



OPEN Dual-directional CIM-based non-volatile SRAM for instant-on/off energy-constrained edge AI devices

Seyed Parsa Hemmasi¹, Mohammad Hossein Moaiyeri^{2✉}, Nima Eslami² & Behzad Ebrahimi¹

Compute-in-Memory (CIM) offers an efficient approach for accelerating DNNs by performing matrix–vector multiplications directly within memory. However, its adoption in edge devices is limited by unstable power supplies and the performance overhead of conventional row- or column-wise computing. This paper presents a two-directional CIM-based nvSRAM cell that performs both row- and column-wise operations, enabling faster and more efficient matrix–vector multiplication. The proposed design stores the CIM outputs within the same computation cycle, referred to as Simultaneous Compute and Write (SCW), thereby reducing latency during complex neural network inference. By integrating a single I-MTJ into each SRAM cell, it also provides reliable data retention and restoration during power failures, making it well-suited for low-power, energy-constrained edge applications. Post-layout simulations were conducted to evaluate the proposed architecture. The detailed post-layout simulation results demonstrate a 31% improvement in write margin, a 40% reduction in PDP in memory mode, and an 85% reduction in energy in backup scenarios, compared to state-of-the-art designs. Furthermore, the proposed design achieves a 39.2% EDP reduction during neural network inference operation under power instability, highlighting its suitability for low-power edge computing.

Keywords NvSRAM, Compute-in-memory, Dual-direction computing, Intermittent computing, Edge devices

The combination of Internet of Things (IoT) and Artificial Intelligence (AI) has opened the door to many smart applications, such as health monitoring, autonomous systems, and speech recognition^{1,2}. These systems often rely on deep neural networks (DNNs), which require constant operation under strict power constraints. This poses challenges for both energy efficiency and system flexibility³.

Compute-in-Memory (CIM) offers a promising solution by performing matrix–vector multiplications (MVMs) directly within memory, significantly reducing the energy consumption of data movement⁴. As a result, CIM is emerging as a key enabler of low-power smart IoT applications, particularly in edge devices that operate under constrained or unstable power conditions⁵. Most CIM designs today are based on either SRAM or non-volatile memory (NVM) technologies.

Although SRAM-based CIM architectures demonstrate promising performance, for example,⁶ achieves high-throughput CIM,⁷ implements 1-bit multiplication using an 8T transistor, and⁸ attains 274 TOPS/W with an SRAM-based macro, they are inherently volatile and require restarting from the first network layer after a power failure⁹. Moreover, static power consumption in SRAM cells becomes a significant contributor to overall power dissipation, especially in deep nanoscale nodes^{10,11}. To overcome power constraints in energy-limited environments, NVM-based CIM architectures have been investigated for their inherent state retention and reliable operation under intermittent power conditions. However, these designs face a critical challenge in memory write operations, which are dominated by frequent weight reloading and activation buffering during AI inference, causing a significant performance bottleneck¹². For instance, in 130 nm technology, the write energy for implementing *TinyML* using NVM is approximately 13 × higher than that of SRAM for 8-bit data and nearly 151 × greater than the energy required for a single multiply–accumulate (MAC) operation, underscoring the significant energy overhead of NVM writes¹³.

To overcome these challenges, non-volatile SRAM (nvSRAM) cell architectures have attracted significant attention by combining the high-speed operation and dynamic energy efficiency of SRAM with a non-volatile element. This integration enables near-zero static power consumption while retaining intermediate data during power interruptions.

¹Department of Electrical Engineering, SR.C., Islamic Azad University, Tehran, Iran. ²Faculty of Electrical Engineering, Shahid Beheshti University, Tehran 1983969411, Iran. ✉email: h_moaiyeri@sbu.ac.ir

In¹⁴, two 9T-2MTJ and 10T-2MTJ nvSRAM cells were proposed to improve hold and write margins while reducing power. In¹⁵, two nvSRAM designs using different MTJ types are presented. The I-MTJ-based design reduces cell area by avoiding transistor upscaling, leveraging efficient switching via spin Hall effect (SHE) currents. In¹⁶, two novel 8T-2MTJ and 9T-2MTJ nvSRAM cells were proposed, in which the store operation is performed using SHE-assisted write in both MTJs.

Although previous designs have aimed to enhance nvSRAM performance, they rely on the conventional 6T SRAM structure, whose stability becomes increasingly limited in submicron technologies, particularly in FinFETs, due to quantized fin dimensions¹⁷. Furthermore, none of the prior designs effectively support CIM, as the 6T structure's shared read/write path constrains CIM performance.

To improve nvSRAM stability,¹⁸ presented a 9T1R cell that leverages the Schmitt-trigger behavior during read operations to enhance read robustness. Similarly,¹³ presented an 8T1R RRAM-based nvSRAM cell with CIM capability. However, the 6T SRAM core used in the design remains constrained by stability issues arising from the intrinsic conflict between read and write operations, as well as read-disturbance problems in submicron technologies¹⁹. This issue is especially critical in AI inference applications, where intermittent read/write operations in the absence of sufficient cell stability can lead to inaccurate outputs.

To the best of our knowledge, all previously proposed nvSRAM cells supporting CIM can perform operations in only one direction. This limitation introduces additional overhead when implementing MVM, which is essential for DNN computations²⁰.

To this end, we propose a novel nvSRAM cell that not only enhances the stability of the SRAM core but also supports CIM operations in both horizontal and vertical directions. This dual-direction capability enables complex matrix transposition via simple write and read operations in opposite directions. Moreover, the proposed design employs a single I-MTJ to store and restore SRAM data during power failures, ensuring non-volatility with minimal area overhead.

The key contributions of the proposed design are as follows:

- A novel SRAM structure that significantly enhances stability parameters, including read and write margins.
- Utilization of a single I-MTJ for efficient data storage and restoration during power failures.
- Support for CIM operations in both horizontal and vertical directions, enabling full MVM functionality essential for DNN applications.
- In-situ storage of CIM results within the same computation cycle, termed Simultaneous Compute and Write (SCW), for both horizontal and vertical directions, enhancing overall computational efficiency.
- Support for content-addressable memory (CAM) operations, enabled by the innovative read path design.

The remainder of this paper is organized as follows: “[Preliminaries](#)” section provides an overview of I-MTJ fundamentals and offers a detailed review of existing nvSRAM cell designs. “[Proposed architecture](#)” section introduces the proposed nvSRAM structure and describes its operation in different modes. “[Evaluation results](#)” section evaluates the proposed design, compares it with existing counterparts, and demonstrates its effectiveness in real-world applications. Finally, “[Conclusion](#)” section concludes the paper by summarizing the key contributions.

Preliminaries

Magnetic tunnel junction

Magnetic Tunnel Junctions (MTJs) are spintronic devices with a multilayer structure, where an ultra-thin insulating barrier is sandwiched between two ferromagnetic layers²¹. The MTJ comprises a fixed layer with a constant magnetic orientation and a free layer that can switch between parallel (*P*) and anti-parallel (*AP*) orientations relative to the fixed layer. The MTJ shows lower resistance (R_P) in the parallel state and higher resistance (R_{AP}) in the anti-parallel state²².

MTJs are categorized by the orientation of their magnetic layers: perpendicular to the substrate, called perpendicular MTJs (P-MTJs), or parallel to the substrate, called in-plane MTJs (I-MTJs)²³.

We have utilized I-MTJ due to its advantages over the P-MTJ, including the mitigation of read disturbance, a higher tunnel magnetoresistance (TMR) ratio for improved data stability, and enhanced reliability and thermal robustness^{24,25}. The TMR reflecting the difference between the R_{AP} and R_P resistances of an I-MTJ is quantified as shown in (1)²⁶. Notably, MTJs with TMR values exceeding 600% have been successfully fabricated at room temperature^{27–29}.

$$TMR = \frac{R_{AP} - R_P}{R_P} \times 100 \quad (1)$$

To switch the resistance of the I-MTJ from R_{AP} to R_P or vice versa, a specific current must pass through the heavy metal, flipping the free-layer orientation. Notably, this current must exceed the critical current (I_C) to write data to the I-MTJ successfully³⁰. The I_C current is dependent on the nanomagnet's volume and its material-specific properties. The critical current can be mathematically calculated using (2).

$$I_C = \frac{2q\alpha M_S V_{MTJ}}{\hbar} (H_k + 2\pi M_S) \quad (2)$$

where q represents the electron charge, α is the magnetic damping coefficient, V_{MTJ} is the volume of the free layer, and M_S , H_k , and \hbar denote the saturation magnetization, anisotropy field, and reduced Planck's constant,

respectively³¹. As illustrated in Fig. 1, the current flow direction between T_2 and T_3 determines whether the I-MTJ is switched to a *P* or *AP* state.

When the current flows from T_2 to T_3 , the device is programmed into the *P* state (Fig. 1a); conversely, when the current flows from T_3 to T_2 , the device is programmed into the *AP* state (Fig. 1b).

Apart from the switching current and TMR, one of the key parameters of the I-MTJ is its thermal stability factor (Δ), which determines the data retention of magnetic memory arrays³². The retention time can be mathematically calculated using (3).

$$\text{Thermal stability factor}(\Delta) = \frac{H_k \times M_s \times V}{2 \times k_B \times T} \quad (3)$$

where V is the volume of the free layer, T is the operating temperature, H_k is the anisotropy field, and M_s is the saturation magnetization.

Moreover, the integration of MTJs with FET transistors has been successfully demonstrated and commercialized within a single-chip platform^{33,34}. Demonstrations by leading institutions, such as ITRI and TSMC, highlight the practical feasibility of this hybrid integration, representing a significant milestone in the advancement of spintronic-FET technology³⁵.

