

### Giacomo Indiveri

Institute of Neuroinformatics
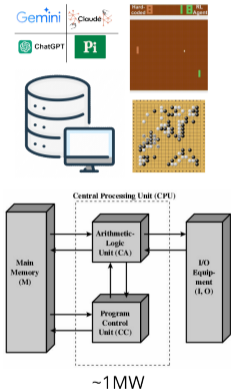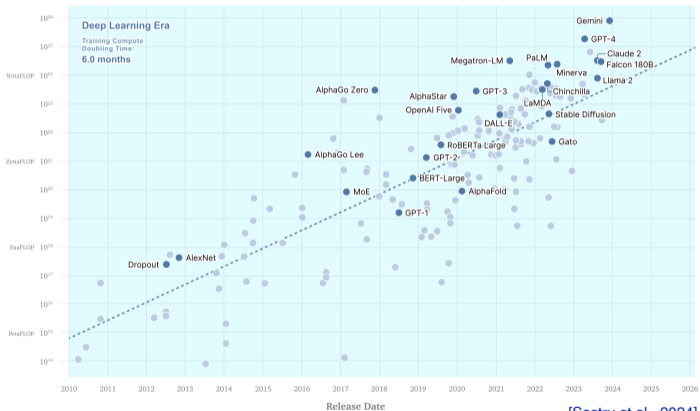University of Zurich and ETH Zurich

**Universität Zürich** [UZH]

**ETH** *zürich*

# 1 Neuromorphic vs Artificial Intelligence

## 2 Building (mixed-signal) neuromorphic systems

## 3 Spike-based learning

## 4 Deploying neuromorphic systems in the real world
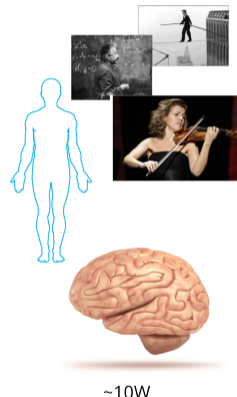
## 5 Conclusions

"Intelligent"



~1MW

**Compute Used for AI Training Runs** (Deep Learning Era)



Deep Learning Era
Training Compute
Doubling Time:
6.0 months

Gemini
GPT-4
Claude 2
Megatron-LM    PaLM    Falcon 180B
Minerva
Llama-2
AlphaGo Zero    AlphaStar    GPT-3    Chinchilla
OpenAI Five    LaMDA
DALL-E    Stable Diffusion
RoBERTa Large    Gato
AlphaGo Lee    GPT-2
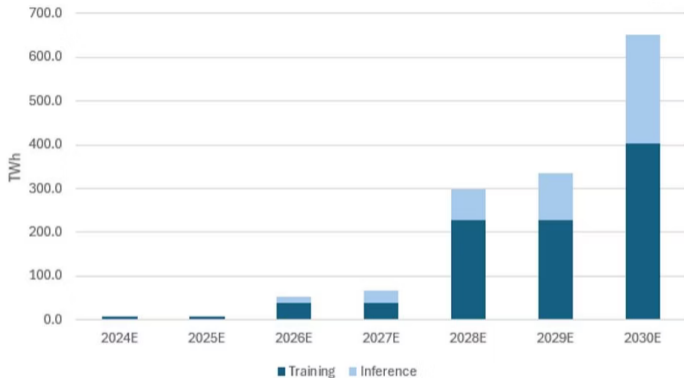BERT-Large
MoE    AlphaFold
GPT-1
Dropout    AlexNet

Release Date

[Sastry et al., 2024]

"Intelligent"



~10W

Generative AI Power Demand, AI Training and Inference

[io-fund.com]

Performance scaling laws

$$L = (C_{min}/2.3 \cdot 10^8)^{-0.050}$$

Compute
Petaflop/s-days*, non-embedding
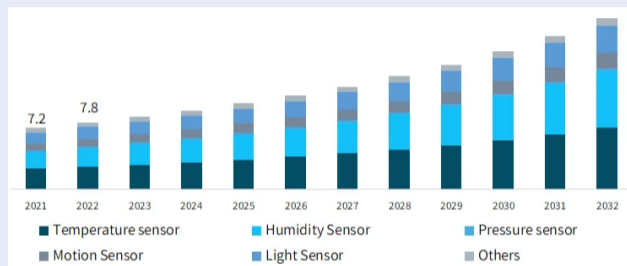
[Sastry et al., 2024]

*AI training is expected to drive the power demand to 402 TWh by 2030 (about the same demand of the whole of France or Germany in 2023)*
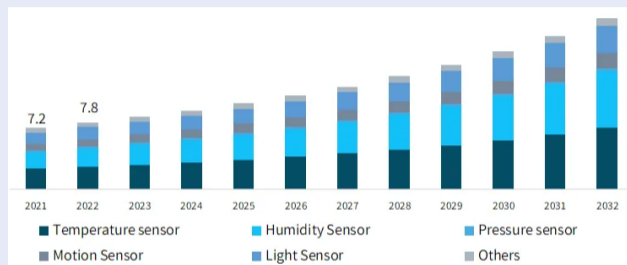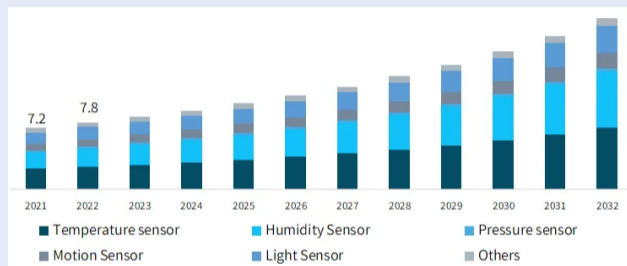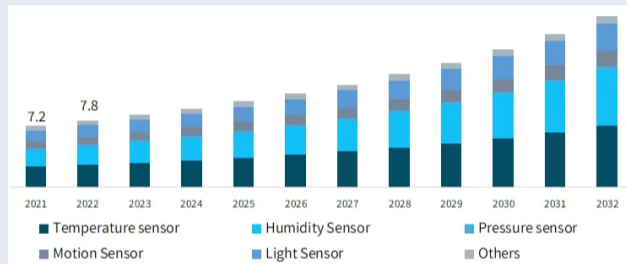
## Global passive sensors market (USD billion)



[gminsights.com]

- More than 50 billion Internet of Things (IoT) devices are expected by 2030
- Embedded devices with sensors and/or actuators are the key components of the IoT
- Local "intelligence" is key to reducing communication, bandwidth and energy consumption.

## Global passive sensors market (USD billion)



7.2  7.8

2021  2022  2023  2024  2025  2026  2027  2028  2029  2030  2031  2032

■ Temperature sensor   ■ Humidity Sensor   ■ Pressure sensor
■ Motion Sensor   ■ Light Sensor   ■ Others

[gminsights.com]

- More than 50 billion Internet of Things (IoT) devices are expected by 2030
- Embedded devices with sensors and/or actuators are the key components of the IoT
- Local "intelligence" is key to reducing communication, bandwidth and energy consumption.

## Global passive sensors market (USD billion)



7.2  7.8

2021 2022 2023 2024 2025 2026 2027 2028 2029 2030 2031 2032

■ Temperature sensor   ■ Humidity Sensor   ■ Pressure sensor
■ Motion Sensor   ■ Light Sensor   ■ Others

[gminsights.com]

- More than 50 billion Internet of Things (IoT) devices are expected by 2030
- Embedded devices with sensors and/or actuators are the key components of the IoT
- Local "intelligence" is key to reducing communication, bandwidth and energy consumption.

