Future Computing Platforms Challenges and Opportunities

Onur Mutlu

omutlu@gmail.com

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

1 December 2021

IEEE Data & Storage Symposium (IEEE Bangalore)

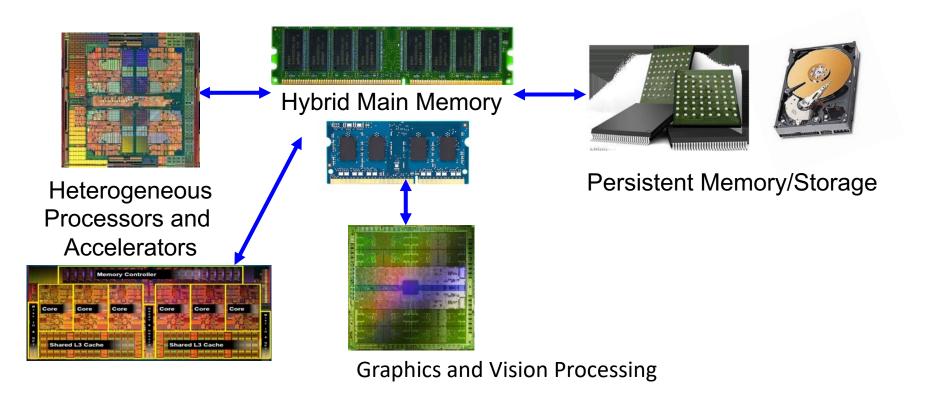




Carnegie Mellon

Current Research Mission

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



Build fundamentally better architectures

Four Key Current Directions

Fundamentally Secure/Reliable/Safe Architectures

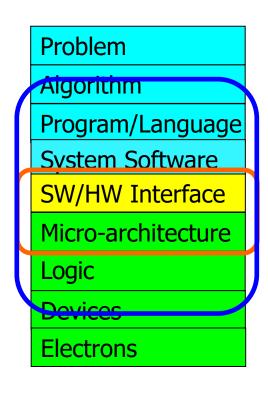
- Fundamentally Energy-Efficient Architectures
 - Memory-centric (Data-centric) Architectures

Fundamentally Low-Latency and Predictable Architectures

Architectures for AI/ML, Genomics, Medicine, Health

The Transformation Hierarchy

Computer Architecture (expanded view)



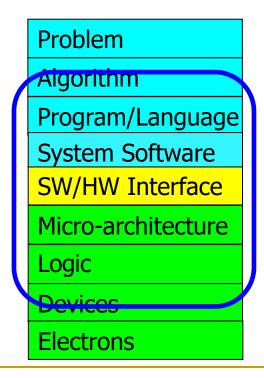
Computer Architecture (narrow view)

Axiom

To achieve the highest energy efficiency and performance:

we must take the expanded view

of computer architecture

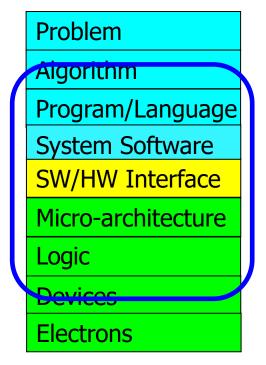


Co-design across the hierarchy:
Algorithms to devices

Specialize as much as possible within the design goals

Current Research Mission & Major Topics

Build fundamentally better architectures



Broad research spanning apps, systems, logic with architecture at the center

- Data-centric arch. for low energy & high perf.
 - Proc. in Mem/DRAM, NVM, unified mem/storage
- Low-latency & predictable architectures
 - Low-latency, low-energy yet low-cost memory
 - QoS-aware and predictable memory systems
- Fundamentally secure/reliable/safe arch.
 - Tolerating all bit flips; patchable HW; secure mem
- Architectures for ML/AI/Genomics/Health/Med
 - Algorithm/arch./logic co-design; full heterogeneity
- Data-driven and data-aware architectures
 - ML/AI-driven architectural controllers and design
 - Expressive memory and expressive systems

Onur Mutlu's SAFARI Research Group

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

https://safari.ethz.ch/safari-newsletter-april-2020/



Think BIG, Aim HIGH!

SAFARI

https://safari.ethz.ch

SAFARI Newsletter January 2021 Edition

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





Newsletter January 2021

Think Big, Aim High, and Have a Wonderful 2021!



Dear SAFARI friends,

Principle: Teaching and Research

Teaching drives Research Research drives Teaching

. . .

Research & Teaching: Some Overview Talks

https://www.youtube.com/onurmutlulectures

- Future Computing Architectures
 - https://www.youtube.com/watch?v=kqiZISOcGFM&list=PL5Q2soXY2Zi8D 5MGV6EnXEJHnV2YFBJI&index=1
- Enabling In-Memory Computation
 - https://www.youtube.com/watch?v=njX 14584Jw&list=PL5Q2soXY2Zi8D 5MGV6EnXEJHnV2YFBJl&index=16
- Accelerating Genome Analysis
 - https://www.youtube.com/watch?v=r7sn41lH-4A&list=PL5Q2soXY2Zi8D_5MGV6EnXEJHnV2YFBJl&index=41
- Rethinking Memory System Design
 - https://www.youtube.com/watch?v=F7xZLNMIY1E&list=PL5Q2soXY2Zi8D_5MGV6EnXEJHnV2YFBJl&index=3
- Intelligent Architectures for Intelligent Machines
 - https://www.youtube.com/watch?v=c6_LgzuNdkw&list=PL5Q2soXY2Zi8D_5MGV6EnXEJHnV2YFBJl&index=25
- The Story of RowHammer
 - https://www.youtube.com/watch?v=sgd7PHQQ1AI&list=PL5Q2soXY2Zi8D_5MGV6EnXEJHnV2YFBJl&index=39

Online Courses & Lectures

First Computer Architecture & Digital Design Course

- Digital Design and Computer Architecture
- Spring 2021 Livestream Edition:
 https://www.youtube.com/watch?v=LbC0EZY8yw4&list=PL5Q2soXY2Zi_uej3aY39YB5pfW4SJ7LIN

Advanced Computer Architecture Course

- Computer Architecture
- Fall 2021 Livestream Edition:
 https://www.youtube.com/watch?v=c3mPdZA-Fmc&list=PL5Q2soXY2Zi9xidyIqBxUz7xRPS-wisBN
- Fall 2020 Edition: https://www.youtube.com/watch?v=4yfkM_5EFgo&list=PL5Q2 soXY2Zi-Mnk1PxjEIG32HAGILkTOF

HOME

VIDEOS

PLAYLISTS

COMMUNITY

CHANNELS

ABOUT

Q

Popular uploads

▶ PLAY ALL

1:33:25

How Computers Work (from the ground up)

Digital Design & Computer

Architecture: Lecture 1:...

49K views • 1 year ago



Computer Architecture -

Lecture 1: Introduction and... 36K views • 3 years ago



Computer Architecture -Lecture 1: Introduction and...

30K views • 8 months ago 31K views • 1 year ago



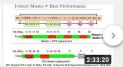
Computer Architecture -Lecture 1: Introduction and...



1:22:29

Design of Digital Circuits -Lecture 1: Introduction and...

22K views • 2 years ago



Computer Architecture -Lecture 2: Fundamentals....

17K views • 3 years ago

First Course in Computer Architecture & Digital Design 2021-2013



Livestream - Digital Design and Digital Design & Computer

Onur Mutlu Lectures VIEW FULL PLAYLIST

38 How Comput =, (from the grou

Architecture - ETH Zürich... Computer Architecture - ETH...

> Onur Mutlu Lectures VIEW FULL PLAYLIST



Design of Digital Circuits - ETH Zürich - Spring 2019

Onur Mutlu Lectures VIEW FULL PLAYLIST



Design of Digital Circuits - ETH Zürich - Spring 2018

Onur Mutlu Lectures VIEW FULL PLAYLIST



Digital Circuits and Computer Architecture - ETH Zurich -...

Onur Mutlu Lectures VIEW FULL PLAYLIST



Spring 2015 -- Computer Architecture Lectures --...

Carnegie Mellon Computer Architec... VIEW FULL PLAYLIST

Advanced Computer Architecture Courses 2020-2012



Computer Architecture - ETH Zürich - Fall 2020

Onur Mutlu Lectures VIEW FULL PLAYLIST



Computer Architecture - ETH Zürich - Fall 2019

Onur Mutlu Lectures VIEW FULL PLAYLIST



Computer Architecture - ETH Zürich - Fall 2018

Onur Mutlu Lectures VIEW FULL PLAYLIST



Computer Architecture - ETH Zürich - Fall 2017

Onur Mutlu Lectures VIEW FULL PLAYLIST



Fall 2015 - 740 Computer Architecture

Carnegie Mellon Computer Architec... VIEW FULL PLAYLIST



Fall 2013 - 740 Computer Architecture - Carnegie Mellon

Carnegie Mellon Computer Architec... VIEW FULL PLAYLIST

Special Courses on Memory Systems



Memory Technology Lectures Onur Mutlu Lectures

VIEW FULL PLAYLIST



Memory Systems and Memory... 2019

Onur Mutlu Lectures VIEW FULL PLAYLIST



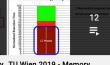
Champéry Winter School 2020 - Perugia NiPS Summer School

Onur Mutlu Lectures VIEW FULL PLAYLIST



SAMOS Tutorial 2019 - Memory TU Wien 2019 - Memory Systems

Onur Mutlu Lectures VIEW FULL PLAYLIST



Systems and Memory-Centric...

Onur Mutlu Lectures VIEW FULL PLAYLIST



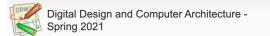
ACACES 2018 Lectures --Memory Systems and Memory...

Onur Mutlu Lectures VIEW FULL PLAYLIST



DDCA (Spring 2021)

- https://safari.ethz.ch/digitaltechnik/ spring2021/doku.php?id=schedule
- https://www.youtube.com/watch?v =LbC0EZY8yw4&list=PL5Q2soXY2Zi uej3aY39YB5pfW4SJ7LIN
- Bachelor's course
 - 2nd semester at ETH Zurich
 - Rigorous introduction into "How Computers Work"
 - Digital Design/Logic
 - Computer Architecture
 - 10 FPGA Lab Assignments



Recent Changes Media Manager Sitemap

schedule

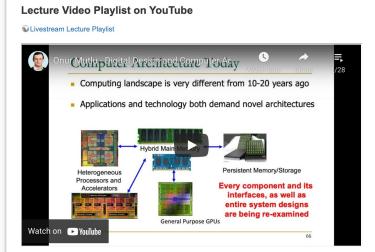
Trace: - schedule

Announcements

- Lectures/Schedule
- Lecture Buzzwords
- Readings Ontional HWs
- Extra Assignments
- Technical Docs

Exams

- Secondary Computer Architecture (CMU)
- SS15: Lecture Videos Computer Architecture (CMU)
- SS15: Course Website
- S Digitaltechnik SS18: Lecture Spigitaltechnik SS18: Course
- Website Specified in the second of the
- Digitaltechnik SS19: Course
- Website Digitaltechnik SS20: Lecture
- Videos Spigitaltechnik SS20: Course
- Website
- Moodle Moodle



Recorded Lecture Playlist



Spring 2021 Lectures/Schedule

Week	Date	Livestream	Lecture	Readings	Lab	HW	
W1	25.02 Thu.	You Tube Live	L1: Introduction and Basics	Required Suggested Mentioned			
	26.02 Fri.	You Tube Live	L2a: Tradeoffs, Metrics, Mindset	Required			
			L2b: Mysteries in Computer Architecture (PDF) (PPT)	Required Mentioned			
W2	04.03 Thu.	You Tube Live	L3a: Mysteries in Computer Architecture II	Required Suggested			



https://www.youtube.com/watch?v=c3 mPdZA-Fmc&list=PL5Q2soXY2Zi9xidyIqBxUz7x **RPS-wisBN**

- Master's level course
 - Taken by Bachelor's/Masters/PhD students
 - Cutting-edge research topics + fundamentals in Computer Architecture
 - 5 Simulator-based Lab Assignments
 - Potential research exploration
 - Many research readings



Q Recent Changes Media Manager Sitemap

schedule

Trace: · start · schedule

Announcements

Materials

- Lectures/Schedule
- Lecture Buzzwords

- Exams Related Courses

- Computer Architecture FS19 Course Webpage
- Computer Architecture FS19:
- Lecture Videos Digitaltechnik SS20: Course
- Webpage Digitaltechnik SS20: Lecture Videos
- Moodle Moodle
- Piazza (Q&A)
- **S** HotCRP
- Verilog Practice Website

Lecture Video Playlist on YouTube



Fall 2020 Lectures & Schedule

Week	Date	Lecture	Readings	Lab	HW
W1	17.09 Thu.	L1: Introduction and Basics (PDF) (PPT) You Video	Described Suggested		HW 0
	18.09 Fri.	L2a: Memory Performance Attacks (PDF) (PPT) Voulton Video	Described Suggested	Lab 1 Out	
		L2b: Data Retention and Memory Refresh (PDF) (PPT) Vou Video	Described Suggested		
		L2c: Course Logistics (PDF) (PPT) You the Video			
W2	24.09 Thu.	L3a: Introduction to Genome Sequence Analysis (PDF) (PPT) (Volume Video	Described Suggested		HW 1 Out
		L3b: Memory Systems: Challenges and Opportunities (PDF) (PPT) Vou Video	Described Suggested		
	25.09 Fri.	L4a: Memory Systems: Solution Directions (PDF) (PPT) You Video	Described Suggested		
		L4b: RowHammer (PDF) (PPT) Vou Video	Described Suggested		
W3	01.10 Thu.	L5a: RowHammer in 2020: TRRespass (PDF) (PPT) Vou Video	Described Suggested		
		L5b: RowHammer in 2020: Revisiting RowHammer (PDF) (PPT) Vou Video	Described Suggested		
		L5c: Secure and Reliable Memory	Described		

Comp Arch (Current)

https://safari.ethz.ch/architecture/fall20 21/doku.php?id=schedule

Youtube Livestream:

https://www.youtube.com/watch?v=4yfk M 5EFgo&list=PL5Q2soXY2Zi-Mnk1PxjEIG32HAGILkTOF

Master's level course

- Taken by Bachelor's/Masters/PhD students
- Cutting-edge research topics + fundamentals in Computer Architecture
- 5 Simulator-based Lab Assignments
- Potential research exploration
- Many research readings



cent Changes Media Manager Sitema

schedule

tooth onangoo moda managor onoma

Trace: · readings · start · schedule

Home

Announcements

Materials

- Lectures/Schedule
- Lecture Buzzwords
- Readings
- HWs
- LabsExams
 - Related Courses
- Tutorials

Pecaurees

- Course Webpage
- Computer Architecture FS20:
- Digitaltechnik SS21: Course
- Digitaltechnik SS21: Lecture Videos
- Moodle
- MotCRP
- S Verilog Practice Website (HDLBits)

Lecture Video Playlist on YouTube



Recorded Lecture Playlist

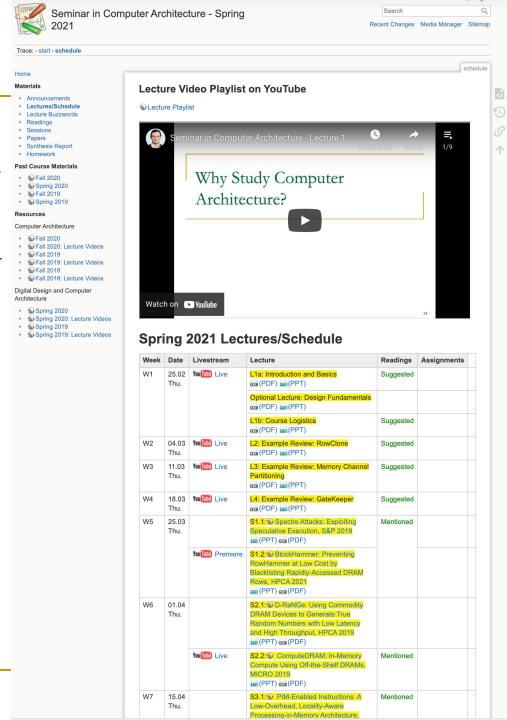


Fall 2021 Lectures & Schedule

Week	Date	Livestream	Lecture	Readings	Lab	HW
W1	30.09 Thu.	You Tube Live	L1: Introduction and Basics	Required Mentioned	Lab 1 Out	HW 0 Out
	01.10 Fri.	You Tube Live	L2: Trends, Tradeoffs and Design Fundamentals (PDF) (PPT)	Required Mentioned		
W2	07.10 Thu.	You Tube Live	L3a: Memory Systems: Challenges and Opportunities	Described Suggested		HW 1 Out
			L3b: Course Info & Logistics			
			L3c: Memory Performance Attacks	Described Suggested		
	08.10 Fri.	You Tube Live	L4a: Memory Performance Attacks (PDF) (PPT)	Described Suggested	Lab 2 Out	
			L4b: Data Retention and Memory Refresh (PDF) (PPT)	Described Suggested		
			L4c: RowHammer	Described Suggested		

Seminar (Spring'21)

- <u>https://safari.ethz.ch/architecture_seminar/spring2021/doku.php?id=schedule</u>
- https://www.youtube.com/watch?v=t3m 93ZpLOyw&list=PL5Q2soXY2Zi awYdjm WVIUegsbY7TPGW4
- Critical analysis course
 - Taken by Bachelor's/Masters/PhD students
 - Cutting-edge research topics + fundamentals in Computer Architecture
 - 20+ research papers, presentations, analyses



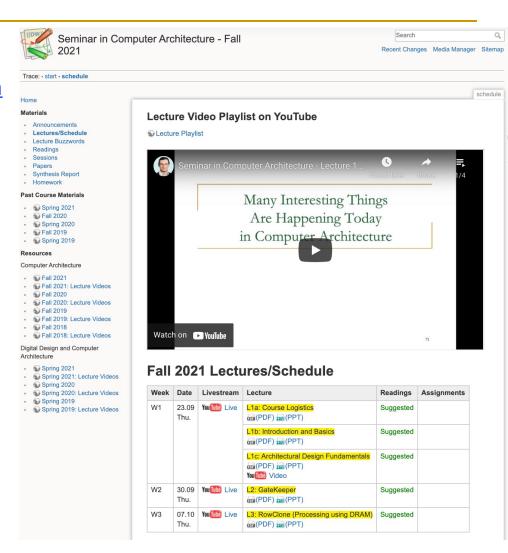


Seminar (Current)

https://safari.ethz.ch/architecture_semin ar/fall2021/doku.php?id=schedule

Youtube Livestream:

- https://www.youtube.com/watch?v=4TcP 297mdsI&list=PL5Q2soXY2Zi 7UBNmC9B 8Yr5JSwTG9yH4
- Critical analysis course
 - Taken by Bachelor's/Masters/PhD students
 - Cutting-edge research topics + fundamentals in Computer Architecture
 - 20+ research papers, presentations, analyses



Hands-On Projects & Seminars Courses

https://safari.ethz.ch/projects_and_seminars/doku.php



Search

start

Recent Changes Media Manager Sitemap

Trace: · start

Projects

Home

- SoftMC
- Ramulator
- Accelerating Genomics
- Mobile Genomics
- Processing-in-Memory
- Heterogeneous Systems
- SSD Simulator

SAFARI Projects & Seminars Courses (Spring 2021)

Welcome to the wiki for Project and Seminar courses SAFARI offers.

Courses we offer:

- Understanding and Improving Modern DRAM Performance, Reliability, and Security with Hands-On **Experiments**
- Designing and Evaluating Memory Systems and Modern Software Workloads with Ramulator
- Accelerating Genome Analysis with FPGAs, GPUs, and New Execution Paradigms
- Genome Sequencing on Mobile Devices
- Exploring the Processing-in-Memory Paradigm for Future Computing Systems
- Hands-on Acceleration on Heterogeneous Computing Systems
- Understanding and Designing Modern NAND Flash-Based Solid-State Drives (SSDs) by Building a **Practical SSD Simulator**



SAFARI Live Seminars (I)



SAFARI Live Seminars (II)



SAFARI Live Seminar: Nastaran Hajinazar 27 Oct 2021

Posted on October 1, 2021 by ewent

Join us for our SAFARI Live Seminar with Nastaran Hajinazar.

Wednesday, October 27 at 7:00 pm Zurich time (CEST)



SAFARI Live Seminar: Gennady Pekhimenko 08 Nov 2021

Posted on November 1, 2021 by ewent

Join us for our SAFARI Live Seminar with Gennady Pekhimenko.

Monday, November 08 at 4:00 pm Zurich time (CET)



SAFARI Live Seminar: Damla Senol Cali 07 Nov 2021

Posted on October 18, 2021 by ewent

Join us for our SAFARI Live Seminar with Damla Senol Cali.

Sunday, November 07 at 6:00 pm Zurich time (CEST)

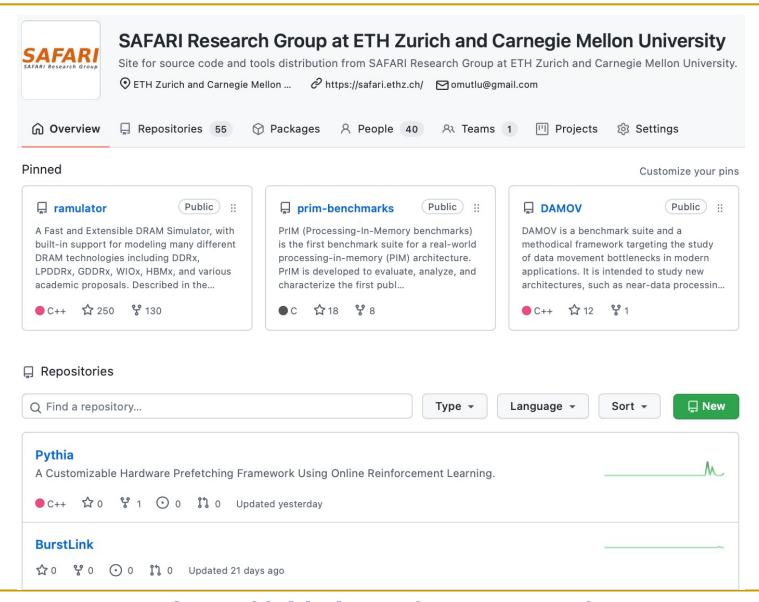


SAFARI Live Seminar: Serghei Mangul 11 Nov 2021
Posted on November 5, 2021 by ewent

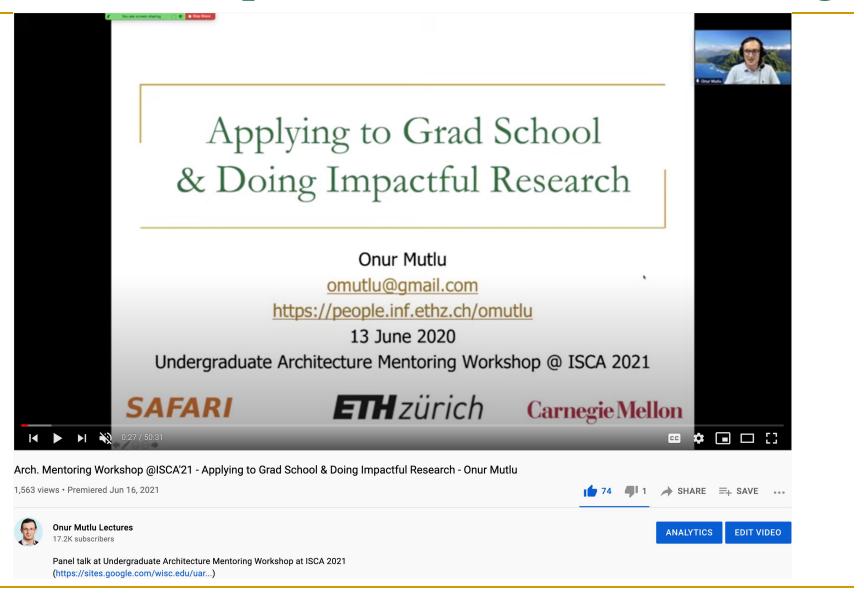
Join us for our SAFARI Live Seminar with Serghei Mangul.

Thursday, November 11 at 11:00 am Zurich time (CET), ETH Zentrum ETZ K91

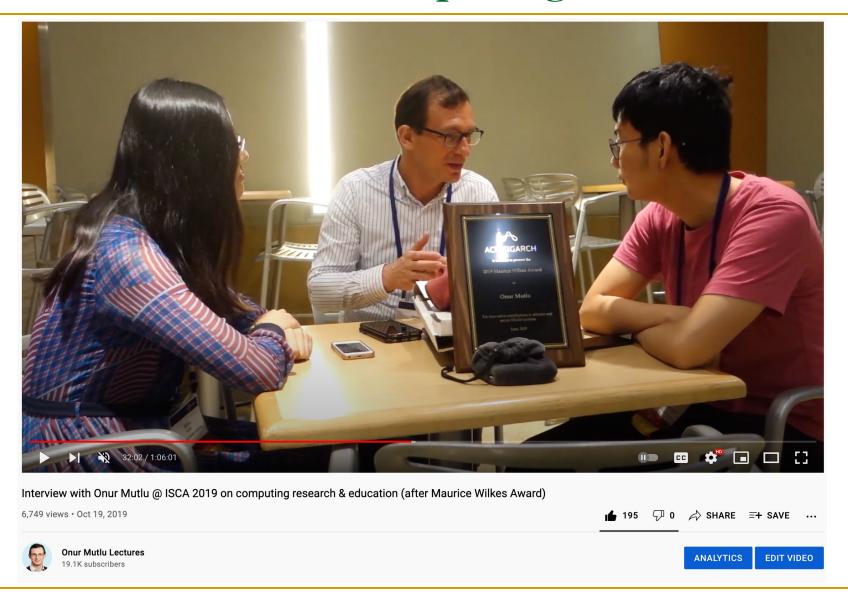
Open Source Tools: SAFARI GitHub



A Talk on Impactful Research & Teaching



An Interview on Computing Futures



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/

Future Computing Platforms Challenges and Opportunities

Why Do We Do Computing?

To Solve Problems

To Gain Insight

To Enable a Better Life & Future

How Does a Computer Solve Problems?

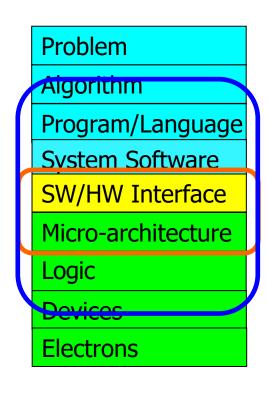
Orchestrating Electrons

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

Different Platforms, Different Goals





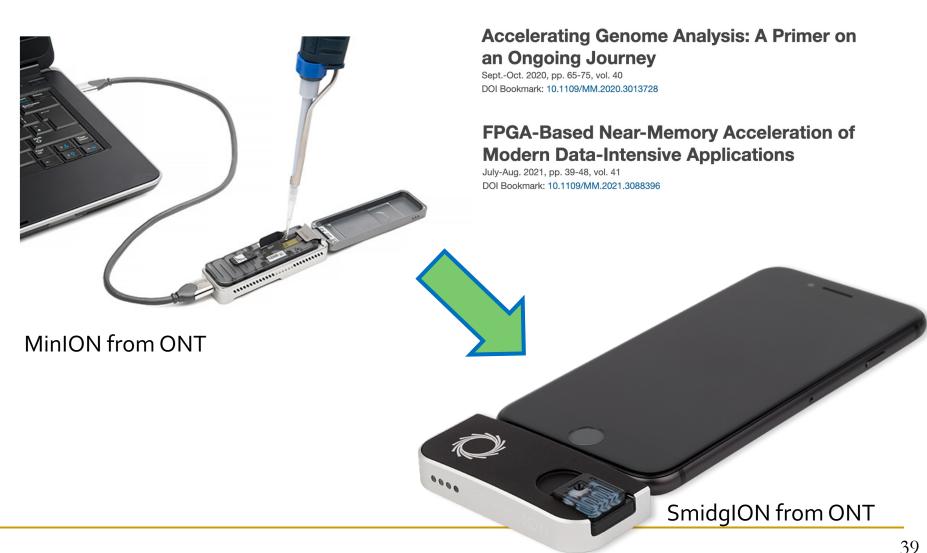
Different Platforms, Different Goals

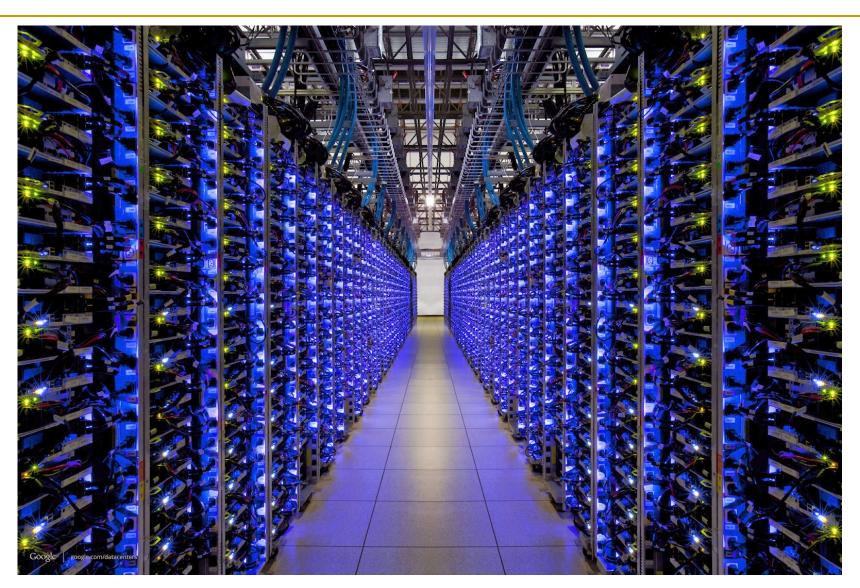






Mohammed Alser, Zülal Bingöl, Damla Senol Cali, Jeremie Kim, Saugata Ghose, Can Alkan, Onur Mutlu "Accelerating Genome Analysis: A Primer on an Ongoing Journey" IEEE Micro, August 2020.







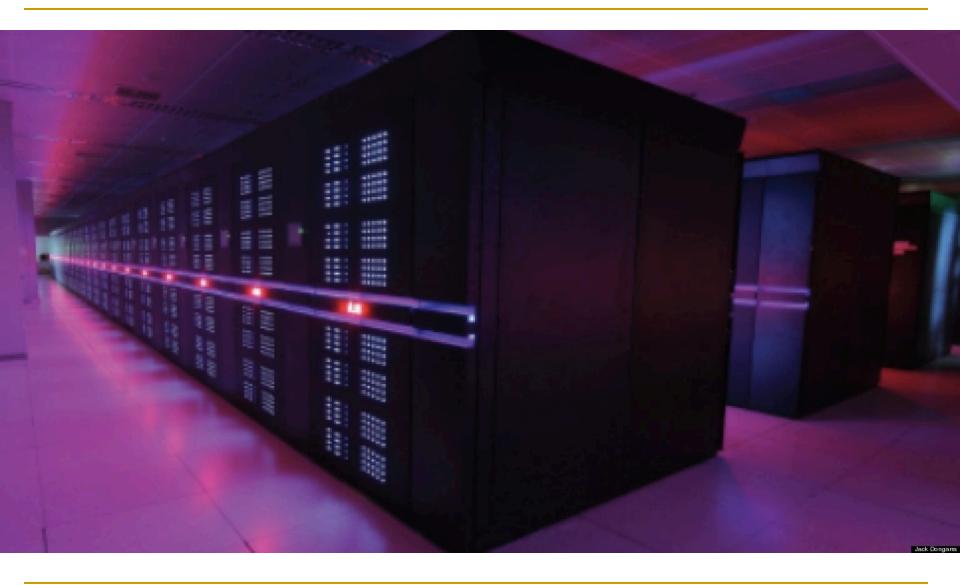




Figure 3. TPU Printed Circuit Board. It can be inserted in the slot for an SATA disk in a server, but the card uses PCIe Gen3 x16.

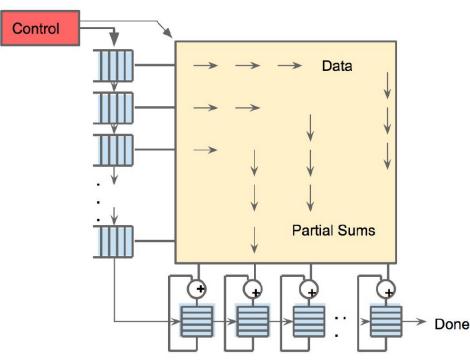
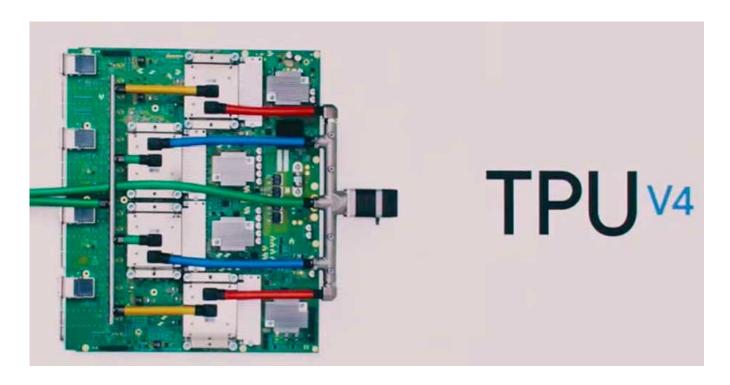


Figure 4. Systolic data flow of the Matrix Multiply Unit. Software has the illusion that each 256B input is read at once, and they instantly update one location of each of 256 accumulator RAMs.

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

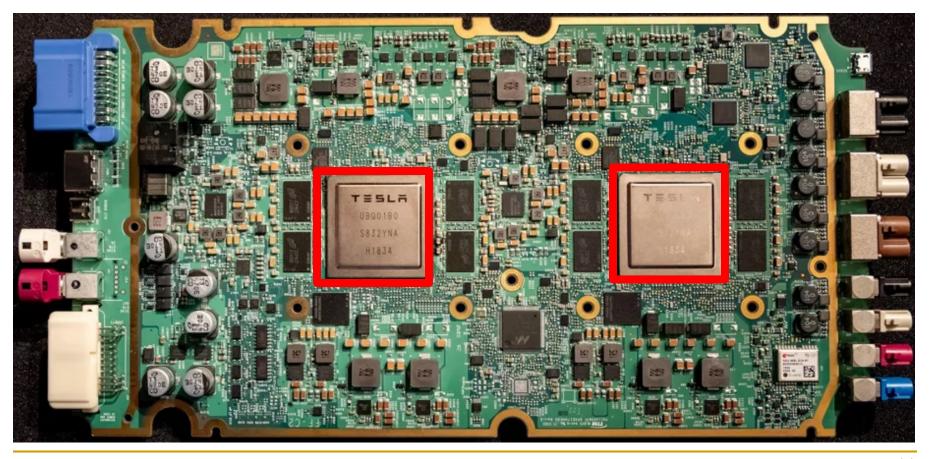


1 ExaFLOPS per board

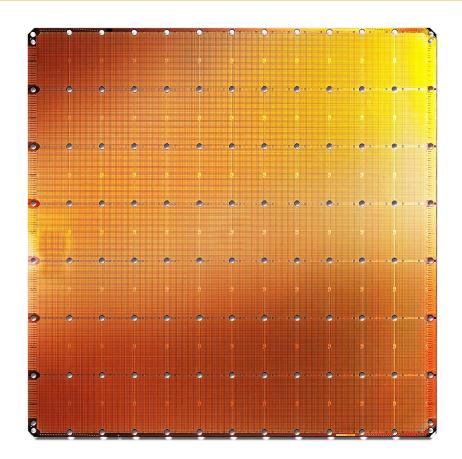
https://spectrum.ieee.org/tech-talk/computing/hardware/heres-how-googles-tpu-v4-ai-chip-stacked-up-in-training-tests

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





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

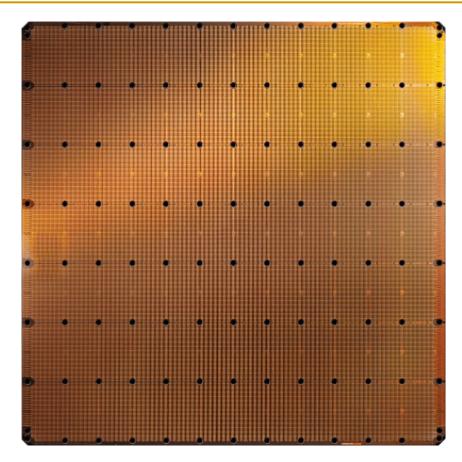
21.1 Billion transistors 815 mm²

NVIDIA TITAN V

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/

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

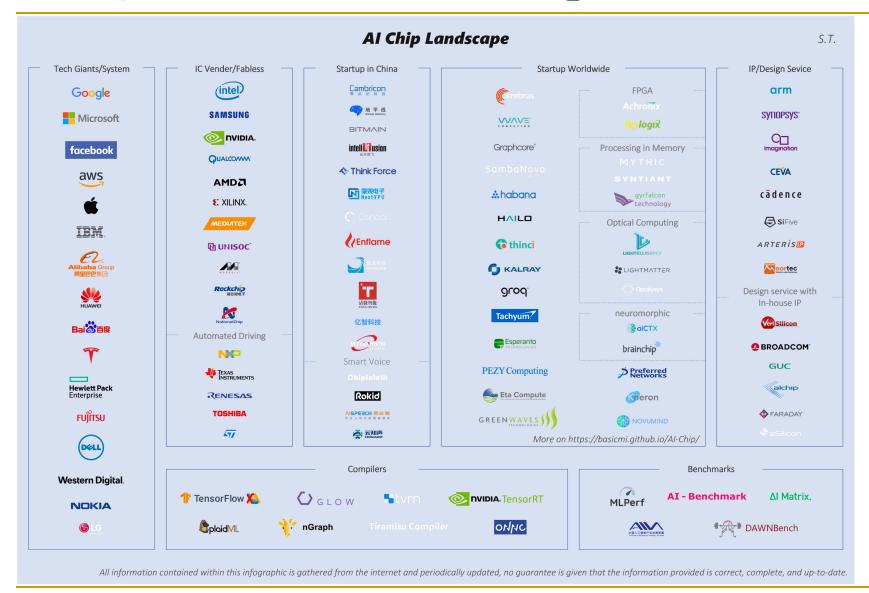
NVIDIA Ampere GA100

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

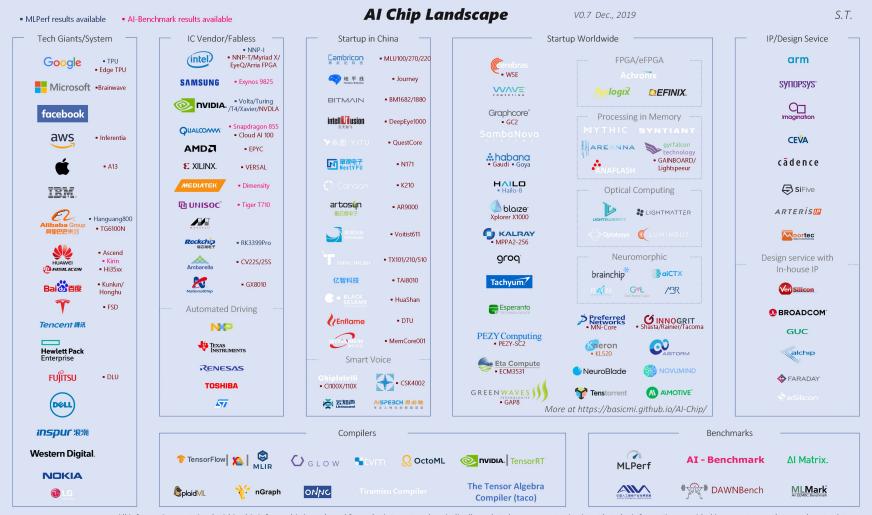
Many (Other) AI/ML Chips

- Alibaba
- Amazon
- Facebook
- Google
- Huawei
- Intel
- Microsoft
- NVIDIA
- Tesla
- Many Others and Many Startups are Building Their Own Chips...
- Many More to Come...

Many (Other) AI/ML Chips (2019)



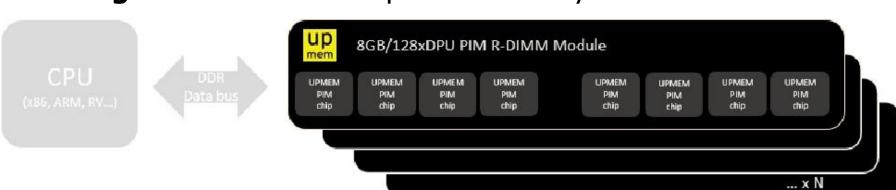
Many (Other) AI/ML Chips (2021)



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.

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





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 3D-stacked 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 state-of-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

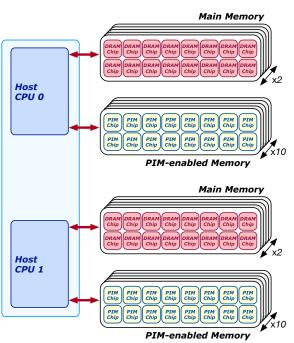
UPMEM Memory Modules

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





2,560-DPU Processing-in-Memory System



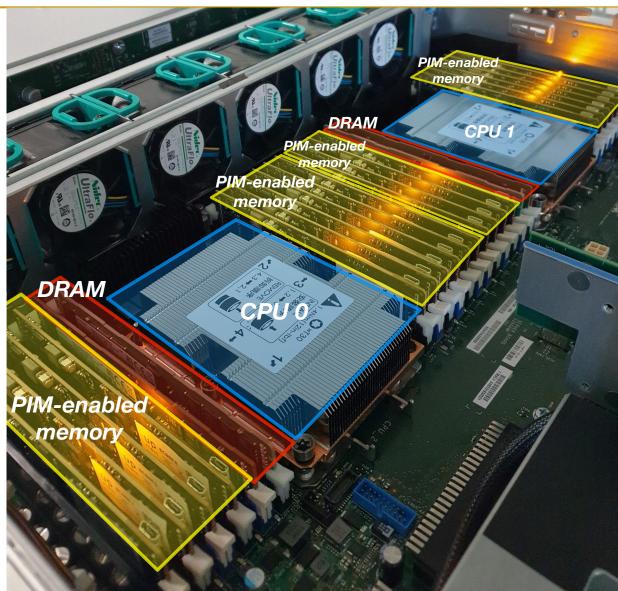
Benchmarking a New Paradigm: An Experimental Analysis of a Real Processing-in-Memory Architecture

JUAN GÓMEZ-LUNA, ETH Zürich, Switzerland
IZZAT EL HAJJ, 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 (PM).

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.

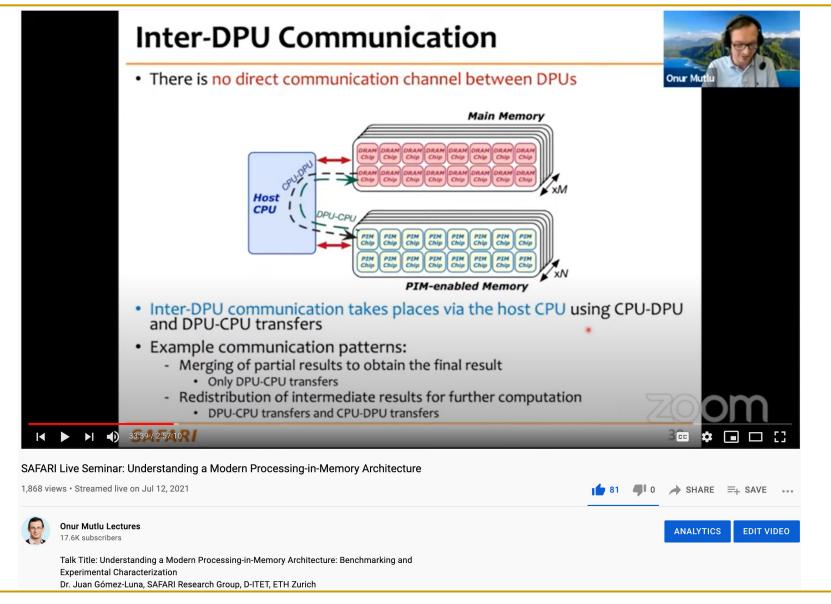
This paper provides the first comprehensive analysis of the first publicly-available real-world PIM architecture. We make two key contributions. First, we conduct an experimental characterization of the UPMEM-based PIM system using microbenchmarks to assess various architecture limits such as compute throughput and memory bandwidth, yielding new insights. Second, we present PIM (Processing,-bendumpy) benchmarks), a benchmark suite of 16 workloads from different application domains (e.g., dense/sparse linear algebra, databases, data analytics, graph processing, which we identify as memory-bound. We evaluate the performance and scaling characteristics of PIM benchmarks on the UPMEM PIM architecture, and compare their performance and energy consumption to their state-of-the-art CPU and CPU counterparts. Our extensive evaluation conducted on two real UPMEM-based PIM systems with 460 and 25.50 DPUs provides new insights about suitability of different workloads to the PIM systems you for the programming recommendations for software designers, and suggestions and hints for hardware and architecture designers of future PIM systems.



