

# How to build a green AI?

Adrien F. Vincent

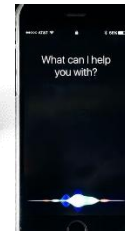
Hardware Artificial Intelligence group

IMS Lab, University of Bordeaux

[adrien.vincent@ims-bordeaux.fr](mailto:adrien.vincent@ims-bordeaux.fr)

# Energy consumption of AI

Artificial Intelligence (AI) may be found in many applications.



- First version: 250 kW
  - One game lasts 2 h
- Second version: 10 kW



~ 850 km!



20h!

# Energy consumption of AI – cont.



- Vehicle battery: 52kW·h (Zoe in 2021)
- Drone Matrice 600 Pro
  - Weight 10kg
  - Payload 6kg during 38'

## ➤ ICT

- 7% of world electricity consumption
- Forecast for digital electronics: 20% or maybe 50% in 2030

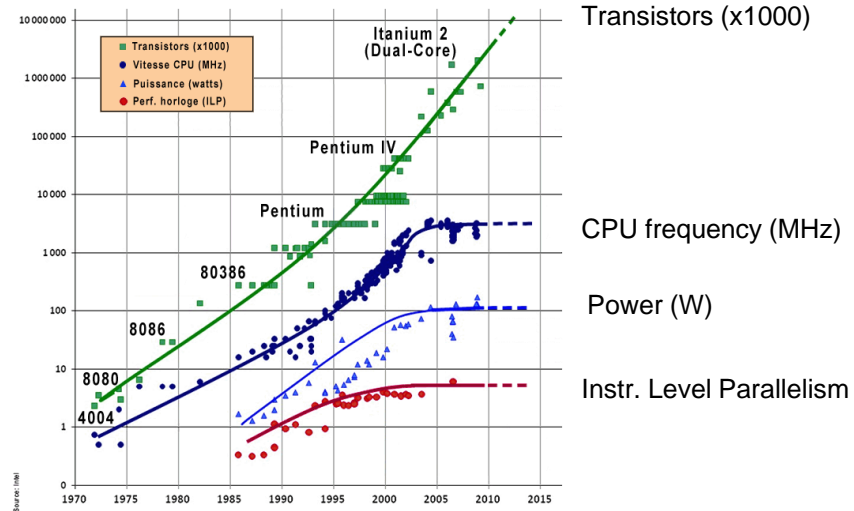


**May not be sustainable!**

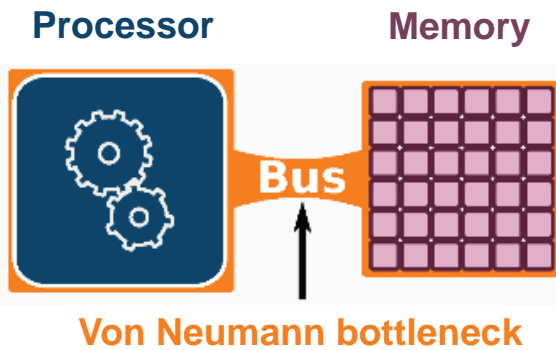
# Processors are facing physical limits

## ➤ Thermal dissipation

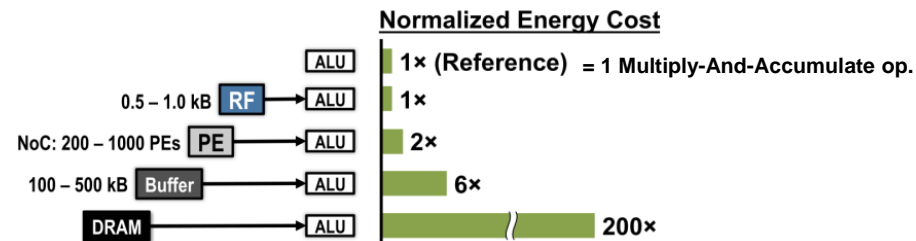
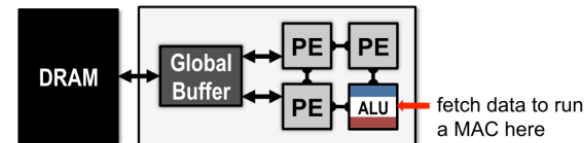
■ 100 W/cm<sup>2</sup>



## ➤ Memory access



**Moving data is expensive!**



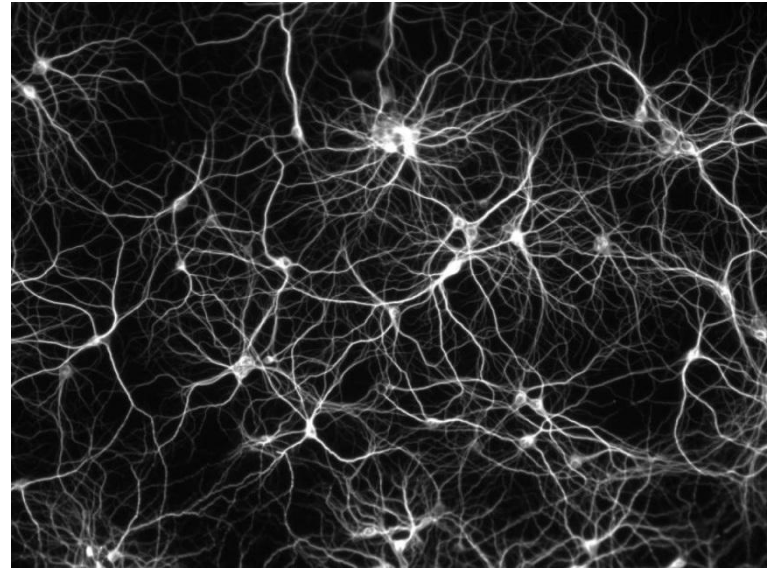
V. Sze & al., Proc. of IEEE, 2017

# An alternative paradigm of computation

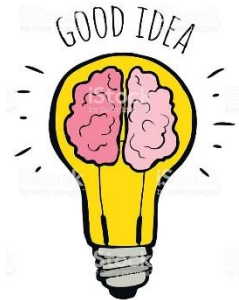
- Take inspiration from biological brains

Human brain:

- ~100 billions of neurons
- ~1000 synapses per neuron
- Neuron frequencies ~10Hz to ~100Hz
- Massively parallel computation
- ~20W
- Efficient for high-level tasks (translation, recognition, synthesis)



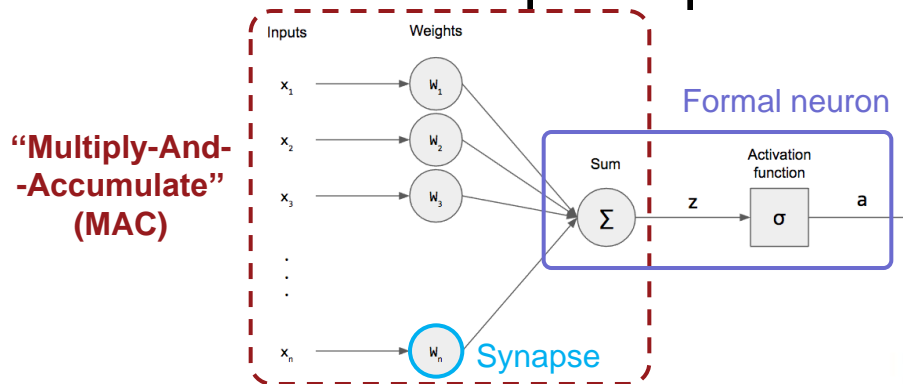
<https://wccftech.com/scientists-artificial-neurons-mimics-human-brain-cells/>



