

Experimental demonstration of magnetic tunnel junction-based computational random-access memory

Yang Lv¹, Brandon R. Zink¹, Robert P. Bloom¹, Hüsrev Cilasun¹, Pravin Khanal², Salonik Resch¹, Zamshed Chowdhury¹, Ali Habiboglu², Weigang Wang², Sachin S. Sapatnekar¹, Ulya Karpuzcu¹, & Jian-Ping Wang^{1*}

¹Department of Electrical and Computer Engineering, University of Minnesota, Minneapolis, Minnesota 55455, USA

²Department of Physics, University of Arizona, Tucson, Arizona 85721, USA

*e-mail: jpwang@umn.edu;

Conventional computing paradigm struggles to fulfill the rapidly growing demands from emerging applications, especially those for machine intelligence, because much of the power and energy is consumed by constant data transfers between logic and memory modules. A new paradigm, called “computational random-access memory (CRAM)” has emerged to address this fundamental limitation. CRAM performs logic operations directly using the memory cells themselves, without having the data ever leave the memory. The energy and performance benefits of CRAM for both conventional and emerging applications have been well established by prior numerical studies. However, there lacks an experimental demonstration and study of CRAM to evaluate its computation accuracy, which is a realistic and application-critical metrics for its technological feasibility and competitiveness. In this work, a CRAM array based on magnetic tunnel junctions (MTJs) is experimentally demonstrated. First, basic memory operations as well as 2-, 3-, and 5-input logic operations are studied. Then, a 1-bit full adder with two different designs is demonstrated. Based on the experimental results, a suite of modeling has been developed to characterize the accuracy of CRAM computation. Further analysis of scalar addition, multiplication, and matrix multiplication shows promising results. These results are then applied to a complete application: a neural network based handwritten digit classifier, as an example to show the connection between the application performance and further MTJ development. The classifier achieved almost-perfect classification accuracy, with reasonable projections of future MTJ development. With the confirmation of MTJ-based CRAM’s accuracy, there is a strong case that this technology will have a significant impact on power- and energy-demanding applications of machine intelligence.

Introduction

Recent advances in machine intelligence^{1,2} for tasks such as recommender systems³, speech recognition⁴, natural language processing⁵, and computer vision⁶, have been placing growing demands on our computing systems, especially for implementations with artificial neural networks. A variety of platforms are used, from general-purpose CPUs and GPUs^{7,8}, to FPGAs⁹, to custom-designed accelerators and processors¹⁰⁻¹³, to mixed- or fully- analog circuits¹⁴⁻²⁰. Most are based on the Von Neumann architecture, with separate logic and memory systems. As shown in Fig. 1a, the inherent segregation of logic and memory requires large amounts of data to be transferred between these modules. In data-intensive scenarios, this transfer becomes a major bottleneck in terms of performance, energy consumption, and cost²¹⁻²³. For example, the data movement consumes about 200 times of the energy used for computation when reading three 64-bit source operands from and writing one 64-bit destination operand to an off-chip main memory²¹. This bottleneck has long been studied. Research aiming at connecting logic and memory more closely has led to new computation paradigms.

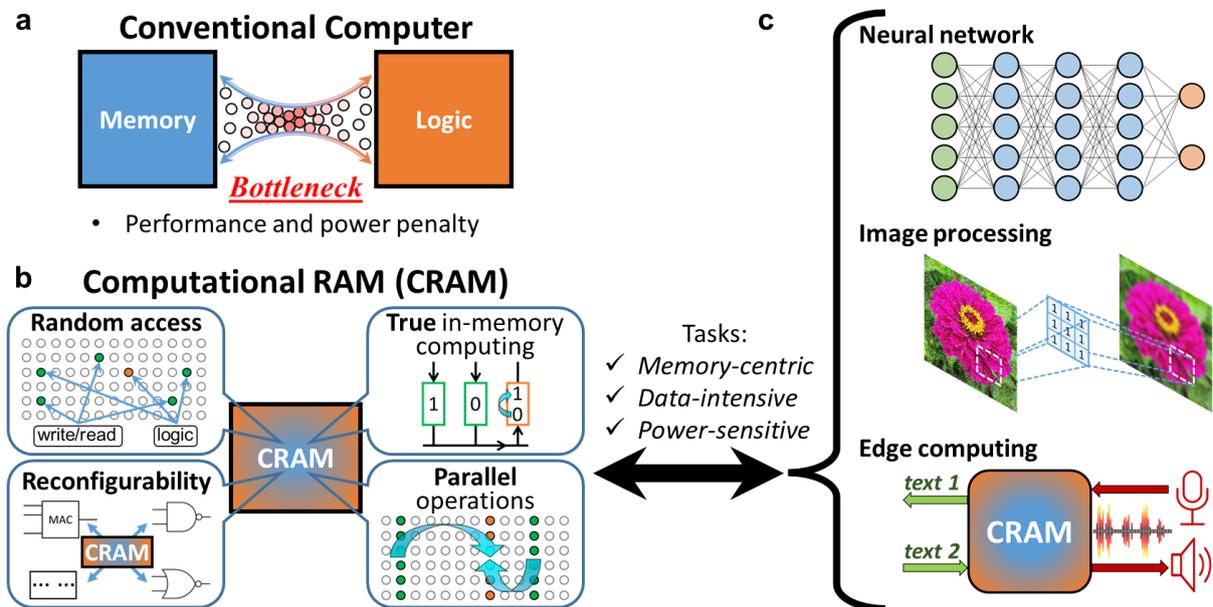


Fig. 1 | Illustrations of CRAM concept, features, and potential applications. **a, b,** Compared to a conventional computer architecture (**a**), which suffers from the memory-logic transfer bottleneck, CRAM (**b**) offers significant power and performance improvements. Its unique architecture allows for computation in memory, as well as, random access, reconfigurability, and parallel operation capability. **c,** The CRAM could excel in (**c**), data-intensive, memory-centric, or power-sensitive applications, such as neural networks, image processing, or edge computing.

Promising paradigms include “near-memory” and “in-memory” computing. Near-memory processing brings logic physically closer to memory by employing 3D-stacked architectures²⁴⁻²⁹. In-memory computing scatters clusters of logic throughout or around the memory banks on a single chip^{14-20,30-35}. Yet another approach is to build systems where the memory itself can perform computation. This has been dubbed “true” in-memory computing³⁶⁻⁴¹. The computational random-access memory (CRAM)^{36,37} is one of the true in-memory computing paradigms. Logic is performed natively by the memory cells; the data for logic operations never has to leave the memory (Fig. 1b). Additionally, CRAM operates in a fully digital fashion, unlike most other reported in-memory computing schemes¹⁴⁻²⁰, which are partially or mostly analog. CRAM promises superior energy efficiency and processing performance for machine intelligence applications. It has unique additional features, such as

random-access of data and operands, massive parallel computing capabilities, and reconfigurability of operations^{36,37}.

The CRAM was initially proposed based on the MTJ device³⁶, an emerging memory device that relies on spin electronics⁴². Such “spintronic” devices, along with other non-volatile emerging memory devices, usually referred as “X” for logic applications, have been intensively investigated over the past several decades for emerging memory and computing applications as “beyond-CMOS” and/or “CMOS+X” technologies. They offer vastly improved speed, energy efficiency, area, and cost. An additional feature that is exploited by CRAM is their non-volatility⁴³. The MTJ device is the most mature of spintronic devices for embedded memory applications, based on endurance⁴⁴, energy efficiency⁴⁵, and speed⁴⁶. We note that CRAM can be implemented based on not only spintronics devices, but also other non-volatile emerging memory devices.

In its simplest form, an MTJ consists of a thin tunneling barrier layer sandwiched by two ferromagnetic (FM) layers. When a voltage is applied between the two layers, electrons tunnel through the barrier resulting in a charge current. The resistance of the MTJ is a function of the magnetic state of the two FM layers, due to the tunneling magnetoresistance (TMR) effect⁴⁷⁻⁴⁹. An MTJ can be engineered to be magnetically bi-stable. Accordingly, it can store information based on its magnetic state. This information can be retrieved by reading the resistance of the device. The MTJ can be electrically switched from one state to the other with a current, due to the spin-transfer torque (STT) effect^{50,51}. In this way, an MTJ can be used as an electrically operated memory device with both read and write functionality. A type of random-access memory, the STT-MRAM⁵²⁻⁵⁵ has been developed commercially, utilizing MTJs as memory cells. A typical STT-MRAM consists of an array of bit cells, each containing one transistor and one MTJ. These are referred to as 1 transistor 1 MTJ (1T1M) cells.

