



# Event-Driven Neuromorphic Solutions for Efficient Dynamic Signal Processing

Dr. Federico Corradi

October 17, 2024

Assistant Professor, Neuromorphic Edge Computing Systems (NECS) Lab

Department of Electrical Engineering, Electronic Systems Group

## The Third Generation of Neural Nets

- Neurons Models

- Neural Coding Schemes

## Training of Spiking Neural Networks

- Conversion: ANN to SNN Conversion

- Time-To-First-Spike Training of SNNs

- Direct Training: Back-propagation Through Time (BPTT)

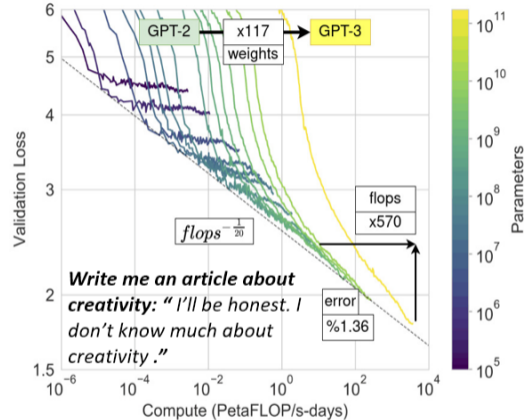
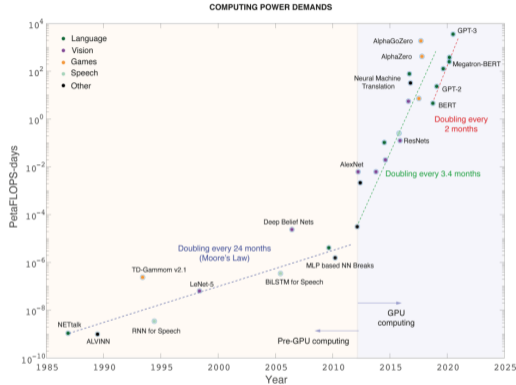
- Direct Training: Forward Propagation Through Time (FPTT)

## Neuromorphic Hardware

- Digital Neuromorphic Processors for The Edge

- Mixed-Signal Neuromorphic Processor ( $\mu$ Brain) for The Extreme Edge

# Deep Learning Seems To Have No Limits



## Energy Expensive (training):

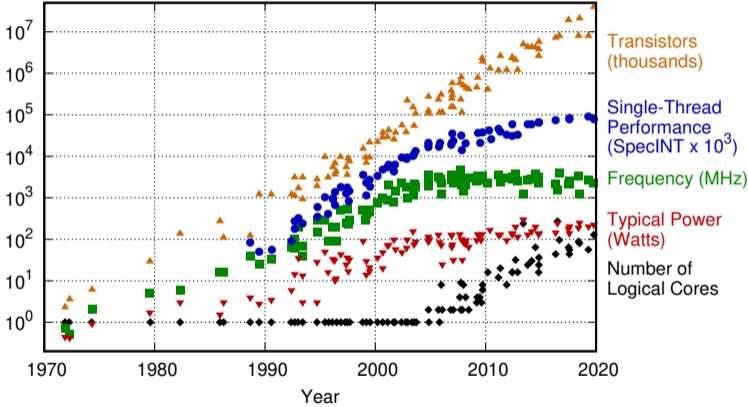
- GPT-3: 1.29 GWh
- GPT4: 50 GWh → 40x

Diminishing returns in deep learning  
(i.e., scaling follows a power law)

[Boahen K., Nature, 2002]

# HW has its Limits: 50 Years of Microprocessors

48 Years of Microprocessor Trend Data

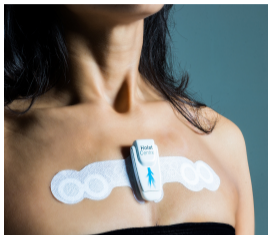


Original data up to the year 2010 collected and plotted by M. Horowitz, F. Labonte, O. Shacham, K. Olukotun, L. Hammond, and C. Batten  
New plot and data collected for 2010-2019 by K. Rupp

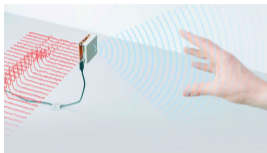
**Physical limits of CMOS & Von Neumann Bottleneck**

- Moore's Law
- Memory Wall
- Heat Wall

# Today, Deep Learning is Not Embedded!



**Biomedical**



**Mobile devices**



**Navigation**



**Fast Autonomy**

## Embedded Deep Learning Main Challenges

- Physical limits of CMOS & Von Neumann bottleneck
- Large volume of data, compute, and energy
- Brittle AI (fixed models)

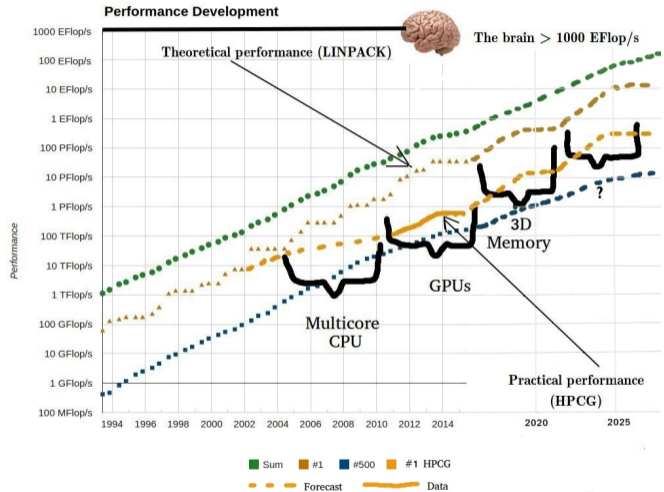
# Brain, Any Limits?

## Brain

- Speed = 1000 Exa Op/s\*
- Power = 20 Watt
- Energy/operation =  
Power / Speed =  
 $2 \times 10^{-5}$  fJ/operation

## Compare to Frontier

- Power = 22.7 MW
- Peak = 1.206 ExaFlop/s
- 8.7 Mcores, 680 m2
- Energy/op. = 19 pJ/op

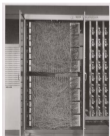
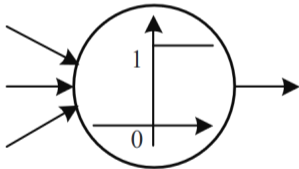


[Dettmers Tim, Blog]

# Evolutionary Journey of Neural Networks Through Bioinspiration

## First Generation

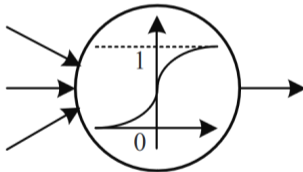
- No weights
- Threshold
- Exc/Inh inputs



[McCulloch & Pitts, 1942]

## Second Generation

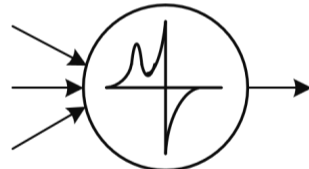
- Sigmoid / ReLU
- Real valued weights
- Differentiable



[Bridle J. 1989]

