

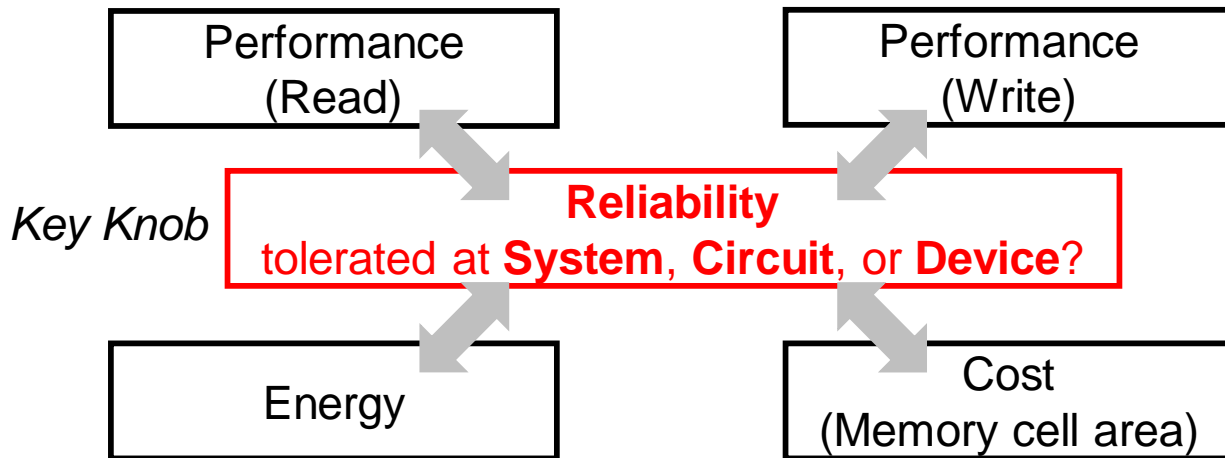


# TaO<sub>x</sub>-based ReRAM for Variability-Aware Approximate Computing

Chihiro Matsui, Shouhei Fukuyama,  
Atsuna Hayakawa, and Ken Takeuchi  
Chuo University, Tokyo, Japan



# Reliability-Aware Approximate Computing in Storage



K. Takeuchi, *IEDM* 2017.

Which hierarchy of ReRAM storage has Error Toleration techniques?  
How to Relax Reliability for Approximate Computing?



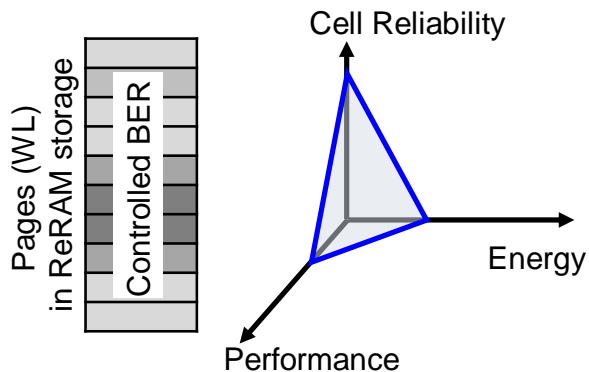
# Outline

- Variability-Aware Approximate Computing (V-AC)
- Application-Induced Variability of TaO<sub>x</sub> ReRAM Cell Errors and V-AC Evaluation Platform
- V-AC Strategies of System, Circuit and Device Co-Design (SCDCD)
- Conclusions

