# Memory Systems in the Multi-Core Era Lecture 2.2: Emerging Technologies and Hybrid Memories

Prof. Onur Mutlu http://www.ece.cmu.edu/~omutlu

onur@cmu.edu

Bogazici University

June 14, 2013



#### SAFARI

## What Will You Learn in Mini Course 2?

- Memory Systems in the Multi-Core Era
   June 13, 14, 17 (1-4pm)
- Lecture 1: Main memory basics, DRAM scaling
- Lecture 2: Emerging memory technologies and hybrid memories
- Lecture 3: Main memory interference and QoS
- 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
- Emerging Memory Technologies
  - http://www.youtube.com/watch? v=LzfOghMKyA0&list=PL5PHm2jkkXmidJOd59REog9jDnPDTG6IJ&index=35
- Multiprocessor Correctness and Cache Coherence
  - <u>http://www.youtube.com/watch?v=U-</u> <u>VZKMgItDM&list=PL5PHm2jkkXmidJOd59REog9jDnPDTG6IJ&index=32</u>

# Readings for Lecture 2.1 (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 Lecture 2.2 (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 Lecture 2.3 (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 Lecture 2.3 (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.
- More to come in next lecture...

# Readings in Flash Memory

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

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

## Online Lectures and More Information

- Online Computer Architecture Lectures
  - <u>http://www.youtube.com/playlist?</u> <u>list=PL5PHm2jkkXmidJOd59REog9jDnPDTG6IJ</u>
- Online Computer Architecture Courses
  - Intro: <u>http://www.ece.cmu.edu/~ece447/s13/doku.php</u>
  - Advanced: <u>http://www.ece.cmu.edu/~ece740/f11/doku.php</u>
  - Advanced: <u>http://www.ece.cmu.edu/~ece742/doku.php</u>
- Recent Research Papers
  - <u>http://users.ece.cmu.edu/~omutlu/projects.htm</u>
  - http://scholar.google.com/citations?
    user=7XyGUGkAAAJ&hl=en

# Emerging Memory Technologies



- 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

#### 



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



- 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)

#### SAFARI

## 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

#### SAFARI



- 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?



- 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





Photo Courtesy: Bipin Rajendran, IBM Slide Courtesy: Moinuddin Qureshi, IBM

## 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=90nm

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

		Table 1. Technology survey.								
	Published prototype									
Parameter*	Horri <sup>6</sup>	Ahn <sup>12</sup>	Bedeschi <sup>13</sup>	Oh14	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 (µm <sup>2</sup> )	••	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 <sup>2</sup>		20.1	9.0	••	12.0		16.6	12.0	5.8	9.0 to 12.0
Access device	**	**	вл	FET	BJT	••	FET	BJT	Diode	BJT
Read time (ns)	**	70	48	68	**	••	62		55	48
Read current (µA)	**	**	40	**	**	••	**		**	40
Read voltage (V)	**	3.0	1.0	1.8	1.6	••	1.8		1.8	1.0
Read power (µW)	**	**	40	**	**	••	••		••	40
Read energy (pJ)	**	**	2.0	**	**	••	••		••	2.0
Set time (ns)	100	150	150	180	**	80	300		400	150
Set current (µA)	200	**	300	200	**	55	••		••	150
Set voltage (V)	**	**	2.0	**	**	1.25	**		**	1.2
Set power (µW)	**	**	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 (µA)	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 (µW)	**	**	1620	**	**	80.4	**		**	480
Reset energy (pJ)	**	**	64.8	**	**	4.8	**	**	**	19.2
Write endurance	107	10 <sup>9</sup>	106	**	10 <sup>8</sup>	104	••	10 <sup>5</sup>	10 <sup>5</sup>	10 <sup>8</sup>

\* 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



## Phase Change Memory Properties

- Dynamic Energy
  - 40 uA Rd, 150 uA Wr
  - 2-43x DRAM, 1x NAND Flash
- Endurance
  - Writes induce phase change at 650C
  - Contacts degrade from thermal expansion/contraction
  - <u>10<sup>8</sup> writes per cell</u>

<sup>10-8</sup>x DRAM, 10<sup>3</sup>x NAND Flash

Cell Size

9-12F<sup>2</sup> using BJT, single-level cells

1.5x DRAM, 2-3x NAND
#### Phase Change Memory: Pros and Cons

- Pros over DRAM
  - Better technology scaling
  - Non volatility
  - Low idle power (no refresh)
- Cons
  - Higher latencies: ~4-15x DRAM (especially write)
  - □ Higher active energy: ~2-50x 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) nonvolatile 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)



- 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



#### 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
  - $\rightarrow$  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?



- 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)





- 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



# 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



# Case 1: RBL-Unaware Policy

#### Access pattern to main memory: A (oldest), B, C, C, C, A, B, D, D, D, A, B (youngest)



# Case 2: RBL-Aware Policy (RBLA)

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



→ Access data at lower row buffer miss latency of DRAM

→ 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)
- 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 unlimitedsized 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

■FREQ ■FREQ-Dyn ■RBLA ■RBLA-Dyn



# Memory Energy Efficiency



## **Thread Fairness**



# Compared to All-PCM/DRAM

RBLA-Dyn □ 16GB DRAM 16GB PCM 2 1.2 Normalized Weighted Speedup 1.6 1.7 1.7 0.8 0 0 0 0 0 7 0 7 0 1 0 0 8 ed Max. Slowdown 1 29% 0.8 31% 0.6 **Our mechanism achieves 31% better performance** than all PCM, within 29% of all DRAM performance 0 0

# 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?
- <u>Observation</u>: PCM array access latency is higher than DRAM's – But peripheral circuit (row buffer) access latencies are similar
- <u>Key Idea</u>: 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** 



- 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

#### SAFARI

#### 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



## TIMBER Tag Management Example (I)

Case 1: TIMBER hit



## TIMBER Tag Management Example (II)

Case 2: TIMBER miss

#### 2. Cache M(Y)



## 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)









# **Dynamic Granularity Performance**



## **TIMBER Performance**



Scalable Hybrid Memories," IEEE Comp. Arch. Letters, 2012.

