

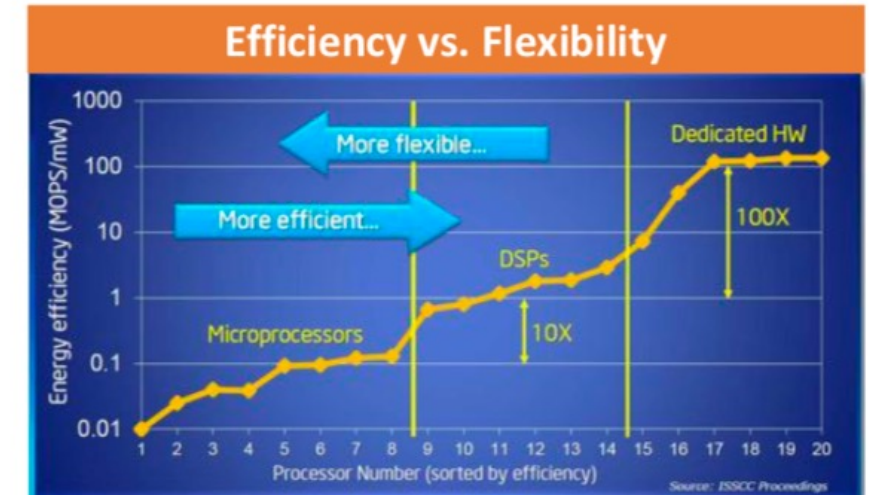
ReRAM-Based Processing-In-Memory(PIM) Accelerators for AI Applications

Siling Yang

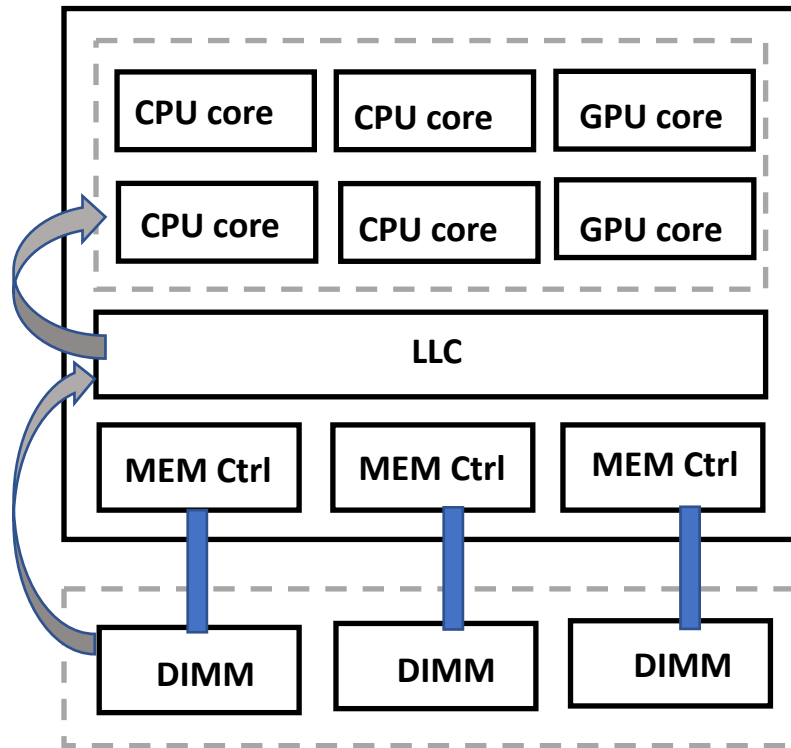
2021.04.17

Existing hardware accelerators for deep learning

- **GPUs**
 - Fast, but high power consumption (~200W)
 - Training DNNs in back-end GPU clusters
- **FPGAs**
 - Massively parallel + low-power (~25W) + reconfigurable
 - Suitable for latency-sensitive real-time inference job
- **ASICs**
 - Fast + energy efficient
 - Long development cycle
- **Novel architectures and emerging devices**