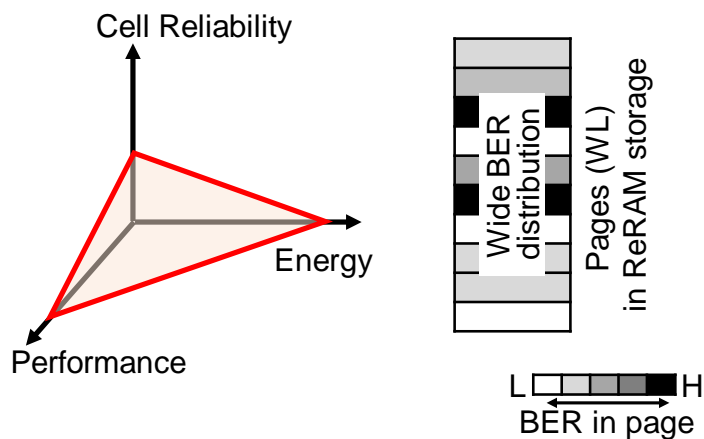


# Variability-Aware Approximate Computing (V-AC)

## Conv. Exact Computing



## V-AC for Machine Learning



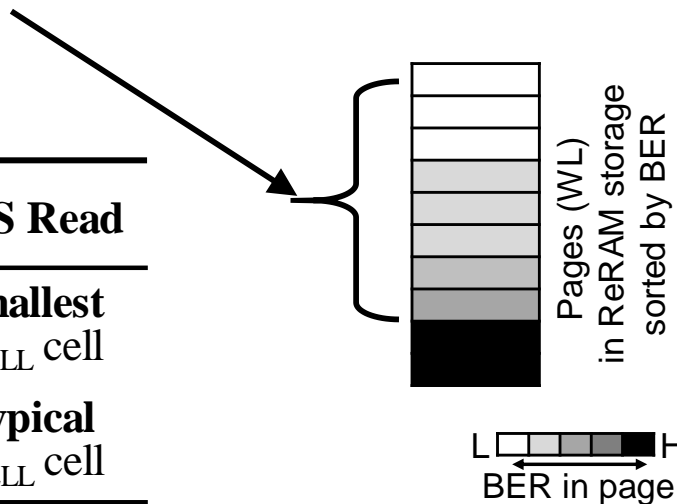
- System, Circuit, and Device in ReRAM-based storage have Variabilities in nature
- By tolerating variability, Performance, Energy, and Cost gain



# Typical Cell Target Strategy of V-AC in ReRAM Storage

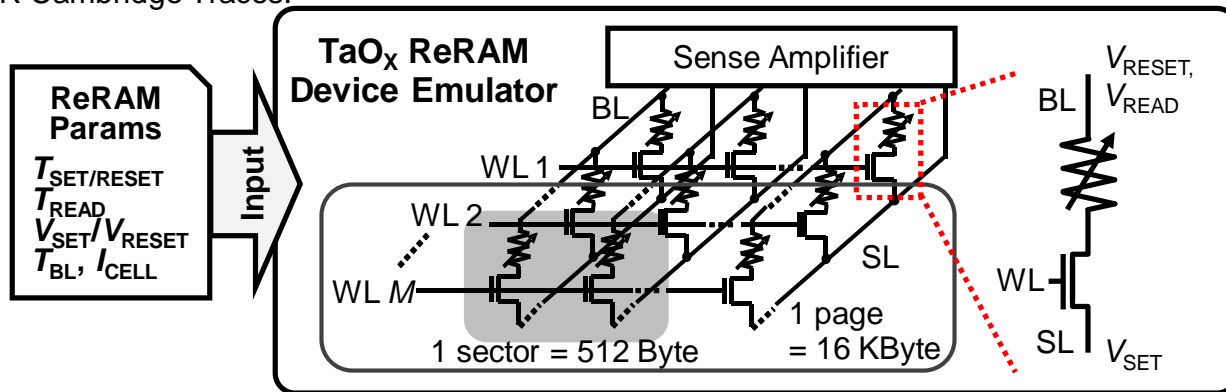
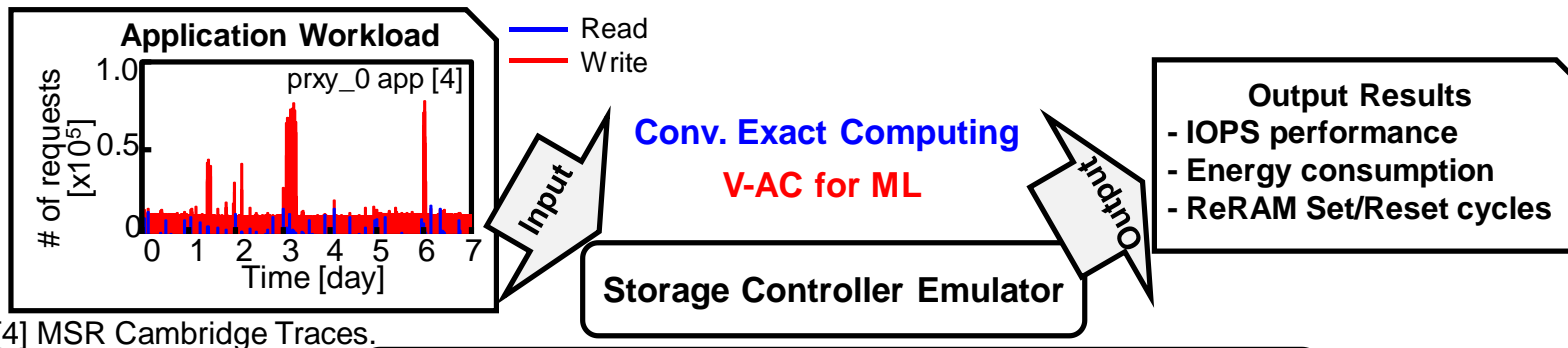
## Target Typical ReRAM Cell, NOT Worst Cell

