

# tinyML<sup>®</sup> On Device Learning Forum

*Enabling Ultra-low Power Machine Learning at the Edge*

“Merging insights from artificial and biological neural networks for neuromorphic intelligence”

May 16, 2023



[www.tinyML.org](http://www.tinyML.org)

# The Dawn of On Device Learning in TinyML



*The goal of On Device Learning (ODL) is to make edge devices “smarter” and more efficient by observing changes in the data collected and self-adjusting/reconfiguring the device’s operating model. Optionally the “knowledge” gained by the device is shared with other deployed devices.*

Danilo Pau, Elias Fallon, Evgeni Gousev, Davis Sawyer, Ira Feldman, Christopher B. Rogers



# tinyML On Device Learning Forum

## 8/31 – 9/1 , 2022 Online

On device learning Forum

- Academia on 8/31/2022

- [On-Device Learning Under 256KB Memory](#), Song HAN, Assistant Professor, MIT EECS
- [Neural Network ODL for Wireless Sensor Nodes](#), Hiroki MATSUTANI, Professor, Keio University
- [Scalable, Heterogeneity-Aware and Trustworthy Federated Learning](#), Yiran CHEN, Professor, Duke University
- [On-Device Learning For Natural Language Processing with BERT](#), Warren J. GROSS, Professor, McGill University
- [Is on-device learning the next “big thing” in TinyML?](#) Manuel ROVERI, Associate Professor, Politecnico di Milano
- [ODL Professors Panel](#)

- Industry on 9/1/2022

- [TinyML ODL in industrial IoT](#), Haoyu REN, PhD Student, Technical University of Munich/Siemens
- [NeuroMem® wearable, hardwired sub milliwatt real time machine learning with wholly parallel access to “neuron memories” fully explainable](#), Guy PAILLET, Co-founder, General Vision
- [Using Coral Dev Board Micro for ODL innovations](#), Bill LUAN, Senior Program Manager, Google
- [Platform for Next Generation Analog AI Hardware Acceleration](#), Kaoutar EL MAGHRAOUI, Principal Research Scientist, IBM T.J Watson Research Center
- [Enabling on-device learning at scale](#), Joseph SORIAGA, Sr. Director of Technology, Qualcomm
- [Training models on tiny edge devices](#), Valeria TOMASELLI, Senior Engineer, STMicroelectronics

# tinyML EMEA Forum - On Device Learning

## 9/12 , 2022 Cyprus, In person



On device learning Forum

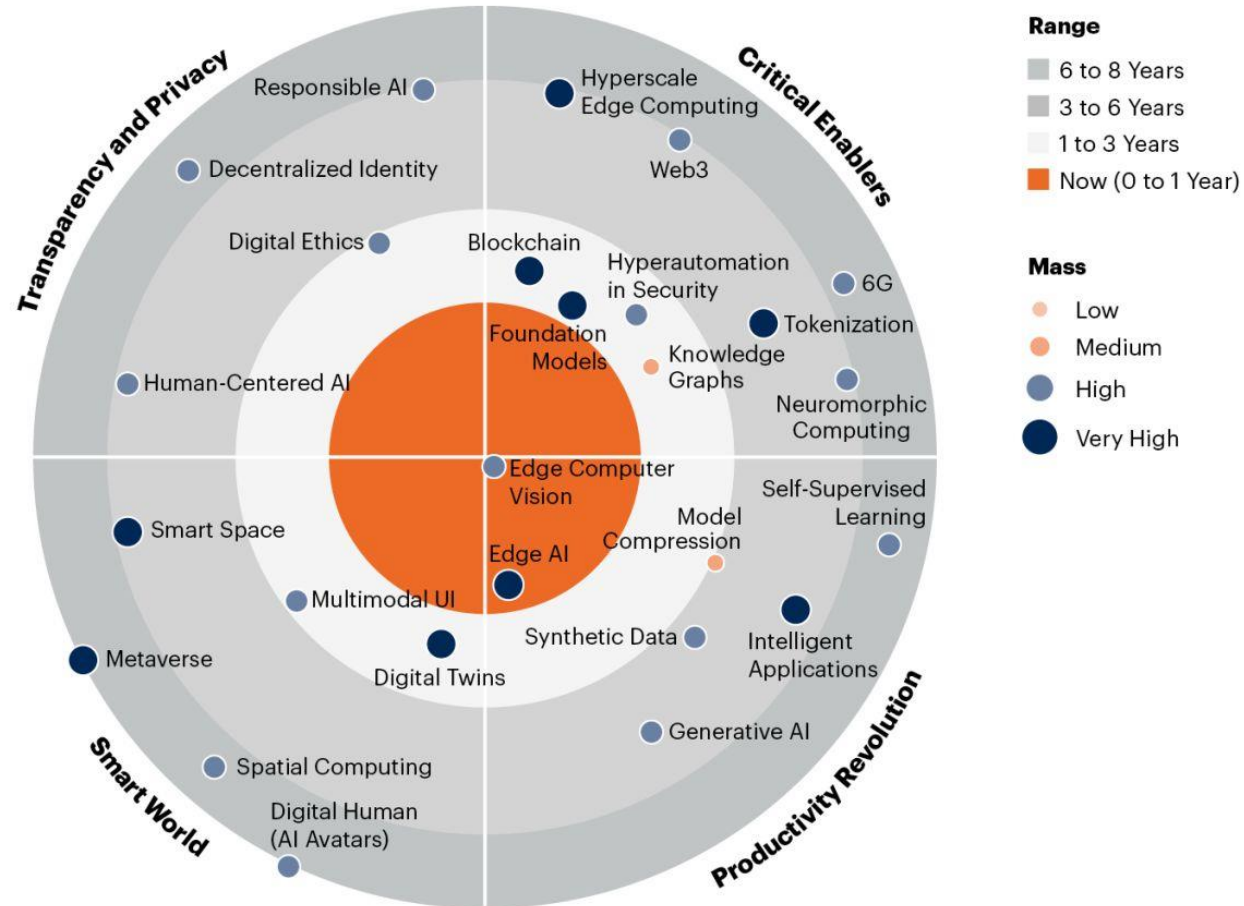
- [A framework of algorithms and associated tool for on-device tiny learning](#), Danilo PAU, Technical Director, IEEE and ST Fellow, STMicroelectronics
- [In Sensor and On-device Tiny Learning for Next Generation of Smart Sensors](#) Michele MAGNO, Head of the Project-based learning Center, ETH Zurich, D-ITET
- [Continual On-device Learning on Multi- Core RISC-V MicroControllers](#) Manuele RUSCI, Embedded Machine Learning Engineer, Greenwaves
- [On-device continuous event-driven deep learning to avoid model drift](#), Bijan MOHAMMADI, CSO, Bondzai

T I N Y



On device learning Forum

# 2023 Gartner Emerging Technologies and Trends Impact Radar



[gartner.com](https://www.gartner.com)

Note: Range measures number of years it will take the technology/trend to cross over from early adopter to early majority adoption. Mass indicates how substantial the impact of the technology or trend will be on existing products and markets.

Source: Gartner  
© 2023 Gartner, Inc. All rights reserved. CM\_GTS\_2034284

**Gartner**

# On Device Learning Forum 2023, May 16 2023

- 8:00 - 8:10 Opening remarks by **Danilo Pau**
- 8:10 - 8:40 **Charlotte Frenkel** "Merging insights from artificial and biological neural networks for neuromorphic edge intelligence"
- 8:40 - 9:40 **Giorgia Dellaferrera** "Forward Learning with Top-Down Feedback: Solving the Credit Assignment Problem without a Backward Pass"
- 9:40 - 10:10 **Guy Paillet** "NeuroMem®, Ultra Low Power hardwired incremental learning and parallel pattern recognition"
- 10:10 - 10:40 **Aida Todri-Sanial** "On-Chip Learning and Implementation Challenges with Oscillatory Neural Networks"
- 10:40 - 11:10 **Eduardo S. Pereira** "Online Learning TinyML for Anomaly Detection Based on Extreme Values Theory"
- 11:10 - 11:15 Closing remarks by Danilo Pau





Thank you, **tinyML Strategic Partners**,  
for committing to take tinyML to the next Level, together

On device learning Forum





On device learning Forum

# Executive Strategic Partners



On device learning Forum



**EDGE IMPULSE**

# The Leading Development Platform for Edge ML

[edgeimpulse.com](https://edgeimpulse.com)

**Qualcomm**  
AI research

# Advancing AI research to make efficient AI ubiquitous

## Power efficiency

Model design, compression, quantization, algorithms, efficient hardware, software tool

## Personalization

Continuous learning, contextual, always-on, privacy-preserved, distributed learning

## Efficient learning

Robust learning through minimal data, unsupervised learning, on-device learning

## A platform to scale AI across the industry



### Perception

Object detection, speech recognition, contextual fusion



### Reasoning

Scene understanding, language understanding, behavior prediction



### Action

Reinforcement learning for decision making



Edge cloud



Cloud



IoT/IIoT



Automotive



Mobile



Accelerate Your Edge Compute

**SYNTIANT**

Making Edge AI A Reality

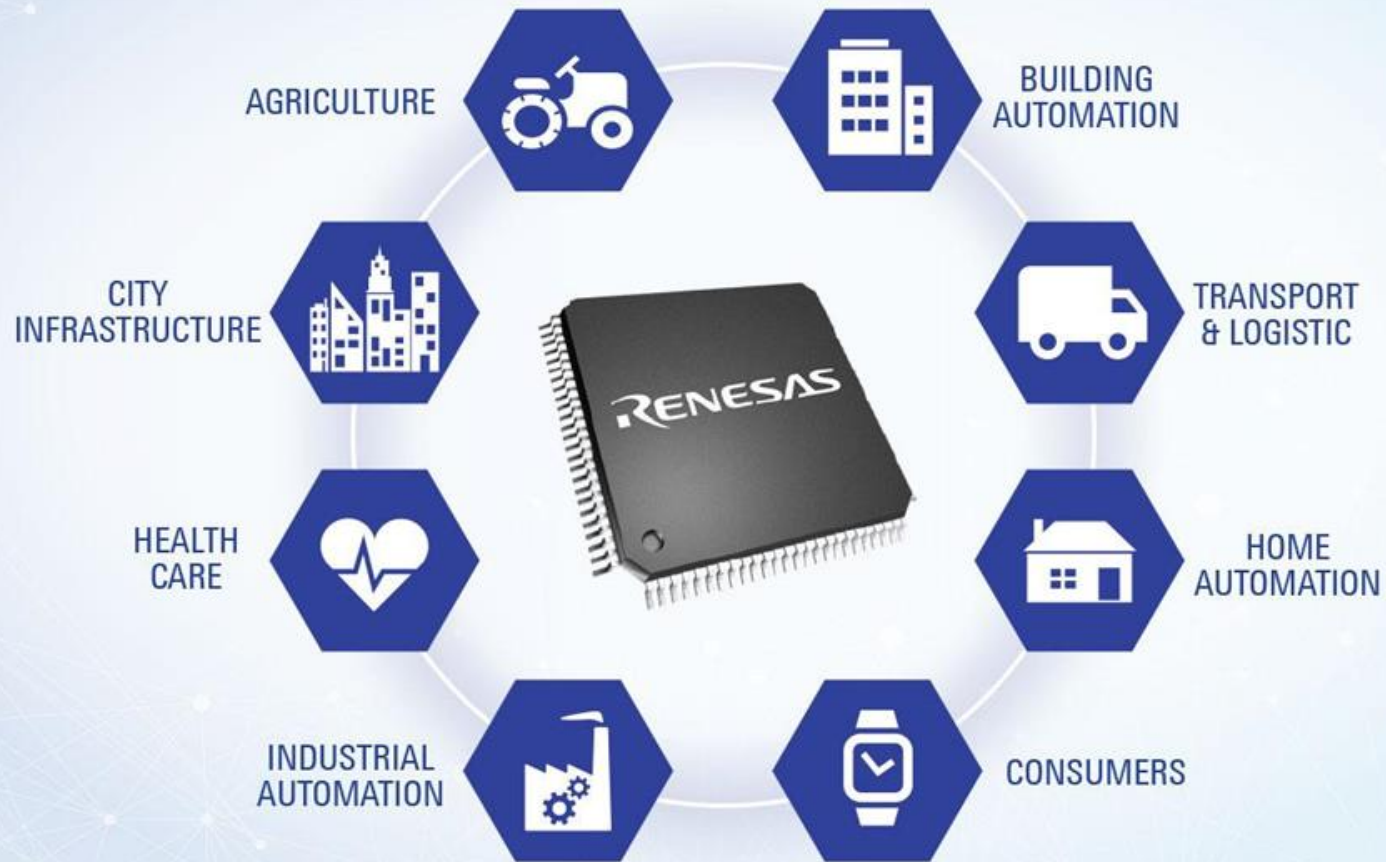
[www.syntiant.com](http://www.syntiant.com)



On device learning Forum

# Platinum Strategic Partners

# Renesas is enabling the next generation of AI-powered solutions that will revolutionize every industry sector.



[renesas.com](https://www.renesas.com)



**DEPLOY VISION AI  
AT THE EDGE AT SCALE**

**SONY**



On device learning Forum

# Gold Strategic Partners



AHEAD OF WHAT'S POSSIBLE™



AHEAD OF WHAT'S POSSIBLE™

Where what if  
becomes what is.

