

Intelligent Architectures for Intelligent Machines

Onur Mutlu
ETH Zurich
omutlu@gmail.com

ABSTRACT

Computing is bottlenecked by data. Large amounts of application data overwhelm storage capability, communication capability, and computation capability of the modern machines we design today. As a result, many key applications' performance, efficiency and scalability are bottlenecked by data movement. In this keynote talk, we describe three major shortcomings of modern architectures in terms of 1) dealing with data, 2) taking advantage of the vast amounts of data, and 3) exploiting different semantic properties of application data. We argue that an intelligent architecture should be designed to handle data well. We show that handling data well requires designing architectures based on three key principles: 1) data-centric, 2) data-driven, 3) data-aware. We give several examples for how to exploit each of these principles to design a much more efficient and high performance computing system. We especially discuss recent research that aims to fundamentally reduce memory latency and energy, and practically enable computation close to data, with at least two promising novel directions: 1) performing massively-parallel bulk operations in memory by exploiting the analog operational properties of memory, with low-cost changes, 2) exploiting the logic layer in 3D-stacked memory technology in various ways to accelerate important data-intensive applications. We discuss how to enable adoption of such fundamentally more intelligent architectures, which we believe are key to efficiency, performance, and sustainability. We conclude with some guiding principles for future computing architecture and system designs.

INTRODUCTION

Existing computing systems process increasingly large amounts of data. Data is key for many modern (and likely even more future) workloads and systems. Important workloads (e.g., machine learning, artificial intelligence, genome analysis, graph analytics, databases, video analytics), whether they execute on cloud servers or mobile systems are all data intensive; they require efficient processing of large amounts of data. Today, we can generate more data than we can process, as exemplified by the rapid increase in the data obtained in astronomy observations and genome sequencing [1].

Unfortunately, the way they are designed, modern computers are not efficient at dealing with large amounts of data: large amounts of application data greatly overwhelm the storage capability, the communication capability, and the computation capability of the modern machines we design today. As such, data becomes a large performance and energy bottleneck, and it greatly impacts system robustness and security as well. As a prime example, we provide evidence that the potential for new genome sequencing technologies, such as nanopore sequencing [2], is greatly limited by how fast and how efficiently we can process the huge amounts of genomic data the underlying technology can provide us with [3, 83].

The processor-centric design paradigm (and the ensuing processor-centric execution model) of modern computing systems is one prime cause of why data overwhelms modern machines [4, 5]. With this paradigm, there is a dichotomy between processing and memory/storage: data has to be brought from storage and memory units to compute units, which are far away from the memory/storage

units. The dichotomy exists at the macro-scale (across the internet) and the micro-scale (within a single compute node). This processor-memory dichotomy leads to large amounts of data movement across the entire system, degrading performance and expending large amounts of energy. For example, a recent work [7] shows that more than 60% of the entire mobile system energy is spent on data movement across the memory hierarchy when executing four major commonly-used consumer workloads, including machine learning inference, video processing and playback, and web browsing. Similarly, due to the current design paradigm, a large fraction of the system resources is dedicated to units that store and move data, and actual computation units constitute only 5-20% of an entire chip [8] – yet, even then, data access is still a major bottleneck due to the large latency and energy costs of accessing large amounts of data.

PRINCIPLES

Our starting axiom for an intelligent architecture is that it should handle (i.e., store, access, and process) data well. But, what does it mean for an architecture to handle data well? We posit (and later demonstrate with examples) that the answer lies in satisfying three major desirable properties (or principles): 1) data-centric, 2) data-driven, and 3) data-aware.

First, the system should ensure that data does not overwhelm its components. Doing so requires effort in intelligent algorithms, intelligent architectures and intelligent whole system designs that co-design across algorithms-architectures-devices, in a manner that puts data and its processing at the center of the design, minimizing data movement and maximizing the efficiency with which data is handled, i.e., stored, accessed, and processed (e.g., as exemplified in [4-38]). We call this first principle *data-centric architectures*.

