

# A Design Framework for Reliable and Energy Efficient Cross-Point Resistive Memory

**Abstract**—With conventional memory technologies approaching their scaling limit, emerging non-volatile memory technologies have attracted considerable attention because of their non-volatility, high access speed, low power consumption, and good scalability. Resistive RAM (ReRAM), with its simple structure, small cell size ( $4F^2$ ), and support for 3D stacking, has been a leading candidate among emerging technologies. A key advantage of ReRAM comes from its non-linear nature, which enable us to build a cross-point RAM array without having a dedicated access transistor in each cell. While cross-point design is effective in improving memory density, it has inherent disadvantages which introduce extra design challenges.

In this work, we analyze these challenges. Based on the circuit characteristics of the cross-point array, we propose a mathematical model to perform a comprehensive analysis of issues of reliability, energy consumption and the area overhead. In addition to the cell-level analysis, different programming schemes are also discussed in detail. The proposed model can enable designers to identify the most energy/area efficient ReRAM organization that meets specific design goals early in the design stage.

## I. INTRODUCTION

The scaling of traditional memory technologies, such as DRAM and FLASH, is approaching its physical limit. In the past few years, emerging non-volatile technologies (NVM), such as Phase Change RAM (PCRAM), Magnetoresistive RAM (STT-RAM), and Resistive RAM (ReRAM) have been widely studied as potential candidates for the next generation memory technologies to meet the need of higher density, faster access time, and lower power consumption. Among all of these emerging memory technologies, ReRAM has many unique characteristics, including simple structure, non-linearity and high resistance ratio, making itself one of the most promising technologies. Researchers have shown that the state-of-the-art single-level-cell ReRAM can achieve sub-8ns random access time for both read and write operations with a resistance ratio larger than 100 [?]. Also, HP labs and Hynix have already announced plans to commercialize the memristor-based ReRAM and predicted that ReRAM could eventually replace traditional memory technologies [?].

Unlike other non-volatile memory technologies, ReRAM can be implemented in a cross-point style structure without any access devices. Specifically, in a nano cross-point array, each bistable ReRAM cell is sandwiched by two orthogonal nanowires, without access devices. Thus the area occupied by each cell is literally the area underneath the intersection of wires, which is  $4F^2$  per bit. However, the simplicity of access device free, cross-point structure introduces challenges to the peripheral circuit design and memory organization.

While there has been prior studies on cross-point ReRAM array [?], [?], [?], [?], they do not consider the effect of voltage drivers and programming methods to cells. In addition, detailed area and energy analysis is also absent. In this work, we address the design challenges of cross-point structure based ReRAM. We build an accurate mathematical model to evaluate memory reliability, energy consumption, and area overhead for different designs and cell parameters. Based on this study, we propose a detailed design methodology which allows for exploring the most energy/area efficient ReRAM design with

TABLE I  
COMPARISON OF EMERGING NON-VOLATILE MEMORY TECHNOLOGIES

Metric	STT-RAM	PCM	FeRAM	ReRAM
Cell Size( $F^2$ )	6 – 20	4 – 8	15	4
Read Latency(ns)	1-10	20-50	20-80	5-50
Write Latency(ns)	2-20	150	100	5-50
Endurance	$10^{15}$	$10^8$	$10^{12}$	$10^{8-10}$

different design constraints and cell parameters at the very beginning of the design stage. On the other hand, the system designers can also leverage the proposed framework to provide valuable feedback to device researchers who will in turn adjust ReRAM cell design. We believe that this kind of collaboration will be very helpful to shorten time-to-market of ReRAM memory.

The rest of this paper is organized as follows. In Section II, the preliminaries of ReRAM technology and cross-point architecture are introduced. Section III discusses the mathematical model we propose for crossbar structure ReRAM and the edge conditions for different write and read schemes. Section IV analyzes different design constraints of write and read operations on the cross-point based ReRAM array. The energy consumption and area overheads are also analyzed in this section. Then in Section V, the effect of non-linearity and write current on the design constraints are evaluated. Finally, the conclusion is presented in Section VI.

## II. PRELIMINARIES

This section provides background of ReRAM technology and cross-point architecture, and discusses their advantages and limitations.

### A. Background of ReRAM technology

Table I compares the properties of state-of-art non-volatile memory technologies. ReRAM and STT-RAM are better than PCM or FeRAM due to their faster access time and high endurance. Although STT-RAM has better read/write delay characteristics, the difference in cell resistance between ON and OFF states is higher in ReRAM. Further more, ReRAM has the best memory density. Hence, ReRAM is a leading candidate for next generation storage.

As implied by the name, a ReRAM cell uses its resistance to represent the stored information. A ReRAM cell can be switched between high resistance state (HRS) and low resistance state (LRS) by applying an external voltage across the cell. In general, a ReRAM cell is built on a Metal-Insulator-Metal (MIM) structure. The resistance switching behaviors have been observed in many MIM nanodevices with different metal oxide materials. For example, a particular  $TiO_2$  based MIM structure ReRAM, named ‘memristor’, was developed by HP Labs in 2008 [?]. The proposed memristor-based ReRAM is considered as the first experimental realization and a theoretical model of the fourth fundamental circuit element, which is predicted by Chua [?] about 40 years ago. The memristor-based ReRAM has very small cell size with an access time of less than 50ns. Another

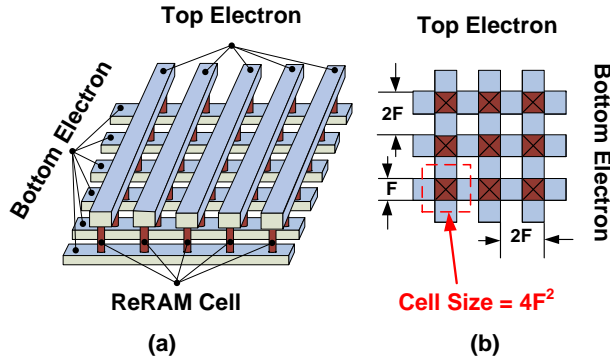


Fig. 1. A schematic view of typical cross-point architecture. (a): Overview of the cross-point architecture. (b): The layout of the cell size of  $4F^2$ .

