SAFARI Group: Research Introduction

Onur Mutlu omutlu@gmail.com https://people.inf.ethz.ch/omutlu

January 31, 2018





Brief Self Introduction

Onur Mutlu



- Professor @ ETH Zurich CS, since September'15, started May'16
- Strecker Professor @ Carnegie Mellon University ECE (CS), 2009-2016, 2016-...
- PhD from UT-Austin, worked @ Google, VMware, Microsoft Research, Intel, AMD
- https://people.inf.ethz.ch/omutlu/
- omutlu@gmail.com (Best way to reach me)
- Publications: <u>https://people.inf.ethz.ch/omutlu/projects.htm</u>
- Research, Education, Consulting in
 - Computer architecture and systems, bioinformatics
 - Memory and storage systems, emerging technologies
 - Many-core systems, heterogeneous systems, core design
 - Interconnects
 - Hardware/software interaction and co-design (PL, OS, Architecture)
 - Predictable and QoS-aware systems
 - Hardware fault tolerance and security
 - Algorithms and architectures for genome analysis

Brief Introduction: SAFARI Group

- 27+ researchers
- Multiple locations: ETH and CMU
- ETH researchers
 - 4 post-PhD (postdocs, visiting faculty)
 - 6 PhD students
 - 2 post-masters students
 - 2+ masters students
- CMU researchers
 - 2 post-PhD (postdoc, research faculty)
 - o 7 PhD student
 - 3 masters

Current Research Focus Areas

<u>**Research Focus:**</u> Computer architecture, HW/SW, bioinformatics

Memory and storage (DRAM, flash, emerging), interconnects

- Heterogeneous & parallel systems, GPUs, systems for data analytics
- System/architecture interaction, new execution models, new interfaces
- Energy efficiency, fault tolerance, hardware security, performance
- Genome sequence analysis & assembly algorithms and architectures
- Biologically inspired systems & system design for bio/medicine



General Purpose GPUs

Four Key Current Directions

Fundamentally Secure/Reliable/Safe Architectures

Fundamentally Energy-Efficient Architectures
 Memory-centric (Data-centric) Architectures

Fundamentally Low-Latency Architectures

Architectures for Genomics, Medicine, Health

In-Memory DNA Sequence Analysis

 Jeremie S. Kim, Damla Senol Cali, Hongyi Xin, Donghyuk Lee, Saugata Ghose, Mohammed Alser, Hasan Hassan, Oguz Ergin, Can Alkan, Onur Mutlu, "GRIM-Filter: Fast Seed Location Filtering in DNA Read Mapping Using Processing-in-Memory Technologies" to appear in <u>BMC Genomics</u>, 2018. to also appear in Proceedings of the <u>16th Asia Pacific Bioinformatics</u> Conference (APBC), Yokohama, Japan, January 2018.

arxiv.org Version (pdf)

GRIM-Filter: Fast Seed Location Filtering in DNA Read Mapping Using Processing-in-Memory Technologies

Jeremie S. Kim^{1,6*}, Damla Senol Cali¹, Hongyi Xin², Donghyuk Lee³, Saugata Ghose¹, Mohammed Alser⁴, Hasan Hassan⁶, Oguz Ergin⁵, Can Alkan^{*4}, and Onur Mutlu^{*6,1}

Memory/Storage Issues We Are Tackling

- (Enable and Exploit) In-Memory Computation
- (Reduce) Latency
- (Reduce) Energy
- Improve) Reliability and Security
- (Make Sense of & Take Advantage of) Heterogeneity
- (Enable and Use) Persistence
- (Enable and Exploit) QoS and Predictability
- (Develop New and Reliable) Infrastructure

Topics Requested (by Email)

1. PIM enabled instructions

2. DRAM Reliability and Performance

3. Processing in Memory (Near-Data Computation)

Three Key Systems Trends

1. Data access is a major bottleneck

Applications are increasingly data hungry

2. Energy consumption is a key limiter

3. Data movement energy dominates compute

Especially true for off-chip to on-chip movement

The Need for More Memory Performance



In-memory Databases

[Mao+, EuroSys'12; Clapp+ (**Intel**), IISWC'15]



In-Memory Data Analytics

[Clapp+ (**Intel**), IISWC'15; Awan+, BDCloud'15]



Graph/Tree Processing [Xu+, IISWC'12; Umuroglu+, FPL'15]



Datacenter Workloads [Kanev+ (**Google**), ISCA'15]

The Performance Perspective (1996-2005)

"It's the Memory, Stupid!" (Richard Sites, MPR, 1996)



Mutlu+, "Runahead Execution: An Alternative to Very Large Instruction Windows for Out-of-Order Processors," HPCA 2003.

The Performance Perspective (Today)

All of Google's Data Center Workloads (2015):



The Performance Perspective (Today)

All of Google's Data Center Workloads (2015):



Figure 11: Half of cycles are spent stalled on caches.

The Performance Perspective

 Onur Mutlu, Jared Stark, Chris Wilkerson, and Yale N. Patt, "Runahead Execution: An Alternative to Very Large Instruction Windows for Out-of-order Processors" Proceedings of the <u>9th International Symposium on High-Performance</u> <u>Computer Architecture</u> (HPCA), pages 129-140, Anaheim, CA, February 2003. <u>Slides (pdf)</u>

Runahead Execution: An Alternative to Very Large Instruction Windows for Out-of-order Processors

Onur Mutlu § Jared Stark † Chris Wilkerson ‡ Yale N. Patt §

§ECE Department The University of Texas at Austin {onur,patt}@ece.utexas.edu †Microprocessor Research Intel Labs jared.w.stark@intel.com

Desktop Platforms Group Intel Corporation chris.wilkerson@intel.com

The Energy Perspective



Data Movement vs. Computation Energy



A memory access consumes ~1000X the energy of a complex addition

Data Movement vs. Computation Energy

Data movement is a major system energy bottleneck

- Comprises 41% of mobile system energy during web browsing [2]
- Costs ~115 times as much energy as an ADD operation [1, 2]



[1]: Reducing data Movement Energy via Online Data Clustering and Encoding (MICRO'16)

[2]: Quantifying the energy cost of data movement for emerging smart phone workloads on mobile platforms (IISWC'14)

Data access is the major performance and energy bottleneck

Our current design principles cause great energy waste (and great performance loss)

Processing of data is performed far away from the data

A Computing System

- Three key components
- Computation
- Communication
- Storage/memory



Burks, Goldstein, von Neumann, "Preliminary discussion of the logical design of an electronic computing instrument," 1946.

Computing System



Image source: https://lbsitbytes2010.wordpress.com/2013/03/29/john-von-neumann-roll-no-15/

A Computing System

- Three key components
- Computation
- Communication
- Storage/memory



Burks, Goldstein, von Neumann, "Preliminary discussion of the logical design of an electronic computing instrument," 1946.

Computing System



Image source: https://lbsitbytes2010.wordpress.com/2013/03/29/john-von-neumann-roll-no-15/

Today's Computing Systems

- Are overwhelmingly processor centric
- All data processed in the processor \rightarrow at great system cost
- Processor is heavily optimized and is considered the master
- Data storage units are dumb and are largely unoptimized (except for some that are on the processor die)



Yet ...

"It's the Memory, Stupid!" (Richard Sites, MPR, 1996)



Mutlu+, "Runahead Execution: An Alternative to Very Large Instruction Windows for Out-of-Order Processors," HPCA 2003.