Second, an intelligent architecture takes advantage of the vast amounts of data and metadata to continuously improve its decision making, by making both its policies and mechanisms better based on online learning. In other words, the architecture should make data-driven, self-optimizing decisions in its components (e.g., as exemplified in [39-51]). We call this second principle *data-driven architectures*.

Third, an intelligent architecture understands and exploits various properties of each piece of data so that it can improve and adapt its algorithms, mechanisms, and policies based on the characteristics of data. In other words, the architecture should make *data-characteristics-aware* decisions in its components and across the entire system (as exemplified in [52-58, 107, 116, 11]). We call this third principle *data-aware architectures*.

COMPUTING ARCHITECTURES TODAY

Based on our qualitative and quantitative analyses, we find that existing computing architectures greatly fall short of handling data well. In particular they violate all of the three major desirable principles. We analyze each briefly next.

First, modern architectures are poor at dealing with data: they are designed to mainly store and move data, as opposed to actually

compute on the data. Most system resources serve the processor (and accelerators) without being capable of processing data. As such, existing architectures are *processor-centric* as opposed to *data-centric*: they place the most value in the processor (not data) and everything else in the system is viewed as secondary serving the processor. We believe this is the wrong mindset and approach in designing a balanced system that handles data well: such a system should be data-centric: i.e., data should be the prime thing that is valued and everything else in the system should be designed to 1) minimize data movement by enabling computation capability at and close to where data resides and 2) maximize the value and efficiency of processing data by enabling low-latency and low-energy access to as well as low-energy and low-cost storage of vast amounts of data.

Second, modern architectures are poor at taking advantage of vast amounts of data (and metadata) available to them during online operation and over time. They are designed to make simple decisions based on fixed policies, ignoring massive amounts of easily-available data. This is because existing architectural policies make *human-driven* decisions as opposed to *data-driven* decisions, and humans, by nature, do not seem capable of designing policies and heuristics that consider hundreds, if not thousands, of different variables that may be useful to examine to dynamically adapt online policies. It is instructive to notice that a memory controller, for example, keeps executing exactly the same fixed policy (e.g., FR-FCFS [59, 60], PAR-BS [61] or some other heuristic based policy [62-73, 117, 118]), during the *entire lifetime* of a system (for many many years!), regardless of the positive or negative impact of the resulting decisions on the system. The controller sees a vast amount of data even in the timeframe of a single millisecond (let alone years), yet it cannot learn from that data and adapt its policy because the policy is rigid and hardcoded by a human. This is clearly not intelligent: for example, as humans, we have the capability to learn from the past and adapt our actions accordingly to not repeat the same mistakes as in the past or to choose the best policy/actions that we believe will provide the highest benefits in the future. Enabling similar intelligence and far-sightedness in controller and system policies in an architecture is necessary for obtaining good performance and efficiency under a variety of system conditions and workloads.

Third, modern architectures are poor at knowing and exploiting different properties of application data. They are designed to treat all data as the same (except for a small set of specialized hints that provide some opportunity to optimize based on data characteristics in a limited manner that is very specific to the particular optimization). As such, the decisions existing architectures make are *component-aware* decisions as opposed to *data-aware* decisions: a component's (e.g., a cache's or a memory controller's) characteristics dominate the policies designed to control that component and the accessed data's characteristics are rarely conveyed to the policy or even known. If the characteristics of the data to be accessed or manipulated were known, the decisions taken could be very different: for example, if we knew the relative compressibility of different types of data [74-81], different components in the entire system could be designed in a manner that adaptively scales their capability to match the compressibility of different data elements, in order to maximize both performance and efficiency. Modifying the architecture and its interface to become richer and more expressive, and to include rich and accurate information on various properties of data that is to be processed, is therefore critical to customizing the architecture to the characteristics of the data and, thus, enabling intelligent adaptation of system policies to data characteristics.

