

A Classification of Memory-Centric Computing

Du Nguyen, H.A.; Yu, J.; Abu Lebdeh, M.F.M.; Taouil, M.; Hamdioui, S.; Catthoor, Francky

DOI

[10.1145/3365837](https://doi.org/10.1145/3365837)

Publication date

2020

Document Version

Final published version

Published in

ACM Journal on Emerging Technologies in Computing Systems

Citation (APA)

Du Nguyen, H. A., Yu, J., Abu Lebdeh, M. F. M., Taouil, M., Hamdioui, S., & Catthoor, F. (2020). A Classification of Memory-Centric Computing. *ACM Journal on Emerging Technologies in Computing Systems*, 16(2), 1-26. Article 13. <https://doi.org/10.1145/3365837>

Important note

To cite this publication, please use the final published version (if applicable).
Please check the document version above.

Copyright

Other than for strictly personal use, it is not permitted to download, forward or distribute the text or part of it, without the consent of the author(s) and/or copyright holder(s), unless the work is under an open content license such as Creative Commons.

Takedown policy

Please contact us and provide details if you believe this document breaches copyrights.
We will remove access to the work immediately and investigate your claim.

Green Open Access added to TU Delft Institutional Repository

'You share, we take care!' - Taverne project

<https://www.openaccess.nl/en/you-share-we-take-care>

Otherwise as indicated in the copyright section: the publisher is the copyright holder of this work and the author uses the Dutch legislation to make this work public.



A Classification of Memory-Centric Computing

HOANG ANH DU NGUYEN, JINTAO YU, MUATH ABU LEBDEH,
MOTTAQIALLAH TAOUIL, and SAID HAMDIOUI, Delft University of Technology
FRANCKY CATTLOOR, Inter-university Micro-Electronics Center (IMEC)

13

Technological and architectural improvements have been constantly required to sustain the demand of faster and cheaper computers. However, CMOS down-scaling is suffering from three technology walls: leakage wall, reliability wall, and cost wall. On top of that, a performance increase due to architectural improvements is also gradually saturating due to three well-known architecture walls: memory wall, power wall, and instruction-level parallelism (ILP) wall. Hence, a lot of research is focusing on proposing and developing new technologies and architectures. In this article, we present a comprehensive classification of memory-centric computing architectures; it is based on three metrics: computation location, level of parallelism, and used memory technology. The classification not only provides an overview of existing architectures with their pros and cons but also unifies the terminology that uniquely identifies these architectures and highlights the potential future architectures that can be further explored. Hence, it sets up a direction for future research in the field.

CCS Concepts: • **Computer systems organization** → *Special purpose systems*; • **Hardware** → *Spintronics and magnetic technologies*;

Additional Key Words and Phrases: Computation-in-memory, resistive computing, memory-centric computer architectures

ACM Reference format:

Hoang Anh Du Nguyen, Jintao Yu, Muath Abu Lebdeh, Mottaqiallah Taouil, Said Hamdioui, and Francky Catthoor. 2020. A Classification of Memory-Centric Computing. *J. Emerg. Technol. Comput. Syst.* 16, 2, Article 13 (January 2020), 26 pages.
<https://doi.org/10.1145/3365837>

1 INTRODUCTION

For several decades, technology scaling has provided a 43% performance gain for each successive node and cheaper computers as a result of a higher operating frequency and lower cost per transistor, respectively [15, 54]. On top of that, smart architectural improvements such as pipelining and cache hierarchies have increased computer performance up to 50% every 2 years [49]. However, CMOS scaling suffers from three main walls: leakage wall, reliability wall, and cost wall [45], while computer architectures also face three walls: memory wall, power wall, and instruction-level

The results presented in this article have been obtained in the framework of the project “Computation-in-memory architecture based on resistive devices” (MNEMOSENE), which has received funding from the European Union’s Horizon 2020 research and innovation programme under grant agreement No 780215.

Authors’ addresses: H. A. D. Nguyen, J. Yu, M. A. Lebdeh, M. Taouil, S. Hamdioui, and F. Catthoor, Mekelweg 4, 2628 CD, Delft, the Netherlands; emails: {H.A.DuNguyen, J.Yu-1, M.F.M.AbuLebdeh, M.Taouil, S.Hamdioui}@tudelft.nl, Francky.Catthoor@imec.be.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

© 2020 Association for Computing Machinery.

1550-4832/2020/01-ART13 \$15.00

<https://doi.org/10.1145/3365837>

parallelism (ILP) wall [100]. In order to address these walls, novel technologies and architectures are under research to improve the performance [54]. As a result, an enormous amount of computer architectures have been proposed recently. Therefore, a complete classification of these architectures is needed, not only to have a useful way of describing and comparing them, but also to have a clear view about what has been explored and what has not been explored yet.

Limited work has addressed this problem. Most of the well-known classifications separate the processors from the memory. Therefore, these classifications often are processor-centric-based architectures, such as Flynn's [35], Skillicorn's [117], and Shami-Hemani's classification [112]. Although these classifications work well for processor-centric architectures proposed in the past decades, they are not applicable to the emerging memory-centric architectures. Other small-scale surveys mostly target a specific type of computer architecture such as vector processors, automata processors, or processing-in-memory architectures [23, 58, 69, 108, 113, 122, 125, 126]. These surveys only discuss a limited part of the computer architecture classification and, in addition, do not contain the complete space of both conventional processor-centric architectures and memory-centric architectures. Therefore, these surveys often make no distinction between processing inside and near the memory. This leads to a confusion in terminology (e.g., processing-in-memory, logic-in-memory, in-memory computing, near-memory computing, etc.). For example, Hybrid Memory Cube is considered to be near-memory computing [101]; however, it is also referred to as processor-in-memory [3]. Some recent classifications and reviews did mention those architectures in the context of technology development [94, 136]. However, these papers mostly targeted the technological feasibility instead of the characteristics and variants of such computer architectures. In addition to the above, there are some architecture-related papers that briefly discussed the features of emerging architectures [89, 91, 106]. However, they are incomplete, focus mostly on relatively narrow aspects, and only classify the architectures based on applications [89] and logic design methods [91, 106]. In short, there is still a lack of systematic and complete classification that focuses on memory-centric computing or computer architectures in general. This is exactly the target of this article.

This article presents a comprehensive classification of memory-centric computing and discusses both conventional and emerging computing architectures. The classification is based on three metrics: computation location, memory technology, and computation parallelism. The computation location indicates where the computations are performed (e.g., near or far from the memory) and provides an insight regarding the severeness of the memory wall. The memory technology, which provides characteristics of the memory, can enable new computer architectures (e.g., resistive computing). The computation parallelism specifies the type of parallelism that can be exploited in an architecture (e.g., task-level parallelism). With these distinct metrics, the classification covers *all* computing architectures in general and memory-centric computing in specific. Note, however, that it does not make previous proposed classifications obsolete, as they typically target specific subclasses. In short, the contributions of this article are the following:

- Unify the terminology for computer architectures such that it is applicable to all computing paradigms including conventional, in-memory, and near-memory computing.
- Propose a complete classification that includes both existing and emerging architectures.
- Explain one representative architecture of each subclass in detail.
- Discuss and evaluate the main advantages and disadvantages of the different classes and selected architectures.
- Highlight the whole space of memory-centric computing, including the nonexplored architectures.

The rest of this article is structured as follows. Section 2 shows the metrics used in the classification, briefly introduces the four classes, and provides a quantitative comparison among them.

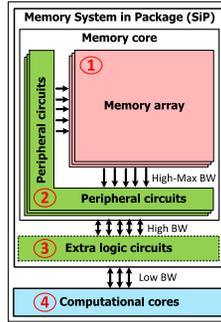


Fig. 1. Computer architecture.

Sections 3, 4, and 5 present the characteristics of the three memory-centric computing classes; the fourth class contains the traditional von Neuman architectures and is out of the scope of this article. Section 6 discusses the pros and cons of this classification and compares it with existing ones. Finally, Section 7 concludes this article.

2 CRITERIA AND CLASSIFICATION

In this section, we first present the set of metrics to classify computer architectures. Thereafter, we map the existing architectures on our classification. Finally, we compare the classes qualitatively based on their most important metrics.

2.1 Classification Metrics

We propose several metrics to classify computer architectures based on the computing resources and memory. A computer architecture or system consists of (one or more) memories and (one or more) computational units as shown in Figure 1. The memories can reside in a core (i.e., memory core) or System in Packages (SiP). A memory core consists of one or more cell arrays (used for storage) and peripheral circuits (used to access the memory cells). Note that register files and caches are not considered as storage here, as they are optimized for speed with relatively small capacity and temporary storage [49]. Hence, the long-term storage of data takes place in the higher layers such as main memory and solid-state disks. Traditionally, the computing takes place in the computational cores. However, recently architectures with computing power in the memory have been proposed [46, 99, 101]. In case the memory contains additional logic circuits such as in Hybrid Memory Cubes (HMCs) [101], we speak of an SiP. With an increasing distance from the main memory array, the available bandwidth (specified by BW in Figure 1) reduces; note that the bandwidth here is related to the memory bottleneck and will be discussed further in Section 2.3. Based on these definitions, the following metrics are used to classify computer architectures: computation location, memory technology, and computation parallelism; they are discussed next.

Computation location: This indicates *where the result of the computation is produced*. A computation is defined here as a primitive logic function (e.g., logical operations) or arithmetic operation (e.g., addition, multiplication). Figure 1 indicates the four possibilities where a computation result can be produced; they can be identified by four circled numbers. If the result is produced *within* the memory core (i.e., the computing takes place within one of the memories), the computer architecture is referred to as *Computation-Inside-Memory (CIM)*. If the result is produced outside the memory core, the architecture is referred to as *Computation-Outside-Memory (COM)*. Both CIM and COM can be further subclassified.

It is worth stressing that **CIM** architectures perform computations *within* the memory core. As already mentioned, the memory consists of a *memory array* and the *peripheral circuits*. Specifically,

depending on *where* the result of the computation is produced, CIM architectures can be divided into two basic subclasses. These subclasses can be combined into many hybrid combinations. We will describe this large space by focusing first on its two extreme sides:

- *CIM-Array (CIM-A)*: In CIM-A, the computing result is produced within the array (noted as position 1 in Figure 1). Note that this is different from a standard write operation. Typical examples of CIM-A architectures use memristive logic designs such as MAGIC and imply [66, 72]. CIM-A architectures always require a redesign of cells to support such logic design, as the conventional memory cell dimensions and their embedding in the bit- and wordline structure do not allow them to be used for logic. A memory cell is namely heavily optimized in terms of processing stack and layout; hence, any changes in the array access require a completely new cell design and characterization process as the material stack of a memory array is specifically optimized for specific control voltages, current, and so forth. In addition, modifications in the periphery are sometimes needed to support the changes in the cell. Therefore, CIM-A architectures can be subdivided into two groups: (1) basic CIM-A, where only changes inside the memory array are required, and (2) hybrid CIM-A, where, in addition to major changes in the memory array, minimal to medium changes are required in the peripheral circuit. An example of basic CIM-A is an architecture that performs computations using implication logic [76]. In this logic style, only one memory row is activated at a time, and a number of columns (bits) are read out through sense amplifiers. Hence, due to the same usage as in normal memory, the peripheral circuits do not require any modifications. An example of hybrid CIM-A is an architecture that performs computations using MAGIC [72]. In this case, multiple memory rows are written simultaneously; due to the high write currents, modifications are required to the cell and medium changes in the peripheral circuits are needed to activate the multiple rows.
- *CIM-Periphery (CIM-P)*: In CIM-P, the computing result is produced within the peripheral circuitry (noted as position 2 in Figure 1). Typical examples of CIM-P architectures contain logical operations and vector-matrix multiplications [21, 81, 134]. CIM-P architectures typically contain dedicated peripheral circuits such as DACs and/or ADCs [37, 111] and customized sense amplifiers [81, 134]. Note that more radical changes in the peripheral circuit can be made in the future (e.g., changing in control voltages leads to radical changes in voltage drivers and sense amplifiers, or including a full functional processor inside memory banks). Even though the computational results are produced in the peripheral circuits for CIM-P, the memory array is a substantial component in the computations. As the peripheral circuits are modified, the currents/voltages applied to the memory array are typically different than in the conventional memory. Hence, similarly to the CIM-A subclasses, the CIM-P architectures are also further divided into two groups: (1) basic CIM-P, where only changes inside the peripheral is required, which means the current levels should not be affected, and (2) hybrid CIM-P, where the majority of the changes take place in the peripheral circuit and minimal to medium changes in the memory array. An example of basic CIM-P is Pinatubo logic [81]. Pinatubo activates two or more (but not many) rows of a memory array simultaneously during read operations for computations; in addition to a customized sense amplifier to perform the logic operation, this architecture also requires modifications in the address decoder to activate several rows. Note, however, that modifications in the cell/array are not required as the total read current is still small. An example of hybrid CIM-P is ISAAC [111]. ISAAC activates all rows of a memory array at the same time during read operations to perform a matrix vector product using an ADC readout circuit. This architecture accumulates currents in the bitline that impose higher electrical loading in the

memory array; hence, not only is the periphery circuit heavily modified but also the cell requires changes due to the high bitline current.

