

# Effortless Locality on Data Systems using Relational Fabric

Tarikul Islam Papon<sup>\*§</sup>, Ju Hyoung Mun<sup>\*§</sup>, Konstantinos Karatsenidis<sup>\*</sup>, Shahin Roozkhosh<sup>\*</sup>, Denis Hoornaert<sup>†</sup>, Ahmed Sanaullah<sup>‡</sup>, Ulrich Drepper<sup>‡</sup>, Renato Mancuso<sup>\*</sup>, Manos Athanassoulis<sup>\*</sup>

<sup>\*</sup>Boston University, email: {papon, jmun, karatse, shahin, rmancuso, mathan}@bu.edu

<sup>†</sup>Technical University of Munich, email: denis.hoornaert@tum.de

<sup>‡</sup>Red Hat, email: {asanauull, drepper}@redhat.com

**Abstract**—A key design decision for data systems is whether they follow the row-store or the column-store paradigm. The former supports transactional workloads, while the latter is better for analytical queries. This decision has a significant impact on the entire data system architecture. The multiple-decade-long journey of these two designs has led to a new family of hybrid transactional/analytical processing (HTAP) architectures. Several efforts have been proposed to reap the benefits of both worlds by proposing systems that maintain multiple copies of data (in different physical layouts) and convert them into the desired layout as required. Due to data duplication, the additional necessary bookkeeping, and the cost of converting data between different layouts, these systems compromise between efficient analytics and data freshness. We depart from existing designs by proposing a radically new approach. We ask the question:

“What if we could access any layout and ship only the relevant data through the memory hierarchy by transparently converting rows to (arbitrary groups of) columns?”

To achieve this functionality, we capitalize on the reinvigorated trend of *hardware specialization* (that has been accelerated due to the tapering of Moore’s law) to propose *Relational Fabric*, a near-data vertical partitioner that allows memory or storage components to perform *on-the-fly* transparent data transformation. By exposing an intuitive API, Relational Fabric pushes vertical partitioning to the hardware, which profoundly impacts the process of designing and building data systems. (A) There is no need for data duplication and layout conversion, making HTAP systems viable using a single layout. (B) It simplifies the memory and storage manager that needs to maintain and update a single data layout. (C) It reduces unnecessary data movement through the memory hierarchy allowing for better hardware utilization and, ultimately, better performance. In this paper, we present Relational Fabric for both memory and storage. We present our initial results on Relational Fabric for in-memory systems and discuss the challenges of building this hardware and the opportunities it brings for simplicity and innovation in the data system software stack, including physical design, query optimization, query evaluation, and concurrency control.

**Index Terms**—HTAP, Data layout, FPGA, Near-data processing

## I. INTRODUCTION

**OLTP vs. OLAP vs. HTAP.** The popularity of large-scale real-time data analytics has soared over the past few years [58], [62], further fueled by the advent of new technological trends like 5G, Internet-of-Things, and cloud computing [20], [32]. Database management systems (DBMS) now need to perform both Online Transactional Processing (OLTP) and Online

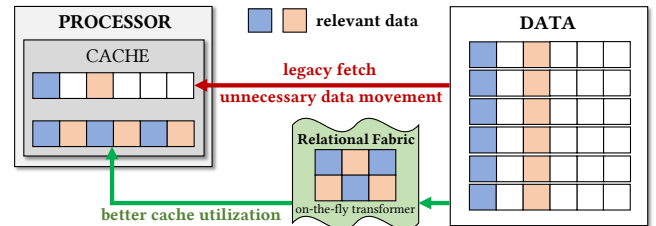


Fig. 1: Relational Fabric removes unnecessary data movement by transparent on-the-fly data transformation.

Analytical Processing (OLAP) since many applications need to analyze fresh data. However, traditional OLTP and OLAP systems are very different – OLTP systems generally use a row-oriented data layout to optimize for write-intensive workloads and point queries. In contrast, OLAP systems use a columnar layout to optimize for read-only analytical queries. Hybrid Transactional/Analytical Processing (HTAP) [62] systems aim to maintain a single data store that offers data freshness and efficient analytics on that same data set. HTAP systems attempt to bridge the OLAP and OLTP requirements by maintaining multiple copies of data in different formats [15], [43] or converting data between different layouts [10], [13], [48], [54], [77]. However, recent research on such hybrid layouts [6], [13], [24], [44] shows the optimal layout of a query is often neither a columnar nor a row-oriented one. Such approaches carry extra complexity, either to convert data between layouts or to manage multiple versions and copies of the data, while not always being able to offer the optimal layout.

**Hardware Specialization On The Rescue.** Due to the historically exponential growth of processor speed, the idea of *hardware specialization*, although recurring every few years, has not been able to attain its true potential. As Moore’s law slows down and the data processing needs keep growing, hardware specialization is becoming a feasible, and more scalable alternative to general-purpose computing [38], [83]. In this *post-Moore* era, this has been further accelerated by the advancements in reconfigurable logic [30], [53]. Recent technological developments (e.g., Google’s TPU [33], Microsoft’s Catapult [57]) show that specialization is here to stay [38]. To quote James Hamilton of AWS “*When there might be several hundred thousand servers all running the same workload, hardware specialization goes from an interesting idea to almost a responsibility*” [36].

<sup>§</sup>These authors contributed equally.

**Idea: Transparent Data Transformation**

*What if we can access any arbitrary data layout using near-data processing via specialized hardware?*

In other words, “*what if the optimal data layout is always physically available?*”. This will eliminate the need to keep multiple layouts and we can perform efficient analytics over the fresh data without converting between different layouts. If the underlying specialized hardware accesses only the relevant data (without accessing unnecessary data and without paying a tuple reconstruction cost) while maintaining a single layout, it will blend the benefits of row-stores and column-stores. It will offer *effortless locality*, alleviating the need for separate row-store and column-store query engines, decoupling the *physical* data layout from the *processing* data layout.

**Our Vision: Relational Fabric**

*A lightweight specialized hardware fabric that allows accessing arbitrary data layouts from memory or storage.*

**Relational Fabric.** We propose *Relational Fabric* [65], a new lightweight specialized hardware fabric that accommodates queries to access arbitrary column groups from memory-resident base data without requiring any data duplication. The base data is stored in a row-oriented physical layout (Figure 1), to allow efficient data ingestion and updates, read-only queries can quickly access only the relevant column groups (or the entire row, if needed) using the underlying machinery. To do this, Relational Fabric exposes a carefully designed API, termed *ephemeral columns* that enables accessing arbitrary *data geometries* (i.e., any subset of data from relational tables) using simple abstractions. This API creates non-materialized aliases of column-groups which, from the cache perspective, pushes arbitrary subsets of columns in dense memory addresses to the memory hierarchy. This, in turn, supports both efficient column- and row-oriented accesses while minimizing CPU cache pollution with unnecessary attributes. Major benefits of Relational Fabric include:

- ✓ **Low Data Complexity:** It allows efficient HTAP processing while maintaining only one layout of the data. There is no need to propagate changes to multiple data copies or convert data among different layouts.
- ✓ **Low Software Complexity:** It reduces data system software complexity by eliminating the requirement to maintain different execution engines. Rather, the execution engine can always assume that only the relevant data will be accessed via *ephemeral columns*.
- ✓ **Efficient Hardware Utilization:** It provides *effortless locality* via shipping only relevant data through the memory hierarchy, alleviating unnecessary data movement, and providing better cache and processor utilization.

As a first instance of Relational Fabric, we have developed *Relational Memory* [59], [70] that utilizes recent advancements in programmable logic [71], and pushes projection to the hardware (§II). The API of Relational Memory is a simple, lightweight programming abstraction, termed *ephemeral vari-*

*ables*, enabling the CPU to access arbitrary data geometries. Relational Memory exploits the inherent parallelism of memory cells to efficiently access data in scattered locations, and uses programmable logic to reorganize and compact it on the fly before pushing it to the CPUs, thus improving locality.

**Simplifying the Data Systems Software Stack.** With Relational Fabric in place, the data system software stack can be significantly simplified [65]. A data system that makes good use of the Relational Fabric would have to drastically simplify its physical design and query optimization components, while the query evaluation engine would now be able to make the most of code generation. Other components would also be significantly affected, especially the transactions manager that would have to implement a multi-version concurrency control (MVCC) approach, and the compression algorithms that should be compatible with scattered data accesses.

In this paper, we present our vision of building *Relational Fabric*, along with the challenges and opportunities of innovation in data systems and some open research questions to fully realize this vision. We begin by discussing how to access arbitrary data geometries and present our initial work on Relational Memory and ephemeral variables (§II). The Relational Fabric vision has several *opportunities for simplicity and innovation* across the data systems stack (§III):

- **Simplify Physical Design:** Relational Fabric will grossly simplify the physical design process. There is no need for creating (physical) vertical partitions anymore. Further, indexes will mostly be useful for workloads with point queries and updates since range queries can be evaluated with column-group accesses very efficiently. Overall, through the Relational Fabric any layout can be achieved via on-the-fly data reorganization.
- **Simplify Query Optimization:** One major challenge in query optimization is its combinatorial nature, that requires to search a vast space to find the optimal query plan. Relational Fabric expands this search space, providing access to any data layout, however, this essentially removes any search constraints. Hence, instead of *solving a combinatorial problem*, we can now *construct* the fastest solution.
- **Efficient MVCC:** The natural way to implement concurrency control using Relational Fabric is MVCC, where there is one source of truth (the base data in row-oriented format), and the ephemeral columns access the correct data using timestamp information associated with MVCC. A big win for Relational Fabric is that it implements timestamp comparisons in hardware, leading to a simple and good-performance implementation of MVCC.