INTELLIGENT COMPUTING ARCHITECTURES

A major chunk of our talk describes in detail the characteristics of an intelligent computing architecture, by concrete examples and their

empirical evaluation. This paper does not go into detail, but provides a brief overview with references to other works that exemplify such architectures. A detailed version of this talk can be found in [82].

Data-Centric

A data-centric architecture has at least four major characteristics. First, it enables processing capability in or near where data resides (i.e., in or near memory structures), as described in detail in [4-6, 8, 38] and exemplified by [7-12, 14, 19, 20, 24, 27, 30, 34, 84, 108-113]. Second, it provides low-latency and low-energy access to data, as exemplified by [11-13, 15-18, 21, 23, 31-33, 84-86]. Third, it enables low-cost data storage and processing (i.e., high capacity memory at low cost, via techniques like new memory technologies, hybrid memory systems and/or compressed memory systems), as exemplified by [22, 87-96, 74, 76, 78, 107, 116]. Fourth, it provides mechanisms for intelligent data management (with intelligent controllers handling robustness, security, cost, etc.), as described in detail in [97-103, 116] and exemplified by, e.g., [104-106, 116].

Data-Driven

A data-driven architecture enables the machine itself to learn the best policies for managing itself and executing programs. Controllers in such an architecture, when needed, are data-driven autonomous agents that automatically learn far-sighted policies. A prime example of such a controller is the *reinforcement learning based self-optimizing memory controllers* [39]. Such controllers can not only improve performance and efficiency under a wide variety of conditions and workloads but also reduce the hardware and system designer's burden in designing sophisticated controllers [39].

Data-Aware

A data-aware architecture understands what it can do with and to each piece of data, and uses this information about data characteristics to maximize system efficiency and performance. In other words, it customizes itself (i.e., its policies and mechanisms) to the characteristics of the data it is dealing with. Such an architecture requires knowledge of various characteristics of different data elements and structures. Many semantic or other characteristics of data (e.g., compressibility, approximability, sparsity, criticality, access and security semantics, locality, latency vs. bandwidth sensitivity) are invisible or unknown to modern hardware and thus need to be communicated or discovered. We believe efficient and expressive software/hardware interfaces and mechanisms, as exemplified by X-Mem [52, 53] and the Virtual Block Interface [56] as well as other works [54, 55, 57, 58, 107, 116, 11], are promising approaches to creating general-purpose data-aware architectures.

REFERENCES

- [1] Z. D. Stevens et al., "Big data: astronomical or genetical?", *PLoS Biology*, 2015.
- [2] D. Senol Cali et al., "Nanopore Sequencing Technology and Tools for Genome Assembly: Computational Analysis of the Current State, Bottlenecks and Future Directions" BIB 2019.
- [3] O. Mutlu, "Accelerating Genome Analysis: A Primer on an Ongoing Journey", Keynote Talk at HiCOMB-17, 2018.
- [4] S. Ghose et al., "Processing-in-Memory: A Workload-Driven Perspective", IBM JRD 2019.
- [5] O. Mutlu et al., "Processing Data Where It Makes Sense: Enabling In-Memory Computation", MICPRO, 2019.
- [6] V. Seshadri and O. Mutlu, "In-DRAM Bulk Bitwise Execution Engine", *Advances in Computers*, 2020.
- [7] A. Boroumand et al., "Google Workloads for Consumer Devices: Mitigating Data Movement Bottlenecks", *ASPLOS* 2018.

