



上海交通大学
SHANGHAI JIAO TONG UNIVERSITY



Self-Terminating Write of Multi-Level Cell ReRAM for Efficient Neuromorphic Computing

Zongwu Wang (Speaker)

Zhezhi He*, Rui Yang, Shiquan Fan, Jie Lin, Fangxin Liu, Yueyang Jia, Chenxi Yuan, Qidong Tang and Li Jiang*

Shanghai Jiao Tong University

2022年4月4日



上海交通大学
SHANGHAI JIAO TONG UNIVERSITY



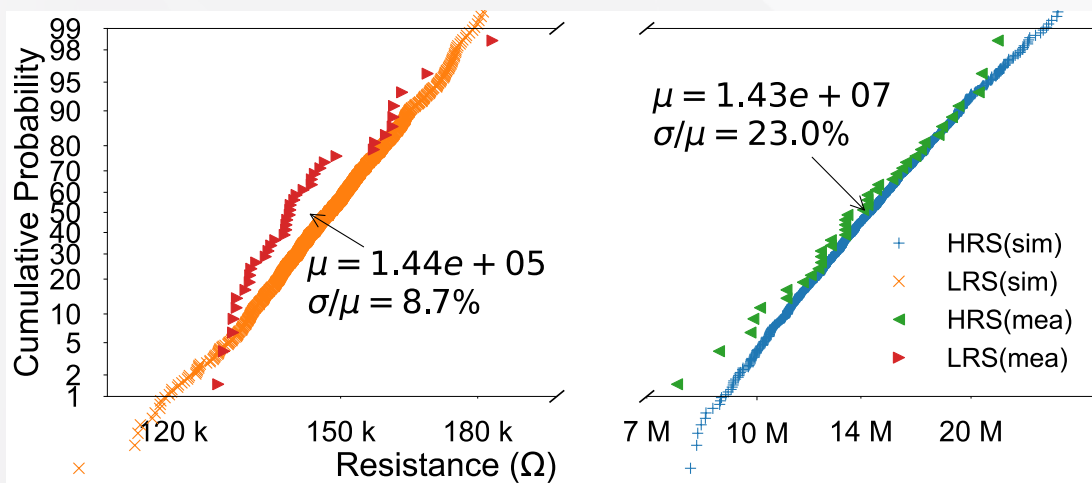


Self-Terminating Write Scheme Overview

Challenges In ReRAM-based PIM

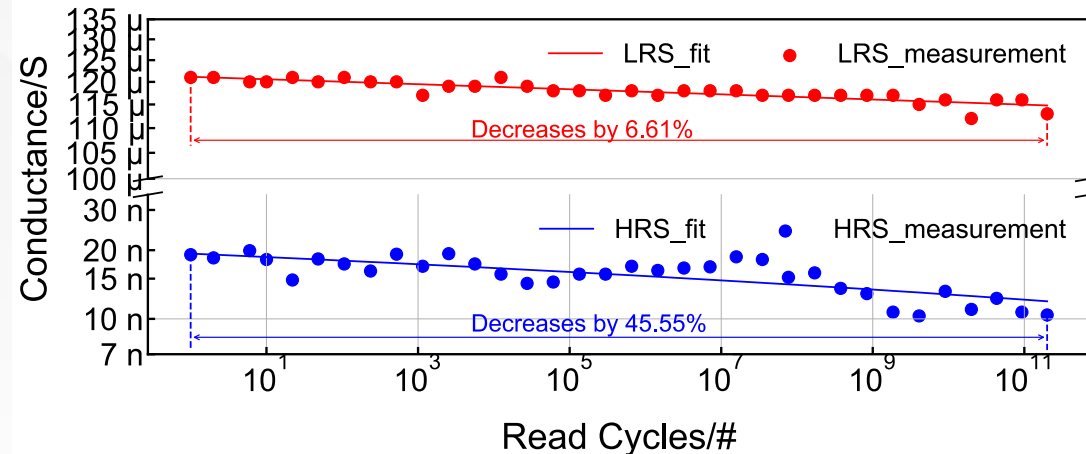
- ReRAM has intrinsic write variation
- Read disturb induces resistance drifting
- Write-verify scheme is relative slow

Write variation (CDF vs resistance)



- Measurement and simulation results comparison
- C2C variation exists in programming
- Set and Reset have significant write variation (8% and 23%, respectively)

Read disturb (conductance vs. #read times)