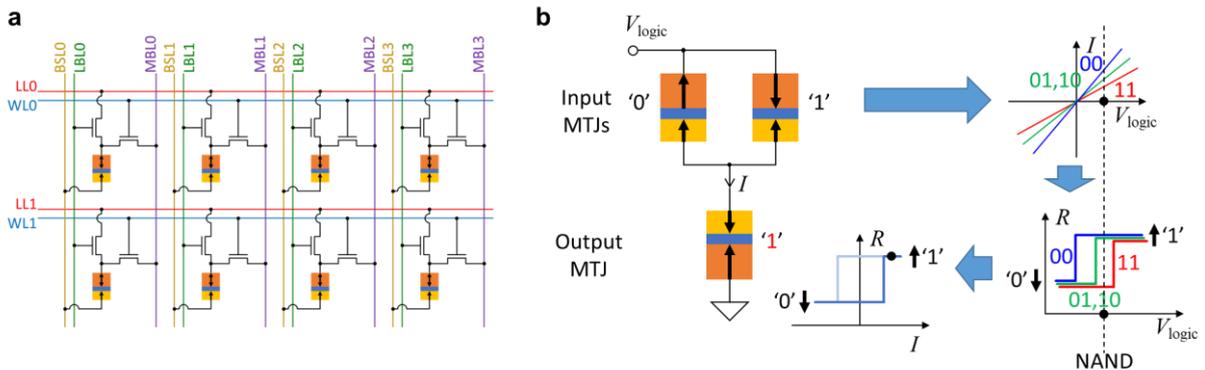


Fig. 2 | Illustrations of the CRAM cell architecture and the working principle of CRAM logic operation. **a**, CRAM adopts the so-called 2 transistor 1 MTJ (2T1M) cell architecture. On top of the 1T1M cell architecture of STT-MRAM, additional transistor, as well as the added logic line (LL), and logic bit line (LBL), allow the CRAM to perform logic operations. During a CRAM logic operation, the transistors and lines are manipulated to form equivalent circuit as shown in **(b)**. Although CRAM can be built based on various emerging memory devices, we use MTJs and MTJ-based CRAM as an example for illustration purposes. **b**, The working principle of CRAM logic operation, the VCL, utilizes the thresholding effect that occurs when switching an MTJ and the TMR effect of MTJ. With an appropriate V_{logic} amplitude, the voltage is translated into different current flowing through the output MTJ by the TMR effect of the input MTJs. Whether the output MTJ switches or not is dependent on the state of input MTJs.

A typical CRAM cell design, as shown in Fig. 2a, is a modification of the 1T1M STT-MRAM architecture⁵⁶. The MTJ, one of the transistors, word line (WL), bit select line (BSL),

and memory bit line (MBL), resemble the 1T1M cell architecture of STT-MRAM, which allow the CRAM to perform memory operations. In order to enable logic operation, a second transistor, as well as a logic line (LL), and a logic bit line (LBL), are added to each memory cell. During a logic operation, corresponding transistors and lines are manipulated so that several MTJs in a row are temporarily connected to a shared LL³⁷. While the LL is left floating, voltage pulses are applied to the lines connecting to input MTJs with that of the output MTJ being grounded. The logic operation is based on a working principle called voltage-controlled logic (VCL)^{57,58}, which utilizes the thresholding effect that occurs when switching an MTJ and the TMR effect of MTJ. As shown in Fig. 2b, when a voltage is applied across the input MTJs, the different resistance values result in different current levels. The current flows through the output MTJ, which may or may not switch its state, depending on the states of the input MTJs. In this way, basic bitwise logic operations, such as AND, OR, NAND, NOR, and MAJ, can be realized. A unique feature of VCL is that the logic operation itself does not require the data in the input MTJs to be read out through sense amplifiers at the edge of the array. Rather, it is used locally within the set of MTJs involved in the computation. This is fundamentally why the CRAM computation represents true-in-memory computing: the computation does not require data to travel out of the memory array. It is always processed locally by nearby cells. We note that this concept would also work with other two-terminal stateful passive memory devices, such as memristors. Accordingly, a CRAM could be implemented with such devices. A CRAM could also be implemented with three-terminal stateful devices, such as spin-orbit torque (SOT) devices. This could result in greater energy efficiency and reliability⁵⁹. As an oversimplified speculation, the performance comparison between CRAMs implemented by various emerging memory devices is expected to roughly follow the comparison between these for memory applications, since CRAM utilizes memory devices in similar manners like in memory application. For example, a CRAM implemented based on MTJs should be expected to offer high endurance and high speed. Also, generally, a CRAM logic operation should consume energy comparable to the energy consumption of a memory write operation, for the same emerging memory device operating at the same speed. However, a careful case-by-case analysis is necessary for CRAMs implemented by each emerging memory device technology. Also note that we do not show a specific circuit design of CRAM peripherals because CRAM does not require significant circuit design change in sensing amplifiers or peripherals compared to 1T1M STT-MRAM. And these in the STT-MRAM are already common and mature. Lastly, the true-in-memory computing characteristic of CRAM is limited to within a continuous CRAM array: any computation that requires access to data across separate CRAM arrays will require additional data access and movement. However, this limitation is true for all other in-memory computing paradigms. CRAM is not at any disadvantage in this scenario.

On top of the potential performance benefits that the emerging memory devices bring, at circuit level, CRAM fundamentally provides several benefits (Fig. 1b): (1) the elimination of the costly performance and energy penalties associated with transferring data between logic and memory; (2) random access of data for the inputs and outputs to operations; (3) the reconfigurability of operations, as any of the logic operation AND, OR, NAND, NOR, and MAJ can be programmed; and (4) the performance gain of massive parallelism, as operations can be performed in parallel in each row of the CRAM array. Based on analysis and benchmarking, CRAM has the potential to deliver significant gains in performance and power efficiency, particularly for data-intensive, memory-centric, or power sensitive applications, such as bioinformatics^{37,60,61}, image⁶² and signal⁶³ processing, neural networks^{62,64}, and edge computing⁶⁵ (Fig. 1c). For example, a CRAM-based machine-learning inference accelerator was estimated to achieve an improvement on the order of 1000× over a state-of-art solution,

in terms of the energy-delay product⁶⁶. And yet, to date, there have been no experimental studies of CRAM.

In this work, we present the first experimental demonstration of a CRAM array. Although based on a small 1×7 array, it successfully shows complete CRAM array operations. We illustrate computation with a 1-bit full adder. This work provides a proof-of-concept as well as a platform with which to study key aspects of the technology experimentally. We provide detailed projections and guidelines for future CRAM design and development. Specifically, based on the experiment results, models and calculations of CRAM logic operations are developed and verified. The results connect the CRAM gate-level accuracy or error rate to MTJ TMR ratio, logic operation pulse width, and other parameters. Then we evaluate the accuracy of a multi-bit adder, a multiplier, and a matrix multiplication unit, which are fundamental building blocks for artificial neural networks. We further evaluate a CRAM-based implementation of the MNIST⁶⁷ handwritten digit classifier. It is used as an example to show the connection between the application performances and further improvements of MTJs. These results link the application-level accuracy to the gate-level error rate of CRAM. They confirm the potential technological relevance and competitiveness of CRAM for applications in conventional domains as well as emerging applications related to machine intelligence.

Experiments

Figure 3 shows the experimental setup, consisting of both hardware and software. The hardware is built with a so-called ‘circuit-around-die’ approach⁶⁸: semiconductor circuitry is built with commercially-available components around the MTJ dies. This approach offers a more rapid development cycle and flexibility needed for exploratory experimental study on CRAM arrays and potential new MTJ technologies, while the major foundries lack the specific process design kit available for making a CRAM array fully integrated with CMOS. The hardware is a 1×7 CRAM array, with the design of cells taken from the 2T1M CRAM cells^{36,37}, modified for simplified memory access. Software on a PC controls the operation. It communicates with the hardware with basic commands: ‘open/close transistors’; ‘apply voltage pulses’ to perform write and logic operations’; and ‘read MTJ cell resistance’. The software collects real-time measurements of the data associated with CRAM operations for analysis and for visualization. All aspects of the 1×7 CRAM array are functional: memory write, memory read, and logic operations.

MTJs with perpendicular interfacial anisotropy are used in the CRAM. They exhibit low resistance-area (RA) product and high TMR ratio – approximately 100% – when sized at 100 nm in diameter.

