

# Bayesian In-Memory Computing

Damien Querlioz

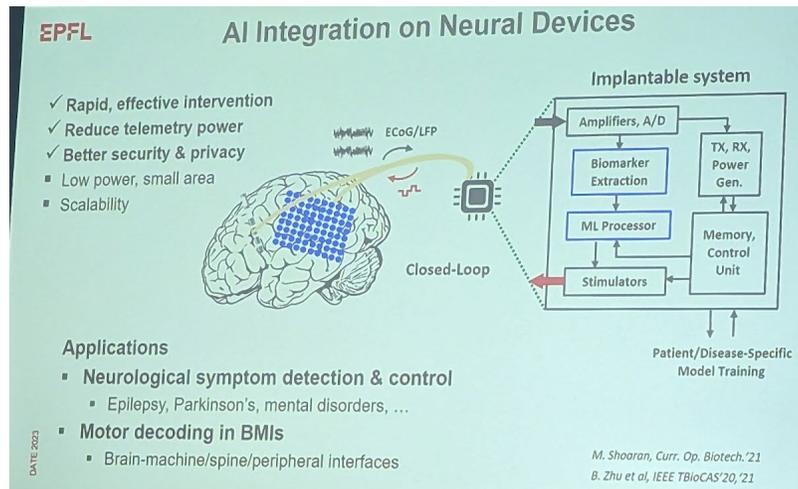
Centre de Nanosciences et de Nanotechnologies  
Université Paris-Saclay, CNRS, Palaiseau, France  
[damien.querlioz@universite-paris-saclay.fr](mailto:damien.querlioz@universite-paris-saclay.fr)



Joint work with the groups of Marc Bocquet and Jean-Michel Portal (Aix-Marseille Univ.) and Elisa Vianello (CEA-LETI)

# Edge AI Has an Incredible Potential for Safety-Critical Applications

**Medical: Predicting epileptic seizures, closed-loop Parkinson DBS...**

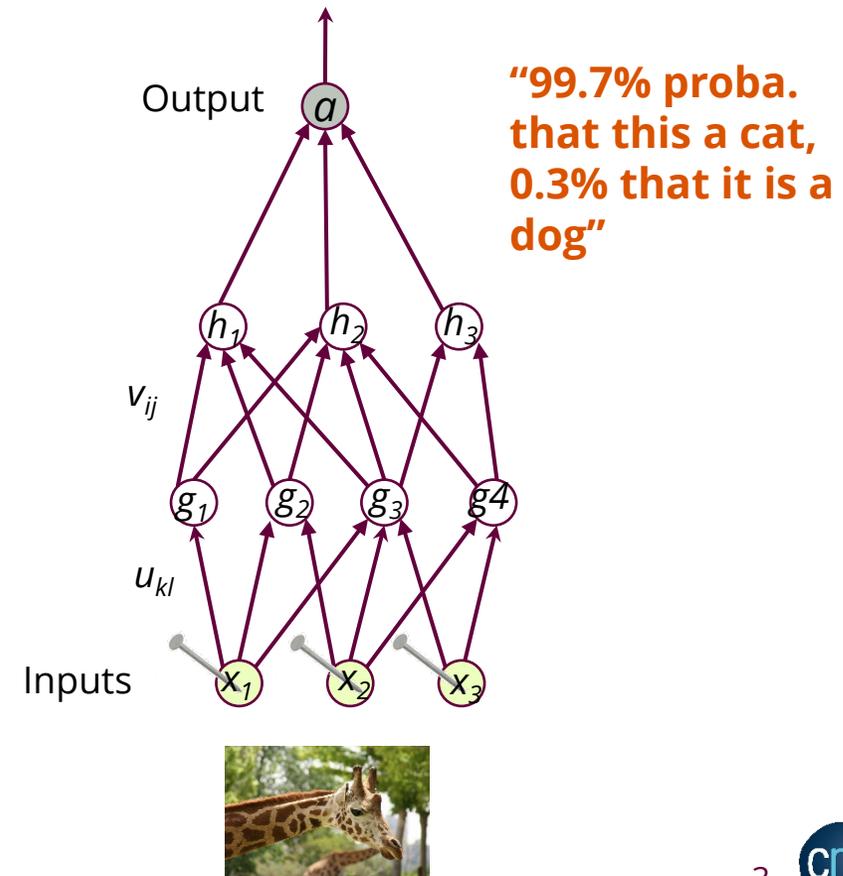
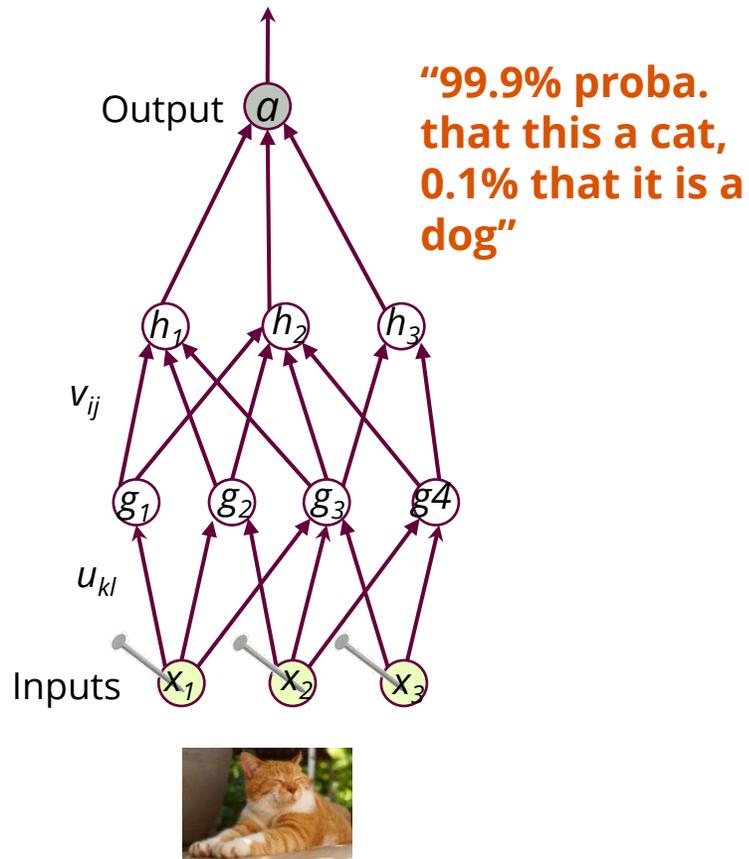


Shoaran Mahsa (EPFL) talk  
DATE 2023, Special Day on Human-AI  
Interaction

**Industrial: Monitoring early drifts/faults to avoid accidents**

# Modern AI Is Amazing But It Has a Confidence Problem

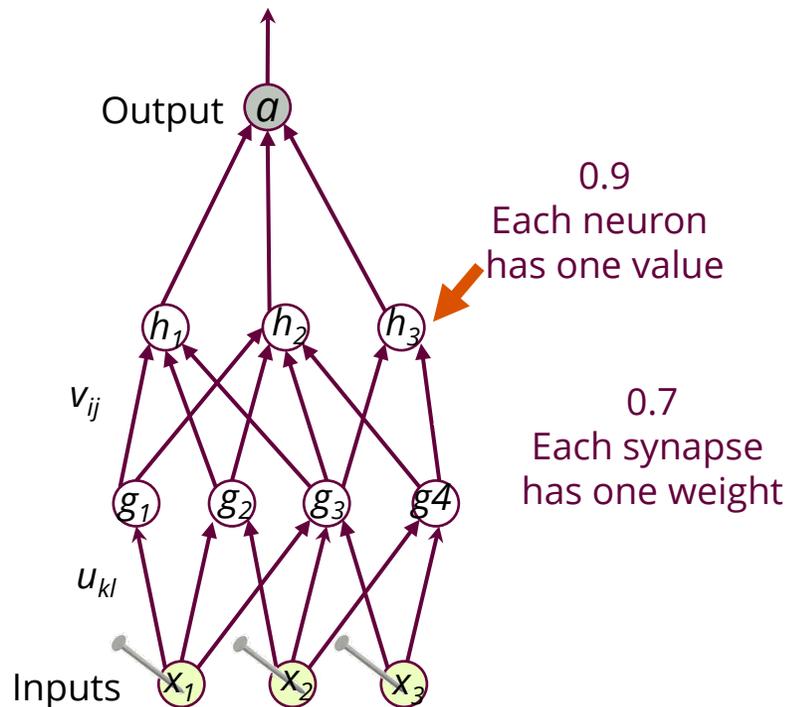
- If a network has been trained to distinguish CATS and DOGS



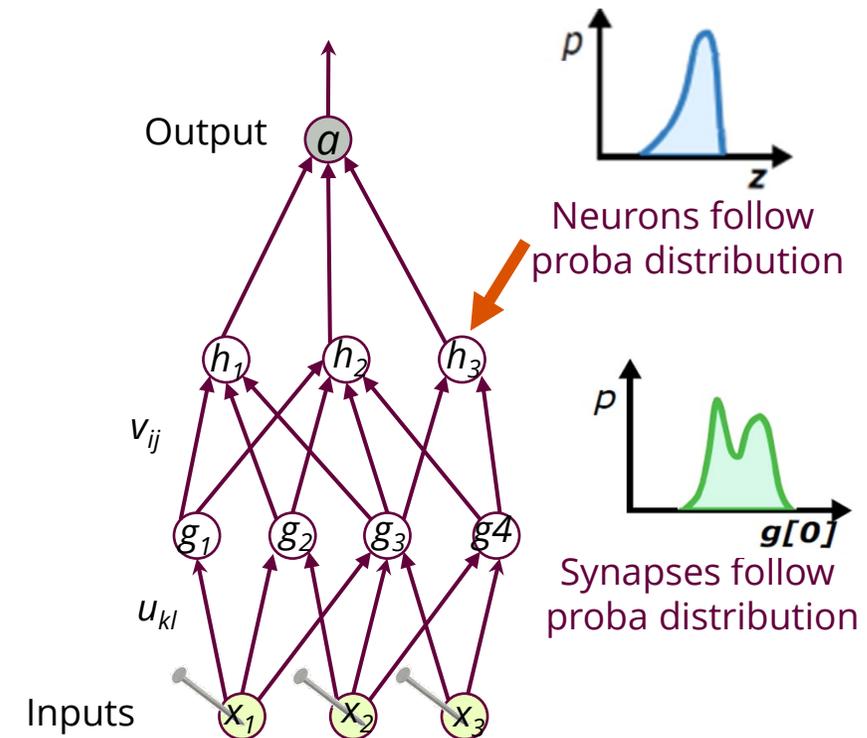
# A Lead: Bayesian Models

- In Bayesian models, everything is considered a random variable that follows specific probability distributions

## Traditional neural network



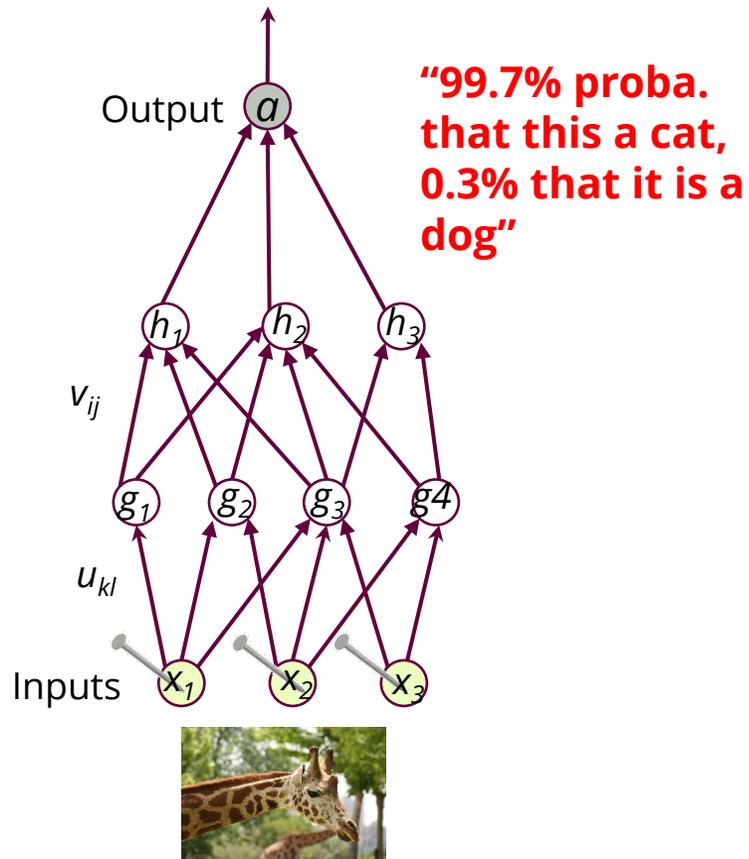
## Bayesian neural network



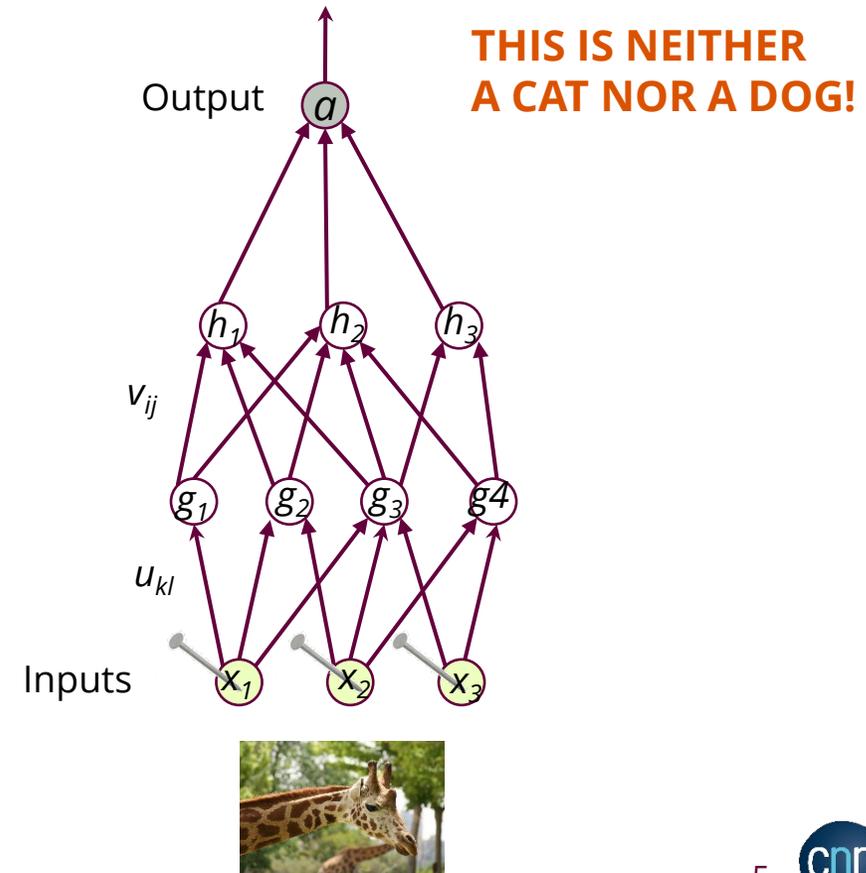
# A Lead: Bayesian Models

- If a network has been trained to distinguish CATS and DOGS

**Traditional neural network**

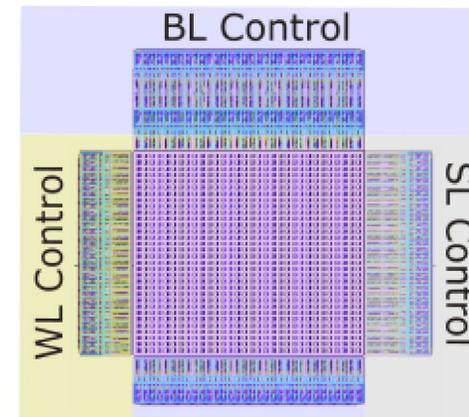
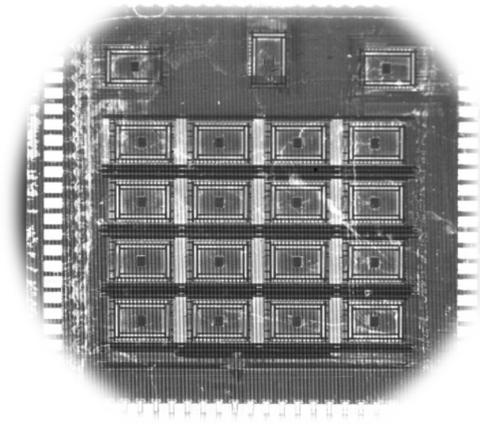
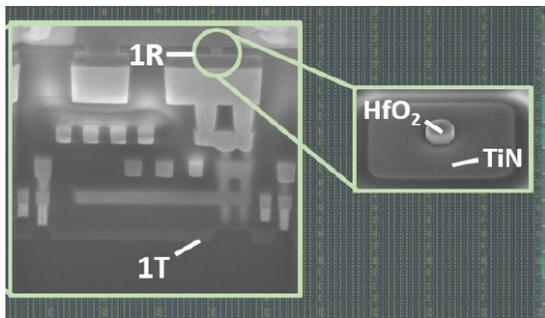


