

# tinyML<sup>®</sup> EMEA

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

## tinyML EMEA Technical Forum 2021 Proceedings

June 7 – 10, 2021

Virtual Event



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

T I N Y



# Bottom-Up and Top-Down Neural Processing Systems Design:

## *Unveiling the Road toward Neuromorphic Intelligence*

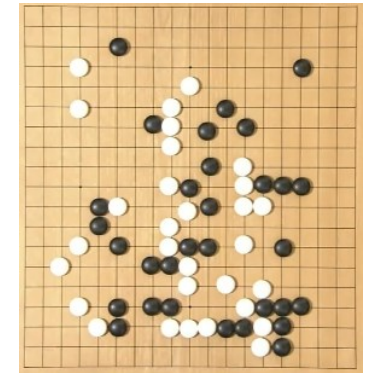
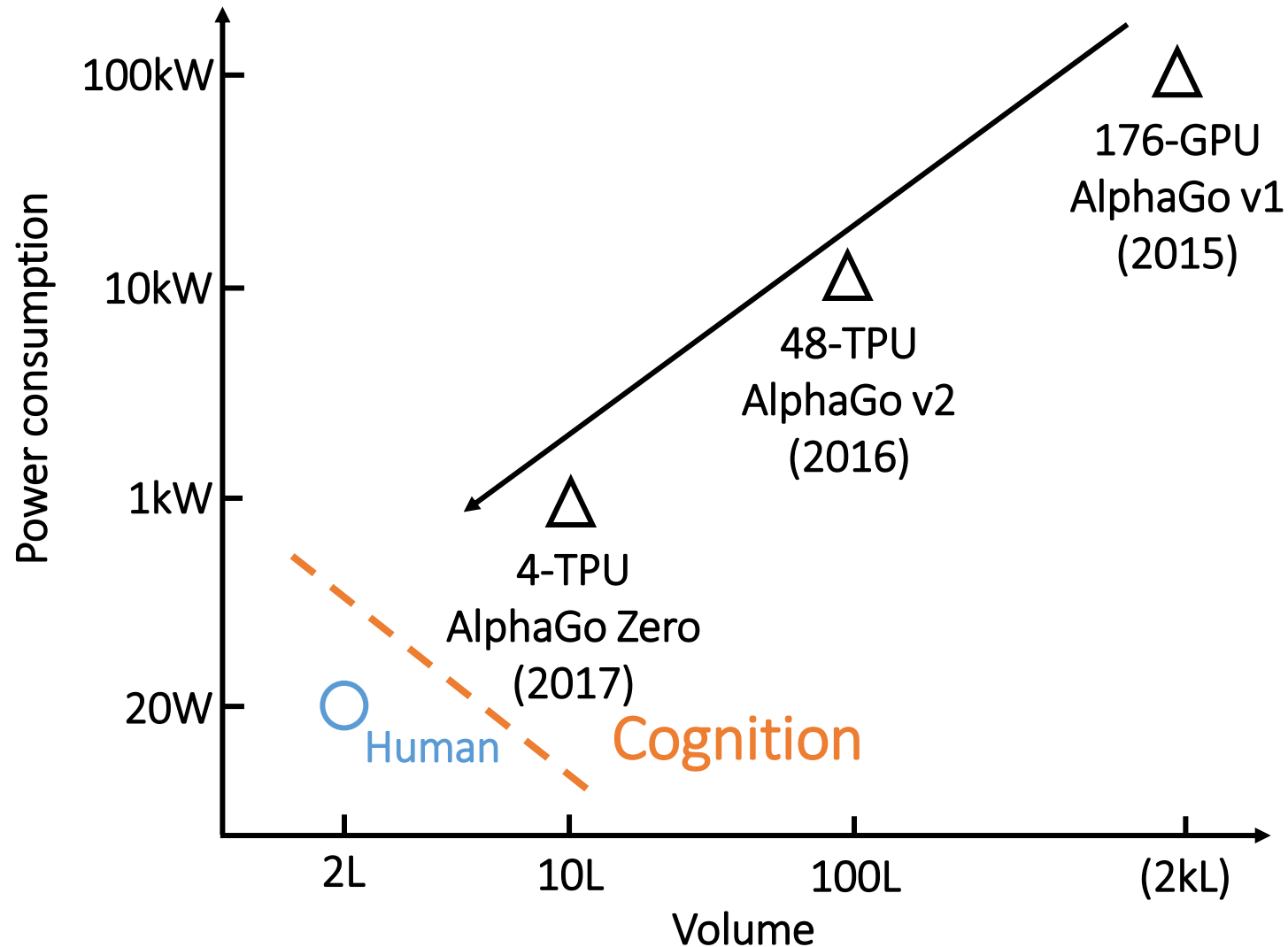
Charlotte Frenkel

Institute of Neuroinformatics, UZH and ETH Zürich, Switzerland  
charlotte@ini.uzh.ch

tinyML EMEA Technical Forum  
Online, June 7-10, 2021

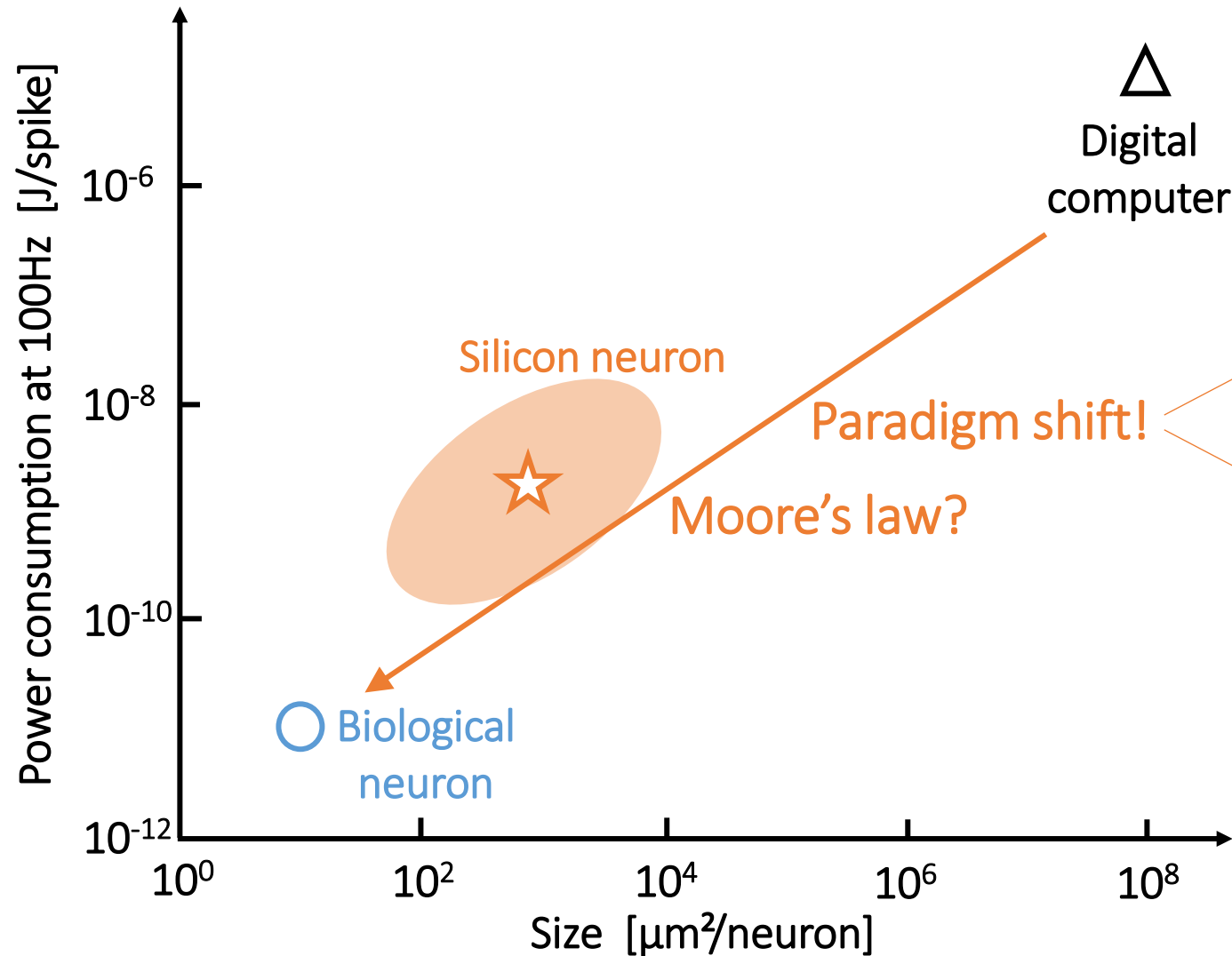
# Neuromorphic Engineering – Why?

*Efficiency of artificial intelligence vs. **natural intelligence**?*



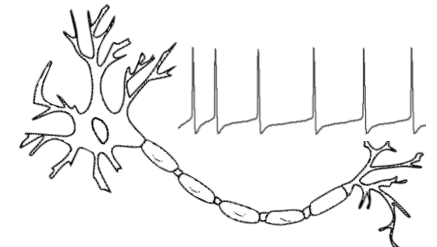
# Neuromorphic Engineering – Why?

*Efficiency of bio-inspired neuromorphic computing?*



Data representation:  
sparse, event-driven spike trains

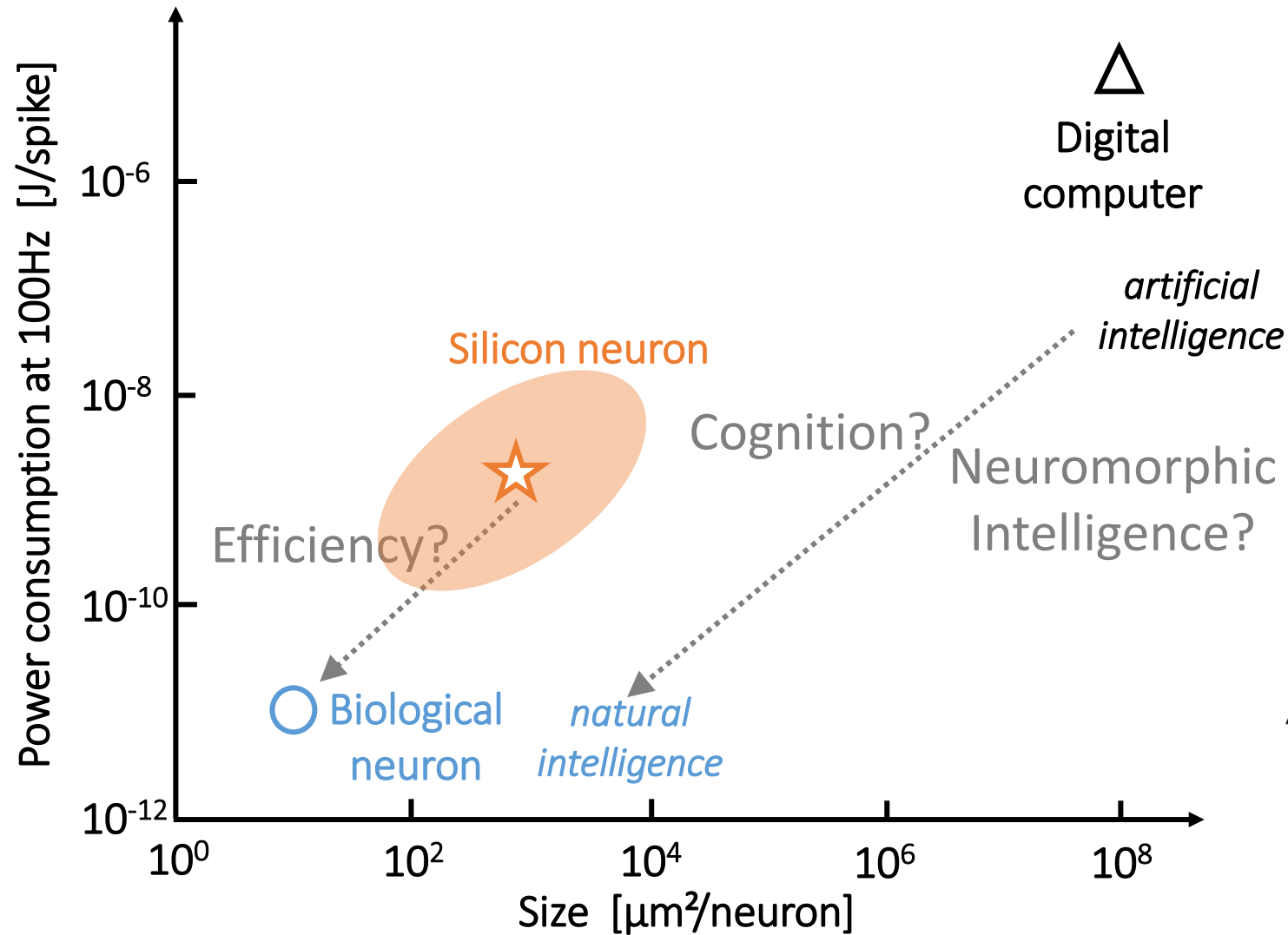
Architecture: distributed processing with co-located neurons and synapses



[Poon & Zhou, *Front. Neurosci.*, 2011]

# Neuromorphic Engineering – How?

*A design strategy toward efficiency and cognition?*



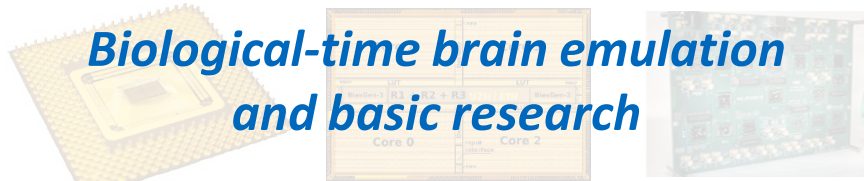
[Poon & Zhou, *Front. Neurosci.*, 2011]

# Neuromorphic Engineering – How?

*A design strategy toward efficiency and cognition?*

*Subthreshold analog (mixed-signal)*

**Biological-time brain emulation  
and basic research**



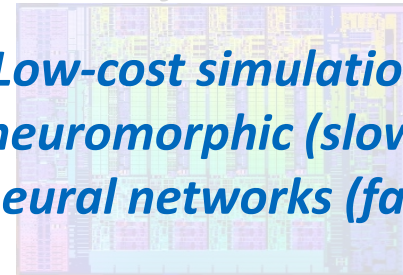
ROLLS  
(UZH/ETHZ)

DYNAPs  
(UZH/ETHZ)

NeuroGrid  
(Stanford)

*Software*

**Low-cost simulation:  
neuromorphic (slow),  
neural networks (fast)**



CPU / GPU

*Dedicated/distributed sim.*

**Simulation acceleration  
for neuroscience  
and neural networks**



FPGA

SpiNNaker 1/2  
(Manchester, TUD)

*Above-threshold analog (mixed-signal)*

**Neuroscience  
simulation acceleration**



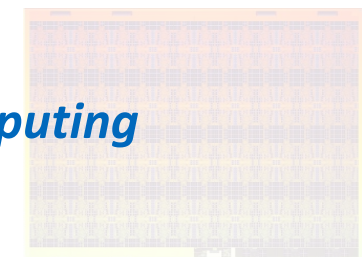
BrainScaleS 1/2 (Heidelberg)

*Large-scale full-custom digital designs*

**Cognitive computing**



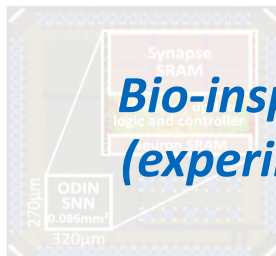
TrueNorth (IBM)



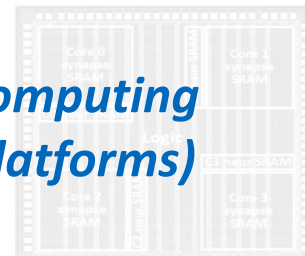
Loihi (Intel)

*Small-scale full-custom digital designs*

**Bio-inspired edge computing  
(experimentation platforms)**

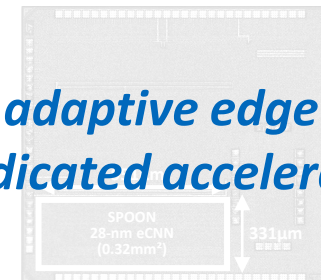


ODIN (UCLouvain)



MorphIC (UCLouvain)