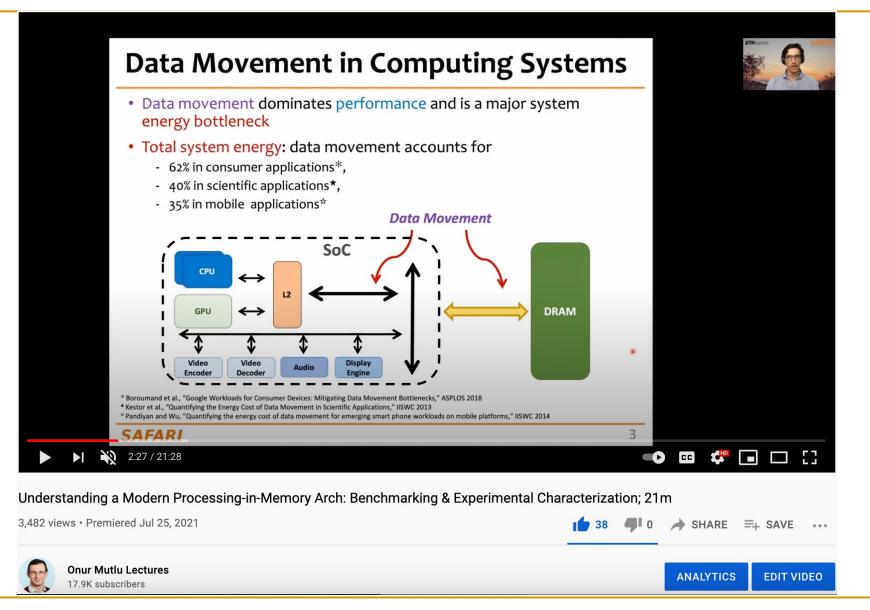
Understanding a Modern PIM Architecture



More on Analysis of the UPMEM PIM Engine



More on Analysis of the UPMEM PIM Engine



FPGA-based Processing Near Memory

Gagandeep Singh, Mohammed Alser, Damla Senol Cali, Dionysios
Diamantopoulos, Juan Gómez-Luna, Henk Corporaal, and Onur Mutlu,

"FPGA-based Near-Memory Acceleration of Modern Data-Intensive

Applications"

IEEE Micro (IEEE MICRO), to appear, 2021.

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

Gagandeep Singh[⋄] Mohammed Alser[⋄] Damla Senol Cali[⋈]
Dionysios Diamantopoulos[▽] Juan Gómez-Luna[⋄]
Henk Corporaal[⋆] Onur Mutlu^{⋄⋈}

[⋄]ETH Zürich [⋈] Carnegie Mellon University *Eindhoven University of Technology [▽]IBM Research Europe

Samsung Function-in-Memory DRAM (2021)

Samsung Newsroom

CORPORATE

PRODUCTS

PRESS RESOURCES

VIEWS

ABOUT US



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

Korea on February 17, 2021

Audio



Share (5)





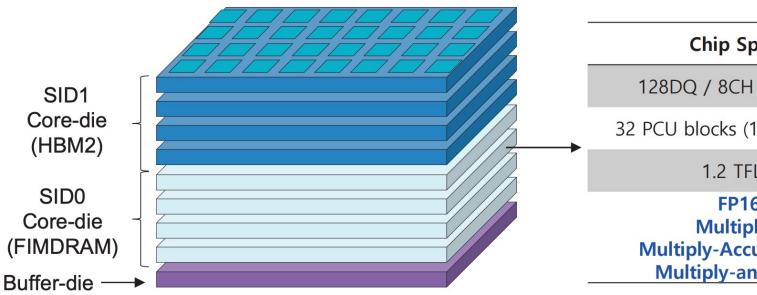
The new architecture will deliver over twice the system performance and reduce energy consumption by more than 70%

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 Al-driven workloads such as HPC, training and inference. We plan to build upon this breakthrough by further collaborating with Al solution providers for even more advanced PIM-powered applications."

Samsung Function-in-Memory DRAM (2021)

FIMDRAM based on HBM2



[3D Chip Structure of HBM with FIMDRAM]

Chip Specification

128DQ / 8CH / 16 banks / BL4

32 PCU blocks (1 FIM block/2 banks)

1.2 TFLOPS (4H)

FP16 ADD /
Multiply (MUL) /
Multiply-Accumulate (MAC) /
Multiply-and- Add (MAD)

ISSCC 2021 / SESSION 25 / DRAM / 25.4

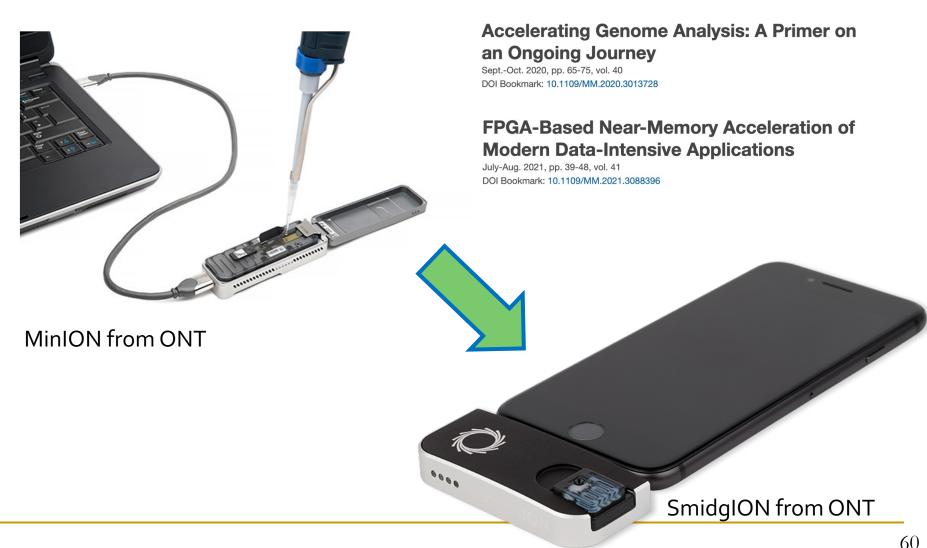
25.4 A 20nm 6GB Function-In-Memory DRAM, Based on HBM2 with a 1.2TFLOPS Programmable Computing Unit Using Bank-Level Parallelism. for Machine Learning Applications

Young-Cheon Kwon', Suk Han Lee', Jaehoon Lee', Sang-Hyuk Kwon', Je Min Ryu', Jong-Pii 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

Future of Genome Sequencing & Analysis

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

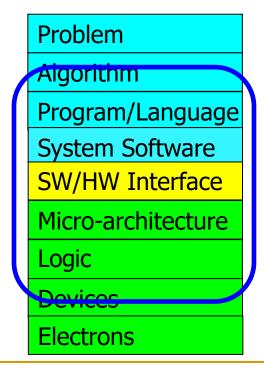




To achieve the highest energy efficiency and performance:

we must take the expanded view

of computer architecture

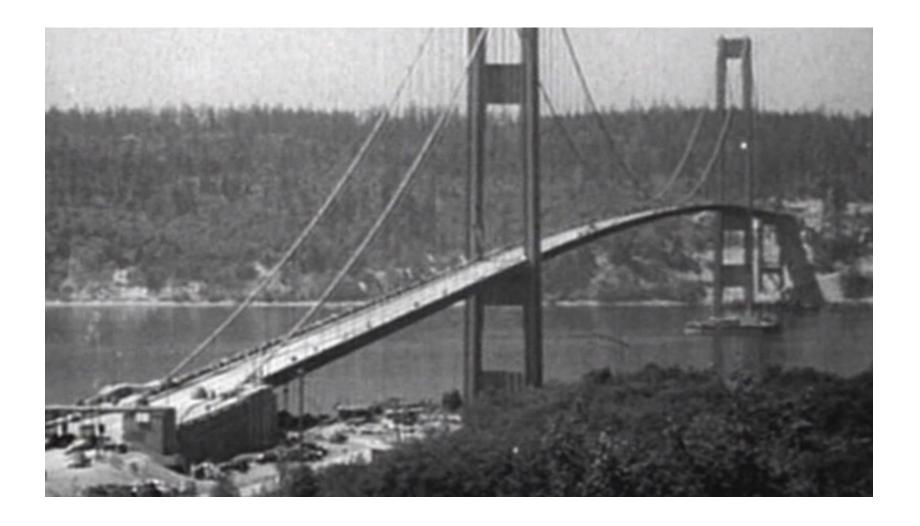


Co-design across the hierarchy:
Algorithms to devices

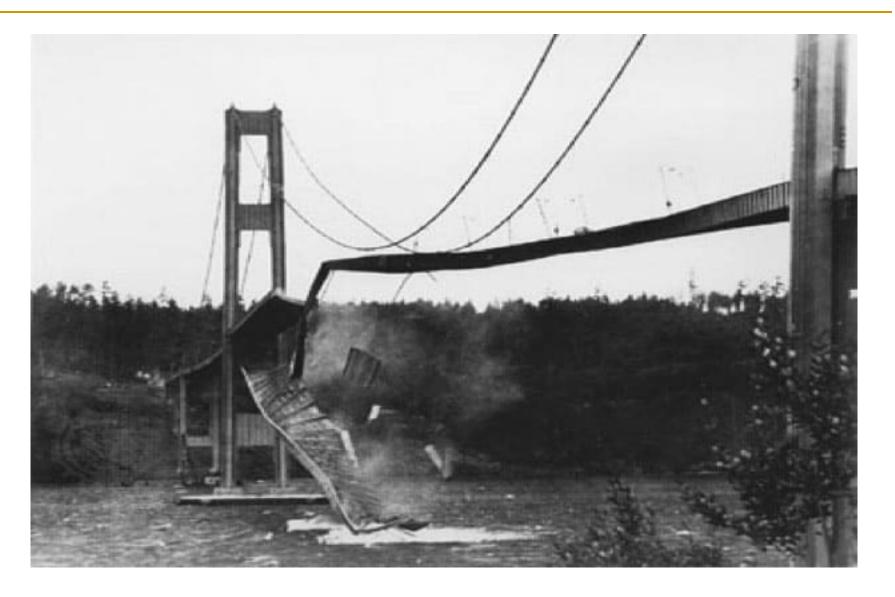
Specialize as much as possible within the design goals

What Kind of a Future Do We Want?

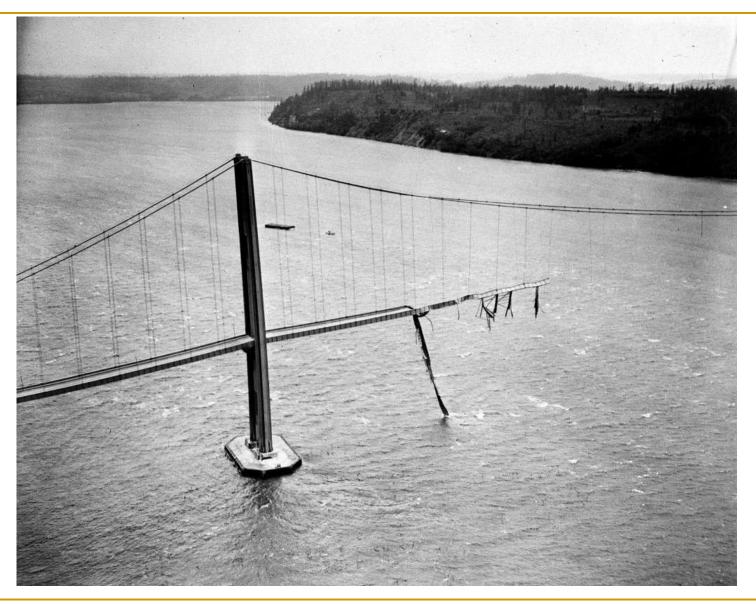
How Reliable/Secure/Safe is This Bridge?



Collapse of the "Galloping Gertie"



Another View



How Secure Are These People?



Security is about preventing unforeseen consequences

How Safe & Secure Is **This** Platform?



Challenge and Opportunity for Future

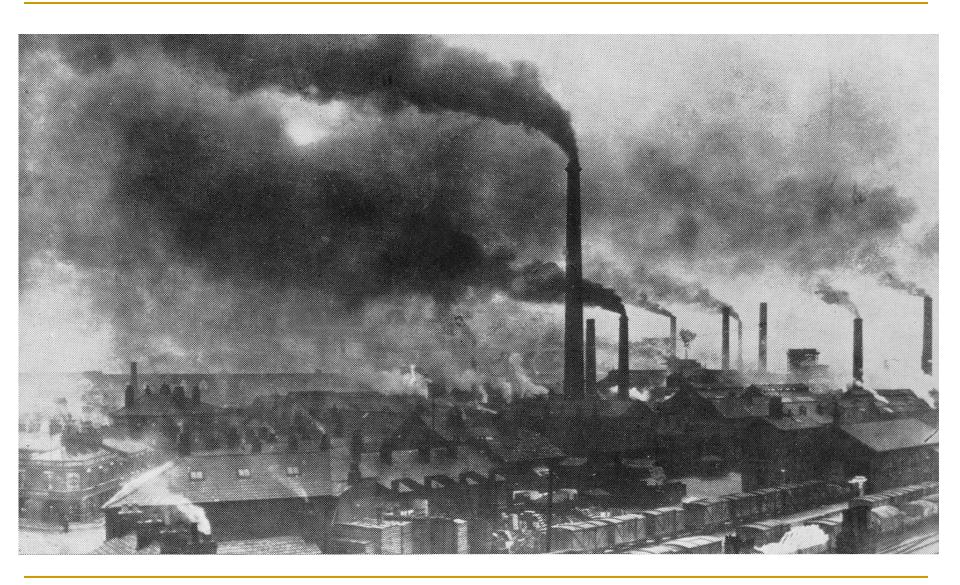
Reliable, Secure, Safe

Do We Want This?





Or This?



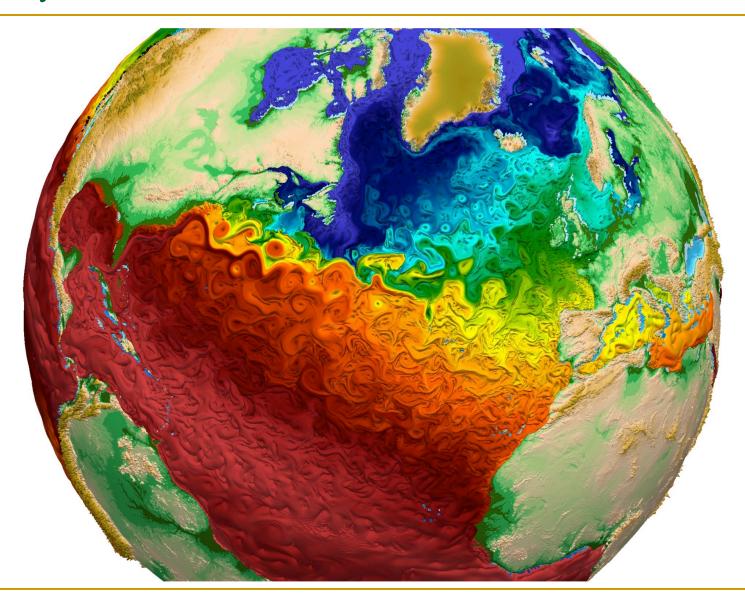
SAFARI

Source: V. Milutinovic 70

Challenge and Opportunity for Future

Sustainable and Energy Efficient

Many Difficult Problems: Climate



Many Difficult Problems: Congestion



Many Difficult Problems: Intelligence

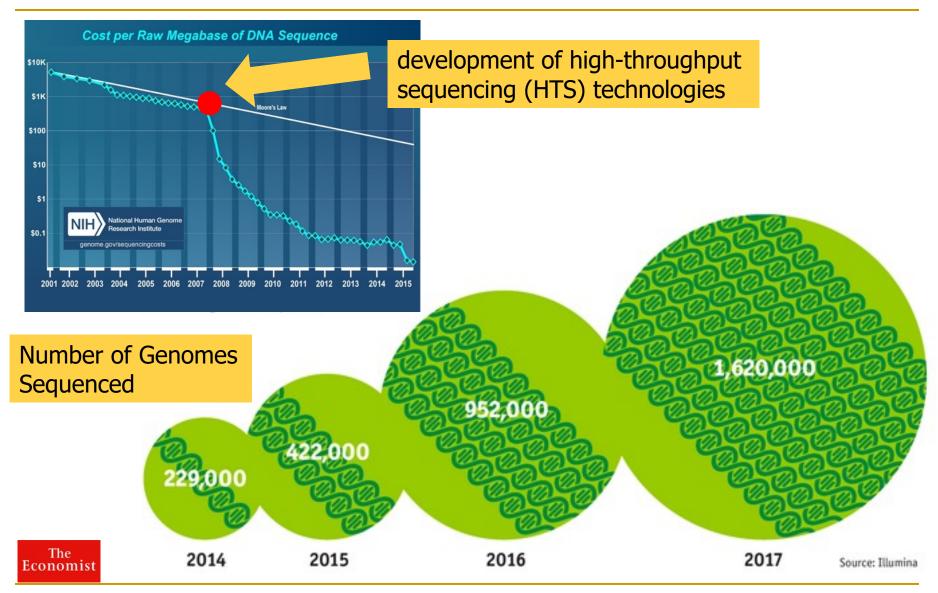




Many Difficult Problems: Public Health



Many Difficult Problems: Genome Analysis



Accelerating Genome Analysis

 Mohammed Alser, Zulal Bingol, Damla Senol Cali, Jeremie Kim, Saugata Ghose, Can Alkan, and Onur Mutlu,

"Accelerating Genome Analysis: A Primer on an Ongoing Journey"

IEEE Micro (IEEE MICRO), Vol. 40, No. 5, pages 65-75, September/October 2020.

[Slides (pptx)(pdf)]

[Talk Video (1 hour 2 minutes)]

Accelerating Genome Analysis: A Primer on an Ongoing Journey

Mohammed Alser

ETH Zürich

Zülal Bingöl

Bilkent University

Damla Senol Cali

Carnegie Mellon University

Jeremie Kim

ETH Zurich and Carnegie Mellon University

Saugata Ghose

University of Illinois at Urbana–Champaign and Carnegie Mellon University

Can Alkan

Bilkent University

Onur Mutlu

ETH Zurich, Carnegie Mellon University, and Bilkent University



GenASM Framework [MICRO 2020]

Damla Senol Cali, Gurpreet S. Kalsi, Zulal Bingol, Can Firtina, Lavanya Subramanian, Jeremie S. Kim, Rachata Ausavarungnirun, Mohammed Alser, Juan Gomez-Luna, Amirali Boroumand, Anant Nori, Allison Scibisz, Sreenivas Subramoney, Can Alkan, Saugata Ghose, and Onur Mutlu, "GenASM: A High-Performance, Low-Power Approximate String Matching Acceleration Framework for Genome Sequence Analysis"
Proceedings of the 53rd International Symposium on Microarchitecture (MICRO), Virtual, October 2020.

[<u>Lighting Talk Video</u> (1.5 minutes)] [<u>Lightning Talk Slides (pptx) (pdf)</u>] [<u>Talk Video</u> (18 minutes)] [<u>Slides (pptx) (pdf)</u>]

GenASM: A High-Performance, Low-Power Approximate String Matching Acceleration Framework for Genome Sequence Analysis

Damla Senol Cali^{†™} Gurpreet S. Kalsi[™] Zülal Bingöl[▽] Can Firtina[⋄] Lavanya Subramanian[‡] Jeremie S. Kim^{⋄†} Rachata Ausavarungnirun[⊙] Mohammed Alser[⋄] Juan Gomez-Luna[⋄] Amirali Boroumand[†] Anant Nori[™] Allison Scibisz[†] Sreenivas Subramoney[™] Can Alkan[▽] Saugata Ghose^{*†} Onur Mutlu^{⋄†▽}

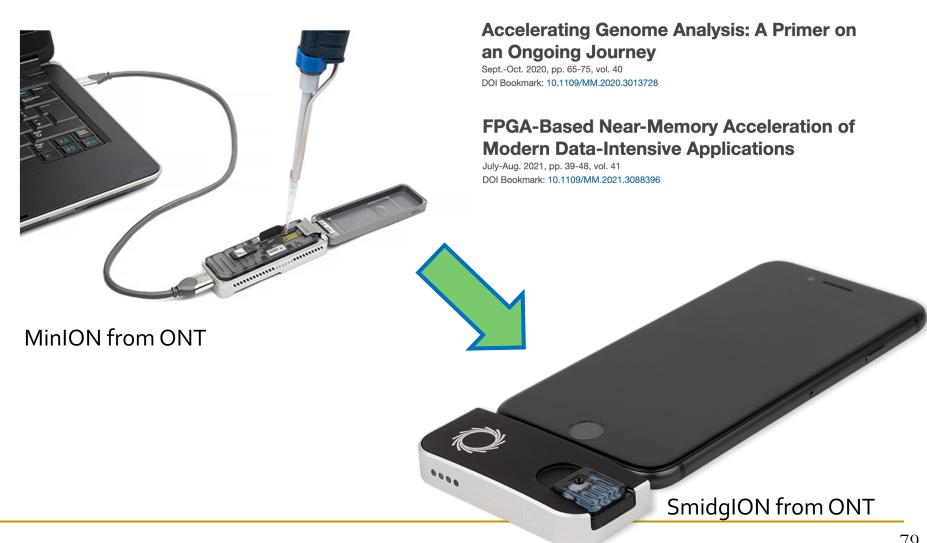
† Carnegie Mellon University [™] Processor Architecture Research Lab, Intel Labs [▽] Bilkent University [⋄] ETH Zürich

‡ Facebook [⊙] King Mongkut's University of Technology North Bangkok ^{*} University of Illinois at Urbana–Champaign

78

Future of Genome Sequencing & Analysis

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



More on Fast & Efficient Genome Analysis

Onur Mutlu,

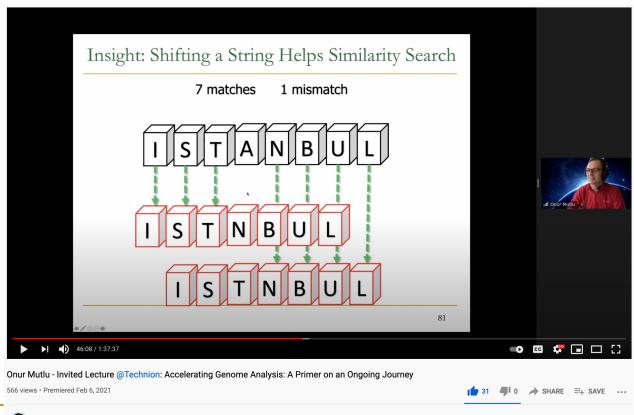
"Accelerating Genome Analysis: A Primer on an Ongoing Journey"

Invited Lecture at <u>Technion</u>, Virtual, 26 January 2021.

[Slides (pptx) (pdf)]

[Talk Video (1 hour 37 minutes, including Q&A)]

[Related Invited Paper (at IEEE Micro, 2020)]





ANALYTICS

Detailed Lectures on Genome Analysis

- Computer Architecture, Fall 2020, Lecture 3a
 - Introduction to Genome Sequence Analysis (ETH Zürich, Fall 2020)
 - https://www.youtube.com/watch?v=CrRb32v7SJc&list=PL5Q2soXY2Zi9xidyIgBxUz7xRPS-wisBN&index=5
- Computer Architecture, Fall 2020, Lecture 8
 - Intelligent Genome Analysis (ETH Zürich, Fall 2020)
 - https://www.youtube.com/watch?v=ygmQpdDTL7o&list=PL5Q2soXY2Zi9xidyIgBxUz7xRPS-wisBN&index=14
- Computer Architecture, Fall 2020, Lecture 9a
 - □ **GenASM: Approx. String Matching Accelerator** (ETH Zürich, Fall 2020)
 - https://www.youtube.com/watch?v=XoLpzmN Pas&list=PL5Q2soXY2Zi9xidyIgBxUz7xRPS-wisBN&index=15
- Accelerating Genomics Project Course, Fall 2020, Lecture 1
 - Accelerating Genomics (ETH Zürich, Fall 2020)
 - https://www.youtube.com/watch?v=rgjl8ZyLsAg&list=PL5Q2soXY2Zi9E2bBVAgCqL gwiDRQDTyId

Challenge and Opportunity for Future

High Performance

(to solve the **toughest** & **all** problems)

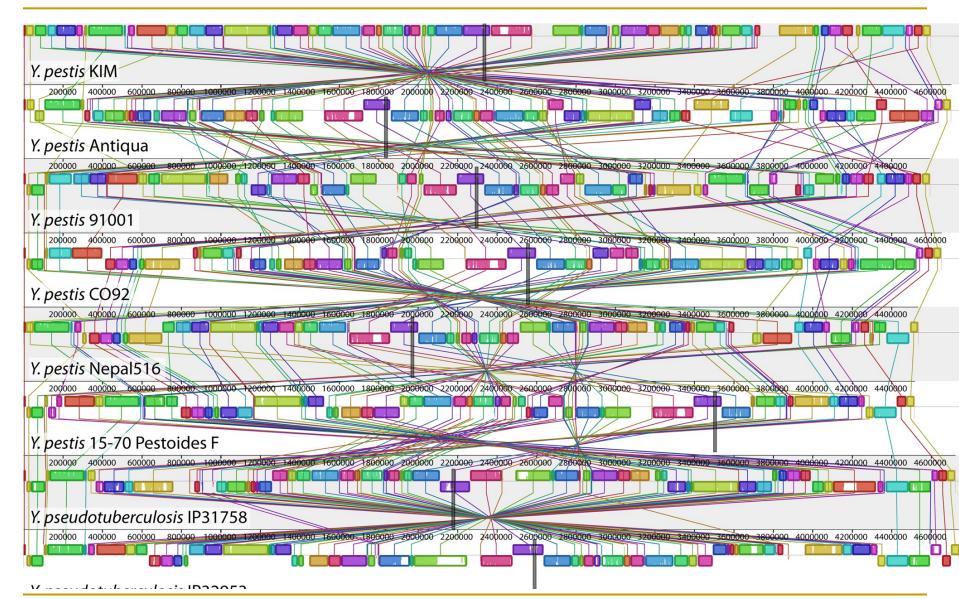
Personalized Medicine





83

Comparative Genomics



Personalized Genomics Technologies

Nanopore sequencing technology and tools for genome assembly: computational analysis of the current state, bottlenecks and future directions

Damla Senol Cali ™, Jeremie S Kim, Saugata Ghose, Can Alkan, Onur Mutlu

Briefings in Bioinformatics, bby017, https://doi.org/10.1093/bib/bby017

Published: 02 April 2018 Article history ▼



Oxford Nanopore MinION

Senol Cali+, "Nanopore Sequencing Technology and Tools for Genome Assembly: Computational Analysis of the Current State, Bottlenecks and Future Directions," Briefings in Bioinformatics, 2018.

[Preliminary arxiv.org version]

Personalized Robotics



Challenge and Opportunity for Future

Personalized and Private

```
(in every aspect of life:
health, medicine,
spaces, devices, robotics, ...)
```

This Lecture is About ...

 Questioning what limits us in designing the best computing architectures for the future

Providing directions for fundamentally better designs

Advocating principled approaches

Increasingly Demanding Applications

Dream...

and, they will come

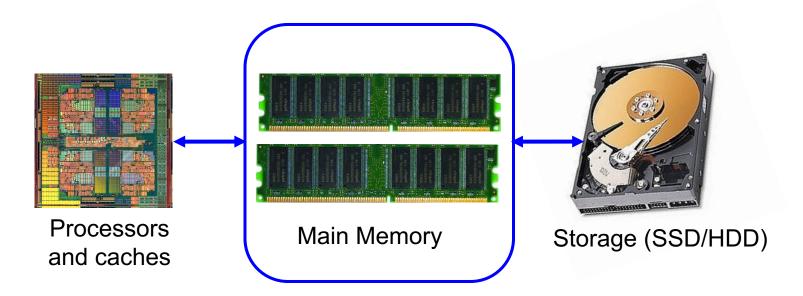
As applications push boundaries, computing platforms will become increasingly strained.

Key Realization

Modern Systems are Bottlenecked by

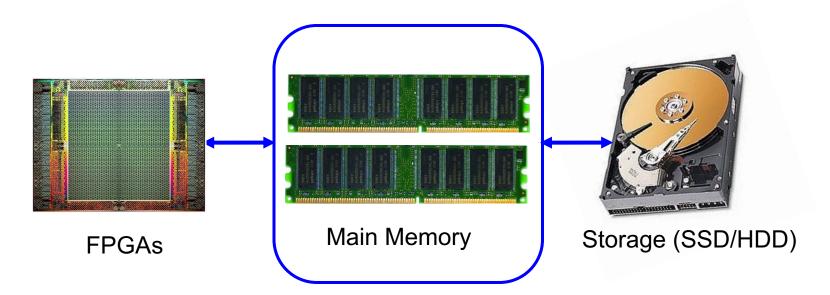
Data Storage and Movement

Focus is on Data Storage Systems (Memory)



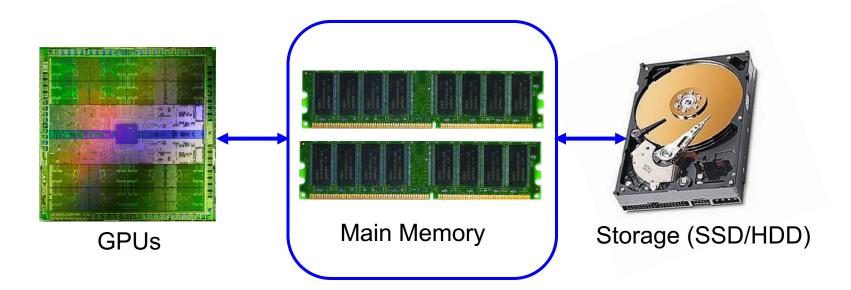
- Main memory is a critical component of all computing systems: server, mobile, embedded, desktop, sensor
- Main memory system must scale (in size, technology, efficiency, cost, and management algorithms) to maintain performance growth and technology scaling benefits

Focus is on Data Storage Systems (Memory)



- 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

Focus is on Data Storage Systems (Memory)



- 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

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

Memory Is Critical for Performance (I)



In-memory Databases

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



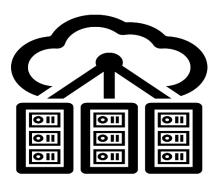
In-Memory Data Analytics

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



Graph/Tree Processing

[Xu+, IISWC'12; Umuroglu+, FPL'15]



Datacenter Workloads

[Kanev+ (Google), ISCA'15]



Memory Is Critical for Performance (I)



In-memory Databases



Graph/Tree Processing

Memory → bottleneck



In-Memory Data Analytics

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



Datacenter Workloads

[Kanev+ (Google), ISCA' 15]



Memory Is Critical for Performance (II)



Chrome

Google's web browser



TensorFlow Mobile

Google's machine learning framework



Google's video codec



Google's video codec

Memory Is Critical for Performance (II)





TensorFlow Mobile

Memory → bottleneck

VP9
VouTube
Video Playback

Google's video codec



Google's video codec

Short Read Alignment Reference Genome

Read Mapping

1 Sequencing

Genome Analysis

reference: TTTATCGCTTCCATGACGCAG

read1: ATCGCATCC read2: TATCGCATC

read3: CATCCATGA

read4: CGCTTCCAT

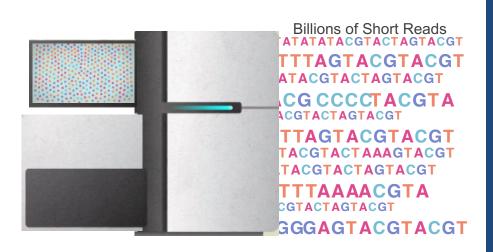
read5: CCATGACGC

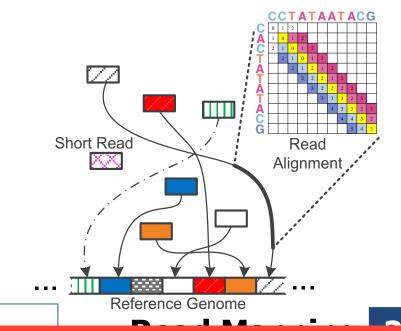
read6: TTCCATGAC



Variant Calling

Scientific Discovery 4





Memory → bottleneck

Tererence. ITTATCGCTTCCATGACGCAG

read1: ATCGCATCC read2: TATCGCATC

read3: CATCCATGA

read4: CGCTTCCAT

read5: CCATGACGC

read6: TTCCATGAC



Scientific Discovery 4

New Genome Sequencing Technologies

Nanopore sequencing technology and tools for genome assembly: computational analysis of the current state, bottlenecks and future directions

Damla Senol Cali ™, Jeremie S Kim, Saugata Ghose, Can Alkan, Onur Mutlu

Briefings in Bioinformatics, bby017, https://doi.org/10.1093/bib/bby017

Published: 02 April 2018 Article history ▼



Oxford Nanopore MinION

Senol Cali+, "Nanopore Sequencing Technology and Tools for Genome Assembly: Computational Analysis of the Current State, Bottlenecks and Future Directions," Briefings in Bioinformatics, 2018.

[Open arxiv.org version]

New Genome Sequencing Technologies

Nanopore sequencing technology and tools for genome assembly: computational analysis of the current state, bottlenecks and future directions

Damla Senol Cali ™, Jeremie S Kim, Saugata Ghose, Can Alkan, Onur Mutlu

Briefings in Bioinformatics, bby017, https://doi.org/10.1093/bib/bby017

Published: 02 April 2018 Article history ▼



Oxford Nanopore MinION

Memory → bottleneck

Memory is Critical for Energy

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 Boroumand¹ Rachata Ausavarungnirun¹ Aki Kuusela³ Allan Knies³

Saugata Ghose¹ Youngsok Kim²

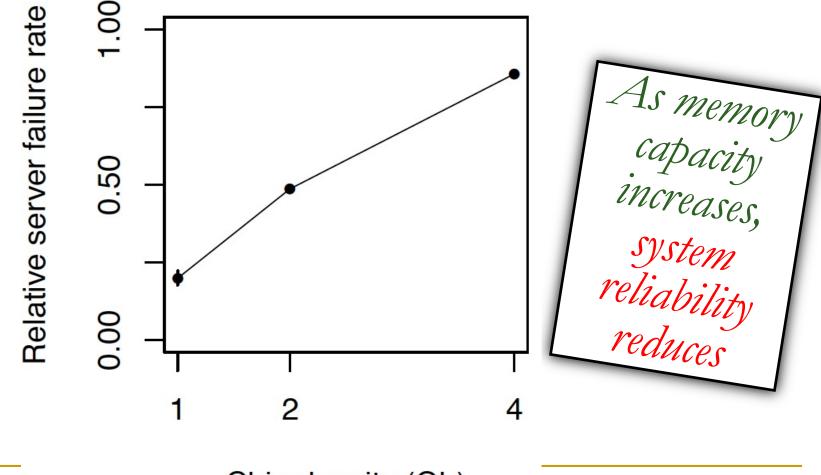
Eric Shiu³ Rahul Thakur³ Daehyun Kim^{4,3}

Parthasarathy Ranganathan³ Onur Mutlu^{5,1}



Memory is Critical for Reliability

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



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 45th Annual IEEE/IFIP International Conference on 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.

107

Modern Systems are Bottlenecked by Memory

An "Early" Overview Paper...

Onur Mutlu,
 "Memory Scaling: A Systems Architecture Perspective"
 Proceedings of the 5th International Memory
 Workshop (IMW), Monterey, CA, May 2013. Slides
 (pptx) (pdf)
 EETimes Reprint

Memory Scaling: A Systems Architecture Perspective

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

Four Key Issues in Future Platforms

Fundamentally Secure/Reliable/Safe Architectures

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

Fundamentally Low-Latency and Predictable Architectures

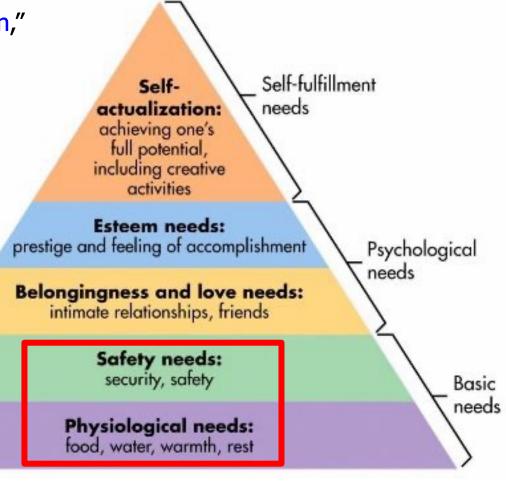
Architectures for AI/ML, Genomics, Medicine, Health

Maslow's (Human) Hierarchy of Needs

Maslow, "A Theory of Human Motivation," Psychological Review, 1943.

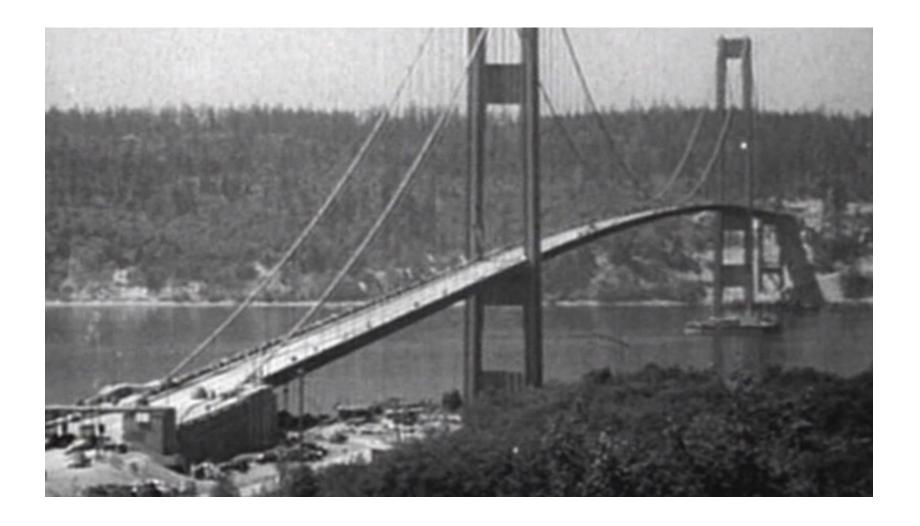
Maslow, "Motivation and Personality," Book, 1954-1970.



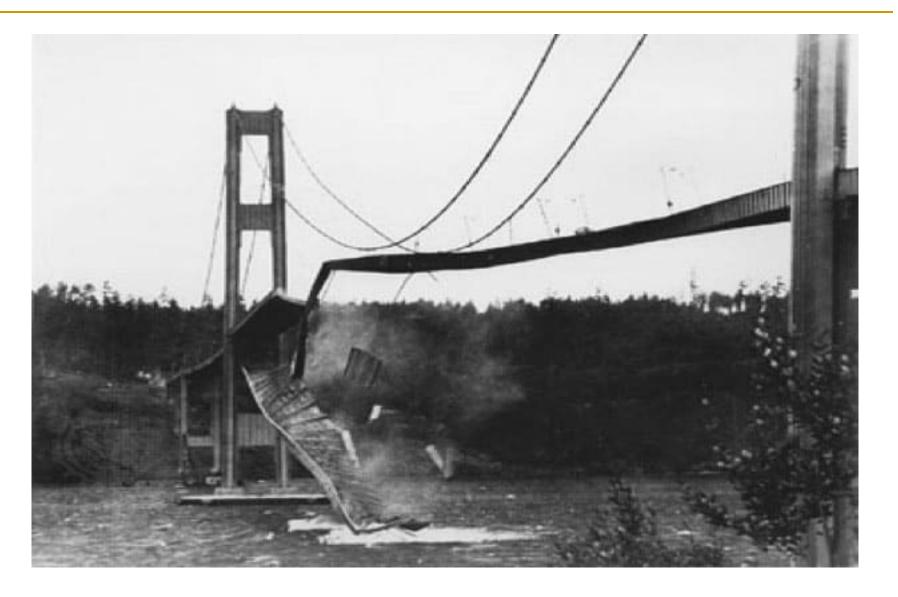


We need to start with reliability, security, safety...

How Reliable/Secure/Safe is This Bridge?



Collapse of the "Galloping Gertie"



How Secure Are These People?

