

## 15.2 A 2.75-to-75.9TOPS/W Computing-in-Memory NN Processor Supporting Set-Associate Block-Wise Zero Skipping and Ping-Pong CIM with Simultaneous Computation and Weight Updating

Jinshan Yue<sup>1,2</sup>, Xiaoyu Feng<sup>1</sup>, Yifan He<sup>1</sup>, Yuxuan Huang<sup>1</sup>, Yipeng Wang<sup>2</sup>, Zhe Yuan<sup>1</sup>, Mingtao Zhan<sup>1</sup>, Jiaxin Liu<sup>1</sup>, Jian-Wei Su<sup>3</sup>, Yen-Lin Chung<sup>3</sup>, Ping-Chun Wu<sup>3</sup>, Li-Yang Hung<sup>3</sup>, Meng-Fan Chang<sup>3</sup>, Nan Sun<sup>1</sup>, Xueqing Li<sup>1</sup>, Huazhong Yang<sup>1</sup>, Yongpan Liu<sup>1</sup>

<sup>1</sup>Tsinghua University, Beijing, China

<sup>2</sup>Pi2star Technology, Beijing, China

<sup>3</sup>National Tsing Hua University, Hsinchu, Taiwan

Computing-in-memory (CIM) is an attractive approach for energy-efficient neural network (NN) processors, especially for low-power edge devices. Previous CIM chips [1-5] have demonstrated macro and system-level design enabling multi-bit operations and sparsity support. However, several challenges exist, as shown in Fig. 15.2.1. First, though a previously proposed block-wise sparsity strategy [5] can power off ADCs, zeros still contributed to storage requirements, and power gating was not applied to computing resources. Second, on-chip SRAM CIM macros are not large enough to hold all weights. Updating weights between computing operations leads to significant performance loss. Finally, the limited sensing margin incurs poor accuracy for large NN models on practical datasets, such as ImageNet. The precision and power of the ADCs should be optimized and adjusted.

To overcome these challenges, this work proposes a CIM NN processor with 2.75-to-75.9TOPS/W system energy efficiency. The main contributions include: 1) A set-associate block-wise zero-skipping architecture (SABZA) that skips zero-valued activations and weights to reduce power, storage and execution time; 2) A ping-pong CIM (PP-CIM) architecture enabling simultaneous computing and weight update operations; 3) An energy-efficient PP-CIM circuit design and a digital-predictor-assisted adaptive 0/2/4b ADC to support accurate inference on practical datasets.

Figure 15.2.2 shows the overall CIM processor, consisting of two global activation/weight SRAMs, four homogeneous CIM cores, an instruction SRAM with top controller, a set-associate adder array and a result SRAM. Each CIM core contains an activation buffer, a sparsity scheduling module (SSM), a weight index SRAM, a PP-CIM macro and a local adder array. Weights are transformed offline into block-wise sparse pattern and stored in a set-associate compact format. The top controller sends weights from the global weight SRAM to CIM macros and executes multiply-accumulation (MAC) operations in weight-stationary mode. The exchangeable CIM banks alternately execute computation and weight update operations, enabling full-time utilization. The SSM utilizes weight indexes and online activation sparsity to select activated CIM columns and ADC precision. The 64 set-associate adders receive CIM results with indexes and accumulate them to the correct results. A 32b accumulation adder dynamically skips zeros and accumulates the partial-sums in the result SRAM.

Figure 15.2.3 presents the SABZA and the digital-predictor-assisted ADC sensing mechanism. In CIM macros, each two adjacent CIM banks compose a set-associate group to calculate results for 8 output channels. Each group of  $N$   $B$ -bit weight data ( $B=1-8b$ ) is accompanied by a 7b index, illustrating sparsity, row and column number. The 112b indexes for 16 CIM banks are sent to SSM, which dynamically detects activation sparsity. The weight/activation sparsity is classified into 3 patterns: all-zero,  $<4$  non-zeros and  $\geq 4$  non-zeros. The all-zero activation group is skipped. When an all-zero weight group appears due to the misalignment or a very sparse weight distribution, the corresponding ADCs are powered off. The  $<4$  and  $\geq 4$  non-zeros patterns are calculated under 2/4b ADC sensing mode.