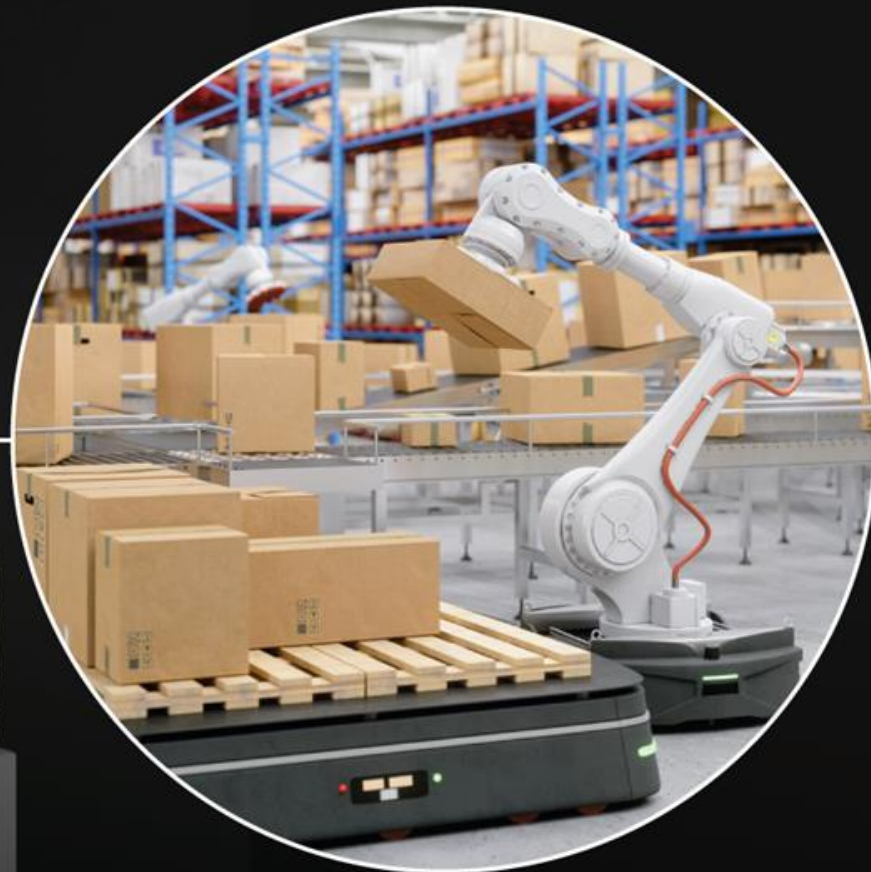
Witness potential made possible at [analog.com](http://analog.com).



PRO™

Easily deploy your  
tinyML solutions with  
Arduino Pro

[arduino.cc/pro](https://arduino.cc/pro)



Made In Italy

arm AI



Powering tinyML Innovation

# Arm AI Virtual Tech Talks

The latest in AI trends, technologies & best practices from Arm and our Ecosystem Partners.

Demos, code examples, workshops, panel sessions and much more!

Fortnightly Tuesday @ 4pm GMT/8am PT

Find out more:

[www.arm.com/techtalks](http://www.arm.com/techtalks)

Decarbonization

Digitalization



Driving decarbonization and digitalization. Together.

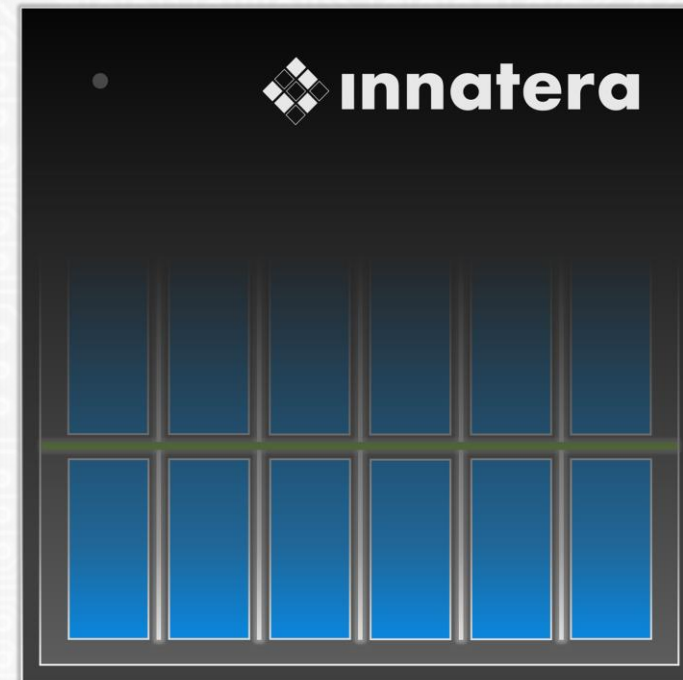
**Infineon serving all target markets as**  
**Leader in Power Systems and IoT**

[www.infineon.com](http://www.infineon.com)





# NEUROMORPHIC INTELLIGENCE FOR THE SENSOR-EDGE



[www.innatera.com](http://www.innatera.com)



Microsoft

The Right Edge AI Tools Can Make or Break Your Next Smart IoT Product



## Analytics Toolkit Suite



[sensiml.com/tinyML](https://sensiml.com/tinyML)





life.augmented

**STMicroelectronics provides extensive solutions to make tiny Machine Learning easy**



# ENGINEERING EXCEPTIONAL EXPERIENCES

We engineer exceptional experiences for consumers in the home, at work, in the car, or on the go.

[www.synaptics.com](http://www.synaptics.com)





On device learning Forum

# Silver Strategic Partners





# Join Growing tinyML Communities:



14.7k members in  
47 Groups in 39 Countries

**tinyML - Enabling ultra-low Power ML at the Edge**

<https://www.meetup.com/tinyML-Enabling-ultra-low-Power-ML-at-the-Edge/>



4k members  
&  
11.6k followers

**The tinyML Community**

<https://www.linkedin.com/groups/13694488/>





On device learning Forum



Subscribe to tinyML YouTube Channel for updates and notifications (including this video)

www.youtube.com/tinyML



tinyML 4.33K subscribers **9.4k subscribers, 559 videos with 327k views**

HOME VIDEOS PLAYLISTS COMMUNITY CHANNELS ABOUT

106 views · 4 days ago	138 views · 4 days ago	54 views · 4 days ago	47 views · 4 days ago	132 views · 4 days ago	137 views · 4 days ago
122 views · 4 days ago	262 views · 2 weeks ago	511 views · 3 weeks ago	229 views · 3 weeks ago	265 views · 3 weeks ago	286 views · 1 month ago
351 views · 1 month ago	462 views · 2 months ago	374 views · 2 months ago	133 views · 2 months ago	287 views · 2 months ago	336 views · 2 months ago
378 views · 2 months ago	214 views · 2 months ago	448 views · 2 months ago	159 views · 2 months ago	190 views · 2 months ago	545 views · 2 months ago



On device learning Forum



FOUNDATION

tinyML EMEA  
Innovation Forum

June 26 -28, 2023

Amsterdam

*EMEA 2023*

<https://www.tinyml.org/event/emea-2023>

More sponsorships are available: [sponsorships@tinyML.org](mailto:sponsorships@tinyML.org)



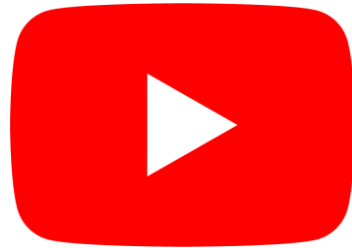
On device learning Forum

# Reminders

Slides & Videos will be posted tomorrow



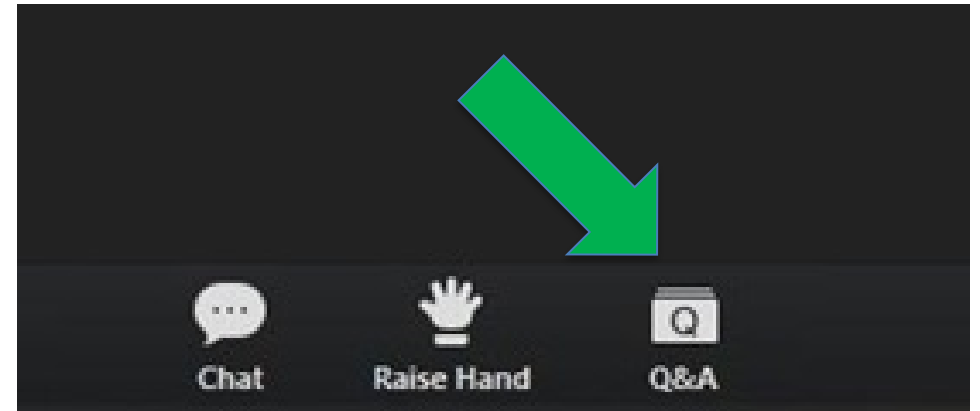
[tinyml.org/forums](https://tinyml.org/forums)



[youtube.com/tinyml](https://youtube.com/tinyml)



Please use the Q&A window for your questions

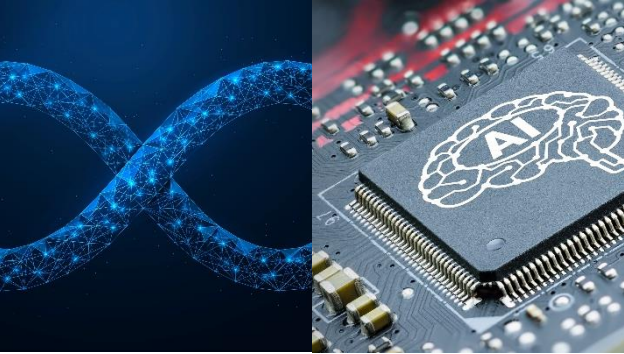




# Charlotte Frenkel



Charlotte Frenkel received the Ph.D. degree from Université catholique de Louvain (UCLouvain), Belgium, in 2020. After a postdoc at the Institute of Neuroinformatics, UZH and ETH Zürich, Switzerland, she joined Delft University of Technology, The Netherlands, as an Assistant Professor in July 2022. Her research focuses on neuromorphic integrated circuit design and learning algorithms for adaptive edge computing. She received a best paper award at the IEEE ISCAS 2020 conference, as well as the FNRS Nokia Bell Labs Scientific Award, the FNRS IBM Innovation Award and the UCLouvain/ICTEAM Best Thesis Award for her Ph.D. thesis. She serves as a TPC member for the tinyML Research Symposium and for the IEEE ESSCIRC, ISLPED, and DATE conferences.



# Merging insights from artificial and biological neural networks for neuromorphic intelligence

Charlotte Frenkel

Delft University of Technology, Microelectronics Department

[c.frenkel@tudelft.nl](mailto:c.frenkel@tudelft.nl)

tinyML On-Device Learning Forum 2023  
Online, May 16<sup>th</sup> 2023

# Outline

① From neuroscience to AI and back again...

...which perspective?

...which starting point?

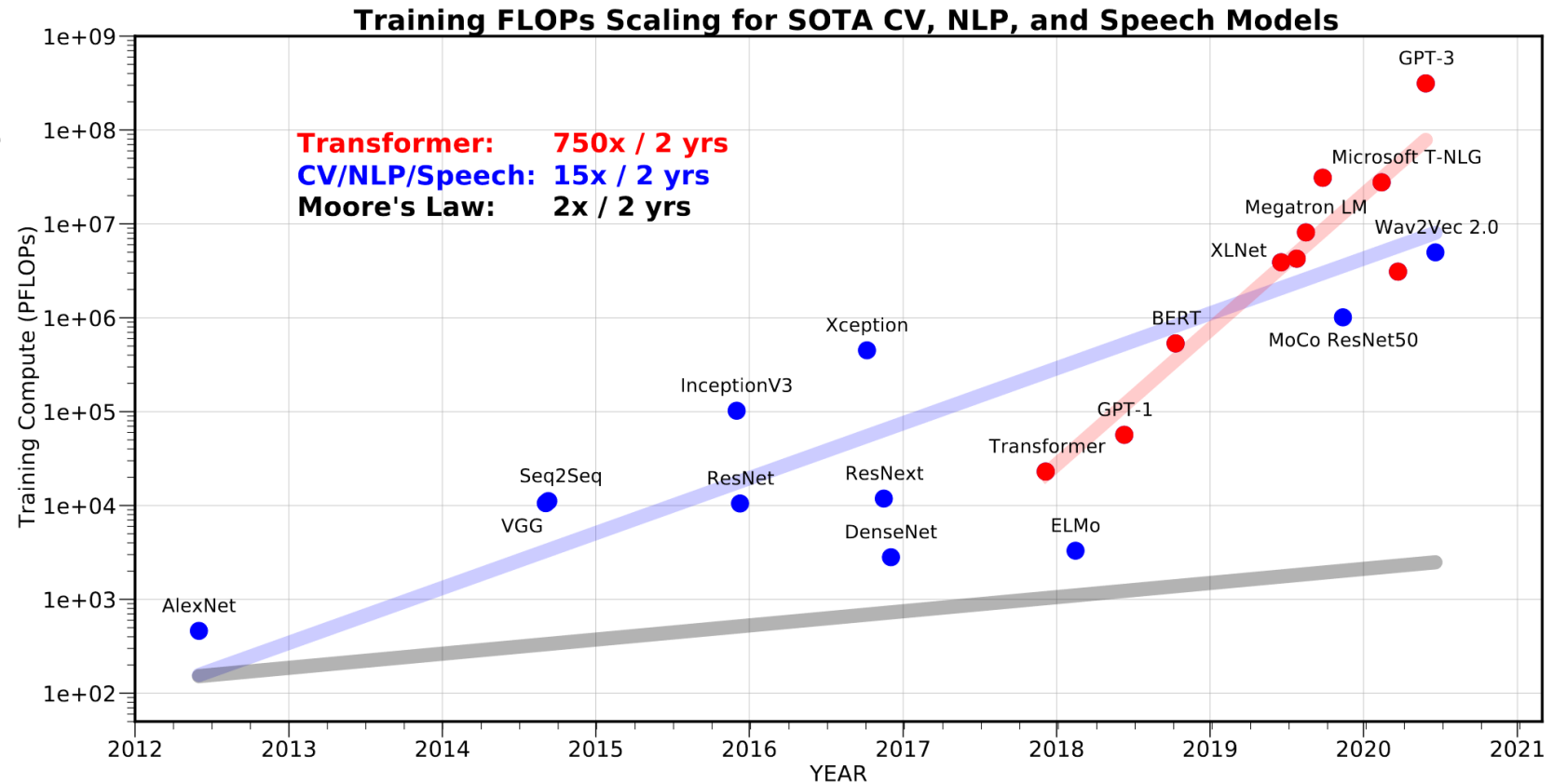
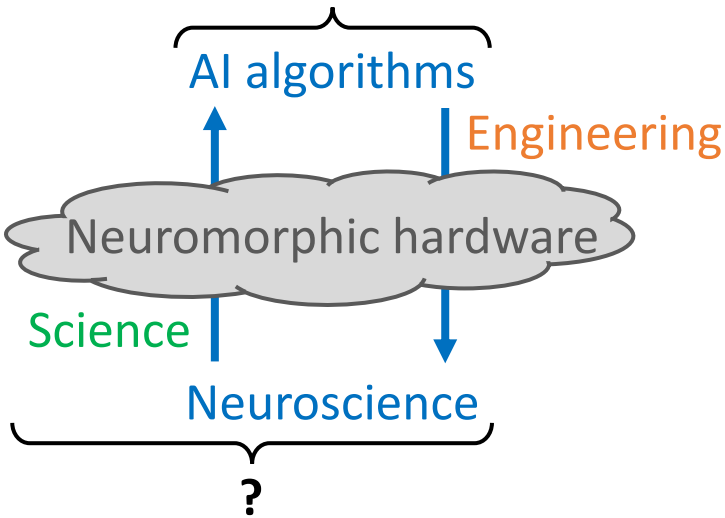
# Outline

