

6.5930/1

Hardware Architectures for Deep Learning

Advanced Technologies

March 9, 2026

Joel Emer and Vivienne Sze

Massachusetts Institute of Technology
Electrical Engineering & Computer Science



Advanced Storage Technology

- **Bring compute and memory closer together** to reduced the cost of data movement
- **Processing/Compute Near Memory** (*a.k.a. near-data processing*)
 - Embedded DRAM (eDRAM)
 - Increase on-chip storage capacity
 - 3D Stacked DRAM (e.g., Hybrid Memory Cube (HMC), High Bandwidth Memory (HBM))
 - Increase memory bandwidth
- **Processing/Compute In Memory** (*a.k.a. in-memory computing*)
 - Static Random Access Memories (SRAM), Dynamic Random Access Memories (DRAM), and Non-Volatile Memories (NVM)
 - Processing ***integrated into*** memory (more aggressive form of processing near memory) versus processing ***using*** memory (using memory device to perform compute)

SRAM (kB – MB)

- SRAM accounts for majority of on-chip storage
 - SRAM is significant portion of chip area
- Usually, 6 transistors per bit-cell

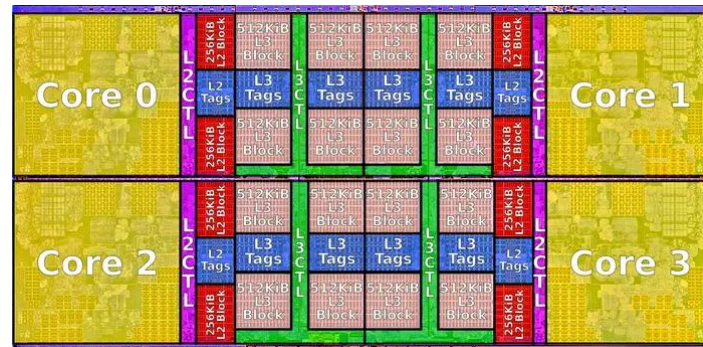
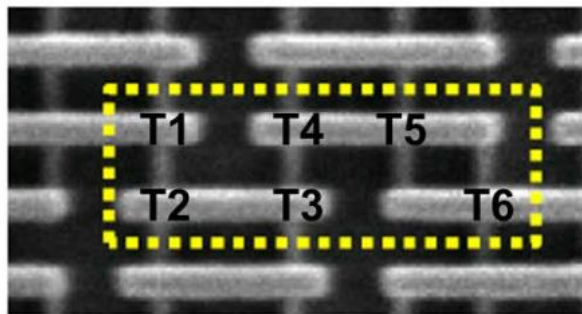
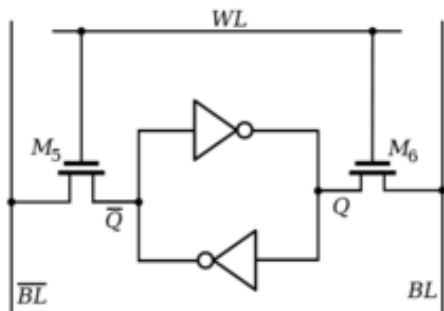


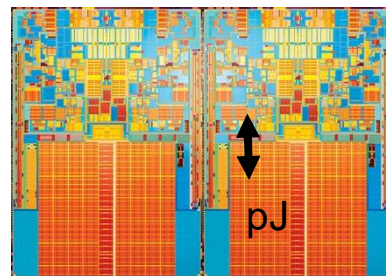
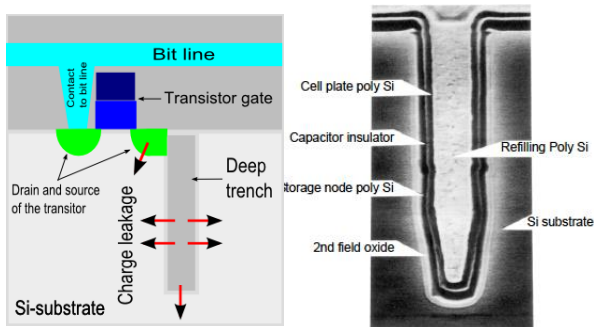
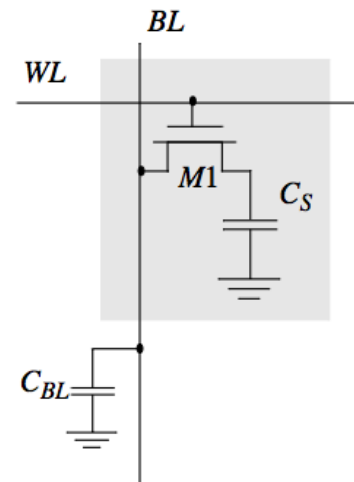
Image Source: <https://en.wikichip.org/wiki/amd/microarchitectures/zen>



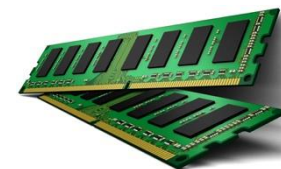
Bit cell size
0.75 μm^2 in 14nm

DRAM (GB)

- Higher density than SRAM
 - 1 transistor per bit-cell
 - Needs periodic refresh
- Special device process
 - Usually off-chip (except eDRAM – which is pricey!)
 - Off-chip interconnect has much higher capacitance

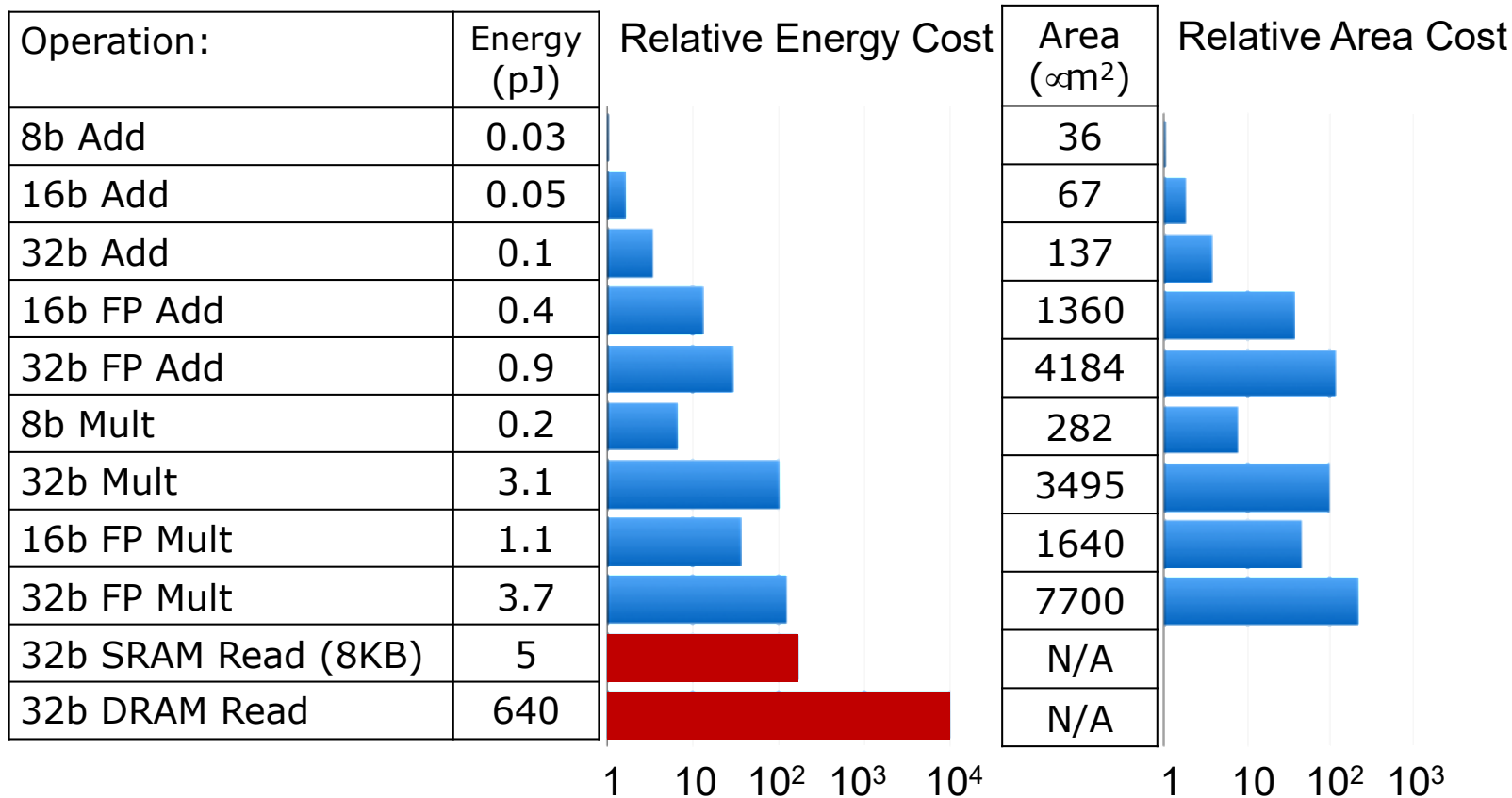


\longleftrightarrow
nJ



Size and Emer

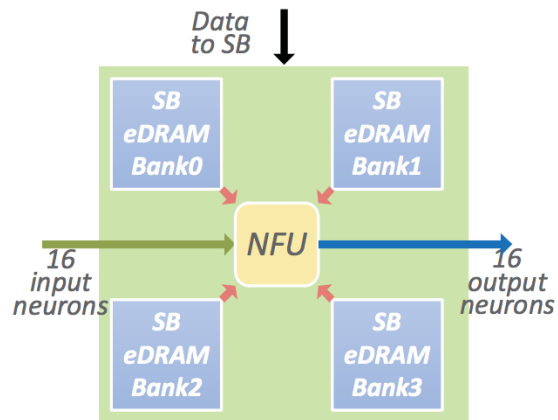
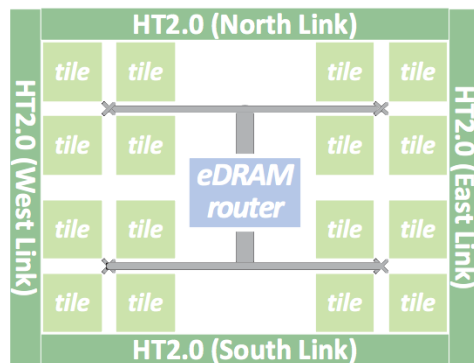
Memory Access Cost for DRAM is Expensive



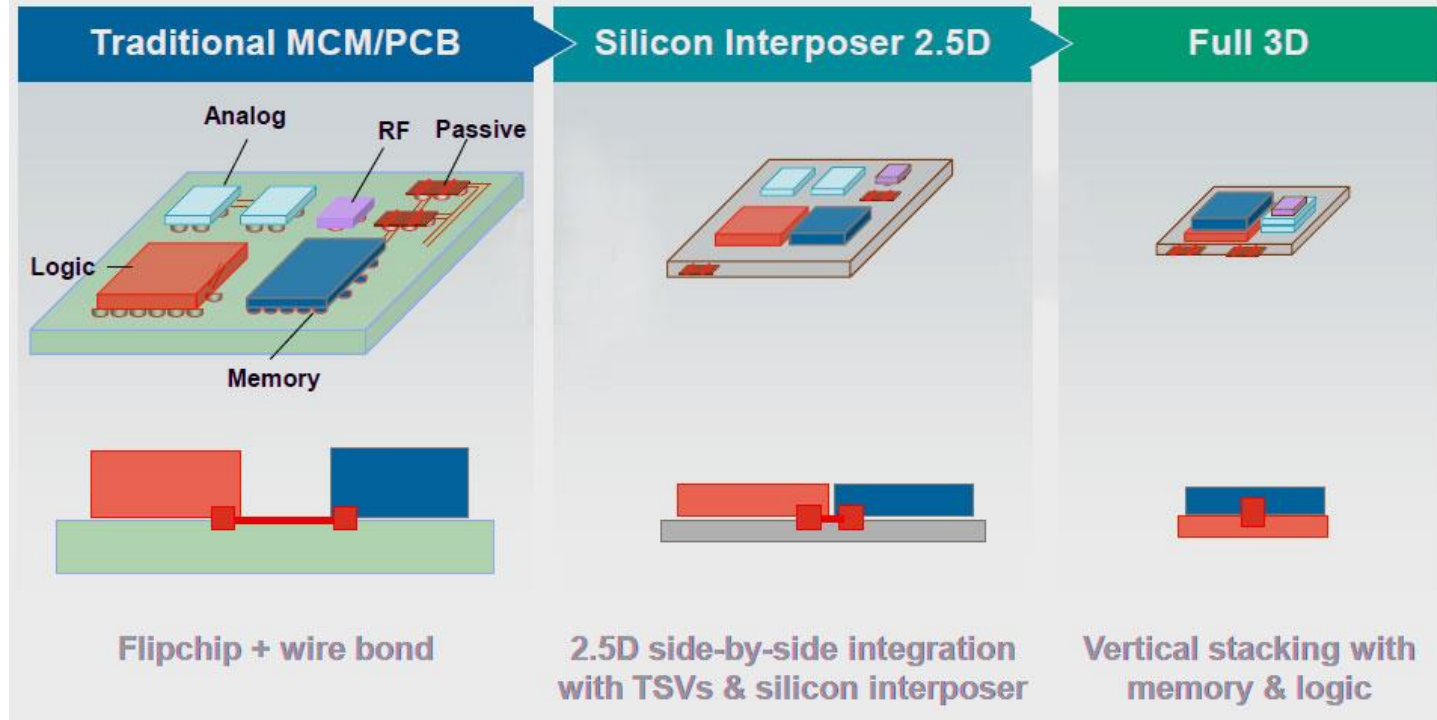
eDRAM (DaDianNao)

- Advantages of eDRAM
 - 2.85x higher density than SRAM
 - 321x more energy-efficient than DRAM (DDR3)
- Store weights in eDRAM (36MB)
 - Target fully connected layers since dominated by weights

16
Parallel
Tiles



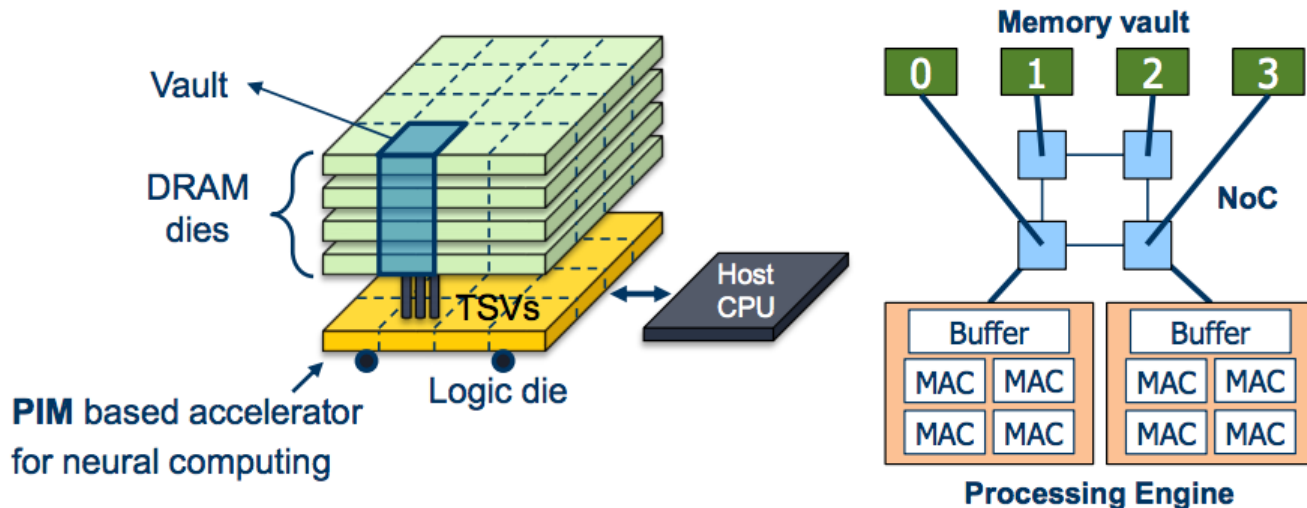
Going 3D!



Source: <https://www.einfochips.com/blog/2-5d-3d-ics-new-paradigms-in-asic/>

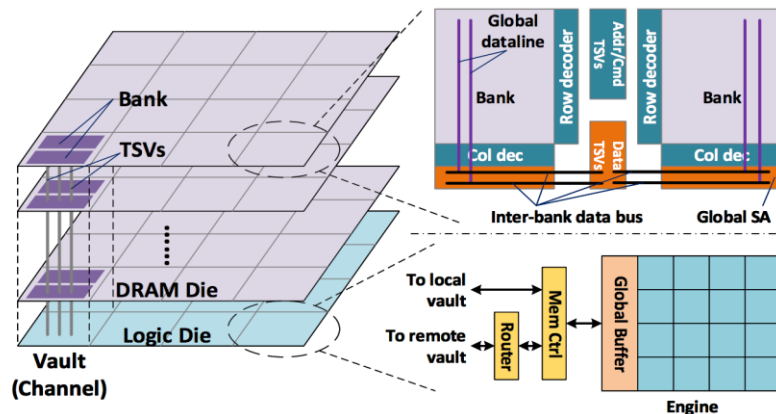
Stacked DRAM (NeuroCube)

- NeuroCube on Hybrid Memory Cube Logic Die
 - 6.25x higher BW than DDR3
 - HMC (16 ch x 10GB/s) > DDR3 BW (2 ch x 12.8GB/s)
 - Computation closer to memory (reduce energy)



Stacked DRAM (Tetris)

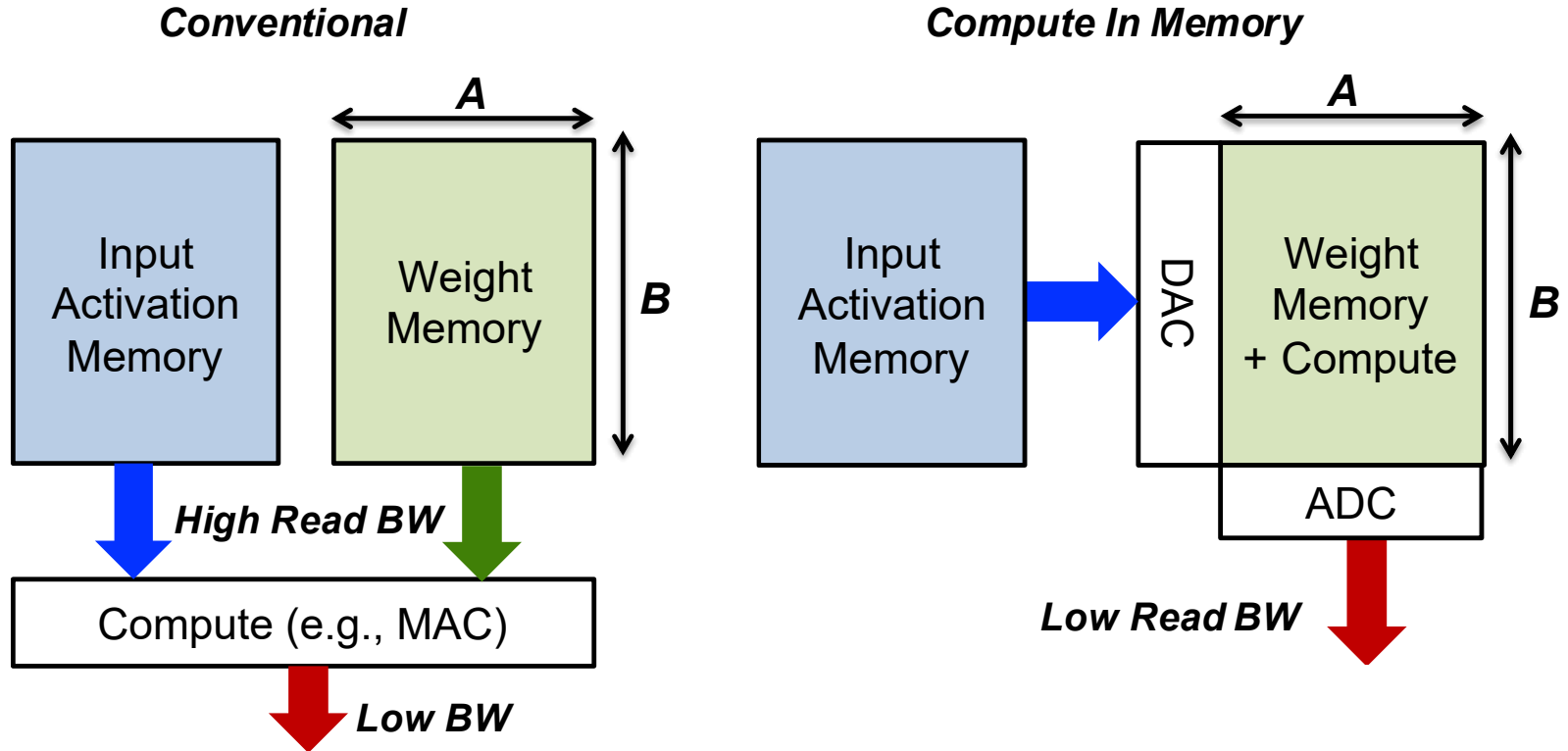
- Explores the use of HMC with the Eyeriss **spatial architecture** and **row stationary dataflow**
- Allocates **more area to the computation (PE array)** than on-chip memory (**global buffer**) to exploit the low energy and high throughput properties of the **HMC**
 - 1.5x energy reduction, 4.1x higher throughput vs. 2-D DRAM



[Gao, ASPLOS 2017]