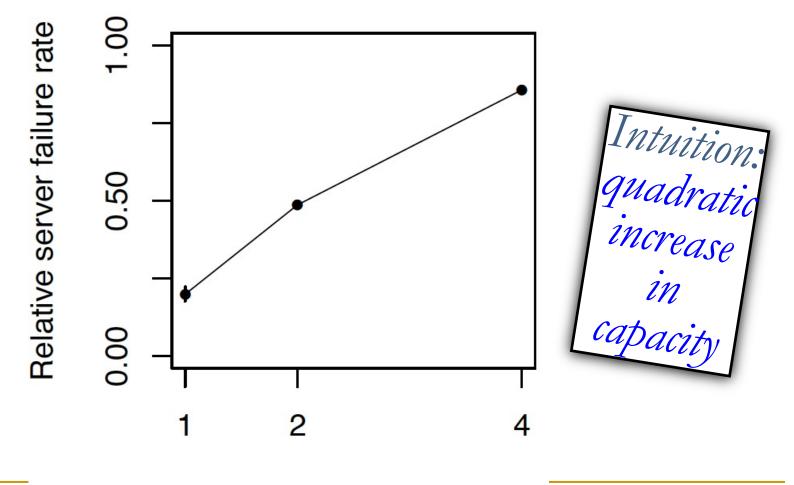


Security is about preventing unforeseen consequences

We do not seem to have design principles for (guaranteeing) reliability and security

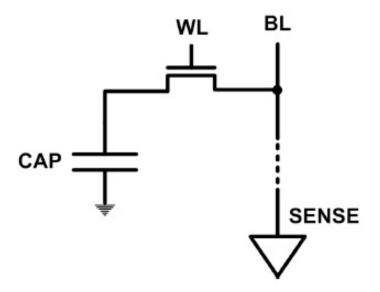
As Memory Scales, It Becomes Unreliable

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



The DRAM Scaling Problem

- DRAM stores charge in a capacitor (charge-based memory)
 - Capacitor must be large enough for reliable sensing
 - Access transistor must be large enough for long data retention time



As DRAM cell becomes smaller, it becomes more vulnerable

Infrastructures to Understand Such Issues



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

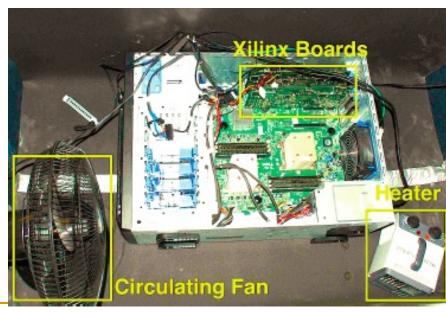
Adaptive-Latency DRAM: Optimizing DRAM
Timing for the Common-Case (Lee et al.,
HPCA 2015)

AVATAR: A Variable-Retention-Time (VRT)

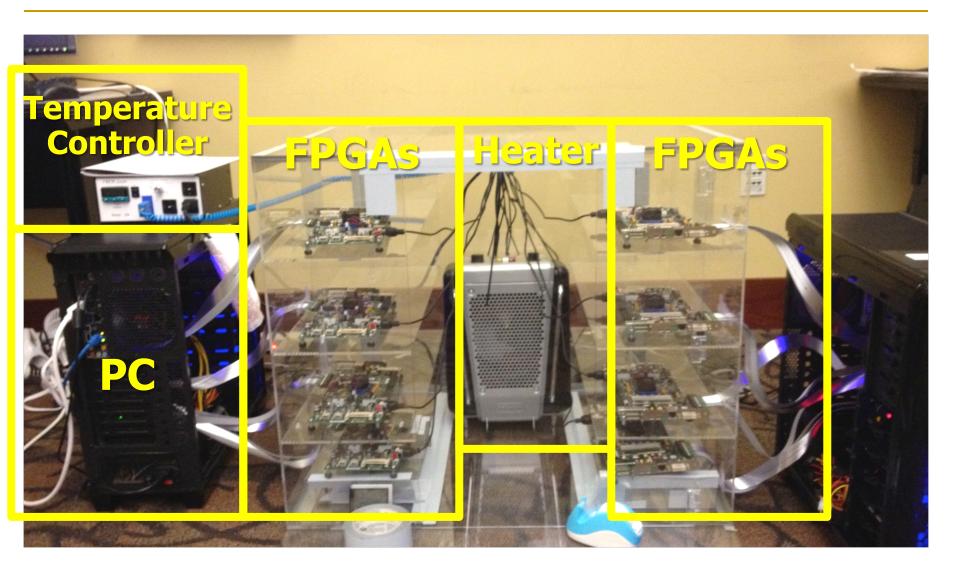
Aware Refresh for DRAM Systems (Qureshi et al., DSN 2015)

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)



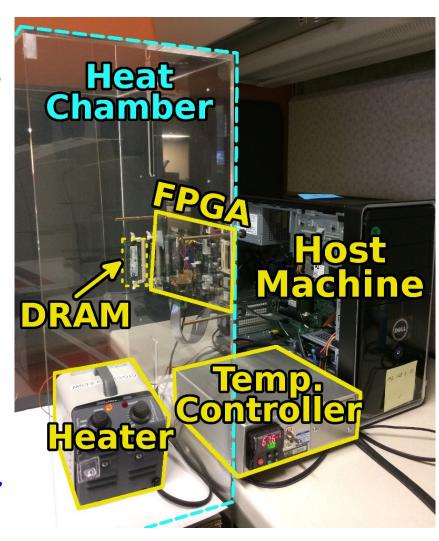
Infrastructures to Understand Such Issues



SoftMC: Open Source DRAM Infrastructure

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

- Flexible
- Easy to Use (C++ API)
- Open-source github.com/CMU-SAFARI/SoftMC



SoftMC

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

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

```
 Hasan Hassan Nandita Vijaykumar Samira Khan Saugata Ghose Kevin Chang Gennady Pekhimenko Donghyuk Lee Gennady Pekhimenko Donghyuk Lee Onur Mutlu Nandita Vijaykumar Samira Khan Saugata Ghose Kevin Chang Gennady Pekhimenko Donghyuk Lee Onur Mutlu Nandita Vijaykumar Samira Khan Saugata Ghose Kevin Chang Gennady Pekhimenko Nandita Vijaykumar Samira Khan Nandita Vijaykumar N
```

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

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

One can predictably induce errors in most DRAM memory chips

DRAM RowHammer

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



Forget Software—Now Hackers Are Exploiting Physics

BUSINESS CULTURE DESIGN GEAR SCIENCE

SHARE

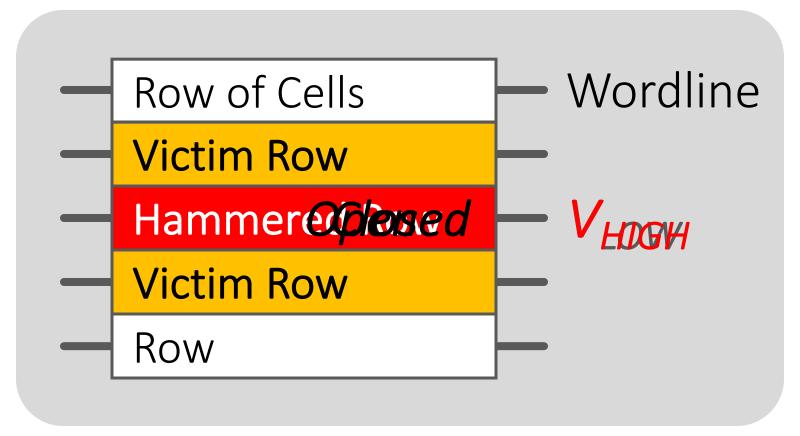




ANDY GREENBERG SECURITY 08.31.16 7:00 AM

FORGET SOFTWARE—NOW HACKERS ARE EXPLOITING PHYSICS

Modern DRAM is Prone to Disturbance Errors



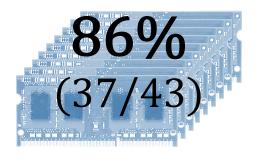
Repeatedly reading a row enough times (before memory gets refreshed) induces disturbance errors in adjacent rows in most real DRAM chips you can buy today

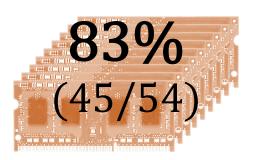
Most DRAM Modules Are Vulnerable

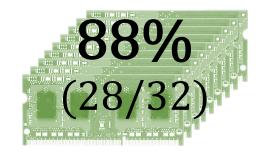
A company

B company

C company







Up to **1.0×10**⁷

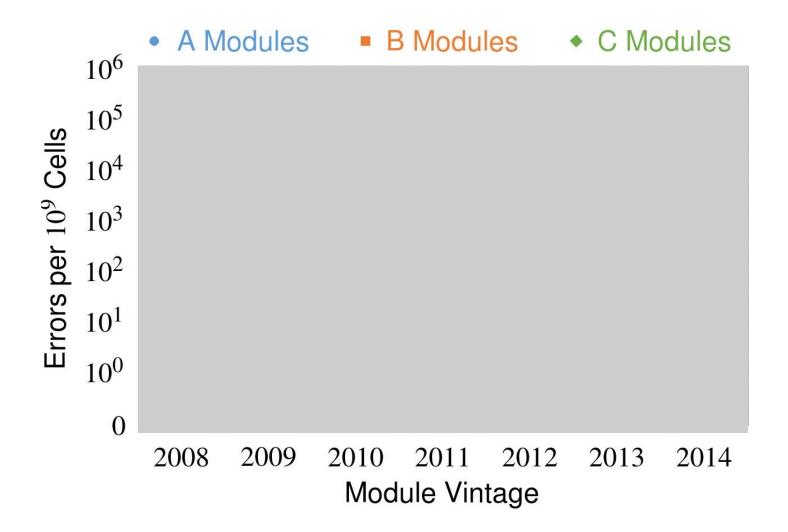
errors

Up to **2.7×10**⁶

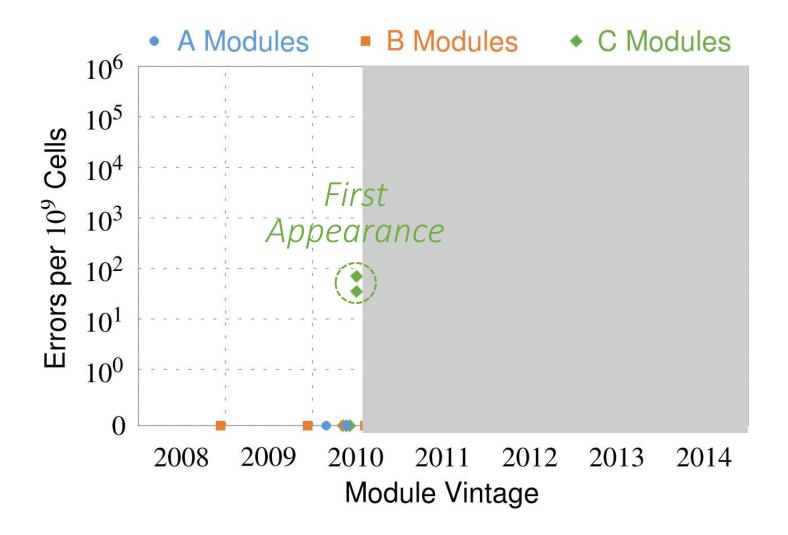
errors

Up to 3.3×10^5 errors

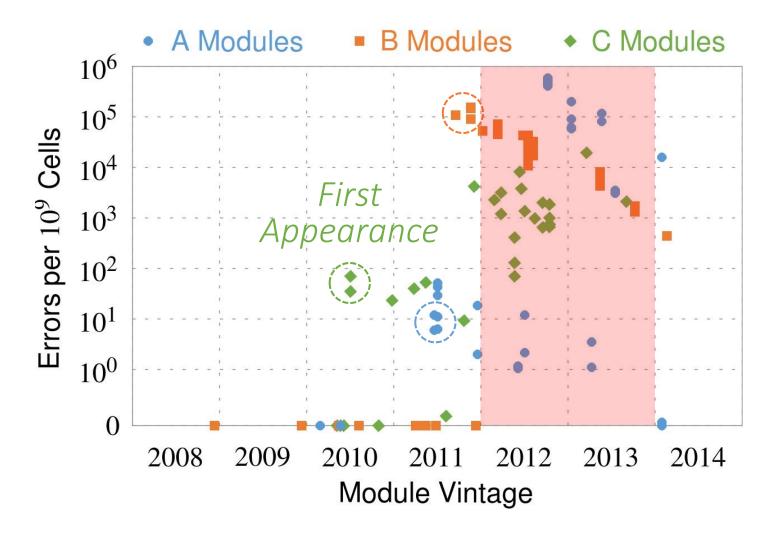
Recent DRAM Is More Vulnerable



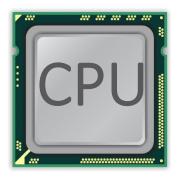
Recent DRAM Is More Vulnerable

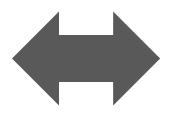


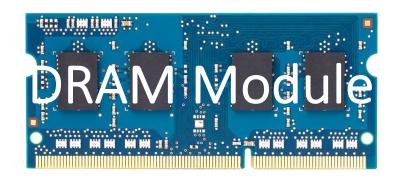
Recent DRAM Is More Vulnerable



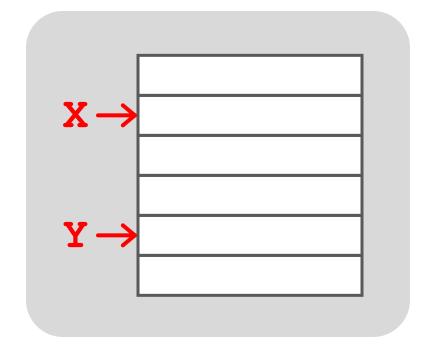
All modules from 2012-2013 are vulnerable

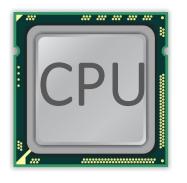


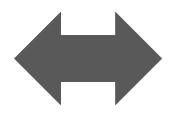


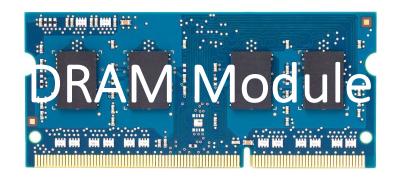


```
loop:
  mov (X), %eax
  mov (Y), %ebx
  clflush (X)
  clflush (Y)
  mfence
  jmp loop
```

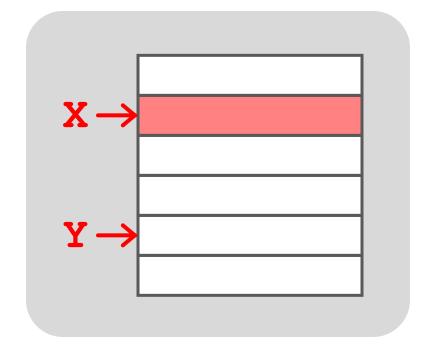


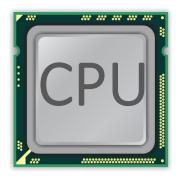


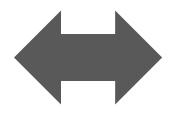




```
loop:
  mov (X), %eax
  mov (Y), %ebx
  clflush (X)
  clflush (Y)
  mfence
  jmp loop
```

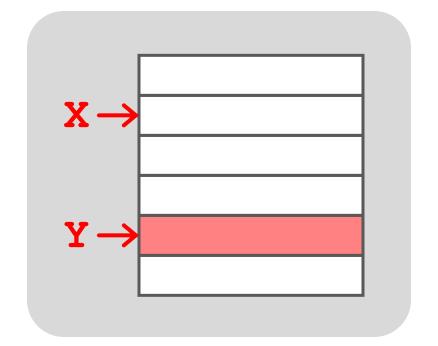


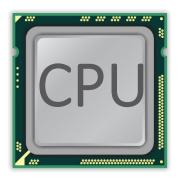


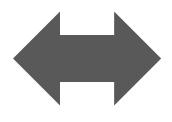


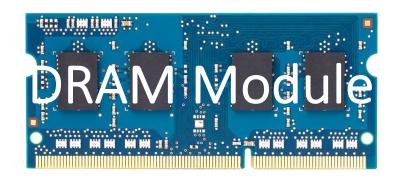


```
loop:
  mov (X), %eax
  mov (Y), %ebx
  clflush (X)
  clflush (Y)
  mfence
  jmp loop
```

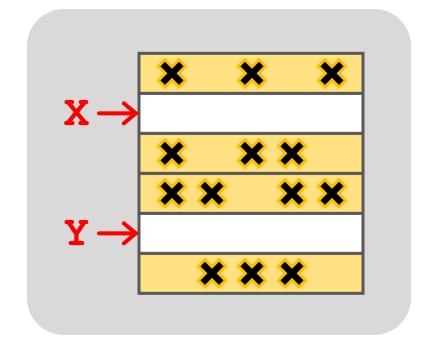








```
loop:
  mov (X), %eax
  mov (Y), %ebx
  clflush (X)
  clflush (Y)
  mfence
  jmp loop
```



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

One Can Take Over an Otherwise-Secure System

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

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

Project Zero

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

News and updates from the Project Zero team at Google

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

Monday, March 9, 2015

Exploiting the DRAM rowhammer bug to gain kernel privileges

RowHammer Security Attack Example

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

Security Implications



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."

www.iaik.tugraz.at

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

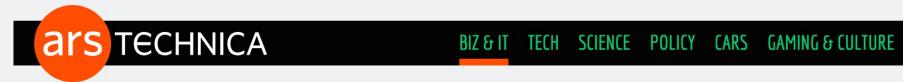
More Security Implications (II)

"Can gain control of a smart phone deterministically" Hammer And Root Millions of Androids

Drammer: Deterministic Rowhammer Attacks on Mobile Platforms, CCS'16¹³⁹

More Security Implications (III)

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



"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)



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 Herbert Bos

VII Amsterdam

Elias Athanasopoulos University of Cyprus

> Kaveh Razavi VU Amsterdam

Cristiano Giuffrida VU Amsterdam

More Security Implications (V)

Rowhammer over RDMA (II)



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

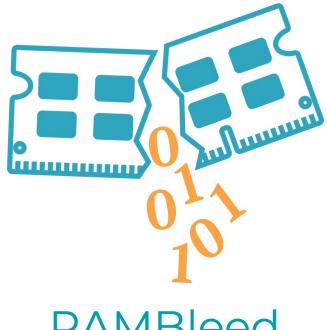
Clémentine Maurice Univ Rennes, CNRS, IRISA

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

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 Dumitraș

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

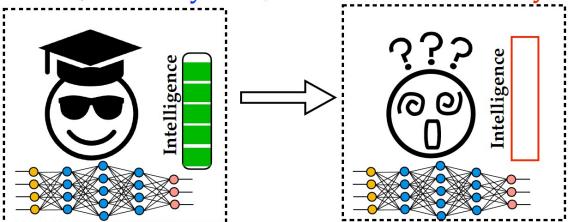
Fan Yao
University of Central Florida
fan.yao@ucf.edu

Adnan Siraj Rakin Deliang Fan Arizona State University asrakin@asu.edu dfan@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)



More Security Implications (IX)

Rowhammer on MLC NAND Flash (based on [Cai+, HPCA 2017])



Security

Rowhammer RAM attack adapted to hit flash storage

Project Zero's two-year-old dog learns a new trick

By Richard Chirgwin 17 Aug 2017 at 04:27

17 🖵

SHARE ▼

From random block corruption to privilege escalation: A filesystem attack vector for rowhammer-like attacks

Anil Kurmus

Nikolas Ioannou

Matthias Neugschwandtner Thomas Parnell

Nikolaos Papandreou

IBM Research – Zurich

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

Solution Direction: Principled Designs

Design fundamentally secure computing architectures

Predict and prevent such safety issues

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

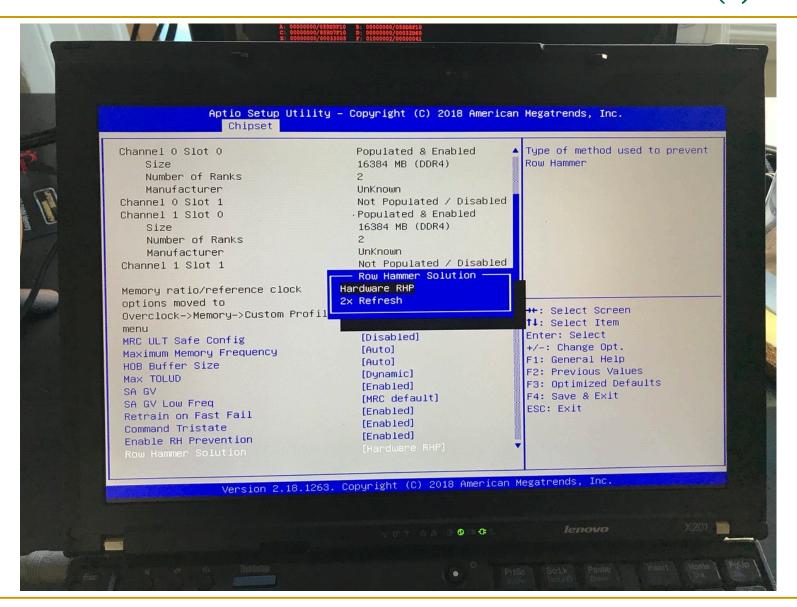
Advantages of PARA

- 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

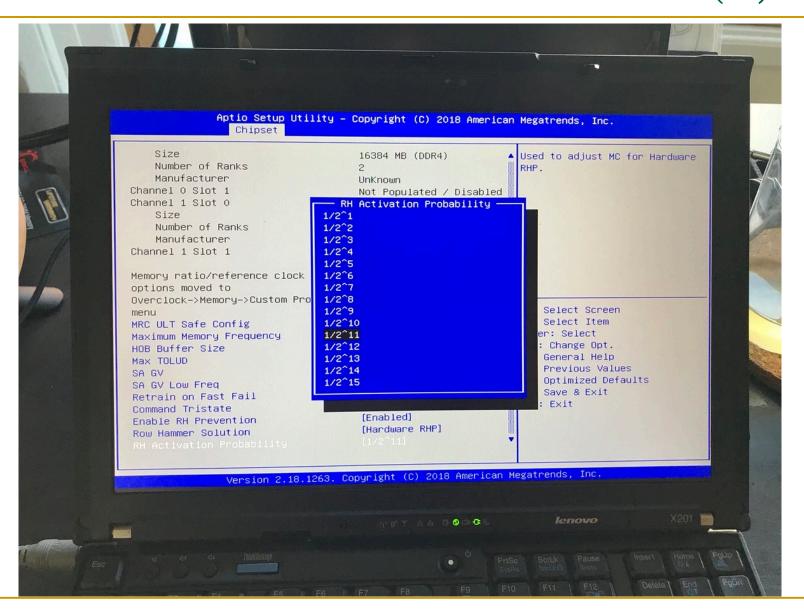
Requirements for PARA

- 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.
 - Kim et al., "Solar-DRAM: Reducing DRAM Access Latency by Exploiting the Variation in Local Bitlines," ICCD 2018.
- If implemented in memory controller (done today)
 - Better coordination between memory controller and DRAM
 - Memory controller should know which rows are physically adjacent

Probabilistic Activation in Real Life (I)



Probabilistic Activation in Real Life (II)



First RowHammer Analysis

 Yoongu Kim, Ross Daly, Jeremie Kim, Chris Fallin, Ji Hye Lee, Donghyuk Lee, Chris Wilkerson, Konrad Lai, and Onur Mutlu,

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

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).

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 155

Retrospective on RowHammer & Future

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.

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

[Slides (pptx) (pdf)]

A More Recent RowHammer Retrospective

Onur Mutlu and Jeremie Kim,

"RowHammer: A Retrospective"

<u>IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems</u> (**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

SAFARI 157

RowHammer in 2020

RowHammer in 2020 (I)

 Jeremie S. Kim, Minesh Patel, A. Giray Yaglikci, Hasan Hassan, Roknoddin Azizi, Lois Orosa, and Onur Mutlu,
 "Revisiting RowHammer: An Experimental Analysis of Modern Devices and Mitigation Techniques"

Proceedings of the <u>47th International Symposium on Computer</u> <u>Architecture</u> (**ISCA**), Valencia, Spain, June 2020.

[Slides (pptx) (pdf)]

[Lightning Talk Slides (pptx) (pdf)]

[Talk Video (20 minutes)]

[Lightning Talk Video (3 minutes)]

Revisiting RowHammer: An Experimental Analysis of Modern DRAM Devices and Mitigation Techniques

```
Jeremie S. Kim^{\S \dagger} Minesh Patel^{\S} A. Giray Yağlıkçı^{\S} Hasan Hassan^{\S} Roknoddin Azizi^{\S} Lois Orosa^{\S} Onur Mutlu^{\S \dagger} ^{\S} ETH Zürich ^{\dagger} Carnegie Mellon University
```

Key Takeaways from 1580 Chips

 Newer DRAM chips are more vulnerable to RowHammer

There are chips today whose weakest cells fail after only
 4800 hammers

• Chips of newer DRAM technology nodes can exhibit RowHammer bit flips 1) in **more rows** and 2) **farther away** from the victim row.

Existing mitigation mechanisms are NOT effective

RowHammer in 2020 (II)

 Pietro Frigo, Emanuele Vannacci, Hasan Hassan, Victor van der Veen, Onur Mutlu, Cristiano Giuffrida, Herbert Bos, and Kaveh Razavi,

"TRRespass: Exploiting the Many Sides of Target Row Refresh"

Proceedings of the <u>41st IEEE Symposium on Security and Privacy</u> (**S&P**), San Francisco, CA, USA, May 2020.

[Slides (pptx) (pdf)]

[Lecture Slides (pptx) (pdf)]

[Talk Video (17 minutes)]

[Lecture Video (59 minutes)]

[Source Code]

[Web Article]

Best paper award.

Pwnie Award 2020 for Most Innovative Research. Pwnie Awards 2020

TRRespass: Exploiting the Many Sides of Target Row Refresh

Pietro Frigo*† Emanuele Vannacci*† Hasan Hassan§ Victor van der Veen¶ Onur Mutlu§ Cristiano Giuffrida* Herbert Bos* Kaveh Razavi*

*Vrije Universiteit Amsterdam

§ETH Zürich

¶Oualcomm Technologies Inc.

RowHammer is still an open problem

Security by obscurity is likely not a good solution

RowHammer in 2020 (III)

Lucian Cojocar, Jeremie Kim, Minesh Patel, Lillian Tsai, Stefan Saroiu,
 Alec Wolman, and Onur Mutlu,

"Are We Susceptible to Rowhammer? An End-to-End Methodology for Cloud Providers"

Proceedings of the <u>41st IEEE Symposium on Security and</u> <u>Privacy</u> (**S&P**), San Francisco, CA, USA, May 2020.

[Slides (pptx) (pdf)]

[Talk Video (17 minutes)]

Are We Susceptible to Rowhammer? An End-to-End Methodology for Cloud Providers

Lucian Cojocar, Jeremie Kim^{§†}, Minesh Patel[§], Lillian Tsai[‡], Stefan Saroiu, Alec Wolman, and Onur Mutlu^{§†} Microsoft Research, [§]ETH Zürich, [†]CMU, [‡]MIT

SAFARI 163

RowHammer in 2020 (IV)

MICRO 2020 Submit Work → Program → Atte

Session 1A: Security & Privacy I 5:00 PM CEST - 5:15 PM CEST Graphene: Strong yet Lightweight Row Hammer Protection Yeonhong Park, Woosuk Kwon, Eojin Lee, Tae Jun Ham, Jung Ho Ahn, Jae W. Lee (Seoul National University) 5:15 PM CEST - 5:30 PM CEST Persist Level Parallelism: Streamlining Integrity Tree Updates for Secure Persistent Memory Alexander Freij, Shougang Yuan, Huiyang Zhou (NC State University); Yan Solihin (University of Central Florida) 5:30 PM CEST - 5:45 PM CEST PThammer: Cross-User-Kernel-Boundary **Rowhammer through Implicit Accesses** Zhi Zhang (University of New South Wales and Data61, CSIRO, Australia); Yueqiang Cheng (Baidu Security); Dongxi Liu, Surya Nepal (Data61, CSIRO, Australia); Zhi Wang (Florida State University); Yuval Yarom (University of Adelaide and Data61, CSIRO, Australia)

RowHammer in 2020 (V)

Session #5: Rowhammer

Room 2

Session chair: Michael Franz (UC Irvine)

RAMBleed: Reading Bits in Memory Without Accessing Them

Andrew Kwong (University of Michigan), Daniel Genkin (University of Michigan), Daniel Gruss Data61)

Are We Susceptible to Rowhammer? An End-to-End Methodology for Cloud Providers Lucian Cojocar (Microsoft Research), Jeremie Kim (ETH Zurich, CMU), Minesh Patel (ETH Zu (Microsoft Research), Onur Mutlu (ETH Zurich, CMU)

Leveraging EM Side-Channel Information to Detect Rowhammer Attacks

Zhenkai Zhang (Texas Tech University), Zihao Zhan (Vanderbilt University), Daniel Balasubrar Peter Volgyesi (Vanderbilt University), Xenofon Koutsoukos (Vanderbilt University)

TRRespass: Exploiting the Many Sides of Target Row Refresh

Pietro Frigo (Vrije Universiteit Amsterdam, The Netherlands), Emanuele Vannacci (Vrije Universiteit Amsterdam, The Netherlands), Cristiano Giuffrida (Vrije Universiteit Amsterdam, The Netherlands)

RowHammer in 2020 (VI)

29TH USENIX SECURITY SYMPOSIUM

ATTEND

PROGRAM

PARTICIPATE

SPONSORS

ABOUT

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

Fan Yao, University of Central Florida: Adnan Sirai Pakin and Deliang Fan Arizona State University

Fan Yao, *University of Central Florida*; Adnan Siraj Rakin and Deliang Fan, *Arizona State University*

AVAILABLE MEDIA 🗋 ず 🕞

Show details >

BlockHammer Solution in 2021

 A. Giray Yaglikci, Minesh Patel, Jeremie S. Kim, Roknoddin Azizi, Ataberk Olgun, Lois Orosa, Hasan Hassan, Jisung Park, Konstantinos Kanellopoulos, Taha Shahroodi, Saugata Ghose, and Onur Mutlu,

"BlockHammer: Preventing RowHammer at Low Cost by Blacklisting Rapidly-Accessed DRAM Rows"

Proceedings of the <u>27th International Symposium on High-Performance</u> <u>Computer Architecture</u> (**HPCA**), Virtual, February-March 2021.

[Slides (pptx) (pdf)]

[Short Talk Slides (pptx) (pdf)]

[Talk Video (22 minutes)]

[Short Talk Video (7 minutes)]

BlockHammer: Preventing RowHammer at Low Cost by Blacklisting Rapidly-Accessed DRAM Rows

A. Giray Yağlıkçı¹ Minesh Patel¹ Jeremie S. Kim¹ Roknoddin Azizi¹ Ataberk Olgun¹ Lois Orosa¹ Hasan Hassan¹ Jisung Park¹ Konstantinos Kanellopoulos¹ Taha Shahroodi¹ Saugata Ghose² Onur Mutlu¹

¹ETH Zürich ²University of Illinois at Urbana–Champaign

SAFARI 167

Two Key RowHammer Papers at MICRO 2021

Lois Orosa, Abdullah Giray Yaglikci, Haocong Luo, Ataberk Olgun, Jisung Park, Hasan Hassan,
 Minesh Patel, Jeremie S. Kim, and Onur Mutlu,

"A Deeper Look into RowHammer's Sensitivities: Experimental Analysis of Real DRAM Chips and Implications on Future Attacks and Defenses"

Proceedings of the <u>54th International Symposium on Microarchitecture</u> (**MICRO**), Virtual, October 2021.

[Slides (pptx) (pdf)]

[Short Talk Slides (pptx) (pdf)]

[Lightning Talk Slides (pptx) (pdf)]

[Talk Video (21 minutes)]

[Lightning Talk Video (1.5 minutes)]

[arXiv version]

A Deeper Look into RowHammer's Sensitivities: Experimental Analysis of Real DRAM Chips and Implications on Future Attacks and Defenses

Lois Orosa* ETH Zürich A. Giray Yağlıkçı* ETH Zürich Haocong Luo ETH Zürich Ataberk Olgun ETH Zürich, TOBB ETÜ Jisung Park ETH Zürich

Hasan Hassan ETH Zürich Minesh Patel ETH Zürich

Jeremie S. Kim ETH Zürich Onur Mutlu ETH Zürich

Two Key RowHammer Papers at MICRO 2021

Hasan Hassan, Yahya Can Tugrul, Jeremie S. Kim, Victor van der Veen, Kaveh Razavi, and Onur Mutlu,

"Uncovering In-DRAM RowHammer Protection Mechanisms: A New Methodology, Custom RowHammer Patterns, and Implications"

Proceedings of the <u>54th International Symposium on Microarchitecture</u> (**MICRO**), Virtual, October 2021.

[Slides (pptx) (pdf)]

[Short Talk Slides (pptx) (pdf)]

[Lightning Talk Slides (pptx) (pdf)]

[Talk Video (25 minutes)]

[<u>Lightning Talk Video</u> (100 seconds)]

arXiv version

Uncovering In-DRAM RowHammer Protection Mechanisms: A New Methodology, Custom RowHammer Patterns, and Implications

Yahya Can Tuğrul^{†‡} Jeremie S. Kim[†] Hasan Hassan[†] Victor van der Veen $^{\sigma}$ Kaveh Razavi[†] Onur Mutlu[†]

 $^\ddagger TOBB\ University\ of\ Economics\ \&\ Technology$ $^\sigma Qualcomm\ Technologies\ Inc.$ †ETH Zürich

More to Come...

Detailed Lectures on RowHammer

- Computer Architecture, Fall 2020, Lecture 4b
 - RowHammer (ETH Zürich, Fall 2020)
 - https://www.youtube.com/watch?v=KDy632z23UE&list=PL5Q2soXY2Zi9xidyIgBxUz 7xRPS-wisBN&index=8
- Computer Architecture, Fall 2020, Lecture 5a
 - RowHammer in 2020: TRRespass (ETH Zürich, Fall 2020)
 - https://www.youtube.com/watch?v=pwRw7QqK_qA&list=PL5Q2soXY2Zi9xidyIgBxUz7xRPS-wisBN&index=9
- Computer Architecture, Fall 2020, Lecture 5b
 - RowHammer in 2020: Revisiting RowHammer (ETH Zürich, Fall 2020)
 - https://www.youtube.com/watch?v=gR7XR-Eepcg&list=PL5Q2soXY2Zi9xidyIgBxUz7xRPS-wisBN&index=10
- Computer Architecture, Fall 2020, Lecture 5c
 - Secure and Reliable Memory (ETH Zürich, Fall 2020)
 - https://www.youtube.com/watch?v=HvswnsfG3oQ&list=PL5Q2soXY2Zi9xidyIgBxUz 7xRPS-wisBN&index=11

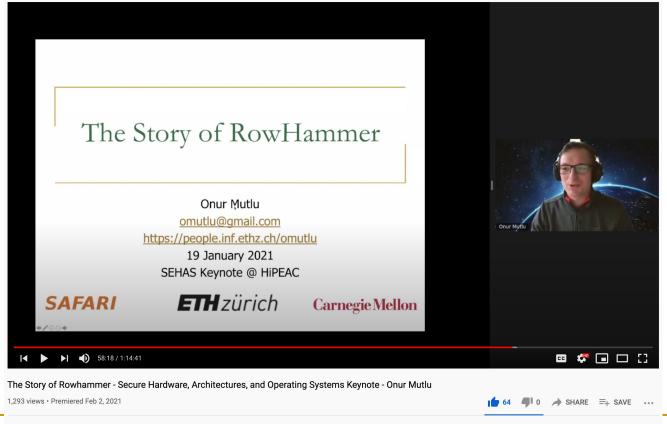
The Story of RowHammer Lecture ...

Onur Mutlu,

"The Story of RowHammer"

Keynote Talk at <u>Secure Hardware, Architectures, and Operating Systems</u>
<u>Workshop</u> (**SeHAS**), held with <u>HiPEAC 2021 Conference</u>, Virtual, 19 January 2021.
[Slides (pptx) (pdf)]

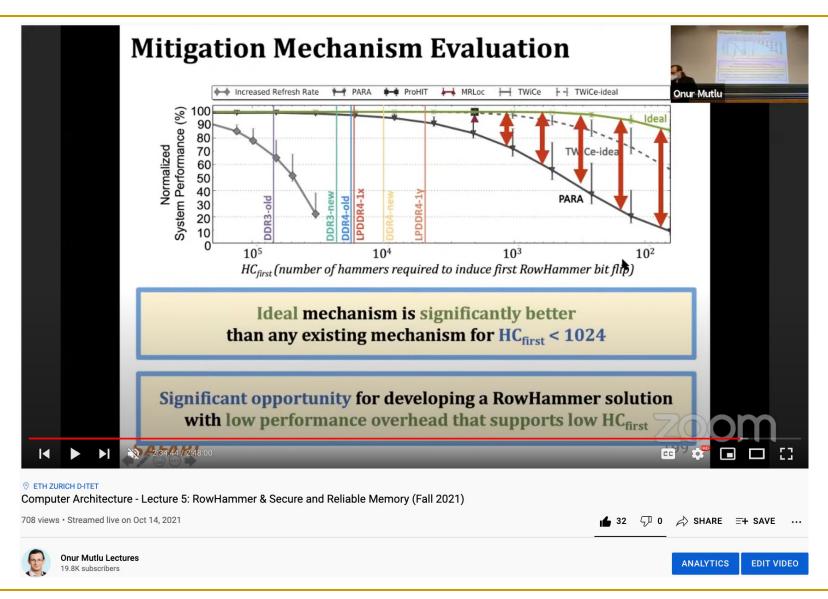
[Talk Video (1 hr 15 minutes, with Q&A)]







Most Recent RowHammer Lecture





Future of Main Memory Reliability

- DRAM is becoming less reliable → more vulnerable
- Due to difficulties in DRAM scaling, other problems may also appear (or they may be going unnoticed)
- Some errors may already be slipping into the field
 - Read disturb errors (Rowhammer)
 - Retention errors
 - Read errors, write errors
 - **...**
- These errors can also pose security vulnerabilities

All Memory Technologies are Vulnerable

- DRAM
- Flash memory
- Emerging Technologies
 - Phase Change Memory
 - STT-MRAM
 - RRAM, memristors
 - **...**

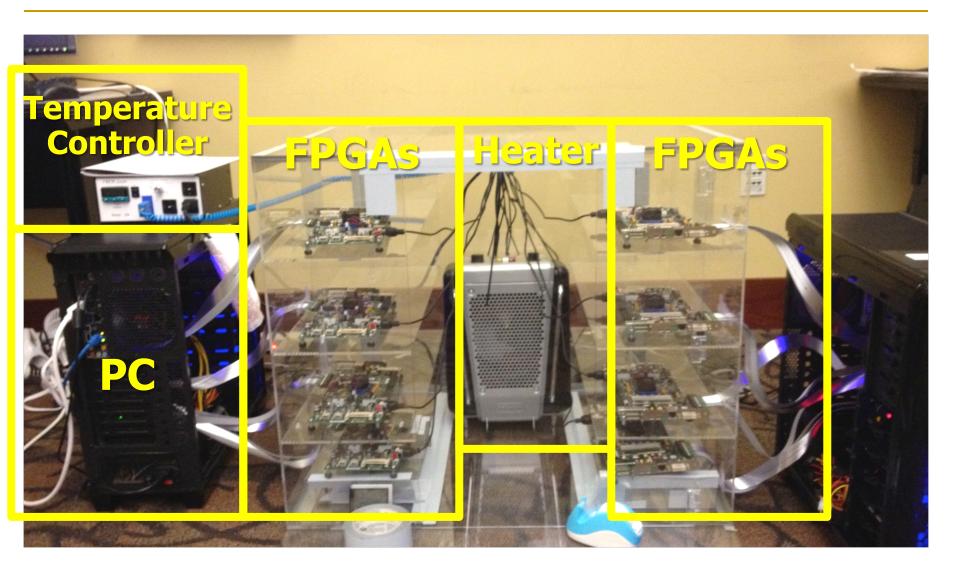
How Do We Keep Memory Secure?

- Understand: Methodologies for failure modeling and discovery
 - Modeling and prediction based on real (device) data

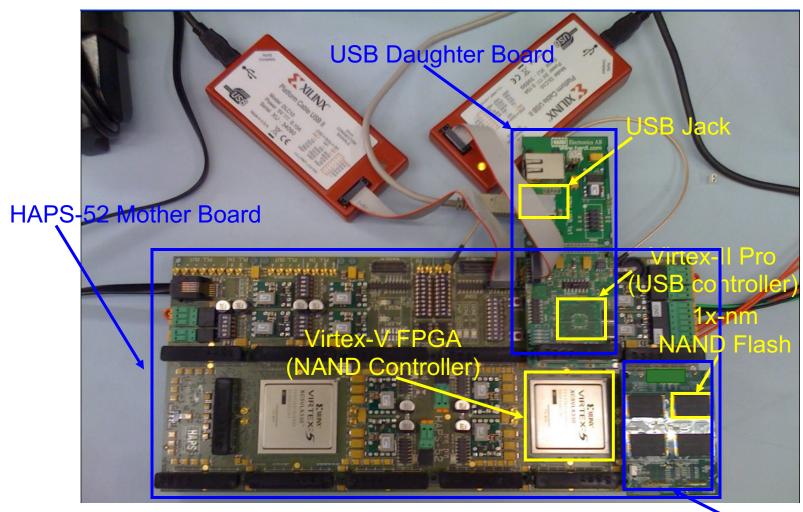
- Architect: Principled co-architecting of system and memory
 - Good partitioning of duties across the stack

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

Understand and Model with Experiments (DRAM)



Understand and Model with Experiments (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

Understanding Flash Memory Reliability



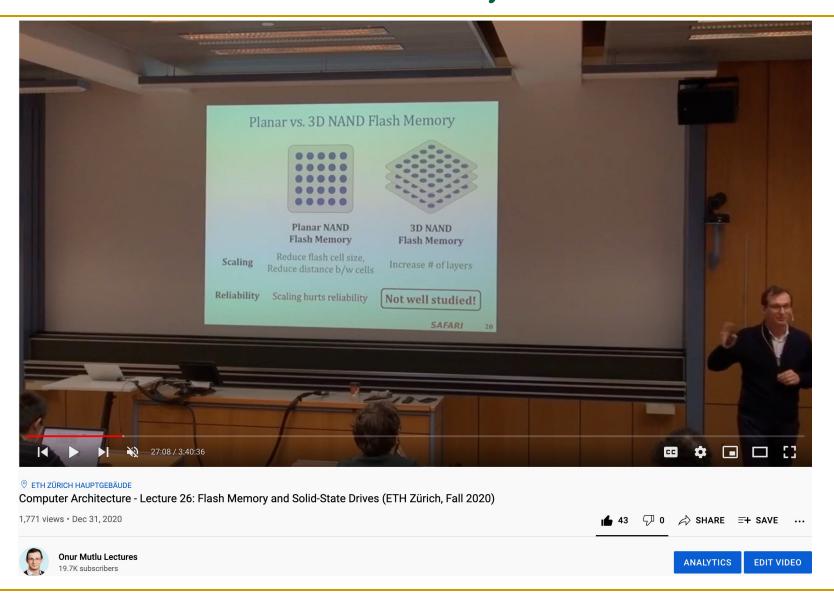
Proceedings of the IEEE, Sept. 2017

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

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

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

Lecture on Flash Memory & SSDs



Special Course on Flash Memory & SSDs



Challenge and Opportunity for Future

Fundamentally Secure, Reliable, Safe Computing Architectures

One Important Takeaway

Main Memory Needs Intelligent Controllers

In-Field Patch-ability (Intelligent Memory) Can Avoid Many Failures

Four Key Issues in Future Platforms

Fundamentally Secure/Reliable/Safe Architectures

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

Fundamentally Low-Latency and Predictable Architectures

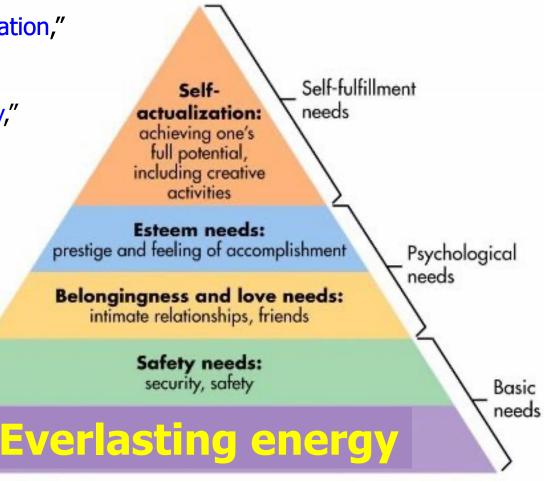
Architectures for AI/ML, Genomics, Medicine, Health

Maslow's (Human) Hierarchy of Needs, Revisited

Maslow, "A Theory of Human Motivation," Psychological Review, 1943.

Maslow, "Motivation and Personality," Book, 1954-1970.





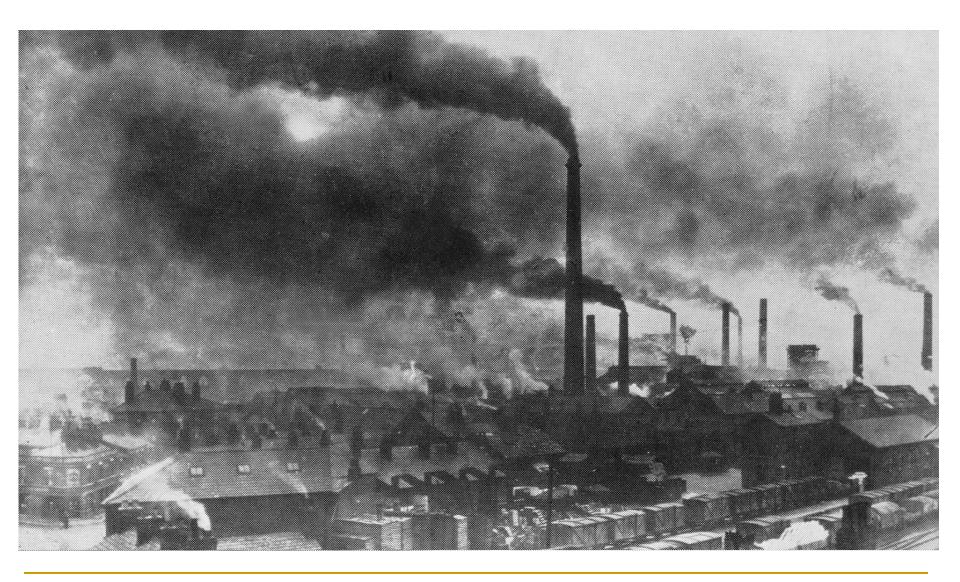
Do We Want This?





188

Or This?



SAFARI

189

Challenge and Opportunity for Future

High Performance, Energy Efficient, Sustainable

The Problem

Data access is the major performance and energy bottleneck

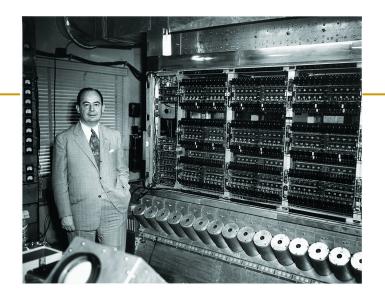
Our current design principles cause great energy waste

(and great performance loss)

Processing of data is performed far away from the data

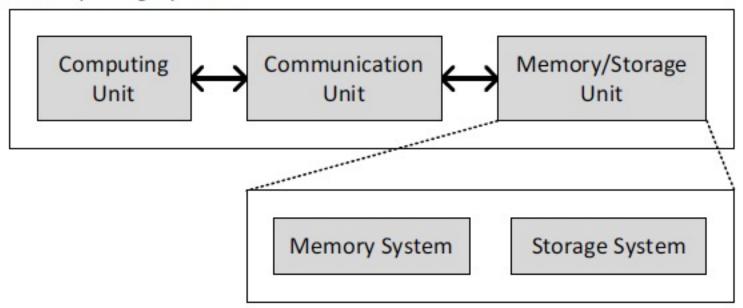
A Computing System

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



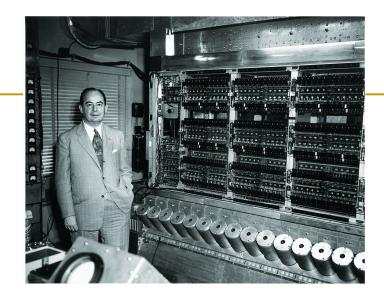
Burks, Goldstein, von Neumann, "Preliminary discussion of the logical design of an electronic computing instrument," 1946.

Computing System



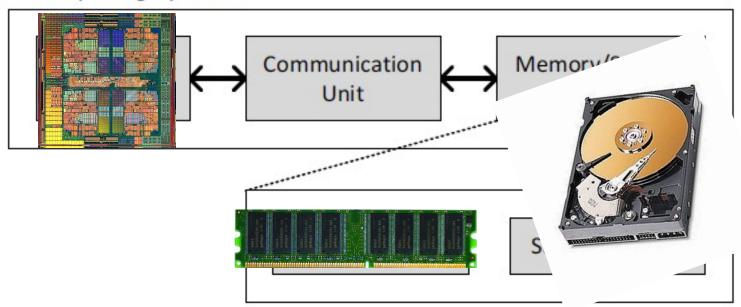
A Computing System

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



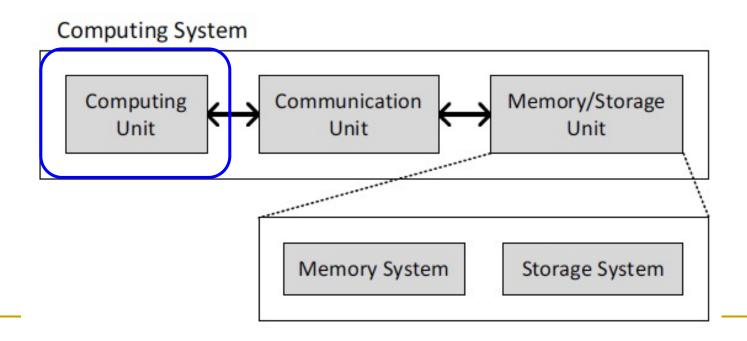
Burks, Goldstein, von Neumann, "Preliminary discussion of the logical design of an electronic computing instrument," 1946.