**Bayesian neural network**

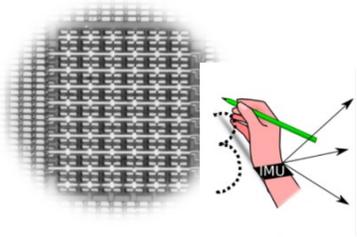


# Bayesian Models Are All About Probabilities

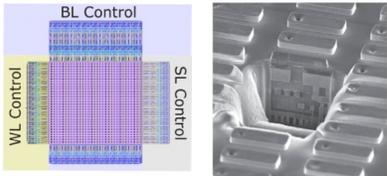
- Probabilities are hard to do on computers
- **Nanoelectronics offers many opportunities**



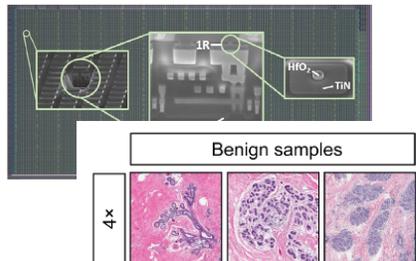
# Bayesian In-Memory Computing



- The Memristor-Based Bayesian Machine

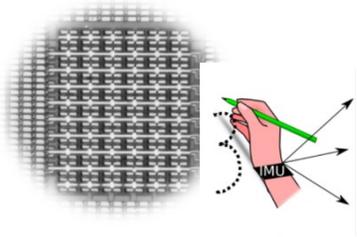


- Bayesian Neural Networks with Memristors

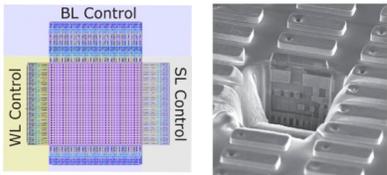


- Bayesian Learning with Memristors

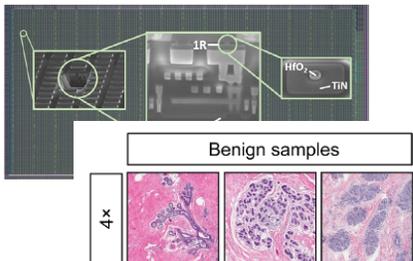
# Bayesian In-Memory Computing



- The Memristor-Based Bayesian Machine

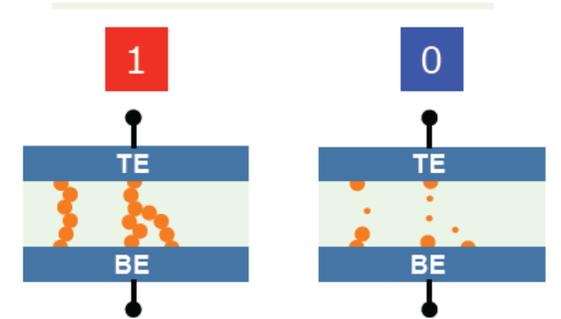
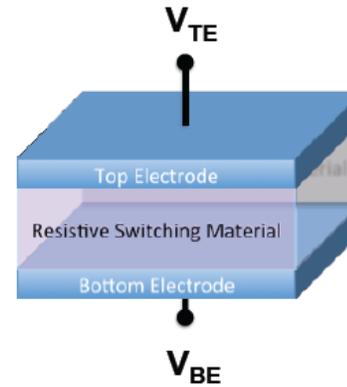
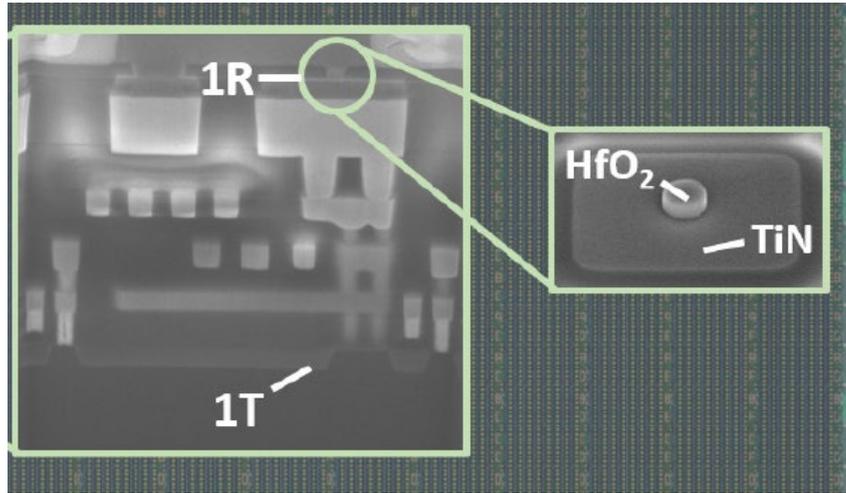


- Bayesian Neural Networks with Memristors



- Bayesian Learning with Memristors

# Memristor/RRAM: A Nanomemory Embeddable at the Core of CMOS



Low resistance

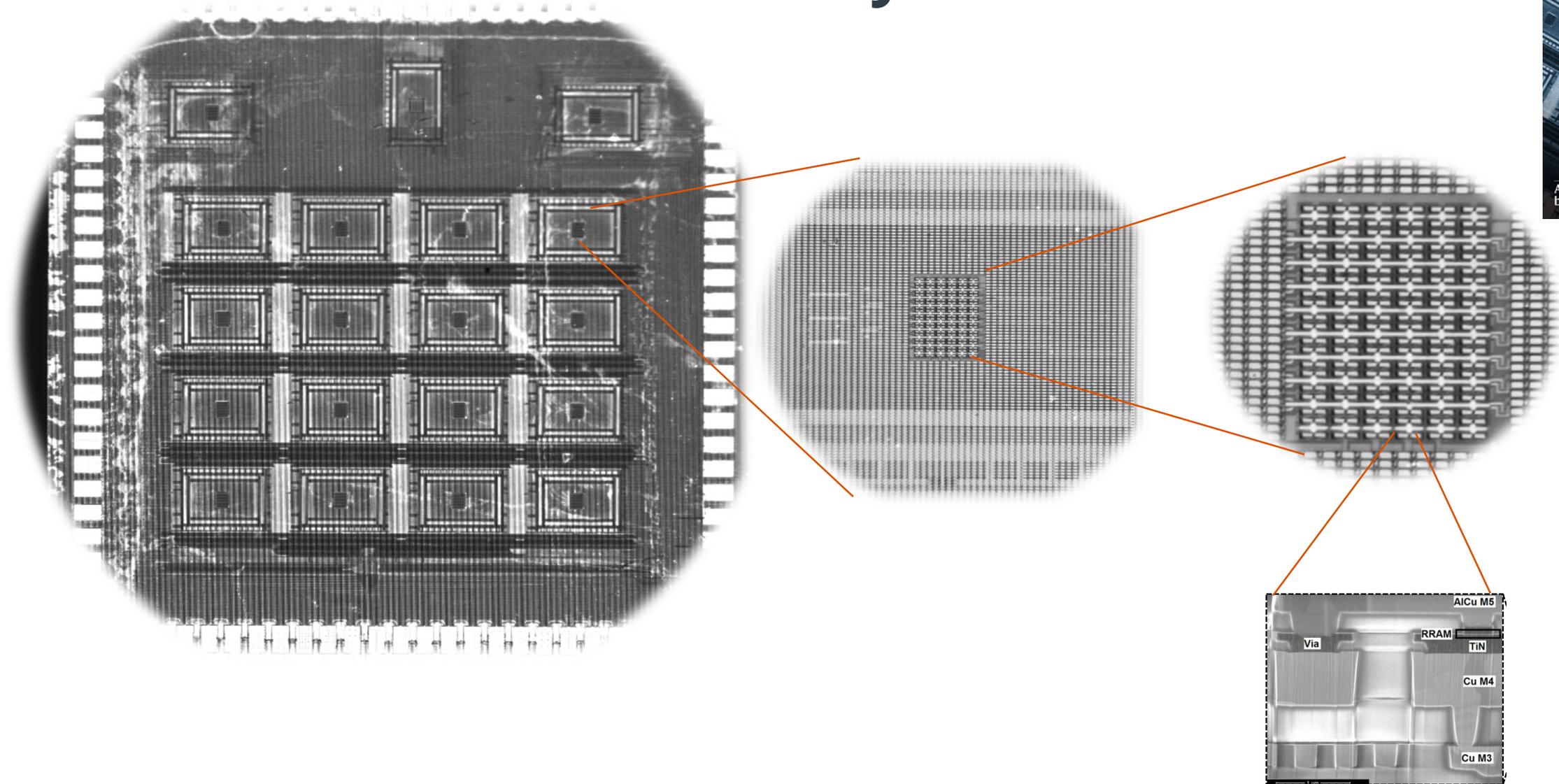
High resistance

TiN/HfO<sub>x</sub>/Ti/TiN stack

**High voltage:** move atoms to switch memristor between low/high resistance

**Low voltage:** allows reading the resistance

# The Memristor-Based Bayesian Machine



Harabi, Hirtzlin, Turck et al, Nature Electronics 6, 53 (2023)

# Bayesian Reasoning: *Better at Small Data*



Thomas Bayes

*Likelihoods*

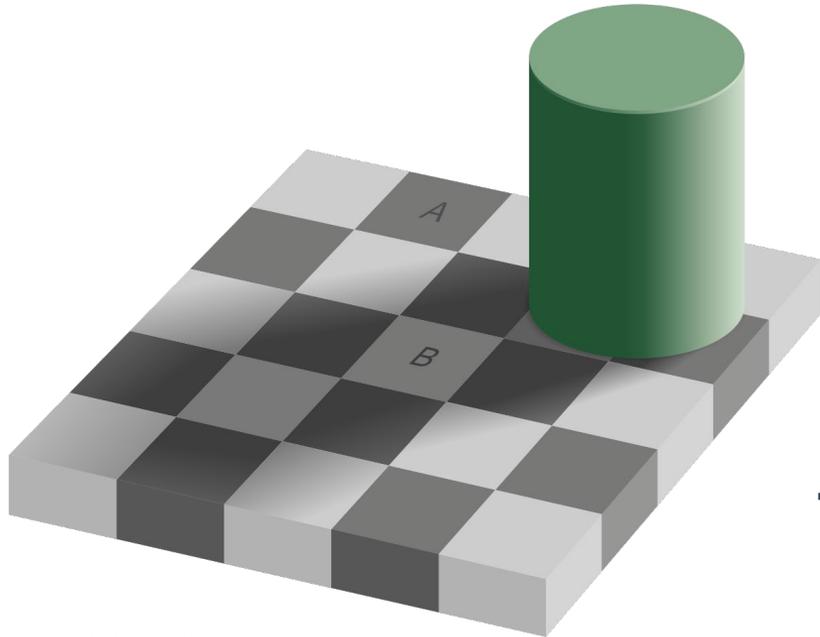
*Prior*

$$p(\text{Disease} \mid \text{Observations}) \propto p(\text{Observations} \mid \text{Disease}) \times p(\text{Disease})$$



