# Tutorial on
# Memory-Centric Computing:
# Processing-Using-Memory

Geraldo F. Oliveira

Prof. Onur Mutlu

ISCA 2024

29 June 2024

**SAFARI**

**ETH** *zürich*

# Agenda

- Introduction to Memory-Centric Computing Systems

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

- Coffee Break

- Real-World Processing-Near-Memory Systems

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

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

- PIM Programming & Infrastructure for PIM Research

- Closing Remarks

*SAFARI*

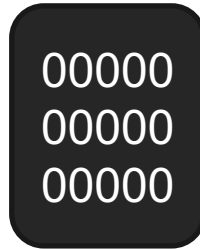# Processing in Memory: Two Approaches

1. Processing near Memory
2. Processing using Memory

# Starting Simple: Data Copy and Initialization

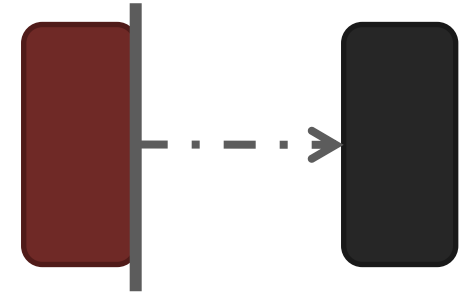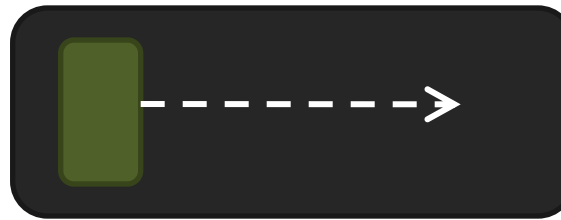*memmove & memcpy:* 5% cycles in Google's datacenter [Kanev+ ISCA'15]

**Forking**

00000
00000
00000

**Zero initialization
(e.g., security)**

**Checkpointing**

**VM Cloning
Deduplication**

**Page Migration**

• • •
Many more

**SAFARI**

# Today's Systems: Bulk Data Copy

3) Cache pollution

1) High latency

**Memory**

CPU   L1   L2   L3   MC
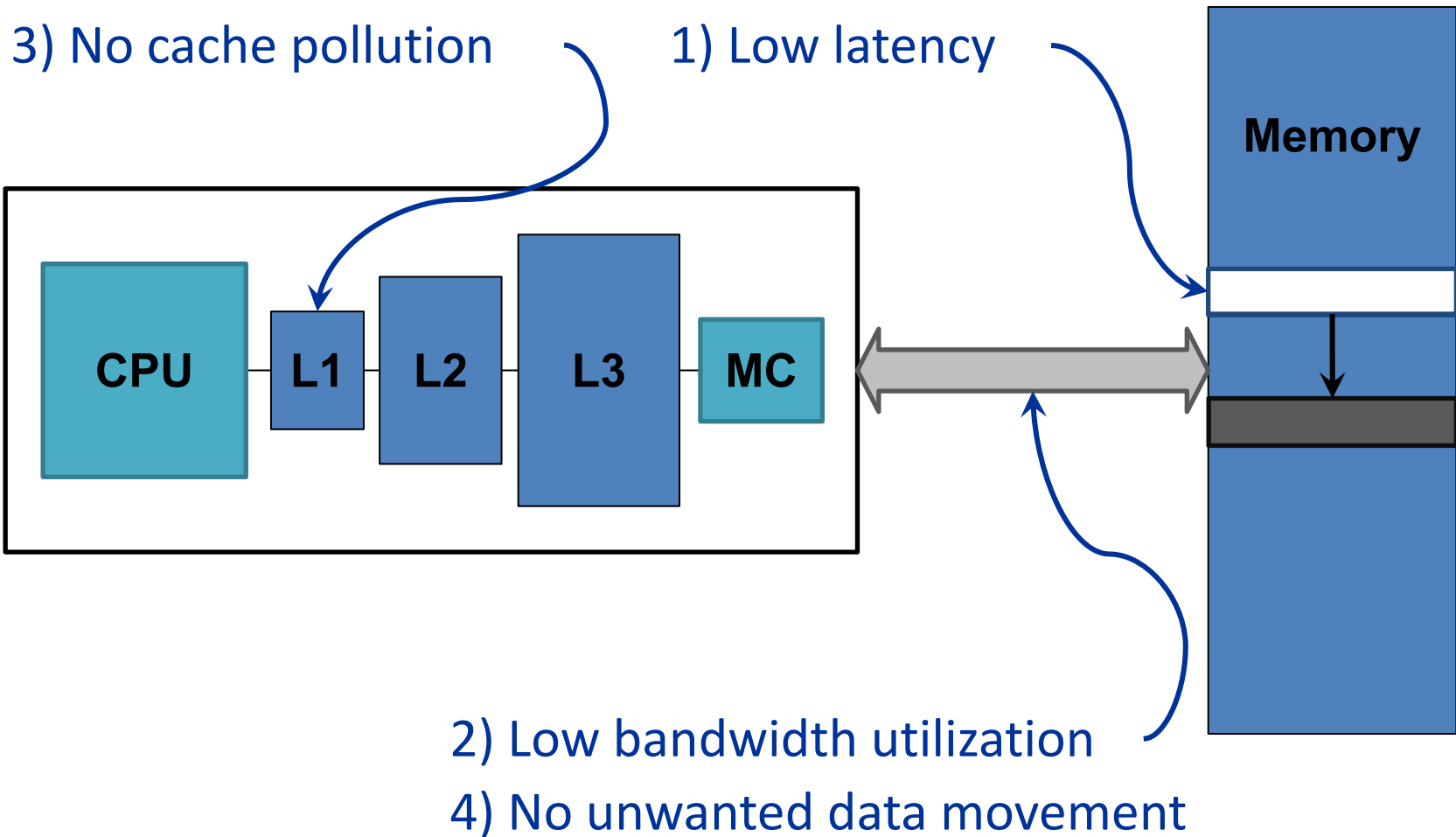
2) High bandwidth utilization

4) Unwanted data movement

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

# Future Systems: In-Memory Copy

3) No cache pollution

1) Low latency

**Memory**

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

2) Low bandwidth utilization
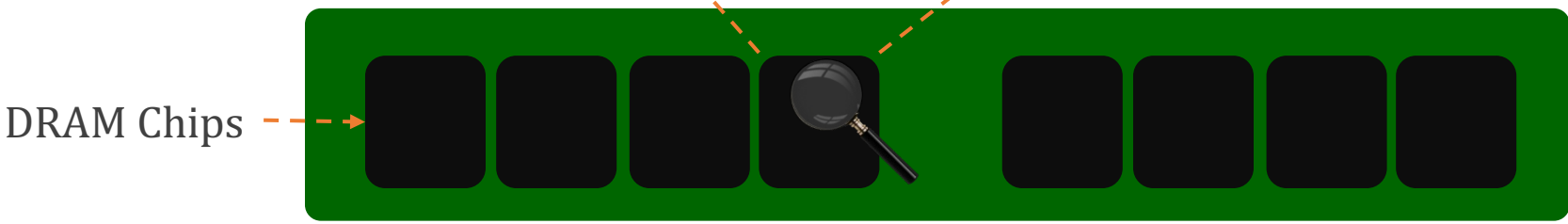
4) No unwanted data movement

1046ns, 3.6uJ  →  90ns, 0.04uJ

# Brief Review:

## Inside A DRAM Chip

# Inside a DRAM Chip



Bitline

Wordline

Subarray
(2D Array
of DRAM Cells)

DRAM Cells

Wordline

Access
Transistor

Bitline

Sense Amplifiers

Row Buffer

DRAM Bank

Storage
Capacitor

DRAM Chips

DRAM Module

SAFARI

8

# Inside a DRAM Chip: Another View



Memory Channel

Chip I/O

Bank

Subarray

Bank I/O

Row of DRAM Cells

Row Buffer

**SAFARI**

# DRAM Cell Operation



wordline

½ V$_{DD}$

access
transistor

storage
capacitor

bitline

enable

sense
amplifier

1. ACTIVATE (ACT)

2. READ/WRITE

3. PRECHARGE (PRE)

SAFARI

# DRAM Cell Operation (1/3)



wordline

1. raise wordline

$\frac{1}{2}V_{DDD}$ + δ

storage
capacitor

access
transistor

bitline

5. capacitor loses charge to bitline
2. capacitor shares charge to bitline

4. amplify deviation
in the bitline

3. enable
sense amplifier

enable

sense
amplifier

1. ACTIVATE (ACT)

2. READ/WRITE

3. PRECHARGE (PRE)

# DRAM Cell Operation (2/3)

wordline

$V_{DD}$

storage capacitor

access transistor

bitline

**1. ACTIVATE (ACT)**

**2. READ/WRITE**

**3. PRECHARGE (PRE)**

enable

sense amplifier

read/write charge
latched in sense amplifier

# DRAM Cell Operation (3/3)



1. lower wordline

wordline

2. precharge bitline for next access

½V$_{DDD}$

bitline

storage capacitor

access transistor

3. disable sense amplifier

enable

sense amplifier

1. ACTIVATE (ACT)

2. READ/WRITE

3. PRECHARGE (PRE)

**SAFARI**

# Future Systems: In-Memory Copy

3) No cache pollution

1) Low latency

**Memory**

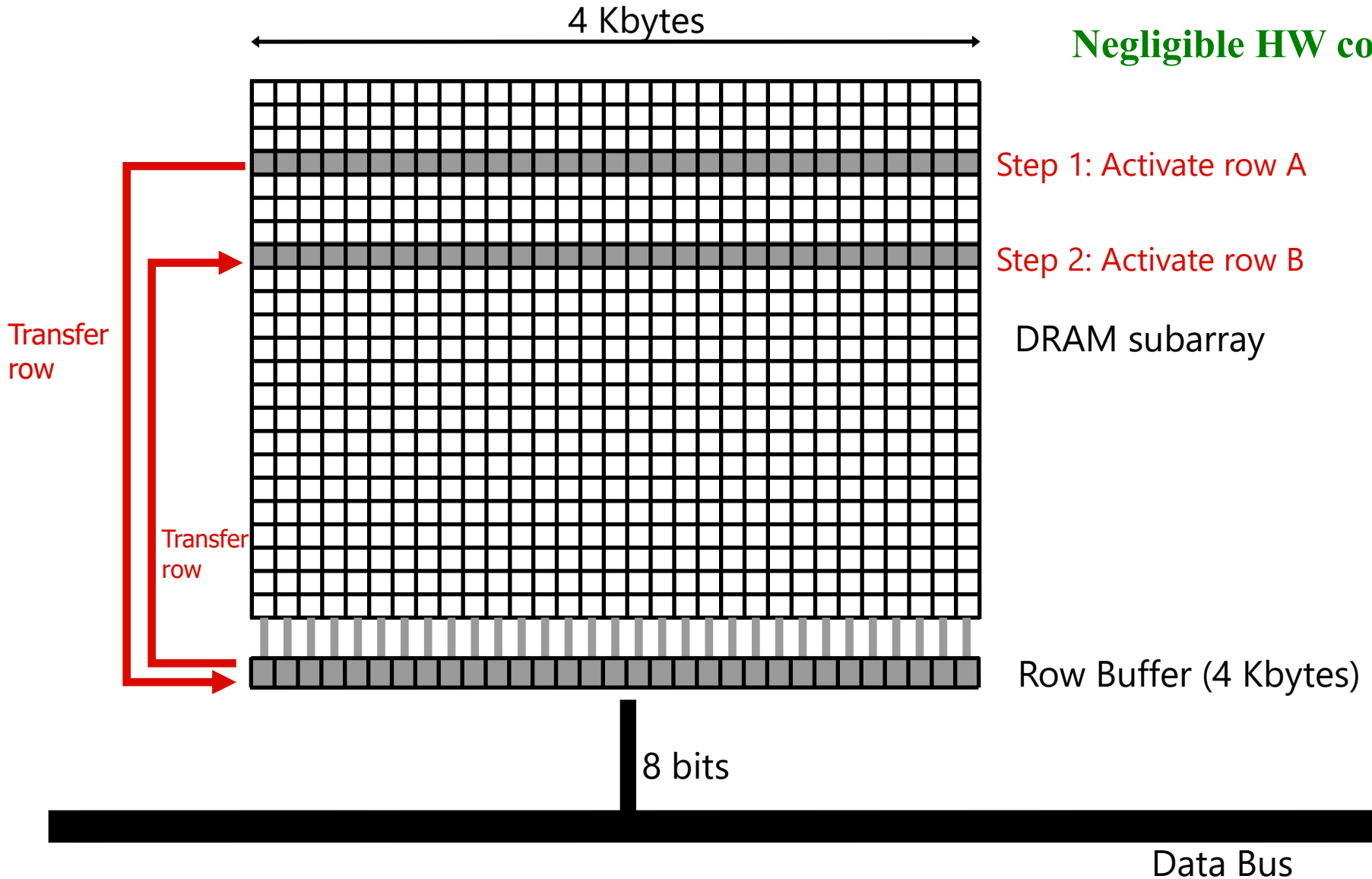**CPU**  **L1**  **L2**  **L3**  **MC**

2) Low bandwidth utilization

4) No unwanted data movement

1046ns, 3.6uJ → 90ns, 0.04uJ

# RowClone: In-DRAM Row Copy

**Idea: Two consecutive ACTivates**

**Negligible HW cost**

4 Kbytes

Step 1: Activate row A

Step 2: Activate row B

DRAM subarray

Transfer row

Transfer row

Row Buffer (4 Kbytes)

8 bits

Data Bus

# RowClone: Intra-Subarray



$V_{DD}/2$ $V_{DD}$ $\delta$

$V_{DD}/2 + \delta$

src    0

dst    0

Data gets
copied

Amplify the
difference

Sense Amplifier
(Row Buffer)

$V_{DD}/2$
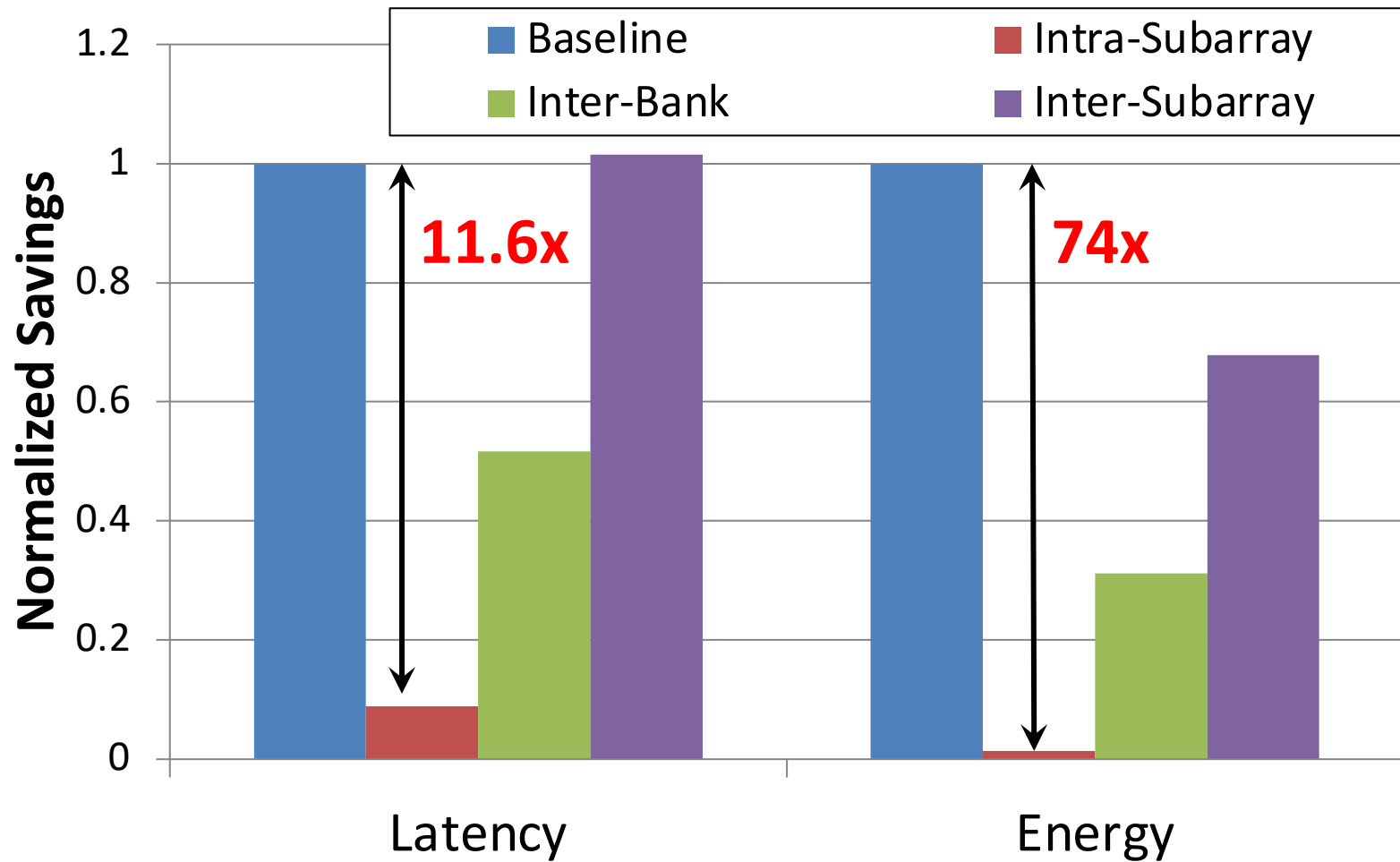0

# RowClone: Latency and Energy Savings



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

# More on RowClone

- Vivek Seshadri, Yoongu Kim, Chris Fallin, Donghyuk Lee, Rachata Ausavarungnirun, Gennady Pekhimenko, Yixin Luo, Onur Mutlu, Michael A. Kozuch, Phillip B. Gibbons, and Todd C. Mowry,
**"RowClone: Fast and Energy-Efficient In-DRAM Bulk Data Copy and Initialization"**
*Proceedings of the 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

# RowClone Extensions and Follow-Up Work

- Can we do faster inter-subarray copy?
  - Yes, see LISA [Chang et al., HPCA 2016]

- Can we enable data movement at smaller granularities within a bank?
  - Yes, see FIGARO [Wang et al., MICRO 2020]

- Can we do better inter-bank copy?
  - Yes, see Network-on-Memory [CAL 2020]

- Can similar ideas and DRAM properties be used to perform computation on data?
  - Yes, see Ambit [Seshadri et al., CAL 2015, MICRO 2017]

# LISA: Increasing Connectivity in DRAM

- Kevin K. Chang, Prashant J. Nair, Saugata Ghose, Donghyuk Lee, Moinuddin K. Qureshi, and Onur Mutlu,
  **"Low-Cost Inter-Linked Subarrays (LISA): Enabling Fast Inter-Subarray Data Movement in DRAM"**
  *Proceedings of the 22nd International Symposium on High-Performance Computer Architecture* (**HPCA**), Barcelona, Spain, March 2016.
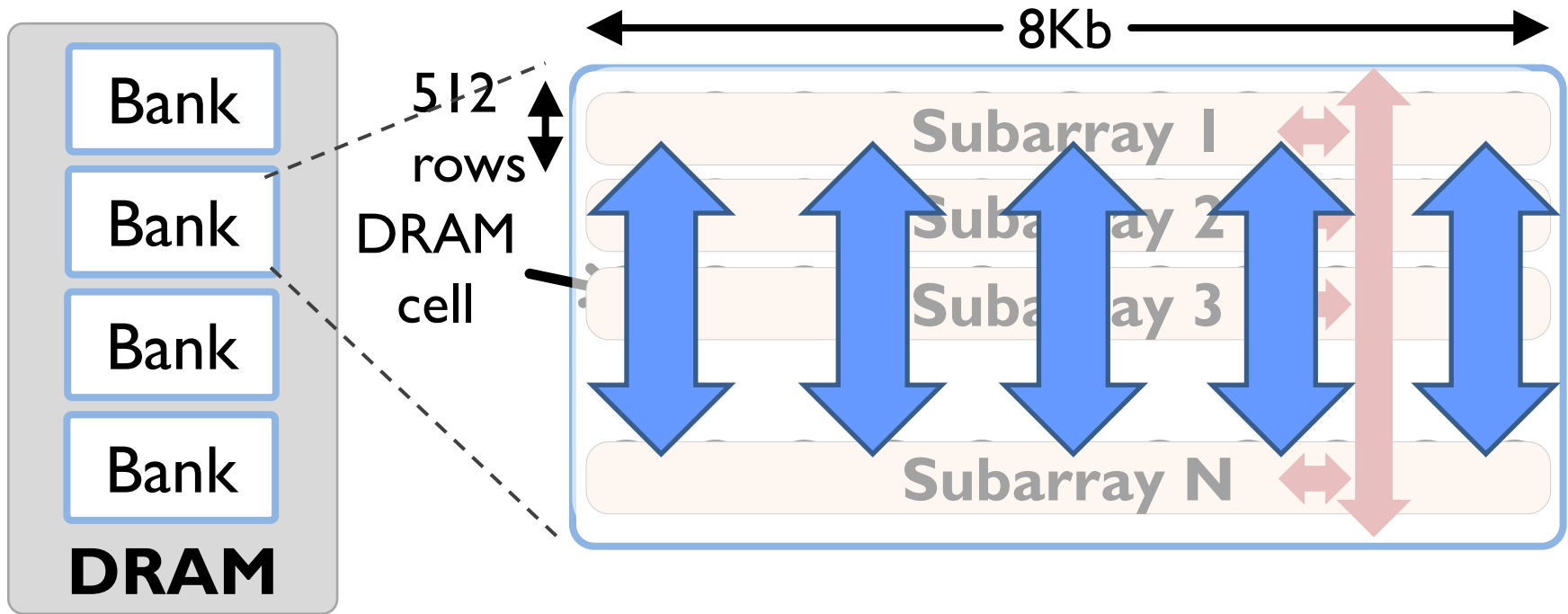  [Slides (pptx) (pdf)]
  [Source Code]