- [8] O. Mutlu, "Enabling Computation with Minimal Data Movement: Changing the Computing Paradigm for High Efficiency", Design Automation Summer School Lecture, DAC 2019. <https://people.inf.ethz.ch/omutlu/pub/onur-DAC-DASS-EnablingInMemoryComputation-June-2-2019.pptx>
- [9] J. Ahn et al., "A Scalable Processing-in-Memory Accelerator for Parallel Graph Processing", ISCA 2015.
- [10] V. Seshadri et al., "Ambit: In-Memory Accelerator for Bulk Bitwise Operations Using Commodity DRAM Technology", MICRO 2017.
- [11] H. Luo et al., "CLR-DRAM: A Low-Cost DRAM Architecture Enabling Dynamic Capacity-Latency Trade-Off", ISCA 2020.
- [12] K. Chang et al., "Low-Cost Inter-Linked Subarrays (LISA): Enabling Fast Inter-Subarray Data Movement in DRAM", HPCA 2016.
- [13] D. Lee et al., "Adaptive-Latency DRAM: Optimizing DRAM Timing for the Common-Case", HPCA 2015.
- [14] K. Hsieh et al., "Accelerating Pointer Chasing in 3D-Stacked Memory: Challenges, Mechanisms, Evaluation", ICCD 2016.
- [15] K. Chang et al., "Understanding Latency Variation in Modern DRAM Chips: Experimental Characterization, Analysis, and Optimization", SIGMETRICS 2016.
- [16] K. Chang et al., "Understanding Reduced-Voltage Operation in Modern DRAM Devices: Experimental Characterization, Analysis, and Mechanisms", SIGMETRICS 2017.
- [17] D. Lee et al., "Design-Induced Latency Variation in Modern DRAM Chips: Characterization, Analysis, and Latency Reduction Mechanisms", SIGMETRICS 2017.
- [18] S. Ghose et al., "What Your DRAM Power Models Are Not Telling You: Lessons from a Detailed Experimental Study", SIGMETRICS 2018.
- [19] K. Hsieh et al., "Transparent Offloading and Mapping (TOM): Enabling Programmer-Transparent Near-Data Processing in GPU Systems", ISCA 2016.
- [20] J. Ahn et al., "PIM-Enabled Instructions: A Low-Overhead, Locality-Aware Processing-in-Memory Architecture", ISCA 2015.
- [21] J. Liu et al., "RAIDR: Retention-Aware Intelligent DRAM Refresh", ISCA 2012.
- [22] B. C. Lee et al., "Architecting Phase Change Memory as a Scalable DRAM Alternative", ISCA 2009.
- [23] D. Lee et al., "Decoupled Direct Memory Access: Isolating CPU and IO Traffic by Leveraging a Dual-Data-Port DRAM", PACT 2015.
- [24] V. Seshadri et al., "Gather-Scatter DRAM: In-DRAM Address Translation to Improve the Spatial Locality of Non-unit Strided Accesses", MICRO 2015.
- [25] D. Lee et al., "Simultaneous Multi-Layer Access: Improving 3D-Stacked Memory Bandwidth at Low Cost", TACO 2016.
- [26] H. Hassan et al., "ChargeCache: Reducing DRAM Latency by Exploiting Row Access Locality", HPCA 2016.
- [27] M. Hashemi et al., "Accelerating Dependent Cache Misses with an Enhanced Memory Controller", ISCA 2016.
- [28] M. Patel et al., "The Reach Profiler (REAPER): Enabling the Mitigation of DRAM Retention Failures via Profiling at Aggressive Conditions", ISCA 2017.
- [29] S. Khan et al., "Detecting and Mitigating Data-Dependent DRAM Failures by Exploiting Current Memory Content", MICRO 2017.
- [30] J. S. Kim et al., "GRIM-Filter: Fast Seed Location Filtering in DNA Read Mapping Using Processing-in-Memory Technologies", BMC Genomics 2018.
- [31] A. Das et al., "VRL-DRAM: Improving DRAM Performance via Variable Refresh Latency", DAC 2018.
- [32] J. S. Kim et al., "Solar-DRAM: Reducing DRAM Access Latency by Exploiting the Variation in Local Bitlines", ICCD 2018.
