

Enabling Practical Processing in and near Memory for Data-Intensive Computing

Onur Mutlu^{a,b} Saugata Ghose^b Juan Gómez-Luna^a Rachata Ausavarungnirun^b

^aETH Zürich ^bCarnegie Mellon University

ABSTRACT

Modern computing systems suffer from the dichotomy between computation on one side, which is performed only in the processor (and accelerators), and data storage/movement on the other, which all other parts of the system are dedicated to. Due to this dichotomy, data moves a lot in order for the system to perform computation on it. Unfortunately, data movement is extremely expensive in terms of energy and latency, much more so than computation. As a result, a large fraction of system energy is spent and performance is lost solely on moving data in a modern computing system.

In this work, we re-examine the idea of reducing data movement by performing Processing in Memory (PIM). PIM places computation mechanisms in or near where the data is stored (i.e., inside the memory chips, in the logic layer of 3D-stacked logic and DRAM, or in the memory controllers), so that data movement between the computation units and memory is reduced or eliminated. While the idea of PIM is not new, we examine two new approaches to enabling PIM: 1) exploiting analog properties of DRAM to perform massively-parallel operations in memory, and 2) exploiting 3D-stacked memory technology design to provide high bandwidth to in-memory logic. We conclude by discussing work on solving key challenges to the practical adoption of PIM.

1 INTRODUCTION

Main memory, built using Dynamic Random Access Memory (DRAM), is a major component in nearly all computing systems, including servers, cloud platforms, mobile/embedded devices, and sensors. Across these systems, the data working set sizes of applications are rapidly growing, while the need for fast analysis of such data is increasing. Thus, main memory is becoming an increasingly significant bottleneck across a wide variety of computing systems and applications [9, 38, 71, 75]. The bottleneck has worsened in recent years, as it has become increasingly difficult to efficiently scale memory capacity, energy, cost, and performance across technology generations [41, 48, 49, 63, 64, 68, 69, 71, 72, 75], as evidenced by the RowHammer problem [48, 72, 74] in recent DRAM chips.

A major reason for the main memory bottleneck is the high energy and latency associated with *data movement*. In today's computers, to perform any operation on data, the processor must retrieve the data from main memory. This requires the memory controller to issue commands to a DRAM module across a relatively slow and power-hungry off-chip bus (known as the *memory channel*). The DRAM module sends the requested data across the memory channel, after which the data is placed in the caches and registers. The CPU can perform computation on the data once the data is in its registers. Data movement from the DRAM to the CPU incurs long latency and consumes significant energy [3, 4, 9, 29, 30]. These costs are often exacerbated by the fact that much of the data

brought into the caches is *not reused* by the CPU [82, 84], providing little benefit in return for the high latency and energy cost.

The cost of data movement is a fundamental issue with the *processor-centric* nature of contemporary computer systems. The CPU is considered the master in the system, and computation is performed only in the processor (and accelerators). In contrast, data storage and communication units, including the main memory, are treated as unintelligent workers that are incapable of computation. As a result of this processor-centric design paradigm, data moves a lot in the system between the computation units and communication/storage units so that computation can be done on it. With the increasingly *data-centric* nature of contemporary and emerging applications, the processor-centric design paradigm leads to great inefficiency in performance, energy and cost: for example, most of the real estate within a single compute node is already dedicated to handling data movement and storage (e.g., large caches, memory controllers, interconnects, and main memory), and our recent work shows that 62% of the entire system energy of a mobile device is spent on data movement between the processor and the memory hierarchy for widely-used mobile workloads [9].

The huge overhead of data movement in modern systems along with technology advances that enable better integration of memory and logic have recently prompted the re-examination of an old idea that we will generally call *Processing in Memory* (PIM). The key idea is to place computation mechanisms in or near where the data is stored (i.e., inside the memory chips, in the logic layer of 3D-stacked DRAM, in the memory controllers, or inside large caches), so that data movement between where the computation is done and where the data is stored is reduced or eliminated, compared to contemporary processor-centric systems.

The idea of PIM has been around for at least four decades [1, 16–18, 24, 42–44, 52, 70, 78, 79, 85, 95, 96]. However, past efforts were *not* widely adopted for various reasons, including 1) the difficulty of integrating processing elements with DRAM, 2) the lack of critical memory-related scaling challenges that current technology and applications face today, and 3) that the data movement bottleneck was not as critical to system cost, energy and performance as it is today. We believe it is crucial to re-examine PIM today with a fresh perspective (i.e., with novel approaches and ideas), by exploiting new memory technologies, with realistic workloads and systems, and with a mindset to ease adoption and feasibility.

