2<sup>nd</sup> Workshop on Memory-Centric Computing: Processing-Using-Memory - Part II

# Dr. Geraldo F. Oliveira https://geraldofojunior.github.io

ICS 2025 08 June 2025



**ETH** zürich

# SIMDRAM Follow-Up: 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 <u>30th International Symposium on High-</u> <u><i>Performance Computer Architecture*</u> (*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<sup>†</sup> Ataberk Olgun<sup>†</sup> Abdullah Giray Yağlıkçı<sup>†</sup> F. Nisa Bostancı<sup>†</sup> Juan Gómez-Luna<sup>†</sup> Saugata Ghose<sup>‡</sup> Onur Mutlu<sup>†</sup>

<sup>†</sup> ETH Zürich <sup>‡</sup> Univ. of Illinois Urbana-Champaign

#### SAFARI

https://arxiv.org/pdf/2402.19080.pdf

# Limitations of PUD Systems: Overview

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

#### SIMD Underutilization

- due to data parallelism variation within and across applications
- leads to throughput and energy waste

#### Limited Computation Support

- due to a lack of low-cost interconnects across columns
- limits PUD operations to only parallel map constructs

#### Challenging Programming Model

- due to a lack of compiler support for PUD systems
- creates a burden on programmers, limiting PUD adoption

# Limitations of PUD Systems: Overview

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

#### SIMD Underutilization

- due to data parallelism variation within and across applications
- leads to throughput and energy waste

#### Limited Computation Support

- due to a lack of low-cost interconnects across columns
- limits PUD operations to only parallel map constructs
- **Q** Challenging Programming Model
  - due to a lack of compiler support for PUD systems
  - creates a burden on programmers, limiting PUD adoption

Limitations of PUD

MIMDRAM Hard

Hardware Overvie w

Software Support

Conclusion

Evaluation

....

# Limitations of PUD Systems: Underutilization of SIMD Lanes (I)

#### **Application Analysis:**

#### quantify the fraction of SIMD parallelism in real applications



#### Ideal maximum vectorization factor = # DRAM columns (e.g., 65,536)

SAFARI

Introduction & Background

Limitations of PUD

. . . . . . . .

MIMDRAM

Hardware Overview

Software Support

Evaluation

....

5

# **Limitations of PUD Systems: Underutilization of SIMD Lanes (II)**

#### **Application Analysis:**

#### quantify the fraction of SIMD parallelism in real applications





Introduction & Background ......

Limitations of PUD . . . . . . . .

MIMDRAM ....

Hardware Overview ....

Evaluation Software Support . . . . . .

•

....

# **Limitations of PUD Systems: Underutilization of SIMD Lanes (II)**

#### **Application Analysis:**

#### quantify the fraction of SIMD parallelism in real applications



Introduction & Background . . . . . . . .

Limitations of PUD . . . . . . . .

Hardware Overview MIMDRAM . . . . ....

Software Support . . . . . .

Conclusion .

....

# Limitations of PUD Systems: Underutilization of SIMD Lanes (III)

#### **Application Analysis:**

#### quantify the fraction of SIMD parallelism in real applications



# 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

#### Limited Computation Support

- due to a lack of low-cost interconnects across columns
- limits PUD operations to only parallel map constructs

#### **Challenging Programming Model**

- due to a lack of compiler support for PUD systems
- creates a burden on programmers, limiting PUD adoption

Limitations of PUD

MIMDRAM Hard

. . . .

Hardware Overvie w

Software Support

Evaluation

....

# **Limitations of PUD Systems: Limited Computation Support**

PUD systems do not support vector reduction at low area cost since data movement is bounded to within a DRAM column



global sense amplifier

no direct communication path across columns

**Directly connecting all DRAM columns** using a custom all-to-all interconnect leads to large (i.e., 21%) area cost

Introduction & Background . . . . . . . .

Limitations of PUD . . . . . . . .

MIMDRAM . . . .

Hardware Overview ....

Software Support . . . . . .

Evaluation ....

Conclusion 10

.

# 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

#### Limited Computation Support

- due to a lack of low-cost interconnects across columns
- limits PUD operations to only parallel map constructs

#### Challenging Programming Model

- due to a lack of compiler support for PUD systems
- creates a burden on programmers, limiting PUD adoption



Limitations of PUD

MIMDRAM Hardwa

. . . .

Software Support

.

Evaluation

. . . . .

