

# Scalable Many-Core Memory Systems

## Topic 2: Emerging Technologies and Hybrid Memories

Prof. Onur Mutlu

<http://www.ece.cmu.edu/~omutlu>

[onur@cmu.edu](mailto:onur@cmu.edu)

HiPEAC ACACES Summer School 2013

July 15-19, 2013

**Carnegie Mellon**

# What Will You Learn in This Course?

---

- Scalable Many-Core Memory Systems
  - July 15-19, 2013
- Topic 1: Main memory basics, DRAM scaling
- Topic 2: Emerging memory technologies and hybrid memories
- Topic 3: Main memory interference and QoS
- Topic 4 (unlikely): Cache management
- Topic 5 (unlikely): Interconnects
- Major Overview Reading:
  - Mutlu, “[Memory Scaling: A Systems Architecture Perspective](#),” IMW 2013.

# Readings and Videos

# Memory Lecture Videos

---

- Memory Hierarchy (and Introduction to Caches)

- <http://www.youtube.com/watch?v=JBdfZ5i21cs&list=PL5PHm2jkkXmidJOd59REog9jDnPDTG6IJ&index=22>

- Main Memory

- <http://www.youtube.com/watch?v=ZLCy3pG7Rc0&list=PL5PHm2jkkXmidJOd59REog9jDnPDTG6IJ&index=25>

- Memory Controllers, Memory Scheduling, Memory QoS

- <http://www.youtube.com/watch?v=ZSotvL3WXmA&list=PL5PHm2jkkXmidJOd59REog9jDnPDTG6IJ&index=26>
- [http://www.youtube.com/watch?v=1xe2w3\\_NzmI&list=PL5PHm2jkkXmidJOd59REog9jDnPDTG6IJ&index=27](http://www.youtube.com/watch?v=1xe2w3_NzmI&list=PL5PHm2jkkXmidJOd59REog9jDnPDTG6IJ&index=27)

- Emerging Memory Technologies

- <http://www.youtube.com/watch?v=LzfOghMKyA0&list=PL5PHm2jkkXmidJOd59REog9jDnPDTG6IJ&index=35>

- Multiprocessor Correctness and Cache Coherence

- <http://www.youtube.com/watch?v=U-VZKMgItDM&list=PL5PHm2jkkXmidJOd59REog9jDnPDTG6IJ&index=32>

# Readings for Topic 1 (DRAM Scaling)

---

- Lee et al., "Tiered-Latency DRAM: A Low Latency and Low Cost DRAM Architecture," HPCA 2013.
- Liu et al., "RAIDR: Retention-Aware Intelligent DRAM Refresh," ISCA 2012.
- Kim et al., "A Case for Exploiting Subarray-Level Parallelism in DRAM," ISCA 2012.
- Liu et al., "An Experimental Study of Data Retention Behavior in Modern DRAM Devices," ISCA 2013.
- Seshadri et al., "RowClone: Fast and Efficient In-DRAM Copy and Initialization of Bulk Data," CMU CS Tech Report 2013.
- David et al., "Memory Power Management via Dynamic Voltage/Frequency Scaling," ICAC 2011.
- Ipek et al., "Self Optimizing Memory Controllers: A Reinforcement Learning Approach," ISCA 2008.

# Readings for Topic 2 (Emerging Technologies)

---

- Lee, Ipek, Mutlu, Burger, “Architecting Phase Change Memory as a Scalable DRAM Alternative,” ISCA 2009, CACM 2010, Top Picks 2010.
- Qureshi et al., “Scalable high performance main memory system using phase-change memory technology,” ISCA 2009.
- Meza et al., “Enabling Efficient and Scalable Hybrid Memories,” IEEE Comp. Arch. Letters 2012.
- Yoon et al., “Row Buffer Locality Aware Caching Policies for Hybrid Memories,” ICCD 2012 Best Paper Award.
- Meza et al., “A Case for Efficient Hardware-Software Cooperative Management of Storage and Memory,” WEED 2013.
- Kultursay et al., “Evaluating STT-RAM as an Energy-Efficient Main Memory Alternative,” ISPASS 2013.

# Readings for Topic 3 (Memory QoS)

---

- Moscibroda and Mutlu, "Memory Performance Attacks," USENIX Security 2007.
- Mutlu and Moscibroda, "Stall-Time Fair Memory Access Scheduling," MICRO 2007.
- Mutlu and Moscibroda, "Parallelism-Aware Batch Scheduling," ISCA 2008, IEEE Micro 2009.
- Kim et al., "ATLAS: A Scalable and High-Performance Scheduling Algorithm for Multiple Memory Controllers," HPCA 2010.
- Kim et al., "Thread Cluster Memory Scheduling," MICRO 2010, IEEE Micro 2011.
- Muralidhara et al., "Memory Channel Partitioning," MICRO 2011.
- Ausavarungnirun et al., "Staged Memory Scheduling," ISCA 2012.
- Subramanian et al., "MISE: Providing Performance Predictability and Improving Fairness in Shared Main Memory Systems," HPCA 2013.
- Das et al., "Application-to-Core Mapping Policies to Reduce Memory System Interference in Multi-Core Systems," HPCA 2013.

# Readings for Topic 3 (Memory QoS)

---

- Ebrahimi et al., “Fairness via Source Throttling,” ASPLOS 2010, ACM TOCS 2012.
- Lee et al., “Prefetch-Aware DRAM Controllers,” MICRO 2008, IEEE TC 2011.
- Ebrahimi et al., “Parallel Application Memory Scheduling,” MICRO 2011.
- Ebrahimi et al., “Prefetch-Aware Shared Resource Management for Multi-Core Systems,” ISCA 2011.

# Readings in Flash Memory

---

- Yu Cai, Gulay Yalcin, Onur Mutlu, Erich F. Haratsch, Adrian Cristal, Osman Unsal, and Ken Mai, **"Error Analysis and Retention-Aware Error Management for NAND Flash Memory"** *Intel Technology Journal (ITJ) Special Issue on Memory Resiliency*, Vol. 17, No. 1, May 2013.
- Yu Cai, Erich F. Haratsch, Onur Mutlu, and Ken Mai, **"Threshold Voltage Distribution in MLC NAND Flash Memory: Characterization, Analysis and Modeling"** *Proceedings of the Design, Automation, and Test in Europe Conference (DATE)*, Grenoble, France, March 2013. [Slides \(ppt\)](#)
- Yu Cai, Gulay Yalcin, Onur Mutlu, Erich F. Haratsch, Adrian Cristal, Osman Unsal, and Ken Mai, **"Flash Correct-and-Refresh: Retention-Aware Error Management for Increased Flash Memory Lifetime"** *Proceedings of the 30th IEEE International Conference on Computer Design (ICCD)*, Montreal, Quebec, Canada, September 2012. [Slides \(ppt\)](#) [\(pdf\)](#)
- Yu Cai, Erich F. Haratsch, Onur Mutlu, and Ken Mai, **"Error Patterns in MLC NAND Flash Memory: Measurement, Characterization, and Analysis"** *Proceedings of the Design, Automation, and Test in Europe Conference (DATE)*, Dresden, Germany, March 2012. [Slides \(ppt\)](#)

# Online Lectures and More Information

---

## ■ Online Computer Architecture Lectures

- <http://www.youtube.com/playlist?list=PL5PHm2jkkXmidJOd59REog9jDnPDTG6IJ>

## ■ Online Computer Architecture Courses

- Intro: <http://www.ece.cmu.edu/~ece447/s13/doku.php>
- Advanced: <http://www.ece.cmu.edu/~ece740/f11/doku.php>
- Advanced: <http://www.ece.cmu.edu/~ece742/doku.php>

## ■ Recent Research Papers

- <http://users.ece.cmu.edu/~omutlu/projects.htm>
- <http://scholar.google.com/citations?user=7XyGUGkAAAAJ&hl=en>

# Emerging Memory Technologies

# Agenda

---

- Major Trends Affecting Main Memory
- Requirements from an Ideal Main Memory System
- Opportunity: Emerging Memory Technologies
- Conclusions
- Discussion

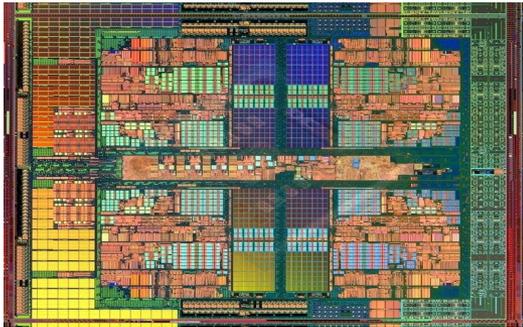
# Major Trends Affecting Main Memory (I)

---

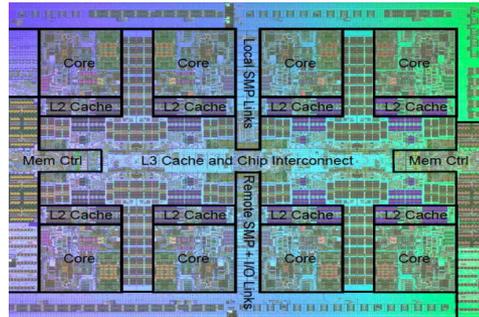
- Need for main memory capacity and bandwidth increasing
- Main memory energy/power is a key system design concern
- DRAM technology scaling is ending

# Demand for Memory Capacity

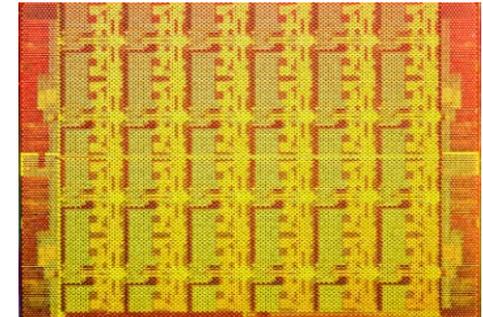
- More cores → More concurrency → Larger working set



AMD Barcelona: 4 cores



IBM Power7: 8 cores

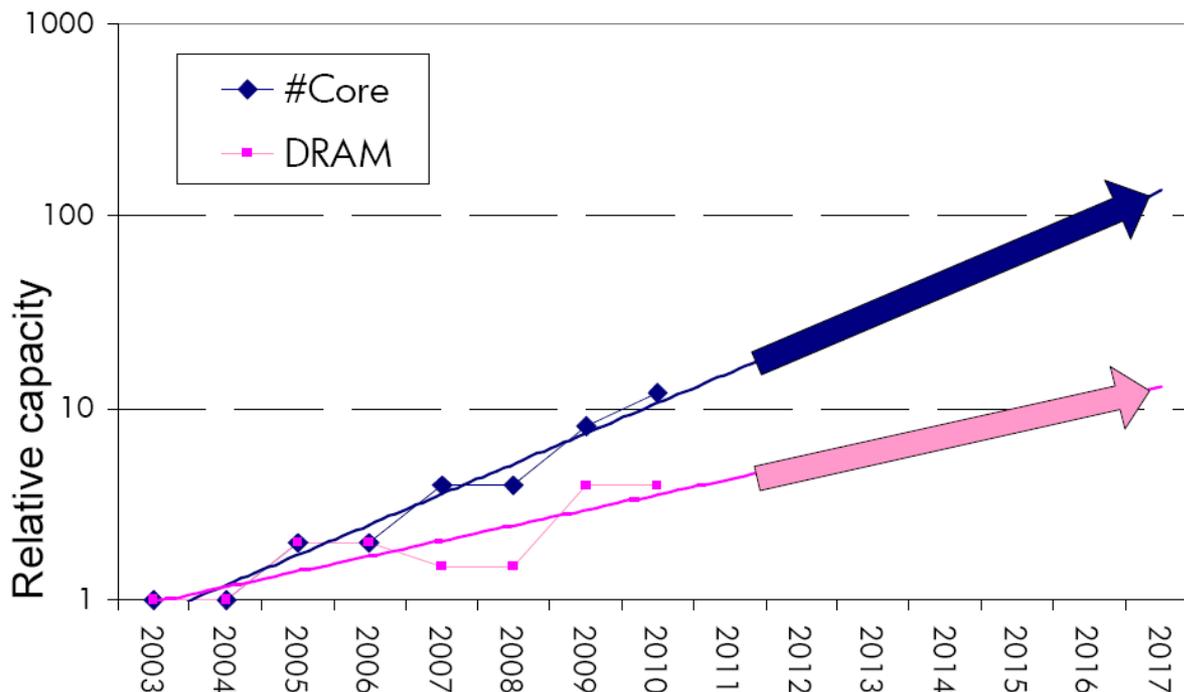


Intel SCC: 48 cores

- Emerging applications are data-intensive
- Many applications/virtual machines (will) share main memory
  - Cloud computing/servers: Consolidation to improve efficiency
  - GP-GPUs: Many threads from multiple parallel applications
  - Mobile: Interactive + non-interactive consolidation

# The Memory Capacity Gap

Core count doubling ~ every 2 years  
DRAM DIMM capacity doubling ~ every 3 years



Source: Lim et al., ISCA 2009.

- Memory capacity per core expected to drop by 30% every two years

# Major Trends Affecting Main Memory (II)

---

- Need for main memory capacity and bandwidth increasing
  - **Multi-core**: increasing number of cores
  - **Data-intensive applications**: increasing demand/hunger for data
  - **Consolidation**: Cloud computing, GPUs, mobile
  
- Main memory energy/power is a key system design concern
  
  
  
  
  
  
  
  
  
  
- DRAM technology scaling is ending

# Major Trends Affecting Main Memory (III)

---

- Need for main memory capacity and bandwidth increasing
- Main memory energy/power is a key system design concern
  - IBM servers: ~50% energy spent in off-chip memory hierarchy [Lefurgy, IEEE Computer 2003]
  - DRAM consumes power when idle and needs periodic refresh
- DRAM technology scaling is ending

# Major Trends Affecting Main Memory (IV)

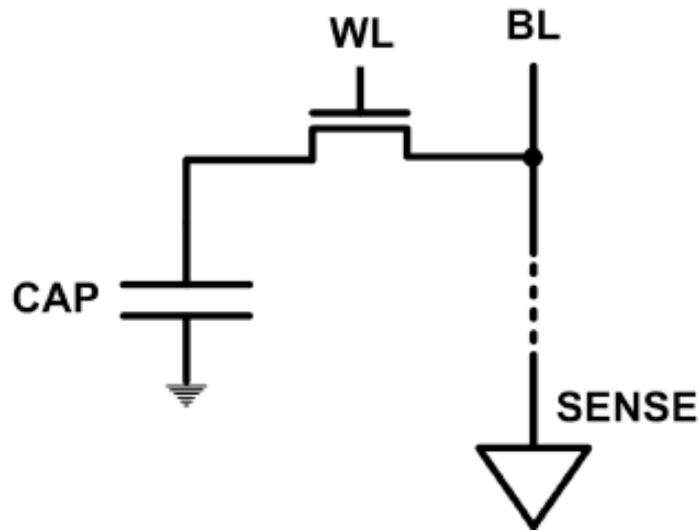
---

- Need for main memory capacity and bandwidth increasing
- Main memory energy/power is a key system design concern
- DRAM technology scaling is ending
  - ITRS projects DRAM will not scale easily below 40nm
  - Scaling has provided many benefits:
    - higher capacity, higher density, lower cost, lower energy

# The DRAM Scaling Problem

---

- DRAM stores charge in a capacitor (charge-based memory)
  - Capacitor must be large enough for reliable sensing
  - Access transistor should be large enough for low leakage and high retention time
  - Scaling beyond 40-35nm (2013) is challenging [ITRS, 2009]



- DRAM capacity, cost, and energy/power hard to scale

# Trends: Problems with DRAM as Main Memory

---

- Need for main memory capacity and bandwidth increasing
  - DRAM capacity hard to scale
  
- Main memory energy/power is a key system design concern
  - DRAM consumes high power due to leakage and refresh
  
- DRAM technology scaling is ending
  - DRAM capacity, cost, and energy/power hard to scale

# Agenda

---

- Major Trends Affecting Main Memory
- Requirements from an Ideal Main Memory System
- Opportunity: Emerging Memory Technologies
- Conclusions
- Discussion

# Requirements from an Ideal Memory System

---

- Traditional
  - Enough capacity
  - Low cost
  - High system performance (high bandwidth, low latency)
  
- New
  - Technology scalability: lower cost, higher capacity, lower energy
  - Energy (and power) efficiency
  - QoS support and configurability (for consolidation)

# Requirements from an Ideal Memory System

---

## ■ Traditional

- ❑ Higher capacity
- ❑ Continuous low cost
- ❑ High system performance (**higher bandwidth**, low latency)

## ■ New

- ❑ Technology scalability: lower cost, higher capacity, lower energy
- ❑ Energy (and power) efficiency
- ❑ QoS support and configurability (for consolidation)

**Emerging, resistive memory technologies (NVM) can help**

---

# Agenda

---

- Major Trends Affecting Main Memory
- Requirements from an Ideal Main Memory System
- Opportunity: Emerging Memory Technologies
- Conclusions
- Discussion

# The Promise of Emerging Technologies

---

- Likely need to replace/augment DRAM with a technology that is
  - Technology scalable
  - And at least similarly efficient, high performance, and fault-tolerant
    - or can be architected to be so
  
- Some emerging resistive memory technologies appear promising
  - Phase Change Memory (PCM)?
  - Spin Torque Transfer Magnetic Memory (STT-MRAM)?
  - Memristors?
  - And, maybe there are other ones
  - Can they be enabled to replace/augment/surpass DRAM?

# Agenda

---

- Major Trends Affecting Main Memory
- Requirements from an Ideal Main Memory System
- Opportunity: Emerging Memory Technologies
  - Background
  - PCM (or Technology X) as DRAM Replacement
  - Hybrid Memory Systems
- Conclusions
- Discussion

# Charge vs. Resistive Memories

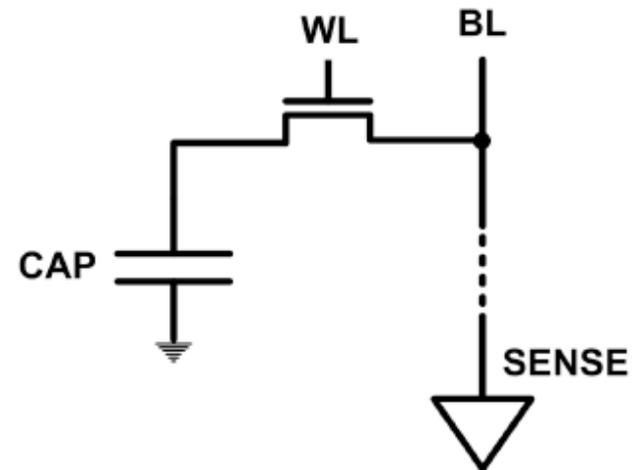
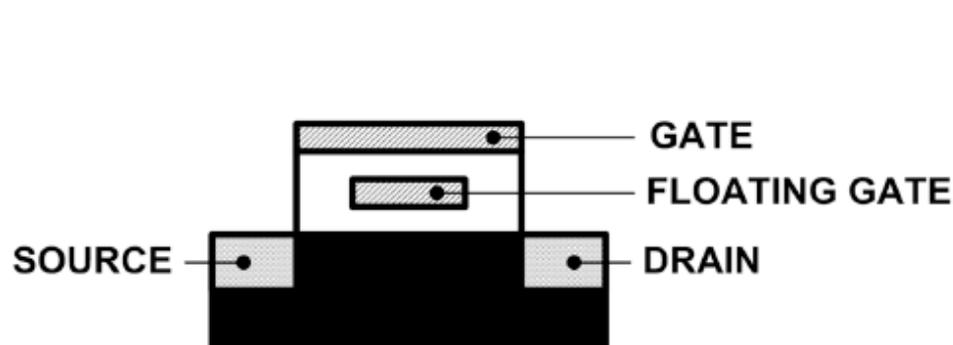
---

- Charge Memory (e.g., DRAM, Flash)
  - Write data by capturing charge  $Q$
  - Read data by detecting voltage  $V$
  
- Resistive Memory (e.g., PCM, STT-MRAM, memristors)
  - Write data by pulsing current  $dQ/dt$
  - Read data by detecting resistance  $R$

# Limits of Charge Memory

---

- Difficult charge placement and control
  - Flash: floating gate charge
  - DRAM: capacitor charge, transistor leakage
- Reliable sensing becomes difficult as charge storage unit size reduces



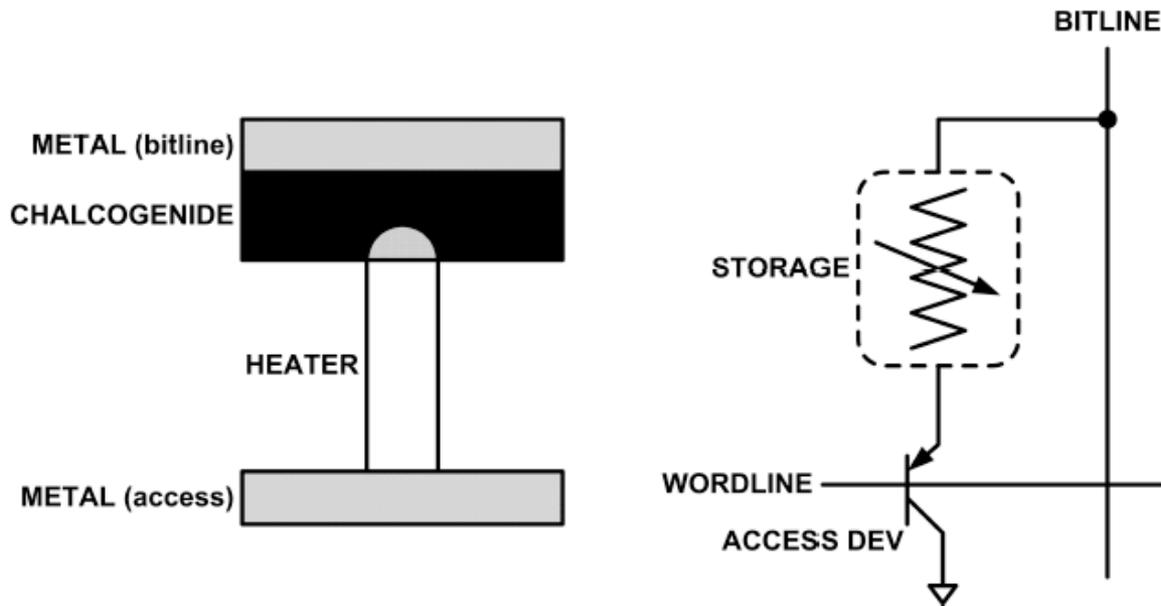
# Emerging Resistive Memory Technologies

---

- PCM
  - Inject current to change material phase
  - Resistance determined by phase
- STT-MRAM
  - Inject current to change magnet polarity
  - Resistance determined by polarity
- Memristors
  - Inject current to change atomic structure
  - Resistance determined by atom distance

# What is Phase Change Memory?

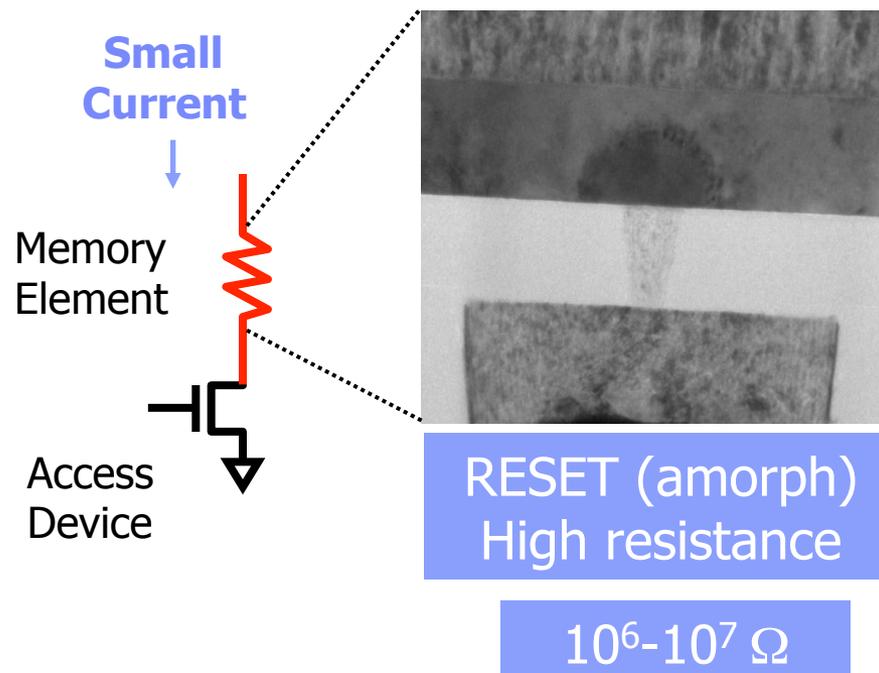
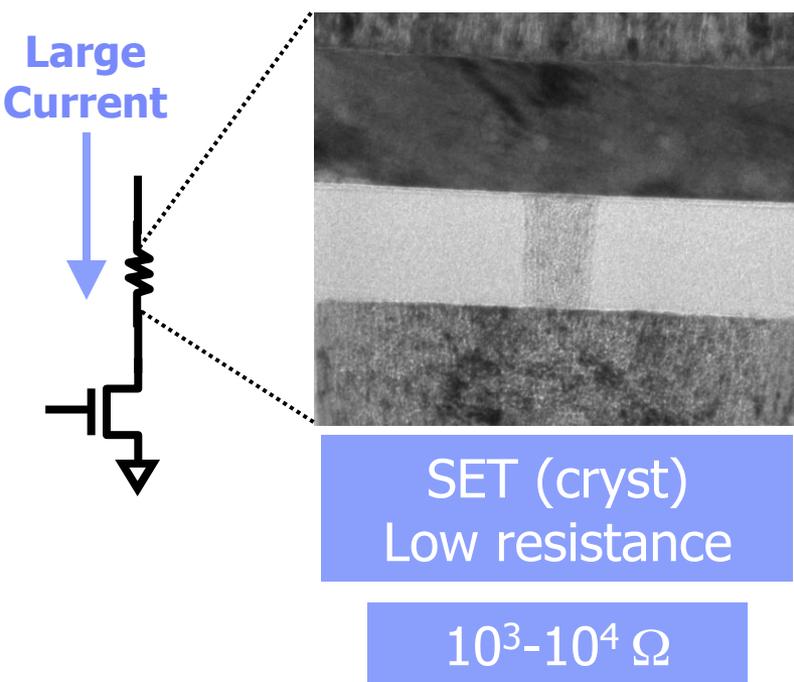
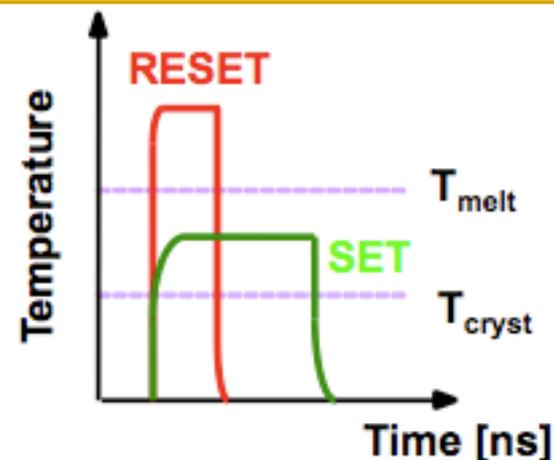
- Phase change material (chalcogenide glass) exists in two states:
  - ❑ Amorphous: Low optical reflexivity and high electrical resistivity
  - ❑ Crystalline: High optical reflexivity and low electrical resistivity



PCM is resistive memory: High resistance (0), Low resistance (1)  
PCM cell can be switched between states reliably and quickly

# How Does PCM Work?

- Write: change phase via current injection
  - SET: sustained current to heat cell above  $T_{cryst}$
  - RESET: cell heated above  $T_{melt}$  and quenched
- Read: detect phase via material resistance
  - amorphous/crystalline



# Opportunity: PCM Advantages

---

- Scales better than DRAM, Flash
  - Requires current pulses, which scale linearly with feature size
  - Expected to scale to 9nm (2022 [ITRS])
  - Prototyped at 20nm (Raoux+, IBM JRD 2008)
- Can be denser than DRAM
  - Can store multiple bits per cell due to large resistance range
  - Prototypes with 2 bits/cell in ISSCC' 08, 4 bits/cell by 2012
- Non-volatile
  - Retain data for >10 years at 85C
- No refresh needed, low idle power

# Phase Change Memory Properties

---

- Surveyed prototypes from 2003-2008 (ITRS, IEDM, VLSI, ISSCC)
- Derived PCM parameters for  $F=90\text{nm}$
  
- Lee, Ipek, Mutlu, Burger, “[Architecting Phase Change Memory as a Scalable DRAM Alternative](#),” ISCA 2009.

Table 1. Technology survey.

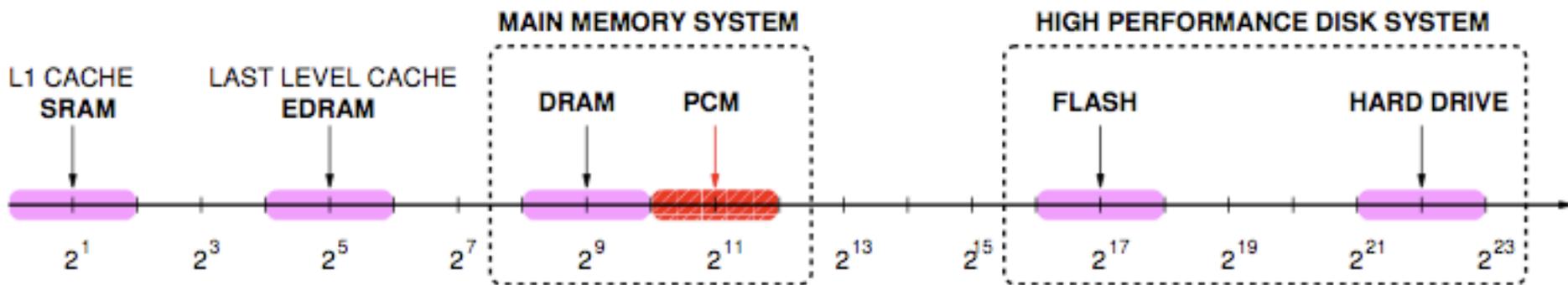
Parameter*	Published prototype									
	Horri <sup>6</sup>	Ahn <sup>12</sup>	Bedeschi <sup>13</sup>	Oh <sup>14</sup>	Pellizer <sup>15</sup>	Chen <sup>5</sup>	Kang <sup>16</sup>	Bedeschi <sup>9</sup>	Lee <sup>10</sup>	Lee <sup>2</sup>
Year	2003	2004	2004	2005	2006	2006	2006	2008	2008	**
Process, $F$ (nm)	**	120	180	120	90	**	100	90	90	90
Array size (Mbytes)	**	64	8	64	**	**	256	256	512	**
Material	GST, N-d	GST, N-d	GST	GST	GST	GS, N-d	GST	GST	GST	GST, N-d
Cell size ( $\mu\text{m}^2$ )	**	0.290	0.290	**	0.097	60 nm <sup>2</sup>	0.166	0.097	0.047	0.065 to 0.097
Cell size, $F^2$	**	20.1	9.0	**	12.0	**	16.6	12.0	5.8	9.0 to 12.0
Access device	**	**	BJT	FET	BJT	**	FET	BJT	Diode	BJT
Read time (ns)	**	70	48	68	**	**	62	**	55	48
Read current ( $\mu\text{A}$ )	**	**	40	**	**	**	**	**	**	40
Read voltage (V)	**	3.0	1.0	1.8	1.6	**	1.8	**	1.8	1.0
Read power ( $\mu\text{W}$ )	**	**	40	**	**	**	**	**	**	40
Read energy (pJ)	**	**	2.0	**	**	**	**	**	**	2.0
Set time (ns)	100	150	150	180	**	80	300	**	400	150
Set current ( $\mu\text{A}$ )	200	**	300	200	**	55	**	**	**	150
Set voltage (V)	**	**	2.0	**	**	1.25	**	**	**	1.2
Set power ( $\mu\text{W}$ )	**	**	300	**	**	34.4	**	**	**	90
Set energy (pJ)	**	**	45	**	**	2.8	**	**	**	13.5
Reset time (ns)	50	10	40	10	**	60	50	**	50	40
Reset current ( $\mu\text{A}$ )	600	600	600	600	400	90	600	300	600	300
Reset voltage (V)	**	**	2.7	**	1.8	1.6	**	1.6	**	1.6
Reset power ( $\mu\text{W}$ )	**	**	1620	**	**	80.4	**	**	**	480
Reset energy (pJ)	**	**	64.8	**	**	4.8	**	**	**	19.2
Write endurance (MLC)	$10^7$	$10^9$	$10^6$	**	$10^6$	$10^4$	**	$10^5$	$10^5$	$10^8$

\* BJT: bipolar junction transistor; FET: field-effect transistor; GST: Ge<sub>2</sub>Sb<sub>2</sub>Te<sub>5</sub>; MLC: multilevel cells; N-d: nitrogen doped.

\*\* This information is not available in the publication cited.

# Phase Change Memory Properties: Latency

- Latency comparable to, but slower than DRAM



Typical Access Latency (in terms of processor cycles for a 4 GHz processor)

- Read Latency
  - 50ns: 4x DRAM,  $10^{-3}$ x NAND Flash
- Write Latency
  - 150ns: 12x DRAM
- Write Bandwidth
  - 5-10 MB/s: 0.1x DRAM, 1x NAND Flash

# Phase Change Memory Properties

---

## ■ Dynamic Energy

- ❑ 40  $\mu\text{A}$  Rd, 150  $\mu\text{A}$  Wr
- ❑ 2-43x DRAM, 1x NAND Flash

## ■ Endurance

- ❑ Writes induce phase change at 650C
- ❑ Contacts degrade from thermal expansion/contraction
- ❑  $10^8$  writes per cell
- ❑  $10^{-8}\text{x}$  DRAM,  $10^3\text{x}$  NAND Flash

## ■ Cell Size

- ❑ 9-12F<sup>2</sup> using BJT, single-level cells
- ❑ 1.5x DRAM, 2-3x NAND (will scale with feature size, MLC)

# Phase Change Memory: Pros and Cons

---

- Pros over DRAM
  - Better technology scaling
  - Non volatility
  - Low idle power (no refresh)
- Cons
  - Higher latencies:  $\sim 4\text{-}15\times$  DRAM (especially write)
  - Higher active energy:  $\sim 2\text{-}50\times$  DRAM (especially write)
  - Lower endurance (a cell dies after  $\sim 10^8$  writes)
- Challenges in enabling PCM as DRAM replacement/helper:
  - Mitigate PCM shortcomings
  - Find the right way to place PCM in the system
  - Ensure secure and fault-tolerant PCM operation

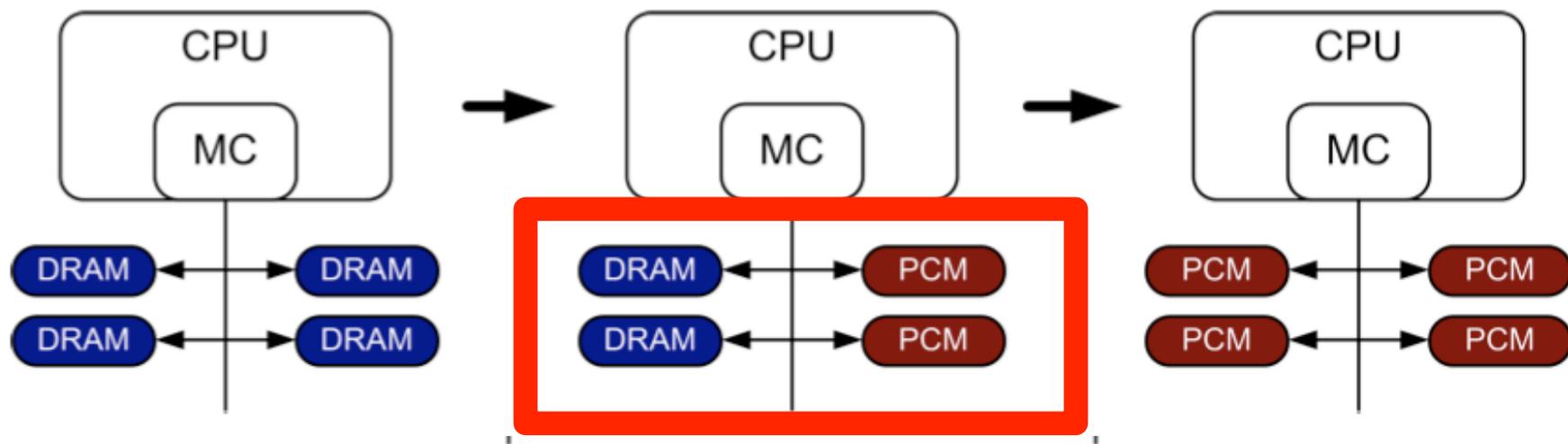
# PCM-based Main Memory: Research Challenges

---

- Where to place PCM in the memory hierarchy?
  - Hybrid OS controlled PCM-DRAM
  - Hybrid OS controlled PCM and hardware-controlled DRAM
  - Pure PCM main memory
- How to mitigate shortcomings of PCM?
- How to minimize amount of DRAM in the system?
- How to take advantage of (byte-addressable and fast) non-volatile main memory?
- Can we design specific-NVM-technology-agnostic techniques?

# PCM-based Main Memory (I)

- How should PCM-based (main) memory be organized?



- **Hybrid PCM+DRAM** [Qureshi+ ISCA'09, Dhiman+ DAC'09, Meza+ IEEE CAL'12]:
  - How to partition/migrate data between PCM and DRAM

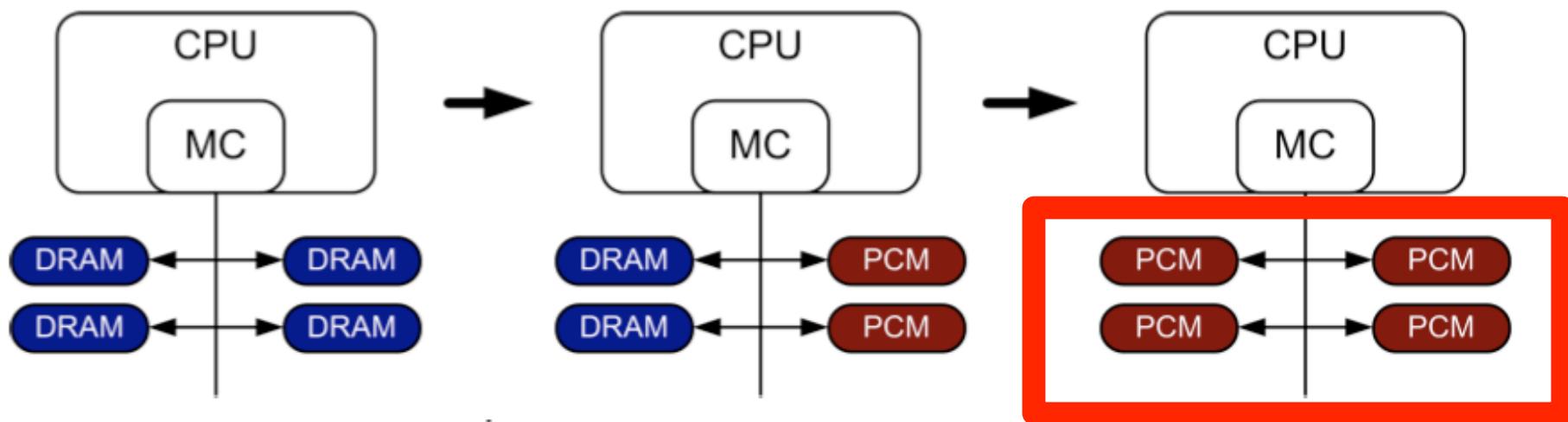
# Hybrid Memory Systems: Challenges

---

- **Partitioning**
  - Should DRAM be a cache or main memory, or configurable?
  - What fraction? How many controllers?
- **Data allocation/movement (energy, performance, lifetime)**
  - Who manages allocation/movement?
  - What are good control algorithms?
  - How do we prevent degradation of service due to wearout?
- **Design of cache hierarchy, memory controllers, OS**
  - Mitigate PCM shortcomings, exploit PCM advantages
- **Design of PCM/DRAM chips and modules**
  - Rethink the design of PCM/DRAM with new requirements

# PCM-based Main Memory (II)

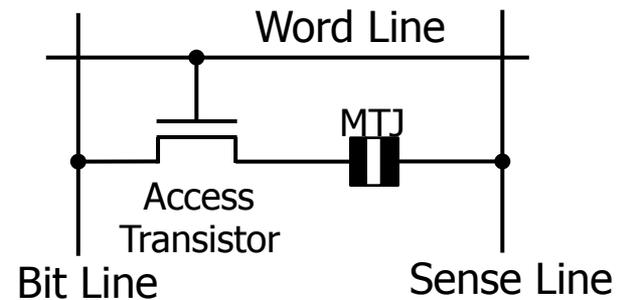
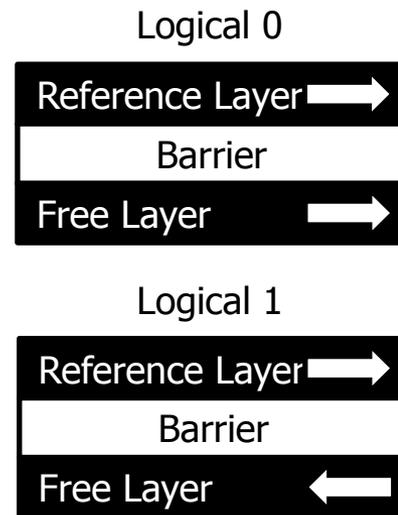
- How should PCM-based (main) memory be organized?