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

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

[†]Carnegie Mellon University    [⋆]Georgia Institute of Technology

# Moving Data Inside DRAM?



**DRAM**

Bank
Bank
Bank
Bank

8Kb

512 rows
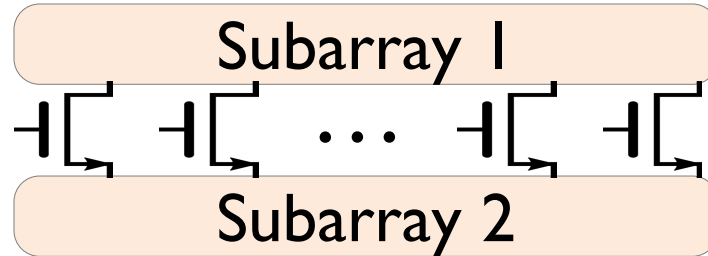
DRAM cell

Subarray 1
Subarray 2
Subarray 3
Subarray N

**Goal: Provide a new substrate to enable wide connectivity between subarrays**

# Key Idea and Applications

- **Low-cost Inter-linked subarrays (LISA)**
  - Fast bulk data movement between subarrays
  - Wide datapath via isolation transistors: 0.8% DRAM chip area



- LISA is a **versatile substrate** → new applications

  Fast bulk data copy: Copy latency 1.363ms→0.148ms (9.2x)
  → 66% speedup, -55% DRAM energy

  In-DRAM caching: Hot data access latency 48.7ns→21.5ns (2.2x)
  → 5% speedup

  Fast precharge: Precharge latency 13.1ns→5.0ns (2.6x)
  → 8% speedup

# More on LISA

- Kevin K. Chang, Prashant J. Nair, Saugata Ghose, Donghyuk Lee, Moinuddin K. Qureshi, and Onur Mutlu,
**"Low-Cost Inter-Linked Subarrays (LISA): Enabling Fast Inter-Subarray Data Movement in DRAM"**
*Proceedings of the 22nd International Symposium on High-Performance Computer Architecture* (**HPCA**), Barcelona, Spain, March 2016.
[Slides (pptx) (pdf)]
[Source Code]

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

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

[†]*Carnegie Mellon University*     [⋆]*Georgia Institute of Technology*

# FIGARO: Fine-Grained In-DRAM Copy

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

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

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

[★]*National University of Defense Technology*   [†]*ETH Zürich*   [⊙]*Chinese University of Hong Kong*
[◇]*University of Illinois at Urbana–Champaign*   [‡]*Carnegie Mellon University*   [§]*Institute of Research in Fundamental Sciences*

# Network-On-Memory: Fast Inter-Bank Copy

- Seyyed Hossein SeyyedAghaei Rezaei, Mehdi Modarressi, Rachata Ausavarungnirun, Mohammad Sadrosadati, Onur Mutlu, and Masoud Daneshtalab,
**"NoM: Network-on-Memory for Inter-Bank Data Transfer in Highly-Banked Memories"**
*IEEE Computer Architecture Letters* (**CAL**), to appear in 2020.

## NoM: NETWORK-ON-MEMORY FOR INTER-BANK DATA TRANSFER IN HIGHLY-BANKED MEMORIES

Seyyed Hossein SeyyedAghaei Rezaei[1]     Mehdi Modarressi[1,3]     Rachata Ausavarungnirun[2]
Mohammad Sadrosadati[3]     Onur Mutlu[4]     Masoud Daneshtalab[5]

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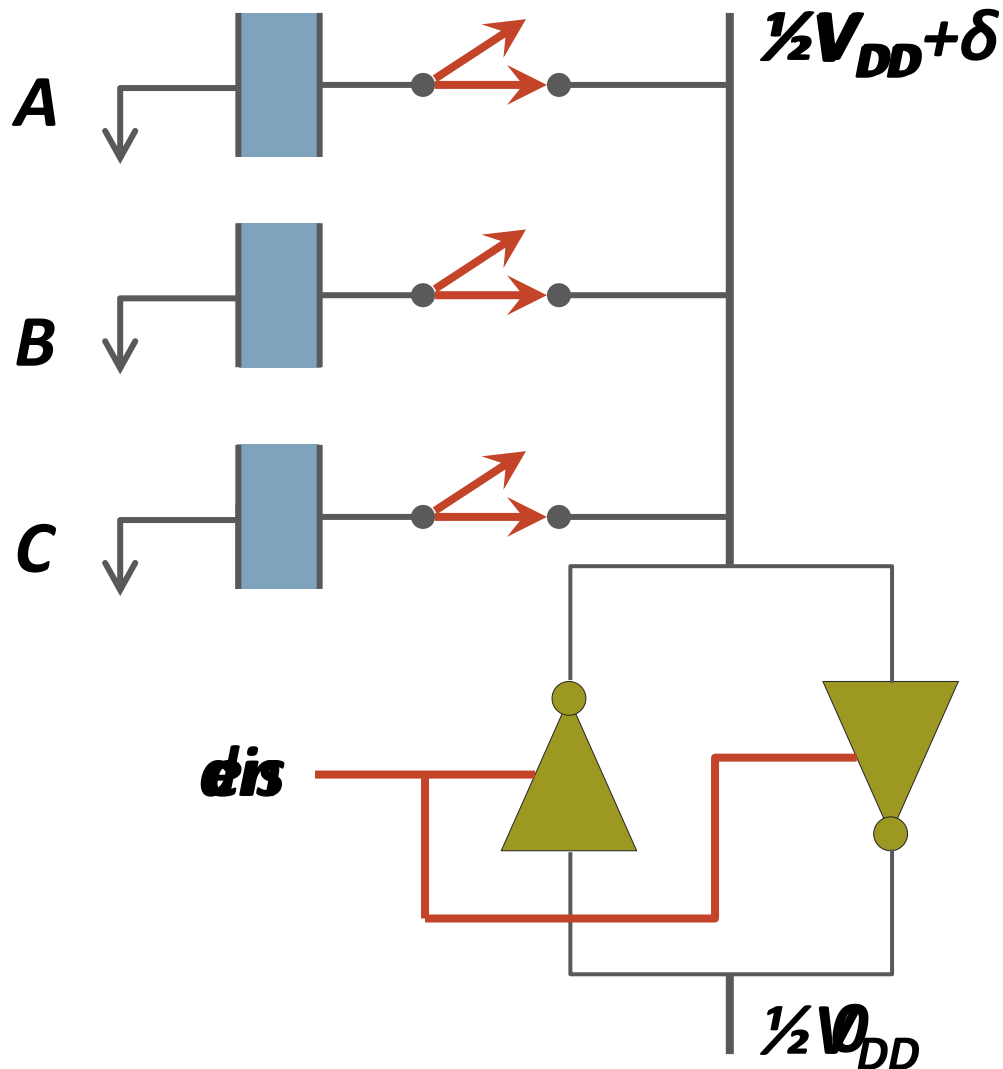[4]ETH Zürich     [5]Mälardalens University

# (Truly) In-Memory Computation

- We can support in-DRAM AND, OR, NOT, MAJ
- At low cost
- Using analog computation capability of DRAM
  - Idea: activating multiple rows performs computation
- 30-60X performance and energy improvement
  - Seshadri+, "Ambit: In-Memory Accelerator for Bulk Bitwise Operations Using Commodity DRAM Technology," MICRO 2017.


- New memory technologies enable even more opportunities
  - Memristors, resistive RAM, phase change mem, STT-MRAM, …
  - Can operate on data with minimal movement

# In-DRAM AND/OR: Triple Row Activation



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

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.

# In-DRAM Bulk Bitwise AND/OR Operation

- **BULKAND A, B → C**

- Semantics: Perform a bitwise AND of two rows A and B and store the result in row C

- R0 – reserved zero row, R1 – reserved one row
- D1, D2, D3 – Designated rows for triple activation

1. RowClone  A  into  D1
2. RowClone  B  into  D2
3. RowClone  R0  into  D3
4. ACTIVATE  D1,D2,D3
5. RowClone  Result  into  C

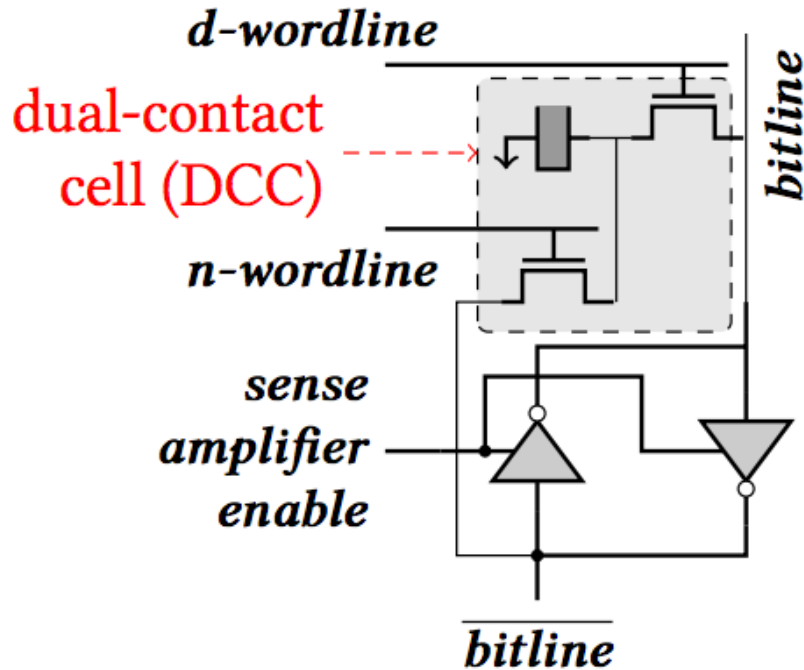**SAFARI**

# More on In-DRAM Bulk AND/OR

- Vivek Seshadri, Kevin Hsieh, Amirali Boroumand, Donghyuk Lee, Michael A. Kozuch, Onur Mutlu, Phillip B. Gibbons, and Todd C. Mowry,
  **"Fast Bulk Bitwise AND and OR in DRAM"**
  *IEEE Computer Architecture Letters* (**CAL**), April 2015.

## Fast Bulk Bitwise AND and OR in DRAM

Vivek Seshadri*, Kevin Hsieh*, Amirali Boroumand*, Donghyuk Lee*,
Michael A. Kozuch[†], Onur Mutlu*, Phillip B. Gibbons[†], Todd C. Mowry*

*Carnegie Mellon University        [†]Intel Pittsburgh

# In-DRAM NOT: Dual Contact Cell



d-wordline

dual-contact cell (DCC)

n-wordline

bitline

sense amplifier enable

$\overline{bitline}$

**Figure 5:** A dual-contact cell connected to both ends of a sense amplifier

Idea:
Feed the negated value in the sense amplifier into a special row

Seshadri+, "Ambit: In-Memory Accelerator for Bulk Bitwise Operations using Commodity DRAM Technology," MICRO 2017.
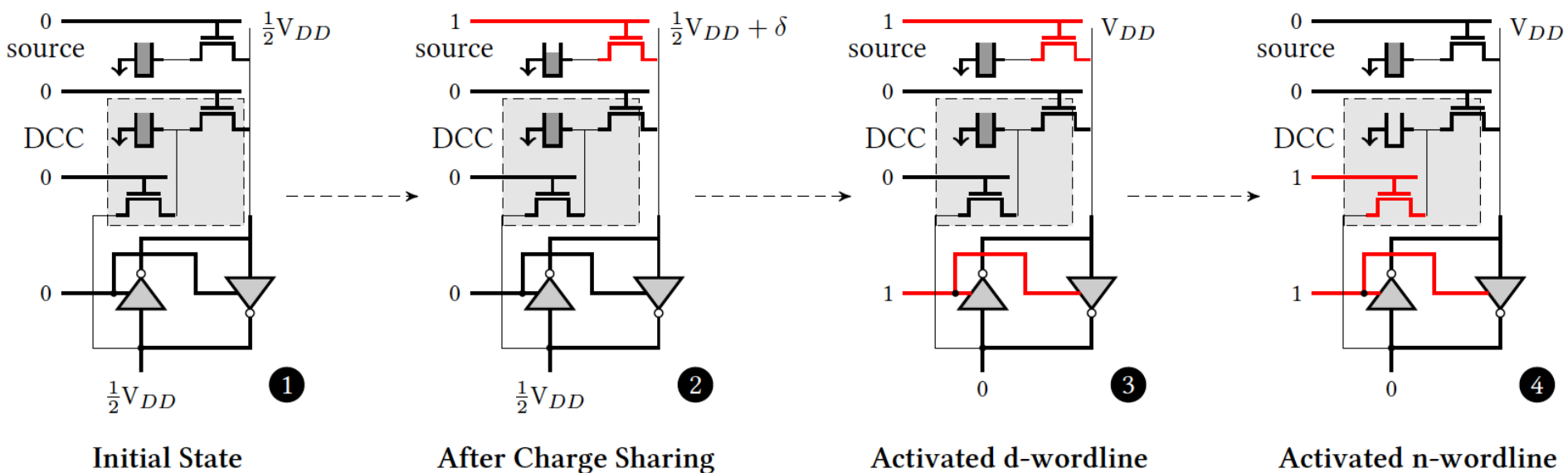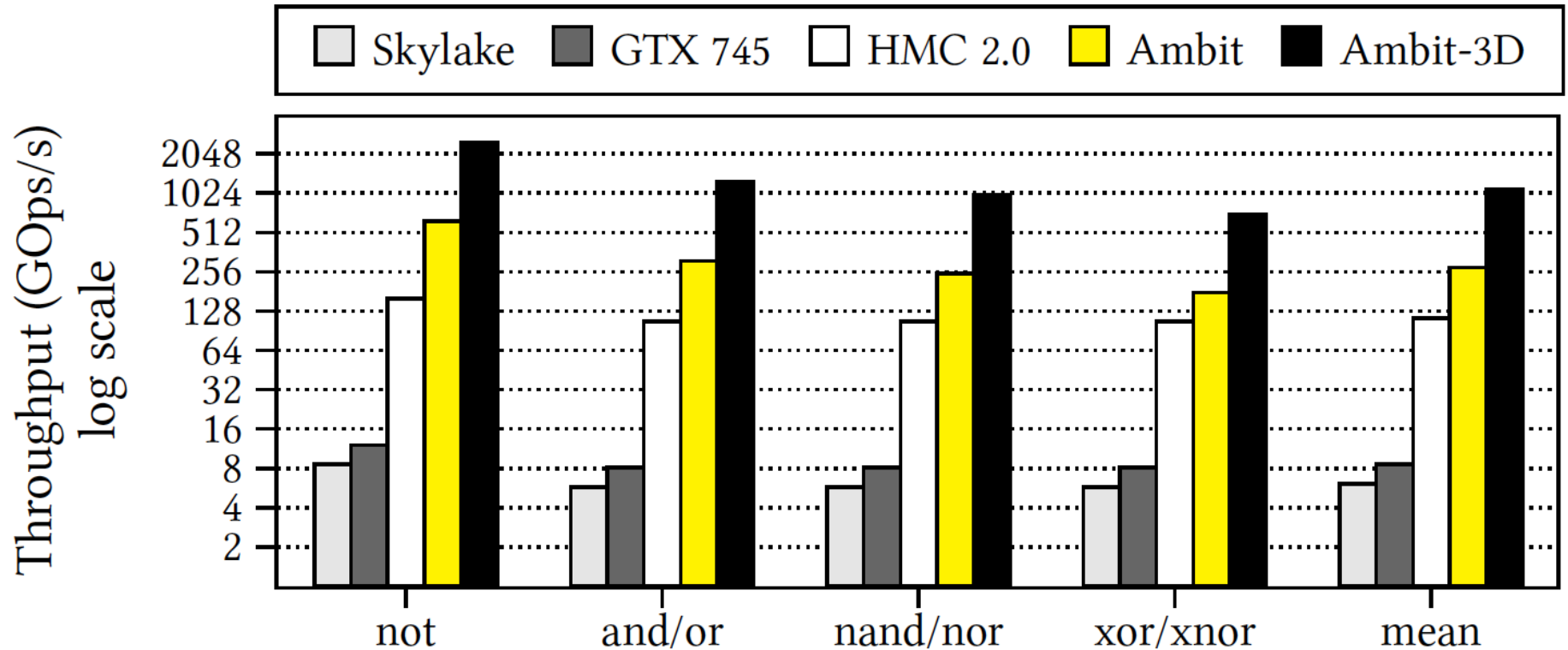
# In-DRAM NOT Operation



**Figure 5: Bitwise NOT using a dual contact capacitor**

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

# Performance: In-DRAM Bitwise Operations



**Figure 9: Throughput of bitwise operations on various systems.**
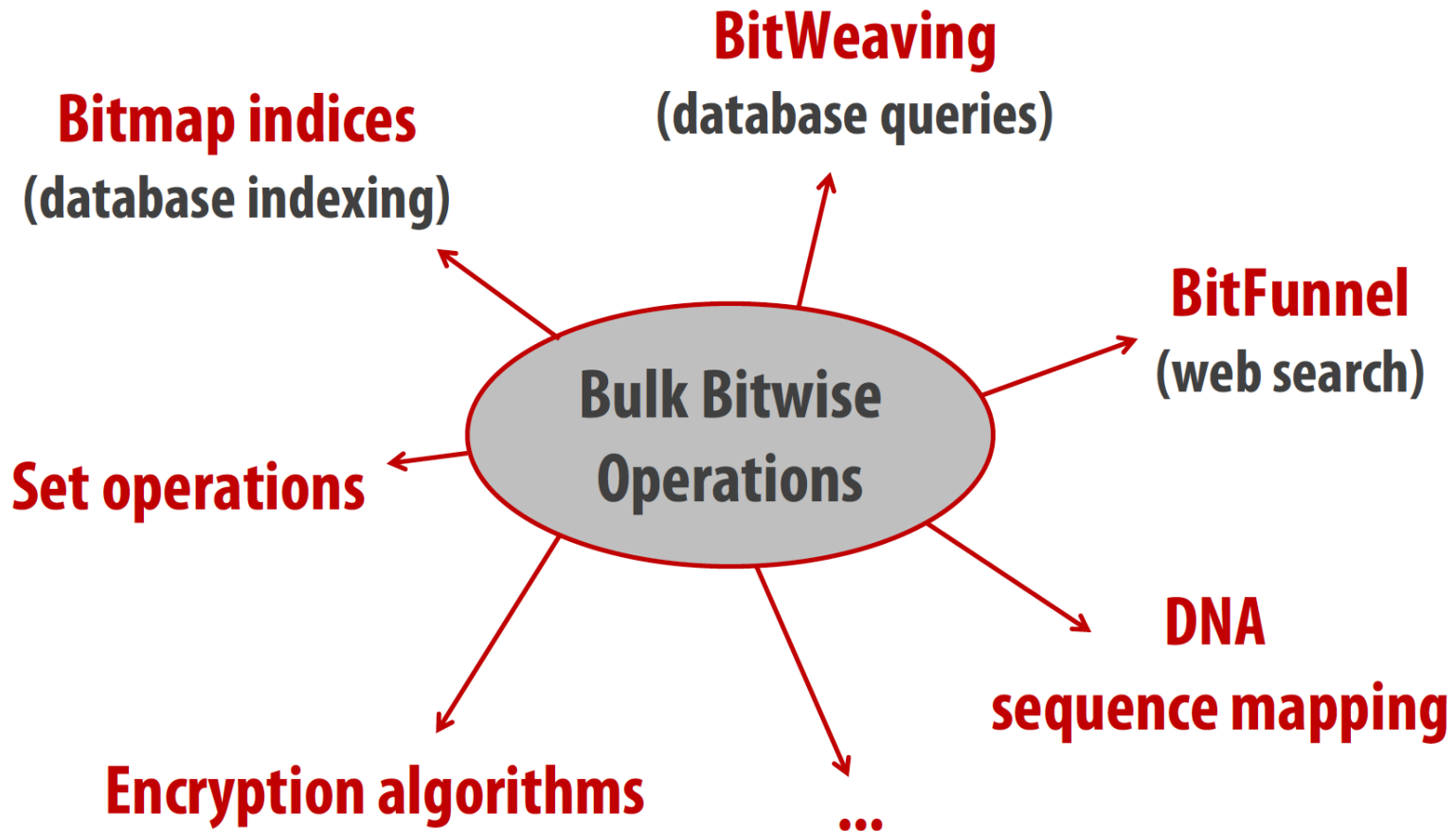
# Energy of In-DRAM Bitwise Operations

|  | Design | not | and/or | nand/nor | xor/xnor |
|---|---|---|---|---|---|
| DRAM & | **DDR3** | 93.7 | 137.9 | 137.9 | 137.9 |
| Channel Energy | **Ambit** | 1.6 | 3.2 | 4.0 | 5.5 |
| (nJ/KB) | (↓) | 59.5X | 43.9X | 35.1X | 25.1X |

