

A review of near-memory computing architectures

Citation for published version (APA):

Singh, G., Chelini, L., Corda, S., Awan, A. J., Stuijk, S., Jordans, R., Corporaal, H., & Boonstra, A. J. (2018). A review of near-memory computing architectures: Opportunities and challenges. In N. Konofaos, M. Novotny, & A. Skavhaug (Eds.), *Proceedings - 21st Euromicro Conference on Digital System Design, DSD 2018* (pp. 608-617). [8491877] Institute of Electrical and Electronics Engineers. <https://doi.org/10.1109/DSD.2018.00106>

DOI:

[10.1109/DSD.2018.00106](https://doi.org/10.1109/DSD.2018.00106)

Document status and date:

Published: 12/10/2018

Document Version:

Accepted manuscript including changes made at the peer-review stage

Please check the document version of this publication:

- A submitted manuscript is the version of the article upon submission and before peer-review. There can be important differences between the submitted version and the official published version of record. People interested in the research are advised to contact the author for the final version of the publication, or visit the DOI to the publisher's website.
- The final author version and the galley proof are versions of the publication after peer review.
- The final published version features the final layout of the paper including the volume, issue and page numbers.

[Link to publication](#)

General rights

Copyright and moral rights for the publications made accessible in the public portal are retained by the authors and/or other copyright owners and it is a condition of accessing publications that users recognise and abide by the legal requirements associated with these rights.

- Users may download and print one copy of any publication from the public portal for the purpose of private study or research.
- You may not further distribute the material or use it for any profit-making activity or commercial gain
- You may freely distribute the URL identifying the publication in the public portal.

If the publication is distributed under the terms of Article 25fa of the Dutch Copyright Act, indicated by the "Taverne" license above, please follow below link for the End User Agreement:

www.tue.nl/taverne

Take down policy

If you believe that this document breaches copyright please contact us at:

openaccess@tue.nl

providing details and we will investigate your claim.

A Review of Near-Memory Computing Architectures: Opportunities and Challenges

Gagandeep Singh*, Lorenzo Chelini*, Stefano Corda*, Ahsan Javed Awan[†], Sander Stuijk*, Roel Jordans*, Henk Corporaal*, Albert-Jan Boonstra[‡]

*Department of Electrical Engineering, Eindhoven University of Technology, Netherlands
 {g.singh, l.chelini, s.corda, s.stuijk, r.jordans, h.corporaal}@tue.nl

[†]Department of Computing, Imperial College London, UK
 ahsan.awan@imperial.ac.uk

[‡]R&D Department, Netherlands Institute for Radio Astronomy, ASTRON, Netherlands
 boonstra@astron.nl

Abstract—The conventional approach of moving stored data to the CPU for computation has become a major performance bottleneck for emerging scale-out data-intensive applications due to their limited data reuse. At the same time, the advancement in integration technologies have made the decade-old concept of coupling compute units close to the memory (called Near-Memory Computing) more viable. Processing right at the home of data can completely diminish the data movement problem of data-intensive applications.

This paper focuses on analyzing and organizing the extensive body of literature on near-memory computing across various dimensions: starting from the memory level where this paradigm is applied, to the granularity of the application that could be executed on the near-memory units. We highlight the challenges as well as the critical need of evaluation methodologies that can be employed in designing these special architectures. Using a case study, we present our methodology and also identify topics for future research to unlock the full potential of near-memory computing.

Index Terms—near-memory computing, data centric computing, modeling of computer architecture, application characterization, survey

I. INTRODUCTION

Over the years memory technology has not been able to keep up with advancements in processor technology in terms of latency and energy consumption which is infamously termed as the *memory wall* [1]. Earlier, system architects tried to bridge this gap by introducing memory hierarchies which mitigated some of the disadvantages of off-chip DRAMs. But the limited number of pins on the memory package is not able to meet today's bandwidth demands of multicore processors. Furthermore, with the demise of Dennard scaling [2], slowing of Moore's law, and dark silicon computer performance has reached a plateau [3].

At the same time, we are witnessing an enormous amount of data being generated across multiple areas like radio astronomy, material science, chemistry, health sciences etc [4]. In radio astronomy for example, the first phase of the Square Kilometre Array (SKA) aims at processing over 100 terabytes of raw data samples per second, yielding of the order of 300 petabytes of SKA data products annually [5]. The SKA

currently is in the design phase with anticipated construction in South Africa and Australia in the first half of the coming decade. These radio-astronomy applications usually exhibit massive data parallelism and low operational intensity with limited locality. On traditional systems, these applications cause frequent data movement between the memory subsystem and the processor which has a severe impact on performance and energy efficiency. Therefore, a lot of current research is being focused on coming up with innovative manufacturing technologies and architectures to overcome these problems.

Today's memory hierarchy usually consists of multiple levels of cache, a main memory, and a storage. The traditional approach is of moving data all the way up to caches from the storage and then processing it. In contrast, *near-memory computing* (NMC) aims at processing close to where the data resides. This data-centric approach couples compute units close to the data and seeks to minimize the expensive data movements. Conceptually, this principle can be applied to any level of memory subsystem and with several technology options, including processing on buffer-on-board (BoB), edge-bonding small processor dies on DRAM chips etc [6]. Notably, three-dimensional stacking is touted as the true enabler of processing close to the memory. It allows the stacking of logic die and memory together by means of through-silicon vias (TSVs) which helps in reducing memory access latency, power consumption, and provides much higher bandwidth [7]. Micron's Hybrid Memory Cube (HMC) [8], High bandwidth Memory (HBM) [9] from AMD and Hynix, and Samsung's Wide I/O [10] are the memory industries competing 3D memory products.

Figure 1 depicts the system evolution based on the information referenced by a program during execution which is referred to as a *working set* [11]. Prior systems were based on a *CPU-centric* approach where data is moved to the core for processing (Fig.1 (a)-(c)), whereas now with near-memory processing (Fig.1 (d)) the processing cores are brought to the place where data resides. Computation in-memory (Fig.1 (e)) is an extreme step of completely reducing data movement by using novel devices (eg. memristors, phase change memory (PCM)) as these memories have inherent compute capability

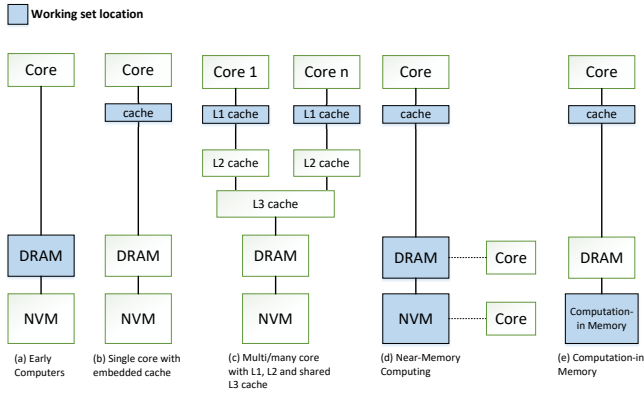


Fig. 1: Classification of computing systems based on working set location

as well. The aim of this paper is to analyze and organize the extensive body of literature related to the novel area of near-memory computing. Figure 2 shows a high level view of our classification which is based on the level in the memory hierarchy and further split into the type of compute implementation (Programmable, Fixed-Function, or Reconfigurable). In our taxonomy we don't include disk-based systems as they can no longer offer timely response due to the high access latency and high failure rate of disks [12]. Nevertheless, there have been research efforts towards providing processing capabilities in the disk. However, it was not adopted widely by the industry due to of the marginal performance improvement that could not justify the associated cost [13], [14]. Instead, we include emerging non-volatile memories termed as storage-class memories (SCMs) [15] which are trying to fill the latency gap between DRAMs and disks. Building up on similar efforts [16], [17], we make the following contributions:

- 1) We analyze and organize the literature related to the landscape of near-memory computing under various dimensions;
- 2) We provide guidelines for design space exploration focusing on near-memory systems;
- 3) We present a case study to illustrate the potential of near-memory computing for memory intensive application with a comparison to the current CPU-centric approach;
- 4) We outline the direction for future research and high-light current challenges in the domain of near-memory computing.