## Global passive sensors market (USD billion)



[gminsights.com]

- More than 50 billion Internet of Things (IoT) devices are expected by 2030
- Embedded devices with sensors and/or actuators are the key components of the IoT
- Local "intelligence" is key to reducing communication, bandwidth and energy consumption.

Clearly it is not possible to use conventional large-scale AI methods to endow IoT devices with intelligence.

## Conventional approaches

1. Make application specific (lose general purpose flexibility)
2. Quantize parameters (reduce bit precision)
3. Minimize resource usage (reduce accuracy)

## Novel approaches

1. Reduce data movement (implement in-memory computing)
2. Reduce clock switching (use asynchronous circuits)
3. Exploit all the physics of the devices (mix analog and digital)

# Reducing energy in embedded systems

## Conventional approaches

1. Make application specific (lose general purpose flexibility)
2. Quantize parameters (reduce bit precision)
3. Minimize resource usage (reduce accuracy)

## Novel approaches

1. Reduce data movement (implement in-memory computing)
2. Reduce clock switching (use asynchronous circuits)
3. Exploit all the physics of the devices (mix analog and digital)

# Reducing energy in embedded systems

## Conventional approaches

1. Make application specific (lose general purpose flexibility)
2. Quantize parameters (reduce bit precision)
3. Minimize resource usage (reduce accuracy)

## Novel approaches

1. Reduce data movement (implement in-memory computing)
2. Reduce clock switching (use asynchronous circuits)
3. Exploit all the physics of the devices (mix analog and digital)

⇓

## Novel computing paradigms: brain-like computation

- Co-localize memory and computation (local processing, local state variables)
- Maximize fine-grain parallelism (massively parallel arrays of memory and processing)
- Use the "physics of computation" (exploit properties of computing substrate)

[Indiveri Sandamirskaya, IEEE Signal Processing Magazine, 2019; Indiveri Liu, Proceedings of IEEE, 2015]

Universität Zürich   **ETH** zürich

## Animal brains

- Slow, noisy and variable processing elements.
- Local connectivity, small world networks.
- Massively parallel distributed computation.
- Closed-loop interaction with the environment.
- Real-time spatio-temporal signal processing.
- Continual always-on learning.

## Existence proof



**Bee brain specs**

| | |
|---|---|
| weight: | 1 mg |
| volume: | 1 mm$^3$ |
| # neurons: | 960'000 |
| energy/op: | $10^{-15}$ J/spike |

## Animal brains

- Slow, noisy and variable processing elements.
- Local connectivity, small world networks.
- Massively parallel distributed computation.
- Closed-loop interaction with the environment.
- Real-time spatio-temporal signal processing.
- Continual always-on learning.

## Time represents itself

The brain uses the time evolution of the physical system to implement its computations. Neural circuits compute by exploiting the natural time evolution of their hardware substrate. [Sterling & Laughlin, 2017]

## Clock speed



Brains outperform faster computing systems in many sensory processing tasks at lower speeds, with less power.

## Neuroscience

- Study the principles of computation in animal brains
- Identify them at the neural circuit level
- Emulate the bio-physics of neurons and synapses using analog electronic circuits
- Validate/invalidate hypotheses of neural computation

## Electronics

- Include novel devices and emerging memory technologies
- Exploit (all) the physics of these nanoscale devices
- Integrate CMOS and memristive devices together
- Engineer efficient "in-memory computing" architectures

## Computing

- Implement "neural processing" systems in custom ASICS
- Integrate processors with sensors and actuators
- Apply them to closed-loop sensory processing tasks
- Develop *cognitive agents* that produce autonomous behavior.

## Neuroscience

- Study the principles of computation in animal brains
- Identify them at the neural circuit level
- Emulate the bio-physics of neurons and synapses using analog electronic circuits
- Validate/invalidate hypotheses of neural computation

## Device physics

- Include novel devices and emerging memory technologies
- Exploit (all) the physics of these nanoscale devices
- Integrate CMOS and memristive devices together
- Engineer efficient "in-memory computing" architectures

## Computing

- Implement "neural processing" systems in custom ASICS
- Integrate processors with sensors and actuators
- Apply them to closed-loop sensory processing tasks
- Develop *cognitive agents* that produce autonomous behavior.

1 Neuromorphic vs Artificial Intelligence

2 Building (mixed-signal) neuromorphic systems

3 Spike-based learning

4 Deploying neuromorphic systems in the real world

5 Conclusions

- Spiking neural networks (SNNs)
- Analog subthreshold circuits.
- Slow temporal, non-linear dynamics.
- Massively parallel operation.
- Compatible with memristive devices
- Inhomogeneous, imprecise, and noisy.

- Spiking neural networks (SNNs)
- Analog subthreshold circuits.
- Slow temporal, non-linear dynamics.
- Massively parallel operation.
- Compatible with memristive devices
- Inhomogeneous, imprecise, and noisy.
- Fast asynchronous digital circuits for routing spikes.
- Reprogrammable network topology

Adaptive Exponential I&F neuron model (beyond HH)

| Work | [12] | [19] | [36] | **This work** |
|---|---|---|---|---|
| Techn. | 180 nm | 28 nm | 28 nm | 22 nm |
| | CMOS | CMOS | FDSOI | FDSOI |
| Type | Mixed | Mixed | Mixed | Mixed |
| $V_{dd}$ | 1.8 V | 0.7-1 V | 1 V | 0.8 V |
| Freq | 30 Hz | - | 30 Hz | 30 Hz |
| Results | Experimental | Experimental | Simulation | Simulation |
| En./spike | 883 pJ | 2.3 nJ-30 nJ | 50 pJ | 14 pJ |

[Brette and Gerstner, 2005, Rubino et al., IEEE TCAS, 2020]

## Device mismatch effects



Time-to-first spike measured across 256 different neurons

## Device mismatch effects



Time-to-first spike measured across 256 different neurons

## How to cope with mismatch?

- Use populations of neurons and average over space and time
- Employ negative feedback, adaptation, and learning mechanisms

## Device mismatch effects



Time-to-first spike measured across 256 different neurons

## How to cope with mismatch?

- Use populations of neurons and average over space and time
- Employ negative feedback, adaptation, and learning mechanisms

## Population codes and averaging

Universität Zürich™ **ETH** zürich

## Device mismatch effects



Time-to-first spike measured across 256 different neurons

## How to cope with mismatch?

- Use populations of neurons and average over space and time
- Employ negative feedback, adaptation, and learning mechanisms

## Choosing bit resolution



| integration time (s) | 2 | 4 | 8 | 16 | 32 | 64 | 128 | |
|---|---|---|---|---|---|---|---|---|
| 5 | 12.7 | 9.0 | 6.4 | 4.2 | 2.7 | 1.7 | 0.9 | -14 bits |
| 2 | 12.7 | 9.0 | 6.4 | 4.3 | 2.7 | 1.8 | 0.9 | -12 bits |
| 1 | 12.7 | 9.1 | 6.4 | 4.3 | 2.7 | 1.8 | 0.9 | |
| 0.5 | 12.8 | 9.1 | 6.5 | 4.3 | 2.8 | 1.8 | 0.9 | -10 bits |
| 0.2 | 12.6 | 9.0 | 6.4 | 4.2 | 2.7 | 1.7 | 1.0 | -8 bits |
| 0.1 | 13.9 | 10.1 | 7.1 | 4.8 | 3.0 | 2.0 | 0.9 | -6 bits |
| 0.05 | 16.1 | 11.0 | 8.2 | 5.5 | 3.4 | 2.1 | 1.2 | |
| 0.02 | 18.3 | 12.6 | 8.7 | 5.9 | 4.2 | 2.9 | 1.2 | -4 bits |

cluster size

Coefficient of variation and Equivalent Number of Bits (ENOB)