- ① From neuroscience to AI and back again...
  - ...which perspective?
  - ...which starting point?
- ② Why should we bother with neuroscience?
- ③ How can we morph these questions into interesting solutions for on-device-learning?

# From neuroscience to AI and back again

*Which starting point? Which perspective?*

**AI without hardware is unsustainable**



[A. Gholami, *RiseLab Medium Post*, 2021]

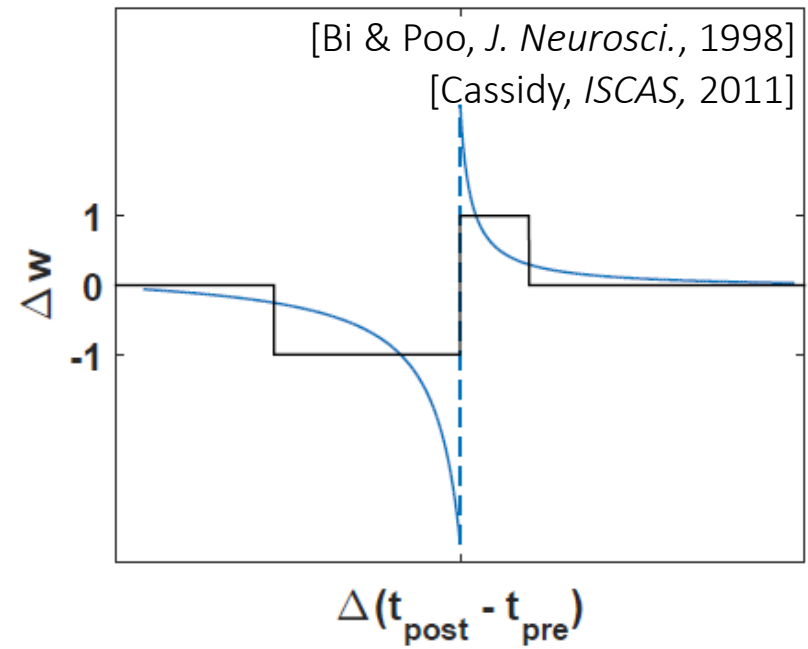
# Outline

- ① From neuroscience to AI and back again...  
...which perspective?  
...which starting point?
- ② Why should we bother with neuroscience?
- ③ How can we morph these questions into interesting solutions for on-device-learning?

# Synaptic plasticity rules – Neuroscience as the starting point

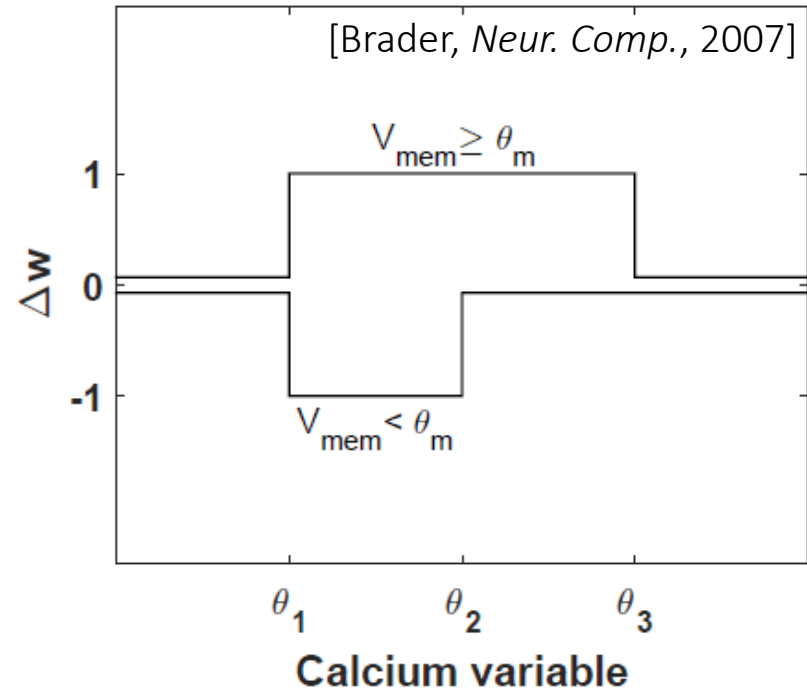
*Synergy with hardware: the perspective of data locality*

## Spike-timing-dependent plasticity (STDP)

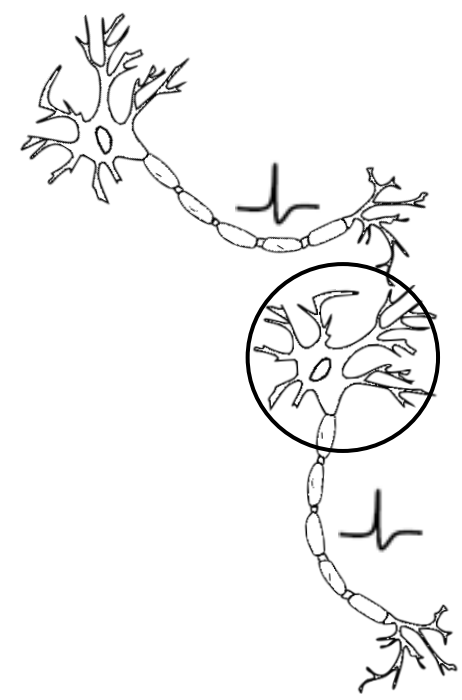


✓ Local

## Spike-dependent synaptic plasticity (SDSP)

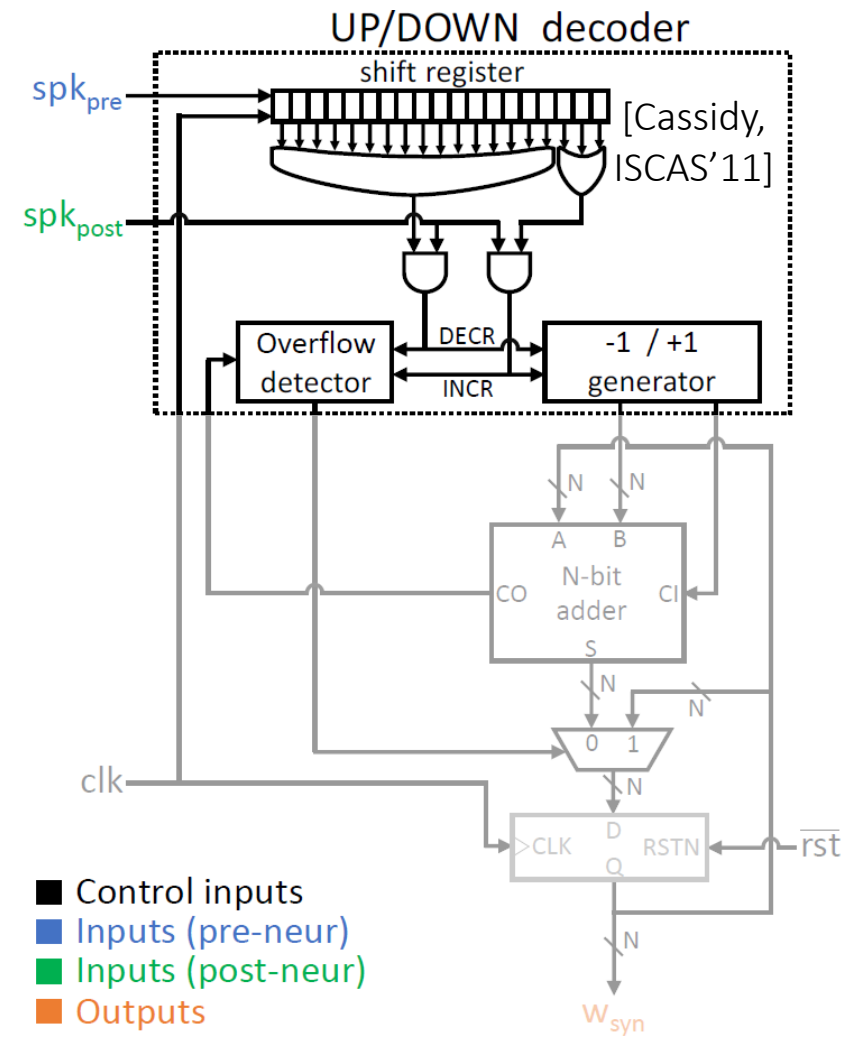


✓ Local

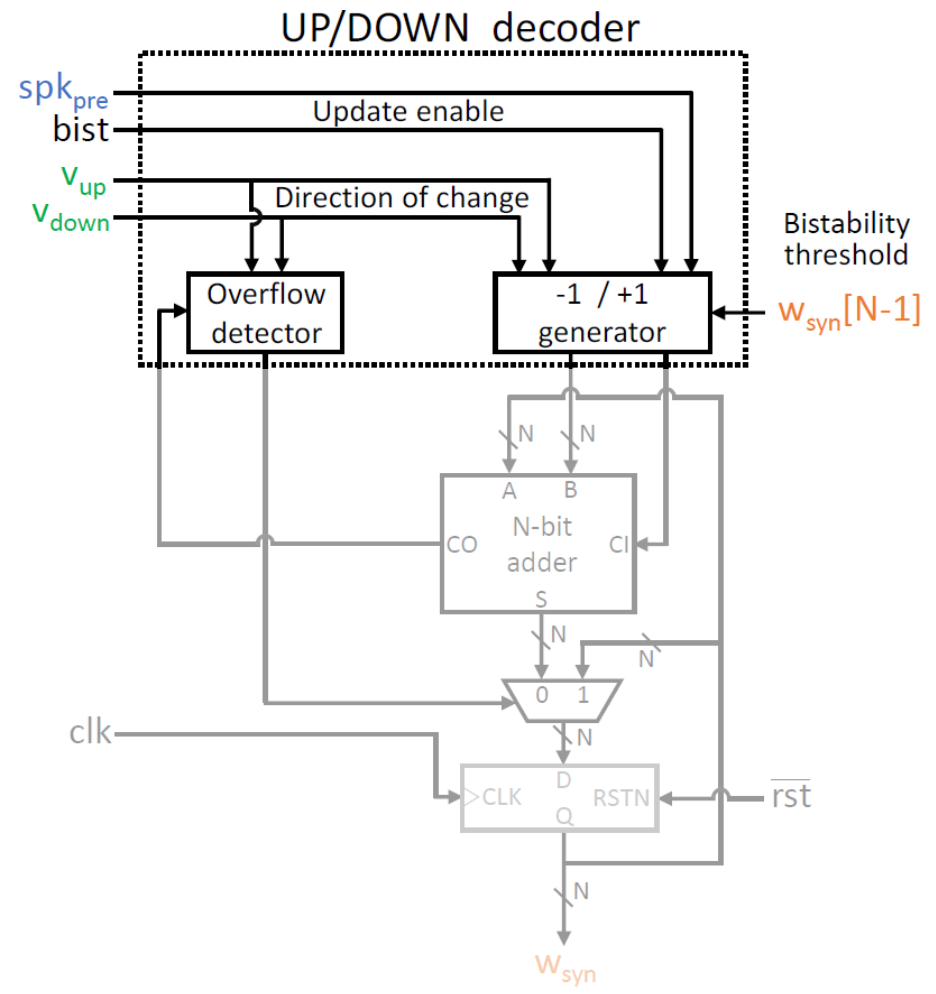


# Synaptic plasticity rules – Neuroscience as the starting point

*Synergy with hardware: the perspective of data locality*



**STDP**



**SDSP**

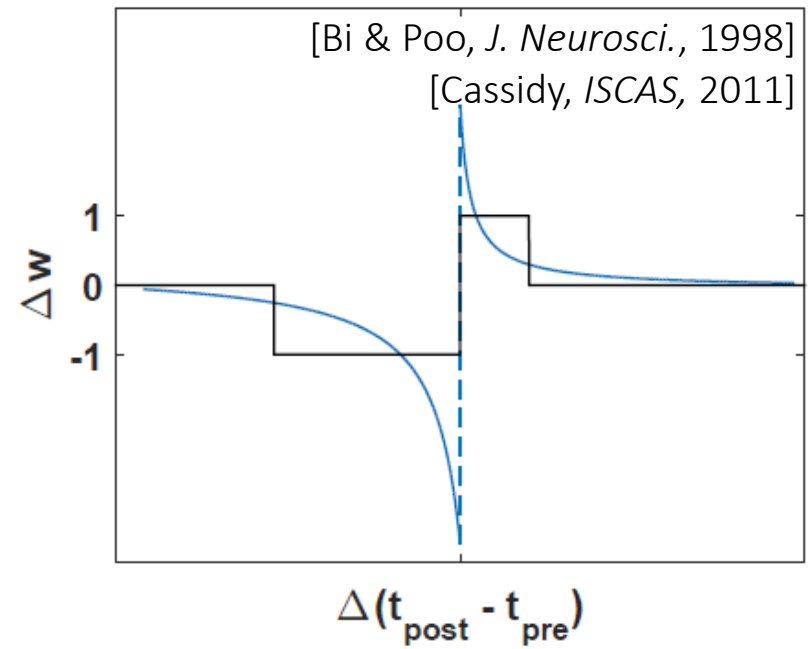
- Control inputs
- Inputs (pre-neur)
- Inputs (post-neur)
- Outputs