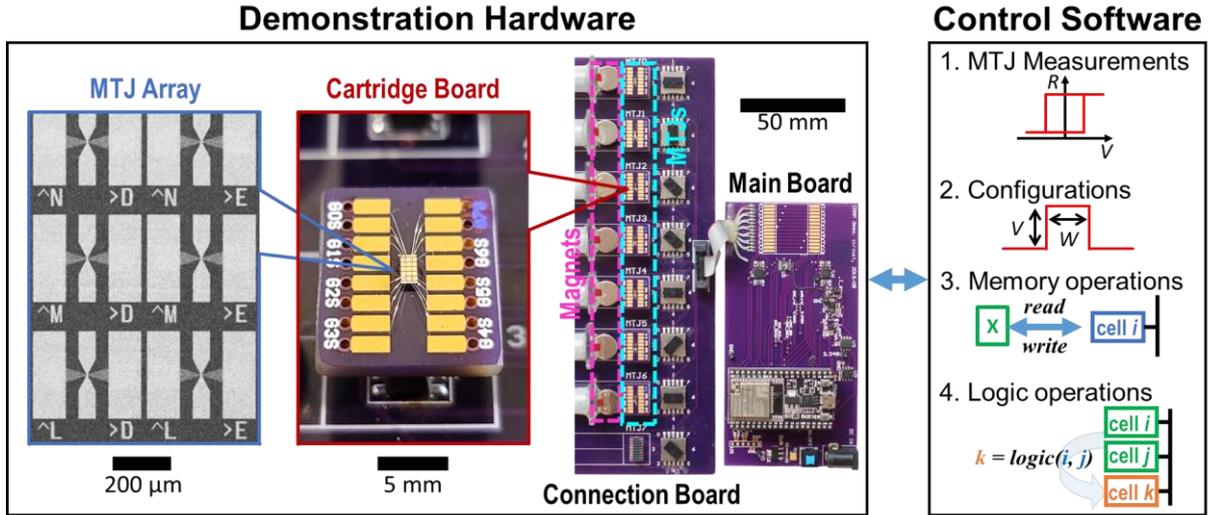


Fig. 3 | CRAM experimental setup. The setup consists of custom-built hardware and a suite of control software. It demonstrates a fully functioning 1×7 CRAM array. The hardware consists of a main board hosting all necessary electronics except for the MTJ devices; a connection board on which passive switches, connectors, and magnetic bias field mechanisms are hosted; and multiple cartridge boards that each have a MTJ array mounted and multiple MTJ devices that are wire bonded. The gray-scale scanning electron microscopy image shows the MTJ array used. The color optical photographs show the cartridge board and the entire hardware setup. The software is responsible for real-time measurements of the MTJs; configuration and execution of CRAM operations: memory write, memory read, and logic; and data collection. It is run on a PC, which communicates wirelessly with the main board.

Results

The experiments begin with measuring the resistance (R)–voltage (V) properties of each MTJ device and of each die. In order to compensate for device-to-device variations, the bias magnetic fields for each MTJ are adjusted so that the R – V properties are as close to each other as possible. As the processes of making CRAM array mature, bias magnetic fields are expected to be no longer needed and all CRAM cells will be able to be operated with uniform parameters and under uniform conditions. The resistance threshold for the MTJs logic states is also determined in this stage.

Then the seven MTJ cells are tested for memory operations with various write pulse amplitudes and widths. Based on the observed error rates for memory write operations, appropriate pulse amplitudes and widths are configured, achieving reliable memory write operations with an average error rate of less than 1.5×10^{-4} . We designate logic ‘0’ and ‘1’ to the parallel (P) low resistance state and anti-parallel (AP) high resistance state of MTJ, respectively.

Two-input logic operations are studied. The output cell is first initialized by writing ‘0’ to it. Then two input cells are connected to the output cell through the LL by turning on corresponding transistors. Voltage pulses of amplitude of V_{logic} , V_{logic} , 0, are simultaneously applied to the two-input cells and the output cell, respectively. This is the same as grounding the output cell while applying a voltage pulse of V_{logic} to the two input cells. Then depending on the input cells states, the output cell will have a certain probability of being switched from ‘0’ to ‘1’. Such a cycle of operations is repeated n times, and the statistical mean of the output logic state, $\langle D_{\text{out}} \rangle$, is obtained. The entire process is repeated for different V_{logic} values and input states. The basis for logic operations in the CRAM is the state-dependent resistance of the input cells. These shift and displace the output cell’s switching probability transfer curve.

As a result, the output cell switches state based on specific input states, therefore, implementing a logic function such as AND, OR, NAND, NOR, or MAJ. A specific initial state of the output cell and V_{logic} value corresponds to one of these logic gates⁶². The time duration or pulse width of the voltage pulse applied during a logic operation is expected to contribute to most of the time required to complete a logic operation. In the following we use the term logic speed to generally refer to the speed of a logic operation. Logic speed is approximately inversely proportional to the time duration of the voltage pulse used during a logic operation.

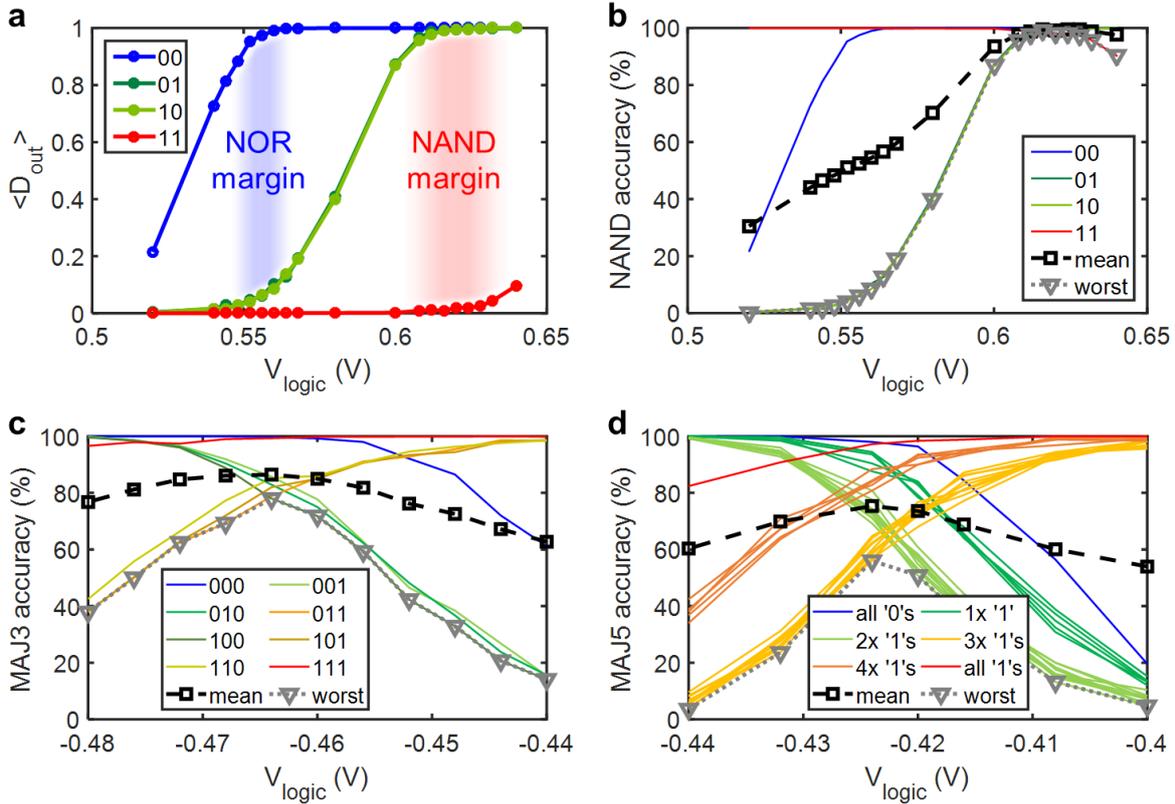


Fig. 4 | Experimental results for CRAM logic operation. **a**, Output logic average, $\langle D_{\text{out}} \rangle$, vs. logic voltage, V_{logic} . In a 2-input logic operation, two input cells and one output MTJ cell are involved. The output cell's terminal is grounded, while the common line is left floating. A logic voltage pulse is applied on the two input cells' terminals for a fixed duration (pulse width) of 1 ms. Before each logic operation, input data is written to the input cells. After each logic operation, the output cell's state is read. Each curve corresponds to a specific input state. Each data point represents the statistical average of the output cell's logic state, $\langle D_{\text{out}} \rangle$, sampled by 1000 repeats ($n = 1000$) of the operations. The separation between the $\langle D_{\text{out}} \rangle$ curves indicates the margins for NOR or NAND operation, highlighted in blue and red, respectively. **b**, Accuracy of 2-input NAND operation vs. logic voltage, V_{logic} . The results in (a) can be converted into a more straightforward metric, accuracy, for the NAND truth table. The curve labeled 'mean', and 'worst' indicate the average and the worst-case accuracy across all input states, respectively. So, for NAND operation, the optimal logic voltage is indicated in such a plot where the mean or worst accuracy is maximized. **c**, **d**, Accuracy of MAJ3 (c) and MAJ5 (d) logic operation vs. logic voltage, V_{logic} . Each curve corresponds to an input state or a group of input states. And each data point represents the statistical average of the output MTJ logic state sampled by $n = 1000$ and $n = 250$, for (c) and (d), respectively.

The experimental results are shown in Figs. 4a and 4b. Generally, for a given input state, $\langle D_{\text{out}} \rangle$ increases with increasing V_{logic} . The $\langle D_{\text{out}} \rangle$ response curves are input-state dependent. The four input states can be divided into three groups:

- The ‘00’ input state yields the lowest resistance at the two input cells, so the output cell switches from ‘0’ to ‘1’ first (with the lowest V_{logic}).
- The ‘11’ input state yields the highest resistance at the two input cells, so the output cell switches from ‘0’ to ‘1’ last (with the highest V_{logic}).
- The ‘01’ and ‘10’ input states both yield resistance that falls in between that of ‘00’ and ‘11’. Accordingly, the output cell’s response curve falls in between that of ‘00’ and ‘11’.

