

Energy Optimizations For Memory Architectures

Mukesh Sekar
Electrical and Computer
Engineering
Colorado State University
Fort Collins, Colorado
msekar@colostate.edu

Abstract— This paper presents a comprehensive survey of recent advancements and challenges in energy-efficient computing, spanning hardware and software domains. It explores the potential of non-volatile memory (NVM) technologies like Racetrack Memory and Resistive Random Access Memory (RRAM) for in-memory computing (IMC) to address computation and memory constraints in edge devices. Furthermore, it discusses optimization strategies in Python-based scientific applications, showcasing efforts to balance abstraction and efficiency, such as integrating SharedArray into the PyCOMPSs programming model for performance gains in distributed environments. The paper also delves into novel job scheduling schemes like Deviation Backfilling in High Performance Computing (HPC) systems, enhancing efficiency while ensuring fairness, and employing job runtime prediction techniques using kNN. Additionally, it investigates approaches to improve energy efficiency and predictability in real-time embedded systems with non-volatile memories, leveraging hybrid hardware-software solutions and optimization techniques to mitigate latency and enhance worst-case execution time (WCET). Lastly, it surveys memory integrity techniques developed over decades to combat memory-related cyber-attacks, assessing their effectiveness and evolution.

Keywords—Energy-efficient computing, Non-volatile memory (NVM), Racetrack Memory, Resistive Random Access Memory (RRAM), In-memory computing (IMC), Python-based scientific applications, SharedArray, PyCOMPSs programming model, Job scheduling, High Performance Computing (HPC), Deviation Backfilling, Real-time embedded systems, Worst-case execution time (WCET), Memory integrity techniques, Cyber-attacks

I. INTRODUCTION

The surge in computational demands across various domains has underscored the critical importance of energy-efficient computing. As traditional architectures primarily focused on performance struggle to cope with escalating energy consumption, a paradigm shift towards energy efficiency becomes imperative. This survey delves into the evolving landscape of energy-efficient computing, examining the challenges posed by escalating energy demands and the strategies devised to mitigate them.

The relentless growth of data and computational complexity has led to an incessant rise in energy requirements within computing systems. Conventional architectures optimized solely for performance often overlook energy efficiency considerations, resulting in substantial power consumption. This issue is particularly pronounced in real-time embedded systems

with constrained battery life and high-performance computing facilities burdened by exorbitant energy costs and environmental implications.

In response to escalating energy demands[1], researchers have devised diverse strategies aimed at curbing power consumption. Energy-efficient processors and memory technologies emerge as pivotal components in reducing overall power requirements. Techniques such as dynamic voltage and frequency scaling (DVFS) dynamically adjust processor speed based on workload, thereby minimizing energy expenditure during less demanding tasks.

The proliferation of deep learning has accentuated concerns regarding model size, particularly for resource-constrained edge devices. Model compression techniques such as pruning and quantization offer viable avenues for significantly reducing the size of deep neural network (DNN) models while preserving acceptable accuracy levels. This translates to diminished memory footprint and lower power consumption during inference on edge devices[2],[3].

Ensuring predictable performance is paramount in real-time embedded systems, necessitating meticulous attention to worst-case execution time (WCET) analysis. Non-volatile memory (NVM) technologies emerge as advantageous for real-time systems owing to their lower leakage power consumption compared to traditional volatile memory. Ongoing research focuses on enhancing memory architecture to optimize WCET and performance while minimizing energy overhead.

Shared memory programming models like SharedArray facilitate efficient parallel computing by enabling rapid data sharing between processes. When integrated with frameworks such as PyCOMPSs[6], these models hold the potential to enhance energy efficiency and performance in high-performance computing (HPC) and scientific computing environments.

Efficient job scheduling algorithms are indispensable for maximizing resource utilization and minimizing energy consumption in HPC systems. Techniques like Deviation Backfilling leverage user prediction deviation for job delays, culminating in more judicious scheduling decisions and reduced resource wastage.

Memory-related vulnerabilities in software pose significant security risks, particularly in systems employing memory-unsafe languages like C and C++. Memory safety techniques, including hardware memory protection and virtualization, are

pivotal in mitigating such vulnerabilities. Nonetheless, these techniques may entail an energy overhead, prompting ongoing research endeavors to achieve memory safety with minimal impact on energy efficiency.

Cache hierarchies play a pivotal role in enhancing processor performance, with optimized cache parameters capable of significantly reducing energy consumption without compromising performance. By scrutinizing cache behavior and tailoring it to specific workloads, researchers aim to strike a delicate balance between performance and energy efficiency.

The pursuit of energy-efficient computing represents a multifaceted endeavor encompassing diverse strategies and technologies. From in-memory computing to model compression and cache optimization, a concerted effort is underway to mitigate the burgeoning energy demands of computational systems. As research endeavors continue to evolve, the prospect of achieving energy efficiency gains while sustaining computational performance remains promising, heralding a more sustainable and environmentally conscious future for computing.

