

Towards Efficient Sparse Matrix Vector Multiplication on Real Processing-In-Memory Architectures

Christina Giannoula

ETH Zürich National Technical University of Athens
christina.giann@gmail.com

Juan Gómez-Luna

ETH Zürich
el1goluj@gmail.com

Georgios Goumas

National Technical University of Athens
goumas@cslab.ece.ntua.gr

Ivan Fernandez

ETH Zürich University of Malaga
ivanferveg@gmail.com

Nectarios Koziris

National Technical University of Athens
nkoziris@cslab.ece.ntua.gr

Onur Mutlu

ETH Zürich
omutlu@gmail.com

Abstract

Sparse Matrix Vector Multiplication (SpMV) has been characterized as one of the most thoroughly studied scientific computation kernels, because it is a fundamental linear algebra kernel for important applications from the scientific computing, machine learning, and graph analytics domains. SpMV performs indirect memory references as a result of storing the sparse matrix in a compressed format, and irregular memory accesses to the input vector due to the sparsity pattern of the input matrix [11]. Therefore, in commodity processor-centric systems, SpMV is a primarily memory-bandwidth-bound kernel for the majority of real sparse matrices, and is bottlenecked by data movement between memory and processors [5, 8, 9].

One promising way to alleviate the data movement bottleneck is the Processing-In-Memory (PIM) paradigm [1, 2, 6–9, 12–14]. PIM moves computation close to application data by equipping memory chips with processing capabilities [1, 12, 13]. To provide large aggregate memory bandwidth for the in-memory processors, several manufacturers have already started to commercialize *near-bank* PIM designs [1, 8, 9, 16, 17]. *Near-bank* PIM designs tightly couple a PIM core with each DRAM bank, exploiting bank-level parallelism to expose high on-chip memory bandwidth of standard DRAM to processors. Three *real* near-bank PIM architectures are Samsung’s FIMDRAM [16], SK hynix’s GDDR6-AiM [17] and the UPMEM PIM system [1, 8, 9].

Most real near-bank PIM architectures [1, 8, 9, 16, 17] support several PIM-enabled memory chips connected to a host CPU via memory channels. Each memory chip comprises multiple low-power PIM cores with relatively low computation capability [8, 9], and each of them is located close to a DRAM bank [8, 9, 16, 17]. Each PIM core can access data located on its local DRAM bank, and typically there is no direct communication channel among PIM cores. Overall, near-bank PIM systems provide high levels of parallelism and very large memory bandwidth, and are thus a very promising computing platform to accelerate the widely-used SpMV kernel.

Our work is the first to efficiently map the SpMV kernel on near-bank PIM systems, and understand its performance implications on a real-world PIM system. We make two key contributions. First, we design efficient SpMV algorithms to accelerate the SpMV kernel in current and future PIM systems, while covering a wide variety of sparse matrices with diverse sparsity patterns. Second, we provide the first comprehensive analysis of SpMV on a real PIM architecture. Specifically, we conduct our rigorous experimental analysis of SpMV kernels in the UPMEM PIM system, the first publicly-available real-world PIM architecture.

We present the freely and openly available *SparseP* library [21] that includes 25 SpMV kernels for real PIM systems. *SparseP* supports (1) the most popular compressed matrix formats (i.e., CSR, COO, BCSR, BCOO formats), (2) a wide range of data types (i.e., 8-bit integer, 16-bit integer, 32-bit integer, 64-bit integer, 32-bit float and 64-bit float data types), (3) two types of well-crafted data partitioning techniques of the sparse matrix to PIM-enabled memory, (4) various load balancing schemes across PIM cores, (5) various load balancing schemes across threads of a multithreaded PIM core, and (6) three synchronization approaches among threads within multithreaded PIM core.

We conduct an extensive characterization and analysis of *SparseP* kernels on the UPMEM PIM system [8, 9]. We analyze the SpMV execution (1) using one single multithreaded PIM core, (2) using thousands of PIM cores, and (3) comparing its performance and energy consumption with that achieved on conventional processor-centric CPU and GPU systems. Our extensive evaluation provides programming recommendations for software designers, and suggestions and hints for hardware and system designers of future PIM systems.

