# A 28nm Nonvolatile AI Edge Processor using 4Mb Analog-Based Near-Memory-Compute ReRAM with 27.2 TOPS/W for Tiny AI Edge Devices

Tai-Hao Wen[1], Je-Min Hung[1], Hung-Hsi Hsu[1], Yuan Wu[1], Fu-Chun Chang[1], Chung-Yuan Li[1], Chih-Han Chien[1], Chin-I Su[2], Win-San Khwa[2], Jui-Jen Wu[2], Chung-Chuan Lo[1], Ren-Shuo Liu[1], Chih-Cheng Hsieh[1], Kea-Tiong Tang[1], Mon-Shu Ho[3], Yu-Der Chih[2], Tsung-Yung Jonathan Chang[2], Meng-Fan Chang[1,2]

[1]National Tsing Hua University, Taiwan; [2]TSMC, Taiwan, [3]National Chung Hsing University, Taiwan

*Email: mfchang@ee.nthu.edu.tw

## Abstract

Tiny AI edge processors prefer using nvCIM to achieve low standby power, high energy efficiency (EF), and short wakeup-to-response latency ($T_{WR}$). Most nvCIMs use in-memory computing for MAC operations; however, this imposes a tradeoff between EF and accuracy, due to MAC accumulation-number ($N_{ACU}$) versus signal margin and readout quantization. To achieve high EF and high accuracy, we developed a system-level nvCIM-friendly control scheme and a nvCIM macro with two analog near-memory computing schemes. The proposed 28nm nonvolatile AI edge processor with 4Mb ReRAM-nvCIM achieved high EF (27.2 TOPS/W), short $T_{WR}$ (3.19ms), and low accuracy loss (<0.5%) The EF of the ReRAM-nvCIM macro was 38.6 TOPS/W.

**Keywords:** multiply-and-accumulation (MAC), nonvolatile memory (NVM), nonvolatile compute-in-memory (nvCIM)

## Introduction

Most AI processors [2]-[4] without nvCIM suffer low EF and long $T_{WR}$ (Fig. 1), due to the need to move weight data from off-chip NVM to on-chip processing units after wakeup. Therefore, AI edge processors [1] prefer using nvCIM [5]-[6] to reduce $T_{WR}$ and improve EF by integrating NVM storage and computation functionalities within the same macro (nvCIM).

The challenges of designing nvCIM-based processors (Fig.2) include (1) limited $N_{ACU}$ in nvCIM due to small and non-linear signal margin associated with the read scheme and parasitic resistance on bitlines (BL); (2) a macro-level tradeoff between ADC precision versus read-yield/energy; (3) a system-level tradeoff between EF/latency versus inference accuracy.

## Overview of proposed Nonvolatile AI Edge Processor

Our objective was to develop a nonvolatile AI processor with high EF, short $T_{WR}$, and high accuracy. This was achieved by developing two macro-level nvCIM schemes and one system-level scheme, including (1) near-memory horizontal accumulation (NMHA) to achieve high $N_{ACU}$ and robust/linear signal margin; (2) non-linear quantization readout (NLQR) to increase precision and sensing margin (read-yield) without increasing ADC energy; (3) MACV-aware LSB truncation (MALT) nvCIM-control scheme to increase system throughput and EF without compromising inference accuracy. The proposed processor (Fig. 3) includes 4Mb ReRAM-nvCIM, a system controller, a 512Kb SRAM buffer, and other neural network functions.

## nvCIM: Near-Memory Horizontal Accumulation

The proposed ReRAM-nvCIM macro (Fig. 4) includes 16 output channels, each of which comprises one ReRAM sub-array, one NMHA unit with 64 global computing cells (GCC), and one NLQR unit. The NMHA operation comprise three phases. Phase-1: When IN=1, each GCC precharges (PRE_EN=1) BL to $V_{BL}$ and pulls CN to $V_{CN}$. When IN=0, the BL and CN remain at VSS. Phase-2: PRE_EN=0 and a wordline (WL) pulse is activated. If the weight stored in a ReRAM cell is 0 (LRS), then the memory-cell current discharges the BL and CN and then turns off N3. If the weight stored in a ReRAM-cell is 1 (HRS), then the BL and CN maintain their voltage and N3 is on ($I_{N3}>0$), indicating that the multiplication result (INxW) of this GCC is "1". Phase-3: NMHA accumulates all $I_{N3}$ horizontally from $GCC_0$ to $GCC_{63}$ and generates dataline current ($I_{DL}$) as the current-domain partial MACV of 64 accumulations (pMACV=$\sum_{k=0}^{63} IN_k \times W_k$).