**Low-cost adaptive edge computing  
(dedicated accelerators)**



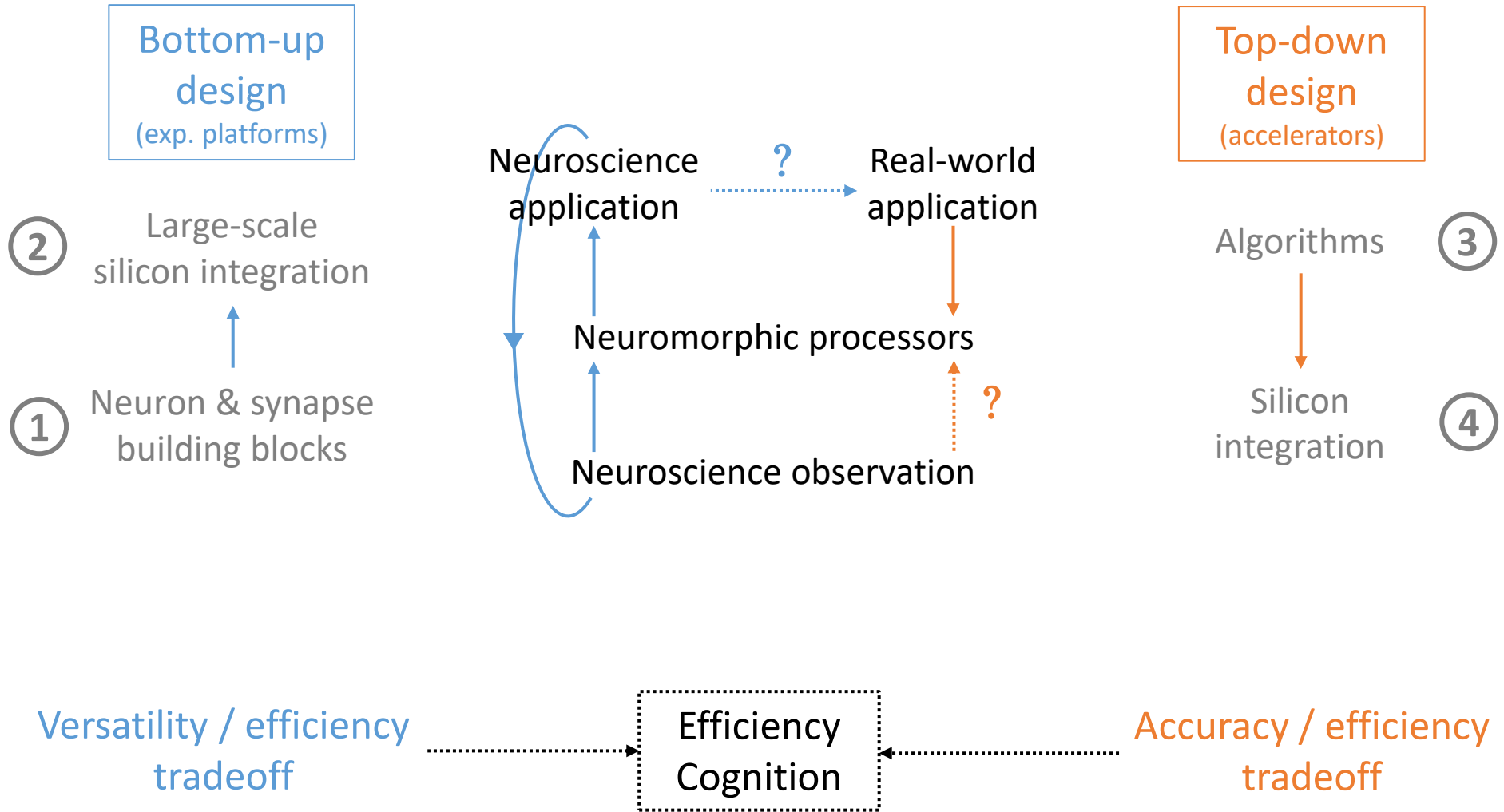
SPOON (UCLouvain)

See also:  
[Seo, CICC'11]  
[Knag, JSSC'15]  
[Park, ISSCC'19]



# Neuromorphic Engineering – How?

*Unveiling roads to embedded cognition*



# Outline

## Part I – Bottom-up neuromorphic design

- Building blocks
- Integration

## Part II – Top-down neuromorphic design

- Algorithms
- Integration

## Conclusion and perspectives



# Outline

## Part I – Bottom-up neuromorphic design

- Building blocks

Neurons and synapses as adaptive processing and memory elements

[Frenkel, *ISCAS*, 2017]

[Frenkel, *BioCAS*, 2017]

- Integration

## Part II – Top-down neuromorphic design

- Algorithms

- Integration

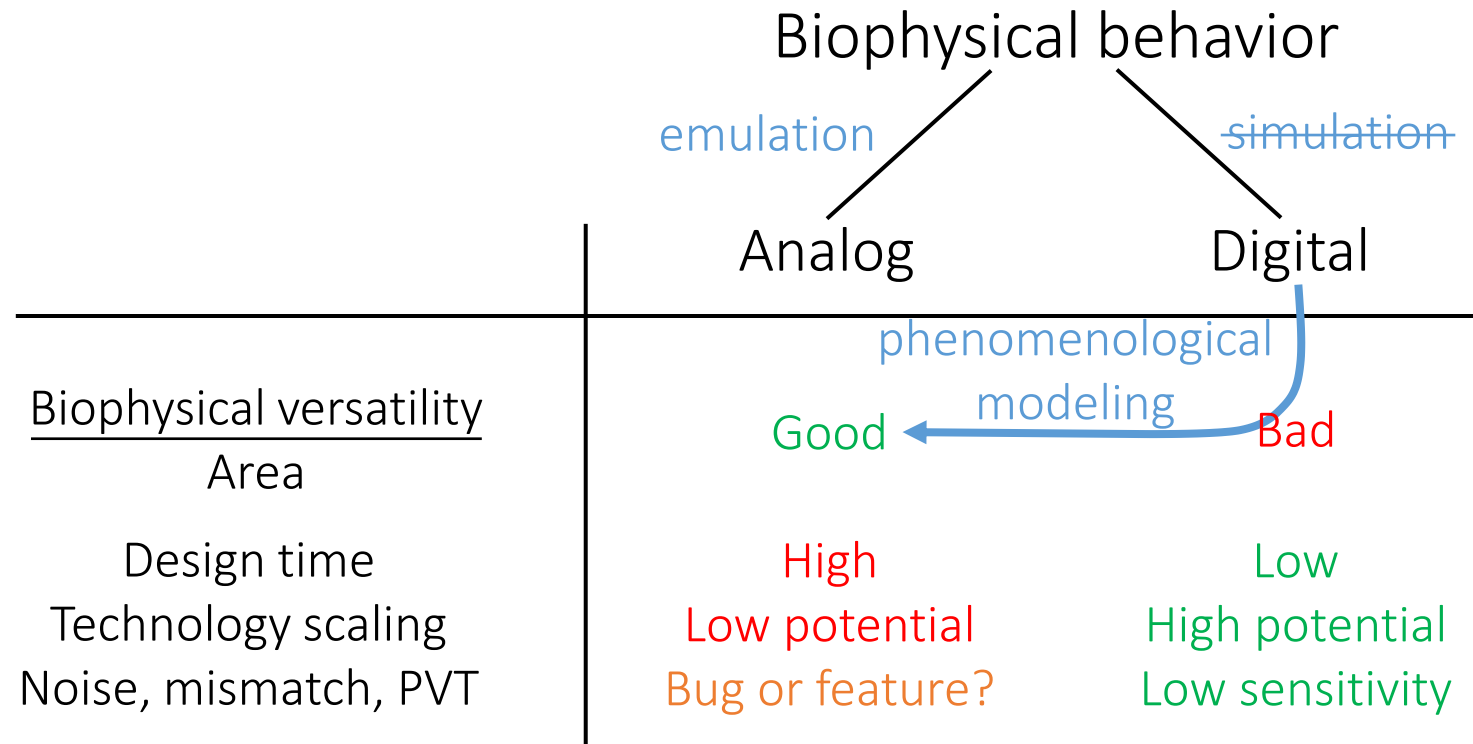
## Conclusion and perspectives

# Design strategy

*Analog or digital?*



We'll come back  
to this. 😊



How can we make the best of both worlds?

# Design strategy

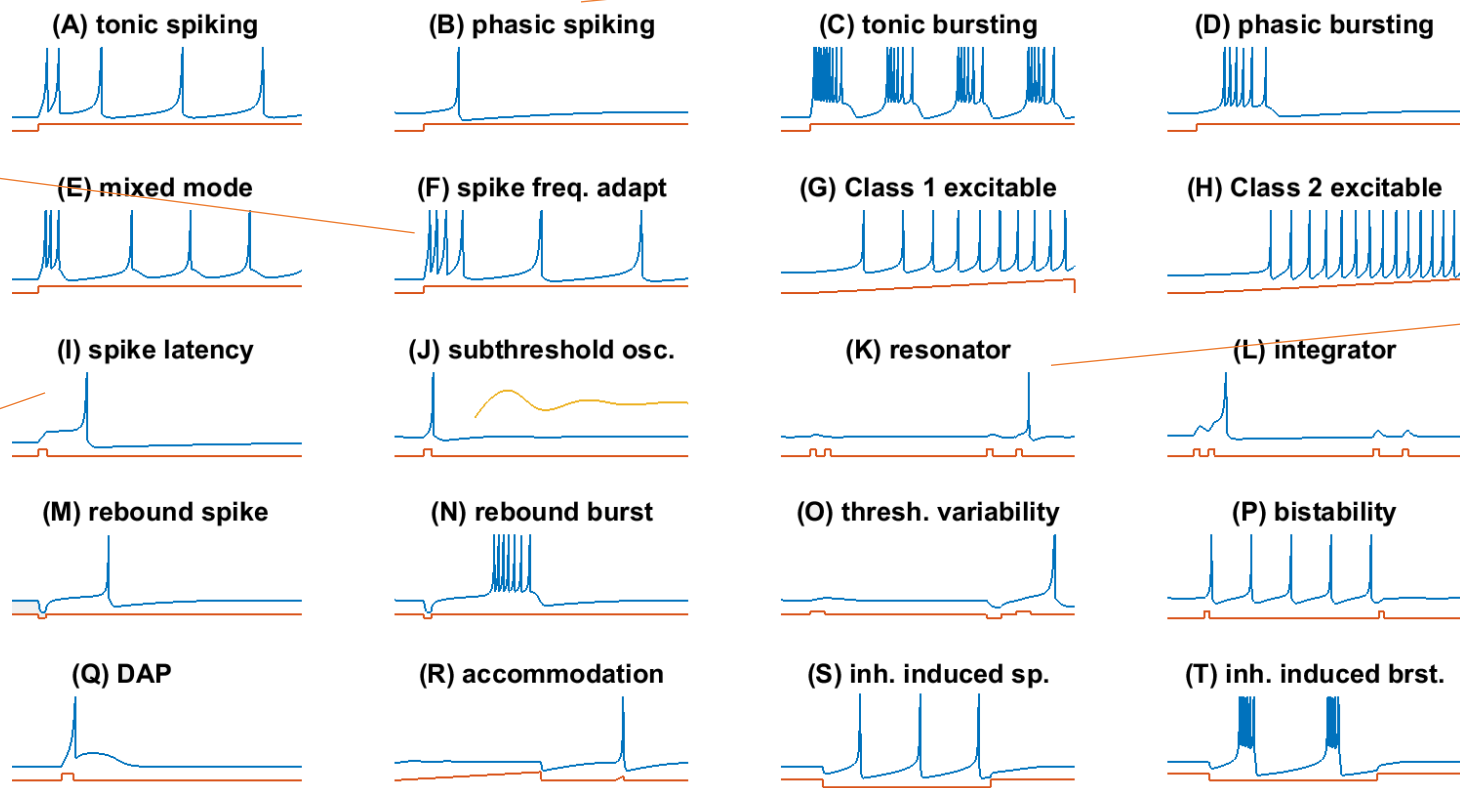
*What should we aim for and phenomenologically implement?*

## Neurons

- 20 Izhikevich behaviors of cortical spiking neurons

Useful for time-to-first-spike encodings

Introduce competition for unsupervised learning in winner-take-all networks  
[Kreiser, BioCAS'17]

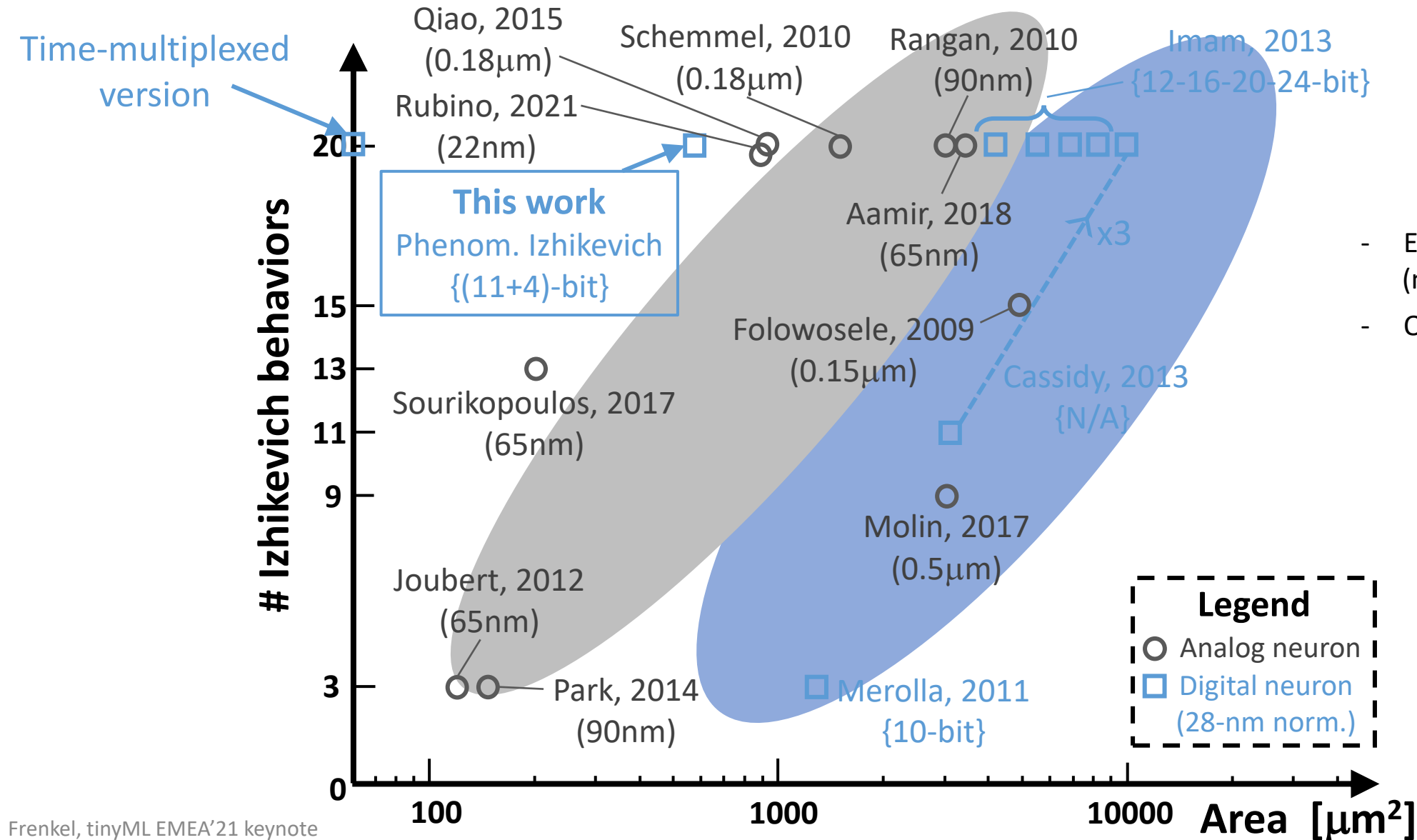


Discrimination of specific frequencies

Stereo sound source localization  
[Schoepe, BioCAS'19]

# Proposed phenomenological digital neuron

*Tackling the versatility/efficiency tradeoff*



## Key features:

- Entirely event-driven (no time-stepped integration)
- Only 4 functions necessary:
  - Threshold adaptation
  - Time window generation
  - Simple template matching
  - Membrane potential sign rotation



