# In-Memory and Error-Immune Differential RRAM Implementation of Binarized Deep Neural Networks

M. Bocquet*[†], T. Hirtzlin*[‡], J.-O. Klein[‡], E. Nowak[§], E. Vianello[§], J.-M. Portal[†] and D. Querlioz[‡]

[†] Aix Marseille Univ, Université de Toulon, CNRS, IM2NP, Marseille, France

[‡] C2N, Univ Paris-Sud, CNRS, Orsay, France

[§] CEA, LETI, Grenoble, France

*These authors contributed equally to the work

# In-Memory and Error-Immune Differential RRAM Implementation of Binarized Deep Neural Networks

M. Bocquet[*†], T. Hirtzlin[*‡], J.-O. Klein[‡], E. Nowak[§], E. Vianello[§], J.-M. Portal[†] and D. Querlioz[‡]

[†] Aix Marseille Univ, Université de Toulon, CNRS, IM2NP, Marseille, France

[‡] C2N, Univ Paris-Sud, CNRS, Orsay, France

[§] CEA, LETI, Grenoble, France

*These authors contributed equally to the work

# The Energy Challenge of Artificial Intelligence

- Deep learning
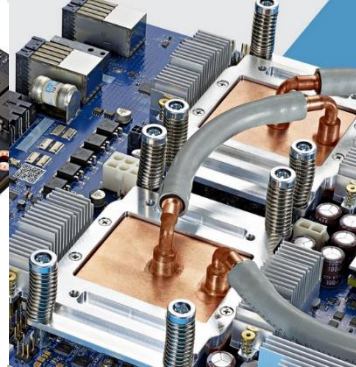

Mask R-CNN

- Energy consumption


Amazon Data-center
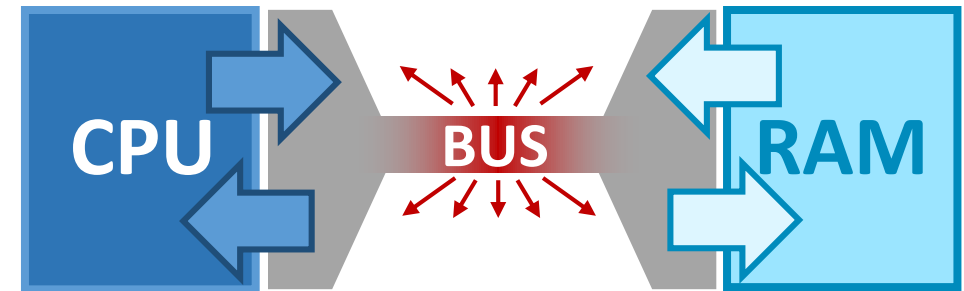
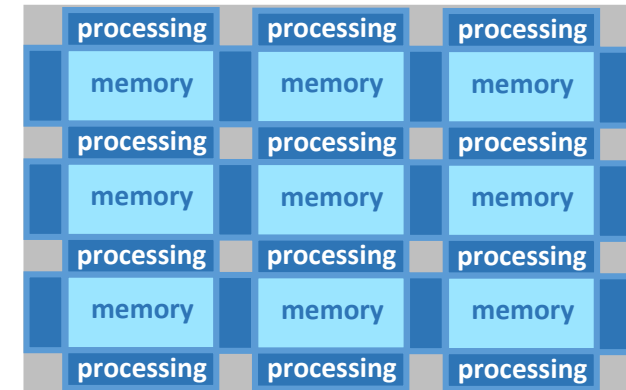- CPU / GPU / TPU > 100W


Intel Xeon CPU   NVIDIA GPU   TPU 3 Google
Titan V100

- von Neumann bottleneck


CPU   BUS   RAM

# Beyond von Neumann

100W → Incompatible with IoT

Microsoft IoT

- In memory computing

| processing | processing | processing |
| memory | memory | memory |
| processing | processing | processing |
| memory | memory | memory |
| processing | processing | processing |
| memory | memory | memory |
| processing | processing | processing |

- New **Non-Volatile** memory technology

**Electrode**
GST
TiN SiO₂
**Electrode**

**PCM**

**MRAM**

Low Resistance State LRS
High Resistance State HRS
TE
BE
TE
BE
**RRAM**

# NEURAL NETWORKS seems Especially Adapted for In Memory Computing



- Operations in Neural Networks :

**Multiplication → Accumulation → Non-linear function**

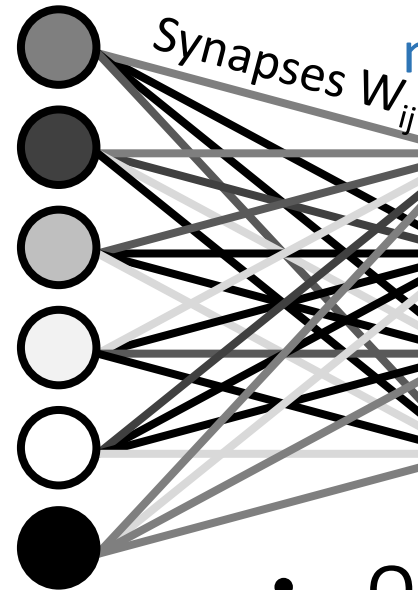Lot of research on RRAM as analog synapse

$$a_i = f\left(\sum_j W_{ij} \cdot a_j\right)$$

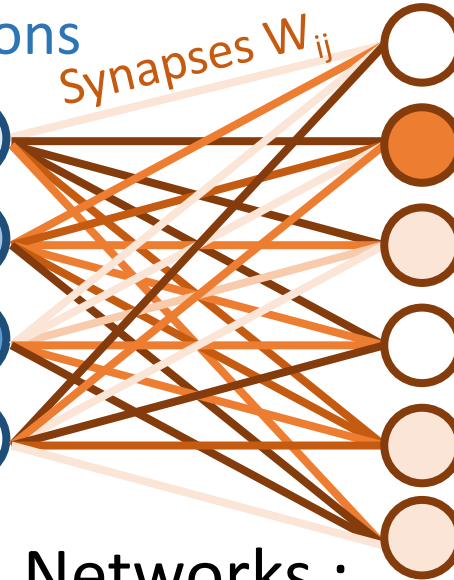# NEURAL NETWORKS seems Especially Adapted for In Memory Computing



- Operations in Neural Networks :

**Multiplication → Accumulation → Non-linear function**
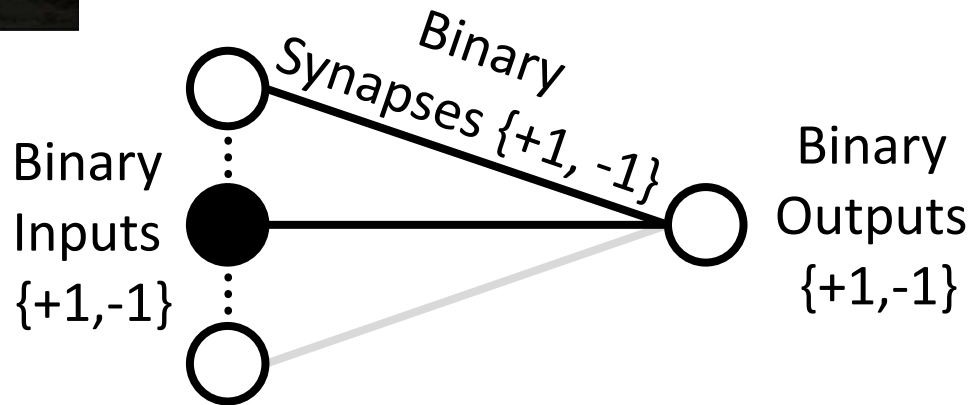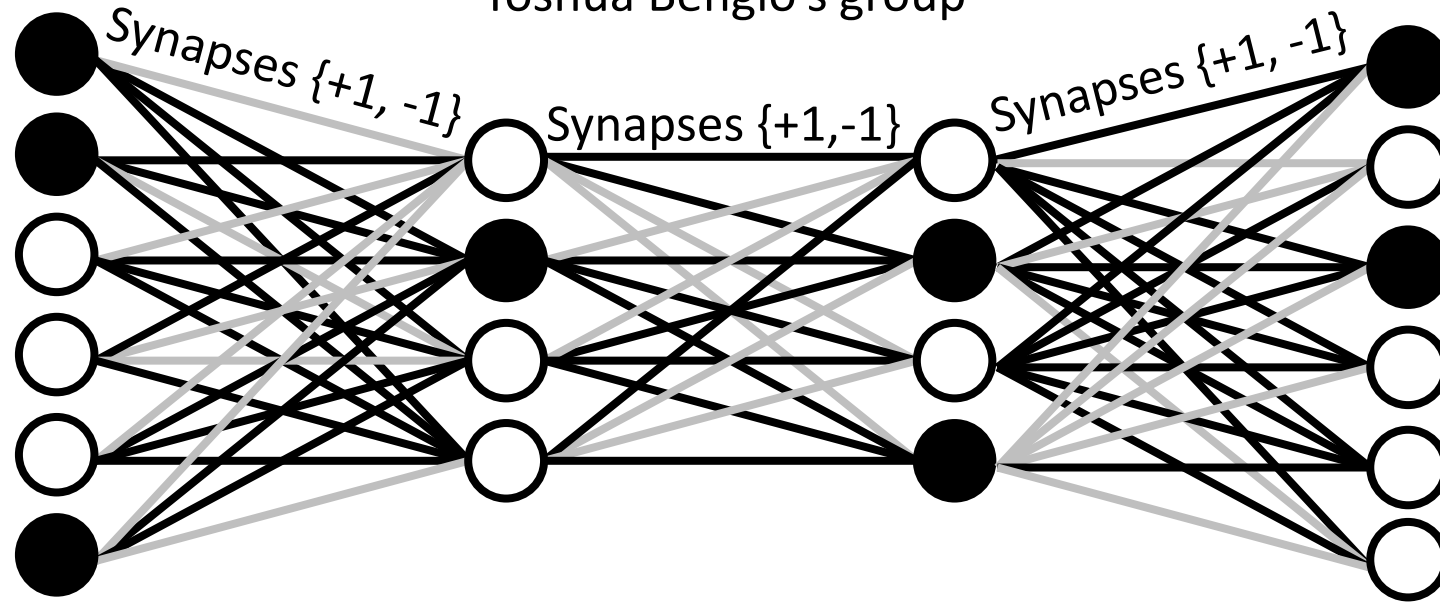
Lot of research on RRAM as analog synapse

