

Machine Learning Training on a Real Processing-in-Memory System

Juan Gómez-Luna¹ Yuxin Guo¹ Sylvan Brocard² Julien Legriel²
Remy Cimadomo² Geraldo F. Oliveira¹ Gagandeep Singh¹ Onur Mutlu¹

¹ETH Zürich ²UPMEM

ABSTRACT

Machine learning (ML) algorithms [1–6] have become ubiquitous in many fields of science and technology due to their ability to learn from and improve with experience with minimal human intervention. These algorithms train by updating their model parameters in an iterative manner to improve the overall prediction accuracy. However, training machine learning algorithms is a computationally intensive process, which requires large amounts of training data. Accessing training data in current processor-centric systems (e.g., CPU, GPU) implies costly data movement between memory and processors, which results in high energy consumption and a large percentage of the total execution cycles. This data movement can become the bottleneck of the training process, if there is not enough computation and locality to amortize its cost.

One way to alleviate the cost of data movement is *processing-in-memory* (PIM) [7–11], a data-centric computing paradigm that places processing elements near or inside the memory arrays. PIM has been explored for decades [9, 12–146]. However, memory technology challenges prevented from its successful materialization in commercial products. For example, the limited number of metal layers in DRAM [147, 148] makes conventional processor designs impractical in commodity DRAM chips [149–152].

Real-world PIM systems have only recently been manufactured and commercialized. The UPMEM company, for example, introduced the first general-purpose commercial PIM architecture [153–157], which integrates small in-order cores near DRAM memory banks. High-bandwidth memory (HBM)-based HBM-PIM [158, 159] and Acceleration DIMM (AxDIMM) [160] are Samsung’s proposals that have been successfully tested via real prototypes. HBM-PIM features *Single Instruction Multiple Data* (SIMD) units, which support multiply-add and multiply-accumulate operations, near the banks in HBM layers [161, 162], and it is designed to accelerate neural network inference. AxDIMM is a near-rank solution that places an FPGA fabric on a DDR module to accelerate specific workloads (e.g., recommendation inference). Accelerator-in-Memory (AiM) [163] is a GDDR6-based PIM architecture from SK Hynix with specialized units for multiply-accumulate and activation functions for deep learning. HB-PNM [164] is a 3D-stacked-based PIM architecture from Alibaba, which stacks a layer of LPDDR4 memory and a logic layer with specialized accelerators for recommendation systems.

These five real-world PIM systems have some important common characteristics, as depicted in Figure 1. First, there is a host processor (CPU or GPU), typically with a deep cache hierarchy, which has access to (1) standard main memory, and (2) PIM-enabled memory (i.e., UPMEM DIMMs, HBM-PIM stacks, AxDIMM DIMMs, AiM GDDR6, HB-PNM LPDDR4). Second, the PIM-enabled memory chip contains multiple PIM processing elements (PIM PEs), which have access to memory (either memory banks or ranks) with higher

bandwidth and lower latency than the host processor. Third, the PIM processing elements (either general-purpose cores, SIMD units, FPGAs, or specialized processors) run at only a few hundred megahertz, and have a small number of registers and relatively small (or no) cache or scratchpad memory. Fourth, processing elements may not be able to communicate directly with each other (e.g., UPMEM DPUs, HBM-PIM PCUs or AiM PUs in different chips), and communication between them happens via the host processor. Figure 1 shows a high-level view of such a state-of-the-art processing-in-memory system.

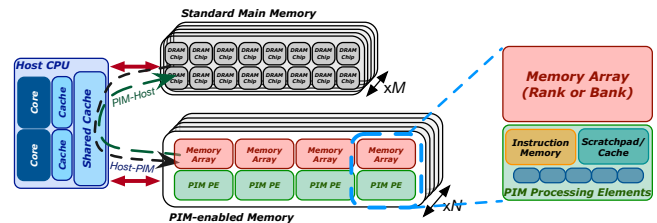


Figure 1: High-level view of a state-of-the-art processing-in-memory system. The host CPU has access to M standard memory modules and N PIM-enabled memory modules.

Our goal in this work is to quantify the potential of general-purpose PIM architectures for training of machine learning algorithms. To this end, we implement four representative classical machine learning algorithms (linear regression [165, 166], logistic regression [165, 167], decision tree [168], K-means clustering [169]) on a general-purpose memory-centric system containing PIM-enabled memory, specifically the UPMEM PIM architecture [153–157]. We do *not* include training of deep learning algorithms in our study, since GPUs and TPUs have a solid position as the preferred and highly optimized accelerators for deep learning training [89, 170–175].