**Table 3: Energy of bitwise operations. (↓) indicates energy reduction of Ambit over the traditional DDR3-based design.**
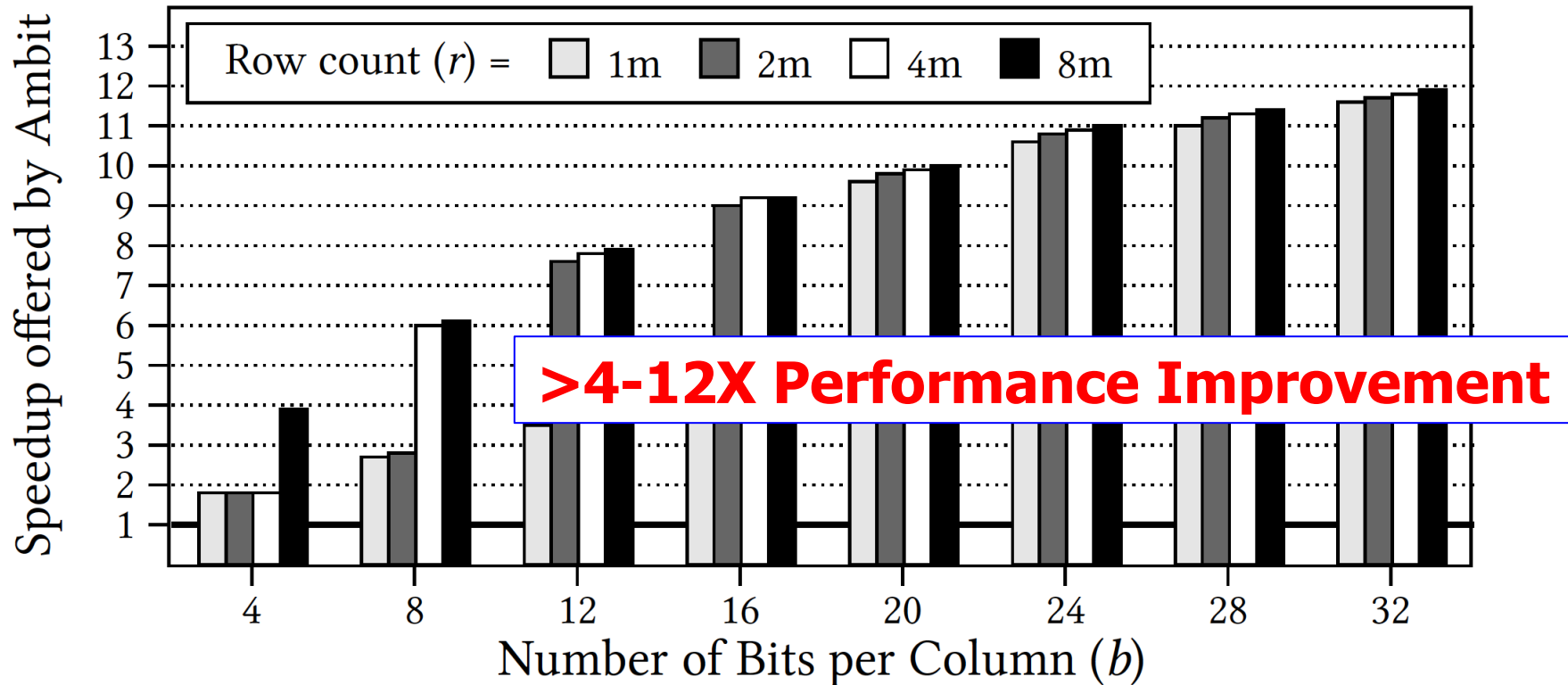
Seshadri+, "Ambit: In-Memory Accelerator for Bulk Bitwise Operations using Commodity DRAM Technology," MICRO 2017.

**SAFARI**

# Bulk Bitwise Operations in Workloads



**BitWeaving**
(database queries)

**Bitmap indices**
(database indexing)

**BitFunnel**
(web search)

**Bulk Bitwise Operations**

**Set operations**

**DNA sequence mapping**

**Encryption algorithms**

**...**

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

# In-DRAM Acceleration of Database Queries



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

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

# More on 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**

# In-DRAM Bulk Bitwise Execution

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

## In-DRAM Bulk Bitwise Execution Engine

Vivek Seshadri
Microsoft Research India
visesha@microsoft.com

Onur Mutlu
ETH Zürich
onur.mutlu@inf.ethz.ch

# SIMDRAM Framework

- Nastaran Hajinazar, Geraldo F. Oliveira, Sven Gregorio, Joao Dinis Ferreira, Nika Mansouri Ghiasi, Minesh Patel, Mohammed Alser, Saugata Ghose, Juan Gomez-Luna, and Onur Mutlu, **"SIMDRAM: An End-to-End Framework for Bit-Serial SIMD Computing in DRAM"** *Proceedings of the 26th International Conference on Architectural Support for Programming Languages and Operating Systems* (**ASPLOS**), Virtual, March-April 2021. [2-page Extended Abstract] [Short Talk Slides (pptx) (pdf)] [Talk Slides (pptx) (pdf)] [Short Talk Video (5 mins)] [Full Talk Video (27 mins)]

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

*Nastaran Hajinazar[1,2]    *Geraldo F. Oliveira[1]    Sven Gregorio[1]    João Dinis Ferreira[1]
Nika Mansouri Ghiasi[1]    Minesh Patel[1]    Mohammed Alser[1]    Saugata Ghose[3]
Juan Gómez-Luna[1]    Onur Mutlu[1]

[1]ETH Zürich    [2]Simon Fraser University    [3]University of Illinois at Urbana–Champaign
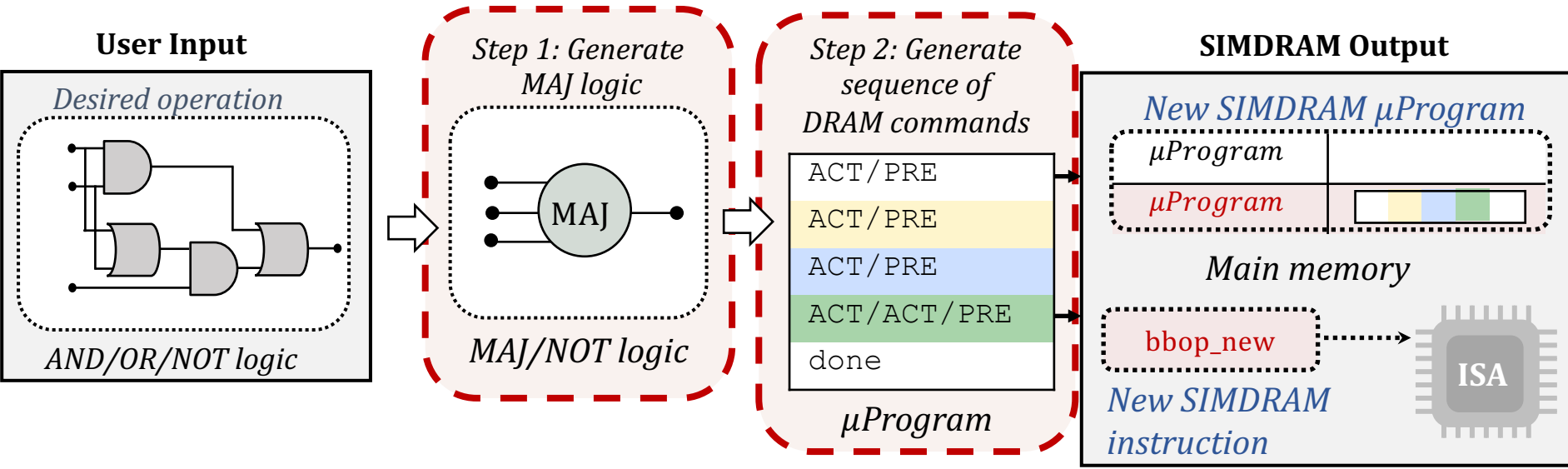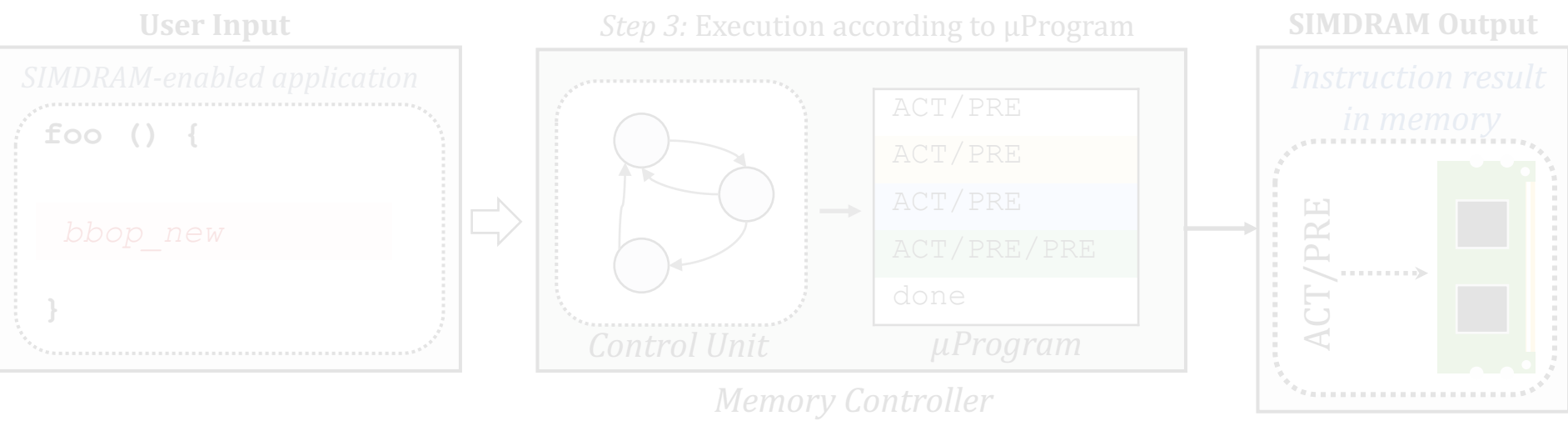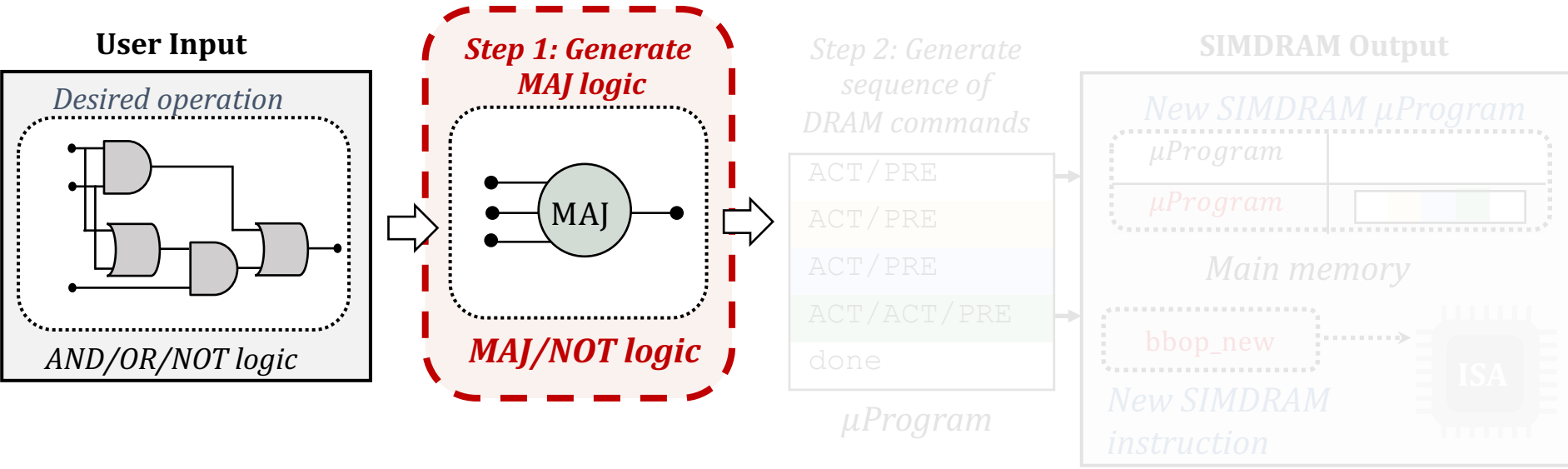
# SIMDRAM Framework: Overview



**User Input**

*Desired operation*

*AND/OR/NOT logic*

**Step 1:** *Generate MAJ logic*

MAJ

*MAJ/NOT logic*

**Step 2:** *Generate sequence of DRAM commands*

| |
|---|
| ACT/PRE |
| ACT/PRE |
| ACT/PRE |
| ACT/ACT/PRE |
| done |

*µProgram*

**SIMDRAM Output**

*New SIMDRAM µProgram*

µProgram

µProgram

*Main memory*

bbop_new → ISA

*New SIMDRAM instruction*

**User Input**

*SIMDRAM-enabled application*

```
foo () {

  bbop_new

}
```

**Step 3:** Execution according to µProgram

Control Unit

| |
|---|
| ACT/PRE |
| ACT/PRE |
| ACT/PRE |
| ACT/PRE/PRE |
| done |

*µProgram*

*Memory Controller*

**SIMDRAM Output**

*Instruction result in memory*

ACT/PRE

# SIMDRAM Framework: Step 1



**User Input**

*Desired operation*

*AND/OR/NOT logic*

**Step 1: Generate MAJ logic**

MAJ

**MAJ/NOT logic**

*Step 2: Generate sequence of DRAM commands*

ACT/PRE
ACT/PRE
ACT/PRE
ACT/ACT/PRE
done

*μProgram*

**SIMDRAM Output**

*New SIMDRAM μProgram*

*μProgram*
*μProgram*

*Main memory*

bbop_new

ISA

*New SIMDRAM instruction*

**User Input**

*SIMDRAM-enabled application*

```
foo () {

    bbop_new

}
```

*Step 3:* Execution according to μProgram

ACT/PRE
ACT/PRE
ACT/PRE
ACT/PRE/PRE
done

*Control Unit*          *μProgram*

*Memory Controller*

**SIMDRAM Output**

*Instruction result in memory*

ACT/PRE

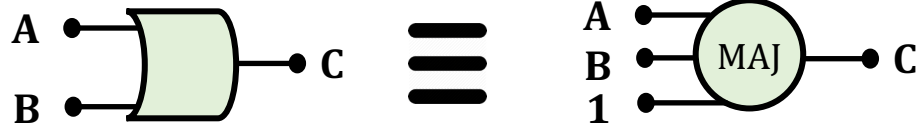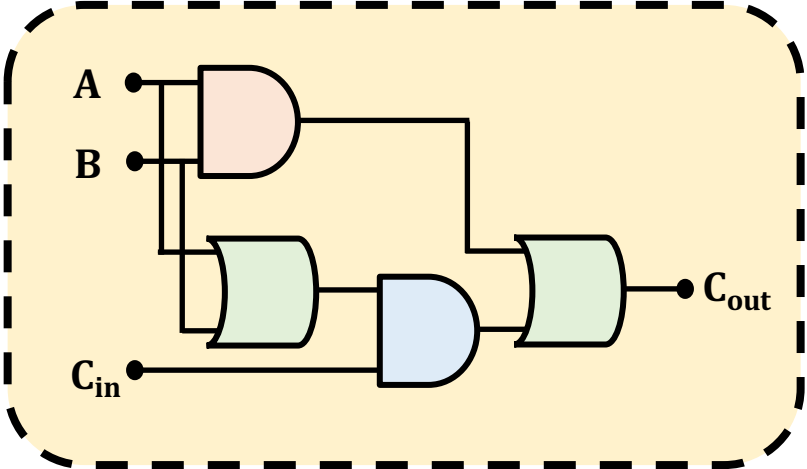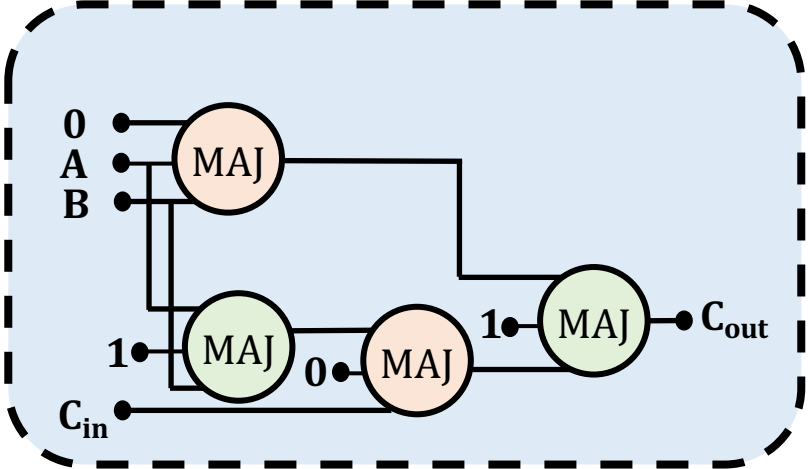**SAFARI**

# Step 1: Naïve MAJ/NOT Implementation



output is "1" only when A = B = "1"

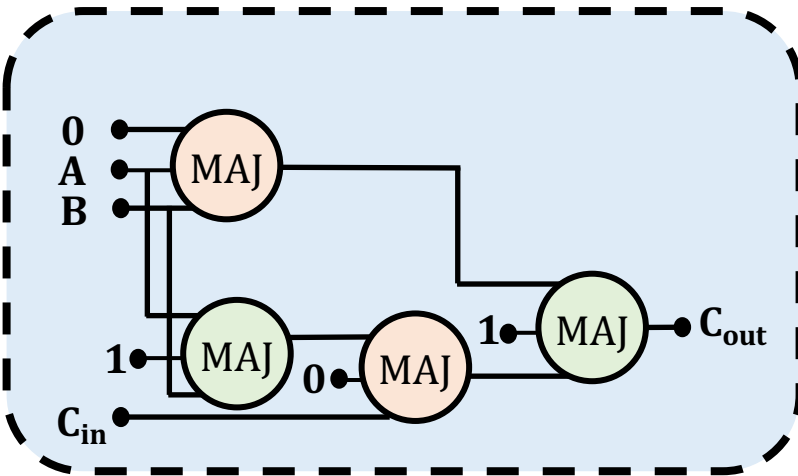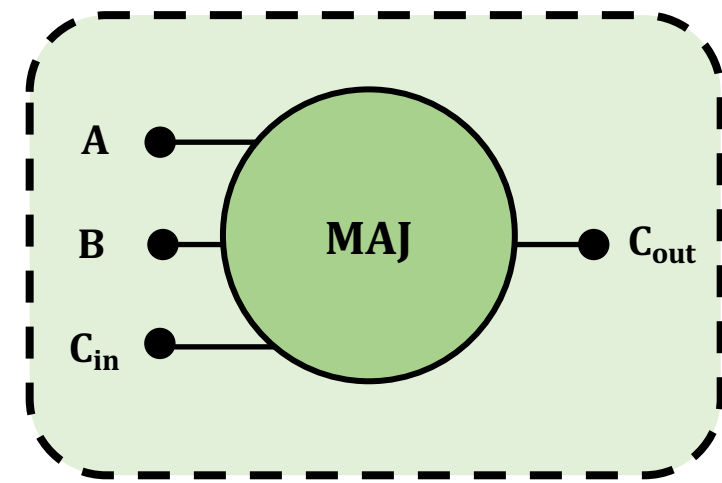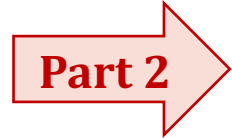output is "0" only when A = B = "0"