- **Pure PCM main memory** [Lee et al., ISCA'09, Top Picks'10]:
  - How to redesign entire hierarchy (and cores) to overcome PCM shortcomings

# Aside: STT-RAM Basics

- Magnetic Tunnel Junction (MTJ)
  - ❑ Reference layer: Fixed
  - ❑ Free layer: Parallel or anti-parallel
- Cell
  - ❑ Access transistor, bit/sense lines
- Read and Write
  - ❑ Read: Apply a small voltage across bitline and senseline; read the current.
  - ❑ Write: Push large current through MTJ. Direction of current determines new orientation of the free layer.
- Kultursay et al., "Evaluating STT-RAM as an Energy-Efficient Main Memory Alternative," ISPASS 2013



# Aside: STT MRAM: Pros and Cons

---

- Pros over DRAM
  - Better technology scaling
  - Non volatility
  - Low idle power (no refresh)
- Cons
  - Higher write latency
  - Higher write energy
  - Reliability?
- Another level of freedom
  - Can trade off non-volatility for lower write latency/energy (by reducing the size of the MTJ)

# Agenda

---

- Major Trends Affecting Main Memory
- Requirements from an Ideal Main Memory System
- Opportunity: Emerging Memory Technologies
  - Background
  - PCM (or Technology X) as DRAM Replacement
  - Hybrid Memory Systems
- Conclusions
- Discussion

# An Initial Study: Replace DRAM with PCM

- Lee, Ipek, Mutlu, Burger, “Architecting Phase Change Memory as a Scalable DRAM Alternative,” ISCA 2009.
  - Surveyed prototypes from 2003-2008 (e.g. IEDM, VLSI, ISSCC)
  - Derived “average” PCM parameters for F=90nm

## Density

- ▷ 9 - 12F<sup>2</sup> using BJT
- ▷ 1.5× DRAM

## Latency

- ▷ 50ns Rd, 150ns Wr
- ▷ 4×, 12× DRAM

## Endurance

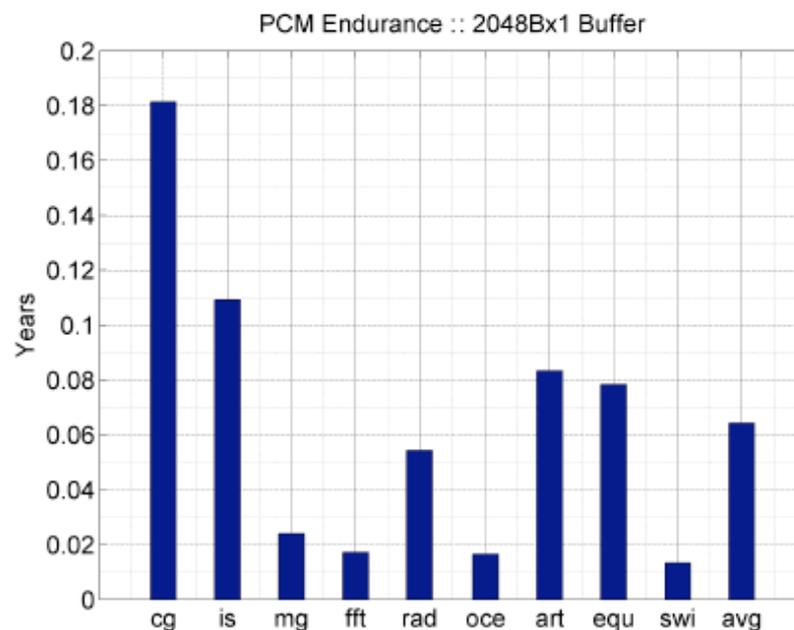
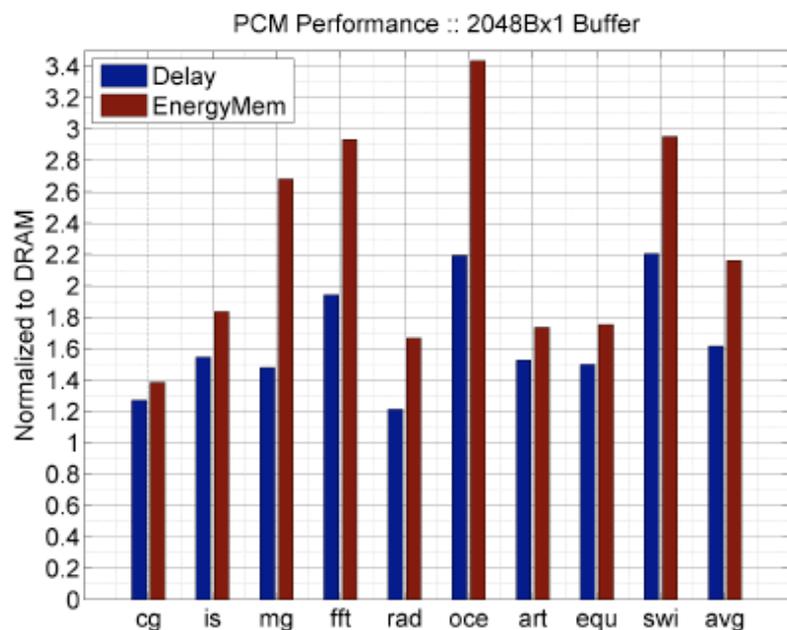
- ▷ 1E+08 writes
- ▷ 1E-08× DRAM

## Energy

- ▷ 40μA Rd, 150μA Wr
- ▷ 2×, 43× DRAM

# Results: Naïve Replacement of DRAM with PCM

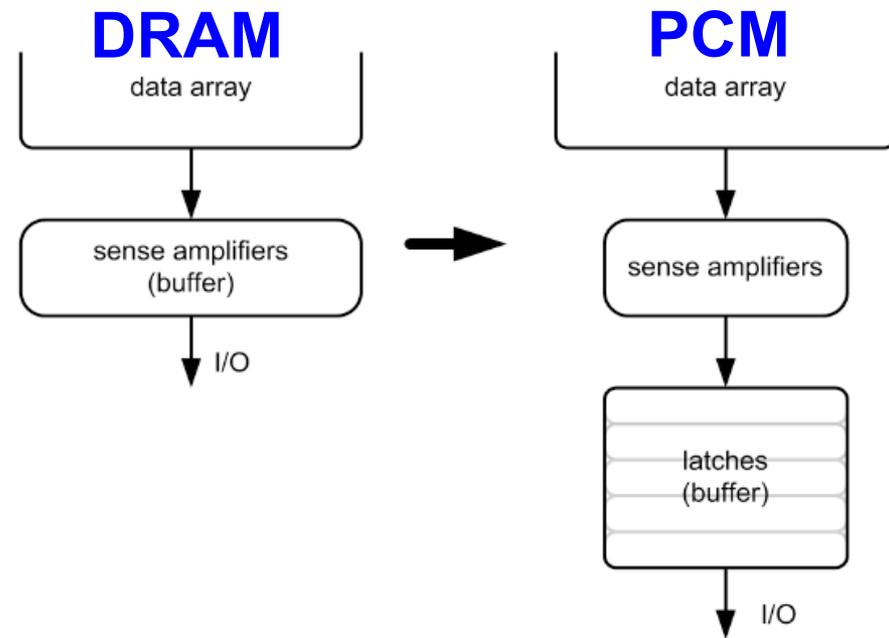
- Replace DRAM with PCM in a 4-core, 4MB L2 system
- PCM organized the same as DRAM: row buffers, banks, peripherals
- 1.6x delay, 2.2x energy, 500-hour average lifetime



- Lee, Ipek, Mutlu, Burger, “Architecting Phase Change Memory as a Scalable DRAM Alternative,” ISCA 2009.

# Architecting PCM to Mitigate Shortcomings

- Idea 1: Use multiple narrow row buffers in each PCM chip  
→ Reduces array reads/writes → better endurance, latency, energy
- Idea 2: Write into array at cache block or word granularity  
→ Reduces unnecessary wear



# Results: Architected PCM as Main Memory

---

- 1.2x delay, 1.0x energy, 5.6-year average lifetime
- Scaling improves energy, endurance, density



- Caveat 1: Worst-case lifetime is much shorter (no guarantees)
- Caveat 2: Intensive applications see large performance and energy hits
- Caveat 3: Optimistic PCM parameters?

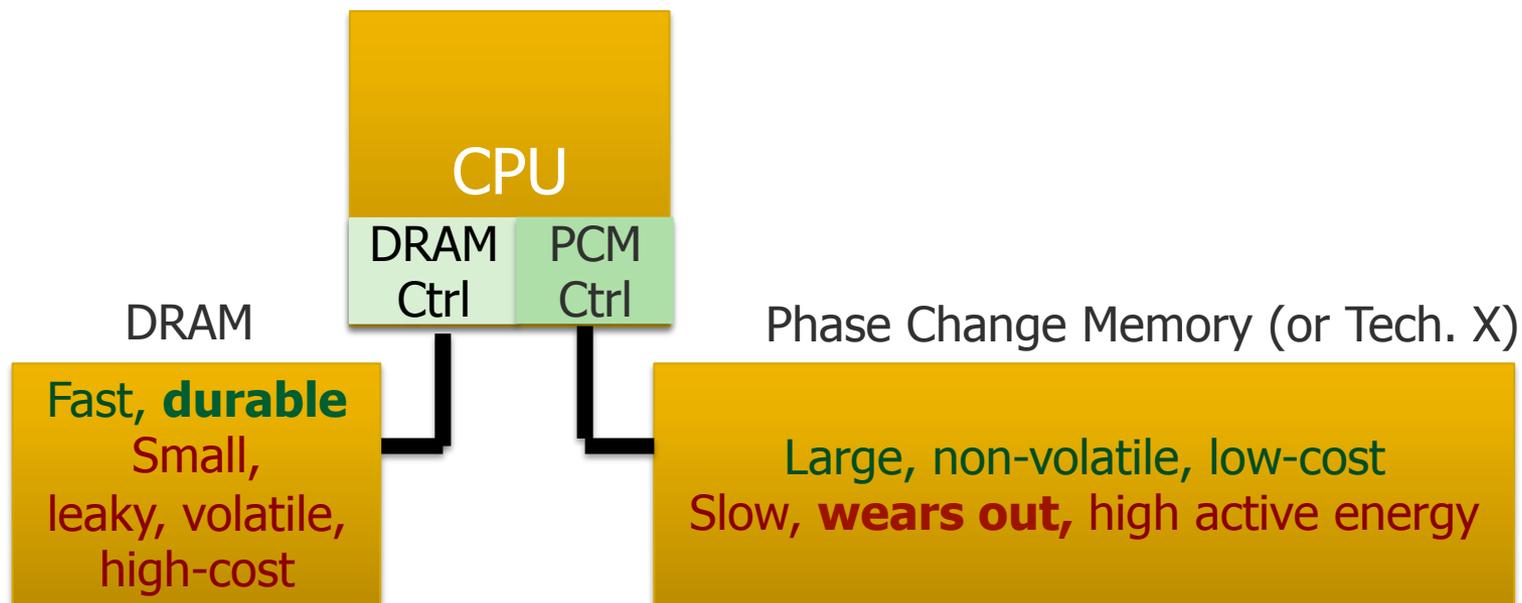
# Agenda

---

- Major Trends Affecting Main Memory
- Requirements from an Ideal Main Memory System
- Opportunity: Emerging Memory Technologies
  - Background
  - PCM (or Technology X) as DRAM Replacement
  - Hybrid Memory Systems
- Conclusions
- Discussion

# Hybrid Memory Systems

---



Hardware/software manage data allocation and movement  
to achieve the best of multiple technologies

Meza, Chang, Yoon, Mutlu, Ranganathan, "Enabling Efficient and Scalable Hybrid Memories,"  
IEEE Comp. Arch. Letters, 2012.

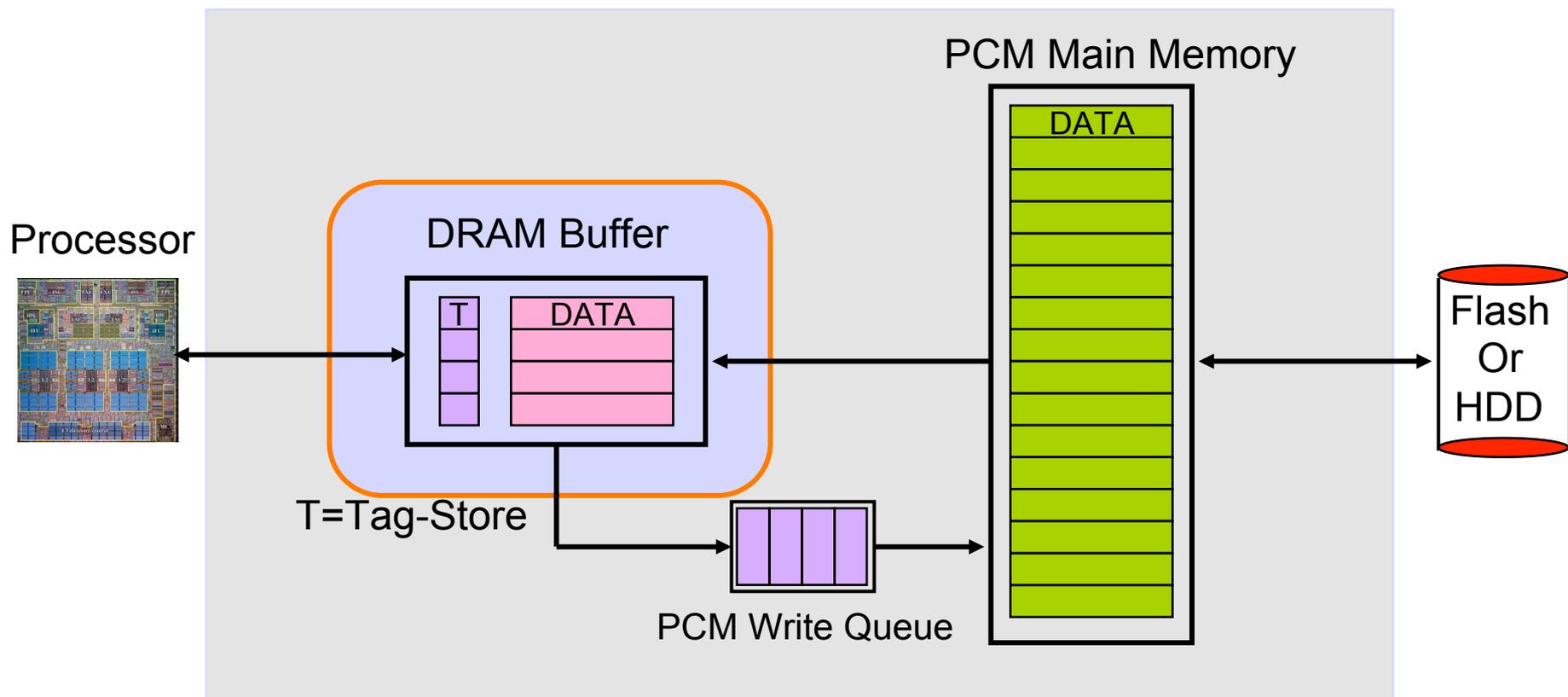
# One Option: DRAM as a Cache for PCM

---

- PCM is main memory; DRAM caches memory rows/blocks
  - Benefits: Reduced latency on DRAM cache hit; write filtering
- Memory controller hardware manages the DRAM cache
  - Benefit: Eliminates system software overhead
- Three issues:
  - What data should be placed in DRAM versus kept in PCM?
  - What is the granularity of data movement?
  - How to design a low-cost hardware-managed DRAM cache?
- Two idea directions:
  - Locality-aware data placement [Yoon+ , ICCD 2012]
  - Cheap tag stores and dynamic granularity [Meza+, IEEE CAL 2012]

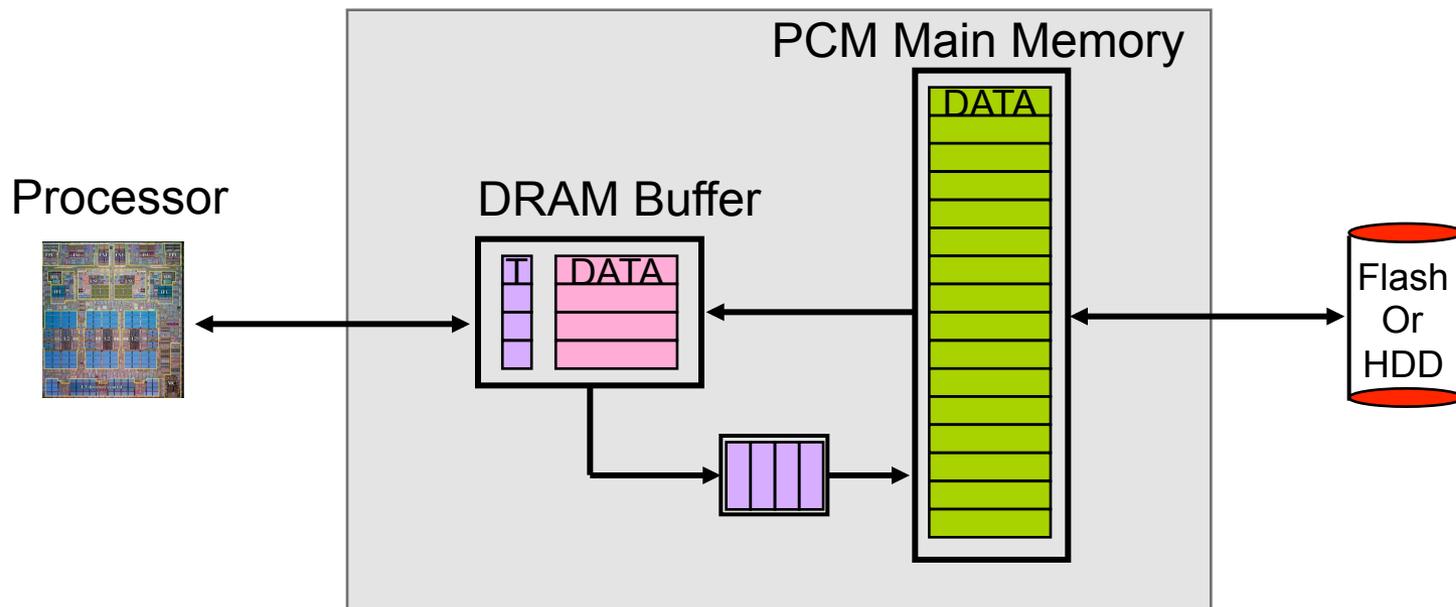
# DRAM as a Cache for PCM

- Goal: Achieve the best of both DRAM and PCM/NVM
  - Minimize amount of DRAM w/o sacrificing performance, endurance
  - DRAM as cache to tolerate PCM latency and write bandwidth
  - PCM as main memory to provide large capacity at good cost and power



# Write Filtering Techniques

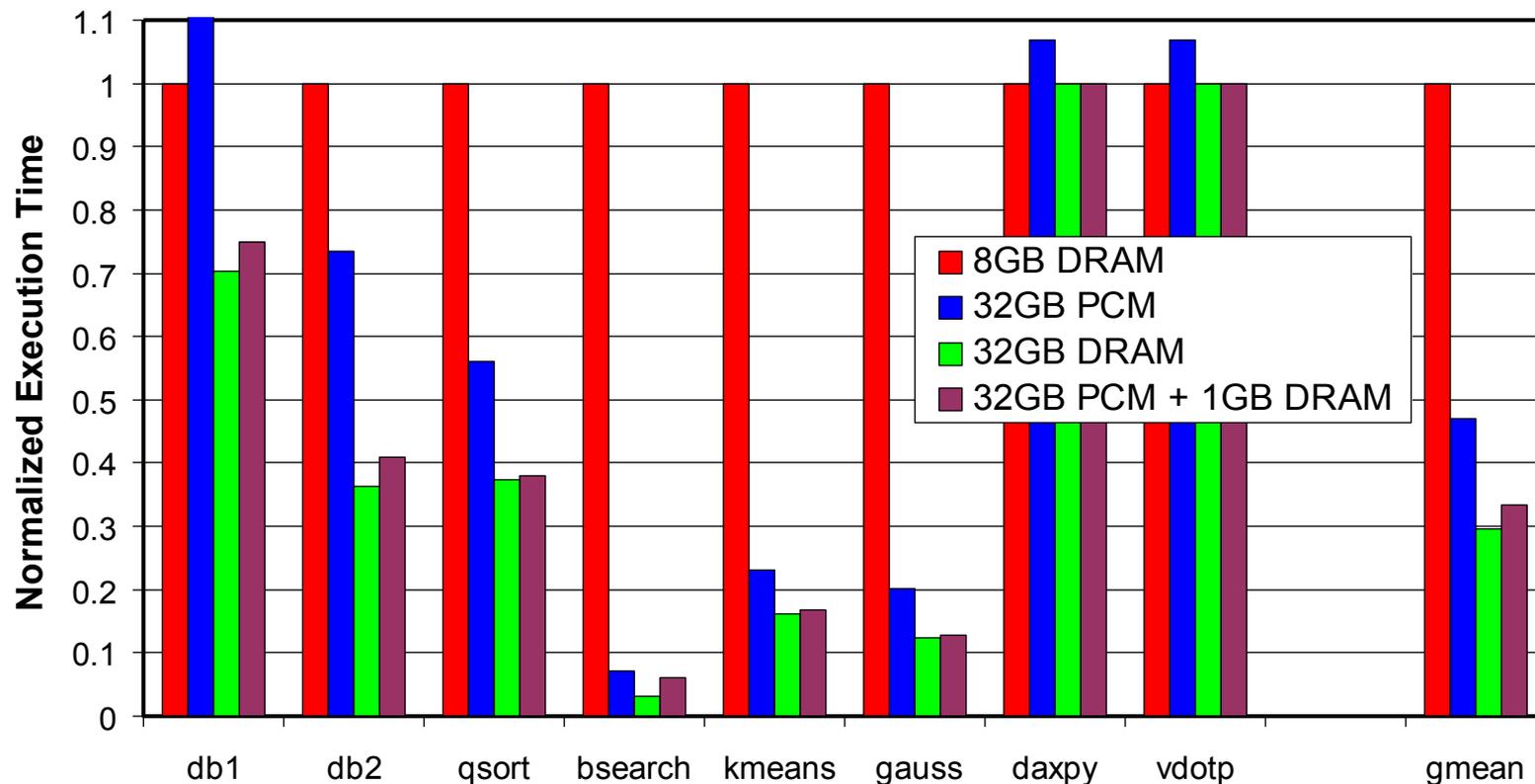
- Lazy Write: Pages from disk installed only in DRAM, not PCM
- Partial Writes: Only dirty lines from DRAM page written back
- Page Bypass: Discard pages with poor reuse on DRAM eviction



- Qureshi et al., “[Scalable high performance main memory system using phase-change memory technology](#),” ISCA 2009.

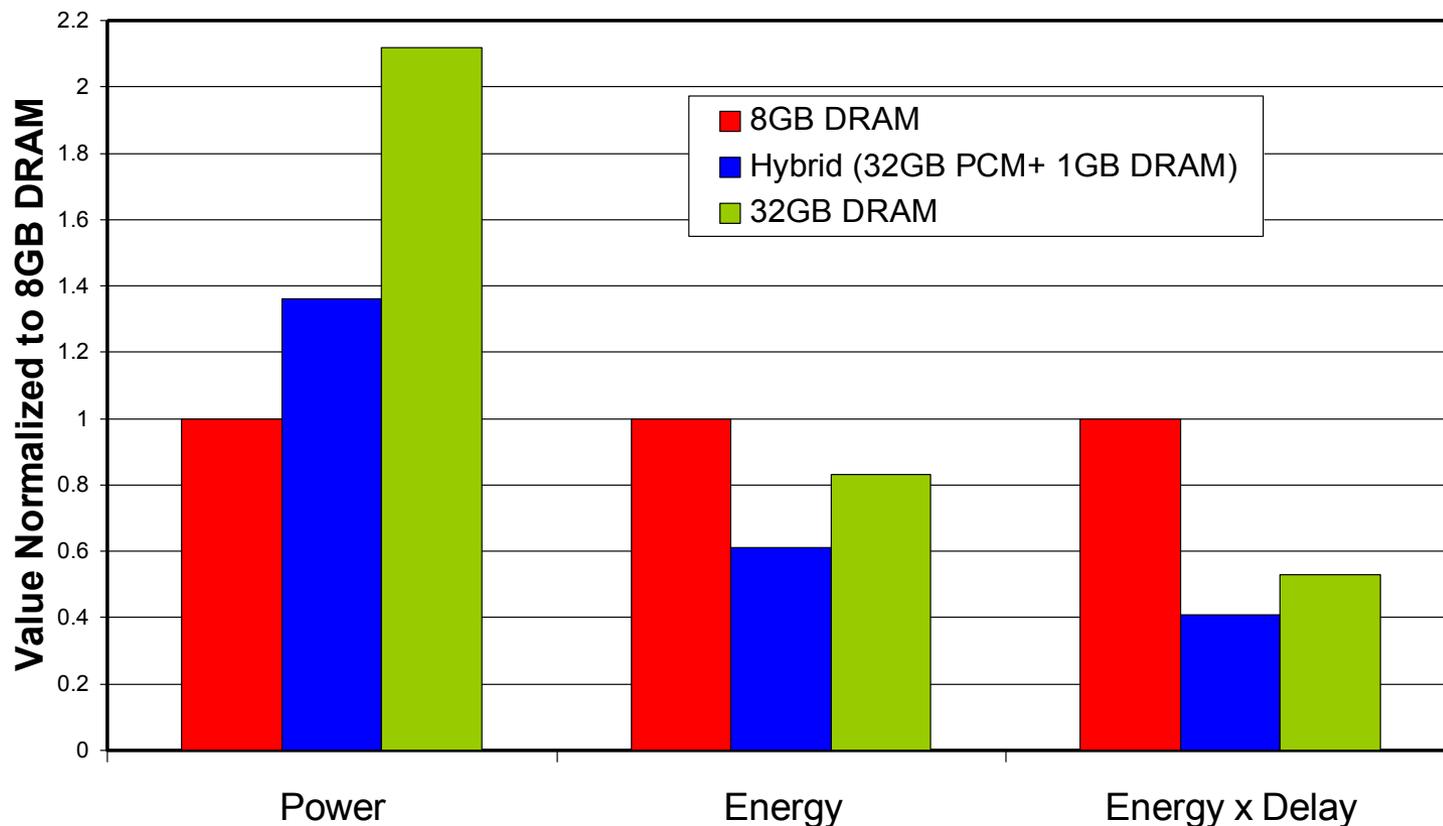
# Results: DRAM as PCM Cache (I)

- Simulation of 16-core system, 8GB DRAM main-memory at 320 cycles, HDD (2 ms) with Flash (32 us) with Flash hit-rate of 99%
- Assumption: PCM 4x denser, 4x slower than DRAM
- DRAM block size = PCM page size (4kB)



# Results: DRAM as PCM Cache (II)

- PCM-DRAM Hybrid performs similarly to similar-size DRAM
- Significant power and energy savings with PCM-DRAM Hybrid
- Average lifetime: 9.7 years (no guarantees)



# Agenda

---

- Major Trends Affecting Main Memory
- Requirements from an Ideal Main Memory System
- Opportunity: Emerging Memory Technologies
  - Background
  - PCM (or Technology X) as DRAM Replacement
  - Hybrid Memory Systems
    - Row-Locality Aware Data Placement
    - Efficient DRAM (or Technology X) Caches
- Conclusions
- Discussion

# Row Buffer Locality Aware Caching Policies for Hybrid Memories

HanBin Yoon

Justin Meza

Rachata Ausavarungnirun

Rachael Harding

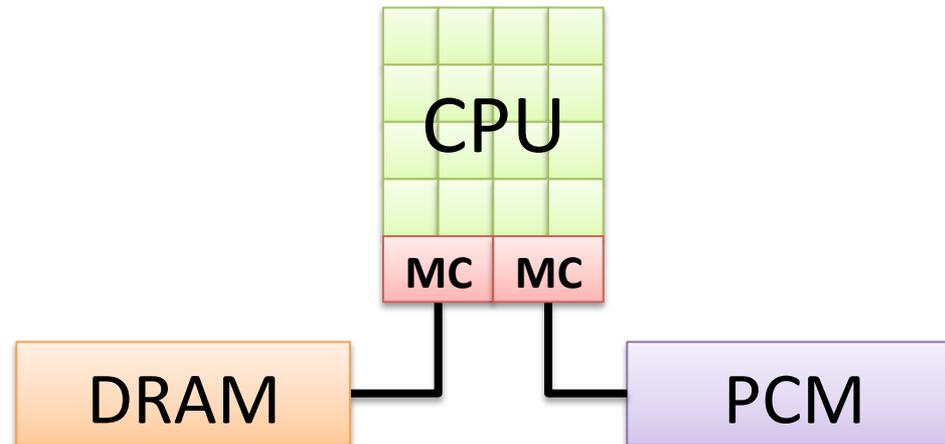
Onur Mutlu

**Carnegie Mellon University**

# Hybrid Memory

---

- Key question: How to place data between the heterogeneous memory devices?



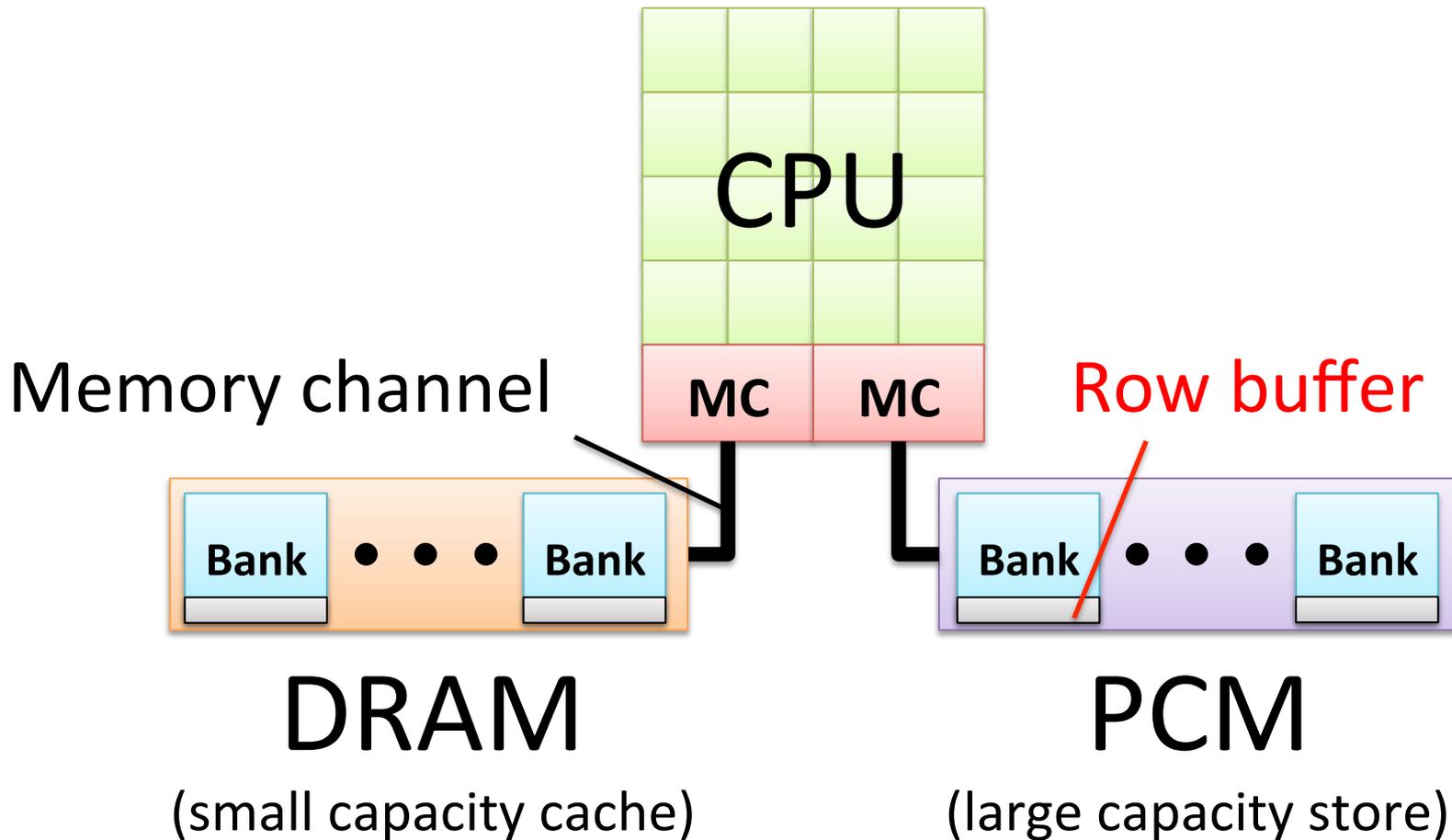
# Outline

---

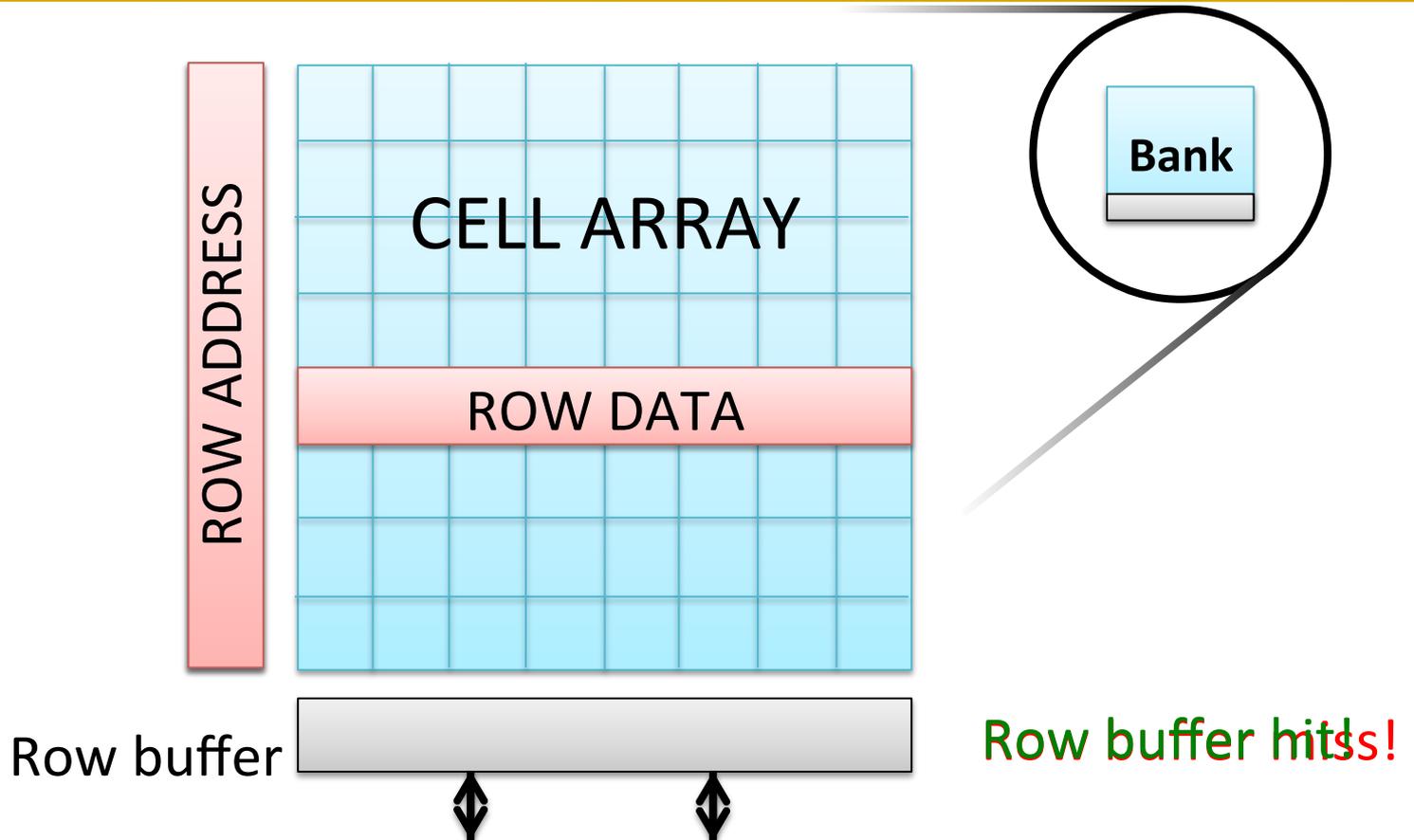
- Background: Hybrid Memory Systems
- **Motivation: Row Buffers and Implications on Data Placement**
- Mechanisms: Row Buffer Locality-Aware Caching Policies
- Evaluation and Results
- Conclusion

# Hybrid Memory: A Closer Look

---



# Row Buffers and Latency



Row (buffer) hit: Access data from row buffer → fast

Row (buffer) miss: Access data from cell array → slow

# Key Observation

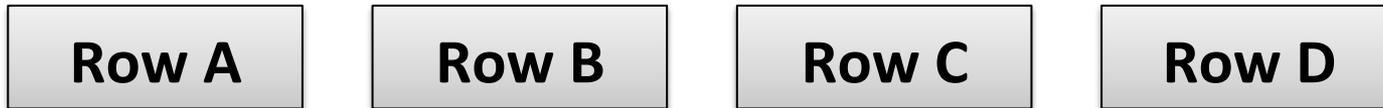
---

- Row buffers exist in both DRAM and PCM
  - Row **hit** latency **similar** in DRAM & PCM [Lee+ ISCA'09]
  - Row **miss** latency **small** in DRAM, **large** in PCM
- Place data in DRAM which
  - is likely to miss in the row buffer (**low row buffer locality**) → miss penalty is smaller in DRAM
  - AND
  - is **reused many times** → cache only the data worth the movement cost and DRAM space

# RBL-Awareness: An Example

---

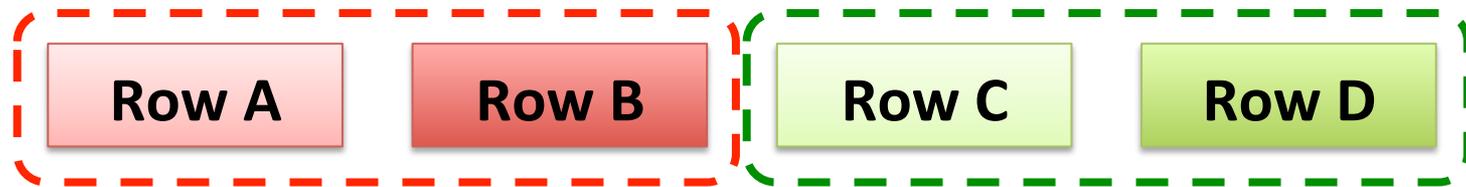
Let's say a processor accesses four rows



# RBL-Awareness: An Example

---

Let's say a processor accesses four rows with different row buffer localities (RBL)



**Low RBL**

(Frequently miss  
in row buffer)

**High RBL**

(Frequently hit  
in row buffer)

Case 1: RBL-*Unaware* Policy (state-of-the-art)

Case 2: RBL-Aware Policy (RBLA)

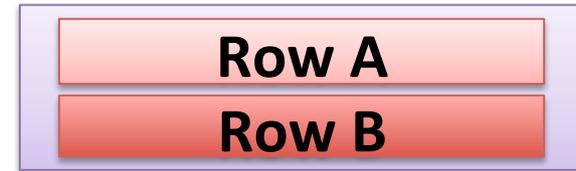
# Case 1: RBL-*Unaware* Policy

---

A **row buffer locality-*unaware*** policy could place these rows in the following manner