**Programmer's Tasks:** 

#### Goal:

Just write my kernel

```
High-level code for

C[i] = (A[i] > pred[i])? A[i] + B[i] : A[i] - B[i]
```

```
for (int i = 0; i < size ; ++ i){
    bool cond = A[i] > pred[i];
    if (cond) C[i] = A[i] + B[i];
    else C[i] = A[i] - B[i];
}
```



Limitations of PUD

MIMDRAM

Hardware Overview

Software Support

Conclusion

12

Programmer's Tasks:	Goal:
Map & align	Just write
data structures	my kernel

High-level code for C[i] = (A[i] > pred[i])? A[i] + B[i] : A[i] - B[i]



Limitations of PUD

MIMDRAM

Hardware Overview

Software Support

Conclusion

Programmer's Tasks:		Goal:	
Map & align	Identify		Just write
data structures	array boundaries		my kernel

#### High-level code for C[i] = (A[i] > pred[i])? A[i] + B[i] : A[i] - B[i]



Limitations of PUD

MIMDRAM

Hardware Overview

view Soft

Software Support

Conclusion

•

Evaluation

....

14

Programmer's Tasks:				Goal:
Map & align	Identify	Manually	Map C to	Just write
data structures	array boundaries	unroll loop	PUD instructions	my kernel

**High-level code for** C[i] = (A[i] > pred[i])? A[i] + B[i] : A[i] - B[i]

```
for (int i = 0; i < size ; ++ i){
   bool cond = A[i] > pred[i];
   if (cond) C[i] = A[i] + B[i];
   else C[i] = A[i] - B[i];
```



Limitations of PUD •••••

MIMDRAM ....

Hardware Overview ....

Software Support

Evaluation .... ....

•

15

Programmer's Tasks:				Goal:	
Map & align	Identify	Manually	Map C to	Orchestrate	Just write
data structures	array boundaries	unroll loop	PUD instructions	data movement	my kernel

**High-level code for** C[i] = (A[i] > pred[i])? A[i] + B[i] : A[i] - B[i]



Limitations of PUD

MIMDRAM

Hardware Overview

w Softwa

Software Support Evaluation

•

Programmer's Tasks:					
Map & align data structures	Identify array boundaries	Manually unroll loop	Map C to PUD instructions	Orchestrate data movement	Just write my kernel
	<b>P</b> C[i] = (	U <b>D's assem</b> A[i] > pred[i]	<b>bly-like code fo</b> )? A[i] + B[i] : A[i]	<b>r</b> — B[i]	
	<pre>bbop_trsp_ini bbop_trsp_ini bbop_trsp_ini</pre>	t(A , si t(B , si t(C , si	ze , elm_siz ze , elm_siz ze , elm_siz	ze); ze); ze);	
	<pre>bbop_add(D , bbop_sub(E ,</pre>	A, B, A, B,	size , elm_s size , elm_s	<pre>size); size);</pre>	

bbop\_greater(F , A , pred , size , elm\_size); bbop\_if\_else(C , D , E , F , size , elm\_size);

SAFARI

Limitations of PUD • • • • • • • •

MIMDRAM Hardware Overview ....

# **Problem & Goal**



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



#### **Fine-Grained DRAM:**

segments the global wordline to access individual DRAM mats

#### **Fine-Grained DRAM:**

#### segments the global wordline to access individual DRAM mats



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



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

#### **1** Improves SIMD utilization

for a single PUD operation, only access the DRAM mats with target data



### **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

### segmented global wordline



global sense amplifier

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

#### 1

#### mproves 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

### **7** Enables low-cost interconnects for vector reduction

- global and local data buses can be used for inter-/intra-mat communication



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

### 1

#### mproves SIMD utilization

- for a single PUD operation, only access the DRAM mats with target data
- for multiple PUD operations, execute independent operations concurrentl
   → multiple instruction, multiple data (MIMD) execution model
- **7** 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 **SAFARI** 

### 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

### Software

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

### 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

### Software

- new compiler support to transparently generate PUD instructions
- system support to map and execute PUD instructions

# MIMDRAM: Modifications to DRAM Chip





### MIMDRAM: Control Unit Design

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



### 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 Hardwar

- 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

- new compiler support to transparently generate PUD instructions
- system support to map and execute PUD instructions

