

# PUMA: Efficient and Low-Cost Memory Allocation and Alignment Support for Processing-Using-Memory Architectures

Geraldo F. Oliveira

Emanuele G. Esposito

Juan Gómez-Luna

Onur Mutlu

ETH Zürich

## 1. Motivation & Problem

Processing-in-memory (PIM) [1–12] is a promising paradigm that aims to alleviate the ever-growing cost of moving data back and forth between computing (e.g., CPU, GPU, accelerators) and memory (e.g., caches, main memory, storage) elements. In PIM architectures, computation is done by adding logic units *near* memory arrays, i.e., processing-near-memory (PNM) [13–98], or by *using* the analog properties of the memory arrays themselves, i.e., processing-using-memory (PUM) [66, 99–139]). Several prior works [66, 101–107, 110, 114–117, 119, 120, 126, 129, 130, 132, 133] have demonstrated the feasibility of processing-using-DRAM (PUD) architectures, which use DRAM cells to implement a variety of PUM operations, including data copy and initialization [104, 116], bitwise Boolean [66, 101, 103, 107, 110], and arithmetic operations [103, 106, 110, 132, 133, 140].

PUD architectures impose a restrictive data layout and alignment for their operands, where source and destination operands (i) *must* reside in the same DRAM subarray (i.e., a group of DRAM rows sharing the same row buffer and row decoder) and (ii) are aligned to the boundaries of a DRAM row. However, standard memory allocation routines (i.e., `malloc`, `posix_memalign`, and `huge pages`-based memory allocation) fail to meet the data layout and alignment requirements for PUD architectures to operate successfully for two main reasons. First, while `malloc` and `posix_memalign` can provide the user application virtually aligned contiguous memory pages, they do *not* guarantee that the allocated virtual pages are also contiguous in physical memory and aligned within a DRAM row. Second, employing `huge pages`-based memory allocation can guarantee that virtual pages are contiguous in physical memory. However, due to its *coarse-grained* page allocation sizes (i.e., Linux-based systems can provide huge pages of 2 MB or 1 GB), a *single* huge page allocation can cover *all* the rows in a DRAM subarray in a single DRAM chip.<sup>1</sup> Therefore, when the PUD instruction requires multiple operands (and thus multiple huge page allocations), it is likely that such operands will resign in different DRAM subarrays, thus imposing extra latency due to inter-subarray data movement [141].

We investigate the potential of using `malloc`, `posix_memalign`, and `huge pages`-based memory allocation for a PUD substrate that can execute AND/OR/NOT Boolean operations (i.e., *Ambit* [101]). We consider that AND/OR/NOT Boolean operations can be executed in the PUD substrate *only* when the data alignment and allocation requirements are met (i.e., source and destination operands are contiguous in physical memory and DRAM row-aligned). We observe that (i) independently of the allocation size for input operands, using `malloc` and `posix_memalign` memory allocators results in 0% of the operations being executed in the PUD substrate due to data misalignment; and (ii) for large-enough allocation sizes (e.g., 32 Kb), *only*

up 60% of the PUD operations that use `huge pages`-based memory allocation can successfully be executed in DRAM. We conclude that traditional memory allocators cannot take full advantage of such PUD techniques since they cannot satisfy the specific memory allocation requirements of PUD substrates. Therefore, our **goal** of this work is to provide a flexible memory allocation mechanism that allows programmers to have control over physical memory allocation and enables PUD execution from the operating system (OS) viewpoint.

## 2. PUMA: Key Idea & Overview

To allow the memory allocation API to influence the OS memory allocator and ensure that memory objects are placed within specific DRAM subarrays, we propose a new *lazy data allocation routine* (in the kernel) for PUD memory objects, called PUMA. The *key idea* of PUMA is to use the internal DRAM mapping information, together with `huge pages`, and then split `huge pages` into *finer-grained* allocation units that are (i) aligned to the page address and size and (ii) virtually contiguous. The PUMA routine has three main components (as Figure 1 illustrates): (i) information regarding the DRAM organization (e.g., row, column, and mat sizes), (ii) the DRAM interleaving scheme, which the memory controller provides via an open firmware device tree [142];<sup>2</sup> and (iii) a `huge pages` pool for PUD memory objects (configured during boot time), which guarantees that virtual addresses assigned to a PUD memory objects are contiguous in the physical address space. The allocation routine uses the DRAM address mapping knowledge to split the `huge pages` into different memory regions. Then, it uses the DRAM interleaving scheme to index each memory region based on their subarray ID (obtained by ORing subarray, bank, channel, and rank mask bits in the DRAM interleaving scheme). PUMA uses an *ordered array* data structure similar to the one used in the Linux Kernel buddy allocator algorithm [146], where each entry represents the number of memory regions in a single subarray. When an application calls the PUD memory allocation API, the allocation routine selects the appropriate memory region that satisfies the memory allocation. PUMA operates by exposing three new memory allocation APIs to the user: (i) `pim_preallocate`, for pre-allocation; (ii) `pim_alloc`, for the first data allocation; and (iii) `pim_alloc_align`, for subsequent aligned allocations.

**Pre-Allocation.** The first step in PUMA is to indicate the number of `huge pages` that are available for PUD allocations using the `pim_preallocate` API (❶ in Figure 1). We left the user the duty to provide the number of `huge pages` used for PUD operations (❷ in Figure 1) because `huge pages` are scarce in the system.

**First Allocation.** PUMA uses the *worst-fit allocation scheme* [147] to manage the allocation of memory regions in the `huge page` pool. The main idea behind this placement

<sup>1</sup>A typical DRAM subarray has 1024 DRAM rows, each with 1024 DRAM columns. Thus, a single DRAM subarray can store 1 MB of data. []

<sup>2</sup>The DRAM interleaving scheme can be obtained by reverse engineering the bit locations of memory addresses [143–145].

