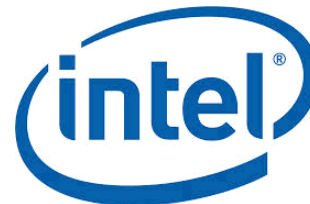# Emerging Non-Volatile Memories- Applications, Challenges and Solutions

## Swaroop Ghosh

School of EECS, The Pennsylvania State University
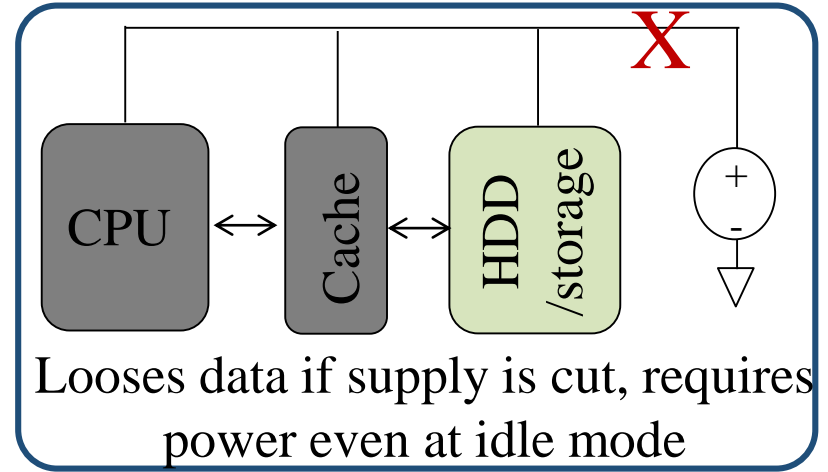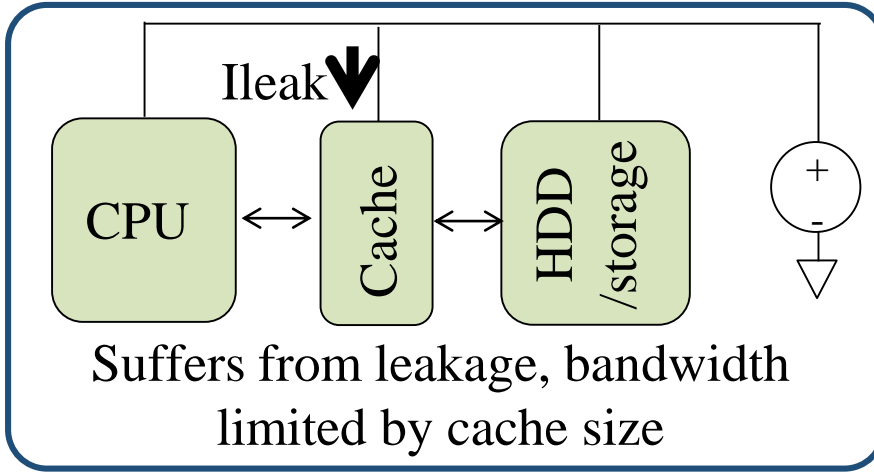
szg212@psu.edu
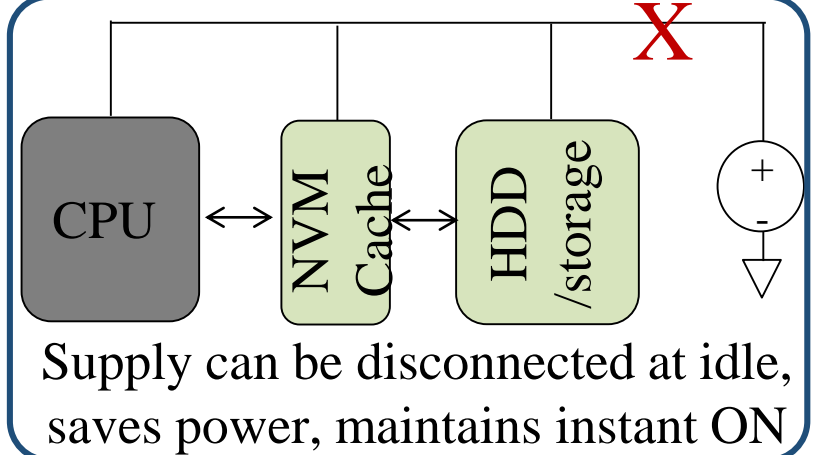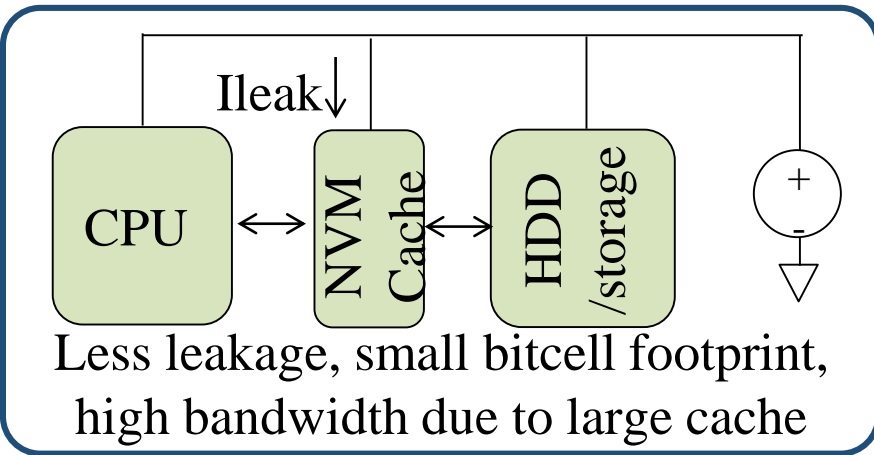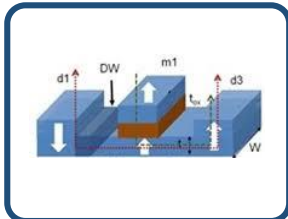
Sponsors

