

BOSS - An Architecture for Database Kernel Composition

Hubert Mohr-Daurat
Imperial College London
h.mohr-daurat19
@imperial.ac.uk

Xuan Sun^{*}
The University of Edinburgh
xuan.sun@ed.ac.uk

Holger Pirk
Imperial College London
hlgr@ic.ac.uk

ABSTRACT

Composable Database System Research has yielded components such as Apache Arrow for Storage, Meta’s Velox for processing and Apache Calcite for query planning. What is lacking, however, is a design for a general, efficient and easy-to-use architecture to connect them. We propose such an architecture. Our proposal is based on the ideas of partial query evaluation and a carefully designed, unified exchange format for query plans and data. We implement the architecture in a system called BOSS¹ that combines the Apache Arrow, the GPU-accelerated compute kernel ArrayFire and the CPU-oriented Velox kernel into a fully-featured relational Data Management System (DMS). We demonstrate that the architecture is general enough to incorporate practically any DMS component, easy-to-use and virtually overhead-free. Based on the architecture, BOSS achieves significant performance improvement over the CPU only Velox kernel and even outperforms the highly-optimized GPU-only DMS HeavyDB for some queries.

1. INTRODUCTION

The use cases of Data Management Systems (DMSs) have evolved from simply storing and retrieving data to managing all aspects of the data lifecycle as well as hardware resources. To this end, DMSs have to account for many kinds of heterogeneity. First, there is increasing workload heterogeneity: classic scenarios like Online Transactional and Analytical Processing are joined by new applications such as data cleaning, data integration, model training or inference. In addition, applications require various heterogeneous data models such as relations, graphs, documents, key-value-pairs and even trained models. Last but not least, systems must manage heterogeneous hardware devices such as CPUs, GPUs, Smart-Storage Devices, Trusted Enclaves, FPGAs, TPUs and other Application-Specific Integrated Circuits (ASICs).

To address the challenge of heterogeneity, systems currently have two options. The first option is to extend the DMS kernel to support heterogeneous devices, workloads and data models. While this approach usually yields good

performance, it requires substantial engineering effort. This effort leads to a fractured ecosystem with many systems supporting one or two “aspects of heterogeneity” (GPUs [10, 13, 28], data cleaning [12], JSON-fields [16], ...) but, to the best of our knowledge, none supporting more. Alternatively, systems can wrap special-purpose libraries in UDFs to support heterogeneity. While this is substantially less effort, it comes at the cost of inferior performance [39, 24].

We propose a third option that combines the best of both: a novel DMS design based on the idea of partial evaluation. Under this paradigm, a query is no longer processed by a single kernel but goes through a sequence of stages, each of which progresses toward the final result. Between stages, the query is passed in a simple, unified format containing all information required to produce the result (data and code). New functionality (such as a kernel or library supporting a specific hardware device) can be integrated as a stage in the evaluation process. As the entire query (including inputs) is passed between kernels, each kernel can opportunistically evaluate as many (or few) operators as deemed beneficial. As long as the last kernel in the pipeline supports the full evaluation of the query, this paradigm guarantees that the query is fully evaluated while giving each kernel complete freedom to implement special data model semantics, support new workloads, exploit hardware-specific features or apply sophisticated optimizations. This insight is reflected in a recent trend towards composable DMSs [37].

There are, however, several challenges when realizing such a design. First, the kernels (and underlying libraries) follow fundamentally different designs, leading to an “impedance” mismatch when combining them: some kernels are “pull-driven” while others are “push-driven”; some provide dynamically typed declarative APIs while others are imperative and statically typed; some are stateful while others maintain no internal state; virtually all of these kernels maintain metadata that is exploited during query evaluation. Efficiently combining these kernels is non-trivial and requires significant amounts of boilerplate/glue code that is not only expensive to write and maintain but also error-prone.

We propose addressing these challenges through a number of technical contributions:

- We introduce the idea of partial database query evaluation as a paradigm for database system composition
- We introduce a unified physical in-memory data representation that is generic enough to support the exchange of data and code between DMS kernels yet simple to use and free from runtime overhead. Where state-of-the-art systems need to combine different interchange formats for data,

^{*}Work done while at Imperial College London

¹project available under MIT license at [booss.1sds.uk](https://github.com/booss/booss)

Copyright is held by the owner/author(s). Publication rights licensed to the VLDB Endowment. This is a minor revision of the paper entitled “BOSS - An Architecture for Database Kernel Composition”, published in PVLDB, Vol. 17, No. 4 ISSN 2150-8097. DOI: <https://doi.org/10.14778/3636218.3636239>

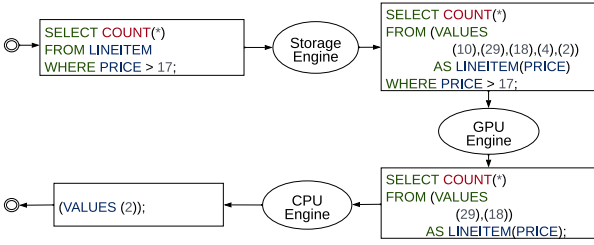


Figure 1: A Pipeline using Partial Query Evaluation

metadata and code, the proposed unified format solves all these problems using a powerful, elegant and highly efficient representation.

- We introduce (in our technical report [32]) a DMS-specific compile-time programming framework (type system, problem-specific language and memory management model) which extends well-known programming language techniques to solve the problem of integrating database kernels with minimal boilerplate code and virtually no performance overhead.
- We demonstrate the practical feasibility of these techniques by presenting the first data management system that composes multiple (unmodified) DMS kernels/libraries into a full-featured data analytics system: the Apache Arrow storage library [1], the CPU-oriented Velox database kernel [36], and the GPU-accelerated tensor compute kernel ArrayFire [29]. The system effectively demonstrates that DMS-components like Arrow, Velox and ArrayFire can be composed to complement their feature sets without performance penalty (compared to a monolithic architecture).

We first describe the general concept of *partial query evaluation* that will drive this system’s objectives, design principles and implementation.

2. PARTIAL QUERY EVALUATION

The goal of our work is to develop a DMS-architecture that allows the composition of separately designed components with minimal engineering effort. We propose for that a simple architecture: query processing is fragmented into stages, each implemented (or integrated from off-the-shelf systems) as “Engines” with more straightforward, smaller kernels rather than a monolithic one. For example, the implementation we will present in Section 5 integrates three engines: the storage engine is solely responsible for resolving the table and column names and loading the data; the GPU engine accelerates selected relational operators; the CPU engine integrates a CPU-based DMS supporting all the relational operators.

To develop this DMS while minimizing the complexity of integrating engines and orchestration logic, we propose the concept of *partial query evaluation*: the queries are passed through a linear pipeline of engines. Each engine accepts an input query and returns a result query. The result query is equivalent to the input query (i.e., a correct interpretation of it) but “closer” to the final result¹. After the last stage, no code shall be left for evaluation in the result query.