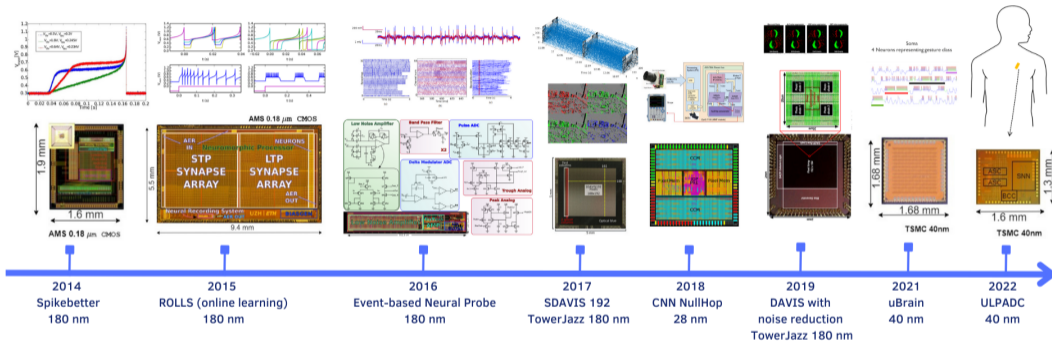
## Third Generation

- Spiking neurons
- In/Out spike trains
- Time-based models



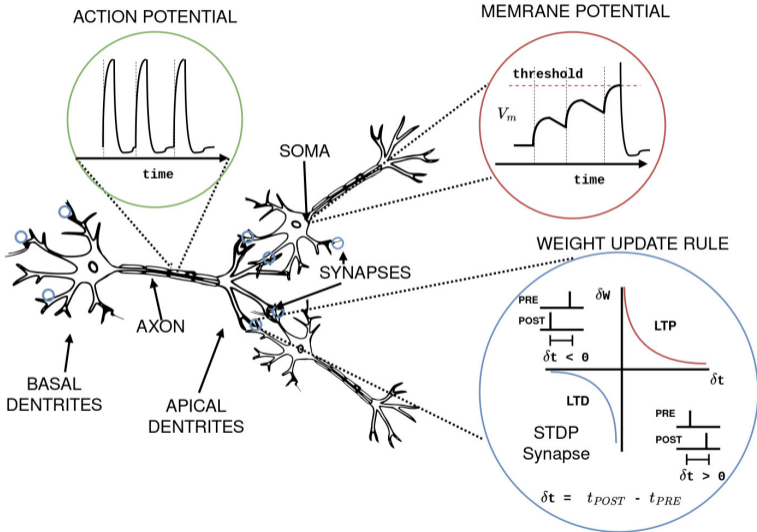
[Maas W., 1997] [Davies M., 2021]

# My Focus: Neuromorphic Sensing and Computing Systems

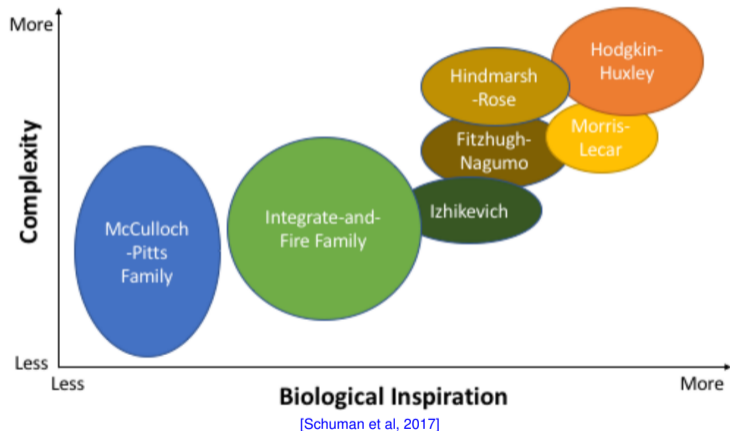


Co-design algorithms and spike-based neuromorphic sensing & computing architectures in CMOS and emerging technologies for edge applications.

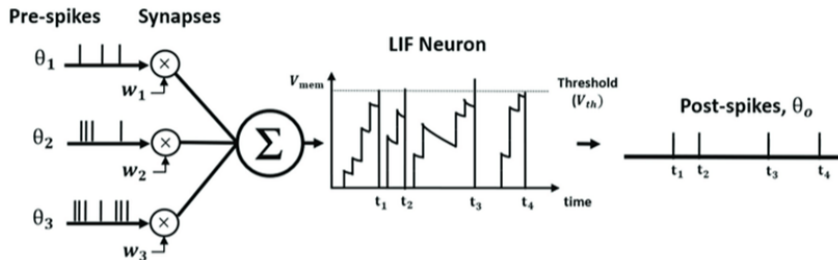
# What is a Spiking Neural Network?



# Neuron models: Complexity of Bio-realism



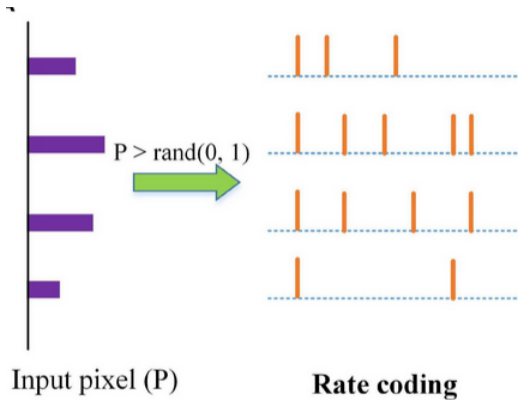
# Leaky-Integrate-and-Fire (LIF) Neuron Model



## LIF

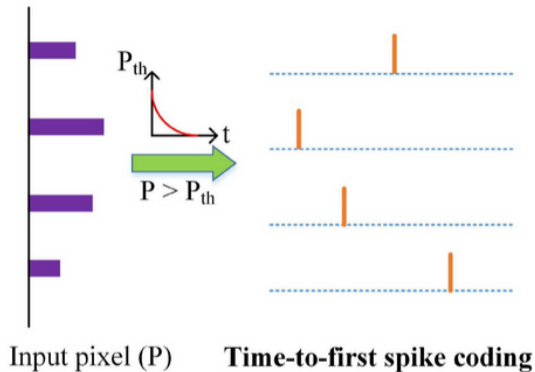
- input spikes, output spikes (digital communication)
- membrane potential decay over time (stateful neuron)
- simple to implement (e.g., RC circuit)

# Neural Coding Schemes



- Rate Coding,  $v = \frac{N_{spikes}}{T}$

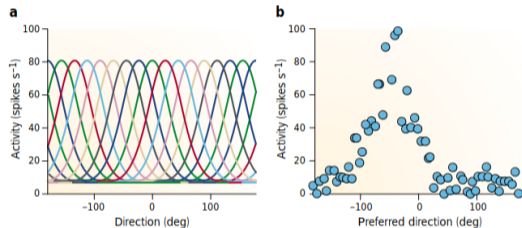
# Neural Coding Schemes



- Rate Coding,  $v = \frac{N_{spikes}}{T}$
- Temporal Coding (TTFS, Inter-Spike-Interval)

[Auge D. et al, 2021] [Stanojevic A. et al, 2024]

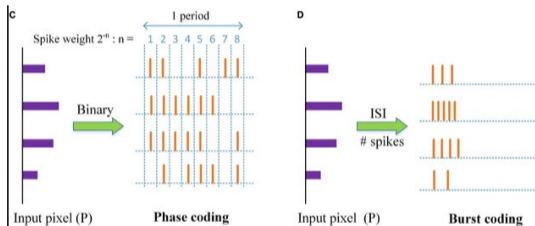
# Neural Coding Schemes



- Rate Coding,  $v = \frac{N_{spikes}}{T}$
- Temporal Coding (TTFS, Inter-Spike-Interval)
- Population Coding (neuron tuning curves)

[Auge D. et al, 2021] [Stanojevic A. et al, 2024] [Pouget et al, 2000]

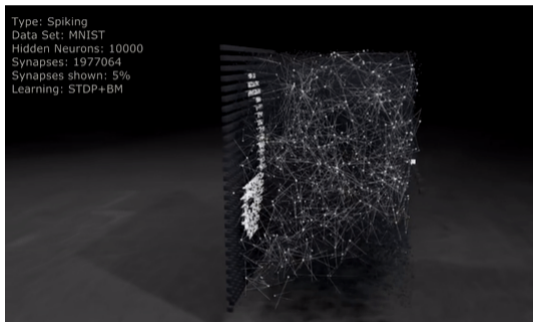
# Neural Coding Schemes



- Rate Coding,  $\nu = \frac{N_{spikes}}{T}$
- Temporal Coding (TTFS, Inter-Spike-Interval)
- Population Coding (neuron tuning curves)
- Sparse Coding, Phase Coding, etc...