Our work : RRAM as binary synapse

# Recent Breakthrough: Binarized Neural network



*Hubara, Courbariaux et al. NIPS 2016*
Yoshua Bengio's group

Synapses {+1, -1}

Synapses {+1,-1}

Synapses {+1, -1}

leopard
jaguar
cheetah
snow leopard
Egyptian cat

motor scooter
go-kart
moped
bumper car
golfcart

Binary Inputs {+1,-1}

Binary Synapses {+1, -1}

Binary Outputs {+1,-1}

Our work :
RRAM as binary synapse

# Binarized Neural network:
# Very Simple Logical Operations

| Multiplication | $\rightarrow$ | Accumulation | $\rightarrow$ | Non-linear function |
|:---:|:---:|:---:|:---:|:---:|
| $W_{ij} . a_j$ | | $\sum_j W_{ij} . a_j$ | | $a_i = f\left(\sum_j W_{ij} . a_j\right)$ |
| XNOR | $\rightarrow$ | Bitcount | $\rightarrow$ | Sign |



| $a_j$ | $W_{ij}^b$ | $W_{ij}^b . a_j$ |
|:---:|:---:|:---:|
| -1 | -1 | 1 |
| -1 | 1 | -1 |
| 1 | -1 | -1 |
| 1 | 1 | 1 |

# RRAM Technology Involved

- **HfO$_2$-based OxRAM** integrated in a 130 nm CMOS logic process
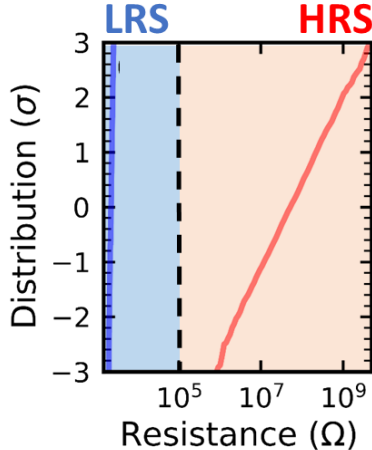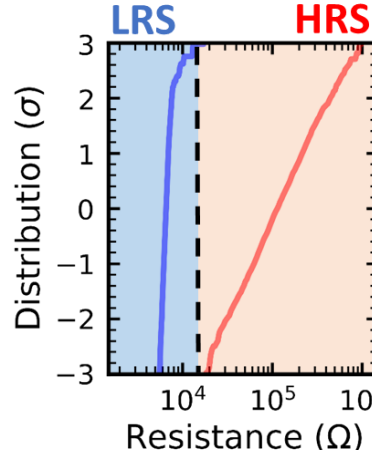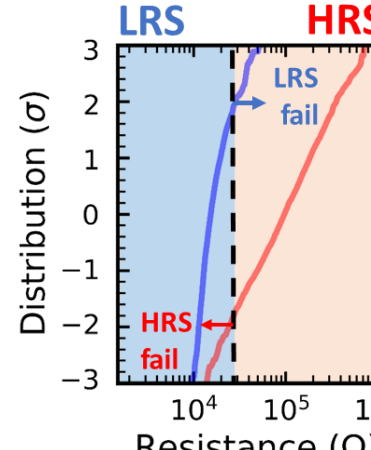- Fast / High retention time / High endurance
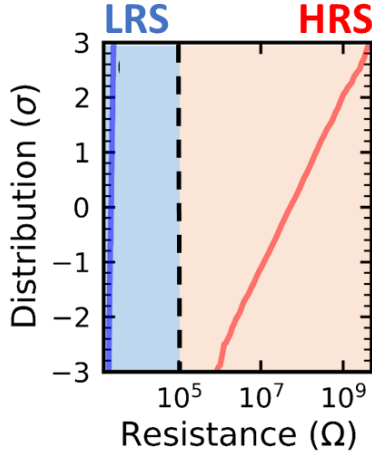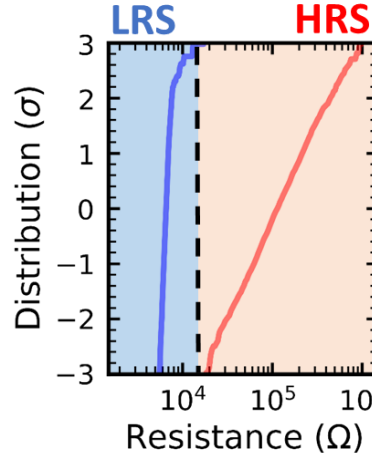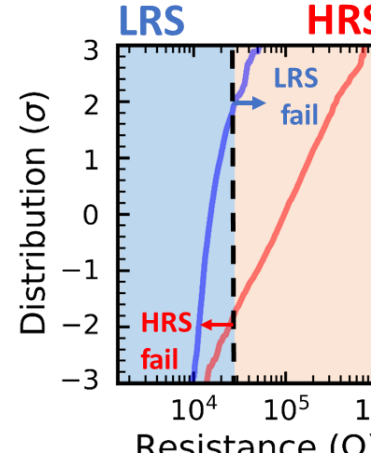


SEM cross-section – TiN/HfO$_2$/Ti/TiN

**LETI MPW Shuttle managed through CMP ( today in 200mm; twice a year; next start may 2019)**

RRAM Challenge → **high variability**

# RRAM Variability Depends Extensively on the Programming Regime

| Programming condition | Very strong | Strong | Weak |
|---|---|---|---|
| SET compliance current | 600µA | 55µA | 20µA |
| RESET voltage | 2.5V | 2.5V | 1.5V |
| Resistance Distribution |  |  |  |
| **Bit error rate (1T1R)** | $< 10^{-6}$ | $9.7 \times 10^{-5}$ | $3.3 \times 10^{-2}$ |
| Prog. Energy (SET/RESET) | $120/150pJ$ | $11/14pJ$ | $4/5pJ$ |
| **Cyclability** | 100 | > 10000 | $> 10^6$ |

# RRAM Variability Depends Extensively on the Programming Regime

| Programming condition | Very strong | Strong | Weak |
|---|---|---|---|
| SET compliance current | 600μA | 55μA | 20μA |
| RESET voltage | 2.5V | 2.5V | 1.5V |
| Resistance Distribution |  |  |  |
| **Bit error rate (1T1R)** | $< 10^{-6}$ | $9.7 \times 10^{-5}$ | $3.3 \times 10^{-2}$ |
| Prog. Energy (SET/RESET) | 120/150$pJ$ | 11/14$pJ$ | 4/5$pJ$ |
| **Cyclability** | 100 | > 10000 | $> 10^6$ |

Exploiting weak programming conditions for BNN

# How to deal with bit errors?

- Classical Approach : Error Correction Code
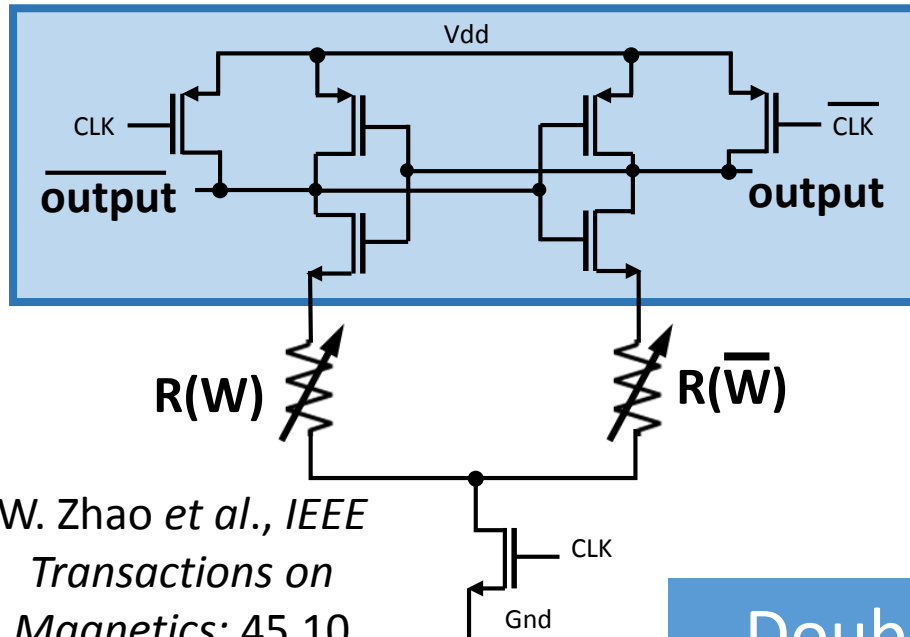
→ **Incompatible with in memory computing**



ECC leads to a big overhead

# Our approach : Two RRAM Devices as One Binary Synapse to reduce bit error rate
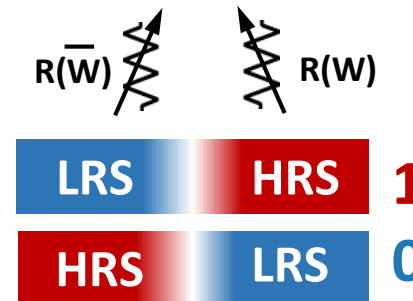


