# Opportunities and Challenges of Emerging Memory Technologies

Onur Mutlu

omutlu@gmail.com

https://people.inf.ethz.ch/omutlu

September 11, 2017 ARM Research Summit



**ETH** zürich

## The Main Memory System



- Main memory is a critical component of all computing systems: server, mobile, embedded, desktop, sensor
- Main memory system must scale (in size, technology, efficiency, cost, and management algorithms) to maintain performance growth and technology scaling benefits

## The Main Memory System



- Main memory is a critical component of all computing systems: server, mobile, embedded, desktop, sensor
- Main memory system must scale (in size, technology, efficiency, cost, and management algorithms) to maintain performance growth and technology scaling benefits

## The Main Memory System



- Main memory is a critical component of all computing systems: server, mobile, embedded, desktop, sensor
- Main memory system must scale (in size, technology, efficiency, cost, and management algorithms) to maintain performance growth and technology scaling benefits

## Memory System: A *Shared Resource* View



Most of the system is dedicated to storing and moving data

### SAFARI

## State of the Main Memory System

- Recent technology, architecture, and application trends
  - lead to new requirements
  - exacerbate old requirements
- DRAM and memory controllers, as we know them today, are (will be) unlikely to satisfy all requirements
- Some emerging non-volatile memory technologies (e.g., PCM) enable new opportunities: memory+storage merging
- We need to rethink the main memory system
   to fix DRAM issues and enable emerging technologies
   to satisfy all requirements

## Major Trends Affecting Main Memory (I)

Need for main memory capacity, bandwidth, QoS increasing

Main memory energy/power is a key system design concern

DRAM technology scaling is ending

## Major Trends Affecting Main Memory (II)

- Need for main memory capacity, bandwidth, QoS increasing
  - Multi-core: increasing number of cores/agents
  - Data-intensive applications: increasing demand/hunger for data
  - Consolidation: cloud computing, GPUs, mobile, heterogeneity

• Main memory energy/power is a key system design concern

DRAM technology scaling is ending

## Example: The Memory Capacity Gap

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



*Memory capacity per core* expected to drop by 30% every two years
Trends worse for *memory bandwidth per core*!

## Major Trends Affecting Main Memory (III)

Need for main memory capacity, bandwidth, QoS increasing

- Main memory energy/power is a key system design concern
  - ~40-50% energy spent in off-chip memory hierarchy [Lefurgy, IEEE Computer'03] >40% power in DRAM [Ware, HPCA'10][Paul,ISCA'15]
  - DRAM consumes power even when not used (periodic refresh)
- DRAM technology scaling is ending

## Major Trends Affecting Main Memory (IV)

Need for main memory capacity, bandwidth, QoS increasing

## Main memory energy/power is a key system design concern

## DRAM technology scaling is ending

- ITRS projects DRAM will not scale easily below X nm
- Scaling has provided many benefits:
  - higher capacity (density), lower cost, lower energy

## Major Trends Affecting Main Memory (V)

- DRAM scaling has already become increasingly difficult
  - Increasing cell leakage current, reduced cell reliability, increasing manufacturing difficulties [Kim+ ISCA 2014], [Liu+ ISCA 2013], [Mutlu IMW 2013], [Mutlu DATE 2017]
  - Difficult to significantly improve capacity, energy

### Emerging memory technologies are promising

## Major Trends Affecting Main Memory (V)

- DRAM scaling has already become increasingly difficult
  - Increasing cell leakage current, reduced cell reliability, increasing manufacturing difficulties [Kim+ ISCA 2014], [Liu+ ISCA 2013], [Mutlu IMW 2013], [Mutlu DATE 2017]
  - Difficult to significantly improve capacity, energy

## Emerging memory technologies are promising

3D-Stacked DRAM	higher bandwidth	smaller capacity	
Reduced-Latency DRAM (e.g., RLDRAM, TL-DRAM)	lower latency	higher cost	
<b>Low-Power DRAM</b> (e.g., LPDDR3, LPDDR4)	lower power	higher latency higher cost	
Non-Volatile Memory (NVM) (e.g., PCM, STTRAM, ReRAM, 3D Xpoint)	larger capacity	higher latency higher dynamic power lower endurance	



- Major Trends Affecting Main Memory
- The Memory Scaling Problem
- Two Complementary Solution Directions
  - New Memory (DRAM) Architectures
  - Enabling Emerging (NVM) Technologies
- Conclusion

## 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



## 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

### SAFARI

## As Memory Scales, It Becomes Unreliable

- Data from all of Facebook's servers worldwide
- Meza+, "Revisiting Memory Errors in Large-Scale Production Data Centers," DSN'15.



Chip density (Gb)

## Large-Scale Failure Analysis of DRAM Chips

- Analysis and modeling of memory errors found in all of Facebook's server fleet
- Justin Meza, Qiang Wu, Sanjeev Kumar, and Onur Mutlu, "Revisiting Memory Errors in Large-Scale Production Data Centers: Analysis and Modeling of New Trends from the Field" Proceedings of the <u>45th Annual IEEE/IFIP International Conference on</u> Dependable Systems and Networks (DSN), Rio de Janeiro, Brazil, June 2015. [Slides (pptx) (pdf)] [DRAM Error Model]

Revisiting Memory Errors in Large-Scale Production Data Centers: Analysis and Modeling of New Trends from the Field

Justin Meza Qiang Wu\* Sanjeev Kumar\* Onur Mutlu

Carnegie Mellon University \* Facebook, Inc.

## Infrastructures to Understand Such Issues



### **SAFARI**

Kim+, "Flipping Bits in Memory Without Accessing Them: An Experimental Study of DRAM Disturbance Errors," ISCA 2014.

## SoftMC: Open Source DRAM Infrastructure

 Hasan Hassan et al., "<u>SoftMC: A</u> <u>Flexible and Practical Open-</u> <u>Source Infrastructure for</u> <u>Enabling Experimental DRAM</u> <u>Studies</u>," HPCA 2017.

