

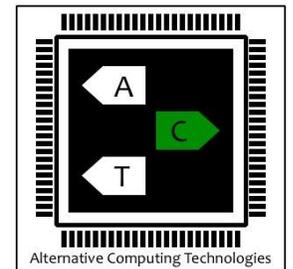
# BiHiWE: Mixed-Signal Charge-Domain Acceleration

**Soroush Ghodrati**

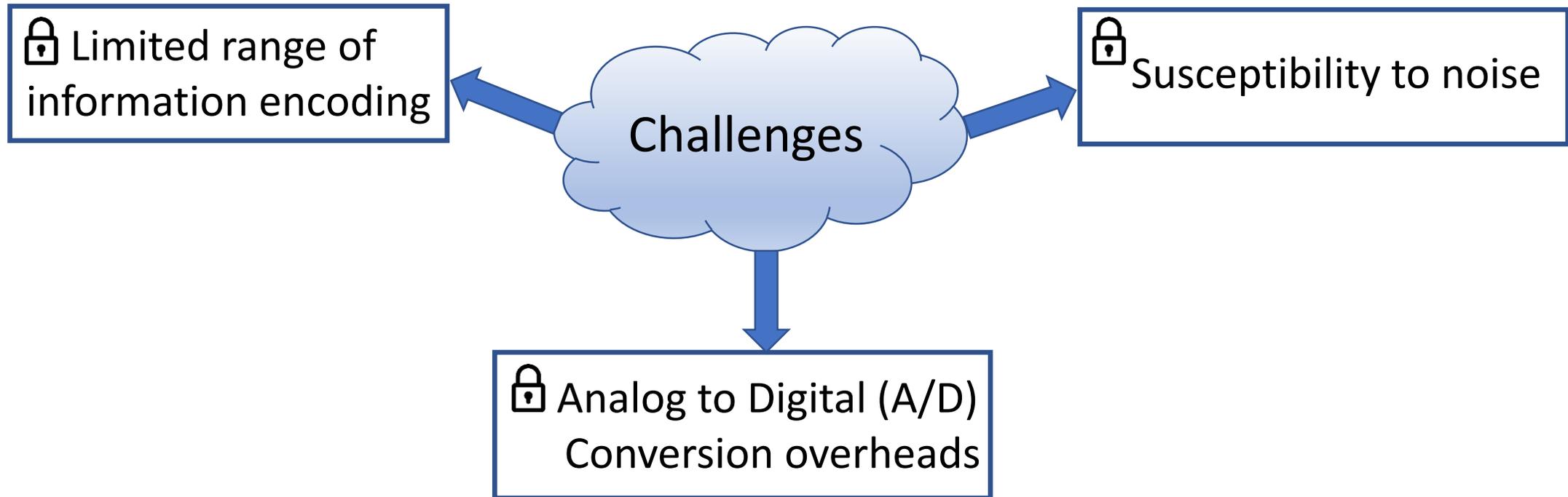
**PhD Student at Computer  
Science and Engineering**

Alternative Computing Technologies (**ACT**) Lab

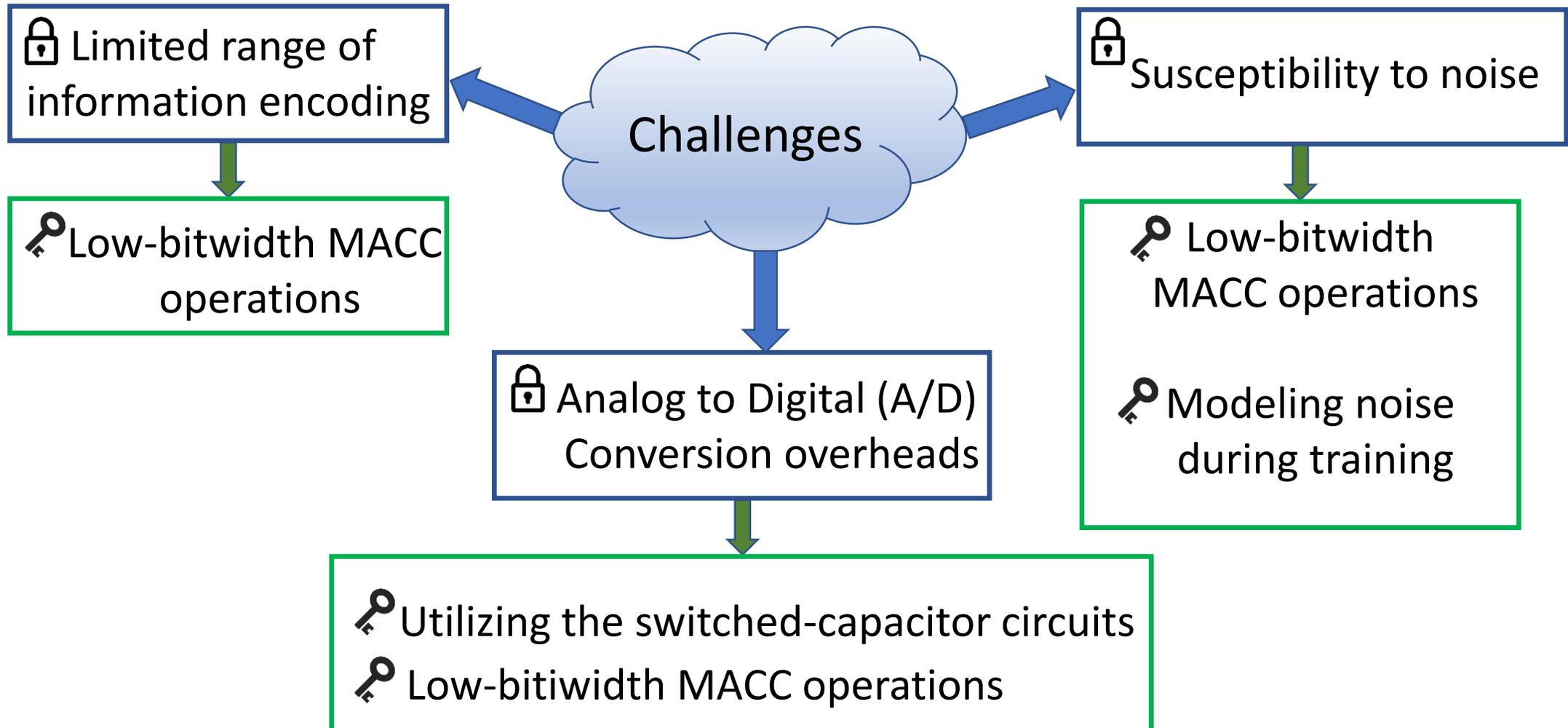
University of California, San Diego



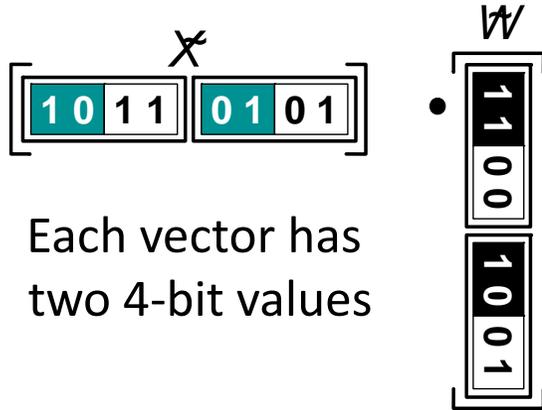
# Challenges in Analog Computing



# Challenges in Analog Computing

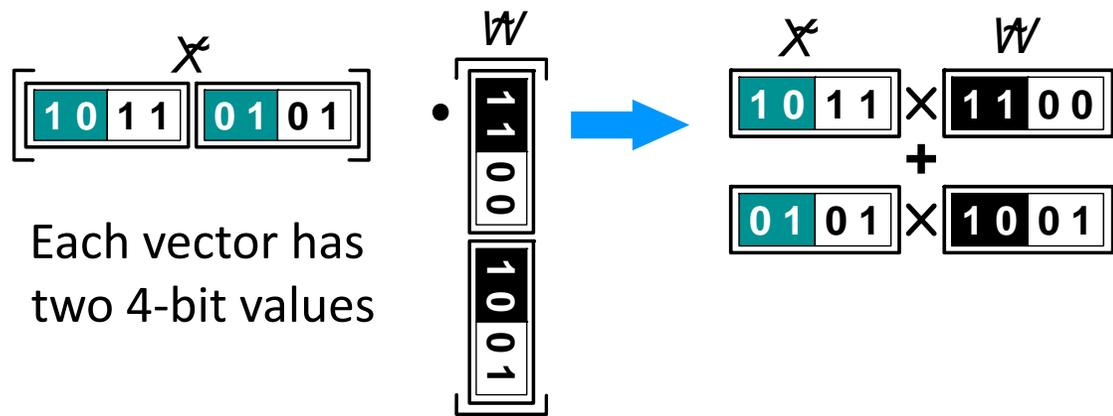


# Our Approach: Wide, Interleaved, and Bit-Partitioned Arithmetic



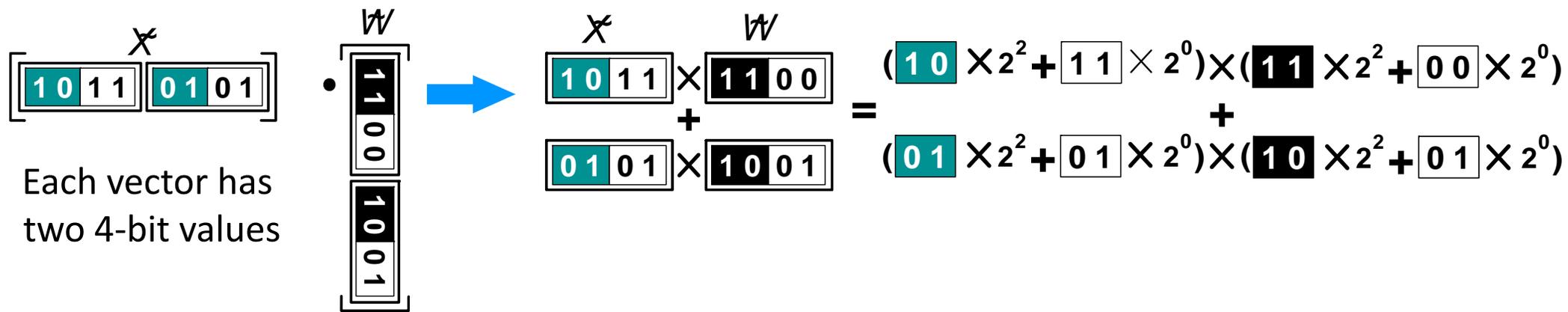
Vector dot-products can be **bit-partitioned** into groups of low bitwidth operations

# Our Approach: Wide, Interleaved, and Bit-Partitioned Arithmetic



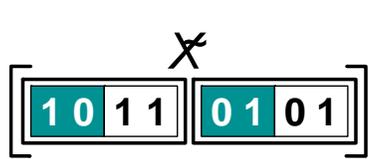
Vector dot-products can be **bit-partitioned** into groups of low bitwidth operations