*Constructed with expert knowledge+Data*

*Hard to Compute*



**The Brain Might Be Using Bayesian Reasoning**

# Stochastic Computing

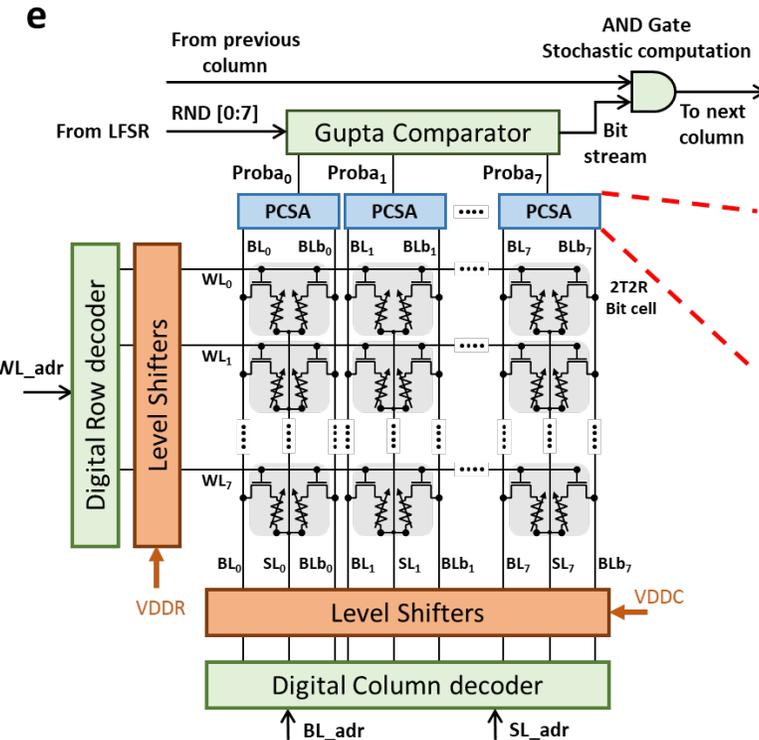
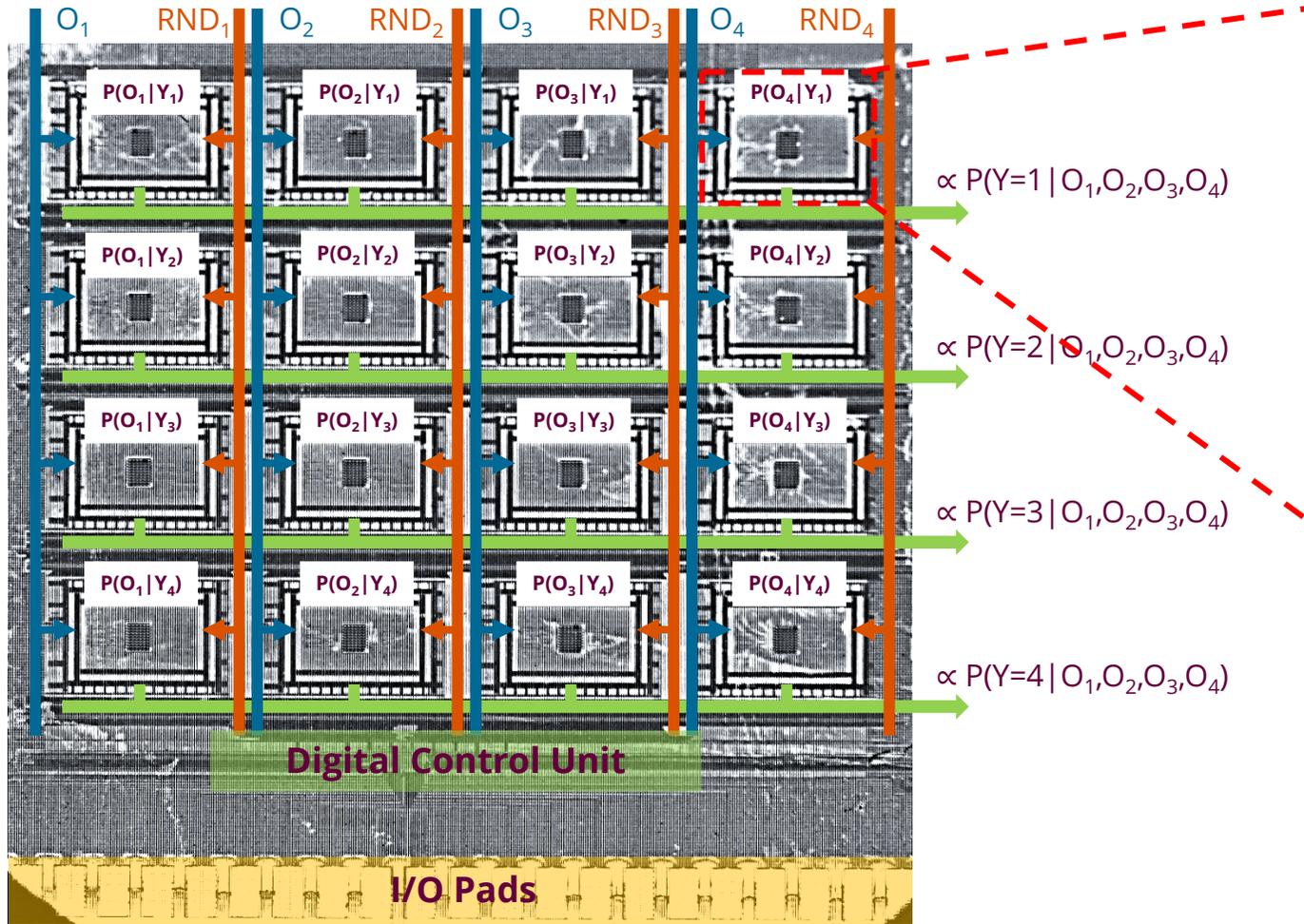


$$p(S_3) = p(S_1)p(S_2)$$

A AND gate implements the multiplication of two probabilities!

Some resemblance with neurons (one wire = one real number)

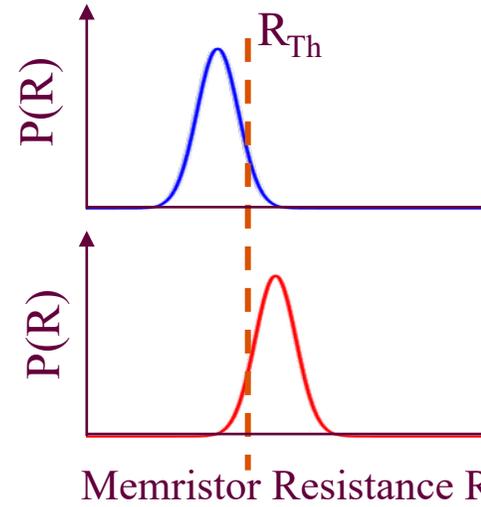
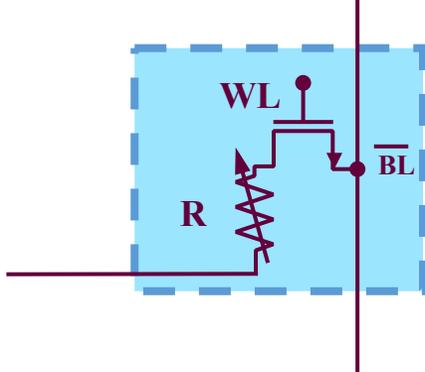
# The Memristor-Based Bayesian Machine



$$P(Y=4 | O_1, O_2, O_3, O_4) \propto P(O_1 | Y=4) \times \dots \times P(O_4 | Y=4)$$

# Reducing the Error Rate without ECC

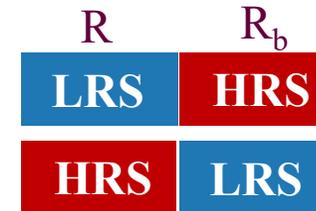
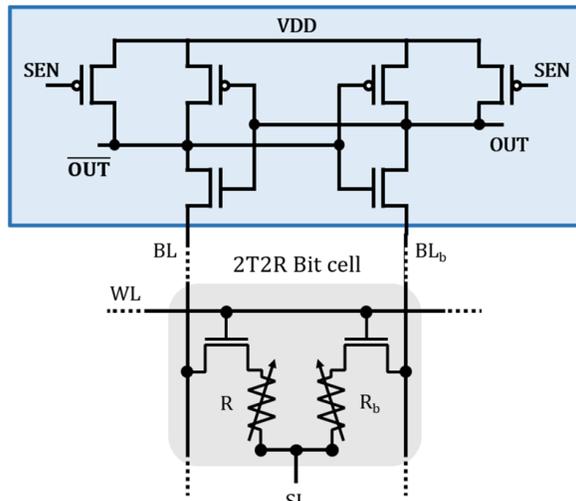
## Traditional approach



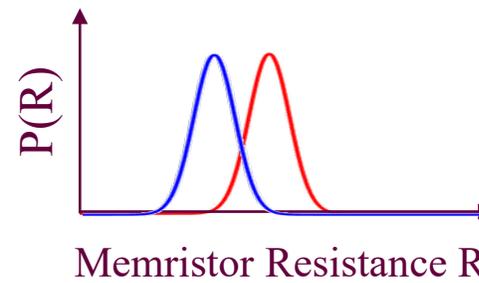
State 0 **LRS**  
Error if  $R(\text{LRS}) > R_{Th}$