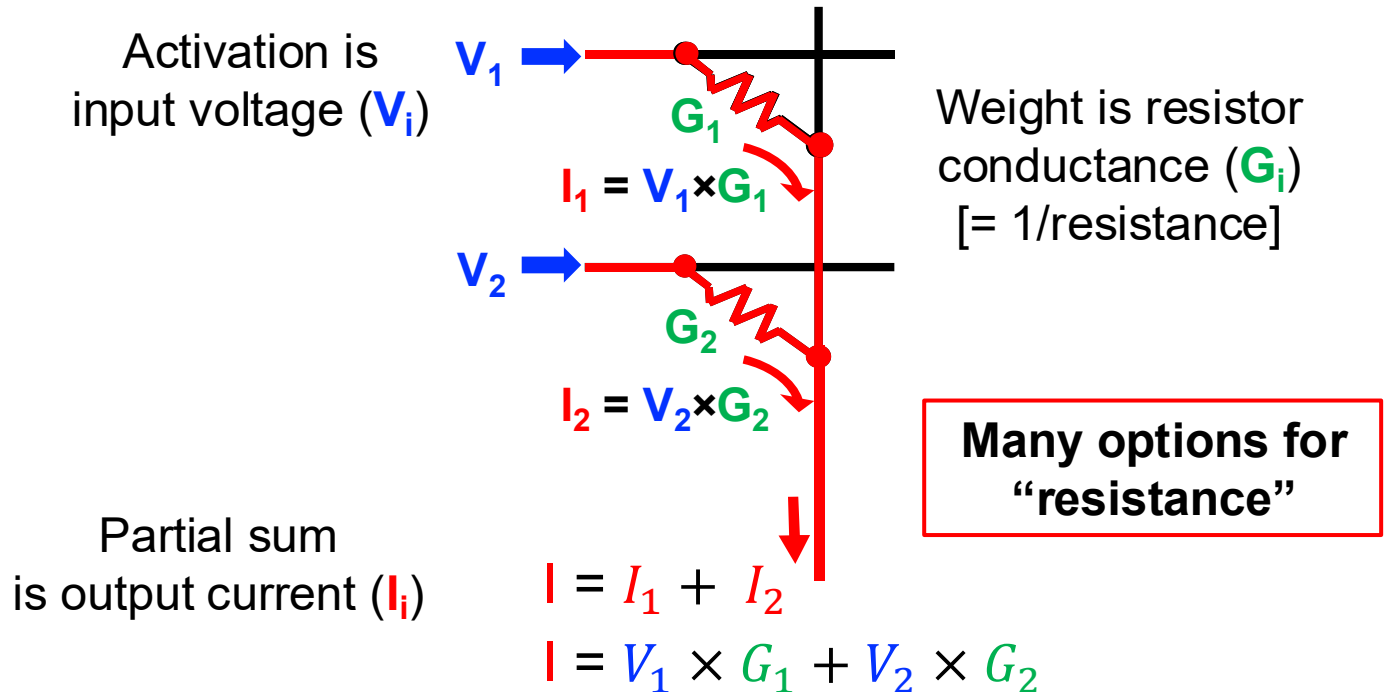
*Eyeriss
design*

Conventional Processing vs. CiM



Analog Computing

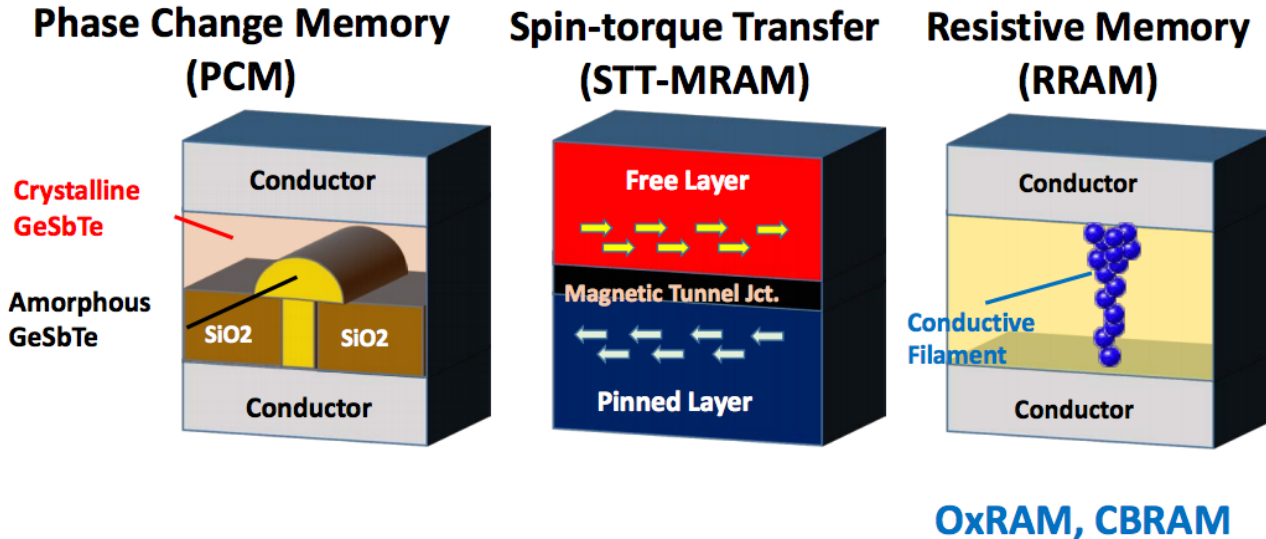
Analog computing is typically required to bring the computation into the array of storage elements or into its peripheral circuits



CiM Using Memristors in Non-Volatile Memory

Use memristors as programmable weights (resistance)

Candidate memristor devices



Source: Darsen Lu, MICRO-49 Tutorial on Emerging Memory

Candidate Devices for Memristor

		PCM	STT-MRAM	RRAM	NAND Flash	DRAM
Power	Ewrite/ Bit	18pJ [15] ¹	1.0pJ [20]	1.0pJ [20]	100pJ [31]	<1.0 pJ [9]
	lwrite	100μA [15]	50μA [33]	1.0μA [18]		
Performance²	Write Lat.	150ns [15]	5ns [6]	50ns [20]	>100μs	5ns [27]
	Read Lat.	80ns [28]	10ns [5]	<10ns [30]	15-50μs [35]	20-80ns [35]
Reliability	Program Window	3 bit/cell [12]	Good [3]	Variable [23]	4 bit/cell [32]	Good
	Endurance	10 ⁸ -10 ⁹ [4,9]	Unlimited [1]	10 ⁵ -10 ¹⁰ [1]	10 ⁵ -10 ⁶ [8]	Unlimited
	Retention	R-drift[12]	Good [6]	RTN [24]	Good	64ms
Density	Cell size	4 F ² [15]	12 F ² [33]	4-6F ² [21]	<4 F ² [15]	7 F ² [15]

Emerging Memory Technologies

(Speed of DRAM; Non-volatility of NAND)

1. Estimated using $I_{\text{reset}} * V_{\text{dd}} * t_{\text{write}}$
2. Required programming pulse duration

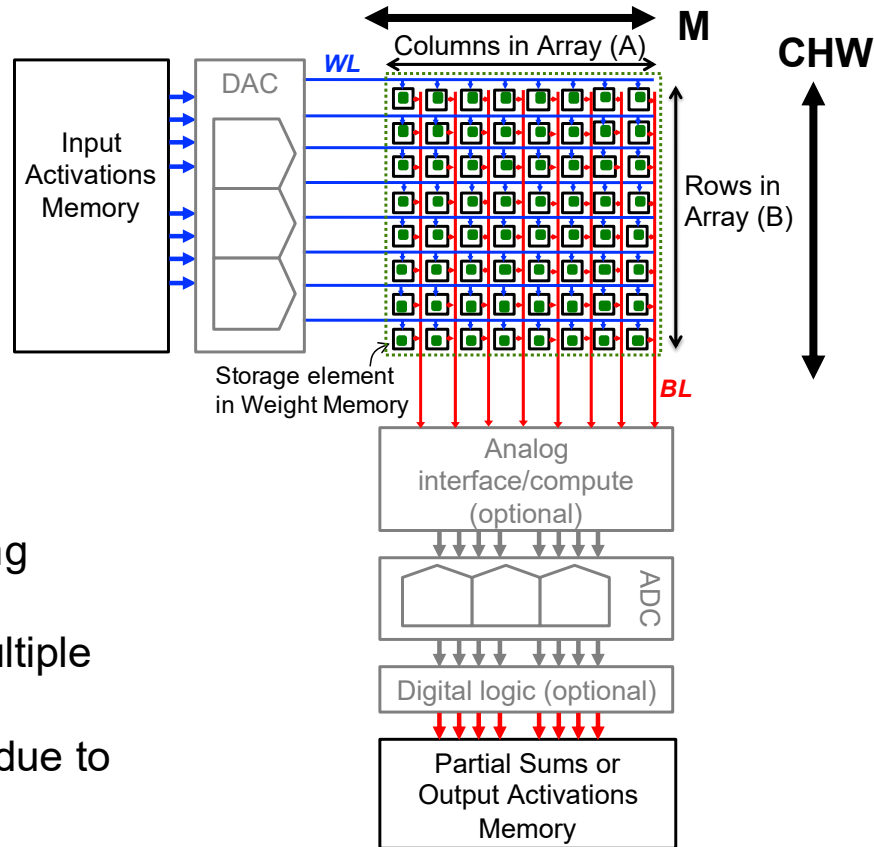
Micro-49 Tutorial on Emerging Memory
Devices (Darsen Lu)

6



Weight Stationary Dataflow for CiM

Transform into **matrix-vector multiply** with weights in matrix, and input activations in vector



Benefits

- Reduces data movement of weights
- Higher memory bandwidth by reading multiple weights in parallel
- Higher throughput by performing multiple computations in parallel
- Lower input activation delivery cost due to increased density of compute

Loop Nest for CiM

$i = \text{Array}(\text{CHW})$ *# Input activations*
 $f = \text{Array}(\text{M}, \text{CHW})$ *# Filter weights*
 $o = \text{Array}(\text{M})$ *# Output partial sums*

For Fully Connected layer

```

parallel-for m in [0, M):
  parallel-for chw in [0, CHW):
    o[m] += i[chw] * f[m, chw]

```

$i = \text{Array}(\text{C}, \text{W})$ *# Input activations*
 $f = \text{Array}(\text{M}, \text{C}, \text{S})$ *# Filter weights*
 $o = \text{Array}(\text{M}, \text{Q})$ *# Output partial sums*

For Convolutional layer
(1-D toy example)

```

parallel-for m in [0, M):
  parallel-for s in [0, S):
    parallel-for c in [0, C):
      for q in [0, Q):
        w = q + s
        o[m, q] += i[c, w] * f[m, c, s]

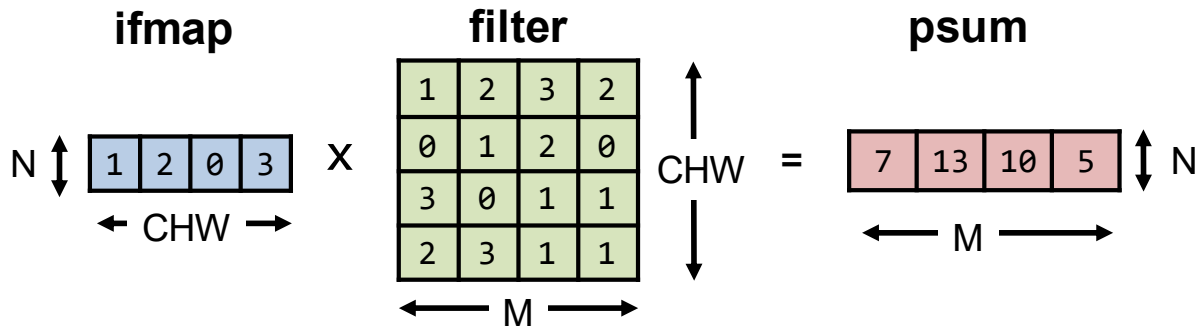
```

Design Considerations for CiM

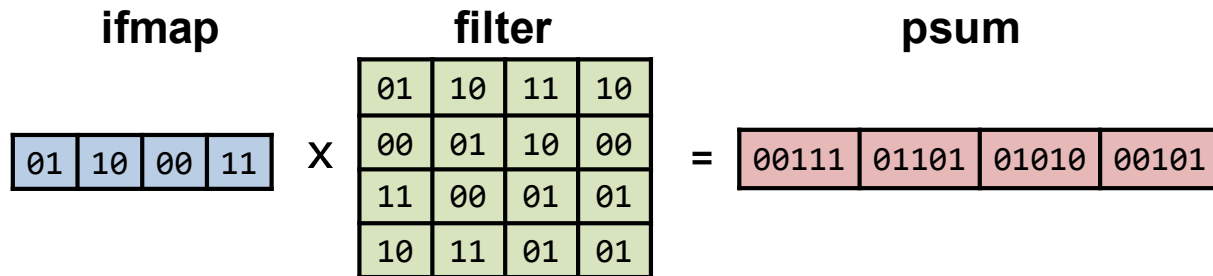
- Analog computing has **increased sensitivity** to circuit and device non-idealities (i.e., non-linearity and process, voltage and temperature variations)
- Requires **trade offs between energy efficiency, throughput, area density, and accuracy**, which reduce the achievable gains over conventional architectures

Toy Example for CiM

Decimal Representation

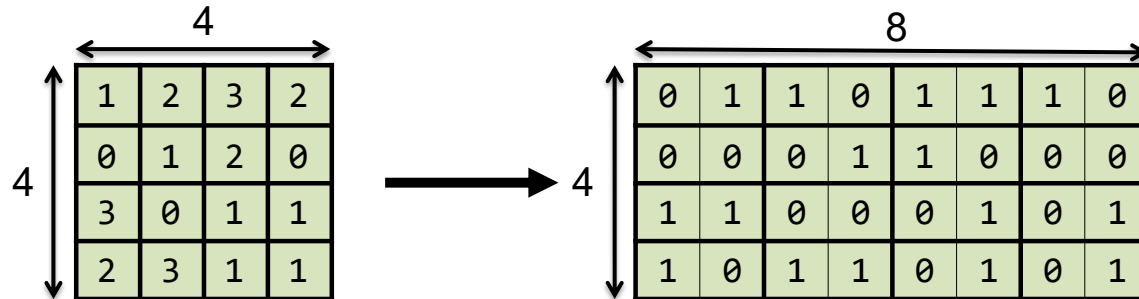


Binary Representation



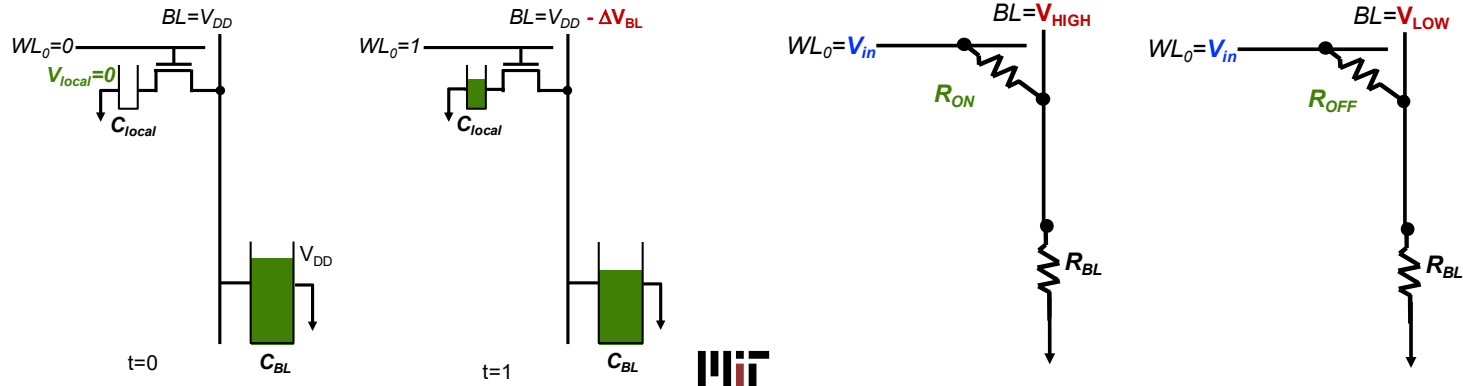
Number of Storage Elements per Weight

- Ideally, it would be desirable to be able to use one storage element (i.e., one device or bit cell) per weight to maximize density.
- In practice, **multiple storage elements are required per weight** due to the limited precision of each device or bit cell (typically on the order of 1 to 4 bits)
 - This is also referred to as **bit slicing** of the weight (**weight slicing**)



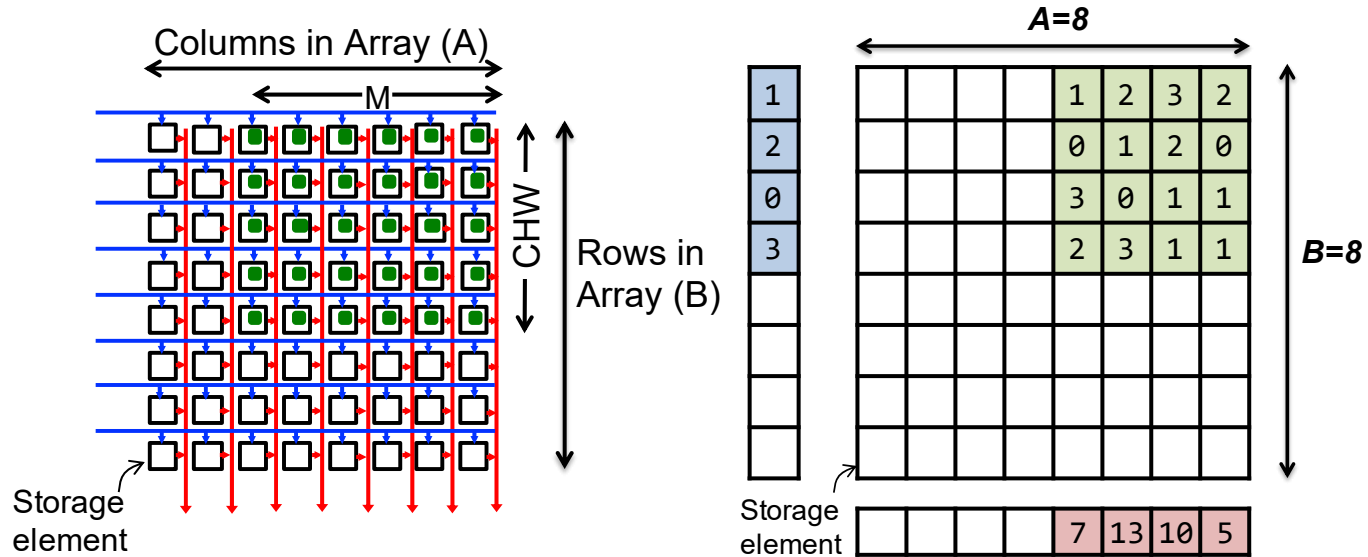
Array Size

- Ideally, it would be **desirable to have a large array size**
 - Increases weight read bandwidth and throughput
 - Increases area density → amortize the cost of the peripheral circuits (e.g., ADC, DAC) which can account for over 50% of Non-Volatile Memory (NVM) designs
- Array size limited by the resistance and capacitance of word line and bit line wires, which impacts robustness, speed and energy consumption
 - If capacitance or resistance of bit line much larger than storage element, it is difficult to sense value in storage element (i.e., change in bit-line voltage/current due to value in storage element is reduced)



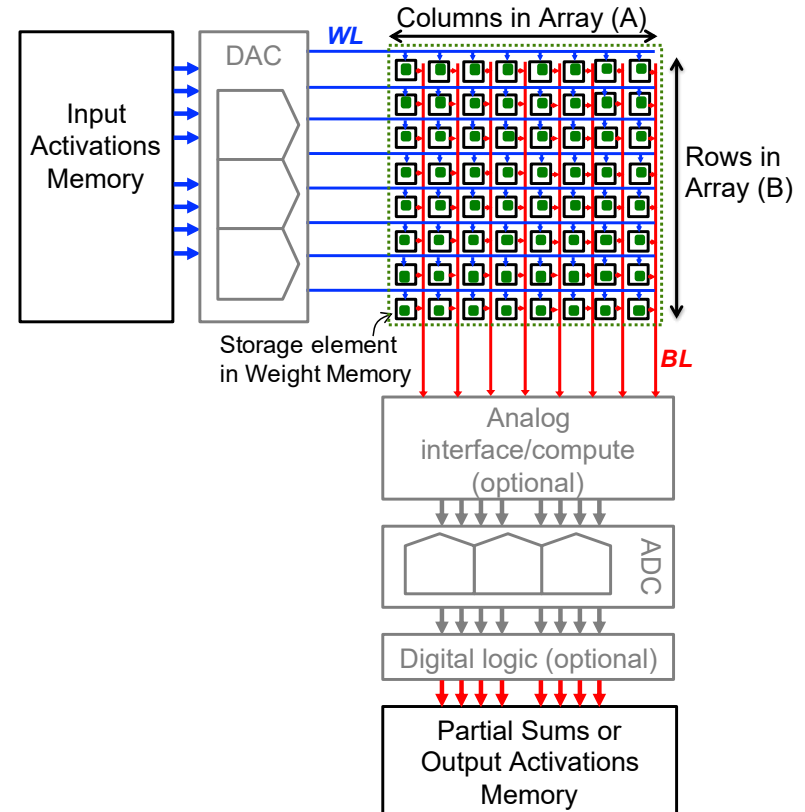
Array Size

Utilization of the array will drop if the workload cannot fill entire column or entire row

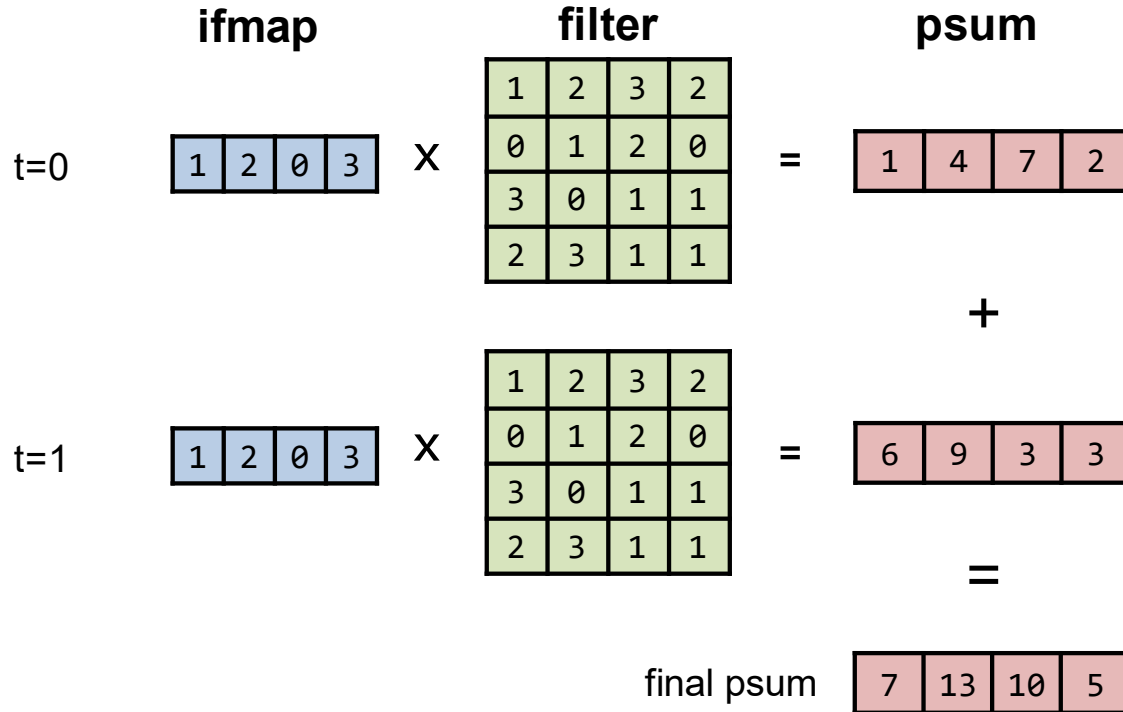


Number of Rows Activated in Parallel

- Ideally, it would be desirable to use all rows at once to maximize parallelism for high bandwidth and high throughput
- In practice, the number of rows that can be used at once is limited due to
 - ADC resolution (number of bits it can resolve)
 - Cumulative effect of the device variations
 - Maximum voltage drop or accumulated current on bit line



