

# Multistate Register Based on Resistive RAM

Ravi Patel, *Student Member, IEEE*, Shahar Kvatinsky, *Student Member IEEE*, Eby G. Friedman, *Fellow IEEE*, and Avinoam Kolodny, *Senior Member IEEE*

**Abstract**—In recent years, memristive technologies, such as resistive RAM (RRAM), have emerged. These technologies are usually considered as replacements to SRAM, DRAM, and Flash. In this paper, a novel digital circuit, the multistate register, is proposed. The multistate register is different than conventional types of memory, and is used to store multiple data bits, where only a single bit is active and the remaining data bits are idle. The active bit is stored within a CMOS flip flop, while the idle bits are stored in an RRAM crossbar co-located with the flip flop. It is demonstrated that additional states require an area overhead of 1.4% per state for a 64 state register. The use of multistate registers as pipeline registers is demonstrated for a novel multithreading architecture – continuous flow multithreading (CFMT), where the total area overhead in the CPU pipeline is only 2.5% for 16 threads as compared to a single thread CMOS pipeline. The use of multistate registers in the CFMT microarchitecture enables higher performance processors (40% average performance improvement) with relatively low energy (6.5% average energy reduction) and area overhead.

**Keywords**—RRAM; memristor; memristive device; flip flop; multithreading

## I. INTRODUCTION

Memristive technologies [1-3] have been proposed to augment existing state-of-the-art CMOS circuits. One interesting memristive technology is resistive RAM (RRAM) [5-9]. RRAM-based memories can be integrated with existing digital circuits to increase functionality and system throughput. RRAM is a two terminal device that exhibits the properties of nonvolatility and high density. Unlike charge-based memories, information in an RRAM is stored by modulating the material state. An RRAM cell dissipates no static power to store a state and provides immunity to radiation and noise induced soft errors. Fabrication of these devices generally requires deposition of a

Manuscript received 17th December 2013; revised 6th May 2014; accepted xxx. This work was partially supported by Hasso Plattner Institute, by the Advanced Circuit Research Center at the Technion, by the Binational Science Foundation under Grant no. 2012139, by the National Science Foundation under Grant no. CCF-1329314, and by grants from Qualcomm, Cisco Systems, and Samsung.

R. Patel and E. G. Friedman are with the Department of Electrical Engineering and Computer Engineering, University of Rochester, Rochester, NY 14627, USA.

S. Kvatinsky and A. Kolodny are with the Department of Electrical Engineering, Technion – Israel Institute of Technology, Haifa 32000, Israel. (S. Kvatinsky corresponding author phone: 972-77887-1923; fax: 972-4829-5757; e-mail: skva@tx.technion.ac.il).

Copyright (c) 2014 IEEE. Personal use of this material is permitted. However, permission to use this material for any other purposes must be obtained from the IEEE by sending an email to pubs-permissions@ieee.org.

thin film material. The integration of these devices with CMOS is constrained primarily by lithographic patterning limits. Thus memristors scale with existing CMOS technologies.

The traditional approach of increasing CPU clock frequency has abated due to constraints on power consumption and density. To increase performance with each CMOS generation, thread level parallelism is exploited with multi-core processors [10]. This approach utilizes an increasing number of CMOS transistors to support additional cores on the same die, rather than increase the frequency of a single processor. This larger number of cores, however, dissipates greater static power. Multithreading is an approach to enhance the performance of an individual core by increasing logic utilization [11], without consuming additional static power. Handling each thread, however, requires duplication of resources (e.g., register files, flags, pipeline registers). This added overhead increases the area, power, and complexity of the processor, potentially increasing on-chip signal delays. The thread count is therefore typically limited to two to four threads per core in modern general purpose processors [12].

The high density, nonvolatility, and soft error immunity exhibited by resistive random access memory enables novel tradeoffs in digital circuits, allowing new mechanisms to increase thread count without increasing the static power. These tradeoffs support innovative memory structures for novel microarchitectures. In this paper, a memristive multistate pipeline register (MPR) is proposed that exploits these properties to enable higher throughput computing. The MPR is compatible with existing digital circuits while leveraging RRAM devices to store multiple machine states within a single register. This behavior enables an individual logic pipeline to be densely integrated with memory while retaining state information for multiple independent on-going operations. The state information for each operation is stored within a local memory and recalled at a later time, allowing the computation to resume without flushing the pipeline.

This functionality is useful in multithreaded processors to store the state of different threads. This situation is demonstrated in the case study of a novel microarchitecture – continuous flow multithreading (CFMT) [13]. It is shown that including an RRAM MPR within the CFMT microarchitecture enhances the performance of a processor, on average, by 40%, while reducing the energy, on average, by 6.5%. The proposed MPR circuit can also be used as a multistate register for applications other than pipeline registers.

Background of RRAM and crosspoint style memories is reviewed in Section II. The operation of the multistate register is presented in Section III. The simulation setup and circuit evaluation process are described in Section IV. A case study examining the multistate register as a pipeline register within a CPU is presented in Section V, followed by some concluding remarks in Section VI.

## II. BACKGROUND

Memristors and memristor-based arrays behave differently than standard CMOS SRAM memory arrays due to the different properties of RRAM devices. Operation of memristive devices and memristor-based crosspoint structures is described in the following section.

### A. Background of memristor and RRAM

Memristors [14] and memristive devices [15] behave as resistors, where the resistance is modulated by an applied bias. Positive and negative biases increase or decrease, respectively, the resistance of the device. In general, a bias applied for a longer duration produces a greater change in resistance. A larger voltage will generally increase the speed of the change in resistance. The device may also exhibit a threshold voltage or current, such that the resistance will change only if the bias exceeds the threshold specific to the device technology [16-18]. Once the bias is removed, the final resistance of the memristor is retained without dissipating any power.

One interesting memristive technology is RRAM, where oxide-based materials (e.g., TaO, TiO, SiO) [19, 20] rely on the migration of dopants to switch the resistance of a tunnel barrier. Dopant chains form through the oxide and reduce the thickness of the tunneling gap. An increase in the gap thickness gives rise to an increase in the resistance of the device while a decrease reduces the resistance. Currently, RRAM is considered a good candidate to replace Flash memory and is being widely investigated both in industry and academia.