# Our Approach: Wide, Interleaved, and Bit-Partitioned Arithmetic

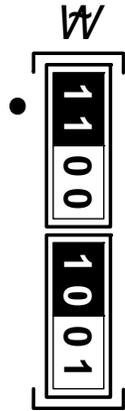


Vector dot-products can be **bit-partitioned** into groups of low bitwidth operations

# Our Approach: Wide, Interleaved, and Bit-Partitioned Arithmetic



Each vector has  
two 4-bit values



$$\begin{array}{l}
 \begin{array}{|c|c|} \hline 1011 \\ \hline \end{array} \times \begin{array}{|c|c|} \hline 1100 \\ \hline \end{array} = (10 \times 2^2 + 11 \times 2^0) \times (11 \times 2^2 + 00 \times 2^0) \\
 + \\
 \begin{array}{|c|c|} \hline 0101 \\ \hline \end{array} \times \begin{array}{|c|c|} \hline 1001 \\ \hline \end{array} = (01 \times 2^2 + 01 \times 2^0) \times (10 \times 2^2 + 01 \times 2^0)
 \end{array}$$

$$\begin{array}{l}
 = \\
 \left( \begin{array}{|c|c|} \hline 10 \times 11 \\ \hline 01 \times 10 \end{array} \right) \ll (2+2) \quad \left( \begin{array}{|c|c|} \hline 10 \times 00 \\ \hline 01 \times 01 \end{array} \right) \ll (2+0) \\
 + \\
 \left( \begin{array}{|c|c|} \hline 11 \times 11 \\ \hline 01 \times 10 \end{array} \right) \ll (0+2) \quad \left( \begin{array}{|c|c|} \hline 11 \times 00 \\ \hline 01 \times 01 \end{array} \right) \ll (0+0)
 \end{array}$$

100110001



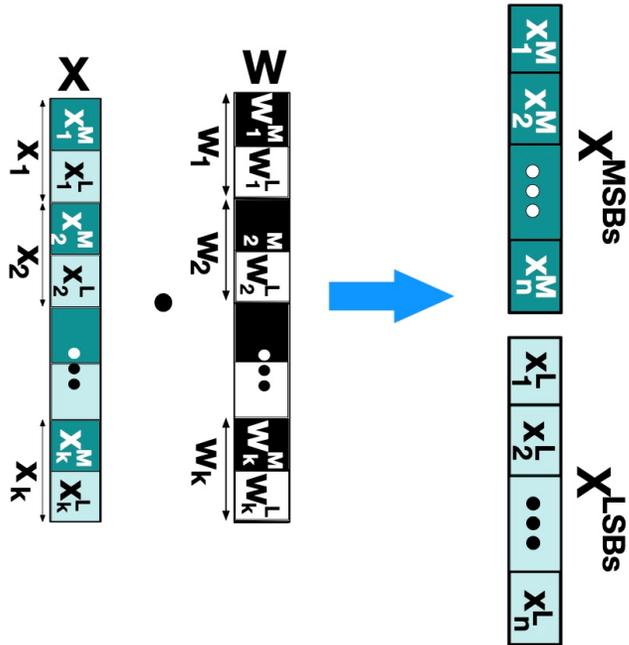
Vector dot-products can be  
**bit-partitioned** into groups  
of low bitwidth operations

# Wide, Interleaved, and Bit-Partitioned Vector Dot-Product

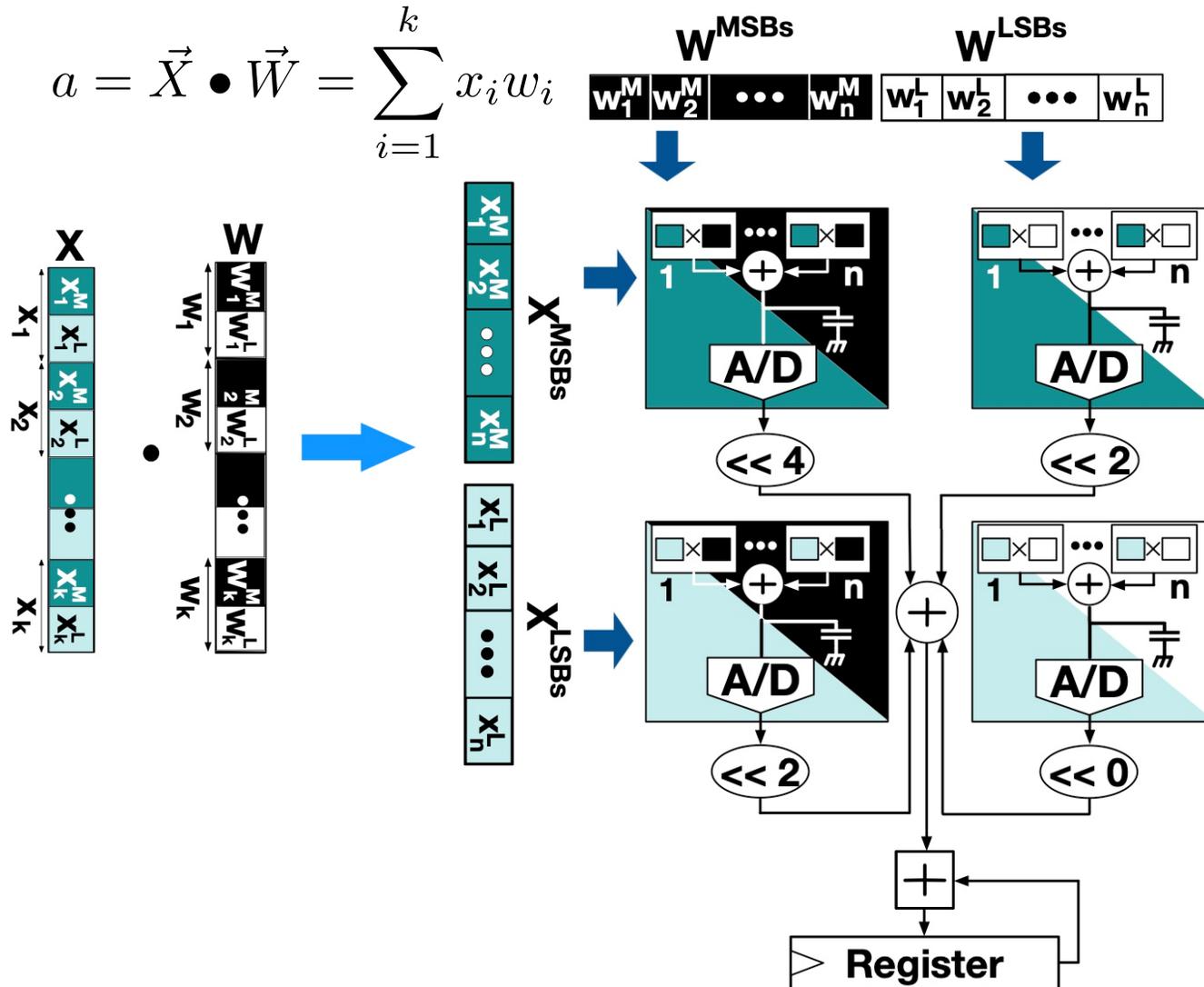
$$a = \vec{X} \bullet \vec{W} = \sum_{i=1}^k x_i w_i$$