The remainder of this article is structured as follows. Section II outlines the evaluation and classification scheme for NMC at main memory (Section III) and storage class memory (Section IV). Section V highlights the challenges with cache coherence, virtual memory, the lack of programming models and data mapping schemes for NMC. In Section VI we look into the tools and techniques used while performing the design space exploration for these systems. This section also illustrates the importance of application characterization for these systems. Section VII presents a high level analytic approach to model NMC systems with a comparison to a traditional CPU-centric approach. Section VIII highlights the lessons learned and future research directions.

	Property	Abbrev	Description
Memory	Hierarchy	MM SCM HM	Main Memory Storage Class Memory Heterogeneous Memory
	Type	C3D DIMM PCM DRAM SSD LRDIMM	Commercial 3D memory Dual in-line memory module Phase Change Memory Dynamic Random-Access Memory Flash SSD Memory Load-Reduce DIMM
	Integration	US S	Conventional Unstacked Stacked using 2.5D or 3D
Processing	NMC/Host Unit	CPU GPU FPGA CGRA ACC	Central Processing Unit Graphics Processing Unit Field Programmable Gate Array Coarse-Grained Reconfigurable Architecture Application Specific Accelerator
	Implementation	P F R	Programmable Unit Fixed function Unit Reconfigurable Unit
	Granularity	I K A	Instruction Kernel Application
	Host Unit		Type of Host Unit
	Tool	Evaluation technique	A S P
Interoperability	Programming Interface	Y/N	Programming Interface support
	Cache Coherence	Y/N	Mechanism for Cache coherence
	Virtual Memory	Y/N	Virtual Memory Support
App. Domain	Workload		Target Application domain for the architecture

Table 1: Classification table, and legend for table 2

II. CLASSIFICATION AND EVALUATION

This section introduces the classification and evaluation that is used in Section III and IV and is summarized in table 1 and table 2. For each architecture five main categories are evaluated and classified:

- *Memory* - The decision of using what kind of memory is one of the most fundamental questions on which the near-memory architecture depends.
- *Processing* - Type of processing unit implemented and the granularity of processing it performs plays a critical role during the design space exploration.
- *Tool* - Any architecture's success depends heavily on the available tool support. As NMC is a fledgling field, there is lack of generic design methodologies.
- *Interoperability* - One of the main challenges faced by NMC architectures is its inadequacy at providing support for a programming model, cache coherence, virtual memory, and efficient data mapping.
- *Application* - Usually these architectures focus on big data applications/algorithms and this massive amount of data is produced across all the markets. Therefore, in our table we include the domain of the application.

III. PROCESSING NEAR MAIN MEMORY

Processing near main memory has been the most researched interpretation. It can be coupled with different processing units ranging from programmable units to fixed-functional units.

A. Programmable Unit

NDC (2014) Pugsley et al. [20] proposed a near-memory computing (NDC) architecture for MapReduce workloads, in which a central host processor with many energy efficient cores is connected to many daisy-chained 3D-stacked memory devices with simple cores in their logic layer; these cores perform Map operations. Reduce operations, however, are executed on the central host processor because it requires random access to data.

TOP-PIM (2014) Zhang et al. [21] proposed GPU-accelerated architecture. It consists of a main APU (GPU

NMC architecture		Memory			Processing			Tool	Interoperability			App. Domain	
Architecture	Year	Hierarchy	Type	Integration	NMC unit	Implementation	Granularity	Host Unit	Evaluation technique	Programming Interface	Cache Coherence	Virtual Memory	Workload
SmartSSD [18]	2013	SCM	SSD	US	CPU	P	A	CPU	P	Y	Y	N	Database
WILLOW [19]	2014	SCM	SSD	US	CPU	P	K	CPU	P	Y	Y	-	Generic
NDC [20]	2014	MM	C3D	S	CPU	P	K	CPU	S	N	R	N	MapReduce
TOP-PIM [21]	2014	MM	C3D	S	GPU	P	K	CPU	S	N	Y	-	Graph and HPC
AMC [4]	2015	MM	C3D	S	CPU	P	K	CPU	S	Y	Y	Y	HPC
JAFAR [22]	2015	MM	DIMM	US	ACC	F	K	CPU	S	Y	-	Y	Database
TESSERACT [23]	2015	MM	C3D	S	CPU	F	A	CPU	S	Y	Y	N	Graph Processing
Gokhale [24]	2015	MM	C3D	S	ACC	F	K	CPU	S	Y	Y	Y	Generic
HRL [25]	2015	MM	C3D	S	CGRA+FPGA	R	A	CPU	S	N	Y	N	MapReduce
ProPRAM [26]	2015	SCM	PCM	US	CPU	P	I	-	S	Y	-	-	Data Analytics
BlueDBM [27]	2015	SCM	SSD	US	FPGA	R	K	-	P	Y	-	-	Data Analytics
NDA [28]	2015	MM	LRDIMM	S	CGRA	R	K	CPU	S	Y	Y	Y	MapReduce
PIM-enabled [29]	2015	MM	C3D	S	ACC	F	I	CPU	S	Y	Y	Y	Generic
IMPICA [30]	2016	MM	C3D	S	ACC	F	K	CPU	S	Y	Y	Y	Pointer chasing
TOM [31]	2016	MM	C3D	S	GPU	P	K	GPU	S	N	Y	Y	Generic
BISCUIT [32]	2016	SCM	SSD	US	ACC	F	K	CPU	P	Y	-	-	Database
CARIBOU [33]	2017	SCM	DRAM	US	FPGA	R	K	CPU	P	-	-	-	Database
Vermij [34]	2017	MM	C3D	S	ACC	F	A	CPU	S	Y	Y	Y	Sorting
SUMMARIZER [35]	2017	SCM	SSD	US	CPU	P	K	CPU	P	Y	-	-	Database
MONDRIAN [36]	2017	MM	C3D	S	CPU	P	K	CPU	A+S	Y	-	Y	Data Analytics

Table 2: Architectures classification and evaluation, refer to table 1 for a legend

+ CPU) interconnected with high-speed serial links with multiple 3D-memory modules. APU allows code portability and easier programmability. The kernels analyzed span across a huge domain from graph processing to fluid and structural dynamics.