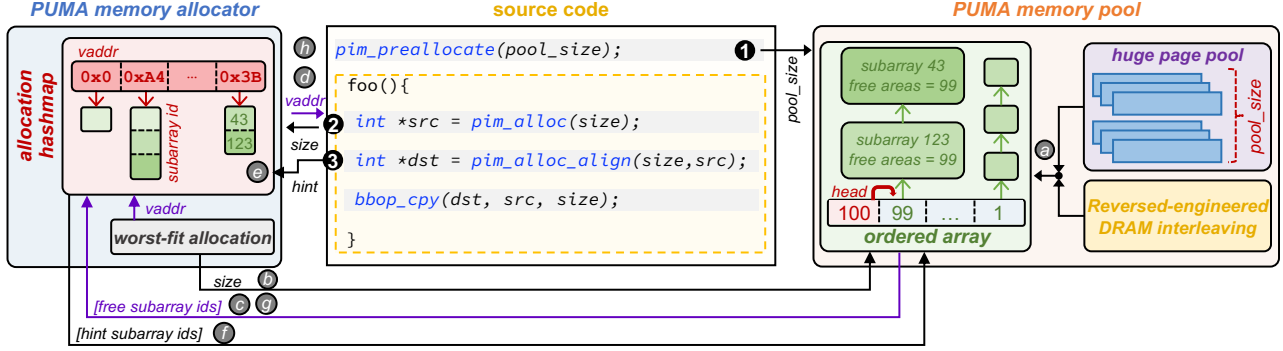


Figure 1: Overview of the PUMA framework.

strategy is to optimize the remaining space post-allocations, thereby increasing the chances of accommodating another process in the remaining memory space. Based on that, for the first PUD memory allocation (using the `pim_alloc` API; ②), PUMA simply scans the *ordered array* to select the subarray with the *largest* amount of memory regions available (③). If the requested memory allocation requires more than one memory region, PUMA interactively scans the *ordered array*, searching for the next largest memory region until the memory allocation is fully satisfied. Once enough space is allocated (④), PUMA creates a new allocation object and inserts it in an *allocation hashmap*, which is indexed (⑤) by the allocation’s virtual address. PUMA needs to keep track of allocations since it might need to find a memory region from the *same* subarray when performing the future aligned allocations (i.e., for the second operand for a Boolean operation).

**Aligned Allocation.** After allocating memory regions for the first operand in a PUD operation, the user can use this memory region as a regular memory object. However, when allocating the remaining operands for a PUD operations (e.g., the second input operand and destination operand in a vector-based Boolean AND operation), PUMA needs to guarantee data alignment for all memory objects within the same DRAM subarray. To this end, we implement a new memory allocation API called `pim_alloc_align`, which takes a `hint` pointer as input (③). Such a pointer indicates a previously allocated memory region to which the current memory allocation must be aligned. The `pim_alloc_align` allocation API works in five main steps. First, PUMA searches the *allocation hashmap* for a match with the address in the `hint` pointer (④). If a match is not found, the allocation fails. Second, if a match is found, PUMA iterates through the `hint`-allocation’s memory regions (⑤). Third, for each memory region, PUMA identifies its source subarray address and tries to allocate another memory region at the same subarray for the new allocation (⑥). Fourth, if the subarray of a given memory region has no free region, PUMA allocates a new memory region from another subarray following the worst-fit allocation scheme (⑦). Since we use a worst-fit allocation scheme, we have a good chance of having a single subarray holding memory regions for multiple allocations. Fifth, since memory regions might come from different huge pages, we must perform `re-mmap` to map such memory regions into contiguous virtual addresses.

### 3. Key Results & Contributions

**Evaluation Methodology.** We implement PUMA as a kernel module using QEMU [148], an open-source emulator and virtualize that can perform hardware virtualization. We emulate a RISC-V machine running Fedora 33 with v5.9.0 Linux Kernel. In our experiments, we evaluate a system with 8 GB DRAM. We emulate the implementation of a PUD system capable of executing row copy operations (as in RowClone [104]) and Boolean AND/OR/NOT operations (as in Ambit [101]). In our experiments, such an operation is performed in the host CPU if a given operation cannot be executed in our PUD substrate (due to data misalignment).

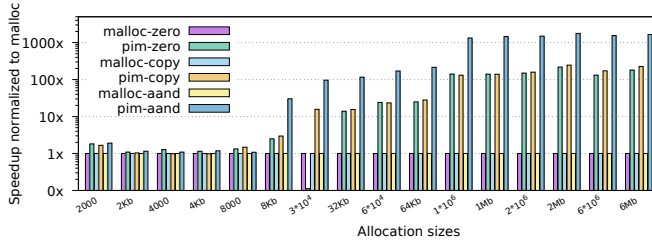
**Baselines & Workloads.** We compare the performance of PUMA to that of using traditional CPU `memcpy` allocation.<sup>3</sup> We use three micro-benchmarks in our analysis: (i) initialize an array with zeros (`*-zero`), (ii) copy data from one array to another (`*-copy`); (iii) perform vector bitwise AND operations `C[i] = A[i] AND B[i]` using Ambit (`*-aand`). For each micro-benchmark, we vary the allocation sizes from 2000 bits to 6 Mb.

**Evaluation Results.** Figure 2 shows PUMA’s performance for each micro-benchmark for different allocation sizes (x-axis) compared to the baseline `malloc` allocator (y-axis). We make two observations for the figure. First, PUMA *significantly* outperforms the baseline memory allocators for all micro-benchmarks and allocation sizes. This is because PUMA increases the likelihood of an operation to be executed in DRAM (due to proper data alignment and allocation), thus increasing overall performance. Second, PUMA’s performance improvements increase as the data allocation sizes increase. This is because the larger the allocation, the more data would need to be moved from DRAM to the CPU in case a PUD operation fails to be executed. Thus severely penalizing overall performance. We conclude that PUMA is a practical and efficient memory allocator for PUD substrates.

We make the following key contributions:

- To our knowledge, this is the first work to propose a practical memory allocation mechanism for PUD substrates.
- We propose PUMA, a data allocation routine for PUD architectures that use the internal DRAM mapping information and huge pages to provide aligned data allocation for PUD instructions.

<sup>3</sup>`posix_mem_align` shows the same performance as `memcpy`.



**Figure 2: PUMA’s performance for three micro-benchmarks and varying data allocation sizes. Values are normalized to the baseline malloc allocator.**

- PUMA does *not* require hardware modifications and operates transparently from the user as a Linux kernel module.
- We evaluate PUMA using three micro-benchmarks, and we observe that PUMA *significantly* increases performance compared to malloc-based memory allocators.