Computing System



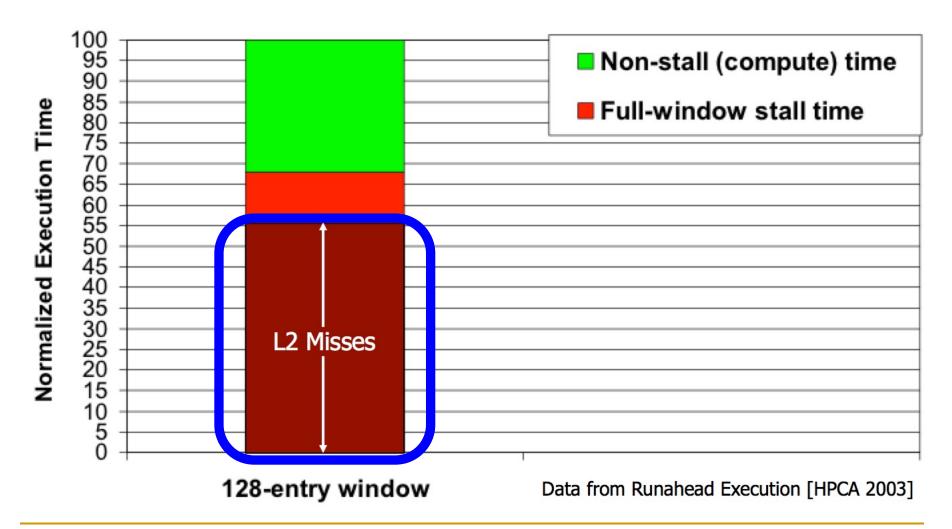
Today's Computing Systems

- Are overwhelmingly processor centric
- All data processed in the processor → 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)



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

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



The Performance Perspective

HPCA Test of Time Award (awarded in 2021).

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 Computer</u>

<u>Architecture</u> (**HPCA**), pages 129-140, Anaheim, CA, February 2003. <u>Slides (pdf)</u>

<u>One of the 15 computer arch. papers of 2003 selected as Top Picks by IEEE Micro.</u>

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

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

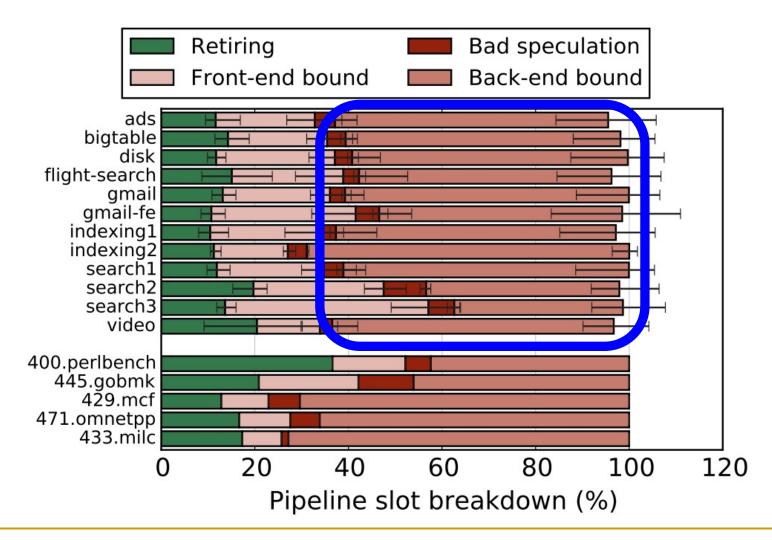
§ECE Department
The University of Texas at Austin
{onur,patt}@ece.utexas.edu

†Microprocessor Research Intel Labs jared.w.stark@intel.com

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

The Performance Perspective (2015)

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



The Performance Perspective (2015)

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

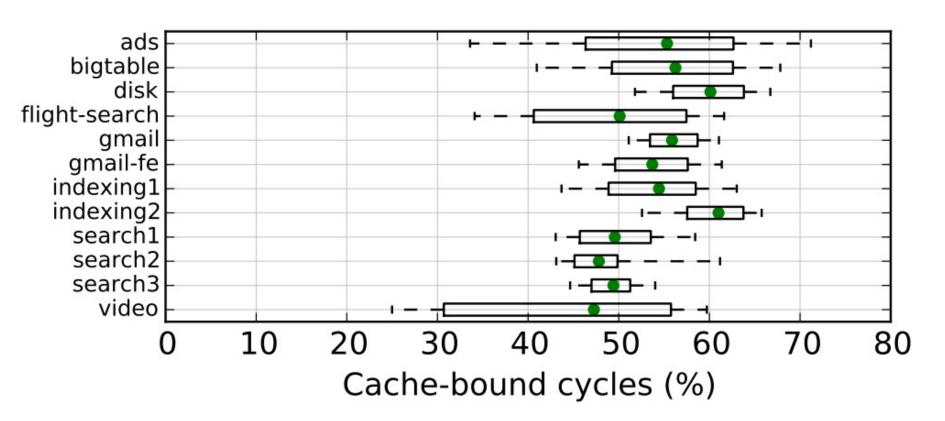
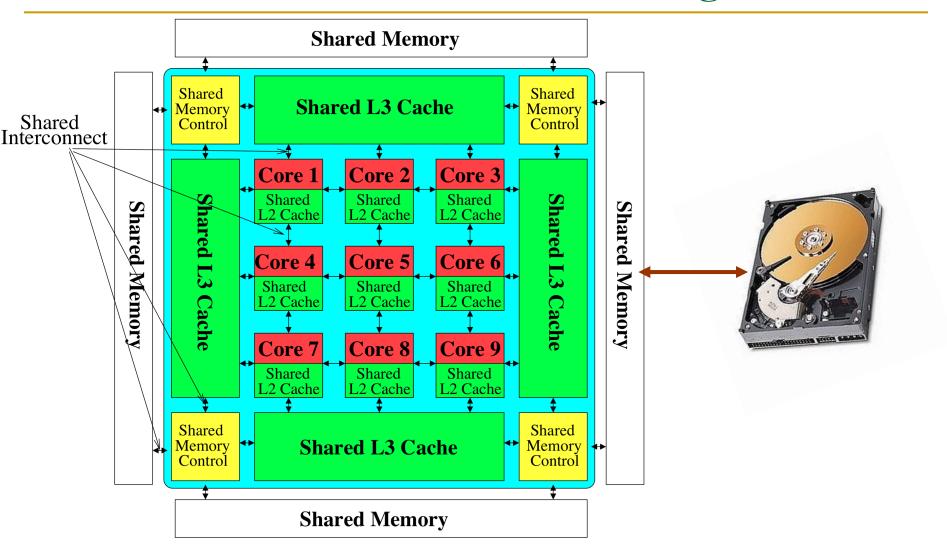


Figure 11: Half of cycles are spent stalled on caches.

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
 - → Energy inefficient
 - → Low performance
 - → Complex
- Overly complex and bloated processor (and accelerators)
 - To tolerate data access from memory
 - Complex hierarchies and mechanisms
 - → Energy inefficient
 - → Low performance
 - → Complex

Perils of Processor-Centric Design



Most of the system is dedicated to storing and moving data

Three Key Systems Trends

1. Data access is a major bottleneck

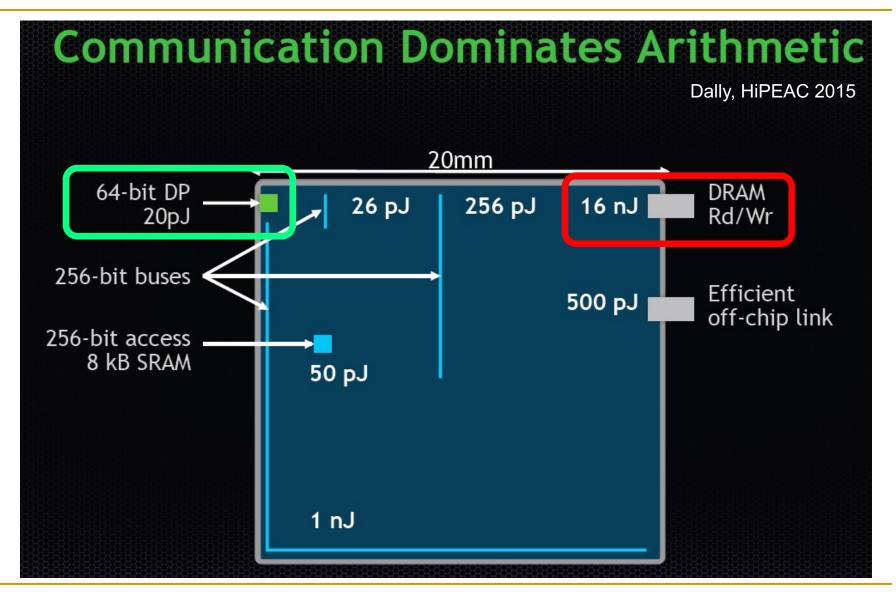
Applications are increasingly data hungry

2. Energy consumption is a key limiter

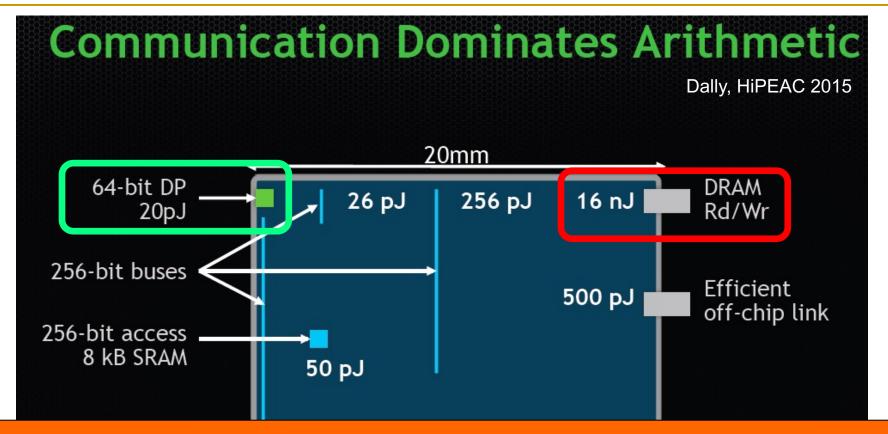
3. Data movement energy dominates compute

Especially true for off-chip to on-chip movement

Data Movement vs. Computation Energy



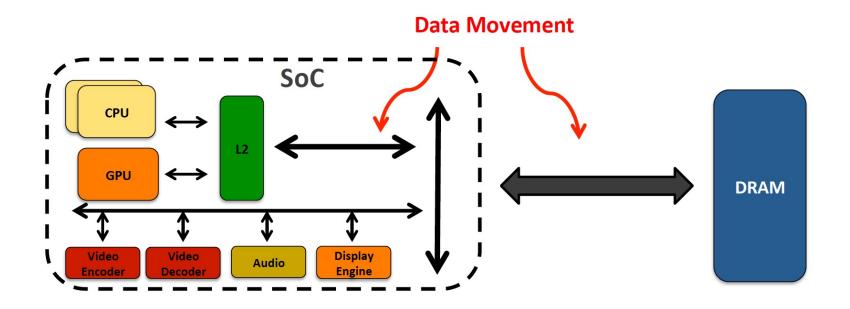
Data Movement vs. Computation Energy



A memory access consumes ~100-1000X 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

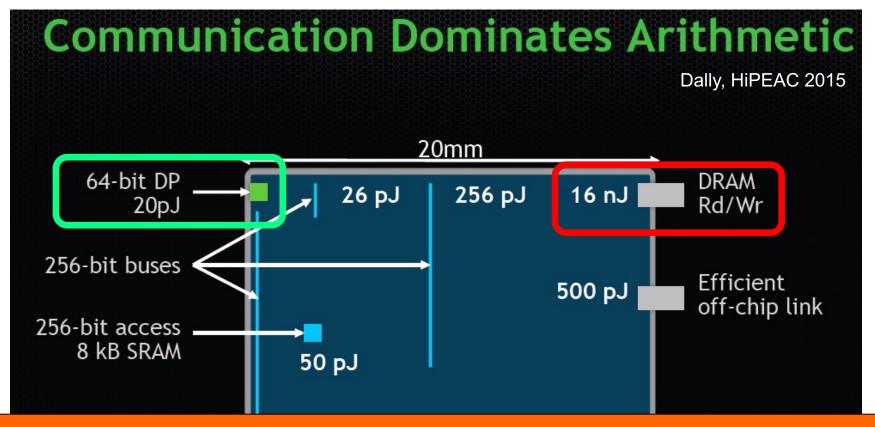
Amirali Boroumand¹ Rachata Ausavarungnirun¹ Aki Kuusela³ Allan Knies³

Saugata Ghose¹ Youngsok Kim²

Eric Shiu³ Rahul Thakur³ Daehyun Kim^{4,3}

Parthasarathy Ranganathan³ Onur Mutlu^{5,1}

We Do Not Want to Move Data!



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

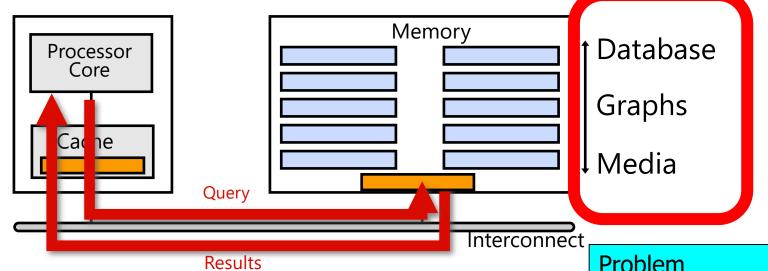
We Need A Paradigm Shift To ...

Enable computation with minimal data movement

Compute where it makes sense (where data resides)

Make computing architectures more data-centric

Goal: Processing Inside Memory



- Many questions ... How do we design the:
 - compute-capable memory & controllers?
 - processor chip and in-memory units?
 - software and hardware interfaces?
 - system software, compilers, languages?
 - algorithms and theoretical foundations?

Problem

Aigorithm

Program/Language

System Software

SW/HW Interface

Micro-architecture

Logic

Electrons

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 Emerging Computing: From Devices to Systems
Looking Beyond Moore and Von Neumann, Springer, to be published in 2021.

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

Abstract

Modern computing systems are overwhelmingly designed to move data to computation. This design choice goes directly against at least three key trends in computing that cause performance, scalability and energy bottlenecks: (1) data access is a key bottleneck as many important applications are increasingly data-intensive, and memory bandwidth and energy do not scale well, (2) energy consumption is a key limiter in almost all computing platforms, especially server and mobile systems, (3) data movement, especially off-chip to on-chip, is very expensive in terms of bandwidth, energy and latency, much more so than computation. These trends are especially severely-felt in the data-intensive server and energy-constrained mobile systems of today.

At the same time, conventional memory technology is facing many technology scaling challenges in terms of reliability, energy, and performance. As a result, memory system architects are open to organizing memory in different ways and making it more intelligent, at the expense of higher cost. The emergence of 3D-stacked memory plus logic, the adoption of error correcting codes inside the latest DRAM chips, proliferation of different main memory standards and chips, specialized for different purposes (e.g., graphics, low-power, high bandwidth, low latency), and the necessity of designing new solutions to serious reliability and security issues, such as the RowHammer phenomenon, are an evidence of this trend.

This chapter discusses recent research that aims to practically enable computation close to data, an approach we call processing-in-memory (PIM). PIM places computation mechanisms in or near where the data is stored (i.e., inside the memory chips, in the logic layer of 3D-stacked memory, or in the memory controllers), so that data movement between the computation units and memory is reduced or eliminated. While the general idea of PIM is not new, we discuss motivating trends in applications as well as memory circuits/technology that greatly exacerbate the need for enabling it in modern computing systems. We examine at least two promising new approaches to designing PIM systems to accelerate important data-intensive applications: (1) processing using memory by exploiting analog operational properties of DRAM chips to perform massively-parallel operations in memory, with low-cost changes, (2) processing near memory by exploiting 3D-stacked memory technology design to provide high memory bandwidth and low memory latency to in-memory logic. In both approaches, we describe and tackle relevant cross-layer research, design, and adoption challenges in devices, architecture, systems, and programming models. Our focus is on the development of in-memory processing designs that can be adopted in real computing platforms at low cost. We conclude by discussing work on solving key challenges to the practical adoption of PIM.

Keywords: memory systems, data movement, main memory, processing-in-memory, near-data processing, computation-in-memory, processing using memory, processing near memory, 3D-stacked memory, non-volatile memory, energy efficiency, high-performance computing, computer architecture, computing paradigm, emerging technologies, memory scaling, technology scaling, dependable systems, robust systems, hardware security, system security, latency, low-latency computing

-	
Con	4 am
Con	цеп

1	Introduction		
2	Major Trends Affecting Main Memory		
3	The Need for Intelligent Memory Controllers		
	to Enhance Memory Scaling		
4	Perils of Processor-Centric Design		
5	Processing-in-Memory (PIM): Technology		
		blers and Two Approaches	12
	5.1	New Technology Enablers: 3D-Stacked	
		Memory and Non-Volatile Memory	12
	5.2	Two Approaches: Processing Using	
		Memory (PUM) vs. Processing Near	
		Memory (PNM)	13
6	Proc	cessing Using Memory (PUM)	14
•	6.1	RowClone	14
	6.2	Ambit	15
	6.3	Gather-Scatter DRAM	17
	6.4	In-DRAM Security Primitives	17
	0.4	III-DRAM Security Frimitives	1/
7	Proc	cessing Near Memory (PNM)	18
	7.1	Tesseract: Coarse-Grained Application-	
		Level PNM Acceleration of Graph Pro-	
		cessing	19
	7.2	Function-Level PNM Acceleration of	
		Mobile Consumer Workloads	20
	7.3	Programmer-Transparent Function-	
		Level PNM Acceleration of GPU	
		Applications	21
	7.4	Instruction-Level PNM Acceleration	
		with PIM-Enabled Instructions (PEI)	21
	7.5	Function-Level PNM Acceleration of	
		Genome Analysis Workloads	22
	7.6	Application-Level PNM Acceleration of	
L		Time Series Analysis	23
8	Enal	bling the Adoption of PIM	24
	8.1	Programming Models and Code Genera-	
		tion for PIM	24
	8.2	PIM Runtime: Scheduling and Data	
		Mapping	25
	8.3	Memory Coherence	27
	8.4	Virtual Memory Support	27
	8.5	Data Structures for PIM	28
	8.6	Benchmarks and Simulation Infrastruc-	
		tures	29
	8.7	Real PIM Hardware Systems and Proto-	
		types	30
	8.8	Security Considerations	30
9	Con	clusion and Future Outlook	31

1. Introduction

Main memory, built using the Dynamic Random Access Memory (DRAM) technology, is a major component in nearly all computing systems, including servers, cloud platforms, mobile/embedded devices, and sensor systems. Across all of these systems, the data working set sizes of modern applications are rapidly growing, while the need for fast analysis of such data is increasing. Thus, main memory is becoming an increasingly significant bottleneck across a wide variety of computing systems and applications [1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16]. Alleviating the main memory bottleneck requires the memory capacity, energy, cost, and performance to all scale in an efficient manner across technology generations. Unfortunately, it has become increasingly difficult in recent years, especially the past decade, to scale all of these dimensions [1, 2, 17, 18, 19, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29, 30, 31, 32, 33, 34, 35, 36, 37, 38, 39, 40, 41, 42, 43, 44, 45, 46, 47, 48, 49], and thus the main memory bottleneck has been worsening.

A major reason for the main memory bottleneck is the high energy and latency cost associated with data movement. In modern computers, to perform any operation on data that resides in main memory, the processor must retrieve the data from main memory. This requires the memory controller to issue commands to a DRAM module across a relatively slow and power-hungry off-chip bus (known as the memory channel). The DRAM module sends the requested data across the memory channel, after which the data is placed in the caches and registers. The CPU can perform computation on the data once the data is in its registers. Data movement from the DRAM to the CPU incurs long latency and consumes a significant amount of energy [7, 50, 51, 52, 53, 54]. These costs are often exacerbated by the fact that much of the data brought into the caches is not reused by the CPU [52, 53, 55, 56], providing little benefit in return for the high latency and energy cost.

The cost of data movement is a fundamental issue with the processor-centric nature of contemporary computer systems. The CPU is considered to be the master in the system, and computation is performed only in the processor (and accelerators). In contrast, data storage and communication units, including the main memory, are treated as unintelligent workers that are incapable of computation. As a result of this processor-centric design paradigm, data moves a lot in the system between the computation units and communication/ storage units so that computation can be done on it. With the increasingly data-centric nature of contemporary and emerging appli-

We Need to Think Differently from the Past Approaches

Processing in Memory: Two Approaches

- 1. Processing using Memory
- 2. Processing near Memory

Two PIM Approaches

5.2. Two Approaches: Processing Using Memory (PUM) vs. Processing Near Memory (PNM)

Many recent works take advantage of the memory technology innovations that we discuss in Section 5.1 to enable and implement PIM. We find that these works generally take one of two approaches, which are categorized in Table 1: (1) processing using memory or (2) processing near memory. We briefly describe each approach here. Sections 6 and 7 will provide example approaches and more detail for both.

Table 1: Summary of enabling technologies for the two approaches to PIM used by recent works. Adapted from [309].

Approach	Enabling Technologies
	SRAM
Processing Using Memory	DRAM
	Phase-change memory (PCM)
	Magnetic RAM (MRAM)
	Resistive RAM (RRAM)/memristors
Processing Near Memory	Logic layers in 3D-stacked memory
	Silicon interposers
	Logic in memory controllers

Processing using memory (PUM) exploits the existing memory architecture and the operational principles of the memory circuitry to enable operations within main memory with minimal changes. PUM makes use

https://people.inf.ethz.ch/omutlu/pub/ModernPrimerOnPIM springer-emerging-computing-bookchapter21.pdf

Onur Mutlu, Saugata Ghose, Juan Gomez-Luna, and Rachata Ausavarungnirun,

<u>"A Modern Primer on Processing in Memory"</u>

Invited Book Chapter in <u>Emerging</u>

<u>Computing: From Devices to Systems - Looking Beyond Moore and Von Neumann</u>,

Springer, to be published in 2021.

[<u>Tutorial Video on "Memory-Centric Computing</u>

Systems" (1 hour 51 minutes)]

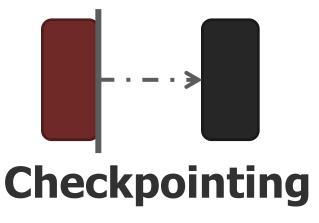
Starting Simple: Data Copy and Initialization

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







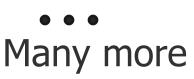




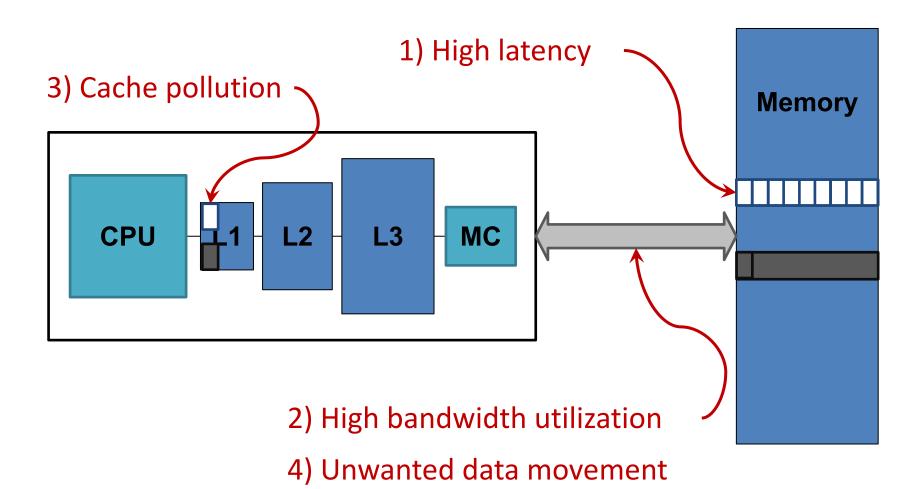




Page Migration

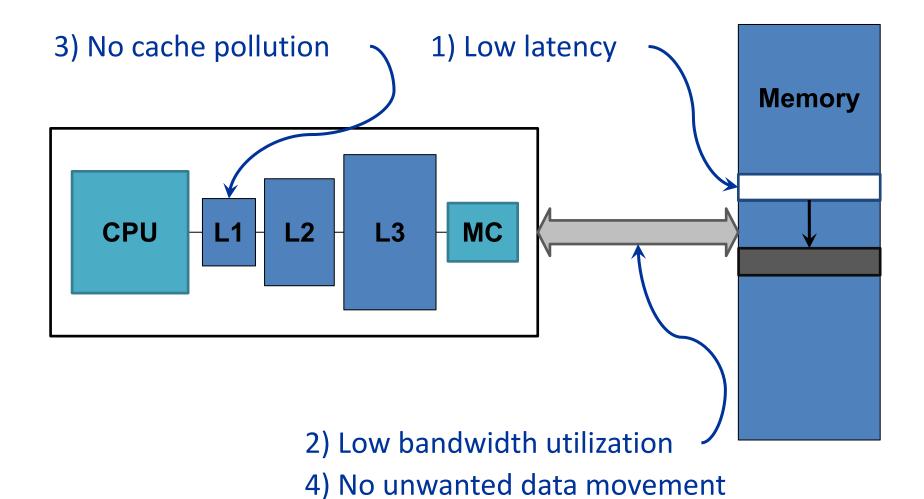


Today's Systems: Bulk Data Copy



1046ns, 3.6uJ (for 4KB page copy via DMA)

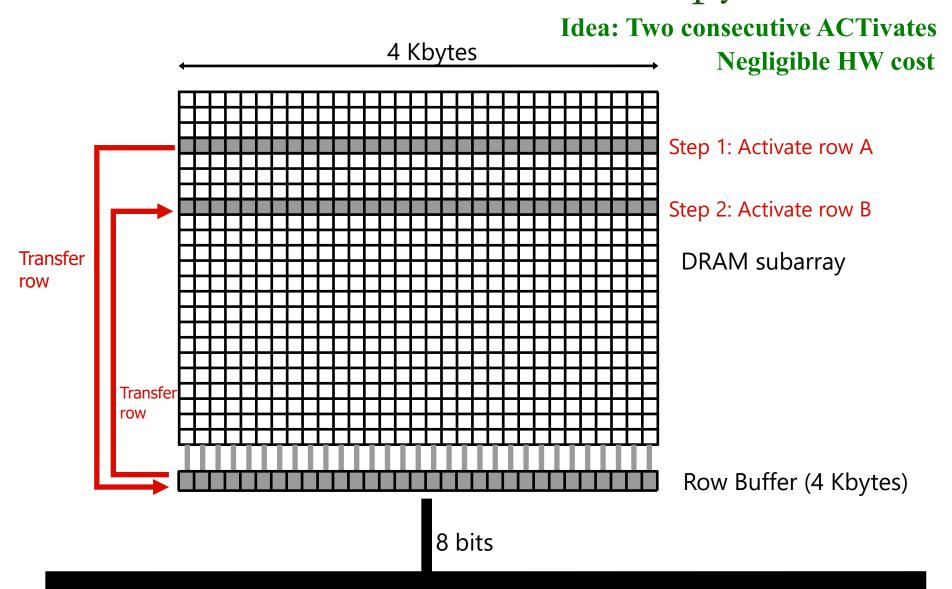
Future Systems: In-Memory Copy



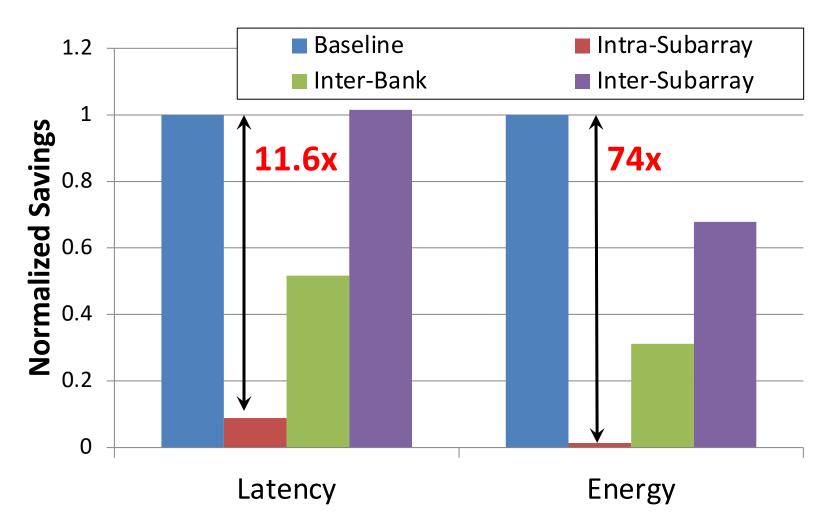
1046ns, 3.6uJ

→ 90ns, 0.04uJ

RowClone: In-DRAM Row Copy



RowClone: Latency and Energy Savings



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

More on RowClone

Vivek Seshadri, Yoongu Kim, Chris Fallin, Donghyuk Lee, Rachata
 Ausavarungnirun, Gennady Pekhimenko, Yixin Luo, Onur Mutlu, Michael A.
 Kozuch, Phillip B. Gibbons, and Todd C. Mowry,

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

Proceedings of the <u>46th International Symposium on Microarchitecture</u> (**MICRO**), Davis, CA, December 2013. [<u>Slides (pptx) (pdf)</u>] [<u>Lightning Session 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 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

RowClone Extensions and Follow-Up Work

- Can this be improved to do faster inter-subarray copy?
 - Yes, see LISA [Chang et al., HPCA 2016]
- Can we enable data movement at smaller granularities within a bank?
 - Yes, see FIGARO [Wang et al., MICRO 2020]
- Can this be improved to do better inter-bank copy?
 - Yes, see Network-on-Memory [CAL 2020]
- Can similar ideas and DRAM properties be used to perform computation on data?
 - Yes, see Ambit [Seshadri et al., CAL 2015, MICRO 2017]

LISA: Increasing Connectivity in DRAM

Kevin K. Chang, Prashant J. Nair, Saugata Ghose, Donghyuk Lee,
 Moinuddin K. Qureshi, and Onur Mutlu,

"Low-Cost Inter-Linked Subarrays (LISA): Enabling Fast Inter-Subarray Data Movement in DRAM"

Proceedings of the <u>22nd International Symposium on High-</u> <u>Performance Computer Architecture</u> (**HPCA**), Barcelona, Spain, March 2016.

[Slides (pptx) (pdf)]
[Source Code]

Low-Cost Inter-Linked Subarrays (LISA): Enabling Fast Inter-Subarray Data Movement in DRAM

Kevin K. Chang[†], Prashant J. Nair*, Donghyuk Lee[†], Saugata Ghose[†], Moinuddin K. Qureshi*, and Onur Mutlu[†]

†Carnegie Mellon University *Georgia Institute of Technology

FIGARO: Fine-Grained In-DRAM Copy

Yaohua Wang, Lois Orosa, Xiangjun Peng, Yang Guo, Saugata Ghose, Minesh Patel, Jeremie S. Kim, Juan Gómez Luna, Mohammad Sadrosadati, Nika Mansouri Ghiasi, and Onur Mutlu, "FIGARO: Improving System Performance via Fine-Grained In-DRAM Data Relocation and Caching"
Proceedings of the <u>53rd International Symposium on</u> Microarchitecture (MICRO), Virtual, October 2020.

FIGARO: Improving System Performance via Fine-Grained In-DRAM Data Relocation and Caching

Yaohua Wang* Lois Orosa[†] Xiangjun Peng[⊙]* Yang Guo* Saugata Ghose^{◇‡} Minesh Patel[†] Jeremie S. Kim[†] Juan Gómez Luna[†] Mohammad Sadrosadati[§] Nika Mansouri Ghiasi[†] Onur Mutlu^{†‡}

*National University of Defense Technology † ETH Zürich $^{\odot}$ Chinese University of Hong Kong $^{\diamond}$ University of Illinois at Urbana–Champaign ‡ Carnegie Mellon University § Institute of Research in Fundamental Sciences

Network-On-Memory: Fast Inter-Bank Copy

 Seyyed Hossein SeyyedAghaei Rezaei, Mehdi Modarressi, Rachata Ausavarungnirun, Mohammad Sadrosadati, Onur Mutlu, and Masoud Daneshtalab,

"NoM: Network-on-Memory for Inter-Bank Data Transfer in Highly-Banked Memories"

<u>IEEE Computer Architecture Letters</u> (CAL), to appear in 2020.

NoM: Network-on-Memory for Inter-bank Data Transfer in Highly-banked Memories

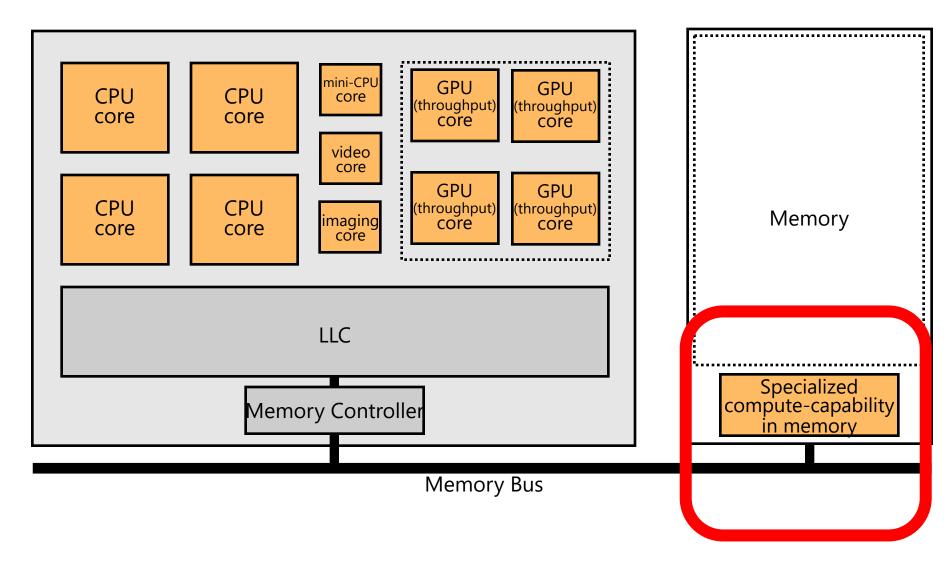
Seyyed Hossein SeyyedAghaei Rezaei¹
Mohammad Sadrosadati³

Mehdi Modarressi^{1,3} Rachata Ausavarungnirun² Onur Mutlu⁴ Masoud Daneshtalab⁵

¹University of Tehran

²King Mongkut's University of Technology North Bangkok ³Institute for Research in Fundamental Sciences ⁴ETH Zürich ⁵Mälardalens University

Mindset: Memory as an Accelerator



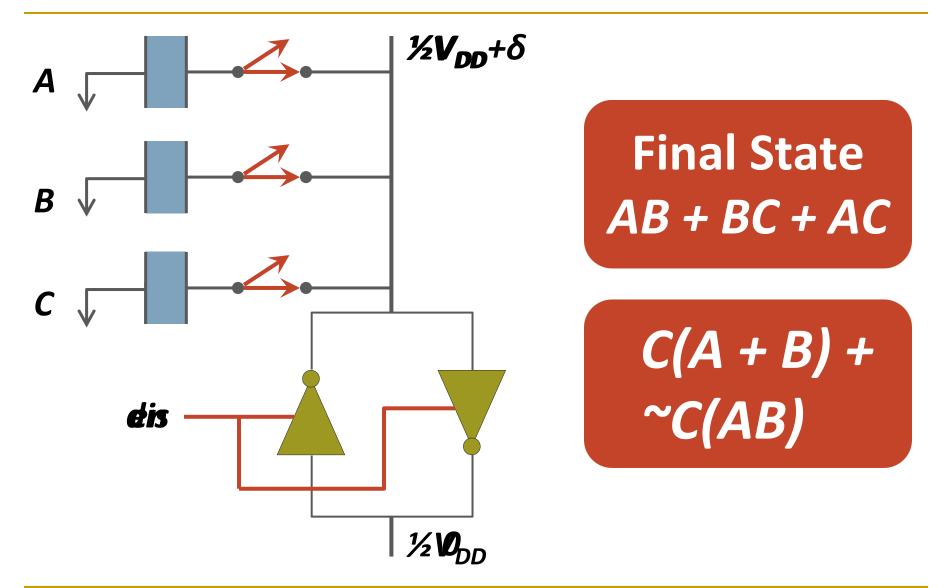
Memory similar to a "conventional" accelerator

In-Memory Bulk Bitwise Operations

- We can also support in-DRAM 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



In-DRAM Acceleration of Database Queries

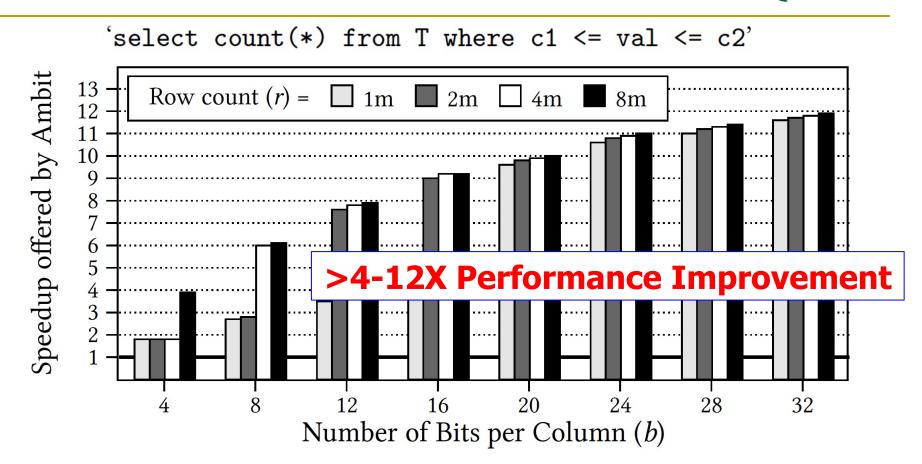


Figure 11: Speedup offered by Ambit over baseline CPU with SIMD for BitWeaving

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,

"Fast Bulk Bitwise AND and OR in DRAM"

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 Hassan 4 Amirali Boroumand 5 Jeremie Kim 4,5 Michael A. Kozuch 3 Onur Mutlu 4,5 Phillip B. Gibbons 5 Todd C. Mowry 5

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

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

SIMDRAM Framework

Nastaran Hajinazar, Geraldo F. Oliveira, Sven Gregorio, Joao Dinis Ferreira, Nika Mansouri Ghiasi, Minesh Patel, Mohammed Alser, Saugata Ghose, Juan Gomez-Luna, and Onur Mutlu, "SIMDRAM: An End-to-End Framework for Bit-Serial SIMD Computing in DRAM" Proceedings of the 26th International Conference on Architectural Support for Programming Languages and Operating Systems (ASPLOS), Virtual, March-April 2021.

[2-page Extended Abstract]

[Short Talk Slides (pptx) (pdf)]

[Talk Slides (pptx) (pdf)]

[Short Talk Video (5 mins)]

[Full Talk Video (27 mins)]

SIMDRAM: A Framework for Bit-Serial SIMD Processing using DRAM

*Nastaran Hajinazar^{1,2}
Nika Mansouri Ghiasi¹

*Geraldo F. Oliveira¹
Minesh Patel¹
Juan Gómez-Luna¹

Sven Gregorio¹ Mohammed Alser¹ Onur Mutlu¹

João Dinis Ferreira¹ Saugata Ghose³

¹ETH Zürich

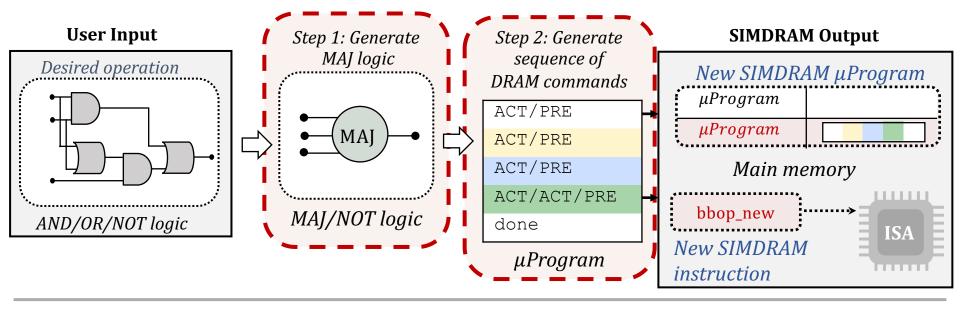
²Simon Fraser University

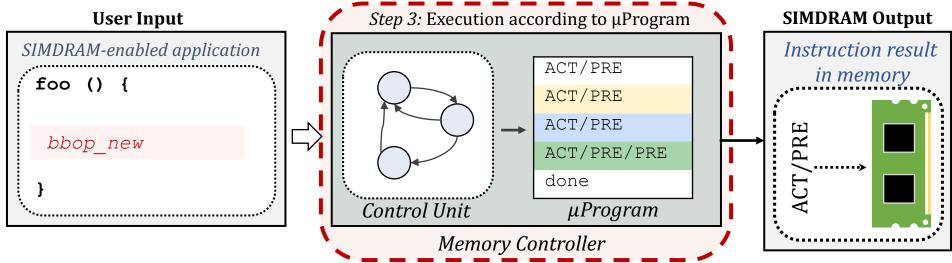
³University of Illinois at Urbana-Champaign

SIMDRAM Key Idea

- **SIMDRAM**: An end-to-end processing-using-DRAM framework that provides the programming interface, the ISA, and the hardware support for:
 - Efficiently computing complex operations in DRAM
 - Providing the ability to implement **arbitrary** operations as required
 - Using an **in-DRAM massively-parallel SIMD substrate** that requires **minimal** changes to DRAM architecture

SIMDRAM Framework: Overview





More on the SIMDRAM Framework

Nastaran Hajinazar, Geraldo F. Oliveira, Sven Gregorio, Joao Dinis Ferreira, Nika Mansouri Ghiasi, Minesh Patel, Mohammed Alser, Saugata Ghose, Juan Gomez-Luna, and Onur Mutlu, "SIMDRAM: An End-to-End Framework for Bit-Serial SIMD Computing in DRAM" Proceedings of the 26th International Conference on Architectural Support for Programming Languages and Operating Systems (ASPLOS), Virtual, March-April 2021.