The example in Figure 1 illustrates these steps: the input query (on the top left) is passed to the storage en-

¹note that it is, in principle, acceptable for an engine to return its input query unmodified

gine; **FROM LINEITEM** is replaced by the content of that table (on the top right). Next, the GPU engine evaluates the **WHERE PRICE > 17** predicate and filters the data but leaves the **SELECT COUNT(*)** unmodified (resulting in the query on the bottom right). Finally, the CPU engine evaluates the aggregation and produces the output (on the bottom left).

Naturally, SQL is an ill-suited representation for this approach. Manipulating such a format incurs implementation and execution overhead as well as scalability issues when exchanging data between engines. Instead, we design a data & code representation as well as a programming framework that suits the requirements.

3. DESIGN PRINCIPLES

Partial Query Evaluation requires designing the DMS-architecture such that engines are implemented and integrated with negligible runtime overhead and minimal boilerplate code. This overarching objective can be achieved by following four guiding principles.

Remain Unopinionated. To ensure that different kernels with different designs can be integrated easily, a composable DMS must not impose design decisions on kernels. In software engineering, this principle is commonly referred to as “unopinionated” design. In the context of DMSs, the principle applies to aspects as diverse as logical data model, data representation, memory model, execution model or concurrency control. A framework must, for example, support the integration of Volcano-style query processors [20], Bulk-processors [8], X100-style [53] processors and just-in-time compiled execution engines [33].

Minimize Boilerplate. Without an appropriate development framework, integrating external database kernels requires substantial boilerplate code. Such code is a productivity hazard for developers and a major source of bugs and performance overhead. Integrating a GPU-coprocessing library, for example, requires code to convert data objects, manage data transfers, schedule kernel execution, control concurrent execution and transfer results back. The amount of boilerplate code can easily exceed that of actual data processing code. The codebase of the Ocelot extension to support GPU-accelerated data processing in MonetDB [23], for example, contains roughly 5,000 lines of OpenCL kernel code for kernel operations and more than 23,000 lines of boilerplate C-code to coordinate execution.

Zero Copies, Minimal Transformations. Creating copies of data is a costly operation for high-performance DMSs. To be performance-competitive with well-designed monolithic systems, a composed DMS must not unnecessarily copy data unless absolutely required.

Interestingly, many kernels have the same internal data representation, making zero-copy data transfer possible: MonetDB, Velox, ArrayFire and Arrow, e.g., share almost the same data representation (the only difference being minor optimizations in string representations). Upon close inspection, this is not surprising, as all of these are designed to maximize CPU efficiency, and modern CPUs are optimized for a specific data representation: that of the ANSI C language. Most kernels store datasets in C-arrays augmented with metadata (such as histograms or sortedness flags). However, all kernels we studied support the extraction of the C-arrays and the construction of the kernel-specific representation without copying data. If exploited effectively, this enables zero-copy data transfers of data between kernels.

Language over Library. This principle is, to a large extent, a consequence of the first and the second but deserves special mention due to its effect on the design. A composition framework that is lightweight and unopinionated must resort to the least common denominator of the kernels it aims to support. That least common denominator is the host programming language (usually C or C++, though Rust and Swift are viable alternatives). This largely precludes the use of libraries to coordinate execution, generate code or transform data. By avoiding libraries, the framework allows a state-of-the-art compiler to perform optimizations such as function inlining, loop unrolling or vectorization even across kernel boundaries. Note, however, that this only applies to the framework and does not prevent kernels from using libraries internally.

In the remainder of this paper, we describe a framework that follows these guidelines and, through careful design and sophisticated metaprogramming, achieves high-performance processing in a surprisingly small codebase (fewer than 5K lines of modern C++ code).

4. CROSS-KERNEL COMMUNICATION

The first challenge to address when connecting multiple kernels without incurring prohibitive copy-overhead is the definition of an appropriate exchange format for data (base tables, indices, histograms, etc.) as well as code (i.e., query plans, selection predicates, UDFs, etc.). While projects like Apache Arrow [1] define an efficient bulk-oriented data exchange format, this format is ill-suited to represent query plans. On the other hand, projects like Substrait [5] define a format for query plans but do not address the data exchange problem. Combining both to orchestrate the operation of multiple database kernels requires substantial boilerplate code and often comes at significant runtime overhead to convert the exchange format to the internal representation. Efficiently connecting multiple kernels requires a unified, lightweight, in-memory representation of data and code that can be interpreted with minimal boilerplate. Such a representation must address several non-trivial challenges. First, it must support the bulk-oriented data format that is shared by virtually all high-performance in-memory DMSs: (statically typed) C-arrays. Next, it must support the arguments, dynamically typed, that come in from frontends like SQL, Python or R. It must also support placeholders (we refer to them as “symbols”) representing table or column references in an execution plan. The representation must, further, be extensible to allow kernels to represent internal datatypes (such as tensors or GPU memory objects) without the need for physical data transformation. Finally, it must provide a concise API that avoids runtime overhead by being amenable to compile-time optimizations like inlining and common subexpression elimination.

While designing such a representation may seem daunting, focusing on data-processing systems provides several opportunities to constrain the problem and, thereby, make it tractable. First, the type system of such systems is “closed” at runtime, i.e., no new types can be defined once the system has been compiled. Next, the type system can be divided into a “core”, i.e., types that every kernel must support and internal “extensions”, i.e., types that cannot be passed between kernels. Lastly, data is usually “typed” in bulk, i.e., the type of the first value in a column (or partition) is the

```

1  template<typename... ExtensionTypes>
2  typedef variant<bool, int64, double, string,
3                ExtensionTypes...> CWDTypedAtom;

```

Figure 2: A (slightly simplified) implementation of CWD using C++ Templates.

same as all other values in that column. Type-interpretation can, therefore, be performed once per-column rather per-value (this concept is referred to as “evidence-typing” in programming language research [25]).

Motivated by this analysis, we propose a unified representation for data and code that addresses the challenges by exploiting these opportunities. It combines several techniques, inspired by established programming language techniques but adapted to the needs of DMSs. We discuss them bottom-up in the rest of this section.

4.1 Extensible Closed-World Dynamic Typing

A unified exchange format’s most fundamental challenge is the efficient and easy-to-use representation of dynamically-typed values (we use the term Atoms). Typically, this is efficiently implemented using “tagged union”, such as C++ variants, but requires the set of types to be statically defined across the kernels. However, a kernel integration might require the declaration of further dynamic types, e.g., one supporting TensorFlow might require a `tf::TensorBuffer` type. As illustrated in Figure 2, our approach extends standard C++ variant by instantiating it with a set of “core” types (`bool`, `int64`, etc.) while allowing the compile-time extension with further types. A core-only type system, e.g., can be declared by instantiating the template without extension types (`CWDAtom<>`) while the type system that also supports TensorFlow would be declared as `CWDAtom<tf::TensorBuffer>`. This approach supports the combination of multiple different type-systems in the same executable. While data can only be exchanged *between* kernels in the standard type system, kernels can use a custom type system internally with *transformation-overhead-free* communication between kernels. As the type system is fully defined at compile time (a.k.a., closed at runtime), the compiler can generate type-specific operators statically.