# Why Emerging NVM?

**Conventional**

CPU ↔ Cache ↔ HDD/storage

Ileak ⬇

Suffers from leakage, bandwidth limited by cache size

CPU ↔ Cache ↔ HDD/storage ✗

Looses data if supply is cut, requires power even at idle mode

**Emerging NVM**

CPU ↔ NVM Cache ↔ HDD/storage

Ileak ↓

Less leakage, small bitcell footprint, high bandwidth due to large cache

CPU ↔ NVM Cache ↔ HDD/storage ✗

Supply can be disconnected at idle, saves power, maintains instant ON
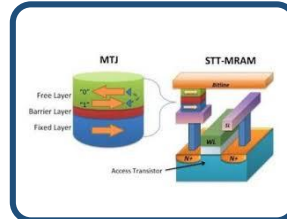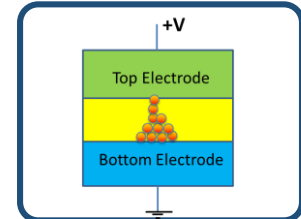
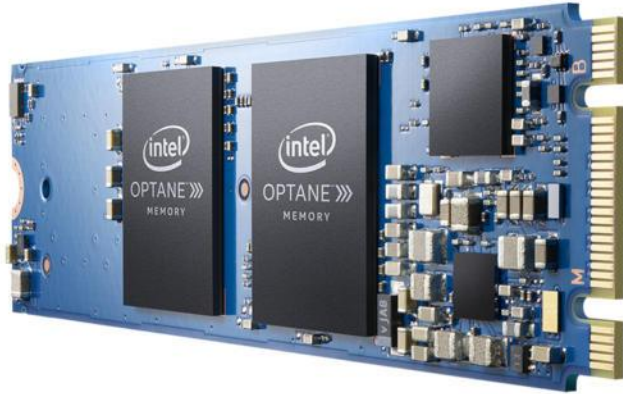**PCRAM**   **DWM**   **FeRAM**   **STTRAM**   **ReRAM**

# Recent Commercialization of Emerging NVMs

## Phase Change RAM*



Intel unveils its Optane hyperfast memory

Intel released few key details around its new non-volatile memory



3D XPoint™ Technology: An Innovative, High-Density Design

## STT- MRAM



Published: *March 9, 2017*

*Everspin unveils a new low latency, PCIe NVMe card based on Spin Torque MRAM*

## ReRAM



3D Resistive RAM as Storage Class Memory

ReRAM is Western Digital's Choice for SCM

Western Digital to Use 3D ReRAM as Storage Class Memory for Special-Purpose SSDs

by Anton Shilov on August 12, 2016 8:00 AM EST

# Emerging Technologies

## STTRAM



- Single bit/cell
- Footprint = ~12-40F$^2$
- Random access
- Read/write latency
  - Read/write of MTJ
- Read: Sensing MTJ resistance
- Write: Flip free layer