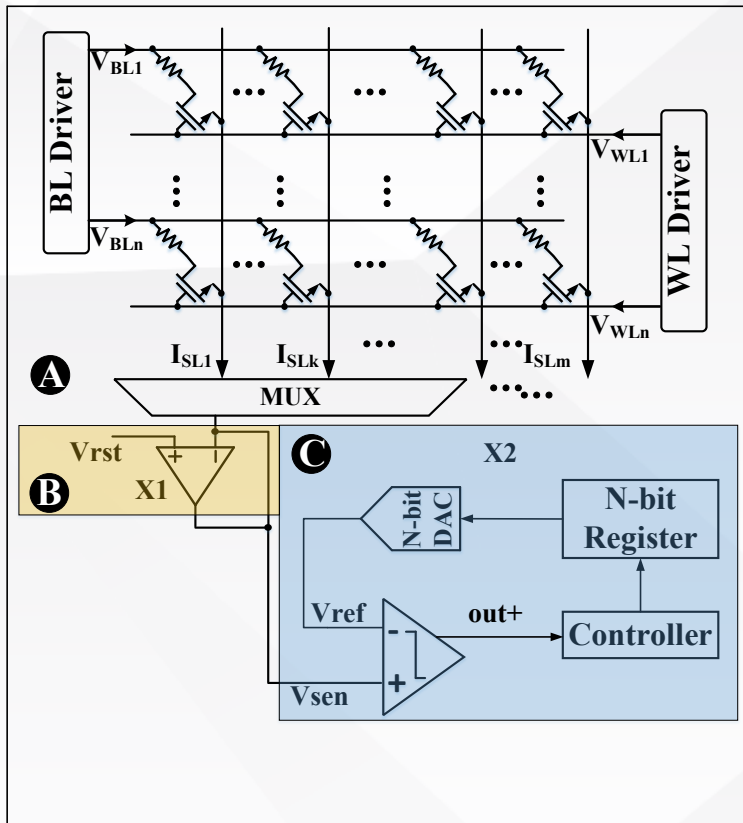


- ReRAM suffers from read induced drifting
- In-Memory computing equivalents to read
- Reliability test shows 6.6% and 45.6% drifting

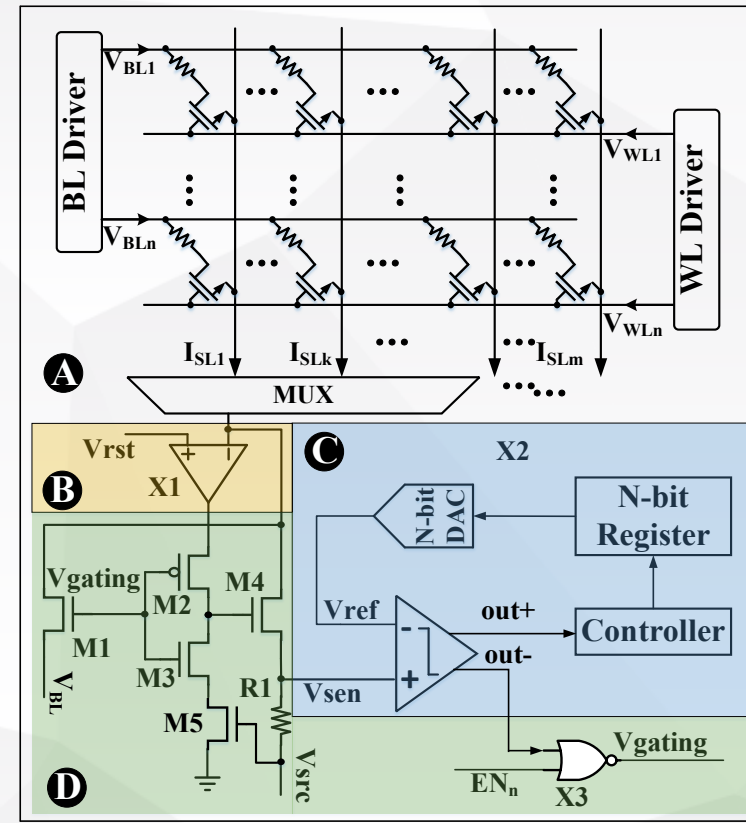
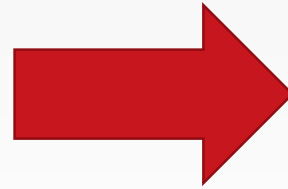
Proposed Solution

- Heavily peripherals reuse achieves precise self-terminating scheme(2-bit)
- Pick appropriate programming range according to circuit design
- Compare to Write-verify scheme, Reduce the latency and energy by 4.7x and 2x, respectively

