# Memristor Overwrite Logic (MOL) for Energy-Efficient In-Memory DNN

Khaled Alhaj Ali , Mostafa Rizk, Amer Baghdadi, Jean-Philippe Diguet, Jalal Jomaah

**2020 IEEE International Symposium on Circuits & Systems**
**Virtual, October 10-21, 2020**

IMT Atlantique
Bretagne-Pays de la Loire
École Mines-Télécom

Lebanese University

Lab-STICC

# OUTLINE

ISCAS 2020

**IMT Atlantique**
Bretagne-Pays de la Loire
École Mines-Télécom

# OUTLINE

**IMT Atlantique**
Bretagne-Pays de la Loire
École Mines-Télécom

**Cloud computing**

Raw data

Decision

Neural networks

Storage
Computation
Analytics

**Cloud**

**Device**

☹ Does not ensure the required real time response

**Edge computing**

Neural networks

Storage
Computation
Analytics

Action overview

**Cloud**

**Edge/Fog**

☹ Slow processing speed and high energy consumption

**IMT Atlantique**
Bretagne-Pays de la Loire
École Mines-Télécom

Memristor Overwrite Logic (MOL) for Energy-Efficient In-Memory DNN
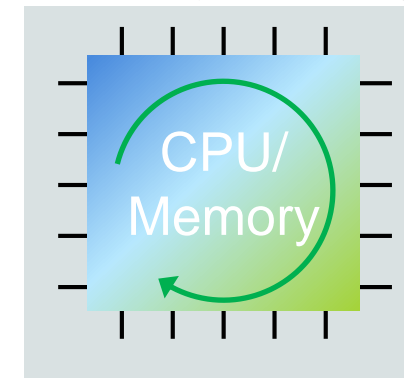
October 10-21, 2020

Lab-STICC

# From Von Neumann to in-memory computing

**Von Neumann model**



▶ Intensive data transfer between memory and CPU/GPU

▶ Large size of the processed data

**In-memory computing**



**Memristor-based logic design** enables true in-memory computing

IMT Atlantique
Bretagne-Pays de la Loire
École Mines-Télécom

Lab-STICC

# OUTLINE

**IMT Atlantique**
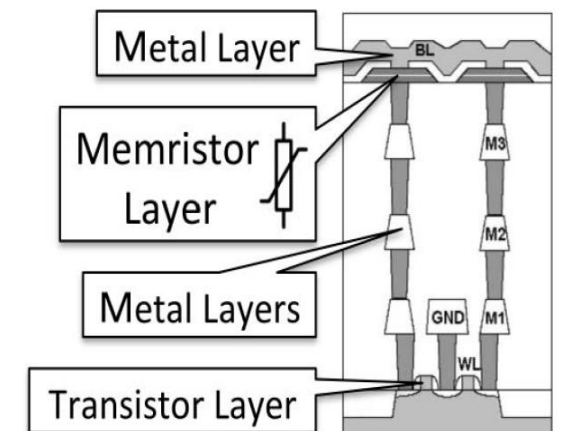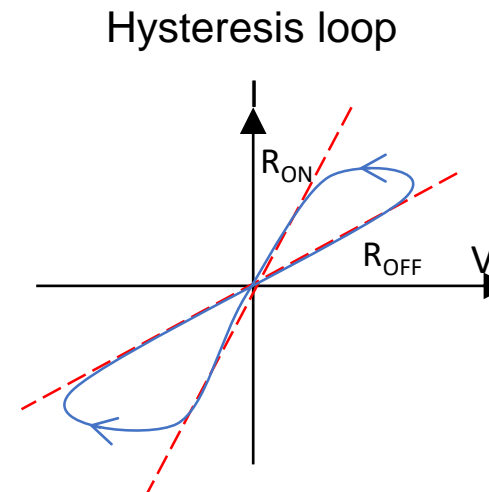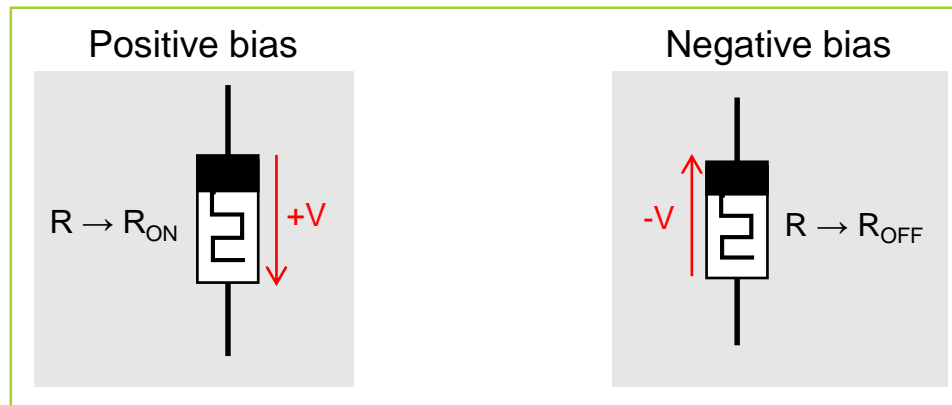Bretagne-Pays de la Loire
École Mines-Télécom

# Basics of memristor technology

► **Recent development of new non-volatile memory technologies (Memristor)**

– Theoretically predicted in 1971 by Chua

– Received attention in 2008 (HP labs)

– Triggered many efforts to explore their use in different applications

► **Basic operation**



Memorize last resistance state

Positive bias
$R \rightarrow R_{ON}$  +V

Negative bias
-V  $R \rightarrow R_{OFF}$

Hysteresis loop
$R_{ON}$
$R_{OFF}$  V

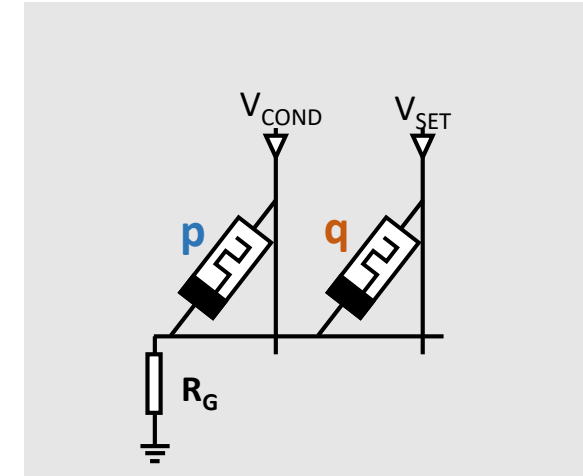Metal Layer
Memristor Layer
Metal Layers
Transistor Layer