[Auge D. et al, 2021] [Stanojevic A. et al, 2024] [Pouget et al, 2000]

# Summary Part I: Spiking Neural Networks



## SNN properties

- Fully parallel
- Spatiotemporal processing
- Asynchronous
- Sparse
- Additive weight operations
- Synaptic Plasticity

## The Third Generation of Neural Nets

- Neurons Models

- Neural Coding Schemes

## Training of Spiking Neural Networks

- Conversion: ANN to SNN Conversion

- Time-To-First-Spike Training of SNNs

- Direct Training: Back-propagation Through Time (BPTT)

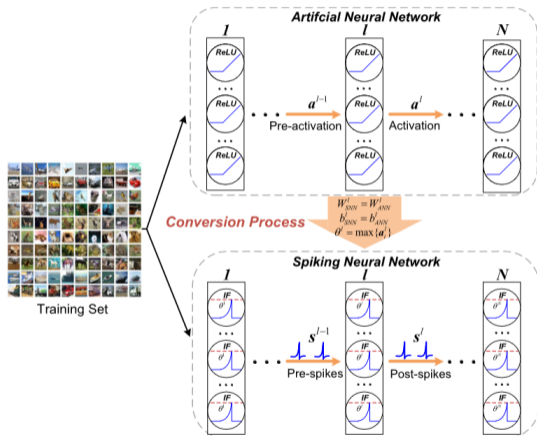
- Direct Training: Forward Propagation Through Time (FPTT)

## Neuromorphic Hardware

- Digital Neuromorphic Processors for The Edge

- Mixed-Signal Neuromorphic Processor ( $\mu$ Brain) for The Extreme Edge

# Analog Neural Networks (ANN) to SNN Conversion



## ANN to SNN

- ReLU (no bias)
- Mean-rate coding
- Train ANN using BP
- Copy weights into SNN
- Infer using SNN

[Diehl P. et al 2015] [Rueckauer B. et al. 2017] [Sengupta A. et al. 2019] [Wang Y. et al. 2024]

# Analog Neural Networks (ANN) to SNN Conversion: Why Bother?

Energy efficiency comparison between ANN and SNN converted by our framework.

		$T = 2$	$T = 4$	$T = 8$	$T = 12$	$T = 16$
CIFAR10	$OP_{ANN}$	256.61M	256.61M	256.61M	256.61M	256.61M
	$OP_{SNN}$	89.07M $\pm$ 0.04M	177.5M $\pm$ 0.07M	351.89M $\pm$ 0.08M	526.09M $\pm$ 0.25M	700.41M $\pm$ 0.26M
	Avg Fire Rate	0.1834	0.1856	0.1857	0.186	0.1864
	$Acc_{ANN}$	95.67	95.67	95.67	95.67	95.67
	$Acc_{SNN}$	90.17 $\pm$ 0.21	94.14 $\pm$ 0.13	95.14 $\pm$ 0.07	95.35 $\pm$ 0.08	95.45 $\pm$ 0.10
	Acc. Drop $\times$ 100%	-5.5% $\pm$ 0.21%	-1.60% $\pm$ 0.13%	-0.55% $\pm$ 0.07%	-0.33% $\pm$ 0.08%	-0.23% $\pm$ 0.10%
	Energy Drop $\times$ 100%	-89.03% $\pm$ 0.01%	-82.43% $\pm$ 0%	-69.72% $\pm$ 0.01%	-57.06% $\pm$ 0.02%	-44.4% $\pm$ 0.02%
CIFAR100	$OP_{ANN}$	256.98M	256.98M	256.98M	256.98M	256.98M
	$OP_{SNN}$	92.48M $\pm$ 0.02M	180.29M $\pm$ 0.03M	352.63M $\pm$ 0.09M	524.96M $\pm$ 0.23M	697.6M $\pm$ 0.12M
	Avg Fire Rate	0.1958	0.1908	0.1883	0.188	0.1882
	$Acc_{ANN}$	77.29	77.29	77.29	77.29	77.29
	$Acc_{SNN}$	64.89 $\pm$ 0.29	73.58 $\pm$ 0.17	76.35 $\pm$ 0.22	77.13 $\pm$ 0.08	77.22 $\pm$ 0.17
	Acc. Drop $\times$ 100%	-12.4% $\pm$ 0.29%	-3.71% $\pm$ 0.17%	-1.22% $\pm$ 0.22%	-0.21% $\pm$ 0.08%	-0.09% $\pm$ 0.17%
	Energy Drop $\times$ 100%	-87.95% $\pm$ 0%	-81.21% $\pm$ 0.01%	-68.76% $\pm$ 0.01%	-56.41% $\pm$ 0.01%	-44.06% $\pm$ 0.01%

## Energy Accuracy Trade-off

SNNs offers a good trade-off among energy and accuracy.

[Wang Y. et al. 2024]

# TTFS Nets Have Exactly The Same Performance of ANN

## Time-To-First-Spikes Nets

- SNN equivalent to feed-forward ReLU networks
- Same training behavior (equivalent gradient)
- Stateless network (frame-based)

nature communications



Article

<https://doi.org/10.1038/s41467-024-51110-5>

## High-performance deep spiking neural networks with 0.3 spikes per neuron

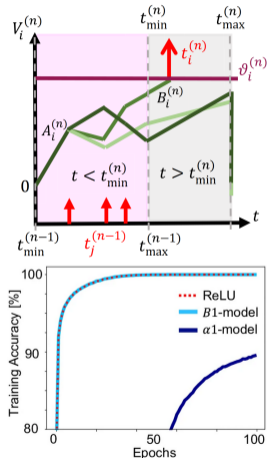
Received: 20 November 2023

Ana Stanojevic<sup>1,2</sup>, Stanisław Woźniak<sup>1</sup>, Guillaume Bellec<sup>2,3</sup>,  
Giovanni Cherubini<sup>1</sup>, Angeliki Pantazi<sup>1</sup> & Wulfram Gerstner<sup>2,3</sup>

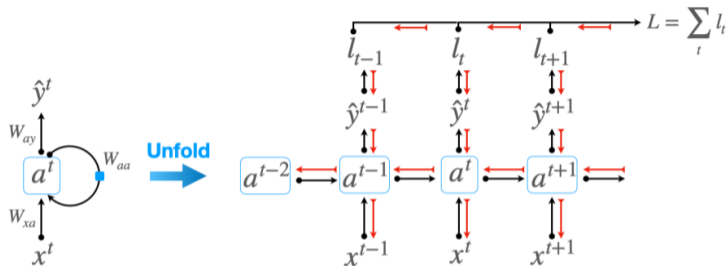
Accepted: 29 July 2024

Published online: 09 August 2024

[Stanojevic et al, 2024]



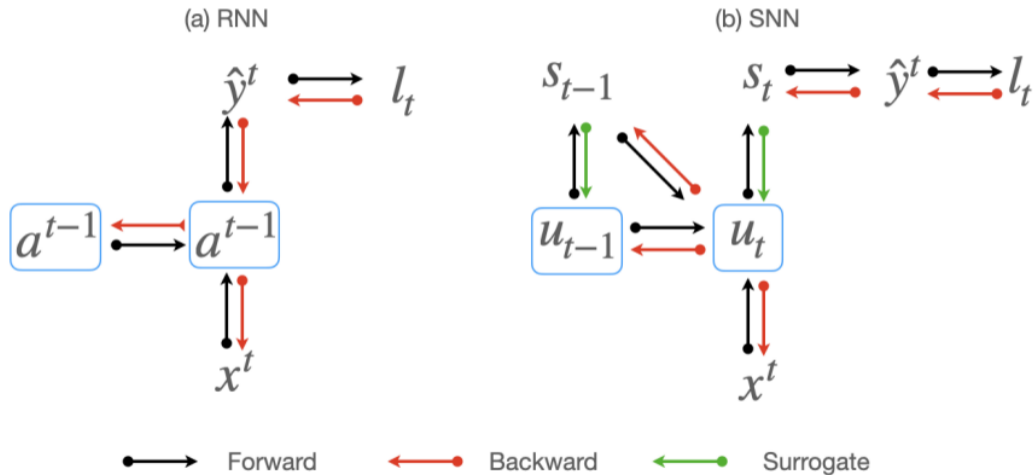
# Spiking LIF Neurons: Backpropagation Through Time (BPTT)



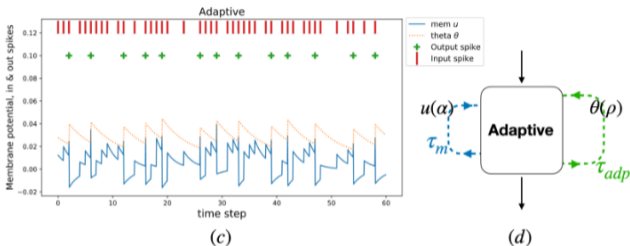
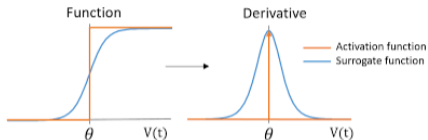
## BPTT

- Exploit Back-propagation
- Store the forward (black arrows)
- Unfold in time backward (red arrows)