## References

- [1] S. Ghose, A. Boroumand, J. S. Kim, J. Gómez-Luna, and O. Mutlu, “Processing-in-Memory: A Workload-Driven Perspective,” *IBM JRD*, 2019.
- [2] O. Mutlu, S. Ghose, J. Gómez-Luna, and R. Ausavarungnirun, “A Modern Primer on Processing in Memory,” in *Emerging Computing: From Devices to Systems — Looking Beyond Moore and Von Neumann*. Springer, 2021.
- [3] G. F. Oliveira, J. Gómez-Luna, L. Orosa, S. Ghose, N. Vijaykumar, I. Fernandez, M. Sadrosadati, and O. Mutlu, “DAMOV: A New Methodology and Benchmark Suite for Evaluating Data Movement Bottlenecks,” *IEEE Access*, 2021.
- [4] S. Ghose, K. Hsieh, A. Boroumand, R. Ausavarungnirun, and O. Mutlu, “The Processing-in-Memory Paradigm: Mechanisms to Enable Adoption,” in *Beyond-CMOS Technologies for Next Generation Computer Design*, 2019.
- [5] O. Mutlu, S. Ghose, J. Gómez-Luna, and R. Ausavarungnirun, “Processing Data Where It Makes Sense: Enabling In-Memory Computation,” *MicPro*, 2019.
- [6] O. Mutlu, S. Ghose, J. Gómez-Luna, and R. Ausavarungnirun, “Enabling Practical Processing in and Near Memory for Data-Intensive Computing,” in *DAC*, 2019.
- [7] O. Mutlu and L. Subramanian, “Research Problems and Opportunities in Memory Systems,” *SUPERFRI*, 2014.
- [8] O. Mutlu, “Memory Scaling: A Systems Architecture Perspective,” in *IMW*, 2013.
- [9] G. H. Loh, N. Jayasena, M. Oskin, M. Nutter, D. Roberts, M. Meswani, D. P. Zhang, and M. Ignatowski, “A Processing in Memory Taxonomy and a Case for Studying Fixed-Function PIM,” in *WoNDP*, 2013.
- [10] R. Balasubramanian, J. Chang, T. Manning *et al.*, “Near-Data Processing: Insights from a MICRO-46 Workshop,” *IEEE Micro*, 2014.
- [11] H. S. Stone, “A Logic-in-Memory Computer,” *IEEE Trans. Comput.*, 1970.
- [12] A. Saulsbury, F. Pong, and A. Nowatzky, “Missing the Memory Wall: The Case for Processor/Memory Integration,” in *ISCA*, 1996.
- [13] A. Farmahini-Farahani, J. H. Ahn, K. Morrow, and N. S. Kim, “NDA: Near-DRAM Acceleration Architecture Leveraging Commodity DRAM Devices and Standard Memory Modules,” in *HPCA*, 2015.
- [14] O. O. Babarinsa and S. Idreos, “JAFAR: Near-Data Processing for Databases,” in *SIGMOD*, 2015.
- [15] F. Devaux, “The True Processing in Memory Accelerator,” in *Hot Chips*, 2019.
- [16] N. M. Ghiasi, J. Park, H. Mustafa, J. Kim, A. Olgun, A. Gollwitzer, D. S. Cali, C. Firtina, H. Mao, N. A. Alserr *et al.*, “GenStore: A High-Performance and Energy-Efficient In-Storage Computing System for Genome Sequence Analysis,” in *ASPLOS*, 2022.
- [17] J. Gómez-Luna, I. El Hajj, I. Fernandez, C. Giannoula, G. F. Oliveira, and O. Mutlu, “Benchmarking Memory-Centric Computing Systems: Analysis of Real Processing-in-Memory Hardware,” in *CUT*, 2021.
- [18] J. Gómez-Luna, I. E. Hajj, I. Fernández, C. Giannoula, G. F. Oliveira, and O. Mutlu, “Benchmarking a New Paradigm: An Experimental Analysis of a Real Processing-in-Memory Architecture,” arXiv:2105.03814 [cs.AR], 2021.
- [19] J. Gómez-Luna, I. El Hajj, I. Fernandez, C. Giannoula, G. F. Oliveira, and O. Mutlu, “Benchmarking a New Paradigm: Experimental Analysis and Characterization of a Real Processing-in-Memory System,” *IEEE Access*, 2022.
- [20] C. Giannoula, N. Vijaykumar, N. Papadopoulou, V. Karakostas, I. Fernandez, J. Gómez-Luna, L. Orosa, N. Koziris, G. Goumas, and O. Mutlu, “SynCron: Efficient Synchronization Support for Near-Data-Processing Architectures,” in *HPCA*, 2021.
- [21] G. Singh, D. Diamantopoulos, C. Hagleitner, J. Gomez-Luna, S. Stuijk, O. Mutlu, and H. Corporaal, “NERO: A Near High-Bandwidth Memory Stencil Accelerator for Weather Prediction Modeling,” in *FPL*, 2020.
- [22] S. Lee, K. Kim, S. Oh, J. Park, G. Hong, D. Ka, K. Hwang, J. Park, K. Kang, J. Kim, J. Jeon, N. Kim, Y. Kwon, K. Vladimir, W. Shin, J. Won, M. Lee, H. Joo *et al.*, “A 1ynm 1.25V 8Gb, 16Gb/s/pin GDDR6-based Accelerator-in-Memory Supporting ITFLOPS MAC Operation and Various Activation Functions for Deep-Learning Applications,” in *ISSCC*, 2022.
- [23] L. Ke, X. Zhang, J. So, J.-G. Lee, S.-H. Kang, S. Lee, S. Han, Y. Cho, J. H. Kim, Y. Kwon *et al.*, “Near-Memory Processing in Action: Accelerating Personalized Recommendation with AxDIMM,” *IEEE Micro*, 2021.
- [24] C. Giannoula, I. Fernandez, J. G. Luna, N. Koziris, G. Goumas, and O. Mutlu, “SparseP: Towards Efficient Sparse Matrix Vector Multiplication on Real Processing-in-Memory Architectures,” in *SIGMETRICS*, 2022.
- [25] H. Shin, D. Kim, E. Park, S. Park, Y. Park, and S. Yoo, “McDRAM: Low Latency and Energy-Efficient Matrix Computations in DRAM,” *IEEE TCADICS*, 2018.
- [26] S. Cho, H. Choi, E. Park, H. Shin, and S. Yoo, “McDRAM v2: In-Dynamic Random Access Memory Systolic Array Accelerator to Address the Large Model Problem in Deep Neural Networks on the Edge,” *IEEE Access*, 2020.
- [27] A. Denzler, R. Bera, N. Hajinazar, G. Singh, G. F. Oliveira, J. Gómez-Luna, and O. Mutlu, “Casper: Accelerating Stencil Computation using Near-Cache Processing,” arXiv:2112.14216 [cs.AR], 2021.
- [28] H. Asghari-Moghaddam, Y. H. Son, J. H. Ahn, and N. S. Kim, “Chameleon: Versatile and Practical Near-DRAM Acceleration Architecture for Large Memory Systems,” in *MICRO*, 2016.
- [29] D. Patterson, T. Anderson, N. Cardwell *et al.*, “A Case for Intelligent RAM,” *IEEE Micro*, 1997.
- [30] D. G. Elliott, M. Stumm, W. M. Snelgrove *et al.*, “Computational RAM: Implementing Processors in Memory,” *Design and Test of Computers*, vol. 16, no. 1, pp. 32–41, Jan. 1999.
- [31] M. A. Z. Alves, P. C. Santos, F. B. Moreira, and others, “Saving memory movements through vector processing in the dram,” in *Int. Conf. on Compilers, Architecture and Synthesis for Embedded Systems*, 2015.
- [32] S. L. Xi, O. Babarinsa, M. Athanassoulis, and S. Idreos, “Beyond the wall: Near-data processing for databases,” in *Int. Workshop on Data Management on New Hardware*, 2015.
- [33] W. Sun, Z. Li, S. Yin, S. Wei, and L. Liu, “ABC-DIMM: Alleviating the Bottleneck of Communication in DIMM-Based Near-Memory Processing with Inter-DIMM Broadcast,” in *ISCA*, 2021.
- [34] K. K. Matam, G. Koo, H. Zha, H.-W. Tseng, and M. Annavaram, “GraphSSD: Graph Semantics Aware SSD,” in *ISCA*, 2019.
- [35] M. Gokhale, B. Holmes, and K. Iobst, “Processing in Memory: The Terasys Massively Parallel PIM Array,” *Computer*, 1995.
- [36] M. Hall, P. Kogge, J. Koller, P. Diniz, J. Chame, J. Draper, J. LaCoss, J. Granacki, J. Brockman, A. Srivastava *et al.*, “Mapping Irregular Applications to DIVA, a PIM-Based Data-Intensive Architecture,” in *SC*, 1999.
- [37] M. A. Z. Alves, P. C. Santos, M. Diener, and L. Carro, “Opportunities and Challenges of Performing Vector Operations Inside the DRAM,” in *MEMSYS*, 2015.
- [38] E. Lockerman, A. Feldmann, M. Bakhshalipour, A. Stanescu, S. Gupta, D. Sanchez, and N. Beckmann, “Livia: Data-Centric Computing Throughout the Memory Hierarchy,” in *ASPLOS*, 2020.
- [39] J. Ahn, S. Hong, S. Yoo, O. Mutlu, and K. Choi, “A Scalable Processing-in-Memory Accelerator for Parallel Graph Processing,” in *ISCA*, 2015.
- [40] L. Nai, R. Hadidi, J. Sim, H. Kim, P. Kumar, and H. Kim, “Graph-PIM: Enabling Instruction-Level PIM Offloading in Graph Computing Frameworks,” in *HPCA*, 2017.
- [41] A. Boroumand, S. Ghose, Y. Kim, R. Ausavarungnirun, E. Shiu, R. Thakur, D. Kim, A. Kussela, A. Knies, P. Ranganathan *et al.*, “Google Workloads for Consumer Devices: Mitigating Data Movement Bottlenecks,” in *ASPLOS*, 2018.