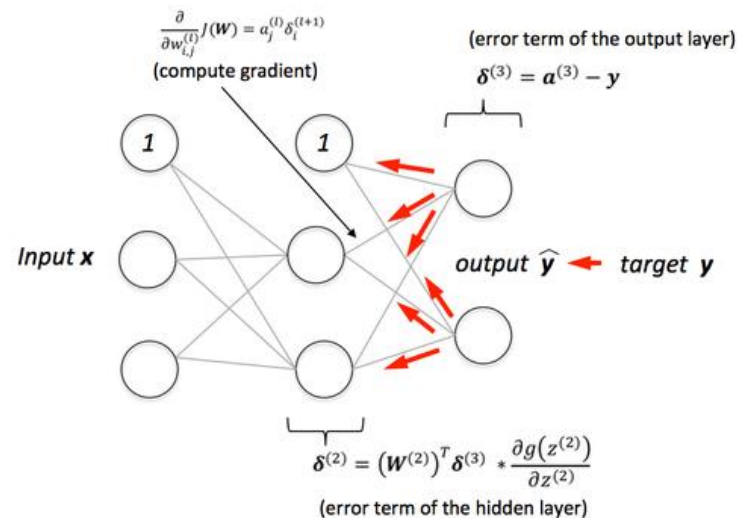
**Original computing architectures to explore!**

# Artificial (and Deep) Neural Networks

- Based on the perceptron concept



- Supervised learning

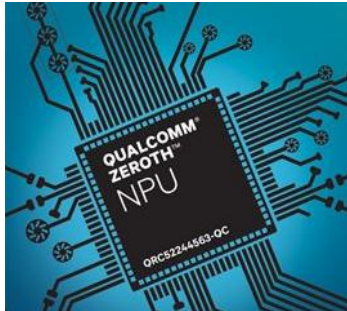


- OK but implemented on Von Neumann computers...



# Dedicated neural network processors

**Qualcomm (2013) :  
Zeroth**



**IBM (2014) :  
TrueNorth**

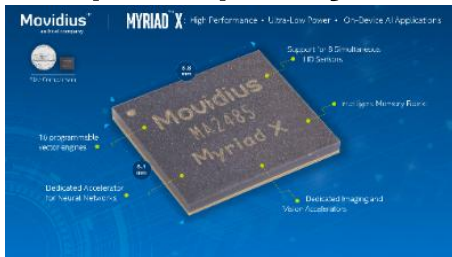


**Spiking neural  
networks**

**Intel (2017) : Loihi**



**Intel (2017) : Myriad**



**Google (2016) : TPU**



**IC purely CMOS  
Learning is not easy**

# Ideas for building a green(er) AI

Let us have a look at “unconventional computing!”

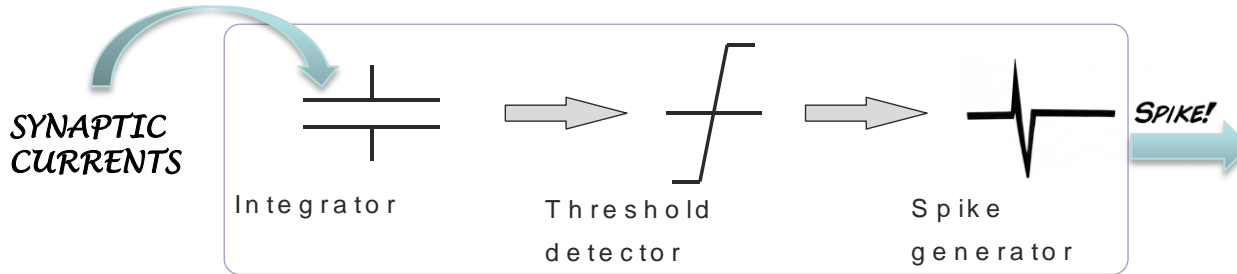
- 1) Event-based computing with spiking neural networks
- 2) Memristor-based disruptive hardware implementations
- 3) Smart System Integration: hardware AI-enhanced sensor  
(ULPEC EU H2020 project)
- 4) Radio-frequency processing with spintronic nanodevices  
(RadioSpin EU H2020 Project)



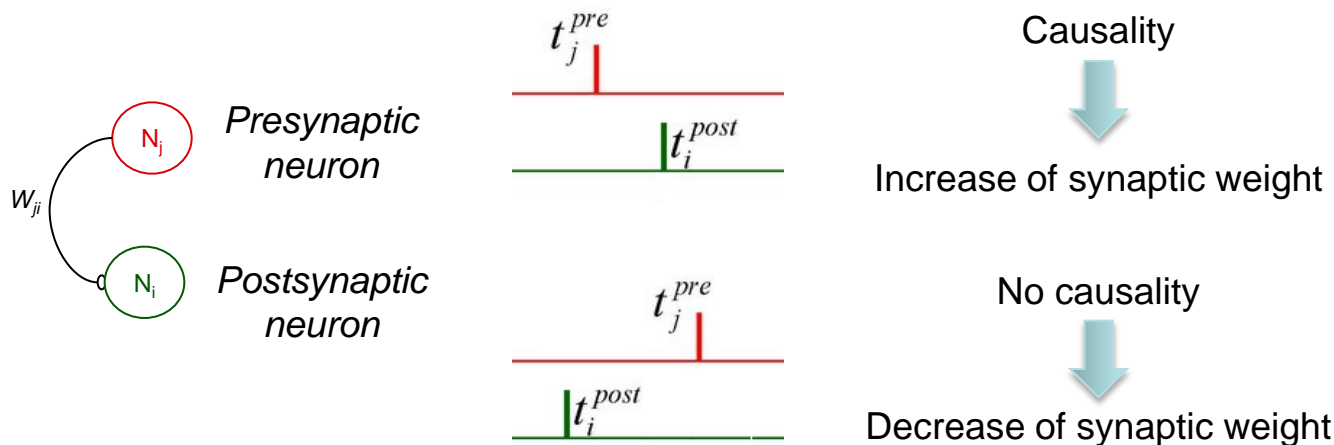
# Event-based computing with spiking neural networks

# Event-based computing & learning rule

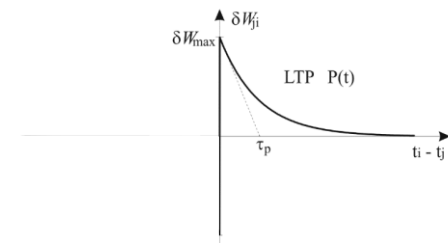
- Spiking neuron = time dependent



- Spike-Timing-Dependent Plasticity (STDP)



**Local and unsupervised learning rule!**



# CMOS implementation of neuromorphic systems



## Neuromorphic silicon neuron circuits

**Giacomo Indiveri<sup>1\*</sup>, Bernabé Linares-Barranco<sup>2</sup>, Tara Julia Hamilton<sup>3</sup>, André van Schaik<sup>4</sup>,  
Ralph Etienne-Cummings<sup>5</sup>, Tobi Delbruck<sup>1</sup>, Shih-Chii Liu<sup>1</sup>, Piotr Dudek<sup>6</sup>, Philipp Häfliger<sup>7</sup>, Sylvie Renaud<sup>8</sup>,  
Johannes Schemmel<sup>9</sup>, Gert Cauwenberghs<sup>10</sup>, John Arthur<sup>11</sup>, Kai Hynna<sup>11</sup>, Fopelola Folowosele<sup>5</sup>,  
Sylvain Saighi<sup>8</sup>, Teresa Serrano-Gotarredona<sup>2</sup>, Jayawan Wijekoon<sup>6</sup>, Yingxue Wang<sup>12</sup> and Kwabena Boahen<sup>11</sup>**