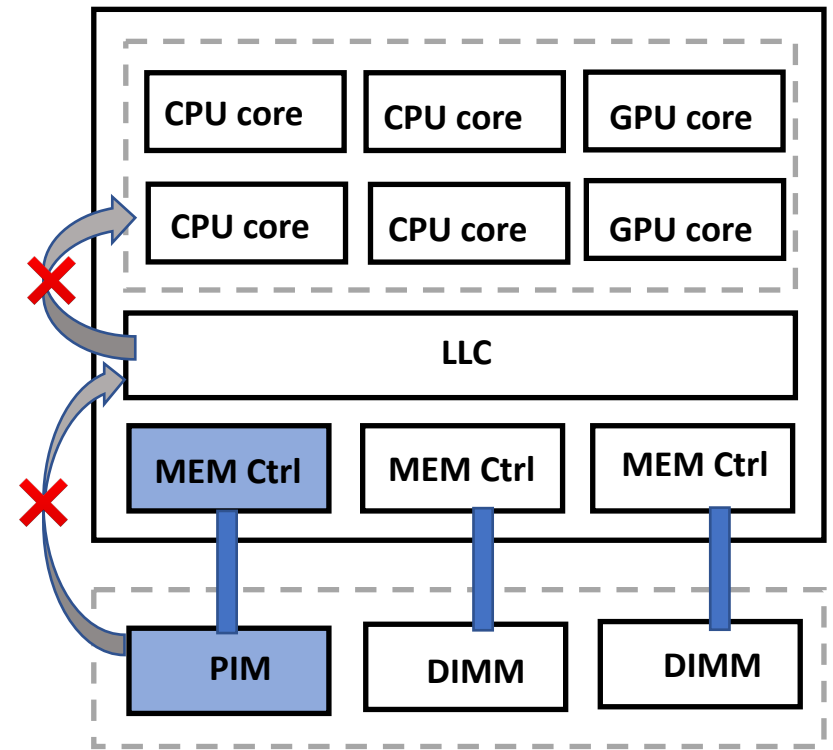
Why PIM?



Von Neumann architecture

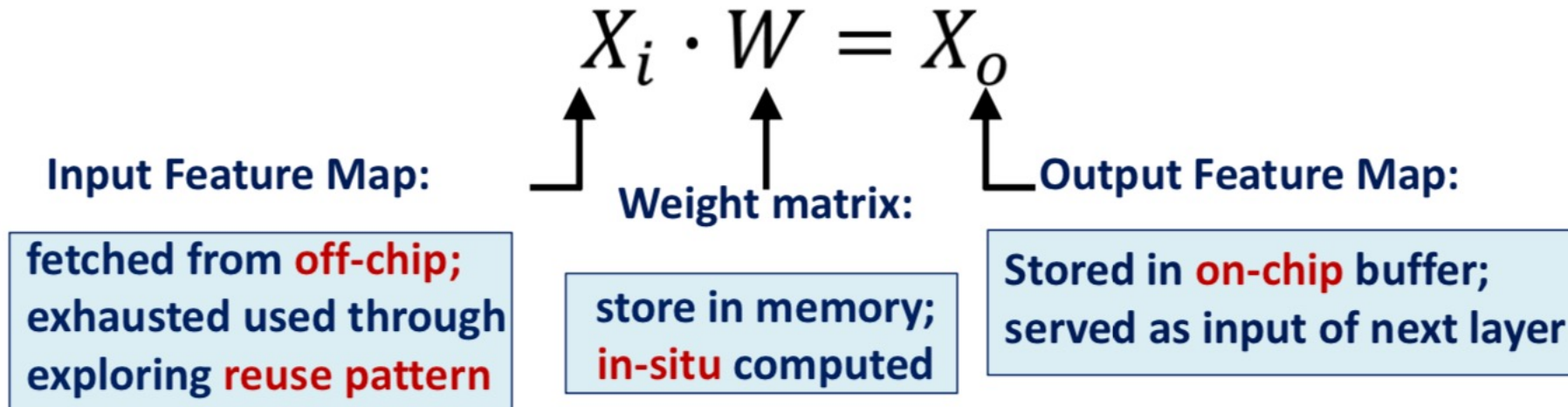
Operation	Energy(pJ)
16b Add	0.05
32b Add	0.1
16b FP Add	0.4
32b FP Add	0.9
32b Mult	3.1
16b FP Mult	1.1
32b FP Mult	3.7
32b SRAM Read(8KB)	5
32b DRAM Read	640