$\mathbf{W}^{\text{MSBs}}$   
 $w_1^M \ w_2^M \ \dots \ w_n^M$

$\mathbf{W}^{\text{LSBs}}$   
 $w_1^L \ w_2^L \ \dots \ w_n^L$

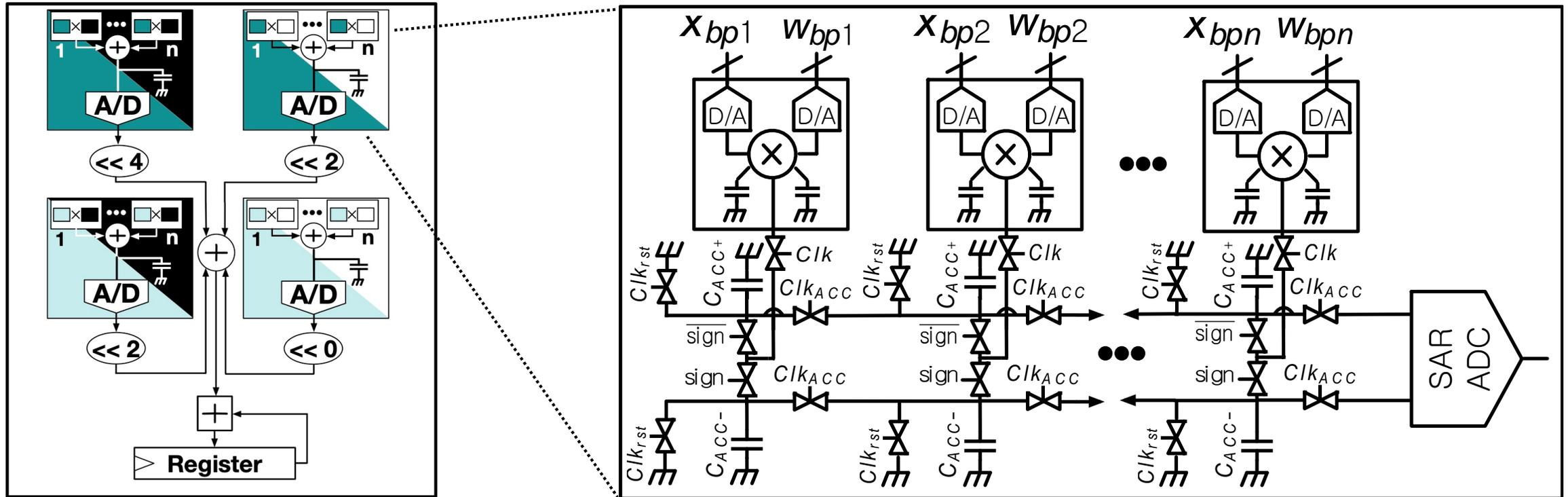


# Wide, Interleaved, and Bit-Partitioned Vector Dot-Product



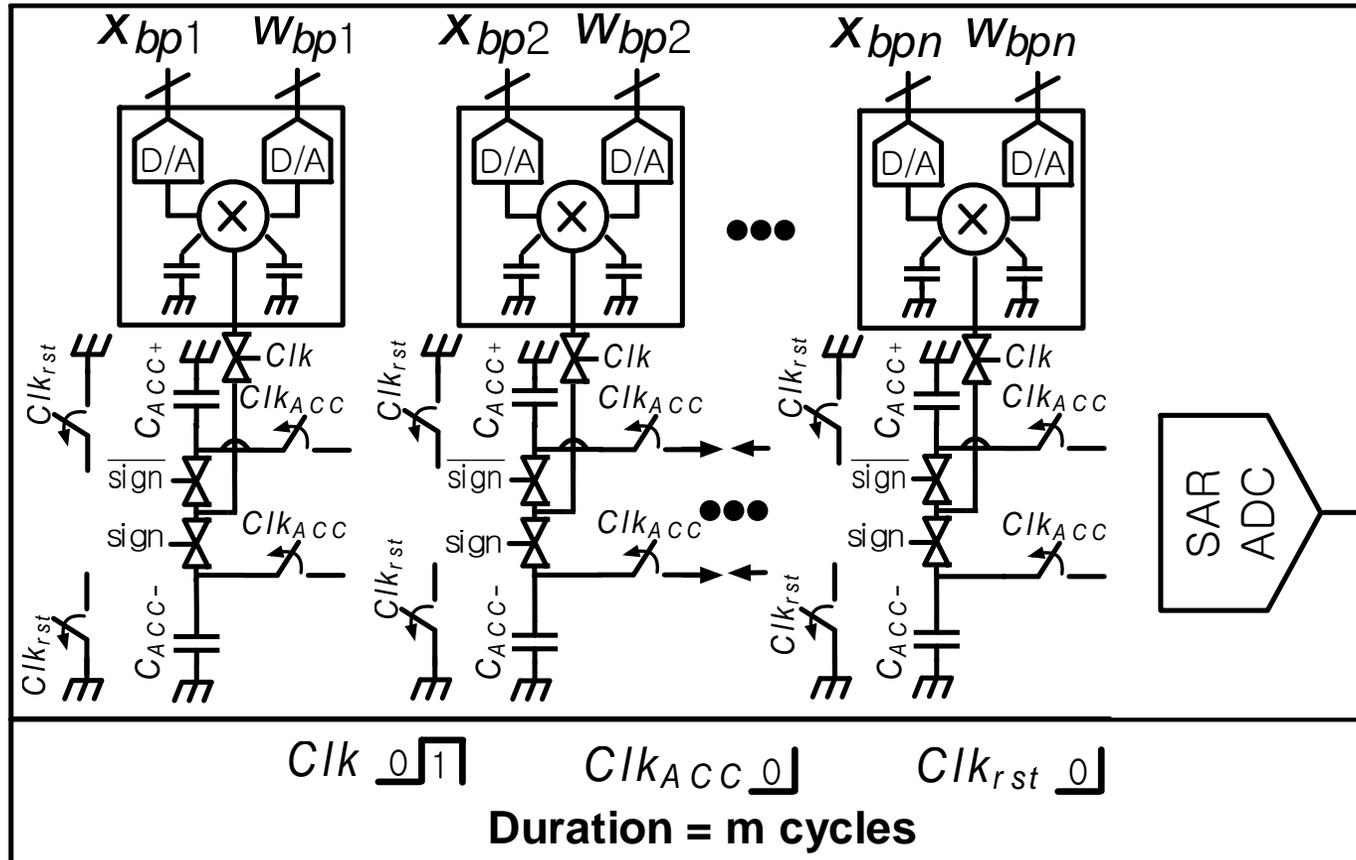
Using low-bitwidth operands provides larger headroom between value encoding in analog domain and reduces the energy/area overhead of A/D and D/A converters

# BiHiWE Microarchitecture



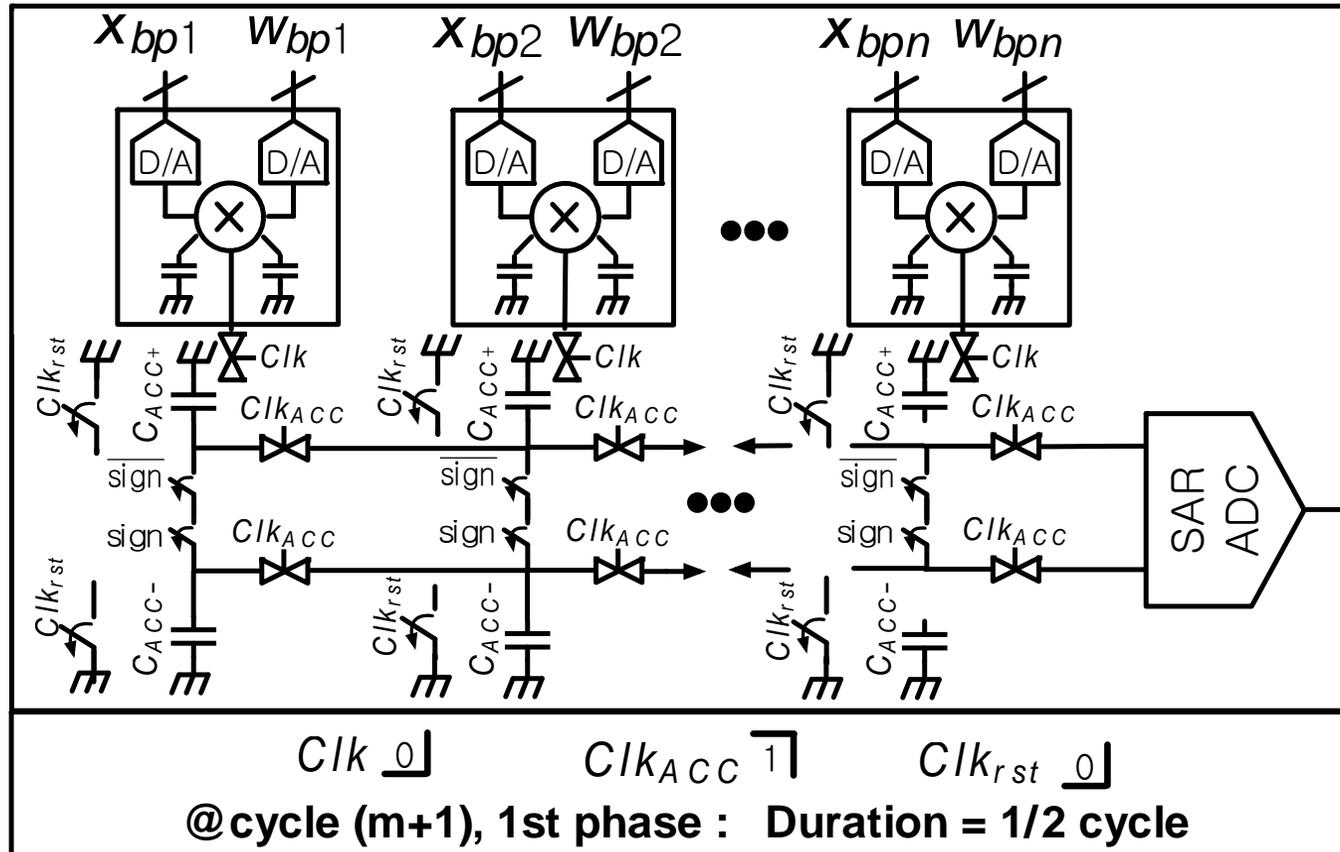
Mixed-Signal Bit-Partitioned MACC Array:  
A **wide** array of **low-bitwidth** MACC units share **single** A/D converter

# BiHIWE Microarchitecture



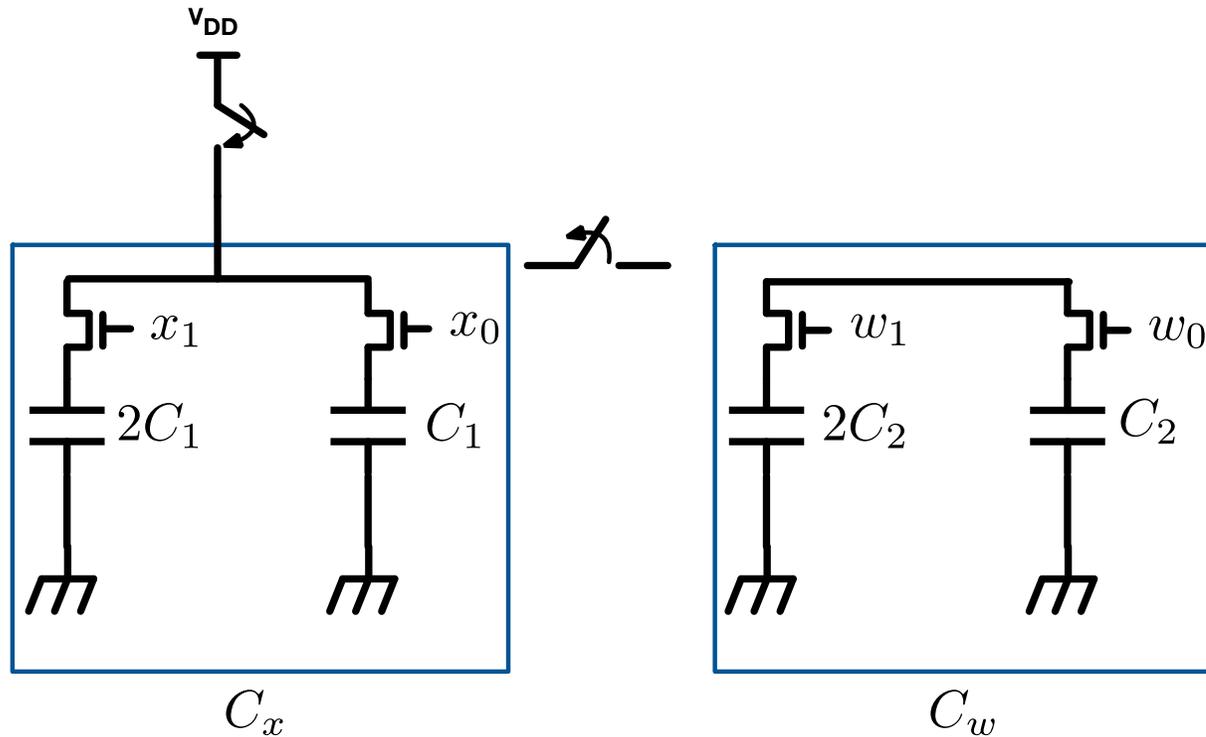
Mixed-Signal Bit-Partitioned MACC Array: **MACC**  
Operations and Private Accumulation

# BiHIWE Microarchitecture



Mixed-Signal Bit-Partitioned MACC Array:  
Accumulating across MACCs and starting A/D conversion