# Limitations of existing logic design styles
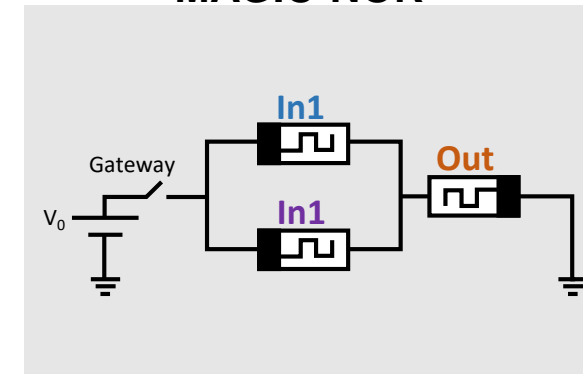
► IMPLY and MAGIC-NOR:

☹ Partial switching in IMPLY

☹ State drift in IMPLY and MAGIC

☹ Requirement of memristive devices with high $R_{OFF}/R_{ON}$ ratio

☹ High number of computational cycles

**IMPLY**



**MAGIC-NOR**



**IMT Atlantique**
Bretagne-Pays de la Loire
École Mines-Télécom

Memristor Overwrite Logic (MOL) for Energy-Efficient In-Memory DNN

October 10-21, 2020

Lab-STICC

# OUTLINE

**IMT Atlantique**
Bretagne-Pays de la Loire
École Mines-Télécom

# Digital representation of memristor

– Non sufficient magnitude or duration of the bias

– $V > Vth$ and $T > Tmin$

**Partial resistance switching**



$R_{ON} < R < R_{OFF}$

**Binary resistance switching**



$T > T_{min}$

$V > V_{th}$

$R \in \{R_{ON}, R_{OFF}\}$

☺ We are able to define the internal state of the memristor in digital domain

**IMT Atlantique**
Bretagne-Pays de la Loire
École Mines-Télécom

Memristor Overwrite Logic (MOL) for Energy-Efficient In-Memory DNN

October 10-21, 2020

Lab-STICC

# Digital representation of memristor

**Current & next state**

**Finite state machine (FSM)**



**State equation**

$$\Rightarrow Q_{n+1} = Q_n A + Q_n \bar{B} + A\bar{B}$$

▶ Derived logic cases

$$Q_{n+1} = Q_n A + Q_n \bar{B} + A\bar{B}$$

➡

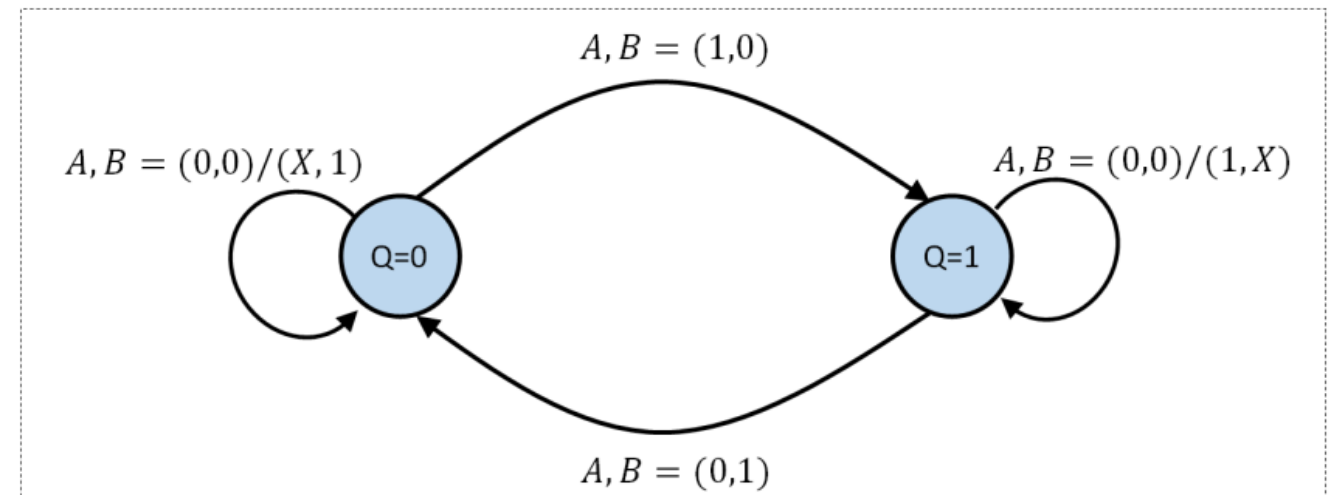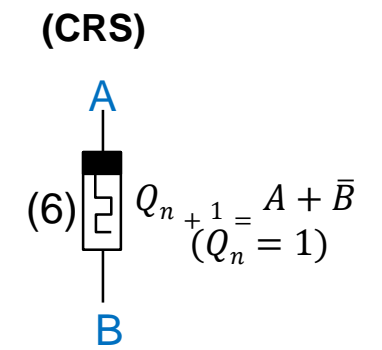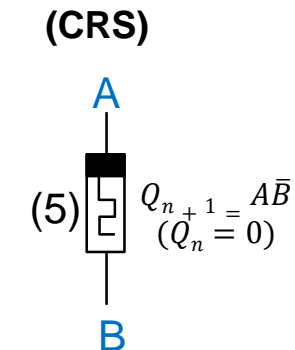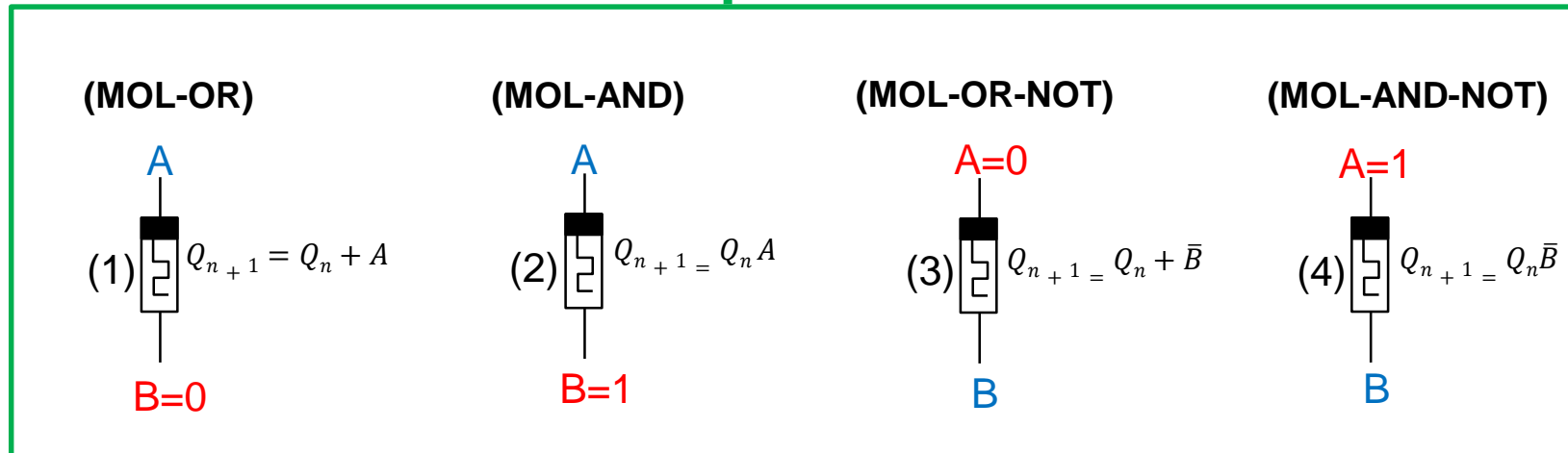$$Q_{n+1} = \begin{cases} Q_n + A, & B = 0, & case : 1 \\ Q_n A, & B = 1, & case : 2 \\ Q_n + \bar{B}, & A = 0, & case : 3 \\ Q_n \bar{B}, & A = 1, & case : 4 \\ A\bar{B}, & Q_n = 0, & case : 5 \\ A + \bar{B}, & Q_n = 1, & case : 6 \end{cases}$$