Every two ADCs in the same CIM bank are for 1b activation (H1b/L1b), and the results are accumulated in the local adder array. The set-associate group contains two CIM banks and eight set-associate adders. Results from the local adders are selected to be added to specific registers by recognizing the corresponding output channel indexes. Fig. 15.2.3 presents a 2-row 8-output set-associate case to illustrate the compact weight format and workflow. As the activation ID and the second weight index in the first column is inconsistent due to the compact weight storage format, one more clock cycle (Cycle 3) is needed. This situation can be optimized during the NN network training stage. In this example with 75% sparsity, a 4 $\times$  macro storage reduction and 3.2 $\times$  cycle reduction are achieved by zero skipping.

Figure 15.2.4 shows the ping-pong CIM (PP-CIM) architecture. The  $\sim 100Kb$  CIM macro cannot store all weight data for large NN models, often requiring  $>1MB$ . The PP-CIM architecture supports simultaneous weight updates and computing operations. In Phase 1, weights are written into the alternate banks of the CIM macros. In Phase 2, computing and weight update banks exchange. The CIM banks with updated weights are used for computing, while the outdated computing banks start weight updates for the next workload. In case the weight transfer is complete while the current CIM operations are not finished, the weight SRAM goes into the retention mode to save power. Fig. 15.2.4 shows the PP-CIM improvement on different NN layers and models, illustrating different performances with varying ratios of weight-update and computing operations. For a ResNet18 model on CIFAR-10, 1.26 $\times$  overall performance improvement is observed by reducing idle weight-update cycles. The improvement can be 1.94 $\times$  in several layers. Over the models on the CIFAR-10 and ImageNet dataset, the PP-CIM architecture shows a 1.02-to-1.30 $\times$  performance gain and 46-to-98% weight SRAM retention time. Compared with a simple double-buffering strategy requiring an additional 100% macro area, the PP-CIM only increases the macro area by approximately 15%.

Figure 15.2.5 presents the details of the PP-CIM macro comprising the SRAM cell array, input driver, normal IO and 32 ADCs. The 128 $\times$ 128 cell array is divided into 16 $\times$ 128 PP-CIM banks. Each bank contains two bitwise multiply units (A-BMU, B-BMU) and one local computing cell (LCC). Each BMU contains four compact 6T SRAM cells. With the additional B-GBL/B-GBLB pair, this macro supports simultaneous computing and weight-update operations. When one BMU is selected to run the computing operation, the other one can perform weight updates. In the SRAM read/write mode, the SRAM cell is accessed through local bitline pairs (LBL/LBLB), and connected to the vertical global bitline (B-GBL) with  $HWL=V_{DD}$ . In the CIM mode for A-BMU, A-HWL turns off M1 and M2. M3-M6 run the multiply operation (IN $\bullet$ W), wherein weight data stored in an SRAM cell is accessed via LBL/LBLB, and two separated 1b activation inputs are rail-to-rail voltage and accessed via A-GBL/A-GBLB. The results go through the horizontal global bitline (HGBL/HGBLB) to two ADCs to get the final results. In each cycle, the CIM macro activates a maximum of 32 columns for a CIM operation, with no overflow ensured by the digital predictor of SSM. The HGBL/HGBLB pair achieves 2.07 $\times$  better sensing margin compared with a previous 5b ADC [4]. Each ADC can switch between 2b/4b sensing and power off (0b) modes to reduce power consumption.

This processor is fabricated in 65nm CMOS technology and measurement results are shown in Fig. 15.2.6. This chip has separate power domains for the CIM macros and other digital circuits. This chip can work at 25-100MHz with 0.62-1.0V digital voltage and 1.0V CIM voltage. The tested ADC current shows a 1.33 $\times$  average power reduction using the digital predictor. The deployed NN models include VGG and ResNet on the CIFAR-10 and ImageNet dataset. This chip achieves a 3.91-to-4.57 $\times$  on-chip storage reduction, with 2.75-to-18.0TOPS/W system energy efficiency for real NN models considering all on-chip power and execution time. The peak system-level energy efficiency is 75.9TOPS/W at 37.5MHz with 0.65V digital and 1.0V CIM voltage. This work demonstrates a CIM chip with both zero-activation and zero-weight skipping and a PP-CIM architecture with simultaneous computing and weight-update capability and shows reasonable accuracy on larger datasets CIFAR-10 (2.24% loss) and ImageNet (3.45% loss) using a digital-predictor-assisted adaptive-precision ADC. An energy-efficiency breakdown shows 6.35 $\times$  improvement compared with state-of-the-art [5]. Fig. 15.2.7 shows the summary table and die photo.