AMC (2015) Nair et al. [4] developed active memory cube (AMC). AMC is built upon the HMC adding several processing elements referred as lanes in the logic layer. The host processor and each AMC are interconnected over a bidirectional link following a daisy-chain topology. For each AMC there are 32 4-slices lanes. Each lane has its own register file, a load/store unit performing read and write operations to the memory contained in the same AMC, and a computational unit. The communications between AMCs should be coordinated by the host processor.

PIM-enabled (2015) Ahn et al. [29] created NMC architectures with the same ISA as the host. They tried to leverage the existing programming model so that the conventional architectures can exploit the PIM concept without changing their programming interface. They implemented it by adding a computation unit composed of computation logic and an SRAM operand buffer, both on host side and in the logic layer of the Hybrid Memory Cube (HMC).

TOM (2016) Hsieh et al. [31] proposed an NMC architecture consisting of main GPU interconnected via high-speed serial links with multiple 3D-stacked memories. Each memory device hosts multiple streaming multiprocessor (SMs) in the logic layer. In line with the philosophy of near-memory computing, blocks of code are offloaded to the SMs in order to reduce off-chip traffic and reduce power consumption due to data movement.

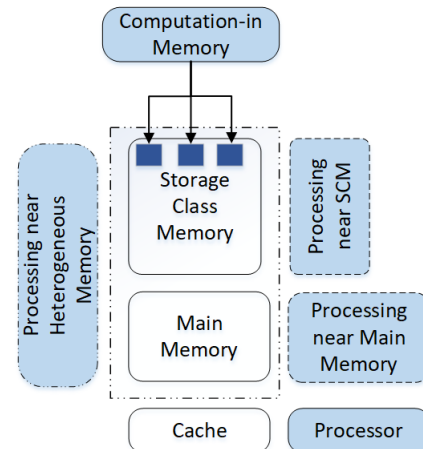


Fig. 2: Processing options in the memory hierarchy

MONDRIAN (2017) Drumond et al. [36] analyzed data analytics operators which are not a fit for NMC. They made a case for optimizing traditional algorithms for NMC system for sequential memory access rather than cache locality.

JAFAR (2015) Xi et al. [22] proposed JAFAR, an NMC accelerator that allows to process filtering operations near the main memory. Thus only relevant data will be pushed up in caches, causing a significant reduction in data movement.

TESSERACT (2015) Ahn et al. [23] focus on graph processing applications. Their architecture consists of one host processor with his own memory and a HMC, with an out-of-order processor mapped to each vault. These cores can see only their own local data partition, but they can communicate to each other using a message passing method. The host

processor as well has access to the entire address space of the HMC. To exploit the high available memory bandwidth in the systems, they developed prefetching mechanisms.

IMPICA (2016) Hsieh et al. [30] accelerated the pointer chasing operations which are ubiquitous in data structures. They proposed adding specialized units that decouple address generation from memory access in the logic layer of 3D stacked memory. These units traverse through the linked data structure in memory and return only the final node found to the CPU. They also proposed completely decoupling the page table of IMPICA from that of CPU to avoid virtual memory related issues.

B. Fixed-Function Unit

Vermij¹ (2017) Vermij et al. [34] evaluated the performance of sorting algorithms in big data workloads running on a 2-socket IBM POWER8 machine. They proposed a heterogeneous system where algorithm phases with high temporal locality are executed on the CPU, while algorithm phases with poor temporal locality are executed on the NMC devices. The architecture proposed consists of a central host processor that relies on two-level memory controllers. A memory-technology agnostic controller located at the CPU side and a memory-specific controller tightly coupled with the memory. The NMC accelerators are placed in the memory-specific controllers and are assisted with an NMC manager. The NMC manager interfaces with the accelerator and the rest of the system and provides support for cache coherency, virtual memory management and communications with the host processor.

C. Re-configurable Unit

Gokhale¹ (2015) Gokhale et al. [24] proposed to place a data rearrangement engine (DRE) in the logic layer of the HMC to accelerate data accesses while still performing the computation on the main CPU. The authors targeted cache unfriendly applications with high memory latency due to irregular access patterns, e.g., sparse matrix multiplication. Each of the DRE engine consists of a scratchpad, a simple controller processor, and a data mover (aka DMA engine). In order to make use of the engines the authors developed an application program interface (API) with different operations. Each operation is issued by the main application running on the central host and served by a control program loaded by the OS on each DRE engine.

HRL (2015) Gao et al. [25] proposed a re-configurable logic architecture called Heterogeneous Re-configurable Logic (HRL) that consists of three main blocks: fine-grained configurable logic blocks (CLBs) for control unit, coarse-grained functional units (FUs) for basic arithmetic and logic operations, and output multiplexer blocks (OMBs) for branch support. Each memory module follows HMC like technology and houses multiple HRL devices in the logic layer. The central host is responsible for data partition and synchronization between NMC units.

NDA (2015) Farmahini et al. [28] proposed three different NMC architectures using coarse-grained reconfigurable arrays

on commodity DRAM modules. This proposal requires minimal change to the DRAM architecture. However, programmer should identify which code will run close to memory. This leads to increased programmer effort for executing compute intensive code on the 3D stack logic layer. Also, it doesn't support direct communication between NMC stacks.

IV. PROCESSING NEAR STORAGE CLASS MEMORY

Flashes and other emerging nonvolatile memories such as phase-change memory (PCM) [37], spin-transfer torque RAM (STT-RAM) [38], memristors [39], etc., are termed as storage-class memories (SCMs) [15]. These memories are trying to fill the latency gap between DRAMs and disks. SCMs like NVRAM are even touted as a future replacement for DRAM [40]. Moving computation in SCM has some of the similar benefits to DRAM in terms of saving in bandwidth, power, latency, and energy but also because of the higher density it allows to work on much larger data-sets as compared to DRAM [41].

A. Programmable Unit

Smart SSD (2013) Kang et al. [18] proposed a model to harness the processing power of the SSD using an object-based communication protocol. They implemented the Smart SSD features in the firmware of a Samsung SSD and modified the Hadoop core and MapReduce framework to use tasklets as a map or a reduce function. To evaluate the prototype, they used a micro-benchmark and log analysis application on both a device and a host. They found that under the current SSD architecture, excessive memory accesses will make the task execution slower than in the host due to the high memory latency and low processing power.

WILLOW (2014) Seshadri et al. [19] proposed a prototype system called Willow, which has processing units called storage processor units (SPU). Each SPU runs a small operating system that manages and enforces security. On the host-side, the Willow driver creates and manages a set of objects that allow the OS and applications to communicate with SPUs. The programmable functionality is provided in the form of SSD Apps. Willow allows programmers to augment and extend the semantics of an SSD with application-specific features without compromising file system protection. The programming model for SSD Apps provides great flexibility, supports the concurrent execution of multiple SSD Apps in Willow, and supports the execution of trusted code in Willow.

ProPRAM (2015) Wang et al. [26] observed that NVM is often naturally incorporated with basic logic like data comparison write or flip-n-write module, and exploited the existing resources inside memory chips to accelerate the key non-compute intensive functions of emerging big data applications.

SUMMARIZER (2017) Koo et al. [35] designed APIs that can be used by the host application to offload a data-intensive task to the inherent ARM based cores inside a SSD processor. Using a fully functional SSD evaluation platform they performed design space exploration of the proposed approach. They evaluate static and dynamic approaches for dividing the work between the host and SSD processor.