# Low-Bitwidth Switched-Capacitor MACC

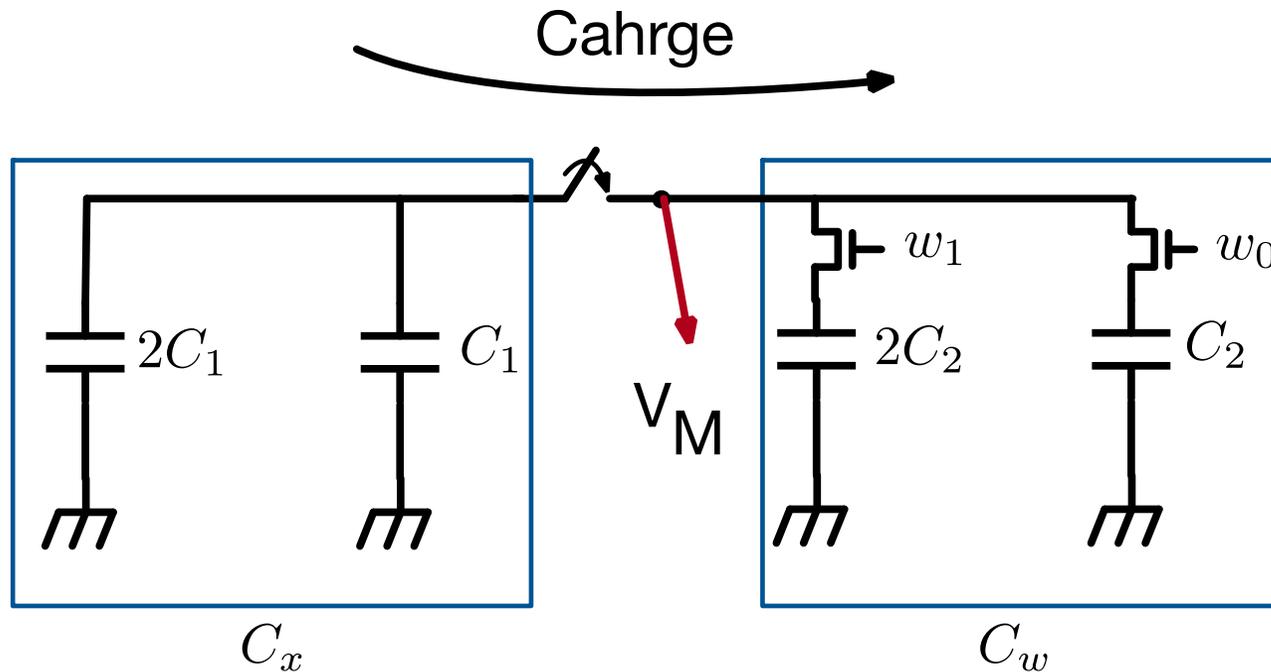


$$Q_x = C_x V_{DD}$$

$$Q_x \propto |X| C_1$$

A charge **proportional** to the magnitude of  $X$  is stored on  $C_x$

# Low-Bitwidth Switched-Capacitor MACC



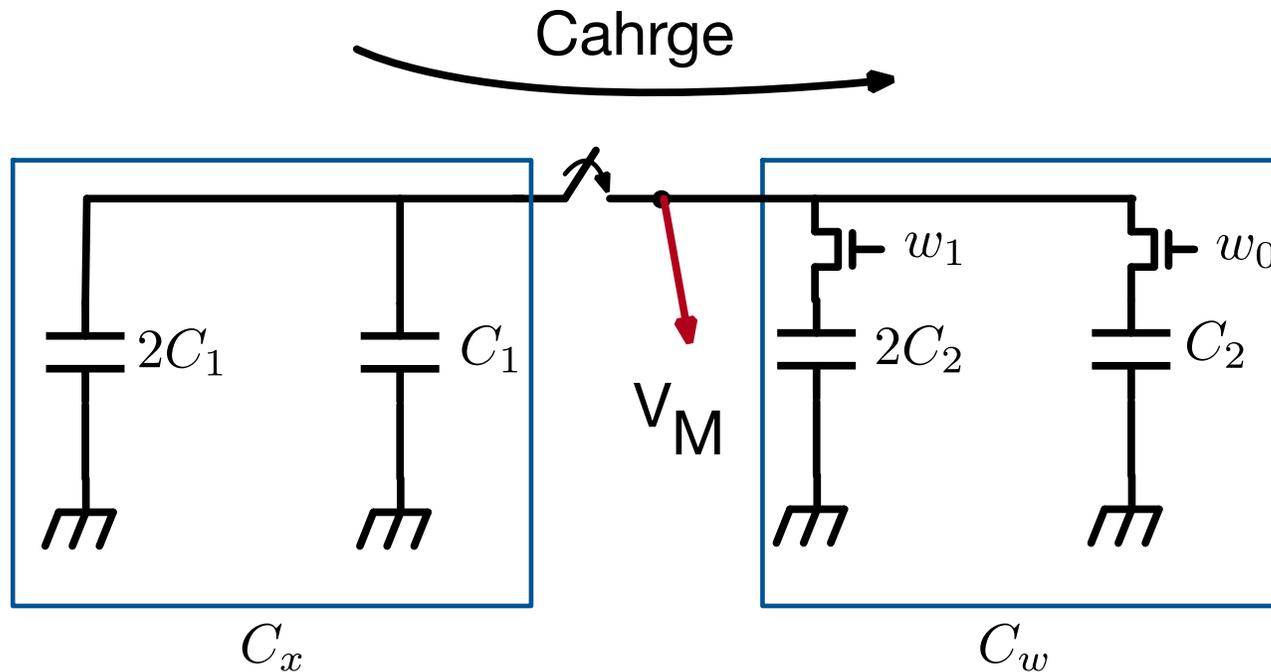
$$V_M = \frac{Q_x}{C_{eq}} = \frac{|X|C_1V_{DD}}{3C_1 + |W|C_2}$$

$$Q_w = C_w V_M = |W|C_2 V_M$$

$$Q_w = |X||W| \frac{C_1 C_2 V_{DD}}{3C_1 + |W|C_2}$$

The sampled charge by  $C_x$  is **shared** by  $C_w$ .

# Low-Bitwidth Switched-Capacitor MACC



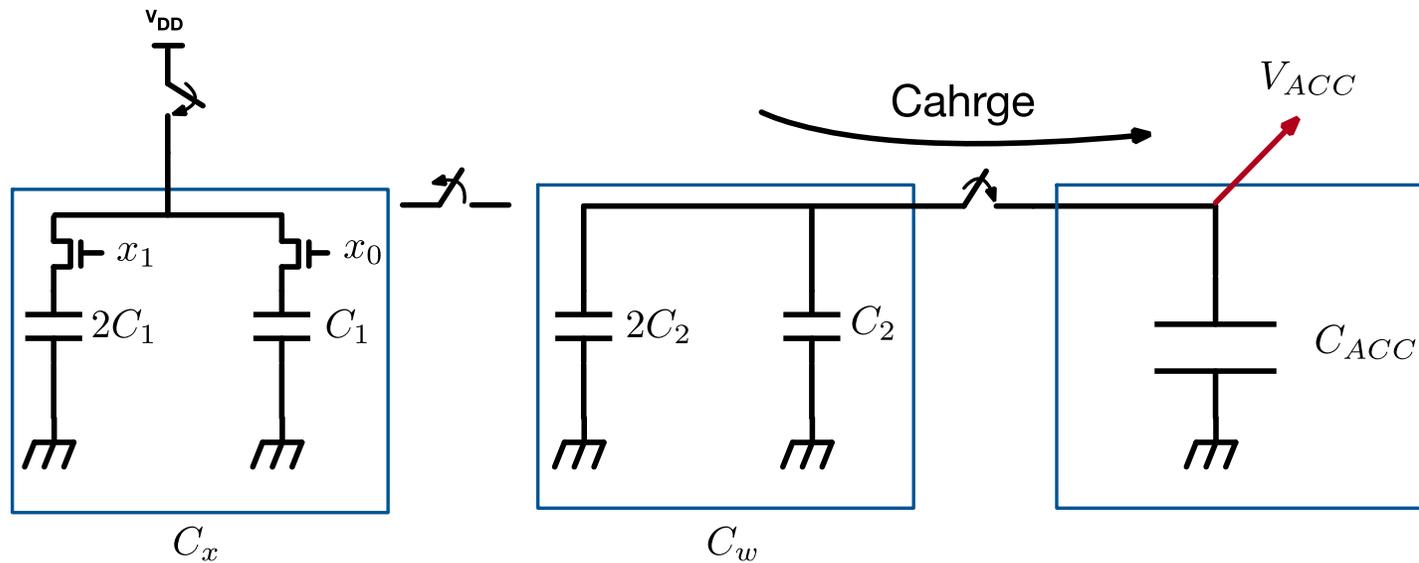
$$C_1 \gg C_2 \quad \frac{C_1}{C_2} \nearrow \infty$$

$$Q_w = |X||W| \frac{C_2 V_{DD}}{3}$$

$$Q_w \propto |X||W|$$

A charge **proportional** to  $|X||W|$  is stored on  $C_w$  and a **multiplication** happens

# Low-Bitwidth Switched-Capacitor MACC

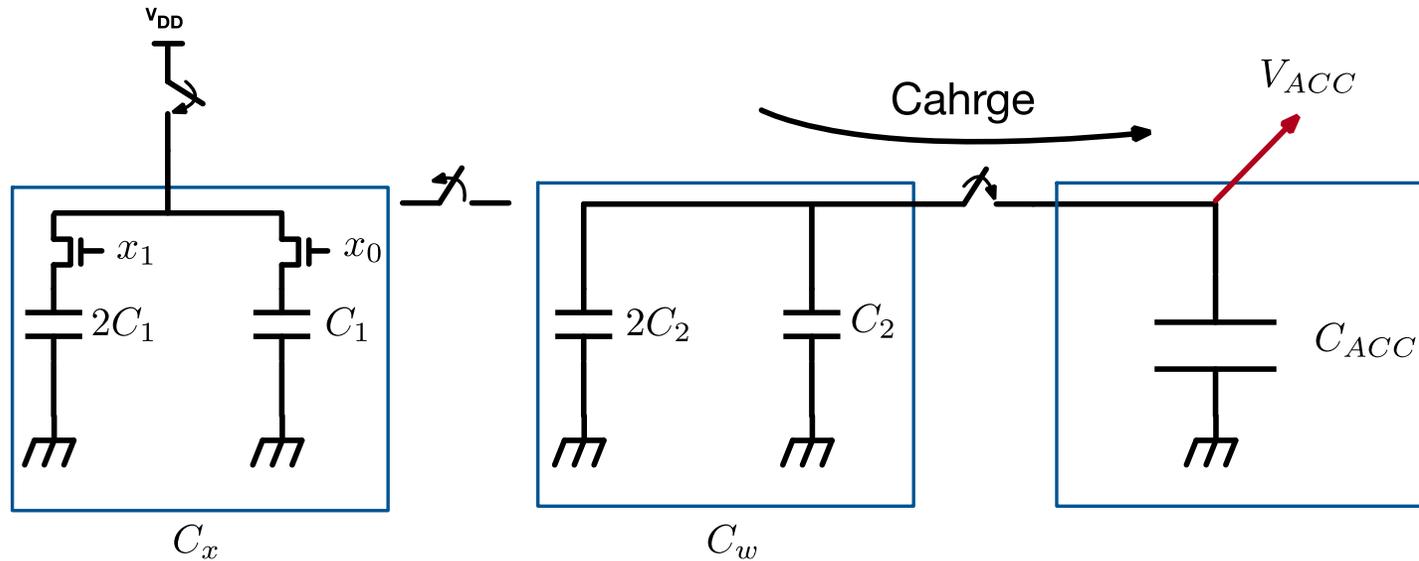


$$C_{ACC} \gg C_2$$

$$Q_{ACC} = \frac{C_{ACC}}{C_{ACC} + 3C_2} Q_w$$

The sampled charge by  $C_w$  is transferred to  $C_{ACC}$  and accumulated there