$HfO_2$ -based bipolar ReRAM was implemented by ITRI this year with an access time as low as 7.2ns [?].

Although there are several variants of ReRAMs, all of them can be classified into two broad categories: unipolar ReRAM and bipolar ReRAM. In an unipolar cell, the resistance switching behaviors do not depend on the polarity of the voltage input across the cell and only relate to magnitude and duration of the voltage input. On the other hand, in a bipolar cell, the voltage polarity for ON-to-OFF switching (RESET operation) is different from OFF-to-ON switching (SET operation). The need for different pulse widths for SET and RESET in unipolar ReRAM means that its write latency will be determined by the longest pulse. Moreover, the control of SET, RESET, and read operations without any disturbance is another crucial design challenge, especially in high speed ReRAM design. For these reasons, most high performance ReRAM studies are dominated by bipolar ReRAM [?], [?], [?]. In this study, we perform a detailed analysis of the design challenges of bipolar ReRAM cross-point arrays.

### B. Cross-Point Architecture

There are two possible memory structures for a bipolar ReRAM implementation: a traditional MOSFET-accessed structure and a cross-point structure. In the MOS-accessed memory array, a MOSFET is added as an access device for each memory cell. As the size of a MOSFET access device is typically much larger than the size of a ReRAM cell, the total area of memory array is primarily dominated by MOSFETs rather than ReRAM cells. Also, in order to prove enough driven current, larger than minimum transistor should be used for writes. Hence, ReRAM's area advantage gets lost because of the access device.

In contrast, the cross-point structure is more area-efficient for the ReRAM memory array [?]. A schematic view of a typical cross-point memory array is shown in Figure 1(a). In a cross-point array, ReRAM is sandwiched between wordlines and bitlines. Figure 1(b) shows the feasibility of  $4F^2$  cells, the theoretical minimum size for a single layer single level memory cell. The memory density can further be improved by using a multi-layer multi-level cross-point ReRAM array [?] [?].

In a cross-point structure, the write operation can either write one bit per access or several bits attached to a wordline at the same time. Although the second scheme has higher bandwidth, it requires a two-step writing operation to prevent unintentional writing [?], significantly increasing the write latency. While writing to a cross-point array, the unselected wordlines and bitlines can either be left floating or half biased. However, while reading a cell, the selected wordline should be biased at read voltage and all other wordlines and bitlines in the array are shunted to ground. The current in each bitline

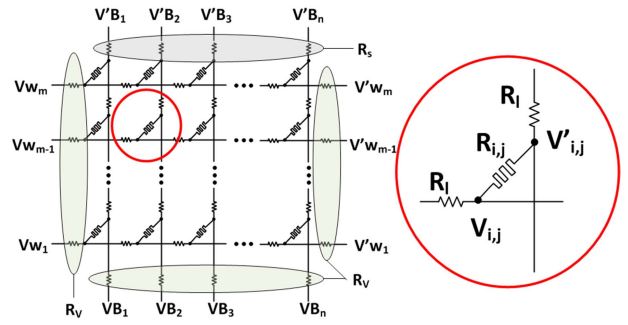


Fig. 2. The basic model of typical cross-point array.

is then sensed and compared to a reference current to determine cell content. However, due to the sneak current existing at the cross-point array, the current in bitline also varies depending upon the values stored in unselected cells. This phenomenon of read disturbance restricts the size of a cross-point array, since sneak current increases with the number of cells attached to wordlines and bitlines. In order to mitigate this problem, a two-step write operation was proposed: in the first step, the background current of the cross-point array will be sensed. Then the total current, comprising of both the background current and current at the selected cell, will be read out. The state of the selected cell can then be determined by computing the difference between the total current and background current. For this scheme to work, the cross-point array should be sized such that the difference between background current and current through the selected cell should be large enough for reliable sensing.

Depending upon the read/write scheme, the size of the array can vary significantly. In this work, we propose a methodology to find the minimum array size that meets specific energy and area constraints based on the worst-case state of the array. This will help designers find an optimal memory organization early in the design flow.

## III. MODELING OF THE CROSS-POINT MEMORY

In this section, we present a detailed mathematical model for cross-point arrays. By using this model, along with specific parameters and edge conditions, the reliability, energy consumption, and area overheads of different read/write schemes can be easily evaluated.

### A. Basic model of Cross-Point Memory

Figure 2 shows the circuit model of an  $M$  by  $N$  cross-point ReRAM array. The model is built upon Kirchhoff's Current Law (KCL) and it's validity can be guaranteed by deductions from the basic circuit theory. The horizontal lines are wordlines and vertical lines represent bitlines. The ReRAM cells are located at each cross-point of wordline and bitlines. The resistance of the ReRAM cell at the cross-point of  $i^{th}$  wordline and  $j^{th}$  bitline is represented by  $R_{i,j}$ . We assume the resistance of the wire connecting two cross-points to be  $R_{line}$ . The input resistance of each wordline and bitline is  $R_v$  and the resistance of sense amplifier is  $R_s$ . In order to set up the KCL equations, the voltage at each cross-point is indicated as  $V_{i,j}$  for wordline and  $V'_{i,j}$  for bitline. A detailed cross-point is also shown in right side of Figure 2. The input voltage for the  $i^{th}$  wordline is  $V_{Wi}$  and the  $i^{th}$  bitline is  $V_{Bi}$ . In the case where a wordline takes input from both the sides, the voltage at the other end of the  $i^{th}$  wordline is represented as  $V_{W1}$ . Finally, the voltage at the sense amplifier is  $V'_{Bi}$  during the read operation.

