# Future of
# Computer Architecture and Hardware Security

Onur Mutlu

omutlu@gmail.com

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

6 March 2025

University of Southern California

**SAFARI**          **ETH**zürich

# Agenda

- **Computer Architecture Today**
  - What is it and where it is going

- **Three Major Hardware Issues That Affect Security**
  - Technology scaling problems
  - Growing system complexity; old methods not keeping up
  - New architectures and technologies

# Why Do We Do Computing?

# Answer

To Solve Problems

# Answer Reworded

# To Gain Insight

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

# Answer Extended

# To Enable
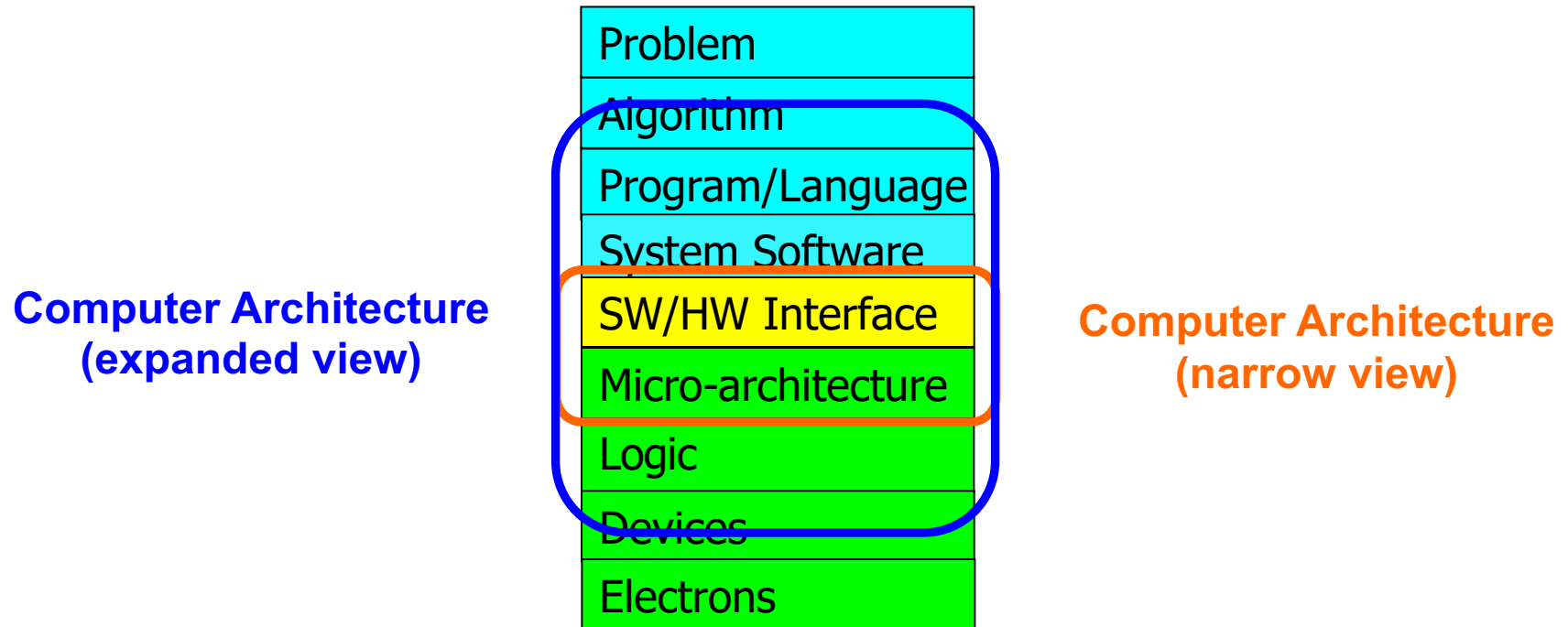# a Better Life & Future

# How Does a Computer Solve Problems?

# Answer

# Orchestrating Electrons

In today's dominant technologies

# How Do Problems
## Get Solved by Electrons?

# The Transformation Hierarchy

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

**Computer Architecture
(expanded view)**

**Computer Architecture
(narrow view)**

# Computer Architecture

- is the science and art of designing computing platforms (hardware, interface, system SW, and programming model)

- to achieve a set of design goals
    - E.g., highest performance on earth on workloads X, Y, Z
    - E.g., longest battery life at a form factor that fits in your pocket with cost < $$$ CHF
    - E.g., best average performance across all known workloads at the best performance/cost ratio
    - ...

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

**SAFARI**

# Different Platforms, Different Goals

**SAFARI**

Source: http://www.sia-online.org (semiconductor industry association)

# Different Platforms, Different Goals

Source: https://iq.intel.com/5-awesome-uses-for-drone-technology/

# Different Platforms, Different Goals

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

# Different Platforms, Different Goals

Source: http://sm.pcmag.com/pcmag_uk/photo/g/google-self-driving-car-the-guts/google-self-driving-car-the-guts_dwx8.jpg

# 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.

**Accelerating Genome Analysis: A Primer on an Ongoing Journey**
Sept.-Oct. 2020, pp. 65-75, vol. 40
DOI Bookmark: 10.1109/MM.2020.3013728

**FPGA-Based Near-Memory Acceleration of Modern Data-Intensive Applications**
July-Aug. 2021, pp. 39-48, vol. 41
DOI Bookmark: 10.1109/MM.2021.3088396

MinION from ONT

SmidgION from ONT

# An Example System in Your Pocket

Sensors

Storage

Main Memory

SoC
with lots of
compute
& caches

Main Memory

Storage

Apple M1 Ultra System (2022)

# Different Platforms, Different Goals

**SAFARI**

Source: http://datacentervoice.com/wp-content/uploads/2015/10/data-center.jpg

# Different Platforms, Different Goals

**SAFARI**

# Different Platforms, Different Goals

Source: https://www.itmagazine.ch/artikel/72401/Fugaku_Der_schnellste_Supercomputer_der_Welt.html

# Different Platforms, Different Goals



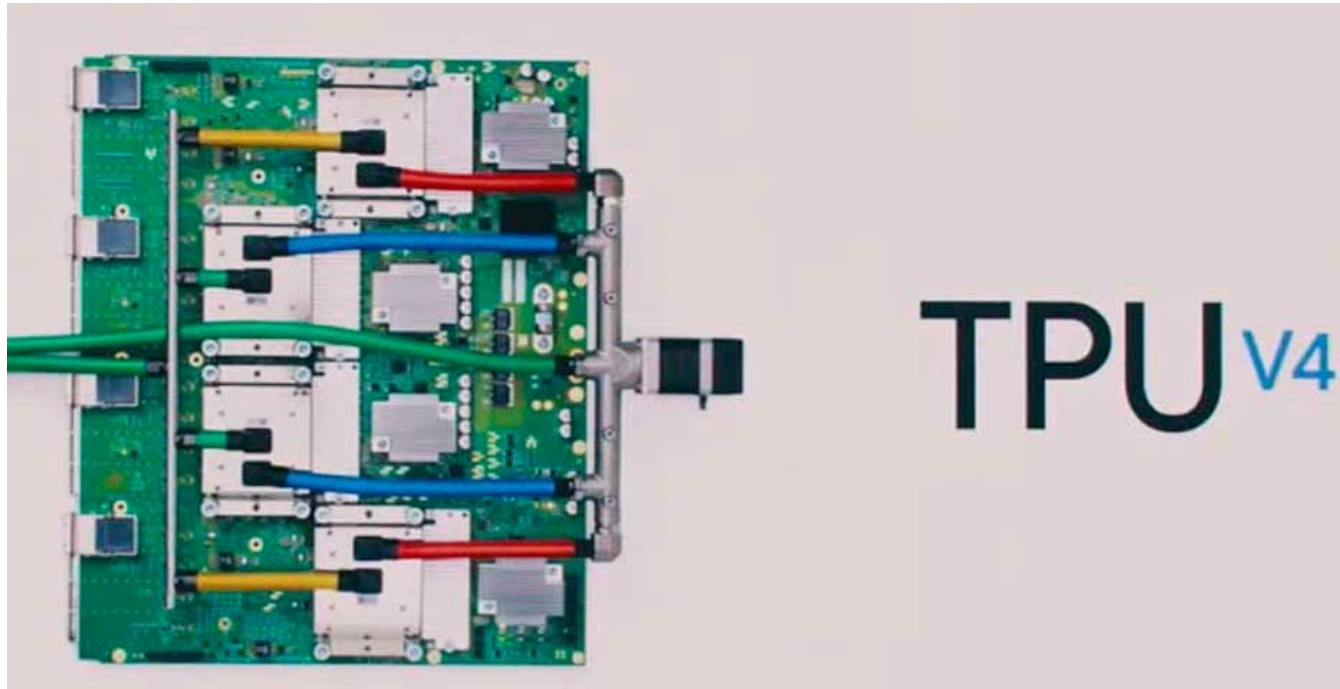**Figure 3.** TPU Printed Circuit Board. It can be inserted in the slot for an SATA disk in a server, but the card uses PCIe Gen3 x16.



**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.
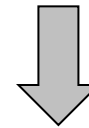
# Different Platforms, Different Goals



**New ML applications (vs. TPU3):**

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

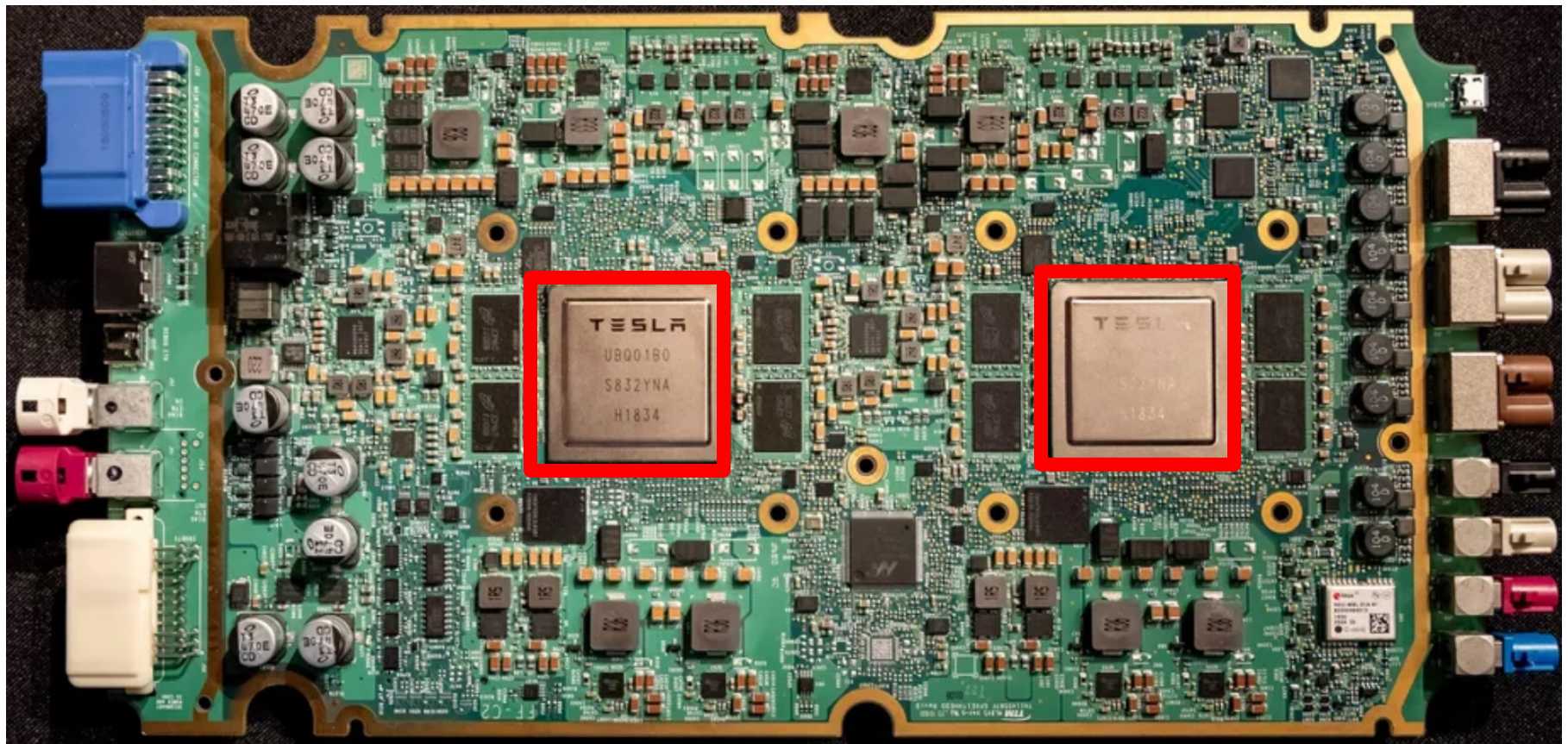**250 TFLOPS per chip in 2021 vs 90 TFLOPS in TPU3**

**1 ExaFLOPS per board**

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

# Different Platforms, Different Goals

- ML accelerator: 260 mm$^2$, 6 billion transistors, 600 GFLOPS GPU, 12 ARM 2.2 GHz CPUs.
- Two redundant chips for better safety.

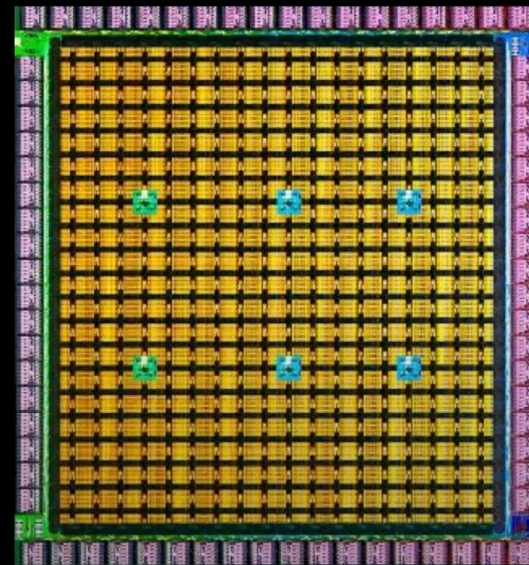# Different Platforms, Different Goals

- **Tesla Dojo Chip & System**

# Different Platforms, Different Goals

- Tesla Dojo Chip & System
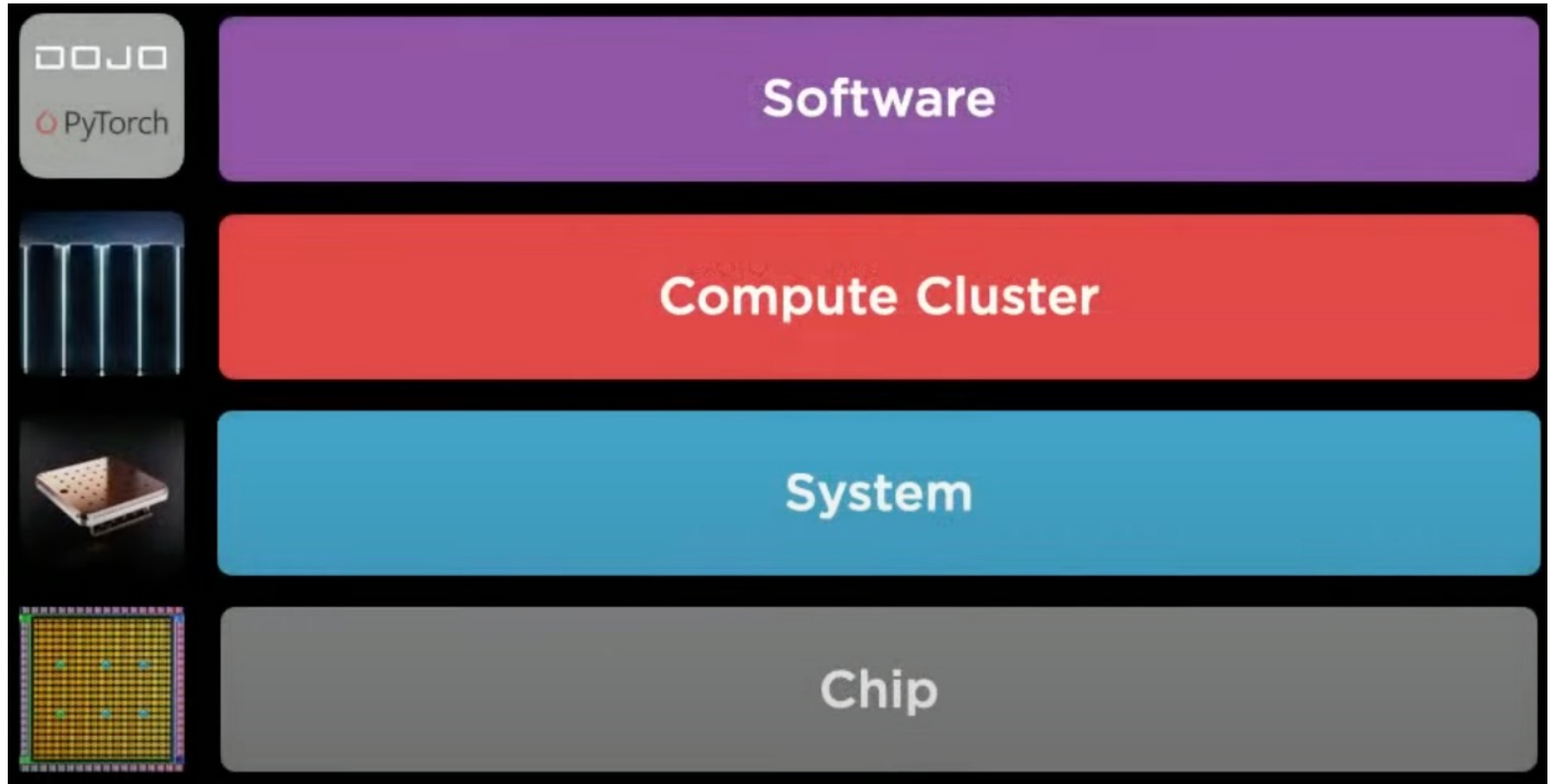
# Different Platforms, Different Goals
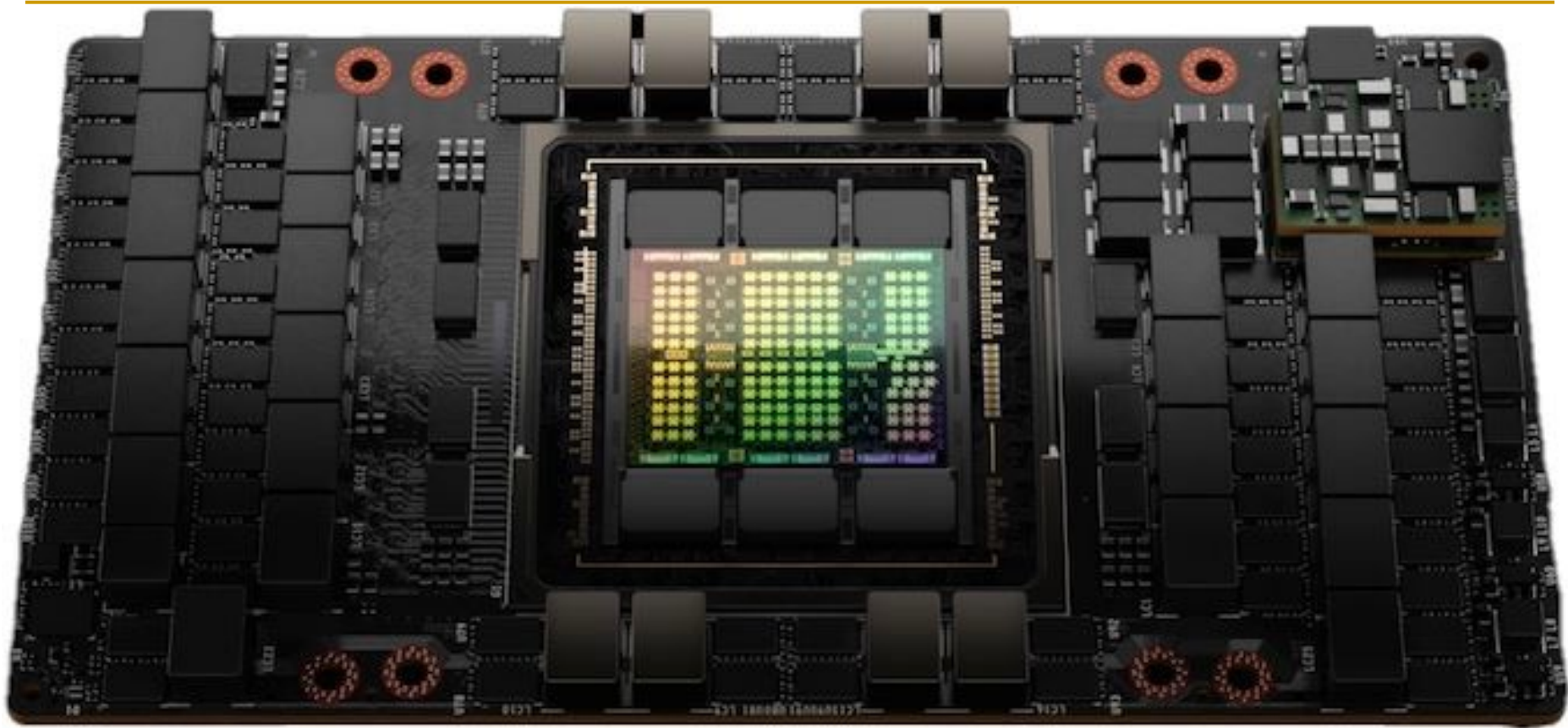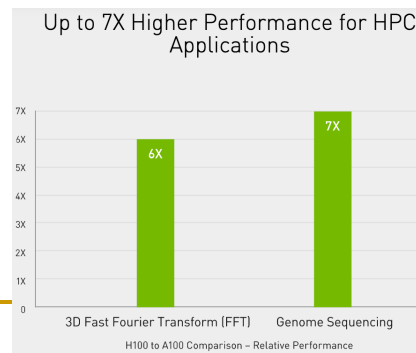
- Tesla Dojo Chip & System

