

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

Journal Article

Author(s):

Olgun, Ataberk; [Gómez Luna, Juan](#) ; Kanellopoulos, Konstantinos; Salami, Behzad; Hassan, Hasan; Ergin, Oguz; Mutlu, Onur

Publication date:

2023-03

Permanent link:

<https://doi.org/10.3929/ethz-b-000601590>

Rights / license:

[Creative Commons Attribution 4.0 International](#)

Originally published in:

ACM Transactions on Architecture and Code Optimization 20(1), <https://doi.org/10.1145/3563697>



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

ATABERK OLGUN, JUAN GÓMEZ LUNA, and KONSTANTINOS KANELLOPOULOS,

ETH Zurich, Switzerland

BEHZAD SALAMI, SAFARI Research Group, Switzerland

HASAN HASSAN, ETH Zurich, Switzerland

OGUZ ERGIN, TOBB University of Economics and Technology, Turkey

ONUR MUTLU, ETH Zurich, Switzerland

Commodity DRAM-based processing-using-memory (PuM) techniques that are supported by off-the-shelf DRAM chips present an opportunity for alleviating the data movement bottleneck at low cost. However, system integration of these techniques imposes non-trivial challenges that are yet to be solved. Potential solutions to the integration challenges require appropriate tools to develop any necessary hardware and software components. Unfortunately, current proprietary computing systems, specialized DRAM-testing platforms, or system simulators do not provide the flexibility and/or the holistic system view that is necessary to properly evaluate and deal with the integration challenges of commodity DRAM-based PuM techniques.

We design and develop Processing-in-DRAM (PiDRAM), the first flexible end-to-end framework that enables system integration studies and evaluation of real, commodity DRAM-based PuM techniques. PiDRAM provides software and hardware components to rapidly integrate PuM techniques across the whole system software and hardware stack. We implement PiDRAM on an FPGA-based RISC-V system. To demonstrate the flexibility and ease of use of PiDRAM, we implement and evaluate two state-of-the-art commodity DRAM-based PuM techniques: (i) in-DRAM copy and initialization (RowClone) and (ii) in-DRAM true random number generation (D-RaNGe). We describe how we solve key integration challenges to make such techniques work and be effective on a real-system prototype, including memory allocation, alignment, and coherence. We observe that end-to-end RowClone speeds up bulk copy and initialization operations by 14.6× and 12.6×, respectively, over conventional CPU copy, even when coherence is supported with inefficient cache flush operations. Over PiDRAM's extensible codebase, integrating both RowClone and D-RaNGe end-to-end on a real RISC-V system prototype takes only 388 lines of Verilog code and 643 lines of C++ code.

CCS Concepts: • **Computer systems organization** → **Architectures**; • **General and reference** → *Performance*; *Design*; *Experimentation*;

Additional Key Words and Phrases: Processing-using-memory, processing-in-memory, RISC-V, FPGA, DRAM, memory controllers

Authors' addresses: A. Olgun, J. G. Luna, K. Kanellopoulos, H. Hassan, and O. Mutlu, ETH Zurich, Gloriastrasse 35, Zurich, Switzerland, 8092; emails: {ataberk.olgun, juan.gomez, konstantinos.kanellopoulos, hasan.hasan, onur.mutlu}@safari.ethz.ch; B. Salami, SAFARI Research Group, Gloriastrasse 35, Zurich, Switzerland, 8092; email: behzadsalami@gmail.com; O. Ergin, TOBB University of Economics and Technology, Sogutozu Caddesi, Ankara, Turkey, 06560; email: oergin@etu.edu.tr.



This work is licensed under a Creative Commons Attribution International 4.0 License.

© 2022 Copyright held by the owner/author(s).

1544-3566/2022/11-ART8

<https://doi.org/10.1145/3563697>

ACM Reference format:

Ataberk Olgun, Juan Gómez Luna, Konstantinos Kanellopoulos, Behzad Salami, Hasan Hassan, Oguz Ergin, and Onur Mutlu. 2022. PiDRAM: A Holistic End-to-end FPGA-based Framework for Processing-in-DRAM. *ACM Trans. Archit. Code Optim.* 20, 1, Article 8 (November 2022), 31 pages.
<https://doi.org/10.1145/3563697>

1 INTRODUCTION

Main memory is a major performance and energy bottleneck in computing systems [48, 120]. One way of overcoming the main memory bottleneck is to move computation into/near memory, a paradigm known as **processing-in-memory (PiM)** [120]. PiM reduces memory latency between the memory units and the compute units, enables the compute units to exploit the large internal bandwidth within memory devices, and reduces the overall power consumption of the system by eliminating the need for transferring data over power-hungry off-chip interfaces [48, 120].

Recent works propose a variety of PiM techniques to alleviate the data movement problem. One set of techniques propose to place compute logic *near* memory arrays (e.g., processing capability in the memory controller, logic layer of three-dimensional- (3D) stacked memory, or near the memory array within the memory chip) [2–4, 12, 20–22, 25, 31, 34, 37, 39, 45–47, 51, 53, 57, 58, 65, 67, 79, 80, 86, 111, 118, 121, 131, 133, 152, 161, 172, 175–177]. These techniques are called **processing-near-memory (PnM)** techniques [120]. Another set of techniques propose to leverage analog properties of memory (e.g., Static Random-Access Memory, Dynamic Random-Access Memory (DRAM), Non-Volatile Memory) operation to perform computation in different ways (e.g., leveraging non-deterministic behavior in memory array operation to generate random numbers, performing bitwise operations within the memory array by exploiting analog charge sharing properties of DRAM operation) [1, 5–9, 17, 19, 24, 28, 32, 36, 42–44, 54–56, 69, 73, 82, 83, 91–93, 102–104, 114, 134, 136, 145, 147, 151, 155, 159, 162, 167, 170, 171]. These techniques are known as **processing-using-memory (PuM)** techniques [120].

A subset of PuM proposals devise mechanisms that enable computation using DRAM arrays [5, 6, 28, 32, 44, 54, 82, 83, 103, 134, 145, 147, 159, 167]. These mechanisms provide significant performance benefits and energy savings by exploiting the high internal bit-level parallelism of DRAM for (1) bulk data copy and initialization operations at row granularity [1, 28, 134, 145, 159], (2) bitwise operations [7–9, 103, 104, 114, 142, 144, 146–148, 167], (3) arithmetic operations [1, 6, 17, 32, 36, 42, 43, 55, 56, 73, 91–93, 102, 103, 151, 162, 170], and (4) security primitives (e.g., true random number generation [83] and physical unclonable functions [82, 126]). Recent works [44, 82, 83] show that some of these PuM mechanisms can already be reliably supported in contemporary, off-the-shelf DRAM chips.¹ Given that DRAM is the dominant main memory technology, these commodity DRAM-based PuM techniques provide a promising way to improve the performance and energy efficiency of existing and future systems at *no additional DRAM hardware cost*.

Integration of these PuM mechanisms in a real system imposes non-trivial challenges that require further research to find appropriate solutions. For example, in-DRAM bulk data copy and initialization techniques [28, 147] require modifications to memory management that affect different parts of the system. First, these techniques have specific memory allocation and alignment requirements (e.g., page-granularity source and destination operand arrays should be allocated and aligned in the same DRAM subarray) that are *not* satisfied by existing memory allocation primitives (e.g., `malloc` [106] and `posix_memalign` [108]). Second, in-DRAM copy requires efficient handling of memory coherence, such that the contents of the source operand in DRAM are up-to-date.

¹We are especially interested in PiM techniques that do *not* require any modification to the DRAM chips or the DRAM interface.

None of these system integration challenges of PuM mechanisms can be efficiently studied in existing general-purpose computing systems (e.g., personal computers, cloud computers, and embedded systems), special-purpose testing platforms (e.g., SoftMC [60]), or system simulators (e.g., gem5 [18, 132], Ramulator [90, 137], Ramulator-PIM [139], zsim [140], DAMOVSim [125, 138], and other simulators [35, 168, 169, 174]). Existing general-purpose computing systems do *not* permit dynamically changing DDRx timing parameters, which is required to integrate many PuM mechanisms into real systems. Although special-purpose testing platforms can be used to dynamically change DDRx timing parameters, these platforms do *not* model an end-to-end computing system where system integration of PuM mechanisms can be studied. System simulators do *not* model DRAM operation that violates manufacturer-recommended timing parameters and do *not* have a way of interacting with real DRAM chips that embody undisclosed and unique characteristics that have implications on how PuM techniques are integrated into real systems.

Our goal is to design and implement a flexible real-system platform that can be used to solve system integration challenges and analyze tradeoffs of end-to-end implementations of commodity DRAM-based PuM mechanisms. To this end, we develop **Processing-in-DRAM (PiDRAM)** framework, the first flexible, end-to-end, and open source framework that enables system integration studies and evaluation of real PuM techniques using real unmodified DRAM devices.

PiDRAM facilitates system integration studies of new commodity DRAM-based PuM mechanisms by providing four customizable hardware and software components that can be used as a common basis to enable system support. PiDRAM contains two main *hardware* components. First, a custom, easy-to-extend *memory controller* allows for implementing new DRAM command sequences that perform PuM operations. For example, the memory controller can be extended with a single state machine in its hardware description to implement a new DDRx command sequence with user-defined timing parameters to implement a new PuM technique (i.e., perform a new PuM operation). Second, an *ISA-transparent controller (PuM Operations Controller (POC))* supervises PuM execution. POC exposes the PuM operations to the software components of PiDRAM over a memory-mapped interface to the processor, allowing the programmer to perform PuM operations using the PiDRAM framework by executing conventional LOAD/STORE instructions. The memory-mapped interface allows PiDRAM to be easily ported to systems that implement different instruction set architectures. PiDRAM contains two main *software* components. First, an *extensible library* allows system designers to implement software support for PuM mechanisms. This library contains customizable functions that communicate with POC to perform PuM operations. Second, a custom *supervisor software* contains the necessary OS primitives (e.g., memory management) to enable end-to-end implementations of commodity DRAM-based PuM techniques.

We demonstrate a prototype of PiDRAM on an FPGA-based RISC-V system [11]. To demonstrate the flexibility and ease of use of PiDRAM, we implement two prominent PuM techniques: (1) *RowClone* [145], an in-DRAM data copy and initialization technique, and (2) an **in-DRAM true random number generation technique (D-RaNGe)** [83] based on activation-latency failures. To support RowClone (Section 5), (i) we customize the PiDRAM memory controller to issue carefully engineered sequences of DRAM commands that perform data copy (and initialization) operations in DRAM, and (ii) we extend the custom supervisor software to implement a new memory management mechanism that satisfies the memory allocation and alignment requirements of RowClone. For D-RaNGe (Section 6), we extend (i) the PiDRAM memory controller with a new state machine that periodically performs DRAM accesses with reduced activation latencies to generate random numbers [83] and a new hardware *random number buffer* that stores the generated random numbers, and (ii) the custom supervisor software with a function that retrieves the random numbers from the hardware buffer to the user program. Our end-to-end evaluation of (i) RowClone demonstrates up to 14.6× speedup for bulk copy and 12.6× initialization operations over CPU copy (i.e.,

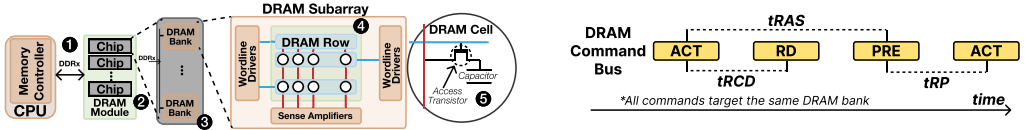


Fig. 1. DRAM organization (left). Timing diagram of DRAM commands (right).

conventional memcpy), even when coherence is satisfied using inefficient cache flush operations, and (ii) D-RaNGe demonstrates that an end-to-end integration of D-RaNGe can provide true random numbers at high throughput (8.30 Mb/s) and low latency (4-bit random number in 220 ns), even without any hardware or software optimizations. Implementing both PuM techniques over the Verilog and C++ codebase provided by PiDRAM requires only 388 lines of Verilog code and 643 lines of C++ code.

Our contributions are as follows:

- We develop PiDRAM, the first flexible framework that enables end-to-end integration and evaluation of PuM mechanisms using real unmodified DRAM chips.
- We develop a prototype of PiDRAM on an FPGA-based platform. To demonstrate the ease-of-use and evaluation benefits of PiDRAM, we implement two state-of-the-art DRAM-based PuM mechanisms, RowClone and D-RaNGe, and evaluate them on PiDRAM's prototype using unmodified DDR3 chips.
- We devise a new memory management mechanism that satisfies the memory allocation and alignment requirements of RowClone. We demonstrate that our mechanism enables RowClone end-to-end in the full system, and provides significant performance improvements over traditional CPU-based copy and initialization operations (memcpy [107] and calloc [105]) as demonstrated on our PiDRAM prototype.
- We implement and evaluate a state-of-the-art D-RaNGe. Our implementation provides a solid foundation for future work on system integration of DRAM-based PuM security primitives (e.g., PUFs [13, 82] and TRNGs [13, 123, 124]), implemented using real unmodified DRAM chips.

2 BACKGROUND

We provide the relevant background on DRAM organization, DRAM operation, and commodity DRAM-based PuM techniques. We refer the reader to prior works for more comprehensive background about DRAM organization and operation [26, 29, 49, 50, 87, 89, 95, 98, 100, 113, 123, 128].

2.1 DRAM Background

DRAM-based main memory is organized hierarchically. Figure 1 (left) depicts this organization. A processor is connected to one or more memory channels (DDR_x in the figure) ①. Each channel has its own command, address, and data buses. Multiple memory modules can be plugged into a single channel. Each module contains several DRAM chips ②. Each chip contains multiple DRAM banks that can be accessed independently ③. Data transfers between DRAM memory modules and processors occur at cache block granularity. The cache block size is typically 64 bytes in current systems.

Inside a DRAM bank, DRAM cells are laid out as a two-dimensional array of wordlines (i.e., DRAM rows) and bitlines (i.e., DRAM columns) ④. Wordlines are depicted in blue and bitlines are depicted in red in Figure 1. Wordline drivers drive the wordlines and sense amplifiers read the values on the bitlines. A DRAM cell is connected to a bitline via an access transistor ⑤. When enabled, an access transistor allows charge to flow between a DRAM cell and the cell's bitline.

DRAM Operation. When all DRAM rows in a bank are closed, DRAM bitlines are precharged to a reference voltage level of $\frac{V_{DD}}{2}$. The memory controller sends an **activate (ACT)** command to the DRAM module to drive a DRAM wordline (i.e., enable a DRAM row). Enabling a DRAM row starts the charge sharing process. Each DRAM cell connected to the DRAM row starts sharing its charge with its bitline. This causes the bitline voltage to deviate from $\frac{V_{DD}}{2}$ (i.e., the charge in the cell perturbs the bitline voltage). The sense amplifier senses the deviation in the bitline and amplifies the voltage of the bitline either to V_{DD} or to 0. As such, an ACT command copies one DRAM row to the sense amplifiers (i.e., row buffer). The memory controller can send READ/WRITE commands to transfer data from/to the sense amplifier array. Once the memory controller needs to access another DRAM row, the memory controller can close the enabled DRAM row by sending a **precharge (PRE)** command on the command bus. The PRE command first disconnects DRAM cells from their bitlines by disabling the enabled wordline and then precharges the bitlines to $\frac{V_{DD}}{2}$.

DRAM Timing Parameters. DRAM datasheets specify a set of timing parameters that define the minimum time window between valid combinations of DRAM commands [26, 27, 81, 97]. The memory controller must wait for tRCD, tRAS, and tRP nanoseconds between successive ACT → RD, ACT → ACT, and PRE → ACT commands, respectively (Figure 1, right). Prior works show that these timing parameters can be violated (e.g., successive ACT → RD commands may be issued with a shorter time window than tRCD) to improve DRAM access latency [26, 27, 81, 96, 97], implement physical unclonable functions [13, 82, 126], generate true random numbers [83, 123, 124], copy data [44, 145], and perform bitwise AND/OR operations [44, 142, 146–148] in commodity DRAM devices.