4.2 Memory-Managed Spans

High-performance data processing kernels store and process data in collections rather than individual atoms. While many kernels store and maintain metadata (histograms, min, max, etc.), the core data structure of virtually all of them is plain C-arrays. While this lack of diversity might seem surprising initially, it is a logical consequence of the hardware support for that format. A universal data exchange format should, therefore, be based on C-arrays.

However, plain C-arrays have no support for memory management, requiring kernels to implement their own mechanism, often based on reference counting or mark-and-sweep garbage collection. As every kernel has its own memory management scheme, unifying them without substantial overhead is non-trivial. However, focusing on (analytical) data processing systems suggests an approach: as allocated memory objects tend to be large (megabytes to gigabytes), small, constant allocation/deallocation overheads are negligible. We propose a simple memory management mechanism based on function pointers and opaque payloads.

```

1  template<typename ElementType>
2  struct Span {
3      ElementType* data;
4      size_t size;
5      void* payload;
6      void (*destructor)(); // function pointer
7  };

```

Figure 3: Spans capture bulk-allocated data

Figure 3 illustrates the implementation of the idea in the form of a *BOSS Span*: inspired by C++20 spans, BOSS Spans are type-generic, thin wrappers around C-arrays their memory deallocation managed in the form of a function pointer and an untyped payload. Upon destruction, the function is called and can, e.g., delete an underlying C-array, unmap a memory-mapped file or not do anything if the span does not own its memory (the default). Note that the payload and the data can point to the same address but do not have to. If, e.g., a span is created from an input that stems from an external library and has a header preceding the data in the same allocated piece of memory, `data` would point to the first data item while `payload` would point to the address of the beginning of the allocated object, i.e., the header.

In addition, we implement Spans as move-only/ non copyable types: only ownership of data is transferred between functions (and engines), and data is never implicitly copied.

4.3 Semi-statically Typed Expressions

While Spans allow representing contiguous in-memory data, this abstraction does not allow to represent complex data structures (e.g., columns, tables, trees, hashtables), query plans/programs or any dynamic behaviour (e.g., on-demand data loading from file). We propose to implement all these abstractions in a single, unified representation, on top of closed-world-typed atoms and spans, with “symbolic expressions (*s-expressions*)”. Classic *s-expressions*, as introduced with LISP [31], are nested (linked) lists of dynamically typed atoms (including symbols). In LISP, *s-expressions* are represented in the iconic notation using parenthesis: `(list 7 2 3)` represents a data structure while `(lambda (x) (+ x 9))` represents a function, i.e., code. In general, the high interpretation overhead makes *s-expressions* manifestly unsuited as a data model for a high-performance DMS. Instead, we propose several novel extensions: most importantly, we statically type some expression arguments at compile-time while allowing some in (closed-world) dynamically typed form. Further, we use the Spans (introduced earlier) to store consecutive, identically-typed arguments of expressions. While Figure 4 illustrates the implementation in C++, let us discuss the extensions in more detail.

Dynamically Typed Arguments

In their simplest form, we propose expressions that are an (almost) faithful reimplement of classic *s-expressions*: they differ from classic *s-expressions* only in that the first element (a.k.a. the head) must be a Symbol (see line 6 in Figure 4). While this has no practical impact on expressivity (one could map expressions with a non-symbolic head to one where the head is a symbol with an empty name), we found it a practical restriction that enables several optimizations we discuss later. The API is simple: an expression representing a column of Extensible Closed-World Dynamic

```

1  typedef CWDTypedAtom<Expression> DynamicArgument;
2  typedef variant<Span<int>,
3              Span<float>, /*...*/> SpanArgument;
4  template <typename... StaticArguments>
5  class Expression {
6      Symbol head;
7      vector<DynamicArgument> dynamicArguments;
8      tuple<StaticArguments...> staticArguments;
9      vector<SpanArgument> spans;
10 public: // API
11     // slow but generic argument access
12     DynamicArgument getArgument(int);
13     // fast and specialized access using specific type
14     template<int i> auto getArgument();
15     SpanArgument getSpanArgument(int);
16 };

```

Figure 4: Semi-statically typed expressions

Typing-typed (CWDT-typed) values would, e.g., be created as: `Expression("Column", 5, 9.2, "seventeen", false)`.

Statically Typed Arguments

While CWDT-typed expressions are very flexible, they come at substantial overhead: on the one hand, a system needs to represent the runtime type using a type tag (usually an integer); on the other hand, operating on dynamically typed values requires dynamic dispatching of values to operators (a.k.a., visitation), which translates into (virtual) function calls. However, most data processing kernels operate on static types. Translating between dynamic and static types requires significant amounts of boilerplate code.

To address these problems, we propose to support static typing of expression arguments (see line 4 in Figure 4). Declaring an expression as `Expression<int, float, string, bool>("Column", 5, 9.2, "seventeen", false)`, e.g., creates an expression with statically typed arguments. Statically typed expressions save memory, CPU cycles and boilerplate code. They are also the basis for the Semi-Static Dispatching technique we introduce in the next section. First, however, let us discuss the last and, arguably, most consequential category of arguments: Spans.

Span Arguments

As discussed earlier, Spans are statically typed collections of atoms. In addition to statically and dynamically typed arguments, we propose to implement expressions that can take arbitrarily many Spans of potentially different types as arguments (see line 9 in Figure 4). The elements of the Spans are “logically” exposed arguments of the expression through the same API as is used to access other arguments (line 12). Physically, however, the Spans remain plain C-arrays that can be accessed directly using a “physical” API (line 15).

To conclude the description of our proposed expression arguments, let us describe the APIs used to access them. There are two APIs with distinct use cases: the “slow and unified” API (line 12 in Figure 4) is designed for non-performance-critical access like rule-based query optimization or printing results to the console while the “fast and specific” API is designed for performance-critical (read only) access like the tight loop of a processing operator (line 14 & 15 show the API to access static arguments and spans).

The proposed abstractions provide convenient *read-only* access to data and code in line with state-of-the-art dataflow

```

1 // decompose() API
2 class Expression {
3     tuple<Symbol, Statics, Dynamics, Spans>
4         decompose() &&;
5 };
6
7 // Decomposition and recomposition example
8 Expression removeDynamicArgs(Expression&& input) {
9     auto [head, statics, dynamics, spans]
10         = move(input).decompose();
11     return Expression(head, statics, {}, spans);
12 }

```

Figure 5: Destructive Decomposition API and Example

systems. However, due to the simplicity of the memory model (i.e., the lack of reference counting), even trivial operators like projections in a column-store would have to perform expensive copies when performing only structural changes to the expressions (i.e., when the data itself is not changed). To address this challenge, we propose a novel design pattern for dataflow systems that requires neither copies nor reference counting: *Destructive Decomposition*.