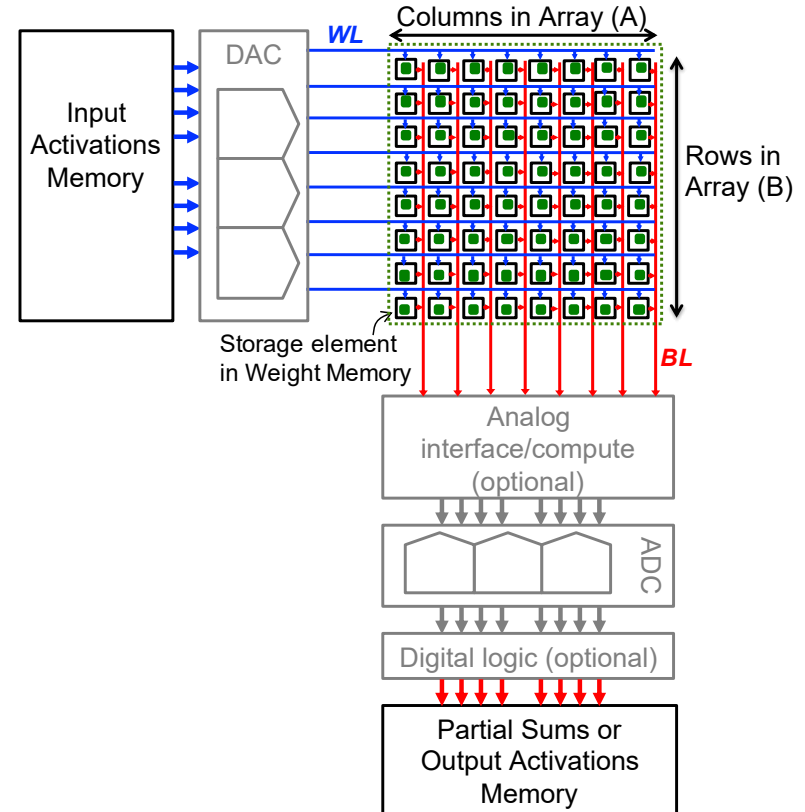
Number of Rows Activated in Parallel



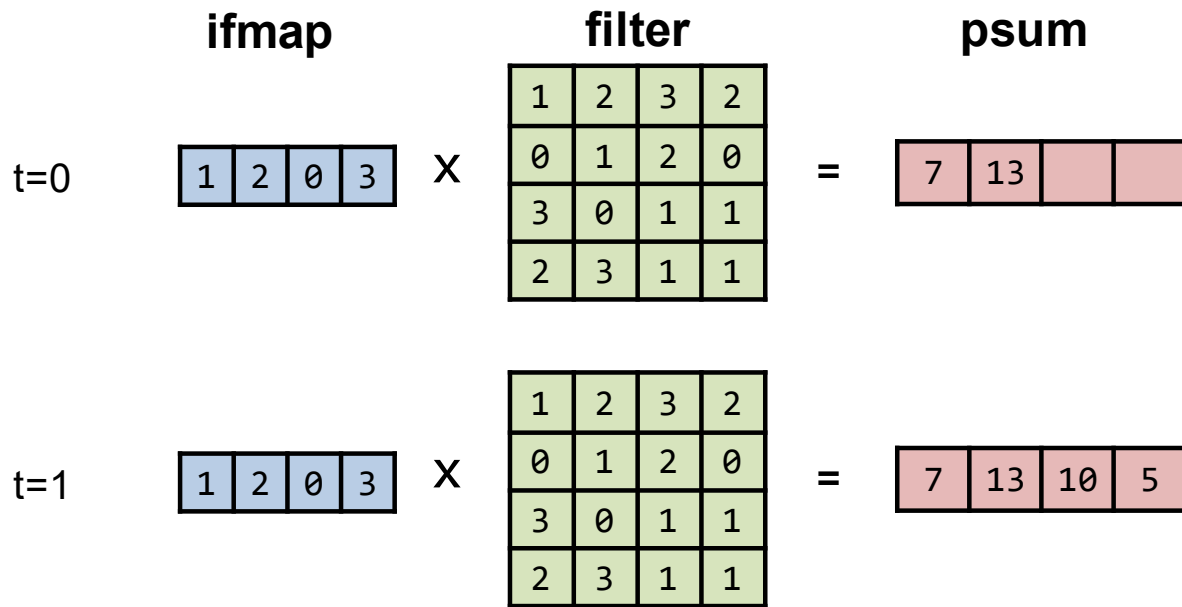
Multiple cycles to complete accumulation across rows → **Reduces Throughput**

Number of Columns Activated in Parallel

- Ideally, it would be desirable to use all columns at once to maximize parallelism for high bandwidth and high throughput
- In practice, the number of columns that can be used is limited by ADC area



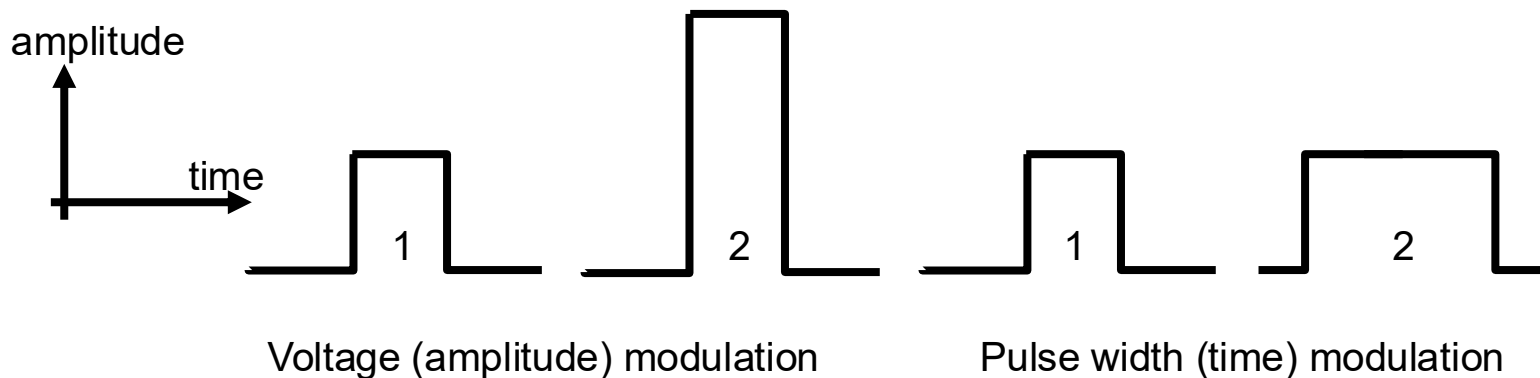
Number of Columns Activated in Parallel



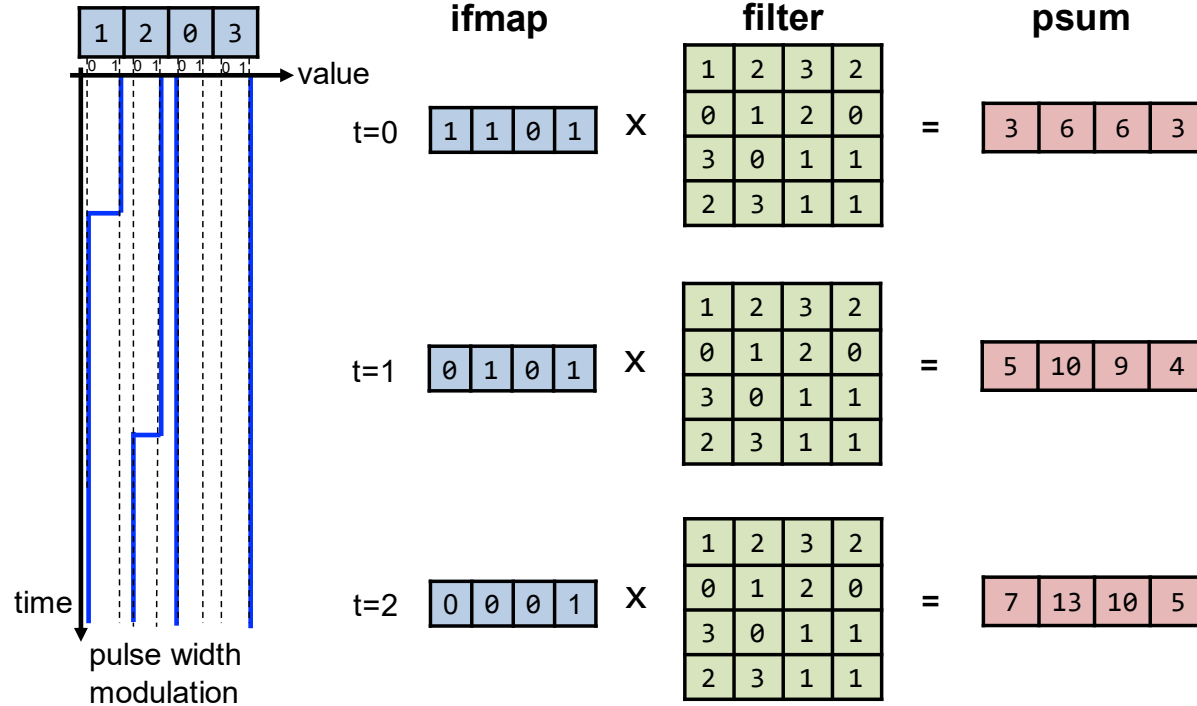
Multiple cycles to complete accumulation for all columns \rightarrow **Reduces Throughput**

Time to Deliver Input

- Ideally, encode input activations onto the word line in the minimum amount of time for maximum throughput (e.g., voltage modulation)
- In practice, challenging due to non-linearity of devices and complexity of DAC → need to encode in time → **Reduces Throughput**

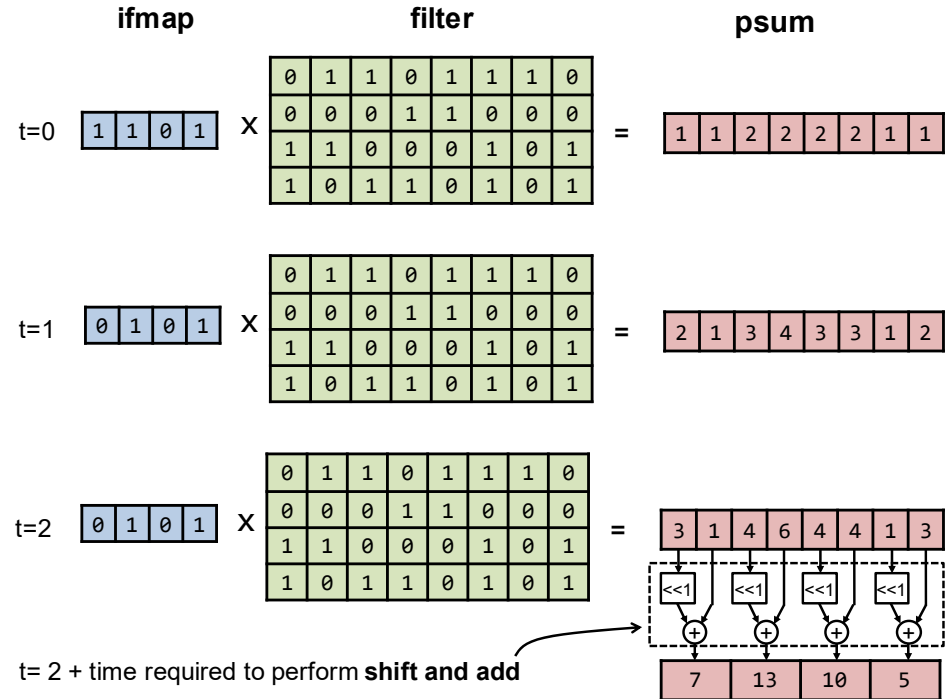


Time to Deliver Input



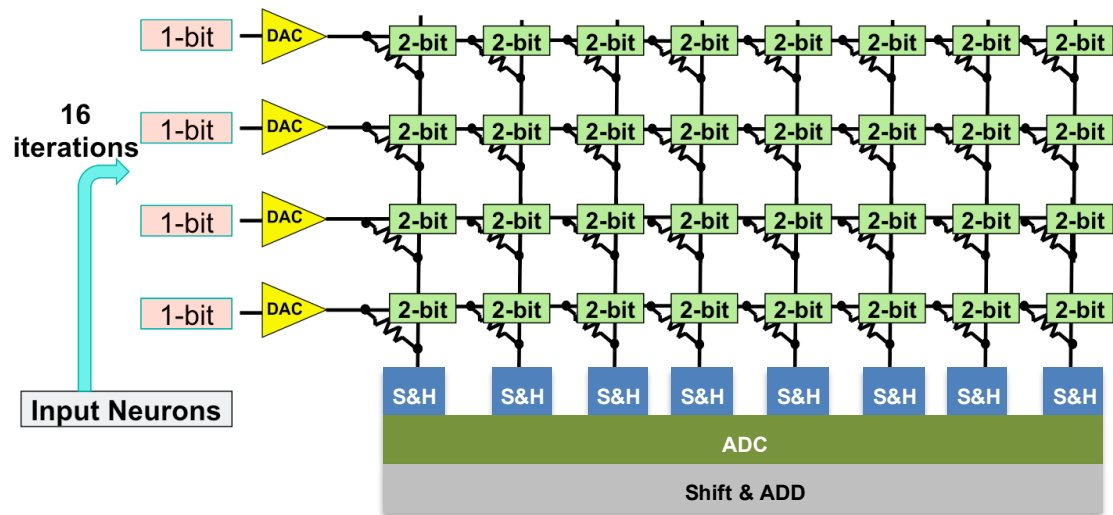
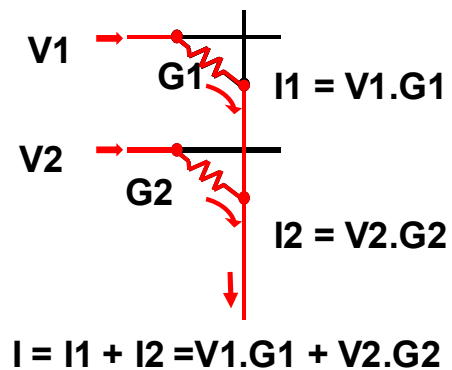
Time to Compute a MAC

- Ideally, compute MAC in a single cycle for high throughput
- In practice, the storage element typically can only perform one-bit or two-bit operations \rightarrow need multiple cycles to build up to a multi-bit operation (temporal accumulation) \rightarrow **Reduces throughput**



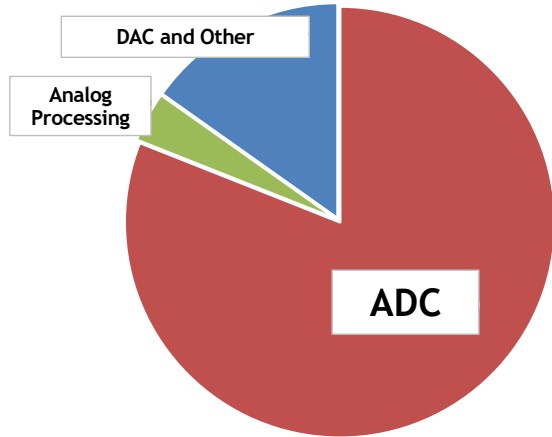
ISAAC

- Replace eDRAM in DaDianNao with memristors
- 16-bit dot-product operation
 - 8 x 2-bits per memristors (weight slicing)
 - 1-bit per cycle computation (input slicing)

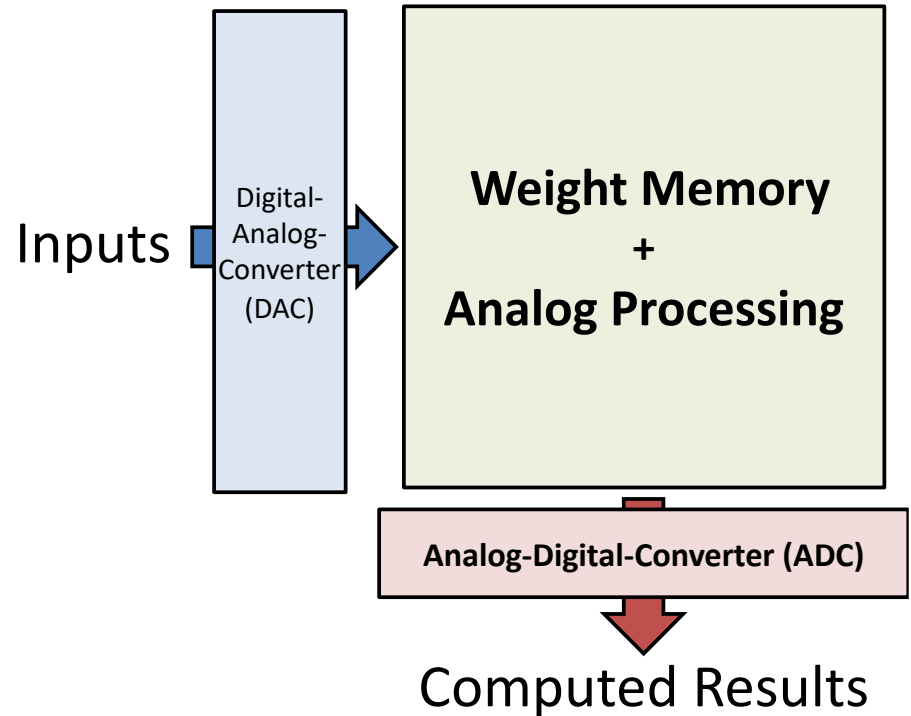


Compute In Memory (CiM) Accelerators

Energy Breakdown



ADC consumes significant energy

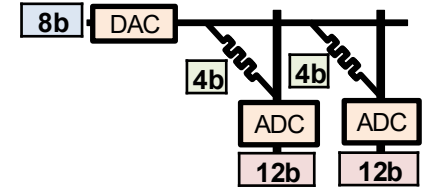
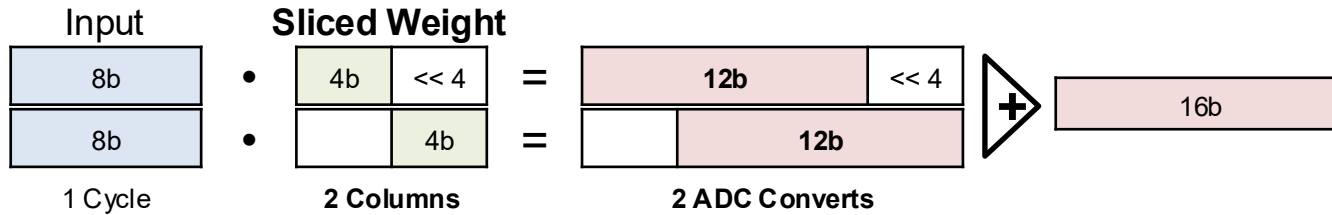
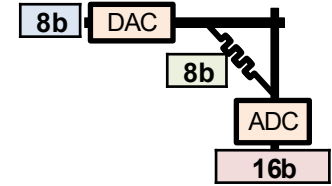
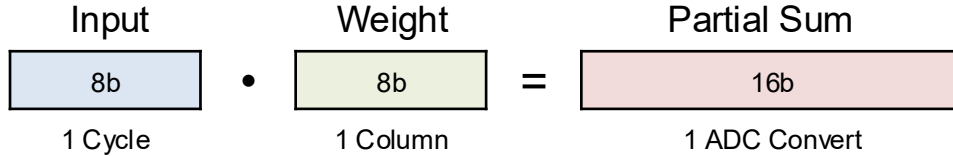


The Titanium Law

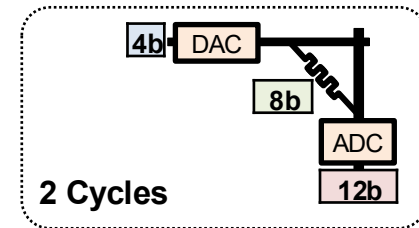
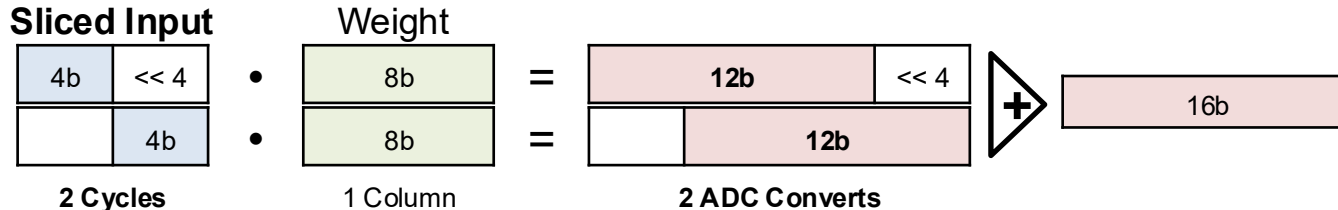
ADC energy is a product of **four** terms

$$\frac{\textit{ADC Energy}}{\textit{DNN}} = \overbrace{\frac{\textit{Energy}}{\textit{Convert}}}^{\uparrow \text{ exponentially with higher ADC resolution}} \times \underbrace{\frac{\textit{Converts}}{\textit{MAC}}}_{\substack{\downarrow \text{ with more rows} \\ \uparrow \text{ with more input/weight slices}}} \times \overbrace{\frac{\textit{MACs}}{\textit{DNN}}}^{\text{set by the DNN Workload}} \times \underbrace{\frac{1}{\textit{Utilization}}}_{\geq 1 \text{ based on row utilization}}$$

Use Bit Slicing to Reduce ADC Resolution



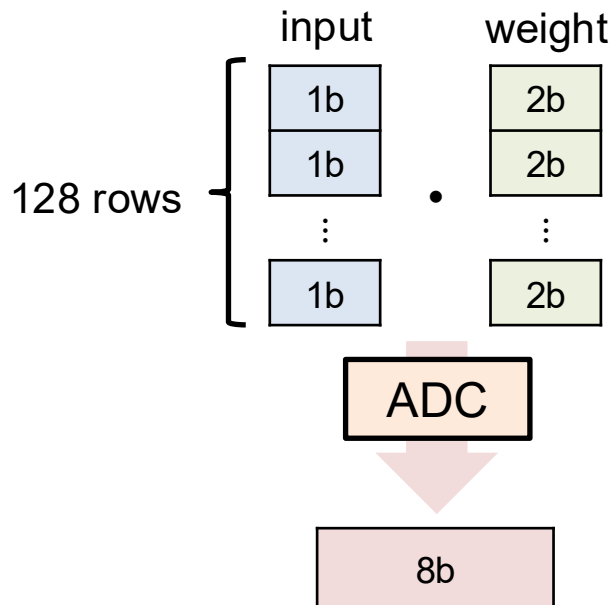
Weight slicing increases area and number of ADC converts



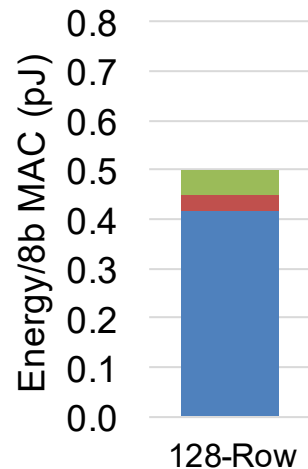
Input slicing increases time and number of ADC converts

The Titanium Law: Revisit ISAAC

$$\frac{\text{ADC Energy}}{\text{DNN}} = \frac{\text{Energy}}{\text{Convert}} \times \frac{\text{Converts}}{\text{MAC}} \times \frac{\text{MACs}}{\text{DNN}} \times \frac{1}{\text{Utilization}}$$



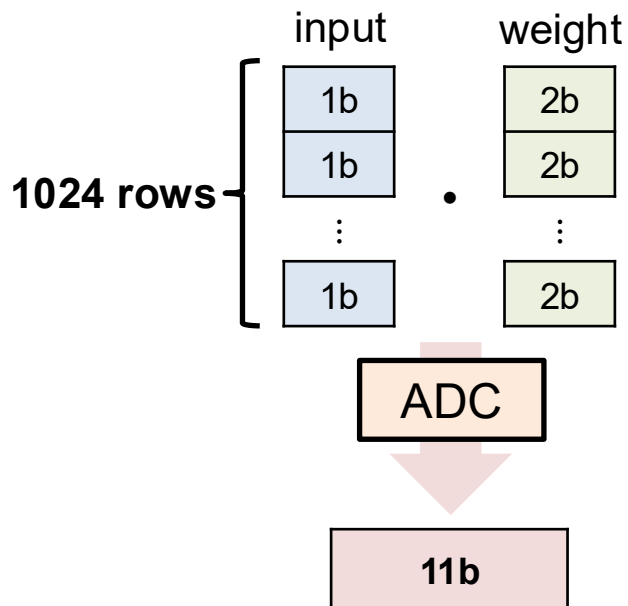
Can we **reduce** ADC energy?