DRAM Internal Address Mapping. DRAM manufacturers use DRAM-internal address mapping schemes [30, 89, 130] to translate from logical (e.g., row, bank, column) DRAM addresses that are used by the memory controller to physical DRAM addresses that are internal to the DRAM chip (e.g., the physical position of a DRAM row within the chip). These schemes allow (i) post-manufacturing row repair techniques to map erroneous DRAM rows to redundant DRAM rows and (ii) DRAM manufacturers to organize DRAM internals in a cost-efficient and reliable way [76, 158]. DRAM-internal address mapping schemes can be substantially different across different DRAM chips [15, 30, 63, 70, 75–77, 88, 96, 110, 127, 129, 130, 141]. Thus, consecutive logical DRAM row addresses might not point to physical DRAM rows in the same subarray.

2.2 PuM Techniques

Prior work proposes a variety of in-DRAM computation mechanisms (i.e., PuM techniques) that (i) have great potential to improve system performance and energy efficiency [6, 28, 40, 54, 144–150] or (ii) can provide low-cost security primitives [13, 14, 82, 83, 124, 126]. A subset of these in-DRAM computation mechanisms are demonstrated on real DRAM chips [13, 44, 82, 83, 124, 126]. We describe the major relevant prior works briefly:

RowClone [145] is a low-cost DRAM architecture that can perform bulk data movement operations (e.g., copy, initialization) inside DRAM chips at high performance and low energy.

Ambit [144, 146, 147, 149, 150] is a new DRAM substrate that can perform (i) bitwise majority (and thus bitwise AND/OR) operations across three DRAM rows by simultaneously activating three DRAM rows and (ii) bitwise NOT operations on a DRAM row using 2-transistor 1-capacitor DRAM cells [72, 112].

ComputeDRAM [44] demonstrates in-DRAM copy (previously proposed by RowClone [145]) and bitwise AND/OR operations (previously proposed by Ambit [147]) on real DDR3 chips. ComputeDRAM performs in-DRAM operations by issuing carefully engineered, valid sequences of DRAM

commands with violated tRAS and tRP timing parameters (i.e., by not obeying manufacturer-recommended timing parameters defined in DRAM chip specifications [116]). By issuing command sequences with violated timing parameters, ComputeDRAM activates two or three DRAM rows in a DRAM bank in quick succession (i.e., performs two or three row activations). ComputeDRAM leverages (i) two row activations to transfer data between two DRAM rows and (ii) three row activations to perform the majority function in real unmodified DRAM chips.

D-RaNGe [83] is a state-of-the-art high-throughput DRAM-based true random number generation technique. D-RaNGe leverages the randomness in DRAM activation (tRCD) failures as its entropy source. D-RaNGe extracts random bits from DRAM cells that fail with 50% probability when accessed with a reduced (i.e., violated) tRCD. D-RaNGe demonstrates high-quality true random number generation on a vast number of real DRAM chips across multiple generations.

QUAC-TRNG [124] demonstrates that four DRAM rows can be activated in a quick succession using an ACT-PRE-ACT command sequence (called QUAC) with violated tRAS and tRP timing parameters in real DDR4 DRAM chips. QUAC-TRNG uses QUAC to generate true random numbers at high throughput and low latency.

3 MOTIVATION

Integrating DRAM-based PuM techniques into a real system requires modifications across the hardware and software stack. End-to-end implementations of PuM techniques require proper tools that (i) are flexible, to enable rapid development of PuM techniques and (ii) support real DRAM devices, to correctly observe the effects of reduced DRAM timing operations that are fundamental to enabling commodity DRAM-based PuM in real unmodified DRAM devices. Existing general-purpose computers, specialized DRAM testing platforms, and simulators (e.g., those mentioned in Section 1) cannot be used to study end-to-end implementations of commodity DRAM-based PuM techniques. We discuss the limitations of such computers, specialized testing platforms, and simulators in detail in Section 8.

Our *goal* is to develop a flexible end-to-end framework that enables rapid system integration of commodity DRAM-based PuM techniques and facilitates studies on end-to-end full-system implementations of PuM techniques using real DRAM devices. To this end, we develop PiDRAM.

4 PIDRAM

Implementing commodity DRAM-based PuM techniques end-to-end requires developing new **hardware (HW)** and **software (SW)** components or augmenting existing components with new functionality (e.g., memory allocation for RowClone requires a new memory allocation routine in the OS, Section 5.1). To ease the process of modifying various components across the hardware and software stack to implement new PuM techniques, PiDRAM provides key HW and SW components. Figure 2 presents an overview of the HW and SW components of the PiDRAM framework. In Section 4.3, we describe the general workflow for executing a PuM operation on PiDRAM.

4.1 Hardware Components

PiDRAM comprises two key hardware components. Both of these components are designed with the goal to provide a flexible and easy to use framework for evaluating PuM techniques.

❶ **PuM Operations Controller.** POC decodes and executes PiDRAM instructions that are used by the programmer to perform PuM operations. POC communicates with the rest of the system over two well-defined interfaces. First, it communicates with the CPU over a memory-mapped interface, where the CPU can send data to or receive data from POC using memory store and load instructions. The CPU accesses the memory-mapped registers (*instruction*, *data*, and *flag*

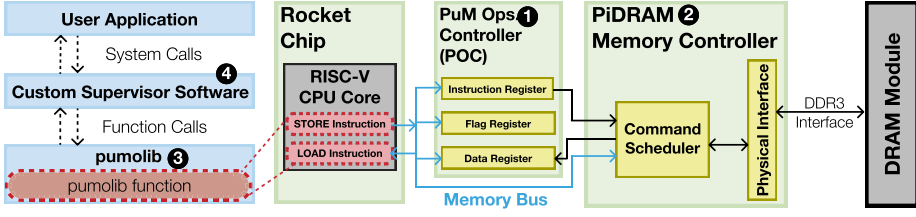


Fig. 2. PiDRAM overview. Modified hardware (in green) and software (in blue) components. Unmodified components are in gray. A pumolib function executes load and store instructions in the CPU to perform PuM operations (in red). We use yellow to highlight the key hardware structures that are controlled by the user to perform PuM operations.

registers) in POC to execute in-DRAM operations. This improves the portability of the framework and facilitates porting the framework to systems that employ different instruction set architectures. Second, POC communicates with the memory controller to perform PuM operations in the DRAM chip over a simple hardware interface. To do so, POC (i) requests the memory controller to perform a PuM operation, (ii) waits until the memory controller performs the operation, and (iii) receives the result of the PuM operation from the memory controller. The CPU can read the result of the operation by executing load instructions that target the *data* register in POC.

Custom Memory Controller. PiDRAM’s memory controller provides an easy-to-extend basis for commodity DRAM-based PuM techniques that require issuing DRAM commands with violated timing parameters [13, 44, 82, 83, 124]. The memory controller is designed modularly and requires easy-to-make modifications to its scheduler to implement new PuM techniques. For instance, RowClone operations (Section 5) is enabled in just 60 lines of Verilog code on top of the baseline custom memory controller’s scheduler that implements conventional DRAM operations (e.g., read, write).

The custom memory controller employs three key sub-modules to facilitate the implementation of new PuM techniques. (i) The *Periodic Operations Module* periodically issues DDR3 refresh [117] and interface maintenance commands [52]. (ii) A simple *DDR3 Command Scheduler* supports conventional DRAM operations (e.g., activate, precharge, read, and write). This scheduler applies an open-bank policy (i.e., DRAM banks are left open following a DRAM row activation) to exploit temporal locality in memory accesses to the DRAM module. LOAD/STORE memory requests are simply handled by the command scheduler in a latency-optimized way. Thus, new modules that are implemented to provide new PuM functionality (e.g., a state machine that controls the execution of a new PuM operation) in the custom memory controller do not compromise the performance of LOAD/STORE memory requests. (iii) The **Configuration Register File (CRF)** comprises 16 user-programmable registers that store the violated timing parameters used for DDRx sequences that trigger PuM operations (e.g., activation latency used in generating true random numbers using D-RaNGe [83], see Section 6) and miscellaneous parameters for PuM implementations (e.g., true random number generation period for D-RaNGe, see Section 6). In our implementation, CRF stores only the timing parameters used for performing PuM operations (e.g., RowClone and D-RaNGe). We do not store every standard DDRx timing parameter (i.e., non-violated, which are used exactly as defined as in DRAM chip specifications) in the CRF. Instead these timings are embedded in the command scheduler.

4.2 Software Components

PiDRAM comprises two key software components that complement and control PiDRAM’s hardware components to provide a flexible and easy to use end-to-end PuM framework.

Table 1. Pumolib Functions

Function	Arguments	Description
set_timings	RowClone_T1, RowClone_T2, tRCD	Updates CRF registers with the timing parameters used in RowClone (T1 and T2) and D-RaNGe (tRCD) operations.
rng_configure	period, address, bit_offsets	Updates CRF registers, configuring the random number generator to access the DRAM cache block at <i>address</i> every <i>period</i> cycles and collect the bits at <i>bit_offsets</i> from the cache block.
copy_row	source_address, destination_address	Performs a RowClone-Copy operation in DRAM from the <i>source_address</i> to the <i>destination_address</i> .
activation_failure	address	Induces an activation failure in a DRAM location pointed by the <i>address</i> .
buf_size	—	Returns the number of random words in the random number buffer.
rand_dram	—	Returns 32 bits (i.e., random words) from the random number buffer.

③ **PuM Operations Library (pumolib).** This extensible library allows system designers to implement software support for PuM techniques. Pumolib contains customizable functions that interface with POC to perform PuM operations in real unmodified DRAM chips. The customizable functions hide the hardware implementation details of PuM techniques from software developers (that use pimolib). For example, although we expose PuM techniques to software via memory LOAD/STORE operations (POC is exposed as a memory-mapped module, Section 4.1), PuM techniques can also be exposed via specialized instructions provided by ISA extensions. Pumolib hides such implementation details from the user of the library and contributes to the modular design of the framework.

We implement a general protocol that defines how programmers express the information required to execute PuM operations to the POC. A typical function in pumolib performs a PuM operation in four steps: It (i) writes a PiDRAM instruction to the POC’s *instruction* register, (ii) sets the *Start* flag in POC’s *flag* register, (iii) waits for the POC to set the *Ack* flag in POC’s *flag* register, and (iv) reads the result of the PuM operation from POC’s *data* register (e.g., the true random number after performing an in-DRAM true random number generation operation, Section 6). We list the currently implemented pumolib functions in Table 1.

④ **Custom Supervisor Software.** PiDRAM provides a custom supervisor software that implements the necessary OS primitives (i.e., virtual memory management, memory allocation, and alignment) for end-to-end implementation of PuM techniques. This facilitates developing end-to-end integration of PuM techniques in the system as these techniques require modifications across the software stack. For example, integrating RowClone end-to-end in the full system requires a new memory allocation mechanism (Section 5.1) that can satisfy the memory allocation constraints of RowClone [145]. Thus, we implement the necessary functions and data structures in the custom supervisor software to implement an allocation mechanism that satisfies RowClone’s constraints. This allows PiDRAM to be extended easily to implement support for new PuM techniques that share similar memory allocation constraints (e.g., Ambit [147], SIMDram [54], and QUAC-TRNG [124], as shown in Table 2).

4.3 Execution of a PuM Operation

We describe the general workflow for a PiDRAM operation (e.g., RowClone-Copy [145]) in Figure 3 over an example `copy_row()` function that is called by the user to perform a RowClone-Copy operation in DRAM.

The user makes a system call to the custom supervisor software ① that in turn calls the `copy_row(source, destination)` function in the pumolib ②. The function executes two store instructions in the RISC-V core ③. The first store instruction updates the *instruction* register with the `copy_row` instruction (that performs a RowClone-Copy operation in DRAM) ④ and the second store instruction sets the *Start* flag in the *flag* register to logic-1 ⑤ in POC. When the *Start* flag is set, POC instructs the PiDRAM memory controller to perform a RowClone-Copy operation using violated timing parameters ⑥. The POC waits until the memory controller starts executing the operation, after which it sets the *Start* flag to logic-0 and the *Ack* flag to logic-1 ⑦, indicating that it started the execution of the PuM operation. The PiDRAM memory controller performs the

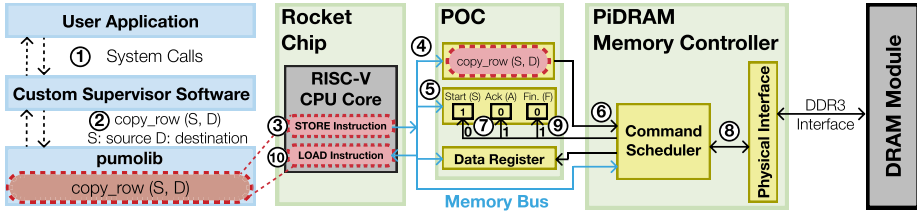


Fig. 3. Workflow for a PiDRAM RowClone-Copy operation.

Table 2. Various Known PuM Techniques That Can Be Studied Using PiDRAM

PuM Technique	Description	Integration Challenges
ComputeDRAM-based [44] RowClone [145]	Bulk data-copy and initialization within DRAM	(i) <i>memory allocation and alignment mechanisms</i> that map source & destination operands of a copy operation into same DRAM subarray; (ii) <i>memory coherence</i> , i.e., source operand must be up-to-date in DRAM.
D-RaNGe [83]	True random number generation using DRAM	(i) periodic generation of true random numbers; (ii) <i>memory scheduling policies</i> that minimize the interference caused by random number requests.
ComputeDRAM-based [44] Ambit [147]	Bitwise operations in DRAM	(i) <i>memory allocation and alignment mechanisms</i> that map operands of a bitwise operation into same DRAM subarray; (ii) <i>memory coherence</i> , i.e., operands of the bitwise operations must be up-to-date in DRAM.
SIMDRAM [54]	Arithmetic operations in DRAM	(i) <i>memory allocation and alignment mechanisms</i> that map operands of an arithmetic operation into same DRAM subarray; (ii) <i>memory coherence</i> , i.e., operands of the arithmetic operations must be up-to-date in DRAM; (iii) <i>bit transposition</i> , i.e., operand bits must be laid out vertically in a single DRAM bitline.
DL-PUF [82]	Physical unclonable functions in DRAM	<i>memory scheduling policies</i> that minimize the interference caused by generating PUF responses.
QUAC-TRNG [123] and Talakder+ [13]	True random number generation using DRAM	(i) periodic generation of true random numbers; (ii) <i>memory scheduling policies</i> that minimize the interference caused by random number requests; (iii) efficient integration of the SHA-256 cryptographic hash function.

PuM techniques we implement in this work are highlighted in bold.

RowClone-Copy operation by issuing a set of DRAM commands with violated timing parameters ⑧. When the last DRAM command is issued, the memory controller sets the Finish flag (denoted as Fin. in Figure 3) in the flag register to logic-1 ⑨, indicating the end of execution for the last PuM operation that the memory controller acknowledged. The copy function periodically checks either the Ack or the Finish flag in the flag register (depending on a user-supplied argument) by executing load instructions that target the flag register ⑩. When the periodically checked flag is set, the copy function returns. This way, the copy function optionally blocks until the start (i.e., the Ack flag is set) or the end (i.e., the Finish flag is set) of the execution of the PuM operation (in this example, RowClone-Copy).²

4.4 Use Cases

Beyond commodity DRAM-based PuM techniques [13, 44, 82, 83, 123], which PiDRAM can be used to study, many prior works propose minor modifications to DRAM arrays to enable various arithmetic [6, 32, 40, 54] and bitwise operations [6, 144, 147, 148, 150] and security primitives [126]. These PuM techniques share common memory allocation and coherence requirements (Section 5.1). PiDRAM facilitates developing new mechanisms that satisfy these requirements. Table 2 describes some of the PuM case studies PiDRAM can enable.³

Other than providing an easy-to-use basis for end-to-end implementations of commodity DRAM-based PuM techniques, PiDRAM can be easily extended with a programmable

²The data register is not used in a RowClone-Copy [145] operation, because the result of the RowClone-Copy operation is stored in *memory* (i.e., the source memory row is copied to the destination memory row). The data register is used in a D-RaNGe [83] operation, as described in Section 6.

