Processing Data Where It Makes Sense in Modern Computing Systems: Enabling In-Memory Computation

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

16 August 2019

APPT Keynote



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

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
 - 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*!

DRAM Capacity, Bandwidth & Latency





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]





In-memory Databases

Graph/Tree Processing

Memory → performance bottleneck



In-Memory Data Analytics

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



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





Chrome

Google's web browser

TensorFlow Mobile

Google's machine learning framework



Google's video codec





Memory → performance bottleneck



Google's video codec



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

Energy Waste in Mobile Devices

 Amirali Boroumand, Saugata Ghose, Youngsok Kim, Rachata Ausavarungnirun, Eric Shiu, Rahul Thakur, Daehyun Kim, Aki Kuusela, Allan Knies, Parthasarathy Ranganathan, and Onur Mutlu, "Google Workloads for Consumer Devices: Mitigating Data Movement Bottlenecks" Proceedings of the <u>23rd International Conference on Architectural Support for Programming</u> <u>Languages and Operating Systems</u> (ASPLOS), Williamsburg, VA, USA, March 2018.

62.7% of the total system energy is spent on data movement

Google Workloads for Consumer Devices: Mitigating Data Movement Bottlenecks

Amirali Boroumand¹Saugata Ghose¹Youngsok Kim²Rachata Ausavarungnirun¹Eric Shiu³Rahul Thakur³Daehyun Kim^{4,3}Aki Kuusela³Allan Knies³Parthasarathy Ranganathan³Onur Mutlu^{5,1}16

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., RL/TL-DRAM, FLY-RAM)	lower latency	higher cost
Low-Power DRAM (e.g., LPDDR3, LPDDR4, Voltron)	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 Trend: Hybrid Main Memory



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.



Main Memory Needs Intelligent Controllers

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



Call for Intelligent Memory Controllers

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



23

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



- Major Trends Affecting Main Memory
- The Need for Intelligent Memory Controllers
 - Bottom Up: Push from Circuits and Devices
 - Top Down: Pull from Systems and Applications
- Processing in Memory: Two Directions
 - Minimally Changing Memory Chips
 - Exploiting 3D-Stacked Memory
- How to Enable Adoption of Processing in Memory
- Conclusion

Maslow's (Human) Hierarchy of Needs



We need to start with reliability and security...

How Reliable/Secure/Safe is This Bridge?



Collapse of the "Galloping Gertie"



How Secure Are These People?



Security is about preventing unforeseen consequences

Source: https://s-media-cache-ak0.pinimg.com/originals/48/09/54/4809543a9c7700246a0cf8acdae27abf.jpg

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

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



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

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

Retention Time Profile of DRAM looks like this:

64-128ms >256ms



Location dependent Stored value pattern dependent Time dependent


Main Memory Needs Intelligent Controllers

More on DRAM Refresh (I)

 Jamie Liu, Ben Jaiyen, Richard Veras, and Onur Mutlu, "RAIDR: Retention-Aware Intelligent DRAM Refresh" Proceedings of the <u>39th International Symposium on</u> <u>Computer Architecture</u> (ISCA), Portland, OR, June 2012. <u>Slides (pdf)</u>

RAIDR: Retention-Aware Intelligent DRAM Refresh

Jamie Liu Ben Jaiyen Richard Veras Onur Mutlu Carnegie Mellon University

More on DRAM Refresh (II)

chris.wilkerson@intel.com

 Jamie Liu, Ben Jaiyen, Yoongu Kim, Chris Wilkerson, and <u>Onur Mutlu</u>,
 "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 bjaiyen@alumni.cmu.edu jamiel@alumni.cmu.edu yoonguk@ece.cmu.edu Chris Wilkerson Onur Mutlu Intel Corporation Carnegie Mellon University 2200 Mission College Blvd. 5000 Forbes Ave. Santa Clara, CA 95054 Pittsburgh, PA 15213

onur@cmu.edu

More on DRAM Refresh (III)

- 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 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

A company B company







C company

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

Flipping Bits in Memory Without Accessing Them: An Experimental Study of DRAM Disturbance Errors, (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 (Y) 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 (Y) mfence jmp loop





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





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

<u>Flipping Bits in Memory Without Accessing Them:</u> <u>An Experimental Study of DRAM Disturbance Errors</u> (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.

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

Security Implications



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 Security Implications (I)

"We can gain unrestricted access to systems of website visitors."

Not there yet, but ...



ROOT privileges for web apps!

Daniel Gruss (@lavados), Clémentine Maurice (@BloodyTangerine), December 28, 2015 - 32c3, Hamburg, Germany





Rowhammer.js: A Remote Software-Induced Fault Attack in JavaScript (DIMVA'16)

Source: https://lab.dsst.io/32c3-slides/7197.html

29

More Security Implications (II)

"Can gain control of a smart phone deterministically"

Hammer And Root

androids Millions of Androids

Drammer: Deterministic Rowhammer

Attacks on Mobile Platforms, CCS'16 59

Source: https://fossbytes.com/drammer-rowhammer-attack-android-root-devices/

More Security Implications (III)

 Using an integrated GPU in a mobile system to remotely escalate privilege via the WebGL interface

ars TECHNICA

BIZ & IT TECH SCIENCE POLICY CARS GAMING & CULTURE

Drive-by Rowhammer attack uses GPU to compromise an Android phone

JavaScript based GLitch pwns browsers by flipping bits inside memory chips.

DAN GOODIN - 5/3/2018, 12:00 PM

Grand Pwning Unit: Accelerating Microarchitectural Attacks with the GPU

Pietro Frigo Vrije Universiteit Amsterdam p.frigo@vu.nl Cristiano Giuffrida Vrije Universiteit Amsterdam giuffrida@cs.vu.nl Herbert Bos Vrije Universiteit Amsterdam herbertb@cs.vu.nl Kaveh Razavi Vrije Universiteit Amsterdam kaveh@cs.vu.nl

More Security Implications (IV)

Rowhammer over RDMA (I)

ars TECHNICA

BIZ & IT TECH SCIENCE POLICY CARS GAMING & CULTURE

THROWHAMMER —

Packets over a LAN are all it takes to trigger serious Rowhammer bit flips

The bar for exploiting potentially serious DDR weakness keeps getting lower.

