

Al ASIC: Design and Practice (ADaP) Fall 2023

Processing-In-Memory

燕博南





- Motivation
- Basic Structure of PIM
- PIM/CIM Behavior Model

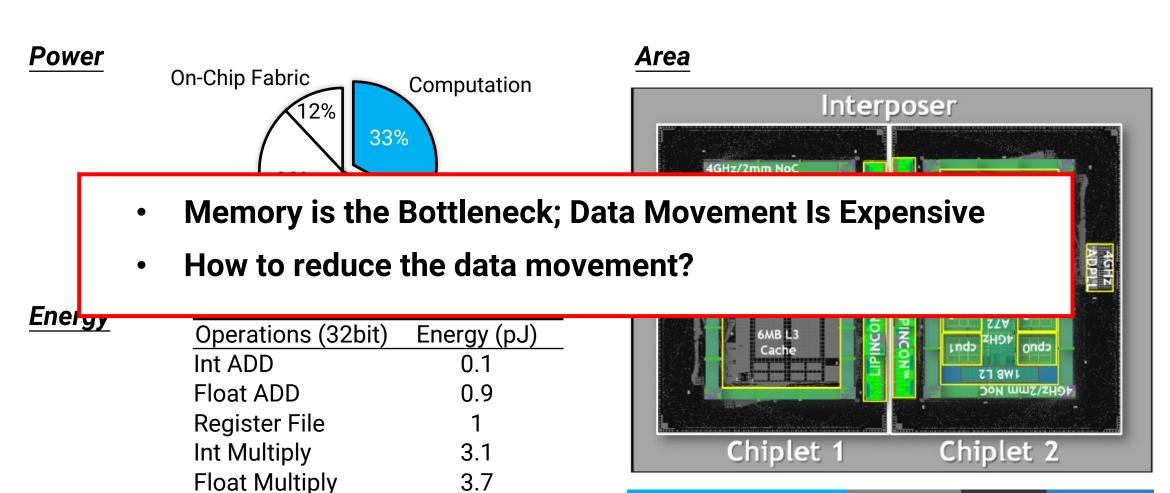


Why Fused Memory & Processing Units?



Why Does Processing-In-Memory (PIM) Rise?





640

On-Chip Mem

CPU

Chiplet

Interface

SRAM Cache

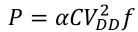
DRAM Memory

Bus &

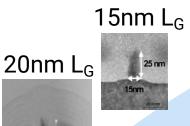


Specialized Hardware Enhances Efficiency

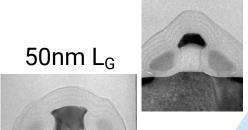




P: Power Dissipation α : Activity Factor C: Load Capacitance V_{DD} : Power Supply f: Clock Frequency



Multicore to Manycore



30nm L_G

MOSFET Scaling
Transistor Size &
Clock Frequency



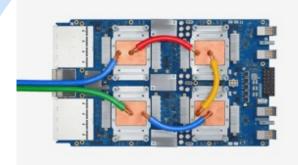




Domain-Specific





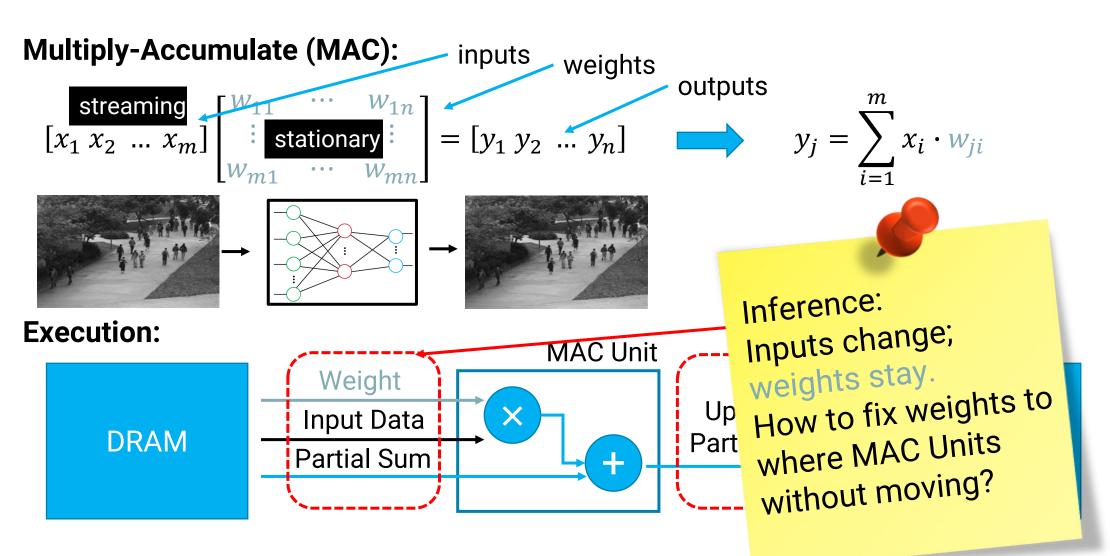


Google TPU v3
~200Watts
(4 TPU to run AlphaGo,
300x less power
consumption)