**DRAM**  
(High RBL)

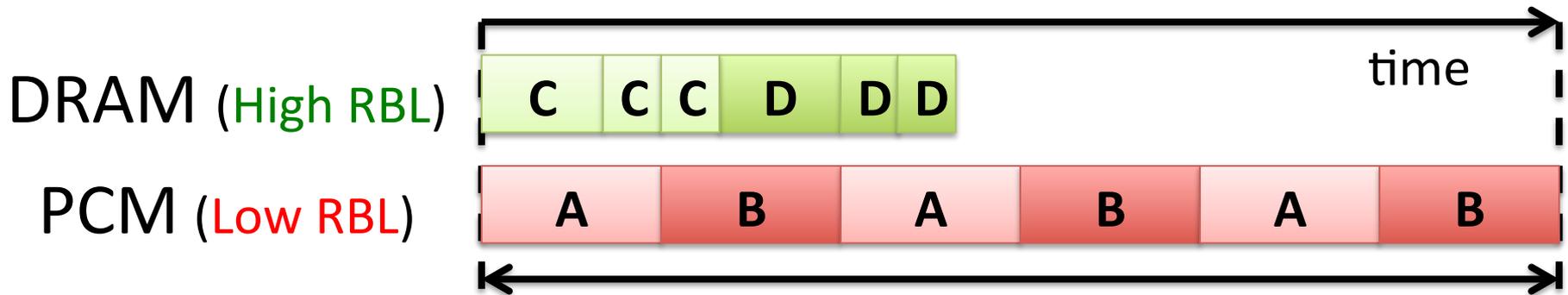


**PCM**  
(Low RBL)

# Case 1: RBL-*Unaware* Policy

Access pattern to main memory:

A (oldest), B, C, C, C, A, B, D, D, D, A, B (youngest)

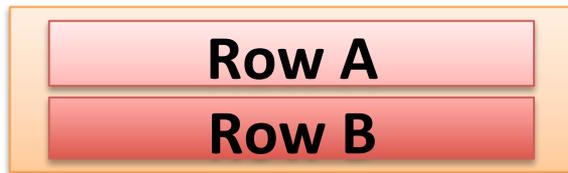


RBL-*Unaware*: Stall time is 6 PCM device accesses

# Case 2: RBL-Aware Policy (RBLA)

---

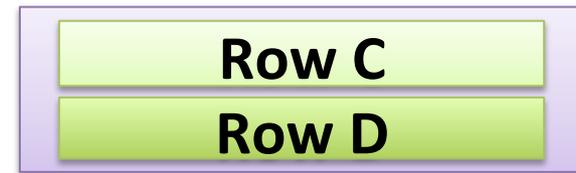
A **row buffer locality-aware** policy would place these rows in the **opposite** manner



**DRAM**

(Low RBL)

→ Access data at lower row buffer **miss** latency of DRAM



**PCM**

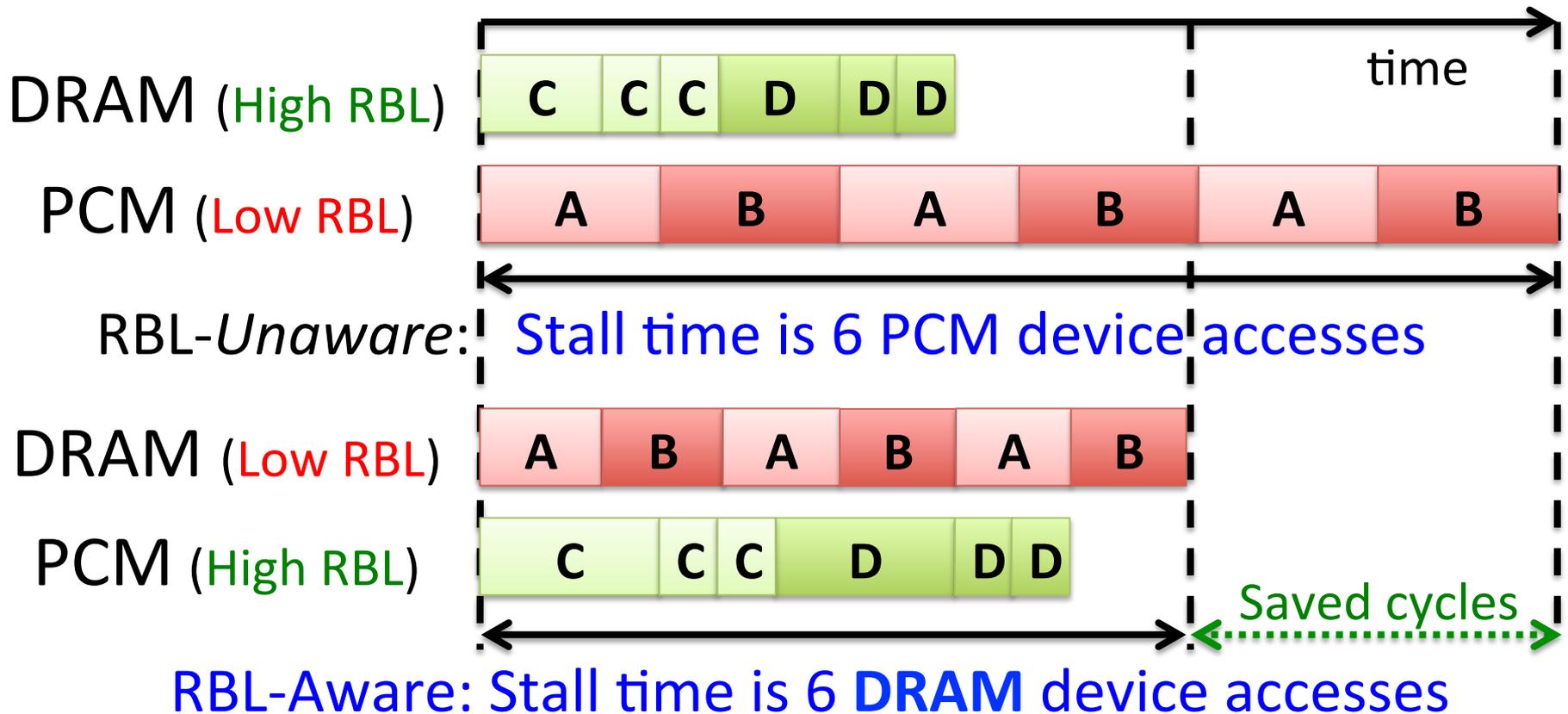
(High RBL)

→ Access data at low row buffer **hit** latency of PCM

# Case 2: RBL-Aware Policy (RBLA)

Access pattern to main memory:

A (oldest), B, C, C, C, A, B, D, D, D, A, B (youngest)



# Outline

---

- Background: Hybrid Memory Systems
- Motivation: Row Buffers and Implications on Data Placement
- **Mechanisms: Row Buffer Locality-Aware Caching Policies**
- Evaluation and Results
- Conclusion

# Our Mechanism: RBLA

---

1. For recently used rows in PCM:
  - Count row buffer **misses** as indicator of row buffer locality (RBL)
2. Cache to DRAM rows with **misses**  $\geq$  **threshold**
  - Row buffer miss counts are periodically reset (only cache rows with high reuse)

# Our Mechanism: RBLA-Dyn

---

1. For recently used rows in PCM:
  - Count row buffer **misses** as indicator of row buffer locality (RBL)
2. Cache to DRAM rows with **misses**  $\geq$  **threshold**
  - Row buffer miss counts are periodically reset (only cache rows with high reuse)
3. Dynamically adjust **threshold** to adapt to workload/system characteristics
  - Interval-based cost-benefit analysis

# Implementation: “Statistics Store”

---

- Goal: To keep count of row buffer misses to recently used rows in PCM
- Hardware structure in memory controller
  - Operation is similar to a cache
    - Input: row address
    - Output: row buffer miss count
  - 128-set 16-way statistics store (9.25KB) achieves system performance within 0.3% of an unlimited-sized statistics store

# Outline

---

- Background: Hybrid Memory Systems
- Motivation: Row Buffers and Implications on Data Placement
- Mechanisms: Row Buffer Locality-Aware Caching Policies
- **Evaluation and Results**
- **Conclusion**

# Evaluation Methodology

---

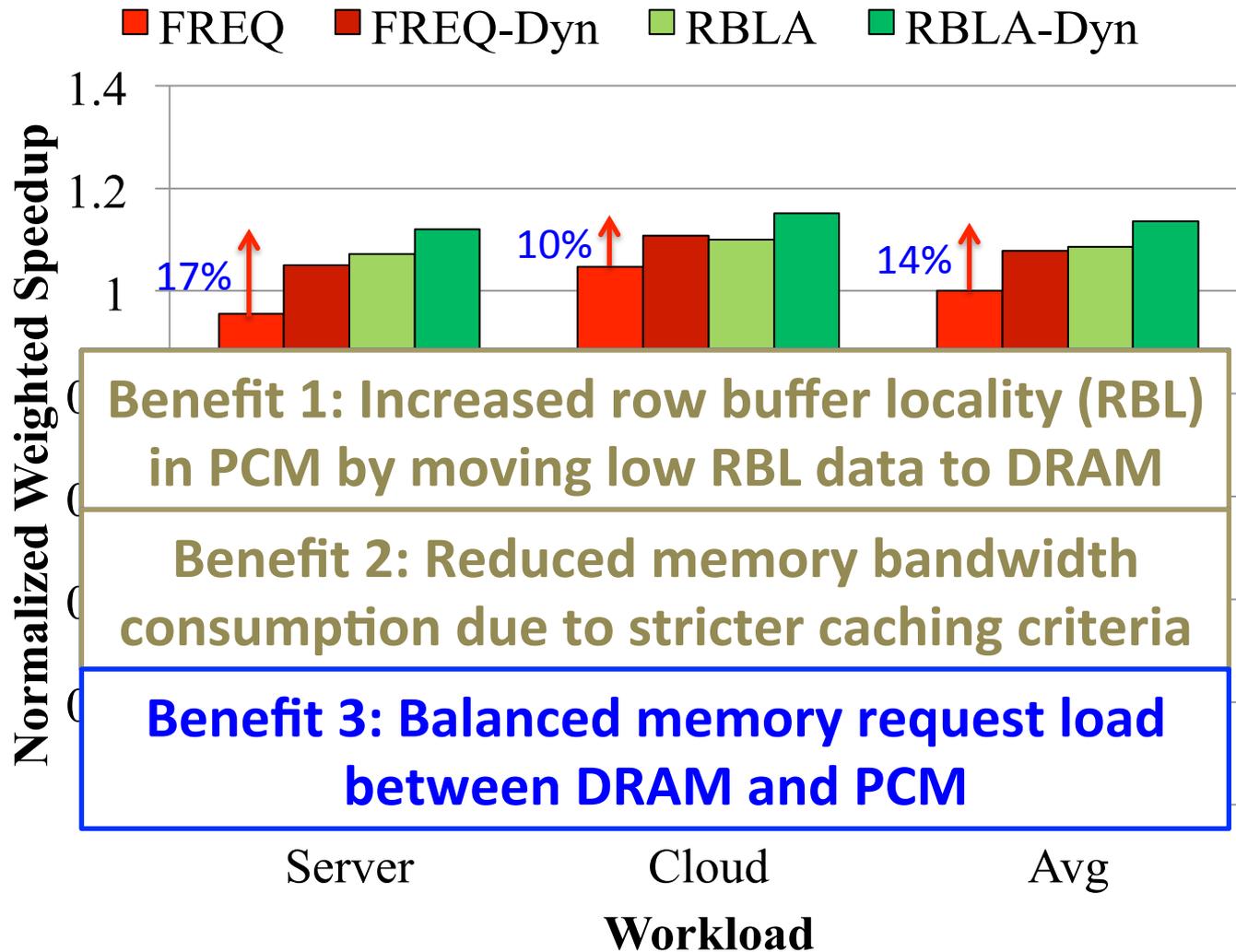
- Cycle-level x86 CPU-memory simulator
  - **CPU**: 16 out-of-order cores, 32KB private L1 per core, 512KB shared L2 per core
  - **Memory**: 1GB DRAM (8 banks), 16GB PCM (8 banks), 4KB migration granularity
- 36 multi-programmed server, cloud workloads
  - Server: TPC-C (OLTP), TPC-H (Decision Support)
  - Cloud: Apache (Webserv.), H.264 (Video), TPC-C/H
- Metrics: Weighted speedup (perf.), perf./Watt (energy eff.), Maximum slowdown (fairness)

# Comparison Points

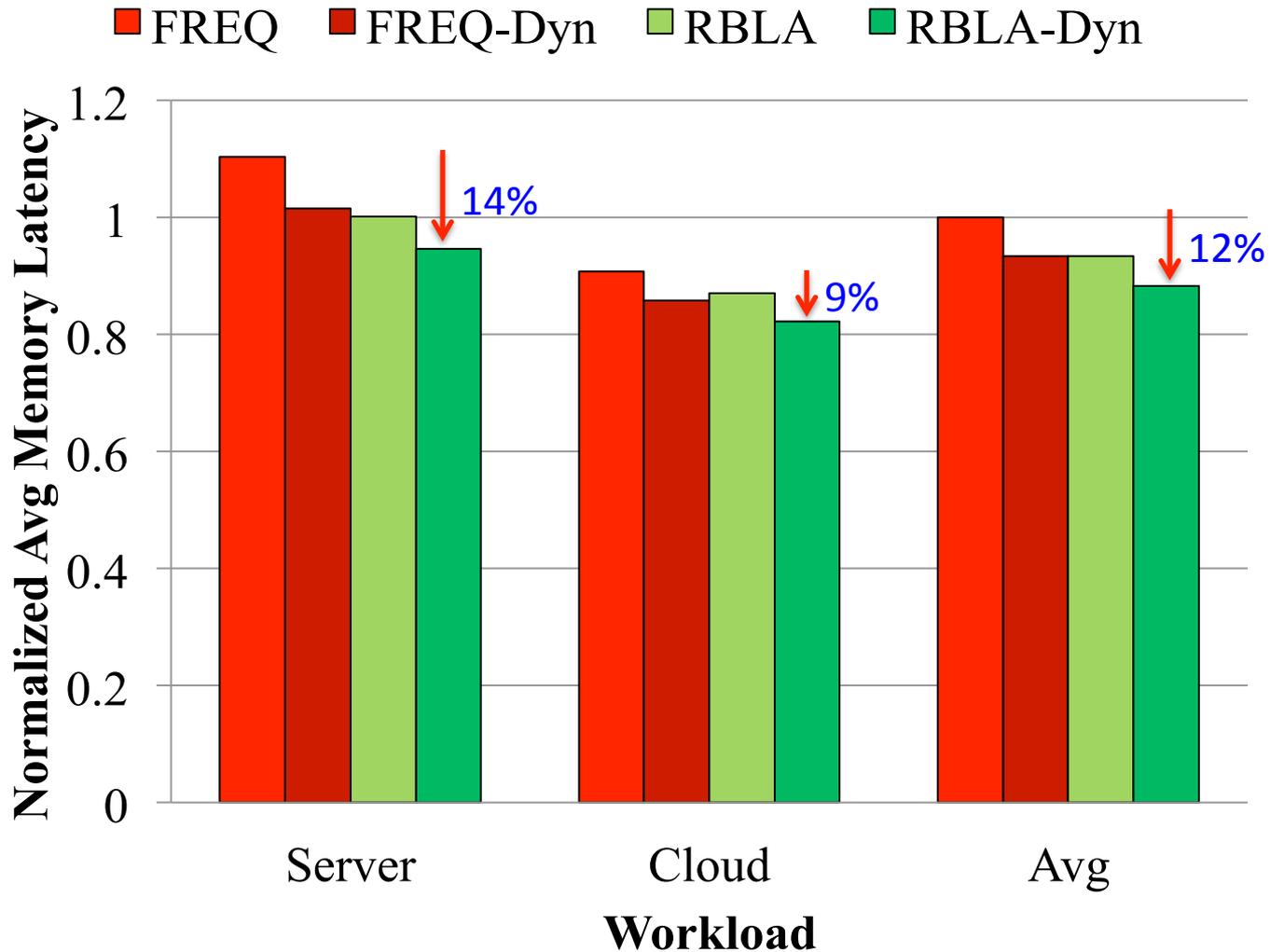
---

- **Conventional LRU Caching**
- **FREQ:** Access-frequency-based caching
  - Places “hot data” in cache [Jiang+ HPCA'10]
  - Cache to DRAM rows with accesses  $\geq$  threshold
  - *Row buffer locality-unaware*
- **FREQ-Dyn:** Adaptive Freq.-based caching
  - **FREQ** + our dynamic threshold adjustment
  - *Row buffer locality-unaware*
- **RBLA:** Row buffer locality-aware caching
- **RBLA-Dyn:** Adaptive RBL-aware caching

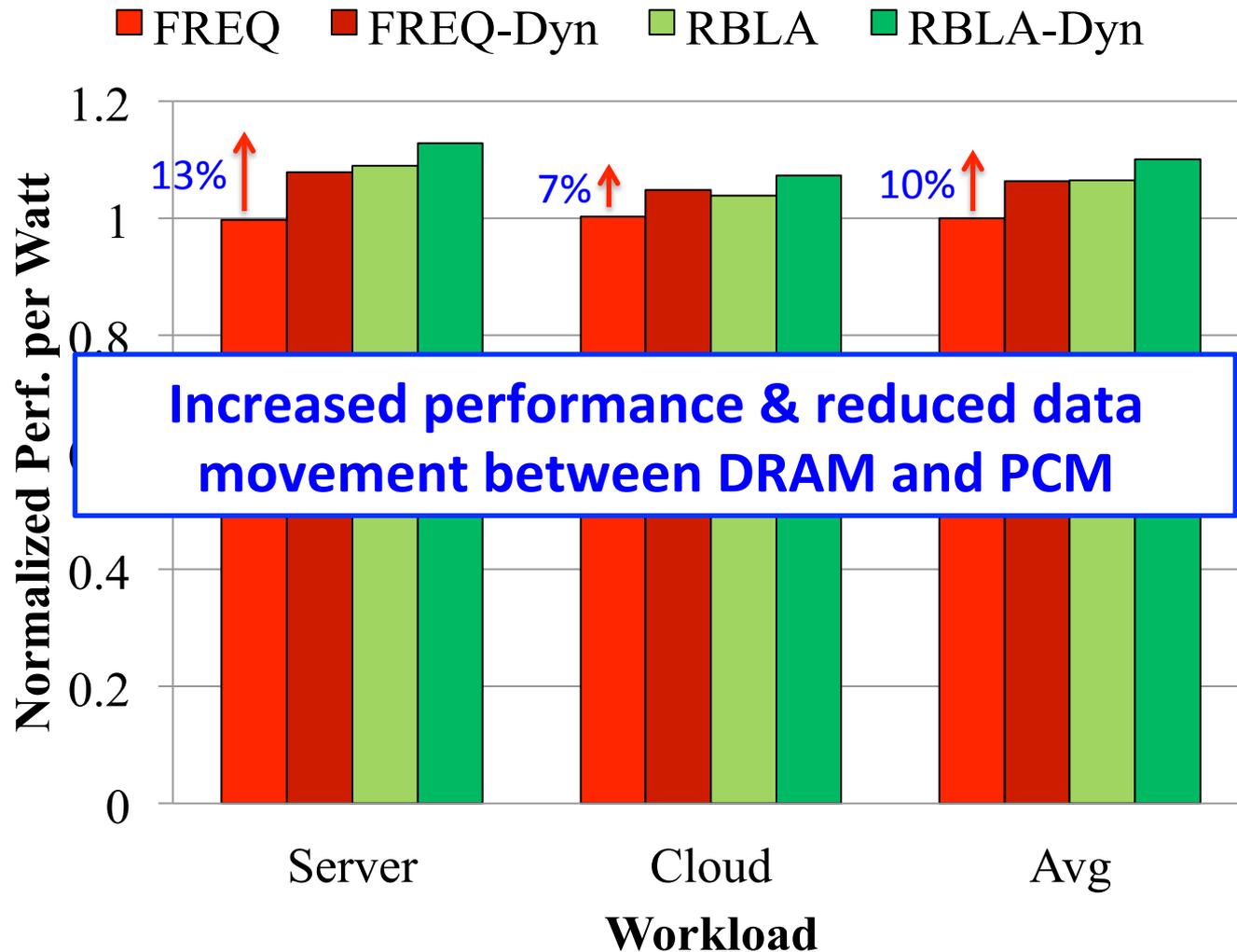
# System Performance



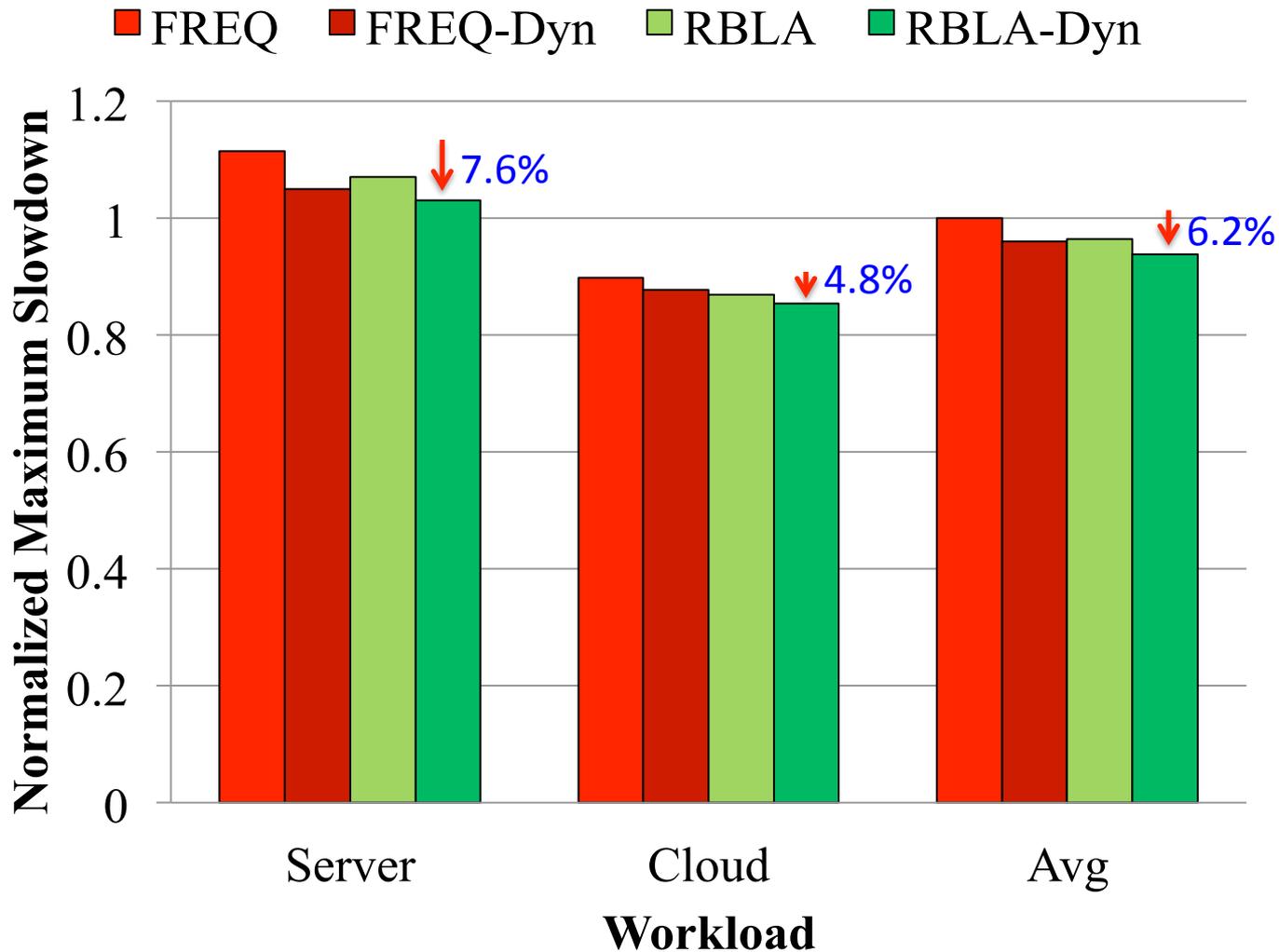
# Average Memory Latency



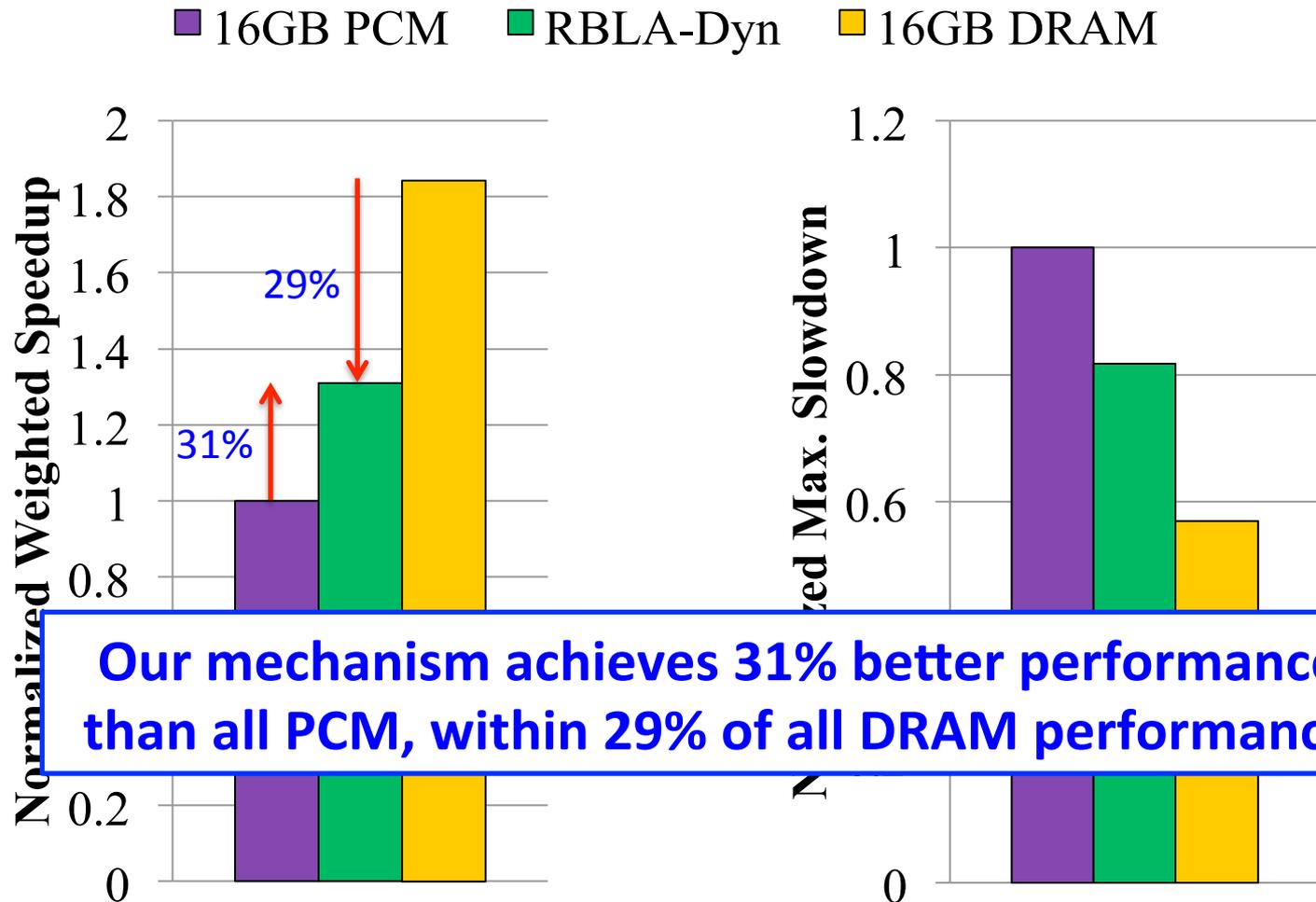
# Memory Energy Efficiency



# Thread Fairness



# Compared to All-PCM/DRAM



# Other Results in Paper

---

- RBLA-Dyn increases the portion of PCM row buffer hit by 6.6 times
- RBLA-Dyn has the effect of balancing memory request load between DRAM and PCM
  - PCM channel utilization increases by 60%.

# Summary

---

- Different memory technologies have different strengths
- A hybrid memory system (DRAM-PCM) aims for best of both
- **Problem:** How to place data between these heterogeneous memory devices?
- **Observation:** PCM array access latency is higher than DRAM's – But peripheral circuit (row buffer) access latencies are similar
- **Key Idea:** Use row buffer locality (RBL) as a key criterion for data placement
- **Solution:** Cache to DRAM rows with low RBL and high reuse
- Improves both performance and energy efficiency over state-of-the-art caching policies

# Row Buffer Locality Aware Caching Policies for Hybrid Memories

HanBin Yoon

Justin Meza

Rachata Ausavarungnirun

Rachael Harding

Onur Mutlu

**Carnegie Mellon University**

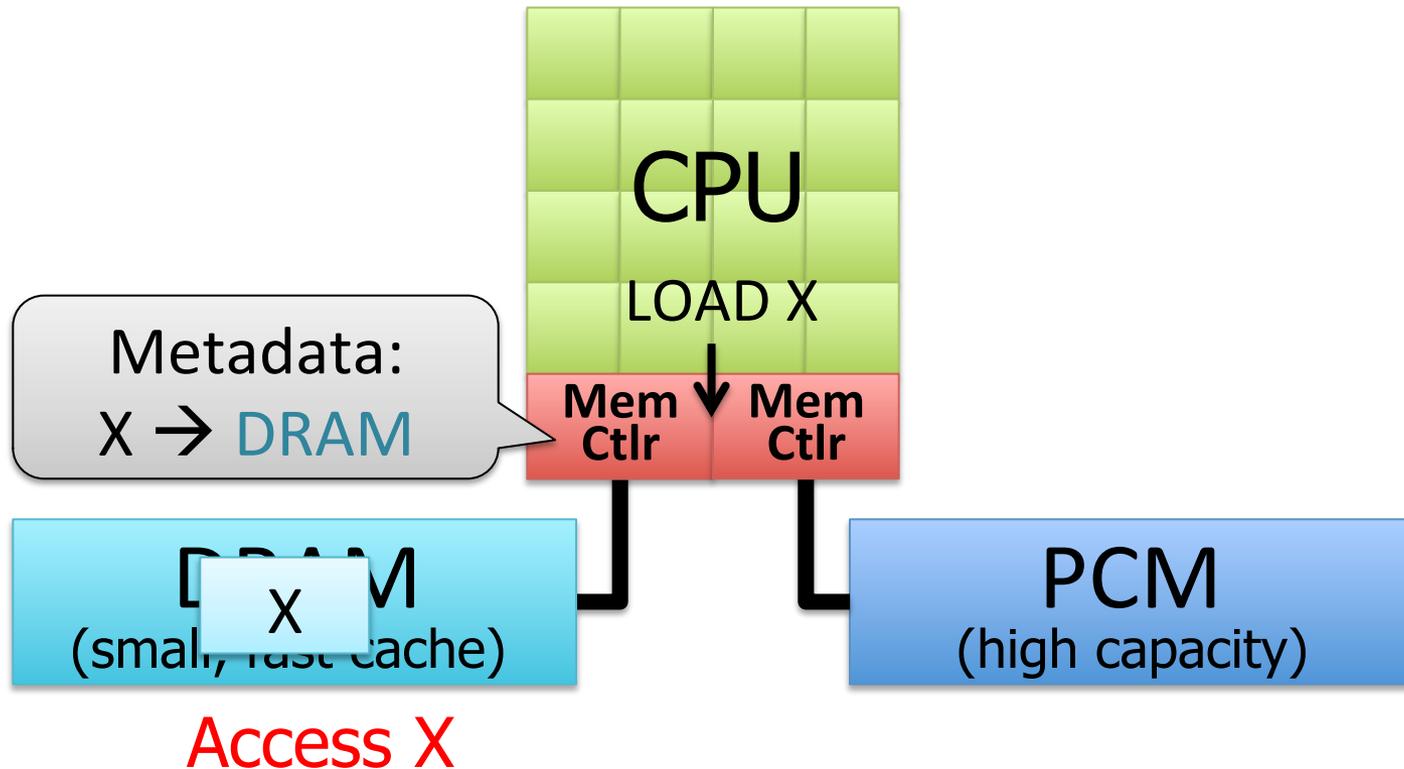
# Agenda

---

- Major Trends Affecting Main Memory
- Requirements from an Ideal Main Memory System
- Opportunity: Emerging Memory Technologies
  - Background
  - PCM (or Technology X) as DRAM Replacement
  - Hybrid Memory Systems
    - Row-Locality Aware Data Placement
    - Efficient DRAM (or Technology X) Caches
- Conclusions
- Discussion

# The Problem with Large DRAM Caches

- A large DRAM cache requires a large metadata (tag + block-based information) store
- How do we design an efficient DRAM cache?



# Idea 1: Tags in Memory

---

- Store tags in the same row as data in DRAM
  - Store metadata in same row as their data
  - Data and metadata can be accessed together



- Benefit: No on-chip tag storage overhead
- Downsides:
  - Cache hit determined only after a DRAM access
  - Cache hit requires two DRAM accesses

# Idea 2: Cache Tags in SRAM

---

- Recall Idea 1: Store all metadata in DRAM
  - To reduce metadata storage overhead
- Idea 2: Cache in on-chip SRAM frequently-accessed metadata
  - Cache only a small amount to keep SRAM size small

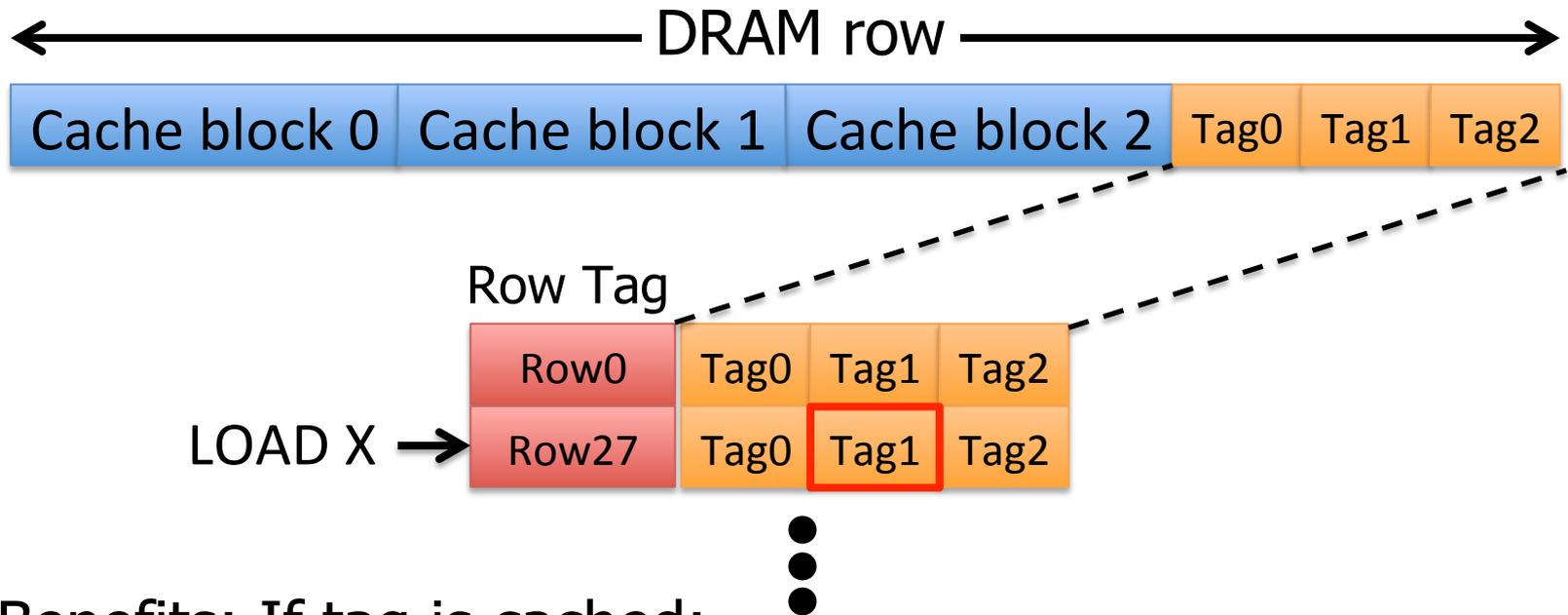
# Idea 3: Dynamic Data Transfer Granularity

---

- Some applications benefit from caching more data
  - They have good spatial locality
- Others do not
  - Large granularity wastes bandwidth and reduces cache utilization
- Idea 3: **Simple dynamic caching granularity policy**
  - Cost-benefit analysis to determine best DRAM cache block size
  - Group main memory into sets of rows
  - Some row sets follow a fixed caching granularity
  - The rest of main memory follows the best granularity
    - Cost-benefit analysis: access latency versus number of cachings
    - Performed every quantum

# TIMBER Tag Management

- A Tag-In-Memory Buffer (TIMBER)
  - Stores recently-used tags in a small amount of SRAM

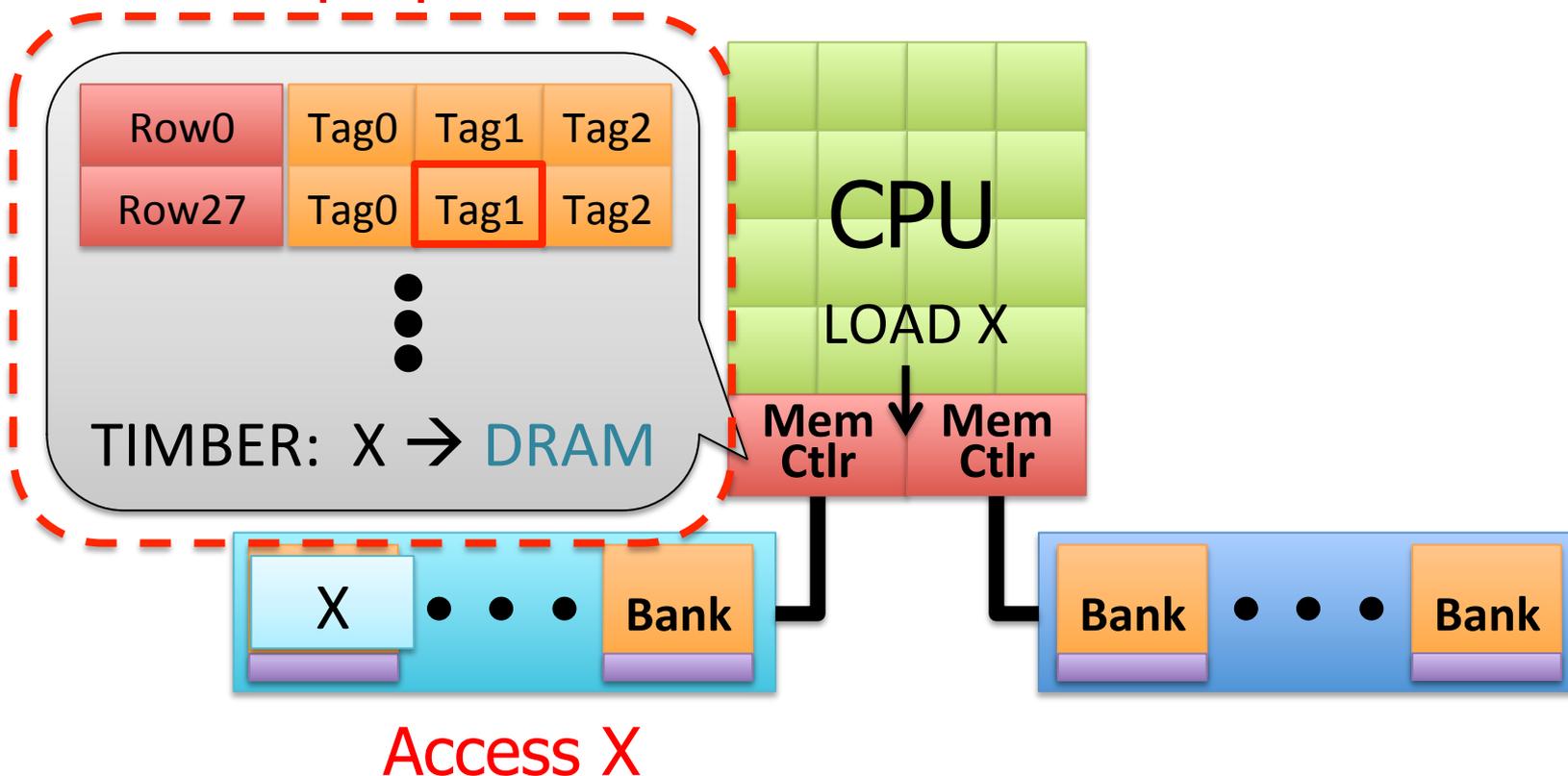


- Benefits: If tag is cached:
  - no need to access DRAM twice
  - cache hit determined quickly