State 1 **HRS**  
Error if  $R(\text{HRS}) < R_{Th}$

## Our approach (Bayesian Machine)



**State 1**  
**State 0**

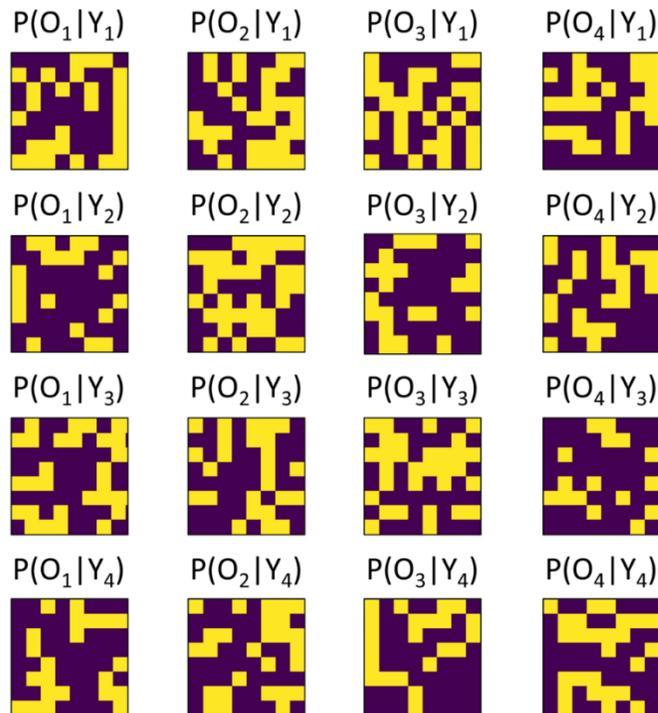


An error happens if  
 $R(\text{LRS}) > R(\text{HRS})$

# The Memristor-Based Bayesian Machine

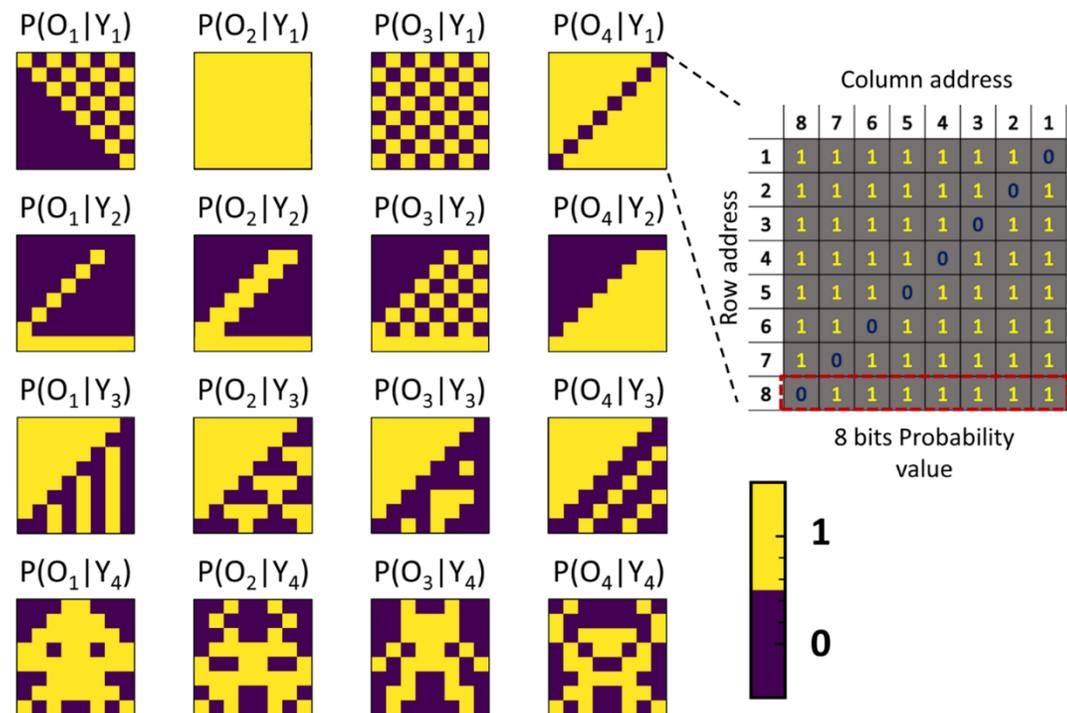
Before programming the memristors

a



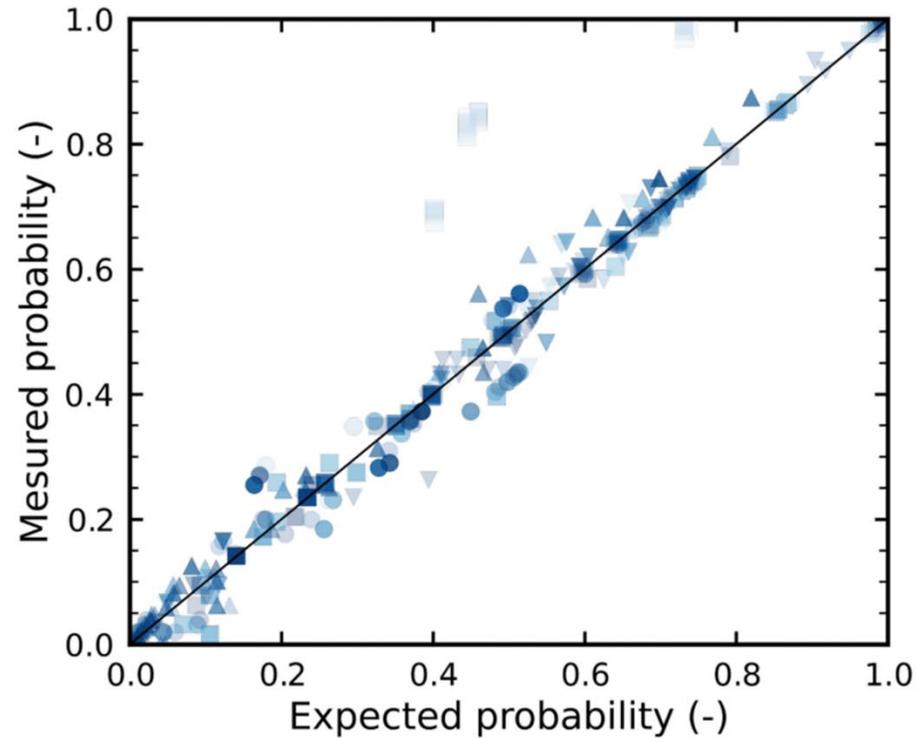
After programming the memristors

b

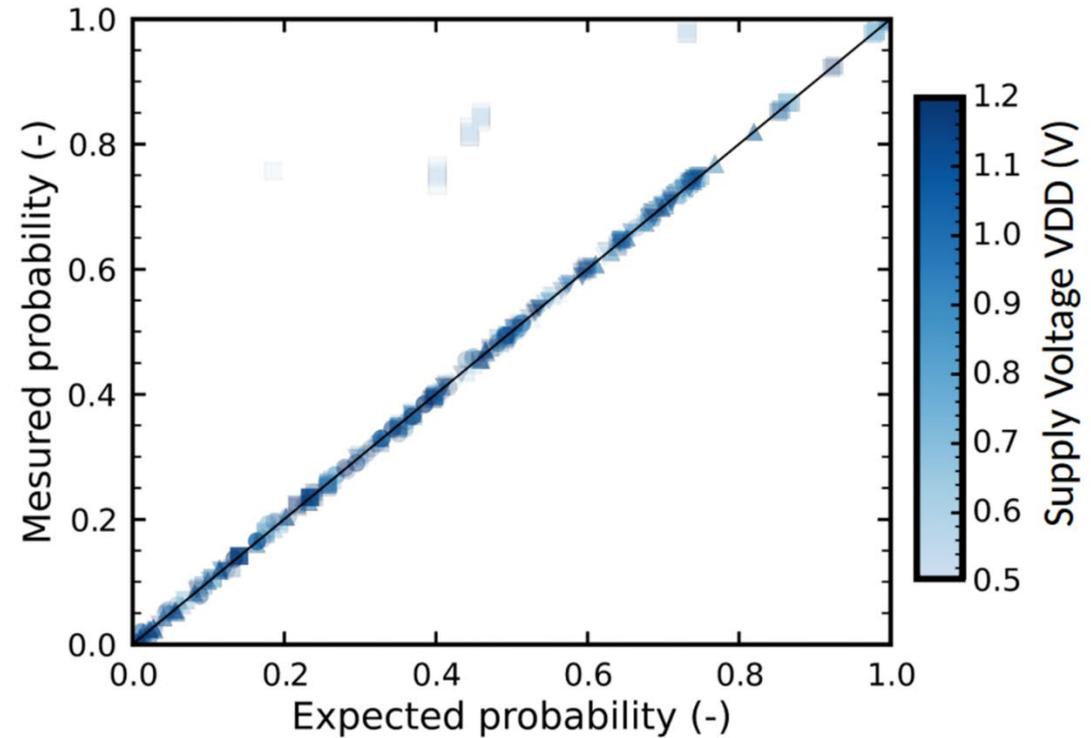


# The Memristor-Based Bayesian Machine

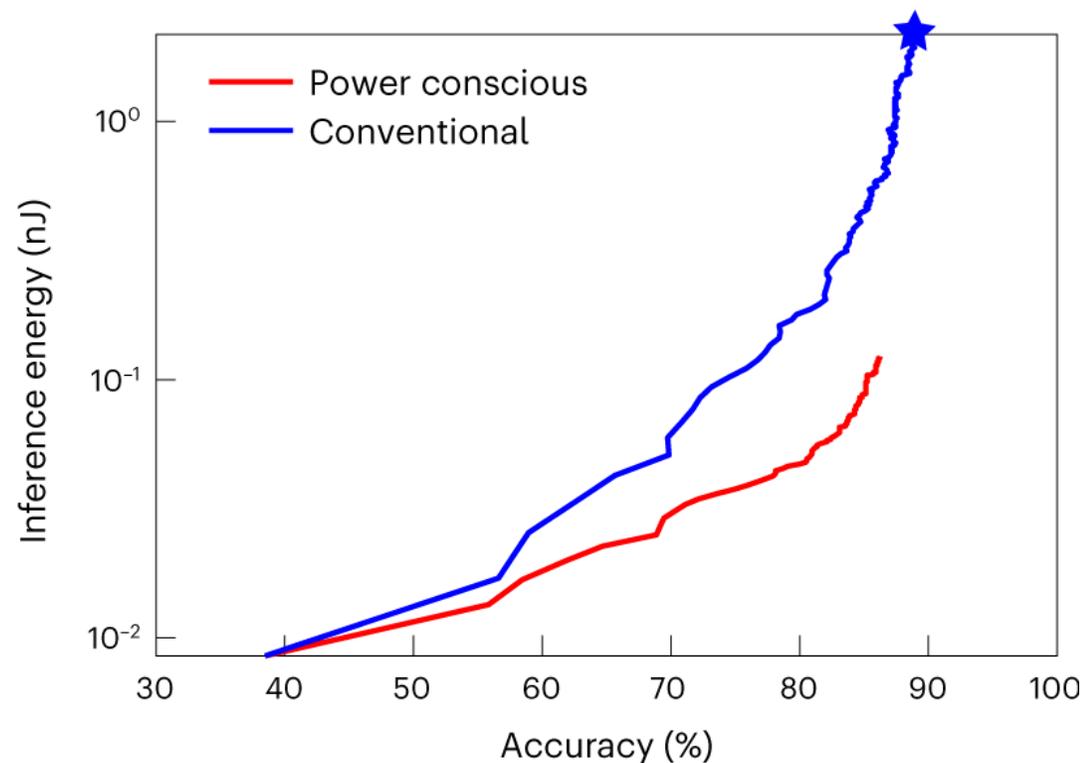
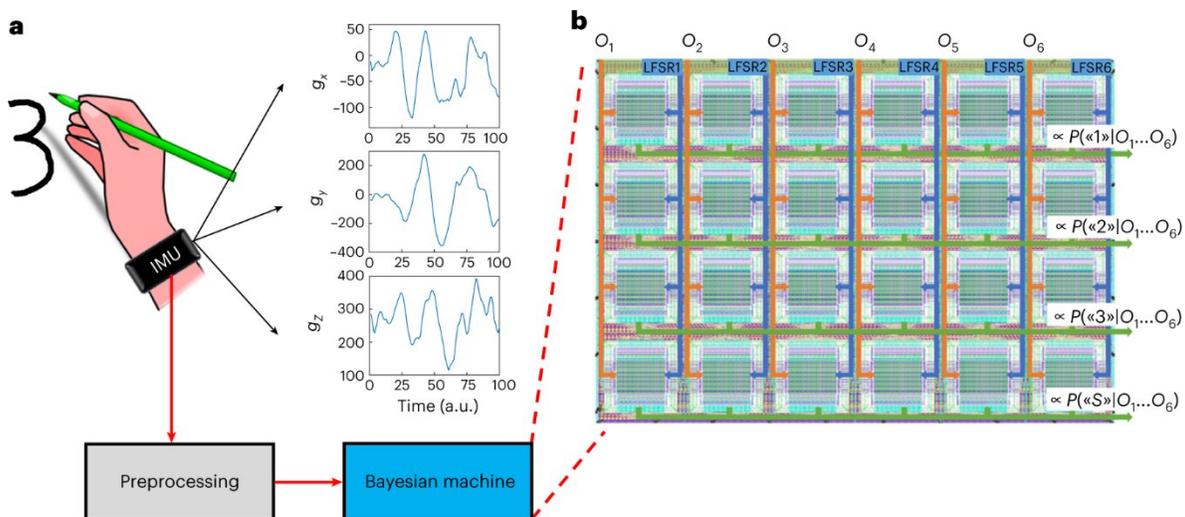
Randomly chosen  
pseudo-RNG seeds



Best choice of  
pseudo-RNG seeds



# Energy Consumption Is Very Low

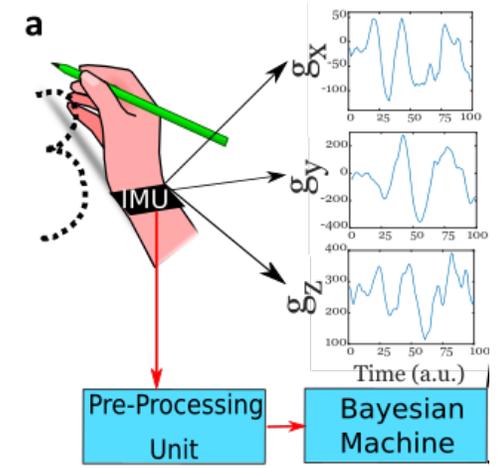
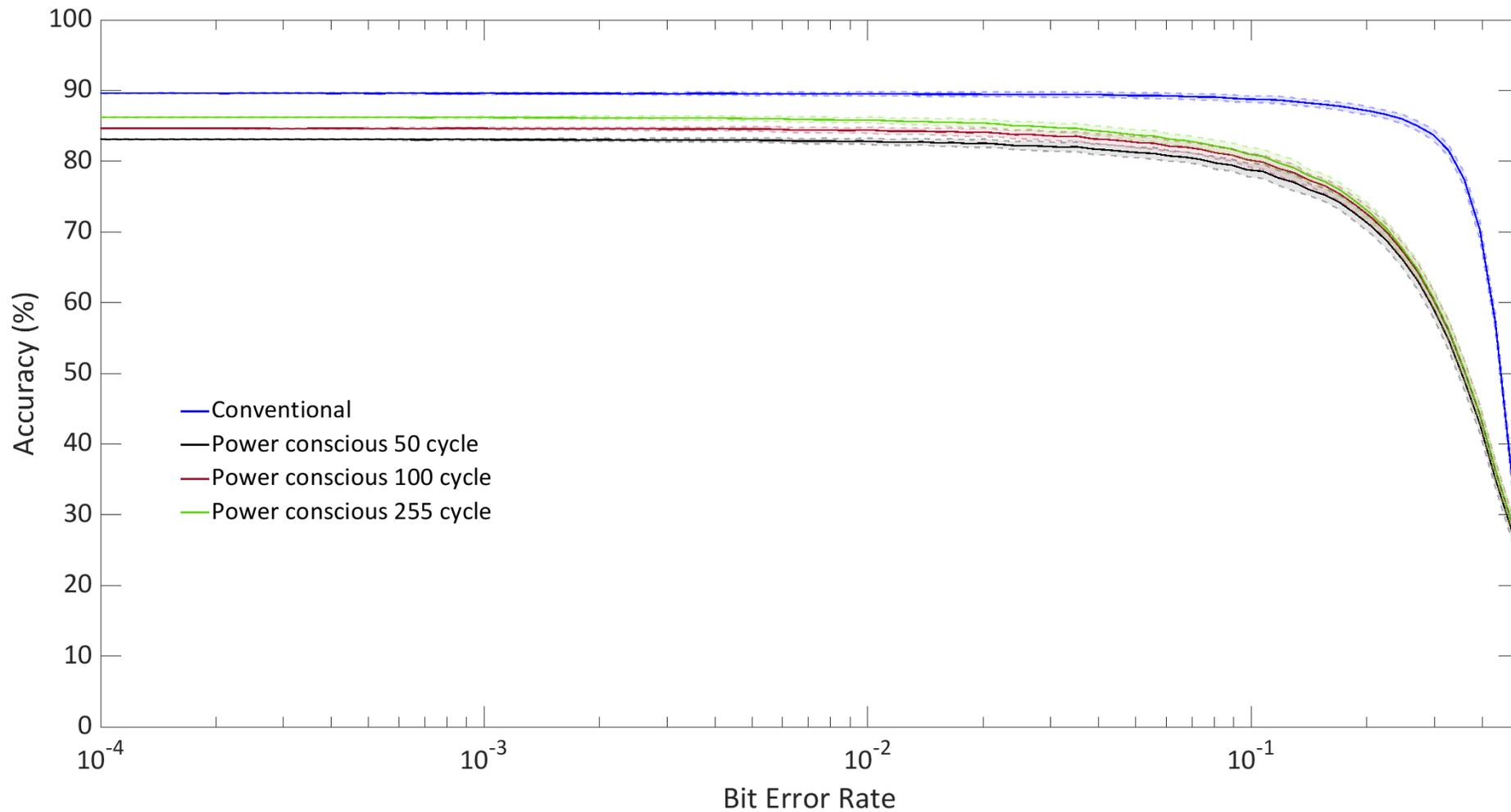


## Two reasons for energy efficiency:

- Near-memory computing
- Stochastic computing

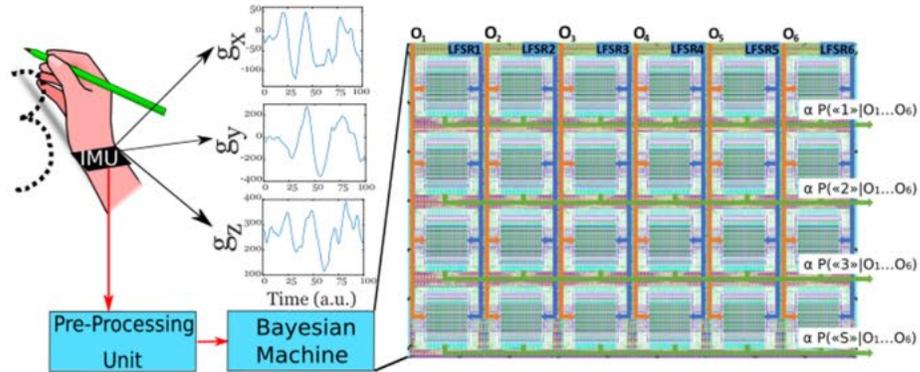
Same task, STM32 microcontroller unit:  $2\mu\text{J}$

# Benefits of Stochastic Computation



# Uncertainty Evaluation

## Signatures of an uncertain Bayesian machine



0 ones

0 ones

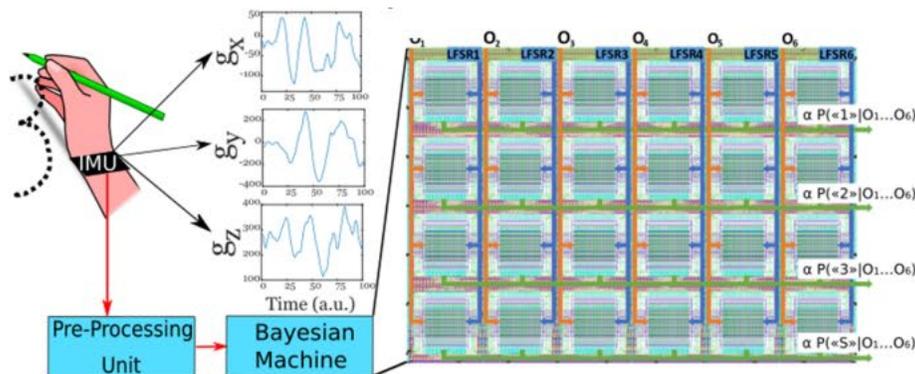
0 ones

3 ones

**Contradictory observations:**  
*All outputs produce a very small number of ones*

**EPISTEMIC UNCERTAINTY**

OR



127 ones

100 ones

93 ones

3 ones

**Ambivalent situation:**  
*Some outputs produce a balanced number of ones*

**ALEATORIC UNCERTAINTY**

# Uncertainty Evaluation

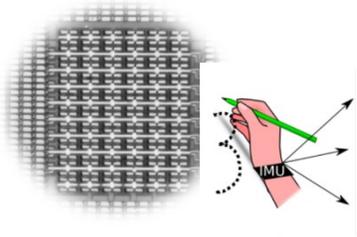
Case when a gesture is presented to the Bayesian machine **corresponding to a different subject**

	T=0	T=5%	T=10%
<b>Bayesian machine was uncertain (desired behavior)</b>	93%	97%	98%
<b>Bayesian machine provided a certain output</b>	7%	3%	2%

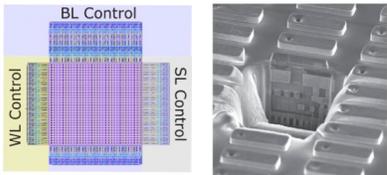
Case when a gesture is presented to the Bayesian machine **corresponding to the appropriate subject**

	T=0	T=5%	T=10%
<b>Bayesian machine was certain about the correct gesture (desired behavior)</b>	88%	80%	71%
<b>Bayesian machine was certain about an incorrect gesture</b>	7%	6%	5%
<b>Bayesian machine was uncertain</b>	5%	14%	24%