The difference between CIM-A and CIM-P classes is the location of producing results. The results of CIM-A architectures are produced inside the memory array, which may sometimes require readout operations to obtain the results for further calculations; instead, in CIM-P the results are obtained directly after the operations and may sometimes need an additional step to write the results back to memory. In order to perform computations, both subclasses impact the design of the memory core. However, in many/most cases both the cell and the peripheral circuitry require changes; i.e., they are hybrids. In case these changes affect mostly the cell, we speak of hybrid CIM-A, otherwise hybrid CIM-P.

In **COM classes**, computations take either place in the extra logic circuits inside the memory SiP (noted as position 3 in Figure 1) or in the traditional computational cores (noted as position 4 in Figure 1) such as CPU, FPGA, and so forth. In case of the former, the computations take place near the memory core and the architecture is referred to as Computation-Outside-Memory Near (COM-N). In case of the latter, the architecture is referred to as Computation-Outside-Memory Far (COM-F). Note that the bandwidth is still high for COM-N as compared to COM-F, but lower than CIM-A and CIM-P.

Note that architectures where the computation takes place in different places (e.g., array and peripheral) are called composite architectures. Hence, they are compositions of the leaf nodes in our classification tree. In addition, an architecture could have multiple primitive functions, each with a different computation location. Also, these architectures are considered to be composite.

In addition to the computation location, which specifies where the results are produced, it is possible to further divide the classes using the computation method by specifying how the computation is performed. For example, CIM-A often uses memristor-based computations such as IMPLY [14, 115], Snider [118], and MAGIC [72]. CIM-P often uses current summations such as Scouting logic [134], Ambit [110], and Pinatubo [81]. However, this metric is not included in the classification for two reasons: (1) it is strongly coupled to the computational location, and (2) it makes the classification too complex and hence loses its simplicity. Nevertheless, including such a metric can be complementary to our work. A further subclassification based on this metric can be based on existing classifications as shown in [26, 106].

Memory technology: This indicates which technology is used to implement the memory array. The technologies are either conventional charge-based memories such as DRAM/SRAM [86, 87, 93] or emerging non-charge-based memories [107]. The non-charge-based memories can be further divided into different types based on their physical mechanism: resistive memories [74, 107], “magnetic” memories [10, 19, 107], molecular memories [41, 78, 79, 103], or mechanical memories [16, 42]. Resistive memories store the data as a resistance value; it includes Resistive RAM (RRAM) [129], phase change memory (PCM) [74], and the like. The resistance in RRAM is determined by the presence or absence of a conductive filament between its two electrodes [107], while the resistance in PCM relies on a change between amorphous and crystalline phases [105, 130]. Magnetic memories, such as Magnetic RAM (MRAM), store the data using the magnetization direction of the free layer with respect to the hard or reference layer; it includes, for example, conventional magnetic RAM [140] and STT-MRAM [36, 51]. The resistive and magnetic memories are organized in crossbars with cells placed at each junction. The other types of memories (i.e., molecular memories, mechanical memories) have not been shown to be useful for computing yet; hence, they are not discussed further in this classification. It is worth mentioning that each of these memory technologies has its own characteristics (read/write latency, endurance, capacity, etc.) and therefore are deployed at different levels in the memory hierarchy [107]. Therefore, the

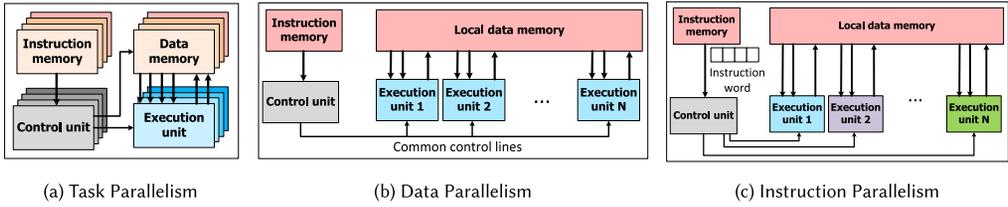


Fig. 2. Three types of computation parallelism.

memory technology dictates not only which CIM operations are technology-wise feasible but also where in the memory hierarchy they take place.

Computation parallelism: This indicates the level of parallelism that can be exploited in a computer system, i.e., task, data, and/or instruction level parallelism, as shown in Figure 2. An architecture supports task-level parallelism when it contains multiple independent control units and multiple data memories (see Figure 2(a)). The independent control units can be used to execute multiple threads or instruction sequences from the same application concurrently; examples are multithreading [30, 124] and multicore systems [40]. In data parallelism, a single control unit is used to apply the *same* instruction concurrently on a collection of data elements (see Figure 2(b)); note that all execution units share the same control signals. The data elements can be processed using constant sizes (e.g., vector and array processor [28, 33]) or varying subword sizes (e.g., SWAR (SIMD Within A Register) processor [34, 102]). In instruction-level parallelism, a single control unit is used to execute *various* instructions concurrently (see Figure 2(c)); hence, the execution units have different control signals. A further distinction can be made based on intra-instruction (e.g., pipelined processor [123]) or interinstruction (e.g., VLIW processor [131]) parallelism, or a combination of them (e.g., speculative processor [88]).

The three above-mentioned metrics (i.e., computation location, memory technology, and computation parallelism) are dependent on each other. The computation location has a big impact on the feasibility of the other two metrics, for example, realizing ILP in CIM-A is quite difficult or realizing COM-N with SRAM is not economically feasible. Also, the parallelism is not fully independent from the computation location and memory technology. For example, data parallelism is often applied straightforwardly in CIM-A and CIM-P [27, 37]; however, it is difficult to realize ILP in CIM-A, while it is much easier in COM-N and COM-F due to the intrinsic pipeline stages in conventional processors. The computation parallelism is also affected by the technology as the technology poses restrictions on the endurance. For example, exploiting ILP in CIM-A architectures demands a high endurance as more writes are required to store immediate stages and, hence, are not attractive for emerging memories like RRAM and PCM with endurance limitations.

2.2 Classification

We classify the existing architectures based on the above discussed metrics; the result is shown in Figure 3. The references of the abbreviated architectures are listed in Table 1.

The classification contains 48 categories. Some categories, the ones located in *red* planes, show that a *lot* of work has been done for that particular class. For the categories in the *pink* planes, a *moderate number* of work has been done. To the best of our knowledge, no architectures exist in the *blue* planes; these fields are currently unexplored as they have received no attention yet from the research community or are nonexistent due to current restrictions of the technology; these blue planes are not further discussed in this article for two main reasons: (1) scope of the article: the technical and economical feasibility of these planes requires an intensive discussion and is by itself a new contribution and (2) space limitations. Later on, we will discuss several architectures from the red and pink planes.

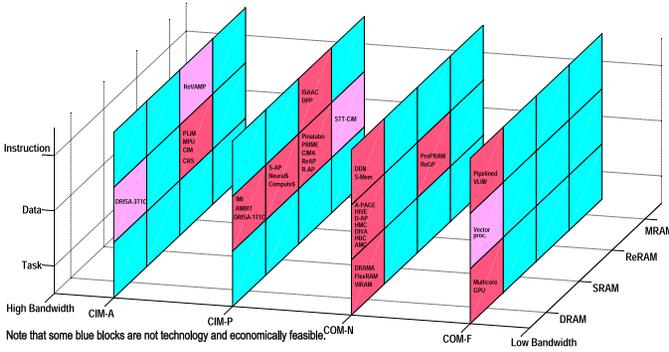


Fig. 3. Memory-centric computing classification.

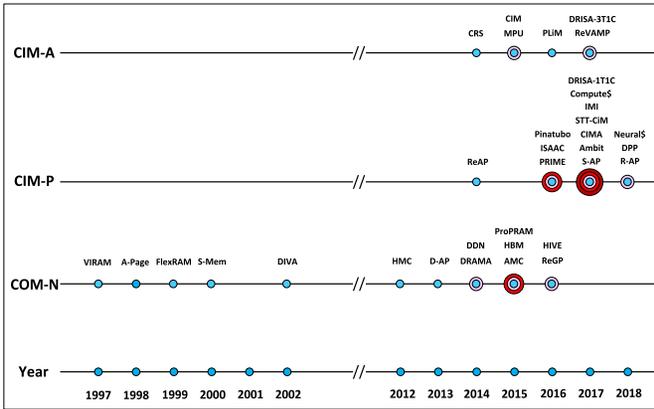


Fig. 4. Memory-centric computing timeline.

Table 1. Abbreviation List

Abbreviation	Reference
DRISA-3T1C	[80]
ReVAMP	[9]
PLiM	[38]
MPU	[52]
CIM	[27, 46]
CRS	[115]
ISAAC	[111]
DPP	[37]
IMI	[32]
AMBIT	[110]
DRISA-1T1C	[80]
S-AP	[121]
Neural\$	[29]
Compute\$	[1]
Pinatubo	[81]
PRIME	[21]
CIMA	[26]
ReAP	[137]
R-AP	[138]
STT-CiM	[56]
DDN	[20]
S-Mem	[85]
A-PAGE	[98]
HIVE	[3]
D-AP	[96]
DIVA	[24]
HMC	[50]
AMC	[95]
HBM	[84]
DRAMA	[31]
FlexRAM	[63]
VIRAM	[70]
ProPRAM	[128]
ReGP	[92]
Pipelined	[49, 123]
VLIW	[83, 131]
Vector Proc.	[22, 33, 102]
Multicore	[53]
GPU	[67, 97]

The developments in memory-centric computing are shown in the timeline of Figure 4; this shows the trend of computing moving from COM-F to COM-N, CIM-A, and CIM-P. In the figure, a larger circle indicates that more work has been proposed in that year. Note that the conventional architectures in COM-F are not memory centric and, hence, are left out. The concept of merging computation and memory was introduced back in 1970 [120]. This concept became popular around

Table 2. Comparison among Architecture Classes in Terms of Data Movement, Computation Requirements, Available Bandwidth, Memory Design Efforts, and Scalability

	Data Movement outside Memory Core	Computation Requirements		Available Bandwidth	Memory Design Efforts			Scalability
		Data Alignment	Complex Function		Cells & Array	Periphery	Controller	
CIM-A	No	Yes	High latency	Max	High	Low/medium	High	Low
CIM-P	No	Yes	High cost	High-Max	Low/medium	High	Medium	Medium
COM-N	Yes	NR	Low cost	High	Low	Low	Low	Medium
COM-F	Yes	NR	Low cost	Low	Low	Low	Low	High

NR: Not required.

1997 in COM-N architectures and was further developed from 2002. These COM-N architectures, such as VIRAM [70] (initially named IRAM), DIVA [24], or FlexRAM [63], never commercialized due to the limitations of embedded DRAM technology (i.e., costly fabrication process and inefficient speed and memory capacity trade-off [55, 64, 65]). After that, a long silent period in academia community was observed from 2002 to 2010. Meanwhile, industrial efforts have been invested to deploy large eDRAM in commercial COM-N systems such as the POWER7 processor [61], PlayStation2 [7], and Intel's top-class CPUs [71]. From 2012 to 2016, new commercial COM-N architectures based on novel 3D stacking technology were proposed such as Hybrid Memory Cube (HMC) [50] and High-Bandwidth Memory (HBM) [84]. In the last several years, with the emergence of resistive technology, CIM-A and CIM-P architectures started to become popular.

Note that many of the architectures are hybrid and/or composite, which means that they can map into multiple classes. In order to simplify Figure 3, the architectures are classified based on their dominant features. For example, DPP exploits both ILP and DLP; however, DPP focuses more on performing various parallel operations using multiple functional units, while it also processes a whole row/column inside the memory; hence, the dominant feature of DPP parallelism is selected as ILP.

2.3 Qualitative Evaluation

In this subsection, we briefly compare the different computing types qualitatively using the most important classification metric, i.e., the computation location. This metric dictates the type of data movements, computation requirements, available bandwidth, memory design efforts, scalability, endurance requirements, and maturity. With respect to the computation requirements, we discuss whether the architectures require a specific data alignment and whether they have the capability to realize complex functions. With respect to available bandwidth, we discuss the capability for data communication between logic and storage units. With respect to the memory design efforts, we discuss the modifications that are required for the cells, array, peripheral circuit, and controller. With respect to the scalability, we discuss the possibility to expand the design to increase the concurrent computing capacity. With respect to the endurance requirement, we indicate the endurance level that the architecture demands in order to execute an application. With respect to maturity, we not only mean the readiness of the memory technology but also the available programming and software support, and current status (i.e., simulations, prototype, or fabrication) for such architectures. With respect to the applications, we roughly indicate the range of applications for each architecture class. The results are shown in Tables 2 and 3; their content will be discussed next.