Previous work

The rising demand for energy-efficient systems, particularly in intermittently powered edge AIoT applications, has spurred extensive research into low-power memory and AI accelerators³⁶. In this context, several nvSRAM cell designs have been proposed to address these challenges.

In¹⁴, two nvSRAM designs are presented, each incorporating two P-MTJs per cell and using SHE-assisted writing. Although SHE-assisted current improves performance and reliability over STT alone, the reliance on a conventional 6T SRAM for the volatile part results in challenges at advanced technology nodes¹⁶. Furthermore, the requirement for two separate voltage levels during store-and-restore operations increases circuit complexity.

In¹⁵ two nvSRAM designs with a shared topology are proposed, differing in the type of MTJ employed. The first design uses two P-MTJs programmed by SHE-assisted STT currents but suffers from high area overhead in advanced technology nodes. The second design adopts two I-MTJs, switchable via SHE current, enabling smaller transistors and improved TMR for enhanced stability and robustness against process variations. However, its coupled read/write path in the SRAM core section increases susceptibility to data corruption, limiting its suitability for in-memory computing applications.

In¹⁶, two nvSRAM cells, 8T-2MTJ and 9T-2MTJ, using dual MTJs are proposed. By leveraging SHE-assisted STT currents to write SRAM data into MTJs, the design eliminates the need for external write circuitry, thereby reducing bit-cell complexity. Array-level evaluations show a 32% improvement in power efficiency. While the cells offer high-speed operation and low dynamic power, the 6T SRAM core limits stability at nanoscale technologies, and the dual-MTJ storage increases power consumption during store operations.

To reduce power consumption during store mode,³⁷ proposes a cost-efficient nvSRAM cell that uses a single P-MTJ combined with a 1T-1MTJ shadow structure for dynamic reconfiguration. However, the store operation involves an additional read step, transferring data from the MTJ to the bit-lines before writing it to the storage nodes, which further decreases efficiency. Furthermore,¹³ presents a novel 8T1R nvSRAM cell capable of supporting CIM operations. The proposed design incorporates neural-architecture-aware and sparsity-aware backup mechanisms, enabling flexible data retention strategies and efficient in-memory multiplication. The design utilizes a conventional 6T SRAM cell as its core volatile storage element.

To improve the scalability and performance of the nvSRAM cell design,¹⁸ proposes a novel architecture that exploits Schmitt-trigger behavior during read operations to enhance overall cell stability. Furthermore, the design performs both store and restore operations using a single RRAM device.

Except for¹⁸, the reviewed works rely on conventional 6T SRAM cells as volatile storage, which face performance limitations at nanoscale nodes due to shared read/write paths. Also, most designs use two non-volatile elements (MTJ or RRAM) for storing and restoring data, which substantially increases power consumption due to their high write energy.

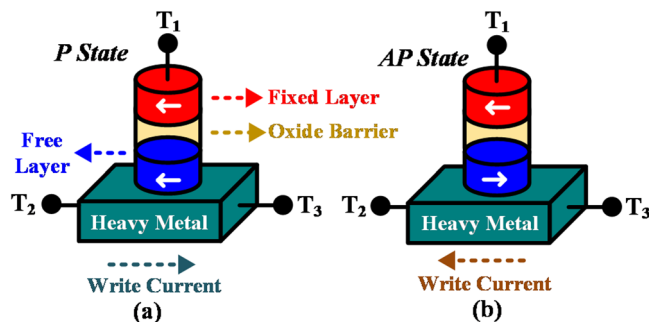


Fig. 1. The structure of I-MTJ (a) parallel (P), (b) anti-parallel (AP).

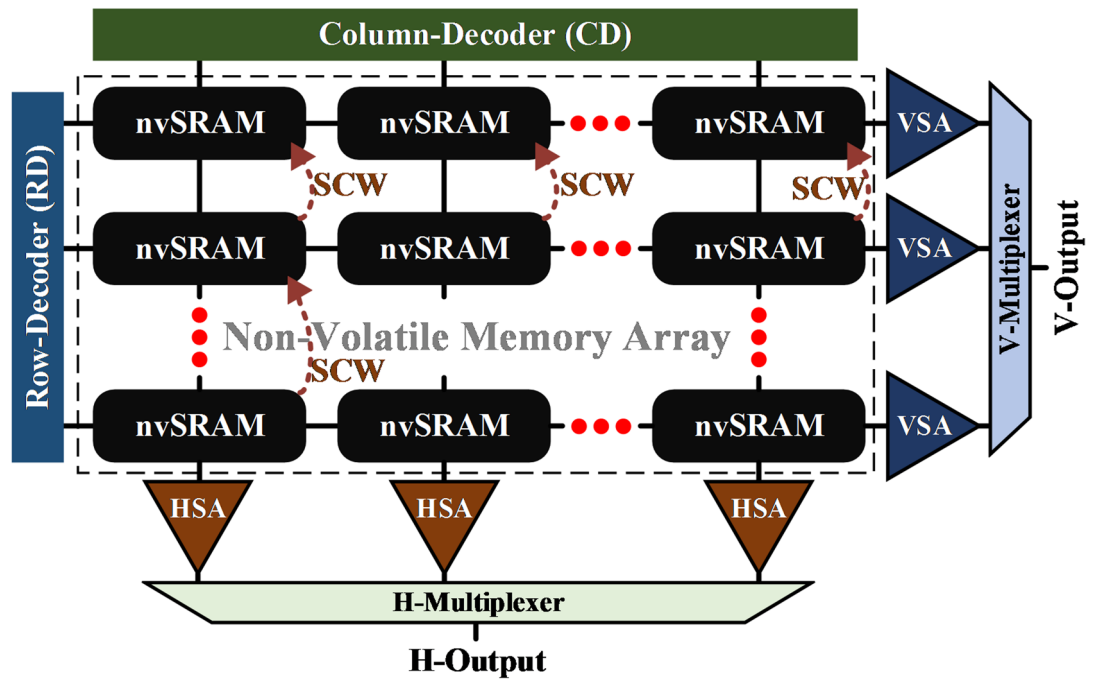


Fig. 2. The proposed dual-directional nvSRAM architecture.

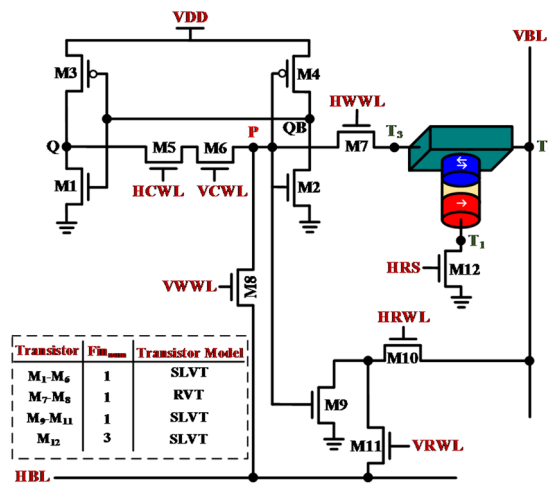


Fig. 3. Schematic of the proposed dual-directional NVSRAM.

Proposed architecture

The proposed nvSRAM architecture is presented in Fig. 2. This architecture supports both row- and column-based read/write operations, as well as bidirectional in-memory computing. The architecture enables in-situ storage of computation results (SCW) within the same cycle in any selected cell along the row or column where the operation occurs. The array comprises a Row Decoder (RD), a Column Decoder (CD), a Vertical Sense Amplifier (VSA), and a Horizontal Sense Amplifier (HSA).

Figure 3 illustrates the schematic of the proposed non-volatile SRAM cell. As depicted in the figure, transistors M_1 - M_4 , together with M_5 and M_6 , are configured to form the cross-coupled inverters, where M_5 and M_6 serve as the horizontal and vertical feedback-cutting transistors, respectively. Transistors M_7 and M_8 , controlled by the write word lines, establish horizontal and vertical access paths. The gates of M_{10} and M_{11} transistors are driven by the horizontal and vertical read word lines, respectively. The Q node is coupled to transistor M_9 to complete the read/compute operation. Finally, transistor M_{12} is governed by the horizontal restore signal. An integrated I-MTJ device serves as the non-volatile element in this design. The proposed dual-direction, single-ended non-volatile SRAM, featuring decoupled read and write paths, operates in two primary modes: memory mode and computing mode. A detailed explanation of these modes is provided in the following subsections.

Memory mode

The proposed cell combines the advantages of SRAM, offering energy- and time-efficient read/write operations, with those of NVM, enabling data retention during power-off. As a result, it can perform the standard hold, read, and write operations of an SRAM cell, as well as store and restore data to and from the I-MTJ connected to the storage node. Table 1 summarizes the overall signaling of the proposed cell across different memory operation modes.

The following sections provide a detailed discussion of each of its five memory operations:

- (1) *Hold mode* During the hold operation, only the VCWL and HCWL signals are activated, which turn on M_5 and M_6 and complete the feedback loop that forms the cross-coupled inverters. Simultaneously, all other signals are deactivated (grounded), thereby isolating the Q and QB nodes and, consequently, the stored data from the bit lines.
- (2) *Read operation* The read operation can be performed either row-wise or column-wise, depending on the appropriate control signals. During the row-wise read operation, activating the horizontal read word line (HRWL) turns on access transistor M_{10} . Subsequently, depending on the logic level stored at node Q, transistor M_9 is either activated or remains off. As a result, the pre-charged vertical bit-line (VBL) is either discharged or remains unaffected, corresponding to the stored value at Q. The column-wise read operation follows a similar mechanism, wherein the vertical read word line (VRWL) is asserted instead of the horizontal read word line (HRWL).
- (3) *Write operation* In the row-wise write operation, asserting the horizontal write word line (HWWL) turns on transistor M_7 , while setting the vertical control word line (VCWL) to ground disables the feedback loop of the cross-coupled inverters, effectively floating the Q node. This condition permits the data applied to the vertical bit-line (VBL) to be reliably written into the Q node without contention.