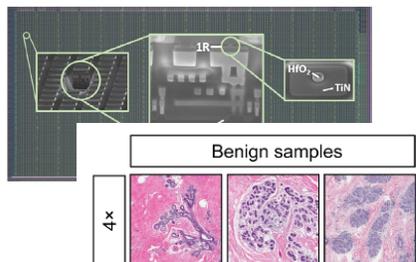
# Bayesian In-Memory Computing



- The Memristor-Based Bayesian Machine



- Bayesian Neural Networks with Memristors

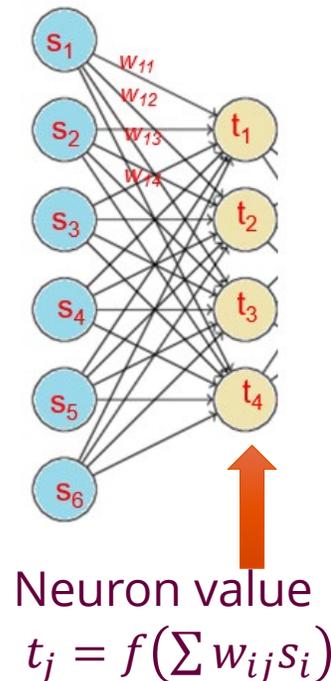
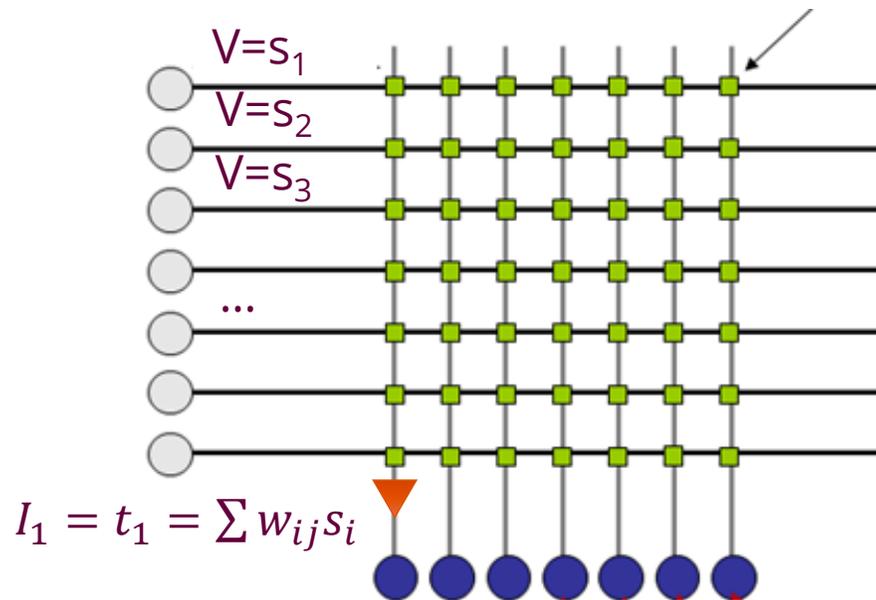


- Bayesian Learning with Memristors

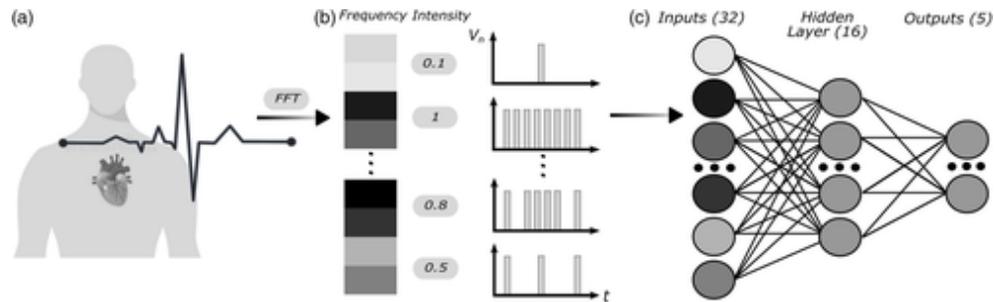
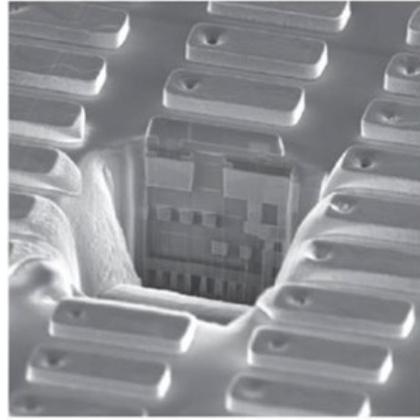
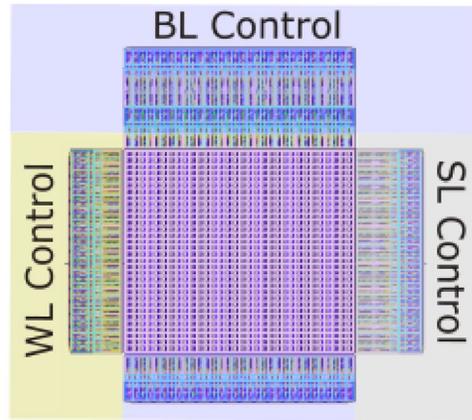
# How to Make a Neural Network with Memristors

- A matrix of analog memristors naturally implements a layer of neural network with Kirchhoff laws!

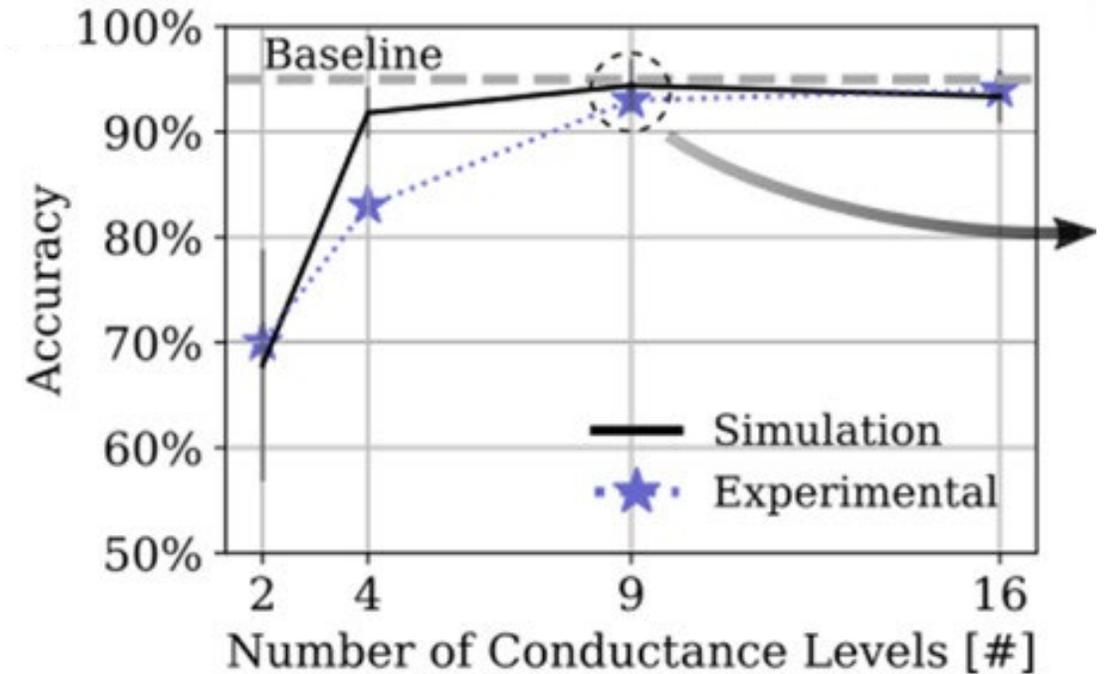
Memristor conductance  $G$  = synaptic weight  $w$



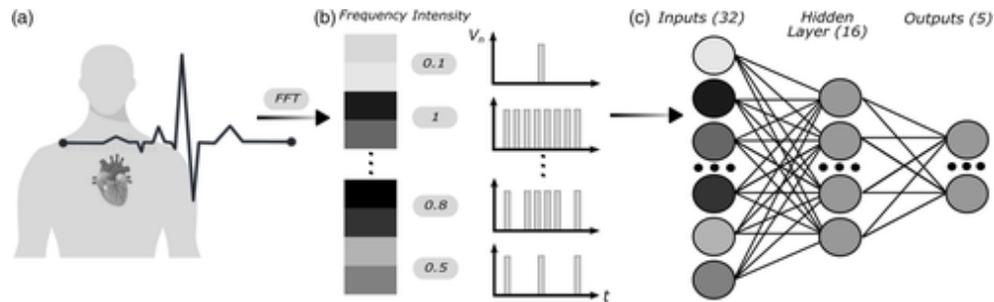
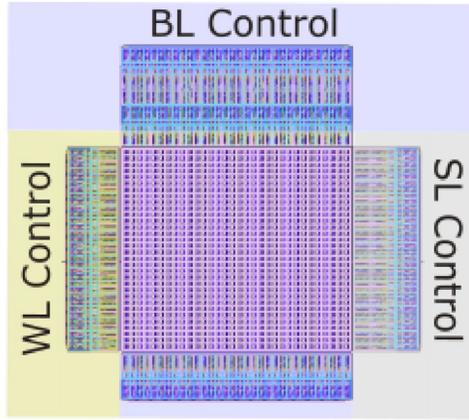
# Experimental Realization with HfO<sub>x</sub> Memristors