# Different Platforms, Different Goals



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



Up to 7X Higher Performance for HPC Applications

H100 to A100 Comparison – Relative Performance

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

# Evolution of Recent GPUs (I)

**Volta**

>21 billion transistors
815mm^2
TSMC 12nm FFN

**Ampere**

>54 billion transistors
826 mm^2
TSMC N7

**Hopper**

>80 billion transistors
814 mm^2
TSMC 4N

**Blackwell**

>208 billion transistors
>1600 mm^2
TSMC 4NP

**Ampere | NVLink3**
12 NVLinks | 50GB/s each
x4@50Gbps-NRZ
600GB/s total

**Hopper | NVLink4**
18 NVLinks | 50GB/s each
x2@100Gbps-PAM4
900GB/s total

**Blackwell | NVLink5**
18 NVLinks | 100GB/s each
x2@200Gbps-PAM4
1800GB/s total

# Multiple GPUs to Tackle Large Workloads

## AI Models Growing Exponentially
### Need for multi-GPU inference at scale



Model Size

**New Capabilities | Trillions of Parameters | 70,000X Growth in a Decade**

# Evolution of Recent GPUs (II)



**2016**
Hybrid Cube Mesh NVLink technology

**2022**
3rd Gen NVLink Switch
All-to-all connection among NVLink domain of 8 GPU

**2024**
4th Gen NVLink Switch Chip
All-to-all connection among NVLink domain of 72 GPU

# Cerebras's Wafer Scale Engine (2019)



- The largest ML accelerator chip

- 400,000 cores

**Cerebras WSE**
1.2 Trillion transistors
46,225 mm$^2$

**Largest GPU**
21.1 Billion transistors
815 mm$^2$
**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 mm$^2$

**Largest GPU**
54.2 Billion transistors
826 mm$^2$

**NVIDIA** Ampere GA100

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

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

# Cerebras's Wafer Scale Engine-3 (2023)



**Cerebras Wafer-Scale Engine**

**The largest chip ever produced**

**46,225 mm²** silicon

**4 trillion** transistors

**900,000** AI cores

**125 Petaflops** of AI compute

**44 Gigabytes** of on-chip memory

**21 PByte/s** memory bandwidth

**214 Pbit/s** fabric bandwidth

**5nm** TSMC process

Wafer-Scale AI: Enabling Unprecedented AI Compute Performance, Hot Chips 2024

# Many (Other) AI/ML Chips (2021)



AI Chip Landscape

V0.7 Dec., 2019

S.T.

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

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.

# Axiom

To achieve the highest efficiency, performance, robustness:

## we must take the expanded view
### of computer architecture

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

**Co-design across the hierarchy:**
**Algorithms to devices**

**Specialize as much as possible**
**within the design goals**

# What Limits Us
## in Computing Today?

# Increasingly Demanding Applications

# Dream…

# and, they will come

As applications push boundaries, computing platforms become increasingly strained

# Many Metrics to Optimize for

- Performance

- Energy/Power

- Correctness

- Robustness (Safety, Security, Reliability, Availability)

- Cost

- Programming Ease

- Usability (Ease of Use)

- Scalability

- Simplicity (Complexity)

- Privacy

- …

Challenging especially with complex systems & hardware

# Three Major Limiters to Computing

- Technology scaling is not going well

- System complexity is increasing; old methods not keeping up

- Processor-centric designs are not keeping up

- These affect all metrics we care about

- These have fundamental impact on security and how we build secure systems

# Technology Scaling

# Technology Scaling Problems

- Circuit size and energy reduction has enabled continuous innovation at all levels of the computing stack

- As circuits become smaller, they become less reliable

- More flaky circuits are a problem for robust (reliable, safe, secure) operation

- If circuits produce wrong results, security can be affected (along with safety, reliability, availability)

**SAFARI**

# 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**

SAFARI

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

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

# How Robust Are **These** Platforms?

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

# Robust
# (Reliable, Secure, Safe)

# An Example: The RowHammer Problem

- One can predictably induce bit flips in commodity DRAM chips
  - All recent DRAM chips are fundamentally vulnerable

- First example of how a simple hardware failure mechanism can create a widespread system security vulnerability

**WIRED**                                    Forget Software—Now Hackers Are Exploiting Physics

| BUSINESS | CULTURE | DESIGN | GEAR | SCIENCE |

ANDY GREENBERG    SECURITY    08.31.16    7:00 AM

SHARE

f  SHARE
   18276

🐦 TWEET

# FORGET SOFTWARE—NOW HACKERS ARE EXPLOITING PHYSICS

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

# One can
# predictably induce errors
# in DRAM memory chips

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

Rowhammer

# Modern Memory is Prone to Disturbance Errors



| Row of Cells | Wordline |
|---|---|
| Victim Row | |
| Hammered~~Opened~~ Row | $V_{HIGH}$ ~~$V_{LOW}$~~ |
| Victim Row | |
| Row | |

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

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

# Recent DRAM Is More Vulnerable



*All modules from 2012–2013 are vulnerable*

# Higher-Level Implications

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

| |
|---|
| Problem |
| Algorithm |
| Program/Language |
| Runtime System (VM, OS, MM) |
| ISA (Architecture) |
| Microarchitecture |
| Logic |
| Devices |
| Electrons |

| |
|---|
| Problem |
| Algorithm |
| Program/Language |

| |
|---|
| Runtime System (VM, OS, MM) |
| ISA |
| Microarchitecture |
| Logic |
| Devices |
| Electrons |

User

# A Simple Program Can Induce Many Errors



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

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

# A Simple Program Can Induce Many Errors

CPU ⟷ DRAM Module

1. Avoid *cache hits*
   – Flush **X** from cache

2. Avoid *row hits* to **X**
   – Read **Y** in another row

**X** →

**Y** →

# A Simple Program Can Induce Many Errors



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

# A Simple Program Can Induce Many Errors



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

# A Simple Program Can Induce Many Errors



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

# Observed Errors in Real Systems

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

## A real reliability, security, safety issue

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

# One Can Take Over an Otherwise-Secure System

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

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

Project Zero

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

# Many RowHammer Security Exploits

- One can exploit RowHammer to

- Take over a system

- Read data they do not have access to

- Break out of virtual machine sandboxes

- Corrupt important data → render ML inference useless

- Steal secret data (e.g., crypto keys & ML model parameters)

# Security Implications

# Security Implications



Rowhammer

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

# Infrastructures to Understand Such Issues



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)

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)

# Infrastructures to Understand Such Issues

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

SAFARI

# 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: Open Source DRAM Infrastructure

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

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

Hasan Hassan[1,2,3]      Nandita Vijaykumar[3]      Samira Khan[4,3]      Saugata Ghose[3]      Kevin Chang[3]
Gennady Pekhimenko[5,3]      Donghyuk Lee[6,3]      Oguz Ergin[2]      Onur Mutlu[1,3]

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

# DRAM Bender: New DRAM Infrastructure

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

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

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

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

# DRAM Bender: FPGA Prototypes

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

## Five out of the box FPGA-based prototypes



Xilinx Alveo U200
FPGA Board
(with DRAM Bender)

DRAM
Module

PCI-e Connection
to the Host Machine

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

# HBM2 DRAM Testing Infrastructure

DRAM Bender on a Bittware XUPVVH



**(3) Cooling Fan**
**(4) Temperature Controller**
**(1) FPGA Board with HBM2**
**(2) Heating Pad**
**(5) PCIe Host Interface**

Fine-grained control over DRAM commands,
timing parameters (±1.67ns), and temperature (±0.5°C )

*Olgun et al., **"DRAM Bender: An Extensible and Versatile FPGA-Based Infrastructure to Easily Test State-of-the-Art DRAM Chips,"** TCAD, 2023. [GitHub: https://github.com/CMU-SAFARI/DRAM-Bender]

# RowHammer [ISCA 2014]

- Yoongu Kim, Ross Daly, Jeremie Kim, Chris Fallin, Ji Hye Lee, Donghyuk Lee, Chris Wilkerson, Konrad Lai, and Onur Mutlu,
  **"Flipping Bits in Memory Without Accessing Them: An Experimental Study of DRAM Disturbance Errors"**
  *Proceedings of the 41st International Symposium on Computer Architecture* (**ISCA**), Minneapolis, MN, June 2014.
  [Slides (pptx) (pdf)] [Lightning Session Slides (pptx) (pdf)] [Source Code and Data] [Lecture Video (1 hr 49 mins), 25 September 2020]
  *One of the 7 papers of 2012-2017 selected as Top Picks in Hardware and Embedded Security for IEEE TCAD (link). Selected to the ISCA-50 25-Year Retrospective Issue covering 1996-2020 in 2023 (Retrospective (pdf) Full Issue). Winner of the 2024 IFIP Jean-Claude Laprie Award in dependable computing (link).*

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

Yoongu Kim[1]    Ross Daly*    Jeremie Kim[1]    Chris Fallin*    Ji Hye Lee[1]
Donghyuk Lee[1]    Chris Wilkerson[2]    Konrad Lai    Onur Mutlu[1]

[1]Carnegie Mellon University    [2]Intel Labs

# Memory Scaling Issues **Are** Real

- Onur Mutlu and Jeremie Kim,
  **"RowHammer: A Retrospective"**
  *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems* (**TCAD**) *Special Issue on Top Picks in Hardware and Embedded Security*, 2019.
  [Preliminary arXiv version]
  [Slides from COSADE 2019 (pptx)]
  [Slides from VLSI-SOC 2020 (pptx) (pdf)]
  [Talk Video (1 hr 15 minutes, with Q&A)]

## RowHammer: A Retrospective

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

# Memory Scaling Issues **Are** Real

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

## Fundamentally Understanding and Solving RowHammer

Onur Mutlu
onur.mutlu@safari.ethz.ch
ETH Zürich
Zürich, Switzerland

Ataberk Olgun
ataberk.olgun@safari.ethz.ch
ETH Zürich
Zürich, Switzerland

A. Giray Yağlıkcı
giray.yaglikci@safari.ethz.ch
ETH Zürich
Zürich, Switzerland

# A Recent PhD Thesis

- A. Giray Yaglikci, **"Enabling Efficient and Scalable DRAM Read Disturbance Mitigation via New Experimental Insights into Modern DRAM Chips,"** PhD Thesis, ETH Zürich, 2024.
  [Slides (pdf) (pptx)]
  [Thesis arXiv (abs) (pdf)]
  [SAFARI News]

## ENABLING EFFICIENT AND SCALABLE DRAM READ DISTURBANCE MITIGATION VIA NEW EXPERIMENTAL INSIGHTS INTO MODERN DRAM CHIPS

### ABDULLAH GİRAY YAĞLIKÇI

**https://arxiv.org/pdf/2408.15044.pdf**

# Main Memory Needs
# Intelligent Controllers

# Industry's Intelligent DRAM Controllers (I)

**ISSCC 2023 / SESSION 28 / HIGH-DENSITY MEMORIES**

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

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

SK hynix Semiconductor, Icheon, Korea

ISSCC 2023

**2023 International Solid-State Circuits Conference**

February 19-23, 2023 San Francisco, CA

SAFARI

# Industry's Intelligent DRAM Controllers (II)

SK hynix Semiconductor, Icheon, Korea

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

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

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

DRAM Design Team, Memory Division, Samsung Electronics

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

# A Solution from Microsoft

## Panopticon: A Complete In-DRAM Rowhammer Mitigation

Tanj Bennett[§], Stefan Saroiu, Alec Wolman, and Lucian Cojocar

Microsoft, [§]Avant-Gray LLC

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

# Recent Improvements in JEDEC (2024)

**JEDEC**
*Global Standards for the Microelectronics Industry*

| STANDARDS & DOCUMENTS | COMMITTEES | NEWS | EVENTS & MEETINGS | JOIN |
|---|---|---|---|---|

**DDR5 SDRAM**                                       JESD79-5C                    Apr 2024

Release Number: Version 1.30

Version 1.30

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

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

# Evaluation of Industry's Recent Solutions

- **Appears at DRAMSec 2024**

## Understanding the Security Benefits and Overheads of Emerging Industry Solutions to DRAM Read Disturbance

Oğuzhan Canpolat[§†]    A. Giray Yağlıkçı[§]    Geraldo F. Oliveira[§]    Ataberk Olgun[§]

Oğuz Ergin[†]    Onur Mutlu[§]

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

https://arxiv.org/pdf/2406.19094

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

# Evaluation of Industry's Recent Solutions

■ Oguzhan Canpolat, Abdullah Giray Yaglikci, Geraldo Francisco de Oliveira, Ataberk Olgun, Nisa Bostanci, Ismail Emir Yuksel, Haocong Luo, Oguz Ergin, and Onur Mutlu,
**"Chronus: Understanding and Securing the Cutting-Edge Industry Solutions to DRAM Read Disturbance"**
*Proceedings of the 31st International Symposium on High-Performance Computer Architecture (**HPCA**)*, Las Vegas, NV, USA, March 2025.
[Chronus Source Code (Officially Artifact Evaluated with All Badges)]
***Officially artifact evaluated as available, functional, and reproduced.***

2025 IEEE International Symposium on High-Performance Computer Architecture (HPCA)

# Chronus: Understanding and Securing the Cutting-Edge Industry Solutions to DRAM Read Disturbance

Oğuzhan Canpolat[§†]     A. Giray Yağlıkçı[§]     Geraldo F. Oliveira[§]     Ataberk Olgun[§]
Nisa Bostancı[§]     Ismail Emir Yuksel[§]     Haocong Luo[§]     Oğuz Ergin[‡†]     Onur Mutlu[§]
[§]*ETH Zürich*     [†]*TOBB University of Economics and Technology*     [‡]*University of Sharjah*

**https://arxiv.org/pdf/2502.12650**

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

84

# Are Solutions Good & Secure?

# Question

Are we now
BitFlip-free
in 2024 and Beyond?

# Are We Now BitFlip Free?

- **Appears at ISCA 2023**

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

## RowPress: Amplifying Read-Disturbance
## in Modern DRAM Chips

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

*ETH Zürich*

# RowPress [ISCA 2023]

- Haocong Luo, Ataberk Olgun, Giray Yaglikci, Yahya Can Tugrul, Steve Rhyner, M. Banu Cavlak, Joel Lindegger, Mohammad Sadrosadati, and Onur Mutlu,
**"RowPress: Amplifying Read Disturbance in Modern DRAM Chips"**
*Proceedings of the 50th International Symposium on Computer Architecture* (**ISCA**), Orlando, FL, USA, June 2023.
[Slides (pptx) (pdf)]
[Lightning Talk Slides (pptx) (pdf)]
[Lightning Talk Video (3 minutes)]
[RowPress Source Code and Datasets (Officially Artifact Evaluated with All Badges)]
*Officially artifact evaluated as available, reusable and reproducible. Best artifact award at ISCA 2023. IEEE Micro Top Pick in 2024.*

# RowPress: Amplifying Read-Disturbance in Modern DRAM Chips

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

*ETH Zürich*

# RowPress vs. RowHammer

Instead of using a high activation count,
☞ increase the time that the aggressor row stays open

**RowHammer Aggressor Row**
Open
Close
**36ns, 47K activations to induce bitflips**

**RowPress Aggressor Row**
Open
Close
**7.8μs, only 5K activations to induce bitflips**

We observe bitflips even with **ONLY ONE activation** in extreme cases where the row stays open for 30ms

**SAFARI**

## RowPress Amplifies Read Disturbance in DRAM

- Reduces the minimum number of row activations needed to induce a bitflip (**ACmin**) by **1-2 orders of magnitude**

- In extreme cases, activating a row **only once** induces bitflips

**Intel Core i5-10400
(Comet Lake)**

**Samsung DDR4 Module
M378A2K43CB1-CTD
(Date Code: 20-10)
w/ TRR RowHammer Mitigation**

**Key Idea:** A proof-of-concept RowPress program keeps a DRAM row open for a longer period by **keeping on accessing different cache blocks in the row**

```
// Sync with Refresh and Loop Below
  for (k = 0; k < NUM_AGGR_ACTS; k++)
    for (j = 0; j < NUM_READS  j++) *AGGRESSOR1[j];
    for (j = 0; j < NUM_READS  j++) *AGGRESSOR2[j];
    for (j = 0; j < NUM_READS  j++)
      clflushopt(AGGRESSOR1[j]);
      clflushopt(AGGRESSOR2[j]);
    mfence();
  activate_dummy_rows();
```

**Number of Cache Blocks Accessed
Per Aggressor Row ACT
(NUM_READS=1 is Rowhammer)**

# Real-System Demonstration (II)

**On 1500 victim rows**



**Leveraging RowPress, our user-level program induces bitflips when RowHammer cannot**

# Combining RowHammer and RowPress

- **Appears at DSN Disrupt 2024**

## An Experimental Characterization of Combined RowHammer and RowPress Read Disturbance in Modern DRAM Chips

Haocong Luo    İsmail Emir Yüksel    Ataberk Olgun    A. Giray Yağlıkçı
Mohammad Sadrosadati    Onur Mutlu
*ETH Zürich*

**https://arxiv.org/pdf/2406.13080v2**

# Combining RowHammer and RowPress

- **Appears at DIMVA 2024**

## Presshammer: Rowhammer and Rowpress without Physical Address Information

Jonas Juffinger[1], Sudheendra Raghav Neela[1], Martin Heckel[2], Lukas Schwarz[1], Florian Adamsky[2], and Daniel Gruss[1]

[1] Graz University of Technology, Graz, Austria
[2] Hof University of Applied Sciences, Hof, Germany

# Understanding RowPress

- **Appears in IEEE TED, 2024**

IEEE TRANSACTIONS ON ELECTRON DEVICES

## Unveiling RowPress in Sub-20 nm DRAM Through Comparative Analysis With Row Hammer: From Leakage Mechanisms to Key Features

Longda Zhou, Sheng Ye, Runsheng Wang, *Member, IEEE*, and Zhigang Ji

SAFARI

# Key Takeaways

Read disturbance is a technology scaling problem

Finding a good solution to read disturbance is difficult (and will become more so)

# More to Come…

# RowHammer Becomes Worse with Aging

Preliminary data on aging via 68-day of continuous hammering

Aging can lead to read disturbance bitflips at **smaller** hammer counts

# RowHammer (Spatial Variation) Analysis (2024)

- **Appears at HPCA 2024**

## Spatial Variation-Aware Read Disturbance Defenses: Experimental Analysis of Real DRAM Chips and Implications on Future Solutions

Abdullah Giray Yağlıkçı    Yahya Can Tuğrul    Geraldo F. Oliveira
İsmail Emir Yüksel    Ataberk Olgun    Haocong Luo    Onur Mutlu
ETH Zürich

**https://arxiv.org/pdf/2402.18652**

# Variable Read Disturbance (2025)

## Key Takeaway

The Read Disturbance Threshold (RDT) of a row
changes randomly and unpredictably over time

Accurately identifying RDT is challenging

# Variable Read Disturbance (2025)

- **Appears at HPCA 2025**

**Variable Read Disturbance:**
**An Experimental Analysis of Temporal Variation in DRAM Read Disturbance**

Ataberk Olgun†    F. Nisa Bostancı†    İsmail Emir Yüksel†    Oğuzhan Canpolat†    Haocong Luo†
Geraldo F. Oliveira†    A. Giray Yağlıkçı†    Minesh Patel‡    Onur Mutlu†

*ETH Zurich*†        *Rutgers University*‡

# Two Major Directions

- **Understanding Bitflips (Hardware errors in general)**
  - Many effects on bitflips still need to be rigorously examined
    - Aging of DRAM Chips
    - Environmental Conditions (e.g., Process, Voltage, Temperature)
    - Memory Access Patterns
    - Memory Controller & System Design Decisions
    - …

- **Solving Bitflips (Hardware errors in general)**
  - Flexible and efficient solutions are necessary
    - In-field patchable / reconfigurable / programmable solutions
  - Co-architecting across the system stack/components is important
    - To avoid performance and denial-of-service problems

# A Recent RowHammer Lecture



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

Stanford Online ✓
529K subscribers

🔔 Subscribed ∨   👍 35  👎   ↪ Share   ↓ Download   ≡+ Save   ...