Data movement outside the memory core indicates whether the data will remain in the memory core during computing or transferred to outside computational units. It affects the

Table 3. Comparison among Architecture Classes in Terms of Endurance Requirement, Maturity, and Applications

	Endurance Requirement	Maturity		Applications
		Software Support and Technology	Development	
CIM-A	High	Emerging	Simulation	Data intensive - Computational complexity (matrix multiplication [48], parallel addition [27])
CIM-P	Medium	Emerging	Simulation	Data intensive - Bitwise operations (database processing [81, 109, 110], graph processing [2], image processing [47, 110], security and biosequencing application [8])
COM-N	Medium	Commercialized	Fabricated	General purpose and Application specific (vector processing [25, 59, 95], automata processing [96], neural computation [20])
COM-F	Low	Common Practise	Fabricated	General purpose

memory bottleneck due to latency and the energy consumption of data transfers. Both CIM-A and CIM-P architectures have a relatively low amount of data movement outside the memory core, as the processing occurs inside the memory core. Therefore, they have the potential to alleviate the memory bottleneck. Instead of moving data from the memory to the computational cores, the instructions are moved and directly applied to the memory; these instructions typically operate on a large dataset, and hence a high level of parallelism can be obtained. Note, however, that the current state of the art typically allows limited functions to be implemented in these architectures. Therefore, complex functions would still require data movements to the computational cores outside the memory. For COM-N and COM-F architectures, data is first read from the memory. Thereafter, they are typically stored in registers before being fed to the processing units. The amount of parallelism is limited here due to constraints in the bandwidth and number of available registers and processing units.

Computation requirements include data alignment and the ability to implement complex functions efficiently. Data alignment is required for all architectures. However, CIM-A and CIM-P classes perform computations directly on the data residing inside the memory, and hence, the robustness and performance are impacted more by data misalignment. Note that performing a data alignment cannot be handled by the host processors in in-memory computing architectures due to a far communication distance, while adding additional logic inside the memory core to handle this is also not trivial. It requires an area overhead to temporarily store operands and do the alignment with CMOS logic. For other classes, the impact of data alignment is less severe; nevertheless, data misalignment can cause a performance degradation in other classes as well.

As the primitive operations in CIM-A and CIM-P are limited, architectures in these classes face challenges in computing complex functions such as arithmetic operations with integer or floating-point numbers. As a result, a lot of primitive operations are required to realize such complex functions, if even possible. For example, a multibit addition in CIM requires multiple single-bit additions as primitive operations and communication operations between these single-bit additions [27]. On top of that, each primitive step that involves a write operation in a memristor-based CIM architecture suffers from a high latency due to its high write time. In addition, current CIM-P architectures require a high cost to implement a diverse set of arithmetic operations as their efficiency today is mainly limited to bitwise logical operations and matrix vector multiplications. Moreover, providing complex functionality using peripheral circuits in CIM-P is difficult, due to the limited available area on the memory core. Note that despite these drawbacks, the performance

can still be high when sufficient parallelism is exploited, e.g., by operating on multiple crossbars in parallel. Furthermore, data doesn't have to be transferred to the main processor, and hence, the energy and performance can be improved. In COM-N and COM-F, computations are performed by CMOS circuits that contain mature, optimized, and, if needed, dedicated functional units. However, the main bottleneck comes from the many additional data transfers through the memory hierarchy.

Available bandwidth specifies how much data can be transferred between the computational and storage units. This metric is important as it affects the amount of parallelism that can be exploited. The available bandwidth is considered similarly to the bandwidth specification of multiple levels in the memory hierarchy; hence, it includes four ranges: max (TBs), high (10GBs), medium (GBs), and low (MBs) [13]. Note that these terms are used for memory technology nowadays as the exact bandwidth values are subject to change with new or different technologies. CIM-A architectures may exploit the maximum bandwidth, as operations happen inside the memory array. CIM-P architectures have a bandwidth range from high to max, depending on the complexity of the peripheral circuitry. Note that the peripheral circuits can be complex, e.g., when large customized sense amplifiers are used. Therefore, the placement of such sense amplifiers may be limited due to area constraints. In such cases, still a relatively high bandwidth can be realized. For COM-N, the bandwidth is bounded by on-chip interconnections between the memory core and extra logic circuits; for example, in Hybrid Memory Cube [101] the bandwidth is limited by the number of TSVs and available registers. This bandwidth for TSV is considered high in comparison with COM-F, where the bandwidth is even lower due to off-chip interconnections [132].

Memory design efforts specify the required efforts needed to modify the memory (as a storage entity) to make it also realize the computing functionality. In some cases, it is very difficult (or may be even practically impossible) to modify the cells, array, periphery, and controller. CIM-A architectures require a redesign of the cell in order to make the computing feasible. Recharacterizing the cell requires a huge effort and induces a huge cost. Other classes, except for hybrid CIM-P, do not require this modification due to the usage of standard memory cells. In terms of changes in the periphery, CIM-P architecture requires complex readout circuits as the output value of two or more accessed cells may end up in multiple levels. Moreover, complex peripheral circuits (i.e., ADC, DAC) limit the scalability when they exhibit internal bottlenecks and could also dominate the area of the memory core when the memory sizes are small. Hence, CIM-P is mainly useful for larger sizes. Other classes, except for hybrid CIM-A, can utilize existing optimized readout circuits and hence do not require modifications in the periphery. In terms of memory controller, the complexity reduces from high to low for CIM-A, CIM-P, COM-N, and COM-F, respectively. CIM-A architectures require a complex controller as it is responsible for both controlling the crossbar (consisting of a large number of states, each controlling different voltage drivers) and handling data transfer (which involves the usage of buffers/registers to store temporary values). CIM-P architectures have relatively simpler controllers as the computations are constructed in a similar manner as for conventional memory (read/write) operations. The difference is that they typically have to deal with more in-memory operations. COM-N and COM-F architectures utilize the memory in a conventional way, and hence, standard memory controllers can be used.

Scalability specifies how easy or hard it is to scale up the architecture in order to maintain the parallelism level. CIM-A has a low scalability due to several reasons such as the lack/complexity of the interconnect network within the memory array it needs and the difficulty in isolating logic units to ensure parallel executing. CIM-P has a medium scalability as the limited amount of resources inside peripheral circuits makes it difficult to fit large and complex logic units; the complexity of the periphery circuits is the main bottleneck. COM-N also has a medium scalability for the same reason; even though the logic layer of memory SiP has more processing resources than

peripheral circuits, it cannot accommodate many complex logic units. COM-F has high scalability due to a mature interconnect network and large space for logic devices.

Endurance requirement specifies how many write operations can be performed before the memory of the architecture starts to fail. A memory that needs a higher number of writes will have a lower lifetime when both have technology-wise the same endurance. Three ranges can be specified for the architectures: a high endurance requirement (i.e., much higher than DRAM endurance 10^{15} [139]), a medium endurance requirement approximately equal to DRAM endurance, and a low endurance requirement much less than the DRAM endurance. CIM-A has in general a high endurance requirement due to the need for multiple write steps to perform simple Boolean functions. CIM-P has a lower endurance requirement as operations are performed during read operations [134]. Nevertheless, results have to be still written back to the memory in order to perform complex functions. As CIM-A and CIM-P architectures are typically based on emerging devices such as memristors, their endurance could be a potential issue. Similarly to CIM-P, COM-N architectures operate closely to the memory and have to write back the results to the main memory due to the absence/limited number of registers and caches. In contrast, COM-F architectures have a much lower endurance requirement as computations are performed using CMOS and the results of the operations are rarely written back to the main memory due to the usage of registers and caches.

Maturity refers not only to how feasible/reliable the memory technology is but also to how much software support exists and the development status of the architectures in the classes. As CIM-A and CIM-P are relatively new concepts and typically resistive based, lots of work still has to be done to realize these architectures from both a hardware and software point of view. Resistive memories and nonvolatile memories in general are typically under research development. For example, the limited endurance puts a constraint on the amount of computations that can be performed in resistive CIM-A architectures. Programming languages, compilers, and simulators still need to be developed in general for these architectures. In the COM-N class, several architectures have been prototyped in the industry, and therefore, they are more mature than CIM-A and CIM-P. Architectures in COM-N also have more software support as they are equipped with tool chains that allow product development on these architectures; for example, Micron's automata processor is already commercialized and is programmed in Automata Network Markup Language (ANML) [96]. COM-F architectures are today's conventional von Neumann architectures. They have the highest maturity from both a technological point and software support. With respect to the development status, CIM-A and CIM-P architectures mostly are verified using simulations, either cycle-accurate simulations [4, 12] or circuit verification simulation (i.e., HSpice). COM-N and COM-F architectures are further developed; they have been demonstrated in prototypes and commercial products [101, 104]. In general, COM architectures are more mature than CIM architectures.

Applications that run effectively on the architectures are also described in Table 3. In general, CIM architectures can be more efficient than COM architectures for certain data-intensive applications as they are less affected by the memory bottleneck. For CIM-A architectures and several CIM-P architectures (e.g., Pinatubo, CIMA, STT-CiM), there are currently limited types of operations that can be efficiently performed on these architectures; hence, a limited range of applications can be mapped on these architectures. For example, CIM-A architectures focus more on arithmetic operations such as matrix multiplication [48] and parallel addition [27]. CIM-P architectures focus on bulk bitwise applications such as database processing, graph processing, image processing, security, and biosequencing application [2, 8, 47, 81, 109, 110]. COM-N architectures are both general purpose and application specific. A limited number of COM-N architectures are considered as general-purpose computers such as FlexRAM [62] and SM [85]. Other COM-N architectures

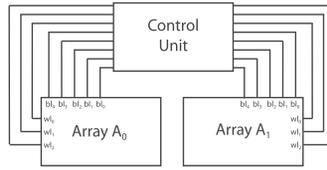


Fig. 5. Complementary resistive switch-logic crossbar array (CRS) [115].

target specific applications such as vector processing (e.g., VIRAM [59], DIVA [25], AMC [95], etc.), automata processing (D-AP [96]), and neural computation (DDN [20]). COM-F architectures are mostly designed for general-purpose applications.

3 COMPUTATION-IN-MEMORY-ARRAY (CIM-A)

The CIM-A class contains mostly resistive computing architectures that use memristive-based logic circuits [26] to perform computations and resistive RAM (RRAM) as memory technology. The resistive logic circuits may implement different design styles such as stateful logic [76], IMPLY [73], MAGIC [72], CRS-based logic [115], and so forth. These design styles can be further classified, as presented in [106]. In addition to resistive computing, computations can be performed using DRAM cells as demonstrated in [80], which will be explained later.

Few architectures have been proposed in this class; they are Complementary Resistive Switch (CRS) [115], Computation-in-Memory (CIM) [27, 43, 46], Memristive Memory Processing Unit (MPU) [52], Programmable Logic-in-Memory Computer (PLiM) [38], ReRAM-based VLIW architecture (ReVAMP) [9], and a DRAM-based Reconfigurable In-Situ Accelerator with a 3T1C cell design (DRISA-3T1C) [80]. Most of the architectures, except for REVAMP, have similar organizations. They contain a memristor crossbar (except for DRISA-3T1C) that is used for both storage and computation and a controller that applies the voltages to the memory array. Each architecture uses a different logic style and controller; for example, CRS, MPU, and PLiM use CRS-based logic, Memristive-Aid loGIC (MAGIC), and majority logic, respectively, while CIM can use any logic style. ReVAMP uses a different architecture and integrates the resistive memory in a pipelined processor in which the memory replaces both the cache and register file. It optimizes traditional pipelined processors by combining the execution, memory, and write-back in a single stage. DRISA-3T1C contains a DRAM memory array and performs NOR instructions by reading two rows simultaneously and writing the results back via the sense amplifier to another row. During the read, the capacitances of the accessed cells will discharge the bitline via a transistor when one or both cell values are high; only when both capacitance values are zero does the bitline remain. As examples, we only describe the CRS and ReVAMP architectures next in more detail; they are the latest proposed architectures that represent a basic CIM-A and hybrid CIM-A architectures, respectively. Due to page limitations, only one representative figure is used to describe each architecture.

3.1 Basic CIM-A Architecture

CRS architecture was proposed in 2014 by Siemon et al., from RWTH Aachen University [115]. It is a memristor-based architecture that exploits data-level parallelism using implication logic. The architecture consists of multiple crossbars and a control unit (as shown in Figure 5 [115]). The crossbar stores and performs logic operations using CRS cells; a CRS cell consists of two resistive switches or resistive RAMs. The control unit distributes signals to the intended addresses (wordlines and bitlines) to perform operations on the crossbars.

