# A New Spin on Neural Networks

**Julie Grollier** 

CNRS/Thales, France













# There is a long standing link between neural networks and spin systems



J. J. Hopfield, PNAS 79, 2554 (1982)

# We are at a convergence point for spin-based neural networks

#### Progress in nanotechnologies



Intel: MRAM integrated into 22nm FinFET CMOS

Physical spin systems for Al C Progress in energy-based Al algorithms



#### Demands of applications: ultra-low power AI

- Human brain : 20 W
- Training the GPT-3 language model: 190,000 kWh
- $\rightarrow$  1,000 years of brain operation!

Neuromorphic chips can reduce the energy consumption of AI by several orders of magnitude



They can solve environmental problem, speed up AI, unlock embedded AI

4

# The dynamical properties of spintronic devices are an asset for neuromorphic computing



J. Grollier, D. Querlioz, K. Y. Camsari, K. Everschor-Sitte, S. Fukami, M. D. Stiles, Nature Electronics 3, 360 (2020)

## Pattern recognition requires separating then grouping data



We can use the complex dynamics of a spintronic nanooscillator to separate input spoken digits, leading to a recognition rate > 99.6%

## Neuron (time-multiplexed)

#### Audio file More complex than input $\rightarrow$ separation 15 200 V<sub>osc</sub> (mV) nput (mV) 0 Λ Voltage Current -200 32.5 33.0 32.0 32.5 33.0 32.0 Time (us) Time (µs) Spintronic oscillator TI-46 database, 5 female Reservoir computing speakers, cochlear pre-processing

J. Torrejon, M. Riou, F. Abreu Araujo et al, Nature 547, 428 (2017)

Synapses (computer)

# Coupled oscillators classify inputs patterns through synchronization



#### Spike sequence: Cheddar



M. Romera, P. Talatchian et al, Nature Com 13, 883 (2022)



#### M. Romera, P. Talatchian et al, Nature 563, 230 (2018) 8

Challenges for building hardware neural networks are the number of components and their connectivity



Image recognition: 10<sup>8</sup> neurons and synapses Brain: 10<sup>11</sup> neurons, 10<sup>15</sup> synapses

 $\rightarrow$  Huge number of nanodevices on chip

Brain: 10<sup>4</sup> synapses/neurons

 $\rightarrow$  Huge connectivity



Motta et al, Science, 366, 6469 (2019)

## **1** – Radio-frequency spintronic neural networks

### 2 – Spin-based neural network learning algorithms

Our idea is to build deep neural networks where components communicate through radio-frequency signals



The spintronic oscillator operates as a neuron that converts DC to RF  $\rightarrow$  we need a synapse that converts RF to DC

N. Leroux, A. Mizrahi et al, Radio-Frequency Multiply-and-Accumulate Operations with Spintronic Synapses, Phys. Rev. Applied 15, 034067 (2021)

# Research performed by:

## CNRS/Thales, France

<u>Nathan Leroux, Andrew Ross</u>, Danijela Marković, Dedalo Sanz Hernandez, Erwan Plouet, Erwann Martin, Teodora Petrisor, Paolo Bortolotti, Juan Trastoy, Alice Mizrahi

## INL, Portugal

<u>Leandro Martins</u>, <u>Alex</u> <u>Jenkins</u> <u>Ricardo Ferreira</u>

# Magnetic tunnel junctions multiply input RF signals by a synaptic weight



Nathan Leroux, Alice Mizrahi et al, Physical Review Applied 15, 034067 (2021)

Weighted sum operation on RF signals is achieved through frequency multiplexing



Nathan Leroux, Alice Mizrahi et al Neuromorph. Comput. Eng. 1 011001 (2021)

## Neurons-to-synapses connection



## The output of neuron 2 is tuned by synaptic weight 2

# Fully spintronic hardware neural network



# Comparison with other nanotechnologies

# We have built the first multilayer fully-nano neural network

Memristor nano-synapses connected by off-the-shelf neuron circuits





Memristor Nanosynapses connected <u>to</u> memristor nano-neurons



 Connecting nano-neurons to nano-synapses is a challenge



Kiani et *al., Sci. Adv.*, 7, 48 (2021)

Oh et al., Nat. Nano. 16 (2021)

Z. Wang et al, Nature Electronics 1, 1737 (2018)

Spintronic RF-to-DC and DC-to-RF conversions enable a nice propagation of signals along the multilayer stack



 $\Rightarrow$  Same device for neurons and synapses

 $\Rightarrow$  Stack layers alternating DC and RF for depth

# Connectivity

# Most physical systems are naturally locally connected

### Spin ice lattice



#### D Wave squid Qubits



J. C. Gartside et al, Nature Nano 13, 53 (2018)

Each magnet is connected to 4 others

8 qubits connected all to all through JJs

2000 spins in total, but only ~50 spins can be coupled all to all through embedding 21

## Our concept is a sneak path free alternative to crossbar arrays

**RF** frequency multiplexing

Crossbar arrays





- Requirement: high OFF/ON ratio devices
- sneak paths
- ~ 100 synapses per neuron

- Requirement: good RF signal propagation
- No sneak paths

 $\sim 1000$  synapses per neuron

## Towards deep RF spintronic neural networks

- First experimental demonstration of communication between multiple layers using nano-devices
- Non-linear classification of two inputs with high accuracy

10 fJ/synapse and 100 fJ/neuron for MTJs with 20 nm diameter: 100 fold energy gain compared to GPUs