1.2K views  1 month ago

# Emerging Memories Also Need Intelligent Controllers

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

## Architecting Phase Change Memory as a Scalable DRAM Alternative

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

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

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

# Intelligent Memory Controllers

# Can Enhance Security & Enable Better Scaling

# Read Disturbance Sessions @ HPCA 2025

**Session 7A (*Acacia A and B*): Hammering the Odds – 1**

Session Chair: *Gururaj Saileshwar (Toronto)*

- **Variable Read Disturbance: An Experimental Analysis of Temporal Variation in DRAM Read Disturbance**
  Ataberk Olgun (ETH Zürich), Nisa Bostanci (ETH Zürich), Ismail Emir Yuksel (ETH Zürich), Giray Yaglikci (ETH Zürich), Geraldo F. Oliveira (ETH Zürich), Haocong Luo (ETH Zürich), Oguzhan Canpolat (ETH Zürich), Minesh Patel (Rutgers University), Onur Mutlu (ETH Zürich)
- **Understanding RowHammer Under Reduced Refresh Latency: Experimental Analysis of Real DRAM Chips and Implications on Future Solutions**
  Yahya Can Tuğrul (TOBB ETÜ & ETH Zürich), Giray Yaglikci (ETH Zürich), Ismail Emir Yuksel (ETH Zürich), Ataberk Olgun (ETH Zürich), Oğuzhan Canpolat (TOBB ETÜ & ETH Zürich), Nisa Bostanci (ETH Zürich), Mohammad Sadrosadati (ETH Zürich), Oguz Ergin (TOBB ETÜ), Onur Mutlu (ETH Zürich)
- **Chronus: Understanding and Securing the Cutting-Edge Industry Solutions to DRAM Read Disturbance**
  Oğuzhan Canpolat (TOBB ETÜ & ETH Zürich), Giray Yaglikci (ETH Zürich), Geraldo Francisco de Oliveira (ETH Zürich), Ataberk Olgun (ETH Zürich), Nisa Bostanci (ETH Zürich), Ismail Emir Yuksel (ETH Zürich), Haocong Luo (ETH Zürich), Oğuz Ergin (TOBB ETÜ), Onur Mutlu (ETH Zürich)

**Session 8A (*Acacia A and B*): Hammering the Odds – 2**

Session Chair: *Sudhanva Gurumurthi (AMD)*

- **AutoRFM: Scaling Low-Cost In-DRAM Trackers to Ultra-Low Rowhammer Thresholds**
  Moinuddin Qureshi (Georgia Tech)
- **DAPPER: A Performance-Attack-Resilient Tracker for RowHammer Defense**
  Jeonghyun Woo (The University of British Columbia (UBC)), Prashant J. Nair (The University of British Columbia (UBC))
- **QPRAC: Towards Secure and Practical PRAC-based Rowhammer Mitigation using Priority Queues**
  Jeonghyun Woo (The University of British Columbia (UBC)), Shaopeng (Chris) Lin (University of Toronto), Prashant J. Nair (The University of British Columbia (UBC)), Aamer Jaleel (NVIDIA), Gururaj Saileshwar (University of Toronto)

**Tuesday, March 4th, 11am and 2pm**

# Data Corruption is in CPU Logic, Too

- Intermittent defects can cause silent data corruption

- They may be hard to detect or replicate

- They may be exploitable

## Silent Data Corruptions at Scale

Harish Dattatraya Dixit
Facebook, Inc.
hdd@fb.com

Sneha Pendharkar
Facebook, Inc.
spendharkar@fb.com

Matt Beadon
Facebook, Inc.
mbeadon@fb.com

Chris Mason
Facebook, Inc.
clm@fb.com

Tejasvi Chakravarthy
Facebook, Inc.
teju@fb.com

Bharath Muthiah
Facebook, Inc.
bharathm@fb.com

Sriram Sankar
Facebook Inc.
sriramsankar@fb.com

## Cores that don't count

Peter H. Hochschild
Paul Turner
Jeffrey C. Mogul
Google
Sunnyvale, CA, US

Rama Govindaraju
Parthasarathy
Ranganathan
Google
Sunnyvale, CA, US

David E. Culler
Amin Vahdat
Google
Sunnyvale, CA, US

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

# Silent Data Corruption In-the-Field (2021)



We have a *new* problem: cores that disobey instructions

CPU cores that
- repeatedly
- but not always
- mis-calculate
- certain computations
- without giving any obvious signal

"Mercurial cores" committing
"Corrupt Execution Errors"

Due to local silicon defects, not eg cosmic rays

Google

0:19 / 9:14 · We have a new probleme cores that disobey instructions

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

## Understanding Silent Data Corruptions in a Large Production CPU Population

Shaobu Wang
Tsinghua University

Guangyan Zhang*
Tsinghua University

Junyu Wei
Tsinghua University

Yang Wang
The Ohio State University

Jiesheng Wu
Alibaba Cloud

Qingchao Luo
Alibaba Cloud

## Understanding and Mitigating Hardware Failures in Deep Learning Training Accelerator Systems

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

Mike Hutton
Google
Sunnyvale, CA, USA
mdhutton@google.com

Steven Chan
Google
Sunnyvale, CA, USA
scchan@google.com

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

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

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

Yanjing Li
University of Chicago
Chciago, IL, USA
yanjingl@uchicago.edu

# How to Detect Hardware Errors? (I)

- Kypros Constantinides, Onur Mutlu, Todd Austin, and Valeria Bertacco, **"Software-Based Online Detection of Hardware Defects: Mechanisms, Architectural Support, and Evaluation"** *Proceedings of the 40th International Symposium on Microarchitecture* (**MICRO**), pages 97-108, Chicago, IL, December 2007. Slides (ppt)

**Software-Based Online Detection of Hardware Defects:**
**Mechanisms, Architectural Support, and Evaluation**

Kypros Constantinides‡        Onur Mutlu†        Todd Austin‡        Valeria Bertacco‡

‡ Advanced Computer Architecture Lab
University of Michigan
Ann Arbor, MI
{kypros, austin, valeria}@umich.edu

† Computer Architecture Group
Microsoft Research
Redmond, WA
onur@microsoft.com

# How to Detect Hardware Errors? (II)

- Kypros Constantinides, Onur Mutlu, and Todd Austin,
  **"Online Design Bug Detection: RTL Analysis, Flexible Mechanisms, and Evaluation"**
  *Proceedings of the 41st International Symposium on Microarchitecture* (**MICRO**), pages 282-293, Lake Como, Italy, November 2008. Slides (ppt)

**Online Design Bug Detection: RTL Analysis, Flexible Mechanisms, and Evaluation**

Kypros Constantinides‡    Onur Mutlu§    Todd Austin‡

‡Advanced Computer Architecture Lab
University of Michigan
{kypros, austin}@umich.edu

§Microsoft Research and Carnegie Mellon University
onur@{microsoft.com,cmu.edu}

# How to Detect Hardware Errors? (III)

- Yanjing Li, Onur Mutlu, and Subhasish Mitra,
  **"Operating System Scheduling for Efficient Online Self-Test in Robust Systems"**
  *Proceedings of the* International Conference on Computer-Aided Design (**ICCAD**), pages 201-208, San Jose, CA, November 2009. Slides (ppt) (pdf)

**Operating System Scheduling for Efficient Online Self-Test in Robust Systems**

| Yanjing Li | Onur Mutlu | Subhasish Mitra |
| --- | --- | --- |
| Stanford University | Carnegie Mellon University | Stanford University |

# How to Detect Hardware Errors? (IV)

- Yanjing Li, Onur Mutlu, Donald S. Gardner, and Subhasish Mitra,
  **"Concurrent Autonomous Self-Test for Uncore Components in System-on-Chips"**
  *Proceedings of the 28th IEEE VLSI Test Symposium* (**VTS**), pages 232-237, Santa Cruz, CA, April 2010. Slides (ppt)
  **Best paper award at VTS 2010.**

**Concurrent Autonomous Self-Test for Uncore Components in System-on-Chips**

Yanjing Li
Stanford University

Onur Mutlu
Carnegie Mellon University

Donald S. Gardner
Intel Corporation

Subhasish Mitra
Stanford University

# How to Detect Hardware Errors? (V)

- Kypros Constantinides, Onur Mutlu, Todd Austin, and Valeria Bertacco, **"A Flexible Software-Based Framework for Online Detection of Hardware Defects"** *IEEE Transactions on Computers* (**TC**), Vol. 58, No. 8, pages 1063-1079, August 2009.

## A Flexible Software-Based Framework for Online Detection of Hardware Defects

Kypros Constantinides, *Student Member, IEEE*, Onur Mutlu, *Member, IEEE*, Todd Austin, *Member, IEEE*, and Valeria Bertacco, *Member, IEEE*

# Takeaways

- **Both memory and logic errors will become worse with technology scaling**

- Hardware errors will create worse robustness problems

- **We cannot afford to ignore data corruption**

# System Complexity

# Complex Systems Cause Many Issues

- **Many hardware components, complex components**

- **Harder to design & verify**

- **Harder to reason about** operational behavior
  - Correctness, performance, energy, security, privacy, …

- **Harder to control interactions** between components and avoid **information leakage**

- **Old methods do not keep up** with new trends and complexity
  - Virtual memory a prime example, also coherence & verification

# Processor Complexity Is Growing



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

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

Our World in Data

SAFARI

119

# Complex CPUs and Memory Hierarchies



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

Intel Alder Lake, 2021

*SAFARI*

# Complex CPUs and Memory Hierarchies



**Core Count:**
8 cores/16 threads

**L1 Caches:**
32 KB per core

**L2 Caches:**
512 KB per core

**L3 Cache:**
32 MB shared

AMD Ryzen 5000, 2020

# Complexity Growing with 3D (2021)

**Zen 3 Layout**

CCD

| CPU Core | 32MB L3 Cache | CPU Core |
| CPU Core | | CPU Core |
| CPU Core | | CPU Core |
| CPU Core | | CPU Core |

https://community.microcenter.com/discussion/5134/comparing-zen-3-to-zen-2

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

**Additional 64 MB L3 cache die stacked on top of the processor die**
- Connected using Through Silicon Vias (TSVs)
- Total of 96 MB L3 cache



Structural silicon

64MB L3 cache die

Direct copper-to-copper bond

Through Silicon Vias (TSVs) for silicon-to-silicon communication

Up to 8-core "Zen 3" CCD

# Processor Complexity and Features

- Leads to <span style="color:red">many (endless) side and covert channels</span>
    - Spectre and Meltdown are prime recent examples
    - These will not go away

- Leads to <span style="color:red">many bugs and unintended behavior</span>
    - Especially with new features or complex interactions
    - Some can be exploitable

- **How to tame processor complexity & resulting issues?**

# Access Control & Protection Mechanisms

- Are based on virtual memory (VM), invented in 1950s

- VM has not changed much even after decades of technology scaling and memory system improvements

- VM causes large performance problems and is responsible for large complexity, power, energy

- VM is poor for fine-grained security and access control

- VM hinders innovation in heterogeneous (e.g., accelerator) systems and new architectures (e.g., processing near data)

- **It is time to rethink virtual memory**

# Virtual Memory: Parting Thoughts

- Virtual Memory is one of the most successful examples of
  - architectural support for programmers
  - how to partition work between hardware and software
  - hardware/software cooperation
  - programmer/architect tradeoff

- Going forward: How does virtual memory fare and scale into the future? Five key trends:
  - Increasing, huge physical memory sizes (local & remote)
  - Hybrid physical memory systems (DRAM + NVM + SSD)
  - Many accelerators in the system accessing physical memory
  - Virtualized systems (hypervisors, software virtualization, local and remote memories)
  - Processing in memory systems – near-data accelerators

# Rethinking Virtual Memory

Nastaran Hajinazar, Pratyush Patel, Minesh Patel, Konstantinos Kanellopoulos, Saugata Ghose, Rachata Ausavarungnirun, Geraldo Francisco de Oliveira Jr., Jonathan Appavoo, Vivek Seshadri, and Onur Mutlu,
**"The Virtual Block Interface: A Flexible Alternative to the Conventional Virtual Memory Framework"**
*Proceedings of the 47th International Symposium on Computer Architecture* (**ISCA**), Virtual, June 2020.
[Slides (pptx) (pdf)]
[Lightning Talk Slides (pptx) (pdf)]
[ARM Research Summit Poster (pptx) (pdf)]
[Talk Video (26 minutes)]
[Lightning Talk Video (3 minutes)]
[Lecture Video (43 minutes)]

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

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

[*]ETH Zürich   [†]Simon Fraser University   [⋈]University of Washington   [‡]Carnegie Mellon University
[⊙]King Mongkut's University of Technology North Bangkok   [◇]Boston University   [▽]Microsoft Research India

# Better Virtual Memory (I)

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

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

Konstantinos Kanellopoulos[1]    Hong Chul Nam[1]    F. Nisa Bostanci[1]    Rahul Bera[1]
Mohammad Sadrosadati[1]    Rakesh Kumar[2]    Davide Basilio Bartolini[3]    Onur Mutlu[1]

[1]ETH Zürich    [2]Norwegian University of Science and Technology    [3]Huawei Zurich Research Center

https://arxiv.org/pdf/2310.04158

# Better Virtual Memory (II)

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

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

Konstantinos Kanellopoulos[1]    Rahul Bera[1]    Kosta Stojiljkovic[1]    Nisa Bostanci[1]    Can Firtina[1]
Rachata Ausavarungnirun[2]    Rakesh Kumar[3]    Nastaran Hajinazar[4]    Mohammad Sadrosadati[1]
Nandita Vijaykumar[5]    Onur Mutlu[1]

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

# Even Better Virtual Memory

- **To Appear at ASPLOS 2025**

## Virtuoso: Enabling Fast and Accurate Virtual Memory Research via an Imitation-based Operating System Simulation Methodology

Konstantinos Kanellopoulos
ETH Zürich
Zürich, Switzerland

Konstantinos Sgouras
ETH Zürich
Zürich, Switzerland

F. Nisa Bostanci
ETH Zürich
Zürich, Switzerland

Andreas Kosmas Kakolyris
ETH Zürich
Zürich, Switzerland

Berkin Kerim Konar
ETH Zürich
Zürich, Switzerland

Rahul Bera
ETH Zürich
Zürich, Switzerland

Mohammad Sadrosadati
ETH Zürich
Zürich, Switzerland

Rakesh Kumar
NTNU
Trondheinm, Norway

Nandita Vijaykumar
University of Toronto
Toronto, Canada

Onur Mutlu
ETH Zürich
Zürich, Switzerland

# New Architectures & Technologies

# New Architectures & Technologies

- Can have large impact on security and robustness
  - Positive or negative

- They need to be designed with system security in mind
  - Ideally as a first-class design goal

- Multiple potentially paradigm-changing new technologies and architectures
  - Processing in memory
  - Accelerator-based computing
  - Quantum computing
  - …

# Processing in Memory

# Problem

## Computing
## is Bottlenecked by Data

# Data is Key for AI, ML, Genomics, …

- Important workloads are all data intensive

- They require rapid and efficient processing of large amounts of data

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

# Huge Demand for Performance & Efficiency

## Exponential Growth of Neural Networks

**Memory and compute requirements**



**1800x more compute**

In just **2 years**

**Tomorrow**, **multi-trillion** parameter models

**~4 orders of magnitude increase in memory requirement in just a few years!**

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

# Huge Demand for Performance & Efficiency



**development of new sequencing technologies**

Oxford Nanopore MinION

**Number of Genomes Sequenced**

229,000 — 2014
422,000 — 2015
952,000 — 2016
1,620,000 — 2017

Source: Illumina

# Do We Want This?

**SAFARI**

# Or This?

SAFARI    Source: V. Milutinovic

# High Performance, Energy Efficient, Sustainable

# (All at the Same Time)

# The Problem

Data access is the major performance and energy bottleneck

## Our current

## design principles

## cause great energy waste
(and great performance loss)

**SAFARI**

# Today's Computing Systems

- Processor centric

- All data processed in the processor → at great system cost

Computing System

# Processor-Centric System Performance

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



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

# Data Movement vs. Computation Energy



**Communication Dominates Arithmetic**

Dally, HiPEAC 2015

- 64-bit DP 20pJ
- 256-bit buses
- 256-bit access 8 kB SRAM
- 20mm
- 26 pJ
- 256 pJ
- 16 nJ — DRAM Rd/Wr
- 500 pJ — Efficient off-chip link
- 50 pJ

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

# Data Movement vs. Computation Energy

Han+, "EIE: Efficient Inference Engine on Compressed Deep Neural Network," ISCA 2016.

**SAFARI**

# Data Movement vs. Computation Energy



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

# Energy Waste in Mobile Devices

- Amirali Boroumand, Saugata Ghose, Youngsok Kim, Rachata Ausavarungnirun, Eric Shiu, Rahul Thakur, Daehyun Kim, Aki Kuusela, Allan Knies, Parthasarathy Ranganathan, and Onur Mutlu, **"Google Workloads for Consumer Devices: Mitigating Data Movement Bottlenecks"** *Proceedings of the 23rd International Conference on Architectural Support for Programming Languages and Operating Systems* (**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 Boroumand[1]    Saugata Ghose[1]    Youngsok Kim[2]

Rachata Ausavarungnirun[1]    Eric Shiu[3]    Rahul Thakur[3]    Daehyun Kim[4,3]

Aki Kuusela[3]    Allan Knies[3]    Parthasarathy Ranganathan[3]    Onur Mutlu[5,1]

# Energy Waste in Accelerators

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

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

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

Amirali Boroumand[†◇]     Saugata Ghose[‡]     Berkin Akin[§]     Ravi Narayanaswami[§]
Geraldo F. Oliveira[★]     Xiaoyu Ma[§]     Eric Shiu[§]     Onur Mutlu[★†]

[†]*Carnegie Mellon Univ.*     [◇]*Stanford Univ.*     [‡]*Univ. of Illinois Urbana-Champaign*     [§]*Google*     [★]*ETH Zürich*

# Energy Wasted on Data Movement



Legend: ■ Total Static  ■ PE  ■ Param Buffer+NoC  ■ Act Buffer+NoC  ■ Off-chip Interconnect  ■ DRAM

Chart x-axis groups (each with Baseline, Base+HB, Mensa): LSTM1, Transd.1, Transd.2, CNN5, CNN9, CNN10, CNN12, RCNN1, RCNN3, Average

Y-axis: Normalized Energy (0, 0.25, 0.5, 0.75, 1)

**In LSTMs and Transducers used by Google, >90% energy spent on off-chip interconnect and DRAM**

# Fundamental Problem

<span style="color:red">Processing</span> of data
is performed
<span style="color:red">far away from the data</span>

# We Need A **Paradigm Shift** To …

- Enable computation with <span style="color:blue">minimal data movement</span>

- <span style="color:blue">Compute where it makes sense</span> (<span style="color:red">where data resides</span>)

- Make computing architectures more <span style="color:blue">data-centric</span>

# Process Data Where It Makes Sense



Sensors

Storage | DRAM | A lot of SRAM | DRAM | Storage

Apple M1 Ultra System (2022)

# Memory as an Accelerator



**Memory similar to a "conventional" accelerator**

# Goal: Processing Inside Memory/Storage



- Many questions … How do we design the:
  - compute-capable memory & controllers?
  - processors & communication units?
  - software & hardware interfaces?
  - system software, compilers, languages?
  - algorithms & theoretical foundations?