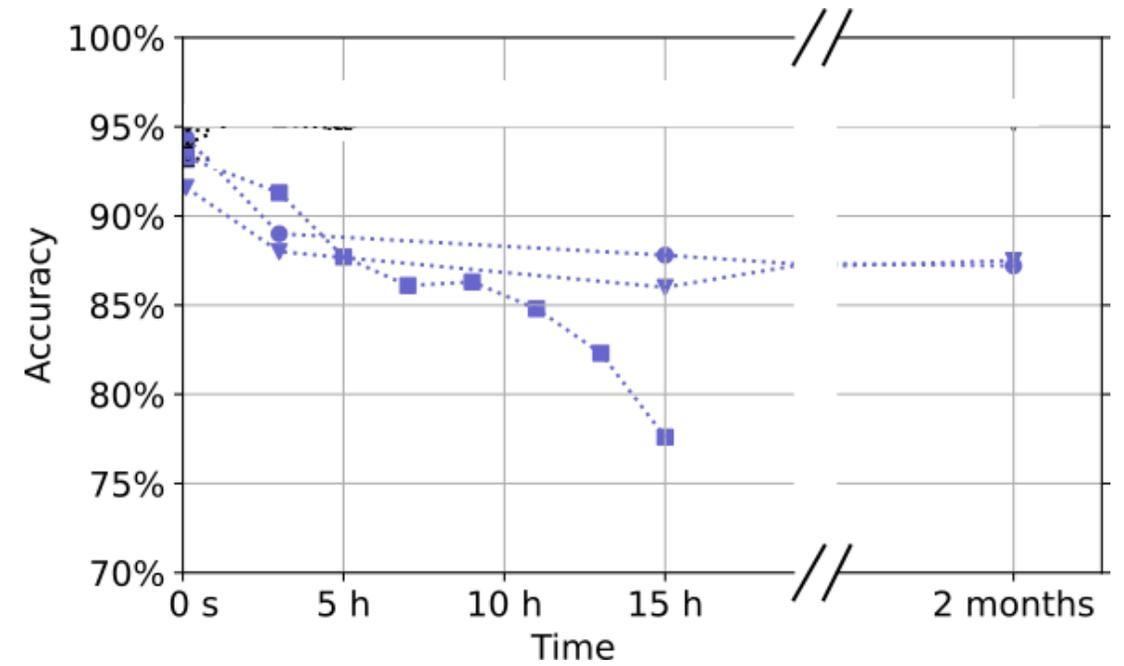
Accuracy on arrhythmia identification



# Experimental Realization with HfO<sub>x</sub> Memristors

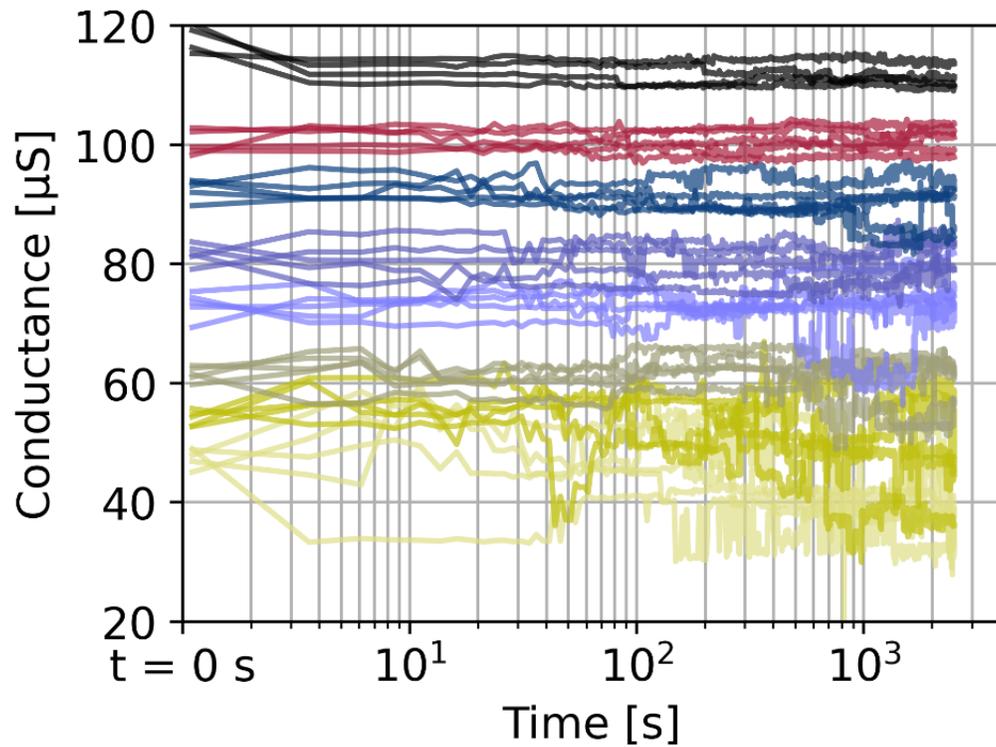


Accuracy on arrhythmia identification

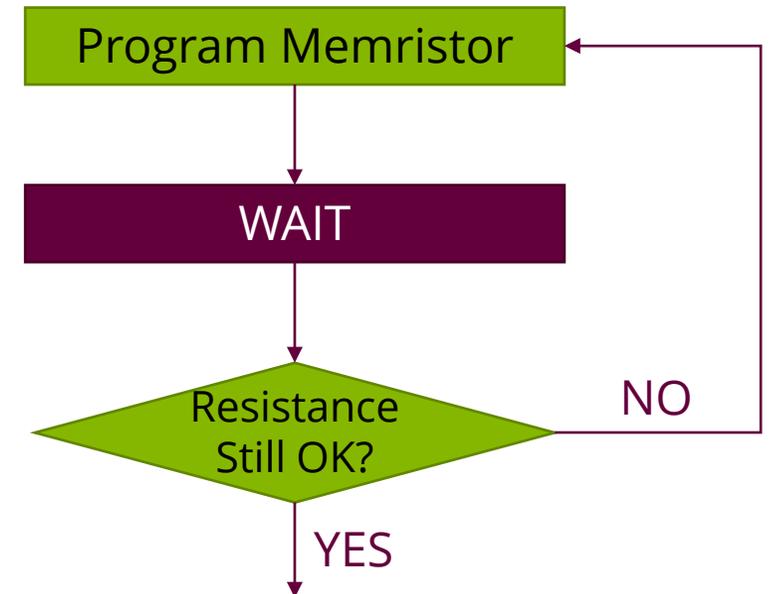


# Writing Memristors Reliably

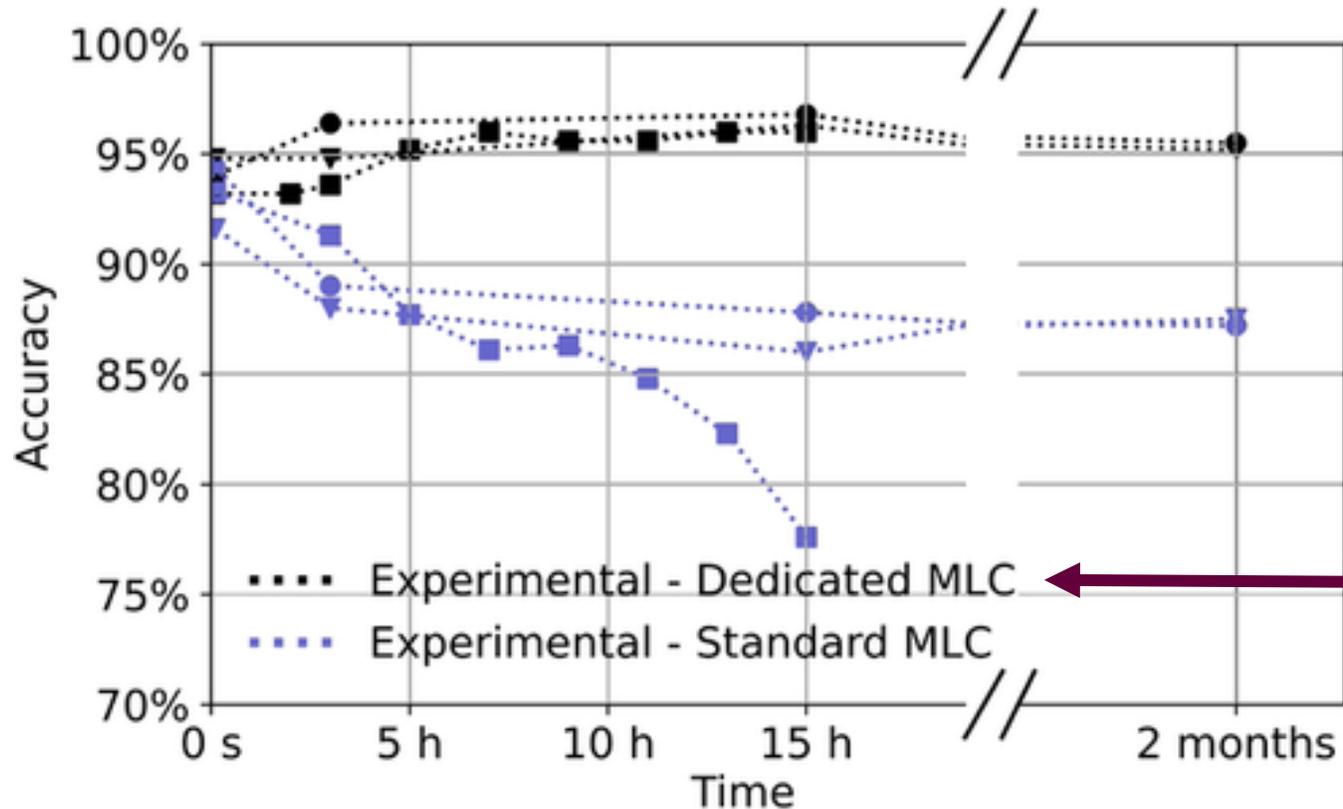
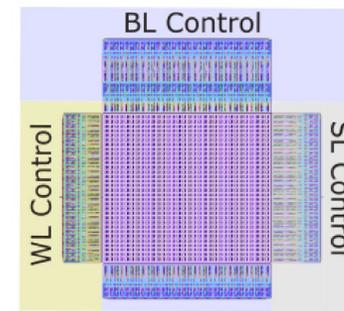
## Problem



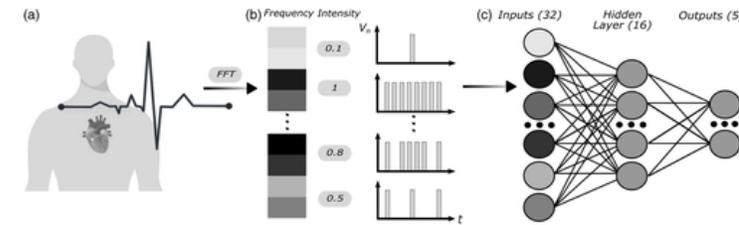
## Solutions



# Writing Memristors Reliably



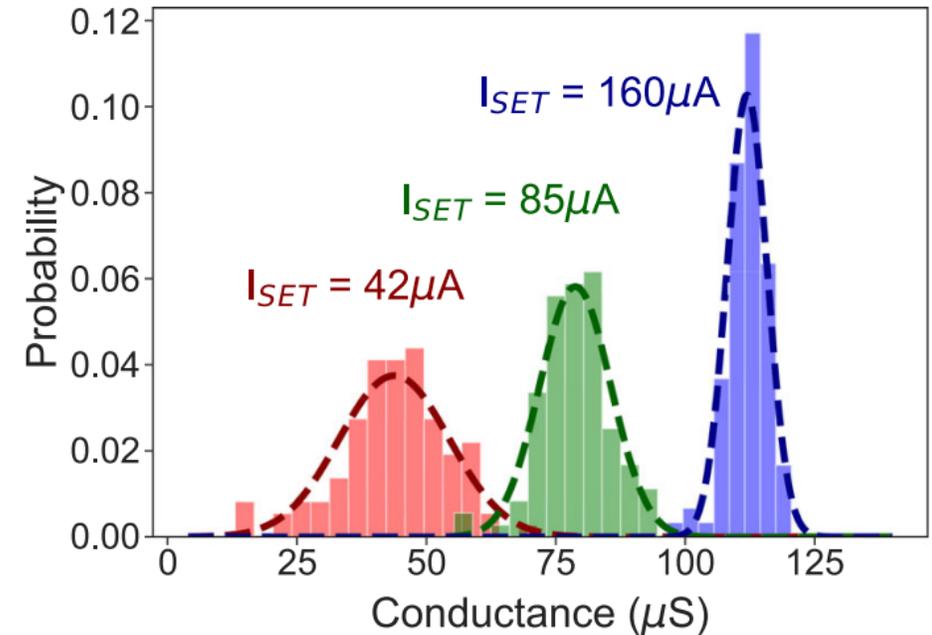
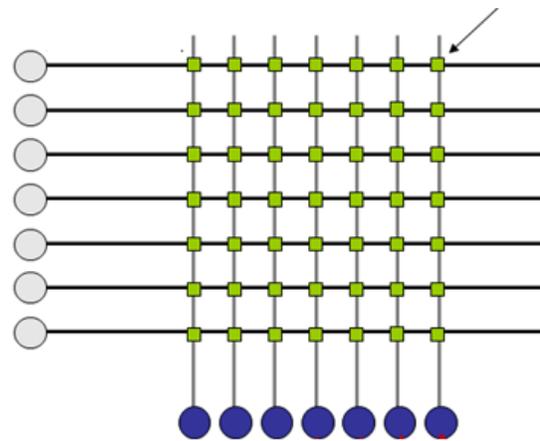
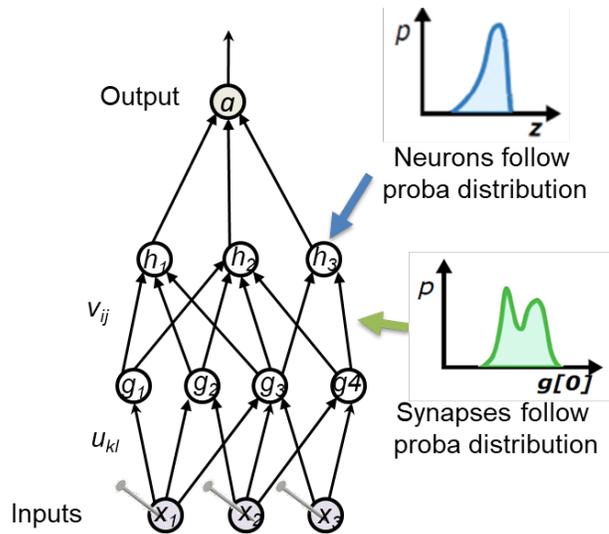
← With the Wait



# Do Not Fight Memristor Imperfections: They Naturally Produce a Bayesian Neural Network!

In Bayesian models, everything is a random variable that follows specific probability distributions

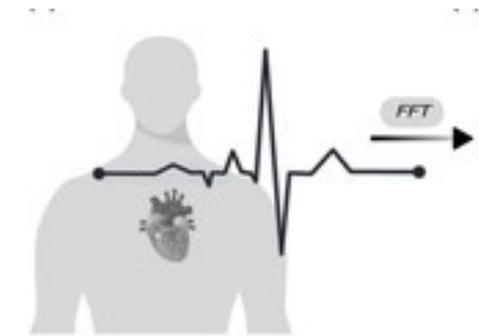
Memristors actually act as a random variable that follow specific probability distributions!



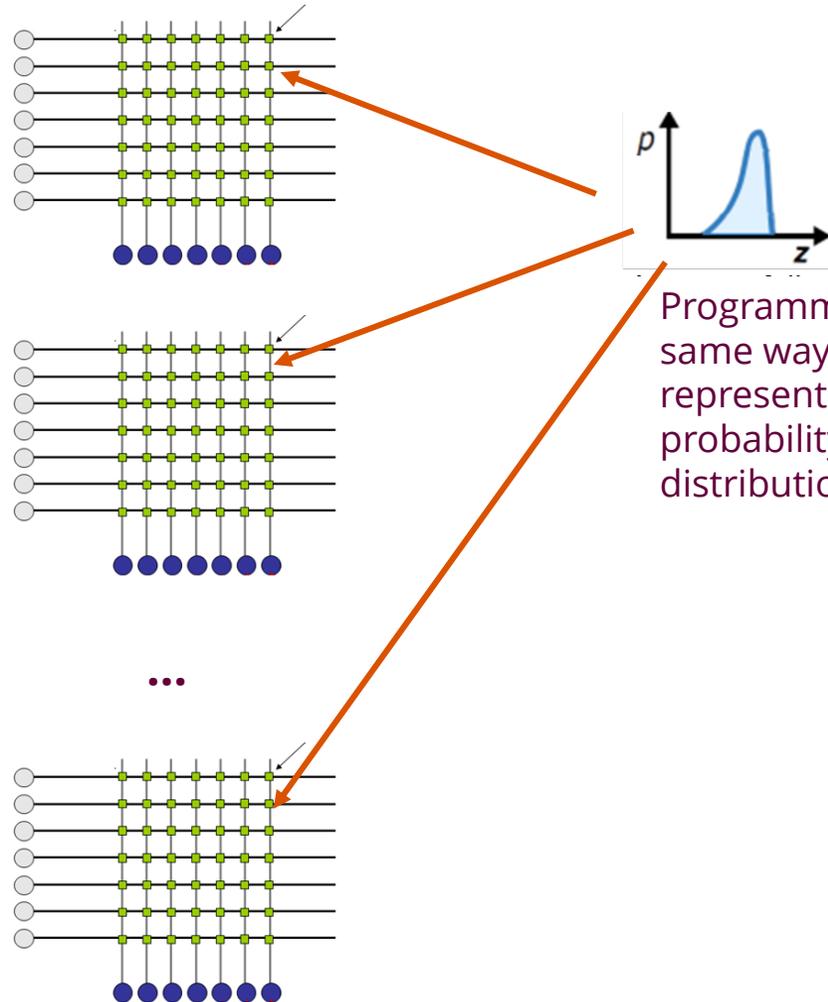
**Our concept: Bayesian models can be a "better" way to exploit memristors**

# Now Let Us Make a Bayesian Version!

We program 50 memristor arrays, and we apply the same input to them



Input: heartbeat to classify

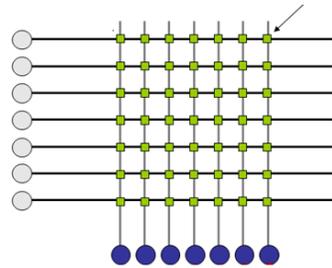


We get 50 outputs:  
Their dispersion tells about  
the *certainty* of the neural  
network

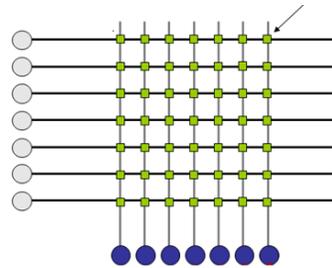
Neural networks needs to have been  
trained in a special way (variational  
inference with technological loss)

# Now Let Us Make a Bayesian Version!

We program 50 memristor arrays, and we apply the same input to them

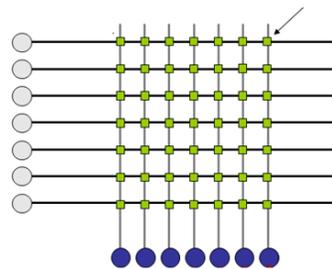


*Hesitates between  
arrhythmia type 1  
and 2*

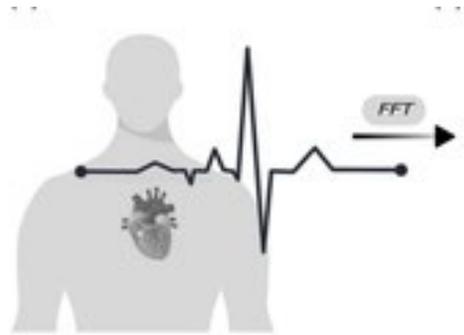


*Hesitates between  
arrhythmia type 1  
and 2*

...