- Classical circuit to differentiate resistance state
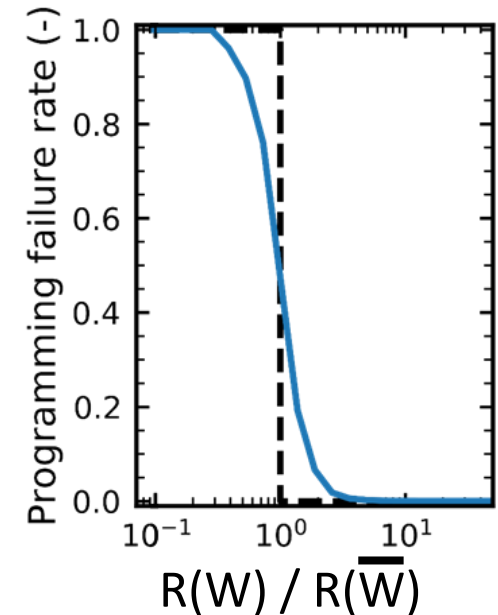
**Pre-Charge Sense Amplifier (PCSA)**



W. Zhao *et al.*, *IEEE Transactions on Magnetics*: 45,10 (2009)

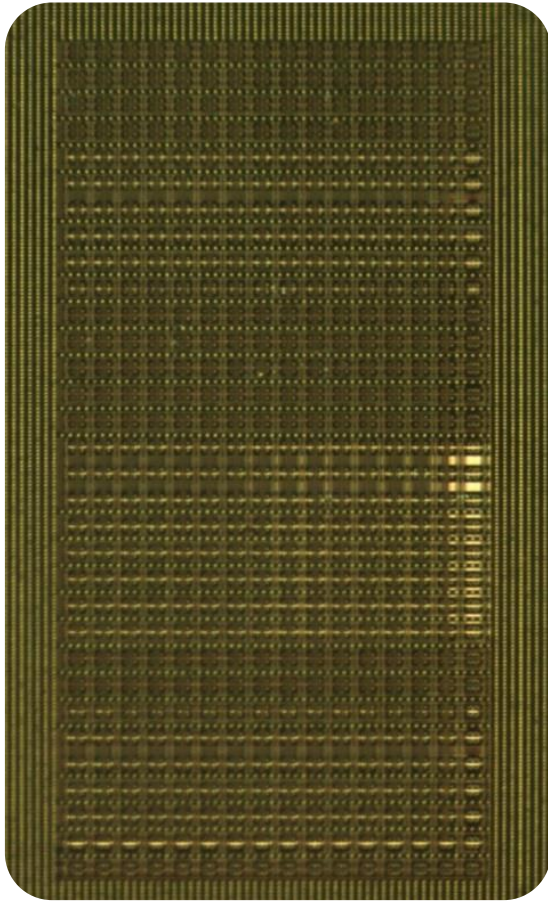- Devices programmed in a complementary fashion :



| LRS | HRS | **1** |
| HRS | LRS | **0** |

- Reading circuit behavior



Double the amount of memory

# Array structure: 2kBits devices
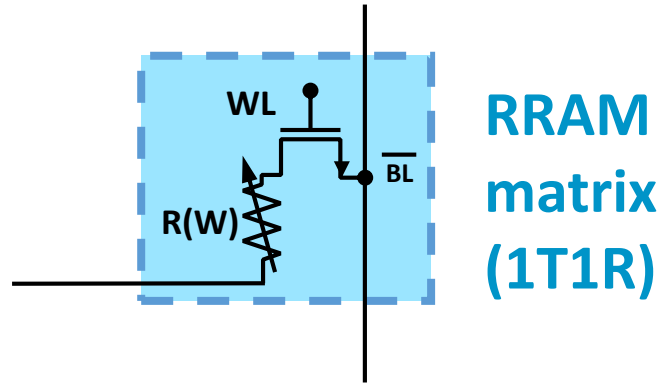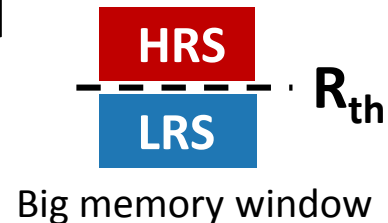
- Photograph of our circuit

- Schematic of the array

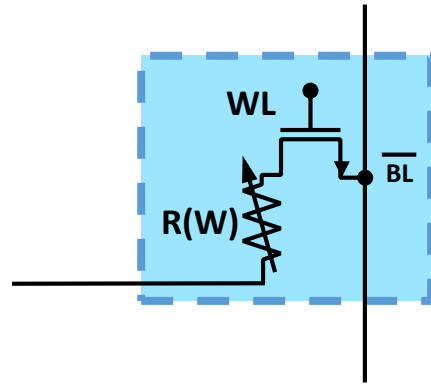# What if we used
# 1 Transistor 1 Resistor (1T1R) array structure ?



RRAM matrix (1T1R)

Devices resistance states :

R(W)

| | |
|---|---|
| HRS | 1 |
| LRS | 0 |

Resistance threshold between the two resistance states :

HRS

LRS

$R_{th}$

Big memory window

# 1 Transistor 1 Resistor (1T1R) is Prone to Errors



**RRAM matrix (1T1R)**

**R(W)**

| | |
|---|---|
| HRS | 1 |
| LRS | 0 |

Devices resistance states :

Resistance threshold between the two resistance states :

$R_{th}$

| | |
|---|---|
| HRS | |
| LRS | |

Resistances States overlap

$R(W_{LRS})$  $R(W_{HRS})$

$R_{th}$

1T1R « 1 » fail    1T1R « 0 » fail

LRS    HRS

LRS fail

HRS fail

$10^4$  $10^5$  $10^6$

Resistance ($\Omega$)

- Errors when :

$R(W_{LRS}) > Rth$ or $R(W_{HRS}) < Rth$
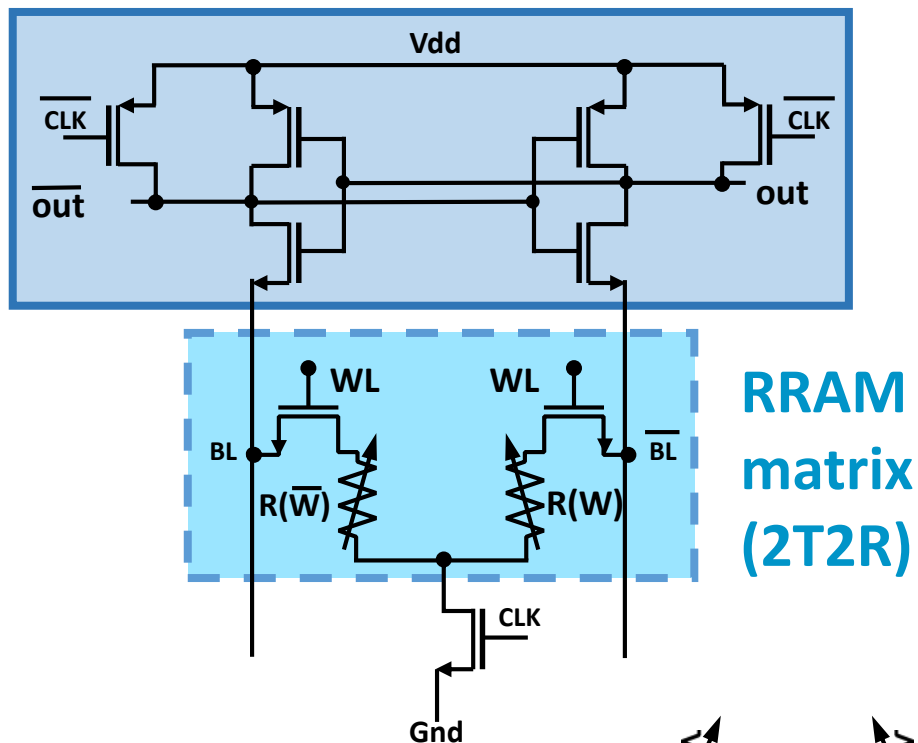
# 2 Transistors 2 Resistors (2T2R) configuration :



**(PCSA)**

**RRAM matrix (2T2R)**

Devices programmed in a complementary fashion :

