

# Accelerating Time Series Analysis via Processing using Non-Volatile Memories

Ivan Fernandez<sup>§†¶</sup> \*Aditya Manglik<sup>†</sup> \*Christina Giannoula<sup>†‡</sup> Ricardo Quisiant<sup>§</sup> Nika Mansouri Ghiasi<sup>†</sup>  
 Juan Gómez-Luna<sup>†</sup> Eladio Gutierrez<sup>§</sup> Oscar Plata<sup>§</sup> Onur Mutlu<sup>†</sup>

<sup>§</sup>University of Malaga <sup>†</sup>ETH Zürich <sup>¶</sup>Barcelona Supercomputing Center <sup>‡</sup>National Technical University of Athens

**Abstract**—*Time Series Analysis (TSA)* is a critical workload for consumer-facing devices. For instance, millions of devices and sensors continuously collect a large amount of data, often stored as a time series (sequential collection of organized data points). Accelerating TSA is vital for many domains as it enables the extraction of valuable information and predict future events, e.g., early detection of earthquakes or heart anomalies. The state-of-the-art algorithm in TSA is the *subsequence Dynamic Time Warping (sDTW)* algorithm. sDTW is a dynamic programming algorithm with a highly effective kernel to detect anomalies and similarities using simple metrics (e.g., absolute difference). However, sDTW’s computation complexity increases quadratically with the time series’ length, resulting in two performance implications. First, the amount of data parallelism available at the application level is significantly higher than the small number of processing units enabled by commodity systems (e.g., CPUs). Second, sDTW is bottlenecked by memory because it 1) has low arithmetic intensity and 2) incurs a large memory footprint, leading to expensive data access costs and data transfers between memory and processors. To tackle these two challenges, we leverage Processing-using-Memory (PuM) by performing in-situ computation where data resides, using the memory cells. PuM provides a promising solution to alleviate data movement bottlenecks and exposes immense parallelism.

In this work, we present *MATSA*, the *first* MRAM-based Accelerator for Time Series Analysis. The key idea is to exploit magneto-resistive memory crossbars to enable energy-efficient and fast time series computation in memory while overcoming endurance issues of other non-volatile memory technologies. MATSA provides the following key benefits: 1) it leverages high levels of parallelism in the memory substrate by exploiting column-wise arithmetic operations, and 2) it significantly reduces the data movement costs performing computation using the memory cells. We evaluate three versions of MATSA to match the requirements of different environments (e.g., embedded, desktop, or HPC computing) based on MRAM technology trends. We perform a design space exploration and demonstrate that our HPC version of MATSA can improve performance by  $7.35\times/6.15\times/6.31\times$  and energy efficiency by  $11.29\times/4.21\times/2.65\times$  over server CPU, GPU and PNM architectures, respectively.

## I. INTRODUCTION

The explosion of the Internet-Of-Things and Big Data era has resulted in the continuous generation of a very large amount of data, which is increasingly difficult to store and analyze [1]. Small sensors and devices produce a significant portion of this data [2], which includes observations (e.g., temperature, voltage, sound) sampled over time. Such a collection of data is also referred to as a time series, a common data representation in almost every scientific discipline and business application [3], for instance, epidemiology, genomics, neuroscience, environmental sciences, and stock markets. Time series analysis (TSA) splits the time series data into *subsequences* of consecutive data points to extract valuable information. This information can be used, for example,

\* Aditya Manglik and Christina Giannoula have equal contribution.

to filter out irrelevant subsequences or to find subsequences of interest. The main goal of this approach is to avoid applying complex and more costly domain-specific algorithms to a large amount of data whenever it is not necessary while accelerating and processing only necessary data. Figure 1 describes the aforementioned process based on an example processing flow.

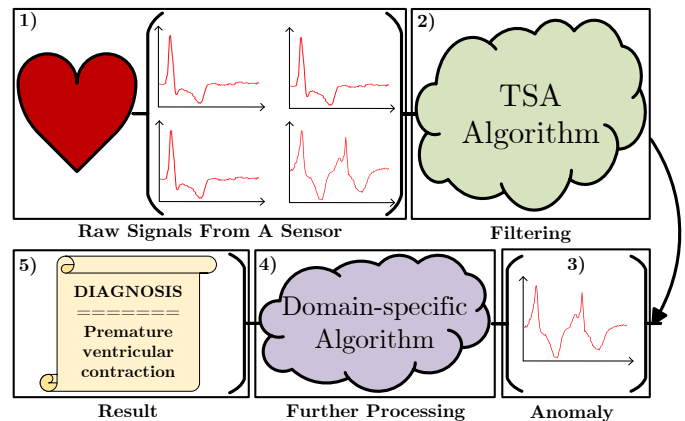


Fig. 1: Example of TSA application. In this processing flow, TSA acts as a filter to avoid most of the computation by selecting the relevant queries (anomalies) and discarding the irrelevant ones (expected behavior).

Assuming that the input is sliced into subsequences, the TSA algorithm can filter out those subsequences that match an expected behavior and runs the computationally expensive domain-specific algorithm (e.g., [4]) only for the anomalies, i.e., only a small amount of critical data that need to be further analyzed. Another example is SquiggleFilter [5], a TSA-based accelerator that filters the Mini-ION sequencer’s output and filters everything except for the sequences of interest, eliminating unnecessary computation that constitutes the 85% of the total workload. However, SquiggleFilter lacks generality to a wide range of application domains.

TSA algorithms define a *distance metric* to find subsequences of interest that can be calculated using different approaches, e.g., by one-to-one or by Dynamic Time Warping (DTW). DTW offers higher precision than one-to-one in common scenarios [6]. Based on the distance metric, the literature classifies the subsequences that have low distance as *motifs* [7] (similarities) and high distance as *discords* [8] (anomalies). Classification is a critical step before further analysis via domain-specific algorithms or human experts.

**Problem.** We characterize the performance of DTW algorithms in commodity hardware platforms in §III-B and find that their performance is bottlenecked by: 1) high data movement costs between memory and processors, and 2) low amount of parallelism provided in commodity hardware. The

DTW algorithm has a large memory footprint and exhibits poor cache locality, resulting in a large number of main memory accesses. Prior work has demonstrated that such data movement significantly impacts the performance and energy efficiency of modern applications [9]. In addition, commodity hardware exposes low levels of parallelism in comparison with a large number of queries at the application level that can be executed in parallel. We observe that DTW is an embarrassingly parallel kernel as each query can be executed as an independent task using simple integer operations.

**Optimization.** Recent works have proposed *Processing-using-Memory* (PuM) to mitigate data movement overheads, which performs operations *in-situ* using the memory cells. Resistive memories are a promising memory technology to implement in-situ computation, where operations are performed based on the voltage difference between activated cells. There are several emerging memory technologies (e.g., ReRAM [10], PCM [11], STT-MRAM [12], and SOT-MRAM [13]). Magneto-Resistive RAM technologies (STT-MRAM and SOT-MRAM) offer low switching energy, non-volatility, superior endurance, high retention time, high integration density, and compatibility with CMOS technology [14]–[16]. Based on these characteristics, we believe that MRAM-based memory substrates offer a promising candidate for designing accelerators for time series analysis while overcoming the bottlenecks in commodity hardware platforms.

**Our Goal.** Our goal in this work is to leverage in-situ computation via MRAM-based crossbars to enable high-performance and energy-efficient time series analysis for a wide range of applications. To this end, 1) we propose an efficient data mapping scheme of sDTW to PuM memories to accelerate its execution using in-situ computation, and 2) we perform a design space exploration in terms of cell latencies and energies based on recent trends of MRAM device characteristics [17], [18]. Our evaluation shows that our HPC version of MATSA can improve performance by  $7.35\times/6.15\times/6.31\times$  and energy efficiency by  $11.29\times/4.21\times/2.65\times$  over server CPU, GPU and PNM architectures, respectively.

This work makes the following *contributions*:

- We characterize subsequence Dynamic Time Warping algorithm using state-of-the-art platforms, demonstrating its performance bottlenecks in commodity platforms and optimization opportunities using a PuM paradigm.
- We propose *MATSA*, the first MRAM-based Accelerator to accelerates TSA via PuM. *MATSA* exploits 1) a data mapping to reduce memory footprint from  $O(MN)$  to  $O(4N)$ , 2) a wavefront scheme to enable per-query parallelism, and 3) the support for computing in the crossbar.
- We perform a design space exploration of *MATSA* and show that it can greatly improve performance and energy efficiency compared to commodity and PNM platforms.

## II. BACKGROUND

### A. Time Series Analysis

A *time series*  $T$  is a sequence of  $n$  data points  $t_i$ , where  $1 \leq i \leq n$ , collected over time. A subsequence of  $T$ , called a *window*, is denoted by  $T_{i,m}$ , where  $i$  is the index of the first

data point, and  $m$  is the number of samples in the subsequence, with  $1 \leq i$ , and  $m \leq n - i$ .

There are two main approaches to perform time series analysis: 1) the self-join, and 2) the query-filtering. In self-join, all sequences of a given time series are compared against the remaining subsequences of the same time series. In contrast, query filtering compares a set of queries against a reference.

Time series analysis algorithms usually define a distance metric to measure the similarity between two subsequences. The state-of-the-art set of tools to perform time series analysis is Matrix Profile [19] (MP). Due to lower computation requirements, prior MP algorithms utilize one-to-one Euclidean Distance as the similarity metric. Recent proposals [20] have started to utilize Dynamic Time Warping (DTW)-based solutions because of higher precision [21]. DTW enables the detection of events of interest in out-of-sync subsequences, e.g., in subsequences that have different sampling rates.

Figure 2 shows the key difference between the one-to-one and the DTW approaches, in which we compare two similar-shape subsequences offset with a fixed constant.

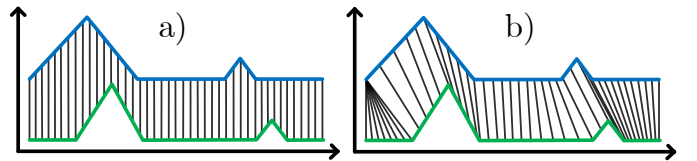


Fig. 2: Example of similarity calculation between two subsequences (blue and green). The one-to-one approach in a) provides a low similarity as it only compares each  $i^{th}$  point of blue with each  $i^{th}$  point of green. In contrast, DTW in b) successfully matches the points of the subsequences.

It can be observed that the DTW algorithm offers better results as it compares a given point with respect to several potential candidates (i.e., determines the best alignment). In contrast, one-to-one executes point-to-point alignment that cannot determine the best alignment in the presence of an offset. One-to-one can be considered as a special case of DTW where the *warping window* is set to '1'. Therefore, we aim to optimize DTW, a more generic and high-precision algorithm, to provide a TSA accelerator for a wide range of applications.

### B. Dynamic Time Warping

The Dynamic Time Warping (DTW) algorithm was first introduced by [22]. The key idea behind DTW is to compute the distance between a certain point in a subsequence and a set of points in another subsequence, considering the minimum distance found. This process is repeated for all the points of the first subsequence, and the sum of all distances provides a similarity measure between the subsequences.

Let's assume that we have two time series subsequences,  $Q$  and  $R$ , of length  $n$  and  $m$ , respectively, where:

$$Q = q_1, q_2, \dots, q_i, \dots, q_n \quad (1)$$

$$R = r_1, r_2, \dots, r_j, \dots, r_m \quad (2)$$

To find how similar those two subsequences are (i.e., find the best alignment between them), DTW constructs an  $n$ -by- $m$  scoring matrix ( $S$ ). We show an example of this matrix in Figure 3 a. Each  $(i^{th}, j^{th})$  cell of the matrix is filled in

two steps. First, the algorithm calculates the distance  $d(q_i, r_j)$  between the two corresponding points of the subsequences. There are several approaches to calculate such distance, while  $d(q_i, c_j) = \text{abs}(q_i - c_j)$  and  $d(q_i, c_j) = (q_i - c_j)^2$  are the most common ones [22]. Second, the distance value is added to the minimum of three neighboring cells, as follows:

$$s_{i,j} = d(q_i, c_j) + \min(s_{i-1,j-1}, s_{i-1,j}, s_{i,j-1}) \quad (3)$$

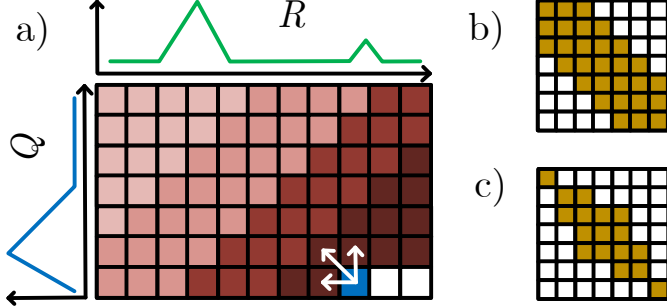


Fig. 3: a) Warping matrix example for a reference time series  $R$  and a query subsequence  $Q$ . The DTW distance between  $R$  and  $Q$  is the minimum value of the last row of the matrix. b) Sakoe-Chiba band. c) Itakura parallelogram.

Once the matrix is computed using dynamic programming, the goal is to find the best alignment (i.e., minimum accumulated cost), known as the *warping path* ( $W$ ). The warping path is a contiguous set of matrix cells that defines the best mapping between  $Q$  and  $R$ , subject to some constraints [22]:

- *Boundary conditions*: the warping path must start and finish in diagonally opposite corner cells of the matrix.
- *Continuity*: the permissible steps in the warping path are restricted to adjacent cells.
- *Monotonicity*: the consecutive points in  $W$  must strictly increase with respect to the time they were collected.

Several DTW implementations try to reduce the computation cost and memory footprint by restricting the DTW score matrix calculation using different approaches, for instance, the Sakoe-Chiba band (Figure 3 b) and the Itakura Parallelogram (Figure 3 c). However, these approaches only apply to scenarios where the end-to-end sequences are quite similar and present small sequence compressions and attenuations (time warps). This fact is important when the goal is to filter out irrelevant subsequences. In this scenario, we usually compare many small queries against a longer reference (set of admissible patterns) and detect the possible anomalies.

**Subsequence Dynamic Time Warping (sDTW).** sDTW is a slightly modified DTW algorithm that allows unbounded subsequence alignment. In the DTW domain, the unbounded alignments are often referred to as *open start* and *open end*, respectively. Using this approach, it is possible to align sequences of different lengths (e.g., a small query over a larger reference). Algorithm 1 presents the pseudocode of sDTW. First, it initializes the matrix  $S$  with zeros. Second, it calculates the distance value of the top-left corner and then the remaining elements of the first row, taking into account the previous values. Third, it fills the remaining elements of the matrix using dynamic programming row by row. Finally, it

returns the minimum element of the last row of the  $S$  matrix, which indicates the similarity between the query and the best alignment with (part of) the reference.

#### Algorithm 1 Subsequence DTW (sDTW)

```

1: procedure sDTW(Q,R)
2:    $S \leftarrow \text{zeros}(N, M)$ ;
3:    $S[0, 0] = \text{dist}(Q[0], R[0])$ ;
4:   for  $i \leftarrow 1$  to  $N$  do
5:      $S[i, 0] \leftarrow S[i-1, 0] + \text{dist}(Q[i], R[0])$ ;
6:   for  $i \leftarrow 1$  to  $N$  do
7:     for  $j \leftarrow 1$  to  $M$  do
8:        $S[i, j] \leftarrow \text{dist}(Q[i], R[j]) +$ 
9:          $\min(S[i-1, j-1], S[i, j-1], S[i-1, j])$ ;
   return  $\min(S[N, :])$ 

```

#### C. MRAM-based PuM Computation

Processing-using-Memory architectures have been proposed to overcome the memory wall challenges in current von-Neumann architectures [23], as they are able to perform computation where data resides. Many prior works demonstrate significant performance and energy efficiency improvements for machine learning workloads using crossbars [24] by exploiting matrix-vector multiplication. Other application domains can exploit bitwise operations [25], [26] to obtain similar benefits. Figure 4 a shows a typical crossbar organization with memory cells connected using bitlines and wordlines.

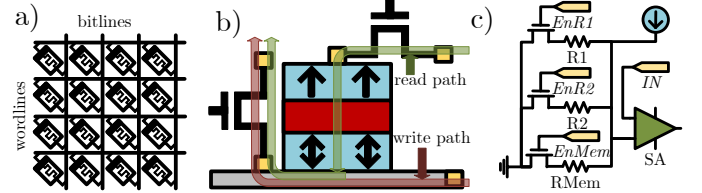


Fig. 4: a) Crossbar organization. b) Magneto-resistive cell. c) Reconfigurable SA that performs in-memory operations based on the voltage across the bitline when two cells are activated.