Building the Relational Fabric hardware requires bringing together software and hardware expertise and benefit from the tight integration of programming systems and programmable logic which is increasingly gaining momentum (§IV). Our preliminary Relational Memory design is implemented in such a tightly integrated platform. We further outline our vision for integrating data transformation with the memory controller and extending the processor’s ISA. The ultimate goal is to be able to use Relational Fabric without needing deep expertise in hardware/software co-design.

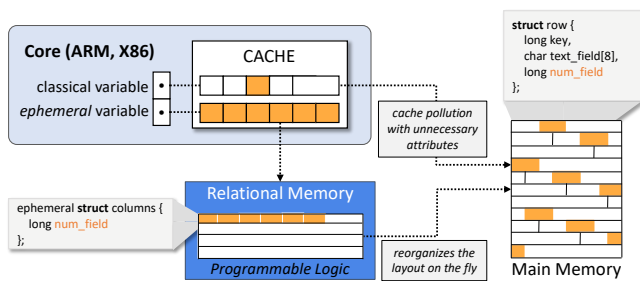


Fig. 2: Relational Memory is an in-memory instance of Relational Fabric that pushes projection closer to data.

## II. ACCESSING ARBITRARY DATA GEOMETRIES

**Motivation.** When considering a row-oriented relational table, accessing a specific (group of) column(s) leads to a strided access pattern, where a few bytes (size of accessed columns) are accessed every hundred or thousand bytes of data (size of the row). This is inefficient since each new row always pulls an entire cache line from memory resulting in lower cache utilization and more data movement than what is required (Figure 1). Further, strided accesses with large strides are not handled well by hardware prefetchers [11]. Column-stores avoid strided accesses by storing each attribute separately. This results in accesses with high locality since only data items strictly required for the final computation are transported from the main memory. However, it also comes at the cost of increased *tuple reconstruction cost* for queries with high projectivity [1], and it is an inefficient layout for inserts and deletes. Instead, we want to achieve *effortless locality* for queries with any projectivity without having unnecessary data movements or any tuple reconstruction cost.

### A. Relational Fabric

To achieve *effortless locality*, we propose to utilize specialized hardware to perform *transparent on-the-fly* data transformation while the base data is stored in a row-oriented layout. This allows efficient updates and inserts on the base data while analytical queries may access the desired columns through *Relational Fabric* [65]. Data transformation and operators like *projection* and *selection* can be pushed closer to data to reduce unnecessary data movement. For memory-based systems, the hardware will sit between memory and the processor, intercept the CPU-originated requests and transparently reorganize the data on-the-fly, making it like the desired data layout already exists in memory. For disk-based systems, the hardware will sit between memory and storage and will be able to transparently reorganize the data on-the-fly, making it like the desired data layout already exists in storage. Operating closer to data enables us to utilize the underlying device parallelism (either within different memory banks [37], [50] or within storage devices like flash-based SSDs [18], [63]) which ensures better overall hardware utilization.

To ensure ease of programmability, the specialized hardware is used via a simple abstraction that allows the data system engineer to request the desired column groups (in case of projection) or rows (in case of selection). For the memory-based design, it will require compiler support to make this

```

1 // layout of the full relational table
2 struct row {
3     long n0;           /* 8 bytes */
4     char s1[12];      /* 12 bytes */
5     char s2[16];      /* 16 bytes */
6     long n1;          /* 8 bytes */
7     long n2;          /* 8 bytes */
8     long n3;          /* 8 bytes */
9     long n4;          /* 8 bytes */
10 };
11
12 // the variable that holds the full relational table
13 struct row table[];
14 // the SQL query to execute
15 char* q = "SELECT SUM(n1 * n4) FROM table WHERE n0 > 10";
16
17 // the ephemeral variable of the SQL query to execute
18 ephemeral struct column_group {
19     long n0;
20     long n1;
21     long n4;
22 };
23 // configuring the ephemeral variable's geometry
24 struct column_group* cg = configure(table, q);
25
26 // executing the query using the ephemeral variable
27 long sum = 0;
28 for (int i = 0; i < rows; i++)
29     if (cg[i].n0 > 10)
30         sum += cg[i].n1 * cg[i].n4;

```

Fig. 3: The ephemeral variable is configured in line 24. Once the CPU accesses it (line 30), the machinery fetches the desired column groups as if they already exist in memory.

fully transparent to the system developer. In addition, in order to make the memory fabric easy to integrate we envision integrating it into the memory controller and modifying the instruction set architecture (ISA). For the disk-based design, the system developer will use a dedicated API that allows to directly request from the proposed machinery to fetch the desired columns from storage. This API will directly feed into the scan operator sending only the relevant data for further query processing. We now discuss our first instance of Relational Fabric and its API that targets in-memory data systems and is termed *Relational Memory*.

### B. Relational Memory

To offer contiguous access to a specific column (or column-group) in memory, we built Relational Memory (RM) [59], [70], which leverages the PLIM paradigm [71] to *create references to data that does not exist in main memory, which the CPU can use as if the data does exist* in main memory. In other words, RM enables accessing the same content in main memory under different strides, but it can be accessed as if it was stored contiguously from the CPU's perspective. Reorganizing data to improve locality minimizes the waste of constrained CPU cache real-estate. In turn, this translates to better efficiency for the query at hand and lower cache pollution. Furthermore, operating closer to the data allows to exploit the inherent parallelism of memory cells. Figure 2 shows a high-level diagram of the proposed design. Relational Memory is located in programmable logic sitting between the memory and the processor. Upon receiving a request, the engine transforms on-the-fly rows to any desired group of columns. The processor directly accesses data in the optimal layout through an abstraction called *ephemeral variables*, which we discuss next.

**Ephemeral Variables.** RM aims to provide the optimal layout for any possible queries; thus, RM executes projection and tuple reconstruction in hardware, and the reconstructed tuples are accessible via *ephemeral variables*, a special type of variables that identifies a specific subset of columns to access. *Ephemeral variables* create memory *aliases* to expose non-contiguous content as if it were contiguous. These transient variables are never instantiated in main memory. Instead, upon accessing such a variable, the underlying machinery generates an on-the-fly projection of the requested columns according to the format that maximizes data locality. We introduce a lightweight software/hardware interface to initialize and deploy the proposed hardware infrastructure. A query needs to perform two tasks to efficiently employ Relational Memory. First, to configure RM at run-time (which needs  $\sim 0.3\mu s$ ), and second to read data from the ephemeral variable sequentially.

Figure 3 shows an example of using ephemeral variables. Imagine a scenario where a complete relational table is loaded into memory and organized as a traditional 2-D array, represented by the array `struct row table[]` (definition of a row provided in lines 2-10). To execute a query (line 16), we create an *ephemeral variable* (lines 18-22), and we register it with RM (line 24) to have direct access to a single column or a group of columns. From the CPU’s viewpoint, accessing this newly created ephemeral variable is equivalent to having direct access to a column group with a packed view of the relevant fields (lines 28-30). This abstraction of *non-materialized in-memory aliases of column-group accesses* grossly simplifies the data management software stack, as we will discuss in Section III. While in our current prototype ephemeral variables are supported by low-level macros, our long-term vision is to add full compiler support by integrating Relational Memory into the memory controller and expanding the ISA, which we discuss in more detail in Section IV. This will allow RM to benefit any data systems with minimal engineering effort.

### III. IMPLICATIONS ON DATA SYSTEMS ARCHITECTURE

We now discuss the opportunities for simplicity and innovation that Relational Fabric brings across the data systems software stack. Figure 4 shows the database components that will be affected by Relational Fabric (green background).

#### A. Physical Design

The physical design process of a database entails (i) the mapping from real-world data to relational data (tables, rows, and columns), and (ii) decisions on how to physically store the relational table regarding normalization, data layout, data partitioning, and indexing. While the first part depends on the nature of the application, the second part aims to optimize performance while respecting some constraints (e.g., space amplification) using workload knowledge. Relational Fabric can radically simplify the second step of physical design because it can allow access to different data geometries without needing to physically store data in that format. We now discuss how the key decisions can be simplified.

**Data Layout and Vertical Partitioning.** One of the most impactful physical design decisions is the data layout: a row-oriented, a column-oriented, or a hybrid layout. This decision

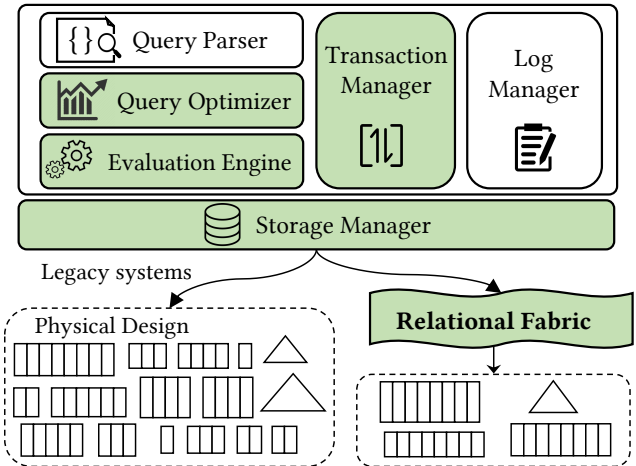


Fig. 4: Relational Fabric enables opportunities for innovation across the entire data system architecture.

is driven by the application requirements since a row-oriented layout is a good fit for transactional applications, while a columnar layout performs well for analytical systems. To bridge the analytical and transactional requirements, hybrid systems often maintain different engines for ingestion and analytics, paying the cost of data layout conversion, duplication, and the necessary additional bookkeeping. Both traditional and hybrid systems use the expected workload to form educated decisions about storing vertically partitioned data. The goal is to collocate columns that are frequently accessed together to reduce unnecessary data movement.