³We acknowledge that PiDRAM's key components require modifications to implement new PuM techniques in PiDRAM and possibly to integrate PiDRAM into other systems. In fact, we quantify the degree of these modifications in our RowClone and D-RaNGe case studies. We show that the key components form a useful and easy-to-extend basis for PuM techniques with our Verilog and C code complexity analyses for both use cases (Sections 5.5.1 and 6.2).

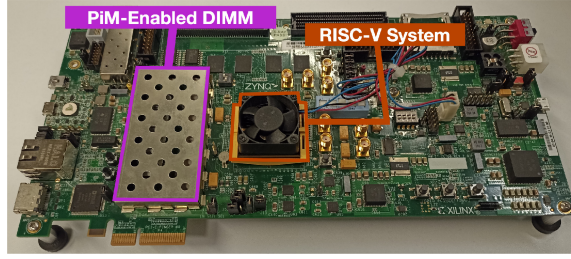


Fig. 4. PiDRAM's FPGA prototype.

microprocessor placed near the memory controller to study system integration challenges of PnM techniques (e.g., efficient pointer chasing [57, 58, 66], general-purpose compute [157], machine learning [74, 84, 94, 101, 122], databases [21, 22, 99], and graph processing [16]).

4.5 PiDRAM Prototype

We develop a prototype of the PiDRAM framework on an FPGA-based platform. We use the Xilinx ZC706 FPGA board [166] to interface with real DDR3 modules. Xilinx provides a DDR3 PHY IP [163] that exposes a low-level “DFI” interface [33] to the DDR3 module on the board. We use this interface to issue DRAM commands to the DDR3 module. We use the existing RISC-V-based SoC generator, Rocket Chip [11], to generate the RISC-V hardware system. Our custom supervisor software extends the RISC-V Proxy Kernel [135] to support the necessary OS primitives on PiDRAM's prototype. Figure 4 shows our prototype.

Simulation Infrastructure. To aid the users in testing the correctness of any modifications they make on top of PiDRAM, we provide the developers with a Verilog simulation environment that injects regular READ/WRITE commands and custom commands (e.g., update the CRF, perform RowClone-Copy, generate random numbers) to the memory controller. When used in conjunction with the Micron DDR3 Verilog model provided by Xilinx [163], the simulation environment can help the developers to easily understand if something unexpected is happening in their implementation (e.g., if timing parameters are violated).

Open Source Repository. We make PiDRAM freely available to the research community as open source software at <https://github.com/CMU-SAFARI/PiDRAM>. Our repository includes the full PiDRAM prototype that has RowClone (Section 5) and D-RaNGe (Section 6) implemented end-to-end on the RISC-V system.

5 CASE STUDY #1: END-TO-END ROWCLONE

We implement support for ComputeDRAM-based RowClone (in-DRAM copy/initialization) operations on PiDRAM to conduct a detailed study of the challenges associated with implementing RowClone end-to-end on a real system. None of the relevant prior works [44, 54, 142, 145, 147, 148, 150, 159] provide a clear description or a real system demonstration of a working memory allocation mechanism that can be implemented in a real operating system to expose RowClone capability to the programmer.

5.1 Implementation Challenges

Data Mapping. RowClone has four data mapping and alignment requirements that cannot be satisfied by current memory allocation mechanisms (e.g., malloc [106]). First, the source and destination operands (i.e., page (4-KiB)-sized arrays) of the copy operation must reside in the same DRAM subarray. We refer to this as the *mapping* requirement. Second, the source and destination

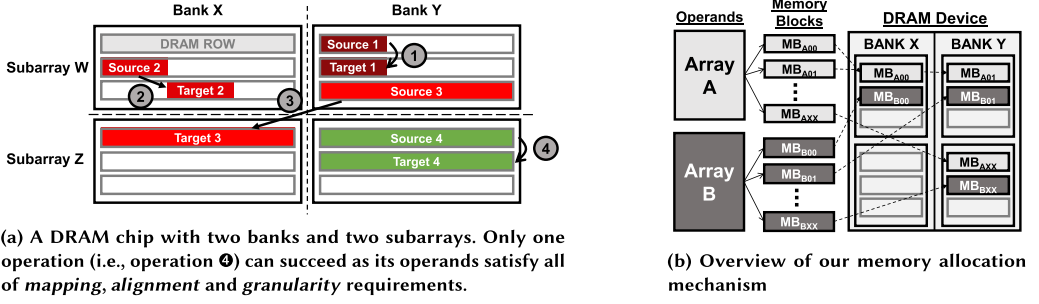


Fig. 5. RowClone memory allocation requirements (left) and memory allocation mechanism overview (right).

operands must be aligned to DRAM rows. We refer to this as the *alignment* requirement. Third, the size of the copied data must be a multiple of the DRAM row size. The size constraint defines the granularity at which we can perform bulk-copy operations using RowClone. We refer to this as the *granularity* requirement. Fourth, RowClone must operate on up-to-date data that resides in main memory. Modern systems employ caches to exploit locality in memory accesses and reduce memory latency. Thus, cache blocks (typically 64 B) of either the source or the destination operands of the RowClone operation may have cache block copies present in the cache hierarchy. Before performing RowClone, the cached copies of pieces of both source and destination operands must be invalidated and written back to main memory. We refer to this as the *memory coherence* requirement.

We explain the data mapping and alignment requirements of RowClone using Figure 5(a). The operand Source 1 cannot be copied to the operand Target 1 as the operands do not satisfy the *granularity* requirement (①). Performing such a copy operation would overwrite the remaining (i.e., non-Target 1) data in Target 1's DRAM row with the remaining (i.e., non-Source 1) data in Source 1's DRAM row. Source 2 cannot be copied to Target 2 as Target 2 is not *aligned* to its DRAM row (②). Source 3 cannot be copied to Target 3, as these operands are not *mapped* to the same DRAM subarray (③). In contrast, Source 4 can be copied to Target 4 using in-DRAM copy, because these operands are (i) *mapped* to the same DRAM subarray, (ii) aligned to their DRAM rows and (iii) occupy their rows completely (i.e., the operands have sizes equal to DRAM row size) (④).

5.2 Memory Allocation Mechanism

Computing systems employ various layers of address mappings that obfuscate the DRAM row-bank-column address mapping from the programmer [30, 61], which makes allocating source and target operands as depicted in Figure 5(a) (④) difficult. Only the virtual addresses are exposed to the programmer. Without control over the virtual address to DRAM address mapping, the programmer *cannot* easily place data in a way that satisfies the mapping and alignment requirements of an in-DRAM copy operation.

We implement a new memory allocation mechanism that can perform memory allocation for RowClone (in-DRAM copy/initialization) operations. This mechanism enables page-granularity RowClone operations (i.e., a virtual page can be copied to another virtual page using RowClone) *without* introducing any changes to the programming model. Figure 5(b) depicts an overview of our memory allocation mechanism.

At a high level, our memory allocation mechanism (i) splits the source and destination operands into page-sized virtually addressed memory blocks, (ii) allocates two physical pages in different DRAM rows in the same DRAM subarray, (iii) assigns these physical pages to virtual pages that correspond to the source and destination memory blocks at the same index such that the source block can be copied to the destination block using RowClone. We repeat this process until we

exhaust the page-sized memory blocks. As the mechanism processes subsequent page-sized memory blocks of the two operands, it allocates physical pages from a different DRAM bank to maximize bank-level parallelism in streaming accesses to these operands.

To overcome the *mapping*, *alignment*, and *granularity* problems, we implement our memory management mechanism in the custom supervisor software of PiDRAM. We expose the allocation mechanism using the `alloc_align(N, ID)` system call. The system call returns a pointer to a contiguous array of N bytes in the virtual address space (i.e., one operand). Multiple calls with the same ID to `alloc_align(N, ID)` place the allocated arrays in the same subarray in DRAM, such that they can be copied from one to another using RowClone. If N is too large such that it exceeds the size of available physical memory, then `alloc_align` fails and causes an exception. Our implementation of RowClone requires application developers to directly use `alloc_align` to allocate data instead of `malloc` and similar function calls.

The custom supervisor software maintains three key structures to make `alloc_align()` work: (i) **Subarray Mapping Table (SAMT)**, (ii) **Allocation ID Table (AIT)**, and (iii) **Initializer Rows Table (IRT)**.

(1) Subarray Mapping Table. We use the SAMT to maintain a list of physical page addresses that point to DRAM rows that are in the same DRAM subarray. `alloc_align()` queries SAMT to find physical addresses that map to rows in one subarray.

SAMT contains the physical pages that point to DRAM rows in each subarray. SAMT is indexed using subarray identifiers in the range $[0, \text{number of subarrays})$. SAMT contains an entry for every subarray. An entry consists of two elements: (i) the number of free physical address tuples and (ii) a list of physical address tuples. Each tuple in the list contains two physical addresses that respectively point to the first and second halves of the same DRAM row. The list of tuples contains all the physical addresses that point to DRAM rows in the DRAM subarray indexed by the SAMT entry. We allocate free physical pages listed in an entry and assign them to the virtual pages (i.e., memory blocks) that make up the row-copy operands (i.e., arrays) allocated by `alloc_align()`. We slightly modify our high-level memory allocation mechanism to allow for two memory blocks (4 KiB virtually addressed pages) of an array to be placed in the same DRAM row, as the page size in our system is 4 KiB and the size of a DRAM row is 8 KiB. We call two memory blocks in the same operand that are placed in the same DRAM row *sibling memory blocks* (also called sibling pages). The parameter N of the `alloc_align()` call defines this relationship: We designate memory blocks that are precisely $N/2$ bytes apart as *sibling memory blocks*.

Finding the DRAM Rows in a Subarray. Finding the DRAM row addresses that belong to the same subarray is not straightforward due to DRAM-internal mapping schemes employed by DRAM manufacturers (Section 2.1). It is extremely difficult to learn which DRAM address (i.e., bank-row-column) is actually mapped to a physical location (e.g., a subarray) in the DRAM device, as these mappings are not exposed through publicly accessible datasheets or standard definitions [71, 116, 130]. We make the key observation that the entire mapping scheme need *not* be available to successfully perform RowClone operations.

We observe that for a set of $\{source, destination\}$ DRAM row address pairs, RowClone operations repeatedly succeed with a 100% probability. We hypothesize that these pairs of DRAM row addresses are mapped to the same DRAM subarray. We identify these row address pairs by conducting a *RowClone success rate* experiment where we repeatedly perform RowClone operations between every *source, destination* row address pair in a DRAM bank. Our experiment works in three steps: We (i) initialize both the source and the destination row with random data, (ii) perform a RowClone operation from the source to the destination row, and (iii) compare the data in the destination row with the source row. RowClone success rate is calculated as the number of bits that differ between the source and destination rows' data divided by the number of bits stored

in a row (8 KiB in our prototype). If there is no difference between the source and the destination rows' data (i.e., the RowClone success rate for the source and the destination row is 100%), then we infer that the RowClone operation was successful. We repeat the experiment for 1000 iterations for each row address pair and if every iteration is successful, we store the address pair in the SAMT, indicating that the row address pair is mapped to different rows in the same DRAM subarray.⁴

(2) Allocation ID Table. To keep track of different operands that are allocated by `alloc_align` using the same *ID* (used to place different arrays in the same subarray), we use the AIT. AIT entries are indexed by *allocation IDs* (the parameter *ID* of the `alloc_align` call). Each AIT entry stores a pointer to an SAMT entry. The SAMT entry pointed by the AIT entry contains the set of physical addresses that were allocated using the same *allocation ID*. AIT entries are used by the `alloc_align` function to find which DRAM subarray can be used to allocate DRAM rows from, such that the newly allocated array can be copied to other arrays allocated using the same *ID*.

(3) Initializer Rows Table. To find which row in a DRAM subarray can be used as the source operand in RowClone-Initialize operations, we maintain the IRT. The IRT is indexed using physical page numbers. RowCopy-Initialize operations query the IRT to obtain the physical address of the DRAM row initialized with zeros and that belongs to the same subarray as the destination operand (i.e., the DRAM row to be initialized with zeros).

Figure 6 describes how `alloc_align()` works over an end-to-end example. Using the RowClone success rate experiment (described above), the custom supervisor software (CSS for short) finds the DRAM rows that are in the same subarray (❶) and initializes the SAMT. The programmer allocates two 128-KiB arrays, A and B, via `alloc_align()` using the same *allocation id* (❷), with the intent to copy from A to B (❸). CSS allocates contiguous ranges of virtual addresses to A and B and then splits the virtual address ranges into page-sized memory blocks (❹). CSS assigns consecutive memory blocks to consecutive DRAM banks and accesses the AIT with the *allocation id* (❺) for each memory block. By accessing the AIT, CSS retrieves the *subarray id* that points to a SAMT entry. The SAMT entry corresponds to the subarray that contains the arrays that are allocated using the *allocation id* (❻). CSS accesses the SAMT entry to retrieve two physical addresses that point to the same DRAM row. CSS maps a memory block and its *sibling memory block* (i.e., the memory block that is $N/2$ bytes away from this memory block, where N is the *size* argument of the `alloc_align()` call) to these two physical addresses, such that they are mapped to the first and the second halves of the same DRAM row (❼). Once allocated, these physical addresses are pinned to main memory and cannot be swapped out to storage. Finally, CSS updates the page table with the physical addresses to map the memory blocks to the same DRAM row (❼).

5.3 Maintaining Memory Coherence

Since memory instructions update the cached copies of data (Section 5.1), a naive implementation of RowClone can potentially operate on stale data, because cached copies of RowClone operands can be modified by CPU store instructions. Thus, we need to ensure memory coherence to prevent RowClone from operating on stale data.

We implement a new custom RISC-V instruction, called *CLFLUSH*, to flush dirty cache blocks to DRAM (RISC-V does not implement any cache management operations [160]) so as to ensure RowClone operates on up-to-date data. A *CLFLUSH* instruction flushes (invalidates) a physically addressed dirty (clean) cache block. *CLFLUSH* or other cache management operations with similar semantics are supported in X86 [68] and ARM architectures [10]. Thus, the *CLFLUSH* instruction

⁴The same RowClone success rate experiment could be conducted in other systems that are based on PiDRAM or in a PiDRAM prototype that uses a different DRAM module. Since the RowClone success rate experiment is a one-time process, its overheads (e.g., time taken to iterate over all DRAM rows using our experiment) are amortized over the lifetime of such a system.

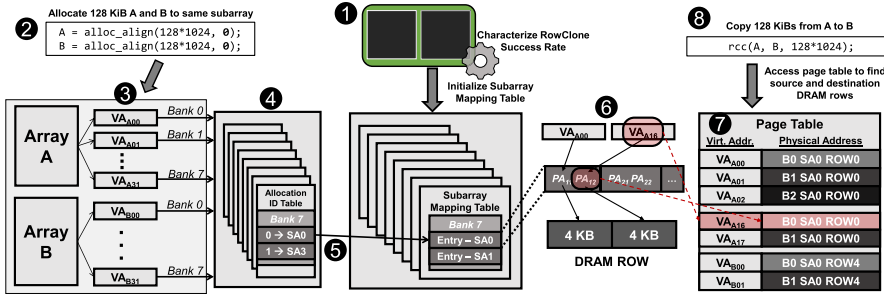


Fig. 6. Alloc_align() and RowClone-Copy (rcc, see Section 5.4) workflow.

(that we implement) provides a minimally invasive solution (i.e., it requires no changes to the specification of commercial ISAs) to the memory coherence problem. Before executing a RowClone Copy or Initialization operation (see Section 5.4), the custom supervisor software flushes (invalidates) the cache blocks of the source (destination) row of the RowClone operation using CLFLUSH.

5.4 RowClone-Copy and RowClone-Initialize

We support the RowClone-Copy and RowClone-Initialize operations in our custom supervisor software via two functions: (i) RowClone-Copy, `rcc(void *dest, void *src, int size)` and (ii) RowClone-Initialize, `rci(void* dest, int size)`. `rcc` copies `size` number of bytes in the virtual address space starting from the `src` memory address to the `dest` memory address. `rci` initializes `size` number of bytes in the virtual address space starting from the `dest` memory address. We expose `rcc` and `rci` to user-level programs using system calls defined in the custom supervisor software.