Self-Terminating Write Scheme Design



Existing ReRAM-based PIM system



ReRAM MLC STW Schematic

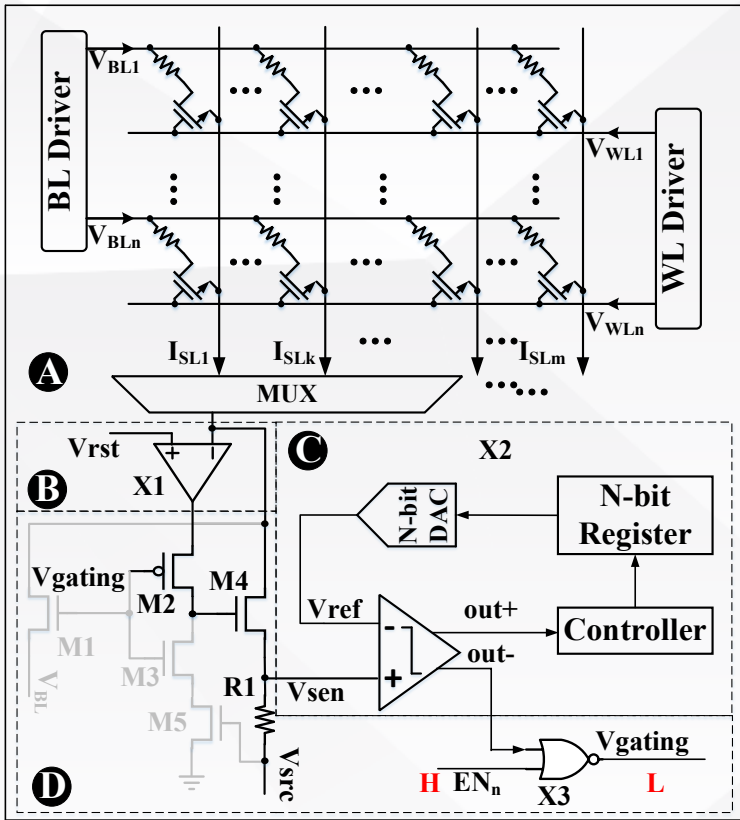
A: ReRAM Array
B: TIA module
C: SAR-ADC
D: Verdict module

Proposed Self-Terminating Write Scheme

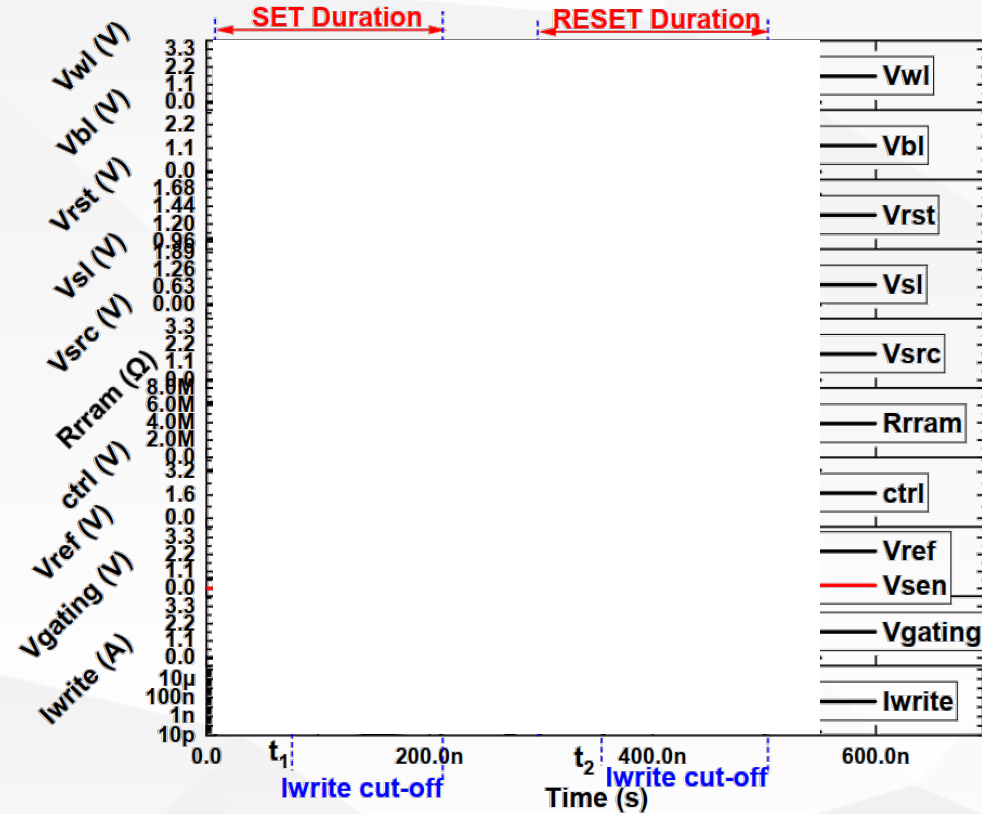
- Heavily reuse the peripherals in ReRAM-based PIM system (**3 + 1 modules**)
- Implementing both Set and Reset termination with circuit sharing (**3 modes**)
- Ultra-compact design contributes to low cost and fast feedback (high precision)