[Zendrikov et al., 2023]

1 Neuromorphic vs Artificial Intelligence

2 Building (mixed-signal) neuromorphic systems

3 Spike-based learning

4 Deploying neuromorphic systems in the real world

5 Conclusions

*Ensemble learning techniques exploit variability of inhomogeneous synapses.*



| | | |
|---|---|---|
| MNIST | deep/CNN (Hinton et al. 2012) | 98.4% |
| | random + bistable synapses | ~ 85% |
| | random + bistable synapses + (mod. protocol) | ~ 96% |
| TIMIT | deep/CNN (Hinton et al. 2012) | 77% |
| | VLSI cochlea + bistable synapses | ~ 60% |

**On-line bagging techniques** *require* **variability**

**AdaBoost theorem:** $1 - \text{error}(H_{final}) \geq 1 - e^{-2\gamma^2 N}$

[Freund and Schapire, 1997]

# Biologically plausible learning rule

## Features

- Compatible with "biological" and "electronic" computing substrate
- Based on latest dendritic multi-compartment models
- Exploits properties of multiple inhibitory cell types
- Makes use of population coding and dynamics
- Implements a "stop learning" mechanism to automatically switch between training and inference

## Cortical motif



Teacher/prediction — Context/attention — VIP — SST — PV — PYR — Plastic weights — Input

## Design team (FDSOI 22 nm, 2024)



The big boss — Modelling and software — Analog hardware — Digital hardware

Spiking inputs arrive from both top-down pathways (context, attention, prediction signals) and bottom-up pathways (sensory signals).

- An E-I balance maintains the population in a proper operating range at all times.
- During training, teacher signals reach the soma and change plastic weights
- During inference, SST inhibitory cells block top-down inputs, the neurons respond only to bottom-up inputs with lower firing rates, and synaptic weights "stop learning".

Spiking inputs arrive from both top-down pathways (context, attention, prediction signals) and bottom-up pathways (sensory signals).

- An E-I balance maintains the population in a proper operating range at all times.
- During training, teacher signals reach the soma and change plastic weights
- During inference, SST inhibitory cells block top-down inputs, the neurons respond only to bottom-up inputs with lower firing rates, and synaptic weights "stop learning".

Spiking inputs arrive from both top-down pathways (context, attention, prediction signals) and bottom-up pathways (sensory signals).

- An E-I balance maintains the population in a proper operating range at all times.
- During training, teacher signals reach the soma and change plastic weights
- During inference, SST inhibitory cells block top-down inputs, the neurons respond only to bottom-up inputs with lower firing rates, and synaptic weights "stop learning".

$$Y_{PYR} = \sigma\left(I_a + I_b\right)$$

$$\theta = \int Y_{PYR} dt$$

$$\Delta w_{up} = \eta(I_a - I_b)\sigma_{LTP} \quad \text{if} \quad I_a \geq I_b$$

$$\Delta w_{dn} = \eta(I_a - I_b)\sigma_{LTD} \quad \text{if} \quad I_a < I_b$$

$$
\begin{aligned}
I_a &= [I_T - I_{SST}]^+ \\
I_b &= (I_{in} + I_{PYR} - I_{PV}) \\
I_{in} &= \int \sum_i w_i \alpha\left(t - \delta(t_i)\right) dt
\end{aligned}
$$

$$\sigma_{LTP} = \begin{cases} 1 & \text{if} \quad \theta_{LTP_-} < \theta < \theta_{LTP_+} \\ 0 & \text{otherwise} \end{cases}$$

$$\sigma_{LTD} = \begin{cases} 1 & \text{if} \quad \theta_{LTD_-} < \theta < \theta_{LTD_+} \\ 0 & \text{otherwise} \end{cases}$$

SPICE circuit simulations of two neurons

NEST SW simulations at full population level (with mismatch)

# Robust spike-based learning mechanisms

There are many hardware-friendly spike-driven learning algorithms that (go beyond STDP).

W. Senn, S. Fusi, N. Brunel, S. Sheik, E. Neftci, R. Zecchina, M. Memmesheimer, etc.

# Robust spike-based learning mechanisms

There are many hardware-friendly spike-driven learning algorithms that (go beyond STDP).

W. Senn, S. Fusi, N. Brunel, S. Sheik, E. Neftci, R. Zecchina, M. Memmesheimer, etc.

All rule have the following *requirements*:

- Redundancy (population codes)
- Bi-stable or multi-stable weights
- Variability and heterogeneity
- Analog, continuous-time state variables



Calcium variable C(t)

Postsynaptic membrane potential $V_{mem}(t)$

Presynaptic spikes

Synaptic state X(t)

# Robust spike-based learning mechanisms

There are many hardware-friendly spike-driven learning algorithms that (go beyond STDP).

W. Senn, S. Fusi, N. Brunel, S. Sheik, E. Neftci, R. Zecchina, M. Memmesheimer, etc.



Calcium variable C(t)

Postsynaptic membrane potential $V_{mem}(t)$

Presynaptic spikes

Synaptic state X(t)

Time (s)

All rule have the following *requirements*:

- Redundancy (population codes)
- Bi-stable or multi-stable weights
- Variability and heterogeneity
- Analog, continuous-time state variables

Many mixed-signal hardware implementations have been demonstrated:

- Supervised learning, mean rates
- Unsupervised learning, precise spike-timing
- Hopfield/attractor networks

- Reservoir computing, liquid state and perceptron
- Ensemble learning (random forest, bagging)

[Khacef et al., 2023]

## Edge intelligence

Mixed-signal neuromorphic systems are optimally suited for extreme-edge computing applications, which require resource constrained electronic systems. They are ideal for always-on in-sensor and in-memory computing applications that need to perform closed-loop interactions the environment, in real-time.



Example: wearables and health monitoring

- Neuromorphic CPG for adaptive pace-makers
  [Abu-Hassan et al., 2019]
- ECG anomaly detection [Bauer et al., 2019,Corradi et al., 2019]
- EMG signal classification [Donati et al., 2019,Ma et al, 2020]
- High-Frequency Oscillation (HFO) detection
  [Sharifhazileh et al., 2021,Burelo et al., 2022]
- Neuromorphic Heart Rate Monitors [Carpegna et al., 2024]

1. Neuromorphic vs Artificial Intelligence

2. Building (mixed-signal) neuromorphic systems

3. Spike-based learning

4. Deploying neuromorphic systems in the real world

5. Conclusions

# Big data needs a hardware revolution

*Artificial intelligence is driving the next wave of innovations in the semiconductor industry.*

Software companies make headlines but research on computer hardware could bring bigger rewards. Credit: Morris MacMatzen/Getty

- Conventional AI increasing power requirements are unsustainable.
- New emerging memory technologies will benefit from massively parallel processing architectures.
- Neuroscience and machine learning are uncovering powerful and robust neural processing methods.
- Hardware implementations of spiking neural networks and sparse event-based sensory-processing systems are starting to show their advantages.
- This is the perfect time to follow the "neuromorphic intelligence" approach for starting a hardware revolution.

# Conclusions

## Big data needs a hardware revolution

*Artificial intelligence is driving the next wave of innovations in the semiconductor industry.*

Software companies make headlines but research on computer hardware could bring bigger rewards. Credit: Morris MacMatzen/Getty

- Conventional AI increasing power requirements are unsustainable.
- New emerging memory technologies will benefit from massively parallel processing architectures.
- Neuroscience and machine learning are uncovering powerful and robust neural processing methods.
- Hardware implementations of spiking neural networks and sparse event-based sensory-processing systems are starting to show their advantages.
- This is the perfect time to follow the "neuromorphic intelligence" approach for starting a hardware revolution.