# TIMBER Tag Management Example (I)

- Case 1: TIMBER hit

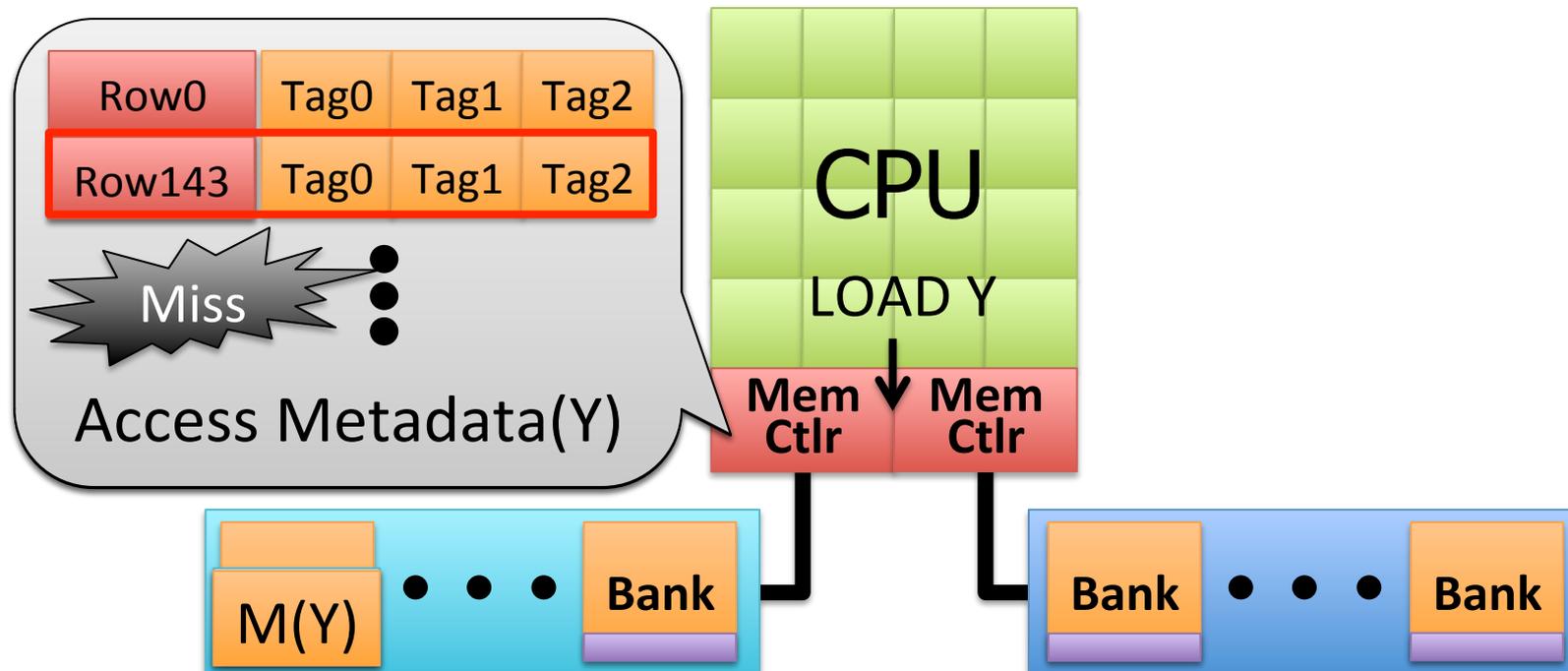
Our proposal



# TIMBER Tag Management Example (II)

- Case 2: TIMBER miss

## 2. Cache M(Y)



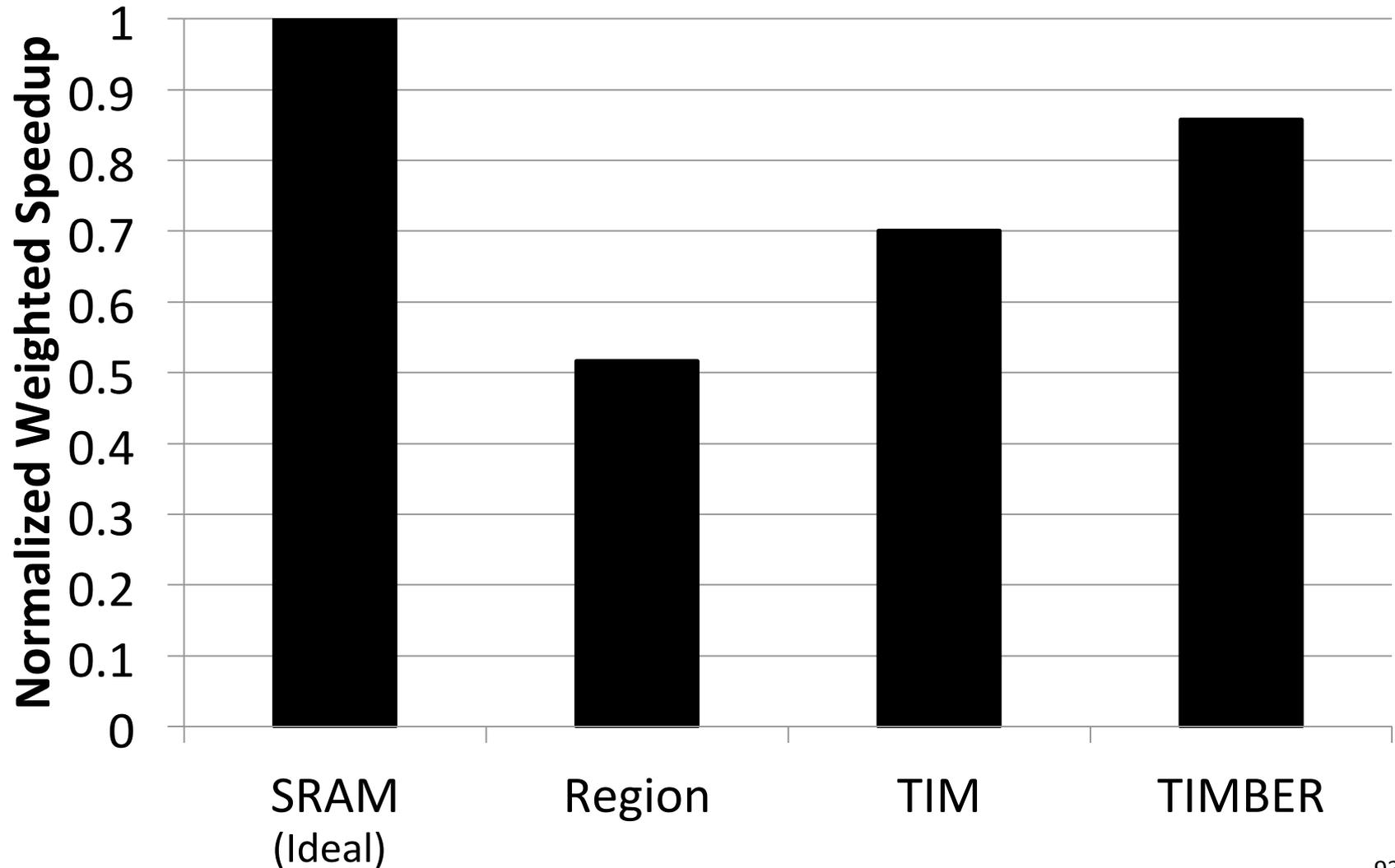
1. Access M(Y)
3. Access Y (row hit)

# Methodology

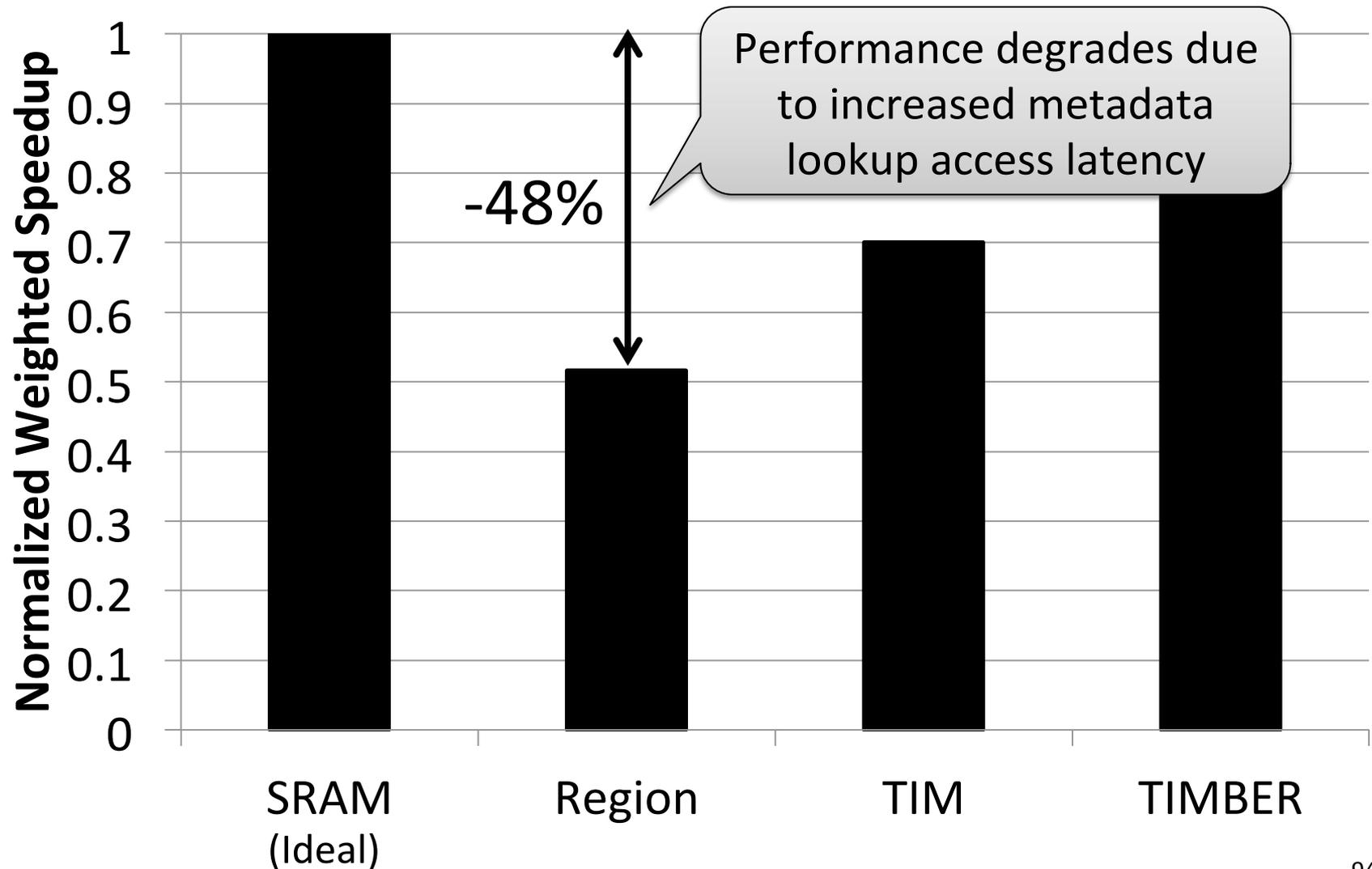
---

- System: 8 out-of-order cores at 4 GHz
- Memory: 512 MB direct-mapped DRAM, 8 GB PCM
  - 128B caching granularity
  - DRAM row hit (miss): 200 cycles (400 cycles)
  - PCM row hit (clean / dirty miss): 200 cycles (640 / 1840 cycles)
- Evaluated metadata storage techniques
  - All SRAM system (8MB of SRAM)
  - Region metadata storage
  - TIM metadata storage (same row as data)
  - TIMBER, 64-entry direct-mapped (8KB of SRAM)

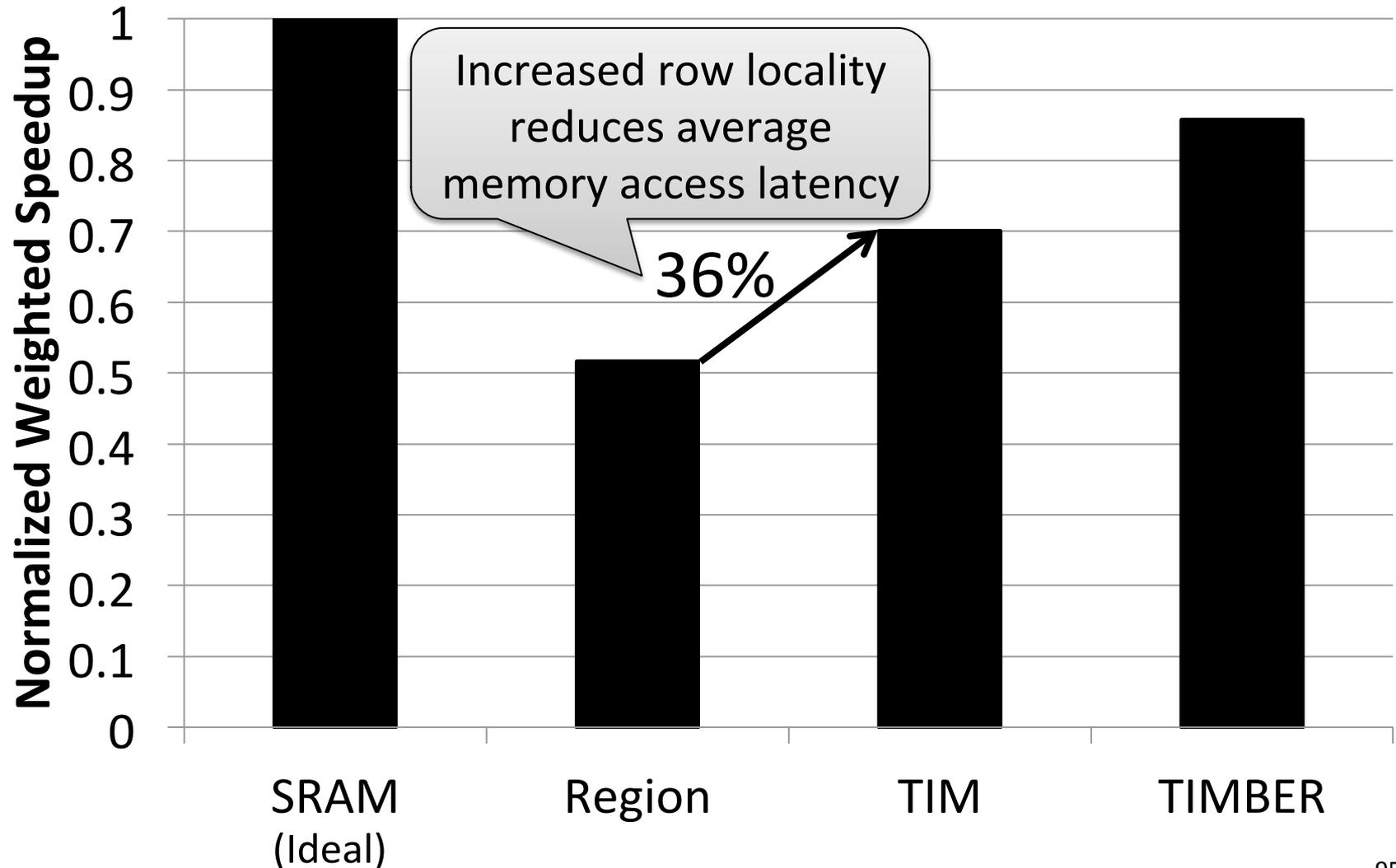
# Metadata Storage Performance



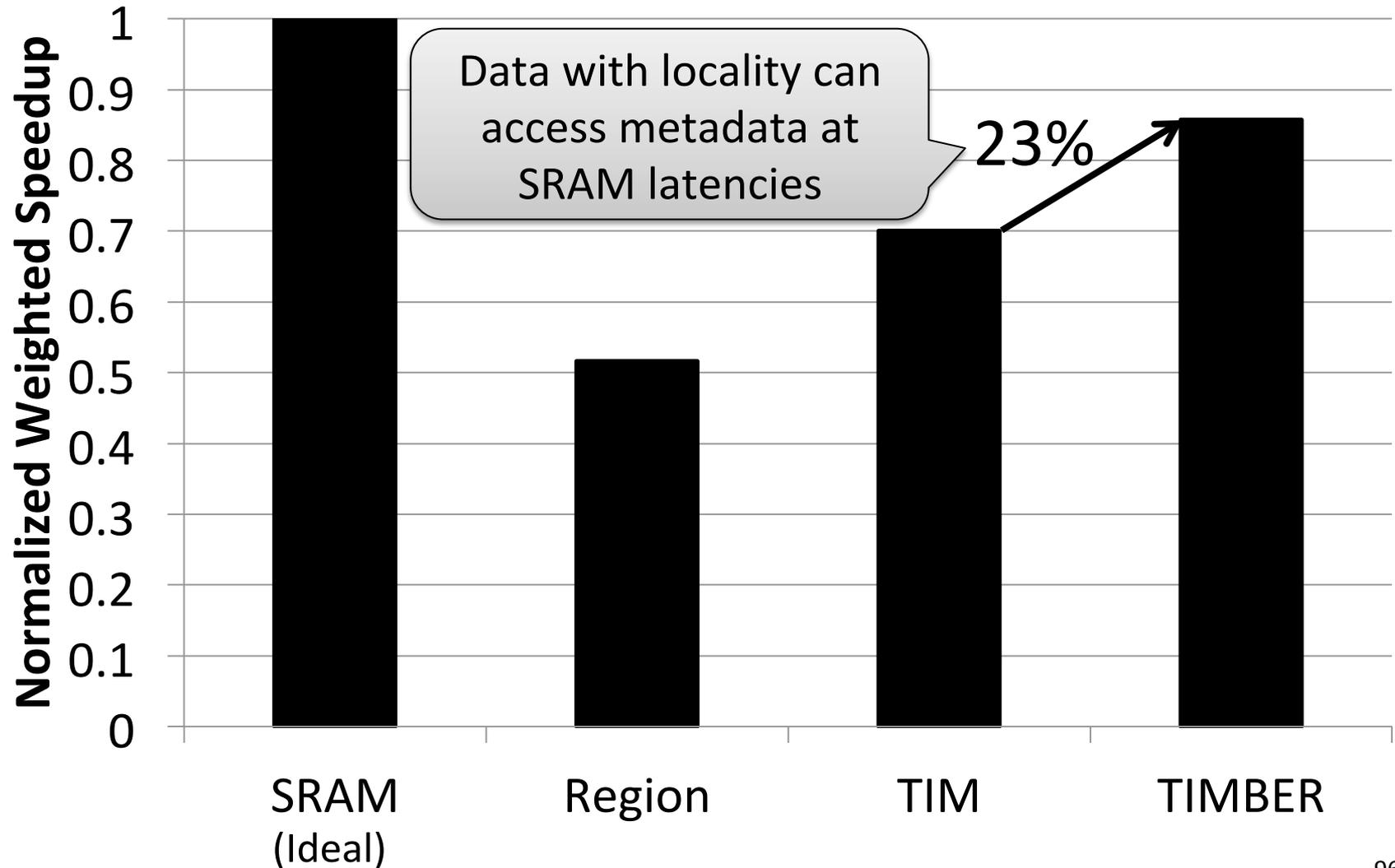
# Metadata Storage Performance



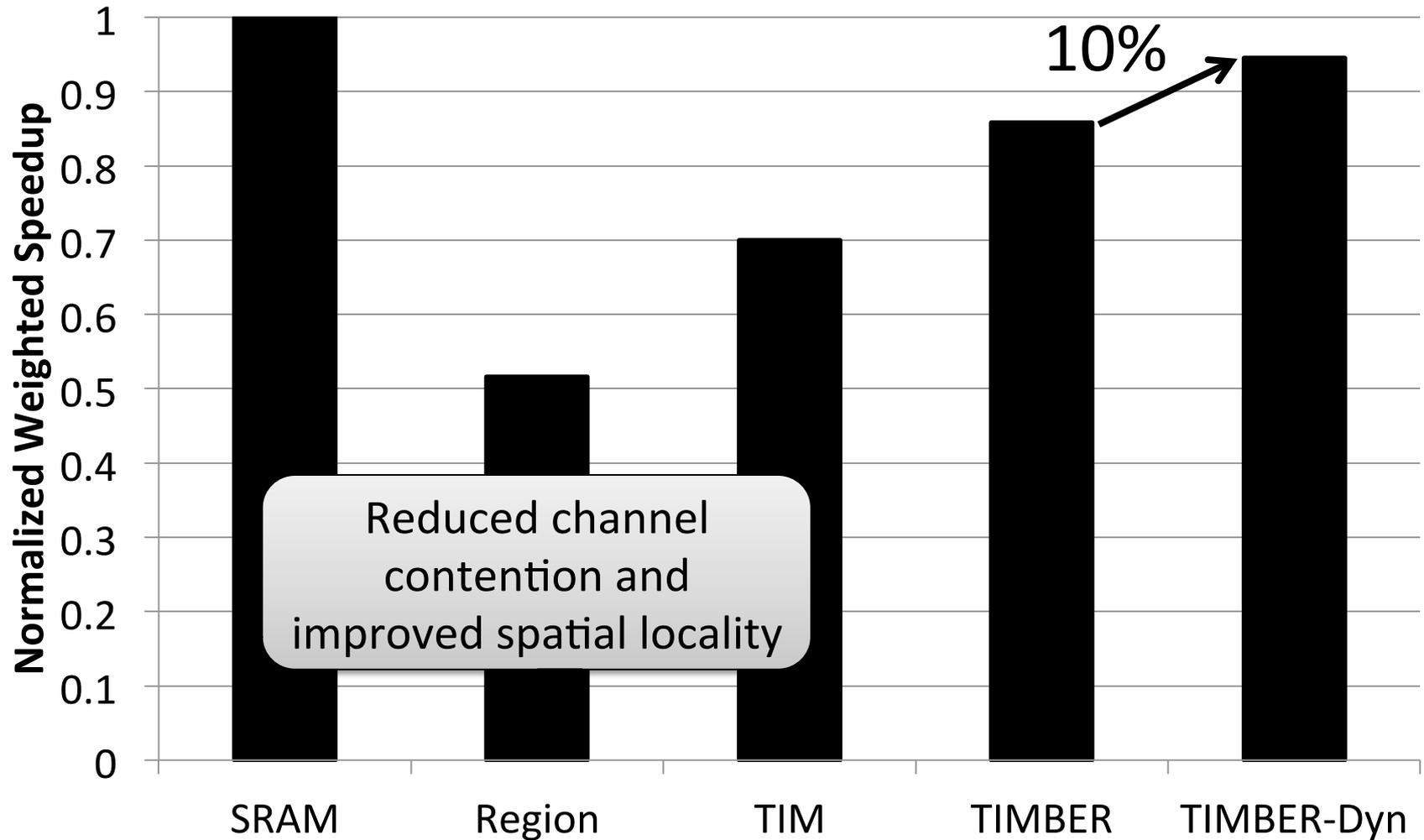
# Metadata Storage Performance



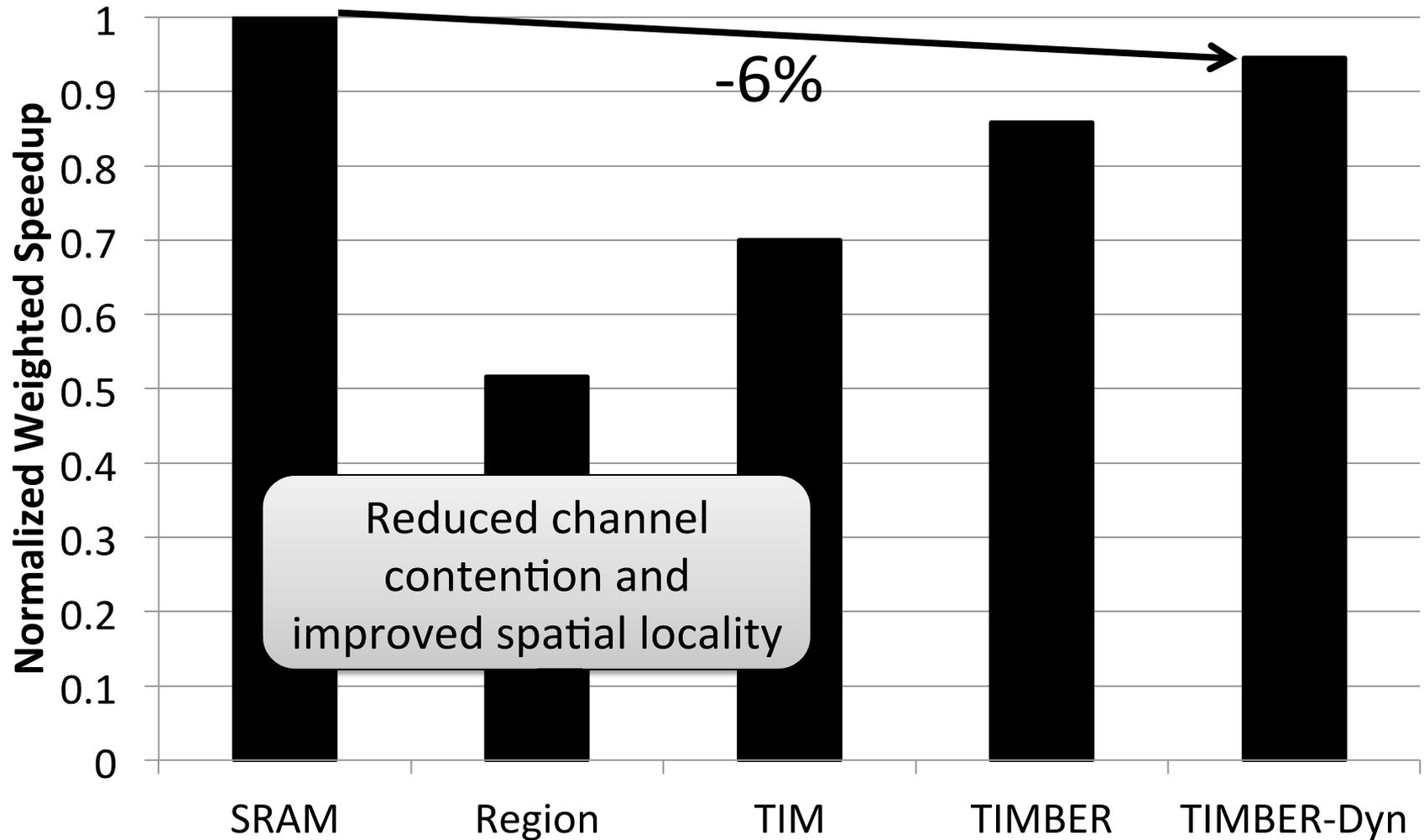
# Metadata Storage Performance



# Dynamic Granularity Performance

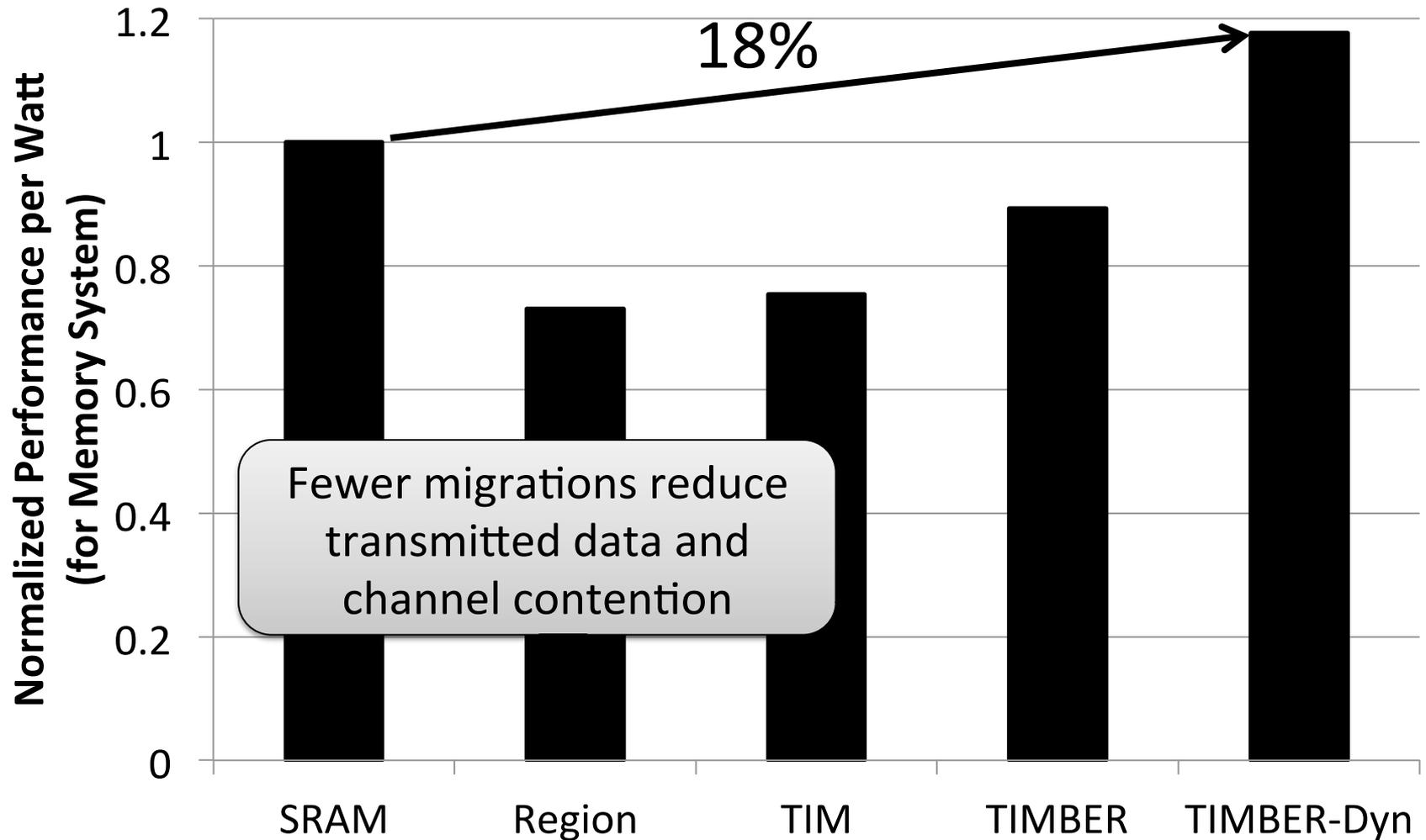


# TIMBER Performance



Meza, Chang, Yoon, Mutlu, Ranganathan, "Enabling Efficient and Scalable Hybrid Memories," IEEE Comp. Arch. Letters, 2012.

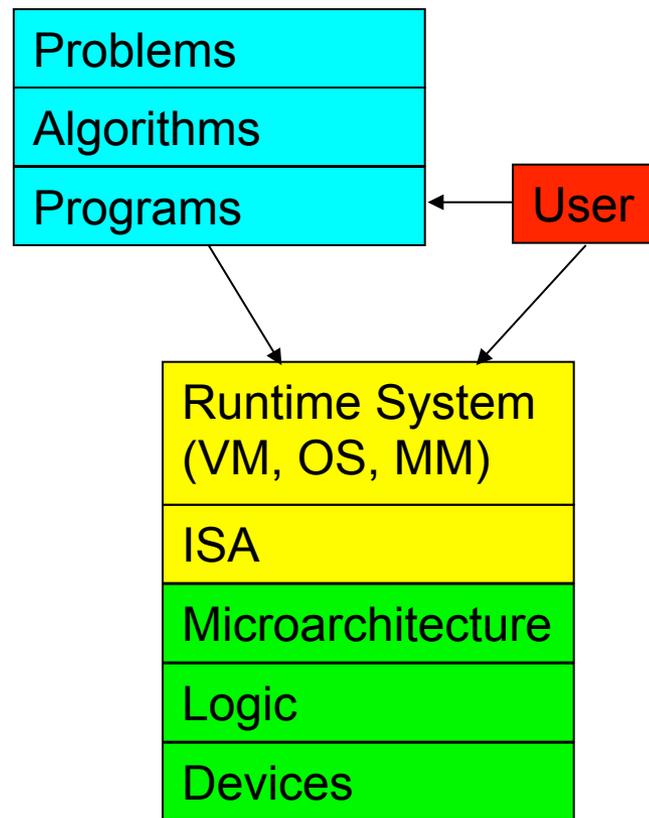
# TIMBER Energy Efficiency



Meza, Chang, Yoon, Mutlu, Ranganathan, "Enabling Efficient and Scalable Hybrid Memories," IEEE Comp. Arch. Letters, 2012.

# Enabling and Exploiting NVM: Issues

- Many issues and ideas from technology layer to algorithms layer
- Enabling NVM and hybrid memory
  - How to **tolerate errors**?
  - How to **enable secure operation**?
  - How to **tolerate performance and power shortcomings**?
  - How to **minimize cost**?
- Exploiting emerging technologies
  - How to **exploit non-volatility**?
  - How to **minimize energy consumption**?
  - How to **exploit NVM on chip**?



# Security Challenges of Emerging Technologies

---

1. Limited endurance → **Wearout attacks**
2. Non-volatility → Data persists in memory after powerdown  
→ **Easy retrieval of privileged or private information**
3. Multiple bits per cell → **Information leakage (via side channel)**

# Securing Emerging Memory Technologies

---

## 1. Limited endurance → **Wearout attacks**

Better architecting of memory chips to absorb writes

Hybrid memory system management

Online wearout attack detection

## 2. Non-volatility → Data persists in memory after powerdown

→ **Easy retrieval of privileged or private information**

Efficient encryption/decryption of whole main memory

Hybrid memory system management

## 3. Multiple bits per cell → **Information leakage (via side channel)**

System design to hide side channel information

# Agenda

---

- Major Trends Affecting Main Memory
- Requirements from an Ideal Main Memory System
- Opportunity: Emerging Memory Technologies
  - Background
  - PCM (or Technology X) as DRAM Replacement
  - Hybrid Memory Systems
- **Conclusions**
- Discussion

# Summary: Memory Scaling (with NVM)

---

- Main memory scaling problems are a critical bottleneck for system performance, efficiency, and usability
- Solution 1: Tolerate DRAM (yesterday)
- Solution 2: Enable emerging memory technologies
  - Replace DRAM with NVM by architecting NVM chips well
  - Hybrid memory systems with automatic data management
- An exciting topic with many other solution directions & ideas
  - Hardware/software/device cooperation essential
  - Memory, storage, controller, software/app co-design needed
  - Coordinated management of persistent memory and storage
  - Application and hardware cooperative management of NVM

# Scalable Many-Core Memory Systems

## Topic 2: Emerging Technologies and Hybrid Memories

Prof. Onur Mutlu

<http://www.ece.cmu.edu/~omutlu>

[onur@cmu.edu](mailto:onur@cmu.edu)

HiPEAC ACACES Summer School 2013

July 15-19, 2013

**Carnegie Mellon**

# Additional Material

# Overview Papers on Two Topics

---

## ■ Merging of Memory and Storage

- Justin Meza, Yixin Luo, Samira Khan, Jishen Zhao, Yuan Xie, and Onur Mutlu,  
**"A Case for Efficient Hardware-Software Cooperative Management of Storage and Memory"**  
*Proceedings of the 5th Workshop on Energy-Efficient Design (WEED)*, Tel-Aviv, Israel, June 2013. [Slides \(pptx\)](#) [Slides \(pdf\)](#)

## ■ Flash Memory Scaling

- Yu Cai, Gulay Yalcin, Onur Mutlu, Erich F. Haratsch, Adrian Cristal, Osman Unsal, and Ken Mai,  
**"Error Analysis and Retention-Aware Error Management for NAND Flash Memory"**  
*Intel Technology Journal (ITJ) Special Issue on Memory Resiliency*, Vol. 17, No. 1, May 2013.

# Merging of Memory and Storage: Persistent Memory Managers

# A Case for Efficient Hardware/Software Cooperative Management of Storage and Memory

Justin Meza<sup>\*</sup>, Yixin Luo<sup>\*</sup>, Samira Khan<sup>\*†</sup>, Jishen Zhao<sup>§</sup>,  
Yuan Xie<sup>§‡</sup>, and **Onur Mutlu<sup>\*</sup>**

<sup>\*</sup>Carnegie Mellon University

<sup>§</sup>Pennsylvania State University

<sup>†</sup>Intel Labs    <sup>‡</sup>AMD Research

# Overview

---

- Traditional systems have a **two-level storage model**
  - Access **volatile** data in memory with a **load/store** interface
  - Access **persistent** data in storage with a **file system** interface
  - Problem: **Operating system (OS) and file system (FS) code and buffering for storage lead to energy and performance inefficiencies**
- Opportunity: New non-volatile memory (NVM) technologies can help provide fast (similar to DRAM), persistent storage (similar to Flash)
  - **Unfortunately, OS and FS code can easily become energy efficiency and performance bottlenecks if we keep the traditional storage model**
- This work: **makes a case for hardware/software cooperative management of storage and memory within a single-level**
  - We describe the idea of a Persistent Memory Manager (PMM) for efficiently coordinating storage and memory, and quantify its benefit
  - And, examine questions and challenges to address to realize PMM

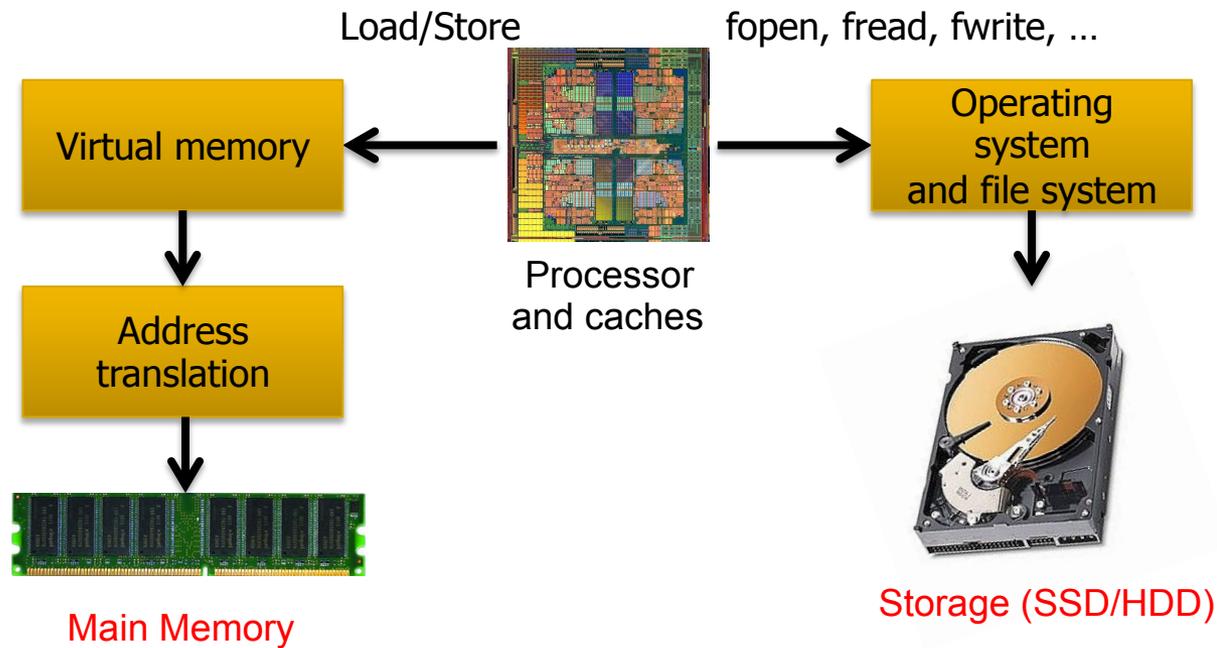
# Talk Outline

---

- Background: Storage and Memory Models
- Motivation: Eliminating Operating/File System Bottlenecks
- Our Proposal: Hardware/Software Coordinated Management of Storage and Memory
  - Opportunities and Benefits
- Evaluation Methodology
- Evaluation Results
- Related Work
- New Questions and Challenges
- Conclusions

# A Tale of Two Storage Levels

- Traditional systems use a two-level storage model
  - Volatile data is stored in DRAM
  - Persistent data is stored in HDD and Flash
- Accessed through two vastly different interfaces



