# A 28nm Hybrid 2T1R RRAM Computing-in-Memory Macro for Energy-efficient AI Edge Inference

Wang Ye[1,3], Chunmeng Dou[1,3], Linfang Wang[1,3], Zhidao Zhou[1,3], Junjie An[1], Weizeng Li[1,3], Hanghang Gao[1,3], Xiaoxin Xu[1,3], Jinshan Yue[1], Jianguo Yang[1,3], Jing Liu[1,3], Dashan Shang[1,3], Jinghui Tian[2], Qi Liu[1,2], Ming Liu[1,2]

[1]Institute of Microelectronics of the Chinese Academy of Sciences, Beijing, China, [2]Frontier Institute of Chip and System, Fudan University, Shanghai, China, [3]University of Chinese Academy of Sciences, Beijing, China

Resistive Memory (RRAM) based computing-in-memory (CIM) can play a key role in intermittently operated AI edge devices and sensors [1-3]. By minimizing the data movement during computation and power switching, it effectively reduces the power and hence improves the real-time performance compared to conventional energy-efficient system. However, as shown in Fig. 1, previous RRAM CIM macros, based on 1T1R or 2T2R cell array, are still facing several critical challenges, including (1) large summation current induced considerable IR drop, which results in the nonlinearity between multiply-and-accumulate (MAC) signals and MAC values (MACV) [4], (2) MAC signal margin degradation due to small resistance-ratio (R-ratio), which makes the leakage of high-resistive-state (HRS) cell become non-negligible compared to that of the low-resistive-state (LRS) one [5], and (3) large hardware-cost for the analogue readout circuitry, which typically requires $2^N-1$ stand-by reference signals for N-bit output precision [6].

To deal with these challenges, this work presents the first hybrid-2T1R (H2T1R) RRAM CIM macro to perform highly efficient and reliable analogue MAC. Compared to the conventional 1T1R cell array, it achieves >13X enhanced R-ratio, >80% reduced summation current, >67% smaller word-line (WL) voltage ($V_{WL}$) for low VDD operation and read disturb suppression, and precise multi-bit weight encoding with signal-level-cell. Besides, the proposed reference-subtracting current sense amplifier (RS-CSA) requires only N reference signals for N-bit output and halves the maximum dynamic range of the current mirror. The test-chip is fabricated using the 28nm high-k/metal-gate (HKMG) process with BEOL integrated TaOx RRAM. It can perform highly-linear analogue MAC over 32 input channels with >63% and >79% power saving for the cell array and the analogue readout circuitry, respectively.

The proposed design, as shown in Fig.2, features (1) a macro structure based on the H2T1R cell with decoupled data path and periphery circuits for memory and CIM mode to meet their different requirements, (2) weighted cell array to encode 3-bit weight using three H2T1R SLC cells and (3) a 4-bit RS-CSA with reduced reference generation cost. The H2T1R cell hybridizes one IO transistor (T1) as the selector for the memory data-path to tolerate the write voltage and another core device (T2) for the CIM data-path to improve the computing performance. The weighted cell array comprises four sub-arrays. Three sub-arrays are used to store 3-bit weight data for analogue MAC, in which different numbers of multiplier (m=1, m=2, and m=4) of T2 in the H2T1R cell are leveraged to represent different bit positions. One additional sub-array with m=4 is used as the redundancy for the most-significant-bit (MSB). In the memory mode (CIMB = 1), the RRAM cells can be addressed by the row and column address (AX and AY) to write/read (WEB = 0/1) the input/output data (DIN/DOUT). In the CIM mode (CIMB = 0), the RRAM cells are row-wisely accessed and the CIM input data (CIN) are applied through the bit-lines (BLs) by the analogue data input driver and the sub-array selector. The CIM output data (COUT) can be read out from the transpose BL (TBL) current ($I_{TBL}$) by RS-CSA.