### Acknowledgement:

This work was supported in part by National Key R&D Program 2018YFA0701500; NSFC Grant 61674094, 61934005, 61720106013; Beijing National Research Center for Information Science and Technology; and Beijing Innovation Center for Future Chips. Corresponding Author: Yongpan Liu.

### References:

- [1] X. Si et al., "A Twin-8T SRAM Computation-In-Memory Macro for Multiple-Bit CNN-Based Machine Learning," *ISSCC*, pp. 396-397, 2019.
- [2] J. Yang et al., "Sandwich-RAM: An Energy-Efficient In-Memory BWN Architecture with Pulse-Width Modulation," *ISSCC*, pp. 394-395, 2019.
- [3] J. Su et al., "A 28nm 64Kb Inference-Training Two-Way Transpose Multibit 6T SRAM Compute-in-Memory Macro for AI Edge Chips," *ISSCC*, pp. 240-241, 2020.
- [4] X. Si et al., "A 28nm 64Kb 6T SRAM Computing-in-Memory Macro with 8b MAC Operation for AI Edge Chips," *ISSCC*, pp. 246-247, 2020.
- [5] J. Yue et al., "A 65nm Computing-in-Memory-Based CNN Processor with 2.9-to-35.8TOPS/W System Energy Efficiency Using Dynamic-Sparsity Performance-Scaling Architecture and Energy-Efficient Inter/Intra-Macro Data Reuse," *ISSCC*, pp. 234-235, 2020.



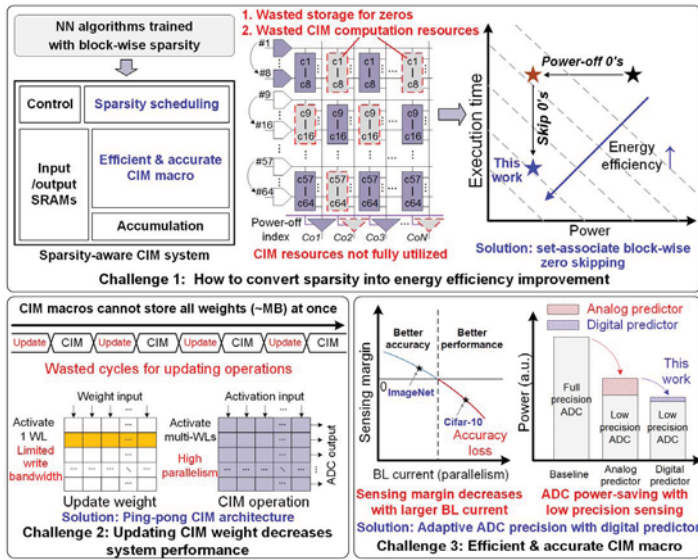


Figure 15.2.1: Design challenges for energy-efficient sparsity-aware CIM neural-network processor.

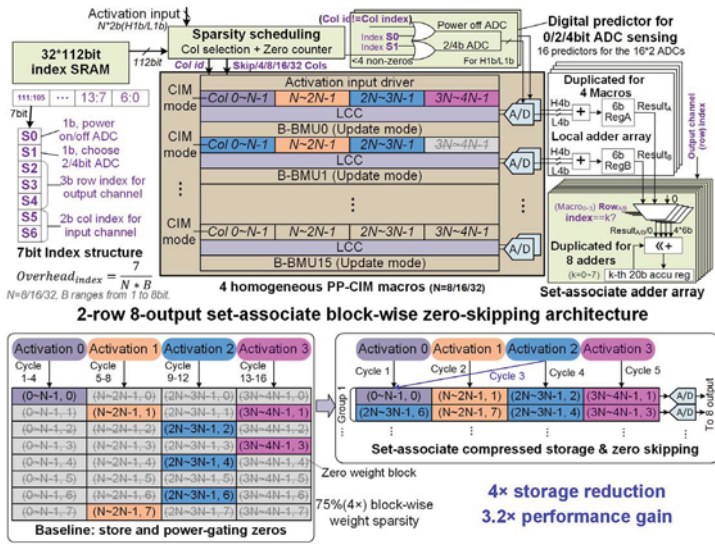


Figure 15.2.3: Workflow of set-associate block-wise zero-skipping and digital-predictor-assisted adaptive 0/2/4bit ADC sensing.