# A Tale of Two Storage Levels

---

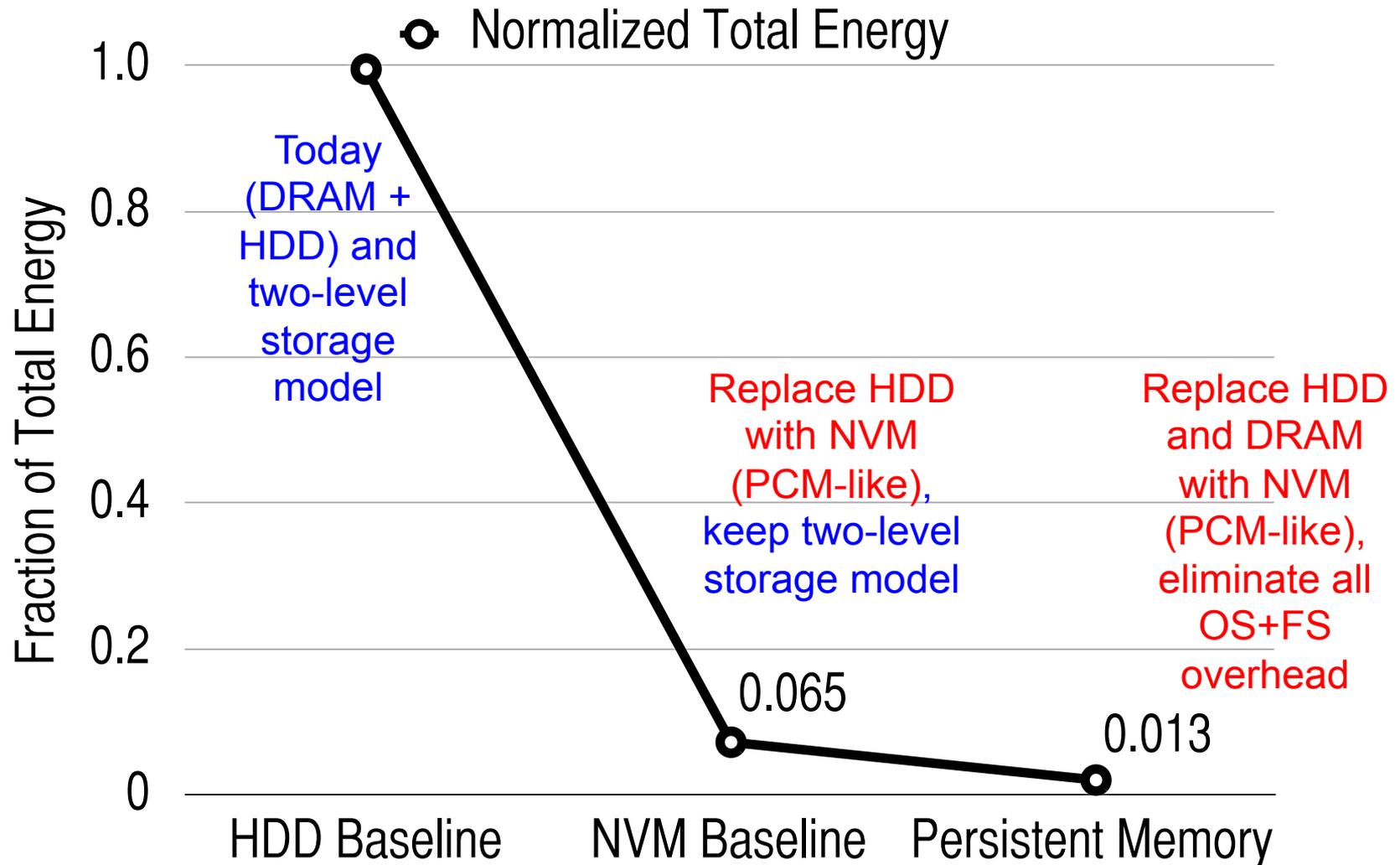
- Two-level storage arose in systems due to the widely different access latencies and methods of the commodity storage devices
  - Fast, low capacity, volatile DRAM → working storage
  - Slow, high capacity, non-volatile hard disk drives → persistent storage
- Data from slow storage media is buffered in fast DRAM
  - After that it can be manipulated by programs → programs cannot directly access persistent storage
  - It is the programmer's job to translate this data between the two formats of the two-level storage (files and data structures)
- Locating, transferring, and translating data and formats between the two levels of storage can waste significant energy and performance

# Opportunity: New Non-Volatile Memories

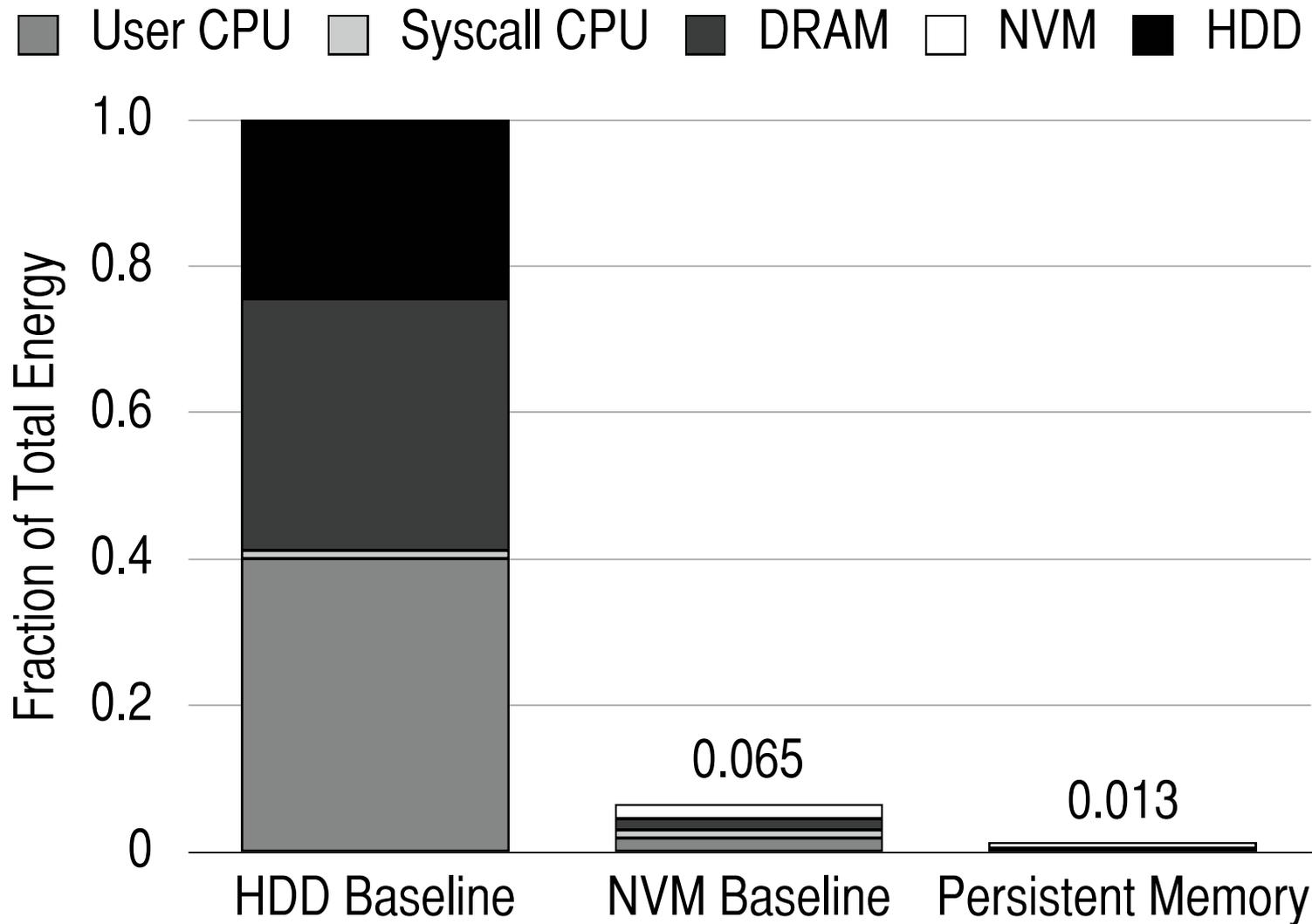
---

- Emerging memory technologies provide the potential for unifying storage and memory (e.g., Phase-Change, STT-RAM, RRAM)
  - **Byte-addressable** (can be accessed like DRAM)
  - **Low latency** (comparable to DRAM)
  - **Low power** (idle power better than DRAM)
  - **High capacity** (closer to Flash)
  - **Non-volatile** (can enable persistent storage)
  - **May have limited endurance** (but, better than Flash)
- Can provide fast access to **both** volatile data and persistent storage
- **Question: if such devices are used, is it efficient to keep a two-level storage model?**

# Eliminating Traditional Storage Bottlenecks

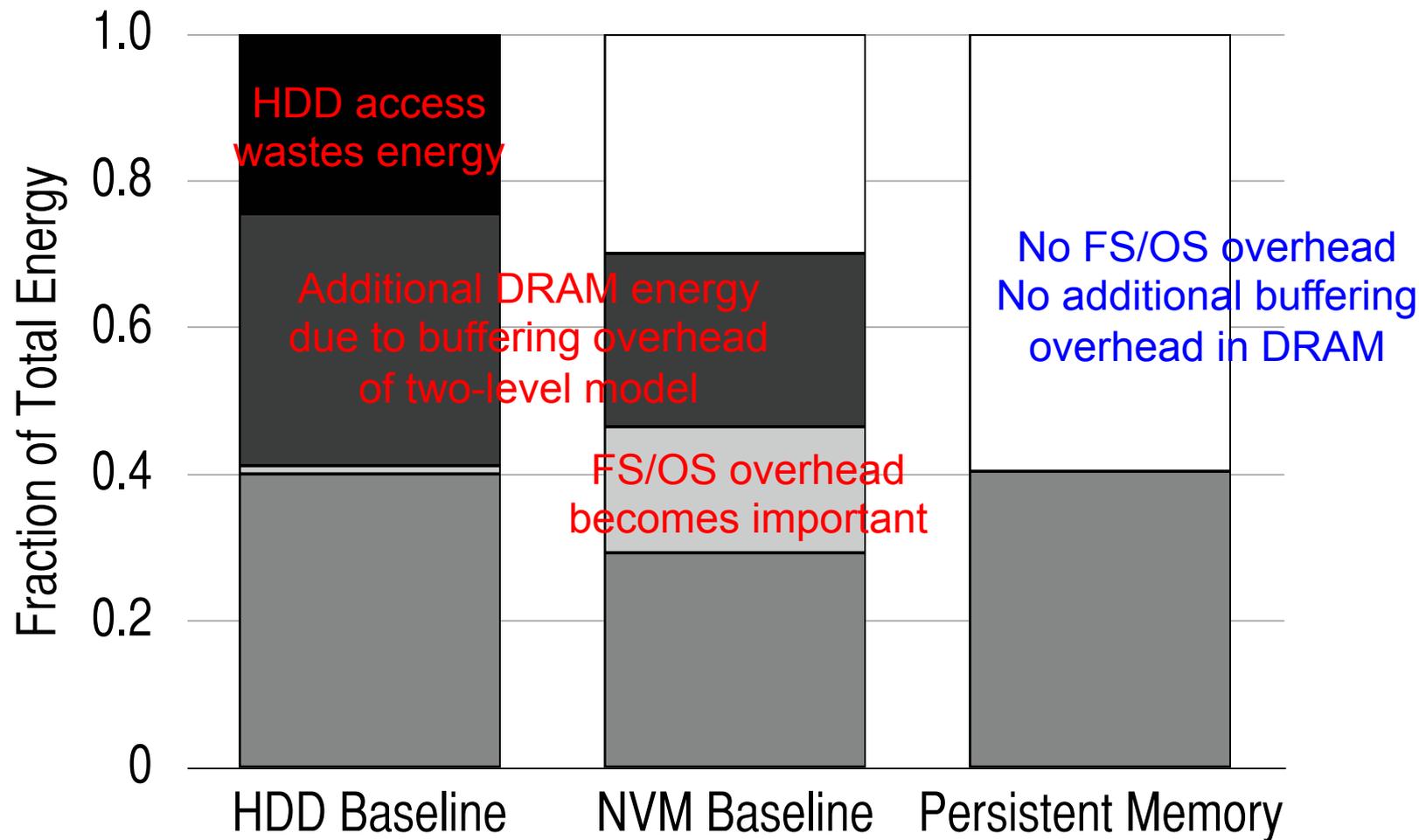


# Eliminating Traditional Storage Bottlenecks



# Where is Energy Spent in Each Model?

■ User CPU   ■ Syscall CPU   ■ DRAM   □ NVM   ■ HDD



# Outline

---

- Background: Storage and Memory Models
- Motivation: Eliminating Operating/File System Bottlenecks
- **Our Proposal: Hardware/Software Coordinated Management of Storage and Memory**
  - Opportunities and Benefits
- Evaluation Methodology
- Evaluation Results
- Related Work
- New Questions and Challenges
- Conclusions

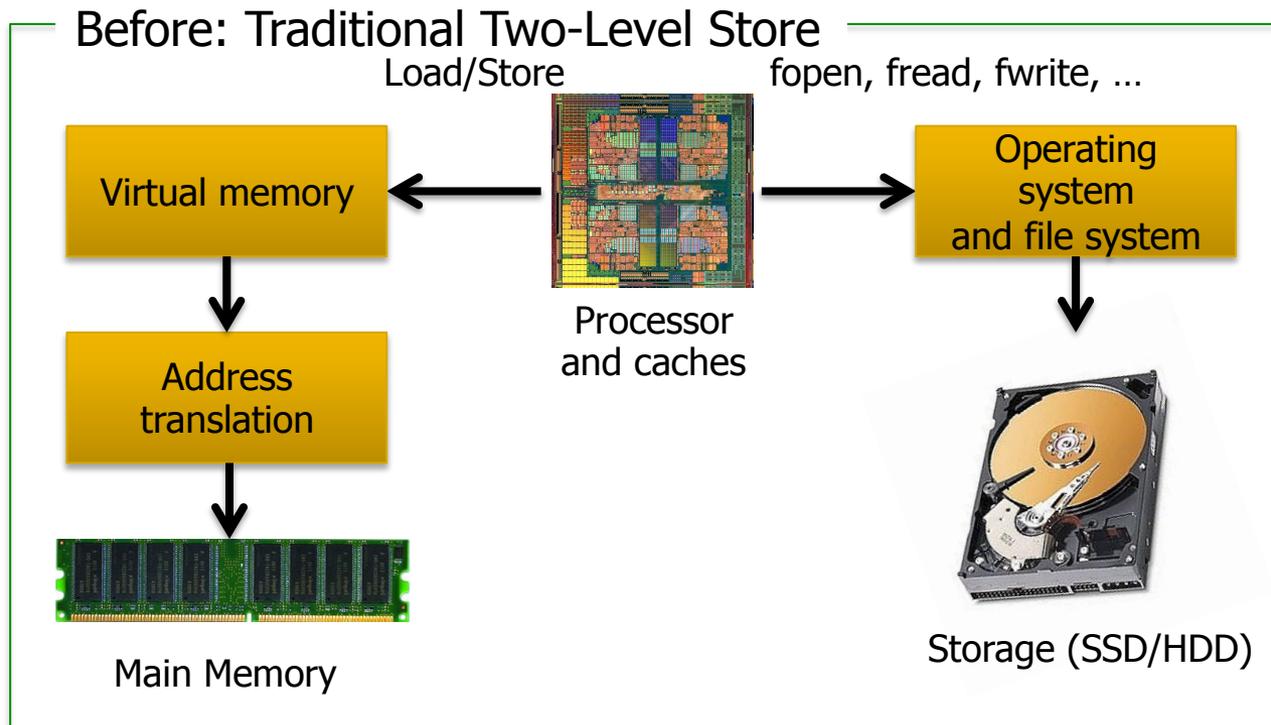
# Our Proposal: Coordinated HW/SW Memory and Storage Management

---

- Goal: Unify memory and storage to eliminate wasted work to locate, transfer, and translate data
  - Improve both energy and performance
  - Simplify programming model as well

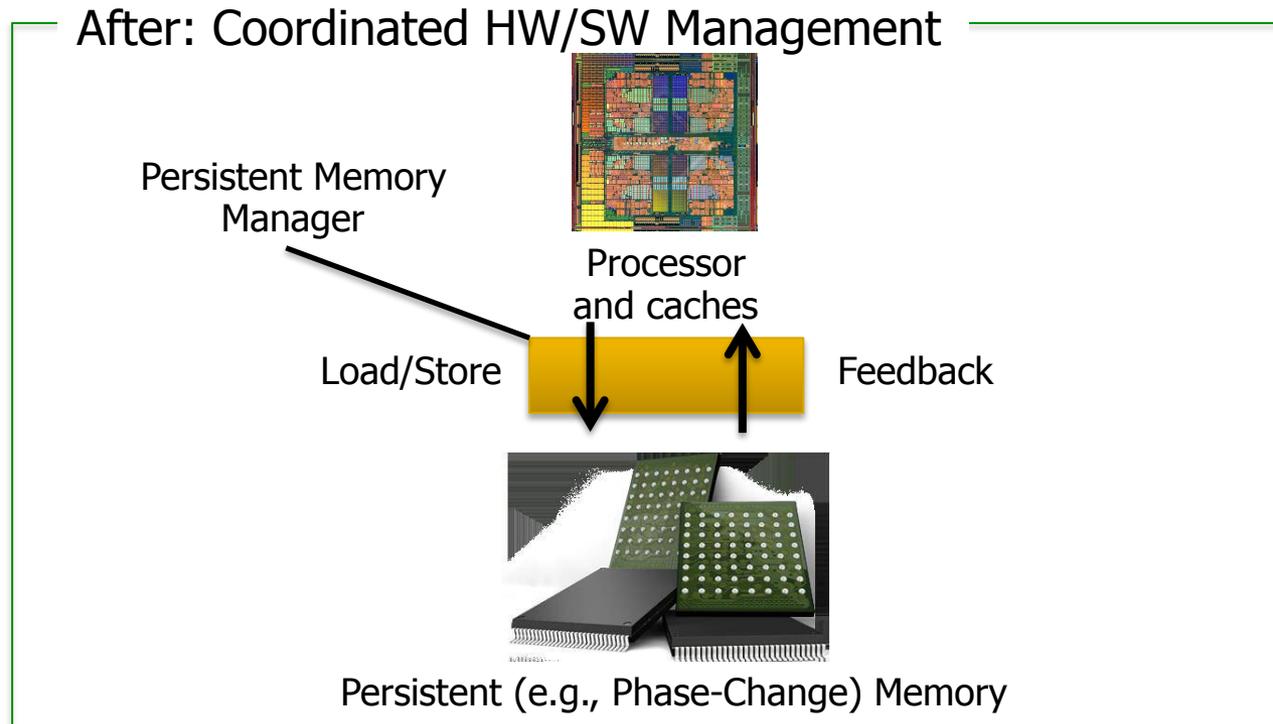
# Our Proposal: Coordinated HW/SW Memory and Storage Management

- Goal: Unify memory and storage to eliminate wasted work to locate, transfer, and translate data
  - Improve both energy and performance
  - Simplify programming model as well



# Our Proposal: Coordinated HW/SW Memory and Storage Management

- Goal: Unify memory and storage to eliminate wasted work to locate, transfer, and translate data
  - Improve both energy and performance
  - Simplify programming model as well



# The Persistent Memory Manager (PMM)

---

- Exposes a load/store interface to access persistent data
  - Applications can directly access persistent memory → no conversion, translation, location overhead for persistent data
- Manages data placement, location, persistence, security
  - To get the best of multiple forms of storage
- Manages metadata storage and retrieval
  - This can lead to overheads that need to be managed
- Exposes hooks and interfaces for system software
  - To enable better data placement and management decisions

# The Persistent Memory Manager

- Persistent Memory Manager

- Exposes a load/store interface to access persistent data
- Manages data placement, location, persistence, security
- Manages metadata storage and retrieval
- Exposes hooks and interfaces for system software

- Example program manipulating a persistent object:

```
1 int main(void) {  
2     // data in file.dat is persistent  
3     FILE myData = "file.dat";  
4     myData = new int [64];  
5 }  
6 void updateValue(int n, int value) {  
7     FILE myData = "file.dat";  
8     myData[n] = value; // value is persistent  
9 }
```

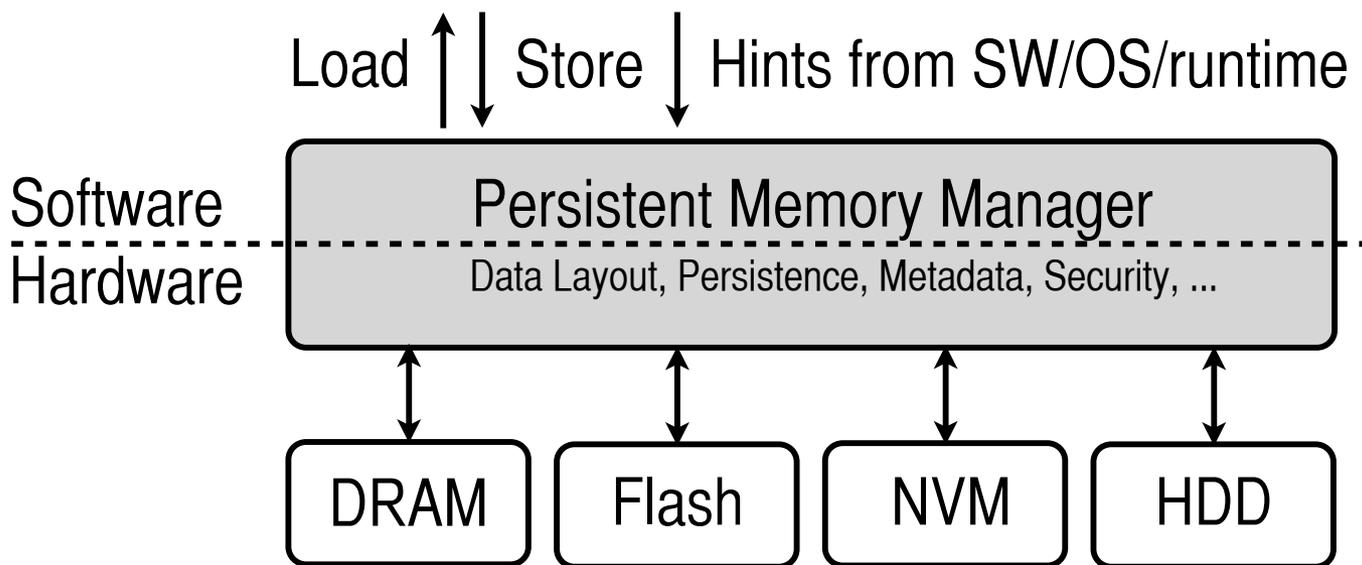
Create persistent object and its handle

Allocate a persistent array and assign

Load/store interface

# Putting Everything Together

```
1 int main(void) {  
2     // data in file.dat is persistent  
3     FILE myData = "file.dat";  
4     myData = new int[64];  
5 }  
6 void updateValue(int n, int value) {  
7     FILE myData = "file.dat";  
8     myData[n] = value; // value is persistent  
9 }
```



**PMM uses access and hint information to allocate, locate, migrate and access data in the heterogeneous array of devices**

# Outline

---

- Background: Storage and Memory Models
- Motivation: Eliminating Operating/File System Bottlenecks
- Our Proposal: Hardware/Software Coordinated Management of Storage and Memory
  - Opportunities and Benefits
- Evaluation Methodology
- Evaluation Results
- Related Work
- New Questions and Challenges
- Conclusions

# Opportunities and Benefits

---

- We've identified at least five opportunities and benefits of a unified storage/memory system that gets rid of the two-level model:
  1. Eliminating system calls for file operations
  2. Eliminating file system operations
  3. Efficient data mapping/location among heterogeneous devices
  4. Providing security and reliability in persistent memories
  5. Hardware/software cooperative data management

# Eliminating System Calls for File Operations

---

- A persistent memory can expose a large, linear, persistent address space
  - Persistent storage objects can be directly manipulated with load/store operations
- This eliminates the need for layers of operating system code
  - Typically used for calls like `open`, `read`, and `write`
- Also eliminates OS file metadata
  - File descriptors, file buffers, and so on

# Eliminating File System Operations

---

- Locating files is traditionally done using a *file system*
  - Runs code and traverses structures in software to locate files
- Existing hardware structures for locating data in virtual memory can be extended and adapted to meet the needs of persistent memories
  - Memory Management Units (MMUs), which map virtual addresses to physical addresses
  - Translation Lookaside Buffers (TLBs), which cache mappings of virtual-to-physical address translations
- Potential to eliminate file system code
- At the cost of additional hardware overhead to handle persistent data storage

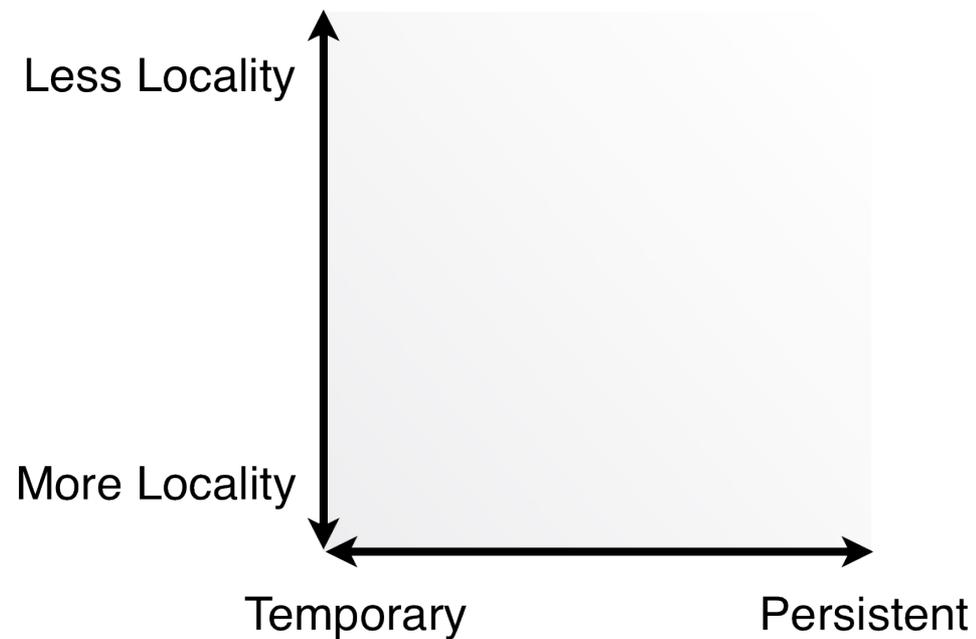
# Efficient Data Mapping among Heterogeneous Devices

---

- A persistent memory exposes a large, persistent address space
  - But it may use many different devices to satisfy this goal
  - From fast, low-capacity volatile DRAM to slow, high-capacity non-volatile HDD or Flash
  - And other NVM devices in between
- Performance and energy can benefit from good placement of data among these devices
  - Utilizing the strengths of each device and avoiding their weaknesses, if possible
  - For example, consider two important application characteristics: locality and persistence

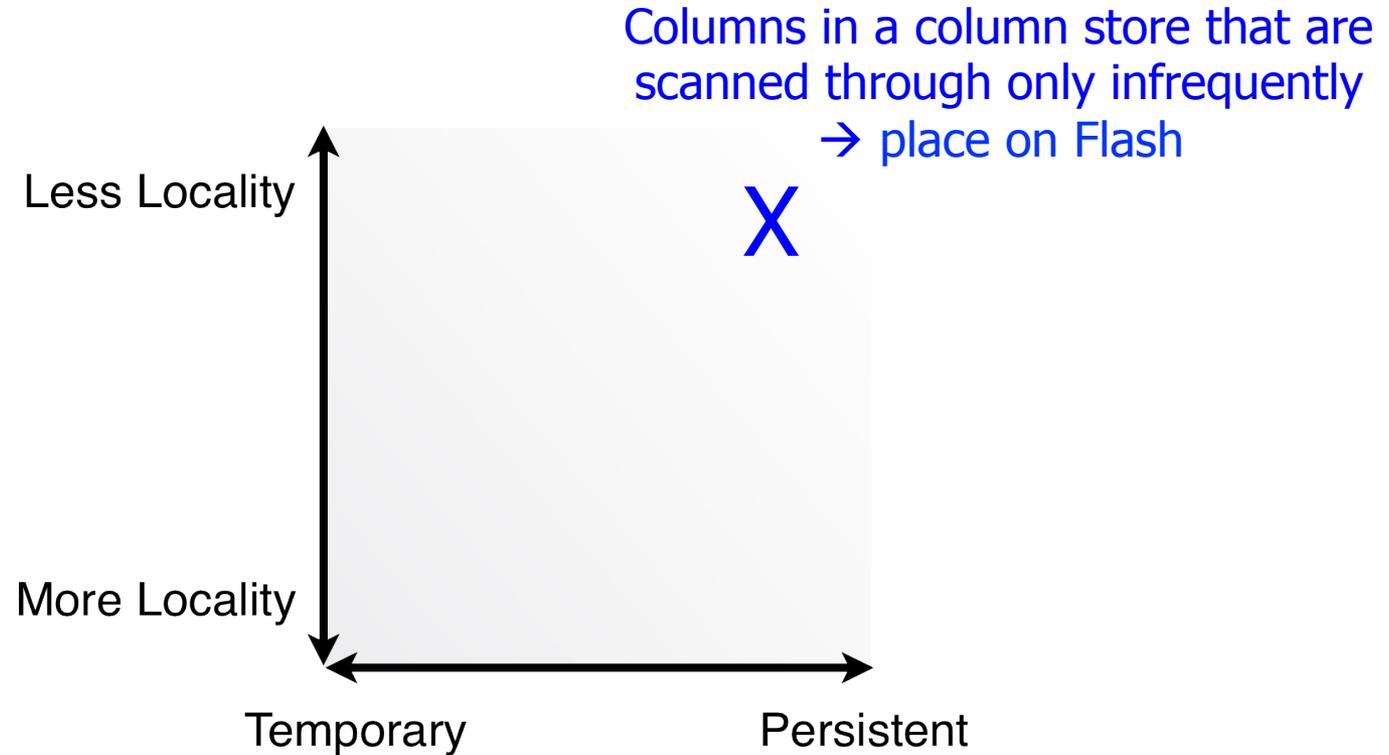
# Efficient Data Mapping among Heterogeneous Devices

---



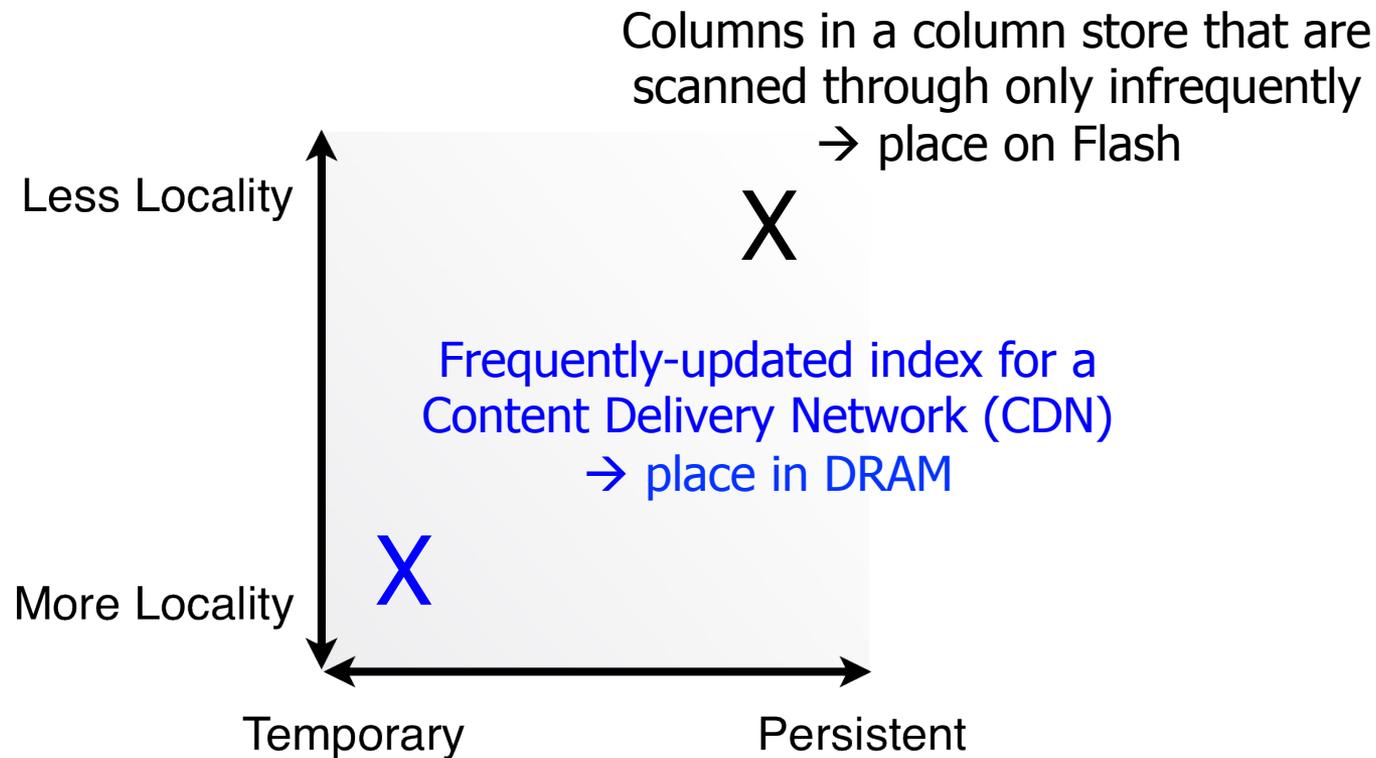
# Efficient Data Mapping among Heterogeneous Devices

---



# Efficient Data Mapping among Heterogeneous Devices

---



**Applications or system software can provide hints for data placement**

# Providing Security and Reliability

---

- A persistent memory deals with data at the granularity of bytes and not necessarily files
  - Provides the opportunity for much finer-grained security and protection than traditional two-level storage models provide/afford
  - Need efficient techniques to avoid large metadata overheads
- A persistent memory can improve application reliability by ensuring updates to persistent data are less vulnerable to failures
  - Need to ensure that changes to copies of persistent data placed in volatile memories become persistent

# HW/SW Cooperative Data Management

---

- Persistent memories can expose hooks and interfaces to applications, the OS, and runtimes
  - Have the potential to provide improved system robustness and efficiency than by managing persistent data with either software or hardware alone
- Can enable fast checkpointing and reboots, improve application reliability by ensuring persistence of data
  - How to redesign availability mechanisms to take advantage of these?
- Persistent locks and other persistent synchronization constructs can enable more robust programs and systems

# Quantifying Persistent Memory Benefits

---

- We have identified several opportunities and benefits of using persistent memories without the traditional two-level store model
- We will next quantify:
  - How do persistent memories affect system performance?
  - How much energy reduction is possible?
  - Can persistent memories achieve these benefits despite additional access latencies to the persistent memory manager?

# Outline

---

- Background: Storage and Memory Models
- Motivation: Eliminating Operating/File System Bottlenecks
- Our Proposal: Hardware/Software Coordinated Management of Storage and Memory
  - Opportunities and Benefits
- Evaluation Methodology
- Evaluation Results
- Related Work
- New Questions and Challenges
- Conclusions

# Evaluation Methodology

---

- Hybrid real system / simulation-based approach
  - System calls are executed on host machine (functional correctness) and timed to accurately model their latency in the simulator
  - Rest of execution is simulated in Multi2Sim (enables hardware-level exploration)
- Power evaluated using McPAT and memory power models
- 16 cores, 4-wide issue, 128-entry instruction window, 1.6 GHz
- Volatile memory: 4GB DRAM, 4KB page size, 100-cycle latency
- Persistent memory
  - HDD (measured): 4ms seek latency, 6Gbps bus rate
  - NVM: (modeled after PCM) 4KB page size, 160-/480-cycle (read/write) latency

# Evaluated Systems

---

- HDD Baseline (HB)
  - Traditional system with volatile DRAM memory and persistent HDD storage
  - Overheads of operating system and file system code and buffering
- HDD without OS/FS (HW)
  - Same as HDD Baseline, but with the ideal elimination of all OS/FS overheads
  - System calls take 0 cycles (but HDD access takes normal latency)
- NVM Baseline (NB)
  - Same as HDD Baseline, but HDD is replaced with NVM
  - Still has OS/FS overheads of the two-level storage model
- Persistent Memory (PM)
  - Uses only NVM (no DRAM) to ensure full-system persistence
  - All data accessed using loads and stores
  - Does not waste energy on system calls
  - Data is manipulated directly on the NVM device

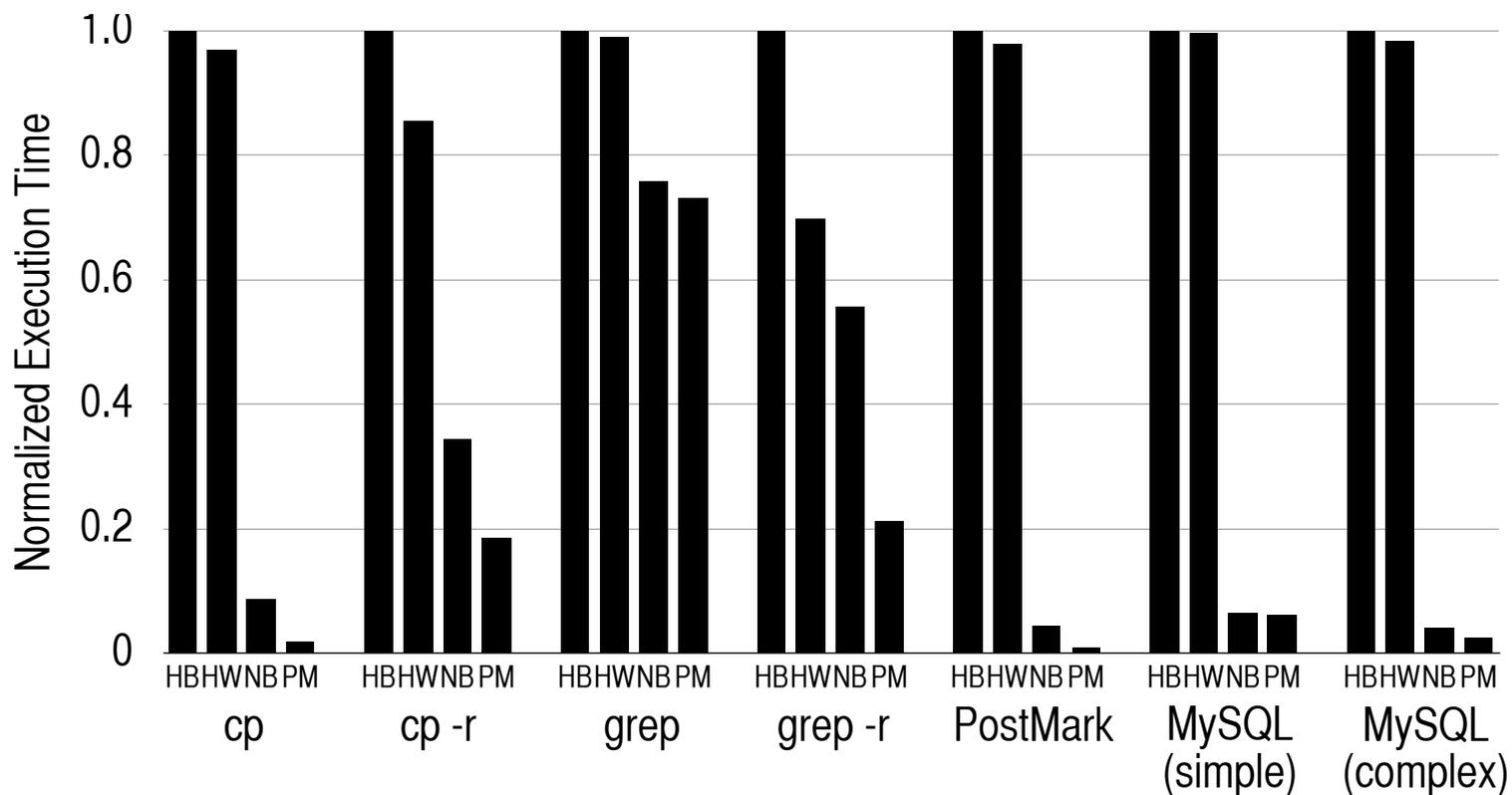
# Evaluated Workloads

---

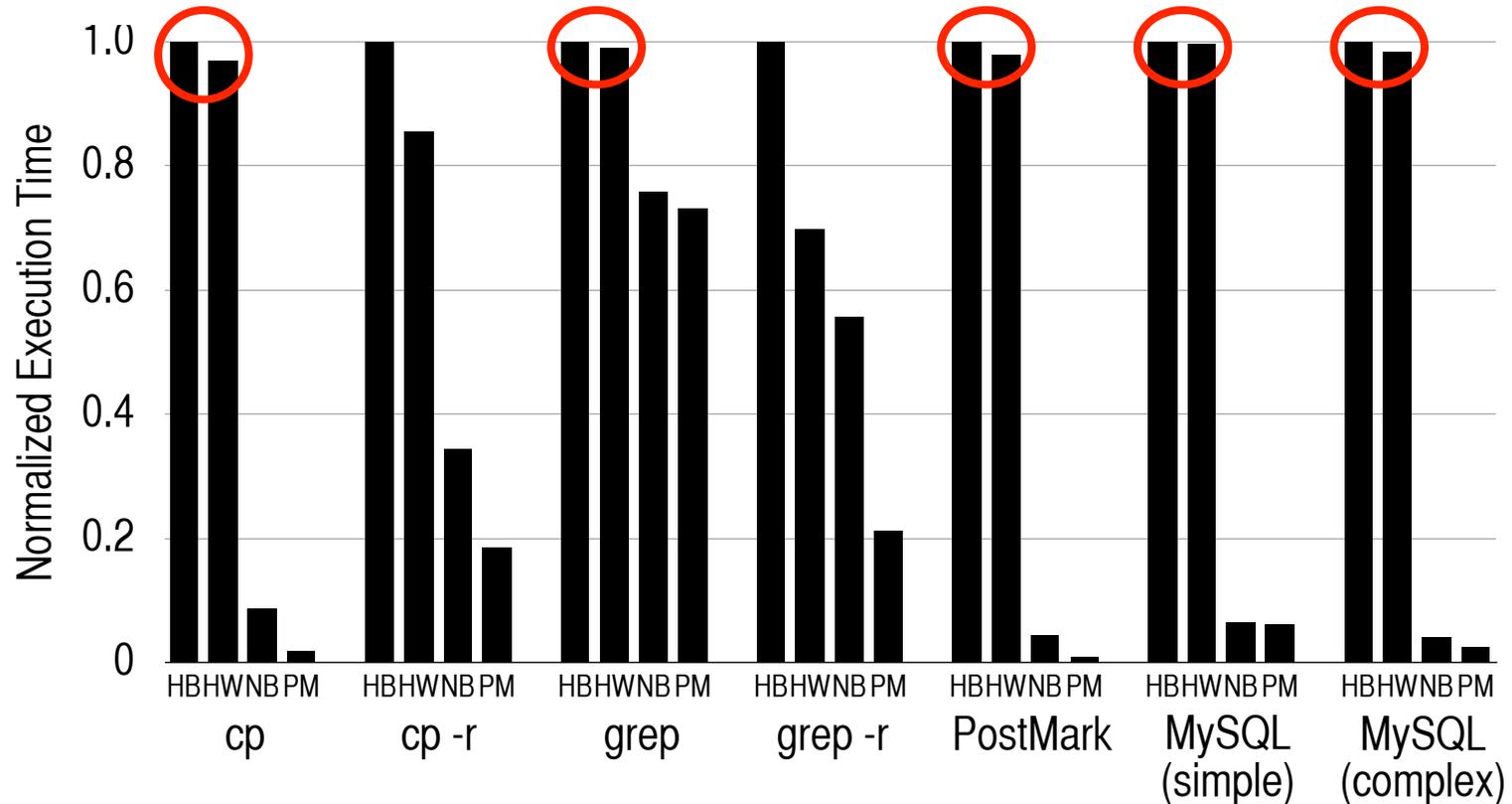
- Unix utilities that manipulate files
  - cp: copy a large file from one location to another
  - cp -r: copy files in a directory tree from one location to another
  - grep: search for a string in a large file
  - grep -r: search for a string recursively in a directory tree
- PostMark: an I/O-intensive benchmark from NetApp
  - Emulates typical access patterns for email, news, web commerce
- MySQL Server: a popular database management system
  - OLTP-style queries generated by Sysbench
  - MySQL (simple): single, random read to an entry
  - MySQL (complex): reads/writes 1 to 100 entries per transaction