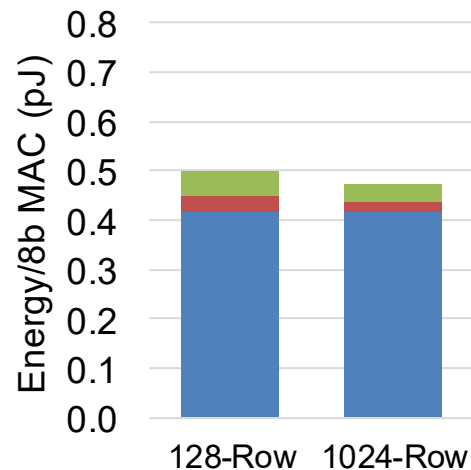
■ ADC ■ Analog Crossbar ■ Other

The Titanium Law: Revisit ISAAC

$$\frac{\text{ADC Energy}}{\text{DNN}} = \uparrow \frac{\text{Energy}}{\text{Convert}} \times \downarrow \frac{\text{Converts}}{\text{MAC}} \times \frac{\text{MACs}}{\text{DNN}} \times \frac{1}{\text{Utilization}}$$



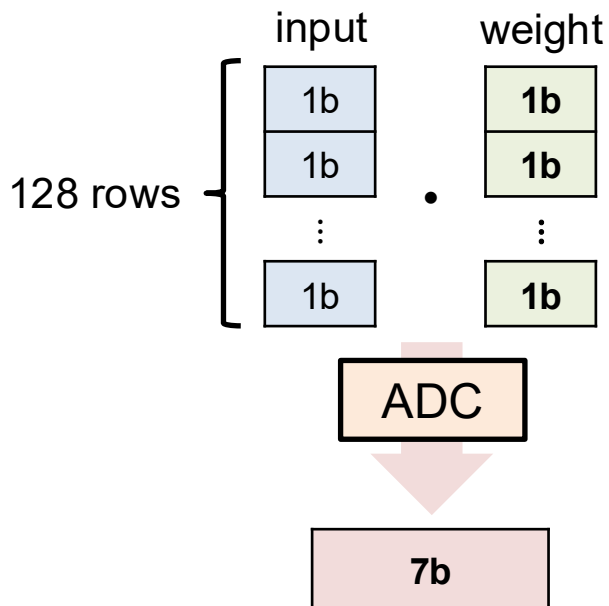
Increase rows
↓
Increase bits



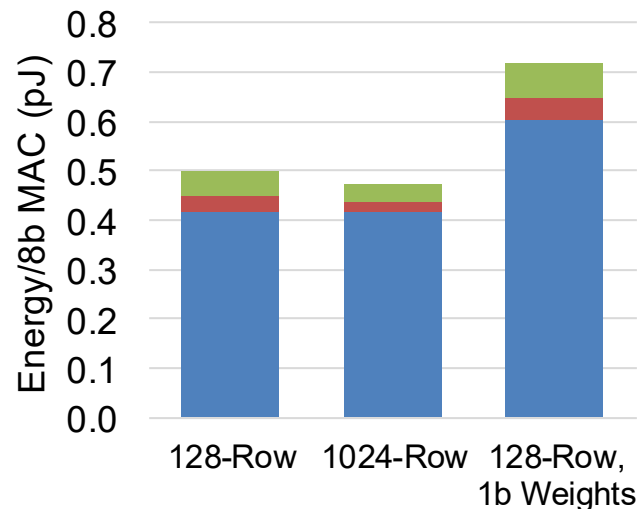
■ ADC ■ Analog Crossbar ■ Other

The Titanium Law: Revisit ISAAC

$$\frac{\text{ADC Energy}}{\text{DNN}} = \downarrow \frac{\text{Energy}}{\text{Convert}} \times \uparrow \frac{\text{Converts}}{\text{MAC}} \times \frac{\text{MACs}}{\text{DNN}} \times \frac{1}{\text{Utilization}}$$



Decrease bits/weight slice
 ↓
Increase weight slices
 ↓
Increase ADC converts

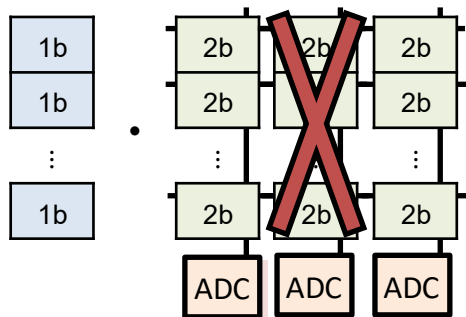


■ ADC ■ Analog Crossbar ■ Other

How Have Prior Works Escaped These Tradeoffs?

$$\frac{\text{ADC Energy}}{\text{DNN}} = \frac{\text{Energy}}{\text{Convert}} \times \frac{\text{Converts}}{\text{MAC}} \times \downarrow \frac{\text{MACs}}{\text{DNN}} \times \frac{1}{\text{Utilization}}$$

Weight-Count-Limited



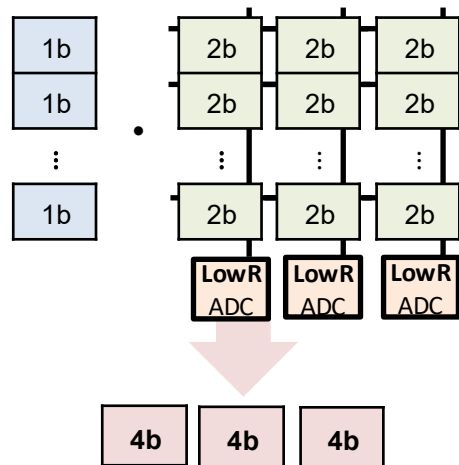
Prune weights
↓
Reduce MACs/DNN



How Have Prior Works Escaped These Tradeoffs?

$$\frac{\text{ADC Energy}}{\text{DNN}} = \downarrow \frac{\text{Energy}}{\text{Convert}} \times \frac{\text{Converts}}{\text{MAC}} \times \frac{\text{MACs}}{\text{DNN}} \times \frac{1}{\text{Utilization}}$$

Sum-Fidelity-Limited

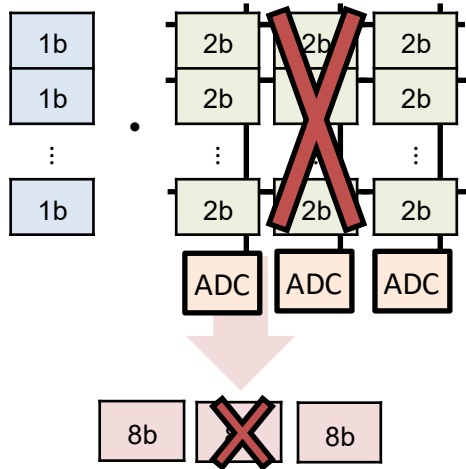


Low-Res ADC
↓
Reduce energy per convert

How Have Prior Works Escaped These Tradeoffs?

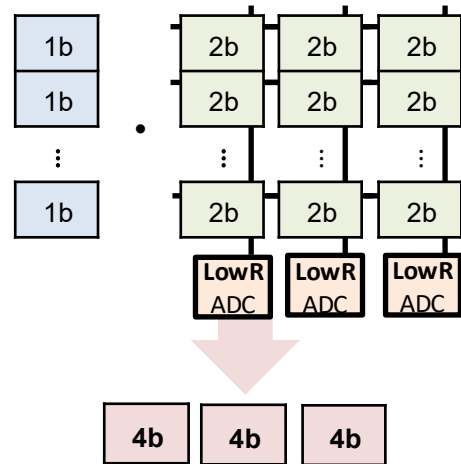
$$\frac{\text{ADC Energy}}{\text{DNN}} = \frac{\text{Energy}}{\text{Convert}} \times \frac{\text{Converts}}{\text{MAC}} \times \frac{\text{MACs}}{\text{DNN}} \times \frac{1}{\text{Utilization}}$$