**EDITORIAL** · 06 FEBRUARY 2018

## Big data needs a hardware revolution

*Artificial intelligence is driving the next wave of innovations in the semiconductor industry.*

Software companies make headlines but research on computer hardware could bring bigger rewards. Credit: Morris MacMatzen/Getty

- Conventional AI increasing power requirements are unsustainable.
- New emerging memory technologies will benefit from massively parallel processing architectures.
- Neuroscience and machine learning are uncovering powerful and robust neural processing methods.
- Hardware implementations of spiking neural networks and sparse event-based sensory-processing systems are starting to show their advantages.
- This is the perfect time to follow the "neuromorphic intelligence" approach for starting a hardware revolution.

**nature**

**EDITORIAL · 06 FEBRUARY 2018**

# Big data needs a hardware revolution

*Artificial intelligence is driving the next wave of innovations in the semiconductor industry.*

Software companies make headlines but research on computer hardware could bring bigger rewards. Credit: Morris MacMatzen/Getty

- Conventional AI increasing power requirements are unsustainable.
- New emerging memory technologies will benefit from massively parallel processing architectures.
- Neuroscience and machine learning are uncovering powerful and robust neural processing methods.
- Hardware implementations of spiking neural networks and sparse event-based sensory-processing systems are starting to show their advantages.
- This is the perfect time to follow the "neuromorphic intelligence" approach for starting a hardware revolution.

**nature**
International journal of science

EDITORIAL · 06 FEBRUARY 2018

## Big data needs a hardware revolution

*Artificial intelligence is driving the next wave of innovations in the semiconductor industry.*

Software companies make headlines but research on computer hardware could bring bigger rewards. Credit: Morris MacMatzen/Getty

- Conventional AI increasing power requirements are unsustainable.
- New emerging memory technologies will benefit from massively parallel processing architectures.
- Neuroscience and machine learning are uncovering powerful and robust neural processing methods.
- Hardware implementations of spiking neural networks and sparse event-based sensory-processing systems are starting to show their advantages.
- This is the perfect time to follow the "neuromorphic intelligence" approach for starting a hardware revolution.

## institute of neuroinformatics

Universität Zürich

**ETH** *zürich*



- Elisa Donati
- Chiara De Luca
- Sapta Ghosh
- Chenxi Wen
- Dmitrii Zendrikov

- Farah Baracat
- Junren Chen
- Yigit Demirag
- Maryada

- Shyam Narayanan
- Arianna Rubino
- Zhe Su

erc

ICT H2020

FNSNF SWISS NATIONAL SCIENCE FOUNDATION

# Thank you for your attention

# Backup slides

# NEUROMORPHIC
Computing and Engineering

OPEN ACCESS

No APCs in 2024

NEUROMORPHIC
Computing and Engineering

iopscience.org/nce

IOP Publishing

A multidisciplinary, open access journal devoted to the application and development of neuromorphic computing, devices, and systems in advancing new scientific discovery and realising emerging new technologies.

**Editor-in-Chief**
**Giacomo Indiveri** University of Zurich, Switzerland

Indexed in Scopus and Web of Science
IMPACT FACTOR COMING IN JUNE **2024**

**IOP** Publishing

Iopscience.org/nce
nce@ioppublishing.org

@IOPneuromorphic

In addition to using populations of neurons and use learning and plasticity to improve robustness of neural processing, it is useful to identify and adopt basic building blocks that implement key principles of computation.

# Neural computational primitives

In addition to using populations of neurons and use learning and plasticity to improve robustness of neural processing, it is useful to identify and adopt basic building blocks that implement key principles of computation.

- Attractor networks
- E-I balanced networks
- Winner-Take-All networks
- Relational networks
- Coupled oscillators
- Neural State Machines

In addition to using populations of neurons and use learning and plasticity to improve robustness of neural processing, it is useful to identify and adopt basic building blocks that implement key principles of computation.

- Attractor networks
- E-I balanced networks
- Winner-Take-All networks
- Relational networks
- Coupled oscillators
- Neural State Machines

  see also poster by Maryada et al.
  (Calcium-based dendritic plasticity)



electrons



A OR Q · A AND Q · A XOR Q
A NOR Q · A NAND Q · A NOT Q̄



500 μm



Cortical

Thalamus · Subcortical

## Edge intelligence

Mixed-signal neuromorphic systems are optimally suited for extreme-edge computing applications, which require resource constrained electronic systems. They are ideal for always-on in-sensor and in-memory computing applications that need to closed-loop interactions the environment, in real-time.

## Edge intelligence

Mixed-signal neuromorphic systems are optimally suited for extreme-edge computing applications, which require resource constrained electronic systems. They are ideal for always-on in-sensor and in-memory computing applications that need to closed-loop interactions the environment, in real-time.



Example: wearables and health monitoring

- Neuromorphic CPG for adaptive pace-makers
  [Abu-Hassan et al., 2019]

- ECG anomaly detection [Bauer et al., 2019,Corradi et al., 2019]

- EMG signal classification [Donati et al., 2019,Ma et al 2020]

- High-Frequency Oscillation (HFO) detection
  [Sharifhazileh et al., 2021,Burelo et al., 2022]

- Neuromorphic Heart Rate Monitors [Carpegna et al., 2024]

Detecting "agitation states" by monitoring monotonic increases in heart rates, over long-time periods.

E  Exc. Population
I  Inh. Population
→  Exc. Coupling
●  Inh. Coupling

SLOW INH SYN (GABA B)
FAST EXC SYN (AMPA)
SLOW EXC SYN (NMDA)
AdExp LIF NEURON
AdExp LIF POPULATION

FILTER BANK · GATING E-I BALANCED · WTA
STATE 3 · ALARM
STATE 2
STATE 1
STATE 0 · RELAXED

# Neuromorphic vs conventional processors

## Pros

- Low-power ($< 1\,\mathrm{mW}$)
- Low latency
- . . .

## Cons

- High area
- High variability, noisy
- Low(er) accuracy

## Pros

- Low-power (< 1 mW)
- Low latency
- . . .

## Cons

- High area
- High variability, noisy
- Low(er) accuracy

## What are they good for?

- Closed-loop sensory-motor processing
- Multi-modal sensory fusion
- Always-on on-line learning

## What are they bad at?

- High precision number crunching
- High accuracy pattern recognition
- Batch processing of large data sets

## Pros

- Low-power (< 1 mW)
- Low latency
- …

## Cons

- High area
- High variability, noisy
- Low(er) accuracy

## What are they good for?

- Closed-loop sensory-motor processing
- Multi-modal sensory fusion
- Always-on on-line learning

## What are they bad at?

- High precision number crunching
- High accuracy pattern recognition
- Batch processing of large data sets

## Open challenges

- How to obtain robust and reliable computation using a noisy and heterogeneous computing substrate.
- How to program networks of spiking neurons (hint: compose *computational primitives* and use *learning*).

## Background

We are building physical, real-time, signal processing systems for real-world sensory data.

### Background

We are building physical, real-time, signal processing systems for real-world sensory data.

### Requirements

1. Robust communication of analog signals across long distances through noisy channels.
2. Local processing, multi-core architectures and distributed computing.
3. Low power and low-latency.

## Background

We are building physical, real-time, signal processing systems for real-world sensory data.

## Requirements

1. Robust communication of analog signals across long distances through noisy channels.
2. Local processing, multi-core architectures and distributed computing.
3. Low power and low-latency.