Perils of Processor-Centric Design

Grossly-imbalanced systems

- Processing done only in **one place**
- Everything else just stores and moves data: data moves a lot
- \rightarrow Energy inefficient
- \rightarrow Low performance
- \rightarrow Complex
- Overly complex and bloated processor (and accelerators)
 - To tolerate data access from memory
 - Complex hierarchies and mechanisms
 - \rightarrow Energy inefficient
 - \rightarrow Low performance
 - \rightarrow Complex

Perils of Processor-Centric Design



Most of the system is dedicated to storing and moving data

We Do Not Want to Move Data!



A memory access consumes ~1000X the energy of a complex addition

We Need A Paradigm Shift To ...

Enable computation with minimal data movement

Compute where it makes sense (where data resides)

Make computing architectures more data-centric

Goal: Processing Inside Memory



Why In-Memory Computation Today?



- Pull from Systems and Applications
 - Data access is a major system and application bottleneck
 - Systems are energy limited
 - Data movement much more energy-hungry than computation

Processing in Memory: Two Approaches

Minimally changing memory chips
 Exploiting 3D-stacked memory

Opportunity: 3D-Stacked Logic+Memory





DRAM Landscape (circa 2015)

Segment	DRAM Standards & Architectures
Commodity	DDR3 (2007) [14]; DDR4 (2012) [18]
Low-Power	LPDDR3 (2012) [17]; LPDDR4 (2014) [20]
Graphics	GDDR5 (2009) [15]
Performance	eDRAM [28], [32]; RLDRAM3 (2011) [29]
3D-Stacked	WIO (2011) [16]; WIO2 (2014) [21]; MCDRAM (2015) [13]; HBM (2013) [19]; HMC1.0 (2013) [10]; HMC1.1 (2014) [11]
Academic	SBA/SSA (2010) [38]; Staged Reads (2012) [8]; RAIDR (2012) [27]; SALP (2012) [24]; TL-DRAM (2013) [26]; RowClone (2013) [37]; Half-DRAM (2014) [39]; Row-Buffer Decoupling (2014) [33]; SARP (2014) [6]; AL-DRAM (2015) [25]

Table 1. Landscape of DRAM-based memory

Kim+, "Ramulator: A Flexible and Extensible DRAM Simulator", IEEE CAL 2015.

Two Key Questions in 3D Stacked PIM

- What is the minimal processing-in-memory support we can provide ?
 - without changing the system significantly
 - while achieving significant benefits of processing in 3Dstacked memory
- How can we accelerate important applications if we use 3D-stacked memory as a coarse-grained accelerator?
 - what is the architecture and programming model?
 - what are the mechanisms for acceleration?

PEI: PIM-Enabled Instructions (Ideas)

- Goal: Develop mechanisms to get the most out of near-data processing with minimal cost, minimal changes to the system, no changes to the programming model
- Key Idea 1: Expose each PIM operation as a cache-coherent, virtually-addressed host processor instruction (called PEI) that operates on only a single cache block
 - e.g., __pim_add(&w.next_rank, value) \rightarrow pim.add r1, (r2)
 - No changes sequential execution/programming model
 - No changes to virtual memory
 - Minimal changes to cache coherence
 - No need for data mapping: Each PEI restricted to a single memory module
- Key Idea 2: Dynamically decide where to execute a PEI (i.e., the host processor or PIM accelerator) based on simple locality characteristics and simple hardware predictors
 - Execute each operation at the location that provides the best performance

Simple PIM Operations as ISA Extensions (II)

```
for (v: graph.vertices) {
value = weight * v.rank;
for (w: v.successors) {
  w.next rank += value;
    Host Processor
      w.next rank
                         64 bytes in
                        64 bytes out
```



Main Memory

Conventional Architecture

Simple PIM Operations as ISA Extensions (III)



In-Memory Addition
Always Executing in Memory? Not A Good Idea





- Executed either in memory or in the processor: dynamic decision
 - Low-cost locality monitoring for a single instruction
- Cache-coherent, virtually-addressed, single cache block only
- Atomic between different PEIs
- Not atomic with normal instructions (use pfence for ordering)
 SAFARI

- Key to practicality: single-cache-block restriction
 - Each PEI can access *at most one last-level cache block*
 - Similar restrictions exist in atomic instructions
- Benefits
 - **Localization**: each PEI is bounded to one memory module
 - Interoperability: easier support for cache coherence and virtual memory
 - Simplified locality monitoring: data locality of PEIs can be identified simply by the cache control logic

Example PEI Microarchitecture



Example PEI uArchitecture

PEI: Initial Evaluation Results

Initial evaluations with 10 emerging data-intensive workloads

- Large-scale graph processing
- In-memory data analytics
- Machine learning and data mining
- Three input sets (small, medium, large) for each workload to analyze the impact of data locality

Table 2: Baseline Simulation Configuration

Component	Configuration
Core	16 out-of-order cores, 4 GHz, 4-issue
L1 I/D-Cache	Private, 32 KB, 4/8-way, 64 B blocks, 16 MSHRs
L2 Cache	Private, 256 KB, 8-way, 64 B blocks, 16 MSHRs
L3 Cache	Shared, 16 MB, 16-way, 64 B blocks, 64 MSHRs
On-Chip Network	Crossbar, 2 GHz, 144-bit links
Main Memory	32 GB, 8 HMCs, daisy-chain (80 GB/s full-duplex)
HMC	4GB, 16 vaults, 256 DRAM banks [20]
– DRAM	FR-FCFS, $tCL = tRCD = tRP = 13.75 \text{ ns} [27]$
 Vertical Links 	64 TSVs per vault with 2 Gb/s signaling rate [23]

Pin-based cycle-level x86-64 simulation

Performance Improvement and Energy Reduction:

- 47% average speedup with large input data sets
- 32% speedup with small input data sets
- 25% avg. energy reduction in a single node with large input data sets

Evaluated Data-Intensive Applications

- Ten emerging data-intensive workloads
 - Large-scale graph processing
 - Average teenage follower, BFS, PageRank, single-source shortest path, weakly connected components
 - In-memory data analytics
 - Hash join, histogram, radix partitioning
 - Machine learning and data mining
 - Streamcluster, SVM-RFE
- Three input sets (small, medium, large) for each workload to show the impact of data locality

PEI Performance Delta: Large Data Sets



PEI Performance: Large Data Sets



PEI Performance Delta: Small Data Sets





PEI Performance: Small Data Sets



PEI Performance Delta: Medium Data Sets

(Medium Inputs, Baseline: Host-Only)



PEI Energy Consumption



PEI: Advantages & Disadvantages

Advantages

- + Simple and low cost approach to PIM
- + No changes to programming model, virtual memory
- + Dynamically decides where to execute an instruction
- Disadvantages
 - Does not take full advantage of PIM potential
 - Single cache block restriction is limiting

More on PIM-Enabled Instructions

 Junwhan Ahn, Sungjoo Yoo, Onur Mutlu, and Kiyoung Choi, "PIM-Enabled Instructions: A Low-Overhead, Locality-Aware Processing-in-Memory Architecture" Proceedings of the <u>42nd International Symposium on</u> <u>Computer Architecture</u> (ISCA), Portland, OR, June 2015. [Slides (pdf)] [Lightning Session Slides (pdf)]

PIM-Enabled Instructions: A Low-Overhead, Locality-Aware Processing-in-Memory Architecture

Junwhan Ahn Sungjoo Yoo Onur Mutlu[†] Kiyoung Choi junwhan@snu.ac.kr, sungjoo.yoo@gmail.com, onur@cmu.edu, kchoi@snu.ac.kr Seoul National University [†]Carnegie Mellon University

Two Key Questions in 3D Stacked PIM

- What is the minimal processing-in-memory support we can provide ?
 - without changing the system significantly
 - while achieving significant benefits of processing in 3Dstacked memory
 - How can we accelerate important applications if we use 3D-stacked memory as a coarse-grained accelerator?
 - what is the architecture and programming model?
 - what are the mechanisms for acceleration?