- Flexible
- Easy to Use (C++ API)
- Open-source

github.com/CMU-SAFARI/SoftMC



### SAFARI



### <u>https://github.com/CMU-SAFARI/SoftMC</u>

## SoftMC: A Flexible and Practical Open-Source Infrastructure for Enabling Experimental DRAM Studies

Hasan Hassan<sup>1,2,3</sup> Nandita Vijaykumar<sup>3</sup> Samira Khan<sup>4,3</sup> Saugata Ghose<sup>3</sup> Kevin Chang<sup>3</sup> Gennady Pekhimenko<sup>5,3</sup> Donghyuk Lee<sup>6,3</sup> Oguz Ergin<sup>2</sup> Onur Mutlu<sup>1,3</sup>

<sup>1</sup>ETH Zürich <sup>2</sup>TOBB University of Economics & Technology <sup>3</sup>Carnegie Mellon University <sup>4</sup>University of Virginia <sup>5</sup>Microsoft Research <sup>6</sup>NVIDIA Research A Curious Discovery [Kim et al., ISCA 2014]

# One can predictably induce errors in most DRAM memory chips

# A simple hardware failure mechanism can create a widespread system security vulnerability



## Modern DRAM is Prone to Disturbance Errors



Repeatedly reading a row enough times (before memory gets refreshed) induces disturbance errors in adjacent rows in most real DRAM chips you can buy today

<u>Flipping Bits in Memory Without Accessing Them: An Experimental Study of DRAM</u> <u>Disturbance Errors</u>, (Kim et al., ISCA 2014)

## More on RowHammer Analysis

Yoongu Kim, Ross Daly, Jeremie Kim, Chris Fallin, Ji Hye Lee, Donghyuk Lee, Chris Wilkerson, Konrad Lai, and Onur Mutlu,
 "Flipping Bits in Memory Without Accessing Them: An

 Experimental Study of DRAM Disturbance Errors"
 Proceedings of the <u>41st International Symposium on Computer</u>
 <u>Architecture</u> (ISCA), Minneapolis, MN, June 2014.

 [Slides (pptx) (pdf)] [Lightning Session Slides (pptx) (pdf)] [Source Code and Data]

## Flipping Bits in Memory Without Accessing Them: An Experimental Study of DRAM Disturbance Errors

Yoongu Kim<sup>1</sup> Ross Daly<sup>\*</sup> Jeremie Kim<sup>1</sup> Chris Fallin<sup>\*</sup> Ji Hye Lee<sup>1</sup> Donghyuk Lee<sup>1</sup> Chris Wilkerson<sup>2</sup> Konrad Lai Onur Mutlu<sup>1</sup> <sup>1</sup>Carnegie Mellon University <sup>2</sup>Intel Labs

## Industry Is Writing Papers About It, Too

### **DRAM Process Scaling Challenges**

#### \* Refresh

- · Difficult to build high-aspect ratio cell capacitors decreasing cell capacitance
- · Leakage current of cell access transistors increasing

### ✤ tWR

- · Contact resistance between the cell capacitor and access transistor increasing
- · On-current of the cell access transistor decreasing
- · Bit-line resistance increasing

### VRT

· Occurring more frequently with cell capacitance decreasing



## Industry Is Writing Papers About It, Too

### **DRAM Process Scaling Challenges**

### \* Refresh

Difficult to build high-aspect ratio cell capacitors decreasing cell capacitance
THE MEMORY FORUM 2014

## Co-Architecting Controllers and DRAM to Enhance DRAM Process Scaling

Uksong Kang, Hak-soo Yu, Churoo Park, \*Hongzhong Zheng, \*\*John Halbert, \*\*Kuljit Bains, SeongJin Jang, and Joo Sun Choi



27

Samsung Electronics, Hwasung, Korea / \*Samsung Electronics, San Jose / \*\*Intel

## How Do We Solve The Memory Problem?



software/hardware/device cooperation



- Major Trends Affecting Main Memory
- The Memory Scaling Problem
- Two Complementary Solution Directions
  - New Memory (DRAM) Architectures
  - Enabling Emerging (NVM) Technologies
- Conclusion

## Solution 1: New Memory Architectures

- Overcome memory shortcomings with
  - Memory-centric system design
  - Novel memory architectures, interfaces, functions
  - Better waste management (efficient utilization)
- Key issues to tackle
  - Enable reliability at low cost  $\rightarrow$  high capacity
  - Reduce energy
  - Improve latency and bandwidth
  - Reduce waste (capacity, bandwidth, latency)
  - Enable computation close to data

