# In-Memory Computing based Machine Learning Accelerators: Opportunities and Challenges

KAUSHIK ROY

DIRECTOR, CENTER FOR BRAIN-INSPIRED COMPUTING

PURDUE UNIVERSITY

KAUSHIK@PURDUE.EDU

SRC
JUMP

C-BRIC

# Machine (Deep) Learning Saga

| 1997 | 2011 | 2016 | 2018 |
|------|------|------|------|



Deep Blue vs. Kasparov
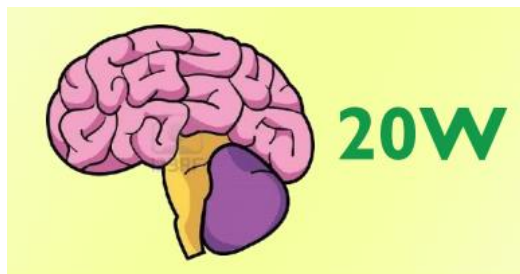**~15000W**

Watson Wins Jeopardy
**~200000W**

AlphaGo vs. Sedol
**~300000W**

Autonomous Driving
**~$$$$$**

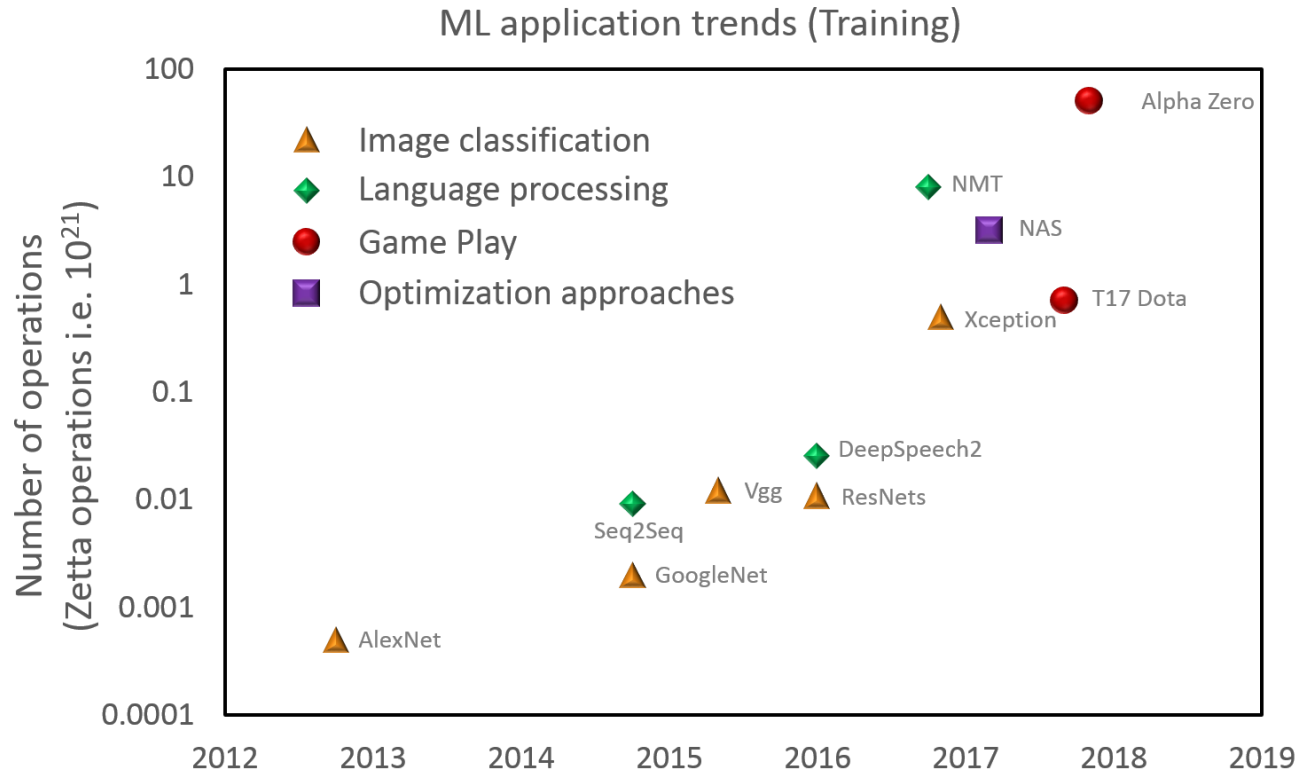- **Advent of Deep Learning, 2012**
- **Fueled by powerful hardware - GPUs**
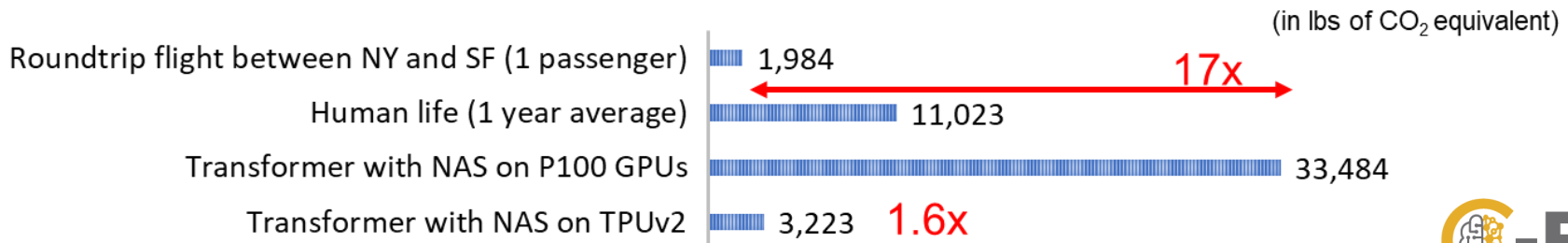
**20W**

Algorithm performance moving closer

Hardware cost moving farther

**"Aspirations have grown faster than the technology available to satisfy them"**

C-BRIC

# AI Compute Demands (Training)

## ML application trends (Training)



Legend:
- ▲ Image classification
- ◆ Language processing
- ● Game Play
- ■ Optimization approaches

Data points: Alpha Zero, NMT, NAS, T17 Dota, Xception, DeepSpeech2, Vgg, ResNets, Seq2Seq, GoogleNet, AlexNet

Y-axis: Number of operations (Zetta operations i.e. $10^{21}$)
X-axis: 2012–2019

## COMMON CARBON FOOTPRINT BENCHMARKS

(in lbs of $CO_2$ equivalent)

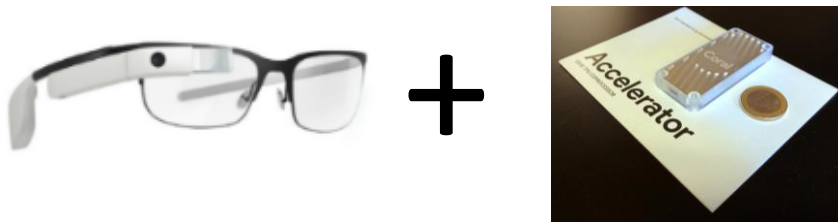| | |
|---|---|
| Roundtrip flight between NY and SF (1 passenger) | 1,984 |
| Human life (1 year average) | 11,023 |
| Transformer with NAS on P100 GPUs | 33,484 |
| Transformer with NAS on TPUv2 | 3,223 |

17x

1.6x

Per Network!

C-BRIC

# Edge Intelligence: Efficiency Gap

➢ Case study: Object recognition in a smart glass with a state-of-the-art accelerator



**+**



Google Edge TPU



Retinanet DNN* on a smart glass

| Performance | |
|---|---|
| Frames/sec | 13.3 |

| Battery Life | |
|---|---|
| Energy/op | 0.5 pJ/op |
| Energy/frame | 0.15 J/frame |
| **Time-to-die (2.1WH)** | **64 mins** |

*300 GOPs/inference

## Where do the in-efficiencies come from?

Algorithms        Hardware Architecture        Circuits and Devices

**Ref**: Venkataramani, S., Roy, K. and Raghunathan, A. "Efficient embedded learning for IoT devices." In *2016 21st Asia and South Pacific Design Automation Conference (ASP-DAC)*, pp. 308-311. IEEE.

C-BRIC

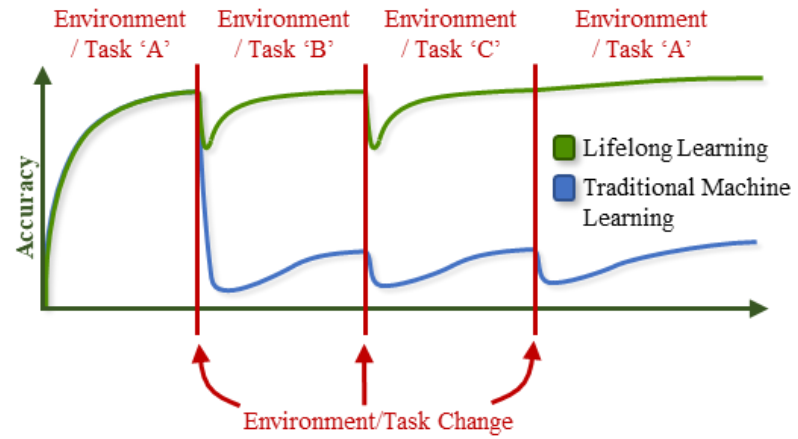# Beyond Compute Efficiency....

➢ Learning with less data

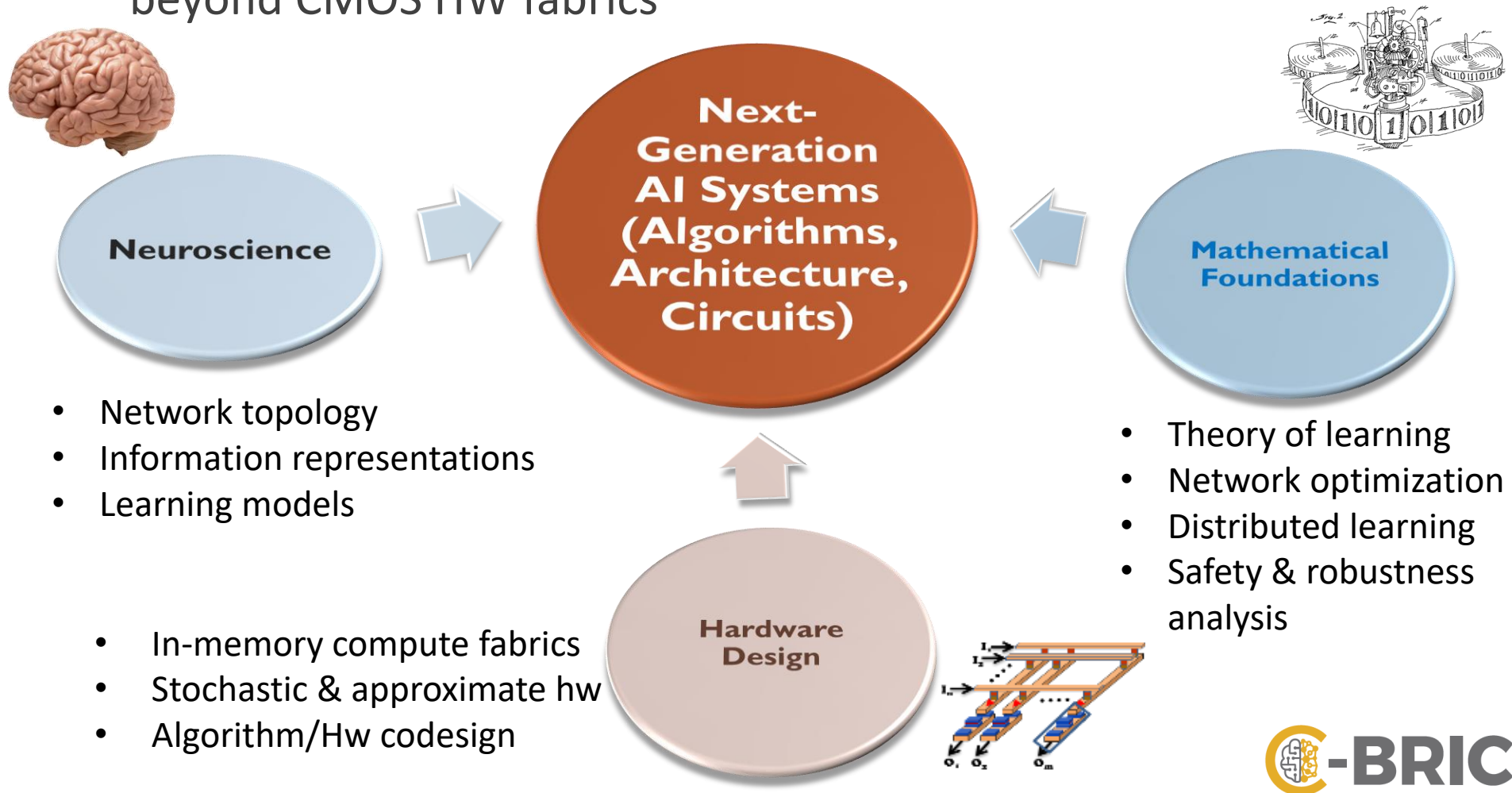➢ Generalization & Robustness/ Security

➢ Lifelong learning



$X_{clean}$

97.3% confidence
Macaw

$+$ 0.005 $\times$

$\Delta$

Adversarial
Perturbation

$=$

$X_{adversary}$

88.9% confidence
Bookcase



Accuracy

Deep Learning

Tradtional
Algorithms

Amount of data



Environment / Task 'A'  Environment / Task 'B'  Environment / Task 'C'  Environment / Task 'A'

Accuracy

■ Lifelong Learning
■ Traditional Machine Learning

Environment/Task Change

C-BRIC

# Center for Brain-Inspired Computing (C-BRIC): Approach

➢ Design next-generation AI systems by drawing from neuroscience, mathematical foundations and using CMOS and beyond CMOS HW fabrics



**Next-Generation AI Systems (Algorithms, Architecture, Circuits)**

**Neuroscience**

- Network topology
- Information representations
- Learning models

**Mathematical Foundations**

- Theory of learning
- Network optimization
- Distributed learning
- Safety & robustness analysis

**Hardware Design**

- In-memory compute fabrics
- Stochastic & approximate hw
- Algorithm/Hw codesign

C-BRIC

# AI Hardware Architecture: Circuits & Devices

➢ Circuits and architectures that can efficiently implement the algorithms (possibly embody computing principles from the brain)
- Near-/In-Memory Computing
- Approximate and stochastic hardware
- Neuromorphic devices and interconnects



Multicores/GPUs

Accelerators

Approximate & Stochastic Hardware

In-memory computing

Neuromorphic Devices

~$10^4$ Energy Gap

C-BRIC

# Background – In-Memory Computing

➢ **Definition:** Design approach that performs computation close to memory to overcome memory bottlenecks – bandwidth, energy

➢ Effective for simple arithmetic - bit-wise operations; fixed-point add, multiply, Truth-tables (ROMs/RAMs)

➢ Typical systems have much higher compute throughput than memory bandwidth(s)

➢ Lots of chip area are memory components (>=50% in TPU)
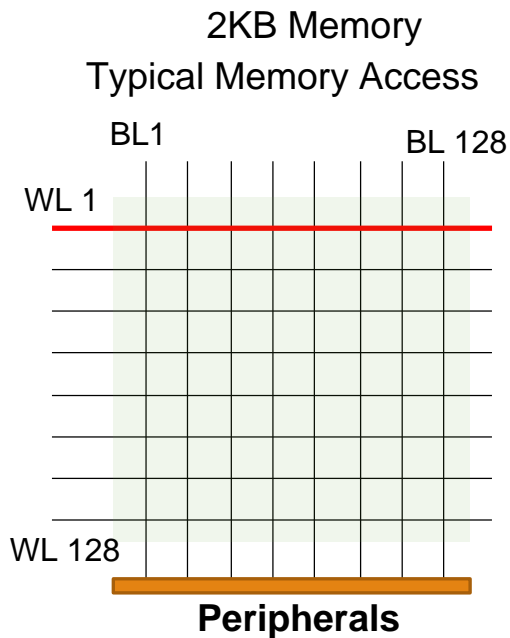  ○ Caches (L1, L2, …), Register File, Scratchpad, Buffers