Uniqueness of Neural Network Execution

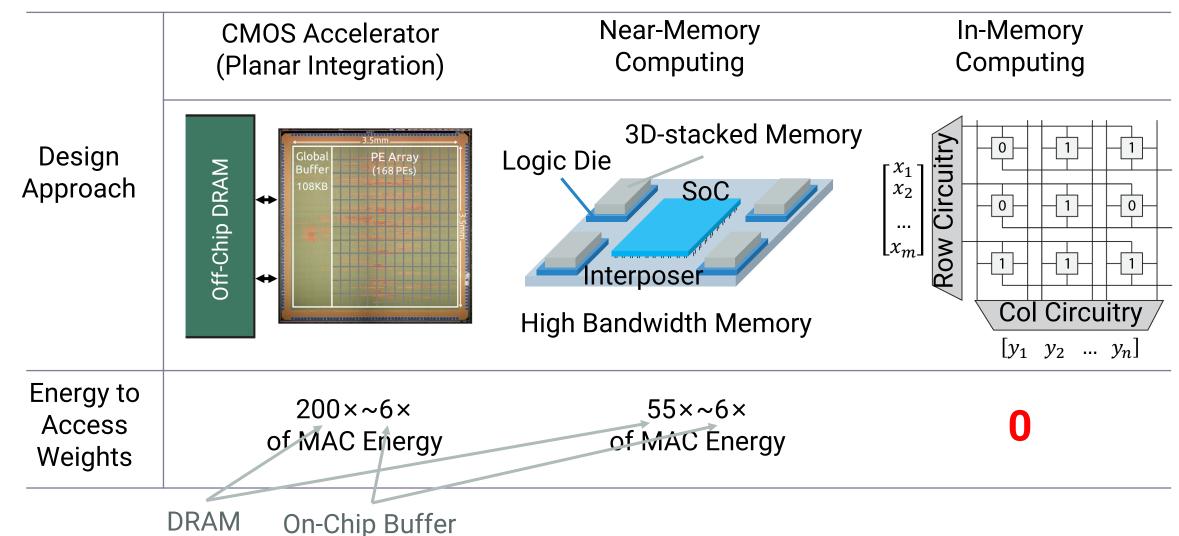






Make Memory Access Less Expensive



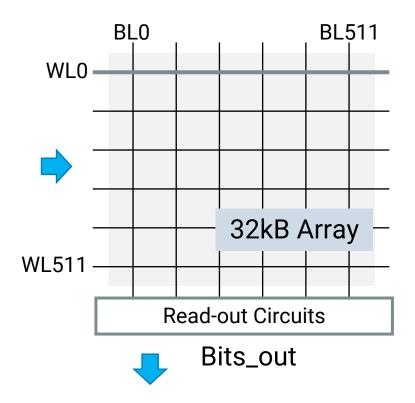


7

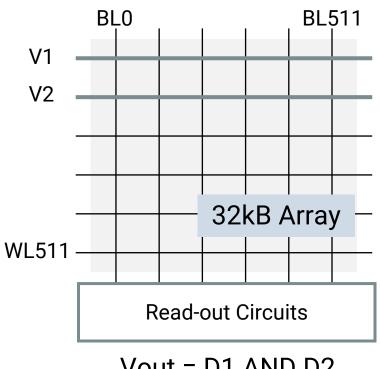
From Memory to In-Memory Computing







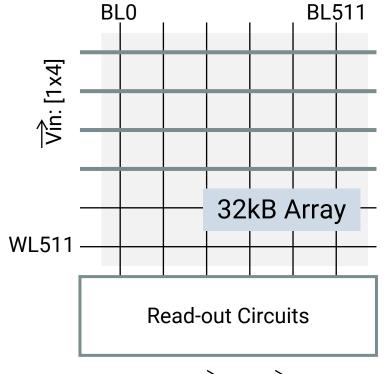
In-Memory Logic



Vout = D1 AND D2

Bits accessed: 1*512=512 Bits accessed: 2*512=1k

In-Memory General Matrix Multiplication (GEMM)



 $\overrightarrow{Vout} = \overrightarrow{Vin} \times M$

Bits accessed: 4*512=2k

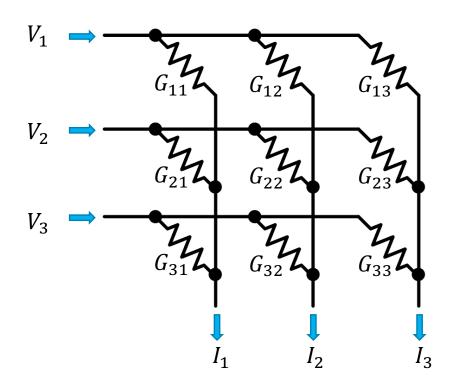


Basic Idea: Memristor PIM

Basic Idea: Resistivity Memory PIM



inputs weights
$$[V_1 \quad V_2 \quad V_3] \begin{bmatrix} G_{11} & G_{12} & G_{13} \\ G_{21} & G_{22} & G_{23} \\ G_{31} & G_{32} & G_{33} \end{bmatrix} = \begin{bmatrix} I_1 & I_2 & I_3 \end{bmatrix}$$