`rcc` (i) splits the source and destination operands into page-aligned, page-sized blocks, (ii) traverses the page table (Figure 6 ⑦) to find the physical address of each block (i.e., the address of a DRAM row), (iii) flushes all cache blocks corresponding to the source operand and invalidates all cache blocks corresponding to the destination operand, and (iv) performs a RowClone operation from the source row to the destination row using `pumolib's copy_row()` function.

`rci` (i) splits the destination operand into page-aligned, page-sized blocks, (ii) traverses the page table to find the physical address of the destination operand, (iii) queries the IRT (see Section 5.2) to obtain the physical address of the initializer row (i.e., source operand), (iv) invalidates the cache blocks corresponding to the destination operand, and (v) performs a RowClone operation from the initializer row to the destination row using `pumolib's copy_row()` function.

5.5 Evaluation

We evaluate our solutions for the challenges in implementing RowClone end-to-end on a real system using PiDRAM. We modify the custom memory controller to implement DRAM command sequences ($ACT \rightarrow PRE \rightarrow ACT$) to trigger RowClone operations. We set the t_{RAS} and t_{RP} parameters to 10 ns (below the manufacturer-recommended 37.5 ns for t_{RAS} and 13.5 ns for t_{RP} [117]).

5.5.1 Experimental Methodology. Table 3 (left) describes the configuration of the components in our system. We use the pipelined and in-order Rocket core with 16-KiB L1 data cache and 4-entry TLB as the main processor of our system. We use the 1-GiB DDR3 module available on the ZC706 board as the main memory where we conduct PuM operations.

Implementing RowClone requires an additional 198 lines of Verilog code over PiDRAM's existing Verilog design. We add 43 and 522 lines of C code to `pumolib` and to our custom supervisor software, respectively, to implement RowClone.

Table 3. PiDRAM System Configuration (Left)

CPU: 50 MHz; in-order Rocket core [11]; TLB 4 entries DTLB; LRU policy	Physical Address	Physical Page Number		Page Offset	
L1 Data Cache: 16 KiB, 4-way; 64 B line; random replacement policy		29	12	11	0
DRAM Memory: 1 GiB DDR3; 800MT/s; single rank; 8 KiB row size	DRAM Address	Row	Bank	Column	Byte Offset
		29	16 15	13 12	3 2 0

Physical address to DRAM address mapping in PiDRAM (right). Byte offset is used to address the byte in the DRAM burst.

Table 3 (right) describes the mapping scheme we use in our custom memory controller to translate from physical to DRAM row-bank-column addresses. We map physical addresses to DRAM columns, banks, and rows from lower-order bits to higher-order bits to exploit the bank-level parallelism in memory accesses to consecutive physical pages. We note that our memory management mechanism is compatible with other physical address \rightarrow DRAM address mappings [62]. For example, for a mapping scheme where page offset bits (physical address (PA) [11:0]) include all or a subset of the bank address bits, a single RowClone operand (i.e., a 4-KiB page) would be split across multiple DRAM banks. This only coarsens the granularity of RowClone operations as the sibling pages that must be copied in unison, to satisfy the granularity constraint, increases. We expect that for other complex or unknown physical address \rightarrow DRAM address mapping schemes, the characterization of the DRAM device for RowClone success rate would take longer. In the worst case, DRAM row addresses that belong to the same DRAM subarray can be found by testing all combinations of physical addresses for their RowClone success rate.

We evaluate `rcc` and `rci` operations under two configurations to understand the copy/initialization throughput improvements provided by `rcc` and `rci` over CPU-copy operations performed by the Rocket core, and to understand the overheads introduced by end-to-end support for commodity DRAM-based PuM operations. We test two configurations:

(1) Bare-Metal. We assume that RowClone operations always target data that are allocated correctly in DRAM (i.e., there is no overhead introduced by address translation, IRT accesses, and CLFLUSH operations). We directly issue RowClone operations via `pumolib` using physical addresses. CPU-copy operations also use physical addresses.

(2) No Flush. We assume that the programmer uses the `alloc_align` function to allocate the operands of RowClone operations. We use a version of `rcc` and `rci` system calls that do not use CLFLUSH to flush cache blocks of source and destination operands of RowClone operations. We run the *No Flush* configuration on our custom supervisor software; `rcc` and `rci` and traditional CPU-copy operations use virtual addresses.

5.5.2 Workloads. For the two configurations, we run a microbenchmark that consists of two programs, *copy* and *init*, on our prototype. Both programs take the argument N , where *copy* copies an N -byte array to another N -byte array and *init* initializes an N -byte array to all zeros. Both programs have two versions: (i) CPU-copy, which copies/initializes data using memory loads and stores, and (ii) RowClone, which uses RowClone operations to perform copy/initialization. All programs use `alloc_align` to allocate data. The performance results we present in this section are the average of a 1,000 runs. To maintain the same initial system state for both CPU-copy and RowClone, we flush all cache blocks before each one of the 1,000 runs. We run each program for array sizes (N) that are powers of two and $8 \text{ KiB} < N < 8 \text{ MiB}$ and find the average copy/initialization throughput across all 1,000 runs (by measuring the # of elapsed CPU cycles to execute copy/initialization operations) for CPU-copy, RowClone-Copy (`rcc`), and RowClone-Initialize (`rci`).⁵

⁵We tested RowClone operations using `alloc_align()` with up to 8 MiB of allocation size, since we observed diminishing returns on performance improvement provided by RowClone operations on larger array sizes.

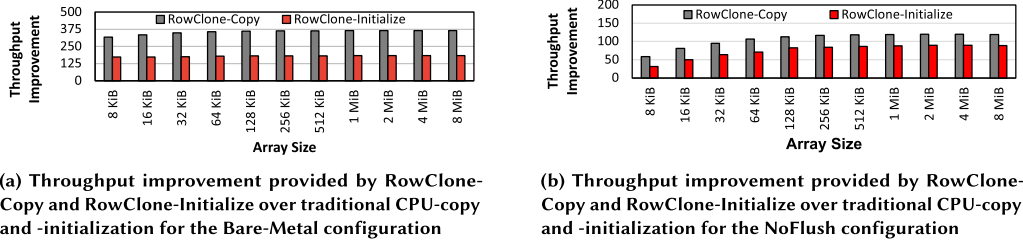


Fig. 7. RowClone-Copy and RowClone-Initialize throughput improvement for the Bare-Metal (left) and the NoFlush (right) configurations.

We analyze the overheads of CLFLUSH operations on copy/initialization throughput that `rcc` and `rci` can provide. We measure the execution time of CLFLUSH operations in our prototype to find how many CPU cycles it takes to flush a (i) dirty and (ii) clean cache block on average across 1,000 measurements. We simulate various scenarios (described in Section 5.5.5) where we assume a certain fraction of the operands of RowClone operations are cached and dirty.

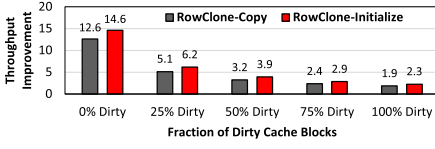
5.5.3 Bare-Metal RowClone. Figure 7(a) shows the throughput improvement provided by `rcc` and `rci` for *copy* and *initialize* over CPU-copy and CPU-initialization for increasing array sizes.

We make two major observations. First, we observe that `rcc` and `rci` provide significant throughput improvement over traditional CPU-copy and CPU-initialization. The throughput improvement provided by `rcc` ranges from $317.5\times$ (for 8 KiB arrays) to $364.8\times$ (for 8 MiB arrays). The throughput improvement provided by `rci` ranges from $172.4\times$ to $182.4\times$. Second, the throughput improvement provided by `rcc` and `rci` increases as the array size increases. This increase saturates when the array size reaches 1 MiB. The load/store instructions used by CPU-copy and CPU-initialization access the operands in a streaming manner. The eviction of dirty cache blocks (i.e., the destination operands of copy and initialization operations) interfere with other memory requests on the memory bus.⁶ We attribute the observed saturation at 1-MiB array size to the interference on the memory bus.

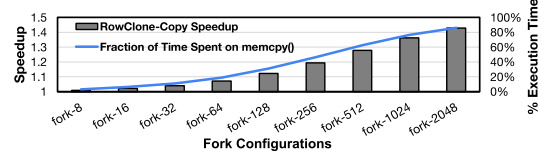
5.5.4 No Flush RowClone. We analyze the overhead in copy/initialization throughput introduced by system support. Figure 7(b) shows the throughput improvement of copy and initialization provided by `rcc` and `rci` operations.

We make two major observations: First, `rcc` improves the copy throughput by $58.3\times$ for 8 KiB and by $118.5\times$ for 8-MiB arrays, whereas `rci` improves initialization throughput by $31.4\times$ for 8 KiB and by $88.7\times$ for 8-MiB arrays. Second, we observe that the throughput improvement provided by `rcc` and `rci` improves *non-linearly* as the array size increases. The execution time (in Rocket core clock cycles) of `rcc` and `rci` operations (not shown in Figure 7(b)) *does not increase linearly* with the array size. For example, the execution time of `rcc` is 397 and 584 cycles at 8-KiB and 16-KiB array sizes, respectively, resulting in a $1.47\times$ increase in execution time between 8-KiB and 16-KiB array sizes. However, the execution time of `rcc` is 92,656 and 187,335 cycles at 4-MiB and 8-MiB array sizes, respectively, resulting in a $2.02\times$ increase in execution time between 4-MiB and 8-MiB array sizes. We make similar observations on the execution time of `rci`. For every RowClone operation, `rcc` and `rci` walk the page table to find the physical addresses corresponding to the source (`rcc`) and the destination (`rcc` and `rci`) operands. We attribute the non-linear increase in `rcc` and `rci`'s execution time to (i) the locality exploited by the Rocket core in accesses to the

⁶Because the data cache in our prototype employs random replacement policy, as the array size increases, the fraction of cache evictions among all memory requests also increases, causing larger interference on the memory bus (i.e., more memory requests to satisfy all cache evictions). The interference saturates at 1-MiB array size.



(a) Throughput improvement provided by rcc and rci with CLFLUSH over Rocket's CPU-copy.



(b) Forkbench speedup (bars, left yaxis) and time spent on memcpy by the CPU baseline (curve, right yaxis)

Fig. 8. Throughput improvement (left) and forkbench speedup (right).

page table and (ii) the diminishing constant cost in the execution time of both rcc and rci due to common instructions executed to perform a system call.

5.5.5 CLFLUSH Overhead. We find that our implementation of CLFLUSH takes 45 Rocket core clock cycles to flush a dirty cache block and 6 Rocket core cycles to invalidate a clean cache block. We estimate the throughput improvement of rcc and rci including the CLFLUSH overhead. We assume that all cache blocks of the source and destination operands are cached and that a fraction of the all cached cache blocks is dirty (quantified on the x axis). We do not include the overhead of accessing the data (e.g., by using *load* instructions) *after* the data gets copied in DRAM. Figure 8(a) shows the estimated improvement in copy and initialization throughput that rcc and rci provide for 8-MiB arrays.

We make three major observations. First, even with inefficient cache flush operations, rcc and rci provide $3.2\times$ and $3.9\times$ higher throughput over the CPU-copy and CPU-initialization operations, assuming 50% of the cache blocks of the 8-MiB source operand are dirty, respectively. Second, as the fraction of dirty cache blocks increases, the throughput improvement provided by both rcc and rci decreases (down to $1.9\times$ for rcc and $2.3\times$ for rci for 100% dirty cache block fraction). Third, we observe that rci can provide better throughput improvement compared to rcc when we include the CLFLUSH overhead. This is because rci flushes cache blocks of one operand (destination), whereas rcc flushes cache blocks of both operands (source and destination).

We do not study the distribution of dirty cache block fractions in real applications as that is not the goal of our CLFLUSH overhead analysis. However, if a large dirty cache block fraction causes severe overhead in a real application, then the system designer or the user of the system would likely decide not to offload the operation to PuM (i.e., performing rcc operations instead of CPU-Copy). PiDRAM's prototype can be useful for studies on different PuM system integration aspects, including such offloading decisions.

We observe that the CLFLUSH operations are inefficient in supporting coherence for RowClone operations. Even so, we see that RowClone-Copy and RowClone-Initialization provides throughput improvements ranging from $1.9\times$ to $14.6\times$. We expect the throughput improvement benefits to increase as coherence between the CPU caches and PIM accelerators become more efficient with new techniques [21, 22, 143].

5.5.6 Real Workload Study. The benefit of rcc and rci on a full application depends on what fraction of execution time is spent on bulk data copy and initialization. We demonstrate the benefit of rcc and rci on *forkbench* [145] and *compile* [145] workloads with varying fractions of time spent on bulk data copy and initialization to show that our infrastructure can enable end-to-end execution and estimation of benefits on real workloads.⁷ We study *forkbench* in detail to demonstrate how the benefits vary with the time spent on data copying in the baseline for this workload.

⁷A full workload study (i.e., with system calls to a full operating system such as Linux) of *forkbench* and *compile* is out of the scope of this article. Our infrastructure currently cannot execute all possible workloads due to the limited library and system call functionality provided by the RISC-V Proxy Kernel [135].

Forkbench first allocates N memory pages and copies data to these pages from a buffer in the process's memory and then accesses 32K random cache blocks within the newly allocated pages to emulate a workload that frequently spawns new processes. We evaluate *forkbench* under varying bulk data copy sizes where we sweep N from 8 to 2,048.

Compile first zero-allocates (`calloc` or `rci`) two pages (8 KiBs) and then executes a number of arithmetic and memory instructions to operate on the zero-allocated data. We carefully develop the *compile* microbenchmark to maintain a realistic ratio between the number of arithmetic and memory instructions executed and zero-allocation function calls made, which we obtain by profiling `gcc` [109]. We use the *No-Flush* configuration of our RowClone implementation for both *forkbench* and *compile*.

Figure 8(b) plots the speedup provided by `rcc` over the CPU-copy baseline, and the proportion of time spent on `memcpy` functions by the CPU-copy baseline, for various configurations of *forkbench* on the x axis.

Forkbench. We observe that RowClone-Copy can significantly improve the performance of *forkbench* by up to 42.9%. RowClone-Copy's performance improvement increases as the number of pages copied increase. This is because the copy operations accelerated by `rcc` contribute a larger amount to the total execution time of the workload. The `memcpy` function calls take 86% of the CPU-copy baseline's time during *forkbench* execution for $N = 2048$.

Compile. RowClone-Initialize improves the performance of *compile* by 9%. Only an estimated 17% of the execution time of *compile* is used for zero-allocation by the CPU-initialization baseline, `rci` reduces the overhead of zero-allocation by (i) performing in-DRAM bulk-initialization and (ii) executing a smaller number of instructions.

Libquantum. To demonstrate that PiDRAM can run real workloads, we run a SPEC2006 [153] workload (libquantum). We modify the `calloc` (allocates and zero initializes memory) function call to allocate data using `alloc_align`, and initialize data using `rci` for allocations that are larger than 8 KiBs.

Using `rci` to bulk initialize data in libquantum improves end-to-end application performance by 1.3% (compared to the baseline that uses CPU-Initialization). This improvement is brought by `rci`, which initializes a total amount of 512 KiBs of memory⁸ using RowClone operations. We note that the proportion of store instructions executed by libquantum to initialize arrays in the CPU-initialization baseline is only 0.2% of all dynamic instructions in the libquantum workload, which amounts to an estimated 2.3% of the total runtime of libquantum. Thus, the 1.3% end-to-end performance improvement provided by `rci` is reasonable, and we expect it to increase with the initialization intensity of workloads.

Summary. We conclude from our evaluation that end-to-end implementations of RowClone (i) can be efficiently supported in real systems by employing memory allocation mechanisms that satisfy the memory *alignment*, *mapping*, *granularity* requirements (Section 5.1) of RowClone operations, (ii) can greatly improve copy/initialization throughput in real systems, and (iii) require cache coherence mechanisms (e.g., PIM-optimized coherence management [21, 22, 143]) that can flush dirty cache blocks of RowClone operands efficiently to achieve optimal copy/initialization throughput improvement. PiDRAM can be used to estimate end-to-end workload execution benefits provided by RowClone operations. Our experiments using libquantum, forkbench, and compile show that (i) PiDRAM can run real workloads, (ii) our end-to-end implementation of RowClone

⁸In libquantum, there are 16 calls to `calloc` that exceed the 8-KiB allocation size. We only bulk initialize data using `rci` for these 16 calls.

operates correctly, and (iii) RowClone can improve the performance of real workloads in a real system, even when inefficient CLFLUSH operations are used to maintain memory coherence.