Self-Terminating Write Transition Waveform



Inference Mode Transition



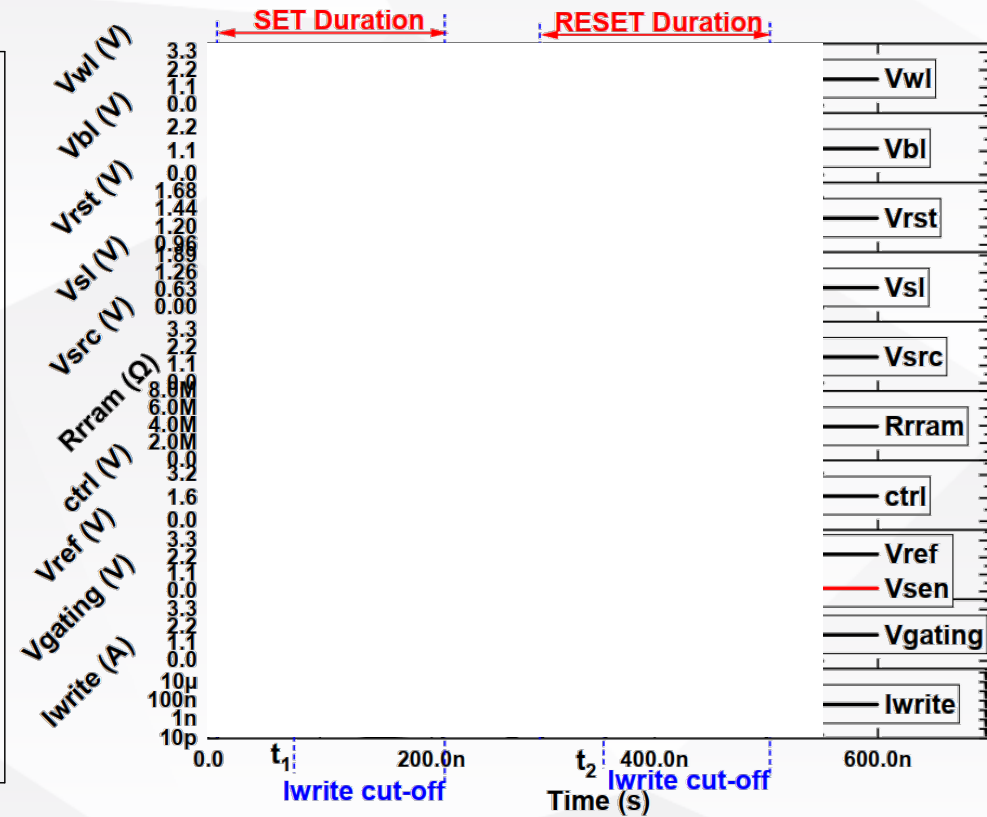
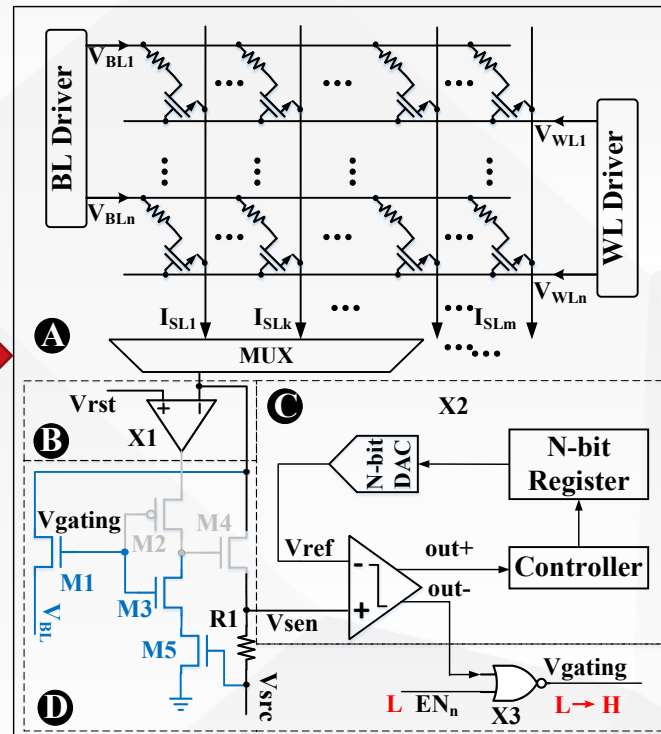
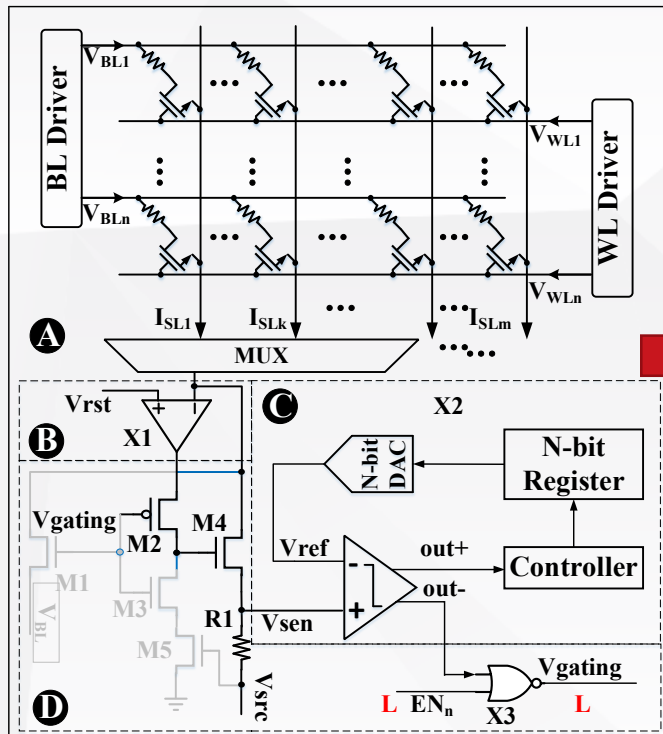
Programming Waveform

