



A Sparsity-Aware Analog-Digital Hybrid eDRAM CIM by Effective Row Activation

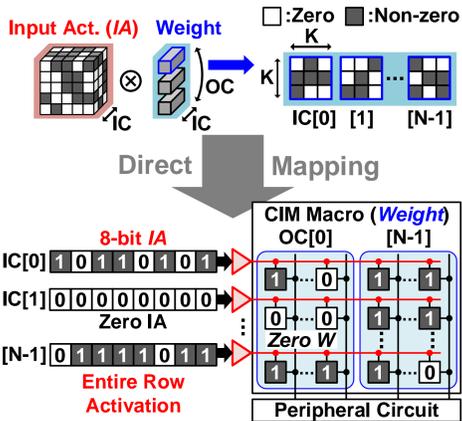
Seungbin Kim¹, Hoichang Jeong¹, and Kyuho Lee²

¹Department of Electrical Engineering, UNIST, Republic of Korea

²Department of Electrical and Electronic Engineering, Yonsei University, Republic of Korea

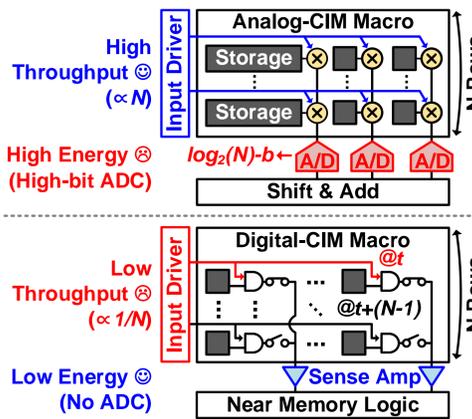
Design Challenges of Previous CIMs

< Waste of Energy Consumption by Zero Computation >



- Input Activations (IA) and Weights in DNN exhibit **Tremendous Sparsity** with random pattern.
- Direct Mapping** of convolution into the CIM Macro and **Entire Row Activation** for MAC operations disables the handling of random sparsity.
- A significant portion of energy is wasted in CIM macro for calculating **83% of zero computation**.

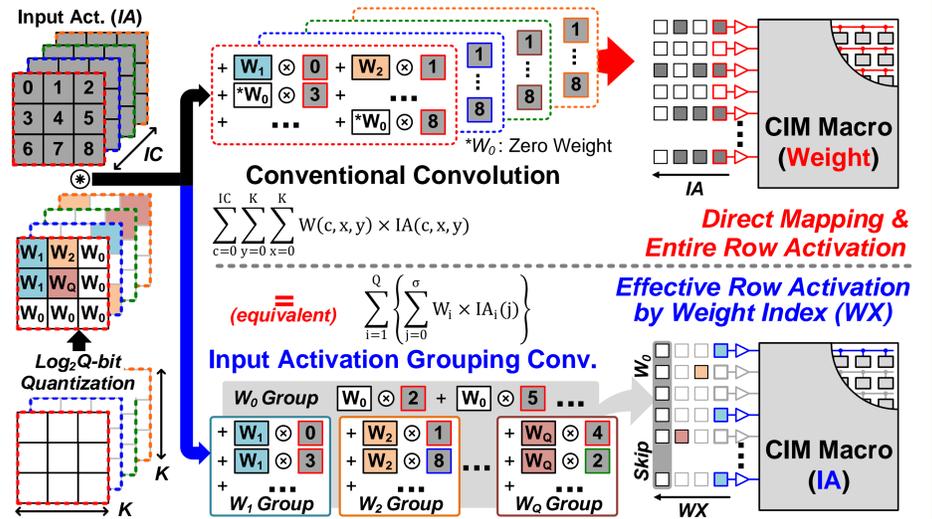
< Distinct Trade-off of Analog CIM and Digital CIM >



- A-CIMs** ensure high throughput by activating entire rows simultaneously. ☹️
- A-CIMs** require high-resolution ADCs to maintain computational accuracy, leading to enormous power usage. ☹️
- D-CIMs** consume less power without utilizing ADC. 😊
- D-CIMs** suffer from low throughput due to sequential row-by-row operations. ☹️