**Part 1**

**Naïvely** converting **AND/OR/NOT-implementation** to **MAJ/NOT-implementation** leads to an **unoptimized circuit**

# Step 1: Efficient MAJ/NOT Implementation



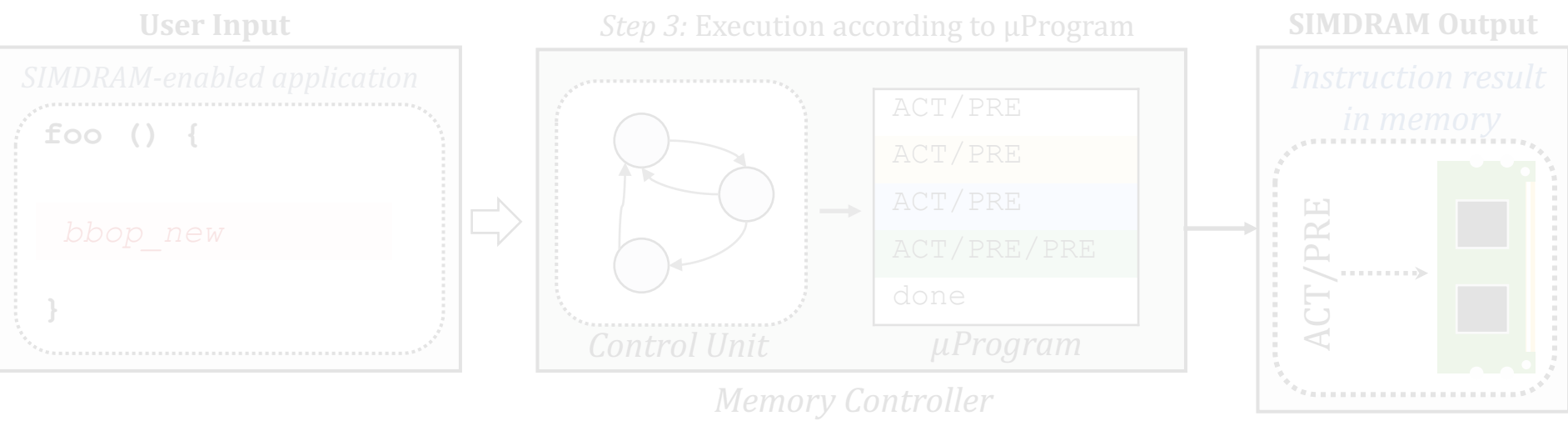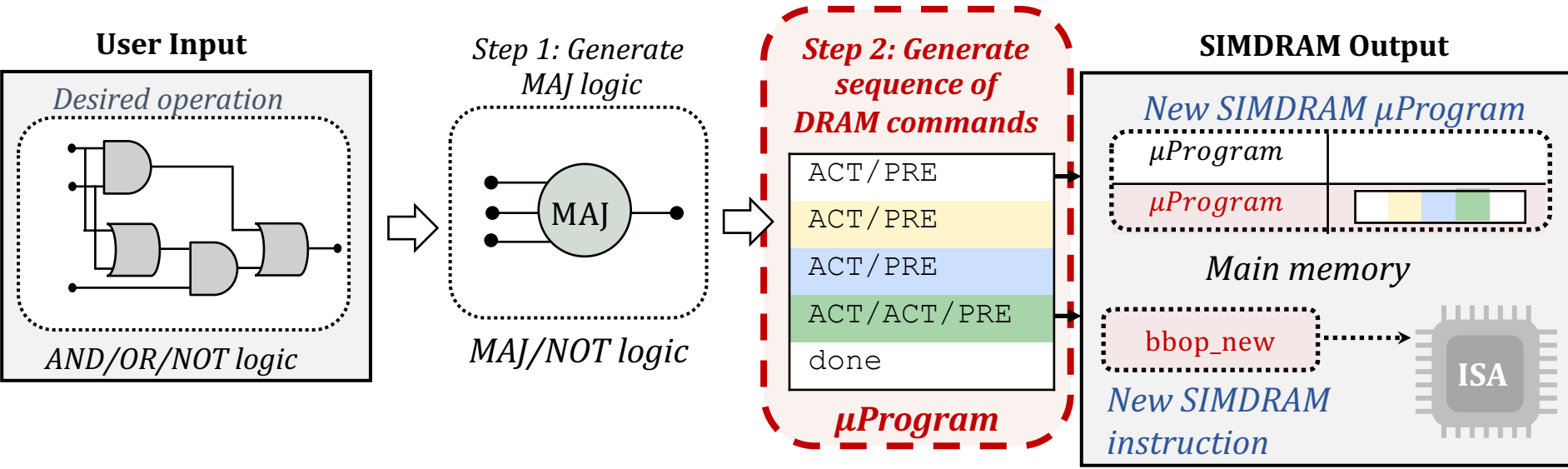Greedy optimization algorithm[4]

Part 2

**Step 1 generates an optimized MAJ/NOT-implementation of the desired operation**

[4] L. Amarù et al, "*Majority-Inverter Graph: A Novel Data-Structure and Algorithms for Efficient Logic Optimization*", DAC, 2014.

# SIMDRAM Framework: Step 2



**User Input**

*Desired operation*

*AND/OR/NOT logic*

*Step 1: Generate MAJ logic*

MAJ

*MAJ/NOT logic*

**Step 2: Generate sequence of DRAM commands**

| |
|---|
| ACT/PRE |
| ACT/PRE |
| ACT/PRE |
| ACT/ACT/PRE |
| done |

*µProgram*

**SIMDRAM Output**

*New SIMDRAM µProgram*

µProgram

*µProgram*

*Main memory*

bbop_new

ISA

*New SIMDRAM instruction*

**User Input**

*SIMDRAM-enabled application*

```
foo () {

    bbop_new

}
```

*Step 3: Execution according to µProgram*

| |
|---|
| ACT/PRE |
| ACT/PRE |
| ACT/PRE |
| ACT/PRE/PRE |
| done |

*Control Unit*  *µProgram*

*Memory Controller*
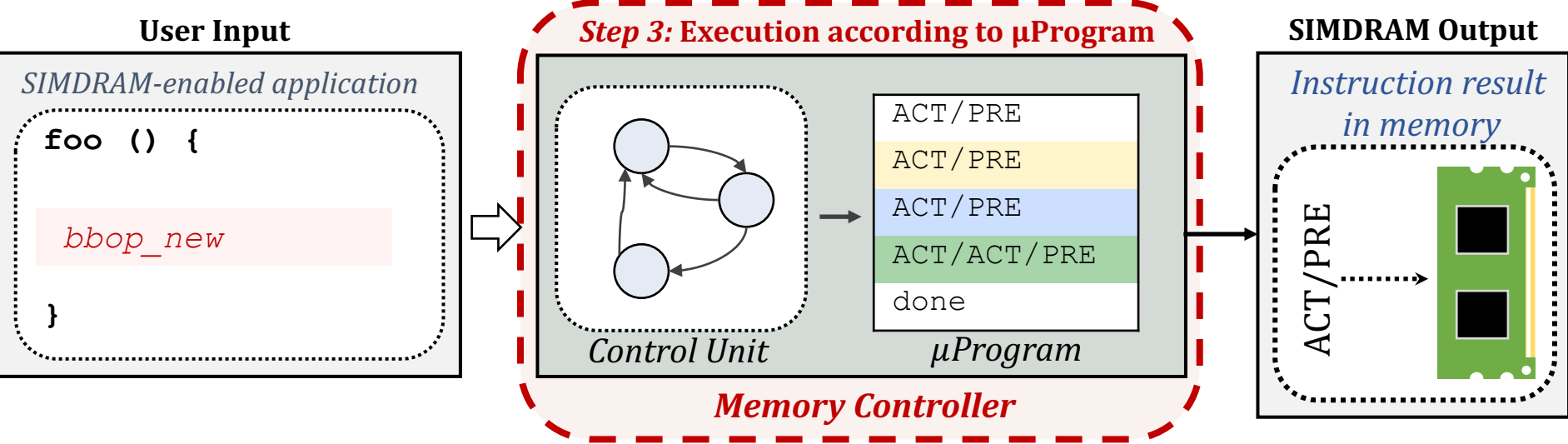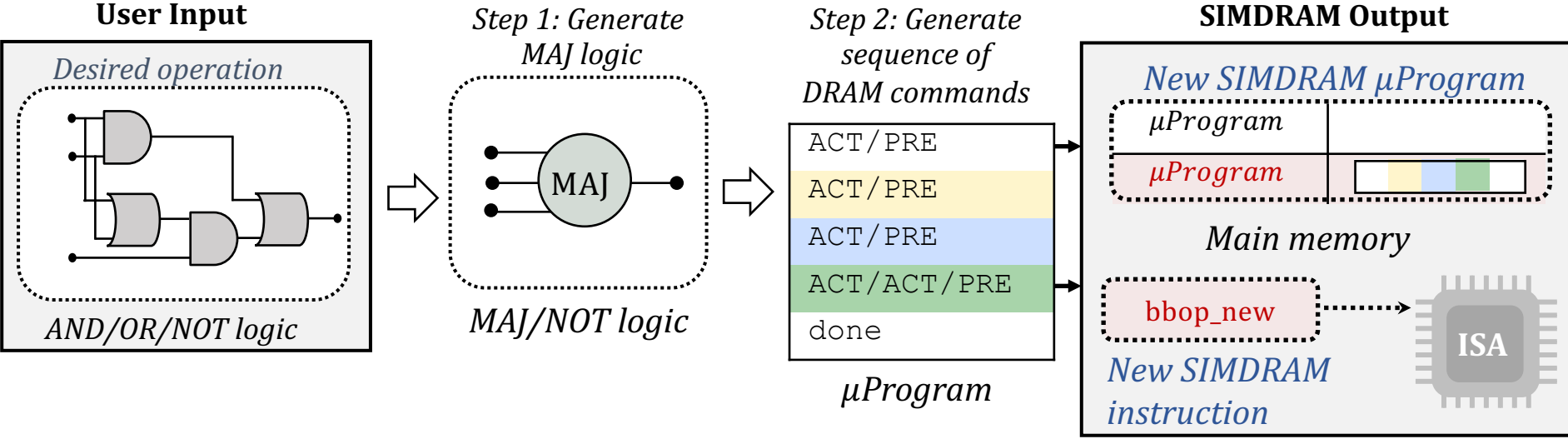
**SIMDRAM Output**

*Instruction result in memory*

ACT/PRE

**SAFARI**

# Step 2: µProgram Generation

- **µProgram:** A series of microarchitectural operations (e.g., ACT/PRE) that SIMDRAM uses to execute SIMDRAM operation in DRAM

- **Goal of Step 2**: To generate the µProgram that executes the desired SIMDRAM operation in DRAM
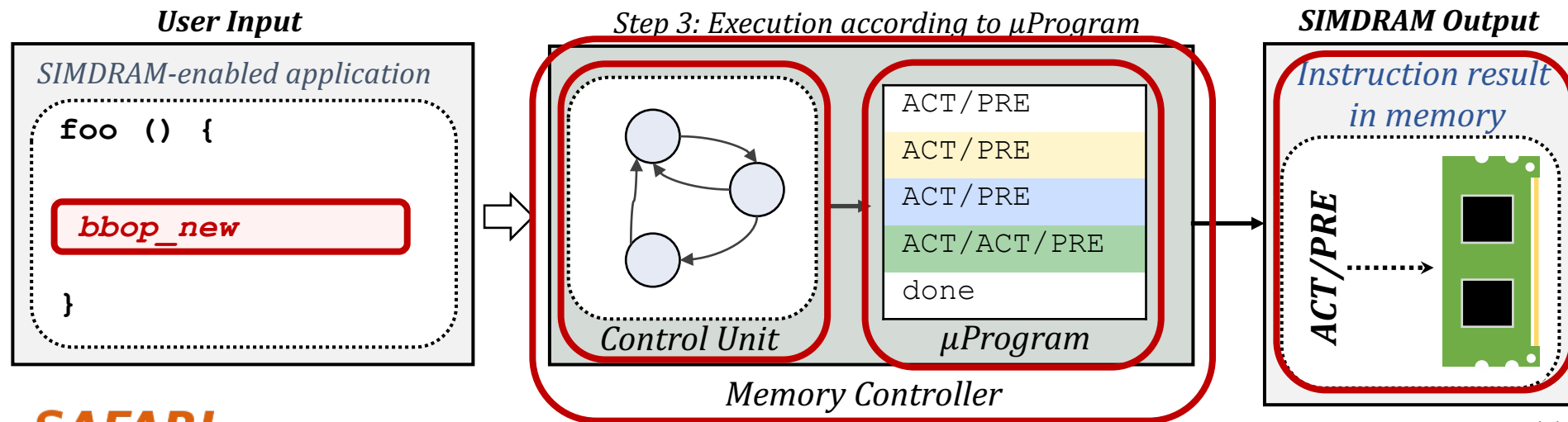
**Task 1: Allocate DRAM rows to the operands**

**Task 2: Generate µProgram**

# SIMDRAM Framework: Step 3

**User Input**

*Desired operation*

*AND/OR/NOT logic*

*Step 1: Generate MAJ logic*

MAJ

*MAJ/NOT logic*

*Step 2: Generate sequence of DRAM commands*

| |
|---|
| ACT/PRE |
| ACT/PRE |
| ACT/PRE |
| ACT/ACT/PRE |
| done |

*µProgram*

**SIMDRAM Output**

*New SIMDRAM µProgram*

| µProgram | |
|---|---|
| µProgram | |

*Main memory*

bbop_new → ISA

*New SIMDRAM instruction*

---

**User Input**

*SIMDRAM-enabled application*

```
foo () {

  bbop_new

}
```

**Step 3:** Execution according to µProgram

| |
|---|
| ACT/PRE |
| ACT/PRE |
| ACT/PRE |
| ACT/ACT/PRE |
| done |

*Control Unit*          *µProgram*

**Memory Controller**

**SIMDRAM Output**

*Instruction result in memory*

ACT/PRE

**SAFARI**

45

# Step 3: μProgram Execution

- **SIMDRAM control unit:** handles the execution of the μProgram at runtime

- Upon receiving a **bbop instruction**, the control unit:

  1. Loads the μProgram corresponding to SIMDRAM operation

  2. Issues the sequence of DRAM commands (ACT/PRE) stored in the μProgram to SIMDRAM subarrays to perform the in-DRAM operation



**User Input**

*SIMDRAM-enabled application*

```
foo () {

    bbop_new

}
```

*Step 3: Execution according to μProgram*

ACT/PRE
ACT/PRE
ACT/PRE
ACT/ACT/PRE
done

*Control Unit*          *μProgram*

*Memory Controller*

**SIMDRAM Output**

*Instruction result in memory*

*ACT/PRE*

# More in the Paper

**https://people.inf.ethz.ch/omutlu/pub/SIMDRAM_asplos21.pdf**

### SIMDRAM: An End-to-End Framework for
### Bit-Serial SIMD Computing in DRAM

*Nastaran Hajinazar[1,2]    *Geraldo F. Oliveira[1]    Sven Gregorio[1]    João Dinis Ferreira[1]
Nika Mansouri Ghiasi[1]    Minesh Patel[1]    Mohammed Alser[1]    Saugata Ghose[3]
Juan Gómez-Luna[1]    Onur Mutlu[1]

[1]ETH Zürich    [2]Simon Fraser University    [3]University of Illinois at Urbana–Champaign

coherence, and interrupts

Handling limited subarray size

Security implications

Limitations of our framework

# SIMDRAM Key Results

Evaluated on:

- 16 complex in-DRAM operations
- 7 commonly-used real-world applications

**SIMDRAM provides:**

- **88×** and **5.8×** the **throughput** of a **CPU** and a **high-end GPU**, respectively, over **16 operations**

- **257×** and **31×** the **energy efficiency** of a **CPU** and a **high-end GPU**, respectively, over **16 operations**

- **21×** and **2.1×** the **performance** of a **CPU** an a **high-end GPU**, over **seven real-world applications**

**SAFARI**

# More on SIMDRAM

- Nastaran Hajinazar, Geraldo F. Oliveira, Sven Gregorio, Joao Dinis Ferreira, Nika Mansouri Ghiasi, Minesh Patel, Mohammed Alser, Saugata Ghose, Juan Gomez-Luna, and Onur Mutlu,
**"SIMDRAM: An End-to-End Framework for Bit-Serial SIMD Computing in DRAM"**
*Proceedings of the 26th International Conference on Architectural Support for Programming Languages and Operating Systems (**ASPLOS**), Virtual, March-April 2021.*
[2-page Extended Abstract]
[Short Talk Slides (pptx) (pdf)]
[Talk Slides (pptx) (pdf)]
[Short Talk Video (5 mins)]
[Full Talk Video (27 mins)]

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

*Nastaran Hajinazar[1,2]    *Geraldo F. Oliveira[1]    Sven Gregorio[1]    João Dinis Ferreira[1]
Nika Mansouri Ghiasi[1]    Minesh Patel[1]    Mohammed Alser[1]    Saugata Ghose[3]
Juan Gómez-Luna[1]    Onur Mutlu[1]

[1]ETH Zürich    [2]Simon Fraser University    [3]University of Illinois at Urbana–Champaign

# SIMDRAM: Follow-Ups

- **Limitations of current substrate?**
  - Computing granularity
  - Data layout conversion
  - High-latency bit-serial operations
  - Assembly-like programming model
  - Application scope
  - ...


- We are working on even better processing-using-memory substrates
  - One step at a time!

# Limitations of PUD Systems:
## Overview

**PUD systems suffer from three sources of inefficiency due to the large and rigid DRAM access granularity**

**1** **SIMD Underutilization**
- due to data parallelism variation within and across applications
- leads to throughput and energy waste

**2** **Limited Computation Support**
- due to a lack of low-cost interconnects across columns
- limits PUD operations to only parallel map constructs

**3** **Challenging Programming Model**
- due to a lack of compiler support for PUD systems
- creates a burden on programmers, limiting PUD adoption

**SAFARI**

# Problem & Goal

**Problem**

Processing-Using-DRAM's <u>large and rigid granularity</u> limits its applicability and efficiency for different applications

**Goal**

Design a <u>flexible PUD system</u> that overcomes the three limitations caused by large and rigid DRAM access granularity

SAFARI

# MIMDRAM:
## Key Idea (I)

**DRAM's hierarchical organization can enable <u>fine-grained access</u>**

DRAM mat

global wordline

row decoder

global sense amplifier

**Key Issue:**
on a DRAM access, the global wordline propagates across all DRAM mats

**Fine-Grained DRAM:**
**segments the global wordline to access individual DRAM mats**

# MIMDRAM:
## Key Idea (II)

**Fine-Grained DRAM:**
**segments the global wordline to access individual DRAM mats**



segmented global wordline

row decoder

global sense amplifier

## Fine-grained DRAM for energy-efficient DRAM access:

**[Cooper-Balis+, 2010]:** Fine-Grained Activation for Power Reduction in DRAM
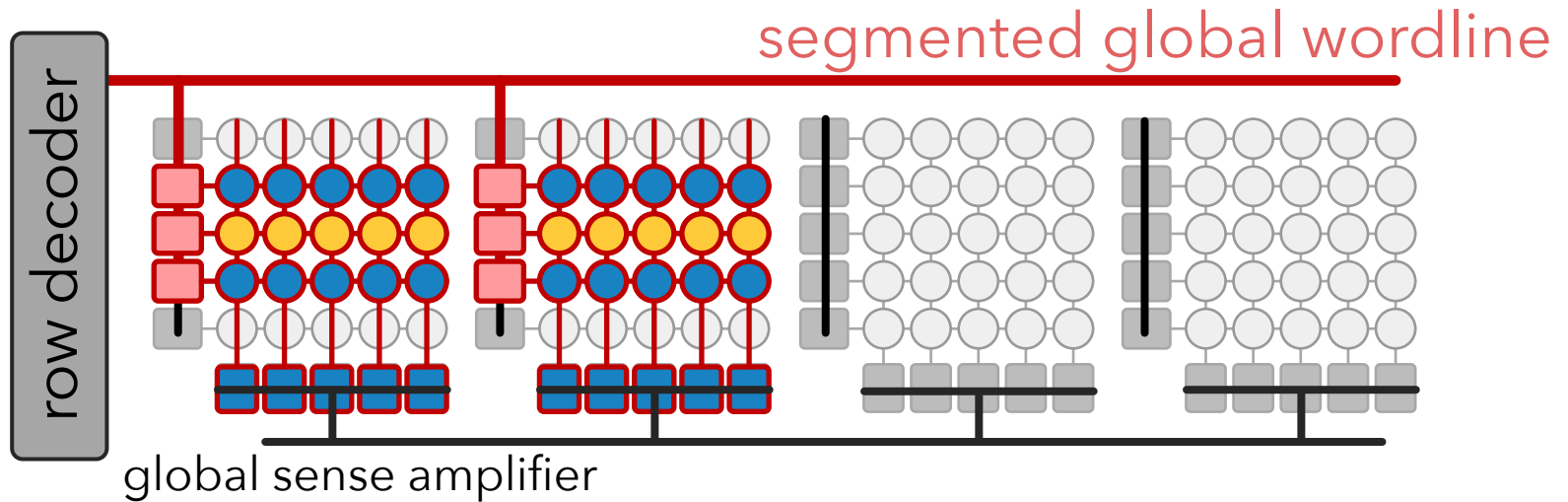**[Udipi+, 2010]:** Rethinking DRAM Design and Organization for Energy-Constrained Multi-Cores
**[Zhang+, 2014]:** Half-DRAM
**[Ha+, 2016]:** Improving Energy Efficiency of DRAM by Exploiting Half Page Row Access
**[O'Connor+, 2017]:** Fine-Grained DRAM
**[Olgun+, 2024]:** Sectored DRAM