**Computer Architecture: A Quantitative Approach**



**TPU Floorplan, ISCA 2017**

# In-Memory Computing for ML

## 2KB Memory
### Typical Memory Access

BL1     BL 128

WL 1

WL 128

**Peripherals**

Bits Accessed = 128
**Operation = Read/Write**

## 2KB Memory
### In-Memory Computing Digital

V1
V2

$V_{out}$

Bits Accessed = 2*128
Operation = Bit-wise AND
**$V_{out}$ = V1 AND V2**

## 2KB Memory
### In-Memory Computing Analog

$V_{in}^T$: 1×4

M: 4×64

$V_{out}$: 1×64

Bits Accessed = 4*128
Operation = Matrix Vector Multiply
**$V_{out}$ = M * $V_{in}$**

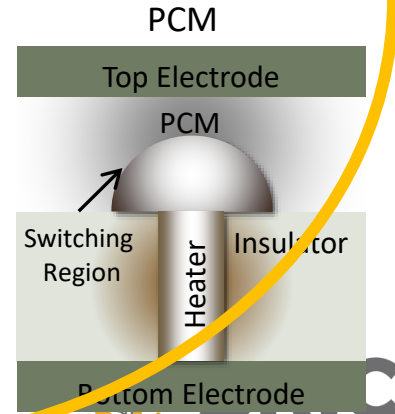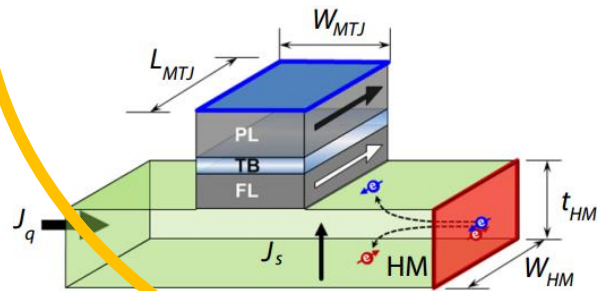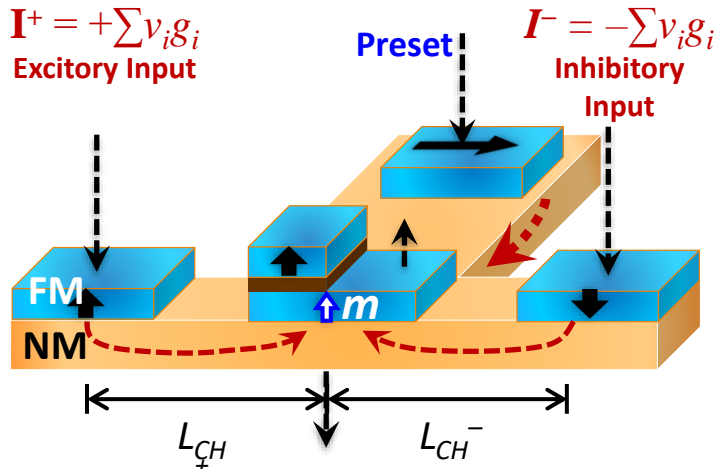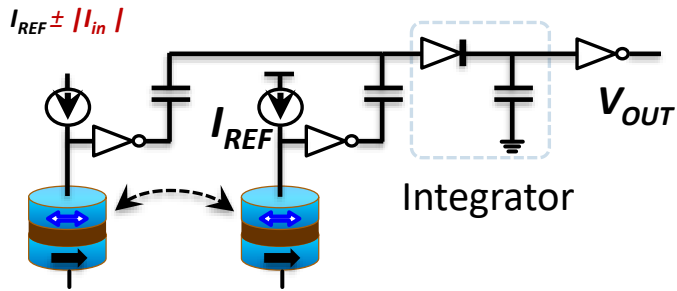| | |
|---|---|
| **Non Volatile Memory** | Jain et al. TVLSI'17, Abbrogio et al. Nature'18, Cai et al. Nature Elec.'19, Xue et al. ISSCC'20, Liu et al., ISSCC'20 |
| **SRAM** | Biswas et al. ISSCC'18, Valavi et al. JSSC'19, Si et al. ISSCC'19, Jaiswal et al. TVLSI'20, Dong et al. ISSCC'20 (TSMC – 7nm) |
| **ROM/RAM** | Lee et. al. EDL 2013, Lee et. al. TVLSI 2013, |

C-BRIC

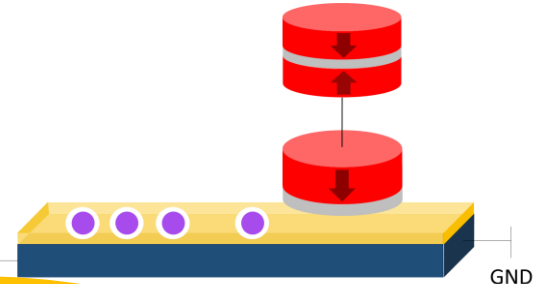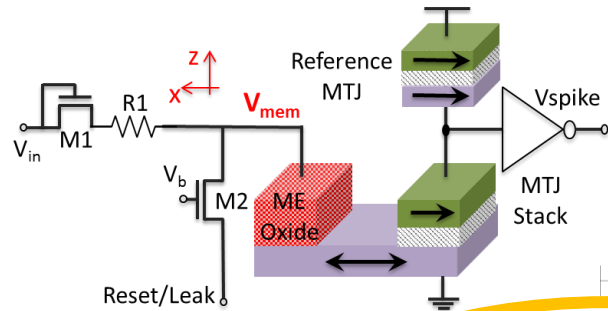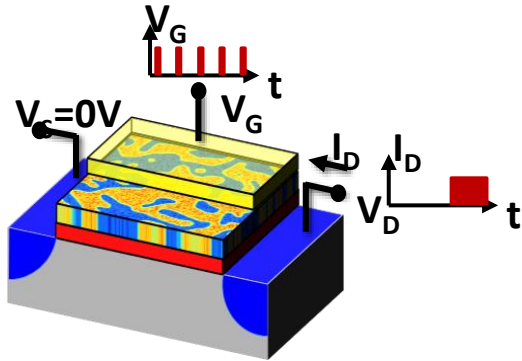# Machine Learning (Deep Learning)

➢ Deep Learning needs – lots of matrix multiplications



*Bruce Fleischer et al, IBM Research, 2018*

➢ Challenge: sustaining deep learning's insatiable compute demands
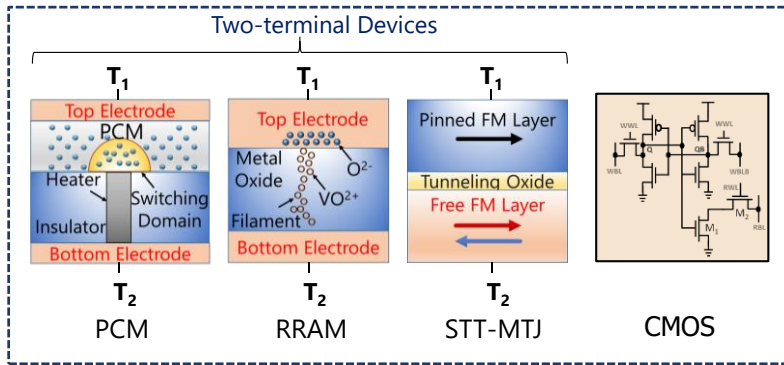
# Technology: Non-Volatile Memories
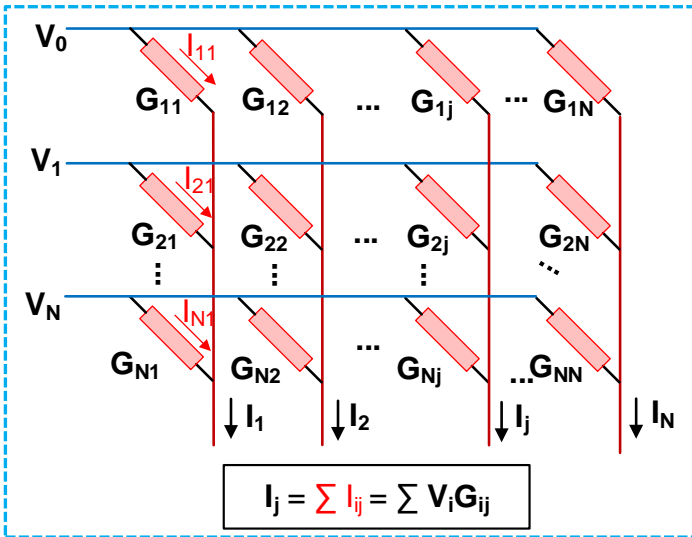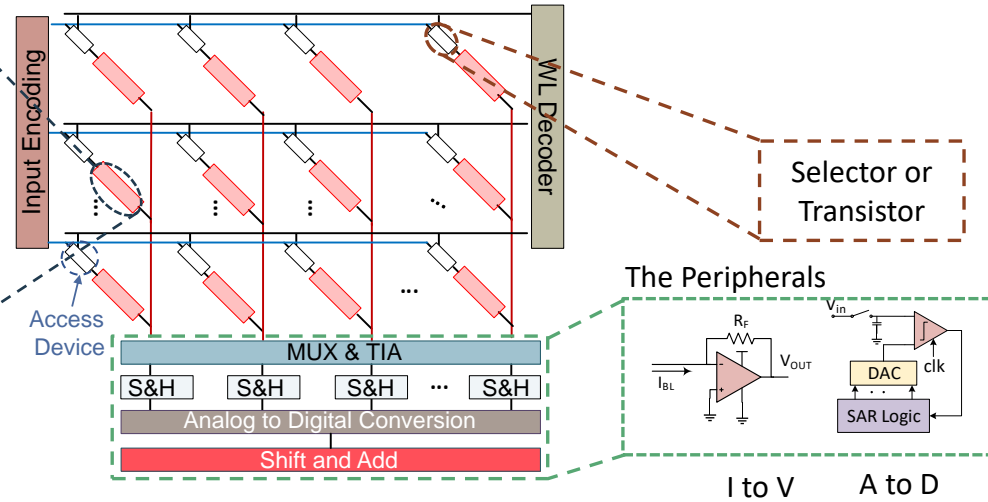
# CMOS SRAM and Non-volatile Memories

| Property | PCM | RRAM | MTJ | CMOS (SRAM) |
|---|---|---|---|---|
| Multi-level cell | Yes | Yes | No | No |
| Storage Density | High | High | High | Low |
| $R_{ON}/R_{OFF}$ | High | High | Low | High |
| Non-volatility | Yes | Yes | Yes | No |
| Leakage | Low | Low | Low | High |
| Cell Area | $16F^2$ | $16F^2$ | $30\text{-}80F^2$ | $160F^2$ (6T), $231F^2$ (8T) |
| Write Energy | 6 nJ | 2 nJ | < 1 nJ | < 0.1 nJ |
| Write Latency | 150 ns | 100 ns | 10 ns | < 1ns |
| Endurance | $10^7$ cycles | $10^5$ cycles | $10^{15}$ cycles | > $10^{16}$ cycles |

C-BRIC

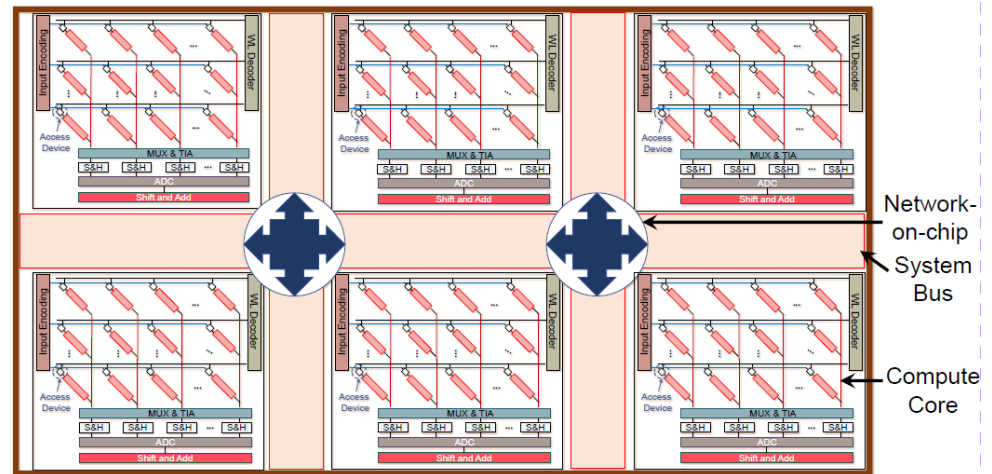# Fundamental building blocks of in-memory computing

In-Memory Computing Memory Devices

Chakraborty et. al. Resistive Crossbars as Approximate Hardware Building Blocks for Machine Learning: Opportunities and Challenges, Proc. of IEEE, 2020