In this paper, we explore two new approaches to enabling PIM in modern systems. The first approach only *minimally changes memory chips* to perform simple yet powerful common operations that the chip is inherently efficient at performing [12, 13, 15, 22, 23, 60, 73, 87–93]. Such solutions take advantage of the existing memory design to perform *bulk operations* (i.e., operations on an entire

row of DRAM cells), such as bulk copy, data initialization, and bitwise operations [13, 88–91]. The second approach enables PIM in a more general-purpose manner by taking advantage of emerging *3D-stacked memory technologies* [3–5, 8–11, 14, 19–21, 26, 27, 31–33, 45–47, 59, 65, 66, 76, 80, 81, 97, 102, 104]. 3D-stacked memory chips have much greater *internal* bandwidth than is available externally on the memory channel [58], and many such chip architectures (e.g., Hybrid Memory Cube [34, 35], High-Bandwidth Memory [37, 58]) include a *logic layer* where designers can add some processing logic (e.g., accelerators, simple cores, reconfigurable logic) that can take advantage of this high internal bandwidth.

Regardless of the approach taken to PIM, there are key practical adoption challenges that system architects and programmers must address to enable the widespread adoption of PIM across the computing landscape and in different domains of workloads. We also briefly discuss these challenges in this paper, along with references to some existing work that addresses these challenges.

2 MINIMALLY CHANGING MEMORY CHIPS

Minimal modifications in existing memory chips can enable simple yet powerful computation capability inside the chip. These modifications take advantage of the existing interconnects in and analog operational behavior of conventional memory chips, e.g., DRAM architectures, without the need for a logic layer and usually without the need for logic processing elements. As a result, the overheads imposed on the memory chip are low. There are a number of mechanisms that use this approach to take advantage of the high internal bandwidth available within each memory cell array [12, 13, 87–91, 93]. We briefly describe one such design, *Ambit*, which enables in-DRAM bulk bitwise operations [88, 90, 91], by building on RowClone, which enables fast and energy-efficient in-DRAM data movement [13, 89].

Ambit: In-DRAM Bulk Bitwise Operations. Many applications use *bulk bitwise operations* [51, 99] (i.e., bitwise operations on large bit vectors), such as bitmap indices, bitwise scan acceleration [62] for databases, accelerated document filtering for web search [25], DNA sequence alignment [6, 7, 47, 100], encryption algorithms [28, 98], graph processing, and networking [99]. Accelerating bulk bitwise operations can thus significantly boost the performance and energy efficiency of a wide range of applications.

We have recently proposed a new **Accelerator-in-Memory** for bulk **Bitwise** operations (*Ambit*) [88, 90, 91]. Unlike prior approaches, *Ambit* uses the analog operation of existing DRAM technology to perform bulk bitwise operations. *Ambit* has two components. The first component, *Ambit-AND-OR*, implements a new operation called *triple-row activation*, where the memory controller simultaneously activates three rows. Triple-row activation uses the charge sharing principles that govern the operation of the DRAM array to perform a bitwise AND or OR on two rows of data, by controlling the initial value on the third row. The second component, *Ambit-NOT*, takes advantage of the two inverters that are connected to each sense amplifier in a DRAM subarray, as the voltage level of one of the inverters represents the negated logical value of the cell. The *Ambit* design adds a special row to the DRAM array to capture this negated value. One possible implementation of the special row [91] is a row of *dual-contact cells* (a 2-transistor 1-capacitor cell [39, 67]), each connected to both inverters inside

a sense amplifier. Even in the presence of process variation (see [91]), *Ambit* can reliably perform AND, OR, and NOT operations completely using DRAM technology, making it functionally (i.e., Boolean logic) complete.

Ambit provides promising performance and energy improvements. Averaged across seven commonly-used bulk bitwise operations (NOT, AND, OR, NAND, NOR, XOR, XNOR), *Ambit* with 8 DRAM banks improves bulk bitwise operation throughput by 44× compared to an Intel Skylake processor [36], and 32× compared to the NVIDIA GTX 745 GPU [77]. Compared to DDR3 DRAM, *Ambit* reduces energy consumption by 35× on average. When integrated directly into the HMC 2.0 device, which has many more banks, *Ambit* improves operation throughput by 9.7× compared to processing in the logic layer of HMC 2.0. Our work evaluates the end-to-end benefits of *Ambit* on real database queries using Bitmap indices and the BitWeaving database [62], showing query latency reductions of 2X to 12X, with larger benefits for larger data set sizes.