# Low-Bitwidth Switched-Capacitor MACC

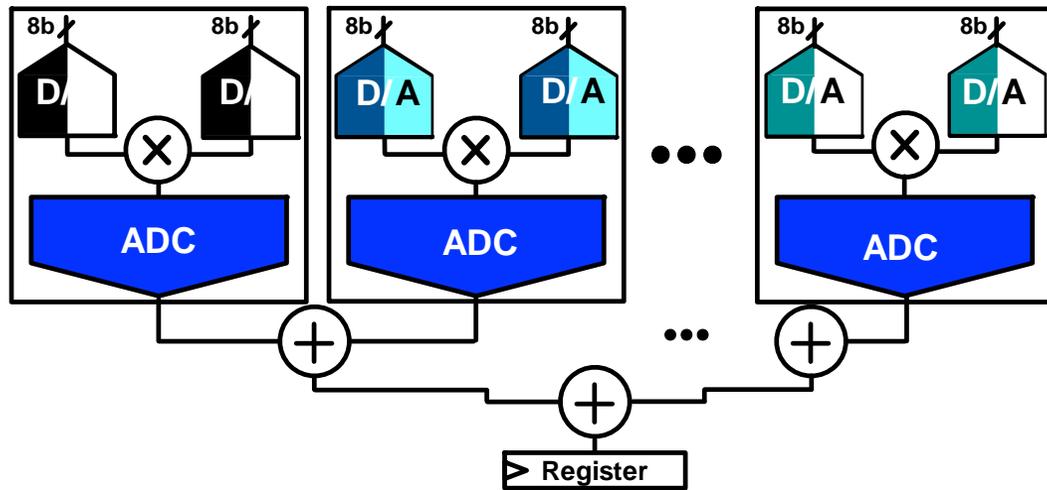


$$V_{ACC} = \frac{Q_{ACC}}{C_{ACC}} \approx \frac{Q_w}{C_{ACC}}$$

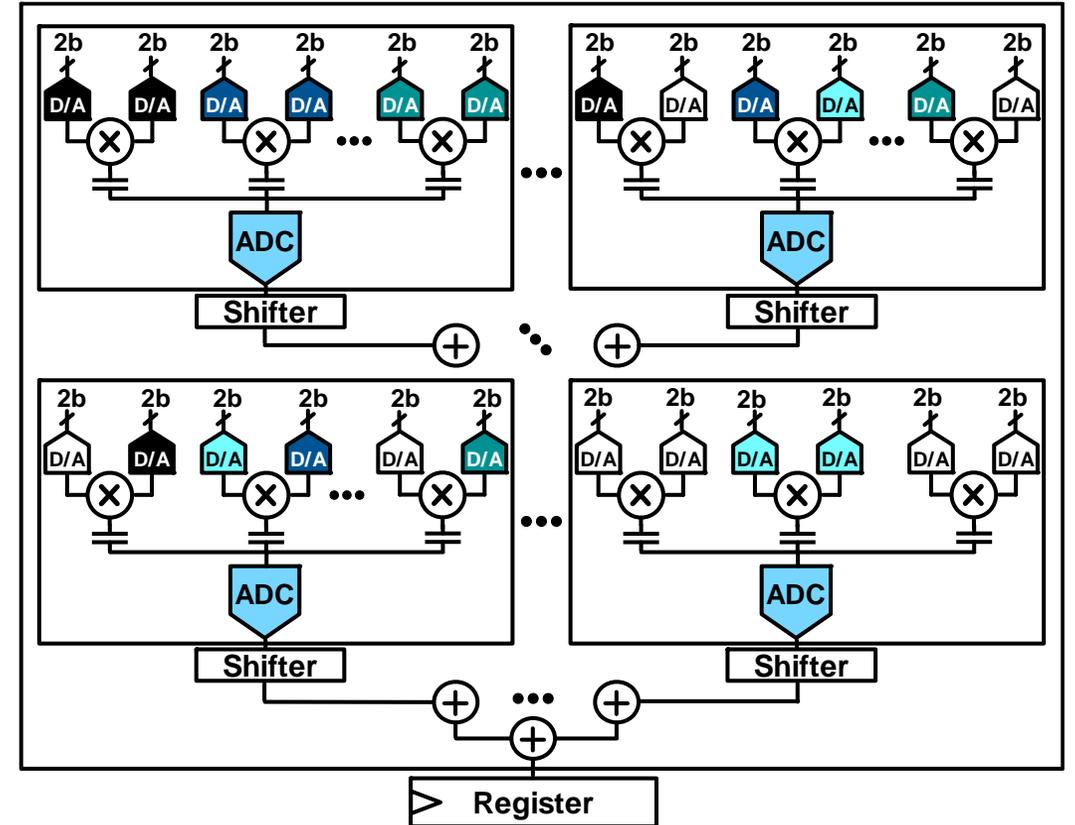
$$V_{ACC} = |X||W| \frac{C_2 V_{DD}}{3C_{ACC}}$$

While the result of the multiplication is being accumulated,  
a new input is sampled and a new round begins