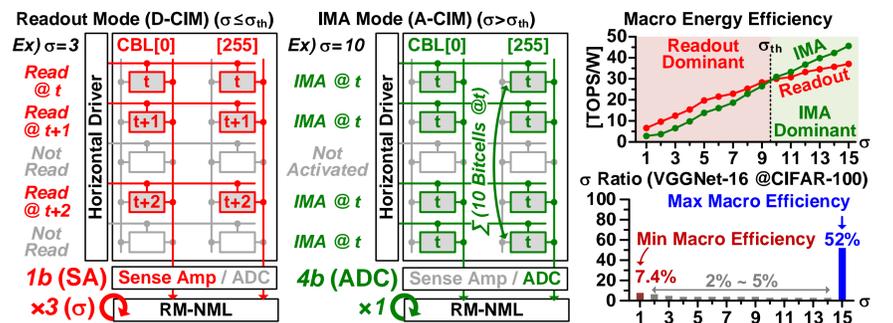
Effective Row Activation w/ Hybrid mode

< Input Activation Grouping Convolution >



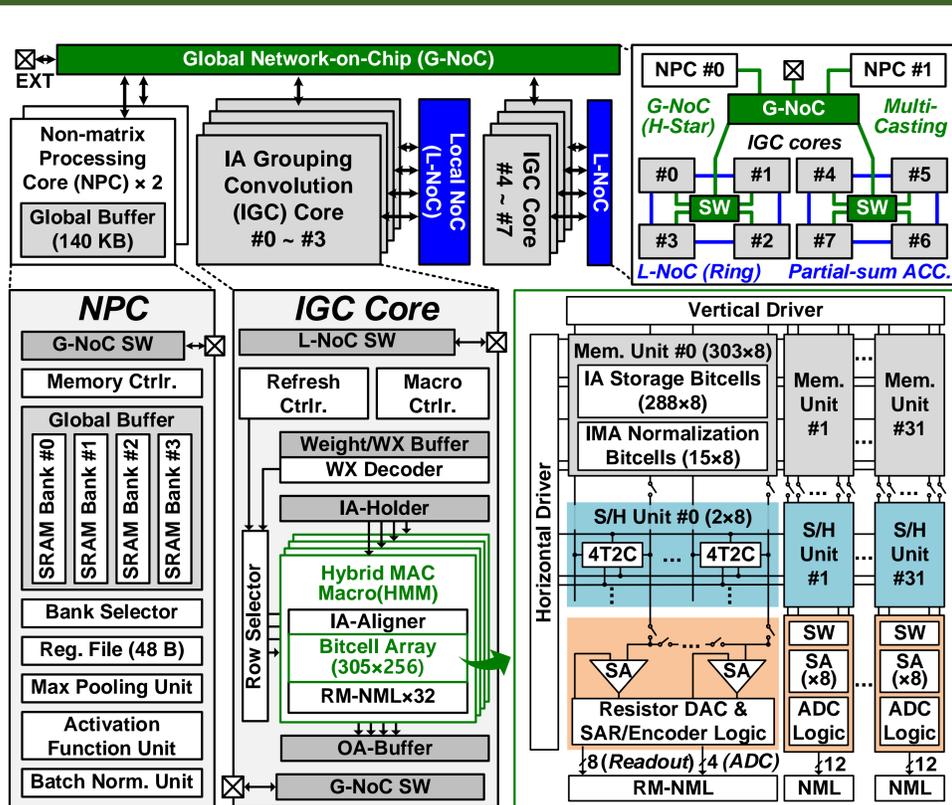
- It **Eliminates Zero Computations** by never activating rows that correspond to the W_0 group, **Retaining Algebraic Equivalence** to conventional convolution.
- IGC improves **Non-zero Computation Ratio** in CIM macro by **1.93x and 4.59x on VGGNet-16 and ResNet-18**.

< Analog-Digital Hybrid MAC Macro >



- Depending on the number of effective row (σ) for W_i group**, HMM decided to operate mode of readout or IMA mode to maximize the energy efficiency.
- The proposed processor achieves benchmark energy efficiency of **24.8 TOPS/W for ResNet-18 and 30.7 TOPS/W for VGGNet-16**.

Overall Architecture



- Hierarchical network-on-chips** consists of global (**G-NoC**) and local (**L-NoC**).
G-NoC: Multicasting data from global buffers to 8 IGC cores w/ H-Star architecture
L-NoC: Facilitating partialsum accumulation from 4 IGC cores w/ Ring architecture
- Hybrid MAC macro** is consists of 305x256 eDRAM bitcell array.
288 rows for storing IAs and perform either in-memory accumulation (IMA) in analog domain or readout in digital domain
15 rows for maintaining computation linearity during IMA
2 rows for sample and hold (S/H) the IMA result and connected to ADC

Chip Photograph & Performance Summary