**Two-terminal Devices**

PCM | RRAM | STT-MTJ | CMOS

Selector or Transistor

The Peripherals

MUX & TIA

Analog to Digital Conversion

Shift and Add

I to V          A to D

Access Device

$$I_j = \sum I_{ij} = \sum V_i G_{ij}$$

Network-on-chip

System Bus

Compute Core

**Efficient MVM**

**Spatially Distributed Cores**

C-BRIC

# Efficient Hardware Architecture: CiM

## In-Memory Computing Memory Devices

Chakraborty et. al. Resistive Crossbars as Approximate Hardware Building Blocks for Machine Learning: Opportunities and Challenges, Proc. of IEEE, 2020
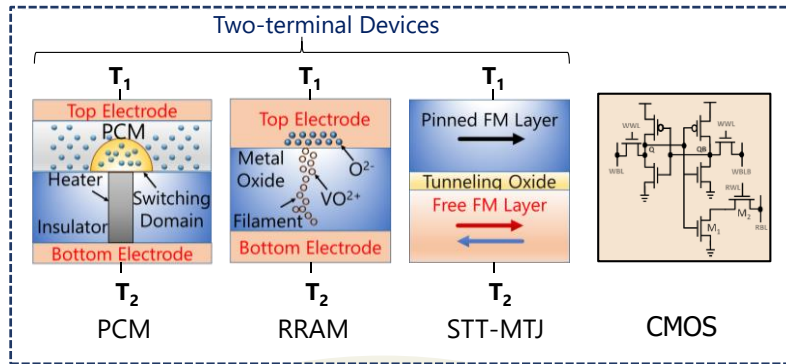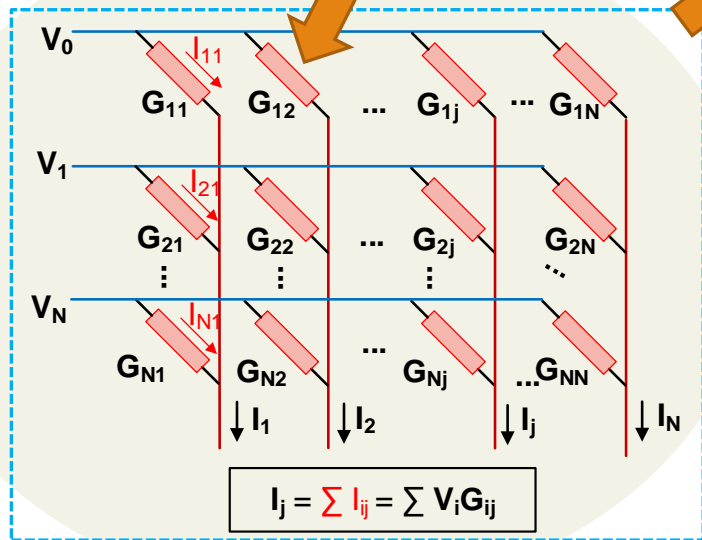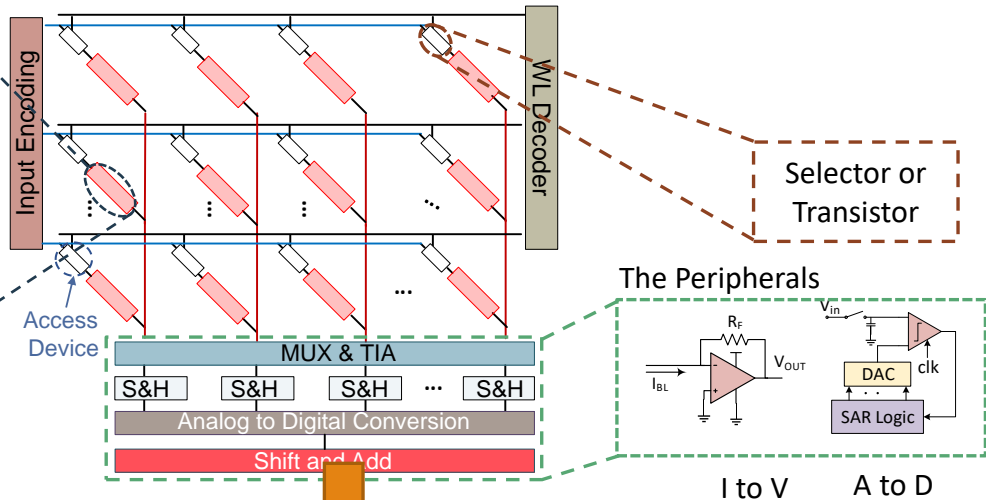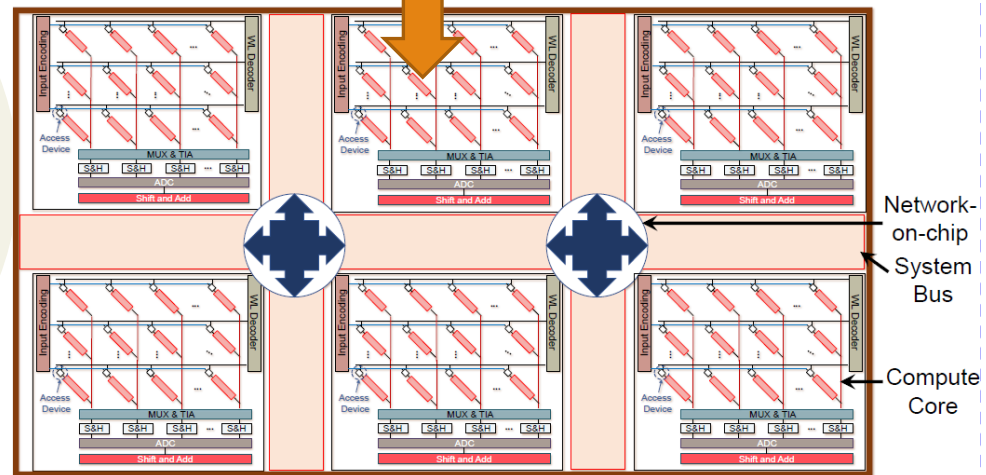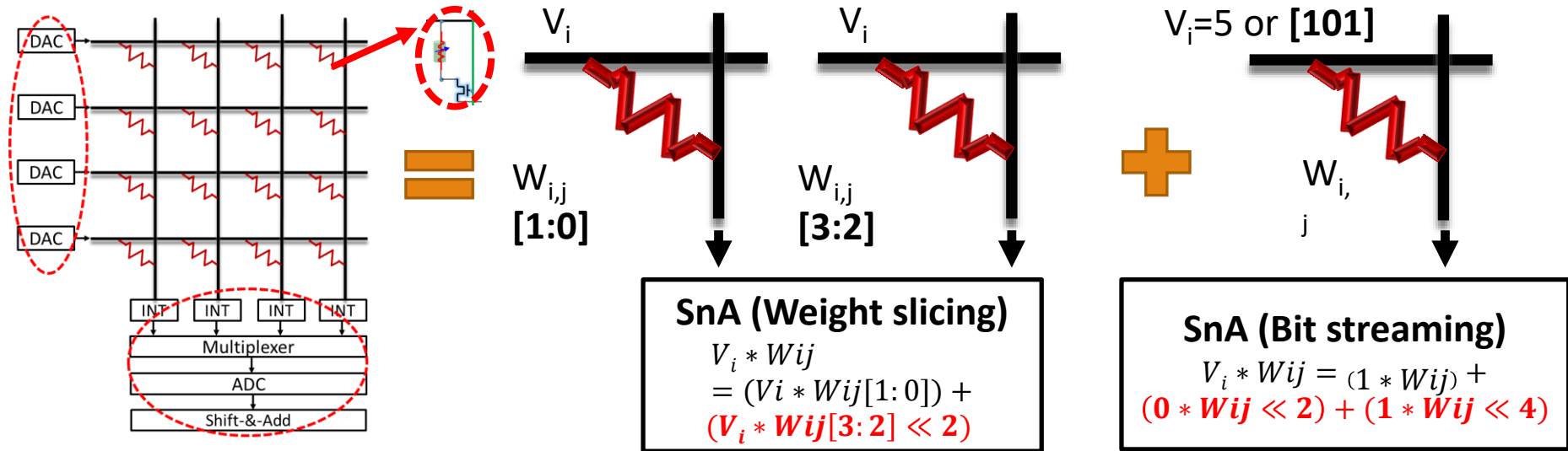


**Two-terminal Devices**

PCM | RRAM | STT-MTJ | CMOS

Selector or Transistor

Access Device

The Peripherals

MUX & TIA
S&H  S&H  S&H  ···  S&H
Analog to Digital Conversion
Shift and Add

I to V      A to D

$$I_j = \sum I_{ij} = \sum V_i G_{ij}$$

**Efficient MVM**

Network-on-chip
System Bus
Compute Core

**Spatially Distributed Cores**

C-BRIC

# Bit-slicing (weights and inputs)



**SnA (Weight slicing)**

$$V_i * Wij = (Vi * Wij[1:0]) + (V_i * Wij[3:2] \ll 2)$$

**SnA (Bit streaming)**

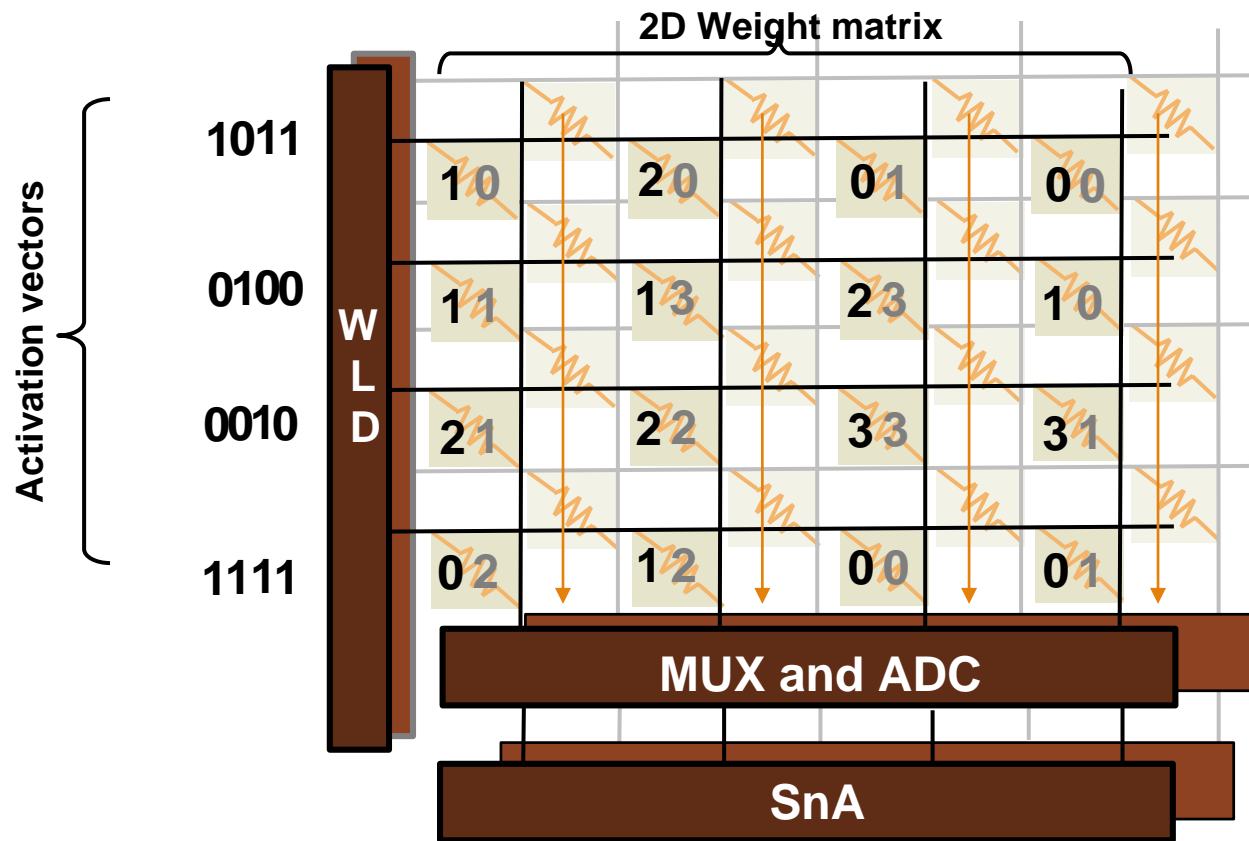$$V_i * Wij = (1 * Wij) + (0 * Wij \ll 2) + (1 * Wij \ll 4)$$

➢ Bit-slicing: weight slicing and input streaming enable using **low precision crossbars** and **low precision DACs** to compose **high precision MVMU**

# Analog CiM : Implementation details

➤ Parallel and efficient GEMV implementation:



**In Memory MVM**

# CiM processing details(1)

➢ Workload Mapping



**Weights**

**Input**

**Output**

**Tiling**

Tile    Tile

Tile    Tile

Tile    Tile

Tile    Tile

Tile    Tile

Divide Weight matrix into tiles based on memory array rows

C-BRIC

# CiM processing details(2)

➢ Workload Mapping



Tiling

Weight bit slice

Intra Tile Accumulation of partial sums using Shift and Add block.

Inter tile Accumulation using NoC

C-BRIC

# Architecture: *Spatial scalability*



**Massively parallel accelerator –> Amenable to Data-Level Parallelism -> Highly efficient ML inference**

Ankit, Roy, et. Al., "PUMA: A Programmable Ultra-efficient Memristor-based Accelerator for Machine Learning Inference", ASPLOS 2019.

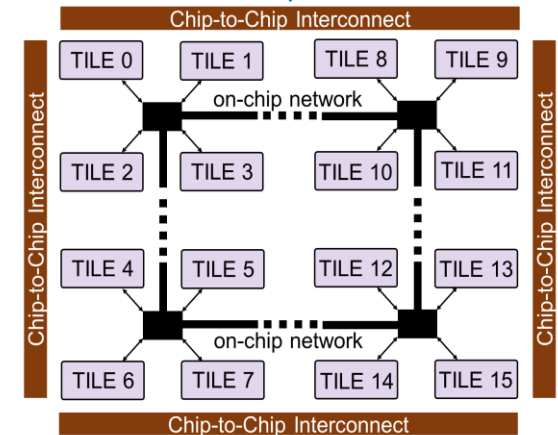# PUMA: Resistive Crossbar based Programmable Architecture

> ## Features

- Analyze the memory-compute characteristics of ML applications
- An ISA-programmable accelerator built with hybrid CMOS-NVM technology
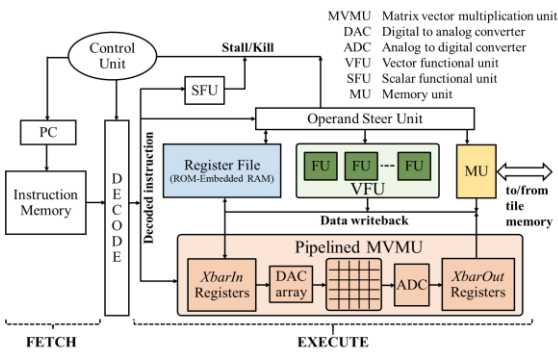


PUMA Chip
(dataflow architecture)

*A. Ankit et al, ASPLOS, 2019*



MVMU  Matrix vector multiplication unit
DAC   Digital to analog converter
ADC   Analog to digital converter
VFU   Vector functional unit
SFU   Scalar functional unit
MU    Memory unit

PUMA Core
(NVM Crossbar + Digital CMOS)

PUMA Tile (multi-core)

# Challenges: NVM devices

**Device**
High ON/OFF Ratio
Device Linearity
High Endurance
Multi-level cells
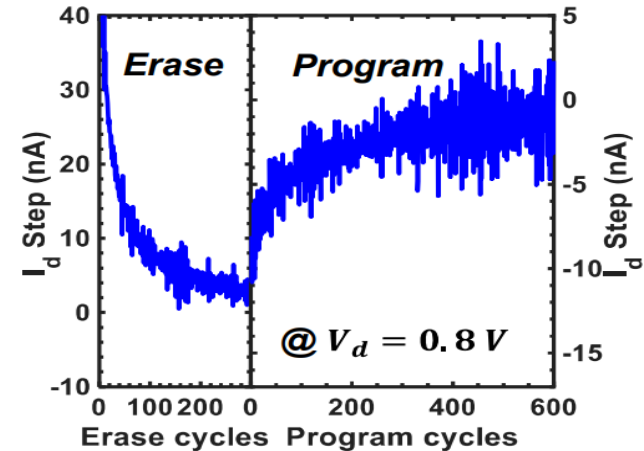
➢ ## Compared to CMOS:
- ✓ Non-volatility
- ✓ High density
- ✓ Low leakage
- ✓ Capable of in-memory compute
- ✗ Write energy/latency

➢ ## Current devices are highly non-linear

➢ Expensive write operations and peripheral circuitry

➢ $R_{ON}/R_{OFF}$ ratios are limited to ~10×

➢ RRAM has poor endurance.

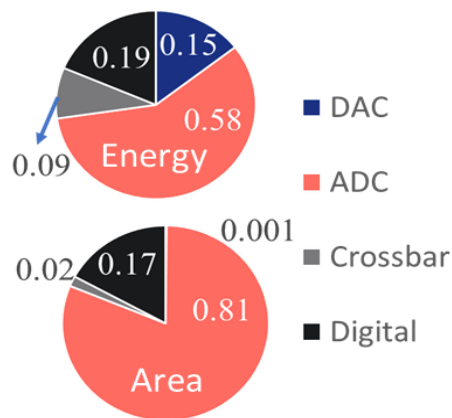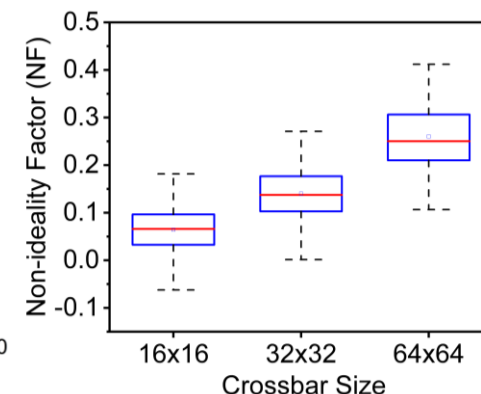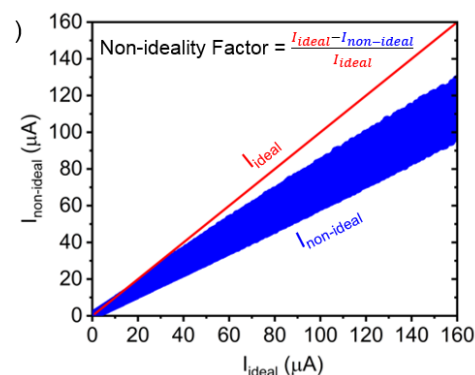➢ More than 4-bits/cell is not reliable yet.

