Future of Computer Architecture and Hardware Security

Onur Mutlu <u>omutlu@gmail.com</u> <u>https://people.inf.ethz.ch/omutlu</u> 16 May 2024

Qualcomm Product Security Summit Keynote Talk







Computer Architecture Today

- What is it and where it is going
- Three Major Hardware Issues That Affect Security
 - Technology scaling problems
 - □ Growing system complexity; old methods not keeping up
 - New architectures and technologies

Why Do We Do Computing?



To Solve Problems



To Gain Insight

SAFARI Hamming, "Numerical Methods for Scientists and Engineers," 1962. ⁵

To Enable a Better Life & Future

How Does a Computer Solve Problems?



Orchestrating Electrons

In today's dominant technologies



How Do Problems Get Solved by Electrons?

The Transformation Hierarchy

Computer Architecture (expanded view)



Computer Architecture (narrow view)

Computer Architecture

- is the science and art of designing computing platforms (hardware, interface, system SW, and programming model)
- to achieve a set of design goals
 - □ E.g., highest performance on earth on workloads X, Y, Z
 - E.g., longest battery life at a form factor that fits in your pocket with cost < \$\$\$ CHF
 - E.g., best average performance across all known workloads at the best performance/cost ratio

• ...

Designing a supercomputer is different from designing a smartphone → But, many fundamental principles are similar







SAFARI so

Source: https://taxistartup.com/wp-content/uploads/2015/03/UK-Self-Driving-Cars.jpg



Mohammed Alser, Zülal Bingöl, Damla Senol Cali, Jeremie Kim, Saugata Ghose, Can Alkan, Onur Mutlu <u>"Accelerating Genome Analysis: A Primer on an Ongoing Journey"</u> IEEE Micro, August 2020.



FPGA-Based Near-Memory Acceleration of Modern Data-Intensive Applications

July-Aug. 2021, pp. 39-48, vol. 41 DOI Bookmark: 10.1109/MM.2021.3088396

MinION from ONT

SmidgION from ONT

An Example System in Your Pocket



SAFARI

https://www.gsmarena.com/apple_announces_m1_ultra_with_20core_cpu_and_64core_gpu-news-53481.php









Control

Figure 3. TPU Printed Circuit Board. It can be inserted in the slot for an SATA disk in a server, but the card uses PCIe Gen3 x16.

Figure 4. Systolic data flow of the Matrix Multiply Unit. Software has the illusion that each 256B input is read at once, and they instantly update one location of each of 256 accumulator RAMs.

Jouppi et al., "In-Datacenter Performance Analysis of a Tensor Processing Unit", ISCA 2017.



New ML applications (vs. TPU3):

- Computer vision
- Natural Language Processing (NLP)
- Recommender system
- Reinforcement learning that plays Go

250 TFLOPS per chip in 2021 vs 90 TFLOPS in TPU3



https://spectrum.ieee.org/tech-talk/computing/hardware/heres-how-googles-tpu-v4-ai-chip-stacked-up-in-training-tests

- ML accelerator: 260 mm², 6 billion transistors, 600 GFLOPS GPU, 12 ARM 2.2 GHz CPUs.
- Two redundant chips for better safety.





Tesla Dojo Chip & System



D1 Chip

362 TFLOPs BF16/CFP8 22.6 TFLOPs FP32

10TBps/dir. On-Chip Bandwidth 4TBps/edge. Off-Chip Bandwidth

400W TDP





50 Billion Transistors

11+ Miles Of Wires

1:53:07 / 3:03:20 · Dojo >

https://www.youtube.com/watch?v=j0z4FweCy4M&t=6340s

Tesla Dojo Chip & System





https://www.youtube.com/watch?v=j0z4FweCy4M&t=6340s

Tesla Dojo Chip & System







NVIDIA is claiming a **7x improvement** in dynamic programming algorithm (**DPX instructions**) performance on a single H100 versus naïve execution on an A100.

https://www.nvidia.com/en-us/data-center/h100/



27

Cerebras's Wafer Scale Engine (2019)



 The largest ML accelerator chip

400,000 cores



Cerebras WSE 1.2 Trillion transistors 46,225 mm² Largest GPU 21.1 Billion transistors 815 mm²

https://www.anandtech.com/show/14758/hot-chips-31-live-blogs-cerebras-wafer-scale-deep-learning

https://www.cerebras.net/cerebras-wafer-scale-engine-why-we-need-big-chips-for-deep-learning⁸

Cerebras's Wafer Scale Engine-2 (2021)



 The largest ML accelerator chip (2021)

850,000 cores



Cerebras WSE-2 2.6 Trillion transistors 46,225 mm² Largest GPU 54.2 Billion transistors 826 mm²

NVIDIA Ampere GA100

https://www.anandtech.com/show/14758/hot-chips-31-live-blogs-cerebras-wafer-scale-deep-learning

https://www.cerebras.net/cerebras-wafer-scale-engine-why-we-need-big-chips-for-deep-learning/

Many (Other) AI/ML Chips (2021)



All information contained within this infographic is gathered from the internet and periodically updated, no guarantee is given that the information provided is correct, complete, and up-to-date.

SAFARI

https://basicmi.github.io/AI-Chip/



To achieve the highest efficiency, performance, robustness:

we must take the expanded view

of computer architecture



What Limits Us in Computing Today?

Increasingly Demanding Applications

Dream...

and, they will come

As applications push boundaries, computing platforms become increasingly strained

Many Metrics to Optimize for

- Performance
- Energy/Power
- Correctness
- Robustness
- Cost
- Programming Ease
- Ease of Use
- Scalability
- Simplicity (Complexity)

Challenging especially with complex systems & hardware

Three Major Limiters to Computing

- Technology scaling is not going well
- System complexity is increasing; old methods not keeping up
- Processor-centric designs are not keeping up

- These affect all metrics we care about
- These have fundamental impact on security and how we build secure systems

Technology Scaling
Technology Scaling Problems

- Circuit size and energy reduction has enabled continuous innovation at all levels of the computing stack
- As circuits become smaller, they become less reliable
- More flaky circuits are a problem for robust (reliable, safe, secure) operation
- If circuits produce wrong results, security can be affected (along with safety, reliability, availability)

How Reliable/Secure/Safe is This Bridge?





Collapse of the "Galloping Gertie"





Another View





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

SAFARI

How Safe & Secure Is This Platform?



SAFARI Source: https://taxistartup.com/wp-content/uploads/2015/03/UK-Self-Driving-Cars.jpg

How Robust Are These Platforms?









SAFARI https://www.kennedyspacecenter.com/explore-attractions/nasa-now https://www.cnet.com/pictures/nasas-wildest-rides-extreme-vehicles-for-earth-and-beyond/7/

Challenge and Opportunity for Future

Robust (Reliable, Secure, Safe)

An Example: The RowHammer Problem

- One can predictably induce bit flips in commodity DRAM chips
 All recent DRAM chips are fundamentally vulnerable
- First example of how a simple hardware failure mechanism can create a widespread system security vulnerability



A Curious Phenomenon [Kim et al., ISCA 2014]

One can predictably induce errors in DRAM memory chips

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



Modern Memory 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

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

Recent DRAM Is More Vulnerable



All modules from 2012–2013 are vulnerable

Higher-Level Implications

This simple circuit level failure mechanism has enormous implications on upper layers of the transformation hierarchy







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



Download from: https://github.com/CMU-SAFARI/rowhammer



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, safety 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

Many RowHammer Security Exploits

- One can exploit RowHammer to
- Take over a system
- Read data they do not have access to
- Break out of virtual machine sandboxes
- Corrupt important data \rightarrow render ML inference useless
- Steal secret data (e.g., crypto keys & ML model parameters)

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

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)



SAFARI

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



SoftMC: Open Source DRAM Infrastructure

 Hasan Hassan, Nandita Vijaykumar, Samira Khan, Saugata Ghose, Kevin Chang, Gennady Pekhimenko, Donghyuk Lee, Oguz Ergin, and Onur Mutlu, "SoftMC: A Flexible and Practical Open-Source Infrastructure for Enabling Experimental DRAM Studies" Proceedings of the 23rd International Symposium on High-Performance Computer Architecture (HPCA), Austin, TX, USA, February 2017.
 [Slides (pptx) (pdf)] [Lightning Session Slides (pptx) (pdf)] [Full Talk Lecture (39 minutes)]

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

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

DRAM Bender: New DRAM Infrastructure

 Ataberk Olgun, Hasan Hassan, A Giray Yağlıkçı, Yahya Can Tuğrul, Lois Orosa, Haocong Luo, Minesh Patel, Oğuz Ergin, and Onur Mutlu, "DRAM Bender: An Extensible and Versatile FPGA-based Infrastructure to Easily Test State-of-the-art DRAM Chips" *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems (TCAD)*, 2023.
 [Extended arXiv version]
 [DRAM Bender Source Code]
 [DRAM Bender Tutorial Video (43 minutes)]

DRAM Bender: An Extensible and Versatile FPGA-based Infrastructure to Easily Test State-of-the-art DRAM Chips

Ataberk Olgun§Hasan Hassan§A. Giray Yağlıkçı§Yahya Can Tuğrul§†Lois Orosa§⊙Haocong Luo§Minesh Patel§Oğuz Ergin†Onur Mutlu§§ETH Zürich†TOBB ETÜ⊙Galician Supercomputing Center

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

DRAM Bender: FPGA Prototypes

Testing Infrastructure	Protocol Support	FPGA Support
SoftMC [134]	DDR3	One Prototype
LiteX RowHammer Tester (LRT) [17]	DDR3/4, LPDDR4	Two Prototypes
DRAM Bender (this work)	DDR3/DDR4	Five Prototypes

Five out of the box FPGA-based prototypes

https://github.com/CMU-SAFARI/DRAM-Bender



SAFARI



66

RowHammer [ISCA 2014]

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 Architecture</u> (ISCA), Minneapolis, MN, June 2014.
 [Slides (pptx) (pdf)] [Lightning Session Slides (pptx) (pdf)] [Source Code and Data] [Lecture Video (1 hr 49 mins), 25 September 2020]
 One of the 7 papers of 2012-2017 selected as Top Picks in Hardware and Embedded Security for IEEE TCAD (link).
 Selected to the ISCA-50 25-Year Retrospective Issue covering 1996-2020 in 2023 (Retrospective (pdf) Full Issue).