## Optimal solution for communication and computation

- The optimal method that minimizes bandwidth and power consumption for achieving this goal, under these constraints, is pulse-frequency modulation. [A. Mortara et al., 1995, K. Boahen, 1998]
- "*Counter to intuition, computing with spikes can be extremely efficient on neuromorphic hardware even when the problem being solved is mathematically formulated in terms of activity rates.*" [M. Davies, Intel, 2019]

## Advantages

- Compute through the dynamics of the circuits
- No need to "count time": avoid use of clock trees
- Avoid large DAC/ADC overhead
- Exploit the full potential of memristors
  - ▸ Exploit intrinsic non-linearities [Brivio et al., 2021]
  - ▸ Exploit intrinsic stochasticity [Gaba et al., 2013, Payvand et al., 2018]

## Disadvantages (?)

- Noisy $\implies$ average across multiple neurons (exploit population coding and heterogeneity)
- Large area requirements $\implies$ employ memristive devices and 3D VLSI (exploit low power)

PCM cross-bar array [Source: IBM]

Device layer

Reconfigurable layer

Device layer

"Dendrocentric learning for synthetic intelligence" [Boahen, 2022]

$$\tau \frac{d}{dt} I_{syn} + I_{syn} = \frac{I_g I_w}{I_\tau}$$

$$\tau = \frac{C U_T}{\kappa I_\tau}$$

[Bartolozzi, Indiveri, NECO 2007]

## Also real neurons are diverse and inhomogeneous

## Balanced HW E-I networks



Maryada

# Reliable dynamics with cross-homeostatic plasticity

## Balanced HW E-I networks



## Cross-homeostatic plasticity induced stability



[Maryada et al., 2023]

# Spike-based backprop approximations

## Local learning rules

- Fusi et al. 2000
- Brader et al. 2007
- Urbanczik, Senn 2014
- Baldassi et al. 2016
- Neftci et al. 2017
- Sacramento et al. 2018
- Bellec et al. 2019
- Zenke, Vogels 2021
- Siddique et al. 2023
- . . .



[Cartiglia et al., AICAS 2019]

*Ensemble and stochastic learning can exploit variability of inhomogeneous synapses.*



| | | |
|---|---|---|
| MNIST | deep/CNN (Hinton et al. 2012) | 98.4% |
| | random + bistable synapses | ~ 85% |
| | random + bistable synapses + (mod. protocol) | ~ 96% |
| TIMIT | deep/CNN (Hinton et al. 2012) | 77% |
| | VLSI cochlea + bistable synapses | ~ 60% |

## On-line bagging techniques

AdaBoost theorem: $1 - \text{error}(H_{final}) \geq 1 - e^{-2\gamma^2 N}$

[Y. Freund And R. E. Schapire, 1995]

- Zebra-finch "Bird's Own Song" classification [Corradi et al., 2015]
- Closed-loop bidirectional brain machine interfaces with in rats and cell-cultures [Boi et al., 2016] [Serb et al. 2020]
- Adaptive pace-maker with neuromorphic CPG network [Abu-Hassan et al., 2019]
- On-line ECG anomaly detection [Bauer et al., 2019]
- On-line classification of EMG signals [Donati et al., 2019]
- Closed-loop obstacle avoidance on roving robot [Milde et al. 2017]
- Closed-loop robot head position control with a neuromorphic processor [Zhao et al., 2020]
- Neuromorphic pattern generation circuits for bioelectronic medicine [Donati et al., 2021]
- Instantaneous stereo depth estimation of real-world stimuli with a stereo-vision setup [Risi et al., 2021]
- On-line detection of vibration anomalies using balanced spiking neural networks [Dennler et al., 2021]
- High-Frequency Oscillation (HFO) detection [Sharifhazileh, Burelo et al., 2021]

## What is an HFO?

Spontaneous EEG events in the frequency range between 80 and 500 Hz consisting of at least four oscillations that clearly stand out from the baseline.



[Fedele et al. 2017]

Universität Zürich

**ETH** *zürich*

HFO are biomarkers for epileptogenic brain tissue.



[Fedele et al., 2017]

**Accuracy: 78 % (vs. 67%)**
**Power consumption: 614 $\mu$W**

[Sharifhazileh, Burelo, et al., 2021]

- Stand-alone sensor+processor
- 256 neurons, 512 synapses
- "Backpropless" two layer network
- One-bit weights
- Inhomogeneous parameters
- Matched time constants
- Power consumption: **614** $\mu$**W**
- Accuracy: **78%** (vs. 67% from s-o-a)

[Sharifhazileh, Burelo et al., Nat. Comms 2021]

## Edge intelligence

We are now entering the era of *neuromorphic intelligence* in which dedicated cognitive "chiplets" will be used to provide intelligence to a multitude of <span style="color:red">extreme edge-computing</span> devices



- Health monitoring
- Wearable sensors
- Environmental sensing
- Industrial monitoring
- Intelligent machine vision
- Consumer applications

**http://capocaccia.cc/**

- Interdisciplinary, international, diverse
- Morning lectures, afternoon hands-on work-groups
- Active and lively discussions (no powerpoint)
- Concrete results, establishment of long-term collaborations

Capo Caccia, Sardinia, Italy. April 28 - May 11, 2024

## Academic/basic research

- Study real brains, start from small neural circuits/systems
- Take into account all properties of neurons and synapses
- Focus on fundamental problems (ignore incremental benchmarks)
- "There's plenty of room at the bottom" (large scale is not all)

## Applied/industrial research

- Choose a specific problem to solve that is not being solved yet
- Consider it's requirements in it's entirety, from end to end
- Be open to using the best of all possible approaches (analog and digital)
- Build the full ecosystem for your solution (devices, software, users)

Early access:

Bottom-Up and Top-Down Approaches for the Design of Neuromorphic Processing Systems: Tradeoffs and Synergies Between Natural and Artificial Intelligence, Frenkel and Indiveri, Proceedings IEEE, 2023.

Universität Zürich

**ETH**zürich

A (Gilbert) normalizer memristive synapse circuit



$$Iw_{pos} = I_b \frac{I_{M1}}{(I_{M1}+I_{M4})}$$

$$Iw_{neg} = I_b \frac{I_{M4}}{(I_{M1}+I_{M4})}$$

[M. Nair et al., Nano Futures, 2017; Payvand et al., Faraday Discuss., 2019]

A (Gilbert) normalizer memristive synapse circuit



Divisive non-linearity "squashes" distributions and <u>reduces mismatch effects</u>



$$Iw_{pos} = I_b \frac{I_{M1}}{(I_{M1}+I_{M4})}$$

$$Iw_{neg} = I_b \frac{I_{M4}}{(I_{M1}+I_{M4})}$$

$CV = 0.429$

$CV = 0.284$

[M. Nair et al., Nano Futures, 2017; Payvand et al., Faraday Discuss., 2019]

DVS

3x3 cxQuad PCB

ROLLS

128X128

Input

32X32

Pooling

4@8X8

Convolution

4@16X16

Pooling

4@8X8

On-line Learning

8@32X1

[Indiveri et al., IEDM 2015]

Experimental setup

← Retinotopic and orientation maps representing the preference of neurons in the visual cortex for the location and orientation of a stimulus on the visual field.

← Retinotopic and orientation maps representing the preference of neurons in the visual cortex for the location and orientation of a stimulus on the visual field.