Weight-Count-Limited



Prune weights
↓
Reduce MACs/DNN

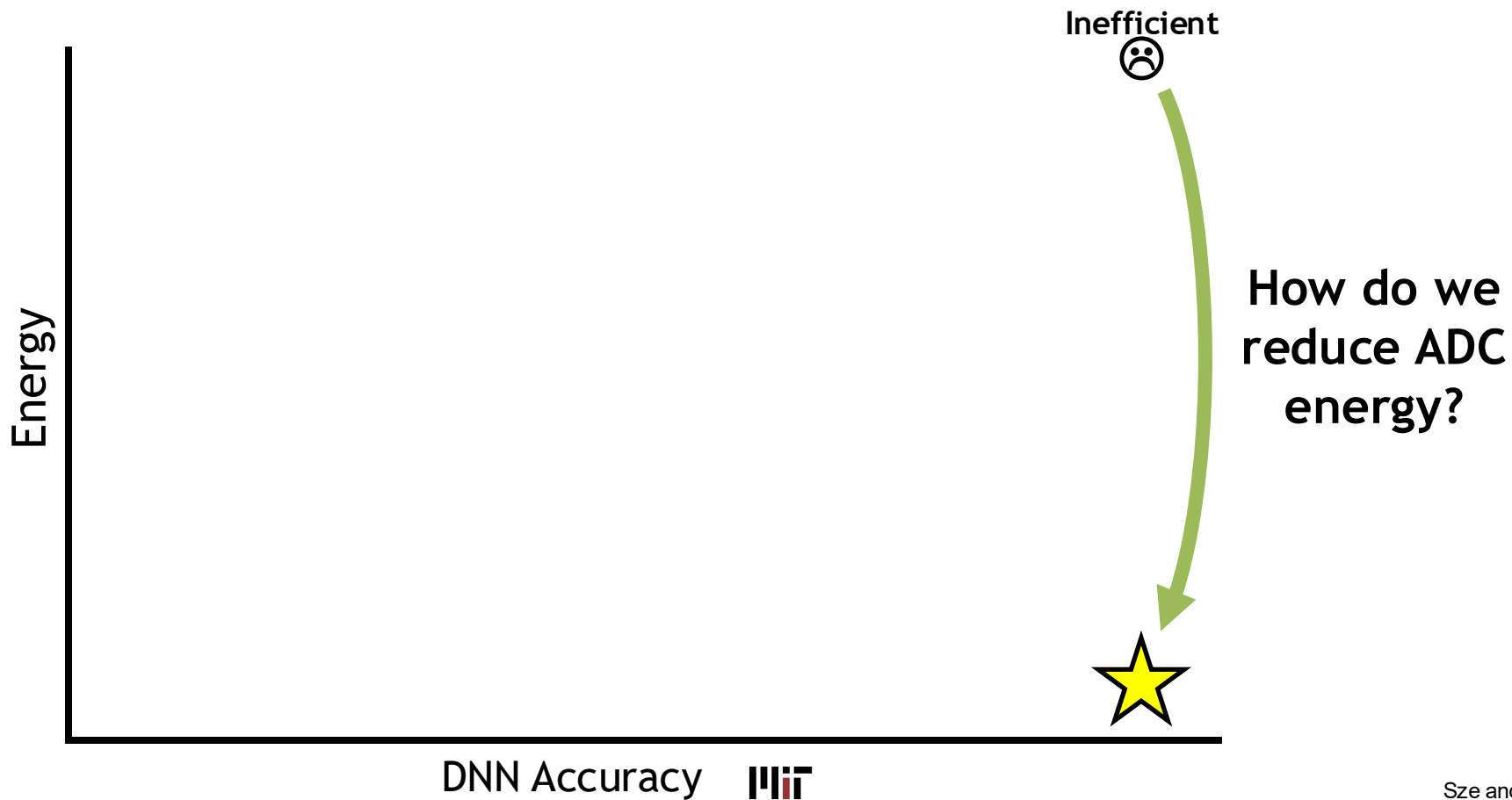
Sum-Fidelity-Limited



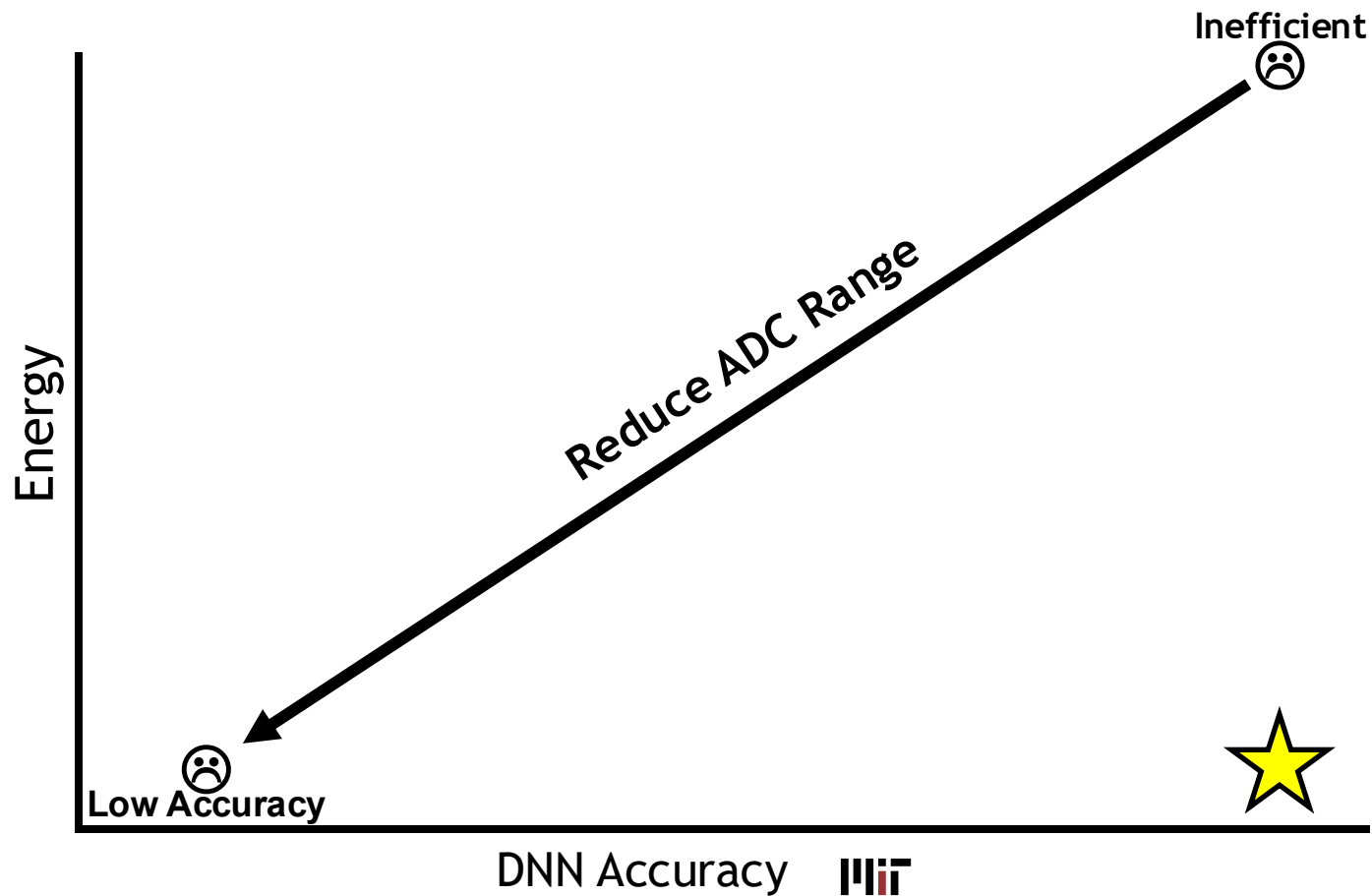
Low-Res ADC
↓
Reduce energy
per convert

Both approaches may require **retraining DNN** to preserve accuracy

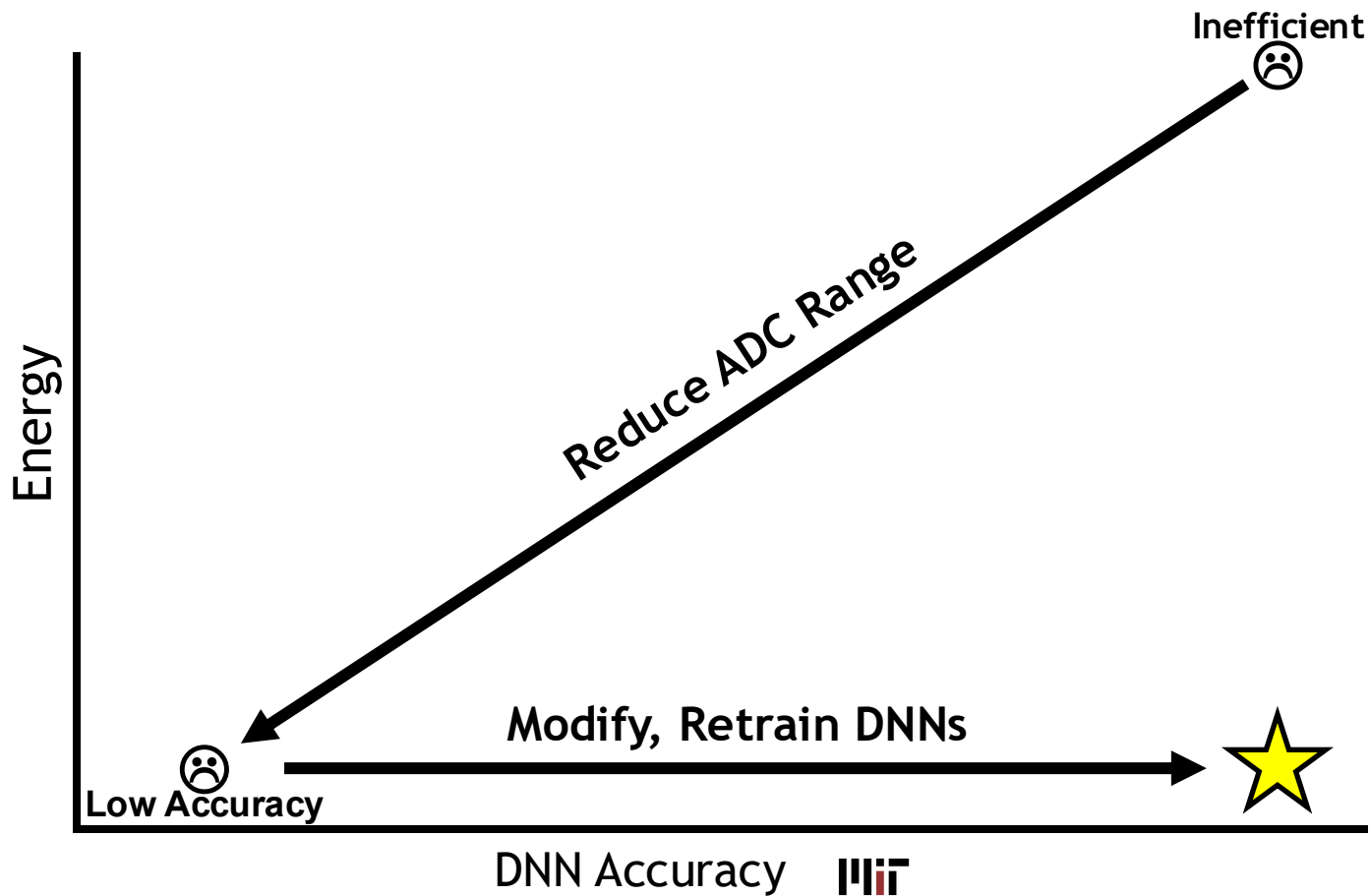
Reducing ADC Energy



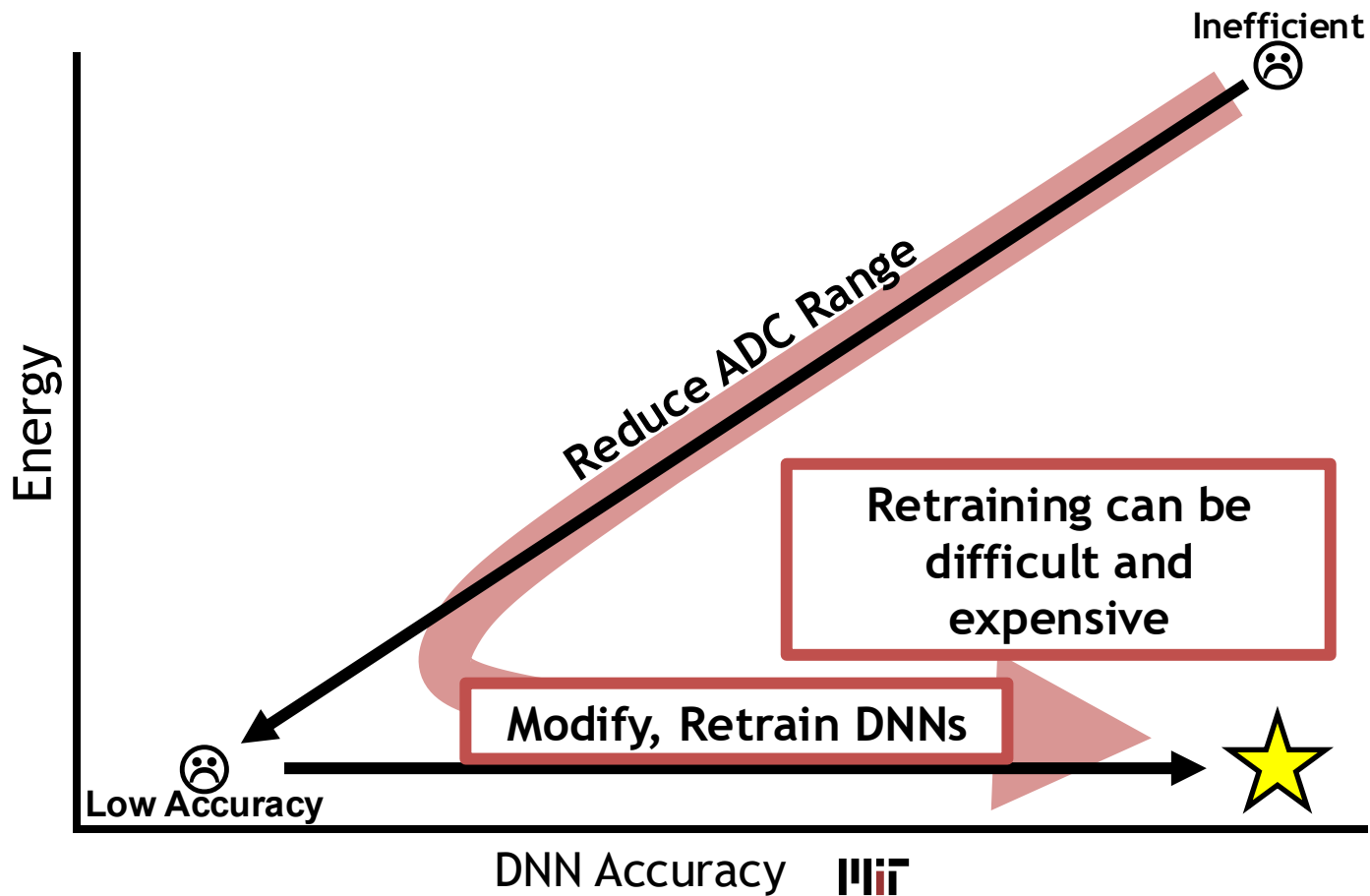
Reducing ADC Energy



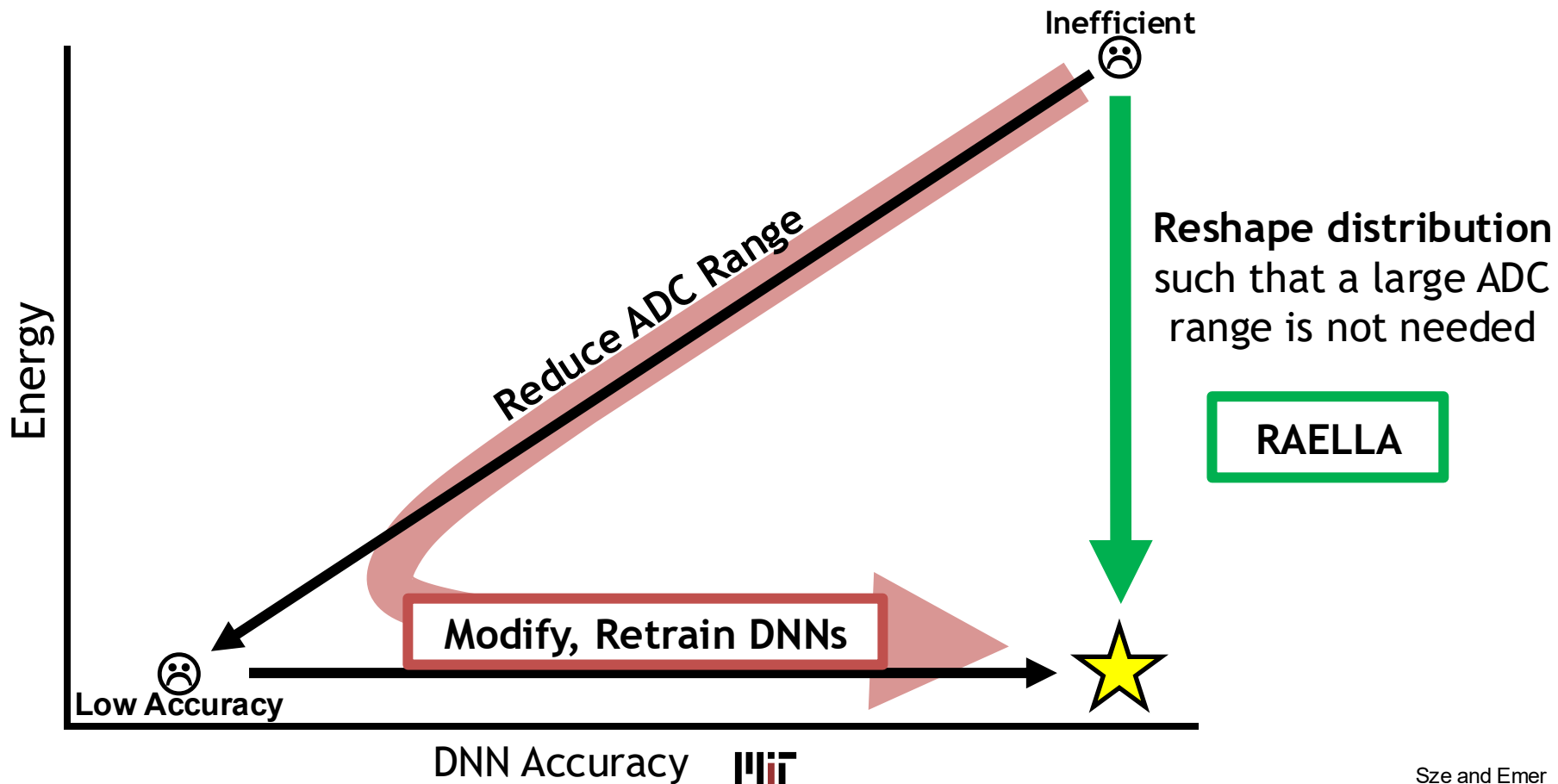
Reducing ADC Energy



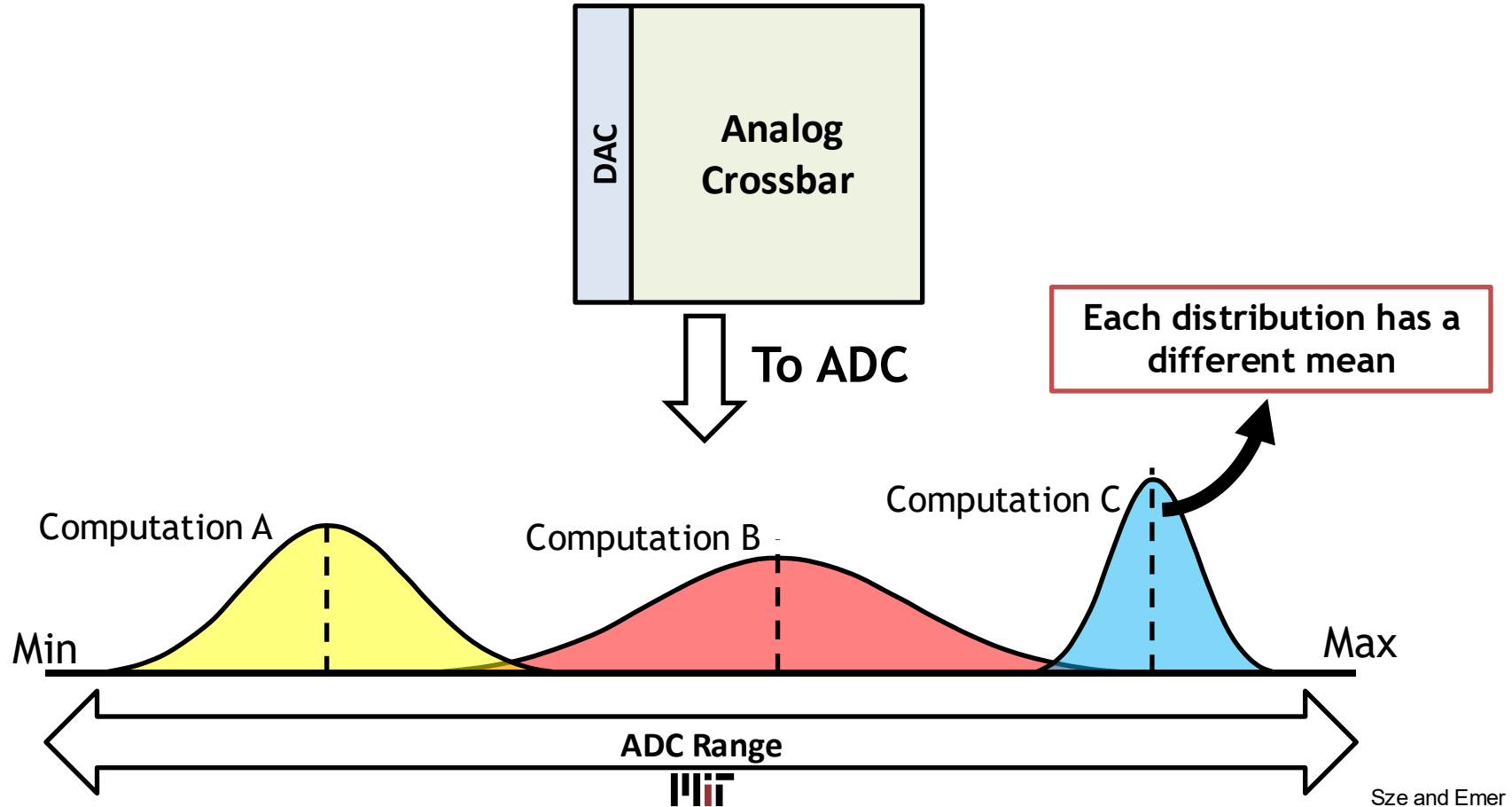
Reducing ADC Energy



Reducing ADC Energy

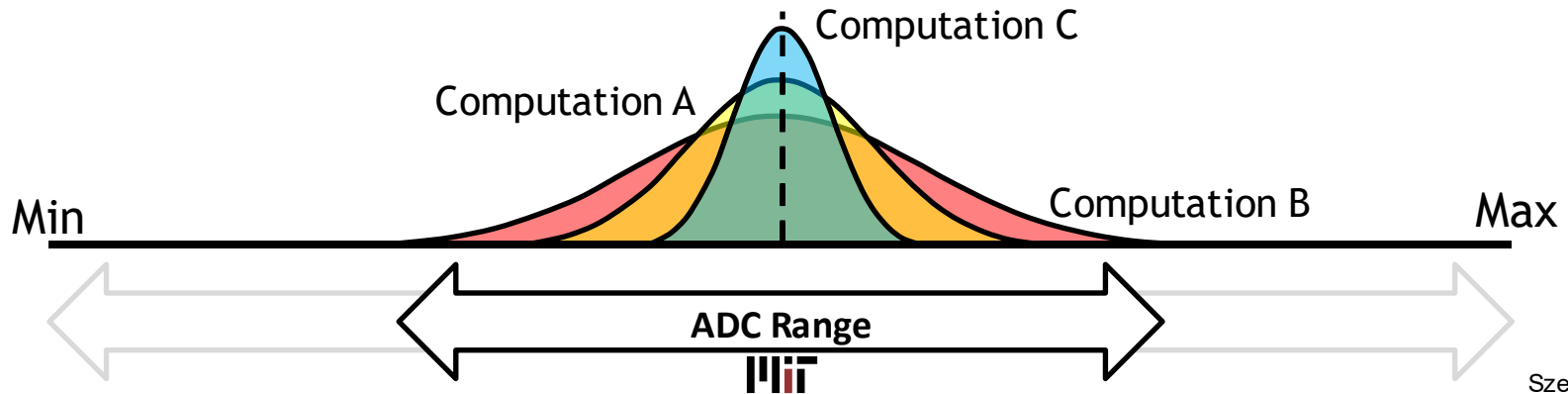
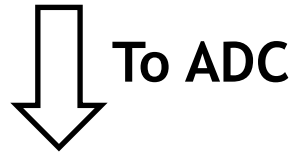
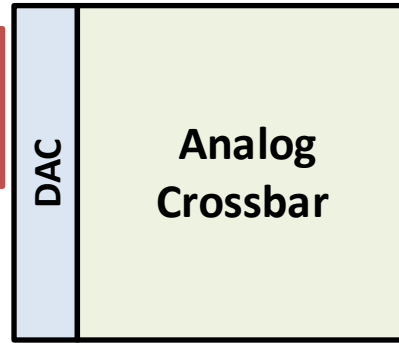


RAELLA: Reshape Distributions to Reduce Input to ADC



RAELLA: Reshape Distributions

1. Shift the mean of each distribution to the center of the ADC range

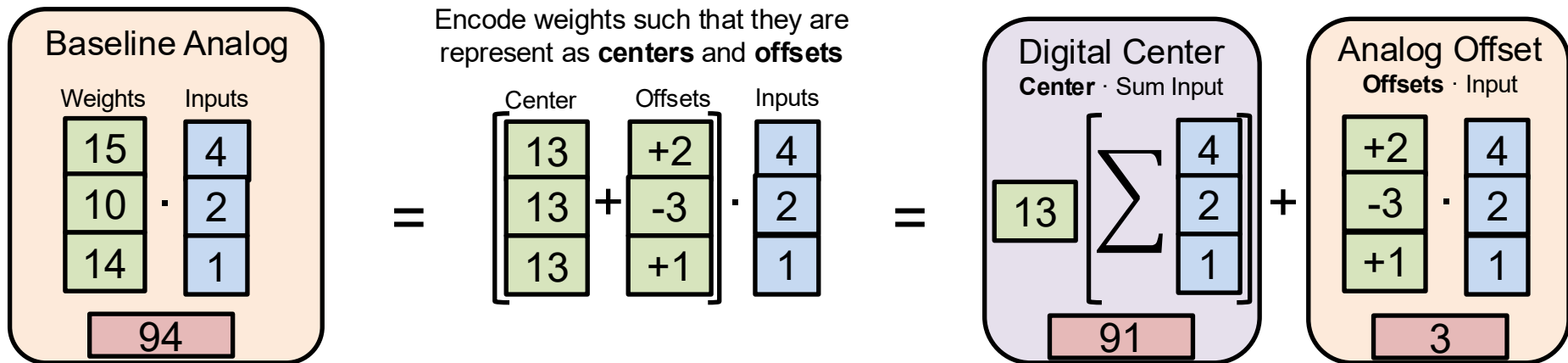


Center + Offset Weight Encoding Zero-Average Analog Results

Partition computation

Digital calculates high-resolution **center** operations

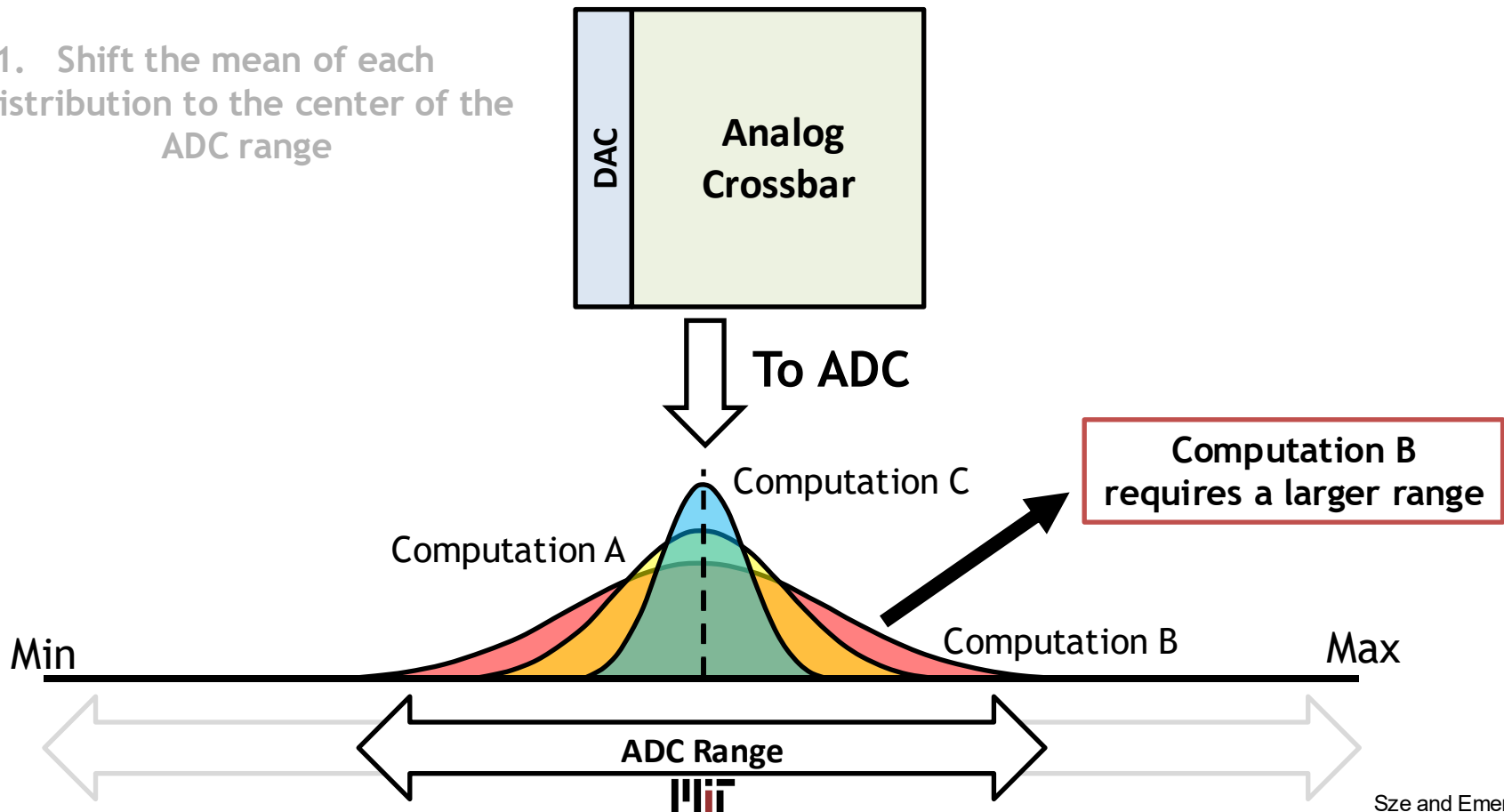
Analog calculates parallel **offset** operations



Encoding allows analog input to ADC to be
smaller and closer to zero

RAELLA: Reshape Distributions

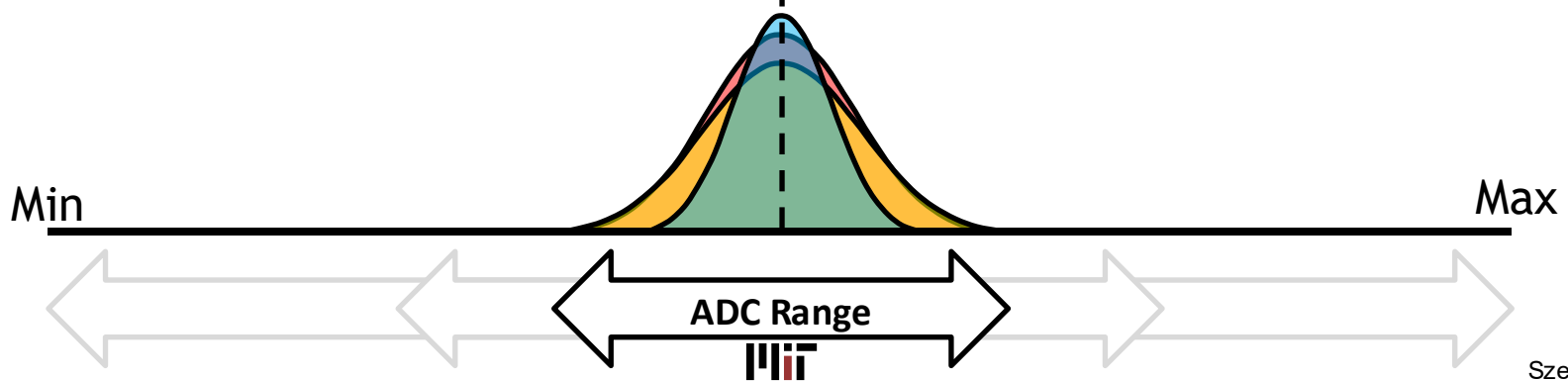
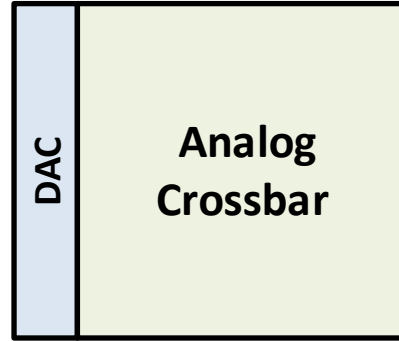
1. Shift the mean of each distribution to the center of the ADC range



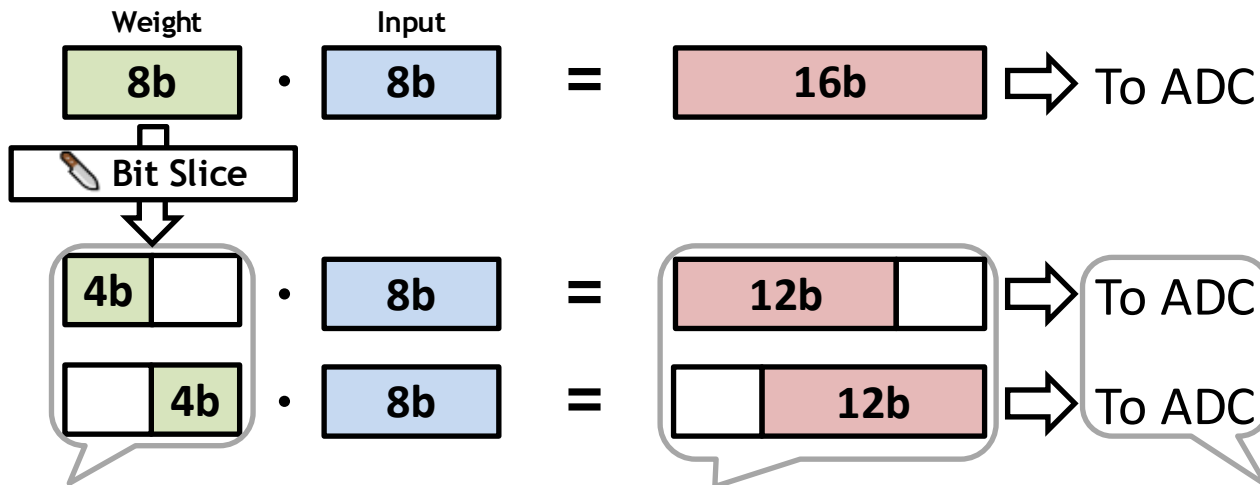
RAELLA: Reshape Distributions

1. Shift the mean of each distribution to the center of the ADC range

2. If a computation produces large results, slice it into smaller pieces



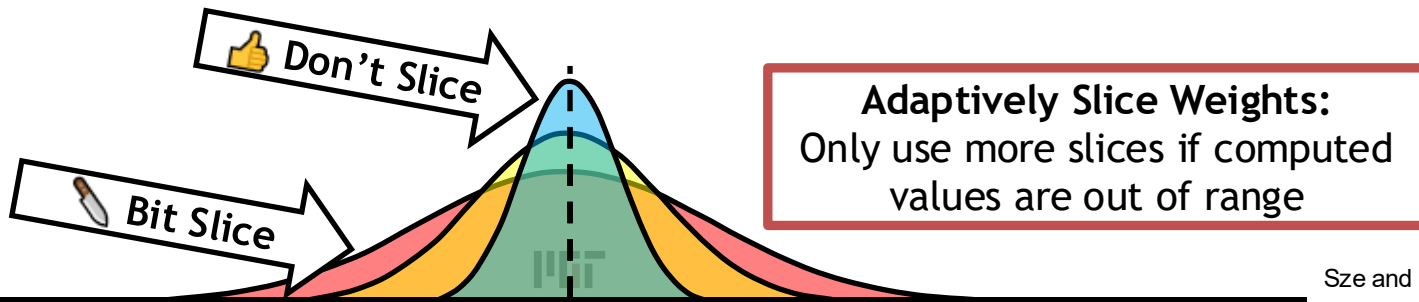
Adaptive Weight Slicing: Slice Large-Result Computations



☹ More Memory (Area)

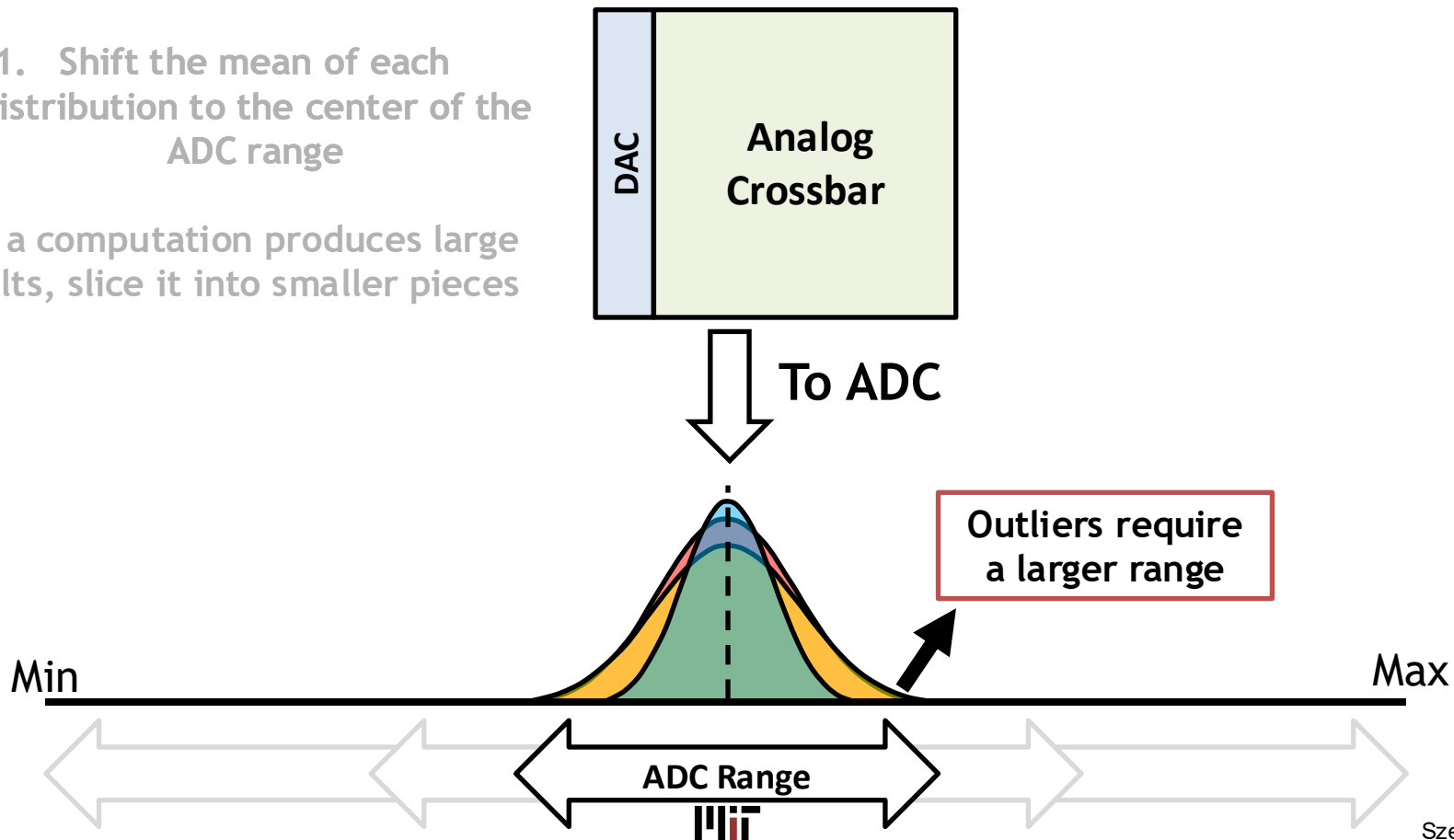
😊 Smaller Range

☹ More ADC Converts (Energy)