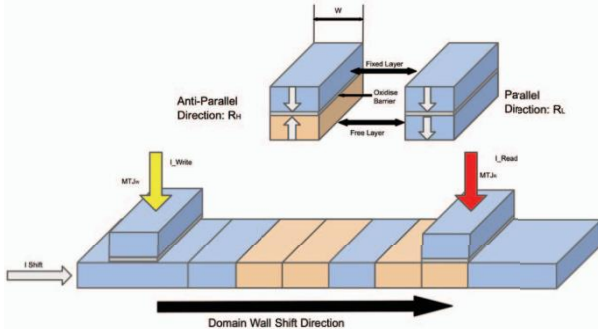


Fig 1.0 Basic structure of the nanowire and vertical magnetic tunnel junction

In this survey, our aim is to comprehensively explore energy optimizations for memory architectures within the realm of energy-efficient computing[4],[6]. We delve into various techniques and technologies aimed at reducing energy consumption while maintaining or improving computational performance. By examining the challenges posed by escalating energy demands and the strategies devised to mitigate them, we seek to identify key avenues for optimizing memory architectures to achieve energy efficiency.

Furthermore, we endeavor to discuss the role of memory optimization techniques, including cache optimization and memory safety measures, in striking a balance between performance and energy efficiency. By analyzing current research trends and emerging technologies in this field, we aim to provide insights into promising avenues for achieving energy optimization in memory architectures.

Ultimately, through this survey, we aspire to contribute to the ongoing discourse on energy-efficient computing by shedding light on the significance of memory optimizations and identifying future directions for research and development in this critical area.

II. RELATED WORKS

This section critically evaluates key papers that substantially contribute to energy optimizations for memory architectures. It specifically focuses on exploring advancements in minimizing energy consumption while enhancing memory performance and efficiency. Each paper undergoes thorough scrutiny regarding its context, the addressed problem[7], proposed optimization techniques, and resultant outcomes. Through this detailed analysis, the objective is to encapsulate the cutting-edge developments in energy-efficient memory architectures and pave the way for identifying potential research gaps and avenues for future exploration in this field[9].

In the realm of energy optimization for memory architectures, a diverse array of approaches has emerged, each tackling unique challenges and offering innovative solutions. Shared memory models like OpenMP and OpenSHMEM have long been pillars of parallel programming, providing scalable solutions across various platforms. However, integrating these models with Python poses challenges due to the language's Global Interpreter Lock (GIL) and limitations in accessing Python objects directly. Python's multiprocessing module attempts to address these issues with its SharedMemory class, but it falls short in providing the fine-grained control necessary for optimal memory management and sharing mechanisms.

Meanwhile, [8][10]the Plasma In-Memory Object Store presents an alternative solution by facilitating efficient storage of immutable objects in shared memory. While promising, Plasma faces its own set of challenges, particularly regarding user-defined objects and integration complexities with frameworks like PyCOMPSs. This underscores the need for versatile and robust memory management solutions, especially in dynamically changing environments where different programming languages and frameworks coexist.

Within the realm of memory subsystem optimization, a plethora of strategies are at play. From heuristic algorithms to Design of Experiments (DoE)[20][21][25] techniques, researchers aim to optimize cache hierarchies and leverage memory low-power modes to enhance energy efficiency. Some approaches focus on dynamic scratchpad memory allocation based on application behavior, while others delve into ILP-based methods for comprehensive optimization, balancing energy savings with performance gains.

Hybrid memory architectures, combining SRAM and NVM, represent another frontier in energy optimization efforts. These architectures seek to strike a delicate balance between performance and energy consumption[12], often through task and data mapping strategies tailored to real-time embedded systems. Moreover, investigations into multi-bank memory architectures aim to mitigate leakage power, with considerations tailored to different memory types and system constraints.

In the realm of real-time systems incorporating non-volatile memories, efforts are directed toward ensuring predictability and performance. Techniques such as data mapping[16], task scheduling, and buffering strategies are honed to minimize energy consumption while meeting stringent timing requirements. Pioneering approaches focusing on fully non-volatile memory architectures offer promising solutions, presenting predictable and high-performance memory solutions crafted specifically for real-time embedded systems.

Overall, the pursuit of energy optimization in memory architectures is a multifaceted endeavor[12], requiring a blend of theoretical insights, practical innovations, and interdisciplinary collaboration. As technology evolves and new challenges emerge, researchers continue to push the boundaries of what's possible, striving to develop solutions that not only optimize energy usage but also enhance overall system performance and reliability..

III. DESIGN AND IMPLEMENTATION

The study focuses on improving memory performance and energy efficiency in multiprocessor systems[17][19] by introducing a software-based distributed memory management system. Departing from conventional structural designs, the authors advocate for a software-centric approach, leveraging software parameters to manage the distributed shared address space more effectively. This approach offers heightened flexibility and adaptability, empowering real-time adjustments and optimizations tailored to the specific requirements of various applications and system configurations[16].

Central to the proposed system is the integration of a sequential consistency mechanism, aimed at governing access to distributed shared memory. This mechanism ensures that read operations consistently yield values in line with the programmer's expectations, thereby fostering a coherent view of shared memory across disparate processes. By enforcing sequential consistency, the system bolsters reliability and predictability, essential for upholding program integrity and mitigating potential errors or discrepancies during parallel execution.

The effectiveness of the approach in enhancing memory performance is validated through comprehensive simulation and evaluation. Detailed analyses spanning cache compression, power consumption, and device utilization provide valuable insights into the realized performance improvements and energy efficiency gains. Additionally, the integration of specialized components such as scratchpad banks, address decoders, and dedicated memory controllers underscores the holistic nature of the optimization endeavors, aimed at maximizing memory utilization, minimizing die area footprint, and optimizing overall system performance.

Furthermore, the study emphasizes the practical implications of the proposed approach, highlighting its potential to address critical challenges in modern multiprocessor systems. By

prioritizing adaptability and real-time optimization, the software-based distributed memory management system offers a versatile solution capable of meeting the evolving demands of diverse applications and system configurations. This versatility extends to potential applications in domains such as the Arrowhead Framework, where standardized engineering data serves as the foundation for various services and applications. As such, the study not only contributes to advancing memory performance and energy efficiency but also underscores the broader impact of innovative memory management solutions in shaping the future of computing architectures.