## DWM



- Multiple bits/cell
- Footprint = 2.5F$^2$
- Tailored for serial access
- Read/write latency
  - Shift + read/write of MTJ
- Read: Sensing MTJ resistance
- Write: Flip NW domain

## ReRAM



- Multiple bit/cell
- Footprint = ~4F$^2$
- Random access
- Read/write latency
  - Read/write of ReRAM
- Read: Sensing ReRAM resistance
- Write: break or make conductive path

# Outline

- Introduction and motivation

- Applications
  - Last level cache
  - Energy efficient computing
  - Security

- Challenges
  - Retention test
  - Long read/write latency
  - High asymmetric read/write current

- Solutions
  - Retention compression
  - Circuit to system synergistic design
  - Attack sensors and prevention
  - Sensing circuit

- Summary

# Last Level Cache



- Performance improvement: 3-33%

- Power reduction: 1.2X-14.4X

# Digital Signal Processing and Neuro-Inspired Computing



DAC'15, TCAS'16

Ongoing (ISLPED under review)

- **DWM based Viterbi decoder**
  - ◆ 66.4 % area and 59.6 % power savings

- **DWM based 8K point FFT processor**
  - ◆ 60.6 % area and 60.3 % power savings

- **Neuro-inspired computing**
  - ◆ 34% energy savings compared to memristor based computing
  - ◆ Bit-width extendibility

# Energy-Efficient Memory Design

## State retentive sequentials



ISLPED'15,
JETC 2017,
MWSCAS'15

## Non volatile CAMs



## Selector for STTRAM



collaborative /w U. Cin

## Write latency sensor



## Read latency sensor



## Retention sensor

# Spintronics for Security


Collaboration with Univ. of Cincinnati


Modeling (DAC, DT*, JETCAS)

PUF (JETC, ISQED, HOST)

Collaboration with iastate

# ReRAM for Security

## PUF



- Sense circuits in conventional memory architecture employed as arbiter

- Number of CRPs increase exponentially with array size

- Minimally invasive

- 0.13% intra HD and 51.3% inter HD with sufficient response randomness

## TRNG



**NIST Tests**



- Energy/bit of the proposed TRNG is 22.8fJ

R. Govindaraj, ICCD, 2016

# Outline

- Introduction and motivation

- Applications
  - Last level cache
  - Energy efficient computing
  - Security

- Challenges
  - Retention test
  - Long read/write latency
  - High asymmetric read/write current

- Solutions
  - Retention compression
  - Circuit to system synergistic design
  - Attack sensors and prevention
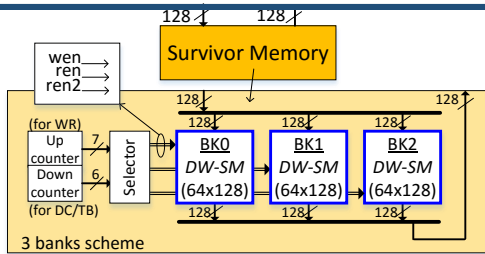  - Sensing circuit

- Summary

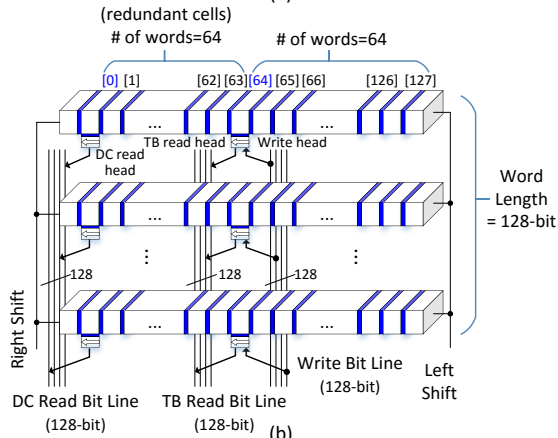# NVM Characteristics-High Retention Time


Retention time Vs Volume

Base MTJ: (40X40X4) nm³
~10 Years
~10 secs


Retention Time Distribution

Min (1) = 69ms
Max (1) = 33.03s
Min (10) = 446ms
Max (10) = 312s
Min (100) = 3.04s
Max (100) = 5357s