**MOL operations**

**(MOL-OR)**

A

(1) $Q_{n+1} = Q_n + A$

B=0

**(MOL-AND)**

A

(2) $Q_{n+1} = Q_n A$

B=1

**(MOL-OR-NOT)**

A=0

(3) $Q_{n+1} = Q_n + \bar{B}$

B

**(MOL-AND-NOT)**

A=1

(4) $Q_{n+1} = Q_n \bar{B}$

B

**(CRS)**

A

(5) $Q_{n+1} = A\bar{B}$ $(Q_n = 0)$

B

**(CRS)**

A

(6) $Q_{n+1} = A + \bar{B}$ $(Q_n = 1)$

B

IMT Atlantique
Bretagne-Pays de la Loire
École Mines-Télécom

Memristor Overwrite Logic (MOL) for Energy-Efficient In-Memory DNN

October 10-21, 2020

Lab-STICC

# MOL logic design

▶ MOL on a vector of bits

**MOL-OR**

– Write

– Overwrite (while selecting by 0)

IMT Atlantique
Bretagne-Pays de la Loire
École Mines-Télécom

Memristor Overwrite Logic (MOL) for Energy-Efficient In-Memory DNN

October 10-21, 2020

Lab-STICC

# MOL logic design

► MOL on a vector of bits

**MOL-AND**

– Write

– Overwrite (while selecting by 1)

$I_k \in \{-1v, 1v\}$
$k \in [0, N-1]$

B=0

$A_k \in \{0v, 1v\}$
$k \in [0, N-1]$

B=1



IMT Atlantique
Bretagne-Pays de la Loire
École Mines-Télécom

Memristor Overwrite Logic (MOL) for Energy-Efficient In-Memory DNN

October 10-21, 2020

Lab-STICC

# MOL logic design

▶ **MOL in crossbar array**

– Write new data bits

– Overwrite the already stored data bits

# OUTLINE

ISCAS 2020

**IMT Atlantique**
Bretagne-Pays de la Loire
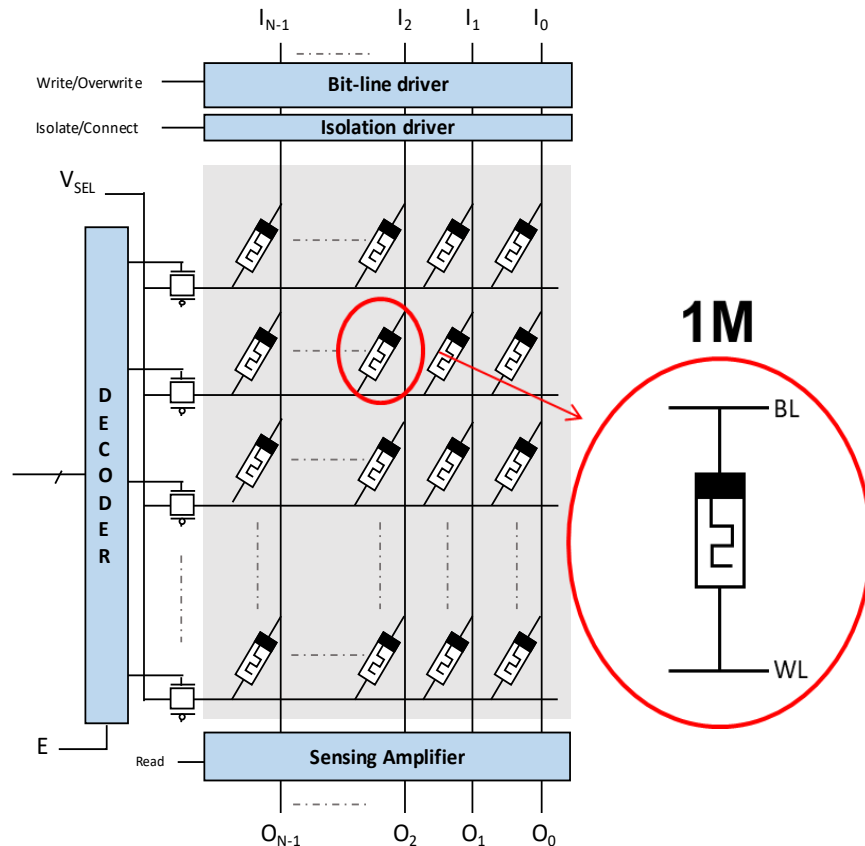École Mines-Télécom

Proposed computational memory
# MOL-based computational memory

► **1M** MOL architecture



► **1T1M** MOL architecture



► Four modes:

1. Write mode

2. Overwrite mode

3. Read mode

4. Idle mode

**IMT Atlantique**
Bretagne-Pays de la Loire
École Mines-Télécom

Memristor Overwrite Logic (MOL) for Energy-Efficient In-Memory DNN

October 10-21, 2020

Lab-STICC

# Proposed computational memory
# MOL-based computational memory

▶ Drivers

– Efficiently shared to all **memristive cells**

– Efficiently shared between **storage** and **computation**

☹ **Supports only logic accumulation** of newly arriving bits



Bit-line driver

Isolation driver

Sensing amplifiers

IMT Atlantique
Bretagne-Pays de la Loire
École Mines-Télécom

Lab-STICC

► Two coupled MOL-memories

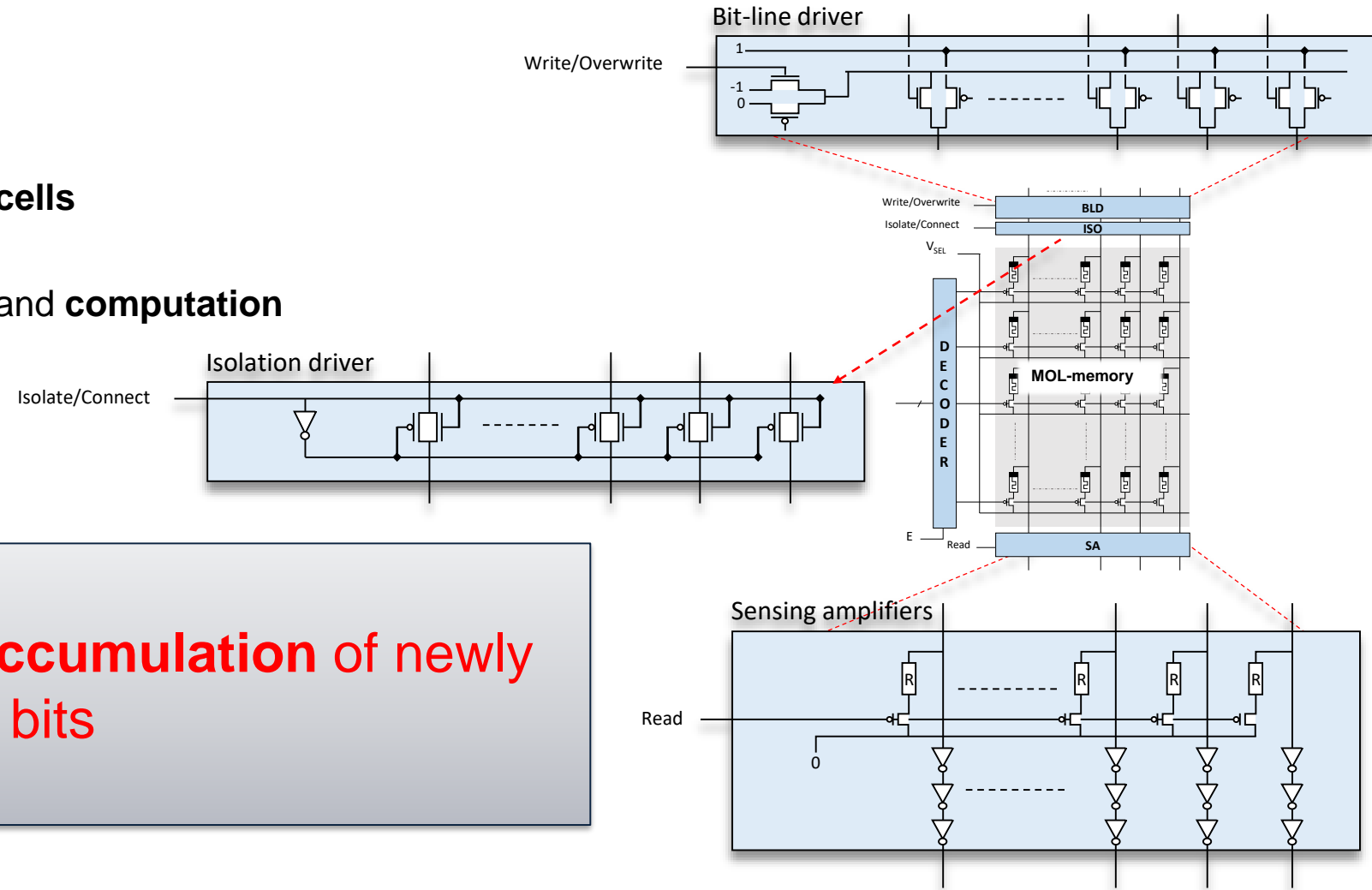– Two interconnected MOL-memory blocks

– Perform MOL **between any two wordlines**

Data in

Control

Address

**MOL-memory A**

Intermediate driver

**MOL-memory B**

Address

Data out

**IMT Atlantique**
Bretagne-Pays de la Loire
École Mines-Télécom

Memristor Overwrite Logic (MOL) for Energy-Efficient In-Memory DNN
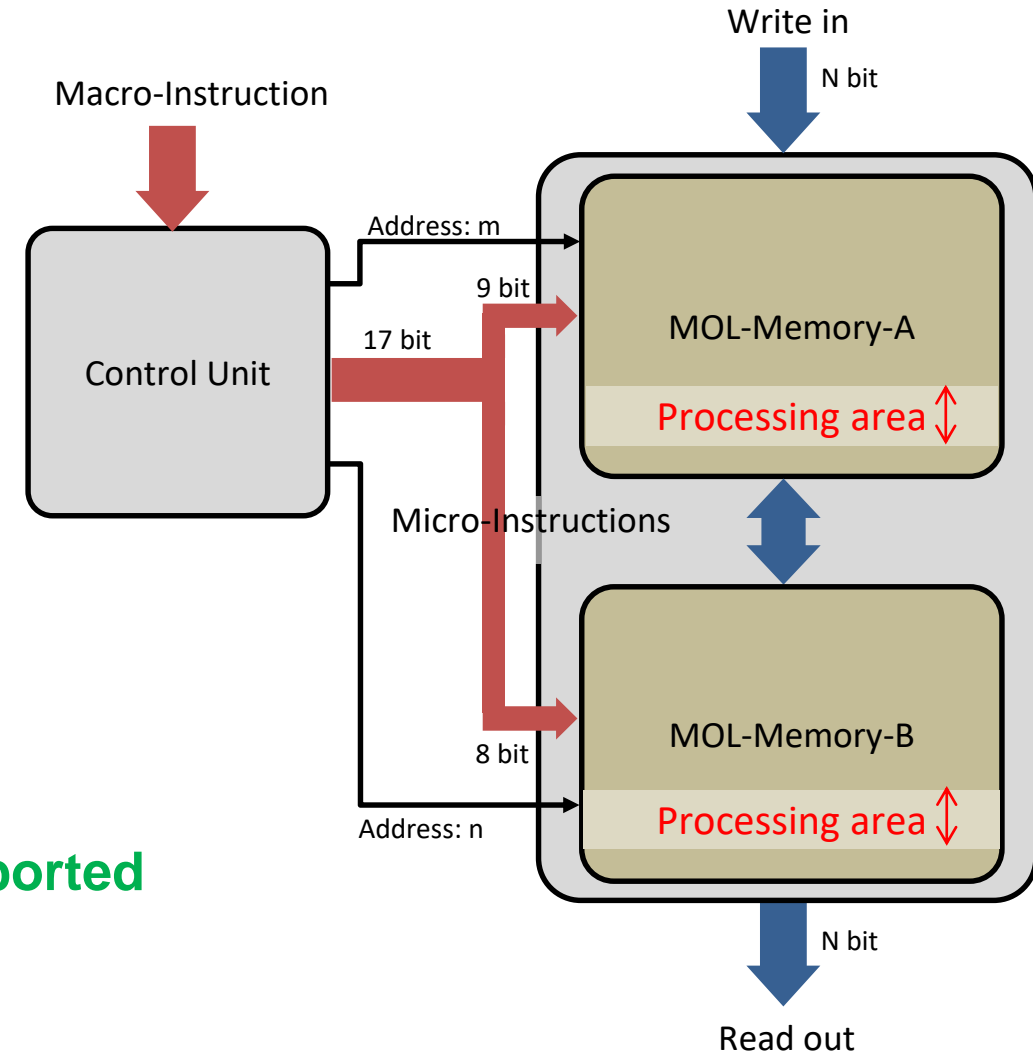
October 10-21, 2020