Relational Fabric simplifies the physical design process. Since the query evaluation engine can always request arbitrary column groups, there is no need to confine the physical design space to a predetermined set of vertical partitions. The base row-oriented data can efficiently facilitate updates and ingestions, and analytical queries can be executed optimally using the ephemeral columns to access arbitrary column group. Thus, Relational Fabric eliminates the need to construct fixed physical vertical partitions since the base data can be on-the-fly vertically partitioned by the specialized hardware. Indexes can still be defined on the original row-wise data. However, the usefulness of indexes is now smaller, since range queries can be efficiently evaluated with columnar accesses, so indexes should be used for point queries and point updates. While Relational Fabric does not necessitate any major changes in the storage manager, it simplifies its operation, as it only needs to maintain a single copy of each relation’s data.

**Horizontal Partitioning.** Contrary to vertical partitioning that can happen on-the-fly using Relational Fabric, horizontal partitioning decisions would still need to be evaluated at physical design time. While Relational Fabric does not affect horizontal partitioning, it can be efficiently combined with it, along with sharding which horizontally partitions data based on a sharding key. Another functionality that Relational Fabric can integrate is to handle the communication with storage devices while exposing its simple *ephemeral columns* API to the query. That way, the data system can request the desired column group on a sharding key range, and the Relational Fabric will directly return the corresponding data to the query.

**Indexing.** Indexes help accelerate access times, ensure uniqueness, and allow sorting and clustering. In general, row-store systems employ indexes that are useful for selection queries and data updates. Column-store systems (and some recent row-store systems) [47], [49] use column projections as a special type of index. The strength of Relational Fabric is that it makes such projections possible without having to materialize them. This has deep implications on the data systems architecture because the query engine may access the data at query time using either the base row-oriented data, an index (if it exists), or the desired columns projected from Relational Fabric. While indexes will be mostly beneficial for workloads with point queries and updates, Relational Fabric allows for efficient range queries due to the arbitrary column-group accesses.

### B. Query Optimization and Query Evaluation

When a query is submitted to the system, it is first processed by the query parser via lexical, syntactic, and semantic analysis. Then the initial query plan is processed by the query optimizer to find the best query plan to execute, which is, in turn, submitted for execution to the query evaluation engine. Relational Fabric will not affect the query parser, however, the optimizer and the evaluation engine components will have a significant impact. Row-oriented query execution models execute query plans one-row-at-a-time, which offers good performance for OLTP queries. In contrast, most column-store systems use vectorized execution that progresses through the query plan via processing batches of column data [45]. While Relational Fabric maintains the base data in row format, it makes columns available to the processor on the fly, and, thus, it can naturally support efficient vectorized query execution. A key challenge with query optimization is the combinatorial process of searching a vast space – the DBMS has to select the most efficient evaluation plan based on the cost of each plan. While Relational Fabric seems to increase the search space by enabling access to arbitrary data layouts, it also allows every query to always use the best layout, and overall simplify the problem, by replacing the *combinatorial search* with the *construction* of the query plan that accesses the optimal data fragments. This opens a whole new avenue of research with the potential to speed up query processing: (i) generate the fastest query plan, (ii) revise existing cost models considering Relational Fabric, and (iii) re-evaluate classical single-column, multi-column partitioning cost models.

**Code Generation.** Adaptive systems like Hyper, H<sub>2</sub>O, Actian Vector, Hekaton, MemSQL and others [6], [23], [44], [55], [77] examine the query and decide how data will be accessed by evaluating alternative access plans. The appropriate code is generated by considering the buffered available data layouts. One disadvantage of this approach is the requirement to compile code on the fly which is alleviated by buffering compiled code fragments. The development of Relational Fabric aids code generation in two ways. First, Relational Fabric does not require buffering different layouts since any arbitrary layout can be accessed on the fly. Second, since data layouts are not buffered, Relational Fabric can buffer more code fragments and reuse previously compiled code fragments more

aggressively which allows for better utilization of memory. The query optimizer can now also consider various factors like which code fragments are buffered and which indexes are available. Through the Relational Fabric design we have the opportunity to develop a novel full-fledged hybrid query engine. This engine would seamlessly switch between row-at-a-time and column-at-a-time processing while working on the same base data. This allows for adaptive query processing, where the system dynamically chooses the most efficient processing method based on the characteristics of the query and the data being processed.

### C. Concurrency Control

The transaction manager is responsible for concurrency control. Relational Fabric can naturally support multi-version concurrency control (MVCC). While the base data is in row format, Relational Fabric offers native access to arbitrary data geometries through ephemeral columns. For example, in our in-memory implementation of Relational Memory, we use ephemeral variables to access column groups. We consider all ephemeral variables (or the respective API) as *read-only* columns or column-groups that accelerate analytical queries. The row-wise base data is marked as *read/write* and updates are handled by appending new rows to this base data. For updates and deletion, Relational Fabric uses two timestamp fields for every row to support multiple versions. The first timestamp is set when a row is inserted to mark the beginning of its validity, while the second timestamp is set upon row deletion or replacement by a newer version, marking the end of its validity. Every time the API is accessed, it generates the column groups that contain the valid rows at the time of the query. A key advantage of this approach is that the timestamp comparison can be implemented in hardware, making this implementation simple and performant. By offering the optimal layout and using the timestamps to ship only valid data, Relational Fabric supports MVCC transactions through snapshot isolation.

### D. Compression

The proposed design stores the base data in a row-oriented format hence it can benefit only from specific types of compression. General compression algorithms of the LZ family [92] are frequently used by row stores, however, they are not a natural fit for Relational Fabric since they require fully decompressing your data before you can access separate columns. Delta, dictionary, and Huffman encoding for compression which are popular among state-of-the-art column stores [2], [3], [93] are easily supported by Relational Fabric. Note that these schemes can be used for row-wise data, and hence, they can benefit any groups of columns requested by ephemeral columns. However, the compression schemes under the run-length encoding family cannot be used out of the box. Unlike dictionary and delta encoding, RLE has an expensive decoding step and relies on data, but it is still quite popular among column stores. More research is required to find compression techniques that can benefit both row-oriented and columnar data and allow for direct operation on compressed data.

#### IV. BUILDING RELATIONAL FABRIC

The key feature of Relational Fabric is that a specialized hardware component transforms data on-the-fly. This section presents the implementation details of our first in-memory Relational Fabric instance, termed Relational Memory (RM). We also discuss extending hardware support to more operators, pushing the logic further into the memory controller, and Relational Fabric for storage devices.

##### A. Implementing Relational Memory

RM is an FPGA-based data transformation engine that sits between the processor and the memory and converts data layouts on the fly as shown in Figure 2. Data transformation in RM is performed in line with the instruction stream via fine-grained information on the exact byte-wise location of data items useful for the computation at hand. Our hardware performs the following four key operations: (1) RM receives the intended access stride of the query (that maps the physical addresses of the columns to be accessed) and then issues parallel main memory requests for the target data, (2) RM communicates with memory via an AXI bus [12] and assembles multiple entries into a single *packed* cache line to be sent to the processor, in the meantime, (3) RM captures the CPU requests and (4) transfers the reorganized data upon availability. This abstraction creates *non-materialized aliases of column-groups* which pushes arbitrary subsets of columns to the memory hierarchy. Hence, RM supports both efficient column- and row-oriented accesses while minimizing CPU cache pollution with unnecessary attributes.

Figure 5 shows the high level overview of RM components and their datapath. RM consists of four modules (Figure 5): the *Trapper*, the *Monitor-Bypass* (MB), the *Requestor*, and the *Fetch Unit*, as well as two Scratch Pad Memories (SPMs) used as a *Metadata Buffer* and a *Reorganization Buffer*. RM interacts with the Processing Subsystem (PS) of the FPGA through two primary and one secondary AXI ports. The configuration port allows the DBMS to specify the location and geometry (tuple width and count, size and positions of the requested columns) of the target table at runtime (① in Figure 5). The Trapper works as the interface between the CPU and RM that intercepts read requests (①) from the CPU. It communicates with the Monitor-Bypass (②) to check the availability of the requested data (③). If the data is already in the reorganization buffer, the Monitor-Bypass sends it to the Trapper (④), and then, the CPU receives the data via RM (⑤). If the requested data is not in the buffer, the Monitor-Bypass informs the Requestor about it (A). The Requestor creates descriptors that identify the location of the desired columns (B) based on the DB geometry. The column-extracting module inside the Fetch Unit reads the bus lines that contain useful data, extracts the relevant part (C), and sends the retrieved parts to the Monitor-Bypass (D) so that it can be stored in the reorganization buffer (E). Thus, RM transforms data into the desired layout and minimizes cache pollution.