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

SAFARI

Memory Scaling Issues Are Real

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]
 [Slides from COSADE 2019 (pptx)]
 [Slides from VLSI-SOC 2020 (pptx) (pdf)]
 [Talk Video (1 hr 15 minutes, with Q&A)]

RowHammer: A Retrospective

Onur Mutlu§‡Jeremie S. Kim‡§§ETH Zürich‡Carnegie Mellon University

Memory Scaling Issues Are Real

 Onur Mutlu, Ataberk Olgun, and A. Giray Yaglikci, <u>"Fundamentally Understanding and Solving RowHammer"</u> *Invited Special Session Paper at the <u>28th Asia and South Pacific Design</u> <u>Automation Conference (ASP-DAC)</u>, Tokyo, Japan, January 2023. [arXiv version] [Slides (pptx) (pdf)] [Talk Video (26 minutes)]*

Fundamentally Understanding and Solving RowHammer

Onur Mutlu onur.mutlu@safari.ethz.ch ETH Zürich Zürich, Switzerland Ataberk Olgun ataberk.olgun@safari.ethz.ch ETH Zürich Zürich, Switzerland A. Giray Yağlıkcı giray.yaglikci@safari.ethz.ch ETH Zürich Zürich, Switzerland

https://arxiv.org/pdf/2211.07613.pdf

The Push from Circuits and Devices

Main Memory Needs Intelligent Controllers

Industry's Intelligent DRAM Controllers (I)

ISSCC 2023 / SESSION 28 / HIGH-DENSITY MEMORIES

28.8 A 1.1V 16Gb DDR5 DRAM with Probabilistic-Aggressor Tracking, Refresh-Management Functionality, Per-Row Hammer Tracking, a Multi-Step Precharge, and Core-Bias Modulation for Security and Reliability Enhancement

Woongrae Kim, Chulmoon Jung, Seongnyuh Yoo, Duckhwa Hong, Jeongjin Hwang, Jungmin Yoon, Ohyong Jung, Joonwoo Choi, Sanga Hyun, Mankeun Kang, Sangho Lee, Dohong Kim, Sanghyun Ku, Donhyun Choi, Nogeun Joo, Sangwoo Yoon, Junseok Noh, Byeongyong Go, Cheolhoe Kim, Sunil Hwang, Mihyun Hwang, Seol-Min Yi, Hyungmin Kim, Sanghyuk Heo, Yeonsu Jang, Kyoungchul Jang, Shinho Chu, Yoonna Oh, Kwidong Kim, Junghyun Kim, Soohwan Kim, Jeongtae Hwang, Sangil Park, Junphyo Lee, Inchul Jeong, Joohwan Cho, Jonghwan Kim

SK hynix Semiconductor, Icheon, Korea



Industry's Intelligent DRAM Controllers (II)

SK hynix Semiconductor, Icheon, Korea

DRAM products have been recently adopted in a wide range of high-performance computing applications: such as in cloud computing, in big data systems, and IoT devices. This demand creates larger memory capacity requirements, thereby requiring aggressive DRAM technology node scaling to reduce the cost per bit [1,2]. However, DRAM manufacturers are facing technology scaling challenges due to row hammer and refresh retention time beyond 1a-nm [2]. Row hammer is a failure mechanism, where repeatedly activating a DRAM row disturbs data in adjacent rows. Scaling down severely threatens reliability since a reduction of DRAM cell size leads to a reduction in the intrinsic row hammer tolerance [2,3]. To improve row hammer tolerance, there is a need to probabilistically activate adjacent rows with carefully sampled active addresses and to improve intrinsic row hammer tolerance [2]. In this paper, row-hammer-protection and refresh-management schemes are presented to guarantee DRAM security and reliability despite the aggressive scaling from 1a-nm to sub 10-nm nodes. The probabilisticaggressor-tracking scheme with a refresh-management function (RFM) and per-row hammer tracking (PRHT) improve DRAM resilience. A multi-step precharge reinforces intrinsic row-hammer tolerance and a core-bias modulation improves retention time: even in the face of cell-transistor degradation due to technology scaling. This comprehensive scheme leads to a reduced probability of failure, due to row hammer attacks, by 93.1% and an improvement in retention time by 17%.
Industry's Intelligent DRAM Controllers (III)



ISSCC 2023 / SESSION 28 / HIGH-DENSITY MEMORIES /

28.8 A 1.1V 16Gb DDR5 DRAM with Probabilistic-Aggressor Tracking, Refresh-Management Functionality, Per-Row Hammer Tracking, a Multi-Step Precharge, and Core-Bias Modulation for Security and Reliability Enhancement

Woongrae Kim, Chulmoon Jung, Seongnyuh Yoo, Duckhwa Hong, Jeongjin Hwang, Jungmin Yoon, Ohyong Jung, Joonwoo Choi, Sanga Hyun, Mankeun Kang, Sangho Lee, Dohong Kim, Sanghyun Ku, Donhyun Choi, Nogeun Joo, Sangwoo Yoon, Junseok Noh, Byeongyong Go, Cheolhoe Kim, Sunil Hwang, Mihyun Hwang, Seol-Min Yi, Hyungmin Kim, Sanghyuk Heo, Yeonsu Jang, Kyoungchul Jang, Shinho Chu, Yoonna Oh, Kwidong Kim, Junghyun Kim, Soohwan Kim, Jeongtae Hwang, Sangil Park, Junphyo Lee, Inchul Jeong, Joohwan Cho, Jonghwan Kim

SK hynix Semiconductor, Icheon, Korea

SAFARI

DSAC: Low-Cost Rowhammer Mitigation Using In-DRAM Stochastic and Approximate Counting Algorithm

Seungki Hong Dongha Kim Jaehyung Lee Reum Oh Changsik Yoo Sangjoon Hwang Jooyoung Lee

DRAM Design Team, Memory Division, Samsung Electronics

https://arxiv.org/pdf/2302.03591v1.pdf

Panopticon: A Complete In-DRAM Rowhammer Mitigation

Tanj Bennett[§], Stefan Saroiu, Alec Wolman, and Lucian Cojocar Microsoft, [§]Avant-Gray LLC

https://stefan.t8k2.com/publications/dramsec/2021/panopticon.pdf

Improvements in JEDEC (2024)

JEDEC ® Global Standards for the Microelectronics Industry				
STANDARDS & DOCUMENTS	COMMITTEES	NEWS	EVENTS & MEETINGS	loc
DDR5 SDRAM	·	JESD79-50	C Apr 2024	
Release Number: Version 1.30				

Version 1.30

This standard defines the DDR5 SDRAM specification, including features, functionalities, AC and DC characteristics, packages, and ball/signal assignments. The purpose of this Standard is to define the minimum set of requirements for JEDEC compliant 8 Gb through 32 Gb for x4, x8, and x16 DDR5 SDRAM devices. This standard was created based on the DDR4 standards (JESD79-4) and some aspects of the DDR, DDR2, DDR3, and LPDDR4 standards (JESD79, JESD79-2, JESD79-3, and JESD209-4).

```
Committee(s): JC-42, JC-42.3
```

Are Solutions Good?



Appears at ISCA 2023

What if there is another phenomenon that **does NOT require high row activation count**?

RowPress: Amplifying Read-Disturbance in Modern DRAM Chips

Haocong Luo Ataberk Olgun A. Giray Yağlıkçı Yahya Can Tuğrul Steve Rhyner Meryem Banu Cavlak Joël Lindegger Mohammad Sadrosadati Onur Mutlu *ETH Zürich*

RowPress [ISCA 2023]



 Haocong Luo, Ataberk Olgun, Giray Yaglikci, Yahya Can Tugrul, Steve Rhyner, M. Banu Cavlak, Joel Lindegger, Mohammad Sadrosadati, and Onur Mutlu, "RowPress: Amplifying Read Disturbance in Modern DRAM Chips" Proceedings of the 50th International Symposium on Computer Architecture (ISCA), Orlando, FL, USA, June 2023.
 [Slides (pptx) (pdf)]
 [Lightning Talk Slides (pptx) (pdf)]
 [Lightning Talk Video (3 minutes)]
 [RowPress Source Code and Datasets (Officially Artifact Evaluated with All Badges)]
 Officially artifact evaluated as available, reusable and reproducible. Best artifact award at ISCA 2023. IEEE Micro Top Pick.

RowPress: Amplifying Read-Disturbance in Modern DRAM Chips

Haocong Luo Ataberk Olgun A. Giray Yağlıkçı Yahya Can Tuğrul Steve Rhyner Meryem Banu Cavlak Joël Lindegger Mohammad Sadrosadati Onur Mutlu *ETH Zürich*

A Recent RowHammer Lecture



Stanford Seminar - RowHammer, RowPress and Beyond: Can We Be Free of Bitflips (Soon)?

SAFARI



https://www.youtube.com/watch?v=0W7YRRhnunw

Emerging Memories Also Need Intelligent Controllers

Benjamin C. Lee, Engin Ipek, Onur Mutlu, and Doug Burger,
"Architecting Phase Change Memory as a Scalable DRAM Alternative"
Proceedings of the <u>36th International Symposium on Computer</u>
<u>Architecture</u> (ISCA), pages 2-13, Austin, TX, June 2009. <u>Slides (pdf)</u>
One of the 13 computer architecture papers of 2009 selected as Top
Picks by IEEE Micro. Selected as a CACM Research Highlight.
2022 Persistent Impact Prize.

Architecting Phase Change Memory as a Scalable DRAM Alternative

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

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

SAFARI

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



Intelligent Memory Controllers **Can Enhance Security** & Enable Better Scaling

Data Corruption is in CPU Logic, Too

- Intermittent defects can cause silent data corruption
- They may be hard to detect or replicate
- They may be exploitable

Silent Data Corruption in Logic (2021)

Silent Data Corruptions at Scale

Harish Dattatraya Dixit Facebook, Inc. hdd@fb.com Sneha Pendharkar Facebook, Inc. spendharkar@fb.com Matt Beadon Facebook, Inc. mbeadon@fb.com Chris Mason Facebook, Inc. clm@fb.com

Tejasvi Chakravarthy Facebook, Inc. teju@fb.com Bharath Muthiah Facebook, Inc. bharathm@fb.com Sriram Sankar Facebook Inc. sriramsankar@fb.com

Cores that don't count

Peter H. Hochschild Paul Turner Jeffrey C. Mogul Google Sunnyvale, CA, US Rama Govindaraju Parthasarathy Ranganathan Google Sunnyvale, CA, US David E. Culler Amin Vahdat Google Sunnyvale, CA, US

Silent Data Corruption In-the-Field (2021)



HotOS 2021: Cores That Don't Count (Fun Hardware)

https://www.youtube.com/watch?v=QMF3rqhjYuM

Silent Data Corruption in Logic (2023)

Understanding Silent Data Corruptions in a Large Production CPU Population

Shaobu Wang Tsinghua University Guangyan Zhang* Tsinghua University Junyu Wei Tsinghua University

Yang Wang The Ohio State University Jiesheng Wu Alibaba Cloud Qingchao Luo Alibaba Cloud

Understanding and Mitigating Hardware Failures in Deep Learning Training Accelerator Systems

Yi He University of Chicago Chciago, IL, USA yiizy@uchicago.edu

Robert de Gruijl Google Sunnyvale, CA, USA rdegruijl@google.com Mike Hutton Google Sunnyvale, CA, USA mdhutton@google.com

Rama Govindaraju Google Sunnyvale, CA, USA govindaraju@google.com

Yanjing Li University of Chicago Chciago, IL, USA yanjingl@uchicago.edu Steven Chan Google Sunnyvale, CA, USA scchan@google.com

Nishant Patil Google Sunnyvale, CA, USA nishantpatil@google.com



 Both memory and logic errors will become worse with technology scaling

Hardware errors will create worse robustness problems

We cannot afford to ignore data corruption

System Complexity

Complex Systems Cause Many Issues

- Many hardware components, complex components
- Harder to design & verify
- Harder to reason about operational behavior
 Correctness, performance, energy, security, privacy, ...
- Harder to control interactions between components and avoid information leakage
- Old methods do not keep up with new trends and complexity
 - Virtual memory is a prime example

SAFARI

CPU Complexity Is Growing

Moore's Law: The number of transistors on microchips doubles every two years

Moore's law describes the empirical regularity that the number of transistors on integrated circuits doubles approximately every two years. This advancement is important for other aspects of technological progress in computing – such as processing speed or the price of computers.