μ = 1s
μ = 1s
μ = 10s
μ = 10s
μ = 100s
μ = 100s

~100X - ~1000X



~ 1.1e4X due to PV
~ 2e6X due to Temp

Max   Min


Test Time: Traditional Vs Proposed

Traditional N=1
Proposed N=1
Traditional N=10
Proposed N=10
Traditional N=100
Proposed N=100

~1000X compression

- High test time due to high retention time of bitcell
  - Cannot waive retention characterization

- Test question
  - How to reduce test time of NVMs?
  - How to characterize retention in presence of variations

A. Iyengar, Swaroop Ghosh, S. Srinivasana, "Retention testing of STTRAM", IEEE Design and Test (D&C), 2016

# NVM Characteristics-High Retention Time (Stochastic Variation)



**Stochastic Variation of Mag**

**Stochastic Retention time CDF**

- ■ Random variation in magnetization
  - ◆ Retention time of the same bit changes over time
  - ◆ Require multiple execution of test to guarantee retention time

- ■ Test question
  - ◆ How to identify the worst case retention time?

A. Iyengar, Swaroop Ghosh, S. Srinivasana, "Retention testing of STTRAM", IEEE Design and Test (D&C), 2016

# Data Privacy Issues in NVM



Contactless tampering & scan · Side channel attack · probing

Unauthorized access

**Secure NVM LLC**

Core — CPU | $
Core — CPU | $

Secure bus

Integrity contr. · Privacy contr.

| mag/T sensor | VS-ECC |
| Forced ret. | Bypass |
| Parity encod. | Encryp. |
| Latency sensor | Erase Arch |

SNVM · Encryp engine

Main mem. · Disk

Protected by encryption (Physically exposed)

Physically protected (contained)

## Thermal tampering



Retention Time CDF — Probability distribution vs Log(Retention time)
Legend: -50C, -25C, 0C, 25C, 50C, 75C, 100C
~1e8 X variation

- Persistent data can be accessed between power cycle

- Short retention bitcells can be used to auto-erase the data in clear
  - Freezing of chip can modulate the retention
  - Encryption is latency sensitive

- New features are needed to secure the data

**Table-I  Simulation Parameters**

| Parameter | Value |
| --- | --- |
| Saturation Magnetization (Ms) | 780 Oe |
| Uniaxial Anisotropy (Ku) | 150150 erg/cc |
| Damping Constant ($\alpha$) | 0.007 |
| $\Delta$ for Tret of 1s, 10s, 100s | 20.73, 23.02 & 25.33 |
| Length and Width | 40nmX40nm |

14

# NVM Characteristics- Sensitivity to Ambient Parameters

## Impact of DC field



Flipping times of different Oe w/o current

**Flip time reduces**

## Impact of DC field



x 10⁻¹⁰ Flipping times of different Oe w/ neg current

**Fast flipping due to assistive STT**

Legend: 50uA, 100uA, 150uA, 200uA, 250uA, noCurr

## Impact of DC field



x 10⁻¹⁰ Flipping times of different Oe w/ pos current

**Slow flipping due to suppressive STT**

Legend: 50uA, 100uA, 150uA, 200uA, 250uA, noCurr

## Impact of AC field



x 10⁻⁹ Flipping times of different Oe w/ diff Freq

**Impact of field magnitude**

**Impact of frequency**

Legend: 2GHz, 1GHz, 750MHz, 500MHz, 450MHz, 350MHz, 250MHz, 150MHz

Jae-won Jang, Jongsun Park, Swaroop Ghosh, Swarup Bhunia, "Self-Correcting STTRAM under Magnetic Field Attacks", IEEE Design Automation Conference (DAC), 2015

- Test challenge
  - How to characterize magnetic tolerance

# NVM Characteristics- Sensitivity to Ambient Parameters

Reference voltage shift





Experimental validation (HOST'16 demo)





http://www.hostsymposium.org/hardware-demo-list.php

- **Test question**
  - ◆ How to characterize NVM under sensitivities
  - ◆ Can we detect security attacks?

# NVM Characteristics-High and Asymmetric Write and Read Current



- **High write current triggers droop**
  - Depends on pattern

- **Test question**
  - Identifying test pattern to validate worst case droop

R. Aluru, Swaroop Ghosh, "Droop mitigating last level STTRAM cache", DATE, 2017

# Security Implications- Privacy

## Asymmetric WR current/latency



Write:0→1    Write:1→0

## 4-bit Write



4 bit without Parity

## Asymmetric RD current



Read current for a 1-bit read operation



Off-chip Voltage Regulator

$I_{DIE}$   +V−

$I_{CPU}$   $I_{LLC}$

Key is recovered!