6 CASE STUDY #2: END-TO-END D-RANGE

Prior work on DRAM-based random number generation techniques [13, 83, 123] do not integrate and evaluate their techniques end-to-end in a real system. We evaluate one DRAM-based true random number generation technique, D-RaNGe [83], end-to-end using PiDRAM. We implement support for D-RaNGe in PiDRAM by enabling access to DRAM with reduced activation latency (i.e., $tRCD$ set to values lower than manufacturer recommendations).

6.1 D-RaNGe Implementation

We implement a simple version of D-RaNGe in PiDRAM. PiDRAM's D-RaNGe controller collects true random numbers from four DRAM cells in the same DRAM cache block inside one DRAM bank. We implement the D-RaNGe controller within the Periodic Operations Module (Section 4.1). The D-RaNGe controller (i) periodically accesses a DRAM cache block with reduced $tRCD$, (ii) reads four of the TRNG DRAM cells in the cache block, (iii) stores the four bits read from the TRNG cells in a 1 KiB random number buffer. We reserve multiple configuration registers in the CRF to configure (i) the TRNG period (in nanoseconds) used by the D-RaNGe controller to periodically generate random numbers by accessing DRAM with reduced activation latency while the buffer is not full (the D-RaNGe controller accesses DRAM every TRNG period), (ii) the timing parameter ($tRCD$) used to induce activation latency failures, and (iii) the physical location (DRAM bank, row, column addresses, and bit offset within the DRAM column) of the TRNG cells to read. We implement two pumolib functions: (i) `buf_size()`, which returns the number of random words (4 bytes) available in the buffer, and (ii) `rand_dram()`, which returns one random word that is read from the buffer. The two functions first execute PiDRAM instructions in the POC that update the data register either with (i) the number of random words available (when `buf_size()` is called) or (ii) a random word read from the random number buffer (when `rand_dram()` is called). The two functions then access the data register using LOAD instructions to retrieve either the size of the random number buffer or a random number. The application developer reads true random numbers using these two functions in pumolib.

Random Cell Characterization. D-RaNGe requires the system designer to characterize the DRAM module for activation latency failures to find DRAM cells that fail with a 50% probability (i.e., randomly) when accessed with reduced $tRCD$. Following the methodology presented in Reference [83], the system designer can characterize a DRAM device or use an automated procedure to find cells that fail with a 50% probability. In PiDRAM, we implement reduced latency access to DRAM by (i) extending the scheduler of the custom memory controller and (ii) adding a pumolib function `activation_failure(address)` that induces a $tRCD$ failure on the DRAM cache block at address.

6.2 Evaluation and Results

Experimental Methodology. We run a microbenchmark to understand the effect of the TRNG period on true random number generation throughput observed by a program running on the Rocket core. The microbenchmark consists of a loop that (i) checks the availability of random numbers using `buf_size()` and (ii) reads a random number from the buffer using `rand_dram()`. We execute the microbenchmark until we read one million bytes of random numbers.

Results. The D-RaNGe controller can perform reduced-latency accesses frequently, every 220 ns. Figure 9 depicts the TRNG throughput observed by the microbenchmark for TRNG periods in the

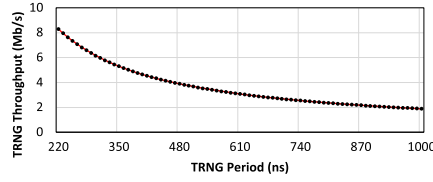


Fig. 9. TRNG throughput observed by our microbenchmark for TRNG periods ranging from 220 to 1,000 ns.

range [220 ns, 1,000 ns] with increments of 10 ns. We observe that the TRNG throughput decreases from 8.30 Mb/s at 220 ns TRNG period to 1.90 Mb/s at 1,000 ns TRNG period. D-RaNGe [83] reports 25.2 Mb/s TRNG throughput using a single DRAM bank when there are four random cells in a cache block. PiDRAM's D-RaNGe controller can be optimized to generate random numbers more frequently to match D-RaNGe's observed maximum throughput.⁹ We leave such optimizations to PiDRAM's D-RaNGe controller for future work.

Including the modifications to the custom memory controller and pumolib, implementing D-RaNGe and reduced-latency DRAM access requires an additional 190 lines of Verilog and 74 lines of C code over PiDRAM's existing codebase. We conclude that our D-RaNGe implementation (i) provides a basis for PiDRAM developers to study end-to-end implementations of DRAM-based true random number generators and (ii) shows that PiDRAM's hardware and software components facilitate the implementation of new commodity DRAM-based PuM techniques, specifically those that are related to security. Our reduced-latency DRAM access implementation provides a basis for other PuM techniques for security purposes, such as the DRAM-latency physical unclonable functions (DL-PUF [82]) and QUAC-TRNG [124] (Section 4.4). We leave further exploration on end-to-end implementations of D-RaNGe, DL-PUF, and QUAC-TRNG, as well as end-to-end analyses of the security benefits they provide using PiDRAM for future work.

7 EXTENDING PIDRAM

We briefly describe the modifications required to extend PiDRAM (i) with new DRAM commands and DRAM timing parameters, (ii) with new case studies, and (iii) to support new FPGA boards.

New DRAM Commands and Timing Parameters. Implementing new DRAM commands or modifying DRAM timing parameters require modifications to PiDRAM's memory controller. This is straightforward as PiDRAM's memory controller's Verilog design is modular and uses well-defined interfaces: It is composed of multiple modules that perform separate tasks. For example, the memory request scheduler comprises two main components: (1) *command timer* and (2) *command scheduler*. To serve LOAD and STORE memory requests, the command scheduler maintains state (e.g., which row is active) for every bank. The command scheduler selects the next DRAM command to satisfy the LOAD or STORE memory request and queries the command timer with the selected DRAM command. The command timer checks for all possible standard DRAM timing constraints and outputs a valid bit if the selected command can be issued in that FPGA clock cycle. To extend the memory controller with a new standard DRAM command (e.g., to implement a newer standard like DDR4 or DDR5), a PiDRAM developer simply needs to (i) add a new timing

⁹D-RaNGe has a smaller true random number generation (TRNG) latency (i.e., takes a smaller amount of time to generate a 4-bit random number) than PiDRAM. PiDRAM has a larger TRNG latency due to (i) discrepancies in the data path (i.e., on-chip interconnect) in D-RaNGe's simulated system and PiDRAM's prototype and (ii) the TRNG period of the D-RaNGe controller (D-RaNGe controller performs a reduced t_{RCD} access only as frequently as one every 220 ns). The D-RaNGe controller can be optimized further to further reduce the TRNG period by down the DRAM row cycle time (t_{RC} standard timing parameter, typically 45 ns [117]).

Table 4. Comparison of PiDRAM with Related State-of-the-art Prototyping and Evaluation Platforms

Platforms	Interface with real DRAM chips	Flexible MC for PuM	System software support	Open-source
Silent-PIM [78]	✗	✗	✓	✗
SoftMC [60]	✓ (DDR3)	✗	✗	✓
ComputeDRAM [44]	✓ (DDR3)	✗	✗	✗
MEG [173]	✓ (HBM)	✗	✓	✓
PiMulator [119]	✗	✓	✗	✓
Commercial platforms (e.g., ZYNQ [165])	✓ (DDR3/4)	✗	✓	✗
Simulators [18, 35, 90, 132, 139, 168, 169, 174]	✗	✓	✓ (potentially)	✓
PiDRAM (this work)	✓ (DDR3)	✓	✓	✓

constraint by replicating the logic in the command timer and (ii) extend the command scheduler to correctly maintain the bank state.

New Case Studies. Implementing new techniques (e.g., those that are listed in Table 2) to perform new case studies requires modifications to PiDRAM’s hardware and software components. We describe the required modifications over an example ComputeDRAM-based in-DRAM bitwise operations case study.

To implement ComputeDRAM-based in-DRAM bitwise operations, the developers need to (i) extend the *custom command scheduler* in PiDRAM’s memory controller with a new state machine that schedules new DRAM command sequences (ACT-PRE-ACT) with an appropriate set of violated timing parameters (our ComputeDRAM-based in-DRAM copy implementation provides a solid basis for this), (ii) expose the functionality to the processor by implementing new PiDRAM instructions in the PuM controller (e.g., by replicating and customizing the existing logic for decoding and executing RowClone operations), and (iii) and make modifications to the software library to expose the new instruction to the programmer (e.g., by replicating the copy_row function’s behavior, described in Table 1).

Porting to New FPGA Boards. Developing new PiDRAM prototypes on different FPGA boards could require modifications to design constraints (e.g., top level input/outputs to physical FPGA pins) and the DDRx PHY IP depending on the FPGA board. Modifying design constraints is a straightforward task involving looking up the FPGA manufacturer datasheets and modifying design constraint files [164]. Manufacturers may provide different DDRx PHY IPs for different FPGAs. Fortunately, these IPs typically expose similar (based on the DFI standard [33]) interfaces to user hardware (in our case, to PiDRAM’s memory controller). Thus, other PiDRAM prototypes on different FPGA boards can be developed with small yet careful modifications to the ZC706 prototype design we provide.

8 RELATED WORK

To our knowledge, this is the first work to develop a flexible, open source framework that enables integration and evaluation of commodity DRAM-based PuM techniques on real DRAM chips by providing the necessary hardware and software components. We demonstrate the first end-to-end implementation of RowClone and D-RaNGe using real DRAM chips. We compare the features of PiDRAM with other state-of-the-art prototyping and evaluation platforms in Table 4 and discuss them below. The four features we use for comparison are as follows: (1) *Interface with real DRAM chips*: The platform allows running experiments using real DRAM chips. (2) *Flexible memory controller for PuM*: The platform provides a flexible memory controller that can easily be extended to perform (e.g., as in PiDRAM) or emulate (e.g., as in PiMulator [119]) new PuM operations. (3) *System software support*: The platform provides support for running system software such as operating systems or supervisor software (e.g., RISC-V PK [135]). (4) *Open source*: The platform is available as open source software.

Silent-PIM [78]. Silent-PIM proposes a new DRAM design that incorporates processing units capable of vector arithmetic computation. Silent-PIM’s goal is to evaluate PIM techniques on a

new, PIM-capable DRAM device using standard DRAM commands (e.g., as defined in DDR4 [71]); it does not provide an evaluation platform or prototype. In contrast, PiDRAM is designed for rapid integration and evaluation of PuM techniques that use *real DRAM devices*. PiDRAM provides key hardware and software components that facilitate end-to-end implementations of PuM techniques.

SoftMC [52, 60]. SoftMC is an FPGA-based DRAM testing infrastructure. SoftMC can issue arbitrary sequences of DDR3 commands to real DRAM devices. SoftMC is widely used in prior work that studies the performance, reliability and security of real DRAM chips [13, 14, 28, 38, 41, 50, 59, 77, 83, 85, 96, 127, 154]. SoftMC is built to test DRAM devices, *not* to study end-to-end implementations of PuM techniques. Thus, SoftMC (i) does *not* support application execution on a real system and (ii) *cannot* use DRAM modules as main memory. While SoftMC is useful in studies that perform exhaustive search on all possible sequences of DRAM commands to potentially uncover undocumented DRAM behavior (e.g., ComputeDRAM [44], QUAC-TRNG [123]), PiDRAM is developed to study end-to-end implementations of PuM techniques. PiDRAM provides an FPGA-based prototype that comprises a RISC-V system and supports using DRAM modules both for storing data (i.e., as main memory) and performing PuM computation.

ComputeDRAM [44]. ComputeDRAM partially demonstrates that two DRAM-based state-of-the-art PuM techniques, RowClone [145] and Ambit [147], are already possible on real off-the-shelf DDR3 chips. ComputeDRAM uses SoftMC to demonstrate in-DRAM copy and bitwise AND/OR operations on real DDR3 chips. ComputeDRAM's goal is *not* to develop a framework to facilitate end-to-end implementations of PuM techniques. Therefore, it does *not* provide (i) a flexible memory controller for PuM or (ii) support for system software. PiDRAM provides the necessary software and hardware components to facilitate end-to-end implementations of PuM techniques.

MEG [173]. MEG is an open source system emulation platform for enabling FPGA-based operation interfacing with **High-Bandwidth Memory (HBM)**. MEG aims to efficiently retrieve data from HBM and perform the computation in the host processor implemented as a soft core on the FPGA. Unlike PiDRAM, MEG does *not* implement a flexible memory controller that is capable of performing PuM operations. We demonstrate the flexibility of PiDRAM by implementing two state-of-the-art PuM techniques [83, 145]. We believe MEG and PiDRAM can be combined to get the functionality and prototyping power of both works.

PiMulator [119]. PiMulator is an open source PiM emulation platform. PiMulator implements a main memory and a PiM model using SystemVerilog, allowing FPGA emulation of PiM architectures. PiMulator enables easy emulation of new PiM techniques. However, it does *not* allow end-to-end execution of workloads that use PiM techniques and it does not provide the user with full control over the DRAM interface.

Commercial Platforms (e.g., ZYNQ [165]). Some commercial platforms implement CPU-FPGA heterogeneous computing systems. A memory controller is provided to access DRAM as the main memory in such systems. However, in such systems, (i) there is *no* support for PuM mechanisms and (ii) the entire hardware-software stack is closed source. PiDRAM can be integrated into these systems, using the closed source computing system as the main processor. Our prototype utilizes an open source system-on-chip (Rocket Chip [11]) as the main processor, which enables developers to study architectural and microarchitectural aspects of PuM techniques (e.g., data allocation and coherence mechanisms). Such studies cannot be conducted using closed source computing systems.

Simulators. Many prior works propose full-system (e.g., References [18, 132]), trace-based (e.g., References [64, 90, 139, 168, 169, 174]), and instrumentation-based (e.g., References [35, 64, 168]) simulators that can be used to evaluate PuM techniques. Although useful, these simulators do

not model DRAM behavior and cannot integrate proprietary device characteristics (e.g., DRAM internal address mapping) into their simulations, without conducting a rigorous characterization study. Moreover, the effects of environmental conditions (e.g., temperature and voltage) on DRAM chips are unlikely to be modeled on accurate, full-system simulators as it would require excessive computation, which would negatively impact the already poor performance (200K instructions per second) of full system simulators [140]. In contrast, PiDRAM interfaces with real DRAM devices and its prototype achieves a 50-MHz clock speed (and can be improved further), which lets PiDRAM execute >10M instructions per second (assuming <5 cycles per instruction). PiDRAM can be used to study end-to-end implementations of PuM techniques and explore solutions that take into account the effects related to the environmental conditions of real DRAM devices. Future versions of PiDRAM could be easily extended (e.g., with real hardware that allows controlling DRAM temperature and voltage [115, 156]) to experiment with different DRAM temperature and voltage levels to better understand the effects of these environmental conditions on the reliability of PuM operations. Using PiDRAM, experiments that require executing real workloads can take an order of magnitude shorter wall clock time compared to using full-system simulators.

Other Related Work. Prior works (see Section 2.2) (i) propose or (ii) demonstrate using real DRAM chips, several DRAM-based PuM techniques that can perform computation [6, 28, 40, 54, 144, 146, 147, 149, 150], move data [145, 159], or implement security primitives [13, 14, 82, 83, 124, 126] in memory. SIMDRAM [54] develops a framework that provides a programming interface to perform in-DRAM computation using the majority operation. DR-STRANGE [23] proposes an end-to-end system design for DRAM-based true random number generators. None of these works provide an end-to-end in-DRAM computation framework that is integrated into a real system using real DRAM chips. We conclude that existing platforms cannot substitute PiDRAM in studying commodity DRAM-based PuM techniques.