## Solution 1: New Memory Architectures

- Liu+, "RAIDR: Retention-Aware Intelligent DRAM Refresh," ISCA 2012.
- Kim+, "A Case for Exploiting Subarray-Level Parallelism in DRAM," ISCA 2012.
- Lee+, "Tiered-Latency DRAM: A Low Latency and Low Cost DRAM Architecture," HPCA 2013.
- Liu+, "An Experimental Study of Data Retention Behavior in Modern DRAM Devices," ISCA 2013.
- Seshadri+, "RowClone: Fast and Efficient In-DRAM Copy and Initialization of Bulk Data," MICRO 2013.
- Pekhimenko+, "Linearly Compressed Pages: A Main Memory Compression Framework," MICRO 2013.
- Chang+, "Improving DRAM Performance by Parallelizing Refreshes with Accesses," HPCA 2014.
- Khan+, "The Efficacy of Error Mitigation Techniques for DRAM Retention Failures: A Comparative Experimental Study," SIGMETRICS 2014.
- Luo+, "Characterizing Application Memory Error Vulnerability to Optimize Data Center Cost," DSN 2014.
- Kim+, "Flipping Bits in Memory Without Accessing Them: An Experimental Study of DRAM Disturbance Errors," ISCA 2014.
- Lee+, "Adaptive-Latency DRAM: Optimizing DRAM Timing for the Common-Case," HPCA 2015.
- Qureshi+, "AVATAR: A Variable-Retention-Time (VRT) Aware Refresh for DRAM Systems," DSN 2015.
- Meza+, "Revisiting Memory Errors in Large-Scale Production Data Centers: Analysis and Modeling of New Trends from the Field," DSN 2015.
- Kim+, "Ramulator: A Fast and Extensible DRAM Simulator," IEEE CAL 2015.
- Seshadri+, "Fast Bulk Bitwise AND and OR in DRAM," IEEE CAL 2015.
- Ahn+, "A Scalable Processing-in-Memory Accelerator for Parallel Graph Processing," ISCA 2015.
- Ahn+, "PIM-Enabled Instructions: A Low-Overhead, Locality-Aware Processing-in-Memory Architecture," ISCA 2015.
- Lee+, "Decoupled Direct Memory Access: Isolating CPU and IO Traffic by Leveraging a Dual-Data-Port DRAM," PACT 2015.
- Seshadri+, "Gather-Scatter DRAM: In-DRAM Address Translation to Improve the Spatial Locality of Non-unit Strided Accesses," MICRO 2015.
- Lee+, "Simultaneous Multi-Layer Access: Improving 3D-Stacked Memory Bandwidth at Low Cost," TACO 2016.
- Hassan+, "ChargeCache: Reducing DRAM Latency by Exploiting Row Access Locality," HPCA 2016.
- Chang+, "Low-Cost Inter-Linked Subarrays (LISA): Enabling Fast Inter-Subarray Data Migration in DRAM," HPCA 2016.
- Chang+, "Understanding Latency Variation in Modern DRAM Chips Experimental Characterization, Analysis, and Optimization," SIGMETRICS 2016.
- Khan+, "PARBOR: An Efficient System-Level Technique to Detect Data Dependent Failures in DRAM," DSN 2016.
- Hsieh+, "Transparent Offloading and Mapping (TOM): Enabling Programmer-Transparent Near-Data Processing in GPU Systems," ISCA 2016.
- Hashemi+, "Accelerating Dependent Cache Misses with an Enhanced Memory Controller," ISCA 2016.
- Boroumand+, "LazyPIM: An Efficient Cache Coherence Mechanism for Processing-in-Memory," IEEE CAL 2016.
- Pattnaik+, "Scheduling Techniques for GPU Architectures with Processing-In-Memory Capabilities," PACT 2016.
- Hsieh+, "Accelerating Pointer Chasing in 3D-Stacked Memory: Challenges, Mechanisms, Evaluation," ICCD 2016.
- Hashemi+, "Continuous Runahead: Transparent Hardware Acceleration for Memory Intensive Workloads," MICRO 2016.
- Khan+, "A Case for Memory Content-Based Detection and Mitigation of Data-Dependent Failures in DRAM"," IEEE CAL 2016.
- Hassan+, "SoftMC: A Flexible and Practical Open-Source Infrastructure for Enabling Experimental DRAM Studies," HPCA 2017.
- Mutlu, "The RowHammer Problem and Other Issues We May Face as Memory Becomes Denser," DATE 2017.
- Lee+, "Design-Induced Latency Variation in Modern DRAM Chips: Characterization, Analysis, and Latency Reduction Mechanisms," SIGMETRICS 2017.
- Chang+, "Understanding Reduced-Voltage Operation in Modern DRAM Devices: Experimental Characterization, Analysis, and Mechanisms," SIGMETRICS 2017.
- Patel+, "The Reach Profiler (REAPER): Enabling the Mitigation of DRAM Retention Failures via Profiling at Aggressive Conditions," ISCA 2017.
- Seshadri and Mutlu, "Simple Operations in Memory to Reduce Data Movement," ADCOM 2017.
- Liu+, "Concurrent Data Structures for Near-Memory Computing," SPAA 2017.
- Khan+, "Detecting and Mitigating Data-Dependent DRAM Failures by Exploiting Current Memory Content," MICRO 2017.
- Seshadri+, "Ambit: In-Memory Accelerator for Bulk Bitwise Operations Using Commodity DRAM Technology," MICRO 2017.
- Avoid DRAM:
  - Seshadri+, "The Evicted-Address Filter: A Unified Mechanism to Address Both Cache Pollution and Thrashing," PACT 2012.
  - Pekhimenko+, "Base-Delta-Immediate Compression: Practical Data Compression for On-Chip Caches," PACT 2012.
  - Seshadri+, "The Dirty-Block Index," ISCA 2014.
  - Pekhimenko+, "Exploiting Compressed Block Size as an Indicator of Future Reuse," HPCA 2015.
  - Vijaykumar+, "A Case for Core-Assisted Bottleneck Acceleration in GPUs: Enabling Flexible Data Compression with Assist Warps," ISCA 2015.
  - Pekhimenko+, "Toggle-Aware Bandwidth Compression for GPUs," HPCA 2016.

# Example: (Truly) In-Memory Computation

- We can support in-DRAM COPY, ZERO, AND, OR, NOT, MAJ
- At low cost
- Using analog computation capability of DRAM
  - Idea: activating multiple rows performs computation
- 30-60X performance and energy improvement
  - Seshadri+, "Ambit: In-Memory Accelerator for Bulk Bitwise Operations Using Commodity DRAM Technology," MICRO 2017.

- New memory technologies enable even more opportunities
  - Memristors, resistive RAM, phase change mem, STT-MRAM, ...
  - Can operate on data with minimal movement

## Performance: In-DRAM Bitwise Operations



Figure 9: Throughput of bitwise operations on various systems.

### SAFARI

	Design	not	and/or	nand/nor	xor/xnor
DRAM &	DDR3	93.7	137.9	137.9	137.9
Channel Energy	Ambit	1.6	3.2	4.0	5.5
(nJ/KB)	$(\downarrow)$	59.5X	43.9X	35.1X	25.1X

Table 3: Energy of bitwise operations. ( $\downarrow$ ) indicates energy reduction of Ambit over the traditional DDR3-based design.

Seshadri+, "Ambit: In-Memory Accelerator for Bulk Bitwise Operations using Commodity DRAM Technology," MICRO 2017.

## More on Ambit

 Vivek Seshadri et al., "<u>Ambit: In-Memory Accelerator</u> for Bulk Bitwise Operations Using Commodity DRAM <u>Technology</u>," MICRO 2017.