MODE	EN_n	VBL	Vsrc	COND	Vsen	Vout-	Vgating
Inference	H	Vread	GND	-	-	-	L
Set_term	L	Vset	GND	↗	↗	↘	↗
Rst_term	L	GND	VDD	↘	↗	↘	↗





Self-Terminating Write Transition Waveform



SET Terminating Transition

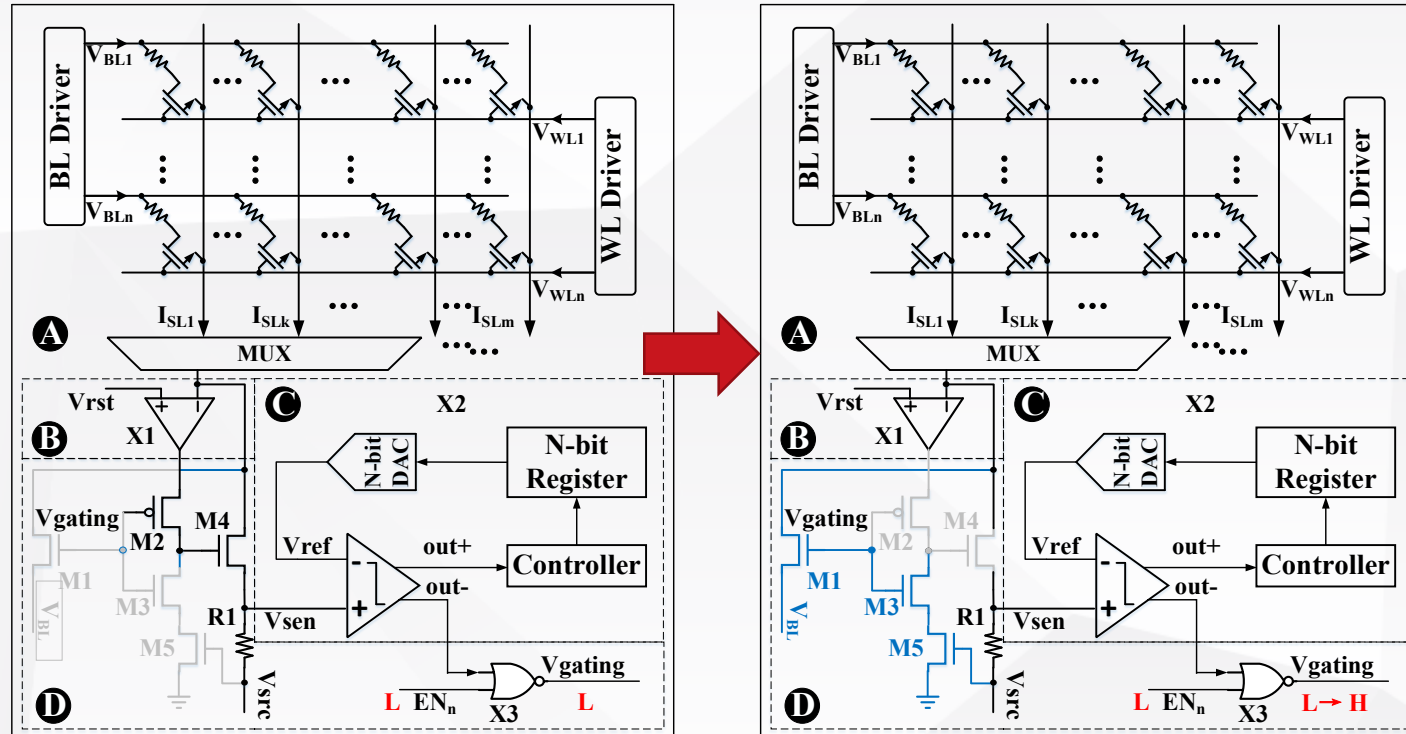
Programming Waveform

MODE	EN_n	VBL	Vsrc	COND	Vsen	Vout-	Vgating
Inference	H	Vread	GND	-	-	-	L