DAN GOODIN - 5/10/2018, 5:26 PM

Throwhammer: Rowhammer Attacks over the Network and Defenses

Andrei Tatar VU Amsterdam Radhesh Krishnan VU Amsterdam Elias Athanasopoulos University of Cyprus

Herbert Bos VU Amsterdam Kaveh Razavi VU Amsterdam Cristiano Giuffrida VU Amsterdam

More Security Implications (V)

Rowhammer over RDMA (II)

Security in a serious way

Nethammer—Exploiting DRAM Rowhammer Bug Through Network Requests



Nethammer: Inducing Rowhammer Faults through Network Requests

Moritz Lipp Graz University of Technology

Daniel Gruss Graz University of Technology Misiker Tadesse Aga University of Michigan

Clémentine Maurice Univ Rennes, CNRS, IRISA

Lukas Lamster Graz University of Technology Michael Schwarz Graz University of Technology

Lukas Raab Graz University of Technology

More Security Implications?



Apple's Patch for RowHammer

https://support.apple.com/en-gb/HT204934

Available for: OS X Mountain Lion v10.8.5, OS X Mavericks v10.9.5

Impact: A malicious application may induce memory corruption to escalate privileges

Description: A disturbance error, also known as Rowhammer, exists with some DDR3 RAM that could have led to memory corruption. This issue was mitigated by increasing memory refresh rates.

CVE-ID

CVE-2015-3693 : Mark Seaborn and Thomas Dullien of Google, working from original research by Yoongu Kim et al (2014)

HP, Lenovo, and other vendors released similar patches

Our Solution to RowHammer

- PARA: <u>Probabilistic Adjacent Row Activation</u>
- Key Idea
 - After closing a row, we activate (i.e., refresh) one of its neighbors with a low probability: p = 0.005
- Reliability Guarantee
 - When p=0.005, errors in one year: 9.4×10^{-14}
 - By adjusting the value of p, we can vary the strength of protection against errors

Advantages of PARA

- PARA refreshes rows infrequently
 - Low power
 - Low performance-overhead
 - Average slowdown: 0.20% (for 29 benchmarks)
 - Maximum slowdown: 0.75%
- PARA is stateless
 - Low cost
 - Low complexity
- PARA is an effective and low-overhead solution to prevent disturbance errors

Requirements for PARA

- If implemented in DRAM chip (done today)
 - Enough slack in timing and refresh parameters
 - Plenty of slack today:
 - Lee et al., "Adaptive-Latency DRAM: Optimizing DRAM Timing for the Common Case," HPCA 2015.
 - Chang et al., "Understanding Latency Variation in Modern DRAM Chips," SIGMETRICS 2016.
 - Lee et al., "Design-Induced Latency Variation in Modern DRAM Chips," SIGMETRICS 2017.
 - Chang et al., "Understanding Reduced-Voltage Operation in Modern DRAM Devices," SIGMETRICS 2017.
 - Ghose et al., "What Your DRAM Power Models Are Not Telling You: Lessons from a Detailed Experimental Study," SIGMETRICS 2018.
- If implemented in memory controller
 - Better coordination between memory controller and DRAM
 - Memory controller should know which rows are physically adjacent

Probabilistic Activation in Real Life (I)

Channel O Slot O Size	Populated & Enabled 16384 MB (DDR4)	Type of method used to prever Row Hammer
Number of Ranks Manufacturer	2 UnKnown	
Channel 0 Slot 1	Not Populated / Disable	d
Size	16384 MB (DDR4)	•
Number of Ranks	2 UnKnown	
Channel 1 Slot 1	Not Populated / Disable	d
Nemery patie/peference clock	Hardware RHP	
options moved to	2x Refresh	the Salast Sanaan
Overclock->Memory->Custom Profi	1	tl: Select Item
MRC ULT Safe Config	[Disabled]	Enter: Select
Maximum Memory Frequency	[Auto]	F1: General Help
Max TOLUD	[Dynamic]	F2: Previous Values
SA GV	[Enabled] [MRC default]	F4: Save & Exit
Retrain on Fast Fail	[Enabled]	ESC: Exit
Command Tristate	[Enabled]	
Row Hammer Solution	[Hardware RHP]	•
Row Hammer Solution		

SAFARI

https://twitter.com/isislovecruft/status/1021939922754723841

Probabilistic Activation in Real Life (II)



SAFARI

https://twitter.com/isislovecruft/status/1021939922754723841

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 71

A RowHammer Retrospective

Onur Mutlu and Jeremie Kim,
 "RowHammer: A Retrospective"
 IEEE Transactions on Computer-Aided Design of Integrated
 Circuits and Systems (TCAD) Special Issue on Top Picks in
 Hardware and Embedded Security, 2019.
 [Preliminary arXiv version]

RowHammer: A Retrospective

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



Call for Intelligent Memory Controllers

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



74

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

Aside: Intelligent Controller for NAND Flash



[DATE 2012, ICCD 2012, DATE 2013, ITJ 2013, ICCD 2013, SIGMETRICS 2014, HPCA 2015, DSN 2015, MSST 2015, JSAC 2016, HPCA 2017, DFRWS 2017, PIEEE 2017, HPCA 2018, SIGMETRICS 2018]

NAND Daughter Board

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

Aside: Intelligent Controller for NAND Flash



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



Main Memory Needs Intelligent Controllers



- Major Trends Affecting Main Memory
- The Need for Intelligent Memory Controllers
 - Bottom Up: Push from Circuits and Devices
 - Top Down: Pull from Systems and Applications
- Processing in Memory: Two Directions
 - Minimally Changing Memory Chips
 - Exploiting 3D-Stacked Memory
- How to Enable Adoption of Processing in Memory
- Conclusion

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]

Do We Want This?



SAFARI Source

Or This?



Maslow's (Human) Hierarchy of Needs, Revisited



Challenge and Opportunity for Future

High Performance, Energy Efficient, Sustainable

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.

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 Performance Perspective (Today)

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



Kanev+, "Profiling a Warehouse-Scale Computer," ISCA 2015.

The Performance Perspective (Today)

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



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

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

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)

Energy Waste in Mobile Devices

 Amirali Boroumand, Saugata Ghose, Youngsok Kim, Rachata Ausavarungnirun, Eric Shiu, Rahul Thakur, Daehyun Kim, Aki Kuusela, Allan Knies, Parthasarathy Ranganathan, and Onur Mutlu, "Google Workloads for Consumer Devices: Mitigating Data Movement Bottlenecks" Proceedings of the <u>23rd International Conference on Architectural Support for Programming</u> <u>Languages and Operating Systems</u> (ASPLOS), Williamsburg, VA, USA, March 2018.