##### B. Pushing Other Relational Operators

Relational Memory is the first instance of a new class of data systems architectures. Implementing projection in hardware

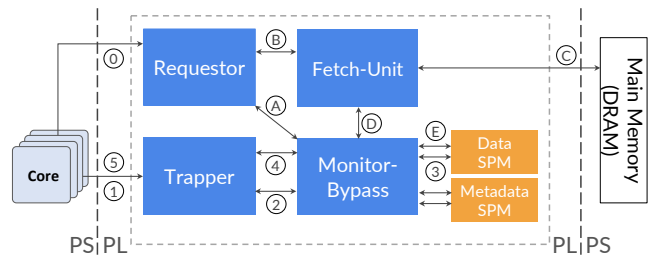


Fig. 5: Abstract overview of Relational Memory components and interconnections with the PS-side.

lays the groundwork for pushing other relational operators to the hardware as well. The Relational Fabric philosophy is that new hardware designs will be adopted if they are simple and general. In other words, a very application-specific design is hard to make its way to mass production. With that in mind, we propose to further reduce unnecessary data movement by transparent on-the-fly data transformation using minimal hardware complexity, by pushing selection and aggregation in the hardware. Both operations are general enough in the sense that they are part of other applications (like operating on matrices and tensors) and have the potential to offer even larger data movement reduction benefits. Implementing selection in Relational Fabric will further alleviate the need for indexes altogether, while aggregation will help both relational and matrix operations. In this design, the ephemeral variables will contain only the required data or the aggregation result, which will be passed through the memory hierarchy ensuring minimal data movement while maintaining the hardware complexity low.

##### C. Pushing RM Further: Relational Memory Controller

Following the aforementioned philosophy, we are developing another instance of Relational Fabric, termed Relational Memory Controller (RMC), where we place RM further closer to memory as shown in Figure 6. Integrating RM into a memory controller (MC) is a game-changer as it will allow for easy adoption of the RM design with minimal development effort. Further, pushing RM into the memory controller maximizes its benefits, since it has low-level access to the physical data placement on the memory DIMMs.

Memory controllers hide the complexity of interfacing with DDR memories, however, they cannot fully exploit the capabilities of DDR memory chips since they have no information about the workload or the application setup. By integrating RM with the memory controller, we pass just enough semantic information about the access patterns that makes it possible to effortlessly offer data locality via exploiting new memory controller commands. This maximizes the memory throughput and reduces unnecessary data movement. To achieve this, we base our implementation on a Red Hat prototype memory controller that models DDR3 DRAM, and expand it to capture first the DDR4 and eventually the DDR5 design. RMC leverages *unaligned* and *interleaved accesses* to implement part of the on-the-fly data transformation in the memory controller while maximizing the parallelism (across memory banks) to further improve performance. Compared to our first prototype of Relational Memory, RMC allows the embedded Relational Memory Engine to see the same memory bandwidth as the

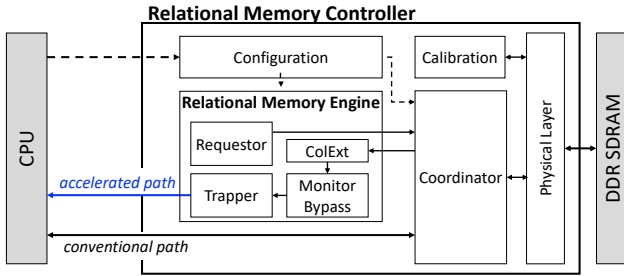


Fig. 6: Architecture of RMC. Pushing RM into the memory controller allows RM to exploit low-level information to maximize memory bandwidth utilization.

CPU, further improving overall efficiency. This is attributed to having direct access to the physical memory and the same clock domain as the memory controller. This is corroborated by our preliminary results that show data transformation through RMC happens at close-to-memory throughput.

Integrating RM into the memory controller has a profound impact. Specifically, it enables rapid research and prototyping of new ideas regarding MCs, something that is increasingly needed with the rapid evolution of interconnects, memory, and reconfigurable hardware.

**Extending the ISA as an RMC Interface.** An Instruction Set Architecture (ISA) is the abstraction between hardware and software. The ISA guarantees that the resulting binary code correctly executes regardless of the toolchain; thus, it helps developers to write and debug software more efficiently [74]. Therefore, integrating RM with ISA, such as RISC-V [14], [86], provides a more valuable interface. The benefits of using RM via ISA are two folds: (1) it will provide a simple interface with no need to understand the details of the underlying hardware and (2) it will simplify the code generation process during the compile time. Thus, RMC along with an ISA extension provides a simple API that can further be beneficial for SQL queries (or other applications benefiting from data transformation) [73], [87].

#### D. Implementing Relational Storage

Following our discussion about building Relational Memory, we propose to develop Relational Fabric in storage devices. Near-storage computation is more challenging than near-memory computation because traditionally storage devices are incapable of performing logic. However, recent modern storage devices like SmartSSD [66] and OpenSSD [16] have processing power that can be exploited to achieve this. We call this approach of pushing computation to storage Relational Storage (RS). RS can be directly implemented in a specialized storage device (i.e., in OpenSSD or SmartSSD) or a programmable logic (i.e., FPGAs), similar to our Relational Memory approach. When implementing *projection* in the custom hardware, a *read-only* analytical query will have access to only a *read-only* version of the optimal layout. In other words, pages will be marked as *read/write* during loading as row-store, while only the columnar pages will be shipped marking as *read-only*. Writes/updates will access the base data in row-format, and the existing *read-only* versions will be marked as *invalid* (similar to an out-of-place data

structure) for the corresponding data. Further, we plan to experiment with a flipped design where the base data is stored in columnar format on storage, allowing for the most efficient compression algorithms like RLE. The processing capabilities of smart SSDs coupled with new custom logic that can be designed in FPGAs that are embedded within the storage device will perform decompression when needed, and tuple reconstruction, removing this burden from the software stack of the database system [19]. This design has great potential since it enables better compression techniques while offloading decompression and tuple reconstruction to the storage. Further research is required to reveal its challenges and opportunities.

Overall, Relational Storage can significantly reduce end-to-end latency by massively reducing data movement. In contrast to RM, it is possible to push operators like selection and aggregation by utilizing the processing power of in-storage custom logic. Exploiting the internal parallelism of the storage device [63], [64] can enhance performance. Furthermore, the software stack will be redesigned to take advantage of near-storage computation for better query processing and optimization in contemporary storage devices.

#### E. Practicality & Feasibility of Relational Fabric

With the tapering of Moore’s Law, hardware specialization is gaining popularity over the past decade [38], [83]. FPGAs, once limited to specific applications, are now ubiquitous across various domains. Our first Relational Memory prototype is developed on a PS-PL platform, however, with reprogrammable and custom hardware being tightly integrated with CPUs increasingly more often, Relational Memory can be deployed to standard servers. For example, Enzian is a research computer built by the Systems Group at ETH Zurich with a big server-class CPU closely coupled to a large FPGA [21]. Further, manufacturers, such as Intel and AMD, are investing substantially in FPGA technology (e.g., with the acquisitions of Altera and Xilinx, respectively), thus driving the momentum towards systems with tightly integrated CPUs and FPGAs [7], [39], [80]. Moreover, recent work on reprogrammable memory controllers [17], [75] further fuels our vision for integrating data reorganization within the memory controller, as discussed in Section IV-C. Finally, the emergence of Smart SSDs and PIM (Processing-in-memory) also adds reprogrammable logic closer to data [16], thus being able to offload processing tasks, like compression [19].

Overall, the current hardware trends (widespread adoption of reprogrammable hardware, slowdown in Moore’s Law, and dark silicon) show great promise for Relational Fabric to improve memory utilization and performance.

## V. EXPERIMENTAL RESULTS ON RELATIONAL MEMORY

We now present selected experimental results of RM showing that it outperforms direct row-wise and direct columnar accesses by offering the optimal layout to any query [70].

**Target Platform.** The full-stack prototype of RM is implemented on a Xilinx Zynq UltraScale+ MPSoC platform [89] which consists of heterogeneous Systems-on-Chip (SoC) where a traditional processing system (PS) is tightly associated with a programmable logic (PL), i.e., an FPGA. The PS

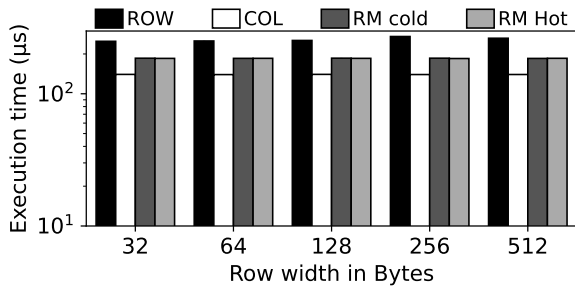


Fig. 7: Execution time for a query projecting 3 columns. RM outperforms row-oriented direct memory accesses while providing comparable performance to column-stores.

equips with 4 Cortex-A53 1.5 GHz cores, each with a private 32+32 KB L1 I+D cache and sharing a unified 1 MB L2 cache. PL side, RM prototype, is constrained to 100 MHz. In order to compare the performance of RM to the row-store (ROW) and the column-store (COL), we custom implement an in-memory row-store following the volcano-style processing model (tuple-at-a-time) and an in-memory column-store following the column-at-a-time processing model. We run two sets of experiments for RM: *hot* (when the targeted data is ready in the Reorganization Buffer) and *cold* (when the targeted data is not yet in the Reorganization Buffer).