¹Architecture has no name, first author's name is shown.

B. Fixed-Function Unit

BISCUIT (2016) Gu et al. [32] presented Biscuit, a near-memory computing framework which allows programmers to write a data-intensive application to run on the host system and the storage system in a distributed manner. The SSD hardware incorporates a pattern matcher IP designed for NMC. When this hardware IP is applied, modified MySQL significantly improves TPC-H performance.

C. Re-configurable Unit

BlueDBM (2015) Jun et al. [27] presented a flash-based platform, called BlueDBM, built of flash storage devices augmented with an application specific FPGA based in-storage processor. The data-sets are stored in the flash array and are read by the FPGA accelerators. Each accelerator implements an array of application-specific distance comparators, used in the high-dimensional nearest-neighbour search algorithms.

CARIBOU (2017) Zsolt et al. [33] explored offloading part of the computation in database engines directly to the storage and deploying FPGAs. They implemented a cuckoo hash table with a slab-based memory allocator to improve the handling of collisions and various values sizes in hardware-based key-value stores. Lookups and scans are performed on the same data to minimize data transferred over the network. They implemented a runtime parametrizable selection operator both for structured and unstructured data in an effort to reduce data movement further.

V. CHALLENGES OF NEAR-MEMORY COMPUTING

One of the biggest challenges in near-memory computing is the lack of interoperability with caches and virtual memory of the host processor due to unconventional programming models. In this section we will explain the current issues with programming model, data mapping, virtual memory support and cache coherency in the context of NMC. Note there are many more design challenges, like design space exploration, reliable simulators which can deal with heterogeneous environments, etc. They are detailed in section VI.

Virtual Memory support

Support for virtual memory is achieved using: *paging* or *direct segment*. In the former approach most NMC systems adopt a software-managed translation lookaside buffer (TLB) to provide a mapping between virtual and physical addresses [6], [42], [31], [43]. Azarkhish et al. in [44] observed that an OS managed TLB may not be the optimal solution and proposed an hardware implementation where a simple controller is responsible for fetching the rule in the page table. In [30] Hsieh et al. notice that for pointer-chasing applications, the accelerator works only on specific data structures that can be mapped onto a contiguous region of the virtual address space. As a consequence, the virtual-to-physical translation for the accelerator is limited to a smaller region. Other works such as [45], [46] decided for a simplified approach where part of the linear virtual address is mapped to physical memory using a direct segment rather than a page. This mechanism allows to remove TLB misses overhead and greatly simplifies the hardware [47].

Cache Coherence

Although in some works the cache coherence has not been treated [46], [48], we classify the approaches proposed in literature as *restricted memory region* and *non-restricted region*.

1) Restricted Memory Regions

Farahani et al. [28] divided the memory in two regions: one for the host processor and one for CGRAs, which is uncacheable. A similar approach has been used by Ahn et al. [23] for graph processing algorithms. Another strategy, proposed by Ahn et al. [29], provides a simple hardware-based solution in which the PIM operations are restricted to only one last level cache block, due to which they can monitor the cache block and request for invalidation or write-back if required.

2) Non-Restricted Memory Regions

Non-restricted memory regions can often lead to a significant amount of coherence traffic. Pattnaik et al. [49] proposed to maintain coherence between the main/host GPU and near-memory ones, flushing the L2 cache in the main GPU after kernel execution. With this approach, applications with highly irregular memory accesses lead to significant amount of overhead. Another way is to implement a look-up based approach as in the case of Hsieh et al. [31]. The NMC units record the cache line address that has been updated by the offloaded block and once when the offloaded block is processed they send this address back to the host system. Subsequently, the host system gets the latest data from memory by invalidating the reported cache lines.

Programming Model

A critical challenge in adoption of NMC is to support a heterogeneous processing environment of a host system and near-memory processing unit. It is not trivial to determine how and who identifies which part of an application should run in the near-memory processing units. Currently, this is left onto the programmer or the compiler [31] to mark sections of the code that must be processed near-memory compute units [17]. In another approach, Ahn et al. [29] use some special set of NMC instructions which invoke NMC logic units. However, this approach calls for sophisticated mechanism (e.g., taking care of the data locality) on the host side to avoid unnecessary exchange of information. Hence, there is still a lot of research required in coming up with an efficient approach to ease the programming burden.

Data mapping

The absence of effective data mapping mechanisms can severely hamper the benefits of processing close to memory. The data should be mapped in such a way that data that will be used by the near-memory processing units should be available in the vicinity. Hence, it is crucial to look into effective data mapping schemes (static and dynamic)

In TOM [31], they proposed a software/hardware cooperative method to predict which pages of the memory will be used by the offloaded code, and they tried to minimize the bandwidth consumption by placing those pages in the memory stack closest to the offloaded code. Yitbarek et al. [46] proposed a data-mapping scheme which they integrated in a

Metrics	Description
Reuse distance (temporal locality)	The number of distinct memory accesses since the last access to a given memory address
Data Stream Stride (spatial locality)	The difference in memory addresses between two consecutive memory accesses
Application bandwidth	The amount of I/O bandwidth utilized by a particular application
Instruction mix	Provides a general overview of the application
Parallelism	Inherent application parallelism such as instruction level parallelism, thread level parallelism, data level parallelism
Branch Entropy	Measures the randomness of branch behavior which is essential to provide speculation support

Table 3: Several metrics for micro-architecture independent analysis of workloads

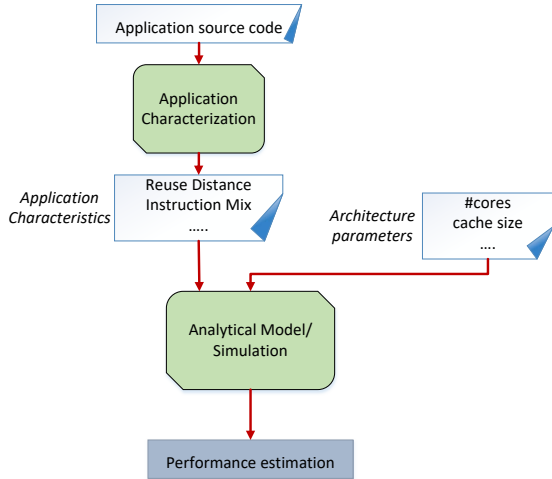


Fig. 3: Design space exploration for system evaluation

platform that used some accelerators near 3D-stacked memory. Near-memory workloads often access data from contiguous addresses for this reason they re-mapped data at page granularity in order to have contiguous addresses in the same memory vaults.

VI. DESIGN SPACE EXPLORATION FOR NMC

As will be clear from the previous classification section, the design space of NMC is huge. In order to understand and evaluate this space, design space exploration (DSE) for NMC system is required, as shown in figure 3.

A. Application Characterization

More than ever, application characterization has taken a crucial role in systems design due to the increasing number of new data-intensive applications. It is used to extract information by using specific metrics, in order to decide which architecture could have the best performance and energy efficiency. In table 3, we have highlighted some of the useful metrics for near-memory architectures. Application characterization can be done in one of the following ways:

Microarchitecture-Dependent Workload Characterization

Traditionally workload characterization is done using hardware performance counters present on our current microprocessors. Awan et al. [50] use hardware performance counters to characterize the scale-out big data processing workloads into CPU-bound, memory-bound, and I/O-bound applications and propose programmable logic near DRAM and NVRAM. However, the use of hardware performance counters is limited by the impact of micro-architecture features like cache size,

issue width etc [51]. To overcome this problem, recently there has been a strong push towards following a ISA independent characterization approach.