## nvCIM: Non-linear Quantization Readout (NLQR)

The proposed NLQR (Fig. 5) uses 5 phases to determine the pMAC precision and sensing margin of the ADC according to the detected pMACV zone. Phase-1: The 2b detector readouts $I_{DL}$ to detect the ZONE of pMACV ($Z_{pMACV}[1:0]$). Phase-2: An adjustable current-mirror (ACM) modifies the mirror-ratio ($R_{CM}$) according to $Z_{pMACV}$. Phase-3: An ACM converts $I_{DL}$ to a voltage ($V_{C0}$) on capacitor C0. Note that the signal margin on $V_{C0}$ is proportional to $R_{CM}$. Phase-4: The ADC reads out 4 bits from $V_{C0}$ and generates scaled-pMACV ($pMACV_S$). Phase-5: The place-value of $pMACV_S$ is recovered according to $Z_{pMACV}$.

## System Scheme: MACV-aware LSB Truncation (MALT)

The proposed MALT (Fig. 6) dynamically truncates the MAC operation for the input-LSBs of nvCIM according to predicted-MACV zone. In each 8b-MAC operation, 64 ($N_{ACU}$) inputs are sent from the input feature map to the nvCIM macro, bitwise from MSB to LSB. The MALT counts the number of "1" ($N_{IN[7]}=\sum_{k=0}^{63} IN_k[7]$) in input-MSB, and then compared it with the threshold values ($N_{TH0-3}$) to determine whether to skip the MAC operation for input-LSBs. If $N_{IN[7]}>N_{TH3}$, the predicted-MACV is in a large-value zone and the MAC operation for IN[3:0] in nvCIM is skipped with small MACV deviation-%. If $N_{IN[7]}<N_{TH0}$, the predicted-MACV is small and the full INPUT is applied to nvCIM without skipping to obtain a matched-MACV without loss of system-level accuracy.

## Measurement Results and Conclusion

A nonvolatile AI edge processor chip is fabricated using foundry provided 28nm ReRAM technology and a demo system is used to verify the chip, as Fig. 7 shows. When applying the ResNet-20 model to the CIFAR-100 dataset for one-shot inference, the measured $T_{WR}$ was 3.19ms. Under 8b precision, the top-1 inference accuracy degradation was 0.45%, while the measured EFs of the processor and macro were 27.2 TOPS/W and 38.6 TOPS/W, respectively. As shown in Table I, the EF of the proposed processor was 15x higher than that of previous nvCIM-based processors.

## References

[1] M. Chang, ISSCC, 2022, pp. 1-3. [2] M. Giordano, VLSI, 2021, pp. 1-2. [3] D. Rossi, ISSCC, 2021, pp. 60-62. [4] V. Jain, VLSI, 2022, pp. 20-21. [5] J. -M. Hung, ISSCC, 2022, pp. 1-3. [6] J. -M. Hung, Nat Electron 4, pp. 921–930, 2021.
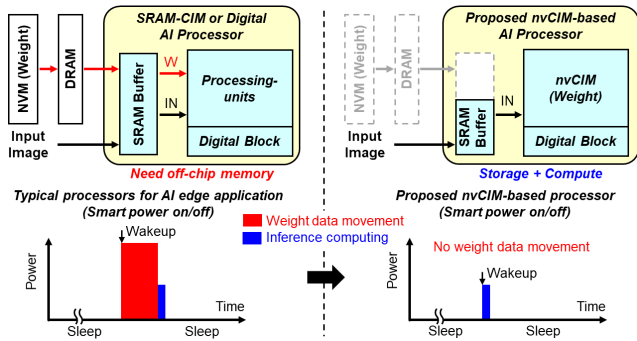
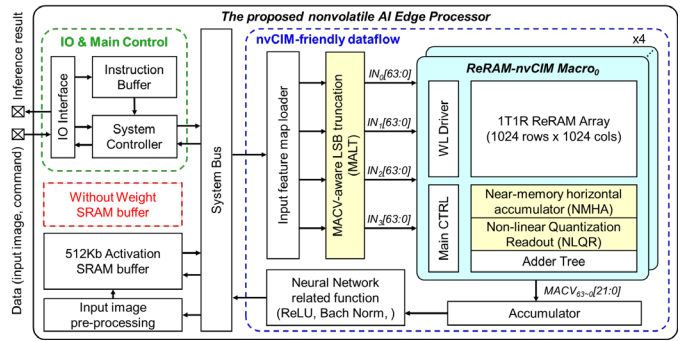Fig.1 Architectures and behavior of typical and nvCIM processor.