### B. Mathematical Model of a Cross-Point Array

Based on the circuit model shown in Figure 2, the current equations for each cross-point can be set following KCL:

$$\sum_{I=1}^k I_k = 0. \quad (1)$$

All of the cross-points have similar structure with no more than three current branches and therefore it is very easy to set up the KCL equations for each cross-point. However, we should treat the cross-points at the edges of the array specially because there are different conditions for different write/read schemes. For example, the unselected wordline for write operation can be either half biased or left floating. Thus, the edge conditions should be adjusted according to each write/read scheme. In particular, all of the cross-points in an array can be classified into three major categories: *normal point*, *activated point* and *floating point*.

The normal points are located inside the memory array. In other words, for all of the nodes with  $1 < i < m$  and  $1 < j < n$ , the KCL equations take the form of

$$R_l^{-1}V_{i,j-1} - (2R_l^{-1} + R_{i,j}^{-1})V_{i,j} + R_l^{-1}V_{i,j+1} + R_{i,j}^{-1}V'_{i,j} = 0, \quad (2)$$

for the node at wordline layer and

$$R_l^{-1}V'_{i-1,j} - (2R_l^{-1} + R_{i,j}^{-1})V'_{i,j} + R_l^{-1}V'_{i+1,j} + R_{i,j}^{-1}V_{i,j} = 0, \quad (3)$$

for the node at bitline layer.

The activated point and floating point represent the nodes at the edge of cross-point array with different conditions: an edge point, which is directly connected to the voltage input or to the ground, can be considered as an activated point. Otherwise, it is a floating point. For example, consider the point located at the intersection of  $i^{th}$  wordline and  $1^{st}$  bitline. If the  $i^{th}$  wordline is activated by an input voltage of  $V_{Wi}$ , this cross-point is an activated point, and the KCL equation for this point is:

$$-(R_v^{-1} + R_l^{-1} + R_{i,1}^{-1})V_{i,1} + R_l^{-1}V_{i,2} + R_{i,1}^{-1}V'_{i,1} = -R_v^{-1}V_{Wi}. \quad (4)$$

Otherwise, it is floating and its KCL equation is

$$-(R_l^{-1} + R_{i,1}^{-1})V_{i,1} + R_l^{-1}V_{i,2} + R_{i,1}^{-1}V'_{i,1} = 0. \quad (5)$$

For clarity, a  $2mn \times 1$  vector  $V$  is defined to represent all of the variables in the KCL equations:

$$V = [V_1^T, V_2^T \dots V_m^T, V'_1{}^T, V'_2{}^T \dots V'_m{}^T]^T, \quad (6)$$

where,

$$V_i = [V_{i,1}, V_{i,2} \dots V_{i,n}]^T, \quad V'_i = [V'_{i,1}, V'_{i,2} \dots V'_{i,n}]^T, \quad (7)$$

for  $i = 1, 2 \dots m$ . Then all of the KCL equations can be considered as a system of linear equations, which has the form,

$$A \cdot V = C. \quad (8)$$

$A$  is a  $2mn \times 2mn$  coefficient matrix, which is determined by Equations(2)-(5).  $C$  is a  $2mn \times 1$  vector, containing the constant terms of these equations. As shown, the KCL equations for each node have very simple structure and are very similar to each other. Therefore, the linear equation system has a relatively fixed format and simple structure, making it very easy to establish and adjust the coefficients and constants according to different design schemes. The characteristics of the linear system can be summarized as:

- 1) As shown in Equation (9), the coefficient matrix  $A$  can be further partitioned into 4 smaller subblocks :

$$A = \begin{bmatrix} A1 & A2 \\ A3 & A4 \end{bmatrix}. \quad (9)$$

All of these subblocks have the same size of  $m \times n$ . Subblock  $A2$  and  $A3$  are diagonal matrixes and have the value of:  $A2_{i,i} = A3_{i,i} = R_{i,i}^{-1}$ .  $A2$  and  $A3$  do not change their values with different schemas. However,  $A1$  and  $A4$  are a little more complex than  $A2$  and  $A3$ .  $A1$  is a tridiagonal matrix and has nonzero elements only in the main diagonal, and the first line below and above the diagonal. Similarly,  $A4$  is a special tridiagonal matrix, which has nonzero elements in the main diagonal, and the  $n^{th}$  line below and above the diagonal, where  $n$  is the number of bitline in the cross-point model. The value of the elements in  $A1$  and  $A4$  can be easily derived from Equation (2) and (3). However, the edge condition varies with different program schemes. Therefore, the coefficients related to the edge condition should be set according to the program schemes. Clearly, the four edges shown in Figure 2 correspond to different coefficients in  $A1$  and  $A4$ . Due to the space limitations, we consider the nodes at the left edge of the array as an example. A similar procedure can be followed to initiate the coefficients of other edge. The coefficients of nodes at the left edge of the array ( $V_{i,1}$ ) can be set as:

$$A1(k, k) = \begin{cases} -(R_l^{-1} + R_{i,1}^{-1}) & \text{if floating} \\ -(R_v^{-1} + R_l^{-1} + R_{i,1}^{-1}) & \text{if activated} \end{cases} \quad (10)$$

where  $k = (n-1)i + 1$  for  $i = 1, 2 \dots m$ .