The column-wise write operation employs the same mechanism, except that the vertical write word line (VWWL) is activated instead of HWWL, and the horizontal control word line (HCWL) is driven low in place of VCWL. It is noteworthy that, to mitigate the half-select issue, the transistors involved in the write operation are activated through both vertical and horizontal control signals.

- (4) *Store operation* In the event of power loss or intermittent power supply, data from the volatile section can be stored in an I-MTJ within the non-volatile section. In this mode, the VBL is set to $V_{DD}/3$, and only the HWWL is activated. If $Q = '1'$, current flows from Q to VBL and then from T_3 to T_2 of the I-MTJ, setting its resistance to the R_{AP} state. Conversely, if $Q = '0'$, current flows from VBL to Q and then from T_2 to T_3 , setting its resistance to the R_p state, as illustrated in Fig. 4.

It is worth noting that the current required to switch the I-MTJ resistance (per Eq. 2) exceeds $6 \mu A$. During the SRAM write operation, however, the current is limited to ~ 200 nA due to the M6 transistor being deactivated, preventing any change in the I-MTJ state.

- (5) *Restore operation* The restore operation transfers data from the non-volatile section to the volatile section when power is restored or stabilized. The operation is based on a voltage divider formed between a reference resistor shared among cells in the same column and the I-MTJ. The value of the reference resistor is chosen to lie between the resistances of the MTJ's R_{AP} and R_p states. During the restore operation, activating the HRS turns on transistor M12, while the horizontal bit-line restore signal HBLRS enables the voltage-divider path. Simultaneously, HWWL is set to '1' and HCWL to '0', allowing the voltage divider output to drive the Q node. If the I-MTJ is in the R_{AP} state, the Q node is driven to a higher voltage (restoring a '1'), while the R_p state drives it to a lower voltage (restoring a '0'). Notably, transistor M_{12} features three fins to minimize its RON impact on the voltage division, as shown in Fig. 5. Notably, this design is the first to use a single I-MTJ

Operations	Hold	Read		Write		Store	Restore
		Row-wise	Column-wise	Row-wise	Column-wise		
HWWL	GND	GND	GND	V_{DD}	GND	V_{DD}	V_{DD}
HRWL	GND	V_{DD}	GND	GND	GND	GND	GND
VCWL	V_{DD}	V_{DD}	V_{DD}	GND	V_{DD}	V_{DD}	V_{DD}
VBL	GND	Pre-Charged	GND	Data ('0' or '1')	GND	$V_{DD}/3$	Float
VWWL	GND	GND	GND	GND	V_{DD}	GND	GND
VRWL	GND	GND	V_{DD}	GND	GND	GND	GND
HCWL	V_{DD}	V_{DD}	V_{DD}	V_{DD}	GND	V_{DD}	GND
HBL	GND	GND	Pre-Charged	GND	Data ('0' or '1')	GND	Float
HRS	GND	GND	GND	GND	GND	GND	V_{DD}

Table 1. Signaling of the proposed design in the memory mode.

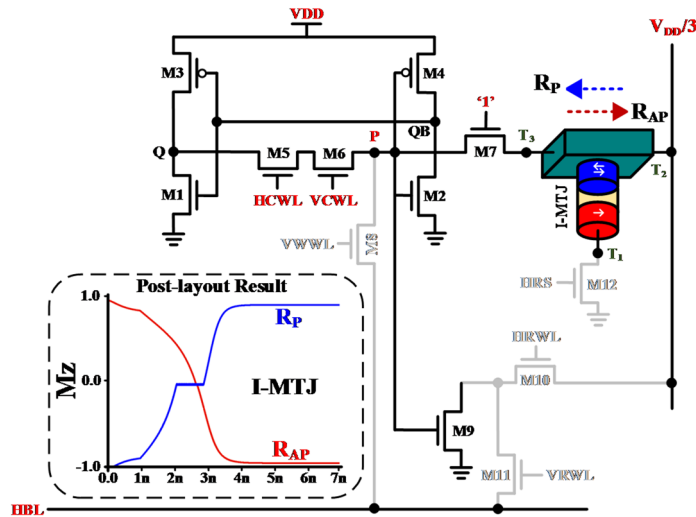


Fig. 4. Store operation in the proposed nvSRAM.

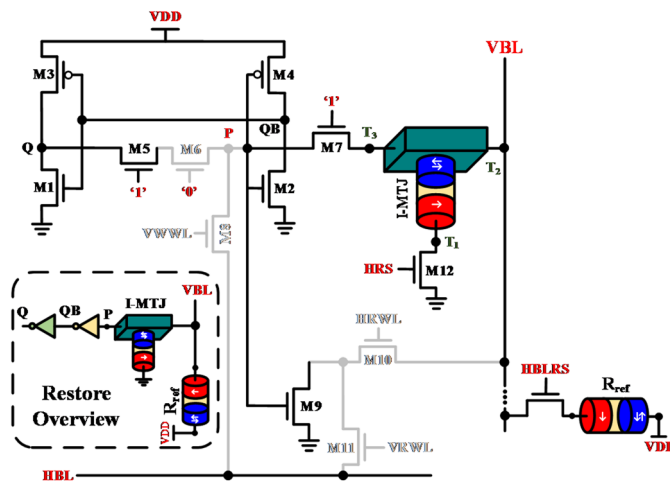


Fig. 5. Restore operation in the proposed nvSRAM.

per cell for both storing and restoring data, whereas previous designs typically use two MTJs per cell. By employing a shared reference resistor rather than a second MTJ, this approach offers a more cost-effective solution than earlier nvSRAM cells.

Computing mode

The proposed cell supports bidirectional Boolean logic operations (both row- and column-based), enabling efficient matrix–vector multiplication with lower latency and power consumption compared to conventional unidirectional (primarily column-based) in-memory computing architectures. In addition, the design facilitates content-addressable memory (CAM) functionality.

- (1) *Boolean Logic Operations* The proposed cell can perform NOR/OR and NAND/AND operations in both row- and column-wise directions. Figure 6 illustrates the implementation of a 2-input column-wise and row-wise NOR/OR operation in the proposed design.

During this operation, the horizontal read word lines (HRWLs) corresponding to the operand cells are activated, and the VBL is pre-charged. Depending on the data stored in the columns, if the cells contain a logic ‘0’, the corresponding M₉ transistors in the nvSRAM cells remain off, allowing VBL to retain its pre-charged value. If one or more cells store a logic ‘1’, the associated M₉ transistors are activated, leading to the discharge of VBL. The sense amplifier connected to the end of VBL then detects the resulting voltage level, thereby completing the NOR/OR operation. The row-wise NOR/OR operation functions similarly but utilizes VRWL and HBL to target data within the same row.

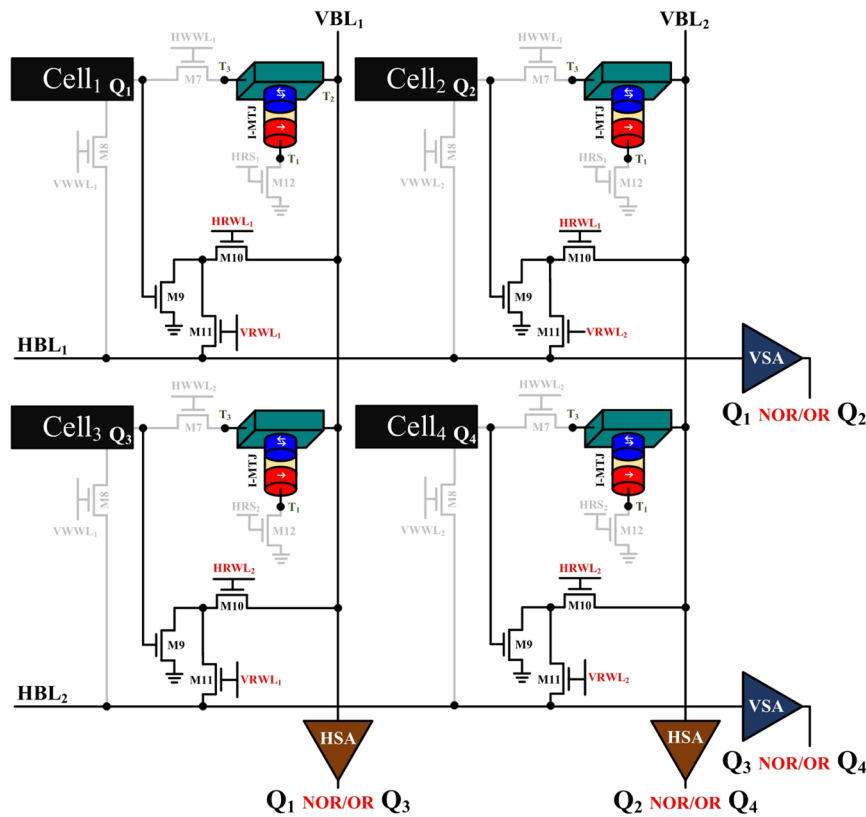


Fig. 6. Implementation of row- and column-wise NOR/OR operations in the proposed nvSRAM array.

The NAND/AND operations closely resemble the NOR/OR operations; however, in this case, the read word lines (VRWLs and HRWLs) are activated for a shorter duration. This limited activation time prevents the bit-line from fully discharging unless both operands are logic ‘1’. Consequently, the cases ‘01’, ‘10’, and ‘11’ can be distinguished, and the resulting bit-line voltage accurately reflects the NAND/AND output for the two selected rows or columns.