## Ambit: In-Memory Accelerator for Bulk Bitwise Operations Using Commodity DRAM Technology

Vivek Seshadri<sup>1,5</sup> Donghyuk Lee<sup>2,5</sup> Thomas Mullins<sup>3,5</sup> Hasan Hassan<sup>4</sup> Amirali Boroumand<sup>5</sup> Jeremie Kim<sup>4,5</sup> Michael A. Kozuch<sup>3</sup> Onur Mutlu<sup>4,5</sup> Phillip B. Gibbons<sup>5</sup> Todd C. Mowry<sup>5</sup>

<sup>1</sup>Microsoft Research India <sup>2</sup>NVIDIA Research <sup>3</sup>Intel <sup>4</sup>ETH Zürich <sup>5</sup>Carnegie Mellon University

### **SAFARI**

## Tesseract System for Graph Processing

Interconnected set of 3D-stacked memory+logic chips with simple cores



**SAFARI** Ahn+, "A Scalable Processing-in-Memory Accelerator for Parallel Graph Processing" ISCA 2015.
### Tesseract Graph Processing Performance

#### >13X Performance Improvement



**SAFARI** Ahn+, "A Scalable Processing-in-Memory Accelerator for Parallel Graph Processing" ISCA 2015.

### Tesseract Graph Processing System Energy



**SAFARI** Ahn+, "A Scalable Processing-in-Memory Accelerator for Parallel Graph Processing" ISCA 2015.

#### More on Tesseract

 Junwhan Ahn, Sungpack Hong, Sungjoo Yoo, Onur Mutlu, and Kiyoung Choi,
 "A Scalable Processing-in-Memory Accelerator for Parallel Graph Processing"
 Proceedings of the <u>42nd International Symposium on</u> <u>Computer Architecture</u> (ISCA), Portland, OR, June 2015.
 [Slides (pdf)] [Lightning Session Slides (pdf)]

#### A Scalable Processing-in-Memory Accelerator for Parallel Graph Processing

Junwhan Ahn Sungpack Hong<sup>§</sup> Sungjoo Yoo Onur Mutlu<sup>†</sup> Kiyoung Choi junwhan@snu.ac.kr, sungpack.hong@oracle.com, sungjoo.yoo@gmail.com, onur@cmu.edu, kchoi@snu.ac.kr Seoul National University <sup>§</sup>Oracle Labs <sup>†</sup>Carnegie Mellon University

### Accelerating Pointer Chasing with PIM

 Kevin Hsieh, Samira Khan, Nandita Vijaykumar, Kevin K. Chang, Amirali Boroumand, Saugata Ghose, and Onur Mutlu, <u>"Accelerating Pointer Chasing in 3D-Stacked Memory:</u> <u>Challenges, Mechanisms, Evaluation"</u> *Proceedings of the <u>34th IEEE International Conference on Computer</u> <u>Design</u> (ICCD), Phoenix, AZ, USA, October 2016.* 

#### Accelerating Pointer Chasing in 3D-Stacked Memory: Challenges, Mechanisms, Evaluation

Kevin Hsieh<sup>†</sup> Samira Khan<sup>‡</sup> Nandita Vijaykumar<sup>†</sup> Kevin K. Chang<sup>†</sup> Amirali Boroumand<sup>†</sup> Saugata Ghose<sup>†</sup> Onur Mutlu<sup>§†</sup> <sup>†</sup>Carnegie Mellon University <sup>‡</sup>University of Virginia <sup>§</sup>ETH Zürich



- Major Trends Affecting Main Memory
- The Memory Scaling Problem
- Two Complementary Solution Directions
  - New Memory (DRAM) Architectures
  - Enabling Emerging (NVM) Technologies
- Conclusion

### 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



### Solution 2: Emerging Memory Technologies

- Some emerging resistive memory technologies seem more scalable than DRAM (and they are non-volatile)
- Example: Phase Change Memory
  - Data stored by changing phase of material
  - Data read by detecting material's resistance
  - Expected to scale to 9nm (2022 [ITRS 2009])
  - Prototyped at 20nm (Raoux+, IBM JRD 2008)
  - Expected to be denser than DRAM: can store multiple bits/cell
- But, emerging technologies have (many) shortcomings
  Can they be enabled to replace/augment/surpass DRAM?



# Solution 2: Emerging Memory Technologies

- Lee+, "Architecting Phase Change Memory as a Scalable DRAM Alternative," ISCA'09, CACM'10, IEEE Micro'10.
- Meza+, "Enabling Efficient and Scalable Hybrid Memories," IEEE Comp. Arch. Letters 2012.
- Yoon, Meza+, "Row Buffer Locality Aware Caching Policies for Hybrid Memories," ICCD 2012.
- Kultursay+, "Evaluating STT-RAM as an Energy-Efficient Main Memory Alternative," ISPASS 2013.
- Meza+, "A Case for Efficient Hardware-Software Cooperative Management of Storage and Memory," WEED 2013.
- Lu+, "Loose Ordering Consistency for Persistent Memory," ICCD 2014.
- Zhao+, "FIRM: Fair and High-Performance Memory Control for Persistent Memory Systems," MICRO 2014.
- Yoon, Meza+, "Efficient Data Mapping and Buffering Techniques for Multi-Level Cell Phase-Change Memories," TACO 2014.
- Ren+, "ThyNVM: Enabling Software-Transparent Crash Consistency in Persistent Memory Systems," MICRO 2015.
- Chauhan+, "NVMove: Helping Programmers Move to Byte-Based Persistence," INFLOW 2016.
- Li+, "Utility-Based Hybrid Memory Management," CLUSTER 2017.
- Yu+, "Banshee: Bandwidth-Efficient DRAM Caching via Software/Hardware Cooperation," MICRO 2017.

# Promising 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/RRAM/ReRAM
  - Inject current to change atomic structure
  - Resistance determined by atom distance

#### **SAFARI**

### 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.
- Lee et al., "Phase Change Technology and the Future of Main Memory," IEEE Micro Top Picks 2010.