Memory Wall

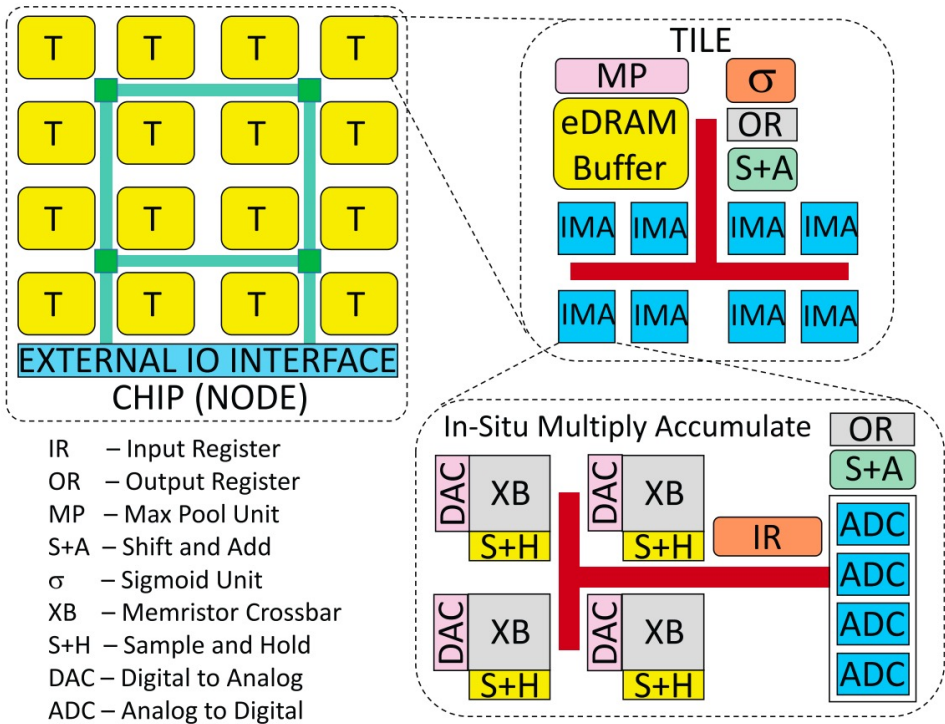


Processing in Memory

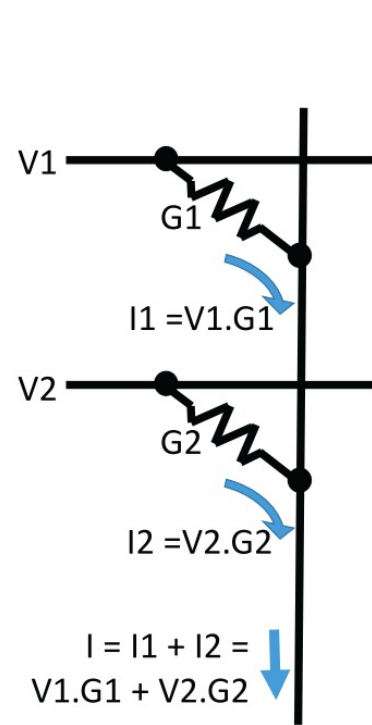
Why PIM benefits NN?