- [42] A. Boroumand, S. Ghose, B. Lucia, K. Hsieh, K. Malladi, H. Zheng, and O. Mutlu, "LazyPIM: An Efficient Cache Coherence Mechanism for Processing-in-Memory," *CAL*, 2017.
- [43] D. Zhang, N. Jayasena, A. Lyashevsky, J. L. Greathouse, L. Xu, and M. Ignatowski, "TOP-PIM: Throughput-Oriented Programmable Processing in Memory," in *HPDC*, 2014.
- [44] M. Gao and C. Kozyrakis, "HRL: Efficient and Flexible Reconfigurable Logic for Near-Data Processing," in *HPCA*, 2016.
- [45] J. S. Kim, D. S. Cali, H. Xin, D. Lee, S. Ghose, M. Alser, H. Hassan, O. Ergin, C. Alkan, and O. Mutlu, "GRIM-Filter: Fast Seed Location Filtering in DNA Read Mapping Using Processing-in-Memory Technologies," *BMC Genomics*, 2018.
- [46] M. Drumond, A. Daglis, N. Mirzadeh, D. Ustiugov, J. Picorel, B. Falsafi, B. Grot, and D. Pneumatikatos, "The Mondrian Data Engine," in *ISCA*, 2017.
- [47] P. C. Santos, G. F. Oliveira, D. G. Tomé, M. A. Z. Alves, E. C. Almeida, and L. Carro, "Operand Size Reconfiguration for Big Data Processing in Memory," in *DATE*, 2017.
- [48] G. F. Oliveira, P. C. Santos, M. A. Alves, and L. Carro, "NIM: An HMC-Based Machine for Neuron Computation," in *ARC*, 2017.
- [49] J. Ahn, S. Yoo, O. Mutlu, and K. Choi, "PIM-Enabled Instructions: A Low-Overhead, Locality-Aware Processing-in-Memory Architecture," in *ISCA*, 2015.
- [50] M. Gao, J. Pu, X. Yang, M. Horowitz, and C. Kozyrakis, "TETRIS: Scalable and Efficient Neural Network Acceleration with 3D Memory," in *ASPLOS*, 2017.
- [51] D. Kim, J. Kung, S. Chai, S. Yalamanchili, and S. Mukhopadhyay, "Neurocube: A Programmable Digital Neuromorphic Architecture with High-Density 3D Memory," in *ISCA*, 2016.
- [52] P. Gu, S. Li, D. Stow, R. Barnes, L. Liu, Y. Xie, and E. Kursun, "Leveraging 3D Technologies for Hardware Security: Opportunities and Challenges," in *GLSVLSI*, 2016.
- [53] A. Boroumand, S. Ghose, M. Patel, H. Hassan, B. Lucia, R. Ausavarungnirun, K. Hsieh, N. Hajinazar, K. T. Malladi, H. Zheng *et al.*, "CoNDA: Efficient Cache Coherence Support for Near-Data Accelerators," in *ISCA*, 2019.
- [54] K. Hsieh, E. Ebrahimi, G. Kim, N. Chatterjee, M. O'Connor, N. Vijaykumar, O. Mutlu, and S. W. Keckler, "Transparent Offloading and Mapping (TOM) Enabling Programmer-Transparent Near-Data Processing in GPU Systems," in *ISCA*, 2016.
- [55] D. S. Cali, G. S. Kalsi, Z. Bingöl, C. Firtina, L. Subramanian, J. S. Kim, R. Ausavarungnirun, M. Alser, J. Gomez-Luna, A. Boroumand *et al.*, "GenASM: A High-Performance, Low-Power Approximate String Matching Acceleration Framework for Genome Sequence Analysis," in *MICRO*, 2020.
- [56] S. H. Pugsley, J. Jests, H. Zhang, R. Balasubramonian *et al.*, "NDC: Analyzing the Impact of 3D-Stacked Memory+Logic Devices on MapReduce Workloads," in *ISPASS*, 2014.
- [57] A. Pattanaik, X. Tang, A. Jog, O. Kayiran, A. K. Mishra, M. T. Kandemir, O. Mutlu, and C. R. Das, "Scheduling Techniques for GPU Architectures with Processing-in-Memory Capabilities," in *PACT*, 2016.
- [58] B. Akin, F. Franchetti, and J. C. Hoe, "Data Reorganization in Memory Using 3D-Stacked DRAM," in *ISCA*, 2015.
- [59] K. Hsieh, S. Khan, N. Vijaykumar, K. K. Chang, A. Boroumand, S. Ghose, and O. Mutlu, "Accelerating Pointer Chasing in 3D-Stacked Memory: Challenges, Mechanisms, Evaluation," in *ICCD*, 2016.
- [60] J. H. Lee, J. Sim, and H. Kim, "BSSync: Processing Near Memory for Machine Learning Workloads with Bounded Staleness Consistency Models," in *PACT*, 2015.
- [61] A. Boroumand, S. Ghose, B. Akin, R. Narayanaswami, G. F. Oliveira, X. Ma, E. Shiu, and O. Mutlu, "Mitigating Edge Machine Learning Inference Bottlenecks: An Empirical Study on Accelerating Google Edge Models," arXiv:2103.00768 [cs.AR], 2021.
- [62] A. Boroumand, S. Ghose, B. Akin, R. Narayanaswami, G. F. Oliveira, X. Ma, E. Shiu, and O. Mutlu, "Google Neural Network Models for Edge Devices: Analyzing and Mitigating Machine Learning Inference Bottlenecks," in *PACT*, 2021.
- [63] A. Boroumand, S. Ghose, G. F. Oliveira, and O. Mutlu, "Polynesia: Enabling High-Performance and Energy-Efficient Hybrid Transactional/Analytical Databases with Hardware/Software Co-Design," in *ICDE*, 2022.
- [64] A. Boroumand, S. Ghose, G. F. Oliveira, and O. Mutlu, "Polynesia: Enabling Effective Hybrid Transactional/Analytical Databases with Specialized Hardware/Software Co-Design," arXiv:2103.00798 [cs.AR], 2021.