# Processing in Memory: Two Types

1. Processing **near** Memory
2. Processing **using** Memory

# Processing-in-Memory Landscape Today



[Samsung 2021]

[Alibaba 2022]

[SK Hynix 2022]

[Samsung 2021]

[UPMEM 2019]

And, many other experimental chips and startups

# Processing-in-Memory Landscape Today

## Computational CXL-Memory Solution for Accelerating Memory-Intensive Applications

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

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

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



Fig. 6. FPGA prototype of proposed CMS card.

SAFARI

# Processing-in-Memory Landscape Today



Samsung Processing in Memory
Technology at Hot Chips 2023

By **Patrick Kennedy** - August 28, 2023

*Samsung PIM PNM For Transformer Based AI HC35_Page_24*

**SAFARI**

# Opportunity: 3D-Stacked Logic+Memory

**Memory**

**Logic**

Other "True 3D" technologies under development

# Tesseract System for Graph Processing

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

# Evaluated Systems



DDR3-OoO    HMC-OoO    HMC-MC    **Tesseract**

8 OoO 4GHz (×4, DDR3-OoO)

8 OoO 4GHz (×4, HMC-OoO)

128 In-Order 2GHz (×4, HMC-MC)

32 Tesseract Cores

102.4GB/s    640GB/s    640GB/s    **8TB/s**

# Tesseract Graph Processing Performance

**>13X Performance Improvement**



On five graph processing algorithms

Speedup (y-axis): 0, 2, 4, 6, 8, 10, 12, 14, 16

- DDR3-OoO
- HMC-OoO: +56%
- HMC-MC: +25%
- Tesseract: 9.0x
- Tesseract-LP: 11.6x
- Tesseract-LP-MTP: 13.8x

# Tesseract Graph Processing System Energy



Legend: Memory Layers, Logic Layers, Cores

Chart comparing HMC-OoO (1.0) with Tesseract with Prefetching (~0.13)

**> 8X Energy Reduction**

# 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 42nd International Symposium on Computer Architecture* (**ISCA**), Portland, OR, June 2015.
  [Slides (pptx) (pdf)] [Lightning Session Slides (pptx) (pdf)]
  ***Top Picks Honorable Mention by IEEE Micro.***
  ***Selected to the ISCA-50 25-Year Retrospective Issue covering 1996-2020 in 2023 (Retrospective (pdf) Full Issue).***

## A Scalable Processing-in-Memory Accelerator for Parallel Graph Processing

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

# A Short Retrospective @ 50 Years of ISCA

**Retrospective:** A Scalable Processing-in-Memory Accelerator for Parallel Graph Processing

Junwhan Ahn[†]    Sungpack Hong[‡]    Sungjoo Yoo[▽]    Onur Mutlu[§]    Kiyoung Choi[▽]
[†]Google DeepMind    [‡]Oracle Labs    [§]ETH Zürich    [▽]Seoul National University

*Abstract*—Our ISCA 2015 paper [1] provides a new programmable processing-in-memory (PIM) architecture and system design that can accelerate key data-intensive applications, with a focus on graph processing workloads. Our major idea was to completely rethink the system, including the programming model, data partitioning mechanisms, system support, instruction set architecture, along with near-memory execution units and their communication architecture, such that an important workload can be accelerated at a maximum level using a distributed system of well-connected near-memory accelerators. We built our accelerator system, Tesseract, using 3D-stacked memories with logic layers, where each logic layer contains general-purpose processing cores and cores communicate with each other using a message-passing programming model. Cores could be specialized for graph processing (or any other application to be accelerated).

To our knowledge, our paper was the first to completely design a near-memory accelerator system from scratch such that it is both generally programmable and specifically customizable to accelerate important applications, with a case study on major graph processing workloads. Ensuing work in academia and industry showed that similar approaches to system design can greatly benefit both graph processing workloads and other applications, such as machine learning, for which ideas from Tesseract seem to have been influential.

This short retrospective provides a brief analysis of our ISCA 2015 paper and its impact. We briefly describe the major ideas and contributions of the work, discuss later works that built on it or were influenced by it, and make some educated guesses on what the future may bring on PIM and accelerator systems.

## I. BACKGROUND, APPROACH & MINDSET

We started our research when 3D-stacked memories (e.g., [2–4]) were viable and seemed to have promise for building effective and practical processing-near-memory systems. Such near-memory processing could lead to improvements, but there was little to no research that examined how an accelerator could be completely (re-)designed using such near-memory technology, from its hardware architecture to its programming model and software system, and what the performance and energy benefits could be of such a re-design. We set out to answer these questions in our ISCA 2015 paper [1].

We followed several major principles to design our accelerator from the ground up. We believe these principles are still important: a major contribution and influence of our work was in putting all of these together in a cohesive full-system design and demonstrating the large performance and energy benefits that can be obtained from such a design. We see a similar approach in many modern large-scale accelerator systems in machine learning today (e.g., [5–9]). Our principles are:

1. *Near-memory execution* to enable/exploit the high data access bandwidth modern workloads (e.g., graph processing) need and to reduce data movement and access latency.

2. *General programmability* so that the system can be easily adopted, extended, and customized for many workloads.

3. *Maximal acceleration capability* to maximize the performance and energy benefits. We set ourselves free from backward compatibility and cost constraints. We aimed to completely re-design the system stack. Our goal was to explore the maximal performance and energy efficiency benefits we can gain from a near-memory accelerator if we had complete freedom to change things as much as we needed. We contrast this approach to the *minimal intrusion* approach we also explored in a separate ISCA 2015 paper [10].

4. *Customizable to specific workloads*, such that we can maximize acceleration benefits. Our focus workload was graph analytics/processing, a key workload at the time and today. However, our design principles are not limited to graph processing and the system we built is customizable to other workloads as well, e.g., machine learning, genome analysis.

5. *Memory-capacity-proportional performance*, i.e., processing capability should proportionally grow (i.e., scale) as memory capacity increases and vice versa. This enables scaling of data-intensive workloads that need both memory and compute.

6. *Exploit new technology (3D stacking)* that enables tight integration of memory and logic and helps multiple above principles (e.g., enables customizable near-memory acceleration capability in the logic layer of a 3D-stacked memory chip).

7. *Good communication and scaling capability* to support scalability to large dataset sizes and to enable memory-capacity-proportional performance. To this end, we provided scalable communication mechanisms between execution cores and carefully interconnected small accelerator chips to form a large distributed system of accelerator chips.

8. *Maximal and efficient use of memory bandwidth* to supply the high-bandwidth data access that modern workloads need. To this end, we introduced new, specialized mechanisms for prefetching and a programming model that helps leverage application semantics for hardware optimization.

## II. CONTRIBUTIONS AND INFLUENCE

We believe the major contributions of our work were 1) complete rethinking of how an accelerator system should be designed to enable maximal acceleration capability, and 2) the design and analysis of such an accelerator with this mindset and using the aforementioned principles to demonstrate its effectiveness in an important class of workloads.

One can find examples of our approach in modern large-scale machine learning (ML) accelerators, which are perhaps the most successful incarnation of scalable near-memory execution architectures. ML infrastructure today (e.g., [5–9]) consists of accelerator chips, each containing compute units and high-bandwidth memory tightly packaged together, and features scale-up capability enabled by connecting thousands of such chips with high-bandwidth interconnection links. The system-wide rethinking that was done to enable such accelerators and many of the principles used in such accelerators resemble our ISCA 2015 paper's approach.

The "memory-capacity-proportional performance" principle we explored in the paper shares similarities with how ML workloads are scaled up today. Similar to how we carefully sharded graphs across our accelerator chips to greatly improve effective memory bandwidth in our paper, today's ML workloads are sharded across a large number of accelerators by leveraging data/model parallelism and optimizing the placement to balance communication overheads and compute scalability [11, 12]. With the advent of large generative models requiring high memory bandwidth for fast training and inference, the scaling behavior where capacity and bandwidth are scaled together has become an essential architectural property to support modern data-intensive workloads.

The "maximal acceleration capability" principle we used in Tesseract provides much larger performance and energy improvements and better customization than the "minimalist" approach that our other ISCA 2015 paper on *PIM-Enabled Instructions* [10] explored: "minimally change" an existing system to incorporate (near-memory) acceleration capability to ease programming and keep costs low. So far, the industry has more widely adopted the maximal approach to overcome the pressing scaling bottlenecks of major workloads. The key enabler that bridges the programmability gap between the maximal approach favoring large performance & energy benefits and the minimal approach favoring ease of programming is compilation techniques. These techniques lower well-defined high-level constructs into lower-level primitives [12, 13]; our ISCA 2015 papers [1,10] and a follow-up work [14] explore them lightly. We believe that a good programming model that enables large benefits coupled with support for it across the entire system stack (including compilers & hardware) will continue to be important for effective near-memory system and accelerator designs [14]. We also believe that the maximal versus minimal approaches that are initially explored in our two ISCA 2015 papers) is a useful way of exploring emerging technologies (e.g., near-memory accelerators) to better understand the tradeoffs of system designs that exploit such technologies.

## III. INFLUENCE ON LATER WORKS

Our paper was at the beginning of a proliferation of scalable near-memory processing systems designed to accelerate key applications (see [15] for many works on the topic). Tesseract has inspired many near-memory system ideas (e.g., [16–28]) and served as the de facto comparison point for such systems, including near-memory graph processing accelerators that built on Tesseract and improved various aspects of Tesseract. Since machine learning accelerators that use high-bandwidth memory (e.g., [5,29]) and industrial PIM prototypes (e.g., [30–41]) are now in the market, near-memory processing is no longer an "eccentric" architecture it used to be when Tesseract was originally published.

Graph processing & analytics workloads remain as an important and growing class of applications in various forms, ranging from large-scale industrial graph analysis engines (e.g., [42]) to graph neural networks [43]. Our focus on large-scale graph processing in our ISCA 2015 paper increased attention to this domain in the computer architecture community, resulting in subsequent research on efficient hardware architectures for graph processing (e.g., [44–46]).

## IV. SUMMARY AND FUTURE OUTLOOK

We believe that our ISCA 2015 paper's principled re-thinking of system design to accelerate an important class of data-intensive workloads provided significant value and enabled/influenced a large body of follow-on works and ideas. We expect that such rethinking of system design for key workloads, especially with a focus on "maximal acceleration capability," will continue to be critical as pressing technology and application scaling challenges increasingly require us to think differently to substantially improve performance and energy (as well as other metrics). We believe the principles exploited in Tesseract are fundamental and they will remain useful and likely become even more important as systems become more constrained due to the continuously-increasing memory access and computation demands of future workloads. We also project that as hardware substrates for near-memory acceleration (e.g., 3D stacking, in-DRAM computation, NVM-based PIM, processing using memory [15]) evolve and mature, systems will take advantage of them even more, likely using principles similar to those used in the design of Tesseract.

## REFERENCES

[1] J. Ahn *et al.*, "A Scalable Processing-in-Memory Accelerator for Parallel Graph Processing," in *ISCA*, 2015.
[2] Hybrid Memory Cube Consortium, "HMC Specification 1.1," 2013.
[3] J. Jeddeloh and B. Keeth, "Hybrid Memory Cube: New DRAM Architecture Increases Density and Performance," in *VLSIT*, 2012.
[4] JEDEC, "High Bandwidth Memory (HBM) DRAM," Standard No. JESD235, 2013.
[5] N. Jouppi *et al.*, "TPU v4: An Optically Reconfigurable Supercomputer for Machine Learning with Hardware Support for Embedding," in *ISCA*, 2023.
[6] J. Fowers *et al.*, "A Configurable Cloud-Scale DNN Processor for Real-Time AI," in *ISCA*, 2018.
[7] S. Lie, "Cerebras Architecture Deep Dive: First Look Inside the Hardware/Software Co-Design for Deep Learning," in *IEEE Micro*, 2023.
[8] E. Talpes *et al.*, "The Microarchitecture of DOJO, Tesla's Exa-Scale Computer," in *IEEE Micro*, 2023.
[9] A. Ishii and R. Wells, "NVLink-Network Switch - NVIDIA's Switch Chip for High Communication-Bandwidth SuperPODs," in *Hot Chips*, 2022.
[10] J. Ahn *et al.*, "PIM-Enabled Instructions: A Low-Overhead, Locality-Aware Processing-in-Memory Architecture," in *ISCA*, 2015.
[11] R. Pope *et al.*, "Efficiently Scaling Transformer Inference," in *MLSys*, 2023.
[12] D. Lepikhin *et al.*, "GShard: Scaling Giant Models with Conditional Computation and Automatic Sharding," in *ICLR*, 2021.
[13] S. Wang *et al.*, "Overlap Communication with Dependent Computation via Decomposition in Large Deep Learning Models," in *ASPLOS*, 2023.
[14] J. Ahn *et al.*, "AIM: Energy-Efficient Aggregation Inside the Memory Hierarchy," in *ACM TACO*, vol. 13, no. 4, 2016.
[15] O. Mutlu *et al.*, "A Modern Primer on Processing in Memory," *Emerging Computing: From Devices to Systems*, 2021, https://arxiv.org/abs/2012.03112.
[16] M. Zhang *et al.*, "GraphP: Reducing Communication for PIM-Based Graph Processing with Efficient Data Partition," in *HPCA*, 2018.
[17] L. Song, "GraphR: Accelerating Graph Processing Using ReRAM," in *HPCA*, 2018.
[18] Y. Zhuo *et al.*, "GraphQ: Scalable PIM-Based Graph Processing," in *MICRO*, 2019.
[19] G. Dai *et al.*, "GraphH: A Processing-in-Memory Architecture for Large-Scale Graph Processing," *IEEE TCAD*, 2018.
[20] G. Li *et al.*, "GraphIA: An In-situ Accelerator for Large-scale Graph Processing," in *MEMSYS*, 2018.
[21] S. Rheindt *et al.*, "NEMESYS: Near-Memory Graph Copy Enhanced System-Software," in *MEMSYS*, 2019.
[22] L. Belayneh and V. Bertacco, "GraphVine: Exploiting Multicast for Scalable Graph Analytics," in *DATE*, 2020.
[23] N. Challapalle *et al.*, "GaaS-X: Graph Analytics Accelerator Supporting Sparse Data Representation using Crossbar Architectures," in *ISCA*, 2020.
[24] M. Zhou *et al.*, "Ultra Efficient Acceleration for De Novo Genome Assembly via Near-Memory Computing," in *PACT*, 2021.
[25] X. Xie *et al.*, "SpaceA: Sparse Matrix Vector Multiplication on Processing-in-Memory Accelerator," in *HPCA*, 2021.
[26] M. Zhou *et al.*, "HyGraph: Accelerating Graph Processing with Hybrid Memory-Centric Computing," in *DATE*, 2021.
[27] M. Lenjani *et al.*, "Gearbox: A Case for Supporting Accumulation Dispatching and Hybrid Partitioning in PIM-based Accelerators," in *ISCA*, 2022.
[28] M. Orenes-Vera *et al.*, "Dalorex: A Data-Local Program Execution and Architecture for Memory-Bound Applications," in *HPCA*, 2023.
[29] J. Choquette, "Nvidia Hopper GPU: Scaling Performance," in *Hot Chips*, 2022.
[30] F. Devaux, "The True Processing In Memory Accelerator," in *Hot Chips*, 2019.
[31] J. Gómez-Luna *et al.*, "Benchmarking a New Paradigm: Experimental Analysis and Characterization of a Real Processing-in-Memory System," *IEEE Access*, 2022.
[32] J. Gomez-Luna *et al.*, "Evaluating Machine Learning Workloads on Memory-Centric Computing Systems," in *ISPASS*, 2023.
[33] S. Lee *et al.*, "Hardware Architecture and Software Stack for PIM Based on Commercial DRAM Technology: Industrial Product," in *ISCA*, 2021.
[34] Y.-C. Kwon *et al.*, "25.4 A 20nm 6GB Function-In-Memory DRAM, Based on HBM2 with a 1.2 TFLOPS Programmable Computing Unit Using Bank-Level Parallelism, for Machine Learning Applications," in *ISSCC*, 2021.
[35] L. Ke *et al.*, "Near-Memory Processing in Action: Accelerating Personalized Recommendation with AxDIMM," *IEEE Micro*, 2021.
[36] D. Lee *et al.*, "Improving In-Memory Database Operations with Acceleration DIMM (AxDIMM)," in *DaMoN*, 2022.
[37] S. Lee *et al.*, "A 1ynm 1.25V 8Gb, 16Gb/s/pin GDDR6-based Accelerator-in-Memory supporting 1TFLOPS MAC Operation and Various Activation Functions for Deep-Learning Applications," in *ISSCC*, 2022.
[38] D. Niu *et al.*, "184QPS/W 64Mb/mm² 3D Logic-to-DRAM Hybrid Bonding with Process-Near-Memory Engine for Recommendation System," in *ISSCC*, 2022.
[39] Y. Kwon, "System Architecture and Software Stack for GDDR6-AiM," in *HCS*, 2022.
[40] G. Singh *et al.*, "FPGA-based Near-Memory Acceleration of Modern Data-Intensive Applications," *IEEE Micro*, 2021.
[41] G. Singh *et al.*, "Accelerating Weather Prediction using Near-Memory Reconfigurable Fabric," *ACM TRETS*, 2021.
[42] S. Hong *et al.*, "PGX.D: A Fast Distributed Graph Processing Engine," in *SC*, 2015.
[43] T. N. Kipf and M. Welling, "Semi-Supervised Classification with Graph Convolutional Networks," in *ICLR*, 2017.
[44] L. Nai *et al.*, "GraphPIM: Enabling Instruction-Level PIM Offloading in Graph Computing Frameworks," in *HPCA*, 2017.
[45] M. Besta *et al.*, "SISA: Set-Centric Instruction Set Architecture for Graph Mining on Processing-in-Memory Systems," in *MICRO*, 2021.
[46] T. J. Ham *et al.*, "Graphicionado: A High-Performance and Energy-Efficient Accelerator for Graph Analytics," in *MICRO*, 2016.

# Processing using DRAM

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

- 30X-257X performance and energy improvements

Seshadri+"RowClone: Fast and Efficient In-DRAM Copy and Initialization of Bulk Data," MICRO 2013.
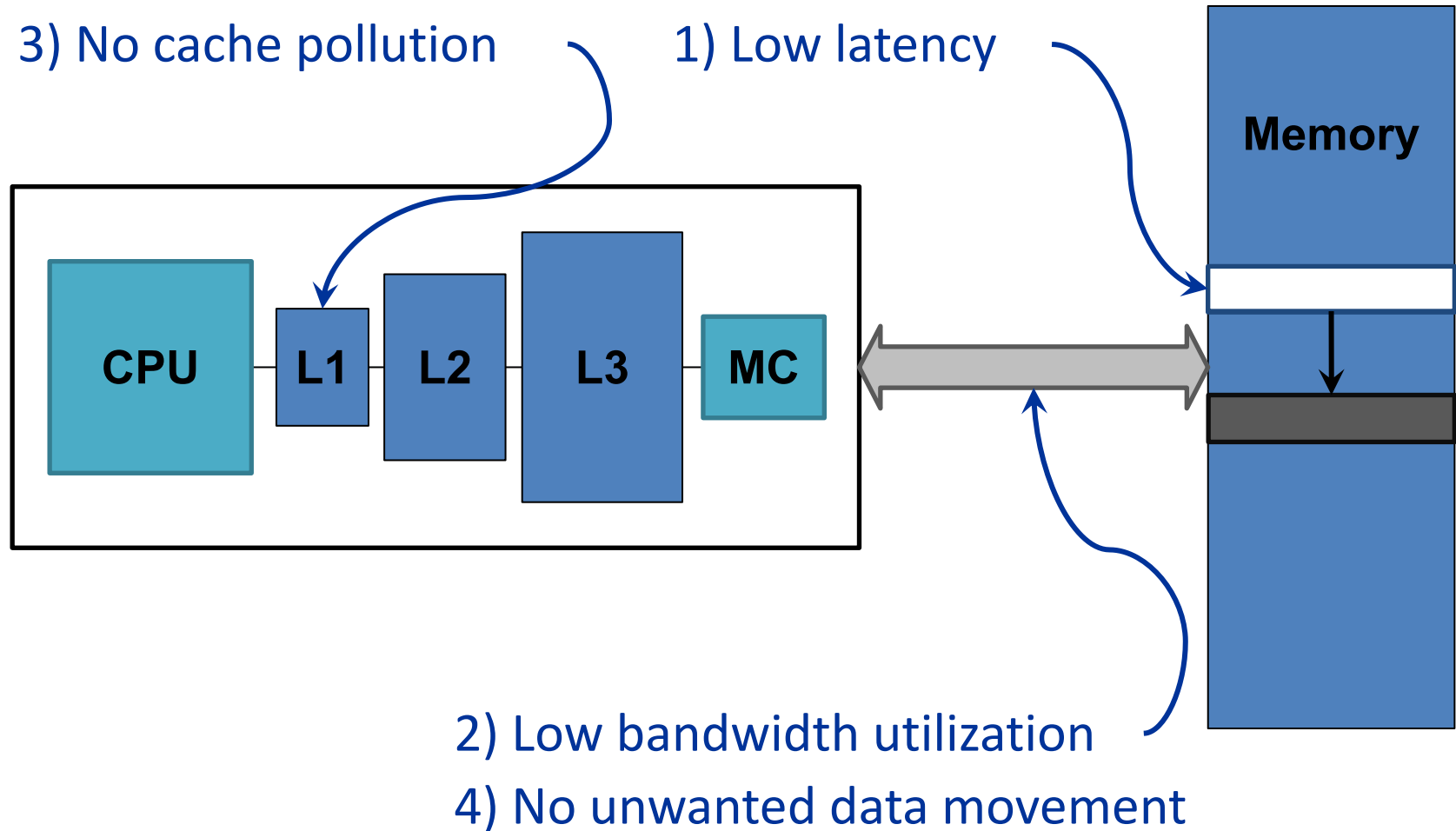
Seshadri+, "Fast Bulk Bitwise AND and OR in DRAM", IEEE CAL 2015.