Set_term	L	Vset	GND	↗	↗	↘	↗
Rst_term	L	GND	VDD	↘	↗	↘	↗

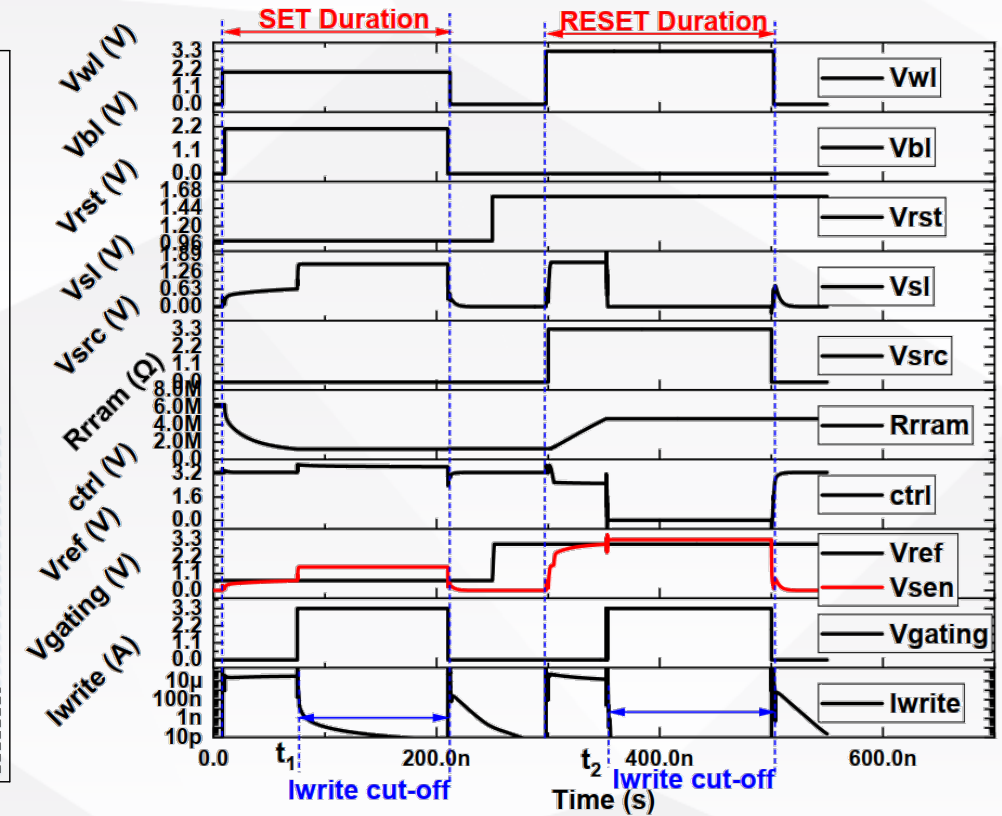




Self-Terminating Write Transition Waveform



RESET Terminating Transition



Programming Waveform

MODE	EN_n	VBL	Vsrc	COND	Vsen	Vout-	Vgating
Inference	H	Vread	GND	-	-	-	L
Set_term	L	Vset	GND	↗	↗	↘	↗
Rst_term	L	GND	VDD	↘	↗	↘	↗