Our PIM implementations of ML algorithms follow PIM programming recommendations in recent literature [154–156, 176]. We apply several optimizations to overcome the limitations of existing general-purpose PIM architectures (e.g., limited instruction set, relatively simple pipeline, relatively low frequency) and take full advantage of the inherent strengths of PIM (e.g., large memory bandwidth, low memory latency).

We evaluate our PIM implementations in terms of training accuracy, performance, and scaling characteristics on a real memory-centric system with PIM-enabled memory [153, 176, 177]. We run our experiments on a real-world PIM system [153] with 2,524 PIM cores running at 425 MHz, and 158 GB of DRAM memory.¹

¹The UPMEM-based PIM system has up to 2,560 PIM cores and 160 GB of DRAM.

Our experimental real system evaluation provides new observations and insights, including the following:

- ML training workloads that show memory-bound behavior in processor-centric systems can greatly benefit from (1) fixed-point data representation, (2) quantization [178, 179], and (3) hybrid precision implementation [163, 180] (without much accuracy loss) in PIM systems, in order to alleviate the lack of native support for floating-point and high-precision (i.e., 32- and 64-bit) arithmetic operations.
- ML training workloads that require complex activation functions (e.g., sigmoid) [181] can take advantage of *lookup tables (LUTs)* [98, 182, 183] in PIM systems instead of function approximation (e.g., Taylor series) [184], when PIM systems lack native support for those activation functions.
- Data can be placed and laid out such that accesses of PIM cores to their nearby memory banks are streaming, which enables better exploitation of the PIM memory bandwidth.
- ML training workloads with large training datasets can greatly benefit from scaling the size of PIM-enabled memory with PIM cores attached to memory banks. Training datasets can remain in memory without being moved to the host processor (e.g., CPU, GPU) in every iteration of the training process. Even if PIM cores need to communicate intermediate results via the host processor, this communication overhead is tolerable with proper overlap of computation and communication.

We compare our PIM implementations of linear regression, logistic regression, decision tree, and K-means clustering to their state-of-the-art CPU and GPU counterparts. We observe that memory-centric systems with PIM-enabled memory can significantly outperform processor-centric systems for memory-bound ML training workloads, when the operations needed by the ML workloads are natively supported by PIM hardware (or can be replaced by efficient LUT implementations).

Our extended paper [185] contains (1) detailed description of our PIM implementations of ML workloads; (2) comprehensive evaluation and comparisons to state-of-the-art CPU and GPU systems; and (3) more insights about the suitability of ML workloads to the PIM system, programming recommendations for ML software developers, and suggestions and hints for future PIM architectures. We aim to open-source all our PIM implementations of ML training workloads, training datasets, and evaluation scripts.

KEYWORDS

machine learning, processing-in-memory, regression, classification, clustering, benchmarking

ACKNOWLEDGMENTS

We acknowledge the generous gifts provided by our industrial partners, including ASML, Facebook, Google, Huawei, Intel, Microsoft, and VMware. We acknowledge support from the Semiconductor Research Corporation and the ETH Future Computing Laboratory.

This extended abstract appears as an invited paper at the 2022 IEEE Computer Society Annual Symposium on VLSI (ISVLSI). It is a summary version of our recent work [185].