Microarchitecture-Independent Workload Characterization

Hoste et al. [52] proposed a Microarchitecture-Independent Workload Characterization (MICA) in which they used Pin [53] as basis. Also, they proposed new methodologies to analyze these characteristics, based on principal components analysis (PCA), which gives a reduction to data set's dimensionality and removes correlation from the data set, and also a genetic algorithm (GA), which reduces the data set's dimensionality, but with a retained dimension easily to understand.

Anghel et al. [54] presented a Platform-Independent Software Analysis (PISA) tool, based on the LLVM compiler. It takes a source code as input and generates a LLVM-IR bitstream which is instrumented with some function calls to a tool library that implements different workload analysis. Characteristics which PISA can extract are: Instruction Mix, Instruction-Level-Parallelism, Memory Access Patterns, Branch Entropy, Communication Patterns. Unlike other tools, PISA is also capable of analyzing multi-threaded programs.

In the same way Caparrós et al. [55] used the LLVM framework for their tool. It's only for sequential code, but there are interesting metrics that are extracted such as: Data-Level-Parallelism (DLP), Thread-Level-Parallelism (TLP) and Task-Level Parallelism. This analysis is made through a Direct-Acyclic-Graph (DAG) that searches for the real data dependencies and at various granularity levels, e.g. instructions, basic blocks, functions.

Using the Intermediate Representation (IR) from the Just-In-Time (JIT) compiler Shao et al. [51] extracted number of opcodes, the value of branch entropy, the value of memory entropy, the unique number of static instructions, and the unique number of data addresses for sequential benchmarks.

B. Performance Evaluation Techniques

Architects often make use of various evaluation techniques to navigate the design-space, avoiding the cost of chip fabrication. Based on the level of detail required, architects make use of analytic models, or more detailed functional or cycle accurate simulators, and even design prototypes. As the field of NMC does not have very mature tools and techniques, researchers often spend a lot of time building the appropriate evaluation environment. [56]

Analytic Modeling

Analytic models abstracts low level system details and provide quick performance estimates at the cost of accuracy. In

Simulator	Year	Category	NMC capabilities
Sinuca [60]	2015	Cycle-Accurate	Yes
HMC-SIM [61]	2016	Cycle-Accurate	Limited
CasHMC [62]	2016	Cycle-Accurate	No
SMC [44]	2016	Cycle-Accurate	Yes
CLAPPS [56]	2017	Cycle-Accurate	Yes

Table 4: Open source simulators

the early design stage, system architects are faced with large design choices which range from semiconductor physics, to circuit level, micro-architectural properties (eg., issue width, branch prediction), cooling concerns, and other system design constrains [57]. During the first stage of design-space exploration analytic models can provide quick estimate number. Mirzadeh et al. [58] studied a join workload, on multiple HMC like 3D-stacked DRAM devices connected together via SerDes links using a first-order analytic model. Zhang et al [59] designed an analytic model using machine learning techniques to estimate the final device performance.

Simulation Based Modeling

To achieve more accurate performance numbers, architects often resort to modeling the entire micro-architecture precisely. However, this approach can be quite slow compared to analytic techniques. Hyeokjun et al. [63] evaluated the potential of NMC for machine learning (ML) using a full-fledged simulator of multi-channel SSD that can execute various ML algorithms on data stored on the SSD. Similarly, Jo et al. [64] developed the iSSD simulator based on the gem5 simulator [65].

Ranganathan et al. [66], [40] for their nano-stores architecture used a hybrid methodology which is a bottom-up approach in which they built a high level model which breaks down applications into various phases, based on compute, network, and I/O sub system activities. This takes input from the low level performance and power models regarding the performance of each phase. They used COTSon [67] for detailed micro-architectural simulation.

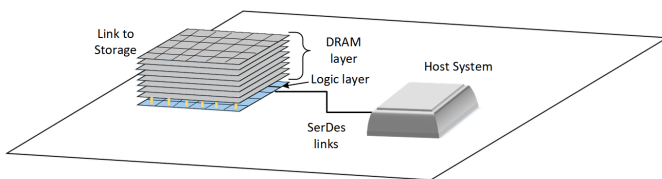


Fig. 4: Our 3D stacked near-memory system connected to a host system

VII. CASE STUDY

We built a high level analytic model which incorporates evaluation methodology as described in figure 3. To see the potential of computing close to the memory and when it is useful to follow a data-centric approach we make a comparison of a CPU-centric approach consisting of a multi-core system which has DDR3 as a main memory, and compared it to a data-centric approach which has HMC like 3D stacked memory instead of DDR3 and is connected to a host system. In the multi-core system each core has its private L1 cache

Description	Symbol	ARM
Cores	n_{cores}	4
Core clock frequency	f_{core}	3 GHz
L1 size	s_{l1}	32 KB
L2 size	s_{l2}	256 KB
DRAM size	s_{DRAM}	4 GB
L1 cache bandwidth	bw_{l1}	137 GB/s
L2 cache bandwidth	bw_{l2}	137 GB/s
DRAM bandwidth	bw_{DRAM}	17 GB/s
L1 cache hit latency	l_{l1}	1 cycles
L2 cache hit latency	l_{l2}	2 cycles
DRAM hit latency	l_{DRAM}	165 cycles
NVM hit latency	l_{ext}	800 cycles

Table 5: System parameters and example set representing ARM Cortex-A9 (ARM) processor

and a shared L2 cache. Figure 4 shows our 3D stacked near-memory system connected with a host system through a serializer/deserializer (SerDes) I/O link. We embed near-memory compute units in the logic layer of 3D stack memory. The near-memory compute units are modeled as ARM Cortex-A9 units (table 5), which are placed in the logic layer of 3D stacked DRAM.

We consider a 4-GB HMC like design, with 8 DRAM layers. Each DRAM layer is split into 16 dual banked partitions with 4 vertical partitions forming a vault. Each vault has its own vault controller in the logic layer which is able to access the DRAM layers at a 10 GB/s bandwidth [4]. In the near-memory system we assume certain parts of computation, which are memory intensive, to be offloaded to the near-memory compute units and the other parts will be executed in the host system. The number of vaults inside the HMC and external SerDes links are kept as varying parameters which have an impact on internal and the external bandwidth of the 3D memory respectively.

Performance calculation: The model estimates the execution time as the ratio between the number of memory access required in order to move the block of data to and from the I/O system and the available I/O bandwidth of each memory subsystem. The execution time takes into account the time to process the instructions and latency due to memory access at different memory levels. The latency of different memory types have been mentioned in table 5.

Energy calculation: A power model for both CPU-centric and near-memory system were built which take into account dynamic as well as static energy. Dynamic energy is the product of the energy per memory access and number of memory accesses at each level of the memory hierarchy. It is assumed to scale linearly with utilization. Static energy is estimated as the product of time and static power for each component. It is scaled based on the number of cores and the cache size.