Figure 3 shows the structure and operation of the H2T1R cell array. The IO T1 and core T2 are vertically aligned to minimize the bit-cell area overhead (30.3%). In order to optimize the signal ratio, T1 is half activated to maximize the change of the voltage at the node X ($\Delta V_X$) at different resistive states. Maximized $\Delta V_X$ can be observed at the point that the square of the resistance of T1 ($R_{T1}$) equals to the product of the cell resistance at HRS ($R_{HRS}$) and LRS ($R_{LRS}$). Furthermore, T2 is sub-$V_{TH}$ operated by controlling the bit-line voltage ($V_{BL}$) to amplify $\Delta V_X$ by its steep sub-threshold swing (S.S.). Besides, the sub-$V_{TH}$ operating T2 can also make $I_{TBL}$ become less subjective to the TBL clamping voltage ($V_{TBL}$). Compared to the conventional 1T1R, the H2T1R read can reduce the LRS current by 88%, reduce the HRS leakage current by 99.8%, lower $V_{WL}$ by 67%, and enhance the signal ratio (R-ratio) by 13.5 times. Notice that lowered $V_{WL}$ can not only support low VDD operation, but also reduce the voltage drop on the RRAM cell by 17% to suppress the read disturb. To perform MAC operation, the H2T1R array is row-wise accessed across different sub-array with analogue BL inputs. The $V_{TBL}$ of the selected row is clamped to 100mV to read out the accumulated $I_{TBL}$ and the rest of them are grounded. The typical H2T1R MAC behavior shows that the accumulated $I_{TBL}$ is well proportional to the number of activated LRS cells with negligible signal margin degradation due to activated HRS cells, which can be attributed to suppressed IR drop by modulating cell current and lowered HRS leakage by R-ratio enhancement.

The proposed RS-CSA (Fig.4) mainly comprises a current mirror to compare the current ($I_{COMP}$) with the reference current ($I_{REF}$), four reference (REF) branches to generate proportional reference currents ($I_{REF}[0:3]$) controlled by SEN[0:3], a leakage current ($I_{LEAK}$) compensation branch to compensate the off-cell leakage in the cell array controlled by SW[0], three reference subtracting (REFS) branches that can subtract $I_{REF}[0]$, $I_{REF}[1]$, or $I_{REF}[2]$ from the data-line current ($I_{DL}$) depending upon the switches (SW[1:3]), RS control for the timing and the subtracting logics, biasing reference voltage generator ($V_{REF\_N}$ for REF and $V_{REF\_P}$ for REFS), buffers and latches to output COUT[0:3]. In this configuration, $I_{COMP}$ equals the difference between $I_{DL}$ and the current of the REFS branches ($I_{REFS}$). A typical readout process, triggered by RD_EN, consists of one phase (PH0) for current stabling and leakage compensating and four phases (PH1-PH4) for current comparing and reference subtracting. In PH0, SW0 is triggered to compensate $I_{LEAK}$ during $I_{DL}$ stabling ($I_{COMP} = I_{DL} - I_{LEAK}$). In PH1-4, $I_{COMP}$ is sequentially compared with $I_{REF}[0:3]$ to output COUT[0:3], in which the RS control logic cumulatively subtracts the $I_{REF}[i-1]$ at the $i^{th}$ phase from $I_{DL}$ if $I_{COMP}>I_{REF}[i-1]$ till the last one. The proposed RS-CSA only needs 4 references and halved the maximum range of $I_{COMP}$, which can reduce the reference power and extend linear dynamic range of the current mirror.

The properties of the H2T1R array and RS-CSA are shown in Fig.5. A clear memory window between the HRS and LRS cells can be observed at $V_{WL}<VDD$. Notice that because 256 core T2 devices are connected to TBL, there is a total $I_{LEAK}$ about 1uA due to the off-cell leakage, which can be compensated by the RS-CSA. While $I_{TBL}$ of the LRS cell linearly increases with increasing m of T2, that of the HRS cell keeps at a negligible level below 30 nA. Thanks to reduced cell currents and IR drop, the H2T1R array performs highly linear MAC operation over channels. Because of the near-$V_{TH}$ operation of T1 and T2, the MAC power of H2T1R array can be reduced by more than 63.4% compared to its 1T1R counterpart. Besides, the RS-CSA achieves more than 79.5% power saving compared to the conventional current-mirror current sense amplifier (CM-CSA) by reducing the number of stand-by reference signals.

Figure 6 shows the measured waveform and power breakdown of the fabricated 28nm 8Kb HKMG H2T1R CIM macro. The measured latency ($T_{AC}$) is 66 ns for 4-bit readout, which can be reduced to 13 ns after sensing phase optimization. It consumes 31.96 uW at VDD = 0.8 V and achieves an energy efficiency of 30.34~154.04 TOPS/W at 1-bit input (IN), 3-bit weight (W), and 4-bit output (O). Compared to the previous work, this work demonstrated the first multi-bit H2T1R CIM macro supports highly linear analogue MAC with reduced reference generation cost.