**Bitwise PuM Mechanism.** Applications based on dynamic programming (e.g., sDTW) require a flexible substrate, thus the matrix-vector approach is not a good fit. Moreover, off-the-shelf crossbar substrates offer limited support for other operations. To overcome this challenge, MAGIC [25], [26] proposes to decompose complex operations into simple Boolean functions (e.g., AND, NOR, XOR) and make them supported in the substrate. The key idea is to vertically map the operands (e.g., 32-bit integers) to the crossbars' columns using (typically) one bit per cell (e.g., each operand value takes 32 bits of a given column). Then, the desired operation (e.g., addition) is decomposed on simple bitwise ones (e.g., NOR) and performed bit by bit by sequentially activating one cell of each operand at the same time. This provokes a difference in the voltage over the bitline depending on the content of the activated cells. Then, a modified sense amplifier generates the result based on that voltage difference and thresholds, storing it in a cell of the same column. While this process is inherently sequential (i.e., the latency is higher than a CMOS-based approach), the 1) independence across columns and 2) the lack of data movement enable an unprecedented parallelism and thus an overall higher throughput. Figure 4 c shows an example of a

modified sense amplifier (SA) with different voltage thresholds and the *regular* memory sensing threshold.

**Technology.** The selection of the cell technology is crucial to ensure feasibility of bitwise based accelerators. While it is possible to build them using DRAM technology (e.g., SIMDRAM [27]), performance scaling is limited by 1) refreshing times and specially by 2) the need of moving data to compute enabled rows. Crossbar-based NVM technologies can overcome those two issues, but it is challenging to support frequent write operations. This is because NVM-based architectures usually suffer from significant latency and energy penalties, and specially low endurance. However, emerging NVM technologies are a possible solution to overcome those drawbacks. Table I presents some NVM characteristics.

Technology	Write/Read Energy (per bit)	Write/Read Time (per bit)	Endurance (Write Cycles)
NAND Flash [28]	470pJ / 46pJ	200 $\mu$ s / 25.2 $\mu$ s	10 <sup>5</sup>
FRAM [29]	1.4nJ / 1.4nJ	120ns / 120ns	10 <sup>15</sup>
STT-MRAM [30]	2nJ / 34pJ	250ns / 10ns	10 <sup>5</sup>
SOT-MRAM [31]	334pJ / 247pJ	1.4ns / 1.1ns	> 10 <sup>15</sup>
ReRAM [32]	1.1nJ / 525fJ	10 $\mu$ s / 5ns	10 <sup>5</sup>
PCM [33]	13.5pJ / 2pJ	150ns / 48ns	10 <sup>7</sup>

TABLE I: Overview of different NVM technologies [34].

To help in the decision of what cell technology use for a potential accelerator, both endurance and latencies have to be taken into consideration. In terms of endurance, we note that both FRAM and SOT-MRAM are good candidates. However, if we add latencies to the equation, we observe that SOT-MRAM is the most promising one, as FRAM latencies are two orders of magnitude larger than the first one. Magneto-Resistive RAM (MRAM)-based substrates (e.g., SOT-MRAM) offer 1) higher density, 2) high endurance, 3) low latencies, and 4) they are CMOS compatible [35]–[39]. Based on that, MRAM technologies are promising architectural candidates for building PuM accelerators [40] where data changes frequently. Figure 4 b shows the basic structure of a Spin-Orbit Torque (SOT)-MRAM device, composed of a stack of Magnetic Tunnel Junctions (MTJs). The key mechanism is based on the change of orientation of one of the layers of the stack, which results in a variation in the device’s electrical resistance. Several companies, like Samsung [41], are raising their interest in these architectures, and we expect that real devices might be commercially available shortly. Therefore, we hypothesize that MRAM-PuM architectures are a good fit to accelerate TSA, and particularly sDTW.

### III. MOTIVATION

In this section we present MATSA’s motivation. Our proposal is motivated by 1) the wide applicability of time series analysis in a wide range of application domains and 2) sDTW’s high data movement overheads in conventional architectures.

#### A. TSA Applications

Time series analysis constitutes one of the most important data mining primitives thanks to its generality in detecting anomalies and similarities for a wide range of applications. Table II presents a few examples for applications of TSA.

In statistics, econometrics, meteorology, and geophysics, the primary goal of time series analysis is prediction and

Field	References	Field	References
Bioinformatics	[42]–[44]	Speech Recognition	[45]
Robotics	[46], [47]	Weather Prediction	[48]
Neuroscience	[42], [49]	Entomology	[50]
Machine Learning	[51]–[53]	Geophysics	[54]–[57]
Econometrics	[58]	Statistics	[59]
Finance	[60], [61]	Control Engineering	[62]–[64]
Signal Processing	[65]	Pattern Recognition	[66]
Communication	[67]–[69]	Medicine	[5], [70]–[72]
Astronomy	[73], [74]	Social Networks	[75], [76]
Clustering	[51], [52]	Classification	[53]
Earthquakes	[54]–[57]	GPS Tracking	[77]
Virtual Reality	[78]	Gesture Recognition	[79], [80]
Trajectories	[81]	Traffic Monitoring	[82]

TABLE II: Time Series Analysis main applications

forecasting. At the same time, in signal processing, control engineering, and communication engineering, it is used for signal detection and estimation. In data mining, pattern recognition, and machine learning, time series motif and discord discovery are used for clustering, classification, anomaly detection, and forecasting. Finally, the most important application of time series motif and discord discovery is clustering seismic data and discovering earthquake pattern clusters from the continuous seismic recording. Consequently, seismic clustering can be applied to earthquake relocation and volcano monitoring to help improve earthquake and volcanic hazard assessments.

**End-to-End Application Benefits of using TSA.** Accelerating TSA is critical for real-life examples of end-to-end applications. For instance, [83] predicts circulatory failure in intensive care units. In this scenario, 90% of the execution time is dominated by TSA preprocessing, while the remainder 10% is used by the machine learning-based application to perform the classification. We also find many other real use case examples that can benefit from this approach, such as:

- **Earthquake detection** [56]. TSA can process the data from a seismograph and detect anomalies that can be further processed with complex algorithms.
- **Electroencephalography** [72]. Assuming an electroencephalograph that is monitoring a patient, TSA can be used to detect anomalies and trigger an alarm about that.
- **Virus Detection** [5]. During the genome assembly process, basecalling is a compute-intensive task that can be skipped for most queries using a TSA-based filter.

#### B. sDTW Bottlenecks

We perform a detailed characterization of parallel versions of sDTW on different commodity computing systems.

**CPU platforms.** We profile the performance of sDTW using a manycore processor (Intel Xeon Phi 7210), analyzing the execution of 16K queries of 8K elements each that are compared against a reference sequence of 32K elements. Figure 5 presents the roofline plot for the experiment.

First, we observe that sDTW only exploits 41% of the system’s integer peak performance, i.e., 59 GINTOPS out of 145 GINTOPS, and it presents low arithmetic intensity (0.55 INTOP/Byte). Second, we find that the memory footprint of the execution is  $\approx$ 570MB. The total memory traffic generated during runtime is  $\approx$ 267 GB, putting the kernel above the DDR4 peak bandwidth most of the execution time. This means

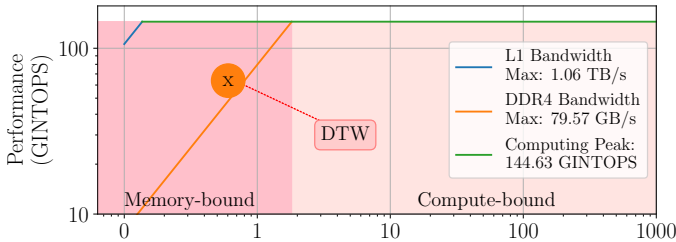


Fig. 5: sDTW’s roofline in an Intel Xeon Phi (64 cores, DDR4). The application is located in the memory-bound zone. that sDTW, while highly parallelizable due to independence across queries, is memory-bound in manycore CPUs.

**GPU platforms.** We find that several prior works [84]–[88] propose to accelerate sDTW using GPUs. However, these works rely on shared-memory optimizations that only perform for certain small query sizes. For large query sizes, these implementations either 1) do not work or 2) use high-latency global-memory to hold the main data structures of sDTW, which results in large performance penalties. We develop a CUDA-based implementation that supports arbitrary sizes.

**FPGA platforms.** Several prior works propose [89]–[91] to accelerate sDTW using FPGAs. However, most of them have very limited onboard memory, and data has to be moved over narrow buses. We develop our own HLS-based FPGA implementation and find that 1) the number of computation units is not enough to exploit the inherent parallelism of sDTW, and 2) the compute units spend most of their time waiting for the memory accesses to be served.

**PNM/PuM platforms.** One way of improving parallelism and reducing data movement costs is to perform computation near or using memory. We analyze several approaches in those directions and detail their benefits and drawbacks.

- *General-purpose PNM.* This approach typically places small CPU cores in the same die as DRAM. The main benefit of this scheme is that the architecture can be potentially used for general-purpose applications. However, this comes at the cost of limited parallelism. We evaluate a general-purpose PNM (upmem baseline) in §V-D, and our evaluations show that this architecture is compute bound when performing the sDTW computation.
- *Specialized PNM.* This approach typically places an ASIC accelerator in the same die as DRAM. The main benefit of this approach is that the processing elements are highly optimized for the target workload. However, the data still needs to be moved from memory to the accelerator. Even though an ASIC accelerator can be used in a PNM architecture, performance would still be bottlenecked by data movement between memory and the accelerator, similarly to general-purpose PNM.
- *SRAM-based PuM.* This approach uses SRAM-based memory arrays to perform in-situ computation (e.g., compute caches [92]). The main benefits of this approach are the high levels of parallelism and reduced data movement. Unfortunately, SRAM suffers from density and scalability issues [93], along with being radiation vulnerable.
- *DRAM-based PuM.* This approach uses DRAM-based memory arrays to perform in-situ computation (e.g., SIM-

DRAM [27]). However, this approach involves internal data movement to perform the operations, as data needs to be moved to specific compute-enabled rows before performing the actual operations. Moreover, DRAM suffers from data volatility and destructive read problems.

### C. Need for an NVM-PuM TSA Accelerator

TSA is widely applicable to many domains [51], while there is *no* previous work to propose a 1) high performance, 2) energy-efficient, and 3) general solution. NVM-PuM architectures are promising candidates to match the requirements.

In contrast to the aforementioned proposals, NVM-PuM architectures 1) enable huge parallelism, 2) eliminate data movement costs, 3) provide better scalability than its competitors thanks to different levels of parallelism, 5) are CMOS-compatible, and 6) overcome the data volatility problem. The main drawback of NVM technologies is the lack of a well-established evaluation methodology. To tackle this issue, we perform an extensive design space exploration study in our evaluation using different NVM technologies.

## IV. MATSA ARCHITECTURE

In this section, we present 1) the mechanism of MATSA, that enables high-performance and energy-efficient TSA, and 2) how MATSA is integrated in commodity systems.

### A. Overview

MATSA is an MRAM-based Accelerator for Time Series Analysis designed to substitute the main processing unit (e.g., CPU or  $\mu$ controller) when performing sDTW. This way, MATSA improves the energy efficiency of the entire system. Figure 6 presents an overview of our proposed architecture. MATSA is composed of one or several chips divided into multiple *banks*. Banks that belong to the same chip share buffers and I/O interfaces and work in a lock-step approach based on an integrated global controller (*Ctrl*) that orchestrates the flow. Each bank is composed of several *Multiple Memory Matrices (MATs)*. The MATs share a *Global Row Buffer (GRB)* and are connected to a *Global Row Decoder (GRD)*. We place a *Local Row Buffer (LRB)* for every pair of subarrays. MATs are composed of several memory *subarrays*.

Each memory subarray is composed of magnetoresistive devices (e.g., SOT-MRAM cells) that are connected to the *Write Word Lines (WWL)*, *Write Bit Lines (WBL)*, *Read Word Lines (RWL)*, *Read Bit Lines (RBL)*, and *Source Lines (SL)*. The memory cells perform the sDTW computation in combination with reconfigurable Sense Amplifiers (SAs). We place a *MATSA controller (Ctrl)* next to each subarray, which is in charge of orchestrating the data flow and activating the subarrays to perform the sDTW computation properly.

### B. MATSA Subarrays

MATSA subarrays are comprised of non-volatile memory cells following a conventional crossbar organization and can work either in regular memory or compute mode. This is a desirable feature since our design consists of 1) subarrays that buffer the data pending to be processed and 2) subarrays that perform the actual computation. Adjacent subarrays are

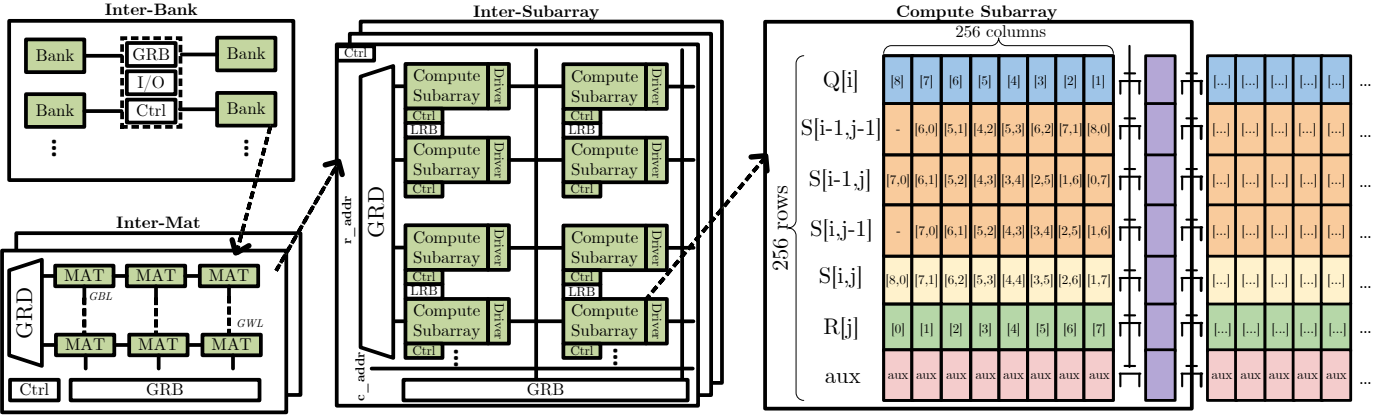


Fig. 6: MATSA’s high-level Architecture and Data Mapping. Note that modifications at *Inter-Bank*, *Inter-Mat*, and *Inter-Subarray* levels with respect to common NVM devices are negligible, easing fabrication and compatibility.

connected using pass gates and auxiliary columns (purple one in Figure 6) to enable the flow at this level of the hierarchy.

**Memory subarrays.** MATSA subarrays in *regular* memory mode support both read and write data operations. To perform a write operation, MATSA first invokes the *Memory Row Decoder* (MRD) to activate the proper WL. Second, the corresponding WD applies the voltage difference needed to switch the cell to the BL and SL (which induces a resistance change in the MJT accordingly). To perform a read operation, MATSA first invokes the MRD to activate the proper WL in the write operation, while the MCD is in charge of connecting the proper BL to its SA. Then, the sense amplifier interprets the voltage value across the path, compares it against a reference voltage, and produces the output result (logic ‘1’ or ‘0’).

**Compute subarrays.** MATSA subarrays working in *compute mode* perform bit-wise operations (e.g., NOR) for operands belonging to the same column. This enables the parallel execution of many operations at the same time since columns in the same subarray work synchronously. The key idea is to select two or three operands simultaneously using the MRD. This produces an equivalent resistance that depends on the content of the selected cells and modifies the sensing voltage across such column accordingly. We include a re-designed reconfigurable sense amplifier (Figure 7) per column, from which MATSA’s Ctrl can select different thresholds depending on the operations. For example, assuming that the desired operation is the majority of three operands, Ctrl activates the three corresponding rows and sets  $EnMaj$  to logic ‘1’. Then, if the equivalent resistance of the activated cells in such column is below  $RMaj$ , the result produced by the SA will be logic ‘1’, or logic ‘0’ otherwise. Most operations are performed in a single cycle thanks to the gates added to SAs.

The execution flow of MATSA is orchestrated by a hierarchy of controllers implemented as finite state machines (FSMs). In particular, MATSA comprises 1) a global controller that orchestrates inter-bank communication, 2) several inter-mat controllers that take care of the inter-mat communications, and 3) subarray controllers that activate the memory rows and drive Reconfigurable SAs according to sDTW’s algorithm.

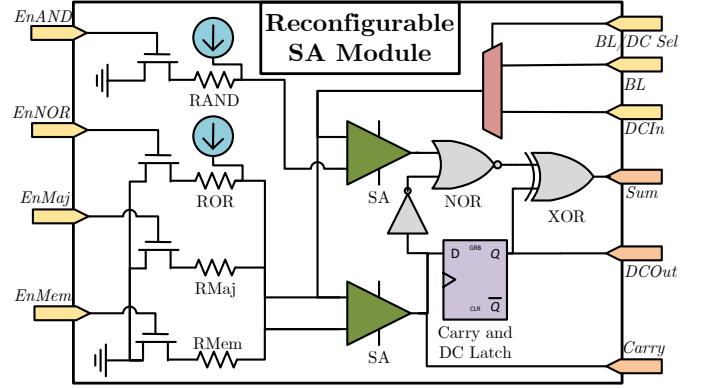


Fig. 7: MATSA’s Reconfigurable Sense Amplifier. Latch register is reused for Addition and Diagonal Copy operations.