**SAFARI**

segmented global wordline

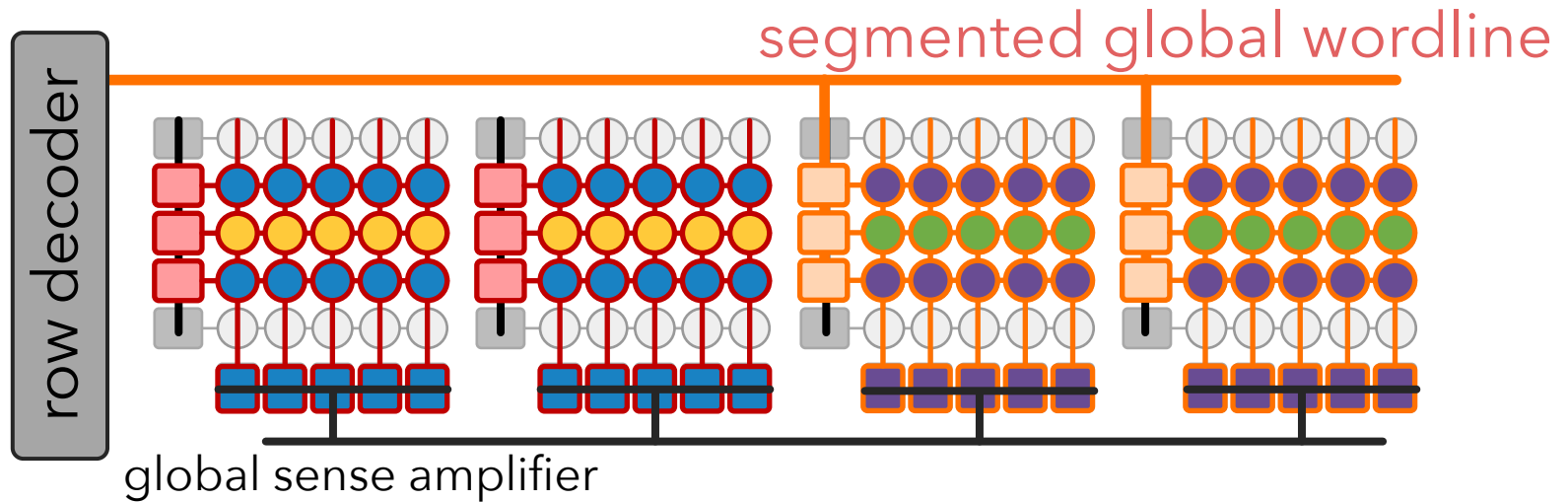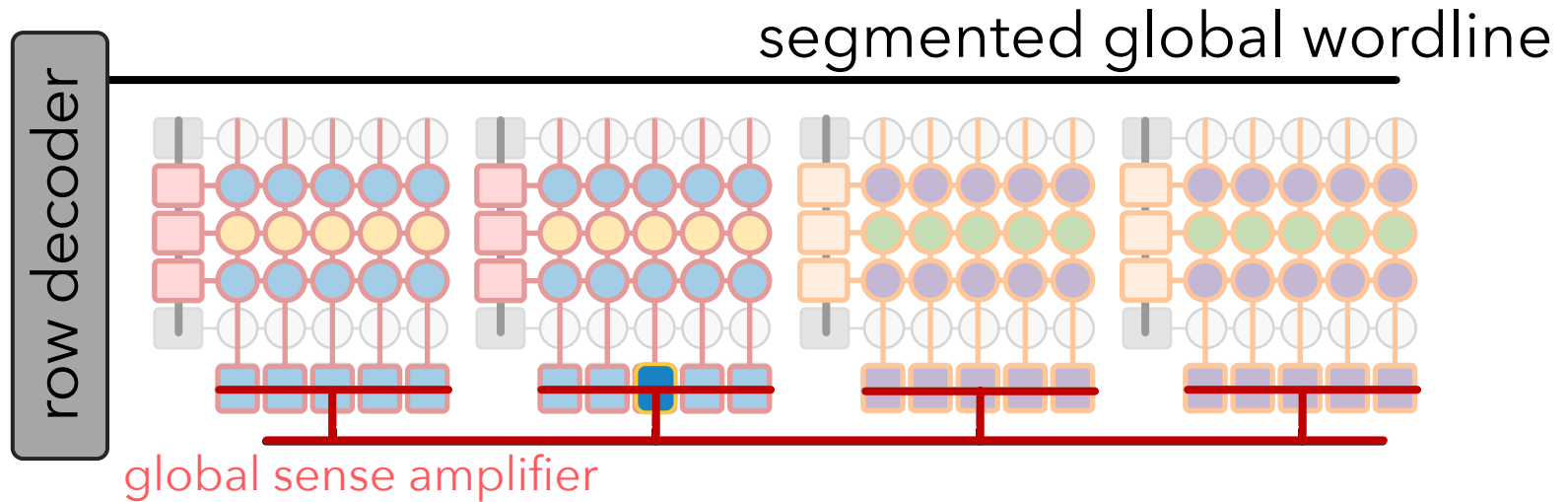row decoder

global sense amplifier

## Fine-grained DRAM for processing-using-DRAM:

**1** **Improves SIMD utilization**
- for a single PUD operation, only access the DRAM mats with target data

segmented global wordline

row decoder

global sense amplifier
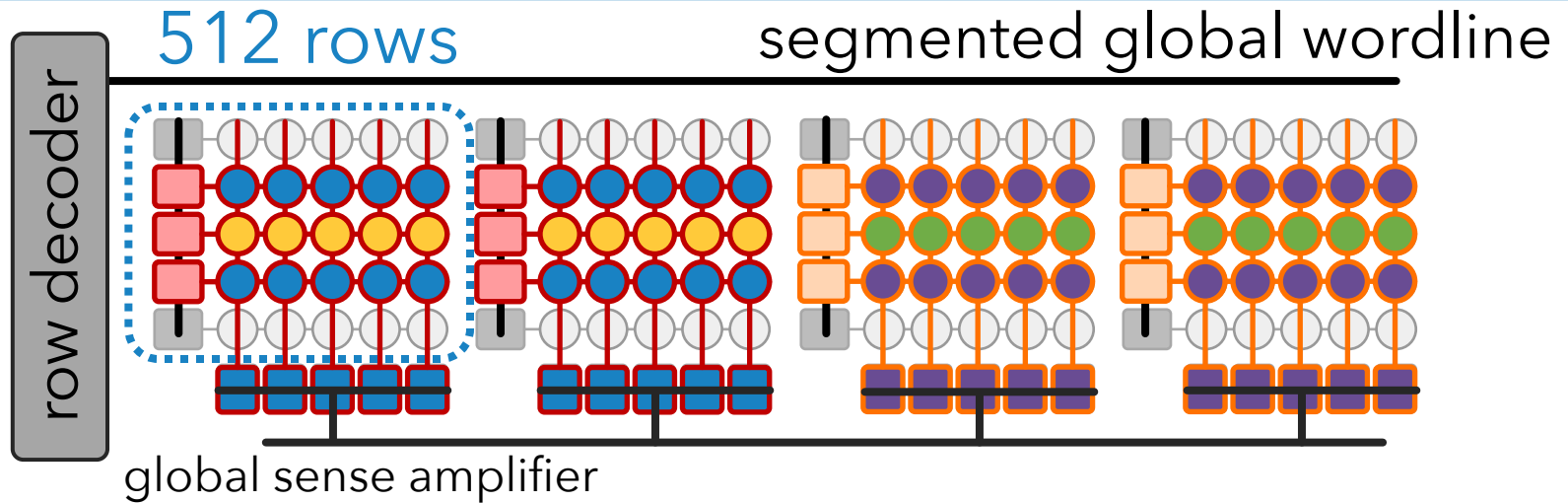
## Fine-grained DRAM for processing-using-DRAM:

**1** **Improves SIMD utilization**
- for a single PUD operation, only access the DRAM mats with target data
- for multiple PUD operations, execute independent operations concurrently
  → **multiple instruction, multiple data (MIMD) execution model**

**SAFARI**

segmented global wordline

row decoder

global sense amplifier

## Fine-grained DRAM for processing-using-DRAM:

**1** **Improves SIMD utilization**
- for a single PUD operation, only access the DRAM mats with target data
- for multiple PUD operations, execute independent operations concurrently
  → **multiple instruction, multiple data (MIMD) execution model**

**2** **Enables low-cost interconnects for vector reduction**
- global and local data buses can be used for inter-/intra-mat communication

**SAFARI**

512 rows — segmented global wordline

row decoder

global sense amplifier

## Fine-grained DRAM for processing-using-DRAM:

**1** **Improves SIMD utilization**
- for a single PUD operation, only access the DRAM mats with target data
- for multiple PUD operations, execute independent operations concurrently
  → **multiple instruction, multiple data (MIMD) execution model**

**2** **Enables low-cost interconnects for vector reduction**
- global and local data buses can be used for inter-/intra-mat communication

**3** **Eases programmability**
- SIMD parallelism in a DRAM mat is on par with vector ISAs' SIMD width

# MIMDRAM:
## Overview

**MIMDRAM is a hardware/software co-designed PUD system that enables fine-grained PUD computation at low cost and programming effort**
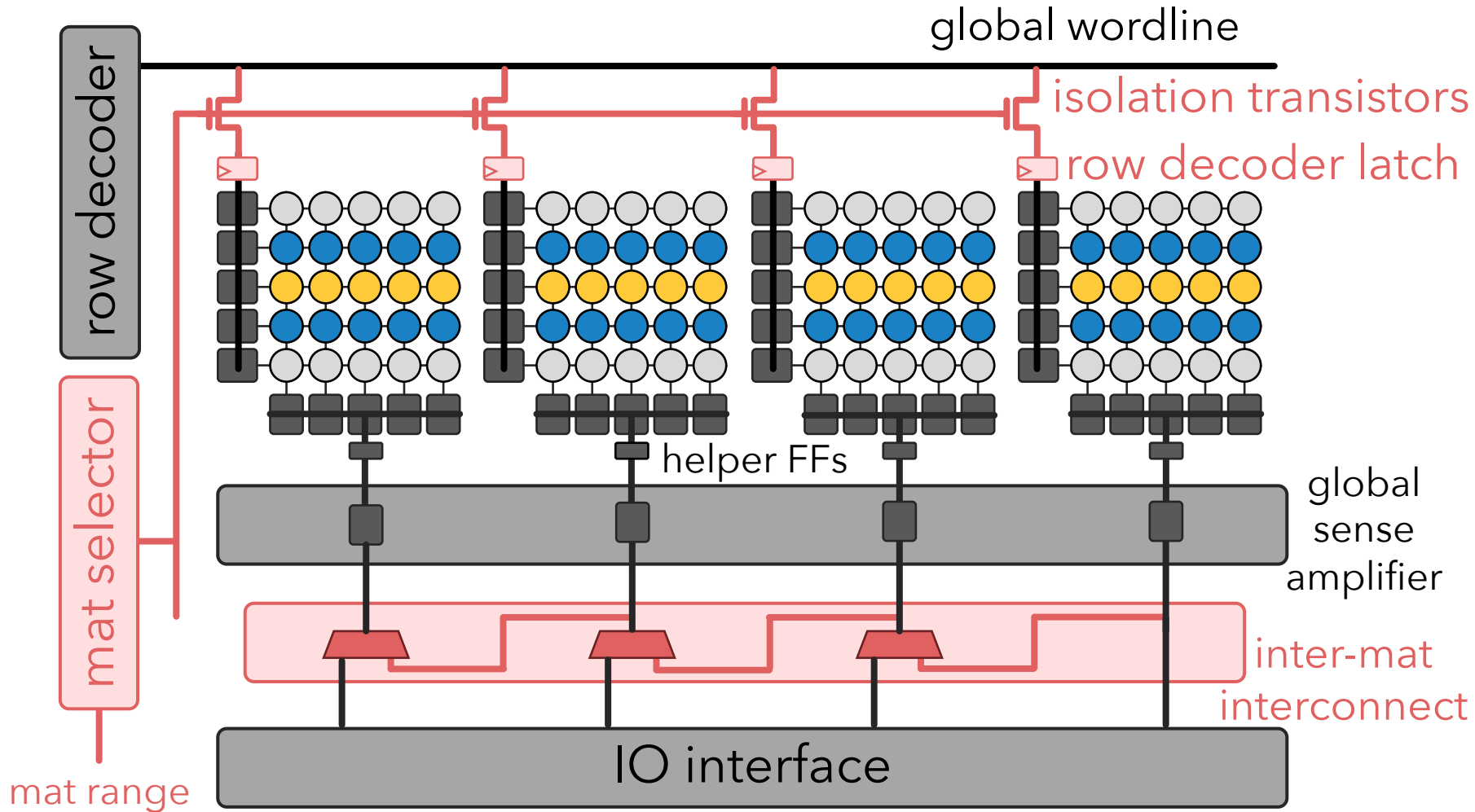
## Main components of MIMDRAM:

**1 Hardware**
- DRAM array modification to enable fine-grained PUD computation
- inter- and intra-mat interconnects to enable PUD vector reduction
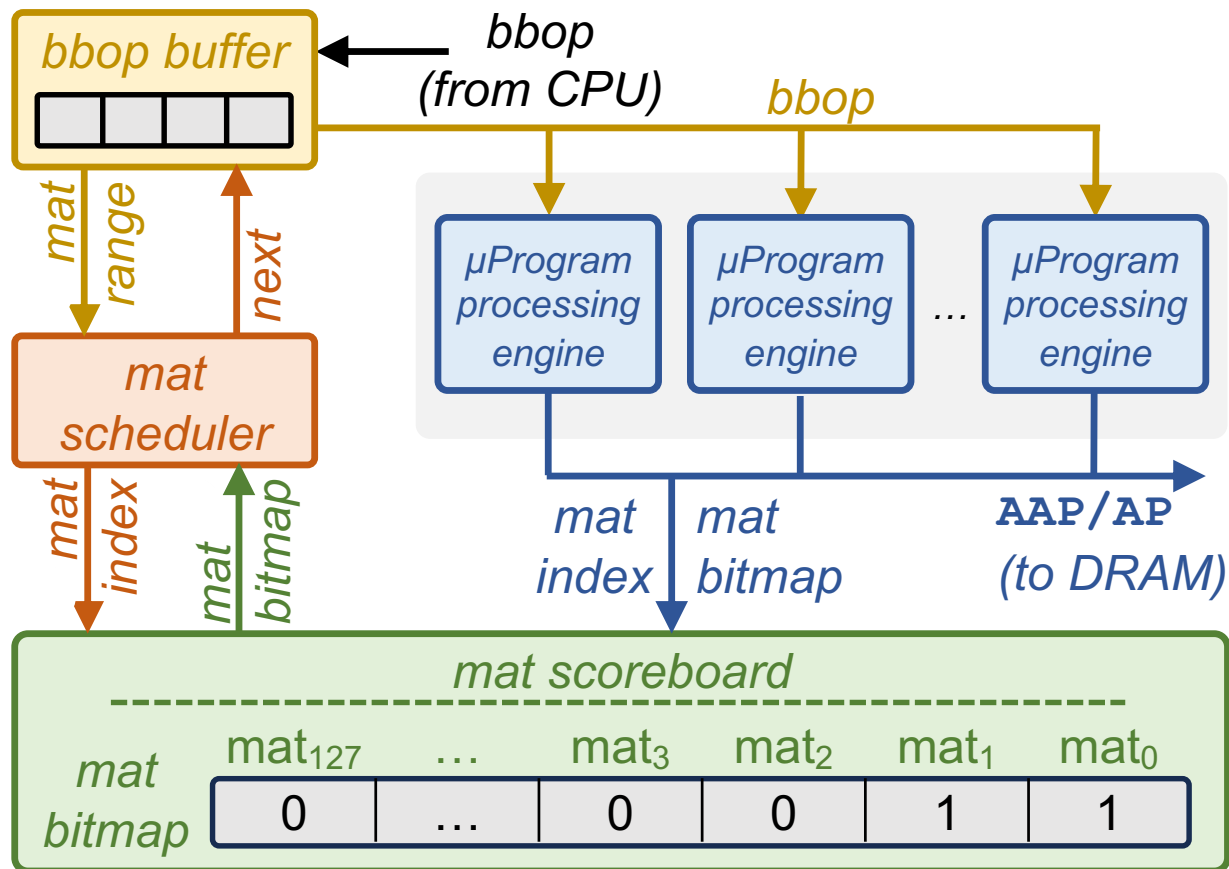- control unit design to orchestrate PUD execution

**2 Software**
- compiler support to transparently generate PUD instructions
- system support to map and execute PUD instructions

# MIMDRAM:
## Modifications to DRAM Chip

**SAFARI**

# MIMDRAM:
## Control Unit Design

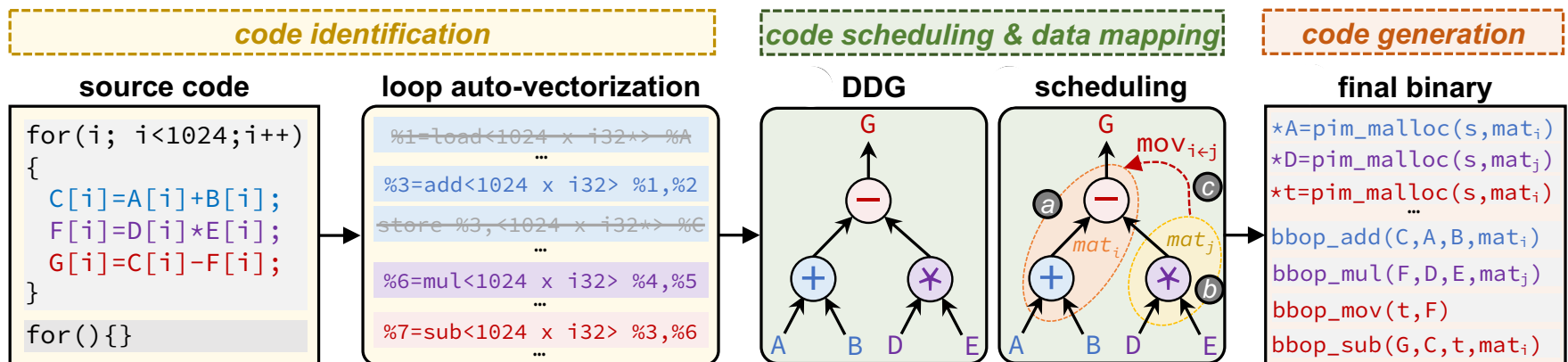**The control unit schedules and orchestrates the execution of multiple PUD operations transparently**

# MIMDRAM:
## Compiler Support

**Transparently:**
**extract SIMD parallelism from an application, and**
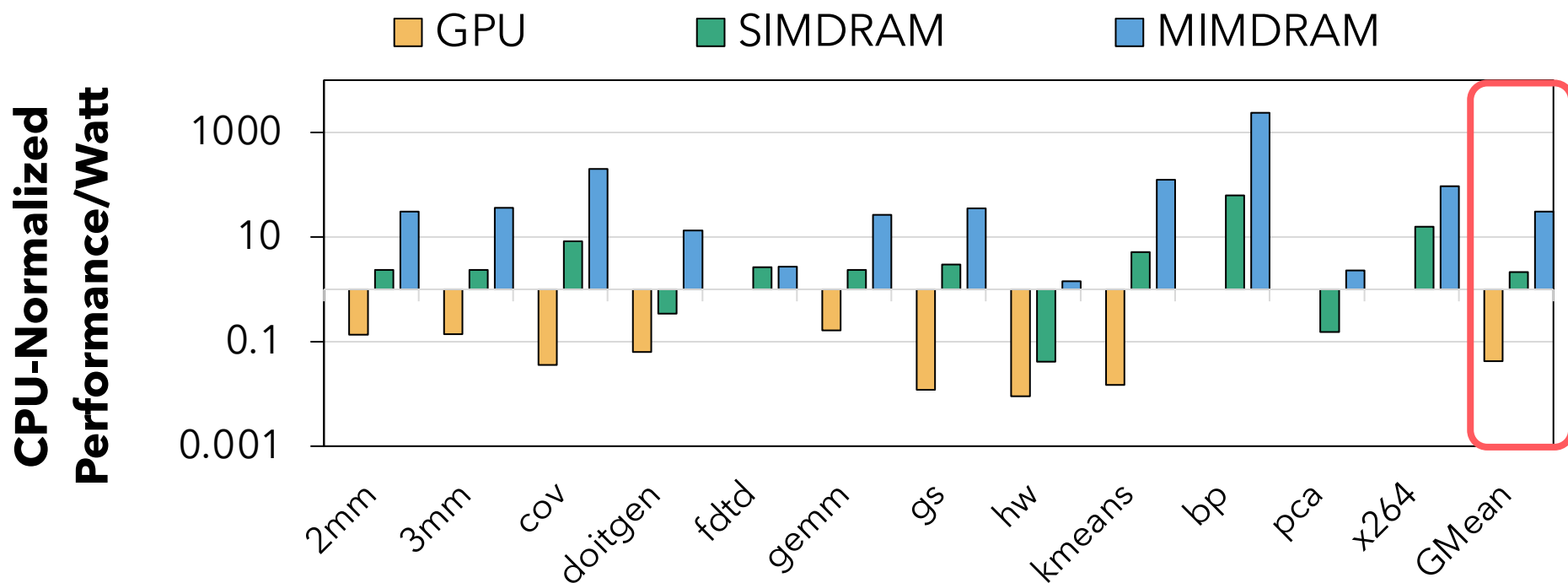**schedule PUD instructions while maximizing utilization**

## Three new LLVM-based passes targeting PUD execution

**CPU-Normalized Performance/Watt** across applications: 2mm, 3mm, cov, doitgen, fdtd, gemm, gs, hw, kmeans, bp, pca, x264, GMean. Legend: GPU, SIMDRAM, MIMDRAM.