We highlight our most significant recommendations for PIM software designers:

- (1) *Design algorithms that provide high load balance across threads of a multithreaded PIM core in terms of computations, loop control iterations, synchronization points and memory accesses.* In SpMV, we find that when the parallelization scheme used causes high disparity in the non-zero elements/blocks/rows processed across threads of a PIM core, or the number of lock acquisitions/lock releases/DRAM memory accesses performed across threads, performance severely degrades in low-area PIM cores with relatively low computation capabilities [8, 9].
- (2) *Design compressed data structures that can be effectively partitioned across DRAM banks, with the goal of providing high*

Permission to make digital or hard copies of part or all of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for third-party components of this work must be honored. For all other uses, contact the owner/author(s).

SIGMETRICS/PERFORMANCE '22 Abstracts, June 6–10, 2022, Mumbai, India.

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

ACM ISBN 978-1-4503-9141-2/22/06.

<https://doi.org/10.1145/3489048.3522661>

computation balance across PIM cores. We observe that the compressed matrix format used to store the input matrix in SpMV determines the data partitioning across DRAM banks of PIM-enabled memory, thereby affecting the load balance across PIM cores with corresponding performance implications.

- (3) *Design adaptive algorithms that trade off computation balance across PIM cores for lower data transfer costs to PIM-enabled memory, and adapt the software strategies to the particular patterns of each input given, as well as the characteristics of the PIM hardware.* Our analysis demonstrates that the best-performing SpMV execution on the UPMEM PIM system is achieved using algorithms that (i) trade off computation for lower data transfer costs, and (ii) select the load balancing strategy and data partitioning policy based on the particular sparsity pattern of the input matrix and the characteristics of the underlying PIM hardware.

We highlight our most significant suggestions for PIM hardware and system designers:

- (1) *Provide low-cost synchronization support and hardware support to enable concurrent memory accesses by multiple threads to the local DRAM bank to increase parallelism in a multithreaded PIM core.* For instance, fine-grained locking approaches in SpMV to increase parallelism in critical sections do not improve performance over coarse-grained approaches. This is because concurrent DRAM accesses performed by multiple threads are serialized by the UPMEM PIM hardware. To improve parallelism, subarray level parallelism [19] or multiple DRAM banks per PIM core could be supported in the PIM hardware, along with lightweight synchronization schemes for PIM cores [2].
- (2) *Optimize the broadcast collective operation in data transfers from main memory to PIM-enabled memory to minimize overheads of copying the input data into all DRAM banks in the PIM system.* When the sparse matrix is horizontally partitioned across PIM cores and the whole input vector is copied into the DRAM bank of each PIM core, SpMV cannot scale up to a large number of PIM cores. This is because it is severely limited by data transfer costs to broadcast the input vector into each DRAM bank of PIM-enabled DIMMs. Such data transfers incur high overheads, because they take place via the narrow off-chip memory bus.
- (3) *Optimize the gather collective operation at DRAM bank granularity for data transfers from PIM-enabled memory to the host CPU to minimize overheads of retrieving the output results.* When the sparse matrix is split in 2D tiles, each of them is assigned to each PIM core, SpMV is severely limited by data transfers to retrieve results for the output vector from DRAM banks of PIM-enabled memory. This is due to two reasons: (i) PIM cores create a large number of partial results that need to be gathered from PIM-enabled memory to the host CPU via the narrow memory bus in order to assemble the final output vector, and (ii) the current implementation of the UPMEM PIM system has the limitation that the transfer sizes from/to all DRAM banks involved in the same parallel transfer need to be the same, and therefore a large amount of padding with empty bytes is performed in such SpMV kernels.
- (4) *Design high-speed communication channels and optimized libraries for data transfers to/from thousands of DRAM banks of PIM-enabled memory.* We find that SpMV execution on the memory-centric UPMEM PIM system achieves a much higher fraction of the machine's peak performance (on average 51.7%

for the 32-bit float data type), compared to that on processor-centric CPU and GPU systems. However, its end-to-end performance is still significantly limited by data transfer overheads on the narrow memory bus. Thus, the software stack of real PIM systems needs to be enhanced with fast data transfers to/from PIM-enabled memory modules, and/or the PIM hardware needs to be enhanced to support efficient direct communication among PIM cores [10, 15, 18, 20].