### C. sDTW Challenges in NVM-PuM

Using an NVM-PuM architecture to accelerate sDTW raises two main challenges: 1) accessing neighboring columns and 2) column-wise parallelism.

**Accessing neighboring columns.** Each sDTW’s calculated value depends on two cells from the  $j-i$  column ( $S[i-1, j-1]$  and  $S[i, j-i]$ ). We address this issue by introducing a novel single cycle and parallel *diagonal copy operation*.

**Column-wise parallelism.** Ensuring that columns can work independently and at the same time is critical for sDTW performance. We address this issue by introducing a novel data mapping and execution flow based on wavefront execution.

### D. Supported Operations

To support the sDTW algorithm, MATSA implements the following in-memory operations:

- **Row Copy.** The key idea behind this mechanism is to perform consecutive memory read and write operations in the same cycle. In the first half cycle, the corresponding subarray’s MRD activates the source row that is read by the LRB. After that, the data is stored in the destination row in the second half cycle. This mechanism works at MAT and bank levels using the Global Row Buffer (GRB) to accelerate the copy operations across the hierarchy.
- **Diagonal Row Copy.** A diagonal copy occurs when the source and destination cells belong to adjacent columns.

To enable such a feature, MATSA leverages the available registers in the sense amplifiers and the interconnections between them. The operation is performed in two steps. First, each SA reads the source value of its own column. Second, each SA takes the value of its left SA neighbor and writes it to its own column. This operation is performed in parallel when the source and destination rows are common for a given set of cells.

- **Addition/Subtraction.** MATSA implements bit-serial addition/subtraction across columns. It computes bit-by-bit operations from the LSB of the two operands until the MSB. Every bit position takes two memory cycles, further divided into four steps. In the first step, the SAs use the two bit cells activated in the same bit lines as input operands and calculate *Sum* taking into account the stored carry in the register available in the SA. In the second step, the SAs write back the *Sum* value to the corresponding cell. In the third step, the SAs calculate the new *Carry* based on a majority function. In step four, SAs write the result in a reserved cell in their register.
- **Absolute value.** To calculate the absolute value, MATSA first checks the sign bit, leading to two possible scenarios: 1) if the number is positive, no change is needed; otherwise, 2) if the number is negative, MATSA inverts the bits of the number and adds '1' to the result using the addition operation (similar to 2's complement).
- **Minimum value.** To calculate the minimum value between three elements, MATSA performs two comparisons based on the subtraction operation. First, it calculates the difference between two of the numbers. Second, it checks the resulting sign from the previous step and selects one of the two numbers for comparison against the third. The final comparison sign determines the overall minimum.

### E. Data Mapping

We design MATSA's data mapping to leverage the parallel column-wise computation support in MRAM. This organization requires that the operands have to be mapped to the same column. There are three basic structures involved in the sDTW computation: 1) reference, 2) queries, and 3) the warping matrix. The size of the warping matrix can be huge:  $O(NM)$ . However, we are interested in the distance value (i.e., no need to store the alignment), which can be computed by iterating over a single vector that holding the current row of the warping matrix. To do so, we define the *s\_vector*, but we note that each *s\_vector* element (which has to be mapped to different crossbar column) requires accessing previous *s\_vector* values that are mapped to another different columns (i.e.,  $S[i-1, j-1]$ ,  $S[i-1, j]$ ). To overcome this challenge, we propose to add three temporal *s\_vector* in the crossbar array, which are updated accordingly for each step of the computation:  $S[i-1, j-1]$ ,  $S[i-1, j]$  and  $S[i, j-i]$  (see Figure 6). Overall, our optimization reduces the memory footprint from  $O(MN)$  to  $O(4N)$ . Each subarray's column is composed of 256 cells sliced in the following way (e.g., using int32 as datatype):

- **Reference elements ( $R[j]$ ).** We vertically map each reference element to 32 cells of a column. If 1) the

number of available columns is bigger than the number of elements in reference, we replicate the reference to multiple columns to increase parallelism (distributing the queries between them). If 2) the number of available columns is lower than the number of elements in reference, we divide the sDTW matrix and complete processing in successive sequential batches.

- **Query elements ( $Q[i]$ ).** We vertically map each query element to 32 cells of a column. New query elements are introduced on the left side of the crossbar, and they are right-shifted in each successive step (see §IV-F).
- **Current *s\_vector* ( $S[i, j]$ ).** We vertically map each element of the *s\_vector* to 32 cells of a column, being aligned with the reference. The query element is computed in a given step of the algorithm (i and j indexes).
- **Temporal *s\_vectors* ( $S[i-1, j-1]$ ,  $S[i-1, j]$  and  $S[i, j-1]$ ).** We vertically map the three temporal vectors that enable independence to 32 cells of a column, in alignment with the reference and query elements that are computed in a given step of the algorithm.
- **Aux Cells.** Each column is provided with a slice of 64 cells that are used to hold the partial results during the execution flow, as explained in the next subsection.

We note that MATSA's design requires internal data movement. For typical reference sizes, data movement happens at the inter-subarray level. If the reference is big enough, a hierarchy of paths enable its computation efficiently.

### F. Execution Flow

MATSA's execution flow follows a wavefront [94] approach, which reflects the computation pattern in dynamic programming applications. The motivation is that sDTW's matrix has to be computed in the wavefront manner due to inter-cell dependencies. Figure 8 shows how we tackle this challenge. The key idea is to perform the computation diagonally by assigning a diagonal element to each processing element (PE), and use the *diagonal row copy* operation to communicate between processing elements (columns) and get the surrounding values. This approach 1) enables parallelism while computing a given query and 2) creates a pipeline where several queries are being processed. The execution flow iterates over the following steps until queries remain.

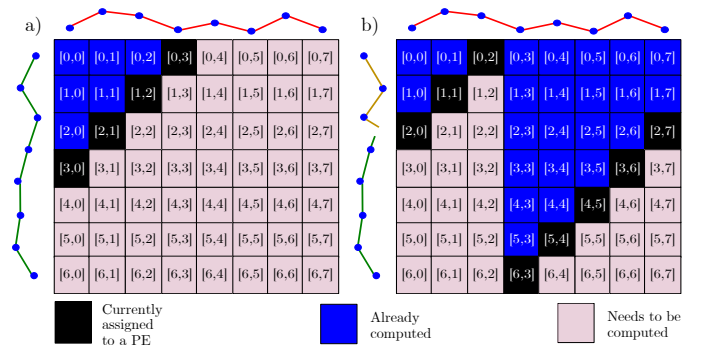


Fig. 8: Wavefront-based sDTW computation. In a, PEs are able to calculate their matrix elements in parallel. In b, the pipeline is full and PEs are also working on different queries.

- 1) **Distance calculation.** Calculation of  $dist(Q[i], R[j])$ , which provides the first partial result  $P1$ . This process implies several substeps depending on the selected distance metric, (e.g., subtraction  $\rightarrow$  absolute value).
- 2) **Minimum.** Calculation without storing the result of  $min(S[i-1, j-1], S[i-1, j], S[i, j-1])$ , which produces the value for the next step  $S1$ .
- 3) **Addition.** Calculation of the addition between the minimum value selected in the previous step ( $S1$ ) and the partial result  $P1$ .
- 4) **Diagonal Copy.** Copying the  $S[i, j]$  vector into the  $S[i, j-1]$  vector shifted by one to the right.
- 5) **Diagonal Copy.** Copying the  $S[i-1, j]$  vector into the  $S[i-1, j-1]$  vector shifted by one to the right.
- 6) **Vertical Copy.** Copying the  $S[i, j]$  vector into the  $S[i-1, j]$  vector.
- 7) **Diagonal Copy.** Copying the  $Q[i]$  vector into the same  $Q[i]$  vector but shifted one position to the right.

### G. System Integration

**Physical Device.** MATSA is designed to work synergistically with the main processing unit (e.g., CPU) to efficiently accelerate TSA. We propose three different versions of MATSA to meet the requirements of each environment:

- 1) **MATSA-Embedded.** A small chip intended to be integrated in edge devices (e.g., sensors).
- 2) **MATSA-Portable.** A USB-based accelerator intended for use in desktops and laptops computers, similar to Intel Neural Compute Stick [95].
- 3) **MATSA-HPC.** A high-performance PCIe-based accelerator intended to be integrated into servers.

**Programming Interface.** To enable efficient use of MATSA by the programmers, we expose an API (Listing 1) that invokes the accelerator based on the input data in a supported DTYPE (int8, int16, int32, int64, fp32 or fp64), the selected mode (*query\_filtering* or *self\_join*) and the distance metric (*abs\_diff* or *square\_diff*). The user can optionally define an anomaly threshold, which provides a boolean array in return with the positions of the anomalies set to true.

```
void matsa(DTYPE * ref, DTYPE * queries, uint32_t *
ref_size, uint32_t * query_sizes, uint32_t n_queries,
char * mode, char * dist_metric, DTYPE anomaly_thres,
bool * anomalies, DTYPE * distances)
```

Listing 1: MATSA’s host interface function.

## V. EVALUATION

In this section, we present the evaluation of MATSA which comprises 1) a self-characterization and 2) a comparison with different baselines in terms of performance and energy.

### A. Methodology

We compare MATSA against several representative (*simulated* and *real*) hardware platforms for the evaluation. MATSA sizes are decided based on the design exploration of MATSA (§V-C), with the goal of making them competitive with their respective baselines. We leave to architects and system integrators the final decision, taking into consideration the tradeoffs between parallelism and area available.

- **(Simulated) CPU-ARM (*cpuarm*):** A portable-class platform based on a 4-core ARM CPU running at 2.5GHz, 32KB L1 caches and an 8GB LPDDR4 memory.
- **CPU-i7 (*cpui7*):** A desktop-class platform based on an 6-core (12 threads) Intel i7 x86 CPU running at 3.2GHz, 64KB L1 caches, 256KB L2 caches, 12MB L3 cache and a 64GB DDR4 memory.
- **CPU-Xeon (*cpuxeon*):** A server-class platform based on two 18-core (36 threads) Intel Xeon Gold x86 CPUs running at 3GHz, 32KB L1 caches, 1MB L2 caches, a 24.75 MB L3 cache and a 768GB DDR4 memory.
- **GPU (*gpu*):** An NVIDIA Tesla V100-SXM2 board equipped with 32GB of HBM memory.
- **FPGA (*fgpa*):** An Alveo U50 board equipped with 872K LUTs and 8GB on-board HBM memory.
- **UPMEM (*upmem*):** A server-class platform based on PNM-enabled memory that is equipped with 2560 DPUs running at 425MHz [96], [97].
- **MATSA-Embedded (*matsa-embedded*):** An ultra low-power version of our accelerator, consisting of 128 compute-enabled crossbars.
- **MATSA-Portable (*matsa-portable*):** A balanced energy/performance version of our accelerator, consisting of 1024 compute-enabled crossbars.
- **MATSA-HPC (*matsa-hpc*):** A high-performance version of our accelerator, consisting of 4096 compute-enabled crossbars.

**Baselines.** To evaluate our *cpuarm* platform, we use ZSim [98] and Ramulator [99], [100] for performance and McPAT [101] for energy consumption evaluations. For the *cpui7* and *cpuxeon* platforms, we obtain the performance from the average of five executions, and energy consumption using rapl-tools [102]. All the above CPU-based platforms execute the same parallel OpenMP C implementation of sDTW. To evaluate the performance of *upmem*, we implement sDTW for this platform. Further, we replicate the reference time series across all DPUs and distribute the queries among them. For the energy measurements, we use an estimation tool provided by UPMEM [103]. To evaluate the performance of *fgpa*, we develop a custom HLS-based implementation of sDTW comprising six compute units that exploit the HBM bandwidth, and obtain the energy consumption using the *xbutil* tool [104]. Finally, we evaluate the performance of *gpu* by developing a custom CUDA-based implementation of sDTW that exploits the HBM bandwidth via coalescing, and obtain the energy consumption using the *nvidia-smi* tool.

Parameter	Values
Crossbar Size (cells)	256x256
Number of Crossbars	128, 256, 512, 1024, 2048, 4096
Read Latency (ns)	1, 3, 5, 10, 20
Write Latency (ns)	1, 3, 5, 10, 20
Read Energy (pJ)	20, 50, 100
Write Energy (pJ)	30, 70, 400

TABLE III: MATSA design exploration parameters.

**MATSA.** To evaluate the performance and energy of MATSA, we perform a sensitivity analysis of the latency and energy of MRAM devices. We create a custom analytical-



based simulator based on sDTW’s execution flow and its required operations. We feed this simulator with the workload parameters along with MATSA characteristics (memory cell latency/energy, etc.) and obtain the execution time and total energy from it. We present the parameter range we use in our evaluation in Table III, going from conservative to optimistic ones based on the MRAM technology trends [105].

### B. Datasets

First, we characterize MATSA using the dataset sizes listed in Table IV, which are extracted from an electrocardiogram (ECG) signal from the European ST-T Database [106]. We cover a wide combination of reference sizes and query sizes to better understand the tradeoffs in MATSA design.

Parameter	Values
Reference Size	64K, 128K, 256K, 512K
Query Size	4K, 8K, 16K, 32K
Number of Queries	4K, 8K, 16K, 64K

TABLE IV: Workloads used in MATSA characterization.

Second, we perform a comparison of MATSA against baselines using real datasets (int32 data type) listed in Table V.

Time Series	Reference Size	Query Size	Num. Queries
Human [107]	7997	120	128K
Song [108]	20234	200	64K
Penguin [109]	109842	800	32K
Seismology [108]	1727990	64	16K
Power [110]	1754985	1536	16K
ECG [106]	1800000	512	16K

TABLE V: Workloads used in MATSA comparison.

### C. MATSA Characterization

We perform a design space exploration of MATSA taking into consideration performance parameters of the cells (i.e., read/write latencies and energies).

**Read/Write Latencies.** We evaluate how changing the read/write latencies affects the execution time and present the results in Figure 9. We observe that, increasing read latency by  $10\times$  incurs a  $4.7\times$  execution time penalty, while increasing the write latency incurs a  $6.5\times$  penalty.

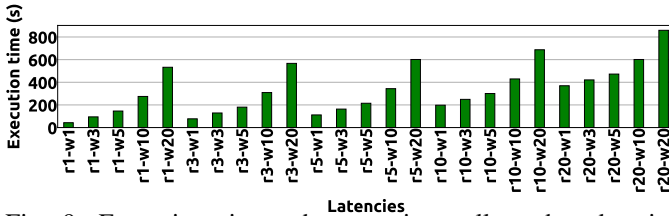


Fig. 9: Execution time when varying cell read and write latencies (ref\_size=128K, query\_size=8K, num\_queries=8K, matsa\_cols=128K).

**Key Observation 1:** using a low write latency memory technology is crucial for MATSA’s design

**Read/Write Energies.** We evaluate how the total execution energy varies with the per word write/read energy, and show the results in Figure 10. We observe here that the contributions of read energy and write energy are similar, thus both of them have to be carefully taken into consideration.

**Key Observation 2:** write energy contributes *only* 19% more than read energy to the total energy

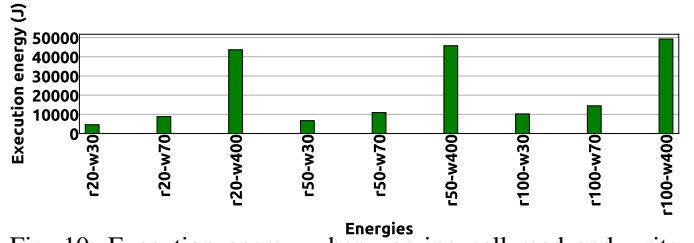


Fig. 10: Execution energy when varying cell read and write energies (ref\_size=128K, query\_size=8K, num\_queries=8K, matsa\_cols=128K).

**Dataset sizes.** First, we evaluate how the execution time varies with different dataset sizes (i.e., ref\_size and query\_size) and present the results in Figure 11. Second, we evaluate how the execution energy varies with different dataset sizes and present the results in Figure 12. We observe that both reference size and query size contribute equally to the execution time and energy. This happens because the total number of operations needed is directly proportional to  $\text{ref\_size} \times \text{query\_size}$ .

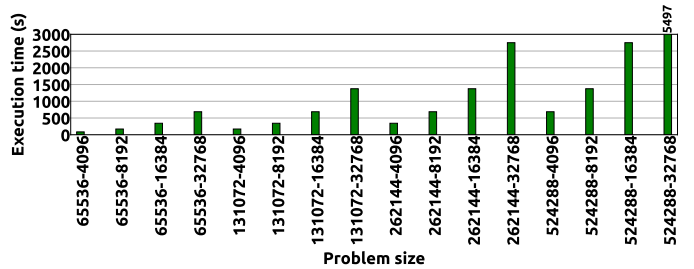


Fig. 11: Execution time when varying dataset sizes (num\_queries=8K, matsa\_cols=128K).