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 System for Graph Processing



Communications In Tesseract (I)

```
for (v: graph.vertices) {
  for (w: v.successors) {
    w.next_rank += weight * v.rank;
                          &w
         V
         w
```

Communications In Tesseract (II)

```
for (v: graph.vertices) {
```

for (w: v.successors) {

w.next_rank += weight * v.rank;



Communications In Tesseract (III)



Remote Function Call (Non-Blocking)

- 1. Send function address & args to the remote core
- 2. Store the incoming message to the message queue
- 3. Flush the message queue when it is full or a synchronization barrier is reached



put(w.id, function() { w.next_rank += value; })

Tesseract System for Graph Processing



Evaluated Systems



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 Performance



Effect of Bandwidth & Programming Model



Tesseract Graph Processing System Energy



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

Tesseract: Advantages & Disadvantages

Advantages

- + Specialized graph processing accelerator using PIM
- + Large system performance and energy benefits
- + Takes advantage of 3D stacking for an important workload

Disadvantages

- Changes a lot in the system
 - New programming model
 - Specialized Tesseract cores for graph processing
- Cost
- Scalability limited by off-chip links or graph partitioning

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[§] Sungjoo Yoo Onur Mutlu[†] Kiyoung Choi junwhan@snu.ac.kr, sungpack.hong@oracle.com, sungjoo.yoo@gmail.com, onur@cmu.edu, kchoi@snu.ac.kr Seoul National University [§]Oracle Labs [†]Carnegie Mellon University

Truly Distributed GPU Processing with PIM?



Accelerating GPU Execution with PIM (I)

 Kevin Hsieh, Eiman Ebrahimi, Gwangsun Kim, Niladrish Chatterjee, Mike O'Connor, Nandita Vijaykumar, Onur Mutlu, and Stephen W. Keckler, <u>"Transparent Offloading and Mapping (TOM): Enabling</u> <u>Programmer-Transparent Near-Data Processing in GPU</u> <u>Systems"</u> *Proceedings of the <u>43rd International Symposium on Computer</u>*

Architecture (ISCA), Seoul, South Korea, June 2016.

[Slides (pptx) (pdf)]

[Lightning Session Slides (pptx) (pdf)]

Transparent Offloading and Mapping (TOM): Enabling Programmer-Transparent Near-Data Processing in GPU Systems

Kevin Hsieh[‡] Eiman Ebrahimi[†] Gwangsun Kim^{*} Niladrish Chatterjee[†] Mike O'Connor[†] Nandita Vijaykumar[‡] Onur Mutlu^{§‡} Stephen W. Keckler[†] [‡]Carnegie Mellon University [†]NVIDIA ^{*}KAIST [§]ETH Zürich

Accelerating GPU Execution with PIM (II)

 Ashutosh Pattnaik, Xulong Tang, Adwait Jog, Onur Kayiran, Asit K. Mishra, Mahmut T. Kandemir, <u>Onur Mutlu</u>, and Chita R. Das, <u>"Scheduling Techniques for GPU Architectures with Processing-</u> <u>In-Memory Capabilities"</u>

Proceedings of the <u>25th International Conference on Parallel</u> <u>Architectures and Compilation Techniques</u> (**PACT**), Haifa, Israel, September 2016.

Scheduling Techniques for GPU Architectures with Processing-In-Memory Capabilities

Ashutosh Pattnaik¹ Xulong Tang¹ Adwait Jog² Onur Kayıran³ Asit K. Mishra⁴ Mahmut T. Kandemir¹ Onur Mutlu^{5,6} Chita R. Das¹ ¹Pennsylvania State University ²College of William and Mary ³Advanced Micro Devices, Inc. ⁴Intel Labs ⁵ETH Zürich ⁶Carnegie Mellon University

Automatic Code and Data Mapping?

 Kevin Hsieh, Eiman Ebrahimi, Gwangsun Kim, Niladrish Chatterjee, Mike O'Connor, Nandita Vijaykumar, Onur Mutlu, and Stephen W. Keckler, <u>"Transparent Offloading and Mapping (TOM): Enabling</u> <u>Programmer-Transparent Near-Data Processing in GPU</u> <u>Systems"</u> *Proceedings of the <u>43rd International Symposium on Computer</u> <u>Architecture</u> (ISCA), Seoul, South Korea, June 2016.*

[Slides (pptx) (pdf)]

[Lightning Session Slides (pptx) (pdf)]

Transparent Offloading and Mapping (TOM): Enabling Programmer-Transparent Near-Data Processing in GPU Systems

Kevin Hsieh[‡] Eiman Ebrahimi[†] Gwangsun Kim^{*} Niladrish Chatterjee[†] Mike O'Connor[†] Nandita Vijaykumar[‡] Onur Mutlu^{§‡} Stephen W. Keckler[†] [‡]Carnegie Mellon University [†]NVIDIA ^{*}KAIST [§]ETH Zürich

Processing in Memory: Two Approaches

Minimally changing memory chips
 Exploiting 3D-stacked memory

Minimally Changing DRAM

- DRAM has great capability to perform bulk data movement and computation internally with small changes
 - Can exploit internal connectivity to move data
 - Can exploit analog computation capability

Examples: RowClone, In-DRAM AND/OR, Gather/Scatter DRAM

- <u>RowClone: Fast and Efficient In-DRAM Copy and Initialization of Bulk Data</u> (Seshadri et al., MICRO 2013)
- □ Fast Bulk Bitwise AND and OR in DRAM (Seshadri et al., IEEE CAL 2015)
- <u>Gather-Scatter DRAM: In-DRAM Address Translation to Improve the Spatial</u> <u>Locality of Non-unit Strided Accesses</u> (Seshadri et al., MICRO 2015)
- "Ambit: In-Memory Accelerator for Bulk Bitwise Operations Using Commodity DRAM Technology" (Seshadri et al., MICRO 2017)

SAFARI

Memory as an Accelerator



Memory similar to a "conventional" accelerator
In-Memory Bulk Operations

- 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

Starting Simple: Data Copy and Initialization

memmove & memcpy: 5% cycles in Google's datacenter [Kanev+ ISCA'15]





---->

Many more

VM Cloning Deduplication

Page Migration

RowClone: In-DRAM Row Copy



RowClone: Latency and Energy Savings



Seshadri et al., "RowClone: Fast and Efficient In-DRAM Copy and Initialization of Bulk Data," MICRO 2013.