- 2) The constant terms  $C$  is a  $2mn \times 1$  vector. Equation(2)-(5) show that only KCL equations of the activated points have constant terms. Therefore, only the following elements in  $C$  may have non-zero value:  $C((i-1)n+1)$ ,  $C(in)$ ,  $C(mn+i)$  and  $C((2m-1)n+i)$  for  $i = 1, 2 \dots m$ , corresponding to the nodes at the four edges respectively. Likewise, as an example, we consider nodes  $V_{i,1}$ . The constant corresponding to these nodes can be defined as:

$$C((i-1)n+1) = \begin{cases} 0 & \text{if floating} \\ -R_v^{-1}V_{Wi} & \text{if activated} \end{cases} \quad (11)$$

Thus, with input parameters such as the resistance of a ReRAM cell, the resistance of interconnect wires, program voltages, and write/read schemes, voltages at various cross points can be obtained by solving Equation (8). With detailed voltage values,  $V_{2mn \times 1}$ , we can analyze the array at a fine granularity. These values are also critical to evaluate reliability, energy consumption, driven current density, and area overheads of a cross-point array.

### IV. ANALYSIS OF DESIGN CONSTRAINTS - A CASE STUDY

In this section, a typical ReRAM cross-point array is analyzed in detail by using our mathematical model.

#### A. Overview

As shown in Figure 2, in order to write or read a cross-point array, proper voltages should be applied to wordlines and bitlines. Since several potential read/write schemes can be used to program a memory array, it is difficult to identify the ideal scheme that meets the design constraints in terms of area, energy, and reliability. In this section, we study the effect of various schemes on cross-point size and reliability. The constraints on array size, energy and area overheads are analyzed in the worst cases scenario. The results of this study will be a useful guide in designing a cross-point array.

Table II shows the circuit parameters of our baseline 32nm design. The data is derived from the recently published studies

on ReRAM [?], [?]. We study reliability, energy consumption, and area overheads for four different write schemes, and discuss the sensitivities of these schemes to the data pattern of HRS and LRS ReRAM cells and cell non-linearity.

Although the goal of a read operation is different from a write operations, both of them are realized by fully biasing the selected cell and floating (or half biasing) unselected cells. Thus, the coefficient matrix  $A$  and the constant vector  $C$  are very similar for both. In addition, their energy consumption and area overhead will also have a similar trend. In the next section, we first study the writing operation comprehensively. After that, for read operation, we mainly focus on the read margin analysis since it is unique to reads.

### B. Write Operation

To write a ReRAM cell, an external voltage is applied across the cell for a certain duration. Intuitively, there are four possible schemes for the write operation:

- 1) According to the location of a selected cell, activate one wordline and one bitline and leave all of other lines floating (FWFB schemes).
- 2) Activate the selected wordline and bitline. Leave all the unselected wordlines floating and half bias the unselected bitlines (FWHB schemes).
- 3) In contrast to the scheme 2), activate the selected wordline and bitline. Leave all the unselected bitlines floating and half bias the unselected word lines (HWFB schemes).
- 4) Activate the selected wordline and bitline. Then half bias all other wordlines and bitlines (HWHB schemes).

Since reliability, energy consumption, and area overheads for these schemes are different, we address these problems separately and finally combine all constraints to provide design guidelines for write operations.

#### Reliable Write Operation.

Write reliability is a serious concern in cross-point arrays. In an ideal condition, the resistance of wires and the sneak currents in unselected cells are negligible. In such a scenario, all the write schemes discussed above will make sure that the write voltage  $V_W(W) - V_B(W)$  is fully applied across the specified cell. However, in reality, both wire resistance and sneak current are non-trivial. Hence, the operation of cross-point array will vary based on the data pattern stored in ReRAM cells. A write is considered reliable if it modifies the content of the selected cells to the new value without disturbing other unselected cells. There are two potential problems with writes: *write failure*, an unsuccessful write on selected cell, and *write disturbance*, an undesirable write on unselected cell. It is

necessary to ensure that a write scheme guarantees reliable operation even in the worst case (w.r.t the location of cells to written and pattern stored in the cross-point array). Otherwise, after several unreliable write operations, the data stored in the cross-point array will become unpredictable.

We first use an example to show the problem with FWFB scheme, which may result in severe write disturbance. Figure 3 shows the voltage drop across each ReRAM cell of a  $64 \times 64$  cross-point array. In this example, in order to write the cell at the cross point of the  $32^{nd}$  wordline and the  $32^{nd}$  bitline, the selected wordline and bitline are biased at 2V and 0V, respectively. All of the other wordlines and bitlines are biased at 1V. The ReRAM cells at the selected bitline are in the HRS, while all of the other cells are in the LRS. It is clear that the voltage drop across the selected cell ( $V_{32,32}$ ) almost has the same magnitude as the unselected cell at the same bitline, resulting the write disturbance to all of the unselected cells at the selected bitline. Actually, for a  $M \times N$  matrix ( $M > N$ ), the worst case voltage drop of the unselected cell can be calculated as:

$$V_{worst} = V_{select} \cdot \left[ 1 - \frac{1}{M + (N - 1)R_{off}/R_{on}} \right]. \quad (12)$$

Considering that the reported On-OFF resistance ratio of ReRAM cell is always  $> 50$  [?], [?], [?], [?], [?], [?], , the worst case voltage drop at the unselected cell is larger than 98% of the voltage at the selected cell, making it is impossible to build a reliable cross-point structure ReRAM with the FWFB scheme. Therefore, in the following discussion, we only compare the results of FWHB, HWFB and HWHB schemes. For each of these three schemes, we can either write the cells at one wordline at the same time or only write one bit per access and separate the write operation to several arrays. In the following discussion, we start from one bit per access write operation, then the results of one wordline per access method are discussed.