Theory weighted sum: N. Leroux, A. Mizrahi *et al. Phys. Rev. Appl.* **15**, 034067 (2021) Experiment weighted sum: N. Leroux, A. Mizrahi *et al.*, Neuromorph. Comput. Eng. **1**, 011001 (2021) Full network experiment: A. Ross, N. Leroux *et al. Manuscript in Preparation* Simulations: N. Leroux *et al Neuromorph. Comput. Eng.* **2**, 034002 (2022) 1 – Radio-frequency spintronic neural networks

### 2 – Spin-based neural network learning algorithms

## Hopfield nets compute by minimizing an Ising energy



for each pattern p, one after the other

J. J. Hopfield, PNAS 79, 2554 (1982)

## Modern algorithms minimize the error at the network output

### Local learning rules



Hebbian learning rules (1949): who fires together wires together

$$\Delta w_{ij} += s_i^p s_j^p$$

### **Backpropagation of errors**

Cost function:  $C = \frac{1}{2}(y_l - t_l)^2$ 



Easy to implement in hardware Low performance on complex tasks Hard to implement in hardware High performance on complex tasks 26 Equilibrium propagation minimizes both the energy and the output error of the system

Inference : free phase



Nudging toward the desired solution



Scellier & Bengio, fnins 2017



Equilibrium propagation directly backpropagates gradients through neuronal dynamics

Inference : free phase

 $\frac{\partial E}{\partial s_i}(s_i^0) = 0$  $\frac{\partial E}{\partial s_i} \left( s_i^\beta \right) + \beta \frac{\partial C}{\partial s_i} = 0$ The perturbation at the output propagates Wait for equilibrium, Wait for equilibrium, through the measure  $s_i^0$ measure  $s_i^{\beta}$ network clamped clamped  $\Delta s_i \propto -\frac{\partial C}{\partial s_i} \implies \Delta w_{ij} \propto -\frac{\partial C}{\partial w_{ij}}$ Difference between the two equilibrium states

Propagation of gradients : nudging phase  $F = E + \beta C$ 

Our recent work: bridging Equilibrium Propagation and hardware

## CNRS/Thales, France

<u>Jérémie Laydevant</u>, <u>Erwann Martin</u>, Dongshu Liu, Shuai Li, <u>Danijela Marković,</u> Julie Grollier

## C2N, France

<u>Maxence Ernoult</u>, <u>Axel Laborieux</u>, Damien Querlioz

## Thales, France

Erwann Martin, Teodora Petrisor

## Mila, Canada

Benjamin Scellier, Yoshua Bengio

The local learning rule of Equilibrium Propagation (EP) leads to high recognition rates on image benchmark tasks



Equivalence of EP local learning rule with the gradients of Backpropagation Through Time (maths + simulation results)

$$\Delta W_{ij} = -\rho(s_i) \,\rho(s_j) + \rho^\beta(s_i) \,\rho^\beta(s_j)$$



EP implements convolutions and scales to CIFAR-10

Modified 3-phase learning rule: EP : 11.68% test error BPTT : 11.1 %

M. Ernoult, J. Grollier, D. Querlioz, Y. Bengio, B. Scellier, NeurIPS (2019) A. Laborieux, M. Ernoult, B. Scellier, Y. Bengio, J. Grollier, D. Querlioz, Frontiers Neuro. 15, 129 (2021) We have implemented a full neural network and trained it to solve MNIST through Equilibrium propagation on the D'Wave quantum computer



#### Ising machine

Work of Jérémie Laydevant and Danijela Marković

1<sup>st</sup> challenge: design a binary version of Equilibrium propagation compatible with Ising machines

• D-wave is an Ising machine based on 2-states qubits s<sub>i</sub>

$$E = -\frac{1}{2} \sum_{i \neq j} w_{ij} s_i s_j - \sum_i b_i s_i$$

• Equilibrium Propagation uses analog neurons  $\rho(s_i)$ 

$$E = -\frac{1}{2} \sum_{i \neq j} w_{ij} \rho(s_i) \rho(s_j) - \sum_i b_i \rho(s_i) + \frac{1}{2} \sum_i s_i^2$$

 $\rightarrow$  We have demonstrated in software that EP trains Binary Neural Networks by using advanced Machine Learning methods

- Binary synapses (CIFAR 10)
- Binary synapses and neurons (MNIST)
- Ternary gradients

J. Laydevant, M. Ernoult, D. Querlioz, J. Grollier, arxiv 2103.08953, published in the Conference on Computer Vision and Pattern CVPR (2021) 2<sup>nd</sup> challenge: mapping a neural network that can classify handwritten digits from the MNIST database to the chip layout



1 qubit  $\neq$  1 neuron due to the chip layout

On chip

3<sup>rd</sup> challenge: hacking the quantum annealing procedure to reach the two different equilibrium points needed for EP



First end-to-end supervised training of a neural network in an Ising machine (D-Wave) with software-equivalent accuracy (85%)



In the future, the accuracy can be improved by increasing:

- the number of presented images in the database
- the number of hardware spins

Spin-based learning algorithms can solve standard AI tasks and achieve SOTA accuracy Spintronic nanodevices communicating through RF signals can implement large scale and low power hardware neural networks







# Learning with high accuracy and low power is the main challenge of neuromorphic chips

Neuromorphic chips today: Limited to Handwritten digit classification





On-chip learning circuits are greedy in energy and area

# Equilibrium propagation replaces gradient computation by gradient measurement



## Can we do better ?

Can we train a hardware neural network through it physics, without additional « artificial circuits », with low power, and with high accuracy despite component variability, like in the brain ?



Eqspike: E Martin, M Ernoult, J Laydevant, S Li, D Querlioz, T Petrisor, J Grollier, iScience 24, 102222 (2021)