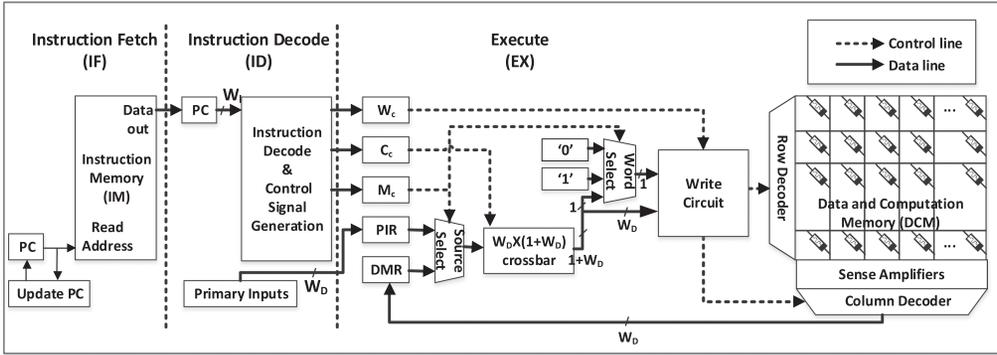


Fig. 6. ReRAM-based VLIW architecture (ReVAMP) [9].

The crossbar is controlled by a sequence of operations including write-in (WI), read-out (RO), write-back (WB), and compute (CP). Before the operations can be performed, the crossbar part used for computation is once entirely reset to a logic value 0. The WI operation writes a logic value into a memristor. The RO operation reads a logic value from a cell; the logic output value is determined by the sense amplifier. The RO operation is destructive and changes the value of the memristor to logic value 1. The task of the WB operation is to recover the destroyed value. Finally, the CP instructions are used to execute the implication logic gates [82, 115]. The data transfer between CRS cells is carried out through the control unit using RO and WB operations; in other words, the control unit reads a value of the source CRS cell and writes this value into the destination cell.

In addition to the general characteristic of CIM-A described in Tables 2 and 3, CRS has the following advantages: (1) It is less impacted by the sneak path currents due to the usage of CRS cells; the cell's resistance is always equivalent to high resistance; hence, sneak path currents are eliminated; however, variations in resistances will make such paths practically unavoidable unless a 1T2R cell is used; (2) CRS logic requires fewer cells to perform computations than Fast Boolean Logic (FBL) [133]. However, it also has the following limitations: (1) the latency of the primitive functions varies and requires readout instructions to determine the voltages that have to be applied; (2) the RO operation is destructive, and hence, a WB operation is required after each RO operation, which increases the latency and energy of computations; (3) the data transfer method is indirect as it is based on the readout and write-back scheme—as all cells have high resistance, direct copying of cells in the crossbar is not applicable; (4) the control unit imposes a high overhead as it is responsible for both controlling the crossbar (requiring a large number of states) and transferring data (which involves the usage of buffers/registers to store temporary values); and (5) the architecture requires additional compiling techniques and tools to convert conventional Boolean logic functions to implication logic. This architecture was only evaluated at circuit level using adders. Therefore, it is hard to make general conclusions on the performance and the applicability of this architecture.

3.2 Hybrid CIM-A Architecture

ReVAMP was proposed in 2017 by Bhattacharjee et al., from Nanyang Technological University [9]. It is a memristor-based architecture that exploits data parallelism using majority logic. The architecture consists of an Instruction Fetch (IF), Instruction Decode (ID), and Execute (EX) stage (as shown in Figure 6 [9]). The IF block fetches instructions from the Instruction Memory using the program counter (PC) as address and puts the resulting instruction in the Instruction

Register (IR). The ID block decodes the instruction and generates control signals, which are placed in the control registers of the EX block. The EX stage finally executes the instruction.

The IF and ID stages are similar to those of the traditional five-pipelined RISC architecture. The IF stage includes an Instruction Memory (IM) and a Program Counter (PC). The ID stage contains registers (IR and Primary Inputs) and an Instruction Decode and Control Signal Generation. The EX stage consists of several registers (i.e., Data Memory Register (DMR), Primary Input Register (PIR), Mux control (M_c) register, Control (C_c) register, Wordline (W_c) register), as well as a crossbar interconnect, wordline select multiplexer, data Source Select multiplexer, and Write circuit to control the crossbar that stores data. Once an instruction is fetched and decoded in IF and ID, respectively, the control registers in EX are filled with suitable values. These values control the multiplexers that are responsible for applying the right control signals to the crossbar. Depending on the operation, primary inputs from PIR or data retrieved from the crossbar stored in DMR can be used for the next operation. The crossbar interconnect permutes the inputs and control signals (indicated by C_c) to generate the voltages that need to be applied to the memory crossbar. The Write circuit applies these voltages to the targeted wordline address (indicated by W_c).

In addition to the general characteristic of CIM-A described in Tables 2 and 3, ReVAMP has the following advantages: (1) the data transfer may include direct (within the crossbar based on copying resistance values) and indirect (based on readout/write-back) schemes; (2) the crossbar is based on only one device per cell, resulting in a more compact architecture as compared with other architectures that make use of two devices per cell (i.e., Complementary Resistive Switch CRS [115]); (3) the architecture does not suffer destructive reads as is the case for CRS architecture [115], and hence the write energy might be less due to the absence of a write-back operation. However, it also has the following limitations: (1) the latency of majority primitive functions varies depending on the functional complexity; in addition, before any operations are applied to the cells, these cells first have to be read out in order to determine the appropriate control voltages; (2) the architecture has to deal with sneak path currents; possible solutions as mentioned above; (3) the EX stage is complex as it integrates both the control signals for memory and computations; therefore, it is not easy to pipeline this architecture, as the EX stage will consume more time than the other stages—i.e., the stages IF, ID, and EX are not balanced; (4) The architecture requires additional compiling techniques and tools to convert conventional Boolean logic functions to majority logic gates. The architecture is simulated and evaluated using EPFL benchmarks [5] and compared against PLiM [38], which is based on a resistive memory with the same logic style.

4 COMPUTATION-IN-MEMORY-PERIPHERY (CIM-P)

The CIM-P class consists of architectures that perform computations during readout operations (i.e., two or more wordlines are activated simultaneously) using special peripheral circuitry. Such operations are typically analog in nature. As there are fewer restrictions on the functionality of the cell, various memory technologies can be used in this category such as DRAM, SRAM, and nonvolatile memory technologies.

A medium number of architectures have been proposed in this class: Resistive Associative Processor (ReAP) [137], a Processing-in-Memory Architecture for Neural Network Computation in ReRAM-based Main Memory (PRIME) [21], a Convolutional Neural Network Accelerator with In Situ Analog Arithmetic (ISAAC) [111], In-Memory Accelerator for Bulk Bitwise Operations Using Commodity DRAM Technology (Ambit) [110], a Processing-in-Memory Architecture for Bulk Bitwise Operations (Pinatubo) [81], In-Memory Intelligence (IMI) [32], Compute Caches (Compute\$) [1], a DRAM-based Reconfigurable In Situ Accelerator with 1T1C design (DRISA-1T1C) [80], Computation-in-Memory Accelerator (CIMA) [26], Computing in Memory Spin-Transfer Torque

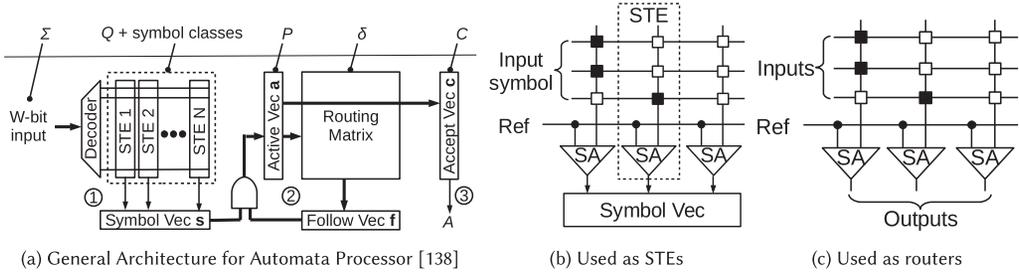


Fig. 7. Resistive RAM automata processor (R-AP) [138].

Magnetic RAM (STT-CiM) [56], Cache Automaton (S-AP) [121], Neural Cache (Neural\$) [29], RRAM Automata Processor (R-AP) [138], and Data Parallel Processor (DPP) [37].

These architectures fundamentally perform computations in the same way by activating multiple rows simultaneously in the memory and using generally specialized sense amplifiers and/or ADC converters to get the results. ReAP performs computations by implementing Content-Addressable-Memory (CAM) operations using look-up tables (LUTs). PRIME and ISAAC perform vector-matrix multiplications for neural applications. Ambit, IMI, Compute\$, DRISA-1T1C, Pinatubo, CIMA, STT-CiM, Neural\$, and DPP perform computations using customized sense amplifiers only; Ambit, IMI, and DRISA-1T1C use DRAM, Compute\$ and Neural\$ use SRAM, and the rest are based on nonvolatile memory. These architectures can only perform logical operations except for IMI, DRISA-1T1C, Neural\$, and DPP, which also perform more complex functions by having additional logic inside the peripheral circuits. S-AP and R-AP implement inner product operations in automata processors. S-AP is implemented using SRAM technology, while R-AP uses nonvolatile memory. As examples, we only describe the R-AP and DPP architectures next in more detail; they are the latest proposed architectures that represent basic CIM-P and hybrid CIM-P architectures, respectively.

4.1 Basic CIM-P Architecture

R-AP was proposed in 2018 by Yu, et al., from Delft University of Technology [138]. The architecture targets an automata processor that exploits data-level parallelism by performing computations using state machines. An automata processor contains two main components: the State Transition Elements (STEs) and the routing matrix; the STE stores the accepting states, while the routing matrix stores the state transitions as shown in Figure 7(a) [138]. The automata processor accepts one input symbol at a time, generates next active states, and decides whether a complete input string is accepted or not.

The architecture consists of STEs and a routing matrix, which are implemented using RRAM technology. Each RRAM column corresponds to an STE that stores the accepting states in RRAM cells, as shown in Figure 7(b) [138]. The input symbol is fed to all the STEs simultaneously. The sense amplifiers collect dot-product results of a vector-matrix multiplication. The output of the STE and the routing matrix are used to determine the next active states, as shown in Figure 7(c) [138]; this process is carried on until all input symbols are processed. In case the one or more final active states are part of the acceptance states, it means that the input string has been matched with the corresponding pattern of the acceptance state. Note that data transfer inside the automata processor is carried out using the routing matrix.

In addition to the general characteristic of CIM-P described in Tables 2 and 3, R-AP has the following advantages: (1) The architecture is used as a read-favored accelerator, which has a positive impact on the endurance due to infrequent use [17, 127]; only when the automata changes do the

STEs and routing matrix have to be updated; (2) automata processing can be used to perform both logical and arithmetic operations in general; data can be transferred using both direct and indirect schemes; (4) the architecture uses nonvolatile memory and hence consumes low energy and has a small footprint; (5) the automata processing techniques and tooling are quite mature; hence, it is feasible to explore many applications using automata processing. However, it also has the following limitations: (1) the modified peripheral circuitry (row drivers) might pose high overhead in the memory system; (2) the architecture requires additional compiling techniques and tools to perform conventional Boolean logic functions using automata processing; the architecture has been validated using circuit-level simulations and evaluated against S-AP [121].

4.2 Hybrid CIM-P Architecture

DPP was proposed in 2018 by Fujiki, et al. from University of Michigan [37]. DPP is a RRAM-based architecture that exploits instruction- and data-level parallelism by performing computations using a combination of RRAM-based dot-product operations and LUTs. The architecture consists of multiple RRAM tiles connected as an H-tree; each tile has multiple clusters and some logic units. Tiles and clusters form a SIMD-like processor that performs the parallel operations. The architecture is considered as a general-purpose architecture as it can perform all primitive functions such as logical, arithmetic, shift, and copy operations.

In addition to clusters, each tile has several units to support the computations including instruction buffer, Shift and Add (S+A), and router. Each cluster additionally has one or more computational units; they are Shift and Add (S+A), Sample and Hold (S+H), DAC and ADC, and a LUT and register file. While reading from the high-latency RRAM, other units are simultaneously used for processing. Therefore, the S+H is used to read data (in the form of a current) from the RRAM array and temporarily store it. Once that data is needed, it is fed to an ADC to convert the analog value to a digital value. The S+A is used to perform carry propagation in a multiple-bit addition. DAC is used to apply a digital value to the RRAM array with an appropriate control voltage. Some complex functions that cannot be realized with these units are performed using LUTs and a register file in each cluster. Data transfer can be performed by enabling two memory rows for direct copy operations or using the buffers and readout operations for indirect copy operations.

In addition to the general characteristic of CIM-P described in Tables 2 and 3, DPP has the following advantages: (1) computations include both logical operations and simple arithmetic operations (i.e., addition, multiplication); (2) data can be transferred using both direct and indirect schemes; (3) the architecture uses nonvolatile memory and hence consumes low energy and has a small footprint; and (4) this architecture is claimed to be general purpose and hence it can exploit the existing instruction set, compiling techniques and tools, and applications. However, it also has the following limitations: (1) the architecture uses nonvolatile memory as main memory, which may impact the lifetime due to limited endurance [17, 127]; (2) as the sense amplifiers are complex, a trade-off between area and bandwidth has to be made. The architecture potential was simulated and evaluated against CPU Intel Xeon E5-2697 using a subset of PARSEC benchmarks [11] and against GPU NVIDIA Titan XP using Rodinia benchmarks [18].