### MIMDRAM: Compiler Support (I)

Transparently: <u>extract</u> SIMD parallelism from an application, and <u>schedule</u> PUD instructions while maximizing <u>utilization</u>

#### Three new LLVM-based passes targeting PUD execution



# **MIMDRAM:** Compiler Support (II)



Identify SIMD parallelism, generate PUD instructions, and set the appropriate vectorization factor

# **MIMDRAM:** Compiler Support (II)



Identify SIMD parallelism, generate PUD instructions, and set the appropriate vectorization factor

#### **Homology SIMD utilization by allowing the distribution of independent PUD** instructions across DRAM mats

# **MIMDRAM:** Compiler Support (III)



Identify SIMD parallelism, generate PUD instructions, and set the appropriate vectorization factor

त्नmprove SIMD utilization by allowing the distribution of independent PUD instructions across DRAM mats

Generate the appropriate binary for data allocation and PUD instructions

# MIMDRAM: System Support

- Instruction set architecture
- Execution & data transposition
- Data coherence
- Address translation
- Data allocation & alignment
- Mat label translation

# **Evaluation:**

### **Single Application Analysis – Energy Efficiency**



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

# **Evaluation:**

### **Multi-Programmed Workload Analysis**


#### **Evaluation:** More in the Paper

- MIMDRAM with subarray and bank-level parallelism
  - MIMDRAM provides significant performance gains compared to the baseline CPU (13.2x) and GPU (2x)
- Comparison to DRISA and Fulcrum for multi-programmed workloads
  - MIMDRAM achieves system throughput on par with DRISA and Fulcrum
- MIMDRAM's SIMD utilization versus SIMDRAM
  - MIMDRAM provides **15.6x** the utilization of SIMDRAM
- Area analysis
  - MIMDRAM adds small area cost to a DRAM chip (1.11%) and CPU die (0.6%)



#### Outline

#### Processing-Using-Memory (PUM) Solutions

- Review of DRAM Background, RowClone, and Ambit
- SIMDRAM & MIMDRAM: Generalizing PUM Computing
- pLUTo: Extending Operation Support with LUTs
- PUM in Off-the-Shelf DRAM Chips
- PUM Beyond Boolean/Arithmetic and DRAM
- Processing-Near-Memory (PNM) Solutions
  - Overview
  - Accelerating Neural Networks & Databases
  - Real PNM Systems
- Barriers for Adoption (and Current Solutions)
- Conclusion

#### 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 <u>55th International Symposium on Microarchitecture</u> (<i>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<sup>§</sup> Gabriel Falcao<sup>†</sup> Juan Gómez-Luna<sup>§</sup> Mohammed Alser§ Mohammad Sadrosadati<sup>§</sup> Lois Orosa<sup>§</sup>∇ Jeremie S. Kim<sup>§</sup> Geraldo F. Oliveira§ Taha Shahroodi<sup>‡</sup> Anant Nori\* Onur Mutlu<sup>§</sup> §ETH Zürich <sup>†</sup>IT, University of Coimbra <sup>∇</sup>*Galicia Supercomputing Center* <sup>‡</sup>TU Delft \*Intel

SAFARI

#### https://arxiv.org/pdf/2104.07699.pdf

#### Limitations of Processing-using-DRAM

Data Movement	RowClone, Seshadri+ 2013 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**



#### 

#### pLUTo: Key Idea



# pLUTo: Key Idea (x) (Insut)

# Replace computation with memory accesses $\rightarrow$ *pLUTo LUT Query* operation

















































?

DRAM array

**....** 

match logic



























DRAM array

**....** 

match logic





match logic

DRAM array

**....** 

#### pLUTo Designs: Tradeoff Space



pLUTo designs cover a *broad design space* and provide *different* performance, energy, and area efficiency

## **System Integration**



#### More in the Paper



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

João Dinis Ferreira§		Gabriel Falcao <sup>†</sup>		Juan Gómez-Luna <sup>§</sup>		Mohammed Alser§	
Lois Orosa <sup>§</sup> ∇		Mohammad Sadrosadati <sup>§</sup>		Jeremie S. Kim§		Geraldo F. Oliveira§	
	Ta	aha Shahroodi‡	Anant	Nori*	Onur Mutlu§		
§ETH Zürich	†IT, Uni	versity of Coimbra	$^{\nabla}Galicia$ Supercomputing Center		• <sup>‡</sup> TU Delft	*Intel	



https://arxiv.org/pdf/2104.07699.pdf



CpLUTo ISA Instructions



С

t-or

eb

ROW

## Performance

#### Average speedup across 7 real-world workloads



pLUTo significantly outperforms CPU, GPU and PnM baselines

## 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

#### SAFARI

### 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	$\hfill \bigtriangleup$ Subscribed $\checkmark$	凸 17	7	A⇒ Share	💥 Clip	∃+ Save	
321 views 9 days ago								
pLUTo: Enabling Massively Parallel Computation in DRAM via Lookup Tables								
Speaker: Geraldo F. Oliveiramore								