<sup>1</sup> Institute of Neuroinformatics, University of Zurich and ETH Zurich, Zurich, Switzerland

<sup>2</sup> National Microelectronics Center, Instituto Microelectronica Sevilla, Sevilla, Spain

<sup>3</sup> School of Electrical Engineering and Telecommunications, University of New South Wales, Sydney, NSW, Australia

<sup>4</sup> School of Electrical and Information Engineering, University of Sydney, Sydney, NSW, Australia

<sup>5</sup> Whiting School of Engineering, Johns Hopkins University, Baltimore, MD, USA

<sup>6</sup> School of Electrical and Electronic Engineering, University of Manchester, Manchester, UK

<sup>7</sup> Department of Informatics, University of Oslo, Oslo, Norway

<sup>8</sup> Laboratoire de l'Intégration du Matériau au Système, Bordeaux University and IMS-CNRS Laboratory, Bordeaux, France

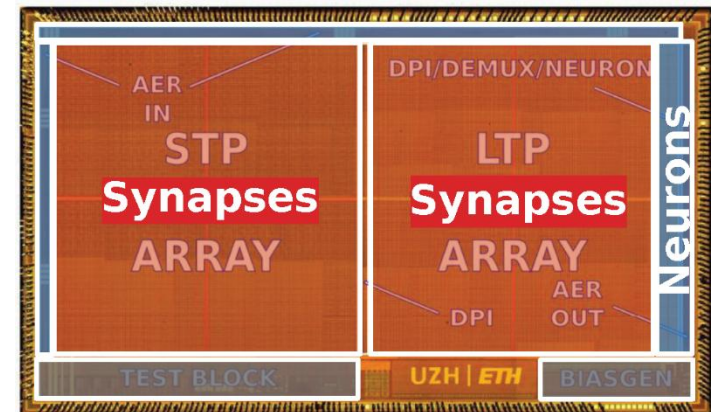
<sup>9</sup> Kirchhoff Institute for Physics, University of Heidelberg, Heidelberg, Germany

<sup>10</sup> Department of Bioengineering and Institute for Neural Computation, University of California San Diego, La Jolla, CA, USA

<sup>11</sup> Stanford Bioengineering, Stanford University, Stanford, CA, USA

<sup>12</sup> Janelia Farm Research Campus, Howard Hughes Medical Institute, Ashburn, VA, USA

**Integrable and learning-capable artificial synapses are still a challenge.**

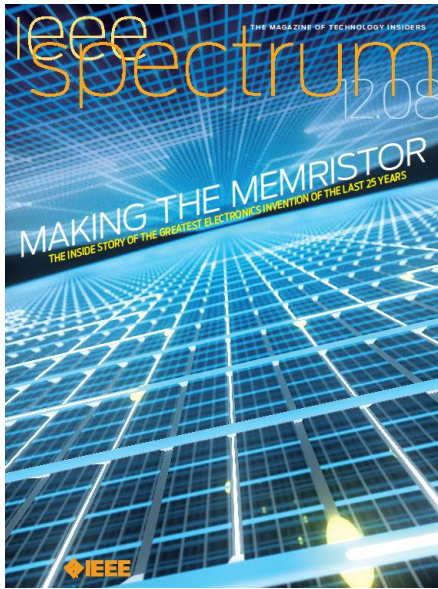


131 072 synapses (~ 65 % surface)    ROLLS [QIAO, 2015]  
256 neurons    (~ 3 % surface)

N. Qiao & al., *Frontiers in Neuroscience*, 2015

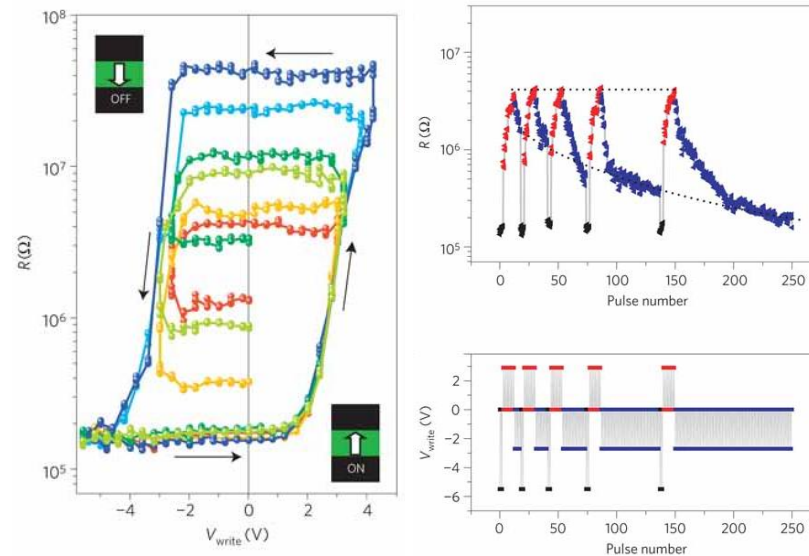
# Memristor-based disruptive hardware implementations

# Outbreak of memristors



- HP's experimental breakthrough in 2008
- What is a “memristor?”
  - Nanodevice ~ 10s nm x 10s nm
  - “Memory resistor,” i.e., a “resistor that learns.”
  - Behavior ~ like a biological synapse...

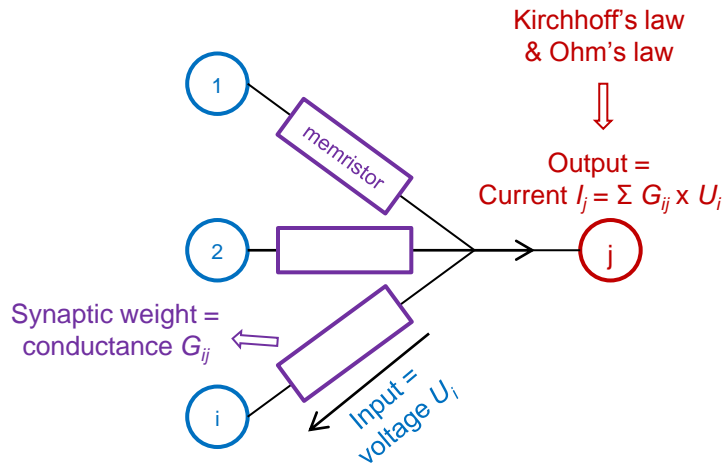
- Non volatile = memory
- Intrinsically plastic
- Several technologies (filamentary, phase change memory, spintronics, ferroelectric...)