5 COMPUTATION-OUT-MEMORY-NEAR (COM-N)

The COM-N class consists of architectures that perform computation using additional logic units outside the memory core but inside the memory SiP. These architectures were proposed in the past and evolved through different memory technologies ranging from conventional DRAM and embedded DRAM to emerging memory technologies such as RRAM.

Many architectures have been proposed in this class: Vector Intelligent RAM (VIRAM) [59, 68, 70, 99], Active Page (A-Page) [98], Advance Intelligent Memory System (FlexRAM) [63],

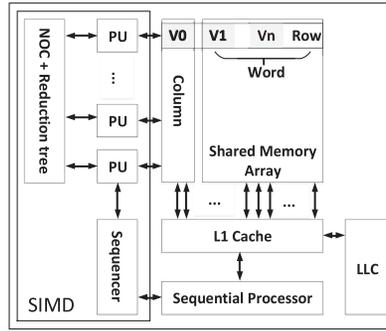
Modular Reconfigurable Smart Memories (S-Mem) [85], Data-Intensive Architecture (DIVA) [24, 25], Hybrid Memory Cube (HMC) [57, 101], Active Memory Cube (AMC) [95], Micron Automata Processor (D-AP) [96], a machine-learning supercomputer (DDN) [20], an Architecture for Accelerated Processing Near Memory (DRAM) [31], High-Bandwidth Memory (HBM) [60, 75, 84, 119], a Near Data Computing Architecture using Non-Volatile Memory (ProPRAM) [128], Resistive GP-SIMD (ReGP) [92], and HMC Instruction Large Vector Extensions (HIVE) [3].

These architectures mainly differ as a consequence of using different technologies. VIRAM, FlexRAM, SM, DIVA, and DRAMA are based on embedded DRAM technology and try to integrate processing units near the main memory; FlexRAM integrates multiple single-core processors with caches; SM a reconfigurable processor; VIRAM and DIVA a vector processor; and DRAMA a coarse-grain reconfigurable accelerator. A-Page is based on reconfigurable DRAM architecture that integrates conventional DRAM into an FPGA; it implements the reconfigurable logic near the main DRAM memory. HMC, AMC, HBM, and HIVE are based on 3D-stacked DRAM; HMC and HBM support general computing, while AMC and HIVE are optimized for VLIW and vector processing, respectively. DaDianNao, D-AP, ReGP, and ProPRAM utilize logic units that are located near the memory. DaDianNao and D-AP implement a neural network and an automata processor, respectively, with very simple logic units inside the conventional DRAM. ReGP integrates a simplified SIMD processor near nonvolatile memory. ProPRAM utilizes existing logic units near the nonvolatile memory to perform simple computations such as addition and logical operations. As an example, we only describe the HIVE and ReGP architectures next in more detail; they are the most recent architectures proposed in the COM-N class.

HIVE was proposed in 2016 by Alves et al., from Federal University of Rio Grande do Sul [3]. HIVE is a Hybrid Memory Cube (HMC) [57, 101]-based architecture that performs large vector operations inside the logic die of an HMC. The architecture consists of a host processor and an HMC module that is extended with a HIVE. The host processor, not shown in the figure, is a pipeline-like architecture with six stages; it fetches, decodes, renames, dispatches, executes, and commits a sequence of instruction. If an instruction fragment has to be executed using in-memory instructions, the processor diverts the instruction fragment to the HMC module. The HMC module executes the fragment and returns the result back to the processor.

The HMC module consists of multiple DRAM layers, logic vaults, a HIVE controller, a crossbar switch, and multiple-lane links to the host processor. The data is stored in multiple DRAM layers and retrieved by the HIVE. The HIVE controller contains a register bank, functional units, and a HIVE sequencer. The logic vaults contains a vault controller, write and read buffer, and DRAM sequencer. Once the HIVE sequencer receives an instruction, it locks the involved memory address space. If the memory has already been locked, the requested instruction returns a fail status to the processor; otherwise, a memory synchronization occurs by flushing related cache data into DRAM. The logic vaults and HIVE subsequently execute the instructions by reading data to read buffers and the register bank, performing operations using functional units, and (optional) storing into memory using write buffers. The operations in HIVE are based on vector operations that operate on 8KB of data at a time executed by the 32 logic vaults and HIVE functional units. As the amount of data is large, a DRAM sequencer and HIVE sequencer schedule these operations accordingly. The results can be collected in register banks and sent back to the host processor through the crossbar switch and links.

In addition to the general characteristic of COM-N described in Tables 2 and 3, HIVE comes with the following advantages: (1) the parallelism is high due to vector processing on 8KB of data; (2) the architecture uses HMC, which is mature, is commercialized, and has some advantages such as high performance, high bandwidth, low power, and high density [57, 101]. However, it also has the limitation that the architecture has a complex HMC module, which has a control, communication,



Resistive GP-SIMD (ReGP) [92]

Fig. 8. Examples of COM-N architectures.

and programming overhead. The architecture is simulated and evaluated using some integer (vector search and memory reset/set operations) and floating-point (vector sum, matrix stencil, and matrix multiplication) kernels against three baseline platforms; both HIVE and baseline platforms are based on the Intel Atom processor. Like HIVE, the three baseline platforms also have additional processing capacities; for the baseline platforms they are as follows: (1) HMC instructions using HMC 2.0 memory [50] (HMC+HMC), (2) 128-bit SSE instructions with DDR-3 1333 modules (SSE+DDR), and (3) 128-bit SSE instructions with HMC 2.0 (SSE+HMC).

ReGP was proposed in 2016 by Morad et al., from Technion-Israel Institute of Technology [92]. ReGP is an RRAM memory-based architecture that exploits data parallelism by attaching a SIMD-like processing unit to the resistive memory, as shown in Figure 8 [92]. The architecture consists of a sequential processor (which is a conventional processor) and its L1 and LLC cache, shared memory array, and SIMD processor. The sequential processor executes traditional code and controls the SIMD processor in a master-slave mode. The SIMD processor executes parallel instructions on the data stored in the shared memory array.

The SIMD processor contains multiple processing units (PUs), a sequencer, and a Network on Chip (NoC) with reduction tree. Each PU contains registers, a single-bit full-adder, and a function generator to perform arithmetic and logical operations. The sequencer receives instructions from the sequential processor and assigns them to PUs. The PUs load data from the shared memory array and perform the requested operations. If required, the NoC and reduction trees are used to perform more complex functions.

In addition to the general characteristic of COM-N described in Tables 2 and 3, ReGP comes with the following advantages: (1) the parallelism is high due to multiple parallel processing units; (2) the architecture uses nonvolatile memory and hence consumes a low amount of energy and has a small footprint; (3) the architecture can reuse compilers, programming languages, and tools from SIMD architectures. However, it also has the limitation that the operations within the processing units are simple; complex functions such as floating-point operations can cause a high overhead. The architecture is simulated and evaluated against CMOS GP-SIMD [91] using a benchmark for dense matrix multiplications [90].

6 DISCUSSION

This section aims to first evaluate the completeness of the proposed classification. Thereafter, we compare it with existing work in the field. Finally, we discuss the limitations of this work and propose directions for future work.

6.1 Completeness

The proposed classification presented in Figure 3 is complete and comprehensive. These points can as follows be proven: (1) theoretically, due to the exploration of all the possible classes derived from the classification metrics, and (2) practically, by mapping all existing memory-centric architectures on the classification.

Theoretically, the classification contains four main classes derived from the “computation location” (first metric); both inside and outside, approximately close or distant from the memory core. Moreover, the second metric consists of both charge-based and non-charge-based memories. Finally, the parallelism metric ranges from instruction to data and task levels. Each metric is in itself complete, and therefore, the entire classification is complete. The classification not only contains the existing solutions but also highlights the potential future solutions that can be further explored (e.g., classes in blue spaces in Figure 3). Note that hybrid architectures are also covered in this classification. For example, a conventional architecture (COM-F) with accelerator in CIM-P class (e.g., ReAP, ISAAC, CIMA) is considered a hybrid architecture, i.e., a COM-F/CIM-P hybrid.

Practically, it contains an overview most of the existing computer architectures and places them in perspective. In addition, the classification can be used to illustrate the past and future trends (see Figure 4). Moreover, it clearly depicts a shift from conventional processor-centric architectures toward memory-centric architectures based on emerging technologies (3D stacking, RRAM, etc.).

6.2 Related Work

Comparison with traditional/processor-centric architecture classifications: Conventional classifications like Flynn’s [35], Skillicorn’s [117], and Shami-Hemani’s [112] classification are quite comprehensive and were considered complete at the time they were published. However, these classifications focus on processor-centric architectures and hence can only be used to classify conventional architectures (i.e., architectures in COM-F class). Aside from the above-mentioned classifications, some publications on the COM-N class have presented intensive architectural reviews [23, 113, 116]. However, they have a restricted focus on near-memory-processing architectures based on 2D, 2.5D, and 3D-stacked DRAM. Signh’s classification [116] is the most recent work that provides a review of near-memory computing architectures, i.e., COM-N architectures. It classifies architectures mainly based on the memory hierarchy and processing type (e.g., programmable unit, fixed functional unit, and reconfigurable unit). Moreover, it evaluates the architectures based on multiple characteristics of memory, processing, evaluation tools, interoperability, and application domains. However, the classification is not easy to use as the metrics are not systematic. Furthermore, it is not clear if the classification is complete and if it covers all ranges of near-computing architectures. Last but not least, in comparison with the aforementioned classifications, our proposed classification goes one step further to cover both conventional and emerging architectures by having the additional classes CIM-A and CIM-P. Moreover, the proposed classification is so broad that several of its classes are not explored yet. New architectures in these unexplored areas can be easily added to the classification. In addition, our proposed classification uses three selective metrics, which create distinctive and easy-to-use terminologies, classes, and subclasses.

Comparison with recent/emerging architecture classifications: Recent surveys and classifications for emerging architectures have been proposed by Mittal [89] and Reuben et al. [106]. Mittal’s classification only tries to link architectures with their applications. Specifically, the classification discusses three unconventional architectures: processing-in-memory, machine learning, and neural-network-based architectures using RRAM. They mostly focus on applications containing dot-product operations in the RRAM crossbar. This classification is not complete, as RRAM in

particular and emerging memory technology in general can also be used to implement other functions such as bitwise logic operations [81, 134], arithmetic operations using implication logic [72, 114], Boolean logic [118, 133], and the like. Reuben's classification classifies existing resistive logic design methods into three classes: in-memory, near-memory, and out-of-memory computing. The near-memory class has three subclasses without identifiers (e.g., they are based on how data moves out of the memory array; this includes data movements (1) between consecutive logic levels, (2) for computing each Sum-of-Products, and (3) for computing each logic gate). This classification, however, tries to redefine the terminologies without defining clear generic metrics for each class. Instead, each class uses different criteria to distinguish between their subclasses. Therefore, it is not a systematic and comprehensive classification, which makes it difficult to use in identifying and exploring architectures. Moreover, it is difficult to judge if the classification is complete. Furthermore, the classification focuses only on resistive memories, while other emerging memory technologies are also promising. Overall, both Mittal and Reuben et al. classifications are not complete and comprehensive enough to classify all architectures.

6.3 Future Directions and Challenges

Memory-entric computing is seen as one of the promising solutions to alleviate (even if partially) the memory bottleneck. Not only the communication between the processing core and the main memory but also the energy consumption will be significantly reduced; the data communication on its own is extremely energy consuming. Implementing CIM based on DRAM or emerging memristive devices seems to be more realistic than using on-chip SRAMs. Although SRAM technology is more CMOS compatible when it comes to manufacturing, the cost per bit for such technology is much higher than that of other memory technologies. Hence, the overall cost of large-capacity SRAMs (which is needed for CIM) is by far much higher; for the same capacity, SRAM consumes much more area/power compared to DRAM and nonvolatile memories. In addition, the two main directions that are currently explored are CIM-A and CIM-P, in which CIM-P is more feasible than CIM-A due to the complex underlining memory technology. CIM-P requires less effort and modification in the memory core (mainly in the periphery). Moreover, CIM architectures do not make conventional architectures obsolete; in fact, multicore architectures with caches are relevant for applications with high data locality, while CIM architectures can only be used efficiently for certain specific applications [44]. Furthermore, building appropriate simulators and tools for CIM architectures based on technology calibrated models will enable a real estimation of the potential of such architectures [6, 39, 77, 135].

It is worth mentioning that the focus of this article is to propose a unified terminology and classification instead of presenting a survey. In our future work, we will present a survey that intensively discusses all architectures.