Data source: Wikipedia (wikipedia.org/wiki/Transistor_count) Year in Which the MicroChip V OurWorldinData.org – Research and data to make progress against the world's largest problems.

Complex CPUs and Memory Hierarchies



10nm ESF=Intel 7 Alder Lake die shot (~209mm²) from Intel: https://www.intel.com/content/www/us/en/newsroom/news/12th-gen-core-processors.html Die shot interpretation by Locuza, October 2021

Intel Alder Lake, 2021

SAFARI

Source: https://twitter.com/Locuza_/status/1454152714930331652

Complex CPUs and Memory Hierarchies



Core Count: 8 cores/16 threads

L1 Caches: 32 KB per core

L2 Caches: 512 KB per core

L3 Cache: 32 MB shared

AMD Ryzen 5000, 2020

Complexity Growing with 3D (2021)



134/comparing-zen-3-to-zen-2

AMD increases the L3 size of their 8-core Zen 3 processors from 32 MB to 96 MB

Additional 64 MB L3 cache die stacked on top of the processor die

- Connected using Through Silicon Vias (TSVs)
- Total of 96 MB L3 cache



CPU Complexity and Features

- Leads to many (endless) side and covert channels
 - Spectre and Meltdown are prime recent examples
 - These will not go away
- Leads to many bugs and unintended behavior
 - Especially with new features or complex interactions
 - Some can be exploitable

How to tame CPU complexity and resulting issues?

Access Control & Protection Mechanisms

- Are based on virtual memory (VM), invented in 1950s
- VM has not changed much even after decades of technology scaling and memory system improvements
- VM causes large performance problems and is responsible for large complexity, power, energy
- VM is poor for fine-grained security and access control
- VM hinders innovation in heterogeneous (accelerator) systems and new architectures (e.g., processing near data)

It is time to rethink virtual memory

Virtual Memory: Parting Thoughts

- Virtual Memory is one of the most successful examples of
 - architectural support for programmers
 - how to partition work between hardware and software
 - hardware/software cooperation
 - programmer/architect tradeoff
- Going forward: How does virtual memory fare and scale into the future? Five key trends:
 - Increasing, huge physical memory sizes (local & remote)
 - Hybrid physical memory systems (DRAM + NVM + SSD)
 - Many accelerators in the system accessing physical memory
 - Virtualized systems (hypervisors, software virtualization, local and remote memories)
 - Processing in memory systems near-data accelerators

Rethinking Virtual Memory

Nastaran Hajinazar, Pratyush Patel, Minesh Patel, Konstantinos Kanellopoulos, Saugata Ghose, Rachata Ausavarungnirun, Geraldo Francisco de Oliveira Jr., Jonathan Appavoo, Vivek Seshadri, and Onur Mutlu, "The Virtual Block Interface: A Flexible Alternative to the Conventional Virtual Memory Framework" *Proceedings of the <u>47th International Symposium on Computer Architecture</u> (<i>ISCA*), Virtual, June 2020. [Slides (pptx) (pdf)] [Lightning Talk Slides (pptx) (pdf)] [ARM Research Summit Poster (pptx) (pdf)] [Talk Video (26 minutes)] [Lightning Talk Video (3 minutes)]

The Virtual Block Interface: A Flexible Alternative to the Conventional Virtual Memory Framework

Nastaran Hajinazar^{*†} Pratyush Patel[⋈] Minesh Patel^{*} Konstantinos Kanellopoulos^{*} Saugata Ghose[‡] Rachata Ausavarungnirun[⊙] Geraldo F. Oliveira^{*} Jonathan Appavoo[◊] Vivek Seshadri[▽] Onur Mutlu^{*‡}

*ETH Zürich [†]Simon Fraser University \bowtie University of Washington [‡]Carnegie Mellon University \odot King Mongkut's University of Technology North Bangkok \diamond Boston University \bigtriangledown Microsoft Research India

SAFARI <u>https://people.inf.ethz.ch/omutlu/pub/VBI-virtual-block-interface_isca20.pdf</u>

Better Virtual Memory (I)

Konstantinos Kanellopoulos, Hong Chul Nam, F. Nisa Bostanci, Rahul Bera, Mohammad Sadrosadati, Rakesh Kumar, Davide Basilio Bartolini, and Onur Mutlu, "Victima: Drastically Increasing Address Translation Reach by Leveraging Underutilized Cache Resources" Proceedings of the <u>56th International Symposium on Microarchitecture</u> (MICRO), Toronto, ON, Canada, November 2023. [Slides (pptx) (pdf)] [arXiv version] [Victima Source Code (Officially Artifact Evaluated with All Badges)] Officially artifact evaluated as available, functional, reusable and reproducible. Distinguished artifact award at MICRO 2023.

Victima: Drastically Increasing Address Translation Reach by Leveraging Underutilized Cache Resources

Konstantinos Kanellopoulos¹ Hong Chul Nam¹ F. Nisa Bostanci¹ Rahul Bera¹ Mohammad Sadrosadati¹ Rakesh Kumar² Davide Basilio Bartolini³ Onur Mutlu¹ ¹ETH Zürich ²Norwegian University of Science and Technology ³Huawei Zurich Research Center

SAFARI

https://arxiv.org/pdf/2310.04158

Better Virtual Memory (II)

Konstantinos Kanellopoulos, Rahul Bera, Kosta Stojiljkovic, Nisa Bostanci, Can Firtina, Rachata Ausavarungnirun, Rakesh Kumar, Nastaran Hajinazar, Mohammad Sadrosadati, Nandita Vijaykumar, and Onur Mutlu, "Utopia: Fast and Efficient Address Translation via Hybrid Restrictive & Flexible Virtual-to-Physical Address Mappings" Proceedings of the <u>56th International Symposium on Microarchitecture</u> (MICRO), Toronto, ON, Canada, November 2023. [Slides (pptx) (pdf)] [arXiv version] [Utopia Source Code]

Utopia: Fast and Efficient Address Translation via Hybrid Restrictive & Flexible Virtual-to-Physical Address Mappings

Konstantinos Kanellopoulos¹ Rahul Bera¹ Kosta Stojiljkovic¹ Nisa Bostanci¹ Can Firtina¹ Rachata Ausavarungnirun² Rakesh Kumar³ Nastaran Hajinazar⁴ Mohammad Sadrosadati¹ Nandita Vijaykumar⁵ Onur Mutlu¹

> ¹ETH Zürich ²King Mongkut's University of Technology North Bangkok ³Norwegian University of Science and Technology ⁴Intel Labs ⁵University of Toronto

https://arxiv.org/abs/2211.12205

New Architectures & Technologies

New Architectures & Technologies

- Can have large impact on security and robustness
 Positive or negative
- They need to be designed with system security in mind
 Ideally as a first-class design goal
- Multiple potentially paradigm-changing new technologies and architectures
 - Processing in memory
 - Accelerator-based computing
 - Quantum computing

Processing in Memory



Computing is Bottlenecked by Data



Data is Key for AI, ML, Genomics, ...

Important workloads are all data intensive

 They require rapid and efficient processing of large amounts of data

- Data is increasing
 - We can generate more than we can process
 - We need to perform more sophisticated analyses on more data

Exponential Growth of Neural Networks



Huge Demand for Performance & Efficiency

105

Source: https://youtu.be/Bh13Idwcb0Q?t=283

SAFAR

Huge Demand for Performance & Efficiency



http://www.economist.com/news/21631808-so-much-genetic-data-so-many-uses-genes-unzipped

106

Do We Want This?



Or This?


Challenge and Opportunity for Future

High Performance, Energy Efficient, Sustainable (All at the Same Time) Data access is the major performance and energy bottleneck

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

Today's Computing Systems

- Processor centric
- All data processed in the processor \rightarrow at great system cost



It's the Memory, Stupid!

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

RICHARD SITES

It's the Memory, Stupid!

When we started the Alpha architecture design in 1988, we estimated a 25-year lifetime and a relatively modest 32% per year compounded performance improvement of implementations over that lifetime (1,000× total). We guestimated about 10× would come from CPU clock improvement, 10× from multiple instruction issue, and 10× from multiple processors.

5, 1996 MICROPROCESSOR REPORT

I expect that over the coming decade memory subsystem design will be the *only* important design issue for microprocessors.

Processor-Centric System Performance



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

Processor-Centric System Performance

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



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

Data Movement vs. Computation Energy



A memory access consumes ~100-1000X the energy of a complex addition

Data Movement vs. Computation Energy



Data Movement vs. Computation Energy



A memory access consumes 6400X the energy of a simple integer addition

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}118

Energy Waste in Accelerators

 Amirali Boroumand, Saugata Ghose, Berkin Akin, Ravi Narayanaswami, Geraldo F. Oliveira, Xiaoyu Ma, Eric Shiu, and Onur Mutlu,
 "Google Neural Network Models for Edge Devices: Analyzing and Mitigating Machine Learning Inference Bottlenecks"
 Proceedings of the <u>30th International Conference on Parallel Architectures and Compilation</u> <u>Techniques</u> (PACT), Virtual, September 2021.
 [Slides (pptx) (pdf)]
 [Talk Video (14 minutes)]

> 90% of the total system energy is spent on memory in large ML models

Google Neural Network Models for Edge Devices: Analyzing and Mitigating Machine Learning Inference Bottlenecks

Amirali Boroumand*Saugata Ghose*Berkin Akin*Ravi Narayanaswami*Geraldo F. Oliveira*Xiaoyu Ma*Eric Shiu*Onur Mutlu**

[†]Carnegie Mellon Univ. [•]Stanford Univ. [‡]Univ. of Illinois Urbana-Champaign [§]Google ^{*}ETH Zürich

Processing of data is performed far away from the data



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

Process Data Where It Makes Sense



SAFARI

https://www.gsmarena.com/apple_announces_m1_ultra_with_20core_cpu_and_64core_gpu-news-53481.php

Memory as an Accelerator



Memory similar to a "conventional" accelerator

Goal: Processing Inside Memory/Storage



Processing in Memory: Two Types

Processing near Memory
 Processing using Memory

Processing-in-Memory Landscape Today



SAFARI

And, many other experimental chips and startups

Processing-in-Memory Landscape Today

IEEE COMPUTER ARCHITECTURE LETTERS, VOL. 22, NO. 1, JANUARY-JUNE

Computational CXL-Memory Solution for Accelerating Memory-Intensive Applications

Joonseop Sim[®], Soohong Ahn[®], Taeyoung Ahn[®], Seungyong Lee[®], Myunghyun Rhee, Jooyoung Kim[®], Kwangsik Shin, Donguk Moon[®], Euiseok Kim, and Kyoung Park[®]

Abstract—CXL interface is the up-to-date technology that enables effective memory expansion by providing a memory-sharing protocol in configuring heterogeneous devices. However, its limited physical bandwidth can be a significant bottleneck for emerging data-intensive applications. In this work, we propose a novel CXL-based memory disaggregation architecture with a real-world prototype demonstration, which overcomes the bandwidth limitation of the CXL interface using near-data processing. The experimental results demonstrate that our design achieves up to 1.9× better performance/power efficiency than the existing CPU system.