Figure 4a shows the experiment results. The two regions highlighted in blue and red that fall in between the three groups of response curves are suitable for NOR and NAND operations, respectively. For example, in the red region, the ‘11’ input has a high probability of yielding a ‘0’ output, while the other three input states have a high probability of yielding a ‘1’ output. This matches the expected truth table for a NAND logic gate. Therefore, if V_{logic} is chosen carefully – within the red region for the CRAM 2-input logic operation – the operation performed is highly likely to be NAND.

The experimental results of $\langle D_{\text{out}} \rangle$ can be converted into a straightforward format representing the accuracy for specified logic function. This translation can be computed by simply subtracting $\langle D_{\text{out}} \rangle$ from 1 for those input states where a ‘0’ output is expected in the truth table of the logic function. Figure 4b shows NAND accuracy of the same 2-input CRAM logic operation. The ‘mean’ and ‘worst’ plots are based on the average value and minimum value of the accuracy, respectively, across all input state combinations at a fixed value for V_{logic} . Based on the experimental results, if $V_{\text{logic}} = 0.624$ or 0.616 V, the CRAM delivers a NAND operation with a best mean and a worst-case accuracy of about 99.4% and 99.0%, respectively. From a circuit perspective, both increasing the effective TMR ratio of input cells and/or making the output cell’s response curve steeper would increase the vertical separation of these input-state-dependent curves, resulting in higher accuracy. For example, higher effective TMR ratio of input cells results in larger contrast of current in the output cell between different input states. Therefore, there is more ‘horizontal’ room to separate the $\langle D_{\text{out}} \rangle$ curves associated with different input states so that for the inputs with which the output is expected to be ‘0’ or ‘1’, the $\langle D_{\text{out}} \rangle$ of output cell is closer to the expected value (‘0’ or ‘1’). Also note that for a logic operation, the ‘accuracy’ and ‘error rate’ are essential two quantities describing the same thing: how true is the logic operation is, statistically. By definition, the sum of accuracy and error rate is always 1. The higher or closer to 1 the accuracy is, the better. The lower or closer to 0 the error rate is, the better. Lastly, to facilitate better visualization of how the resistance changes of different input cell states are translated into voltage differences on the output cell resulting in it being switched or unswitched, we list the equivalent resistance of the two input cells combined in parallel and the resulted voltage on the output cell in the following. With $V_{\text{logic}} = 0.620$ V, the equivalent resistance of input cells and the resulted voltage on the output cell are 0.4133 V and 1120 Ω , 0.3753 V and 1461 Ω , and 0.3248 V and 2037 Ω , for input states ‘00’, ‘01’ or ‘10’, and ‘11’, respectively. Note that these values are estimated by the experiment-based modeling, which is introduced in the later part of this paper.

With more input cells, we also studied 3-input and 5-input majority logic operations. As the number of input cells involved in a CRAM logic operation increases, two mechanisms impact the accuracy. The first is that there are more input states that must be distinguished in the limited range of the relative resistance changes. This makes distinguishing input states more difficult, therefore, impacting accuracy. Note that the input cells are logically capable of producing 2^N states, where N is the number of input cells. However, in terms of the CRAM logic operations, the number of different effective resistance values usually falls into fewer

distinguishable groups (assuming all cells have identical electrical properties). For example, three input cells generate eight logic states, but there are only four distinguishable levels of resistance, since there are 0, 1, 2, or 3 cells among all three input cells to be in low or high resistance state. The second mechanism is that, with more input cells, the absolute amplitude of the logic voltage, V_{logic} , required to switch the output cell generally decreases due to the increased number of current paths from the input cells. Therefore, the portion of V_{logic} on the input cells also decreases. Because the effective TMR ratio of MTJs increases with smaller voltage bias, it leads to an increase in the effective TMR ratio for each input cell compared to a logic operation with fewer input cells. Figure 4c shows the accuracy of a 3-input MAJ3 logic operation. At $V_{\text{logic}} = -0.464$ V, both the optimal mean and the worst-case accuracy are observed to be 86.5% and 78.0%, respectively. Similarly, for a 5-input MAJ5 logic operation, shown in Fig. 4d, both the optimal mean and the worst-case accuracy are observed to be 75% and 56%, respectively. As expected, comparing 2-input, 3-input, and 5-input logic operation, the accuracy decreases with an increasing number of inputs.

Having demonstrated fundamental elements of CRAM – memory write operations, memory read operations, and logic operations – we turn to more complex operations. We demonstrate a 1-bit full adder. This device takes two 1-bit operands, A and B, as well as a 1-bit carry-in, C, as inputs, and outputs a 1-bit sum, S, and a 1-bit carry-out, C_{out} . A variety of implementations exist. We investigate two common designs: (1) one that uses a combination of majority and inversion logic gates, which we will refer to as a ‘MAJ+NOT’ design; and (2) one that uses only NAND gates, which we will refer to as an ‘all-NAND’ design. Figures 5a and 5b illustrate these designs.

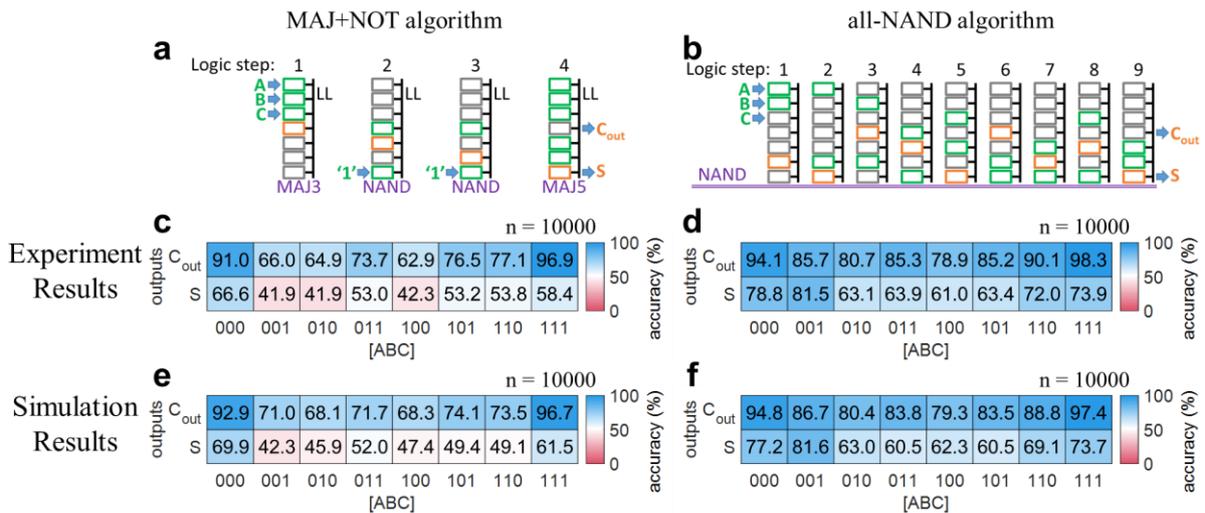


Fig. 5 | CRAM 1-bit full adder demonstration results. **a, b**, Illustrations of the ‘MAJ+NOT’ and ‘all-NAND’ 1-bit full adder designs. Green and orange letter symbols indicate input and output data for the full adder, respectively. From left to right, numbered by ‘logic step’, each drawing shows the intended input (green rectangle) and output (orange rectangle) cells involved in the logic operation. The text in purple under each drawing indicates the intended function of the logic operation (MAJ3, NAND, or MAJ5). **c, d, e, f**, Experimental (**c, d**) and simulation (**e, f**) results of the output accuracy of 1-bit full adder operations by CRAM with the MAJ+NOT (**c, e**) and all-NAND (**d, f**) designs. The CRAM adder’s outputs, S and C_{out} , are assessed against the expected values, i.e., their truth table, for all input states of A, B, and C. The accuracy of each result for each input state is shown by the numerical value in black font, as well as, represented by the color of the box with red (or blue) indicating wrong (or correct), or accuracy of 0% (100%). The accuracy is calculated based on the statistical average of outputs obtained by repeating the full adder execution n times, for $n = 10000$. The experimental results for the MAJ+NOT (**c**) and all-NAND (**d**) designs are obtained by repeatedly executing the operation for all input states and

observing the output states. The simulation results for the MAJ+NOT (e) and all-NAND (f) designs are obtained with probabilistic modeling, using Monte Carlo methods. The accuracy of individual logic operation is set to what was observed experimentally.

Figures 5c-f show both the experimental and simulation results for the MAJ+NOT and the all-NAND designs, respectively. Each plot is a colormap that lists the accuracy of the output bits S and C_{out} , with each input state coded as $[ABC]$. The blue (red) indicates good/desired (bad/undesired) accuracy. In the boxes of colormap, results in saturated blue are the most desirable. The numerical values of accuracy are also labeled accordingly.