[2-page Extended Abstract]

[Short Talk Slides (pptx) (pdf)]

[Talk Slides (pptx) (pdf)]

[Short Talk Video (5 mins)]

[Full Talk Video (27 mins)]

SIMDRAM: A Framework for Bit-Serial SIMD Processing using DRAM

*Nastaran Hajinazar^{1,2} Nika Mansouri Ghiasi¹ *Geraldo F. Oliveira¹
Minesh Patel¹
Juan Gómez-Luna¹

Sven Gregorio¹ Mohammed Alser¹ Onur Mutlu¹

João Dinis Ferreira¹ Saugata Ghose³

¹ETH Zürich

²Simon Fraser University

³University of Illinois at Urbana–Champaign

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 Devices"

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> <u>Architecture</u> (**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

SAFARI 238

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,

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

Proceedings of the <u>48th International Symposium on Computer Architecture</u> (**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 Olgun^{§†} Minesh Patel[§] A. Giray Yağlıkçı[§] Haocong Luo[§] Jeremie S. Kim[§] F. Nisa Bostancı^{§†} Nandita Vijaykumar^{§⊙} Oğuz Ergin[†] Onur Mutlu[§]

§ETH Zürich † TOBB University of Economics and Technology $^{\odot}$ University of Toronto

SAFARI 239

Processing in Memory: Two Approaches

- 1. Processing using Memory
- 2. Processing near Memory

Another Example: In-Memory Graph Processing

Large graphs are everywhere (circa 2015)



36 Million Wikipedia Pages



1.4 Billion Facebook Users

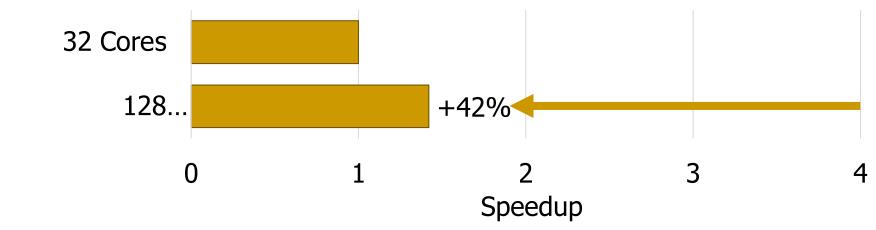


300 Million Twitter Users



30 Billion Instagram Photos

Scalable large-scale graph processing is challenging

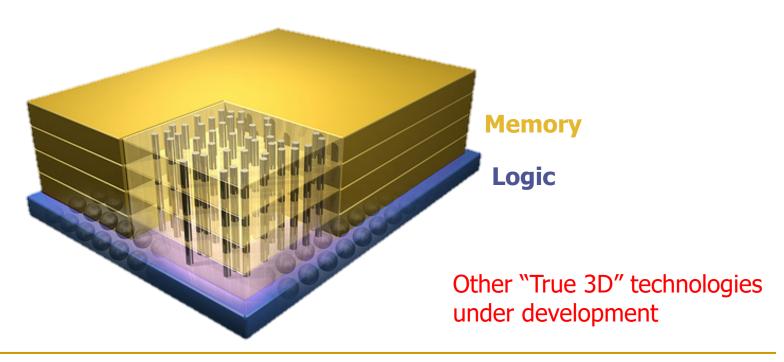


Key Bottlenecks in Graph Processing

```
for (v: graph.vertices) {
     for (w: v.successors) {
       w.next rank += weight * v.rank;
                       1. Frequent random memory accesses
                                   &w
            V
 w.rank
w.next rank
                              weight * v.rank
 w.edges
            W
                              2. Little amount of computation
```

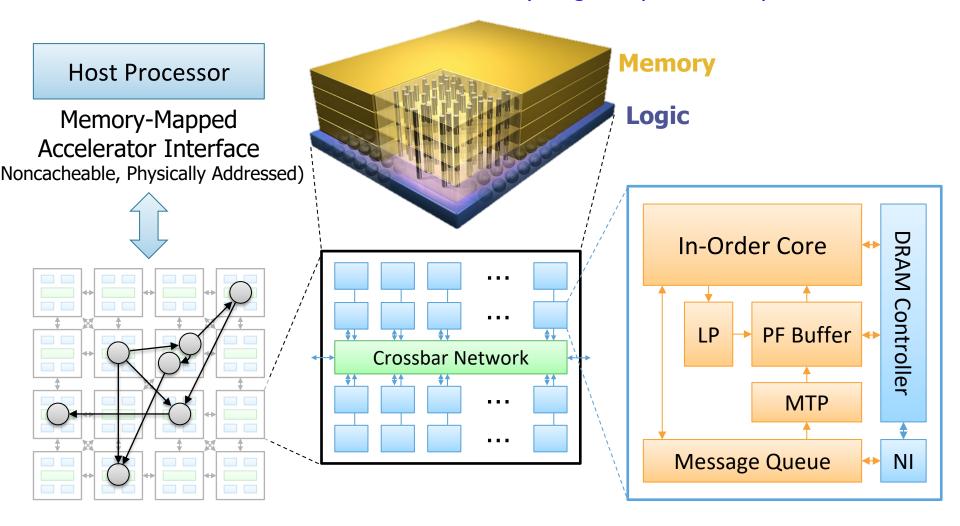
Opportunity: 3D-Stacked Logic+Memory



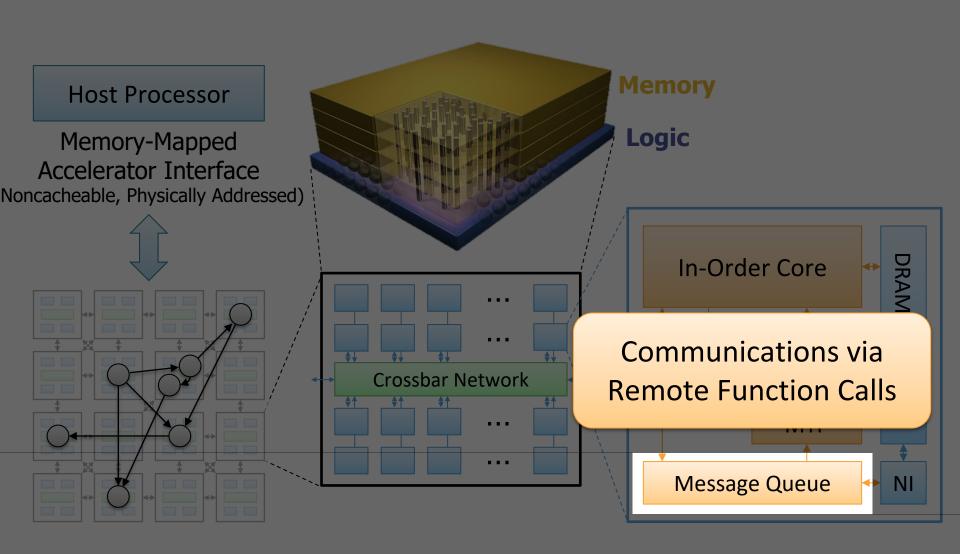


Tesseract System for Graph Processing

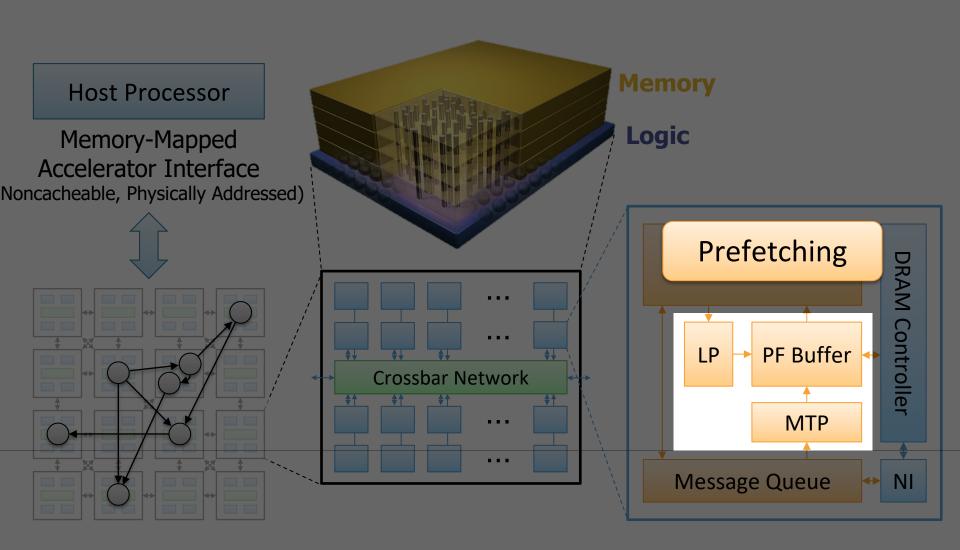
Interconnected set of 3D-stacked memory+logic chips with simple cores



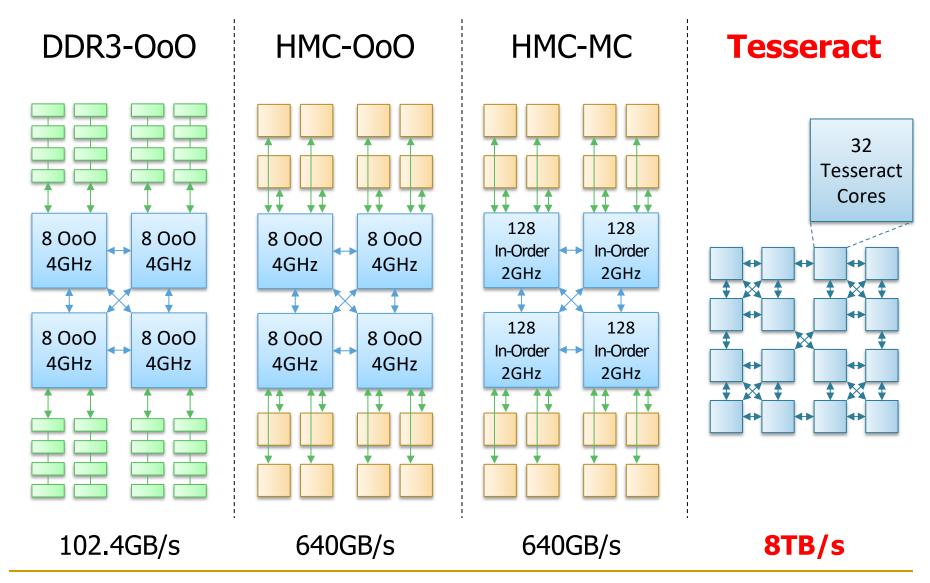
Tesseract System for Graph Processing



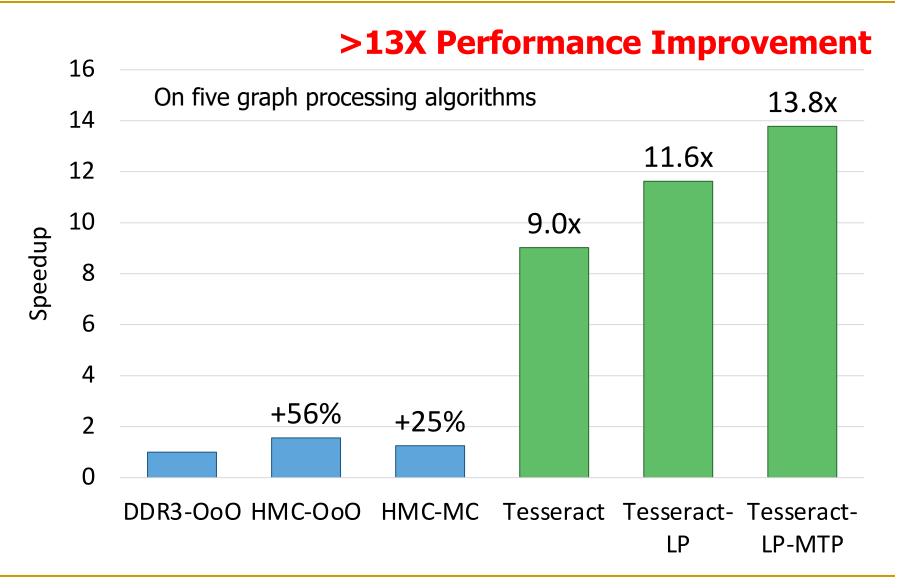
Tesseract System for Graph Processing



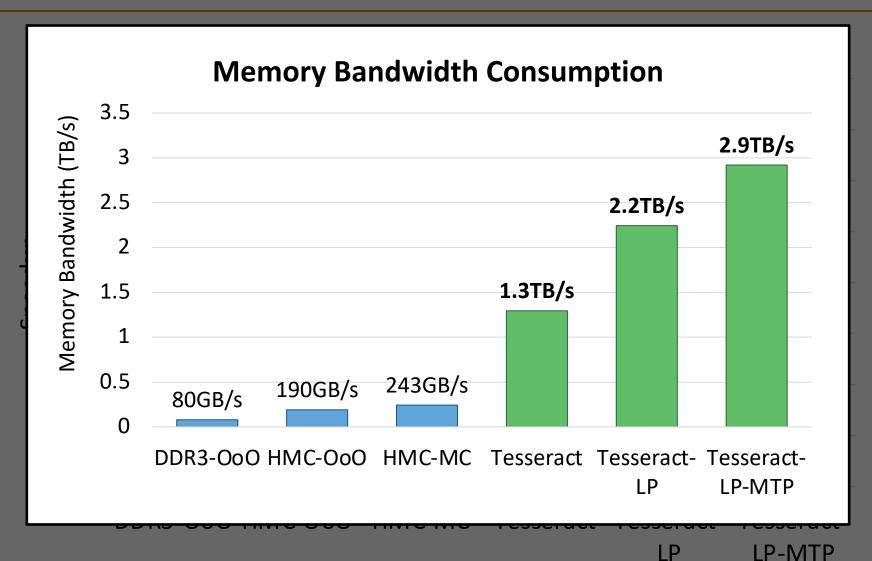
Evaluated Systems



Tesseract Graph Processing Performance

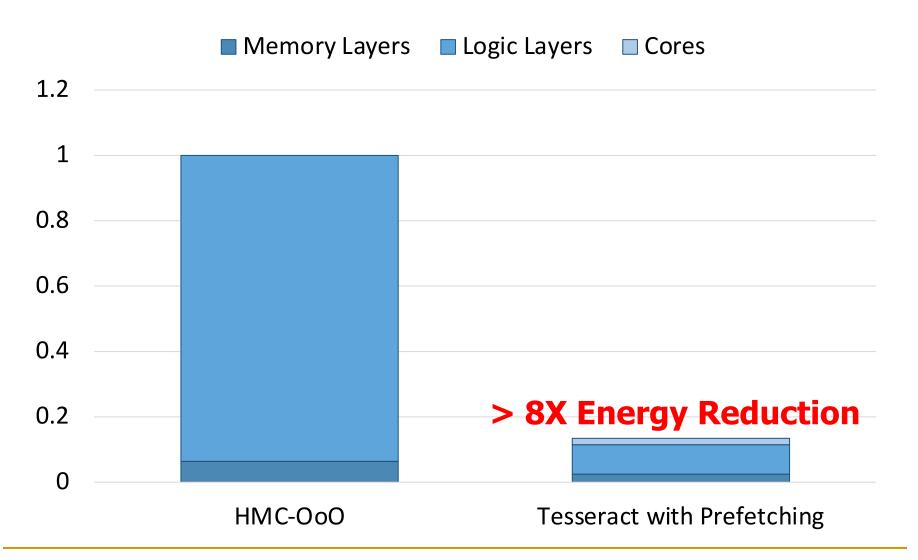


Tesseract Graph Processing Performance



249

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

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 Boroumand¹ Saugata Ghose¹ Youngsok Kim² Rachata Ausavarungnirun¹ Eric Shiu³ Rahul Thakur³ Daehyun Kim^{4,3} Aki Kuusela³ Allan Knies³ Parthasarathy Ranganathan³ Onur Mutlu^{5,1}

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













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



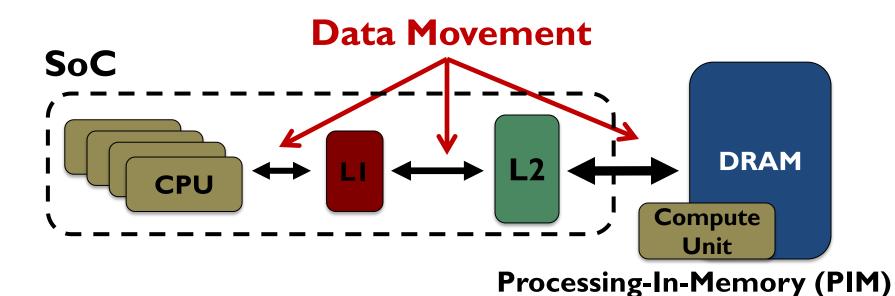
Google's video codec



Google's video codec

Energy Cost of Data Movement

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



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 2.3X and 2.2X

Workload Analysis



Chrome

Google's web browser



TensorFlow Mobile

Google's machine learning framework

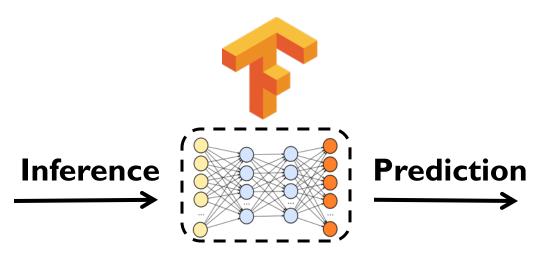


Google's video codec



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

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</u>
<u>Programming Languages and Operating Systems</u> (**ASPLOS**), Williamsburg, VA, USA, March 2018.

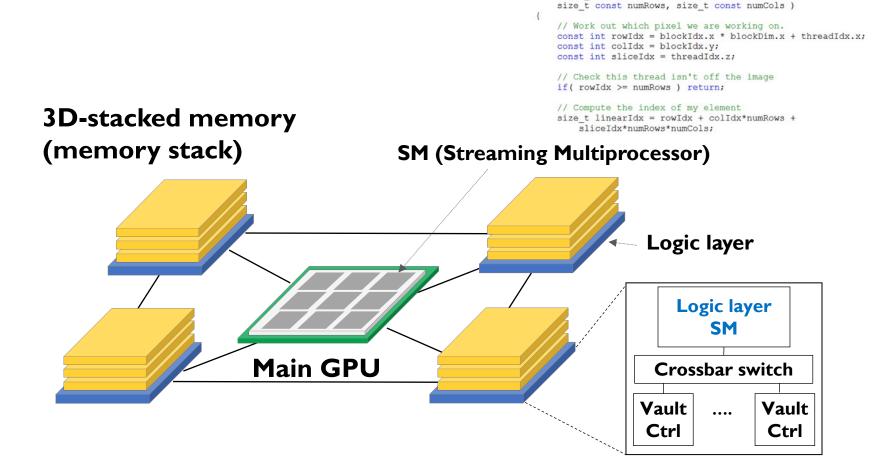
[Slides (pptx) (pdf)] [Lightning Session Slides (pptx) (pdf)] [Poster (pptx) (pdf)] [Lightning Talk Video (2 minutes)] [Full Talk Video (21 minutes)]

Google Workloads for Consumer Devices: Mitigating Data Movement Bottlenecks

Amirali Boroumand¹ Saugata Ghose¹ Youngsok Kim² Rachata Ausavarungnirun¹ Eric Shiu³ Rahul Thakur³ Daehyun Kim^{4,3} Aki Kuusela³ Allan Knies³ Parthasarathy Ranganathan³ Onur Mutlu^{5,1}

SAFARI

Truly Distributed GPU Processing with PIM



void applyScaleFactorsKernel(uint8_T * const out, uint8_T const * const in, const double *factor,

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> <u>Architecture</u> (**ISCA**), Seoul, South Korea, June 2016. [Slides (pptx) (pdf)]

Lightning Session Slides (pptx) (pdf)

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

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

Accelerating GPU Execution with PIM (II)

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

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,
 "Accelerating Pointer Chasing in 3D-Stacked Memory:
 Challenges, Mechanisms, Evaluation"
 Proceedings of the 34th IEEE International Conference on Computer
 Design (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 Memory Controller"

Proceedings of the <u>43rd International Symposium on Computer</u> <u>Architecture</u> (**ISCA**), Seoul, South Korea, June 2016. [Slides (pptx) (pdf)]

[Lightning Session Slides (pptx) (pdf)]

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

Accelerating Runahead Execution

Milad Hashemi, Onur Mutlu, and Yale N. Patt,
 "Continuous Runahead: Transparent Hardware Acceleration for Memory Intensive Workloads"
 Proceedings of the 49th International Symposium on Microarchitecture (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

Accelerating Climate Modeling

 Gagandeep Singh, Dionysios Diamantopoulos, Christoph Hagleitner, Juan Gómez-Luna, Sander Stuijk, Onur Mutlu, and Henk Corporaal, "NERO: A Near High-Bandwidth Memory Stencil Accelerator for Weather Prediction Modeling"

Proceedings of the <u>30th International Conference on Field-Programmable Logic</u> <u>and Applications</u> (**FPL**), Gothenburg, Sweden, September 2020.

[Slides (pptx) (pdf)]

[Lightning Talk Slides (pptx) (pdf)]

[Talk Video (23 minutes)]

Nominated for the Stamatis Vassiliadis Memorial Award.

NERO: A Near High-Bandwidth Memory Stencil Accelerator for Weather Prediction Modeling

Gagandeep Singh a,b,c Dionysios Diamantopoulos c Christoph Hagleitner c Juan Gómez-Luna b Sander Stuijk a Onur Mutlu b Henk Corporaal a Eindhoven University of Technology b ETH Zürich c IBM Research Europe, Zurich

Accelerating Approximate String Matching

Damla Senol Cali, Gurpreet S. Kalsi, Zulal Bingol, Can Firtina, Lavanya Subramanian, Jeremie S. Kim, Rachata Ausavarungnirun, Mohammed Alser, Juan Gomez-Luna, Amirali Boroumand, Anant Nori, Allison Scibisz, Sreenivas Subramoney, Can Alkan, Saugata Ghose, and Onur Mutlu, "GenASM: A High-Performance, Low-Power Approximate String Matching Acceleration Framework for Genome Sequence Analysis"
Proceedings of the 53rd International Symposium on Microarchitecture (MICRO), Virtual, October 2020.

[<u>Lighting Talk Video</u> (1.5 minutes)] [<u>Lightning Talk Slides (pptx) (pdf)</u>] [<u>Talk Video</u> (18 minutes)] [<u>Slides (pptx) (pdf)</u>]

GenASM: A High-Performance, Low-Power Approximate String Matching Acceleration Framework for Genome Sequence Analysis

Damla Senol Cali^{†™} Gurpreet S. Kalsi[™] Zülal Bingöl[▽] Can Firtina[⋄] Lavanya Subramanian[‡] Jeremie S. Kim^{⋄†} Rachata Ausavarungnirun[⊙] Mohammed Alser[⋄] Juan Gomez-Luna[⋄] Amirali Boroumand[†] Anant Nori[™] Allison Scibisz[†] Sreenivas Subramoney[™] Can Alkan[▽] Saugata Ghose^{*†} Onur Mutlu^{⋄†▽}

† Carnegie Mellon University [™] Processor Architecture Research Lab, Intel Labs [▽] Bilkent University [⋄] ETH Zürich

‡ Facebook [⊙] King Mongkut's University of Technology North Bangkok ^{*} University of Illinois at Urbana–Champaign

268

Accelerating Time Series Analysis

Ivan Fernandez, Ricardo Quislant, Christina Giannoula, Mohammed Alser, Juan Gómez-Luna, Eladio Gutiérrez, Oscar Plata, and Onur Mutlu,
 "NATSA: A Near-Data Processing Accelerator for Time Series Analysis"
 Proceedings of the <u>38th IEEE International Conference on Computer</u>
 <u>Design</u> (ICCD), Virtual, October 2020.
 [Slides (pptx) (pdf)]
 [Talk Video (10 minutes)]

NATSA: A Near-Data Processing Accelerator for Time Series Analysis

Ivan Fernandez § Ricardo Quislant § Christina Giannoula † Mohammed Alser ‡ Juan Gómez-Luna ‡ Eladio Gutiérrez § Oscar Plata § Onur Mutlu ‡ § University of Malaga † National Technical University of Athens ‡ ETH Zürich

DAMOV Methodology & Workloads

DAMOV: A New Methodology and Benchmark Suite for Evaluating Data Movement Bottlenecks

GERALDO F. OLIVEIRA, ETH Zürich, Switzerland
JUAN GÓMEZ-LUNA, ETH Zürich, Switzerland
LOIS OROSA, ETH Zürich, Switzerland
SAUGATA GHOSE, University of Illinois at Urbana-Champaign, USA
NANDITA VIJAYKUMAR, University of Toronto, Canada
IVAN FERNANDEZ, University of Malaga, Spain & ETH Zürich, Switzerland
MOHAMMAD SADROSADATI, Institute for Research in Fundamental Sciences (IPM), Iran & ETH
Zürich, Switzerland
ONUR MUTLU, ETH Zürich, Switzerland

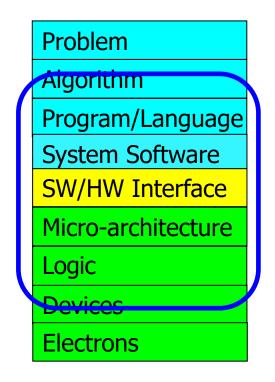
Data movement between the CPU and main memory is a first-order obstacle against improving performance, scalability, and energy efficiency in modern systems. Computer systems employ a range of techniques to reduce overheads tied to data movement, spanning from traditional mechanisms (e.g., deep multi-level cache hierarchies, aggressive hardware prefetchers) to emerging techniques such as Near-Data Processing (NDP), where some computation is moved close to memory. Prior NDP works investigate the root causes of data movement bottlenecks using different profiling methodologies and tools. However, there is still a lack of understanding about the key metrics that can identify different data movement bottlenecks and their relation to traditional and emerging data movement mitigation mechanisms. Our goal is to methodically identify potential sources of data movement over a broad set of applications and to comprehensively compare traditional compute-centric data movement mitigation techniques (e.g., caching and prefetching) to more memory-centric techniques (e.g., NDP), thereby developing a rigorous understanding of the best techniques to mitigate each source of data movement.

With this goal in mind, we perform the first large-scale characterization of a wide variety of applications, across a wide range of application domains, to identify fundamental program properties that lead to data movement to/from main memory. We develop the first systematic methodology to classify applications based on the sources contributing to data movement bottlenecks. From our large-scale characterization of 77K functions across 345 applications, we select 144 functions to form the first open-source benchmark suite (DAMOV) for main memory data movement studies. We select a diverse range of functions that (1) represent different types of data movement bottlenecks, and (2) come from a wide range of application domains. Using NDP as a case study, we identify new insights about the different data movement bottlenecks and use these insights to determine the most suitable data movement mitigation mechanism for a particular application. We open-source DAMOV and the complete source code for our new characterization methodology at https://github.com/CMU-SAFARI/DAMOV.

SAFARI

https://arxiv.org/pdf/2105.03725.pdf

We Need to Revisit the Entire Stack



We can get there step by step

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 Emerging Computing: From Devices to Systems
Looking Beyond Moore and Von Neumann, Springer, to be published in 2021.

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

Abstract

Modern computing systems are overwhelmingly designed to move data to computation. This design choice goes directly against at least three key trends in computing that cause performance, scalability and energy bottlenecks: (1) data access is a key bottleneck as many important applications are increasingly data-intensive, and memory bandwidth and energy do not scale well, (2) energy consumption is a key limiter in almost all computing platforms, especially server and mobile systems, (3) data movement, especially off-chip to on-chip, is very expensive in terms of bandwidth, energy and latency, much more so than computation. These trends are especially severely-felt in the data-intensive server and energy-constrained mobile systems of today.

At the same time, conventional memory technology is facing many technology scaling challenges in terms of reliability, energy, and performance. As a result, memory system architects are open to organizing memory in different ways and making it more intelligent, at the expense of higher cost. The emergence of 3D-stacked memory plus logic, the adoption of error correcting codes inside the latest DRAM chips, proliferation of different main memory standards and chips, specialized for different purposes (e.g., graphics, low-power, high bandwidth, low latency), and the necessity of designing new solutions to serious reliability and security issues, such as the RowHammer phenomenon, are an evidence of this trend.

This chapter discusses recent research that aims to practically enable computation close to data, an approach we call processing-in-memory (PIM). PIM places computation mechanisms in or near where the data is stored (i.e., inside the memory chips, in the logic layer of 3D-stacked memory, or in the memory controllers), so that data movement between the computation units and memory is reduced or eliminated. While the general idea of PIM is not new, we discuss motivating trends in applications as well as memory circuits/technology that greatly exacerbate the need for enabling it in modern computing systems. We examine at least two promising new approaches to designing PIM systems to accelerate important data-intensive applications: (1) processing using memory by exploiting analog operational properties of DRAM chips to perform massively-parallel operations in memory, with low-cost changes, (2) processing near memory by exploiting 3D-stacked memory technology design to provide high memory bandwidth and low memory latency to in-memory logic. In both approaches, we describe and tackle relevant cross-layer research, design, and adoption challenges in devices, architecture, systems, and programming models. Our focus is on the development of in-memory processing designs that can be adopted in real computing platforms at low cost. We conclude by discussing work on solving key challenges to the practical adoption of PIM.

Keywords: memory systems, data movement, main memory, processing-in-memory, near-data processing, computation-in-memory, processing using memory, processing near memory, 3D-stacked memory, non-volatile memory, energy efficiency, high-performance computing, computer architecture, computing paradigm, emerging technologies, memory scaling, technology scaling, dependable systems, robust systems, hardware security, system security, latency, low-latency computing

-	
Car	iteni
	пеш

1	Introduction		
2	Major Trends Affecting Main Memory		
3	The Need for Intelligent Memory Controllers		
	to Enhance Memory Scaling	6	
4	Perils of Processor-Centric Design	9	
5	Processing-in-Memory (PIM): Technology		
	Enablers and Two Approaches	12	
	5.1 New Technology Enablers: 3D-Stacked		
	Memory and Non-Volatile Memory	12	
	5.2 Two Approaches: Processing Using		
	Memory (PUM) vs. Processing Near		
	Memory (PNM)	13	
6	Processing Using Memory (PUM)	14	
U	6.1 RowClone	14	
	6.2 Ambit	15	
		17	
	6.3 Gather-Scatter DRAM		
	6.4 In-DRAM Security Primitives	17	
7	Processing Near Memory (PNM)	18	
	7.1 Tesseract: Coarse-Grained Application-		
	Level PNM Acceleration of Graph Pro-		
	cessing	19	
	7.2 Function-Level PNM Acceleration of		
	Mobile Consumer Workloads	20	
	7.3 Programmer-Transparent Function-		
	Level PNM Acceleration of GPU		
	Applications	21	
	7.4 Instruction-Level PNM Acceleration		
	with PIM-Enabled Instructions (PEI)	21	
	7.5 Function-Level PNM Acceleration of		
	Genome Analysis Workloads	22	
_	7.6 Application-Level PNM Acceleration of		
L	Time Series Analysis	23	
8	Enabling the Adoption of PIM	24	
	8.1 Programming Models and Code Genera-		
	tion for PIM	24	
	8.2 PIM Runtime: Scheduling and Data		
	Mapping	25	
	8.3 Memory Coherence	27	
	8.4 Virtual Memory Support	27	
	8.5 Data Structures for PIM	28	
	8.6 Benchmarks and Simulation Infrastruc-		
	tures	29	
	8.7 Real PIM Hardware Systems and Proto-	-	
	types	30	
	8.8 Security Considerations	30	
9	Conclusion and Future Outlook	31	

1. Introduction

Main memory, built using the Dynamic Random Access Memory (DRAM) technology, is a major component in nearly all computing systems, including servers, cloud platforms, mobile/embedded devices, and sensor systems. Across all of these systems, the data working set sizes of modern applications are rapidly growing, while the need for fast analysis of such data is increasing. Thus, main memory is becoming an increasingly significant bottleneck across a wide variety of computing systems and applications [1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16]. Alleviating the main memory bottleneck requires the memory capacity, energy, cost, and performance to all scale in an efficient manner across technology generations. Unfortunately, it has become increasingly difficult in recent years, especially the past decade, to scale all of these dimensions [1, 2, 17, 18, 19, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29, 30, 31, 32, 33, 34, 35, 36, 37, 38, 39, 40, 41, 42, 43, 44, 45, 46, 47, 48, 49], and thus the main memory bottleneck has been worsening.

A major reason for the main memory bottleneck is the high energy and latency cost associated with data movement. In modern computers, to perform any operation on data that resides in main memory, the processor must retrieve the data from main memory. This requires the memory controller to issue commands to a DRAM module across a relatively slow and power-hungry off-chip bus (known as the memory channel). The DRAM module sends the requested data across the memory channel, after which the data is placed in the caches and registers. The CPU can perform computation on the data once the data is in its registers. Data movement from the DRAM to the CPU incurs long latency and consumes a significant amount of energy [7, 50, 51, 52, 53, 54]. These costs are often exacerbated by the fact that much of the data brought into the caches is not reused by the CPU [52, 53, 55, 56], providing little benefit in return for the high latency and energy cost.

The cost of data movement is a fundamental issue with the processor-centric nature of contemporary computer systems. The CPU is considered to be the master in the system, and computation is performed only in the processor (and accelerators). In contrast, data storage and communication units, including the main memory, are treated as unintelligent workers that are incapable of computation. As a result of this processor-centric design paradigm, data moves a lot in the system between the computation units and communication/ storage units so that computation can be done on it. With the increasingly data-centric nature of contemporary and emerging appli-

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

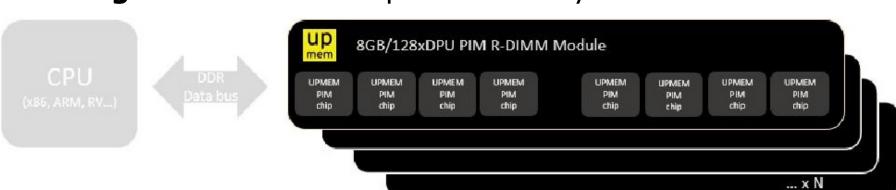
Hardware for Artificial Intelligence, to appear in November 2019.

[Preliminary arXiv version]

SAFARI

UPMEM Processing-in-DRAM Engine (2019)

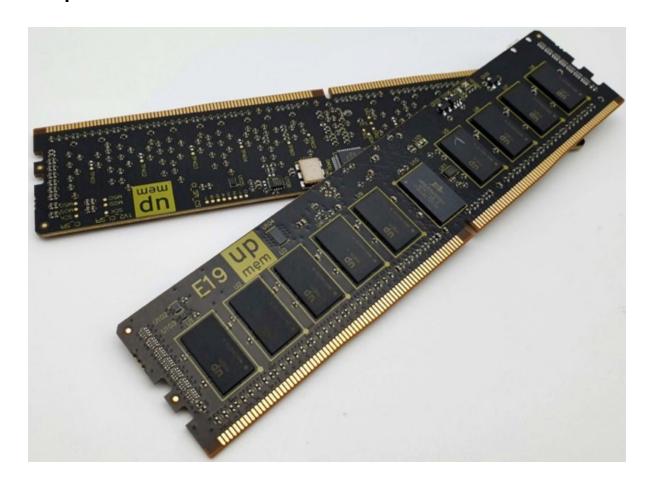
- Processing in DRAM Engine
- Includes standard DIMM modules, with a large number of DPU processors combined with DRAM chips.
- Replaces standard DIMMs
 - DDR4 R-DIMM modules
 - 8GB+128 DPUs (16 PIM chips)
 - Standard 2x-nm DRAM process
 - Large amounts of compute & memory bandwidth





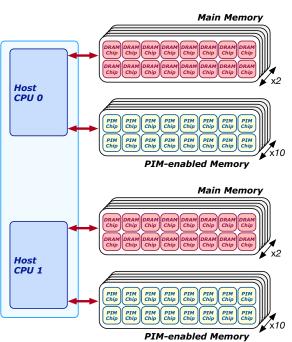
UPMEM Memory Modules

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





2,560-DPU Processing-in-Memory System



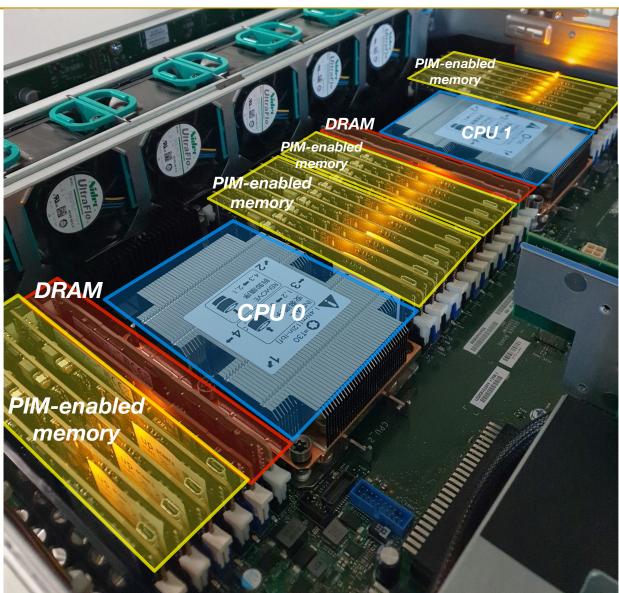
Benchmarking a New Paradigm: An Experimental Analysis of a Real Processing-in-Memory Architecture

JUAN GÓMEZ-LUNA, ETH Zürich, Switzerland
IZZAT EL HAJJ, 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

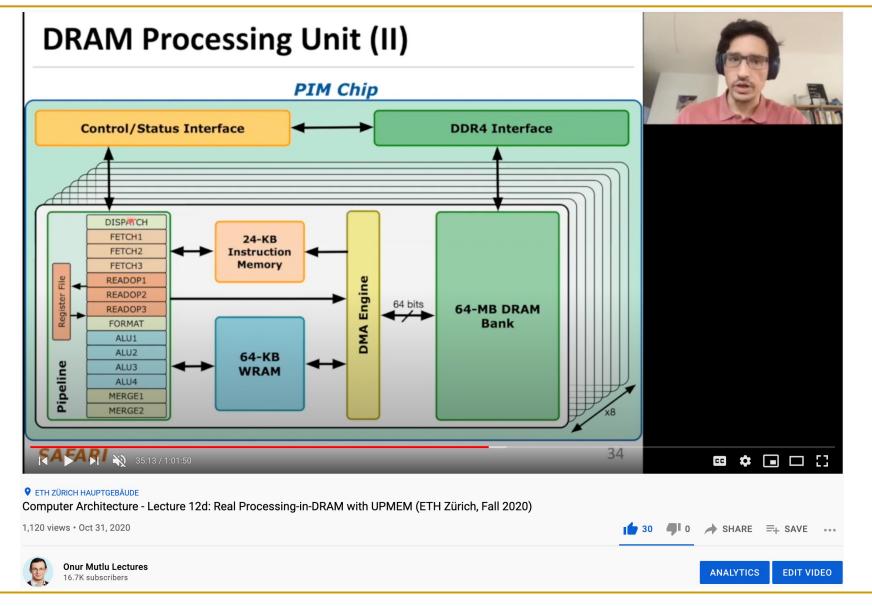
Many modern workloads, such as neural networks, databases, and graph processing, are fundamentally memory-bound. For such workloads, the data novement between main memory and CPU cores imposes a significant overhead in terns of both latency and energy. A major reason is that this communication happens 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 (PM).

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 PIM (Processing,-bendumpy) benchmarks), a benchmark suite of 16 workloads from different application domains (e.g., dense/sparse linear algebra, databases, data analytics, graph processing, which we identify as memory-bound. We evaluate the performance and scaling characteristics of PIM benchmarks on the UPMEM PIM architecture, and compare their performance and energy consumption to their state-of-the-art CPU and CPU counterparts. Our extensive evaluation conducted on two real UPMEM-based PIM systems with 460 and 25.50 DPUs provides new insights about suitability of different workloads to the PIM systems you for the programming recommendations for software designers, and suggestions and hints for hardware and architecture designers of future PIM systems.



More on the UPMEM PIM System



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 3D-stacked 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 state-of-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.

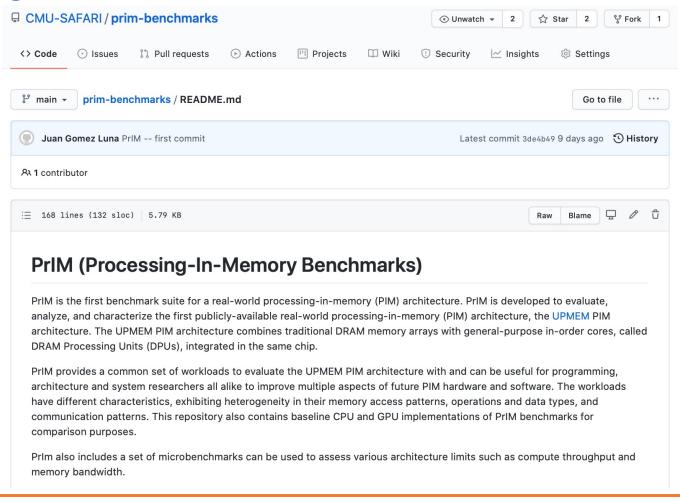
https://arxiv.org/pdf/2105.03814.pdf

PrIM Benchmarks: Application Domains

Domain	Benchmark	Short name
Dance linear algebra	Vector Addition	VA
Dense linear algebra	Matrix-Vector Multiply	GEMV
Sparse linear algebra	Sparse Matrix-Vector Multiply	SpMV
	Select	SEL
Databases	Unique	UNI
Data analytica	Binary Search	BS
Data analytics	Time Series Analysis	TS
Graph processing	Breadth-First Search	BFS
Neural networks	Multilayer Perceptron	MLP
Bioinformatics	Needleman-Wunsch	NW
luna da mua assalin d	Image histogram (short)	HST-S
Image processing	Image histogram (large)	HST-L
	Reduction	RED
Devallal maioriticas	Prefix sum (scan-scan-add)	SCAN-SSA
Parallel primitives	Prefix sum (reduce-scan-scan)	SCAN-RSS
	Matrix transposition	TRNS

PrIM Benchmarks are Open Source

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



Understanding a Modern PIM Architecture

Understanding a Modern Processing-in-Memory Architecture: Benchmarking and Experimental Characterization

```
Juan Gómez-Luna<sup>1</sup> Izzat El Hajj<sup>2</sup> Ivan Fernandez<sup>1,3</sup> Christina Giannoula<sup>1,4</sup> Geraldo F. Oliveira<sup>1</sup> Onur Mutlu<sup>1</sup>
```

¹ETH Zürich ²American University of Beirut ³University of Malaga ⁴National Technical University of Athens

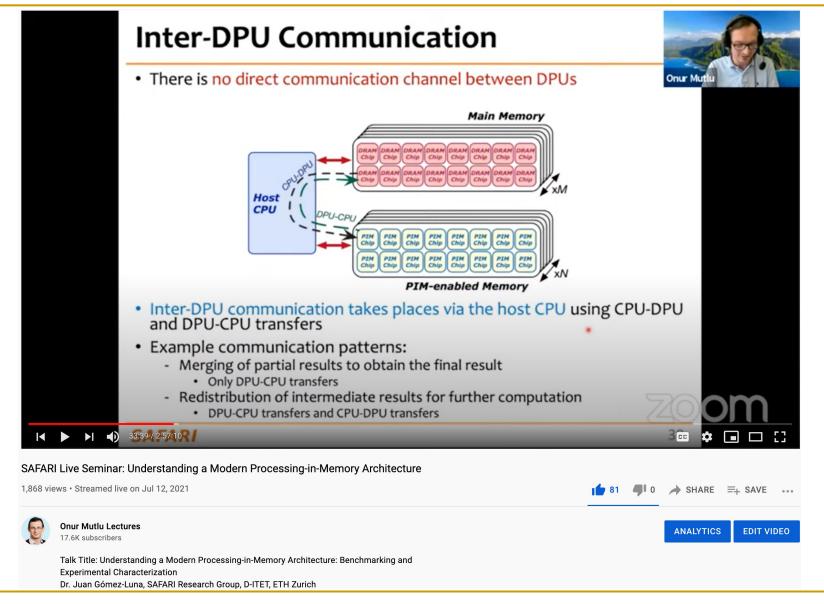
https://arxiv.org/pdf/2105.03814.pdf

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

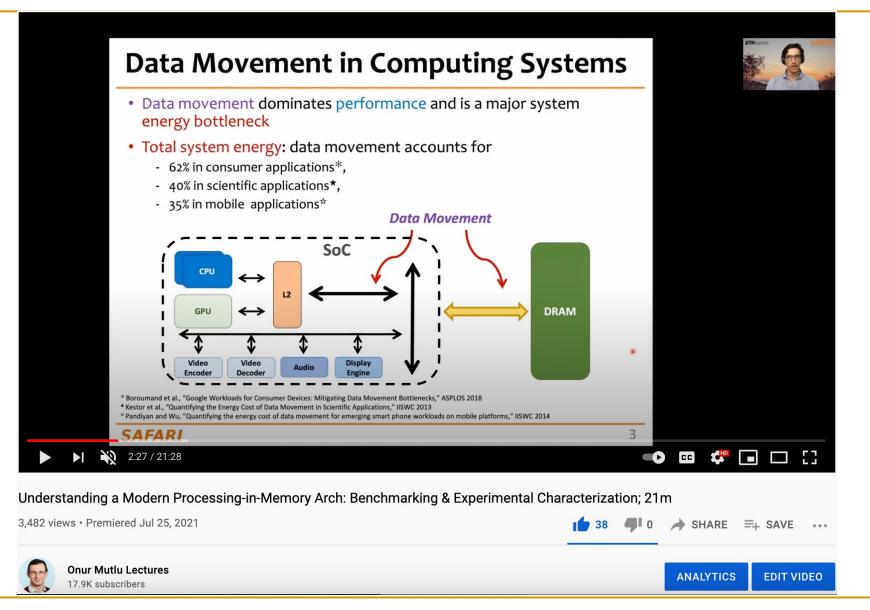
Understanding a Modern PIM Architecture



More on Analysis of the UPMEM PIM Engine



More on Analysis of the UPMEM PIM Engine



FPGA-based Processing Near Memory

Gagandeep Singh, Mohammed Alser, Damla Senol Cali, Dionysios
Diamantopoulos, Juan Gómez-Luna, Henk Corporaal, and Onur Mutlu,

"FPGA-based Near-Memory Acceleration of Modern Data-Intensive

Applications"

IEEE Micro (IEEE MICRO), to appear, 2021.