	Error tolerance	ECC	Set/Reset	LRS Read
<b>Conv. Computing</b>	No error	Worst cell	Slowest cell	Smallest $I_{CELL}$ cell
<b>V-AC</b>	1-10 % errors [*]	Typical cell	Typical cell	Typical $I_{CELL}$ cell



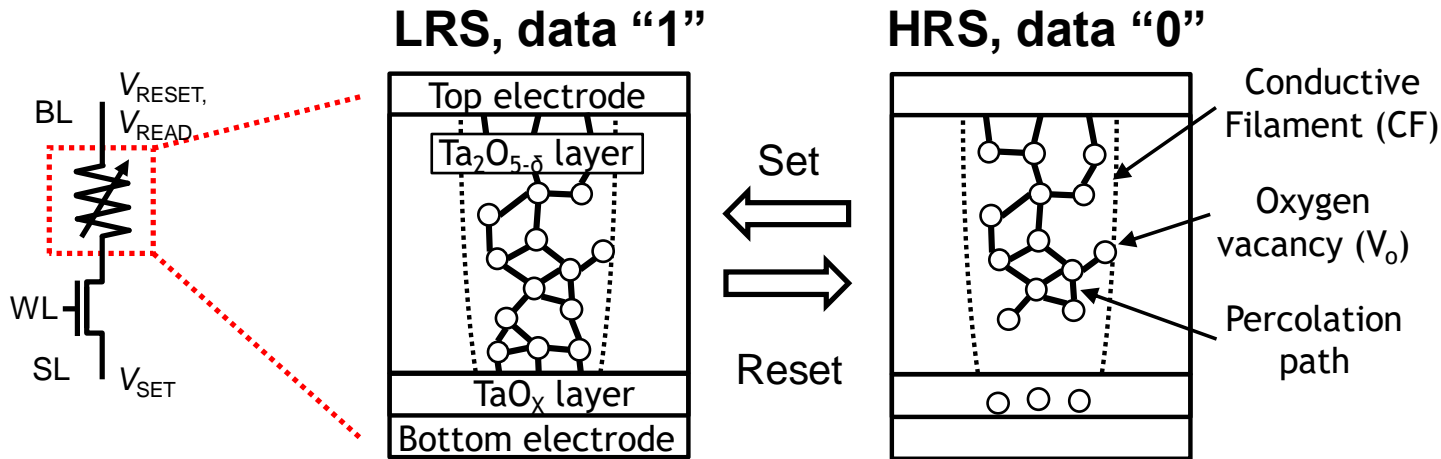


# System, Circuit and Device Co-Design (SCDCCD) Platform [\*]





# Set/Reset in TaO<sub>x</sub>-based ReRAM Cell [\*]

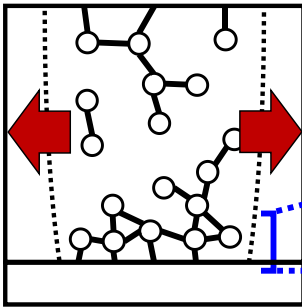


- In Set (Reset) operation, LRS (HRS) is formed by moving O<sup>2-</sup> to TaO<sub>x</sub> layer (CF)
- Percolation paths connect (disconnect) between V<sub>o</sub>s in LRS (HRS)

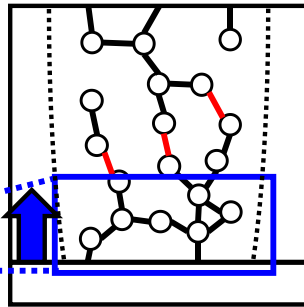


# TaO<sub>x</sub> ReRAM Conductive Filament (CF) Model [\*]

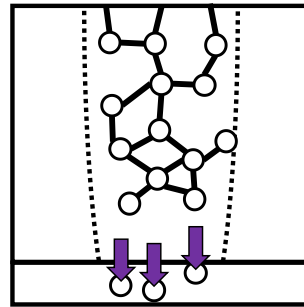
Set/Reset error by **Write-hot data**



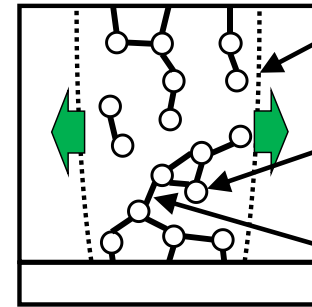
Relaxation effect by long write interval