**RM Enables Native Columnar Accesses.** Our initial experiment illustrates that RM can efficiently access and retrieve individual columns while the base data is stored as a row-store. We evaluate a query that projects three non-contiguous columns having the same width, while varying the row size from 32B to 512B. Figure 7 shows the execution time where we compare the query time when data is accessed directly (i) from the in-memory row-store (ROW), (ii) through RM *Hot* and (iii) RM *Cold* accesses, and (iv) using the in-memory columnar format (COL). We observe that RM outperforms ROW regardless of accesses being *cold* or *hot* which shows RM can provide optimal data layout. In fact, RM achieves an average latency that is comparable to columnar accesses. Further, RM’s performance remains virtually unchanged irrespective of whether the data is ready in the Reorganization Buffer, meaning that RM achieves virtually the same benefit for the cold accesses. We also experiment with varying column width where we observe RM to outperform ROW consistently and have comparable performance as COL for smaller column width while RM outperforms COL for larger column-widths [70]. This experiment shows RM’s ability to efficiently access column groups without unnecessary data movement.

**RM Shines for Queries with High Projectivity.** In this experiment, we vary the projectivity from 1 to 11 columns for 4-byte wide columns and 64-byte wide rows, as shown in Figure 8. For any projectivity, RM outperforms direct row-wise accesses since RM provides the optimal layout that minimizes cache pollution. When the projectivity is low ( $\leq 4$ ), columnar accesses are faster since the tuple materialization cost is still small and the prefetcher can efficiently support up to four parallel sequential accesses. As projectivity becomes larger than four columns, however, RM starts to outperform direct columnar accesses since the columnar accesses exceed the capacity of prefetcher as observed in Figure 7.

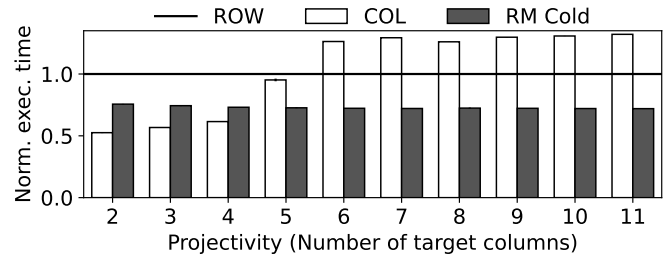
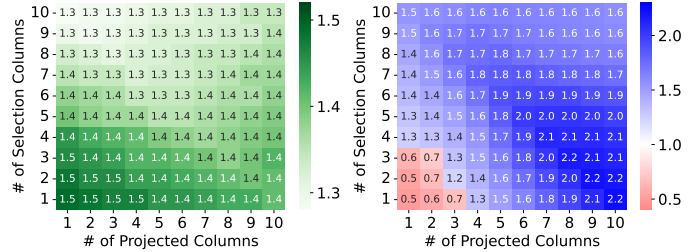


Fig. 8: RM outperforms row-wise memory accesses irrespective of projectivity, while RM shows better performance than columnar accesses for projecting more than 5 columns.



(a) Speedup - RM vs Row (b) Speedup - RM vs Columnar

Fig. 9: (a) RM always outperforms row access. (b) RM dominates column store for larger number of columns ( $> 4$ ).

**RM vs Col-Store vs Row-Store.** The experiment shown in Figures 9a and 9b compares the performance of RM with direct row-wise and columnar accesses while varying both the number of projected columns ( $x$ -axis) and the columns used for the selection ( $y$ -axis). They both range from 1 to 10 columns. Figure 9a shows the speedup of RM compared to the direct row-wise accesses. Similarly to the previous experiment, RM consistently outperforms the direct row-wise access by 1.3-1.5 $\times$ . In contrast, direct columnar access achieves better performance than RM when the number of columns used for projection and selection is less than four, as shown in the lower left corner of Figure 9b. As the number of columns in a query increases, RM outperforms the columnar accesses due to the tuple materialization cost and the diminished prefetching benefits of column-store. Overall, RM achieves better performance than direct row-wise accesses for any number of target columns, while RM outperforms a columnar layout only when the number of target columns is large enough ( $> 4$ ).

**RM Shows Stable Performance for Practical Queries.** In order to evaluate RM in a practical environment, we execute  $Q_1$  and  $Q_6$  from TPC-H [84] while varying the data size. RM supports arbitrary data sizes even with a small data memory of 2 MB on the FPGA by refilling it whenever it is full. Figure 10 shows the running time of  $Q_1$  and  $Q_6$  on tables from 11 MB to 692 MB. Since we choose the data size based on the size of target columns (shown in the parentheses of  $x$ -axis), the range of data sizes varies for  $Q_1$  and  $Q_6$ . For  $Q_1$ , the execution time is similar for all layouts, as shown in Figure 10a. This is because executing CPU-intensive operations in  $Q_1$  dominates the data movement cost. On the other hand, for queries such as  $Q_6$  where data movement is the bottleneck, RM accelerates the execution time by offering the optimal layout (Figure 10b).



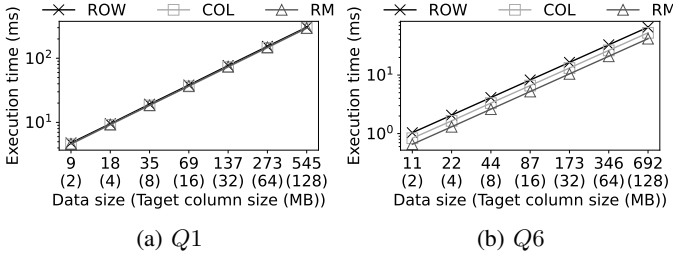


Fig. 10: RM shows better performance than direct row-wise or columnar accesses in practical queries such as TPC-H Q1 and Q6 regardless of the data size.

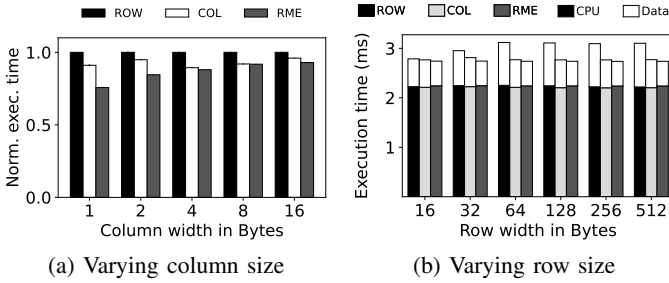


Fig. 11: RM performs join faster than traditional row-store join by minimizing data movement.

**RM Reduces Data Movement for Joins.** Queries that involve multiple tables (e.g., join), RM projects only the relevant columns, encompassing the join attributes and the target columns for projection or selection. This reduces unnecessary data movement and minimizes cache pollution. In this particular experiment, we perform a join operation using a state-of-the-art hash-based join with a single-pass hash table generation. This generated hash table is subsequently probed by the second relation where half of the entries of the outer relation have a match in the inner relation. Figure 11a shows the normalized query latency while varying target column sizes. We observe that RM gives a benefit between 5% and 10% compared to row-wise access when performing joins. The figure also shows that RM outperforms the columnar join as well, providing up to 10% improvement. Figure 11b compares the execution time of this query while varying the row sizes. RM reduces the total runtime by up to 12% depending on the row size. A key takeaway from the figure is that the CPU overhead (depicted by the solid portion of the bars) of hashing represents the predominant portion of the runtime which is constant across RM, ROW, and COL. However, as the row size increases, RM can optimize the data movement by up to 41% because of its lower cache misses, better-strided accesses, and higher cache utilization.

**RM Offers Any Layout via Data Reorganization.** Figure 12a highlights the capability of RM to on-the-fly vertically partition. In this experiment, we use a table with 128 columns of 4B each (512-byte rows). We use four columns for selection and another four columns for the final projection. We compare six approaches: a full row-store (ROW), a full column-store (COL), RM that accesses exactly all 8 columns (RM 8), RM that accesses separately the four columns for selection first, and then the four columns of projection (RM 4+4), and the optimal layout where the 8 and the 4+4 column-groups

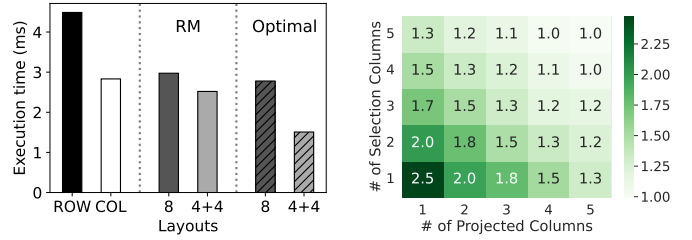


Fig. 12: (a) RM offers any layout efficiently, thus simplifying physical design. (b) Relational Storage pushes projection to storage, thus improving query processing.

are already created as projections (Optimal 8 and Optimal 4+4). RM 8 outperforms ROW and matches COL, while RM 4+4 further reduces the data movement avoiding some more unnecessary accesses from the projected columns. When compared with having the optimal layout readily available, RM offers similar performance with Optimal 8, while some additional performance optimizations seem to be possible for the RM 4+4 vs Optimal 4+4 case. Overall, we see that RM can virtually mimic having any layout available having significant impact on simplifying physical design.