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

Gagandeep Singh[⋄] Mohammed Alser[⋄] Damla Senol Cali[⋈]
Dionysios Diamantopoulos[▽] Juan Gómez-Luna[⋄]
Henk Corporaal[⋆] Onur Mutlu^{⋄⋈}

[⋄]ETH Zürich [⋈] Carnegie Mellon University *Eindhoven University of Technology [▽]IBM Research Europe

DAMOV Analysis Methodology & Workloads

DAMOV: A New Methodology and Benchmark Suite for Evaluating Data Movement Bottlenecks

GERALDO F. OLIVEIRA, ETH Zürich, Switzerland
JUAN GÓMEZ-LUNA, ETH Zürich, Switzerland
LOIS OROSA, ETH Zürich, Switzerland
SAUGATA GHOSE, University of Illinois at Urbana-Champaign, USA
NANDITA VIJAYKUMAR, University of Toronto, Canada
IVAN FERNANDEZ, University of Malaga, Spain & ETH Zürich, Switzerland
MOHAMMAD SADROSADATI, Institute for Research in Fundamental Sciences (IPM), Iran & ETH
Zürich, Switzerland
ONUR MUTLU, ETH Zürich, Switzerland

Data movement between the CPU and main memory is a first-order obstacle against improving performance, scalability, and energy efficiency in modern systems. Computer systems employ a range of techniques to reduce overheads tied to data movement, spanning from traditional mechanisms (e.g., deep multi-level cache hierarchies, aggressive hardware prefetchers) to emerging techniques such as Near-Data Processing (NDP), where some computation is moved close to memory. Prior NDP works investigate the root causes of data movement bottlenecks using different profiling methodologies and tools. However, there is still a lack of understanding about the key metrics that can identify different data movement bottlenecks and their relation to traditional and emerging data movement mitigation mechanisms. Our goal is to methodically identify potential sources of data movement over a broad set of applications and to comprehensively compare traditional compute-centric data movement mitigation techniques (e.g., caching and prefetching) to more memory-centric techniques (e.g., NDP), thereby developing a rigorous understanding of the best techniques to mitigate each source of data movement.

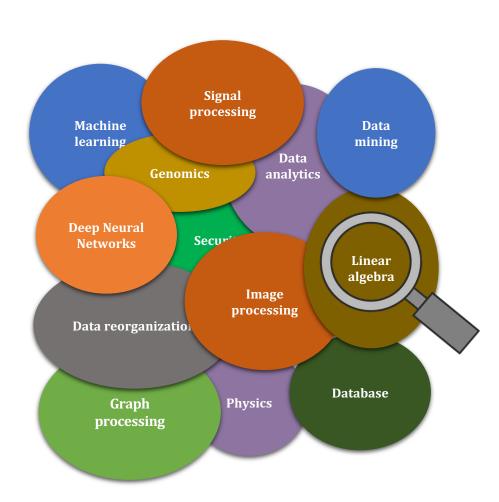
With this goal in mind, we perform the first large-scale characterization of a wide variety of applications, across a wide range of application domains, to identify fundamental program properties that lead to data movement to/from main memory. We develop the first systematic methodology to classify applications based on the sources contributing to data movement bottlenecks. From our large-scale characterization of 77K functions across 345 applications, we select 144 functions to form the first open-source benchmark suite (DAMOV) for main memory data movement studies. We select a diverse range of functions that (1) represent different types of data movement bottlenecks, and (2) come from a wide range of application domains. Using NDP as a case study, we identify new insights about the different data movement bottlenecks and use these insights to determine the most suitable data movement mitigation mechanism for a particular application. We open-source DAMOV and the complete source code for our new characterization methodology at https://github.com/CMU-SAFARI/DAMOV.

SAFARI

https://arxiv.org/pdf/2105.03725.pdf

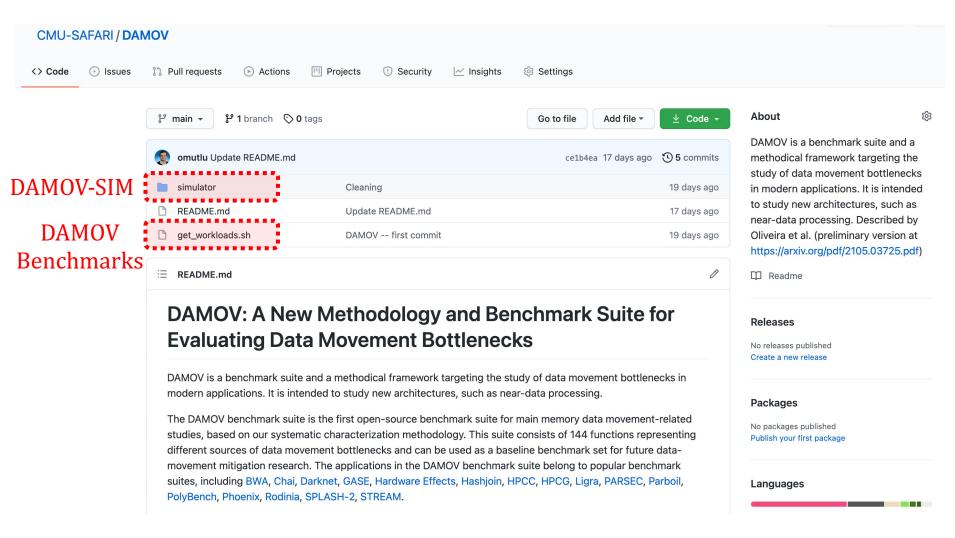
Step 1: Application Profiling

- We analyze 345 applications from distinct domains:
- Graph Processing
- Deep Neural Networks
- Physics
- High-Performance Computing
- Genomics
- Machine Learning
- Databases
- Data Reorganization
- Image Processing
- Map-Reduce
- Benchmarking
- Linear Algebra



DAMOV is Open Source

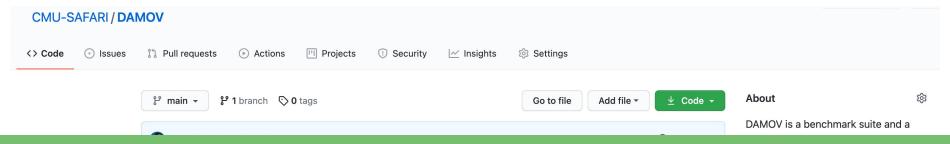
We open-source our benchmark suite and our toolchain





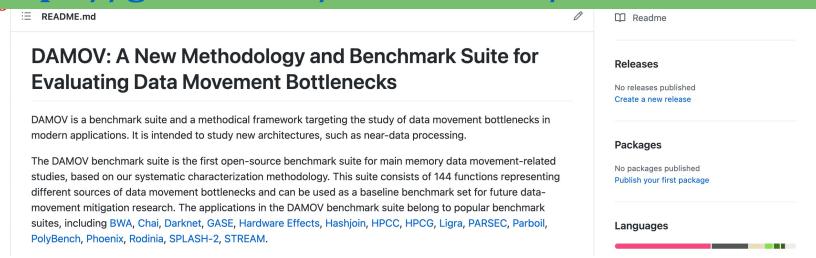
DAMOV is Open Source

We open-source our benchmark suite and our toolchain



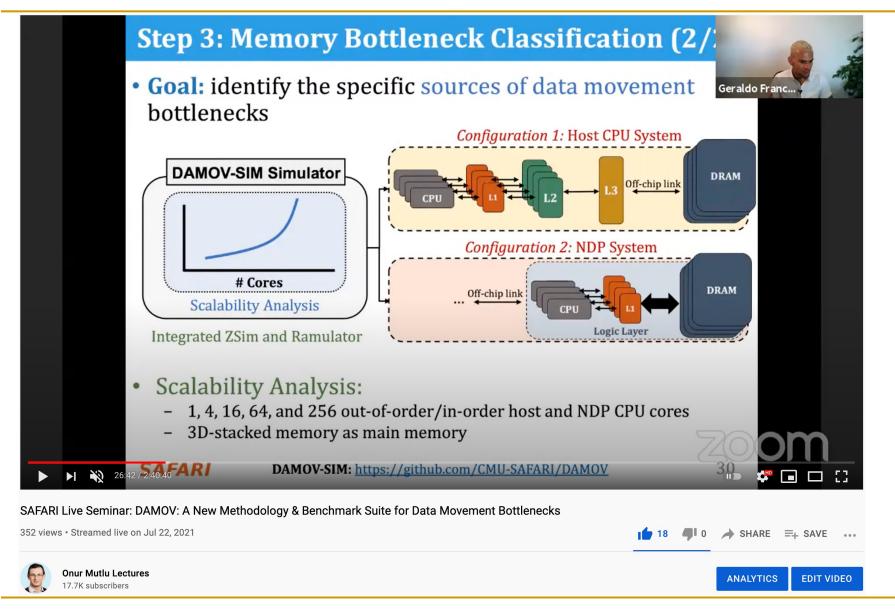
Get DAMOV at:

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





More on DAMOV Analysis Methodology & Workloads



More on DAMOV

Geraldo F. Oliveira, Juan Gomez-Luna, Lois Orosa, Saugata Ghose, Nandita Vijaykumar, Ivan fernandez, Mohammad Sadrosadati, and Onur Mutlu,
 "DAMOV: A New Methodology and Benchmark Suite for Evaluating Data Movement Bottlenecks"

Preprint in <u>arXiv</u>, 8 May 2021.

[arXiv preprint]

[DAMOV Suite and Simulator Source Code]

[SAFARI Live Seminar Video (2 hrs 40 mins)]

ONUR MUTLU, ETH Zürich, Switzerland

[Short Talk Video (21 minutes)]

DAMOV: A New Methodology and Benchmark Suite for Evaluating Data Movement Bottlenecks

GERALDO F. OLIVEIRA, ETH Zürich, Switzerland
JUAN GÓMEZ-LUNA, ETH Zürich, Switzerland
LOIS OROSA, ETH Zürich, Switzerland
SAUGATA GHOSE, University of Illinois at Urbana-Champaign, USA
NANDITA VIJAYKUMAR, University of Toronto, Canada
IVAN FERNANDEZ, University of Malaga, Spain & ETH Zürich, Switzerland
MOHAMMAD SADROSADATI, ETH Zürich, Switzerland

Samsung Function-in-Memory DRAM (2021)

Samsung Newsroom

CORPORATE

PRODUCTS

PRESS RESOURCES

VIEWS

ABOUT US

Q

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

Korea on February 17, 2021

Audio



Share (5





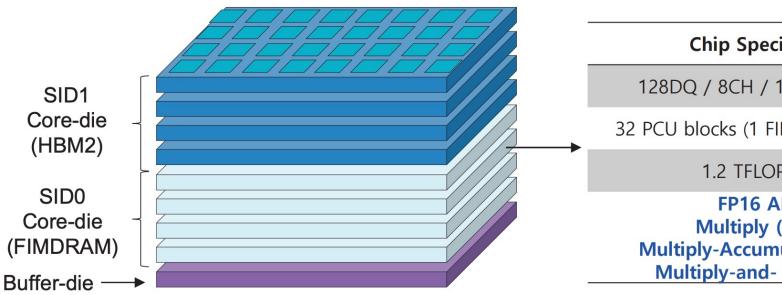
The new architecture will deliver over twice the system performance and reduce energy consumption by more than 70%

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 Al-driven workloads such as HPC, training and inference. We plan to build upon this breakthrough by further collaborating with Al solution providers for even more advanced PIM-powered applications."

Samsung Function-in-Memory DRAM (2021)

FIMDRAM based on HBM2



[3D Chip Structure of HBM with FIMDRAM]

Chip Specification

128DQ / 8CH / 16 banks / BL4

32 PCU blocks (1 FIM block/2 banks)

1.2 TFLOPS (4H)

FP16 ADD / Multiply (MUL) / Multiply-Accumulate (MAC) / Multiply-and- Add (MAD)

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

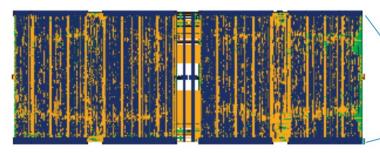
Young-Cheon Kwon¹, Suk Han Lee¹, Jaehoon Lee¹, Sang-Hyuk Kwon¹, Je Min Ryu1, Jong-Pil Son1, Seongil O1, Hak-Soo Yu1, Haesuk Lee1, Soo Young Kim¹, Youngmin Cho¹, Jin Guk Kim¹, Jongyoon Choi¹, Hyun-Sung Shin¹, Jin Kim¹, BengSeng Phuah¹, HyoungMin Kim¹, Myeong Jun Song¹, Ahn Choi¹, Daeho Kim¹, SooYoung Kim¹, Eun-Bong Kim¹, David Wang², Shinhaeng Kang¹, Yuhwan Ro³, Seungwoo Seo³, JoonHo Song³, Jaeyoun Youn1, Kyomin Sohn1, Nam Sung Kim1

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

Samsung Function-in-Memory DRAM (2021)

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 Kwont, Suk Han Lee!, Jaehoon Lee! Sang-Hruk Kwon', Je Min Ryu', Jong-Pil Son', Seongil O', Hak-Soo Yu', Haesuk Lee', Soo Young Kim', Youngmin Cho', Jin Guk Kim', Jongyoon Cho'r, Hyun-Sung Shin', Jin Kim', BengSeng Phuah', HyoungMin Kong, Ahn Choi, Jaehoo Kim', Soo'Young Kim', Eun-Bong Kim', David Wang', Shinhaeng Kang', Yuhwan Ro', Seungwoo Seo', JoonHo Song', Jaeyoun Youn', Kyomin Sohn', Man Sung Kim'

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		
Cell array for bank3	Cell array for bank7	Cell array for bank3	Cell array for bank7		
		TSV &	Peri C	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 bank9	Cell array for bank14 Cell array for bank13	Cell array for bank10 Cell array for bank9	Cell array for bank14 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

Detailed Lectures on PIM (I)

- Computer Architecture, Fall 2020, Lecture 6
 - Computation in Memory (ETH Zürich, Fall 2020)
 - https://www.youtube.com/watch?v=oGcZAGwfEUE&list=PL5Q2soXY2Zi9xidyIgBxUz 7xRPS-wisBN&index=12
- Computer Architecture, Fall 2020, Lecture 7
 - Near-Data Processing (ETH Zürich, Fall 2020)
 - https://www.youtube.com/watch?v=j2GIigqn1Qw&list=PL5Q2soXY2Zi9xidyIgBxUz7xRPS-wisBN&index=13
- Computer Architecture, Fall 2020, Lecture 11a
 - Memory Controllers (ETH Zürich, Fall 2020)
 - https://www.youtube.com/watch?v=TeG773OgiMQ&list=PL5Q2soXY2Zi9xidyIgBxUz 7xRPS-wisBN&index=20
- Computer Architecture, Fall 2020, Lecture 12d
 - Real Processing-in-DRAM with UPMEM (ETH Zürich, Fall 2020)
 - https://www.youtube.com/watch?v=Sscy1Wrr22A&list=PL5Q2soXY2Zi9xidyIgBxUz7xRPS-wisBN&index=25

Detailed Lectures on PIM (II)

- Computer Architecture, Fall 2020, Lecture 15
 - Emerging Memory Technologies (ETH Zürich, Fall 2020)
 - https://www.youtube.com/watch?v=AlE1rD9G_YU&list=PL5Q2soXY2Zi9xidyIgBxUz 7xRPS-wisBN&index=28
- Computer Architecture, Fall 2020, Lecture 16a
 - Opportunities & Challenges of Emerging Memory Technologies
 (ETH Zürich, Fall 2020)
 - https://www.youtube.com/watch?v=pmLszWGmMGQ&list=PL5Q2soXY2Zi9xidyIgBx Uz7xRPS-wisBN&index=29
- Computer Architecture, Fall 2020, Guest Lecture
 - In-Memory Computing: Memory Devices & Applications (ETH Zürich, Fall 2020)
 - https://www.youtube.com/watch?v=wNmqQHiEZNk&list=PL5Q2soXY2Zi9xidyIgBxUz7xRPS-wisBN&index=41

A Tutorial on PIM

Onur Mutlu,

"Memory-Centric Computing Systems"

Invited Tutorial at <u>66th International Electron Devices</u>

Meeting (IEDM), Virtual, 12 December 2020.

[Slides (pptx) (pdf)]

[Executive Summary Slides (pptx) (pdf)]

[Tutorial Video (1 hour 51 minutes)]

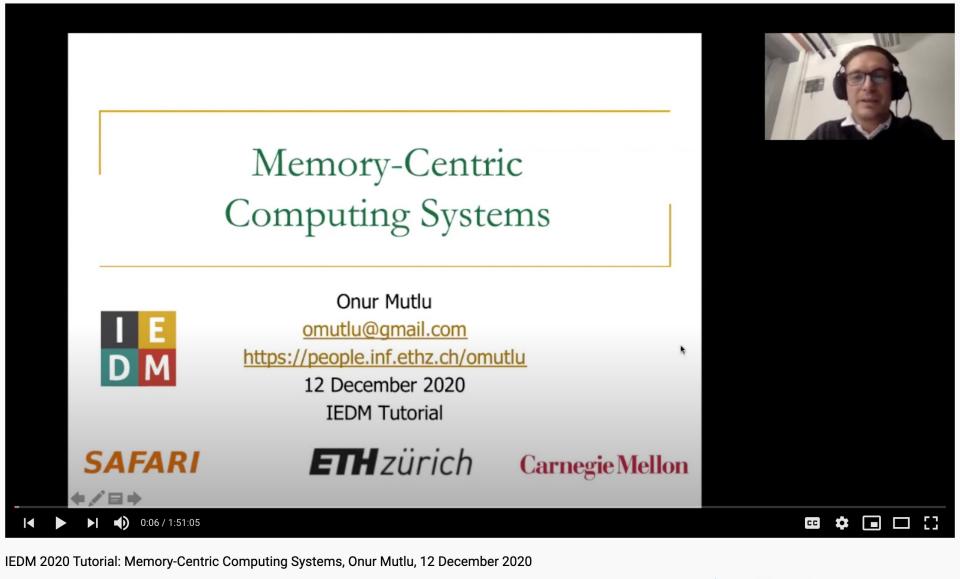
[Executive Summary Video (2 minutes)]

[Abstract and Bio]

[Related Keynote Paper from VLSI-DAT 2020]

[Related Review Paper on Processing in Memory]

https://www.youtube.com/watch?v=H3sEaINPBOE



1,641 views • Dec 23, 2020 ♣ SHARE =+ SAVE •



ANALYTICS

EDIT VIDEO

Challenge and Opportunity for Future

Computing Architectures with Minimal Data Movement

Challenge and Opportunity for Future

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

Challenge and Opportunity for Future

Fundamentally High-Performance (Data-Centric) Computing Architectures

Four Key Issues in Future Platforms

Fundamentally Secure/Reliable/Safe Architectures

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

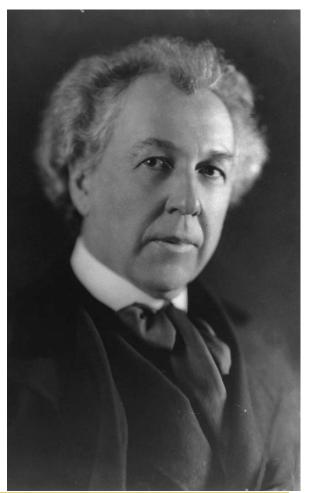
Fundamentally Low-Latency and Predictable Architectures

Architectures for AI/ML, Genomics, Medicine, Health

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"

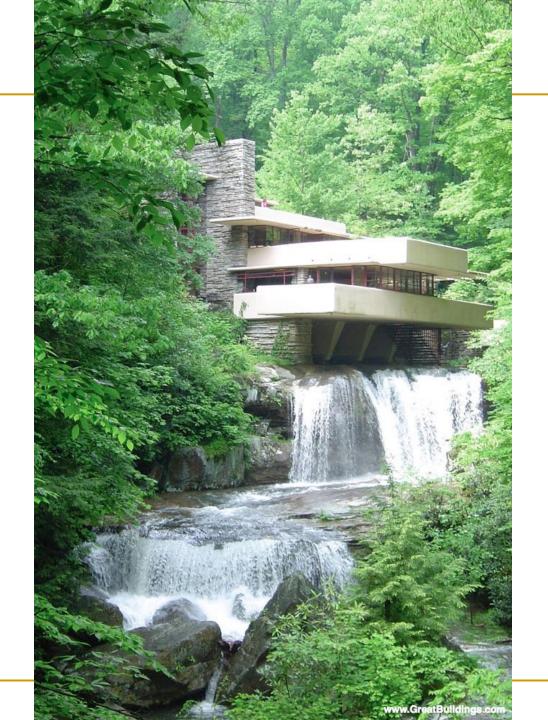


Principled Design

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



308



Another Example: Precedent-Based Design



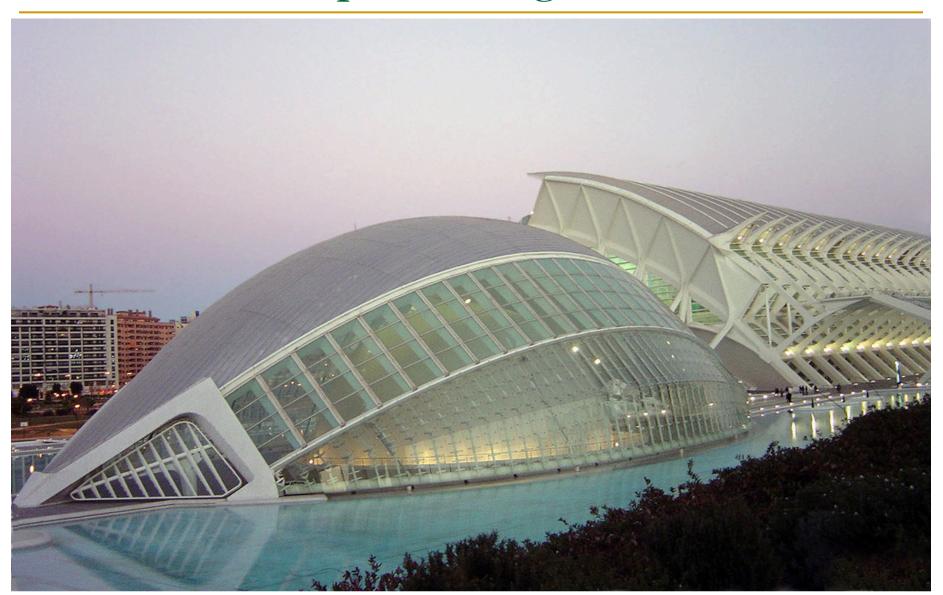
Principled Design



Another Principled Design



Another Principled Design



Principle Applied to Another Structure





314

Source: By 準建築人手札網站 Forgemind ArchiMedia - Flickr: IMG_2489.JPG, CC BY 2.0, SOURCE: https://www.dezeen.gom/2016/2016/2016-09.jpd.cc BY 2.0, Source: By 準建築人手札網站 Forgemind ArchiMedia - Flickr: IMG_2489.JPG, CC BY 2.0, Source: By 準建築人手札網站 Forgemind ArchiMedia - Flickr: IMG_2489.JPG, CC BY 2.0, Source: By 準建築人手札網站 Forgemind ArchiMedia - Flickr: IMG_2489.JPG, CC BY 2.0, Source: By 準建築人手札網站 Forgemind ArchiMedia - Flickr: IMG_2489.JPG, CC BY 2.0, Source: By 準建築人手札網站 Forgemind ArchiMedia - Flickr: IMG_2489.JPG, CC BY 2.0, Source: By 準建築人手札網站 Forgemind ArchiMedia - Flickr: IMG_2489.JPG, CC BY 2.0, Source: B

Overarching Principles for Computing?



Fundamentally Better Architectures

Data-centric

Data-driven

Data-aware

A Blueprint for Fundamentally Better Architectures

Onur Mutlu,

"Intelligent Architectures for Intelligent Computing Systems"

Invited Paper in Proceedings of the <u>Design, Automation, and Test in</u> <u>Europe Conference</u> (**DATE**), Virtual, February 2021.

[Slides (pptx) (pdf)]

[IEDM Tutorial Slides (pptx) (pdf)]

[Short DATE Talk Video (11 minutes)]

[Longer IEDM Tutorial Video (1 hr 51 minutes)]

Intelligent Architectures for Intelligent Computing Systems

Onur Mutlu ETH Zurich omutlu@gmail.com

We Need to Exploit Good Principles

- Data-centric design
- All components intelligent
- Good cross-layer communication, expressive interfaces
- Better-than-worst-case design
- Heterogeneity
- Flexibility, adaptability

Open minds

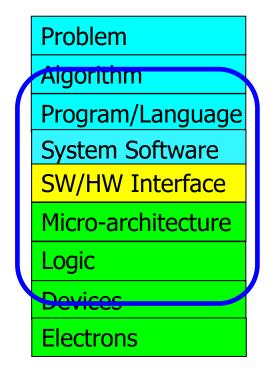
Concluding Remarks

- It is time to design principled computing architectures to achieve the highest security, performance, and efficiency
- Discover design principles for fundamentally secure and reliable computer architectures
- Design complete systems to be balanced and energy-efficient,
 i.e., data-centric (or memory-centric) and low-latency
- Enable new platforms for genomics, medicine, health, AI/ML
- This can
 - Lead to orders-of-magnitude improvements
 - Enable new applications & computing platforms
 - Enable better understanding of nature

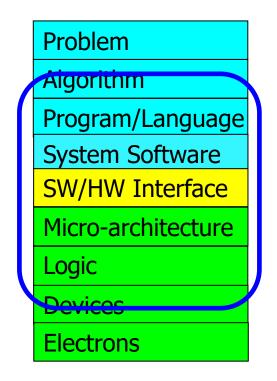
319

The Future is Very Bright

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



We Need to Think and Act Across the Stack



We can get there step by step

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>

Looking Beyond Moore and Von Neumann, Springer, to be published in 2021.

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

Abstract

Modern computing systems are overwhelmingly designed to move data to computation. This design choice goes directly against at least three key trends in computing that cause performance, scalability and energy bottlenecks: (1) data access is a key bottleneck as many important applications are increasingly data-intensive, and memory bandwidth and energy do not scale well, (2) energy consumption is a key limiter in almost all computing platforms, especially server and mobile systems, (3) data movement, especially off-chip to on-chip, is very expensive in terms of bandwidth, energy and latency, much more so than computation. These trends are especially severely-felt in the data-intensive server and energy-constrained mobile systems of today.

At the same time, conventional memory technology is facing many technology scaling challenges in terms of reliability, energy, and performance. As a result, memory system architects are open to organizing memory in different ways and making it more intelligent, at the expense of higher cost. The emergence of 3D-stacked memory plus logic, the adoption of error correcting codes inside the latest DRAM chips, proliferation of different main memory standards and chips, specialized for different purposes (e.g., graphics, low-power, high bandwidth, low latency), and the necessity of designing new solutions to serious reliability and security issues, such as the RowHammer phenomenon, are an evidence of this trend.

This chapter discusses recent research that aims to practically enable computation close to data, an approach we call processing-in-memory (PIM). PIM places computation mechanisms in or near where the data is stored (i.e., inside the memory chips, in the logic layer of 3D-stacked memory, or in the memory controllers), so that data movement between the computation units and memory is reduced or eliminated. While the general idea of PIM is not new, we discuss motivating trends in applications as well as memory circuits/technology that greatly exacerbate the need for enabling it in modern computing systems. We examine at least two promising new approaches to designing PIM systems to accelerate important data-intensive applications: (1) processing using memory by exploiting analog operational properties of DRAM chips to perform massively-parallel operations in memory, with low-cost changes, (2) processing near memory by exploiting 3D-stacked memory technology design to provide high memory bandwidth and low memory latency to in-memory logic. In both approaches, we describe and tackle relevant cross-layer research, design, and adoption challenges in devices, architecture, systems, and programming models. Our focus is on the development of in-memory processing designs that can be adopted in real computing platforms at low cost. We conclude by discussing work on solving key challenges to the practical adoption of PIM.

Keywords: memory systems, data movement, main memory, processing-in-memory, near-data processing, computation-in-memory, processing using memory, processing near memory, 3D-stacked memory, non-volatile memory, energy efficiency, high-performance computing, computer architecture, computing paradigm, emerging technologies, memory scaling, technology scaling, dependable systems, robust systems, hardware security, system security, latency, low-latency computing

Contents

1	Introduction	2
2	Major Trends Affecting Main Memory	4
_		
3	The Need for Intelligent Memory Controllers	_
	to Enhance Memory Scaling	6
4	Perils of Processor-Centric Design	9
5	Processing-in-Memory (PIM): Technology	
	Enablers and Two Approaches	12
	5.1 New Technology Enablers: 3D-Stacked	
	Memory and Non-Volatile Memory	12
_	5.2 Two Approaches: Processing Using	
	Memory (PUM) vs. Processing Near	12
_	Memory (PNM)	13
6	Processing Using Memory (PUM)	14
O	6.1 RowClone	14
	6.2 Ambit	15
		17 17
	6.4 In-DRAM Security Primitives	1/
7	Processing Near Memory (PNM)	18
	7.1 Tesseract: Coarse-Grained Application-	
	Level PNM Acceleration of Graph Pro-	
	cessing	19
	7.2 Function-Level PNM Acceleration of	
	Mobile Consumer Workloads	20
	7.3 Programmer-Transparent Function-	
	Level PNM Acceleration of GPU	
	Applications	21
	7.4 Instruction-Level PNM Acceleration	
	with PIM-Enabled Instructions (PEI)	21
	7.5 Function-Level PNM Acceleration of	
	Genome Analysis Workloads	22
	7.6 Application-Level PNM Acceleration of	
	Time Series Analysis	23
	E III di Ali di ADITA	•
8	Enabling the Adoption of PIM	24
	8.1 Programming Models and Code Genera-	24
	tion for PIM	24
	8.2 PIM Runtime: Scheduling and Data	25
_	Mapping	27
		27
	8.5 Data Structures for PIM	28
	8.6 Benchmarks and Simulation Infrastruc-	29
	P. 7 Paul DIM Handware Systems and Prote	29
	8.7 Real PIM Hardware Systems and Proto-	30
	types	30
	6.6 Security Considerations	30

Conclusion and Future Outlook

1. Introduction

31

Main memory, built using the Dynamic Random Access Memory (DRAM) technology, is a major component in nearly all computing systems, including servers, cloud platforms, mobile/embedded devices, and sensor systems. Across all of these systems, the data working set sizes of modern applications are rapidly growing, while the need for fast analysis of such data is increasing. Thus, main memory is becoming an increasingly significant bottleneck across a wide variety of computing systems and applications [1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16]. Alleviating the main memory bottleneck requires the memory capacity, energy, cost, and performance to all scale in an efficient manner across technology generations. Unfortunately, it has become increasingly difficult in recent years, especially the past decade, to scale all of these dimensions [1, 2, 17, 18, 19, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29, 30, 31, 32, 33, 34, 35, 36, 37, 38, 39, 40, 41, 42, 43, 44, 45, 46, 47, 48, 49], and thus the main memory bottleneck has been worsening.

A major reason for the main memory bottleneck is the high energy and latency cost associated with data movement. In modern computers, to perform any operation on data that resides in main memory, the processor must retrieve the data from main memory. This requires the memory controller to issue commands to a DRAM module across a relatively slow and power-hungry off-chip bus (known as the memory channel). The DRAM module sends the requested data across the memory channel, after which the data is placed in the caches and registers. The CPU can perform computation on the data once the data is in its registers. Data movement from the DRAM to the CPU incurs long latency and consumes a significant amount of energy [7, 50, 51, 52, 53, 54]. These costs are often exacerbated by the fact that much of the data brought into the caches is not reused by the CPU [52, 53, 55, 56], providing little benefit in return for the high latency and energy cost.

The cost of data movement is a fundamental issue with the processor-centric nature of contemporary computer systems. The CPU is considered to be the master in the system, and computation is performed only in the processor (and accelerators). In contrast, data storage and communication units, including the main memory, are treated as unintelligent workers that are incapable of computation. As a result of this processor-centric design paradigm, data moves a lot in the system between the computation units and communication/ storage units so that computation can be done on it. With the increasingly data-centric nature of contemporary and emerging appli-

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]

A Tutorial on Memory-Centric Systems

Onur Mutlu,

"Memory-Centric Computing Systems"

Invited Tutorial at <u>66th International Electron Devices</u>

Meeting (IEDM), Virtual, 12 December 2020.

[Slides (pptx) (pdf)]

[Executive Summary Slides (pptx) (pdf)]

[Tutorial Video (1 hour 51 minutes)]

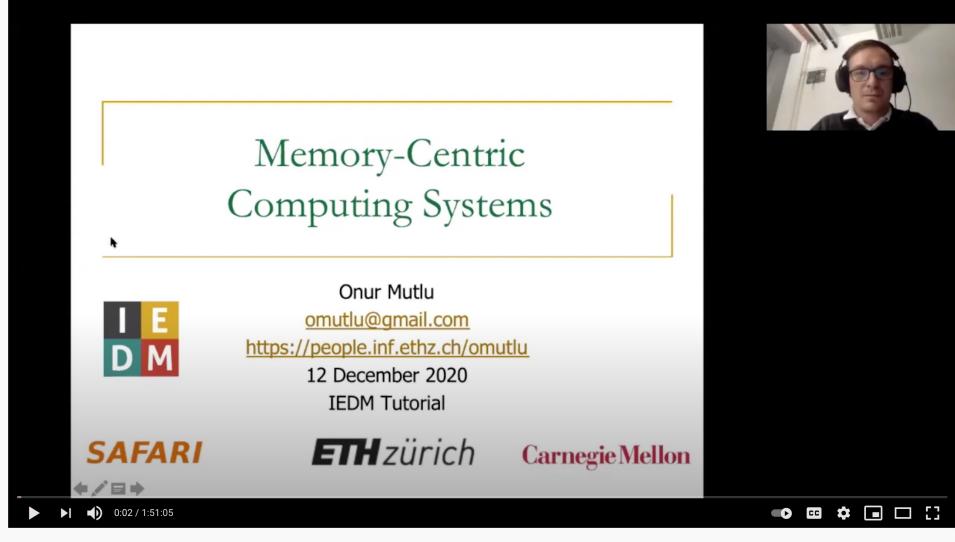
[Executive Summary Video (2 minutes)]

[Abstract and Bio]

[Related Keynote Paper from VLSI-DAT 2020]

[Related Review Paper on Processing in Memory]

https://www.youtube.com/watch?v=H3sEaINPBOE



IEDM 2020 Tutorial: Memory-Centric Computing Systems, Onur Mutlu, 12 December 2020



Onur Mutlu Lectures 15.2K subscribers

Speaker: Professor Onur Mutlu (https://people.inf.ethz.ch/omutlu/)

Date: December 12, 2020
Abstract and Bio: https://ieee-iedm.org/wp-content/uplo...

ANALYTICS

EDIT VIDEO

Funding Acknowledgments

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

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, ...

Acknowledgments



Think BIG, Aim HIGH!

https://safari.ethz.ch

Onur Mutlu's SAFARI Research Group

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

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



Think BIG, Aim HIGH!

SAFARI

https://safari.ethz.ch

SAFARI Newsletter April 2020 Edition

https://safari.ethz.ch/safari-newsletter-april-2020/





View in your browser

Think Big, Aim High



Dear SAFARI friends,

SAFARI Newsletter January 2021 Edition

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





Think Big, Aim High, and Have a Wonderful 2021! Newsletter January 2021



Dear SAFARI friends,

Future Computing Platforms Challenges and Opportunities

Onur Mutlu

omutlu@gmail.com

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

1 December 2021

IEEE Data & Storage Symposium (IEEE Bangalore)





Carnegie Mellon

Four Key Issues in Future Platforms

Fundamentally Secure/Reliable/Safe Architectures

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

Fundamentally Low-Latency and Predictable Architectures

Architectures for AI/ML, Genomics, Medicine, Health

Low Latency Communication is Critical

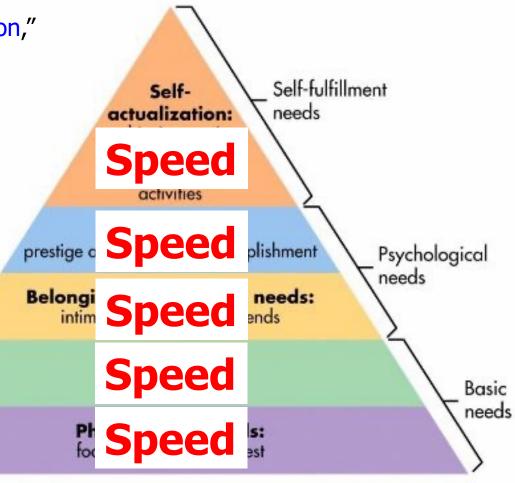


Maslow's Hierarchy of Needs, A Third Time

Maslow, "A Theory of Human Motivation," Psychological Review, 1943.

Maslow, "Motivation and Personality," Book, 1954-1970.

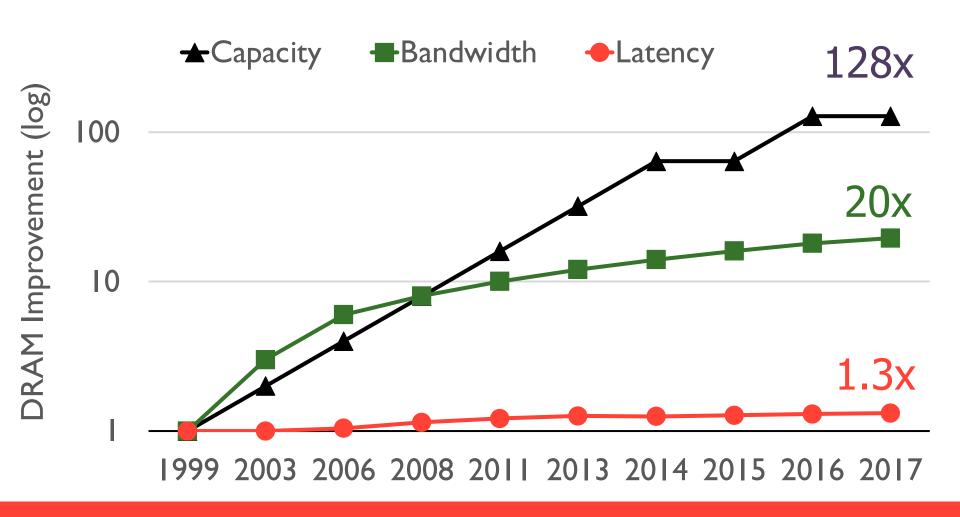




Challenge and Opportunity for Future

Fundamentally Low-Latency Computing Architectures

Main Memory Latency Lags Behind



Memory latency remains almost constant

The Memory Latency Problem

- High memory latency is a significant limiter of system performance and energy-efficiency
- It is becoming increasingly so with higher memory contention in multi-core and heterogeneous architectures
 - Exacerbating the bandwidth need
 - Exacerbating the QoS problem
- It increases processor design complexity due to the mechanisms incorporated to tolerate memory latency

Retrospective: Conventional Latency Tolerance Techniques

- Caching [initially by Wilkes, 1965]
 - Widely used, simple, effective, but inefficient, passive
 - Not all applications/phases exhibit temporal or spatial locality
- Prefetching Γinitially in IRM 360/91 19671

None of These Fundamentally Reduce Memory Latency

ongoing research effort

- Out-of-order execution [initially by Tomasulo, 1967]
 - Tolerates cache misses that cannot be prefetched
 - Requires extensive hardware resources for tolerating long latencies



Truly Reducing Memory Latency

Why the Long Memory Latency?

- Reason 1: Design of DRAM Micro-architecture
 - Goal: Maximize capacity/area, not minimize latency
- Reason 2: "One size fits all" approach to latency specification
 - Same latency parameters for all temperatures
 - Same latency parameters for all DRAM chips
 - Same latency parameters for all parts of a DRAM chip
 - Same latency parameters for all supply voltage levels
 - Same latency parameters for all application data
 - **...**

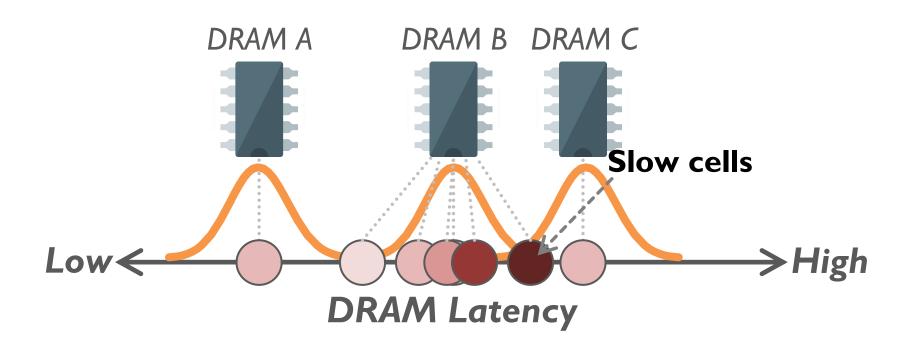
Tackling the Fixed Latency Mindset

- Reliable operation latency is actually very heterogeneous
 - Across temperatures, chips, parts of a chip, voltage levels, ...
- Idea: Dynamically find out and use the lowest latency one can reliably access a memory location with
 - Adaptive-Latency DRAM [HPCA 2015]
 - Flexible-Latency DRAM [SIGMETRICS 2016]
 - Design-Induced Variation-Aware DRAM [SIGMETRICS 2017]
 - Voltron [SIGMETRICS 2017]
 - DRAM Latency PUF [HPCA 2018]
 - DRAM Latency True Random Number Generator [HPCA 2019]
 - **-** ...
- We would like to find sources of latency heterogeneity and exploit them to minimize latency

344

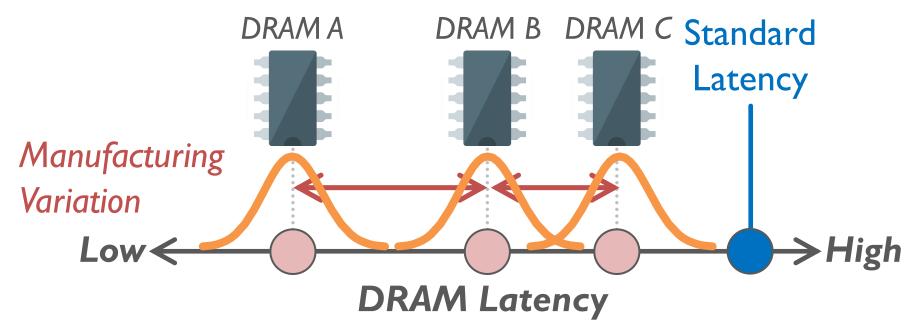
Latency Variation in Memory Chips

Heterogeneous manufacturing & operating conditions → latency variation in timing parameters



Why is Latency High?

- DRAM latency: Delay as specified in DRAM standards
 - Doesn't reflect true DRAM device latency
- Imperfect manufacturing process → latency variation
- High standard latency chosen to increase yield



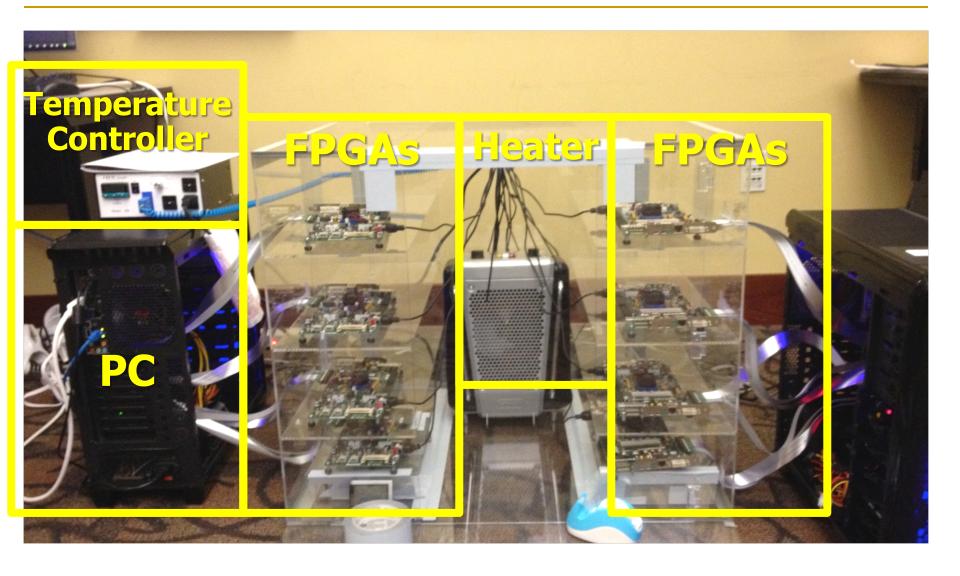
What Causes the Long Memory Latency?

Conservative timing margins!