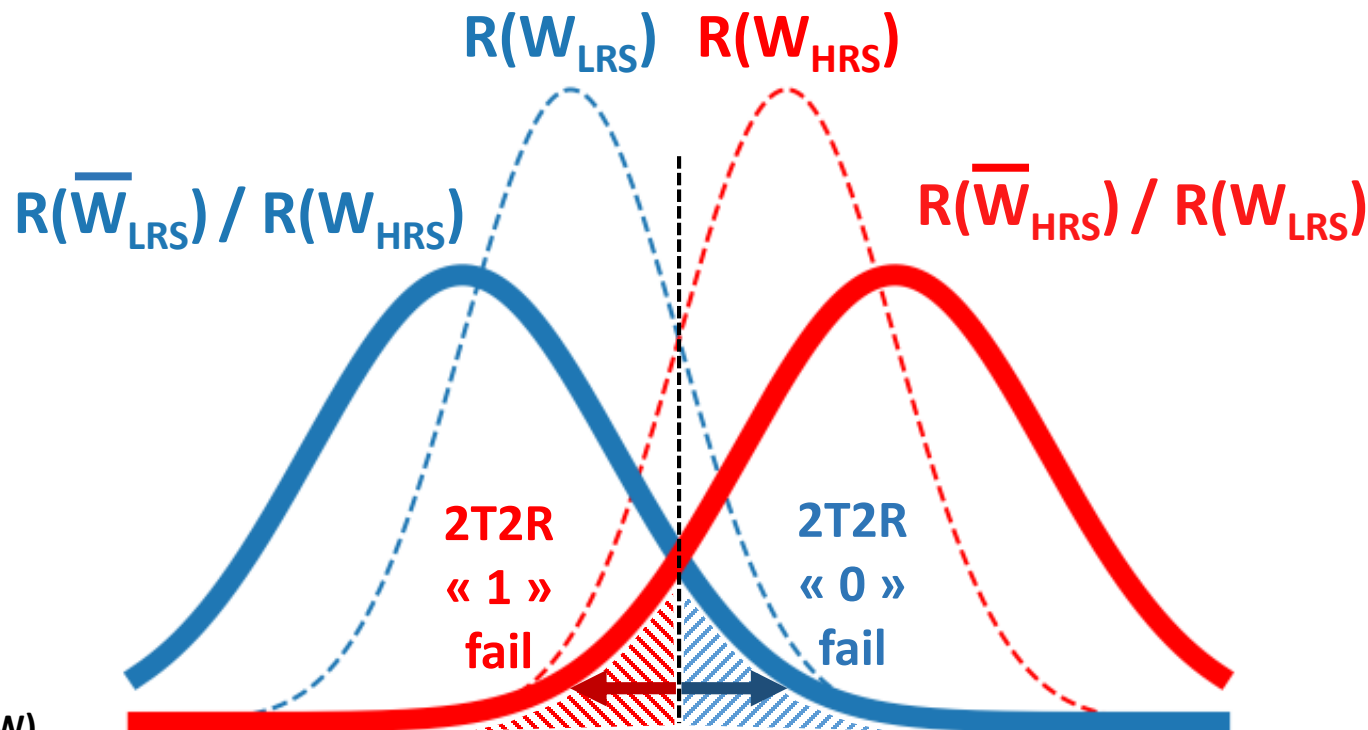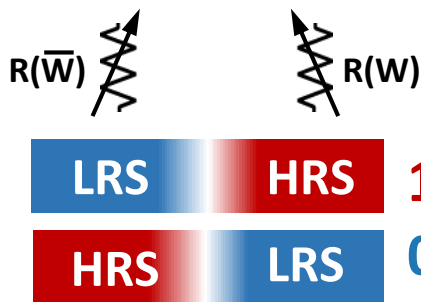| | | |
|---|---|---|
| LRS | HRS | 1 |
| HRS | LRS | 0 |

# 2 Transistors 2 Resistors (2T2R): Decreasing Error Rate



**(PCSA)**

**RRAM matrix (2T2R)**
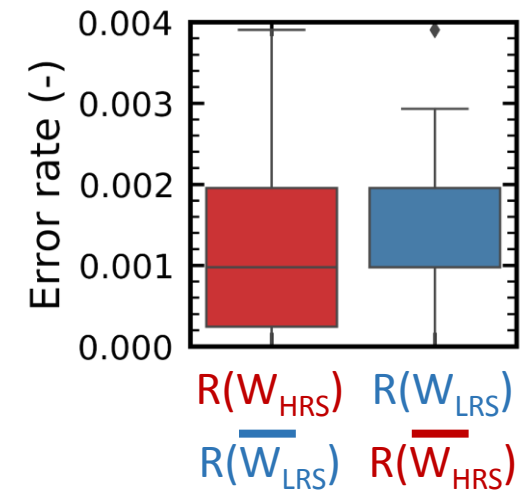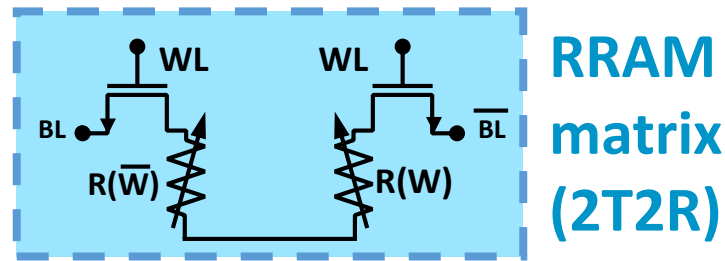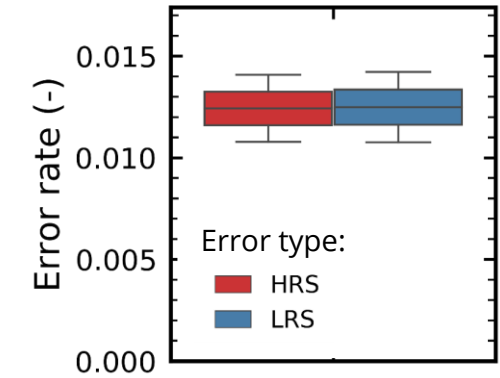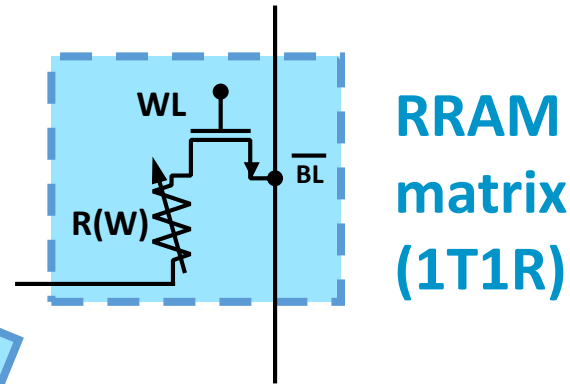
Devices programmed in a complementary fashion :

| LRS | HRS | 1 |
| HRS | LRS | 0 |

$R(W_{LRS})$ $R(\overline{W}_{HRS})$

$R(\overline{W}_{LRS}) / R(W_{HRS})$

$R(\overline{W}_{HRS}) / R(W_{LRS})$

2T2R « 1 » fail

2T2R « 0 » fail

- Errors when :

$$R(\overline{W}_{LRS}) > R(W_{HRS}) \text{ or } R(\overline{W}_{HRS}) < R(W_{LRS})$$
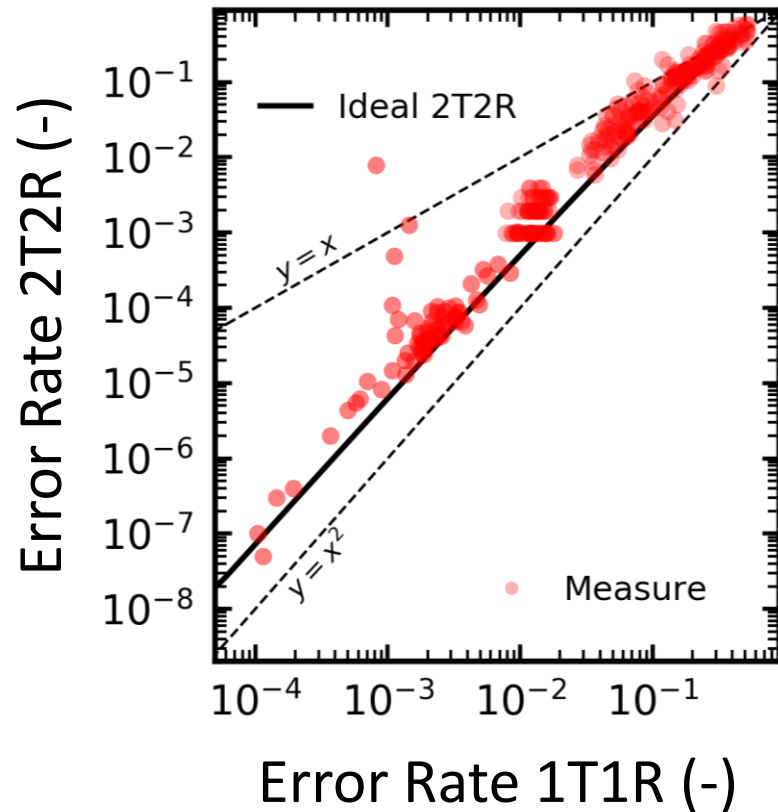
18

# Comparison between 1T1R & 2T2R
# Bit error rates extraction
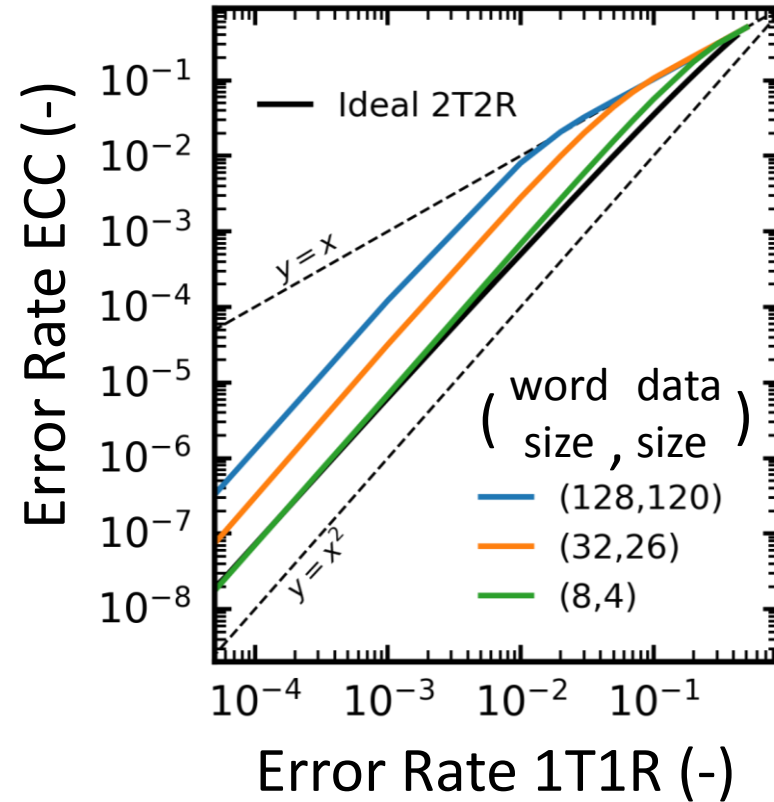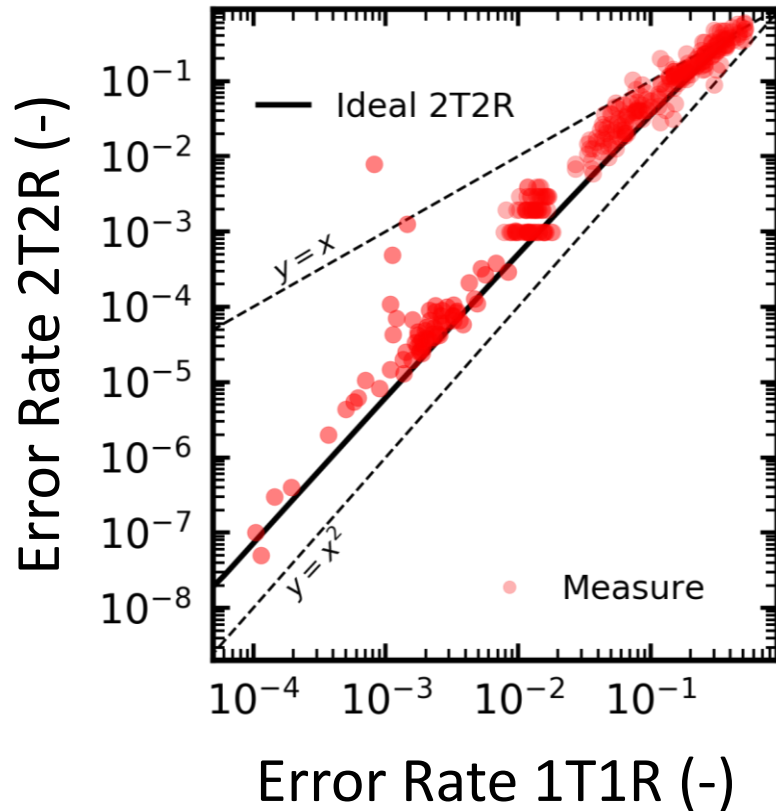
# Comparison between 1T1R & 2T2R