		Table 1. Technology survey.								
Parameter*	Published prototype									
	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 (MLC)	10 <sup>7</sup>	10 <sup>9</sup>	10 <sup>6</sup>	••	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: Pros and Cons

#### Pros over DRAM

- Better technology scaling (capacity and cost)
- □ Non volatile  $\rightarrow$  Persistent
- 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)
- Reliability issues (resistance drift)
- Challenges in enabling PCM as DRAM replacement/helper:
  - Mitigate PCM shortcomings
  - □ Find the right way to place PCM in the system

#### SAFARI

## PCM-based Main Memory (I)

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



Hybrid PCM+DRAM [Qureshi+ ISCA'09, Dhiman+ DAC'09]:

How to partition/migrate data between PCM and DRAM

## 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



# 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

 $\triangleright$  9 - 12 $F^2$  using BJT

▷ 1.5× DRAM

#### Endurance

▷ 1E+08 writes



#### Latency

50ns Rd, 150ns Wr

⊳ 4×, 12× DRAM

#### Energy

▷ 40µA Rd, 150µA Wr

 $\triangleright$  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
  - $\rightarrow$  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?

#### SAFARI

### More on PCM As Main Memory

 Benjamin C. Lee, Engin Ipek, Onur Mutlu, and Doug Burger, <u>"Architecting Phase Change Memory as a Scalable DRAM</u> <u>Alternative"</u>

Proceedings of the <u>36th International Symposium on Computer</u> <u>Architecture</u> (**ISCA**), pages 2-13, Austin, TX, June 2009. <u>Slides</u> (pdf)

#### Architecting Phase Change Memory as a Scalable DRAM Alternative

Benjamin C. Lee† Engin Ipek† Onur Mutlu‡ Doug Burger†

†Computer Architecture Group Microsoft Research Redmond, WA {blee, ipek, dburger}@microsoft.com

SAFARI

‡Computer Architecture Laboratory Carnegie Mellon University Pittsburgh, PA onur@cmu.edu

### More on PCM As Main Memory (II)

 Benjamin C. Lee, Ping Zhou, Jun Yang, Youtao Zhang, Bo Zhao, Engin Ipek, Onur Mutlu, and Doug Burger,
 "Phase Change Technology and the Future of Main Memory" IEEE Micro, Special Issue: Micro's Top Picks from 2009 Computer Architecture Conferences (MICRO TOP PICKS), Vol. 30, No. 1, pages 60-70, January/February 2010.

# Phase-Change Technology and the Future of Main Memory

### STT-MRAM as Main Memory

- Magnetic Tunnel Junction (MTJ) device
  - Reference layer: Fixed magnetic orientation
  - □ Free layer: Parallel or anti-parallel
- Magnetic orientation of the free layer determines logical state of device
  - High vs. low resistance
- Write: Push large current through MTJ to change orientation of free layer
- Read: Sense current flow

SAFARI

 Kultursay et al., "Evaluating STT-RAM as an Energy-Efficient Main Memory Alternative," ISPASS 2013.





### STT-MRAM: Pros and Cons

#### Pros over DRAM

- Better technology scaling (capacity and cost)
- □ Non volatile  $\rightarrow$  Persistent
- Low idle power (no refresh)

#### Cons

- Higher write latency
- Higher write energy
- Poor density (currently)
- Reliability?
- Another level of freedom
  - Can trade off non-volatility for lower write latency/energy (by reducing the size of the MTJ)

### Architected STT-MRAM as Main Memory

- 4-core, 4GB main memory, multiprogrammed workloads
- ~6% performance loss, ~60% energy savings vs. DRAM



Kultursay+, "Evaluating STT-RAM as an Energy-Efficient Main Memory Alternative," ISPASS 2013.

#### **SAFARI**

### More on STT-MRAM as Main Memory

 Emre Kultursay, Mahmut Kandemir, Anand Sivasubramaniam, and Onur Mutlu,
 "Evaluating STT-RAM as an Energy-Efficient Main Memory Alternative"
 Proceedings of the 2013 IEEE International Symposium on Performance Analysis of Systems and Software (ISPASS),

Austin, TX, April 2013. Slides (pptx) (pdf)

# Evaluating STT-RAM as an Energy-Efficient Main Memory Alternative

Emre Kültürsay\*, Mahmut Kandemir\*, Anand Sivasubramaniam\*, and Onur Mutlu<sup>†</sup> \*The Pennsylvania State University and <sup>†</sup>Carnegie Mellon University

#### A More Viable Approach: Hybrid Memory Systems



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

Meza+, "Enabling Efficient and Scalable Hybrid Memories," IEEE Comp. Arch. Letters, 2012. Yoon+, "Row Buffer Locality Aware Caching Policies for Hybrid Memories," ICCD 2012 Best Paper Award.

#### SAFARI

# Providing the Best of Multiple Metrics with

# **Multiple Memory Technologies**



Heterogeneous, Configurable, Programmable **Memory Systems** 

### Hybrid Memory Systems: Issues

- Cache vs. Main Memory
- Granularity of Data Move/Manage-ment: Fine or Coarse
- Hardware vs. Software vs. HW/SW Cooperative
- When to migrate data?
- How to design a scalable and efficient large cache?

### Data Placement in Hybrid Memory



# Which memory do we place each page in, to maximize system performance?

- Memory A is fast, but small
- Load should be balanced on both channels
- Page migrations have performance and energy overhead
  SAFARI

### Data Placement Between DRAM and PCM

- Idea: Characterize data access patterns and guide data placement in hybrid memory
- Streaming accesses: As fast in PCM as in DRAM
- Random accesses: Much faster in DRAM
- Idea: Place random access data with some reuse in DRAM; streaming data in PCM
- Yoon+, "Row Buffer Locality-Aware Data Placement in Hybrid Memories," ICCD 2012 Best Paper Award.

# Hybrid vs. All-PCM/DRAM [ICCD'12]

■ 16GB PCM ■ RBLA-Dyn ■ 16GB DRAM



Yoon+, "Row Buffer Locality-Aware Data Placement in Hybrid Memories," ICCD 2012 Best Paper Award.