A number of *Ambit*-like bitwise operation substrates have been proposed in recent years, making use of emerging resistive memory technologies, e.g., phase-change memory (PCM) [55–57, 83, 101, 103], SRAM, or specialized DRAM. These substrates can perform bulk bitwise operations in a special DRAM array augmented with computational circuitry [60] and in PCM [61]. Similar substrates can perform simple arithmetic operations in SRAM [2, 40] and arithmetic and logical operations in memristors [53, 54, 94]. Resistive memory technologies are amenable to in-place updates, and can thus incorporate *Ambit*-like operations with even less data movement than DRAM. Thus, we believe it is extremely important to continue exploring low-cost *Ambit*-like substrates, as well as more sophisticated computational substrates, for all types of memory technologies, old and new.

3 PIM USING 3D-STACKED MEMORY

Several works propose to place some form of processing logic (typically accelerators, simple cores, or reconfigurable logic) inside the logic layer of 3D-stacked memory [58]. This *PIM processing logic*, which we also refer to as *PIM cores*, can execute portions of applications (from individual instructions to functions) or entire threads and applications, depending on the design of the architecture. The PIM cores connect to the memory stacks that are on top of them using vertical *through-silicon vias* [58], which provide high-bandwidth and low-latency access to data. In this section, we discuss examples of how systems can make use of relatively simple PIM cores to avoid data movement and thus obtain significant performance and energy improvements for a variety of application domains.

Tesseract: Graph Processing. A popular modern application is large-scale graph processing/analytics. Graph processing has broad applicability and use in many domains, from social networks to machine learning, from data analytics to bioinformatics. Graph analysis workloads put large pressure on memory bandwidth due to 1) frequent random memory accesses across large memory regions (leading to limited cache efficiency and unnecessary data

transfer on the memory bus) and 2) small amount of computation per data item fetched from memory (leading to limited ability to hide long memory latencies and exercising the memory energy bottleneck). These two characteristics make it very challenging to scale up such workloads despite their inherent parallelism, especially with conventional architectures based on large on-chip caches and relatively scarce off-chip memory bandwidth for random access.

To overcome the limitations of conventional architectures, we design Tesseract, a programmable PIM accelerator for large-scale graph processing [3]. Tesseract consists of 1) simple in-order PIM cores that exploit the high memory bandwidth available in the logic layer of 3D-stacked memory, where each core manipulates data only on the memory partition it is assigned to control, 2) an efficient communication interface that allows a PIM core to request computation on data elements that reside in the memory partition controlled by another core, and 3) a message-passing based programming interface, similar to how modern distributed systems are programmed, which enables remote function calls on data that resides in each memory partition. Tesseract moves functions to data rather than moving data elements across different memory partitions and cores. Our comprehensive evaluations using five state-of-the-art graph processing workloads with large graphs show that Tesseract improves average system performance by 13.8× and reduces average system energy by 87% over a state-of-the-art conventional system.

Consumer Workloads. A popular domain of computing is consumer devices, including smartphones, tablets, web-based computers (e.g., Chromebooks), and wearable devices. In such devices, energy efficiency is a first-class concern due to the limited battery capacity and the stringent thermal power budget. We find that *data movement* is a major contributor to energy (and execution time) in modern consumer devices: across four popular workloads (described next), 62.7% of the total system energy, on average, is spent on data movement across the memory hierarchy [9].

We comprehensively analyze the energy and performance impact of data movement for several widely-used Google consumer workloads [9]: 1) the Chrome web browser, 2) TensorFlow Mobile (Google’s machine learning framework), 3) the VP9 video playback engine, and 4) the VP9 video capture engine. We find that offloading key functions (called *target functions*) of these workloads to PIM logic greatly reduces data movement. However, consumer devices are extremely stringent in terms of the extra area and energy they can accommodate. As a result, it is important to identify what kind of PIM logic can both 1) maximize energy efficiency and 2) be implemented at minimum possible area and energy costs.