[1] IEDM, 2019

| Property | PCM | RRAM | MTJ | CMOS |
|---|---|---|---|---|
| Multi-level cell | Yes | Yes | No | No |
| Storage Density | High | High | High | Low |
| $R_{ON}/R_{OFF}$ | High | High | Low | High |
| Non-volatility | Yes | Yes | Yes | No |
| Leakage | Low | Low | Low | High |
| Write Energy | 6 nJ | 2 nJ | < 1 nJ | < 0.1 nJ |
| Write Latency | 150 ns | 100 ns | 10 ns | < 1ns |
| Endurance | $10^7$ cycles | $10^5$ cycles | $10^{15}$ cycles | > $10^{16}$ cycles |

C-BRIC

# Challenges: NVM Compute Macro

**Macro**

Large Crossbar Size
Low cost peripheral overhead
Good selector device
Source/Sink/Line resistances

- NVM crossbars can have various non-idealities (parasitics, non-ideal devices)

- Such non-idealities can introduce varying amounts of functional errors based on different voltage and conductance

- Errors increase with higher crossbar sizes

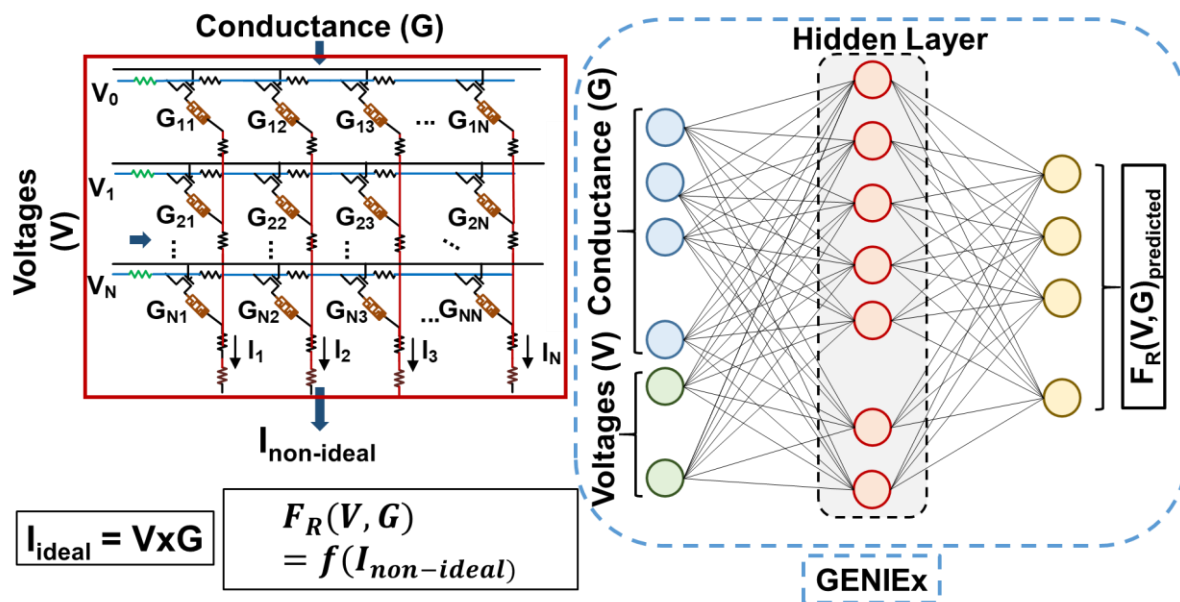- ADCs consume 58% and >80% of the total energy and area, respectively

Non-ideality Factor = $\frac{I_{ideal} - I_{non-ideal}}{I_{ideal}}$

| Analog Operations | Energy (pJ) | Dig. Operations | Energy (pJ) |
|---|---|---|---|
| MVM Energy | 3.84 | ALU | 25.6 |
| ADC Energy | 128 | Access (FMA) | 480 |
| Other peripherals | 12.8 | | |
| Total | 144.6 | Total | 505.6 |

8-bit MVM

C-BRIC

# GENIEx: A Generalized Approach to Emulating Non-Ideality in Memristive X-bars

$$I_{column} = f(V_i, G_{ij}(V), R_{source}, R_{sink}, R_{wire})$$

- ➢ $f$ is a data-dependent non-linear function.

- ➢ Neural networks are efficient tools for capturing the close inter-dependence of its inputs.

- ➢ Neural network to model the behavior of non-ideal crossbars



$I_{ideal} = V \times G$

$$F_R(V, G) = f(I_{non-ideal})$$

Chakraborty et al, DAC 2020,
https://arxiv.org/pdf/2003.06902.pdf

- ➢ GENIEx provides modeling capability for different non-idealities

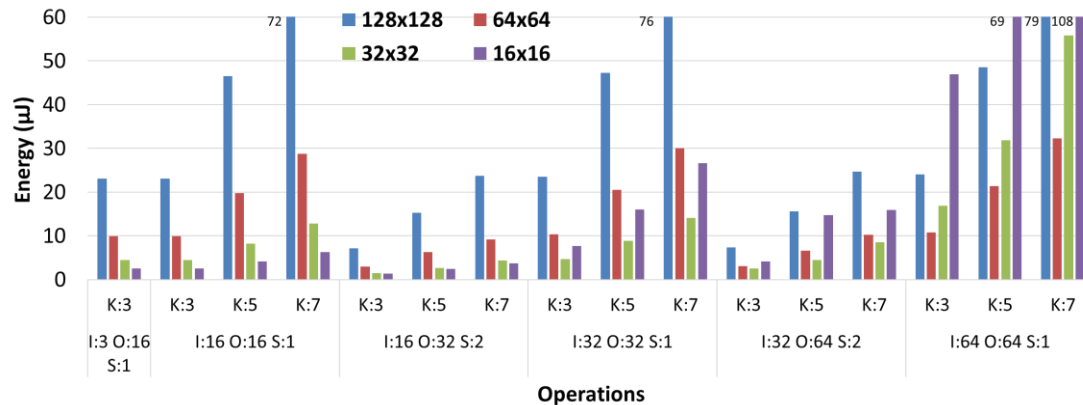C-BRIC

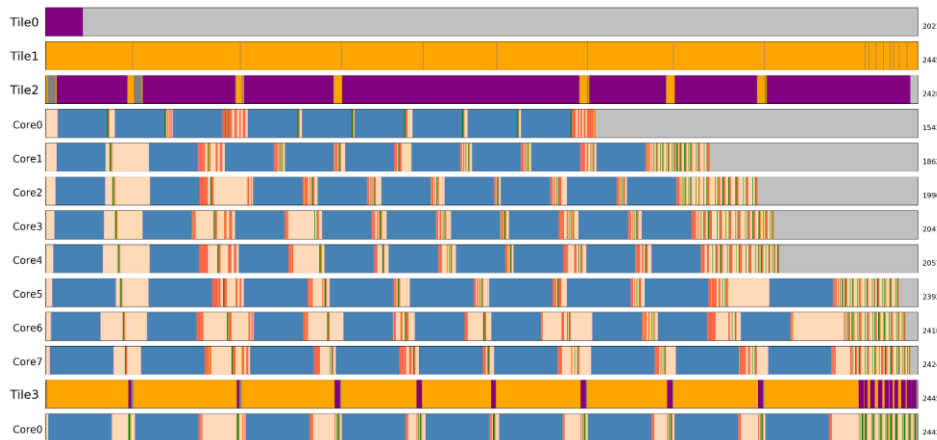# Resistive Crossbar Based Accelerator Design Flow

# Performance Simulator – Scope



- **Design space exploration of ML kernels**

- Efficiency depends on multiple parameters
  - Workload properties
  - Architecture configuration
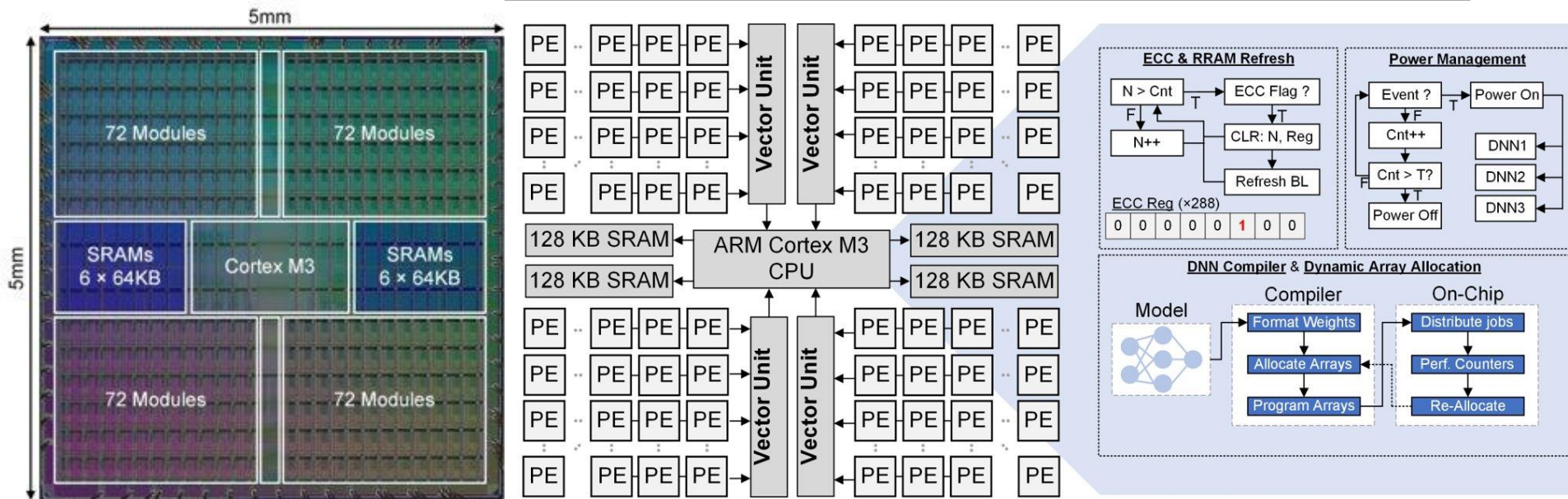  - Runtime Utilization

- **Performance Bottleneck Analysis**

- Runtime characteristics has complex dependency of workload and hardware properties

- CNNs show upto 13.0× reduction (least). **High weight reuse, even at batch-size 1.**

- MLPs show upto 80.1× reduction. **No weight reuse, small models.**

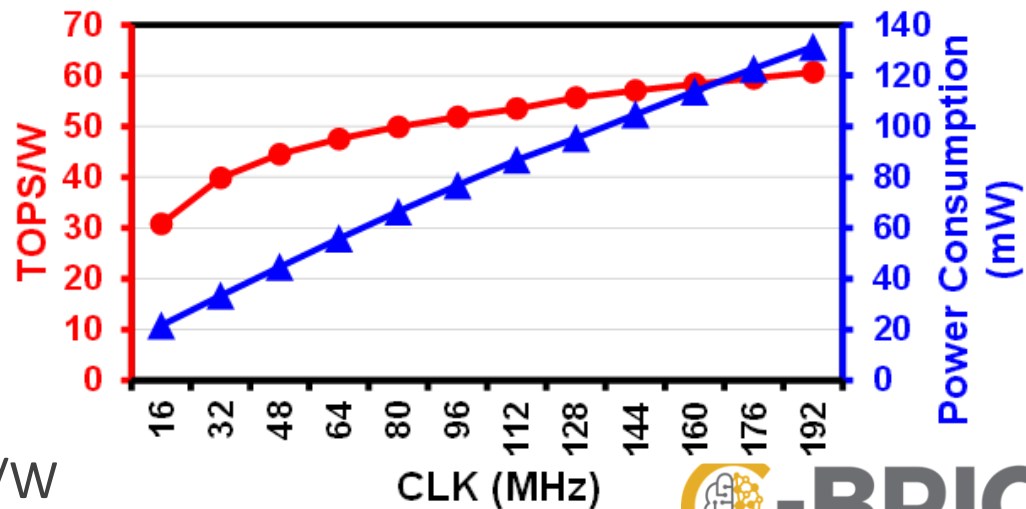- LSTMs show upto 2446× reduction. **Little Weight reuse, large models (billions of parameters).**

# Gen2 N40 RRAM CIM with Embedded Processor
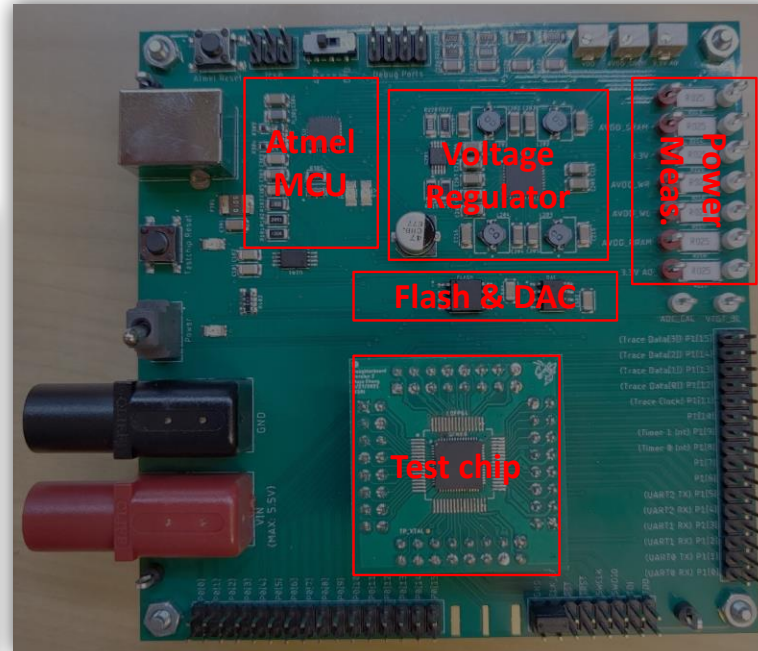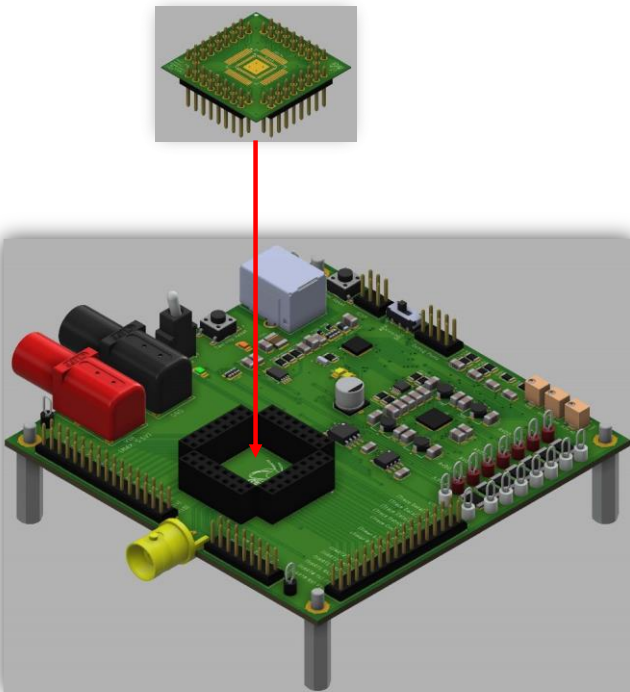


Technology: TSMC 40nm RRAM

**Key Innovation**

➢ Full system demonstration with Embedded Cortex M3 processor

➢ Highest effective RRAM density with >3X improvement of array density w.r.t. SOTA and >50 TOPS/W

Raychowdhury, GaTech, ISSCC 2022

# Evaluation Board



| Block | Model |
|---|---|
| MCU | Atmel Atmega32u2 |
| Voltage Regulator | Analog Device LTC3676 |
| External Flash | Adesto Tech. AT45DB321E |
| DAC | Burr-Brown DAC7612 |

https://muyachang.github.io/rram-pyterminal/

➤ Full python programmability and OS support.

➤ Currently being used as a test-vehicle for both research and undergraduate teaching

➤ Planning to share the evaluation board with CBRIC PIs so that we can use this as a test-bed for algorithmic and embedded system research
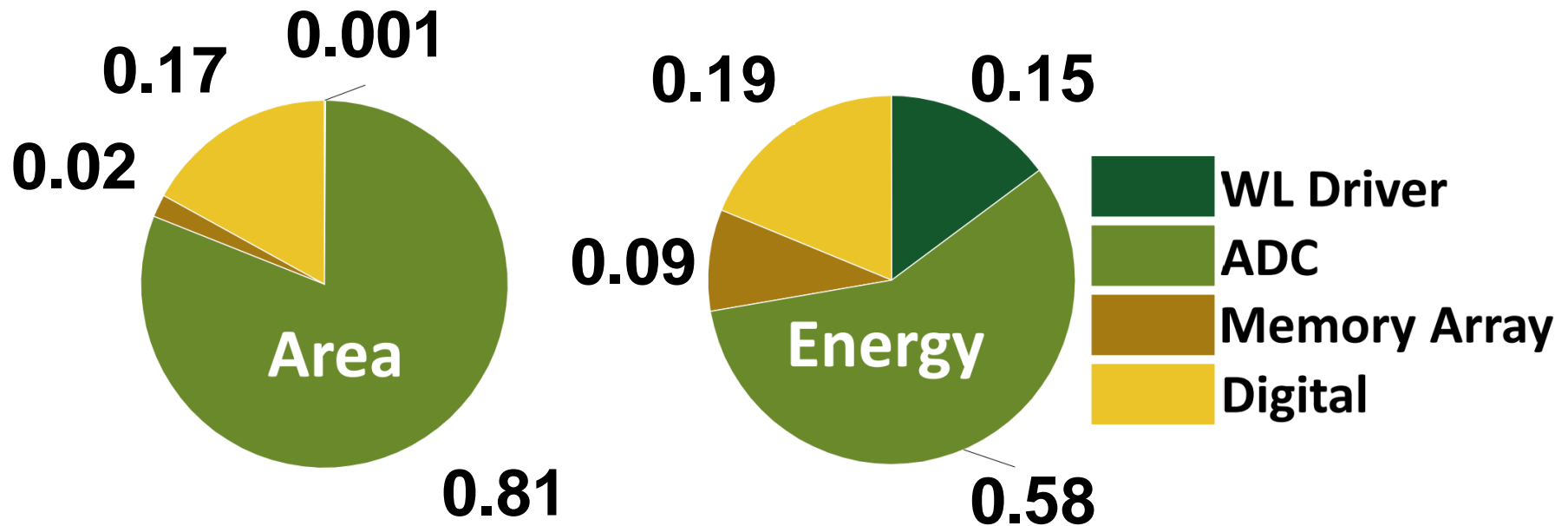
C-BRIC