- Experimental bit error rate



2T2R reduces error rate

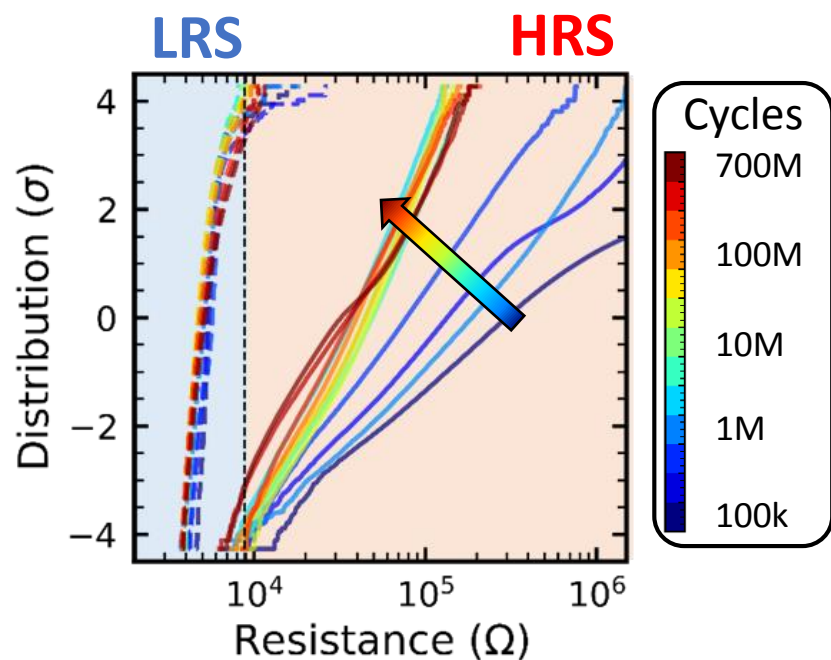# Comparison with Error Correction Code

- Experimental bit error rate
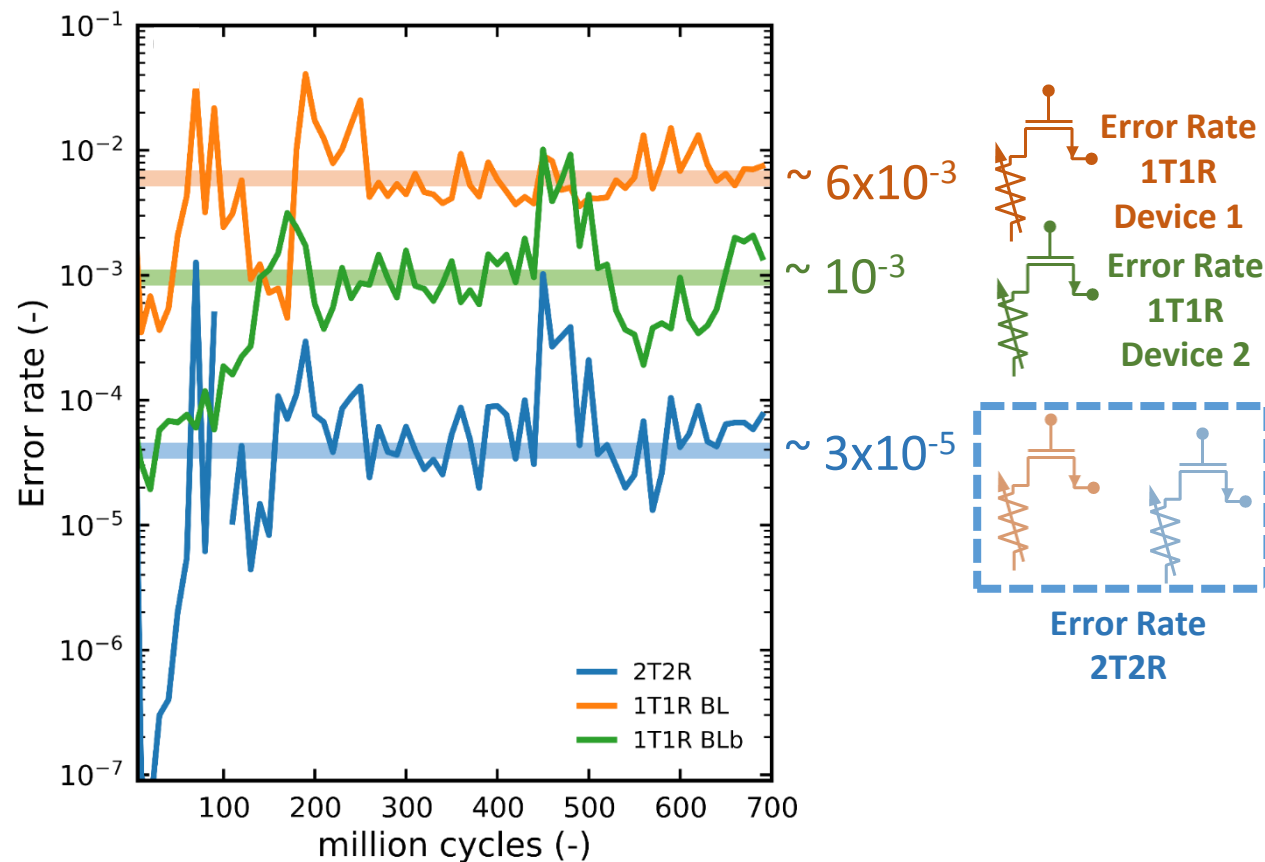
- Error Correction Code SECDED