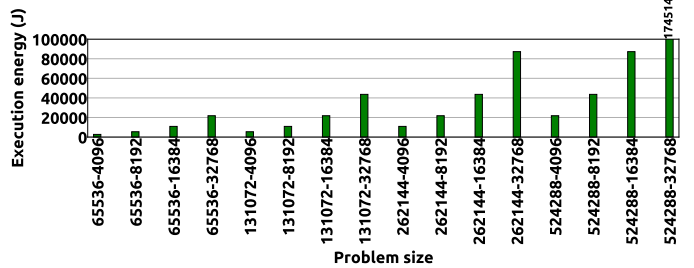


Fig. 12: Execution energy when varying dataset sizes (num\_queries=8K, matsa\_cols=128K).

**Key Observation 3:** total execution time and energy are proportional both to the ref\_size and the query\_size

**MATSA sizes.** We evaluate how the execution time varies when changing the number of MATSA’s compute-enabled columns in Figure 13. MATSA provides almost-ideal scaling.

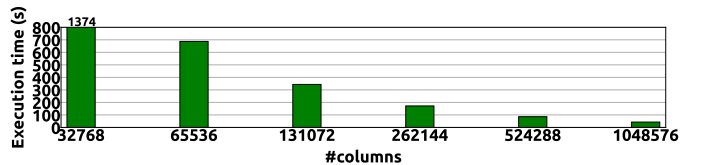


Fig. 13: Execution time when varying MATSA sizes.

**Key Observation 4:** inherent independence across columns enables almost-ideal scaling when increasing # of columns

**Endurance.** Assuming that MATSA is built using 5ns rd/wr cells and runs 24/7 for 10 years, we estimate that (on average) each cell will be written on the order of  $4 \times 10^9$  times. Based on Table I, limited-endurance technologies (e.g. NAND-Flash or ReRAM) would make MATSA to stop working properly after a single day. However, magnetoresistive technologies as SOT-MRAM largely satisfy the endurance requirements and thus are good candidates to build a reliable MATSA accelerator.

**Hardware Overheads.** MATSA introduces hardware overheads in two components: 1) Reconfigurable SAs and 2) MATSA controllers. Reconfigurable SAs add 13 transistors to a traditional SA, thus taking into consideration typical SA and cell areas [111], [112], our design increases the overall crossbar area by less than 1%. MATSA controllers are implemented as small finite-state machines whose area is negligible with respect to the memory arrays.

#### D. MATSA Comparison

In this section, we compare the performance and energy of three versions of MATSA with the corresponding baselines.

**MATSA-Embedded.** We first compare the performance of MATSA-Embedded (32K compute-enabled columns) with `cpuarm`, `cpui7` and `fpga` and present the results in Figure 14. MATSA’s design provides larger parallelism and only implement the required operations, which enables MATSA-Embedded to provide significantly lower execution times than `cpuarm`, `cpui7` and even `fpga`<sup>1</sup>. Second, we compare the energy consumption of MATSA-Embedded with the baselines in Figure 15. We observe that MATSA-Embedded reduces the energy consumption with respect to the other platforms, explained by the lack of off-chip data movement in MATSA.

**MATSA-Portable.** We first compare the performance of MATSA-Portable (256K compute-enabled columns) with `cpuarm`, `cpui7` and `fpga` and show the results in Figure 14. We observe that MATSA-Portable provides even lower execution times than MATSA-Embedded, thanks to the increased parallelism by the number of columns. Second, we compare the energy consumption of MATSA-Portable with the other platforms in Figure 15. We observe a similar energy reduction as MATSA-Embedded, (execution time is decreased in proportion to the number of columns, but instantaneous energy is increased by a similar proportion).

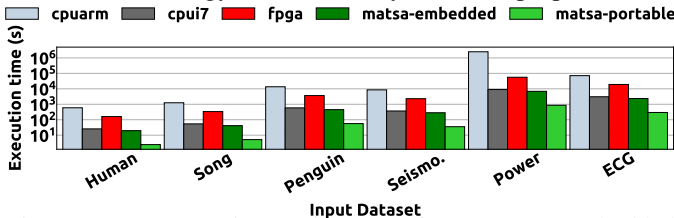


Fig. 14: Execution times of MATSA-Embedded (num\_cols=32K) and MATSA-Portable (num\_cols=256K) versus baselines (rd\_latency=5ns, wr\_latency=10ns).

<sup>1</sup>The reader may wonder how the performance of the `cpui7` baseline can be better than the `fpga` one. The main reason for that is the cache memory in the CPU baseline, which is large enough to hold the main data structures of the algorithm, whereas the FPGA’s scratchpad memory is smaller and the data structures have to be HBM located. Thus, HBM latency lags the FPGA behind the CPU baseline.

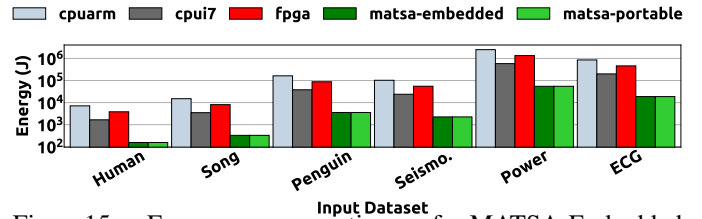


Fig. 15: Energy consumption of MATSA-Embedded (num\_cols=32K) and MATSA-Portable (num\_cols=256K) versus baselines (rd\_energy=50nJ, wr\_energy=70nJ)

**MATSA-HPC.** We first compare the performance of MATSA-HPC with `cpuxeon`, `gpu` and `upmem` and present the results in Figure 16. We observe that MATSA-HPC outperforms all of them thanks to its huge parallelism (1M columns). Second, we compare the energy consumption of MATSA-HPC with the baselines in Figure 17. We note that MATSA-HPC is the most energy-efficient for the same reason as MATSA-Embedded and MATSA-Portable are. We note that `cpuxeon` is bottlenecked by 1) the limited parallelism (# of cores) and 2) memory hierarchy. Our `gpu` baseline provides high parallelism, but it is limited by data movement from and to memory. The PNM-based `upmem` baseline provides high parallelism and closeness to memory, but it is compute-bound for this application, as it is based on general purpose cores.

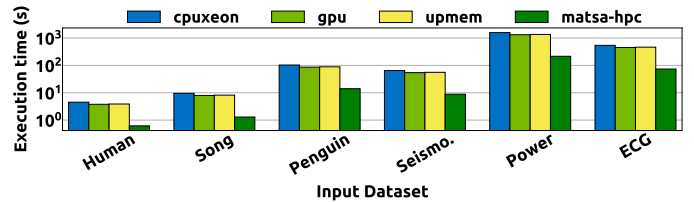


Fig. 16: Execution times of MATSA-HPC versus baselines (rd\_latency=5ns, wr\_latency=10ns, num\_cols=1M)

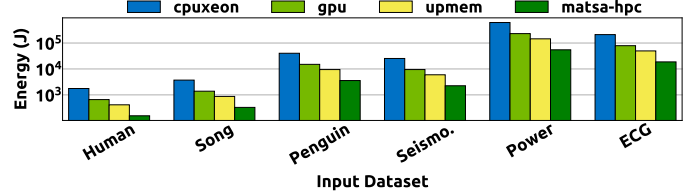


Fig. 17: Energy consumption of MATSA-HPC versus baselines (rd\_energy=50nJ, wr\_energy=70nJ, num\_cols=1M)

**Key Observation 5:** MATSA-Embedded, MATSA-Portable and MATSA-HPC reduce the energy consumption in a similar way with respect to a given baseline. This fact is explained because the main difference between the three MATSA versions is the increment in the number of columns (i.e., SIMD lanes). As a result, while execution time is reduced, the instantaneous power is incremented by the same factor, thus total energy consumed is similar.

**Overall MATSA Benefits.** Finally, we present a summary of MATSA speedups and energy savings compared to baselines in Table VI. We observe that MATSA benefits are larger in embedded environments, where reduced power consumption and dissipation are of great interest.

## VI. DISCUSSION

**End-To-End benefits of MATSA.** MATSA is a general

MATSA Version	Baseline	Speedup	Energy Savings
Embedded	cpuarm	30.20×	45.67×
Portable	cpui7	10.41×	10.65×
	fpga	65.01×	24.58×
HPC	Xeon	7.35×	11.29×
	UPMEM	6.31×	2.65×
	GPU	6.15×	4.21×

TABLE VI: MATSA’s Speedup and Energy over baselines.

TSA accelerator based on sDTW, thus it does not infer domain-specific information from the data that is being analyzed. Based on that, MATSA is likely to be used in a dataflow where our accelerator feeds a domain-specific algorithm with the queries of interest. Based on real examples [5], [83], we conservatively assume that the sDTW’s execution will take at least 75% of the execution time. Thus, we expect the end-to-end speedup of MATSA-HPC to be  $5.51 \times / 4.61 \times / 4.73 \times$  over server CPU, GPU and PNM architectures, respectively.

**Generality of MATSA.** While MATSA is tailored for the important and widely-applicable sDTW, the techniques that we propose are potentially a good fit for other DP algorithms with -mostly minor- modifications. However, we leave the experimentation of other DP algorithms as future work.

## VII. RELATED WORK

**Dynamic Time Warping.** Dynamic time warping is a widely-used method to find patterns in time series [22]. Over the last decade, there has been interest in accelerating DTW via software and hardware-based approaches [113], [114].

Several works explore many-core and GPU architectures to accelerate DTW [114]–[117], while FPGA and ASIC-based solutions can also provide significant performance improvements [114]. Li and Tai [118] explore the computational intensive characteristics of DTW. These characteristics make DTW unsuitable for real-time series analysis using commodity systems. This motivates prior works to propose a FPGA-based single-element processing unit as a building block for DTW time series mining. Loftian et al. [119] propose a wearable ASIC for human activity recognition. They propose a battery-less solution by leveraging a DTW algorithm based on a block-tiered architecture. The low-power DTW block (with lower sampling frequency/bit resolution) is placed in the first stage for filtering majority of the non-target movements. Power-consuming DTW blocks are placed in the later stages, taking advantage of power gating. Koul et al. [120] focus on the multiplication operation, which constitutes the largest computational portion of the DTW algorithm. They compare the performance of different FPGA multiplier designs. Chen and Gu [121] propose an accelerator that exploits DTW pipelining using a specialized designed time flip-flop. Xu et al. [122] propose a memristor-based DTW accelerator for real-time time series mining on data centers. Compared to our proposal, even though they also use memristors for computation, they do not leverage on processing-using-memory. The time series samples must be moved in and out the memristor-based accelerator to compute the DTW distance (i.e., the memristor devices do not store data). To the best of our knowledge, this is the first work that analyzes the memory-bound nature of DTW time series mining and proposes a PuM solution to accelerate it.

**Processing using Memory.** There has been significant interest in Processing-using-Memory-based solutions for overcoming the von Neumann bottleneck in modern computation platforms. AlignS [123] is an energy-efficient and parallel PuM accelerator which implements a DNA short read alignment algorithm. DRISA [124] is a DRAM-based accelerator and shows important benefits when running CNNs. MRIMA [125] proposes hardware-friendly bit-line computing in STT-MRAM arrays, with reconfigurable sense amplifiers implementing complete Boolean logic functions in one cycle. MRIMA proposal improves energy and performance over ASIC counterparts when accelerating Convolutional Neural Networks (CNNs). CMP-PIM [126] uses SOT-MRAM with modified sense amplifiers, as well as a Digital Processing Unit (DPU) with three ancillary units to accelerate CNNs. PIM-Aligner [127] approaches the DNA short read alignment problem such as AlignS does. However, they use three SAs per bit-line to perform computations in a single cycle. In contrast, the AlignS architecture has two SAs and a two-cycle addition scheme. However, none of these prior works explore the widely-used DTW kernel in PuM.

**Processing Near Memory.** Recent works explore Near Data Processing [23], [103], [103], [126], [128], [129], [129], [130], [130], [131], [131], [132], [132], [133], [133]–[265], [265], [266], [266], [267], [267], [268], [268] for various applications using accelerators or general-purpose cores. In [269], ARM cores are used as NDP compute units to improve data analytics operators (e.g., group, join, sort). IMPICA [270] is an NDP pointer chasing accelerator. Tesseract [271] is a scalable NDP accelerator for parallel graph processing. TETRIS [207] is an NDP neural network accelerator. Lee et al. [272] propose an NDP accelerator for similarity search. GRIM-Filter [273] is an NDP accelerator for pre-alignment filtering [274]–[278] in genome analysis [279]. Boroumand et al. [9] analyze the energy and performance impact of data movement for several widely-used Google consumer workloads, providing NDP accelerators for them. CoNDA [129] provides efficient cache coherence support for NDP accelerators. SparseP [280], [281] provides efficient data partitioning/mapping techniques of the SpMV kernel tailored for near-bank NDP architectures. Finally, an NDP architecture [203] has been proposed for MapReduce-style applications. However, none of these works explores DTW, while they focus on DRAM-based NDP.

## VIII. CONCLUSIONS

This paper presents MATSA, the first MRAM-based Accelerator for Time Series Analysis. The key idea is to exploit magnetoresistive crossbars to enable energy efficient and fast time series computation in-memory. MATSA provides the following key benefits: 1) significantly increases parallelism exploiting column-level bitwise operations, and 2) reduces the overheads of data movement by performing computation in the memory cells. We evaluate three versions of MATSA to match the requirements of different environments and perform a design space exploration. MATSA significantly improves performance and energy over commodity and PNM platforms.