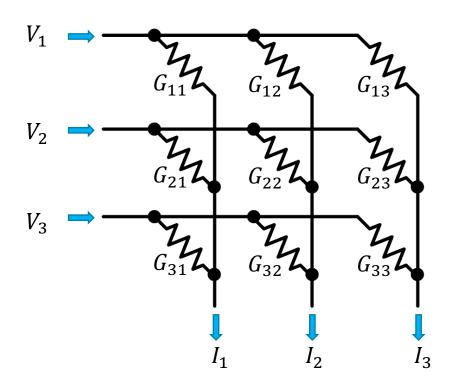
outputs

- Weight Matrix Stored as Conductance G
- Rely on Analog Computation (Kirchhoff's Current Law) for "almost free"
 - Multiplication: $I = V \cdot G$
 - Addition: $I^{column} = I_1^{row} + I_2^{row} + I_2^{row}$
- Ideal Nanoscale Devices for G:
 - Programmable Conductance
 - Multi-Level Cell
 - Small Footprint/High Density
 - Compatible with Existing CMOS Process

Basic Idea: Resistivity Memory PIM



$$[V_1 \quad V_2 \quad V_3] \begin{bmatrix} G_{11} & G_{12} & G_{13} \\ G_{21} & G_{22} & G_{23} \\ G_{31} & G_{32} & G_{33} \end{bmatrix} = \begin{bmatrix} I_1 & I_2 & I_3 \end{bmatrix}$$
 outputs



Q: Can we use currents as the inputs?

$$\begin{bmatrix} I_1 & I_2 & I_3 \end{bmatrix} \begin{bmatrix} R_{11} & R_{12} & R_{13} \\ R_{21} & R_{22} & R_{23} \\ R_{31} & R_{32} & R_{33} \end{bmatrix} = \begin{bmatrix} V_1 & V_2 & V_3 \end{bmatrix}$$

Q: What circuits do we need? Why?



How A SoC Designer Sees PIM/CIM?