It is worth noting that the proposed design fully supports in-situ storage of computation results within the target cell during the same cycle as the logic operation (SCW). This is achieved by configuring the target row or column in write mode concurrently with the logic operation. As a result, there is no need for a separate write-back, and the overhead associated with additional decoders and control signals is eliminated.

Consequently, the proposed design’s SCW capability greatly enhances energy and time efficiency for intermediate computations, such as those in neural network hidden layers, compared to prior architectures.

(2) *CAM Operation* The proposed architecture supports both row- and column-wise CAM operations. In column-wise CAM implementation, the search data and its complement are applied to the HWWL and HRWL, respectively. Meanwhile, the VBL functions as the match line (ML), delivering the output corresponding to the match results. In this operation, the ML is pre-charged to the supply voltage. If the search data matches the stored data, the ML retains its pre-charged level; otherwise, it is discharged, indicating a mismatch (Fig. 7a). The row-wise CAM operation follows a similar principle; however, the VWWL, VRWL, and HBL are used to apply the search data, its complement, and serve as the ML, respectively, enabling data processing within the same row (Fig. 7b).

It is worth noting that, during parallel CAM searches, data corruption may occur if the VBL or HBL drops too rapidly. The worst-case scenario arises when all search data bits are logic ‘1’ (e.g., 111...1), but only a single stored bit is logic ‘1’ (e.g., 100...0). To mitigate this issue, the proposed cell employs deliberately weakened access transistors with reduced threshold voltages, thereby limiting their drive strength and preventing abrupt voltage drops on the bit lines. Additionally, during CAM operations, these transistors can be biased at a lower voltage than the supply voltage to further suppress the risk of data corruption.

Evaluation results

To evaluate the proposed nvSRAM architecture, the design was simulated in Cadence Virtuoso using the 7nm FinFET technology developed by the Arizona State University in cooperation with ARM^{38,39}, along with an experimentally approved Verilog-A compact model of an in-plane (I-MTJ)^{40,41}. The post-layout results were

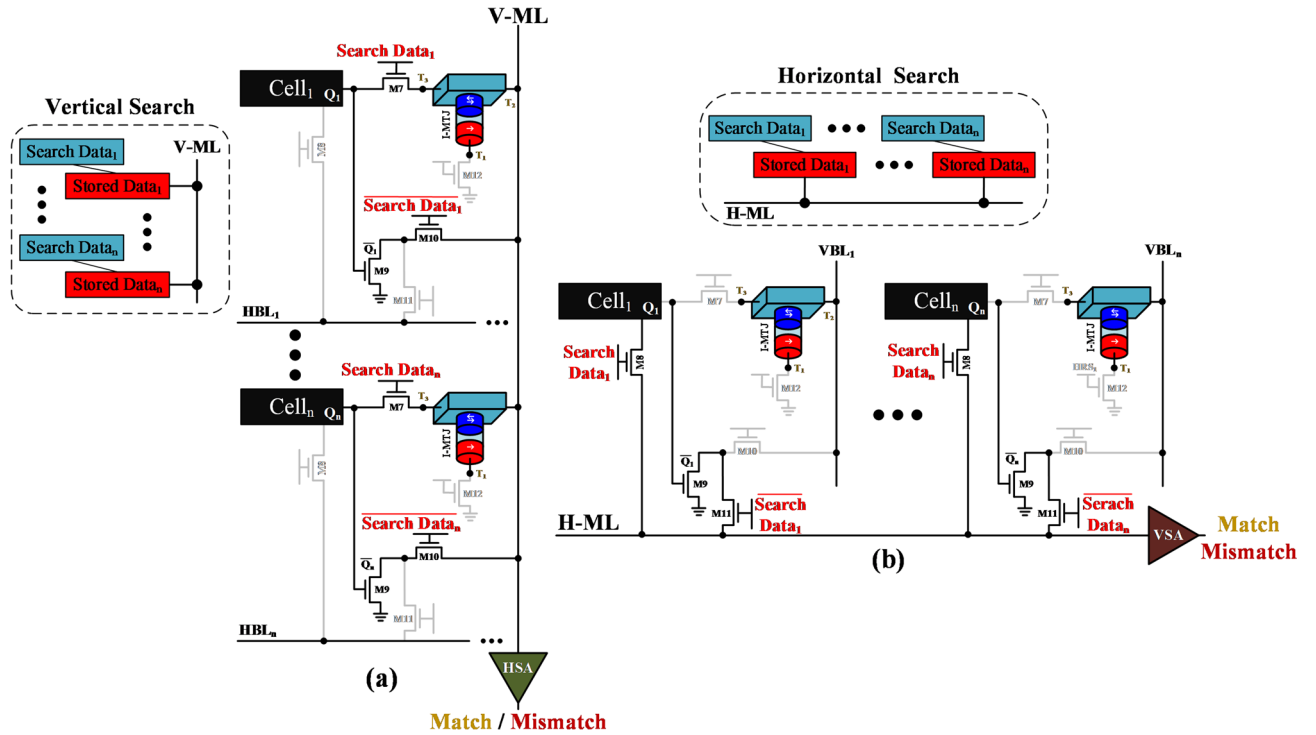


Fig. 7. CAM Operation (a) column-wise (b) row-wise.

Parameter	Value
I-MTJ	
Free layer thickness (T_p)	1.4 nm
The thickness of the MgO barrier layer (T_b)	1.2 nm
The elliptical surface of the free layer (S)	$110 \times 60 \times \frac{\pi}{4} \text{ nm}^2$
Dimensions of the metal strip (l, w, d)	$120 \times 70 \times 1.8 \text{ nm}^3$
MTJ Resistance-Area product (RAP)	$10 \text{ } \Omega \text{ } \mu\text{m}^2$
Tunnel magnetoresistance (TMR)	500%
Parallel (P) resistance (LRS)	1.98 k Ω
Potential barrier height of MgO oxide (ϕ)	0.4 V
The saturation magnetization (M_s)	$10.8 \times 10^5 \text{ A m}^{-1}$
Anisotropy field (H_k)	$3.17 \times 10^4 \text{ A m}^{-1}$
Spin-hall angle (η)	0.3
Thermal stability	37.96
Retention time	~ One year
FinFET 7nm	
Fin thickness (T_{fin})	6.5 nm
Fin height (H_{fin})	34 nm
Fin pitch (P_{fin})	27 nm
Gate length (L)	21 nm
Oxide thickness (T_{ox})	2.1 nm
Equivalent oxide thickness (EOT)	1 nm

Table 2. The key parameters of I-MTJ and FinFET technologies.

extracted for a 4 Kb array configuration. Table 2 summarizes the key parameters of both the I-MTJ and the 7nm FinFET technology used in this work.

Notably, a TMR of 500% is assumed for the I-MTJs. It is supported by^{27–29}, which reported the fabrication of I-MTJs with TMRs exceeding 600% at room temperature, whereas the value considered in this work is lower.

To provide a comprehensive performance evaluation, comparisons were conducted with state-of-the-art nonvolatile SRAM cells, including the 10T-2MTJ architecture reported in¹⁴, the I-MTJ- and P-MTJ-based designs from¹⁵, and the 8T-1R design from¹³. All designs were post-layout simulated using Cadence under the same 7nm technology node and array size to ensure a fair comparison. Additionally, the parameters of all I-MTJ and P-MTJ devices were carefully calibrated to achieve equivalent retention times across the evaluated designs¹³. Furthermore, Fig. 8 presents the layout of the proposed cell within a 64×64 array, along with an individual-cell view, implemented in 7nm FinFET technology.

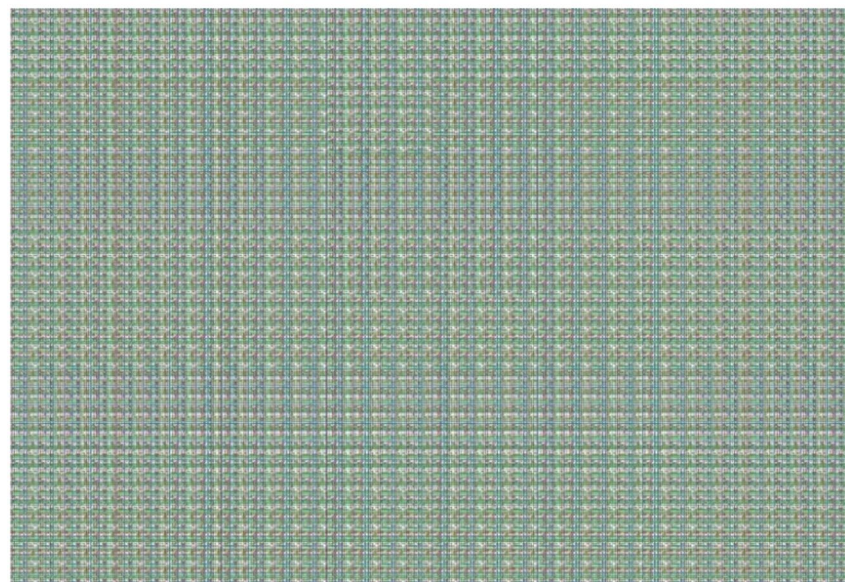
Memory mode analysis

This section evaluates the proposed nvSRAM cell under Memory Mode operation.