$K_{0,0}$ .... $K_{0,3}$
$K_{3,0}$ .... $K_{3,3}$

Plain Text
Key → AES → Cipher Text → STTRAM LLC
Processor



Attack Success Rate Comparison

- Test question
  - ◆ Characterize asymmetric write current

N. Rathi, S Ghosh, H. Naeimi, "Side Channel Attacks on STTRAM and Low-Overhead Countermeasures ", DFT 2016

# NVM Characteristics-High and Asymmetric Write and Read Latency



High Write and Read current

Long Write and Read Latency

- Long tail of read and write latency

- Test question
  - How to characterize write and read latency at fast test time?

# Outline

- Introduction and motivation

- Applications
  - Last level cache
  - Energy efficient computing
  - Security

- Challenges
  - Retention test
  - Long read/write latency
  - High asymmetric read/write current

- Solutions
  - Retention compression
  - Circuit to system synergistic design
  - Attack sensors and prevention
  - Sensing circuit

- Summary

# Retention Testing using Test Time Compression



Retention Time Vs Thermal Barrier

Retention Time CDF

$\sim 10^{13} X$ variation

Volume$_{MTJ}$ 1.04*10$^{-17}$ cm$^3$

0Oe, 50 Oe, 100 Oe, 150 Oe, 200 Oe, 220 Oe

Retention Time CDF

25C, 50C, 75C, 100C, 125C

$\sim 10^4 X$ variation

Base MTJ 1.04*10$^{-17}$ cm$^3$ 3$\sigma$ of 5%

Retention Time Distribution

MBI+BI, MBI, Base

- Exploit STTRAM sensitivity to compress retention time

- Test time with lower retention is low

M. N.,I. Khan A. Iyengar, Swaroop Ghosh, "Magnetic burn-in for STTRAM retention testing", DATE, 2017

# Retention Testing using Test Time Compression



M. N.,I. Khan A. Iyengar, Swaroop Ghosh, "Magnetic burn-in for STTRAM retention testing", DATE, 2017

# Energy-Efficient Memory Design

| Layout |
|---|

Bit-Cell Layout ⟷ Head Positioning ⟷ Utilization Factor ⟷ Sharing of diffusion, bitlines and shift lines

| Circuit |
|---|

Merged Read-Write Head ⟷ Access transistor sizing | Shift gating & WL strapping ⟷ Shift Circuit & Write Driver with three operating modes

| uArch |
|---|

Cache Organization ⟷ Cache segregation with novel replacement policy

| System |
|---|

Workload aware current boosting

ISLPED'14, TNANO'15, DATE'15

## Layout/circuit opt

② Area    Notch wid/dep    ③ d    t
$t_{ox}$ ①    TMR ④    h ⑤    ⑥ $V_{th}$

TMR(Tunnel Magnetic Ratio)=(RH-RL)/RL

Latency Vs Power

Fast    Med    Slow    Slow    Med    Fast

$\overline{FSR}$  MSR  $\overline{SSR}$    Shift Circuit    $\overline{SSL}$  MSL  $\overline{FSL}$

SL    SR

## Cache segregation

L2 Access

Hit? — No → Fetch block from main memory and replace with LRU block

Yes

Fast Way? — No

Yes → Access granted, Mark it as MRU

Medium Way? — No → Replace block with LRU block in Medium way and mark it as MRU

Yes → Replace block with LRU block in Fast ways and mark it as MRU

## Workload adaptive modulation

Number of L2 accesses (×10⁵)

Boost-up

Workload
$Th_H$
$Th_L$
Comparator
Boost enable
L2-Cache

$TH_H$
$TH_L$    No Boost

Boost-down

Run Time

- **~30% perf improvement**
- **>10X power saving**

23

# Energy-Efficient Memory Design

## Slope sensing circuit



ISLPED'15
(collaborative /w Intel)

## Sensing with column voltage boosting



CICC'17 (under review, collaborative /w Intel)

# Magnetic Field Sensor



**Flip Times for diff Volumes w/o Current**

100Oe, 150Oe, 200Oe, 220Oe, 230Oe, 280Oe, 350Oe, 400Oe

-33%, 2X

Flipping time (ps) vs Volume (um³)