Fig. 3 Architecture of proposed nonvolatile AI edge processor.
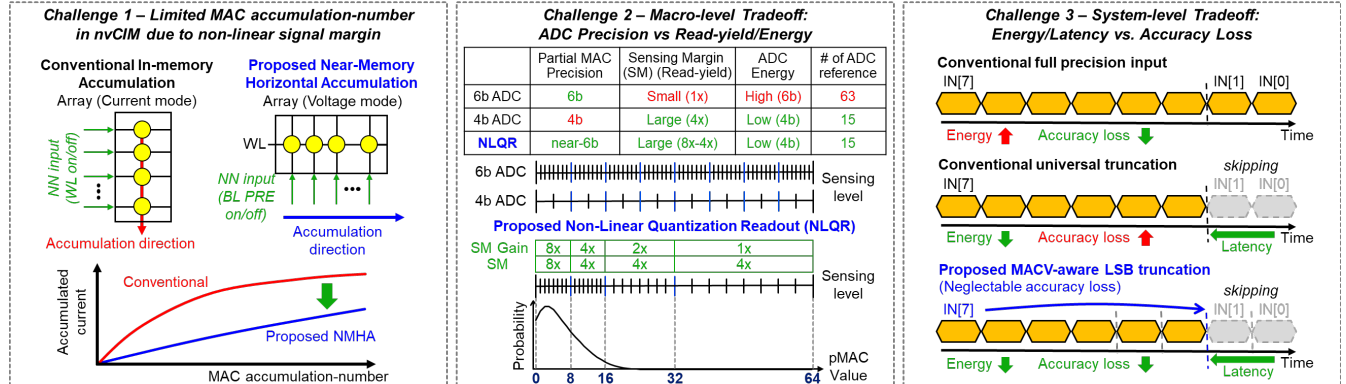


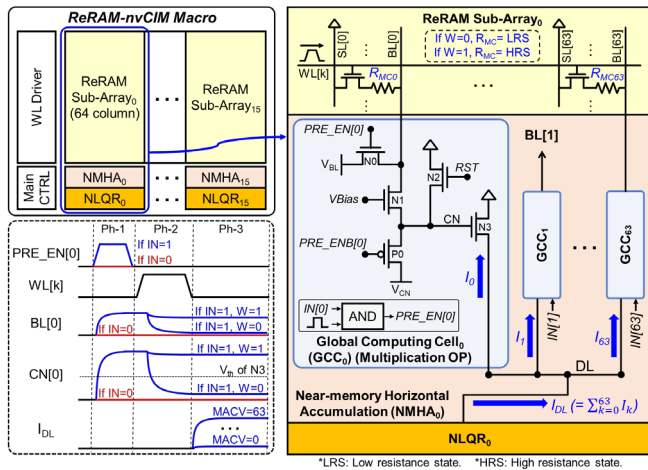Fig. 2 Challenges in developing an nvCIM-based processor and the proposed schemes.