[Frenkel, *Trans. BioCAS*, 2019]

# Synaptic plasticity rules – Neuroscience as the starting point

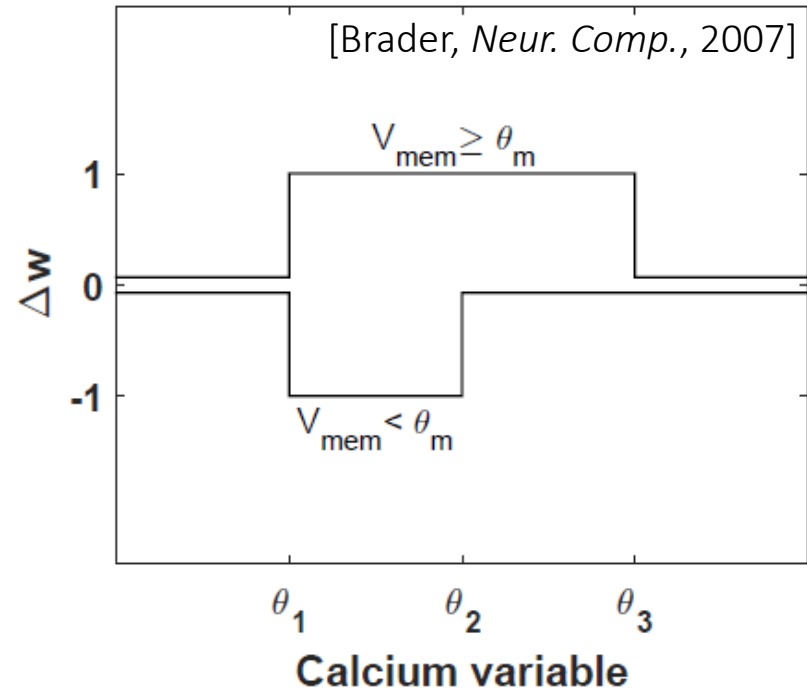
*Synergy with hardware: the perspective of data locality*

## Spike-timing-dependent plasticity (STDP)

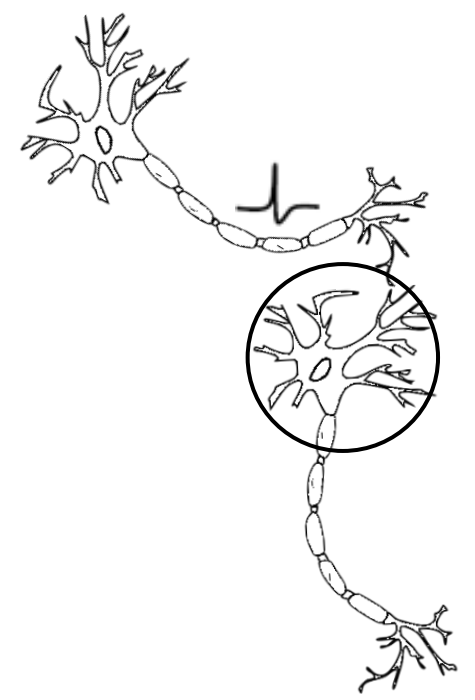


- ✓ **Local** in space
- ✗ **Non-local** in time

## Spike-dependent synaptic plasticity (SDSP)



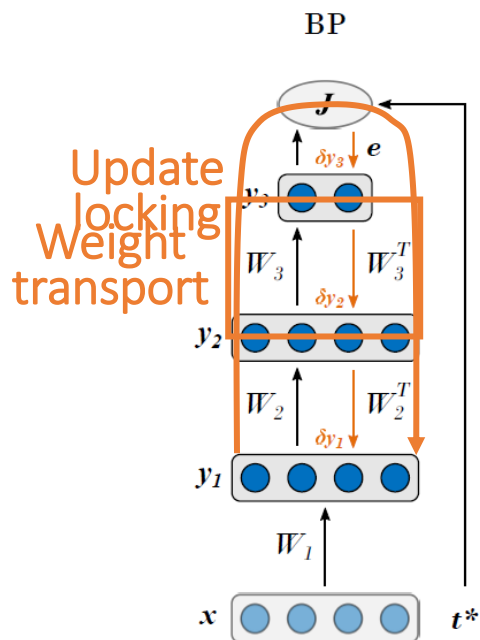
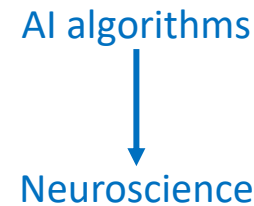
- ✓ **Local** in space
- ✓ **Local** in time



[Clopath and Gerstner, *Front. Syn. Neuro.*, 2010]

# Neural network training – Bio-plausibility as the end goal

*Synergy with hardware: latency, memory access patterns*



(signed-based)

sign(e)

$$e_c = \begin{cases} 1 - y_{3c} > 0 \\ 0 - y_{3c} < 0 \end{cases}$$

Aligns with output-independent target signals in the dendritic instructive pathways of cortical intermediate-layer neurons?  
[Magee & Grienberger, Ann. Rev. Neuro., 2020]

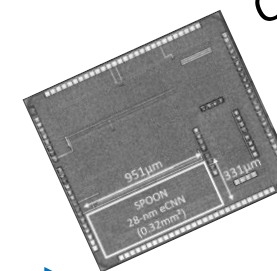
$\delta y_k$

$$\frac{\partial J}{\partial y_k} = W_{k+1}^T \delta z_{k+1}$$



(for classification problems)

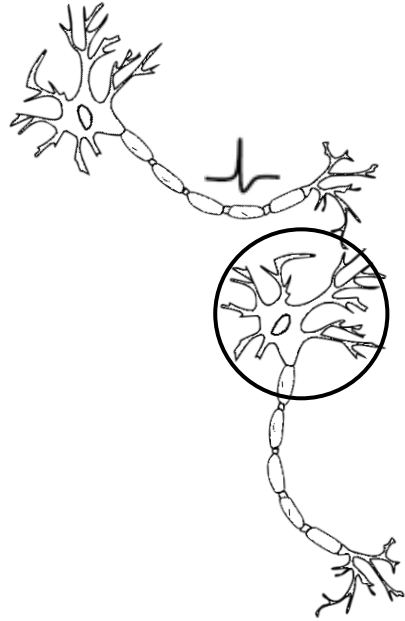
Computational and memory cost



Only ~15% overhead in power and area [Frenkel, ISCAS'20] (🏆 Best paper award)

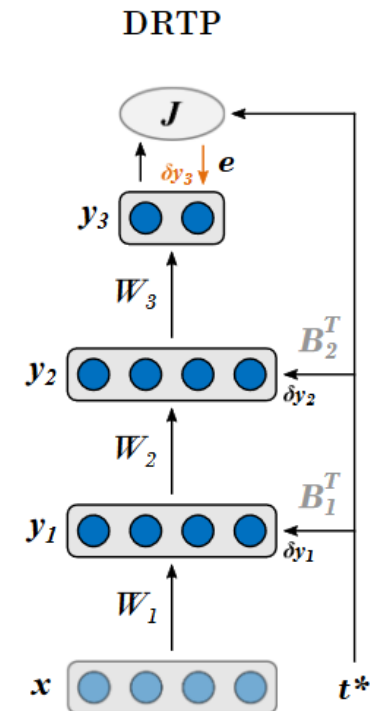
# HW efficiency and bio-plausibility are often two sides of the same coin!

*Many more examples: quantization, stochastic computing, event-driven computation,...*



Designing efficient hardware hints toward bio-plausible mechanisms

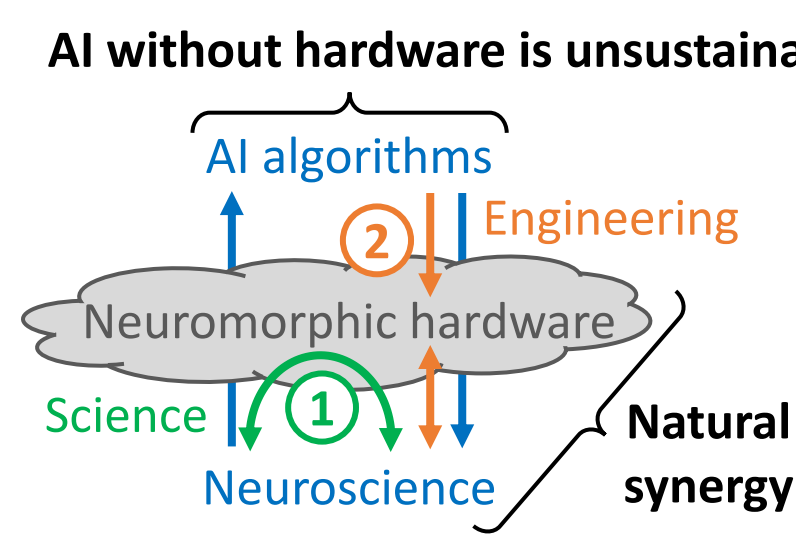
Bringing AI closer to neuroscience leads to hardware efficiency



# From neuroscience to AI and back again

*Which starting point? Which perspective?*

**AI without hardware is unsustainable**



## ① Bottom-up science-driven approach

- ✓ Analysis-by-synthesis
- ✗ Difficult to scale efficiently to real-world problems

## ② Top-down engineering-driven approach

- ✓ Starts from working solutions to real-world problems
- ✗ Which “salt & pepper” from neuroscience?

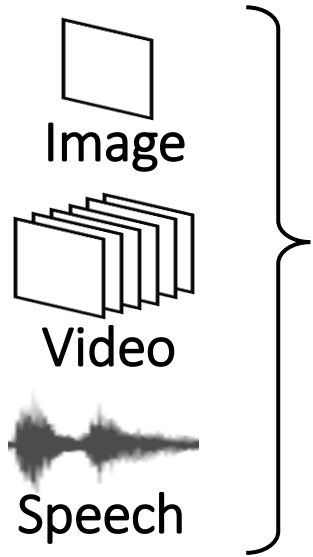
**Neuromorphic intelligence:**

② should be fed by ①

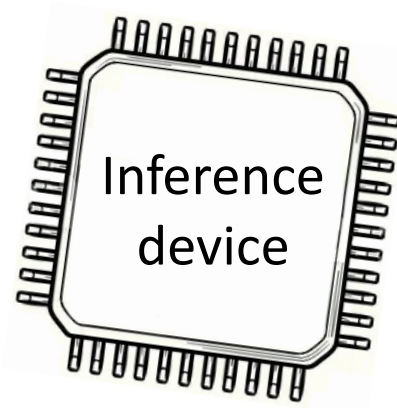
# Outline

- ① From neuroscience to AI and back again...
  - ...which perspective?
  - ...which starting point?
- ② Why should we bother with neuroscience?
- ③ How can we morph these questions into interesting solutions for on-device-learning?