Read-disturb error by **Read-hot data**



Data-retention error by **Cold data**



Conductive Filament (CF)  
Oxygen vacancy (V<sub>O</sub>)  
Percolation path

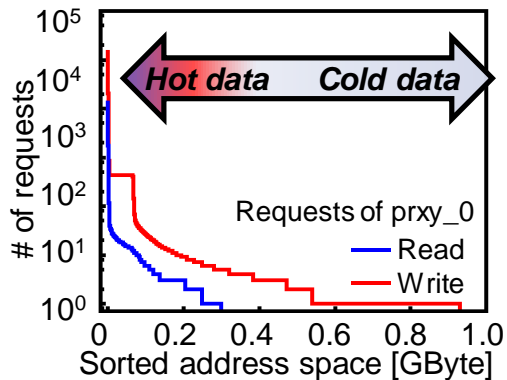
- Write-hot data decrease V<sub>O</sub> density in CF by horizontal diffusion
- Relaxation effect reconnects percolation paths by interface V<sub>O</sub> diffusion to CF
- Read-hot data cause weak reset by vertical V<sub>O</sub> diffusion
- Data retention of cold data causes horizontal V<sub>O</sub> diffusion



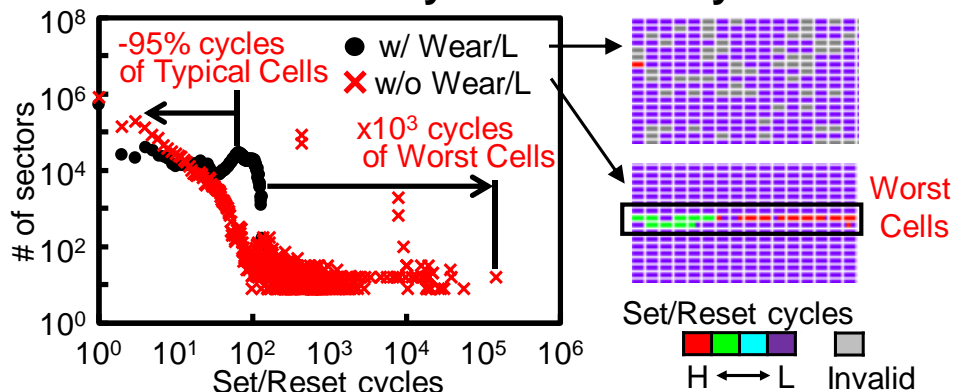


# Storage System Variability

## Non-uniformity of Data Access



## Set/Reset Cycles Variability

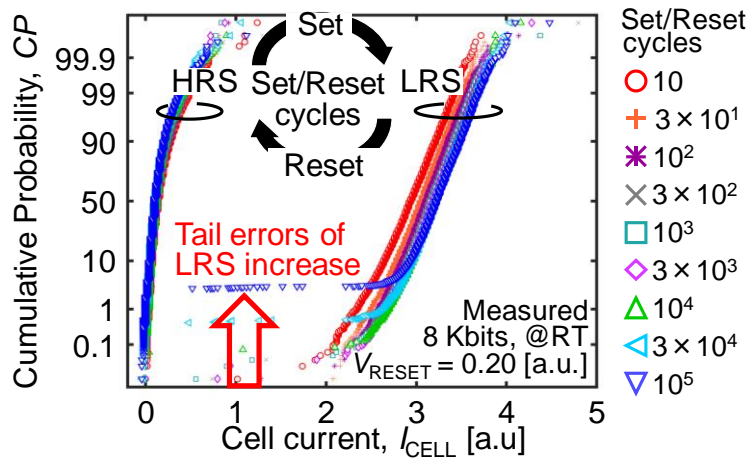


- Write-hot data induce large Set/Reset cycle difference in cells without smoothing by Wear-leveling (Wear/L) [\*]
- Set/Reset cycles of Typical Cells reduce by 95% while those of Worst Cells increase by  $\times 10^3$

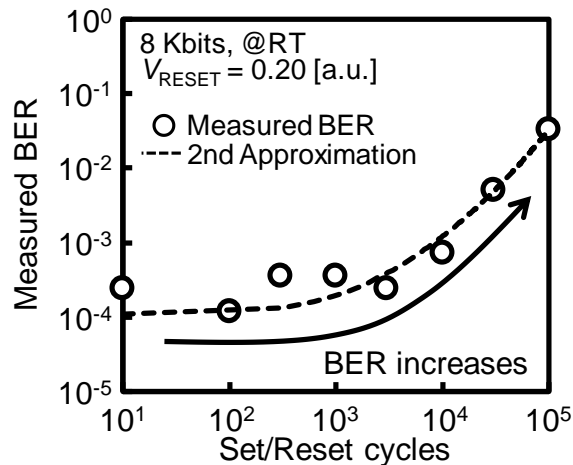


# ReRAM Device-induced BER Variability

### Tail Errors of LRS by Endurance



### Set/Reset BER increase

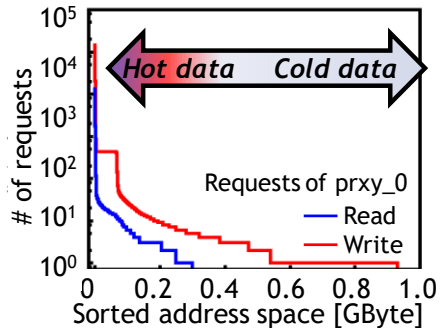


- Measured LRS show tail error cells at high Set/Reset cycles [\*]
- BER increases with Set/Reset cycles