*Hesitates between  
arrhythmia type 1  
and 2*



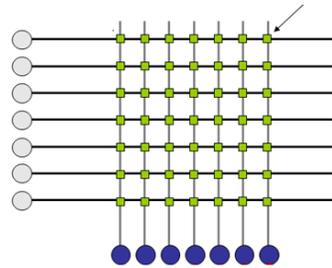
Input: heartbeat to  
classify

**Ambivalence between  
two types or  
arrhythmias:**

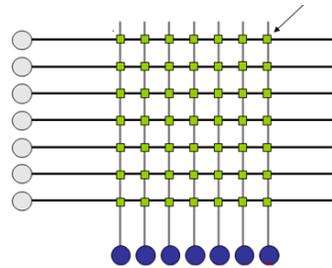
**ALEATORIC  
UNCERTAINTY**

# Now Let Us Make a Bayesian Version!

We program 50 memristor arrays, and we apply the same input to them

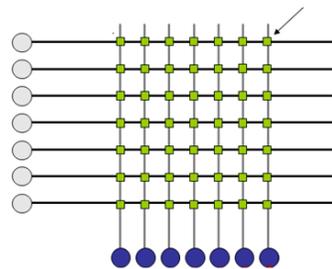


*Hesitates between  
arrythmia type 1  
and 2*

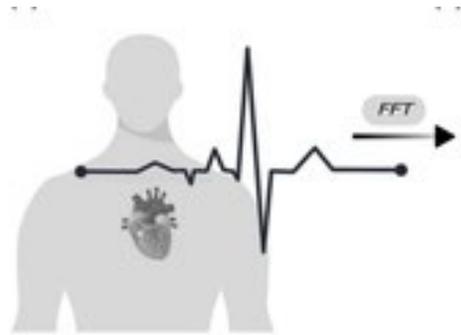


*Hesitates between  
arrythmia type 3  
and 4*

...



*Hesitates between  
arrythmia type 1  
and 5*



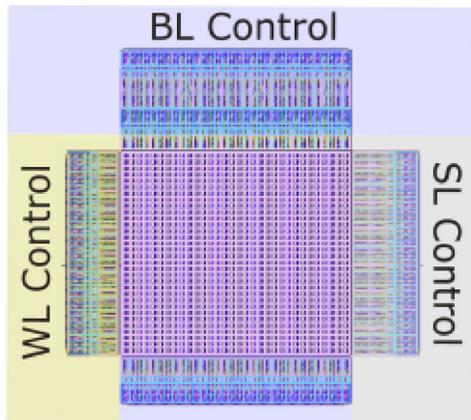
Input: heartbeat to  
classify

**Unknown data:**

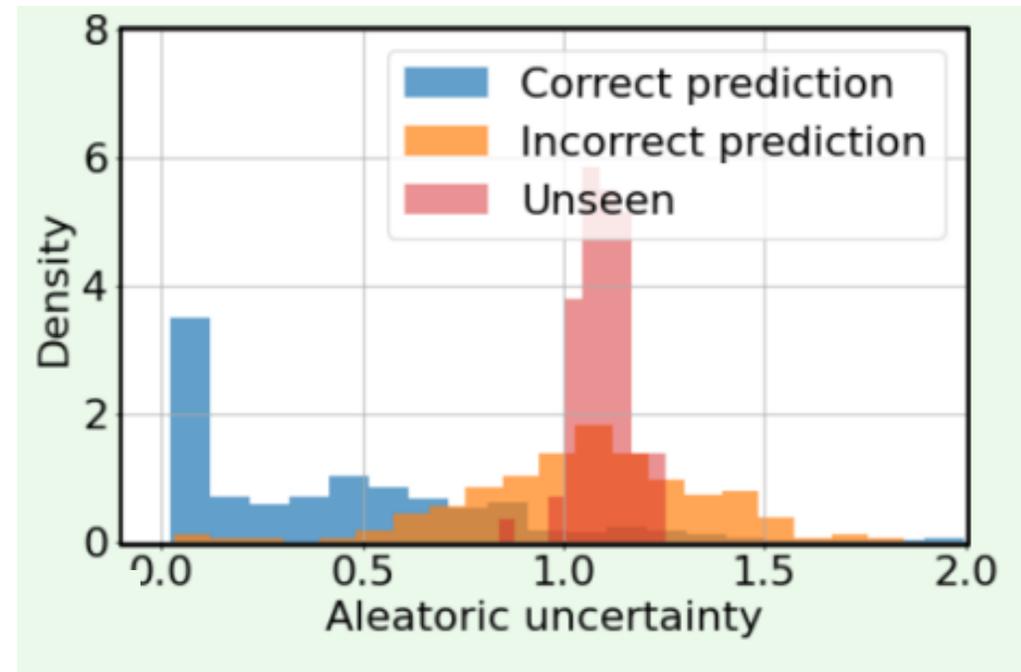
**EPISTEMIC  
UNCERTAINTY**

# Fully Experimental Arrhythmia Recognition with Uncertainty Evaluation

- **79% raw accuracy (software: 80%)**
- **Unknown types of arrhythmia are easily recognized**

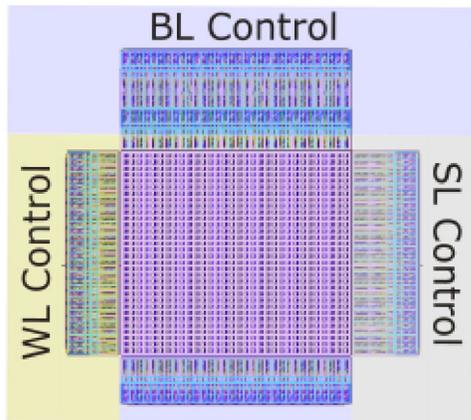


Using 50 memristors neural networks  
« Wait and Verify » Programming not needed

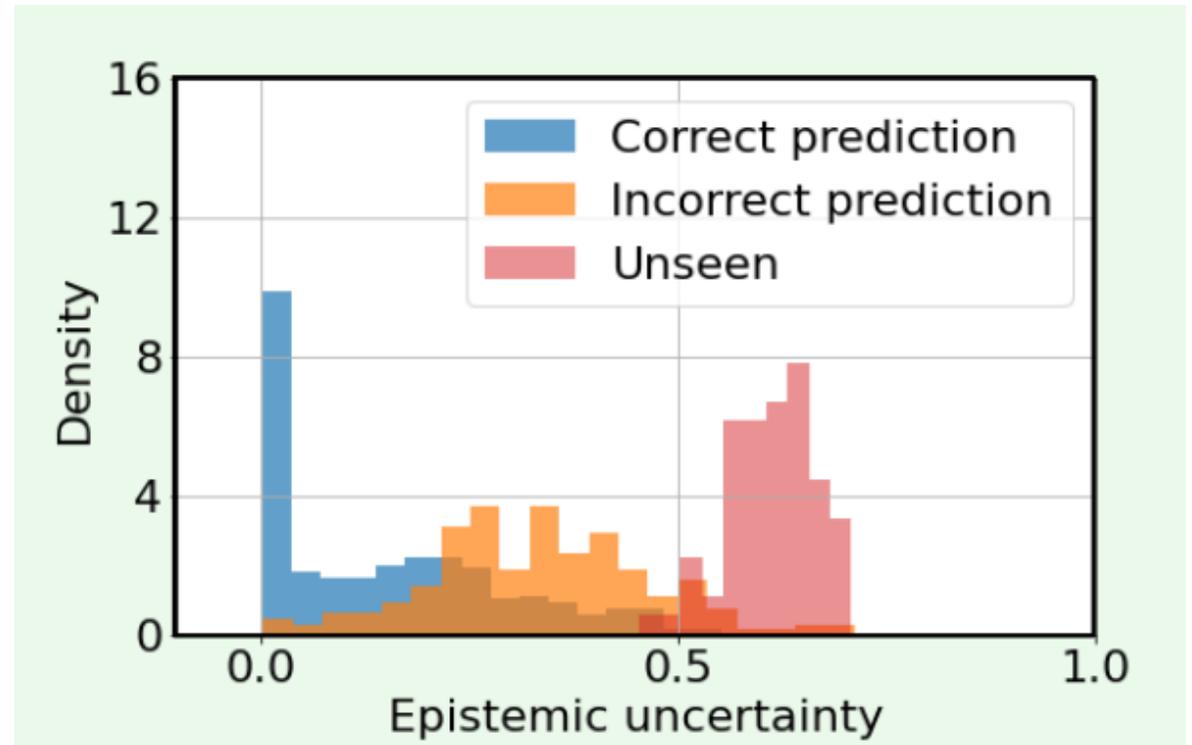


# Fully Experimental Arrhythmia Recognition with Uncertainty Evaluation

- **79% raw accuracy (software: 80%)**
- **Unknown types of arrhythmia are easily recognized**

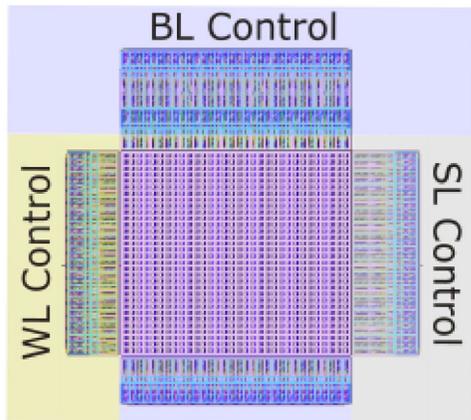


Using 50 memristors neural networks  
« Wait and Verify » Programming not needed

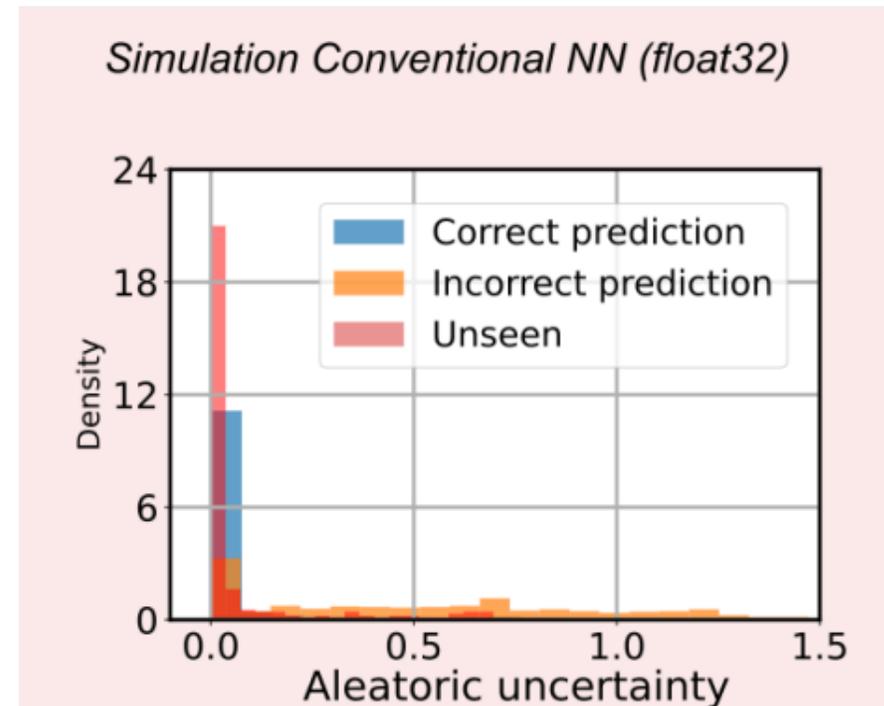


# Fully Experimental Arrhythmia Recognition with Uncertainty Evaluation

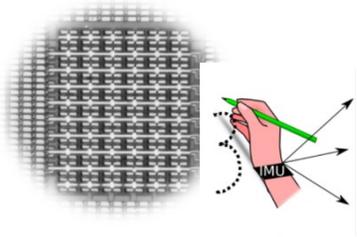
- **79% raw accuracy (software: 80%)**
- **Unknown types of arrhythmia are easily recognized**



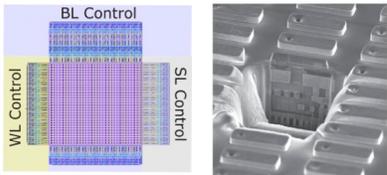
Using 50 memristors neural networks  
« Wait and Verify » Programming not needed



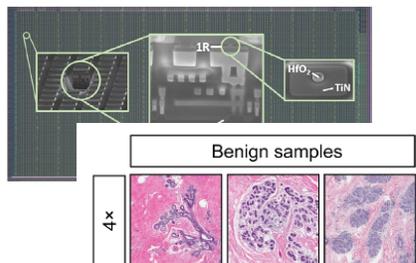
# Bayesian In-Memory Computing



- The Memristor-Based Bayesian Machine

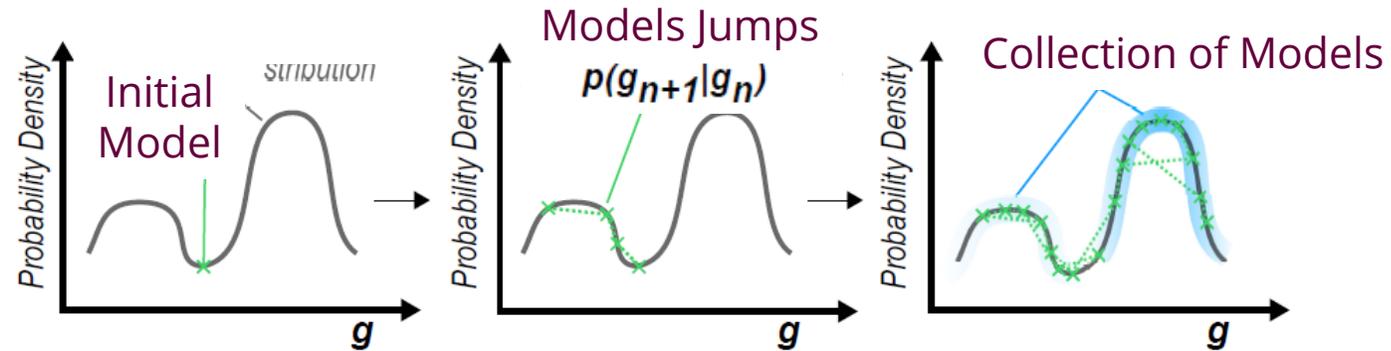


- Bayesian Neural Networks with Memristors

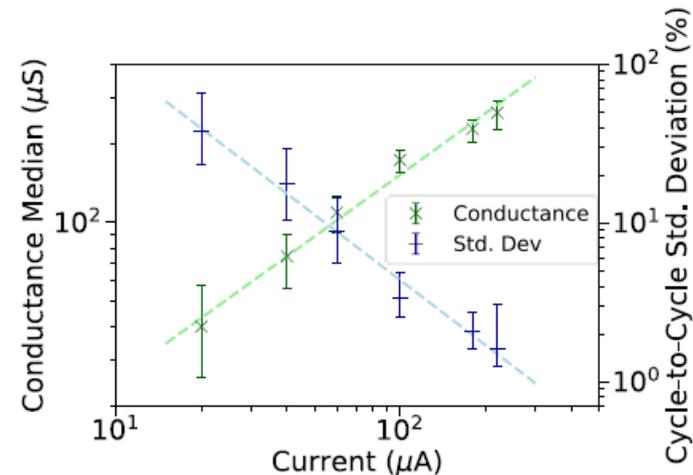
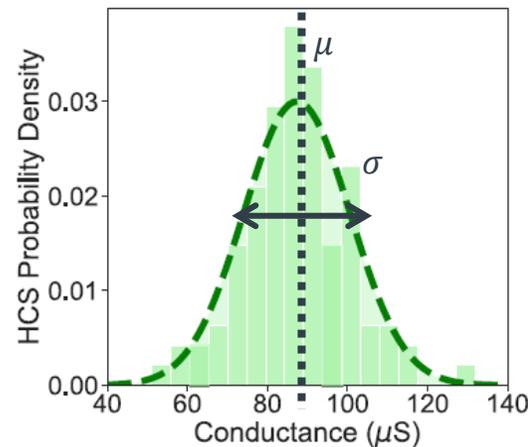
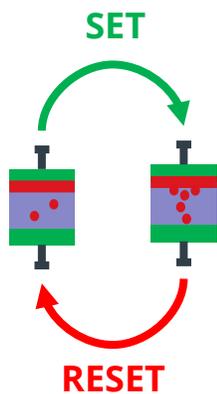