A. Chanthbouala & al., Nature Materials, 2012

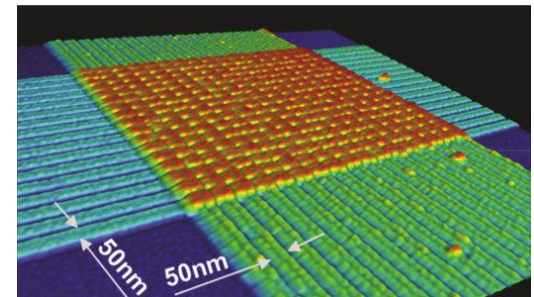
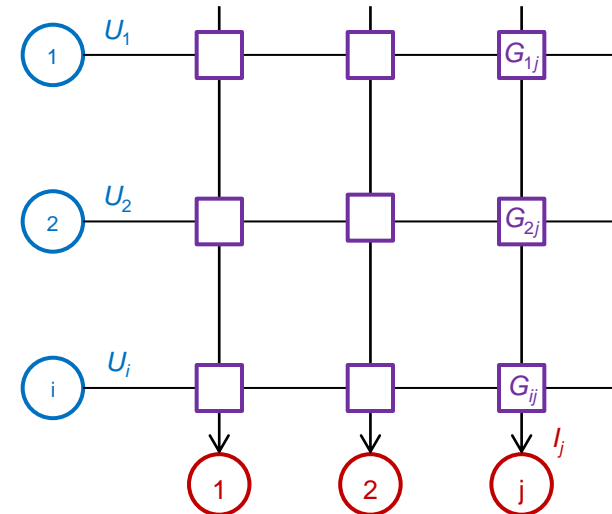
# Memristors can help for the inference

## Analog Multiply-And-Accumulate operation



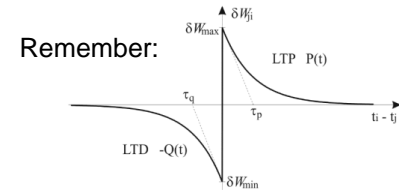
**The physics is doing the computing for us!**

## Memristive crossbar array $\Rightarrow$ Matrix-by-Vector Multiplication



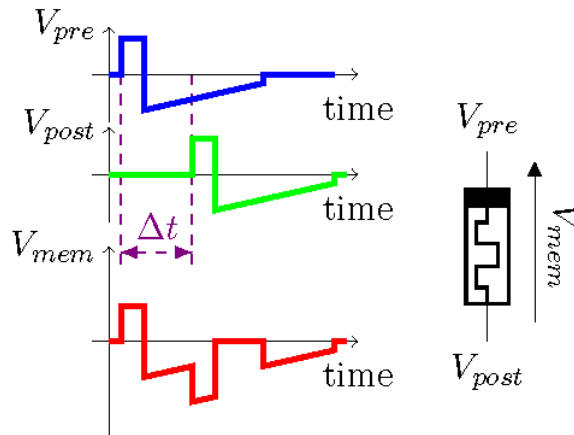
Y. Yang & al., *Appl. Phys. Lett.*, 2012

# Memristors can also help for the learning

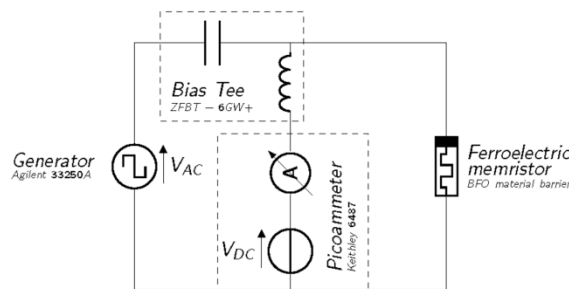


## STDP-like behavior in ferroelectric memristors

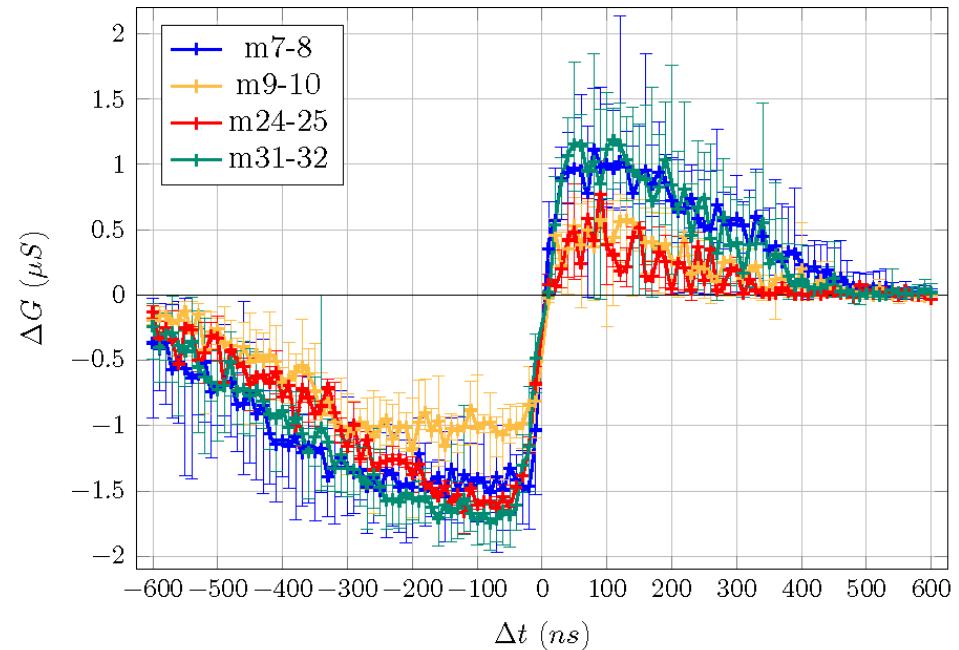
Voltage waveforms



Experimental setup



Measurements

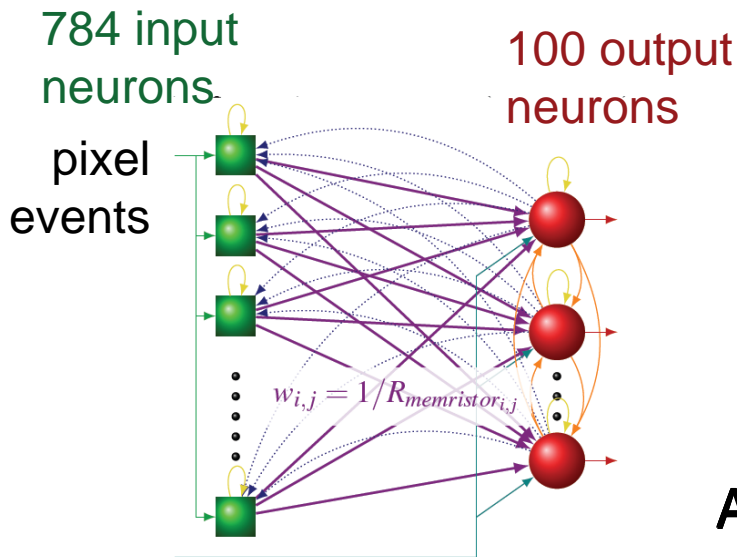
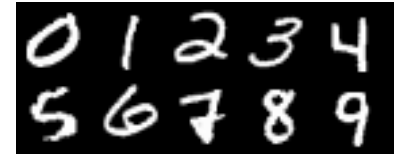