## REFERENCES

- [1] Heiner Lasi, Peter Fettke, Hans-Georg Kemper, Thomas Feld, and Michael Hoffmann. Industry 4.0. *Business & Information Systems Engineering*, 6(4):239–242, 2014.
- [2] Feichi Zhou and Yang Chai. Near-sensor and in-sensor Computing. *Nature Electronics*, 3(11):664–671, 2020.
- [3] Hairong Song. Review of time series analysis and its applications with r examples , by robert h. shumway & david s. stoffer: New york, ny: Springer, 2011, 596 pp., 2017.
- [4] Nguyen Cong Long, Phayung Meesad, and Herwig Unger. A highly accurate firefly based algorithm for heart disease prediction. *Expert Systems with Applications*, 42(21):8221–8231, 2015.
- [5] Tim Dunn, Harisankar Sadasivan, Jack Wadden, Kush Goliya, Kuan-Yu Chen, David Blaauw, Reetuparna Das, and Satish Narayanasamy. Squigglefilter: An accelerator for portable virus detection. In *MICRO*, pages 535–549, 2021.
- [6] Eamonn Keogh and Chotirat Ann Ratanamahatana. Exact indexing of dynamic time warping. *Knowledge and information systems*, 7(3):358–386, 2005.
- [7] Pranav Patel, Eamonn Keogh, Jessica Lin, and Stefano Lonardi. Mining Motifs in Massive Time Series Databases. In *ICDM*, 2002.
- [8] Eamonn Keogh, Jessica Lin, Sang-Hee Lee, and Helga Van Herle. Finding the most unusual time series subsequence: algorithms and applications. *Knowledge and Information Systems*, 11(1):1–27, 2007.
- [9] Amirali Boroumand, Saugata Ghose, Youngsok Kim, Rachata Ausavarungnirun, Eric Shiu, Rahul Thakur, Daehyun Kim, Aki Kususela, Allan Knies, Parthasarathy Ranganathan, and Onur Mutlu. Google Workloads for Consumer Devices: Mitigating Data Movement Bottlenecks. *ASPLOS*, 2018.
- [10] H-S Philip Wong, Heng-Yuan Lee, Shimeng Yu, Yu-Sheng Chen, Yi Wu, Pang-Shiu Chen, Byoungil Lee, Frederick T Chen, and Ming-Jinn Tsai. Metal-oxide RRAM. *Proceedings of the IEEE*, 2012.
- [11] H-S Philip Wong, Simone Raoux, SangBum Kim, Jiale Liang, John P Reifenberg, Bipin Rajendran, Mehdi Asheghi, and Kenneth E Goodson. Phase change memory. *Proceedings of the IEEE*, 98(12):2201–2227, 2010.
- [12] Dmytro Apalkov, Alexey Khvalkovskiy, Steven Watts, Vladimir Nikitin, Xueti Tang, Daniel Lottis, Kiseok Moon, Xiao Luo, Eugene Chen, Adrian Ong, et al. Spin-transfer torque magnetic random access memory (stt-mram). *ACM Journal on Emerging Technologies in Computing Systems (JETC)*, 9(2):1–35, 2013.
- [13] Fabian Oboril, Rajendra Bishnoi, Mojtaba Ebrahimi, and Mehdi B Tahoori. Evaluation of hybrid memory technologies using sot-mram for on-chip cache hierarchy. *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems*, 34(3):367–380, 2015.
- [14] Tetsuo Endoh, Hiroaki Honjo, Koichi Nishioka, and Shoji Ikeda. Recent progresses in stt-mram and sot-mram for next generation mram. In *2020 IEEE Symposium on VLSI Technology*, pages 1–2. IEEE, 2020.
- [15] Guillaume Prenat, Kotb Jabeur, Gregory Di Pendina, Olivier Bouille, and Gilles Gaudin. Beyond stt-mram, spin orbit torque ram sot-mram for high speed and high reliability applications. In *Spintronics-based Computing*, pages 145–157. Springer, 2015.
- [16] Guillaume Prenat, Kotb Jabeur, Pierre Vanhauwaert, Gregory Di Pendina, Fabian Oboril, Rajendra Bishnoi, Mojtaba Ebrahimi, Nathalie Lamard, Olivier Bouille, Kevin Garello, et al. Ultra-fast and high-reliability sot-mram: From cache replacement to normally-off computing. *IEEE Transactions on Multi-Scale Computing Systems*, 2(1):49–60, 2015.
- [17] Shimeng Yu and Pai-Yu Chen. Emerging memory technologies: Recent trends and prospects. *IEEE Solid-State Circuits Magazine*, 8(2):43–56, 2016.
- [18] William J Gallagher, Eric Chien, Tien-Wei Chiang, Jian-Cheng Huang, Meng-Chun Shih, CY Wang, Christine Bair, George Lee, Yi-Chun Shih, Chia-Fu Lee, et al. Recent progress and next directions for embedded mram technology. In *2019 Symposium on VLSI Circuits*, pages T190–T191. IEEE, 2019.
- [19] Chin-Chia Michael Yeh, Yan Zhu, Liudmila Ulanova, Nurjahan Begum, Yifei Ding, Hoang Anh Dau, Diego Furtado Silva, Abdullah Mueen, and Eamonn Keogh. Matrix Profile I: All Pairs Similarity Joins for Time Series: A Unifying View That Includes Motifs, Discords and Shapelets. In *ICDM*, 2016.
- [20] Sara Alaei, Ryan Mercer, Kaveh Kamgar, and Eamonn Keogh. Time series motifs discovery under dtw allows more robust discovery of conserved structure. *Data Mining and Knowledge Discovery*, 35(3):863–910, 2021.
- [21] Chotirat Ann Ratanamahatana and Eamonn Keogh. Making time-series classification more accurate using learned constraints. In *Proceedings of the 2004 SIAM international conference on data mining*, pages 11–22. SIAM, 2004.
- [22] Donald J Berndt and James Clifford. Using dynamic time warping to find patterns in time series. In *KDD workshop*, volume 10, pages 359–370. Seattle, WA, USA:, 1994.
- [23] Onur Mutlu, Saugata Ghose, Juan Gómez-Luna, and Rachata Ausavarungnirun. Processing Data Where it Makes Sense: Enabling In-Memory Computation. *Microprocessors and Microsystems*, 2019.
- [24] Sparsh Mittal. A survey of reram-based architectures for processing-in-memory and neural networks. *Machine learning and knowledge extraction*, 1(1):75–114, 2018.
- [25] Jeffrey Louis, Barak Hoffer, and Shahar Kvatinisky. Performing memristor-aided logic (magic) using stt-mram. In *2019 26th IEEE International Conference on Electronics, Circuits and Systems (ICECS)*, pages 787–790. IEEE, 2019.
- [26] Shahar Kvatinisky, Dmitry Belousov, Slavik Liman, Guy Satat, Nimrod Wald, Eby G Friedman, Avinoam Kolodny, and Uri C Weiser. Magic—memristor-aided logic. *IEEE Transactions on Circuits and Systems II: Express Briefs*, 61(11):895–899, 2014.
- [27] Nastaran Hajinazar, Geraldo F Oliveira, Sven Gregorio, João Dinis Ferreira, Nika Mansouri Ghiasi, Minesh Patel, Mohammed Alser, Saugata Ghose, Juan Gómez-Luna, and Onur Mutlu. SimDRAM: a framework for bit-serial simDRAM processing using dram. In *Proceedings of the 26th ACM International Conference on Architectural Support for Programming Languages and Operating Systems*, pages 329–345, 2021.
- [28] Laura M Grupp, Adrian M Caulfield, Joel Coburn, Steven Swanson, Eitan Yaakobi, Paul H Siegel, and Jack K Wolf. Characterizing flash memory: Anomalies, observations, and applications. In *Proceedings of the 42nd Annual IEEE/ACM International Symposium on Microarchitecture*, pages 24–33, 2009.
- [29] Texas Instruments. FRAM – New Generation of Non-Volatile Memory. [https://www.antaos.fr/IMG/pdf/web\\_site\\_sot\\_whitepaper.pdf](https://www.antaos.fr/IMG/pdf/web_site_sot_whitepaper.pdf). 2009.
- [30] Yu-Der Chih, Yi-Chun Shih, Chia-Fu Lee, Yen-An Chang, Po-Hao Lee, Hon-Jarn Lin, Yu-Lin Chen, Chieh-Pu Lo, Meng-Chun Shih, Kuei-Hung Shen, et al. 13.3 a 22nm 32mb embedded stt-mram with 10ns read speed, 1m cycle write endurance, 10 years retention at 150 c and high immunity to magnetic field interference. In *2020 IEEE International Solid-State Circuits Conference (ISSCC)*, pages 222–224. IEEE, 2020.
- [31] Antaios. Spin-Orbit Torque MRAM. Technical Report. Antaios, 51 Avenue Jean Kuntzmann, 38830 Montbonnot – France. 3 pages. [https://www.antaos.fr/IMG/pdf/web\\_site\\_sot\\_whitepaper.pdf](https://www.antaos.fr/IMG/pdf/web_site_sot_whitepaper.pdf). 2020.
- [32] Pulkit Jain, Umut Arslan, Meenakshi Sekhar, Blake C Lin, Liqiong Wei, Tanaya Sahu, Juan Alzate-Vinasco, Ajay Vangapaty, Mesut Metereliyoz, Nathan Strutt, et al. 13.2 a 3.6 mb 10.1 mb/mm<sup>2</sup> embedded non-volatile reram macro in 22nm finfet technology with adaptive forming/set/reset schemes yielding down to 0.5 v with sensing time of 5ns at 0.7 v. In *2019 IEEE International Solid-State Circuits Conference (ISSCC)*, pages 212–214. IEEE, 2019.
- [33] Benjamin C Lee, Engin Ipek, Onur Mutlu, and Doug Burger. Architecting phase change memory as a scalable dram alternative. In *Proceedings of the 36th annual international symposium on Computer architecture*, pages 2–13, 2009.
- [34] Tim Daulby, Anand Savanth, Alex S Weddell, and Geoff V Merrett. Comparing nvm technologies through the lens of intermittent computation. In *Proceedings of the 8th International Workshop on Energy Harvesting and Energy-Neutral Sensing Systems*, pages 77–78, 2020.
- [35] Alessandro Grossi, Cristian Zambelli, Piero Olivo, Paolo Pellati, Michele Ramponi, Christian Wenger, Jeremy Alvarez-Herault, and Ken Mackay. An automated test equipment for characterization of emerging mram and rram arrays. *IEEE Transactions on Emerging Topics in Computing*, 6(2):269–277, 2016.
- [36] H-S Philip Wong and Sayeef Salahuddin. Memory leads the way to better computing. *Nature nanotechnology*, 10(3):191–194, 2015.
- [37] Haitong Li, Mudit Bhargava, Paul N Whatmough, and H-S Philip Wong. On-chip memory technology design space explorations for mobile deep neural network accelerators. In *2019 56th ACM/IEEE design automation conference (DAC)*, pages 1–6. IEEE, 2019.
- [38] Mohamed M Sabry Aly, Mingyu Gao, Gage Hills, Chi-Shuen Lee, Greg Pitner, Max M Shulaker, Tony F Wu, Mehdi Asheghi, Jeff Bokor, Franz Franchetti, et al. Energy-efficient abundant-data computing: The n3xt 1,000 x. *Computer*, 48(12):24–33, 2015.

- [39] Yiming Huai, Yuchen Zhou, Ioan Tudosa, Roger Malmhall, Rajiv Ranjan, and Jing Zhang. Progress and outlook for stt-mram. In *2011 IEEE/ACM International Conference on Computer-Aided Design (ICCAD)*, pages 235–235. IEEE, 2011.
- [40] Hao Yan, Hebin R. Cherian, Ethan C. Ahn, Xuehai Qian, and Lide Duan. icelia: A full-stack framework for stt-mram-based deep learning acceleration. *IEEE Transactions on Parallel and Distributed Systems*, 31(2):408–422, 2020.
- [41] Seungchul Jung, Hyungwoo Lee, Sungmeen Myung, Hyunsoo Kim, Seung Keun Yoon, Soon-Wan Kwon, Yongmin Ju, Minje Kim, Wooseok Yi, Shinhee Han, et al. A crossbar array of magnetoresistive memory devices for in-memory computing. *Nature*, 601(7892):211–216, 2022.
- [42] André E X Brown, Eviatar Yemini, Laura J Grundy, Tadas Jucikas, and William Schafer. A dictionary of behavioral motifs reveals clusters of genes affecting *Caenorhabditis elegans* locomotion. *Proceedings of the National Academy of Sciences of the United States of America*, 110, 2012.
- [43] Martin Straume. Dna microarray time series analysis: automated statistical assessment of circadian rhythms in gene expression patterning. In *Methods in enzymology*, volume 383, pages 149–166. Elsevier, 2004.
- [44] Ziv Bar-Joseph. Analyzing time series gene expression data. *Bioinformatics*, 20(16):2493–2503, 2004.
- [45] Arvind Balasubramanian, Jun Wang, and Balakrishnan Prabhakaran. Discovering Multidimensional Motifs in Physiological Signals for Personalized Healthcare. *JSTSP*, 2016.
- [46] Yoshiki Tanaka, Kazuhisa Iwamoto, and Kuniaki Uehara. Discovery of Time-Series Motif from MultiDimensional Data Based on MDL Principle. *Machine Learning*, 58:269–300, 2005.
- [47] Alireza Vahdatpour, Navid Amini, and Majid Sarrafzadeh. Toward unsupervised activity discovery using multi-dimensional motif detection in time series. In *Int'l. Jont Conf. on Artificial Intelligence*, pages 1261–1266, 2009.
- [48] Amy McGovern, Derek H. Rosendahl, Rodger A. Brown, and Kelvin K. Droegemeier. Identifying predictive multi-dimensional time series motifs: an application to severe weather prediction. *Data Mining and Knowledge Discovery*, 22(1):232–258, 2011.
- [49] Ilya Kolb, Giovanni Talei Franzesi, Michael Wang, Suhasa B. Kandaramaiah, Craig R. Forest, Edward S. Boyden, and Annabelle C. Singer. Evidence for long-timescale patterns of synaptic inputs in ca1 of awake behaving mice. *Journal of Neuroscience*, 38(7):1821–1834, 2018.
- [50] Balázs Szigeti, Ajinkya Deogade, and Barbara Webb. Searching for motifs in the behaviour of larval *drosophila melanogaster* and *caenorhabditis elegans* reveals continuity between behavioural states. *Journal of The Royal Society Interface*, 12(113):20150899, 2015.
- [51] T. Warren Liao. Clustering of time series data - A survey. *Pattern Recogn.*, 38(11):1857–1874, 2005.
- [52] Saeed Aghabozorgi, Ali Seyed Shirkhorshidi, and Teh Ying Wah. Time-series clustering - A decade review. *Inf. Syst.*, 53(C):16–38, 2015.
- [53] Eamonn Keogh and Shrutli Kasetty. On the need for time series data mining benchmarks: A survey and empirical demonstration. *Data Min. Knowl. Discov.*, 7(4):349–371, 2003.
- [54] Daniel T. Trugman and Peter M Shearer. GrowClust: A hierarchical clustering algorithm for relative earthquake relocation, with application to the Spanish Springs and Sheldon, Nevada, earthquake sequences. *Seismological Research Letters*, 88(2A), 2017.
- [55] A double-difference earthquake location algorithm: Method and application to the northern Hayward fault. *Bull. Seism. Soc. Am.*, 2000.
- [56] Annemarie Christophersen, Natalia I. Deligne, Anca M. Hanea, Lauriane Chardot, Nicolas Fournier, and Willy P. Aspinall. Bayesian network modeling and expert elicitation for probabilistic eruption forecasting: Pilot study for Whakaari/White Island, New Zealand. *Frontiers in Earth Science*, 6:211, 2018.
- [57] Chris McKee, Ima Itikarai, and Hugh Davies. Instrumental volcano surveillance and community awareness in the lead-up to the 1994 eruptions at rabaul, papua new guinea. In *Observing the Volcano World: Volcano Crisis Communication*, pages 205–233. Springer International Publishing, 2018.
- [58] E. Philip Howrey. The role of time series analysis in econometric model evaluation. *Evaluation of Econometric Models*, pages 275–307, 1980.
- [59] R.H. Shumway. Applied statistical time series analysis. *Prentice-Hall, Englewood Cliffs.*, 1988.
- [60] Ruey S. Tsay. *Analysis of financial time series*. Wiley series in probability and statistics. Wiley-Interscience, 2005.
- [61] Kei Nakagawa, Mitsuyoshi Imamura, and Kenichi Yoshida. Stock price prediction using k-medoids clustering with indexing dynamic time warping. *Electronics and Communications in Japan*, 2019.
- [62] George E. P. Box, Gwilym M. Jenkins, Gregory C. Reinsel, and Greta M. Ljung. *Time Series Analysis: Forecasting and Control*. Wiley-Interscience, 2015.
- [63] Alex T Nelson and Eric A Wan. A two-observation kalman framework for maximum-likelihood modeling of noisy time series. In *1998 IEEE International Joint Conference on Neural Networks Proceedings. IEEE World Congress on Computational Intelligence (Cat. No. 98CH36227)*, volume 3, pages 2489–2494. IEEE, 1998.
- [64] David Barber, A Taylan Cemgil, and Silvia Chiappa. *Bayesian time series models*. Cambridge University Press, 2011.
- [65] S. A. P. Kumar and P. K. Bora. Time series analysis and signal processing. In *Conf. on Computational Intelligence and Signal Processing*, pages 24–24, 2012.
- [66] Berlin Wu. Pattern recognition and classification in time series analysis. *Applied Mathematics and Computation*, 62(1):29 – 45, 1994.
- [67] Anukool Lakhina, Mark Crovella, and Christiphe Diot. Characterization of network-wide anomalies in traffic flows. In *IMC*, pages 201–206. ACM, 2004.
- [68] Raj Jain and Shawn Routhier. Packet trains—measurements and a new model for computer network traffic. *IEEE journal on selected areas in Communications*, 4(6):986–995, 1986.
- [69] Paul Barford and David Plonka. Characteristics of network traffic flow anomalies. In *Internet Measurement Workshop*, pages 69–73. Citeseer, 2001.
- [70] Lal Hussain, Wajid Aziz, Jalal S. Alowibdi, Nazneen Habib, Muhammad Rafique, Sharjil Saeed, and Syed Zaki Hassan Kazmi. Symbolic time series analysis of electroencephalographic (EEG) epileptic seizure and brain dynamics with eye-open and eye-closed subjects during resting states. *Journal of Physiological Anthropology*, 36(1), 2017.
- [71] T. Balli and R. Palaniappan. EEG time series analysis with exponential autoregressive modelling. In *Canadian Conf. on Electrical and Computer Engineering*, 2008.
- [72] Guangyan Chen, Guoliang Lu, Zhaohong Xie, and Wei Shang. Anomaly detection in eeg signals: a case study on similarity measure. *Computational intelligence and neuroscience*, 2020.
- [73] R Vio, Niels Kristensen, Henrik Madsen, and W Wamsteker. Time series analysis in astronomy: limits and potentialities. *Astronomy and Astrophysics*, 435, 10 2004.
- [74] Jeffrey D. Scargle. Studies in astronomical time series analysis. V. Bayesian blocks, a new method to analyze structure in photon counting data. *The Astrophysical Journal*, 504(1):405–418, 1998.
- [75] Kajal Nusratullah, Shoab Ahmad Khan, Asadullah Shah, and Wasi Haider Butt. Detecting changes in context using time series analysis of social network. In *2015 SAI Intelligent Systems Conference (IntelliSys)*, pages 996–1001. IEEE, 2015.
- [76] Sitaram Asur and Bernardo A Huberman. Predicting the future with social media. In *Proceedings of the 2010 IEEE/WIC/ACM International Conference on Web Intelligence and Intelligent Agent Technology-Volume 01*, pages 492–499. IEEE Computer Society, 2010.
- [77] Anna Klos, Machiel S Bos, and Janusz Bogusz. Detecting time-varying seasonal signal in gps position time series with different noise levels. *GPS solutions*, 22(1):1–11, 2018.
- [78] Robert Stoermer, Ralph Mager, Andreas Roessler, Franz Mueller-Spahn, and Alex H Bullinger. Monitoring human-virtual reality interaction: a time series analysis approach. *CyberPsychology & Behavior*, 3(3):401–406, 2000.
- [79] Shah Muhammed Abid Hussain and ABM Harun-ur Rashid. User independent hand gesture recognition by accelerated dtw. In *ICIEV*, pages 1033–1037. IEEE, 2012.
- [80] Alina Delia Calin. Gesture recognition on kinect time series data using dynamic time warping and hidden markov models. In *SYNASC*, pages 264–271. IEEE, 2016.
- [81] Samet Ayhan and Hanan Samet. Time series clustering of weather observations in predicting climb phase of aircraft trajectories. In *IWCTS*, pages 25–30, 2016.
- [82] Li Li, Xiaonan Su, Yi Zhang, Yuetong Lin, and Zhiheng Li. Trend modeling for traffic time series analysis: An integrated study. *IEEE Transactions on Intelligent Transportation Systems*, 16(6):3430–3439, 2015.
- [83] Stephanie L Hyland, Martin Faltsy, Matthias Hüser, Xinrui Lyu, Thomas Gumbsch, Cristóbal Esteban, Christian Bock, Max Horn, Michael Moor, Bastian Rieck, et al. Early prediction of circulatory failure in the intensive care unit using machine learning. *Nature medicine*, 26(3):364–373, 2020.