We find that many of the target functions for PIM in consumer workloads are comprised of simple operations (e.g., *memcpy/memset*, basic arithmetic and bitwise operations), and can be implemented easily in the logic layer using either 1) a small low-power general-purpose core or 2) small fixed-function accelerators. Our analysis shows that the area of a PIM core and a PIM accelerator take up no more than 9.4% and 35.4%, respectively, of the area available for PIM logic in an HMC-like [35] 3D-stacked memory architecture. Both the PIM core and PIM accelerator eliminate a large amount of data movement, and thereby significantly reduce total

system energy (by an average of 55.4% across all the workloads) and execution time (by an average of 54.2%).

4 ENABLING PIM ADOPTION

Pushing computation from the CPU into memory introduces new challenges for system architects and programmers to overcome. Many of these challenges must be addressed for PIM to be adopted in a wide variety of systems of workloads, without placing a heavy burden on most programmers [22, 73]. These challenges include 1) how to easily program PIM systems (with good programming model, library, compiler and tools support) [4, 32]; 2) how to design runtime systems and system software that can take advantage of PIM (e.g., runtime scheduling of code on PIM logic, data mapping) [4, 9, 32, 80]; 3) how to efficiently enable coherence between PIM logic and CPU/accelerator cores that operate on shared data [4, 10, 11]; 4) how to efficiently enable virtual memory support on the PIM logic [33]; 5) how to design high-performance data structures for PIM whose performance is better than concurrent data structures on multi-core machines [65]; 6) how to accurately assess the benefits and shortcomings of PIM using realistic workload suites, rigorous analysis methodologies, and accurate and flexible simulation infrastructures [50, 86].

We believe these challenges provide exciting cross-layer research opportunities. Fundamentally solving the data movement problem requires a paradigm shift to a data-centric computing system design, where computation happens in or near memory, with minimal data movement. We argue that research enabled towards such a paradigm shift would be very useful for both PIM as well as other potential ideas that can reduce data movement.

ACKNOWLEDGMENTS

We thank members of the SAFARI Research Group and collaborators at Carnegie Mellon, ETH Zurich, and other universities, who have contributed to the various works we describe in this paper. Thanks also to our research group’s industrial sponsors over the past ten years, especially Alibaba, Google, Huawei, Intel, Microsoft, NVIDIA, Samsung, and VMware. This work was also partially supported by the Semiconductor Research Corporation and NSF.

REFERENCES

- [1] A. Acharya et al. 1998. Active Disks: Programming Model, Algorithms and Evaluation. In *ASPLOS*.
- [2] S. Aga et al. 2017. Compute Caches. In *HPCA*.
- [3] J. Ahn et al. 2015. A Scalable Processing-in-Memory Accelerator for Parallel Graph Processing. In *ISCA*.
- [4] J. Ahn et al. 2015. PIM-Enabled Instructions: A Low-Overhead, Locality-Aware Processing-in-Memory Architecture. In *ISCA*.
- [5] B. Akin et al. 2015. Data Reorganization in Memory Using 3D-Stacked DRAM. In *ISCA*.
- [6] M. Alser et al. Shouji: A Fast and Efficient Pre-Alignment Filter for Sequence Alignment. *Bioinformatics* (2019).
- [7] M. Alser et al. GateKeeper: A New Hardware Architecture for Accelerating Pre-Alignment in DNA Short Read Mapping. *Bioinformatics* (2017).
- [8] O. O. Babarinsa et al. 2015. JAFAR: Near-Data Processing for Databases. In *SIGMOD*.
- [9] A. Boroumand et al. 2018. Google Workloads for Consumer Devices: Mitigating Data Movement Bottlenecks. In *ASPLOS*.
- [10] A. Boroumand et al. 2019. CoNDA: Enabling Efficient Near-Data Accelerator Communication by Optimizing Data Movement. In *ISCA*.
- [11] A. Boroumand et al. LazyPIM: An Efficient Cache Coherence Mechanism for Processing-in-Memory. *CAL* (2016).
- [12] K. K. Chang. 2017. *Understanding and Improving the Latency of DRAM-Based Memory Systems*. Ph.D. Dissertation. Carnegie Mellon Univ.