Conclusion

➤ No difference between set & reset



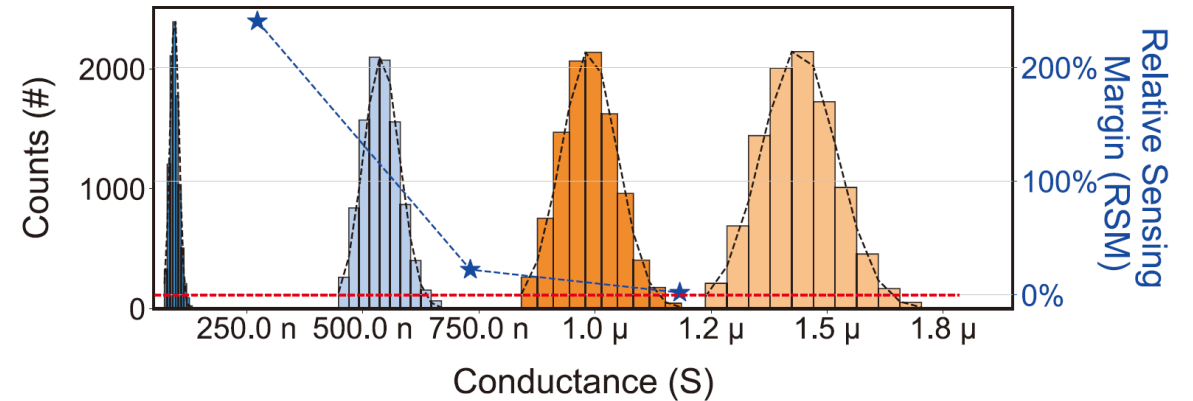


Self-Terminating Write Scheme Evaluation

	Structure	Area	Terminate	Precision
This work	2Amp+5T+NOR	Medium	both	2 bits
JSSC-2013 [10]	2Amp+R+30T +DelayUnit+others	Large	both	1 bit
ISSCC-2014 [24]	4T	Small	set	1 bit
IEDM-2017 [6]	RESET: Amp+4SW+6T SET: 5T	Medium	both	1 bit
ISSCC-2021 [25]	2Amp+R+5T+3INV +AND+Delay Unit	Large	set	1 bit

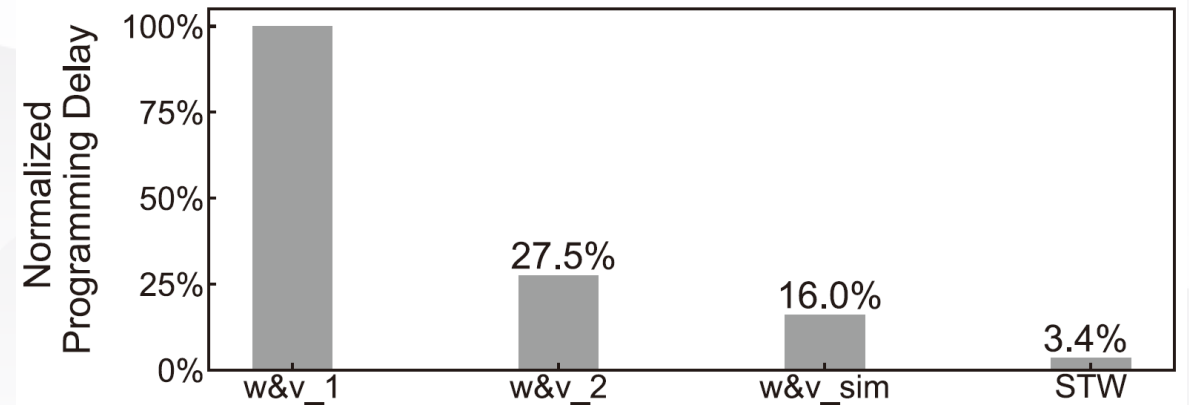
Comparison with previous works (area, programming polarity and precision) :

- Reduces area overhead by peripherals reuse
- Supports both Set and Reset termination
- Achieves 2-bit MLC self-terminating



10^4 trials MC simulation with range selection algorithm

- the proposed STW scheme achieves 2-bit precision

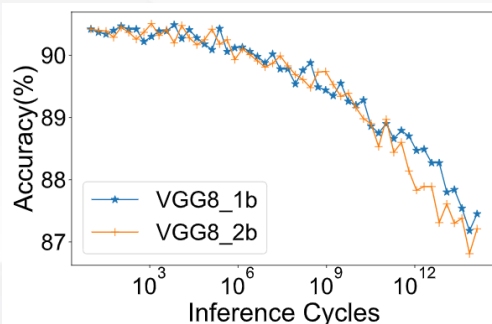


Latency comparison between different schemes

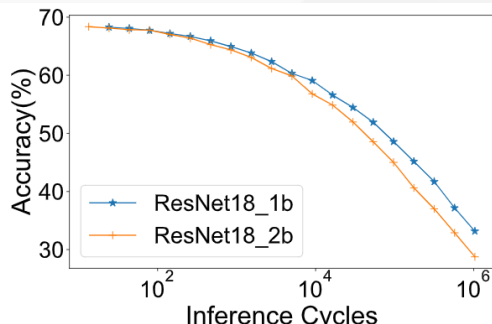
- STW scheme shows 4.7x speedup (conservative)