62.7% of the total system energy is spent on data movement

Google Workloads for Consumer Devices: Mitigating Data Movement Bottlenecks

Amirali Boroumand¹Saugata Ghose¹Youngsok Kim²Rachata Ausavarungnirun¹Eric Shiu³Rahul Thakur³Daehyun Kim^{4,3}Aki Kuusela³Allan Knies³Parthasarathy Ranganathan³Onur Mutlu^{5,1}99

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



- Major Trends Affecting Main Memory
- The Need for Intelligent Memory Controllers
 - Bottom Up: Push from Circuits and Devices
 - Top Down: Pull from Systems and Applications
- Processing in Memory: Two Directions
 - Minimally Changing Memory Chips
 - Exploiting 3D-Stacked Memory
- How to Enable Adoption of Processing in Memory
- Conclusion

Processing in Memory: Two Approaches

Minimally changing memory chips
 Exploiting 3D-stacked memory

Approach 1: 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

Starting Simple: Data Copy and Initialization

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





VM Cloning Deduplication



Many more

Page Migration

Today's Systems: Bulk Data Copy


Future Systems: In-Memory Copy



RowClone: In-DRAM Row Copy



Data Bus

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. [<u>Slides (pptx) (pdf)</u>] [<u>Lightning Session</u> <u>Slides (pptx) (pdf)</u>] [<u>Poster (pptx) (pdf)</u>]

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

Memory as an Accelerator



Memory similar to a "conventional" accelerator

In-Memory Bulk Bitwise 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

In-DRAM AND/OR: Triple Row Activation



Seshadri+, "Fast Bulk Bitwise AND and OR in DRAM", IEEE CAL 2015.

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.



Performance: In-DRAM Bitwise Operations



Figure 9: Throughput of bitwise operations on various systems.

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)	(\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.

Ambit vs. DDR3: Performance and Energy



Bulk Bitwise Operations in Workloads



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

Example Data Structure: Bitmap Index

- Alternative to B-tree and its variants
- Efficient for performing *range queries* and *joins*
- Many bitwise operations to perform a query



Performance: Bitmap Index on Ambit



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

>5.4-6.6X Performance Improvement

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> *IEEE Computer Architecture Letters (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

More on In-DRAM Bulk Bitwise Execution

 Vivek Seshadri and Onur Mutlu,
"In-DRAM Bulk Bitwise Execution Engine" Invited Book Chapter in Advances in Computers, to appear in 2020.
[Preliminary arXiv version]

In-DRAM Bulk Bitwise Execution Engine

Vivek Seshadri Microsoft Research India visesha@microsoft.com Onur Mutlu ETH Zürich onur.mutlu@inf.ethz.ch Challenge and Opportunity for Future

Computing Architectures with

Minimal Data Movement



Challenge: Intelligent Memory Device

Does memory have to be dumb?



- Major Trends Affecting Main Memory
- The Need for Intelligent Memory Controllers
 - Bottom Up: Push from Circuits and Devices
 - Top Down: Pull from Systems and Applications
- Processing in Memory: Two Directions
 - Minimally Changing Memory Chips
 - Exploiting 3D-Stacked Memory
- How to Enable Adoption of Processing in Memory
- Conclusion

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 are the performance and energy benefits of using 3D-stacked memory as a coarse-grained accelerator?
 - By changing the entire system
 - By performing simple function offloading

- What is the minimal processing-in-memory support we can provide?
 - With minimal changes to system and programming

Graph Processing

Large graphs are everywhere (circa 2015)



Scalable large-scale graph processing is challenging



Key Bottlenecks in Graph Processing



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



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



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[§] 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

Two Key Questions in 3D-Stacked PIM

- What are the performance and energy benefits of using 3D-stacked memory as a coarse-grained accelerator?
 - By changing the entire system
 - By performing simple function offloading

- What is the minimal processing-in-memory support we can provide?
 - With minimal changes to system and programming
PIM on Mobile Devices

 Amirali Boroumand, Saugata Ghose, Youngsok Kim, Rachata Ausavarungnirun, Eric Shiu, Rahul Thakur, Daehyun Kim, Aki Kuusela, Allan Knies, Parthasarathy Ranganathan, and Onur Mutlu, "Google Workloads for Consumer Devices: Mitigating Data Movement Bottlenecks"

Proceedings of the <u>23rd International Conference on Architectural</u> <u>Support for Programming Languages and Operating</u> <u>Systems</u> (**ASPLOS**), Williamsburg, VA, USA, March 2018.

Google Workloads for Consumer Devices: Mitigating Data Movement Bottlenecks

Amirali Boroumand1Saugata Ghose1Youngsok Kim2Rachata Ausavarungnirun1Eric Shiu3Rahul Thakur3Daehyun Kim4,3Aki Kuusela3Allan Knies3Parthasarathy Ranganathan3Onur Mutlu^{5,1}

Consumer Devices



Consumer devices are everywhere!

Energy consumption is a first-class concern in consumer devices



Four Important Workloads





Chrome

Google's web browser

TensorFlow Mobile

Google's machine learning framework



Google's video codec



Energy Cost of Data Movement

Ist key observation: 62.7% of the total system energy is spent on data movement



Processing-In-Memory (PIM)

Potential solution: move computation close to data

Challenge: limited area and energy budget

Using PIM to Reduce Data Movement

2nd key observation: a significant fraction of the data movement often comes from simple functions

We can design lightweight logic to implement these <u>simple functions</u> in <u>memory</u>

Small embedded low-power core

> PIM Core

Small fixed-function accelerators



Offloading to PIM logic reduces energy and improves performance, on average, by 55.4% and 54.2%

Workload Analysis





Chrome Google's web browser



TensorFlow Mobile

Google's machine learning framework



Google's video codec



TensorFlow Mobile



57.3% of the inference energy is spent on data movement

54.4% of the data movement energy comes from packing/unpacking_and quantization



Packing



A simple data reorganization process that requires simple arithmetic

Quantization



A simple data conversion operation that requires shift, addition, and multiplication operations

Normalized Energy





PIM core and PIM accelerator reduce <u>energy consumption</u> on average by 49.1% and 55.4% SAFARI

Normalized Runtime

S CPU-Only ■ PIM-Core ■ PIM-Acc



Offloading these kernels to PIM core and PIM accelerator improves performance on average by 44.6% and 54.2%

More on PIM for Mobile Devices

 Amirali Boroumand, Saugata Ghose, Youngsok Kim, Rachata Ausavarungnirun, Eric Shiu, Rahul Thakur, Daehyun Kim, Aki Kuusela, Allan Knies, Parthasarathy Ranganathan, and Onur Mutlu, "Google Workloads for Consumer Devices: Mitigating Data Movement Bottlenecks" Proceedings of the <u>23rd International Conference on Architectural Support for Programming</u> <u>Languages and Operating Systems</u> (ASPLOS), Williamsburg, VA, USA, March 2018.