The experimental results for the MAJ+NOT design of the full adder are shown in Fig. 5c. We make two observations:

- The first is that the accuracy of C_{out} is generally higher than that of S . This is likely due to the fact that C_{out} is directly produced by the first MAJ3 operation from inputs A , B , and C , while S is produced after multiple logic operations, including a MAJ5. Accordingly, we note that C_{out} is produced earlier than S so that it is less impacted by error propagation and accumulation during each step; and the MAJ5 involved in producing S is inherently less accurate than the MAJ3.
- The second observation is that both C_{out} and S have higher accuracy when the input $[ABC] = 000$ and when $[ABC]=111$ than in the other cases. This is expected as the input states of all '0's and all '1's yield higher accuracy than these with mixed numbers of '0's and '1's for both MAJ3 and MAJ5.

The experimental results for the all-NAND design are shown in Fig. 5d. The same observations regarding accuracy vs. inputs as the MAJ+NOT design apply. However, it is clear that the accuracy of the all-NAND full adder, at 78.5%, is higher than that of the MAJ+NOT full adder, at 63.8%. This is likely due to the fact that 2-input NAND operations are inherently more accurate than MAJ3 and MAJ5 operations. This offsets the impact of the additional steps required in the all-NAND design. We note that the accuracy of all computation blocks will improve as the underlying MTJ technology evolves. Accordingly, the relative accuracy of the all-NAND versus the MAJ+NOT designs may change⁶².

To understand the origin of errors, how they accumulate, and how they propagate, we performed numerical simulations of the full adder designs. These are based on probabilistic models of logic operations, implemented by Monte Carlo methods. Figures 5e and 5f show the simulation results for the MAJ+NOT and all-NAND design, respectively. In these, the accuracy of individual logic operations was set to match what was experimentally observed. The simulation results for the overall designs of the full adders correspond well to what was observed experimentally for these, which confirms the validity of the proposed probabilistic models.

We note that beyond the inherent inaccuracy of logic operations, other factors such as device drift and device-to-device variation in MTJ devices will contribute to error in a CRAM. Shifts in the overall output responses will occur. Specifically, shifts in temperature, external magnetic field, MTJ anisotropy, and MTJ resistance can lead to shifts of the response curve, $\langle D_{out} \rangle$. Most likely, any such shift will result in a reduction of the accuracy.

There could be several impacts of device-to-device variation. If the input MTJ cells' high or low resistance values are not closely grouped, then the logic states will separate into more levels. This will shrink the margin of accuracy in the worst cases. Another impact of device-to-device variation is disturbances to the data stored in the input cells during a logic operation. Ideally, a logic operation should only alter the state of the output cell while the data in input

cells should not be disturbed or destroyed. However, when input cells are not identical, some input cells could carry currents closer to their critical switching current so that they are more likely to be disturbed. For example, if the input and/or output cells happen to be these that are easier and/or harder to switch, the input cells are more likely to be disturbed. And disturbance of data in input cells may result in more errors. Therefore, the reduction or control of device-to-device variation, which can be achieved by utilizing industry-level production processes that are developed and matured for STT-MRAM, is critical and should bring major improvement (reduction) of CRAM accuracy (error rate).

The accuracy of logic operations will significantly benefit from TMR ratio improvements as the technology for MTJs evolves. To project what the accuracy of CRAM operations might be in the future, we apply several types of physical modeling, informed by existing experimental results.

Three sets of assumptions on the accuracies (or error rates) of NAND logic operations underlie the following studies.

- The ‘experimental’ assumptions are based on the best accuracy experimentally observed among the 9 NAND steps involved with the all-NAND 1-bit full adder. These are adjusted linearly to ensure that the error for inputs ‘01’ and ‘10’ equals that for input ‘11’. In reality, as supported by the experimental results shown in fig. 4a, such a condition can be reached by properly tuning the V_{logic} . Therefore, assuming the gate-level error rate is already optimized by tuning the V_{logic} , then the per-input-state NAND accuracies can be further simplified so that an error rate, δ ($0 \leq \delta \leq 1$), can be used to characterize the error, accuracy, and probabilistic truth table of NAND operations in a CRAM. The NAND accuracy is $[1, 1-\delta, 1-\delta, 1-\delta]$, and the NAND probabilistic truth table is $[1, 1-\delta, 1-\delta, \delta]$, both being a function of δ . Through the above-mentioned modeling and calculations, the ‘experimental’ assumptions yield $\delta = 0.0076$, which corresponds to a TMR ratio of approximately 109%, based on experiments.
- Two additional sets of assumptions, labelled as ‘production’, and ‘improved’ assume MTJ TMR ratios of 200%, and 300%, respectively. These two assumptions yield $\delta = 2.1 \times 10^{-4}$, and $\delta = 7.6 \times 10^{-6}$, respectively, based on modeling and calculations. The ‘production’ assumptions represent the current industry-level TMR ratios developed for STT-MRAM technologies. The ‘improved’ assumptions present reasonable expectations for near-future MTJ developments.

CRAM NAND error rates vs. TMR ratio with various logic voltage pulse widths are shown in Fig. 6a. Higher TMR ratios and faster logic speed – so shorter V_{logic} pulse widths – lead to smaller error rates. Note that, for all subsequent results, we will use the NAND error rate at the assumed TMR ratios, with pulse widths of 1 ms. This is more conservative but is consistent with the pulse widths used in the experimental results reported above.

With these defined sets of assumptions, we provide projections of CRAM accuracy at a larger scale for meaningful applications. First we evaluate ripple-carry adders and array multipliers⁶⁹ operating on scalar operands, with up to 6 bits. Then, with these primitives, we evaluate dot-product operations, which form the basis of matrix multiplication. Dot products consist of element-wise multiplication of two unsigned integer vectors, followed by addition. We perform additions with binary trees to maintain smaller circuit depth.

Finally, we apply the analysis to an image classification problem, using a multilayer perceptron based neural network. Here we use the MNIST digit recognition benchmark set⁶⁷.

The entire 10000 test images in the MNIST set are used for our study. We performed the parameter sweeps shown in Fig. 6c. The inference is performed as a sequence of matrix to vector multiplications with the serialized images, followed by sigmoid computations. The main component of the inference, unsigned fixed matrix-vector multiplications, was performed by a CRAM. The sigmoid calculation as well as the sign handling was handled off the CRAM. Figure 6b shows the multilayer perceptron architecture. It receives serialized 28×28 images, passes these through two layers of computation, and assigns them to one of 10 output classes. As shown in Fig. 6c, the classification accuracy on the MNIST test set increases and then saturates as the bit width of the input data increases for both the ‘experimental’ and the ‘production’ assumptions. But the accuracy is almost constant vs. bit width for ‘improved’ assumptions, as it gets very close to the ‘baseline’, which is the ceiling for MNIST accuracy with our specific neural network implementation ran on a perfect, error-free hardware. At bit widths of 6, the MNIST accuracy is 19.74%, 90.13%, 97.09%, and 97.26%, for the ‘experimental’, ‘production’, ‘improved’, and ‘baseline’ assumptions, respectively. And the best accuracy observed is 19.74% (6-bit), 90.13% (6-bit), 97.48% (4-bit), and 97.64% (4-bit), for the ‘experimental’, ‘production’, ‘improved’, and ‘baseline’ assumptions, respectively. Additionally, we also evaluate a simpler implementation. The only difference is the number of hidden layer neurons is reduced to 64. The ‘baseline’ accuracy is, as expected, lower, with a value of 92.03% (6-bit). Both the ‘production’ and ‘improved’ accuracy levels are also lower, with values of 86.16% (6-bit) and 91.81% (6-bit), respectively. Interestingly, the ‘experimental’ accuracy, being 27.25% (6-bit), is higher than that from the more complex implementation (presented here) with 300 hidden layer neurons. As the implementation becomes bigger or more complex, the baseline accuracy increases. For those assumptions that yield sufficiently low gate-level error rate, the accuracy drop from baseline due to the gate-level error is less than the improvement of baseline due to more complex neural network implementation. However, when the gate-level error rate is larger, the less complex implementation performs better. Note that different neural network training and implementation⁷⁰ results in different ‘baseline’ MNIST accuracy in MNIST benchmark. The specific implementation of the classifier is relatively simple and does not produce the highest MNIST accuracy if compared to other more complex implementations. These results are not intended to demonstrate the highest absolute value of MNIST accuracy, but rather to show that the CRAM with finite gate-level error rate can deliver minimal MNIST accuracy drop from ‘baseline’ or accuracy degradation. These results confirm the technological relevance of CRAM for larger scale, meaningful applications.