- [65] A. Boroumand, "Practical Mechanisms for Reducing Processor-Memory Data Movement in Modern Workloads," Ph.D. dissertation, Carnegie Mellon University, 2020.
- [66] M. Besta, R. Kanakagiri, G. Kwasniewski, R. Ausavarungnirun, J. Beránek, K. Kanellopoulos, K. Janda, Z. Vonarburg-Shmaria, L. Gianinazzi, I. Stefan *et al.*, "SISA: Set-Centric Instruction Set Architecture for Graph Mining on Processing-in-Memory Systems," in *MICRO*, 2021.
- [67] I. Fernandez, R. Quisilant, E. Gutiérrez, O. Plata, C. Giannoula, M. Alser, J. Gómez-Luna, and O. Mutlu, "NATSA: A Near-Data Processing Accelerator for Time Series Analysis," in *ICCD*, 2020.
- [68] G. Singh, G. F. Oliveira, S. Corda, S. Stuijk, O. Mutlu, and H. Corporaal, "NAPEL: Near-Memory Computing Application Performance Prediction via Ensemble Learning," in *DAC*, 2019.
- [69] 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 *et al.*, "A 20nm 6GB Function-in-Memory DRAM, Based on HBM2 with a 1.2 TFLOPS Programmable Computing Unit using Bank-Level Parallelism, for Machine Learning Applications," in *ISSCC*, 2021.
- [70] S. Lee, S.-h. Kang, J. Lee, H. Kim, E. Lee, S. Seo, H. Yoon, S. Lee, K. Lim, H. Shin *et al.*, "Hardware Architecture and Software Stack for PIM Based on Commercial DRAM Technology: Industrial Product," in *ISCA*, 2021.
- [71] D. Niu, S. Li, Y. Wang, W. Han, Z. Zhang, Y. Guan, T. Guan, F. Sun, F. Xue, L. Duan *et al.*, "184QPS/W 64Mb/mm<sup>2</sup> 3D Logic-to-DRAM Hybrid Bonding with Process-Near-Memory Engine for Recommendation System," in *ISSCC*, 2022.
- [72] Q. Zhu, T. Graf, H. E. Sumbul, L. Pileggi, and F. Franchetti, "Accelerating Sparse Matrix-Matrix Multiplication with 3D-Stacked Logic-in-Memory Hardware," in *HPEC*, 2013.
- [73] E. Azarkhish, C. Pfister, D. Rossi, I. Loi, and L. Benini, "Logic-Base Interconnect Design for Near Memory Computing in the Smart Memory Cube," *IEEE VLSI*, 2016.
- [74] E. Azarkhish, D. Rossi, I. Loi, and L. Benini, "Neurostream: Scalable and Energy Efficient Deep Learning with Smart Memory Cubes," *TPDS*, 2018.
- [75] Q. Guo, N. Alachiotis, B. Akin, F. Sadi, G. Xu, T. M. Low, L. Pileggi, J. C. Hoe, and F. Franchetti, "3D-Stacked Memory-Side Acceleration: Accelerator and System Design," in *WoNDP*, 2014.
- [76] J. P. C. de Lima, P. C. Santos, M. A. Alves, A. Beck, and L. Carro, "Design Space Exploration for PIM Architectures in 3D-Stacked Memories," in *CF*, 2018.
- [77] B. Akin, J. C. Hoe, and F. Franchetti, "HAMLeT: Hardware Accelerated Memory Layout Transform within 3D-Stacked DRAM," in *HPEC*, 2014.
- [78] Y. Huang, L. Zheng, P. Yao, J. Zhao, X. Liao, H. Jin, and J. Xue, "A Heterogeneous PIM Hardware-Software Co-Design for Energy-Efficient Graph Processing," in *IPDPS*, 2020.
- [79] G. Dai, T. Huang, Y. Chi, J. Zhao, G. Sun, Y. Liu, Y. Wang, Y. Xie, and H. Yang, "GraphH: A Processing-in-Memory Architecture for Large-Scale Graph Processing," *TCAD*, 2018.
- [80] J. Liu, H. Zhao, M. A. Ogleari, D. Li, and J. Zhao, "Processing-in-Memory for Energy-Efficient Neural Network Training: A Heterogeneous Approach," in *MICRO*, 2018.
- [81] P.-A. Tsai, C. Chen, and D. Sanchez, "Adaptive Scheduling for Systems with Asymmetric Memory Hierarchies," in *MICRO*, 2018.
- [82] P. Gu, X. Xie, Y. Ding, G. Chen, W. Zhang, D. Niu, and Y. Xie, "iPIM: Programmable In-Memory Image Processing Accelerator using Near-Bank Architecture," in *ISCA*, 2020.
- [83] A. Farmahini-Farahani, J. H. Ahn, K. Compton, and N. S. Kim, "DRAMA: An Architecture for Accelerated Processing Near Memory," *Computer Architecture Letters*, 2014.
- [84] H. Asghari-Moghaddam, A. Farmahini-Farahani, K. Morrow *et al.*, "Near-DRAM Acceleration with Single-ISA Heterogeneous Processing in Standard Memory Modules," *IEEE Micro*, 2016.
- [85] J. Huang, R. R. Puli, P. Majumder, S. Kim, R. Boyapati, K. H. Yum, and E. J. Kim, "Active-Routing: Compute on the Way for Near-Data Processing," in *HPCA*, 2019.
- [86] C. D. Kersey, H. Kim, and S. Yalamanchili, "Lightweight SIMT Core Designs for Intelligent 3D Stacked DRAM," in *MEMSYS*, 2017.
- [87] J. Li, X. Wang, A. Tumeo, B. Williams, J. D. Leidel, and Y. Chen, "PIMS: A Lightweight Processing-in-Memory Accelerator for Stencil Computations," in *MEMSYS*, 2019.
- [88] J. S. Kim, D. Senol, H. Xin, D. Lee, S. Ghose, M. Alser, H. Hassan, O. Ergin, C. Alkan, and O. Mutlu, "GRIM-Filter: Fast Seed Filtering in Read Mapping using Emerging Memory Technologies," arXiv:1708.04329 [q-bio.GN], 2017.



