# Tutorial on
# Memory-Centric Computing:
## Conclusion Remarks

Geraldo F. Oliveira

Prof. Onur Mutlu

ISCA 2024

29 June 2024

SAFARI

ETH zürich

# Agenda

- Introduction to Memory-Centric Computing Systems

- Invited Talk by Prof. Minsoo Rhu:
  "*Memory-Centric Computing Systems – For AI and Beyond*"

- Coffee Break

- Real-World Processing-Near-Memory Systems

- Processing-Using-Memory Architectures for Bulk Bitwise Op.

- Invited Talk by Prof. Saugata Ghose:
  "*RACER and ReRAM PUM*"

- PIM Programming & Infrastructure for PIM Research

- Closing Remarks

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

# Fundamentally High-Performance (Data-Centric) Computing Architectures

# Computing Architectures with Minimal Data Movement

# Concluding Remarks

- We must design systems to be **balanced**, **high-performance**, **energy-efficient** (all at the same time) → intelligent systems
  - **Data-centric, data-driven, data-aware**

- Enable computation capability inside and close to memory

- This can
  - Lead to **orders-of-magnitude** improvements
  - **Enable new applications & computing platforms**
  - **Enable better understanding of nature**
  - **...**

- Future of **truly memory-centric computing** is bright
  - We need to do research & design across the computing stack

# Fundamentally Better Architectures

**Data-centric**

**Data-driven**

**Data-aware**

# We Need to Revisit the Entire Stack

| Problem |
| --- |
| Algorithm |
| Program/Language |
| System Software |
| SW/HW Interface |
| Micro-architecture |
| Logic |
| Devices |
| Electrons |

**We can get there step by step**

# We Need to Exploit Good Principles

- Data-centric system design

- All components intelligent

- Better (cross-layer) communication, better interfaces

- Better-than-worst-case design

- Heterogeneity

- Flexibility, adaptability

**Open minds**

# 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*

[a]*ETH Zürich*
[b]*Carnegie Mellon University*
[c]*University of Illinois at Urbana-Champaign*
[d]*King 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.*

# Special Research Sessions & Courses (I)

- Special Session at ISVLSI 2022: 9 cutting-edge talks



**SAFARI**

**https://www.youtube.com/watch?v=qeukNs5XI3g**

# Special Research Sessions & Courses (II)

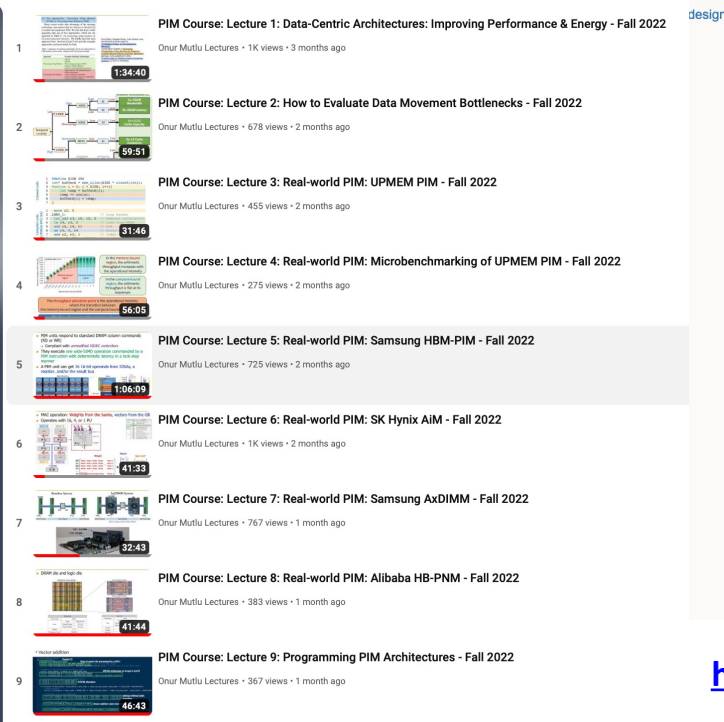- **Special Session at ISVLSI 2022: 9 cutting-edge talks**



19 — **GenStore: In-Storage Filtering for High-Performance and Energy-Efficient Genome Analysis**
Onur Mutlu Lectures • • Premieres 3/12/23, 7:00 PM
UPCOMING

20 — **Introduction to the ISVLSI 2022 Special Session on Processing-in-Memory**
Onur Mutlu Lectures • 286 views • 2 days ago

21 — **Heterogeneous Data-Centric Architectures for Data-Intensive Applications: Case Studies in ML and DB**
Onur Mutlu Lectures • 2 waiting • Premieres 3/10/23, 7:00 PM
UPCOMING

22 — **Machine Learning Training on a Real Processing-In-Memory System**
Onur Mutlu Lectures • • Premieres 3/14/23, 7:00 PM
UPCOMING