A. Enhancing Memory Performance and Energy Efficiency in

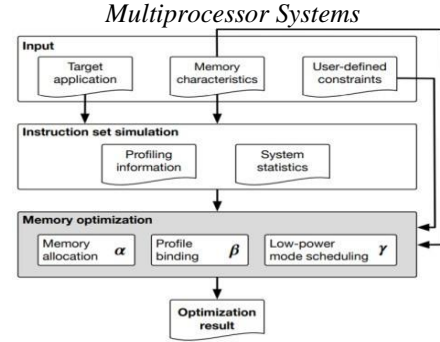


Figure 1.1 Time Memory Optimization work flow

B. Co-Optimization of Device and Architecture for RRAM-Based In-Memory Computing Systems

The research conducted by the authors' group at Peking University focuses on co-optimizing device engineering and architecture design for Resistive Random Access Memory (RRAM)-based In-Memory Computing (IMC) systems[23][24][25]. The methodology employed involves a multi-faceted approach to address various challenges encountered in RRAM-based IMC systems. Firstly, the authors delve into device optimization, where they meticulously investigate and model non-ideal behaviors of RRAM devices, such as random telegraph noise (RTN), early-stage fluctuation (ERF), and retention degradation (RD). Through modeling and analysis, they assess the impacts of these non-idealities on the accuracy of IMC systems[29][27]. Additionally, thermal effects in RRAM arrays, including self-heating and thermal crosstalk, are studied, and mitigation techniques such as 3D vertical RRAM structures and resistance temperature compensation are proposed.

In the realm of architecture optimization, the authors introduce innovative solutions tailored to enhance the efficiency and performance of RRAM-based IMC systems. They propose the Lattice architecture, an ADC/DAC-less RRAM[30] IMC architecture, which eliminates costly analog-digital converters and offers effective data mapping, leading to heightened area and energy efficiency. Furthermore, the development of the MobiLattice block, a hybrid digital/analog IMC block, enables efficient support for both analog and digital mode operations on RRAM crossbars. This advancement significantly accelerates

the processing of depth-wise convolutional neural networks, contributing to improved performance in AI and neural network applications. Additionally, the authors introduce PIMulator-NN, an end-to-end simulation framework for IMC-based accelerators, facilitating flexible architecture descriptions and accurate architectural modeling, thus aiding in the optimization process.

The research methodology adopted in this study takes a comprehensive approach, covering device engineering and architecture design to co-optimize [24]RRAM-based In-Memory Computing (IMC) systems. This approach aims to tackle a range of challenges, including non-idealities, thermal effects, hardware overhead, and architectural flexibility. Through the proposed techniques and architectures, the goal is to improve the accuracy, reliability, energy efficiency, and overall performance of RRAM-based IMC systems, with a particular focus on applications in AI and neural networks.

C. Optimizing Two-Level Cache Memory Hierarchies for Energy Efficiency Using TECH-DoE

Optimizing energy efficiency in embedded applications, this paper introduces TECH-DoE, a novel mechanism focused on tuning two-level cache memory hierarchies[33][36]. TECH-DoE employs Design of Experiments (DoE) to systematically analyze various cache parameters and their interactions, aiming to minimize energy consumption while maintaining performance. Through 2^k factorial design in DoE, a comprehensive set of experiments is planned to evaluate the effects of cache parameters, including size, line size, and associativity, on energy efficiency. Subsequently, a tuning sequence is derived based on the degree of influence of each parameter obtained from the DoE analysis. The TECH-DoE heuristic iteratively adjusts cache parameters in accordance with this sequence, retaining configurations that reduce energy consumption without significantly compromising performance. Implementation involves the utilization of tools like SimpleScalar for energy consumption data and the eCACTI[18] tool for performance statistics, facilitating efficient exploration and evaluation of cache configurations.

Implemented with a focus on simplicity and computational efficiency, TECH-DoE[32] offers significant advantages over existing cache tuning heuristics. By employing a systematic approach and leveraging DoE techniques, TECH-DoE achieves an average energy reduction of approximately 6% for data caches compared to prior methods. Moreover, it demonstrates the capability to identify configurations close to the optimal Pareto front while being computationally less demanding than alternatives such as genetic algorithms. Through evaluation on representative applications from the MiBench benchmark suite, TECH-DoE[31] showcases its effectiveness in enhancing energy efficiency in embedded systems, highlighting its potential as a valuable tool for cache optimization in real-world scenarios.

D. Predictable and High-Performance Fully Non-Volatile Data Memory Architecture for Real-Time Embedded Systems

The methodology presented in this paper revolves around the design and implementation of a hybrid hardware-software solution aimed at enhancing the predictability and performance of fully non-volatile data memory architecture for real-time embedded systems. The hardware component of the solution entails the partitioning of the non-volatile memory [38][39](NVM) into multiple banks, supplemented by the incorporation of a write buffer situated between the processor and memory. Additional hardware components such as registers, multiplexers, and decoders are integrated to efficiently manage parallel write operations. This hardware configuration enables the masking of long write latencies inherent in NVMs like STT-RAM, allowing the processor to continue execution without waiting for write operations to complete.