For modeling the HMC we referred to [7] and [20]. We consider energy per bit of 3.7 pJ/b for accessing the DRAM layers and 1.5 pJ/b for the operations and data movements in the HMC logic layer. Also a static power of 0.96 W [20] was considered to model the additional components in the logic layer.

Design Space Exploration: We vary the cache miss rate, which encapsulates application characteristics, both at L1 and L2 level to see the impact it has on performance of the entire

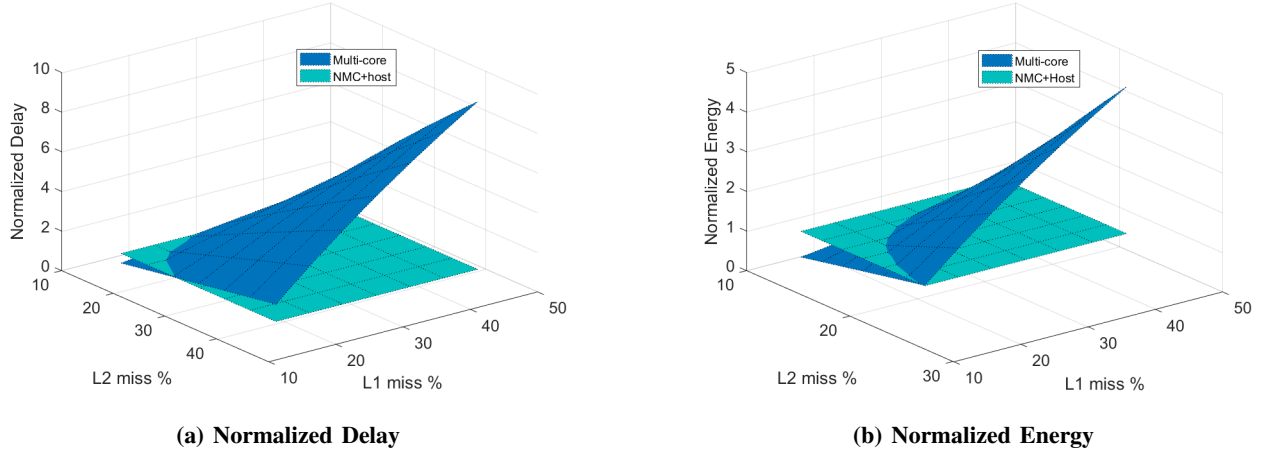


Fig. 5: Performance and energy comparison between multi-core and NMC with a host system using high-level analytic model

system. In case of a near-memory system, we considered the scenario when all access are done to the different vaults inside the HMC memory which can exploit the inherent parallelism offered by the 3D memory. In figure 5, the x and y axis represent the miss rate at L1 and L2 respectively. Results in figure 5 are normalized to NMC+host system. In both the figures we can see a similar trend if the miss rate (L1 and L2) increases the multi core system performance degrades compared to the NMC systems this because of increase in data movement. The data has to be fetched from the off-chip DRAM. If the data brought into the caches is not reused by the CPU, the use of caching becomes inefficient. Therefore, the application (or function) with low locality should take advantage of the near-memory approach whereas the other application (or function) with high locality will benefit more from the traditional approach.

VIII. GUIDELINES FOR FUTURE DIRECTION

Based on analysis of many existing NMC systems we identify the following key topics for future research, which we regard as essential to unlock the full potential of processing close to memory:

- Currently, it is unclear which emerging memory technology best supports near-memory architectures.
- The field requires a generic set of tools and techniques for these novel systems, as often researchers have to spend a significant amount of time and effort in building the required simulation environment.
- Although many architectures feature some sort of tool-flow, very few are open source. Often, the tool-flow is very architecture specific, thereby limiting reproducibility.
- Most of the evaluated architectures focus on the compute aspect. Few architectures focus on providing cache coherency and virtual memory support.
- More exploration is required for interconnect network between the near-memory compute units and also between the host and near-memory compute system.

- At application level, algorithms need to be adapted according to the location of the data for energy efficient processing.
- DRAM and NVM have different memory attributes. A hybrid design can revolutionize our current systems. Processing near heterogeneous memories is a new research topic and in future, we expect a lot of interest in this direction.
- 3D stacking also needs novel power and thermal solutions as traditional heat sink technology will not suffice if we want to add more computation close to the memory.
- Static and dynamic decision support for offloading processing and data to near-memory systems.
- Appropriate data (re-)mappings are required for the near-memory architectures.

IX. CONCLUSION

The improvements in performance and energy have stimulated a great amount of research in processing close to the memory. It can be the right approach for certain big data applications. However, to embrace this paradigm we need to provide a complete ecosystem from hardware to software. In this paper, we analyzed and organized the extensive literature of placing compute units close to the memory using different synonyms (eg. *Processing-in memory*, *Near-data Processing*) under the umbrella of *Near-Memory Computing*. This is done to distinguish it from the in-situ *Computation-in Memory* (CIM) by means of novel non-volatile memories like memristors. It provides systematization for positioning of new proposals within the existing body of work. We stressed on the demand of mature tools and techniques and outlined a methodology for the design space exploration for these novel architectures. Furthermore, based on our analysis of existing NMC systems we identified key topics for future research, which we regard as essential to unlock the full potential of processing close to memory.

ACKNOWLEDGMENT

This work was performed in the framework of Horizon 2020 program and is funded by European Commission under Marie Skłodowska-Curie Innovative Training Networks European Industrial Doctorate (Project ID: 676240).