More on RowClone

 Vivek Seshadri, Yoongu Kim, Chris Fallin, Donghyuk Lee, Rachata Ausavarungnirun, Gennady Pekhimenko, Yixin Luo, Onur Mutlu, Michael A. Kozuch, Phillip B. Gibbons, and Todd C. Mowry,
 "RowClone: Fast and Energy-Efficient In-DRAM Bulk Data Copy and Initialization" Proceedings of the <u>46th International Symposium on Microarchitecture</u>

(*MICRO*), Davis, CA, December 2013. [Slides (pptx) (pdf)] [Lightning Session Slides (pptx) (pdf)] [Poster (pptx) (pdf)]

RowClone: Fast and Energy-Efficient In-DRAM Bulk Data Copy and Initialization

Vivek Seshadri Yoongu Kim Chris Fallin* Donghyuk Lee vseshadr@cs.cmu.edu yoongukim@cmu.edu cfallin@c1f.net donghyuk1@cmu.edu Rachata Ausavarungnirun Gennady Pekhimenko Yixin Luo rachata@cmu.edu gpekhime@cs.cmu.edu yixinluo@andrew.cmu.edu Onur Mutlu Phillip B. Gibbons† Michael A. Kozuch† Todd C. Mowry onur@cmu.edu phillip.b.gibbons@intel.com michael.a.kozuch@intel.com tcm@cs.cmu.edu Carnegie Mellon University †Intel Pittsburgh

Bulk Bitwise Operations in Workloads



SAFARI

[1] Li and Patel, BitWeaving, SIGMOD 2013[2] Goodwin+, BitFunnel, SIGIR 2017

In-DRAM AND/OR: Triple Row Activation



In-DRAM NOT: Dual Contact Cell



Figure 5: A dual-contact cell connected to both ends of a sense amplifier Idea: Feed the negated value in the sense amplifier into a special row

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

	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)	(↓)	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.

Ambit vs. DDR3: Performance and Energy



82

Performance: Bitmap Index on Ambit



Figure 10: Bitmap index performance. The value above each bar indicates the reduction in execution time due to Ambit.

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

Performance: BitWeaving on Ambit

`select count(*) from T where c1 <= val <= c2'</pre>



Figure 11: Speedup offered by Ambit over baseline CPU with SIMD for BitWeaving

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

More on In-DRAM Bulk AND/OR

 Vivek Seshadri, Kevin Hsieh, Amirali Boroumand, Donghyuk Lee, Michael A. Kozuch, Onur Mutlu, Phillip B. Gibbons, and Todd C. Mowry,
 <u>"Fast Bulk Bitwise AND and OR in DRAM"</u> <u>IEEE Computer Architecture Letters</u> (CAL), April 2015.

Fast Bulk Bitwise AND and OR in DRAM

Vivek Seshadri*, Kevin Hsieh*, Amirali Boroumand*, Donghyuk Lee*, Michael A. Kozuch[†], Onur Mutlu*, Phillip B. Gibbons[†], Todd C. Mowry* *Carnegie Mellon University [†]Intel Pittsburgh

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^{1,5} Donghyuk Lee^{2,5} Thomas Mullins^{3,5} Hasan Hassan⁴ Amirali Boroumand⁵ Jeremie Kim^{4,5} Michael A. Kozuch³ Onur Mutlu^{4,5} Phillip B. Gibbons⁵ Todd C. Mowry⁵

¹Microsoft Research India ²NVIDIA Research ³Intel ⁴ETH Zürich ⁵Carnegie Mellon University

Eliminating the Adoption Barriers

How to Enable Adoption of Processing in Memory

Barriers to Adoption of PIM

1. Functionality of and applications for PIM

- 2. Ease of programming (interfaces and compiler/HW support)
- 3. System support: coherence & virtual memory

4. Runtime systems for adaptive scheduling, data mapping, access/sharing control

5. Infrastructures to assess benefits and feasibility

We Need to Revisit the Entire Stack

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

Key Challenge 1: Code Mapping

• Challenge 1: Which operations should be executed in memory vs. in CPU?



Key Challenge 2: Data Mapping

• Challenge 2: How should data be mapped to different 3D memory stacks?



How to Do the Code and Data Mapping?

- Kevin Hsieh, Eiman Ebrahimi, Gwangsun Kim, Niladrish Chatterjee, Mike O'Connor, Nandita Vijaykumar, Onur Mutlu, and Stephen W. Keckler, <u>"Transparent Offloading and Mapping (TOM): Enabling</u> <u>Programmer-Transparent Near-Data Processing in GPU</u> <u>Systems"</u> *Proceedings of the <u>43rd International Symposium on Computer</u> <u>Architecture</u> (ISCA), Seoul, South Korea, June 2016.*
 - [<u>Slides (pptx) (pdf)</u>]

[Lightning Session Slides (pptx) (pdf)]

Transparent Offloading and Mapping (TOM): Enabling Programmer-Transparent Near-Data Processing in GPU Systems

Kevin Hsieh[‡] Eiman Ebrahimi[†] Gwangsun Kim^{*} Niladrish Chatterjee[†] Mike O'Connor[†] Nandita Vijaykumar[‡] Onur Mutlu^{§‡} Stephen W. Keckler[†] [‡]Carnegie Mellon University [†]NVIDIA ^{*}KAIST [§]ETH Zürich

How to Schedule Code?

 Ashutosh Pattnaik, Xulong Tang, Adwait Jog, Onur Kayiran, Asit K. Mishra, Mahmut T. Kandemir, <u>Onur Mutlu</u>, and Chita R. Das, <u>"Scheduling Techniques for GPU Architectures with Processing-</u> <u>In-Memory Capabilities"</u>

Proceedings of the <u>25th International Conference on Parallel</u> <u>Architectures and Compilation Techniques</u> (**PACT**), Haifa, Israel, September 2016.

Scheduling Techniques for GPU Architectures with Processing-In-Memory Capabilities

Ashutosh Pattnaik¹ Xulong Tang¹ Adwait Jog² Onur Kayıran³ Asit K. Mishra⁴ Mahmut T. Kandemir¹ Onur Mutlu^{5,6} Chita R. Das¹ ¹Pennsylvania State University ²College of William and Mary ³Advanced Micro Devices, Inc. ⁴Intel Labs ⁵ETH Zürich ⁶Carnegie Mellon University

Challenge: Coherence for Hybrid CPU-PIM Apps



How to Maintain Coherence?

 Amirali Boroumand, Saugata Ghose, Minesh Patel, Hasan Hassan, Brandon Lucia, Kevin Hsieh, Krishna T. Malladi, Hongzhong Zheng, and Onur Mutlu,
 "LazyPIM: An Efficient Cache Coherence Mechanism for Processing-in-Memory" IEEE Computer Architecture Letters (CAL), June 2016.

LazyPIM: An Efficient Cache Coherence Mechanism for Processing-in-Memory

Amirali Boroumand[†], Saugata Ghose[†], Minesh Patel[†], Hasan Hassan^{†§}, Brandon Lucia[†], Kevin Hsieh[†], Krishna T. Malladi^{*}, Hongzhong Zheng^{*}, and Onur Mutlu^{‡†} [†]Carnegie Mellon University *Samsung Semiconductor, Inc. [§]TOBB ETÜ [‡]ETH Zürich



How to Support Virtual Memory?

 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[†] Samira Khan[‡] Nandita Vijaykumar[†] Kevin K. Chang[†] Amirali Boroumand[†] Saugata Ghose[†] Onur Mutlu^{§†} [†]Carnegie Mellon University [‡]University of Virginia [§]ETH Zürich

How to Design Data Structures for PIM?