- [13] K. K. Chang et al. 2016. Low-Cost Inter-Linked Subarrays (LISA): Enabling Fast Inter-Subarray Data Movement in DRAM. In *HPCA*.
- [14] P. Chi et al. 2016. PRIME: A Novel Processing-in-Memory Architecture for Neural Network Computation in ReRAM-Based Main Memory. In *ISCA*.
- [15] Q. Deng et al. 2018. DrAcc: a DRAM Based Accelerator for Accurate CNN Inference. In *DAC*.
- [16] J. Draper et al. 2002. The Architecture of the DIVA Processing-in-Memory Chip. In *SC*.
- [17] D. Elliott et al. Computational RAM: Implementing Processors in Memory. *IEEE Design & Test* (1999).
- [18] D. G. Elliott et al. 1992. Computational RAM: A Memory-SIMD Hybrid and Its Application to DSP. In *CICC*.
- [19] A. Farmahini-Farahani et al. 2015. NDA: Near-DRAM Acceleration Architecture Leveraging Commodity DRAM Devices and Standard Memory Modules. In *HPCA*.
- [20] M. Gao et al. 2015. Practical Near-Data Processing for In-Memory Analytics Frameworks. In *PACT*.
- [21] M. Gao et al. 2016. HRL: Efficient and Flexible Reconfigurable Logic for Near-Data Processing. In *HPCA*.
- [22] S. Ghose et al. 2018. Enabling the Adoption of Processing-in-Memory: Challenges, Mechanisms, Future Research Directions. arXiv:1802.00320 [cs:AR].
- [23] S. Ghose et al. 2019. The Processing-in-Memory Paradigm: Mechanisms to Enable Adoption. In *Beyond-CMOS Technologies for Next Generation Computer Design*.
- [24] M. Gokhale et al. Processing in Memory: The Terasys Massively Parallel PIM Array. *IEEE Computer* (1995).
- [25] B. Goodwin et al. 2017. BitFunnel: Revisiting Signatures for Search. In *SIGIR*.
- [26] B. Gu et al. 2016. Biscuit: A Framework for Near-Data Processing of Big Data Workloads. In *ISCA*.
- [27] Q. Guo et al. 2014. 3D-Stacked Memory-Side Acceleration: Accelerator and System Design. In *WoNDP*.
- [28] J.-W. Han et al. Optical Image Encryption Based on XOR Operations. *SPIE OE* (1999).
- [29] M. Hashemi et al. 2016. Accelerating Dependent Cache Misses with an Enhanced Memory Controller. In *ISCA*.
- [30] M. Hashemi et al. 2016. Continuous Runahead: Transparent Hardware Acceleration for Memory Intensive Workloads. In *MICRO*.
- [31] S. M. Hassan et al. 2015. Near Data Processing: Impact and Optimization of 3D Memory System Architecture on the Uncore. In *MEMSYS*.
- [32] K. Hsieh et al. 2016. Transparent Offloading and Mapping (TOM): Enabling Programmer-Transparent Near-Data Processing in GPU Systems. In *ISCA*.
- [33] K. Hsieh et al. 2016. Accelerating Pointer Chasing in 3D-Stacked Memory: Challenges, Mechanisms, Evaluation. In *ICCD*.
- [34] Hybrid Memory Cube Consortium. 2013. HMC Specification 1.1.
- [35] Hybrid Memory Cube Consortium. 2014. HMC Specification 2.0.
- [36] Intel Corp. 2018. 6th Generation Intel Core Processor Family Datasheet.
- [37] JEDEC. 2013. High Bandwidth Memory (HBM) DRAM. Standard No. JESD235.
- [38] S. Kanev et al. 2015. Profiling a Warehouse-Scale Computer. In *ISCA*.
- [39] H. Kang et al. 2009. One-Transistor Type DRAM. US Patent 7701751.
- [40] M. Kang et al. 2014. An Energy-Efficient VLSI Architecture for Pattern Recognition via Deep Embedding of Computation in SRAM. In *ICASSP*.
- [41] U. Kang et al. 2014. Co-Architecting Controllers and DRAM to Enhance DRAM Process Scaling. In *The Memory Forum*.
- [42] Y. Kang et al. 1999. FlexRAM: Toward an Advanced Intelligent Memory System. In *ICCD*.
- [43] S. Kaxiras et al. 1997. Distributed Vector Architecture: Beyond a Single Vector-IRAM. In *First Workshop on Mixing Logic and DRAM: Chips that Compute and Remember*.