- [84] Huanzhou Zhu, Zhuoer Gu, Haiming Zhao, Keyang Chen, Chang-Tsun Li, and Ligang He. Developing a pattern discovery method in time series data and its gpu acceleration. *Big Data Mining and Analytics*, 1(4):266–283, 2018.
- [85] Jia Liu, Yong Xue, Kaijun Ren, Junqiang Song, Christopher Windmill, and Patrick Merritt. High-performance time-series quantitative retrieval from satellite images on a gpu cluster. *IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing*, 12(8):2810–2821, 2019.
- [86] Kai-Wei Chang, Biplab Deka, Wen-Mei W Hwu, and Dan Roth. Efficient pattern-based time series classification on gpu. In *2012 IEEE 12th International Conference on Data Mining*, pages 131–140. IEEE, 2012.
- [87] Igor M Coelho, Vitor N Coelho, Eduardo J da S Luz, Luiz S Ochi, Frederico G Guimarães, and Eyder Rios. A gpu deep learning metaheuristic based model for time series forecasting. *Applied Energy*, 201:412–418, 2017.
- [88] Taban Eslami and Fahad Saeed. Fast-gpu-pcc: A gpu-based technique to compute pairwise pearson’s correlation coefficients for time series data—fmri study. *High-throughput*, 7(2):11, 2018.
- [89] Amin Kalantar, Zachary Zimmerman, and Philip Brisk. Fa-lamp: fpga-accelerated learned approximate matrix profile for time series similarity prediction. In *2021 IEEE 29th Annual International Symposium on Field-Programmable Custom Computing Machines (FCCM)*, pages 40–49. IEEE, 2021.
- [90] Seoungyoung Kang, Jinyeong Moon, and Sang-Woo Jun. Fpga-accelerated time series mining on low-power iot devices. In *2020 IEEE 31st International Conference on Application-specific Systems, Architectures and Processors (ASAP)*, pages 33–36. IEEE, 2020.
- [91] Miquel L Alomar, Vincent Canals, Nicolas Perez-Mora, Víctor Martínez-Moll, and Josep L Rosselló. Fpga-based stochastic echo state networks for time-series forecasting. *Computational intelligence and neuroscience*, 2016, 2016.
- [92] Shaizeen Aga, Supreet Jeloka, Arun Subramanian, Satish Narayanasamy, David Blaauw, and Reetuparna Das. Compute caches. In *HPCA*, 2017.
- [93] Saransh Gupta, Mohsen Imani, and Tajana Rosing. Exploring processing in-memory for different technologies. In *Proceedings of the 2019 on Great Lakes Symposium on VLSI*, pages 201–206, 2019.
- [94] Antonio J Dios, Angeles Navarro, Rafael Asenjo, Francisco Corbera, and Emilio L Zapata. A case study of the task-based parallel wavefront pattern. In *Applications, Tools and Techniques on the Road to Exascale Computing*, pages 65–72. IOS Press, 2012.
- [95] Oleksandra Aleksandrova and Yevgen Bashkov. Face recognition systems based on neural compute stick 2, cpu, gpu comparison. In *ATTI*, pages 104–107. IEEE, 2020.
- [96] Fabrice Devaux. The true processing in memory accelerator. In *2019 IEEE Hot Chips 31 Symposium (HCS)*, pages 1–24. IEEE Computer Society, 2019.
- [97] Juan Gómez-Luna, Izzat El Hajj, Ivan Fernandez, Christina Giannoula, Geraldo F Oliveira, and Onur Mutlu. Benchmarking a new paradigm: An experimental analysis of a real processing-in-memory architecture. *arXiv preprint arXiv:2105.03814*, 2021.
- [98] Daniel Sanchez and Christos Kozyrakis. ZSim: Fast and Accurate Microarchitectural Simulation of Thousand-Core Systems. In *ISCA*, 2013.
- [99] Yoongu Kim, Weikun Yang, and Onur Mutlu. Ramulator: A Fast and Extensible DRAM Simulator. *CAL*, 2015.
- [100] SAFARI Research Group. Ramulator Source Code. <https://github.com/CMU-SAFARI/ramulator>. Accessed 23 September 2020.
- [101] Sheng Li, Jung Ho Ahn, Richard D Strong, Jay B Brockman, Dean M Tullsen, and Norman P Jouppi. McPAT: An Integrated Power, Area, and Timing Modeling Framework for Multicore and Manycore Architectures. In *MICRO*, 2009.
- [102] Intel Open Source. RAPL Power Meter. <https://01.org/rapl-power-meter>, 2014.
- [103] UPMEM. Introduction to UPMEM PIM. Processing-in-memory (PIM) on DRAM Accelerator (White Paper), 2018.
- [104] Xilinx. Xilinx runtime (xrt) documentation. <https://xilinx.github.io/XRT/>. Accessed Oct 10, 2022.
- [105] Rajesh Saha, Yogendra Pratap Pundir, and Pankaj Kumar Pal. Comparative analysis of stt and sot based mrams for last level caches. *Journal of Magnetism and Magnetic Materials*, 551:169161, 2022.
- [106] A. Taddei, G. Distanto, M. Emdin, P. Pisani, G. B. Moody, C. Zee-lenberg, and C. Marchesi. The European ST-T Database: Standard for Evaluating Systems for the Analysis of ST-T Changes in Ambulatory Electrocardiography. *European Heart Journal*, 1992.
- [107] Ashok Veeraraghavan, Rama Chellappa, and Amit K Roy-Chowdhury. The function space of an activity. In *2006 IEEE computer society conference on computer vision and pattern recognition (CVPR’06)*, volume 1, pages 959–968. IEEE, 2006.
- [108] C M Yeh, H V Herle, and E Keogh. Matrix Profile III: The Matrix Profile Allows Visualization of Salient Subsequences in Massive Time Series. In *ICDM*, 2016.
- [109] Penguin Data. <https://www.cs.ucr.edu/~eamonn/MatrixProfile.html>, 2022.
- [110] David Murray, Jing Liao, Lina Stankovic, Vladimir Stankovic, Richard Hauxwell-Baldwin, Charlie Wilson, Michael Coleman, Tom Kane, and Steven Firth. A data management platform for personalised real-time energy feedback. 2015.
- [111] Mesbah Uddin and Garrett S Rose. A practical sense amplifier design for memristive crossbar circuits (puf). In *2018 31st IEEE International System-on-Chip Conference (SOCC)*, pages 209–214. IEEE, 2018.
- [112] Yeongkyo Seo, Kon-Woo Kwon, and Kaushik Roy. Area-efficient sot-mram with a schottky diode. *IEEE Electron Device Letters*, 37(8):982–985, 2016.
- [113] Hui Ding, Goce Trajcevski, Peter Scheuermann, Xiaoyue Wang, and Eamonn Keogh. Querying and mining of time series data: Experimental comparison of representations and distance measures. *Proc. VLDB Endow.*, 1(2):1542–1552, aug 2008.
- [114] Doruk Sart, Abdullah Mueen, Walid Najjar, Eamonn Keogh, and Vit Niennattrakul. Accelerating dynamic time warping subsequence search with gpus and fpgas. In *ICDM*, pages 1001–1006, 2010.
- [115] Gustavo Poli, Joao F. Mari, Jose Hiroki Saito, and Alexandre L. M. Levada. Voice command recognition with dynamic time warping (dtw) using graphics processing units (gpu) with compute unified device architecture (cuda). In *SBAC-PAD*, pages 19–25, 2007.
- [116] Aleksandr Movchan and Mikhail Zymbler. Time series subsequence similarity search under dynamic time warping distance on the intel many-core accelerators. In *Similarity Search and Applications*, pages 295–306. Springer International Publishing, 2015.
- [117] Bertil Schmidt and Christian Hundt. cuDTW++: Ultra-fast dynamic time warping on CUDA-enabled GPUs. In *Euro-Par*, pages 597–612. Springer International Publishing, 2020.
- [118] Kin Fun Li and James Shueyen Tai. Dynamic time warping in hardware. In *IWAS*, page 132–137, 2012.
- [119] Reza Lotfian and Roozbeh Jafari. An ultra-low power hardware accelerator architecture for wearable computers using dynamic time warping. In *DATE*, pages 913–916, 2013.
- [120] Ritik Koul, Mukul Yadav, and Kriti Suneja. Comparative analysis of fpga based hardware design of dynamic time warping algorithm using different multiplier architectures. In *GUCON*, pages 599–603, 2020.
- [121] Zhengyu Chen and Jie Gu. High-throughput dynamic time warping accelerator for time-series classification with pipelined mixed-signal time-domain computing. *IEEE Journal of Solid-State Circuits*, 56(2):624–635, 2021.
- [122] Xiaowei Xu, Feng Lin, Aosen Wang, Xinwei Yao, Qing Lu, Wenyao Xu, Yiyu Shi, and Yu Hu. Accelerating dynamic time warping with memristor-based customized fabrics. *IEEE Trans. on Computer-Aided Design of Integrated Circuits and Systems*, 37(4):729–741, 2018.
- [123] Shaahin Angizi, Jiao Sun, Wei Zhang, and Deliang Fan. Aligns: A processing-in-memory accelerator for dna short read alignment leveraging sot-mram. In *DAC*, pages 1–6. IEEE, 2019.
- [124] Shuangchen Li, Dimin Niu, Krishna T Malladi, Hongzhong Zheng, Bob Brennan, and Yuan Xie. Drisa: A dram-based reconfigurable in-situ accelerator. In *2017 50th Annual IEEE/ACM International Symposium on Microarchitecture (MICRO)*, pages 288–301. IEEE, 2017.
- [125] Shaahin Angizi, Zhezhi He, Amro Awad, and Deliang Fan. Mrima: An mram-based in-memory accelerator. *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems*, 39(5):1123–1136, 2019.
- [126] S. Angizi, A. S. Rakin, and D. Fan. CMP-PIM: An Energy-efficient Comparator-based Processing-in-Memory Neural Network Accelerator. In *DAC*, 2018.
- [127] Shaahin Angizi, Jiao Sun, Wei Zhang, and Deliang Fan. Pim-aligner: A processing-in-mram platform for biological sequence alignment. In *DATE*, pages 1265–1270. IEEE, 2020.
- [128] Saugata Ghose, Amirali Boroumand, Jeremie S Kim, Juan Gómez-Luna, and Onur Mutlu. Processing-in-Memory: A Workload-driven Perspective. *IBM JRD*, 2019.