- [89] A. Boroumand, S. Ghose, M. Patel, H. Hassan, B. Lucia, N. Hajinazar, K. Hsieh, K. T. Malladi, H. Zheng, and O. Mutlu, "LazyPIM: Efficient Support for Cache Coherence in Processing-in-Memory Architectures," arXiv:1706.03162 [cs.AR], 2017.
- [90] Y. Zhuo, C. Wang, M. Zhang, R. Wang, D. Niu, Y. Wang, and X. Qian, "GraphQ: Scalable PIM-Based Graph Processing," in *MICRO*, 2019.
- [91] M. Zhang, Y. Zhuo, C. Wang, M. Gao, Y. Wu, K. Chen, C. Kozyrakis, and X. Qian, "GraphP: Reducing Communication for PIM-Based Graph Processing with Efficient Data Partition," in *HPCA*, 2018.
- [92] H. Lim and G. Park, "Triple Engine Processor (TEP): A Heterogeneous Near-Memory Processor for Diverse Kernel Operations," *TACO*, 2017.
- [93] E. Azarkhish, D. Rossi, I. Loi, and L. Benini, "A Case for Near Memory Computation Inside the Smart Memory Cube," in *EMS*, 2016.
- [94] M. A. Z. Alves, M. Diener, P. C. Santos, and L. Carro, "Large Vector Extensions Inside the HMC," in *DATE*, 2016.
- [95] J. Jang, J. Heo, Y. Lee, J. Won, S. Kim, S. J. Jung, H. Jang, T. J. Ham, and J. W. Lee, "Charon: Specialized Near-Memory Processing Architecture for Clearing Dead Objects in Memory," in *MICRO*, 2019.
- [96] R. Nair, S. F. Antao, C. Bertolli, P. Bose *et al.*, "Active Memory Cube: A Processing-in-Memory Architecture for Exascale Systems," *IBM JRD*, 2015.
- [97] R. Hadidi, L. Nai, H. Kim, and H. Kim, "CAIRO: A Compiler-Assisted Technique for Enabling Instruction-Level Offloading of Processing-in-Memory," *TACO*, 2017.
- [98] P. C. Santos, G. F. Oliveira, J. P. Lima, M. A. Alves, L. Carro, and A. C. Beck, "Processing in 3D Memories to Speed Up Operations on Complex Data Structures," in *DATE*, 2018.
- [99] P. Chi, S. Li, C. Xu, T. Zhang, J. Zhao, Y. Liu, Y. Wang, and Y. Xie, "PRIME: A Novel Processing-in-Memory Architecture for Neural Network Computation in ReRAM-Based Main Memory," in *ISCA*, 2016.
- [100] A. Shafiee, A. Nag, N. Muralimanohar, R. Balasubramonian, J. P. Strachan, M. Hu, R. S. Williams, and V. Srikumar, "ISAAC: A Convolutional Neural Network Accelerator with In-Situ Analog Arithmetic in Crossbars," in *ISCA*, 2016.
- [101] V. Seshadri, D. Lee, T. Mullins, H. Hassan, A. Boroumand, J. Kim, M. A. Kozuch, O. Mutlu, P. B. Gibbons, and T. C. Mowry, "Ambit: In-Memory Accelerator for Bulk Bitwise Operations Using Commodity DRAM Technology," in *MICRO*, 2017.
- [102] V. Seshadri and O. Mutlu, "In-DRAM Bulk Bitwise Execution Engine," arXiv:1905.09822 [cs.AR], 2019.
- [103] S. Li, D. Niu, K. T. Malladi, H. Zheng, B. Brennan, and Y. Xie, "DRISA: A DRAM-Based Reconfigurable In-Situ Accelerator," in *MICRO*, 2017.
- [104] V. Seshadri, Y. Kim, C. Fallin, D. Lee, R. Ausavarungnirun, G. Pekhimenko, Y. Luo, O. Mutlu, P. B. Gibbons, M. A. Kozuch *et al.*, "RowClone: Fast and Energy-Efficient In-DRAM Bulk Data Copy and Initialization," in *MICRO*, 2013.
- [105] V. Seshadri and O. Mutlu, "The Processing Using Memory Paradigm: In-DRAM Bulk Copy, Initialization, Bitwise AND and OR," arXiv:1610.09603 [cs.AR], 2016.
- [106] Q. Deng, L. Jiang, Y. Zhang, M. Zhang, and J. Yang, "DrAcc: A DRAM Based Accelerator for Accurate CNN Inference," in *DAC*, 2018.
- [107] X. Xin, Y. Zhang, and J. Yang, "ELP2IM: Efficient and Low Power Bitwise Operation Processing in DRAM," in *HPCA*, 2020.
- [108] L. Song, Y. Zhuo, X. Qian, H. Li, and Y. Chen, "GraphR: Accelerating Graph Processing Using ReRAM," in *HPCA*, 2018.
- [109] L. Song, X. Qian, H. Li, and Y. Chen, "PipeLayer: A Pipelined ReRAM-Based Accelerator for Deep Learning," in *HPCA*, 2017.
- [110] F. Gao, G. Tziantzioulis, and D. Wentzlaff, "ComputeDRAM: In-Memory Compute Using Off-the-Shelf DRAMs," in *MICRO*, 2019.
- [111] C. Eckert, X. Wang, J. Wang, A. Subramaniyan, R. Iyer, D. Sylvester, D. Blaauw, and R. Das, "Neural Cache: Bit-Serial In-Cache Acceleration of Deep Neural Networks," in *ISCA*, 2018.
- [112] S. Aga, S. Jeloka, A. Subramaniyan, S. Narayanasamy, D. Blaauw, and R. Das, "Compute Caches," in *HPCA*, 2017.
- [113] D. Fujiki, S. Mahlke, and R. Das, "Duality Cache for Data Parallel Acceleration," in *ISCA*, 2019.
- [114] V. Seshadri, D. Lee, T. Mullins, H. Hassan, A. Boroumand, J. Kim, M. A. Kozuch, O. Mutlu, P. B. Gibbons, and T. C. Mowry, "BuddyRAM: Improving the Performance and Efficiency of Bulk Bitwise Operations Using DRAM," arXiv:1611.09988 [cs.AR], 2016.
- [115] V. Seshadri and O. Mutlu, "Simple Operations in Memory to Reduce Data Movement," in *Advances in Computers, Volume 106*, 2017.