Index Terms—Compute express link (CXL), near-data-processing (NDP)



SAFARI

Fig. 6. FPGA prototype of proposed CMS card.

Processing-in-Memory Landscape Today

-

Samsung Processing in Memory Technology at Hot Chips 2023

By Patrick Kennedy - August 28, 2023





Samsung PIM PNM For Transformer Based AI HC35_Page_24

Opportunity: 3D-Stacked Logic+Memory



Hybrid Memory Cube



SAFARI

Tesseract System for Graph Processing

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



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

Tesseract Graph Processing Performance

>13X Performance Improvement



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

Tesseract Graph Processing System Energy



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

Processing using DRAM

- We can support
 - Bulk bitwise AND, OR, NOT, MAJ
 - Bulk bitwise COPY and INIT/ZERO
 - True Random Number Generation; Physical Unclonable Functions
 - More complex computation using Lookup Tables
- At low cost
- Using analog computation capability of DRAM
 - Idea: activating (multiple) rows performs computation
 - Even in commodity off-the-shelf DRAM chips!

30X-257X performance and energy improvements

Seshadri+"RowClone: Fast and Efficient In-DRAM Copy and Initialization of Bulk Data," MICRO 2013. Seshadri+, "Fast Bulk Bitwise AND and OR in DRAM", IEEE CAL 2015. Seshadri+, "Ambit: In-Memory Accelerator for Bulk Bitwise Operations Using Commodity DRAM Technology," MICRO 2017. Hajinazar+, "SIMDRAM: A Framework for Bit-Serial SIMD Processing using DRAM," ASPLOS 2021.

Oliveira+, "MIMDRAM: An End-to-End Processing-Using-DRAM System for High-Throughput, Energy-Efficient and Programmer-Transparent Multiple-Instruction Multiple-Data Processing," HPCA 2024.

Future Systems: In-Memory Copy



Capabilities of Off-The-Shelf Memory

Existing DRAM Chips Are Already Quite Capable



Real Processing Using Memory Prototype

- End-to-end RowClone & TRNG using off-the-shelf DRAM chips
- Idea: Violate DRAM timing parameters to mimic RowClone

PiDRAM: A Holistic End-to-end FPGA-based Framework for Processing-in-DRAM

Ataberk Olgun^{§†} Juan Gómez Luna[§] Konstantinos Kanellopoulos[§] Behzad Salami^{§*} Hasan Hassan[§] Oğuz Ergin[†] Onur Mutlu[§] [§]ETH Zürich [†]TOBB ETÜ ^{*}BSC

https://arxiv.org/pdf/2111.00082.pdf https://github.com/cmu-safari/pidram https://www.youtube.com/watch?v=qeukNs5XI3g&t=4192s

Real Processing-using-Memory Prototype



https://arxiv.org/pdf/2111.00082.pdf https://github.com/cmu-safari/pidram https://www.youtube.com/watch?v=qeukNs5XI3g&t=4192s

Real Processing-using-Memory Prototype

E README.md

0

Building a PiDRAM Prototype

To build PiDRAM's prototype on Xilinx ZC706 boards, developers need to use the two sub-projects in this directory. fpga-zynq is a repository branched off of UCB-BAR's fpga-zynq repository. We use fpga-zynq to generate rocket chip designs that support end-to-end DRAM PuM execution. controller-hardware is where we keep the main Vivado project and Verilog sources for PiDRAM's memory controller and the top level system design.

Rebuilding Steps

- Navigate into fpga-zynq and read the README file to understand the overall workflow of the repository

 Follow the readme in fpga-zynq/rocket-chip/riscv-tools to install dependencies
- Create the Verilog source of the rocket chip design using the ZynqCopyFPGAConfig

 Navigate into zc706, then run make rocket C0NFIG=ZynqCopyFPGAConfig -j<number of cores>
- 3. Copy the generated Verilog file (should be under zc706/src) and overwrite the same file in controllerhardware/source/hdl/impl/rocket-chip
- 4. Open the Vivado project in controller-hardware/Vivado_Project using Vivado 2016.2
- 5. Generate a bitstream
- 6. Copy the bitstream (system_top.bit) to fpga-zynq/zc706
- 7. Use the ./build_script.sh to generate the new boot.bin under fpga-images-zc706, you can use this file to program the FPGA using the SD-Card
 - For details, follow the relevant instructions in fpga-zynq/README.md

You can run programs compiled with the RISC-V Toolchain supplied within the fpga-zynq repository. To install the toolchain, follow the instructions under fpga-zynq/rocket-chip/riscv-tools.

Generating DDR3 Controller IP sources

We cannot provide the sources for the Xilinx PHY IP we use in PiDRAM's memory controller due to licensing issues. We describe here how to regenerate them using Vivado 2016.2. First, you need to generate the IP RTL files:

1- Open IP Catalog 2- Find "Memory Interface Generator (MIG 7 Series)" IP and double click

https://arxiv.org/pdf/2111.00082.pdf https://github.com/cmu-safari/pidram

https://www.youtube.com/watch?v=qeukNs5XI3g&t=4192s

Microbenchmark Copy/Initialization Throughput



In-DRAM Copy and Initialization improve throughput by 119x and 89x

SAFARI @kasırga

More on PiDRAM

 Ataberk Olgun, Juan Gomez Luna, Konstantinos Kanellopoulos, Behzad Salami, Hasan Hassan, Oguz Ergin, and Onur Mutlu,
 "PiDRAM: A Holistic End-to-end FPGA-based Framework for Processing-in-DRAM"
 ACM Transactions on Architecture and Code Optimization (TACO), March 2023.
 [arXiv version]
 Presented at the 18th HiPEAC Conference, Toulouse, France, January 2023.
 [Slides (pptx) (pdf)]
 [Longer Lecture Slides (pptx) (pdf)]
 [Lecture Video (40 minutes)]
 [PiDRAM Source Code]

PiDRAM: A Holistic End-to-end FPGA-based Framework for <u>Processing-in-DRAM</u>

Ataberk Olgun§

Juan Gómez Luna
§ Konstantinos Kanellopoulos
§
Hasan Hassan
§ Oğuz Ergin
† Onur Mutlu
§

[§]ETH Zürich [†]TOBB University of Economics and Technology

Behzad Salami[§]

In-DRAM Physical Unclonable Functions

 Jeremie S. Kim, Minesh Patel, Hasan Hassan, and Onur Mutlu, "The DRAM Latency PUF: Quickly Evaluating Physical Unclonable Functions by Exploiting the Latency-Reliability Tradeoff in Modern DRAM <u>Devices"</u> Proceedings of the <u>24th International Symposium on High-Performance Computer</u> <u>Architecture</u> (HPCA), Vienna, Austria, February 2018.

 [Lightning Talk Video]
 [Slides (pptx) (pdf)] [Lightning Session Slides (pptx) (pdf)]
 [Full Talk Lecture Video (28 minutes)]

The DRAM Latency PUF:

Quickly Evaluating Physical Unclonable Functions by Exploiting the Latency-Reliability Tradeoff in Modern Commodity DRAM Devices

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

In-DRAM True Random Number Generation

 Jeremie S. Kim, Minesh Patel, Hasan Hassan, Lois Orosa, and Onur Mutlu, "D-RaNGe: Using Commodity DRAM Devices to Generate True Random Numbers with Low Latency and High Throughput" Proceedings of the <u>25th International Symposium on High-Performance Computer</u> Architecture (HPCA), Washington, DC, USA, February 2019. [Slides (pptx) (pdf)] [Full Talk Video (21 minutes)] [Full Talk Lecture Video (27 minutes)] Top Picks Honorable Mention by IEEE Micro.

D-RaNGe: Using Commodity DRAM Devices to Generate True Random Numbers with Low Latency and High Throughput

Jeremie S. Kim^{‡§}

Minesh Patel[§] Hasan Hassan[§] Lois Orosa[§] Onur Mutlu^{§‡} [‡]Carnegie Mellon University [§]ETH Zürich

In-DRAM True Random Number Generation

 Ataberk Olgun, Minesh Patel, A. Giray Yaglikci, Haocong Luo, Jeremie S. Kim, F. Nisa Bostanci, Nandita Vijaykumar, Oguz Ergin, and Onur Mutlu, <u>"QUAC-TRNG: High-Throughput True Random Number Generation Using</u> <u>Quadruple Row Activation in Commodity DRAM Chips"</u> *Proceedings of the <u>48th International Symposium on Computer Architecture</u> (<i>ISCA*), Virtual, June 2021.
 [Slides (pptx) (pdf)]
 [Short Talk Slides (pptx) (pdf)]
 [Talk Video (25 minutes)]
 [SAFARI Live Seminar Video (1 hr 26 mins)]

QUAC-TRNG: High-Throughput True Random Number Generation Using Quadruple Row Activation in Commodity DRAM Chips

Ataberk OlgunMinesh PatelA. Giray YağlıkçıHaocong LuoJeremie S. KimF. Nisa BostancıNandita VijaykumarOğuz ErginOnur Mutlu§ETH Zürich†TOBB University of Economics and TechnologyOUniversity of Toronto

In-DRAM True Random Number Generation

F. Nisa Bostanci, Ataberk Olgun, Lois Orosa, A. Giray Yaglikci, Jeremie S. Kim, Hasan Hassan, Oguz Ergin, and Onur Mutlu,
 "DR-STRaNGe: End-to-End System Design for DRAM-based True Random Number Generators"
 Proceedings of the <u>28th International Symposium on High-Performance Computer</u>
 <u>Architecture</u> (HPCA), Virtual, April 2022.
 [Slides (pptx) (pdf)]
 [Short Talk Slides (pptx) (pdf)]

DR-STRaNGe: End-to-End System Design for DRAM-based True Random Number Generators

F. Nisa Bostanci^{†§} Ataberk Olgun^{†§} Lois Orosa[§] A. Giray Yağlıkçı[§]
 Jeremie S. Kim[§] Hasan Hassan[§] Oğuz Ergin[†] Onur Mutlu[§]

[†]TOBB University of Economics and Technology

nur Mutlu^s [§]ETH Zürich

SAFARI

https://arxiv.org/pdf/2201.01385.pdf
In-Flash Bulk Bitwise Execution

 Jisung Park, Roknoddin Azizi, Geraldo F. Oliveira, Mohammad Sadrosadati, Rakesh Nadig, David Novo, Juan Gómez-Luna, Myungsuk Kim, and Onur Mutlu, "Flash-Cosmos: In-Flash Bulk Bitwise Operations Using Inherent Computation Capability of NAND Flash Memory" Proceedings of the <u>55th International Symposium on Microarchitecture</u> (MICRO), Chicago, IL, USA, October 2022.
[Slides (pptx) (pdf)]
[Longer Lecture Slides (pptx) (pdf)]
[Lecture Video (44 minutes)]
[arXiv version]

Flash-Cosmos: In-Flash Bulk Bitwise Operations Using Inherent Computation Capability of NAND Flash Memory

Jisung Park^{§∇} Roknoddin Azizi[§] Geraldo F. Oliveira[§] Mohammad Sadrosadati[§] Rakesh Nadig[§] David Novo[†] Juan Gómez-Luna[§] Myungsuk Kim[‡] Onur Mutlu[§]

[§]ETH Zürich ∇ POSTECH [†]LIRMM, Univ. Montpellier, CNRS [‡]Kyungpook National University

SAFARI https://arxiv.org/pdf/2209.05566.pdf

PIM Review and Open Problems

A Modern Primer on Processing in Memory

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

SAFARI Research Group

^aETH Zürich ^bCarnegie Mellon University ^cUniversity of Illinois at Urbana-Champaign ^dKing Mongkut's University of Technology North Bangkok

Onur Mutlu, Saugata Ghose, Juan Gomez-Luna, and Rachata Ausavarungnirun, "A Modern Primer on Processing in Memory" *Invited Book Chapter in <u>Emerging Computing: From Devices to Systems -</u> <i>Looking Beyond Moore and Von Neumann*, Springer, to be published in 2021.