- [129] Amirali Boroumand, Saugata Ghose, Minesh Patel, Hasan Hassan, Brandon Lucia, Rachata Ausavarungnirun, Kevin Hsieh, Nastaran Hajinazar, Krishna T Malladi, Hongzhong Zheng, et al. CoNDA: Efficient Cache Coherence Support for Near-Data Accelerators. In *ISCA*, 2019.
- [130] Damla Senol Cali, Gurpreet S Kalsi, Zülal Bingöl, Can Firtina, Lavanya Subramanian, Jeremie S Kim, Rachata Ausavarungnirun, Mohammed Alser, Juan Gomez-Luna, Amirali Boroumand, et al. GenASM: A High-Performance, Low-Power Approximate String Matching Acceleration Framework for Genome Sequence Analysis. In *MICRO*, 2020.
- [131] Gagandeep Singh, Juan Gómez-Luna, Giovanni Mariani, Geraldo F Oliveira, Stefano Corda, Sander Stuijk, Onur Mutlu, and Henk Corporaal. NAPEL: Near-memory Computing Application Performance Prediction Via Ensemble Learning. In *DAC*, 2019.
- [132] Gagandeep Singh, Dionysios Diamantopoulos, Christoph Hagleitner, Juan Gomez-Luna, Sander Stuijk, Onur Mutlu, and Henk Corporaal. NERO: A Near High-Bandwidth Memory Stencil Accelerator for Weather Prediction Modeling. In *FPL*, 2020.
- [133] Daichi Fujiki, Scott Mahlke, and Reetuparna Das. Duality Cache for Data Parallel Acceleration. In *ISCA*, 2019.
- [134] Christina Giannoula, Nandita Vijaykumar, Nikola Papadopoulou, Vasileios Karakostas, Ivan Fernandez, Juan Gómez-Luna, Lois Orosa, Nectarios Koziris, Georgios Goumas, and Onur Mutlu. Syncron: Efficient synchronization support for near-data-processing architectures. In *HPCA*, pages 263–276. IEEE, 2021.
- [135] Harold S Stone. A Logic-in-Memory Computer. *IEEE TC*, 1970.
- [136] W. H. Kautz. Cellular Logic-in-Memory Arrays. *IEEE TC*, 1969.
- [137] David Elliot Shaw, Salvatore J. Stolfo, Hussein Ibrahim, Bruce Hillyer, Gio Wiederhold, and JA Andrews. The NON-VON Database Machine: A Brief Overview. *IEEE Database Eng. Bull.*, 1981.
- [138] P. M. Kogge. EXECUBE - A New Architecture for Scaleable MPPs. In *ICPP*, 1994.
- [139] Maya Gokhale, Bill Holmes, and Ken Iobst. Processing in Memory: The Terasys Massively Parallel PIM Array. *IEEE Computer*, 1995.
- [140] David Patterson, Thomas Anderson, Neal Cardwell, Richard Fromm, Kimberly Keeton, Christoforos Kozyrakis, Randi Thomas, and Katherine Yelick. A Case for Intelligent RAM. *IEEE Micro*, 1997.
- [141] Mark Oskin, Frederic T. Chong, and Timothy Sherwood. Active Pages: A Computation Model for Intelligent Memory. In *ISCA*, 1998.
- [142] Yi Kang, Wei Huang, Seung-Moon Yoo, Diana Keen, Zhenzhou Ge, Vinh Lam, Pratap Pattanaik, and Josep Torrellas. FlexRAM: Toward an Advanced Intelligent Memory System. In *ICCD*, 1999.
- [143] Ken Mai, Tim Paaske, Nuwan Jayasena, Ron Ho, William J. Dally, and Mark Horowitz. Smart Memories: A Modular Reconfigurable Architecture. In *ISCA*, 2000.
- [144] Richard C Murphy, Peter M Kogge, and Arun Rodrigues. The Characterization of Data Intensive Memory Workloads on Distributed PIM Systems. In *Intelligent Memory Systems*. Springer, 2001.
- [145] Jeff Draper, Jacqueline Chame, Mary Hall, Craig Steele, Tim Barrett, Jeff LaCoss, John Granacki, Jaewook Shin, Chun Chen, Chang Woo Kang, Ihn Kim, and Gokhan Daglikoca. The Architecture of the DIVA Processing-in-Memory Chip. In *SC*, 2002.
- [146] Shaizeen Aga, Supreet Jeloka, Arun Subramanian, Satish Narayanasamy, David Blaauw, and Reetuparna Das. Compute Caches. In *HPCA*, 2017.
- [147] Charles Eckert, Xiaowei Wang, Jingcheng Wang, Arun Subramanian, Ravi Iyer, Dennis Sylvester, David Blaauw, and Reetuparna Das. Neural Cache: Bit-serial In-cache Acceleration of Deep Neural Networks. In *ISCA*, 2018.
- [148] Mingu Kang, Min-Sun Keel, Naresh R Shanbhag, Sean Eilert, and Ken Cuiwewitz. An Energy-Efficient VLSI Architecture for Pattern Recognition via Deep Embedding of Computation in SRAM. In *ICASSP*, 2014.
- [149] V. Seshadri, D. Lee, T. Mullins, H. Hassan, A. Boroumand, J. Kim, M. A. Kozuch, O. Mutlu, P. B. Gibbons, and T. C. Mowry. Ambit: In-Memory Accelerator for Bulk Bitwise Operations Using Commodity DRAM Technology. In *MICRO*, 2017.
- [150] V. Seshadri, D. Lee, T. Mullins, H. Hassan, A. Boroumand, J. Kim, M. A. Kozuch, O. Mutlu, P. B. Gibbons, and T. C. Mowry. BuddyRAM: Improving the Performance and Efficiency of Bulk Bitwise Operations Using DRAM. arXiv:1611.09988 [cs:AR], 2016.
- [151] Vivek Seshadri, K. Hsieh, A. Boroumand, D. Lee, M. A. Kozuch, O. Mutlu, P. B. Gibbons, and T. C. Mowry. Fast Bulk Bitwise AND and OR in DRAM. *CAL*, 2015.
- [152] Vivek Seshadri, Yoongu Kim, Chris Fallin, Donghyuk Lee, Rachata Ausavarungnirun, Gennady Pekhimenko, Yixin Luo, Onur Mutlu, Michael A Kozuch, Phillip B Gibbons, and Todd C. Mowry. RowClone: Fast and Energy-Efficient In-DRAM Bulk Data Copy and Initialization. In *MICRO*, 2013.
- [153] Shaahin Angizi and Deliang Fan. Graphide: A Graph Processing Accelerator Leveraging In-DRAM-computing. In *GLSVLSI*, 2019.
- [154] J. Kim, M. Patel, H. Hassan, and O. Mutlu. The DRAM Latency PUF: Quickly Evaluating Physical Unclonable Functions by Exploiting the Latency–Reliability Tradeoff in Modern DRAM Devices. In *HPCA*, 2018.
- [155] J. Kim, M. Patel, H. Hassan, L. Orosa, and O. Mutlu. D-RaNGe: Using Commodity DRAM Devices to Generate True Random Numbers with Low Latency and High Throughput. In *HPCA*, 2019.
- [156] Fei Gao, Georgios Tziantzioulis, and David Wentzloff. ComputeDRAM: In-Memory Compute Using Off-the-Shelf DRAMs. In *MICRO*, 2019.
- [157] Kevin K. Chang, Prashant J. Nair, Donghyuk Lee, Saugata Ghose, Moinuddin K. Qureshi, and Onur Mutlu. Low-Cost Inter-Linked Subarrays (LISA): Enabling Fast Inter-Subarray Data Movement in DRAM. In *HPCA*, 2016.
- [158] Xin Xin, Youtao Zhang, and Jun Yang. ELP2IM: Efficient and Low Power Bitwise Operation Processing in DRAM. In *HPCA*, 2020.
- [159] S. Li, D. Niu, K. T. Malladi, H. Zheng, B. Brennan, and Y. Xie. DRISA: A DRAM-Based Reconfigurable In-Situ Accelerator. In *MICRO*, 2017.
- [160] Q. Deng, L. Jiang, Y. Zhang, M. Zhang, and J. Yang. DrAcc: A DRAM Based Accelerator for Accurate CNN Inference. In *DAC*, 2018.
- [161] Nastaran Hajinazar, Geraldo F Oliveira, Sven Gregorio, João Dinis Ferreira, Nika Mansouri Ghiasi, Minesh Patel, Mohammed Alser, Saugata Ghose, Juan Gómez-Luna, and Onur Mutlu. SIMDRAM: A Framework for Bit-Serial SIMD Processing Using DRAM. In *ASPLOS*, 2021.
- [162] S. H. S. Rezaei, M. Modarressi, R. Ausavarungnirun, M. Sadrosadati, O. Mutlu, and M. Daneshtalab. NoM: Network-on-Memory for Inter-Bank Data Transfer in Highly-Banked Memories. *CAL*, 2020.
- [163] Yaohua Wang, Lois Orosa, Xiangjun Peng, Yang Guo, Saugata Ghose, Minesh Patel, Jeremie S Kim, Juan Gómez Luna, Mohammad Sadrosadati, Nika Mansouri Ghiasi, et al. FIGARO: Improving System Performance via Fine-Grained In-DRAM Data Relocation and Caching. In *MICRO*, 2020.
- [164] Mustafa F Ali, Akhilesh Jaiswal, and Kaushik Roy. In-Memory Low-Cost Bit-Serial Addition Using Commodity DRAM Technology. In *TCAS-I*, 2019.
- [165] Shuangchen Li, Cong Xu, Qiaosha Zou, Jishen Zhao, Yu Lu, and Yuan Xie. Pinatubo: A Processing-in-Memory Architecture for Bulk Bitwise Operations in Emerging Non-Volatile Memories. In *DAC*, 2016.
- [166] S. Angizi, Z. He, and D. Fan. PIMA-Logic: A Novel Processing-in-Memory Architecture for Highly Flexible and Energy-efficient Logic Computation. In *DAC*, 2018.
- [167] S. Angizi, J. Sun, W. Zhang, and D. Fan. AlignS: A Processing-in-Memory Accelerator for DNA Short Read Alignment Leveraging SOT-MRAM. In *DAC*, 2019.
- [168] Yifat Levy, Jehoshua Bruck, Yuval Cassuto, Eby G. Friedman, Avinoam Kolodny, Eitan Yaakobi, and Shahar Kvatinsky. Logic Operations in Memory Using a Memristive Akers Array. *Microelectronics Journal*, 2014.
- [169] S. Kvatinsky, D. Belousov, S. Liman, G. Satat, N. Wald, E. G. Friedman, A. Kolodny, and U. C. Weiser. MAGIC—Memristor-Aided Logic. *IEEE TCAS II: Express Briefs*, 2014.
- [170] Ali Shafiee, Anirban Nag, Naveen Muralimanohar, Rajeev Balasubramanian, John Paul Strachan, Miao Hu, R Stanley Williams, and Vivek Srikumar. ISAAC: A Convolutional Neural Network Accelerator with In-situ Analog Arithmetic in Crossbars. *ISCA*, 2016.
- [171] S. Kvatinsky, A. Kolodny, U. C. Weiser, and E. G. Friedman. Memristor-Based IMPLY Logic Design Procedure. In *ICCD*, 2011.
- [172] S. Kvatinsky, G. Satat, N. Wald, E. G. Friedman, A. Kolodny, and U. C. Weiser. Memristor-Based Material Implication (IMPLY) Logic: Design Principles and Methodologies. *TVLSI*, 2014.
- [173] P.-E. Gaillardon, L. Amaru, A. Siemon, and et al. The Programmable Logic-in-Memory (PLiM) Computer. In *DATE*, 2016.
- [174] D. Bhattacharjee, R. Devadoss, and A. Chattopadhyay. ReVAMP: ReRAM based VLIW Architecture for In-memory Computing. In *DATE*, 2017.
- [175] S. Hamdioui, L. Xie, H. A. D. Nguyen, and et al. Memristor Based Computation-in-Memory Architecture for Data-intensive Applications. In *DATE*, 2015.
- [176] L. Xie, H. A. D. Nguyen, M. Taouil, and et al. Fast Boolean Logic Papped on Memristor Crossbar. In *ICCD*, 2015.

- [177] S. Hamdioui, S. Kvatinsky, and et al. G. Cauwenberghs. Memristor for Computing: Myth or Reality? In *DATE*, 2017.
- [178] J. Yu, H. A. D. Nguyen, L. Xie, and et al. Memristive Devices for Computation-in-Memory. In *DATE*, 2018.
- [179] Christina Giannoula, Nandita Vijaykumar, Nikela Papadopoulou, Vasileios Karakostas, Ivan Fernandez, Juan Gómez-Luna, Lois Orosa, Nectarios Koziris, Georgios Goumas, and Onur Mutlu. SynCron: Efficient Synchronization Support for Near-Data-Processing Architectures. In *HPCA*, 2021.
- [180] Ivan Fernandez, Ricardo Quisilant, Eladio Gutiérrez, Oscar Plata, Christina Giannoula, Mohammed Alser, Juan Gómez-Luna, and Onur Mutlu. NATSA: A Near-data Processing Accelerator for Time Series Analysis. In *ICCD*, 2020.
- [181] J. S. Kim, D. Senol, H. Xin, D. Lee, S. Ghose, M. Alser, H. Hassan, O. Ergin, C. Alkan, and O. Mutlu. GRIM-Filter: Fast Seed Location Filtering in DNA Read Mapping Using Processing-in-Memory Technologies. *BMC Genomics*, 2018.
- [182] Junwhan Ahn, Sungjoo Yoo, Onur Mutlu, and Kiyoun Choi. PIM-Enabled Instructions: A Low-Overhead, Locality-Aware Processing-in-Memory Architecture. In *ISCA*, 2015.
- [183] Junwhan Ahn, Sungpack Hong, Sungjoo Yoo, Onur Mutlu, and Kiyoun Choi. A Scalable Processing-in-Memory Accelerator for Parallel Graph Processing. In *ISCA*, 2015.
- [184] A. Boroumand, S. Ghose, Y. Kim, R. Ausavarungnirun, E. Shiu, R. Thakur, D. Kim, A. Kuusela, A. Knies, P. Ranganathan, and O. Mutlu. Google Workloads for Consumer Devices: Mitigating Data Movement Bottlenecks. In *ASPLOS*, 2018.
- [185] Hadi Aghari-Moghaddam, Young Hoon Son, Jung Ho Ahn, and Nam Sung Kim. Chameleon: Versatile and Practical Near-DRAM Acceleration Architecture for Large Memory Systems. In *MICRO*, 2016.
- [186] Oreoluwatomiwa O. Babarinsa and Stratos Idreos. JAFAR: Near-Data Processing for Databases. In *SIGMOD*, 2015.
- [187] P. Chi, S. Li, C. Xu, T. Zhang, J. Zhao, Y. Liu, Y. Wang, and Y. Xie. PRIME: A Novel Processing-In-Memory Architecture for Neural Network Computation In ReRAM-Based Main Memory. In *ISCA*, 2016.
- [188] Amin Farmahini-Farahani, Jung Ho Ahn, Katherine Morrow, and Nam Sung Kim. NDA: Near-DRAM acceleration architecture leveraging commodity DRAM devices and standard memory modules. In *HPCA*, 2015.
- [189] Mingyu Gao, Grant Ayers, and Christos Kozyrakis. Practical Near-Data Processing for In-Memory Analytics Frameworks. In *PACT*, 2015.
- [190] Mingyu Gao and Christos Kozyrakis. HRL: Efficient and Flexible Reconfigurable Logic for Near-Data Processing. In *HPCA*, 2016.
- [191] Boncheol Gu, A. S. Yoon, D.-H. Bae, I. Jo, J. Lee, J. Yoon, J.-U. Kang, M. Kwon, C. Yoon, S. Cho, J. Jeong, and D. Chang. Biscuit: A Framework for Near-Data Processing of Big Data Workloads. In *ISCA*, 2016.
- [192] Q. Guo, N. Alachiotis, B. Akin, F. Sadi, G. Xu, T. M. Low, L. Pileggi, J. C. Hoe, and F. Franchetti. 3D-Stacked Memory-Side Acceleration: Accelerator and System Design. In *WoNDP*, 2014.
- [193] Milad Hashemi, Khubaib, Eiman Ebrahimi, Onur Mutlu, and Yale N. Patt. Accelerating Dependent Cache Misses with an Enhanced Memory Controller. In *ISCA*, 2016.
- [194] M. Hashemi, O. Mutlu, and Y. N. Patt. Continuous Runahead: Transparent Hardware Acceleration for Memory Intensive Workloads. In *MICRO*, 2016.
- [195] Kevin Hsieh, Eiman Ebrahimi, Gwangsun Kim, Niladri Chatterjee, Mike O’Conner, Nandita Vijaykumar, Onur Mutlu, and Stephen Keckler. Transparent Offloading and Mapping (TOM): Enabling Programmer-Transparent Near-Data Processing in GPU Systems. In *ISCA*, 2016.
- [196] Duckhwan Kim, J. Kung, S. Chai, S. Yalamanchili, and S. Mukhopadhyay. Neurocube: A Programmable Digital Neuromorphic Architecture with High-Density 3D Memory. In *ISCA*, 2016.
- [197] G. Kim, N. Chatterjee, M. O’Connor, and K. Hsieh. Toward Standardized Near-Data Processing with Unrestricted Data Placement for GPUs. In *SC*, 2017.
- [198] Joo Hwan Lee, Jaewoong Sim, and Hyesoon Kim. BSSync: Processing Near Memory for Machine Learning Workloads with Bounded Staleness Consistency Models. In *PACT*, 2015.
- [199] Zhiyu Liu, Irina Calciu, Maurice Herlihy, and Onur Mutlu. Concurrent Data Structures for Near-Memory Computing. In *SPAA*, 2017.
- [200] Amir Morad, Leonid Yavits, and Ran Ginosar. GP-SIMD Processing-in-Memory. *ACM TACO*, 2015.
- [201] Lifeng Nai, Ramyad Hadidi, Jaewoong Sim, Hoyjong Kim, Pranith Kumar, and Hyesoon Kim. GraphPIM: Enabling Instruction-Level PIM Offloading in Graph Computing Frameworks. In *HPCA*, 2017.
- [202] Ashutosh Pattnaik, Xulong Tang, Adwait Jog, Onur Kayiran, Asit K. Mishra, Mahmut T. Kandemir, Onur Mutlu, and Chita R. Das. Scheduling Techniques for GPU Architectures with Processing-in-Memory Capabilities. In *PACT*, 2016.
- [203] Seth H. Pugsley, Jeffrey Jests, Huihui Zhang, Rajeev Balasubramanian, Vijayalakshmi Srinivasan, Alper Buyuktosunoglu, Al Davis, and Feifei Li. NDC: Analyzing the Impact of 3D-Stacked Memory+Logic Devices on MapReduce Workloads. In *ISPASS*, 2014.
- [204] D. P. Zhang, N. Jayasena, A. Lyashevsky, J. L. Greathouse, L. Xu, and M. Ignatowski. TOP-PIM: Throughput-Oriented Programmable Processing in Memory. In *HPDC*, 2014.
- [205] Qiuling Zhu, Tobias Graf, H Ekin Sumbul, Larry Pileggi, and Franz Franchetti. Accelerating Sparse Matrix-Matrix Multiplication with 3D-Stacked Logic-in-Memory Hardware. In *HPEC*, 2013.
- [206] Berkin Akin, Franz Franchetti, and James C. Hoe. Data Reorganization in Memory Using 3D-Stacked DRAM. In *ISCA*, 2015.
- [207] Mingyu Gao, Jing Pu, Xuan Yang, Mark Horowitz, and Christos Kozyrakis. Tetris: Scalable and Efficient Neural Network Acceleration with 3D Memory. In *ASPLOS*, 2017.
- [208] Mario Drumond, Alexandros Daglis, Nooshin Mirzadeh, Dmitrii Ustiugov, Javier Picorel Obando, Babak Falsafi, Boris Grot, and Dionisios Pnevmatikatos. The Mondrian Data Engine. In *ISCA*, 2017.
- [209] Guohao Dai, Tianhao Huang, Yuze Chi, Jishen Zhao, Guangyu Sun, Yongpan Liu, Yu Wang, Yuan Xie, and Huazhong Yang. GraphH: A Processing-in-Memory Architecture for Large-scale Graph Processing. *IEEE TCAD*, 2018.
- [210] Mingxing Zhang, Youwei Zhuo, Chao Wang, Mingyu Gao, Yongwei Wu, Kang Chen, Christos Kozyrakis, and Xuehai Qian. GraphP: Reducing Communication for PIM-based Graph Processing with Efficient Data Partition. In *HPCA*, 2018.
- [211] Yu Huang, Long Zheng, Pengcheng Yao, Jiashan Zhao, Xiaofei Liao, Hai Jin, and Jingling Xue. A Heterogeneous PIM Hardware-Software Co-Design for Energy-Efficient Graph Processing. In *IPDPS*, 2020.
- [212] Youwei Zhuo, Chao Wang, Mingxing Zhang, Rui Wang, Dimin Niu, Yanzhi Wang, and Xuehai Qian. GraphQ: Scalable PIM-based Graph Processing. In *MICRO*, 2019.
- [213] Paulo C Santos, Geraldo F Oliveira, Diego G Tomé, Marco AZ Alves, Eduardo C Almeida, and Luigi Carro. Operand Size Reconfiguration for Big Data Processing in Memory. In *DATE*, 2017.
- [214] Saugata Ghose, Amirali Boroumand, Jeremie S Kim, Juan Gómez-Luna, and Onur Mutlu. Processing-in-Memory: A Workload-Driven Perspective. *IBM JRD*, 2019.
- [215] Wen-Mei Hwu, Izzat El Hajj, Simon Garcia De Gonzalo, Carl Pearson, Nam Sung Kim, Deming Chen, Jinjun Xiong, and Zehra Sura. Rebooting the Data Access Hierarchy of Computing Systems. In *ICRC*, 2017.
- [216] Maciej Besta, Raghavendra Kanakagiri, Grzegorz Kwasniewski, Rachata Ausavarungnirun, Jakub Beránek, Konstantinos Kanellopoulos, Kacper Janda, Zur Vonarburg-Shmria, Lukas Gianinazzi, Ioana Stefan, et al. SISA: Set-Centric Instruction Set Architecture for Graph Mining on Processing-in-Memory Systems. In *MICRO*, 2021.
- [217] João Dinis Ferreira, Gabriel Falcao, Juan Gómez-Luna, Mohammed Alser, Lois Orosa, Mohammad Sadrosadati, Jeremie S. Kim, Geraldo F. Oliveira, Taha Shahroodi, Anant Nori, and Onur Mutlu. pLUTO: In-DRAM Lookup Tables to Enable Massively Parallel General-Purpose Computation. *arXiv:2104.07699 [cs.AR]*, 2021.
- [218] Ataberk Olgun, Minesh Patel, Abdullah Giray Yağlıkcı, Haocong Luo, Jeremie S. Kim, F. Nisa Bostancı, Nandita Vijaykumar, Oğuz Ergin, and Onur Mutlu. QUAC-TRNG: High-Throughput True Random Number Generation Using Quadruple Row Activation in Commodity DRAMs. In *ISCA*, 2021.
- [219] Scott Lloyd and Maya Gokhale. In-memory Data Rearrangement for Irregular, Data-intensive Computing. *Computer*, 2015.
- [220] Duncan G Elliott, Michael Stumm, W Martin Snelgrove, Christian Cojocar, and Robert McKenzie. Computational RAM: Implementing Processors in Memory. *IEEE Design & Test of Computers*, 1999.
- [221] Le Zheng, Sangho Shin, Scott Lloyd, Maya Gokhale, Kyungmin Kim, and Sung-Mo Kang. RRAM-based TCAMs for pattern search. In *ISCAS*, 2016.
- [222] Joshua Landgraf, Scott Lloyd, and Maya Gokhale. Combining Emulation and Simulation to Evaluate a Near Memory Key/Value Lookup Accelerator, 2021.
- [223] Arun Rodrigues, Maya Gokhale, and Gwendolyn Voskuilen. Towards a Scatter-Gather Architecture: Hardware and Software Issues. In *MEMSYS*, 2019.