# Spiking LIF Neurons: Backpropagation Through Time (BPTT)



# Spiking Recurrent Neural Networks

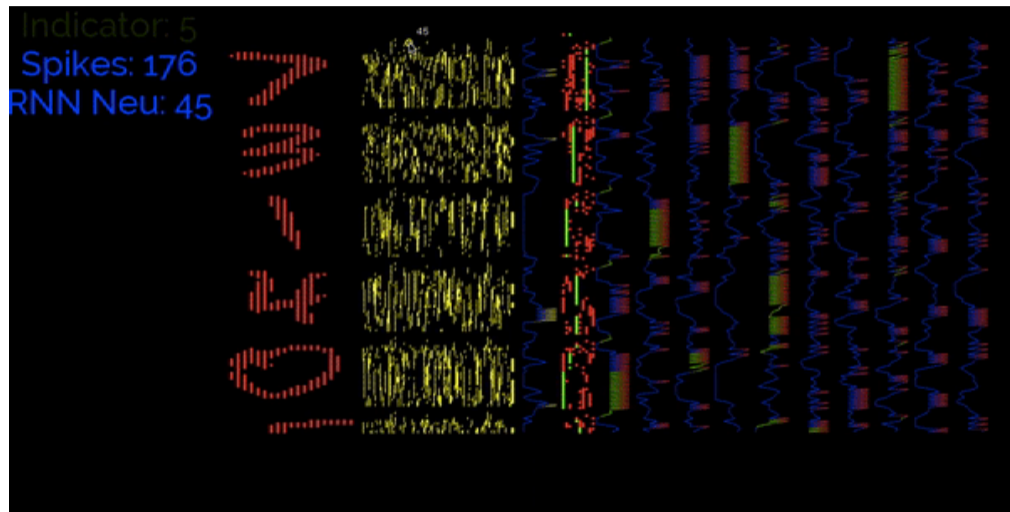


## Spiking Recurrent Neural Networks (SRNN)

SRNNs achieve **competitive** performance of classical recurrent neural networks (RNNs) while exhibiting **sparse** activity.

[Yin B., Corradi F., Bothe S., ICONS, 2021]

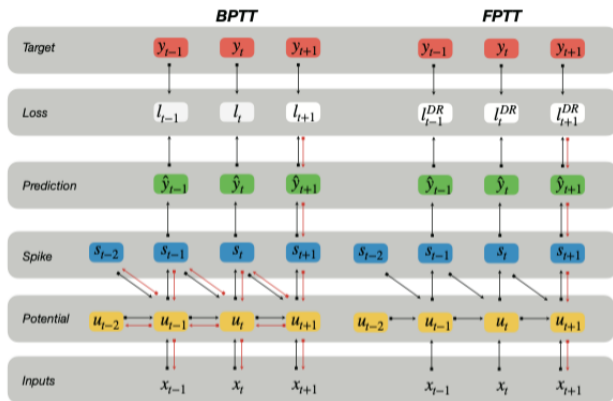
# Spiking Neural Networks of LIF Neurons: SMNIST example



## Algorithmic limitations of BPTT

- Cumbersome gradient computation path (via STE), and the **approximation error surrogate gradient accumulates** along time.
- Difficulty of back-propagating the gradient in deep spiking neural network models, **vanishing gradient**.
- **Memory requirements** for network training increases with the sequence length. It is challenging to train deep spiking neural network models on long sequences or when timestep become small.

# Forward Propagation Through Time (FPTT)



[Yin B., Corradi F., Bothe S., Nature Machine Intelligence, 2023]

## Intuition

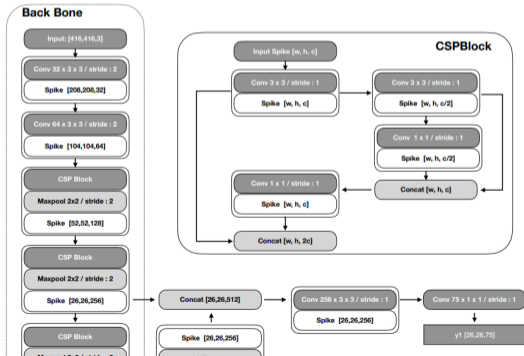
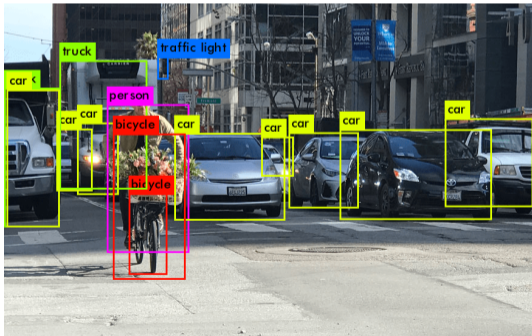
FPTT optimizes an instantaneous risk function, which includes a loss (as BPTT), and a dynamic regularization penalty (prev. observed losses).

## FPTT (online)

- Lessens surrogate gradient cumulative effect
- Lessens memory requirements
- Cost more computation

# Deep Spiking Neural Networks Models (SpyYoloV4)

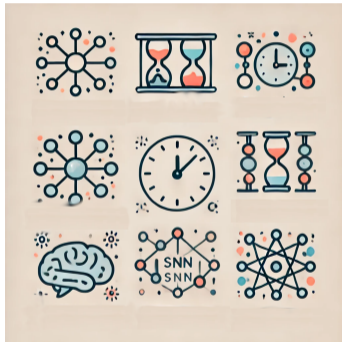
- 6.2 M spiking neurons
- Conv, fully connected, cross-stage partial subnets
- 14 M parameters
- Approaching DNNs



# Deep Spiking Neural Networks Models (SpyYoloV4)



[Yin B., Corradi F., Bothe S., Nature Machine Intelligence, 2023]



Training **competitive and deep SNN** today is possible!  
Interesting properties for HW:

- Sparse communication
- Event-driven computation
- Temporal information processing
- Trade-off energy and accuracy
- Same accuracy as ANN!

## The Third Generation of Neural Nets

- Neurons Models

- Neural Coding Schemes

## Training of Spiking Neural Networks

- Conversion: ANN to SNN Conversion

- Time-To-First-Spike Training of SNNs

- Direct Training: Back-propagation Through Time (BPTT)

- Direct Training: Forward Propagation Through Time (FPTT)

## Neuromorphic Hardware

- Digital Neuromorphic Processors for The Edge

- Mixed-Signal Neuromorphic Processor ( $\mu$ Brain) for The Extreme Edge

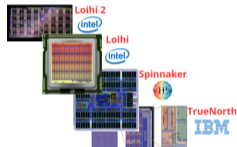
## Neuromorphic Computing Systems

- **Digital Neuromorphic Hardware**
  - Software models.
  - Easy to design, scalable.
- **Mixed-Signal Neuromorphic Hardware**
  - Bio-realistic models.
  - Hard to design, noise, mismatch.
- **Emerging Technologies**
  - Non-volatile memories.
  - Spintronics.
  - New design styles (e.g., in-memory computing).

### Analog neuromorphic architectures



### Digital neuromorphic architectures



### Emerging neuromorphic architectures

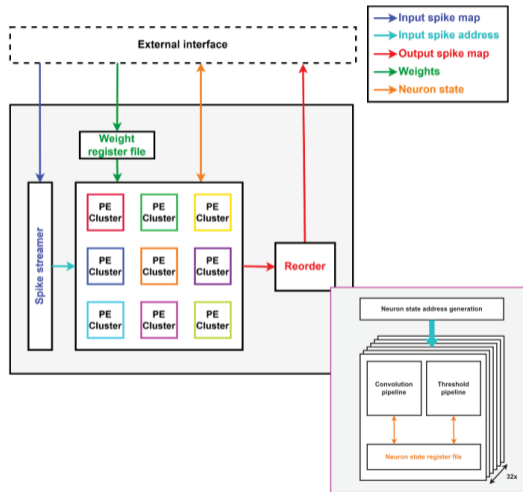


# Mega: Digital SNN Accelerator for SpyYoloV4 (TU/e)

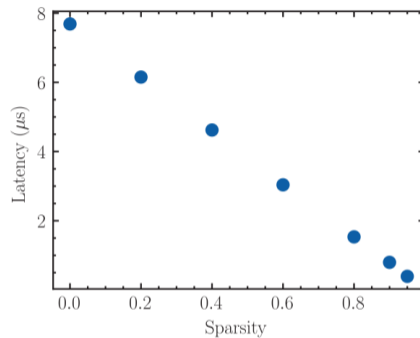
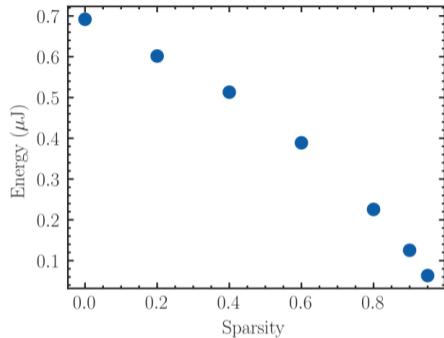
- 3x3 spiking convolution
  - Event-based: operations happen when there is a spike
  - Each PE cluster updates 32 neurons in parallel
- Memory hierarchy
  - Small, fast memories near the PEs
  - Larger memory for flexibility
- 0.501 pJ/SOP
  - Post-synthesis in 22 nm
  - 625 MHz @ 0.8V



[Luiken R., unpublished]



# Mega: What Happens if Spiking is Sparse?



- Energy and latency for a 96x48 spike map.
- SpyYoloV4 application.