62.7% of the total system energy is spent on data movement

Google Workloads for Consumer Devices: Mitigating Data Movement Bottlenecks

Amirali Boroumand¹Saugata Ghose¹Youngsok Kim²Rachata Ausavarungnirun¹Eric Shiu³Rahul Thakur³Daehyun Kim^{4,3}Aki Kuusela³Allan Knies³Parthasarathy Ranganathan³Onur Mutlu^{5,1}156

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

Accelerating Linked Data Structures

 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

Accelerating Dependent Cache Misses

 Milad Hashemi, Khubaib, Eiman Ebrahimi, Onur Mutlu, and Yale N. Patt, "Accelerating Dependent Cache Misses with an Enhanced <u>Memory Controller"</u> *Proceedings of the <u>43rd International Symposium on Computer</u> <i>Architecture (ISCA)*, Seoul, South Korea, June 2016. [Slides (pptx) (pdf)] [Lightning Session Slides (pptx) (pdf)]

Accelerating Dependent Cache Misses with an Enhanced Memory Controller

Milad Hashemi^{*}, Khubaib[†], Eiman Ebrahimi[‡], Onur Mutlu[§], Yale N. Patt^{*}

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

Two Key Questions in 3D-Stacked PIM

- What are the performance and energy benefits of using 3D-stacked memory as a coarse-grained accelerator?
 - By changing the entire system
 - By performing simple function offloading

What is the minimal processing-in-memory support we can provide?

With minimal changes to system and programming

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



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

Example (Abstract) PEI uArchitecture



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	4 GB, 16 vaults, 256 DRAM banks [20]
– DRAM	FR-FCFS, tCL = tRCD = tRP = 13.75 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

PEI Performance Delta: Large Data Sets



PIM-Only Locality-Aware

PEI Energy Consumption





Simpler PIM: 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

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>*

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

Automatic Offloading of Critical Code

 Milad Hashemi, Khubaib, Eiman Ebrahimi, Onur Mutlu, and Yale N. Patt, "Accelerating Dependent Cache Misses with an Enhanced <u>Memory Controller"</u> *Proceedings of the <u>43rd International Symposium on Computer</u> <i>Architecture (ISCA)*, Seoul, South Korea, June 2016. [Slides (pptx) (pdf)] [Lightning Session Slides (pptx) (pdf)]

Accelerating Dependent Cache Misses with an Enhanced Memory Controller

Milad Hashemi^{*}, Khubaib[†], Eiman Ebrahimi[‡], Onur Mutlu[§], Yale N. Patt^{*}

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

Automatic Offloading of Prefetch Mechanisms

 Milad Hashemi, Onur Mutlu, and Yale N. Patt,
 "Continuous Runahead: Transparent Hardware Acceleration for Memory Intensive Workloads"
 Proceedings of the <u>49th International Symposium on</u> <u>Microarchitecture</u> (MICRO), Taipei, Taiwan, October 2016.
 [Slides (pptx) (pdf)] [Lightning Session Slides (pdf)] [Poster (pptx) (pdf)]

Continuous Runahead: Transparent Hardware Acceleration for Memory Intensive Workloads

Milad Hashemi^{*}, Onur Mutlu[§], Yale N. Patt^{*}

* The University of Texas at Austin §ETH Zürich

Efficient Automatic Data Coherence Support

 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"

<u>IEEE Computer Architecture Letters</u> (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



Efficient Automatic Data Coherence Support

Amirali Boroumand, Saugata Ghose, Minesh Patel, Hasan Hassan, Brandon Lucia, Kevin Hsieh, Krishna T. Malladi, Hongzhong Zheng, and Onur Mutlu, "CoNDA: Efficient Cache Coherence Support for Near-**Data Accelerators**" Proceedings of the <u>46th International Symposium on Computer</u>

Architecture (ISCA), Phoenix, AZ, USA, June 2019.

CoNDA: Efficient Cache Coherence Support for Near-Data Accelerators

Saugata Ghose[†] Minesh Patel^{*} Hasan Hassan^{*} Amirali Boroumand[†] Brandon Lucia[†] Rachata Ausavarungnirun^{†‡} Kevin Hsieh[†] Nastaran Hajinazar^{\phi†} Krishna T. Malladi[§] Hongzhong Zheng[§] Onur Mutlu^{*†} [†]Carnegie Mellon University *****ETH Zürich [‡]KMUTNB

*Simon Fraser University

[§]Samsung Semiconductor, Inc.

Challenge and Opportunity for Future

Computing Architectures with

Minimal Data Movement





- Major Trends Affecting Main Memory
- The Need for Intelligent Memory Controllers
 - Bottom Up: Push from Circuits and Devices
 - Top Down: Pull from Systems and Applications
- Processing in Memory: Two Directions
 - Minimally Changing Memory Chips
 - Exploiting 3D-Stacked Memory
- How to Enable Adoption of Processing in Memory
- Conclusion

Eliminating the Adoption Barriers

How to Enable Adoption of Processing in Memory

Barriers to Adoption of PIM

1. Functionality of and applications & software for PIM

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

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

5. Infrastructures and models to assess benefits and feasibility

All can be solved with change of mindset
We Need to Revisit the Entire Stack

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

We can get there step by step

SAFARI

PIM Review and Open Problems

Processing Data Where It Makes Sense: Enabling In-Memory Computation

Onur Mutlu^{a,b}, Saugata Ghose^b, Juan Gómez-Luna^a, Rachata Ausavarungnirun^{b,c}

^aETH Zürich ^bCarnegie Mellon University ^cKing Mongkut's University of Technology North Bangkok

Onur Mutlu, Saugata Ghose, Juan Gomez-Luna, and Rachata Ausavarungnirun, "Processing Data Where It Makes Sense: Enabling In-Memory Computation" Invited paper in Microprocessors and Microsystems (MICPRO), June 2019. [arXiv version]

https://arxiv.org/pdf/1903.03988.pdf

SAFAR

PIM Review and Open Problems (II)

A Workload and Programming Ease Driven Perspective of Processing-in-Memory

Saugata Ghose†Amirali Boroumand†Jeremie S. Kim†§Juan Gómez-Luna§Onur Mutlu§††Carnegie Mellon University§ETH Zürich