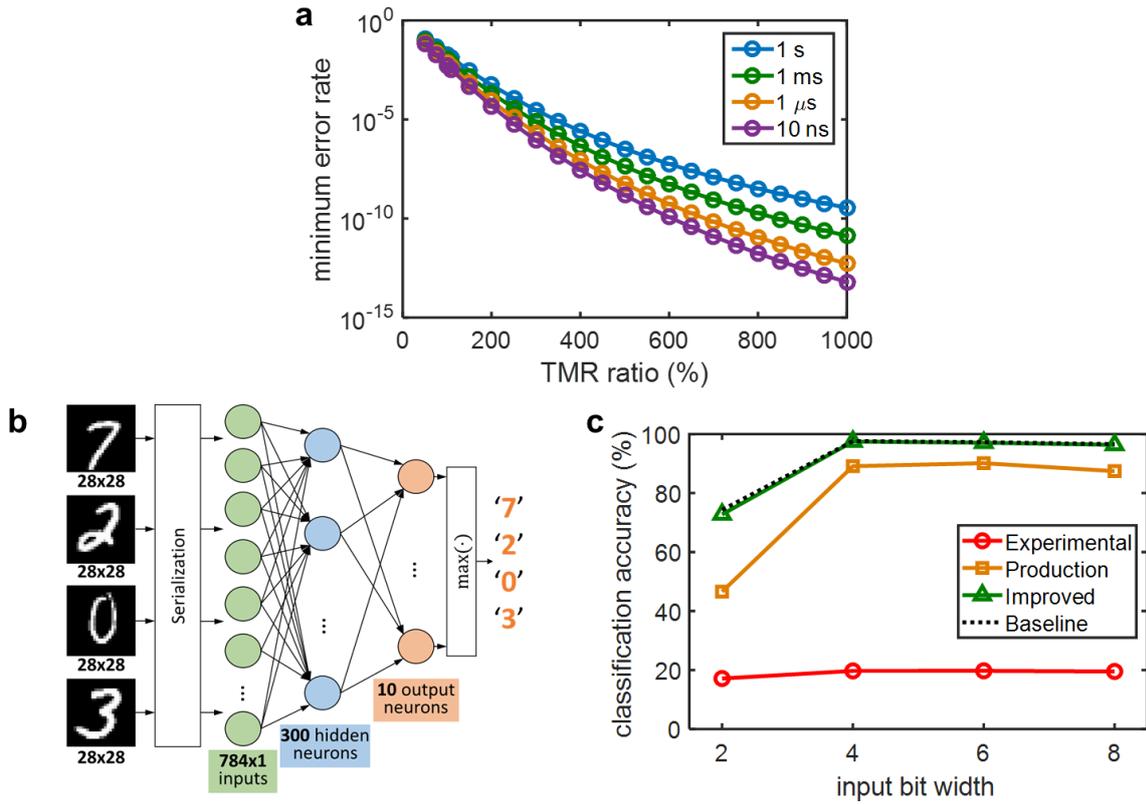


Fig. 6 | CRAM accuracy projections. **a**, NAND gate minimum error rate vs. MTJ TMR ratio with various V_{logic} pulse widths. For a given TMR ratio, the error rate is a function of V_{logic} . So, the ‘minimum error rate’ represents the minimum error rate achievable with an appropriate V_{logic} value. All subsequent studies use the error rates observed with 1 ms pulse widths (to be consistent with the earlier experimental studies) at assumed TMR ratios. **b**, Illustration of the architecture of the CRAM-based handwritten digit classifier. **c**, Accuracy of the CRAM-based MNIST handwritten digit classifier with various assumptions on the gate-level error rate. Results are obtained with the full 10000 test images from the MNIST dataset. The calculated NAND gate error rates for pulse widths of 1 ms from (a) are used.

Discussions

To summarize the experimental work, an MTJ-based 1×7 CRAM array hardware was experimentally demonstrated and systematically evaluated. The basic memory write and read operations of CRAM were achieved with high reliability. The study on CRAM logic operations began with 2-input logic operations. It was found that a 2-input NAND operation could be performed with accuracy as high as 99.4%. As the number of input cells was increased, for example, for 3-input MAJ3 and 5-input MAJ5 operations, the accuracy decreased to 86.5% and 75%, respectively. The decrease was attributed to having too many levels corresponding to the input states crowding a limited operating margin. Next, two versions of a 1-bit full adder were experimentally demonstrated using the 1×7 CRAM array: an all-NAND version and a MAJ+NOT version. The all-NAND design achieved an accuracy of 78.5% while the seemingly simpler MAJ+NOT, which involves 3- and 5-input MAJ operations, only achieved an accuracy of 63.8%. Note that although each type of logic operation achieves optimal accuracy performance with a specific voltage value, the value is expected to only need to be static or constant. Therefore, only a finite number of power rails is needed to accommodate the logic operations of the CRAM array. Also, if the multiple logic pulse duration is allowed by a peripheral design, it is possible to operate the CRAM array with a single set of power rails for both memory write and logic operations.

A probabilistic model was proposed that accounts for the origin of errors, their propagation, and their accumulation during a multi-step CRAM operation. The model was shown to be effective when matched with the experimental results for the 1-bit full adder. The working principles of this model were adopted for the rest of the studies.

A suite of MTJ device circuit models were fitted to the existing experimental data and used to project CRAM NAND gate-level accuracy in the form of error rates. The gate-level error rates were shown to be 7.6×10^{-6} , with reasonable expectations of TMR ratio improvement as MTJ technology develops. Other device properties, such as the switching probability transfer curve, could also significantly affect the CRAM gate-level error rate. This calls for improvements or new discoveries of the physical mechanisms for memory read-out, and memory write. Error is an inherent property of any physical hardware, including CMOS logic components, which are commonly perceived as deterministic. As the development of CRAM proceeds, the gate-level error rate of CRAM will be further reduced over time. For now, while the error rate of CRAM is still higher compared to that of CMOS logic circuits, CRAM is naturally more suitable for applications that require less precision but can still benefit from the true-in-memory computing features and advantages of CRAM, instead of those that require high precision and determinism. Additionally, logic operations with many inputs, such as majority, may be desirable in certain scenarios. And yet, these were shown to have lower accuracy than 2-input operations. So, a tradeoff might exist.

Building on the experimental demonstration and evaluation of the 1-bit full adder designs, simulation and analysis was performed for larger functional circuits: scalar addition and multiplication up to 6 bits and matrix multiplication up to 5 bits, with an input size up to 4×4 . The parameters for the simulations were experimentally measured values as well as reasonable projections for future MTJ technology. The results show promising accuracy performance of the CRAM at a functional level.

Finally, a detailed large-scale simulation and analysis of a meaningful application was performed: an MNIST handwritten digit classifier. The results reveal that, with the experimental TMR ratio set to 109%, the classification accuracy is unsatisfactory. But when the TMR ratio is set to 200% and 300%, satisfactory and almost-perfect classification accuracy are achieved, respectively. Improvement of CRAM performances by improving TMR ratio can be reasonably expected since more than 600% TMR ratio was recently achieved experimentally at room temperature⁷¹. It was noted that with shorter logic pulses, so higher speed, the gate-level error improves. Besides the TMR ratio and logic speed, changes of various other design parameters and switching mechanisms, such as reducing the barrier thickness and increasing thermal stability, could be used to improve the gate-level error rate. Accordingly, satisfactory classification accuracy could be achieved with a TMR ratio that is significantly lower than 200%.

The MNIST⁶⁷ classifier implemented by CRAM is relatively small, compared to other representative designs of neural networks^{6,72,81,82,73-80}, in terms of number of parameters and number of layers. However, with the similar methodology employed in a recent work on crossbar in-memory computing²⁰, it is possible to further extend the architectural support to a larger neural network, as the computational primitives needed are similar to those studied in this work. The expected accuracy of a neural network could be highly nonlinear or weakly related to the size of it if it were to be implemented by CRAM with finite gate-level error rate. Therefore, it is possible that the currently ‘improved’ CRAM could perform satisfactorily even if a larger neural network is implemented. Nevertheless, a careful case-by-case analysis would be needed to determine the exact sensitivity to error rates. However, we note that, if the CRAM achieves reasonable low error rates, these applications may become feasible with the

technology, especially if error-correction is implemented. Lastly, extension to a bigger real-life application would still require bigger CRAM arrays, although the specific size is subject to the specific implementation.

In summary, this work serves as the first step in experimentally demonstrating the viability, feasibility, and realistic properties of MTJ-based CRAM hardware. Through modeling and simulation, it also lays out the foundation for a coherent view of CRAM, from the device physics level up to the application level. Prior work had established the potential of CRAM through numerical simulation only. Accordingly, there had been considerable interest in the unique features, speed, power, and energy benefits of the technology. This study puts the earlier work on a firm experimental footing, providing application-critical metrics of gate-level accuracy or error rate and linking it to the application accuracy. It paves the way for future work on large scale applications, in conventional domains as well as new ones emerging in machine intelligence. It also indicates the possibility of making competitive large-scale CMOS-integrated CRAM hardware.