# Why on-chip learning?

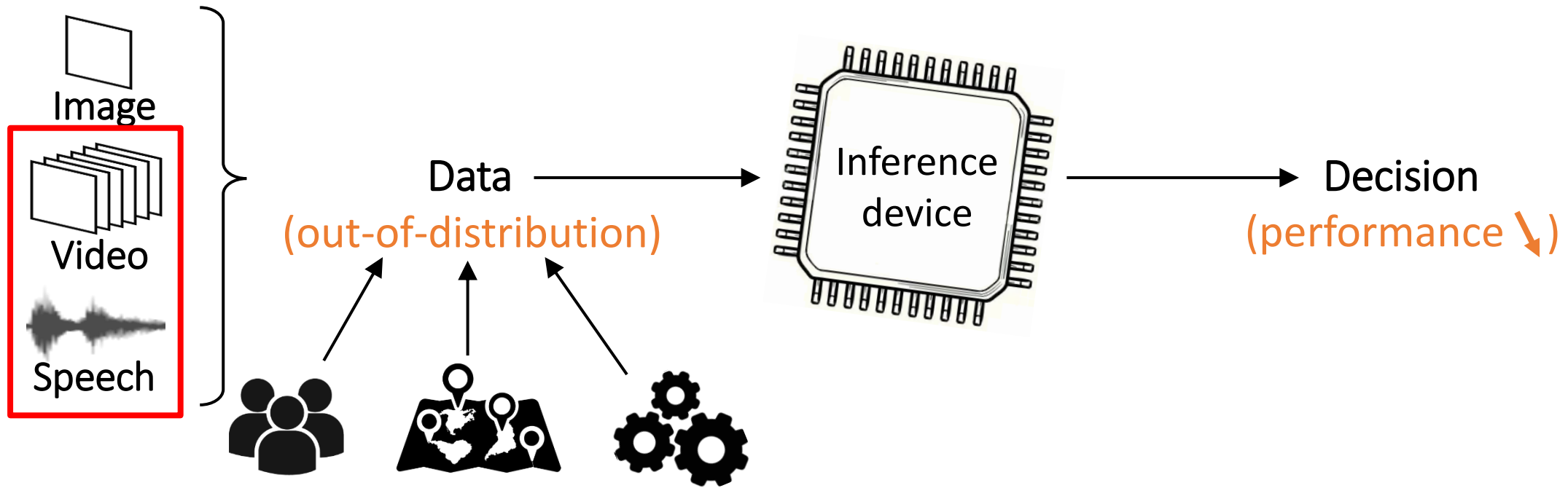


Data  
(in-distribution)



Decision  
(performance within specs)

# Why on-chip learning?



Different users, environments, task requirements

## More training data before deployment?

Issues: cost, robustness, flexibility

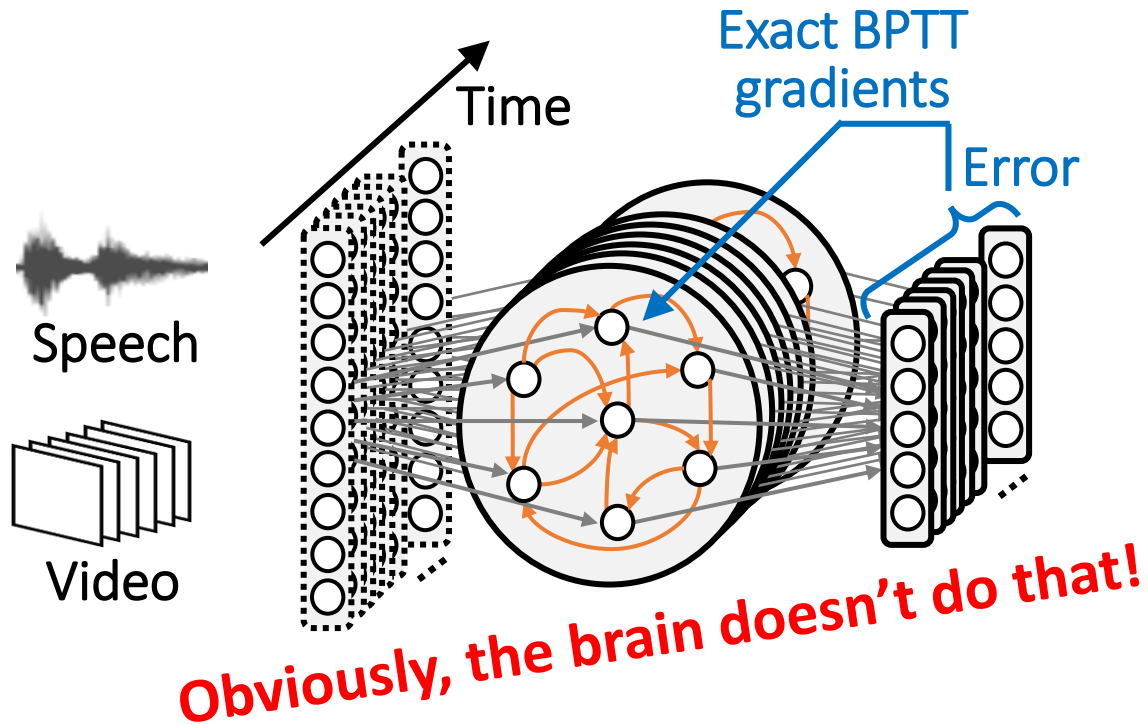
## Data exchange with the cloud?

Issues: power budget, privacy

**On-chip training**  
**(bio-inspired, end-to-end)**

# Why is on-chip learning over second-long timescales difficult?

*Let's solve a yet unsolved engineering challenge!*



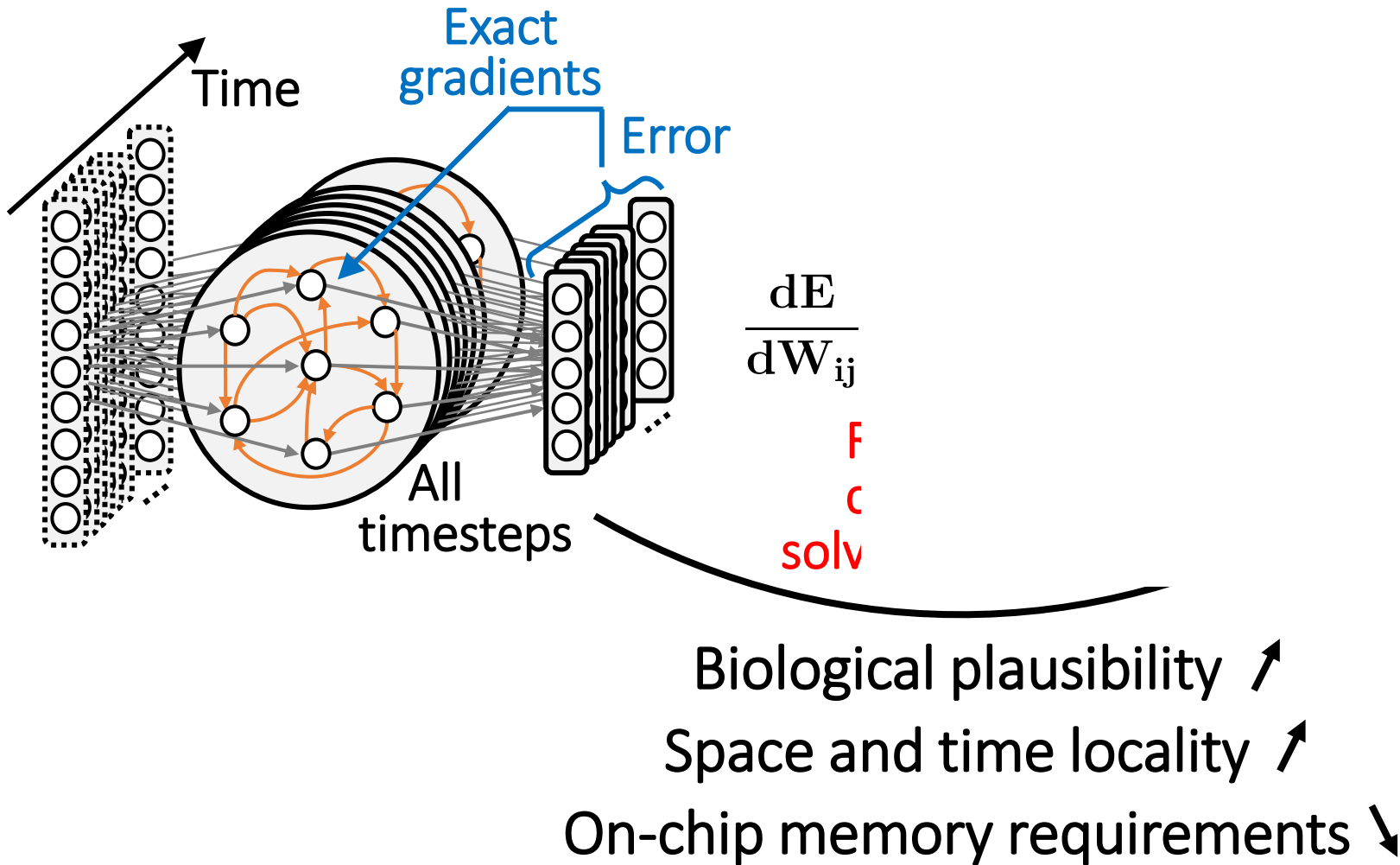
- Unrolling in time: very deep network (current learning ICs for static stimuli:  $\leq 3$  layers)
- Intractable memory/latency requirements
- No end-to-end on-chip solution to date

**Key challenge: On-chip learning over long timescales while keeping a fine-grained temporal resolution**

# The bio-inspired solution

*Backward- vs. forward-mode training*

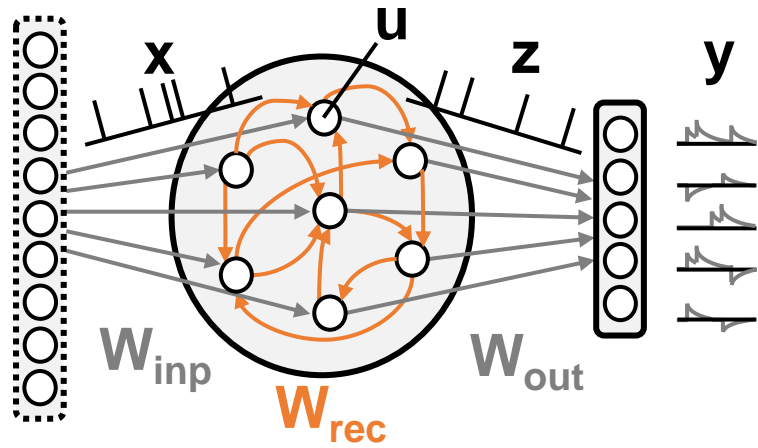
## Backprop through time (BPTT, backward)



# Algorithmic developments toward efficient long-term on-chip training

*Network definitions and evaluation task*

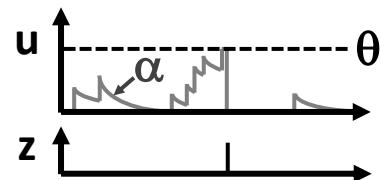
## Network model



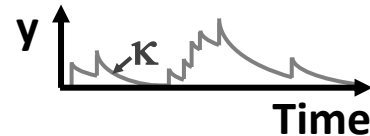
Sample input



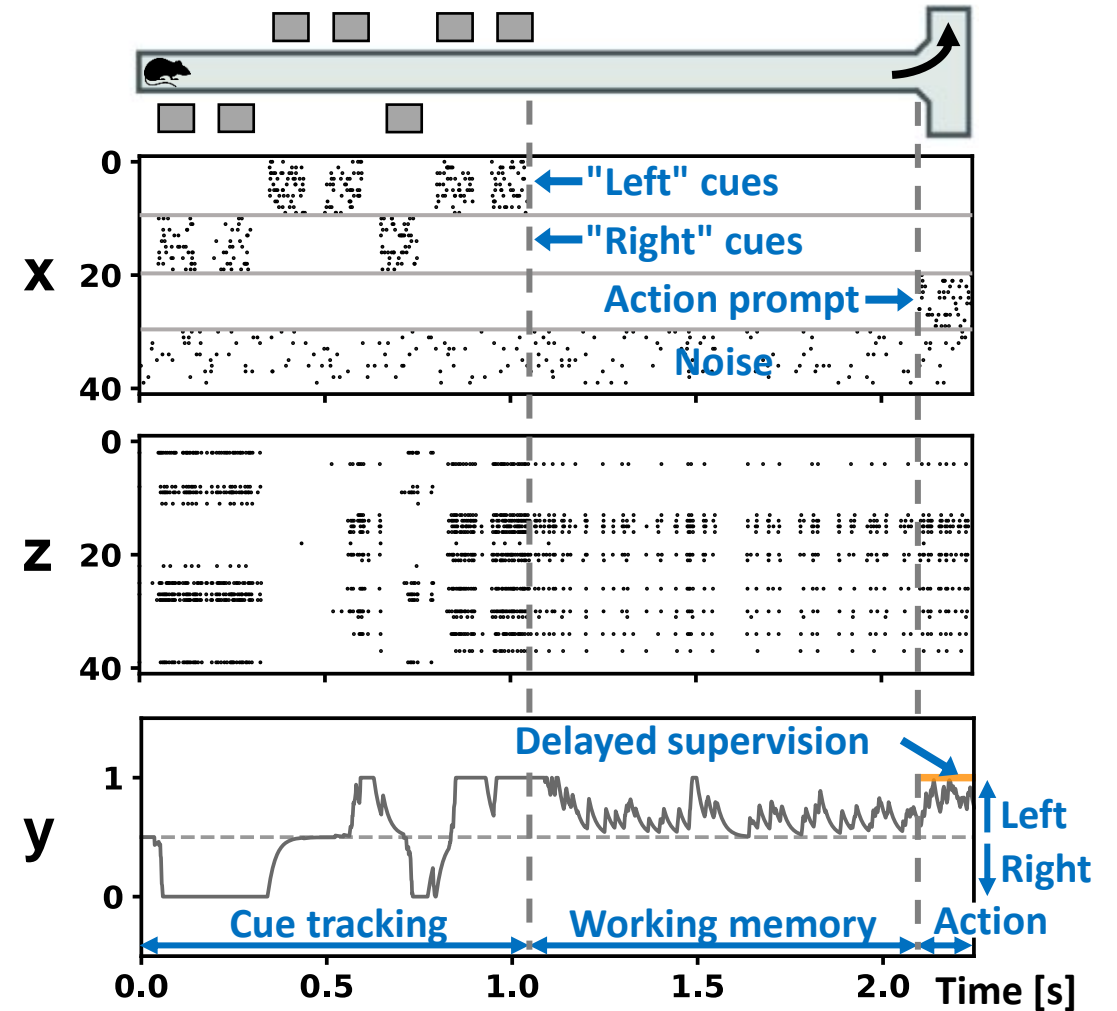
Leaky integrate-and-fire (LIF)



Leaky integrator (LI)