RAELLA: Reshape Distributions

1. Shift the mean of each distribution to the center of the ADC range
2. If a computation produces large results, slice it into smaller pieces

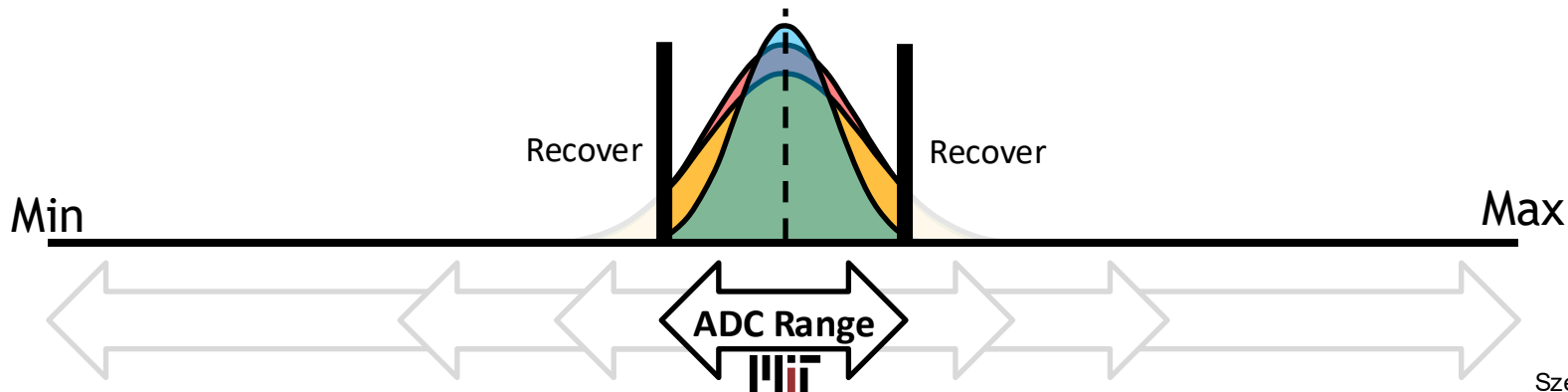
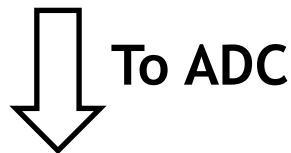
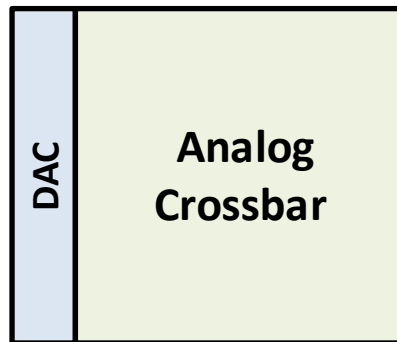


RAELLA: Reshape Distributions

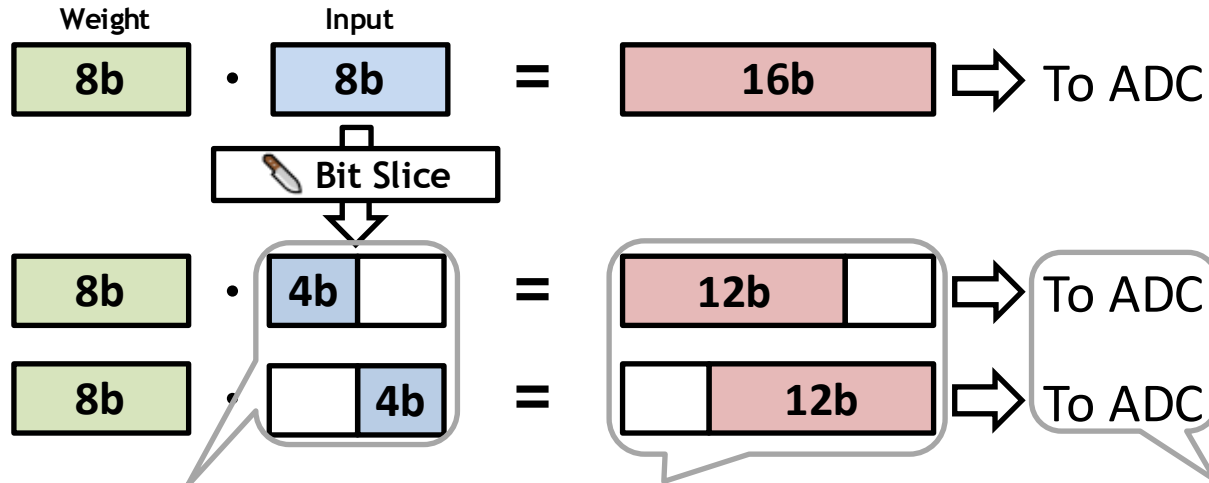
1. Shift the mean of each distribution to the center of the ADC range

2. If a computation produces large results, slice it into smaller pieces

3. Speculate that results are in-range, recover out-of-range results



Dynamic Input Slicing: Try Again with Smaller Slices

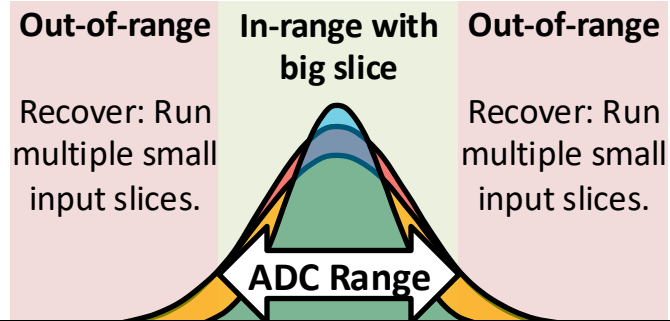


☹ More Cycles (Time)

😊 Smaller Range

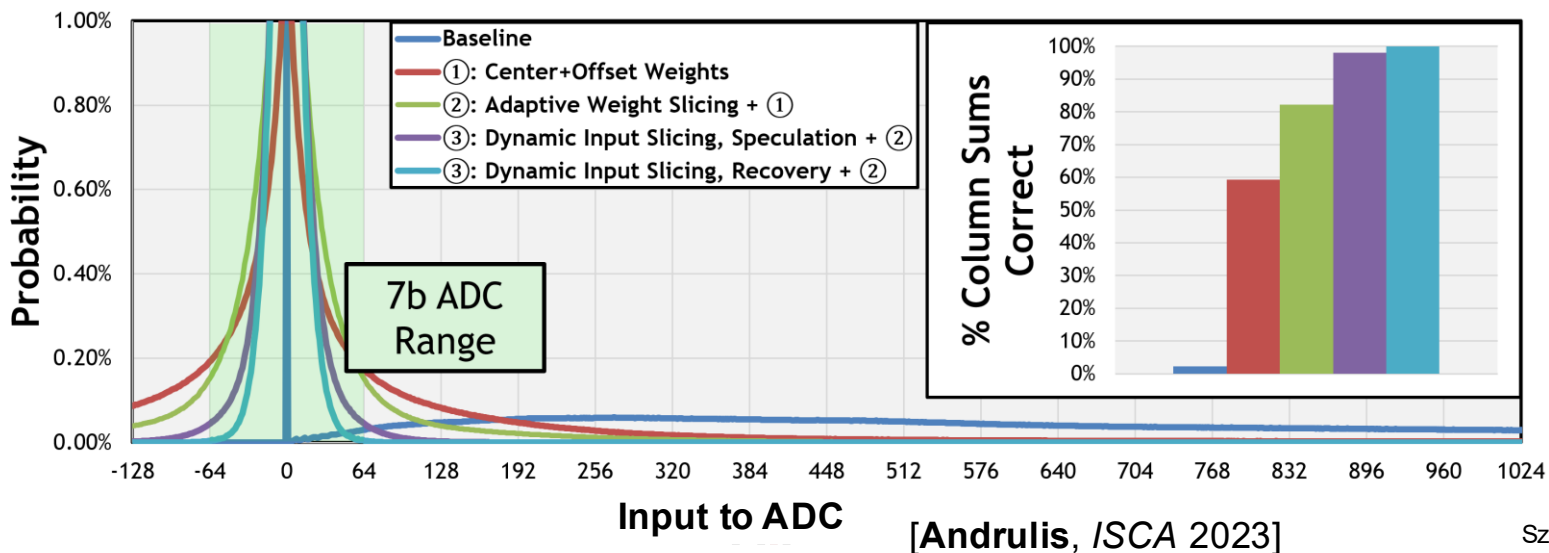
☹ More ADC Converts (Energy)

1. Speculate with big input slice
2. Recover out-of-range results with multiple smaller input slices



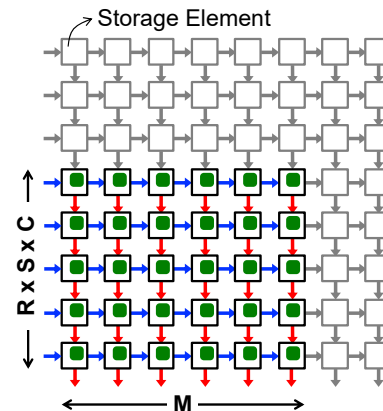
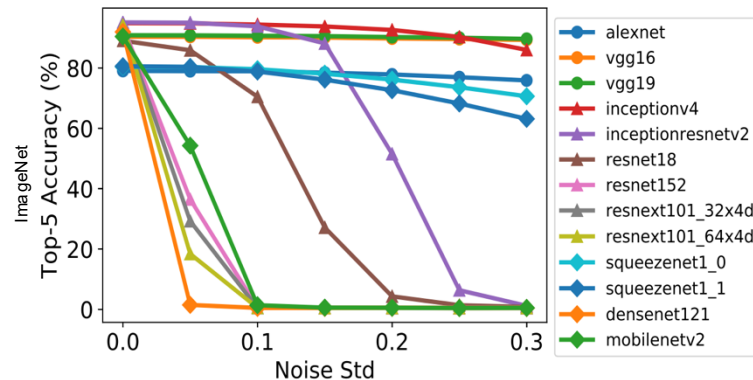
RAELLA: Reshape Distributions of Input to ADC

- Makes analog operations produce low-resolution results
 - 1024x reduction of input to ADC
- Enables more compute per ADC convert while using lower-resolution ADCs
 - Improves energy efficiency by 3.9x and throughput by 1.8x compared to iso-area ISAAC
- Maintains DNN accuracy without changing DNN or retraining



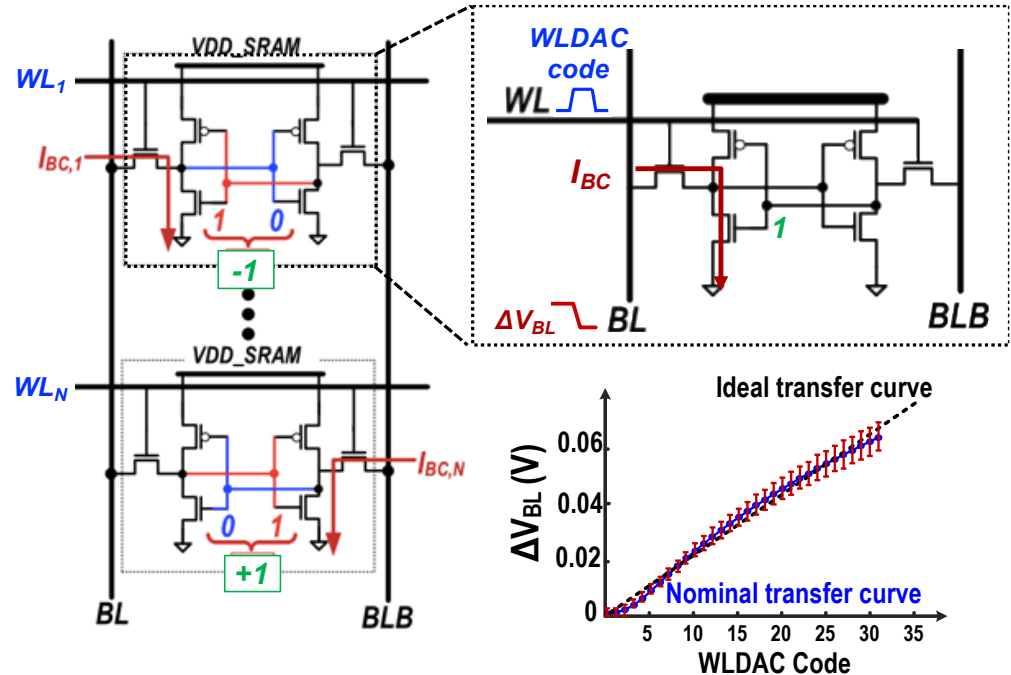
Designing DNN Models for CiM

- Designing DNNs for CiM may differ from DNNs for digital processors
- Highest accuracy DNN on digital processor may be different on CiM
 - Accuracy drops based on robustness to non-idealities
- Reducing number of weights is less desirable
 - Since CiM is weight stationary, may be better to reduce number of activations
 - CiM tend to have larger arrays \rightarrow fewer weights may lead to low utilization on CiM
- Current trend is deeper and smaller filters
 - CiM may prefer to have shallower and larger filters



CiM Using SRAM Bit Cell

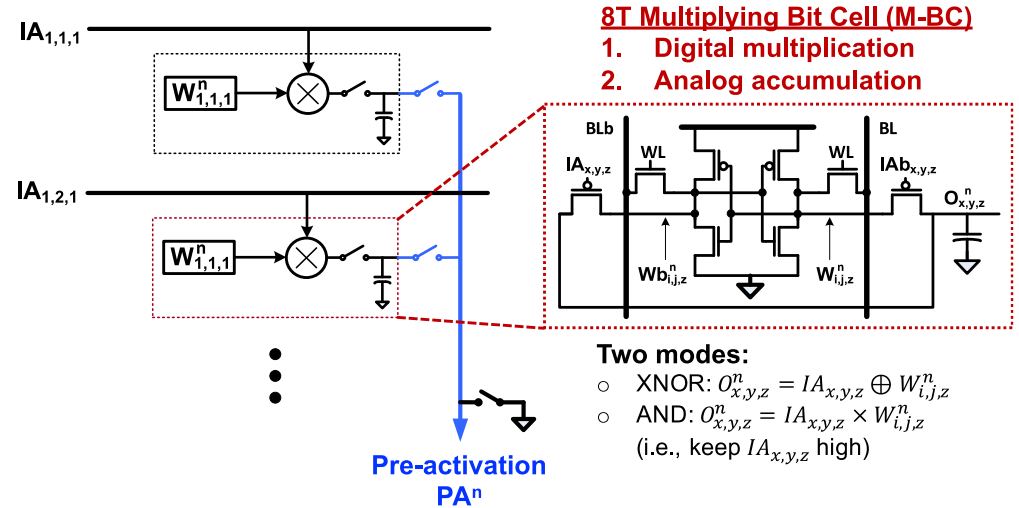
- Multiplication uses I-V relationship of access transistor (WL) and stored value in bit-cell
 - Assumes binary weights and multi-bit input activation
- Addition using **current addition** on bit line (BL)
 - Limited by nonlinearity and sensitivity to variations



[Verma, SSCS 2019]

CiM Using SRAM Bit Cell

- Binary multiplication (AND or XNOR) using access transistor (WL) and stored value in bit-cell
 - Explicit capacitor to store charge
- Addition using **charge sharing** on bit line (BL)
 - Better linearity and matching



Using Charge Sharing for Addition

Capacitive charge sharing Analog snippets

$Q_{INIT} = C_1(V_1 - V_1') + C_2(V_2 - V_2')$
 $Q_{FINAL} = C_1(V_F - V_1') + C_2(V_F - V_2')$

$Q_{INIT} = Q_{FINAL}$
 $\Rightarrow C_1 V_1 + C_2 V_2 = V_F (C_1 + C_2)$
 $\Rightarrow \boxed{V_F = \frac{C_1 V_1 + C_2 V_2}{C_1 + C_2}}$

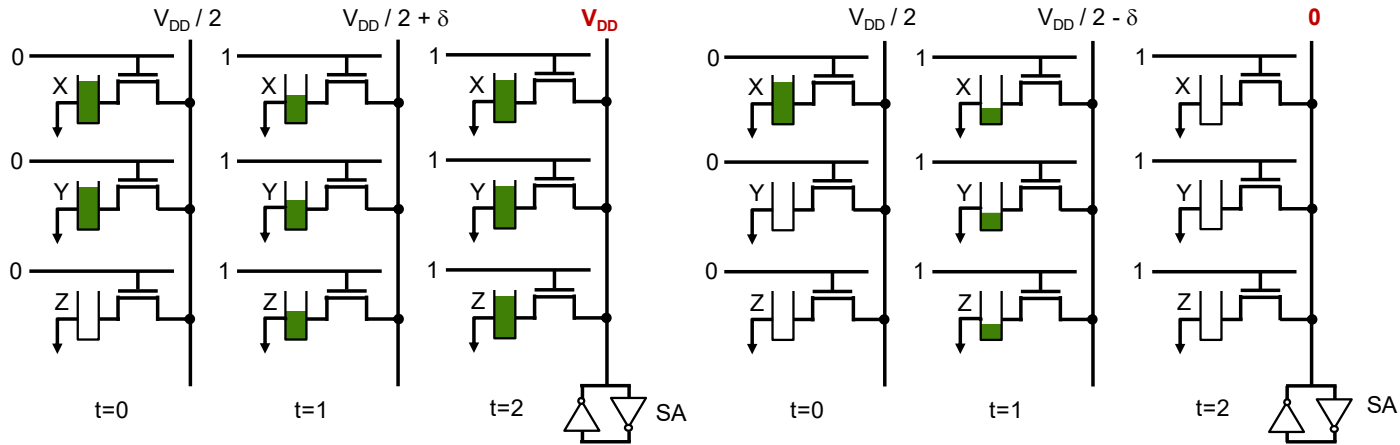
Image Source: https://www.youtube.com/watch?v=XRQ_Xldr2nk

If $C_1 = C_2$, $V_f = \frac{1}{2}(V_1 + V_2)$, which is a scaled value of the sum (addition)

CiM Using DRAM

Performs bit-wise operations using charge sharing

If $Z=0$, perform $\text{AND}(X, Y)$



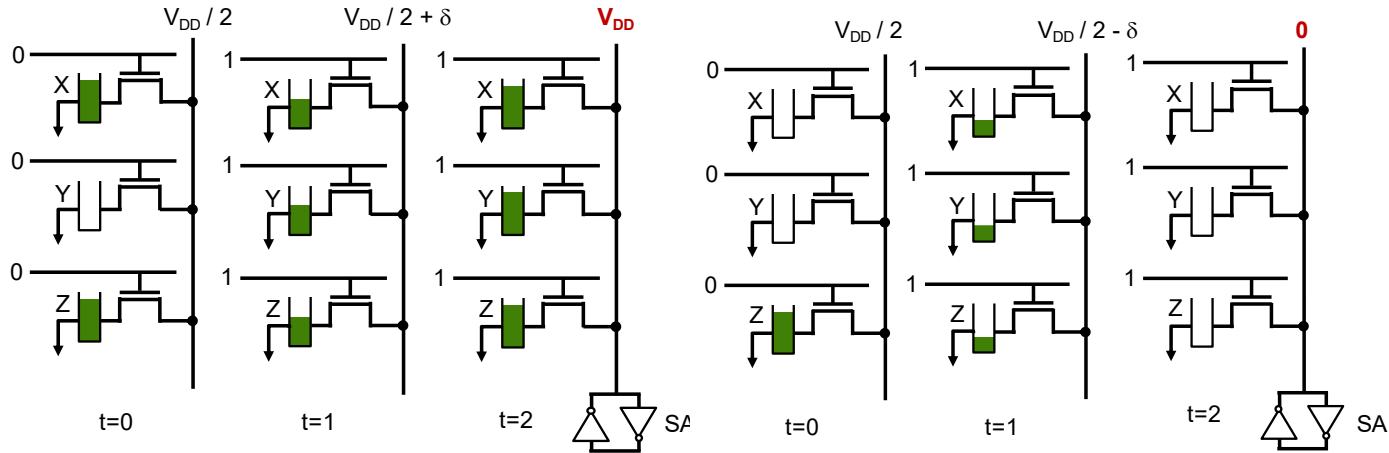
$$\text{AND}(X=1, Y=1) = 1$$

$$\text{AND}(X=1, Y=0) = 0$$

CiM Using DRAM

Performs bit-wise operations using charge sharing

If $Z=1$, perform $OR(X, Y)$



$OR(X=1, Y=1) = 1$

$OR(X=0, Y=0) = 0$

Takes multiple cycles to built up to a multiplication.
However, can perform many operations in parallel (bus width of DRAM)

Lots of Compute In Memory Research Across the Stack

SESSION 34 Wednesday, February 21st, 1:30 PM Compute-In-Memory

- Session Chair:** Anker Agrawal, IBM T. J. Watson Research Center, Yorktown Heights, NY
- Session Co-Chair:** Eric Wang, TSMC, Hsinchu City, Taiwan
- 1:30 PM
- 34.1 A 28nm 63.2Tflops/W POSIT-Based Compute-in-Memory Macro for High-Accuracy AI Applications**
Y. Wang, X. Yang, Y. Qin, Z. Zhao, R. Guo, Z. Yue, H. Han, S. Wu, Y. Hu, S. Yin
Tsinghua University, Beijing, China
- 1:55 PM
- 34.2 A 16nm 90k Integer-Floating-Point Dual-Mode Gain-Cell-Computing-In-Memory Macro Achieving 71.3-163.2TOPS/W and 21.2-21.2Tflops/W for AI-Edge Devices**
W.S. Ahsa, P.C. Wu*, J.J. Wu, J.H. Su*, H.Y. Chen, Z.E. Ke, F.C. Chou, J.H. Hsu*, D.F. Cheng, Y.C. Chen, C.C. Lo, R.S. Liu, S.C. Heem, K.T. Tang, M.H. Chang*
*TSMC Corporate Research, Hsinchu, Taiwan; *National Tsing Hua University, Hsinchu, Taiwan
*Equally Credited Authors (ESCA)
- 2:20 PM
- 34.3 A 22nm 64k Lightning-Like Hybrid Computing-in-Memory Macro with Compressed Aider Tree and Analog-Surrogate Quantizers for Transformer and CNNs**
A. Guo, X. Chen, F. Ding, J. Chen, Z. Yuan*, H. Xu, Y. Zhang, J. Zhang, Y. Tang, Z. Zhang, G. Chen, D. Yang, Z. Zhang, L. Ren, T. Xiong, B. Wang, B. Liu, W. Shan, X. Liu, H. Cai, G. Sun, J. Yang, X. Gu
Southeast University, Nanjing, China
Peking University, Beijing, China
NIOUM, Beijing, China
- 2:45 PM
- 34.4 A 3nm 32.5TOPS/W, 55.0TOPS/mm² and 3.78Mbit/mm² Fully-Digital Computing-in-Memory Macro Supporting INT12 x INT12 with a Parallel-MAC Architecture and Friendly 6T SRAM 6H1 Cell**
H. Fujiwara, H. Mori, W.C. Zhao, K. Akashi, C.E. Lee, X. Peng, V. Joshi, C.H. Chuang, S.H. Hsu, T. Hattori, T. Kobayashi, D.H. Yen, H.Y. Liu, Y.C. Lai, C.C. Lee, F.H. Chou, K. Akamatsu, S. Adachi, Y. Wang, Y.D. Chou*, H.H. Chen, H.J. Liao, F.Y.J. Chang*
*TSMC, Hsinchu, Taiwan; *TSMC, San Jose, CA
*TSMC, Ottawa, Canada; *TSMC, Yokohama, Japan
- 3:00 PM
- 34.5 A 618-4094TOPS/W Capacitor-Reconfigured CIM Macro for Unified Acceleration of CNNs and Transformers**
K. Yoshida, Keio University, Yokohama, Japan
- Break 3:15 PM
- 3:30 PM
- 34.6 A 28nm 72.12Tflops/W Hybrid-Domain Outer-Product Based Floating-Point SRAM Computing-in-Memory Macro with Logarithm Bit-Width Residual ADC**
Y. Yuan*, Y. Song, X. Wang, J. Li, C. Ma*, G. Chen, M. Tang, X. Ren, Z. Hou, J. Zhu*, H. Wu*, G. Ren*, G. Xing*, P.H. Mak, F. Zhang*
*Institute of Microelectronics of the Chinese Academy of Sciences, Beijing, China
*University of Chinese Academy of Sciences, Beijing, China
*Beijing Institute of Technology, Beijing, China; *University of Macau, Macau, China
- 4:00 PM
- 34.7 A 28nm 2.4Mbit/mm² 6.9-16.3TOPS/mm² eDRAM-LUT-Based Digital-Computing-in-Memory Macro with In-Memory Encoding and Refreshing**
Y. Wu, S. Guo, X. Li, L. Liu, W. An, G. Ding, Y. Li, Z. Huang, Z. Gu, C. Han, X. Li, H. Yang, H. Xu, Y. Liu*

Session C3: AI/ML Accelerators and CiM 10:15 AM, Honolulu HI

VLSI 2024

Co-Chairs: Jaydeep Kulkarni, University of Texas at Austin
Masanao Yamaoka, Hitachi, Ltd.