REFERENCES

- [1] A. Géron, *Hands-on Machine Learning With Scikit-Learn, Keras, and TensorFlow: Concepts, tools, and techniques to build intelligent systems*, 2019.
- [2] E. Alpaydin, *Introduction to Machine Learning*, 2020.
- [3] I. Goodfellow *et al.*, *Deep Learning*, 2016.
- [4] M. Mohri *et al.*, *Foundations of Machine Learning*, 2018.
- [5] S. Shalev-Shwartz and S. Ben-David, *Understanding Machine Learning: From Theory to Algorithms*, 2014.
- [6] S. Raschka and V. Mirjalili, *Python Machine Learning: Machine Learning and Deep Learning with Python, Scikit-Learn, and TensorFlow 2*, 2019.
- [7] O. Mutlu *et al.*, “Processing Data Where It Makes Sense: Enabling In-Memory Computation,” *MicPro*, 2019.
- [8] O. Mutlu *et al.*, “A Modern Primer on Processing in Memory,” *Emerging Computing: From Devices to Systems - Looking Beyond Moore and Von Neumann*, 2021, <https://arxiv.org/pdf/2012.03112.pdf>.
- [9] S. Ghose *et al.*, “Processing-in-Memory: A Workload-Driven Perspective,” *IBM JRD*, 2019.
- [10] V. Seshadri and O. Mutlu, “In-DRAM Bulk Bitwise Execution Engine,” arXiv:1905.09822 [cs.AR], 2020.
- [11] O. Mutlu *et al.*, “Enabling Practical Processing in and near Memory for Data-Intensive Computing,” in *DAC*, 2019.
- [12] H. S. Stone, “A Logic-in-Memory Computer,” *IEEE TC*, 1970.
- [13] W. H. Kautz, “Cellular Logic-in-Memory Arrays,” *IEEE TC*, 1969.
- [14] D. E. Shaw *et al.*, “The NON-VON Database Machine: A Brief Overview,” *IEEE Database Eng. Bull.*, 1981.
- [15] P. M. Kogge, “EXECUBE - A New Architecture for Scalable MPPs,” in *ICPP*, 1994.
- [16] M. Gokhale *et al.*, “Processing in Memory: The Terasys Massively Parallel PIM Array,” *IEEE Computer*, 1995.
- [17] D. Patterson *et al.*, “A Case for Intelligent RAM,” *IEEE Micro*, 1997.
- [18] M. Oskin *et al.*, “Active Pages: A Computation Model for Intelligent Memory,” in *ISCA*, 1998.
- [19] Y. Kang *et al.*, “FlexRAM: Toward an Advanced Intelligent Memory System,” in *ICCD*, 1999.
- [20] K. Mai *et al.*, “Smart Memories: A Modular Reconfigurable Architecture,” in *ISCA*, 2000.
- [21] R. C. Murphy *et al.*, “The Characterization of Data Intensive Memory Workloads on Distributed PIM Systems,” in *Intelligent Memory Systems*. Springer.
- [22] J. Draper *et al.*, “The Architecture of the DIVA Processing-in-Memory Chip,” in *SC*, 2002.
- [23] S. Aga *et al.*, “Compute Caches,” in *HPCA*, 2017.
- [24] C. Eckert *et al.*, “Neural Cache: Bit-serial In-cache Acceleration of Deep Neural Networks,” in *ISCA*, 2018.
- [25] D. Fujiki *et al.*, “Duality Cache for Data Parallel Acceleration,” in *ISCA*, 2019.
- [26] M. Kang *et al.*, “An Energy-Efficient VLSI Architecture for Pattern Recognition via Deep Embedding of Computation in SRAM,” in *ICASSP*, 2014.
- [27] V. Seshadri *et al.*, “Ambit: In-Memory Accelerator for Bulk Bitwise Operations Using Commodity DRAM Technology,” in *MICRO*, 2017.
- [28] V. Seshadri *et al.*, “Buddy-RAM: Improving the Performance and Efficiency of Bulk Bitwise Operations Using DRAM,” arXiv:1611.09988 [cs.AR], 2016.
- [29] V. Seshadri *et al.*, “Fast Bulk Bitwise AND and OR in DRAM,” *CAL*, 2015.
- [30] V. Seshadri *et al.*, “RowClone: Fast and Energy-Efficient In-DRAM Bulk Data Copy and Initialization,” in *MICRO*, 2013.
- [31] S. Angizi and D. Fan, “Graphide: A Graph Processing Accelerator Leveraging In-DRAM-computing,” in *GLSVLSI*, 2019.
- [32] J. Kim *et al.*, “The DRAM Latency PUF: Quickly Evaluating Physical Unclonable Functions by Exploiting the Latency-Reliability Tradeoff in Modern DRAM Devices,” in *HPCA*, 2018.
- [33] J. Kim *et al.*, “D-RaNGe: Using Commodity DRAM Devices to Generate True Random Numbers with Low Latency and High Throughput,” in *HPCA*, 2019.
- [34] F. Gao *et al.*, “ComputeDRAM: In-Memory Compute Using Off-the-Shelf DRAMs,” in *MICRO*, 2019.
- [35] K. K. Chang *et al.*, “Low-Cost Inter-Linked Subarrays (LISA): Enabling Fast Inter-Subarray Data Movement in DRAM,” in *HPCA*, 2016.
- [36] X. Xin *et al.*, “ELP2IM: Efficient and Low Power Bitwise Operation Processing in DRAM,” in *HPCA*, 2020.
- [37] S. Li *et al.*, “DRISA: A DRAM-Based Reconfigurable In-Situ Accelerator,” in *MICRO*, 2017.
- [38] Q. Deng *et al.*, “DrAcc: A DRAM Based Accelerator for Accurate CNN Inference,” in *DAC*, 2018.
- [39] N. Hajinazar *et al.*, “SIMDRAM: A Framework for Bit-Serial SIMD Processing Using DRAM,” in *ASPLOS*, 2021.
- [40] S. H. S. Rezaei *et al.*, “NoM: Network-on-Memory for Inter-Bank Data Transfer in Highly-Banked Memories,” *CAL*, 2020.
- [41] Y. Wang *et al.*, “FIGARO: Improving System Performance via Fine-Grained In-DRAM Data Relocation and Caching,” in *MICRO*, 2020.
- [42] M. F. Ali *et al.*, “In-Memory Low-Cost Bit-Serial Addition Using Commodity DRAM Technology,” in *TCAS-I*, 2019.