**Flip time difference 3nm to 2nm**

No Curr, 50uA, 100uA

400us, 100us, 1us, 2X, 80X

Flipping time difference (s) vs Oersted Fields (Oe), 50Oe

- **Key requirements**
  - ◆ Proactive sensing
  - ◆ Sense magnitude and polarity

- **Sensor design**
  - ◆ Small volume for early sensing
  - ◆ Weak write of sensor array to fail early

- **Challenges**
  - ◆ Identifying false alarms
  - ◆ Power consumption in sensor



SA0, SA1, Midlogic, SA2, SA3

Sensor array

Control logic for sensor array
(address /read/write/biasing)

Write driver, rden, BL, S, L, W, L, Senseamp, wren, rden, Wordline (WL) driver, WL bias

# Prevention (1)- Stalling

- Stall the CPU and wait till the attack is over

- For gradually ramping attack
  - Better than shutting down the entire system
  - Will not work for sudden attack since dirty data is corrupted

| Attack launched | Attack sensed | CPU Stall | Attack ended | Resume |
|---|---|---|---|---|

Writeback done     Discard LLC

- For sudden attack
  - Processor is restarted after the attack
  - Applications can resume from application level checkpointing

| Attack launched | Attack sensed | CPU Stall | Attack ended | Restart |
|---|---|---|---|---|

- Both approaches disable computations during attack
  - Attacker can exploit these features to drain the battery

# Prevention (2)- Cache Bypassing



- **Key ideas**
  - ◆ Since LLC is under attack, bypass it
  - ◆ Perform computation seamlessly without LLC
  - ◆ Update the main memory before starting bypass
  - ◆ Invalidate LLC before exiting bypass

# Prevention (3)- Checkpointing



■ **Key ideas**
  ◆ Save processor state and update main memory periodically
  ◆ If attack, go back to last saved state & start with LLC bypass
  ◆ Can handle sudden corruption of memory

■ **Challenges**
  ◆ Need to stop main memory writeback between checkpoints
  ◆ Performance loss due to checkpoint which depends on
    ⇨ Epoch
    ⇨ LLC full

# Protecting Data Privacy



**CLK**

**Reset** → **CPU**    $\overline{WL} = 1$

**Address**

**Tag | Index**    **MUX1**    ②    **Tag Array**    **Data Array**

**Counter**    **Tag | Index**    **Index** Valid Tag    **Index**
ER → **EN**                            0                    0
**RST** **END**    ①    ⑤    1    ⑤    1
ER    ⑦    2    2

**Counter**    ⑥    ⑧    **Data Pattern**
ER → **EN**
**RST**

**Write Controller**    **=**    **Hit/Miss**    **DECODER**    ④
ER → **EN**    ③ $\overline{ER}$
**RST**    **MUX2**    **ER** →    **Main Memory**

**Valid Erasure**

① **ER = 1**    ⑤ **Tag = 0, Valid = 1**
② **Tag & Index made 0**    ⑥ **Valid erasure**
③ $\overline{ER}$**=0. Force Miss**    ⑦ **Location selected**
④ **MM is bypassed**    ⑧ **Valid erased**