Complementing the hardware aspect, the software component introduces a novel "flush" instruction designed to ensure the completion of all pending write operations in the buffer. This contributes to the predictability of worst-case execution time (WCET) analysis, a critical consideration in real-time embedded systems. Furthermore, a WCET analysis model and algorithm are developed to estimate the WCET considering the proposed memory architecture with the write buffer. This analysis model can be seamlessly integrated into existing WCET estimation tools, facilitating accurate performance prediction and system optimization.

To evaluate the efficacy of the proposed design, extensive experimental evaluation is conducted[16][18], comparing the performance of the hybrid memory architecture against a baseline configuration across various parameters such as the number of banks and write buffer size. The evaluation metrics include WCET, system performance, and energy consumption. The results demonstrate significant improvements in WCET and performance compared to the baseline configuration, validating the effectiveness of the hybrid hardware-software solution in enhancing the predictability and performance of fully non-volatile data memory architecture for real-time embedded systems.

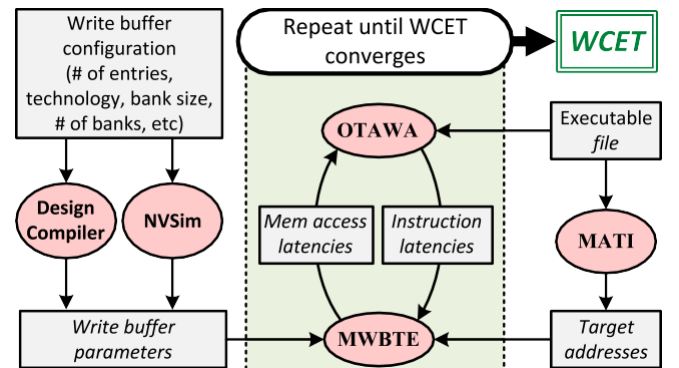


Figure 1.2 The WCET estimation flow includes estimating memory subsystem parameters, analyzing application instructions, and iterating to achieve stable WCET.

E. Enhancing Job Scheduling Efficiency with Deviation Backfilling in HPC Systems

The methodology of the Deviation Backfilling scheme revolves around optimizing job scheduling efficiency in High-Performance Computing (HPC)[36],[39] systems through innovative approaches to delay threshold determination and runtime prediction. Initially, the scheme employs a backfilling approach, allowing shorter jobs to precede longer ones without significantly delaying the latter, thus enhancing overall system throughput. What sets Deviation Backfilling apart is its unique determination of the delay threshold for the first job in the queue, which relies on the deviation between user-provided runtime estimates and system-predicted runtimes. This deviation is calculated using a k-Nearest Neighbors (kNN) algorithm, providing a dynamic and adaptive threshold mechanism that considers both user inputs and system predictions.

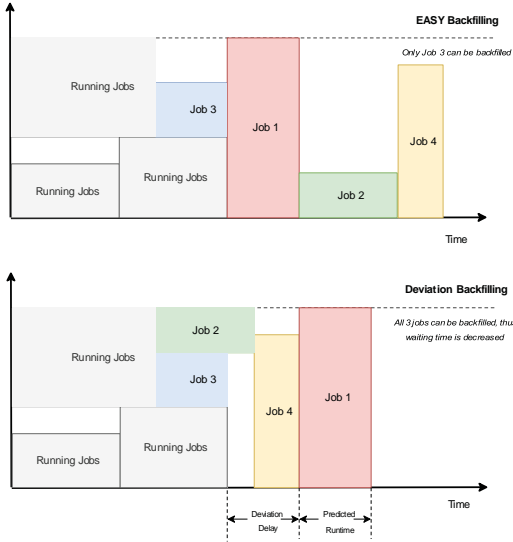


Figure 1.3 Comparison between EASY Backfilling and Deviation Backfilling.

A key aspect of the methodology lies in the incorporation of a kNN-based runtime prediction method to estimate job runtimes. This prediction method utilizes historical data to forecast the execution time of submitted jobs, facilitating more accurate scheduling decisions. By comparing these predictions with user-provided estimates, the scheme calculates the deviation, which serves as the delay threshold for determining job backfilling eligibility. Detailed algorithms provided in the paper delineate the step-by-step processes involved in calculating the delay threshold, identifying backfillable jobs, and predicting job runtimes, thereby offering a comprehensive understanding of the Deviation Backfilling scheme's implementation[40][41].

F. Enabling System-wide Shared Memory in Python for Distributed Computing with PyCOMPSs

The methodology of this study primarily focuses on extending the Python language to support system-wide shared memory and integrating this functionality into the PyCOMPSs programming model, tailored for distributed computing platforms. The initial step involved extending the SharedArray library, which provides an API for creating NumPy arrays backed by shared memory segments provided by the system. By enabling the creation of shared NumPy[37][38] arrays, the extension facilitates efficient data reuse across multiple processes, mitigating the need for costly serialization and deserialization operations. This extension enhances the interoperability of Python with system-level shared memory, fostering a more seamless integration of parallel processing paradigms into Python-based distributed computing environments.

Following the extension of the SharedArray library, the authors proceeded to integrate the shared memory functionality into the PyCOMPSs programming model, which facilitates task-based parallelism in distributed computing environments. This integration involved modifying the task decorator in PyCOMPSs to include a flag for marking task arguments as "recurrent read-only". If an argument is flagged as such and is a NumPy array, it is either loaded into the system shared memory or retrieved from it if previously added. Additionally, an internal dictionary is utilized to manage the shared memory segments for each worker process efficiently. This integration enhances the usability of PyCOMPSs[5] for data-intensive applications by enabling seamless access to shared memory resources within the Python programming environment.