Orientation tuning: ⟶ Non-human primate response to moving bars (top); Neuromorphic processor response to flashing bars (bottom)

Feature tuning via populations of neurons

[Dayan & Abbott, 2005]

[Chicca et al., 2007]

[R. Krause et al., 2021]

[E. Donati et al., 2021]

Université Zürich ETH zürich

**Synaptic matrix**



Inhibitory neurons

Excitatory neurons

Excitatory neurons

Excitatory neurons

**Output spikes**



Freq (Hz)

Time (s)

Input

Input



Methods and tools:

- Mean Field Theory
- Effective Response Function
- Self consistency condition

[M. Giulioni et al., 2012, M. Giulioni et al., 2015]

$$\tau \frac{d}{dt}\mu = -\mu + W\nu_{in} - \beta$$

$$\nu_{out} = \Phi(\mu(\nu_{in}), \sigma)$$

$$\nu_{out} = \nu_{in} = \nu$$



Effective Response Function

mean rate 100 Hz!

[Neftci et al., 2013]

[Zendrikov et al., 2023]

Universität Zürich<sup>UZH</sup>

**ETH**_zürich_

Binary variables V0,V1: V0≠V1



In absence of external input (evidence), the network settles to the lowest energy state (all constraints satisfied).

[Mostafa et al., 2015]

Can be applied to all Boolean satisfiability problems, such as graph coloring problem, SAT, etc.

[Mostafa et al., 2015]

Exploiting the device mismatch in the neuron's refractory period.



[Bias et al., 2016]

[J. Zhao et al., 2020]

16 neurons

16 neurons

31 neurons

### Finite State Machines vs Neural State Machines

A finite-state machine (FSM) is a mathematical model of computation used to design both computer programs and sequential logic circuits. It is conceived as an abstract machine that can be in one of a finite number of states. [Wikipedia]



- Recognizes regular expression $B^*[AB^*A]^*$



▷ Initial

◯ Accept

[Minsky, 1967]

Universität Zürich

**ETH** *zürich*

## Single WTA

Inhibitory neurons    Excitatory neurons



Global Inhibition    Nearest-N Excitation



## Coupled WTAs



E — Exc. Population
I — Inh. Population
→ Exc. Coupling
•— Inh. Coupling

[R. Rutishauser & R.J. Douglas, 2009, R. Rutishauser et al., 2011, E. Neftci et al., 2013]

## Single WTA

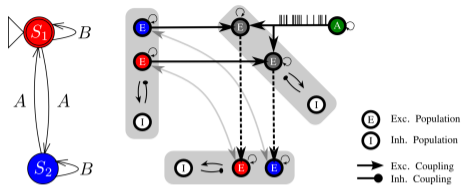Inhibitory neurons    Excitatory neurons



Global Inhibition    Nearest-N Excitation



## Coupled WTAs



E — Exc. Population
I — Inh. Population
→ Exc. Coupling
●— Inh. Coupling

[R. Rutishauser & R.J. Douglas, 2009, R. Rutishauser et al., 2011, E. Neftci et al., 2013]

## Single WTA

Inhibitory neurons    Excitatory neurons

Global Inhibition    Nearest-N Excitation

## Coupled WTAs

E Exc. Population
I Inh. Population
→ Exc. Coupling
•→ Inh. Coupling

[R. Rutishauser & R.J. Douglas, 2009, R. Rutishauser et al., 2011, E. Neftci et al., 2013]

(a) State Machine

(b) Network Architecture

(c) Example Run

Chip Output

[E. Neftci et al., 2013]

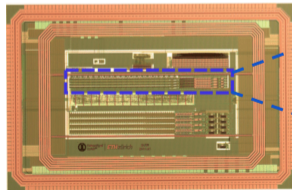# Synthesizing neuromorphic cognitive systems
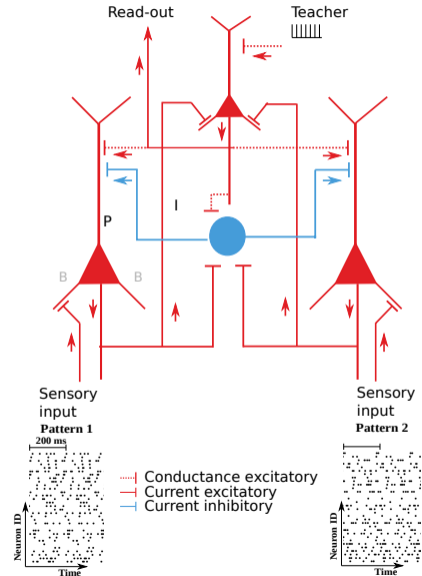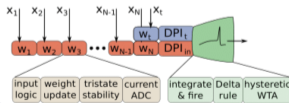
Università Zürich · ETH zürich

- Bi-stable synapses with STDP circuits

  [Indiveri et al, 2006]

- Spike-driven synaptic plasticity with stop-learning

  [Mitra et al, 2009,Qiao et al., 2015]

- Error-propagation with local learning

  [Cartiglia et al., 2020]

- Dendritic Hebbian synaptic plasticity with stop-learning

  [Rubino et al, 2023]
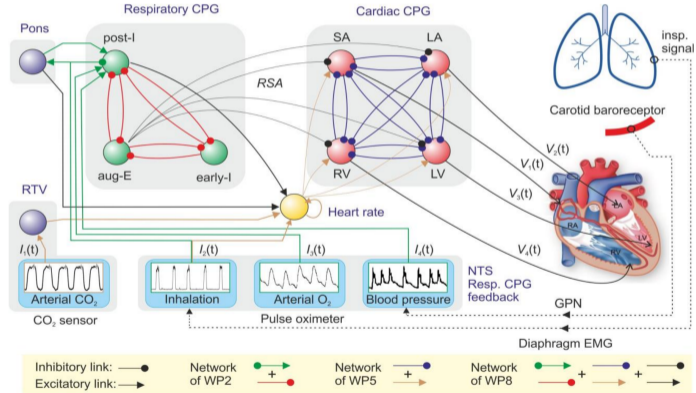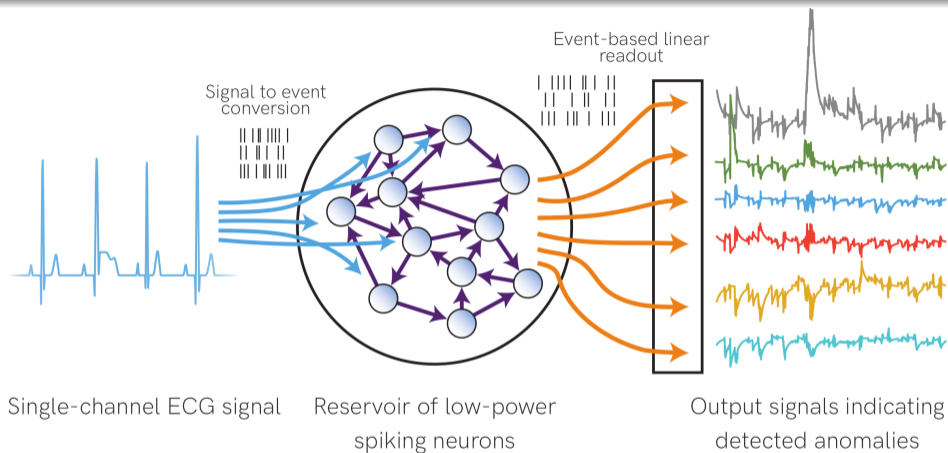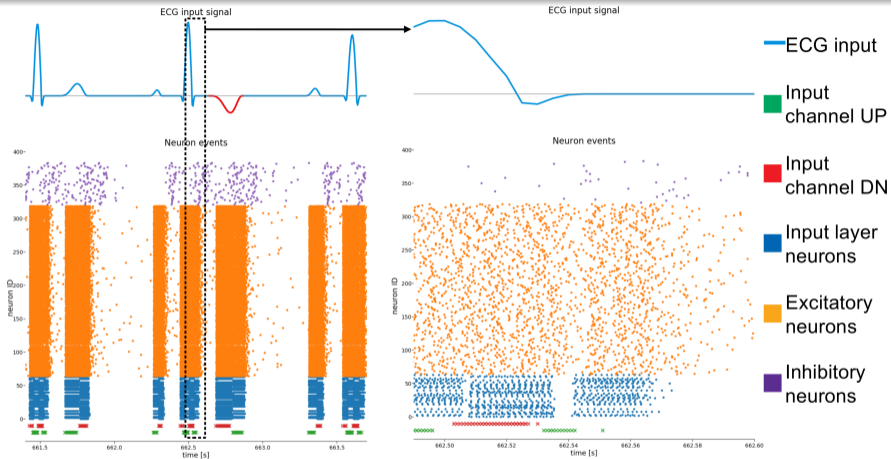
Build an adaptive pacemaker that responds to physiological feedback in real time to recover heart rate adaptation functionality.
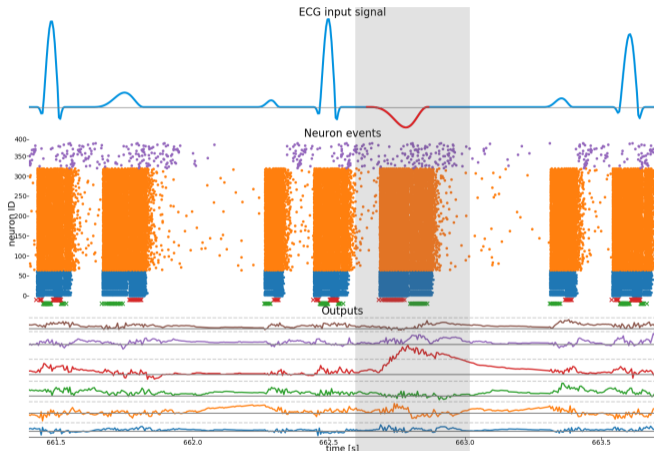