Saugata Ghose, Amirali Boroumand, Jeremie S. Kim, Juan Gomez-Luna, and Onur Mutlu, "Processing-in-Memory: A Workload-Driven Perspective" *Invited Article in IBM Journal of Research & Development, Special Issue on Hardware for Artificial Intelligence*, to appear in November 2019. [Preliminary arXiv version]

SAFARI

https://arxiv.org/pdf/1907.12947.pdf

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>*

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

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

SAFARI

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-pim</u>
 - <u>https://github.com/CMU-SAFARI/ramulator</u>
 - Source Code for Ramulator-PIM

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



Simulation Infrastructures for PIM (in SSDs)

 Arash Tavakkol, Juan Gomez-Luna, Mohammad Sadrosadati, Saugata Ghose, and <u>Onur Mutlu</u>,
 "MQSim: A Framework for Enabling Realistic Studies of Modern Multi-Queue SSD Devices"
 Proceedings of the <u>16th USENIX Conference on File and Storage</u> <u>Technologies</u> (FAST), Oakland, CA, USA, February 2018.
 [Slides (pptx) (pdf)]
 [Source Code]

MQSim: A Framework for Enabling Realistic Studies of Modern Multi-Queue SSD Devices

Arash Tavakkol[†], Juan Gómez-Luna[†], Mohammad Sadrosadati[†], Saugata Ghose[‡], Onur Mutlu^{†‡} [†]*ETH Zürich* [‡]*Carnegie Mellon University*

Performance & Energy Models for PIM

- Gagandeep Singh, Juan Gomez-Luna, Giovanni Mariani, Geraldo F. Oliveira, Stefano Corda, Sander Stujik, <u>Onur Mutlu</u>, and Henk Corporaal, <u>"NAPEL: Near-Memory Computing Application Performance</u> <u>Prediction via Ensemble Learning"</u> *Proceedings of the <u>56th Design Automation Conference</u> (<i>DAC*), Las Vegas, NV, USA, June 2019.
 [Slides (pptx) (pdf)]
 [Poster (pptx) (pdf)]
 - Source Code for Ramulator-PIM

NAPEL: Near-Memory Computing Application Performance Prediction via Ensemble Learning

Gagandeep Singh^{*a,c*} Juan Gómez-Luna^{*b*} Stefano Corda^{*a,c*} Sander Stuijk^{*a*} ^{*a*}Eindhoven University of Technology ^{*b*}E

Juan Gómez-Luna^bGiovanni Mariani^cGeraldo F. Oliveira^bSander Stuijk^aOnur Mutlu^bHenk Corporaal^aiversity of Technology^bETH Zürich^cIBM Research - Zurich

New Applications and Use Cases for PIM

 Jeremie S. Kim, Damla Senol Cali, Hongyi Xin, Donghyuk Lee, Saugata Ghose, Mohammed Alser, Hasan Hassan, Oguz Ergin, Can Alkan, and Onur Mutlu,
 "GRIM-Filter: Fast Seed Location Filtering in DNA Read Mapping Using Processing-in-Memory Technologies"
 <u>BMC Genomics</u>, 2018.

Proceedings of the <u>16th Asia Pacific Bioinformatics Conference</u> (**APBC**), Yokohama, Japan, January 2018. <u>arxiv.org Version (pdf)</u>

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*}

From The Sixteenth Asia Pacific Bioinformatics Conference 2018 Yokohama, Japan. 15-17 January 2018

SAFARI



Genome Read In-Memory (GRIM) Filter:

Fast Seed Location Filtering in DNA Read Mapping using Processing-in-Memory Technologies

Jeremie Kim,

Damla Senol, Hongyi Xin, Donghyuk Lee, Saugata Ghose, Mohammed Alser, Hasan Hassan, Oguz Ergin, Can Alkan, and Onur Mutlu







TOBB UNIVERSITY OF ECONOMICS AND TECHNOLOGY



Executive Summary

- Genome Read Mapping is a very important problem and is the first step in many types of genomic analysis
 Could lead to improved health care, medicine, quality of life
 - Could lead to improved health care, medicine, quality of life
- Read mapping is an **approximate string matching** problem
 - □ Find the best fit of 100 character strings into a 3 billion character dictionary
 - Alignment is currently the best method for determining the similarity between two strings, but is very expensive
- We propose an in-memory processing algorithm GRIM-Filter for accelerating read mapping, by reducing the number of required alignments
- We implement GRIM-Filter using in-memory processing within 3Dstacked memory and show up to 3.7x speedup.

Google Workloads for Consumer Devices: Mitigating Data Movement Bottlenecks

Amirali Boroumand

Saugata Ghose, Youngsok Kim, Rachata Ausavarungnirun, Eric Shiu, Rahul Thakur, Daehyun Kim, Aki Kuusela, Allan Knies, Parthasarathy Ranganathan, Onur Mutlu









SEOUL NATIONAL UNIVERSITY





Open Problems: PIM Adoption

Enabling the Adoption of Processing-in-Memory: Challenges, Mechanisms, Future Research Directions

SAUGATA GHOSE, KEVIN HSIEH, AMIRALI BOROUMAND, RACHATA AUSAVARUNGNIRUN

Carnegie Mellon University

ONUR MUTLU ETH Zürich and Carnegie Mellon University

Saugata Ghose, Kevin Hsieh, Amirali Boroumand, Rachata Ausavarungnirun, Onur Mutlu, "Enabling the Adoption of Processing-in-Memory: Challenges, Mechanisms, Future Research Directions" Invited Book Chapter, to appear in 2018. [Preliminary arxiv.org version]

PIM Review and Open Problems

Processing Data Where It Makes Sense: Enabling In-Memory Computation

Onur Mutlu^{a,b}, Saugata Ghose^b, Juan Gómez-Luna^a, Rachata Ausavarungnirun^{b,c}

^aETH Zürich ^bCarnegie Mellon University ^cKing Mongkut's University of Technology North Bangkok

Onur Mutlu, Saugata Ghose, Juan Gomez-Luna, and Rachata Ausavarungnirun, "Processing Data Where It Makes Sense: Enabling In-Memory Computation" Invited paper in Microprocessors and Microsystems (MICPRO), June 2019. [arXiv version]

https://arxiv.org/pdf/1903.03988.pdf

SAFAR

PIM Review and Open Problems (II)