## Adding dummy bits



4 bit with Even Parity

1-1-1-1 / 0
0-1-1-1 / 1
0-0-0-1 / 1
0-0-1-1 / 0

0-0-0-0 / 0

Current (mA): 0, -2, -4, -6
Write Latency (ns): 0, 0.4, 0.8



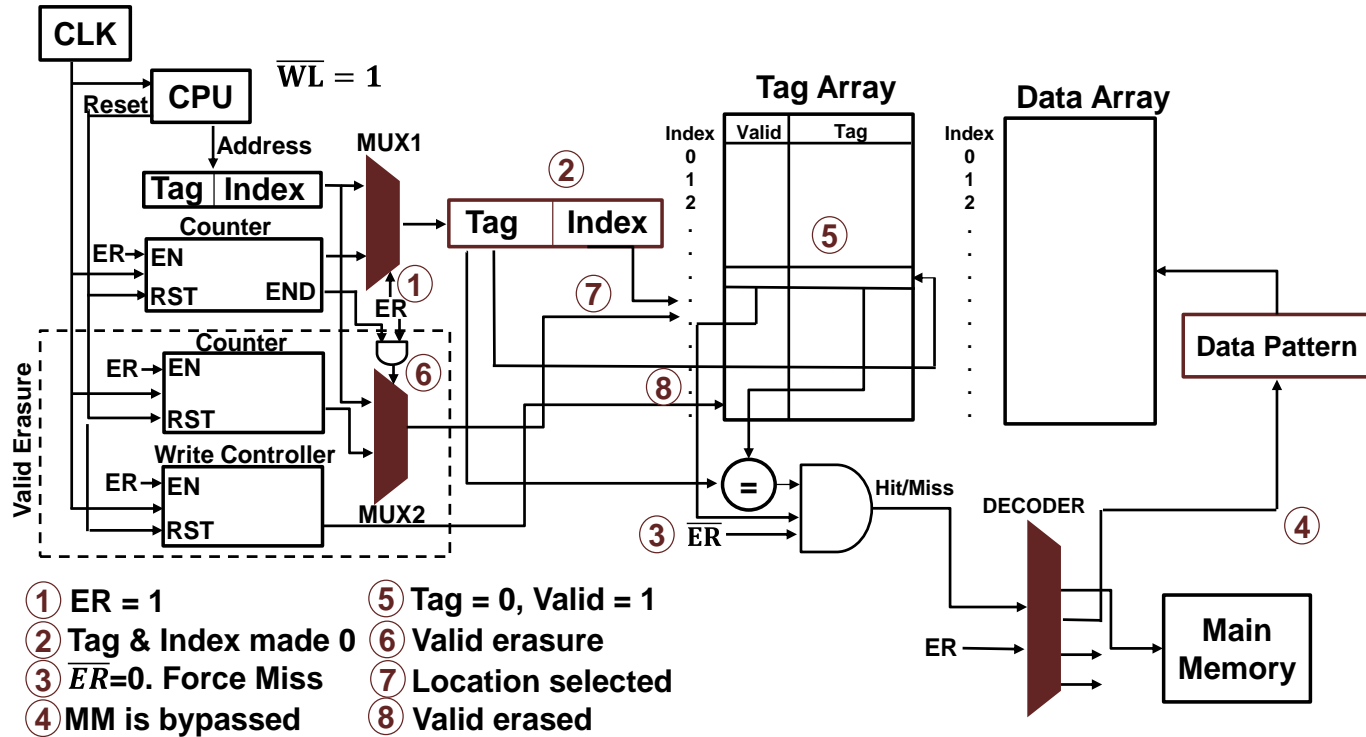% reduction in states vs Word Size (bits)
Legend: 2-bit, 3-bit, 4-bit

## Constant Current Write



Current Mirror

Write '1'    mismatch (0.4ns)

Magnetization

Write '0'

MTJ    $M_C$    $V_B$

Write Latency (ns): 0, 0.4, 0.8, 1.2

Nitin Rathi, Swaroop Ghosh, Anirudh Iyengar and Helia Naeimi, "Data Privacy in Non-Volatile Cache: Challenges, Attack Models and Solutions", ASPDAC, 2016.

# Conclusions

- Emerging NVMs are promising for broad range of applications

- NVMs possess unique challenges that could be design and security issues

- We proposed novel techniques to solve the challenges

- Proposed solutions are also applicable to other NVMs

# **Thank You!**

## **Acknowledgements**

LOGICS lab students, collaborators from Intel, Nanyang Tech Univ (NTU), Univ. of Florida, Iowa State, Univ. of Cincinnati and Korea Univ.

# Graduate Students

**PhD**

A. Iyengar('13)  H. Motaman('13)     A. Saki ('17)     R. Govindaraj('14)     J. Jang('15)

P:7, Pat:3          P:6, Pat:1          P:0, Pat:0          P:3+2*, Pat:1          P:4, Pat:1

D. Vontela    C. Lin         R. Aluru         I. Reddy         Md. N. Khan ('16)     A. De ('16)

**MS**

P:1+2*, Pat:0  P:3, Pat:2  P:1*, Pat:0     P:2, Pat:1          P:1*, Pat:0          P:1*, Pat:0

- **Graduated**
  - ◆ 5 MS: Kenneth Ramclam, Jae-won Jang, Radha Aluru, Deepak Vontela, Ithihasa Reddy
  - ◆ Published more than 40 IEEE papers in last 5 years

*LOGICS: Lab. Of Green and secure Integrated Circuits and Systems