(Elisa Donati, Renate Krause)

Signal to event conversion

Event-based linear readout

Single-channel ECG signal

Reservoir of low-power spiking neurons

Output signals indicating detected anomalies

[H. Jaeger, 2003] [W. Maass et al., 2002] [F. Bauer and D. Muir, SynSense]

Universität Zürich

**ETH** zürich



ECG input

Input channel UP

Input channel DN

Input layer neurons

Excitatory neurons

Inhibitory neurons

[H. Jaeger, 2003] [W. Maass et al., 2002] [F. Bauer and D. Muir, SynSense]

- Generic, single-led ECG
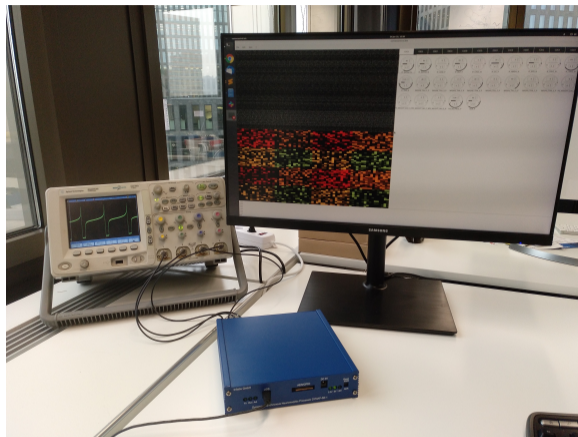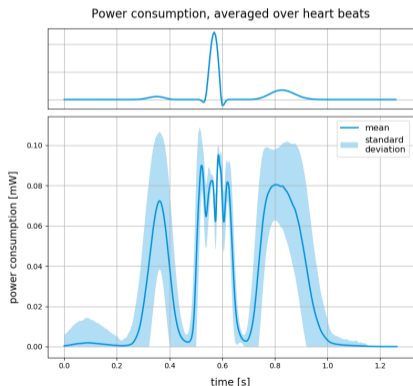- Six different anomaly types
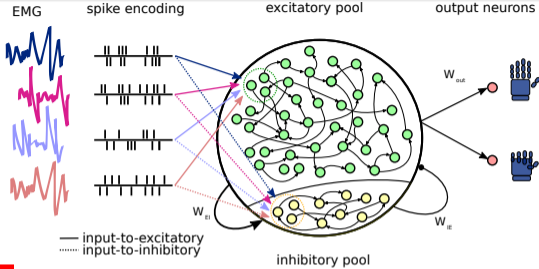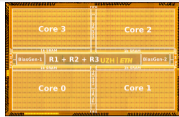- One read-out unit per anomaly



True positives rate (specificity): 91.3%
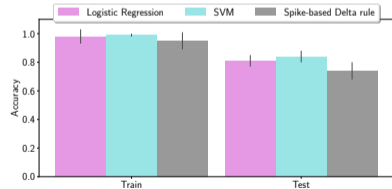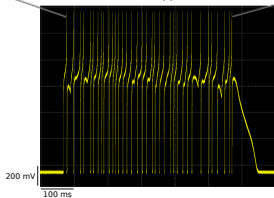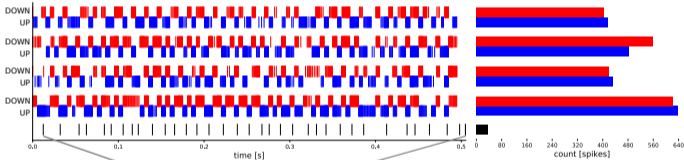True negative rate (sensitivity): 97.6%

[F. Bauer et al., 2019]

Mean neural event rate:     $14.8 \cdot 10^3$/s
Mean synaptic event rate:     $787.6 \cdot 10^3$/s
Energy per neural event:     100 pJ
Energy per synaptic event:     40 pJ
Mean power consumption:     $< 500\,\mu$W
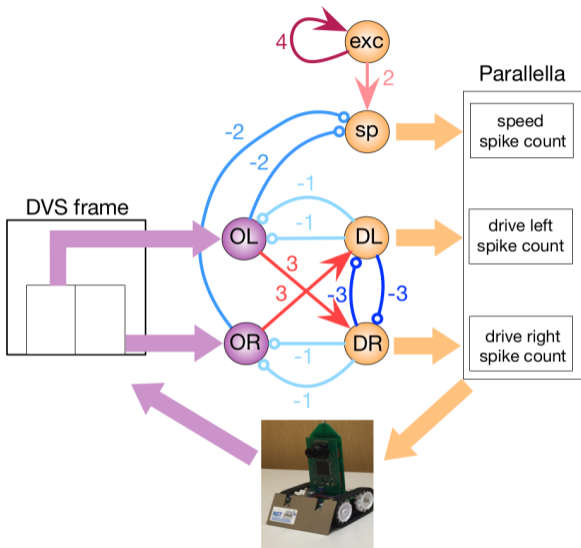


Power consumption, averaged over heart beats

open



EMG   spike encoding   excitatory pool   output neurons

— input-to-excitatory
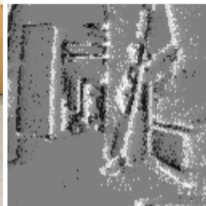···· input-to-inhibitory

inhibitory pool

[Donati et al., 2019] ‹ ›

Robot driving in the office

DVS events

[R. Kreiser et al, Frontiers in Neuromorphic Eng., 2018]

[M. Milde et al., Frontiers in Neurorobotics, 2017]

[H. Blum et al., RSS, 2017]

