





# Computation-oriented Fault-tolerance Schemes for RRAM-based Computing Systems

Wenqin Huangfu<sup>1</sup>, <u>Lixue Xia<sup>1</sup></u>, Ming Cheng<sup>1</sup>, Xilin Yin<sup>1</sup>, Tianqi Tang<sup>1</sup>, Boxun Li<sup>1</sup>, Krishnendu Chakrabarty<sup>2</sup>, Yuan Xie<sup>3</sup>, Yu Wang<sup>1</sup>, Huazhong Yang<sup>1</sup>

1 Dept. of E.E., Tsinghua National Laboratory for Information Science and Technology (TNList), Tsinghua University, Beijing, China

2 Department of Electrical and Computer Engineering, Duke University

North Carolina, USA

3 Department of Electrical and Computer Engineering, University of California at Santa Barbara, California, USA



# **RRAM-based Computation**



 Emerging Devices, such as RRAM devices, provide a promising solution to realize better implementation of brain inspired circuits and systems in a high energy efficiency way



# **Related Work about RRAM-based NNs**

- Architecture and circuit designs
  - RRAM-based Deep Neural Networks
     [Hu\_DAC\_2012,Chi\_ISCA\_2016,Liu\_IEEE TCAS I\_2016, etc]
  - RRAM-based Spiking NNs [Tang\_DATE\_2015,Hu\_IEEE TCAD\_2016]
  - RRAM-based Convolutional NNs [Wang\_ISCAS\_2016, Shafiee\_ISCA\_2016]
  - Training NN on RRAM [Hasan\_IJCNN\_2014]
- EDA tools
  - Simulator [Xia\_DATE\_2016]
  - Design optimization tool [Gu\_ASPDAC\_2015]
- Demonstration
  - NN/Hopfield Network [Prezioso\_Nature\_2015,Lee\_IEDM\_2016]
  - Convolver [Gao\_IEEE Electron Device Letters\_2016]

#### Fault Classification of RRAM Device:

- Hard Fault (Stuck-At Faults)
  - The resistance is unchangeable
- Soft Fault
  - The resistance is changeable but not correct
- Static Fault
  - Already happened before using
- Dynamic Fault
  - Happened during using
- Soft faults can be tolerant by inner fault tolerance of NN [Gu\_ASPDAC\_2015]
- Previous designs ignore the influence of hard faults



# **Impact of Hard Faults**

- The yield (the percentage of available cells) of RRAM device varies from 60% to more than 90% in different materials and technologies [Chen\_IEEE TC\_2015]
- We test the classification accuracy of a 784x100x10 NN on MNIST test bench with random stuck-at-0 and stuck-at-1 faults:

Yield	Ideal	95%	90%	80%
Accuracy	97.8%	26.7~60.4%	15.5~38.6%	10.6~28.0%
Reduction	-	>37%	>59%	>69%

 Hard faults can obviously influence the performance of RRAM-based computation



# **Fault-tolerant Research**

- Memory-oriented methods use redundant cells and arrays to substitute faulty cells [Wang\_DAC\_12,Koh\_ICCD\_09]
- These methods are not available for RRAM-based computation
  - Redundant cell is helpless: the basic unit is the whole column instead of single cell
  - Unable to directly shut down faulty cells: if the size of RRAM crossbar array is M \* N, M \* N control lines are needed to control every cell independently, which is unacceptable
- We need computation-oriented research to tolerant the hard faults in RRAM-based computation system



# Fault-tolerant Method for RRAM-based Neural Computing System

- Background
  - NNs on RRAM
  - RRAM Faults
- Proposed Methods
  - Fault-tolerant mapping and redundant scheme
- Future Work

# **Proposed Methods**



- Mapping algorithm with inner fault-tolerant ability
- Redundant schemes and circuits
  - Redundant crossbars
  - Independent Redundant Columns

# **Mapping Algorithm**

• Traditionally, two RRAM crossbars are needed to represent the positive and negative values of matrix:  $\vec{V_o} = C^+ \cdot \vec{V_i} + C^- \cdot - \vec{V_i}$ 

 $= (C^+ - C^-) \cdot \vec{V}_i$ 

Original mapping algorithm accords to the sign of weights to determine the RRAM values:  $\int c_{k,i} = c_{k,i} > 0$ 

 $= C \cdot \vec{V}_i$ 

$$c_{k,j}^{+} = \begin{cases} c_{k,j} & c_{k,j} > 0\\ 0 & c_{k,j} \le 0 \end{cases}$$
$$c_{k,j}^{-} = \begin{cases} -c_{k,j} & c_{k,j} < 0\\ 0 & c_{k,j} \ge 0 \end{cases}$$