Lab-STICC

# MOL-based computational memory

▶ Performing arithmetic tasks

– At each time step, either **storage** or **computation**

– An arithmetic task (**Macro-Instruction**) is broken into several MOL operations (**Micro-instructions**)

– A **processing area** is reserved for computation

– **Dynamically changed** to maintain uniform endurance of cells

## More than 30 micro-instructions are supported



Write in
N bit

Macro-Instruction

Address: m

9 bit

17 bit

Control Unit

MOL-Memory-A

Processing area

Micro-Instructions

MOL-Memory-B

8 bit

Processing area

Address: n

N bit

Read out

IMT Atlantique
Bretagne-Pays de la Loire
École Mines-Télécom

Memristor Overwrite Logic (MOL) for Energy-Efficient In-Memory DNN

October 10-21, 2020

Lab-STICC

# OUTLINE

ISCAS 2020

**IMT Atlantique**
Bretagne-Pays de la Loire
École Mines-Télécom

# Convolution process

## Convolution process in neural networks
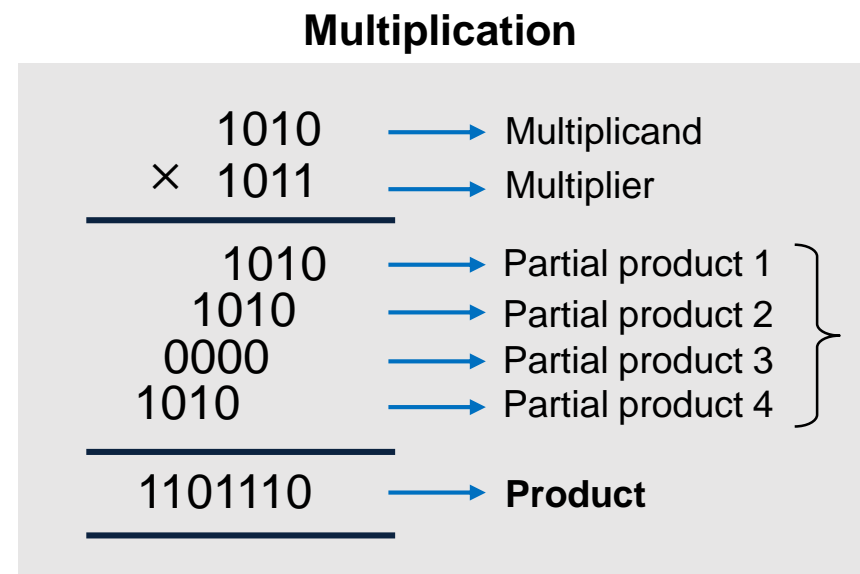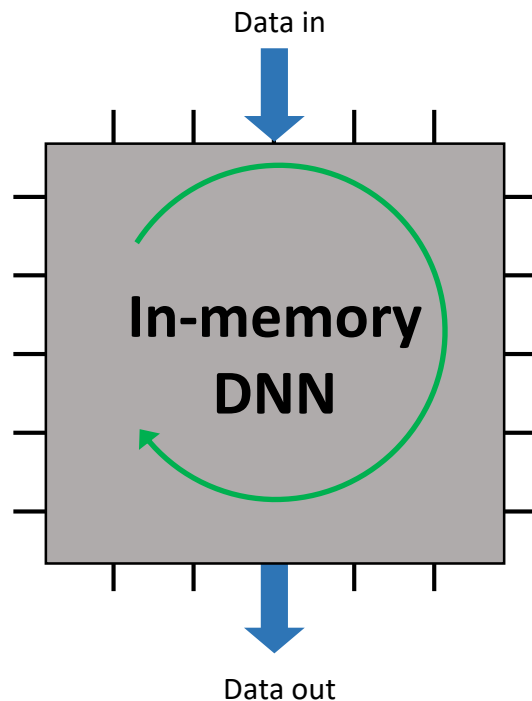


**Input Data**     **Neuron Weights**     **Output Equations**

$$[X_0 \ X_1 \ldots \ X_{K-1}] * \begin{bmatrix} A_0 & B_0 & C_0 \\ A_1 & B_1 & C_1 \\ \ldots & \ldots & \ldots \\ A_{K-1} & B_{K-1} & C_{K-1} \end{bmatrix} = \begin{bmatrix} Y_A = X_0 A_0 + X_1 A_1 + \cdots X_{K-1} A_{K-1} \\ Y_B = X_0 B_0 + X_1 B_1 + \cdots X_{K-1} B_{K-1} \\ Y_C = X_0 C_0 + X_1 C_1 + \cdots X_{K-1} C_{K-1} \end{bmatrix}$$

–   Number of MAC is proportional to the size of the network

–   **For large DNN**, implies intensive data movement between memory and processing cores