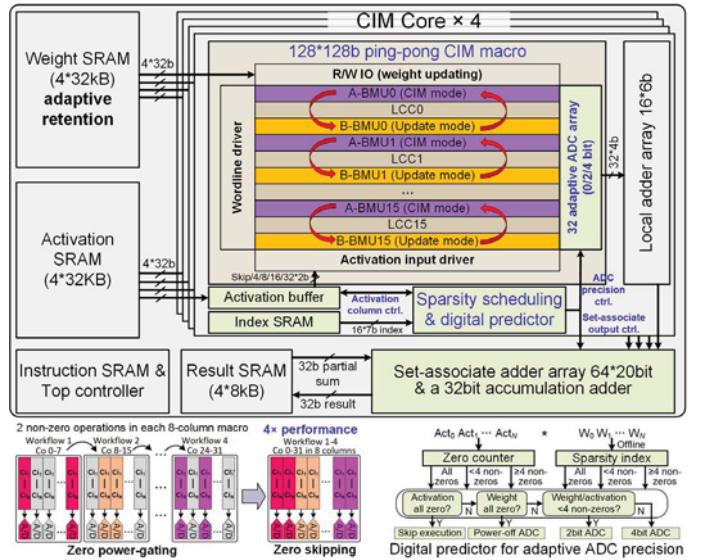


Figure 15.2.2: Overall CIM processor architecture with set-associate block-wise zero-skipping and ping-pong CIM.

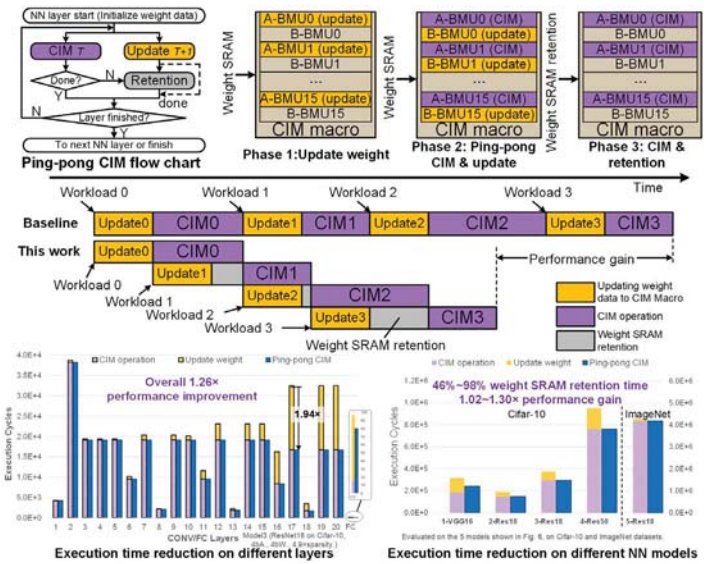


Figure 15.2.4: Ping-pong CIM architecture and performance improvement on different NN layers and models.