# System Demonstration

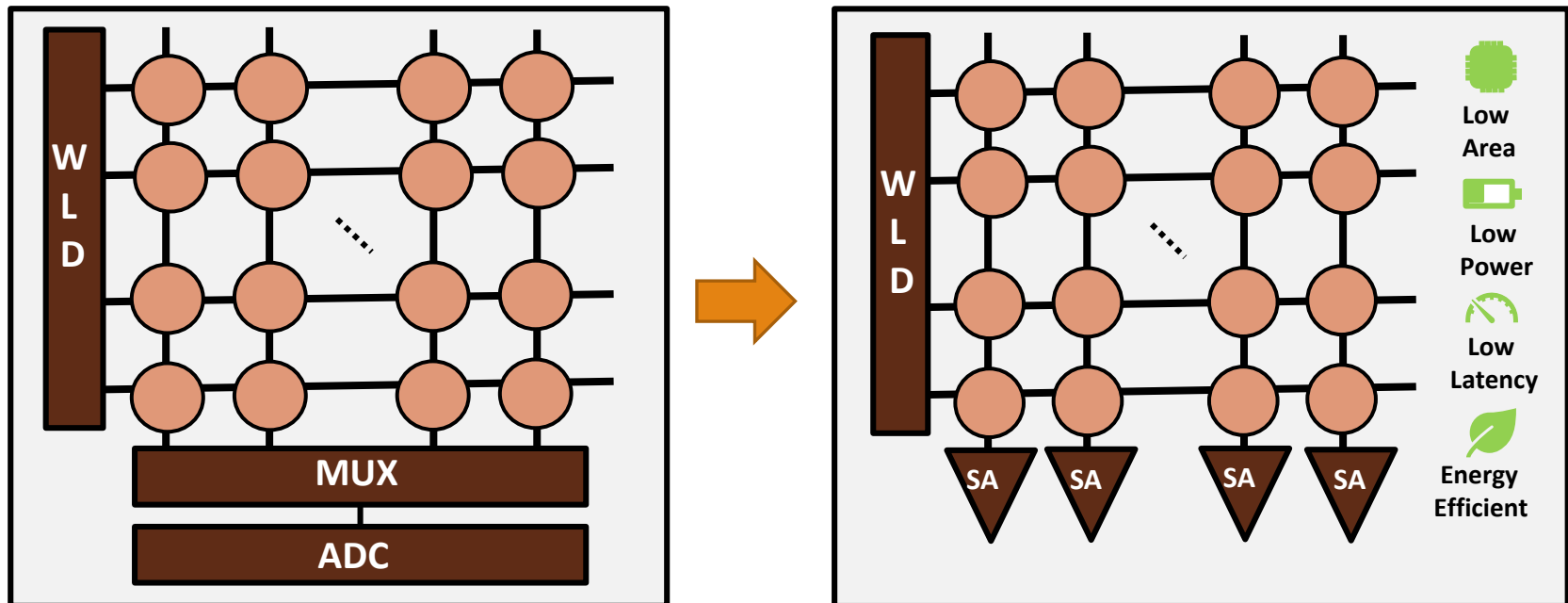# Revisit ADC: Near ADC-less CiM

# ADC overhead in CiM accelerators

➤ Large percentage of area and energy profile dominated by ADCs.

# Mitigating Overhead with ADC-Less Design

➤ Area and energy profile dominated by ADCs in CiM accelerators.

➤ ADC-Less Design: Use Sense Amplifiers for Analog to Digital conversion of Array output
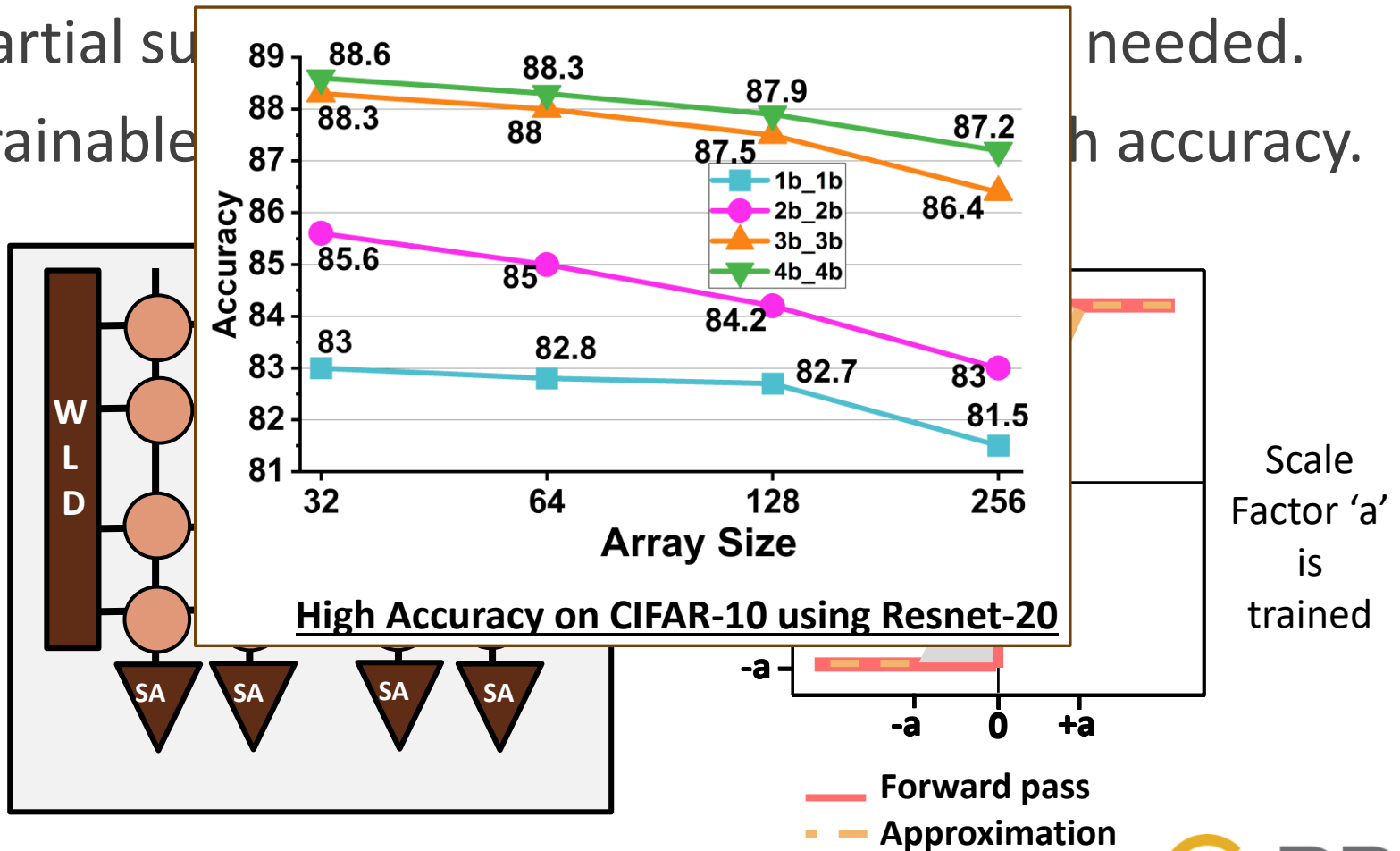


Saxena, Chakraborty, Roy, Towards ADC-Less Compute-In-Memory Accelerators for Energy Efficient Deep Learning, DATE 2022

# SW co-design for ADC-Less CiM Accelerators

➤ ADC-Less CiM accelerators have 1b partial sums.

➤ Partial su[...] needed.

➤ Trainable [...] accuracy.



**High Accuracy on CIFAR-10 using Resnet-20**

Scale Factor 'a' is trained

Forward pass

Approximation

# CMOS SRAM and Non-volatile Memories

| Property | PCM | RRAM | MTJ | CMOS (SRAM) |
|---|---|---|---|---|
| Multi-level cell | Yes | Yes | No | No |
| Storage Density | High | High | High | Low |
| $R_{ON}/R_{OFF}$ | High | High | Low | High |
| Non-volatility | Yes | Yes | Yes | No |
| Leakage | Low | Low | Low | High |
| Cell Area | $16F^2$ | $16F^2$ | $30\text{-}80F^2$ | $160F^2$ (6T), $231F^2$ (8T) |
| Write Energy | 6 nJ | 2 nJ | < 1 nJ | < 0.1 nJ |
| Write Latency | 150 ns | 100 ns | 10 ns | < 1ns |
| Endurance | $10^7$ cycles | $10^5$ cycles | $10^{15}$ cycles | > $10^{16}$ cycles |

C-BRIC

# In-Memory Bit-wise Vector Boolean Operations

[ Modifying the Peripheral Sensing Circuits to Read-out A Vector Boolean Function ]



Simultaneous Activated WLs

Processing Core

Memory Peripherals

Modified Memory Sensing Peripherals
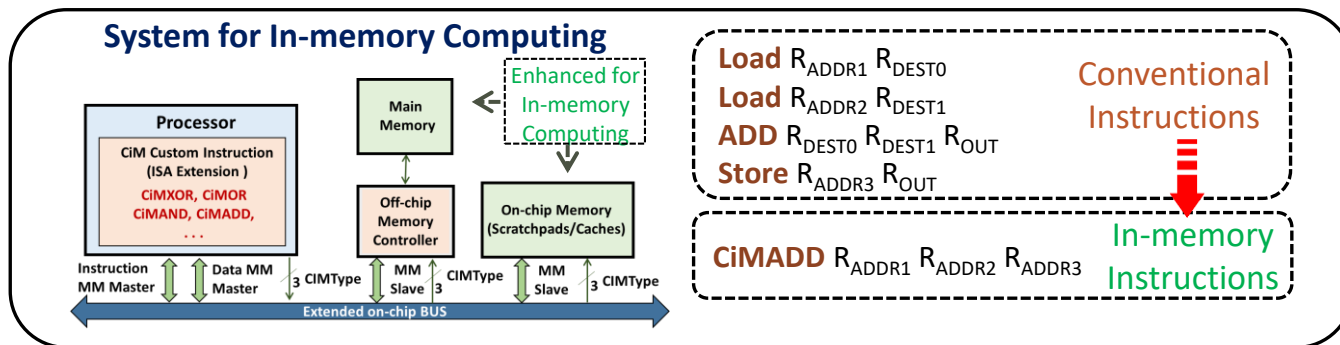
C-BRIC

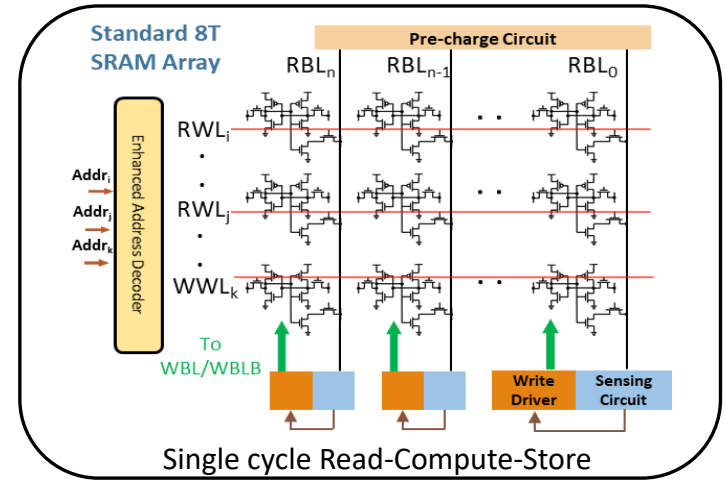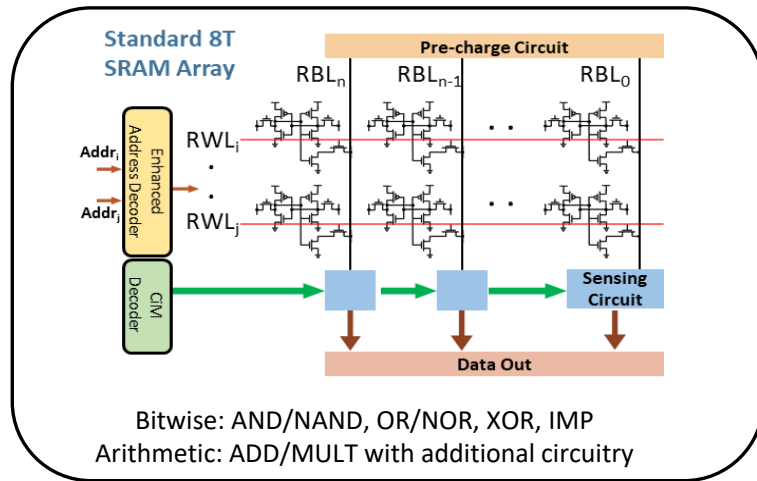# SRAM: In-Memory Bit-wise Boolean Operations
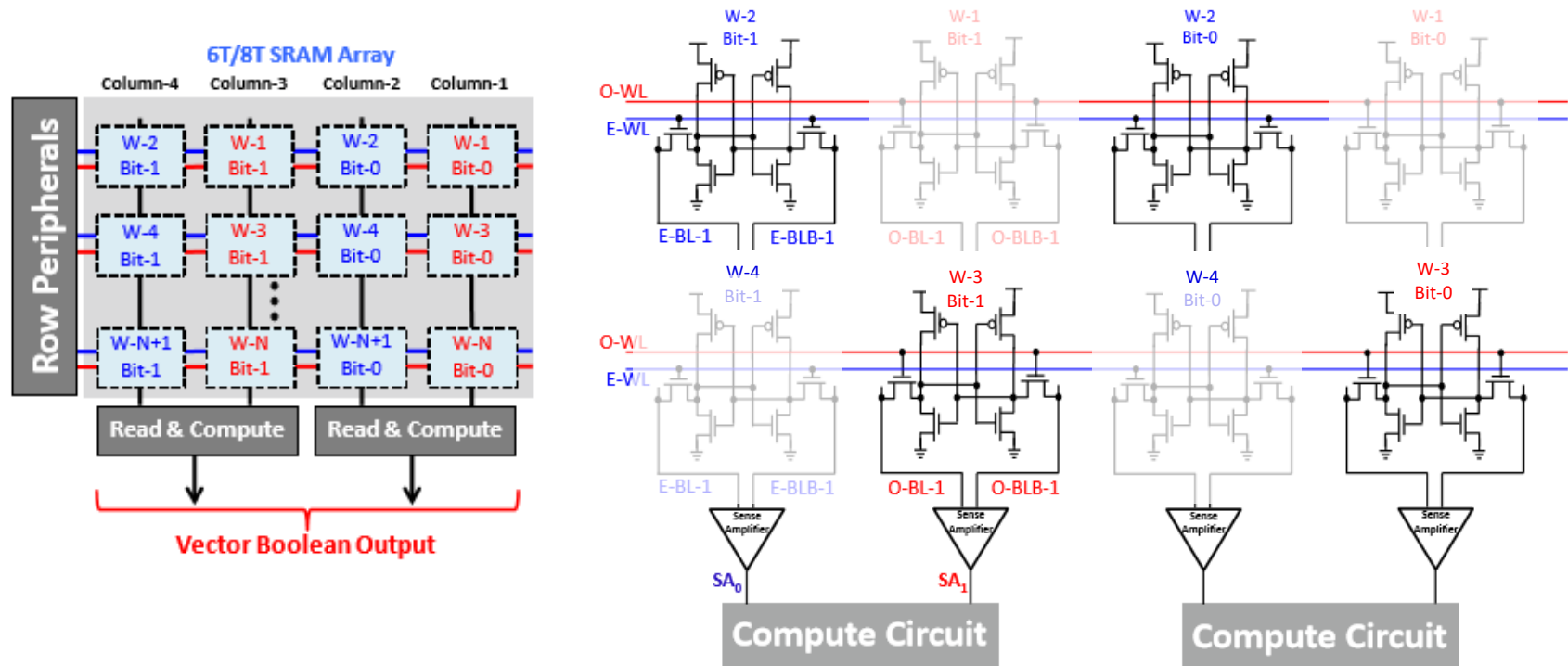


*Shanbhag et. al.*

## 6T-SRAM
➤ Staggered WL activation to avoid short circuits between cells.
➤ Asymmetric SAs help detect bitwise NAND/NOR/XORs

# X-SRAM: Bit-Wise Vector Boolean Operations



Bitwise: AND/NAND, OR/NOR, XOR, IMP
Arithmetic: ADD/MULT with additional circuitry