[Luiken R., unpublished]

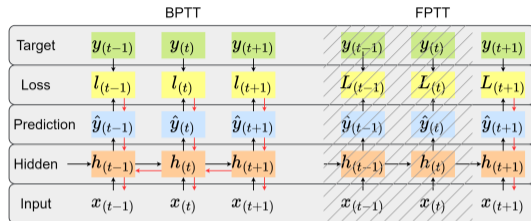
# Recurrent Neural Network Training at the Edge

- traditionally in one go: update the weight **once** by all network states generated after processing the entire sequence

$$W_{new} = W_{old} - \eta \frac{\partial \sum_{t=1}^T l(t)}{\partial W_{old}} \quad (1)$$

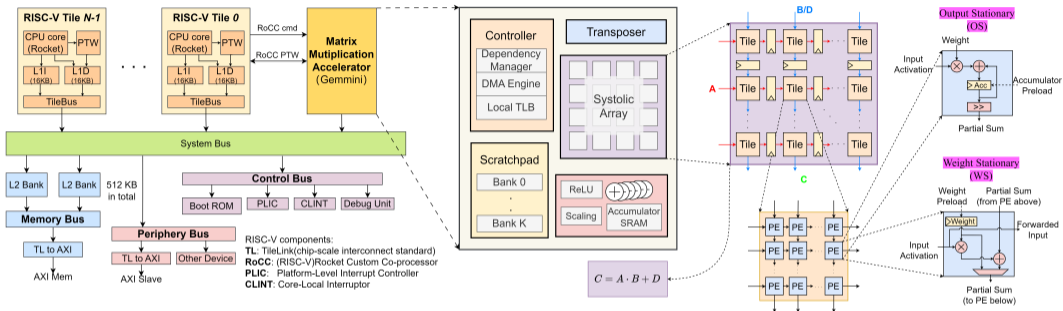
- Forward Propagation Through Time: update by the current network state **after every time step**; application of regularization term  $R(t)$  for stabilization

$$W_{(t+1)} = W_{(t)} - \eta \frac{\partial (l_{(t+1)} + R_{(t+1)})}{\partial W_{(t)}} \quad (2)$$



- How** to efficiently accelerate FPTT computations by designing an edge hardware architecture?
- Explore** several level of partitioning for the sequence.
- Sequential** MNIST Divide it into K parts: 1,2,7,14,28,56

# Customized Embedded Platform for FPTT Training



## Architecture

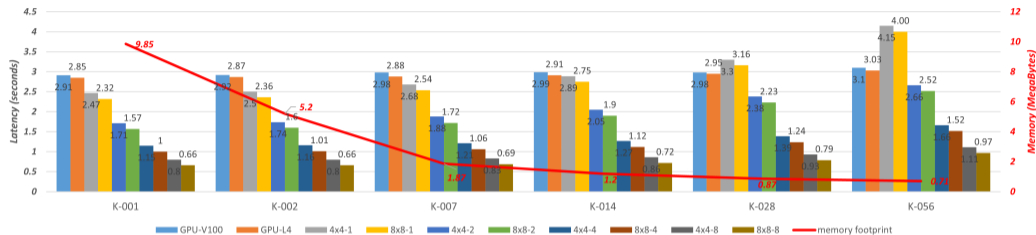
- RISC-V & Systolic Arrays
- Brain Float 16 (BF16), low memory req.
- Weight Stationary only

## Explore System Scaling

- 2 mesh sizes of the systolic array in Gemmini: 4x4, 8x8
- 4 types of CPU cores: 1, 2, 4, 8

[Zhang Y. Gomony D.M., Corporral H., Corradi F., VLSISOC, 2024]

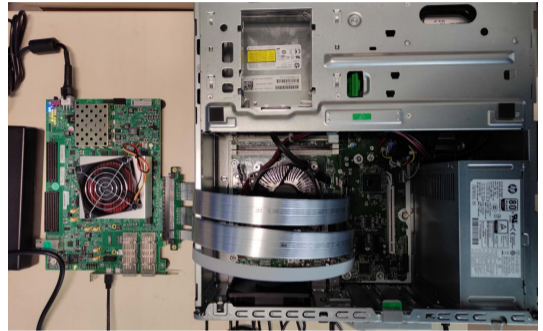
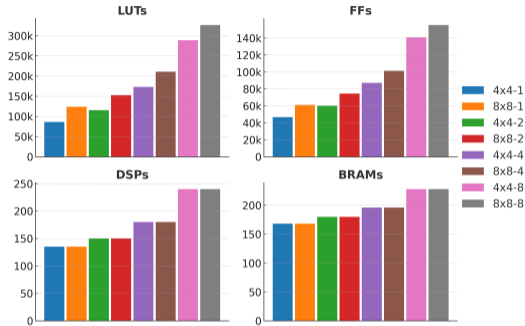
# Results: Training at the Edge (FPTT) Architecture



- For most clusters of K, customized architectures outperform two GPU platforms in latency.
- Trade-off between latency and resource utilization, given any cluster of K.
- Degree of partition: for the cluster, a benchmark of larger K, i.e., finer partition over a sequence leads to a slight latency increase but significant memory saving.

[Zhang Y. Gomony D.M., Corporral H., Corradi F., VLSISOC, 2024]

# Results: Training at the Edge (FPTT) Architecture



Cycle-accurate hardware simulation in FPGA.

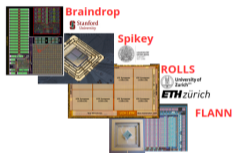
Latency and Memory Savings when exploiting FPTT with a Compressed BF16 architecture.

[Zhang Y. Gomony D.M., Corporral H., Corradi F., VLSISOC, 2024]

## Neuromorphic Computing Systems

- **Digital Neuromorphic Hardware**
  - Software models.
  - Easy to design, scalable.
- **Mixed-Signal Neuromorphic Hardware**
  - Bio-realistic models.
  - Hard to design, noise, mismatch.
- **Emerging Technologies**
  - Non-volatile memories.
  - Spintronics.
  - New design styles (e.g., in-memory computing).