- [224] Scott Lloyd and Maya Gokhale. Design Space Exploration of Near Memory Accelerators. In *MEMSYS*, 2018.
- [225] Scott Lloyd and Maya Gokhale. Near Memory Key/Value Lookup Acceleration. In *MEMSYS*, 2017.
- [226] Maya Gokhale, Scott Lloyd, and Chris Hajas. Near Memory Data Structure Rearrangement. In *MEMSYS*, 2015.
- [227] Ravi Nair, Samuel F Antao, Carlo Bertolli, Pradip Bose, Jose R Brunheroto, Tong Chen, C-Y Cher, Carlos HA Costa, Jun Doi, Constantinos Evangelinos, et al. Active Memory Cube: A Processing-in-Memory Architecture for Exascale Systems. *IBM JRD*, 2015.
- [228] Arpith C Jacob, Zehra Sura, Tong Chen, Carlo Bertolli, Samuel Antao, Olivier Sallenave, Kevin O'Brien, Hans Jacobson, Ravi Nair, Jose R Brunheroto, et al. Compiling for the Active Memory Cube. Technical report, Tech. rep. RC25644 (WAT1612-008). IBM Research Division, 2016.
- [229] Zehra Sura, Arpith Jacob, Tong Chen, Bryan Rosenburg, Olivier Sallenave, Carlo Bertolli, Samuel Antao, Jose Brunheroto, Yoonho Park, Kevin O'Brien, et al. Data Access Optimization in a Processing-in-Memory System. In *CF*, 2015.
- [230] Ravi Nair. Evolution of Memory Architecture. *Proceedings of the IEEE*, 2015.
- [231] Rajeev Balasubramonian, Jichuan Chang, Troy Manning, Jaime H Moreno, Richard Murphy, Ravi Nair, and Steven Swanson. Near-Data Processing: Insights from a MICRO-46 Workshop. *IEEE Micro*, 2014.
- [232] Yue Xi, Bin Gao, Jianshi Tang, An Chen, Meng-Fan Chang, Xiaobo Sharon Hu, Jan Van Der Spiegel, He Qian, and Huaqiang Wu. In-Memory Learning With Analog Resistive Switching Memory: A Review and Perspective. *Proceedings of the IEEE*, 2020.
- [233] K. Hsieh, S. Khan, N. Vijaykumar, K. K. Chang, A. Boroumand, S. Ghose, and O. Mutlu. Accelerating Pointer Chasing in 3D-Stacked Memory: Challenges, Mechanisms, Evaluation. In *ICCD*, 2016.
- [234] Amiral Boroumand, Saugata Ghose, Minesh Patel, Hasan Hassan, Brandon Lucia, Kevin Hsieh, Krishna T. Malladi, Hongzhong Zheng, and Onur Mutlu. LazyPIM: An Efficient Cache Coherence Mechanism for Processing-in-Memory. *CAL*, 2016.
- [235] Christina Giannoula, Ivan Fernandez, Juan Gómez-Luna, Nectarios Koziris, Georgios Goumas, and Onur Mutlu. SparseP: Towards Efficient Sparse Matrix Vector Multiplication on Real Processing-In-Memory Systems. *arXiv preprint arXiv:2201.05072*, 2022.
- [236] Christina Giannoula, Ivan Fernandez, Juan Gómez-Luna, Nectarios Koziris, Georgios Goumas, and Onur Mutlu. Towards Efficient Sparse Matrix Vector Multiplication on Real Processing-in-Memory Architectures. In *SIGMETRICS*, 2022.
- [237] Alain Denzler, Rahul Bera, Nastaran Hajinazar, Gagandeep Singh, Geraldo F Oliveira, Juan Gómez-Luna, and Onur Mutlu. Casper: Accelerating stencil computation using near-cache processing. *arXiv preprint arXiv:2112.14216*, 2021.
- [238] Amiral Boroumand, Saugata Ghose, Geraldo F Oliveira, and Onur Mutlu. Polynesia: Enabling Effective Hybrid Transactional/Analytical Databases with Specialized Hardware/Software Co-Design. *arXiv:2103.00798 [cs.AR]*, 2021.
- [239] Amiral Boroumand, Saugata Ghose, Geraldo F Oliveira, and Onur Mutlu. Polynesia: Enabling Effective Hybrid Transactional Analytical Databases with Specialized Hardware Software Co-Design. In *ICDE*, 2022.
- [240] Gagandeep Singh, Mohammed Alser, Damla Senol Cali, Dionysios Diamantopoulos, Juan Gómez-Luna, Henk Corporaal, and Onur Mutlu. FPGA-based Near-Memory Acceleration of Modern Data-Intensive Applications. *IEEE Micro*, 2021.
- [241] Gagandeep Singh, Dionysios Diamantopoulos, Juan Gómez-Luna, Christoph Hagleitner, Sander Stuijk, Henk Corporaal, and Onur Mutlu. Accelerating Weather Prediction using Near-Memory Reconfigurable Fabric. *ACM TRETS*, 2021.
- [242] Jose M Herruzo, Ivan Fernandez, Sonia González-Navarro, and Oscar Plata. Enabling Fast and Energy-Efficient FM-Index Exact Matching Using Processing-Near-Memory. *The Journal of Supercomputing*, 2021.
- [243] Leonid Yavits, Roman Kaplan, and Ran Ginosar. GIRAF: General Purpose In-Storage Resistive Associative Framework. *IEEE TPDS*, 2021.
- [244] Bahar Asgari, Ramyad Hadidi, Jiashen Cao, Da Eun Shim, Sung-Kyu Lim, and Hyesoon Kim. FAFNIR: Accelerating Sparse Gathering by Using Efficient Near-Memory Intelligent Reduction. In *HPCA*, 2021.
- [245] Amiral Boroumand, Saugata Ghose, Berkin Akin, Ravi Narayanaswami, Geraldo F Oliveira, Xiaoyu Ma, Eric Shiu, and Onur Mutlu. Google Neural Network Models for Edge Devices: Analyzing and Mitigating Machine Learning Inference Bottlenecks. *arXiv preprint arXiv:2109.14320*, 2021.
- [246] Amiral Boroumand, Saugata Ghose, Berkin Akin, Ravi Narayanaswami, Geraldo F. Oliveira, Xiaoyu Ma, Eric Shiu, and Onur Mutlu. Google Neural Network Models for Edge Devices: Analyzing and Mitigating Machine Learning Inference Bottlenecks. In *PACT*, 2021.
- [247] Amiral Boroumand. *Practical Mechanisms for Reducing Processor-Memory Data Movement in Modern Workloads*. PhD thesis, Carnegie Mellon University, 2020.
- [248] Vivek Seshadri and Onur Mutlu. Simple Operations in Memory to Reduce Data Movement. In *Advances in Computers, Volume 106*. 2017.
- [249] Safaa Diab, Amir Nassereldine, Mohammed Alser, Juan Gómez Luna, Onur Mutlu, and Izzat El Hajj. High-throughput Pairwise Alignment with the Wavefront Algorithm using Processing-in-Memory. *arXiv preprint arXiv:2204.02085*, 2022.
- [250] Safaa Diab, Amir Nassereldine, Mohammed Alser, Juan Gómez Luna, Onur Mutlu, and Izzat El Hajj. High-throughput Pairwise Alignment with the Wavefront Algorithm using Processing-in-Memory. In *HICOMB*, 2022.
- [251] Daichi Fujiki, Scott Mahlke, and Reetuparna Das. In-Memory Data Parallel Processor. In *ASPLOS*, 2018.
- [252] Yue Zha and Jing Li. Hyper-AP: Enhancing Associative Processing Through A Full-Stack Optimization. In *ISCA*, 2020.
- [253] Onur Mutlu. Memory Scaling: A Systems Architecture Perspective. *IMW*, 2013.
- [254] Onur Mutlu and Lavanya Subramanian. Research Problems and Opportunities in Memory Systems. *SUPERFRI*, 2014.
- [255] Hameeza Ahmed, Paulo C Santos, João PC Lima, Rafael F Moura, Marco AZ Alves, Antônio CS Beck, and Luigi Carro. A Compiler for Automatic Selection of Suitable Processing-in-Memory Instructions. In *DATE*, 2019.
- [256] Shubham Jain, Sachin Sapatnekar, Jian-Ping Wang, Kaushik Roy, and Anand Raghunathan. Computing-in-memory with spintronics. In *2018 Design, Automation & Test in Europe Conference & Exhibition (DATE)*, pages 1640–1645. IEEE, 2018.
- [257] Nika Mansouri Ghiasi, Jisung Park, Harun Mustafa, Jeremie Kim, Ataberk Olgun, Arvid Gollwitzer, Damla Senol Cali, Can Firtina, Haiyu Mao, Nour Almadhoun Alser, et al. GenStore: A High-Performance and Energy-Efficient In-Storage Computing System for Genome Sequence Analysis. In *ASPLOS*, 2022.
- [258] Geraldo F. Oliveira, Juan Gómez-Luna, Lois Orosa, Saugata Ghose, Nandita Vijaykumar, Ivan Fernandez, Mohammad Sadrosadati, and Onur Mutlu. DAMOV: A New Methodology and Benchmark Suite for Evaluating Data Movement Bottlenecks. *IEEE Access*, 2021.
- [259] Geraldo F. Oliveira, Juan Gómez-Luna, Lois Orosa, Saugata Ghose, Nandita Vijaykumar, Ivan Fernandez, Mohammad Sadrosadati, and Onur Mutlu. DAMOV: A New Methodology and Benchmark Suite for Evaluating Data Movement Bottlenecks. *arXiv:2105.03725 [cs.AR]*, 2021.
- [260] Seunghwan Cho, Haerang Choi, Eunhyeok Park, Hyunsung Shin, and Sungjoo Yoo. McDRAM v2: In-Dynamic Random Access Memory Systolic Array Accelerator to Address the Large Model Problem in Deep Neural Networks on the Edge. *IEEE Access*, 2020.
- [261] Hyunsung Shin, Dongyoung Kim, Eunhyeok Park, Sungho Park, Yongsik Park, and Sungjoo Yoo. McDRAM: Low latency and energy-efficient matrix computations in DRAM. *IEEE TCADICS*, 2018.
- [262] Peng Gu, Xinfeng Xie, Yufei Ding, Guoyang Chen, Weifeng Zhang, Dimin Niu, and Yuan Xie. iPIM: Programmable In-Memory Image Processing Accelerator using Near-Bank Architecture. In *ISCA*, 2020.
- [263] D. Lavenier, R. Cimadomo, and R. Jodin. Variant Calling Parallelization on Processor-in-Memory Architecture. In *BIBM*, 2020.
- [264] Vasileios Zois, Divya Gupta, Vassilis J. Tsotras, Walid A. Najjar, and Jean-Francois Roy. Massively Parallel Skyline Computation for Processing-in-Memory Architectures. In *PACT*, 2018.
- [265] UPMEM. UPMEM Website. <https://www.upmem.com>, 2020.
- [266] Juan Gómez-Luna, Izzat El Hajj, Ivan Fernández, Christina Giannoula, Geraldo F. Oliveira, and Onur Mutlu. Benchmarking a New Paradigm: An Experimental Analysis of a Real Processing-in-Memory Architecture. *arXiv:2105.03814 [cs.AR]*, 2021.
- [267] Juan Gómez-Luna, Izzat El Hajj, Ivan Fernandez, Christina Giannoula, Geraldo F. Oliveira, and Onur Mutlu. Benchmarking a New Paradigm: Experimental Analysis and Characterization of a Real Processing-in-Memory System. *IEEE Access*, 2022.

- [268] Juan Gómez-Luna, Izzat El Hajj, Ivan Fernandez, Christina Giannoula, Geraldo F. Oliveira, and Onur Mutlu. Benchmarking Memory-Centric Computing Systems: Analysis of Real Processing-In-Memory Hardware. In *IGSC*, 2021.
- [269] Mario Paulo Drumond, Alexandros Daglis, Nooshin Mirzadeh, Dmitrii Ustiugov, Javier Picorel Obando, Babak Falsafi, Boris Grot, and Dionisios Pneumatikatos. The Mondrian Data Engine. In *ISCA*, 2017.
- [270] Kevin Hsieh, Samira Khan, Nandita Vijaykumar, Kevin Chang, Amirali Boroumand, Saugata Ghose, and Onur Mutlu. Accelerating Pointer Chasing in 3D-stacked Memory: Challenges, Mechanisms, Evaluation. In *ICCD*, 2016.
- [271] Junwhan Ahn, Sungpack Hong, Sungjoo Yoo, Onur Mutlu, and Kiyoung Choi. A scalable processing-in-memory accelerator for parallel graph processing. In *ISCA*, pages 105–117, 2015.
- [272] V T Lee, A Mazumdar, C C del Mundo, A Alaghi, L Ceze, and M Oskin. Application Codesign of NDP for Similarity Search. In *IPDPS*, 2018.
- [273] Jeremie S. Kim, Damla Senol Cali, Hongyi Xin, Donghyuk Lee, Saugata Ghose, Mohammed Alser, Hasan Hassan, Oguz Ergin, Can Alkan, and Onur Mutlu. GRIM-Filter: Fast seed Location Filter. in DNA Read Mapping Using PIM Technologies. *BMC Genomics*, 2018.
- [274] Hongyi Xin, Donghyuk Lee, Farhad Hormozdiari, Can Alkan, and Onur Mutlu. FastHASH: A New GPU-friendly Algorithm for Fast and Comprehensive Next-Generation Sequence Mapping. In *BMC Genomics*, 2013.
- [275] Mohammed Alser, Hasan Hassan, Akash Kumar, Onur Mutlu, and Can Alkan. Shouji: a Fast and Efficient Pre-alignment Filter for Sequence Alignment. *Bioinformatics*, 2019.
- [276] Hongyi Xin, John Greth, John Emmons, Gennady Pekhimenko, Carl Kingsford, Can Alkan, and Onur Mutlu. Shifted Hamming Distance: A Fast and Accurate SIMD-friendly Filter to Accelerate Alignment Verification in Read Mapping. *Bioinformatics*, 2015.
- [277] Mohammed Alser, Taha Shahroodi, Juan Gomez-Luna, Can Alkan, and Onur Mutlu. SneakySnake: A Fast and Accurate Universal Genome Pre-Alignment Filter for CPUs, GPUs, and FPGAs. *arXiv*, 2019.
- [278] Mohammed Alser, Hasan Hassan, Hongyi Xin, Oğuz Ergin, Onur Mutlu, and Can Alkan. GateKeeper: A New Hardware Arch. for Accelerating Pre-alignment in DNA Short Read Mapping. *Bioinformatics*, 2017.
- [279] Mohammed Alser, Zülal Bingöl, Damla Senol Cali, Jeremie Kim, Saugata Ghose, Can Alkan, and Onur Mutlu. Accelerating Genome Analysis: A Primer on an Ongoing Journey. *IEEE Micro*, 2020.
- [280] Christina Giannoula, Ivan Fernandez, Juan Gómez-Luna, Nectarios Koziris, Georgios Goumas, and Onur Mutlu. SparseP: Towards Efficient Sparse Matrix Vector Multiplication on Real Processing-In-Memory Architectures. In *Proc. ACM Meas. Anal. Comput. Syst.*, 2022.
- [281] Christina Giannoula, Ivan Fernandez, Juan Gómez-Luna, Nectarios Koziris, Georgios Goumas, and Onur Mutlu. Towards Efficient Sparse Matrix Vector Multiplication on Real Processing-In-Memory Architectures. In *SIGMETRICS*, 2022.