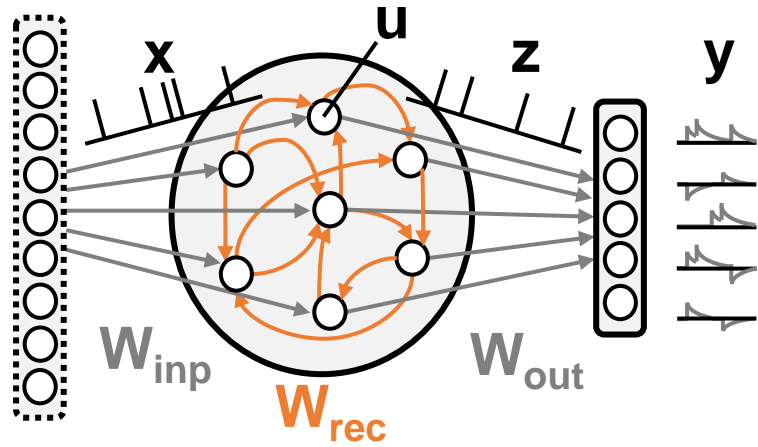
## Delayed-supervision navigation task



# Algorithmic developments – Step 1

*Neuron model selection*

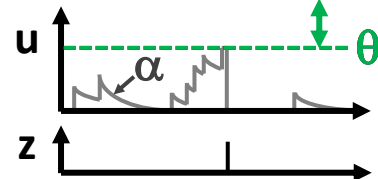
## Network model



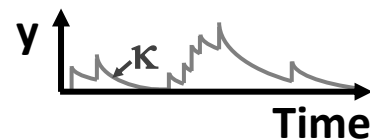
Sample input



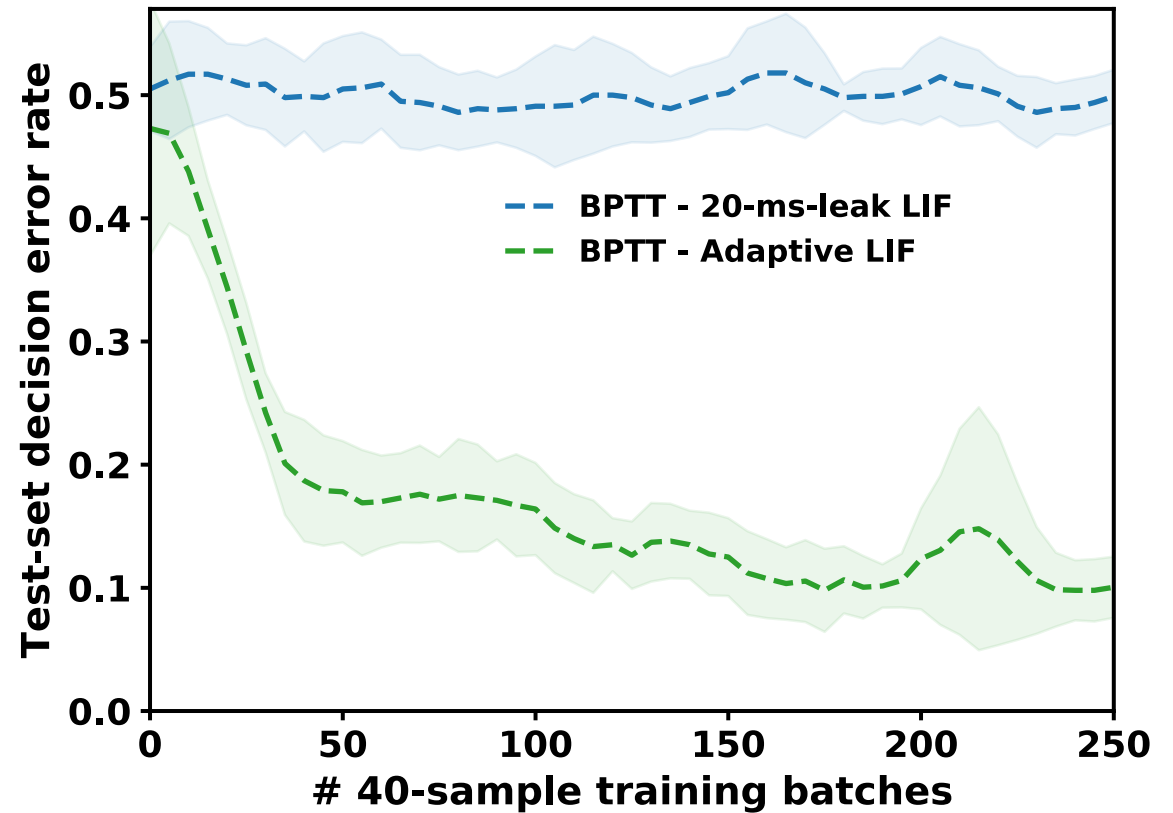
Leaky integrate-and-fire (LIF)  
or **adaptive LIF (ALIF)**



Leaky integrator (LI)



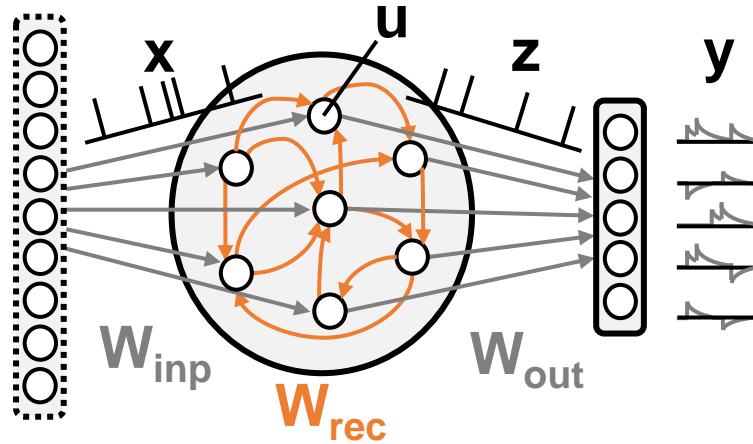
## Task performance



# Algorithmic developments – Step 1

*Neuron model selection*

## Network model



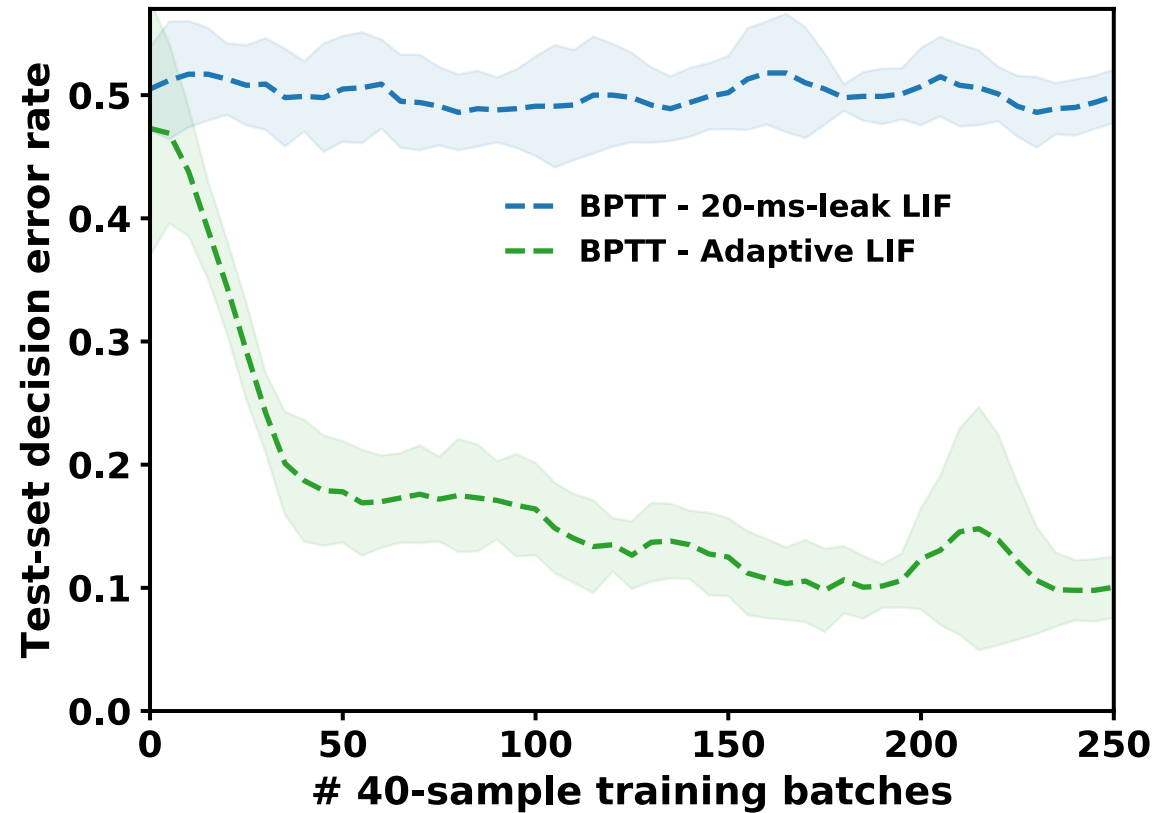
Leaky integrate-and-fire (LIF):

✗ Only a short time constant (~20-ms leak)

Adaptive LIF (ALIF):

✓ Embeds threshold adaptation over 100s of ms

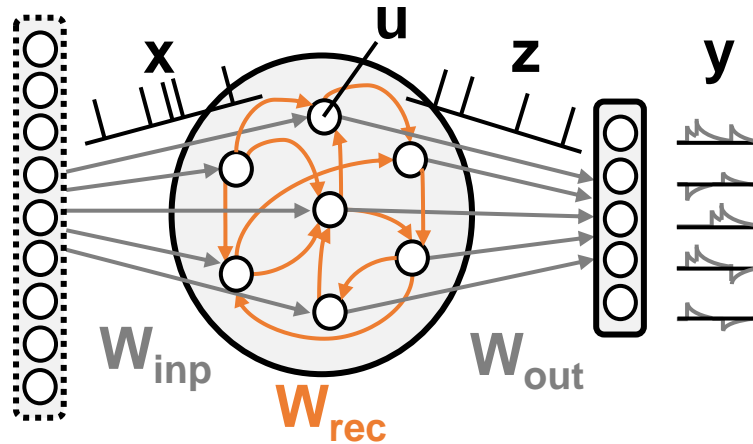
## Task performance



# Algorithmic developments – Step 1

## *Neuron model selection*

### Network model



### Leaky integrate-and-fire (LIF):

- ✗ Only a short time constant (~20-ms leak)
- ✓ Eligibility traces: simple activity LPF

### Adaptive LIF (ALIF):

- ✓ Embeds threshold adaptation over 100s of ms
- ✗ ET: Complex per-synapse multi-scale filtering

### LIF with configurable leak:

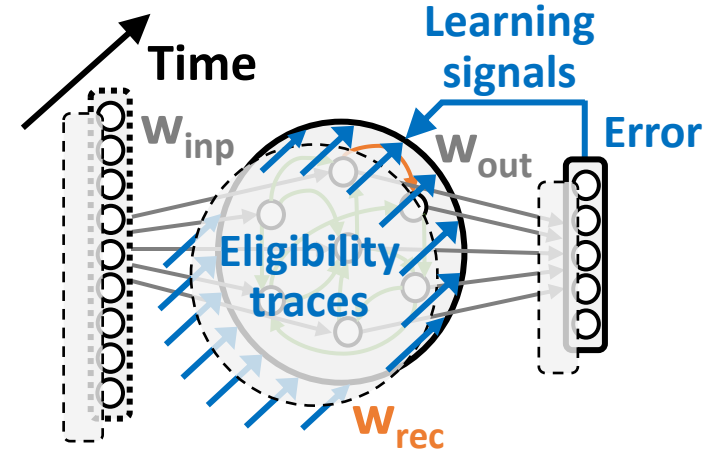
- ✓ Flexible time constant (ms to sec)
- ✓ Eligibility traces: simple activity LPF
- ≈ Less biologically plausible