4.4 Destructive Decomposition

Like Spans, we propose implementing Expressions as non-copyable (move-only in C++-nomenclature) types. In fact, they must be move-only types because one of their components (the Span arguments) is move-only. However, while they cannot be copied, they can be decomposed into components which can be reassembled to implement operations like projections. Figure 5, lines 3& 4, illustrates the API: a `decompose()` function that returns the expression components (head, statics, dynamic arguments and spans) as a tuple. The `&&`-suffix to the function indicates that the function can only be applied to an rvalue-reference, i.e., an object that is about to be destroyed. It is the very semantics that is required for destructive decomposition.

Lines 8 to 11 in Figure 5 illustrate how destructive decomposition can be used to structurally modify an expression (in this example, to remove the dynamic arguments): according to the established ownership model, the function owns the expression object; it can safely destroy the object and extract its arguments. To do so, the function must mark `input` for destruction (by calling `move`) and call `decompose()`. Just like for Spans, the use of `move` is required and enforced by the compiler. Using the input expression after it has been moved results in a warning or compilation error (depending on compiler flags). The elements of the returned tuple are assigned directly to variables using the `auto []`-syntax. They are then reassembled into a new Expression object. Using destructive decomposition this way avoids expensive copies and (opinionated) reference counting.

5. PARTIALLY-EVALUATING ENGINES

To demonstrate the feasibility of multi-kernel DMSs and assess the utility of the abstractions we introduced in Sections 4, we implemented them in a system named BOSS, short for **Bulk-Oriented Symbol-Store** [2]. BOSS supports CPU/GPU co-processing of relational queries by pipelining a storage & loading kernel, a GPU-accelerated relational kernel and a CPU-based relational kernel that supports a subset of what is required for relational data analytics. Let us start by introducing the processing pipeline at a high level

```

1 ;; input query:
2 (Select Lineitem (Where Mode != "Truck"))
3 ;; after Storage engine substitutes symbols/Strings:
4 (Select (Table (Key 1 2 3 4 5) (Mode 0 0 1 2 0))
5         (Where Mode != 0))
6 ;; after ArrayFire partially-evaluates the query:
7 (Gather (Table (Key 1 2 3 4 5) (Mode 0 0 1 2 0))
8         2 3)
9 ;; after Velox evaluates the query:
10 (Table (Key 3 4) (Mode 1 2))
11 ;; after Storage engine re-write the Strings back:
12 (Table (Key 3 4) (Mode "Mail" "Air"))

```

Figure 6: Query evaluation in a staged pipeline

before presenting the different kernels.

5.1 Engine Pipelining

BOSS combines three engines: the Apache Arrow persistence engine [1], the CPU-oriented relational processing engine Velox [36], and the ArrayFire parallel computing kernel [29] to add GPU co-processing support. The example in Figure 6 illustrates the role of each engine in the execution pipeline. The initial query (line 2) is, first, evaluated by the Storage engine (lines 4 & 5) to substitute symbols referring to tables and columns with expressions containing the data of those columns. Next, the ArrayFire engine opportunistically evaluates the parts of the plan that benefit from GPU acceleration (lines 7& 8). Finally, the Velox engine evaluates all unevaluated operators and returns an evaluated result (line 10). There is no need for explicit operator assignment to kernels to decide which operators are evaluated by the ArrayFire or the Velox engine: the only orchestration logic is the order of the engines in the pipeline.

Let us, now, discuss the technical details of each of these engines and how they achieve zero-overhead interaction.

5.2 Storage & Loading

The storage engine implemented in BOSS has two functions: loading relational data from files into memory and substituting the table identifier symbols with the stored relational expressions.

Loading Relational Data

We implement the storage engine on top of Apache Arrow [1], which provides in-memory columnar storage in the form of *Arrow Arrays*. Arrow Arrays are wrappers around C-arrays. They can, therefore, be converted to BOSS Spans with zero-copy by passing the C-array pointer to the Span's constructor. Because the Arrow Array owns the C-array, to ensure that the C-array is released only when the Span is deallocated, the Arrow Array's shared pointer is passed as the payload for the Span's destructor.

To partition column data into smaller arrays that fit in the CPU's last-level cache, both Apache Arrow and BOSS support partitioning but with a different layout. As shown on the left of Figure 7, Apache Arrow uses a *relation-level partitioning layout*, i.e., relations are sets of partitions which are sets of column arrays. We adopted in BOSS the *column-level partitioning layout*, i.e., columns are composed of multiple partitions, shown on the right of Figure 7. This layout, supported by the Expression API by storing multiple Spans in the expressions (see Section 4.3) is simpler for manipulating data per column compared with Apache Arrow's layout

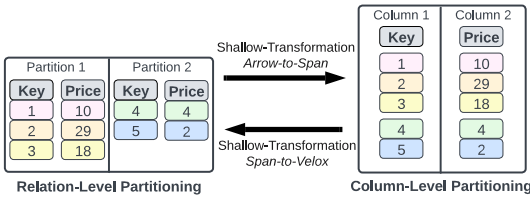


Figure 7: Conversion between relation-level (Arrow & Velox) and column-level (BOSS Span & ArrayFire) partitioning

and has less structural redundancy (e.g., column metadata is stored only once).

Data in one layout is converted into the other representation with minimal overhead using *shallow-transformation*, i.e., the target data structure is created, but the data is transferred without copies.

Symbol Substitution

As illustrated in Figure 6 lines 1 to 5, table identifier symbols are substituted by the storage engine with stored relational expressions, such as replacing the `Lineitem` symbol in line 2 with `(Table (Key 1 2 3 4 5) (Mode 0 0 1 2 0))` in line 4. To efficiently perform this substitution, table identifier symbols are stored in a dictionary of table identifiers to relational expressions. During query execution, the storage engine traverses query plan expressions and replaces table identifier symbols with stored relational expressions. As the storage engine owns persistent relational expressions (i.e., the `Table` expressions), they are shallow-copied into the plan: the storage engine instantiates new spans referencing the input C-arrays without passing a destructor function or payload. This effectively makes the spans non-owning. As the storage engine outlives the execution of the query, handing out non-owning references to spans is safe.

Compressed String Dictionaries

When a database contains string columns with low cardinality (i.e., few unique strings), storing duplicates of the strings in memory would be wasteful. In addition, executing equality predicates on a string is inefficient when only an integer comparison suffices. DMSs, therefore, typically implement a compressed string dictionary, i.e., strings are stored with an integer column whose values are indexed into a list of unique strings. The following illustrates how the partial query evaluation paradigm allows the implementation of this feature in the storage engine without changes in the other engines.

During loading, low-cardinality string columns are dictionary-compressed. During the query evaluation, the Storage engine evaluates string predicates on the dictionary and rewrites the query plan into one that operates directly on the (integer) keys (lines 4 & 5 in Figure 6). This allows “downstream” engines to efficiently process compressed strings without having to implement the functionality. To turn keys back to strings in the result, the storage engine wraps the query in a `StringLookup` operator (turning line 10 into 12).

5.3 Relational Processing on CPU

We implement a general-purpose CPU execution engine for the Multi-kernel DMS to handle relational processing based on Meta’s Velox kernel [36]. Integrating Velox into BOSS, requires transferring data from BOSS to Velox and

back and translating BOSS query plans to Velox. For that purpose, we implemented a thin wrapper around the Velox kernel we call the *Velox Engine*.