- DRAM timing parameters are set to cover the worst case
- Worst-case temperatures
 - 85 degrees vs. common-case
 - to enable a wide range of operating conditions
- Worst-case devices
 - DRAM cell with smallest charge across any acceptable device
 - to tolerate process variation at acceptable yield
- This leads to large timing margins for the common case

Understanding and Exploiting Variation in DRAM Latency

DRAM Characterization Infrastructure



Adaptive-Latency DRAM

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



Latency Reduction Summary of 115 DIMMs

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

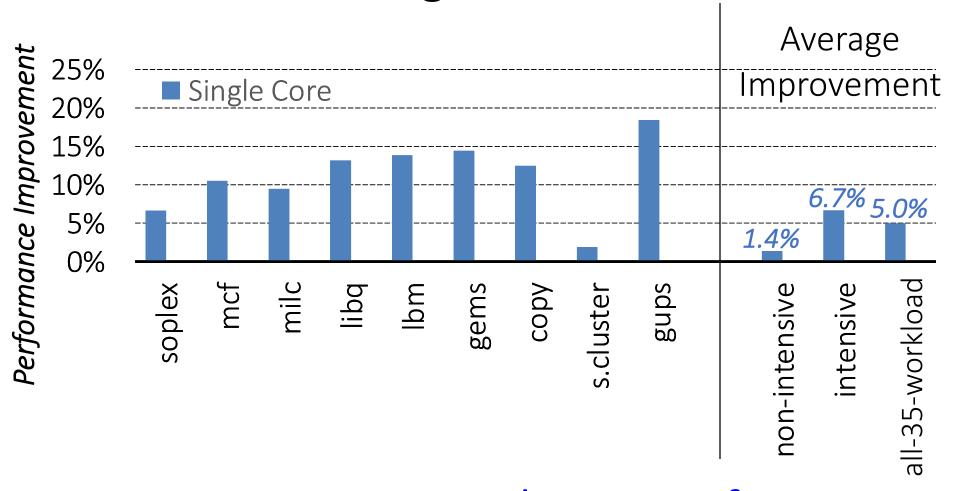


AL-DRAM: Real System Evaluation

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

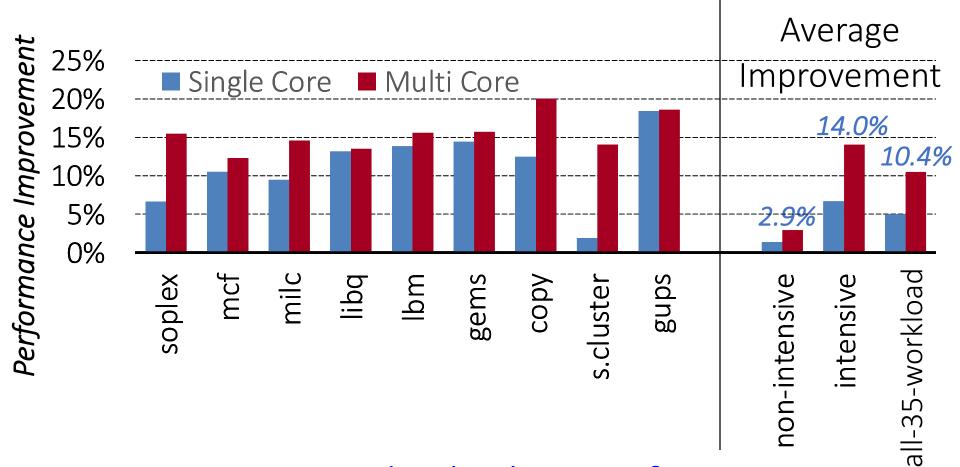
D18F2x200 dct[0] mp[1:0] DDR3 DRAM Timing 0 Reset: 0F05 0505h. See 2.9.3 [DCT Configuration Registers]. Description Bits Reserved 31:30 29:24 Tras: row active strobe. Read-write. BIOS: See 2.9.7.5 [SPD ROM-Based Configuration]. Specifies the minimum time in memory clock cycles from an activate command to a precharge command, both to the same chip select bank. Description Bits 07h-00h Reserved 2Ah-08h <Tras> clocks 3Fh-2Bh Reserved 23:21 Reserved 20:16 Trp: row precharge time. Read-write. BIOS: See 2.9.7.5 [SPD ROM-Based Configuration]. Specifies the minimum time in memory clock cycles from a precharge command to an activate command or auto refresh command, both to the same bank.

AL-DRAM: Single-Core Evaluation



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

AL-DRAM: Multi-Core Evaluation



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

Reducing Latency Also Reduces Energy

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

More on Adaptive-Latency DRAM

 Donghyuk Lee, Yoongu Kim, Gennady Pekhimenko, Samira Khan, Vivek Seshadri, Kevin Chang, and Onur Mutlu,
 "Adaptive-Latency DRAM: Optimizing DRAM Timing for the Common-Case"

Proceedings of the <u>21st International Symposium on High-</u> <u>Performance Computer Architecture</u> (**HPCA**), Bay Area, CA, February 2015.

[Slides (pptx) (pdf)] [Full data sets]

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

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

356

CLR-DRAM: Capacity-Latency Reconfigurability

Haocong Luo, Taha Shahroodi, Hasan Hassan, Minesh Patel, A. Giray Yaglikci, Lois Orosa, Jisung Park, and Onur Mutlu,
 "CLR-DRAM: A Low-Cost DRAM Architecture Enabling Dynamic Capacity-Latency Trade-Off"

Proceedings of the <u>47th International Symposium on Computer</u> <u>Architecture</u> (**ISCA**), Valencia, Spain, June 2020.

[Slides (pptx) (pdf)]

[Lightning Talk Slides (pptx) (pdf)]

[Talk Video (20 minutes)]

[Lightning Talk Video (3 minutes)]

CLR-DRAM: A Low-Cost DRAM Architecture Enabling Dynamic Capacity-Latency Trade-Off

Haocong Luo^{§†} Taha Shahroodi[§] Hasan Hassan[§] Minesh Patel[§] A. Giray Yağlıkçı[§] Lois Orosa[§] Jisung Park[§] Onur Mutlu[§]

§ETH Zürich †ShanghaiTech University



Analysis of Latency Variation in DRAM Chips

Kevin Chang, Abhijith Kashyap, Hasan Hassan, Samira Khan, Kevin Hsieh, Donghyuk Lee, Saugata Ghose, Gennady Pekhimenko, Tianshi Li, and Onur Mutlu,

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

Proceedings of the <u>ACM International Conference on Measurement and</u> Modeling of Computer Systems (SIGMETRICS), Antibes Juan-Les-Pins, France, June 2016.

[Slides (pptx) (pdf)]

Source Code

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

Kevin K. Chang¹ Abhijith Kashyap¹ Hasan Hassan^{1,2} Saugata Ghose¹ Kevin Hsieh¹ Donghyuk Lee¹ Tianshi Li^{1,3} Gennady Pekhimenko¹ Samira Khan⁴ Onur Mutlu^{5,1}

¹Carnegie Mellon University ²TOBB ETÜ ³Peking University ⁴University of Virginia ⁵ETH Zürich SAFARI

Design-Induced Latency Variation in DRAM

 Donghyuk Lee, Samira Khan, Lavanya Subramanian, Saugata Ghose, Rachata Ausavarungnirun, Gennady Pekhimenko, Vivek Seshadri, and Onur Mutlu,

"Design-Induced Latency Variation in Modern DRAM Chips:
Characterization, Analysis, and Latency Reduction Mechanisms"
Proceedings of the ACM International Conference on Measurement and
Modeling of Computer Systems (SIGMETRICS), Urbana-Champaign, IL,
USA, June 2017.

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

Donghyuk Lee, NVIDIA and Carnegie Mellon University
Samira Khan, University of Virginia
Lavanya Subramanian, Saugata Ghose, Rachata Ausavarungnirun, Carnegie Mellon University
Gennady Pekhimenko, Vivek Seshadri, Microsoft Research
Onur Mutlu, ETH Zürich and Carnegie Mellon University

Solar-DRAM: Putting It Together

Jeremie S. Kim, Minesh Patel, Hasan Hassan, and Onur Mutlu,
 "Solar-DRAM: Reducing DRAM Access Latency by
 Exploiting the Variation in Local Bitlines"
 Proceedings of the 36th IEEE International Conference on Computer Design (ICCD), Orlando, FL, USA, October 2018.
 [Slides (pptx) (pdf)]
 [Talk Video (16 minutes)]

Solar-DRAM: Reducing DRAM Access Latency by Exploiting the Variation in Local Bitlines

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

360

Challenge and Opportunity for Future

Fundamentally Low-Latency Computing Architectures

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

<u>Jeremie S. Kim</u> Minesh Patel Hasan Hassan Lois Orosa Onur Mutlu

SAFARI



Carnegie Mellon

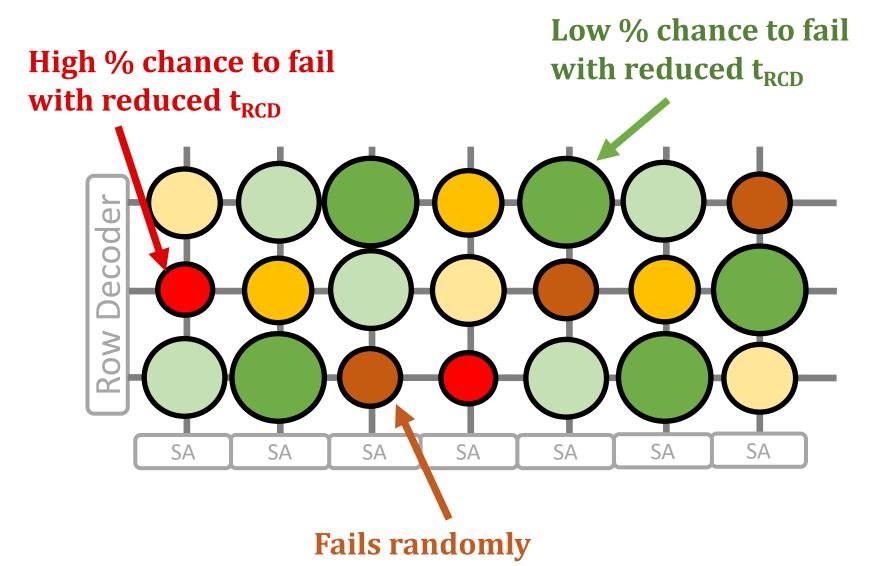
DRAM Latency Characterization of 282 LPDDR4 DRAM Devices

 Latency failures come from accessing DRAM with reduced timing parameters.

Key Observations:

- 1. A cell's **latency failure** probability is determined by **random process variation**
- 2. Some cells fail **randomly**

D-RaNGe Key Idea



with reduced t_{RCD}

SAFARI

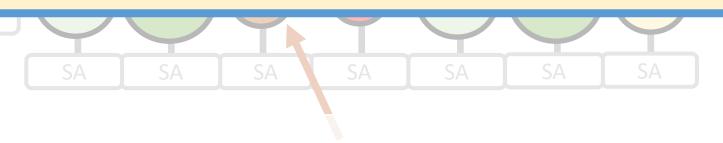
364/4

D-RaNGe Key Idea

High % chance to fail with reduced t_{RCD}

Low % chance to fail with reduced t_{RCD}

We refer to cells that fail randomly when accessed with a reduced t_{RCD} as RNG cells



Fails randomly with reduced t_{rcp}

Our D-RaNGe Evaluation

- We generate random values by repeatedly accessing RNG cells and aggregating the data read
- The random data satisfies the NIST statistical test suite for randomness
- The D-RaNGE generates random numbers
 - **Throughput**: 717.4 Mb/s
 - **Latency**: 64 bits in <1us
 - **Power**: 4.4 nJ/bit

DRAM Latency True Random Number Generator

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> <u>Architecture</u> (**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

SAFARI 367

DRAM Latency 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 Devices"

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

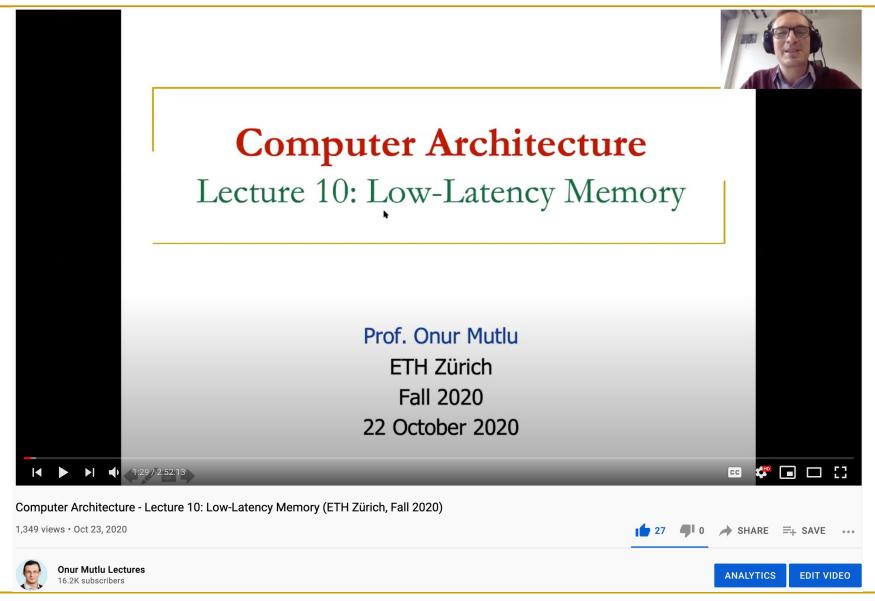
Jeremie S. Kim^{†§} Minesh Patel[§] Hasan Hassan[§] Onur Mutlu^{§†}

[†]Carnegie Mellon University [§]ETH Zürich

Lectures on Low-Latency Memory

- Computer Architecture, Fall 2020, Lecture 10
 - Low-Latency Memory (ETH Zürich, Fall 2020)
 - https://www.youtube.com/watch?v=vQd1YgOH1Mw&list=PL5Q2soXY2Zi9xidyIgBx Uz7xRPS-wisBN&index=19
- Computer Architecture, Fall 2020, Lecture 12b
 - Capacity-Latency Reconfigurable DRAM (ETH Zürich, Fall 2020)
 - https://www.youtube.com/watch?v=DUtPFW3jxq4&list=PL5Q2soXY2Zi9xidyIgBxUz 7xRPS-wisBN&index=23
- Computer Architecture, Fall 2019, Lecture 11a
 - DRAM Latency PUF (ETH Zürich, Fall 2019)
 - https://www.youtube.com/watch?v=7gqnrTZpjxE&list=PL5Q2soXY2Zi-DyoI3HbqcdtUm9YWRR_z-&index=15
- Computer Architecture, Fall 2019, Lecture 11b
 - DRAM True Random Number Generator (ETH Zürich, Fall 2020)
 - https://www.youtube.com/watch?v=Y3hPv1I5f8Y&list=PL5Q2soXY2Zi-DyoI3HbqcdtUm9YWRR_z-&index=16

A Tutorial on Low-Latency Memory



https://www.youtube.com/onurmutlulectures

Challenge and Opportunity for Future

Fundamentally Low-Latency Computing Architectures

Four Key Issues in Future Platforms

Fundamentally Secure/Reliable/Safe Architectures

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

Fundamentally Low-Latency and Predictable Architectures

Architectures for AI/ML, Genomics, Medicine, Health

Intel Optane Persistent Memory (2019)

- Non-volatile main memory
- Based on 3D-XPoint Technology



PCM as Main Memory: Idea in 2009

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</u> (pdf)

Architecting Phase Change Memory as a Scalable DRAM Alternative

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

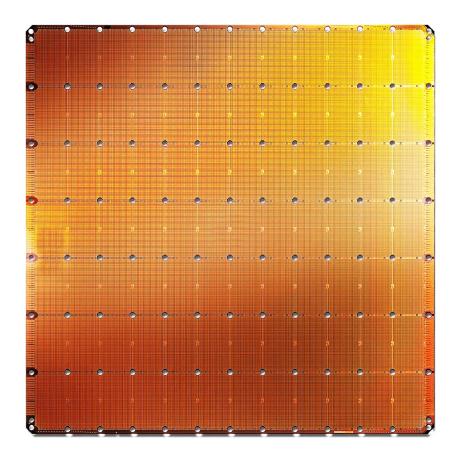
†Computer Architecture Group Microsoft Research Redmond, WA {blee, ipek, dburger}@microsoft.com ‡Computer Architecture Laboratory Carnegie Mellon University Pittsburgh, PA onur@cmu.edu

PCM as Main Memory: Idea in 2009

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

PHASE-CHANGE TECHNOLOGY AND THE FUTURE OF MAIN MEMORY

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

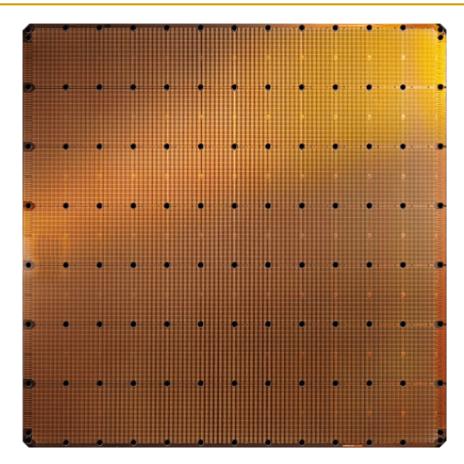
21.1 Billion transistors 815 mm²

NVIDIA TITAN V

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/

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

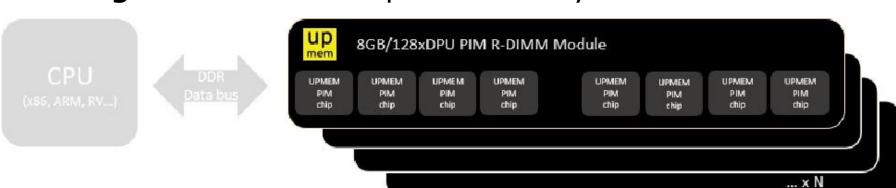
54.2 Billion transistors 826 mm²

NVIDIA Ampere GA100

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

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





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 3D-stacked 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 state-of-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

Samsung Newsroom

CORPORATE

PRODUCTS

PRESS RESOURCES

VIEWS

ABOUT US



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

Korea on February 17, 2021

Audio



Share (5



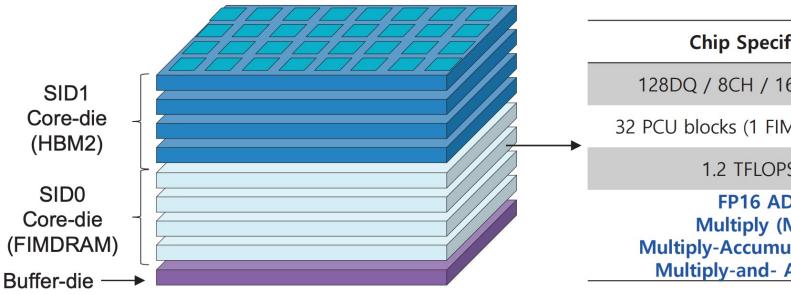


The new architecture will deliver over twice the system performance and reduce energy consumption by more than 70%

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 Al-driven workloads such as HPC, training and inference. We plan to build upon this breakthrough by further collaborating with Al solution providers for even more advanced PIM-powered applications."

FIMDRAM based on HBM2



[3D Chip Structure of HBM with FIMDRAM]

Chip Specification

128DQ / 8CH / 16 banks / BL4

32 PCU blocks (1 FIM block/2 banks)

1.2 TFLOPS (4H)

FP16 ADD / Multiply (MUL) / Multiply-Accumulate (MAC) / Multiply-and- Add (MAD)

ISSCC 2021 / SESSION 25 / DRAM / 25.4

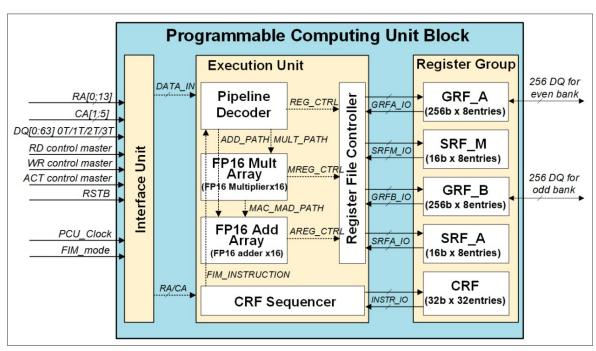
25.4 A 20nm 6GB Function-In-Memory DRAM, Based on HBM2 with a 1.2TFLOPS Programmable Computing Unit Using Bank-Level Parallelism, for Machine Learning Applications

Young-Cheon Kwon¹, Suk Han Lee¹, Jaehoon Lee¹, Sang-Hyuk Kwon¹, Je Min Ryu1, Jong-Pil Son1, Seongil O1, Hak-Soo Yu1, Haesuk Lee1, Soo Young Kim¹, Youngmin Cho¹, Jin Guk Kim¹, Jongyoon Choi¹, Hyun-Sung Shin¹, Jin Kim¹, BengSeng Phuah¹, HyoungMin Kim¹, Myeong Jun Song¹, Ahn Choi¹, Daeho Kim¹, SooYoung Kim¹, Eun-Bong Kim¹, David Wang², Shinhaeng Kang¹, Yuhwan Ro³, Seungwoo Seo³, JoonHo Song³, Jaeyoun Youn1, Kyomin Sohn1, Nam Sung Kim1

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 Let', Jaehoon Let', Sang-Hyuk Kwon', Ja Min Ryu', Jong-Pi Son', Seongli O', Hak Soo Yu', Hesauk Let', Soo Young Kim', Youngmin Cho', Jin Guk Kim', Jongyoon Choi', Hyun-Sung Shin', Jin Kim', BengSeng Phuah', HyoungMin Kim', Hyeong Jun Song', Alm Choi', Daeho Kim', Soo Young Kim', Eun-Bong Kim', David Wang', Shinhaend 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

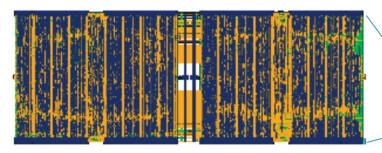
ISSCC 2021 / SESSION 25 / DRAM / 25.4

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-Cheon Kwon', Suk Han Let', Jaehoon Let', Sang-Hyuk Kwon', Ja Min Ryu', Jong-Pi Son', Seongli O', Hak Soo Yu', Hesauk Let', Soo Young Kim', Youngmin Cho', Jin Guk Kim', Jongyoon Choi', Hyun-Sung Shin', Jin Kim', BengSeng Phuah', HyoungMin Kim', Hyeong Jun Song', Alm Choi', Daeho Kim', Soo Young Kim', Eun-Bong Kim', David Wang', Shinhaend Kang', Yuhwan Ro', Seungwoo Seo', JoonHo Song', Jaeyoun Youn', Kyomin Sohn', Man Sung Kim'

Chip Implementation

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

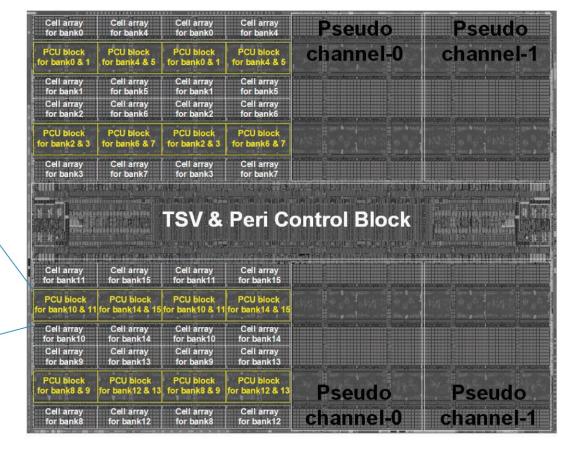


[Digital RTL design for PCU block]

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 Ler', Jaehoon Ler', 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', Hyeong Jun Song', Am Choi', Deach Kim', Soo'Qung Kim', Eun-Bong Kim', David Wang', Shinhaeng Kang', Yuhwan Ro', Seungwoo Seo', JoonHo Song', Jaeyoun Youn', Kyomin Sohn', Man Sung Kim'



More on Processing in Memory (I)