## **TIMBER Energy Efficiency**



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?

SAFARI

- 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  $\rightarrow$  Wearout attacks

2. Non-volatility  $\rightarrow$  Data persists in memory after powerdown  $\rightarrow$  Easy retrieval of privileged or private information

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

#### Securing Emerging Memory Technologies

- Limited endurance → Wearout attacks
   Better architecting of memory chips to absorb writes
   Hybrid memory system management
   Online wearout attack detection
- 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
  SAFARI



- 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
- We are examining many other solution directions and 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

#### Flash Memory Scaling

## Readings in Flash Memory

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

 Yu Cai, Erich F. Haratsch, <u>Onur Mutlu</u>, and Ken Mai, <u>"Error Patterns in MLC NAND Flash Memory: Measurement, Characterization,</u> <u>and Analysis"</u>

Proceedings of the <u>Design, Automation, and Test in Europe Conference</u> (**DATE**), Dresden, Germany, March 2012. <u>Slides (ppt)</u>



## **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

#### SAFARI

#### 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<sup>-15</sup>) is increasing exponentially to reach *less* endurance

UBER: Uncorrectable bit error rate. Fraction of erroneous bits after error correction. SAFARI
Carnegie Mellon
## Future NAND Flash Storage Architecture



### Need to understand NAND flash error patterns

### SAFARI

#### **Carnegie Mellon**

## **Test System Infrastructure**



#### SAFARI

#### **Carnegie Mellon**

## NAND Flash Testing Platform



NAND Daughter Board

#### **Carnegie Mellon**

#### SAFARI

## NAND Flash Usage and Error Model



#### SAFARI

#### **Carnegie Mellon**

## 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

### SAFARI

#### **Carnegie Mellon**

## 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

### SAFARI

#### **Carnegie Mellon**

# **Retention Error Value Dependency**



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

### SAFARI

#### **Carnegie Mellon**

### More Details on Flash Error Analysis

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



# **Threshold Voltage Distribution Shifts**



**Carnegie Mellon** 

As P/E cycles increase ... Distribution shifts to the right Distribution becomes wider

### **SAFARI**

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

### SAFARI

### Carnegie Mellon<sup>121</sup>

### 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  $\rightarrow$  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

123

**Carnegie Melle** 

### 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<sup>-15</sup>) is increasing exponentially to reach *less* endurance

124

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

### 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
    - $\rightarrow$  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

### SAFARI

## 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
    - $\rightarrow$  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 lengt (n)	h	Correctable Errors (t)	Acceptable Raw BER	Norm. Power	Norm. Area
512		7	1.0x10 <sup>-4</sup> (1x)	1	1
1024		12	4.0x10 <sup>-4</sup> (4x)	2	2.1
2048		22	1.0x10 <sup>-3</sup> (10x)	4.1	3.9
4096		40	1.7x10 <sup>-3</sup> (17x)	8.6	10.3
8192		74	2.2x10 <sup>-3</sup> (22x)	17.8	21.3
32768		259	2.6x10 <sup>-3</sup> (26x)	71	85
	Code lengt (n) 512 1024 2048 4096 8192 32768	Code length (n)512102420484096819232768	Code length (n)Correctable Errors (t)512710241220482240964081927432768259	Code length (n)Correctable Errors (t)Acceptable Raw BER51271.0x10-4 (1x)1024124.0x10-4 (4x)2048221.0x10-3 (10x)4096401.7x10-3 (17x)8192742.2x10-3 (22x)327682592.6x10-3 (26x)	Code length (n)Correctable Errors (t)Acceptable Raw BERNorm. Power51271.0x10-4 (1x)11024124.0x10-4 (4x)22048221.0x10-3 (10x)4.14096401.7x10-3 (17x)8.68192742.2x10-3 (22x)17.8327682592.6x10-3 (26x)71

Carnegie Mellon<sup>131</sup>

## Resulting Flash Lifetime with Strong ECC

Lifetime improvement comparison of various BCH codes



Strong ECC is very inefficient at improving lifetime

### SAFARI

### 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

### SAFARI

### FCR Intuition

	Errors with No refresh	Errors with Periodic refresh
Program Page	×	×
After time T	<b>× × ×</b>	××××
After time 2T	× × × × ×	× × × ×
After time 3T	$\times \times \times \times \times \times \times$	××××

× Retention Error × Program Error

### **SAFARI**

### **Carnegie Mellon**<sup>136</sup>

## 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)
  Select next Block
  - 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



140

**Carnegie Mell** 

• Problem: Causes additional erase operations  $\rightarrow$  more wearout

- Bad for read-intensive workloads (few erases really needed)
- Lifetime degrades for such workloads (see paper)

### SAFARI

### 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



Pro: No remapping needed  $\rightarrow$  no additional erase operations

Con: Increases the occurrence of program errors

### SAFARI

### **Carnegie Mellon**<sup>143</sup>

## Program Errors in Flash Memory

- When a cell is being programmed, voltage level of a neighboring cell changes (unintentionally) due to parasitic capacitance coupling
  - $\rightarrow$  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



SAFARI

### **Carnegie Mellon**<sup>145</sup>

# 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

- 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

#### SAFARI

### Carnegie Mellon<sup>149</sup>

- 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

- 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





# Normalized Flash Memory Lifetime



Carnegie Mellon<sup>155</sup>

# 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

Remapping-based Refresh Hybrid Refresh



 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



- 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**