- [116] V. Seshadri, Y. Kim, C. Fallin, D. Lee, R. Ausavarungnirun, G. Pekhimenko, Y. Luo, O. Mutlu, P. B. Gibbons, M. A. Kozuch *et al.*, "RowClone: Accelerating Data Movement and Initialization Using DRAM," arXiv:1805.03502 [cs.AR], 2018.
- [117] V. Seshadri, K. Hsieh, A. Boroum, D. Lee, M. A. Kozuch, O. Mutlu, P. B. Gibbons, and T. C. Mowry, "Fast Bulk Bitwise AND and OR in DRAM," *CAL*, 2015.
- [118] S. Li, C. Xu, Q. Zou, J. Zhao, Y. Lu, and Y. Xie, "Pinatubo: A Processing-in-Memory Architecture for Bulk Bitwise Operations in Emerging Non-Volatile Memories," in *DAC*, 2016.
- [119] J. D. Ferreira, G. Falcao, J. Gómez-Luna, M. Alser, L. Orosa, M. Sadrosadati, J. S. Kim, G. F. Oliveira, T. Shahroodi, A. Nori *et al.*, "pLUTO: In-DRAM Lookup Tables to Enable Massively Parallel General-Purpose Computation," arXiv:2104.07699 [cs.AR], 2021.
- [120] J. D. Ferreira, G. Falcao, J. Gómez-Luna, M. Alser, L. Orosa, M. Sadrosadati, J. S. Kim, G. F. Oliveira, T. Shahroodi, A. Nori *et al.*, "pLUTO: Enabling Massively Parallel Computation in DRAM via Lookup Tables," in *MICRO*, 2022.
- [121] M. Imani, S. Gupta, Y. Kim, and T. Rosing, "FloatPIM: In-Memory Acceleration of Deep Neural Network Training with High Precision," in *ISCA*, 2019.
- [122] Z. He, L. Yang, S. Angizi, A. S. Rakin, and D. Fan, "Sparse BD-Net: A Multiplication-Less DNN with Sparse Binarized Depth-Wise Separable Convolution," *JETC*, 2020.
- [123] J. Park, R. Azizi, G. F. Oliveira, M. Sadrosadati, R. Nadig, D. Novo, J. Gómez-Luna, M. Kim, and O. Mutlu, "Flash-Cosmos: In-Flash Bulk Bitwise Operations Using Inherent Computation Capability of NAND Flash Memory," in *MICRO*, 2022.
- [124] M. S. Truong, L. Shen, A. Glass, A. Hoffmann, L. R. Carley, J. A. Bain, and S. Ghose, "Adapting the RACER Architecture to Integrate Improved In-ReRAM Logic Primitives," *JETCAS*, 2022.
- [125] M. S. Truong, E. Chen, D. Su, L. Shen, A. Glass, L. R. Carley, J. A. Bain, and S. Ghose, "RACER: Bit-Pipelined Processing Using Resistive Memory," in *MICRO*, 2021.
- [126] A. Olgun, M. Patel, A. G. Yağlıkcı, H. Luo, J. S. Kim, F. N. Bostancı, N. Vijaykumar, O. Ergin, and O. Mutlu, "QUAC-TRNG: High-Throughput True Random Number Generation Using Quadruple Row Activation in Commodity DRAMs," in *ISCA*, 2021.
- [127] J. S. Kim, M. Patel, H. Hassan, L. Orosa, and O. Mutlu, "D-RaNGe: Using Commodity DRAM Devices to Generate True Random Numbers With Low Latency and High Throughput," in *HPCA*, 2019.
- [128] J. S. Kim, M. Patel, H. Hassan, and O. Mutlu, "The DRAM Latency PUF: Quickly Evaluating Physical Unclonable Functions by Exploiting the Latency-Reliability Tradeoff in Modern Commodity DRAM Devices," in *HPCA*, 2018.
- [129] F. N. Bostancı, A. Olgun, L. Orosa, A. G. Yağlıkcı, J. S. Kim, H. Hassan, O. Ergin, and O. Mutlu, "DR-STRaNGe: End-to-End System Design for DRAM-Based True Random Number Generators," in *HPCA*, 2022.
- [130] A. Olgun, J. G. Luna, K. Kanellopoulos, B. Salami, H. Hassan, O. Ergin, and O. Mutlu, "PiDRAM: A Holistic End-to-End FPGA-Based Framework for Processing-in-DRAM," *TACO*, 2022.
- [131] M. F. Ali, A. Jaiswal, and K. Roy, "In-Memory Low-Cost Bit-Serial Addition Using Commodity DRAM Technology," in *TCAS-I*, 2019.
- [132] S. Angizi and D. Fan, "GraphiDe: A Graph Processing Accelerator Leveraging In-DRAM-Computing," in *GLSVLSI*, 2019.
- [133] S. Li, A. O. Glova, X. Hu, P. Gu, D. Niu, K. T. Malladi, H. Zheng, B. Brennan, and Y. Xie, "SCOPE: A Stochastic Computing Engine for DRAM-Based In-Situ Accelerator," in *MICRO*, 2018.
- [134] A. Subramaniyan and R. Das, "Parallel Automata Processor," in *ISCA*, 2017.
- [135] Y. Zha and J. Li, "Hyper-AP: Enhancing Associative Processing Through A Full-Stack Optimization," in *ISCA*, 2020.
- [136] D. Fujiki, S. Mahlke, and R. Das, "In-Memory Data Parallel Processor," in *ASPLOS*, 2018.
- [137] L. Orosa, Y. Wang, M. Sadrosadati, J. Kim, M. Patel, I. Puddu, H. Luo, K. Razavi, J. Gómez-Luna, H. Hassan, N. M. Ghiasi, S. Ghose, and O. Mutlu, "CODIC: A Low-Cost Substrate for Enabling Custom In-DRAM Functionalities and Optimizations," in *ISCA*, 2021.
- [138] M. Sharad, D. Fan, and K. Roy, "Ultra Low Power Associative Computing with Spin Neurons and Resistive Crossbar Memory," in *DAC*, 2013.
- [139] S. H. S. Rezaei, M. Modarressi, R. Ausavarungnirun, M. Sadrosadati, O. Mutlu, and M. Daneshlab, "NoM: Network-on-Memory for Inter-Bank Data Transfer in Highly-Banked Memories," *CAL*, 2020.
- [140] N. Hajinazar, G. F. Oliveira, S. Gregorio, J. D. Ferreira, N. M. Ghiasi, M. Patel, M. Alser, S. Ghose, J. Gómez-Luna, and O. Mutlu, "SIM-DRAM: A Framework for Bit-Serial SIMD Processing Using DRAM," in *ASPLOS*, 2021.