No need of ECC overhead

# Characteristics over aging devices with weak programming conditions
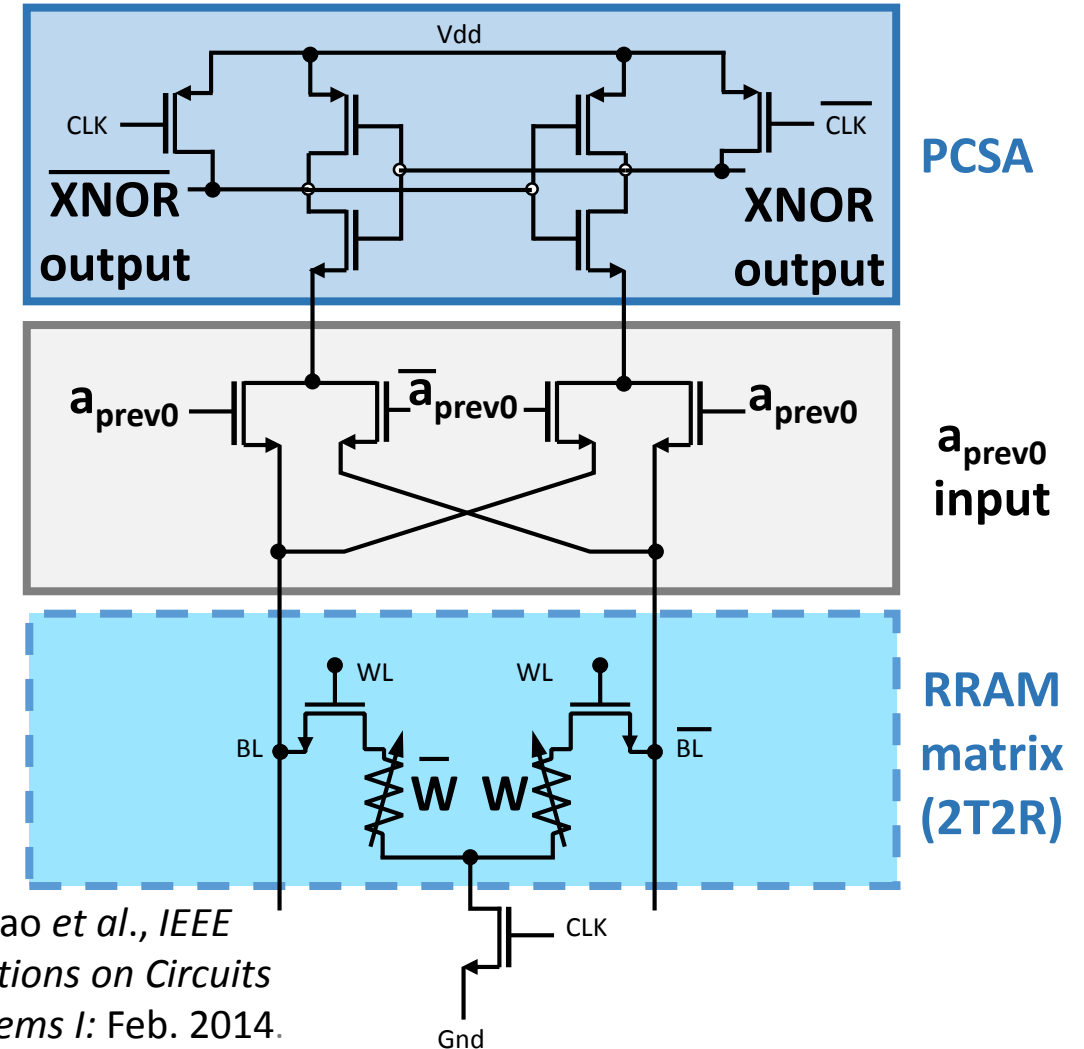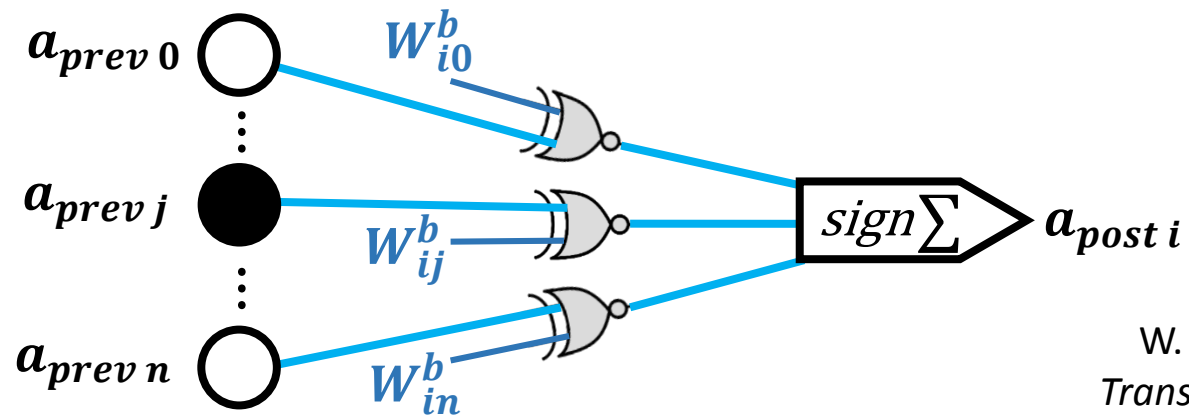
- Distribution



LRS      HRS

- Error rate for 2 devices



~ $6\times10^{-3}$ — Error Rate 1T1R Device 1

~ $10^{-3}$ — Error Rate 1T1R Device 2

~ $3\times10^{-5}$ — Error Rate 2T2R

# Logic in Memory Reading Circuit

- No ECC offers opportunity for in memory operation
- XNOR operation directly in PCSA circuit

**XNOR → Bitcount → Sign**
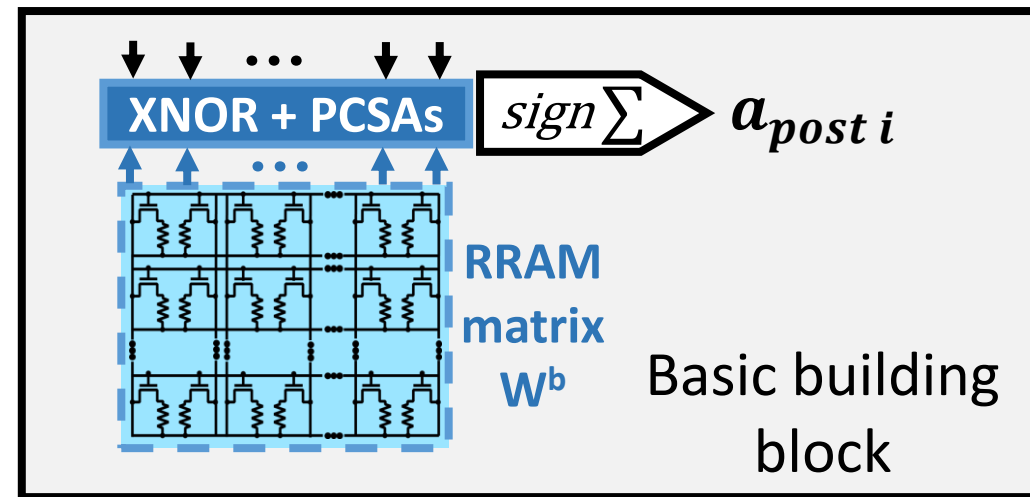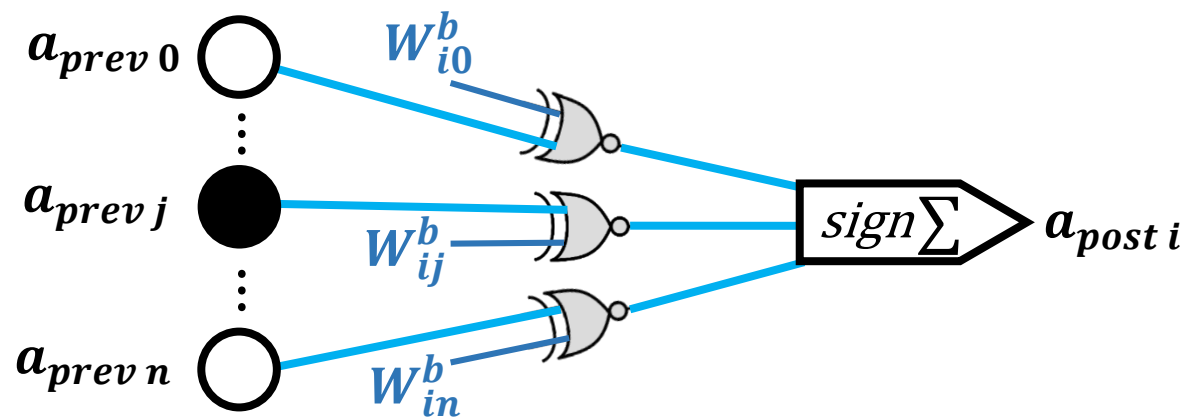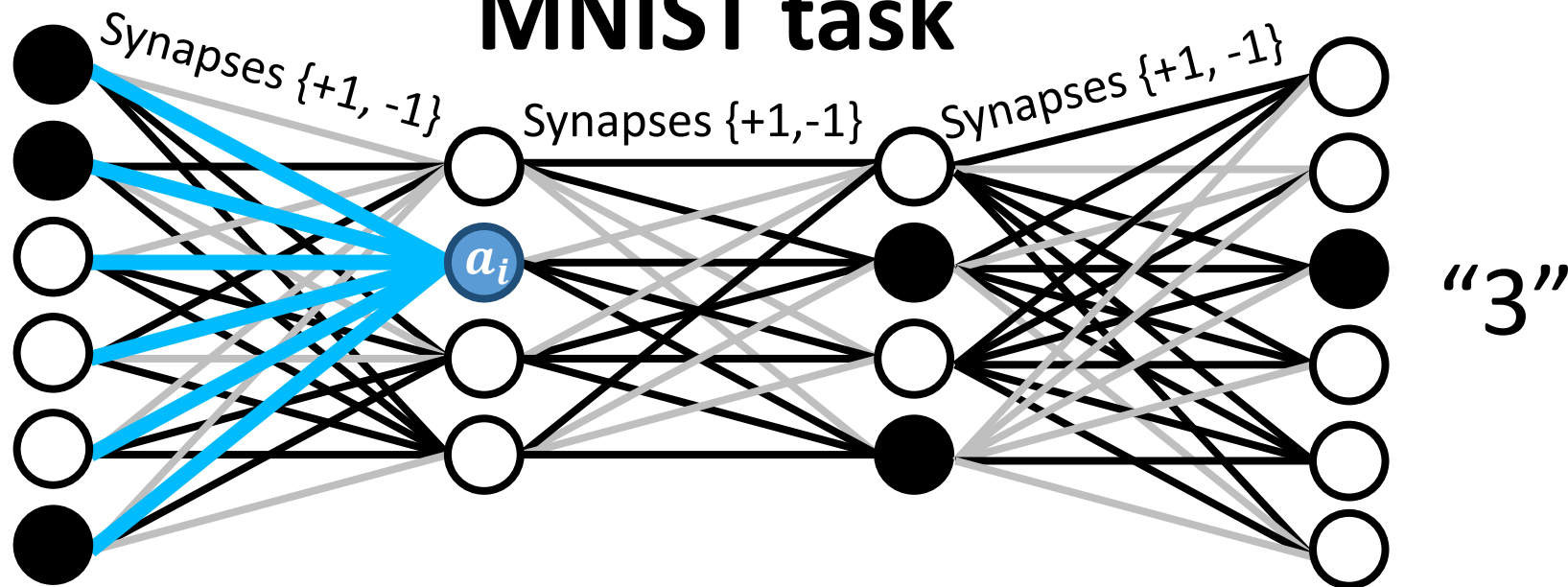


$a_{prev\,0}$  $W^b_{i0}$

$a_{prev\,j}$

$W^b_{ij}$

$a_{prev\,n}$  $W^b_{in}$

$sign\sum$  $a_{post\,i}$

W. Zhao *et al., IEEE Transactions on Circuits and Systems I:* Feb. 2014.

# Binarized Neural Network: MNIST task

Handwritten Digit example :

Synapses {+1, -1}

Synapses {+1,-1}

Synapses {+1, -1}

$a_i$

"3"