**SAFARI** 

#### https://youtu.be/9t1FJQ6nNw4?si=bhylWCLZde2DC7os

### Outline

#### Processing-Using-Memory (PUM) Solutions

- Review of DRAM Background, RowClone, and Ambit
- SIMDRAM & MIMDRAM: Generalizing PUM Computing
- pLUTo: Extending Operation Support with LUTs
- PUM in Off-the-Shelf DRAM Chips
- PUM Beyond Boolean/Arithmetic and DRAM
- Processing-Near-Memory (PNM) Solutions
  - Overview
  - Accelerating Neural Networks & Databases
  - Real PNM Systems
- Barriers for Adoption (and Current Solutions)
- Conclusion

### 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 <u>30th International Symposium on High-</u> <i>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**:

Can simultaneously activate up to48 rows in two neighboring subarrays

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

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

2

3

### **DRAM Testing Infrastructure**

- Developed from DRAM Bender [Olgun+, TCAD'23]\*
- Fine-grained control over DRAM commands, timings, and temperature



**SAFAR** \*Olgun et al., "<u>DRAM Bender: An Extensible and Versatile FPGA-based Infrastructure</u> 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 <sup>a</sup>	Chip Density	Chip Org.	Speed Rate
	9 (72)	М	N/A	4Gb	x8	2666MT/s
	5 (40)	А	N/A	4Gb	x8	2133MT/s
CV Uunin	1 (16)	А	N/A	8Gb	x8	2666MT/s
SK HYIIIX	1 (32)	А	18-14	4Gb	x4	2400MT/s
	1 (32)	А	16-49	8Gb	x4	2400MT/s
	1 (32)	М	16-22	8Gb	x4	2666MT/s
	1 (8)	F	21-02	4Gb	x8	2666MT/s
Samsung	2 (16)	D	21-10	8Gb	x8	2133MT/s
	1 (8)	А	22-12	8Gb	x8	3200MT/s

### The Capability of COTS DRAM Chips

### We demonstrate that COTS DRAM chips:

# Can simultaneously activate up to48 rows in two neighboring subarrays

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

### 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 **32 rows** 

A total of **48 rows** 

#### A total of **32 rows SAFARI**

### The Capability of COTS DRAM Chips

We demonstrate that COTS DRAM chips:



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

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



2

### **Characterization Methodology**

Sweep Row A and Row B addresses



• Sweep DRAM chip temperature

![](_page_83_Figure_4.jpeg)

![](_page_83_Picture_5.jpeg)

### **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

![](_page_84_Picture_5.jpeg)

### The Capability of COTS DRAM Chips

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

Can simultaneously activate up to 48 rows in two neighboring subarrays

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

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

3

### Key Idea

### Manipulate the bitline voltage to express a wide variety of functions using

multiple-row activation in neighboring subarrays

![](_page_86_Figure_3.jpeg)

### **Two-Input AND and NAND Operations**

![](_page_87_Figure_1.jpeg)

![](_page_87_Figure_2.jpeg)

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

### **Two-Input AND and NAND Operations**

![](_page_88_Figure_1.jpeg)

![](_page_88_Figure_2.jpeg)

 $V_{DD} = 1 \& GND = 0$ COM **REF** Х  $\left( \right)$ 0 1  $\mathbf{0}$ 1  $\mathbf{O}$  $\left( \right)$ 1 1 1 0 1 AND NAND

### **Characterization Methodology**

• Sweep Row A and Row B addresses

![](_page_89_Figure_2.jpeg)

### **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

#### SAFARI

### Real Processing Using Memory Prototype

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

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

Ataberk Olgun<sup>§†</sup> Juan Gómez Luna<sup>§</sup> Konstantinos Kanellopoulos<sup>§</sup> Behzad Salami<sup>§\*</sup> Hasan Hassan<sup>§</sup> Oğuz Ergin<sup>†</sup> Onur Mutlu<sup>§</sup> <sup>§</sup>ETH Zürich <sup>†</sup>TOBB ETÜ <sup>\*</sup>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