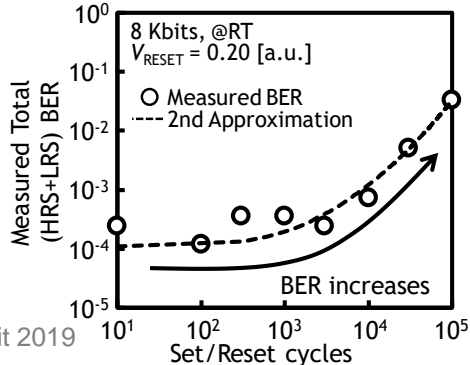


# System-induced BER Variability

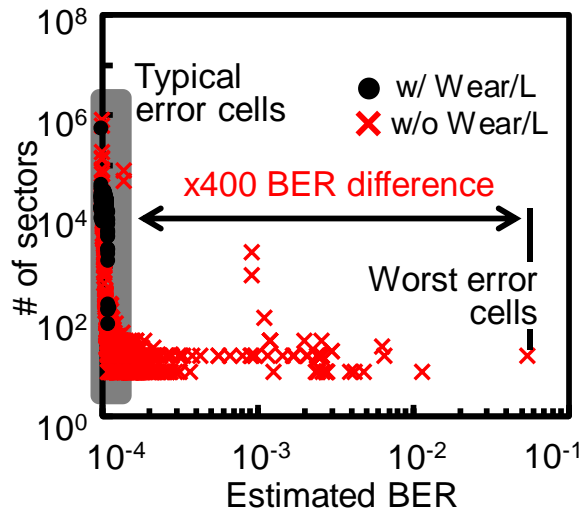
## Application-induced Variability



## Device-induced Variability



## System-induced ReRAM Cell BER Variability



For Variability-Aware  
Approximate Computing (V-AC)

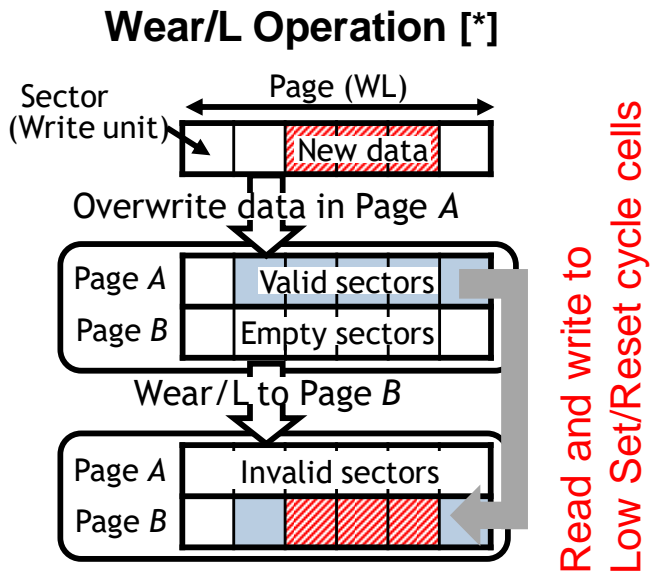


# V-AC Error Tolerant Strategies [\*]

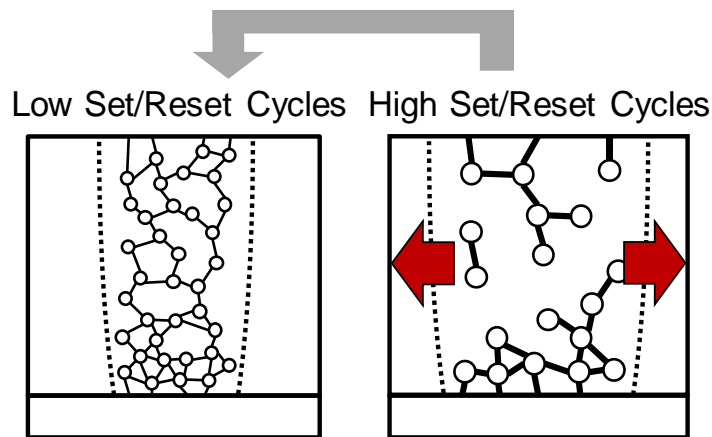
Hierarchy	Operation	Conv. computing	Proposed V-AC Strategy	V-AC Technique	Data characteristic		
					Write-hot	Read-hot	Cold
System	Wear-leveling (Wear/L)	w/ Wear/L	<b>I</b>	w/o Wear/L	✓	✓	
	ECC	Worst-error target (35-bit correction)	<b>II</b>	Typical-error target (5-bit correction)	✓	✓	
Circuit	Read	NA	<b>III</b>	Adaptive Read		✓	
Device	Set/Reset	Verify	<b>IV</b>	w/o Verify	✓		
		NA		Lower $V_{SET}/V_{RESET}$	✓		✓