Write failure typically results from the voltage drop at the interconnect wires along the wordline and bitline. It has been shown that [?], for one bit per access write operation, the worst case voltage drop occurs when

$$\begin{cases} R_{M,N} = R_{on} \\ V_{WM} = V_W(W) \\ V_{BN} = V_B(W). \end{cases} \quad (13)$$

In order to avoid the write failure and successfully program the selected ReRAM cell, the driven voltage should be boosted to a higher level, making sure that the voltage across the cell exceeds the threshold voltage even at the worst case. Figure 4 shows the lower bound of the driven voltage for different sizes of cross-point array. The minimum word/bitline voltage increases from 2.01 V for

TABLE II  
PARAMETERS OF THE BASELINE CROSS-POINT ARRAY

Metric	Description	Values
$S_{cell}$	Cell Size	$4F^2$
$R_l$	Interconnection Resistance	$1.25\Omega$
$V_{RESET}$	Threshold voltage for RESET	2.0V
$V_{SET}$	Threshold voltage for SET	-2.0V
$V_{READ}$	Read Voltage of Cell	0.5V
$R_{off}$	HRS Resistance	$500K\Omega$
$R_{on}$	LRS Resistance	$10K\Omega$
$V_W(R)$	Wordline Voltage during Read	0.4V
$V_W(W)$	Wordline Voltage during Write	$\pm 2V$
$V_W(H)$	Half Selected wordline Voltage	1V
$V_B(R)$	Bitline Voltage during Read	0V
$V_B(W)$	Bitline Voltage during Write	0V
$V_B(H)$	Half Selected bitline Voltage	1V
$M$	Number of wordlines	64
$N$	Number of bitlines	64

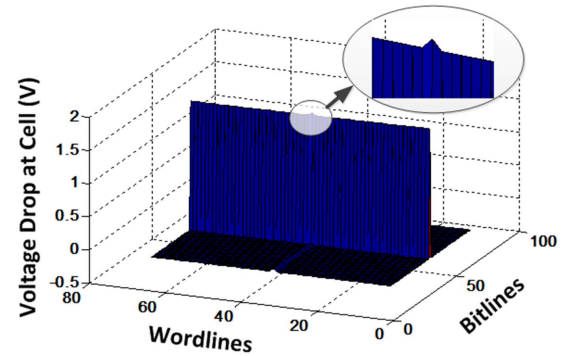


Fig. 3. Write disturbance for FWFB schemes. ( $V_{W32} = 2V$ ,  $V_{B32} = 0V$ .  $R_{x,32}$  at HRS, others at LRS.)

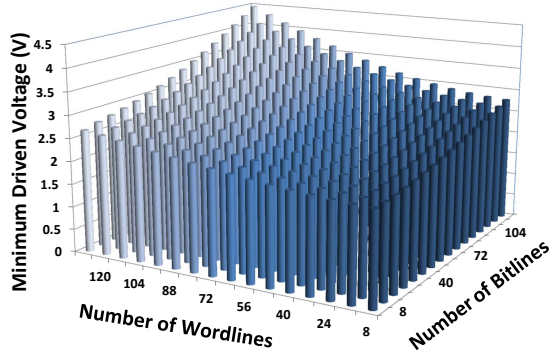


Fig. 4. Write voltage requirement (Threshold voltage = 2V).

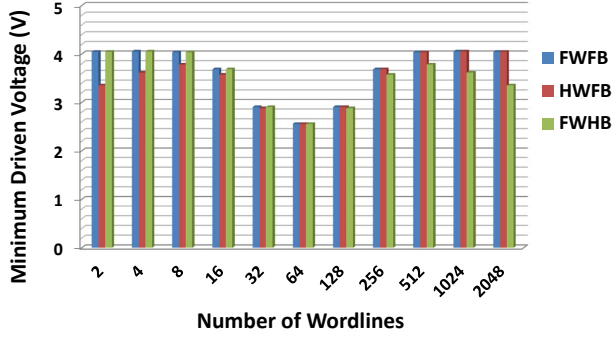


Fig. 5. Write voltage requirement with different memory shape. (Array capacity = 4Kbits, Activated wordline voltage = 2V, Activated bitline voltage = 0V.)

a  $8 \times 8$  array to 4.47 V for a  $128 \times 128$  cross-point array. In addition, for a memory capability, the cross-point array can be organized with different number of wordlines and bitlines. For example, a 4K bits cross-point array can be implemented either by a  $64 \times 64$  array or by a  $32 \times 128$  array. In the latter case, the voltage drops along the wordline will be much more serious than along the bitline. Figure 5 examines the voltage requirement for different array organizations with different write schemes. The result shows that from a reliability point of view, a cross-point array with same numbers of wordlines and bitlines is the best choice. Furthermore, we also notice that when the array has the same number of wordlines and bitlines, FWFB, HWFB and FWHB schemes have the same minimum driven voltage.

However, boosting the driven voltage also introduces other potential problems for the array. In particular, increasing the driven voltage also increases the voltage applied at unselected cells. Therefore, a write disturbance may occur when the voltage applied at an unselected cell exceeds the threshold voltage for SET or RESET operation. Figure 6 shows the maximum voltage applied at unselected cells with the minimum driven voltage, which is determined in Figure 4. Since the threshold voltage of the ReRAM cell is 2V, only array sizes with worst case voltage less than 2V are allowable. Otherwise, the array is unreliable because it can not avoid write failure and write disturbance at the same time. Therefore, Figure 6 provides a hard constraint on array size, and all of the following energy and area tradeoffs should be bounded by this constraint.