# Performance Results

---

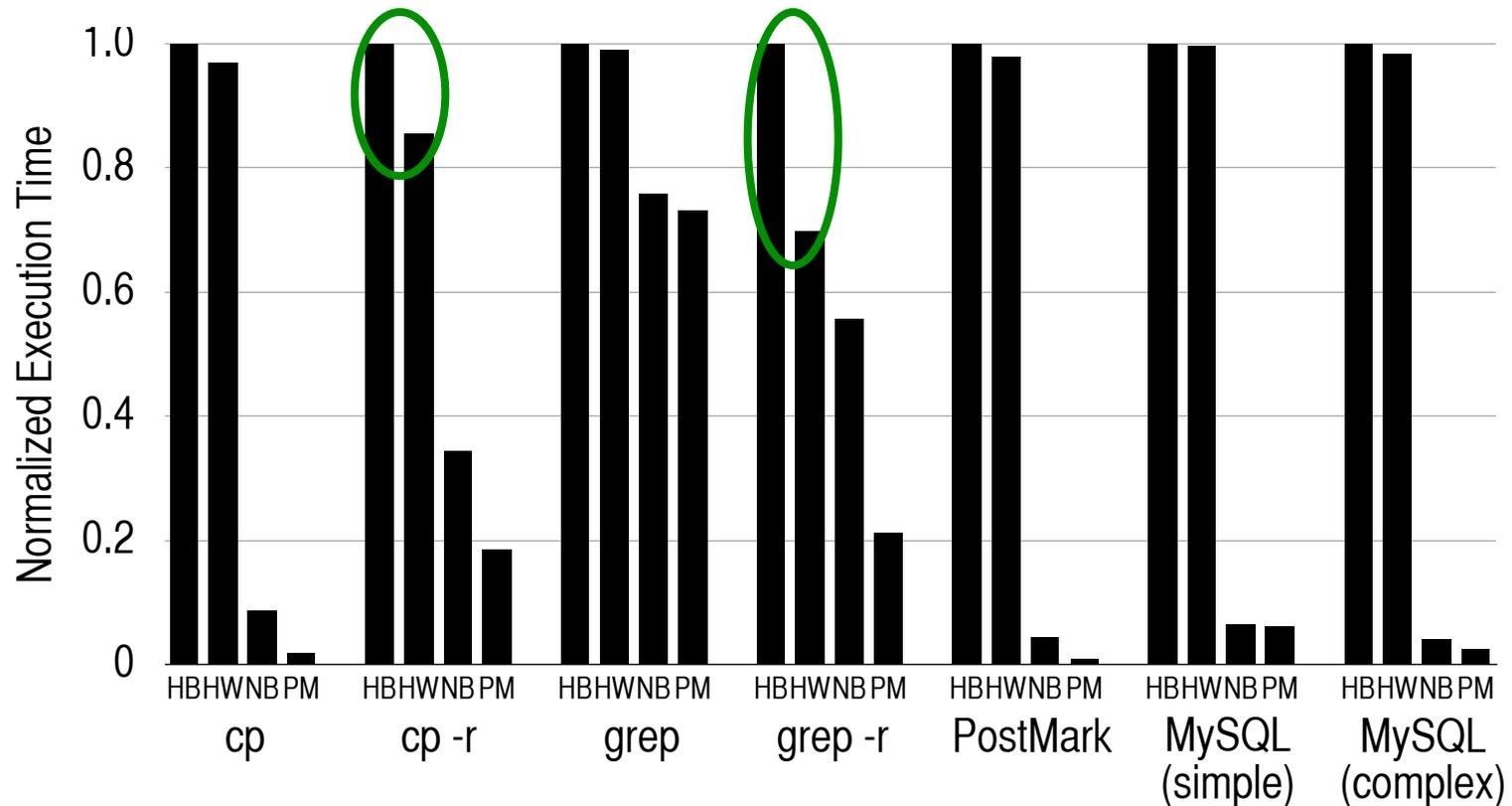


# Performance Results: HDD w/o OS/FS



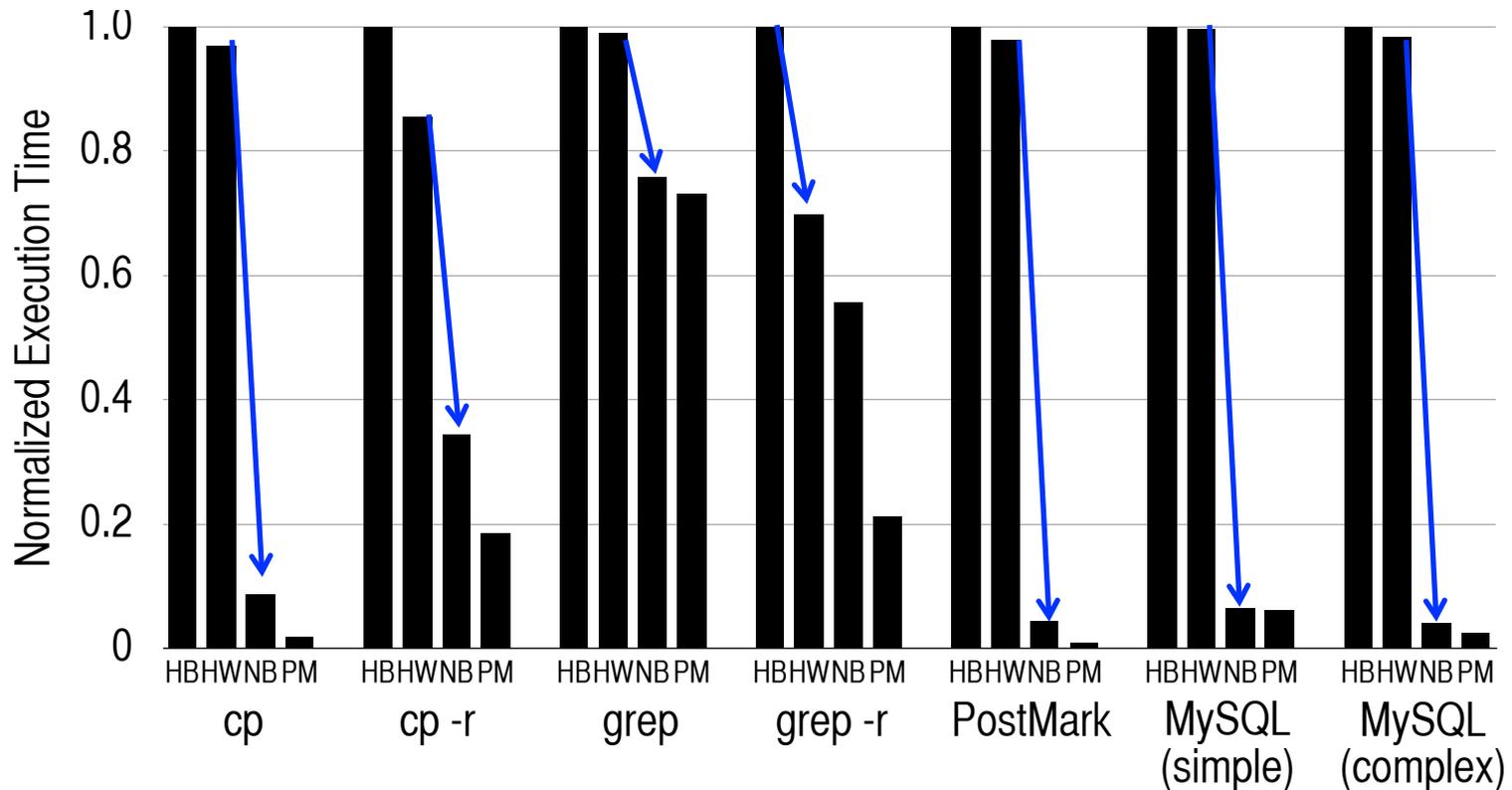
For HDD-based systems, eliminating OS/FS overheads typically leads to small performance improvements → execution time dominated by HDD access latency

# Performance Results: HDD w/o OS/FS



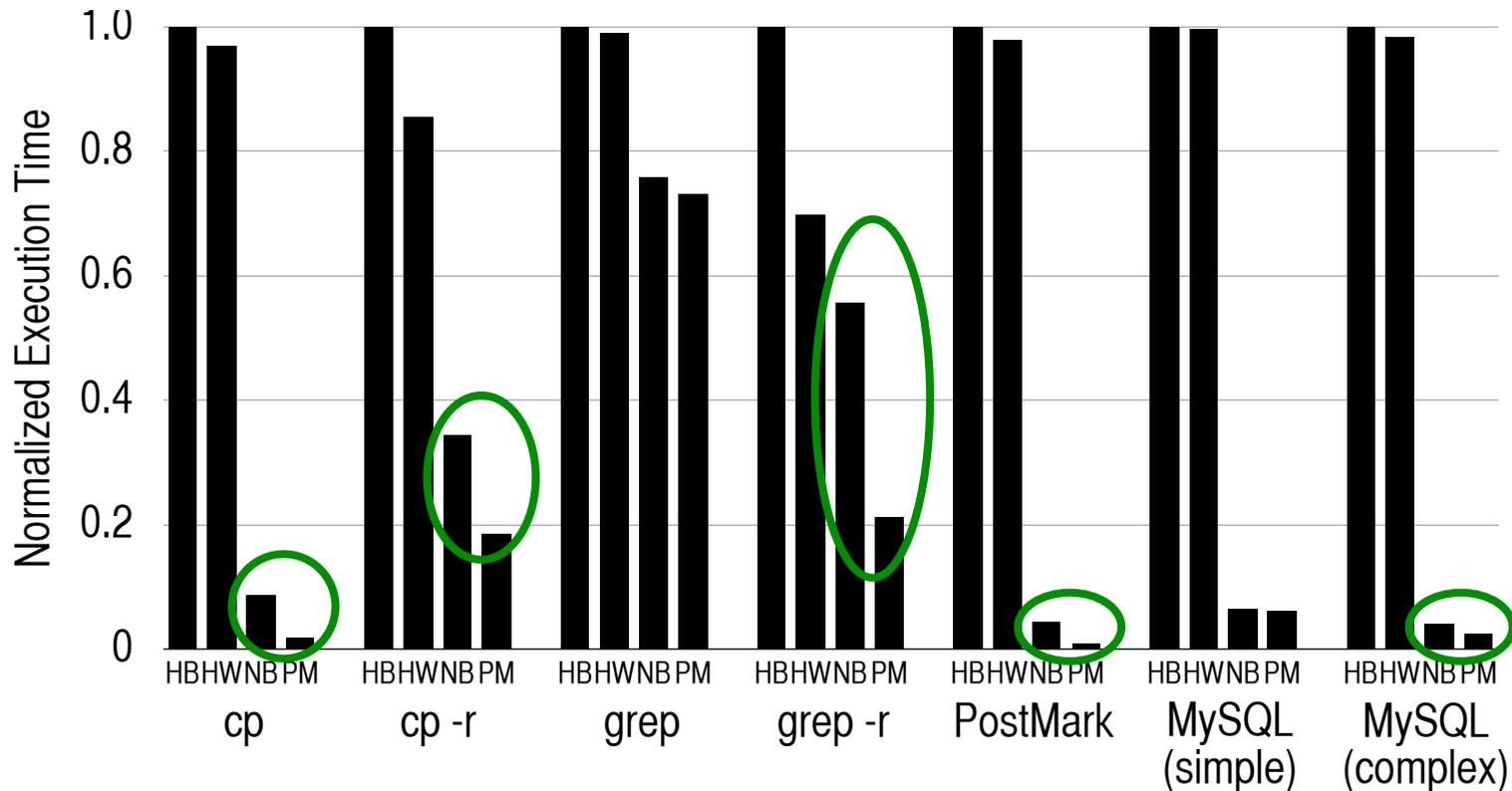
Though, for more complex file system operations like directory traversal (seen with `cp -r` and `grep -r`), eliminating the OS/FS overhead improves performance

# Performance Results: HDD to NVM



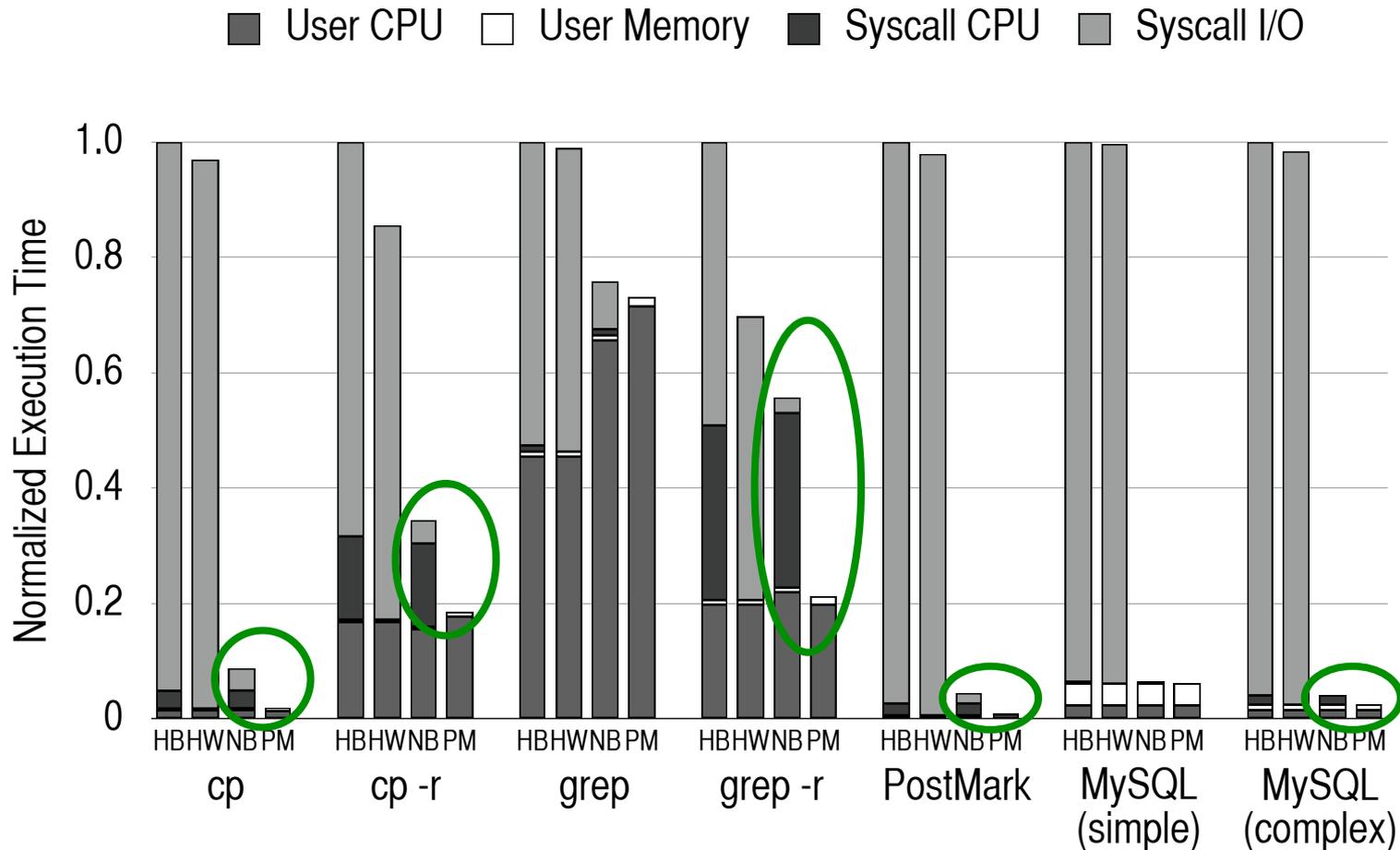
Switching from an HDD to NVM greatly reduces execution time due to NVM's much faster access latencies, especially for I/O-intensive workloads (cp, PostMark, MySQL)

# Performance Results: NVM to PMM



For most workloads, eliminating OS/FS code and buffering improves performance greatly on top of the NVM Baseline system (even when DRAM is eliminated from the system)

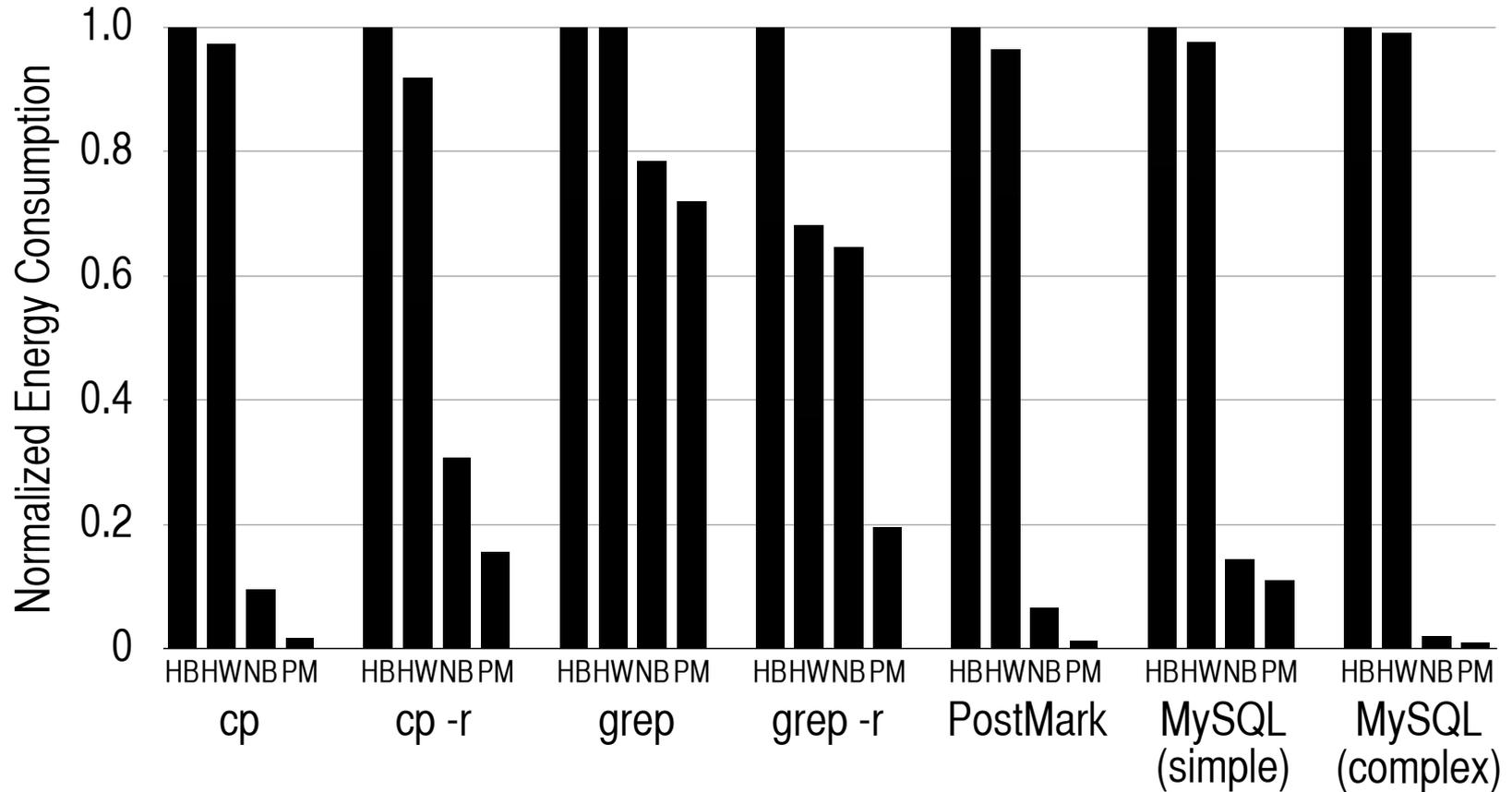
# Performance Results



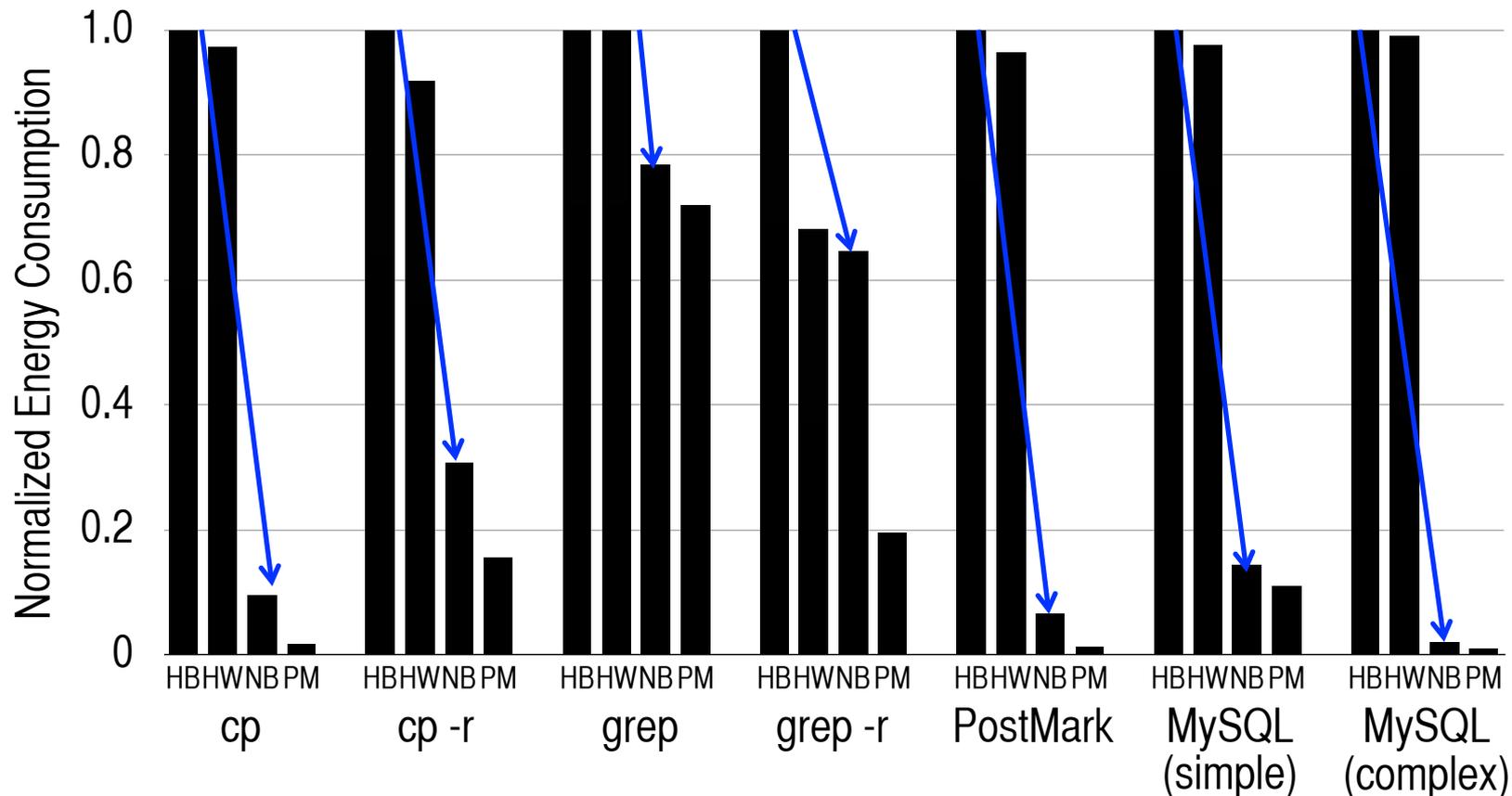
The workloads that see the greatest improvement from using a Persistent Memory are those that spend a large portion of their time executing system call code due to the two-level storage model

# Energy Results

---

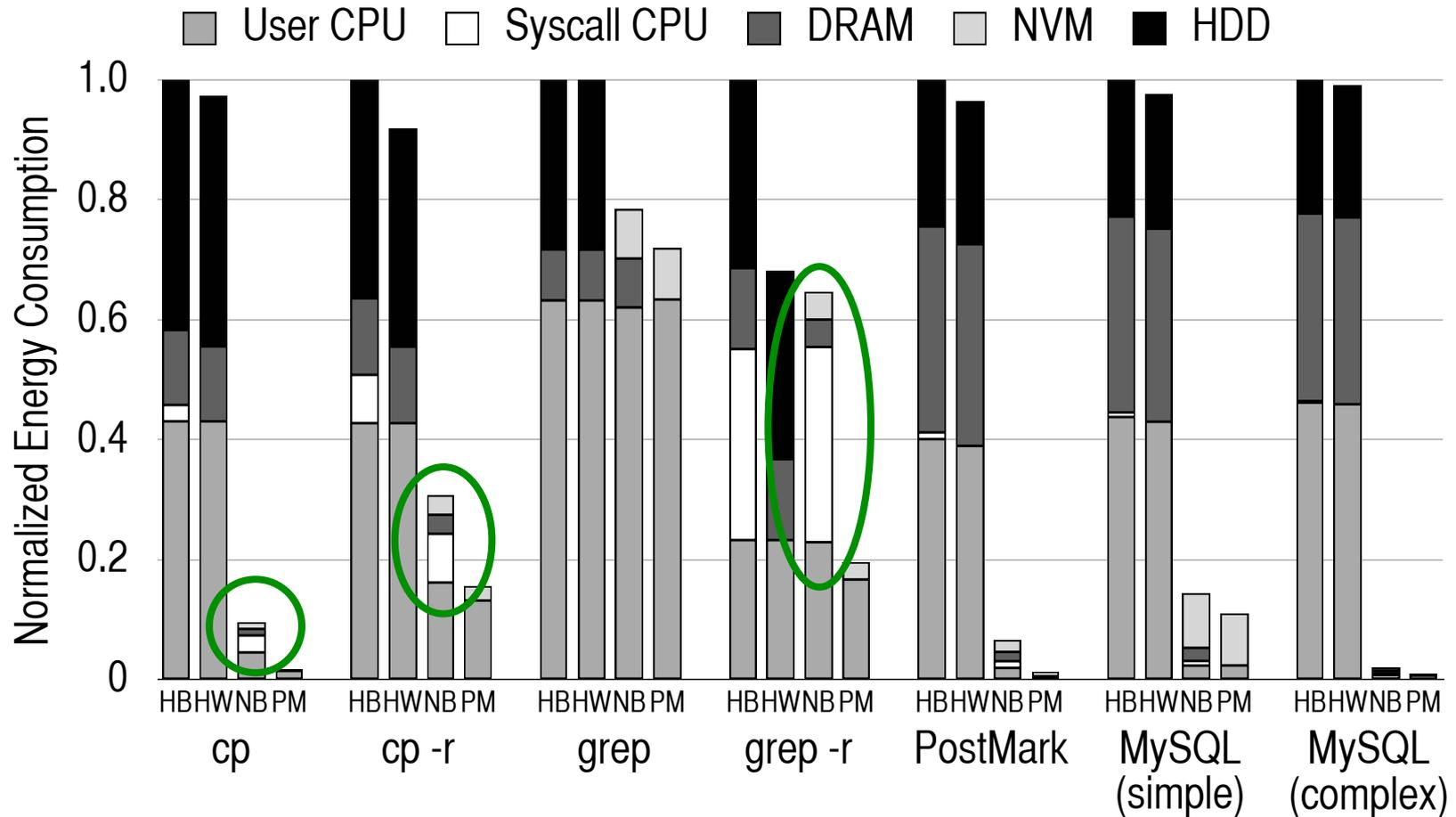


# Energy Results: HDD to NVM



Between HDD-based and NVM-based systems, lower NVM energy leads to greatly reduced energy consumption

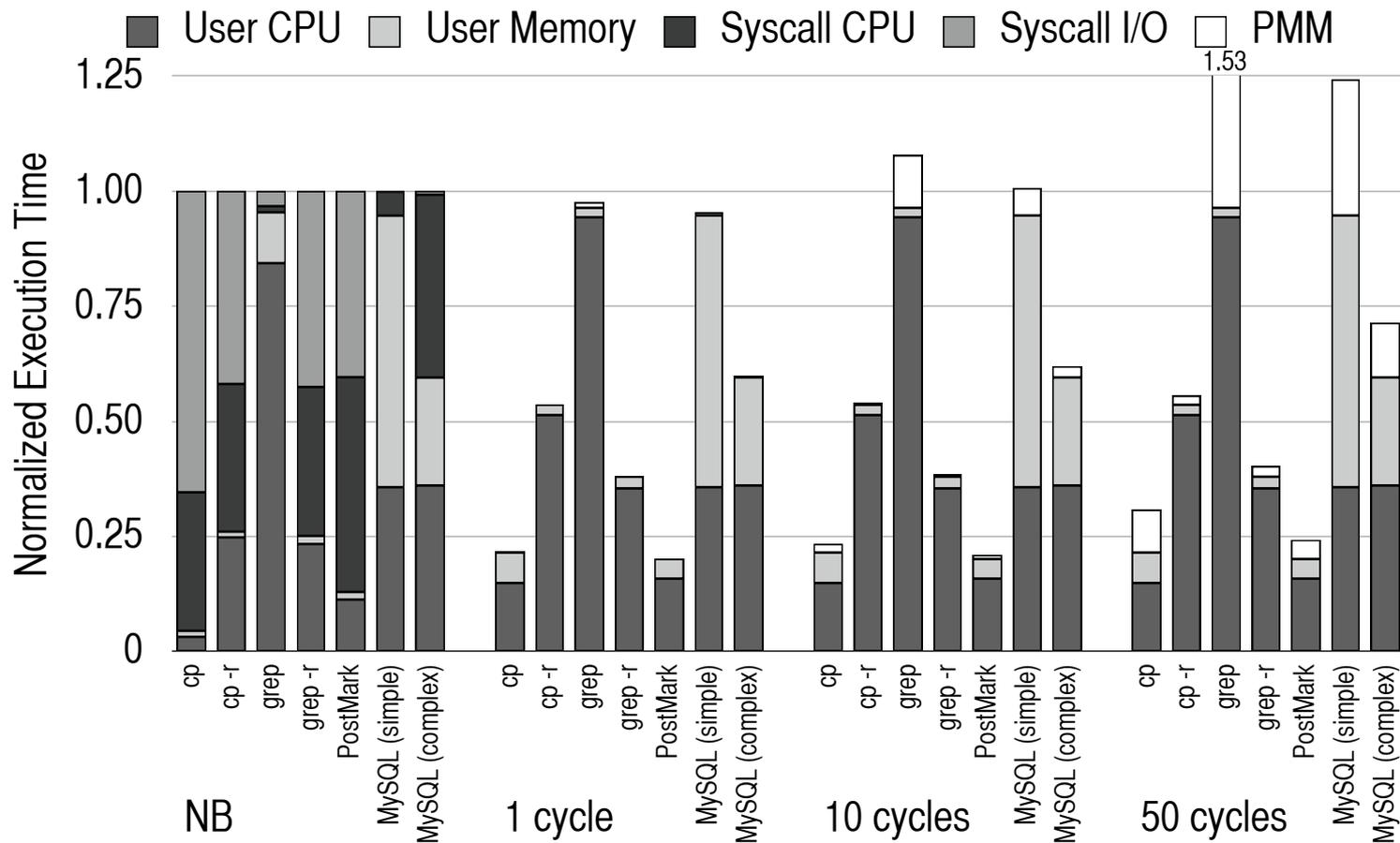
# Energy Results: NVM to PMM



Between systems with and without OS/FS code, energy improvements come from:  
1. reduced code footprint, 2. reduced data movement

**Large energy reductions with a PMM over the NVM based system**

# Scalability Analysis: Effect of PMM Latency



Even if each PMM access takes a non-overlapped 50 cycles (conservative), PMM still provides an overall improvement compared to the NVM baseline

**Future research should target keeping PMM latencies in check**

# Outline

---

- Background: Storage and Memory Models
- Motivation: Eliminating Operating/File System Bottlenecks
- Our Proposal: Hardware/Software Coordinated Management of Storage and Memory
  - Opportunities and Benefits
- Evaluation Methodology
- Evaluation Results
- **Related Work**
- New Questions and Challenges
- Conclusions

# Related Work

---

- We provide a comprehensive overview of past work related to single-level stores and persistent memory techniques
  1. Integrating file systems with persistent memory
    - ❑ Need optimized hardware to fully take advantage of new technologies
  2. Programming language support for persistent objects
    - ❑ Incurs the added latency of indirect data access through software
  3. Load/store interfaces to persistent storage
    - ❑ Lack efficient and fast hardware support for address translation, efficient file indexing, fast reliability and protection guarantees
  4. Analysis of OS overheads with Flash devices
    - ❑ Our study corroborates findings in this area and shows even larger consequences for systems with emerging NVM devices
- The goal of our work is to provide cheap and fast hardware support for memories to enable high energy efficiency and performance

# Outline

---

- Background: Storage and Memory Models
- Motivation: Eliminating Operating/File System Bottlenecks
- Our Proposal: Hardware/Software Coordinated Management of Storage and Memory
  - Opportunities and Benefits
- Evaluation Methodology
- Evaluation Results
- Related Work
- **New Questions and Challenges**
- Conclusions

# New Questions and Challenges

---

- We identify and discuss several open research questions
  - Q1. How to tailor applications for systems with persistent memory?
  - Q2. How can hardware and software cooperate to support a scalable, persistent single-level address space?
  - Q3. How to provide efficient backward compatibility (for two-level stores) on persistent memory systems?
  - Q4. How to mitigate potential hardware performance and energy overheads?