–   Inefficient in terms of **time** and **energy**
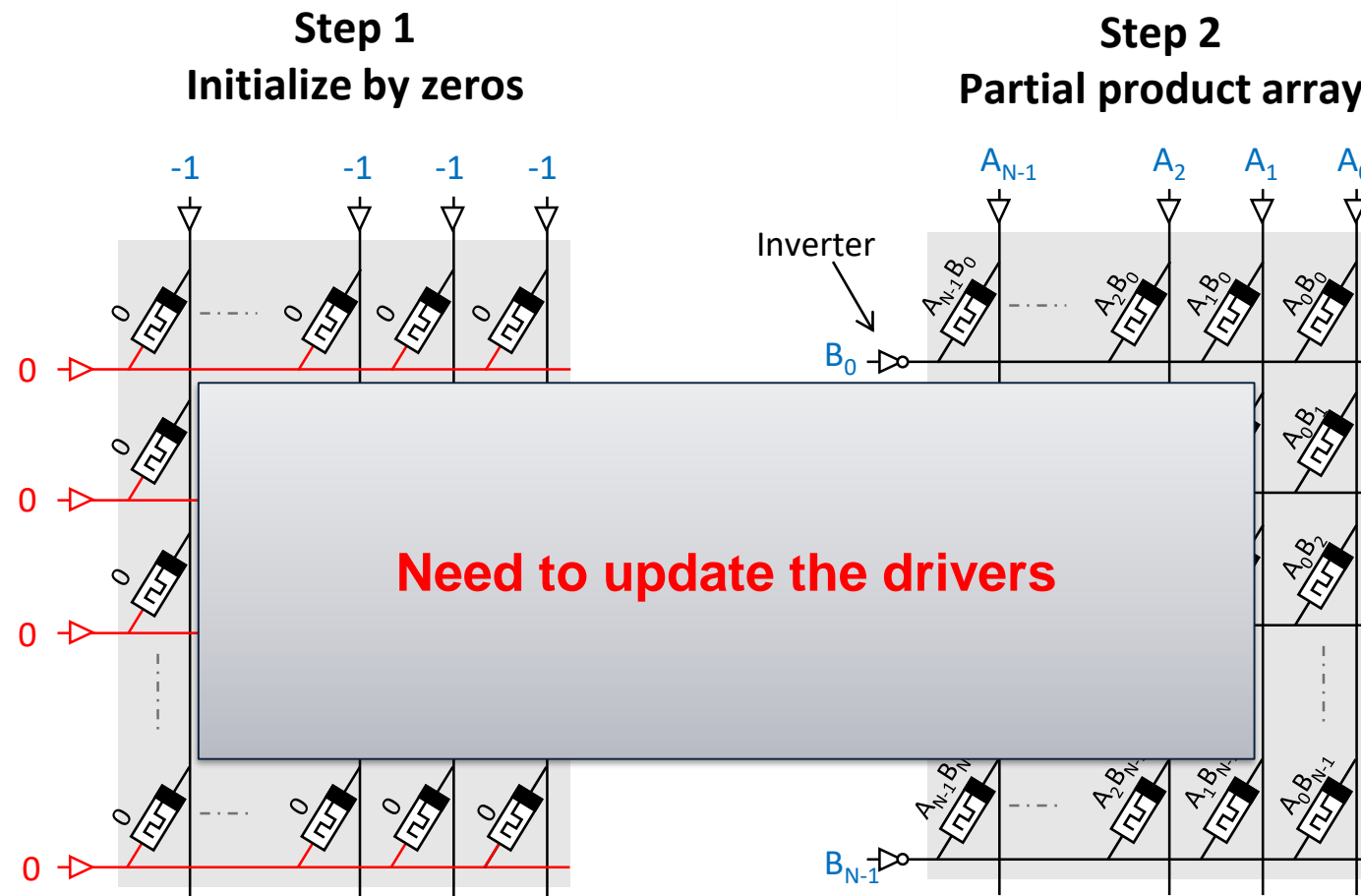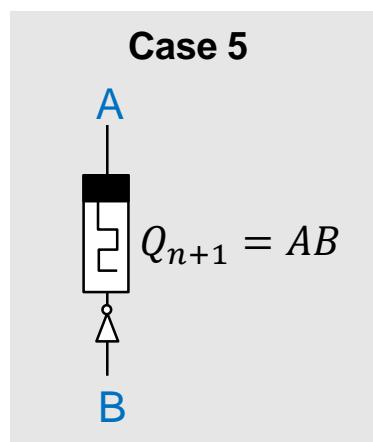
► Multiply and Accumulate (MAC) process

Data in

**In-memory DNN**

Data out

**Multiplication**

| | | |
|---|---|---|
| 1010 | → | Multiplicand |
| × 1011 | → | Multiplier |

| | | |
|---|---|---|
| 1010 | → | Partial product 1 |
| 1010 | → | Partial product 2 |
| 0000 | → | Partial product 3 |
| 1010 | → | Partial product 4 |

| | | |
|---|---|---|
| 1101110 | → | **Product** |

Significant number of

