18-447 Computer Architecture Lecture 34: Emerging Memory Technologies

Prof. Onur Mutlu Carnegie Mellon University Spring 2014, 5/2/2014

Lab 5 Statistics

- MAX 100
 MIN 67.54
 MEDIAN 100
 MEAN 93.30
- STD 10.96

Lab 5 Grade Distribution



Lab 5 Extra Credit (Cache Sweep)

1. Aaron Reyes, Bailey Forrest, Max Regan, Mengzhe Li, Xiang Lin, John Greth

2. Chang Sheng Loh, Fazle Sadi, Jacquelyn Harris, Jeremie Kim, Nicolas Mellis

- 3. Erik Pintar, Albert Cho
- 4. Teng Fei Liao

5. Doci Mou, Jonathan Leung

16 Extra Credit Winners for a Lab: A Record for 447!

Lab 6 Statistics

- MAX 100
 MIN 67.54
 MEDIAN 100
 MEAN 93.30
- STD 10.96

Lab 6 Grade Distribution



Lab 6 Extra Credit

- Albert Cho (best performance with prefetcher)
- Fazle Sadi (second best performance with prefetcher)
- Bailey Forrest (stride prefetcher)
- Are we missing anyone else?

Overall Extra Credit Champions

- All labs (5/5)
 - Bailey Forrest
- All minus one (4/5)
 Albert Cho
- All minus two (3/5)
 - John Greth
 - Jeremie Kim
 - Teng Fei Liao
 - Xiang Lin
 - Chang Sheng Loh

Final Exam: May 6

- May 6, 8:30-11:30am, Hamerschlag Hall B103
- Comprehensive (over all topics in course)
- Three cheat sheets allowed
- We might have a review session
- Remember this is 25% of your grade
 - □ I will take into account your improvement over the course
 - □ Know all concepts, especially the previous midterm concepts
 - Same advice as before for Midterms I and II

A Note on 742, Research, Jobs

- I am teaching Parallel Computer Architecture next semester (Fall 2014)
 - Deep dive into many topics we covered
 - And, many topics we did not cover
 - Research oriented with an open-ended research project
 - Cutting edge research and topics in HW/SW interface
- If you are enjoying 447 and are doing well, you can take it
 - \rightarrow no need to have taken 640/740
 - \rightarrow talk with me
- If you are excited about Computer Architecture research or looking for a job/internship in this area
 - \rightarrow talk with me

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



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
- Major Solution Directions
- Requirements from an Ideal Main Memory System
- Opportunity: Emerging Memory Technologies

Major Trends Affecting Main Memory (I)

Need for main memory capacity and bandwidth increasing

Main memory energy/power is a key system design concern

DRAM technology scaling is ending

Demand for Memory Capacity

■ More cores → More concurrency → Larger working set



AMD Barcelona: 4 cores



IBM Power7: 8 cores



Intel SCC: 48 cores

Emerging applications are data-intensive

Many applications/virtual machines (will) share main memory

- Cloud computing/servers: Consolidation to improve efficiency
- GP-GPUs: Many threads from multiple parallel applications
- Mobile: Interactive + non-interactive consolidation

The Memory Capacity Gap

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



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 X nm
- 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
- Major Solution Directions
- Requirements from an Ideal Main Memory System
- Opportunity: Emerging Memory Technologies

Solutions to the DRAM Scaling Problem

- Two potential solutions
 - Tolerate DRAM (by taking a fresh look at it)
 - Enable emerging memory technologies to eliminate/minimize DRAM
- Do both
 - Hybrid memory systems

Solution 1: Tolerate DRAM

- Overcome DRAM shortcomings with
 - System-DRAM co-design
 - Novel DRAM architectures, interface, functions
 - Better waste management (efficient utilization)
- Key issues to tackle
 - Reduce energy
 - Enable reliability at low cost
 - Improve bandwidth and latency
 - Reduce waste

Solution 1: Tolerate DRAM

- Liu, Jaiyen, Veras, Mutlu, "RAIDR: Retention-Aware Intelligent DRAM Refresh," ISCA 2012.
- Kim, Seshadri, Lee+, "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.