The exact behavior of RRAM devices varies for different oxide materials. To simulate the behavior of memristive circuits, a general device model is used – the TEAM model [20]. In the TEAM model, the behavior of the resistive device is represented by the following expressions,

$$\frac{dx(t)}{dt} = \begin{cases} k_{off} \cdot \left(\frac{i(t)}{i_{off}} - 1\right)^{\alpha_{off}} \cdot f_{off}(x), & 0 < i_{off} < i, & (1a) \\ 0, & i_{on} < i < i_{off}, & (1b) \\ k_{on} \cdot \left(\frac{i(t)}{i_{on}} - 1\right)^{\alpha_{on}} \cdot f_{on}(x), & i < i_{on} < 0, & (1c) \end{cases}$$

where  $k_{off}$  and  $k_{on}$  are fitting parameters,  $\alpha_{on}$  and  $\alpha_{off}$  are adaptive nonlinearity parameters,  $i_{off}$  and  $i_{on}$  are current threshold parameters,  $f_{on}(x)$  and  $f_{off}(x)$  are window functions,

$$v(t) = \left[ R_{ON} + \frac{R_{OFF} - R_{ON}}{x_{off} - x_{on}} (x - x_{on}) \right] \cdot i(t), \quad (2)$$

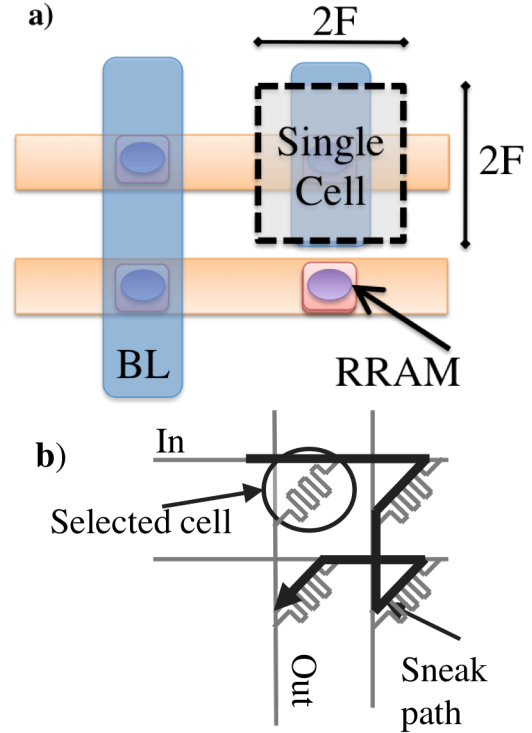


Figure 1. RRAM crosspoint (a) structure, and (b) an example of a parasitic sneak path within a 2 x 2 crosspoint array.

$R_{ON}$  and  $R_{OFF}$  are, respectively, the minimum and maximum resistance of the memristor, and  $x_{on}$  and  $x_{off}$  are, respectively, the minimum and maximum allowed value of the internal state variable  $x$ . The window function returns a value between zero and one and describes the rate at which the change of the state variable becomes nonlinear near the minimum and maximum resistance of a memristor. A Joglekar window function is used with a p-coefficient of two [22]. An I-V curve of a memristive device based on the TEAM model is shown in Figure 2a, exhibiting a pinched hysteresis loop.

### B. Crosspoints and nonlinearity

RRAM exhibits high density when utilized in a crosspoint array configuration. In this structure, a thin film is sandwiched between two sets of parallel interconnects. Each set of interconnects is orthogonal, allowing any individual memristive device to be selected by biasing one vertical and one horizontal metal line. In this configuration, the circuit density is only limited by the available metal pitch. The structure of a crosspoint is shown in Figure 1a.

Crosspoint arrays have the inherent problem of sneak path currents where currents propagate between two selected lines through unselected memristors. The sneak path phenomenon is illustrated in Figure 1b. The nonlinear I-V characteristic of certain memristive devices lessens the sneak path phenomenon [23]. This nonlinearity can be achieved by depositing additional materials above or below the memristive thin film. Depending on the material system used for RRAM, the nonlinearity can result from an insulator to metal transition or

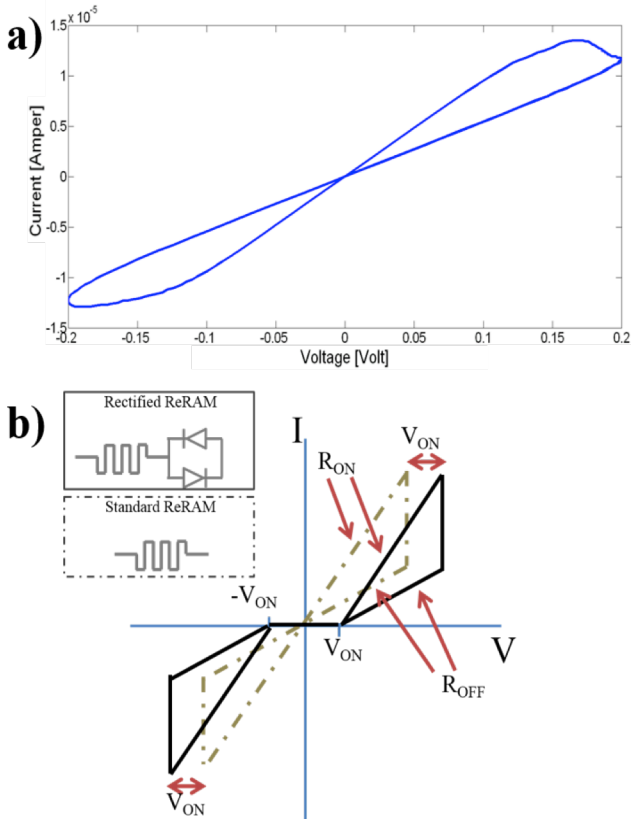


Figure 2. I-V characteristic of a memristor for (a) a ThrEshold Adaptive Memristor (TEAM) model with a 0.2 volt sinusoidal input operating at a frequency of 2 GHz, and (b) resistive devices with and without ideal cross-coupled diodes. The parameters of the TEAM models are listed in Table II.  $V_{ON}$  is the on-voltage of the diode, and  $R_{ON}$  and  $R_{OFF}$  are, respectively, the minimum and maximum resistance of the memristor.

a negative differential resistance [23]. From a circuits perspective, the combined device can be modeled as a pair of cross coupled diodes in series with a memristor, as shown in the inset of Figure 2b. Since the rectifying structure requires an additional thin film layer, there is no effect on the area of the crosspoint structure.

An I-V curve of a memristive device with cross coupled diodes is shown in Figure 2b. The high resistance of the unselected devices reduces sneak currents and ensures that the leakage power of the array is relatively small. Reducing sneak currents ensures that the leakage power of the array is relatively small. A DC analysis of the on and off crosspoint currents is listed in Table I, where a 4 x 4 crosspoint array with RRAM devices is DC biased at 0.8 volts. These RRAM devices exhibit an on/off current ratio of 30. In an unrectified crosspoint, the observed current ratio drops to less than two. The rectified crosspoint displays a current ratio of 28.5, only 5% less than the ideal ratio of an RRAM device. Furthermore, the total power consumption is reduced by almost an order of magnitude.