# BiHiWE Microarchitecture

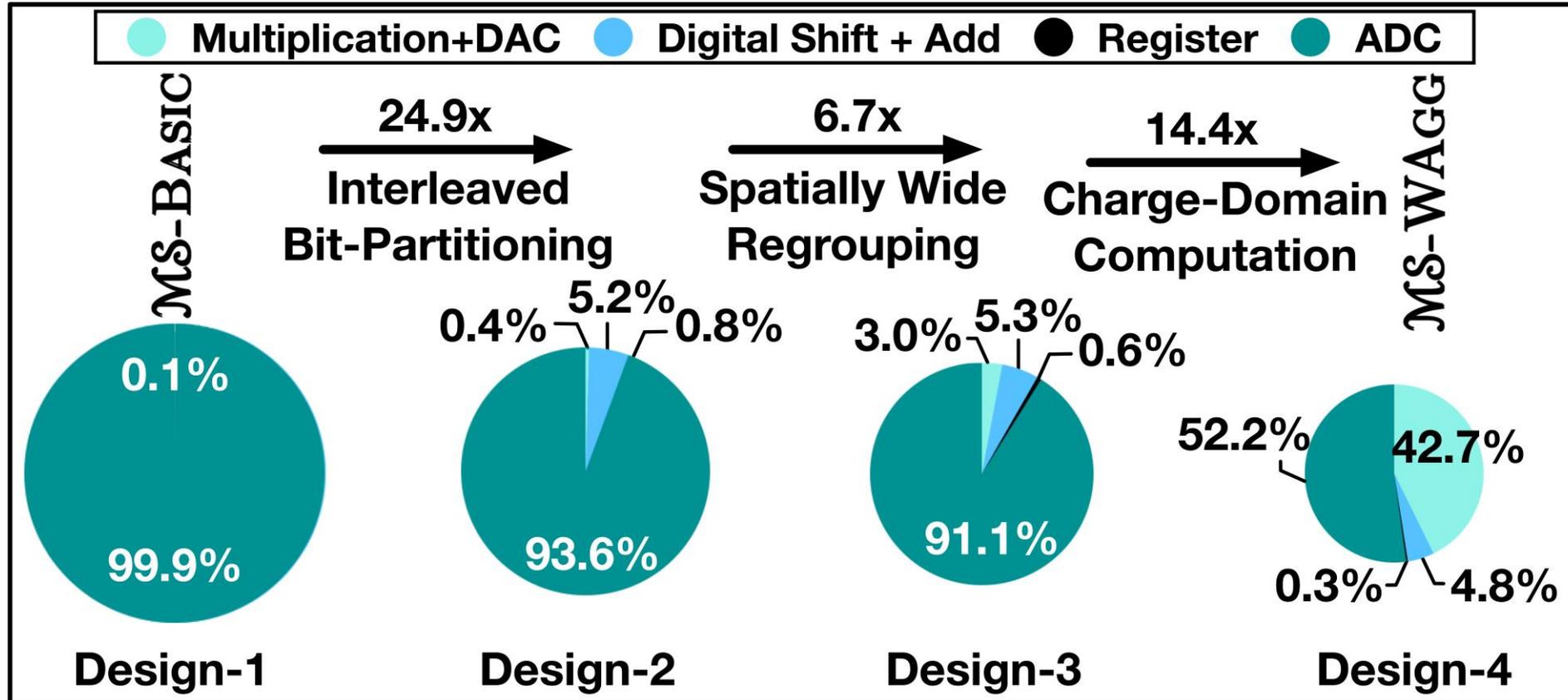


Basic Dot-Product Engine



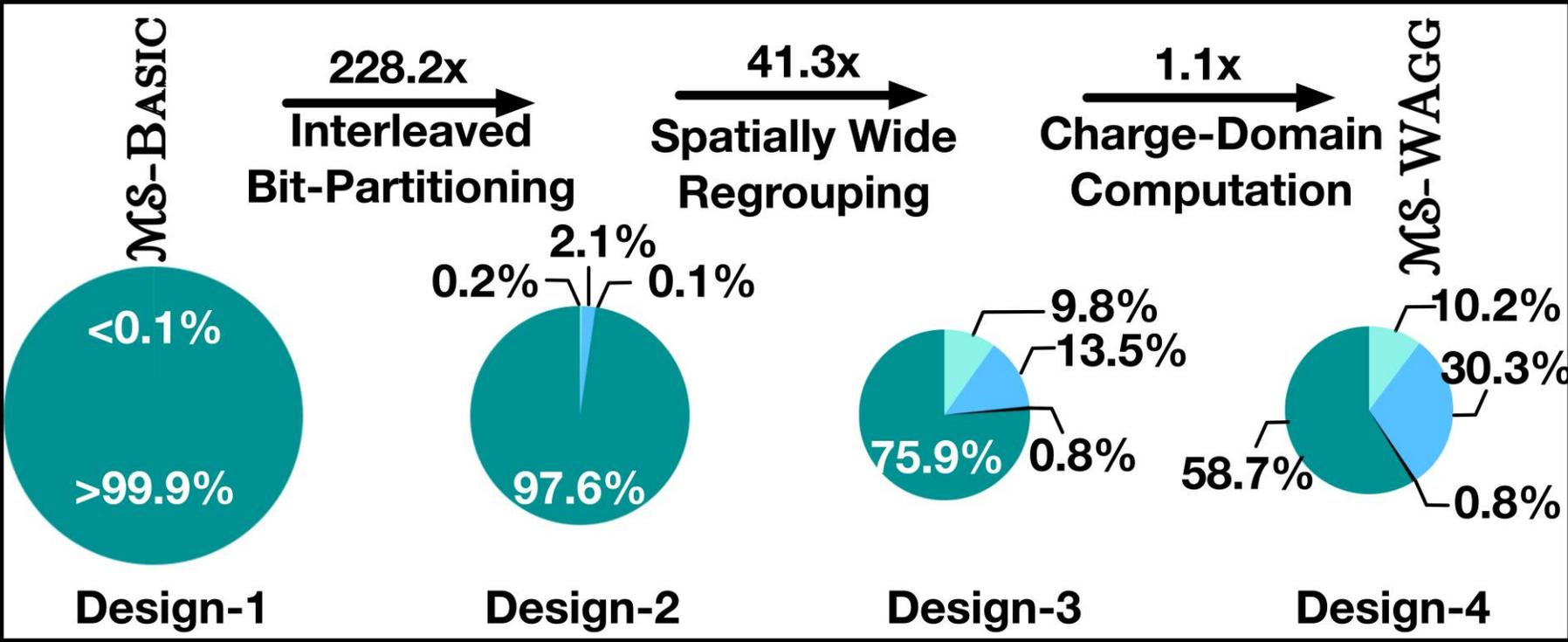
Our Dot-Product Engine:  
MS-WAGG

# BiHiWE Microarchitecture: Design Decisions & Tradeoffs



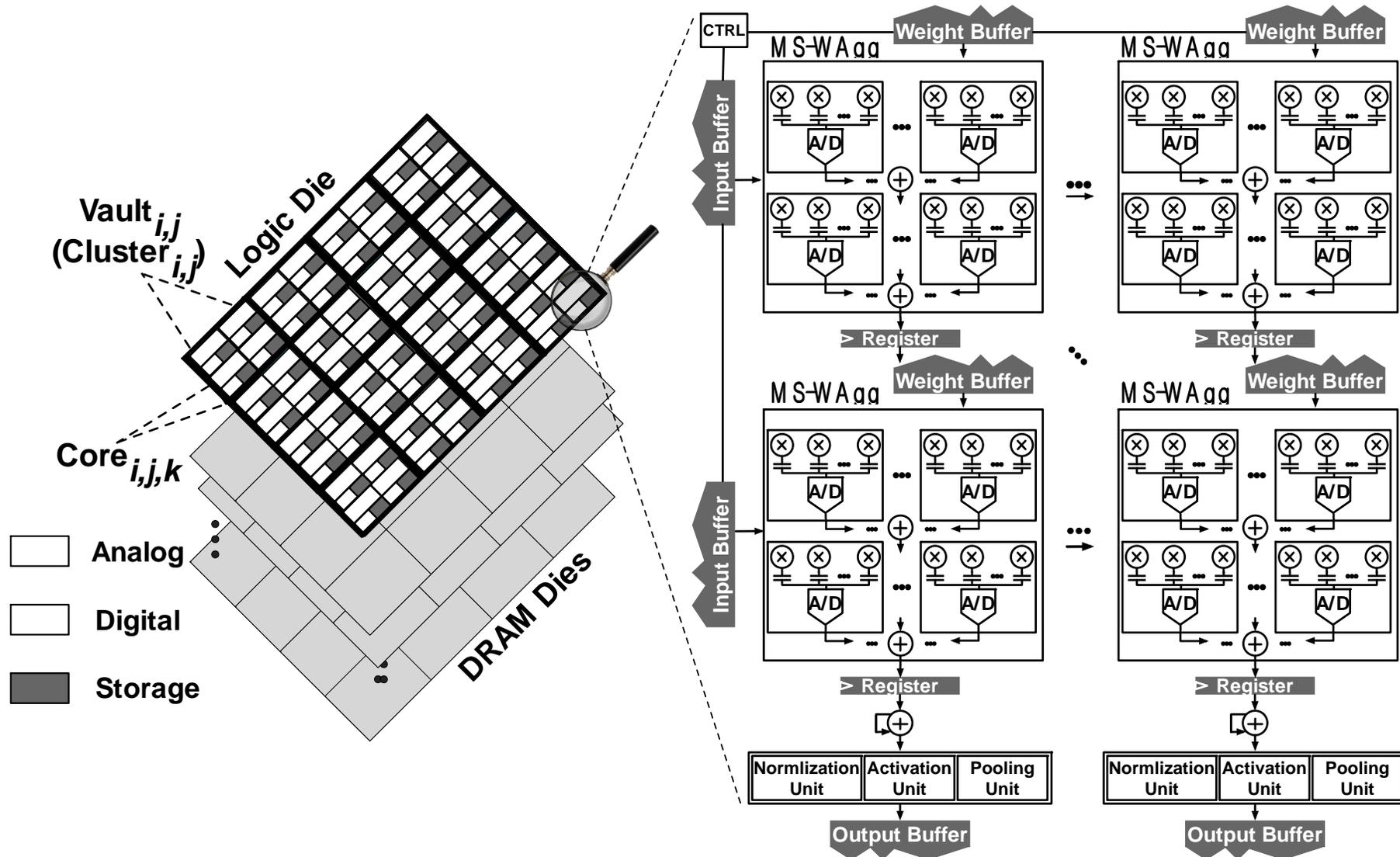
Improvement in Power; Step-by-Step Analysis

# BiHiWE Microarchitecture: Design Decisions & Tradeoffs

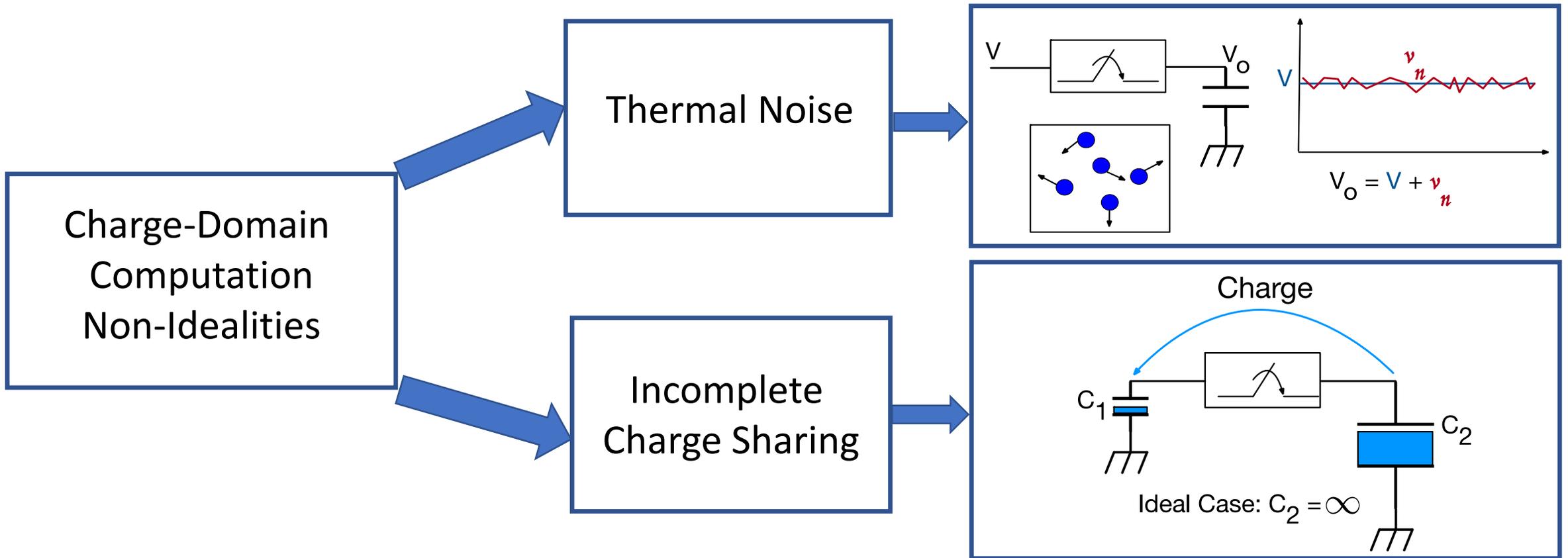


Improvement in Area; Step-by-Step Analysis

# BiHiWE Hierarchical Clustered Architecture

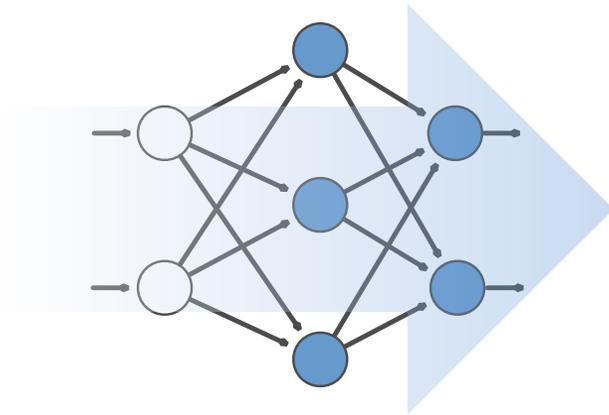


# Mixed-Signal Non-Idealities and Their Mitigation

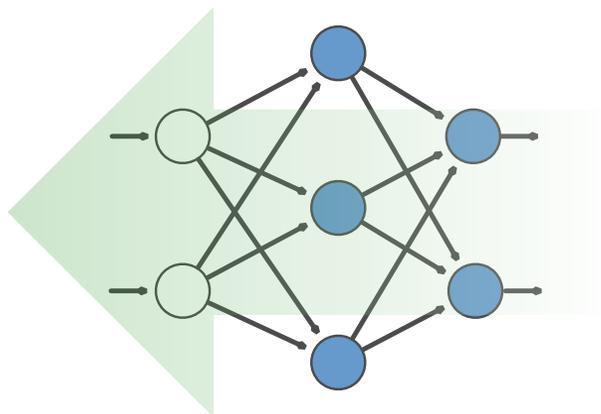


🔑 Injecting the non-idealities to the model and fine-tuning the parameters of the model by retraining the network

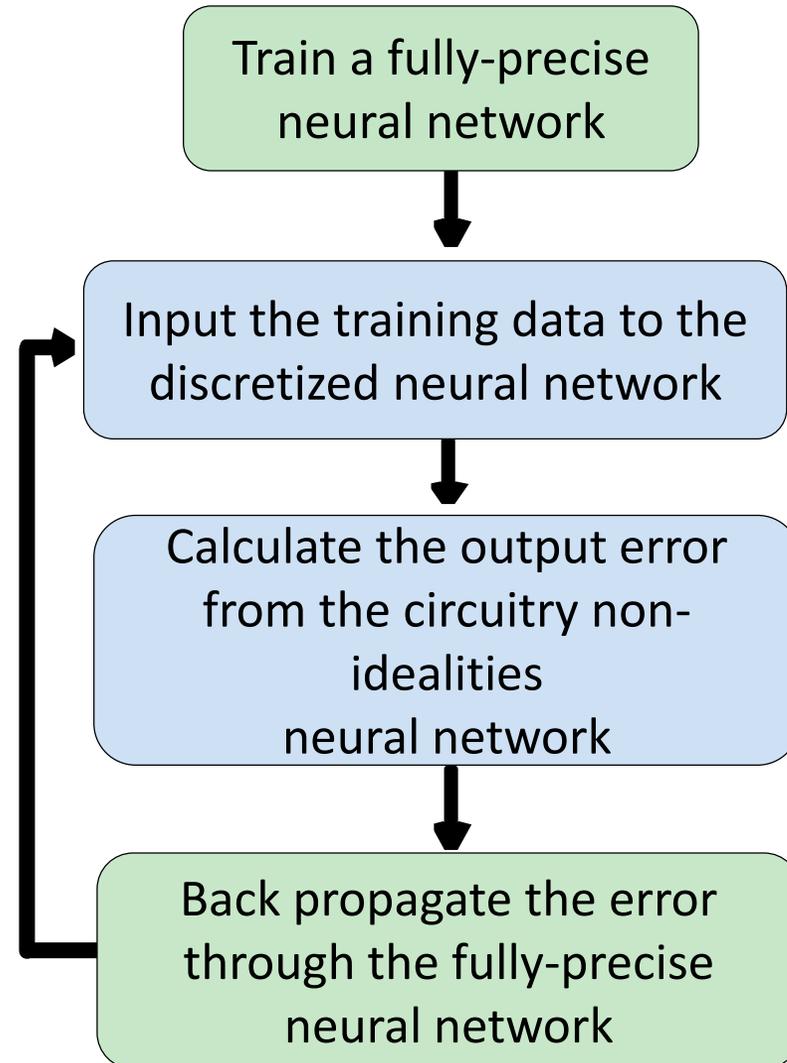
# Mixed-Signal Non-Idealities and Their Mitigation