→ perspectives

# Design strategy

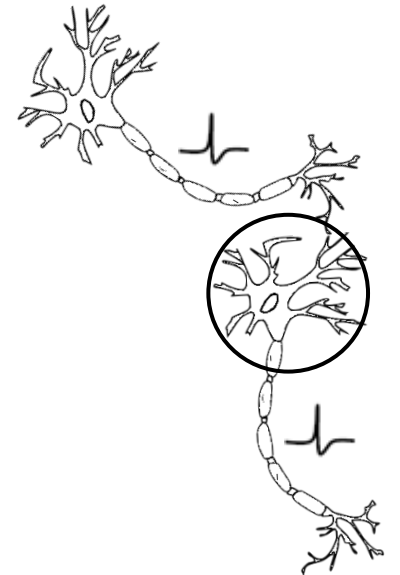
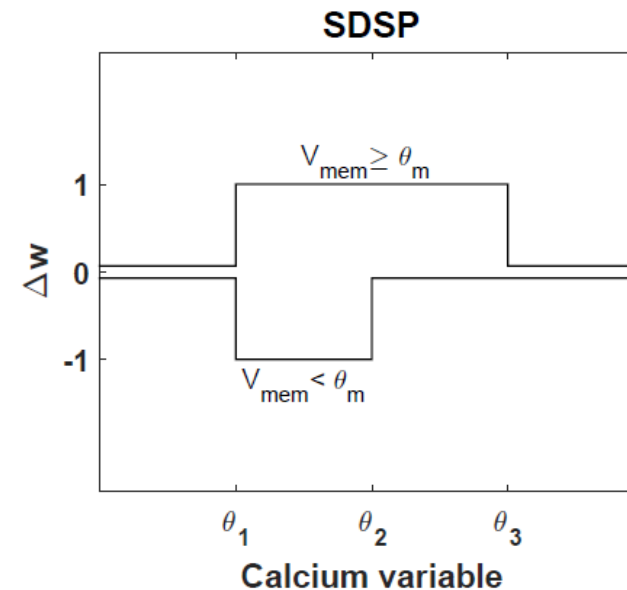
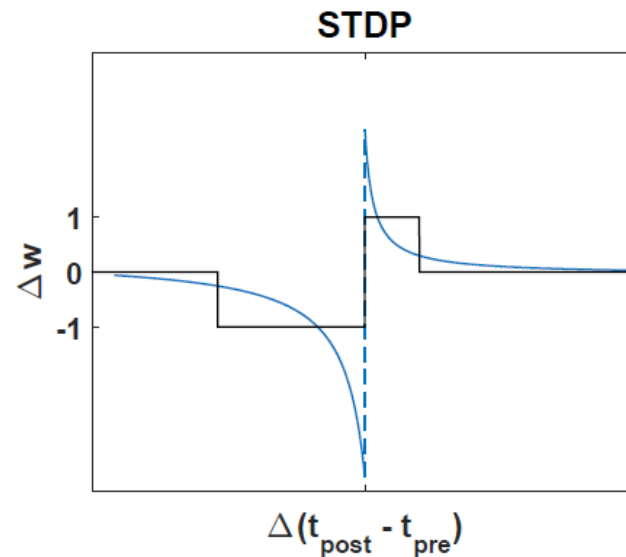
*What should we aim for and phenomenologically implement?*

## Neurons

- 20 Izhikevich behaviors of cortical spiking neurons

## Synapses

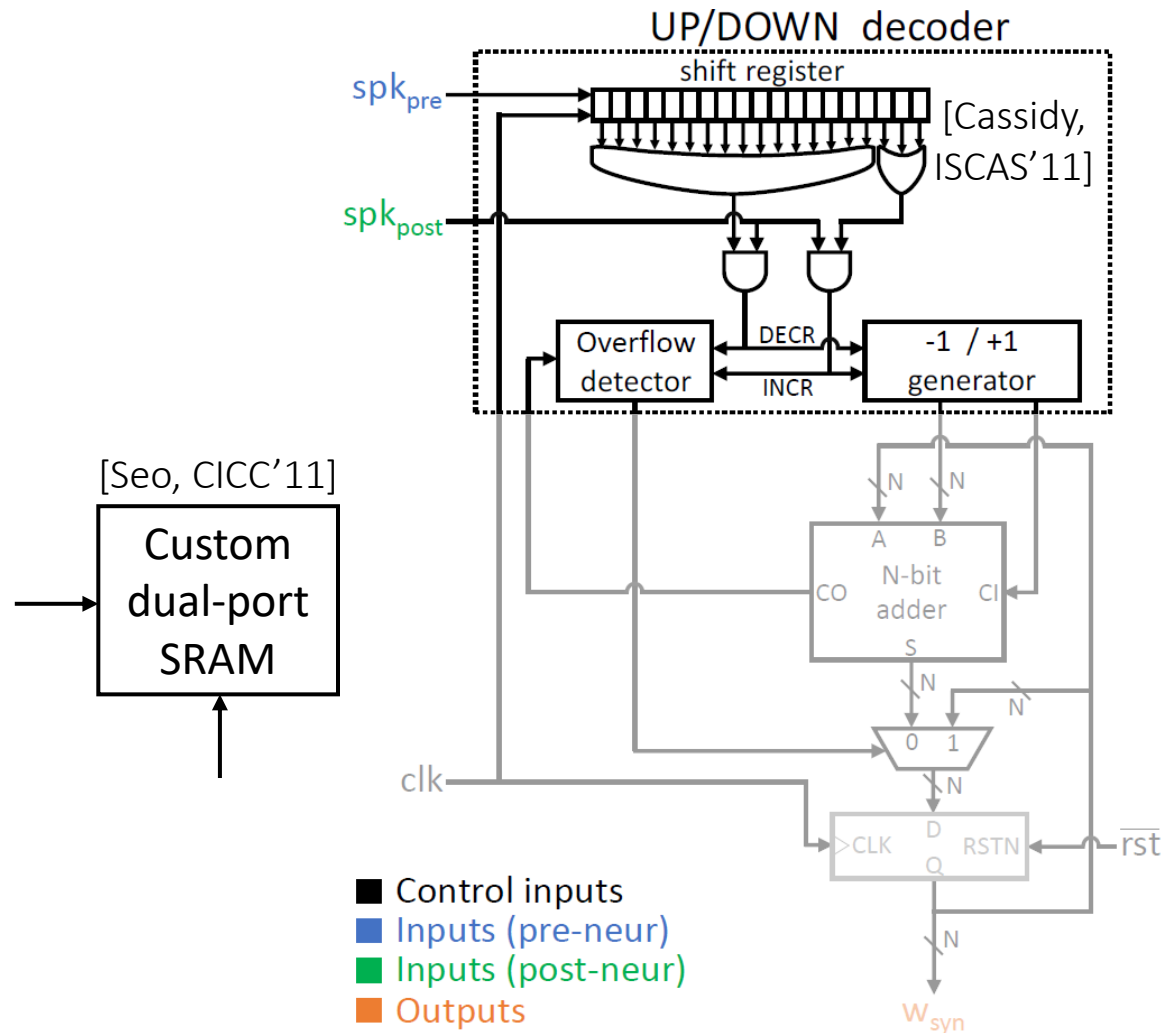
- Spike-based online learning



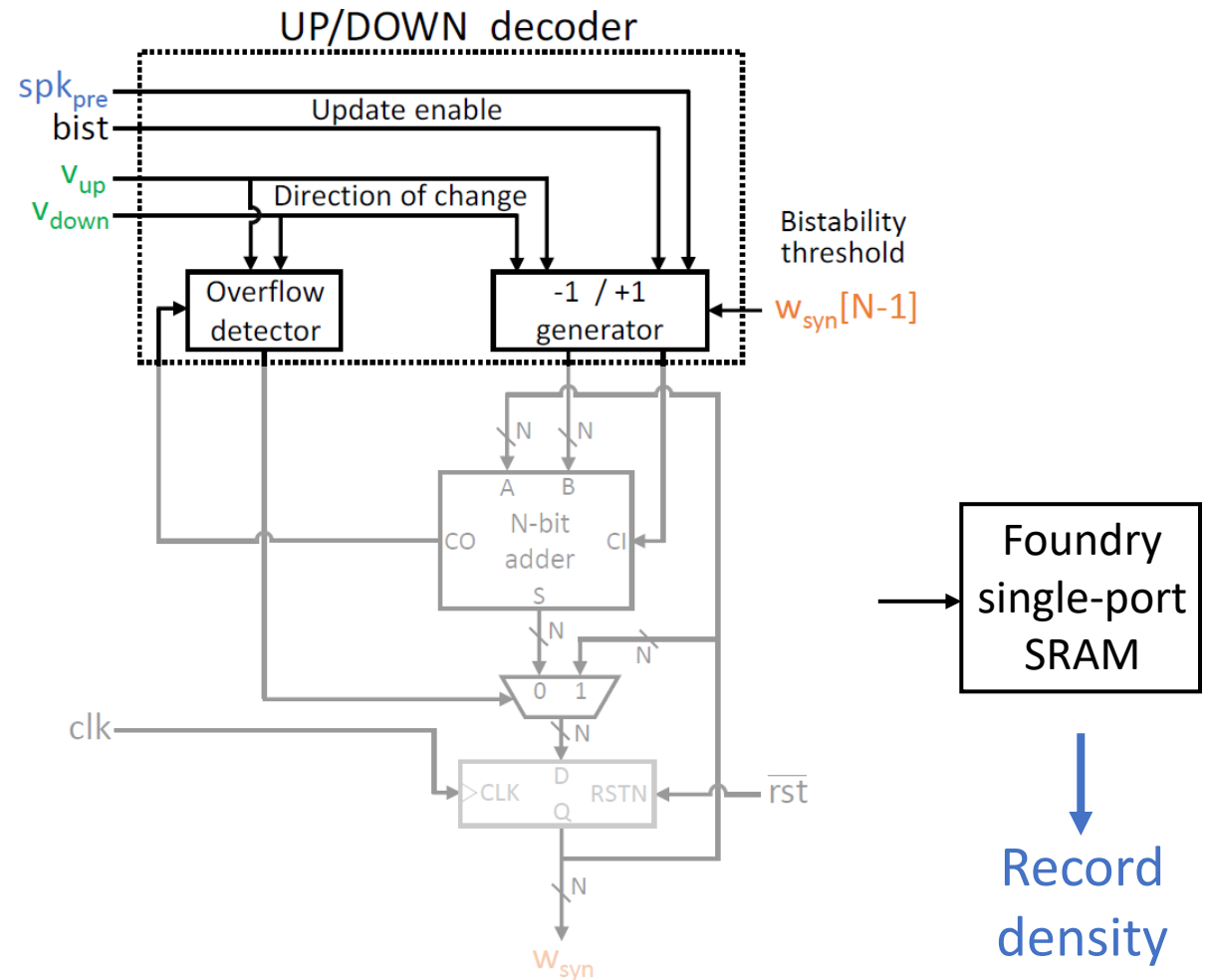
# Proposed digital synapse

*Tackling the versatility/efficiency tradeoff*

Key challenge – Fan-in = 100-10000 synapses/neuron



STDP



SDSP

Record density

# Outline

## Part I – Bottom-up neuromorphic design

- Building blocks
- Integration

Proposed neuromorphic experimentation platforms

[Frenkel, *Trans. BioCAS*, 2019a]

[Frenkel, *Trans. BioCAS*, 2019b]

## Part II – Top-down neuromorphic design

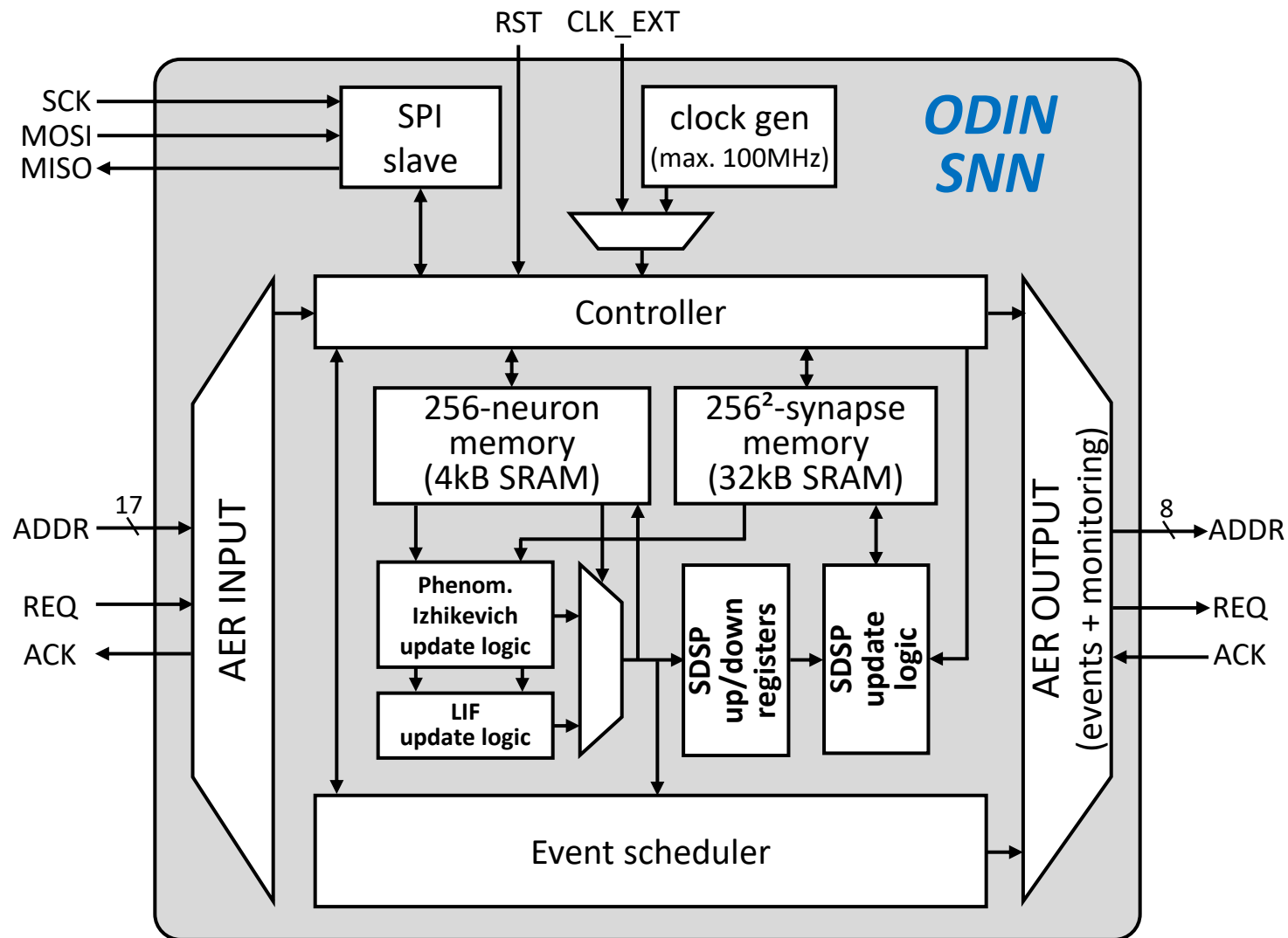
- Algorithms
- Integration

## Conclusion and perspectives

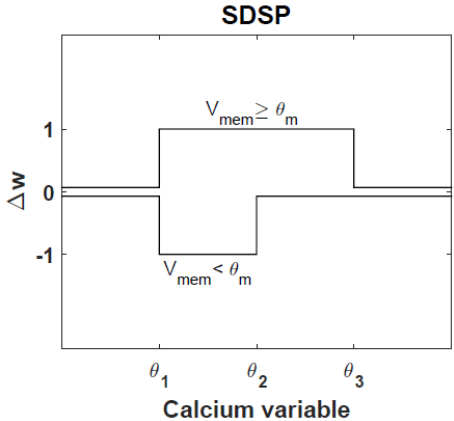
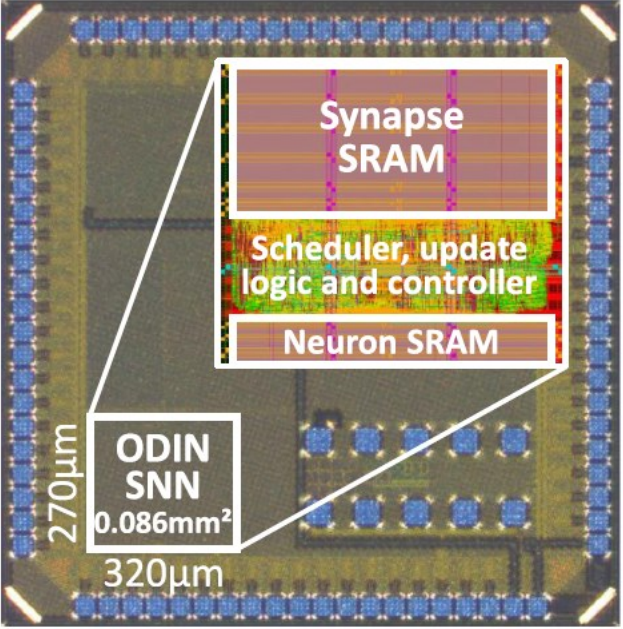