REFERENCES

- [1] W. A. Wulf and S. A. McKee, "Hitting the memory wall: Implications of the obvious," *SIGARCH Comput. Archit. News*, vol. 23, no. 1, pp. 20–24, Mar. 1995.
- [2] R. H. Dennard, F. H. Gaensslen, V. L. Rideout, E. Bassous, and A. R. LeBlanc, "Design of ion-implanted MOSFET's with very small physical dimensions," *IEEE Journal of Solid-State Circuits*, vol. 9, no. 5, pp. 256–268, Oct 1974.
- [3] H. Esmaeilzadeh *et al.*, "Dark silicon and the end of multicore scaling," *IEEE Micro*, vol. 32, no. 3, pp. 122–134, May 2012.
- [4] R. Nair *et al.*, "Active memory cube: A processing-in-memory architecture for exascale systems," *IBM Journal of Research and Development*, vol. 59, no. 2/3, pp. 17–1, 2015.
- [5] R. Jongerius, S. Wijnholds, R. Nijboer, and H. Corporaal, "An end-to-end computing model for the square kilometre array," *Computer*, vol. 47, no. 9, pp. 48–54, Sept 2014.
- [6] M. Gao, G. Ayers, and C. Kozyrakis, "Practical near-data processing for in-memory analytics frameworks," in *2015 International Conference on Parallel Architecture and Compilation (PACT)*, Oct 2015, pp. 113–124.
- [7] J. Jeddelloh and B. Keeth, "Hybrid memory cube new DRAM architecture increases density and performance," in *2012 Symposium on VLSI Technology (VLSIT)*, June 2012, pp. 87–88.
- [8] J. T. Pawlowski, "Hybrid memory cube (hmc)," in *2011 IEEE Hot Chips 23 Symposium (HCS)*, Aug 2011, pp. 1–24.
- [9] D. U. Lee *et al.*, "25.2 a 1.2v 8Gb 8-channel 128GB/s high-bandwidth memory (HBM) stacked DRAM with effective microbump I/O test methods using 29nm process and TSV," in *2014 IEEE International Solid-State Circuits Conference Digest of Technical Papers (ISSCC)*, Feb 2014, pp. 432–433.
- [10] J. S. Kim *et al.*, "A 1.2 v 12.8 GB/s 2 Gb mobile wide-I/O DRAM with 4 times 128 I/Os using TSV based stacking," *IEEE Journal of Solid-State Circuits*, vol. 47, no. 1, pp. 107–116, Jan 2012.
- [11] S. Hamdioui, L. Xie, H. A. D. Nguyen, M. Taouil, K. Bertels, H. Corporaal, H. Jiao, F. Cathoor, D. Wouters, L. Eike, and J. van Lunteren, "Memristor based computation-in-memory architecture for data-intensive applications," in *2015 Design, Automation Test in Europe Conference Exhibition (DATE)*, March 2015, pp. 1718–1725.
- [12] H. Zhang, G. Chen, B. C. Ooi, K. L. Tan, and M. Zhang, "In-memory big data management and processing: A survey," *IEEE Transactions on Knowledge and Data Engineering*, vol. 27, no. 7, pp. 1920–1948, July 2015.
- [13] K. Keeton, D. A. Patterson, and J. M. Hellerstein, "A case for intelligent disks (IDISKS)," *SIGMOD Rec.*, vol. 27, no. 3, pp. 42–52, Sep. 1998.
- [14] E. Riedel, G. A. Gibson, and C. Faloutsos, "Active storage for large-scale data mining and multimedia," in *Proceedings of the 24rd International Conference on Very Large Data Bases*, ser. VLDB '98. San Francisco, CA, USA: Morgan Kaufmann Publishers Inc., 1998, pp. 62–73.
- [15] R. Nair, "Evolution of memory architecture," *Proceedings of the IEEE*, vol. 103, no. 8, pp. 1331–1345, Aug 2015.
- [16] P. Siegl, R. Buchty, and M. Berekovic, "Data-centric computing frontiers: A survey on processing-in-memory," in *Proceedings of the Second International Symposium on Memory Systems*. ACM, 2016, pp. 295–308.
- [17] S. Ghose, K. Hsieh, A. Boroumand, R. Ausavarungnirun, and O. Mutlu, "Enabling the adoption of processing-in-memory: Challenges, mechanisms, future research directions," *arXiv preprint arXiv:1802.00320*, 2018.
- [18] Y. Kang *et al.*, "Enabling cost-effective data processing with smart SSD," in *Mass Storage Systems and Technologies (MSST), 2013 IEEE 29th Symposium on*. IEEE, 2013, pp. 1–12.
- [19] S. Seshadri *et al.*, "Willow: A user-programmable SSD," in *OSDI*, 2014, pp. 67–80.
- [20] S. H. Pugsley *et al.*, "NDC: Analyzing the impact of 3d-stacked memory+logic devices on mapreduce workloads," in *2014 IEEE International Symposium on Performance Analysis of Systems and Software (ISPASS)*, March 2014, pp. 190–200.
- [21] D. Zhang *et al.*, "TOP-PIM: throughput-oriented programmable processing in memory," in *Proceedings of the 23rd international symposium on High-performance parallel and distributed computing*. ACM, 2014, pp. 85–98.
- [22] S. L. Xi *et al.*, "Beyond the wall: Near-data processing for databases," in *Proceedings of the 11th International Workshop on Data Management on New Hardware*, ser. DaMoN'15. New York, NY, USA: ACM, 2015, pp. 2:1–2:10.
- [23] J. Ahn *et al.*, "A scalable processing-in-memory accelerator for parallel graph processing," in *2015 ACM/IEEE 42nd Annual International Symposium on Computer Architecture (ISCA)*, June 2015, pp. 105–117.
- [24] M. Gokhale, S. Lloyd, and C. Hajas, "Near memory data structure rearrangement," in *Proceedings of the 2015 International Symposium on Memory Systems*. ACM, 2015, pp. 283–290.
- [25] M. Gao and C. Kozyrakis, "HRL: Efficient and flexible reconfigurable logic for near-data processing," in *2016 IEEE International Symposium on High Performance Computer Architecture (HPCA)*, March 2016, pp. 126–137.
- [26] Y. Wang *et al.*, "ProPRAM: exploiting the transparent logic resources in non-volatile memory for near data computing," in *Proceedings of the 52nd Annual Design Automation Conference*. ACM, 2015, p. 47.
- [27] S. W. Jun *et al.*, "BlueDBM: An appliance for big data analytics," *SIGARCH Comput. Archit. News*, vol. 43, no. 3, pp. 1–13, Jun. 2015.
- [28] A. Farmahini-Farahani, J. H. Ahn, K. Morrow, and N. S. Kim, "NDA: Near-DRAM acceleration architecture leveraging commodity DRAM devices and standard memory modules," in *2015 IEEE 21st International Symposium on High Performance Computer Architecture (HPCA)*, Feb 2015, pp. 283–295.
- [29] J. Ahn, S. Yoo, O. Mutlu, and K. Choi, "PIM-enabled instructions: a low-overhead, locality-aware processing-in-memory architecture," in *Proceedings of the 42nd Annual International Symposium on Computer Architecture*. ACM, 2015, pp. 336–348.
- [30] K. Hsieh *et al.*, "Accelerating pointer chasing in 3D-stacked memory: Challenges, mechanisms, evaluation," in *Computer Design (ICCD), 2016 IEEE 34th International Conference on*. IEEE, 2016, pp. 25–32.
- [31] K. Hsieh, E. Ebrahim *et al.*, "Transparent offloading and mapping (TOM): Enabling programmer-transparent near-data processing in GPU systems," in *2016 ACM/IEEE 43rd Annual International Symposium on Computer Architecture (ISCA)*, June 2016, pp. 204–216.
- [32] B. Gu *et al.*, "Biscuit: A framework for near-data processing of big data workloads," in *Computer Architecture (ISCA), 2016 ACM/IEEE 43rd Annual International Symposium on*. IEEE, 2016, pp. 153–165.
- [33] Z. István, D. Sidler, and G. Alonso, "Caribou: intelligent distributed storage," *Proceedings of the VLDB Endowment*, vol. 10, no. 11, pp. 1202–1213, 2017.
- [34] E. Vermij, L. Fiorin, C. Hagleitner, and K. Bertels, "Sorting big data on heterogeneous near-data processing systems," in *Proceedings of the Computing Frontiers Conference*. ACM, 2017, pp. 349–354.
- [35] G. Koe *et al.*, "Summarizer: trading communication with computing near storage," in *Proceedings of the 50th Annual IEEE/ACM International Symposium on Microarchitecture*. ACM, 2017, pp. 219–231.
- [36] D. L. De Oliveira, M. Paulo, A. Daglis, N. Mirzadeh, D. Ustiugov, J. Picorel Obando, B. Falsafi, B. Grot, and D. Pnevmatikatos, "The Mondrian Data Engine," in *Proceedings of the 44th International Symposium on Computer Architecture*, no. EPFL-CONF-227947, 2017.
- [37] M. K. Qureshi, V. Srinivasan, and J. A. Rivers, "Scalable high performance main memory system using Phase-Change Memory technology," *SIGARCH Comput. Archit. News*, vol. 37, no. 3, pp. 24–33, Jun. 2009.
- [38] M. Hosomi *et al.*, "A novel nonvolatile memory with spin torque transfer magnetization switching: spin-ram," in *IEEE International Electron Devices Meeting, 2005. IEDM Technical Digest.*, Dec 2005, pp. 459–462.
- [39] L. Chua, "Memristor-the missing circuit element," *IEEE Transactions on Circuit Theory*, vol. 18, no. 5, pp. 507–519, September 1971.
- [40] P. Ranganathan, "From microprocessors to nanostores: Rethinking data-centric systems (vol 44, pg 39, 2010)," *COMPUTER*, vol. 44, no. 3, pp. 6–6, 2011.
- [41] L. C. Quero *et al.*, "Self-sorting SSD: Producing sorted data inside active SSDs," in *Mass Storage Systems and Technologies (MSST), 2015 31st Symposium on*. IEEE, 2015, pp. 1–7.
- [42] Z. Sura *et al.*, "Data access optimization in a processing-in-memory system," in *Proceedings of the 12th ACM International Conference on Computing Frontiers*, ser. CF '15. New York, NY, USA: ACM, 2015, pp. 6:1–6:8.
- [43] M. Wei, M. Snir, J. Torrellas, and R. B. Tremaine, "A near-memory processor for vector, streaming and bit manipulation workloads," in *In*