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

Hajinazar+, "SIMDRAM: A Framework for Bit-Serial SIMD Processing using DRAM," ASPLOS 2021.

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

# Future Systems: In-Memory Copy

3) No cache pollution

1) Low latency

**Memory**

**CPU** **L1** **L2** **L3** **MC**

2) Low bandwidth utilization

4) No unwanted data movement

1046ns, 3.6uJ → 90ns, 0.04uJ

168

# 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 46th International Symposium on Microarchitecture (MICRO)*, Davis, CA, December 2013. [Slides (pptx) (pdf)] [Lightning Session Slides (pptx) (pdf)] [Poster (pptx) (pdf)]

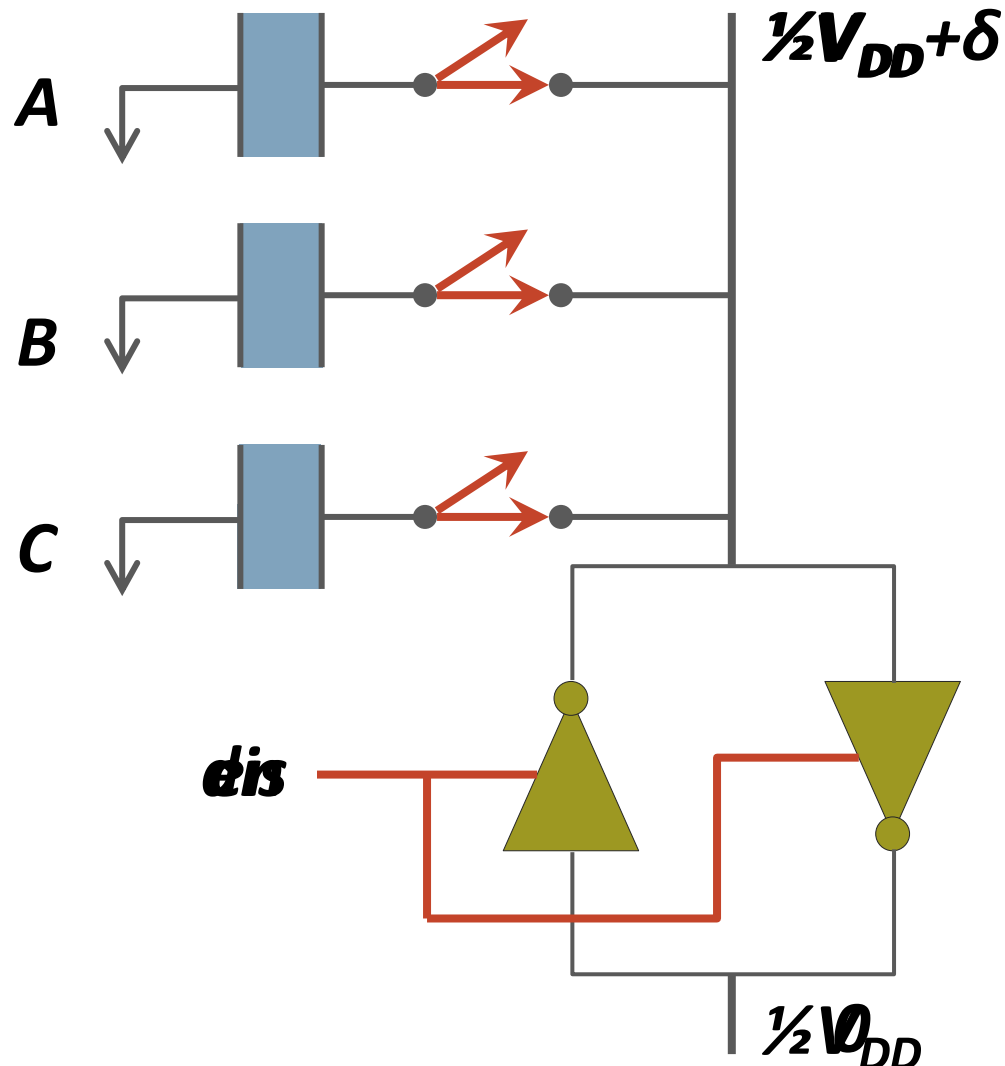## RowClone: Fast and Energy-Efficient In-DRAM Bulk Data Copy and Initialization

Vivek Seshadri
vseshadr@cs.cmu.edu

Yoongu Kim
yoongukim@cmu.edu

Chris Fallin*
cfallin@c1f.net

Donghyuk Lee
donghyuk1@cmu.edu

Rachata Ausavarungnirun
rachata@cmu.edu

Gennady Pekhimenko
gpekhime@cs.cmu.edu

Yixin Luo
yixinluo@andrew.cmu.edu

Onur Mutlu
onur@cmu.edu

Phillip B. Gibbons†
phillip.b.gibbons@intel.com

Michael A. Kozuch†
michael.a.kozuch@intel.com

Todd C. Mowry
tcm@cs.cmu.edu

Carnegie Mellon University    †Intel Pittsburgh

# In-DRAM AND/OR: Triple Row Activation



$\frac{1}{2}V_{DD}+\delta$

A

B

C

dis

$\frac{1}{2}V_{DD}$

**Final State**
*AB + BC + AC*

*C(A + B) +*
*~C(AB)*

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

# More on Ambit

- Vivek Seshadri, Donghyuk Lee, Thomas Mullins, Hasan Hassan, Amirali Boroumand, Jeremie Kim, Michael A. Kozuch, Onur Mutlu, Phillip B. Gibbons, and Todd C. Mowry,
**"Ambit: In-Memory Accelerator for Bulk Bitwise Operations Using Commodity DRAM Technology"**
*Proceedings of the 50th International Symposium on Microarchitecture* (**MICRO**), Boston, MA, USA, October 2017.
[Slides (pptx) (pdf)] [Lightning Session Slides (pptx) (pdf)] [Poster (pptx) (pdf)]

## 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**

# Capabilities of Off-The-Shelf Memory

# Existing DRAM Chips

# Are Already Quite Capable

# Real Processing Using Memory Prototype

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

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

Ataberk Olgun[§†]   Juan Gómez Luna[§]   Konstantinos Kanellopoulos[§]   Behzad Salami[§*]
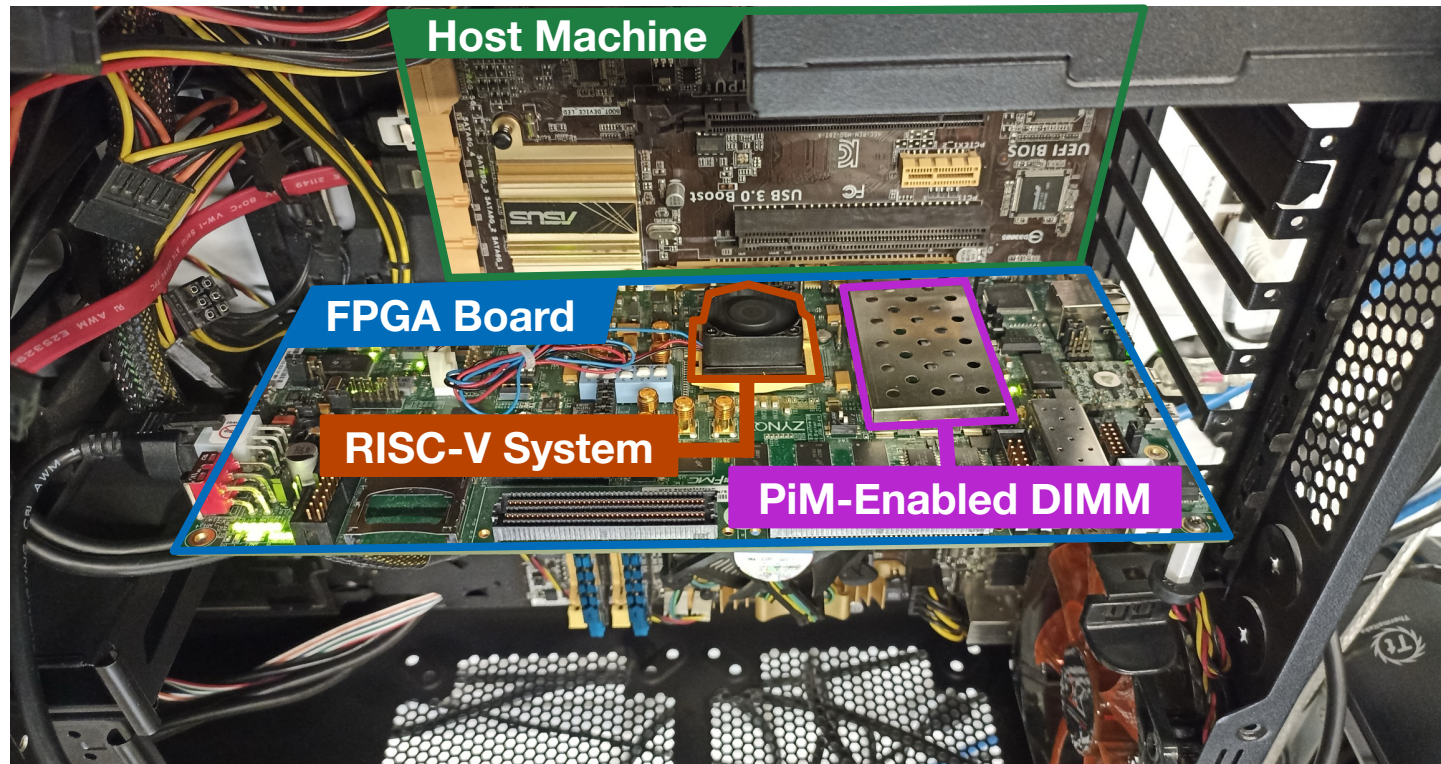Hasan Hassan[§]   Oğuz Ergin[†]   Onur Mutlu[§]

[§]ETH Zürich   [†]TOBB ETÜ   [*]BSC

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

# Real Processing-using-Memory Prototype



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

# Real Processing-using-Memory Prototype

**README.md**

## Building a PiDRAM Prototype

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

### Rebuilding Steps

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

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

### Generating DDR3 Controller IP sources

We cannot provide the sources for the Xilinx PHY IP we use in PiDRAM's memory controller due to licensing issues. We describe here how to regenerate them using Vivado 2016.2. First, you need to generate the IP RTL files:
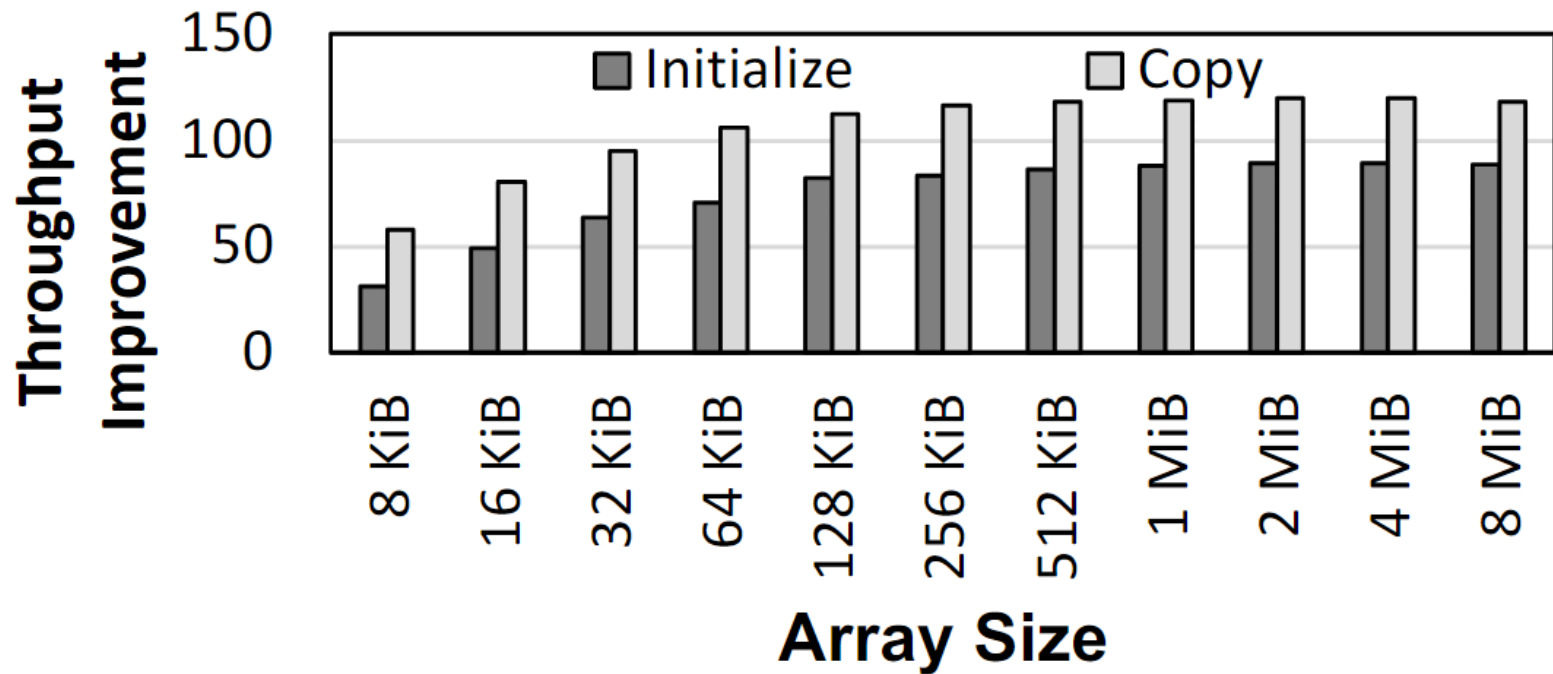
1- Open IP Catalog
2- Find "Memory Interface Generator (MIG 7 Series)" IP and double click

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

# Microbenchmark Copy/Initialization Throughput



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

# More on PiDRAM

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

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

Ataberk Olgun[§]    Juan Gómez Luna[§]    Konstantinos Kanellopoulos[§]    Behzad Salami[§]
Hasan Hassan[§]    Oğuz Ergin[†]    Onur Mutlu[§]

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

# DRAM Chips Are Already (Quite) Capable!

- **Appears at HPCA 2024**  **https://arxiv.org/pdf/2402.18736.pdf**

## Functionally-Complete Boolean Logic in Real DRAM Chips: Experimental Characterization and Analysis

İsmail Emir Yüksel    Yahya Can Tuğrul    Ataberk Olgun    F. Nisa Bostancı    A. Giray Yağlıkçı
Geraldo F. Oliveira    Haocong Luo    Juan Gómez-Luna    Mohammad Sadrosadati    Onur Mutlu

ETH Zürich

*We experimentally demonstrate that COTS DRAM chips are capable of performing 1) functionally-complete Boolean operations: NOT, NAND, and NOR and 2) many-input (i.e., more than two-input) AND and OR operations. We present an extensive characterization of new bulk bitwise operations in 256 off-the-shelf modern DDR4 DRAM chips. We evaluate the reliability of these operations using a metric called success rate: the fraction of correctly performed bitwise operations. Among our 19 new observations, we highlight four major results. First, we can perform the NOT operation on COTS DRAM chips with 98.37% success rate on average. Second, we can perform up to 16-input NAND, NOR, AND, and OR operations on COTS DRAM chips with high reliability (e.g., 16-input NAND, NOR, AND, and OR with average success rate of 94.94%, 95.87%, 94.94%, and 95.85%, respectively). Third, data pattern only slightly*

# The Capability of COTS DRAM Chips

We **demonstrate** that **COTS DRAM chips:**

**1** Can **copy one row into up to 31 other rows with >99.98% success rate**

**2** Can perform **NOT operation** with up to **32 output operands**

**3** Can perform up to **16-input AND, NAND, OR, and NOR** operations

# 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 24th International Symposium on High-Performance Computer Architecture* (**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 _25th International Symposium on High-Performance Computer Architecture_ (**HPCA**), Washington, DC, USA, February 2019.
  [Slides (pptx) (pdf)]
  [Full Talk Video (21 minutes)]
  [Full Talk Lecture Video (27 minutes)]
  **_Top Picks Honorable Mention by IEEE Micro._**

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

Jeremie S. Kim[‡§]    Minesh Patel[§]    Hasan Hassan[§]    Lois Orosa[§]    Onur Mutlu[§‡]

[‡]Carnegie Mellon University        [§]ETH Zürich

# In-DRAM True Random Number Generation

- Ataberk Olgun, Minesh Patel, A. Giray Yaglikci, Haocong Luo, Jeremie S. Kim, F. Nisa Bostanci, Nandita Vijaykumar, Oguz Ergin, and Onur Mutlu,
  **"QUAC-TRNG: High-Throughput True Random Number Generation Using Quadruple Row Activation in Commodity DRAM Chips"**
  *Proceedings of the 48th International Symposium on Computer Architecture* (**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*    [⊙]*University of Toronto*

# In-DRAM True Random Number Generation

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

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

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

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

# In-Flash Bulk Bitwise Execution

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

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

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

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

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

# 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.*

# How to Enable Adoption of Processing in Memory

# Potential Barriers to Adoption of PIM

1. **Applications** & **software** for PIM

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

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

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

5. **Infrastructures** to assess benefits and feasibility

**All can be solved with change of mindset**

# We Need to Revisit the Entire Stack

- With a **memory-centric mindset**

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

**We can get there step by step**

# Security Issues in Processing in Memory

- Does PIM make security better or easier?

- Does PIM make security worse?

- Many interesting & important questions here

- Some recent papers:
    - Evaluating Homomorphic Operations on a Real-World Processing-In-Memory System **[IISWC 2023]**
    - CIPHERMATCH: Accelerating Homomorphic Encryption based String Matching via Memory-Efficient Data Packing and In-Flash Processing **[ASPLOS 2025]**
    - Amplifying Main Memory-Based Timing Covert and Side Channels using Processing-in-Memory Operations **[arxiv 2024]**

# Potential Security Issues & Benefits (I)

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

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

*SAFARI*

# Potential Security Issues & Benefits (II)

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

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

**SAFARI**

# PIM Helps Security: Many Examples (I)

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

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

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

*ETH Zürich*

**https://arxiv.org/pdf/2309.06545**

# PIM Helps Security: Many Examples (II)

- Mayank Kabra, Rakesh Nadig, Harshita Gupta, Rahul Bera, Manos Frouzakis, Vamanan Arulchelvan, Yu Liang, Haiyu Mao, Mohammad Sadrosadati, and Onur Mutlu,
**"CIPHERMATCH: Accelerating Homomorphic Encryption based String Matching via Memory-Efficient Data Packing and In-Flash Processing,"**
*Proceedings of the [30th International Conference on Architectural Support for Programming Languages and Operating Systems](#)* (**ASPLOS**), Rotterdam, Netherlands, April 2025.