#### Energy Consumption of Write Operation.

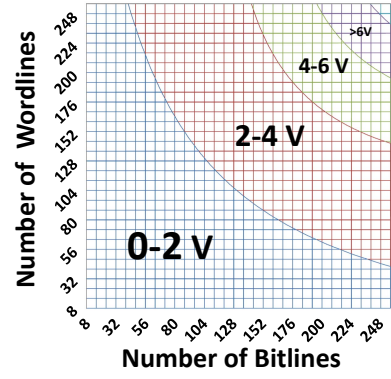


Fig. 6. The maximum voltage applied at unselected cells with the minimum driven voltage.

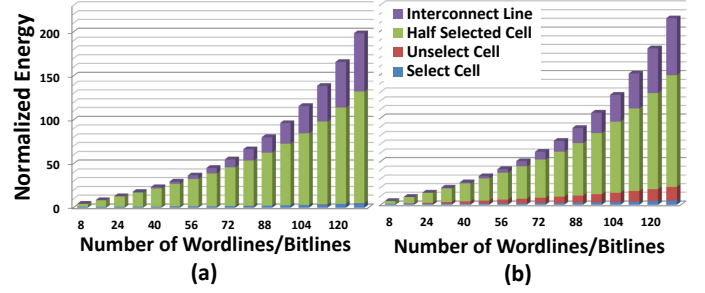


Fig. 7. The normalized energy consumption. (a): HWFB scheme (b): FWHB and HWFB schemes.

The energy consumption of a write operation for a cross-point array can be calculated as:

$$E_{write} = E_{select} + E_{unselect} + E_{halfselect} + E_{line}, \quad (14)$$

where the  $E_{select}$  is the energy consumed to change the state of the selected cell, the  $E_{unselect}$  and  $E_{halfselect}$  are the undesired energy wasted at the half selected and unselected cells. The energy consumed by the interconnect lines are represented by  $E_{line}$ . Figure 7 shows the decomposed energy consumption for the cross-point array. Note that, the  $E_{line}$  and  $E_{halfselect}$  take a large amount of the total energy consumption. Also, this part of energy wasted during the write operation takes greater part of the total energy for larger array sizes. For example, the undesired energy consumption for writing a  $128 \times 128$  array is more than 1000 times larger than the  $8 \times 8$  array. We also notice that, since the impact of sneak paths for floating schemes (FWHB and HWFB) is more serious, the energy consumed at unselected cells for floating schemes is larger than the half-biased scheme. Due to this reason, the total energy consumptions for FWHB and HWFB schemes are at least 10% larger than that of HWFB scheme.

#### Area cost of Write Operation.

The write operation for a  $M \times N$  array requires totally  $M + N$  voltage drivers. Therefore, the average number of ReRAM cells per voltage driver can be calculated as  $mn/(m + n)$ . Given the array capacity of  $C_{array}$ , it is easy to find that the optimal array organization can be achieved when  $M = N = \sqrt{C_{array}}$  and the maximum number of cells per voltage driver is:  $\sqrt{C_{array}}/2$ . However, the area overhead of a voltage driver is also related to its current drive capability. Figure 8(a) shows the maximum write current with different ReRAM array sizes. According to the current requirement, the area of the voltage driver can be directly calculated,



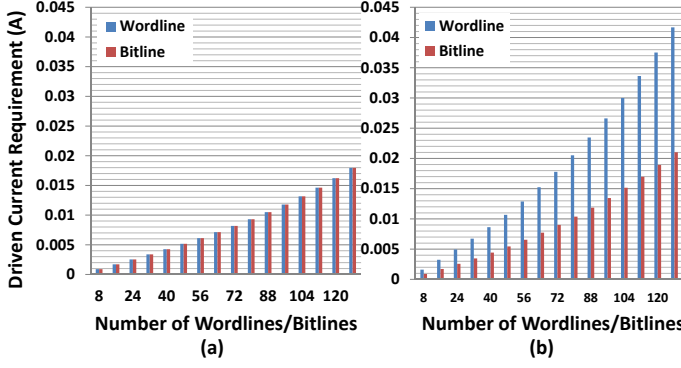


Fig. 8. The driven current requirements for wordlines and bitlines. (a) One bit writing; (B) Multi-bit writing.

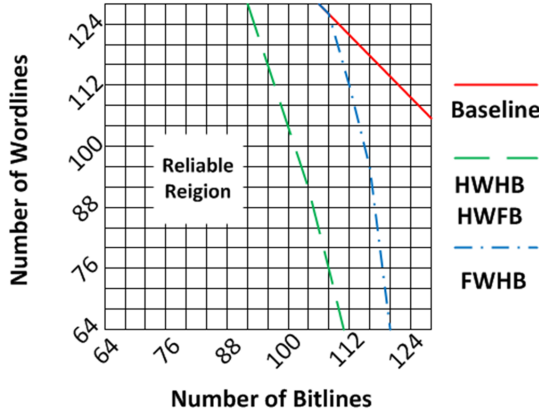


Fig. 9. The array size requirement for the cross-point array with different write schemes. (Baseline: one bit per access. HWHB, HWFB and FWFB: one wordline per access.

which will be shown in the following discussion.

#### Discussion on Multi-Bits Write Operation.

So far, we have only discussed the write operation with one bit per access. In this section, we consider the difference between one bit per access and one wordline per access write operations. Firstly, writing a wordline at a time will worsen the voltage drop along the wordline. Therefore, as shown in Figure 9, the reliable size of the cross-point array will be further reduced. The maximum array size reduces from  $116 \times 116$  to  $100 \times 100$  for HWFB and HWFB schemes.