 Zhiyu Liu, Irina Calciu, Maurice Herlihy, and Onur Mutlu, "Concurrent Data Structures for Near-Memory Computing" Proceedings of the <u>29th ACM Symposium on Parallelism in Algorithms</u> <u>and Architectures</u> (SPAA), Washington, DC, USA, July 2017. [Slides (pptx) (pdf)]

Concurrent Data Structures for Near-Memory Computing

Zhiyu Liu Computer Science Department Brown University zhiyu_liu@brown.edu

Maurice Herlihy Computer Science Department Brown University mph@cs.brown.edu Irina Calciu VMware Research Group icalciu@vmware.com

Onur Mutlu Computer Science Department ETH Zürich onur.mutlu@inf.ethz.ch

Simulation Infrastructures for PIM

- Ramulator extended for PIM
 - Flexible and extensible DRAM simulator
 - Can model many different memory standards and proposals
 - Kim+, "Ramulator: A Flexible and Extensible DRAM Simulator", IEEE CAL 2015.
 - <u>https://github.com/CMU-SAFARI/ramulator</u>

Ramulator: A Fast and Extensible DRAM Simulator

Yoongu Kim¹ Weikun Yang^{1,2} Onur Mutlu¹ ¹Carnegie Mellon University ²Peking University

An FPGA-based Test-bed for PIM?

 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



Topics Requested (by Email)

- 1. PIM enabled instructions
- 2. DRAM Reliability and Performance
- 3. Processing in Memory (Near-Data Computation)

Memory/Storage Issues We Are Tackling

- (Enable and Exploit) In-Memory Computation
- (Reduce) Latency
- (Reduce) Energy
- (Improve) Reliability and Security
- (Make Sense of & Take Advantage of) Heterogeneity
- (Enable and Use) Persistence
- (Enable and Exploit) QoS and Predictability
- (Develop New and Reliable) Infrastructure

DRAM Reliability, Security, Refresh

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]



As DRAM cell becomes smaller, it becomes more vulnerable

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)

Testing DRAM Scaling Issues ...



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

Adaptive-Latency DRAM: Optimizing DRAM Timing for the Common-Case (Lee et al., HPCA 2015)

AVATAR: A Variable-Retention-Time (VRT) Aware Refresh for DRAM Systems (Qureshi et al., DSN 2015) An Experimental Study of Data Retention Behavior in Modern DRAM Devices: Implications for Retention Time Profiling Mechanisms (Liu et al., ISCA 2013)

The Efficacy of Error Mitigation Techniques for DRAM Retention Failures: A Comparative Experimental Study (Khan et al., SIGMETRICS 2014)



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





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

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

Hasan Hassan^{1,2,3} Nandita Vijaykumar³ Samira Khan^{4,3} Saugata Ghose³ Kevin Chang³ Gennady Pekhimenko^{5,3} Donghyuk Lee^{6,3} Oguz Ergin² Onur Mutlu^{1,3}

¹ETH Zürich ²TOBB University of Economics & Technology ³Carnegie Mellon University ⁴University of Virginia ⁵Microsoft Research ⁶NVIDIA 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)

Most DRAM Modules Are Vulnerable









Up to	Up to	Up to
1.0×10 ⁷	2.7×10 ⁶	3.3×10 ⁵
errors	errors	errors

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

Recent DRAM Is More Vulnerable



Recent DRAM Is More Vulnerable



Recent DRAM Is More Vulnerable



All modules from 2012–2013 are vulnerable



loop: mov (X), %eax mov (Y), %ebx clflush (X) clflush (X) mfence jmp loop





- Avoid *cache hits* Flush X from cache
- Avoid *row hits* to X
 Read Y in another row





loop: mov (X), %eax mov (Y), %ebx clflush (X) clflush (X) mfence jmp loop





loop: mov (X), %eax mov (Y), %ebx clflush (X) clflush (X) mfence jmp loop





loop: mov (X), %eax mov (Y), %ebx clflush (X) clflush (X) mfence jmp loop



Observed Errors in Real Systems

CPU Architecture	Errors	Access-Rate
Intel Haswell (2013)	22.9K	12.3M/sec
Intel Ivy Bridge (2012)	20.7K	11.7M/sec
Intel Sandy Bridge (2011)	16.1K	11.6M/sec
AMD Piledriver (2012)	59	6.1M/sec

A real reliability & security issue

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

One Can Take Over an Otherwise-Secure System

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

Abstract. Memory isolation is a key property of a reliable and secure computing system — an access to one memory address should not have unintended side effects on data stored in other addresses. However, as DRAM process technology

Project Zero

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

News and updates from the Project Zero team at Google

Exploiting the DRAM rowhammer bug to gain kernel privileges (Seaborn, 2015)

Monday, March 9, 2015

Exploiting the DRAM rowhammer bug to gain kernel privileges

RowHammer Security Attack Example

- "Rowhammer" is a problem with some recent DRAM devices in which repeatedly accessing a row of memory can cause bit flips in adjacent rows (Kim et al., ISCA 2014).
 - Flipping Bits in Memory Without Accessing Them: An Experimental Study of DRAM Disturbance Errors (Kim et al., ISCA 2014)
- We tested a selection of laptops and found that a subset of them exhibited the problem.
- We built two working privilege escalation exploits that use this effect.
 - Exploiting the DRAM rowhammer bug to gain kernel privileges (Seaborn, 2015)
- One exploit uses rowhammer-induced bit flips to gain kernel privileges on x86-64 Linux when run as an unprivileged userland process.
- When run on a machine vulnerable to the rowhammer problem, the process was able to induce bit flips in page table entries (PTEs).
- It was able to use this to gain write access to its own page table, and hence gain read-write access to all of physical memory.

Security Implications



It's like breaking into an apartment by repeatedly slamming a neighbor's door until the vibrations open the door you were after

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 Kim¹ Ross Daly^{*} Jeremie Kim¹ Chris Fallin^{*} Ji Hye Lee¹ Donghyuk Lee¹ Chris Wilkerson² Konrad Lai Onur Mutlu¹ ¹Carnegie Mellon University ²Intel Labs

Future of Memory Reliability

Onur Mutlu, "The RowHammer Problem and Other Issues We May Face as Memory Becomes Denser" Invited Paper in Proceedings of the Design, Automation, and Test in Europe Conference (DATE), Lausanne, Switzerland, March 2017. [Slides (pptx) (pdf)]

The RowHammer Problem and Other Issues We May Face as Memory Becomes Denser

Onur Mutlu ETH Zürich onur.mutlu@inf.ethz.ch https://people.inf.ethz.ch/omutlu

SAFARI https://people.inf.ethz.ch/omutlu/pub/rowhammer-and-other-memory-issues_date17.pdf 126

Future of Main Memory

• DRAM is becoming less reliable \rightarrow more vulnerable

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.

SAFARI

DRAM Reliability Reducing



Chip density (Gb)

Meza+, "Revisiting Memory Errors in Large-Scale Production Data Centers," DSN'15.

Future of Main Memory

- DRAM is becoming less reliable \rightarrow more vulnerable
- Due to difficulties in DRAM scaling, other problems may also appear (or they may be going unnoticed)
- Some errors may already be slipping into the field
 - Read disturb errors (Rowhammer)
 - Retention errors
 - Read errors, write errors

```
• ...
```

These errors can also pose security vulnerabilities

DRAM Data Retention Time Failures

- Determining the data retention time of a cell/row is getting more difficult
- Retention failures may already be slipping into the field

More on DRAM Data Retention

DRAM Refresh

- DRAM capacitor charge leaks over time
- The memory controller needs to refresh each row periodically to restore charge
 - Activate each row every N ms
 - Typical N = 64 ms
- Downsides of refresh
 - -- Energy consumption: Each refresh consumes energy
 - -- Performance degradation: DRAM rank/bank unavailable while refreshed
 - -- QoS/predictability impact: (Long) pause times during refresh
 - -- Refresh rate limits DRAM capacity scaling



Refresh Overhead: Performance



Refresh Overhead: Energy



Data Retention in Memory [Liu et al., ISCA 2013]

Retention Time Profile of DRAM looks like this:

64-128ms >256ms **Location** dependent 128-256ms Stored value pattern dependent **Time** dependent

Liu et al., "RAIDR: Retention-Aware Intelligent DRAM Refresh," ISCA 2012.