# Strategy I: Wear-Leveling (Wear/L) Elimination



## Endurance Error Reduction by Wear/L

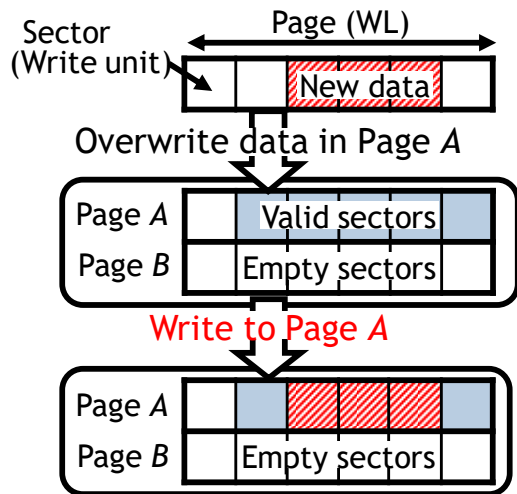


- Wear/L reduces BER of ReRAM storage by smoothing Set/Reset cycles. However, Total Set/Reset cycles increase by extra data copy

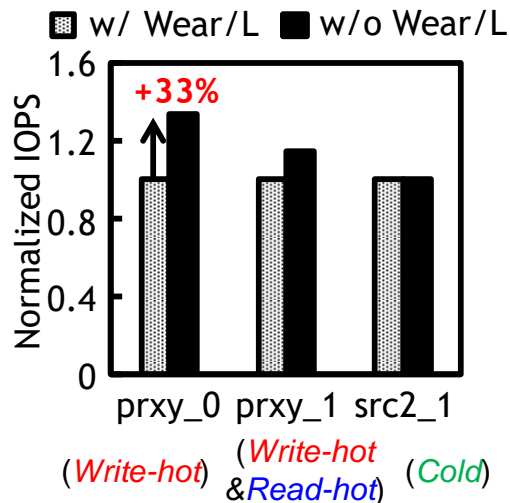


# Strategy I: Wear-Leveling (Wear/L) Elimination

## Wear/L Elimination



## ReRAM Storage Performance



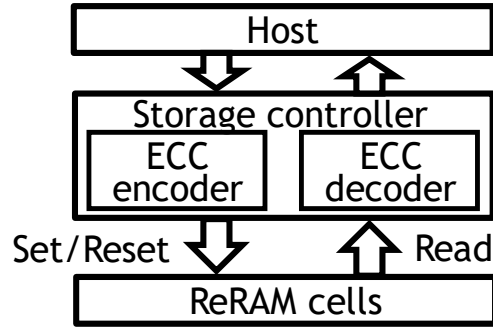
- Strategy I eliminates Wear/L to remove extra data copy and improves storage performance by 33%