SAFARI

Eliminating the Adoption Barriers

How to Enable Adoption of Processing in Memory

Potential Barriers to Adoption of PIM

1. Applications & software for PIM

2. Ease of **programming** (interfaces and compiler/HW support)

3. **System** and **security** support: coherence, synchronization, virtual memory, isolation, communication interfaces, ...

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

5. Infrastructures to assess benefits and feasibility

All can be solved with change of mindset

We Need to Revisit the Entire Stack

• With a **memory-centric mindset**

Problem		
Aigo	orithm	
Prog	gram/Language	
Syst	tem Software	
SW,	HW Interface	
Mic	ro-architecture	l
Log	ic	J
Dev	vices	
Elec	trons	

We can get there step by step

SAFARI

Potential Security Issues & Benefits (I)

Can PIM worsen security?

- Worsened or easier-to-induce physical issues (e.g., RowHammer)?
- Worsened or new side channels?
- Hardware bugs?
- New threat models?

] ...

Can PIM enhance security?

- Less exposure of data (& keys?)
- □ In-memory (homomorphic) encryption & cryptographic hashing
- Execution of security functions; trusted execution in memory
- Support for security primitives (TRNGs, PUFs, encryption, ...)
- More or better isolation, virtualization, containerization?

Potential Security Issues & Benefits (II)

- Security analysis of PIM Systems
 - Different types of PIM: PnM vs. PuM
 - Different locations: cache, MC, DRAM, NVM, storage, remote, ...
 - □ General-purpose vs. special-purpose PIM?
 - Multi tenancy vs. single workload?
 - Concurrent host and PIM access?
 - Memory bus protection; memory wire(s) protection?
 - Robustness issues like RowHammer, RowPress, ...

• ...

- Can PIM support (more) secure execution of workloads?
 - What is needed to do so?
 - Secure PIM enclaves?

PIM Helps Security: Many Examples

Harshita Gupta, Mayank Kabra, Juan Gómez-Luna, Konstantinos Kanellopoulos, and Onur Mutlu, "Evaluating Homomorphic Operations on a Real-World Processing-In-Memory System" *Proceedings of the 2023 IEEE International Symposium on Workload Characterization Poster Session (IISWC)*, Ghent, Belgium, October 2023. [arXiv version] [Lightning Talk Slides (pptx) (pdf)] [Poster (pptx) (pdf)]

Evaluating Homomorphic Operations on a Real-World Processing-In-Memory System

Harshita Gupta* Mayank Kabra* Juan Gómez-Luna Konstantinos Kanellopoulos Onur Mutlu ETH Zürich

https://arxiv.org/pdf/2309.06545

SAFAR

Amplifying Main Memory-Based Timing Covert and Side Channels using Processing-in-Memory Operations

Konstantinos Kanellopoulos^{†*} F. Nisa Bostancı^{†*} Ataberk Olgun[†] A. Giray Yağlıkçı[†] İsmail Emir Yüksel[†] Nika Mansouri Ghiasi[†] Zülal Bingöl^{†‡} Mohammad Sadrosadati[†] Onur Mutlu[†] [†]ETH Zürich [‡]Bilkent University

Concluding Remarks

Summary: Three Major Limiters

- Technology scaling is not going well
- System complexity is increasing; old methods not keeping up
- Processor-centric designs are not keeping up
- These affect all metrics we care about
- These have fundamental impact on security and how we build secure systems
- We need to revisit how we build architectures and how we secure them

Funding Acknowledgments

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

Thank you!

Acknowledgments

SAFARI Research Group safari.ethz.ch



SAFARI Newsletter June 2023 Edition

<u>https://safari.ethz.ch/safari-newsletter-june-2023/</u>



Think Big, Aim High



y in D

View in your browser June 2023



Referenced Papers, Talks, Artifacts

All are available at

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

https://www.youtube.com/onurmutlulectures

https://github.com/CMU-SAFARI/

Open Source Tools: SAFARI GitHub



https://github.com/CMU-SAFARI/

Future of Computer Architecture and Hardware Security

Onur Mutlu <u>omutlu@gmail.com</u> <u>https://people.inf.ethz.ch/omutlu</u> 16 May 2024

Qualcomm Product Security Summit Keynote Talk

SAFARI





Onur Mutlu's SAFARI Research Group

Computer architecture, HW/SW, systems, bioinformatics, security, memory

https://safari.ethz.ch/safari-newsletter-january-2021/



SAFARI Newsletter December 2021 Edition

<u>https://safari.ethz.ch/safari-newsletter-december-2021/</u>



Think Big, Aim High



f y in 🛛

View in your browser December 2021



SAFARI Newsletter June 2023 Edition

<u>https://safari.ethz.ch/safari-newsletter-june-2023/</u>



Think Big, Aim High



y in D

View in your browser June 2023



SAFARI Introduction & Research

Computer architecture, HW/SW, systems, bioinformatics, security, memory



Seminar in Computer Architecture - Lecture 5: Potpourri of Research Topics (Spring 2023)



SAFARI PhD and Post-Doc Alumni

<u>https://safari.ethz.ch/safari-alumni/</u>

- Hasan Hassan (Rivos), EDAA Outstanding Dissertation Award 2023; S&P 2020 Best Paper Award, 2020 Pwnie Award, IEEE Micro TP HM 2020
- Christina Giannoula (Univ. of Toronto), NTUA Best Dissertation Award 2023
- Minesh Patel (Rutgers, Asst. Prof.), DSN Carter Award Best Thesis 2022; ETH Medal 2023; MICRO'20 & DSN'20 Best Paper Awards; ISCA HoF 2021
- Damla Senol Cali (Bionano Genomics), SRC TECHCON 2019 Best Student Presentation Award; RECOMB-Seq 2018 Best Poster Award
- Nastaran Hajinazar (Intel)
- Gagandeep Singh (AMD/Xilinx), FPL 2020 Best Paper Award Finalist
- Amirali Boroumand (Stanford Univ → Google), SRC TECHCON 2018 Best Presentation Award
- Jeremie Kim (Apple), EDAA Outstanding Dissertation Award 2020; IEEE Micro Top Picks 2019; ISCA/MICRO HoF 2021
- Nandita Vijaykumar (Univ. of Toronto, Assistant Professor), ISCA Hall of Fame 2021
- Kevin Hsieh (Microsoft Research, Senior Researcher)
- Justin Meza (Facebook), HiPEAC 2015 Best Student Presentation Award; ICCD 2012 Best Paper Award
- Mohammed Alser (ETH Zurich), IEEE Turkey Best PhD Thesis Award 2018
- Yixin Luo (Google), HPCA 2015 Best Paper Session
- Kevin Chang (Facebook), SRC TECHCON 2016 Best Student Presentation Award
- Rachata Ausavarungnirun (KMUNTB, Assistant Professor), NOCS 2015 and NOCS 2012 Best Paper Award Finalist
- Gennady Pekhimenko (Univ. of Toronto, Assistant Professor), ISCA Hall of Fame 2021; ASPLOS 2015 SRC Winner
- Vivek Seshadri (Microsoft Research)
- Donghyuk Lee (NVIDIA Research, Senior Researcher), HPCA Hall of Fame 2018
- Yoongu Kim (Software Robotics → Google), TCAD'19 Top Pick Award; IEEE Micro Top Picks'10; HPCA'10 Best Paper Session
- Lavanya Subramanian (Intel Labs → Facebook)
- Samira Khan (Univ. of Virginia, Assistant Professor), HPCA 2014 Best Paper Session
- Saugata Ghose (Univ. of Illinois, Assistant Professor), DFRWS-EU 2017 Best Paper Award
- Jawad Haj-Yahya (Huawei Research Zurich, Principal Researcher)
- Lois Orosa (Galicia Supercomputing Center, Director)
- Jisung Park (POSTECH, Assistant Professor)
- Gagandeep Singh (AMD/Xilinx, Researcher)
- Juan Gomez-Luna (NVIDIA, Researcher), ISPASS 2023 Best Paper Session

SAFARI

Processing in Memory: Evaluation Methods

Simulators (Open Source)

- Ramulator 2.0 & Ramulator-PIM
- DAMOVSim
- UPMEMSim (UPMEM)
- AiMSim (SK Hynix)

Ramulator + Gem5

 Haocong Luo, Yahya Can Tugrul, F. Nisa Bostanci, Ataberk Olgun, A. Giray Yaglikci, and Onur Mutlu, "Ramulator 2.0: A Modern, Modular, and Extensible DRAM Simulator" *Preprint on arxiv*, August 2023.
[arXiv version]
[Ramulator 2.0 Source Code]

Ramulator 2.0: A Modern, Modular, and Extensible DRAM Simulator

Haocong Luo, Yahya Can Tuğrul, F. Nisa Bostancı, Ataberk Olgun, A. Giray Yağlıkçı, and Onur Mutlu

https://arxiv.org/pdf/2308.11030.pdf

SAFARI https://github.com/CMU-SAFARI/ramulator2

Opportunity: 3D-Stacked Logic+Memory



Hybrid Memory Cube



Tesseract System for Graph Processing

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



Tesseract System for Graph Processing



SAFARI

Tesseract System for Graph Processing



SAFARI

Simulated Systems



Tesseract Graph Processing Performance

>13X Performance Improvement



Tesseract Graph Processing System Energy



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 Computer</u> Architecture (ISCA), Portland, OR, June 2015. [Slides (pptx) (pdf)] [Lightning Session Slides (pptx) (pdf)] Top Picks Honorable Mention by IEEE Micro. Selected to the ISCA-50 25-Year Retrospective Issue covering 1996-2020 in 2023 (Retrospective (pdf) Full Issue).

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



Processing-in-Memory in the Real World



Processing-in-Memory Landscape Today



SAFARI

And, many other experimental chips and startups

Processing-in-Memory Landscape Today

IEEE COMPUTER ARCHITECTURE LETTERS, VOL. 22, NO. 1, JANUARY-JUNE

Computational CXL-Memory Solution for Accelerating Memory-Intensive Applications

Joonseop Sim[®], Soohong Ahn[®], Taeyoung Ahn[®], Seungyong Lee[®], Myunghyun Rhee, Jooyoung Kim[®], Kwangsik Shin, Donguk Moon[®], Euiseok Kim, and Kyoung Park[®]

Abstract—CXL interface is the up-to-date technology that enables effective memory expansion by providing a memory-sharing protocol in configuring heterogeneous devices. However, its limited physical bandwidth can be a significant bottleneck for emerging data-intensive applications. In this work, we propose a novel CXL-based memory disaggregation architecture with a real-world prototype demonstration, which overcomes the bandwidth limitation of the CXL interface using near-data processing. The experimental results demonstrate that our design achieves up to 1.9× better performance/power efficiency than the existing CPU system.