23 — **Exploiting Near-Data Processing to Accelerate Time Series Analysis**
Onur Mutlu Lectures • • Premieres 3/11/23, 7:00 PM
UPCOMING

24 — **PiDRAM: An FPGA-Based Framework for End-To-End Evaluation of Processing-In-DRAM Techniques**
Onur Mutlu Lectures • • Premieres 3/9/23, 7:00 PM
UPCOMING

25 — **The Road to Widely Deploying Processing-In-Memory: Challenges and Opportunities**
Onur Mutlu Lectures • 399 views • 1 day ago

26 — **SparseP: Efficient Sparse Matrix Vector Multiplication on Real Processing-In-Memory Architectures**
Onur Mutlu Lectures • 1 waiting • Premieres 3/13/23, 7:00 PM
UPCOMING

27 — **HPCA 2023 Tutorial: Real-World Processing-in-Memory Architectures**
Onur Mutlu Lectures • 1.6K views • Streamed 10 days ago

# Processing-in-Memory Course (Fall 2022)

- Short weekly lectures
- Hands-on projects



https://safari.ethz.ch/projects_and_seminars/fall2022/
doku.php?id=processing_in_memory

https://youtube.com/playlist?list=PL5Q2soXY2Zi8KzG2CQYRNQOVD0GOBrnKy

**SAFARI**

# PIM Course (Fall 2022)



- **Fall 2022 Edition:**
  - https://safari.ethz.ch/projects_and_seminars/fall2022/doku.php?id=processing_in_memory

- **Spring 2022 Edition:**
  - https://safari.ethz.ch/projects_and_seminars/spring2022/doku.php?id=processing_in_memory

- **Youtube Livestream (Fall 2022):**
  - https://www.youtube.com/watch?v=QLL0wQ9I4Dw&list=PL5Q2soXY2Zi8KzG2CQYRNQOVD0GOBrnKy

- **Youtube Livestream (Spring 2022):**
  - https://www.youtube.com/watch?v=9e4Chnwdovo&list=PL5Q2soXY2Zi-841fUYYUK9EsXKhQKRPyX

- Project course
  - Taken by Bachelor's/Master's students
  - Processing-in-Memory lectures
  - Hands-on research exploration
  - Many research readings

**https://www.youtube.com/onurmutlulectures**

**SAFARI**

Spring 2022 Meetings/Schedule

| Week | Date | Livestream | Meeting | Learning Materials | Assignments |
|------|------|-----------|---------|-------------------|-------------|
| W1 | 10.03 Thu. | Live | M1: P&S PIM Course Presentation (PDF) (PPT) | Required Materials Recommended Materials | HW 0 Out |
| W2 | 15.03 Tue. | | Hands-on Project Proposals | | |
| | 17.03 Thu. | Premiere | M2: Real-world PIM: UPMEM PIM (PDF) (PPT) | | |
| W3 | 24.03 Thu. | Live | M3: Real-world PIM: Microbenchmarking of UPMEM PIM (PDF) (PPT) | | |
| W4 | 31.03 Thu. | Live | M4: Real-world PIM: Samsung HBM-PIM (PDF) (PPT) | | |
| W5 | 07.04 Thu. | Live | M5: How to Evaluate Data Movement Bottlenecks (PDF) (PPT) | | |
| W6 | 14.04 Thu. | Live | M6: Real-world PIM: SK Hynix AiM (PDF) (PPT) | | |
| W7 | 21.04 Thu. | Premiere | M7: Programming PIM Architectures (PDF) (PPT) | | |
| W8 | 28.04 Thu. | Premiere | M8: Benchmarking and Workload Suitability on PIM (PDF) (PPT) | | |
| W9 | 05.05 Thu. | Premiere | M9: Real-world PIM: Samsung AxDIMM (PDF) (PPT) | | |
| W10 | 12.05 Thu. | Premiere | M10: Real-world PIM: Alibaba HB-PNM (PDF) (PPT) | | |
| W11 | 19.05 Thu. | Live | M11: SpMV on a Real PIM Architecture (PDF) (PPT) | | |
| W12 | 26.05 Thu. | Live | M12: End-to-End Framework for Processing-using-Memory (PDF) (PPT) | | |
| W13 | 02.06 Thu. | Live | M13: Bit-Serial SIMD Processing using DRAM (PDF) (PPT) | | |
| W14 | 09.06 Thu. | Live | M14: Analyzing and Mitigating ML Inference Bottlenecks (PDF) (PPT) | | |
| W15 | 15.06 Thu. | Live | M15: In-Memory HTAP Databases with HW/SW Co-design (PDF) (PPT) | | |
| W16 | 23.06 Thu. | Live | M16: In-Storage Processing for Genome Analysis (PDF) (PPT) | | |
| W17 | 18.07 Mon. | Premiere | M17: How to Enable the Adoption of PIM? (PDF) (PPT) | | |
| W18 | 09.08 Tue. | Premiere | SS1: ISVLSI 2022 Special Session on PIM (PDF & PPT) | | |