- [43] S. Li *et al.*, "Pinatubo: A Processing-in-Memory Architecture for Bulk Bitwise Operations in Emerging Non-Volatile Memories," in *DAC*, 2016.
- [44] S. Angizi *et al.*, "PIMA-Logic: A Novel Processing-in-Memory Architecture for Highly Flexible and Energy-efficient Logic Computation," in *DAC*, 2018.
- [45] S. Angizi *et al.*, "CMP-PIM: An Energy-efficient Comparator-based Processing-in-Memory Neural Network Accelerator," in *DAC*, 2018.
- [46] S. Angizi *et al.*, "AlignS: A Processing-in-Memory Accelerator for DNA Short Read Alignment Leveraging SOT-MRAM," in *DAC*, 2019.
- [47] Y. Levy *et al.*, "Logic Operations in Memory Using a Memristive Akers Array," *Microelectronics Journal*, 2014.
- [48] S. Kvatinisky *et al.*, "MAGIC—Memristor-Aided Logic," *IEEE TCAS II: Express Briefs*, 2014.
- [49] A. Shafiee *et al.*, "ISAAC: A Convolutional Neural Network Accelerator with In-situ Analog Arithmetic in Crossbars," in *ISCA*, 2016.
- [50] S. Kvatinisky *et al.*, "Memristor-Based IMPLY Logic Design Procedure," in *ICCD*, 2011.
- [51] S. Kvatinisky *et al.*, "Memristor-Based Material Implication (IMPLY) Logic: Design Principles and Methodologies," *TVLSI*, 2014.
- [52] P.-E. Gaillardon *et al.*, "The Programmable Logic-in-Memory (PLiM) Computer," in *DATE*, 2016.
- [53] D. Bhattacharjee *et al.*, "ReVAMP: ReRAM based VLIW Architecture for In-memory Computing," in *DATE*, 2017.
- [54] S. Hamdioui *et al.*, "Memristor Based Computation-in-Memory Architecture for Data-intensive Applications," in *DATE*, 2015.
- [55] L. Xie *et al.*, "Fast Boolean Logic Papped on Memristor Crossbar," in *ICCD*, 2015.
- [56] S. Hamdioui *et al.*, "Memristor for Computing: Myth or Reality?" in *DATE*, 2017.
- [57] J. Yu *et al.*, "Memristive Devices for Computation-in-Memory," in *DATE*, 2018.
- [58] C. Giannoula *et al.*, "SynCron: Efficient Synchronization Support for Near-Data-Processing Architectures," in *HPCA*, 2021.
- [59] I. Fernandez *et al.*, "NATSA: A Near-Data Processing Accelerator for Time Series Analysis," in *ICCD*, 2020.
- [60] D. S. Cali *et al.*, "GenASM: A High-Performance, Low-Power Approximate String Matching Acceleration Framework for Genome Sequence Analysis," in *MICRO*, 2020.
- [61] J. S. Kim *et al.*, "GRIM-Filter: Fast Seed Location Filtering in DNA Read Mapping Using Processing-in-Memory Technologies," *BMC Genomics*, 2018.
- [62] J. Ahn *et al.*, "PIM-Enabled Instructions: A Low-Overhead, Locality-Aware Processing-in-Memory Architecture," in *ISCA*, 2015.
- [63] J. Ahn *et al.*, "A Scalable Processing-in-Memory Accelerator for Parallel Graph Processing," in *ISCA*, 2015.
- [64] A. Boroumand *et al.*, "Google Workloads for Consumer Devices: Mitigating Data Movement Bottlenecks," in *ASPLOS*, 2018.
- [65] A. Boroumand *et al.*, "CoNDA: Efficient Cache Coherence Support for near-Data Accelerators," in *ISCA*, 2019.
- [66] G. Singh *et al.*, "NAPEL: Near-memory Computing Application Performance Prediction via Ensemble Learning," in *DAC*, 2019.
- [67] H. Asghari-Moghaddam *et al.*, "Chameleon: Versatile and Practical Near-DRAM Acceleration Architecture for Large Memory Systems," in *MICRO*, 2016.
- [68] O. O. Babarinsa and S. Idreos, "JAFAR: Near-Data Processing for Databases," in *SIGMOD*, 2015.