- [33] Y. Wang et al., "Reducing DRAM Latency via Charge-Level-Aware Look-Ahead Partial Restoration", MICRO 2018.
- [34] J. S. Kim et al., "D-RaNGe: Using Commodity DRAM Devices to Generate True Random Numbers with Low Latency and High Throughput", HPCA 2019.
- [35] H. Hassan et al., "CROW: A Low-Cost Substrate for Improving DRAM Performance, Energy Efficiency, and Reliability", ISCA 2019.
- [36] S. Song et al., "Improving Phase Change Memory Performance with Data Content Aware Access", ISMM 2020.
- [37] G. Singh et al., "NAPEL: Near-Memory Computing Application Performance Prediction via Ensemble Learning", DAC 2019.
- [38] O. Mutlu et al., "Enabling Practical Processing in and near Memory for Data-Intensive Computing", DAC 2019.
- [39] E. Ipek et al., "Self Optimizing Memory Controllers: A Reinforcement Learning Approach", ISCA 2008.
- [40] D. A. Jimenez and C. Lin, "Dynamic Branch Prediction with Perceptrons", HPCA 2001.
- [41] D. A. Jimenez, "Fast Path-Based Neural Branch Prediction", MICRO 2003.
- [42] D. A. Jimenez, "Piecewise Linear Branch Prediction", ISCA 2005.
- [43] D. A. Jimenez, "An optimized scaled neural branch predictor", ICCD 2011.
- [44] E. Teran et al., "Perceptron learning for reuse prediction", MICRO 2016.
- [45] E. Garza et al., "Bit-level perceptron prediction for indirect branches", ISCA 2019.
- [46] E. Bhatia et al., "Perceptron-based prefetch filtering", ISCA 2019.
- [47] L. Peled et al., "A Neural Network Prefetcher for Arbitrary Memory Access Patterns", TACO 2020.
- [48] L. Peled et al., "Semantic locality and context-based prefetching using reinforcement learning", ISCA 2015.
- [49] R. Bitirgen et al., "Coordinated management of multiple interacting resources in chip multiprocessors: A machine learning approach", MICRO 2008.
- [50] J. Mukundan and J. F. Martinez, "MORSE: Multi-objective reconfigurable self-optimizing memory scheduler", HPCA 2012.
- [51] J. F. Martinez and E. Ipek, "Dynamic Multicore Resource Management: A Machine Learning Approach", IEEE Micro 2009.
- [52] N. Vijaykumar et al., "A Case for Richer Cross-layer Abstractions: Bridging the Semantic Gap with Expressive Memory", ISCA 2018.
- [53] N. Vijaykumar et al., "The Locality Descriptor: A Holistic Cross-Layer Abstraction to Express Data Locality in GPUs", ISCA 2018.
- [54] S. Koppula et al., "EDEN: Enabling Energy-Efficient, High-Performance Deep Neural Network Inference Using Approximate DRAM", MICRO 2019.
- [55] K. Kanellopoulos et al., "SMASH: Co-designing Software Compression and Hardware-Accelerated Indexing for Efficient Sparse Matrix Operations", MICRO 2019.
- [56] N. Hajinazar et al., "The Virtual Block Interface: A Flexible Alternative to the Conventional Virtual Memory Framework", ISCA 2020.
- [57] Z. Yu et al., "Labeled RISC-V: A New Perspective on Software-Defined Architecture", CARRV 2017.
- [58] J. Ma et al., "Supporting Differentiated Services in Computers via Programmable Architecture for Resourcing-on-Demand (PARD)", ASPLOS 2015.
- [59] S. Rixner et al., "Memory access scheduling", ISCA 2000.
- [60] W. K. Zuravleff and T. Robinson, "Controller for a synchronous DRAM that maximizes throughput by allowing memory requests and commands to be issued out of order", U.S. Patent Number 5,630,096, May 1997.
- [61] O. Mutlu and T. Moscibroda, "Parallelism-Aware Batch Scheduling: Enhancing both Performance and Fairness of Shared DRAM Systems", ISCA 2008.