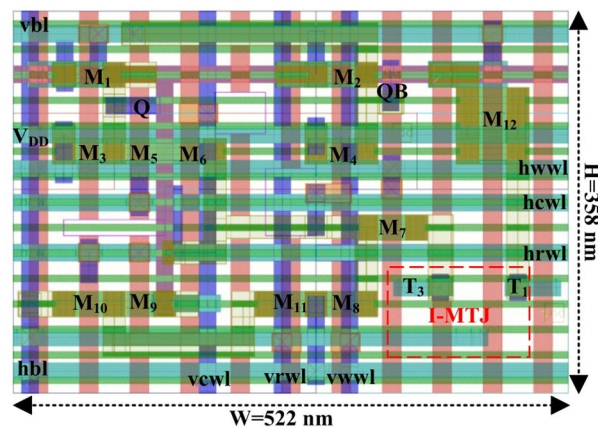
- (1) *Transient waveform analysis* To assess the functionality of the proposed design operating in memory mode, a transient waveform was obtained from post-layout simulations, as illustrated in Fig. 9. The operation sequence begins with a hold phase, followed by a write operation to the SRAM cell.

Subsequently, the stored data is read to verify the integrity of the write operation. Before power-down, the SRAM content is backed up in situ to the I-MTJ. Upon power restoration, the stored data in the I-MTJ is successfully recovered and restored to the SRAM cell.

- (2) *SRAM mode analysis* To benchmark the proposed nvSRAM against existing designs in SRAM mode, the cell's stability is assessed using the static noise margin (SNM), while its delay, power, and energy consumption are evaluated within a 4 Kb memory array. The corresponding post-layout simulation results are summarized in Table 3. Notably, since the proposed nvSRAM cell supports memory-mode operations in both



(a)



(b)

Fig. 8. Layout of the proposed NVSRAM (a) 4 Kb array (b) cell.

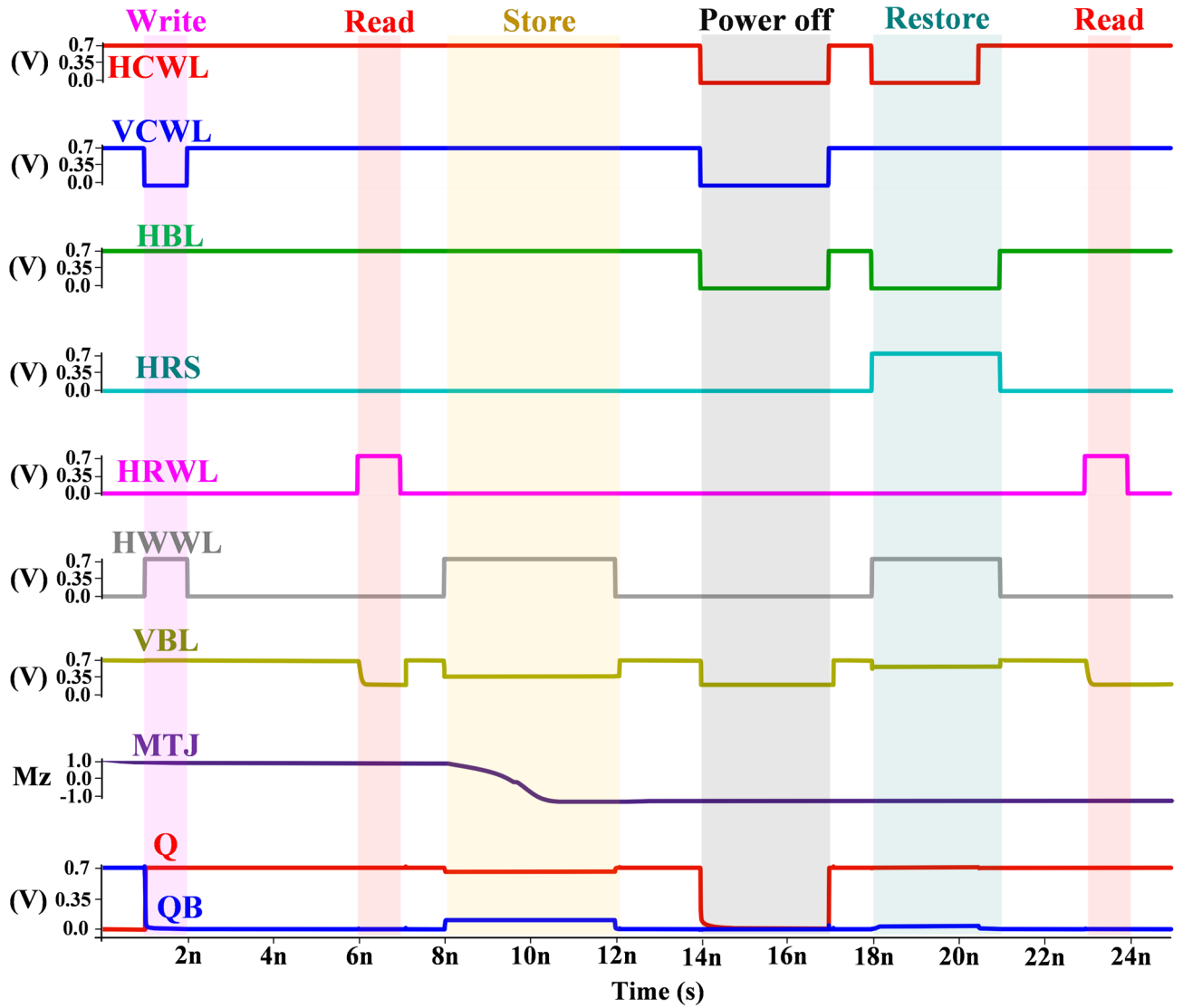


Fig. 9. Transient waveform of the proposed design in different operational modes.

Parameters	¹⁴	I-MTJ ¹⁵	P-MTJ ¹⁵	¹³	Proposed
HSNM (mV)	217	221	219	219	222
Half-select (mV)	54	57	51	56	110
RSNM (mV)	54	57	51	56	222
WSNM (mV)	208	217	224	225	372
Read delay (ps)	10.9	10.5	10.7	10.4	38.2
Read power (μ W)	572.7	76.4	500.3	150.0	36.4
Read PDP (fJ)	6.2	0.8	5.4	1.6	1.4
Write delay (ps)	11.9	10.9	11.6	8.8	31.8
Write power (μ W)	333.8	44.6	292.5	91.9	11.0
Write PDP (fJ)	4.0	0.5	3.4	0.8	0.3
Leakage power (μ W)	813.4	105.6	668.4	214.1	164.2
Bit cell area (μm^2)	0.218	0.141	0.260	0.076	0.186

Table 3. The post-layout simulation results of different nvSRAM cells in the 4 Kb memory array.

row- and column-wise directions, the reported parameters represent the worst-case values obtained from evaluations in both directions.

As shown, the proposed design exhibits higher stability parameters (HSNM, RSNM, WSNM, Half-selected SNM) than the referenced cells. This improvement is attributed to the separation of read and write paths and to the application of the feedback-cutting technique during the write process. The proposed design exhibits approximately 74.3% and 39.5% improvements in overall RSNM and WSNM, respectively, compared to the existing designs. The proposed design shows a slightly higher read delay than its counterpart, primarily due to the use of a minimum-size read transistor and a single-ended read operation.

However, the overall PDP of the read operation in the proposed design is lower than that of the compared cells, attributed to reduced power consumption. A similar pattern is observed during the write operation. The proposed design incurs a higher write delay, primarily due to its non-differential structure and relatively weak access transistors. However, it demonstrates significantly lower write power. During the write process, the feedback-cutting transistors isolate node Q, leaving it floating; as a result, the access transistor does not contend with the cross-coupled inverter transistors. This mechanism leads to the lowest write power consumption among all evaluated designs, thereby achieving the minimum write PDP. The proposed design achieves an overall PDP approximately 40% lower than that of the compared cells.

In terms of static power during the hold mode, the proposed design achieves the second-lowest leakage power among the evaluated cells. The analysis confirms that the proposed nvSRAM cell offers robust performance and reliability, positioning it as a promising solution for low-power, resource-constrained edge applications.

Moreover, as shown in Table 3, the per-bit area overhead of the proposed design and the compared designs has been evaluated and extracted using Cadence. It is worth noting that the MTJ is fabricated on top of the transistors using a back-end-of-line (BEOL) process and, therefore, does not contribute to the overall circuit area⁴². The results indicate that the proposed design has a larger area compared to the designs in¹³ and¹⁵. However, despite its higher per-bit area, the proposed design's bidirectional IMC capability enables matrix transposition by simply reading data in the reverse direction, without requiring additional memory. Consequently, although the proposed design occupies more area at the cell level, it achieves a lower overall area at the system level.

- (3) *In-situ store and restore analysis* The proposed design enables the storage and restoration of volatile data within the I-MTJ, offering significant power savings, particularly beneficial for edge devices operating under limited or unstable power conditions. Table 4 summarizes the non-volatile characteristics, namely, the store and restore operations, of the evaluated nvSRAM cells.

As shown in the table, the proposed design exhibits the lowest store power, attributed to the use of a single I-MTJ rather than a pair, coupled with a circuit topology optimized using minimum-sized transistors. Moreover, the proposed design achieves approximately 85.7% lower PDP during the store operation compared to the other evaluated designs.

Moreover, during the restore operation, although the proposed design ranks second in PDP after the¹³ design, it achieves a 48% faster restore time compared to this design. In addition, the RRAM-based design in¹³ requires erasing the RRAM cell before each store operation, implying that storing a logic '1' must be preceded by a logic '0' store cycle.

Computing mode analysis

The NOR/OR and NAND/AND operations are performed by activating two rows or two columns of cells, enabling bidirectional computing. Figure 10 shows the delay and energy consumption for these operations at different supply voltages.

The evaluation involves executing the logic operation alongside writing the computed CIM data into the memory cell within the same cycle, enabling seamless data.

reuse in subsequent CIM operations. As illustrated in the figure, the proposed design reliably performs logic operations over a wide supply voltage range.