7 CONCLUSION

In this article, we have proposed a classification using three metrics: computational location, memory technology, and level of parallelism. We have used the most important metric, i.e., computational location, to describe and evaluate the four main classes (and the selected architectures therein). The work shows that architectures are required to be not only memory bottleneck free but also energy and area efficient. In order to accomplish that, the architectures must be implemented with the right technologies. The relationship and dependency between the architecture and technologies becomes stronger for memory-centric computing architectures. This work also showed that new architectures typically emerge after new technology developments (e.g., introduction of 3D stacking and RRAM). Our classification unifies the prior work and aims to provide a comprehensive and unique terminology for memory-centric computing architectures. Finally,

the classification not only presents an overview of existing architectures but also predicts the potential of future architecture variants, including hybrid architectures that may combine different strengths of the different classes.

REFERENCES

- [1] Shaizeen Aga, Supreet Jeloka, Arun Subramaniyan, Satish Narayanasamy, David Blaauw, and Reetuparna Das. 2017. Compute caches. In *2017 IEEE International Symposium on High Performance Computer Architecture (HPCA '17)*. IEEE, 481–492.
- [2] Junwhan Ahn, Sungpack Hong, Sungjoo Yoo, Onur Mutlu, and Kiyong Choi. 2016. A scalable processing-in-memory accelerator for parallel graph processing. *ACM SIGARCH Computer Architecture News* 43, 3 (2016), 105–117.
- [3] Marco A. Z. Alves, Matthias Diener, Paulo C. Santos, and Luigi Carro. 2016. Large vector extensions inside the HMC. In *Design, Automation & Test in Europe Conference & Exhibition (DATE '16)*. IEEE, 1249–1254.
- [4] Marco Antonio Zanata Alves, Carlos Villavieja, Matthias Diener, Francis Birck Moreira, and Philippe Olivier Alexandre Navaux. 2015. SiNUCA: A validated micro-architecture simulator. In *Proceeding of International Conference on High Performance Computing and Communications (HPCC), International Symposium on Cyberspace Safety and Security (CSS), and International Conference on Embedded Software and Systems (ICSS)*. 605–610.
- [5] Luca Amarú, Pierre-Emmanuel Gaillardon, and Giovanni De Micheli. 2015. The EPFL combinational benchmark suite. In *Proceedings of the 24th International Workshop on Logic & Synthesis (IWLS '15)*.
- [6] Ali BanaGozar, Kanishkan Vadivel, Sander Stuijk, Henk Corporaal, Stephan Wong, Muath Abu Lebdeh, Jintao Yu, and Said Hamdioui. 2019. CIM-SIM: Computation in memory SIMUlator. In *International Workshop on Software and Compilers for Embedded Systems*. ACM, 1–4.
- [7] John Barth, Don Plass, Erik Nelson, Charlie Hwang, Gregory Fredeman, Michael Sperling, Abraham Mathews, Toshiaki Kirihata, William R. Reohr, Kavita Nair, and Nianzheng Cao. 2010. A 45nm SOI embedded DRAM macro for the POWER™ processor 32 MByte on-chip L3 cache. *IEEE Journal of Solid-State Circuits* 46, 1 (2010), 64–75.
- [8] Gary Benson, Yozen Hernandez, and Joshua Loving. 2013. A bit-parallel, general integer-scoring sequence alignment algorithm. In *Annual Symposium on Combinatorial Pattern Matching*. Springer, 50–61.
- [9] Debjyoti Bhattacharjee, Rajeswari Devadoss, and Anupam Chattopadhyay. 2017. ReVAMP: ReRAM based VLIW architecture for in-memory computing. In *2017 Design, Automation & Test in Europe Conference & Exhibition (DATE '17)*. IEEE, 782–787.
- [10] Sabpreet Bhatti, Rachid Sbiaa, Atsufumi Hirohata, Hideo Ohno, Shunsuke Fukami, and S. N. Piramanayagam. 2017. Spintronics based random access memory: A review. *Materials Today* 20, 9 (2017), 530–548.
- [11] Christian Bienia, Sanjeev Kumar, Jaswinder Pal Singh, and Kai Li. 2008. The PARSEC benchmark suite: Characterization and architectural implications. In *Proceedings of the 17th International Conference on Parallel Architectures and Compilation Techniques*. ACM, 72–81.
- [12] Nathan Binkert, Bradford Beckmann, Gabriel Black, Steven K. Reinhardt, Ali Saidi, Arkaprava Basu, Joel Hestness, Derek R. Hower, Tushar Krishna, Somayeh Sardashti, et al. 2011. The gem5 simulator. *ACM SIGARCH Computer Architecture News* 39, 2 (2011), 1–7.
- [13] Evgeny Bolotin, David Nellans, Oreste Villa, Mike O'Connor, Alex Ramirez, and Stephen W. Keckler. 2015. Designing efficient heterogeneous memory architectures. *IEEE Micro* 35, 4 (2015), 60–68.
- [14] Julien Borghetti, Gregory S. Snider, Philip J. Kuekes, J. Joshua Yang, Duncan R. Stewart, and R. Stanley Williams. 2010. Memristive switches enable stateful logic operations via material implication. *Nature* 464, 7290 (2010), 873–876.
- [15] S. Borkar. 1999. Design challenges of technology scaling. *IEEE Micro* 19, 4 (July 1999), 23–29. DOI : <https://doi.org/10.1109/40.782564>
- [16] Rafmag Cabrera, Emmanuelle Merced, and Nelson Sepúlveda. 2013. A micro-electro-mechanical memory based on the structural phase transition of VO2. *Physica Status Solidi (a)* 210, 9 (2013), 1704–1711.
- [17] Meng-Fan Chang, Ching-Hao Chuang, Min-Ping Chen, Lai-Fu Chen, Hiroyuki Yamauchi, Pi-Feng Chiu, and Shyh-Shyuan Sheu. 2012. Endurance-aware circuit designs of nonvolatile logic and nonvolatile SRAM using resistive memory (memristor) device. In *2012 17th Asia and South Pacific Design Automation Conference (ASP-DAC '12)*. IEEE, 329–334.
- [18] Shuai Che, Michael Boyer, Jiayuan Meng, David Tarjan, Jeremy W. Sheaffer, Sang-Ha Lee, and Kevin Skadron. 2009. Rodinia: A benchmark suite for heterogeneous computing. In *IEEE International Symposium on Workload Characterization, 2009 (IISWC '09)*. IEEE, 44–54.
- [19] E. Chen, D. Apalkov, Z. Diao, A. Driskill-Smith, D. Druist, D. Lottis, V. Nikitin, X. Tang, S. Watts, S. Wang, et al. 2010. Advances and future prospects of spin-transfer torque random access memory. *IEEE Transactions on Magnetics* 46, 6 (2010), 1873–1878.

- [20] Yunji Chen, Tao Luo, Shaoli Liu, Shijin Zhang, Liqiang He, Jia Wang, Ling Li, Tianshi Chen, Zhiwei Xu, Ninghui Sun, et al. 2014. Dadiannao: A machine-learning supercomputer. In *IEEE/ACM International Symposium on Microarchitecture*. IEEE Computer Society, 609–622.
- [21] Ping Chi, Shuangchen Li, Cong Xu, Tao Zhang, Jishen Zhao, Yongpan Liu, Yu Wang, and Yuan Xie. 2016. PRIME: A novel processing-in-memory architecture for neural network computation in ReRAM-based main memory. In *ACM SIGARCH Computer Architecture News*, Vol. 44. IEEE Press, 27–39.
- [22] Gianni Conte, Stefano Tommesani, and Francesco Zanichelli. 2000. The long and winding road to high-performance image processing with MMX/SSE. In *Proceedings of the 5th IEEE International Workshop on Computer Architectures for Machine Perception, 2000*. IEEE, 302–310.
- [23] Joao Paulo C. de Lima, Paulo Cesar Santos, Marco A. Z. Alves, Antonio C. S. Beck, and Luigi Carro. 2018. Design space exploration for PIM architectures in 3D-stacked memories. In *Computer Frontier*. ACM, 295–308.
- [24] Jaffrey Draper, J. Tim Barrett, Jeff Sonddeen, Sumit Mediratta, Chang Woo Kang, Ihn Kim, and Gokhan Daglikoca. 2005. A prototype processing-in-memory (PIM) chip for the data-intensive architecture (DIVA) system. *Journal of VLSI Signal Processing Systems for Signal, Image and Video Technology* 40, 1 (2005), 73–84.
- [25] Jeff Draper, Jacqueline Chame, Mary Hall, Craig Steele, Tim Barrett, Jeff LaCoss, John Granacki, Jaewook Shin, Chun Chen, Chang Woo Kang, et al. 2002. The architecture of the DIVA processing-in-memory chip. In *Proceedings of the 16th International Conference on Supercomputing*. ACM, 14–25.
- [26] H. A. Du Nguyen, Jintao Yu, Lei Xie, Mottaqiallah Taouil, Said Hamdioui, and Dietmar Fey. 2017. Memristive devices for computing: Beyond CMOS and beyond von Neumann. In *2017 IFIP/IEEE International Conference on Very Large Scale Integration (VLSI-SoC'17)*. IEEE, 1–10.
- [27] Hoang Anh Du Nguyen, Lei Xie, Mottaqiallah Taouil, Razvan Nane, Said Hamdioui, and Koen Bertels. 2017. On the implementation of computation-in-memory parallel adder. *IEEE Transactions on Very Large Scale Integration (VLSI) Systems* 25, 8 (2017), 2206–2219.
- [28] P. Dudek and S. J. Carey. 2006. General-purpose 128/spl times/128 SIMD processor array with integrated image sensor. *Electronics Letters* 42, 12 (2006), 678–679.
- [29] Charles Eckert, Xiaowei Wang, Jingcheng Wang, Arun Subramaniyan, Ravi Iyer, Dennis Sylvester, David Blaauw, and Reetuparna Das. 2018. Neural cache: Bit-serial in-cache acceleration of deep neural networks. *arXiv preprint arXiv:1805.03718* (2018).
- [30] Susan J. Eggers, Joel S. Emer, Henry M. Levy, Jack L. Lo, Rebecca L. Stamm, and Dean M. Tullsen. 1997. Simultaneous multithreading: A platform for next-generation processors. *IEEE Micro* 17, 5 (1997), 12–19.
- [31] Amin Farmahini-Farahani, Jung Ho Ahn, Katherine Morrow, and Nam Sung Kim. 2015. DRAMA: An architecture for accelerated processing near memory. *IEEE Computer Architecture Letters* 14, 1 (2015), 26–29.
- [32] Tim Finkbeiner, Glen Hush, Troy Larsen, Perry Lea, John Leidel, and Troy Manning. 2017. In-memory intelligence. *IEEE Micro* 37, 4 (2017), 30–38.
- [33] Nadeem Firasta, Mark Buxton, Paula Jinbo, Kaveh Nasri, and Shihjong Kuo. 2008. Intel AVX: New frontiers in performance improvements and energy efficiency. *Intel White Paper* 19 (2008), 20.
- [34] Randall James Fisher. 2003. General-purpose SIMD within a register: Parallel processing on consumer microprocessors. Doctoral Dissertation.
- [35] M. Flynn. 1966. Very high-speed computing systems. *Proceedings of the IEEE* 54, 12 (Dec. 1966), 1901–1909. DOI : <https://doi.org/10.1109/PROC.1966.5273>
- [36] G. D. Fuchs, N. C. Emley, I. N. Krivorotov, P. M. Braganca, E. M. Ryan, S. I. Kiselev, J. C. Sankey, D. C. Ralph, R. A. Buhrman, and J. A. Katine. 2004. Spin-transfer effects in nanoscale magnetic tunnel junctions. *Applied Physics Letters* 85, 7 (2004), 1205–1207.
- [37] Daichi Fujiki, Scott Mahlke, and Reetuparna Das. 2018. In-memory data parallel processor. In *Proceedings of the 23rd International Conference on Architectural Support for Programming Languages and Operating Systems*. ACM, 1–14.
- [38] Pierre-Emmanuel Gaillardon, Luca Amar, Anne Siemon, Eike Linn, Rainer Waser, Anupam Chattopadhyay, and Giovanni De Micheli. 2016. The programmable logic-in-memory (PLiM) computer. In *2016 Design, Automation & Test in Europe Conference & Exhibition (DATE'16)*. IEEE, 427–432.
- [39] Mingyu Gao, Grant Ayers, and Christos Kozyrakis. 2015. Practical near-data processing for in-memory analytics frameworks. In *2015 International Conference on Parallel Architecture and Compilation (PACT'15)*. IEEE, 113–124.
- [40] Simcha Gochman, Avi Mendelson, Alon Naveh, and Efraim Rotem. 2006. Introduction to Intel core duo processor architecture. *Intel Technology Journal* 10, 2 (2006), 89–97.
- [41] Jonathan E. Green, Jang Wook Choi, Akram Boukai, Yuri Bunimovich, Ezekiel Johnston-Halperin, Erica DeIonno, Yi Luo, Bonnie A. Sheriff, Ke Xu, Young Shik Shin, et al. 2007. A 160-kilobit molecular electronic memory patterned at 10 11 bits per square centimetre. *Nature* 445, 7126 (2007), 414.
- [42] Beat Halg. 1990. On a micro-electro-mechanical nonvolatile memory cell. *IEEE Transactions on Electron Devices* 37, 10 (1990), 2230–2236.

