2025 International Conference on Supercomputing (ICS)



Achieving High-Performance Processing-Using-DRAM with Dynamic Bit-Precision, Adaptive Data Representation, and Flexible Arithmetic

Geraldo F. Oliveira

Mayank KabraYuxin GuoKangqi ChenGiray YağlıkçıMelina SoysalMohammad SadrJoaquin O. BuenoSaugata GhoseJuan Gómez-LunaOnur Mutlu













Executive Summary

Problem: Processing-Using-DRAM (PUD) suffers from three limitations caused naive employing a bit-serial execution model

- <u>Rigid and static data representation</u> leading to unnecessary computation
- <u>Throughput-oriented execution</u> leading to limited latency tolerance
- <u>Scalability challenge</u>s for high-precision operations

Goal: Design a PUD system that overcomes the three limitations caused by naively employing a bit-serial execution model

Key Mechanism: Proteus, a data-aware runtime PUD framework that

- 1. <u>dynamically</u> adjust the bit-precision, and based on that,
- 2. <u>chooses</u> and <u>uses</u> the most appropriate data representation and arithmetic algorithm implementation

for a given PUD operation

Key Results: Proteus achieves

- 17x (90.3x), 7.3x (21x), and 10.2x (8.1x) the performance per mm² (energy reduction) of a high-end CPU, GPU, and state-of-the-art PUD system
- Small area cost to a DRAM chip (1.6%) and CPU die (0.03%)

SAFARI

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

Outline

1	Introduction
2	Background
3	Limitations of PUD Systems
4	Proteus: Key Ideas
5	Proteus: Overview
6	Evaluation
7	Conclusion

Outline

1	Introduction
2	Background
3	Limitations of PUD Systems
4	Proteus: Key Ideas
5	Proteus: Overview
6	Evaluation
7	Conclusion

Processing-in-Memory: Overview

Two main approaches for Processing-in-Memory:

- Processing-<u>Near</u>-Memory: PIM logic is added to the same die as memory of to the logic layer of 3D-stacked memory
- 2 Processing-Using-Memory: uses the operational principles of memory cells to perform computation



Outline

1	Introduction
2	Background
3	Limitations of PUD Systems
4	Proteus: Key Ideas
5	Proteus: Overview
6	Evaluation
7	Conclusion

Processing-in-Memory: Overview

Two main approaches for Processing-in-Memory:

- **1** Processing-<u>Near</u>-Memory: PIM logic is added to the same die as memory of to the logic layer of 3D-stacked memory
- 2 Processing-Using-Memory: uses the operational principles of memory cells to perform computation



Background: DRAM Hierarchical Organization



Background: DRAM Hierarchical Organization



Background: DRAM Operation – Row Access (ACTIVATE)



Background: DRAM Operation – Column Access (READ)



Background: In-DRAM Majority Operations

In-DRAM majority is performed by simultaneously activating three DRAM rows



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

Bulk bitwise arithmetic can be performed by orchestrating in-DRAM row copy and majority operations





Bulk bitwise arithmetic can be performed by orchestrating in-DRAM row copy and majority operations





Bulk bitwise arithmetic can be performed by orchestrating in-DRAM row copy and majority operations





Bulk bitwise arithmetic can be performed by orchestrating in-DRAM row copy and majority operations





Bulk bitwise arithmetic can be performed by orchestrating in-DRAM row copy and majority operations





Bulk bitwise arithmetic can be performed by orchestrating in-DRAM row copy and majority operations



Oliveira, Geraldo F., et al. "SIMDRAM: An End-to-End Framework for Bit-Serial SIMD Computing in DRAM," in ASPLOS, 2021

Bulk bitwise arithmetic can be performed by orchestrating in-DRAM row copy and majority operations



Processing-Using-DRAM architectures (e.g., SIMDRAM) are very-wide (e.g., 65,536 wide) bit-serial SIMD engines