- The Second Watson Conference on Interaction between Architecture, Circuits, and Compilers*, 2005.
- [44] E. Azarkhish, D. Rossi, I. Loi, and L. Benini, "Design and evaluation of a processing-in-memory architecture for the smart memory cube," in *Proceedings of the 29th International Conference on Architecture of Computing Systems – ARCS 2016 - Volume 9637*. New York, NY, USA: Springer-Verlag New York, Inc., 2016, pp. 19–31.
- [45] J. Lee, J. H. Ahn, and K. Choi, "Buffered compares: Excavating the hidden parallelism inside dram architectures with lightweight logic," in *2016 Design, Automation Test in Europe Conference Exhibition (DATE)*, March 2016, pp. 1243–1248.
- [46] S. F. Yitbarek, T. Yang, R. Das, and T. Austin, "Exploring specialized near-memory processing for data intensive operations," in *2016 Design, Automation Test in Europe Conference Exhibition (DATE)*, March 2016, pp. 1449–1452.
- [47] A. Basu *et al.*, "Efficient virtual memory for big memory servers," *SIGARCH Comput. Archit. News*, vol. 41, no. 3, pp. 237–248, Jun. 2013.
- [48] A. Farmahini-Farahani *et al.*, "Drama: An architecture for accelerated processing near memory," *IEEE Computer Architecture Letters*, vol. 14, no. 1, pp. 26–29, Jan 2015.
- [49] A. Pattnaik *et al.*, "Scheduling techniques for GPU architectures with processing-in-memory capabilities," in *2016 International Conference on Parallel Architecture and Compilation Techniques (PACT)*, Sept 2016, pp. 31–44.
- [50] A. J. Awan *et al.*, "Identifying the potential of near data processing for apache spark," in *Proceedings of the International Symposium on Memory Systems, MEMSYS 2017, Alexandria, VA, USA, October 02 - 05, 2017*, 2017, pp. 60–67.
- [51] Y. S. Shao and D. Brooks, "ISA-independent workload characterization and its implications for specialized architectures," in *2013 IEEE International Symposium on Performance Analysis of Systems and Software (ISPASS)*, April 2013, pp. 245–255.
- [52] K. Hoste and L. Eeckhout, "Microarchitecture-independent workload characterization," *IEEE Micro*, vol. 27, no. 3, pp. 63–72, May 2007.
- [53] C. K. Luk *et al.*, "Pin: Building customized program analysis tools with dynamic instrumentation," *SIGPLAN Not.*, vol. 40, no. 6, pp. 190–200, Jun. 2005.
- [54] A. Anghel *et al.*, "An instrumentation approach for hardware-agnostic software characterization," *International Journal of Parallel Programming*, vol. 44, no. 5, pp. 924–948, Oct 2016.
- [55] V. C. Cabezas and P. Stanley-Marbell, "Parallelism and data movement characterization of contemporary application classes," in *SPAA 2011: Proceedings of the 23rd Annual ACM Symposium on Parallelism in Algorithms and Architectures, San Jose, CA, USA, June 4-6, 2011 (Co-located with FCRC 2011)*, 2011, pp. 95–104.
- [56] O. Junior *et al.*, "A generic processing in memory cycle accurate simulator under hybrid memory cube architecture," 2017.
- [57] R. Jongerius *et al.*, "Analytic multi-core processor model for fast design-space exploration," *IEEE Transactions on Computers*, 2017.
- [58] N. Mirzadeh, Y. O. Koçberber, B. Falsafi, and B. Grot, "Sort vs. hash join revisited for near-memory execution," in *5th Workshop on Architectures and Systems for Big Data (ASBD 2015)*, no. EPFL-CONF-209121, 2015.
- [59] L. Xu, D. P. Zhang, and N. Jayasena, "Scaling deep learning on multiple in-memory processors," 2015.
- [60] M. A. Z. Alves *et al.*, "SiNUCA: A validated micro-architecture simulator," in *2015 IEEE 17th International Conference on High Performance Computing and Communications, 2015 IEEE 7th International Symposium on Cyberspace Safety and Security, and 2015 IEEE 12th International Conference on Embedded Software and Systems*, Aug 2015, pp. 605–610.
- [61] J. D. Leidel and Y. Chen, "HMC-Sim-2.0: A simulation platform for exploring custom memory cube operations," in *2016 IEEE International Parallel and Distributed Processing Symposium Workshops (IPDPSW)*, May 2016, pp. 621–630.
- [62] D. I. Jeon and K. S. Chung, "CasHMC: A cycle-accurate simulator for hybrid memory cube," *IEEE Computer Architecture Letters*, vol. 16, no. 1, pp. 10–13, Jan 2017.
- [63] H. Choe *et al.*, "Near-data processing for machine learning," *CoRR*, vol. abs/1610.02273, 2016.
- [64] Y.-Y. Jo *et al.*, "Data mining in intelligent SSD: Simulation-based evaluation," in *Big Data and Smart Computing (BigComp), 2016 International Conference on*. IEEE, 2016, pp. 123–128.
- [65] N. Binkert *et al.*, "The gem5 simulator," *SIGARCH Comput. Archit. News*, vol. 39, no. 2, pp. 1–7, Aug. 2011.
- [66] J. Chang, P. Ranganathan, T. Mudge, D. Roberts, M. A. Shah, and K. T. Lim, "A limits study of benefits from nanostore-based future data-centric system architectures," in *Proceedings of the 9th conference on Computing Frontiers*. ACM, 2012, pp. 33–42.
- [67] E. Argollo, A. Falcón, P. Faraboschi, M. Monchiero, and D. Ortega, "COTSon: infrastructure for full system simulation," *SIGOPS Oper. Syst. Rev.*, vol. 43, no. 1, pp. 52–61, 2009.