Tolerating DRAM: Example Techniques

Retention-Aware DRAM Refresh: Reducing Refresh Impact

Refresh Access Parallelization: Reducing Refresh Impact

Tiered-Latency DRAM: Reducing DRAM Latency

RowClone: Accelerating Page Copy and Initialization

Subarray-Level Parallelism: Reducing Bank Conflict Impact

Linearly Compressed Pages: Efficient Memory Compression

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, Ipek, Mutlu, Burger, "Architecting Phase Change Memory as a Scalable DRAM Alternative," ISCA 2009, CACM 2010, Top Picks 2010.
- Meza, Chang, Yoon, Mutlu, Ranganathan, "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.
- 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.

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.

An Orthogonal Issue: Memory Interference



Cores' interfere with each other when accessing shared main memory

An Orthogonal Issue: Memory Interference

- Problem: Memory interference between cores is uncontrolled
 - \rightarrow unfairness, starvation, low performance
 - \rightarrow uncontrollable, unpredictable, vulnerable system
- Solution: QoS-Aware Memory Systems
 - Hardware designed to provide a configurable fairness substrate
 - Application-aware memory scheduling, partitioning, throttling
 - Software designed to configure the resources to satisfy different QoS goals
- QoS-aware memory controllers and interconnects can provide predictable performance and higher efficiency

Designing QoS-Aware Memory Systems: Approaches

- Smart resources: Design each shared resource to have a configurable interference control/reduction mechanism
 - QoS-aware memory controllers [Mutlu+ MICRO'07] [Moscibroda+, Usenix Security'07] [Mutlu+ ISCA'08, Top Picks'09] [Kim+ HPCA'10] [Kim+ MICRO'10, Top Picks'11] [Ebrahimi+ ISCA'11, MICRO'11] [Ausavarungnirun+, ISCA'12][Subramanian+, HPCA'13] [Kim+, RTAS'14]
 - QoS-aware interconnects [Das+ MICRO'09, ISCA'10, Top Picks '11] [Grot+ MICRO'09, ISCA'11, Top Picks '12]
 - QoS-aware caches
- Dumb resources: Keep each resource free-for-all, but reduce/control interference by injection control or data mapping
 - Source throttling to control access to memory system [Ebrahimi+ ASPLOS'10, ISCA'11, TOCS'12] [Ebrahimi+ MICRO'09] [Nychis+ HotNets'10] [Nychis+ SIGCOMM'12]
 - □ QoS-aware data mapping to memory controllers [Muralidhara+ MICRO'11]
 - QoS-aware thread scheduling to cores [Das+ HPCA'13]



- Major Trends Affecting Main Memory
- Major Solution Directions
- Requirements from an Ideal Main Memory System
- Opportunity: Emerging Memory Technologies

Requirements from an Ideal Memory System

Traditional

- Enough capacity
- Low cost
- High system performance (high bandwidth, low latency)

New

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

Requirements from an Ideal Memory System

Traditional

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

New

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

Emerging, resistive memory technologies (NVM) can help



- Major Trends Affecting Main Memory
- Requirements from an Ideal Main Memory System
- Opportunity: Emerging Memory Technologies
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
 - Other Opportunities with Emerging Technologies

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/RRAM/ReRAM
 - 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 Tcryst
 - 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.								
Parameter*	Published prototype									
	Horri ⁶	Ahn ¹²	Bedeschi ¹³	Oh14	Pellizer ¹⁵	Chen ⁵	Kang ¹⁶	Bedeschi ⁹	Lee ¹⁰	Lee ²
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 ²)	••	0.290	0.290	••	0.097	60 nm ²	0.166	0.097	0.047	0.065 to 0.097
Cell size, F ²		20.1	9.0	••	12.0		16.6	12.0	5.8	9.0 to 12.0
Access device	**	**	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 ⁷	10 ⁹	10 ⁶	••	10 ⁸	104		10 ⁵	10 ⁵	10 ⁸

* BJT: bipolar junction transistor; FET: field-effect transistor; GST: Ge₂Sb₂Te₅; MLC: multilevel cells; N-d: nitrogen doped. ** This information is not available in the publication cited.