Index Terms—Compute express link (CXL), near-data-processing (NDP)



SAFARI

Fig. 6. FPGA prototype of proposed CMS card.
Processing-in-Memory Landscape Today

-

Samsung Processing in Memory Technology at Hot Chips 2023

By Patrick Kennedy - August 28, 2023





Samsung PIM PNM For Transformer Based AI HC35_Page_24

Samsung AxDIMM (2021)

DDRx-PIM

DLRM recommendation system





AxDIMM System





Ske Eta Rear-Memory Processing in Action: Accelerating Personalized Recommendation with AxDIMM", IEEE Micro (2021)

Samsung Newsroom

CORPORATE | PRODUCTS | PRESS RESOURCES | VIEWS | ABOUT US

Audio

Share (🔊

Samsung Develops Industry's First High Bandwidth Memory with AI Processing Power

Korea on February 17, 2021

The new architecture will deliver over twice the system performance and reduce energy consumption by more than 70%

Samsung Electronics, the world leader in advanced memory technology, today announced that it has developed the industry's first High Bandwidth Memory (HBM) integrated with artificial intelligence (AI) processing power – the HBM-PIM The new processing-in-memory (PIM) architecture brings powerful AI computing capabilities inside high-performance memory, to accelerate large-scale processing in data centers, high performance computing (HPC) systems and AI-enabled mobile applications.

Kwangil Park, senior vice president of Memory Product Planning at Samsung Electronics stated, "Our groundbreaking HBM-PIM is the industry's first programmable PIM solution tailored for diverse AI-driven workloads such as HPC, training and inference. We plan to build upon this breakthrough by further collaborating with AI solution providers for even more advanced PIM-powered applications."





[3D Chip Structure of HBM with FIMDRAM]

ISSCC 2021 / SESSION 25 / DRAM / 25.4

25.4 A 20nm 6GB Function-In-Memory DRAM, Based on HBM2 with a 1.2TFLOPS Programmable Computing Unit Using Bank-Level Parallelism, for Machine Learning Applications

Young-Cheon Kwon', Suk Han Lee', Jaehoon Lee', Sang-Hyuk Kwon', Je Min Ryu', Jong-Pil Son', Seongil O', Hak-Soo Yu', Haesuk Lee', Soo Young Kim', Youngmin Cho', Jin Guk Kim', Jongyoon Choi', Hyun-Sung Shin', Jin Kim', BengSeng Phuah', HyoungMin Kim', Myeong Jun Song', Ahn Choi', Daeho Kim', SooYoung Kim', Eun-Bong Kim', David Wang', Shinhaeng Kang', Yuhwan Ro³, Seungwoo Seo³, JoonHo Song³, Jaeyoun Youn', Kyomin Sohn', Nam Sung Kim'

¹Samsung Electronics, Hwaseong, Korea ²Samsung Electronics, San Jose, CA ³Samsung Electronics, Suwon, Korea

Programmable Computing Unit

- Configuration of PCU block
 - Interface unit to control data flow
 - Execution unit to perform operations
 - Register group
 - 32 entries of CRF for instruction memory
 - 16 GRF for weight and accumulation
 - 16 SRF to store constants for MAC operations



[Block diagram of PCU in FIMDRAM]

ISSCC 2021 / SESSION 25 / DRAM / 25.4

25.4 A 20nm 6GB Function-In-Memory DRAM, Based on HBM2 with a 1.2TFLOPS Programmable Computing Unit Using Bank-Level Parallelism, for Machine Learning Applications

Young-Cheon Kwon', Suk Han Lee', Jaehoon Lee', Sang-Hyuk Kwon', Je Min Ryu', Jong-Pil Son', Seongi O, Hak-Soo Yu', Haesuk Lee', Soo Young Kim', Youngmin Che', In Guk Kim', Jongyoon Choi', Hyun-Sung Shin', Jin Kim', BengSeng Phuah', HyoungMin Kim', Myeong Jun Song, Ahn Choi', Daeho Kim', SooYoung Kim', Eun-Bong Kim', David Wang', Shinhaeng Kang', Yuhwan Ro', Seungwoo Seo', JoonHo Song', Jaeyoun Youn', Kyomin Sohn', Man Sung Kim'

[Available instruction list for FIM operation]

Туре	CMD	Description	
Floating Point	ADD	FP16 addition	
	MUL	FP16 multiplication	
	MAC	FP16 multiply-accumulate	
	MAD	FP16 multiply and add	
Data Path	MOVE	Load or store data	
	FILL	Copy data from bank to GRFs	
Control Path	NOP	Do nothing	
	JUMP	Jump instruction	
	EXIT	Exit instruction	

ISSCC 2021 / SESSION 25 / DRAM / 25.4

25.4 A 20nm 6GB Function-In-Memory DRAM, Based on HBM2 with a 1.2TFLOPS Programmable Computing Unit Using Bank-Level Parallelism, for Machine Learning Applications

Young-Theon Kwon', Suk Han Lee', Jaehoon Lee', Sang-Hyuk Kwon', Ja Min Ryu', Jong-Hi Son, Seongi Ol, Hak-Soo Yu', Hasauk Lee', Soo Young Kim', Youngmin Cho', Jin Guk Kim', Jongyoon Choi', Hyun-Sung Shin', Jan Kim', BengSeng Pinair', HyoungMin Kim', Myeong Jun Song', Ahn Choi, Daeho Kim', SooYoong Kim', Eun-Bong Kim', David Wang', Shinhaeng Kang', Yuhwan Ro', Seungwoo Seo', JoonHo Song', Jaeyoun Youn', Kyomin Sohn', Mam Sung Kim'

Chip Implementation

- Mixed design methodology to implement FIMDRAM
 - Full-custom + Digital RTL



[Digital RTL design for PCU block]

ISSCC 2021 / SESSION 25 / DRAM / 25.4

25.4 A 20nm 6GB Function-In-Memory DRAM, Based on HBM2 with a 1.2TFLOPS Programmable Computing Unit Using Bank-Level Parallelism, for Machine Learning Applications

Young-Cheon Kwon', Suk Han Lee', Jaehoon Lee', Sang-Hyuk Kwon', Je Min Ryu', Jong-Pil Son', Seongil O', Hak-Soo Yu', Haesuk Lee', Soo Young Kim', Youngmin Cho', Jin Guk Kim', Jongyoon Cho', Hyun-Sung Shin', Jin Kim', BengSeng Phuah', HyoungMin Kim', Myeong Jun Soong', Ahn Cho', Deaho Kim', Sooryoung Kim', Euro-Bong Kim', David Wang', Shinhaeng Kang', Yuhwan Ro', Seungwoo See', JoonHo Song', Jaeyoun Youn', Kiyomin Sohn', Nam Sung Kim'

¹Samsung Electronics, Hwaseong, Korea ²Samsung Electronics, San Jose, CA ³Samsung Electronics, Suwon, Korea

Cell array for bank0	Cell array for bank4	Cell array for bank0	Cell array for bank4	Pseudo	Pseudo
PCU block for bank0 & 1	PCU block for bank4 & 5	PCU block for bank0 & 1	PCU block for bank4 & 5	channel-0	channel-1
Cell array for bank1 Cell array for bank2	Cell array for bank5 Cell array for bank6	Cell array for bank1 Cell array for bank2	Cell array for bank5 Cell array for bank6		
PCU block for bank2 & 3	PCU block for bank6 & 7	PCU block for bank2 & 3	PCU block for bank6 & 7		and the part of th
Cell array for bank3	Cell array for bank7	Cell array for bank3	Cell array for bank7		
	n an				
		TSV &	Peri Co	ontrol Block	
Cell array for bank11	Cell array for bank15	Cell array for bank11	Cell array for bank15		
PCU block for bank10 & 11	PCU block for bank14 & 15	PCU block for bank10 & 11	PCU block for bank14 & 15		
Cell array for bank10	Cell array for bank14	Cell array for bank10	Cell array for bank14		
Cell array for bank9	Cell array for bank13	Cell array for bank9	Cell array for bank13		
PCU block for bank8 & 9	PCU block for bank12 & 13	PCU block for bank8 & 9	PCU block for bank12 & 13	Pseudo	Pseudo
Cell array for bank8	Cell array for bank12	Cell array for bank8	Cell array for bank12	channel-0	channel-1

SK Hynix Accelerator-in-Memory (2022)

SKhynix NEWSROOM

SK hvnix STORY

INSIGHT

PRESS CENTER

Search

🌐 ENG 🗸

Q

SK hynix Develops PIM, Next-Generation AI Accelerator

MULTIMEDIA

February 16, 2022

Seoul, February 16, 2022

SK hynix (or "the Company", www.skhynix.com) announced on February 16 that it has developed PIM*, a nextgeneration memory chip with computing capabilities.

*PIM(Processing In Memory): A next-generation technology that provides a solution for data congestion issues for AI and big data by adding computational functions to semiconductor memory

It has been generally accepted that memory chips store data and CPU or GPU, like human brain, process data. SK hynix, following its challenge to such notion and efforts to pursue innovation in the next-generation smart memory, has found a breakthrough solution with the development of the latest technology.

SK hynix plans to showcase its PIM development at the world's most prestigious semiconductor conference, 2022 ISSCC*, in San Francisco at the end of this month. The company expects continued efforts for innovation of this technology to bring the memory-centric computing, in which semiconductor memory plays a central role, a step closer In Paper 11.1, SK Hynix describes an 1ynm, GDDR6-based accelerator-in-memory with a command set for deep-learning operation. The to the reality in devices such as smartphones.

*ISSCC: The International Solid-State Circuits Conference will be held virtually from Feb. 20 to Feb. 24 this year with a theme of "Intelligent Silicon for a Sustainable World"

For the first product that adopts the PIM technology, SK hynix has developed a sample of GDDR6-AiM (Accelerator* in memory). The GDDR6-AiM adds computational functions to GDDR6* memory chips, which process data at 16Gbps. A combination of GDDR6-AiM with CPU or GPU instead of a typical DRAM makes certain computation speed 16 times faster. GDDR6-AiM is widely expected to be adopted for machine learning, high-performance computing, and big data computation and storage.



11.1 A 1ynm 1.25V 8Gb, 16Gb/s/pin GDDR6-based Accelerator-in-Memory supporting 1TFLOPS MAC Operation and Various Activation Functions for Deep-Learning Applications

Seongiu Lee, SK hynix, Icheon, Korea

8Gb design achieves a peak throughput of 1TFLOPS with 1GHz MAC operations and supports major activation functions to improve accuracy.

https://news.skhynix.com/sk-hynix-develops-pim-next-generation-ai-accelerator/

SK Hynix Accelerator-in-Memory (2022)



ASPLOS 2023 Tutorial: Real-world Processing-in-Memory Systems for Modern Workloads



1,146 views Streamed live on Mar 26, 2023 Livestream - Data-Centric Architectures: Fundamentally Improving Performance and Energy (Spring 2023) ASPLOS 2023 Tutorial: Real-world Processing-in-Memory Systems for Modern Workloads https://events.safari.ethz.ch/asplos-...