10:15 AM

C3.1 122.7 TOPS/W Stdcell-Based DNN Accelerator Based on Transition Density Data Representation, Clock-Less MAC Operation, Pseudo-Sparsity Exploitation in 40 nm, Animesh Gupta¹, Japesh Vohra¹, Viveka Konandur Rajanna¹, Massimo Aliotti¹¹National University of Singapore

A DNN whose activation magnitude is represented by digital transition density is introduced for low energy, under the proposed Dyadic Digital Transition Modulation (DDTM). MAC operations are simplified into transition counting, enabling 1) activation pseudo-sparsity for lower energy, 2) clock-less neuron operation via simple un-down asynchronous counters. >100 TOPS/W in 40 nm stdcell design is on par with

Session 23: Neuromorphic Computing (NC) - Compute-in-Memory for Deep Learning 2:15 PM, Continental 6

Co-Chairs: Martin Frank, IBM and Dugyu Kuzum, University of California San Diego

This session describes advances in compute-in-memory (CiM) technologies and 3D integration for deep learning. Such circuits hold promise for deep learning inference and training by enabling faster and more energy-efficient neural network operations than digital CMOS. The session will be opened with a report on an in-memory compute chip fabricated in 14 nm CMOS technology that employs a carbon-based liner underneath the phase-change material to improve temporal stability and inference accuracy. The second paper reports on low-temperature monolithic 3D integration of carbon nanotube FETs and resistive RAM (ReRAM) arrays to

3:45 PM - 4:00 PM

Session 4A: PIM ACCELERATORS

Location: Pacifico B (floor map)

Session Chair: Jongse Park

3:45 PM - 4:00 PM

PreSto: An In-Storage Data Preprocessing System for Training Recommendation Models

Y. Lee, H. Kim, M. Rhu

4:00 PM - 4:15 PM

pSyncPIM: Partially Synchronous Execution of Sparse Matrix Operations for All-bank PIM Architectures

D. Baek, S. Hwang, J. Huh

4:15 PM - 4:30 PM

NDSearch: Accelerating Graph-Traversal-Based Approximate Nearest Neighbor Search through Near Data Processing

Y. Wang, S. Li, Q. Zheng, L. Song, Z. Li, A. Chang, H. Li, Y. Chen

4:30 PM - 4:45 PM

Enabling Efficient Large Recommendation Model Training with Near CXL Memory Processing

H. Liu, L. Zheng, Y. Huang, J. Zhou, C. Liu, R. Wang, X. Liao, H. Jin, J. Xue

4:45 PM - 5:00 PM

Exploiting Similarity Opportunity of Emerging AI Models on 3D Hybrid Bonding Architecture

ISSCC 2024

IEDM 2024

ISCA 2024

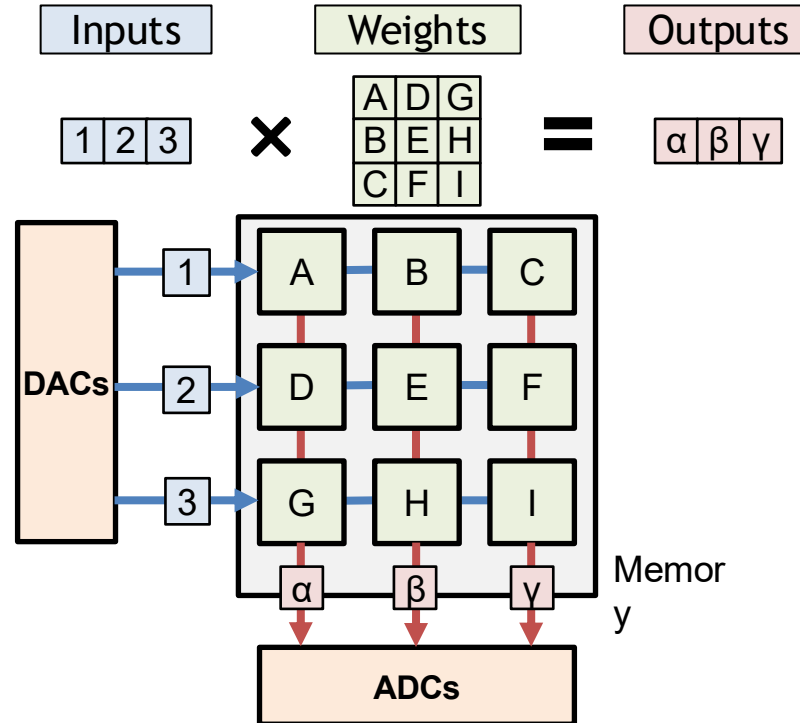
Many dedicated sessions on CiM at architecture, circuits, and devices conferences

CiM Research Spans Full Stack

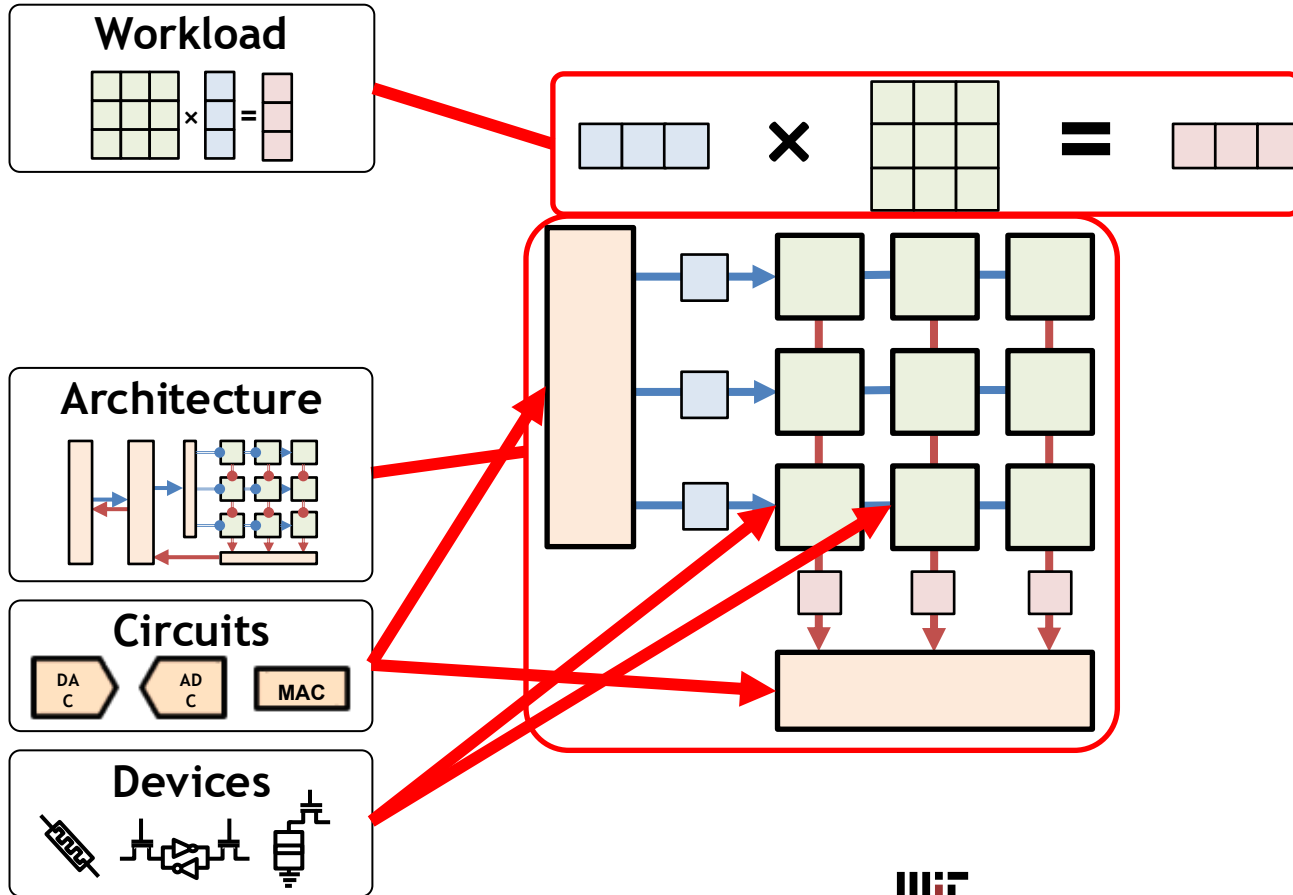
- **Devices:** The components forming each memory cell (e.g., SRAM, DRAM, ReRAM, STT-RAM)
- **Circuits:** The components performing computation, analog/digital conversion, storage, data movement, and other actions
- **Architecture:** The organization of components into a larger system (e.g., the number of each component and how components are connected)
- **Workload:** The DNN to be run, which we model as a series of extended-Einsum operations with tensors of varying shapes and values
- **Mapping:** The temporal and spatial scheduling of the workload onto the system

Need for modeling tool to enable apple-to-apple comparison and design space exploration → **CiMLoop** (used in Lab 5)

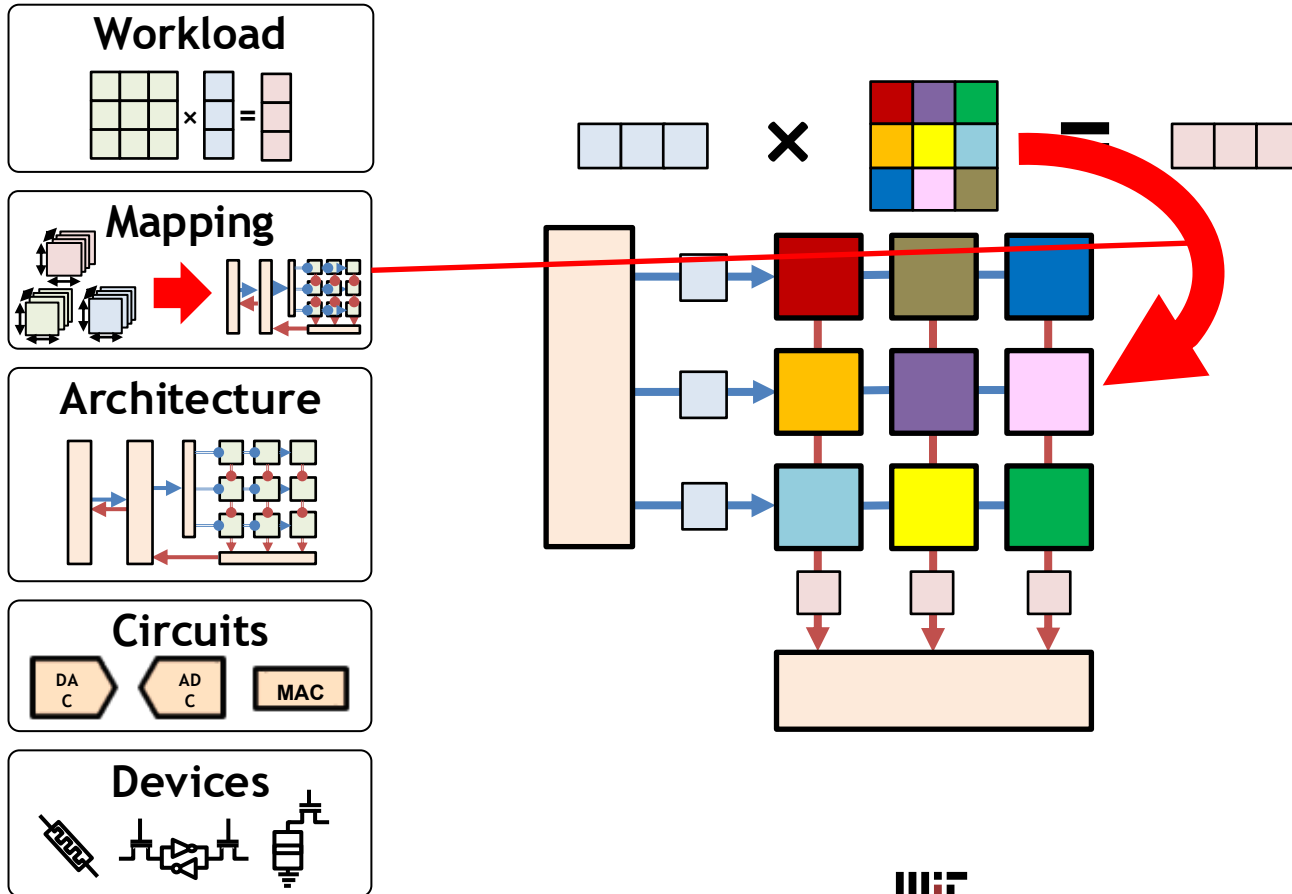
Compute In Memory (CiM) Accelerators



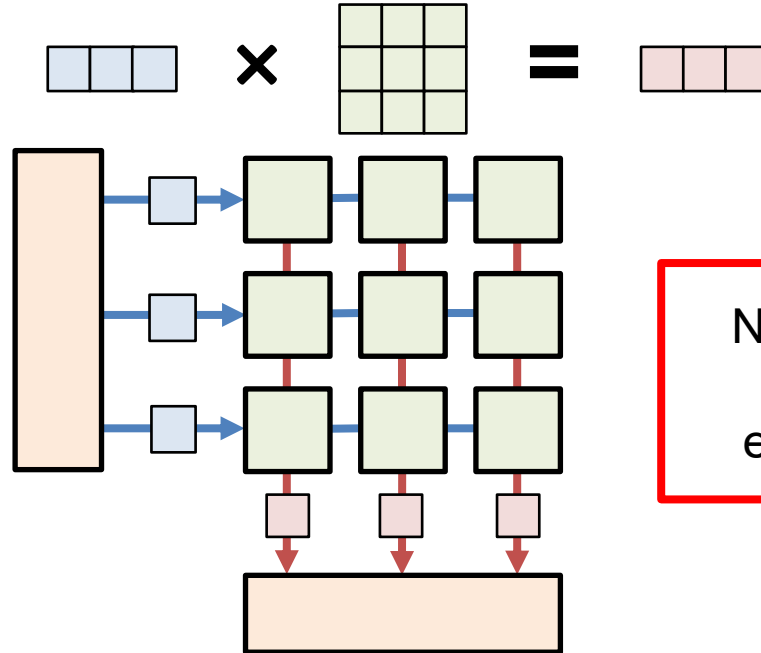
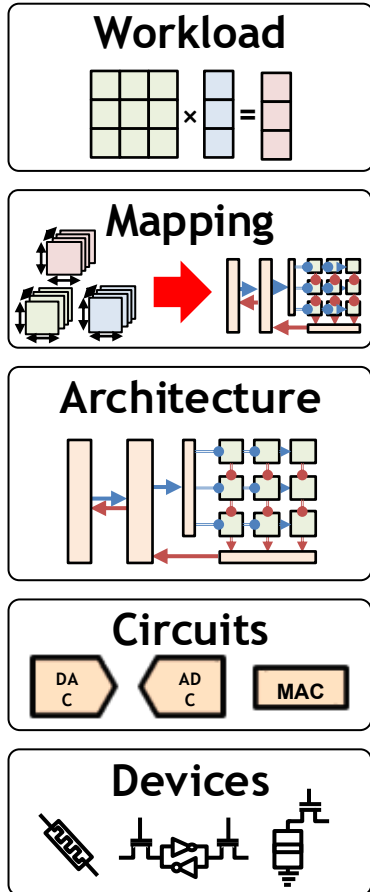
The CiM Stack: Large Design Space



The CiM Stack: Large Design Space



The CiM Stack: Large Design Space

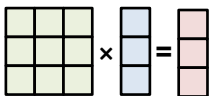


Need for modeling tool to enable design space exploration → **CiMLoop**

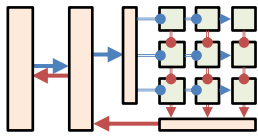
CiMLoop: A Flexible, Accurate, and Fast CiM Modeling Tool

Inputs

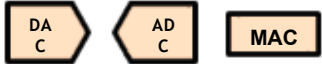
Workload



Architecture



Circuits



Devices



CiMLoop

Data-Distribution-Dependent Component Model

Data Distribution Calculation

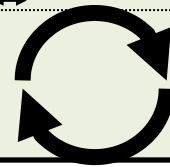
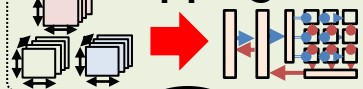
Data Distributions



Component Model Library

Timeloop + Accelergy

Mapping



Full-Stack Model

Outputs

System Energy,
Area,
Throughput

Code available at

<https://emze.csail.mit.edu/cimloop>

CiMLoop: A Flexible, Accurate, and Fast CiM Modeling Tool

- Flexibility to represent co-design space
 - **Challenge:** There are diverse choices at each level
 - **Solution:** Flexible user-defined specifications

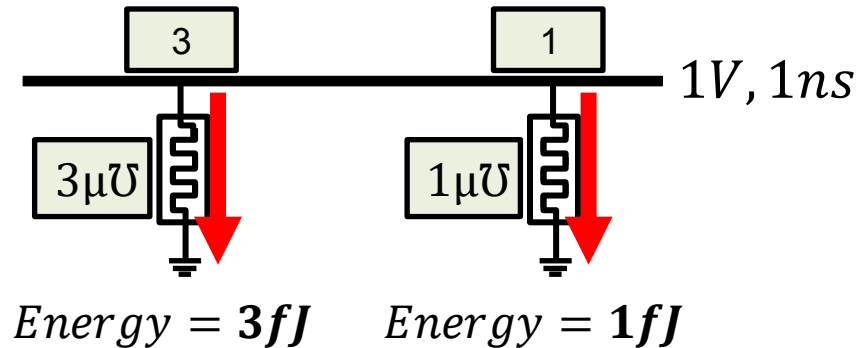
CiMLoop: A Flexible, Accurate, and Fast CiM Modeling Tool

- Flexibility to represent co-design space
 - **Challenge:** There are diverse choices at each level
 - **Solution:** Flexible user-defined specifications
- Accurately model energy
 - **Challenge:** Workload values and representation affect component energy
 - **Solution:** Models capture these cross-stack interactions (error within 8%)

Accurately Modeling Energy: Data-Value-Dependence

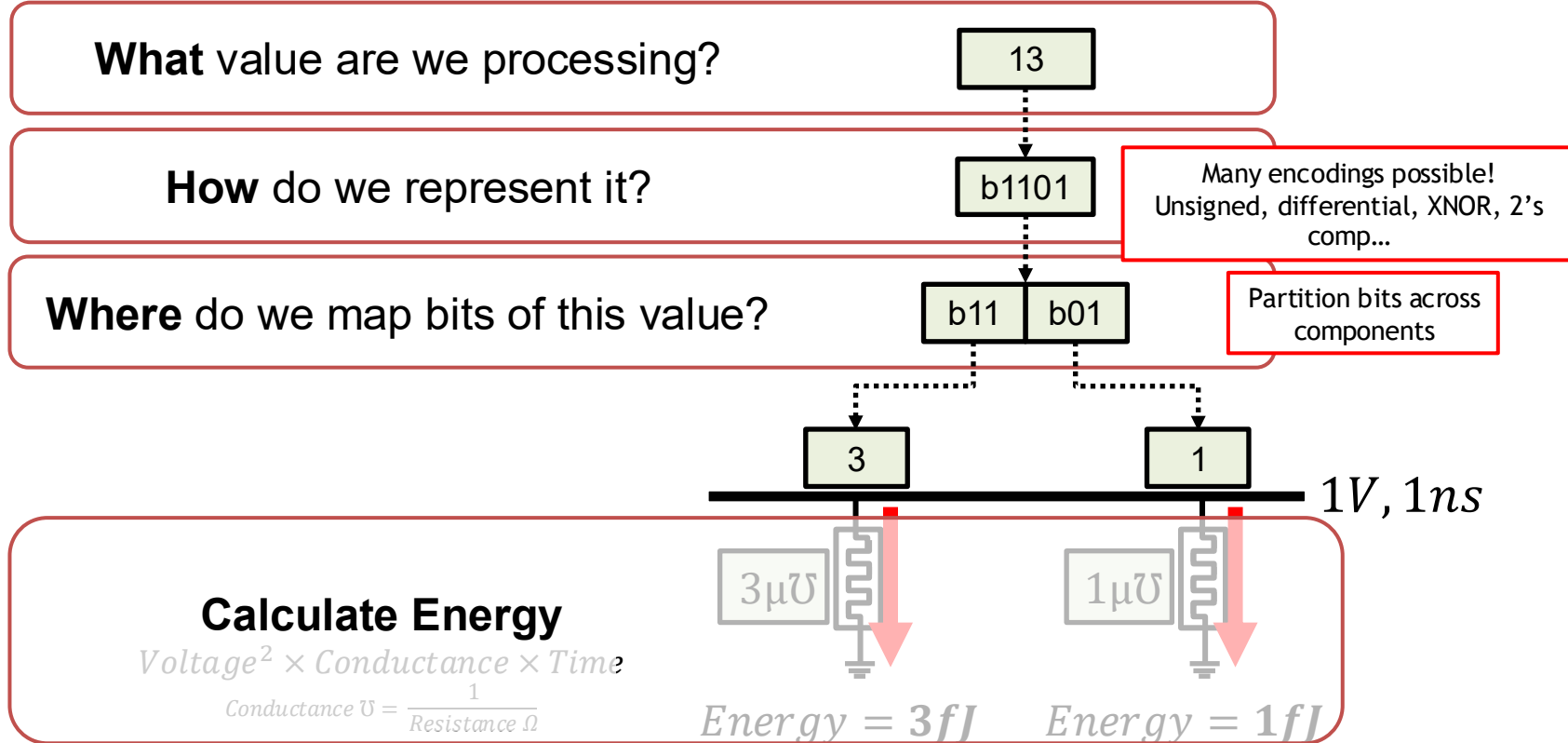
$$\text{Voltage}^2 \times \text{Conductance} \times \text{Time}$$

$$\text{Conductance } \mathcal{U} = \frac{1}{\text{Resistance } \Omega}$$



Data-value-dependence significantly impacts device and circuit energy
 Prior works assume fixed energy → significant error