RAIDR: Eliminating Unnecessary Refreshes

- Observation: Most DRAM rows can be refreshed much less often without losing data [Kim+, EDL'09][Liu+ ISCA'13]
- Key idea: Refresh rows containing weak cells more frequently, other rows less frequently

1. Profiling: Profile retention time of all rows



50%

64 Gb

32 Gb

16 Gb Device capacity

2. Binning: Store rows into bins by retention time in memory controller

Efficient storage with Bloom Filters (only 1.25KB for 32GB memory)

14(

80

60

4020

4 Gb

8 Gb

Energy per

3. Refreshing: Memory controller refreshes rows in different bins at different rates 160

- access (n) 120 100 80 Results: 8-core, 32GB, SPEC, TPC-C, TPC-H
 - 74.6% refresh reduction @ 1.25KB storage
 - ~16%/20% DRAM dynamic/idle power reduction
 - ~9% performance improvement
 - Benefits increase with DRAM capacity

Liu et al., "RAIDR: Retention-Aware Intelligent DRAM Refresh," ISCA 2012.

Analysis of Data Retention Failures [ISCA'13]

 Jamie Liu, Ben Jaiyen, Yoongu Kim, Chris Wilkerson, and Onur Mutlu,
 "An Experimental Study of Data Retention Behavior in Modern DRAM Devices: Implications for Retention Time Profiling Mechanisms" Proceedings of the <u>40th International Symposium on Computer Architecture</u> (ISCA), Tel-Aviv, Israel, June 2013. <u>Slides (ppt)</u> <u>Slides (pdf)</u>

An Experimental Study of Data Retention Behavior in Modern DRAM Devices: Implications for Retention Time Profiling Mechanisms

Jamie Liu* Ben Jaiyen^{*} Yoongu Kim Carnegie Mellon University Carnegie Mellon University Carnegie Mellon University 5000 Forbes Ave. 5000 Forbes Ave. 5000 Forbes Ave. Pittsburgh, PA 15213 Pittsburgh, PA 15213 Pittsburgh, PA 15213 jamiel@alumni.cmu.edu bjaiyen@alumni.cmu.edu yoonguk@ece.cmu.edu Chris Wilkerson Onur Mutlu Intel Corporation Carnegie Mellon University 5000 Forbes Ave. 2200 Mission College Blvd. Santa Clara, CA 95054 Pittsburgh, PA 15213

onur@cmu.edu

chris.wilkerson@intel.com

Two Challenges to Retention Time Profiling

Data Pattern Dependence (DPD) of retention time

Variable Retention Time (VRT) phenomenon

Two Challenges to Retention Time Profiling

- Challenge 1: Data Pattern Dependence (DPD)
 - Retention time of a DRAM cell depends on its value and the values of cells nearby it
 - □ When a row is activated, all bitlines are perturbed simultaneously



Data Pattern Dependence

- Electrical noise on the bitline affects reliable sensing of a DRAM cell
- The magnitude of this noise is affected by values of nearby cells via
 - □ Bitline-bitline coupling \rightarrow electrical coupling between adjacent bitlines
 - Bitline-wordline coupling → electrical coupling between each bitline and the activated wordline



Data Pattern Dependence

- Electrical noise on the bitline affects reliable sensing of a DRAM cell
- The magnitude of this noise is affected by values of nearby cells via
 - □ Bitline-bitline coupling \rightarrow electrical coupling between adjacent bitlines
 - Bitline-wordline coupling → electrical coupling between each bitline and the activated wordline

- Retention time of a cell depends on data patterns stored in nearby cells
 - \rightarrow need to find the worst data pattern to find worst-case retention time
 - \rightarrow this pattern is location dependent

Two Challenges to Retention Time Profiling

- Challenge 2: Variable Retention Time (VRT)
 - Retention time of a DRAM cell changes randomly over time
 - a cell alternates between multiple retention time states
 - Leakage current of a cell changes sporadically due to a charge trap in the gate oxide of the DRAM cell access transistor
 - When the trap becomes occupied, charge leaks more readily from the transistor's drain, leading to a short retention time
 - Called *Trap-Assisted Gate-Induced Drain Leakage*
 - This process appears to be a random process [Kim + IEEE TED'11]
 - Worst-case retention time depends on a random process
 → need to find the worst case despite this

Modern DRAM Retention Time Distribution



Newer device families have more weak cells than older ones Likely a result of technology scaling
An Example VRT Cell



Variable Retention Time



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



147

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



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

Mitigation of Retention Issues [SIGMETRICS'14]

Samira Khan, Donghyuk Lee, Yoongu Kim, Alaa Alameldeen, Chris Wilkerson, and Onur Mutlu, "The Efficacy of Error Mitigation Techniques for DRAM Retention Failures: A Comparative Experimental Study" Proceedings of the ACM International Conference on Measurement and *Modeling of Computer Systems (SIGMETRICS)*, Austin, TX, June 2014. [Slides (pptx) (pdf)] [Poster (pptx) (pdf)] [Full data sets]

The Efficacy of Error Mitigation Techniques for DRAM Retention Failures: A Comparative Experimental Study

Samira Khan[†]* samirakhan@cmu.edu

Donghyuk Lee[†] donghyuk1@cmu.edu

Chris Wilkerson*

Yoongu Kim[†] yoongukim@cmu.edu

Alaa R. Alameldeen* alaa.r.alameldeen@intel.com chris.wilkerson@intel.com

Onur Mutlu[†] onur@cmu.edu

*Intel Labs [†]Carnegie Mellon University

Towards an Online Profiling System

Key Observations:

- Testing alone cannot detect all possible failures
- Combination of ECC and other mitigation techniques is much more effective
 - But degrades performance
- Testing can help to reduce the ECC strength
 - Even when starting with a higher strength ECC

Khan+, "The Efficacy of Error Mitigation Techniques for DRAM Retention Failures: A Comparative Experimental Study," SIGMETRICS 2014.