SAFARI

https://www.youtube.com/watch?v=oYCaLcT0Kmo

AliBaba PIM Recommendation System (2022)



Figure 29.1.2: Illustration of 3D-stacked chip, cross-illustration of package, DRAM array layout and design blocks on logic die.

Neural Engine Region DRAM Die Photo (36Gb) JE I Match Engine Region ME IF 7.02 mr DRAM Die Logic Die Technology 25nm Technology 55nm Total* 602 22 mm Total Neural Engine 32 mm² Area Neural Engine Area Match Engine 32 mm² Match Engine # of MC Voltage 1.1 V 16 per IF Frequency (max) 150 MHz Voltage 1.2 V 300 mW per 1Gb Frequency 300 MHz Power 153.60 GB/s / 1.38 TB/s 977.70 mW Bandwidth* Power * Total die size of test chip Precision INT8 ** Memory bandwidth of NE and ME / Total bandwidth of DRAM die

Figure 29.1.7: Die micrographs of DRAM die, NE and ME. Detailed specifications of DRAM die and logic die.

29.1 184QPS/W 64Mb/mm² 3D Logic-to-DRAM Hybrid Bonding with Process-Near-Memory Engine for Recommendation System

Dimin Niu¹, Shuangchen Li¹, Yuhao Wang¹, Wei Han¹, Zhe Zhang², Yijin Guan², Tianchan Guan³, Fei Sun¹, Fei Xue¹, Lide Duan¹, Yuanwei Fang¹, Hongzhong Zheng¹, Xiping Jiang⁴, Song Wang⁴, Fengguo Zuo⁴, Yubing Wang⁴, -**SAFARI** Bing Yu⁴, Qiwei Ren⁴, Yuan Xie¹

602 22 mm

5.90 mm²

7.02 mm²

UPMEM Processing-in-DRAM Engine (2019)

Processing in DRAM Engine

 Includes standard DIMM modules, with a large number of DPU processors combined with DRAM chips.

Replaces standard DIMMs

- DDR4 R-DIMM modules
 - 8GB+128 DPUs (16 PIM chips)
 - Standard 2x-nm DRAM process



Large amounts of compute & memory bandwidth



https://www.upmem.com/video-upmem-presenting-its-true-processing-in-memory-solution-hot-chips-2019/

UPMEM Memory Modules

- E19: 8 chips DIMM (1 rank). DPUs @ 267 MHz
- P21: 16 chips DIMM (2 ranks). DPUs @ 350 MHz



SAFARI

2,560-DPU Processing-in-Memory System



Benchmarking a New Paradigm: An Experimental Analysis of a Real Processing-in-Memory Architecture

JUAN GÓMEZ-LUNA, ETH Zürich, Switzerland IZZAT EL HAJJ, Amerian University of Betrut, Lebanon IVAN FERNANDEZ, ETH Zürich, Switzerland and University of Malaga, Spain CHRISTINA GIANNOULA, ETH Zürich, Switzerland and NTUA, Greece GERALDO F. OLIVEIRA, ETH Zürich, Switzerland ONUR MUTLU, ETH Zürich, Switzerland

Many modern workloads, such as neural networks, databases, and graph processing, are fundamentally memory-bound for such workloads, the data movement between main memory and CPU cores imposes a significant overhead in terms of both latency and energy. A major reason is that this communication happense through a narrow bus with high latency and limited bandwidth, and the low data reuse in memory-bound workloads is insufficient to amorize the cost of main memory ancess. Fundamentally addressing this data movement builteneck requires a paradigm where the memory system assumes an active role in computing by integrating processing capabilities. This paradigm is known as *processing-in-memory* (PdM).

Recent research explores different forms of PIM architectures, motivated by the emergence of new 3Dstacked memory technologies that integrate memory with a logic layer where processing elements can be easily placed. Past works evaluate these architectures in simulation or, at best, with simplified hardware prototypes. In contrast, the UPMEM company has designed and manufactured the first publicly-available real-world PIM architecture. The UPMEM PIM architecture combines traditional DRAM memory arrays with general-purpose in-order cores, called DRAM Processing Units (DPUs), integrated in the same chip.

This paper provides the first comprehensive analysis of the first publicly-available real-world PM architecture. We make two evolution string, we conduct an experimental characterization of the UPMRM-based PIM system using microbenchmarks to assess various architecture limits such as compute throughput and memory bandwisht, yielding new insights. Second, we present PIM (*Coressing in-Phaemory benchmarks*), a benchmark suite of 16 worldoads from different application domains (e.g., dense/sparse linear algebra, dathases, data analytics, graph processing, neural networks, bioinformatics, image processing), which we identify as memory-bound. We evaluate the performance and acaling characteristics of PIM benchmarks on the UPMRM PIM architecture, and compare their performance and energy consumption to their stateof-the-art CPU and CPU counterparts. Our extensive evaluation conducted on two real UPMRM-based PIM systems with 64 and 2550 PUP sproides new insights about suitability of different worldoads to the PIM systems reares of future PIM systems.



https://arxiv.org/pdf/2105.03814.pdf

More on the UPMEM PIM System



https://www.youtube.com/watch?v=Sscy1Wrr22A&list=PL5Q2soXY2Zi9xidyIgBxUz7xRPS-wisBN&index=26

Experimental Analysis of the UPMEM PIM Engine

Benchmarking a New Paradigm: An Experimental Analysis of a Real Processing-in-Memory Architecture

JUAN GÓMEZ-LUNA, ETH Zürich, Switzerland IZZAT EL HAJJ, American University of Beirut, Lebanon IVAN FERNANDEZ, ETH Zürich, Switzerland and University of Malaga, Spain CHRISTINA GIANNOULA, ETH Zürich, Switzerland and NTUA, Greece GERALDO F. OLIVEIRA, ETH Zürich, Switzerland ONUR MUTLU, ETH Zürich, Switzerland

Many modern workloads, such as neural networks, databases, and graph processing, are fundamentally memory-bound. For such workloads, the data movement between main memory and CPU cores imposes a significant overhead in terms of both latency and energy. A major reason is that this communication happens through a narrow bus with high latency and limited bandwidth, and the low data reuse in memory-bound workloads is insufficient to amortize the cost of main memory access. Fundamentally addressing this *data movement bottleneck* requires a paradigm where the memory system assumes an active role in computing by integrating processing capabilities. This paradigm is known as *processing-in-memory (PIM*).

Recent research explores different forms of PIM architectures, motivated by the emergence of new 3Dstacked memory technologies that integrate memory with a logic layer where processing elements can be easily placed. Past works evaluate these architectures in simulation or, at best, with simplified hardware prototypes. In contrast, the UPMEM company has designed and manufactured the first publicly-available real-world PIM architecture. The UPMEM PIM architecture combines traditional DRAM memory arrays with general-purpose in-order cores, called *DRAM Processing Units* (*DPUs*), integrated in the same chip.

This paper provides the first comprehensive analysis of the first publicly-available real-world PIM architecture. We make two key contributions. First, we conduct an experimental characterization of the UPMEM-based PIM system using microbenchmarks to assess various architecture limits such as compute throughput and memory bandwidth, yielding new insights. Second, we present *PrIM* (*Processing-In-Memory benchmarks*), a benchmark suite of 16 workloads from different application domains (e.g., dense/sparse linear algebra, databases, data analytics, graph processing, neural networks, bioinformatics, image processing), which we identify as memory-bound. We evaluate the performance and scaling characteristics of PrIM benchmarks on the UPMEM PIM architecture, and compare their performance and energy consumption to their stateof-the-art CPU and GPU counterparts. Our extensive evaluation conducted on two real UPMEM-based PIM systems with 640 and 2,556 DPUs provides new insights about suitability of different workloads to the PIM system, programming recommendations for software designers, and suggestions and hints for hardware and architecture designers of future PIM systems.

https://arxiv.org/pdf/2105.03814.pdf

UPMEM PIM System Summary & Analysis

Juan Gomez-Luna, Izzat El Hajj, Ivan Fernandez, Christina Giannoula, Geraldo F. Oliveira, and Onur Mutlu, "Benchmarking Memory-Centric Computing Systems: Analysis of Real **Processing-in-Memory Hardware**" Invited Paper at Workshop on Computing with Unconventional Technologies (CUT), Virtual, October 2021. [arXiv version] PrIM Benchmarks Source Code [Slides (pptx) (pdf)] [Talk Video (37 minutes)] [Lightning Talk Video (3 minutes)]

Benchmarking Memory-Centric Computing Systems: Analysis of Real Processing-in-Memory Hardware

Juan Gómez-Luna ETH Zürich

Izzat El Hajj American University of Beirut

University of Malaga

Ivan Fernandez Christina Giannoula Geraldo F. Oliveira Onur Mutlu National Technical University of Athens

ETH Zürich ETH Zürich

PrIM Benchmarks: Application Domains

Domain	Benchmark	Short name
Donco linear algobra	Vector Addition	VA
Dense linear algebra	Matrix-Vector Multiply	GEMV
Sparse linear algebra	Sparse Matrix-Vector Multiply	SpMV
Databases	Select	SEL
DataDases	Unique	UNI
Data analytics	Binary Search	BS
Data analytics	Time Series Analysis	TS
Graph processing	Breadth-First Search	BFS
Neural networks	Multilayer Perceptron	MLP
Bioinformatics	Needleman-Wunsch	NW
Image processing	Image histogram (short)	HST-S
inage processing	Image histogram (large)	HST-L
	Reduction	RED
Darallal primitivas	Prefix sum (scan-scan-add)	SCAN-SSA
r ar aner prinnitives	Prefix sum (reduce-scan-scan)	SCAN-RSS
	Matrix transposition	TRNS

SAFARI

PrIM Benchmarks are Open Source

- All microbenchmarks, benchmarks, and scripts
- <u>https://github.com/CMU-SAFARI/prim-benchmarks</u>

G CMU-SAFARI / prim-benchmarks	ⓒ Unwatch v 2 \overleftrightarrow Star 2 $\frac{29}{6}$ Fork 1				
<> Code ⊙ Issues 10 Pull requests ⊙ Actions ⊡ Projects □ Wiki	🕕 Security 🗁 Insights 🕸 Settings				
ি main 👻 prim-benchmarks / README.md	Go to file				
Juan Gomez Luna PrIM first commit Latest commit 3de4b49 9 days ago 3 History					
At 1 contributor					
i≡ 168 lines (132 sloc) 5.79 KB	Raw Blame 🖵 🖉 🖞				
PrIM (Processing-In-Memory Benchmarks)					
analyze, and characterize the first publicly-available real-world processing-in-memory (PIM) architecture, the UPMEM PIM architecture. The UPMEM PIM architecture combines traditional DRAM memory arrays with general-purpose in-order cores, called DRAM Processing Units (DPUs), integrated in the same chip.					
PrIM provides a common set of workloads to evaluate the UPMEM PIM architecture with and can be useful for programming, architecture and system researchers all alike to improve multiple aspects of future PIM hardware and software. The workloads have different characteristics, exhibiting heterogeneity in their memory access patterns, operations and data types, and communication patterns. This repository also contains baseline CPU and GPU implementations of PrIM benchmarks for comparison purposes.					
PrIm also includes a set of microbenchmarks can be used to assess various architecture limits such as compute throughput and					

memory bandwidth.