Analog Dot Product in crossbars

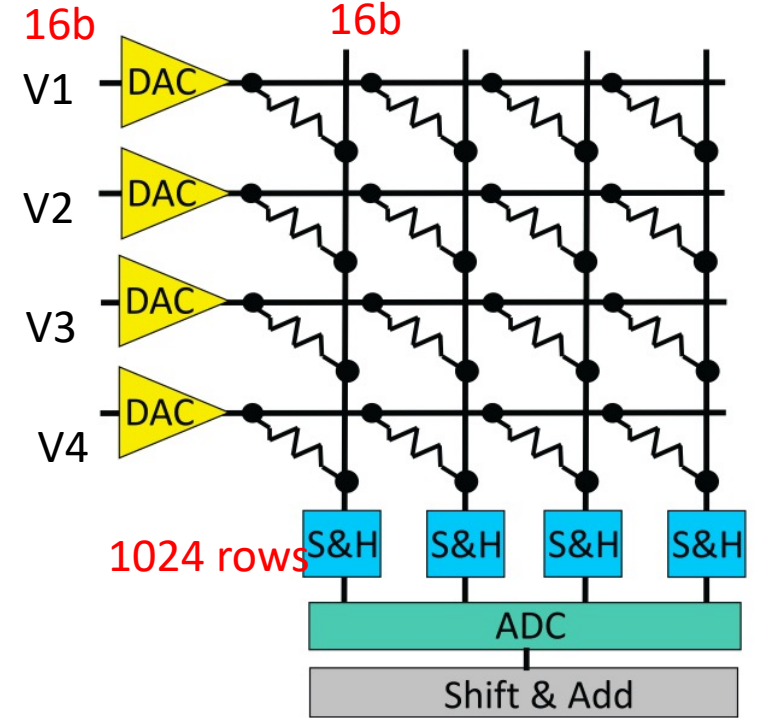
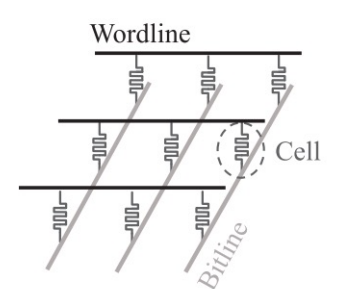
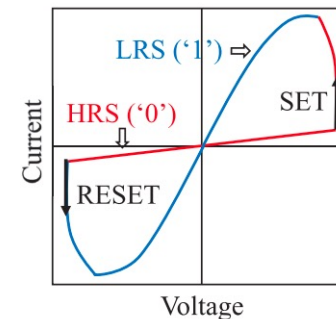
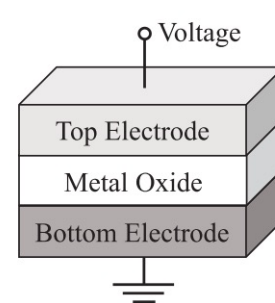


ISAAC architecture hierarchy. [ISCA16]