The implementation details of the extended SharedArray library include the provision of a copy constructor for NumPy arrays, allowing the creation of new arrays with contiguous buffers pointing to shared memory regions. Furthermore, the behavior of the library was modified to support read-on-write operations without synchronization, thereby optimizing performance for scenarios with frequent data updates. The authors also ensured compatibility with two shared memory models: file-backed shared memory using mmap[2] and system-provided shared memory using shm_open[2][3], with the latter chosen for its superior performance characteristics. These implementation details underscore the robustness and versatility of the proposed Python extension for system-wide shared memory.

Algorithm 2 Main algorithm for blocked matrix multiplication

INPUT
dim Matrices dimension
A, B Square 2D matrices, shared

C Square 2D matrix, not shared, zero initialized

OUTPUT

C Result matrix

```

← 1: for r  0 to dim  1 do
2:   for c  0 to dim  1 do
3:     ← for i  0 to dim  1 do
4:       Crc matmul_block(Ari, Bic, Crc)
5:       end for
6:     end for
7: end for

```

To assess the performance impact of their extension, the authors conducted a comprehensive evaluation using two representative applications: k-means clustering and blocked matrix multiplication. The results demonstrated notable performance improvements, particularly in scenarios with high data reuse, such as k-means clustering, where performance gains of up to 40% were observed. However, the evaluation also revealed mixed results for other applications, with some cases showing improvements and others indicating slight performance degradation[41][43][47]. Overall, the performance evaluation underscores the potential benefits of leveraging system-wide shared memory for enhancing the efficiency of data-intensive applications in distributed computing environments.

G. Efficient In-Memory Integer Multiplication Design Using Racetrack Memory Technology

The methodology of this study revolves around designing an energy-efficient in-memory integer multiplication scheme based on racetrack memory technology[48][50]. The primary approach involves proposing a weight-based parallel Booth multiplier architecture, leveraging the unique properties of racetrack memory. Unlike conventional designs that rely on adder arrays, this architecture parallelizes the addition of partial products using racetrack memory-based adders. This innovative approach exploits the high density of racetrack memory, allowing for more efficient processing of integer multiplication operations.

The weight-based parallel addition technique employed in the architecture utilizes $2n$ adders, comprising $n-1$ full adders and 1 half adder, where n represents the bit-width of the operation. Bits of the same weight from different partial products are grouped and fed into the designated adder in a pipelined manner using FIFOs[55][58] of varying lengths, implemented with racetrack stripes. This pipelining strategy enhances the throughput of the multiplication process and optimizes resource utilization within the architecture.

A key contribution of the study is the proposal of a novel write optimization method aimed at transforming energy-intensive write operations in racetrack memory into more efficient shift operations. This optimization technique involves pre-storing 0s

ENERGY OPTIMIZATION FOR MEMORY ARCHITECTURE.

and 1s in stripes connected to input Magnetic Tunnel Junctions (MTJs) and employing a finite state machine control logic. By minimizing the need for write operations and substituting them with energy-efficient shift operations, significant reductions in energy consumption and processing time are achieved, enhancing the overall efficiency of the design.

The combined implementation of the weight-based parallel architecture and the write optimization technique results in notable improvements in throughput and energy efficiency. Through SPICE simulations using a 45nm CMOS design kit and a validated racetrack memory model, the design parameters such as radix[51][57][60], architectures, and input bit-widths are evaluated comprehensively. The simulations enable a thorough analysis of the trade-offs between throughput, area overhead, and energy consumption for different multiplier configurations, providing insights into the performance characteristics of the proposed design.

The methodology involves crafting and assessing an energy-efficient in-memory integer multiplication scheme specifically crafted for racetrack memory technology[48][49]. By introducing innovative architectural principles and optimization methods, the research endeavors to optimize racetrack memory's dense parallel capabilities while mitigating its energy-related drawbacks. Through meticulous simulation-driven assessments, the study systematically scrutinizes the efficacy and effectiveness of the proposed design, underscoring its prospective role in advancing energy-conscious computing systems.

IV. FUTURE SCOPE AND ITS CHALLENGES

Looking ahead, advancing energy optimizations for memory architectures requires addressing the evolving landscape of computing workloads and applications. As new computing paradigms emerge, such as edge computing and Internet of Things (IoT) [51][53][56] devices, the demand for energy-efficient memory solutions tailored to these applications will grow. Future research could focus on developing specialized memory architectures optimized for low-power, resource-constrained environments, while maintaining high performance and reliability. This involves exploring novel design methodologies and techniques that strike a balance between energy efficiency, performance, and scalability in diverse computing environments.

Additionally, ensuring the security and integrity of energy-efficient memory architectures is paramount[13][15] in the face of evolving cyber threats and vulnerabilities. Future research efforts could explore innovative approaches to integrating security features directly into memory architectures, thereby minimizing the energy overhead associated with traditional security mechanisms. This includes investigating hardware-based security solutions such as secure enclaves and trusted execution environments, as well as exploring software-driven approaches for enhancing memory security without compromising energy efficiency. By embedding security

directly into memory architectures, future systems can better withstand cyber attacks while maintaining optimal energy efficiency.

Furthermore, collaboration and interdisciplinary research efforts will be essential for advancing the field of energy-efficient memory architectures[57][59]. Collaborations between hardware designers, software developers, materials scientists, and other domain experts can foster innovation and drive the development of holistic solutions that address the multifaceted challenges facing energy optimization in memory systems. Interdisciplinary research initiatives can explore synergies between different technologies and disciplines, leading to breakthroughs in energy efficiency, performance, and sustainability. By fostering collaboration and knowledge exchange across diverse fields, researchers can unlock new insights and approaches to tackle the complex challenges of energy-efficient memory architectures.