- [43] Said Hamdioui, Koenraad Laurent Maria Bertels, and Mottaqiallah Taouil. 2017. Computing Device for Big Data Applications Using Memristors. US Patent 9,824,753.
- [44] Said Hamdioui, Hoang Anh Du Nguyen, Mottaqiallah Taouil, Abu Sebastian, Manuel Le Gallo, Sandeep Pande, Siebren Schaafsma, Francky Catthoor, Shidhartha Das, Fernando G. Redondo, et al. 2019. Applications of computation-in-memory architectures based on memristive devices. In *2019 Design, Automation & Test in Europe Conference & Exhibition (DATE'19)*. IEEE, 486–491.
- [45] Said Hamdioui, Shahar Kvatinsky, Gert Cauwenberghs, Lei Xie, Nimrod Wald, Siddharth Joshi, Hesham Mostafa Elsayed, Henk Corporaal, and Koen Bertels. 2017. Memristor for computing: Myth or reality? In *Proceedings of the Conference on Design, Automation & Test in Europe*. European Design and Automation Association, 722–731.
- [46] Said Hamdioui, Lei Xie, Hoang Anh Du Nguyen, Mottaqiallah Taouil, Koen Bertels, Henk Corporaal, Hailong Jiao, Francky Catthoor, Dirk Wouters, Linn Eike, et al. 2015. Memristor based computation-in-memory architecture for data-intensive applications. In *Proceedings of the 2015 Design, Automation & Test in Europe Conference & Exhibition*. EDA Consortium, 1718–1725.
- [47] JongWook Han, Choon-Sik Park, Dae-Hyun Ryu, and Eun-Soo Kim. 1999. Optical image encryption based on XOR operations. *Optical Engineering* 38, 1 (1999), 47–55.
- [48] Adib Haron, Jintao Yu, Razvan Nane, Mottaqiallah Taouil, Said Hamdioui, and Koen Bertels. 2016. Parallel matrix multiplication on memristor-based computation-in-memory architecture. In *2016 International Conference on High Performance Computing & Simulation (HPCS'16)*. IEEE, 759–766.
- [49] John L. Hennessy and David A. Patterson. 2011. *Computer Architecture: A Quantitative Approach*. Elsevier.
- [50] HMC. 2018. Hybrid Memory Cube Specification 2.1. Retrieved from <http://hybridmemorycube.org/>.
- [51] M. Hosomi, H. Yamagishi, T. Yamamoto, K. Bessho, Y. Higo, K. Yamane, H. Yamada, M. Shoji, H. Hachino, C. Fukumoto, et al. 2005. A novel nonvolatile memory with spin torque transfer magnetization switching: Spin-RAM. In *IEEE International Electron Devices Meeting, 2005. IEDM Technical Digest*. IEEE, 459–462.
- [52] Rotem Ben Hur and Shahar Kvatinsky. 2016. Memristive memory processing unit (MPU) controller for in-memory processing. In *IEEE International Conference on the Science of Electrical Engineering (ICSEE'16)*. IEEE, 1–5.
- [53] IBM. 2014. Power 4 - The First Multi-Core, 1GHz Processor.
- [54] ITRS. 2010. ITRS ERD Report. Retrieved from <http://www.itrs.net>.
- [55] Subramanian S. Iyer and Howard L. Kalter. 1999. Embedded DRAM technology: Opportunities and challenges. *IEEE Spectrum* 36, 4 (1999), 56–64.
- [56] Shubham Jain, Ashish Ranjan, Kaushik Roy, and Anand Raghunathan. 2017. Computing in memory with spin-transfer torque magnetic RAM. *arXiv preprint arXiv:1703.02118* (2017).
- [57] Joe Jeddelloh and Brent Keeth. 2012. Hybrid memory cube new DRAM architecture increases density and performance. In *2012 Symposium on VLSI Technology (VLSIT'12)*. IEEE, 87–88.
- [58] Zhang Jianwu, Zhao Danying, et al. 2008. Survey on microprocessor architecture and development trends. In *11th IEEE International Conference on Communication Technology, 2008 (ICCT'08)*. IEEE, 297–300.
- [59] David Judd, Katherine Yelick, Christoforos Kozyrakis, David Martin, and David Patterson. 2001. Exploiting on-chip memory bandwidth in the VIRAM compiler. In *Intelligent Memory Systems*. Springer, 122–134.
- [60] Hongshin Jun, Jinhee Cho, Kangseol Lee, Ho-Young Son, Kwiwook Kim, Hanho Jin, and Keith Kim. 2017. HBM (high bandwidth memory) DRAM technology and architecture. In *2017 IEEE International Memory Workshop (IMW'17)*. IEEE, 1–4.
- [61] Ron Kalla, Balam Sinharoy, William J. Starke, and Michael Floyd. 2010. Power7: IBM's next-generation server processor. *IEEE Micro* 30, 2 (2010), 7–15.
- [62] Yi Kang, Wei Huang, Seung-Moon Yoo, D. Keen, Zhenzhou Ge, V. Lam, P. Pattnaik, and J. Torrellas. [n.d.]. FlexRAM: Toward an advanced intelligent memory system. In *2012 IEEE 30th International Conference on Computer Design (ICCD'12)*. 5–14. DOI : <https://doi.org/10.1109/ICCD.2012.6378608>
- [63] Yi Kang, Wei Huang, Seung-Moon Yoo, Diana Keen, Zhenzhou Ge, Vinh Lam, Pratap Pattnaik, and Josep Torrellas. 2012. FlexRAM: Toward an advanced intelligent memory system. In *2012 IEEE 30th International Conference on Computer Design (ICCD'12)*. IEEE, 5–14.
- [64] Doris Keitel-Schulz and Norbert Wehn. 1998. Issues in embedded DRAM development and applications. In *Proceedings of the 11th International Symposium on System Synthesis*. IEEE Computer Society, 23–31.
- [65] Doris Keitel-Schulz and Norbert Wehn. 2001. Embedded DRAM development: Technology, physical design, and application issues. *IEEE Design & Test of Computers* 18, 3 (2001), 7–15.
- [66] Kyosun Kim, Sangho Shin, and Sung-Mo Kang. 2011. Stateful logic pipeline architecture. In *2011 IEEE International Symposium of Circuits and Systems (ISCAS'11)*. IEEE, 2497–2500.
- [67] David Kirk et al. 2007. NVIDIA CUDA software and GPU parallel computing architecture. In *ISMM*, Vol. 7. 103–104.
- [68] Christoforos Kozyrakis. 2002. *Scalable Vector Media-Processors for Embedded Systems*. Technical Report. California University Berkeley Computer Science Division.

- [69] Christoforos Kozyrakis and David Patterson. 2002. Vector vs. superscalar and VLIW architectures for embedded multimedia benchmarks. In *Proceedings of the 35th Annual ACM/IEEE International Symposium on Microarchitecture*. IEEE Computer Society Press, 283–293.
- [70] Christoforos E. Kozyrakis, Stylianos Perissakis, David Patterson, Thomas Anderson, Krste Asanovic, Neal Cardwell, Richard Fromm, Jason Golbus, Benjamin Gribstad, Kimberly Keeton, et al. 1997. Scalable processors in the billion-transistor era: IRAM. *Computer* 30, 9 (1997), 75–78.
- [71] Nasser Kurd, Muntaquim Chowdhury, Edward Burton, Thomas P. Thomas, Christopher Mozak, Brent Boswell, Praveen Mosalikanti, Mark Neidengard, Anant Deval, Ashish Khanna, et al. 2014. Haswell: A family of IA 22nm processors. *IEEE Journal of Solid-State Circuits* 50, 1 (2014), 49–58.
- [72] Shahar Kvatinsky, Dmitry Belousov, Slavik Liman, Guy Satat, Nimrod Wald, Eby G. Friedman, Avinoam Kolodny, and Uri C. Weiser. 2014. MAGIC—Memristor-aided logic. *IEEE Transactions on Circuits and Systems II: Express Briefs* 61, 11 (2014), 895–899.
- [73] Shahar Kvatinsky, Guy Satat, Nimrod Wald, Eby G. Friedman, Avinoam Kolodny, and Uri C. Weiser. 2014. Memristor-based material implication (IMPLY) logic: Design principles and methodologies. *IEEE Transactions on Very Large Scale Integration (VLSI) Systems* 22, 10 (2014), 2054–2066.
- [74] Benjamin C. Lee, Engin Ipek, Onur Mutlu, and Doug Burger. 2010. Phase change memory architecture and the quest for scalability. *Communications of the ACM* 53, 7 (2010), 99–106.
- [75] Jong Chern Lee, Jihwan Kim, Kyung Whan Kim, Young Jun Ku, Dae Suk Kim, Chunseok Jeong, Tae Sik Yun, Hongjung Kim, Ho Sung Cho, Yeon Ok Kim, et al. 2016. 18.3 A 1.2 V 64Gb 8-channel 256GB/s HBM DRAM with peripheral-base-die architecture and small-swing technique on heavy load interface. In *2016 IEEE International Solid-State Circuits Conference (ISSCC'16)*. IEEE, 318–319.
- [76] Eero Lehtonen, Jussi H. Poikonen, and Mika Laiho. 2014. Memristive stateful logic. In *Memristor Networks*. Springer, 603–623.
- [77] John D. Leidel and Yong Chen. 2016. Hmc-sim-2.0: A simulation platform for exploring custom memory cube operations. In *2016 IEEE International Parallel and Distributed Processing Symposium Workshops (IPDPSW'16)*. IEEE, 621–630.
- [78] Chao Li, Wendy Fan, Bo Lei, Daihua Zhang, Song Han, Tao Tang, Xiaolei Liu, Zuqin Liu, Sylvia Asano, Meyya Meyyappan, et al. 2004. Multilevel memory based on molecular devices. *Applied Physics Letters* 84, 11 (2004), 1949–1951.
- [79] Chao Li, Daihua Zhang, Xiaolei Liu, Song Han, Tao Tang, Chongwu Zhou, Wendy Fan, Jessica Koehne, Jie Han, Meyya Meyyappan, et al. 2003. Fabrication approach for molecular memory arrays. *Applied Physics Letters* 82, 4 (2003), 645–647.
- [80] Shuangchen Li, Dimin Niu, Krishna T. Malladi, Hongzhong Zheng, Bob Brennan, and Yuan Xie. 2017. DRISA: A DRAM -based reconfigurable in-situ accelerator. In *Proceedings of the 50th Annual IEEE/ACM International Symposium on Microarchitecture*. ACM, 288–301.
- [81] Shuangchen Li, Cong Xu, Qiaosha Zou, Jishen Zhao, Yu Lu, and Yuan Xie. 2016. Pinatubo: A processing-in-memory architecture for bulk bitwise operations in emerging non-volatile memories. In *Proceeding of ACM/EDAC/IEEE Design Automation Conference (DAC)*. IEEE, 173–178.
- [82] E. Linn, R. Rosezin, S. Tappertzhofen, R. Waser, et al. 2012. Beyond von Neumann—logic operations in passive crossbar arrays alongside memory operations. *Nanotechnology* 23, 30 (2012), 305205.
- [83] Andrea Lodi, Mario Toma, Fabio Campi, Andrea Cappelli, Roberto Canegallo, and Roberto Guerrieri. 2003. A VLIW processor with reconfigurable instruction set for embedded applications. *IEEE Journal of Solid-state Circuits* 38, 11 (2003), 1876–1886.
- [84] Joe Macri. 2015. AMD’s next generation GPU and high bandwidth memory architecture: FURY. In *2015 IEEE Hot Chips 27 Symposium (HCS'15)*. IEEE, 1–26.
- [85] Ken Mai, Tim Paaske, Nuwan Jayasena, Ron Ho, William J. Dally, and Mark Horowitz. 2000. Smart memories: A modular reconfigurable architecture. *ACM SIGARCH Computer Architecture News* 28, 2 (2000), 161–171.
- [86] Ariel Maislos et al. 2011. A new era in embedded Flash memory. In *Flash Memory Summit*.
- [87] Jack A. Mandelman, Robert H. Dennard, Gary B. Bronner, John K. DeBrosse, Rama Divakaruni, Yujun Li, and Carl J. Radens. 2002. Challenges and future directions for the scaling of dynamic random-access memory (DRAM). *IBM Journal of Research and Development* 46, 2.3 (2002), 187–212.
- [88] Pedro Marcuello, Antonio González, and Jordi Tubella. 1998. Speculative multithreaded processors. In *Proceedings of the 12th International Conference on Supercomputing*. ACM, 77–84.
- [89] Sparsh Mittal. 2018. A survey of ReRAM-based architectures for processing-in-memory and neural networks. *Machine Learning and Knowledge Extraction* 1, 1 (2018), 75–114. DOI : <https://doi.org/10.3390/make1010005>
- [90] Amir Morad, Leonid Yavits, and Ran Ginosar. 2014. Efficient dense and sparse Matrix multiplication on GP-SIMD. In *2014 24th International Workshop on Power and Timing Modeling, Optimization and Simulation (PATMOS'14)*. IEEE, 1–8.