- Bayesian Learning with Memristors

# Bayesian Models Can Learn Using Metropolis-Hastings Markov Chain Monte Carlo (MCMC)

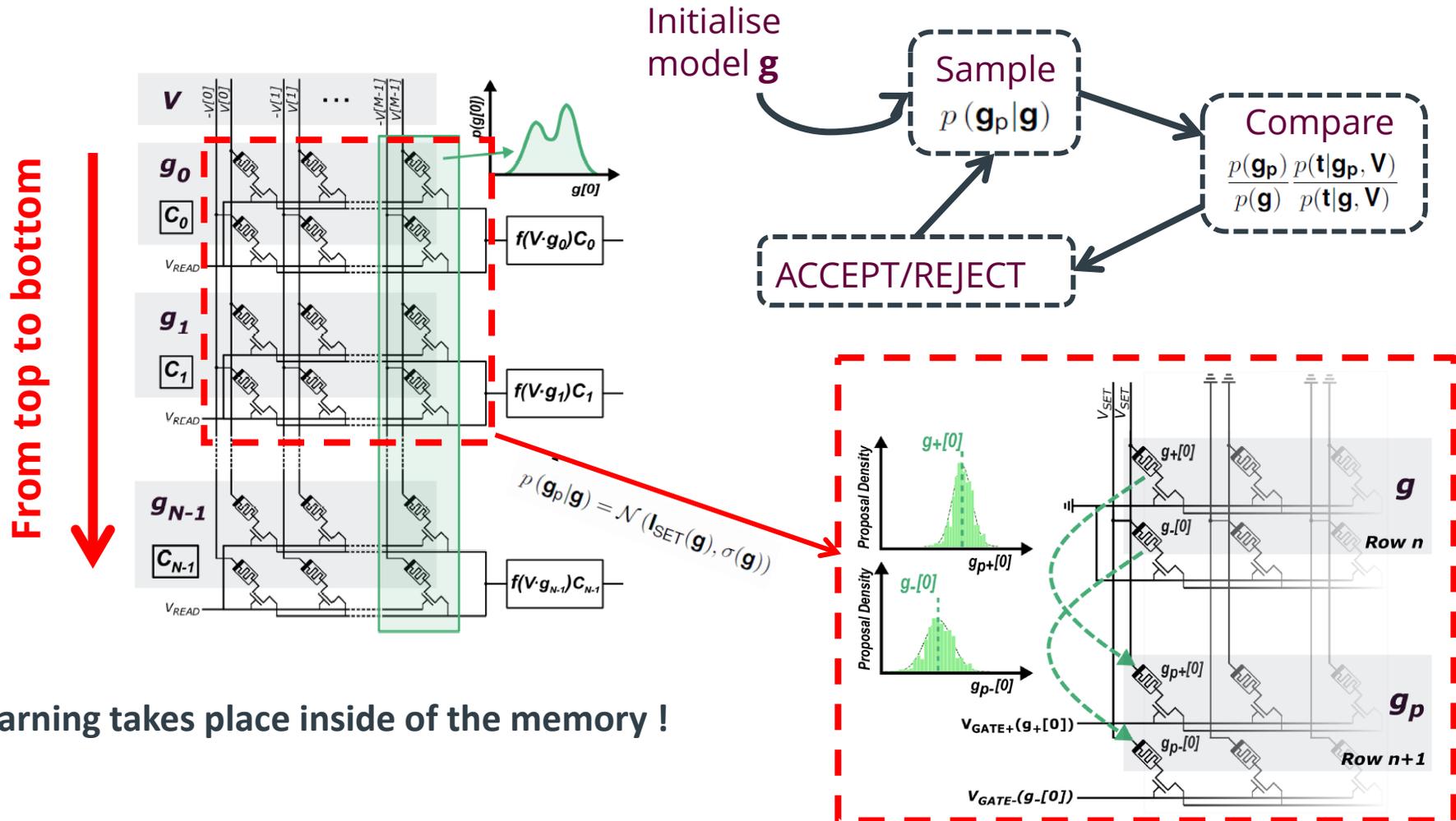


The jumps  $p(g_{n+1}|g_n)$  can be performed easily using the statistical behavior of memristors!



Memristors are ideal for MCMC-based learning!

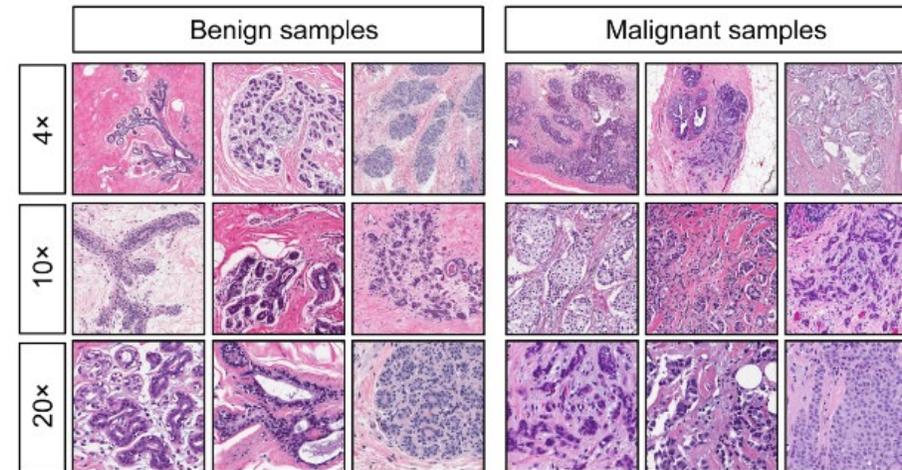
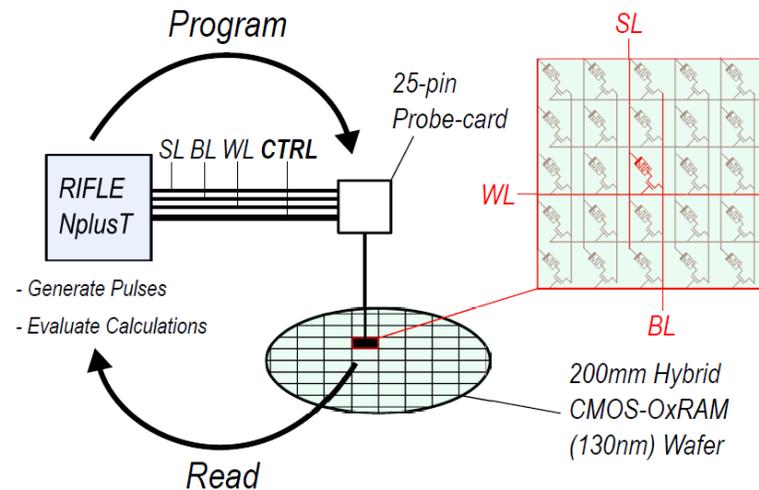
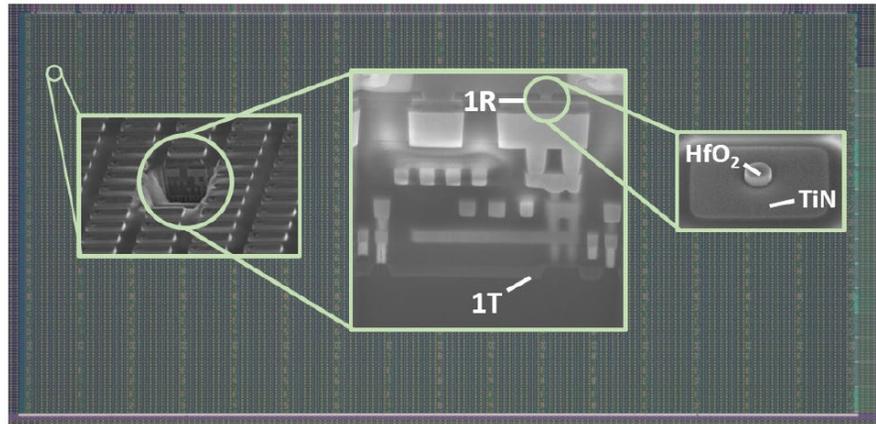
# Memristor-Based MCMC in Practice



Learning takes place inside of the memory !

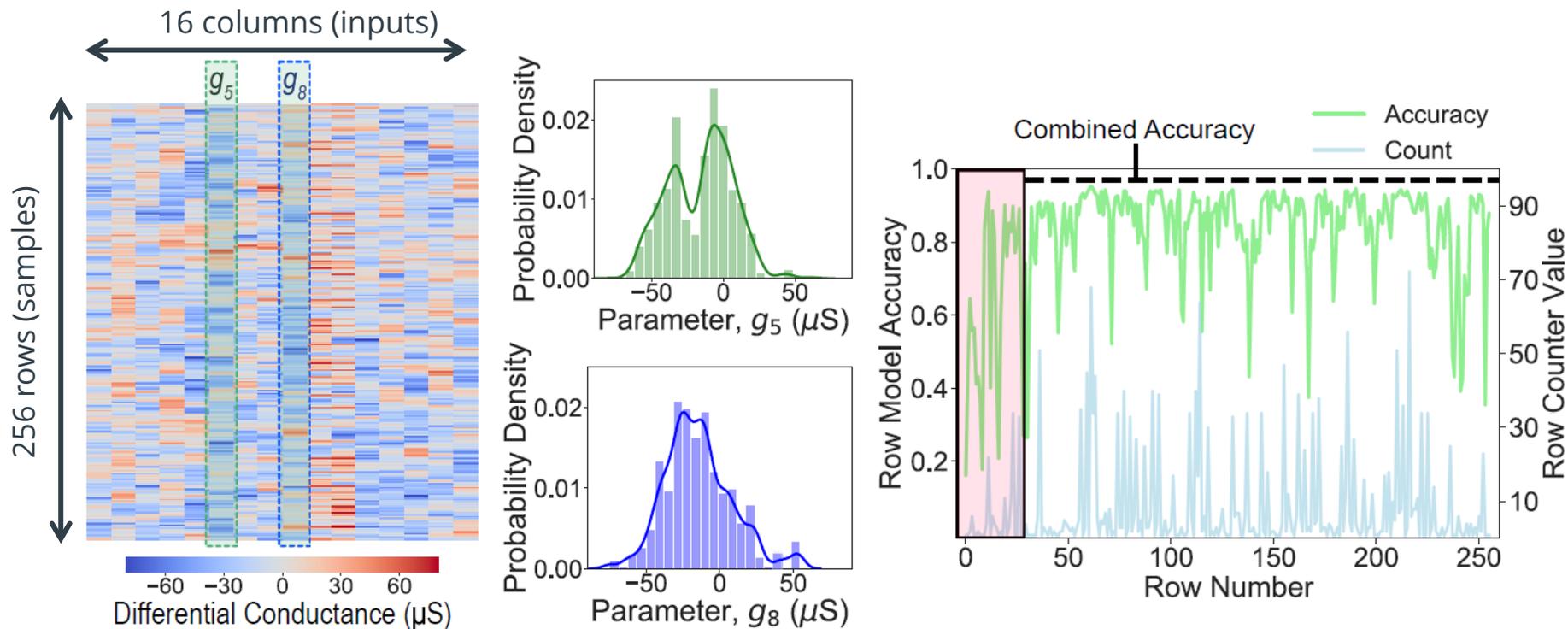
# Memristor-Based MCMC in Practice

Computer-in-the-loop experiment with an array of 16,384 memristors



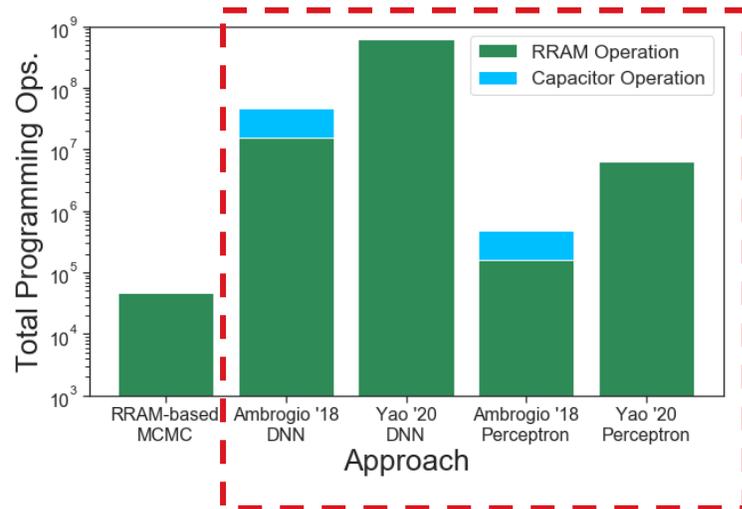
Mangasarian, O. L., et al (1995). *Operations Research*, 43(4), 570-577.

# Supervised Learning with Memristor-Based MCMC



**The experimental system was able to detect malignant tissue with 98% accuracy**

# MCMC Learning Is Highly Energy-Efficient for Small Data



Memristor-based backpropagation

Preliminary mixed-signal design results (full model training)

Intel Xeon processor (7nm) implementation of MCMC sampling required **600mj**

	Step 1 (Model evaluation)	Step 2 (Model acceptance/rejection)	Step 3 (RRAM programming)	Total
Number of repetitions	500 × 10 × 512	10 × 512	10 × 512	
Total energy (130nm)	5.8μJ	120nJ	1.1μJ	<b>6.9μJ</b>
Total energy (28nm)	2.5μJ	34nJ	1.1μJ	<b>3.6μJ</b>

# Conclusion

- Nanoelectronics enables a wide range of Bayesian concepts
- The probabilistic nature of memristors can be exploited for probabilistic machine learning, i.e., Bayesian models
- This approach can be used for both learning and inference
- Particularly appropriate for “small data”/ high uncertainty situations where wrong answers have dramatic impact, e.g., medical tasks

# Acknowledgments



- Kamel-Eddine Harabi
- Cl ment Turck
- Tifenn Hirtzlin
- Atreya Majumdar
- Marie Drouhin
- Jacques-Olivier Klein



- Elisa Vianello
- Djohan Bonnet
- Thomas Dalgaty
- Tifenn Hirtzlin
- Eduardo Esmanhotto
- Niccolo Castellani
- Fran ois Andrieu



- Jacques Droulez
- Pierre Bessi re



- Rapha el Laurent



- Jean-Michel Portal
- Jean-Pierre Walder
- Marc Bocquet
- Eloi Muhr
- Fadi Jebali
- Mathieu-Coumba Faye



NEURONIC



NANOINFER

# Thank you for your attention!

 @DamienQuerlioz  
damien.querlioz@universite-paris-saclay.fr  
<https://sites.google.com/site/damienquerlioz/>