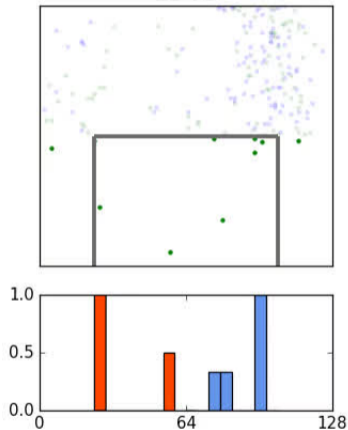[R. Kreiser et al., ISCAS, 2017]

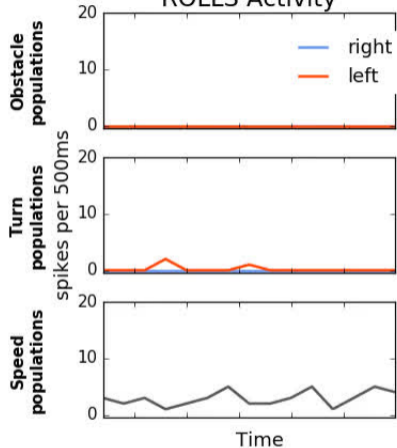[M. Milde et al., ISCAS 2017]

[R. Kreiser et al., IROS, 2018]

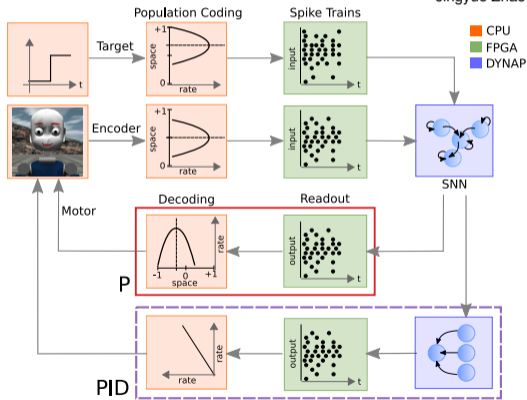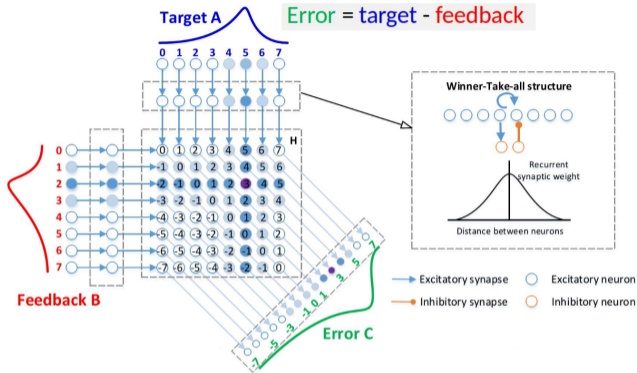[S. Glatz et al., arXiv:1810.10801, 2018]

Camera

EDVS

ROLLS Activity

Target A

Error = target - feedback

Winner-Take-all structure

Recurrent synaptic weight

Distance between neurons

Feedback B

Error C

Excitatory synapse | Excitatory neuron
Inhibitory synapse | Inhibitory neuron

Jingyue Zhao

CPU
FPGA
DYNAP

Population Coding | Spike Trains

Target

Encoder

SNN

Decoding | Readout

Motor

P

PID

Stay still

Yes

Target $(x^*, y^*)$ → SNN → Target $(\theta_1^*, \theta_2^*)$ → Close enough? → No → Command $(\theta_1^*, \theta_2^*)$

Encoder $(\theta_1, \theta_2)$

[Zhao et al. 2023 (in press, npj Robotics)]

[Zhao et al. 2023 (in press, npj Robotics)]

## Industrial Predictive Maintenance (PM)

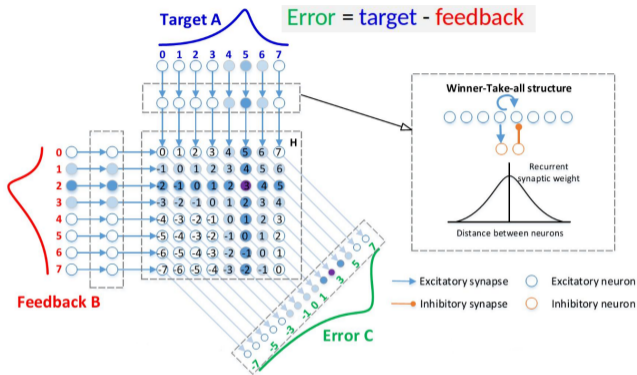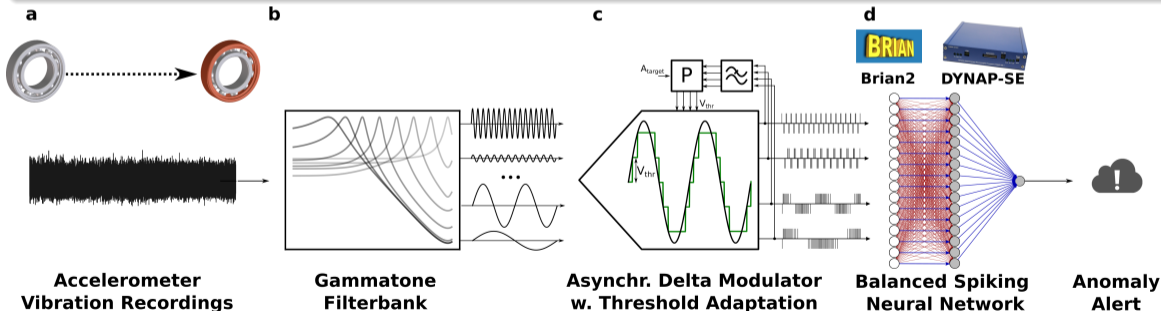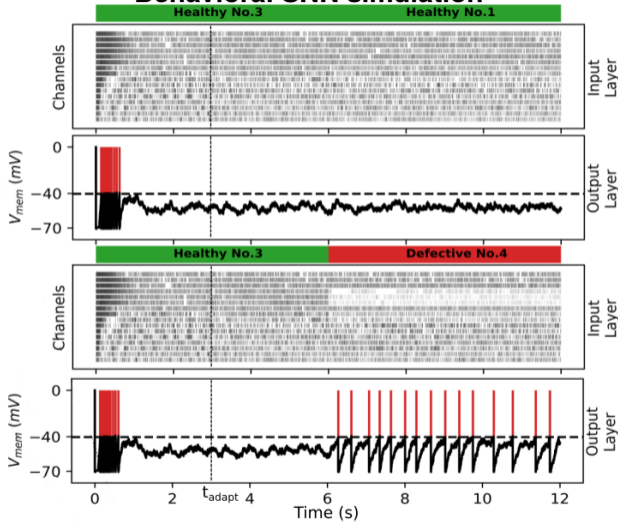- Predictive Maintenance involves the health monitoring of a degrading system.
- Vibration patterns yield valuable information about the health state of a running machine.
- PM is typically applied to large industrial tasks, but could be useful for small appliances and robots as well.



**a** | **b** | **c** | **d**

**Accelerometer Vibration Recordings** — **Gammatone Filterbank** — **Asynchr. Delta Modulator w. Threshold Adaptation** — **Balanced Spiking Neural Network** — **Anomaly Alert**

## Behavioral SNN simulation



## Validation with the DYNAP-SE chip



DETECTION TIMES (DATAPOINT) FOR RUN-TO-FAILURE DATASET

|           | b1      | b2      | b3      | b4      |
|-----------|---------|---------|---------|---------|
| LSSVM     | **533** | **823** | 893     | 700     |
| AEC       | 547     | -       | -       | -       |
| This work | 543     | 890     | **873** | **683** |

[Dennler et al., 2021] ‹ ›