Single cycle Read-Compute-Store



Load $R_{ADDR1}$ $R_{DEST0}$
Load $R_{ADDR2}$ $R_{DEST1}$
ADD $R_{DEST0}$ $R_{DEST1}$ $R_{OUT}$
Store $R_{ADDR3}$ $R_{OUT}$

Conventional Instructions

CiMADD $R_{ADDR1}$ $R_{ADDR2}$ $R_{ADDR3}$

In-memory Instructions

*Agrawal, A et al., 2018. X-sram,. IEEE TCAS-I*

C-BRIC

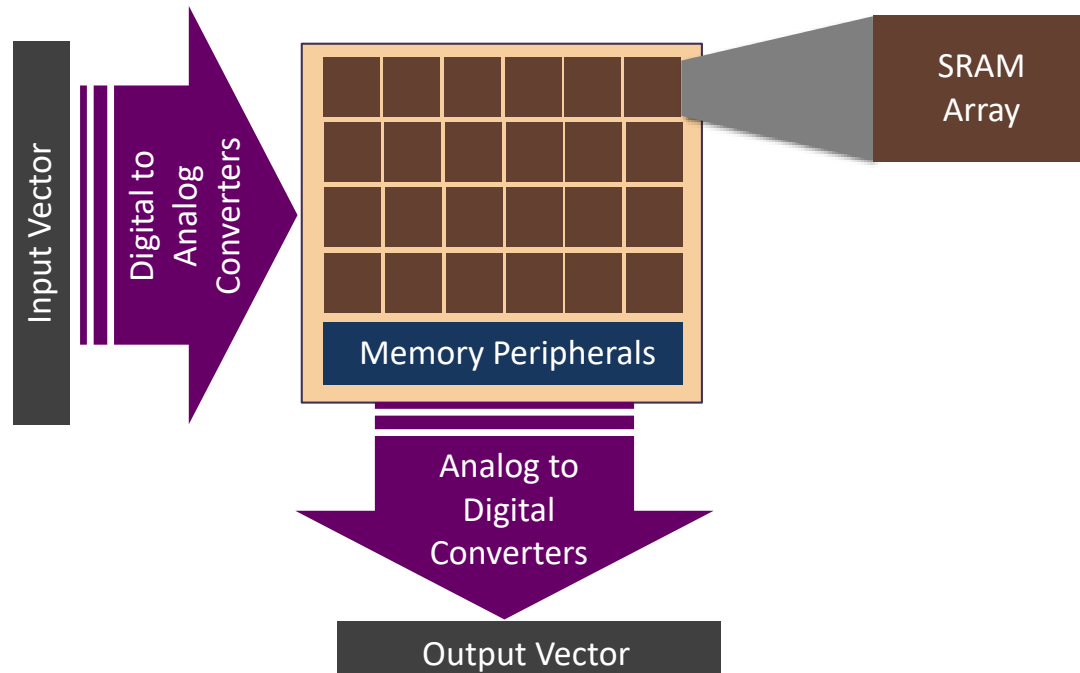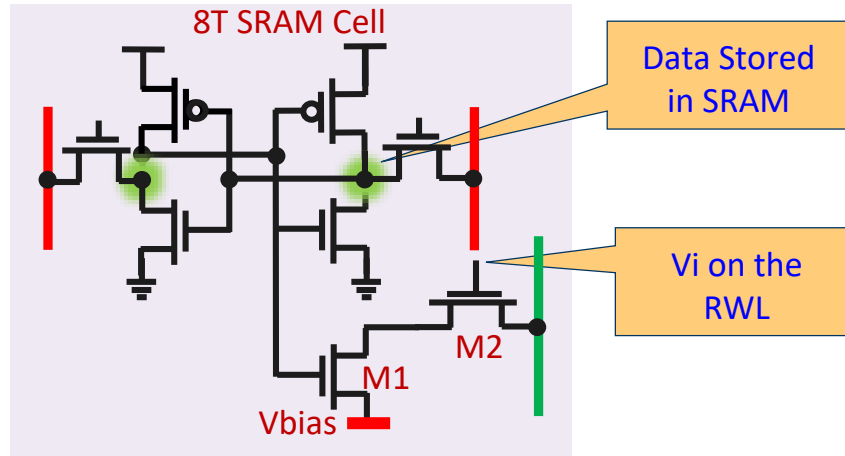# i-SRAM: Interleaved Wordlines for In-Memory Vector Boolean Operations



## 8T-SRAM

➢ Each row has two read word-lines, and bit-lines are connected to read and compute blocks

➢ Interleaved 6T cells with bit-lines connected to sense-amplifiers then compute circuits.

➢ The circuit schematic for NAND(AND)

➢ The circuit scheme for NOR(OR)

Jaiswal, et al. "i-SRAM: Interleaved Wordlines for Vector Boolean Operations Using SRAMs." *IEEE Trans. CAS I:* (2020).

# In-memory Dot Product Computations

In-Memory Dot Product Acceleration by use of Current-Mode Computations in SRAMs
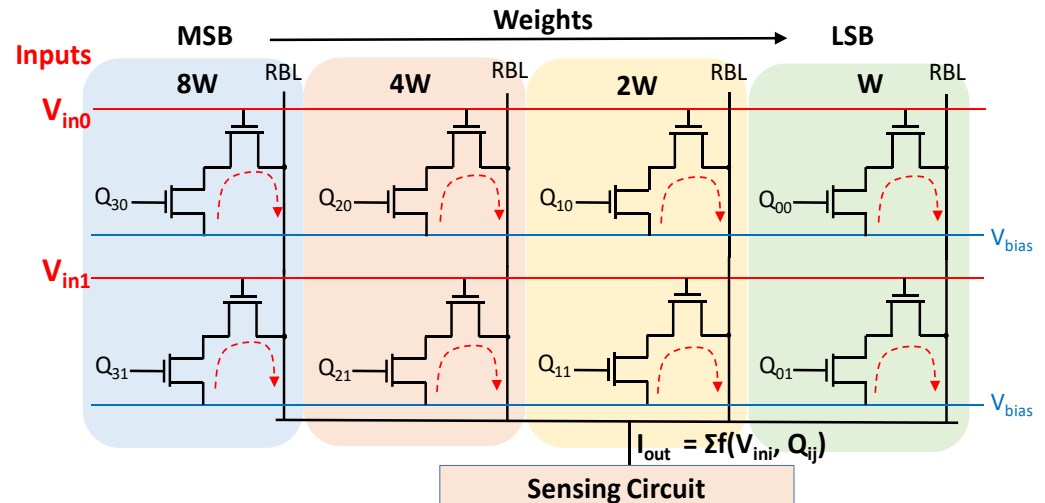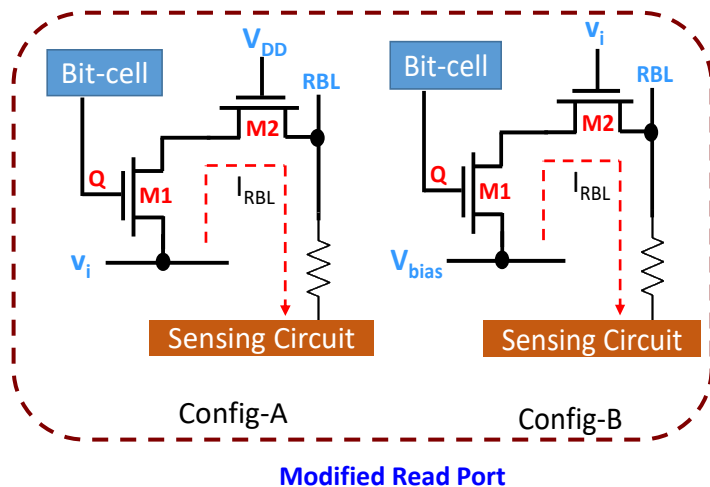
# 8T SRAM as a Multi-Bit Dot Product Engine



8T SRAM Cell

Data Stored in SRAM
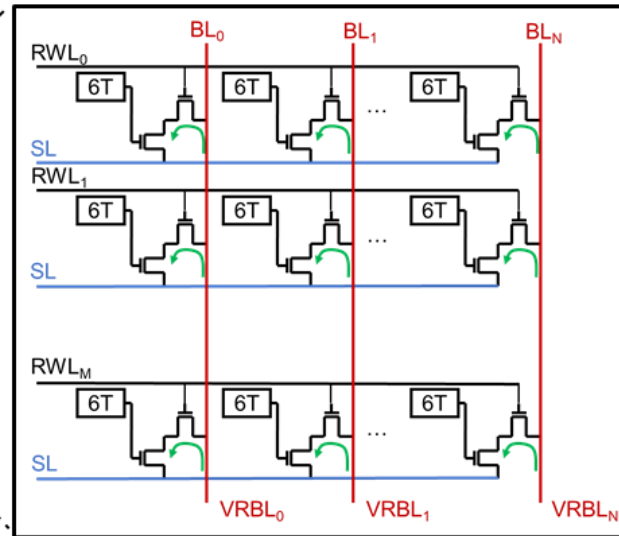
Vi on the RWL

M2

M1

Vbias

The Dot Product = $\Sigma V_i \cdot W_i$

$V_i$: Analog Voltages on RWL or SL

$W_i$: Data Stored in SRAM

Summation: KCL addition



Bit-cell

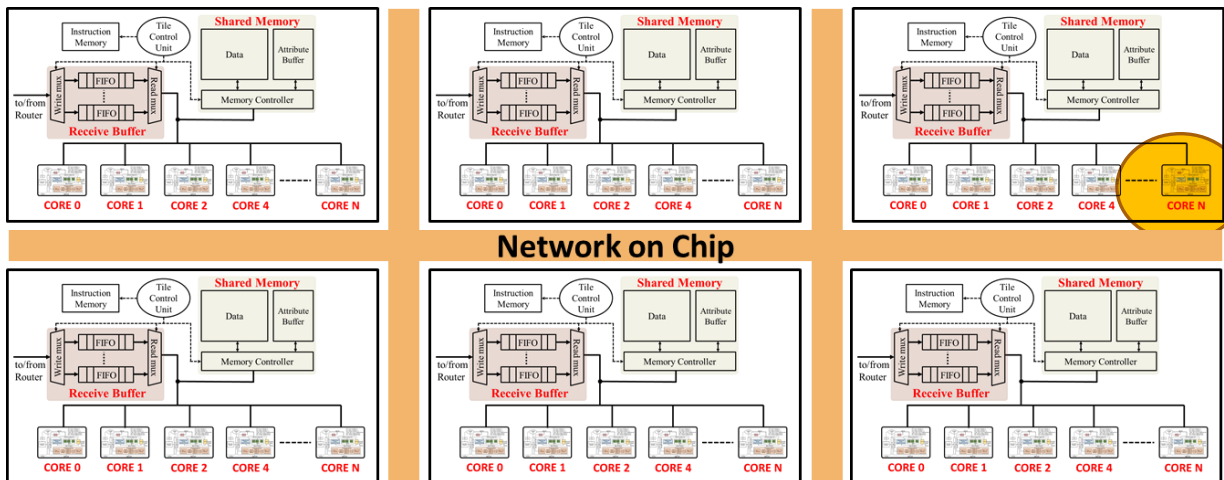$V_{DD}$

RBL

M2

Q

M1

$I_{RBL}$

$V_i$

Sensing Circuit

Config-A

Bit-cell

$V_i$

RBL

M2

Q

M1

$I_{RBL}$

$V_{bias}$

Sensing Circuit

Config-B

Modified Read Port

MSB — Weights — LSB

Inputs

$V_{in0}$

$V_{in1}$

8W    4W    2W    W

RBL   RBL   RBL   RBL

$Q_{30}$   $Q_{20}$   $Q_{10}$   $Q_{00}$

$Q_{31}$   $Q_{21}$   $Q_{11}$   $Q_{01}$

$V_{bias}$

$V_{bias}$

$I_{out} = \Sigma f(V_{ini}, Q_{ij})$

Sensing Circuit

C-BRIC

# CiM Macro



1b IN -1b W compute logic

| RWL | Q | Current |
|---|---|---|
| 0 | 0 | No BL Discharge |
| 0 | 1 | No BL Discharge |
| 1 | 0 | No BL Discharge |
| 1 | 1 | BL Discharge |

ADC_OUT[9:0]     ADC_OUT[9:0]

Step 1: Precharge BL to VDD
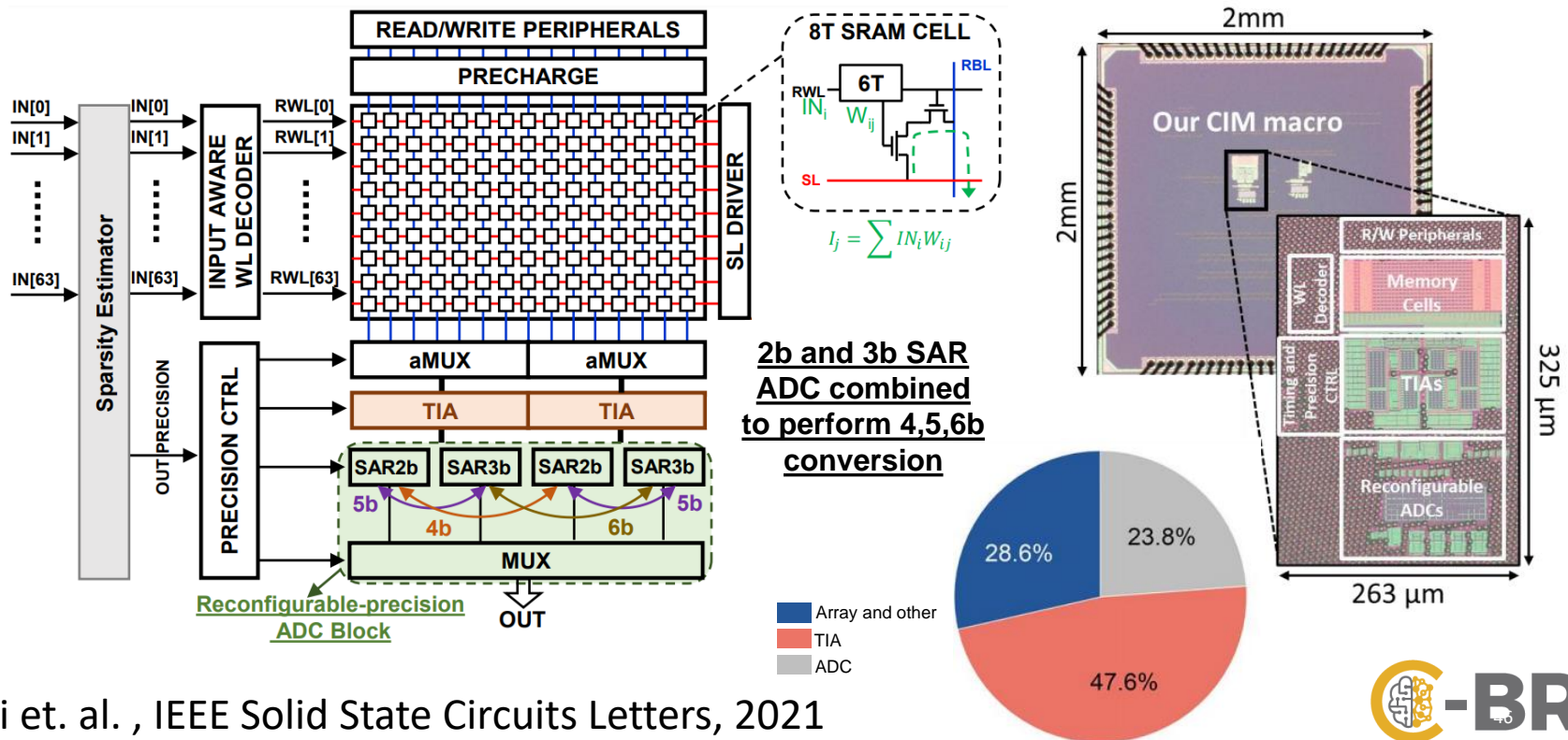
Step 2: Discharge BL based on RWL and cell content

➤ More the number of ones in input or weights in a column, more discharge of BL.