### Analog neuromorphic architectures



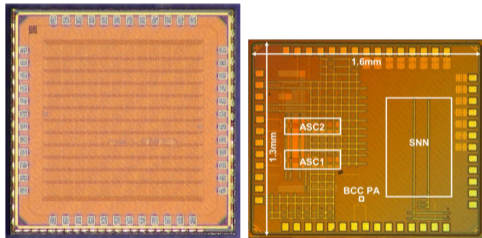
### Digital neuromorphic architectures



### Emerging neuromorphic architectures



## Synthesizable mixed-signal

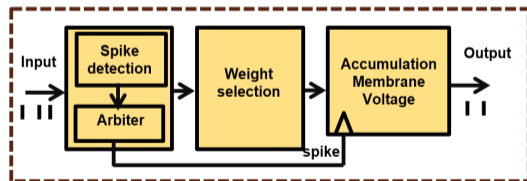
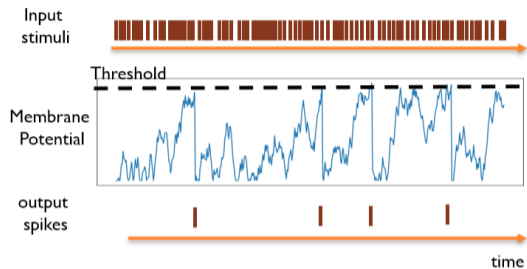


[Stuijt J. ... Corradi F., Front. Neurosc., 2021] [He Y., JSSC, 2022]

## Synthesizable mixed-signal:

- Benefits:
  - Fully asynchronous
  - Fast development-to-deployment cycle
  - Ultra-low-power
- Downsides:
  - Larger area
  - Difficult to train

# $\mu$ Brain: Fully Synthesizable Mixed-signal Neuromorphic Processor

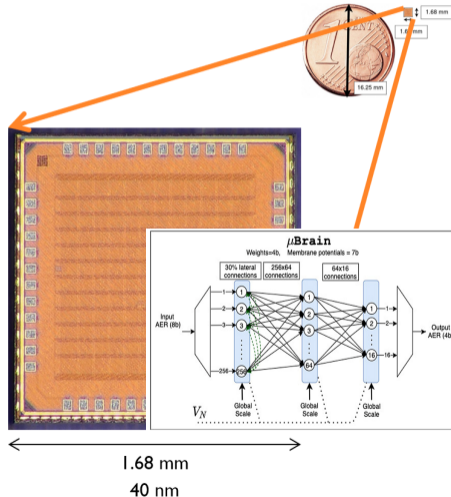


## LIF

- Self-timed (asynchronous)
- Programmable synapses
- Scalable
- Assumption: arbitration delays are in the order of  $ns$ , while incoming spikes are spaced in  $\mu s$  or even  $ms$

[Stuijt J., ... Corradi F., Front. Neurosc., 2021]

# $\mu$ Brain: Fully Synthesizable Mixed-signal Neuromorphic Processor

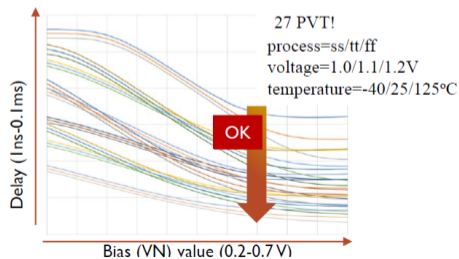
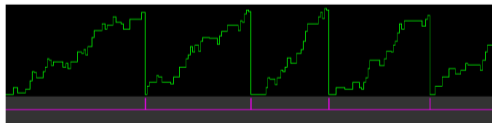
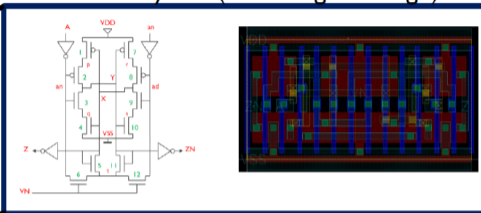
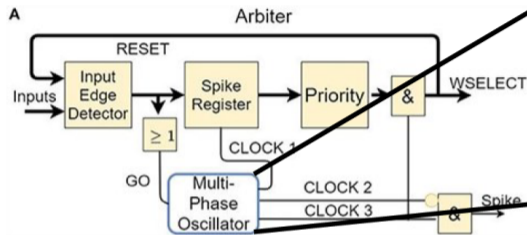


## $\mu$ Brain

- $1.68\text{mm}^2$  (60% syn, 4% neu, 36% arb)
- 300 neu, 200000 syn
- $\leq 100\ \mu\text{W}$
- Programmable, Fully parallel, Asynchronous
- High-fidelity (SW vs. HW)

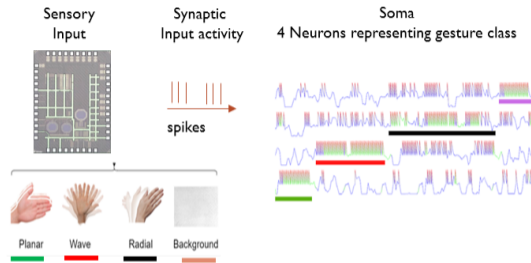
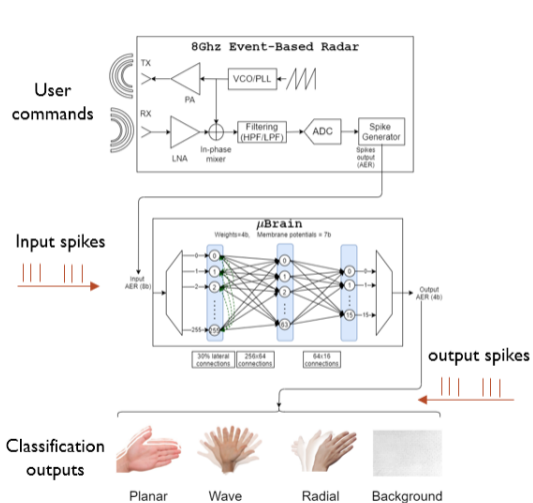
# $\mu$ Brain: Self-timed Neurons

## Neurons delay cell (mixed signal design)



[Stuijt J., ..., Corradi F., Front. Neurosc., 2021]

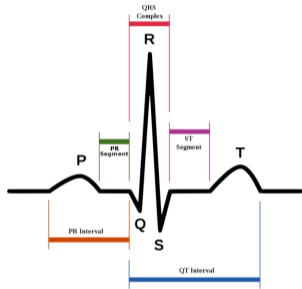
# $\mu$ Brain: Radar Signal Processing



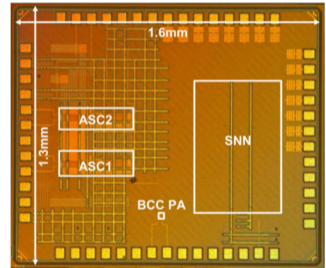
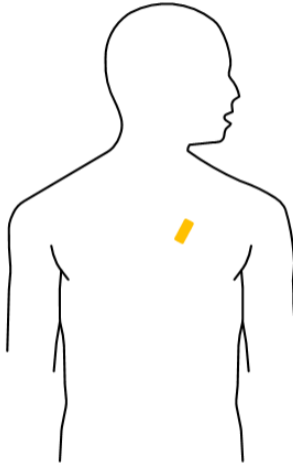
## Performance

- Always-on ( $< 100\mu\text{W}$  compute)
- $> 93\%$  accuracy
- High-fidelity SW vs HW
- Near- and in-sensor computing

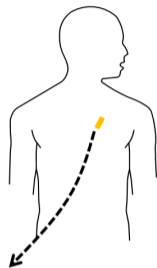
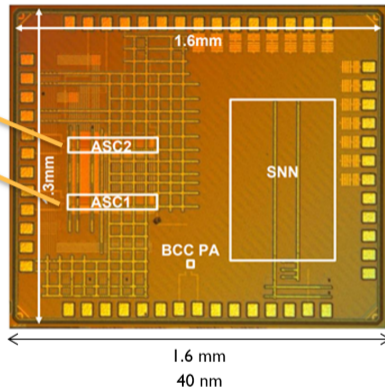
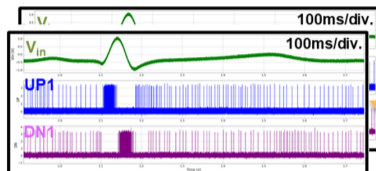
[Stuijt J., ... , Corradi F., Front. Neurosc., 2021]



Always-on monitoring  
Temporal feature extraction  
for anomaly detection



# $\mu$ Brain: Implantable Biomedical Devices

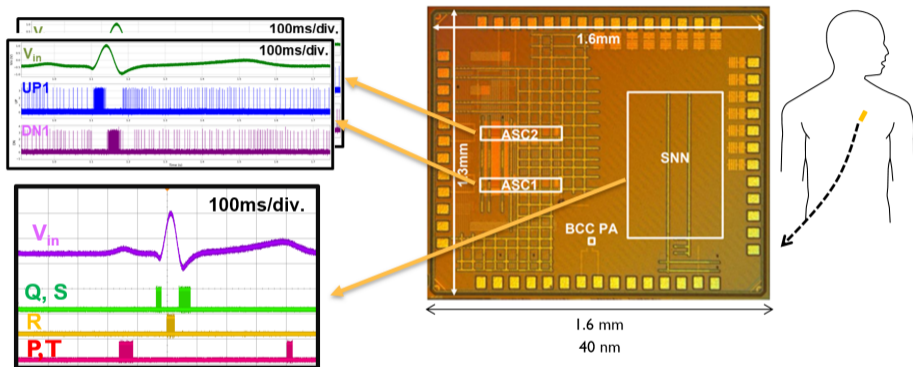


- Analog to Spike Converter: bio-inspired!
  - Lower data rate (compression)
  - Active on demand