9 CONCLUSION

We develop PiDRAM, a flexible and open source prototyping framework for integrating and evaluating end-to-end commodity DRAM-based PuM techniques. PiDRAM comprises the necessary hardware and software structures to facilitate end-to-end implementation of PuM techniques. We build an FPGA-based prototype of PiDRAM along with an open source RISC-V system and enable computation on real DRAM chips. Using PiDRAM, we implement and evaluate RowClone (in-DRAM data copy and initialization) and D-RaNGe end-to-end in the entire real system. Our results show that RowClone significantly improves data copy and initialization throughput in a real system on real workloads, and efficient cache coherence mechanisms are needed to maximize RowClone's potential benefits. Our implementation of D-RaNGe requires small additions to PiDRAM's codebase and provides true random numbers at high throughput and with low latency. We conclude that unlike existing prototyping and evaluation platforms, PiDRAM enables (i) easy integration of existing and new PuM techniques end-to-end in a real system and (ii) novel studies on end-to-end implementations of PuM techniques using real DRAM chips. PiDRAM is freely available as an open source tool for researchers and designers in both academia and industry to experiment with and build on.

ACKNOWLEDGMENTS

We thank the reviewers of MICRO 2021, HPCA 2022, and TACO for feedback. We thank the SAFARI Research Group members for valuable feedback and the stimulating intellectual environment they provide. We acknowledge the generous gifts provided by our industrial partners, including Google, Huawei, Intel, Microsoft, and VMware. This research was also supplied in part by the Semiconductor Research Corporation and the ETH Future Computing Laboratory.

REFERENCES

- [1] Shaizeen Aga, Supreet Jeloka, Arun Subramaniyan, Satish Narayanasamy, David Blaauw, and Reetuparna Das. 2017. Compute caches. In *HPCA*.
- [2] Junwhan Ahn, Sungpack Hong, Sungjoo Yoo, Onur Mutlu, and Kiyoun Choi. 2015. A scalable processing-in-memory accelerator for parallel graph processing. In *ISCA*.
- [3] Junwhan Ahn, Sungjoo Yoo, Onur Mutlu, and Kiyoun Choi. 2015. PIM-enabled instructions: A low-overhead, locality-aware processing-in-memory architecture. In *ISCA*.
- [4] Berkin Akin, Franz Franchetti, and James C. Hoe. 2015. Data reorganization in memory using 3D-stacked DRAM. In *ISCA*.
- [5] Mustafa F. Ali, Akhilesh Jaiswal, and Kaushik Roy. 2019. In-memory low-cost bit-serial addition using commodity DRAM technology. *IEEE Transactions on Circuits and Systems I: Regular Papers* 67, 1 (2020), 155–165. <https://doi.org/10.1109/TCSL.2019.2945617>
- [6] Shaahin Angizi and Deliang Fan. 2019. Graphide: A graph processing accelerator leveraging in-dram-computing. In *GLSVLSI*.
- [7] S. Angizi, Z. He, and D. Fan. 2018. PIMA-logic: A novel processing-in-memory architecture for highly flexible and energy-efficient logic computation. In *DAC*.
- [8] S. Angizi, A. S. Rakin, and D. Fan. 2018. CMP-PIM: An energy-efficient comparator-based processing-in-memory neural network accelerator. In *DAC*.
- [9] S. Angizi, J. Sun, W. Zhang, and D. Fan. 2019. AlignS: A processing-in-memory accelerator for DNA short read alignment leveraging SOT-MRAM. In *DAC*.
- [10] ARM. 2021. Cache Maintenance Operations. Retrieved from <https://developer.arm.com/documentation/ddi0246/h/programmers-model/register-descriptions/cache-maintenance-operations>.
- [11] Krste Asanović, Rimas Avizienis, Jonathan Bachrach, Scott Beamer, David Biancolin, Christopher Celio, Henry Cook, Palmer Dabbelt, John R. Hauser, Adam M. Izraelevitz, Sagar Karandikar, Benjamin Keller, Donggyu Kim, John Koenig, Yunsup Lee, Eric Love, Martin Maas, Albert Magyar, Howard Mao, Miquel Moretó, Albert Ou, David A. Patterson, B. H. Richards, Colin Schmidt, Stephen M. Twigg, Huy Vo, and Andrew Waterman. 2016. The rocket chip generator. Technical Report No. UCB/EECS-2016-17.
- [12] Hadi Asghari-Moghaddam, Young Hoon Son, Jung Ho Ahn, and Nam Sung Kim. 2016. Chameleon: Versatile and practical near-DRAM acceleration architecture for large memory systems. In *MICRO*.
- [13] B. M. S. Bahar Talukder, J. Kerns, B. Ray, T. Morris, and M. T. Rahman. 2019. Exploiting DRAM latency variations for generating true random numbers. In *ICCE*.
- [14] B. M. S. Bahar Talukder, Biswajit Ray, Domenic Forte, and Md Tauhidur Rahman. 2019. PreLatPUF: Exploiting DRAM latency variations for generating robust device signatures. *IEEE Access* 7 (2019), 81106–81120. <https://doi.org/10.1109/ACCESS.2019.2923174>
- [15] Alessandro Barengi, Luca Breveglieri, Niccolò Izzo, and Gerardo Pelosi. 2018. Software-only reverse engineering of physical DRAM mappings for rowhammer attacks. In *IVSW*.
- [16] Maciej Besta, Raghavendra Kanakagiri, Grzegorz Kwasniewski, Rachata Ausavarungnirun, Jakub Beránek, Konstantinos Kanellopoulos, Kacper Janda, Zur Vonarburg-Shmaria, Lukas Gianinazzi, Ioana Stefan, Juan Gómez Luna, Jakub Golinowski, Marcin Copik, Lukas Kapp-Schwoerer, Salvatore Di Girolamo, Nils Blach, Marek Konieczny, Onur Mutlu, and Torsten Hoefer. 2021. SISA: Set-centric instruction set architecture for graph mining on processing-in-memory systems. In *MICRO*.
- [17] D. Bhattacharjee, R. Devadoss, and A. Chattopadhyay. 2017. ReVAMP: ReRAM based VLIW architecture for in-memory computing. In *DATE*.
- [18] N. Binkert, B. Beckman, A. Saidi, G. Black, and A. Basu. 2011. The gem5 simulator. *SIGARCH Comput. Archit. News* 39, 2 (aug 2011), 1–7. <https://doi.org/10.1145/2024716.2024718>
- [19] Julien Borghetti, Gregory Snider, Philip Kuekes, Jianhua Joshua Yang, Duncan Stewart, and Stan Williams. 2010. Memristive switches enable stateful logic operations via material implication. *Nature* 464 (04 2010), 873–6. <https://doi.org/10.1038/nature08940>
- [20] A. Boroumand, S. Ghose, Y. Kim, R. Ausavarungnirun, E. Shiu, R. Thakur, D. Kim, A. Kuusela, A. Knies, P. Ranganathan, and O. Mutlu. 2018. Google workloads for consumer devices: Mitigating data movement bottlenecks. In *ASPLOS*.
- [21] Amirali Boroumand, Saugata Ghose, Minesh Patel, Hasan Hassan, Brandon Lucia, Rachata Ausavarungnirun, Kevin Hsieh, Nastaran Hajinazar, Krishna T. Malladi, Hongzhong Zheng, and Onur Mutlu. 2019. CoNDA: Efficient cache coherence support for near-data accelerators. In *ISCA*.
- [22] Amirali Boroumand, Saugata Ghose, Minesh Patel, Hasan Hassan, Brandon Lucia, Kevin Hsieh, Krishna T. Malladi, Hongzhong Zheng, and Onur Mutlu. 2016. LazyPIM: An efficient cache coherence mechanism for processing-in-memory. *IEEE Computer Architecture Letters* 16, 1 (2017), 46–50. <https://doi.org/10.1109/LCA.2016.2577557>

- [23] F. Bostanci, A. Olgun, L. Orosa, A. Yaglikci, J. S. Kim, H. Hassan, O. Ergin, and O. Mutlu. 2022. DR-STraNGe: End-to-end system design for DRAM-based true random number generators. In *HPCA*.
- [24] Geoffrey W. Burr, Robert M. Shelby, Abu Sebastian, Sangbum Kim, Seyoung Kim, Severin Sidler, Kumar Virwani, Masatoshi Ishii, Pritish Narayanan, Alessandro Fumarola, Lucas L. Sanches, Irem Boybat, Manuel Le Gallo, Kibong Moon, Jiyoo Woo, Hyunsang Hwang, and Yusuf Leblebici. 2017. Neuromorphic computing using non-volatile memory. *Advances in Physics: X* 2, 1 (2017), 89–124. <https://doi.org/10.1080/23746149.2016.1259585> arXiv:<https://doi.org/10.1080/23746149.2016.1259585>
- [25] Damla Senol Cali, Gurpreet S. Kalsi, Zülal Bingöl, Can Firtina, Lavanya Subramanian, Jeremie S. Kim, Rachata Ausavarungrinun, Mohammed Alser, Juan Gomez-Luna, Amirali Boroumand, Anant Nori, Allison Scibisz, Sreenivas Subramoney, Can Alkan, Saugata Ghose, and Onur Mutlu. 2020. GenASM: A high-performance, low-power approximate string matching acceleration framework for genome sequence analysis. In *MICRO*.
- [26] K. Chang. 2017. *Understanding and Improving the Latency of DRAM-based Memory Systems*. Ph.D. Dissertation. Carnegie Mellon University.
- [27] Kevin K. Chang, Abhijith Kashyap, Hasan Hassan, Saugata Ghose, Kevin Hsieh, Donghyuk Lee, Tianshi Li, Gennady Pekhimenko, Samira Khan, and Onur Mutlu. 2016. Understanding latency variation in modern DRAM chips: Experimental characterization, analysis, and optimization. In *SIGMETRICS*.
- [28] Kevin K. Chang, Prashant J. Nair, Donghyuk Lee, Saugata Ghose, Moinuddin K. Qureshi, and Onur Mutlu. 2016. Low-cost inter-linked subarrays (LISA): Enabling fast inter-subarray data movement in DRAM. In *HPCA*.
- [29] Kevin K. Chang, Abdullah Giray Yağlıkçı, Saugata Ghose, Aditya Agrawal, Niladrish Chatterjee, Abhijith Kashyap, Donghyuk Lee, Mike O'Connor, Hasan Hassan, and Onur Mutlu. 2017. Understanding reduced-voltage operation in modern DRAM devices: Experimental characterization, analysis, and mechanisms. In *SIGMETRICS*.
- [30] Lucian Cojocar, Jeremie Kim, Minesh Patel, Lillian Tsai, Stefan Saroiu, Alec Wolman, and Onur Mutlu. 2020. Are we susceptible to rowhammer? An end-to-end methodology for cloud providers. In *S&P*.
- [31] Guohao Dai, Tianhao Huang, Yuze Chi, Jishen Zhao, Guangyu Sun, Yongpan Liu, Yu Wang, Yuan Xie, and Huazhong Yang. 2018. GraphH: A processing-in-memory architecture for large-scale graph processing. *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems* 38, 4 (2019), 640–653. <https://doi.org/10.1109/TCAD.2018.2821565>
- [32] Q. Deng, L. Jiang, Y. Zhang, M. Zhang, and J. Yang. 2018. DrAcc: A DRAM based accelerator for accurate CNN inference. In *DAC*.
- [33] DFI Group. 2018. DFI 5.0 Specification. <https://www.ddd-phy.org/>.
- [34] Mario Paulo Drumond Lages De Oliveira, Alexandros Daglis, Nooshin Mirzadeh, Dmitrii Ustiugov, Javier Picorel Obando, Babak Falsafi, Boris Grot, and Dionisios Pnevmatikatos. 2017. The mondrian data engine. In *ISCA*.
- [35] Bruno E. Forlin, Paulo C. Santos, Augusto E. Becker, Marco A. Z. Alves, and Luigi Carro. 2022. Sim2PIM: A complete simulation framework for processing-in-memory. *J. Syst. Archit.* 128, C (Jul 2022), 11 pages. <https://doi.org/10.1016/j.sysarc.2022.102528>
- [36] Charles Eckert, Xiaowei Wang, Jingcheng Wang, Arun Subramaniyan, Ravi Iyer, Dennis Sylvester, David Blaauw, and Reetuparna Das. 2018. Neural cache: Bit-serial in-cache acceleration of deep neural networks. In *ISCA*.
- [37] Amin Farmahini-Farahani, Jung Ho Ahn, Katherine Morrow, and Nam Sung Kim. 2015. NDA: Near-DRAM acceleration architecture leveraging commodity DRAM devices and standard memory modules. In *HPCA*.
- [38] Mohammad Farmani, Mark Tehranipoor, and Fahim Rahman. 2021. RHAT: Efficient rowhammer-aware test for modern DRAM modules. In *ETS*.
- [39] Ivan Fernandez, Ricardo Quisilant, Christina Giannoula, Mohammed Alser, Juan Gomez-Luna, Eladio Gutierrez, Oscar Plata, and Onur Mutlu. 2020. NATSA: A near-data processing accelerator for time series analysis. In *ICCD*.
- [40] João Dinis Ferreira, Gabriel Falcao, Juan Gómez-Luna, Mohammed Alser, Lois Orosa, Mohammad Sadrosadati, Jeremie S. Kim, Geraldo F. Oliveira, Taha Shahroodi, Anant Nori, et al. 2021. pLUTo: In-DRAM lookup tables to enable massively parallel general-purpose computation. arXiv:2104.07699. Retrieved from <https://arxiv.org/abs/2104.07699>.
- [41] Pietro Frigo, Emanuele Vannacci, Hasan Hassan, Victor van der Veen, Onur Mutlu, Cristiano Giuffrida, Herbert Bos, and Kaveh Razavi. 2020. TRRespass: Exploiting the many sides of target row refresh. In *S&P*.
- [42] Daichi Fujiki, Scott Mahlke, and Reetuparna Das. 2019. Duality cache for data parallel acceleration. In *ISCA*.
- [43] Pierre-Emmanuel Gaillardon, Luca Amarú, Anne Siemon, Eike Linn, Rainer Waser, Anupam Chattopadhyay, and Giovanni De Micheli. 2016. The programmable logic-in-memory (PLiM) computer. In *DATe*.
- [44] Fei Gao, Georgios Tziantzioulis, and David Wentzlaff. 2019. ComputeDRAM: In-memory compute using off-the-shelf DRAMs. In *MICRO*.
- [45] Mingyu Gao, Grant Ayers, and Christos Kozyrakis. 2015. Practical near-data processing for in-memory analytics frameworks. In *PACT*.
- [46] Mingyu Gao and Christos Kozyrakis. 2016. HRL: Efficient and flexible reconfigurable logic for near-data processing. In *HPCA*.