## CIPHERMATCH:
## Accelerating Homomorphic Encryption-Based String Matching via Memory-Efficient Data Packing and In-Flash Processing

Mayank Kabra
ETH Zürich
Zürich, Switzerland

Rakesh Nadig
ETH Zürich
Zürich, Switzerland

Harshita Gupta
ETH Zürich
Zürich, Switzerland

Rahul Bera
ETH Zürich
Zürich, Switzerland

Manos Frouzakis
ETH Zürich
Zürich, Switzerland

Vamanan Arulchelvan
ETH Zürich
Zürich, Switzerland

Yu Liang
ETH Zürich
Zürich, Switzerland

Haiyu Mao
King's College London
London, United Kingdom

Mohammad Sadrosadati
ETH Zürich
Zürich, Switzerland

Onur Mutlu
ETH Zürich
Zürich, Switzerland

**SAFARI**

# PIM Worsens Side Channels: An Example

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

Konstantinos Kanellopoulos[†*]    F. Nisa Bostancı[†*]    Ataberk Olgun[†]

A. Giray Yağlıkçı[†]    İsmail Emir Yüksel[†]    Nika Mansouri Ghiasi[†]

Zülal Bingöl[†‡]    Mohammad Sadrosadati[†]    Onur Mutlu[†]

[†]ETH Zürich    [‡]Bilkent University

# A Short Talk on Security of PIM Systems



Security of PIM Systems: Invited Talk at Dagstuhl MAD Seminar - 30.11.2023

# Concluding Remarks

# Summary: Three Major Limiters

- Technology scaling is not going well

- System complexity is increasing; old methods not keeping up

- Processor-centric designs are not keeping up

- These affect all metrics we care about

- These have fundamental impact on security and how we build secure systems

- **We need to revisit how we build architectures and how we secure them**

# Funding Acknowledgments

Thank you!

# Acknowledgments



SAFARI
SAFARI Research Group
safari.ethz.ch

Think BIG, Aim HIGH!

https://safari.ethz.ch

# SAFARI Newsletter June 2023 Edition

- https://safari.ethz.ch/safari-newsletter-june-2023/

# SAFARI Newsletter July 2024 Edition

- **https://safari.ethz.ch/safari-newsletter-july-2024/**

# 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.

👥 **440** followers   ⊙ ETH Zurich and Carnegie Mellon U...   🔗 https://safari.ethz.ch/   ✉ omutlu@gmail.com

🏠 **Overview**   📖 Repositories **80**   ⊞ Projects   📦 Packages   👤 People **13**

---

📖 **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++   ☆ 558   ⅄ 207

📖 **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   ☆ 131   ⅄ 48

📖 **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++   ☆ 271   ⅄ 149

📖 **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   ☆ 214   ⅄ 42

📖 **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   ☆ 122   ⅄ 28

📖 **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++   ☆ 112   ⅄ 34

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

# Future of
# Computer Architecture and Hardware Security

Onur Mutlu

omutlu@gmail.com

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

6 March 2025

University of Southern California

**SAFARI**          **ETH** zürich

# ISMM 2025

🏠 PLDI 2025 (series) / ISMM 2025 (series) /

# International Symposium on Memory Management

ISMM 2025

| About | Call for Papers |
|---|---|

**Update 2025-03-06:** One week submission extension, submissions now due 2025-03-18 (AoE UTC-12h)

Welcome to the home page of the 2025 ACM SIGPLAN International Symposium on Memory Management (ISMM 2025)! ISMM is the premier forum dedicated to research in memory management, covering the areas of memory performance, allocator design, garbage collection, architectural support for memory management, persistent memories, emerging memory technologies, and more.

ISMM'25 will held in-person as part of PLDI'25, sharing the venue and activities.

## Code of Conduct

ISMM follows the *ACM Policy Against Harassment at ACM Activities*. Please familiarize yourself with the policy and guide for reporting unacceptable behavior.

### Important Dates   🌐🕐 AoE (UTC-12h)

Tue 18 Mar 2025    **updated**
Paper submission

Wed 23 - Fri 25 Apr 2025
Author response period

Sat 3 May 2025
Author notification

Tue 13 May 2025
Camera-ready papers due

Tue 17 Jun 2025
Conference day

**SAFARI**          **https://conf.researchr.org/home/ismm-2025**          205

# DRAMSec 2025

## Fifth Workshop on DRAM Security (DRAMSec)
## June 21, 2025, co-located with ISCA 2025

DRAM is the most prevalent memory technology used in laptops, mobile phones, workstations and servers. As such, its security is paramount, yet DRAM attacks remain as viable as ever despite many attempts to resolve its security problems. The toolkit of DRAM disturb attacks has expanded with the introduction of new techniques such as Half-Double and RowPress, and it is likely that additional form of disturbance errors (and reliability and security issues) will emerge as we scale DRAM devices to smaller feature sizes. DRAM is also plagued by additional forms of attack, including side-channel, Denial-of-Service (DoS), and cold-boot attacks.

Against this backdrop, the industry is introducing new DRAM security solutions that require independent scrutiny from the academia. Academia continues to propose novel fixes for Rowhammer, often without the benefit of insight into constraints faced by the industry.

# PIM Tutorial November 2024 Edition



**MICRO 2024 - Tutorial on Memory-Centric Computing Systems**

Saturday, November 2nd, Austin, Texas, USA

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

**Program:** https://events.safari.ethz.ch/micro24-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
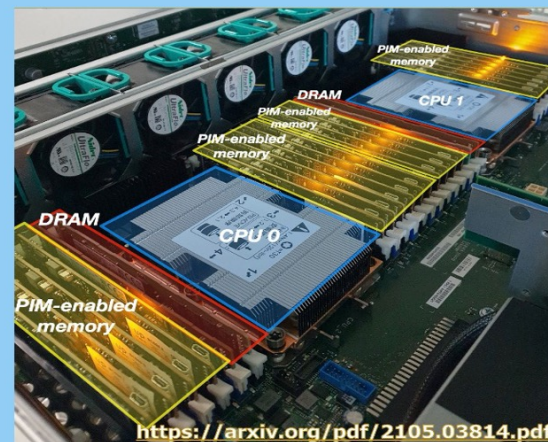Research challenges & opportunities

**https://www.youtube.com/watch?v=KV2MXvcBgb0**

**https://events.safari.ethz.ch/micro24-memorycentric-tutorial/**

# PIM Tutorial @ PPoPP/HPCA/CGO/CC



**PPoPP 2025 - Tutorial on Memory-Centric Computing Systems**

March 1st, Las Vegas, Nevada, USA

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

**Program:** https://events.safari.ethz.ch/ppopp25-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
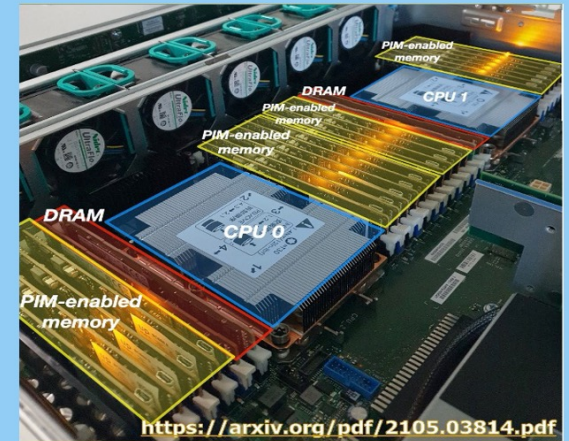Research challenges & opportunities

**https://www.youtube.com/live/NkDY6osus6g**

**https://events.safari.ethz.ch/ppopp25-memorycentric-tutorial/** 208

# Upcoming PIM Tutorials/Workshops (I)

## ASPLOS 2025 – 1ˢᵗ Workshop on Memory-Centric Computing Systems

Sunday, March 30ᵗʰ, Rotterdam, The Netherlands

**Organizers:** Geraldo F. Oliveira, Dr. Mohammad Sadrosadati,
Ataberk Olgun, Professor Onur Mutlu
**Program:** https://events.safari.ethz.ch/asplos25-MCCSys/doku.php

- 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



ASPLOS 2025

ACM International Conference on Architectural Support for Programming Languages and Operating Systems

Rotterdam, The Netherlands

March 30 - April 3, 2025



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

**https://events.safari.ethz.ch/asplos25-MCCSys/doku.php**

# Upcoming PIM Tutorials/Workshops (II)

## ICS 2025 – 2nd Workshop on Memory-Centric Computing Systems

Sunday, June 8th, Salt Lake City, USA

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

**Program:** https://events.safari.ethz.ch/ics25-MCCSys/doku.php

- 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



ACM International Conference on Supercomputing 2025

June 8-11, 2025
Salt Lake City, U.S.A.



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

**https://events.safari.ethz.ch/ics25-MCCSys/doku.php**

# Upcoming PIM Tutorials/Workshops (III)

**ISCA 2025 – 3rd Workshop on Memory-Centric Computing Systems**

Saturday, 21st June, 2025, Tokyo, Japan

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

**Program:** https://events.safari.ethz.ch/isca25-MCCSys/doku.php

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 2025

June 21–25, 2025

Tokyo, Japan



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

**https://events.safari.ethz.ch/isca25-MCCSys/doku.php**

# Backup Slides – Longer Version

# How Do We Make Sure Solution is Good?

- Many challenges (some below)
- Security by obscurity (as done in JEDEC DDR5 spec) unhelpful

- How do we guarantee we use **correct thresholds**?
  - Determining RH threshold is not easy
  - Many factors & conditions: RowPress, temperature, voltage, spatial variation, aging, voltage, and the unknowns

- How do we guarantee the **correct mitigating actions**?
  - In the presence of blast radius, address remapping, …

- How do we guarantee we perform **accurate bookkeeping**?
  - Updating counters properly? Row open time? When to reset? Worst-case access patterns?, …

# Good at What Cost?

- Even if we **magically assume** the solution prevents bitflips…

- Does the problem turn into a memory performance attack?

- How do we avoid large performance and energy losses?

- Is there a smarter way of handling things more holistically?
    - Better partitioning of responsibilities between CPU and memory
    - Our memory interface is terrible
        - DRAM has "no real freedom" to do things internally
        - DRAM should be more self-managing – we need more open minds

# What is RowPress?

Keeping a DRAM row **open for a long time**
causes bitflips in adjacent rows

These bitflips do **NOT** require many row activations

**Only one activation** is enough in some cases!

## DRAM chips tested

- 164 DDR4 chips from all 3 major DRAM manufacturers
- Covers different die densities and revisions

| Mfr. | #DIMMs | #Chips | Density | Die Rev. | Org. | Date |
|---|---|---|---|---|---|---|
| Mfr. S (Samsung) | 2 | 8 | 8Gb | B | x8 | 20-53 |
| | 1 | 8 | 8Gb | C | x8 | N/A |
| | 3 | 8 | 8Gb | D | x8 | 21-10 |
| | 2 | 8 | 4Gb | F | x8 | N/A |
| Mfr. H (SK Hynix) | 1 | 8 | 4Gb | A | x8 | 19-46 |
| | 1 | 8 | 4Gb | X | x8 | N/A |
| | 2 | 8 | 16Gb | A | x8 | 20-51 |
| | 2 | 8 | 16Gb | C | x8 | 21-36 |
| Mfr. M (Micron) | 1 | 16 | 8Gb | B | x4 | N/A |
| | 2 | 4 | 16Gb | B | x16 | 21-26 |
| | 1 | 16 | 16Gb | E | x4 | 20-14 |
| | 2 | 4 | 16Gb | E | x16 | 20-46 |
| | 1 | 4 | 16Gb | F | x16 | 21-50 |

RowPress significantly **amplifies** DRAM's vulnerability to read disturbance

RowPress has a **different** underlying error **mechanism** from RowHammer

# Reported HC$_{first}$ Values (2012 – Now)

**HC$_{first}$ : Number of hammers for the first bitflip**



HC$_{first}$ = ∞
(all good)

HC$_{first}$ = 1
(DRAM doomed?)

DDR3 @ **139K**
[Kim+, ISCA'14]

DDR3 @ **24K**
[Kim+, ISCA'20]

HBM2 @ **14K**
[Olgun+, DSN'24]

DDR4 @ **10K**
[Kim+, ISCA'20]

LPDDR4 @ **4.8K**
[Kim+, ISCA'20]

*Not shown: Significant variance in HC$_{first}$ across vendors and die variations

**SAFARI**

# RowPress at $t_{AggON}$ = Refresh Interval

**$HC_{first}$ : Number of hammers for the first bitflip**[*]

**$HC_{first} = \infty$**
**(all good)**

**$HC_{first} = 1$**
**(DRAM is doomed)**

*row on time = 7.2 us*
DDR4-new @ **380**
[Luo+, ISCA'23]

*row on time = 3.9 us*
HBM2-20nm @ **335**
[Olgun+, DSN'24]

*Not shown: Significant variance in $HC_{first}$ across vendors and die variations

# RowPress at $t_{AggON}$ = 9 * Refresh Interval

HC_{first} = ∞
(all good)

**HC$_{first}$ Scale*$^*$**

HC$_{first}$ = 1
(DRAM is doomed)

*row on time = 70.2 us*
DDR4-new @ **51**
[Luo+, ISCA'23]

*row on time = 35.1 us*
HBM2-20nm @ **123**
[Olgun+, DSN'24]

*Not shown: Significant variance in HC$_{first}$ across vendors and die variations

**SAFARI**

# Key Idea: NOT Operation

**Connect rows in neighboring subarrays** through **a NOT gate** by simultaneously activating rows

# Key Idea: NAND, NOR, AND, OR

**Manipulate the bitline voltage** to express
**a wide variety of functions** using
multiple-row activation in neighboring subarrays



$V_{REF}$

A

B

Multiple Row ACT

$V_{REF}$

X

Y

$V_{(A,B)}$

A

B

sense amp.
compares
$V_{(A,B)}$ and $V_{(X,Y)}$

X

Y

$V_{(X,Y)}$

SAFARI

222

# DRAM Chips Tested

- 256 DDR4 chips from two major DRAM manufacturers
- Covers different die revisions and chip densities

| Chip Mfr. | #Modules (#Chips) | Die Rev. | Mfr. Date[a] | Chip Density | Chip Org. | Speed Rate |
|---|---|---|---|---|---|---|
| SK Hynix | 9 (72) | M | N/A | 4Gb | x8 | 2666MT/s |
| | 5 (40) | A | N/A | 4Gb | x8 | 2133MT/s |
| | 1 (16) | A | N/A | 8Gb | x8 | 2666MT/s |
| | 1 (32) | A | 18-14 | 4Gb | x4 | 2400MT/s |
| | 1 (32) | A | 16-49 | 8Gb | x4 | 2400MT/s |
| | 1 (32) | M | 16-22 | 8Gb | x4 | 2666MT/s |
| Samsung | 1 (8) | F | 21-02 | 4Gb | x8 | 2666MT/s |
| | 2 (16) | D | 21-10 | 8Gb | x8 | 2133MT/s |
| | 1 (8) | A | 22-12 | 8Gb | x8 | 3200MT/s |

# Performing AND, NAND, OR, and NOR



**COTS DRAM chips can perform
{2, 4, 8, 16}-input AND, NAND, OR, and NOR operations**

# Performing AND, NAND, OR, and NOR



COTS DRAM chips can perform
16-input AND, NAND, OR, and NOR operations
with very high success rate (>94%)

# Other Backup Slides

# Onur Mutlu's SAFARI Research Group

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

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



40+ Researchers

**Think BIG, Aim HIGH!**

SAFARI

https://safari.ethz.ch

# SAFARI Newsletter December 2021 Edition

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

# SAFARI Newsletter June 2023 Edition

- https://safari.ethz.ch/safari-newsletter-june-2023/

# SAFARI Introduction & Research

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



# Think BIG, Aim HIGH!

**https://www.youtube.com/watch?v=mV2OuB2djEs**

# SAFARI PhD and Post-Doc Alumni

- https://safari.ethz.ch/safari-alumni/

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

- Samira Khan (Univ. of Virginia, Assistant Professor), HPCA 2014 Best Paper Session
- Saugata Ghose (Univ. of Illinois, Assistant Professor), DFRWS-EU 2017 Best Paper Award
- Jawad Haj-Yahya (Huawei Research Zurich, Principal Researcher)
- Lois Orosa (Galicia Supercomputing Center, Director)
- Jisung Park (POSTECH, Assistant Professor)
- Gagandeep Singh (AMD/Xilinx, Researcher)
- Juan Gomez-Luna (NVIDIA, Researcher), ISPASS 2023 Best Paper Session

# Processing in Memory: Evaluation Methods

# Simulators (Open Source)

- Ramulator 2.0 & Ramulator-PIM

- DAMOVSim

- UPMEMSim (UPMEM)

- AiMSim (SK Hynix)

- ...

*SAFARI*

# Ramulator + Gem5

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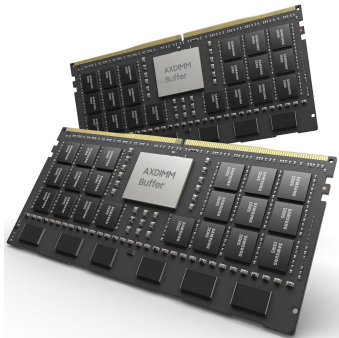  [arXiv version]
  [Ramulator 2.0 Source Code]

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

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

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

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

# Processing-in-Memory
# in the Real World

# Processing-in-Memory Landscape Today



[Samsung 2021]

[Alibaba 2022]

[SK Hynix 2022]

[Samsung 2021]

[UPMEM 2019]

**SAFARI** And, many other experimental chips and startups

# Processing-in-Memory Landscape Today

## Computational CXL-Memory Solution for Accelerating Memory-Intensive Applications

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

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

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



Fig. 6. FPGA prototype of proposed CMS card.