- [62] H. Usui et al., “DASH: Deadline-Aware High-Performance Memory Scheduler for Heterogeneous Systems with Hardware Accelerators”, TACO 2016.
- [63] O. Mutlu and T. Moscibroda, “Stall-Time Fair Memory Access Scheduling for Chip Multiprocessors”, MICRO 2007.
- [64] Y. Kim et al., “ATLAS: A Scalable and High-Performance Scheduling Algorithm for Multiple Memory Controllers”, HPCA 2010.
- [65] Y. Kim et al., “Thread Cluster Memory Scheduling: Exploiting Differences in Memory Access Behavior”, MICRO 2010.
- [66] I. Hur and C. Lin, “Adaptive History-Based Memory Schedulers”, MICRO 2004.
- [67] I. Hur and C. Lin, “Memory scheduling for modern microprocessors”, ACM TOCS 2007.
- [68] C. Natarajan et al., “A study of performance impact of memory controller features in multi-processor server environment”, WMPI 2004.
- [69] S. Rixner, “Memory controller optimizations for web servers”, MICRO 2004.
- [70] L. Subramanian et al., “BLISS: Balancing Performance, Fairness and Complexity in Memory Access Scheduling”, IEEE TPDS 2016.
- [71] L. Subramanian et al., “The Blacklisting Memory Scheduler: Achieving High Performance and Fairness at Low Cost”, ICCD 2014.
- [72] R. Ausavarungnirun et al., “Staged Memory Scheduling: Achieving High Performance and Scalability in Heterogeneous Systems”, ISCA 2012.
- [73] K. J. Nesbit et al., “Fair Queuing Memory Systems”, MICRO 2006.
- [74] G. Pekhimenko et al., “Base-Delta-Immediate Compression: Practical Data Compression for On-Chip Caches”, PACT 2012.
- [75] N. Vijaykumar et al., “A Case for Core-Assisted Bottleneck Acceleration in GPUs: Enabling Flexible Data Compression with Assist Warps”, ISCA 2015.
- [76] G. Pekhimenko et al., “Linearly Compressed Pages: A Low-Complexity, Low-Latency Main Memory Compression Framework”, MICRO 2013.
- [77] G. Pekhimenko et al., “A Case for Toggle-Aware Compression for GPU Systems”, HPCA 2016.
- [78] M. Ekman and P. Stenstrom, “A Robust Main-Memory Compression Scheme”, ISCA 2005.
- [79] A. Arelakis et al., “HyComp: a hybrid cache compression method for selection of data-type-specific compression methods”, MICRO 2015.
- [80] A. Arelakis and P. Stenstrom, “SC²: A statistical compression cache scheme”, ISCA 2014.
- [81] G. Pekhimenko et al., “Exploiting Compressed Block Size as an Indicator of Future Reuse”, HPCA 2015.
- [82] O. Mutlu, “Intelligent Architectures for Intelligent Machines”, Keynote Talk at 17th ChinaSys Workshop, December 2019. https://www.youtube.com/watch?v=n8Aj_A0WSg8
- [83] M. Alser et al., “Shouji: A Fast and Efficient Pre-Alignment Filter for Sequence Alignment”, Bioinformatics 2019.
- [84] V. Seshadri et al., “RowClone: Fast and Energy-Efficient IN-DRAM Bulk Data Copy and Initialization”, MICRO 2013.
- [85] D. Lee et al., “Tiered-Latency DRAM: A Low Latency and Low Cost DRAM Architecture”, HPCA 2013.
- [86] Y. Kim et al., “A Case for Exploiting Subarray-Level Parallelism (SALP) in DRAM”, ISCA 2012.
- [87] B. C Lee et al., “Phase Change Technology and the Future of Main Memory”, IEEE Micro 2010.
- [88] Y. Li et al., “Utility-Based Hybrid Memory Management”, CLUSTER 2017.
- [89] H. Yoon et al., “Row Buffer Locality Aware Caching Policies for Hybrid Memories”, ICCD 2012.
- [90] C. Wang et al., “Panthera: Holistic Memory Management for Big Data Processing over Hybrid Memories”, PLDI 2019.