Figure 7 shows a die photo and summary table .

**References:**

[1] W.-H. Chen *et al.*, " CMOS-integrated memristive non-volatile computing-in-memory for AI edge processors," Nature Electronics, 2019:420-428.
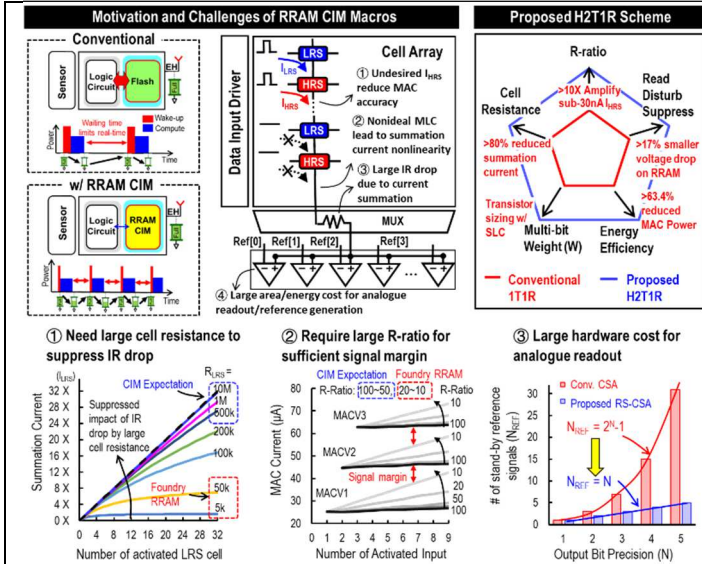
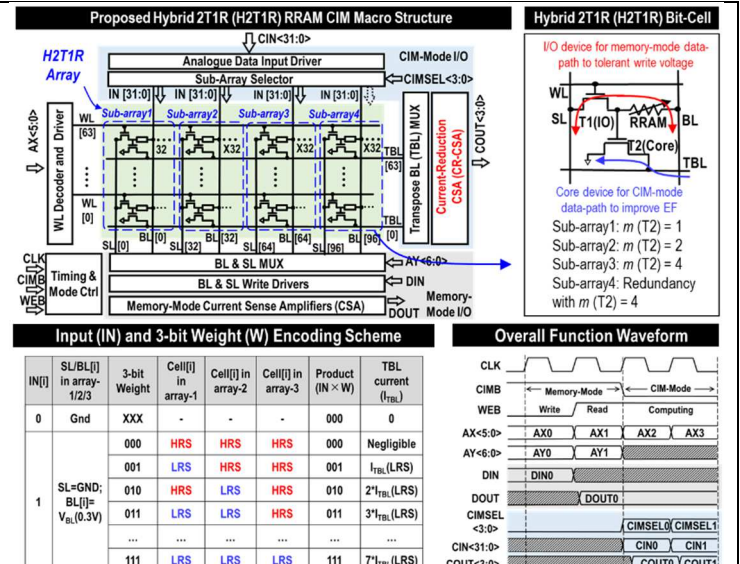Fig. 1. Motivation, challenges and advantages of the proposed H2T1R RRAM CIM macro.



Fig. 2. Proposed H2T1R RRAM CIM macro, input and weight encoding method, and its operational waveforms.
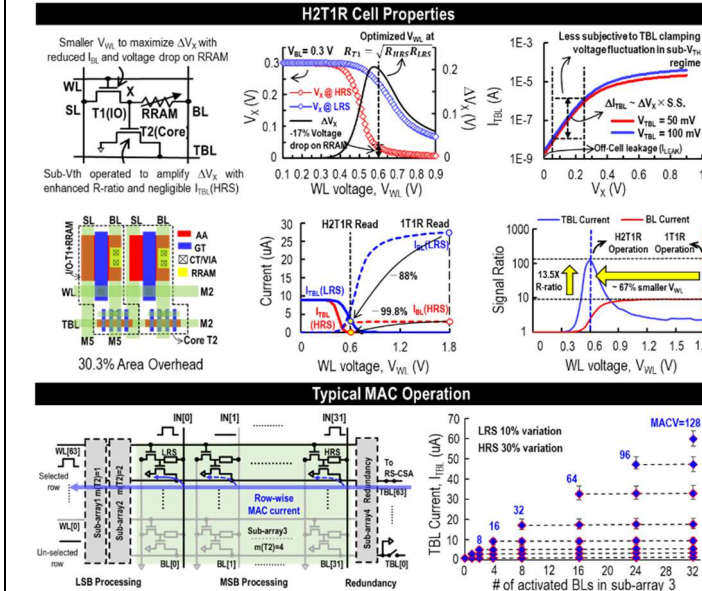