**Takeaway**

MIMDRAM significantly improves energy efficiency compared to CPU (30.6x), GPU (6.8x), and SIMDRAM (14.3x)

**SAFARI**

63

# More on MIMDRAM

- Geraldo F. Oliveira, Ataberk Olgun, Abdullah Giray Yağlıkçı, F. Nisa Bostancı, Juan Gómez-Luna, Saugata Ghose, and Onur Mutlu
**" MIMDRAM: An End-to-End Processing-Using-DRAM System for High-Throughput, Energy-Efficient and Programmer-Transparent Multiple-Instruction Multiple-Data Computing"**
*Proceedings of the 30th International Symposium on High-Performance Computer Architecture* (**HPCA**), Edinburgh, Scotland, March 2024.

**MIMDRAM: An End-to-End Processing-Using-DRAM System for High-Throughput, Energy-Efficient and Programmer-Transparent Multiple-Instruction Multiple-Data Processing**

Geraldo F. Oliveira[†]   Ataberk Olgun[†]   Abdullah Giray Yağlıkçı[†]   F. Nisa Bostancı[†]
Juan Gómez-Luna[†]   Saugata Ghose[‡]   Onur Mutlu[†]

[†] *ETH Zürich*   [‡] *Univ. of Illinois Urbana-Champaign*

# In-DRAM Physical Unclonable Functions

- Jeremie S. Kim, Minesh Patel, Hasan Hassan, and Onur Mutlu,
  **"The DRAM Latency PUF: Quickly Evaluating Physical Unclonable Functions by Exploiting the Latency-Reliability Tradeoff in Modern DRAM Devices"**
  Proceedings of the *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-DRAM Lookup-Table Based Execution

João Dinis Ferreira, Gabriel Falcao, Juan Gómez-Luna, Mohammed Alser, Lois Orosa, Mohammad Sadrosadati, Jeremie S. Kim, Geraldo F. Oliveira, Taha Shahroodi, Anant Nori, and Onur Mutlu,
**"pLUTo: Enabling Massively Parallel Computation in DRAM via Lookup Tables"**
*Proceedings of the* 55th International Symposium on Microarchitecture (**MICRO**), Chicago, IL, USA, October 2022.
[Slides (pptx) (pdf)]
[Longer Lecture Slides (pptx) (pdf)]
[Lecture Video (26 minutes)]
[arXiv version]
[Source Code (Officially Artifact Evaluated with All Badges)]
***Officially artifact evaluated as available, reusable and reproducible.***



## pLUTo: Enabling Massively Parallel Computation in DRAM via Lookup Tables

João Dinis Ferreira[§]   Gabriel Falcao[†]   Juan Gómez-Luna[§]   Mohammed Alser[§]
Lois Orosa[§▽]   Mohammad Sadrosadati[§]   Jeremie S. Kim[§]   Geraldo F. Oliveira[§]
Taha Shahroodi[‡]   Anant Nori[*]   Onur Mutlu[§]

[§]*ETH Zürich*   [†]*IT, University of Coimbra*   [▽]*Galicia Supercomputing Center*   [‡]*TU Delft*   [*]*Intel*

# Limitations of Processing-using-DRAM

| Data Movement | *RowClone, Seshadri+ 2013* <br> *LISA, Chang+ 2013* |
|---|---|
| Bitwise Operations | *Ambit, Seshadri+ 2017* |
| Bit Shifting | *DRISA, Li+ 2017* |
| Arithmetic Operations | *SIMDRAM, Hajinazar & Oliveira+ 2021* |

**Existing Processing-using-DRAM architectures only support a limited range of operations**
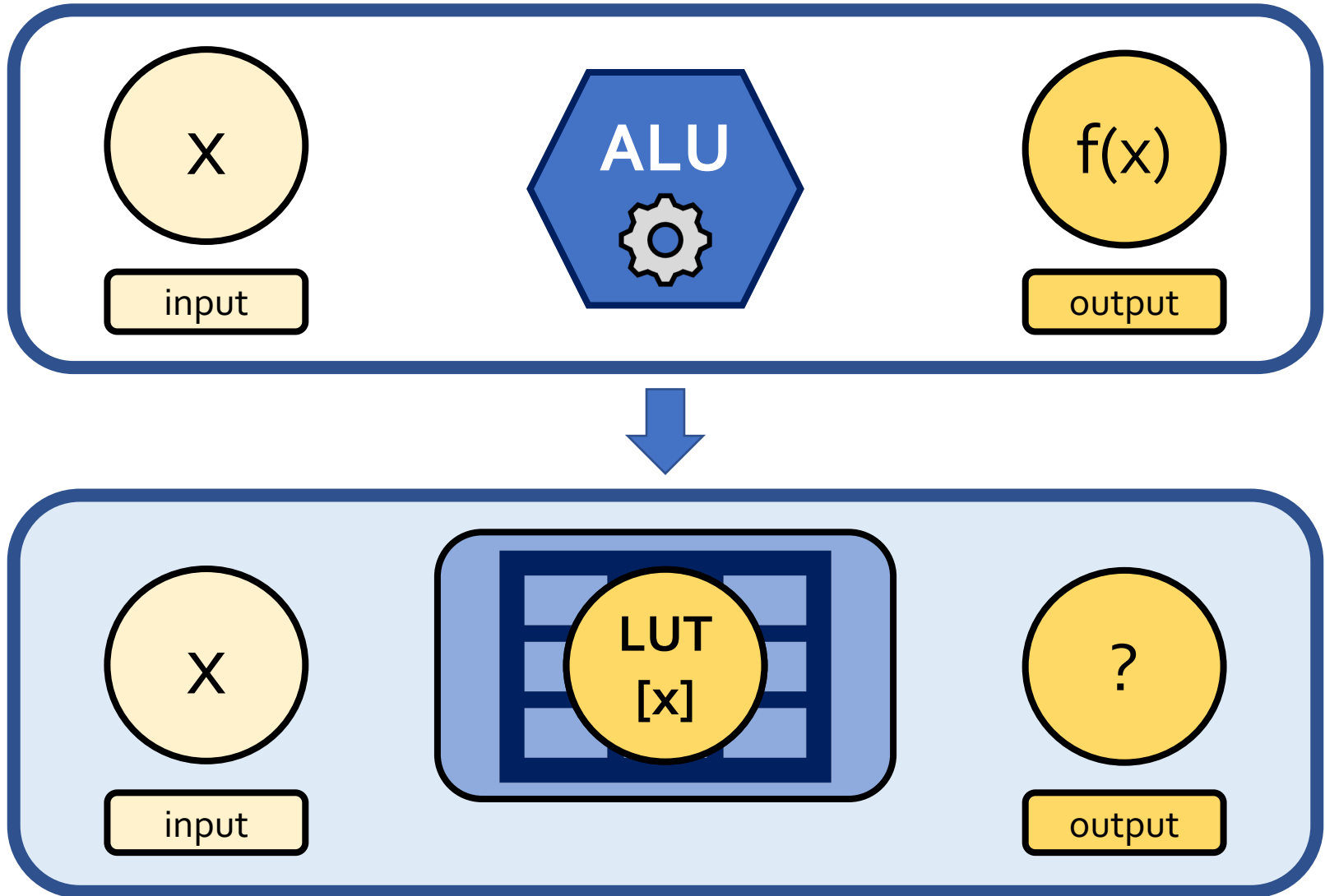
# The Goal of pLUTo

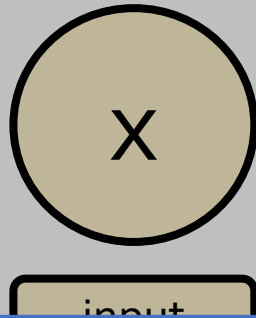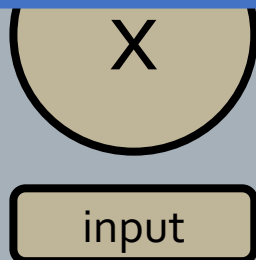*Extend* **Processing-using-DRAM to support the execution of** *arbitrarily complex operations*

**SAFARI**

# pLUTo: Key Idea

X
input

ALU

f(x)
output

# pLUTo: Key Idea

**SAFARI**

# pLUTo: Key Idea



**Replace computation with memory accesses**
**→ *pLUTo LUT Query* operation**

X

ALU

f(x)

input

output

X

LUT
[x]

input

output

# System Integration



| pLUTo Compiler | | pLUTo Controller |

C-Like Code with **pLUTo API** calls → Assembly Code with **pLUTo ISA Extensions** → Execution in the **DRAM Substrate**

api_pluto_mul

pluto_subarray_alloc
pluto_bit_shift_l
pluto_or
pluto_op

ACT
PRE
ACT
ACT
PRE
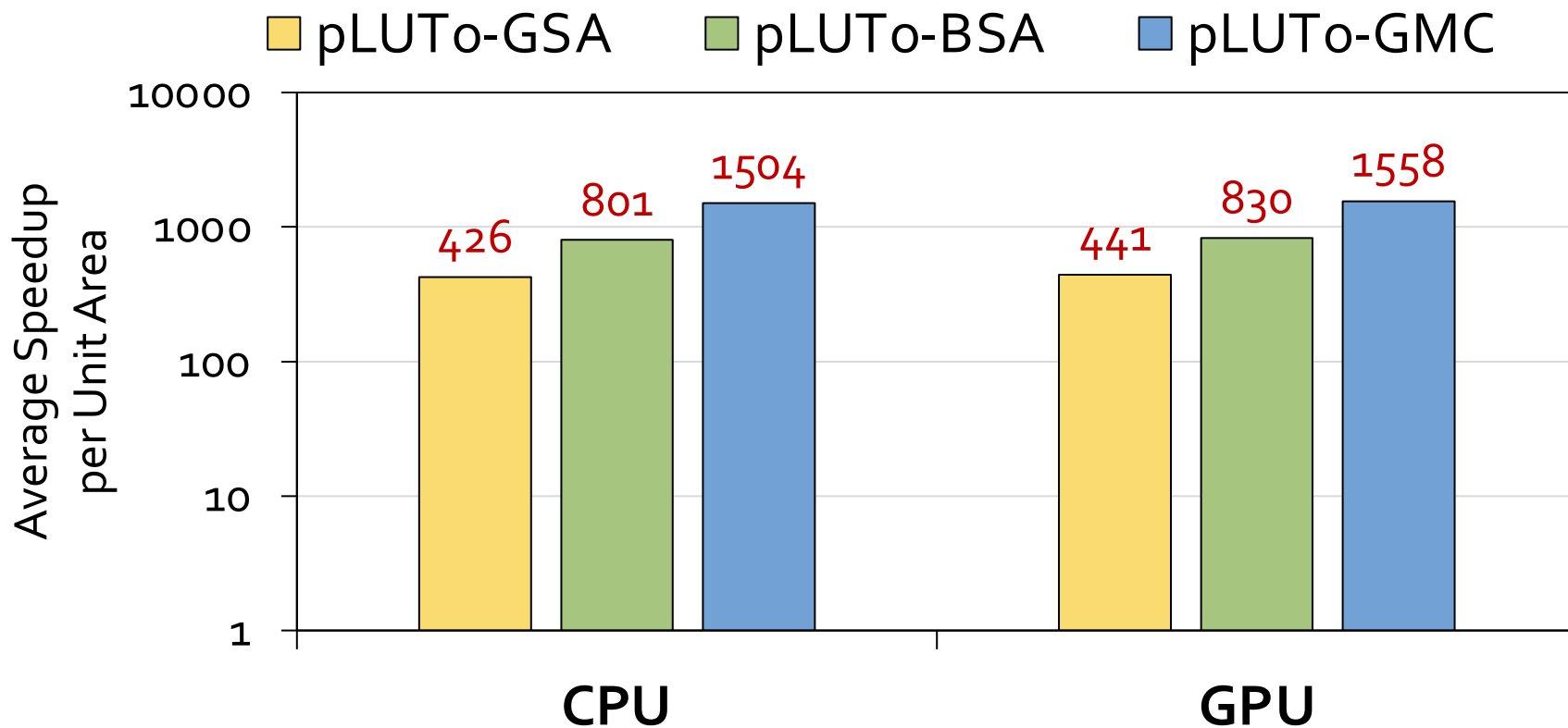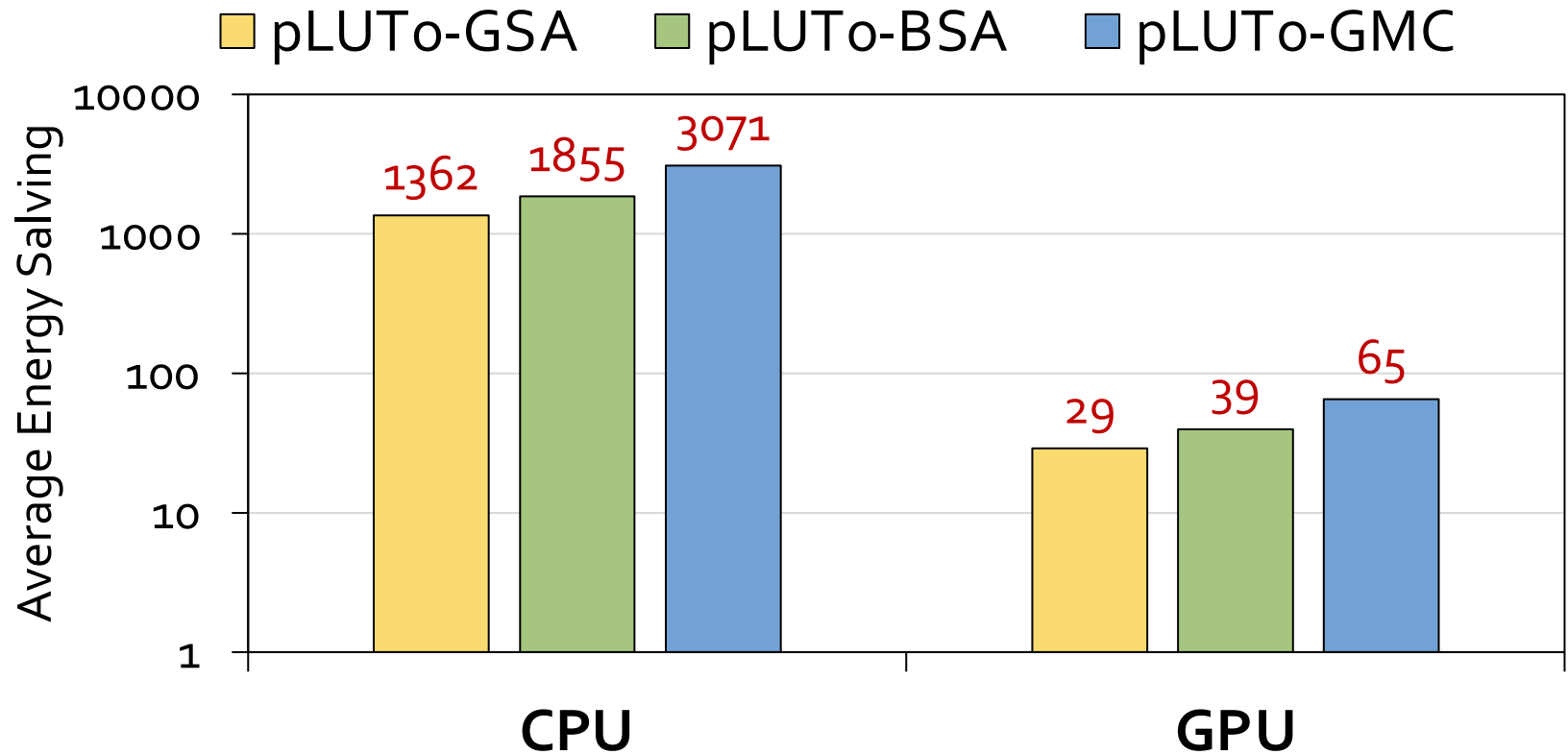...

# Performance (normalized to area)

Average speedup normalized to area across 7 real-world workloads



pLUTo provides *substantially higher* performance per unit area than *both* the CPU and the GPU

# Energy Consumption

Average energy consumption across 7 real-world workloads



pLUTo *significantly reduces energy consumption* compared to processor-centric architectures for various workloads

# More Results in the Paper

- **Comparison with FPGA**
- **Area Overhead Analysis**
- **Circuit-Level Reliability & Correctness**

- **Subarray-Level Parallelism**
- **LUT Loading Overhead**
- **Range of Supported Operations**



## pLUTo: Enabling Massively Parallel Computation in DRAM via Lookup Tables

João Dinis Ferreira[§]    Gabriel Falcao[†]    Juan Gómez-Luna[§]    Mohammed Alser[§]
Lois Orosa[§▽]    Mohammad Sadrosadati[§]    Jeremie S. Kim[§]    Geraldo F. Oliveira[§]
Taha Shahroodi[‡]    Anant Nori[⋆]    Onur Mutlu[§]

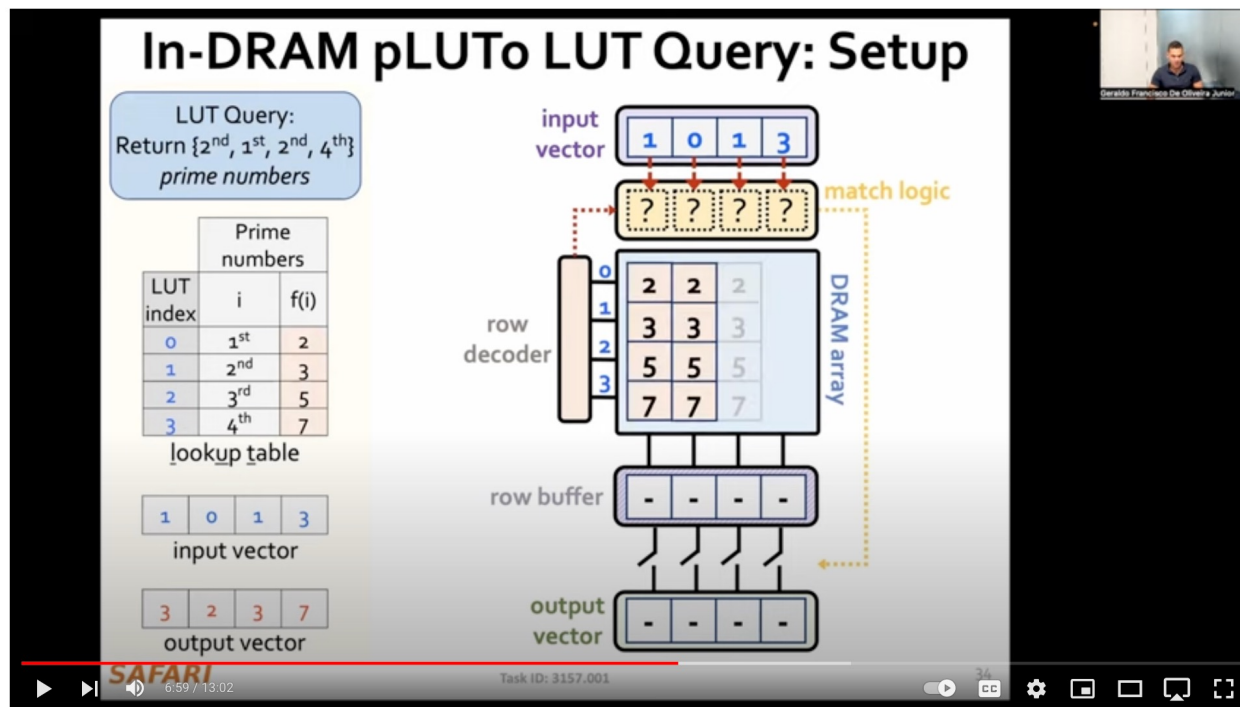[§]*ETH Zürich*    [†]*IT, University of Coimbra*    [▽]*Galicia Supercomputing Center*    [‡]*TU Delft*    [⋆]*Intel*

# SRC TECHCON Presentation

- **Geraldo F. Oliveira**
  - ❑ pLUTo: Enabling Massively Parallel Computation in DRAM via Lookup Tables
  - ❑ https://arxiv.org/pdf/2104.07699.pdf



pLUTo: Enabling Massively Parallel Computation in DRAM via Lookup Tables, SRC TECHCON 2023

**Onur Mutlu Lectures**
35.5K subscribers

321 views 9 days ago
pLUTo: Enabling Massively Parallel Computation in DRAM via Lookup Tables
Speaker: Geraldo F. Oliveira ...more

SAFARI

# Bulk Bitwise Operations in Real DRAM Chips

- Ismail Emir Yüksel, Yahya Can Tugrul Ataberk Olgun, F. Nisa Bostancı, A. Giray Yaglıkçı, Geraldo F. Oliveira, Haocong Luo, Juan Gómez-Luna, Mohammad Sadrosadati, Onur Mutlu, "**Functionally-Complete Boolean Logic in Real DRAM Chips: Experimental Characterization and Analysis**," *Proceedings of the 30th International Symposium on High-Performance Computer Architecture* (**HPCA**), Edinburgh, Scotland, March 2024.

## 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

# The Capability of COTS DRAM Chips

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

**1** Can **simultaneously activate** up to **48 rows** in **two neighboring subarrays**

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

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

# DRAM Testing Infrastructure

- Developed from DRAM Bender [Olgun+, TCAD'23]*

- Fine-grained control over DRAM commands, timings, and temperature

*Olgun et al., "DRAM Bender: An Extensible and Versatile FPGA-based Infrastructure to Easily Test State-of-the-art DRAM Chips," TCAD, 2023.