Data Transfer

Velox is a vectorized execution engine that processes data in batches of tuples in decomposed format. The fundamental materialization unit for decomposed data in Velox is a `FlatVector`, which represents a column or column-partition in a single memory block. Consequently, data As can be converted between Velox’ `FlatVector` and BOSS Spans without copy. Composed from multiple `FlatVectors`, `RowVectors` are used to represent a group of columns that are being passed between operators.

A Velox `FlatVector` can assume ownership of a memory object, evaluate operations on it and return ownership without requiring copying. BOSS’ Velox Engine adopts this Owned mode, by providing `spanToVelox` and `veloxToSpan` functions to transform Span to a `FlatVector` and vice-versa.

Illustrated in Figure 7, *Velox-to-Span* implements a *shallow transformation* between BOSS Span and Velox’ `RowVectors`. A vector of `RowVector` represents multiple Spans generated based on the *relation-level partitioning layout*. For example, the first batch of column `Key` and column `Price` is transformed to `RowVector 1`, while the second batch is transformed to `RowVector 2`.

Opportunities

When integrating Velox kernel into BOSS, we found that the Velox engine does not yet exploit all optimization opportunities. The `join` operator provided by Velox, e.g., does not exploit available indices, such as primary/foreign keys. Furthermore, the current version of Velox does not support GPU acceleration. A secondary engine can improve performance by exploiting these opportunities without modifying the Velox kernel itself. Let us in the following illustrate how BOSS offloads compute-intensive operators onto a GPU by using the ArrayFire Tensor compute kernel.

5.4 Acceleration on GPU

Due to limited capacity of GPU memory, transfer cost and massively parallel architecture, not all operators can be GPU accelerated [21]. Consequently, we implemented a GPU-accelerated engine that executes only some parts of the query. Taking advantage of partial query evaluation, the queries are traversed during execution to opportunistically evaluate some operators, leaving the rest of the query unevaluated (for a CPU-kernel to evaluate).

To demonstrate the use of off-the-shelf kernels, we build on the ArrayFire Tensor Processing Kernel [29] to implement the GPU-accelerated engine. ArrayFire provides a unified API to manipulate tensors, called *af::arrays*, with backends for CUDA, OpenCL and classic CPUs. Our engine implements relational operations on top.

Relational Operators Implementation

We limit our efforts to cases that can efficiently be implemented using only ArrayFire functions.

Projection. In its simple form, this operator is a no-op, i.e., only renaming columns using destructive decomposition. However, some projections evaluate arithmetic expressions. These are straightforward to implement using ArrayFire’s arithmetic operations.

Selection. This operator is implemented using ArrayFire by evaluating the predicate(s) to a bit array using ArrayFire’s boolean operations on `af::arrays` and transforming the bit array into a position list (i.e., an array of indices) using ArrayFire’s `af::where` operator. The kernel opportunistically evaluates as many predicates as possible given the GPU memory constraints.

As illustrated in Figure 6 (lines 7& 8), the ArrayFire Engine leaves some selections unevaluated to minimize cross-device data transfer. Those that are evaluated are transformed to a `Gather` operator, taking an (unmodified) relation and the position list as arguments to be evaluated by the Velox Engine.

Join. The typical implementation of a join on the GPU is a partitioned hash join [45, 28] which is non-trivial to implement with ArrayFire. However, ArrayFire supports the case of indexed primary/foreign key joins, which are, arguably, the most common joins. The join is implemented using ArrayFire’s `af::lookup` function to resolve the probe side of the joined relation. If the input data is partitioned into multiple spans as discussed in 5.2, they are combined into a single contiguous array in the GPU memory. This operation is performed at no extra cost when the data is transferred to the GPU.

6. EVALUATION

To assess the benefits of GPU-accelerating relational operators with our approach, we evaluate the performance of BOSS [2] with and without GPU acceleration and compare it with CPU-based DMSs and a GPU-based DMS. Then, we verify the benefits of avoiding transfer between devices (CPU to GPU) and between engines.

6.1 Experimental setup

Systems. We integrated Velox (git rev. fb33fbfec5895) [36] for the implementation of the CPU-based kernel and ArrayFire v3.8.3 [29], compiled with CUDA v12.0, for the GPU-accelerated kernel. To evaluate the performance of co-processing in BOSS in comparison with CPU-based DMS, we compare BOSS with MonetDB Jun2023_SP2_release [53] and DuckDB v0.8.1 [38]. To evaluate the performance with a GPU-based DMS, we compare BOSS with HeavyDB v6.4.0 [22]. While HeavyDB is a GPU-only system and we, therefore, do not expect to outperform it in terms of runtime, *there exists, to the best of our knowledge, no system that is designed for general-purpose co-processing.*

Hardware. All experiments are performed on a server with two Intel Xeon Silver 4114 2.20 GHz CPUs, each with 10 physical cores, a 14 MB LLC cache and 196 GB of memory. The GPU is a GeForce GTX Titan Xp with 12 GB of memory with NVIDIA driver v525.105.17. We use Ubuntu 18.04 with Linux kernel 4.15.0-209 and compile all code with Clang version 14 using the compiler flags `-O3 -mavx2`.

Workload. In all the experiments, we use the TPC-H benchmark [48]. We store numerical values with double-precision floating-point type for BOSS (due to better performance for GPU-acceleration). Following established practice [26], we evaluate the five queries that capture the benchmark’s choke points [9]. None of these queries performs integer arithmetic besides the aggregations, which ensures fair comparison since the ArrayFire engine does not check for arithmetic overflow yet (and does not evaluate aggregations). Due to the lack of support for non-dictionary com-

pressed string columns in the storage engine yet, Q9 is modified to filter the `PART` table on `P_RETAILPRICE` rather than `P_NAME` and Q18 to group by `C_CUSTKEY` rather than `C_NAME`. However, the cardinalities are preserved. Q1, Q3 and Q6 are the original TPC-H queries.

Query Plans. To ensure a fair comparison of BOSS’s hand-written query plans with other DMSs, we apply fixed heuristics that we plan to integrate into an automated query optimizer in the future. To decide the join order, we iteratively pick the two smallest among all the pairs of relations that can be joined next and always use the smaller of the two on the build side. To ensure that the query plan maximizes the operations that the GPU-accelerated kernel evaluates, we apply two more rules when using this kernel. First, we choose, as the next joined relations, the ones having indexed PK/FK keys over smaller relations. Second, we push the selections on joined tables down the query plan and order the predicates based on their independent selectivity (estimated by sampling the tables). For example, for Q9, `ORDERS` and `LINEITEM` are joined first because they have indexed PK/FK and the selection on `P_RETAILPRICE` is applied only after joining the `PART` table, so the indexed FK/PK on `PARTSUPP` and `PART` can be exploited.