# Algorithmic developments – Step 2

*Space and time locality*

$$\frac{dE}{dW_{ij}} \approx \sum_t L_j^t e_{ji}^t$$

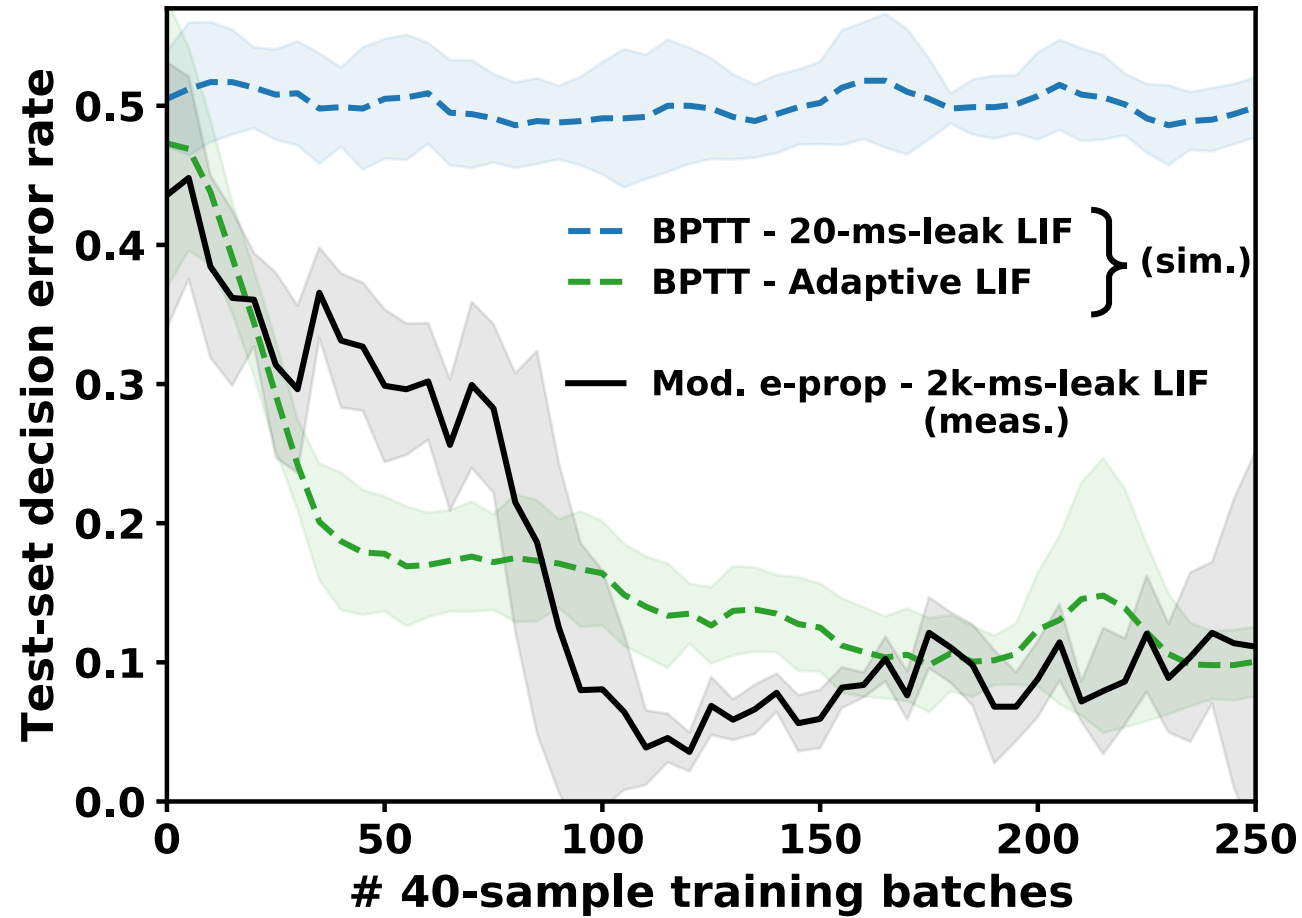
Local decoupling of  
space and time



Step 3 – Stochastic weight updates allow reducing weight resolution to 8 bits.