![](_page_92_Picture_4.jpeg)

![](_page_92_Picture_5.jpeg)

Easy-to-extend Memory Controller

2 ISA-transparent PuM Controller

### Software

![](_page_92_Picture_9.jpeg)

![](_page_92_Picture_10.jpeg)

2 Custom Supervisor Software

#### SAFARI

### Real Processing Using Memory Prototype

![](_page_93_Picture_1.jpeg)

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

## **PiDRAM Workflow**

![](_page_94_Figure_1.jpeg)

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

#### SAFARI

### Real Processing Using Memory Prototype

#### E README.md

Ø

#### **Building a PiDRAM Prototype**

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

#### **Rebuilding Steps**

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

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

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

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

#### **Generating DDR3 Controller IP sources**

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

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

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

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

#### **Microbenchmark Copy/Initialization Throughput**

![](_page_96_Figure_1.jpeg)

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

SAFARI @kasırga

### **PiDRAM is Open Source**

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

□       CMU-SAFARI / PiDRAM (Public)         □       □						
<> Code O Issues 11 Pull requests	🕑 Actions 🗄 Projects 🕮 Wiki	🛈 Security 🗠 Insights	<b>段</b> Settings			
양 master → 양 2 branches ा 0 tag	S	Go to file Add	file ▼ Code ▼	About	鐐	
olgunataberk Fix small mistake in REAL	iii olgunataberk Fix small mistake in README       46522cc on Dec 5, 2021       Iii commits					
controller-hardware	Add files via upload		7 months ago	Processing-using-Memory techniques.		
📘 fpga-zynq	fpga-zynq Adds instructions to reproduce two key results 7 months a			Prototype on a RISC-V rocket chip system		
README.md		7 months ago	our preprint: https://arxiv.org/abs/2111.00082			
i≡ README.md			Ø	C Readme		
PiDRAM	<ul> <li>☆ 21 stars</li> <li>③ 3 watching</li> <li>♀ 2 forks</li> </ul>					
PiDRAM is the first flexible end-to-er	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**

SAFARI (

kasırga

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

Sear $V_1 V > c_5 > arXiv:2111.00082$	ch All fields 🗸 Search
	Help   Advanced Search
Computer Science > Hardware Architecture	Download:
[Submitted on 29 Oct 2021 (v1), last revised 19 Dec 2021 (this version, v3)]	• PDF
PiDRAM: A Holistic End-to-end FPGA-based Framework for Processing-in-DRAM	Other formats
Ataberk Olgun, Juan Gómez Luna, Konstantinos Kanellopoulos, Behzad Salami, Hasan Hassan, Oğuz Ergin, Onur Mutlu	Current browse context: cs.AR
Processing-using-memory (PuM) techniques leverage the analog operation of memory cells to perform computation. Several recent works have demons	trated < prev   next >
techniques represent an opportunity for alleviating the data movement bottleneck at very low cost. However, system integration of PuM techniques impo non-trivial challenges that are yet to be solved. Design space exploration of potential solutions to the PuM integration challenges requires appropriate to	ses Change to browse by: ols to cs
develop necessary hardware and software components. Unfortunately, current specialized DRAM-testing platforms, or system simulators do not provide 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.,	the References & Citations <ul> <li>NASA ADS</li> <li>Google Scholar</li> <li>Semantic Scholar</li> </ul>
necessary modifications in the operating system, memory controller). We implement PiDRAM on an FPGA-based platform along with an open-source RI 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 built 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 c using PiDRAM's software and hardware components.	SC-V DBLP - CS Bibliography listing   bibtex ulk Juan Gómez-Luna Behzad Salami Hasan Hassan Oguz Ergin Onur Muthu
Comments: 15 pages, 12 figures	Export Bibtex Citation
Subjects: Hardware Architecture (cs.AR)	
Cite as: arXiv:2111.00082 [CS.AR]	Bookmark
https://doi.org/10.48550/arXiv.2111.00082	💥 🔛 😥 🔤

### Long Talk + Tutorial on Youtube

### https://youtu.be/s z S6FYpC8

![](_page_99_Figure_2.jpeg)