A Workload and Programming Ease Driven Perspective of Processing-in-Memory

Saugata Ghose†Amirali Boroumand†Jeremie S. Kim†§Juan Gómez-Luna§Onur Mutlu§††Carnegie Mellon University§ETH Zürich

Saugata Ghose, Amirali Boroumand, Jeremie S. Kim, Juan Gomez-Luna, and Onur Mutlu, "Processing-in-Memory: A Workload-Driven Perspective" *Invited Article in IBM Journal of Research & Development, Special Issue on Hardware for Artificial Intelligence*, to appear in November 2019. [Preliminary arXiv version]

SAFARI

https://arxiv.org/pdf/1907.12947.pdf



- Major Trends Affecting Main Memory
- The Need for Intelligent Memory Controllers
 - Bottom Up: Push from Circuits and Devices
 - Top Down: Pull from Systems and Applications
- Processing in Memory: Two Directions
 - Minimally Changing Memory Chips
 - Exploiting 3D-Stacked Memory
- How to Enable Adoption of Processing in Memory
- Conclusion

Challenge and Opportunity for Future

Fundamentally **Energy-Efficient** (Data-Centric) **Computing Architectures**

Challenge and Opportunity for Future

Fundamentally Low-Latency (Data-Centric) **Computing Architectures**

Challenge and Opportunity for Future

Computing Architectures with

Minimal Data Movement





Main Memory Needs Intelligent Controllers



Concluding Remarks

A Quote from A Famous Architect

 "architecture [...] based upon principle, and not upon precedent"



Precedent-Based Design?

"architecture [...] based upon principle, and not upon precedent"



Principled Design

"architecture [...] based upon principle, and not upon precedent"





The Overarching Principle

Organic architecture

From Wikipedia, the free encyclopedia

Organic architecture is a philosophy of architecture which promotes harmony between human habitation and the natural world through design approaches so sympathetic and well integrated with its site, that buildings, furnishings, and surroundings become part of a unified, interrelated composition.

A well-known example of organic architecture is Fallingwater, the residence Frank Lloyd Wright designed for the Kaufmann family in rural Pennsylvania. Wright had many choices to locate a home on this large site, but chose to place the home directly over the waterfall and creek creating a close, yet noisy dialog with the rushing water and the steep site. The horizontal striations of stone masonry with daring cantilevers of colored beige concrete blend with native rock outcroppings and the wooded environment.

Another Example: Precedent-Based Design



Principled Design



Another Principled Design



Source: By Martín Gómez Tagle - Lisbon, Portugal, CC BY-SA 3.0, https://commons.wikimedia.org/w/index.php?curid=13764903 Source: http://www.arcspace.com/exhibitions/unsorted/santiago-calatrava/
Another Principled Design



Principle Applied to Another Structure



The Overarching Principle

Zoomorphic architecture

From Wikipedia, the free encyclopedia

Zoomorphic architecture is the practice of using animal forms as the inspirational basis and blueprint for architectural design. "While animal forms have always played a role adding some of the deepest layers of meaning in architecture, it is now becoming evident that a new strand of biomorphism is emerging where the meaning derives not from any specific representation but from a more general allusion to biological processes."^[1]

Some well-known examples of Zoomorphic architecture can be found in the TWA Flight Center building in New York City, by Eero Saarinen, or the Milwaukee Art Museum by Santiago Calatrava, both inspired by the form of a bird's wings.^[3]

Overarching Principle for Computing?



Source: http://spectrum.ieee.org/image/MjYzMzAyMg.jpeg

Concluding Remarks

- It is time to design principled system architectures to solve the memory problem
- Design complete systems to be balanced, high-performance, and energy-efficient, i.e., data-centric (or memory-centric)
- Enable computation capability inside and close to memory
- This can
 - Lead to orders-of-magnitude improvements
 - Enable new applications & computing platforms
 - **D** Enable better understanding of nature

The Future of Processing in Memory 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

We Need to Revisit the Entire Stack

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

We can get there step by step

SAFARI

If In Doubt, See Other Doubtful Technologies

- 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

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

PIM Review and Open Problems

Processing Data Where It Makes Sense: Enabling In-Memory Computation

Onur Mutlu^{a,b}, Saugata Ghose^b, Juan Gómez-Luna^a, Rachata Ausavarungnirun^{b,c}

^aETH Zürich ^bCarnegie Mellon University ^cKing Mongkut's University of Technology North Bangkok

Onur Mutlu, Saugata Ghose, Juan Gomez-Luna, and Rachata Ausavarungnirun, "Processing Data Where It Makes Sense: Enabling In-Memory Computation" Invited paper in Microprocessors and Microsystems (MICPRO), June 2019. [arXiv version]

https://arxiv.org/pdf/1903.03988.pdf

SAFAR

PIM Review and Open Problems (II)

A Workload and Programming Ease Driven Perspective of Processing-in-Memory

Saugata Ghose†Amirali Boroumand†Jeremie S. Kim†§Juan Gómez-Luna§Onur Mutlu§††Carnegie Mellon University§ETH Zürich

Saugata Ghose, Amirali Boroumand, Jeremie S. Kim, Juan Gomez-Luna, and Onur Mutlu, "Processing-in-Memory: A Workload-Driven Perspective" *Invited Article in IBM Journal of Research & Development, Special Issue on Hardware for Artificial Intelligence*, to appear in November 2019. [Preliminary arXiv version]

SAFARI

https://arxiv.org/pdf/1907.12947.pdf

Acknowledgments

My current and past students and postdocs

Rachata Ausavarungnirun, Abhishek Bhowmick, Amirali Boroumand, Rui Cai, Yu Cai, Kevin Chang, Saugata Ghose, Kevin Hsieh, Tyler Huberty, Ben Jaiyen, Samira Khan, Jeremie Kim, Yoongu Kim, Yang Li, Jamie Liu, Lavanya Subramanian, Donghyuk Lee, Yixin Luo, Justin Meza, Gennady Pekhimenko, Vivek Seshadri, Lavanya Subramanian, Nandita Vijaykumar, HanBin Yoon, Jishen Zhao, ...

My collaborators

 Can Alkan, Chita Das, Phil Gibbons, Sriram Govindan, Norm Jouppi, Mahmut Kandemir, Mike Kozuch, Konrad Lai, Ken Mai, Todd Mowry, Yale Patt, Moinuddin Qureshi, Partha Ranganathan, Bikash Sharma, Kushagra Vaid, Chris Wilkerson, ...