TABLE I. COMPARISON OF DC ON/OFF CURRENT FOR 4 X 4 CROSSPOINT ARRAY

	$I_{on}$ [mA]	$I_{off}$ [mA]	Ratio	Average Active Power [mW]
Unrectified	2.3	0.132	1.7	1.45
Rectified	0.486	0.017	28.5	0.201

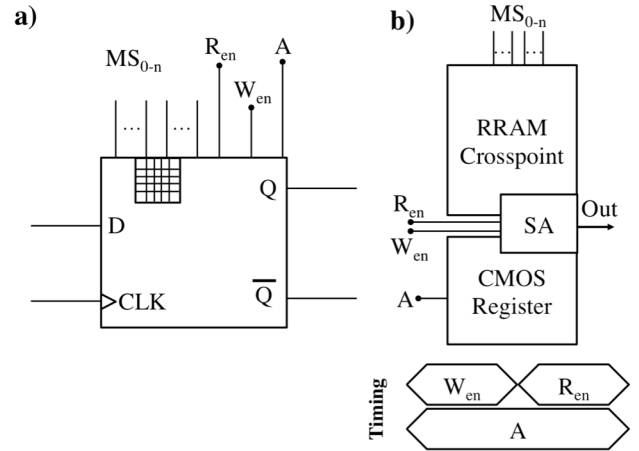


Figure 3. Multistate register element. (a) Symbol of the multistate register, and (b) block diagram with control signal timing. The symbol is similar to a standard CMOS D flip flop with the addition of a crosspoint array symbol.

III. RRAM MULTISTATE REGISTER

The multistate register is a novel circuit used to store multiple bits within a single logic gate. The multistate register is "drop-in" compatible with existing CMOS based flip flops. The element utilizes a clocked CMOS register augmented by additional sense circuitry (SA) and global memristor select (MS) lines. The symbol and topology of the multistate register are shown in Figure 3. Multistate registers can be used as pipeline registers within a processor pipeline, as shown in Figure 4 and further explained in Section V.

The MS lines select individual RRAM devices within the crosspoint memory co-located with the CMOS register. A schematic of the proposed RRAM multistate register is shown in Figure 5a. The signals  $W_{en}$  and  $R_{en}$  are global control signals that, respectively, write and read within the local crosspoint memory. Signal  $A$  sets the CMOS register into an intermediate state that facilitates writes and reads from the crosspoint. An individual RRAM device is selected using a set of global MS lines. Local writes to the RRAM crosspoint array are controlled by the master stage within the register. The gates within the slave stage of the CMOS register are reconfigured to provide a built-in sense amplifier to read the RRAM crosspoint array [24]. The overhead of the additional circuitry (shown in Figure 3) is relatively small (see Section IV.B).

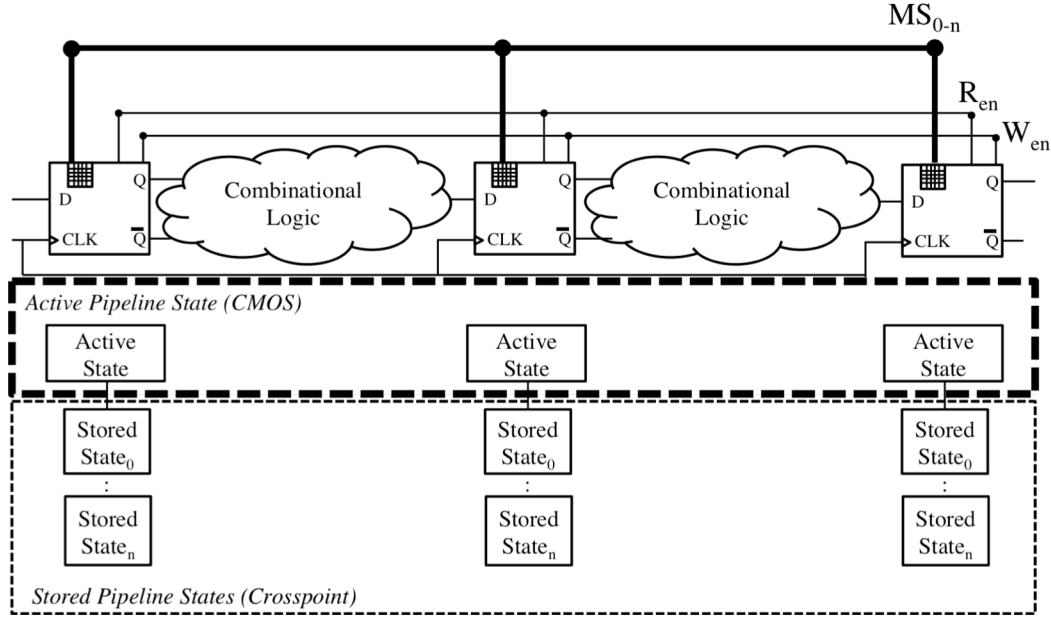


Figure 4. Multistate pipeline register (MPR) based pipeline and logic diagram of active and stored pipeline states. The MPR replaces a conventional pipeline register and time multiplexes the stored states.

The multistate register primarily operates as a CMOS register. In this mode, the structure behaves as a standard D flip flop, where a single bit is stored and is active while the idle states are stored within the RRAM crosspoint array. When global control circuitry triggers a change of the pipeline state (e.g., for a pipeline stall or context switch), the circuit stores the current bit of the register and reads out the value of the next active bit from the internal RRAM-based storage. Switching between active bits consists of two phases. In the first half of the cycle, an RRAM write operation stores the current state of the register. During a write operation, the transmission gate  $A$  disconnects the first stage from the following stage, isolating the structure into two latches. The input latch stores the currently evaluated state, while the output latch stores the data of the previous state. Once  $W_{en}$  goes high, the input latch drives a pair of multiplexers that write the currently stored state into the RRAM cell selected by the global MS lines. The active devices during the write phase are shown in Figure 5b. The write phase may require more than half a cycle depending upon the switching time of the RRAM technology. During the second half of the clock cycle, the new active bit is selected within the resistive crosspoint array and sensed by the output stage of the CMOS D flip flop. During a read operation, the globally selected row is grounded through the common node  $N_{in}$ . The voltage on the common line  $N_{out}$  is set by the state of the RRAM cell. To bias the RRAM cell, the common line is connected through a PMOS transistor to the supply voltage  $V_{DD}$ . The voltage is sensed at the output of  $M_1$ . If  $R_{en}$  is set high,  $M_1$  to  $M_5$  reconfigure the last inverter stage as a single ended sense amplifier [13], and the crosspoint array is read. The active devices during the read phase are shown in Figure 5c.