# Outline

---

- Background: Storage and Memory Models
- Motivation: Eliminating Operating/File System Bottlenecks
- Our Proposal: Hardware/Software Coordinated Management of Storage and Memory
  - Opportunities and Benefits
- Evaluation Methodology
- Evaluation Results
- Related Work
- New Questions and Challenges
- **Conclusions**

# Summary and Conclusions

---

- Traditional two-level storage model is inefficient in terms of performance and energy
  - Due to OS/FS code and buffering needed to manage two models
  - Especially so in future devices with NVM technologies, as we show
- New non-volatile memory based persistent memory designs that use a single-level storage model to unify memory and storage can alleviate this problem
- We quantified the performance and energy benefits of such a single-level persistent memory/storage design
  - Showed significant benefits from reduced code footprint, data movement, and system software overhead on a variety of workloads
- Such a design requires more research to answer the questions we have posed and enable efficient persistent memory managers
  - can lead to a fundamentally more efficient storage system

# A Case for Efficient Hardware/Software Cooperative Management of Storage and Memory

Justin Meza<sup>\*</sup>, Yixin Luo<sup>\*</sup>, Samira Khan<sup>\*†</sup>, Jishen Zhao<sup>§</sup>,  
Yuan Xie<sup>§‡</sup>, and **Onur Mutlu<sup>\*</sup>**

<sup>\*</sup>Carnegie Mellon University

<sup>§</sup>Pennsylvania State University

<sup>†</sup>Intel Labs    <sup>‡</sup>AMD Research

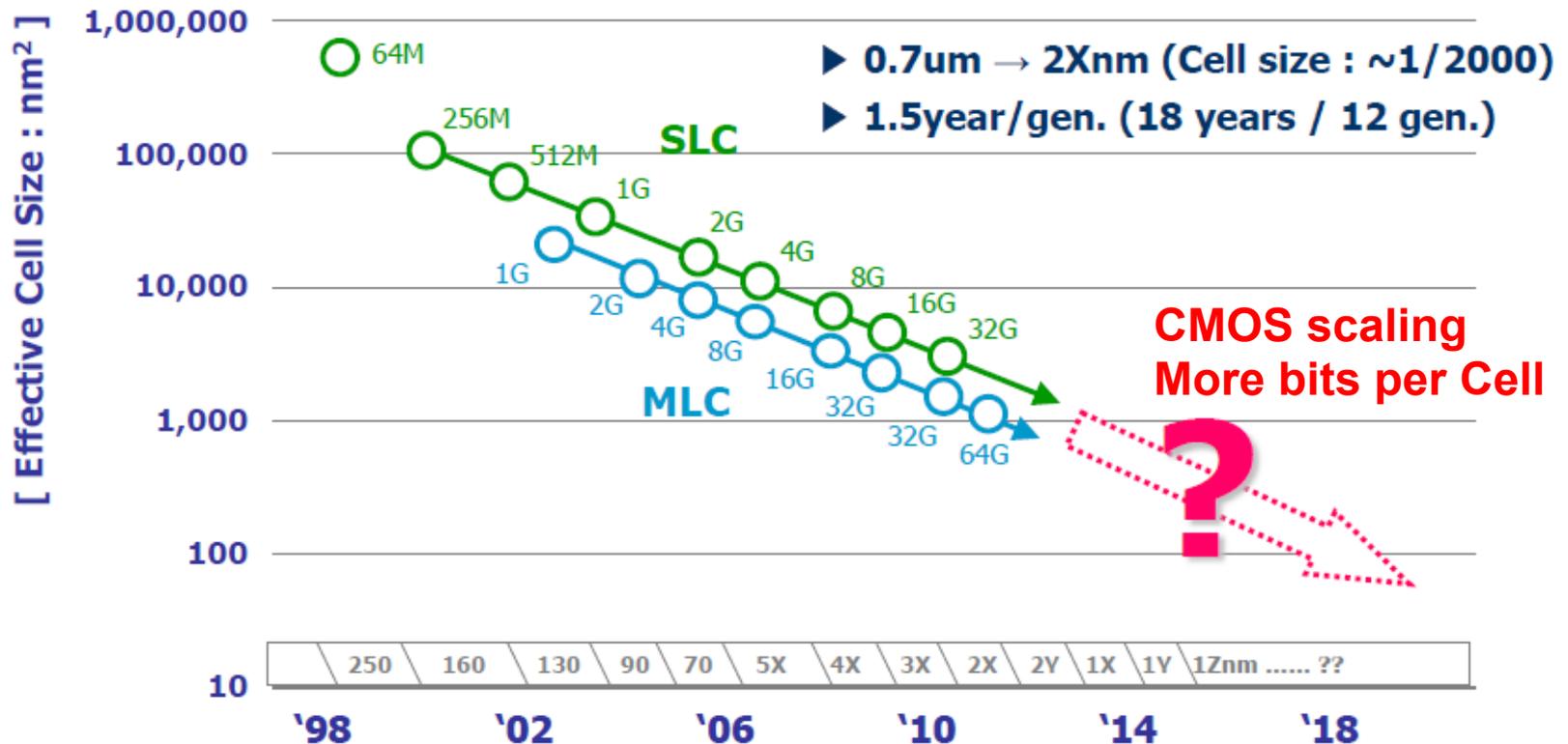
# Flash Memory Scaling

# Readings in Flash Memory

---

- Yu Cai, Gulay Yalcin, Onur Mutlu, Erich F. Haratsch, Adrian Cristal, Osman Unsal, and Ken Mai, **"Error Analysis and Retention-Aware Error Management for NAND Flash Memory"** *Intel Technology Journal (ITJ) Special Issue on Memory Resiliency*, Vol. 17, No. 1, May 2013.
- Yu Cai, Erich F. Haratsch, Onur Mutlu, and Ken Mai, **"Threshold Voltage Distribution in MLC NAND Flash Memory: Characterization, Analysis and Modeling"** *Proceedings of the Design, Automation, and Test in Europe Conference (DATE)*, Grenoble, France, March 2013. [Slides \(ppt\)](#)
- Yu Cai, Gulay Yalcin, Onur Mutlu, Erich F. Haratsch, Adrian Cristal, Osman Unsal, and Ken Mai, **"Flash Correct-and-Refresh: Retention-Aware Error Management for Increased Flash Memory Lifetime"** *Proceedings of the 30th IEEE International Conference on Computer Design (ICCD)*, Montreal, Quebec, Canada, September 2012. [Slides \(ppt\)](#) [\(pdf\)](#)
- Yu Cai, Erich F. Haratsch, Onur Mutlu, and Ken Mai, **"Error Patterns in MLC NAND Flash Memory: Measurement, Characterization, and Analysis"** *Proceedings of the Design, Automation, and Test in Europe Conference (DATE)*, Dresden, Germany, March 2012. [Slides \(ppt\)](#)

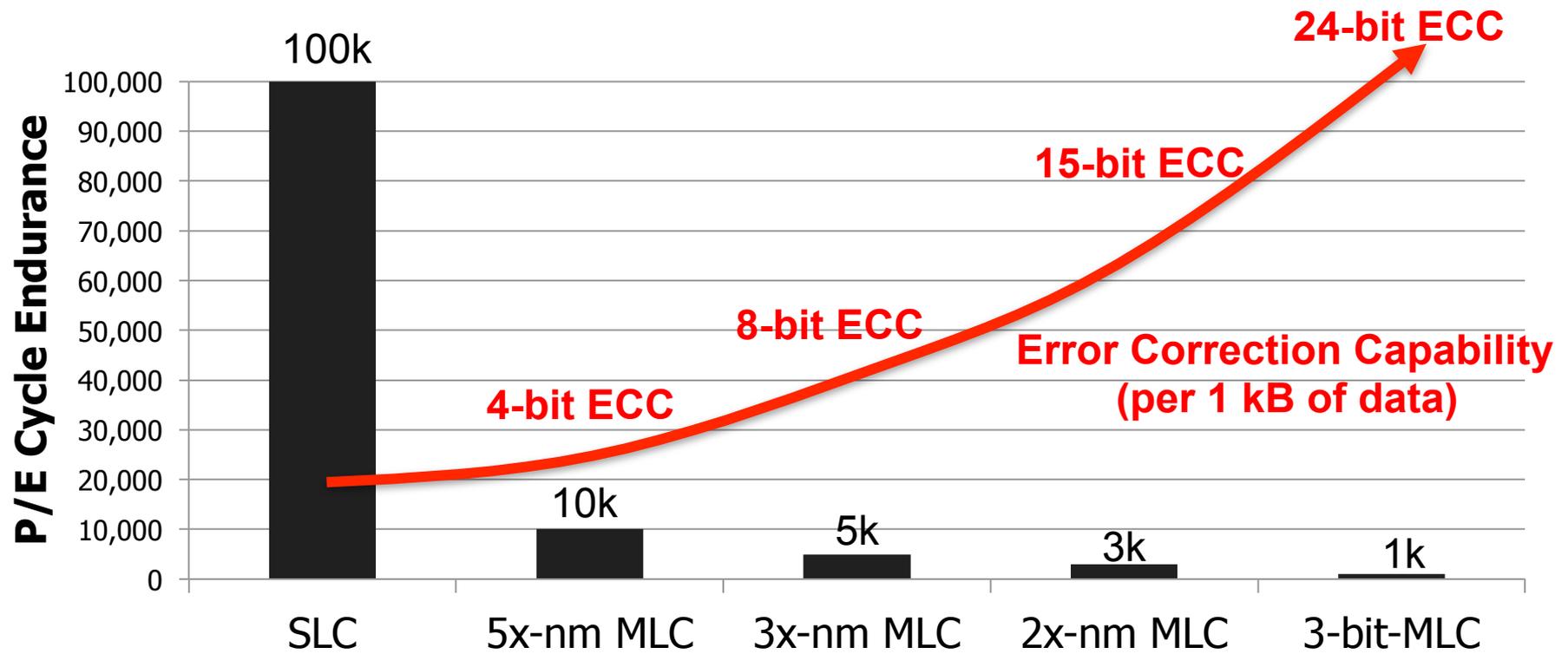
# Evolution of NAND Flash Memory



Seaung Suk Lee, "Emerging Challenges in NAND Flash Technology", Flash Summit 2011 (Hynix)

- Flash memory widening its range of applications
  - Portable consumer devices, laptop PCs and enterprise servers

# Decreasing Endurance with Flash Scaling



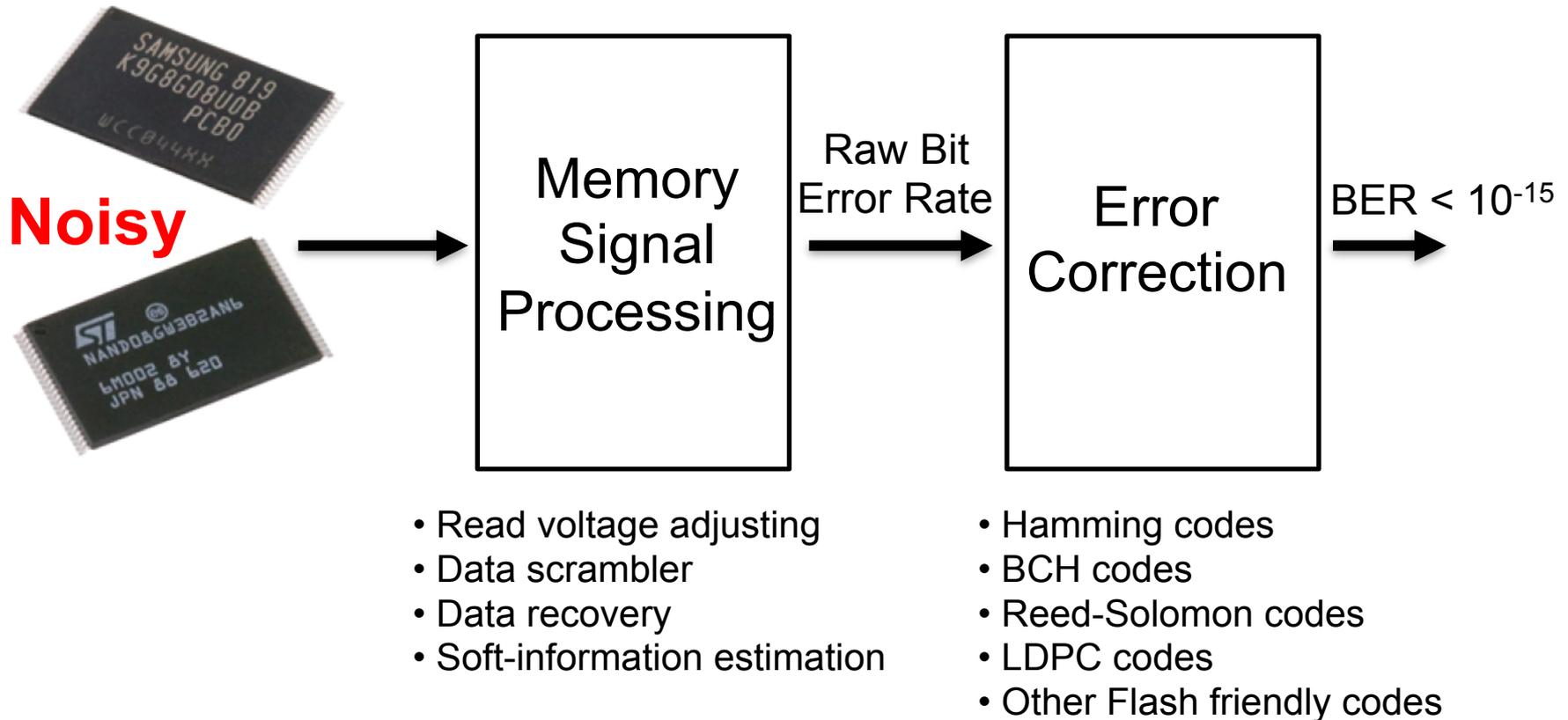
Ariel Maislos, "A New Era in Embedded Flash Memory", Flash Summit 2011 (Anobit)

- Endurance of flash memory decreasing with scaling and multi-level cells
- Error correction capability required to guarantee storage-class reliability (UBER <  $10^{-15}$ ) is increasing exponentially to reach less endurance

UBER: Uncorrectable bit error rate. Fraction of erroneous bits after error correction.

# Future NAND Flash Storage Architecture

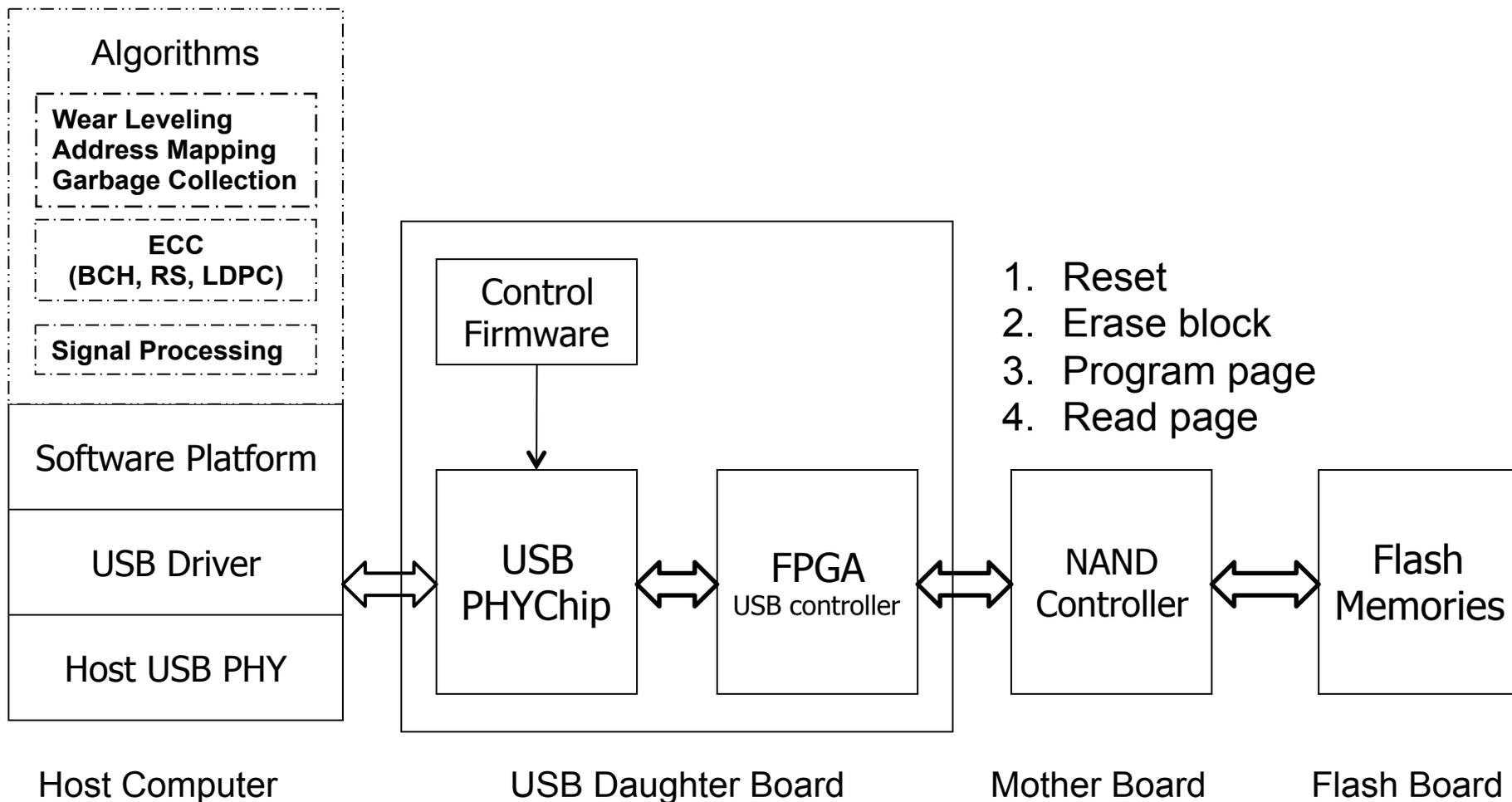
---



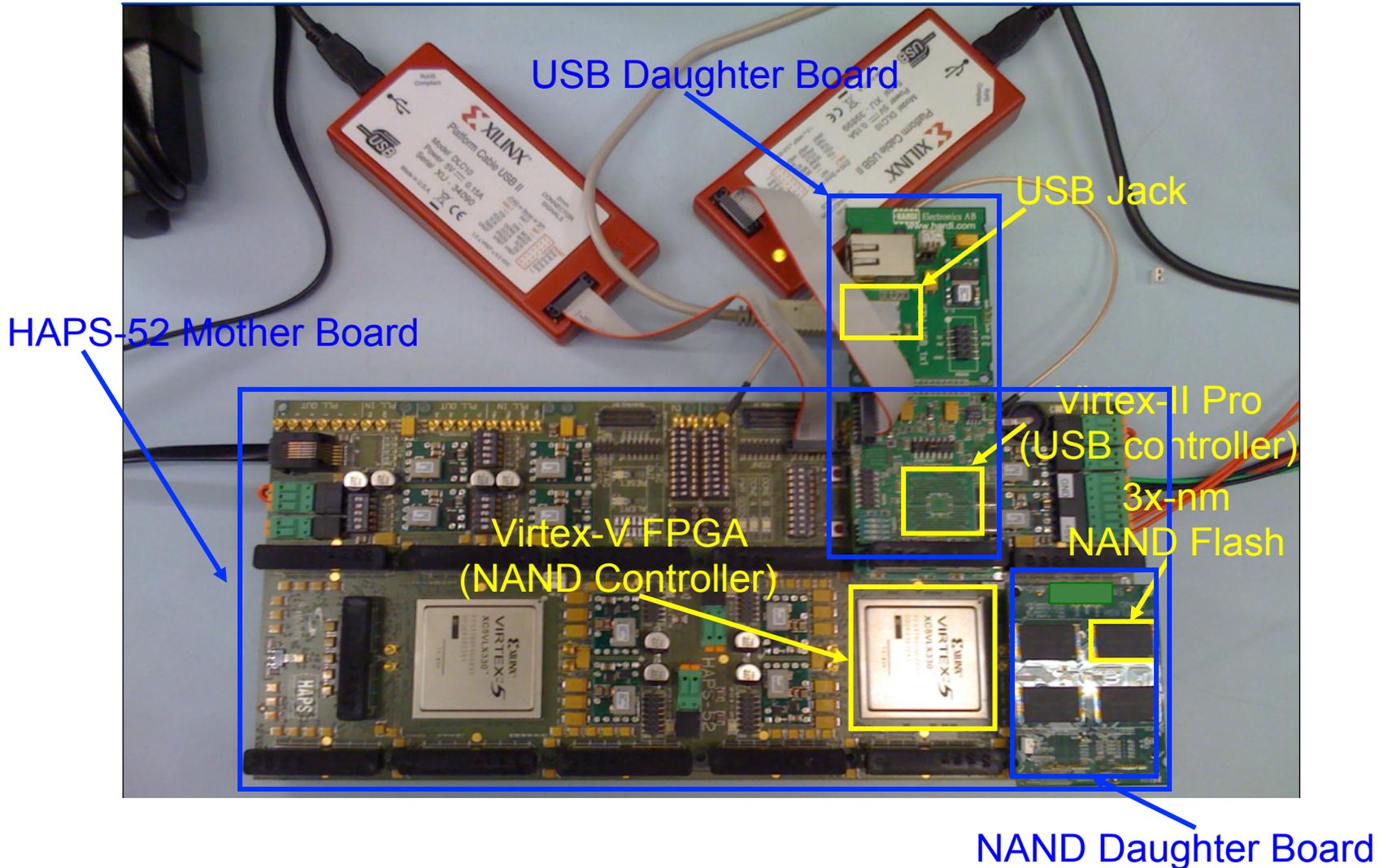
**Need to understand NAND flash error patterns**

---

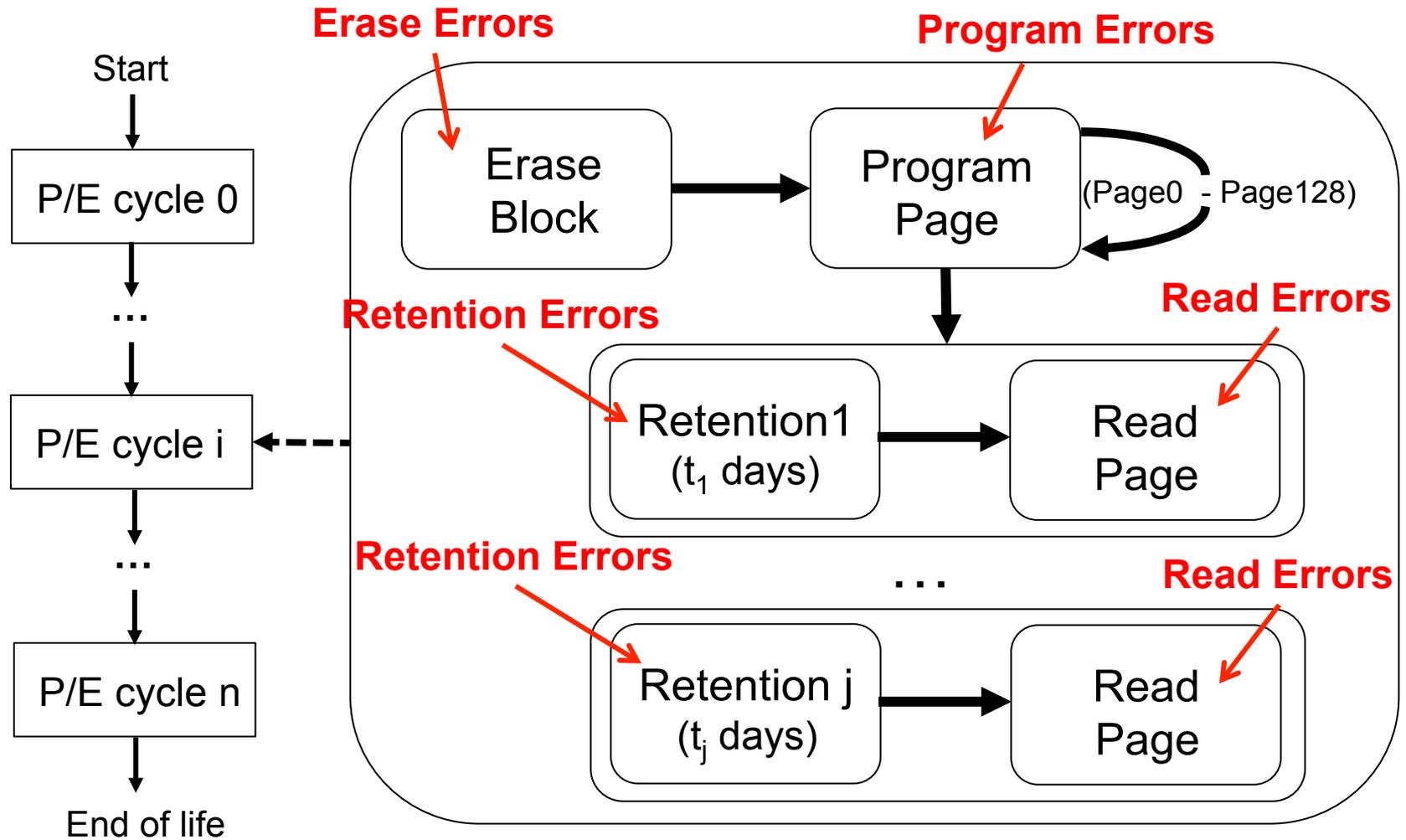
# Test System Infrastructure



# NAND Flash Testing Platform



# NAND Flash Usage and Error Model

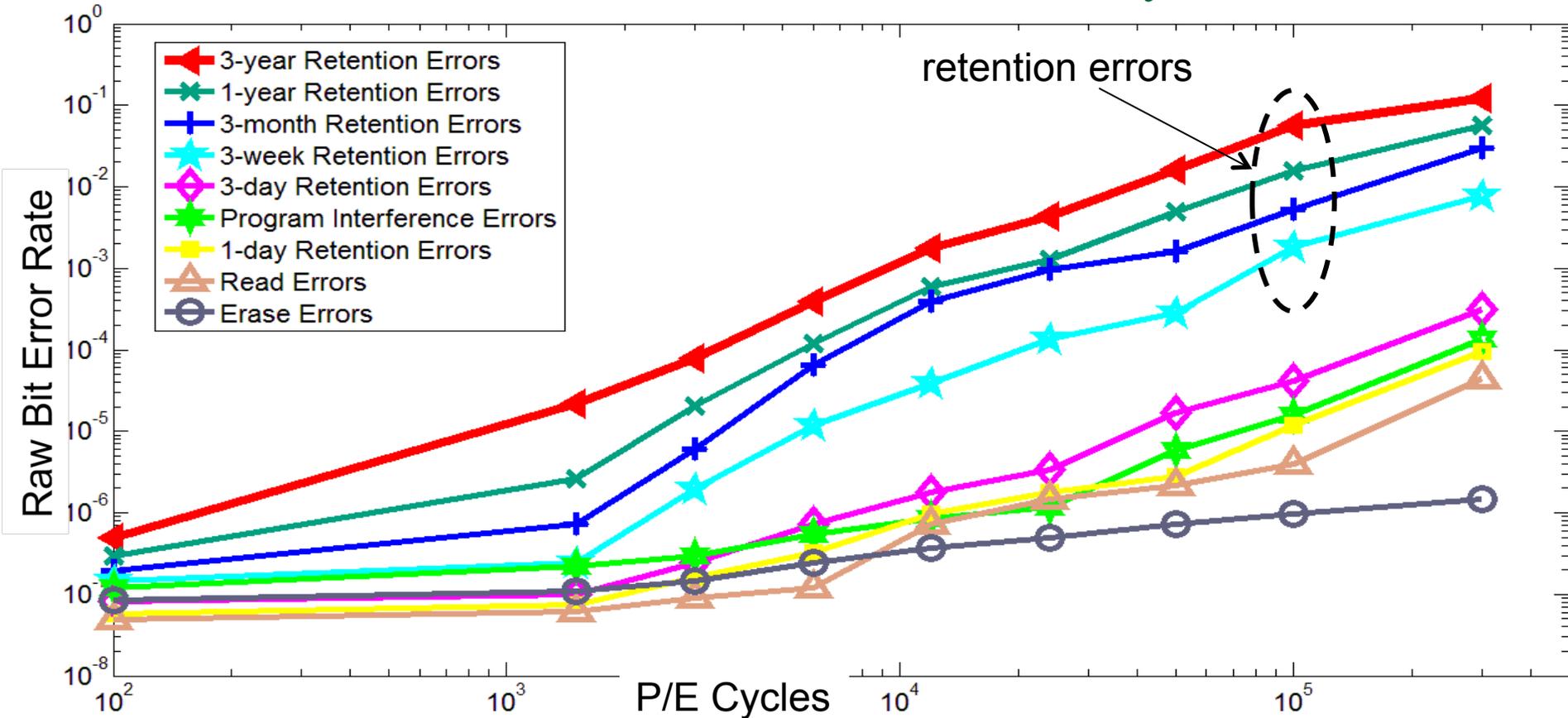


# Error Types and Testing Methodology

---

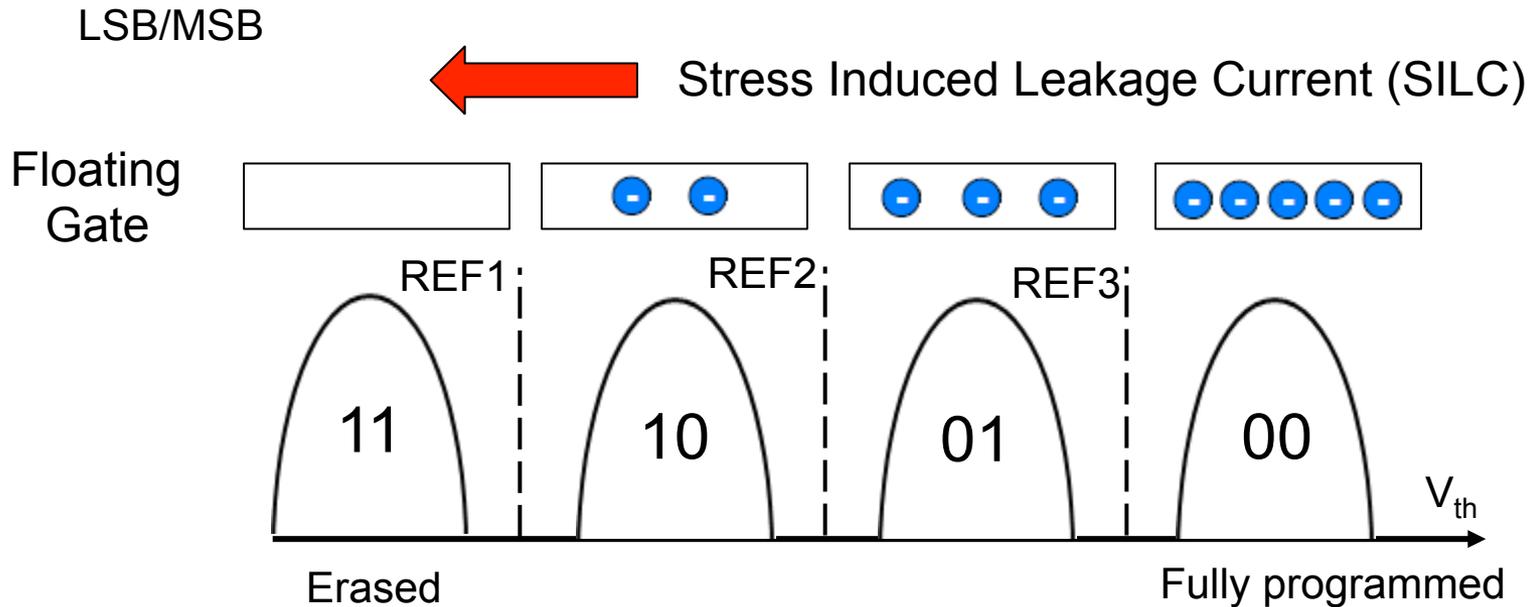
- Erase errors
  - Count the number of cells that fail to be erased to “11” state
  
- Program interference errors
  - Compare the data immediately after page programming and the data after the whole block being programmed
  
- Read errors
  - Continuously read a given block and compare the data between consecutive read sequences
  
- Retention errors
  - Compare the data read after an amount of time to data written
    - Characterize short term retention errors under room temperature
    - Characterize long term retention errors by baking in the oven under 125°C

# Observations: Flash Error Analysis



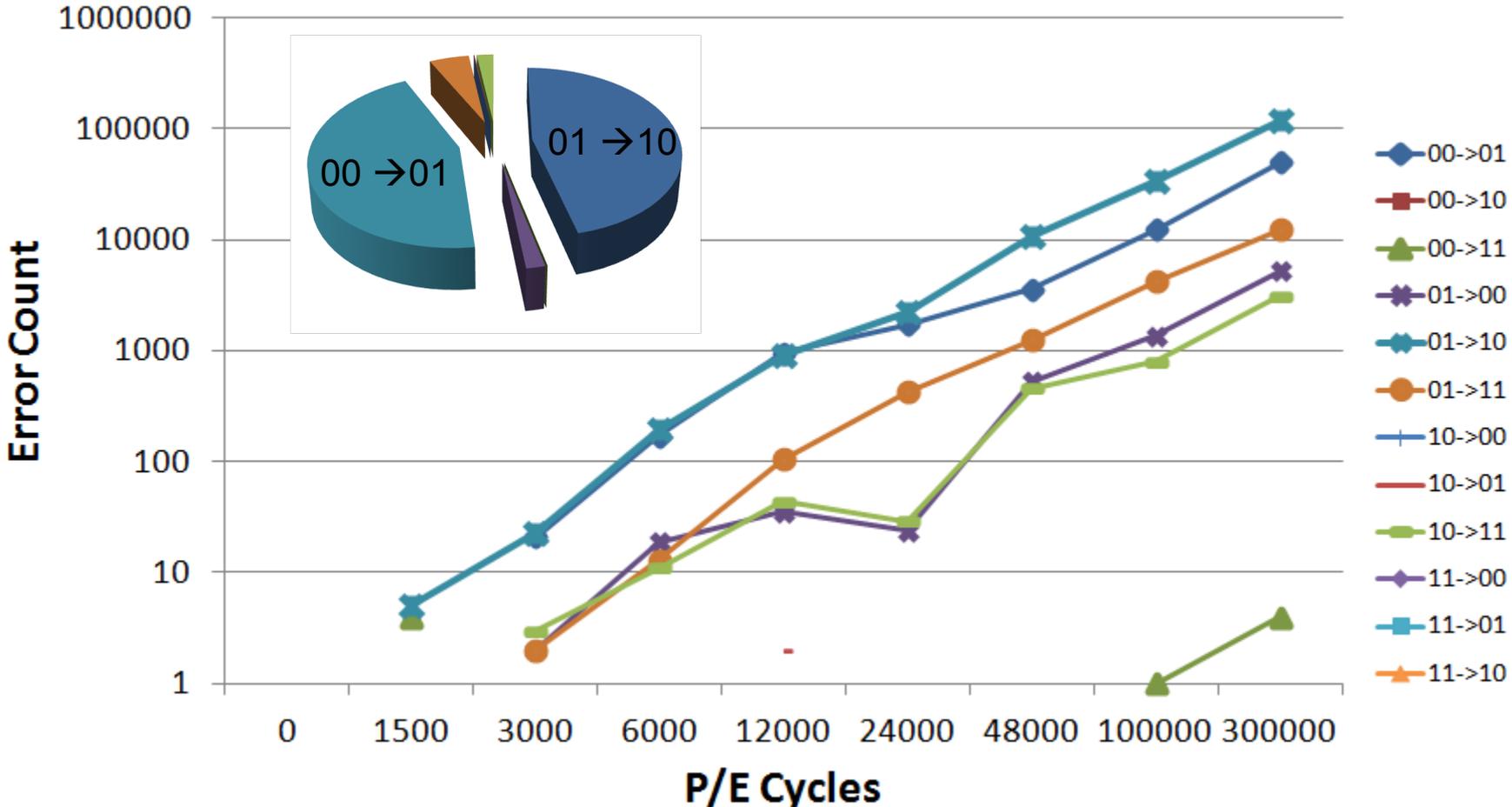
- Raw bit error rate increases exponentially with P/E cycles
- Retention errors are dominant (>99% for 1-year ret. time)
- Retention errors increase with retention time requirement

# Retention Error Mechanism



- Electron loss from the floating gate causes retention errors
  - ❑ Cells with more programmed electrons suffer more from retention errors
  - ❑ Threshold voltage is more likely to shift by one window than by multiple

# Retention Error Value Dependency



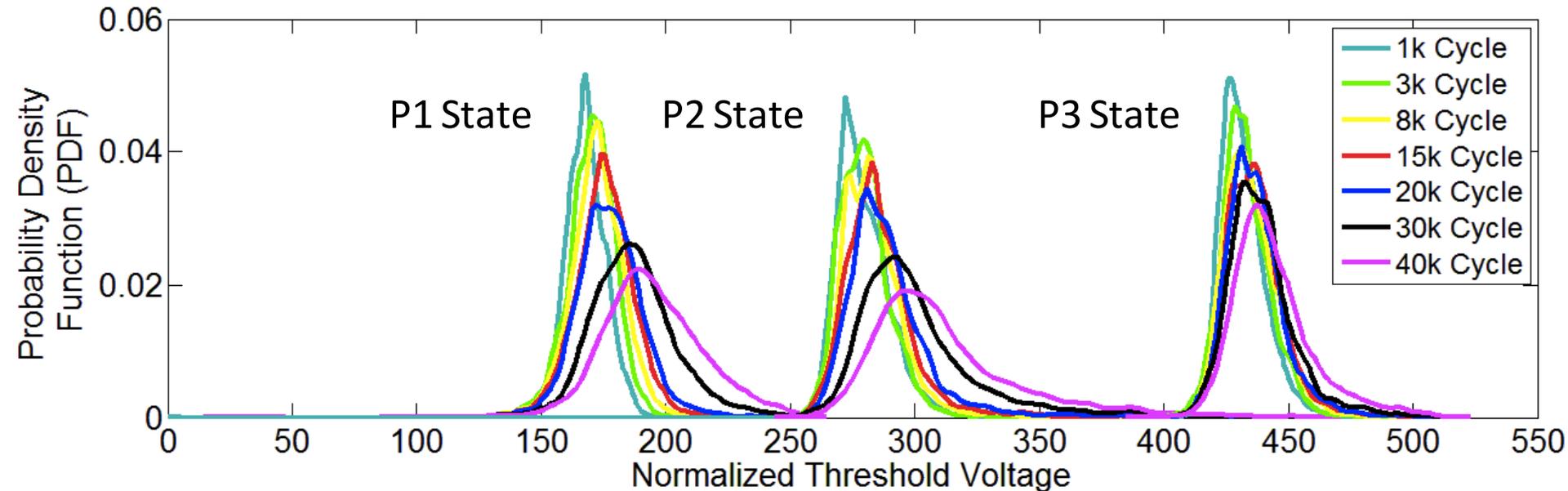
- Cells with more programmed electrons tend to suffer more from retention noise (i.e. 00 and 01)

# More Details on Flash Error Analysis

---

- Yu Cai, Erich F. Haratsch, Onur Mutlu, and Ken Mai, **"Error Patterns in MLC NAND Flash Memory: Measurement, Characterization, and Analysis"** *Proceedings of the Design, Automation, and Test in Europe Conference (DATE)*, Dresden, Germany, March 2012. Slides (ppt)

# Threshold Voltage Distribution Shifts



As P/E cycles increase ...

- Distribution shifts to the right
- Distribution becomes wider

# More Detail

---

- Yu Cai, Erich F. Haratsch, Onur Mutlu, and Ken Mai, **"Threshold Voltage Distribution in MLC NAND Flash Memory: Characterization, Analysis and Modeling"** *Proceedings of the Design, Automation, and Test in Europe Conference (DATE)*, Grenoble, France, March 2013. Slides (ppt)

# Flash Correct-and-Refresh

## Retention-Aware Error Management for Increased Flash Memory Lifetime

Yu Cai<sup>1</sup> Gulay Yalcin<sup>2</sup> Onur Mutlu<sup>1</sup> Erich F. Haratsch<sup>3</sup>  
Adrian Cristal<sup>2</sup> Osman S. Unsal<sup>2</sup> Ken Mai<sup>1</sup>

<sup>1</sup> Carnegie Mellon University

<sup>2</sup> Barcelona Supercomputing Center

<sup>3</sup> LSI Corporation



**SAFARI**

**Carnegie Mellon**

# Executive Summary

---

- NAND flash memory has low endurance: a flash cell dies after 3k P/E cycles vs. 50k desired → Major scaling challenge for flash memory
  - Flash error rate increases exponentially over flash lifetime
  - **Problem:** Stronger error correction codes (ECC) are ineffective and undesirable for improving flash lifetime due to
    - diminishing returns on lifetime with increased correction strength
    - prohibitively high power, area, latency overheads
  - **Our Goal:** Develop techniques to tolerate high error rates w/o strong ECC
  - **Observation:** Retention errors are the dominant errors in MLC NAND flash
    - flash cell loses charge over time; retention errors increase as cell gets worn out
  - **Solution:** Flash Correct-and-Refresh (FCR)
    - Periodically read, correct, and reprogram (in place) or remap each flash page before it accumulates more errors than can be corrected by simple ECC
    - Adapt “refresh” rate to the severity of retention errors (i.e., # of P/E cycles)
  - **Results:** FCR improves flash memory lifetime by 46X with no hardware changes and low energy overhead; outperforms strong ECCs
-

# Outline

---

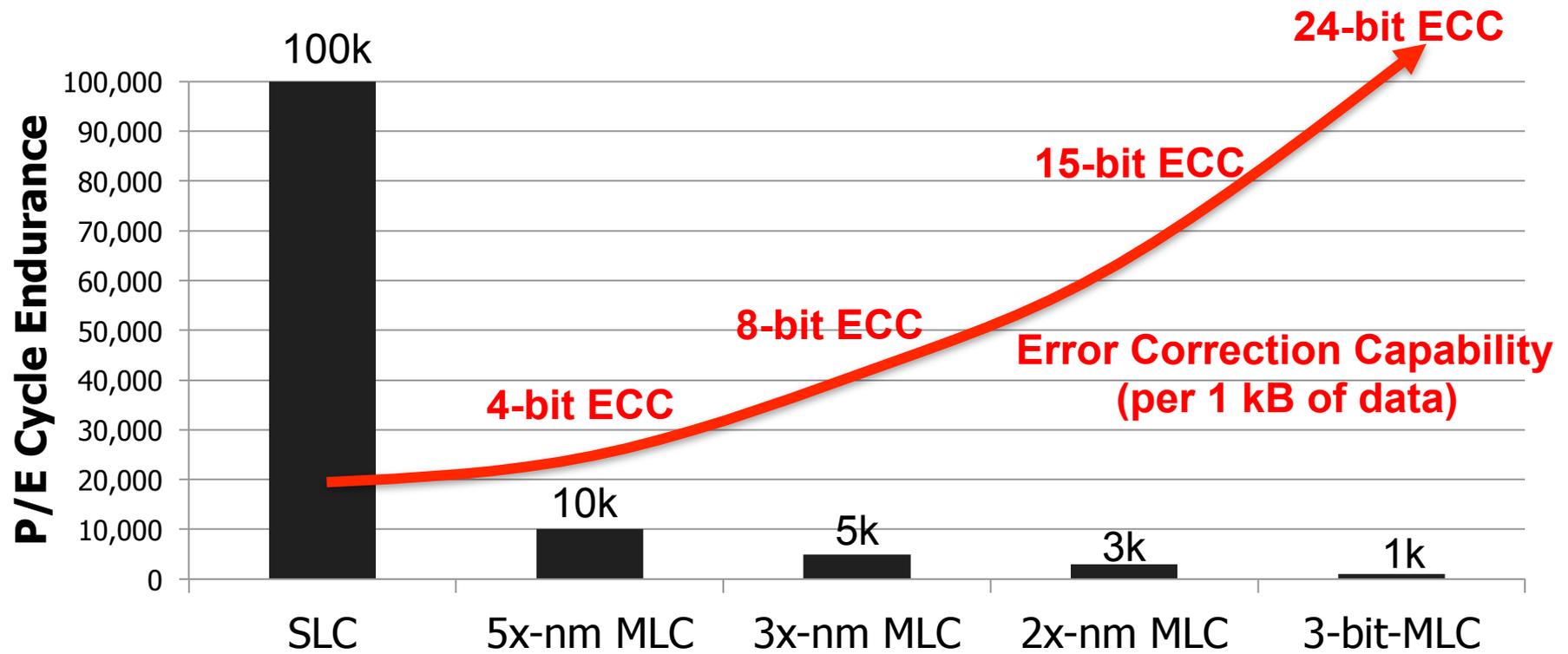
- Executive Summary
- **The Problem: Limited Flash Memory Endurance/Lifetime**
- Error and ECC Analysis for Flash Memory
- Flash Correct and Refresh Techniques (FCR)
- Evaluation
- Conclusions

# Problem: Limited Endurance of Flash Memory

---

- **NAND flash has limited endurance**
    - A cell can tolerate a small number of Program/Erase (P/E) cycles
    - 3x-nm flash with 2 bits/cell → 3K P/E cycles
  - Enterprise data storage requirements demand very high endurance
    - >50K P/E cycles (10 full disk writes per day for 3-5 years)
  - **Continued process scaling and more bits per cell will reduce flash endurance**
  - One potential solution: stronger error correction codes (ECC)
    - **Stronger ECC not effective enough and inefficient**
-

# Decreasing Endurance with Flash Scaling



Ariel Maislos, "A New Era in Embedded Flash Memory", Flash Summit 2011 (Anobit)

- Endurance of flash memory decreasing with scaling and multi-level cells
- Error correction capability required to guarantee storage-class reliability (UBER <  $10^{-15}$ ) is increasing exponentially to reach less endurance

UBER: Uncorrectable bit error rate. Fraction of erroneous bits after error correction.