- [47] Mingyu Gao, Jing Pu, Xuan Yang, Mark Horowitz, and Christos Kozyrakis. 2017. Tetris: Scalable and efficient neural network acceleration with 3D memory. In *ASPLOS*.
- [48] Saugata Ghose, Amirali Boroumand, Jeremie S. Kim, Juan Gómez-Luna, and Onur Mutlu. 2019. Processing-in-memory: A workload-driven perspective. *IBM Journal of Research and Development* 63, 6 (2019), 3:1–3:19. <https://doi.org/10.1147/JRD.2019.2934048>
- [49] Saugata Ghose, Tianshi Li, Nastaran Hajinazar, Damla Senol Cali, and Onur Mutlu. 2019. Demystifying complex workload-DRAM interactions: An experimental study. In *SIGMETRICS*.
- [50] Saugata Ghose, Abdullah Giray Yaglikçi, Raghav Gupta, Donghyuk Lee, Kais Kudrolli, William X. Liu, Hasan Hassan, Kevin K. Chang, Niladri Chatterjee, Aditya Agrawal, Mike O'Connor, and Onur Mutlu. 2018. What your DRAM power models are not telling you: Lessons from a detailed experimental study. In *SIGMETRICS*.
- [51] Christina Giannoula, Nandita Vijaykumar, Nikela Papadopoulou, Vasileios Karakostas, Ivan Fernandez, Juan Gómez-Luna, Lois Orosa, Nectarios Koziris, Georgios Goumas, and Onur Mutlu. 2021. SynCron: Efficient synchronization support for near-data-processing architectures. In *HPCA*.
- [52] SAFARI Research Group. 2021. SoftMC v1.0—GitHub Repository. Retrieved from <https://github.com/CMU-SAFARI/SoftMC>.
- [53] Boncheol Gu, A. S. Yoon, D.-H. Bae, I. Jo, J. Lee, J. Yoon, J.-U. Kang, M. Kwon, C. Yoon, S. Cho, J. Jeong, and D. Chang. 2016. Biscuit: A framework for near-data processing of big data workloads. In *ISCA*.
- [54] Nastaran Hajinazar, Geraldo F. Oliveira, Sven Gregorio, João Dinis Ferreira, Nika Mansouri Ghiasi, Minesh Patel, Mohammed Alser, Saugata Ghose, Juan Gómez-Luna, and Onur Mutlu. 2021. SIMDRAM: A framework for bit-serial SIMD processing using DRAM. In *ASPLOS*.
- [55] S. Hamdioui, S. Kvatinsky, and et al. G. Cauwenberghs. 2017. Memristor for computing: Myth or reality? In *DATE*.
- [56] Said Hamdioui, Lei Xie, Hoang Anh Du Nguyen, Mottaqiallah Taouil, Koen Bertels, Henk Corporaal, Hailong Jiao, Francky Catthoor, Dirk Wouters, Linn Eike, and Jan van Lunteren. 2015. Memristor based computation-in-memory architecture for data-intensive applications. In *DATE*.
- [57] Milad Hashemi, Khubaib, Eiman Ebrahimi, Onur Mutlu, and Yale N. Patt. 2016. Accelerating dependent cache misses with an enhanced memory controller. In *ISCA*.
- [58] M. Hashemi, O. Mutlu, and Y. N. Patt. 2016. Continuous runahead: Transparent hardware acceleration for memory intensive workloads. In *MICRO*.
- [59] Hasan Hassan, Yahya Can Tugrul, Jeremie S. Kim, Victor van der Veen, Kaveh Razavi, and Onur Mutlu. 2021. Uncovering In-DRAM rowhammer protection mechanisms: A new methodology, custom rowhammer patterns, and implications. In *MICRO*.
- [60] Hasan Hassan, Nandita Vijaykumar, Samira Khan, Saugata Ghose, Kevin Chang, Gennady Pekhimenko, Donghyuk Lee, Oguz Ergin, and Onur Mutlu. 2017. SoftMC: A flexible and practical open-source infrastructure for enabling experimental DRAM studies. In *HPCA*.
- [61] C. Helm, S. Akiyama, and K. Taura. 2020. Reliable reverse engineering of intel DRAM addressing using performance counters. In *MASCOTS*.
- [62] Marius Hillenbrand. 2017. Physical Address Decoding in Intel Xeon v3/v4 CPUs: A Supplemental Datasheet.
- [63] M. Horiguchi. 1997. Redundancy techniques for high-density DRAMs. In *ISIS*.
- [64] HPS Research Group. 2022. Scarab—Github Repository. Retrieved from <https://github.com/hpsresearchgroup/scarab>.
- [65] Kevin Hsieh, Eiman Ebrahimi, Gwangsun Kim, Niladri Chatterjee, Mike O'Conner, Nandita Vijaykumar, Onur Mutlu, and Stephen Keckler. 2016. Transparent offloading and mapping (TOM): Enabling programmer-transparent near-data processing in GPU systems. In *ISCA*.
- [66] K. Hsieh, S. Khan, N. Vijaykumar, K. K. Chang, A. Boroumand, S. Ghose, and O. Mutlu. 2016. Accelerating pointer chasing in 3D-stacked memory: Challenges, mechanisms, evaluation. In *ICCD*.
- [67] Yu Huang, Long Zheng, Pengcheng Yao, Jieshan Zhao, Xiaofei Liao, Hai Jin, and Jingling Xue. 2020. A heterogeneous PIM hardware-software co-design for energy-efficient graph processing. In *IPDPS*.
- [68] Intel. 2011. Intel 64 and IA-32 Architectures Software Developer Manuals. Retrieved from <http://www.intel.com/content/www/us/en/processors/architectures-software-developer-manuals.html>.
- [69] Intel. 2022. Taking Neuromorphic Computing to the Next Level with Loihi 2. Technology Brief.
- [70] K. Itoh. 2001. *VLSI Memory Chip Design*. Springer.
- [71] JEDEC. 2012. DDR4. JEDEC Standard JESD79–4 (2012).
- [72] Hee Bok Kang and Suk Kyoung Hong. 2009. One-Transistor Type DRAM. US Patent 7701751.
- [73] Mingu Kang, Min-Sun Keel, Naresh R. Shanbhag, Sean Eilert, and Ken Curewitz. 2014. An energy-efficient VLSI architecture for pattern recognition via deep embedding of computation in SRAM. In *ICASSP*.
- [74] Liu Ke, Xuan Zhang, Jinin So, Jong-Geon Lee, Shin-Haeng Kang, Sukhan Lee, Songyi Han, Yeongon Cho, Jin Hyun Kim, Yongsuk Kwon, et al. 2021. Near-memory processing in action: Accelerating personalized recommendation with AxDIMM. *IEEE Micro* 42, 1 (2022), 116–127. <https://doi.org/10.1109/MM.2021.3097700>

- [75] B. Keeth and R. J. Baker. 2001. *DRAM Circuit Design: A Tutorial*. Wiley.
- [76] Samira Khan, Donghyuk Lee, and Onur Mutlu. 2016. PARBOR: An efficient system-level technique to detect data dependent failures in DRAM. In *DSN*.
- [77] Samira Khan, Chris Wilkerson, Z. Wang, Alaa Alameldeen, Donghyuk Lee, and Onur Mutlu. 2017. Detecting and mitigating data-dependent DRAM failures by exploiting current memory content. In *MICRO*.
- [78] C. H. Kim, W. J. Lee, Y. Paik, K. Kwon, S. Y. Kim, I. Park, and S. W. Kim. 2021. Silent-PIM: Realizing the processing-in-memory computing with standard memory requests. *IEEE Transactions on Parallel and Distributed Systems* 33, 2 (2022), 251–262. <https://doi.org/10.1109/TPDS.2021.3065365>
- [79] Duckhwan Kim, J. Kung, S. Chai, S. Yalamanchili, and S. Mukhopadhyay. 2016. Neurocube: A programmable digital neuromorphic architecture with high-density 3D memory. In *ISCA*.
- [80] G. Kim, N. Chatterjee, M. O'Connor, and K. Hsieh. 2017. Toward standardized near-data processing with unrestricted data placement for GPUs. In *SC*.
- [81] J. Kim, M. Patel, H. Hassan, and O. Mutlu. 2018. Solar-DRAM: Reducing DRAM access latency by exploiting the variation in local bitlines. In *ICCD*.
- [82] J. Kim, M. Patel, H. Hassan, and O. Mutlu. 2018. The DRAM latency PUF: Quickly evaluating physical unclonable functions by exploiting the latency–reliability tradeoff in modern DRAM devices. In *HPCA*.
- [83] J. Kim, M. Patel, H. Hassan, L. Orosa, and O. Mutlu. 2019. D-RaNGe: Using commodity DRAM devices to generate true random numbers with low latency and high throughput. In *HPCA*.
- [84] Jin Hyun Kim, Shin-haeng Kang, Sukhan Lee, Hyeonsu Kim, Woongjae Song, Yuhwan Ro, Seungwon Lee, David Wang, Hyunsung Shin, Bengseng Phuah, et al. 2021. Aquabolt-XL: Samsung HBM2-PIM with in-memory processing for ML accelerators and beyond. In *Hot Chips*.
- [85] Jeremie S. Kim, Minesh Patel, A. Giray Yağlıkçı, Hasan Hassan, Roknoddin Azizi, Lois Orosa, and Onur Mutlu. 2020. Revisiting RowHammer: An experimental analysis of modern DRAM devices and mitigation techniques. In *ISCA*.
- [86] J. S. Kim, D. Senol, H. Xin, D. Lee, S. Ghose, M. Alser, H. Hassan, O. Ergin, C. Alkan, and O. Mutlu. 2018. GRIM-Filter: Fast seed location filtering in DNA read mapping using processing-in-memory technologies. 19, 2 (2018), 89.
- [87] Yoongu Kim. 2015. *Architectural Techniques to Enhance DRAM Scaling*. Ph.D. Dissertation. Carnegie Mellon University.
- [88] Yoongu Kim, Ross Daly, Jeremie Kim, Chris Fallin, Ji Hye Lee, Donghyuk Lee, Chris Wilkerson, Konrad Lai, and Onur Mutlu. 2014. Flipping bits in memory without accessing them: An experimental study of DRAM disturbance errors. In *ISCA*.
- [89] Yoongu Kim, Vivek Seshadri, Donghyuk Lee, Jamie Liu, and Onur Mutlu. 2012. A case for exploiting subarray-level parallelism (SALP) in DRAM. In *ISCA*.
- [90] Yoongu Kim, Weikun Yang, and Onur Mutlu. 2015. Ramulator: A fast and extensible DRAM simulator. *IEEE Computer Architecture Letters* 15, 1 (2016), 45–49. <https://doi.org/10.1109/LCA.2015.2414456>
- [91] S. Kvatinsky, D. Belousov, S. Liman, G. Satat, N. Wald, E. G. Friedman, A. Kolodny, and U. C. Weiser. 2014. MAGIC—memristor-aided logic. In *IEEE TCAS II: Express Briefs*. 61 (2014), 895–899.
- [92] S. Kvatinsky, A. Kolodny, U. C. Weiser, and E. G. Friedman. 2011. Memristor-based IMPLY logic design procedure. In *ICCD*.
- [93] S. Kvatinsky, G. Satat, N. Wald, E. G. Friedman, A. Kolodny, and U. C. Weiser. 2014. Memristor-based material implication (IMPLY) logic: Design principles and methodologies. *IEEE Transactions on Very Large Scale Integration (VLSI) Systems* 22 (2014), 2054–2066.
- [94] Y.-C. Kwon, S. H. Lee, J. Lee, S.-H. Kwon, J. M. Ryu, J.-P. Son, O. Seongil, H.-S. Yu, H. Lee, S. Y. Kim, Y. Cho, J. G. Kim, J. Choi, H.-S. Shin, J. Kim, B. Phuah, H. Kim, M. J. Song, A. Choi, D. Kim, S. Kim, E.-B. Kim, D. Wang, S. Kang, Y. Ro, S. Seo, J. Song, J. Youn, K. Sohn, and N. S. Kim. 2021. 25.4 A 20nm 6GB function-in-memory dram, based on HBM2 with a 1.2TFLOPS programmable computing unit using bank-level parallelism, for machine learning applications. In *ISSCC*.
- [95] D. Lee. 2016. *Reducing DRAM Latency at Low Cost by Exploiting Heterogeneity*. Ph.D. Dissertation. Carnegie Mellon University.
- [96] D. Lee, S. Khan, L. Subramanian, S. Ghose, R. Ausavarungnirun, G. Pekhimenko, V. Seshadri, and O. Mutlu. 2017. Design-induced latency variation in modern DRAM chips: Characterization, analysis, and latency reduction mechanisms. In *SIGMETRICS*.
- [97] Donghyuk Lee, Yoongu Kim, Gennady Pekhimenko, Samira Khan, Vivek Seshadri, Kevin Chang, and Onur Mutlu. 2015. Adaptive-latency DRAM: Optimizing DRAM timing for the common-case. In *HPCA*.
- [98] Donghyuk Lee, Yoongu Kim, Vivek Seshadri, Jamie Liu, Lavanya Subramanian, and Onur Mutlu. 2013. Tiered-latency DRAM: A low latency and low cost DRAM architecture. In *HPCA*.

- [99] Donghun Lee, Jinin So, Minseon Ahn, Jong-Geon Lee, Jungmin Kim, Jeonghyeon Cho, Rebholz Oliver, Vishnu Charan Thummala, Ravi Shankar JV, Sachin Suresh Upadhyaya, et al. 2022. Improving in-memory database operations with acceleration DIMM (AxDIMM). In *DaMoN*.
- [100] Donghyuk Lee, Lavanya Subramanian, Rachata Ausavarungnirun, Jongmoo Choi, and Onur Mutlu. 2015. Decoupled direct memory access: Isolating CPU and IO traffic by leveraging a dual-data-port DRAM. In *PACT*.
- [101] Seongju Lee, Kyuyoung Kim, Sanghoon Oh, Joonhong Park, Gimoon Hong, Dongyoon Ka, Kyudong Hwang, Jeongje Park, Kyeongpil Kang, Jungyeon Kim, Junyeol Jeon, Nahsung Kim, Yongkee Kwon, Kornijuk Vladimir, Woojae Shin, Jongsoo Won, Minkyu Lee, Hyunha Joo, Haerang Choi, Jaewook Lee, Donguc Ko, Younggun Jun, Keewon Cho, Ilwoong Kim, Choungki Song, Chunseok Jeong, Daehan Kwon, Jieun Jang, Il Park, Junhyun Chun, and Joohwan Cho. 2022. A 1nm 1.25V 8Gb, 16Gb/s/pin GDDR6-based accelerator-in-memory supporting 1TFLOPS MAC operation and various activation functions for deep-learning applications. In *ISSCC*.
- [102] Yifat Levy, Jehoshua Bruck, Yuval Cassuto, Eby G. Friedman, Avinoam Kolodny, Eitan Yaakobi, and Shahar Kvatinsky. 2014. Logic operations in memory using a memristive akers array. *Microelectr. J.* 45, 11 (2014), 1429–1437. <https://doi.org/10.1016/j.mejo.2014.06.006>
- [103] S. Li, D. Niu, K. T. Malladi, H. Zheng, B. Brennan, and Y. Xie. 2017. DRISA: A DRAM-based reconfigurable in-situ accelerator. In *MICRO*.
- [104] Shuangchen Li, Cong Xu, Qiaosha Zou, Jishen Zhao, Yu Lu, and Yuan Xie. 2016. Pinatubo: A processing-in-memory architecture for bulk bitwise operations in emerging non-volatile memories. In *DAC*.
- [105] Linux man-pages Project. 2022. `calloc(3p)`—Linux Manual Page. Retrieved from <https://man7.org/linux/man-pages/man3/calloc.3p.html>.
- [106] Linux man-pages Project. 2022. `malloc(3)`—Linux Manual Page. Retrieved from <https://man7.org/linux/man-pages/man3/malloc.3.html>.
- [107] Linux man-pages Project. 2022. `memcpy(3)`—Linux Manual Page. Retrieved from <https://man7.org/linux/man-pages/man3/memcpy.3.html>.
- [108] Linux man-pages Project. 2022. `posix_memalign(3)`—Linux Manual Page. Retrieved from https://man7.org/linux/man-pages/man3/posix_memalign.3.html.
- [109] Linux Wiki. 2021. perf: Linux Profiling with Performance Counters. Retrieved from https://perf.wiki.kernel.org/index.php/Main_Page.
- [110] Jamie Liu, Ben Jaiyen, Yoongu Kim, Chris Wilkerson, and Onur Mutlu. 2013. An experimental study of data retention behavior in modern DRAM devices: Implications for retention time profiling mechanisms. In *ISCA*.
- [111] Zhiyu Liu, Irina Calciu, Maurice Herlihy, and Onur Mutlu. 2017. Concurrent data structures for near-memory computing. In *SPAA*.
- [112] Shih-Lien Lu, Ying-Chen Lin, and Chia-Lin Yang. 2015. Improving DRAM latency with dynamic asymmetric subarray. In *MICRO*.
- [113] Haocong Luo, Taha Shahroodi, Hasan Hassan, Minesh Patel, A. Giray Yaglikci, Lois Orosa, Jisung Park, and Onur Mutlu. 2020. CLR-DRAM: A low-cost DRAM architecture enabling dynamic capacity-latency trade-off. In *ISCA*.
- [114] Jack A. Mandelman, Robert H. Dennard, Gary B. Bronner, John K. DeBrosse, Rama Divakaruni, Yujun Li, and Carl J. Radens. 2002. Challenges and future directions for the scaling of dynamic random-access memory (DRAM). *IBM Journal of Research and Development* 46, 2–3 (2002), 187–212. <https://doi.org/10.1147/rd.462.0187>
- [115] Maxwell. 2022. FT20X. Retrieved from <https://www.maxwell-fa.com/upload/files/base/8/m/311.pdf>.
- [116] Micron. 2016. DDR4 SDRAM Datasheet.
- [117] Micron. 2018. DDR3 SDRAM: MT41J128M8. Data Sheet.
- [118] Amir Morad, Leonid Yavits, and Ran Ginosar. 2015. GP-SIMD processing-in-memory. *ACM Trans. Arch. Code Optim.* 11, 4, Article 53 (Jan 2015), 26 pages. <https://doi.org/10.1145/2686875>
- [119] Sergiu Mosanu, Mohammad Nazmus Sakib, Tommy II, Ersin Cukurtas, Alif Ahmed, Preslav Ivanov, Samira Khan, Kevin Skadron, and Mircea Stan. 2022. PiMulator: A fast and flexible processing-in-memory emulation platform. In *DATE*.
- [120] Onur Mutlu, Saugata Ghose, Juan Gómez-Luna, and Rachata Ausavarungnirun. 2021. A modern primer on processing in memory. In *Emerging Computing: From Devices to Systems—Looking Beyond Moore and Von Neumann*. Springer, Singapore, 171–243.
- [121] Lifeng Nai, Ramyad Hadidi, Jaewoong Sim, Hyojong Kim, Pranith Kumar, and Hyesoon Kim. 2017. GraphPIM: Enabling instruction-level PIM offloading in graph computing frameworks. In *HPCA*.
- [122] Dimin Niu, Shuangchen Li, Yuhao Wang, Wei Han, Zhe Zhang, Yijin Guan, Tianchan Guan, Fei Sun, Fei Xue, Lide Duan, et al. 2022. 184QPS/W 64Mb/mm 2 3D Logic-to-DRAM hybrid bonding with process-near-memory engine for recommendation system. In *ISSCC*.
- [123] Ataberk Olgun, Minesh Patel, A. Giray Yaglikci, Haocong Luo, Jeremie S. Kim, Nisa Bostanci, Nandita Vijaykumar, Oguz Ergin, and Onur Mutlu. 2021. QUAC-TRNG: High-Throughput True Random Number Generation Using