# 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 |

# The Capability of COTS DRAM Chips

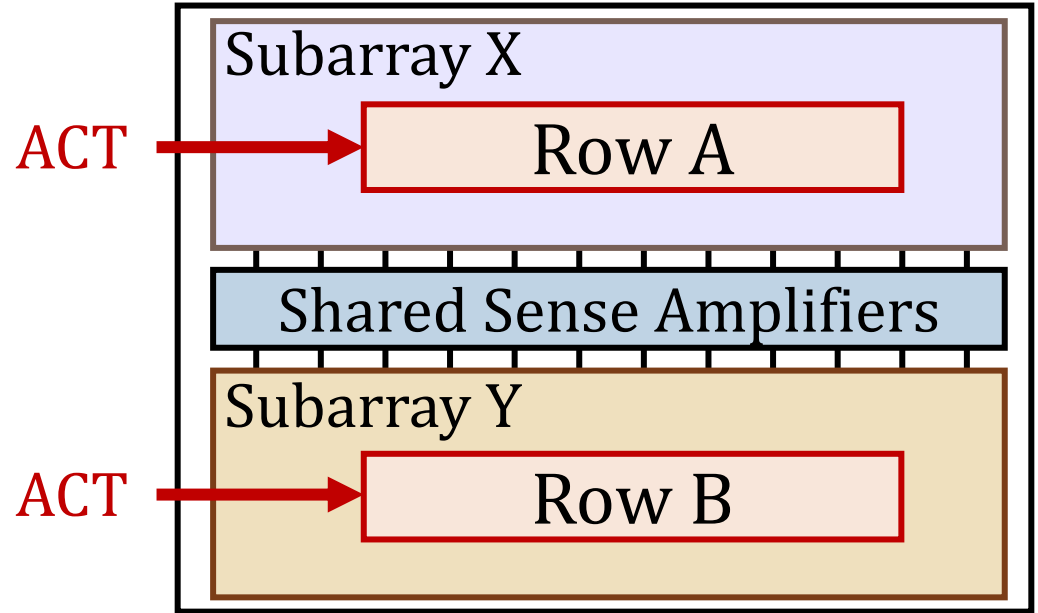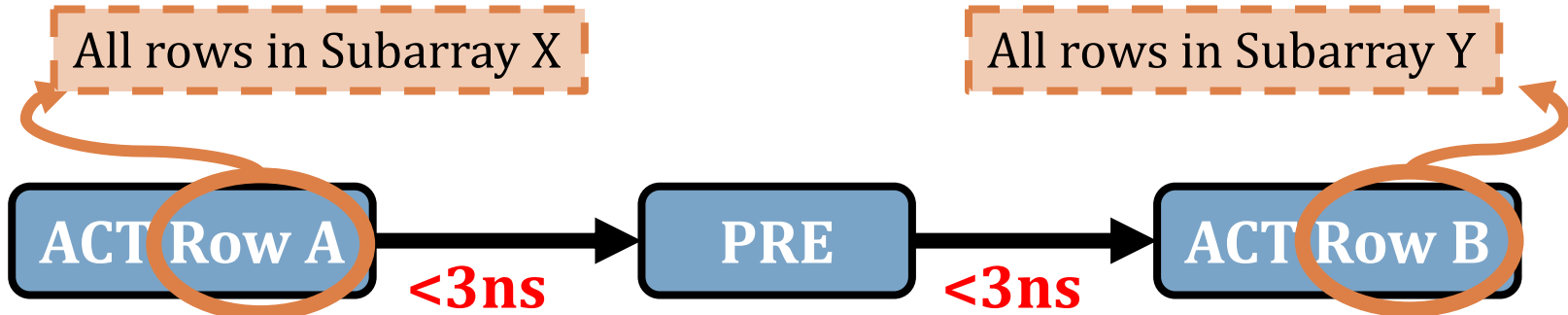We **demonstrate** that **COTS DRAM chips:**

**1** Can **simultaneously activate** up to **48 rows** in **two neighboring subarrays**

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

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

# Characterization Methodology

- To understand **which and how many** rows are simultaneously activated
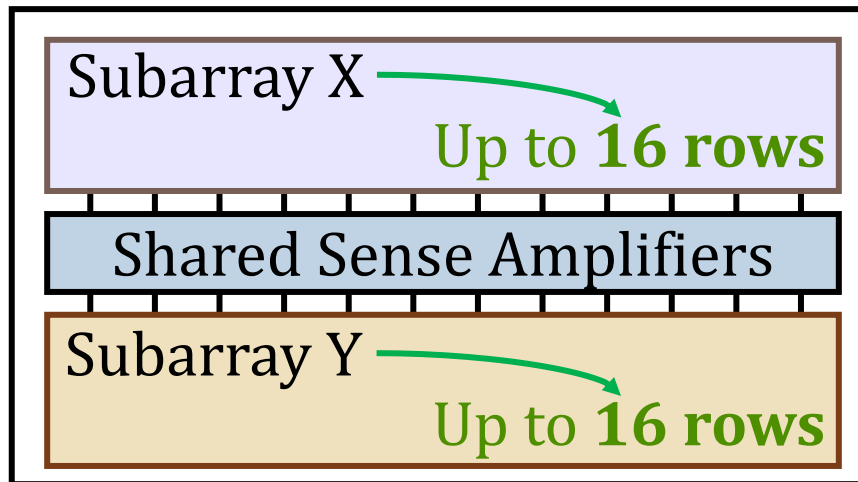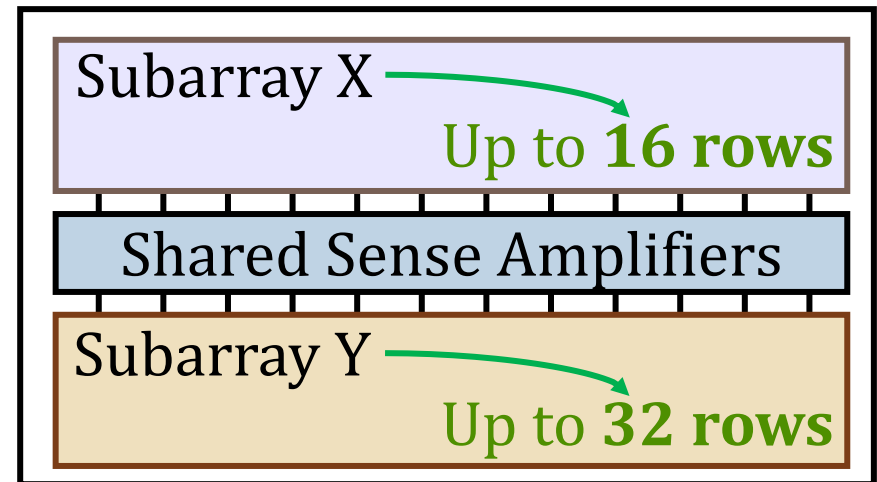  - **Sweep** Row A and Row B addresses

# Key Results

COTS DRAM chips have **two distinct** sets of activation patterns in **neighboring subarrays** when two rows are activated with **violated timings**

| **Exactly the same number** of rows in each subarray are activated | **Twice as many** rows in one subarray **compared to its neighbor subarray** are activated |
|---|---|

Subarray X — Up to **16 rows**

Shared Sense Amplifiers

Subarray Y — Up to **16 rows**

A total of **32 rows**

Subarray X — Up to **16 rows**

Shared Sense Amplifiers

Subarray Y — Up to **32 rows**

A total of **48 rows**

# The Capability of COTS DRAM Chips

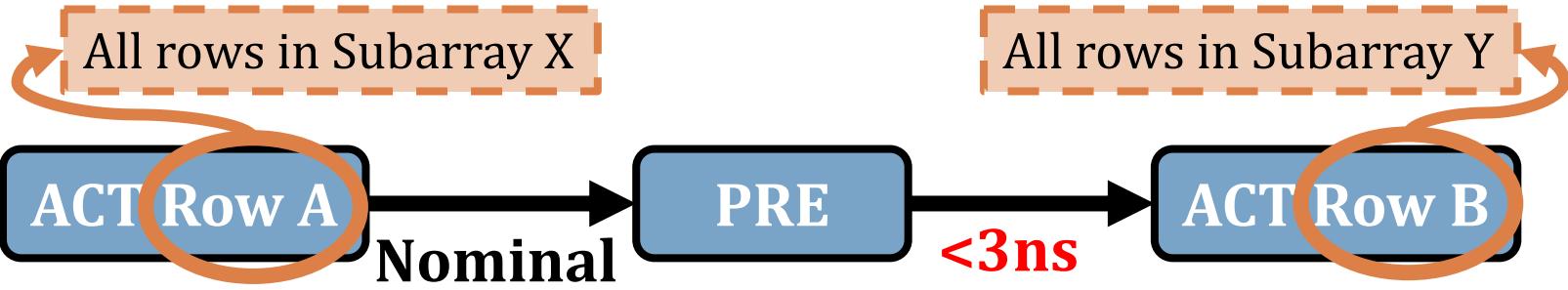**1** Can **simultaneously activate** up to **48 rows** in **two neighboring subarrays**

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

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

# Characterization Methodology

- Sweep **Row A and Row B addresses**

All rows in Subarray X

All rows in Subarray Y

**ACT Row A** → **Nominal** → **PRE** → **<3ns** → **ACT Row B**

- Sweep **DRAM chip temperature**

95°C

50°C

Temperature

# Key Takeaways from In-DRAM NOT Operation

| Key Takeaway 1 |
|:---:|
| **COTS DRAM chips can perform NOT operations with up to 32 destination rows** |

| Key Takeaway 2 |
|:---:|
| **Temperature has a small effect on the reliability of NOT operations** |

# The Capability of COTS DRAM Chips

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

**1** Can **simultaneously activate** up to 48 rows in **two neighboring subarrays**
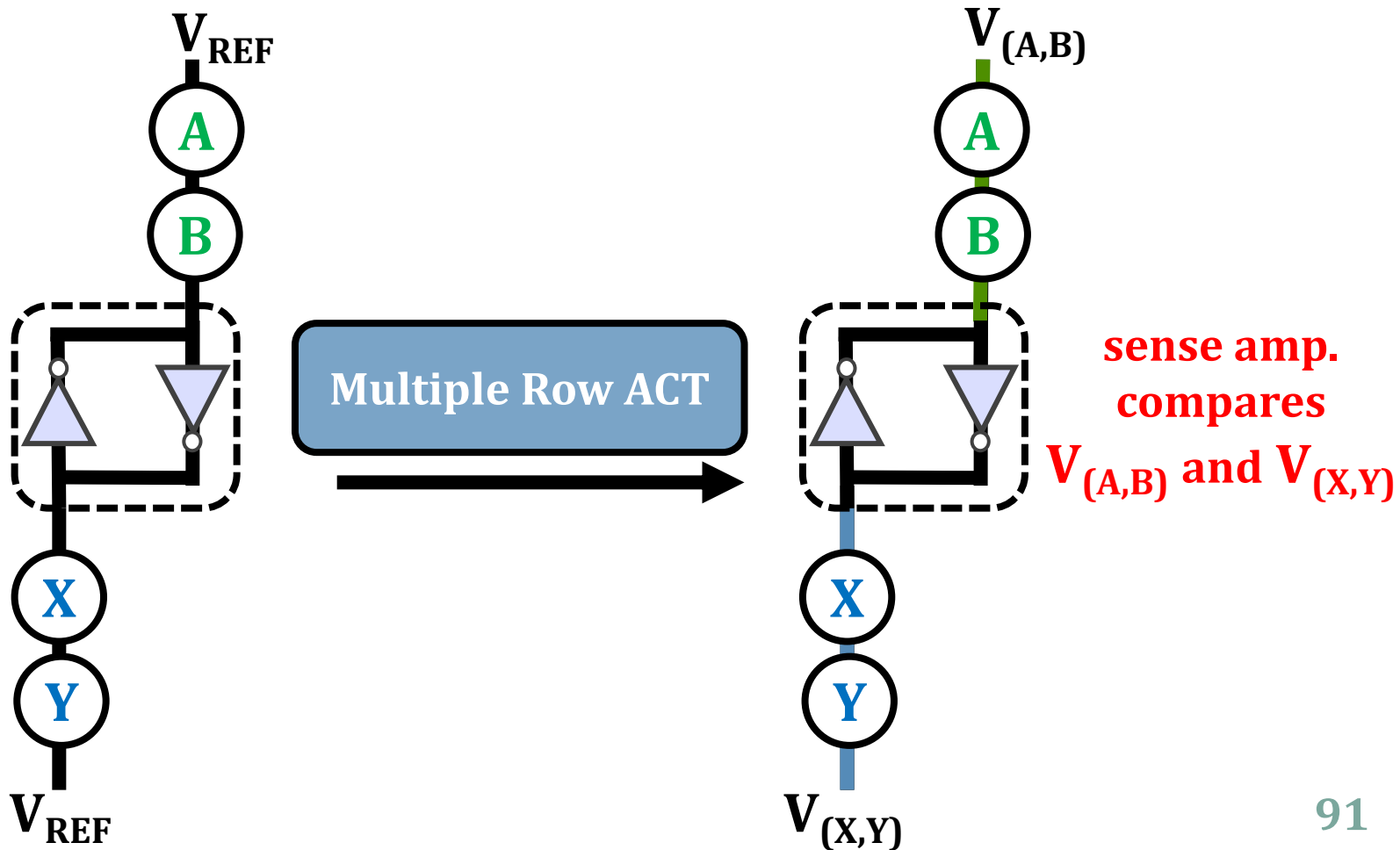
**2** Can perform **NOT operation** with **up to 32** output operands
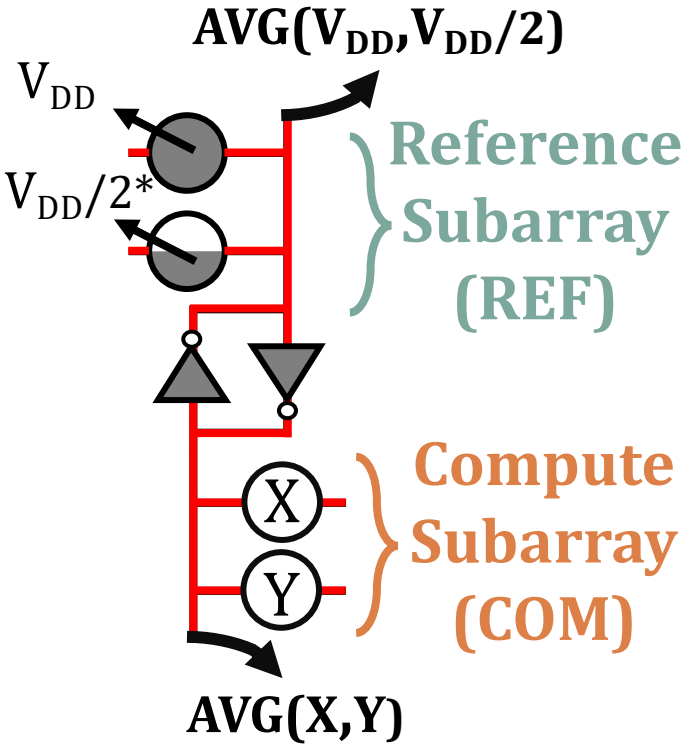
**3** Can perform **up to 16-input AND, NAND, OR, and NOR** operations

# Key Idea

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



$V_{REF}$

A

B

$V_{(A,B)}$

A

B

Multiple Row ACT

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

X

Y

X

Y

$V_{REF}$

$V_{(X,Y)}$

SAFARI

# Two-Input AND and NAND Operations



$V_{DD}=1$ & GND = 0

| X | Y | COM | REF |
|---|---|-----|-----|
| 0 | 0 | 0 | 1 |
| 0 | 1 | 0 | 1 |
| 1 | 0 | 0 | 1 |
| 1 | 1 | 1 | 0 |
| | | AND | NAND |

*Gao et al., "FracDRAM: Fractional Values in Off-the-Shelf DRAM," in MICRO, 2022.

# Key Takeaways from In-DRAM Operations

## Key Takeaway 1

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

## Key Takeaway 2

**COTS DRAM chips can perform
AND, NAND, OR, and NOR operations
with very high reliability**

## Key Takeaway 3

**Data pattern slightly affects
the reliability of AND, NAND, OR, and NOR operations**

# 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 <u>P</u>rocessing-<u>i</u>n-<u>DRAM</u>

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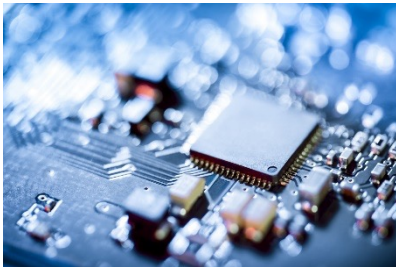**https://www.youtube.com/watch?v=qeukNs5XI3g&t=4192s**

# PiDRAM

**Goal:** Develop a **flexible** platform to explore **end-to-end** implementations of PuM techniques

- Enable rapid integration via key components

## Hardware



**1** Easy-to-extend Memory Controller
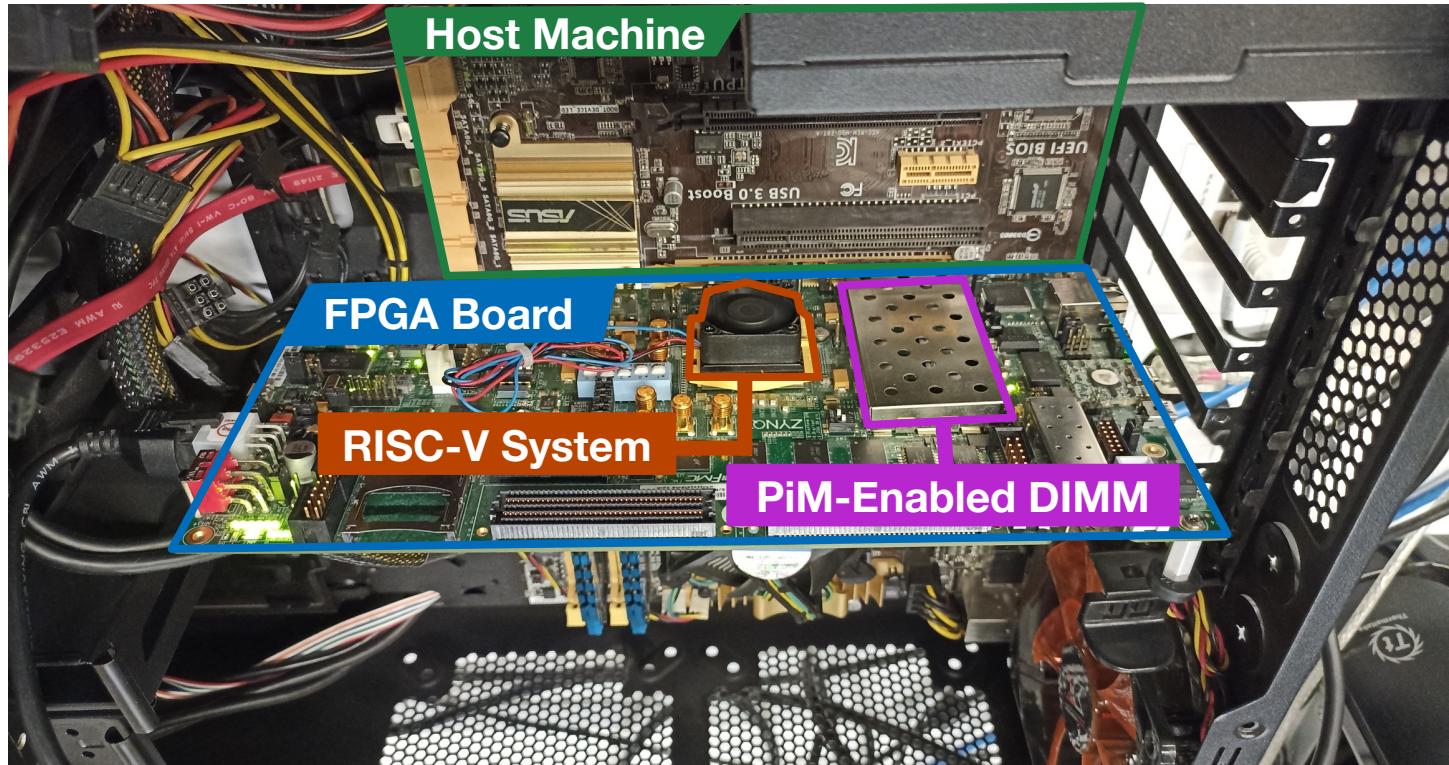
**2** ISA-transparent PuM Controller

## Software



**1** Extensible Software Library

**2** Custom Supervisor Software

# Real Processing Using Memory Prototype



**https://arxiv.org/pdf/2111.00082.pdf**

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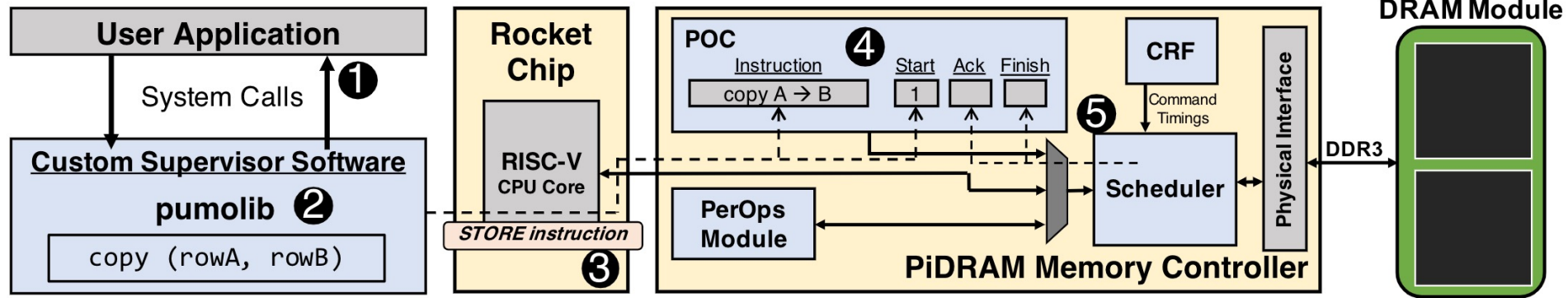
**https://www.youtube.com/watch?v=qeukNs5XI3g&t=4192s**