computational steps (bottleneck)

Need for an optimized **MAC** process inside memory

IMT Atlantique
Bretagne-Pays de la Loire
École Mines-Télécom

Memristor Overwrite Logic (MOL) for Energy-Efficient In-Memory DNN

October 10-21, 2020

Lab-STICC

# MOL – based in-memory DNN
# DNN in-memory

▶ **MOL-based In-memory DNN computing**

**Realization of Partial product**

– Initialize all cells to the logic zero

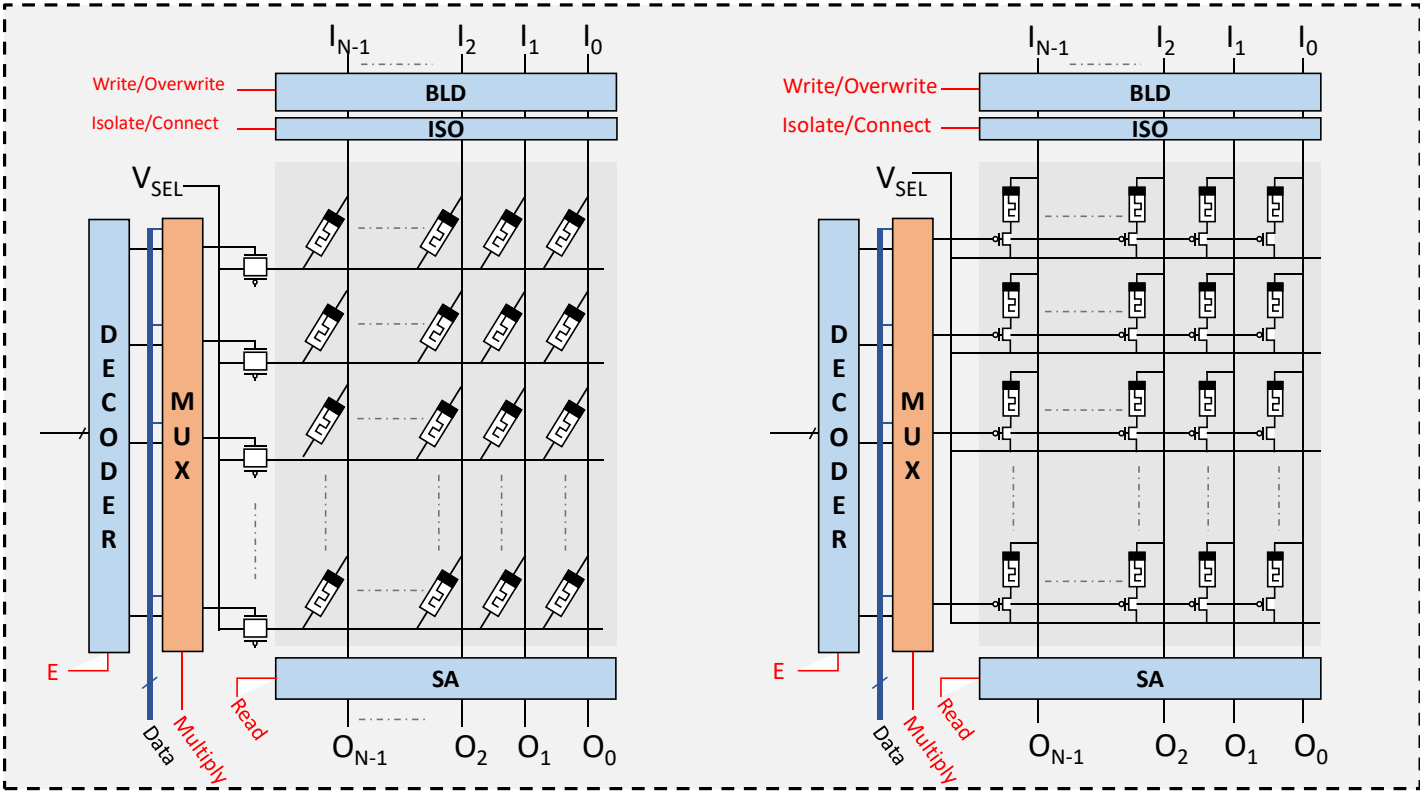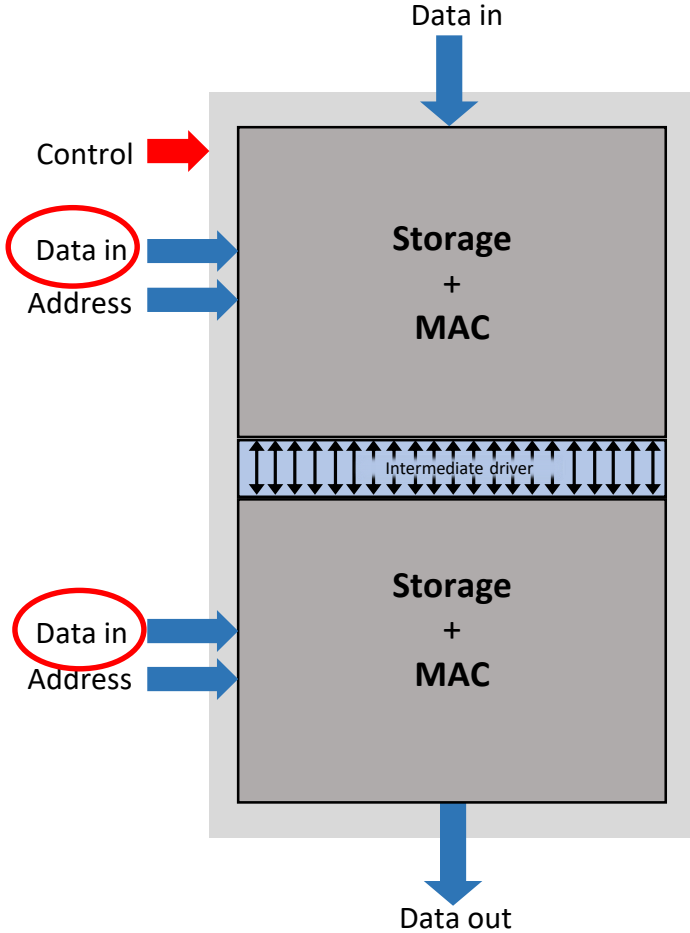– Vector A and an inverted vector B are fed to the columns and rows