Towards an Online Profiling System



Handling Variable Retention Time [DSN'15]

 Moinuddin Qureshi, Dae Hyun Kim, Samira Khan, Prashant Nair, and Onur Mutlu, "AVATAR: A Variable-Retention-Time (VRT) Aware Refresh for DRAM <u>Systems</u>" Proceedings of the <u>45th Annual IEEE/IFIP International Conference on</u> <u>Dependable Systems and Networks</u> (DSN), Rio de Janeiro, Brazil, June 2015. [Slides (pptx) (pdf)]

AVATAR: A Variable-Retention-Time (VRT) Aware Refresh for DRAM Systems

Moinuddin K. Qureshi[†] Dae-Hyun Kim[†] [†]Georgia Institute of Technology {*moin, dhkim, pnair6*}@*ece.gatech.edu* Samira Khan[‡]

Prashant J. Nair[†] Onur Mutlu[‡] [‡]Carnegie Mellon University

{samirakhan, onur}@cmu.edu

Handling Data-Dependent Failures [DSN'16]

 Samira Khan, Donghyuk Lee, and Onur Mutlu,
 "PARBOR: An Efficient System-Level Technique to Detect Data-Dependent Failures in DRAM"
 Proceedings of the <u>45th Annual IEEE/IFIP International Conference on</u> Dependable Systems and Networks (DSN), Toulouse, France, June 2016.
 [Slides (pptx) (pdf)]

PARBOR: An Efficient System-Level Technique to Detect Data-Dependent Failures in DRAM

Samira Khan^{*} Donghyuk Lee^{†‡} Onur Mutlu^{*†} *University of Virginia [†]Carnegie Mellon University [‡]Nvidia ^{*}ETH Zürich

Handling Data-Dependent Failures [MICRO'17]

 Samira Khan, Chris Wilkerson, Zhe Wang, Alaa R. Alameldeen, Donghyuk Lee, and Onur Mutlu,
 "Detecting and Mitigating Data-Dependent DRAM Failures by Exploiting Current Memory Content"
 Proceedings of the 50th International Symposium on Microarchitecture (MICRO), Boston, MA, USA, October 2017.
 [Slides (pptx) (pdf)] [Lightning Session Slides (pptx) (pdf)] [Poster (pptx) (pdf)]

Detecting and Mitigating Data-Dependent DRAM Failures by Exploiting Current Memory Content

Samira Khan^{*} Chris Wilkerson[†] Zhe Wang[†] Alaa R. Alameldeen[†] Donghyuk Lee[‡] Onur Mutlu^{*} ^{*}University of Virginia [†]Intel Labs ^{*}Nvidia Research ^{*}ETH Zürich

Handling Both DPD and VRT [ISCA'17]

- Minesh Patel, Jeremie S. Kim, and Onur Mutlu,
 "The Reach Profiler (REAPER): Enabling the Mitigation of DRAM Retention Failures via Profiling at Aggressive Conditions"
 Proceedings of the <u>44th International Symposium on Computer</u> Architecture (ISCA), Toronto, Canada, June 2017.
 [Slides (pptx) (pdf)]
 [Lightning Session Slides (pptx) (pdf)]
- First experimental analysis of (mobile) LPDDR4 chips
- Analyzes the complex tradeoff space of retention time profiling
- Idea: enable fast and robust profiling at higher refresh intervals & temperatures

The Reach Profiler (REAPER): Enabling the Mitigation of DRAM Retention Failures via Profiling at Aggressive Conditions

Minesh Patel^{§‡} Jeremie S. Kim^{‡§} Onur Mutlu^{§‡} [§]ETH Zürich [‡]Carnegie Mellon University

How Do We Keep Memory Secure?

- Understand: Solid methodologies for failure modeling and discovery
 - Modeling based on real device data small scale and large scale

- Architect: Principled co-architecting of system and memory
 - Good partitioning of duties across the stack

- Design & Test: Principled electronic design, automation, testing
 - High coverage and good interaction with system reliability methods

Understand and Model with Experiments (DRAM)



SAFARI

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

Understand and Model with Experiments (Flash)



[DATE 2012, ICCD 2012, DATE 2013, ITJ 2013, ICCD 2013, SIGMETRICS 2014, NAND Daughter Board HPCA 2015, DSN 2015, MSST 2015, JSAC 2016, HPCA 2017, DFRWS 2017, PIEEE'17]

Cai+, "Error Characterization, Mitigation, and Recovery in Flash Memory Based Solid State Drives," Proc. IEEE 2017.

NAND Flash Errors and Reliability



Proceedings of the IEEE, Sept. 2017

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



This paper reviews the most recent advances in solid-state drive (SSD) error characterization, mitigation, and data recovery techniques to improve both SSD's reliability and lifetime.

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

https://arxiv.org/pdf/1706.08642

Reducing Memory Latency

Main Memory Latency Lags Behind



Memory latency remains almost constant

A Closer Look ...



Figure 1: DRAM latency trends over time [20, 21, 23, 51].

Chang+, "<u>Understanding Latency Variation in Modern DRAM Chips: Experimental</u> Characterization, Analysis, and Optimization"," SIGMETRICS 2016.

DRAM Latency Is Critical for Performance



In-memory Databases

[Mao+, EuroSys'12; Clapp+ (**Intel**), IISWC'15]



In-Memory Data Analytics

[Clapp+ (**Intel**), IISWC'15; Awan+, BDCloud'15]



Graph/Tree Processing [Xu+, IISWC'12; Umuroglu+, FPL'15]



Datacenter Workloads [Kanev+ (**Google**), ISCA'15]

DRAM Latency Is Critical for Performance





In-memory Databases

Graph/Tree Processing

Long memory latency \rightarrow performance bottleneck



In-Memory Data Analytics

[Clapp+ (**Intel**), IISWC'15; Awan+, BDCloud'15]



Datacenter Workloads [Kanev+ (**Google**), ISCA'I 5]

Design of DRAM uArchitecture

• Goal: Maximize capacity/area, not minimize latency

"One size fits all" approach to latency specification

- Same latency parameters for all temperatures
- Same latency parameters for all DRAM chips (e.g., rows)
- Same latency parameters for all parts of a DRAM chip
- Same latency parameters for all supply voltage levels
- Same latency parameters for all application data

□ ..

Latency Variation in Memory Chips

Heterogeneous manufacturing & operating conditions → latency variation in timing parameters



DRAM Characterization Infrastructure



SAFARI

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

DRAM Characterization 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



SoftMC: Open Source DRAM Infrastructure

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

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

Hasan Hassan^{1,2,3} Nandita Vijaykumar³ Samira Khan^{4,3} Saugata Ghose³ Kevin Chang³ Gennady Pekhimenko^{5,3} Donghyuk Lee^{6,3} Oguz Ergin² Onur Mutlu^{1,3}

¹ETH Zürich ²TOBB University of Economics & Technology ³Carnegie Mellon University ⁴University of Virginia ⁵Microsoft Research ⁶NVIDIA Research

Tackling the Fixed Latency Mindset

- Reliable operation latency is actually very heterogeneous
 Across temperatures, chips, parts of a chip, voltage levels, ...
- Idea: Dynamically find out and use the lowest latency one can reliably access a memory location with
 - Adaptive-Latency DRAM [HPCA 2015]
 - Flexible-Latency DRAM [SIGMETRICS 2016]
 - Design-Induced Variation-Aware DRAM [SIGMETRICS 2017]
 - Voltron [SIGMETRICS 2017]
 - ••••
- We would like to find sources of latency heterogeneity and exploit them to minimize latency

Adaptive-Latency DRAM

- Key idea
 - Optimize DRAM timing parameters online
- Two components
 - DRAM manufacturer provides multiple sets of reliable DRAM timing parameters at different temperatures for each DIMM
 - System monitors DRAM temperature & uses appropriate DRAM timing parameters

SAFARI Lee+, "Adaptive-Latency DRAM: Optimizing DRAM Timing for the Common-Case," HPCA 171 2015.

Latency Reduction Summary of 115 DIMMs

- Latency reduction for read & write (55°C)
 - Read Latency: **32.7%**
 - Write Latency: **55.1%**
- Latency reduction for each timing parameter (55°C)
 - Sensing: **17.3%**
 - Restore: **37.3%** (read), **54.8%** (write)
 - Precharge: **35.2%**

SAFARI Lee+, "Adaptive-Latency DRAM: Optimizing DRAM Timing for the Common-Case," HPCA 172 2015.