#### More on Hybrid Memory Data Placement

 HanBin Yoon, Justin Meza, Rachata Ausavarungnirun, Rachael Harding, and Onur Mutlu,
 "Row Buffer Locality Aware Caching Policies for Hybrid Memories"

Proceedings of the <u>30th IEEE International Conference on</u> <u>Computer Design</u> (**ICCD**), Montreal, Quebec, Canada, September 2012. <u>Slides (pptx) (pdf)</u>

#### Row Buffer Locality Aware Caching Policies for Hybrid Memories

HanBin Yoon, Justin Meza, Rachata Ausavarungnirun, Rachael A. Harding and Onur Mutlu Carnegie Mellon University {hanbinyoon,meza,rachata,onur}@cmu.edu, rhardin@mit.edu

## Weaknesses of Existing Solutions

- They are all heuristics that consider only a *limited part* of memory access behavior
- Do not *directly* capture the overall system performance impact of data placement decisions

- Example: None capture memory-level parallelism (MLP)
  - Number of *concurrent* memory requests from the same application when a page is accessed
  - Affects how much page migration helps performance
#### Importance of Memory-Level Parallelism



#### Page migration decisions need to consider MLP

A **generalized** mechanism that

1. Directly estimates the **performance benefit of migrating a page** between **any two types of memory** 

2. Places **only** the **performance-critical data** in the fast memory

# Utility-Based Hybrid Memory Management

A memory manager that works for *any* hybrid memory
 e.g., DRAM-NVM, DRAM-RLDRAM

#### Key Idea

- For each page, use comprehensive characteristics to calculate estimated *utility* (i.e., performance impact) of migrating page from one memory to the other in the system
- Migrate only pages with the highest utility (i.e., pages that improve system performance the most when migrated)
- Li+, "Utility-Based Hybrid Memory Management", CLUSTER 2017.
   SAFARI

#### Key Mechanisms of UH-MEM

For each page, estimate utility using a performance model

#### Application stall time reduction

How much would migrating a page benefit the performance of the application that the page belongs to?

Application performance sensitivity

How much does the improvement of a single application's performance increase the *overall* system performance?

 $Utility = \Delta StallTime_i \times Sensitivity_i$ 

 Migrate only pages whose utility exceed the migration threshold from slow memory to fast memory

Periodically adjust migration threshold

#### SAFARI

#### Results: System Performance



UH-MEM improves system performance over the best state-of-the-art hybrid memory manager



#### Results: Sensitivity to Slow Memory Latency



UH-MEM improves system performance for a wide variety of hybrid memory systems

#### More on UH-MEM

 Yang Li, Saugata Ghose, Jongmoo Choi, Jin Sun, Hui Wang, and Onur Mutlu,
 <u>"Utility-Based Hybrid Memory Management"</u> *Proceedings of the <u>19th IEEE Cluster Conference</u> (CLUSTER), Honolulu, Hawaii, USA, September 2017.
 [Slides (pptx) (pdf)]* 

#### **Utility-Based Hybrid Memory Management**

Yang Li<sup>†</sup> Saugata Ghose<sup>†</sup> <sup>†</sup>Carnegie Mellon University Jongmoo Choi<sup>‡</sup> Jin Sun<sup>†</sup> Hui Wang<sup>\*</sup> <sup>‡</sup>Dankook University <sup>\*</sup>Beihang University Onur Mutlu<sup> $+\dagger$ </sup> +*ETH Zürich*  Challenge and Opportunity

Enabling an Emerging Technology to Augment DRAM

# Managing Hybrid Memories



# Designing Effective Large (DRAM) Caches

#### One 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

# On Large DRAM Cache Design

 Justin Meza, Jichuan Chang, HanBin Yoon, Onur Mutlu, and Parthasarathy Ranganathan,
 <u>"Enabling Efficient and Scalable Hybrid Memories</u> <u>Using Fine-Granularity DRAM Cache Management"</u> <u>IEEE Computer Architecture Letters</u> (CAL), February 2012.

#### Enabling Efficient and Scalable Hybrid Memories Using Fine-Granularity DRAM Cache Management

Justin Meza\* Jichuan Chang<sup>†</sup> HanBin Yoon\* Onur Mutlu\* Parthasarathy Ranganathan<sup>†</sup> \*Carnegie Mellon University <sup>†</sup>Hewlett-Packard Labs {meza,hanbinyoon,onur}@cmu.edu {jichuan.chang,partha.ranganathan}@hp.com

#### DRAM Caches: Many Recent Options

**Table 1: Summary of Operational Characteristics of Different State-of-the-Art DRAM Cache Designs** – We assume perfect way prediction for Unison Cache. Latency is relative to the access time of the off-package DRAM (see Section 6 for baseline latencies). We use different colors to indicate the high (dark red), medium (white), and low (light green) overhead of a characteristic.

Scheme	DRAM Cache Hit	DRAM Cache Miss	Replacement Traffic	<b>Replacement Decision</b>	Large Page Caching
Unison [32]	In-package traffic: 128 B	In-package traffic: 96 B	On every miss	Hardware managed,	Yes
	(data + tag read and up-	(spec. data + tag read)	Footprint size [31]	set-associative,	
	date)	Latency: ~2x		LRU	
	Latency: ~1x				
Alloy [50]	In-package traffic: 96 B	In-package traffic: 96 B	On some misses	Hardware managed,	Yes
	(data + tag read)	(spec. data + tag read)	Cacheline size (64 B)	direct-mapped,	
	Latency: ~1x	Latency: ~2x		stochastic [20]	
TDC [38]	In-package traffic: 64 B	In-package traffic: 0 B	On every miss	Hardware managed,	No
	Latency: ~1x	Latency: ~1x	Footprint size [28]	fully-associative,	
	TLB coherence	TLB coherence		FIFO	
HMA [44]	In-package traffic: 64 B	In-package traffic: 0 B	Software managed, high replacement cost		Yes
	Latency: ~1x	Latency: ~1x			
Banshee	In-package traffic: 64 B	In-package traffic: 0 B	Only for hot pages	Hardware managed,	Yes
(This work)	Latency: ~1x	Latency: ~1x	Page size (4 KB)	set-associative,	
				frequency based	

Yu+, "Banshee: Bandwidth-Efficient DRAM Caching via Software/Hardware Cooperation," MICRO 2017.