The physical design of the multistate register can be achieved by two approaches. RRAM devices can be integrated between the first two metals, as illustrated in Figure 6a, or the RRAM can be integrated on the middle level metal layers, as shown in Figure 6b. The middle metal layer approach allows the RRAM to be integrated above the CMOS circuitry, saving area. A standard cell floorplan is shown in Figure 7b, where a dedicated track is provided for the RRAM interface circuitry. This dedicated track runs parallel to the CMOS track. The addition of this track wastes area in those cases where multistate registers are sparsely located among the CMOS gates. Additional routing overhead increases the area required to pass signals around the crosspoint array.

The approach illustrated in Figure 7a, where the RRAM is integrated on the lower metal layers, requires slightly more area but is compatible with standard cell CMOS layout rules. Fabrication on the lower levels maintains standard routing conventions, where the lower metal layers are dedicated to routing within the gates, and the middle metal layers are used to route among the gates.

#### IV. SIMULATION SETUP AND CIRCUIT EVALUATION

The multistate register has been evaluated for use within a high performance microprocessor pipeline. The latency, energy, and area of the register are described in this section as well as the sensitivity to process variations.

##### A. Latency and energy

The latency and energy of an MPR are dependent on the parameters of an RRAM device and the CMOS sensing circuitry built into the MPR. The RRAM device is modeled using the TEAM model [21] based on the parameters listed in Table II. The parameters of the resistive device are chosen to

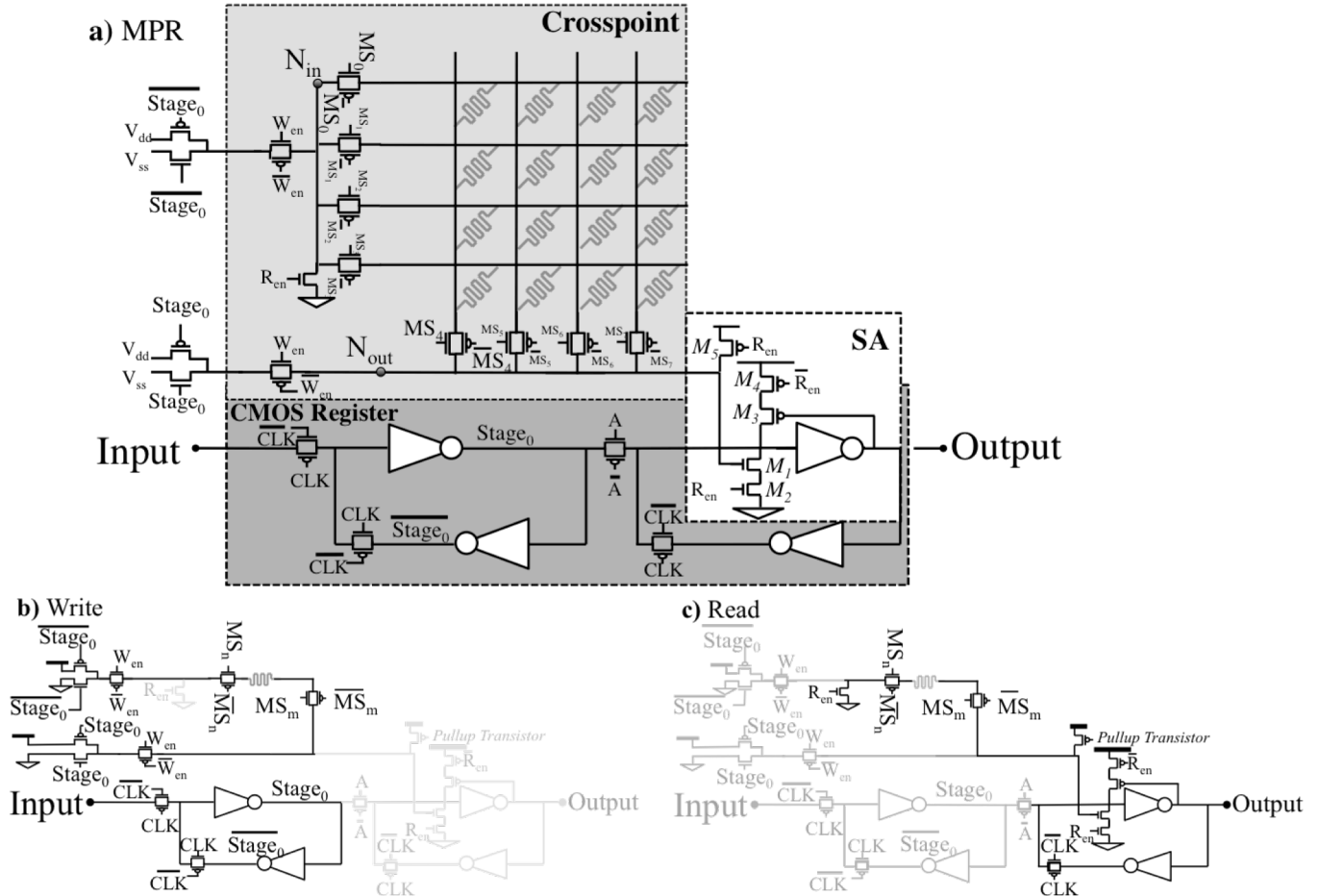


Figure 5. Proposed RRAM multistate pipeline register. (a) The complete circuit consists of a RRAM-based crosspoint array above a CMOS-based flip flop, where the second stage (the slave) also behaves as a sense amplifier. The (b) write and (c) read operations of the proposed circuit.

incorporate device nonlinearity into the I-V characteristic, as shown in Figure 2b and described in Section IIA. The multistate register is evaluated across a range of internal cross point array sizes (e.g., different number of states per register). The resistance of the device is extracted from [23]. The transistor and cell track sizing information is from the FREEPDK45 Standard Cell Library [25] and scaled to a 22 nm technology. Circuit simulations utilize the 22 nm PTM CMOS transistor model [26]. The RRAM and diode device parameters are listed in Table II. Standard CMOS timing information for the register is listed in Table III. The read operation requires 28.6 ps, equivalent to a 16 GHz clock frequency (the read operation is less than half a clock cycle). The register operates primarily as a CMOS register and only accesses the RRAM crosspoint array to switch between idle and active pipeline states. Note that the eight row by eight column crosspoint array is small as compared to large scale memory crosspoint arrays, and therefore places a small electrical load on the sensing circuitry. Hence, the read operation is relatively fast and does not limit the operation of the multistate register.

The performance of the multistate register is limited by the switching characteristics of the RRAM device. To maintain high performance, the desired RRAM devices must be relatively fast [31]. These characteristics are chosen to achieve a target write latency of a 3 GHz CPU. As mentioned in Section II, the RRAM write operation occurs sequentially prior to the read operation. Due to the sequential nature of the multistate register access to the RRAM array, a half cycle is devoted to the read operation.