- [69] P. Chi *et al.*, "PRIME: A Novel Processing-In-Memory Architecture for Neural Network Computation In ReRAM-Based Main Memory," in *ISCA*, 2016.
- [70] A. Farmahini-Farahani *et al.*, "NDA: Near-DRAM acceleration architecture leveraging commodity DRAM devices and standard memory modules," in *HPCA*, 2015.
- [71] M. Gao *et al.*, "Practical Near-Data Processing for In-Memory Analytics Frameworks," in *PACT*, 2015.
- [72] M. Gao and C. Kozyrakis, "HRL: Efficient and Flexible Reconfigurable Logic for Near-Data Processing," in *HPCA*, 2016.
- [73] B. Gu *et al.*, "Biscuit: A Framework for Near-Data Processing of Big Data Workloads," in *ISCA*, 2016.
- [74] Q. Guo *et al.*, "3D-Stacked Memory-Side Acceleration: Accelerator and System Design," in *WoNDP*, 2014.
- [75] M. Hashemi *et al.*, "Accelerating Dependent Cache Misses with an Enhanced Memory Controller," in *ISCA*, 2016.
- [76] M. Hashemi *et al.*, "Continuous Runahead: Transparent Hardware Acceleration for Memory Intensive Workloads," in *MICRO*, 2016.
- [77] K. Hsieh *et al.*, "Transparent Offloading and Mapping (TOM): Enabling Programmer-Transparent Near-Data Processing in GPU Systems," in *ISCA*, 2016.
- [78] D. Kim *et al.*, "Neurocube: A Programmable Digital Neuromorphic Architecture with High-Density 3D Memory," in *ISCA*, 2016.
- [79] G. Kim *et al.*, "Toward Standardized Near-Data Processing with Unrestricted Data Placement for GPUs," in *SC*, 2017.
- [80] J. H. Lee *et al.*, "BSSync: Processing Near Memory for Machine Learning Workloads with Bounded Staleness Consistency Models," in *PACT*, 2015.
- [81] Z. Liu *et al.*, "Concurrent Data Structures for Near-Memory Computing," in *SPAA*, 2017.
- [82] A. Morad *et al.*, "GP-SIMD Processing-in-Memory," *ACM TACO*, 2015.
- [83] L. Nai *et al.*, "GraphPIM: Enabling Instruction-Level PIM Offloading in Graph Computing Frameworks," in *HPCA*, 2017.
- [84] A. Pattnaik *et al.*, "Scheduling Techniques for GPU Architectures with Processing-in-Memory Capabilities," in *PACT*, 2016.
- [85] S. H. Pugsley *et al.*, "NDC: Analyzing the Impact of 3D-Stacked Memory+Logic Devices on MapReduce Workloads," in *ISPASS*, 2014.
- [86] D. P. Zhang *et al.*, "TOP-PIM: Throughput-Oriented Programmable Processing in Memory," in *HPDC*, 2014.
- [87] Q. Zhu *et al.*, "Accelerating Sparse Matrix-Matrix Multiplication with 3D-Stacked Logic-in-Memory Hardware," in *HPEC*, 2013.
- [88] B. Akin *et al.*, "Data Reorganization in Memory Using 3D-Stacked DRAM," in *ISCA*, 2015.
- [89] M. Gao *et al.*, "Tetris: Scalable and Efficient Neural Network Acceleration with 3D Memory," in *ASPLOS*, 2017.
- [90] M. Drummond *et al.*, "The Mondrian Data Engine," in *ISCA*, 2017.
- [91] G. Dai *et al.*, "GraphH: A Processing-in-Memory Architecture for Large-scale Graph Processing," *IEEE TCAD*, 2018.
- [92] M. Zhang *et al.*, "GraphP: Reducing Communication for PIM-based Graph Processing with Efficient Data Partition," in *HPCA*, 2018.
- [93] Y. Huang *et al.*, "A Heterogeneous PIM Hardware-Software Co-Design for Energy-Efficient Graph Processing," in *IPDPS*, 2020.
- [94] Y. Zhuo *et al.*, "GraphQ: Scalable PIM-based Graph Processing," in *MICRO*, 2019.
- [95] P. C. Santos *et al.*, "Operand Size Reconfiguration for Big Data Processing in Memory," in *DATE*, 2017.
- [96] W.-M. Hwu *et al.*, "Rebooting the Data Access Hierarchy of Computing Systems," in *ICRC*, 2017.