**Exploiting the intrinsic physics of the device to implement the learning rule.**



# Memristor-based learning: proof of concept

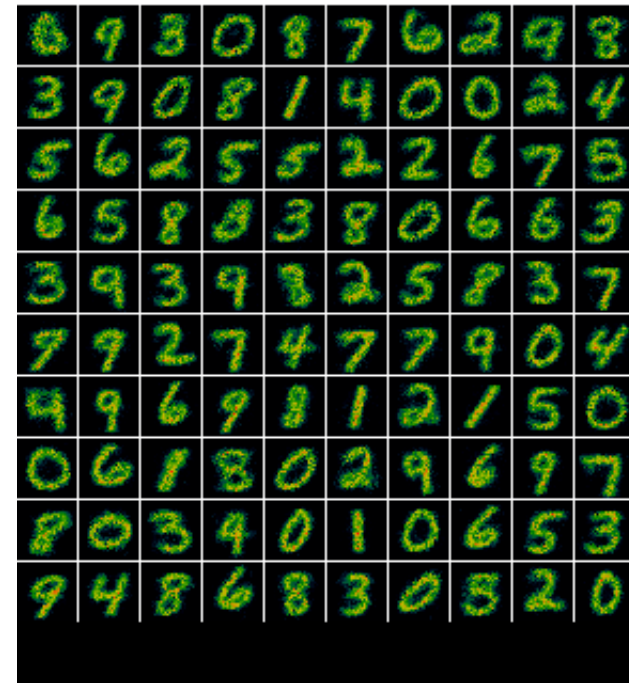
- M-NIST dataset  
(60k pictures for learning and 10k pictures for testing)
- 28 x 28 pixels = 784 input neurons
- 100 output neurons
- Memristor variability is taken into account
- STDP learning rule



- 100 boxes = 100 neurons
- 1 pixel = 1 synapse

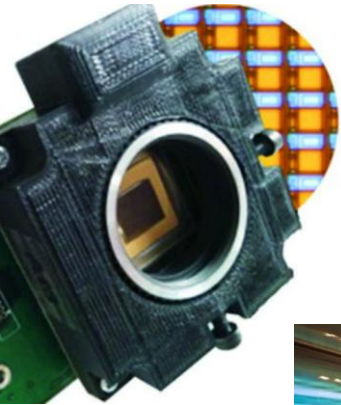
After 60000 pictures

Simulation results  
(synaptic conductance maps)

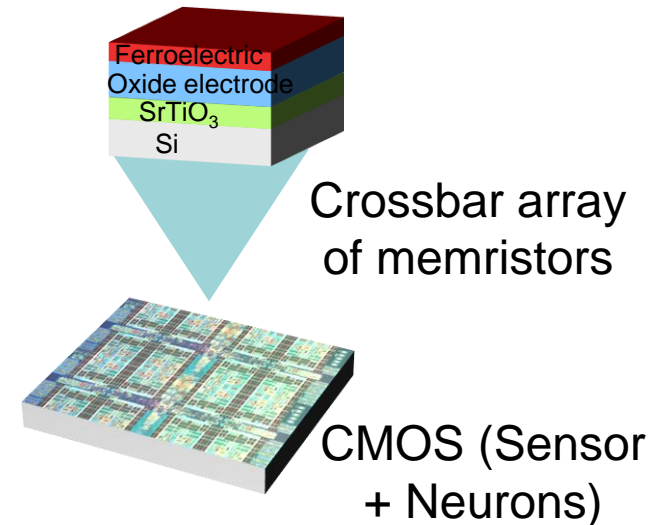


# Smart System Integration: hardware AI-enhanced sensor (ULPEC EU H2020 project, 2017-2021)

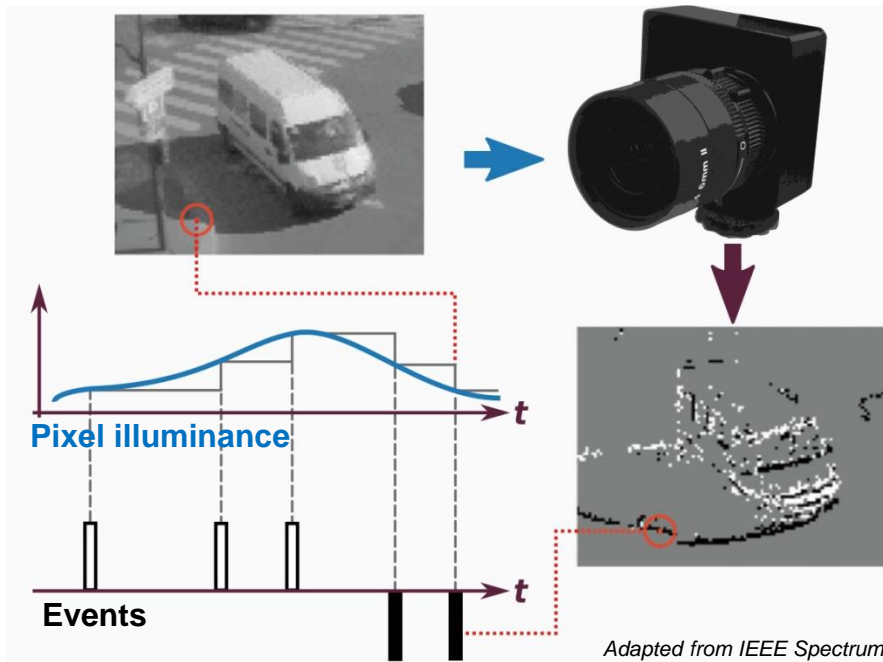




- **ULPEC : Ultra-Low Power Event based Camera**
  - Smart System Integration call
  - From sensor to decision taking for driving aid



# Event-based camera: principle & main features



- Event temporal accuracy: sub-millisecond
- Pixel individual event rates: 0 to tens of kHz
- High dynamic range response: >120dB
- Intrinsic data compression

<https://www.prophesee.ai/event-based-vision-applications>  
Example of high-speed detection and tracking for automotive

**Detection of changes:  
outputs sparse data => low power consumption.**

# Learning can be difficult in real world...

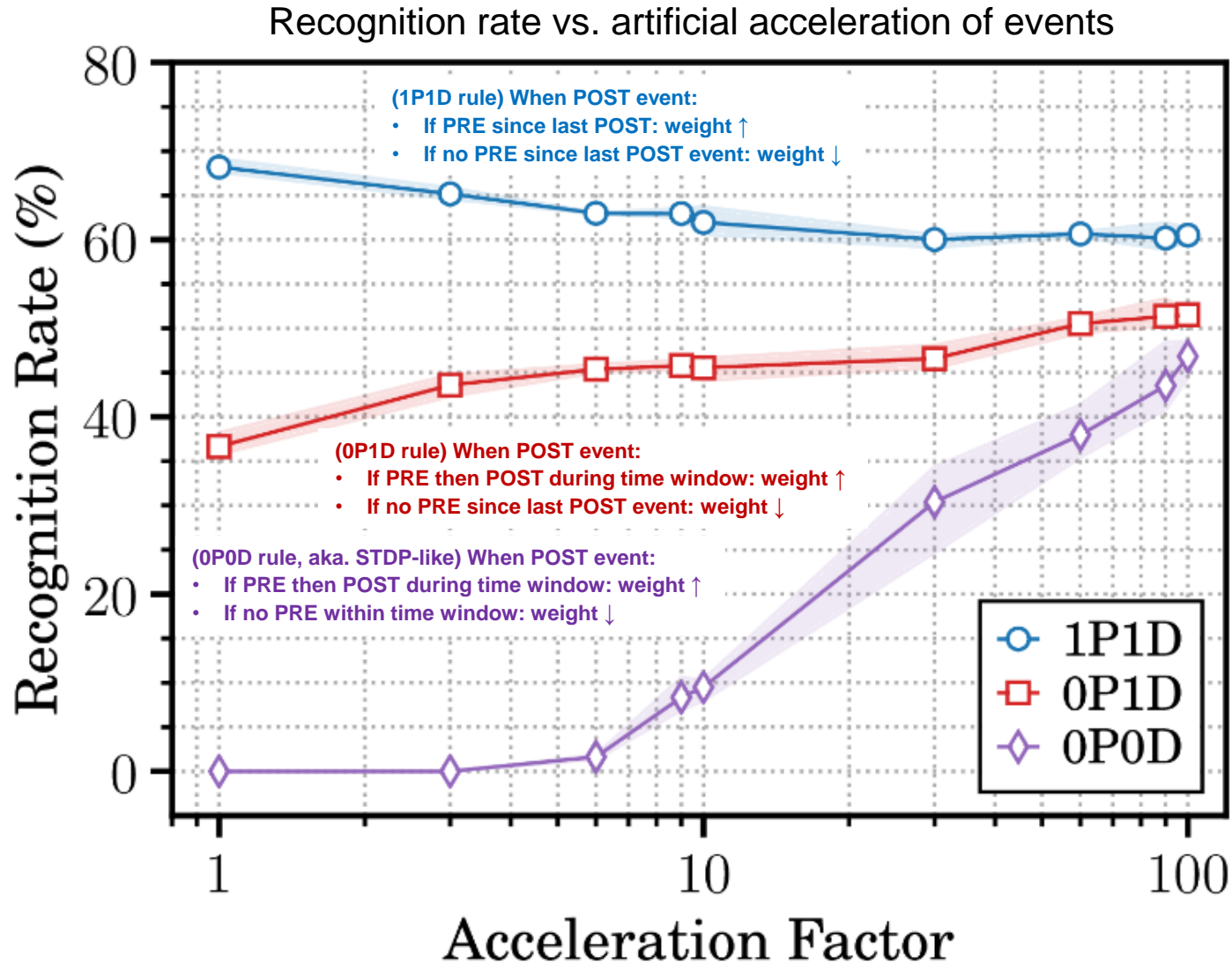
If one performs simulations with:

- MNIST inputs filmed with an event-based camera => N-MNIST
  - Camera movement in front of digits (100 ms-saccades)
- STDP-like rule:
  - Time window for STDP: 10  $\mu$ s (due to analog design constraints)
  - Presynaptic event and then a postsynaptic one in the time window: synaptic weight  $\uparrow$
  - Postsynaptic event without a presynaptic one in the time window: synaptic weight  $\downarrow$

**0 % recognition rate**

**Not enough events in the STDP time window for learning.**

# One may need to use other learning rules



P. Lewden & al.,  
IJCNN, 2020

# Radio-frequency processing with spintronic nanodevices

(RadioSpin EU H2020 Project, 2021-2025)





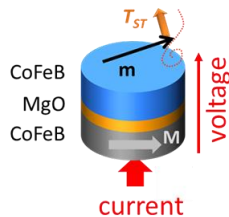
# Magnetic Tunnel Junction (MTJ)

## A multifunctional spintronic nanodevice

DC: "direct current," i.e., 0 Hz  
RF: radio-frequency

*N. Leroux & al., Physical Review Applied 15, 034067, 2021*

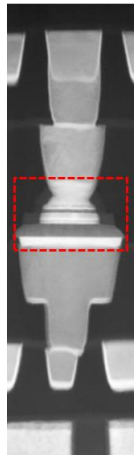
Magnetic tunnel junction



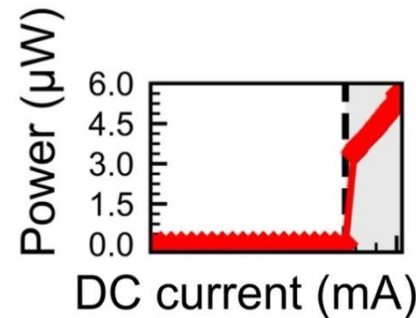
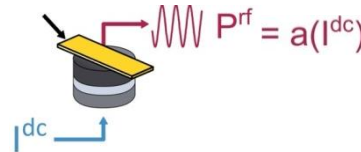
Non-volatility  
Endurance  
Reproducibility

*J. Grollier & al., Nature Electronics 3, 360, 2021*

Samsung IEDM  
2018, 28nm  
Logic



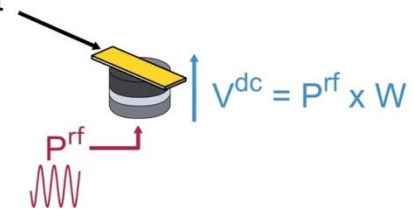
Spin-torque nano-oscillator



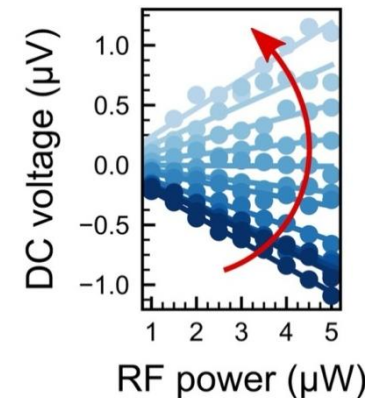
Neuron  
(DC → RF)

Spin-diode

Tuning weight  
by DC current



Tunable  
weight

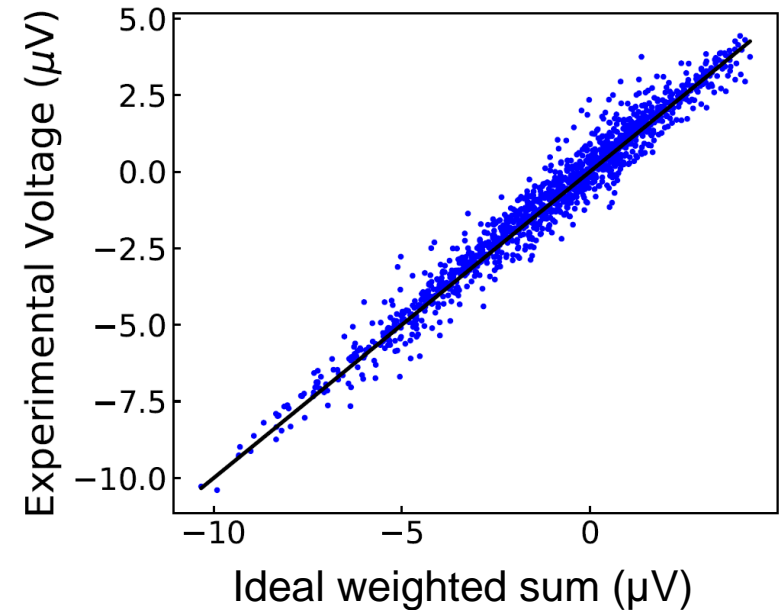
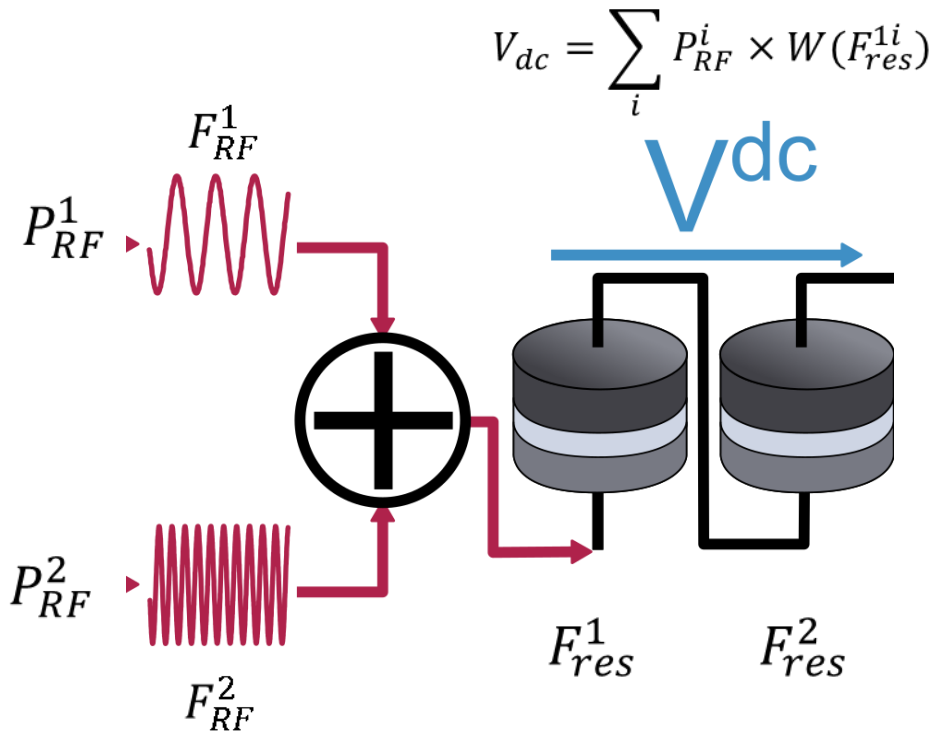