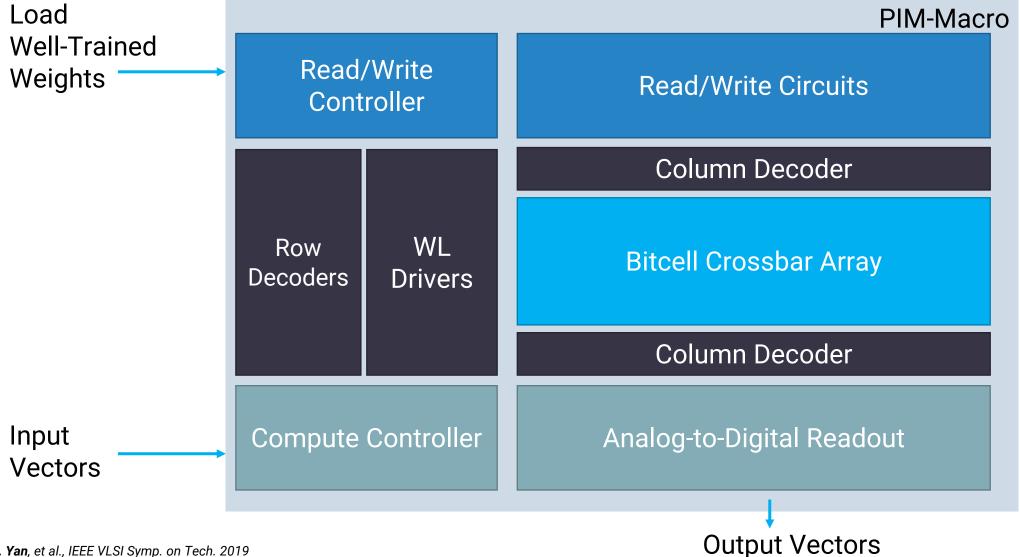


CIM Behavior Model



In-Memory Computing Solution Stack-Macro





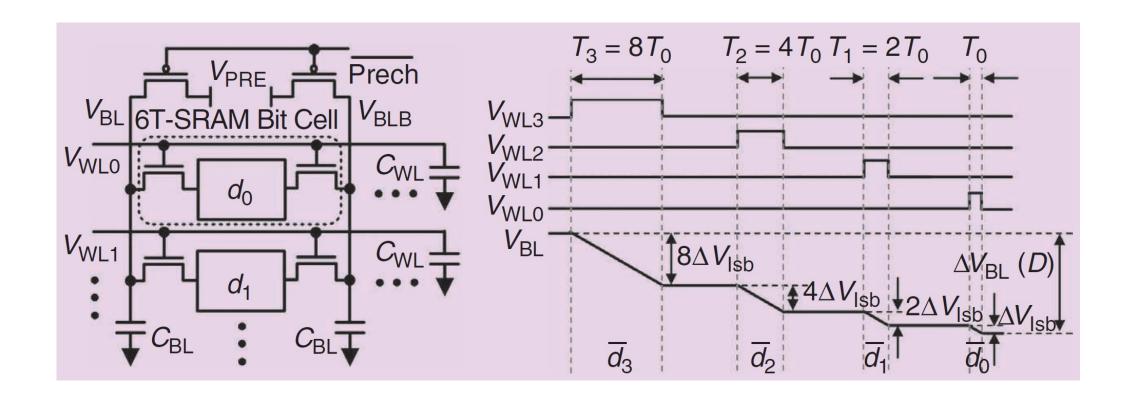


Basic Idea: SRAM PIM



Basic Idea: SRAM PIM



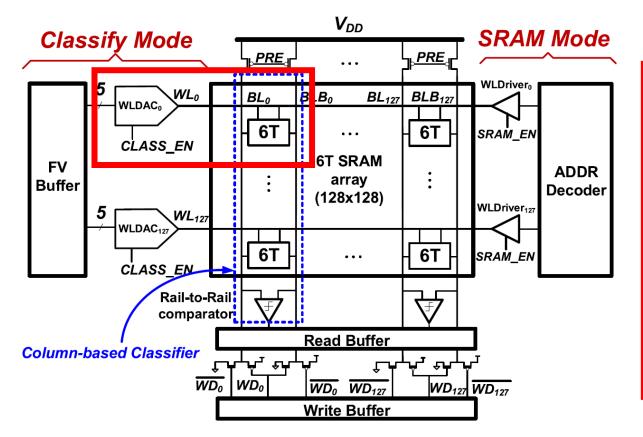


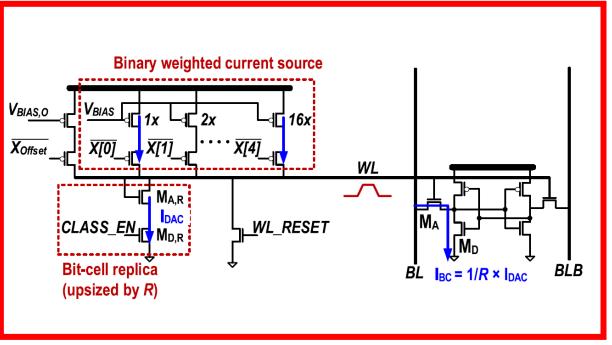
- Precharge BL
- · Control the opening time windows of each WL