#### Banshee [MICRO 2017]

- Tracks presence in cache using TLB and Page Table
  - No tag store needed for DRAM cache
  - Enabled by a new lightweight lazy TLB coherence protocol
- New bandwidth-aware frequency-based replacement policy



#### More on Banshee

 Xiangyao Yu, Christopher J. Hughes, Nadathur Satish, Onur Mutlu, and Srinivas Devadas,
 "Banshee: Bandwidth-Efficient DRAM Caching via Software/Hardware Cooperation"
 Proceedings of the 50th International Symposium on Microarchitecture (MICRO), Boston, MA, USA, October 2017.

# Banshee: Bandwidth-Efficient DRAM Caching via Software/Hardware Cooperation

Xiangyao Yu<sup>1</sup> Christopher J. Hughes<sup>2</sup> Nadathur Satish<sup>2</sup> Onur Mutlu<sup>3</sup> Srinivas Devadas<sup>1</sup> <sup>1</sup>MIT <sup>2</sup>Intel Labs <sup>3</sup>ETH Zürich

#### SAFARI

#### Other Opportunities with Emerging Technologies

#### Merging of memory and storage

- e.g., a single interface to manage all data
- New applications
  - e.g., ultra-fast checkpoint and restore
- More robust system design
  - e.g., reducing data loss
- Processing tightly-coupled with memory
   e.g., enabling efficient search and filtering

# **TWO-LEVEL STORAGE MODEL**



# **TWO-LEVEL STORAGE MODEL**



Non-volatile memories combine characteristics of memory and storage

# Two-Level Memory/Storage Model

- The traditional two-level storage model is a bottleneck with NVM
  - Volatile data in memory  $\rightarrow$  a load/store interface
  - **Persistent** data in storage  $\rightarrow$  a **file system** interface
  - Problem: Operating system (OS) and file system (FS) code to locate, translate, buffer data become performance and energy bottlenecks with fast NVM stores



## Unified Memory and Storage with NVM

- Goal: Unify memory and storage management in a single unit to eliminate wasted work to locate, transfer, and translate data
  - Improves both energy and performance
  - Simplifies programming model as well



93

#### The Persistent Memory Manager (PMM)



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

#### Performance Benefits of a Single-Level Store



**SAFARI** Meza+, "A Case for Efficient Hardware-Software Cooperative Management of 95 Storage and Memory," WEED 2013.

#### Energy Benefits of a Single-Level Store



**SAFARI** Meza+, "A Case for Efficient Hardware-Software Cooperative Management of 96 Storage and Memory," WEED 2013.

#### On Persistent Memory Benefits & Challenges

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

#### 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> Onur Mutlu<sup>\*</sup> \*Carnegie Mellon University <sup>†</sup>Pennsylvania State University <sup>‡</sup>Intel Labs <sup>§</sup>AMD Research

#### SAFARI

## Challenge and Opportunity

# Combined Memory & Storage

## Challenge and Opportunity

# A Unified Interface to All Data

## One Key Challenge in Persistent Memory

- How to ensure consistency of system/data if all memory is persistent?
- Two extremes
  - Programmer transparent: Let the system handle it
  - Programmer only: Let the programmer handle it
- Many alternatives in-between...

# **CHALLENGE: CRASH CONSISTENCY**



#### **Persistent Memory System**

# System crash can result in permanent data corruption in NVM

## **CRASH CONSISTENCY PROBLEM**

Example: Add a node to a linked list



System crash can result in inconsistent memory state

# **CURRENT SOLUTIONS**

#### **Explicit interfaces to manage consistency**

- NV-Heaps [ASPLOS'11], BPFS [SOSP'09], Mnemosyne [ASPLOS'11]

# AtomicBegin { Insert a new node; } AtomicEnd;

#### Limits adoption of NVM Have to rewrite code with clear partition between volatile and non-volatile data

# **Burden on the programmers**

# **OUR APPROACH: ThyNVM**

# Goal: Software transparent consistency in persistent memory systems

# **ThyNVM: Summary**

A new hardware-based checkpointing mechanism

- Checkpoints at multiple granularities to reduce both checkpointing latency and metadata overhead
- Overlaps checkpointing and execution to reduce checkpointing latency
- Adapts to DRAM and NVM characteristics

Performs within **4.9%** of an *idealized DRAM* with zero cost consistency

### More About ThyNVM

 Jinglei Ren, Jishen Zhao, Samira Khan, Jongmoo Choi, Yongwei Wu, and <u>Onur Mutlu</u>,
 "ThyNVM: Enabling Software-Transparent Crash Consistency in Persistent Memory Systems"
 Proceedings of the <u>48th International Symposium on</u> <u>Microarchitecture</u> (MICRO), Waikiki, Hawaii, USA, December 2015.
 [Slides (pptx) (pdf)] [Lightning Session Slides (pptx) (pdf)] [Poster (pptx) (pdf)]
 [Source Code]

#### ThyNVM: Enabling Software-Transparent Crash Consistency in Persistent Memory Systems

Jinglei Ren<sup>\*†</sup> Jishen Zhao<sup>‡</sup> Samira Khan<sup>†</sup> Jongmoo Choi<sup>+†</sup> Yongwei Wu<sup>\*</sup> Onur Mutlu<sup>†</sup>

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

<sup>‡</sup>University of California, Santa Cruz <sup>′</sup>University of Virginia <sup>+</sup>Dankook University

#### Another Key Challenge in Persistent Memory

# Programming Ease to Exploit Persistence



# Tools/Libraries to Help Programmers

 Himanshu Chauhan, Irina Calciu, Vijay Chidambaram, Eric Schkufza, Onur Mutlu, and Pratap Subrahmanyam,
 "NVMove: Helping Programmers Move to Byte-Based Persistence"

Proceedings of the <u>4th Workshop on Interactions of NVM/Flash</u> <u>with Operating Systems and Workloads</u> (**INFLOW**), Savannah, GA, USA, November 2016. [<u>Slides (pptx) (pdf)</u>]

#### **NVMOVE: Helping Programmers Move to Byte-Based Persistence**

Himanshu Chauhan \* UT Austin Irina Calciu VMware Research Group Vijay Chidambaram UT Austin

Eric Schkufza VMware Research Group Onur Mutlu ETH Zürich Pratap Subrahmanyam VMware