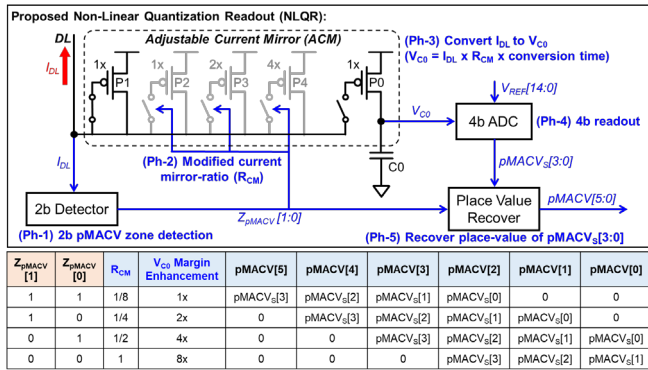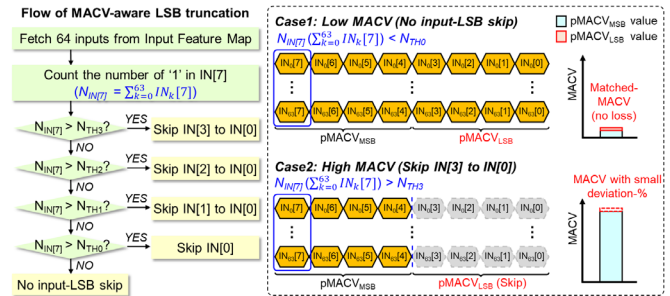Fig. 4 Structure of ReRAM-CIM macro and near-memory horizontal accumalation (NMHA) scheme.



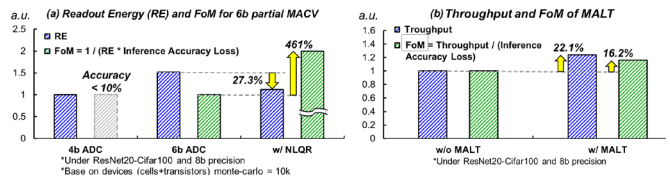Fig. 6 Flow chart and illustration of proposed MACV-aware LSB truncation (MALT) scheme.



Fig. 7 Simulation results of proposed schemes: (a) NLQR decreased energy by 27.3% and increase FoM by 461%; (b) MALT increased throughput by 22.1% and FoM by 16.5%.
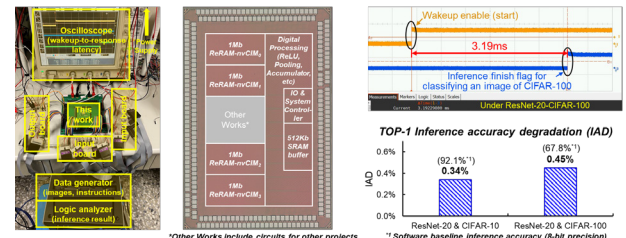


Fig.5 Circuitry and operation of proposed non-linear quantization readout (NLQR) scheme.



Fig.8 Die photo, demo system, and measured results.

| System works | This work | ISSCC2022 [1] | VLSI2021 [2] | ISSCC2021 [3] | VLSI2022 [4] |
|---|---|---|---|---|---|
| Technology | 28nm | 40nm | 40nm | 22nm | 22nm |
| On-chip NVM Device | ReRAM | ReRAM | ReRAM | MRAM | MRAM |
| NVM type | Storage + Compute | Storage + Compute | Storage | Storage | Storage |
| Computing approach | Analog NMC | IMC | Digital | Digital | Digital |
| Chip area | $6^{*3}$ | 25 | 29.2 | 12 | 6.25 |
| Supply voltage (V) | 0.7 - 0.8 | 0.9 | 1.1 | 0.5 - 0.8 | 0.4 - 0.9 |
| Frequency (MHz) | 50 - 200 | 200 | 200 | 0.032 - 450 | 150 |
| System level IN/W precision | INT1 to INT8 | INT1 to INT8 | INT 8, FP16 | INT8, 16, 32, FP, BF | INT2, 4, 8 |
| Throughput (TOPS) *1 | $0.34^{*4}$ | N/A | 0.92 | 0.032 | 0.0176 |
| Energy efficiency (TOPS/W) *1 | $27.2^{*4}$ | 1.8 | 2.2 | 1.3 | 2.47 |
| Computing density (TOPS/mm²)*1 (Include on-chip NVM device) | $0.056^{*4}$ | N/A | $0.031^{*2}$ | $0.0026^{*2}$ | $0.0028^{*2}$ |

*1 Under 8b input-weight precision. *2 Estimated by (Throughput / Active area) *3 One operation (OP) represents one multiplication or one addition. *4 Measured at VDD=0.8V, 200MHz, 8-bit precision. *5 Include 277 IO PADs for general purpose and testmode.

Table I Comparison table of nonvalatile AI processors.