- [44] K. Keeton et al. A Case for Intelligent Disks (IDISKS). *SIGMOD Rec.* (1998).
- [45] D. Kim et al. 2016. Neurocube: A Programmable Digital Neuromorphic Architecture with High-Density 3D Memory. In *ISCA*.
- [46] G. Kim et al. 2017. Toward Standardized Near-Data Processing with Unrestricted Data Placement for GPUs. In *SC*.
- [47] J. S. Kim et al. GRIM-Filter: Fast Seed Location Filtering in DNA Read Mapping Using Processing-in-Memory Technologies. *BMC Genomics* (2018).
- [48] Y. Kim et al. 2014. Flipping Bits in Memory Without Accessing Them: An Experimental Study of DRAM Disturbance Errors. In *ISCA*.
- [49] Y. Kim et al. 2012. A Case for Exploiting Subarray-Level Parallelism (SALP) in DRAM. In *ISCA*.
- [50] Y. Kim et al. Ramulator: A Fast and Extensible DRAM Simulator. *CAL* (2015).
- [51] D. E. Knuth. 2009. The Art of Computer Programming, Volume 4 Fascicle 1: Bitwise Tricks & Techniques; Binary Decision Diagrams.
- [52] P. M. Kogge. 1994. EXECUBE—A New Architecture for Scaleable MPPs. In *ICPP*.
- [53] S. Kvatinsky et al. MAGIC—Memristor-Aided Logic. *IEEE TCAS II: Express Briefs* (2014).
- [54] S. Kvatinsky et al. Memristor-Based Material Implication (IMPLY) Logic: Design Principles and Methodologies. *TVLSI* (2014).
- [55] B. C. Lee et al. 2009. Architecting Phase Change Memory as a Scalable DRAM Alternative. In *ISCA*.
- [56] B. C. Lee et al. Phase Change Memory Architecture and the Quest for Scalability. *CACM* (2010).
- [57] B. C. Lee et al. Phase-Change Technology and the Future of Main Memory. *IEEE Micro* (2010).
- [58] D. Lee et al. Simultaneous Multi-Layer Access: Improving 3D-Stacked Memory Bandwidth at Low Cost. *TACO* (2016).
- [59] J. H. Lee et al. 2015. BSSync: Processing Near Memory for Machine Learning Workloads with Bounded Staleness Consistency Models. In *PACT*.
- [60] S. Li et al. 2017. DRISA: A DRAM-Based Reconfigurable In-Situ Accelerator. In *MICRO*.
- [61] S. Li et al. 2016. Pinatubo: A Processing-in-Memory Architecture for Bulk Bitwise Operations in Emerging Non-Volatile Memories. In *DAC*.
- [62] Y. Li et al. 2013. BitWeaving: Fast Scans for Main Memory Data Processing. In *SIGMOD*.
- [63] J. Liu et al. 2013. An Experimental Study of Data Retention Behavior in Modern DRAM Devices: Implications for Retention Time Profiling Mechanisms. In *ISCA*.
- [64] J. Liu et al. 2012. RAIDR: Retention-Aware Intelligent DRAM Refresh. In *ISCA*.
- [65] Z. Liu et al. 2017. Concurrent Data Structures for Near-Memory Computing. In *SPAA*.
- [66] G. H. Loh et al. 2013. A Processing in Memory Taxonomy and a Case for Studying Fixed-Function PIM. In *WoNDP*.
- [67] S.-L. Lu et al. 2015. Improving DRAM Latency with Dynamic Asymmetric Subarray. In *MICRO*.
- [68] Y. Luo et al. 2017. Using ECC DRAM to Adaptively Increase Memory Capacity. arXiv:1706.08870 [cs:AR].
- [69] Y. Luo et al. 2014. Characterizing Application Memory Error Vulnerability to Optimize Datacenter Cost via Heterogeneous-Reliability Memory. In *DSN*.
- [70] K. Mai et al. 2000. Smart Memories: A Modular Reconfigurable Architecture. In *ISCA*.
- [71] O. Mutlu. Memory Scaling: A Systems Architecture Perspective. *IMW* (2013).
- [72] O. Mutlu. 2017. The RowHammer Problem and Other Issues We May Face as Memory Becomes Denser. In *DATE*.
- [73] O. Mutlu et al. Processing Data Where It Makes Sense: Enabling In-Memory Computation. *Microprocessors and Microsystems* (2019).
- [74] O. Mutlu et al. 2019. RowHammer: A Retrospective. In *IEEE TCAD*.