$a_{prev\ 0}$ $W_{i0}^b$

$a_{prev\ j}$

$a_{prev\ n}$ $W_{ij}^b$

$W_{in}^b$

$sign\sum$ $a_{post\ i}$

XNOR + PCSAs

$sign\sum$ $a_{post\ i}$

RRAM matrix $W^b$

Basic building block

# Hardware implementation

- Fully-connected BNN:
2 layers with 1024 neurons each

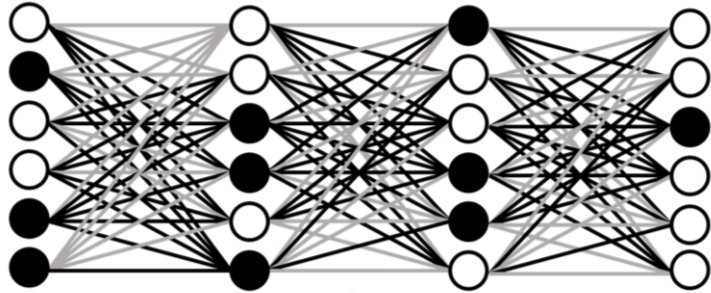| RRAM In-memory BNN (this work) | 25nJ |
|---|---|
| RRAM In-memory 8-bit fixed point | 80nJ |
| Analog Phase Change Memory* | ~56nJ |
| GPU (Tesla V100) | ~µJ |
| CPU (Xeon E5) | ~mJ |

- Projection to 28 nm technology
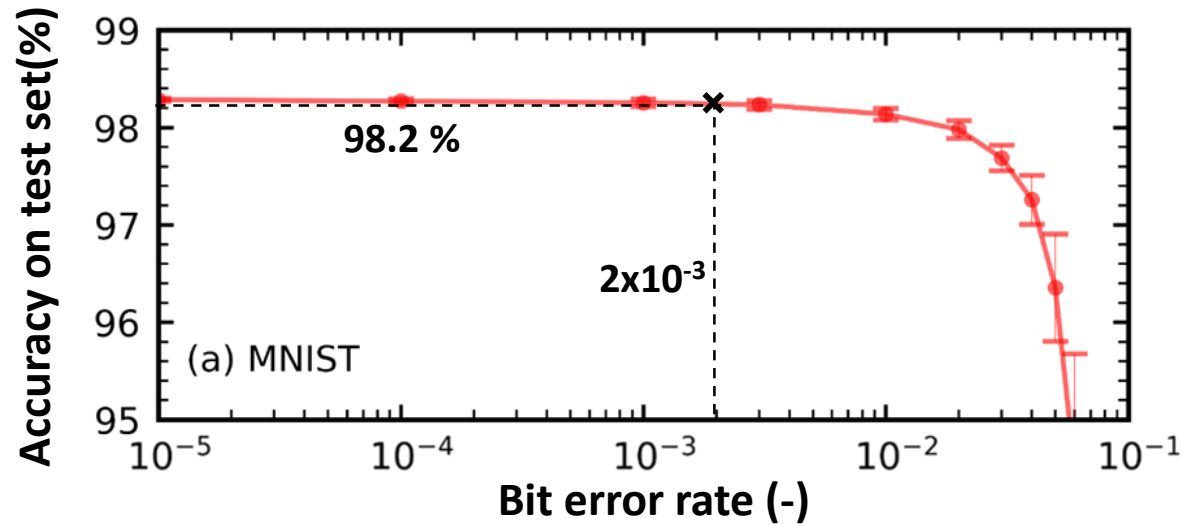
**A technology ready**
↳ low energy and low area

*Ambrogio, Stefano, et al. "Equivalent-accuracy accelerated neural-network training using analogue memory." *Nature* (2018)

# Error evaluation on two different tasks
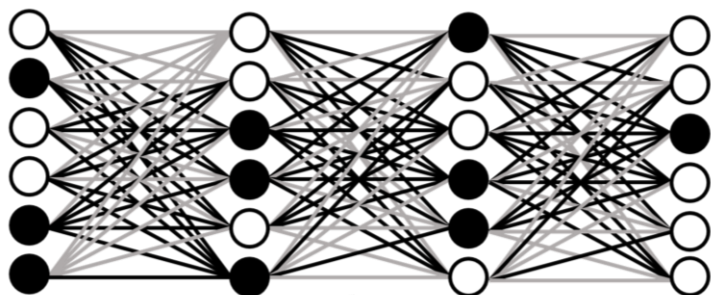
- MNIST with fully-connected NN



(a) MNIST

Accuracy on test set(%)
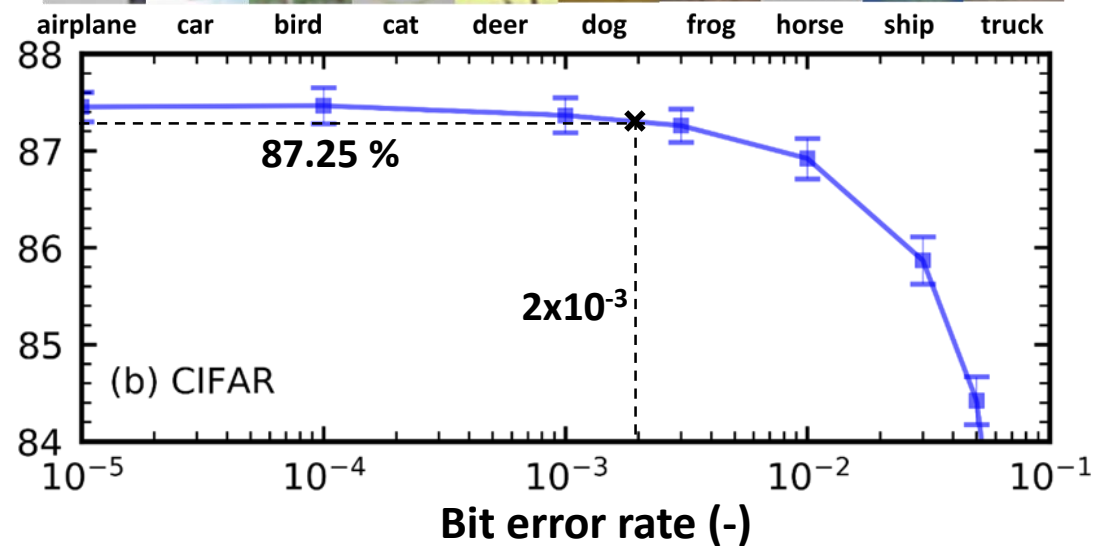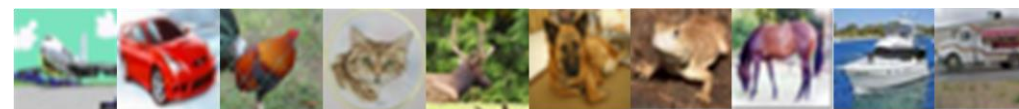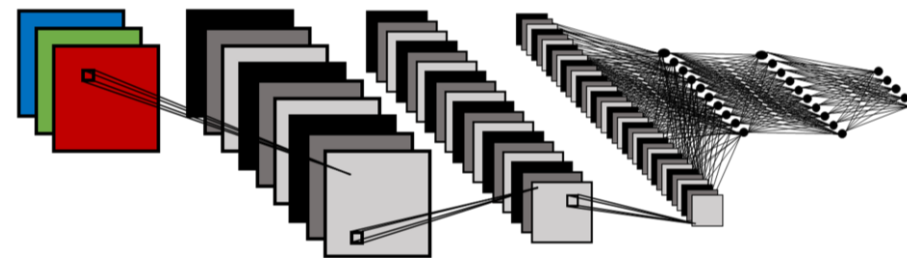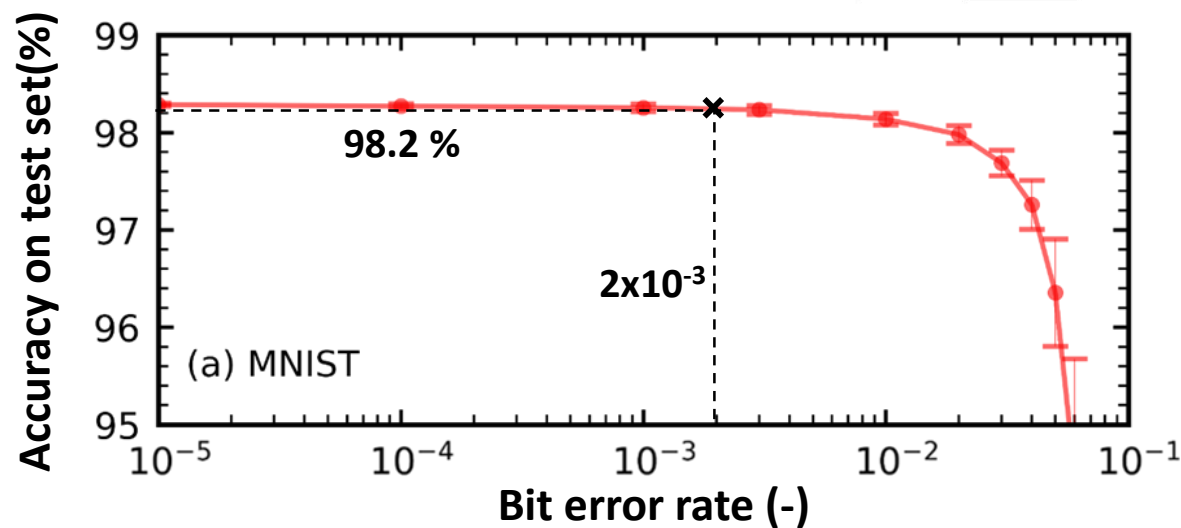
Bit error rate (-)