The energy of the multistate register depends upon the RRAM switching latency, as listed in Table IV.  $E_{Low-High}$  and  $E_{High-Low}$  are the energy required to switch, respectively, to  $R_{off}$  and  $R_{on}$  for a single device write to the multistate register crosspoint array. Since the switching time of the memristor dominates the delay of a write to the multistate register,  $E_{Low-High}$  and  $E_{High-Low}$  increase linearly as the switching time increases. Note that the read energy only depends on  $R_{ON}$  and  $R_{OFF}$  and is therefore constant for different switching times. The read energy, however, depends on the size of the crosspoint array, as listed in Table V.

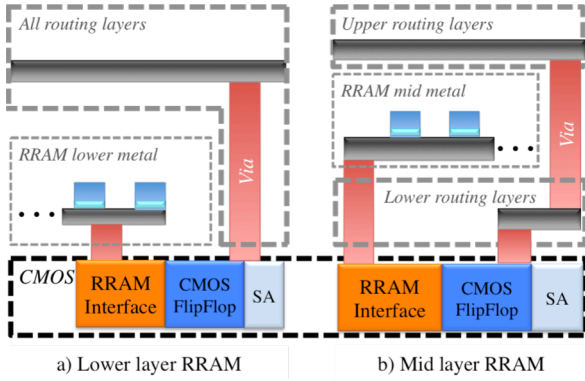


Figure 6. Vertical layout of RRAM in MPR circuit for a) lower level, and b) mid-layer crosspoint RRAM array.

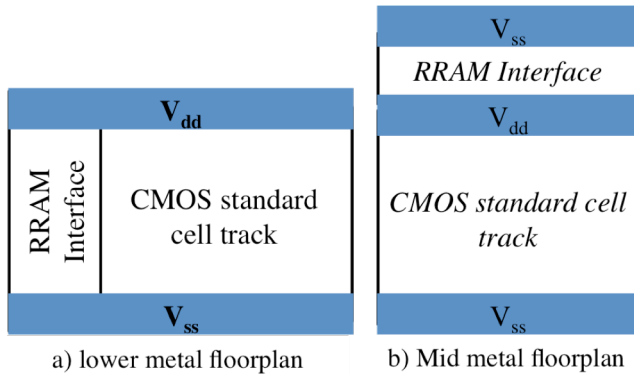


Figure 7. Planar floorplan of MPR with lower metal and upper metal RRAM layers. The RRAM array is not marked in this figure since it is located above the CMOS layer and has a smaller area footprint.

**B. Layout and physical area**

The energy and latency of an MPR are dependent on both the parameters of an RRAM device and on the CMOS sensing circuitry built into the MPR. An individual crosspoint RRAM cell is  $0.001934 \mu\text{m}^2$  ( $4F^2$ , where  $F$  is the feature size). The layout of the proposed RRAM multistate register is shown in Figure 8. The layout of the multistate register is based on 45 nm design rules and scaled to the target technology of 22 nm. The number of RRAM devices within a crosspoint array is scaled from four devices to 64 devices. The MPR is evaluated for both the middle metal and lower metal approaches, as described in Section III. The physical area is listed in Table VI.

The transistors required to access the crosspoint, as shown in Figure 8, dominate the area overhead of both the lower metal and middle metal multistate register. Due to the relatively small on-resistance of the RRAM devices, the access transistor needs to be sufficiently large to facilitate a write operation. Additionally, CMOS transmission gates are used to ensure that there is no threshold voltage drop across the pass transistors. As a result, the area of the crosspoint memory array is only a small fraction of the area overhead of

TABLE II. MEMRISTOR AND DIODE PARAMETERS

$R_{on}$ [k $\Omega$ ]	0.5
$R_{off}$ [k $\Omega$ ]	30
$k_{on}$	-0.021-0.07
$k_{off}$	0.0021-0.007
$\alpha_{on,off}$	3
$i_{on}$ [ $\mu\text{A}$ ]	-1
$i_{off}$ [ $\mu\text{A}$ ]	1
$V_{ON}$ (diode) [V]	0.5
$R_{out}$ (diode) [ $\Omega$ ]	1

TABLE III. ACCESS LATENCY OF A 16 BIT MPR

Clock to Q [ps]	11.2
Setup Time [ps]	13.2
RRAM Read [ps]	28.6

TABLE IV. WRITE LATENCY AND ENERGY OF A 16-BIT MULTISTATE REGISTER

Write Time [cycles @ 3 GHz]	0.5	1.5	2.5	3.5	4.5
$E_{Low-High}$ [fJ]	2.24	5.26	8.3	10.49	13.23
$E_{High-Low}$ [fJ]	3.78	10.33	16.89	23.5	30.08

TABLE V. READ ACCESS ENERGY OF RRAM

States per Multistate Register	4 States	16 States	64 States
$E_{read,Off}$ [fJ]	1.6	2.2	3.5
$E_{read,On}$ [fJ]	0.33	0.41	0.71

the multistate register. Note that alternative RRAM technologies with a higher  $R_{on}$  supports smaller transistors and reduced area. Under these constraints, the most area efficient structure is a 64 bit array, as the overhead per state is, respectively,  $0.08 \mu\text{m}^2$  for the lower metal approach and  $3.75 \mu\text{m}^2$  for the middle metal approach.

As shown, the middle metal register requires less area than a lower metal multistate register. As described in Section III and depicted in Figure 8b, the middle metal register requires an additional track dedicated to the control transistors within the crosspoint array. Positioning the crosspoint array over the register also adds complexity as the upper metal layers can no longer be used to route signals above the multistate register.

**C. Sensitivity and device variations**

The built-in sense amplifier circuit senses the RRAM based on a threshold voltage. Any voltage above the threshold of the registers produces a logical zero at the output, and any voltage below the threshold produces a logical one. Similar to digital CMOS circuits, the structure is tolerant to variability in the RRAM resistance. To evaluate the sensitivity of the circuit

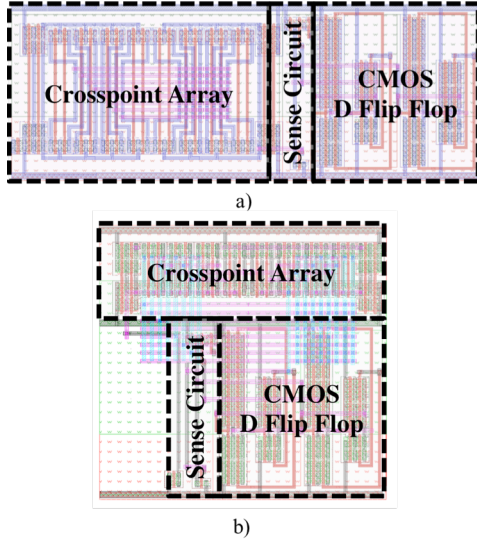


Figure 8. Physical layout of 64 state MPR within the crosspoint array on (a) lower metal layers (M1 and M2), and (b) upper metal layers (M2 and M3) above the D flip flop.