- [75] O. Mutlu et al. Research Problems and Opportunities in Memory Systems. *SUPERFRI* (2014).
- [76] L. Nai et al. 2017. GraphPIM: Enabling Instruction-Level PIM Offloading in Graph Computing Frameworks. In *HPCA*.
- [77] NVIDIA Corp. 2014. GeForce GTX 745 Specification.
- [78] M. Oskin et al. 1998. Active Pages: A Computation Model for Intelligent Memory. In *ISCA*.
- [79] D. Patterson et al. A Case for Intelligent RAM. *IEEE Micro* (1997).
- [80] A. Pattnaik et al. 2016. Scheduling Techniques for GPU Architectures with Processing-in-Memory Capabilities. In *PACT*.
- [81] S. H. Pugsley et al. 2014. NDC: Analyzing the Impact of 3D-Stacked Memory+Logic Devices on MapReduce Workloads. In *ISPASS*.
- [82] M. K. Qureshi et al. 2007. Adaptive Insertion Policies for High-Performance Caching. In *ISCA*.
- [83] M. K. Qureshi et al. 2009. Scalable High Performance Main Memory System Using Phase-Change Memory Technology. In *ISCA*.
- [84] M. K. Qureshi et al. 2007. Line Distillation: Increasing Cache Capacity by Filtering Unused Words in Cache Lines. In *HPCA*.
- [85] E. Riedel et al. 1998. Active Storage for Large-scale Data Mining and Multimedia Applications. In *Vldb*.
- [86] SAFARI Research Group. 2015. Ramulator: A DRAM Simulator – GitHub Repository. <https://github.com/CMU-SAFARI/ramulator/>.
- [87] V. Seshadri. 2016. *Simple DRAM and Virtual Memory Abstractions to Enable Highly Efficient Memory Systems*. Ph.D. Dissertation, Carnegie Mellon Univ.
- [88] V. Seshadri et al. Fast Bulk Bitwise AND and OR in DRAM. *CAL* (2015).
- [89] V. Seshadri et al. 2013. RowClone: Fast and Energy-Efficient In-DRAM Bulk Data Copy and Initialization. In *MICRO*.
- [90] V. Seshadri et al. 2016. Buddy-RAM: Improving the Performance and Efficiency of Bulk Bitwise Operations Using DRAM. arXiv:1611.09988 [cs:AR].
- [91] V. Seshadri et al. 2017. Ambit: In-Memory Accelerator for Bulk Bitwise Operations Using Commodity DRAM Technology. In *MICRO*.
- [92] V. Seshadri et al. 2015. Gather-Scatter DRAM: In-DRAM Address Translation to Improve the Spatial Locality of Non-Unit Strided Accesses. In *MICRO*.
- [93] V. Seshadri et al. 2017. Simple Operations in Memory to Reduce Data Movement. In *Advances in Computers, Volume 106*.
- [94] A. Shafiee et al. 2016. ISAAC: A Convolutional Neural Network Accelerator with In-Situ Analog Arithmetic in Crossbars. In *ISCA*.
- [95] D. E. Shaw et al. The NON-VON Database Machine: A Brief Overview. *IEEE Database Eng. Bull.* (1981).
- [96] H. S. Stone. A Logic-in-Memory Computer. *TC* (1970).

- [97] Z. Sura et al. 2015. Data Access Optimization in a Processing-in-Memory System. In *CF*.
- [98] P. Tuyls et al. XOR-Based Visual Cryptography Schemes. *Designs, Codes and Cryptography* (2005).
- [99] H. S. Warren. 2012. *Hacker's Delight* (2nd ed.). Addison-Wesley Professional.
- [100] H. Xin et al. Shifted Hamming Distance: A Fast and Accurate SIMD-Friendly Filter to Accelerate Alignment Verification in Read Mapping. *Bioinformatics* (2015).
- [101] H. Yoon et al. Efficient Data Mapping and Buffering Techniques for Multilevel Cell Phase-Change Memories. *ACM TACO* (2014).
- [102] D. P. Zhang et al. 2014. TOP-PIM: Throughput-Oriented Programmable Processing in Memory. In *HPDC*.
- [103] P. Zhou et al. 2009. A Durable and Energy Efficient Main Memory Using Phase Change Memory Technology. In *ISCA*.
- [104] Q. Zhu et al. 2013. Accelerating Sparse Matrix-Matrix Multiplication with 3D-Stacked Logic-in-Memory Hardware. In *HPEC*.