# The Problem with Stronger Error Correction

---

- Stronger ECC detects and corrects more raw bit errors → increases P/E cycles endured
- Two shortcomings of stronger ECC:
  1. High implementation complexity
    - Power and area overheads increase super-linearly, but correction capability increases sub-linearly with ECC strength
  2. Diminishing returns on flash lifetime improvement
    - Raw bit error rate increases exponentially with P/E cycles, but correction capability increases sub-linearly with ECC strength

# Outline

---

- Executive Summary
- The Problem: Limited Flash Memory Endurance/Lifetime
- Error and ECC Analysis for Flash Memory
- Flash Correct and Refresh Techniques (FCR)
- Evaluation
- Conclusions

# Methodology: Error and ECC Analysis

---

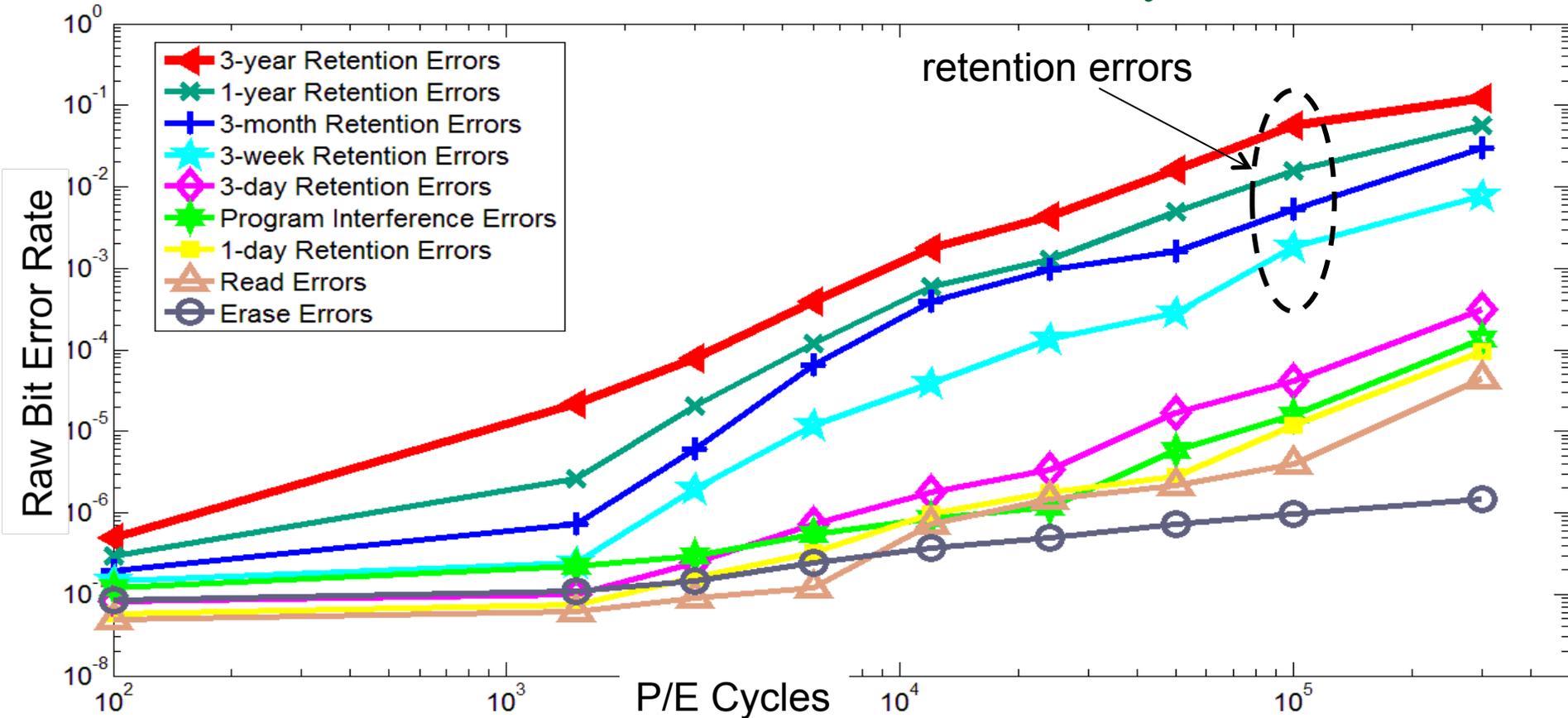
- **Characterized errors and error rates** of 3x-nm MLC NAND flash using an experimental FPGA-based flash platform
  - Cai et al., "Error Patterns in MLC NAND Flash Memory: Measurement, Characterization, and Analysis," DATE 2012.
- **Quantified Raw Bit Error Rate (RBER) at a given P/E cycle**
  - Raw Bit Error Rate: Fraction of erroneous bits without any correction
- **Quantified error correction capability** (and area and power consumption) of various BCH-code implementations
  - Identified how much RBER each code can tolerate
    - how many P/E cycles (flash lifetime) each code can sustain

# NAND Flash Error Types

---

- Four types of errors [Cai+, DATE 2012]
- Caused by **common flash operations**
  - **Read** errors
  - **Erase** errors
  - **Program** (interference) errors
- Caused by flash **cell losing charge over time**
  - **Retention** errors
    - Whether an error happens depends on required retention time
    - Especially problematic in MLC flash because voltage threshold window to determine stored value is smaller

# Observations: Flash Error Analysis



- Raw bit error rate increases exponentially with P/E cycles
- Retention errors are dominant (>99% for 1-year ret. time)
- Retention errors increase with retention time requirement

# Methodology: Error and ECC Analysis

---

- **Characterized errors and error rates** of 3x-nm MLC NAND flash using an experimental FPGA-based flash platform
  - Cai et al., "Error Patterns in MLC NAND Flash Memory: Measurement, Characterization, and Analysis," DATE 2012.
- **Quantified Raw Bit Error Rate (RBER) at a given P/E cycle**
  - Raw Bit Error Rate: Fraction of erroneous bits without any correction
- **Quantified error correction capability** (and area and power consumption) of various BCH-code implementations
  - Identified how much RBER each code can tolerate
    - how many P/E cycles (flash lifetime) each code can sustain

# ECC Strength Analysis

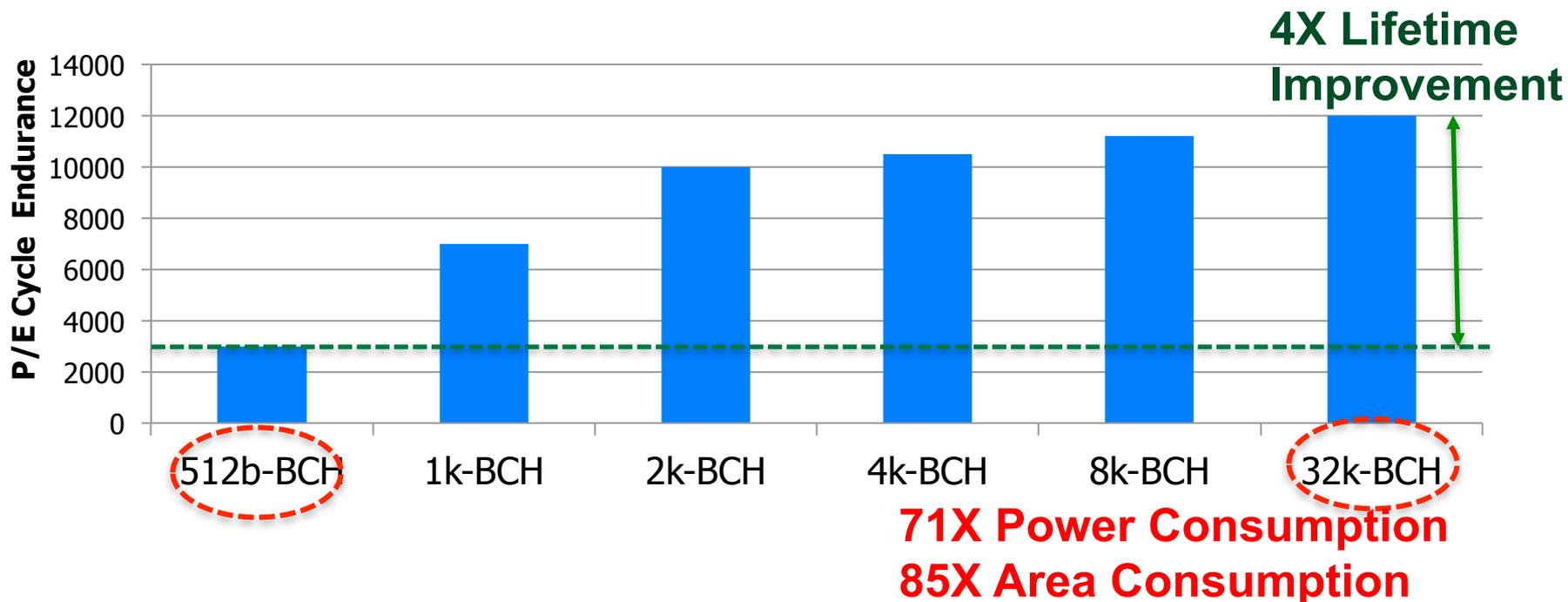
Error correction capability increases sub-linearly

Power and area overheads increase super-linearly

Code length (n)	Correctable Errors (t)	Acceptable Raw BER	Norm. Power	Norm. Area
512	7	$1.0 \times 10^{-4}$ (1x)	1	1
1024	12	$4.0 \times 10^{-4}$ (4x)	2	2.1
2048	22	$1.0 \times 10^{-3}$ (10x)	4.1	3.9
4096	40	$1.7 \times 10^{-3}$ (17x)	8.6	10.3
8192	74	$2.2 \times 10^{-3}$ (22x)	17.8	21.3
32768	259	$2.6 \times 10^{-3}$ (26x)	71	85

# Resulting Flash Lifetime with Strong ECC

- Lifetime improvement comparison of various BCH codes



Strong ECC is very inefficient at improving lifetime

# Our Goal

---

Develop new techniques  
to improve flash lifetime  
without relying on stronger ECC

# Outline

---

- Executive Summary
- The Problem: Limited Flash Memory Endurance/Lifetime
- Error and ECC Analysis for Flash Memory
- Flash Correct and Refresh Techniques (FCR)
- Evaluation
- Conclusions

# Flash Correct-and-Refresh (FCR)

---

- Key Observations:

- Retention errors are the dominant source of errors in flash memory [Cai+ DATE 2012][Tanakamaru+ ISSCC 2011]  
→ limit flash lifetime as they increase over time
- Retention errors can be corrected by “refreshing” each flash page periodically

- Key Idea:

- Periodically read each flash page,
  - Correct its errors using “weak” ECC, and
  - Either remap it to a new physical page or reprogram it in-place,
  - Before the page accumulates more errors than ECC-correctable
  - Optimization: Adapt refresh rate to endured P/E cycles
-

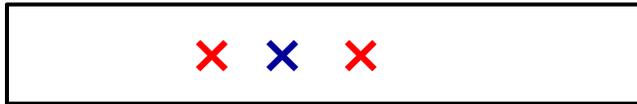
# FCR Intuition

Errors with  
No refresh

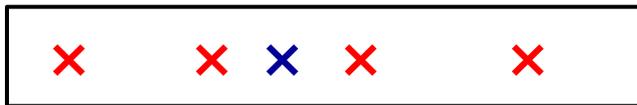
Program  
Page



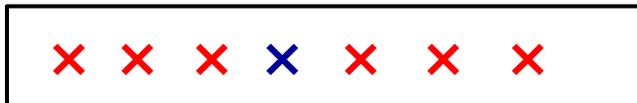
After  
time T



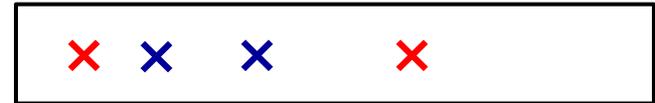
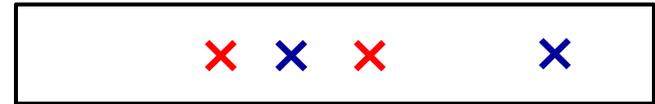
After  
time 2T



After  
time 3T



Errors with  
Periodic refresh



**X Retention Error**    **X Program Error**

# FCR: Two Key Questions

---

- How to refresh?
  - Remap a page to another one
  - Reprogram a page (in-place)
  - Hybrid of remap and reprogram
  
- When to refresh?
  - Fixed period
  - Adapt the period to retention error severity

# Outline

---

- Executive Summary
- The Problem: Limited Flash Memory Endurance/Lifetime
- Error and ECC Analysis for Flash Memory
- Flash Correct and Refresh Techniques (FCR)
  1. Remapping based FCR
  2. Hybrid Reprogramming and Remapping based FCR
  3. Adaptive-Rate FCR
- Evaluation
- Conclusions

# Outline

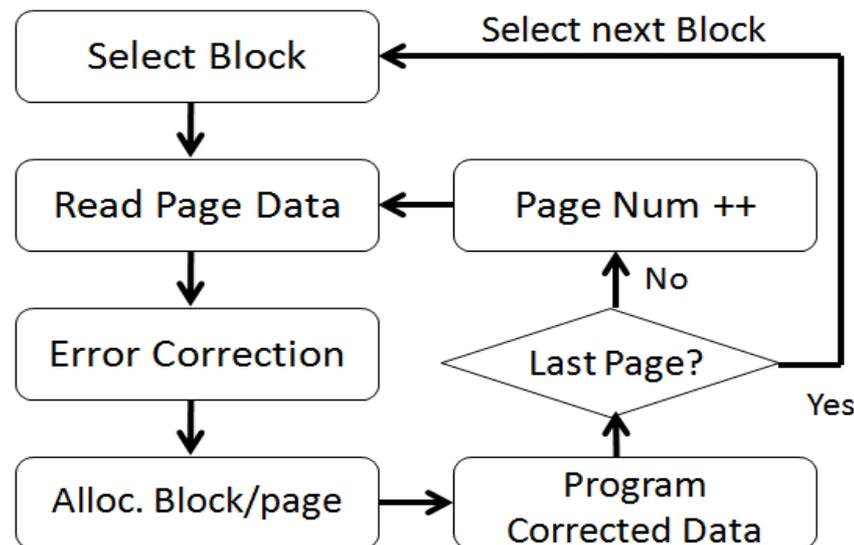
---

- Executive Summary
- The Problem: Limited Flash Memory Endurance/Lifetime
- Error and ECC Analysis for Flash Memory
- Flash Correct and Refresh Techniques (FCR)
  1. Remapping based FCR
  2. Hybrid Reprogramming and Remapping based FCR
  3. Adaptive-Rate FCR
- Evaluation
- Conclusions

# Remapping Based FCR

- Idea: Periodically remap each page to a different physical page (after correcting errors)

- Also [Pan et al., HPCA 2012]
- FTL already has support for changing logical → physical flash block/page mappings
- Deallocated block is erased by garbage collector



- Problem: Causes additional erase operations → more wearout
  - Bad for read-intensive workloads (few erases really needed)
  - Lifetime degrades for such workloads (see paper)

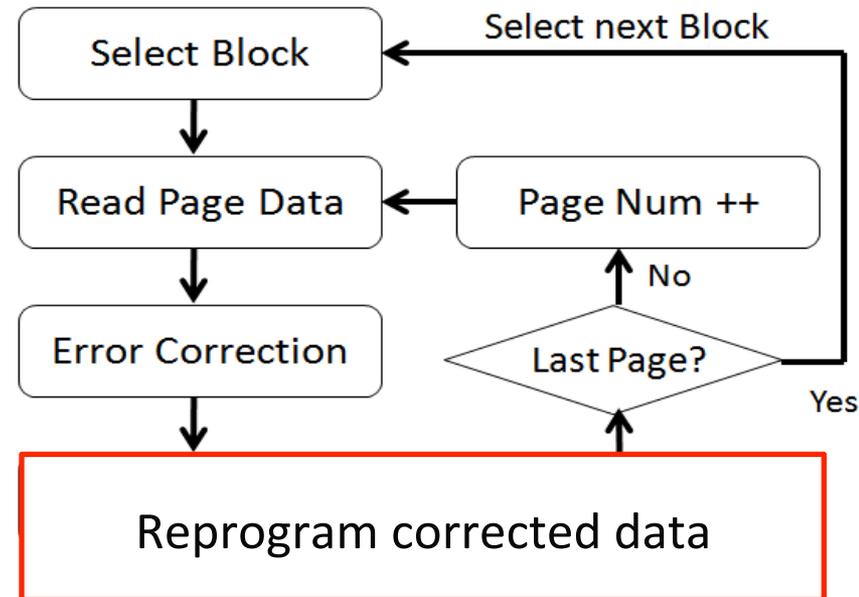
# Outline

---

- Executive Summary
- The Problem: Limited Flash Memory Endurance/Lifetime
- Error and ECC Analysis for Flash Memory
- Flash Correct and Refresh Techniques (FCR)
  1. Remapping based FCR
  2. Hybrid Reprogramming and Remapping based FCR
  3. Adaptive-Rate FCR
- Evaluation
- Conclusions

# In-Place Reprogramming Based FCR

- Idea: Periodically reprogram (in-place) each physical page (after correcting errors)
  - Flash programming techniques (ISPP) can correct retention errors in-place by recharging flash cells



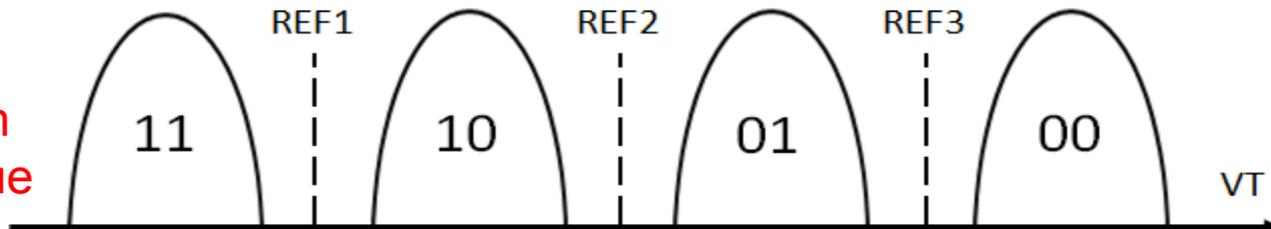
- Problem: Program errors accumulate on the same page → may not be correctable by ECC after some time

# In-Place Reprogramming of Flash Cells

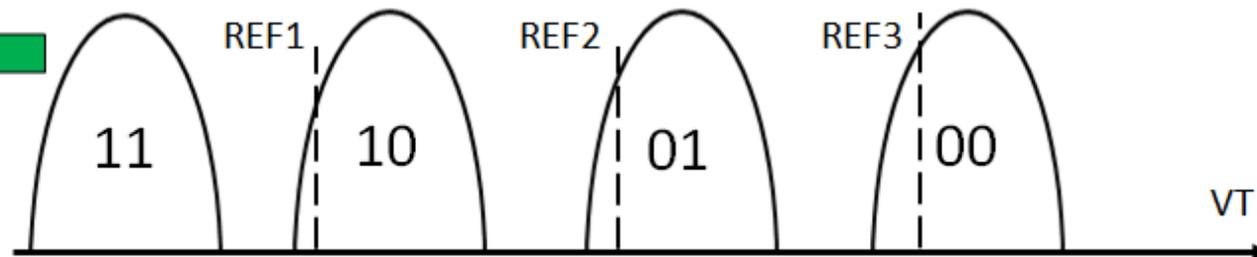
Floating Gate



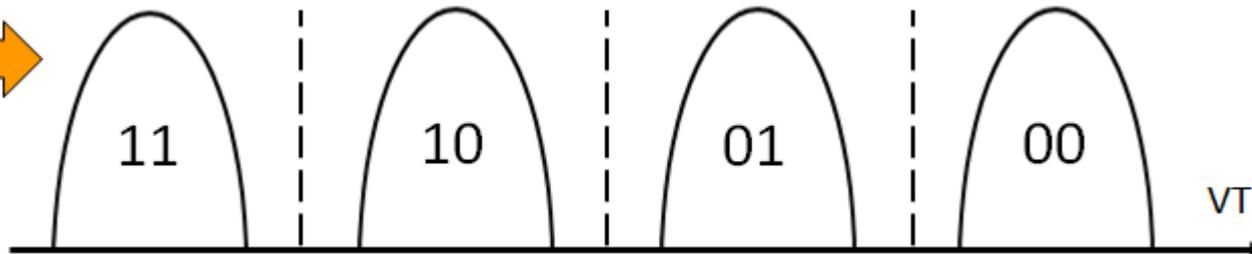
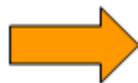
Floating Gate  
Voltage Distribution  
for each Stored Value



Retention errors are  
caused by cell voltage  
shifting to the left



ISPP moves cell  
voltage to the right;  
fixes retention errors



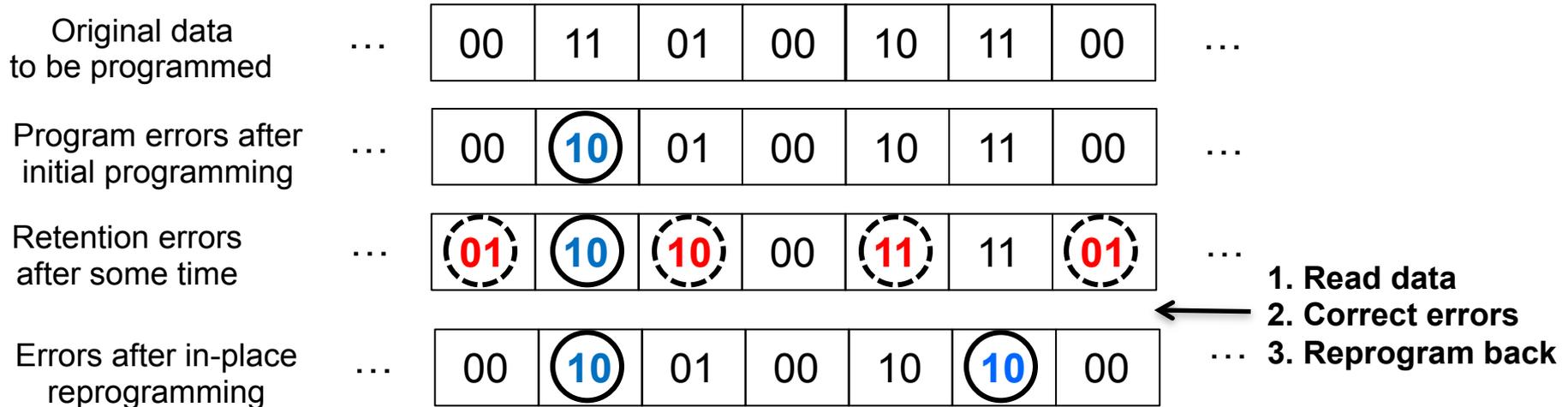
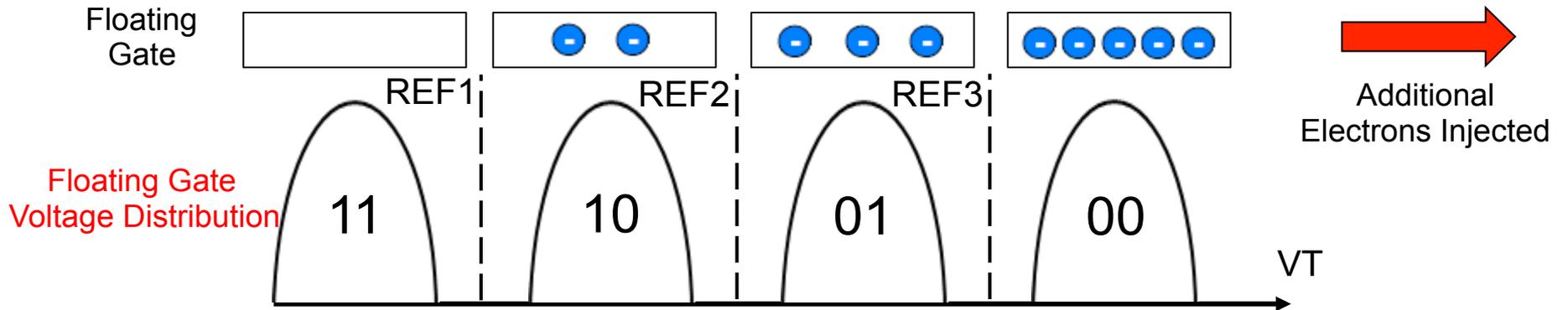
- Pro: No remapping needed → no additional erase operations
- Con: Increases the occurrence of program errors

# Program Errors in Flash Memory

---

- When a cell is being programmed, **voltage level of a neighboring cell changes** (unintentionally) due to parasitic capacitance coupling
  - **can change the data value stored**
- Also called program interference error
- Program interference causes neighboring cell voltage to shift to the right

# Problem with In-Place Reprogramming



**Problem: Program errors can accumulate over time**

# Hybrid Reprogramming/Remapping Based FCR

---

- Idea:
  - Monitor the count of right-shift errors (after error correction)
  - If count < threshold, in-place reprogram the page
  - Else, remap the page to a new page
- Observation:
  - Program errors much less frequent than retention errors → Remapping happens only infrequently
- Benefit:
  - Hybrid FCR greatly reduces erase operations due to remapping

# Outline

---

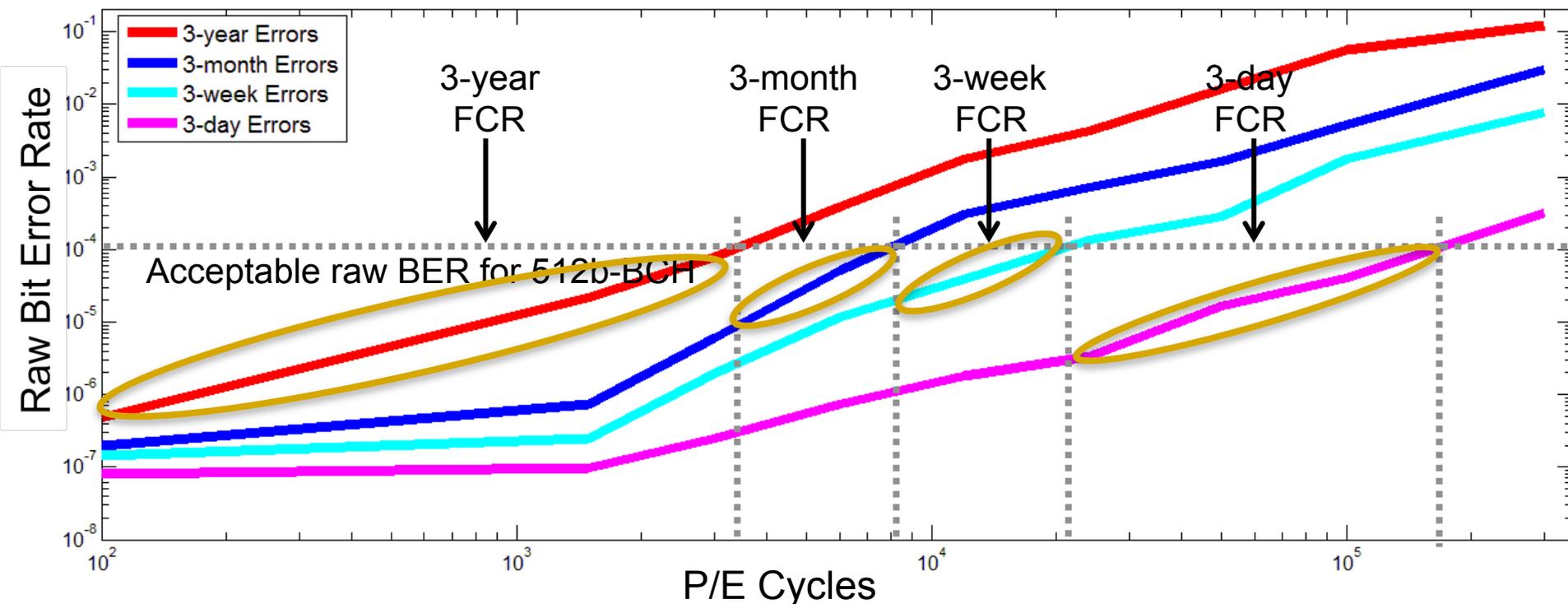
- Executive Summary
- The Problem: Limited Flash Memory Endurance/Lifetime
- Error and ECC Analysis for Flash Memory
- Flash Correct and Refresh Techniques (FCR)
  1. Remapping based FCR
  2. Hybrid Reprogramming and Remapping based FCR
  3. Adaptive-Rate FCR
- Evaluation
- Conclusions

# Adaptive-Rate FCR

---

- Observation:
  - Retention error rate strongly depends on the P/E cycles a flash page endured so far
  - No need to refresh frequently (at all) early in flash lifetime
- Idea:
  - Adapt the refresh rate to the P/E cycles endured by each page
  - Increase refresh rate gradually with increasing P/E cycles
- Benefits:
  - Reduces overhead of refresh operations
  - Can use existing FTL mechanisms that keep track of P/E cycles

# Adaptive-Rate FCR (Example)



Select refresh frequency such that error rate is below acceptable rate

# Outline

---

- Executive Summary
- The Problem: Limited Flash Memory Endurance/Lifetime
- Error and ECC Analysis for Flash Memory
- Flash Correct and Refresh Techniques (FCR)
  1. Remapping based FCR
  2. Hybrid Reprogramming and Remapping based FCR
  3. Adaptive-Rate FCR
- Evaluation
- Conclusions

# FCR: Other Considerations

---

- Implementation cost
  - No hardware changes
  - FTL software/firmware needs modification
- Response time impact
  - FCR not as frequent as DRAM refresh; low impact
- Adaptation to variations in retention error rate
  - Adapt refresh rate based on, e.g., temperature [Liu+ ISCA 2012]
- FCR requires power
  - Enterprise storage systems typically powered on

# Outline

---

- Executive Summary
- The Problem: Limited Flash Memory Endurance/Lifetime
- Error and ECC Analysis for Flash Memory
- Flash Correct and Refresh Techniques (FCR)
- Evaluation
- Conclusions

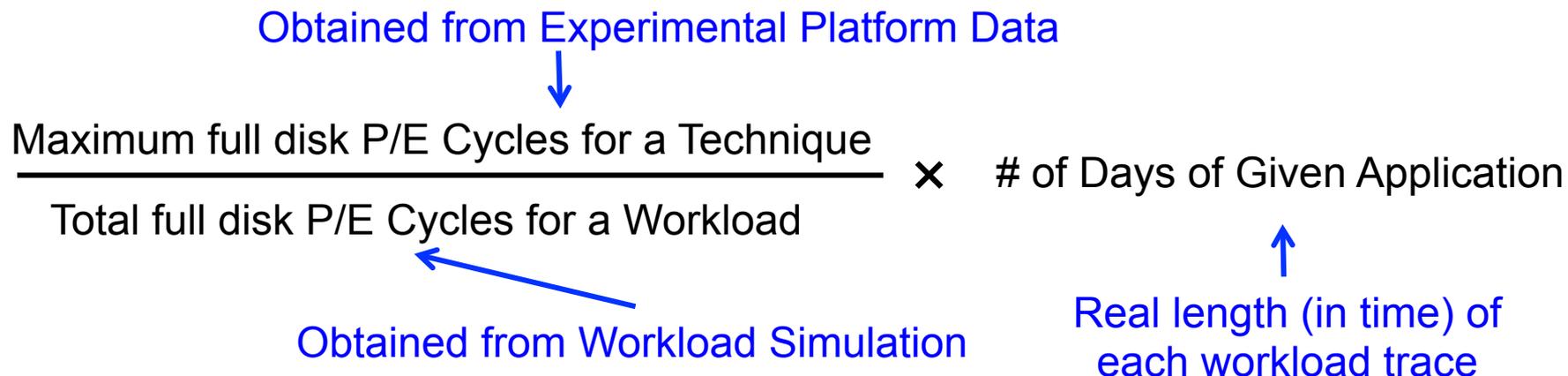
# Evaluation Methodology

---

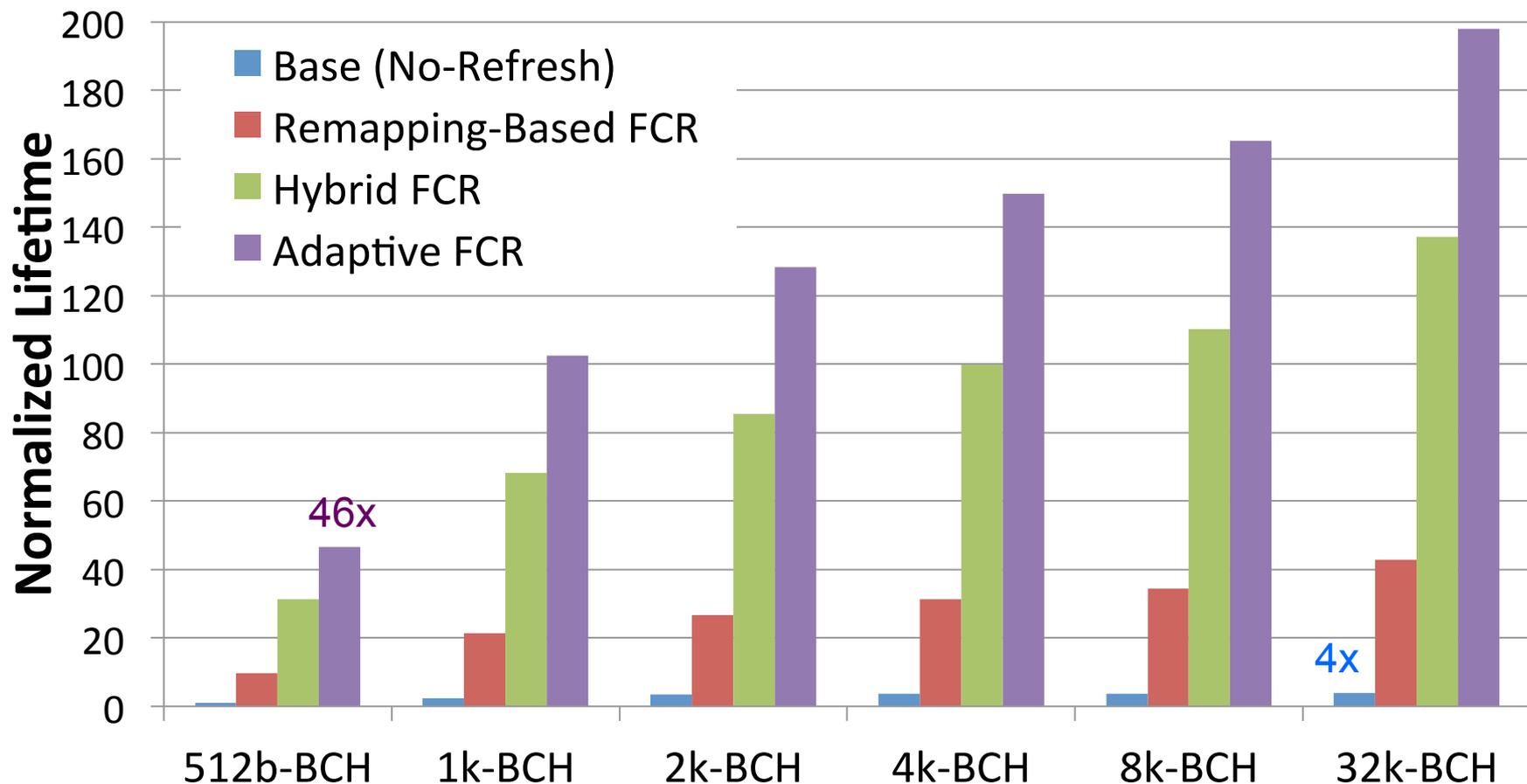
- Experimental flash platform to obtain error rates at different P/E cycles [Cai+ DATE 2012]
- Simulation framework to obtain P/E cycles of real workloads: DiskSim with SSD extensions
- Simulated system: 256GB flash, 4 channels, 8 chips/channel, 8K blocks/chip, 128 pages/block, 8KB pages
- Workloads
  - File system applications, databases, web search
  - Categories: Write-heavy, read-heavy, balanced
- Evaluation metrics
  - Lifetime (extrapolated)
  - Energy overhead, P/E cycle overhead

# Extrapolated Lifetime

---



# Normalized Flash Memory Lifetime



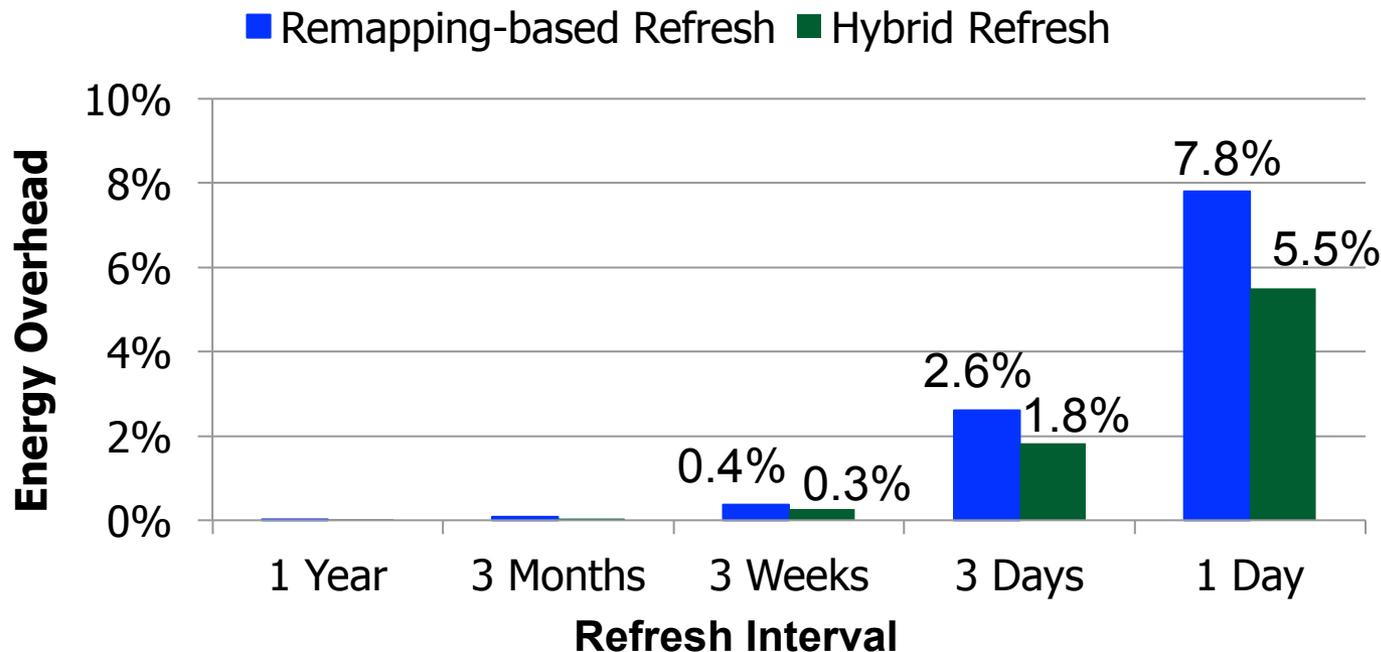
**Lifetime of FCR much higher than lifetime of stronger ECC**

# Lifetime Evaluation Takeaways

---

- Significant average lifetime improvement over no refresh
  - Adaptive-rate FCR: 46X
  - Hybrid reprogramming/remapping based FCR: 31X
  - Remapping based FCR: 9X
- FCR lifetime improvement larger than that of stronger ECC
  - 46X vs. 4X with 32-kbit ECC (over 512-bit ECC)
  - FCR is less complex and less costly than stronger ECC
- Lifetime on all workloads improves with Hybrid FCR
  - Remapping based FCR can degrade lifetime on read-heavy WL
  - Lifetime improvement highest in write-heavy workloads

# Energy Overhead



- Adaptive-rate refresh: <1.8% energy increase until daily refresh is triggered

# Overhead of Additional Erases

---

- Additional erases happen due to remapping of pages
- Low (2%-20%) for write intensive workloads
- High (up to 10X) for read-intensive workloads
- Improved P/E cycle lifetime of all workloads largely outweighs the additional P/E cycles due to remapping

# More Results in the Paper

---

- Detailed workload analysis
- Effect of refresh rate

# Outline

---

- Executive Summary
- The Problem: Limited Flash Memory Endurance/Lifetime
- Error and ECC Analysis for Flash Memory
- Flash Correct and Refresh Techniques (FCR)
- Evaluation
- Conclusions

# Conclusion

---

- NAND flash memory lifetime is limited due to uncorrectable errors, which increase over lifetime (P/E cycles)
- **Observation: Dominant source of errors in flash memory is retention errors** → retention error rate limits lifetime
- **Flash Correct-and-Refresh (FCR) techniques reduce retention error rate to improve flash lifetime**
  - **Periodically read, correct, and remap or reprogram each page** before it accumulates more errors than can be corrected
  - Adapt refresh period to the severity of errors
- **FCR improves flash lifetime by 46X at no hardware cost**
  - More effective and efficient than stronger ECC
  - Can enable better flash memory scaling

# Flash Correct-and-Refresh

## Retention-Aware Error Management for Increased Flash Memory Lifetime

Yu Cai<sup>1</sup> Gulay Yalcin<sup>2</sup> Onur Mutlu<sup>1</sup> Erich F. Haratsch<sup>3</sup>  
Adrian Cristal<sup>2</sup> Osman S. Unsal<sup>2</sup> Ken Mai<sup>1</sup>

<sup>1</sup> Carnegie Mellon University

<sup>2</sup> Barcelona Supercomputing Center

<sup>3</sup> LSI Corporation



**SAFARI**

**Carnegie Mellon**