Basic idea: SRAM PIM



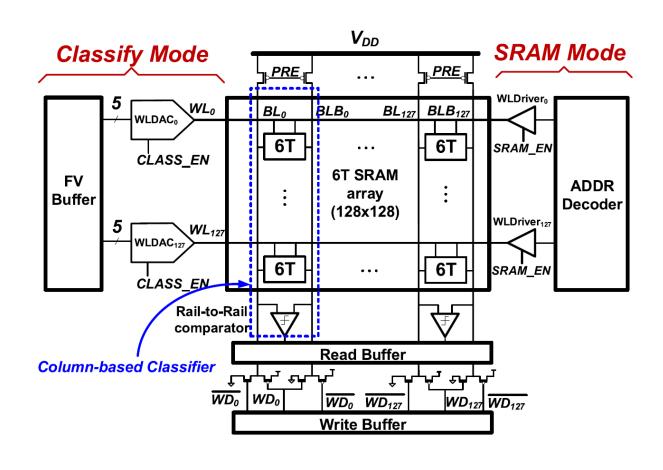


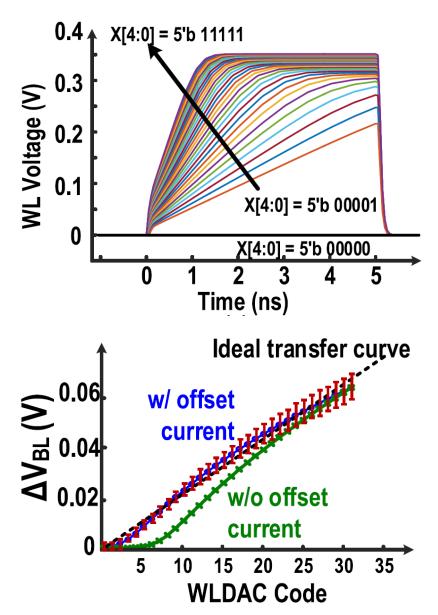




Basic idea: SRAM PIM

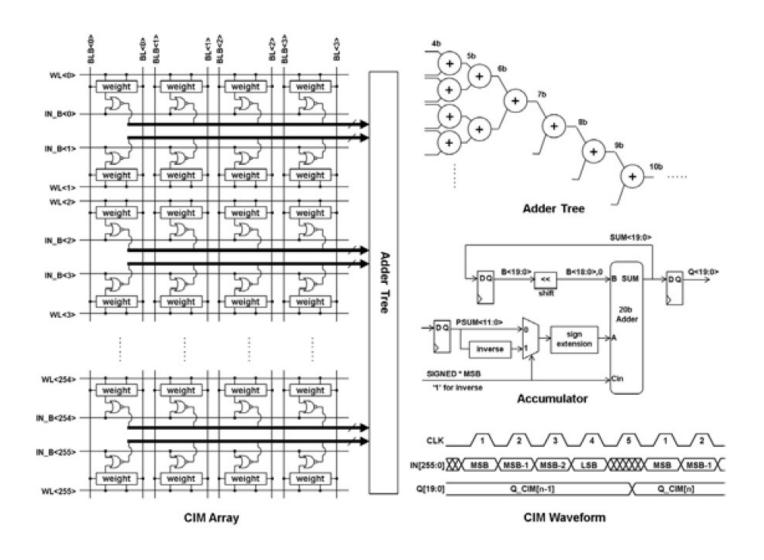






Basic Idea: Digital SRAM PIM





Why Digital PIM?

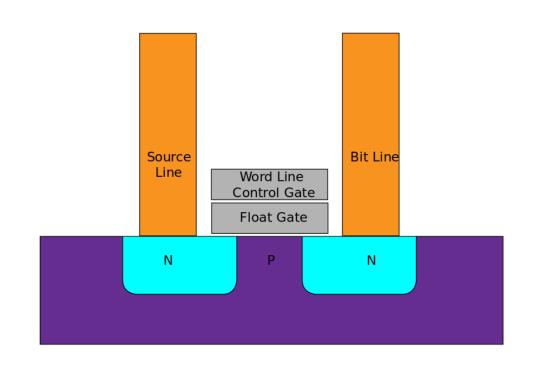
- Strength: High Noise Margin
- Weakness: Low parallelism

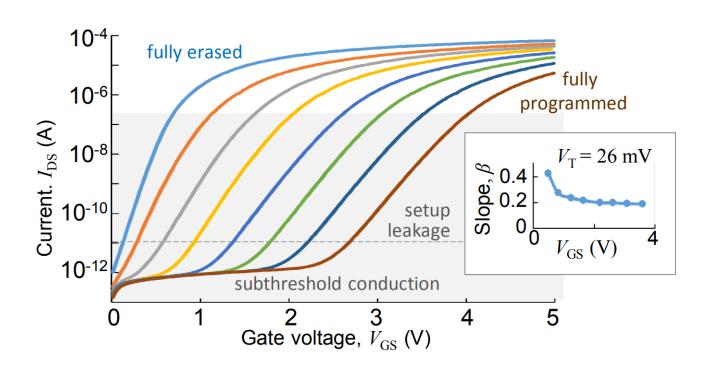


Basic Idea: Flash PIM







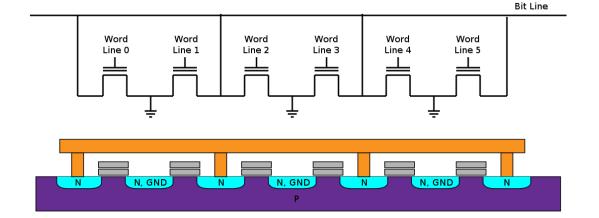




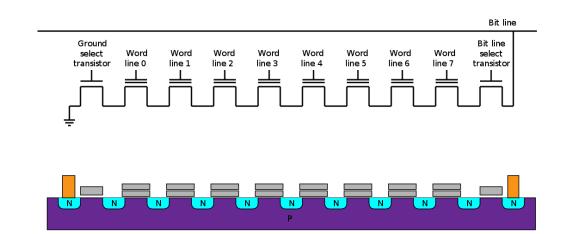
Basic Idea: Flash PIM: 2 Structures





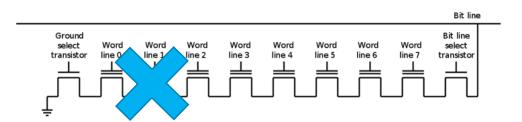


NAND Flash



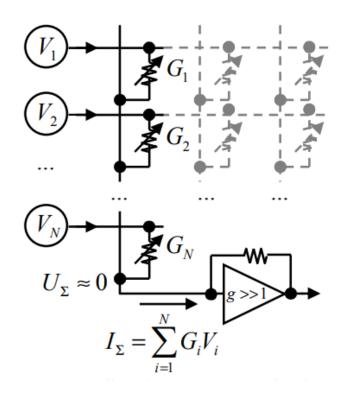
Q: Which one can be used for PIM-VMM?