- This mapping method doesn't consider RRAM faults
  - Example
  - Target Value: -5
  - Range of value represent by RRAM devices: [1, 10]



# **Mapping Algorithm**

 However, in some situations, the influence of hard faults can be directly tolerant in mapping phase



- Proposed Mapping Method
  - Target: minimize the difference between the ideal value of matrix and the practical value represented by RRAM
  - Step 1: Initializing
  - Step 2: Greedily mapping the weights cell-by-cell
  - Cannot solve all situations



Algorithm 1: Mapping Algorithms with Inner Fault-Tolerant Ability

```
Input: C, G_{max}, G_{min} g_s, StuckAtMap,
          Available\_Rdundant\_Cells
  Output: G^+_{Output}, G^-_{Output}
for i = 1 : C.cells do
      if G^+_{Output}(i) doesn't get stuck then
           G^+_{Output} = Gmax
       end
      if G_{Output}^{-}(i) doesn't get stuck then
           G^-_{Output} = Gmin
      end
8 end
9 for i = 1 : C.cells do
       for j = 1: Available_Rdundant_Cells(i) do
           Diff = C(i) - (G^+_{Output}(i,:)/Gs - G^-_{Output}(i,:)/Gs);
           Adjust G^+_{Output}(j) and G^-_{Output}(j) to minimize Diff;
           Diff = C(i) - (G^+_{Output}(i,:)/Gs - G^-_{Output}(i,:)/Gs);
           if Diff == 0 then
                break;
           end
       end
   end
17 return G^+_{Output}, G^-_{Output}
```

# **Proposed Method**

- Mapping algorithm with inner fault-tolerant ability
- Redundant schemes and circuits



### **Hardware Structure**

- Two Redundant schemes and corresponding redundant circuits are proposed
  - Redundant crossbars
  - Independent Redundant Columns





Mapping algorithm: Co-mapping different crossbars to get the best available results

- Target application matrix: C
- Original computing process:



Computing process with redundant crossbars:



#### **Cell-by-cell Mapping Process** Target Value: -5 Range of value represent by RRAM devices: [1, 10] Redundancy Ratio (Number of Redundant Crossbars): 1 Initialization: 3 4 10 🄀 🕂 🔞 — 🍫 — 🏷 5 Step 1 (Adjust the first cell): skip 7 Step 2 (Adjust the second cell): 10 11 12 13 Step 3 (Adjust the third cell): 14 15

• Step 4 (Adjust the fourth cell):



Algorithm 1: Mapping Algorithms with Inner Fault-Tolerant Ability Input: C,  $G_{max}$ ,  $G_{min}$   $g_s$ , StuckAtMap,  $Available_Rdundant_Cells$ **Output**:  $G^+_{Output}$ ,  $G^-_{Output}$ 1 for i = 1 : C.cells do if  $G^+_{Output}(i)$  doesn't get stuck then  $G^+_{Output} = Gmax$ end if  $G_{Output}^{-}(i)$  doesn't get stuck then  $G^-_{Output} = Gmin$ end 8 end 9 for i = 1 : C.cells do for j = 1: Available\_Rdundant\_Cells(i) do  $Diff = C(i) - (G^+_{Output}(i,:)/Gs - G^-_{Output}(i,:)/Gs);$ Adjust  $G^+_{Output}(j)$  and  $G^-_{Output}(j)$  to minimize Diff;  $Diff = C(i) - (G^+_{Output}(i, :)/Gs - G^-_{Output}(i, :)/Gs);$ if Diff == 0 then break; end end 16 end 17 return  $G^+_{Output}$ ,  $G^-_{Output}$ 

- Computing Process: merge the results from the original RRAM crossbar and the results from redundant crossbars together
- Definition: R is the redundant ratio, standing of extra RRAM crossbar arrays that used as redundant structure



**Positive Parts** 

**Negative Parts** 

# **Independent Redundant Columns**

- Given the fault ratio of a single device is *P*, the size of RRAM crossbar array is *M* x *N*, the expectation of faults in each column is *P* x *M* 
  - Solution: A smaller redundant unit which only contains *P* x *M* cells
  - Divide each column into P x M parts (called 'cut'), while the expectation of faults in each cut is 1
  - R indicates the number of redundant devices in a redundant column for one cut



# **Independent Redundant Columns**



# **Independent Redundant Columns**

Mapping algorithm: Co-mapping corresponding columns to get best available result

- Target application matrix: C
- Original Mapping process:



Modified Mapping process:



### **Hardware Overhead**

- The proposed redundant methods requires additional RRAM cells and corresponding function/control modules
- The detailed comparison between original method and proposed two methods is as follows:

TABLE III Hardware overhead of original RCS design and two redundant schemes we proposed

	Original	Redundant Xbar	Redundant Columns		
RRAM	$2 \times M \times N$	• $2 \ge (R+1) \times M \times N$	$2 \times M \times N + 2 \times 2 \times \lceil R \times P \rceil \times M \times N$		
ADC	2  imes N	$2 \times (R+1) \times N$	N		
DAC	M	M	M		
Op Amps	$2 \times N$	$2 \times (R+1) \times N$	2  imes 2  imes N		
Decoder	$2 \times M$ outputs +	$2 \times (R+1) \times M$ outputs +	$2 \times M$ outputs + $2 \times N$ outputs +		
Decoder	$2 \times N$ outputs	$2 \times (R+1) \times N$ outputs	$2 \times (2 \times \lceil R \times P \rceil \times M \text{ outputs})$		
Mux	0	• <sub>s</sub> 0	$2 \times 2 \times [R \times P] \times M \times N \times (\lfloor \frac{1}{n} \rfloor \times \text{inputs})$		
Adder	0	R  imes N	2  imes N		
Sub	$2 \times N$	2  imes N	$2 \times N$		

Three case studies

- 1. Mapping an 128x128 matrix onto RRAM-based computing systems
- 2. 128x128 Matrix-vector multiplication on RRAM-based computing system
- 3. 784x100x10 NN for MNIST on RRAM-based computing system;

Three evaluation metrics

- 1. Mapping Error: I<sub>2</sub> norm distance between ideal matrix and mapped matrix
- 2. Computing Error: I<sub>2</sub> norm distance between ideal output vector and RRAM computed output vector
- 3. MNIST Error: error rate of MNIST classification accuracy. The error rate of MNIST on CPU platform (ideal case) is 2.17%



- Performance of Mapping Algorithm
  - Great improvement for all three evaluation metrics under different faults
  - Especially effective for testing cases with low (below 5%) faults
  - For 5% SAFs, MAO reduce MNIST error from 52.11% to 4.01%
  - No extra overhead is introduced by MAO

SAFs	Mapping Error (%)			Computational Error (%)			MNIST Error (%)		
(%)	Origin [10]	MAO	Percentage Improvement (%)	Origin [10]	MAO	Percentage Improvement (%)	Origin [10]	MAO	Percentage Improvement (%)
1	16.60	10.10	6.50	16.56	10.14	6.42	12.42	2.24	10.18
3	28.82	17.71	11.11	28.93	17.53	11.40	34.60	2.89	31.71
5	37.04	23.11	13.93	36.80	23.11	13.69	52.11	4.01	48.10
7	43.89	27.70	16.19	43.40	27.21	16.19	62.31	7.55	54.76
8	46.83	29.88	16.95	46.88	29.71	17.17	66.52	10.99	55.53
10	52.83	34.81	18.02	52.39	34.88	17.51	74.17	15.38	58.79
15	63.86	42.80	21.06	63.89	42.60	21.29	79.67	41.57	38.10
20	73.72	53.15	20.57	73.70	53.31	20.39	82.25	58.19	24.06

• Performance and overhead under some typical cases

TABLE IV Performance and Energy overhead of RX and IRC on MNIST Simulate 100 times with random SAFs

SAFs (%)	Design of RCS	R	Range of MNIST Error (%)	Average MNIST Error (%)	Average Percentage Improvement (%)	Energy Overhead (%)
	Original	0	[39.56, 73.42]	52.11	-	-
	MAO	0	[2.88, 9.16]	4.01	48.10	-
5	KX	1	[2.13, 6.75]	2.22	49.89	37.58
5	RX	2	[2.13, 2.27]	2.18	49.93	75.17
	IRC	1	[2.46, 4.07]	3.32	48.79	20.85
	IRC	2	[2.14, 2.85]	2.32	49.79	23.55
	Original	0	[61.44, 84.47]	74.17	-	-
10	MAO	0	[8.53, 36.15]	15.38	58.79	-
	RX	1	[2.37, 7.85]	2.83	71.34	37.58
	RX	2	[2.15, 7.02]	2.70	71.47	75.17
	IRC	1	[4.97, 26.06]	11.44	60.03	22.68
	IRC	2	[2.25, 9.87]	3.87	70.30	27.22
20	Original	0	[71.99, 89.38]	82.25	-	-
	MAO	0	[49.82, 70.46]	58.19	24.06	-
	RX	1	[10.84, 37.8]	19.30	57.50	37.58
	RX	2	[3.33, 13.37]	5.91	77.95	75.17
	RX	3	[2.38, 6.46]	2.65	79.60	112.75
	IRC	1	[35.96, 69.55]	54.21	28.04	26.46
	IRC	2	[6.00, 29.82]	13.52	68.73	34.78
	IRC	3	[2.36, 4.00]	3.65	78.60	43.09



• Energy-Performance Trade-Off (Task: MNIST)



Fig. 4. (a).*R*-Energy relation of RX. (b).*R*-Energy relation of IRC. (c).Energy-Error trade-off of RX (10% SAFs). (d).Energy-Error trade-off of IRC (10% SAFs).