# PiDRAM Workflow



1- User application interfaces with the OS via system calls

2- OS uses PuM Operations Library (pumolib) to convey operation related information to the hardware *using*

    3- STORE instructions that target the memory mapped registers of the PuM Operations Controller (POC)

4- POC oversees the execution of a PuM operation (e.g., RowClone, bulk bitwise operations)

5- Scheduler arbitrates between regular (load, store) and PuM operations and issues DRAM commands with custom timings

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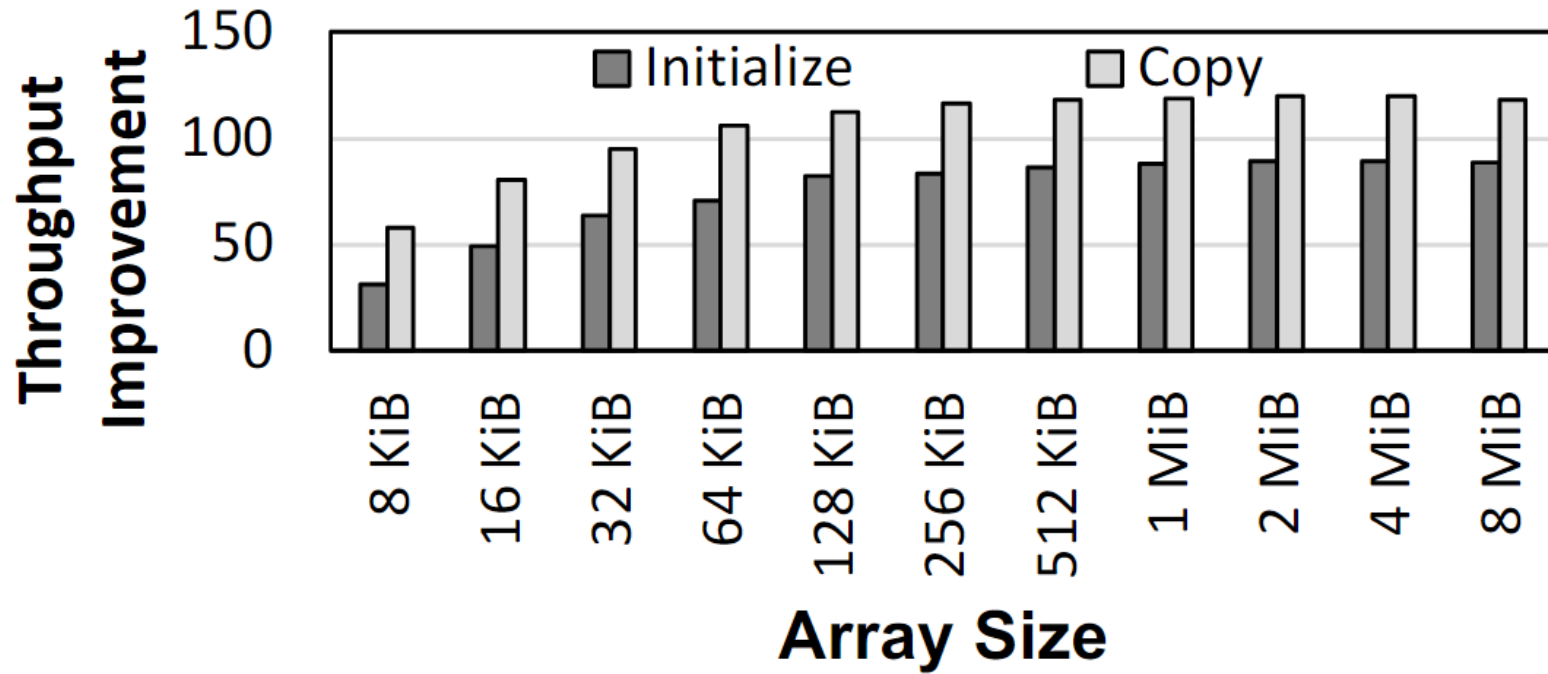
1- Open IP Catalog
2- Find "Memory Interface Generator (MIG 7 Series)" IP and double click

**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**

# PiDRAM is Open Source

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



□ CMU-SAFARI / **PiDRAM** ( Public )

⚡ Edit Pins ▾    👁 Watch 3 ▾    ⑂ Fork 2    ☆ Star 21    ▾

<> Code   ⊙ Issues   ⑃ Pull requests   ⊙ Actions   ⊞ Projects   📖 Wiki   ⊘ Security   📈 Insights   ⚙ Settings

⑂ master ▾    ⑂ 2 branches   ⬡ 0 tags    Go to file    Add file ▾    Code ▾

**About**

PiDRAM is the first flexible end-to-end framework that enables system integration studies and evaluation of real Processing-using-Memory techniques. Prototype on a RISC-V rocket chip system implemented on an FPGA. Described in our preprint:

https://arxiv.org/abs/2111.00082

📖 Readme

☆ 21 stars

👁 3 watching

⑂ 2 forks

🐙 olgunataberk Fix small mistake in README    46522cc  on Dec 5, 2021   ⏱ 11 commits

📁 controller-hardware    Add files via upload    7 months ago
📁 fpga-zynq    Adds instructions to reproduce two key results    7 months ago
📄 README.md    Fix small mistake in README    7 months ago

☰ README.md    ✎

## PiDRAM

PiDRAM is the first flexible end-to-end framework that enables system integration studies and evaluation of real Processing-using-Memory (PuM) techniques. PiDRAM, at a high level, comprises a RISC-V system and a custom memory controller that can perform PuM operations in real DDR3 chips. This repository contains all sources required to build PiDRAM and develop its prototype on the Xilinx ZC706 FPGA boards.

**Releases**

No releases published
Create a new release

**SAFARI**   ⬭ **kasırga**

# Extended Version on ArXiv

## https://arxiv.org/abs/2111.00082

---

arXiv > cs > arXiv:2111.00082

Search... | All fields | Search

**Computer Science > Hardware Architecture**

# 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

Processing-using-memory (PuM) techniques leverage the analog operation of memory cells to perform computation. Several recent works have demonstrated PuM techniques in off-the-shelf DRAM devices. Since DRAM is the dominant memory technology as main memory in current computing systems, these PuM techniques represent an opportunity for alleviating the data movement bottleneck at very low cost. However, system integration of PuM techniques imposes non-trivial challenges that are yet to be solved. Design space exploration of potential solutions to the PuM integration challenges requires appropriate tools to develop necessary hardware and software components. Unfortunately, current specialized DRAM-testing platforms, or system simulators do not provide the flexibility and/or the holistic system view that is necessary to deal with PuM integration challenges.

We design and develop PiDRAM, the first flexible end-to-end framework that enables system integration studies and evaluation of real PuM techniques. PiDRAM provides software and hardware components to rapidly integrate PuM techniques across the whole system software and hardware stack (e.g., necessary modifications in the operating system, memory controller). We implement PiDRAM on an FPGA-based platform along with an open-source RISC-V system. Using PiDRAM, we implement and evaluate two state-of-the-art PuM techniques: in-DRAM (i) copy and initialization, (ii) true random number generation. Our results show that the in-memory copy and initialization techniques can improve the performance of bulk copy operations by 12.6x and bulk initialization operations by 14.6x on a real system. Implementing the true random number generator requires only 190 lines of Verilog and 74 lines of C code using PiDRAM's software and hardware components.

Comments: 15 pages, 12 figures
Subjects: **Hardware Architecture (cs.AR)**
Cite as: arXiv:2111.00082 **[cs.AR]**
(or arXiv:2111.00082v3 **[cs.AR]** for this version)
https://doi.org/10.48550/arXiv.2111.00082

**Download:**
- PDF
- Other formats

(cc) BY

Current browse context:
**cs.AR**
< prev | next >
new | recent | 2111
Change to browse by:
cs

**References & Citations**
- NASA ADS
- Google Scholar
- Semantic Scholar

**DBLP** - CS Bibliography
listing | bibtex

Juan Gómez-Luna
Behzad Salami
Hasan Hassan
Oguz Ergin
Onur Mutlu

**Export Bibtex Citation**

**Bookmark**

# Long Talk + Tutorial on Youtube

## https://youtu.be/s_z_S6FYpC8



Processing in Memory Course: Meeting 6: End-to-end Framework for Processing-using-Memory - Fall'21

615 views • Streamed live on 9 Nov 2021 • Project & Seminar, ETH Zürich, Fall 2021 Show more

# 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)]

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

# Pinatubo: RowClone and Bitwise Ops in PCM

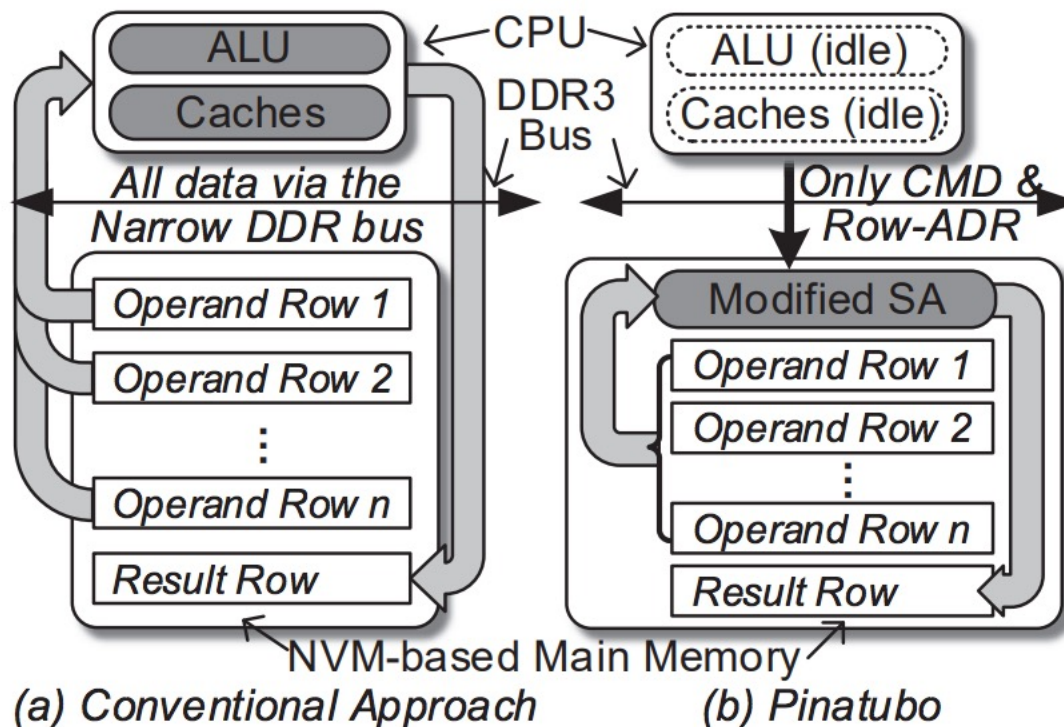## Pinatubo: A Processing-in-Memory Architecture for Bulk Bitwise Operations in Emerging Non-volatile Memories

Shuangchen Li[1]*, Cong Xu[2], Qiaosha Zou[1,5], Jishen Zhao[3], Yu Lu[4], and Yuan Xie[1]

University of California, Santa Barbara[1], Hewlett Packard Labs[2]
University of California, Santa Cruz[3], Qualcomm Inc.[4], Huawei Technologies Inc.[5]
{shuangchenli, yuanxie}ece.ucsb.edu[1]

# Pinatubo: RowClone and Bitwise Ops in PCM



Figure 2: Overview: (a) Computing-centric approach, moving tons of data to CPU and write back. (b) The proposed Pinatubo architecture, performs $n$-row bitwise operations inside NVM in one step.

# 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*

# Aside: In-Memory Crossbar Computation



(a) Multiply-Accumulate operation

(b) Vector-Matrix Multiplier
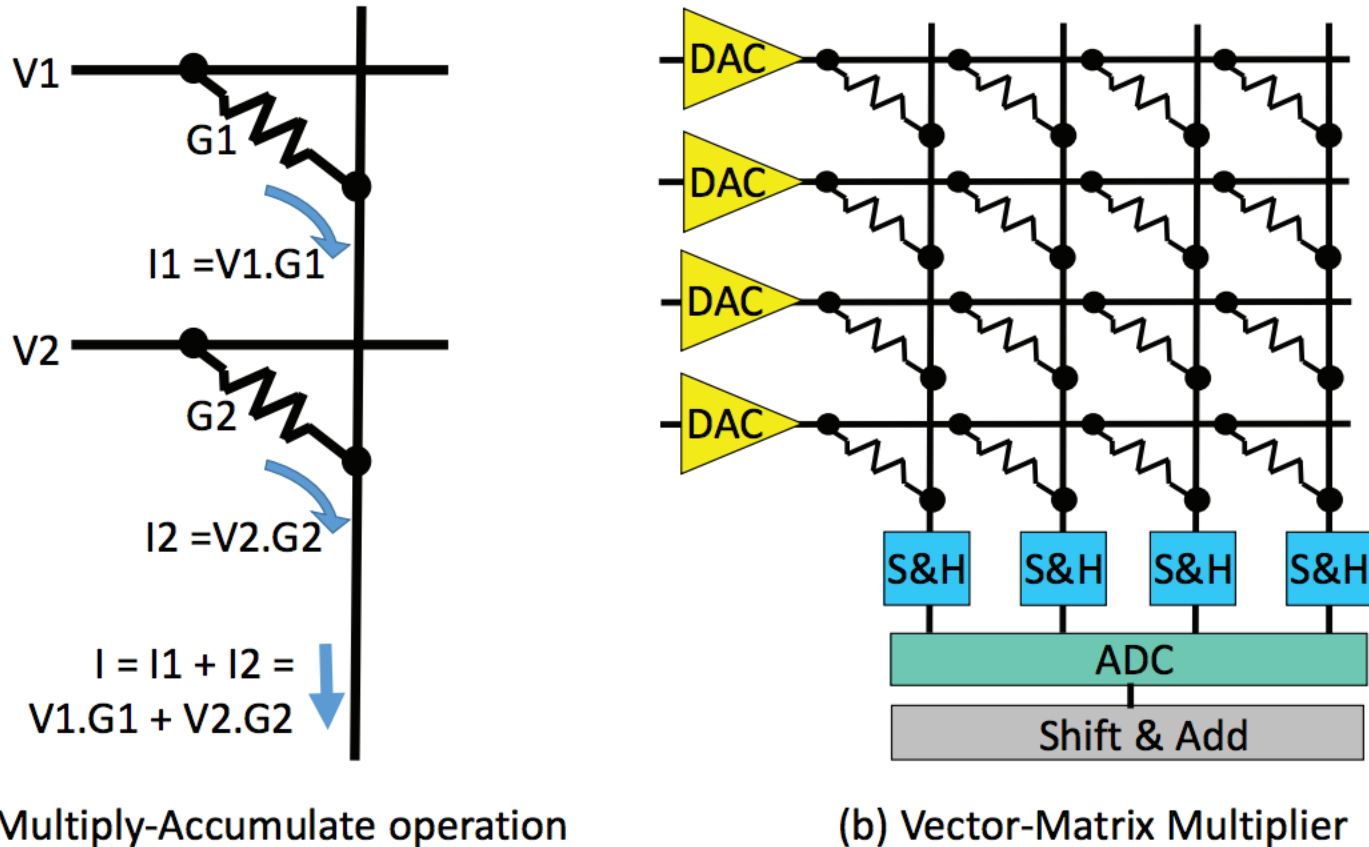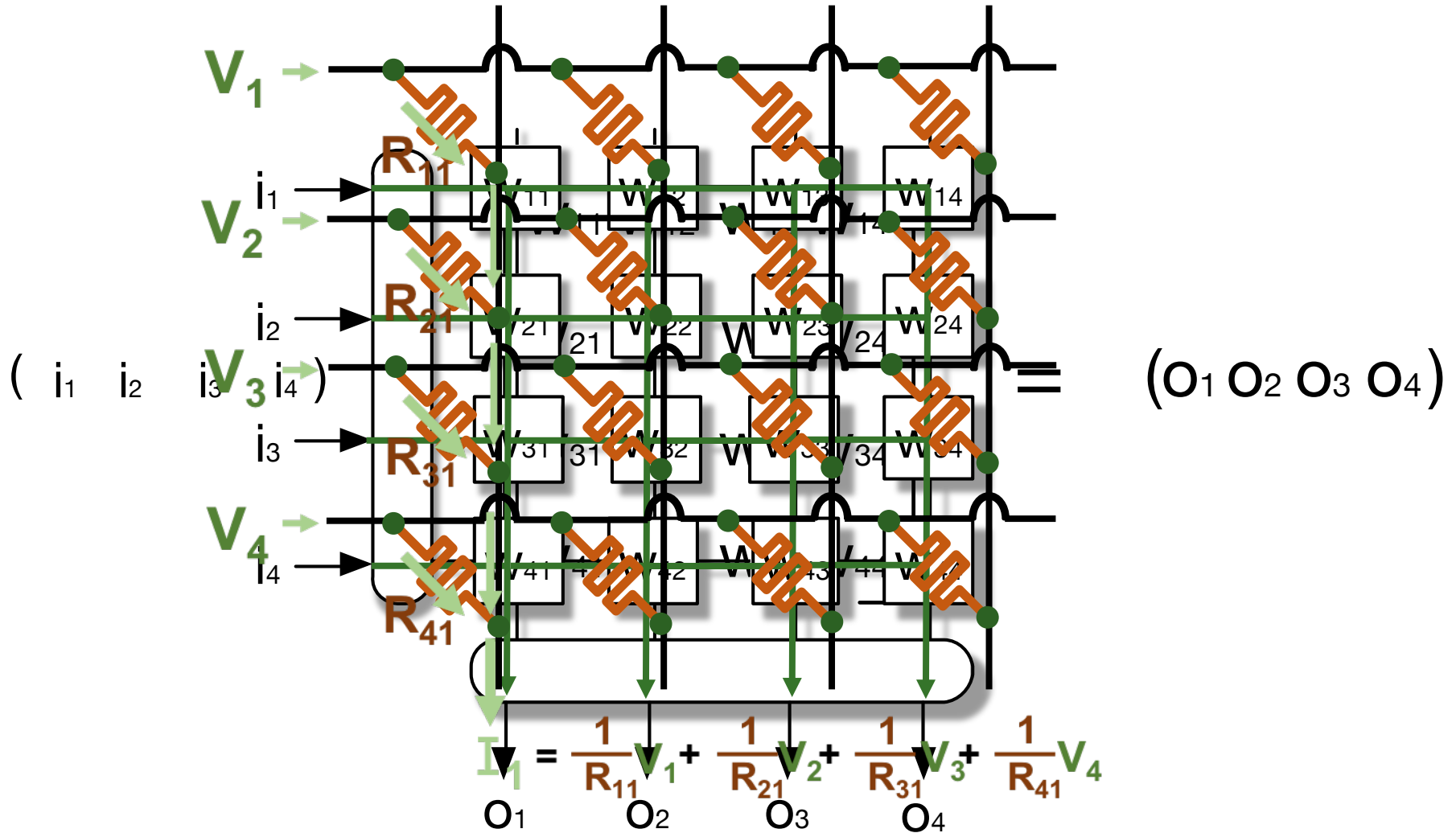
Fig. 1. (a) Using a bitline to perform an analog sum of products operation. (b) A memristor crossbar used as a vector-matrix multiplier.

Shafiee+, "ISAAC: A Convolutional Neural Network Accelerator with In-Situ Analog Arithmetic in Crossbars", ISCA 2016.

# Aside: In-Memory Crossbar Computation



$$( \ i_1 \quad i_2 \quad i_3 \quad i_4 \ ) = (o_1 \ o_2 \ o_3 \ o_4)$$

$$I_1 = \frac{1}{R_{11}}V_1 + \frac{1}{R_{21}}V_2 + \frac{1}{R_{31}}V_3 + \frac{1}{R_{41}}V_4$$

# Tutorial on Memory-Centric Computing:
## Processing-Using-Memory

Geraldo F. Oliveira

Prof. Onur Mutlu

ISCA 2024

29 June 2024

**SAFARI**

**ETH** *zürich*

# Agenda

- Introduction to Memory-Centric Computing Systems

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

- Coffee Break

- Real-World Processing-Near-Memory Systems

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

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

- PIM Programming & Infrastructure for PIM Research

- Closing Remarks

**SAFARI**