Acknowledgements

This work was supported in part by the Defense Advanced Research Projects Agency (DARPA) via No. HR001117S0056-FP-042 “Advanced MTJs for computation in and near random access memory” and by the National Institute of Standards and Technology. This work was supported in part by NSF SPX grant no. 1725420. The authors also thank Cisco Inc. for the support. Portions of this work were conducted in the Minnesota Nano Center, which was supported by the National Science Foundation through the National Nanotechnology Coordinated Infrastructure (NNCI) under Award No. ECCS-2025124. The authors acknowledge the Minnesota Supercomputing Institute (MSI, URL: <http://www.msi.umn.edu>) at the University of Minnesota for providing resources that contributed to the research results reported within this paper. The authors thank Prof. Marc Riedel and Prof. John Sartori from the Department of Electrical and Computer Engineering at University of Minnesota for proofreading the manuscript. Yang Lv, Brandon Zink, and Hüsrev Cilasun are CISCO Fellow.

Author Contributions

J.-P.W. conceived the CRAM research and coordinated the entire project. Y.L. and J.-P.W. designed the experiments. Y.L. and R.P.B. designed and developed the demonstration hardware and software. P.K., A.H. and W.W. grew part of the perpendicular MTJ stacks. B.R.Z. fabricated the MTJ nano devices. Y.L. conducted the CRAM demonstration experiments and analyzed the results. Y.L. studied the probabilistic model of CRAM operations and conducted simulations of 1-bit full adder. Y.L., B.R.Z., and R.P.B. developed the device physics modeling of CRAM logic operations and gate-level error rate and conducted related calculations. H.C., S.R., Z.C., and U.K. carried out the simulation studies of multi-bit adder, multiplier, matrix multiplication, and the MNIST handwritten digit classification. S.S. participated in discussions of modeling and simulation. All authors reviewed and discussed the results. Y.L. and J.-P.W. wrote the draft manuscript. All authors contributed to the completion of the manuscript.

References

1. Lecun, Y., Bengio, Y. & Hinton, G. Deep learning. *Nature* **521**, 436–444 (2015).
2. Jordan, M. I. & Mitchell, T. M. Machine learning: trends, perspectives, and prospects. *Science* **349**, 255–260 (2015).
3. Adomavicius, G. & Tuzhilin, A. Toward the next generation of recommender systems: A survey of the state-of-the-art and possible extensions. *IEEE Trans. Knowl. Data Eng.*

- 17, 734–749 (2005).
4. Hinton, G. *et al.* Deep neural networks for acoustic modeling in speech recognition: The shared views of four research groups. *IEEE Signal Process. Mag.* **29**, 82–97 (2012).
 5. Collobert, R. *et al.* Natural language processing (almost) from scratch. *J. Mach. Learn. Res.* **12**, 2493–2537 (2011).
 6. Krizhevsky, A., Sutskever, I. & Hinton, G. E. ImageNet classification with deep convolutional neural networks. *Commun. ACM* **60**, 84–90 (2017).
 7. Oh, K. S. & Jung, K. GPU implementation of neural networks. *Pattern Recognit.* **37**, 1311–1314 (2004).
 8. Strigl, D., Kofler, K. & Podlipnig, S. Performance and scalability of GPU-based convolutional neural networks. in *2010 18th Euromicro Conference on Parallel, Distributed and Network-based Processing* 317–324 (IEEE, 2010).
 9. Nurvitadhi, E. *et al.* Accelerating binarized neural networks: comparison of FPGA, CPU, GPU, and ASIC. in *2016 International Conference on Field-Programmable Technology (FPT)* 77–84 (IEEE, 2017).
 10. Sawada, J. *et al.* TrueNorth ecosystem for brain-inspired computing: scalable systems, software, and applications. in *SC '16: Proceedings of the International Conference for High Performance Computing, Networking, Storage and Analysis* 130–141 (IEEE, 2016).
 11. Jouppi, N. P. *et al.* In-datacenter performance analysis of a tensor processing unit. in *Proceedings of the 44th Annual International Symposium on Computer Architecture* 1–12 (ACM, 2017).
 12. Chen, Y. H., Krishna, T., Emer, J. S. & Sze, V. Eyeriss: an energy-efficient reconfigurable accelerator for deep convolutional neural networks. *IEEE J. Solid-State Circuits* **52**, 127–138 (2017).
 13. Yin, S. *et al.* A high energy efficient reconfigurable hybrid neural network processor for deep learning applications. *IEEE J. Solid-State Circuits* **53**, 968–982 (2018).
 14. Borghetti, J. *et al.* Memristive switches enable stateful logic operations via material implication. *Nature* **464**, 873–876 (2010).
 15. Chi, P. *et al.* PRIME: a novel processing-in-memory architecture for neural network computation in ReRAM-based main memory. *ACM SIGARCH Comput. Archit. News* **44**, 27–39 (2016).
 16. Shafiee, A. *et al.* ISAAC: a convolutional neural network accelerator with in-situ analog arithmetic in crossbars. *ACM SIGARCH Comput. Archit. News* **44**, 14–26 (2016).
 17. Hu, M. *et al.* Dot-product engine for neuromorphic computing: programming 1T1M crossbar to accelerate matrix-vector multiplication. in *2016 53rd ACM/EDAC/IEEE Design Automation Conference (DAC)* 1–6 (IEEE, 2016).
 18. Seshadri, V. *et al.* Ambit: in-memory accelerator for bulk bitwise operations using commodity DRAM technology. in *2017 50th Annual IEEE/ACM International Symposium on Microarchitecture (MICRO)* 273–287 (IEEE, 2017).
 19. Yao, P. *et al.* Fully hardware-implemented memristor convolutional neural network. *Nature* **577**, 641–646 (2020).
 20. Jung, S. *et al.* A crossbar array of magnetoresistive memory devices for in-memory computing. *Nature* **601**, 211–216 (2022).
 21. Keckler, S. W., Dally, W. J., Khailany, B., Garland, M. & Glasco, D. GPUs and the future of parallel computing. *IEEE Micro* **31**, 7–17 (2011).
 22. Bergman, K. *et al.* *ExaScale Computing Study: Technology Challenges in Achieving Exascale Systems*. www.cse.nd.edu/Reports/2008/TR-2008-13.pdf (2008).
 23. Horowitz, M. Computing’s energy problem (and what we can do about it). in *2014*

- IEEE International Solid-State Circuits Conference Digest of Technical Papers (ISSCC)* 10–14 (IEEE, 2014).
24. Kim, D., Kung, J., Chai, S., Yalamanchili, S. & Mukhopadhyay, S. Neurocube: a programmable digital neuromorphic architecture with high-density 3D memory. *ACM SIGARCH Comput. Archit. News* **44**, 380–392 (2016).
 25. Huang, J. *et al.* Active-routing: compute on the way for near-data processing. in *2019 IEEE International Symposium on High Performance Computer Architecture (HPCA)* 674–686 (IEEE, 2019).
 26. Nair, R. *et al.* Active memory cube: a processing-in-memory architecture for exascale systems. *IBM J. Res. Dev.* **59**, 17:1-17:14 (2015).
 27. Pawlowski, J. T. Hybrid memory cube (HMC). in *2011 IEEE Hot Chips 23 Symposium (HCS)* 1–24 (IEEE, 2011).
 28. Gao, M., Ayers, G. & Kozyrakis, C. Practical near-data processing for in-memory analytics frameworks. in *2015 International Conference on Parallel Architecture and Compilation (PACT)* 113–124 (IEEE, 2015).
 29. Gao, M., Pu, J., Yang, X., Horowitz, M. & Kozyrakis, C. TETRIS: scalable and efficient neural network acceleration with 3D memory. *SIGARCH Comput. Arch. News* **45**, 751–764 (2017).
 30. Matsunaga, S. *et al.* Fabrication of a Nonvolatile Full Adder Based on Logic-in-Memory Architecture Using Magnetic Tunnel Junctions. *Appl. Phys. Express* **1**, 091301 (2008).
 31. Hanyu, T. *et al.* Standby-power-free integrated circuits using MTJ-based VLSI computing. *Proc. IEEE* **104**, 1844–1863 (2016).
 32. Li, S. *et al.* Pinatubo: a processing-in-memory architecture for bulk bitwise operations in emerging non-volatile memories. in *2016 53rd ACM/EDAC/IEEE Design Automation Conference (DAC)* 1–6 (IEEE, 2016).
 33. Aga, S. *et al.* Compute caches. in *2017 IEEE International Symposium on High Performance Computer Architecture (HPCA)* 481–492 (IEEE, 2017).
 34. Zidan, M. A., Strachan, J. P. & Lu, W. D. The future of electronics based on memristive systems. *Nat. Electron.* *2017 11* **1**, 22–29 (2018).
 35. Jeon, K., Ryu, J. J., Jeong, D. S. & Kim, G. H. Dot-Product Operation in Crossbar Array Using a Self-Rectifying Resistive Device. *Adv. Mater. Interfaces* **9**, 2200392 (2022).
 36. Wang, J.-P. & Harms, J. D. General structure for computational random access memory (CRAM). US patent 14/259,568 (2015).
 37. Chowdhury, Z. *et al.* Efficient in-memory processing using spintronics. *IEEE Comput. Archit. Lett.* **17**, 42–46 (2018).
 38. Kvatinsky, S. *et al.* MAGIC—memristor-aided logic. *IEEE Trans. Circuits Syst. II Express Briefs* **61**, 895–899 (2014).
 39. Kvatinsky, S. *et al.* Memristor-based material implication (IMPLY) logic: Design principles and methodologies. *IEEE Trans. Very Large Scale Integr. Syst.* **22**, 2054–2066 (2014).
 40. Gupta, S., Imani, M. & Rosing, T. FELIX: fast and energy-efficient logic in memory. in *2018 IEEE/ACM International Conference on Computer-Aided Design (ICCAD)* 1–7 (IEEE, 2018).
 41. Gao, F., Tziantzioulis, G. & Wentzlaff, D. ComputeDRAM: in-memory compute using off-the-shelf DRAMs. in *2019 52nd Annual IEEE/ACM International Symposium on Microarchitecture (MICRO)* 100–113 (IEEE, 2019).
 42. Žutić, I., Fabian, J. & Sarma, S. Das. Spintronics: fundamentals and applications. *Rev. Mod. Phys.* **76**, 323–410 (2004).
 43. Nikonov, D. E. & Young, I. A. Benchmarking of beyond-CMOS exploratory devices