**Potential of Relational Storage.** In our last experiment, we quantify the headroom of improvement for Relational Storage. We vary the number of projection and selection columns, and we compare the performance of a row store with 10 columns (4-bytes each) against the optimal data layout containing only the relevant columns. Figure 12b illustrates the potential benefit (up to 2.5× for low projectivity). Building on-the-fly data transformation capabilities within storage (e.g., using Smart SSDs) will unlock this performance benefit for any query that can benefit without compromising the performance of other queries.

Overall, our proof-of-concept prototype Relational Memory shows better or comparable performance compared to in-memory row- or column-store for various queries by offering data transformation on the fly (more experimental results are available in our conference paper [70]). These experimental results of RME further fuel our vision for Relational Fabric.

## VI. RELATED WORK

**Hybrid Layouts.** Many HTAP systems like SAP HANA [31], Oracle TimesTen [48], MemSQL [77], BatchDB [52], and L-store [72] follow the *one size does not fit all* rule [81], hence, they use the row-format to ingest data and then convert it to columnar-format for analytical processing [62]. The *optimal* layout is more often neither a column-store nor a row-store [6]. On the other hand, systems like H<sub>2</sub>O [6], Hyper [44], Peloton [13], and OctopusDB [24] use adaptive layouts depending on the query patterns. All these systems need to store multiple layouts of the data and convert between formats which increases the complexity, materialization overhead, and maintenance cost.

**Hardware Specialization.** There have been many efforts to utilize specialized hardware for data management systems [30], [40]. We categorize the developed specialized hard-

ware by its objectives. The first line of specialized hardware is to accelerate particular DBMS operators such as selection [82], aggregation [22], compression [67], decompression [29], data partitioning [42], sort [91], group by [4], and join [35], [90]. Secondly, we classify attempts to offload the SQL query itself or the subset of queries [60], [61], [78], [88]. Due to the inflexible nature of hardware, these approaches' main limitation lies in supporting ad-hoc queries. A third class is query accelerators accessing non-local memory aiming to reduce data movement [5], [9], [46], [68], [79].

Contrary to the aforementioned related work or the Processing-In-Memory (PIM) approach [51], [76], the Relational Fabric paradigm does not aim to implement complex logic near memory/storage, nor to change the physical memory/storage hardware (e.g., memory or flash cells). Rather, Relational Fabric sits between the query execution engine and the data, and offers a light-weight layer that performs on-the-fly transparent data transformation into the optimal layout for the query in question without materializing it. Our first Relational Fabric instance, RM, sits between the CPU and memory and transparently transforms data into the optimal layout that does not exist in main memory. Therefore, any ad hoc queries can be accelerated with no data duplication. Furthermore, RM does not require any modification of the memory hierarchy unlike PIM and is fully implemented on commercially available platforms [8], [28], [56], [89].

To develop Relational Fabric for storage devices, we capitalize on recent advancements in computational SSDs (OpenSSD [16], SmartSSD [66]). These SSDs have processing power in the flash controller that allows programmability which can be utilized to enable highly efficient SSD execution [26]. There have been several works on performing near-data processing in SSDs [25], [34], [41], [85] leveraging their computational capability which can also aid the development of Relational Fabric in modern storage devices.

## VII. DISCUSSION & FUTURE WORK

### A. Relational Fabric on Disaggregated Storage

With the network interfaces being upgraded to 50Gbps or 100Gbps in many (private or public) cloud deployments and many applications being storage space limited [27], it is increasingly common to leverage disaggregated storage. In this context, Relational Fabric can offer ideal data movement over the network, ensuring that despite not having the optimal layout on storage, only the desired columns are shipped from the storage node to the compute node. This can be implemented using the computational capabilities of SSDs that are attached to the storage node as discussed in Section IV-D or through a software approach that is implemented in the storage node. Depending on the specific setup (network interface speed, local storage throughput, data geometry), the Relational Fabric approach can deliver superior performance along with better device and network utilization, all while using hardware capabilities that are widely available. It is worth experimenting with both the hardware-based and the software-based solutions to uncover the tradeoffs between the two approaches, which is left for future work.

### B. Generalized Data Transformation

Moving further than relational data, we plan to build a generalized *Data Transformation Unit* (DTU) capable of enhancing the locality of spatiotemporal data in arbitrarily high-dimensional objects, with *tensors* being a prime example. This will allow ML-based applications to enhance their data access locality and, ultimately, their efficiency. To achieve this, we are developing a resource-efficient data transformation algebra that can handle complex, multi-dimensional access patterns. The DTU we are currently building uses the requested address to identify and retrieve an  $\mathcal{N}$ -dimensional memory object, extracting details about the access pattern. It then accesses the desired (potentially scattered) data points within the tensor and consolidates them into a *restructured cache line*, ready for delivery to CPUs in response to the cache-line refill request. During writes, the DTU starts from an altered cache line aiming to distribute the (again, potentially scattered) updates to their positions in the base data. With more high-throughput cache coherent interfaces between CPUs and FPGA being available, integrating DTUs is now efficient and feasible as a practical solution. Recent findings [69] indicate a 30% higher data rate compared to the original PLIM [71] approach by leveraging cache coherence ports for cache line refills.

## VIII. OPEN QUESTIONS

In addition to the opportunities for simplicity and innovation discussed in this paper, the Relational Fabric vision has several *open research challenges* that require further investigation.

**Q1.** *Is data transformation (projection) enough?* Relational Fabric is a layer that offers transparent and efficient projection that leads to the benefits we discussed above. Further, data transformation has great potential for other data-intensive applications over multi-dimensional data (matrix/tensor slicing and vectorized operations on matrix/tensor slices). In addition, there have been several recent efforts to implement more complex logic near or within memory. We purposefully avoid this path because it increases the hardware complexity and specialization, making it less general and, thus, to our understanding, less appealing for real-life use and deployment. However, it remains an open question whether more logic can be implemented between the memory and the processor. Overall, our thesis is that any added logic should benefit many different applications to be ultimately viable.

**Q2.** *How does Relational Fabric interact with compression?* While delta and dictionary compression schemes can be used as a starting point, we also believe it is worth investigating new compression schemes that can be applied to row-oriented data and allow for on-the-fly vertical partitioning and potentially allow for operating on compressed data.

**Q3.** *Can you have Relational Fabric both on storage and in memory?* The vision we outline assumes that the Relational Fabric is implemented either in memory on storage, depending on the use-case. However, a scheme that uses Relational Fabric in both storage and memory may also be interesting. Consider that the two fabrics may play different roles. For example, the storage one can convert from compressed columns to rows in memory, and the in-memory one can allow the processor

to access arbitrary column groups. We believe that more investigation in this direction is warranted.

## IX. CONCLUSION

In this paper, we present our vision of Relational Fabric, a new lightweight specialized hardware fabric that offers *effortless locality* by accessing arbitrary data layouts from row-oriented base data without any data duplication. Relational Fabric will simplify data and software complexity, and it will enable efficient hardware utilization and true HTAP processing. We outline the principles, goals, and impact of Relational Fabric, and as a proof-of-concept, we present its first instance, Relational Memory that uses reprogrammable hardware to implement logic between the memory and the processor. Relational Memory on-the-fly converts rows to arbitrary column groups, alleviating the need to vertically partition data. We also outline the necessary steps toward building Relational Fabric in memory, discuss its opportunities for innovation in data systems architecture in physical design, query processing, and concurrency control, and some open questions that require further research. We further discuss building Relational Fabric in computational SSDs by developing Relational Storage. Developing Relational Fabric in memory and storage has the potential to be a paradigm shift where different specialized hardware components (in memory and storage) can synergistically turn data processing more efficient, scalable, and resource-efficient for data-intensive applications.

## ACKNOWLEDGMENT

We thank the reviewers for their constructive feedback and Teona Bagashvili for her assistance in rerunning some experiments. This work is funded by a RedHat Research Incubation Award, a RedHat Research Award, a Cisco gift, and partially supported by the National Science Foundation (NSF) under grant number IIS-2144547, CCF-2008799 and CNS-2238476. The opinions expressed in this publication are those of the authors and not necessarily reflective of NSF views. Denis Hoornaert was supported by the Chair for Cyber-Physical Systems in Production Engineering at TUM and the Alexander von Humboldt Foundation.