Accurately Modeling Energy: Data-Value-Dependence



Capture data-value-dependence:

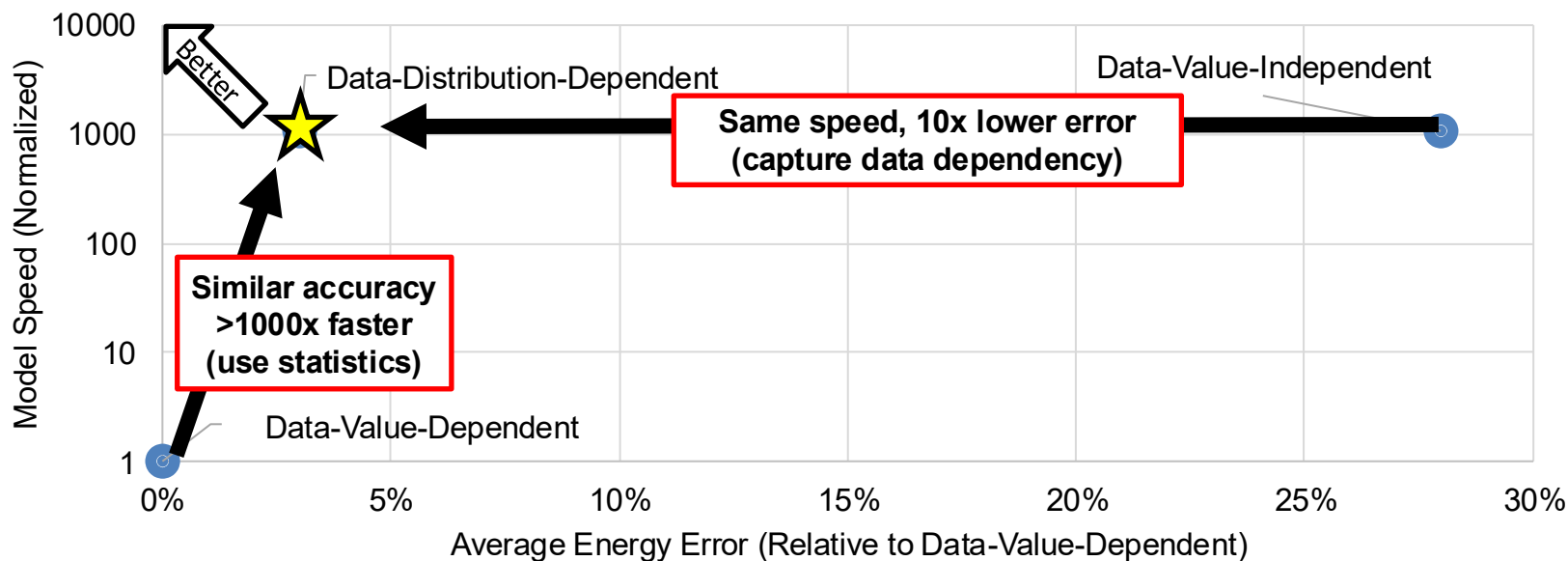
What values are there? **How** do we represent them? **Where** do we map their bits?

CiMLoop: A Flexible, Accurate, and Fast CiM Modeling Tool

- Flexibility to represent co-design space
 - **Challenge:** There are diverse choices at each level
 - **Solution:** Flexible user-defined specifications
- Accurately model energy
 - **Challenge:** Workload values and representation affect energy
 - **Solution:** Models capture these cross-stack interactions (error within 8%)
- Fast exploration of co-design space
 - **Challenge:** Accurate energy models may simulate many ($>10^{12}$) values
 - **Solution:** Statistical models that are 1000x faster than prior accurate models

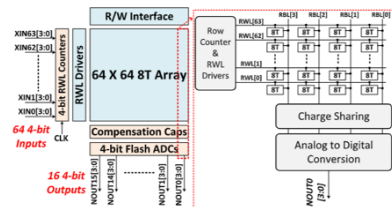
Fast and Accurate Statistical Energy Modeling

	Data-Value-Independent Timeloop [Parashar, <i>ISPASS</i> 2019]	Data-Value-Dependent NeuroSim [Peng, <i>TCAD</i> 2021]	Data-Distribution-Dependent CiMLoop [Andrulis, <i>ISPASS</i> 2024]
Model Accuracy	Low	High	High
Model Speed	High	Low	High

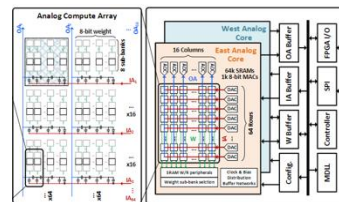


Example: Apples-to-Apples Comparison

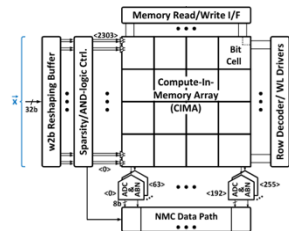
Macro



[Sinangil, JSSC 2021]



[Wang, VLSI 2022]



[Jia, JSSC 2020]

Technology Node

7nm

22nm

65nm

ADC Type

4b Flash

8b SAR

8b SAR

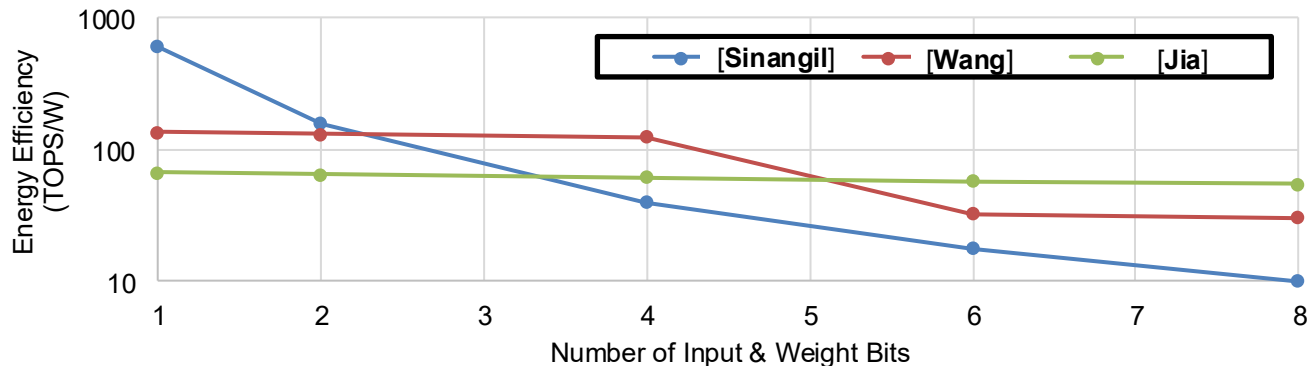
Memory Device

6T SRAM

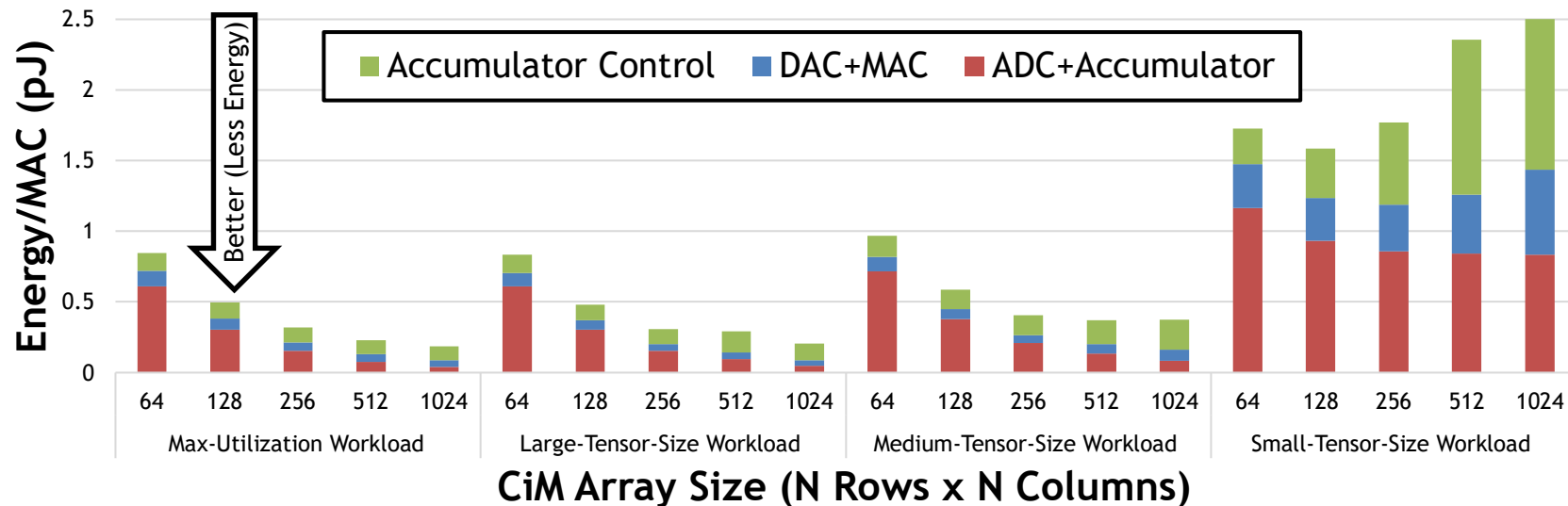
8T SRAM + Capacitor

6T SRAM

Compare Designs:
Same technology, ADC,
device for all macros



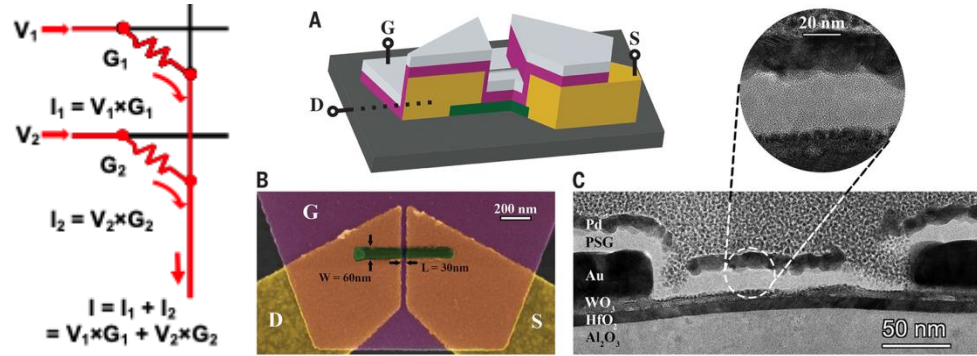
Example: Design Space Exploration



Explore array size (architecture) and DNN shapes (workload)

CiMLoop Enabling Collaborations Across Stack at MIT

Computing in memory with programmable resistive devices [w/ Jesus del Alamo]



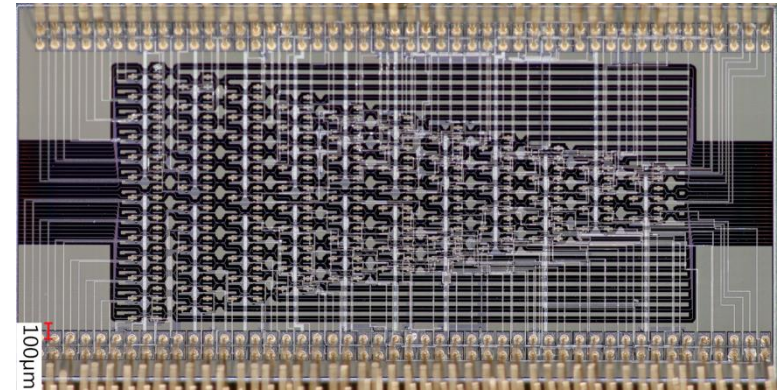
Ultra-low Power Superconducting Electronics for Deep Learning Accelerator Architectures: Evaluating Energy Efficiency and Scalability

L. Camron Blackburn, Evan Golden, Tanner Andrusis, Vivienne Sze, Joel Emer, Neil Gershenfeld, Karl K. Berggren

Sponsorship: MIT Lincoln Laboratory, the MIT AI Hardware Program

9.22

Computing with superconducting electronics
[w/ Karl Berggren and Neil Gershenfeld]
(Started from 6.5930/1 final project!)

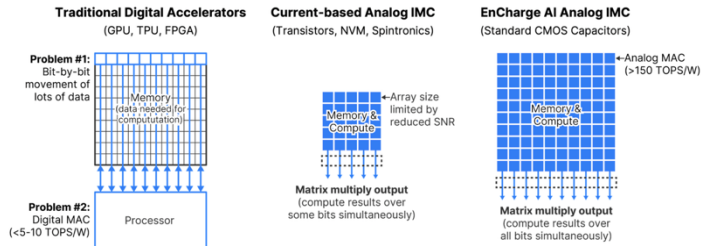


Modeling helps identify level of abstraction and facilitates communication between teams

Companies doing Analog CiM

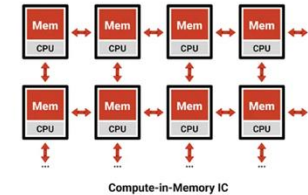
In-Memory Computing (IMC)

In-memory computing greatly enhances compute efficiency and reduces data movement.



Compute-in-Memory

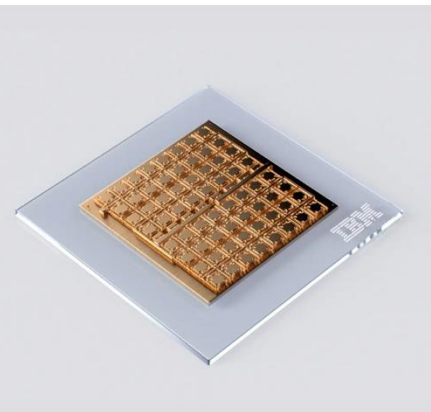
Boosting memory capacity and processing speed



Today's most common computing architectures are built on assumptions about how memory is accessed and used. These systems assume that the full memory space is too large to fit on-chip near the processor, and that we do not know what memory will be needed at what time. To address

Analog AI

Making Deep Neural Network systems more capable and energy-efficient.



Compute with Light

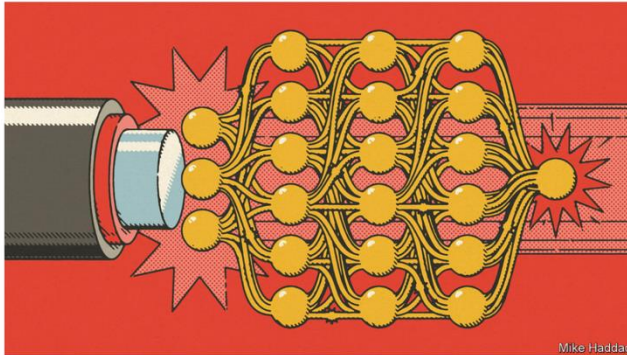
The Economist

Menu Weekly edition Search

Science & technology | Information technology

Artificial intelligence and the rise of optical computing

Photonic data-processing is well-suited to the age of deep learning



Mike Haddad

Dec 20th 2022

Save Share Give

MODERN INFORMATION technology (IT) relies on division of labour. Photons carry data around the world and electrons process them. But, before optical fibres, electrons did both—and some people hope to complete the transition by having photons process data as well as carrying them.

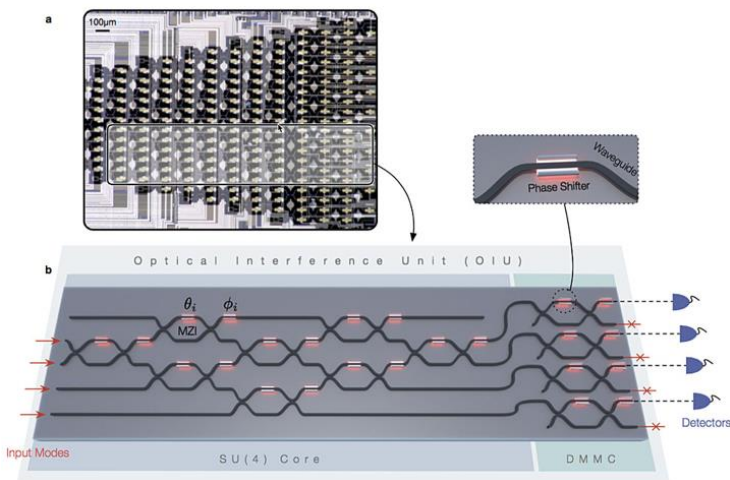
“Unlike electrons, photons (which are electrically neutral) can cross each others’ paths without interacting, so glass fibres can handle many simultaneous signals in a way that copper wires cannot. **An optical computer could likewise do lots of calculations at the same time.** Using photons reduces power consumption, too. **Electrical resistance generates heat, which wastes energy. The passage of photons through transparent media is resistance-free.**”

<https://www.economist.com/science-and-technology/2022/12/20/artificial-intelligence-and-the-rise-of-optical-computing>

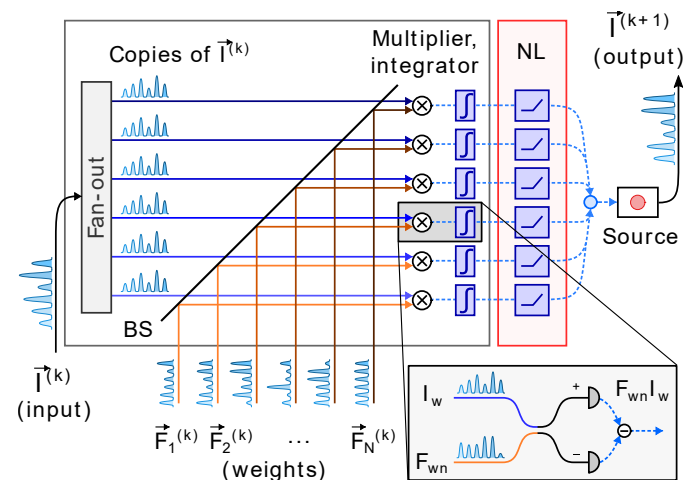
Compute with Light

Matrix Multiplication in the Optical Domain

- Cost of moving a photon can be **independent** of distance
- Multiplication can be performed **passively**



[Shen, *Nature Photonics* 2017]



[Bernstein, *CLEO* 2020]

Compute with Light

WILL KNIGHT BUSINESS 03.10.2021 07:00 AM

This Chip for AI Works Using Light, Not Electrons

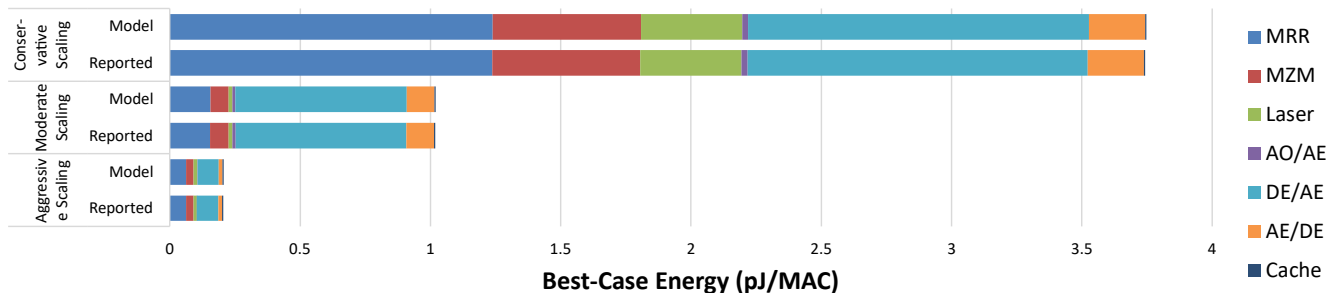
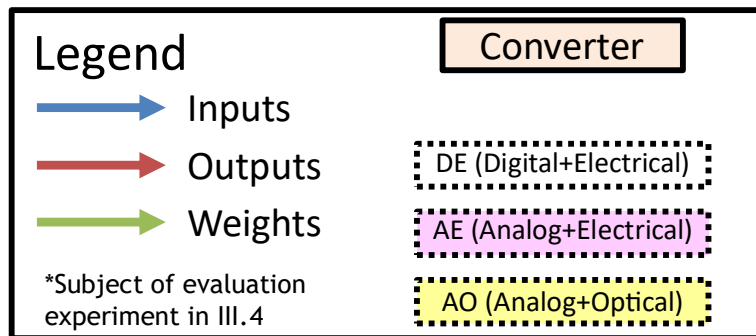
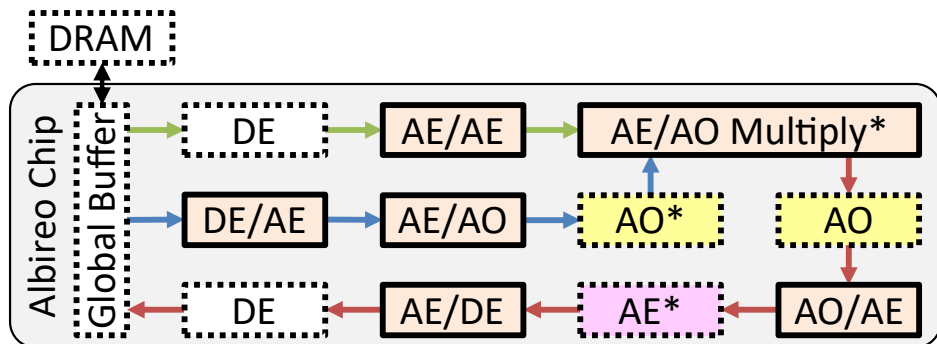
Lightmatter says the computing and power demands of complex neural networks need new technologies like these to keep up.



“...chip runs **1.5 to 10 times** faster than a top-of-the-line Nvidia A100 AI chip, Running a natural language model called BERT, for example, Lightmatter says Enviser is **five times faster** than the Nvidia chip; it also consumes **one-sixth of the power**”

<https://www.wired.com/story/chip-ai-works-using-light-not-electrons/>

CiMLoop for Photonics Modeling



Previous final projects in 6.5930 involved modeling and validation

Summary

- Cross-layer design critical for providing additional efficiency improvements
- For DNN processing using Advanced Technologies, it is important to factor in device and circuit limitations into the architecture
- Textbook Chapter 10
 - <https://doi.org/10.1007/978-3-031-01766-7>
- Other References
 - Y. N. Wu, V. Sze, J. S. Emer, “An Architecture-Level Energy and Area Estimator for Processing-In-Memory Accelerator Designs,” *IEEE International Symposium on Performance Analysis of Systems and Software (ISPASS)*, April 2020 [paper [PDF](#) | code [github](#)]
 - T. Andrulis, J. Emer, V. Sze, “RAELLA: Reforming the Arithmetic for Efficient, Low-Resolution, and Low-Loss Analog PIM: No Retraining Required!,” *International Symposium on Computer Architecture (ISCA)*, June 2023 [[PDF](#)]
 - T. Andrulis, J. Emer, V. Sze, “CiMLoop: A Flexible, Accurate, and Fast Compute-In-Memory Modeling Tool,” *IEEE International Symposium on Performance Analysis of Systems and Software (ISPASS)*, May 2024