### Outline

#### Processing-Using-Memory (PUM) Solutions

- Review of DRAM Background, RowClone, and Ambit
- SIMDRAM & MIMDRAM: Generalizing PUM Computing
- pLUTo: Extending Operation Support with LUTs
- PUM in Off-the-Shelf DRAM Chips

#### **PUM Beyond Boolean/Arithmetic and DRAM**

- Processing-Near-Memory (PNM) Solutions
  - Overview
  - Accelerating Neural Networks & Databases
  - Real PNM Systems
- Barriers for Adoption (and Current Solutions)
- Conclusion

#### **SAFARI**

### In-DRAM Physical Unclonable Functions

- Jeremie S. Kim, Minesh Patel, Hasan Hassan, and Onur Mutlu, "The DRAM Latency PUF: Quickly Evaluating Physical Unclonable Functions by Exploiting the Latency-Reliability Tradeoff in Modern DRAM <u>Devices"</u> Proceedings of the <u>24th International Symposium on High-Performance Computer</u> <u>Architecture</u> (HPCA), Vienna, Austria, February 2018. [Lightning Talk Video]
  - [Slides (pptx) (pdf)] [Lightning Session Slides (pptx) (pdf)]
  - [Full Talk Lecture Video (28 minutes)]

#### The DRAM Latency PUF:

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

Jeremie S. Kim<sup>†§</sup> Minesh Patel<sup>§</sup> Hasan Hassan<sup>§</sup> Onur Mutlu<sup>§†</sup> <sup>†</sup>Carnegie Mellon University <sup>§</sup>ETH Zürich

### In-DRAM True Random Number Generation

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

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

Jeremie S. Kim<sup>‡§</sup>

Minesh Patel<sup>§</sup> Hasan Hassan<sup>§</sup> Lois Orosa<sup>§</sup> Onur Mutlu<sup>§‡</sup> <sup>‡</sup>Carnegie Mellon University <sup>§</sup>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 <u>48th International Symposium on Computer Architecture</u> (ISCA), Virtual, June 2021.
 [Slides (pptx) (pdf)]
 [Short Talk Slides (pptx) (pdf)]
 [Talk Video (25 minutes)]
 [SAFARI Live Seminar Video (1 hr 26 mins)]

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

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

### In-DRAM True Random Number Generation

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

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

F. Nisa Bostanci<sup>†§</sup> Ataberk Olgun<sup>†§</sup> Lois Orosa<sup>§</sup> Jeremie S. Kim<sup>§</sup> Hasan Hassan<sup>§</sup> Oğuz Ergin<sup>†</sup>

<sup>†</sup>TOBB University of Economics and Technology

A. Giray Yağlıkçı<sup>§</sup> Onur Mutlu<sup>§</sup> <sup>§</sup>ETH Zürich

SAFARI

https://arxiv.org/pdf/2201.01385.pdf

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

Shuangchen Li<sup>1</sup>; Cong Xu<sup>2</sup>, Qiaosha Zou<sup>1,5</sup>, Jishen Zhao<sup>3</sup>, Yu Lu<sup>4</sup>, and Yuan Xie<sup>1</sup>

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

#### **SAFARI** <u>https://cseweb.ucsd.edu/~jzhao/files/Pinatubo-dac2016.pdf</u> <sup>106</sup>

### Pinatubo: RowClone and Bitwise Ops in PCM

![](_page_106_Figure_1.jpeg)

proach, 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 <u>Onur Mutlu</u>,
 "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<sup>§∇</sup> Roknoddin Azizi<sup>§</sup> Geraldo F. Oliveira<sup>§</sup> Mohammad Sadrosadati<sup>§</sup> Rakesh Nadig<sup>§</sup> David Novo<sup>†</sup> Juan Gómez-Luna<sup>§</sup> Myungsuk Kim<sup>‡</sup> Onur Mutlu<sup>§</sup>

<sup>§</sup>ETH Zürich <sup>¬</sup>POSTECH <sup>†</sup>LIRMM, Univ. Montpellier, CNRS <sup>‡</sup>Kyungpook National University

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

SAFARI
# 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.

### SAFARI

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

## Aside: In-Memory Crossbar Computation



#### SAFARI

2<sup>nd</sup> Workshop on Memory-Centric Computing: Processing-Using-Memory - Part II

# Dr. Geraldo F. Oliveira https://geraldofojunior.github.io

ICS 2025 08 June 2025



**ETH** zürich