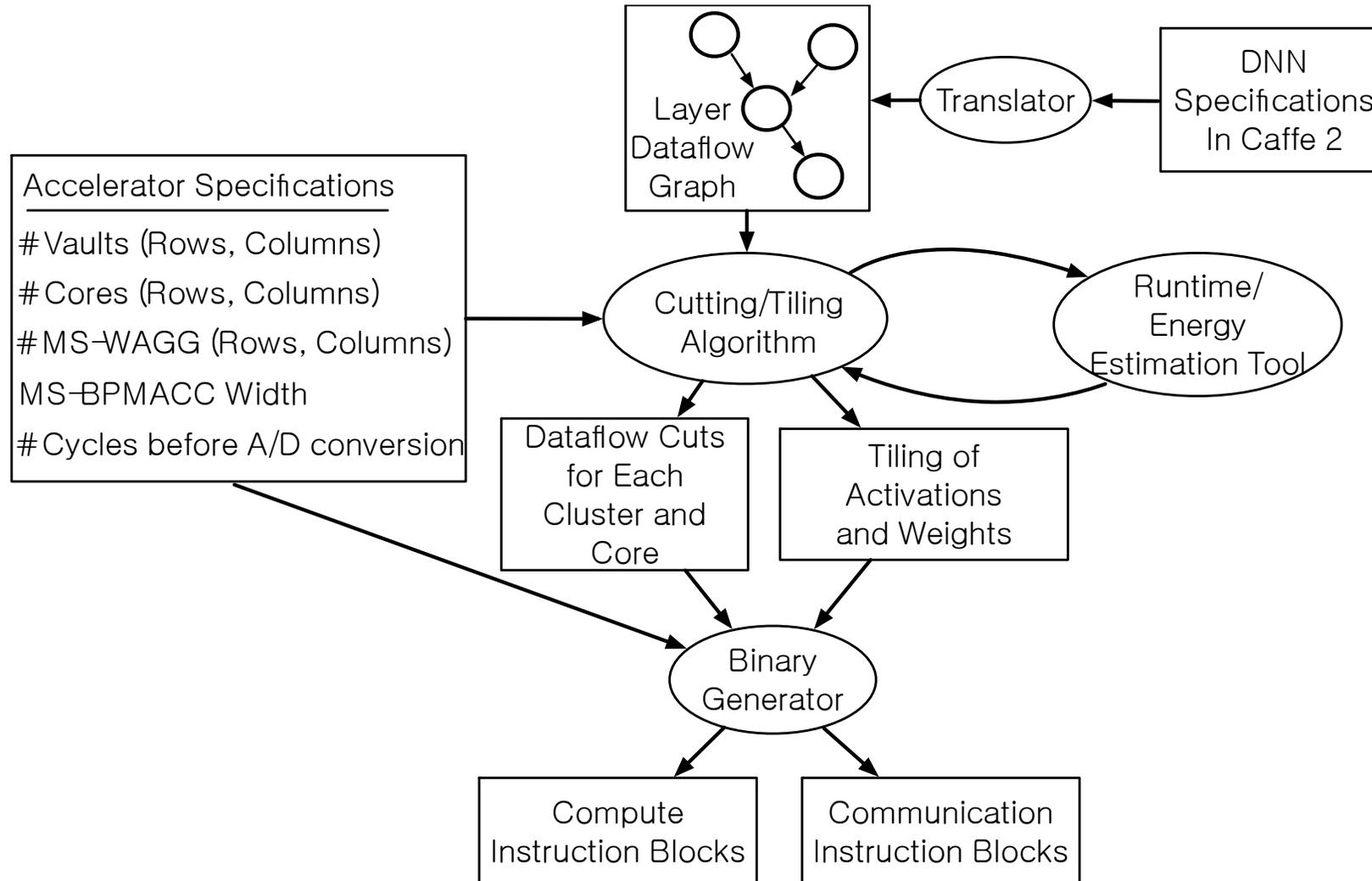
Noisy Network



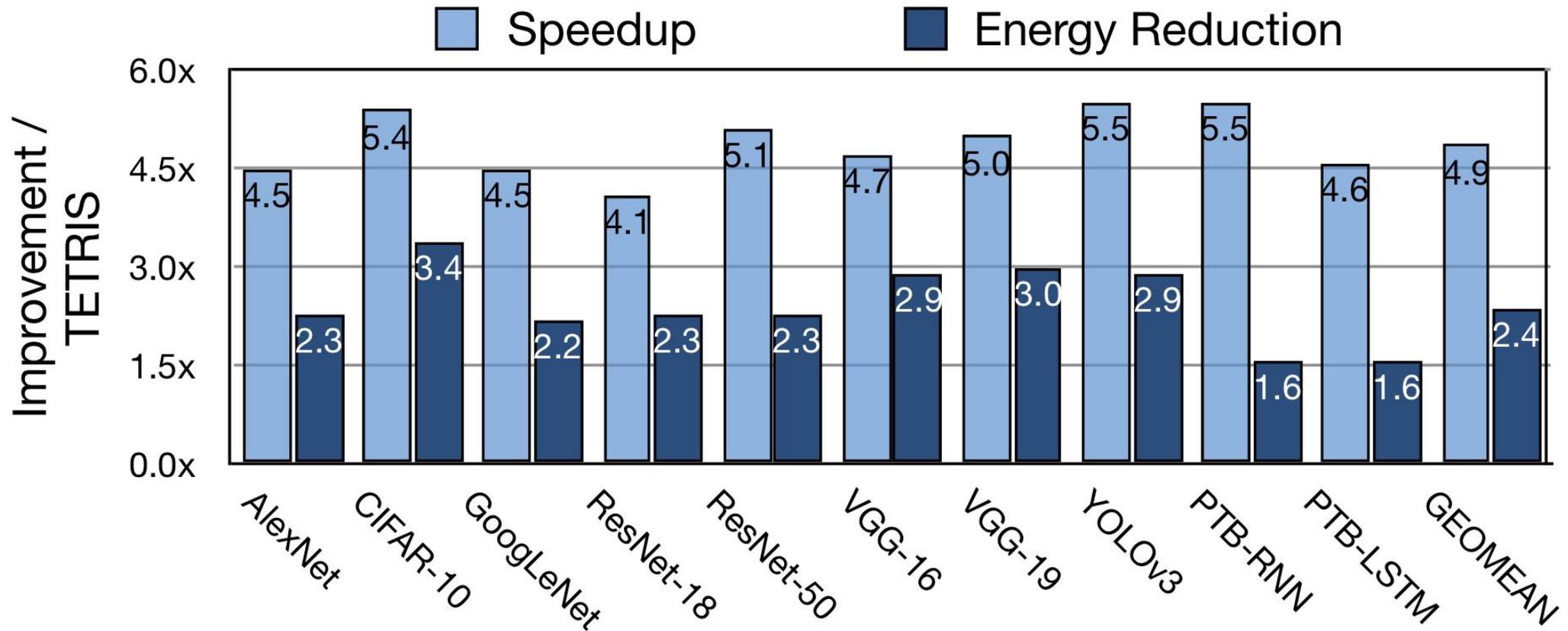
Ideal Network



# BiHiWE Compilation Stack

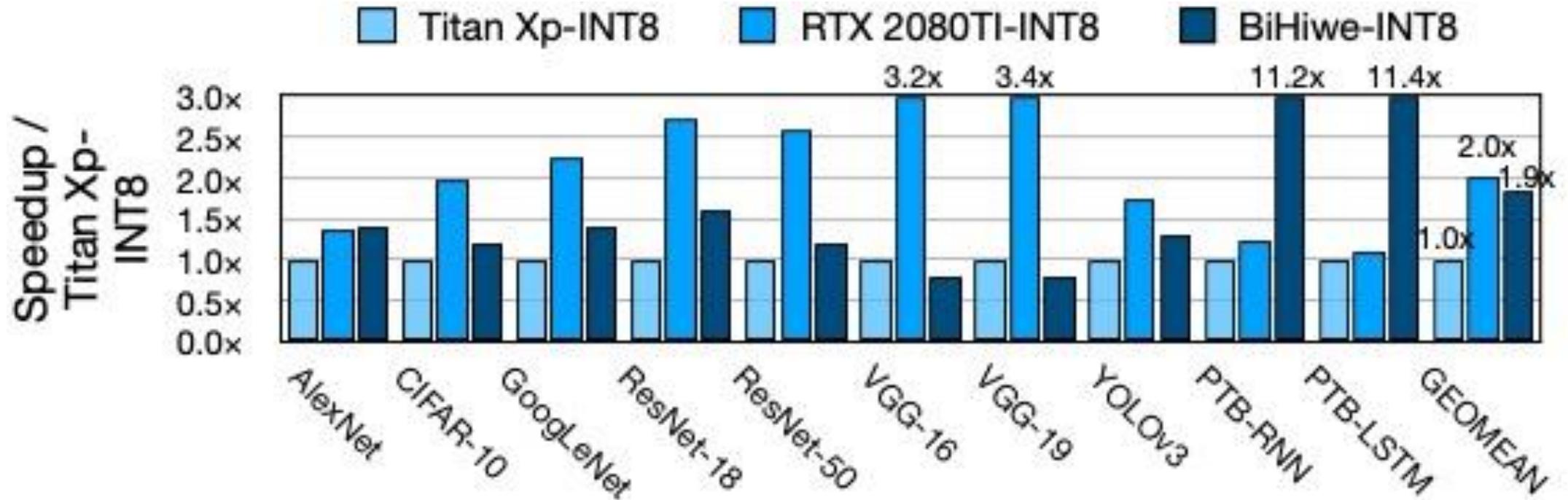


# Comparison with TETRIS



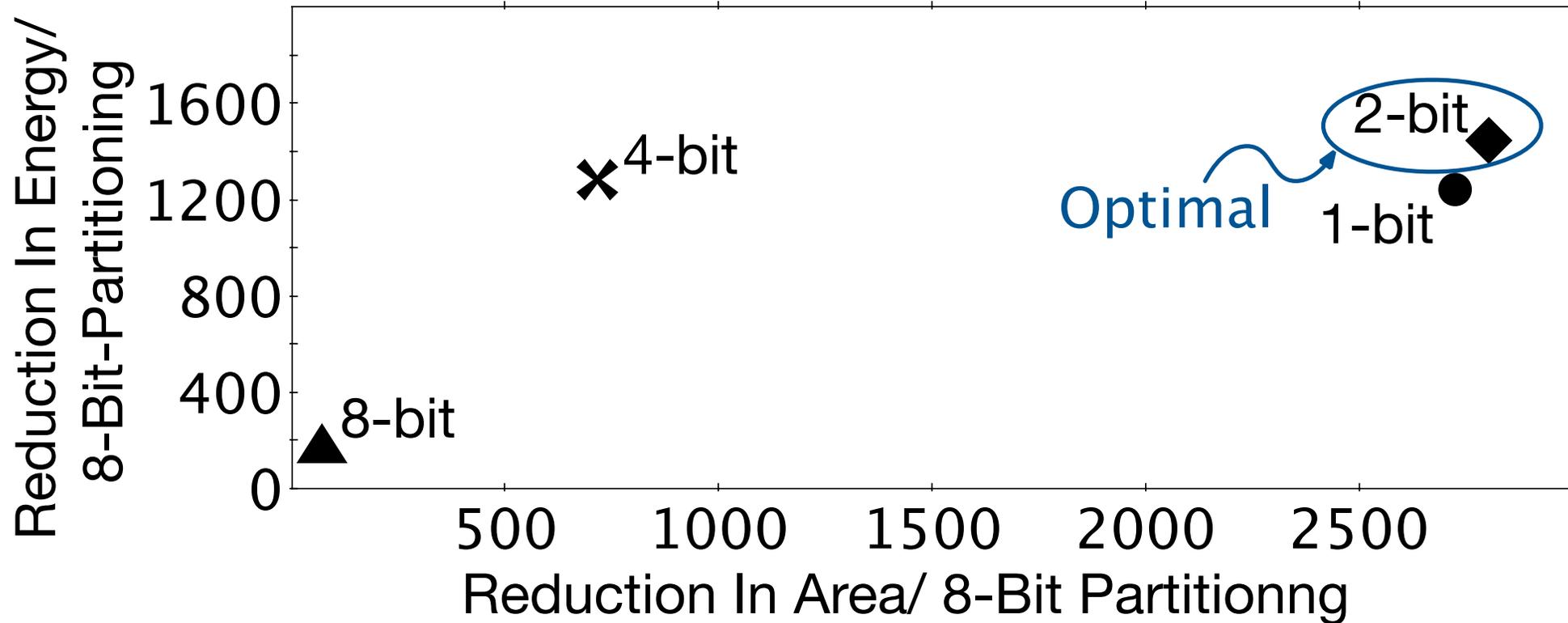
4.9x speedup and 2.4x energy reduction over TETRIS,  
an optimized 3D-stacked fully-digital accelerator for DNNs

# Comparison with GPUs



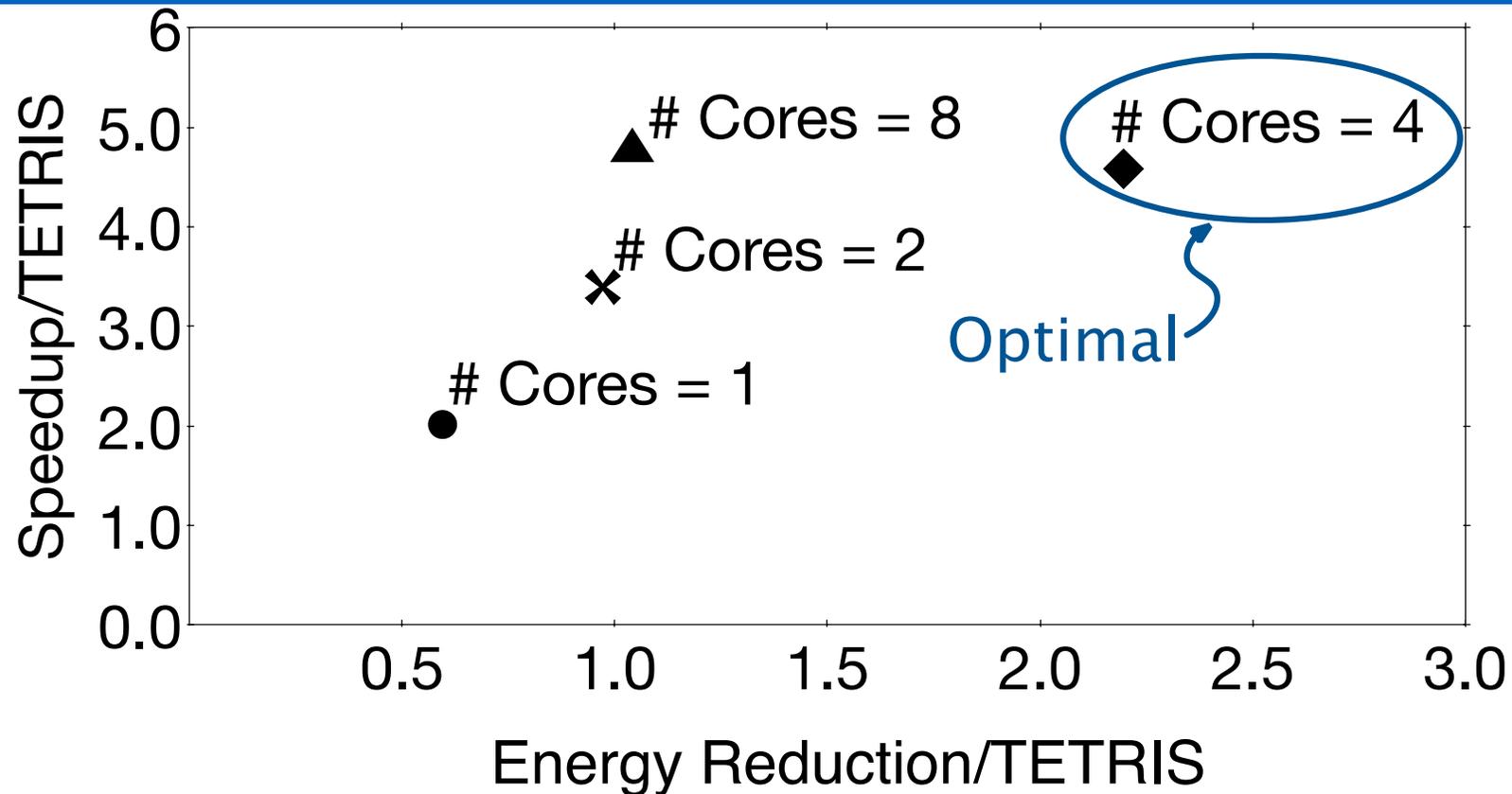
BiHiwe delivers 70.1x and 35.4x higher Performance-per-Watt compared to Nvidia Titan Xp and RTX 2080 TI

# Design Space Exploration for Bit-Partitioning



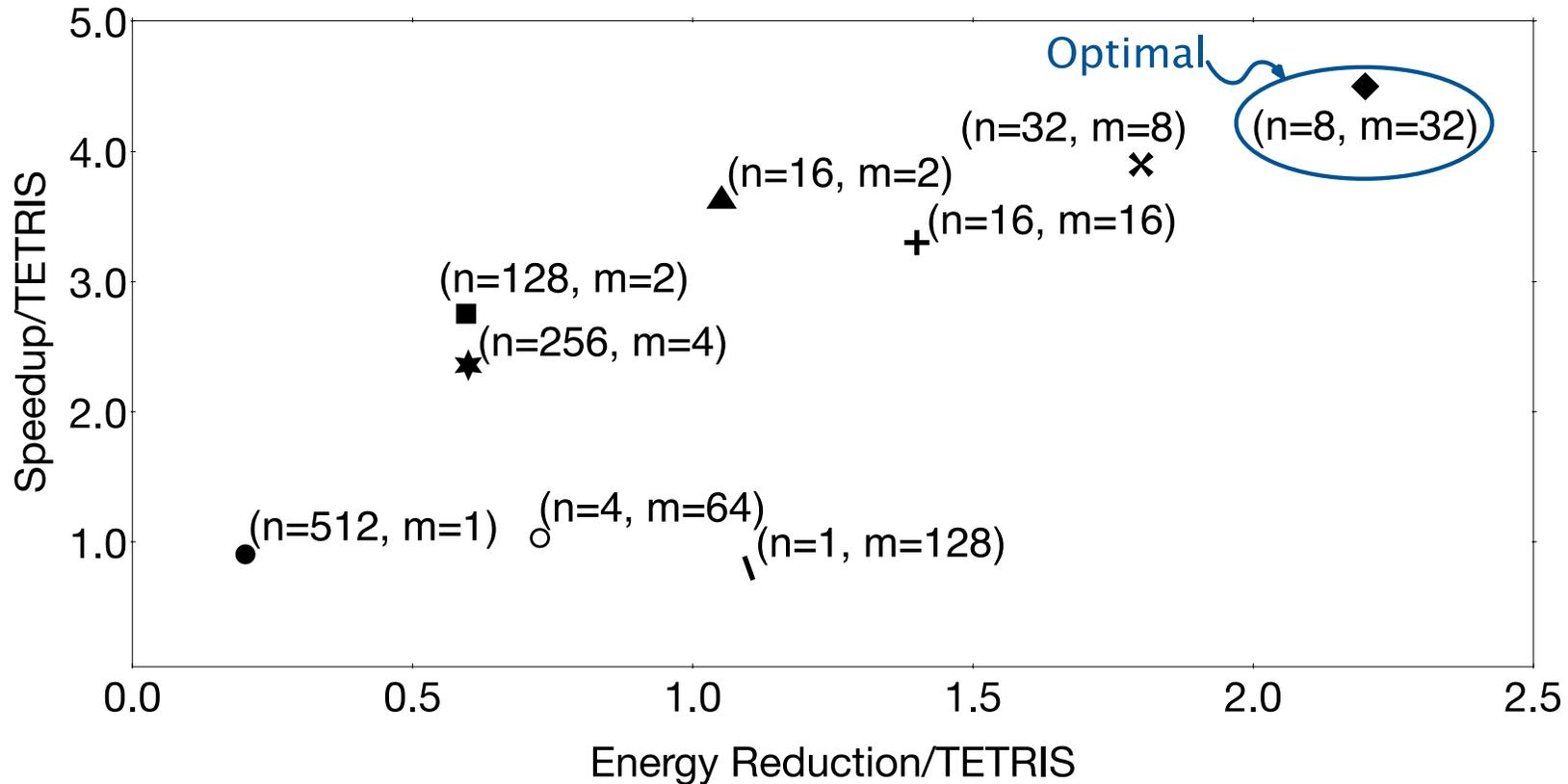
2-bit bit-partitioning is the optimal choice based on this design style and technology node

# Design Space Exploration for # of cores



Each cluster (vault) of the BIHIWE consists of four accelerator cores