- Quadruple Row Activation in Commodity DRAM Chips. arXiv:2105.08955. Retrieved from <https://arxiv.org/abs/2105.08955>.
- [124] Ataberk Olgun, Minesh Patel, A. Giray Yağlıkçı, Haocong Luo, Jeremie S. Kim, F. Nisa Bostancı, Nandita Vijaykumar, Oğuz Ergin, and Onur Mutlu. 2021. QUAC-TRNG: High-throughput true random number generation using quadruple row activation in commodity DRAM chips. In *ISCA*.
 - [125] Geraldo F. Oliveira, Juan Gómez-Luna, Lois Orosa, Saugata Ghose, Nandita Vijaykumar, Ivan Fernandez, Mohammad Sadrosadati, and Onur Mutlu. 2021. DAMOV: A new methodology and benchmark suite for evaluating data movement bottlenecks. *IEEE Access* 9 (2021), 134457–134502. <https://doi.org/10.1109/ACCESS.2021.3110993>
 - [126] Lois Orosa, Yaohua Wang, Mohammad Sadrosadati, Jeremie S. Kim, Minesh Patel, Ivan Puddu, Haocong Luo, Kaveh Razavi, Juan Gómez-Luna, Hasan Hassan, Nika Mansouri-Ghiasi, Saugata Ghose, and Onur Mutlu. 2021. CODIC: A low-cost substrate for enabling custom In-DRAM functionalities and optimizations. In *ISCA*.
 - [127] Lois Orosa, Abdullah Giray Yaglikci, Haocong Luo, Ataberk Olgun, Jisung Park, Hasan Hassan, Minesh Patel, Jeremie S. Kim, and Onur Mutlu. 2021. A deeper look into RowHammer's sensitivities: Experimental analysis of real DRAM chips and implications on future attacks and defenses. In *MICRO*.
 - [128] Minesh Patel, Jeremie S. Kim, and Onur Mutlu. 2017. The reach profiler (REAPER): Enabling the mitigation of DRAM retention failures via profiling at aggressive conditions. In *ISCA*.
 - [129] Minesh Patel, Jeremie S. Kim, Taha Shahroodi, Hasan Hassan, and Onur Mutlu. 2020. Bit-exact ECC recovery (BEER): Determining DRAM On-Die ECC functions by exploiting DRAM data retention characteristics. In *MICRO*.
 - [130] Minesh Patel, Taha Shahroodi, Aditya Manglik, A. Giray Yaglikci, Ataberk Olgun, Haocong Luo, and Onur Mutlu. 2022. A case for transparent reliability in DRAM systems. arXiv:2204.10378. Retrieved from <https://arxiv.org/abs/2204.10378>.
 - [131] Ashutosh Pattnaik, Xulong Tang, Adwait Jog, Onur Kayiran, Asit K. Mishra, Mahmut T. Kandemir, Onur Mutlu, and Chita R. Das. 2016. Scheduling techniques for GPU architectures with processing-in-memory capabilities. In *PACT*.
 - [132] J. Power, J. Hestness, M. S. Orr, M. D. Hill, and D. A. Wood. 2015. gem5-gpu: A heterogeneous CPU-GPU simulator. *IEEE Computer Architecture Letters* 14, 1 (2015), 34–36. <https://doi.org/10.1109/LCA.2014.2299539>
 - [133] Seth H. Pugsley, Jeffrey Jestes, Huihui Zhang, Rajeev Balasubramonian, Vijayalakshmi Srinivasan, Alper Buyuktosunoglu, Al Davis, and Feifei Li. 2014. NDC: Analyzing the impact of 3D-stacked memory+logic devices on mapreduce workloads. In *ISPASS*.
 - [134] S. H. S. Rezaei, M. Modarressi, R. Ausavarungnirun, M. Sadrosadati, O. Mutlu, and M. Daneshtalab. 2020. NoM: Network-on-memory for inter-bank data transfer in highly-banked memories. *IEEE Computer Architecture Letters* 19, 1 (2020), 80–83. <https://doi.org/10.1109/LCA.2020.2990599>
 - [135] RISC-V. 2022. RISC-V Proxy Kernel. Retrieved from <https://github.com/riscv/riscv-pk>.
 - [136] Ronny Ronen, Adi Eliahu, Orian Leitersdorf, Natan Peled, Kunal Korgaonkar, Anupam Chattopadhyay, Ben Perach, and Shahar Kvatinsky. 2022. The bitlet model: A parameterized analytical model to compare PIM and CPU systems. *ACM Journal on Emerging Technologies in Computing Systems* 18, 2, Article 43 (2022), 29 pages. <https://doi.org/10.1145/3465371>
 - [137] SAFARI Research Group. 2015. Ramulator: A DRAM Simulator–GitHub Repository. Retrieved from <https://github.com/CMU-SAFARI/ramulator/>.
 - [138] SAFARI Research Group. 2021. DAMOV–GitHub Repository. Retrieved from <https://github.com/CMU-SAFARI/DAMOV>.
 - [139] SAFARI Research Group. 2021. Ramulator-PIM: A Processing-in-Memory Simulation Framework–GitHub Repository. Retrieved from <https://github.com/CMU-SAFARI/ramulator-pim>.
 - [140] Daniel Sanchez and Christos Kozyrakis. 2013. ZSim: Fast and accurate microarchitectural simulation of thousand-core systems. In *ISCA*.
 - [141] Stefan Saroiu, Alec Wolman, and Lucian Cojocar. 2022. The price of secrecy: How hiding internal DRAM topologies hurts rowhammer defenses. In *IRPS*.
 - [142] V. Seshadri. 2016. *Simple DRAM and Virtual Memory Abstractions to Enable Highly Efficient Memory Systems*. Ph.D. Dissertation. Carnegie Mellon University.
 - [143] Vivek Seshadri, Abhishek Bhowmick, Onur Mutlu, Phillip B. Gibbons, Michael A. Kozuch, and Todd C. Mowry. 2014. The dirty-block index. In *ISCA*.
 - [144] Vivek Seshadri, K. Hsieh, A. Boroumand, D. Lee, M. A. Kozuch, O. Mutlu, P. B. Gibbons, and T. C. Mowry. 2015. Fast bulk bitwise AND and OR in DRAM. In *CAL*.
 - [145] 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. 2013. RowClone: Fast and energy-efficient In-DRAM bulk data copy and initialization. In *MICRO*.

- [146] V. Seshadri, D. Lee, T. Mullins, H. Hassan, A. Boroumand, J. Kim, M. A. Kozuch, O. Mutlu, P. B. Gibbons, and T. C. Mowry. 2016. Buddy-RAM: Improving the performance and efficiency of bulk bitwise operations using DRAM. arXiv:1611.09988. Retrieved from <https://arxiv.org/abs/1611.09988>.
- [147] V. Seshadri, D. Lee, T. Mullins, H. Hassan, A. Boroumand, J. Kim, M. A. Kozuch, O. Mutlu, P. B. Gibbons, and T. C. Mowry. 2017. Ambit: In-memory accelerator for bulk bitwise operations using commodity DRAM technology. In *MICRO*.
- [148] Vivek Seshadri and Onur Mutlu. 2016. The processing using memory paradigm: In-DRAM bulk copy, initialization, bitwise AND and OR. arXiv:1610.09603. Retrieved from <https://arxiv.org/abs/1610.09603>.
- [149] Vivek Seshadri and Onur Mutlu. 2017. Simple operations in memory to reduce data movement. In *Advances in Computers, Volume 106*.
- [150] Vivek Seshadri and Onur Mutlu. 2020. In-DRAM bulk bitwise execution engine. arXiv:1905.09822. Retrieved from <https://arxiv.org/abs/1905.09822>.
- [151] Ali Shafiee, Anirban Nag, Naveen Muralimanohar, Rajeev Balasubramonian, John Paul Strachan, Miao Hu, R. Stanley Williams, and Vivek Srikumar. 2016. ISAAC: A convolutional neural network accelerator with in-situ analog arithmetic in crossbars. In *ISCA*.
- [152] Gagandeep Singh, Juan Gomez-Luna, Giovanni Mariani, Geraldo F. Oliveira, Stefano Corda, Sander Stujik, Onur Mutlu, and Henk Corporaal. 2019. NAPEL: Near-memory computing application performance prediction via ensemble learning. In *DAC*.
- [153] Standard Performance Evaluation Corp. 2006. SPEC CPU 2006. Retrieved from <http://www.spec.org/cpu2006>.
- [154] B. S. Bahar Talukder, V. Menon, B. Ray, T. Neal, and M. Rahman. 2020. Towards the avoidance of counterfeit memory: Identifying the DRAM origin. In *HOST*.
- [155] Eleonora Testa, Mathias Soeken, Odysseas Zografos, Luca Amaru, Praveen Raghavan, Rudy Lauwereins, Pierre-Emmanuel Gaillardon, and Giovanni De Micheli. 2016. Inversion optimization in majority-inverter graphs. In *NANOARCH*.
- [156] T.Ti. 2022. PL & PL-P Series DC Power Supplies Data Sheet—Issue 5. Retrieved from https://resources.aimtti.com/datasheets/AIM-PL+PL-P_series_DC_power_supplies_data_sheet-Iss5.pdf.
- [157] UPMEM. 2018. Introduction to UPMEM PIM. Processing-in-memory (PIM) on DRAM Accelerator.
- [158] A. J. van de Goor and I. Schanstra. 2002. Address and data scrambling: Causes and impact on memory tests. In *IEEE International Workshop on Electronic Design, Test and Applications*.
- [159] Yaohua Wang, Lois Orosa, Xiangjun Peng, Yang Guo, Saugata Ghose, Minesh Patel, Jeremie S. Kim, Juan Gómez Luna, Mohammad Sadosadati, Nika Mansouri Ghiasi, and Onur Mutlu. 2020. FIGARO: Improving system performance via fine-grained In-DRAM data relocation and caching. In *MICRO*.
- [160] Andrew Waterman and Krste Asanovic. 2021. The RISC-V Instruction Set Manual. Retrieved from <https://riscv.org/wp-content/uploads/2019/06/riscv-spec.pdf>.
- [161] Sam (Likun) Xi, Oreoluwa Babarinsa, Manos Athanassoulis, and Stratos Idreos. 2015. Beyond the wall: Near-data processing for databases. In *DaMoN*.
- [162] Lei Xie, Hoang Anh Du Nguyen, Mottaqiallah Taouil, Said Hamdioui, and Koen Bertels. 2015. Fast boolean logic mapped on memristor crossbar. In *ICCD*.
- [163] Xilinx. 2011. *7 Series FPGAs Memory Interface Solutions*.
- [164] Xilinx. 2021. *Vivado Design Suite: Using Constraints*.
- [165] Xilinx. 2021. Xilinx Ultrascale+ MPSoC. Retrieved from <https://www.xilinx.com/products/silicon-devices/soc/zynq-ultrascale-mpsoc.html>.
- [166] Xilinx. 2021. Xilinx Zynq-7000 SoC ZC706 Evaluation Kit. Retrieved from <https://www.xilinx.com/products/boards-and-kits/ek-z7-zc706-g.html>.
- [167] Xin Xin, Youtao Zhang, and Jun Yang. 2020. ELP2IM: Efficient and low power bitwise operation processing in DRAM. In *HPCA*.
- [168] Sheng Xu, Xiaoming Chen, Ying Wang, Yinhe Han, Xuehai Qian, and Xiaowei Li. 2019. PIMSim: A flexible and detailed processing-in-memory simulator. *IEEE Computer Architecture Letters* 18, 1 (2019), 6–9. <https://doi.org/10.1109/LCA.2018.2885752>
- [169] Chao Yu, Sihang Liu, and Samira Khan. 2021. MultiPIM: A detailed and configurable multi-stack processing-in-memory simulator. *IEEE Computer Architecture Letters* 20, 1 (2021), 54–57. <https://doi.org/10.1109/LCA.2021.3061905>
- [170] Jintao Yu, Hoang Anh Du Nguyen, Lei Xie, Mottaqiallah Taouil, and Said Hamdioui. 2018. Memristive devices for computation-in-memory. In *DATE*.
- [171] Yue Zha, Etienne Nowak, and Jing Li. 2019. Liquid silicon: A nonvolatile fully programmable processing-in-memory processor with monolithically integrated ReRAM for big data/machine learning applications. In *VLSIC*.
- [172] D. P. Zhang, N. Jayasena, A. Lyashevsky, J. L. Greathouse, L. Xu, and M. Ignatowski. 2014. TOP-PIM: Throughput-oriented programmable processing in memory. In *HPDC*.

- [173] Jialiang Zhang, Yue Zha, Nicholas Beckwith, Bangya Liu, and Jing Li. 2020. MEG: A RISC-V-based system emulation infrastructure for near-data processing using FPGAs and high-bandwidth memory. *ACM Transactions on Reconfigurable Technology and Systems* 13, 4, Article 19 (Sep 2020), 24 pages. <https://doi.org/10.1145/3409114>
- [174] Liang Zhang and Li Shen. 2022. PIM-HBMSim: A processing in memory simulator based on high bandwidth memory. In *CICA*.
- [175] Mingxing Zhang, Youwei Zhuo, Chao Wang, Mingyu Gao, Yongwei Wu, Kang Chen, Christos Kozyrakis, and Xuehai Qian. 2018. GraphP: Reducing communication for PIM-based graph processing with efficient data partition. In *HPCA*.
- [176] Qiuling Zhu, Tobias Graf, H Ekin Sumbul, Larry Pileggi, and Franz Franchetti. 2013. Accelerating sparse matrix-matrix multiplication with 3D-stacked logic-in-memory hardware. In *HPEC*.
- [177] Youwei Zhuo, Chao Wang, Mingxing Zhang, Rui Wang, Dimin Niu, Yanzhi Wang, and Xuehai Qian. 2019. GraphQ: Scalable PIM-based graph processing. In *MICRO*.

Received 20 December 2021; revised 15 June 2022; accepted 14 July 2022