**SAFARI**

# Processing-in-Memory Landscape Today

## Samsung Processing in Memory Technology at Hot Chips 2023

By **Patrick Kennedy** - August 28, 2023

SAMSUNG

**CXL-PNM**

Industry's 1st CXL-PNM (Processing-near-Memory)

*Samsung PIM PNM For Transformer Based AI HC35_Page_24*

# Samsung AxDIMM (2021)

- ## DDRx-PIM

  - DLRM recommendation system



**Baseline System**

**AxDIMM System**

CH1: AxDIMM
CH0: RDIMM

Intel Broadwell Server

Ke et al. "Near-Memory Processing in Action: Accelerating Personalized Recommendation with AxDIMM", IEEE Micro (2021)

SAFARI

# 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

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

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

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

# Samsung Function-in-Memory DRAM (2021)

- **FIMDRAM based on HBM2**



**[3D Chip Structure of HBM with FIMDRAM]**

Labels on diagram:
- SID1 Core-die (HBM2)
- SID0 Core-die (FIMDRAM)
- Buffer-die

**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)**

Young-Cheon Kwon[1], Suk Han Lee[1], Jaehoon Lee[1], Sang-Hyuk Kwon[1],
Je Min Ryu[1], Jong-Pil Son[1], Seongil O[1], Hak-Soo Yu[1], Haesuk Lee[1],
Soo Young Kim[1], Youngmin Cho[1], Jin Guk Kim[1], Jongyoon Choi[1],
Hyun-Sung Shin[1], Jin Kim[1], BengSeng Phuah[1], HyoungMin Kim[1],
Myeong Jun Song[1], Ahn Choi[1], Daeho Kim[1], SooYoung Kim[1], Eun-Bong Kim[1],
David Wang[2], Shinhaeng Kang[1], Yuhwan Ro[3], Seungwoo Seo[3], JoonHo Song[3],
Jaeyoun Youn[1], Kyomin Sohn[1], Nam Sung Kim[1]

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

# Samsung Function-in-Memory DRAM (2021)

## 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]

# Samsung Function-in-Memory DRAM (2021)

**[Available instruction list for FIM operation]**

| Type | 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 |

# Chip Implementation

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



**[Digital RTL design for PCU block]**



Cell array for bank0 | Cell array for bank4 | Cell array for bank0 | Cell array for bank4 | **Pseudo channel-0** | **Pseudo channel-1**

PCU block for bank0 & 1 | PCU block for bank4 & 5 | PCU block for bank0 & 1 | PCU block for bank4 & 5

Cell array for bank1 | Cell array for bank5 | Cell array for bank1 | Cell array for bank5

Cell array for bank2 | Cell array for bank6 | Cell array for bank2 | 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 Control Block**

Cell array for bank11 | Cell array for bank15 | Cell array for bank11 | Cell array for bank15

PCU block for bank10 & 11 | PCU block for bank14 & 15 | PCU block for bank10 & 11 | PCU block for bank14 & 15

Cell array for bank10 | Cell array for bank14 | Cell array for bank10 | Cell array for bank14

Cell array for bank9 | Cell array for bank13 | Cell array for bank9 | Cell array for bank13

PCU block for bank8 & 9 | PCU block for bank12 & 13 | PCU block for bank8 & 9 | PCU block for bank12 & 13 | **Pseudo channel-0** | **Pseudo channel-1**

Cell array for bank8 | Cell array for bank12 | Cell array for bank8 | Cell array for bank12

Young-Cheon Kwon[1], Suk Han Lee[1], Jaehoon Lee[1], Sang-Hyuk Kwon[1], Je Min Ryu[1], Jong-Pil Son[1], Seongil O[1], Hak-Soo Yu[1], Haesuk Lee[1], Soo Young Kim[1], Youngmin Cho[1], Jin Guk Kim[1], Jongyoon Choi[1], Hyun-Sung Shin[1], Jin Kim[1], BengSeng Phuah[1], HyoungMin Kim[1], Myeong Jun Song[1], Ahn Choi[1], Daeho Kim[1], SooYoung Kim[1], Eun-Bong Kim[1], David Wang[2], Shinhaeng Kang[1], Yuhwan Ro[3], Seungwoo Seo[3], JoonHo Song[3], Jaeyoun Youn[1], Kyomin Sohn[1], Nam Sung Kim[1]
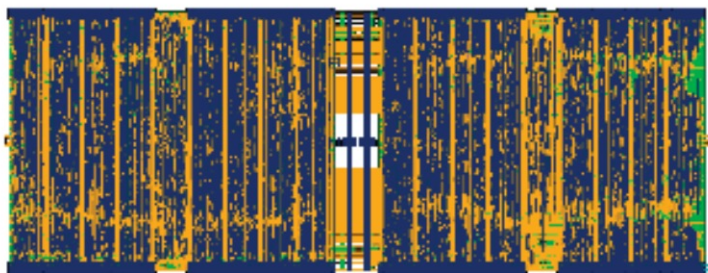
[1]Samsung Electronics, Hwaseong, Korea
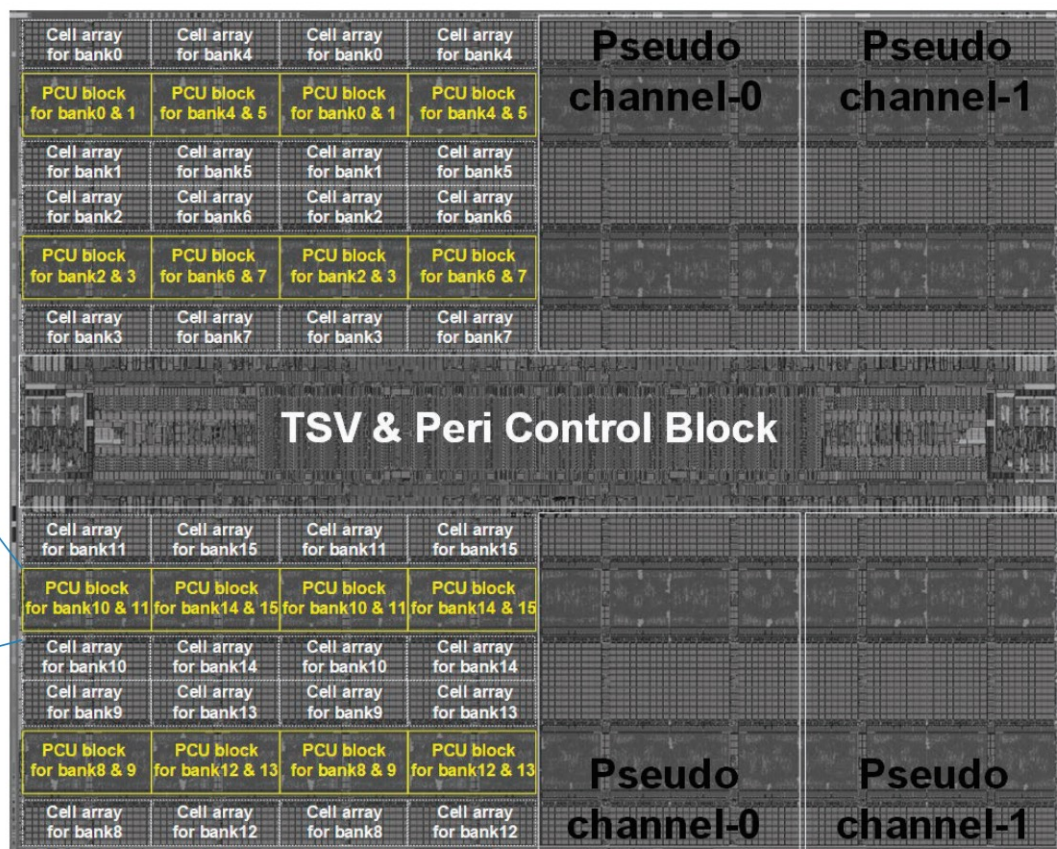[2]Samsung Electronics, San Jose, CA
[3]Samsung Electronics, Suwon, Korea

# SK Hynix Accelerator-in-Memory (2022)

## SK hynix Develops PIM, Next-Generation AI Accelerator

February 16, 2022                                    (in) (f) (y) (🔗)

**Seoul, February 16, 2022**

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

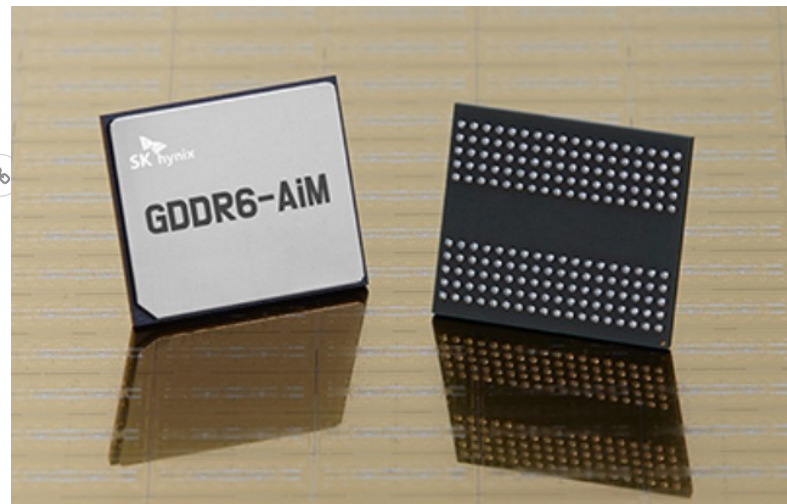*PIM(Processing In Memory): A next-generation technology that provides a solution for data congestion issues for AI and big data by adding computational functions to semiconductor memory

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

SK hynix plans to showcase its PIM development at the world's most prestigious semiconductor conference, 2022 ISSCC*, in San Francisco at the end of this month. The company expects continued efforts for innovation of this technology to bring the memory-centric computing, in which semiconductor memory plays a central role, a step closer to the reality in devices such as smartphones.

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

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

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

Seongju Lee, SK hynix, Icheon, Korea

In Paper 11.1, SK Hynix describes an 1ynm, GDDR6-based accelerator-in-memory with a command set for deep-learning operation. The 8Gb design achieves a peak throughput of 1TFLOPS with 1GHz MAC operations and supports major activation functions to improve accuracy.

# SK Hynix Accelerator-in-Memory (2022)

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

# AliBaba PIM Recommendation System (2022)



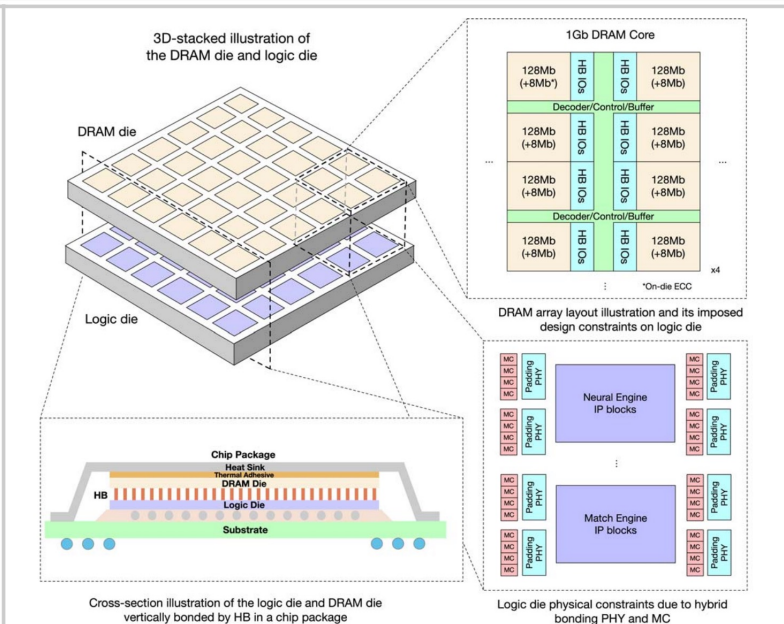ISSCC 2022 / February 24, 2022 / 8:30 AM

Figure 29.1.2: Illustration of 3D-stacked chip, cross-illustration of package, DRAM array layout and design blocks on logic die.
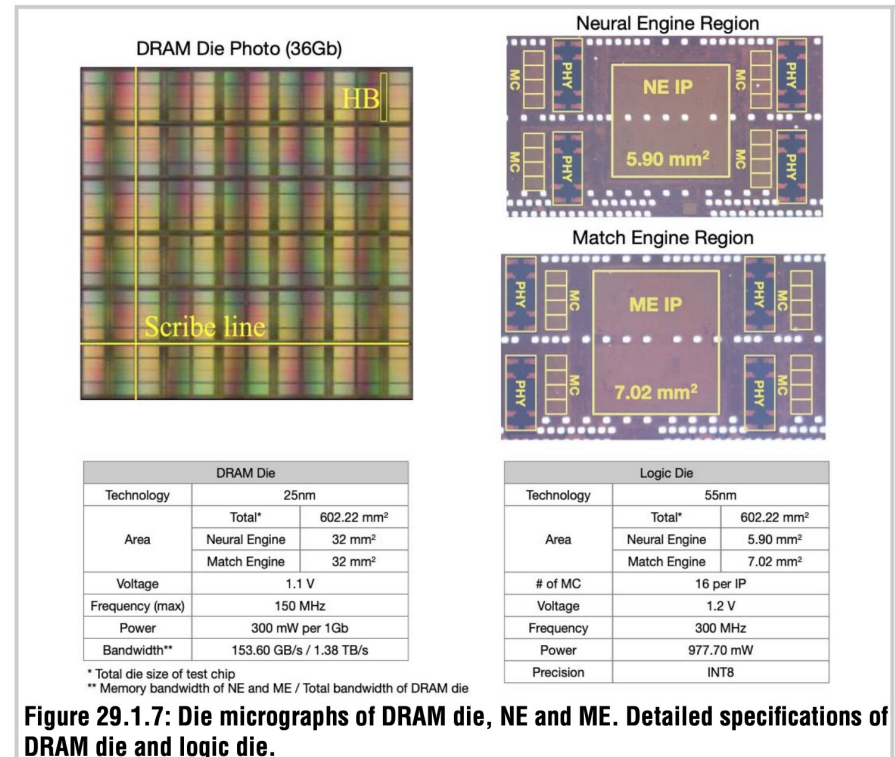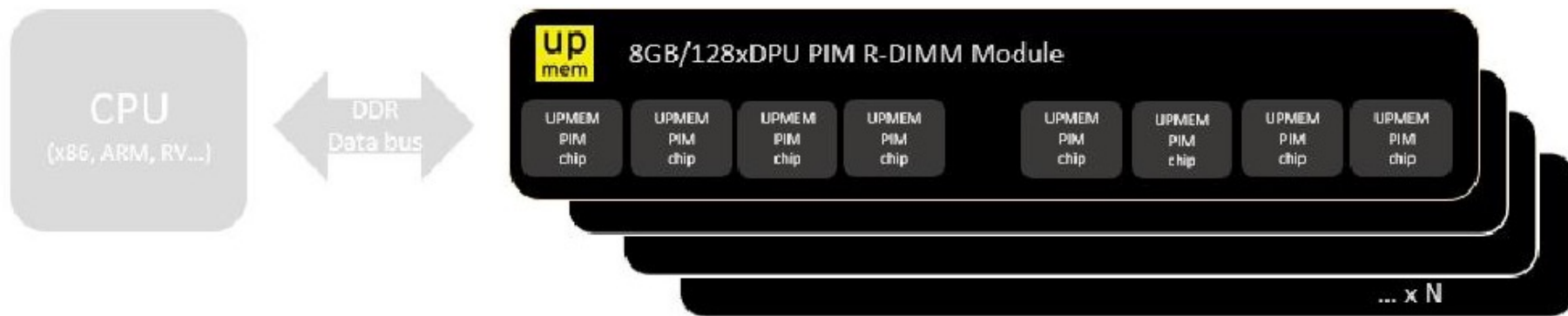
Figure 29.1.7: Die micrographs of DRAM die, NE and ME. Detailed specifications of DRAM die and logic die.

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

Dimin Niu[1], Shuangchen Li[1], Yuhao Wang[1], Wei Han[1], Zhe Zhang[2], Yijin Guan[2], Tianchan Guan[3], Fei Sun[1], Fei Xue[1], Lide Duan[1], Yuanwei Fang[1], Hongzhong Zheng[1], Xiping Jiang[4], Song Wang[4], Fengguo Zuo[4], Yubing Wang[4], Bing Yu[4], Qiwei Ren[4], Yuan Xie[1]
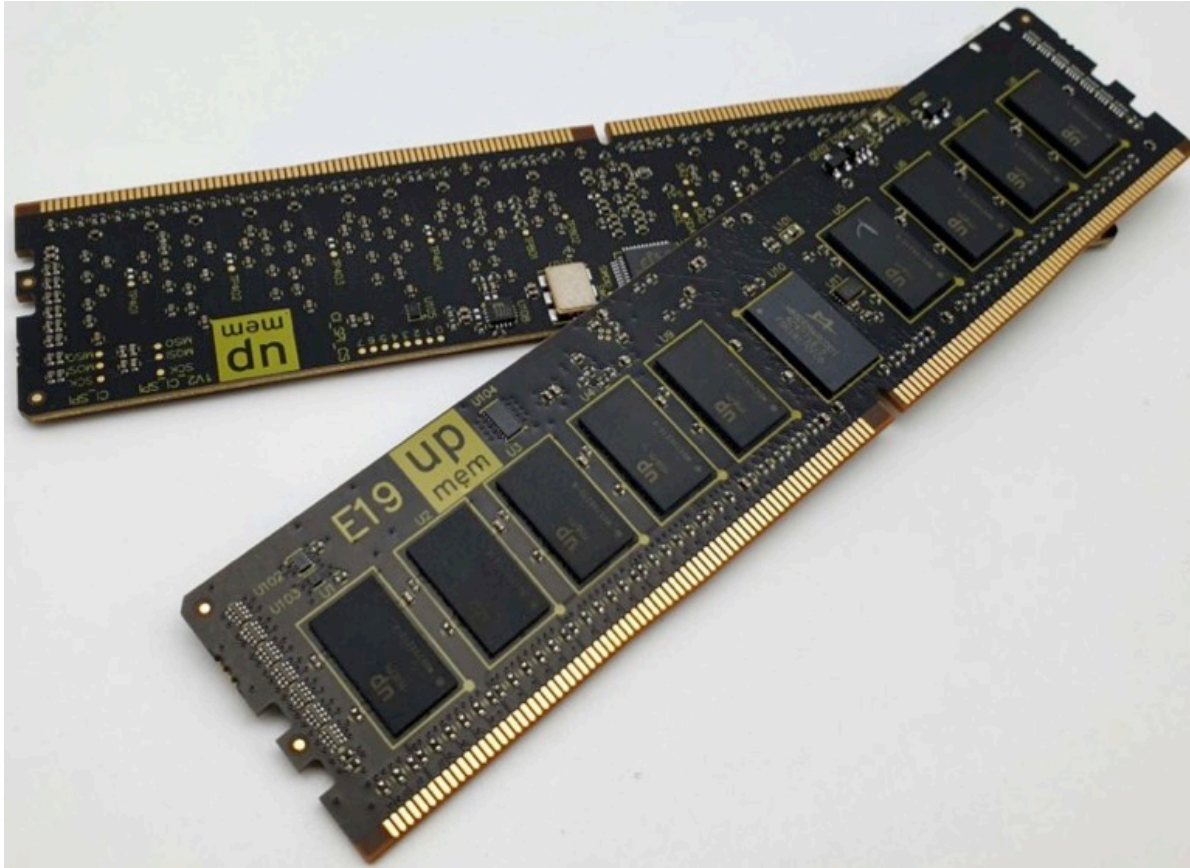
# UPMEM Processing-in-DRAM Engine (2019)

- <span style="color:red">Processing in DRAM Engine</span>

- 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
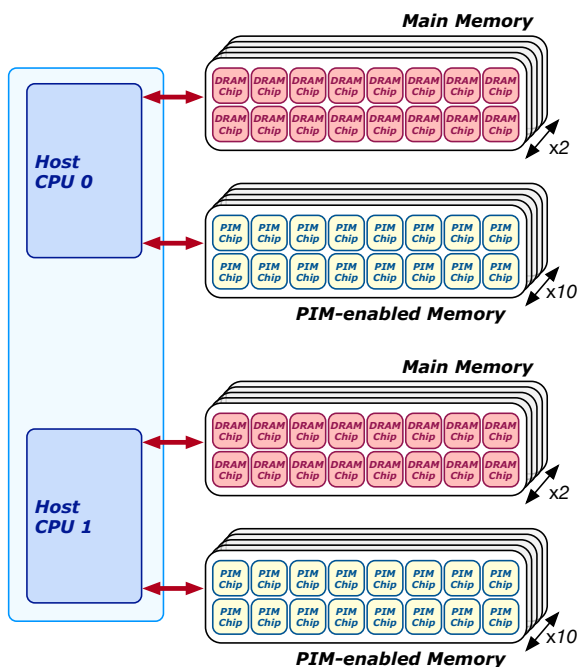  - **Large amounts of** compute & memory bandwidth

# UPMEM Memory Modules

- E19: 8 chips DIMM (1 rank). DPUs @ 267 MHz
- P21: 16 chips DIMM (2 ranks). DPUs @ 350 MHz
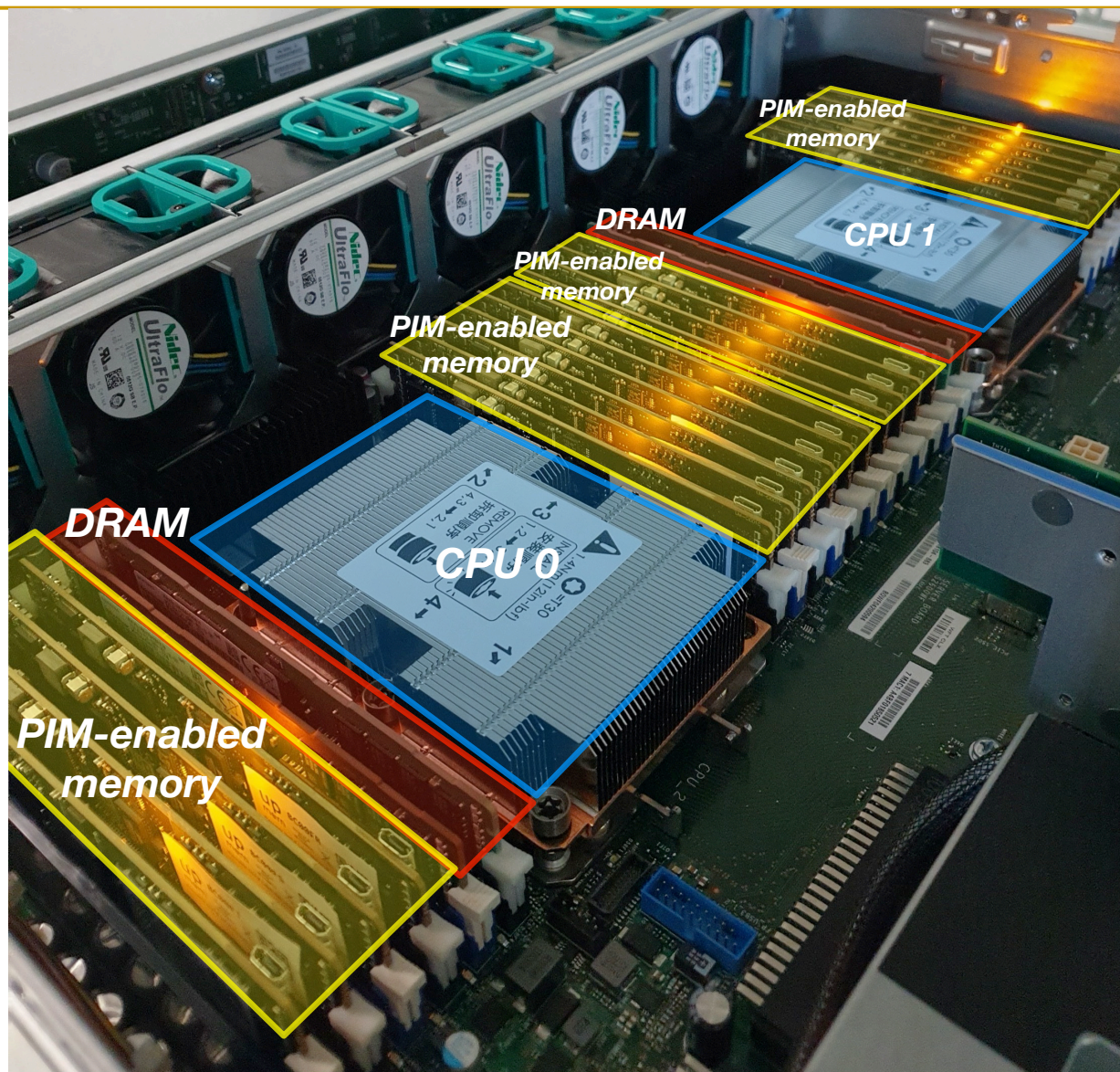
# 2,560-DPU Processing-in-Memory System