98.2 %

2x10⁻³

# Error evaluation on two different tasks

- MNIST with fully-connected NN
- CIFAR 10 with Convolutional NN


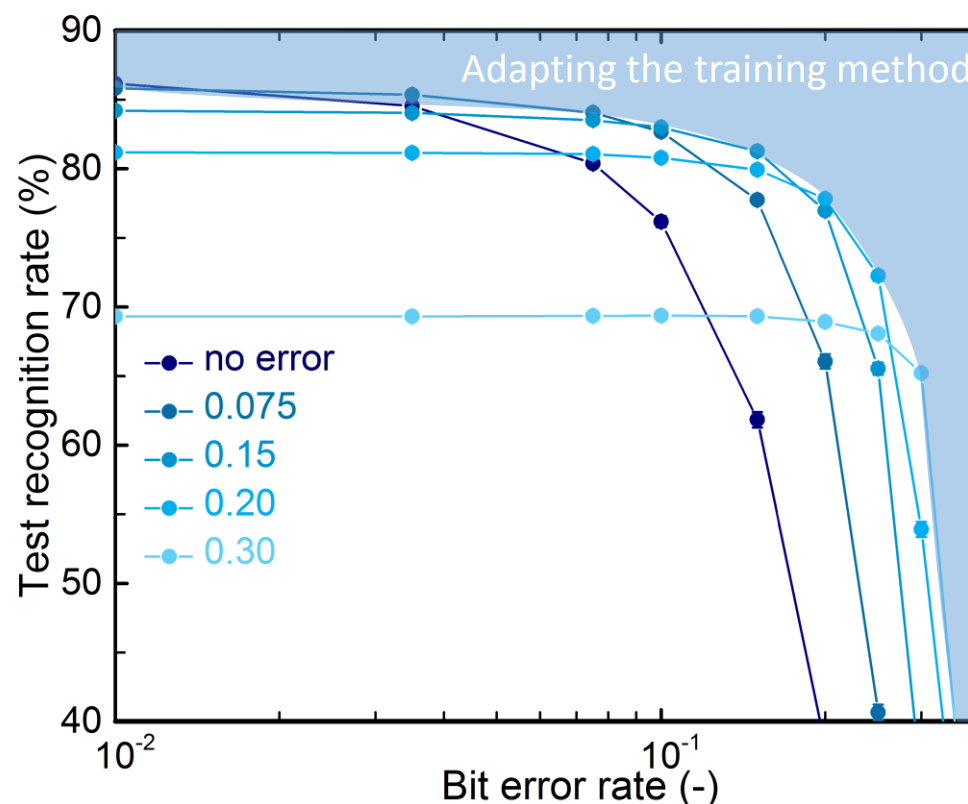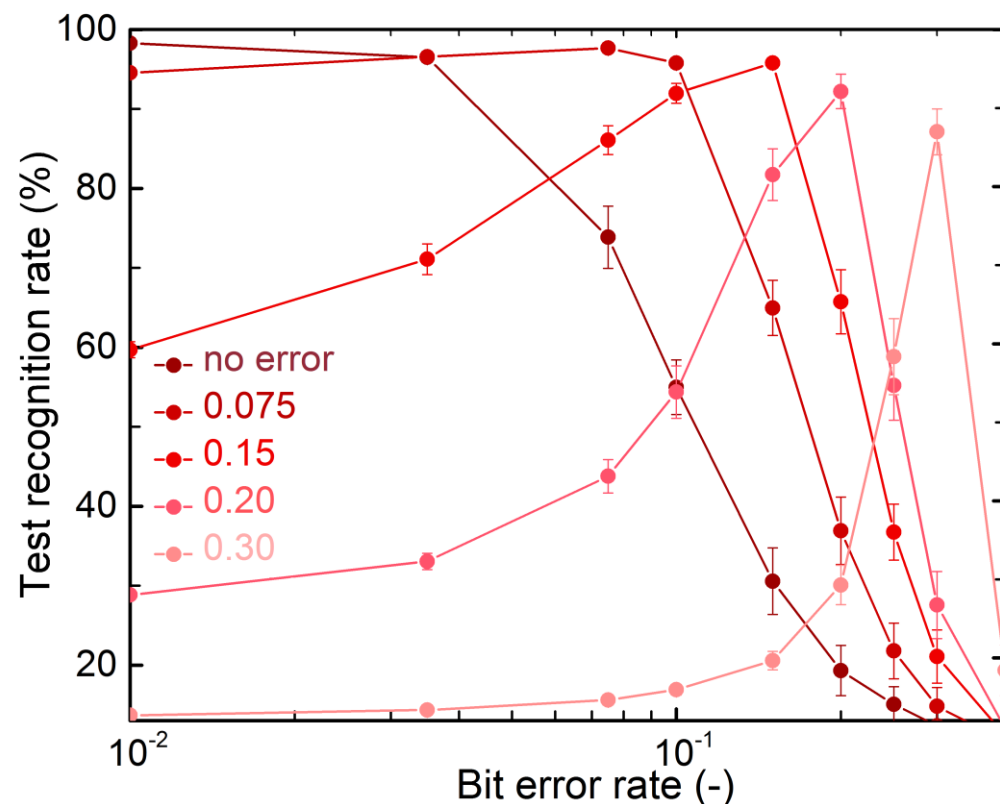
airplane   car   bird   cat   deer   dog   frog   horse   ship   truck



(a) MNIST

98.2 %

$2 \times 10^{-3}$

Accuracy on test set(%)

Bit error rate (-)

(b) CIFAR

87.25 %

$2 \times 10^{-3}$

Bit error rate (-)

**Binarized Neural Networks are resilient to errors**
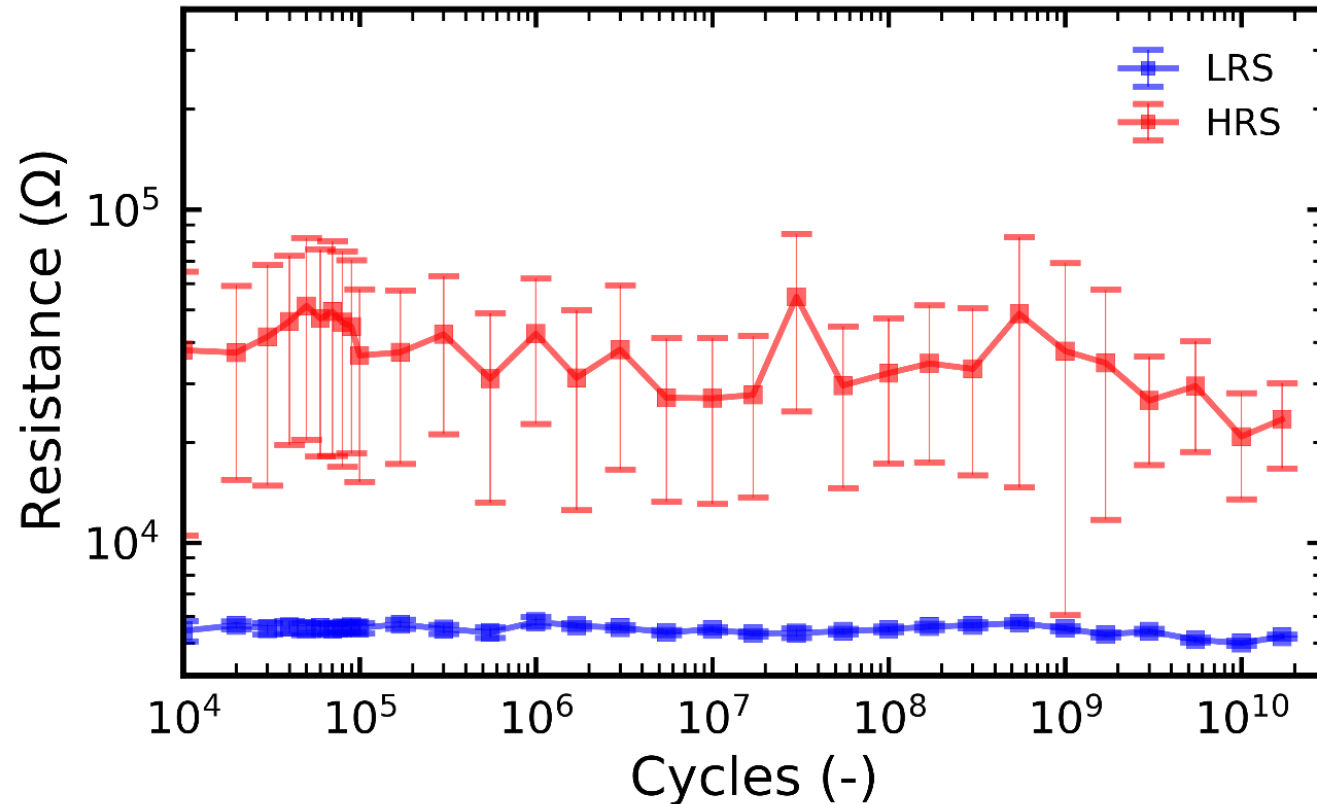
# Including bit errors during the training process

- MNIST with fully-connected NN
- CIFAR 10 with Convolutional NN



Adapting the training method can extend the bit error tolerance

# Reliance on 2T2R + Tolerance to Errors Allows Using Weak Programming Condition



- $V_{reset}$ = 1.5V
- Compliance current = 200 µA
- Error rate 1T1R = **~10$^{-2}$**
- Error rate 2T2R = **2x10$^{-3}$**

Endurance > **10$^{10}$** cycles

# To sum up

- Binarized Neural Network ideal for in memory computing

- 2T2R decreases bit error in comparison to 1T1R and ECC
  - Avoid the use of ECC
  - Logic operation integrated in reading circuit

- Binarized Neural Network are resilient to errors
  - Low voltage/current for reading/programming
  ↳ leads to : Low Energy & high endurance

- Technology ready today

**Hardware & learning algorithm co-development**

# Thank you for your attention !