The optimal operating point is observed to be approximately 0.6 V. However, it is worth noting that to ensure sufficient current for the store and restore operations in the I-MTJ with minimum-sized transistors, a supply voltage of 0.7 V is selected. This voltage level provides the required current while maintaining compact transistor dimensions and remains close to the optimal point for efficient CIM computation.

Desings		14	I-MTJ ¹⁵	P-MTJ ¹⁵	13	Proposed
Store	Delay (ns)	6.3	4.0	2.8	2.6	3.4
	Power (μ W)	193.0	46.9	49.9	147.4	5.9
	PDP (pJ)	1.21	0.19	0.14	0.38	0.02
Restore	Delay (ps)	2008	1002	17.2	88.5	46
	Power (μ W)	556.8	65.5	406.3	2.8	91.6
	PDP (fJ)	1118.0	65.6	6.98	0.25	4.21

Table 4. The post-layout simulation results of the store and restore operations in the 4 Kb memory array.

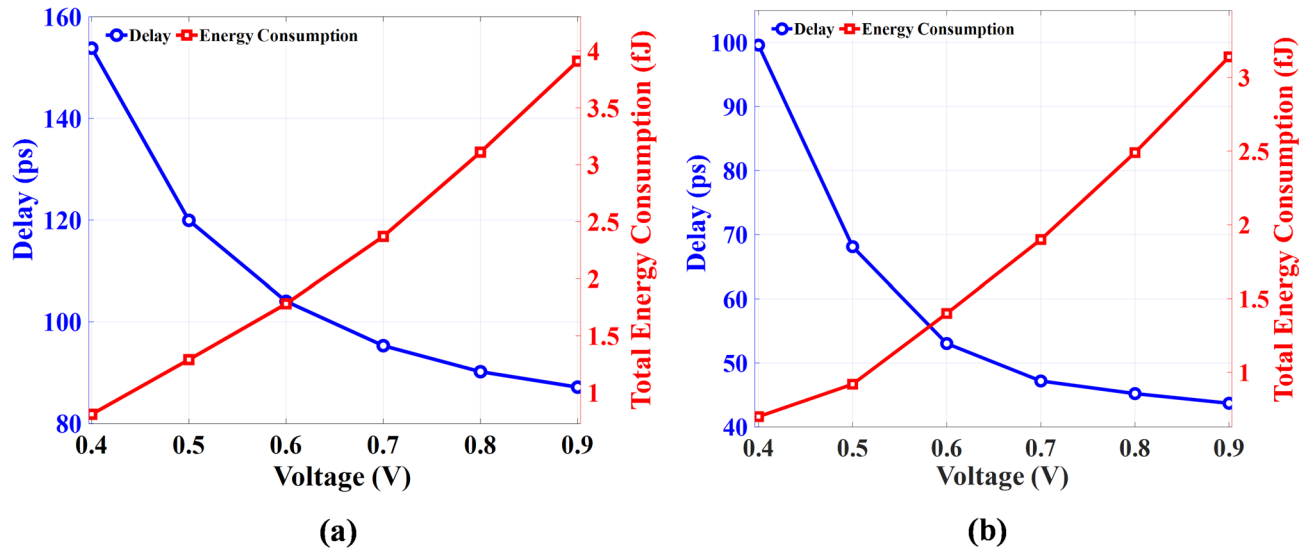


Fig. 10. CIM evaluation in different supply voltages (a) delay and energy of NOR/OR, (b) delay and energy of NAND/AND.

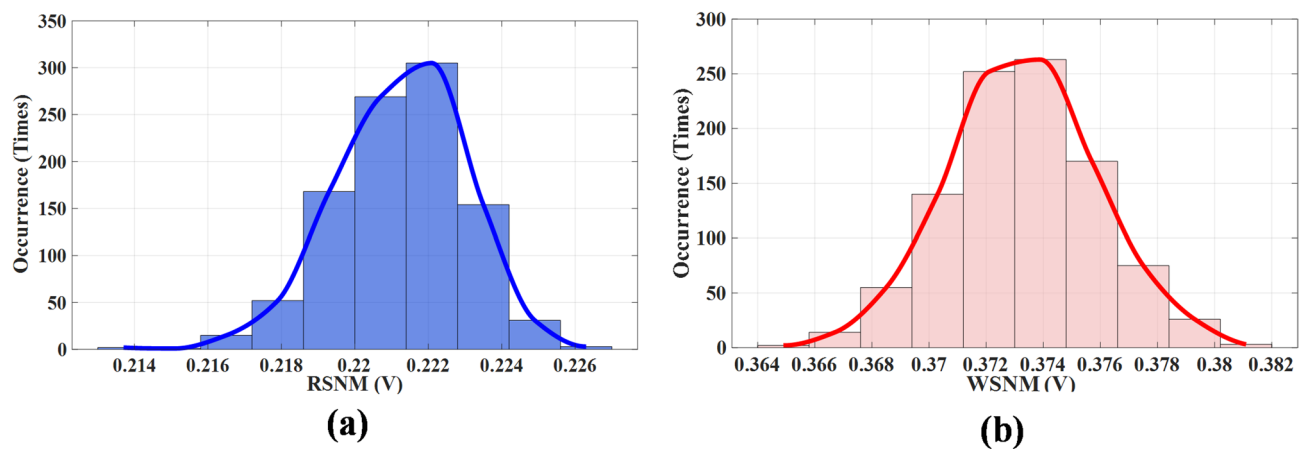


Fig. 11. Post-Layout yield results under SRAM mode operation for a 4 Kb memory array (a) RSNM (b) WSNM.

- (1) *Variation analysis* To assess the practical performance of the proposed design, 1000 Monte-Carlo simulations were conducted, targeting its functional behavior under variations in critical and sensitive parameters. In the context of FinFET technology, a 10% variation was applied to key process parameters, including fin height (H_{Fin}), silicon fin thickness (T_{Si}), gate length (L_g), and oxide thickness (T_{ox})⁴³. Furthermore, parameter variations of 10%, 5%, 15%, 5%, 15%, and 15% were applied to the TMR, barrier thickness (t_b), resistance-area product (RAP), free layer thickness (T_F), and the length and width of the free layer in the I-MTJs, respectively^{44,45}. Figure 11 presents the yield of the proposed design under SRAM mode operation (stability analysis) for a 4 Kb array configuration, confirming the practical robustness of the design. Moreover, the store and restore operations of the proposed design were thoroughly evaluated, and no failures were found during these processes. Figure 12 presents the transient waveforms of the proposed design during various memory operations, taking process variations into account.

Application

As previously discussed, a significant advantage of the proposed design is its ability to perform bidirectional CIM operations while concurrently storing the computation data within the same cycle. This integration enables efficient MVM execution without requiring additional memory management or internal data transfers. Moreover, because the cell can store intermediate data in its non-volatile component, the system can effectively resume operation after a power failure in edge devices running DNNs once power is restored and stabilized. To assess the effectiveness of the proposed design in real-world scenarios, AlexNet is implemented as a representative neural network⁴⁶.

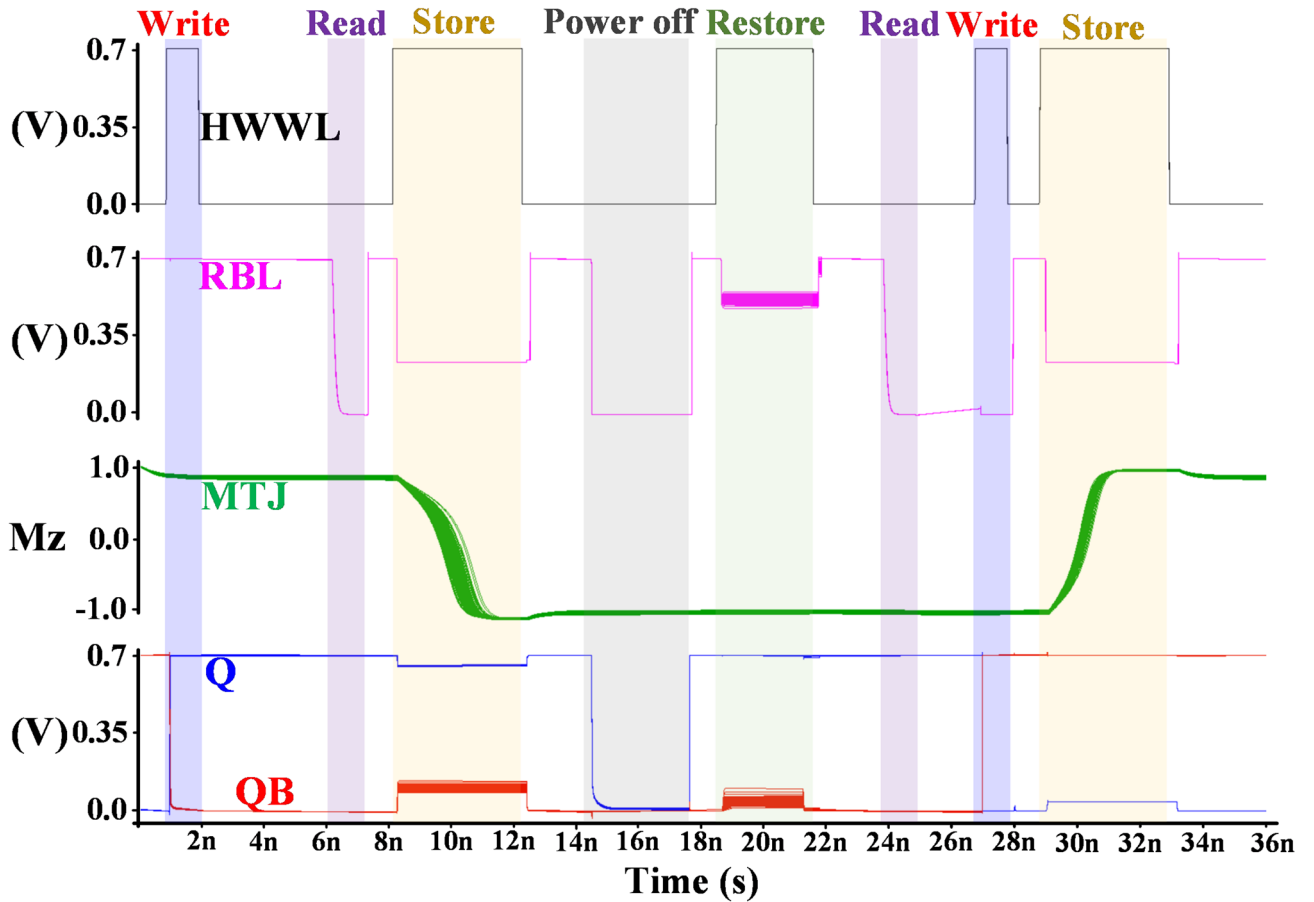


Fig. 12. Transient waveform of the proposed design in different operations, considering process variation.