Funding Acknowledgments

- Alibaba, AMD, Google, Facebook, HP Labs, Huawei, IBM, Intel, Microsoft, Nvidia, Oracle, Qualcomm, Rambus, Samsung, Seagate, VMware
- NSF
- NIH
- GSRC
- SRC
- CyLab

Acknowledgments

SAFARI Research Group safari.ethz.ch

Think BIG, Aim HIGH! https://safari.ethz.ch

Processing Data Where It Makes Sense in Modern Computing Systems: Enabling In-Memory Computation

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

16 August 2019

APPT Keynote



ETH zürich



Slides Not Covered But Could Be Useful

Readings, Videos, Reference Materials

Accelerated Memory Course (~6.5 hours)

ACACES 2018

- Memory Systems and Memory-Centric Computing Systems
- Taught by Onur Mutlu July 9-13, 2018
- ~6.5 hours of lectures
- Website for the Course including Videos, Slides, Papers
 - https://safari.ethz.ch/memory_systems/ACACES2018/
 - https://www.youtube.com/playlist?list=PL5Q2soXY2Zi-HXxomthrpDpMJm05P6J9x
- All Papers are at:
 - <u>https://people.inf.ethz.ch/omutlu/projects.htm</u>
 - Final lecture notes and readings (for all topics)

Longer Memory Course (~18 hours)

Tu Wien 2019

- Memory Systems and Memory-Centric Computing Systems
- Taught by Onur Mutlu June 12-19, 2019
- ~18 hours of lectures
- Website for the Course including Videos, Slides, Papers
 - https://safari.ethz.ch/memory_systems/TUWien2019
 - https://www.youtube.com/playlist?list=PL5Q2soXY2Zi_gntM55
 VoMIKIw7YrXOhbl
- All Papers are at:
 - https://people.inf.ethz.ch/omutlu/projects.htm
 - Final lecture notes and readings (for all topics)

Some Overview Talks

https://www.youtube.com/watch?v=kgiZISOcGFM&list=PL5Q2soXY2Zi8D_5MGV6EnXEJHnV2YFBJI

Future Computing Architectures

- https://www.youtube.com/watch?v=kgiZlSOcGFM&list=PL5Q2soXY2Zi8D_5MG V6EnXEJHnV2YFBJl&index=1
- Enabling In-Memory Computation
 - https://www.youtube.com/watch?v=oHqsNbxgdzM&list=PL5Q2soXY2Zi8D_5M GV6EnXEJHnV2YFBJl&index=7

Accelerating Genome Analysis

<u>https://www.youtube.com/watch?v=hPnSmfwu2-</u> <u>A&list=PL5Q2soXY2Zi8D_5MGV6EnXEJHnV2YFBJl&index=9</u>

Rethinking Memory System Design

https://www.youtube.com/watch?v=F7xZLNMIY1E&list=PL5Q2soXY2Zi8D_5MG V6EnXEJHnV2YFBJl&index=3

Reference Overview Paper I

ΔΓΔΠ

Processing Data Where It Makes Sense: Enabling In-Memory Computation

Onur Mutlu^{a,b}, Saugata Ghose^b, Juan Gómez-Luna^a, Rachata Ausavarungnirun^{b,c}

^aETH Zürich ^bCarnegie Mellon University ^cKing Mongkut's University of Technology North Bangkok

Onur Mutlu, Saugata Ghose, Juan Gomez-Luna, and Rachata Ausavarungnirun, "Processing Data Where It Makes Sense: Enabling In-Memory Computation" Invited paper in Microprocessors and Microsystems (MICPRO), June 2019. [arXiv version]

https://arxiv.org/pdf/1903.03988.pdf

Reference Overview Paper II

Enabling the Adoption of Processing-in-Memory: Challenges, Mechanisms, Future Research Directions

SAUGATA GHOSE, KEVIN HSIEH, AMIRALI BOROUMAND, RACHATA AUSAVARUNGNIRUN

Carnegie Mellon University

ONUR MUTLU ETH Zürich and Carnegie Mellon University

Saugata Ghose, Kevin Hsieh, Amirali Boroumand, Rachata Ausavarungnirun, Onur Mutlu, "Enabling the Adoption of Processing-in-Memory: Challenges, Mechanisms, Future Research Directions" Invited Book Chapter, to appear in 2018. [Preliminary arxiv.org version]

SAFARI

Reference Overview Paper III

 Onur Mutlu and Lavanya Subramanian, <u>"Research Problems and Opportunities in Memory</u> <u>Systems"</u> *Invited Article in <u>Supercomputing Frontiers and Innovations</u> (SUPERFRI), 2014/2015.*

Research Problems and Opportunities in Memory Systems

Onur Mutlu¹, Lavanya Subramanian¹

https://people.inf.ethz.ch/omutlu/pub/memory-systems-research_superfri14.pdf

Reference Overview Paper IV

Onur Mutlu, **"The RowHammer Problem and Other Issues We May Face as Memory Becomes Denser"** *Invited Paper in Proceedings of the <u>Design, Automation, and Test in</u> <i>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

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

Reference Overview Paper V

 Onur Mutlu, <u>"Memory Scaling: A Systems Architecture</u> <u>Perspective"</u> *Technical talk at <u>MemCon 2013</u> (MEMCON)*, Santa Clara, CA, August 2013. [Slides (pptx) (pdf)] [Video] [Coverage on StorageSearch]

Memory Scaling: A Systems Architecture Perspective

Onur Mutlu Carnegie Mellon University onur@cmu.edu http://users.ece.cmu.edu/~omutlu/

https://people.inf.ethz.ch/omutlu/pub/memory-scaling_memcon13.pdf

Reference Overview Paper VI



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

Reference Overview Paper VII

Onur Mutlu and Jeremie Kim,
 "RowHammer: A Retrospective"
 <u>IEEE Transactions on Computer-Aided Design of Integrated</u>
 <u>Circuits and Systems</u> (TCAD) Special Issue on Top Picks in
 Hardware and Embedded Security, 2019.
 [Preliminary arXiv version]

RowHammer: A Retrospective

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

Reference Overview Paper VIII

A Workload and Programming Ease Driven Perspective of Processing-in-Memory

Saugata Ghose†Amirali Boroumand†Jeremie S. Kim†§Juan Gómez-Luna§Onur Mutlu§††Carnegie Mellon University§ETH Zürich

Saugata Ghose, Amirali Boroumand, Jeremie S. Kim, Juan Gomez-Luna, and Onur Mutlu, "Processing-in-Memory: A Workload-Driven Perspective" *Invited Article in <u>IBM Journal of Research & Development</u>, Special Issue on Hardware for Artificial Intelligence*, to appear in November 2019. [Preliminary arXiv version]