# Processing-in-Memory Course (Spring 2023)

- Short weekly lectures
- Hands-on projects

**SAFARI**

# Real PIM Tutorials [ISCA'23, ASPLOS'23, HPCA'23]

- June, March, Feb : Lectures + Hands-on labs + Invited talks



**https://events.safari.ethz.ch/isca-pim-tutorial/**

# Real PIM Tutorial [ISCA 2023]

- June 18: Lectures + Hands-on labs + Invited talks



## Tutorial Materials

| Time | Speaker | Title | Materials |
|---|---|---|---|
| 8:55am-9:00am | Dr. Juan Gómez Luna | Welcome & Agenda | (PDF) (PPT) |
| 9:00am-10:20am | Prof. Onur Mutlu | Memory-Centric Computing | (PDF) (PPT) |
| 10:20am-11:00am | Dr. Juan Gómez Luna | Processing-Near-Memory: Real PNM Architectures / Programming General-purpose PIM | (PDF) (PPT) |
| 11:20am-11:50am | Prof. Izzat El Hajj | High-throughput Sequence Alignment using Real Processing-in-Memory Systems | (PDF) (PPT) |
| 11:50am-12:30pm | Dr. Christina Giannoula | SparseP: Towards Efficient Sparse Matrix Vector Multiplication for Real Processing-In-Memory Systems | (PDF) (PPT) |
| 2:00pm-2:45pm | Dr. Sukhan Lee | Introducing Real-world HBM-PIM Powered System for Memory-bound Applications | (PDF) (PPT) |
| 2:45pm-3:30pm | Dr. Juan Gómez Luna / Ataberk Olgun | Processing-Using-Memory: Exploiting the Analog Operational Properties of Memory Components / PUM Prototypes: PiDRAM | (PDF) (PPT) (PDF) (PPT) |
| 4:00pm-4:40pm | Dr. Juan Gómez Luna | Accelerating Modern Workloads on a General-purpose PIM System | (PDF) (PPT) |
| 4:40pm-5:20pm | Dr. Juan Gómez Luna | Adoption Issues: How to Enable PIM? | (PDF) (PPT) |
| 5:20pm-5:30pm | Dr. Juan Gómez Luna | Hands-on Lab: Programming and Understanding a Real Processing-in-Memory Architecture | (Handout) (PDF) (PPT) |



https://www.youtube.com/live/GIb5EgSrWk0

https://events.safari.ethz.ch/isca-pim-tutorial/

# Real PIM Tutorial [ASPLOS 2023]

■ **March 26: Lectures + Hands-on labs + Invited talks**



## Tutorial Materials

| Time | Speaker | Title | Materials |
|------|---------|-------|-----------|
| 9:00am-10:20am | Prof. Onur Mutlu | Memory-Centric Computing | (PDF) (PPT) |
| 10:40am-12:00pm | Dr. Juan Gómez Luna | Processing-Near-Memory: Real PNM Architectures Programming General-purpose PIM | (PDF) (PPT) |
| 1:40pm-2:20pm | Prof. Alexandra (Sasha) Fedorova (UBC) | Processing in Memory in the Wild | (PDF) (PPT) |
| 2:20pm-3:20pm | Dr. Juan Gómez Luna & Ataberk Olgun | Processing-Using-Memory: Exploiting the Analog Operational Properties of Memory Components | (PDF) (PPT) (PDF) (PPT) |
| 3:40pm-4:10pm | Dr. Juan Gómez Luna | Adoption issues: How to enable PIM? Accelerating Modern Workloads on a General-purpose PIM System | (PDF) (PPT) (PDF) (PPT) |
| 4:10pm-4:50pm | Dr. Yongkee Kwon & Eddy (Chanwook) Park (SK Hynix) | System Architecture and Software Stack for GDDR6-AiM | (PDF) (PPT) |
| 4:50pm-5:00pm | Dr. Juan Gómez Luna | Hands-on Lab: Programming and Understanding a Real Processing-in-Memory Architecture | (Handout) (PDF) (PPT) |