**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 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*).
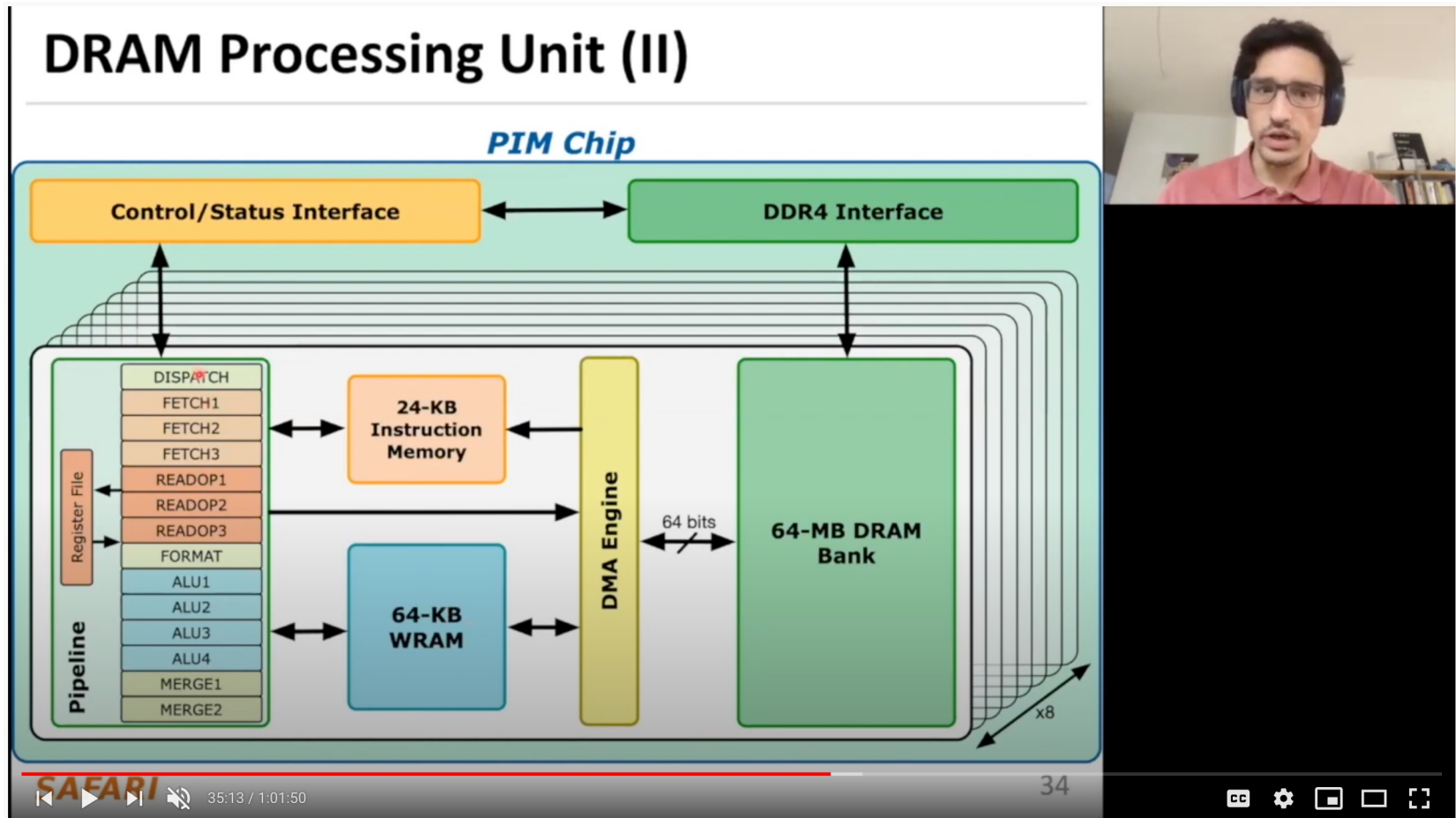
Recent research explores different forms of PIM architectures, motivated by the emergence of new 3D-stacked memory technologies that integrate memory with a logic layer where processing elements can be easily placed. Past works evaluate these architectures in simulation or, at best, with simplified hardware 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**

# More on the UPMEM PIM System



Computer Architecture - Lecture 12d: Real Processing-in-DRAM with UPMEM (ETH Zürich, Fall 2020)

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

# Experimental Analysis of the UPMEM PIM Engine

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

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

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

Recent research explores different forms of PIM architectures, motivated by the emergence of new 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 PIM System Summary & Analysis

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

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

Juan Gómez-Luna
*ETH Zürich*

Izzat El Hajj
*American University of Beirut*

Ivan Fernandez
*University of Malaga*

Christina Giannoula
*National Technical University of Athens*
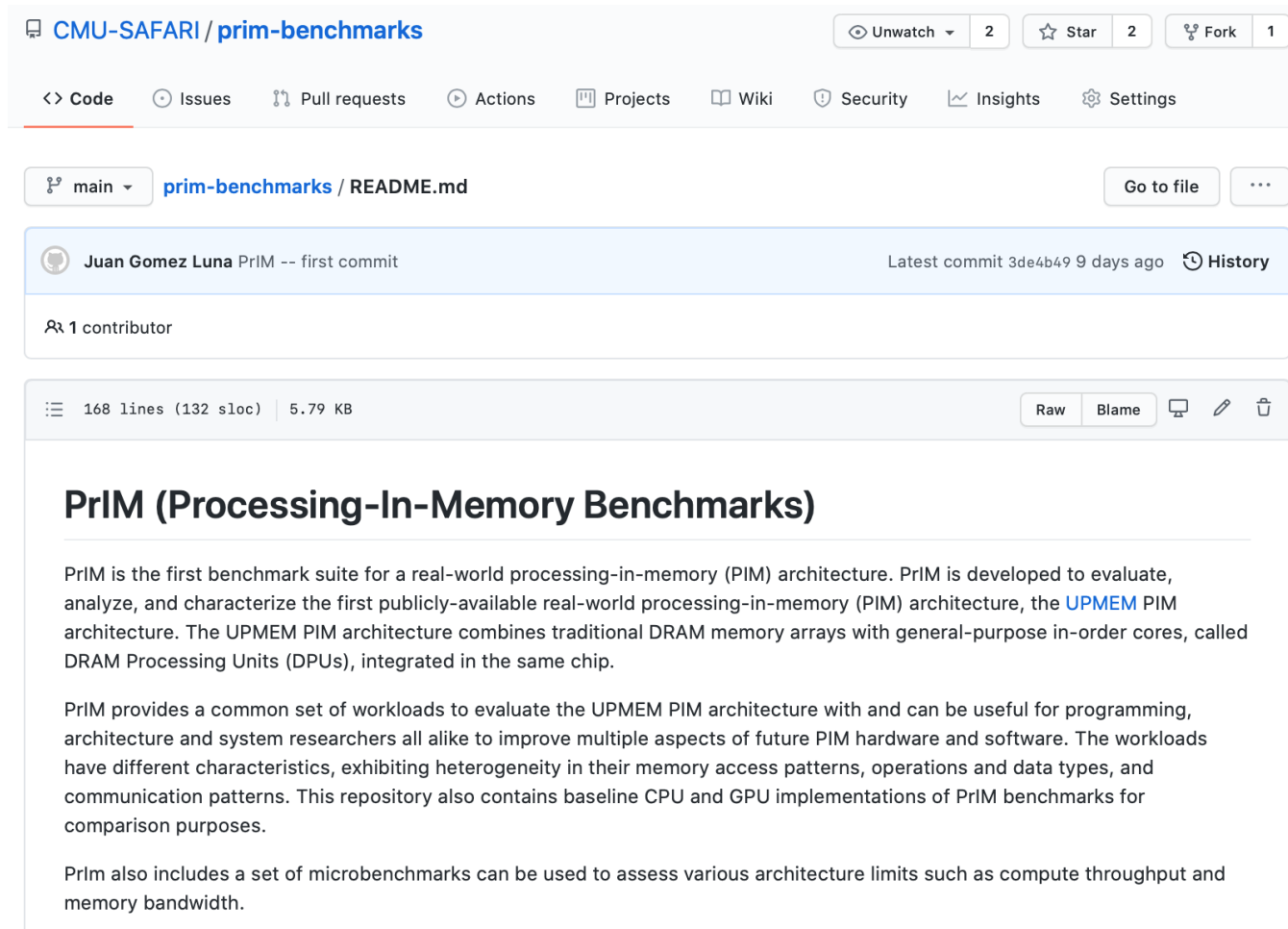
Geraldo F. Oliveira
*ETH Zürich*

Onur Mutlu
*ETH Zürich*

# PrIM Benchmarks: Application Domains

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



CMU-SAFARI / prim-benchmarks

Unwatch ▾  2    ☆ Star  2    ⑂ Fork  1

<> Code  ⊙ Issues  ⇅ Pull requests  ▷ Actions  ⊞ Projects  ▭ Wiki  ⊘ Security  ⮑ Insights  ⚙ Settings

⑂ main ▾    prim-benchmarks / README.md    Go to file    ···

Juan Gomez Luna PrIM -- first commit    Latest commit 3de4b49 9 days ago   ⟲ History

⋀ 1 contributor

≡   168 lines (132 sloc)   5.79 KB    Raw   Blame   🖥 ✏ 🗑

## PrIM (Processing-In-Memory Benchmarks)

PrIM is the first benchmark suite for a real-world processing-in-memory (PIM) architecture. PrIM is developed to evaluate, analyze, and characterize the first publicly-available real-world processing-in-memory (PIM) architecture, the UPMEM PIM architecture. The UPMEM PIM architecture combines traditional DRAM memory arrays with general-purpose in-order cores, called DRAM Processing Units (DPUs), integrated in the same chip.

PrIM provides a common set of workloads to evaluate the UPMEM PIM architecture with and can be useful for programming, architecture and system researchers all alike to improve multiple aspects of future PIM hardware and software. The workloads have different characteristics, exhibiting heterogeneity in their memory access patterns, operations and data types, and communication patterns. This repository also contains baseline CPU and GPU implementations of PrIM benchmarks for comparison purposes.

PrIm also includes a set of microbenchmarks can be used to assess various architecture limits such as compute throughput and memory bandwidth.

# Understanding a Modern PIM Architecture

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

JUAN GÓMEZ-LUNA[1], IZZAT EL HAJJ[2], IVAN FERNANDEZ[1,3], CHRISTINA GIANNOULA[1,4],
GERALDO F. OLIVEIRA[1], AND ONUR MUTLU[1]

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

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

https://arxiv.org/pdf/2105.03814.pdf
https://github.com/CMU-SAFARI/prim-benchmarks

# More Security Implications (I)

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



Not there yet, but ...

ROWHAMMERJS

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"**

Drammer: Deterministic Rowhammer Attacks on Mobile Platforms, CCS'16

258

# More Security Implications (III)

- Using an integrated GPU in a mobile system to remotely escalate privilege via the WebGL interface. IEEE S&P 2018

# Grand Pwning Unit: Accelerating Microarchitectural Attacks with the GPU

Pietro Frigo
Vrije Universiteit
Amsterdam
p.frigo@vu.nl

Cristiano Giuffrida
Vrije Universiteit
Amsterdam
giuffrida@cs.vu.nl

Herbert Bos
Vrije Universiteit
Amsterdam
herbertb@cs.vu.nl

Kaveh Razavi
Vrije Universiteit
Amsterdam
kaveh@cs.vu.nl

# More Security Implications (IV)

- Rowhammer over RDMA (I) USENIX ATC 2018

*THROWHAMMER* —

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

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

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

## Throwhammer: Rowhammer Attacks over the Network and Defenses

Andrei Tatar
*VU Amsterdam*

Radhesh Krishnan
*VU Amsterdam*

Elias Athanasopoulos
*University of Cyprus*

Cristiano Giuffrida
*VU Amsterdam*

Herbert Bos
*VU Amsterdam*

Kaveh Razavi
*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

Misiker Tadesse Aga
University of Michigan

Michael Schwarz
Graz University of Technology

Daniel Gruss
Graz University of Technology

Clémentine Maurice
Univ Rennes, CNRS, IRISA

Lukas Raab
Graz University of Technology

Lukas Lamster
Graz University of Technology

# More Security Implications (VI)

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)

- <span style="color:blue">USENIX Security 2019</span>

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

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

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

**A Single Bit-flip Can Cause Terminal Brain Damage to DNNs**

*One specific bit-flip in a DNN's representation leads to accuracy drop over 90%*

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

**Read More**

# More Security Implications (VIII)

- **USENIX Security 2020**

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

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

Adnan Siraj Rakin
*Arizona State University*
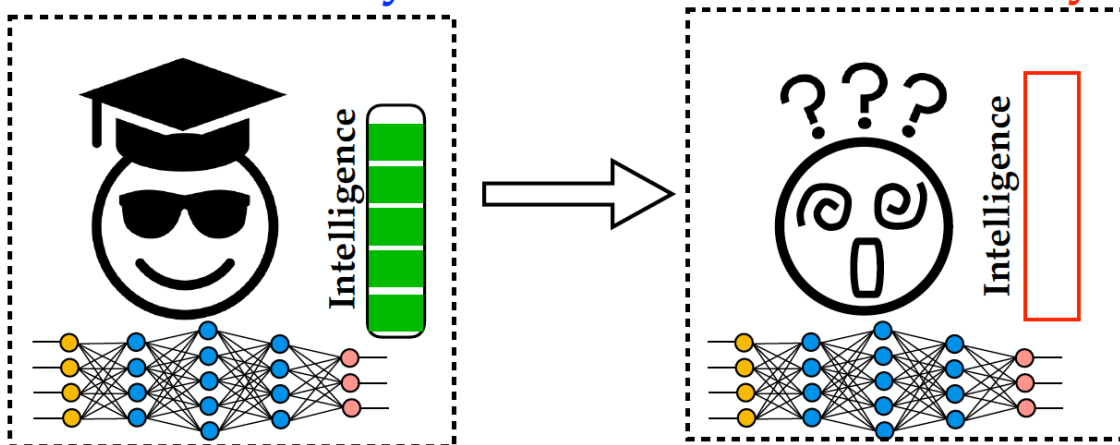*asrakin@asu.edu*

Deliang Fan
*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)**

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

**Google** Security Blog

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

---

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

May 25, 2021

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

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

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

# More Security Implications (VIII)

- **USENIX Security 2022**

- Google's Half-Double RowHammer Attack

## Half-Double: Hammering From the Next Row Over

Andreas Kogler[1]    Jonas Juffinger[1,2]    Salman Qazi[3]    Yoongu Kim[3]    Moritz Lipp[4]*
Nicolas Boichat[3]    Eric Shiu[5]    Mattias Nissler[3]    Daniel Gruss[1]
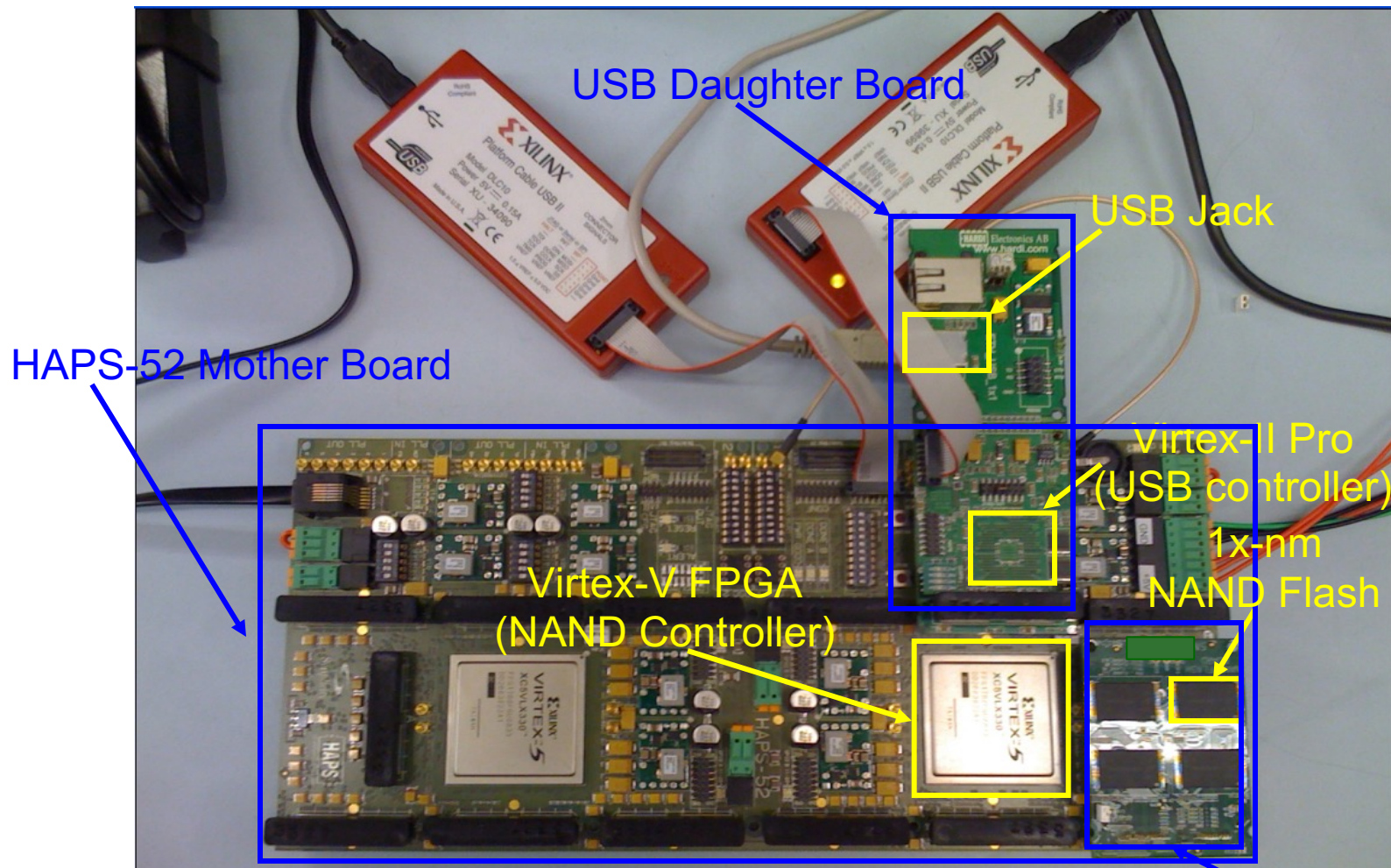
[1]*Graz University of Technology*    [2]*Lamarr Security Research*    [3]*Google*
[4]*Amazon Web Services*    [5]*Rivos*

SAFARI

# More Security Implications?

# Aside: Intelligent Controller for NAND Flash



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

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

# Intelligent Flash Controllers [PIEEE'17]

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

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

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

https://arxiv.org/pdf/1706.08642

269