# Design Space Exploration for MS-BPMAcc



Each MS-BPMAcc in BiHiWE has an array of 8 low-bitwidth MACC units which perform operations for 32 cycles before A/D conversion

# Evaluating Circuitry Non-Idealities

DNN Model	Dataset	Top-1 Accuracy (With Non-Idealities)	Top-1 Accuracy (After Fine-Tuning)	Top-1 Accuracy (Ideal)	Final Accuracy Loss
AlexNet	Imagenet	53.12%	56.64%	57.11%	0.47 %
CIFAR-10	CIFAR-10	90.82%	91.01%	91.03%	0.02 %
GoogLeNet	Imagenet	67.15%	68.39%	68.72%	0.33 %
ResNet-18	Imagenet	66.91%	68.96%	68.98%	0.02 %
ResNet-50	Imagenet	74.5%	75.21%	75.25%	0.04 %
VGG-16	Imagenet	70.31%	71.28%	71.46%	0.18 %
VGG-19	Imagenet	73.24%	74.20%	74.52%	0.32 %
YOLOv3	Imagenet	75.92%	77.1%	77.22%	0.21 %
PTB-RNN	Penn TreeBank	1.1 BPC	1.6 BPC	1.1 BPC	0.0 BPC
PTB-LSTM	Penn TreeBank	97 PPW	170 PPW	97 PPW	0.0 PPW

BiHiWE has **no virtual impact** on the classification accuracy of the DNN models