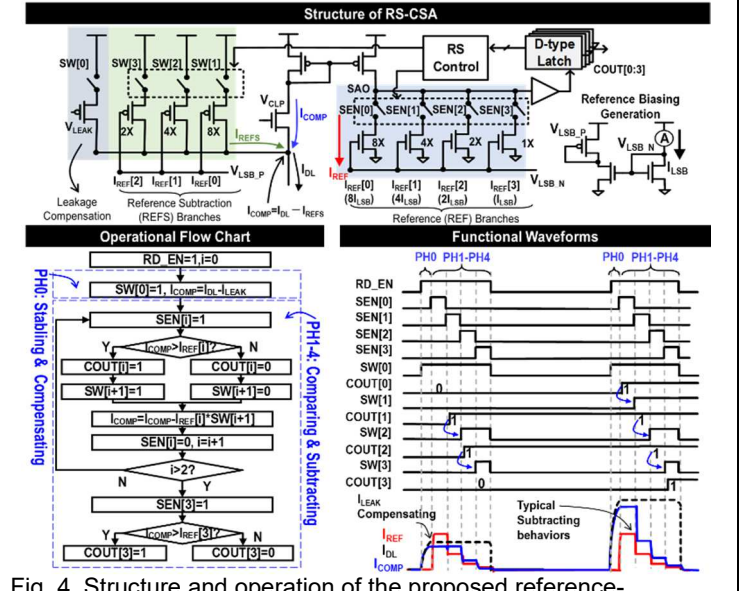Fig. 3. H2T1R cell properties and typical MAC behaviors.



Fig. 4. Structure and operation of the proposed reference-subtracting current sense amplifier (RS-CSA).
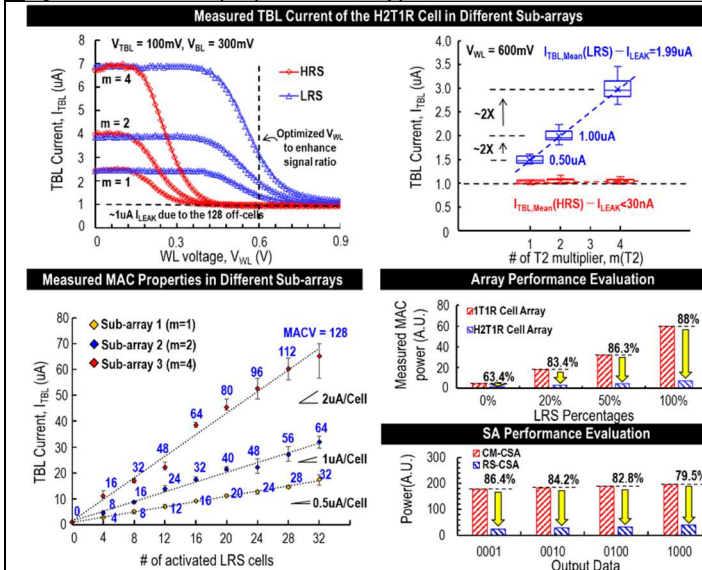


Fig. 5. Measured H2T1R cell array characteristics and performance evaluation of proposed array and sense amplifier.
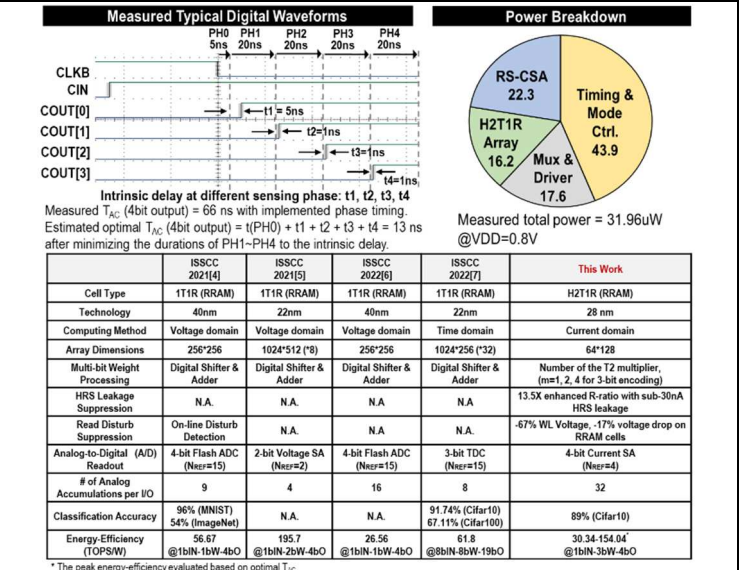