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

**https://events.safari.ethz.ch/asplos-pim-tutorial/**

# Real PIM Tutorial [HPCA 2023]

- ## February 26: Lectures + Hands-on labs + Invited Talks





https://www.youtube.com/watch?v=f5-nT1tbz5w

https://events.safari.ethz.ch/real-pim-tutorial/

| Time | Speaker | Title | Materials |
|------|---------|-------|-----------|
| 8:00am-8:40am | Prof. Onur Mutlu | Memory-Centric Computing | (PDF) (PPT) |
| 8:40am-10:00am | Dr. Juan Gómez Luna | Processing-Near-Memory: Real PNM Architectures Programming General-purpose PIM | (PDF) (PPT) |
| 10:20am-11:00am | Dr. Dimin Niu | A 3D Logic-to-DRAM Hybrid Bonding Process-Near-Memory Chip for Recommendation System | |
| 11:00am-11:40am | Dr. Christina Giannoula | SparseP: Towards Efficient Sparse Matrix Vector Multiplication on Real Processing-In-Memory Architectures | (PDF) (PPT) |
| 1:30pm-2:10pm | Dr. Juan Gómez Luna | Processing-Using-Memory: Exploiting the Analog Operational Properties of Memory Components | (PDF) (PPT) |
| 2:10pm-2:50pm | Dr. Manuel Le Gallo | Deep Learning Inference Using Computational Phase-Change Memory | |
| 2:50pm-3:30pm | Dr. Juan Gómez Luna | PIM Adoption Issues: How to Enable PIM Adoption? | (PDF) (PPT) |
| 3:40pm-5:40pm | Dr. Juan Gómez Luna | Hands-on Lab: Programming and Understanding a Real Processing-in-Memory Architecture | (Handout) (PDF) (PPT) |

# Real PIM Tutorial [MICRO 2023]

- October 29: Lectures + Hands-on labs + Invited talks





https://www.youtube.com/live/ohUooNSIxOI

https://events.safari.ethz.ch/micro-pim-tutorial

# PIM Tutorial at HEART 2024



**HEART 2024 Memory-Centric Computing Systems Tutorial**

Friday, June 21, Porto, Portugal

**Organizers:** Geraldo F. Oliveira, Dr. Mohammad Sadrosadati,
Ataberk Olgun, Professor Onur Mutlu
**Program:** https://events.safari.ethz.ch/heart24-memorycentric-tutorial/

Overview of PIM | PIM taxonomy
PIM in memory & storage
Real-world PNM systems
PUM for bulk bitwise operations
Programming techniques & tools
Infrastructures for PIM Research
Research challenges & opportunities

# This PIM Tutorial at ISCA 2024



## ISCA 2024 Memory-Centric Computing Systems Tutorial

Saturday, June 29, Buenos Aires, Argentina

**Organizers:** Geraldo F. Oliveira, Dr. Mohammad Sadrosadati, Ataberk Olgun, Professor Onur Mutlu

**Program:** https://events.safari.ethz.ch/isca24-memorycentric-tutorial/

Overview of PIM | PIM taxonomy

PIM in memory & storage

Real-world PNM systems

PUM for bulk bitwise operations

Programming techniques & tools

Infrastructures for PIM Research

Research challenges & opportunities

ISCA 2024
June 29 - July 3, 2024
Buenos Aires, Argentina

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

**https://events.safari.ethz.ch/isca24-memorycentric-tutorial**

# 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/**

*SAFARI*

# Open Source Tools: SAFARI GitHub



**SAFARI Research Group at ETH Zurich and Carnegie Mellon University**

Site for source code and tools distribution from SAFARI Research Group at ETH Zurich and Carnegie Mellon University.

241 followers · ETH Zurich and Carnegie Mellon U... · https://safari.ethz.ch/ · omutlu@gmail.com

Overview · Repositories 80 · Projects · Packages · People 13

Pinned · Customize pins

**ramulator** (Public)

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...

C++ · ☆ 442 · ⑂ 195

**prim-benchmarks** (Public)

PrIM (Processing-In-Memory benchmarks) is the first benchmark suite for a real-world processing-in-memory (PIM) architecture. PrIM is developed to evaluate, analyze, and characterize the first publ...

C · ☆ 100 · ⑂ 38

**MQSim** (Public)

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 implement...

C++ · ☆ 213 · ⑂ 120

**rowhammer** (Public)

Source code for testing the Row Hammer error mechanism in DRAM devices. Described in the ISCA 2014 paper by Kim et al. at http://users.ece.cmu.edu/~omutlu/pub/dram-row-hammer_isca14.pdf.

C · ☆ 208 · ⑂ 41

**SoftMC** (Public)

SoftMC is an experimental FPGA-based memory controller design that can be used to develop tests for DDR3 SODIMMs using a C++ based API. The design, the interface, and its capabilities and limitatio...

Verilog · ☆ 104 · ⑂ 26

**Pythia** (Public)

A customizable hardware prefetching framework using online reinforcement learning as described in the MICRO 2021 paper by Bera et al. (https://arxiv.org/pdf/2109.12021.pdf).

C++ · ☆ 85 · ⑂ 25

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

24

# Tutorial on
# Memory-Centric Computing:
## Conclusion Remarks

Geraldo F. Oliveira

Prof. Onur Mutlu

ISCA 2024

29 June 2024

**SAFARI**

**ETH** *zürich*