Benchmarking a New Paradigm: Experimental Analysis and Characterization of a Real Processing-in-Memory System

JUAN GÓMEZ-LUNA¹, IZZAT EL HAJJ², IVAN FERNANDEZ^{1,3}, CHRISTINA GIANNOULA^{1,4}, GERALDO F. OLIVEIRA¹, AND ONUR MUTLU¹

¹ETH Zürich

²American University of Beirut

³University of Malaga

⁴National Technical University of Athens

Corresponding author: Juan Gómez-Luna (e-mail: juang@ethz.ch).

https://arxiv.org/pdf/2105.03814.pdf

https://github.com/CMU-SAFARI/prim-benchmarks

SAFARI

Real Processing Using Memory Prototype

- End-to-end RowClone & TRNG using off-the-shelf DRAM chips
- Idea: Violate DRAM timing parameters to mimic RowClone

PiDRAM: A Holistic End-to-end FPGA-based Framework for Processing-in-DRAM

Ataberk Olgun^{§†} Juan Gómez Luna[§] Konstantinos Kanellopoulos[§] Behzad Salami^{§*} Hasan Hassan[§] Oğuz Ergin[†] Onur Mutlu[§] [§]ETH Zürich [†]TOBB ETÜ ^{*}BSC

https://arxiv.org/pdf/2111.00082.pdf https://github.com/cmu-safari/pidram https://www.youtube.com/watch?v=qeukNs5XI3g&t=4192s

Real Processing Using Memory Prototype



https://arxiv.org/pdf/2111.00082.pdf https://github.com/cmu-safari/pidram https://www.youtube.com/watch?v=qeukNs5XI3g&t=4192s

Real Processing Using Memory Prototype

E README.md

0

Building a PiDRAM Prototype

To build PiDRAM's prototype on Xilinx ZC706 boards, developers need to use the two sub-projects in this directory. fpga-zynq is a repository branched off of UCB-BAR's fpga-zynq repository. We use fpga-zynq to generate rocket chip designs that support end-to-end DRAM PuM execution. controller-hardware is where we keep the main Vivado project and Verilog sources for PiDRAM's memory controller and the top level system design.

Rebuilding Steps

- Navigate into fpga-zynq and read the README file to understand the overall workflow of the repository

 Follow the readme in fpga-zynq/rocket-chip/riscv-tools to install dependencies
- Create the Verilog source of the rocket chip design using the ZynqCopyFPGAConfig

 Navigate into zc706, then run make rocket C0NFIG=ZynqCopyFPGAConfig -j<number of cores>
- 3. Copy the generated Verilog file (should be under zc706/src) and overwrite the same file in controllerhardware/source/hdl/impl/rocket-chip
- 4. Open the Vivado project in controller-hardware/Vivado_Project using Vivado 2016.2
- 5. Generate a bitstream
- 6. Copy the bitstream (system_top.bit) to fpga-zynq/zc706
- 7. Use the ./build_script.sh to generate the new boot.bin under fpga-images-zc706 , you can use this file to program the FPGA using the SD-Card
 - For details, follow the relevant instructions in fpga-zynq/README.md

You can run programs compiled with the RISC-V Toolchain supplied within the fpga-zynq repository. To install the toolchain, follow the instructions under fpga-zynq/rocket-chip/riscv-tools.

Generating DDR3 Controller IP sources

We cannot provide the sources for the Xilinx PHY IP we use in PiDRAM's memory controller due to licensing issues. We describe here how to regenerate them using Vivado 2016.2. First, you need to generate the IP RTL files:

1- Open IP Catalog 2- Find "Memory Interface Generator (MIG 7 Series)" IP and double click

https://arxiv.org/pdf/2111.00082.pdf https://github.com/cmu-safari/pidram

https://www.youtube.com/watch?v=qeukNs5XI3g&t=4192s

Microbenchmark Copy/Initialization Throughput



In-DRAM Copy and Initialization improve throughput by 119x and 89x

SAFARI @kasırga

More on PiDRAM

 Ataberk Olgun, Juan Gomez Luna, Konstantinos Kanellopoulos, Behzad Salami, Hasan Hassan, Oguz Ergin, and Onur Mutlu,
 "PiDRAM: A Holistic End-to-end FPGA-based Framework for Processing-in-DRAM"
 ACM Transactions on Architecture and Code Optimization (TACO), March 2023.
 [arXiv version]
 Presented at the 18th HiPEAC Conference, Toulouse, France, January 2023.
 [Slides (pptx) (pdf)]
 [Longer Lecture Slides (pptx) (pdf)]
 [Lecture Video (40 minutes)]
 [PiDRAM Source Code]

PiDRAM: A Holistic End-to-end FPGA-based Framework for <u>P</u>rocessing-<u>i</u>n-<u>DRAM</u>

Ataberk Olgun[§]

Juan Gómez Luna
§ Konstantinos Kanellopoulos
§
Hasan Hassan
§ Oğuz Ergin
† Onur Mutlu
§

[§]ETH Zürich [†]TOBB University of Economics and Technology

Behzad Salami[§]

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

www.iaik.tugraz.at





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'16206

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. IEEE S&P 2018

ars technica

BIZ & IT TECH SCIENCE POLICY CARS GAMING & CULTURE

"GRAND PWNING UNIT" --

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) USENIX ATC 2018

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 (VI)

IEEE S&P 2020

RAMBleed: Reading Bits in Memory Without Accessing Them

Andrew Kwong University of Michigan ankwong@umich.edu Daniel Genkin University of Michigan genkin@umich.edu Daniel Gruss Graz University of Technology daniel.gruss@iaik.tugraz.at Yuval Yarom University of Adelaide and Data61 yval@cs.adelaide.edu.au

More Security Implications (VII)

USENIX Security 2019

Terminal Brain Damage: Exposing the Graceless Degradation in Deep Neural Networks Under Hardware Fault Attacks

Sanghyun Hong, Pietro Frigo[†], Yiğitcan Kaya, Cristiano Giuffrida[†], Tudor Dumitraş

University of Maryland, College Park [†]Vrije Universiteit Amsterdam



A Single Bit-flip Can Cause Terminal Brain Damage to DNNs One specific bit-flip in a DNN's representation leads to accuracy drop over 90%

Our research found that a specific bit-flip in a DNN's bitwise representation can cause the accuracy loss up to 90%, and the DNN has 40-50% parameters, on average, that can lead to the accuracy drop over 10% when individually subjected to such single bitwise corruptions...

Read More

More Security Implications (VIII)

USENIX Security 2020

DeepHammer: Depleting the Intelligence of Deep Neural Networks through Targeted Chain of Bit Flips

Fan Yao A University of Central Florida fan.yao@ucf.edu a

Adnan Siraj RakinDeliang FanArizona State Universityasrakin@asu.edudfan@asu.edu

Degrade the **inference accuracy** to the level of **Random Guess**

Example: ResNet-20 for CIFAR-10, 10 output classes

Before attack, Accuracy: 90.2% After attack, Accuracy: ~10% (1/10)



Google's Half-Double RowHammer Attack (May 2021)

Google Security Blog

The latest news and insights from Google on security and safety on the Internet

Introducing Half-Double: New hammering technique for DRAM Rowhammer bug

May 25, 2021

Research Team: Salman Qazi, Yoongu Kim, Nicolas Boichat, Eric Shiu & Mattias Nissler

Today, we are sharing details around our discovery of Half-Double, a new Rowhammer technique that capitalizes on the worsening physics of some of the newer DRAM chips to alter the contents of memory.

Rowhammer is a DRAM vulnerability whereby repeated accesses to one address can tamper with the data stored at other addresses. Much like speculative execution vulnerabilities in CPUs, Rowhammer is a breach of the security guarantees made by the underlying hardware. As an electrical coupling phenomenon within the silicon itself, Rowhammer allows the potential bypass of hardware and software memory protection policies. This can allow untrusted code to break out of its sandbox and take full control of the system.

More Security Implications (VIII)

USENIX Security 2022

 Google's Half-Double RowHammer Attack

Google Security Blog

The latest news and insights from Google on security and safety on the Internet

Introducing Half-Double: New hammering technique for DRAM Rowhammer bug May 25, 2021

Research Team: Salman Qazi, Yoongu Kim, Nicolas Boichat, Eric Shiu & Mattias Nissler

Today, we are sharing details around our discovery of Half-Double, a new Rowhammer technique that capitalizes on the worsening physics of some of the newer DRAM chips to alter the contents of memory.

Rowhammer is a DRAM vulnerability whereby repeated accesses to one address can tamper with the data stored at other addresses. Much like speculative execution vulnerabilities in CPUs, Rowhammer is a breach of the security guarantees made by the underlying hardware. As an electrical coupling phenomenon within the silicon itself, Rowhammer allows the potential bypass of hardware and software memory protection policies. This can allow untrusted code to break out of its sandbox and take full control of the system.

Half-Double: Hammering From the Next Row Over

Andreas Kogler¹ Jonas Juffinger^{1,2} Salman Qazi³ Yoongu Kim³ Moritz Lipp^{4*} Nicolas Boichat³ Eric Shiu⁵ Mattias Nissler³ Daniel Gruss¹

¹Graz University of Technology ²Lamarr Security Research ³Google ⁴Amazon Web Services ⁵Rivos

SAFARI

https://www.usenix.org/system/files/sec22-kogler-half-double.pdf

More Security Implications?



A RowHammer Survey Across the Stack

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]
 [Slides from COSADE 2019 (pptx)]
 [Slides from VLSI-SOC 2020 (pptx) (pdf)]
 [Talk Video (1 hr 15 minutes, with Q&A)]

RowHammer: A Retrospective

Onur Mutlu§‡Jeremie S. Kim‡§§ETH Zürich‡Carnegie Mellon University
A RowHammer Survey: Recent Update

 Onur Mutlu, Ataberk Olgun, and A. Giray Yaglikci, "Fundamentally Understanding and Solving RowHammer" Invited Special Session Paper at the <u>28th Asia and South Pacific Design</u> Automation Conference (ASP-DAC), Tokyo, Japan, January 2023. [arXiv version] [Slides (pptx) (pdf)] [Talk Video (26 minutes)]

Fundamentally Understanding and Solving RowHammer

Onur Mutlu onur.mutlu@safari.ethz.ch ETH Zürich Zürich, Switzerland Ataberk Olgun ataberk.olgun@safari.ethz.ch ETH Zürich Zürich, Switzerland A. Giray Yağlıkcı giray.yaglikci@safari.ethz.ch ETH Zürich Zürich, Switzerland

https://arxiv.org/pdf/2211.07613.pdf

SAFARI



Main Memory Needs **Intelligent Controllers** for Security, Safety, Reliability, Scaling

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.

Intelligent Flash Controllers [PIEEE'17]



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

Two Major RowHammer Directions

Understanding RowHammer

- Many effects still need to be rigorously examined
 - Aging of DRAM Chips
 - Environmental Conditions (e.g., Process, Voltage, Temperature)
 - Memory Access Patterns
 - Memory Controller & System Design Decisions

...

Solving RowHammer

- Flexible and efficient solutions are necessary
 - In-field patchable / reconfigurable / programmable solutions
- Co-architecting System and Memory is important
 - To avoid performance and denial-of-service problems