(a) Multiply-Accumulate operation

Word line
Bit line



(b) Vector-Matrix Multiplier

Mapping Filter Weights of DNNs in crossbar-based Accelerators

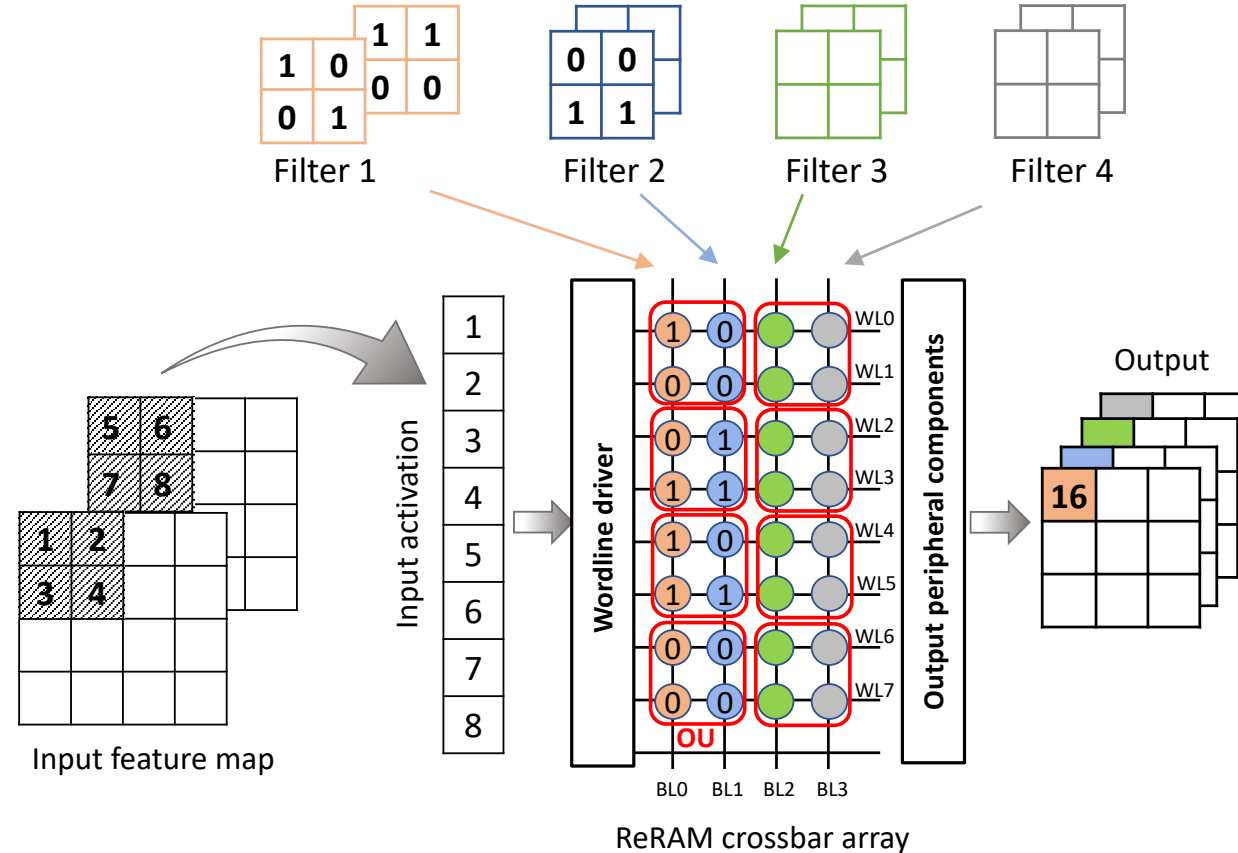


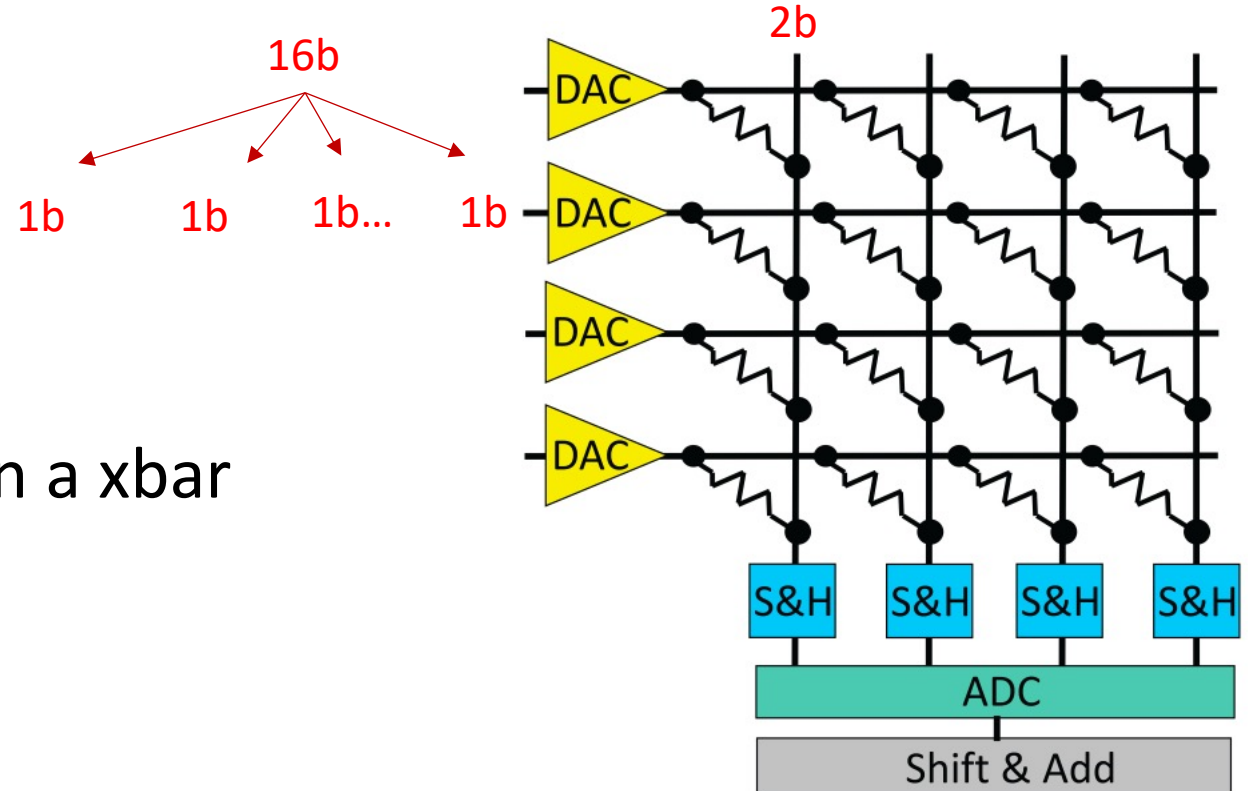
Illustration of mapping filter weights to a crossbar array used in the architecture of ReRAM-based accelerators. BL: bitline; WL: wordline; OU: operation unit.