6.2 Performance with TPC-H Benchmark

To start the evaluation, we study the typical case for GPU acceleration by executing the TPC-H queries with scale factor (SF) 10 (medium-sized dataset) and SF 100 (where data required for the GPU processing does not entirely fit into the GPU memory). The results in Figure 8 show, overall, a significant performance benefit from co-processing but with variations depending on the query characteristics.

At SF 10, Velox is outperformed by MonetDB and DuckDB for Q6, due to Velox’s handling of multiple predicates², and Q3, due to Velox not taking advantage of the indexed primary/foreign key present for the join relations of this query³. For Q9, Velox takes advantage of multi-threading for the high-cardinality join calculation and outperforms DuckDB but not MonetDB (due to inferior sorting and aggregation).

BOSS improves Velox’s runtime by a factor ranging from 1.3x (due to large data transfer back to CPU impeding the benefits for Q1) to 16x (Q6, benefiting from GPU-accelerating the selection). Q3, Q9 and Q18 also benefit from GPU-accelerated indexed primary/foreign key joins. Combining the efficiency of the CPU-based operators and the GPU-accelerated operators allows BOSS to outperform the two CPU-based baselines, MonetDB and DuckDB, for Q6 and Q18 and to reduce the gap for Q1, Q3 and Q9.

BOSS also outperforms the GPU-only baseline HeavyDB for Q3 (3x), Q6 (25x) and Q18 (4x) due to BOSS’s ability to GPU-accelerate only the part of the query that benefits from the GPU but is outperformed by HeavyDB for Q1 (1.4x) and Q9 (1.5x) due to HeavyDB transferring only a small fraction of (aggregated) data back to the CPU.

At SF 100, HeavyDB fails to process Q3, Q9 and Q18 (hence the missing bars) because the implementation does not provide a fallback method when the GPU memory is insufficient. Unlike HeavyDB, BOSS’ GPU-accelerated ker-

²Velox uses dictionary vectors to avoid the materialization of intermediate position lists for successive filter operators, but the nesting of index indirections affects the performance

³The Velox team considers taking advantage of an index structure outside their project’s scope

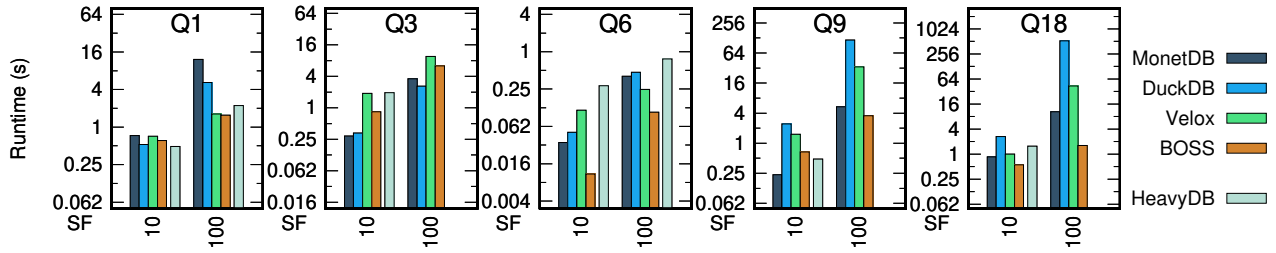


Figure 8: Runtime for TPC-H on five representative queries with scale factor 10 (i.e., 10GB)

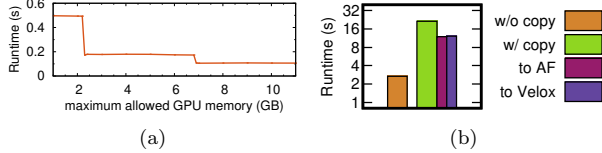


Figure 9: (a) Q6 runtime (SF 100) for various GPU size
(b) TPC-H query runtime (SF 10) with and without data copy between engines

nel can partially evaluate the query and leave the remaining evaluation to the CPU kernel, as explained in Section 5.4. The results show that, while fewer operations are GPU-accelerated, BOSS always benefits from GPU acceleration and outperforms the Velox implementation for all queries (although generally to a lesser extent than SF 10). For Q9 and Q18, BOSS outperforms Velox by factors of 8x and 32x, taking advantage of indexes and GPU-accelerated joins.

6.3 Avoiding Transfer

We designed BOSS to minimize data transfer between devices (CPU to GPU) with the partial evaluation strategy and between engines with the zero-copy data transfer.

To assess the effects of the partial evaluation strategy, we measure how much the data not fitting into the GPU memory affects the runtime execution. For this experiment, we evaluate TPC-H Q6 with SF 100 and vary the maximum size of data allowed to be transferred to the GPU from 0GB (i.e., no GPU acceleration) to 11GB (i.e., maximum GPU transfer allowing to GPU-accelerate two of the three selections). As expected, the results in Figure 9a show a significant performance improvement at 2.3GB when the memory is sufficient to evaluate the first (low-cardinality) selection and a smaller gain at 6.9GB (when the kernel can evaluate the second, higher-cardinality, selection). The partial evaluation allows to maximize GPU acceleration’s benefits for any available GPU memory resources.

To assess the benefits of zero-copy data transfer between engines, we measure the runtime overhead with a modified version of BOSS with data copy between the kernels, evaluated for TPC-H queries at SF 10. The results in Figure 9b show a high cost for data movement: between four and 250 times the execution runtime, depending on the query. This result confirms the significant performance benefit from avoiding data transfer between the kernels, which is made possible by the proposed data exchange format used to integrate DMSs into a single data processing pipeline.

We report more detailed results, insights and additional experiments in our technical report [32].

7. RELATED WORK

Composable DMSs is a relatively recent trend [36, 3, 4] even though DMS researchers have studied extensible designs for a long time [20, 14]. All of these focus on improving the reusability of the components. We, however, focus on a framework to compose kernels.

GPU Acceleration has been widely used in data analytics. Numerous contributions have been made to utilize GPU resources to speed up classic algorithms, such as sorting [46, 44, 42, 30] and nearest neighbor search [18, 19, 49]. DBMS acceleration using GPU has also received attention. Two main options exist: offloading operators [41, 45, 47, 40, 35, 52, 10] or building a standalone GPU DMS [22, 6, 23].

Various optimizations are proposed to improve the efficiency of GPU pipelines [50, 27, 17, 15, 34] and CPU/GPU co-processing [10, 11, 43, 7, 52, 51]. All these techniques are applicable but orthogonal to our work.

8. CONCLUSION

Composable DMS design promises a new generation of DMSs that are more extensible, easier to develop, cheaper, safer and more efficient while taking advantage of the latest techniques and technologies. However, current DMS architectures obstruct this goal by requiring substantial boilerplate code and expensive copies between components. To address this problem, we propose a fundamentally new architecture developed around the idea of partial query evaluation and a unified exchange format for data and executable query plans. We demonstrate the feasibility and utility of our approach by implementing BOSS [2], a system that integrates the Apache Arrow storage library, the GPU-accelerated ArrayFire compute kernel and the CPU-oriented Velox kernel. We demonstrated that BOSS is overhead-free and achieves significant performance improvement over the CPU-only Velox kernel and even outperforms the highly-optimized GPU-only DMS HeavyDB for some queries.