Continuing research is needed to overcome the scalability challenges associated with energy-efficient memory architectures. As memory capacities continue to increase, maintaining energy efficiency becomes increasingly challenging. Future efforts could explore innovative memory hierarchies[13][10], such as hybrid memory systems that integrate volatile and non-volatile memory technologies, to optimize energy usage while meeting the demands of large-scale computing applications. Additionally, advancements in materials science and nanotechnology may offer new opportunities for developing energy-efficient memory devices with higher densities and lower power consumption.

Moreover, addressing the environmental impact of energy-efficient memory architectures is crucial for achieving sustainable computing. Future research could focus on developing recycling and disposal methods for memory components to minimize electronic waste and promote circular economy principles. Additionally, exploring renewable energy sources for powering memory-intensive computing systems can further reduce their environmental footprint. By adopting a holistic approach that considers not only energy efficiency but also environmental sustainability, future memory architectures can contribute to building a more environmentally conscious computing infrastructure.

The future direction of energy-efficient memory architectures hinges on innovation, collaboration, and interdisciplinary cooperation. By customizing memory solutions to align with the evolving computing landscape, bolstering security measures, overcoming scalability hurdles, and prioritizing environmental sustainability, researchers can steer towards an era characterized by energy-conscious and ecologically responsible computing. Through collective endeavors and cross-disciplinary alliances, the field stands poised to advance, ushering in a new epoch of energy-efficient memory architectures primed to address the needs of future computing models while minimizing their ecological footprint.

CONCLUSION

The paper showcases a range of techniques and methodologies tailored to address challenges in energy optimizations for memory architectures across diverse computing domains. By leveraging non-volatile memory technologies like Racetrack Memory and RRAM for in-memory computing, optimizing cache hierarchies, enhancing memory management in real-time embedded systems, optimizing job scheduling in HPC systems, and integrating shared memory functionality into Python-based [61]distributed computing frameworks, the paper demonstrates how these techniques can collectively improve energy efficiency, performance, and scalability while overcoming existing challenges.

Through the utilization of these techniques, various hurdles such as energy inefficiencies, resource wastage, timing guarantees, memory management complexities, and inefficient scheduling practices can be effectively mitigated. These approaches offer promising avenues for reducing energy consumption, improving performance predictability, enhancing data reuse, and optimizing specialized operations, thereby paving the way for more sustainable and efficient computing systems across different application domains.

Looking ahead, continued research and development efforts in these areas can further refine and optimize these techniques, enabling the creation of energy-efficient memory architectures that meet the evolving demands of computing paradigms while minimizing their environmental footprint. By embracing innovation, collaboration, and interdisciplinary approaches, the field can continue to advance towards a future where energy-efficient computing is not only achievable but also integral to the development of sustainable and high-performance computing systems[1][6].

REFERENCES

1. H. Kwon, A. Samajdar, and T. Krishna, "Rethinking NoCs for spatial neural network accelerators," in *Proc. 11th IEEE/ACM Int. Symp. Netw.-Chip*, Oct. 2017, pp. 1–8.
2. N. Jiang et al., "A detailed and flexible cycle-accurate network-on-chip simulator," in *Proc. IEEE Int. Symp. Perform. Anal. Syst. Softw. (ISPASS)*, Apr. 2013,
3. X. Qiao et al., "Atomlayer: A universal reRAM-based CNN accelerator with atomic layer computation," in *Proc. IEEE/ACM DAC*, Jun. 2018,.
4. Z. Zhu et al., "Mixed size crossbar based RRAM CNN accelerator with overlapped mapping method," in *Proc. IEEE/ACM ICCAD*, Nov. 2018, pp. 1–8.
5. A. Lottarini et al., "Master of none acceleration:
6. A comparison of accelerator architectures for analytical query processing," in *Proc. ACM/IEEE ISCA*, Jun. 2019, pp. 762–773.