Design in ISAAC: simplifying the ADC

[Never drop bits]

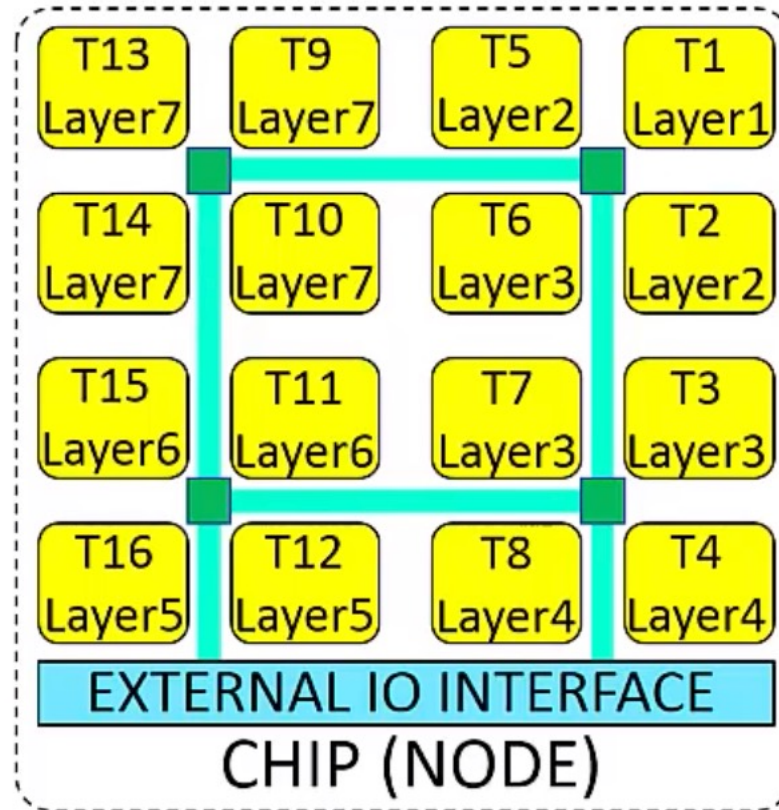
Distribute the computation:

- Across time
- Across multiple columns in a xbar
- Across multiple xbars
- Use an encoding trick