# Architecture of ODIN

ODIN – A 256-neuron 64k-synapse Online-learning Digital Neurosynaptic core



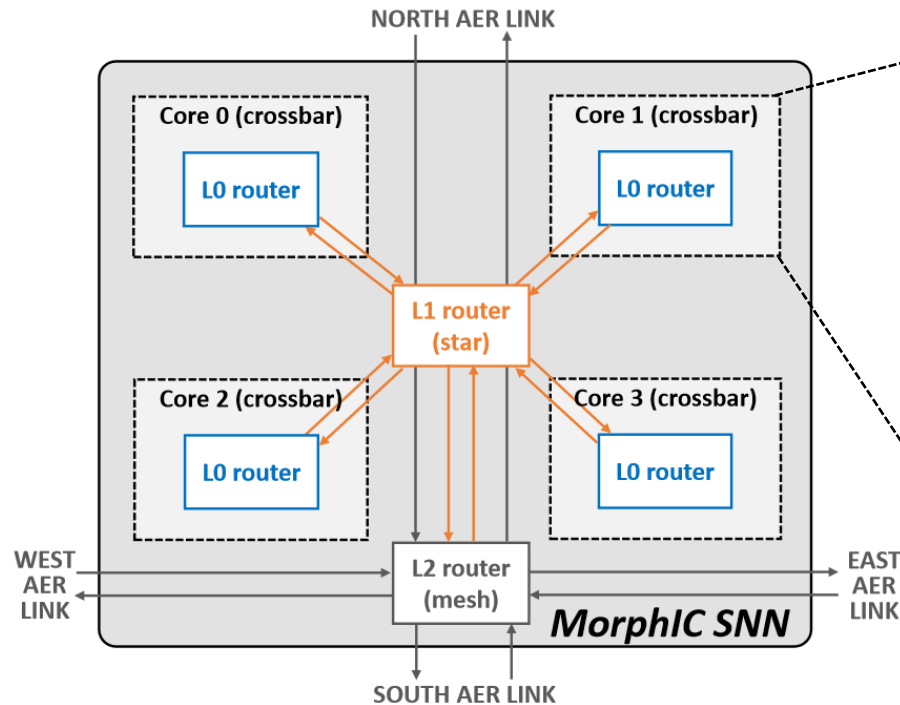
# ODIN – Chip microphotograph and specifications



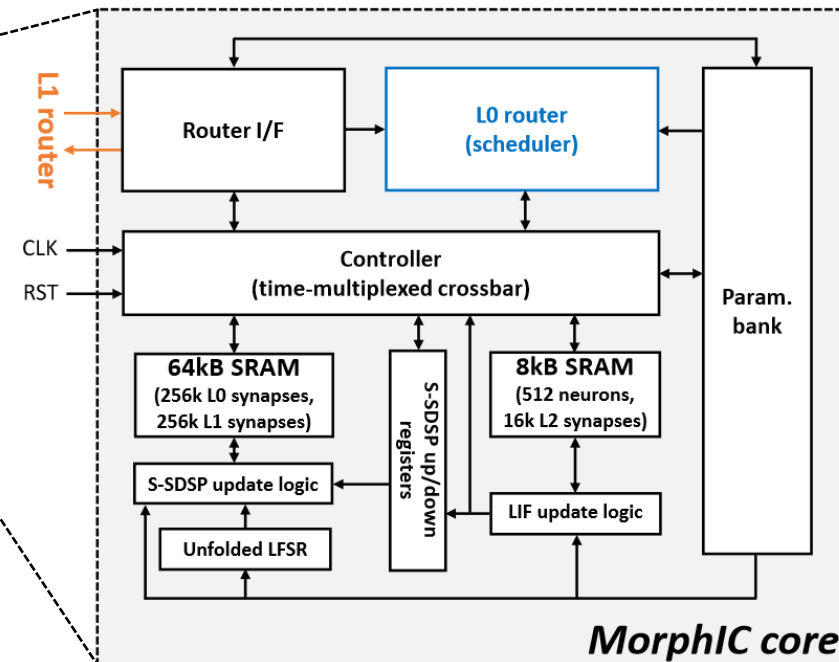
Technology	28nm FDSOI
Implementation	Digital
Area	0.086mm <sup>2</sup>
# neurons	256
# synapses	64k
# Izhikevich behav.	20
Online learning	SDSP, (3+1)-bit weight
Time constant	Biological to accelerated
Supply voltage	0.55V – 1.0V
Leakage power ( $P_{leak}$ )	27.3μW @0.55V
Idle power ( $P_{idle}$ )	1.78μW/MHz @0.55V
Incr. energy/SOP ( $E_{SOP}$ )	8.43pJ @0.55V
Global energy/SOP ( $E_{tot.SOP}$ )	>12.7pJ @0.55V
Routing flexibility/efficiency	☹ (AER)
Fan-in	256
Fan-out	256

# Architecture of MorphIC

Chip-level architecture



Core architecture

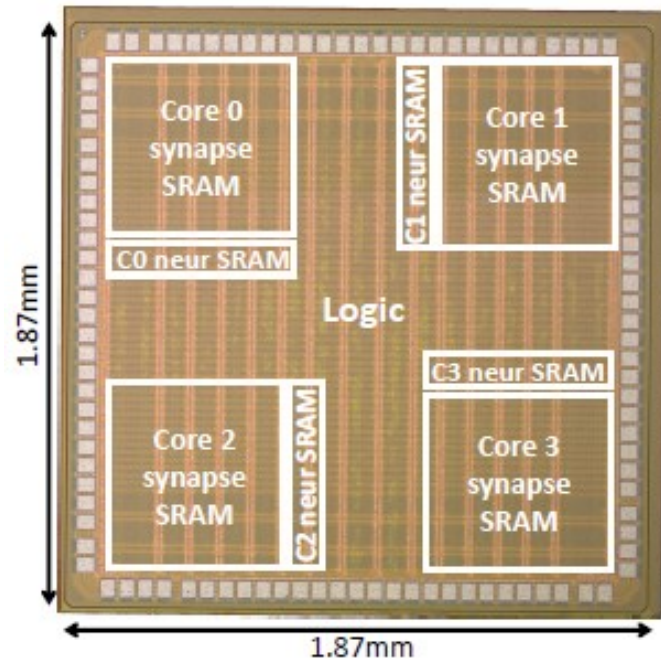


Neurons/core 512  
Synapses/core 528k

Fan-in 1k  
Fan-out 2k

Stochastic SDSP (S-SDSP)  
on binary synapses

# MorphIC – Chip microphotograph and specifications



Technology	65nm LP CMOS
Implementation	Digital
Area	3.5mm <sup>2</sup> (incl. pads) 2.86mm <sup>2</sup> (excl. pads)
Number of cores	4
Total # neurons (type)	2048 (LIF)
Total # synapses (hier.)	1M (L0), 1M (L1), 64k (L2)
Fan-in (hier.)	512 (L0), 512 (L1), 32 (L2)
Fan-out (hier.)	512 (L0), 3x512 (L1), 4 (L2)
Online learning	Stochastic SDSP, 1-bit weight
Time constant	Biological to accelerated
Supply voltage	0.8V – 1.2V
Max. clock frequency	55MHz (0.8V) – 210MHz (1.2V)
Leakage power ( $P_{leak}$ )	45μW @0.8V
Idle power ( $P_{idle}$ )	41.3μW/MHz @0.8V
Energy/SOP ( $E_{SOP}$ )	30pJ @0.8V

# Comparison with SoA experimentation platforms

## Mixed-signal

## Digital

Author	Schemmel	Benjamin	Qiao	Moradi	Painkras	Akopyan	Davies	Frenkel	Frenkel
Publication	ISCAS, 2010	PIEEE, 2014	Front. NS, 2015	TBioCAS, 2017	JSSC, 2013	TCAD, 2015	IEEE Micro, 2018	TBCAS, 2019a	TBCAS, 2019b
Chip name	HICANN	Neurogrid	ROLLS	DYNAPs	SpiNNaker	TrueNorth	Loihi	ODIN	MorphIC
Implementation	Mixed-signal	Mixed-signal	Mixed-signal	Mixed-signal	Digital	Digital	Digital	Digital	Digital
Technology	0.18 $\mu$ m	0.18 $\mu$ m	0.18 $\mu$ m	0.18 $\mu$ m	0.13 $\mu$ m	28nm	14nm FinFET	28nm FDSOI	65nm LP
# cores	1	16	1	4	18	4096	128	1	4
Neurosynaptic core area [mm <sup>2</sup> ]	49	168	51.4	7.5	3.75	0.095	0.4	0.086	0.715
# Izhikevich behaviors	(20)	N/A	(20)	(20)	Programmable	11 (3 neur: 20)	(6)	20	3
# neurons per core	512	64k	256	256	max. 1000	256	max. 1024	256	512
Synaptic weight storage	4-bit (SRAM)	Off-chip	Capacitor	12-bit (CAM)	Off-chip	1-bit (SRAM)	1- to 9-bit (SRAM)	(3+1)-bit (SRAM)	1-bit (SRAM)
Embedded online learning	STDP	No	SDSP	No	Programmable	No	Programmable	SDSP	S-SDSP
# synapses per core	112k	—	128k	16k	—	64k	1M to 114k (1-9 bits)	64k	528k
Time constant	Accelerated	Biological	Biological	Biological	Bio. to accel.	Biological	N/A	Bio. to accel.	Bio. to accel.
Flexibility routing	Medium	Medium	Low	Medium	High	Medium	High	Low	Medium
Flexibility learning	Low	—	Low	Low	—	—	High	Low	Low
Neuron core density [neur/mm <sup>2</sup> ]	raw: 10.5 norm. —	390	5	34	max. 267 max. 5.8k	2.6k 2.6k	max. 2.5k max. 1k	3.0k 3.0k	716 3.9k
Synapse core density [syn/mm <sup>2</sup> ]	raw: 2.3k norm. —	—	2.5k	2.1k	—	674k 674k	2.5M to 282k 1M to 113k	741k 741k	738k 4M
Supply voltage	1.8V	3.0V	1.8V	1.3V-1.8V	1.2V	0.7V-1.05V	0.5V-1.25V	0.55V-1.0V	0.8V-1.2V
Energy per SOP	raw: N/A norm. —	(941pJ) <sup>▲</sup>	>77fJ <sup>▲</sup>	134fJ <sup>▲</sup> /30pJ <sup>▲</sup> (1.3V)	>11.3nJ <sup>▲</sup> /26.6nJ <sup>▲</sup> (>2.4nJ <sup>▲</sup> /5.7nJ <sup>▲</sup> )	26pJ <sup>▲</sup> (0.775V) 26pJ <sup>▲</sup>	>23.6pJ <sup>▲</sup> (0.75V) (66.1pJ <sup>▲</sup> )	8.4pJ <sup>▲</sup> /12.7pJ <sup>▲</sup> (0.55V) 8.4pJ <sup>▲</sup> /12.7pJ <sup>▲</sup>	30pJ <sup>▲</sup> /51pJ <sup>▲</sup> (0.8V) 12.9pJ <sup>▲</sup> /22pJ <sup>▲</sup>

Most direct comparison: IBM TrueNorth core vs. ODIN (same technology node, same number of neurons and synapses per neurosynaptic core, same area).

	ODIN	TrueNorth
Synapses	✓ 4-bit <b>with</b> learning	1-bit <b>without</b> learning ✗
Neurons	✓ 20 Izh. beh.	11 Izh. beh. ✗
Energy/SOP	✓ 12.7pJ @0.55V	26pJ @0.775V ✗
Connectivity	✗ AER	large-scale mesh ✓ → MorphIC

# Comparison with SoA experimentation platforms

## Mixed-signal

## Digital

Author	Schemmel	Benjamin	Qiao	Moradi	Painkras	Akopyan	Davies	Frenkel	Frenkel
Publication	ISCAS, 2010	PIEEE, 2014	Front. NS, 2015	TBioCAS, 2017	JSSC, 2013	TCAD, 2015	IEEE Micro, 2018	TBCAS, 2019a	TBCAS, 2019b
Chip name	HICANN	Neurogrid	ROLLS	DYNAPs	SpiNNaker	TrueNorth	Loihi	ODIN	MorphIC
Implementation	Mixed-signal	Mixed-signal	Mixed-signal	Mixed-signal	Digital	Digital	Digital	Digital	Digital
Technology	0.18 $\mu$ m	0.18 $\mu$ m	0.18 $\mu$ m	0.18 $\mu$ m	0.13 $\mu$ m	28nm	14nm FinFET	28nm FDSOI	65nm LP
# cores	1	16	1	4	18	4096	128	1	4
Neurosynaptic core area [mm <sup>2</sup> ]	49	168	51.4	7.5	3.75	0.095	0.4	0.086	0.715
# Izhikevich behaviors	(20)	N/A	(20)	(20)	Programmable	11 (3 neur: 20)	(6)	20	3
# neurons per core	512	64k	256	256	max. 1000	256	max. 1024	256	512
Synaptic weight storage	4-bit (SRAM)	Off-chip	Capacitor	12-bit (CAM)	Off-chip	1-bit (SRAM)	1- to 9-bit (SRAM)	(3+1)-bit (SRAM)	1-bit (SRAM)
Embedded online learning	STDP	No	SDSP	No	Programmable	No	Programmable	SDSP	S-SDSP
# synapses per core	112k	–	128k	16k	–	64k	1M to 114k (1-9 bits)	64k	528k
Time constant	Accelerated	Biological	Biological	Biological	Bio. to accel.	Biological	N/A	Bio. to accel.	Bio. to accel.
Flexibility routing	Medium	Medium	Low	Medium	High	Medium	High	Low	Medium
Flexibility learning	Low	–	Low	Low	–	–	High	Low	Low
Neuron core density [neur/mm <sup>2</sup> ]	raw 10.5 norm. –	390	5	34	max. 267 max. 5.8k	2.6k 2.6k	max. 2.5k max. 1k	3.0k 3.0k	716 3.9k
Synapse core density [syn/mm <sup>2</sup> ]	raw 2.3k norm. –	–	2.5k	2.1k	–	674k 674k	2.5M to 282k 1M to 113k	741k 741k	738k 4M
Supply voltage	1.8V	3.0V	1.8V	1.3V-1.8V	1.2V	0.7V-1.05V	0.5V-1.25V	0.55V-1.0V	0.8V-1.2V
Energy per SOP raw	N/A	(941pJ) <sup>▲</sup>	>77fJ <sup>▲</sup>	134fJ <sup>▲</sup> /30pJ <sup>▲</sup> (1.3V)	>11.3nJ <sup>▲</sup> /26.6nJ <sup>▲</sup>	26pJ <sup>▲</sup> (0.775V)	>23.6pJ <sup>▲</sup> (0.75V)	8.4pJ <sup>▲</sup> /12.7pJ <sup>▲</sup> (0.55V)	30pJ <sup>▲</sup> /51pJ <sup>▲</sup> (0.8V)
Energy per SOP norm.	–	–	–	–	>2.4nJ <sup>▲</sup> /5.7nJ <sup>▲</sup>	26pJ <sup>▲</sup>	(66.1pJ <sup>▲</sup> )	8.4pJ <sup>▲</sup> /12.7pJ <sup>▲</sup>	12.9pJ <sup>▲</sup> /22pJ <sup>▲</sup>

## Area

ODIN and MorphIC have the highest neuron and synapse densities among all SNNs with embedded synaptic weight storage