- [97] M. Besta *et al.*, "SISA: Set-Centric Instruction Set Architecture for Graph Mining on Processing-in-Memory Systems," in *MICRO*, 2021.
- [98] J. D. Ferreira *et al.*, "pLUTO: In-DRAM Lookup Tables to Enable Massively Parallel General-Purpose Computation," *arXiv:2104.07699 [cs.AR]*, 2021.
- [99] A. Olgun *et al.*, "QUAC-TRNG: High-Throughput True Random Number Generation Using Quadruple Row Activation in Commodity DRAMs," in *ISCA*, 2021.
- [100] S. Lloyd and M. Gokhale, "In-memory Data Rearrangement for Irregular, Data-intensive Computing," *Computer*, 2015.
- [101] D. G. Elliott *et al.*, "Computational RAM: Implementing Processors in Memory," *IEEE Design & Test of Computers*, 1999.
- [102] L. Zheng *et al.*, "RRAM-based TCAMs for pattern search," in *ISCAS*, 2016.
- [103] J. Landgraf *et al.*, "Combining Emulation and Simulation to Evaluate a Near Memory Key/Value Lookup Accelerator," 2021.
- [104] A. Rodrigues *et al.*, "Towards a Scatter-Gather Architecture: Hardware and Software Issues," in *MEMSYS*, 2019.
- [105] S. Lloyd and M. Gokhale, "Design Space Exploration of Near Memory Accelerators," in *MEMSYS*, 2018.
- [106] S. Lloyd and M. Gokhale, "Near Memory Key/Value Lookup Acceleration," in *MEMSYS*, 2017.
- [107] M. Gokhale *et al.*, "Near Memory Data Structure Rearrangement," in *MEMSYS*, 2015.
- [108] R. Nair *et al.*, "Active Memory Cube: A Processing-in-Memory Architecture for Exascale Systems," *IBM JRD*, 2015.
- [109] A. C. Jacob *et al.*, "Compiling for the Active Memory Cube," Tech. rep. RC25644 (WAT1612-008). IBM Research Division, Tech. Rep., 2016.
- [110] Z. Sura *et al.*, "Data Access Optimization in a Processing-in-Memory System," in *CF*, 2015.
- [111] R. Nair, "Evolution of Memory Architecture," *Proceedings of the IEEE*, 2015.
- [112] R. Balasubramonian *et al.*, "Near-Data Processing: Insights from a MICRO-46 Workshop," *IEEE Micro*, 2014.
- [113] Y. Xi *et al.*, "In-Memory Learning With Analog Resistive Switching Memory: A Review and Perspective," *Proceedings of the IEEE*, 2020.
- [114] K. Hsieh *et al.*, "Accelerating Pointer Chasing in 3D-Stacked Memory: Challenges, Mechanisms, Evaluation," in *ICCD*, 2016.
- [115] A. Boroumand *et al.*, "LazyPIM: An Efficient Cache Coherence Mechanism for Processing-in-Memory," *CAL*, 2016.
- [116] C. Giannoula *et al.*, "SparseP: Towards Efficient Sparse Matrix Vector Multiplication on Real Processing-In-Memory Systems," *arXiv preprint arXiv:2201.05072*, 2022.
- [117] C. Giannoula *et al.*, "Towards Efficient Sparse Matrix Vector Multiplication on Real Processing-in-Memory Architectures," in *SIGMETRICS*, 2022.
- [118] A. Denzler *et al.*, "Casper: Accelerating stencil computation using near-cache processing," *arXiv preprint arXiv:2112.14216*, 2021.
- [119] A. Boroumand *et al.*, "Polynesia: Enabling Effective Hybrid Transactional/Analytical Databases with Specialized Hardware/Software Co-Design," *arXiv:2103.00798 [cs.AR]*, 2021.
- [120] A. Boroumand *et al.*, "Polynesia: Enabling Effective Hybrid Transactional Analytical Databases with Specialized Hardware Software Co-Design," in *ICDE*, 2022.
- [121] G. Singh *et al.*, "FPGA-based Near-Memory Acceleration of Modern Data-Intensive Applications," *IEEE Micro*, 2021.