## REFERENCES

- [1] D. J. Abadi, P. A. Boncz, and S. Harizopoulos, "Column-oriented Database Systems," *PVLDB*, vol. 2, no. 2, pp. 1664–1665, 2009.
- [2] D. J. Abadi, P. A. Boncz, S. Harizopoulos, S. Idreos, and S. Madden, "The Design and Implementation of Modern Column-Oriented Database Systems," *Found. Trends Databases*, vol. 5, no. 3, pp. 197–280, 2013.
- [3] D. J. Abadi, S. Madden, and M. Ferreira, "Integrating Compression and Execution in Column-oriented Database Systems," *SIGMOD*, 2006.
- [4] I. Absalyamov, P. Budhkar, S. Windh, R. J. Halstead, W. A. Najjar, and V. J. Tsotras, "FPGA-accelerated group-by aggregation using synchronizing caches," *DAMON*, 2016.
- [5] M. K. Aguilera, K. Keeton, S. Novakovic, and S. Singhal, "Designing Far Memory Data Structures: Think Outside the Box," *HotOS*, 2019.
- [6] I. Alagiannis, S. Idreos, and A. Ailamaki, "H2O: A Hands-free Adaptive Store," *SIGMOD*, 2014.
- [7] P. Alcorn, "AMD to Fuse FPGA AI Engines Onto EPYC Processors," 2022. [Online]. Available: <https://www.tomshardware.com/news/amd-to-fuse-fpga-ai-engines-onto-epyc-processors-arrives-in-2023>
- [8] G. Alonso, T. Roscoe, D. Cock, M. Ewaida, K. Kara, D. Korolija, D. Sidler, and Z. Wang, "Tackling Hardware/Software co-design from a database perspective," *CIDR*, 2020.
- [9] E. Amaro, C. Branner-Augmon, Z. Luo, A. Ousterhout, M. K. Aguilera, A. Panda, S. Ratnasamy, and S. Shenker, "Can far memory improve job throughput?" *EuroSys*, 2020.
- [10] R. Appuswamy, M. Karpathiotakis, D. Porobic, and A. Ailamaki, "The Case For Heterogeneous HTAP," *CIDR*, 2017.
- [11] ARM, "Arm Cortex-A53 MPCore Processor Technical Reference Manual," Tech. Rep., 2018. [Online]. Available: <https://developer.arm.com/documentation/ddi0500/j>
- [12] —, "AMBA AXI and ACE Protocol Specification," [https://static.docs.arm.com/ihi0022/g/IHI0022G\\_amba\\_axi\\_protocol\\_spec.pdf](https://static.docs.arm.com/ihi0022/g/IHI0022G_amba_axi_protocol_spec.pdf), 2019.
- [13] J. Arulraj, A. Pavlo, and P. Menon, "Bridging the Archipelago between Row-Stores and Column-Stores for Hybrid Workloads," *SIGMOD*, 2016.
- [14] K. Asanović and D. A. Patterson, "Instruction sets should be free: The case for risc-v," *EECS Department, University of California, Berkeley, Tech. Rep. UCB/EECS-2014-146*, 2014.
- [15] R. Barber, G. M. Lohman, V. Raman, R. Sidle, S. Lightstone, and B. Schiefer, "In-Memory BLU Acceleration in IBM's DB2 and dashDB: Optimized for Modern Workloads and Hardware Architectures," *ICDE*, 2015.
- [16] M. Bjørling, J. González, and P. Bonnet, "LightNVM: The Linux open-channel SSD subsystem," *FAST*, 2019.
- [17] M. N. Bojnordi and E. Ipek, "PARDIS: A programmable memory controller for the DDRx interfacing standards," *ISCA 2012*, 2012.
- [18] F. Chen, B. Hou, and R. Lee, "Internal Parallelism of Flash Memory-Based Solid-State Drives," *TOS*, vol. 12, no. 3, pp. 13:1–13:39, 2016.
- [19] X. Chen, N. Zheng, S. Xu, Y. Qiao, Y. Liu, J. Li, and T. Zhang, "KallaxDB: A Table-less Hash-based Key-Value Store on Storage Hardware with Built-in Transparent Compression," *DAMON*, 2021.
- [20] Cisco, "Cisco Global Cloud Index: Forecast and Methodology, 2016–2021," *White Paper*, 2018.
- [21] D. Cock, A. Ramdas, D. Schwyn, M. Giardino, A. Turowski, Z. He, N. Hossle, D. Korolija, M. Licciardello, K. Martsenko, R. Achermann, G. Alonso, and T. Roscoe, "Enzian: an open, general, CPU/FPGA platform for systems software research," *ASPLOS*, 2022.
- [22] C. Dennl, D. Ziener, and J. Teich, "Acceleration of SQL Restrictions and Aggregations through FPGA-Based Dynamic Partial Reconfiguration," *FCCM*, 2013.
- [23] C. Diaconu, C. Freedman, E. Ismert, P.-A. Larson, P. Mittal, R. Stonecipher, N. Verma, and M. Zwillig, "Hekaton: SQL server's memory-optimized OLTP engine," *SIGMOD*, 2013.
- [24] J. Dittrich and A. Jindal, "Towards a One Size Fits All Database Architecture," *CIDR*, 2011.
- [25] J. Do, Y.-S. Kee, J. M. Patel, C. Park, K. Park, and D. J. DeWitt, "Query processing on smart SSDs: opportunities and challenges," *SIGMOD*, 2013.
- [26] J. Do, S. Sengupta, and S. Swanson, "Programmable solid-state storage in future cloud datacenters," *CACM*, vol. 62, no. 6, pp. 54–62, 2019.
- [27] S. Dong, S. S. P. S. Pan, A. Ananthabhotla, D. Ekambaram, A. Sharma, S. Dayal, N. V. Parikh, Y. Jin, A. Kim, S. Patil, J. Zhuang, S. Dunster, A. Mahajan, A. Chelluri, C. Datye, L. V. Santana, N. Garg, and O. Gawde, "Disaggregating RocksDB: A Production Experience," *Proc. ACM Manag. Data*, vol. 1, no. 2, pp. 192:1—192:24, 2023.
- [28] ETHZ, "Enzian Systems," <http://enzian.systems/>, 2021.
- [29] J. Fang, J. Chen, J. Lee, Z. Al-Ars, and H. P. Hofstee, "A Fine-Grained Parallel Snappy Decompressor for FPGAs Using a Relaxed Execution Model," *FCCM*, 2019.
- [30] J. Fang, Y. T. B. Mulder, J. Hidders, J. Lee, and H. P. Hofstee, "In-memory database acceleration on FPGAs: a survey," *VLDBJ*, vol. 29, no. 1, pp. 33–59, 2020.
- [31] F. Färber, N. May, W. Lehner, P. Große, I. Müller, H. Rauhe, and J. Dees, "The SAP HANA Database – An Architecture Overview," *IEEE DEBULL*, vol. 35, no. 1, pp. 28–33, 2012.
- [32] Gartner, "Gartner Says 8.4 Billion Connected "Things" Will Be in Use in 2017, Up 31 Percent From 2016," <https://tinyurl.com/Gartner2020>, 2017.
- [33] Google, "Cloud TPU," <https://cloud.google.com/tpu/>, 2017.
- [34] B. 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," *ISCA*, 2016.
- [35] R. J. Halstead, I. Absalyamov, W. A. Najjar, and V. J. Tsotras, "FPGA-based Multithreading for In-Memory Hash Joins," *CIDR*, 2015.
- [36] J. Hamilton, "Tensor Processing Unit," <https://perspectives.mvdirona.com/2017/04/tensor-processing-unit/>.
- [37] M. Hassan, "Reduced latency DRAM for multi-core safety-critical real-time systems," *Real-Time Systems*, vol. 56, no. 2, pp. 171–206, 2020.