- [141] K. K. Chang, P. J. Nair, D. Lee, S. Ghose, M. K. Qureshi, and O. Mutlu, "Low-Cost Inter-Linked Subarrays (LISA): Enabling Fast Inter-Subarray Data Movement in DRAM," in *HPCA*, 2016.
- [142] T. K. D. Community, "Linux and the Devicetree The Linux Kernel Documentation," <https://www.kernel.org/doc/html/latest/devicetree/usage-model.html>.
- [143] J. S. Kim, M. Patel, A. G. Yağlıkçı, H. Hassan, R. Azizi, L. Orosa, and O. Mutlu, "Revisiting RowHammer: An Experimental Analysis of Modern DRAM Devices and Mitigation Techniques," in *ISCA*, 2020.
- [144] L. Orosa, A. G. Yaglikci, H. Luo, A. Olgun, J. Park, H. Hassan, M. Patel, J. S. Kim, and O. Mutlu, "A Deeper Look into RowHammer's Sensitiv-  
ities: Experimental Analysis of Real DRAM Chips and Implications on Future Attacks and Defenses"," in *MICRO*, 2021.
- [145] A. G. Yağlıkçı, H. Luo, G. F. De Oliveira, A. Olgun, M. Patel, J. Park, H. Hassan, J. S. Kim, L. Orosa, and O. Mutlu, "Understanding RowHammer Under Reduced Wordline Voltage: An Experimental Study Using Real DRAM Devices," in *DSN*, 2022.
- [146] K. C. Knowlton, "A Programmer's Description of L6," *CACM*, 1966.
- [147] D. S. Johnson, "Near-Optimal Bin Packing Algorithms," Ph.D. dissertation, Massachusetts Institute of Technology, 1973.
- [148] F. Bellard, "QEMU, a Fast and Portable Dynamic Translator," in *USENIX ATC*, 2005.