- [91] J. Meza et al., “Enabling Efficient and Scalable Hybrid Memories Using Fine-Granularity DRAM Cache Management”, IEEE CAL 2012.
- [92] M. K. Qureshi et al., “Scalable high performance main memory system using phase-change memory technology”, ISCA 2009.
- [93] M. K. Qureshi et al., “Morphable memory system: a robust architecture for exploiting multi-level phase change memories”, ISCA 2010.
- [94] C-C. Chou et al., “CAMEO: A Two-Level Memory Organization with Capacity of Main Memory and Flexibility of Hardware-Managed Cache”, MICRO 2014.
- [95] V. Young et al., “Enabling Transparent Memory-Compression for Commodity Memory Systems”, HPCA 2014.
- [96] X. Yu et al., “Banshee: Bandwidth-Efficient DRAM Caching via Software/Hardware Cooperation”, MICRO 2017.
- [97] O. Mutlu, “Memory Scaling: A Systems Architecture Perspective”, IMW 2013.
- [98] O. Mutlu and L. Subramanian, “Research Problems and Opportunities in Memory Systems”, SUPERFRI 2015.
- [99] O. Mutlu and J. Kim, “RowHammer: A Retrospective”, IEEE TCAD 2019.
- [100] Y. Cai et al., “Error Characterization, Mitigation, and Recovery in Flash Memory Based Solid State Drives”, Proc. IEEE 2017.
- [101] Y. Cai et al., “Errors in Flash-Memory-Based Solid-State Drives: Analysis, Mitigation, and Recovery”, Inside Solid-State Drives, 2018.
- [102] O. Mutlu, “The RowHammer Problem and Other Issues We May Face as Memory Becomes Denser”, DATE 2017.
- [103] O. Mutlu et al., “Recent Advances in DRAM and Flash Memory Architectures”, IPSI TIR, July 2018.
- [104] Y. Kim et al., “Flipping Bits in Memory Without Accessing Them: An Experimental Study of DRAM Disturbance Errors”, ISCA 2014.
- [105] J. S. Kim et al., “Revisiting RowHammer: An Experimental Analysis of Modern Devices and Mitigation Techniques”, ISCA 2020.
- [106] P. Frigo et al., “TRRespass: Exploiting the Many Sides of Target Row Refresh”, S&P 2020.
- [107] Y. Luo et al., “Characterizing Application Memory Error Vulnerability to Optimize Data Center Cost via Heterogeneous-Reliability Memory”, DSN 2014.
- [108] A. Boroumand et al., “CoNDA: Efficient Cache Coherence Support for Near-Data Accelerators”, ISCA 2019.
- [109] V. Seshadri et al., “Fast Bulk Bitwise AND and OR in DRAM”, IEEE CAL 2015.
- [110] A. Boroumand et al., “LazyPIM: An Efficient Cache Coherence Mechanism for Processing-in-Memory”, IEEE CAL 2016.
- [111] M. Hashemi et al., “Continuous Runahead: Transparent Hardware Acceleration for Memory Intensive Workloads”, MICRO 2016.
- [112] A Pattnaik et al., “Scheduling Techniques for GPU Architectures with Processing-In-Memory Capabilities”, PACT 2016.
- [113] A. Boroumand et al., “LazyPIM: An Efficient Cache Coherence Mechanism for Processing-in-Memory”, IEEE CAL 2016.
- [114] H. Hassan et al., “SoftMC: A Flexible and Practical Open-Source Infrastructure for Enabling Experimental DRAM Studies”, HPCA 2017.
- [115] Y. Kim et al., “Ramulator: A Fast and Extensible DRAM Simulator”, IEEE CAL 2015.
- [116] J. Meza et al., “A Case for Efficient Hardware-Software Cooperative Management of Storage and Memory”, WEED 2013.
- [117] L. Subramanian et al., “MISE: Providing Performance Predictability and Improving Fairness in Shared Main Memory Systems”, HPCA 2013.
- [118] C. J. Lee et al., “Prefetch-Aware DRAM Controllers”, MICRO 2008.