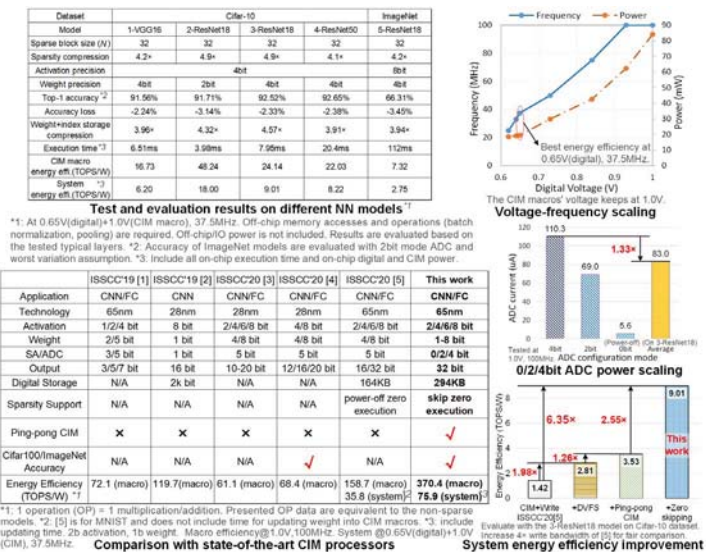


Figure 15.2.5: Circuit diagram of energy efficient PP-CIM macro and digital-predictor-assisted adaptive 0/2/4bit ADC.

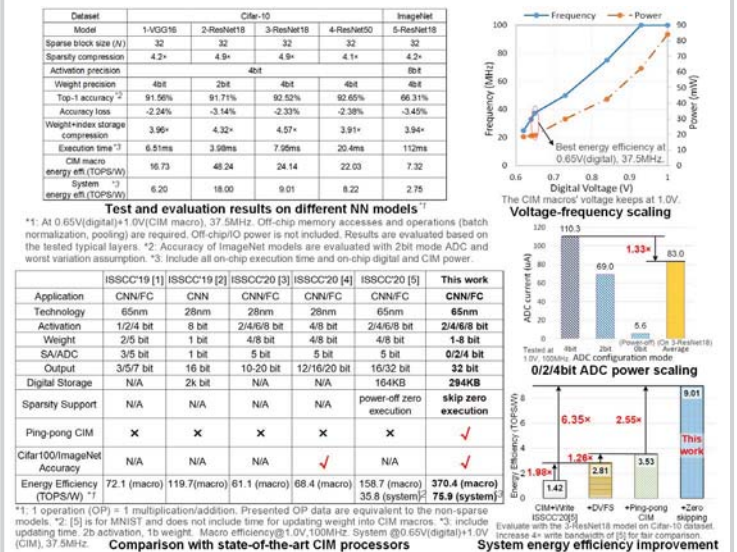


Figure 15.2.6: Measurement results and comparison with state-of-the-art CIM processors.

Technology	TSMC 65nm GP
Chip Area	3.0mm×4.0mm
Core Area	2.44mm×3.41mm
CIM Macro Area	0.69mm×0.62mm×4
Application	CNN / FC
Weight Bit-precision	1-8 bit
Activation Bit-precision	2/4/6/8 bit
Total Digital SRAM	294kB
CIM SRAM	16kbit × 4
Supply Voltage (Digital)	0.62-1.0V
Supply Voltage (Analog)	1.0V
Frequency	25-100MHz
CIM Macro Power <sup>*1</sup>	2.55mW × 4
System Power	18.6-84.1mW
Performance <sup>*2</sup>	0.10-3.16TOPS
<b>CIM Macro<sup>*3</sup></b>	<b>7.32-370TOPS/W</b>
<b>System<sup>*3</sup></b>	<b>2.75-75.9TOPS/W</b>

\*1: Test at 1.0V, 100MHz.

\*2: 25-100MHz, from non-sparse models to 4× sparse models, 2bit activation, 1bit weight.

\*3: Off-chip and IO power is not included. Off-chip operation time (batch normalization, pooling) is not included.

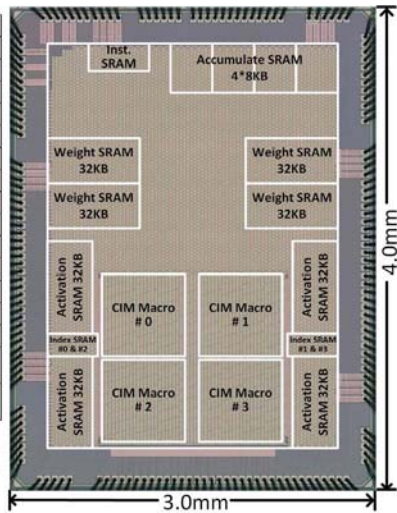


Figure 15.2.7: Die photo and metrics.