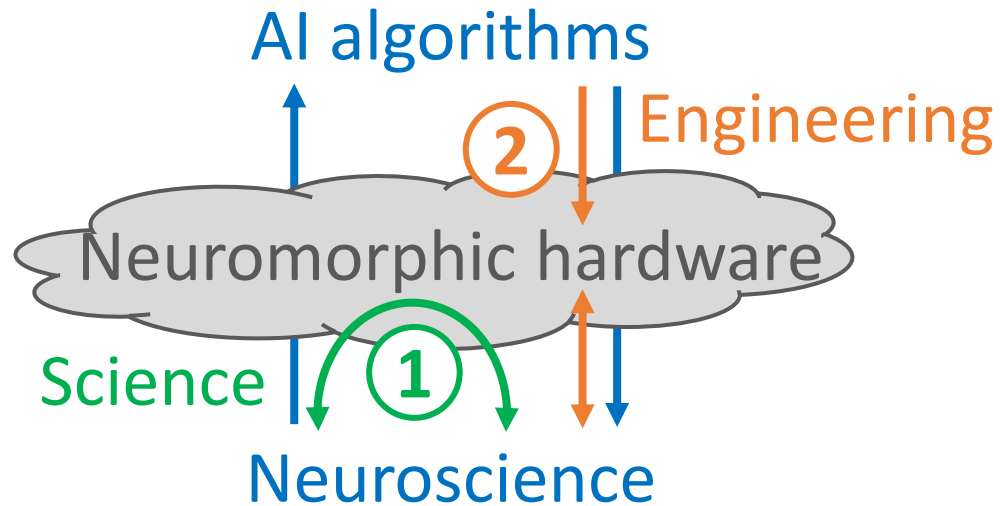
# Algorithmic developments

*Final performance*

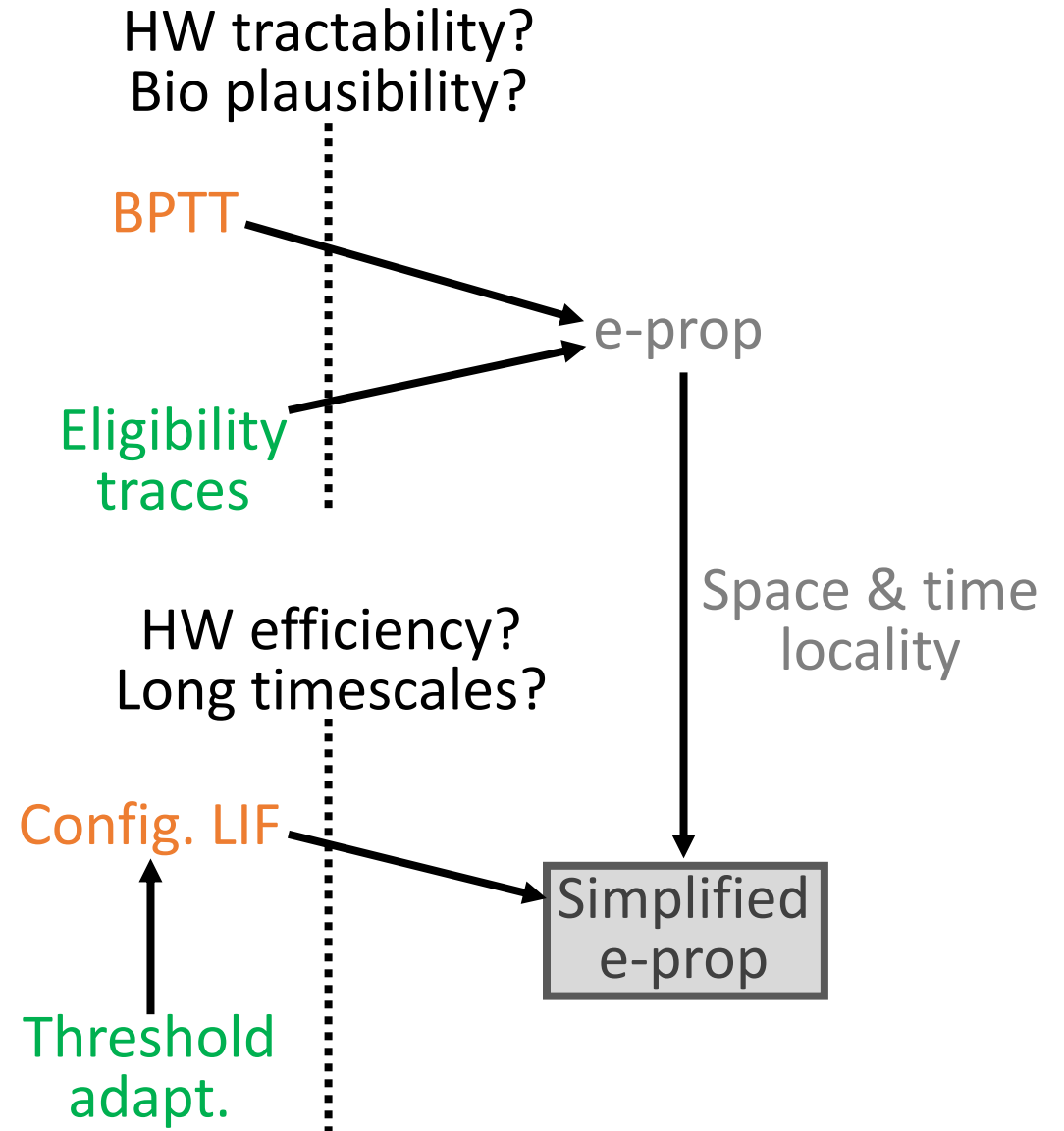


# From neuroscience to AI and back again

*Which starting point? Which perspective?*



**Neuromorphic intelligence:**  
② should be fed by ①

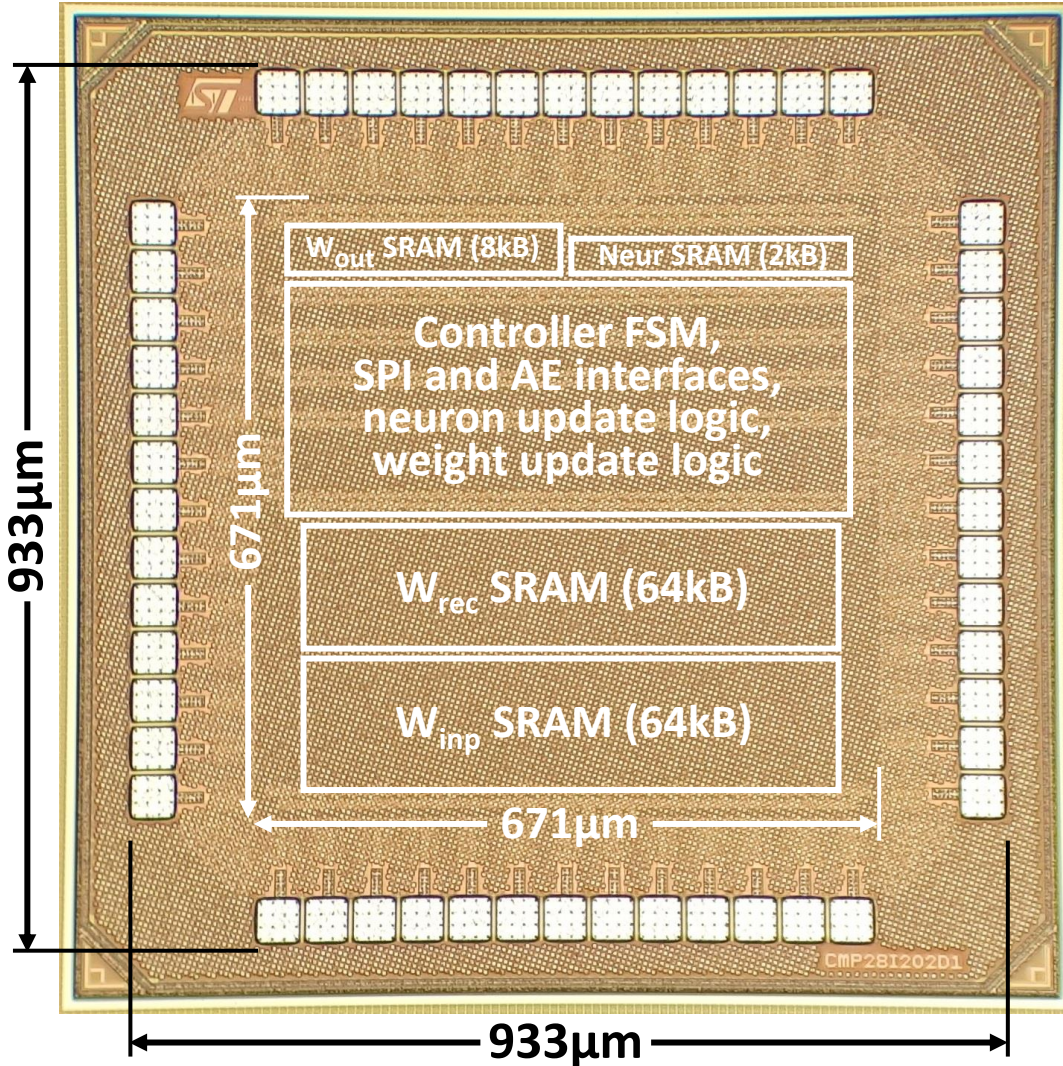


# ReckOn – Neuroscience and AI meet in efficient hardware

*Chip microphotograph and summary*



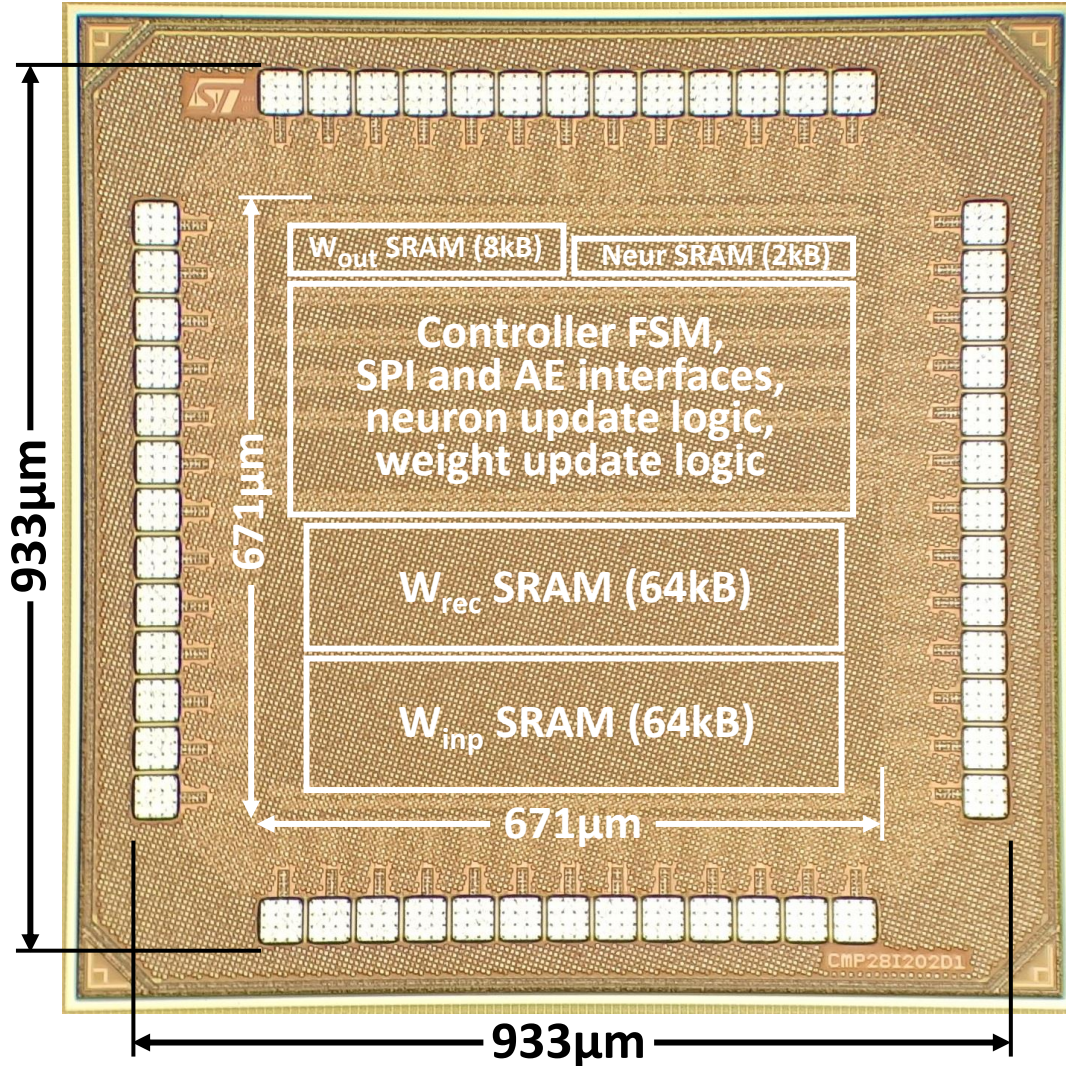
open source hardware



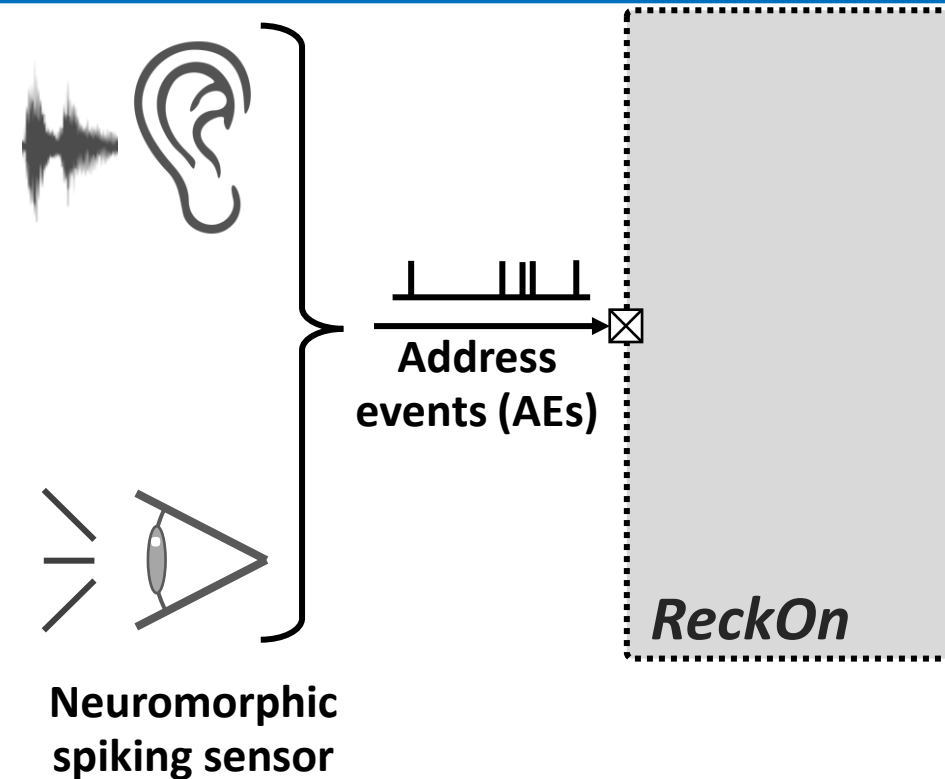
Technology	28nm FDSOI CMOS	
Core size	0.67 x 0.67 mm <sup>2</sup>	0.45mm <sup>2</sup>
Die size	0.93 x 0.93 mm <sup>2</sup>	
SRAM	138kB	+ 0kB ext. DRAM!
Network	Spiking RNN	
Training timespan	Max. 32k steps	

# ReckOn – Neuroscience and AI meet in efficient hardware

*Chip microphotograph and summary*



- Event-driven / sparsity-aware computation
- Sensor-agnostic raw-data processing
- **Task-agnostic processing and learning**



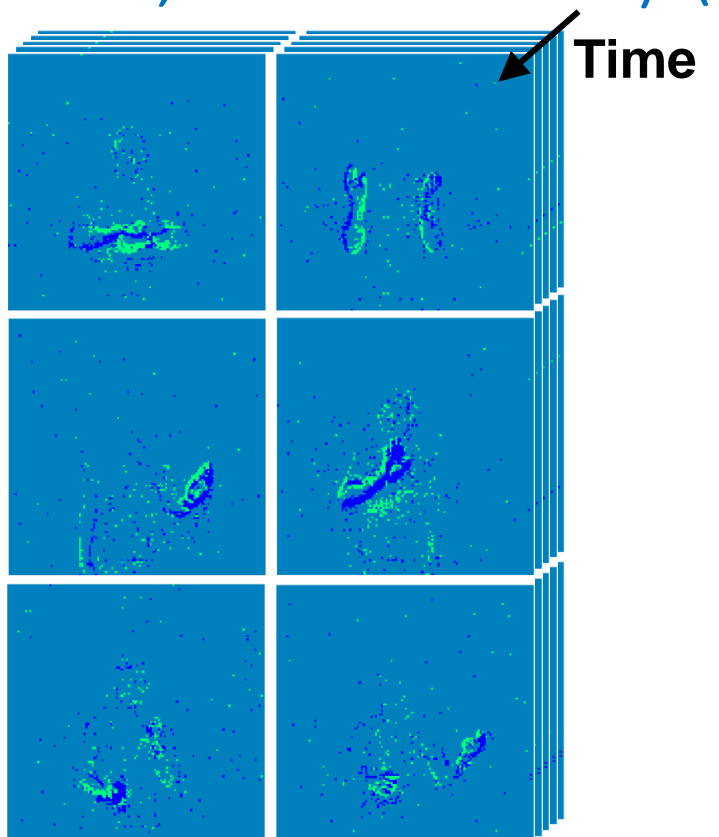
# ReckOn – Measurements and benchmarking

*Three benchmarks that demonstrate task-agnostic learning*



## Vision

IBM DVS Gestures dataset  
(10 classes, shrunk to 16x16)

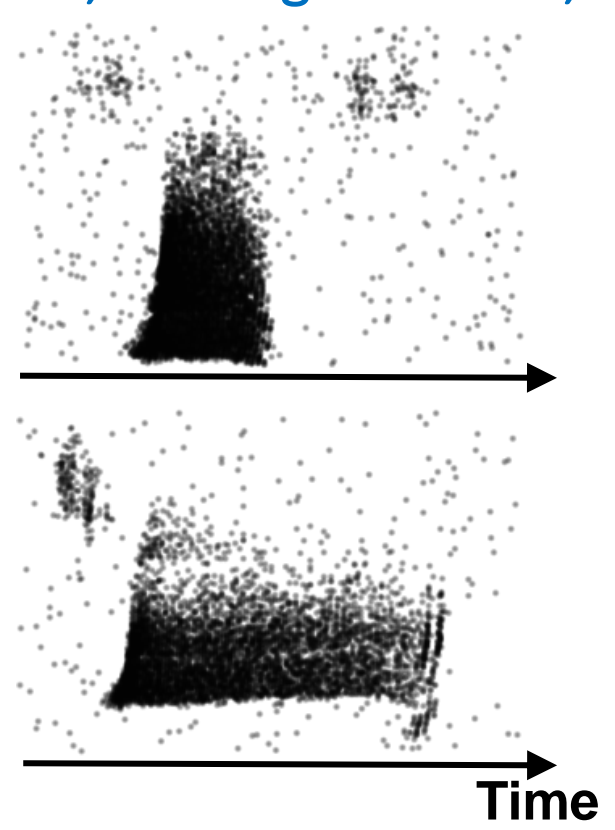


**Accuracy: 87.3% (28μW @0.5V)**



## Audition

Spiking Heidelberg Digits (EN) dataset  
(1-word KWS, 1:1 target vs. filler, 1:3 sub)

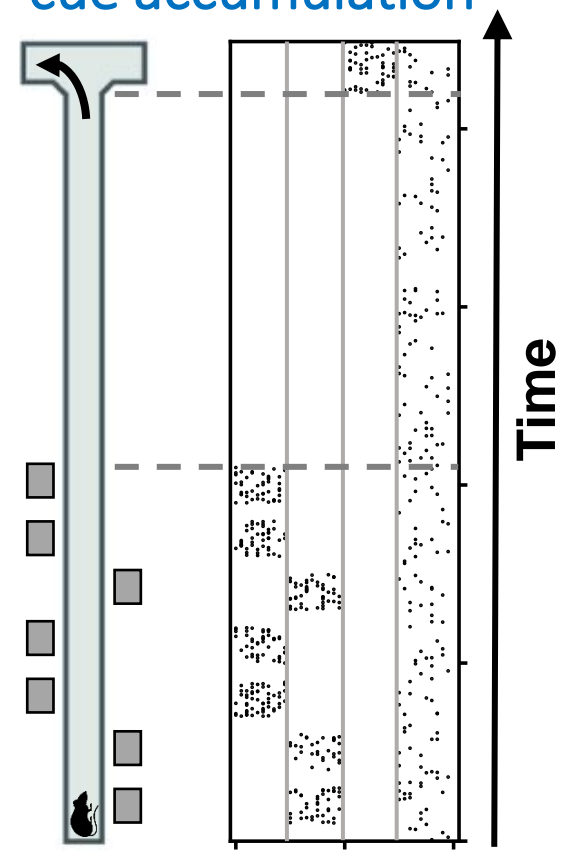


**Accuracy: 90.7% (46μW @0.5V)**



## Navigation

Delayed-supervision  
cue accumulation



**Accuracy: 96.4% (14μW @0.5V)**

# ReckOn – What you should remember

*Key elements toward neuromorphic edge intelligence*

ReckOn merges AI, neuroscience and hardware to

Neuromorphic intelligence outlines an exciting future for tiny on-device learning!

# The *Cognitive Sensor Nodes and Systems* (CogSys) Team

We bridge the bottom-up (bio-inspired) and top-down (engineering-driven) design approaches toward neuromorphic intelligence.

Things we like (non-exhaustive list!):

- designing neuromorphic ICs and tinyML accelerators (mostly digital, going mixed-signal)
- bio-plausible training algorithms and synaptic plasticity mechanisms
- system-level optimization for autonomous sensorimotor agents, from sensing to decision

Positions will open soon!



# Questions?



@C\_Frenkel



cfrenkel



ChFrenkel



Charlotte-Frenkel



c.frenkel@tudelft.nl



chfrenkel.github.io

## Main references:

- ODIN: [C. Frenkel et al., “A 0.086-mm<sup>2</sup> 12.7-pJ/SOP 64k-synapse 256-neuron online-learning digital spiking neuromorphic processor in 28nm CMOS,” *IEEE Trans. BioCAS*, 2019]
- DRTP: [C. Frenkel, M. Lefebvre et al., “Learning without feedback: Fixed random learning signals allow for feedforward training of deep neural networks,” *Frontiers in Neuroscience*, 2021]
- SPOON: [C. Frenkel et al., “A 28-nm convolutional neuromorphic processor enabling online learning with spike-based retinas,” *IEEE ISCAS*, 2020]
- **Review:** [C. Frenkel, D. Bol and G. Indiveri, “Bottom-up and top-down approaches for the design of neuromorphic processing systems: Tradeoffs and synergies between natural and artificial intelligence,” *Proceedings of the IEEE* (to appear), 2023]
- **ReckOn:** [C. Frenkel and G. Indiveri, “ReckOn: A 28-nm Sub-mm<sup>2</sup> Task-Agnostic Spiking Recurrent Neural Network Processor Enabling On-Chip Learning over Second-Long Timescales,” *IEEE International Solid-State Circuits Conference (ISSCC)*, 2022]

*Open-sourced!*

[github.com/ChFrenkel/ODIN](https://github.com/ChFrenkel/ODIN)

*Open-sourced!*

[github.com/ChFrenkel/DirectRandomTargetProjection](https://github.com/ChFrenkel/DirectRandomTargetProjection)

*Already available on*

[arxiv.org/abs/2106.01288](https://arxiv.org/abs/2106.01288)

*Open-sourced!*

[github.com/ChFrenkel/ReckOn](https://github.com/ChFrenkel/ReckOn)



On device learning Forum

# Copyright Notice

This multimedia file is copyright © 2023 by tinyML Foundation. All rights reserved. It may not be duplicated or distributed in any form without prior written approval.

tinyML<sup>®</sup> is a registered trademark of the tinyML Foundation.

[www.tinyml.org](http://www.tinyml.org)



On device learning Forum

# Copyright Notice

This presentation in this publication was presented as a tinyML® Talks webcast. The content reflects the opinion of the author(s) and their respective companies. The inclusion of presentations in this publication does not constitute an endorsement by tinyML Foundation or the sponsors.

There is no copyright protection claimed by this publication. However, each presentation is the work of the authors and their respective companies and may contain copyrighted material. As such, it is strongly encouraged that any use reflect proper acknowledgement to the appropriate source. Any questions regarding the use of any materials presented should be directed to the author(s) or their companies.

tinyML is a registered trademark of the tinyML Foundation.

**[www.tinyml.org](http://www.tinyml.org)**