7. S. S. Parkin, M. Hayashi, and L. Thomas, "Magnetic domain-wall racetrack memory," *Science*, vol. 320, no. 5873, pp. 190–194, 2008.
8. Zhao, Xia, Yao Guo, Qing Feng, and Xiangqun Chen. "A system context-aware approach for battery lifetime prediction in smartphones." In *Proceedings of the 2011 ACM Symposium on Applied Computing*, pp. 641–646. 2011.
9. S. S. Parkin, M. Hayashi, and L. Thomas, "Magnetic domain-wall racetrack memory," *Science*, vol. 320, no. 5873, pp. 190–194, 2008.
10. T. Luo, W. Zhang, B. He, and D. Maskell, "A racetrack memory based in-memory booth multiplier for cryptography application," in *2016 21st Asia and South Pacific Design Automation Conference (ASP-DAC)*. IEEE, 2016, pp. 286–291.
11. D. Baran, M. Aktan, and V. G. Oklobdzija, "Energy efficient implementation of parallel cmos multipliers with improved compressors," in *Proceedings of the 16th ACM/IEEE international symposium on Low power electronics and design*. ACM, 2010, pp. 147–152.
12. S. Nangate, "California (2008). 45nm open cell library," URL; <http://www.nangate.com>, 2008.
13. Y. Zhang, W. Zhao, D. Ravelosona, J.-O. Klein, J. Kim, and C. Chappert, "Perpendicular-magnetic-anisotropy cofeb racetrack memory," *Journal of Applied Physics*, vol. 111, no. 9, p. 093925, 2012.
14. C. Zhang, G. Sun, W. Zhang, F. Mi, H. Li, and W. Zhao, "Quantitative modeling of racetrack memory, a tradeoff among area, performance, and power," in *Design Automation Conference (ASP-DAC)*, 2015 20th Asia and South Pacific. IEEE, 2015, pp. 100–105.
15. X. Dong, C. Xu, Y. Xie, and N. P. Jouppi, "Nvsim: A circuit-level performance, energy, and area model for emerging nonvolatile memory," *IEEE TCAD*, 2012.
16. T. Luo, B. He, W. Zhang, and D. L. Maskell, "A novel two-stage modular multiplier based on racetrack memory for asymmetric cryptography," in *2017 IEEE/ACM International Conference on Computer-Aided Design (ICCAD)*. IEEE, 2017, pp. 276–282.
17. Zhang Z, Wang Z, Shi T, et al. Memory materials and devices: From concept to application[J]. *InfoMat*, 2020, 2(2): 261-290.
18. L. Han, P. Huang, Z. Zhou, Y. Chen, X. Liu and J. Kang, "A Convolution Neural Network Accelerator Design with Weight Mapping and Pipeline Optimization," *2023 60th ACM/IEEE Design Automation Conference (DAC)*, San Francisco, CA, USA, 2023, pp. 1-6, doi: 10.1109/DAC56929.2023.10247977.
19. Wang, Yuan, et al. "A stable rhombohedral phase in ferroelectric Hf(Zr)1+ xO2 capacitor with ultralow coercive field." *Science* 381.6657 (2023): 558-563.
20. Zhang, Wenbin, et al. "Edge learning using a fully integrated neuro-inspired memristor chip." *Science* 381.6663 (2023): 1205-1211.
21. Jung, Seungchul, et al. "A crossbar array of magnetoresistive memory devices for in-memory computing." *Nature* 601.7892 (2022): 211-216.
22. Rasch, Malte J., et al. "Hardware-aware training for large-scale and diverse deep learning inference workloads using in-memory computing-based accelerators." *Nature Communications* 14.1 (2023): 5282.
23. Jiang, Mingrui, et al. "Efficient combinatorial optimization by quantum-inspired parallel annealing in analogue memristor crossbar." *Nature Communications* 14.1 (2023): 5927.
24. J. Kang, Z. Yu, L. Wu, Z. Wang, Y. Cai, R. Huang, et al., "Time-dependent variability in RRAM-based analog neuromorphic system for pattern recognition," *2017 IEEE International Electron Devices Meeting (IEDM)*, San Francisco, CA, USA, 2017, pp. 6.4.1-6.4.4.
25. Yu, Zhizhen, et al. "Early-stage fluctuation in low-power analog resistive memory: impacts on neural network and mitigation approach." *IEEE Electron Device Letters* 41.6 (2020): 940-943.
26. F. Lordan, E. Tejedor, J. Ejarque, R. Rafanell, J. Álvarez, F. Marozzo,
27. D. Lezzi, R. Sirvent, D. Talia, and R. M. Badia, "Servicess: An interoperable programming framework for the cloud," *Journal of grid computing*, vol. 12, no. 1, pp. 67–91, 2014.
28. E. Tejedor, Y. Becerra, G. Alomar, A. Queralt, R. M. Badia, J. Torres,
29. T. Cortes, and J. Labarta, "Pycompss: Parallel computational workflows in python," *The International Journal of High Performance Computing Applications*, vol. 31, no. 1, pp. 66–82, 2017.
30. M. Mirmont, "Python extension: Sharedarray," <https://pypi.org/project/SharedArray/>, 2019.
31. OpenMP Architecture Review Board, "OpenMP Application Programming Interface Specification," Web page at <http://www.openmp.org/specifications/>, (Date of last access: 3rd May, 2017).
32. L. Dagum and R. Menon, "Openmp: an industry standard api for shared-memory programming," *Computational Science & Engineering, IEEE*, vol. 5, no. 1, pp. 46–55, 1998.
33. B. Chapman, T. Curtis, S. Pophale, S. Poole, J. Kuehn, C. Koelbel, and
34. L. Smith, "Introducing openshmem: Shmem for the pgas community," in *Proceedings of the Fourth Conference on Partitioned Global AddressSpace Programming Model*, 2010, pp. 1–3.
35. M. Driscoll, A. Kamil, S. Kamil, Y. Zheng, and K. Yelick, "PyGAS: A partitioned global address space extension for python," Poster abstract, 2012. [Online]. Available: <http://web.eecs.umich.edu/~akamil/papers/pgas12.pdf>
36. J. Álvarez Cid-Fuentes, P. Álvarez, R. Amela, K. Ishii, R. K. Morizawa, and R. M. Badia, "Efficient development of high performance data analytics in python," *Future Generation Computer Systems*, 2019. [Online]. Available: <http://www.sciencedirect.com/science/article/pii/S0167739X18321393>
37. S. J. van der Walt, S. C. Colbert, and G. Varoquaux, "The numpy array: A structure for efficient numerical computation," *Computing in Science Engineering*, vol. 13, no. 2, pp. 22–30, March 2011.
38. BSC - CNS. (2012) MareNostrum 3. [Online]. Available: <https://www.bsc.es/marenostrum/marenostrum/mn3>
39. D. Computadors, V. Pillet, J. Labarta, T. Cortes, and S. Girona, "Paraver: A tool to visualize and analyze parallel code," *WoTUG-18*, vol. 44, 03 1995.
40. J. D. Ullman, "NP-complete scheduling problems," *Journal of Computer and System Sciences*, vol. 10, no. 3, pp. 384–393, jun 1975.
41. P.-F. Dutot, M. Mercier, M. Poquet, and O. Richard, "Batsim: a Realistic Language-Independent Resources and Jobs Management Systems Simulator," in *20th Workshop on Job Scheduling Strategies for Parallel Processing*, Chicago, United States, May 2016.
42. H. Casanova, "Simgrid: A toolkit for the simulation of application scheduling," in *Proceedings of the 1st International Symposium on Cluster Computing and the Grid*, ser. CCGRID '01. USA: IEEE Computer Society, 2001, p. 430.
43. D. G. Feitelson, D. Tsafir, and D. Krakov, "Experience with using the parallel workloads archive," *Journal of Parallel and Distributed Computing*, vol. 74, no. 10, pp. 2967 – 2982, 2014.
44. D. G. Feitelson, "Metrics for parallel job scheduling and their convergence," in *Revised Papers from the 7th International Workshop on Job Scheduling Strategies for Parallel Processing*, ser. JSSPP '01. Berlin, Heidelberg: Springer-Verlag, 2001, p. 188–206.
45. S. Mittal, J. S. Vetter, and D. Li, "A Survey of architectural approaches for managing embedded DRAM and non-volatile on-chip caches," *IEEE Trans. Parallel Distrib. Syst.*, vol. 26, no. 6, pp. 1524–1537, Jun. 2015.
46. G. C. Buttazzo, *Hard Real-Time Computing Systems: Predictable Scheduling Algorithms and Applications*, vol. 24, Berlin, Germany: Springer, 2011.