We argue that Partial Query Evaluation is the right framework for designing and implementing composable DMSs and will enable the integration of many data-oriented techniques and technologies. This includes the integration of new hardware, such as smart storage devices or application-specific circuits, as well as cloud services, such as serverless computation and smart cloud object stores. It will also accelerate the integration of new data management techniques, such as learned indices and new data models. Finally, this new paradigm holds the potential to simplify existing DMS architectures: kernels could, e.g., apply a mix of processing and optimization, effectively implement an adaptive query optimizer. We will explore such opportunities in future work.

9. REFERENCES

- [1] Apache Arrow, 2023. URL: <https://arrow.apache.org>.
- [2] BOSS, 2023. URL: <http://boss.lids.ksu.edu>.
- [3] Delta Lake, 2023. URL: <https://delta.io/>.
- [4] Google SQL, 2023. URL: <https://cloud.google.com/spanner/docs/reference/standard-sql/overview>.
- [5] Substrait, 2023. URL: <https://substrait.io>.
- [6] BlazingDB. Blazing SQL, 2023. URL: <https://github.com/BlazingDB/blazingsql>.
- [7] N. Boesch and C. Binnig. GaccO - A GPU-Accelerated OLTP DBMS. In *Proceedings of the 2022 International Conference on Management of Data*, SIGMOD '22, pages 1003–1016, New York, NY, USA, 2022. Association for Computing Machinery. doi:10.1145/3514221.3517876.
- [8] P. Boncz and M. L. Kersten. *Monet: A next-Generation DBMS Kernel for Query-Intensive Applications*. PhD thesis, Universiteit van Amsterdam, 2002.
- [9] P. Boncz, T. Neumann, and O. Erling. TPC-H analyzed: Hidden messages and lessons learned from an influential benchmark. In *Technology Conference on Performance Evaluation and Benchmarking*, pages 61–76. Springer, 2013.
- [10] S. Breß. The design and implementation of CoGaDB: A column-oriented GPU-accelerated DBMS. *Datenbank-Spektrum*, 14:199–209, 2014.
- [11] S. Breß, H. Funke, and J. Teubner. Robust query processing in co-processor-accelerated databases. In *Proceedings of the 2016 International Conference on Management of Data*, SIGMOD '16, pages 1891–1906, New York, NY, USA, 2016. Association for Computing Machinery. doi:10.1145/2882903.2882936.
- [12] J. Cambroneiro, J. K. Feser, M. J. Smith, and S. Madden. Query optimization for dynamic imputation. *Proceedings of the VLDB Endowment*, 10(11):1310–1321, 2017.
- [13] J. Cao, R. Sen, M. Interlandi, J. Arulraj, and H. Kim. Revisiting query performance in GPU database systems, 2023. arXiv:2302.00734.
- [14] M. J. Carey, D. J. DeWitt, G. Graefe, D. M. Haight, J. E. Richardson, D. T. Schuh, E. J. Shekita, and S. L. Vandenberg. The EXODUS extensible DBMS project: An overview. 1988.
- [15] P. Chrysogelos, M. Karpathiotakis, R. Appuswamy, and A. Ailamaki. HetExchange: Encapsulating heterogeneous CPU-GPU parallelism in JIT compiled engines. *Proc. VLDB Endow.*, 12(5):544–556, Jan. 2019. doi:10.14778/3303753.3303760.
- [16] D. Durner, V. Leis, and T. Neumann. JSON tiles: Fast analytics on semi-structured data. In *Proceedings of the 2021 International Conference on Management of Data*, pages 445–458, 2021.
- [17] H. Funke, S. Breß, S. Noll, V. Markl, and J. Teubner. Pipelined query processing in coprocessor environments. In *Proceedings of the 2018 International Conference on Management of Data*, SIGMOD '18, pages 1603–1618, New York, NY, USA, 2018. Association for Computing Machinery. doi:10.1145/3183713.3183734.
- [18] V. Garcia, E. Debreuve, and M. Barlaud. Fast k nearest neighbor search using GPU. In *2008 IEEE Computer Society Conference on Computer Vision and Pattern Recognition Workshops*, pages 1–6, 2008. doi:10.1109/CVPRW.2008.4563100.
- [19] V. Garcia, É. Debreuve, F. Nielsen, and M. Barlaud. K-nearest neighbor search: Fast GPU-based implementations and application to high-dimensional feature matching. In *2010 IEEE International Conference on Image Processing*, pages 3757–3760, 2010. doi:10.1109/ICIP.2010.5654017.
- [20] G. Graefe. Volcano—an extensible and parallel query evaluation system. *IEEE Transactions on Knowledge and Data Engineering*, 6(1):120–135, Feb. 1994. doi:10.1109/69.273032.
- [21] C. Gregg and K. Hazelwood. Where is the data? Why you cannot debate CPU vs. GPU performance without the answer. In *(IEEE ISPASS) IEEE International Symposium on Performance Analysis of Systems and Software*, pages 134–144. IEEE, 2011.
- [22] HEAVY.AI. HeavyDB, 2023. URL: <https://www.heavy.ai/product/heavydb>.
- [23] M. Heimel, M. Saecker, H. Pirk, S. Manegold, and V. Markl. Hardware-oblivious parallelism for in-Memory column-stores. *Proc. VLDB Endow.*, 6(9):709–720, July 2013. doi:10.14778/2536360.2536370.
- [24] D. Hirn and T. Grust. One WITH RECURSIVE is Worth Many GOTOs. In *Proceedings of the 2021 International Conference on Management of Data*, pages 723–735, Virtual Event China, June 2021. ACM. doi:10.1145/3448016.3457272.
- [25] W. Jones, T. Field, and T. Allwood. Deconstraining DSLs. In *Proceedings of the 17th ACM SIGPLAN International Conference on Functional Programming*, pages 299–310, Copenhagen Denmark, Sept. 2012. ACM. doi:10.1145/2364527.2364571.
- [26] T. Kersten, V. Leis, A. Kemper, T. Neumann, A. Pavlo, and P. Boncz. Everything you always wanted to know about compiled and vectorized queries but were afraid to ask. *Proceedings of the VLDB Endowment*, 11(13):2209–2222, Sept. 2018. doi:10.14778/3275366.3275370.
- [27] J. Li, H.-W. Tseng, C. Lin, Y. Papakonstantinou, and S. Swanson. HippogriffDB: Balancing I/O and GPU bandwidth in big data analytics. *Proc. VLDB Endow.*, 9(14):1647–1658, Oct. 2016. doi:10.14778/3007328.3007331.
- [28] C. Lutz, S. Breß, S. Zeuch, T. Rabl, and V. Markl. Pump Up the Volume: Processing Large Data on GPUs with Fast Interconnects. In *Proceedings of the 2020 ACM SIGMOD International Conference on Management of Data*, pages 1633–1649, Portland OR USA, June 2020. ACM. doi:10.1145/3318464.3389705.
- [29] J. Malcolm, P. Yalamanchili, C. McClanahan, V. Venugopalakrishnan, K. Patel, and J. Melonakos. ArrayFire: A GPU acceleration platform. In E. J. Kelmelis, editor, *SPIE Defense, Security, and Sensing*, page 84030A, Baltimore, Maryland, USA, May 2012. doi:10.1117/12.921122.
- [30] T. Maltnerberger, I. Ilic, I. Tolovski, and T. Rabl. Evaluating multi-GPU sorting with modern