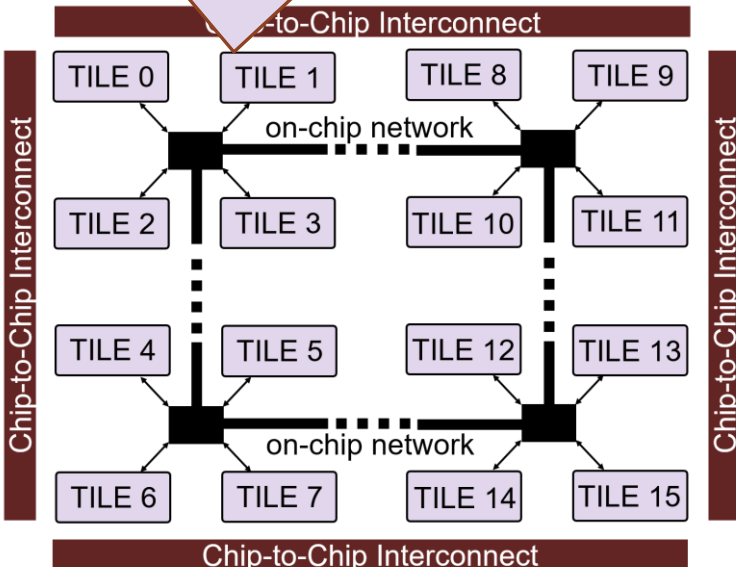➤ $V_{RBL}$ proportional to MAC output between inputs and weights.
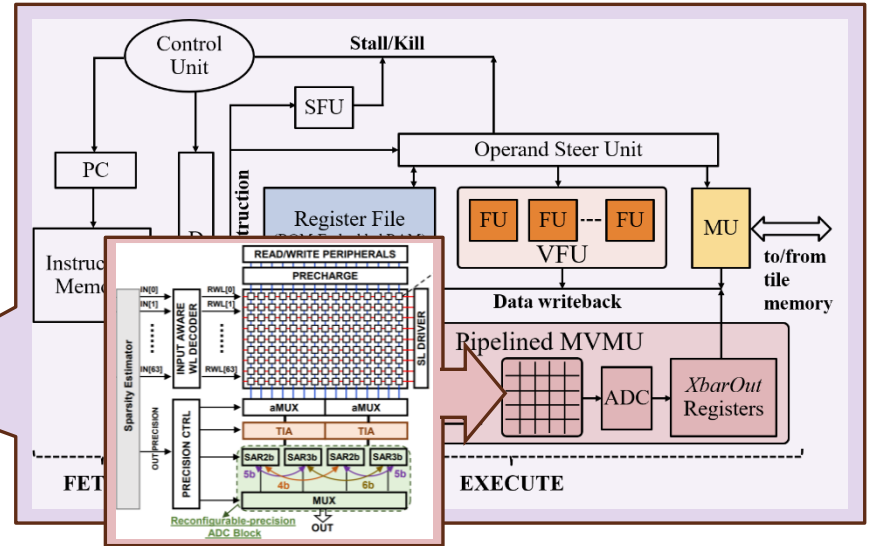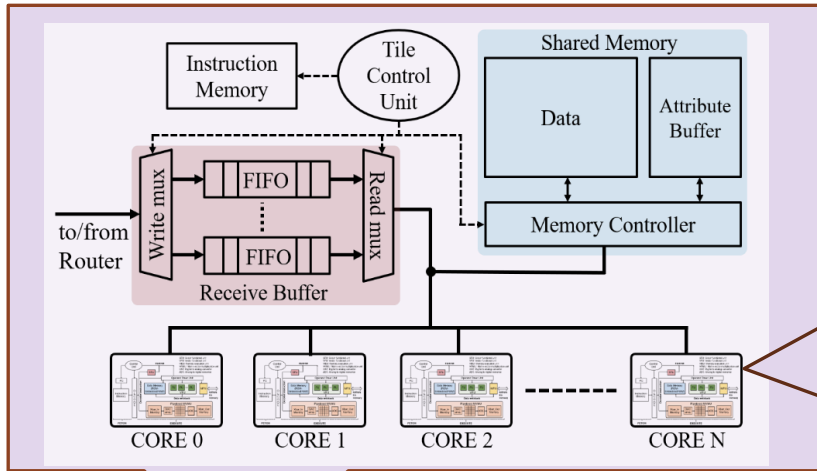
44

# IMC CMOS Cores (ANNs)
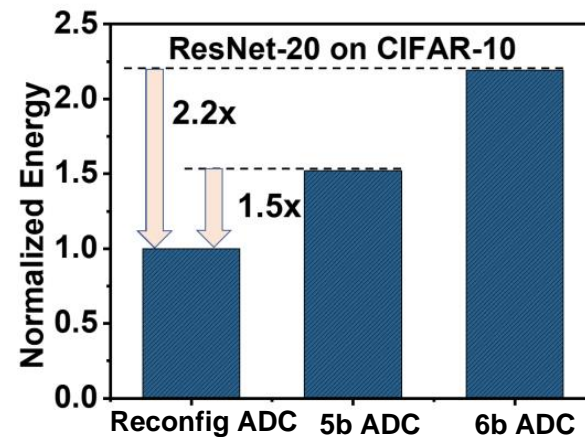
# CMOS Array: Sparsity Aware CiM Macro design

➤ Bit serial CiM acceleration provides opportunities to leverage abundant bit level sparsity.

➤ Different levels of input sparsity can be leveraged by dynamically reconfiguring ADC precision (2b-6b).
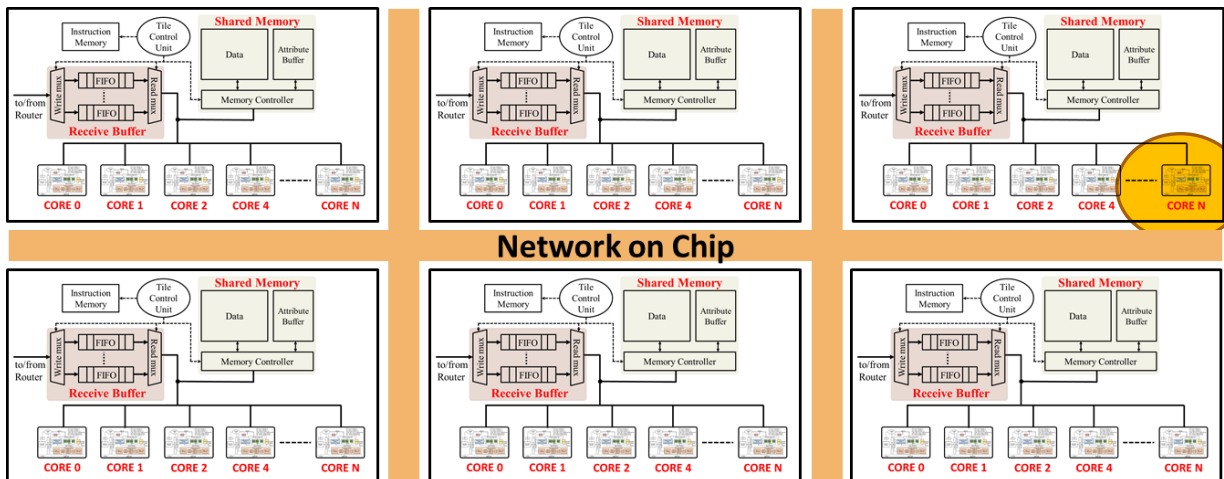


**2b and 3b SAR ADC combined to perform 4,5,6b conversion**

$I_j = \sum IN_i W_{ij}$

A. Ali et. al. , IEEE Solid State Circuits Letters, 2021

# System with sparsity aware Core



Sparsity aware macro

2.2X improvement in inference energy

TILE

CROE

NODE

CORE 0  CORE 1  CORE 2  CORE N

Instruction Memory

Tile Control Unit

Shared Memory

Data  Attribute Buffer

Write mux  FIFO  Read mux

to/from Router

Receive Buffer

Memory Controller

Control Unit  Stall/Kill

SFU

PC

Operand Steer Unit

Instruction Memory

Register File

FU FU FU  MU

VFU

Data writeback

to/from tile memory

Pipelined MVMU

ADC  XbarOut Registers

FETCH  EXECUTE

Chip-to-Chip Interconnect

TILE 0  TILE 1  TILE 8  TILE 9

on-chip network

TILE 2  TILE 3  TILE 10  TILE 11

TILE 4  TILE 5  TILE 12  TILE 13

on-chip network

TILE 6  TILE 7  TILE 14  TILE 15

Chip-to-Chip Interconnect

Chip-to-Chip Interconnect

Chip-to-Chip Interconnect

ResNet-20 on CIFAR-10

Normalized Energy

2.2x

1.5x

Reconfig ADC  5b ADC  6b ADC
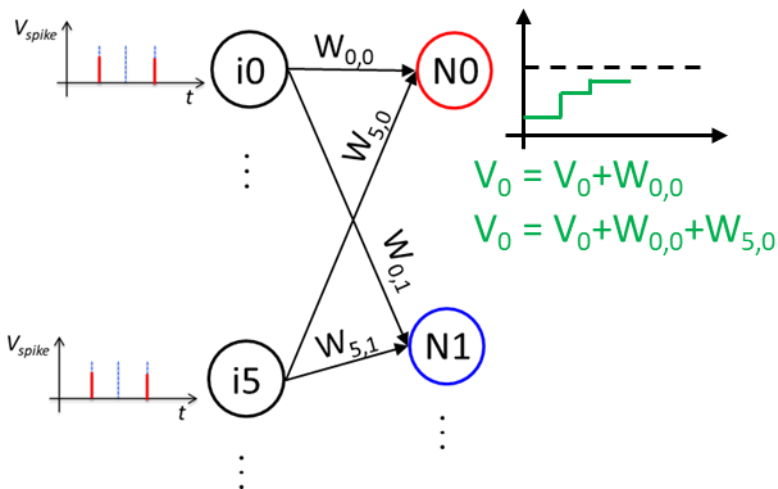
C-BRIC

# IMC CMOS Cores (SNNs)

# SNNs

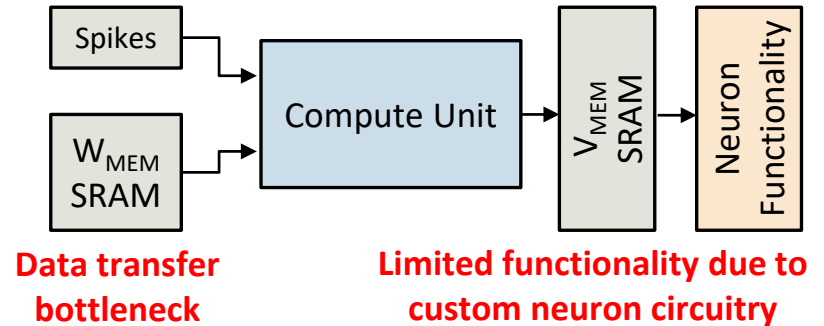➢ **Challenges in current digital SNN hardware:**

o Data transfer bottleneck to/from memory.

o Additional SNN-specific data movements for processing $V_{MEM}$ for multiple timesteps.

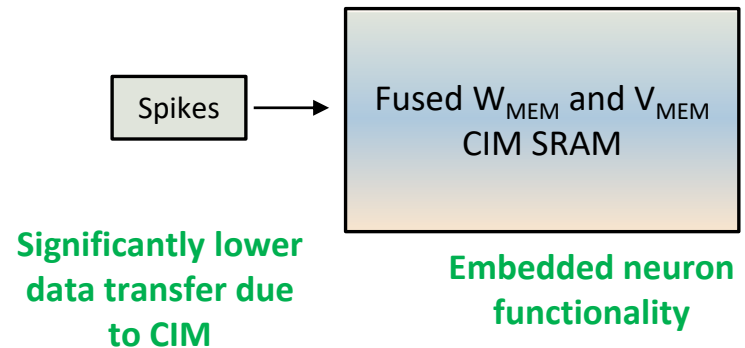o Limited functionality due to area and power expensive custom neuron circuitry.

➢ **Approach:**

✓ Fused $W_{MEM}$ and $V_{MEM}$ CIM Array integrating all processing modes required for SNN inference – accumulate, threshold, reset etc.
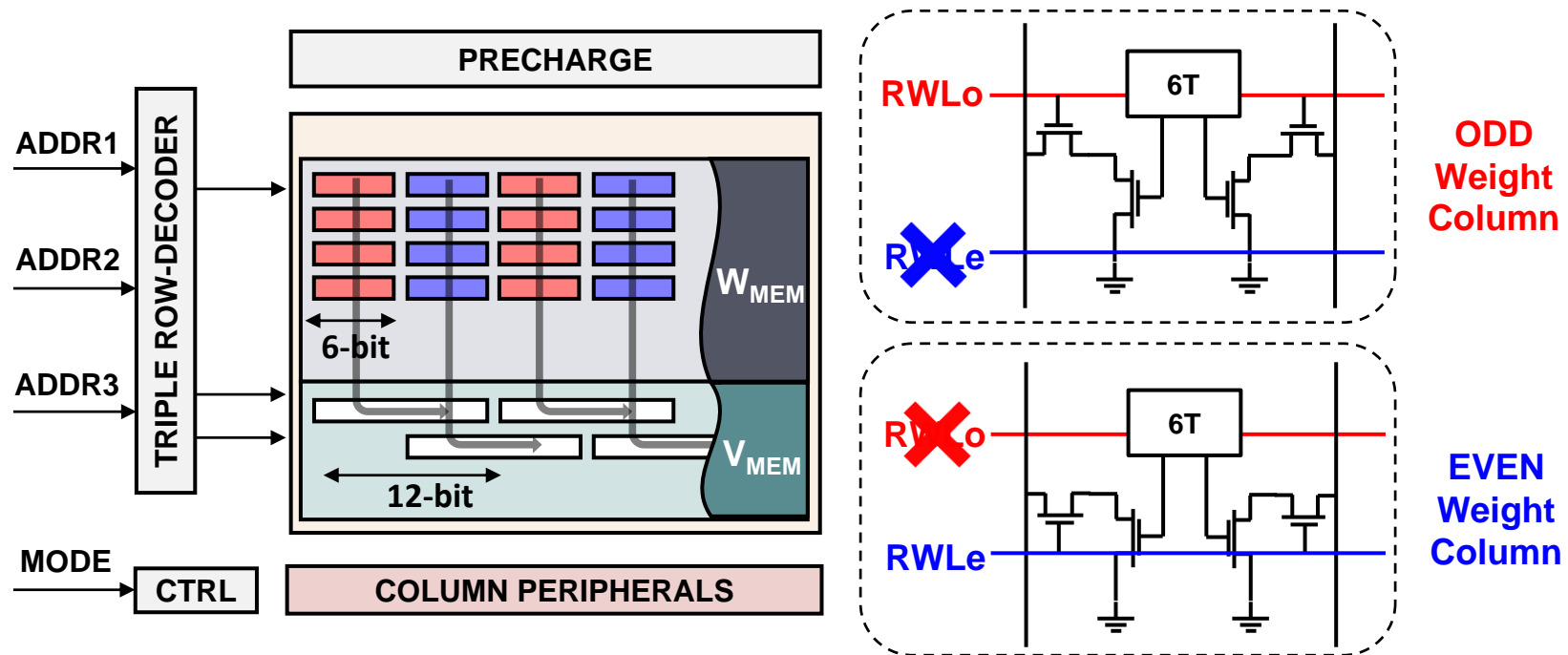
**Conventional Approach:**



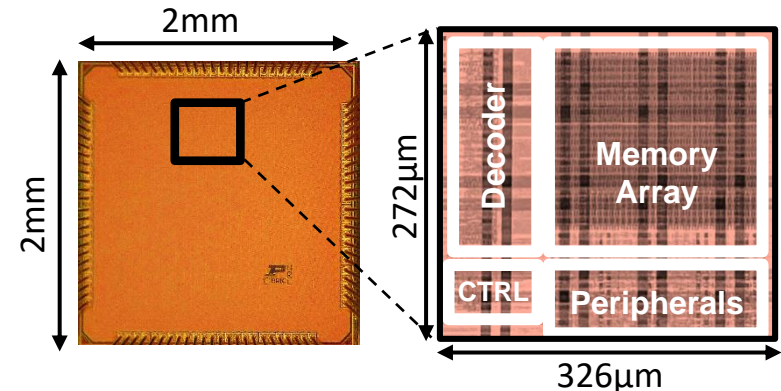**Data transfer bottleneck**

**Limited functionality due to custom neuron circuitry**

**Proposed Approach:**



**Significantly lower data transfer due to CIM**

**Embedded neuron functionality**



$$V_0 = V_0 + W_{0,0}$$
$$V_0 = V_0 + W_{0,0} + W_{5,0}$$

C-BRIC

# Organization of the Fused $W_{MEM}/V_{MEM}$ Array



- ➢ **Different bit-precision** requirements for weights and $V_{MEM}$.
- ➢ We propose a mapping strategy to **efficiently use the SRAM area** with **minimal peripheral complexity**:
  - o Fit more weights in each row using odd/even RWLs.
  - o Staggered alignment of corresponding $V_{MEM}$ data.
  - o Same peripherals are used by reconfiguring them in odd/even cycles.