Synapse  
(RF → DC)

**The same technology can be used both for neurons and for synapses.**

# Combining several synaptic MTJs

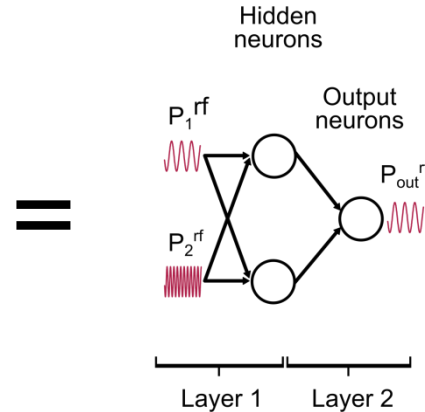
N. Leroux & al., Neuromorph. Comput. Eng. 1 011001, 2021



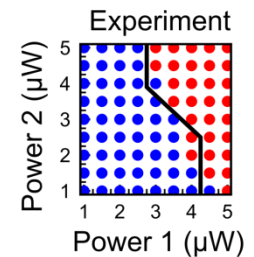
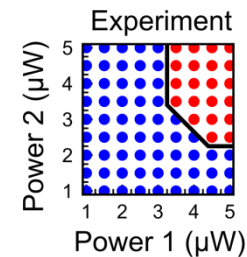
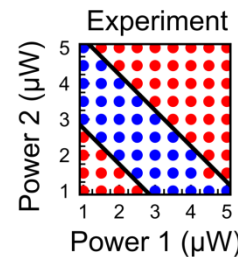
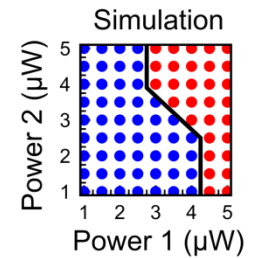
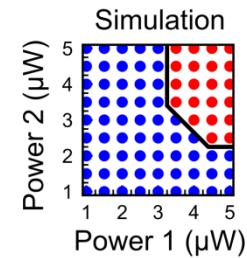
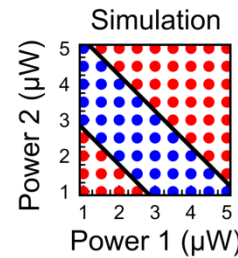
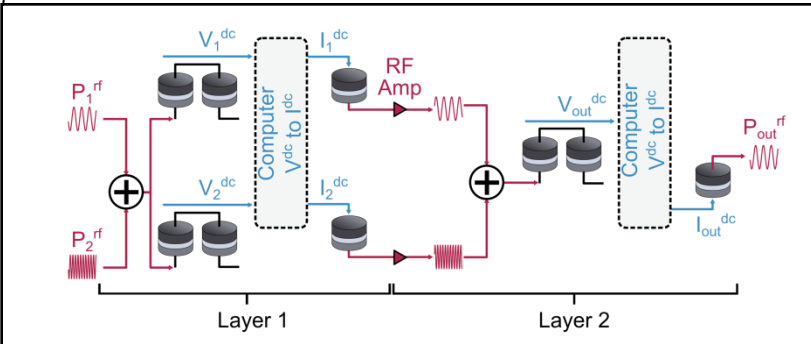
**Multiply-And-Accumulate operations on RF signals are possible in hardware!**

# Fully spintronic hardware neural network

A. Ross & al., ArXiv 2211.03659, under review, 2022

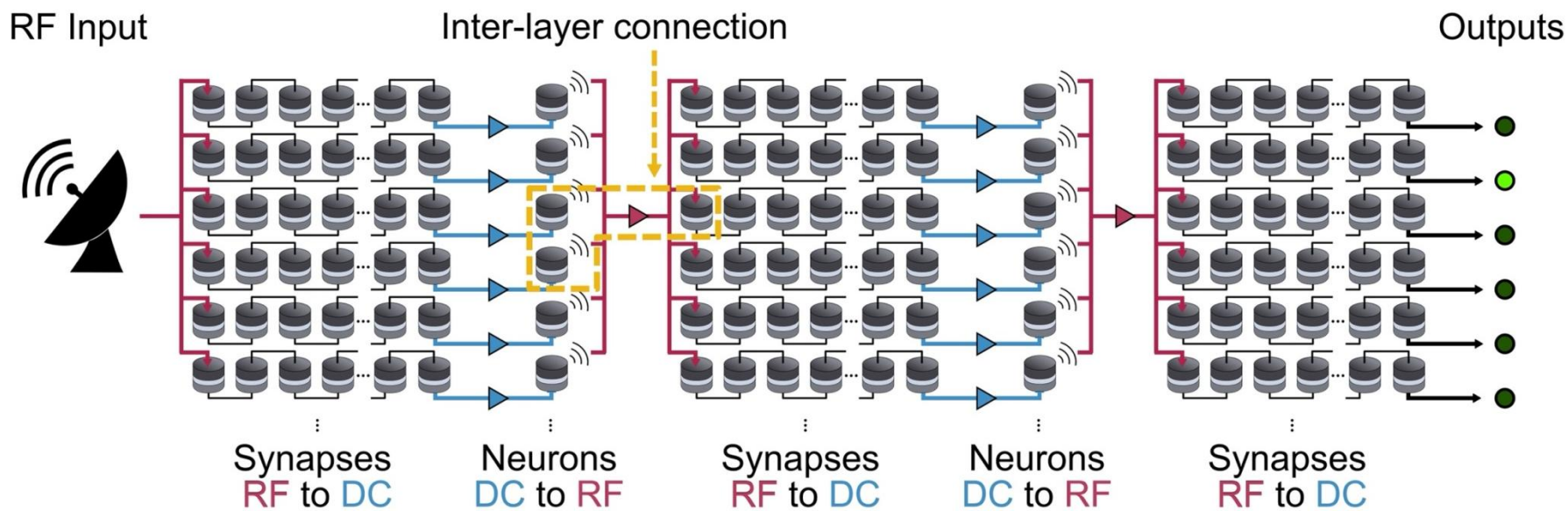


**97.7 % accuracy**



# Towards deep RF spintronics neural networks

A. Ross & al., ArXiv 2211.03659, under review, 2022



**Energy consumption estimate:**  
~ 10 fJ/synapse and ~ 100 fJ/neuron for MTJs with 20 nm diameter:  
**100 fold energy gain compared to GPUs.**

# Conclusion

# Conclusions and outlook

- Need for energy sustainability of AI before great disillusion.
- Event-based computation (neuromorphic) is a good candidate: mainly for edge-computing.
- CMOS will be not enough: unconventional nanoelectronics is likely to play a major role in Green AI.
- It will likely be a long way before industrial products... however that is the purpose of research!
- In any case: best energy savings are not generating irrelevant data!