In order to fairly compare the energy consumption, we compare the energy-per-bit instead of the total energy. For example, in order to write a wordline with size of 128, the energy-per-bit can be calculated as:  $E_{ave} = E_{total}/128/2$ . Figure 10 shows the energy-per-bit of the multi-bit write operation. The energy shown in this figure is normalized to the same unit as Figure 7 for easier comparison. The results show that for large cross point array sizes, the multi-bit write operation is much more energy efficient. This is because the energy wasted at the unselected and half-selected cells are shared by multiple bits and the average energy for one each bit is therefore reduced. However, although the multi-bit write operation has the advantage of lower energy consumption, the maximum current requirement for each wordline also increases. As shown in Figure 8(b), the maximum driven current for each bitline is almost the same as when writing one bit, however the drive capability for each word line is almost doubled for multi-bit writing. Since the area of the voltage driver

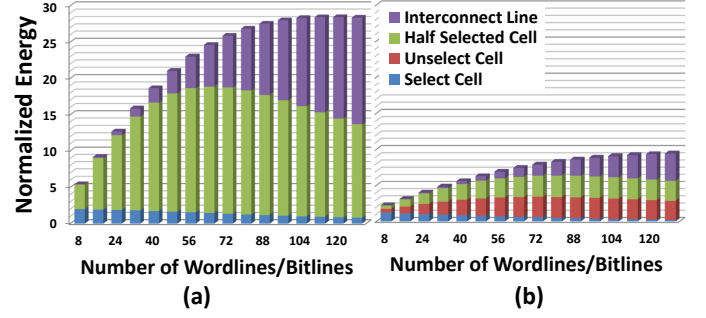


Fig. 10. The normalized energy consumption per bit for multi-bits write operation. (a): HWHB and FWFB schemes; (b): HWFB scheme.

increases proportionally with its drive capability, the area overhead for multi-bit writing is about 50% larger than for one bit writing.

#### Non-linearity of the ReRAM Cell.

One of the most distinct features of ReRAM is its non-linearity. For example, the non-linearity of memristor-based ReRAM is observed at LRS when the resistance of the memristor cell is not constant but varies with the applied voltage. The non-linearity coefficient is defined as:  $K_r(p, V) = p \times R(V/p)/R(V)$ , where  $R(V/p)$  and  $R(V)$  are the equivalent resistance of the memristor biased at  $V/p$  and  $V$  [?]. Normally, the  $K_r(p, V)$  value for memristor-based ReRAM is larger than 20, meaning that the resistance of a half-biased cell is at least 10 times larger than a full-biased cell. Clearly, the ReRAM cell with a larger non-linearity coefficient results in a better memory cell since the current in the sneak path will be significantly reduced. In addition, the increased resistance at half-selected and unselected cells can also mitigate the voltage drop along the activated wordline and bitline. Figure 11 shows the influence of different non-linearity coefficients on the array size requirements for one bit HWHB writing scheme. In this figure, the maximum array size increases from  $112 \times 112$  to  $340 \times 340$  when the non-linearity coefficient  $K_r$  increases from 1 to 10. Similarly, the non-linearity can also increase the maximum array size for other write schemes.

Moreover, the non-linearity can also reduce the energy consumption and area overheads of the cross-point array. For example, consider a  $128 \times 128$  array. As shown in Figure 12, the energy consumption for the write operation decreases dramatically with the increase of non-linearity coefficient  $K_r$ . As  $K_r$  increases from 1 to 40, the write energy is reduced by 98.3%. The driven current requirement is shown in Figure 13(a), and the corresponding area overheads of the voltage drivers are compared to the array size at Figure 13(b). The baseline design is unacceptable because the area

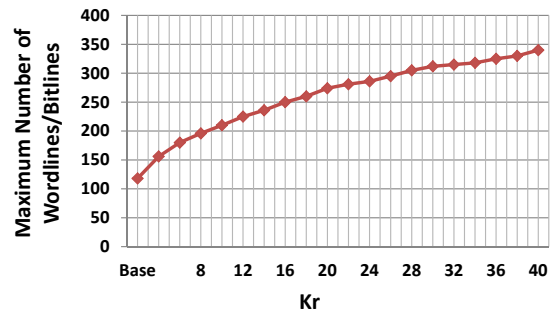


Fig. 11. The maximum array size with different non-linearity coefficient.

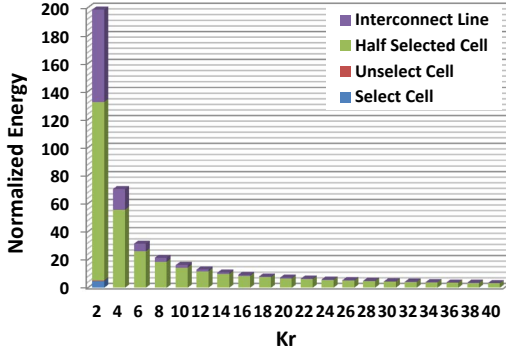


Fig. 12. The normalized energy consumption with non-linear ReRAM cells.