## Re-synthesized $\mu$ Brain [He Y. et al, JSSC, 2022]

- 100 neurons, 7000 syn,  $\leq 50 \mu\text{W}$
- Analog-to-spike converters (asynchronous)

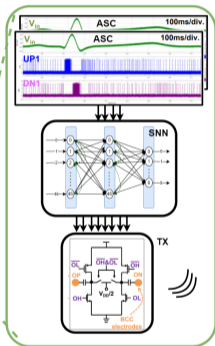
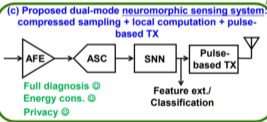
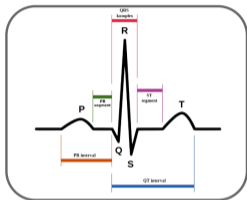
# $\mu$ Brain: Implantable Biomedical Devices



## Re-synthesized $\mu$ Brain [He Y. et al, JSSC, 2022]

- Temporal feature extraction
- With an accuracy of less than 1 millisecond!

# μBrain: Implantable Biomedical Devices



	<b>This work</b>		[2] Luo, TBIOCAS'19	[3] Kim, TBIOCAS'14	Reveal LINQ [10]	
Tech. (nm)	40		130	180	N.A.	
Supply (V)	0.9/1/1.1		0.9	1.2	N.A.	
System	ADC+SNN+TX		AFE+ADC+DSP+TX	AFE+ADC+DSP+TX	N.A.	
Applications	<b>Insertable</b>		Wearable	Wearable		Insertable
	Daily monit.	Diagnostic		Daily monit.	Diagnostic	
ADC architecture	LC		SAR	SAR		
ADC sample rate × nr. of bits (bps)	~ 140 (event-driven) × 2b =280		256 × 12b =3072	512/64 (Adaptive) × 12b =877 (7× compression)		
ADC ENOB (bits)	5	9.6	10.32	N.A.		
TX technology	Pulse based BCC		QPSK/BFSK	BLE		
Power con. (μW)						
TX	1.2	1.5	69.4~317.4 <sup>B</sup>	1000	13302	-
DSP/NN/Memory	13	15	-	42	0	-
ADC	14	34	4.6	18	38	-
AFE	-	-	4.6	21	21	-
Total (excl. AFE)	<b>28.2</b>	<b>50.5</b>	74~322	1060	13340	N.A. <sup>C</sup>
Nr. of ECG chan.	2		2	3		N.A.
Core area (mm <sup>2</sup> )	<b>0.32<sup>A</sup></b>		6.25	23.5		N.A.
System area (mm <sup>2</sup> )	<b>72</b>		N.A.	~500		322 <sup>D</sup>
On-chip labeling	<b>PQRST</b>		No	R		R
Label precision (ms)	<b>&lt;1</b>		>4 <sup>E</sup>	>2 <sup>E</sup>		>4 <sup>E</sup>

<sup>A</sup> No AFE. <sup>B</sup> 0.38% TX duty cycling. <sup>C</sup> 3-yr battery life. <sup>D</sup> Include ECG electrodes. <sup>E</sup> Limited by sampling rate.

Accurate and efficient spike-based electrocardiogram monitoring.

[He Y. et al, JSCC 2022]

## Summary

- **Diverse solutions** for different requirements (Digital, Mixed-Signals, In-Memory Computing)
  - Trade-off: Flexibility, Energy, Memory
- **Asynchronous tiny SNNs** chips for IoT applications achieve competitive performance!
  - Always-on (on-demand), ULP ( $\leq 100 \mu\text{W}$ )
  - Consume data locally, ASIC for specific tasks

## Summary

- **Diverse solutions** for different requirements (Digital, Mixed-Signals, In-Memory Computing)
  - Trade-off: Flexibility, Energy, Memory
- **Asynchronous tiny SNNs** chips for IoT applications achieve competitive performance!
  - Always-on (on-demand), ULP ( $\leq 100 \mu\text{W}$ )
  - Consume data locally, ASIC for specific tasks

## What's next?

- **Scaling** to large SNN models.
- **Heterogenous** computing architectures.
- **In-memory computing** with emerging technologies.

# Neuromorphic Edge Computing Systems: Summary

## Summary

- **Diverse solutions** for different requirements (Digital, Mixed-Signals, In-Memory Computing)
  - Trade-off: Flexibility, Energy, Memory
- **Asynchronous tiny SNNs** chips for IoT applications achieve competitive performance!
  - Always-on (on-demand), ULP ( $\leq 100 \mu\text{W}$ )
  - Consume data locally, ASIC for specific tasks

## What's next?

- **Scaling** to large SNN models.
- **Heterogenous** computing architectures.
- **In-memory computing** with emerging technologies.

## Collaborators

- Sander Bothe (CWI)
- Yao-Hong Liu (IMEC)
- Manil Dev Gomony (TU/e)

## Students

- Yuming He - PhD (IMEC)
- Rick Luiken - PhD (TU/e)
- Bojian Yin - PostDoc (TU/e)



## Neuromorphic Edge Computing Systems Lab

**Federico Corradi**

Eindhoven University of Technology

f.corradi@tue.nl, +31(0)402 472 556

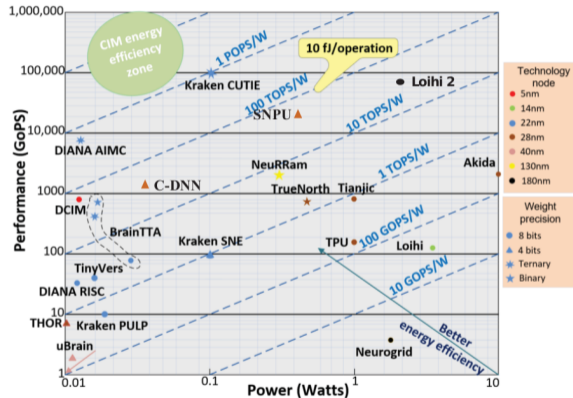
Neuromorphic Edge Computing Systems Lab



# SOTA in Edge-AI Processor/Accelerator HW (ANNs and SNNs)

## SOTA

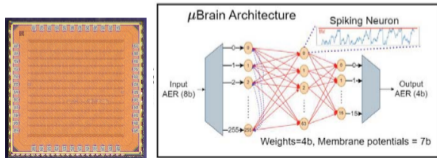
- Mostly Artificial Neural Nets (ANN)
- Some Spiking Neural Nets (SNN)
- Weight Precision (8/4/Ternary/1)
- Technology Nodes (5-180nm)



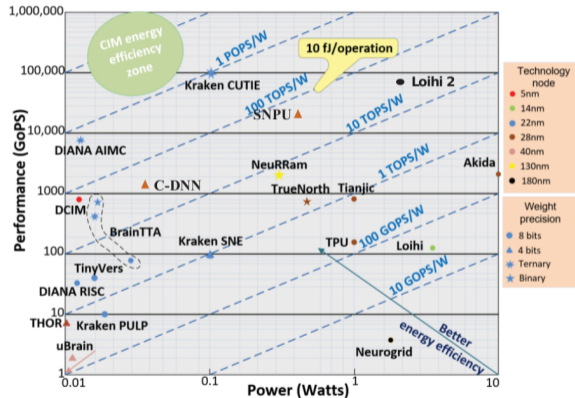
[Gomony D. M. et al, 2023]

## $\mu$ Brain (IMEC)

- 40 nm,  $1.4\text{mm}^2$
- Event based architecture
- Asynchronous design (no clock)
- Without schedules, clocks, state machines
- Extreme low power, not en-eff.