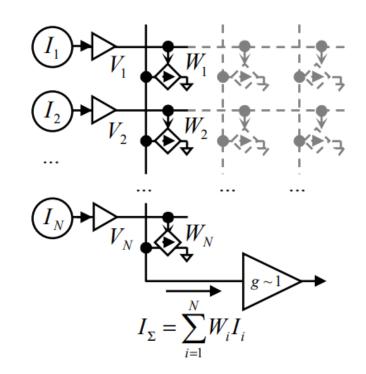
A: NOR Flash



Basic Idea: Flash PIM







Memristor PIM

Flash PIM



Memory Technologies

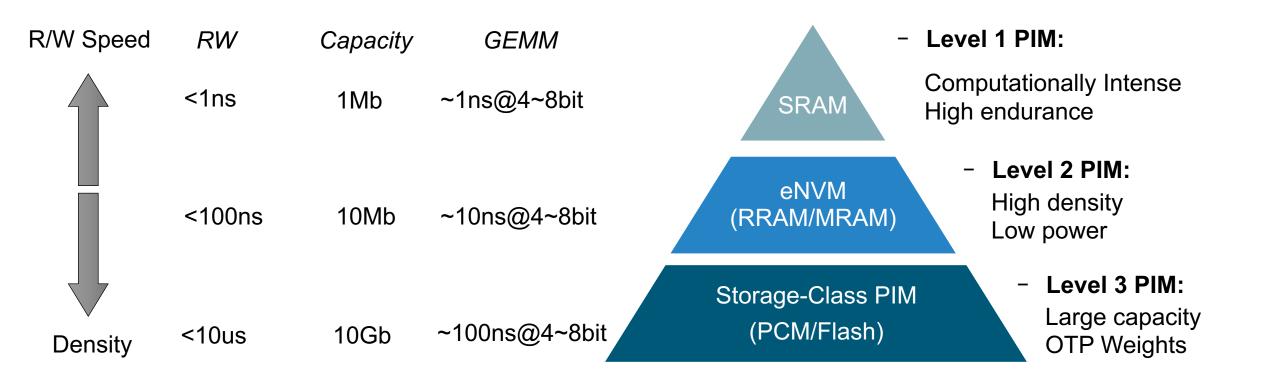


Property	РСМ	RRAM	MRAM	CMOS-SRAM
Multilevel Cell	Yes	Which is better?	No	No
Storage Density	High		High	Low
R _{OFF} /R _{ON}	High		Low	High
Nonvolatility	Yes		Yes	No
Leakage	Low		Low	High
Cell Area	16F ²		30~80F ²	160F ² (6T) 231F ² (8T)
Write Energy	6nJ	2nJ	<1nJ	<0.1nJ
Write Latency	150ns	100ns	10ns	<1ns
Endurance	10 ⁷ cycles	10 ⁶ cycles	10 ¹⁵ cycles	10 ¹⁶ cycles