TABLE VI. MPR AREA

		Area [μm <sup>2</sup> ]	Overhead [%]	Overhead per State [%]
	CMOS Register (1 state)	2.8	-	-
Lower Metal	MPR 4 states	5.5	96.2%	24%
	MPR 16 states	6.3	126.5%	8%
	MPR 64 states	8.1	192.5%	3%
Middle Metal	MPR 4 states	3.9	39.3%	9.8%
	MPR 16 states	4.3	53.6%	3.3%
	MPR 64 states	5.2	85.7%	1.3%

to variations, the nominal  $R_{on}$  is varied from 0.35 to 0.65 kΩ. This range produces a maximum and minimum change of  $\pm 2$  mV in the voltage input of the sense amplifier. For  $21 \text{ k}\Omega > R_{off} < 39 \text{ k}\Omega$ , a voltage ranging from -40 mV to +26 mV is produced. Both ranges represent a 30% variation in the device resistance of  $R_{on}$  and  $R_{off}$ . In these cases, the correct output state is read out, indicating a high degree of tolerance to variations in the RRAM resistance.

The RRAM circuit can tolerate an  $R_{on}$  of up to 12 kΩ before the circuit produces an incorrect output. In a 64 bit multistate register, this behavior corresponds to an increase in the RRAM read delay from 78 ps to 476 ps. With increasing  $R_{on}$ , the sense amplifier no longer generates a full range signal at the output, dissipating static energy. Much of this increased delay is due to the device operating near the switching threshold of the sense amplifier.

As  $R_{off}$  varies from 30 kΩ to 300 MΩ, the performance of the circuit improves due to two effects. As the resistance increases, the voltage at the sense amplifier input also increases, placing the transistor into a higher bias state, which lowers the delay of the sense amplifier. Additionally, the large

TABLE VII. SOE MT AND CFMT PROCESSOR CONFIGURATIONS

	Switch on Event	RRAM-based CFMT
Number of pipeline stages	10	
CMOS process	22 nm	
Clock frequency [GHz]	3	
Switch penalty [cycles]	7	1 to 5
L1 read/write latency [cycles]	0	
L1 miss penalty [cycles]	200	
Data L1 cache configuration	32 kB, 4 way set associative	
Instruction L1 cache configuration	32 kB, 4 way set associative	
Branch predictor	Tournament, Ishare 18kB/gshare 8kB	

TABLE VIII. PERFORMANCE SPEEDUP FOR DIFFERENT MPR WRITE LATENCIES AS COMPARED TO SWITCH-ON-EVENT MULTITHREADING PROCESSOR FOR CPU SPEC 2006

Benchmark	MPR Write Latency [clock cycles]				
	1	2	3	4	5
<i>libquantum</i>	1.35	1.28	1.21	1.15	1.09
<i>bwaves</i>	1.22	1.15	1.08	1.04	1
<i>milc</i>	1.47	1.26	1.18	1.11	1.06
<i>zeusmp</i>	1.85	1.59	1.40	1.29	1.21
<i>gromacs</i>	1.53	1.32	1.21	1.17	1.14
<i>leslie3d</i>	1.67	1.48	1.33	1.22	1.15
<i>namd</i>	1.40	1.24	1.15	1.08	1.04
<i>soplex.pds-50</i>	1.35	1.28	1.21	1.16	1.1
<i>lbm</i>	1.5	1.31	1.2	1.12	1.08
<i>bzip2.combined</i>	1.13	1.1	1.08	1.05	1.03
<i>gcc.166</i>	1.35	1.28	1.21	1.15	1.09
<i>gobmk.trevorc</i>	1.3	1.24	1.19	1.14	1.09
<i>h264ref.foreman_baseline</i>	1.06	1.02	1	1	1
<i>GemsFDTD</i>	1.45	1.3	1.18	1.08	1.04
<i>hammer.nph3</i>	1.18	1.14	1.11	1.07	1.04
<i>soplex.ref</i>	1.7	1.42	1.29	1.19	1.1
<i>gcc.c-typeck</i>	1.33	1.26	1.21	1.15	1.1
<i>gobmk.trevord</i>	1.29	1.23	1.18	1.13	1.08
<b>Average</b>	1.40	1.27	1.19	1.13	1.08

TABLE IX. ENERGY AND AREA EVALUATION FOR CFMT TEST CASE

	Switch on Event	RRAM-based CFMT	Difference
Thread switch energy [pJ]	109.9	9.1 @ 1 cycle penalty	-91.7%
		19.1 @ 2 cycle penalty	-82.6%
		29.2 @ 3 cycle penalty	-73.4%
		38.4 @ 4 cycle penalty	-65.1%
		48.2 @ 5 cycle penalty	-56.1%
Processor area [mm <sup>2</sup> ]	123.276	126.426	2.55%

resistance of the sensed RRAM device prevents the sense line within the crosspoint array from dissipating charge, maintaining a high voltage at the input of the sense amplifier. Counterintuitively, this effect lowers the delay when  $R_{off}$  is greater than 30 MΩ. Due to the interplay of  $R_{on}$  and  $R_{off}$ , a

TABLE X. ENERGY PER INSTRUCTION FOR VARIOUS CPU SPEC 2006 BENCHMARK APPLICATIONS

Benchmark	SoE MT [pJ/inst.]	CFMT				
		RRAM MPR – various thread switch latencies				
		1 cycle [pJ/inst.]	2 cycles [pJ/inst.]	3 cycles [pJ/inst.]	4 cycles [pJ/inst.]	5 cycles [pJ/inst.]
<i>libquantum</i>	15.17	14.12	14.29	14.46	14.63	14.80
<i>bwaves</i>	19.63	18.83	19.03	19.25	19.42	19.42
<i>milc</i>	24.51	22.61	23.23	23.47	23.74	24.11
<i>zeusmp</i>	21.10	18.04	18.62	19.19	19.18	19.95
<i>gromacs</i>	30.16	27.94	28.62	29.05	29.23	29.34
<i>leslie3d</i>	27.27	24.72	25.20	25.68	26.08	26.39
<i>namd</i>	22.90	21.42	21.91	22.21	22.50	22.65
<i>soplex.pds-50</i>	17.62	16.52	16.71	16.88	17.03	17.20
<i>lbm</i>	22.54	20.29	20.90	21.36	21.76	21.94
<i>bzip2.combined</i>	21.86	21.44	21.51	21.65	21.65	21.72
<i>gcc.166</i>	19.37	18.32	18.49	18.66	18.83	19.01
<i>gobmk.trevorc</i>	23.05	22.15	22.28	22.71	22.56	22.71
<i>h264ref.foreman_baseline</i>	25.95	25.27	25.35	25.50	25.69	25.76
<i>GemsFDTD</i>	23.89	21.88	22.43	22.99	23.36	23.49
<i>hmmer.nph3</i>	24.27	23.65	23.75	23.84	23.84	24.04
<i>soplex.ref</i>	21.92	19.47	20.04	20.44	20.80	21.17
<i>gcc.c-typeck</i>	19.94	19.16	19.12	19.27	19.43	19.58
<i>gobmk.trevord</i>	22.73	21.71	21.87	22.40	22.25	22.40
Average	22.44	20.97	21.30	21.61	21.78	21.98

delay tradeoff exists between the average resistance of the RRAM technology and the resistive ratio of the device.