SAFARI

https://arxiv.org/pdf/1907.12947.pdf

Related Videos and Course Materials (I)

- <u>Undergraduate Computer Architecture Course Lecture</u> <u>Videos (2015, 2014, 2013)</u>
- <u>Undergraduate Computer Architecture Course</u> <u>Materials</u> (2015, 2014, 2013)
- Graduate Computer Architecture Course Lecture Videos (2018, 2017, 2015, 2013)
- Graduate Computer Architecture Course Materials (2018, 2017, 2015, 2013)
- Parallel Computer Architecture Course Materials (Lecture Videos)

Related Videos and Course Materials (II)

- Freshman Digital Circuits and Computer Architecture Course Lecture Videos (2018, 2017)
- Freshman Digital Circuits and Computer Architecture Course Materials (2018)
- <u>Memory Systems Short Course Materials</u> (<u>Lecture Video on Main Memory and DRAM Basics</u>)

Some Open Source Tools (I)

- Rowhammer Program to Induce RowHammer Errors
 - <u>https://github.com/CMU-SAFARI/rowhammer</u>
- Ramulator Fast and Extensible DRAM Simulator
 - https://github.com/CMU-SAFARI/ramulator
- MemSim Simple Memory Simulator
 - https://github.com/CMU-SAFARI/memsim
- NOCulator Flexible Network-on-Chip Simulator
 - <u>https://github.com/CMU-SAFARI/NOCulator</u>
- SoftMC FPGA-Based DRAM Testing Infrastructure
 - https://github.com/CMU-SAFARI/SoftMC
- Other open-source software from my group
 - https://github.com/CMU-SAFARI/

<u>http://www.ece.cmu.edu/~safari/tools.html</u>
SAFARI

Some Open Source Tools (II)

- MQSim A Fast Modern SSD Simulator
 - <u>https://github.com/CMU-SAFARI/MQSim</u>
- Mosaic GPU Simulator Supporting Concurrent Applications
 - https://github.com/CMU-SAFARI/Mosaic
- IMPICA Processing in 3D-Stacked Memory Simulator
 - https://github.com/CMU-SAFARI/IMPICA
- SMLA Detailed 3D-Stacked Memory Simulator
 - https://github.com/CMU-SAFARI/SMLA
- HWASim Simulator for Heterogeneous CPU-HWA Systems
 <u>https://github.com/CMU-SAFARI/HWASim</u>
- Other open-source software from my group
 - https://github.com/CMU-SAFARI/

<u>http://www.ece.cmu.edu/~safari/tools.html</u>
SAFARI

More Open Source Tools (III)

- A lot more open-source software from my group
 - https://github.com/CMU-SAFARI/
 - http://www.ece.cmu.edu/~safari/tools.html

SAFARI Research Group at ETH Zurich and Carnegie Mellon University							
Site for source code and tools distribution from SAFARI Research Group at ETH Zurich © ETH Zurich and Carnegi © http://www.ece.cmu.ed 🖂 omutlu@gmail.com Repositories 30	ch and Carnegie Mellon University.						
Search repositories Type: All Language: All	Customize pinned repositories						
MQSim MQSim is a fast and accurate simulator modeling the performance of modern multi-queue (MQ) SSDs as well as traditional SATA based SSDs. MQSim faithfully models new high-bandwidth protocol implementations, steady-state SSD conditions, and the full end-to-end latency of	Top languages ● C++ ● C ● C# ● AGS Script ● Verilog						
equests in modern SSDs. It is described in detail in the FAST 2018 paper by A ● C++ ★ 14	Most used topics Manage						

SAFARI



All are available at

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

http://scholar.google.com/citations?user=7XyGUGkAAAAJ&hl=en

https://people.inf.ethz.ch/omutlu/acaces2018.html

Ramulator: A Fast and Extensible DRAM Simulator [IEEE Comp Arch Letters'15]

Ramulator Motivation

- DRAM and Memory Controller landscape is changing
- Many new and upcoming standards
- Many new controller designs
- A fast and easy-to-extend simulator is very much needed

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				

Ramulator

- Provides out-of-the box support for many DRAM standards:
 - DDR3/4, LPDDR3/4, GDDR5, WIO1/2, HBM, plus new proposals (SALP, AL-DRAM, TLDRAM, RowClone, and SARP)
- ~2.5X faster than fastest open-source simulator
- Modular and extensible to different standards

Simulator	ulator Cycles (10 ⁶)		Runtime (sec.)		<i>Req/sec</i> (10 ³)		Memory	
(clang -03)	Random	Stream	Random	Stream	Random	Stream	(<i>MB</i>)	
Ramulator	652	411	752	249	133	402	2.1	
DRAMSim2	645	413	2,030	876	49	114	1.2	
USIMM	661	409	1,880	750	53	133	4.5	
DrSim	647	406	18,109	12,984	6	8	1.6	
NVMain	666	413	6,881	5,023	15	20	4,230.0	

Table 3. Comparison of five simulators using two traces
Case Study: Comparison of DRAM Standards

Star	ndard	Rate (MT/s)	Timing (CL-RCD-RP)	Data-Bus (Width×Chan.)	Rank-per-Chan	BW (GB/s)
DD	R3	1,600	11-11-11	64 -bit $\times 1$	1	11.9
DD	R4	2,400	16-16-16	64 -bit $\times 1$	1	17.9
SAI	P^{\dagger}	1,600	11-11-11	64 -bit $\times 1$	1	11.9
LPI	DDR3	1,600	12 - 15 - 15	64 -bit $\times 1$	1	11.9
LPI	DDR4	2,400	22-22-22	32 -bit $\times 2^*$	1	17.9
GD	DR5 [12]	6,000	18-18-18	64 -bit $\times 1$	1	44.7
HB	M	1,000	7-7-7	128 -bit $\times 8^*$	1	119.2
WIG)	266	7-7-7	128 -bit $\times 4^*$	1	15.9
WIG	02	1,066	9-10-10	128-bit \times 8*	1	127.2



Ramulator Paper and Source Code

- Yoongu Kim, Weikun Yang, and Onur Mutlu,
 "Ramulator: A Fast and Extensible DRAM Simulator"
 <u>IEEE Computer Architecture Letters</u> (CAL), March 2015.
 [Source Code]
- Source code is released under the liberal MIT License
 <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