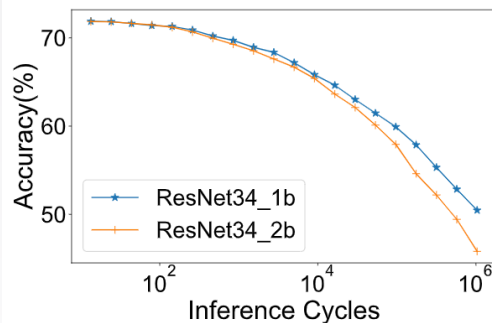
Self-Terminating Write Scheme Evaluation



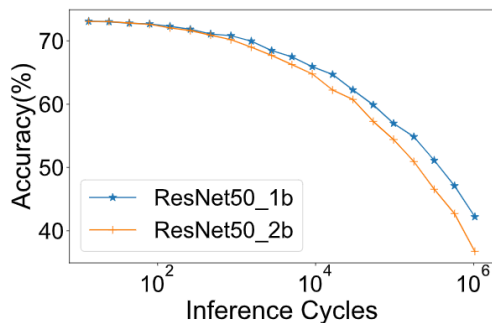
(a) VGG8 on CIFAR-10



(b) ResNet-18 on ImageNet



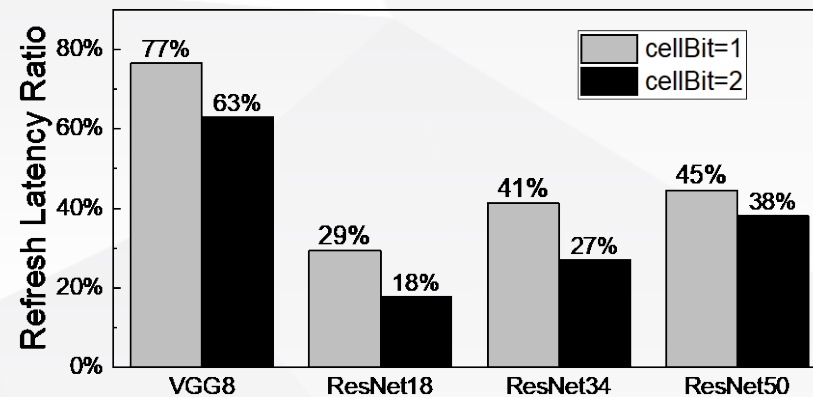
(c) ResNet-34 on ImageNet



(d) ResNet-50 on ImageNet

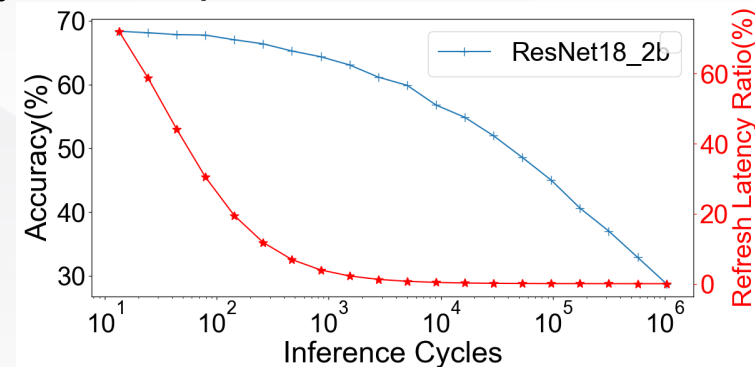
Impact of Read disturb on inference accuracy:

- Accuracy loss with the continuous inference after the network deployed
- MLC can reduce the storage/computation cost, but it is more vulnerable to read disturb



Proportion of delay on different networks

- Ratio of refresh latency is low on compact networks
- From the perspective of deployment cost, programming delay is an important factor



Refresh Frequency and Accuracy Balancing:

- The lower the refresh frequency, the lower the proportion of refresh delay, but the lower the accuracy



1. An auto-calibrate Framework

- Provides an easy-use and confident ReRAM compact model

2. A valid self-terminated programming scheme for MLC

- Heavily reuses the original peripheral
- Compact design achieves low cost and high precision
- Reduce the latency and energy by 4.7x and 2x, respectively

3. Cross-layer simulation (device/circuit/system) to validate the design





Thanks for Listening

饮水思源 爱国荣校