The gain and offset of the sense amplifier have a small effect on the circuit performance. A higher sense amplifier gain improves the tolerance of the sense circuit to variations of the RRAM device. An offset voltage shifts the reference threshold voltage, but must be comparable to the supply voltage ( $0.3V_{DD}$  or more) before the circuit performance is affected.

#### V. MULTISTATE REGISTERS AS MULTISTATE PIPELINE REGISTER FOR MULTITHREAD PROCESSORS – A TEST CASE

Replacing CMOS memory (e.g., register file and caches) with non-volatile memristors significantly reduces power consumption. Multithreaded machines can exploit the high density and CMOS compatibility of memristors to store the state of the in-flight instructions with fine granularity within a CPU. Hence, using memristive technology can dramatically increase the number of threads running within a single core. This approach is demonstrated in this test case, where RRAM multistate registers store the state of multiple threads within a CPU pipeline.

In continuous flow multithreading [13], the multistate registers are used as MPRs to store the state of multiple threads. A single thread is active within the pipeline and the instructions from the other threads are stored within the MPRs. The MPRs therefore eliminate the need to flush instructions within the pipeline, significantly improving the performance of the processor, as illustrated in Figure 9.

To exemplify this behavior, the performance and energy of a CFMT processor with the proposed RRAM-based MPRs have been evaluated [27]. To evaluate the performance, the GEM5 simulator [28] is extended to support CFMT. The

energy has been evaluated by the McPAT simulator [29]. The simulated processor is a ten stage single scalar ARM processor, where the execution stage operates at the eighth stage. The performance and energy of the CFMT processor are compared to a switch-on-event (SoE) multithreading processor [30], where a thread switch occurs for each long latency instruction (e.g., L1 cache miss, floating point instructions), causing the pipeline to flush. The characteristics of the evaluated processors are listed in Table VII. The energy is compared to a 16 thread processor (i.e., with an MPR storing 16 states) which is a sufficient number of threads to achieve the maximum performance for most benchmark applications.

The performance of the processors is measured by the average number of instructions per clock cycle (IPC), as listed in Table VIII. The average speedup in performance is 40%. A comparison of the thread switch energy is listed in Table IX. The average energy per instruction for various CPU SPEC 2006 benchmarks is listed in Table X, where the average reduction in energy is 6.5%. The area overhead for a 16 thread CFMT as compared to an SoE is approximately 2.5%, as listed in Table IX.

For the CFMT configuration described herein, the simulations show that 16 threads are sufficient to achieve the maximum performance for the vast majority of SPEC CPU 2006 benchmarks. Alternate configurations with many long latency events or different machines may benefit from additional states.

Physically, a linear increase in the number of rows and columns within a crosspoint array generates a quadratically increasing number of states and physical area, increasing the efficiency of the crosspoint array. A small increase in the number of rows and columns supports many more threads. However, as previously mentioned, 64 states is sufficient for most applications.



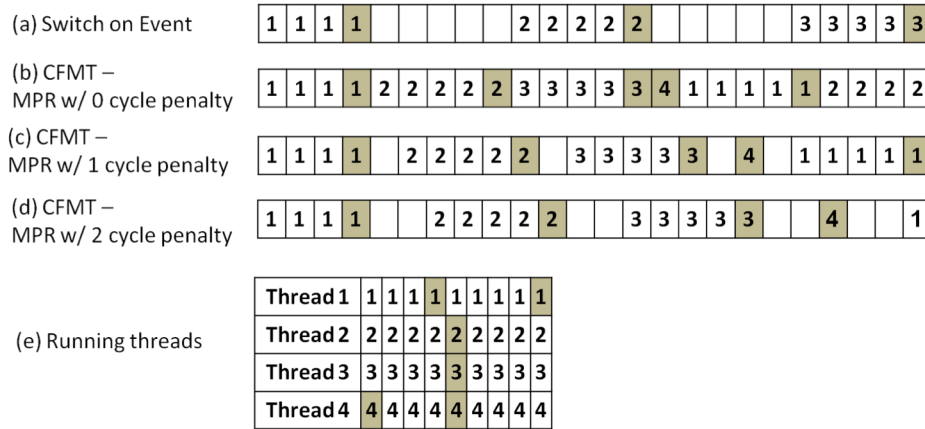


Figure 9. Illustration of different multithreading techniques with four threads (marked by the number). All processors run the same four threads as shown in (e). The latency of the 'white' and 'shaded' instructions is, respectively, a single clock cycle and ten clock cycles. For (a) switch-on-event multithreaded processor, the long latency instruction that triggers a thread switch is the shaded instruction and the thread switch penalty is five clock cycles due to the pipeline flush. For continuous flow multithreading, the thread switch penalty depends upon the read and write times of the multistate register, and is lower than traditional switch-on-event processors. The different thread switch penalties illustrated in this example are (b) zero for an ideal multistate register, (c) one clock cycle, and (d) two clock cycles. The performance (measured by the instructions per cycles) in this example is (a) 7/12, (b) 1 (71% improvement as compared to switch-on-event), (c) 0.83 (43% improvement), and (d) 0.67 (14% improvement) [12].

For the MPR to enhance performance, the cost of a thread switch must be smaller than the latency of a cache miss or other long latency events. This situation is typical for all practical thread switching events.

### VI. CONCLUSIONS

Emerging memory technologies, such as RRAM, are more than just a drop-in replacement to existing memory technologies. In this paper, a RRAM-based multistate register is proposed using an embedded array of memristive memory cells within a single flip flop. The multistate register can store additional data that is not conventionally contained within a computational pipeline.

The proposed multistate register is relatively fast due to the physical closeness of the CMOS and RRAM devices. A 16 state multistate register requires only 54% additional area as compared to a single state standard register. The multistate register is also relatively low power due to the non-volatility of the resistive devices.

As an example, the proposed multistate register has been applied to a continuous flow multithreaded processor, exhibiting a significant 40% performance improvement with low energy as compared to a conventional switch-on-event processor. An RRAM-based MPR therefore enables novel microarchitectures, such as the CFMT. The proposed multistate register significantly improves performance and reduces energy with small area overhead.

### ACKNOWLEDGMENTS

The authors thank Yoav Etsion, Yuval H. Nacson, and Uri C. Weiser for their contributions.