- interconnects. In *Proceedings of the 2022 International Conference on Management of Data*, SIGMOD '22, pages 1795–1809, New York, NY, USA, 2022. Association for Computing Machinery. doi:10.1145/3514221.3517842.
- [31] J. McCarthy. Recursive functions of symbolic expressions and their computation by machine, Part I. *Communications of the ACM*, 3(4):184–195, Apr. 1960. doi:10.1145/367177.367199.
- [32] H. Mohr-Daurat, X. Sun, and H. Pirk. BOSS - An Architecture for Database Kernel Composition. *Proceedings of the VLDB Endowment*, 17(4):877–890, Dec. 2023. doi:10.14778/3636218.3636239.
- [33] T. Neumann. Efficiently compiling efficient query plans for modern hardware. *Proceedings of the VLDB Endowment*, 4(9):539–550, June 2011. doi:10.14778/2002938.2002940.
- [34] J. Paul, J. He, and B. He. GPL: A GPU-Based pipelined query processing engine. In *Proceedings of the 2016 International Conference on Management of Data*, SIGMOD '16, pages 1935–1950, New York, NY, USA, 2016. Association for Computing Machinery. doi:10.1145/2882903.2915224.
- [35] J. Paul, S. Lu, B. He, and C. T. Lau. MG-Join: A scalable join for massively parallel multi-GPU architectures. In *Proceedings of the 2021 International Conference on Management of Data*, SIGMOD '21, pages 1413–1425, New York, NY, USA, 2021. Association for Computing Machinery. doi:10.1145/3448016.3457254.
- [36] P. Pedreira, O. Erling, M. Basmanova, K. Wilfong, L. Sakka, K. Pai, W. He, and B. Chattopadhyay. Velox: Meta’s unified execution engine. *Proc. VLDB Endow.*, 15(12):3372–3384, Aug. 2022. doi:10.14778/3554821.3554829.
- [37] P. Pedreira, O. Erling, K. Karanasos, S. Schneider, W. McKinney, S. R. Valluri, M. Zait, and J. Nadeau. The composable data management system manifesto. *Proceedings of the VLDB Endowment*, 16(10):2679–2685, 2023.
- [38] M. Raasveldt and H. Mühleisen. Data Management for Data Science Towards Embedded Analytics. 2020.
- [39] K. Ramachandra, K. Park, K. V. Emani, A. Halverson, C. Galindo-Legaria, and C. Cunningham. Froid: Optimization of imperative programs in a relational database. *Proceedings of the VLDB Endowment*, 11(4):432–444, Dec. 2017. doi:10.1145/3186728.3164140.
- [40] R. Rui, H. Li, and Y.-C. Tu. Efficient join algorithms for large database tables in a multi-GPU environment. *Proc. VLDB Endow.*, 14(4):708–720, Dec. 2020. doi:10.14778/3436905.3436927.
- [41] R. Rui and Y.-C. Tu. Fast equi-join algorithms on GPUs: Design and implementation. In *Proceedings of the 29th International Conference on Scientific and Statistical Database Management*, SSDBM '17, New York, NY, USA, 2017. Association for Computing Machinery. doi:10.1145/3085504.3085521.
- [42] N. Samardzic, W. Qiao, V. Aggarwal, M.-C. F. Chang, and J. Cong. Bonsai: High-performance adaptive merge tree sorting. In *2020 ACM/IEEE 47th Annual International Symposium on Computer Architecture (ISCA)*, pages 282–294, 2020. doi:10.1109/ISCA45697.2020.00033.
- [43] A. Shanbhag, S. Madden, and X. Yu. A study of the fundamental performance characteristics of GPUs and CPUs for database analytics. In *Proceedings of the 2020 ACM SIGMOD International Conference on Management of Data*, SIGMOD '20, pages 1617–1632, New York, NY, USA, 2020. Association for Computing Machinery. doi:10.1145/3318464.3380595.
- [44] A. Shanbhag, H. Pirk, and S. Madden. Efficient top-K query processing on massively parallel hardware. SIGMOD '18, pages 1557–1570, New York, NY, USA, 2018. Association for Computing Machinery. doi:10.1145/3183713.3183735.
- [45] P. Sioulas, P. Chrysogelos, M. Karpathiotakis, R. Appuswamy, and A. Ailamaki. Hardware-Conscious Hash-Joins on GPUs. In *2019 IEEE 35th International Conference on Data Engineering (ICDE)*, pages 698–709, Macao, Macao, Apr. 2019. IEEE. doi:10.1109/ICDE.2019.00068.
- [46] E. Stehle and H.-A. Jacobsen. A memory bandwidth-efficient hybrid radix sort on GPUs. In *Proceedings of the 2017 ACM International Conference on Management of Data*, SIGMOD '17, pages 417–432, New York, NY, USA, 2017. Association for Computing Machinery. doi:10.1145/3035918.3064043.
- [47] Y. Tao, X. He, A. Machanavajjhala, and S. Roy. Computing local sensitivities of counting queries with joins. In *Proceedings of the 2020 ACM SIGMOD International Conference on Management of Data*, SIGMOD '20, pages 479–494, New York, NY, USA, 2020. Association for Computing Machinery. doi:10.1145/3318464.3389762.
- [48] TPC. TPC-H benchmark (revision 2.16.0), 2013. URL: <http://www.tpc.org/tpch/>.
- [49] P. Wieschollek, O. Wang, A. Sorkine-Hornung, and H. Lensch. Efficient large-scale approximate nearest neighbor search on the gpu. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2027–2035, 2016.
- [50] H. Wu, G. Damos, S. Cadambi, and S. Yalamanchili. Kernel weaver: Automatically fusing database primitives for efficient GPU computation. In *2012 45th Annual IEEE/ACM International Symposium on Microarchitecture*, pages 107–118, 2012. doi:10.1109/MICRO.2012.19.
- [51] B. W. Yogatama, W. Gong, and X. Yu. Orchestrating data placement and query execution in heterogeneous CPU-GPU DBMS. *Proc. VLDB Endow.*, 15(11):2491–2503, July 2022. doi:10.14778/3551793.3551809.
- [52] Y. Yuan, R. Lee, and X. Zhang. The yin and yang of processing data warehousing queries on GPU devices. *Proc. VLDB Endow.*, 6(10):817–828, Aug. 2013. doi:10.14778/2536206.2536210.
- [53] M. Zukowski, P. A. Boncz, N. Nes, and S. Héman. MonetDB/X100-A DBMS in the CPU cache. *IEEE Data Eng. Bull.*, 28(2):17–22, 2005.