#### References

[1] Prakash, Aravind, and S. Maikap. "Improved resistance memory characteristics and switching mechanism using TiN electrode on TaOx/W structure." *Proc. Winter Simul. Conf.* 2013.

[2] Chen, C.; Shih, H.; Wu, C.; Lin, C.; Chiu, P.; Sheu, S.; Chen, F.T., "RRAM Defect Modeling and Failure Analysis Based on March Test and a Novel Squeeze-Search Scheme," *Computers, IEEE Transactions on*, vol.64, no.1, pp.1,11, Jan. 1 2015

[3] Lee D, Park J, Moon K, et al. Oxide based nanoscale analog synapse device for neural signal recognition system[C]//2015 IEEE International Electron Devices Meeting (IEDM). IEEE, 2015: 4.7. 1-4.7. 4.

[4] Prezioso M, Merrikh-Bayat F, Hoskins B D, et al. Training and operation of an integrated neuromorphic network based on metal-oxide memristors[J]. Nature, 2015, 521(7550): 61-64.

[5] Wang, Jue, Xiangyu Dong, and Yuan Xie. "Point and discard: a hard-error-tolerant architecture for non-volatile last level caches." *Proceedings of the 49th Annual Design Automation Conference*. ACM, 2012.

[6] Cheng-Kok Koh; Weng-Fai Wong; Yiran Chen; Hai Li, "The salvage cache: A fault-tolerant cache architecture for next-generation memory technologies," *Computer Design*, 2009. ICCD 2009. IEEE International Conference on , vol., no., pp.268,274, 4-7 Oct. 2009

[7] Hasan R, Taha T M. Enabling back propagation training of memristor crossbar neuromorphic processors[C]//2014 International Joint Conference on Neural Networks (IJCNN). IEEE, 2014: 21-28.

[8] Degraeve R, Fantini A, Raghavan N, et al. Causes and consequences of the stochastic aspect of filamentary RRAM[J]. Microelectronic Engineering, 2015, 147: 171-175.

[9] Shih, Hsiu-Chuan, et al. "Training-based forming process for RRAM yield improvement." VLSI Test Symposium (VTS), 2011 IEEE 29th. IEEE, 2011.

[10] Chen Y Y, Roelofs R, Redolfi A, et al. Tailoring switching and endurance/retention reliability characteristics of HfO 2/Hf RRAM with Ti, Al, Si dopants[C]//VLSI Technology (VLSI-Technology): Digest of Technical Papers, 2014 Symposium on. IEEE, 2014: 1-2.

[11] Shafiee A, Nag A, Muralimanohar N, et al. ISAAC: A Convolutional Neural Network Accelerator with In-Situ Analog Arithmetic in Crossbars[C]//Proc. ISCA. 2016.

[12] Liu X, Mao M, Liu B, et al. Harmonica: A framework of heterogeneous computing systems with memristor-based neuromorphic computing accelerators[J]. IEEE Transactions on Circuits and Systems I: Regular Papers, 2016, 63(5): 617-628.

[13] Xia L, Li B, Tang T, et al. MNSIM: Simulation platform for memristor-based neuromorphic computing system[C]//2016 Design, Automation & Test in Europe Conference & Exhibition (DATE). IEEE, 2016: 469-474.

[14] Gu P, Li B, Tang T, et al. Technological exploration of rram crossbar array for matrix-vector multiplication[C]//The 20th Asia and South Pacific Design Automation Conference. IEEE, 2015: 106-111.

[15] Tang T, Xia L, Li B, et al. Spiking neural network with rram: Can we use it for real-world application?[C]//Proceedings of the 2015 Design, Automation & Test in Europe Conference & Exhibition. EDA Consortium, 2015: 860-865.

[16] Chi P, Li S, Qi Z, et al. PRIME: A Novel Processing-In-Memory Architecture for Neural Network Computation in ReRAM-based Main Memory[C]//Proceedings of ISCA. 2016, 43.



# Thank you!