# Comparison with SoA experimentation platforms

## Mixed-signal

## Digital

Author	Schemmel	Benjamin	Qiao	Moradi	Painkras	Akopyan	Davies	Frenkel	Frenkel
Publication	ISCAS, 2010	PIEEE, 2014	Front. NS, 2015	TBioCAS, 2017	JSSC, 2013	TCAD, 2015	IEEE Micro, 2018	TBCAS, 2019a	TBCAS, 2019b
Chip name	HICANN	Neurogrid	ROLLS	DYNAPs	SpiNNaker	TrueNorth	Loihi	ODIN	MorphIC
Implementation	Mixed-signal	Mixed-signal	Mixed-signal	Mixed-signal	Digital	Digital	Digital	Digital	Digital
Technology	0.18 $\mu$ m	0.18 $\mu$ m	0.18 $\mu$ m	0.18 $\mu$ m	0.13 $\mu$ m	28nm	14nm FinFET	28nm FDSOI	65nm LP
# cores	1	16	1	4	18	4096	128	1	4
Neurosynaptic core area [mm <sup>2</sup> ]	49	168	51.4	7.5	3.75	0.095	0.4	0.086	0.715
# Izhikevich behaviors	(20)	N/A	(20)	(20)	Programmable	11 (3 neur: 20)	(6)	20	3
# neurons per core	512	64k	256	256	max. 1000	256	max. 1024	256	512
Synaptic weight storage	4-bit (SRAM)	Off-chip	Capacitor	12-bit (CAM)	Off-chip	1-bit (SRAM)	1- to 9-bit (SRAM)	(3+1)-bit (SRAM)	1-bit (SRAM)
Embedded online learning	STDP	No	SDSP	No	Programmable	No	Programmable	SDSP	S-SDSP
# synapses per core	112k	–	128k	16k	–	64k	1M to 114k (1-9 bits)	64k	528k
Time constant	Accelerated	Biological	Biological	Biological	Bio. to accel.	Biological	N/A	Bio. to accel.	Bio. to accel.
Flexibility routing	Medium	Medium	Low	Medium	High	Medium	High	Low	Medium
Flexibility learning	Low	–	Low	Low	–	–	High	Low	Low
Neuron core density [neur/mm <sup>2</sup> ]	raw 10.5	390	5	34	max. 267	2.6k	max. 2.5k	3.0k	716
	norm. –	–	–	–	max. 5.8k	2.6k	max. 1k	3.0k	3.9k
Synapse core density [syn/mm <sup>2</sup> ]	raw 2.3k	–	2.5k	2.1k	–	674k	2.5M to 282k	741k	738k
	norm. –	–	–	–	–	674k	1M to 113k	741k	4M
Supply voltage	1.8V	3.0V	1.8V	1.3V-1.8V	1.2V	0.7V-1.05V	0.5V-1.25V	0.55V-1.0V	0.8V-1.2V
Energy per SOP	raw N/A	(941pJ) <sup>▲</sup>	>77fJ <sup>▲</sup>	134fJ <sup>▲</sup> /30pJ <sup>▲</sup> (1.3V)	>11.3nJ <sup>▲</sup> /26.6nJ <sup>▲</sup>	26pJ <sup>▲</sup> (0.775V)	>23.6pJ <sup>▲</sup> (0.75V)	8.4pJ <sup>▲</sup> /12.7pJ <sup>▲</sup> (0.55V)	30pJ <sup>▲</sup> /51pJ <sup>▲</sup> (0.8V)
	norm. –	–	–	–	>2.4nJ <sup>▲</sup> /5.7nJ <sup>▲</sup>	26pJ <sup>▲</sup>	(66.1pJ <sup>▲</sup> )	8.4pJ <sup>▲</sup> /12.7pJ <sup>▲</sup>	12.9pJ <sup>▲</sup> /22pJ <sup>▲</sup>

## Power

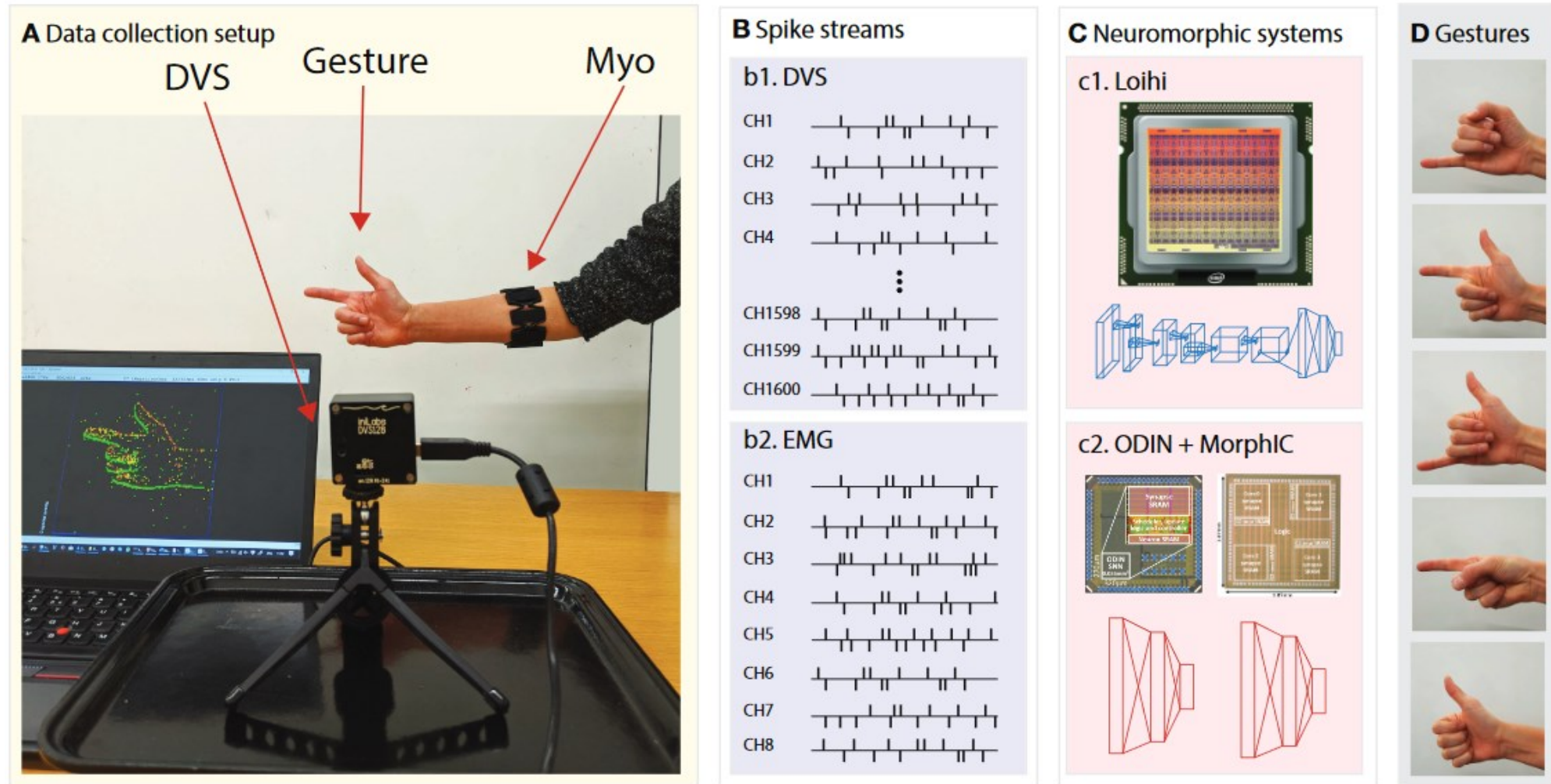
ODIN has the lowest energy per synaptic event among all digital SNNs, MorphIC keeps a competitive energy efficiency.

They outperform subthreshold analog SNNs in accelerated time, but not for biological-time processing.



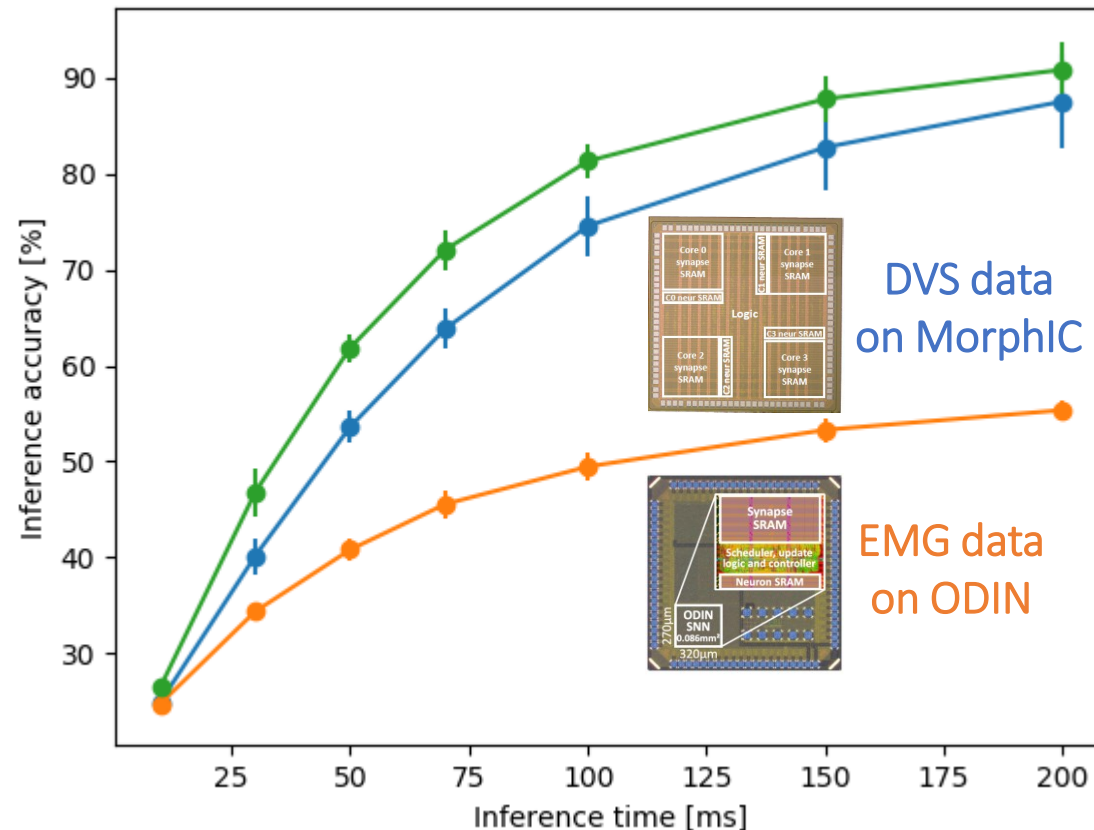
# Results on the spiking EMG/DVS sensor fusion benchmark

[Ceolini, Frenkel, Shrestha et al., *Front. Neurosci.*, 2020]



# Results on the spiking EMG/DVS sensor fusion benchmark

[Ceolini, Frenkel, Shrestha et al., *Front. Neurosci.*, 2020]



Sensor fusion

ODIN+MorphIC 89.4% / 37.4µJ

Loihi 96% / 1105µJ

Software 95.4% / 32100µJ

Accuracy / Energy tradeoff

Neuromorphic designs are more efficient than GPUs, as would be expected from dedicated hardware. But are they more efficient than conventional accelerators?



→ perspectives

See the ODIN and MorphIC papers for more benchmarking, incl. online- and offline-trained MNIST.

# Outline

## Part I – Bottom-up neuromorphic design

- Building blocks
- Integration

## Part II – Top-down neuromorphic design

- Algorithms

Minimizing the training cost of neural networks for adaptive edge computing

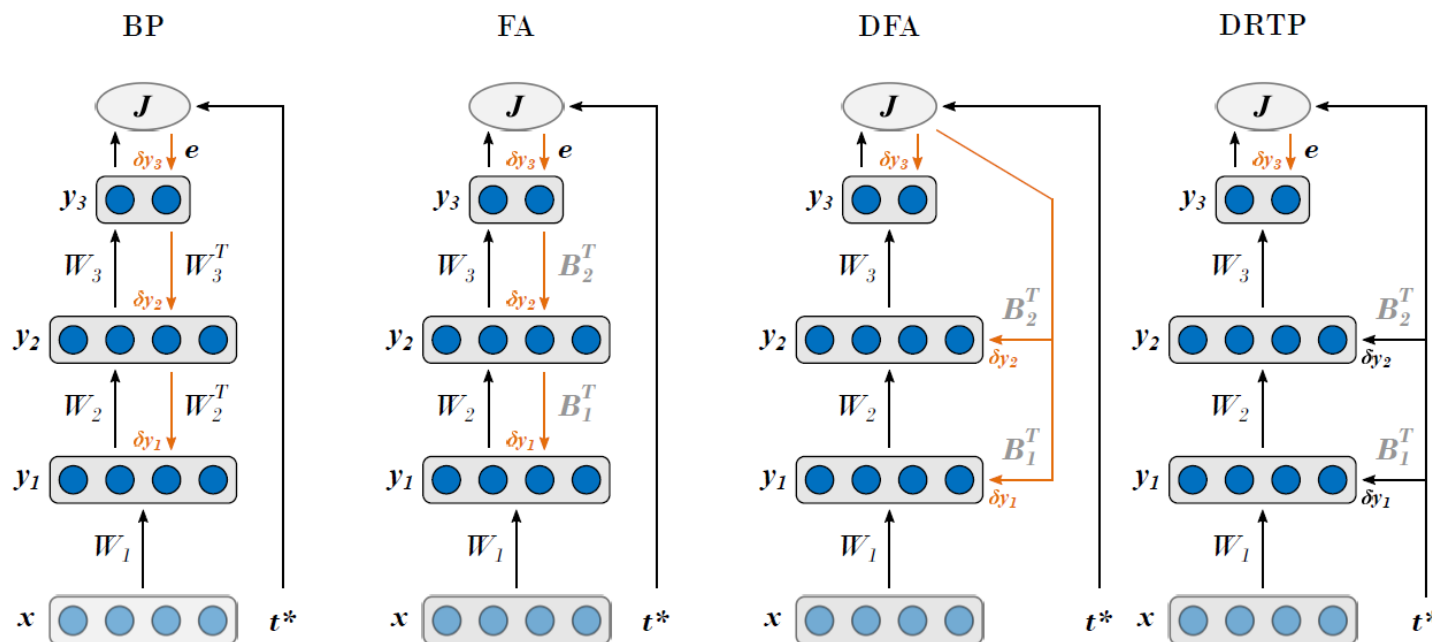
[Frenkel & Lefebvre, *Front. Neurosci.*, 2021]

- Integration

## Conclusion and perspectives

# Learning without feedback

*Releasing the weight transport and update locking of backprop*



	$\delta y_k$	$\frac{\partial J}{\partial y_k} = W_{k+1}^T \delta z_{k+1}$	$B_k^T \delta z_{k+1}$	$B_k^T e$	$B_k^T t^*$
Weight-transport-free	×	×	✓	✓	✓
Update-unlocked	×	×	×	×	✓

Feedforward  
local training

↘ Computational and memory cost ↘

# Direct Random Target Projection (D RTP)

*Ideal use cases?*

## Adaptive edge computing

- Very low power and area overheads can be expected for an on-chip implementation.
- Datasets representative of the complexity associated to autonomous smart sensors: MNIST or CIFAR-10.

→ We'll verify these claims *in silico*.

Disclaimer: whether D RTP scales to ImageNET is probably **not** the right question. 😊

## Neuroscience

D RTP could come in line with recent findings in cortical areas that reveal the existence of output-independent target signals in the dendritic instructive pathways of intermediate-layer neurons.