 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 Hassan^4 Amirali Boroumand^5 Jeremie Kim^{4,5} Michael A. Kozuch^3 Onur Mutlu^{4,5} Phillip B. Gibbons^5 Todd C. Mowry^5
```

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

More on Processing in Memory (II)

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

More on Processing in Memory (III)

Nastaran Hajinazar, Geraldo F. Oliveira, Sven Gregorio, Joao Dinis Ferreira, Nika Mansouri Ghiasi, Minesh Patel, Mohammed Alser, Saugata Ghose, Juan Gomez-Luna, and Onur Mutlu, "SIMDRAM: An End-to-End Framework for Bit-Serial SIMD Computing in DRAM" Proceedings of the 26th International Conference on Architectural Support for Programming Languages and Operating Systems (ASPLOS), Virtual, March-April 2021.

[2-page Extended Abstract]

[Short Talk Slides (pptx) (pdf)]

[Talk Slides (pptx) (pdf)]

[Short Talk Video (5 mins)]

[Full Talk Video (27 mins)]

SIMDRAM: A Framework for Bit-Serial SIMD Processing using DRAM

*Nastaran Hajinazar^{1,2} Nika Mansouri Ghiasi¹ *Geraldo F. Oliveira¹ Minesh Patel¹ Juan Gómez-Luna¹

Sven Gregorio¹ Mohammed Alser¹ Onur Mutlu¹ João Dinis Ferreira¹ Saugata Ghose³

¹ETH Zürich

²Simon Fraser University

³University of Illinois at Urbana–Champaign

More on Processing in Memory (IV)

 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

More on Processing in Memory (V)

 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 Boroumand¹ Saugata Ghose¹ Youngsok Kim² Rachata Ausavarungnirun¹ Eric Shiu³ Rahul Thakur³ Daehyun Kim^{4,3} Aki Kuusela³ Allan Knies³ Parthasarathy Ranganathan³ Onur Mutlu^{5,1}

More on Processing in Memory (VI)

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>
 Computer Architecture (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

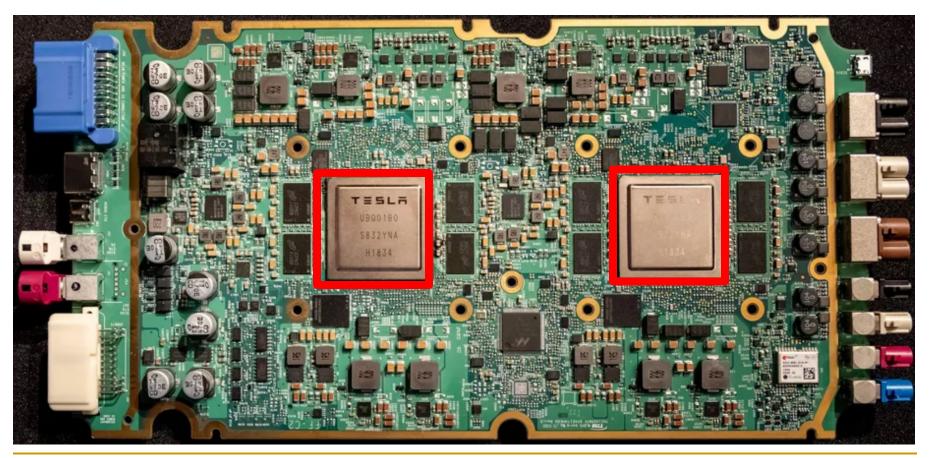
Seoul National University [†]Carnegie Mellon University

SAFARI

TESLA Full Self-Driving Computer (2019)

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





Google TPU Generation I (~2016)



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.

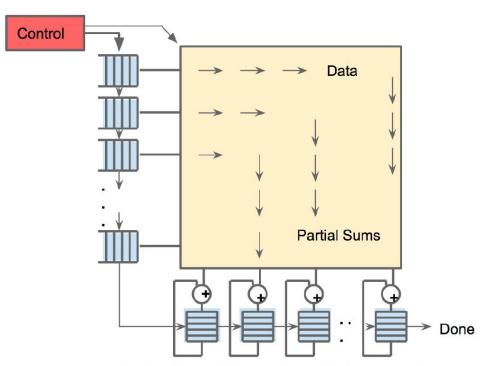


Figure 4. Systolic data flow of the Matrix Multiply Unit. Software has the illusion that each 256B input is read at once, and they instantly update one location of each of 256 accumulator RAMs.

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

Google TPU Generation II (2017)



https://www.nextplatform.com/2017/05/17/first-depth-look-googles-new-second-generation-tpu/

4 TPU chips vs 1 chip in TPU1

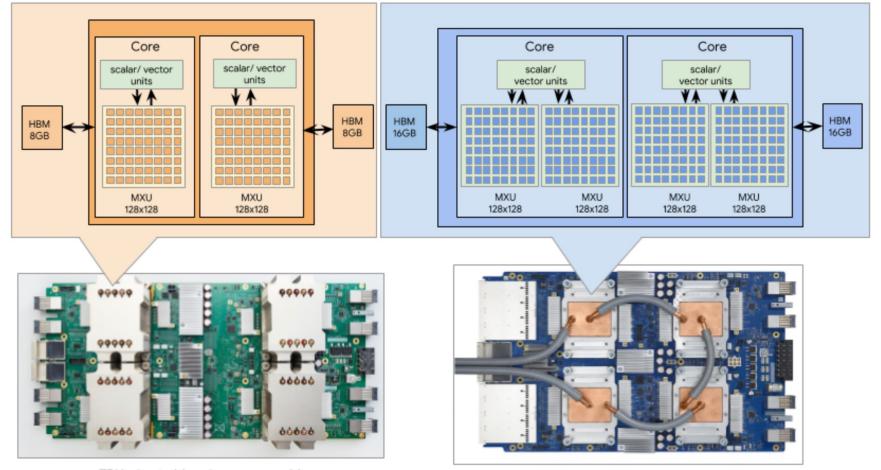
High Bandwidth Memory vs DDR3

Floating point operations vs FP16

45 TFLOPS per chip vs 23 TOPS

Designed for training and inference vs only inference

Google TPU Generation III



TPU v2 - 4 chips, 2 cores per chip

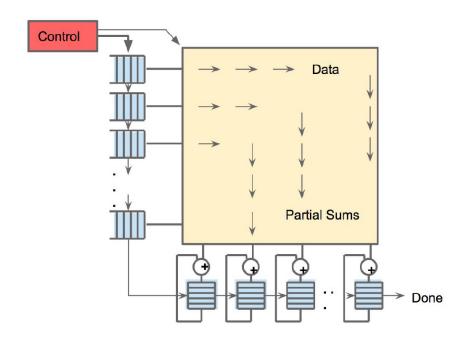
TPU v3 - 4 chips, 2 cores per chip

32GB HBM per chip vs 16GB HBM in TPU2

4 Matrix Units per chip 90 TFLOPS per chip vs 2 Matrix Units in TPU2 vs 45 TFLOPS in TPU2

An Example Modern Systolic Array: TPU (II)

As reading a large SRAM uses much more power than arithmetic, the matrix unit uses systolic execution to save energy by reducing reads and writes of the Unified Buffer [Kun80][Ram91][Ovt15b]. Figure 4 shows that data flows in from the left, and the weights are loaded from the top. A given 256-element multiply-accumulate operation moves through the matrix as a diagonal wavefront. The weights are preloaded, and take effect with the advancing wave alongside the first data of a new block. Control and data are pipelined to give the illusion that the 256 inputs are read at once, and that they instantly update one location of each of 256 accumulators. From a correctness perspective, software is unaware of the systolic nature of the matrix unit, but for performance, it does worry about the latency of the unit.



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

An Example Modern Systolic Array: TPU (III)

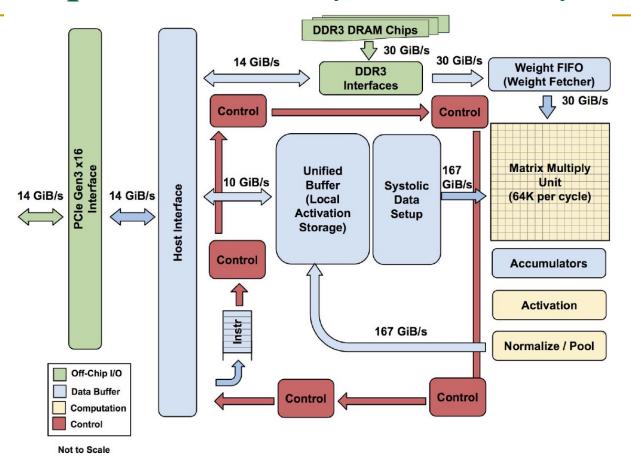
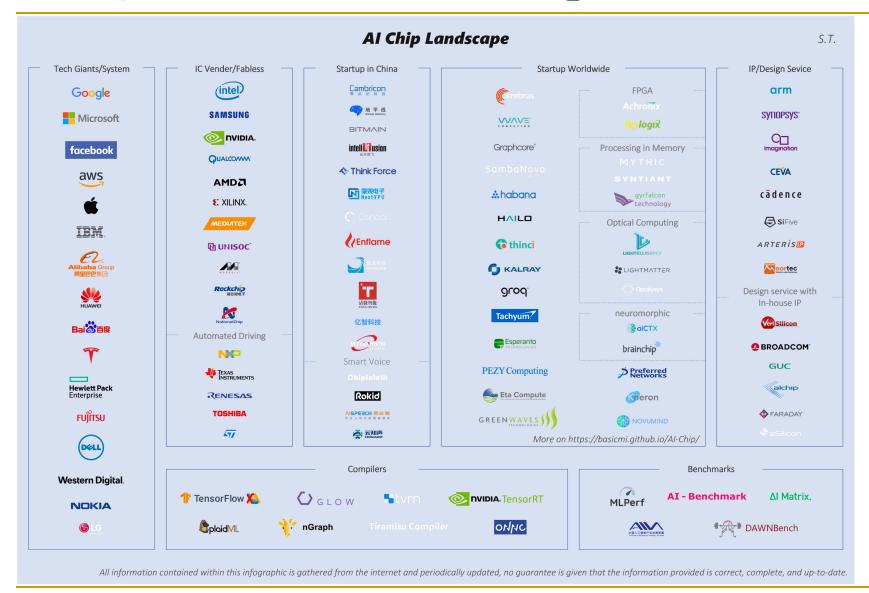


Figure 1. TPU Block Diagram. The main computation part is the yellow Matrix Multiply unit in the upper right hand corner. Its inputs are the blue Weight FIFO and the blue Unified Buffer (UB) and its output is the blue Accumulators (Acc). The yellow Activation Unit performs the nonlinear functions on the Acc, which go to the UB.

Many (Other) AI/ML Chips

- Alibaba
- Amazon
- Facebook
- Google
- Huawei
- Intel
- Microsoft
- NVIDIA
- Tesla
- Many Others and Many Startups...
- Many More to Come...

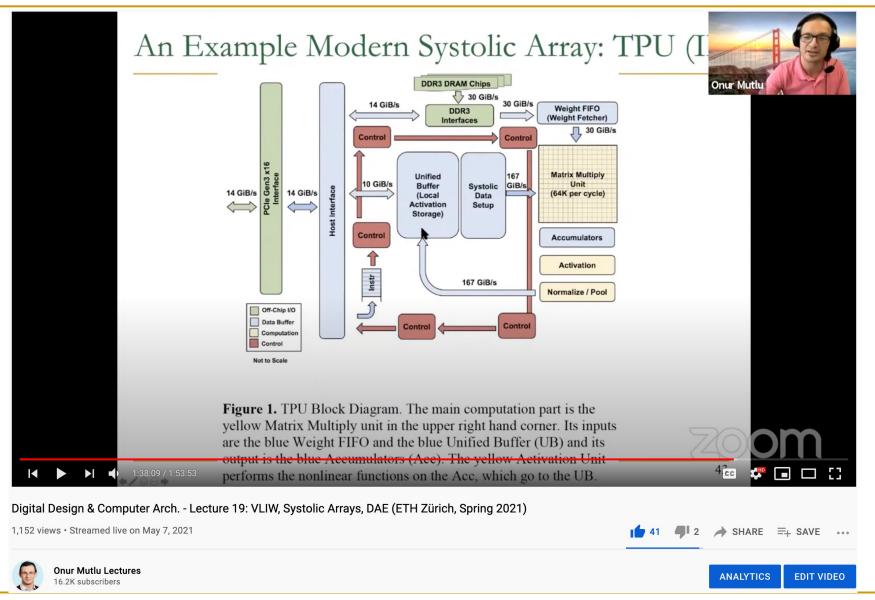
Many (Other) AI/ML Chips



Lectures on Systolic Arrays & ML/AI Acceleration

- Digital Design and Computer Architecture, Spring 2021, Lecture 19
 - VLIW, Systolic Arrays, DAE (ETH Zürich, Spring 2021)
 - https://www.youtube.com/watch?v=UtLy4Yagdys&list=PL5Q2soXY2Zi_uej3aY39YB 5pfW4SJ7LlN&index=21
- Computer Architecture, Fall 2020, Lecture 9b
 - □ **EDEN: Efficient DNN Inference** (ETH Zürich, Fall 2020)
 - https://www.youtube.com/watch?v=HmB32OXMKMY&list=PL5Q2soXY2Zi9xidyIgBx Uz7xRPS-wisBN&index=16
- Computer Architecture, Fall 2020, Lecture 9c
 - SMASH: Accelerating Sparse Matrix Operations (ETH Zürich, Fall 2020)
 - https://www.youtube.com/watch?v=78aikMbxkGc&list=PL5Q2soXY2Zi9xidyIgBxUz7 xRPS-wisBN&index=17

Lecture on Systolic Arrays & ML Acceleration



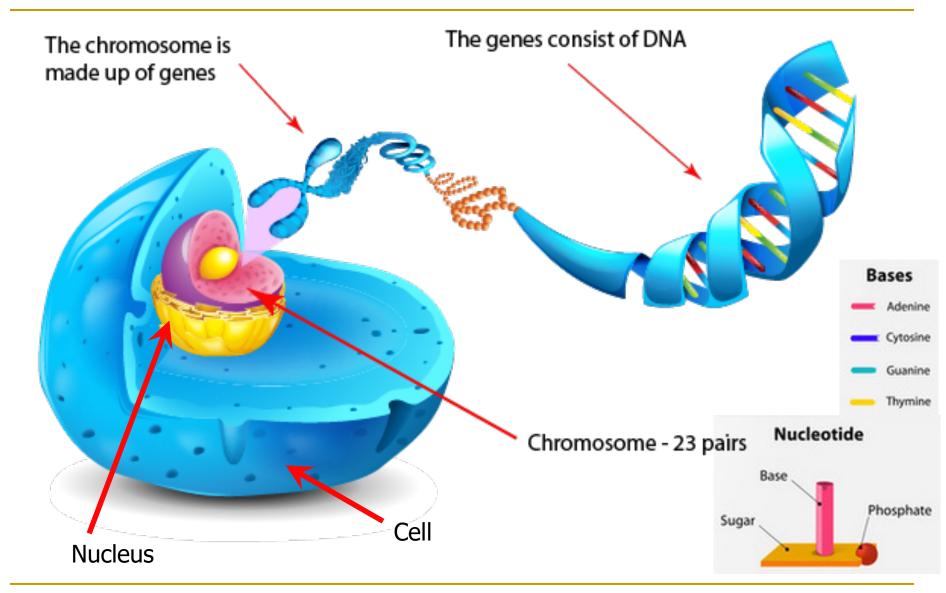
Accelerating Genome Analysis

Our Dream (circa 2007)

- An embedded device that can perform comprehensive genome analysis in real time (within a minute)
 - Which of these DNAs does this DNA segment match with?
 - What is the likely genetic disposition of this patient to this drug?
 - What disease/condition might this particular DNA/RNA piece associated with?

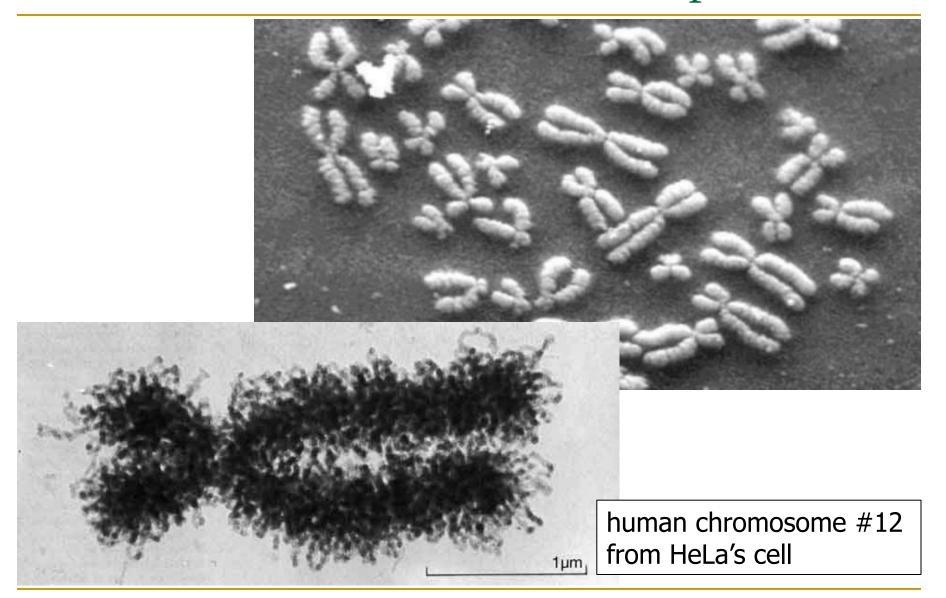
- . . .

What Is a Genome Made Of?



403

DNA Under Electron Microscope



DNA Sequencing

Goal:

Find the complete sequence of A, C, G, T's in DNA.

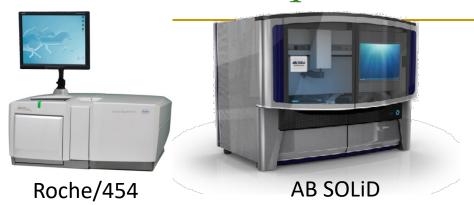
Challenge:

- There is no machine that takes long DNA as an input, and gives the complete sequence as output
- All sequencing machines chop DNA into pieces and identify relatively small pieces (but not how they fit together)

Untangling Yarn Balls & DNA Sequencing



Genome Sequencers





Illumina HiSeq2000



Pacific Biosciences RS



Ion Torrent Proton



Illumina MiSeq



Complete Genomics



Oxford Nanopore MinION



Illumina NovaSeq 6000

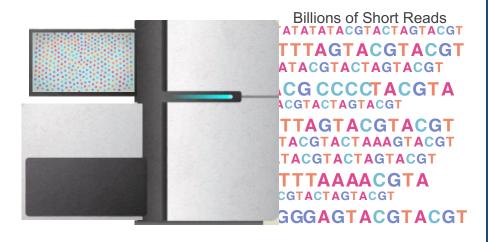


Oxford Nanopore GridION

... and more! All produce data with different properties.



Ion Torrent PGM



Short Read Alignment Reference Genome

Read Mapping

1 Sequencing

Genome Analysis

reference: TTTATCGCTTCCATGACGCAG

read1: ATCGCATCC read2: TATCGCATC

read3: CATCCATGA

read4: CGCTTCCAT

read5: CCATGACGC

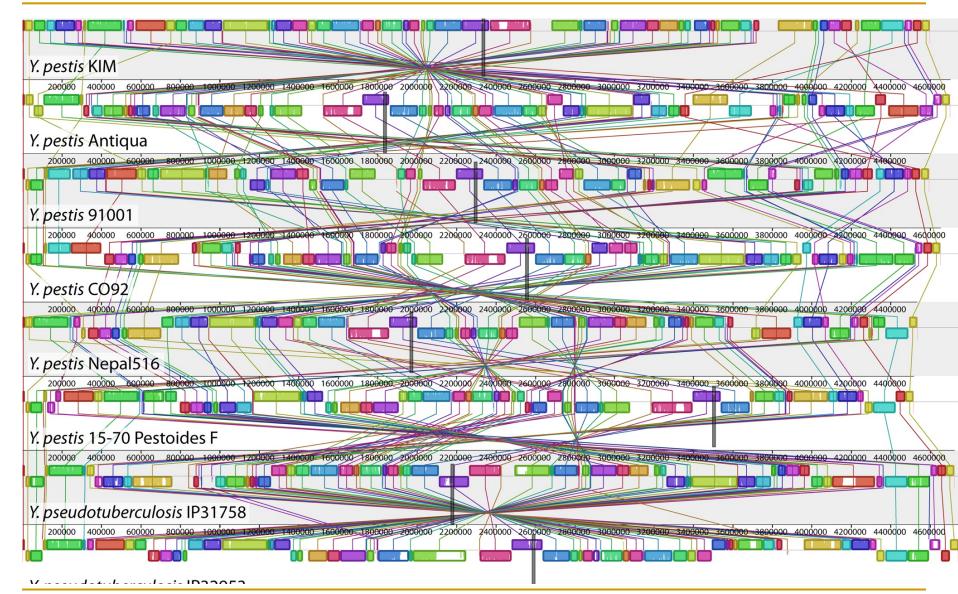
read6: TTCCATGAC



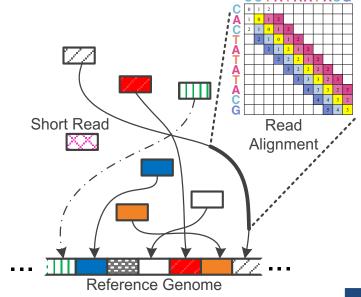
Variant Calling

Scientific Discovery 4

Genome Sequence Alignment: Example

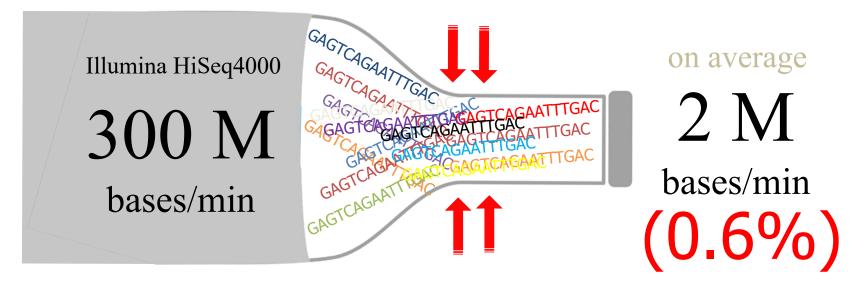






Read Mapping

Bottlenecked in Mapping!!



Hash Table Based Read Mappers

- + Guaranteed to find all mappings → sensitive
- + Can tolerate up to e errors

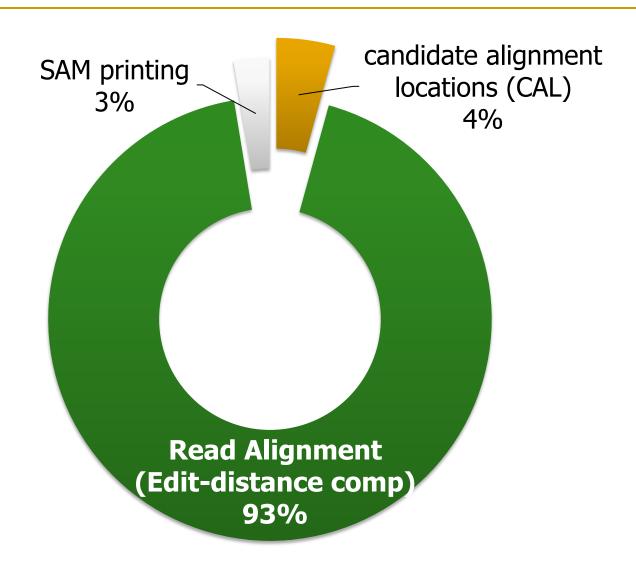
nature genetics

http://mrfast.sourceforge.net/

Personalized copy number and segmental duplication maps using next-generation sequencing

Can Alkan^{1,2}, Jeffrey M Kidd¹, Tomas Marques-Bonet^{1,3}, Gozde Aksay¹, Francesca Antonacci¹, Fereydoun Hormozdiari⁴, Jacob O Kitzman¹, Carl Baker¹, Maika Malig¹, Onur Mutlu⁵, S Cenk Sahinalp⁴, Richard A Gibbs⁶ & Evan E Eichler^{1,2}

Read Mapping Execution Time Breakdown



Filter fast before you align

Minimize costly "approximate string comparisons"

Our First Filter: Pure Software Approach

- Download the source code and try for yourself
 - Download link to FastHASH

Xin et al. BMC Genomics 2013, **14**(Suppl 1):S13 http://www.biomedcentral.com/1471-2164/14/S1/S13



PROCEEDINGS

Open Access

Accelerating read mapping with FastHASH

Hongyi Xin¹, Donghyuk Lee¹, Farhad Hormozdiari², Samihan Yedkar¹, Onur Mutlu^{1*}, Can Alkan^{3*}

From The Eleventh Asia Pacific Bioinformatics Conference (APBC 2013) Vancouver, Canada. 21-24 January 2013

Shifted Hamming Distance: SIMD Acceleration

https://github.com/CMU-SAFARI/Shifted-Hamming-Distance

Bioinformatics, 31(10), 2015, 1553-1560

doi: 10.1093/bioinformatics/btu856

Advance Access Publication Date: 10 January 2015

Original Paper



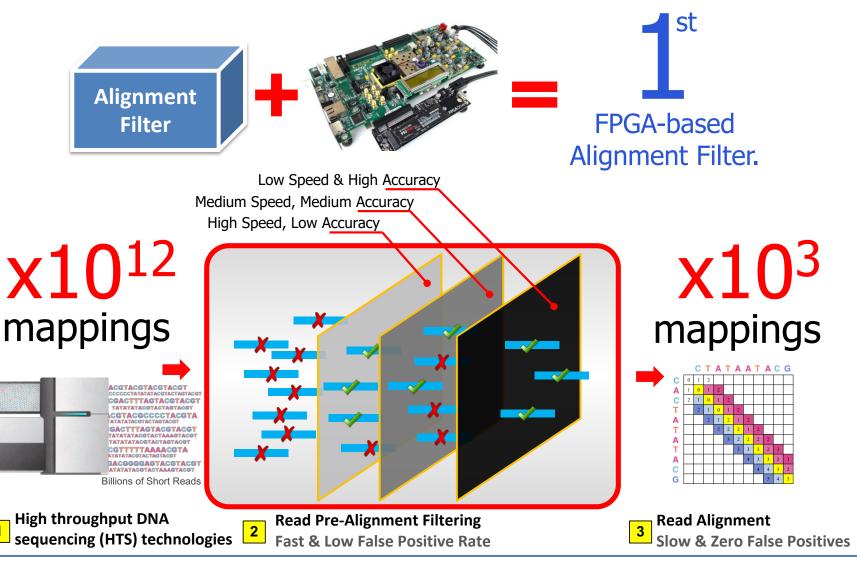
Sequence analysis

Shifted Hamming distance: a fast and accurate SIMD-friendly filter to accelerate alignment verification in read mapping

Hongyi Xin^{1,*}, John Greth², John Emmons², Gennady Pekhimenko¹, Carl Kingsford³, Can Alkan^{4,*} and Onur Mutlu^{2,*}

Xin+, "Shifted Hamming Distance: A Fast and Accurate SIMD-friendly Filter to Accelerate Alignment Verification in Read Mapping", Bioinformatics 2015.

GateKeeper: FPGA-Based Alignment Filtering



GateKeeper: FPGA-Based Alignment Filtering

 Mohammed Alser, Hasan Hassan, Hongyi Xin, Oguz Ergin, Onur Mutlu, and Can Alkan

"GateKeeper: A New Hardware Architecture for Accelerating Pre-Alignment in DNA Short Read Mapping" Bioinformatics, [published online, May 31], 2017.

[Source Code]

[Online link at Bioinformatics Journal]

GateKeeper: a new hardware architecture for accelerating pre-alignment in DNA short read mapping

Mohammed Alser ™, Hasan Hassan, Hongyi Xin, Oğuz Ergin, Onur Mutlu ™, Can Alkan ™

Bioinformatics, Volume 33, Issue 21, 1 November 2017, Pages 3355–3363,

https://doi.org/10.1093/bioinformatics/btx342

Published: 31 May 2017 Article history ▼

SAFARI

DNA Read Mapping & Filtering

- Problem: Heavily bottlenecked by Data Movement
- GateKeeper FPGA performance limited by DRAM bandwidth [Alser+, Bioinformatics 2017]
- Ditto for SHD on SIMD [Xin+, Bioinformatics 2015]
- Solution: Processing-in-memory can alleviate the bottleneck
- However, we need to design mapping & filtering algorithms to fit processing-in-memory

In-Memory DNA Sequence Analysis

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

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

BMC Genomics, 2018.

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

[Slides (pptx) (pdf)]

Source Code

[arxiv.org Version (pdf)]

Talk Video at AACBB 2019

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

Shouji (障子) [Alser+, Bioinformatics 2019]

Mohammed Alser, Hasan Hassan, Akash Kumar, Onur Mutlu, and Can Alkan, "Shouji: A Fast and Efficient Pre-Alignment Filter for Sequence Alignment" Bioinformatics, [published online, March 28], 2019.

Source Code

Online link at Bioinformatics Journal

Bioinformatics, 2019, 1–9 doi: 10.1093/bioinformatics/btz234 Advance Access Publication Date: 28 March 2019 Original Paper



Sequence alignment

Shouji: a fast and efficient pre-alignment filter for sequence alignment

Mohammed Alser^{1,2,3,*}, Hasan Hassan¹, Akash Kumar², Onur Mutlu^{1,3,*} and Can Alkan^{3,*}

¹Computer Science Department, ETH Zürich, Zürich 8092, Switzerland, ²Chair for Processor Design, Center For Advancing Electronics Dresden, Institute of Computer Engineering, Technische Universität Dresden, 01062 Dresden, Germany and ³Computer Engineering Department, Bilkent University, 06800 Ankara, Turkey

Associate Editor: Inanc Birol

SAFARI

^{*}To whom correspondence should be addressed.

SneakySnake [Alser+, Bioinformatics 2020]

Mohammed Alser, Taha Shahroodi, Juan-Gomez Luna, Can Alkan, and Onur Mutlu, "SneakySnake: A Fast and Accurate Universal Genome Pre-Alignment Filter for CPUs, GPUs, and FPGAs"

Bioinformatics, to appear in 2020.

[Source Code]

Online link at Bioinformatics Journal

Bioinformatics

doi.10.1093/bioinformatics/xxxxxx

Advance Access Publication Date: Day Month Year

Manuscript Category



Subject Section

SneakySnake: A Fast and Accurate Universal Genome Pre-Alignment Filter for CPUs, GPUs, and FPGAs

Mohammed Alser ^{1,2,*}, Taha Shahroodi ¹, Juan Gómez-Luna ^{1,2}, Can Alkan ^{4,*}, and Onur Mutlu ^{1,2,3,4,*}

¹Department of Computer Science, ETH Zurich, Zurich 8006, Switzerland

²Department of Information Technology and Electrical Engineering, ETH Zurich, Zurich 8006, Switzerland

³Department of Electrical and Computer Engineering, Carnegie Mellon University, Pittsburgh 15213, PA, USA

Department of Computer Engineering, Bilkent University, Ankara 06800, Turkey

GenASM Framework [MICRO 2020]

Damla Senol Cali, Gurpreet S. Kalsi, Zulal Bingol, Can Firtina, Lavanya Subramanian, Jeremie S. Kim, Rachata Ausavarungnirun, Mohammed Alser, Juan Gomez-Luna, Amirali Boroumand, Anant Nori, Allison Scibisz, Sreenivas Subramoney, Can Alkan, Saugata Ghose, and Onur Mutlu, "GenASM: A High-Performance, Low-Power Approximate String Matching Acceleration Framework for Genome Sequence Analysis"
Proceedings of the 53rd International Symposium on Microarchitecture (MICRO), Virtual, October 2020.

[<u>Lighting Talk Video</u> (1.5 minutes)]
[<u>Lightning Talk Slides (pptx) (pdf)</u>]
[<u>Talk Video</u> (18 minutes)]
[<u>Slides (pptx) (pdf)</u>]

GenASM: A High-Performance, Low-Power Approximate String Matching Acceleration Framework for Genome Sequence Analysis

Damla Senol Cali^{†™} Gurpreet S. Kalsi[™] Zülal Bingöl[▽] Can Firtina[⋄] Lavanya Subramanian[‡] Jeremie S. Kim^{⋄†} Rachata Ausavarungnirun[⊙] Mohammed Alser[⋄] Juan Gomez-Luna[⋄] Amirali Boroumand[†] Anant Nori[™] Allison Scibisz[†] Sreenivas Subramoney[™] Can Alkan[▽] Saugata Ghose^{*†} Onur Mutlu^{⋄†▽}

† Carnegie Mellon University [™] Processor Architecture Research Lab, Intel Labs [▽] Bilkent University [⋄] ETH Zürich

‡ Facebook [⊙] King Mongkut's University of Technology North Bangkok ^{*} University of Illinois at Urbana–Champaign

422

Quick Note: Key Principles and Results

- Two key principles:
 - Exploit the structure of the genome to minimize computation
 - Morph and exploit the structure of the underlying hardware to maximize performance and efficiency

- Algorithm-architecture co-design for DNA read mapping
 - Speeds up read mapping by ~100-1000X
 - Improves accuracy of read mapping in the presence of errors

New Genome Sequencing Technologies

Nanopore sequencing technology and tools for genome assembly: computational analysis of the current state, bottlenecks and future directions

Damla Senol Cali ™, Jeremie S Kim, Saugata Ghose, Can Alkan, Onur Mutlu

Briefings in Bioinformatics, bby017, https://doi.org/10.1093/bib/bby017

Published: 02 April 2018 Article history ▼



Oxford Nanopore MinION

Senol Cali+, "Nanopore Sequencing Technology and Tools for Genome Assembly: Computational Analysis of the Current State, Bottlenecks and Future Directions," Briefings in Bioinformatics, 2018.

[Preliminary arxiv.org version]

Nanopore Genome Assembly Pipeline

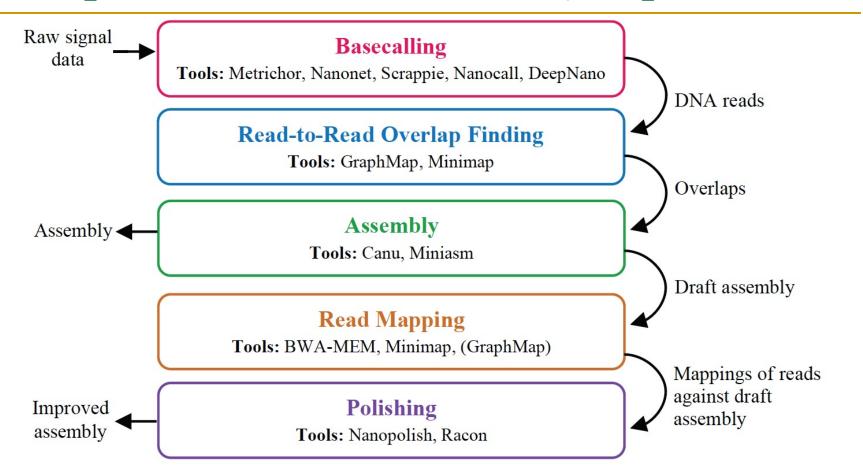


Figure 1. The analyzed genome assembly pipeline using nanopore sequence data, with its five steps and the associated tools for each

___step.

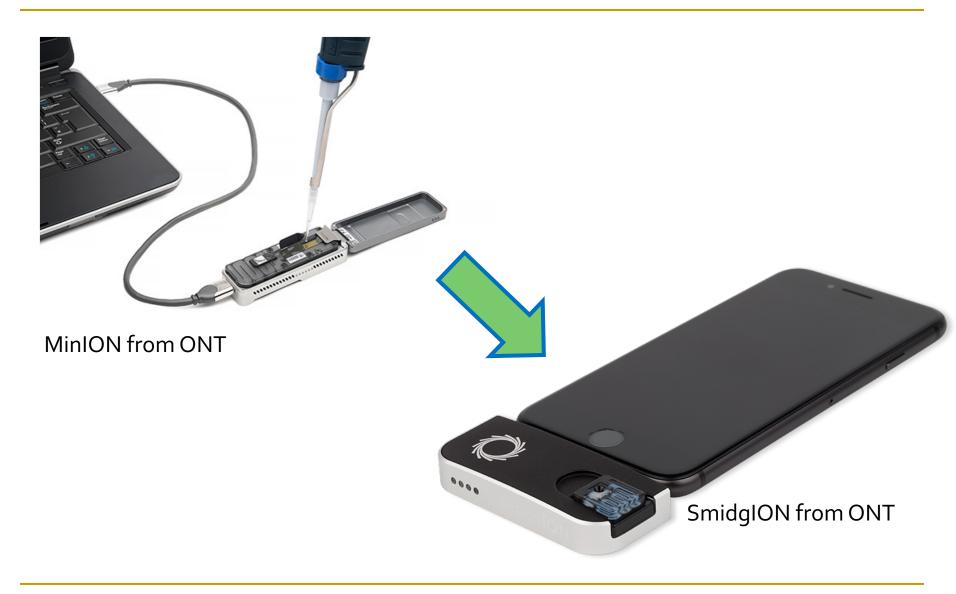
Senol Cali+, "Nanopore Sequencing Technology and Tools for Genome Assembly," Briefings in Bioinformatics, 2018.

425

Recall Our Dream

- An embedded device that can perform comprehensive genome analysis in real time (within a minute)
- Still a long ways to go
 - Energy efficiency
 - Performance (latency)
 - Security
 - Huge memory bottleneck

Future of Genome Sequencing & Analysis



Why Do We Care? An Example from 2020

200 Oxford Nanopore sequencers have left UK for China, to support rapid, near-sample coronavirus sequencing for outbreak surveillance

Fri 31st January 2020

Following extensive support of, and collaboration with, public health professionals in China, Oxford Nanopore has shipped an additional 200 MinION sequencers and related consumables to China. These will be used to support the ongoing surveillance of the current coronavirus outbreak, adding to a large number of the devices already installed in the country.



Each MinION sequencer is approximately the size of a stapler, and can provide rapid sequence information about the coronavirus.





700Kg of Oxford Nanopore sequencers and consumables are on their way for use by Chinese scientists in understanding the current coronavirus outbreak.



Sequencing of COVID-19

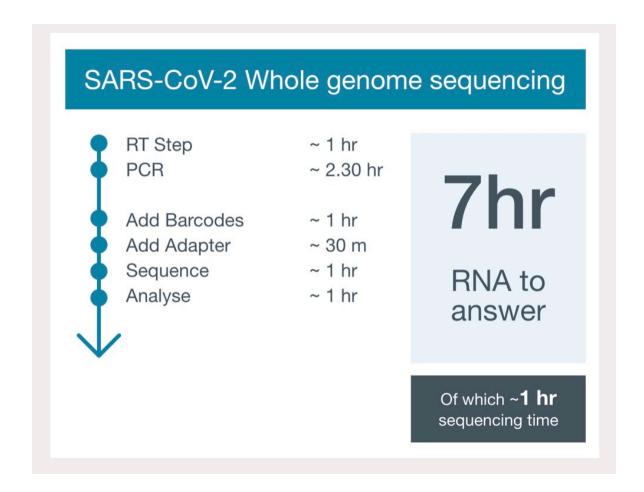
Whole genome sequencing (WGS) and sequence data analysis are important

- To detect the virus from a human sample such as saliva,
 Bronchoalveolar fluid etc.
- To understand the sources and modes of transmission of the virus
- To discover the genomic characteristics of the virus, and compare with better-known viruses (e.g., 02-03 SARS epidemic)
- To design and evaluate the diagnostic tests and deep-dive studies

Two key areas of COVID-19 genomic research

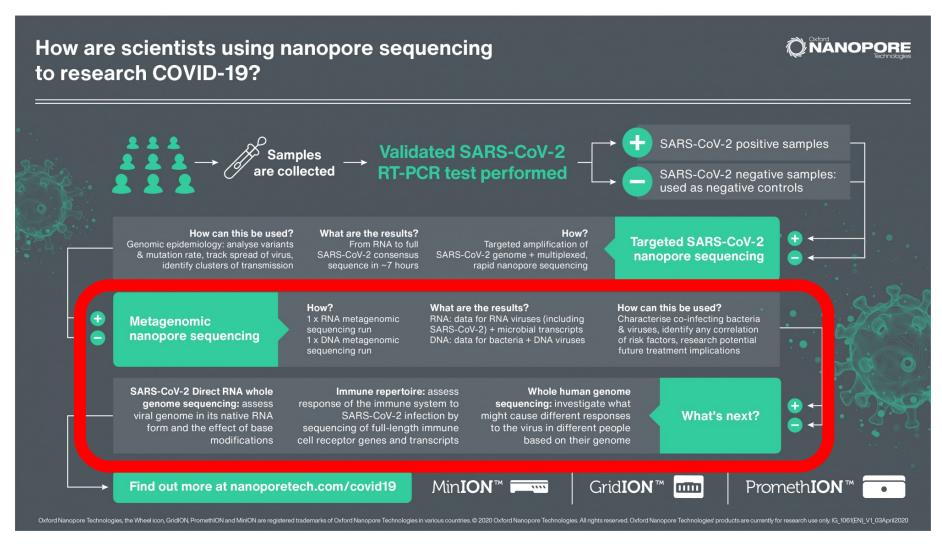
- To sequence the genome of the virus itself, COVID-19, in order to track the mutations in the virus.
- To explore the genes of infected patients. This analysis can be used to understand why some people get more severe symptoms than others, as well as, help with the development of new treatments in the future.

COVID-19 Nanopore Sequencing (I)



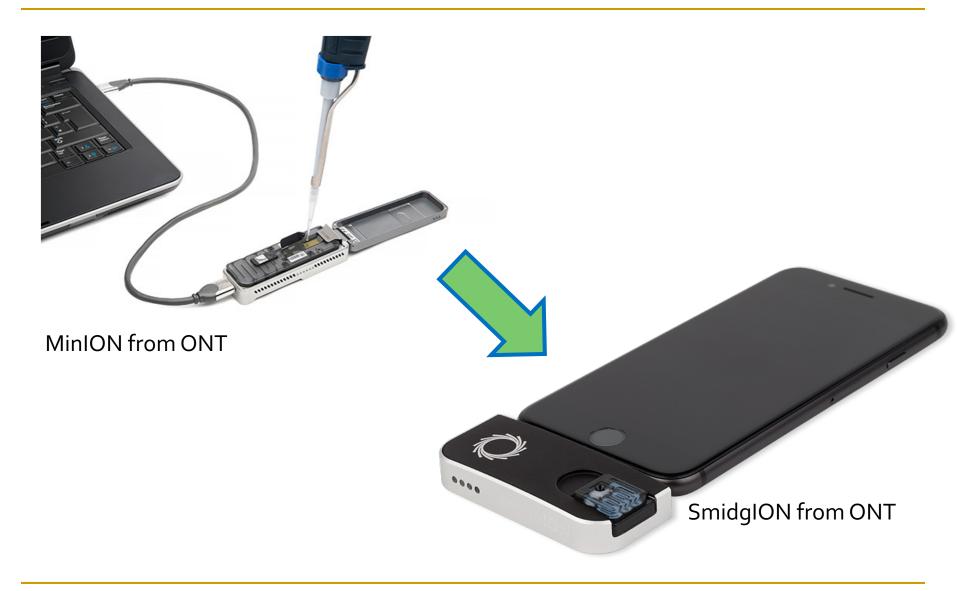
From ONT (https://nanoporetech.com/covid-19/overview)

COVID-19 Nanopore Sequencing (II)



From ONT (https://nanoporetech.com/covid-19/overview)

Future of Genome Sequencing & Analysis



Accelerating Genome Analysis: Overview

 Mohammed Alser, Zulal Bingol, Damla Senol Cali, Jeremie Kim, Saugata Ghose, Can Alkan, and Onur Mutlu,

"Accelerating Genome Analysis: A Primer on an Ongoing Journey"

IEEE Micro (IEEE MICRO), Vol. 40, No. 5, pages 65-75, September/October 2020.

[Slides (pptx)(pdf)]

[Talk Video (1 hour 2 minutes)]

Accelerating Genome Analysis: A Primer on an Ongoing Journey

Mohammed Alser

ETH Zürich

Zülal Bingöl

Bilkent University

Damla Senol Cali

Carnegie Mellon University

Jeremie Kim

ETH Zurich and Carnegie Mellon University

Saugata Ghose

University of Illinois at Urbana–Champaign and Carnegie Mellon University

Can Alkan

Bilkent University

Onur Mutlu

ETH Zurich, Carnegie Mellon University, and Bilkent University

More on Fast Genome Analysis ...

Onur Mutlu,

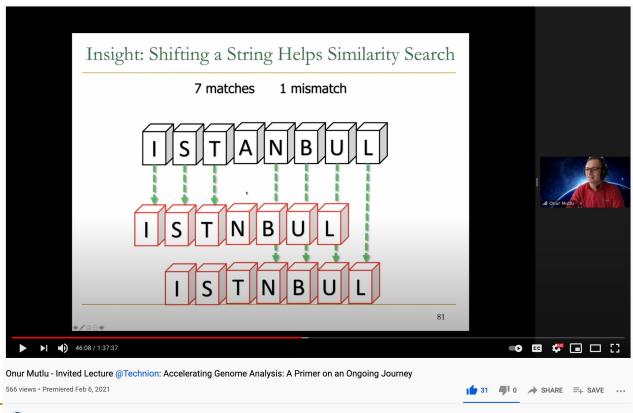
"Accelerating Genome Analysis: A Primer on an Ongoing Journey"

Invited Lecture at <u>Technion</u>, Virtual, 26 January 2021.

[Slides (pptx) (pdf)]

[Talk Video (1 hour 37 minutes, including Q&A)]

[Related Invited Paper (at IEEE Micro, 2020)]





Detailed Lectures on Genome Analysis

- Computer Architecture, Fall 2020, Lecture 3a
 - Introduction to Genome Sequence Analysis (ETH Zürich, Fall 2020)
 - https://www.youtube.com/watch?v=CrRb32v7SJc&list=PL5Q2soXY2Zi9xidyIgBxUz7xRPS-wisBN&index=5
- Computer Architecture, Fall 2020, Lecture 8
 - □ **Intelligent Genome Analysis** (ETH Zürich, Fall 2020)
 - https://www.youtube.com/watch?v=ygmQpdDTL7o&list=PL5Q2soXY2Zi9xidyIgBxUz7xRPS-wisBN&index=14
- Computer Architecture, Fall 2020, Lecture 9a
 - □ **GenASM: Approx. String Matching Accelerator** (ETH Zürich, Fall 2020)
 - https://www.youtube.com/watch?v=XoLpzmN Pas&list=PL5Q2soXY2Zi9xidyIgBxUz7xRPS-wisBN&index=15
- Accelerating Genomics Project Course, Fall 2020, Lecture 1
 - Accelerating Genomics (ETH Zürich, Fall 2020)
 - https://www.youtube.com/watch?v=rgjl8ZyLsAg&list=PL5Q2soXY2Zi9E2bBVAgCqL gwiDRQDTyId

Challenge and Opportunity for Future

High Performance

(to solve the **toughest** & **all** problems)

Challenge and Opportunity for Future

Personalized and Private

```
(in every aspect of life:
health, medicine,
spaces, devices, robotics, ...)
```

More on My Research & Teaching

Brief Self Introduction



Onur Mutlu

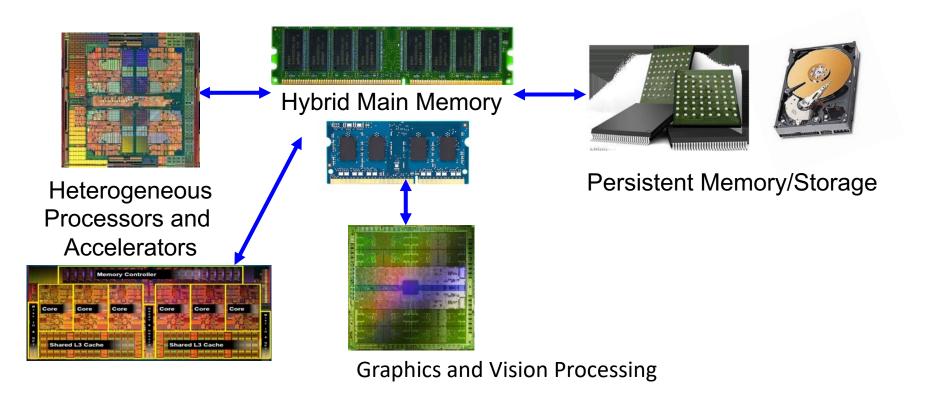
- □ Full Professor @ ETH Zurich ITET (INFK), since September 2015
- □ Strecker Professor @ Carnegie Mellon University ECE/CS, 2009-2016, 2016-...
- PhD from UT-Austin, worked at Google, VMware, Microsoft Research, Intel, AMD
- https://people.inf.ethz.ch/omutlu/
- omutlu@gmail.com (Best way to reach me)
- https://people.inf.ethz.ch/omutlu/projects.htm

Research and Teaching in:

- Computer architecture, computer systems, hardware security, bioinformatics
- Memory and storage systems
- Hardware security, safety, predictability
- Fault tolerance
- Hardware/software cooperation
- Architectures for bioinformatics, health, medicine
- **-** ...

Current Research Mission

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



Build fundamentally better architectures

Four Key Current Directions

Fundamentally Secure/Reliable/Safe Architectures

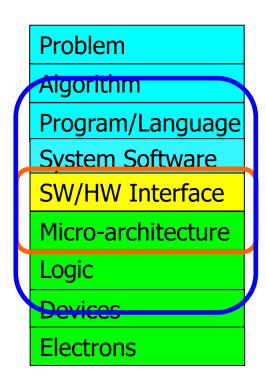
- Fundamentally Energy-Efficient Architectures
 - Memory-centric (Data-centric) Architectures

Fundamentally Low-Latency and Predictable Architectures

Architectures for AI/ML, Genomics, Medicine, Health

The Transformation Hierarchy

Computer Architecture (expanded view)



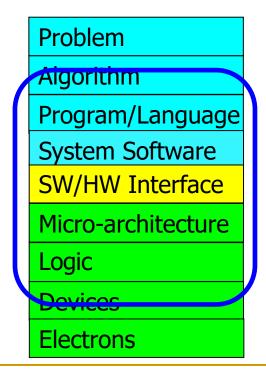
Computer Architecture (narrow view)

Axiom

To achieve the highest energy efficiency and performance:

we must take the expanded view

of computer architecture

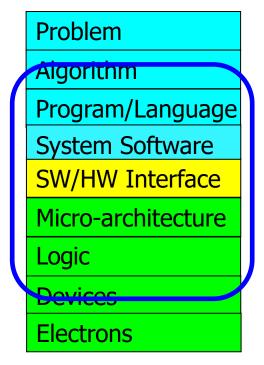


Co-design across the hierarchy:
Algorithms to devices

Specialize as much as possible within the design goals

Current Research Mission & Major Topics

Build fundamentally better architectures



Broad research spanning apps, systems, logic with architecture at the center

- Data-centric arch. for low energy & high perf.
 - Proc. in Mem/DRAM, NVM, unified mem/storage
- Low-latency & predictable architectures
 - Low-latency, low-energy yet low-cost memory
 - QoS-aware and predictable memory systems
- Fundamentally secure/reliable/safe arch.
 - Tolerating all bit flips; patchable HW; secure mem
- Architectures for ML/AI/Genomics/Health/Med
 - Algorithm/arch./logic co-design; full heterogeneity
- Data-driven and data-aware architectures
 - ML/AI-driven architectural controllers and design
 - Expressive memory and expressive systems

Onur Mutlu's SAFARI Research Group

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

https://safari.ethz.ch/safari-newsletter-april-2020/



Think BIG, Aim HIGH!

SAFARI

https://safari.ethz.ch

SAFARI Newsletter January 2021 Edition

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





Newsletter January 2021

Think Big, Aim High, and Have a Wonderful 2021!



Dear SAFARI friends,

Principle: Teaching and Research

Teaching drives Research Research drives Teaching

•

Focus on Insight Encourage New Ideas

Research & Teaching: Some Overview Talks

https://www.youtube.com/onurmutlulectures

- Future Computing Architectures
 - https://www.youtube.com/watch?v=kgiZISOcGFM&list=PL5Q2soXY2Zi8D_5MGV6EnXEJHnV2YFBJI&index=1
- Enabling In-Memory Computation
 - https://www.youtube.com/watch?v=njX 14584Jw&list=PL5Q2soXY2Zi8D 5MGV6EnXEJHnV2YFBJl&index=16
- Accelerating Genome Analysis
 - https://www.youtube.com/watch?v=r7sn41lH-4A&list=PL5Q2soXY2Zi8D_5MGV6EnXEJHnV2YFBJl&index=41
- Rethinking Memory System Design
 - https://www.youtube.com/watch?v=F7xZLNMIY1E&list=PL5Q2soXY2Zi8D_5MGV6EnXEJHnV2YFBJl&index=3
- Intelligent Architectures for Intelligent Machines
 - https://www.youtube.com/watch?v=c6_LgzuNdkw&list=PL5Q2soXY2Zi8D_5MGV6EnXEJHnV2YFBJl&index=25
- The Story of RowHammer
 - https://www.youtube.com/watch?v=sgd7PHQQ1AI&list=PL5Q2soXY2Zi8D_5MGV6EnXEJHnV2YFBJl&index=39

An Interview on Research and Education

- Computing Research and Education (@ ISCA 2019)
 - https://www.youtube.com/watch?v=8ffSEKZhmvo&list=PL5Q2 soXY2Zi_4oP9LdL3cc8G6NIjD2Ydz

- Maurice Wilkes Award Speech (10 minutes)
 - https://www.youtube.com/watch?v=tcQ3zZ3JpuA&list=PL5Q2 soXY2Zi8D_5MGV6EnXEJHnV2YFBJl&index=15

More Thoughts and Suggestions

Onur Mutlu,

"Some Reflections (on DRAM)"

Award Speech for <u>ACM SIGARCH Maurice Wilkes Award</u>, at the **ISCA** Awards Ceremony, Phoenix, AZ, USA, 25 June 2019.

[Slides (pptx) (pdf)]

[Video of Award Acceptance Speech (Youtube; 10 minutes) (Youku; 13 minutes)]

[Video of Interview after Award Acceptance (Youtube; 1 hour 6 minutes) (Youku;

1 hour 6 minutes)

[News Article on "ACM SIGARCH Maurice Wilkes Award goes to Prof. Onur Mutlu"]

Onur Mutlu,

"How to Build an Impactful Research Group"

57th Design Automation Conference Early Career Workshop (DAC), Virtual, 19 July 2020.

[Slides (pptx) (pdf)]

Referenced Papers, Talks, Artifacts

All are available at

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

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

https://www.youtube.com/onurmutlulectures

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

Readings, Videos, Reference Materials

Research & Teaching: Some Overview Talks

https://www.youtube.com/onurmutlulectures

- Future Computing Architectures
 - https://www.youtube.com/watch?v=kgiZISOcGFM&list=PL5Q2soXY2Zi8D_5MGV6EnXEJHnV2YFBJI&index=1
- Enabling In-Memory Computation
 - https://www.youtube.com/watch?v=njX 14584Jw&list=PL5Q2soXY2Zi8D 5MGV6EnXEJHnV2YFBJl&index=16
- Accelerating Genome Analysis
 - https://www.youtube.com/watch?v=r7sn41lH-4A&list=PL5Q2soXY2Zi8D_5MGV6EnXEJHnV2YFBJl&index=41
- Rethinking Memory System Design
 - https://www.youtube.com/watch?v=F7xZLNMIY1E&list=PL5Q2soXY2Zi8D_5MGV6EnXEJHnV2YFBJl&index=3
- Intelligent Architectures for Intelligent Machines
 - https://www.youtube.com/watch?v=c6_LgzuNdkw&list=PL5Q2soXY2Zi8D_5MGV6EnXEJHnV2YFBJl&index=25
- The Story of RowHammer
 - https://www.youtube.com/watch?v=sqd7PHQQ1AI&list=PL5Q2soXY2Zi8D 5MGV6EnXEJHnV2YFBJl&index=39

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://people.inf.ethz.ch/omutlu/acaces2018.html
 - https://www.youtube.com/playlist?list=PL5Q2soXY2Zi-HXxomthrpDpMJm05P6J9x

All Papers are at:

- https://people.inf.ethz.ch/omutlu/projects.htm
- 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 VoMlKlw7YrXOhbl

All Papers are at:

- https://people.inf.ethz.ch/omutlu/projects.htm
- Final lecture notes and readings (for all topics)

An Interview on Research and Education

- Computing Research and Education (@ ISCA 2019)
 - https://www.youtube.com/watch?v=8ffSEKZhmvo&list=PL5Q2 soXY2Zi_4oP9LdL3cc8G6NIjD2Ydz

- Maurice Wilkes Award Speech (10 minutes)
 - https://www.youtube.com/watch?v=tcQ3zZ3JpuA&list=PL5Q2 soXY2Zi8D_5MGV6EnXEJHnV2YFBJl&index=15

More Thoughts and Suggestions

Onur Mutlu,

"Some Reflections (on DRAM)"

Award Speech for <u>ACM SIGARCH Maurice Wilkes Award</u>, at the **ISCA** Awards Ceremony, Phoenix, AZ, USA, 25 June 2019.

[Slides (pptx) (pdf)]

[Video of Award Acceptance Speech (Youtube; 10 minutes) (Youku; 13 minutes)]

[Video of Interview after Award Acceptance (Youtube; 1 hour 6 minutes)] (Youku;

1 hour 6 minutes)

[News Article on "ACM SIGARCH Maurice Wilkes Award goes to Prof. Onur Mutlu"]

Onur Mutlu,

"How to Build an Impactful Research Group"

57th Design Automation Conference Early Career Workshop (DAC), Virtual, 19 July 2020.

[Slides (pptx) (pdf)]

Reference Overview Paper

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 Emerging Computing: From Devices to Systems
Looking Beyond Moore and Von Neumann, Springer, to be published in 2021.

Reference Overview Paper I

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 <u>Microprocessors and Microsystems</u> (**MICPRO**), June 2019. [arXiv version]

SAFARI

Reference Overview Paper 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]

Reference Overview Paper III

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]

Reference Overview Paper IV

Onur Mutlu and Lavanya Subramanian,
 "Research Problems and Opportunities in Memory Systems"

Invited Article in <u>Supercomputing Frontiers and Innovations</u> (**SUPERFRI**), 2014/2015.

Research Problems and Opportunities in Memory Systems

Onur Mutlu¹, Lavanya Subramanian¹

Reference Overview Paper V

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> <u>Europe Conference</u> (**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

Reference Overview Paper VI

Onur Mutlu,
 "Memory Scaling: A Systems Architecture
 Perspective"

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/

Reference Overview Paper VII



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

Reference Overview Paper VIII

Onur Mutlu and Jeremie Kim,

"RowHammer: A Retrospective"

<u>IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems</u> (**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 (30 minutes)]

RowHammer: A Retrospective

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

SAFARI 467

Related Videos and Course Materials (I)

- Undergraduate Digital Design & Computer
 Architecture Course Lecture
 Videos (2020, 2019, 2018, 2017, 2015, 2014, 2013)
- Undergraduate Digital Design & Computer
 Architecture Course
 Materials (2020, 2019, 2018, 2015, 2014, 2013)
- Graduate Computer Architecture Course Lecture
 Videos (2019, 2018, 2017, 2015, 2013)
- Graduate Computer Architecture Course
 Materials (2019, 2018, 2017, 2015, 2013)
- Parallel Computer Architecture Course Materials (Lecture Videos)

Related Videos and Course Materials (II)

- Seminar in Computer Architecture Course Lecture
 Videos (Spring 2020, Fall 2019, Spring 2019, 2018)
- Seminar in Computer Architecture Course
 Materials (Spring 2020, Fall 2019, Spring 2019, 2018)
- Memory Systems Course Lecture Videos (Sept 2019, July 2019, June 2019, October 2018)
- Memory Systems Short Course Lecture Materials (Sept 2019, July 2019, June 2019, October 2018)
- ACACES Summer School Memory Systems Course Lecture Videos (2018, 2013)
- ACACES Summer School Memory Systems Course Materials (2018, 2013)

Some Open Source Tools (I)

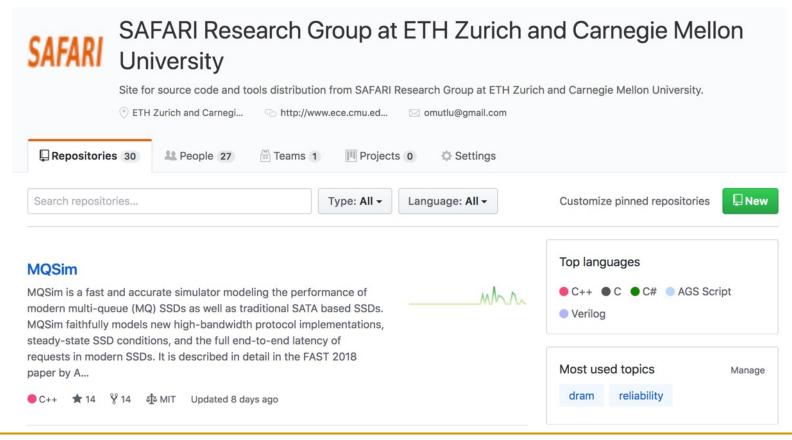
- Rowhammer Program to Induce RowHammer Errors
 - https://github.com/CMU-SAFARI/rowhammer
- Ramulator Fast and Extensible DRAM Simulator
 - https://github.com/CMU-SAFARI/ramulator
- MemSim Simple Memory Simulator
 - https://github.com/CMU-SAFARI/memsim
- NOCulator Flexible Network-on-Chip Simulator
 - https://github.com/CMU-SAFARI/NOCulator
- SoftMC FPGA-Based DRAM Testing Infrastructure
 - https://github.com/CMU-SAFARI/SoftMC
- Other open-source software from my group
 - https://github.com/CMU-SAFARI/
 - http://www.ece.cmu.edu/~safari/tools.html

Some Open Source Tools (II)

- MQSim A Fast Modern SSD Simulator
 - https://github.com/CMU-SAFARI/MQSim
- 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
 - https://github.com/CMU-SAFARI/HWASim
- Other open-source software from my group
 - https://github.com/CMU-SAFARI/
 - http://www.ece.cmu.edu/~safari/tools.html

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



ramulator-pim

A fast and flexible simulation infrastructure for exploring general-purpose processing-in-memory (PIM) architectures. Ramulator-PIM combines a widely-used simulator for out-of-order and in-order processors (ZSim) with Ramulator, a DRAM simulator with memory models for DDRx, LPDDRx, GDDRx, WIOx, HBMx, and HMCx. Ramulator is described in the IEEE ...

● C++ ♀ 11 ☆ 29 ① 6 ┆ 0 Updated 19 days ago

SMASH

SMASH is a hardware-software cooperative mechanism that enables highly-efficient indexing and storage of sparse matrices. The key idea of SMASH is to compress sparse matrices with a hierarchical bitmap compression format that can be accelerated from hardware.

Described by Kanellopoulos et al. (MICRO '19) https://people.inf.ethz.ch/omutlu/pub/SMA...

●C ೪1 ☆6 ①0 ♯0 Updated on May 17

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 requests in modern SSDs. It is described in detail in the FAST 2018 paper by A...

●C++ គ្ MIT ೪ 54 ☆62 ①10 រឿ 1 Updated on May 15

Apollo

Apollo is an assembly polishing algorithm that attempts to correct the errors in an assembly. It can take multiple set of reads in a single run and polish the assemblies of genomes of any size. Described in the Bioinformatics journal paper (2020) by Firtina et al. at https://people.inf.ethz.ch/omutlu/pub/apollotechnology-independent-genome-asse...

ramulator

A Fast and Extensible DRAM Simulator, with built-in support for modeling many different DRAM technologies including DDRx, LPDDRx, GDDRx, WIOx, HBMx, and various academic proposals. Described in the IEEE CAL 2015 paper by Kim et al. at

http://users.ece.cmu.edu/~omutlu/pub/ramulator_dram_simulator-ieee-cal15.pdf

●C++ Φ MIT 및 93 ☆ 170 ① 37 % 2 Updated on Apr 13

Shifted-Hamming-Distance

Source code for the Shifted Hamming Distance (SHD)
filtering mechanism for sequence alignment. Described
in the Bioinformatics journal paper (2015) by Xin et al. at
http://users.ece.cmu.edu/~omutlu/pub/shiftedhamming-distance_bioinformatics15_proofs.pdf

SneakySnake

The first and the only pre-alignment filtering algorithm that works on all modern high-performance computing architectures. It works efficiently and fast on CPU, FPGA, and GPU architectures and that greatly (by more than two orders of magnitude) expedites sequence alignment calculation. Described by Alser et al. (preliminary version at https://a...

AirLift

AirLift is a tool that updates mapped reads from one reference genome to another. Unlike existing tools, It accounts for regions not shared between the two reference genomes and enables remapping across all parts of the references. Described by Kim et al. (preliminary version at http://arxiv.org/abs/1912.08735)

●C ♀O ☆3 ①O ₺ O Updated on Feb 19

GPGPUSim-Ramulator

The source code for GPGPUSim+Ramulator simulator. In this version, GPGPUSim uses Ramulator to simulate the DRAM. This simulator is used to produce some of the

Referenced Papers, Talks, Artifacts

All are available at

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

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

https://www.youtube.com/onurmutlulectures

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

An Interview on Research and Education

- Computing Research and Education (@ ISCA 2019)
 - https://www.youtube.com/watch?v=8ffSEKZhmvo&list=PL5Q2 soXY2Zi_4oP9LdL3cc8G6NIjD2Ydz

- Maurice Wilkes Award Speech (10 minutes)
 - https://www.youtube.com/watch?v=tcQ3zZ3JpuA&list=PL5Q2 soXY2Zi8D_5MGV6EnXEJHnV2YFBJl&index=15

More Thoughts and Suggestions

Onur Mutlu,

"Some Reflections (on DRAM)"

Award Speech for <u>ACM SIGARCH Maurice Wilkes Award</u>, at the **ISCA** Awards Ceremony, Phoenix, AZ, USA, 25 June 2019.

[Slides (pptx) (pdf)]

[Video of Award Acceptance Speech (Youtube; 10 minutes) (Youku; 13 minutes)]

[Video of Interview after Award Acceptance (Youtube; 1 hour 6 minutes) (Youku;

1 hour 6 minutes)

[News Article on "ACM SIGARCH Maurice Wilkes Award goes to Prof. Onur Mutlu"]

Onur Mutlu,

"How to Build an Impactful Research Group"

57th Design Automation Conference Early Career Workshop (DAC), Virtual, 19 July 2020.

[Slides (pptx) (pdf)]

End of Backup Slides