- [122] G. Singh *et al.*, "Accelerating Weather Prediction using Near-Memory Reconfigurable Fabric," *ACM TRETS*, 2021.
- [123] J. M. Herruzo *et al.*, "Enabling Fast and Energy-Efficient FM-Index Exact Matching Using Processing-Near-Memory," *The Journal of Supercomputing*, 2021.
- [124] L. Yavits *et al.*, "GIRAF: General Purpose In-Storage Resistive Associative Framework," *IEEE TPDS*, 2021.
- [125] B. Asgari *et al.*, "FAFNIR: Accelerating Sparse Gathering by Using Efficient Near-Memory Intelligent Reduction," in *HPCA*, 2021.
- [126] A. Boroumand *et al.*, "Google Neural Network Models for Edge Devices: Analyzing and Mitigating Machine Learning Inference Bottlenecks," *arXiv preprint arXiv:2109.14320*, 2021.
- [127] A. Boroumand *et al.*, "Google Neural Network Models for Edge Devices: Analyzing and Mitigating Machine Learning Inference Bottlenecks," in *PACT*, 2021.
- [128] A. Boroumand, "Practical Mechanisms for Reducing Processor-Memory Data Movement in Modern Workloads," Ph.D. dissertation, Carnegie Mellon University, 2020.
- [129] G. Singh *et al.*, "NERO: A Near High-Bandwidth Memory Stencil Accelerator for Weather Prediction Modeling," in *FPL*, 2020.
- [130] V. Seshadri and O. Mutlu, "Simple Operations in Memory to Reduce Data Movement," in *Advances in Computers, Volume 106*, 2017.
- [131] S. Diab *et al.*, "High-throughput Pairwise Alignment with the Wavefront Algorithm using Processing-in-Memory," *arXiv preprint arXiv:2204.02085*, 2022.
- [132] S. Diab *et al.*, "High-throughput Pairwise Alignment with the Wavefront Algorithm using Processing-in-Memory," in *HICOMB*, 2022.
- [133] D. Fujiki *et al.*, "In-Memory Data Parallel Processor," in *ASPLOS*, 2018.
- [134] Y. Zha and J. Li, "Hyper-AP: Enhancing Associative Processing Through A Full-Stack Optimization," in *ISCA*, 2020.
- [135] O. Mutlu, "Memory Scaling: A Systems Architecture Perspective," *IMW*, 2013.
- [136] O. Mutlu and L. Subramanian, "Research Problems and Opportunities in Memory Systems," *SUPERFRI*, 2014.
- [137] H. Ahmed *et al.*, "A Compiler for Automatic Selection of Suitable Processing-in-Memory Instructions," in *DATE*, 2019.
- [138] S. Jain *et al.*, "Computing-in-Memory with Spintronics," in *DATE*, 2018.
- [139] N. M. Ghiassi *et al.*, "GenStore: A High-Performance and Energy-Efficient In-Storage Computing System for Genome Sequence Analysis," in *ASPLOS*, 2022.
- [140] G. F. Oliveira *et al.*, "DAMOV: A New Methodology and Benchmark Suite for Evaluating Data Movement Bottlenecks," *IEEE Access*, 2021.
- [141] G. F. Oliveira *et al.*, "DAMOV: A New Methodology and Benchmark Suite for Evaluating Data Movement Bottlenecks," *arXiv:2105.03725 [cs.AR]*, 2021.
- [142] S. Cho *et al.*, "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.
- [143] H. Shin *et al.*, "McDRAM: Low latency and energy-efficient matrix computations in DRAM," *IEEE TCADICS*, 2018.
- [144] P. Gu *et al.*, "iPIM: Programmable In-Memory Image Processing Accelerator using Near-Bank Architecture," in *ISCA*, 2020.
- [145] D. Lavenier *et al.*, "Variant Calling Parallelization on Processor-in-Memory Architecture," in *BIBM*, 2020.
- [146] V. Zois *et al.*, "Massively Parallel Skyline Computation for Processing-in-Memory Architectures," in *PACT*, 2018.
- [147] D. Weber *et al.*, "Current and Future Challenges of DRAM Metallization," in *IITC*, 2005.
- [148] Y. Peng *et al.*, "Design, Packaging, and Architectural Policy Co-optimization for DC Power Integrity in 3D DRAM," in *DAC*, 2015.
- [149] F. Devaux, "The True Processing In Memory Accelerator," in *Hot Chips*, 2019.
- [150] M. Yuffe *et al.*, "A Fully Integrated Multi-CPU, GPU and Memory Controller 32nm processor," in *ISSCC*, 2011.
- [151] R. Christy *et al.*, "8.3 A 3GHz ARM Neoverse N1 CPU in 7nm FinFET for Infrastructure Applications," in *ISSCC*, 2020.
- [152] T. Singh *et al.*, "3.2 Zen: A Next-generation High-performance x86 Core," in *ISSCC*, 2017.
- [153] UPMEM, "UPMEM Website," <https://www.upmem.com>, 2020.