[Magee & Grienberger,  
*Annual Review of  
Neuroscience*, 2020]

# Outline

## Part I – Bottom-up neuromorphic design

- Building blocks
- Integration

## Part II – Top-down neuromorphic design

- Algorithms
- Integration

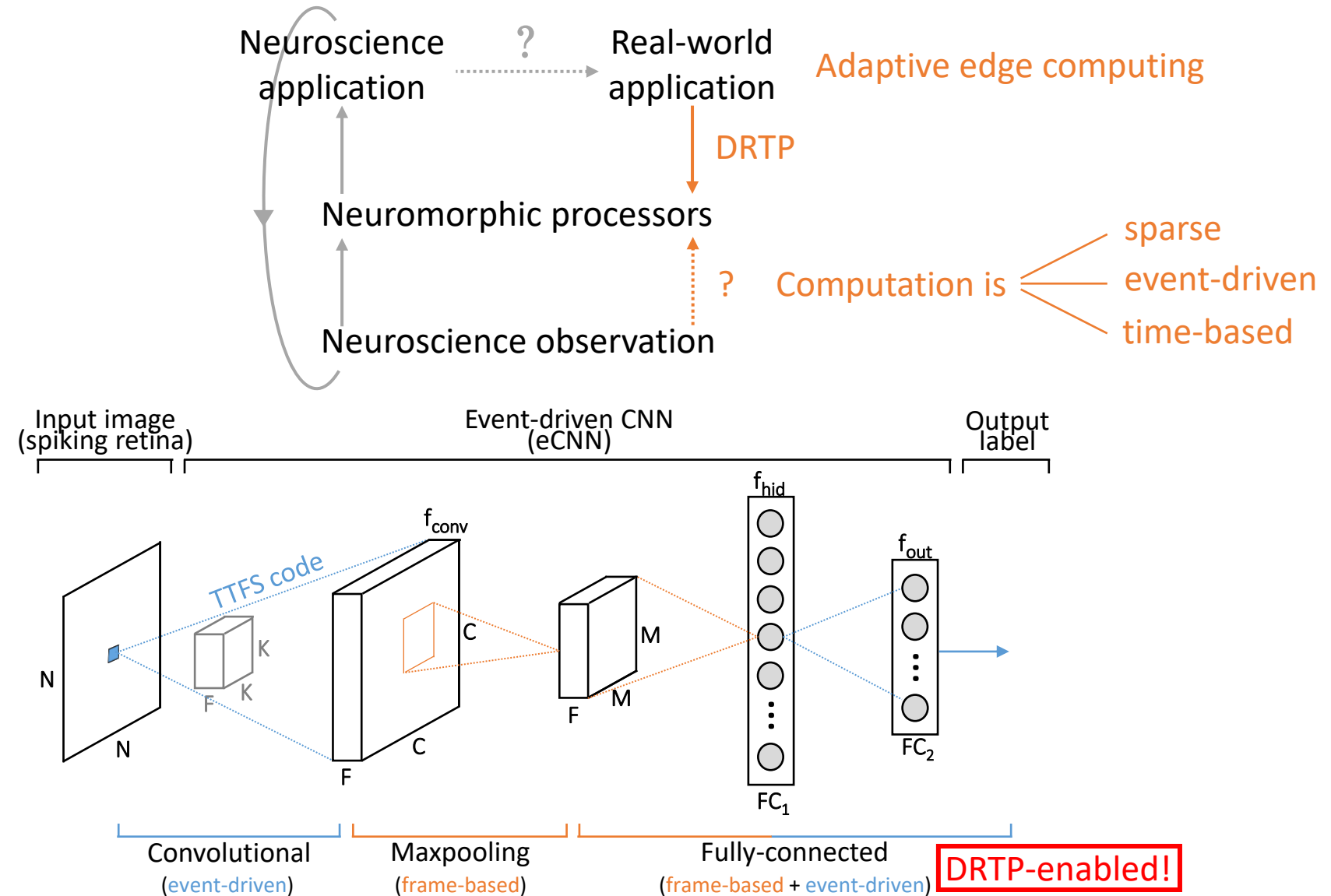
Neuromorphic accelerators

[Frenkel, *ISCAS*, 2020]  
(*Best paper award* 🏆)

Conclusion and perspectives

# Which bio-inspired elements?

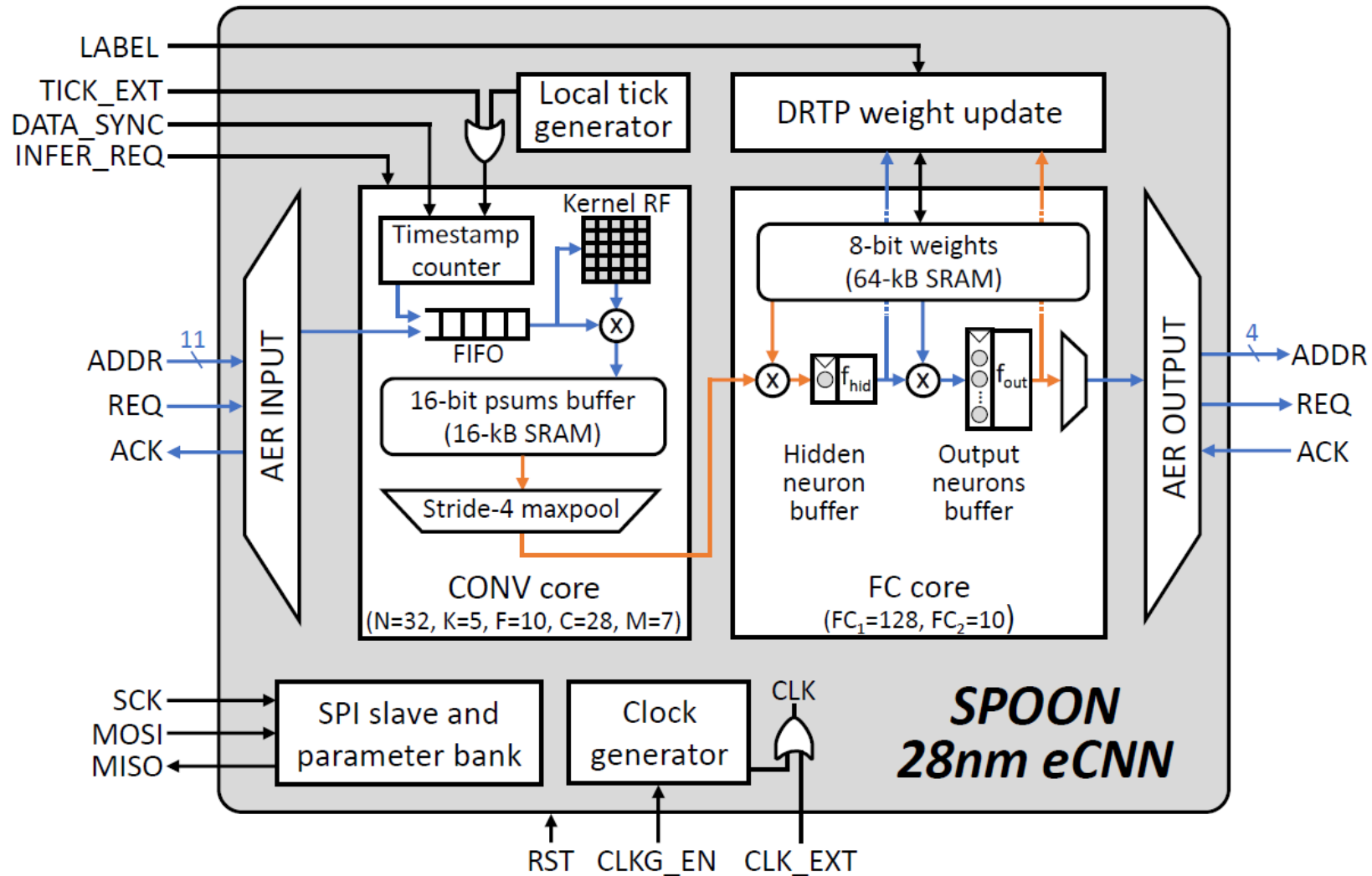
*Taking a step back with the top-down design strategy*





# Architecture of SPOON

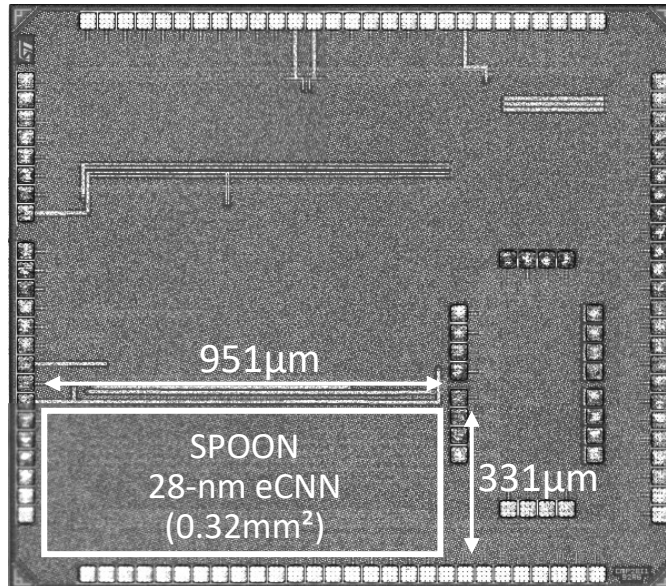
SPOON – A Spiking One-Learning Convolutional Neuromorphic Processor



# SPOON – Chip microphotograph and specifications



→ perspectives



*(pre-silicon numbers, not yet updated)*

Technology	28nm FDSOI CMOS
Implementation	Digital
Area	0.32mm <sup>2</sup> (0.26mm <sup>2</sup> excl. rails)
Topology	C5×5@10–FC128–FC10
Online learning	Stochastic DRTP, 8-bit weights
Time constant	Biological to accelerated
Supply voltage	0.6V – 1.0V
Max. clock frequency	150MHz
Leakage power	61μW at 0.6V
Energy for CONV core	1.7nJ/event at 0.6V
Energy for FC core	55nJ/inference at 0.6V
Online learning overhead	16.8% in power, 11.8% in area

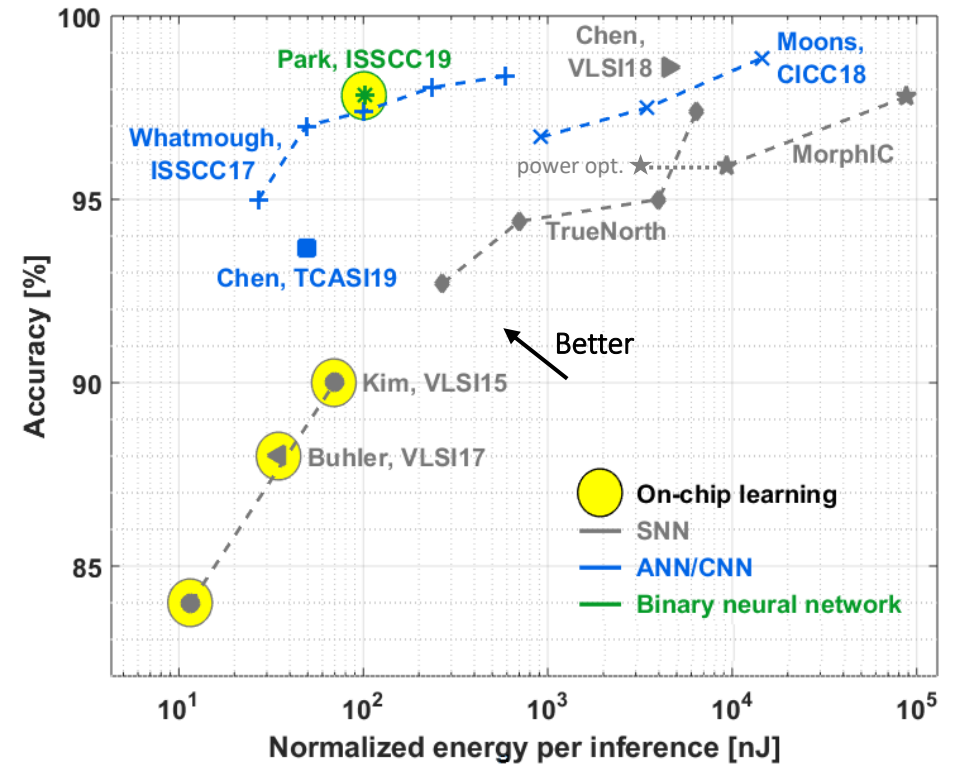
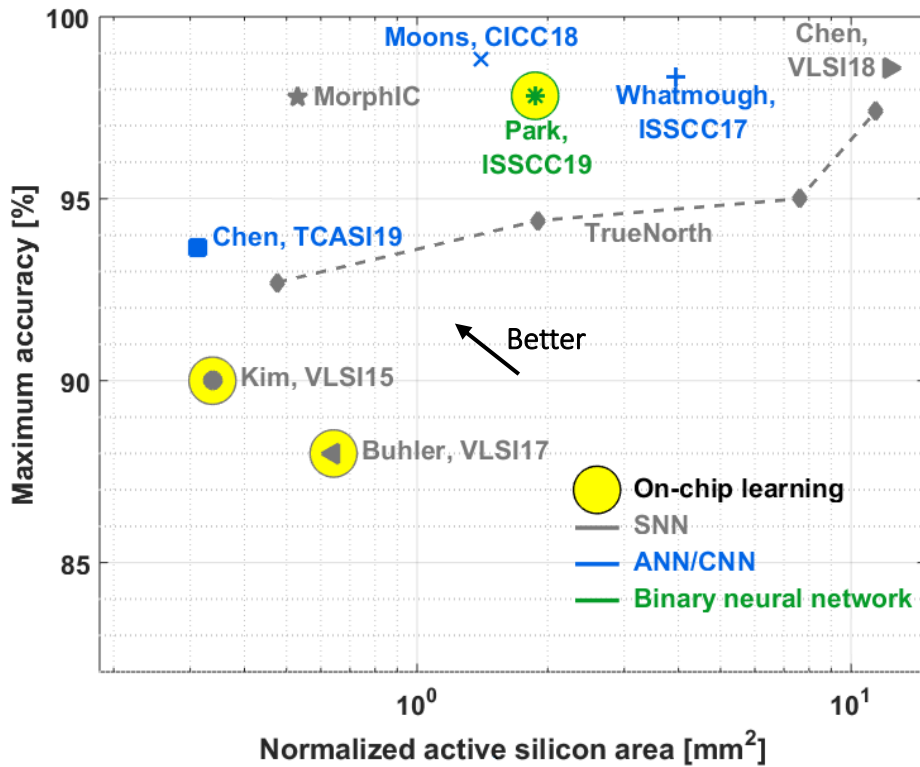
Stay tuned for the journal extension!

DRTP can be implemented on-chip at a very low cost!

Benchmarking: **MNIST** and N-MNIST

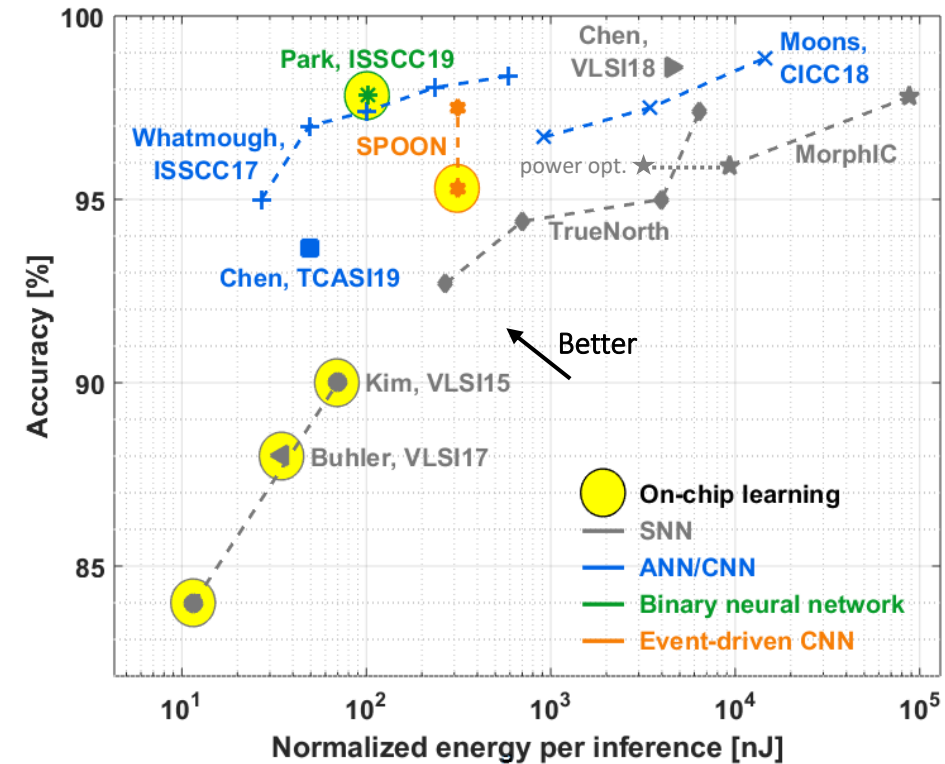
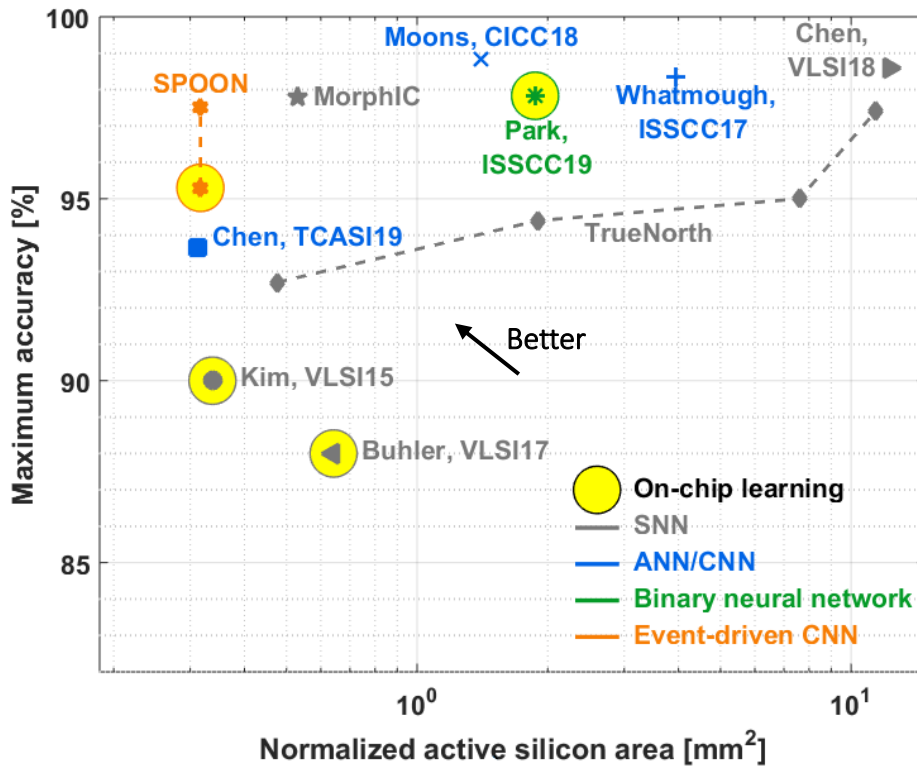
# SPOON benchmarking

*Against SoA spiking neural networks on MNIST*



# SPOON benchmarking

*Against SoA spiking neural networks on MNIST*



Only SPOON allows reaching the efficiency of ANN/CNN/BNN accelerators while enabling online learning with event-based sensors.