- [38] J. L. Hennessy and D. A. Patterson, "A new golden age for computer architecture," *CACM*, vol. 62, no. 2, pp. 48–60, 2019.
- [39] Intel, "Intel Completes Acquisition of Altera," 2015. [Online]. Available: <https://www.intel.com/news-events/press-releases/detail/302/intel-completes-acquisition-of-altera>
- [40] Z. István, "The Glass Half Full: Using Programmable Hardware Accelerators in Analytics," *IEEE DEBULL*, vol. 42, no. 1, pp. 49–60, 2019.
- [41] Y. Jin, H.-W. Tseng, Y. Papakonstantinou, and S. Swanson, "KAML: A Flexible, High-Performance Key-Value SSD," *2017 IEEE International Symposium on High Performance Computer Architecture, HPCA 2017, Austin, TX, USA, February 4-8, 2017*, 2017.
- [42] K. Kara, J. Giceva, and G. Alonso, "FPGA-based Data Partitioning," *SIGMOD*, 2017.
- [43] M. Karpathiotakis, M. Branco, I. Alagiannis, and A. Ailamaki, "Adaptive Query Processing on RAW Data," *PVLDB*, vol. 7, no. 12, pp. 1119–1130, 2014.
- [44] A. Kemper and T. Neumann, "HyPer: A Hybrid OLTP & OLAP Main Memory Database System Based on Virtual Memory Snapshots," *ICDE*, 2011.
- [45] T. Kersten, V. Leis, A. Kemper, T. Neumann, A. Pavlo, and P. A. Boncz, "Everything You Always Wanted to Know About Compiled and Vectorized Queries But Were Afraid to Ask," *PVLDB*, vol. 11, no. 13, pp. 2209–2222, 2018.
- [46] D. Korolija, D. Koutsoukos, K. Keeton, K. Taranov, D. S. Milojicic, and G. Alonso, "Farview: Disaggregated Memory with Operator Off-loading for Database Engines," *CIDR*, 2022.
- [47] T. Lahiri, S. Chavan, M. Colgan, D. Das, A. Ganesh, M. Gleeson, S. Hase, A. Holloway, J. Kamp, T.-H. Lee, J. Loaiza, N. Macnaughton, V. Marwah, N. Mukherjee, A. Mullick, S. Muthulingam, V. Raja, M. Roth, E. Soylemez, and M. Zait, "Oracle Database In-Memory: A Dual Format In-Memory Database," *ICDE*, 2015.
- [48] T. Lahiri, M.-A. Neimat, and S. Folkman, "Oracle TimesTen: An In-Memory Database for Enterprise Applications," *IEEE DEBULL*, vol. 36, no. 2, pp. 6–13, 2013.
- [49] A. Lamb, M. Fuller, and R. Varadarajan, "The Vertica Analytic Database: C-Store 7 Years Later," *PVLDB*, vol. 5, no. 12, pp. 1790–1801, 2012.
- [50] C. J. Lee, V. Narasiman, O. Mutlu, and Y. N. Patt, "Improving memory bank-level parallelism in the presence of prefetching," *MICRO*, 2009.
- [51] G. Loh, N. Jayasena, M. Oskin, M. Nutter, D. Roberts, M. Meswani, D. P. Zhang, and M. Ignatowski, "A Processing in Memory Taxonomy and a Case for Studying Fixed-function PIM," *WoNDP*, 2013.
- [52] D. Makreshanski, J. Giceva, C. Barthels, and G. Alonso, "BatchDB: Efficient Isolated Execution of Hybrid OLTP+OLAP Workloads for Interactive Applications," *SIGMOD*, 2017.
- [53] R. Mancuso, S. Roozkhosh, D. Hoornaert, J. H. Mun, T. I. Papon, and M. Athanassoulis, "Software-Shaped Platforms," *Proceedings of Cyber-Physical Systems and Internet of Things Week 2023, CPS-IoT Week 2023 Workshops, San Antonio, TX, USA, May 9-12, 2023*, 2023.
- [54] N. May, A. Böhm, and W. Lehner, "SAP HANA - The Evolution of an In-Memory DBMS from Pure OLAP Processing Towards Mixed Workloads," *BTW*, 2017.
- [55] P. Menon, A. Pavlo, and T. C. Mowry, "Relaxed Operator Fusion for In-Memory Databases: Making Compilation, Vectorization, and Prefetching Work Together At Last," *PVLDB*, vol. 11, no. 1, pp. 1–13, 2017.
- [56] Microsemi — Microchip Technology Inc., "PolarFire SoC - Lowest Power, Multi-Core RISC-V SoC FPGA," July 2020. [Online]. Available: <https://www.microsemi.com/product-directory/soc-fpgas/5498-polarfire-soc-fpga>
- [57] Microsoft, "Project Catapult," <https://www.microsoft.com/en-us/research/project/project-catapult/>, 2017.
- [58] C. Mohan, "Hybrid Transaction and Analytics Processing (HTAP): State of the Art," *BIRTE*, 2016.
- [59] J. H. Mun, K. Karatsenidis, T. I. Papon, S. Roozkhosh, D. Hoornaert, U. Drepper, A. Sanaullah, R. Mancuso, and M. Athanassoulis, "On-the-fly Data Transformation in Action," *PVLDB*, vol. 16, no. 12, pp. 3950–3953, 2023.
- [60] M. Najafi, M. Sadoghi, and H.-A. Jacobsen, "Flexible Query Processor on FPGAs," *PVLDB*, vol. 6, no. 12, pp. 1310–1313, 2013.
- [61] Oracle, "DAX," <https://blogs.oracle.com/linux/post/oracle-data-analytics-accelerator-dax-for-sparc>, 2021.
- [62] F. Özcan, Y. Tian, and P. Tözün, "Hybrid Transactional/Analytical Processing: A Survey," *SIGMOD*, 2017.
- [63] T. I. Papon and M. Athanassoulis, "A Parametric I/O Model for Modern Storage Devices," *DAMON*, 2021.
- [64] —, "The Need for a New I/O Model," *CIDR*, 2021.
- [65] T. I. Papon, J. H. Mun, S. Roozkhosh, D. Hoornaert, A. Sanaullah, U. Drepper, R. Mancuso, and M. Athanassoulis, "Relational Fabric: Transparent Data Transformation," *ICDE*, 2023.
- [66] K. Park, Y.-S. Kee, J. M. Patel, J. Do, C. Park, and D. J. DeWitt, "Query Processing on Smart SSDs," *IEEE DEBULL*, vol. 37, no. 2, pp. 19–26, 2014.
- [67] W. Qiao, J. Du, Z. Fang, M. Lo, M.-C. F. Chang, and J. Cong, "High-Throughput Lossless Compression on Tightly Coupled CPU-FPGA Platforms," *FCCM*, 2018.
- [68] A. Redshift, "Aqua (advanced query accelerator) for amazon redshift," 2021. [Online]. Available: <https://aws.amazon.com/redshift/features/aqua/>
- [69] S. Roozkhosh, D. Hoornaert, and R. Mancuso, "CAESAR: Coherence-Aided Elective and Seamless Alternative Routing via on-chip FPGA," in *Proceedings of the 43rd IEEE Real-Time Systems Symposium (RTSS)*, Houston, TX, USA, 2022.
- [70] S. Roozkhosh, D. Hoornaert, J. H. Mun, T. I. Papon, A. Sanaullah, U. Drepper, R. Mancuso, and M. Athanassoulis, "Relational Memory: Native In-Memory Accesses on Rows and Columns," *EDBT*, 2023.
- [71] S. Roozkhosh and R. Mancuso, "The Potential of Programmable Logic in the Middle: Cache Bleaching," *RTAS*, 2020.
- [72] M. Sadoghi, S. Bhattacharjee, B. Bhattacharjee, and M. C. Anim, "L-Store: A Real-time OLTP and OLAP System," *EDBT*, 2018.
- [73] B. Salami, G. A. Malazgirt, O. Arcas-Abella, A. Yurdakul, and N. Sönmez, "AxleDB: A novel programmable query processing platform on FPGA," *Microprocessors and Microsystems*, vol. 51, pp. 142–164, 2017.
- [74] A. Sanaullah, "Risc-v for fpgas: benefits and opportunities," *Red Hat Research Quarterly*, no. May, 2022.
- [75] B. C. Schwedock, P. Yoovidhya, J. Seibert, and N. Beckmann, "täkö: a polymorphic cache hierarchy for general-purpose optimization of data movement," *ISCA*, 2022.
- [76] V. Seshadri, T. Mullins, A. Boroumand, O. Mutlu, P. B. Gibbons, M. A. Kozuch, and T. C. Mowry, "Gather-scatter DRAM: in-DRAM address translation to improve the spatial locality of non-unit strided accesses," *MICRO*, 2015.
- [77] N. Shangunov, "The MemSQL In-Memory Database System," *IMDM*, 2014.
- [78] D. Sidler, M. Owaida, Z. István, K. Kara, and G. Alonso, "doppioDB: A hardware accelerated database," *FPL*, 2017.
- [79] D. Sidler, Z. Wang, M. Chiosa, A. Kulkarni, and G. Alonso, "StRoM: smart remote memory," *EuroSys*, 2020.
- [80] A. D. C. Solution, "Xilinx Deep Learning Solution on AMD EPYC™ Processors," 2019. [Online]. Available: <https://www.amd.com/content/dam/amd/en/documents/epyc-business-docs/white-papers/xilinx-deep-learning-solution-on-amd-epyc-processors.pdf>
- [81] M. Stonebraker and U. Cetintemel, "'One Size Fits All': An Idea Whose Time Has Come and Gone," *ICDE*, 2005.
- [82] X. Sun, C. J. Xue, J. Yu, T.-W. Kuo, and X. Liu, "Accelerating data filtering for database using FPGA," *Journal of Systems Architecture*, vol. 114, p. 101908, 2021.
- [83] N. C. Thompson and S. Spanuth, "The decline of computers as a general purpose technology," *CACM*, vol. 64, no. 3, pp. 64–72, 2021.
- [84] TPC, "TPC-H benchmark," <http://www.tpc.org/tpch/>, 2021.
- [85] J. Wang, D. Park, Y. Papakonstantinou, and S. Swanson, "SSD In-Storage Computing for Search Engines," *IEEE TC*, p. 1, 2016.
- [86] A. S. Waterman, *Design of the RISC-V instruction set architecture*. University of California, Berkeley, 2016.
- [87] L. Wu, A. Lottarini, T. K. Paine, M. A. Kim, and K. A. Ross, "Q100: the architecture and design of a database processing unit," *ASPLOS*, 2014.
- [88] —, "The Q100 Database Processing Unit," *IEEE Micro*, vol. 35, no. 3, pp. 34–46, 2015.
- [89] Xilinx, Inc., "Zynq UltraScale+ MPSoC - All Programmable Heterogeneous MPSoC," August 2016. [Online]. Available: <https://www.xilinx.com/products/silicon-devices/soc/zynq-ultrascale-mpsoc.html>
- [90] M. Xue, Q. Xing, C. Feng, F. Yu, and Z.-G. Ma, "FPGA-Accelerated Hash Join Operation for Relational Databases," *IEEE Transactions on Circuits and Systems II: Express Briefs*, vol. 67-II, no. 10, pp. 1919–1923, 2020.
- [91] C. Zhang, R. Chen, and V. K. Prasanna, "High Throughput Large Scale Sorting on a CPU-FPGA Heterogeneous Platform," *IPDPS Workshops*, 2016.
- [92] J. Ziv and A. Lempel, "A Universal Algorithm for Sequential Data Compression," *TIT*, vol. 23, no. 3, pp. 337–343, 1977.
- [93] M. Zukowski, S. Héman, N. Nes, and P. A. Boncz, "Super-Scalar RAM-CPU Cache Compression," *ICDE*, 2006.