- [91] Amir Morad, Leonid Yavits, and Ran Ginosar. 2015. GP-SIMD processing-in-memory. *ACM Transactions on Architecture and Code Optimization (TACO)* 11, 4 (2015), 53.
- [92] Amir Morad, Leonid Yavits, Shahar Kvatinsky, and Ran Ginosar. 2016. Resistive GP-SIMD processing-in-memory. *ACM Transactions on Architecture and Code Optimization (TACO)* 12, 4 (2016), 57.
- [93] Onur Mutlu. 2013. Memory scaling: A systems architecture perspective. In *2013 5th IEEE International Memory Workshop (IMW'13)*. IEEE, 21–25.
- [94] Ravi Nair. 2015. Evolution of memory architecture. *Proceedings of the IEEE* 103, 8 (2015), 1331–1345.
- [95] Ravi Nair, Samuel F. Antao, Carlo Bertolli, Pradip Bose, Jose R. Brunheroto, Tong Chen, C.-Y. Cher, Carlos H. A. Costa, Jun Doi, Constantinos Evangelinos, et al. 2015. Active memory cube: A processing-in-memory architecture for exascale systems. *IBM Journal of Research and Development* 59, 2/3 (2015), 17–1.
- [96] H. Noyes et al. 2014. Micron’s automata processor architecture: Reconfigurable and massively parallel automata processing. In *Proceedings of 5th International Symposium on Highly-Efficient Accelerators and Reconfigurable Technologies*.
- [97] NVIDIA. 2012. Tesla K20X GPU Accelerator Board Specification.
- [98] Mark Oskin, Frederic T. Chong, and Timothy Sherwood. 1998. *Active Pages: A Computation Model for Intelligent Memory*. Vol. 26. IEEE Computer Society.
- [99] David Patterson, Thomas Anderson, Neal Cardwell, Richard Fromm, Kimberly Keeton, Christoforos Kozyrakis, Randi Thomas, and Katherine Yelick. 1997. A case for intelligent RAM. *IEEE Micro* 17, 2 (1997), 34–44.
- [100] David A. Patterson. 2006. Future of computer architecture. In *Berkeley EECS Annual Research Symposium (BEARS), College of Engineering*, UC Berkeley, US.
- [101] J. Thomas Pawlowski. 2011. Hybrid memory cube (HMC). In *2011 IEEE Hot Chips 23 Symposium (HCS'11)*. IEEE, 1–24.
- [102] Alex Peleg and Uri Weiser. 1996. MMX technology extension to the Intel architecture. *IEEE Micro* 16, 4 (1996), 42–50.
- [103] M. Radosavljević, M. Freitag, K. V. Thadani, and A. T. Johnson. 2002. Nonvolatile molecular memory elements based on ambipolar nanotube field effect transistors. *Nano Letters* 2, 7 (2002), 761–764.
- [104] R. M. Ramanathan. 2006. Intel® multi-core processors. In *Making the Move to Quad-Core and Beyond*.
- [105] Simone Raoux, Feng Xiong, Matthias Wuttig, and Eric Pop. 2014. Phase change materials and phase change memory. *MRS Bulletin* 39, 8 (2014), 703–710.
- [106] John Reuben, Rotem Ben-Hur, Nimrod Wald, Nishil Talati, Ameer Haj Ali, Pierre-Emmanuel Gaillardon, and Shahar Kvatinsky. 2017. Memristive logic: A framework for evaluation and comparison. In *2017 27th International Symposium on Power and Timing Modeling, Optimization and Simulation (PATMOS'17)*. IEEE, 1–8.
- [107] Gurtej S. Sandhu. 2013. Emerging memories technology landscape. In *2013 13th Non-Volatile Memory Technology Symposium (NVMTS'13)*. IEEE, 1–5.
- [108] Karthikeyan Sankaralingam, Ramadass Nagarajan, Haiming Liu, Changkyu Kim, Jaehyuk Huh, Doug Burger, Stephen W. Keckler, and Charles R. Moore. 2003. Exploiting ILP, TLP, and DLP with the polymorphous TRIPS architecture. In *ACM SIGARCH Computer Architecture News*, Vol. 31. ACM, 422–433.
- [109] Vivek Seshadri, Kevin Hsieh, Amirali Boroum, Donghyuk Lee, Michael A. Kozuch, Onur Mutlu, Phillip B. Gibbons, and Todd C. Mowry. 2015. Fast bulk bitwise AND and OR in DRAM. *IEEE Computer Architecture Letters* 14, 2 (2015), 127–131.
- [110] Vivek Seshadri, Donghyuk Lee, Thomas Mullins, Hasan Hassan, Amirali Boroumand, Jeremie Kim, Michael A. Kozuch, Onur Mutlu, Phillip B. Gibbons, and Todd C. Mowry. 2017. Ambit: In-memory accelerator for bulk bitwise operations using commodity DRAM technology. In *Proceedings of the 50th Annual IEEE/ACM International Symposium on Microarchitecture*. ACM, 273–287.
- [111] Ali Shafiee, Anirban Nag, Naveen Muralimanohar, Rajeev Balasubramonian, John Paul Strachan, Miao Hu, R. Stanley Williams, and Vivek Srikumar. 2016. ISAAC: A convolutional neural network accelerator with in-situ analog arithmetic in crossbars. *ACM SIGARCH Computer Architecture News* 44, 3 (2016), 14–26.
- [112] M. A. Shami and A. Hemani. 2012. Classification of massively parallel computer architectures. In *2012 IEEE 26th International Parallel and Distributed Processing Symposium Workshops PhD Forum (IPDPSW'12)*. 344–351. DOI: <https://doi.org/10.1109/IPDPSW.2012.42>
- [113] Patrick Siegl, Rainer Buchty, and Mladen Berekovic. 2016. Data-centric computing frontiers: A survey on processing-in-memory. In *Proceedings of the 2nd International Symposium on Memory Systems*. ACM, 295–308.
- [114] A. Siemon, S. Menzel, A. Chattopadhyay, R. Waser, and E. Linn. 2015. In-memory adder functionality in 1S1R arrays. In *2015 IEEE International Symposium on Circuits and Systems (ISCAS'15)*. IEEE, 1338–1341.
- [115] Anne Siemon, Stephan Menzel, Rainer Waser, and Eike Linn. 2015. A complementary resistive switch-based crossbar array adder. *IEEE Journal on Emerging and Selected Topics in Circuits and Systems* 5, 1 (2015), 64–74.
- [116] Gagandeep Singh, Lorenzo Chelini, Stefano Corda, Ahsan Javed Awan, Sander Stuijk, Roel Jordans, Henk Corporaal, and Albert-Jan Boonstra. 2018. A review of near-memory computing architectures: Opportunities and challenges. In *Proceedings of the 21st Euromicro Conference on Digital System Design (DSD'18)*.

- [117] D. B. Skillicorn. 1988. A taxonomy for computer architectures. *Computer* 21, 11 (Nov. 1988), 46–57. DOI : <https://doi.org/10.1109/2.86786>
- [118] G. Snider. 2005. Computing with hysteretic resistor crossbars. *Applied Physics A: Materials Science & Processing* 80, 6 (2005), 1165–1172.
- [119] Kyomin Sohn, Won-Joo Yun, Reum Oh, Chi-Sung Oh, Seong-Young Seo, Min-Sang Park, Dong-Hak Shin, Won-Chang Jung, Sang-Hoon Shin, Je-Min Ryu, et al. 2017. A 1.2 V 20nm 307GB/s HBM DRAM with at-speed wafer-level IO test scheme and adaptive refresh considering temperature distribution. *IEEE Journal of Solid-State Circuits* 52, 1 (2017), 250–260.
- [120] Harold S. Stone. 1970. A logic-in-memory computer. *IEEE Transactions on Computing* 100, 1 (1970), 73–78.
- [121] Arun Subramaniyan, Jingcheng Wang, Ezhil R. M. Balasubramanian, David Blaauw, Dennis Sylvester, and Reetuparna Das. 2017. Cache automaton. In *Proceedings of the 50th Annual IEEE/ACM International Symposium on Microarchitecture (MICRO-50'17)*. ACM, New York, NY, 259–272. DOI : <https://doi.org/10.1145/3123939.3123986>
- [122] Jinwoo Suh, Eun-Gyu Kim, Stephen P. Crago, Lakshmi Srinivasan, and Matthew C. French. 2003. A performance analysis of PIM, stream processing, and tiled processing on memory-intensive signal processing kernels. In *ACM SIGARCH Computer Architecture News*, Vol. 31. ACM, 410–421.
- [123] Mark R. Thistle and Burton J. Smith. 1988. A processor architecture for Horizon. In *Proceedings of Supercomputing '88. Vol. 1*. IEEE, 35–41.
- [124] Dean M. Tullsen, Susan J. Eggers, and Henry M. Levy. 1995. Simultaneous multithreading: Maximizing on-chip parallelism. In *ACM SIGARCH Computer Architecture News*, Vol. 23. ACM, 392–403.
- [125] Mario Vestias and Horácio Neto. 2014. Trends of CPU, GPU and FPGA for high-performance computing. In *2014 24th International Conference on Field Programmable Logic and Applications (FPL'14)*. IEEE, 1–6.
- [126] Borui Wang, Martin Torres, Dong Li, Jishen Zhao, and Florin Rusu. 2016. Performance implications of processing-in-memory designs on data-intensive applications. In *2016 45th International Conference on Parallel Processing Workshops (ICPPW'16)*. IEEE, 115–122.
- [127] Jue Wang, Xiangyu Dong, Yuan Xie, and Norman P. Jouppi. 2014. Endurance-aware cache line management for non-volatile caches. *ACM Transactions on Architecture and Code Optimization (TACO)* 11, 1 (2014), 4.
- [128] Ying Wang, Yinhe Han, Lei Zhang, Huawei Li, and Xiaowei Li. 2015. ProPRAM: Exploiting the transparent logic resources in non-volatile memory for near data computing. In *Proceedings of the 52nd Annual Design Automation Conference*. ACM, 47.
- [129] Rainer Waser. 2012. Redox-based resistive switching memories. *Journal of Nanoscience and Nanotechnology* 12, 10 (2012), 7628–7640.
- [130] Rainer Waser and Masakazu Aono. 2007. Nanoionics-based resistive switching memories. *Nature Materials* 6, 11 (2007), 833.
- [131] Stephan Wong, Thijs Van As, and Geoffrey Brown. 2008. ρ -VEX: A reconfigurable and extensible softcore VLIW processor. In *International Conference on ICECE Technology, 2008 (FPT'08)*. IEEE, 369–372.
- [132] Wm A. Wulf and Sally A. McKee. 1995. Hitting the memory wall: Implications of the obvious. *ACM SIGARCH Computer Architecture News* 23, 1 (1995), 20–24.
- [133] Lei Xie, Hoang Anh Du Nguyen, Mottaqiallah Taouil, and Koen Bertels Said Hamdioui. 2015. Fast Boolean logic mapped on memristor crossbar. In *2015 33rd IEEE International Conference on Computer Design (ICCD'15)*. IEEE, 335–342.
- [134] Lei Xie, Hoang Anh Du Nguyen, Jintao Yu, Ali Kaichouhi, Mottaqiallah Taouil, Mohammad AlFailakawi, and Said Hamdioui. 2017. Scouting logic: A novel memristor-based logic design for resistive computing. In *IEEE Computer Society Annual Symposium on VLSI (ISVLSI'17)*. IEEE, 335–340.
- [135] Sheng Xu, Xiaoming Chen, Ying Wang, Yinhe Han, Xuehai Qian, and Xiaowei Li. 2018. PIMSim: A flexible and detailed processing-in-memory simulator. *IEEE Computer Architecture Letters* 18, 1 (2018), 6–9.
- [136] J. Joshua Yang, Dmitri B. Strukov, and Duncan R. Stewart. 2013. Memristive devices for computing. *Nature Nanotechnology* 8, 1 (2013), 13–24.
- [137] Leonid Yavits, Shahar Kvatinsky, Amir Morad, and Ran Ginosar. 2015. Resistive associative processor. In *CAL*.
- [138] Jintao Yu, Lei Xie, Mottaqiallah Taouil, and Said Hamdioui. 2018. Memristive devices for computation-in-memory. In *Design, Automation and Test in Europe (DATE'18)*.
- [139] Shimeng Yu and Pai-Yu Chen. 2016. Emerging memory technologies: Recent trends and prospects. *IEEE Solid-State Circuits Magazine* 8, 2 (2016), 43–56.
- [140] Jian-Gang Zhu. 2008. Magnetoresistive random access memory: The path to competitiveness and scalability. *Proceedings of the IEEE* 96, 11 (2008), 1786–1798.

Received December 2018; revised September 2019; accepted October 2019