For more information about our thorough characterization on the SpMV PIM execution, results, insights and the open-source *SparseP* software package [21], we refer the reader to the full version of the paper [3, 4]. We hope that our work can provide valuable insights to programmers in the development of efficient sparse linear algebra kernels and other irregular kernels from different application domains tailored for real PIM systems, and enlighten architects and system designers in the development of future memory-centric computing systems. The *SparseP* software package is publicly and freely available at <https://github.com/CMU-SAFARI/SparseP>.

ACM Reference Format:

Christina Giannoula, Ivan Fernandez, Juan Gómez-Luna, Nectarios Koziris, Georgios Goumas, and Onur Mutlu. 2022. Towards Efficient Sparse Matrix Vector Multiplication on Real Processing-In-Memory Architectures. In *Abstract Proceedings of the 2022 ACM SIGMETRICS/IFIP PERFORMANCE Joint International Conference on Measurement and Modeling of Computer Systems (SIGMETRICS/PERFORMANCE '22 Abstracts)*, June 6–10, 2022, Mumbai, India. 2 pages. <https://doi.org/10.1145/3489048.3522661>

References

- [1] F. Devaux. 2019. The True Processing In Memory Accelerator. In *Hot Chips*.
- [2] Christina Giannoula et al. 2021. SynCron: Efficient Synchronization Support for Near-Data-Processing Architectures. In *HPCA*.
- [3] Christina Giannoula et al. 2022. SparseP: Towards Efficient Sparse Matrix Vector Multiplication on Real Processing-In-Memory Architectures. *Proc. ACM Meas. Anal. Comput. Syst.*
- [4] Christina Giannoula et al. 2022. SparseP: Towards Efficient Sparse Matrix Vector Multiplication on Real Processing-In-Memory Systems. In *CoRR*. <https://arxiv.org/abs/2201.05072>
- [5] Georgios Goumas et al. 2009. Performance Evaluation of the Sparse Matrix-Vector Multiplication on Modern Architectures. In *J. Supercomput.*
- [6] Ivan Fernandez et al. 2020. NATSA: A Near-Data Processing Accelerator for Time Series Analysis. In *ICCD*.
- [7] Junwhan Ahn et al. 2015. A Scalable Processing-In-Memory Accelerator for Parallel Graph Processing. In *ISCA*.
- [8] Juan Gómez-Luna et al. 2021. Benchmarking a New Paradigm: An Experimental Analysis of a Real Processing-in-Memory Architecture. In *CoRR*.
- [9] Juan Gómez-Luna et al. 2021. Benchmarking Memory-Centric Computing Systems: Analysis of Real Processing-In-Memory Hardware. In *JGSC*.
- [10] Kevin Chang et al. 2016. Low-Cost Inter-Linked Subarrays (LISA): Enabling Fast Inter-Subarray Data Movement in DRAM. In *HPCA*.
- [11] Konstantinos Kanellopoulos et al. 2019. SMASH: Co-Designing Software Compression and Hardware-Accelerated Indexing for Efficient Sparse Matrix Operations. In *MICRO*.
- [12] Onur Mutlu et al. 2019. Processing Data Where It Makes Sense: Enabling In-Memory Computation. In *MICPRO*.
- [13] Onur Mutlu et al. 2021. A Modern Primer on Processing in Memory. In *Emerging Computing: From Devices to Systems - Looking Beyond Moore and Von Neumann*.
- [14] Saugata Ghose et al. 2019. Processing-in-Memory: A Workload-Driven Perspective. In *IBM JRD*.
- [15] Seyyed Hossein SeyyedAghaei Rezaei et al. 2020. NoM: Network-on-Memory for Inter-Bank Data Transfer in Highly-Banked Memories. In *IEEE CAL*.
- [16] Sukhan Lee et al. 2021. Hardware Architecture and Software Stack for PIM Based on Commercial DRAM Technology: Industrial Product. In *ISCA*.
- [17] Seongju Lee et al. 2022. A 1ynm 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*.
- [18] Vivek Seshadri et al. 2013. RowClone: Fast and energy-efficient in-DRAM bulk data copy and initialization. In *MICRO*.
- [19] Yoongu Kim et al. 2012. A Case for Exploiting Subarray-Level Parallelism (SALP) in DRAM. In *ISCA*.
- [20] Yaohua Wang et al. 2020. FIGARO: Improving System Performance via Fine-Grained In-DRAM Data Relocation and Caching. In *MICRO*.
- [21] SAFARI Research Group. 2022. *SparseP* Software Package. <https://github.com/CMU-SAFARI/SparseP>