Future PIM Memory Hierarchy







Place PIM into System



In-Memory Computing Solution Stack-Macro

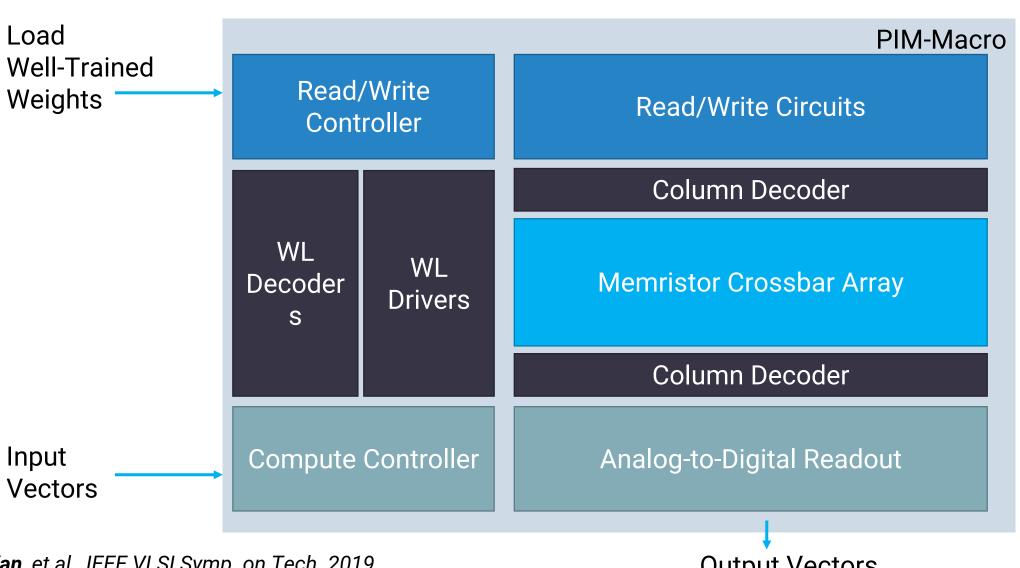


SW

Arch

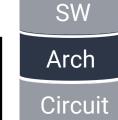
Circuit

Device

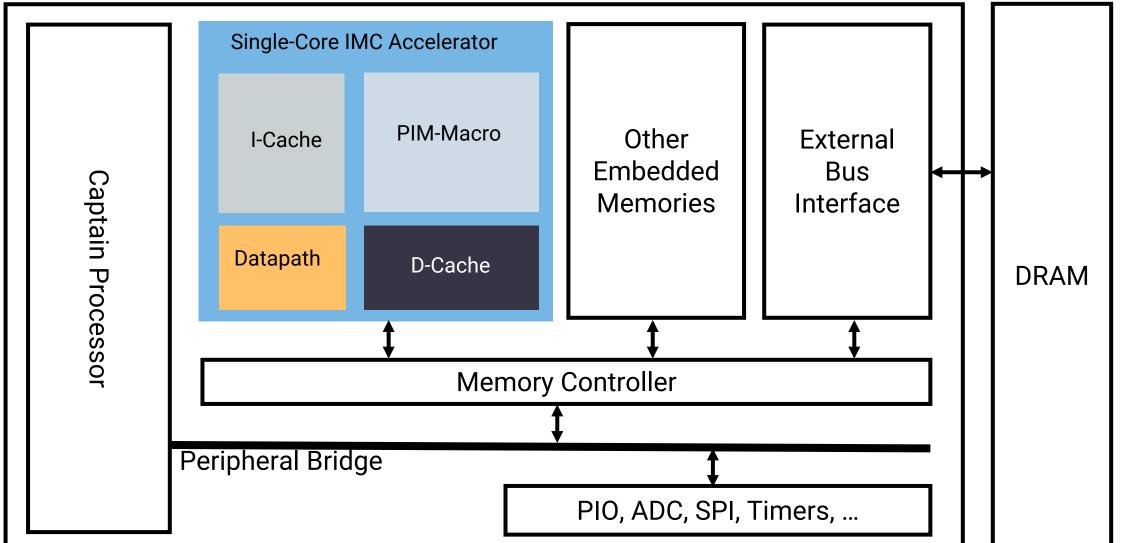


In-Memory Computing Solution Stack-SoC





Device





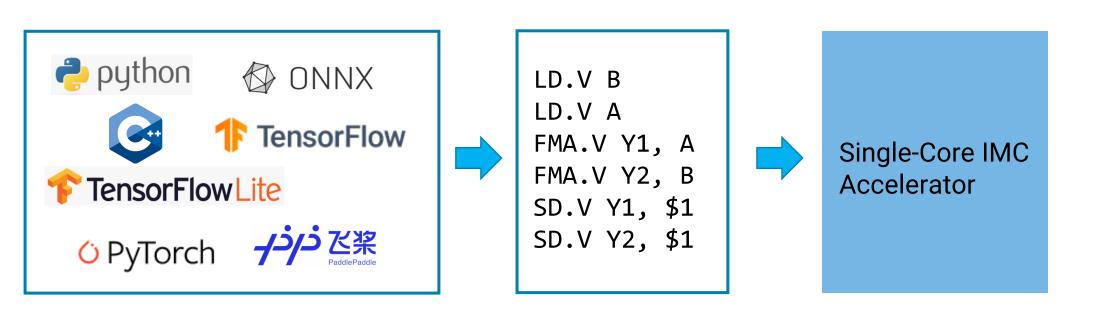
In-Memory Computing Solution Stack-Software



Arch Circuit

SW

Device



Hardware-Software codesign is required for best use of ASICs:

- Various ASICs need specialized compilers
- Various ML Models need specially optimized scheduling/dataflow



Not Only PIM



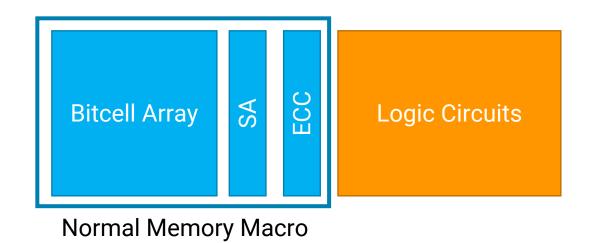
Compromise of PIM: Processing-Near Memory



Challenges:

A/D converter based compute readout circuits are

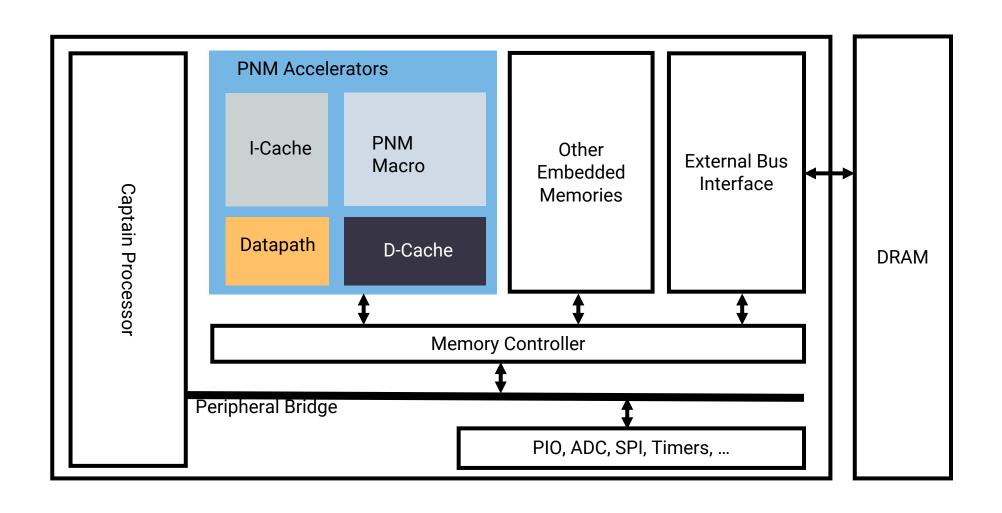
- sooooo hard to design
- sooooo large
- sooooo prune to noise/process variations



Processing-Near Memory Merging into Systems



At what level does PNM shares behaves the same as PIM?





Different Level of Memory-Compute Fusion

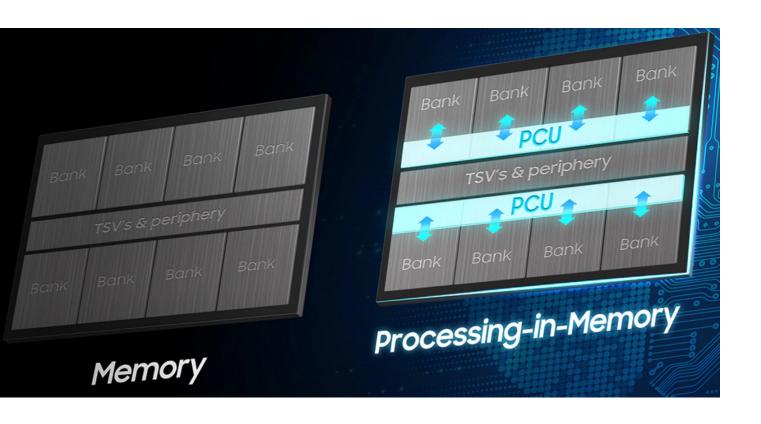


- In-Memory Computing
 - Memristor PIM
 - SRAM PIM
 - Flash PIM
- Near-Memory Computing
 - HBM-PIM
 - Embedded MRAM+SA+Computation
- In-Storage Computing / Computational Storage
 - Solid-State Drive (SSD) + ARM CPU = In-Storage Computing Card



Examples: HBM-PIM ... PNM





Near-Memory Computing

3D-stacked Memory
Logic Die
SoC
Interposer
High Bandwidth Memory

For a system designer, the processing is actually happening in the "memory" chips



Examples: In-Storage Computing



SmartSSD® CSD





- SSD + FPGA
- SSD + CPU
- Major purpose: storage



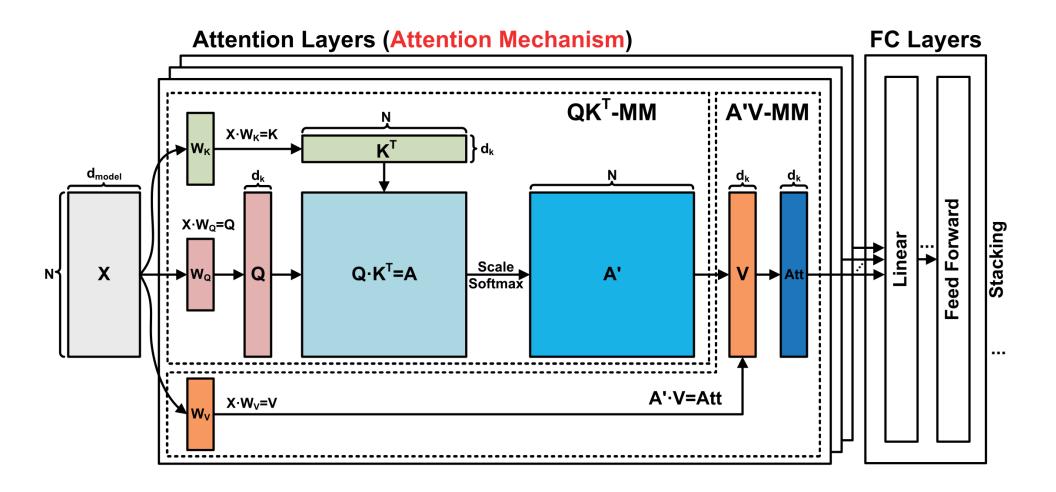


• PISLIB: PIM macro model



Accelerator Computing Task





Courtesy: TranCim Paper