- [154] UPMEM, "Introduction to UPMEM PIM. Processing-in-memory (PIM) on DRAM Accelerator (White Paper)," 2018.
- [155] J. Gómez-Luna *et al.*, "Benchmarking a New Paradigm: An Experimental Analysis of a Real Processing-in-Memory Architecture," *arXiv:2105.03814 [cs.AR]*, 2021.
- [156] J. Gómez-Luna *et al.*, "Benchmarking a New Paradigm: Experimental Analysis and Characterization of a Real Processing-in-Memory System," *IEEE Access*, 2022.
- [157] J. Gómez-Luna *et al.*, "Benchmarking Memory-Centric Computing Systems: Analysis of Real Processing-In-Memory Hardware," in *IGSC*, 2021.
- [158] Y.-C. Kwon *et al.*, "25.4 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.
- [159] S. Lee *et al.*, "Hardware Architecture and Software Stack for PIM Based on Commercial DRAM Technology: Industrial Product," in *ISCA*, 2021.
- [160] L. Ke *et al.*, "Near-Memory Processing in Action: Accelerating Personalized Recommendation with AxDIMM," *IEEE Micro*, 2021.
- [161] JEDEC, "High Bandwidth Memory (HBM) DRAM," Standard No. JESD235, 2013.
- [162] D. Lee *et al.*, "Simultaneous Multi-Layer Access: Improving 3D-Stacked Memory Bandwidth at Low Cost," *TACO*, 2016.
- [163] S. Lee *et al.*, "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*, 2022.
- [164] D. Niu *et al.*, "184QPS/W 64Mb/mm2 3D Logic-to-DRAM Hybrid Bonding with Process-Near-Memory Engine for Recommendation System," in *ISSCC*, 2022.
- [165] D. A. Freedman, *Statistical Models: Theory and Practice*, 2009.
- [166] X. Yan and X. Su, *Linear Regression Analysis: Theory and Computing*, 2009.
- [167] D. W. Hosmer Jr *et al.*, *Applied Logistic Regression*, 2013.
- [168] S. Suthaharan, "Decision Tree Learning," in *Machine Learning Models and Algorithms for Big Data Classification*, 2016.
- [169] S. P. Lloyd, "Least Squares Quantization in PCM," *IEEE Transactions on Information Theory*, 1982.
- [170] D. B. Kirk *et al.*, *Programming Massively Parallel Processors, 3rd Edition, Chapter 16 - Application Case Study: Machine Learning*. Morgan Kaufmann, 2017.
- [171] N. P. Jouppi *et al.*, "cuDNN: Efficient Primitives for Deep Learning," *arXiv preprint arXiv:1410.0759*, 2014.
- [172] M. Abadi *et al.*, "Tensorflow: A System for Large-scale Machine Learning," in *OSDI*, 2016.
- [173] Run:AI, "Best GPU for Deep Learning," <https://www.run.ai/guides/gpu-deep-learning/best-gpu-for-deep-learning/>, 2021.
- [174] N. P. Jouppi *et al.*, "In-Datacenter Performance Analysis of a Tensor Processing Unit," in *ISCA*, 2017.
- [175] N. P. Jouppi *et al.*, "Ten Lessons from Three Generations Shaped Google's TPUV4i: Industrial Product," in *ISCA*, 2021.
- [176] UPMEM, "UPMEM User Manual. Version 2021.3.0," 2021.
- [177] UPMEM, "UPMEM Software Development Kit (SDK)." <https://sdk.upmem.com>, 2021.
- [178] N. Zmora *et al.*, "Achieving FP32 Accuracy for INT8 Inference Using Quantization Aware Training with NVIDIA TensorRT," <https://developer.nvidia.com/blog/achieving-fp32-accuracy-for-int8-inference-using-quantization-aware-training-with-tensorrt/>.
- [179] A. Gholami *et al.*, "A Survey of Quantization Methods for Efficient Neural Network Inference," in *Low-Power Computer Vision*.
- [180] NVIDIA, "NVIDIA H100 Tensor Core GPU Architecture. White Paper," <https://nvdam.widen.net/s/9bz6dw7dqr/gtc22-whitepaper-hopper>, 2022.
- [181] J. Han and C. Moraga, "The Influence of the Sigmoid Function Parameters on the Speed of Backpropagation Learning," in *IWANN*, 1995.
- [182] Q. Deng *et al.*, "LAcc: Exploiting Lookup Table-based Fast and Accurate Vector Multiplication in DRAM-based CNN Accelerator," in *DAC*, 2019.
- [183] M. Gao *et al.*, "DRAF: A Low-power DRAM-based Reconfigurable Acceleration Fabric," in *ISCA*, 2016.
- [184] E. W. Weisstein, "Taylor Series," <https://mathworld.wolfram.com/TaylorSeries.html>, 2004.
- [185] J. Gómez-Luna *et al.*, "An Experimental Evaluation of Machine Learning Training on a Real Processing-in-Memory System," *arXiv preprint arXiv:2207.07886*, 2022.