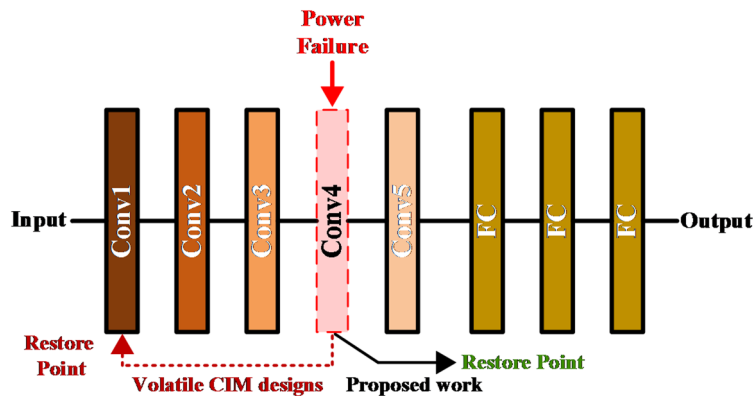


Fig. 13. Comparison of the backup scheme in the proposed design and volatile CIM under power failure scenarios in the implementation of AlexNet.

In this evaluation, we compare the proposed design with existing SRAM-based CIM architectures under a power-failure condition during network execution. To this end, the delay and power consumption associated with the full adder, multiplier, and read/write operations of both the proposed and baseline designs are simulated in Cadence, and post-layout results are extracted. These extracted circuit-level parameters are then incorporated into a Python-based implementation of the full AlexNet architecture to estimate the required MAC operations and overall system behavior. For instance, in the scenario where a failure occurs in Layer 4 of the network, as illustrated in Fig. 13, the SRAM-based CIM counterparts must reload all weights and restart computation from the first layer due to the loss of intermediate results. In contrast, the proposed design preserves critical and intermediate data in its non-volatile memory subsystem. Consequently, the system can resume operation

Designs	Delay (ms)	Power (mW)	EDP ($\mu\text{J s}$)
⁴⁷	100.7	1.17	117.82
⁴⁸	98.7	1.32	130.28
⁴⁹	100.4	1.31	131.52
Proposed	108.4	0.66	71.54

Table 5. Post-layout evaluation of the designs for implementing sAlexNet under power failure conditions on layer 4.

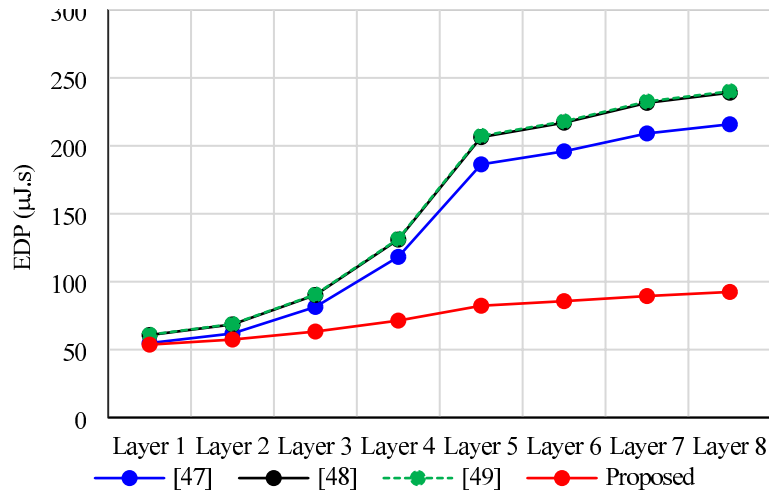


Fig. 14. EDP of the proposed design in different power failure layers of the AlexNet.

directly from the layer at which the power failure occurred, significantly improving energy efficiency and reducing recovery latency.

Table 5 summarizes the results of the AlexNet implementation with 8-bit network precision for the proposed design and its existing counterparts under a power-failure scenario in Layer 4. As shown, the proposed design outperforms the compared architecture by supporting bidirectional computation, enabling intermediate data writes within the same computation cycle, and retaining critical data during power failures. Together, these capabilities prevent unnecessary re-computation of previous layers and significantly improve the overall efficiency of the system. The proposed design achieves approximately 39.2% lower Energy Delay Product (EDP) consumption compared to the design presented in⁴⁷.

Moreover, we evaluate power failures occurring at different layers of the network to assess the performance of the proposed design across various scenarios, as illustrated in Fig. 14. As shown, when the failure occurs within the first two layers, the required re-computation of weights and intermediate data in the SRAM-based design^{47–49} is comparable to the power and delay overhead of the store-and-restore process in the proposed design. Consequently, the proposed architecture's performance improvement in these early layers is negligible. However, when the power failure occurs in deeper layers, the amount of re-computation required by the SRAM-based design increases substantially. In such cases, the proposed design clearly outperforms the compared designs, as it preserves and restores all critical data, allowing the network to resume from the failure point rather than recomputing from scratch.

Conclusion

This paper introduces a novel two-directional nvSRAM cell with integrated CIM functionality, tailored for edge devices with unstable power supplies. Unlike conventional CIM designs that require additional cycles for intermediate data storage, the proposed architecture enables direct storage of CIM outputs within the same computation cycle. Supporting both horizontal and vertical CIM operations facilitates efficient matrix-vector multiplication without data movement or extra storage overhead, addressing key limitations of one-directional CIM in DNN execution. The integration of an I-MTJ element ensures reliable data retention and restoration during power-off periods, significantly reducing power consumption. Detailed post-layout simulations demonstrate a 31% improvement in write margin, a 40% reduction in PDP in memory mode, and an 85% reduction in energy in backup scenarios compared to baseline designs. Notably, the proposed design achieves a 39.2% EDP reduction during neural network operation under power instability, underscoring its effectiveness for low-power edge computing.

Data availability

Data related to the current study are available from the corresponding author upon reasonable request.