AL-DRAM: Real System Evaluation

System

- CPU: AMD 4386 (8 Cores, 3.1GHz, 8MB LLC)

D18F2x200_dct[0]_mp[1:0] DDR3 DRAM Timing 0		
Reset: 0F05_0505h. See 2.9.3 [DCT Configuration Registers].		
Bits	s Description	
31:30	0 Reserved.	
29:24	4 Tras: row active strobe. Read-write. BIOS: S the minimum time in memory clock cycles from to the same chip select bank. <u>Bits</u> <u>Description</u> 07h-00h Reserved 2Ah-08h <tras> clocks 3Fh-2Bh Reserved</tras>	ee 2.9.7.5 [SPD ROM-Based Configuration]. Specifies om an activate command to a precharge command, both
23:21	1 Reserved.	
20:16	Trp: row precharge time . Read-write. BIOS: See 2.9.7.5 [SPD ROM-Based Configuration]. Specifies the minimum time in memory clock cycles from a precharge command to an activate command or auto refresh command, both to the same bank.	



AL-DRAM *improves single-core performance* on a real system

AL-DRAM: Multi-Core Evaluation



AL-DRAM provides higher performance on multi-programmed & multi-threaded workloads SAFARI

Reducing Latency Also Reduces Energy

- AL-DRAM reduces DRAM power consumption by 5.8%
- Major reason: reduction in row activation time

More on Adaptive-Latency DRAM

 Donghyuk Lee, Yoongu Kim, Gennady Pekhimenko, Samira Khan, Vivek Seshadri, Kevin Chang, and Onur Mutlu,
 <u>"Adaptive-Latency DRAM: Optimizing DRAM Timing for</u> <u>the Common-Case"</u>
 Proceedings of the <u>21st International Symposium on High-</u> <i>Performance Computer Architecture (HPCA), Bay Area, CA, February 2015.
 [Slides (pptx) (pdf)] [Full data sets]

Adaptive-Latency DRAM: Optimizing DRAM Timing for the Common-Case

Donghyuk Lee Yoongu Kim Gennady Pekhimenko Samira Khan Vivek Seshadri Kevin Chang Onur Mutlu

Carnegie Mellon University

Heterogeneous Latency within A Chip



40 Workloads

Chang+, "<u>Understanding Latency Variation in Modern DRAM Chips: Experimental</u> Characterization, Analysis, and Optimization"," SIGMETRICS 2016.

Analysis of Latency Variation in DRAM Chips

- Kevin Chang, Abhijith Kashyap, Hasan Hassan, Samira Khan, Kevin Hsieh, Donghyuk Lee, Saugata Ghose, Gennady Pekhimenko, Tianshi Li, and Onur Mutlu,
 - "Understanding Latency Variation in Modern DRAM Chips: Experimental Characterization, Analysis, and Optimization" Proceedings of the ACM International Conference on Measurement and
 - Proceedings of the <u>ACM International Conference on Measurement and</u> <u>Modeling of Computer Systems</u> (**SIGMETRICS**), Antibes Juan-Les-Pins, France, June 2016. [<u>Slides (pptx) (pdf)</u>]
 - [Source Code]

Understanding Latency Variation in Modern DRAM Chips: Experimental Characterization, Analysis, and Optimization

Kevin K. Chang¹ Abhijith Kashyap¹ Hasan Hassan^{1,2} Saugata Ghose¹ Kevin Hsieh¹ Donghyuk Lee¹ Tianshi Li^{1,3} Gennady Pekhimenko¹ Samira Khan⁴ Onur Mutlu^{5,1} ¹Carnegie Mellon University ²TOBB ETÜ ³Peking University ⁴University of Virginia ⁵ETH Zürich SAFARI



Systematic variation in cell access times caused by the *physical organization* of DRAM


Profile only slow regions to determine min. latency → Dynamic & low cost latency optimization



Combine error-correcting codes & online profiling → Reliably reduce DRAM latency

DIVA-DRAM Reduces Latency



DIVA-DRAM *reduces latency more aggressively* and uses ECC to correct random slow cells

Design-Induced Latency Variation in DRAM

- Donghyuk Lee, Samira Khan, Lavanya Subramanian, Saugata Ghose, Rachata Ausavarungnirun, Gennady Pekhimenko, Vivek Seshadri, and <u>Onur Mutlu</u>,
 - "Design-Induced Latency Variation in Modern DRAM Chips: Characterization, Analysis, and Latency Reduction Mechanisms" Proceedings of the <u>ACM International Conference on Measurement and</u> <u>Modeling of Computer Systems</u> (SIGMETRICS), Urbana-Champaign, IL, USA, June 2017.

Design-Induced Latency Variation in Modern DRAM Chips: Characterization, Analysis, and Latency Reduction Mechanisms

Donghyuk Lee, NVIDIA and Carnegie Mellon University

Samira Khan, University of Virginia

Lavanya Subramanian, Saugata Ghose, Rachata Ausavarungnirun, Carnegie Mellon University Gennady Pekhimenko, Vivek Seshadri, Microsoft Research

Onur Mutlu, ETH Zürich and Carnegie Mellon University

Voltron: Exploiting the Voltage-Latency-Reliability Relationship

Executive Summary

- DRAM (memory) power is significant in today's systems
 - Existing low-voltage DRAM reduces voltage **conservatively**
- <u>Goal</u>: Understand and exploit the reliability and latency behavior of real DRAM chips under *aggressive reduced-voltage operation*

• Key experimental observations:

- Huge voltage margin -- Errors occur beyond some voltage
- Errors exhibit spatial locality
- Higher operation latency mitigates voltage-induced errors
- Voltron: A new DRAM energy reduction mechanism
 - Reduce DRAM voltage without introducing errors
 - Use a **regression model** to select voltage that does not degrade performance beyond a chosen target \rightarrow 7.3% system energy reduction

Analysis of Latency-Voltage in DRAM Chips

 Kevin Chang, A. Giray Yaglikci, Saugata Ghose, Aditya Agrawal, Niladrish Chatterjee, Abhijith Kashyap, Donghyuk Lee, Mike O'Connor, Hasan Hassan, and <u>Onur Mutlu</u>, <u>"Understanding Reduced-Voltage Operation in Modern DRAM</u> <u>Devices: Experimental Characterization, Analysis, and</u> <u>Mechanisms"</u> *Proceedings of the <u>ACM International Conference on Measurement and</u> <u>Modeling of Computer Systems</u> (<i>SIGMETRICS*), Urbana-Champaign, IL, USA, June 2017.

Understanding Reduced-Voltage Operation in Modern DRAM Chips: Characterization, Analysis, and Mechanisms

Kevin K. Chang[†] Abdullah Giray Yağlıkçı[†] Saugata Ghose[†] Aditya Agrawal[¶] Niladrish Chatterjee[¶] Abhijith Kashyap[†] Donghyuk Lee[¶] Mike O'Connor^{¶,‡} Hasan Hassan[§] Onur Mutlu^{§,†}

[†]Carnegie Mellon University [¶]NVIDIA [‡]The University of Texas at Austin [§]ETH Zürich

And, What If ...

we can sacrifice reliability of some data to access it with even lower latency?

SAFARI Group: Research Introduction

Onur Mutlu <u>omutlu@gmail.com</u> <u>https://people.inf.ethz.ch/omutlu</u> January 31, 2018