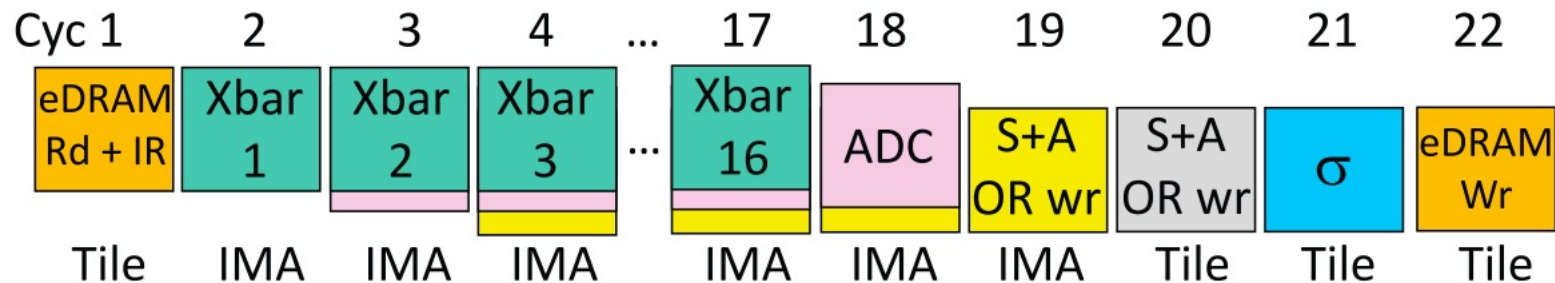
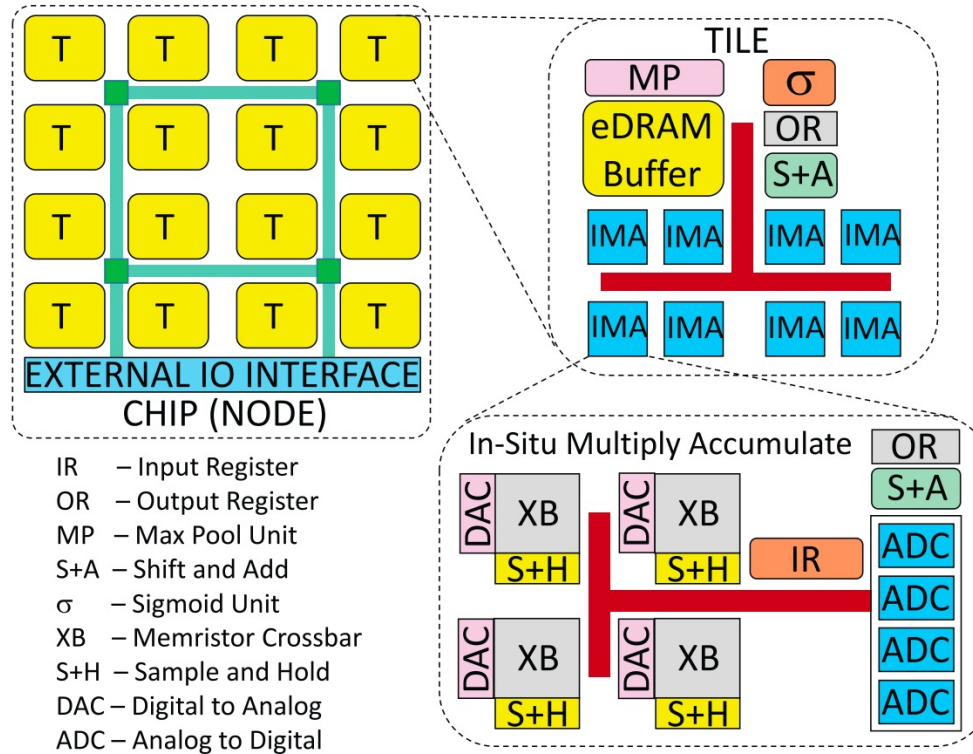


(b) Vector-Matrix Multiplier

ISAAC



ISAAC: Pipeline



ISAAC Design space

Balanced use of xbars and ADCs;

Peak Thruput: 45 TOPs/s;

Storage: 63MB

Optimize:

- **CE: Computational Efficiency** is represented by the number of 16-bit operations performed per second per mm² (GOPS/s \times mm²).
- **PE: Power Efficiency** is represented by the number of 16-bit operations performed per watt (GOPS/W).
- **SE: Storage Efficiency** is the on-chip capacity for synaptic weights per unit area (MB/ mm²).



EIE(ISCA'16):

- Peak Thruput: 3TOPs/s
- Storage: 80MB (8MB)

Many large xbars; few ADCs;

Peak Thruput: 10 TOPs/s;

Storage: 1.7 GB



Research on ReRAM-based PIM

ReRAM-based ANN Architectures:

- Mapping NN to ReRAM
- Architectures for NN inference
- Architectures for NN training
- **MCA-aware pruning strategy**
- Reconfigurable architectures
- Reducing overhead of analog implementation

ReRAM-based PIM Techniques

- Arithmetic and logical operations
- Data search operations
- Graph-processing operations
- Approximate computing approaches

Thanks