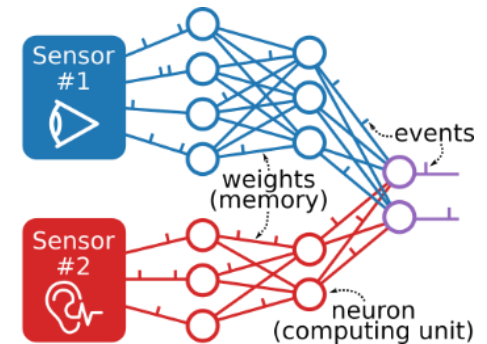
# Currently in Bordeaux: Green AI project

- GrAI: chair on energy-efficient hardware AI supported by the French National Research Agency



- Objective: designing better sensors and computing architectures for low-power spiking neural networks.

- Building new event-based sensors
- Optimize spiking neural network architectures
- Study sensor fusion
- Digital or analog fashion



- Led by Sylvain Saïghi ([sylvain.saighi@ims-bordeaux.fr](mailto:sylvain.saighi@ims-bordeaux.fr))



# To go further... at the French level

- GDR BioComp
  - Hardware implementation of natural computation
  - 5 CNRS institutes

[www.gdr-biocomp.fr](http://www.gdr-biocomp.fr)

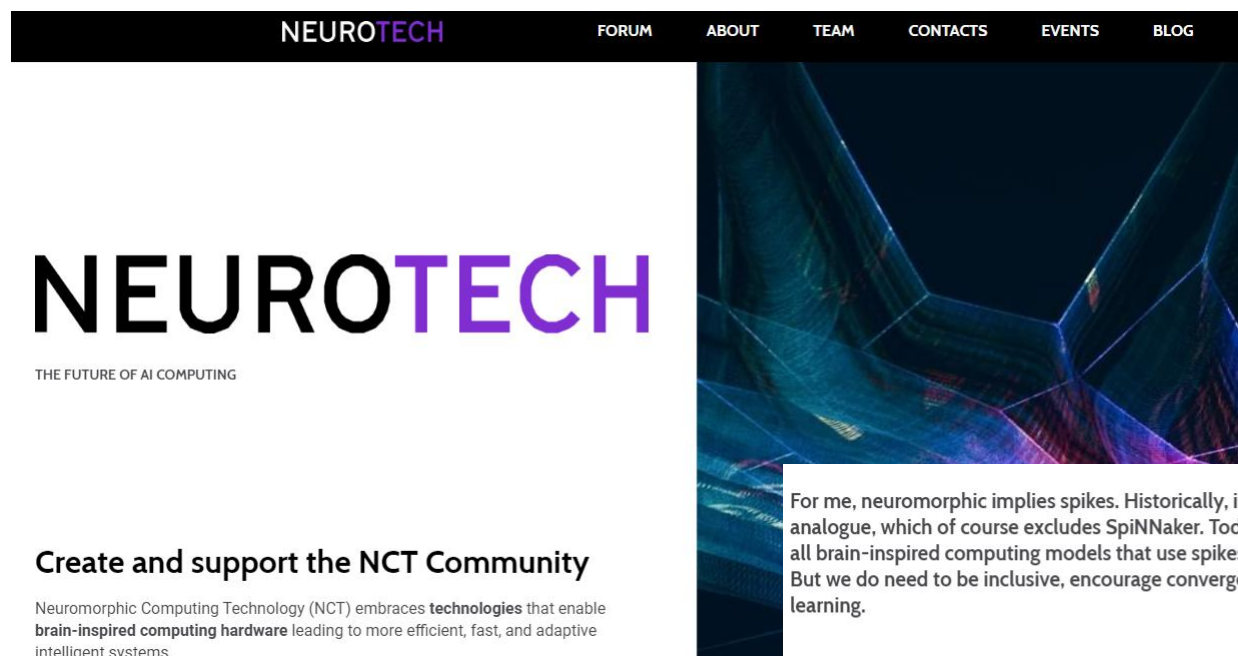
The screenshot shows the GDR BioComp website with the following content:

- Header:** GDR BioComp *implémentations matérielles du calcul naturel*
- Navigation:** Accueil, Le GDR BioComp, Colloques, Journées thématiques, Annonces, Inscription, Liens utiles
- Section: GDR BIOCOMP**
- Workshop Announcement:**
  - processing**
  - EVÉNEMENTS:** 24 NOV., 2015
  - Workshop Dynamical systems and brain-inspired information processing**
  - Description:** Le workshop Dynamical systems and Brain-inspired Information Processing organisé en Novembre par le laboratoire de mathématiques de Besançon, et soutenu par le GDR BioComp a été un grand succès. La suite en Belgique en...
- Participants List:**
  - Bruno Garbin (Université Nonlinéaire de Nice)
  - Julie Grollier (CNRS/Thales)
  - Michiel Hermans (Université Libre de Bruxelles)
  - Lars Keuninckx (Vrije Universiteit Brussel)
  - Andrew Katumba (Université Gent)
  - Serge Massar (Université Libre de Bruxelles)
  - Simon Morando (FEMTO-ST, Beloit)
  - Michel Salomon (FEMTO-ST, Besançon)
  - Ilia Degenova
- Image:** A photograph of a complex electronic circuit board with various components and glowing green LEDs.

On the right side of the screenshot, there is a sidebar with social media links and a 'FOLLOW:' section containing tweets and a photo of a group of people.

# To go further... at the European level

- NeuroTech EU project
  - Aims to gather all academic and industrial NCT stakeholders
  - Started on Fall 2019 for 3 years (just finished!)



NEUROTECH

FORUM ABOUT TEAM CONTACTS EVENTS BLOG

# NEUROTECH

THE FUTURE OF AI COMPUTING

## Create and support the NCT Community

Neuromorphic Computing Technology (NCT) embraces **technologies** that enable **brain-inspired computing hardware** leading to more efficient, fast, and adaptive intelligent systems.

For me, neuromorphic implies spikes. Historically, it is restricted to sub-threshold analogue, which of course excludes SpiNNaker. Today I think it is wider, and covers all brain-inspired computing models that use spikes for primary communication. But we do need to be inclusive, encourage convergence with ANNs and machine learning.

[www.neurotechai.eu](http://www.neurotechai.eu)



– Steve Furber

# Thank you for your attention!

[adrien.vincent@ims-bordeaux.fr](mailto:adrien.vincent@ims-bordeaux.fr)

This work was supported by the European Union's Horizon 2020 research and innovation programme under grant RadioSpin No 101017098.

This work has received funding from the European Union's Horizon 2020 research and innovation programme under grant agreement 732642 (ULPEC project).

This work was supported by a grant overseen by the French National Research Agency (ANR) as part of the "Chaires IA" Program (GrAI Project, ANR-19-CHIA-0003).

Financial support from the French Agence Nationale de la Recherche (ANR) through MIRA project is acknowledged.