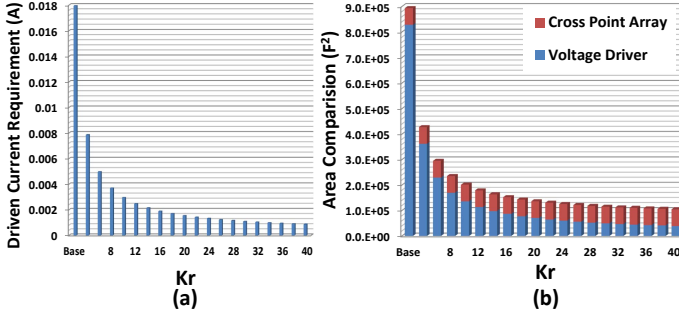


Fig. 13. The driven current requirements and area overheads with different non-linearity coefficients

of voltage drivers is about 11.6 times larger than the area of the cross-point array. In this case, the area efficiency of ReRAM's  $4F^2$  cell size will be offset by the extremely huge area overhead of the voltage drivers. However, with the increase of non-linearity, the area of voltage drivers becomes comparable to the array area. Therefore, we can conclude that, the ReRAM cells with a small non-linearity coefficient are not suitable for the cross-point structure based memory array. Next, we study the area overhead of multi-bit write. Figure 14 shows the normalized areas of the voltage drivers for one bit and multi-bit write operations. As mentioned, multi-bit write operations require larger driven current. Therefore, the area of voltage drivers for multi-bit write operations are much larger than that for one bit write operations. Finally, normalized areas of the one bit and multi-bit write operations have opposite trends as the array size increases. Normalized area for one bit write operation increases with the array size. On the contrary, normalized area for multi-bit write decreases as the array size increase.

### C. Read Operation

In this section we applied the similar sensing scheme as [?] and [?]. In order to read cell  $R_{i,j}$ , the  $i^{th}$  wordline is biased at  $V_{READ}$  and all of the other wordlines and bitlines are grounded. Then the state of the selected cell is read out by measuring the voltage across  $R_s$ . The energy consumption for read operation can be analyzed by the same way as that of the write operation. Since the read voltage is much smaller than write voltage, the read energy is expected at least one order smaller than write operation. Additionally, since the read voltage/current is much lower than the write, we believe that the voltage drivers can always provide enough current for the read operation if they meet the current requirement for write operation. Therefore, we can conclude that the area overhead of voltage drivers is determined by the write current. However, the reliability of read operation is different from the write operation. The read reliability

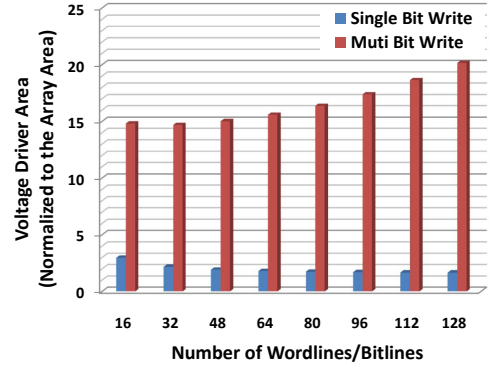


Fig. 14. The normalized area overhead of voltage drivers ( $K_r = 20$ , the areas are normalized to the area of cross-point array).

is determined by the voltage swing for reading HRS and LRS cells. Figure 15 (a) shows the voltage swing with different array sizes and  $K_r$  values. Large array sizes and large non-linearity are harmful to the voltage swing: on the one hand, a larger array has more sneak paths, making the output voltage very sensitive to the data pattern of unselected cells; on the other hand, the non-linearity increases the resistance of LRS and therefore the resistance difference between HRS and LRS cells is reduced. In order to improve the reliability of the read operation, a two-step sensing scheme can be applied, which senses the current of an unselected cell first, then the overall current is sensed, and after that the current difference is converted to the output voltage. The voltage swing of this two-step sensing scheme is shown in Figure 15 (b). By using this two-step sensing schemes, the voltage swing for a given array size and non-linearity coefficient is doubled.

## V. NON-LINEARITY AND WRITE CURRENT

## VI. CONCLUSION

ReRAM is a promising candidate for next-generation non-volatile memory technology. The area efficient cross-point structure is the most attractive memory organization for ReRAM memory design. However, problems inherent in the cross-point structure, such as the existence of sneak current and voltage drops along the nanowires introduce challenges to the design of reliable ReRAM cross-point array. In this paper, we first establish a mathematical model for cross-point arrays. We show that the proposed model has a simple structure and is

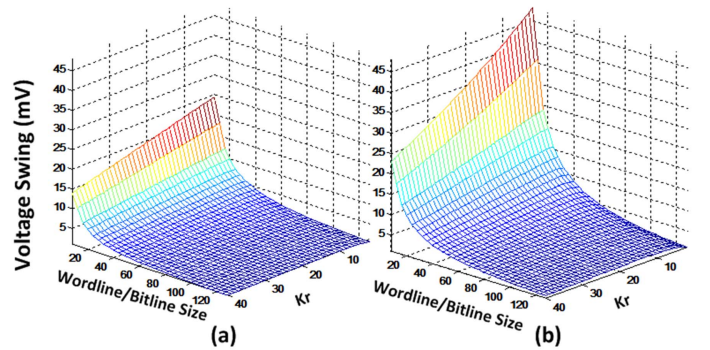


Fig. 15. Relationships among the voltage swing, array size and non-linearity. (a) Normal sensing scheme; (b) Two-step sensing scheme

flexible to evaluate different write/read schemes. By using this model, we study in detail how reliability affects the array organization, size, energy consumption, and area overheads in designing cross-point arrays. The simulation results show that, the multi-bit write operation is more energy efficient than one bit write operation and therefore is more suitable for energy-constrained design. However, from an area-constrained design, one bit write operation is better. Also, we point out that the non-linearity of the ReRAM cell can reduce the energy consumption and area overhead significantly, and it is favorable for large, energy efficient ReRAM design. Based on the results of our study, a detailed design methodology is proposed to help designers identify the optimal organization that meets the design constraints early in the design stage.