Fig. 6. Measured waveforms and power breakdown of the fabricated macro and the comparison table.
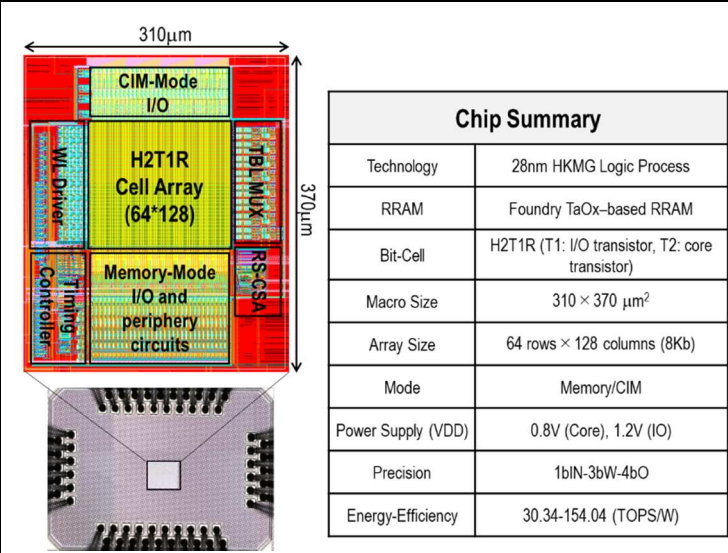
| Chip Summary | |
|---|---|
| Technology | 28nm HKMG Logic Process |
| RRAM | Foundry TaOx–based RRAM |
| Bit-Cell | H2T1R (T1: I/O transistor, T2: core transistor) |
| Macro Size | $310 \times 370 \ \mu m^2$ |
| Array Size | 64 rows $\times$ 128 columns (8Kb) |
| Mode | Memory/CIM |
| Power Supply (VDD) | 0.8V (Core), 1.2V (IO) |
| Precision | 1bIN-3bW-4bO |
| Energy-Efficiency | 30.34-154.04 (TOPS/W) |

Fig. 7. Die photo and chip summary.

**Additional References:**

[2] C.-X. Xue *et al.*, " A 22nm 2Mb ReRAM Compute-in-Memory Macro with 121-28TOPS/W for Multibit MAC Computing for Tiny AI Edge Devices," ISSCC 2020:244-245.

[3] C.-X. Xue *et al.*, " A CMOS-integrated compute-in-memory macro based on resistive random-access memory for AI edge devices," Nature Electronics 2021:81-90.

[4] J.-H. Yoon *et al.*, " A 40nm 64Kb 56.67TOPS/W Read-Disturb-Tolerant Compute-in-Memory/Digital RRAM Macro with Active-Feedback-Based Read and In-Situ Write Verification," ISSCC 2021:404-405.

[5] C.-X. Xue *et al.*, " A 22nm 4Mb 8b-Precision ReRAM Computing-in-Memory Macro with 11.91 to 195.7TOPS/W for Tiny AI Edge Devices," ISSCC 2021:246-247.

[6] S D. Spetalnick *et al.*, " A 40nm 64kb 26.56TOPS/W 2.37Mb/mm$^2$ RRAM Binary/Compute-in-Memory Macro with 4.23× Improvement in Density and >75% Use of Sensing Dynamic Range," ISSCC 2022:268-269.

[7] J.-M. Hung *et al.*, " An 8-Mb DC-Current-Free Binary-to-8b Precision ReRAM Nonvolatile Computing-in-Memory Macro using Time-Space-Readout with 1286.4 - 21.6TOPS/W for Edge-AI Devices," ISSCC 2022:182-183.