Received: 19 October 2025; Accepted: 31 December 2025

Published online: 30 January 2026

References

- Sun, Y. et al. Energy-efficient nonvolatile SRAM design based on resistive switching multi-level cells. *IEEE Trans. Circuits Syst. II: Express Briefs* **66**, 753–757 (2019).
- Lin, Z. et al. High restore yield NVSRAM structures with dual complementary RRAM devices for high-speed applications. *IEEE Trans Very Large Scale Integr. VLSI Syst.* **31**, 522–531 (2023).
- Cheng, L. et al. Advancements in accelerating deep neural network inference on AIoT devices: a survey. *IEEE Trans. Sustain. Comput.* **9**, 830–847 (2024).
- Liu, J. et al. A SOT-MRAM-based CIM design with multi-bit resistance-sum paradigm and non-idealities tuning mechanism. *IEEE Trans. Magn.* **61**, 1–6 (2025).
- Jhang, C. J., Xue, C. X., Hung, J. M., Chang, F. C. & Chang, M. F. Challenges and trends of SRAM-based computing-in-memory for AI edge devices. *IEEE Trans. Circuits Syst. I: Regular Pap.* **68**, 1773–1786 (2021).
- Liu, J. et al. HiT-CIM: A high-throughput compute-in-memory SRAM architecture with simultaneous weight loading/computing and balance capabilities. *IEEE Trans. Emerg. Top. Comput.* **13**, 1396–1409 (2025).
- Lou, M. et al. Area-efficient and low-power 8T compute-SRAM bitcell design for digital compute-in-memory macros in 22nm CMOS. *IEEE Trans. Circuits Syst. II Express Briefs* **72**, 1605–1609 (2025).
- Yun, S.-J. et al. 5-nm high-efficiency and high-density digital SRAM in-memory-computing macros for AI accelerators. *IEEE Solid State Circuits Lett.* **8**, 269–272 (2025).
- Mu, C. et al. A 28-nm RRAM/SRAM collaborative CIM accelerator supporting RRAM-endurance-latency awareness for edge fine-tuning. *IEEE J Solid-State Circuits* 1–13 (2025) <https://doi.org/10.1109/JSSC.2025.3577335>.
- Han, S. & Jiang, Y. RISC-V-based evaluation and strategy exploration of MRAM triple-level hybrid cache systems. *IEEE Trans. Very Large Scale Integr. VLSI Syst.* **31**, 980–992 (2023).
- Choi, S., Son, J., Cho, K. & Kim, S. One-transistor static random-access memory cell array comprising single-gated feedback field-effect transistors. *Sci Rep* **11**, 17983 (2021).
- Jia, H. et al. Scalable and programmable neural network inference accelerator based on in-memory computing. *IEEE J. Solid-State Circuits* **57**, 198–211 (2022).
- Wu, T. et al. A Heterogeneous Microprocessor for Intermittent AI Inference using Nonvolatile-SRAM-based Compute-In-Memory. *IEEE Trans. Circuits Syst. II Express Briefs* <https://doi.org/10.1109/TCSII.2023.3289493> (2023).
- Tripathi, S., Choudhary, S. & Misra, P. K. A Novel STT-SOT MTJ-based nonvolatile SRAM for power gating applications. *IEEE Trans. Electron. Devices* **69**, 1058–1064 (2022).
- Kang, W., Lv, W., Zhang, Y. & Zhao, W. Low store power high-speed high-density nonvolatile SRAM design with spin hall effect-driven magnetic tunnel junctions. *IEEE Trans. Nanotechnol.* **16**, 148–154 (2017).
- Tripathi, S., Choudhary, S. & Misra, P. K. Highly reliable, stable, and store energy efficient 8T/9T-2D-2MTJ NVSRAMs. *IEEE Trans. Nanotechnol.* **23**, 89–94 (2024).
- Wang, X. et al. 29.2 A 0.021 μm^2 High-Density SRAM in Intel-18A-RibbonFET Technology with PowerVia-Backside Power Delivery. in *2025 IEEE International Solid-State Circuits Conference (ISSCC)* 494–496 (IEEE, 2025). <https://doi.org/10.1109/ISSCC49661.2025.10904657>.
- Fazili, M. M., Ahmad, S. & Iqbal, B. A half-Schmitt trigger-based 9T1R nonvolatile robust SRAM cell for instant on-off application. *IEEE Trans. Circuits Syst. I Regul. Pap.* **72**, 169–179 (2025).
- Lu, W. et al. High-reliability and high-throughput CIM 10T-SRAM for multiplication and accumulation operations with 274.3 GOPS and 200–237.5 TOPS/W. *IEEE Trans. Very Large Scale Integr. VLSI Syst.* **33**, 1073–1081 (2025).
- Wei, F., Cui, X., Zhang, S. & Zhang, X. An 11T SRAM Cell for the dual-direction in-array logic/CAM operations. *IEEE Trans. Circuits Syst. II: Express Briefs* **71**, 2329–2333 (2024).
- Joshi, V. K. Spintronics: A contemporary review of emerging electronics devices. *Eng. Sci. Technol. Int. J.* **19**, 1503–1513. <https://doi.org/10.1016/j.jestch.2016.05.002> (2016).
- Li, H., Yao, E., Qin, P. & Jiang, S. USCIM: Computing-in-memory with unipolar switching SOT-MRAM. *IEEE Trans. Magn.* **61**, 1–6 (2025).
- Mohseni, A., Moaiyeri, M. H. & Adel, M. J. A novel obfuscation method based on majority logic for preventing unauthorized access to binary deep neural networks. *Sci. Rep.* **15**, 24416 (2025).
- Wang, C. et al. Technically feasible robust complementary SOT-MRAM design for improving the area and energy efficiency. *IEEE Trans. Circuits Syst. I Regul. Pap.* **72**, 2327–2340 (2025).
- Ashtari Gargari, M., Eslami, N. & Moaiyeri, M. H. A reconfigurable nonvolatile memory architecture for prolonged wearable health monitoring devices. *IEEE Trans. Consum. Electron.* **70**, 4717–4728 (2024).
- Shreya, S. & Kaushik, B. K. Differential spin hall effect-based nonvolatile static random access memory for energy-efficient and fast data restoration application. *IEEE Trans. Magn.* **55**, 1–11 (2019).
- Scheike, T., Wen, Z., Sukegawa, H. & Mitani, S. 631% room temperature tunnel magnetoresistance with large oscillation effect in CoFe/MgO/CoFe(001) junctions. *Appl. Phys. Lett.* **122**, 112404 (2023).
- Ikeda, S. et al. Tunnel magnetoresistance of 604% at 300 K by suppression of Ta diffusion in CoFeBMgOCoFeB pseudo-spin-valves annealed at high temperature. *Appl. Phys. Lett.* **93**, 082508 (2008).
- He, Q. et al. High-dose X-ray radiation induced MgO degradation and breakdown in spin transfer torque magnetic tunnel junctions. *Sci. Rep.* **12**, 18620 (2022).
- Razi, F., Moaiyeri, M. H., Rajaei, R. & Mohammadi, S. A variation-aware ternary spin-hall assisted STT-RAM based on hybrid MTJ/GAA-CNTFET logic. *IEEE Trans. Nanotechnol.* **18**, 598–605 (2019).
- Tserkovnyak, Y., Brataas, A. & Bauer, G. E. W. Enhanced Gilbert damping in thin ferromagnetic films. *Phys. Rev. Lett.* **88**, 117601 (2002).
- Sivabalan, K. et al. Evaluation of effective thermal stability factor for patterned magnetic tunnel junction array. *IEEE Trans. Magn.* **57**, 1–5 (2021).
- Lee, K. et al. 1Gbit High Density Embedded STT-MRAM in 28nm FDSOI Technology. in *Technical Digest - International Electron Devices Meeting, IEDM* vols 2019-December (2019).
- Naik, V. B. et al. Manufacturable 22nm FD-SOI Embedded MRAM Technology for Industrial-grade MCU and IOT Applications. in *Technical Digest - International Electron Devices Meeting, IEDM* vols 2019-December (2019).
- Wei, L. et al. 13.3 A 7Mb STT-MRAM in 22FFL FinFET Technology with 4ns Read Sensing Time at 0.9V Using Write-Verify-Write Scheme and Offset-Cancellation Sensing Technique. in *Digest of Technical Papers - IEEE International Solid-State Circuits Conference* vols 2019-February (2019).

36. Sang, H., Xie, W., Park, G. & Yoo, H. J. An 2.31uJ/Inference ultra-low power always-on event-driven AI-IoT SoC with switchable nvSRAM compute-in-memory macro. *IEEE Trans. Circuits Syst. II: Express Briefs* **71**, 2534–2538 (2024).
37. Bagheriye, L., Toofan, S., Saeidi, R., Zeinali, B. & Moradi, F. A reduced store/restore energy MRAM-based SRAM cell for a non-volatile dynamically reconfigurable FPGA. *IEEE Trans. Circuits Syst. II: Express Briefs* **65**, 1708–1712 (2018).
38. Bahador, A., Moaiyeri, M. H. & Ghaderi, R. Algorithmically-enhanced design of spintronic-based tunable true random number generator for dependable stochastic computing. *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems* **1–1** (2024) <https://doi.org/10.1109/TCAD.2024.3446715>.
39. Clark, L. T. et al. ASAP7: A 7-nm finFET predictive process design kit. *Microelectron. J.* **53**, 105–115 (2016).
40. Kang, W. et al. Spintronic logic design methodology based on spin Hall effect-driven magnetic tunnel junctions. *J. Phys. D Appl. Phys.* **49**, 065008 (2016).
41. Wang, Y. et al. Compact model of dielectric breakdown in spin-transfer torque magnetic tunnel junction. *IEEE Trans. Electron. Devices* **63**, 1762–1767 (2016).
42. Huang, Y. et al. Top spin-orbit-torque switching of magnetic tunnel junction with in-situ efficiency quantification. *IEEE Electron. Device Lett.* **46**, 1409–1412 (2025).
43. Eslami, N., Moaiyeri, M. H. & Yousefi, A. Hybrid ternary RRAM-based in-memory computing architecture for energy-efficient data-driven applications. *Arab. J. Sci. Eng.* <https://doi.org/10.1007/s13369-025-10164-3> (2025).
44. Wang, Y. et al. Compact model of magnetic tunnel junction with stochastic spin transfer torque switching for reliability analyses. *Microelectron. Reliab.* **54**, 1774–1778 (2014).
45. Fathollahi, A., Amirany, A. & Moaiyeri, M. H. Ternary computing using a novel spintronic multi-operator logic-in-memory architecture. *Res. Eng.* **25**, 104011 (2025).
46. Lammie, C. et al. LIONHEART: A Layer-based mapping framework for heterogeneous systems with analog in-memory computing tiles. *IEEE Trans. Emerg. Top Comput.* 1–13. <https://doi.org/10.1109/TETC.2025.3546128> (2025).
47. Chen, J. et al. A reliable 8T SRAM for high-speed searching and logic-in-memory operations. *IEEE Trans. Very Large Scale Integr. VLSI Syst.* **30**, 769–780 (2022).
48. Lin, Z. et al. In-memory computing with double word lines and three read ports for four operands. *IEEE Trans. Very Large Scale Integr. VLSI Syst.* **28**, 1316–1320 (2020).
49. Agrawal, A., Jaiswal, A., Lee, C. & Roy, K. X-SRAM: Enabling in-memory boolean computations in CMOS static random access memories. *IEEE Trans. Circuits Syst. I Regul. Pap.* **65**, 4219–4232 (2018).

Author contributions

S.P.H. proposed the main idea, performed the experiments, and wrote the main paper. M.H.M. helped develop the concept, reviewed the simulations, discussed and analyzed the results, edited the paper, supervised the study, and proofread the paper. N.E. helped develop the idea, conducted some simulations, discussed the results, and proofed the paper. B.E. discussed the results and proofed the paper.

Declarations

Competing interests

The authors declare no competing interests.

Additional information

Correspondence and requests for materials should be addressed to M.H.M.

Reprints and permissions information is available at www.nature.com/reprints.

Publisher's note Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

Open Access This article is licensed under a Creative Commons Attribution-NonCommercial-NoDerivatives 4.0 International License, which permits any non-commercial use, sharing, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons licence, and indicate if you modified the licensed material. You do not have permission under this licence to share adapted material derived from this article or parts of it. The images or other third party material in this article are included in the article's Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons licence and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this licence, visit <http://creativecommons.org/licenses/by-nc-nd/4.0/>.

© The Author(s) 2026