# In-memory Computing Macro for SNNs: Chip Summary

- We propose a 10T SRAM based CIM macro for SNN inference.

- The macro consists of a fused $W_{MEM}$ and $V_{MEM}$ and supports all processing modes required for SNN inference - accumulate, threshold, spike-check, reset etc.

- The macro also supports multiple neuron functionalities through various instruction sequences.

- The macro also leverages sparsity in the input spikes for energy-efficiency.

- The prototype chip was fabricated in 65nm LP CMOS process, achieves an energy-efficiency of 0.99 TOPS/W @ 0.85V, 200MHz, for signed 11-bit operations.

- We demonstrate sentiment classification using the intrinsic dynamics of SNNs achieving competitive accuracies with LSTMs.



| Technology | 65nm |
|---|---|
| Macro Area | 0.089 mm² |
| Cell Type | 10T |
| Memory | $W_{MEM}$ : 9kb $V_{MEM}$ : 2.25kb |
| Weight/Vmem bits | 6-bit/11-bit |
| Supply Voltage | 0.7 ~ 1.2 V |
| Max. Frequency | 500 MHz |
| Energy Efficiency | 0.99 (TOPS/W) @200MHz, 0.85V (signed 11-bit op) |

A. Agrawal et. al. , "IMPULSE….",  IEEE Solid State Circuits Letters, 2021

C-BRIC

# Embedding ROM in a RAM

## Embedding ROM in CMOS and 1T-1R Arrays Enabling Near-Memory Computing through Lookup Tables



Processing Core

Memory Peripherals

ROM & RAM

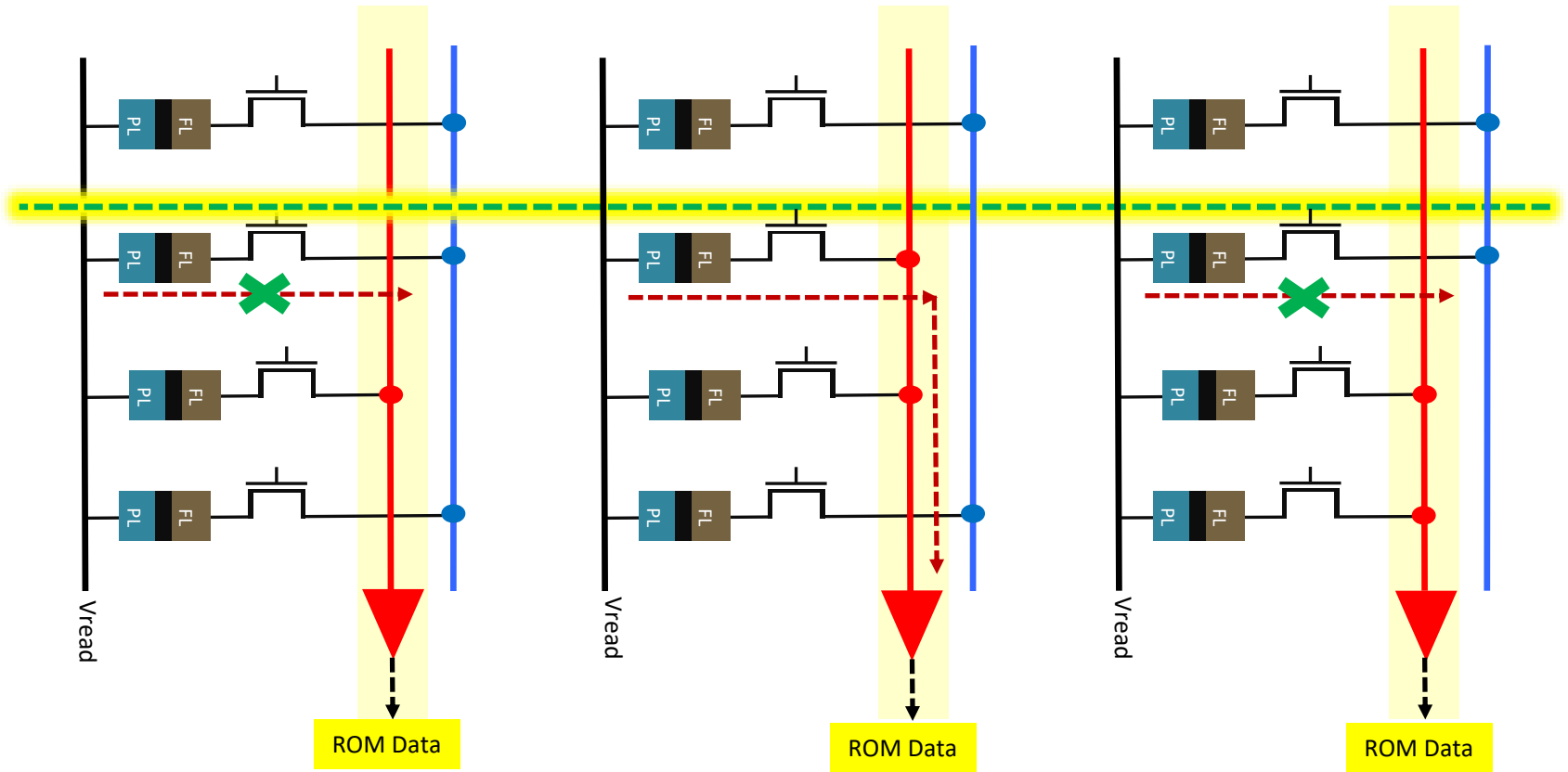Modified Bit-Cell Storing Both ROM and RAM Data

C-BRIC

# Embedding ROMs in RAMs (NVMs)



Both ROM and RAM data are stored in the same bit-cell. The read cycle determines whether the ROM or the RAM data is being read.

*D. Lee, K. Roy, IEEE EDL 2013*

# Embedding ROMs in RAMs (STT-MRAM)



*D. Lee, K. Roy, IEEE EDL 2013*

# Computing-In-DRAM

# In-DRAM Low-cost Bit-serial Addition

➤ **Majority-based vector addition**

$$C_{out} = Maj(A, B, C_{in})$$

$$S = Maj(A, B, C_{in}, \overline{C_{out}}, \overline{C_{out}})$$

➤ **Store data vectors in column-based fashion**

➤ **Same subarray peripheral circuits**

➤ **Add 9 reserved rows for compute (<1% area overhead)**



➤ **The result of A+B addition is stored in the same column**

➤ **Massively-parallel vector additions in bit-serial mode**

➤ **No need for carry shifts across bitlines!**

Ali, Jaiswal, Roy, "In-Memory Low-Cost Bit-Serial Addition Using Commodity DRAM Technology," TCAS-I, 2019

# In-DRAM Low-cost Bit-serial Addition

## Multiple row activation to calculate Maj(A,B,C,D,E)



(1) Initial state    (2) Enable WLs    (3) Enable Sense amp

# An example of 1-bit addition of two vectors A and B



**For n-bit vector addition, 4n+1 operations are needed**

(0) Initial state  (1) Copy A  (2) Copy B  (3) Calculate $C_{out}/\overline{C_{out}}$  (4) Calculate S

# Case Study: Compute-in-DRAM based k-NN Acceleration

- Checks *k* closest samples, Query vector assigned to the group that holds majority among those k samples
- Mainly consists of two computation stages: Distance computation and Global top-k sort
- The k value: The number of closest samples to be checked
- One-to-one distance computation: a large # of memory accesses



| Processor | X86, 2GHz |
|---|---|
| L1 Cache | 32KB I- and D-Cache |
| L2 Cache | 2 MB |
| Main memory | 1024 MB, DDR3-1600, 1 channel, 1 rank, 8 banks |

Using modified gem5 simulator from *S.Xu et. al,ICAL'19*



**Host Processor**

(Top-*k* Sort)

*k*NN Output

Add/Sub requests

Sum/Diff data

**Main Memory &
In-DRAM Computing**

(Manhattan Distance)

Test Data

Train Data

**11.7x performance improvement**

Ali, Jaiswal, Roy, "In-Memory Low-Cost Bit-Serial Addition Using Commodity DRAM Technology," TCAS-I, 2019

C-BRIC

# Multiplication Primitive



|  |  | | $A_3$ | $A_2$ | $A_1$ | $A_0$ |
|---|---|---|---|---|---|---|
|  |  | | $B_3$ | $B_2$ | $B_1$ | $B_0$ |
|  |  | | $A_3B_0$ | $A_2B_0$ | $A_1B_0$ | $A_0B_0$ |
|  |  | $A_3B_1$ | $A_2B_1$ | $A_1B_1$ | $A_0B_1$ | |
|  | $A_3B_2$ | $A_2B_2$ | $A_1B_2$ | $A_0B_2$ | | |
| Carry out | $A_3B_3$ | $A_2B_3$ | $A_1B_3$ | $A_0B_3$ | | |
| $P_7$ | $P_6$ | $P_5$ | $P_4$ | $P_3$ | $P_2$ | $P_1$ | $P_0$ |

**AND** ➕ **ADD**

**compute rows**

**DRAM subarray**

- ➤ Multiplication can be broken down in terms of AND and ADD operations.
- ➤ Multiplication operation reserves 9 compute rows in the DRAM subarray.

# PIM-DRAM Architecture



Roy, Ali, Raghunathan, PIM-DRAM: Accelerating Machine Learning Workloads using Processing in Commodity DRAM, IEEE JETCAS, 2021

*4-bit act. and W*

- ➢ Special Function Units consist of ReLu, Pooling, Batch normalization and Quantize units.

- ➢ The Multiplication operation happens in parallel across different DRAM subarrays for data in transposed layout.

- ➢ Every DRAM bank is allocated to a layer in the Neural Network.

- ➢ P1-P5 represents the degree of parallelism while mapping DNN layers on DRAM banks.

# C-BRIC Artifacts: Chip Gallery #1
## *In-Memory Computing, Digital DNN Accelerators*

**SRAM IMC**



65nm: XNOR-SRAM
S. VLSI'18, JSSC'20



65nm: IMPULSE
SSCL'21



65nm: Dyn. Sparsity
SSCL'21

**IMC Macros**



28nm: Prog. IMC accelerator
S. VLSI'21

**RRAM IMC**



Winbond 90nm:
XNOR-RRAM, 2-bit-per-cell
Micro'19, TED'20, SSCL'20



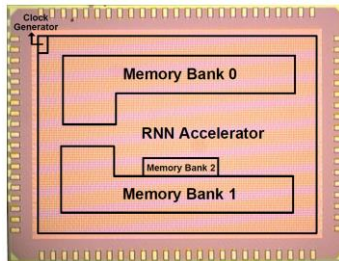TSMC 40nm: Binary RRAM,
multi-bit encoding
ISSCC'21, CICC'21, JSSC'22
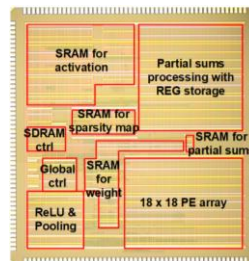


SUNY 65nm:
RRAM/SRAM
based hybrid IMC

**IMC Systems**



TSMC 40nm: RRAM IMC System
with embedded processor

**Digital Accelerators**



65nm: LSTM inference w/
hierarchical structured sparsity
ESSCIRC'19, JSSC'20



40nm: CNN inference w/
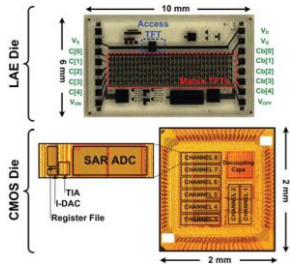conditional computing
CICC'20, JSSC'21



65nm: 16-bit fixed-point
CNN training accelerator
SSCL'20

*ASU*
*Purdue*
*Georgia Tech*
*(w/ ASCENT,*
*Samsung, TSMC)*

**C-BRIC**

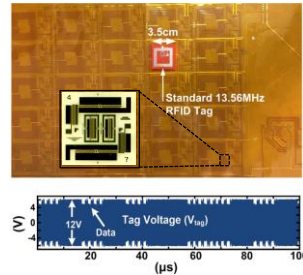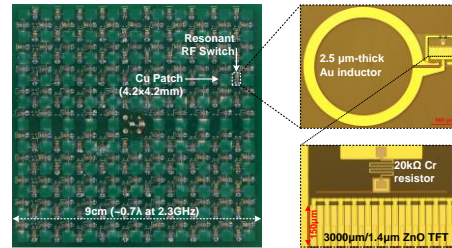# C-BRIC Artifacts: Chip Gallery #2
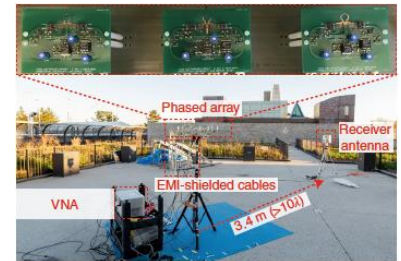## *Large-Area Sensing, RL, SNN, Optimization, Robotics*



**Hybrid 130nm CMOS / LAE tactile sensing array**
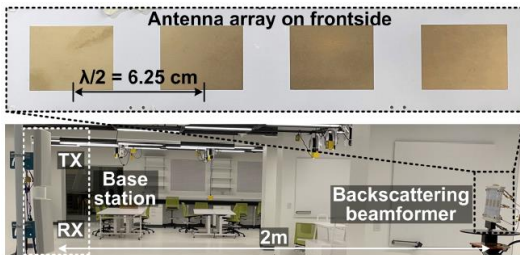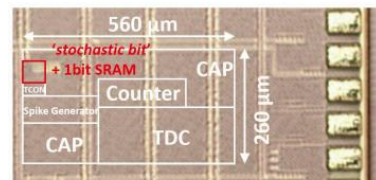**ISSCC'19**



**LAE 13.56 MHz RFID reader array**
**SSCL'18**



**LAE 2.4 GHz monolithically-integrable reconfigurable antenna**
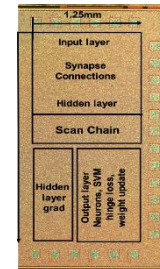**IEDM'20**



**LAE 1 GHz Phased Array Nature Electronics'21**
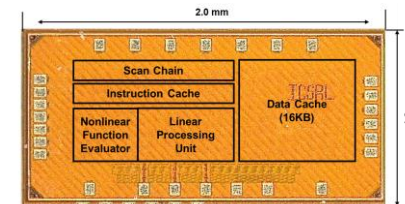


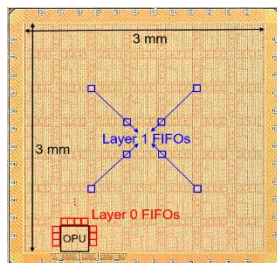**LAE 2.4 GHz passive backscattering beamformer for event-driven sensing**



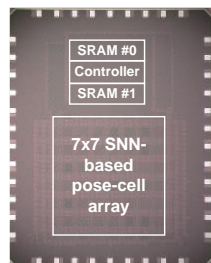**Stochastic binary SNN 90nm: TCAS-I'20**



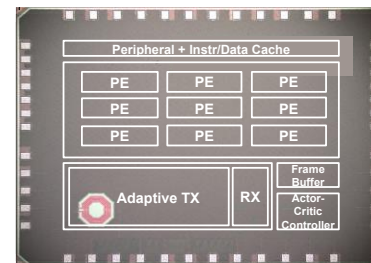**Tiny-RL, analog compute 130nm: JSSC'18**



**Edge AI for swarm robotics 130nm: ISSCC'18, JSSC'19**



**65nm: OPTIMO 49-core optimization processor CICC'19, JSSC'19**



**65nm: NeuroSLAM SNN-based visual SLAM acc. ISSCC'20, JSSC'21**



**65nm: Edge SoC with edge-cloud load balancing S' VLSI'20, JSSC (review)**

*Princeton*
*Georgia Tech*
*Purdue (w/ Intel)*

C-BRIC

# Parting Thoughts…

➢ While possibilities of achieving large improvements in inference latency and energy is possible…

➢ There are several challenges…
- Array efficiency
- Cross-bar non-idealities: Device non-linearity, access transistor/selector device, circuit non-idealities (line resistance, source/sink resistance), process variability
- Reliability and endurance of non-volatile devices
- High write cost for NVMs
- A/D and D/A converters
- Data movements from partial sums
- Training/mapping to the hardware – degradation over time for some NVM technologies
- Need for vector operations and floating point operations
- SRAMs are large compared to emerging NVMs

C-BRIC

# Neuro-electronics Research Laboratory

# Acknowledgement



**Vannevar Bush
Faculty Fellow**