### REFERENCES

- [1] D. B. Strukov, G. S. Snider, D. R. Stewart, and R. S. Williams, "The Missing Memristor Found," *Nature*, Vol. 453, pp. 80-83, May 2008.
- [2] M. Hosomi et al., "A Novel Nonvolatile Memory with Spin Torque Transfer Magnetization Switching: Spin-RAM," *Proceedings of the IEEE International Electron Devices Meeting*, pp. 459-462, December 2005.
- [3] L. Chua, "Resistance Switching Memories are Memristors," *Applied Physics A*, Vol. 102, No. 4, pp. 765-783, March 2011.
- [4] E. Salman and E. G. Friedman, *High Performance Integrated Circuit Design*, McGraw-Hill Publishers, 2012.
- [5] H. S. Wong et al., "Metal-Oxide RRAM," *Proceedings of the IEEE*, Vol. 100, No. 6, pp. 1951-1970, June 2012.
- [6] Y. Ho, G. M. Huang, and P. Li. "Nonvolatile Memristor Memory: Device Characteristics and Design Implications," *Proceedings of the IEEE/ACM International Conference on Computer-Aided Design*, pp. 485-490, November 2009.
- [7] D. R. Lamb and P. C. Rundle, "A Non-Filamentary Switching Action in Thermally Grown Silicon Dioxide Films," *British Journal of Applied Physics*, Vol. 18, No.1, pp. 29-32, January 1967.
- [8] G. M. Ribeiro et al., "Designing Memristors: Physics, Materials Science and Engineering," *Proceedings of the IEEE International Symposium on Circuits and Systems*, pp. 2513-2516, May 2012.
- [9] J. J. Yang et al., "Memristive Switching Mechanism for Metal/Oxide/Metal Nanodevices," *Nature Nanotechnology*, Vol. 3, No. 7, pp. 429-433, June 2008.
- [10] J. Li and J. F. Martinez, "Power-Performance Implications of Thread-level Parallelism on Chip Multiprocessors," *Proceedings of the IEEE International Symposium on Performance Analysis of Systems and Software*, pp.124-134, March 2005.
- [11] D. M. Tullsen, S. J. Eggers, and H. M. Levy. "Simultaneous Multithreading: Maximizing On-Chip Parallelism," *Proceedings of the IEEE/ACM International Symposium on Computer Architecture*, pp. 392-403, May 1995.
- [12] Intel Ivy Bridge Specifications (two threads per core) <http://ark.intel.com/>
- [13] S. Kvatinsky, Y. H. Nacson, Y. Etsion, E. G. Friedman, A. Kolodny, and U. C. Weiser, "Memristor-Based Multithreading," *IEEE Computer Architecture Letters*, 2013 (in press).

- [14] L. O. Chua, "Memristor – The Missing Circuit Element," *IEEE Transactions on Circuit Theory*, Vol. 18, No. 5, pp. 507-519, September 1971.
- [15] L. O. Chua and S. M. Kang, "Memristive Devices and Systems," *Proceedings of the IEEE*, Vol. 64, No. 2, pp. 209- 223, February 1976.
- [16] T. Prodromakis, K. Michelakisy, and C. Toumazou. "Fabrication and Electrical Characteristics of Memristors with  $\text{TiO}_2/\text{TiO}_{2+x}$  Active Layers," *Proceedings of the IEEE International Symposium on Circuits and Systems*, pp.1520-1522, May 2010.
- [17] Z. Biolek, D. Biolek, and V. Biolkova., "SPICE Model of Memristor with Nonlinear Dopant Drift," *Radioengineering*, Vol. 18, No .2, pp. 210-214, June 2009.
- [18] G. M. Ribeiro *et al.*, "Designing Memristors: Physics, Materials Science and Engineering," *Proceedings of the IEEE International Symposium on Circuits and Systems*, pp. 2513–2516, May 2012.
- [19] H.Y. Lee *et al.*, "Low Power and High Speed Bipolar Switching with a Thin Reactive Ti Buffer Layer in Robust HfO<sub>2</sub> Based RRAM," *Proceedings of the IEEE International Electron Devices Meeting*, pp.1–4, December 2008.
- [20] Y. F. Chang *et al.*, "Study of SiO<sub>x</sub>-Based Complementary Resistive Switching Memristor," *Proceedings of the Annual Device Research Conference*, pp. 49–50, June 2012.
- [21] S. Kvatinsky, E. G. Friedman, A. Kolodny, and U. C. Weiser, "TEAM: ThrEshold Adaptive Memristor Model," *IEEE Transactions on Circuits and Systems I: Regular Papers*, Vol. 60, No. 1, pp. 211–221, January 2013.
- [22] Y. N. Joglekar, and S. J. Wolf, "The Elusive Memristor: Properties of Basic Electrical Circuits," *European Journal of Physics*, Vol. 30, No. 4, pp. 661–675, July 2009.
- [23] J. J. Yang *et al.*, "Engineering Nonlinearity into Memristors for Passive Crossbar Applications," *Applied Physics Letters*, Vol. 100, No. 11, pp. 113501–113501, March 2012.
- [24] K. Pagiamtzis and A. Sheikholeslami, "Content-Addressable Memory (CAM) Circuits and Architectures: A Tutorial and Survey," *IEEE Journal of Solid-State Circuits*, Vol. 41, No. 3, pp. 712–727, March 2006.
- [25] FreePDK45 User Guide, April 2011, <http://www.eda.ncsu.edu/wiki/FreePDK45>.
- [26] W. Zhao and Y. Cao, "New Generation of Predictive Technology Model for Sub-45 nm Early Design Exploration," *IEEE Transactions on Electron Devices*, Vol. 53, No. 11, pp. 2816–2823, January 2006.
- [27] S. Kvatinsky, Y. H. Nacson, R. Patel, Y. Etsion, E. G. Friedman, A. Kolodny, and U. C. Weiser, "Multithreading with Emerging Technologies – Dense Integration of Memory within Logic," (in submission).
- [28] The gem5 Simulator System, May 2012, <http://www.m5sim.org/>.
- [29] S. Li *et al.*, "McPAT: an Integrated Power, Area, and Timing Modeling Framework for Multicore and Manycore Architectures," *Proceedings of the IEEE/ACM International Symposium on Microarchitecture*, pp. 469-480, December 2009.
- [30] S. Kvatinsky, E. G. Friedman, A. Kolodny, and U. C. Weiser, "The Desired Memristor for Circuit Designers," *IEEE Circuits and Systems Magazine*, Vol. 13, No. 2, pp. 17-22, Second Quarter 2013.
- [31] M. K. Farrens and A. R. Pleszkun, "Strategies for Achieving Improved Processor Throughput," *Proceedings of the ACM International Symposium on Computer Architecture*, pp. 362-369, May 1991.