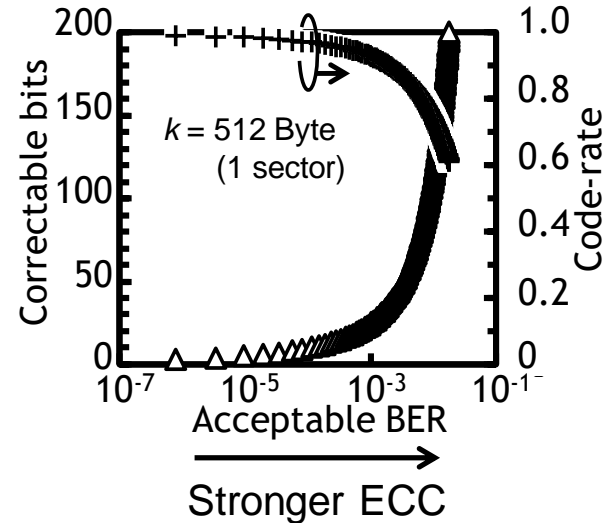
# Strategy II: Typical-Error Target ECC

## ECC Structure

Code word =  $n$ , Code-rate =  $k/n$



## BCH ECC Trade-off

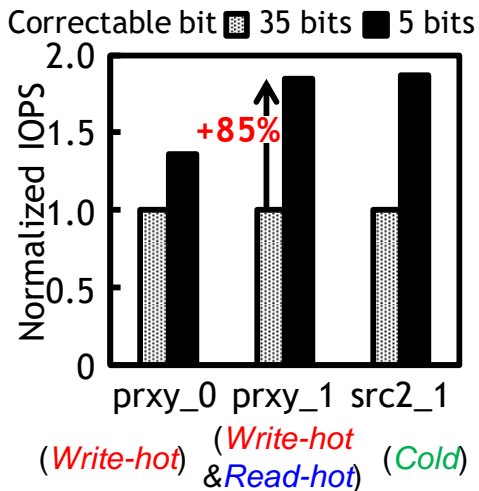


- ECC has trade-off between error correction capability, decoding time and code-rate (cell area)

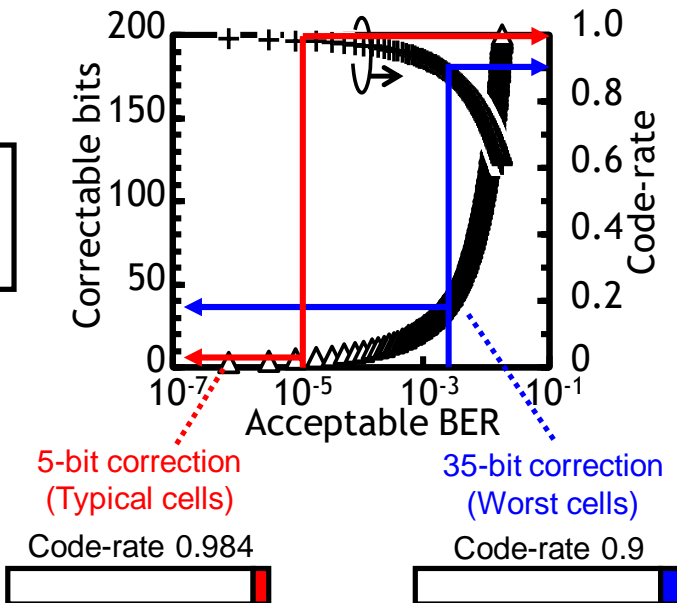


# Strategy II: Typical-Error Target ECC

### ReRAM Storage Performance



### Relaxed Correction Capability



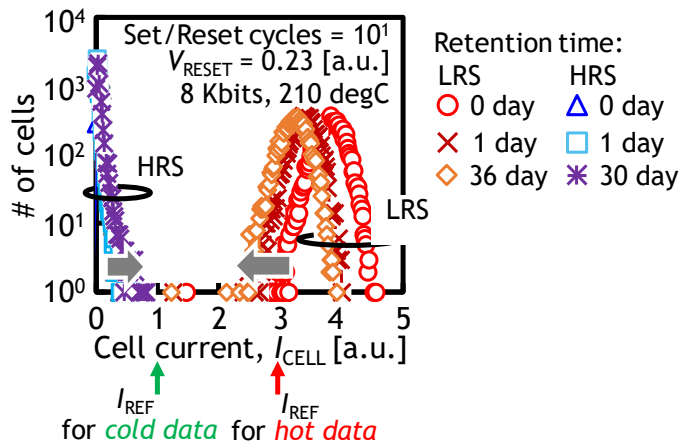
- ECC code-rate increase and performance improves by 85%





# Strategy III & IV: Error Tolerance in Circuit & Device

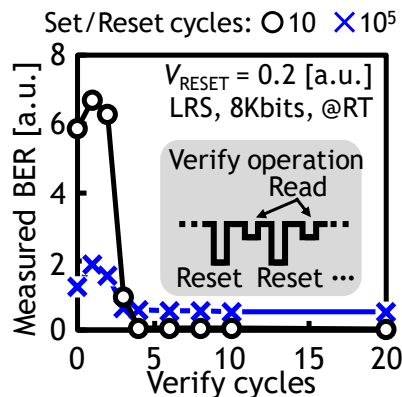
## Strategy III (Circuit): Adaptive Read



**Read Performance**

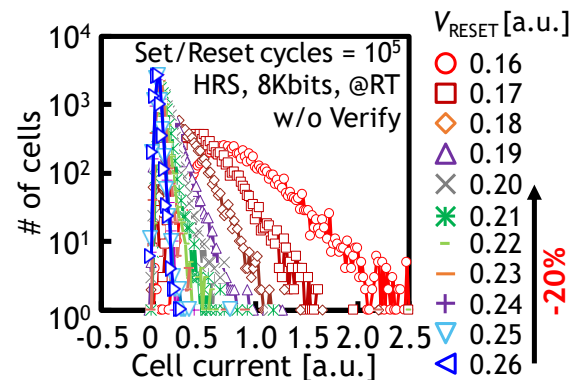
## Strategy IV (Device):