**Case 5**

A

$Q_{n+1} = AB$

B

**Step 1**
**Initialize by zeros**

**Step 2**
**Partial product array**

Inverter

**Need to update the drivers**

IMT Atlantique
Bretagne-Pays de la Loire
École Mines-Télécom

Memristor Overwrite Logic (MOL) for Energy-Efficient In-Memory DNN

October 10-21, 2020

Lab-STICC

# MOL – based in-memory DNN
# **DNN in-memory**

▶ MOL-based In-memory DNN computing

MOL-memory with **Updated drivers**

**Computational memory (CMEM)**

# DNN in-memory
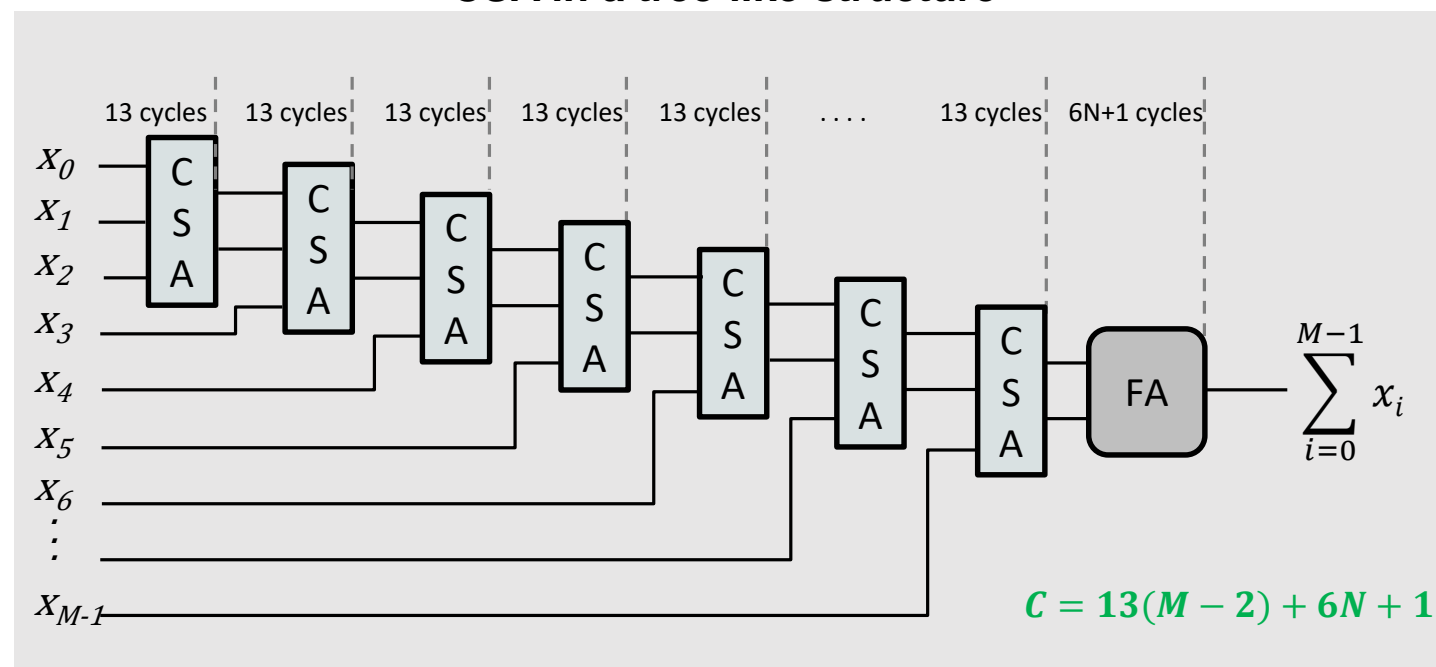
► MOL-based In-memory DNN computing

**Addition of Partial product**

– $M \times N$ multiplication requires $(M-1)(6N+1)$ cycles

**Use the method of carry save adder (CSA):**

– Provides a 3:2 operands reduction

– Fixed latency

– Number of cycles reduced to $C = 13(M-2) + 6N + 1$

**CSA in a tree-like structure**



$$C = 13(M-2) + 6N + 1$$

IMT Atlantique
Bretagne-Pays de la Loire
École Mines-Télécom

Memristor Overwrite Logic (MOL) for Energy-Efficient In-Memory DNN

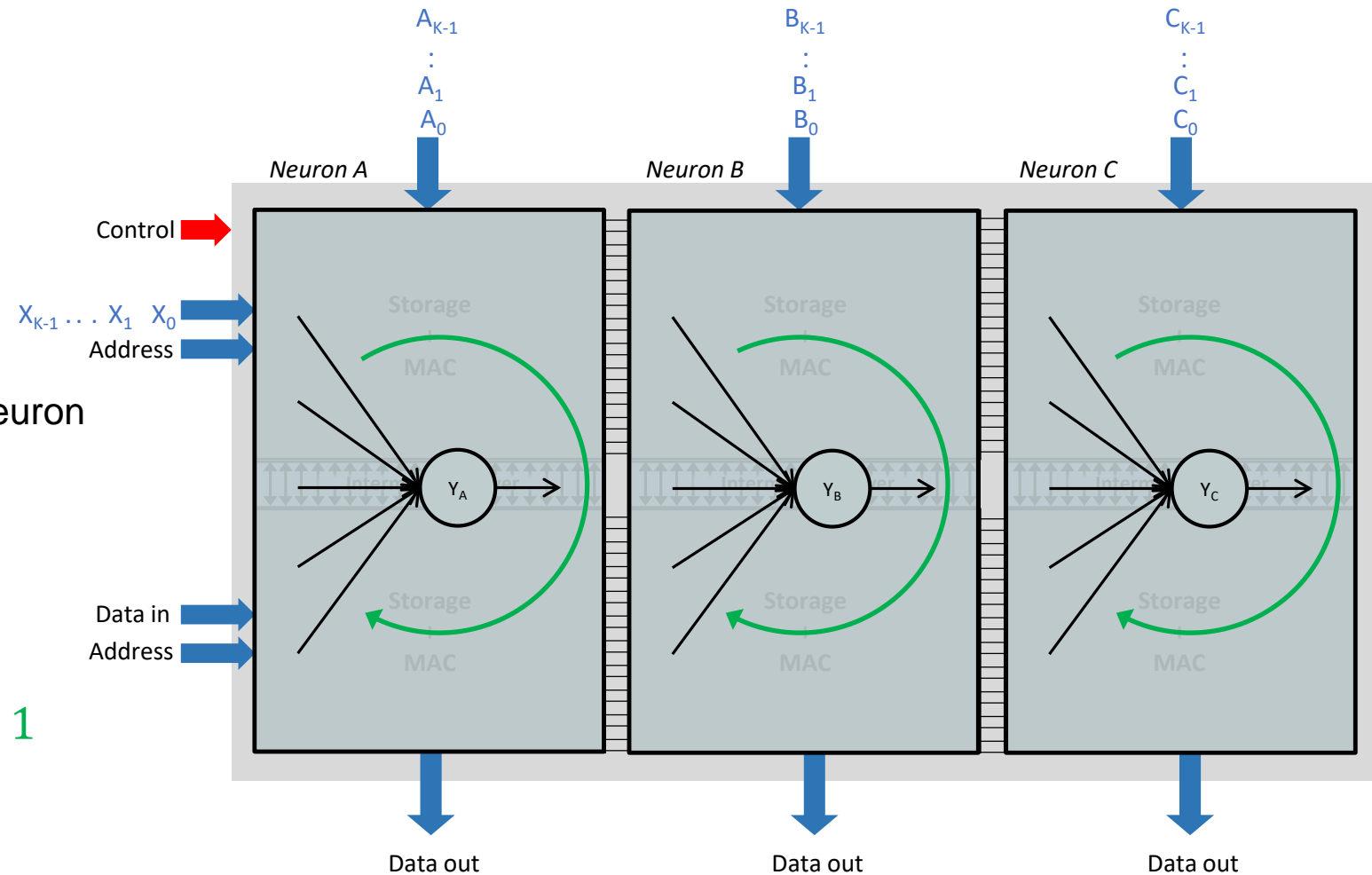October 10-21, 2020

Lab-STICC

# DNN in-memory

▶ MOL-based In-memory DNN computing

**CMEM-based DNN architecture**

– Interconnected CMEM blocks

– Each CMEM compute the output of a single neuron

– Parallel execution is possible
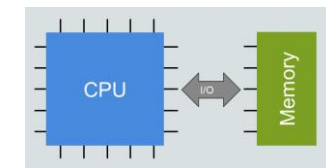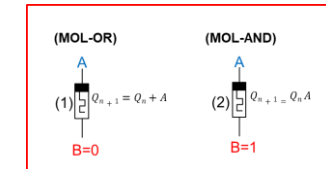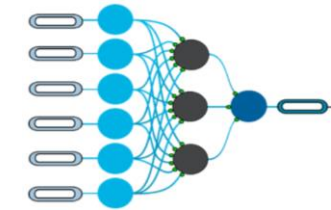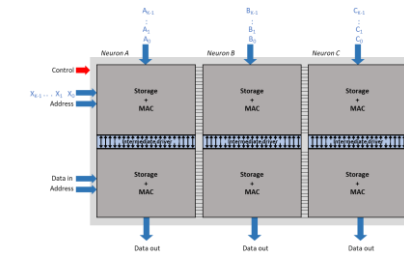
$$Total\ latency = (C + 15)K + 6(N + M) + 1$$

IMT Atlantique
Bretagne-Pays de la Loire
École Mines-Télécom

Lab-STICC

# OUTLINE

ISCAS 2020

**IMT Atlantique**
Bretagne-Pays de la Loire
École Mines-Télécom

# Conclusion

☑ Novel architecture design for **in-memory DNN** applications



☑ Programmable and allows to execute any sequence of arithmetic tasks including **MAC**



☑ Computation is performed based on our proposed **MOL design style**



☑ Addresses the inefficiency of moving data between memory and processing cores which is **time** and **energy** consuming.



**IMT Atlantique**
Bretagne-Pays de la Loire
École Mines-Télécom

Memristor Overwrite Logic (MOL) for Energy-Efficient In-Memory DNN

October 10-21, 2020

Lab-STICC

# Future works

► **Control unit** and the associated micro-operations synthesis tool

► Perform DNN on **real dataset** such as CIFAR-10 or MNIST

► Performance evaluation and **comparison** with state-of the art CPU and GPU implementations.

**IMT Atlantique**
Bretagne-Pays de la Loire
École Mines-Télécom

Lab-STICC

# Thank you for your attention

IMT Atlantique
Bretagne-Pays de la Loire
École Mines-Télécom

Lab-STICC