# Outline

## Part I – Bottom-up neuromorphic design

- Building blocks
- Integration

## Part II – Top-down neuromorphic design

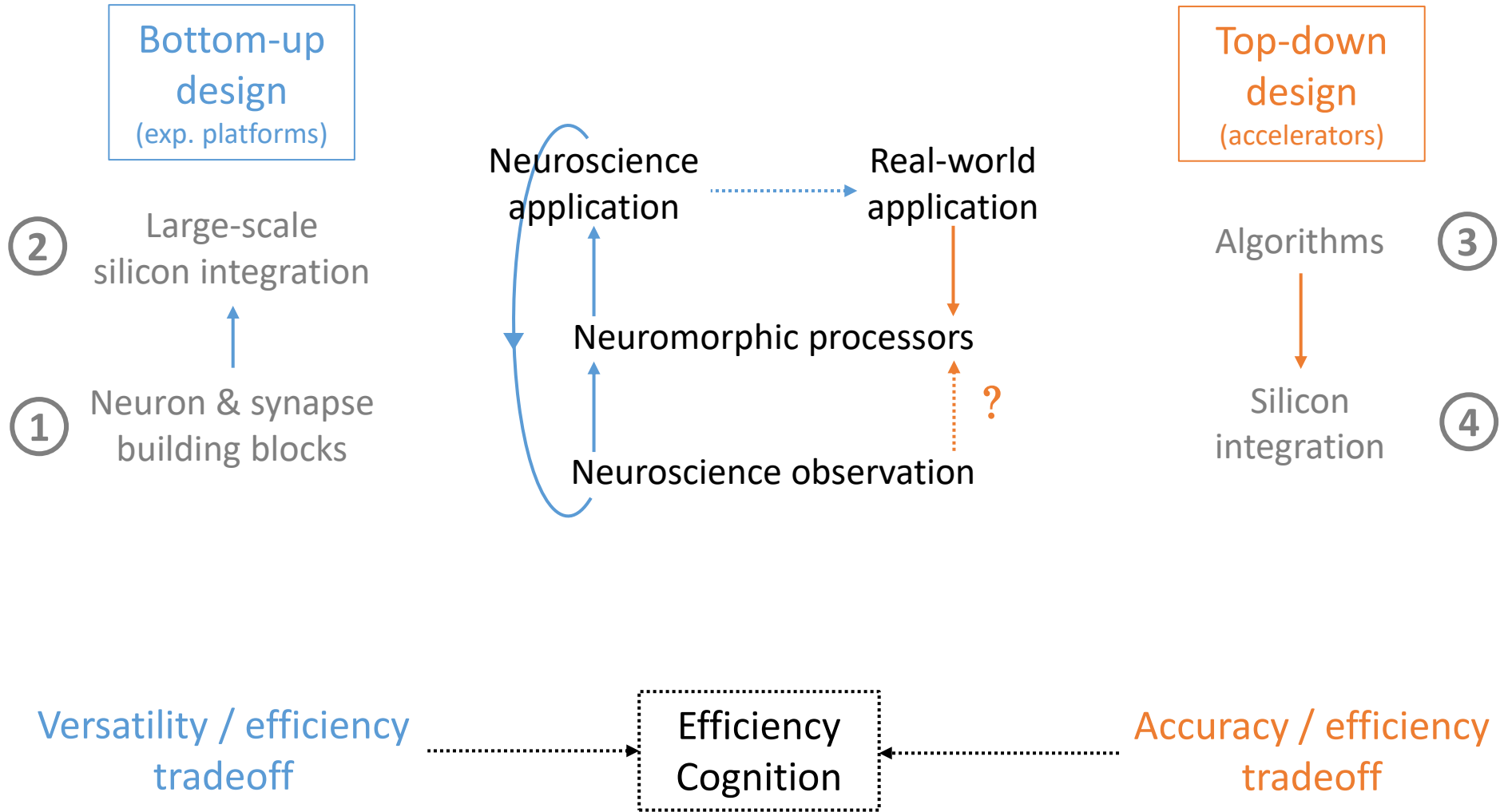
- Algorithms
- Integration

## Conclusion and perspectives

Summary of the key messages, next directions

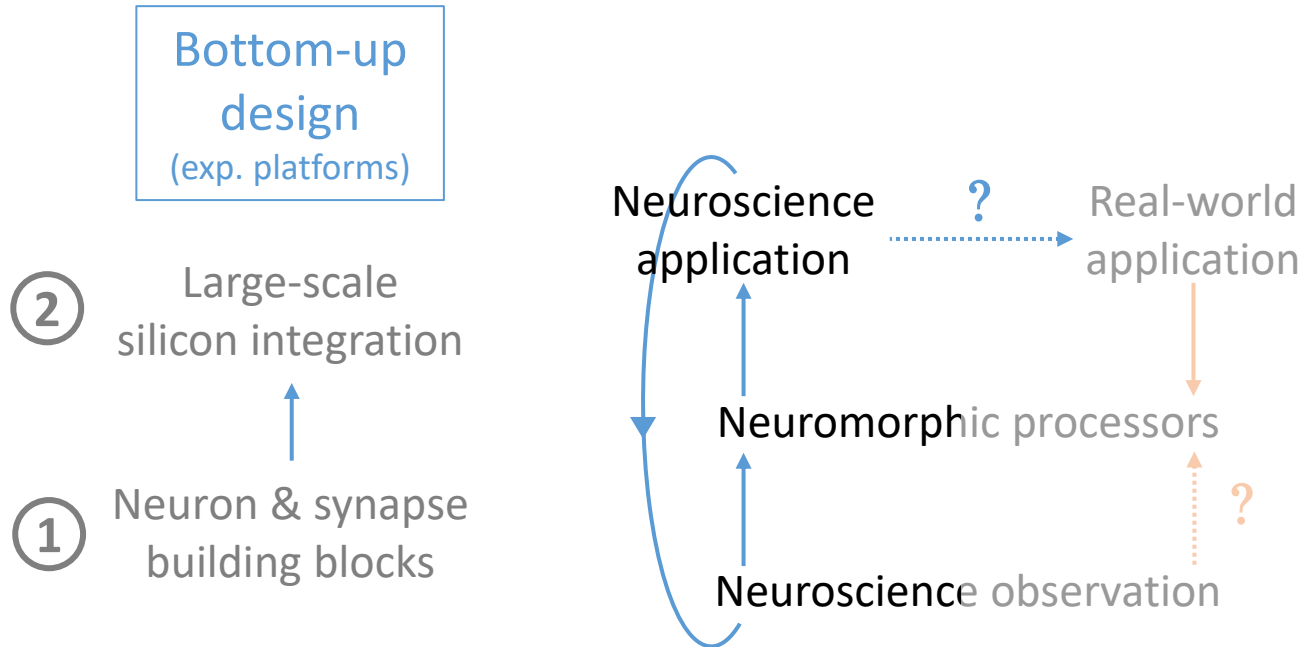
# Neuromorphic Engineering – Key Claims

*Unveiling roads to embedded cognition*



# Neuromorphic Engineering – Key Claims

*Unveiling roads to embedded cognition*



Versatility / efficiency  
tradeoff

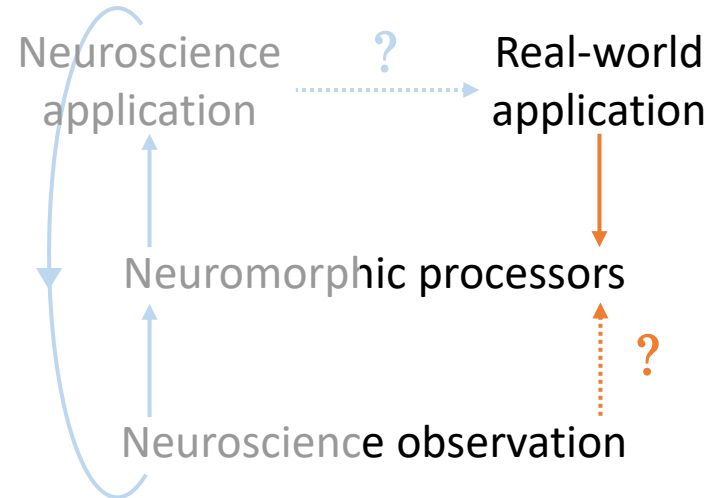
## Claim 1

Hardware-aware neuroscience model design and selection allows reaching record neuron and synapse densities with low-power operation for large-scale integration *in silico*.



# Neuromorphic Engineering – Key Claims

*Unveiling roads to embedded cognition*



Top-down  
design  
(accelerators)

Algorithms

③

Silicon  
integration

④

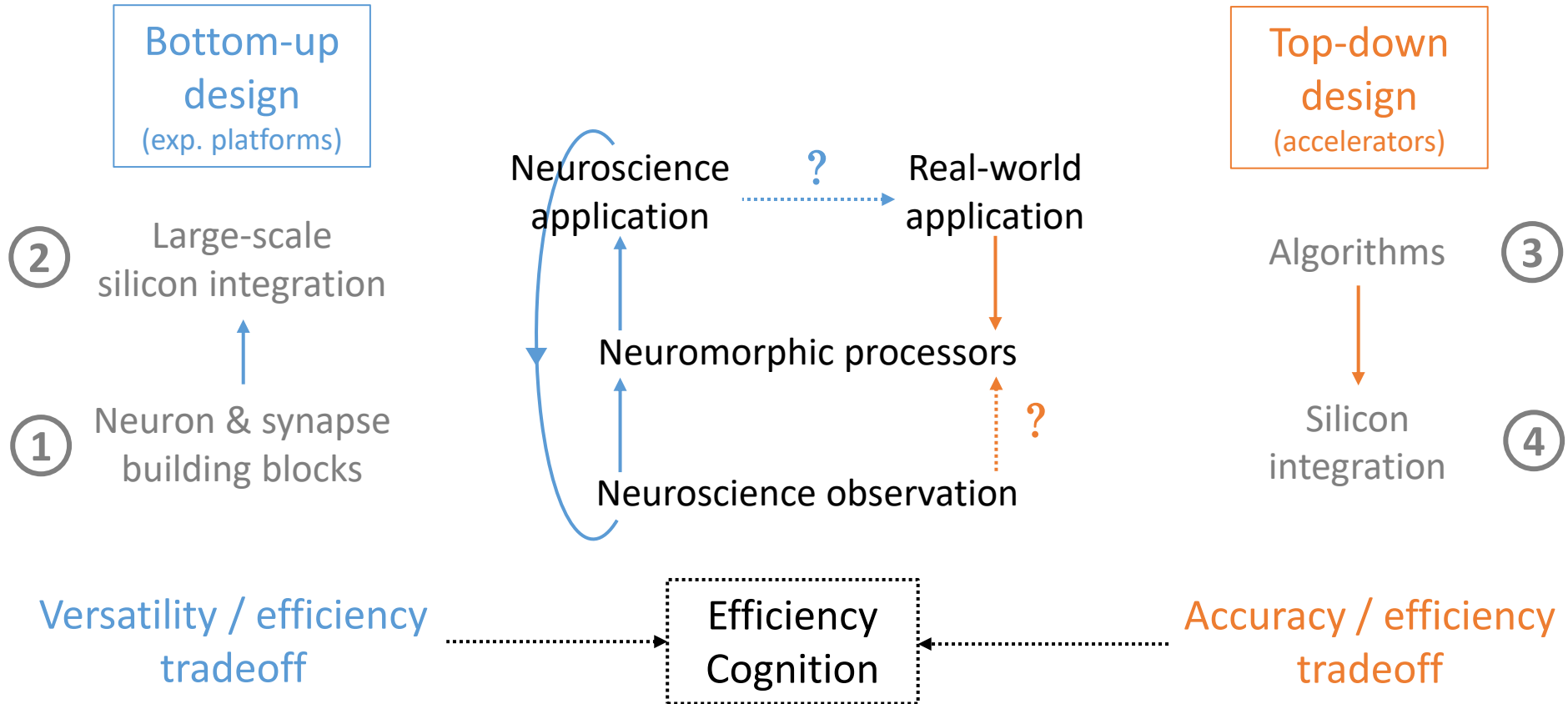
## Claim 2

Combining event-driven and frame-based processing with weight-transport-free update-unlocked training supports low-cost adaptive edge computing with spike-based sensors.

Accuracy / efficiency  
tradeoff

# Neuromorphic Engineering – Key Claims

*Unveiling roads to embedded cognition*

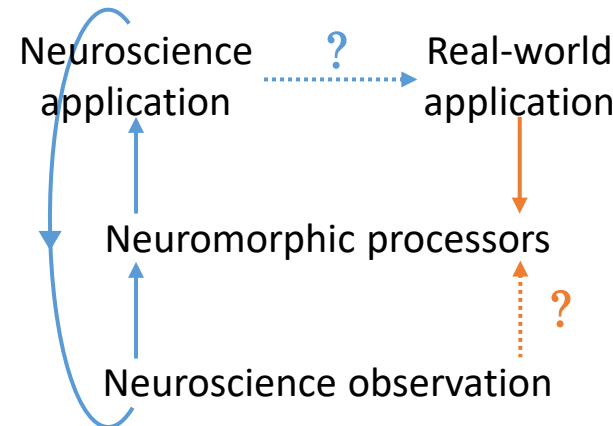


## Claim 3

Top-down guidance helps pushing bottom-up neuron and synapse integration beyond the purpose of neuroscience experimentation platforms, while bottom-up guidance supports top-down design toward brain reverse-engineering.