6

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>10⁸ writes per cell</u>

 \sim 10⁻⁸x DRAM, 10³x NAND Flash

Cell Size

9-12F² 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 ~10⁸ 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: Research 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 DCM shortcomings, explait DCM advantages
 - 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

SAFARI

- 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)

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
 - Other Opportunities with Emerging Technologies

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





Latency

50ns Rd, 150ns Wr

 \triangleright 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?



- 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
 - Other Opportunities with Emerging Technologies

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

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 huge (DRAM) cache at low cost?
- Two solutions:

Locality-aware data placement [Yoon+, ICCD 2012]

Cheap tag stores and dynamic granularity [Meza+, IEEE CAL 2012]

SAFARI

DRAM vs. PCM: An Observation

- Row buffers are the same in DRAM and PCM
- Row buffer hit latency same in DRAM and PCM
- Row buffer miss latency small in DRAM, large in PCM



- Accessing the row buffer in PCM is fast
- What incurs high latency is the PCM array access \rightarrow avoid this

Row-Locality-Aware Data Placement

- Idea: Cache in DRAM only those rows that
 - □ Frequently cause row buffer conflicts → because row-conflict latency is smaller in DRAM
 - □ Are reused many times → to reduce cache pollution and bandwidth waste
- Simplified rule of thumb:
 - Streaming accesses: Better to place in PCM
 - Other accesses (with some reuse): Better to place in DRAM

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

Row-Locality-Aware Data Placement: Results





Hybrid vs. All-PCM/DRAM

■ 16GB PCM ■ RBLA-Dyn ■ 16GB 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
 - Other Opportunities with Emerging Technologies

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

Coordinated Memory and Storage with NVM (I)

- 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



Coordinated Memory and Storage with NVM (II)

- 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



71

The Persistent Memory Manager (PMM)

- Exposes a load/store interface to access persistent data
 - □ Applications can directly access persistent memory → no conversion, translation, location overhead for persistent data
- Manages data placement, location, persistence, security
 - To get the best of multiple forms of storage
- Manages metadata storage and retrieval
 - This can lead to overheads that need to be managed
- Exposes hooks and interfaces for system software
 - To enable better data placement and management decisions
- Meza+, "A Case for Efficient Hardware-Software Cooperative Management of Storage and Memory," WEED 2013.
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



Energy Benefits of a Single-Level Store



Enabling and Exploiting NVM: Issues

- Many issues and ideas from technology layer to algorithms layer
- Enabling NVM and hybrid memory
 - How to tolerate errors?
 - How to enable secure operation?
 - How to tolerate performance and power shortcomings?
 - How to minimize cost?
- Exploiting emerging tecnologies
 - 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 \rightarrow 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
- 2. Non-volatility \rightarrow Data persists in memory after powerdown
 - → Easy retrieval of privileged or private information
 Efficient encryption/decryption of whole main memory
 Hybrid memory system management
- Multiple bits per cell → Information leakage (via side channel)
 System design to hide side channel information



- 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
 - Other Opportunities with Emerging Technologies
- Summary

Summary

- Key trends affecting main memory
 - End of DRAM scaling (cost, capacity, efficiency)
 - Need for high capacity
 - Need for energy efficiency
- Emerging NVM technologies can help
 - PCM more scalable than DRAM and non-volatile
 - But, it has critical shortcomings: latency, active energy, endurance
- We need to enable promising NVM technologies by overcoming their shortcomings
- Many exciting opportunities to reinvent main memory at all layers of computing stack