Background: PUD in Commodity Off-the-Shelf DRAM

Commodity off-the-shelf DRAM chips can

perform bulk bitwise operations without hardware modifications

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

Processing-using-DRAM (PuD) is an emerging paradigm that leverages the analog operational properties of DRAM circuitry to enable massively parallel in-DRAM computation. PuD has the potential to significantly reduce or eliminate costly data movement between processing elements and main memory. A common approach for PuD architectures is to make use of bulk bitwise computation (e.g., AND, OR, NOT). Prior works experimentally demonstrate three-input MAJ (i.e., MAJ3) and two-input AND and OR operations in commercial off-the-shelf (COTS) DRAM chips. Yet, demonstrations on COTS DRAM chips do not provide a functionally complete set of operations (e.g., NAND or AND and NOT).

systems and applications [12, 13]. Processing-using-DRAM (PuD) [29-32] is a promising paradigm that can alleviate the data movement bottleneck. PuD uses the analog operational properties of the DRAM circuitry to enable massively parallel in-DRAM computation. Many prior works [29-53] demonstrate that PuD can greatly reduce or eliminate data movement.

A widely used approach for PuD is to perform bulk bitwise operations, i.e., bitwise operations on large bit vectors. To perform bulk bitwise operations using DRAM, prior works propose modifications to the DRAM circuitry [29-31, 33, 35, 36, 43, 44, 46, 48-58]. Recent works [38, 41, 42, 45] experimentally demonstrate the feasibility of executing data copy & initializa-

https://arxiv.org/pdf/2402.18736.pdf

Outline

1	Introduction
2	Background
3	Limitations of PUD Systems
4	Proteus: Key Ideas
5	Proteus: Overview
6	Evaluation
7	Conclusion

Limitations of Bit-Serial PUD Systems: Overview

PUD systems suffer from three sources of inefficiency due to the naive use of a bit-serial execution model

- **1** Rigid and Static Data Representation
 - leading to a subpar performance in the presence of narrow values

- **7** Throughput-Oriented Execution with Low Latency Tolerance
 - failing to reduce the latency of a single PUD operation

- **3** Scalability Challenges for High-Precision Operations
 - due to the linear/quadratically scaling nature of bit-serial algorithms

Limitations of Bit-Serial PUD Systems: Overview

PUD systems suffer from three sources of inefficiency due to the naive use of a bit-serial execution model

- **1** Rigid and Static Data Representation
 - leading to a subpar performance in the presence of narrow values

- **7** Throughput-Oriented Execution with Low Latency Tolerance
 - failing to reduce the latency of a single PUD operation

- **3** Scalability Challenges for High-Precision Operations
 - due to the linear/quadratically scaling nature of bit-serial algorithms



















Application Analysis:

quantify the amount required bit-precision in real applications



---> minimum number of bits required to represent input operands

Application Analysis:

quantify the amount required bit-precision in real applications



Applications display a significant amount of <u>narrow values</u> \rightarrow bit-precision can be <u>reduced</u> from 32 bits to 20 bits, on average

Limitations of Bit-Serial PUD Systems: Overview

PUD systems suffer from three sources of inefficiency due to the naive use of a bit-serial execution model

1 Rigid and Static Data Representation

leading to a subpar performance in the presence of narrow values

7 Throughput-Oriented Execution with Low Latency Tolerance

failing to reduce the latency of a single PUD operation

- **3** Scalability Challenges for High-Precision Operations
 - due to the linear/quadratically scaling nature of bit-serial algorithms

Limitations of Bit-Serial PUD Systems: Throughput Oriented Execution (I)




























SAFARI

41





SAFARI

42





SAFARI

43







Limitations of Bit-Serial PUD Systems: Overview

PUD systems suffer from three sources of inefficiency due to the naive use of a bit-serial execution model

- **1** Rigid and Static Data Representation
 - leading to a subpar performance in the presence of narrow values

- **7** Throughput-Oriented Execution with Low Latency Tolerance
 - failing to reduce the latency of a single PUD operation

- **3** Scalability Challenges for High-Precision Operations
 - due to the linear/quadratically scaling nature of bit-serial algorithms

Limitations of Bit-Serial PUD Systems: Scalability Challenges (I)



Problem & Goal

Naively employing a bit-serial execution model for Processing-
Using-DRAM architectures
limits its applicability and
efficiency for different applications