#### SAFARI


- Major Trends Affecting Main Memory
- The Memory Scaling Problem
- Two Complementary Solution Directions
  - New Memory (DRAM) Architectures
  - Enabling Emerging (NVM) Technologies
- Conclusion

#### SAFARI

#### The Future of Emerging Technologies is Bright

- Regardless of challenges
  - in underlying technology and overlying problems/requirements

Can enable:

- Orders of magnitude improvements

- New applications and computing systems

Problem	
Aigorithm	
Program/Language	
System Software	
SW/HW Interface	
Micro-architecture	
Logic	
Devices	
Electrons	

Yet, we have to

- Think across the stack
- Design enabling systems

## If In Doubt, Refer to Flash Memory

A very "doubtful" emerging technology
 for at least two decades

Proceedings of the IEEE, Sept. 2017 Error Characterization, Mitigation, and Recovery in Flash-Memory-Based Solid-State Drives

By YU CAI, SAUGATA GHOSE, ERICH F. HARATSCH, YIXIN LUO, AND ONUR MUTLU

**ABSTRACT** | NAND flash memory is ubiquitous in everyday life today because its capacity has continuously increased and

**KEYWORDS** | Data storage systems; error recovery; fault tolerance; flash memory; reliability; solid-state drives

SAFARI

https://arxiv.org/pdf/1706.08642

## Opportunities and Challenges of Emerging Memory Technologies

Onur Mutlu

omutlu@gmail.com

https://people.inf.ethz.ch/omutlu

September 11, 2017 ARM Research Summit



**ETH** zürich

### Overview Paper on Flash Reliability

- Yu Cai, Saugata Ghose, Erich F. Haratsch, Yixin Luo, and Onur Mutlu,
  - "Error Characterization, Mitigation, and Recovery in Flash Memory Based Solid State Drives"

to appear in <u>Proceedings of the IEEE</u>, 2017.

#### Error Characterization, Mitigation, and Recovery in Flash Memory Based Solid State Drives

Yu Cai, Saugata Ghose, Erich F. Haratsch, Yixin Luo, and Onur Mutlu

## Solution 2: Emerging Memory Technologies

- Some emerging resistive memory technologies seem more scalable than DRAM (and they are non-volatile)
- Example: Phase Change Memory
  - Expected to scale to 9nm (2022 [ITRS])
  - Expected to be denser than DRAM: can store multiple bits/cell
- But, emerging technologies have shortcomings as well
  Can they be enabled to replace/augment/surpass DRAM?
- Lee+, "Architecting Phase Change Memory as a Scalable DRAM Alternative," ISCA'09, CACM'10, IEEE Micro'10.
- Meza+, "Enabling Efficient and Scalable Hybrid Memories," IEEE Comp. Arch. Letters 2012.
- Yoon, Meza+, "Row Buffer Locality Aware Caching Policies for Hybrid Memories," ICCD 2012.
- Kultursay+, "Evaluating STT-RAM as an Energy-Efficient Main Memory Alternative," ISPASS 2013.
- Meza+, "A Case for Efficient Hardware-Software Cooperative Management of Storage and Memory," WEED 2013.
- Lu+, "Loose Ordering Consistency for Persistent Memory," ICCD 2014.
- Zhao+, "FIRM: Fair and High-Performance Memory Control for Persistent Memory Systems," MICRO 2014.
- Yoon, Meza+, "Efficient Data Mapping and Buffering Techniques for Multi-Level Cell Phase-Change Memories," TACO 2014.
- Ren+, "ThyNVM: Enabling Software-Transparent Crash Consistency in Persistent Memory Systems," MICRO 2015.
- Chauhan+, "NVMove: Helping Programmers Move to Byte-Based Persistence," INFLOW 2016.
- Li+, "Utility-Based Hybrid Memory Management," CLUSTER 2017.
- Yu+, "Banshee: Bandwidth-Efficient DRAM Caching via Software/Hardware Cooperation," MICRO 2017.

#### **SAFARI**

## Combination: Hybrid Memory Systems



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

Meza+, "Enabling Efficient and Scalable Hybrid Memories," IEEE Comp. Arch. Letters, 2012. Yoon, Meza et al., "Row Buffer Locality Aware Caching Policies for Hybrid Memories," ICCD 2012 Best Paper Award.

#### SAFARI

### Exploiting Memory Error Tolerance with Hybrid Memory Systems



On Microsoft's Web Search workload Reduces server hardware cost by 4.7 % Achieves single server availability target of 99.90 % Heterogeneous-Reliability Memory [DSN 2014]

## More on Heterogeneous Reliability Memory

Yixin Luo, Sriram Govindan, Bikash Sharma, Mark Santaniello, Justin Meza, Aman Kansal, Jie Liu, Badriddine Khessib, Kushagra Vaid, and Onur Mutlu,
 <u>"Characterizing Application Memory Error Vulnerability to Optimize</u>
 <u>Data Center Cost via Heterogeneous-Reliability Memory"</u>
 *Proceedings of the <u>44th Annual IEEE/IFIP International Conference on</u>
 <u>Dependable Systems and Networks</u> (DSN), Atlanta, GA, June 2014. [Summary]
 [Slides (pptx) (pdf)] [Coverage on ZDNet]* 

#### Characterizing Application Memory Error Vulnerability to Optimize Datacenter Cost via Heterogeneous-Reliability Memory

Yixin Luo Sriram Govindan<sup>\*</sup> Bikash Sharma<sup>\*</sup> Mark Santaniello<sup>\*</sup> Justin Meza Aman Kansal<sup>\*</sup> Jie Liu<sup>\*</sup> Badriddine Khessib<sup>\*</sup> Kushagra Vaid<sup>\*</sup> Onur Mutlu Carnegie Mellon University, yixinluo@cs.cmu.edu, {meza, onur}@cmu.edu \*Microsoft Corporation, {srgovin, bsharma, marksan, kansal, jie.liu, bkhessib, kvaid}@microsoft.com

# Providing the Best of Multiple Metrics with

# **Multiple Memory Technologies**



Heterogeneous, Configurable, Programmable **Memory Systems**