[Stuijt J. et al, Front. Neurosc., 2021]

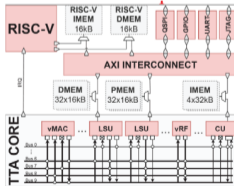


[Gomony D. M. et al, 2023]

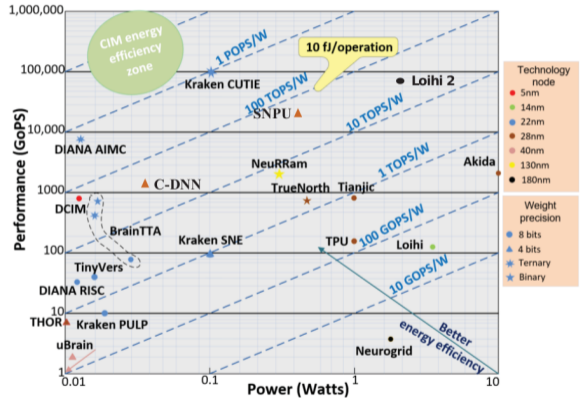
# BrainTTA - Flexible & Mixed-precision (TU/e)

## BrainTTA (TU/e)

- TTA-based accelerator, 22nm
- Fully-programmable (C-compiler)
- SNPU
- Flexible precision:
  - INT8, ternary, binary
  - 405 / 67 / 35 fJ/op



[Molendijk M. et al, 2023]

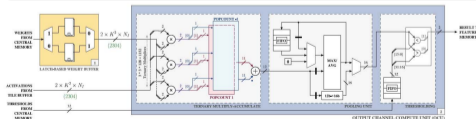


[Gomony D. M. et al, 2023]

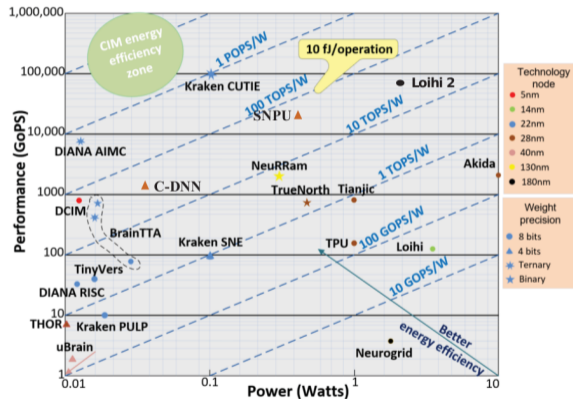
# CUTIE: Ternary DNN (ETH)

## CUTIE (ETH)

- ANN: Ternary (-1,0,1) Inf.
  - exploit zero suppression
  - ternary compression
- 1.6 bit / symbol
- complete unrolled loops in HW
- $3 \times 3 \times 128 = 1152$  conv. size
- 22 nm, 2.5-4 fJ/op (%Sparsity dependent)



[Scherer M, et al., 2021]

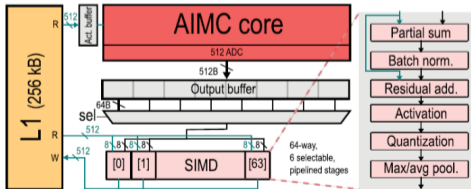


[Gomony D. M. et al, 2023]

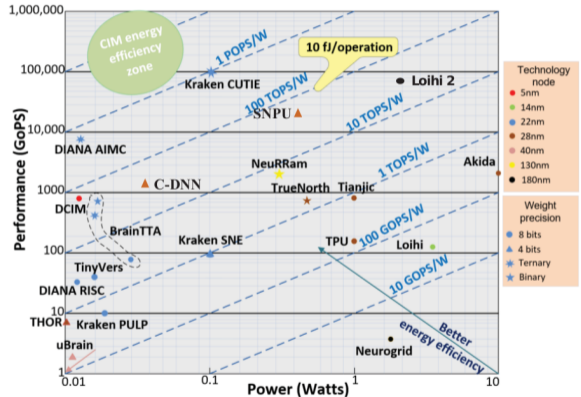
# DIANA: Mixed-Signal In-memory Compute (IMEC)

## DIANA (IMEC)

- Precision-scalable digital core
- Analog In-Memory Core (AIMC)
  - INT 8/4/2
  - Ternary weights, 7-b activation
  - Programmable SIMD
  - 22nm 1.7 fJ/Op (I/W/O= 7/1.5/6-bit)



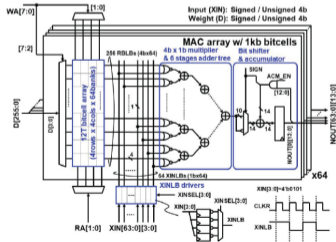
[Houshmand P. et al. , 2022]



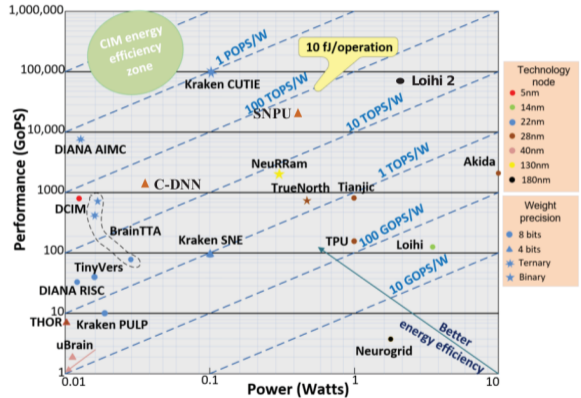
[Gomony D. M. et al, 2023]

## Digital (CIM)

- 12T bitcell based architecture
- Flexible precision
- INT 8/4, 5nm, 15.8 / 3.9 fJ/Op
- MAC using 1-b at a time



[Fujiwara H., et al, 2022]

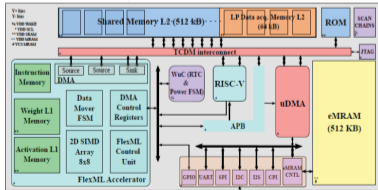


[Gomony D. M. et al, 2023]

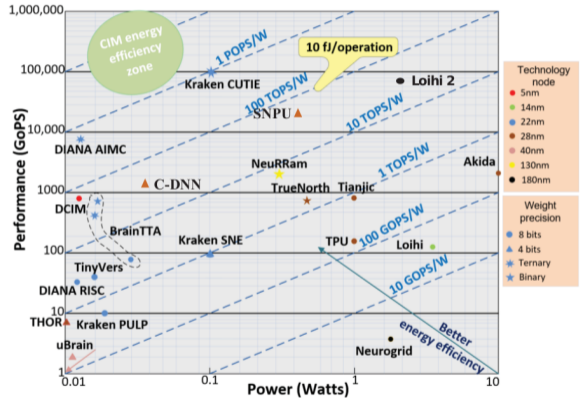
# TinyVers - Embedding MRAM

## TinyVers

- Supports various DNN layers, to traditional ML models like SVM
- RISC-V + ML Acc.
- Flexible-precision scalable digital accelerator, max 17 TOPS/W  $\sim$  59 fJ/Op (Int2)



[Vikram J et al, 2022]

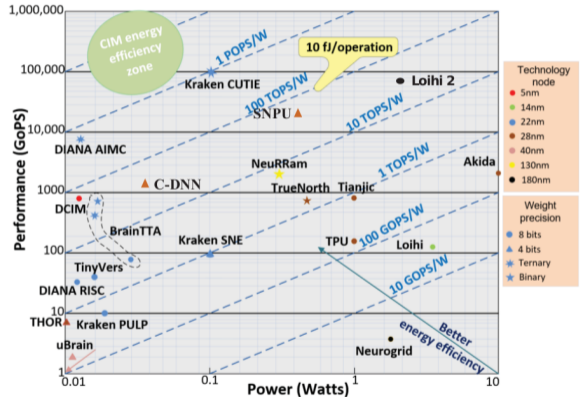


[Gomony D. M. et al, 2023]

# SOTA in Edge-AI Processor/Accelerator HW: Summary

## SOTA

- ANN chips approaching 1 fJ/Op
  - However:
    - Requires complete unrolling (CUTIE) and/or AIMC (DIANA)
    - System overhead often neglected
- Flexibility has its price:
  - E.g. BrainTTA suffers at least one-order in energy efficiency
  - DRAM overhead not included
- SNNs have high potential
  - However:
    - at best in the 35 fJ/Op range
    - need to catch-up



[Gomony D. M. et al, 2023]

# Interested? Reach out!

## Collaborators

- Henk Corporaal (TU/e)
- Manil Dev Gomony (TU/e)



CONVOLVE



imec

## Neuromorphic Edge Computing Systems Lab

**dr. Federico Corradi**

Eindhoven University of Technology

f.corradi@tue.nl, +31(0)402 472 556

Neuromorphic Edge Computing Systems Lab