Design a <u>flexible PUD system</u> that overcomes the three limitations caused by naively employing a bit-serial execution model



Goal

Outline

1	Introduction
2	Background
3	Limitations of PUD Systems
4	Proteus: Key Ideas
4 5	Proteus: Key Ideas Proteus: Overview
4 5 5	Proteus: Key Ideas Proteus: Overview Evaluation

Opportunities for Bit-Serial PUD: Overview

PUD systems suffer from three sources of inefficiency due to the naive use of a bit-serial execution model

1 Rigid and Static Data Representation

- leading to a subpar performance in the presence of narrow values
- <u>opportunity 1</u>: narrow values for PUD computation

7 Throughput-Oriented Execution with Low Latency Tolerance

- failing to reduce the latency of a single PUD operation
- <u>opportunity 2</u>: DRAM parallelism for latency-oriented execution

3 Scalability Challenges for High-Precision Operations

- due to the linear/quadratically scaling nature of bit-serial algorithms
- <u>opportunity 3</u>: alternative data representation for high-precision computation

Opportunities for Bit-Serial PUD: Overview

PUD systems suffer from three sources of inefficiency due to the naive use of a bit-serial execution model

Rigid and Static Data Representation

- leading to a subpar performance in the presence of narrow values
- <u>opportunity 1</u>: narrow values for PUD computation
- **7** Throughput-Oriented Execution with Low Latency Tolerance
 - failing to reduce the latency of a single PUD operation
 - <u>opportunity 2</u>: DRAM parallelism for latency-oriented execution
- **3** Scalability Challenges for High-Precision Operations
 - due to the linear/quadratically scaling nature of bit-serial algorithms
 - <u>opportunity 3</u>: alternative data representation for high-precision computation

Opportunities for Bit-Serial PUD: Narrow Values for PUD Operations

Narrow Values:



Opportunities for Bit-Serial PUD: Overview

PUD systems suffer from three sources of inefficiency due to the naive use of a bit-serial execution model

- **1** Rigid and Static Data Representation
 - leading to a subpar performance in the presence of narrow values
 - <u>opportunity 1</u>: narrow values for PUD computation

7 Throughput-Oriented Execution with Low Latency Tolerance

- failing to reduce the latency of a single PUD operation
- <u>opportunity 2</u>: DRAM parallelism for latency-oriented execution
- **3** Scalability Challenges for High-Precision Operations
 - due to the linear/quadratically scaling nature of bit-serial algorithms
 - <u>opportunity 3</u>: alternative data representation for high-precision computation



Proteus' one-bit per-subarray data mapping →

scatter bits of a data word across DRAM subarrays to expose bit-level parallelism



one-bit per-subarray → exposes bit-level parallelism















Opportunities for Bit-Serial PUD: Overview

PUD systems suffer from three sources of inefficiency due to the naive use of a bit-serial execution model

- **1** Rigid and Static Data Representation
 - leading to a subpar performance in the presence of narrow values
 - <u>opportunity 1</u>: narrow values for PUD computation
- **7** Throughput-Oriented Execution with Low Latency Tolerance
 - failing to reduce the latency of a single PUD operation
 - <u>opportunity 2</u>: DRAM parallelism for latency-oriented execution

3 Scalability Challenges for High-Precision Operations

- due to the linear/quadratically scaling nature of bit-serial algorithms
- <u>opportunity 3</u>: alternative data representation for high-precision computation

Opportunities for Bit-Serial PUD: Alternative Data Representation for High-Precision

Use a positional number system representation, i.e., the <u>redundant binary representation (RBR)</u> for high-precision PUD operations:

(1) the operation <u>no longer needs to propagate</u> carry bits
(2) the operation latency is <u>independent</u> of the data precision

SAFARI

Key Idea

Outline

1	Introduction
2	Background
3	Limitations of PUD Systems
4	Proteus: Key Ideas
5	Proteus: Overview
5	Evaluation
6	Conclusion

Proteus: Design Overview (I)

Proteus, a data-aware runtime mechanism that <u>dynamically</u> and <u>transparently</u> adapts the in-DRAM implementation, bit-precision, and data format based on the input data for efficient PUD execution