48. D. Gao, D. Reis, X. S. Hu, and C. Zhuo, "Eva-CiM: A system-level energy evaluation framework for computing-in-memory architectures," 2019, *arXiv: 1901.09348*.
49. Y. Xiang and S. Pasricha, "A hybrid framework for application allocation and scheduling in multicore systems with energy harvesting," in *Proc. ACM Great Lakes Symp. VLSI*, 2014, pp. 163–168.
50. J. Zhou, T. Wei, M. Chen, J. Yan, X. S. Hu, and Y. Ma, "Thermal-aware task scheduling for energy minimization in heterogeneous real-time MPSoC systems," *IEEE Trans. Comput.-Aided Design Integr. Circuits Syst.*, vol. 35, no. 8, pp. 1269–1282, Aug. 2016.
51. E. C. R. Shin, D. Song, and R. Moazzezi, "Recognizing functions in binaries with neural networks," in *Proc. 24th USENIX Secur. Symp.*, 2015, pp. 611–626.
52. Z. L. Chua, S. Shen, P. Saxena, and Z. Liang, "Neural nets can learn function type signatures from binaries," in *Proc. 26th USENIX Secur. Symp.*, 2017, pp. 99–116.
53. X. Xu, C. Liu, Q. Feng, H. Yin, L. Song, and D. Song, "Neural network-based graph embedding for cross-platform binary code similarity detection," in *Proc. ACM SIGSAC Conf. Comput. Commun. Secur.*, Oct. 2017, pp. 363–376.
54. W. Guo, D. Mu, J. Xu, P. Su, G. Wang, and X. Xing, "LEMNA: Explaining deep learning based security applications," in *Proc. ACM SIGSAC Conf. Comput. Commun. Secur.*, Oct. 2018, pp. 364–379.
55. R. L. Rivest, A. Shamir, and L. Adleman, "A method for obtaining digital signatures and public-key cryptosystems," *Commun. ACM*, vol. 26, no. 1, pp. 96–99, Jan. 1983.
56. P. W. Shor, "Algorithms for quantum computation: Discrete logarithms and factoring," in *Proc. 35th Annu. Symp. Found. Comput. Sci.*, 1994, pp. 124–134.
57. D. J. Bernstein and L. Tanja, "Post-quantum cryptography," *Nature*, vol. 549, no. 7671, pp. 188–194, 2017.
58. L. Chen, L. Chen, S. Jordan, Y.-K. Liu, D. Moody, R. Peralta, R. A. Perlner, and D. Smith-Tone, "Report on post-quantum cryptography, volume 12," US Dept. Commerce, Nat. Inst. Standards Technol., Gaithersburg, MD, USA, Tech. Rep. NIST IR 8105, 2016.
59. S. Lapshev and S. Hasan, "New Low Glitch and Low Power DET Flip-Flops Using Multiple C-Elements", *IEEE Transactions on Circuits and Systems I: Regular Papers*, vol. 63, no. 10, pp. 1673–1681, 2016.
60. Y. Li, H. Wang, R. Liu, L. Chen, I. Nofal, Q. Chen, A. He, G. Guo, S. Baeg, S. Wen, R. Wong, Q. Wu and M. Chen, "A 65 nm Temporally Hardened Flip-Flop Circuit", *IEEE Transactions on Nuclear Science*, vol. 63, no. 6, pp. 2934–2940, 2016.
61. N. Kulkarni, J. Yang, J. Seo and S. Vrudhula, "Reducing Power, Leakage, and Area of Standard-Cell ASICs Using Threshold Logic Flip-Flops", *IEEE Transactions on Very Large Scale Integration (VLSI) Systems*, vol. 24, no. 9, pp. 2873–2886, 2016.