# Perspectives

- Neuromorphic engineering and spiking neural networks:  
“Can we make it work?” —→ “Will it bring a competitive advantage?” (not only against GPUs)  
Need something better than MNIST —→ Audio (KWS) and bio-signal processing (time, biological-time)  
*[Davies, Nat. Mach. Intel., 2019]*
- Phenomenological digital design: pragmatic short-to-midterm approach.  
Promising avenues: leveraging the variability of subthreshold analog design; fine-grained mixed-signal design.
- Bottom-up trend: dendrites
- Top-down trend: new wave of training algorithms mapping onto bio-plausible primitives  
*[Sacramento, NeurIPS’18]*  
*[Payeur, bioRxiv, 2020]*  
*[Bellec, Nat. Comms., 2020]*
- Cognition: a case for neuromorphic robots?  
*[Man & Damasio, Nat. Mach. Intel., 2019]*



# Acknowledgments

PhD



Postdoc

Institution:



University of  
Zurich <sup>UZH</sup>

**ETH**

Eidgenössische Technische Hochschule Zürich  
Swiss Federal Institute of Technology Zurich

Funding:



LE FONDS EUROPEEN DE DEVELOPPEMENT REGIONAL  
ET LA WALLONIE INVESTISSENT DANS VOTRE AVENIR



Supervisors:



*Profs. David Bol, Jean-Didier Legat*



*Prof. Giacomo Indiveri*

Key colleagues:



*Martin Lefebvre*



# Questions?



@C\_Frenkel



cfrenkel



Charlotte-Frenkel



ChFrenkel



charlotte@ini.uzh.ch

## 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]
- MorphIC: [C. Frenkel et al. “MorphIC: A 65-nm 738k-synapse/mm<sup>2</sup> quad-core binary-weight digital neuromorphic processor with stochastic spike-driven online learning,” *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 neural processing systems design: Neuromorphic intelligence as the convergence of natural and artificial intelligence”, *arXiv preprint arXiv:2106.01288*, 2021]

*Open-sourced!*

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

*Open-sourced!*

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

*Journal extension coming soon*

***Just released!***

# Premier Sponsor



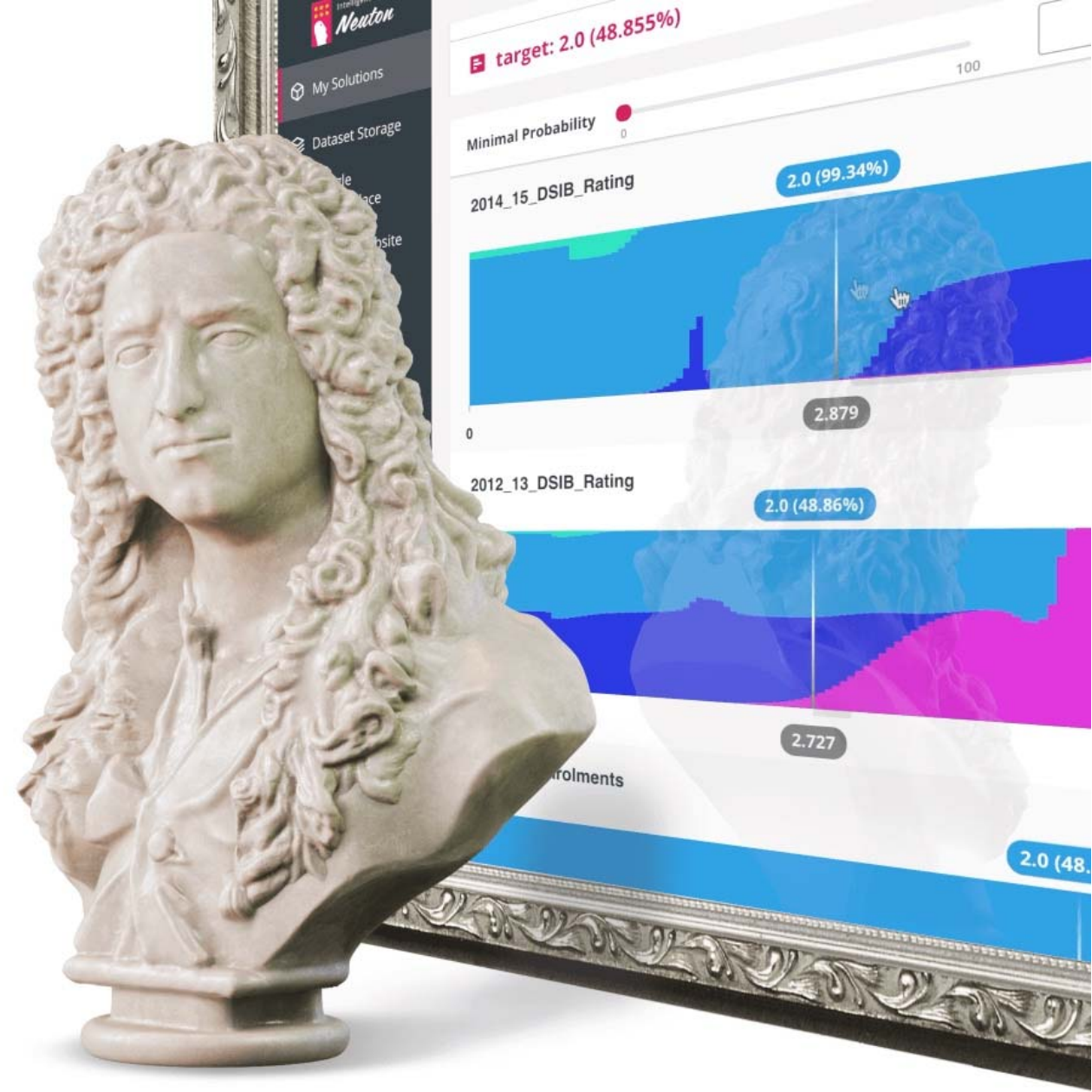
# Automated TinyML

Zero-code SaaS solution

**Create tiny models, ready for embedding,  
in just a few clicks!**

Compare the benchmarks of our compact models to those of TensorFlow and other leading neural network frameworks.

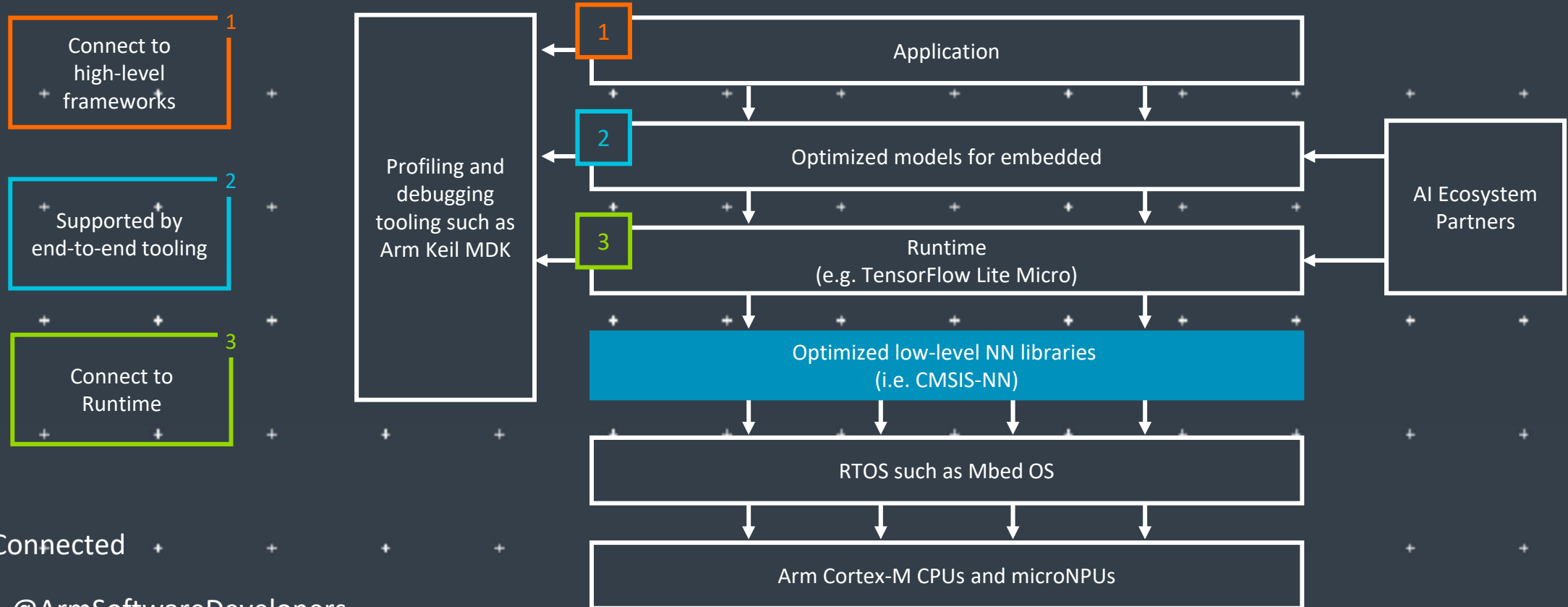
***Build Fast. Build Once. Never Compromise.***



# Executive Sponsors



# Arm: The Software and Hardware Foundation for tinyML



Stay Connected



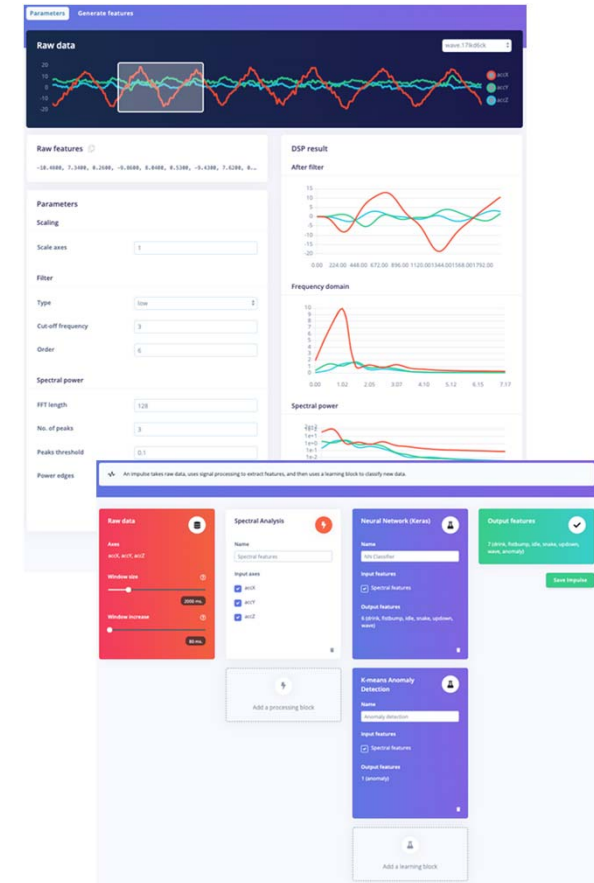
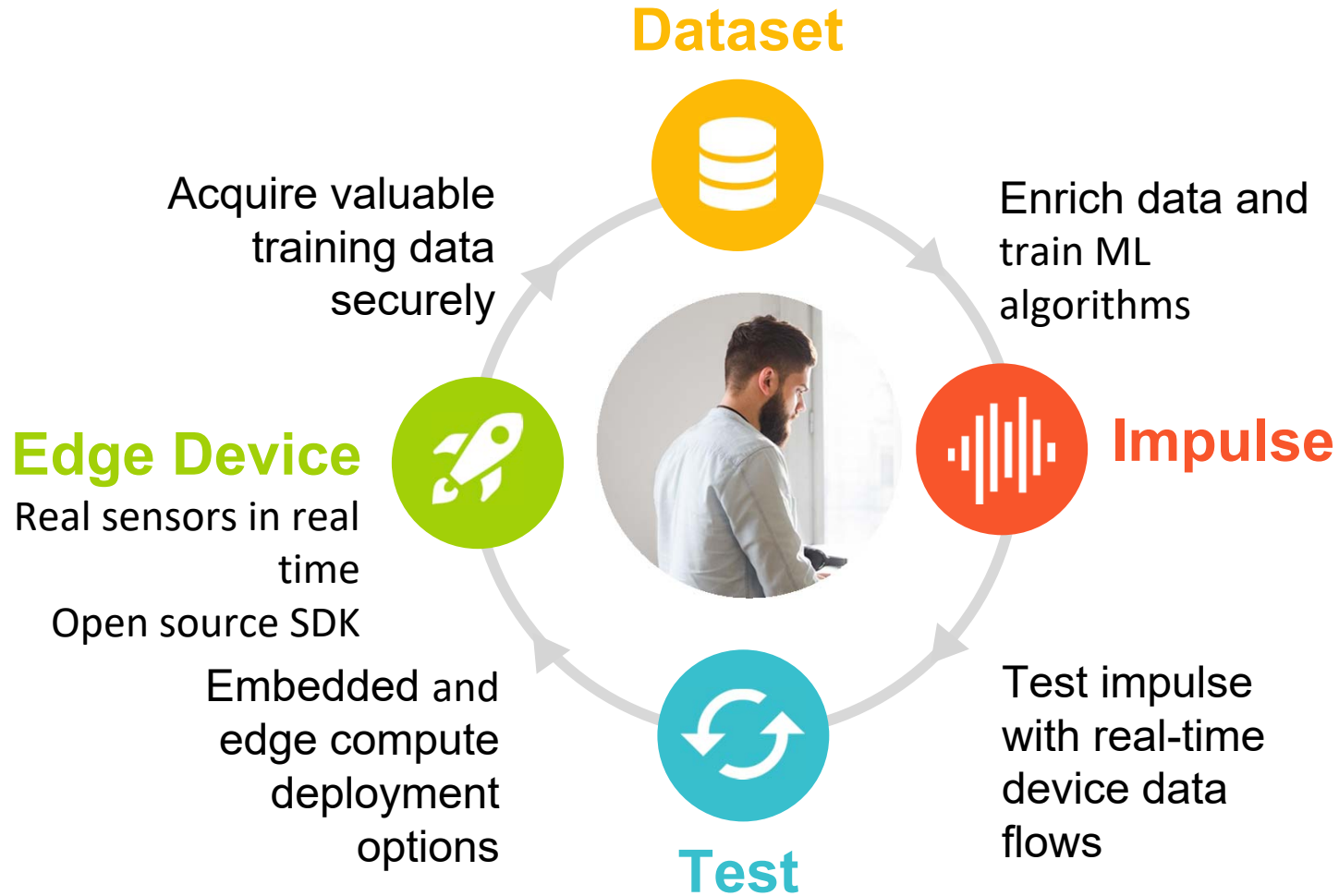
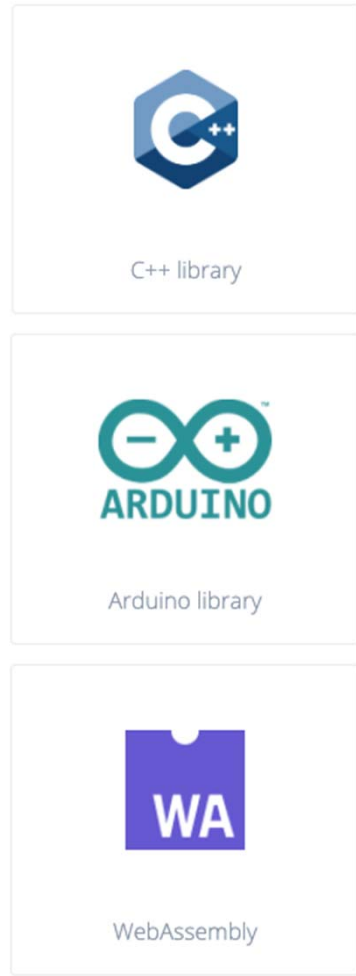
@ArmSoftwareDevelopers



@ArmSoftwareDev

Resources: [developer.arm.com/solutions/machine-learning-on-arm](https://developer.arm.com/solutions/machine-learning-on-arm)

# TinyML for all developers



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



# 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

# SYNTIANT

[Syntiant Corp.](#) is moving artificial intelligence and machine learning from the cloud to edge devices. Syntiant's chip solutions merge deep learning with semiconductor design to produce ultra-low-power, high performance, deep neural network processors. These network processors enable always-on applications in battery-powered devices, such as smartphones, smart speakers, earbuds, hearing aids, and laptops. Syntiant's Neural Decision Processors™ offer wake word, command word, and event detection in a chip for always-on voice and sensor applications.

Founded in 2017 and headquartered in Irvine, California, the company is backed by Amazon, Applied Materials, Atlantic Bridge Capital, Bosch, Intel Capital, Microsoft, Motorola, and others. Syntiant was recently named a [CES® 2021 Best of Innovation Awards Honoree](#), [shipped over 10M units worldwide](#), and [unveiled the NDP120](#) part of the NDP10x family of inference engines for low-power applications.

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



@Syntiantcorp

# Platinum Sponsors



Part of your life. Part of tomorrow.

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



# Reality AI<sup>®</sup>

## Add Advanced Sensing to your Product with Edge AI / TinyML

<https://reality.ai>



[info@reality.ai](mailto:info@reality.ai)



[@SensorAI](https://twitter.com/SensorAI)



[Reality AI](#)

## Pre-built Edge AI sensing modules, plus tools to build your own

### Reality AI solutions

Prebuilt sound recognition models for  
indoor and outdoor use cases

Solution for industrial anomaly detection

Pre-built automotive solution that lets cars  
“see with sound”

### Reality AI Tools<sup>®</sup> software

Build prototypes, then turn them into  
real products

Explain ML models and relate the function  
to the physics

Optimize the hardware, including  
sensor selection and placement

# Gold Sponsors





# LatentAI

## Adaptive AI for the Intelligent Edge

[Latentai.com](https://latent.ai)



# Build Smart IoT Sensor Devices From Data

SensiML pioneered TinyML software tools that auto generate AI code for the intelligent edge.

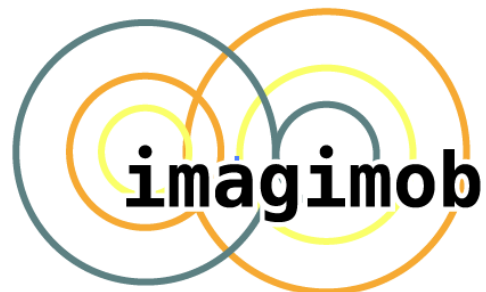
- End-to-end AI workflow
- Multi-user auto-labeling of time-series data
- Code transparency and customization at each step in the pipeline

We enable the creation of production-grade smart sensor devices.



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

# Silver Sponsors



# Copyright Notice

The presentation(s) in this publication comprise the proceedings of tinyML® EMEA Technical Forum 2021. The content reflects the opinion of the authors and their respective companies. This version of the presentation may differ from the version that was presented at tinyML EMEA. 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)