### Verify Elimination



**Write Performance**

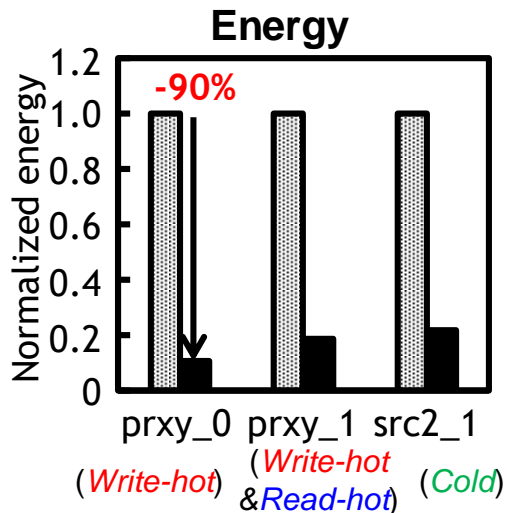
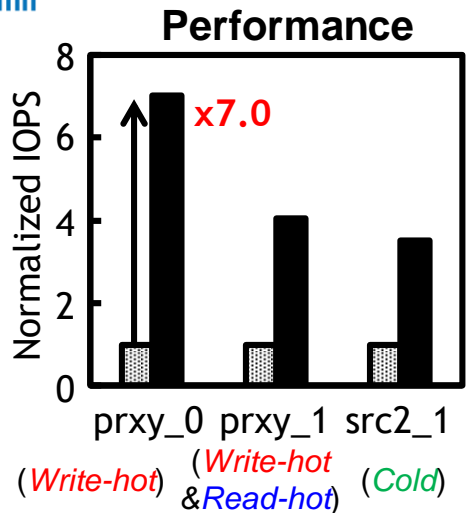
### Lower $V_{SET}/V_{RESET}$



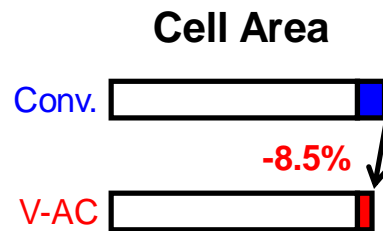
**Selector Scalability**



# Conclusions



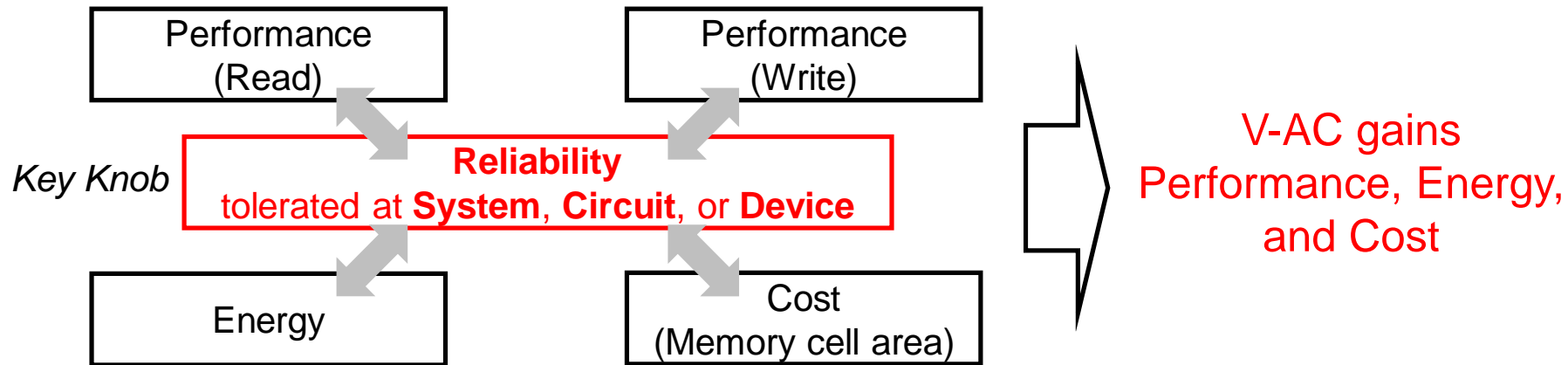
■ Conv. Computing  
■ Proposed V-AC



- Application-induced Variability-aware Approximate Computing (V-AC) is proposed with System, Circuit and Device Co-Design (SCDCD)
- Performance, Energy, and Cell Area of ReRAM storage improve by x7.0, 90%, and 8.5%



# Conclusions



## Thank you for your attention

This presentation is based on results obtained from a project commissioned by the New Energy and Industrial Technology Development Organization (NEDO)