Proteus is composed of three main components:

- **1** Parallelism-Aware μProgram Library
- **2** Dynamic Bit-Precision Engine
- **3** μProgram Select Unit

Proteus: Design Overview (II)



Proteus is composed of three main components:

- 1 Parallelism-Aware μProgram Library hand-optimized implementations of different PUD operations with different performance vs. bit-precision trade-offs
- **2** Dynamic Bit-Precision Engine
- **3** μProgram Select Unit

Proteus: Design Overview (III)



Proteus is composed of three main components:

- **1 Parallelism-Aware μProgram Library** hand-optimized implementations of different PUD operations with different performance vs. bit-precision trade-offs
- 2 Dynamic Bit-Precision Engine identifies the dynamic range of evicted cache lines belonging to the target PUD operation
- **3** μProgram Select Unit

Proteus: Design Overview (IV)



Proteus is composed of three main components:

- **1 Parallelism-Aware μProgram Library** hand-optimized implementations of different PUD operations with different performance vs. bit-precision trade-offs
- 2 Dynamic Bit-Precision Engine identifies the dynamic range of evicted cache lines belonging to the target PUD operation
- **3** μProgram Select Unit

identifies the appropriate bit-precision based on the input operations of the target PUD operation

Proteus: Design Overview (IV)



Proteus is composed of three main components:

- **1 Parallelism-Aware μProgram Library** hand-optimized implementations of different PUD operations with different performance vs. bit-precision trade-offs
- 2 Dynamic Bit-Precision Engine identifies the dynamic range of evicted cache lines belonging to the target PUD operation
- **3** μProgram Select Unit

identifies the appropriate bit-precision based on the input operations of the target PUD operation

Outline

1	Introduction
2	Background
3	Limitations of PUD Systems
4	Proteus: Key Ideas
5	Proteus: Overview
6	Evaluation
7	Conclusion

Evaluation: Methodology Overview

- Evaluation Setup
 - CPU: Intel Skylake CPU
 - GPU: NVIDIA A100 GPU (with and without Tensor Cores)
 - PUD: SIMDRAM [Oliveira+, 2021]
 - https://github.com/CMU-SAFARI/Proteus
- Workloads
 - 12 workloads from Polybench, Rodinia, Phoenix, and SPEC2017
 - Linear algebra, data mining, video processing, machine learning
- Metrics
 - CPU-Normalized **performance per area**
 - Energy reduction (compared to baseline CPU)

Evaluation: Performance Analysis



Proteus significantly improves performance/mm² compared to CPU (17x), GPU (7.3x), and SIMDRAM (10.2x)

Evaluation: Energy Analysis



Proteus significantly improves energy-efficiency compared to CPU (90.3x), GPU (21x), and SIMDRAM (8.1x)

SAFARI

CPU-Normalized Energy
Outline

SAFARI		73
7	Conclusion	
6	Evaluation	
5	Proteus: Overview	
4	Proteus: Key Ideas	
3	Limitations of PUD Systems	
2	Background	
1	Introduction	

Conclusion

We introduce *Proteus,* a data-aware hardware runtime PUD framework that

that <u>dynamically</u> adjusts the bit-precision and, based on that, <u>chooses</u> and <u>uses</u> the most appropriate data representation and arithmetic algorithm implementation for a given PUD operation

Our extensive evaluation shows that *Proteus* achieves

17x (90.3x), 7.3x (21x), and 10.2x (8.1x) the performance per mm² (energy reduction) of a CPU, GPU, and state-of-the-art PUD system, small area cost to a DRAM chip (1.6%) and CPU die (0.03%)

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

SAFARI

2025 International Conference on Supercomputing (ICS)



Achieving High-Performance Processing-Using-DRAM with Dynamic Bit-Precision, Adaptive Data Representation, and Flexible Arithmetic

Geraldo F. Oliveira

Mayank KabraYuxin GuoKangqi ChenGiray YağlıkçıMelina SoysalMohammad SadrJoaquin O. BuenoSaugata GhoseJuan Gómez-LunaOnur Mutlu