- for logic integrated circuits. *IEEE J. Explor. Solid-State Comput. Devices Circuits* **1**, 3–11 (2015).
44. Shiokawa, Y. *et al.* High write endurance up to 1012 cycles in a spin current-type magnetic memory array. *AIP Adv.* **9**, 035236 (2019).
 45. Jan, G. *et al.* Demonstration of ultra-low voltage and ultra low power STT-MRAM designed for compatibility with 0x node embedded LLC applications. in *2018 IEEE Symposium on VLSI Technology* 65–66 (IEEE, 2018).
 46. Zhao, H. *et al.* Sub-200 ps spin transfer torque switching in in-plane magnetic tunnel junctions with interface perpendicular anisotropy. *J. Phys. D. Appl. Phys.* **45**, 025001 (2012).
 47. Julliere, M. Tunneling between ferromagnetic films. *Phys. Lett. A* **54**, 225–226 (1975).
 48. Parkin, S. S. P. *et al.* Giant tunnelling magnetoresistance at room temperature with MgO (100) tunnel barriers. *Nat. Mater.* **3**, 862–867 (2004).
 49. Yuasa, S., Nagahama, T., Fukushima, A., Suzuki, Y. & Ando, K. Giant room-temperature magnetoresistance in single-crystal Fe/MgO/Fe magnetic tunnel junctions. *Nat. Mater.* **3**, 868–871 (2004).
 50. Berger, L. Emission of spin waves by a magnetic multilayer traversed by a current. *Phys. Rev. B* **54**, 9353–9358 (1996).
 51. Slonczewski, J. C. Current-driven excitation of magnetic multilayers. *J. Magn. Magn. Mater.* **159**, L1–L7 (1996).
 52. Wei, L. *et al.* 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 *2019 IEEE International Solid- State Circuits Conference - (ISSCC)* 214–216 (IEEE, 2019).
 53. Gallagher, W. J. *et al.* 22nm STT-MRAM for reflow and automotive uses with high yield, reliability, and magnetic immunity and with performance and shielding options. in *2019 IEEE International Electron Devices Meeting (IEDM)* 2.7.1-2.7.4 (IEEE, 2019).
 54. Chih, Y. Der *et al.* A 22nm 32Mb embedded STT-MRAM with 10ns read speed, 1M cycle write endurance, 10 years retention at 150°C and high immunity to magnetic field interference. in *2020 IEEE International Solid- State Circuits Conference - (ISSCC)* 222–224 (IEEE, 2020).
 55. Edelstein, D. *et al.* A 14 nm embedded STT-MRAM CMOS technology. in *2020 IEEE International Electron Devices Meeting (IEDM)* 11.5.1-11.5.4 (IEEE, 2020).
 56. Chun, K. C. *et al.* A scaling roadmap and performance evaluation of in-plane and perpendicular MTJ based STT-MRAMs for high-density cache memory. *IEEE J. Solid-State Circuits* **48**, 598–610 (2013).
 57. Lilja, D. J. *et al.* Systems and methods for direct communication between magnetic tunnel junctions. US patent 13/475,544 (2014).
 58. Lyle, A. *et al.* Direct communication between magnetic tunnel junctions for nonvolatile logic fan-out architecture. *Appl. Phys. Lett.* **97**, 152504 (2010).
 59. Zabihi, M. *et al.* Using spin-Hall MTJs to build an energy-efficient in-memory computation platform. in *20th International Symposium on Quality Electronic Design (ISQED)* 52–57 (IEEE, 2019).
 60. Chowdhury, Z. I. *et al.* A DNA read alignment accelerator based on computational RAM. *IEEE J. Explor. Solid-State Comput. Devices Circuits* **6**, 80–88 (2020).
 61. Chowdhury, Z. *et al.* CRAM-Seq: accelerating RNA-Seq abundance quantification using computational RAM. *IEEE Trans. Emerg. Top. Comput.* **10**, 2055–2071 (2022).
 62. Zabihi, M. *et al.* In-memory processing on the spintronic CRAM: from hardware design to application mapping. *IEEE Trans. Comput.* **68**, 1159–1173 (2019).
 63. Cilasun, H. *et al.* CRAFFT: High resolution FFT accelerator in spintronic

- computational RAM. in *2020 57th ACM/IEEE Design Automation Conference (DAC)* 1–6 (IEEE, 2020).
64. Resch, S. *et al.* PIMBALL: Binary neural networks in spintronic memory. *ACM Trans. Archit. Code Optim.* **16**, 41 (2019).
 65. Chowdhury, Z. I. *et al.* CAMEleon: reconfigurable B(T)CAM in computational RAM. in *Proceedings of the 2021 on Great Lakes Symposium on VLSI* 57–63 (ACM, 2021).
 66. Resch, S. *et al.* MOUSE: inference in non-volatile memory for energy harvesting applications. in *2020 53rd Annual IEEE/ACM International Symposium on Microarchitecture (MICRO)* 400–414 (IEEE, 2020).
 67. LeCun, Y., Cortes, C. & Burges, C. J. C. The MNIST database of handwritten digits. <http://yann.lecun.com/exdb/mnist/> (1998).
 68. Lv, Y., Bloom, R. P. & Wang, J.-P. Experimental demonstration of probabilistic spin logic by magnetic tunnel junctions. *IEEE Magn. Lett.* **10**, 1–5 (2019).
 69. Subathradevi, S. & Vennila, C. Systolic array multiplier for augmenting data center networks communication link. *Cluster Comput.* **22**, 13773–13783 (2019).
 70. Guo, X. *et al.* Fast, energy-efficient, robust, and reproducible mixed-signal neuromorphic classifier based on embedded NOR flash memory technology. in *2017 IEEE International Electron Devices Meeting (IEDM)* 6.5.1–6.5.4 (IEEE, 2017).
 71. 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).
 72. Simonyan, K. & Zisserman, A. Very deep convolutional networks for large-scale image recognition. *arXiv Prepr.* arXiv:1409.1556 (2015) doi:10.48550/arxiv.1409.1556.
 73. Szegedy, C. *et al.* Going deeper with convolutions. in *2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)* 1–9 (IEEE, 2015).
 74. He, K., Zhang, X., Ren, S. & Sun, J. Deep residual learning for image recognition. in *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)* 770–778 (IEEE, 2016).
 75. Bojarski, M. *et al.* End to end learning for self-driving cars. *arXiv Prepr.* arXiv:1604.07316 (2016) doi:10.48550/arxiv.1604.07316.
 76. Zagoruyko, S. & Komodakis, N. Wide residual networks. in *Proceedings of the British Machine Vision Conference (BMVC)* 87.1–87.12 (BMVA Press, 2016).
 77. Devlin, J., Chang, M. W., Lee, K. & Toutanova, K. BERT: pre-training of deep bidirectional transformers for language understanding. *arXiv Prepr.* arXiv:1810.04805 (2018) doi:10.48550/arxiv.1810.04805.
 78. Kocić, J., Jovičić, N. & Drndarević, V. An end-to-end deep neural network for autonomous driving designed for embedded automotive platforms. *Sensors* **19**, 2064 (2019).
 79. Brown, T. B. *et al.* Language models are few-shot learners. *Adv. Neural Inf. Process. Syst.* **33**, 1877–1901 (2020).
 80. Raffel, C. *et al.* Exploring the limits of transfer learning with a unified text-to-text transformer. *J. Mach. Learn. Res.* **21**, 1–67 (2020).
 81. Ramesh, A. *et al.* Zero-shot text-to-image generation. in *Proceedings of the 38th International Conference on Machine Learning* 8821–8831 (PMLR, 2021).
 82. Fedus, W., Zoph, B. & Shazeer, N. Switch transformers: scaling to trillion parameter models with simple and efficient sparsity. *J. Mach. Learn. Res.* **23**, 1–40 (2022).