



Apache Arrow DataFusion: A Fast, Embeddable, Modular Analytic Query Engine

Andrew Lamb
InfluxData
Boston, MA, USA
alamb@apache.org

Yijie Shen
Space and Time
Irvine, CA, USA
yjshen@apache.org

Daniël Heres
Coralogix
Utrecht, The Netherlands
dheres@apache.org

Jayjeet Chakraborty
UC Santa Cruz
Santa Cruz, CA, USA
jayjeetc@ucsc.edu

Mehmet Ozan Kabak
Synnada
Austin, TX, USA
ozankabak@apache.org

Liang-Chi Hsieh
Apple
Seattle, WA, USA
viirya@apache.org

Chao Sun
Apple
Cupertino, CA, USA
sunchao@apache.org

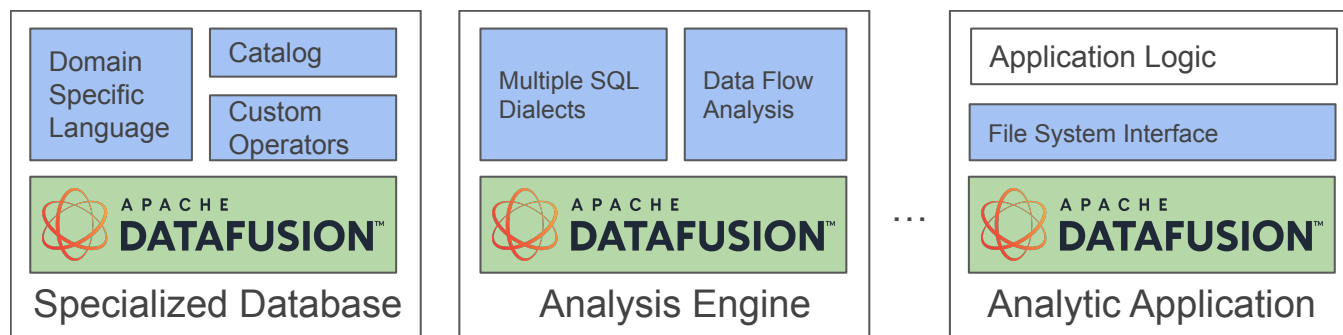


Figure 1: When building with DataFusion, system designers implement domain-specific features via extension APIs (blue), rather than re-implementing standard OLAP query engine technology (green).

ABSTRACT

Apache Arrow DataFusion[25] is a fast, embeddable, and extensible query engine written in Rust[76] that uses Apache Arrow[24] as its memory model. In this paper we describe the technologies on which it is built, and how it fits in long-term database implementation trends. We then enumerate its features, optimizations, architecture and extension APIs to illustrate the breadth of requirements of modern OLAP engines as well as the interfaces needed by systems built with them. Finally, we demonstrate open standards and extensible design do not preclude state-of-the-art performance using a series of experimental comparisons to DuckDB[66].

While the individual techniques used in DataFusion have been previously described many times, it differs from other industrial strength engines by providing competitive performance *and* an open architecture that can be customized using more than 10 major extension APIs. This flexibility has led to use in many commercial and open source databases, machine learning pipelines, and other

data-intensive systems. We anticipate that the accessibility and versatility of DataFusion, along with its competitive performance, will further the proliferation of high-performance custom data infrastructures tailored to specific needs assembled from modular components[18, 61].

CCS CONCEPTS

• **Information systems** → **Database management system engines; Online analytical processing engines; DBMS engine architectures;** Relational database model; *Database query processing*; • **Software and its engineering** → *Abstraction, modeling and modularity; Software performance; Software usability.*

KEYWORDS

Database Systems, Modular Query Engines, Column Stores, OLAP, Vectorized Execution, Parallel Execution, API Design

ACM Reference Format:

Andrew Lamb, Yijie Shen, Daniël Heres, Jayjeet Chakraborty, Mehmet Ozan Kabak, Liang-Chi Hsieh, and Chao Sun. 2024. Apache Arrow DataFusion: A Fast, Embeddable, Modular Analytic Query Engine. In *Companion of the 2024 International Conference on Management of Data (SIGMOD-Companion '24)*, June 9–15, 2024, Santiago, AA, Chile. ACM, New York, NY, USA, 13 pages. <https://doi.org/10.1145/3626246.3653368>

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

SIGMOD-Companion '24, June 9–15, 2024, Santiago, AA, Chile

© 2024 Copyright held by the owner/author(s). Publication rights licensed to ACM.

ACM ISBN 979-8-4007-0422-2/24/06.

<https://doi.org/10.1145/3626246.3653368>

1 INTRODUCTION

Traditionally, the realm of high-performance analytic query engines has been dominated by tightly integrated systems such as Vertica [47], Spark [77], and DuckDB [66]. This approach optimizes the interfaces between the file format, in-memory layout, and processing engine to achieve peak performance. However, building such a system is expensive and requires substantial commercial and/or research funding, given the extensive software engineering required.

As new requirements such as elastically scaling compute in public clouds and supplying AI/Machine Learning pipelines require new data systems, it has become clear that continually reimplementing such query engines is unnecessary. Due to the number analytic systems that have been built in industry and studied in academia, we know the best boundaries between subsystems such as file format, catalog, language front-ends and execution engine [18, 61].

DataFusion is designed with these API boundaries in mind, permitting assembly of end-to-end systems from open, reusable and high-quality components, using standards such as Apache Arrow (Section 2.1) and Apache Parquet (Section 2.2). Its competitive performance while being extensible demonstrates that a modern OLAP engine does not need to have a tight-knit architecture. Its mere existence demonstrates that permissive licensing and an open organizational structure [27] are capable of creating and maintaining this level of technology.

This paper makes the following contributions:

- (1) Describes the ecosystem of technologies that power DataFusion and that will likely power the majority of successful analytic systems in the next decade.
- (2) Describes several systems built with DataFusion, illustrating the possibilities of commodity OLAP engines.
- (3) Describes DataFusion’s architecture, feature set, and optimizations, illustrating the breadth of features required of modern analytic engines and quantifying the effort necessary to implement one.
- (4) Defines DataFusion’s extension APIs, outlining key module boundaries in an analytic stack.
- (5) Evaluates DataFusion’s performance, demonstrating that state-of-the-art performance is achievable using modular components and open standards.

The rest of this paper is organized as follows: **Section 2** reviews the technologies on which DataFusion is built. **Section 3** describes use cases and examples of real-world adoption. **Section 4** explores the trend towards modular databases. **Section 5** describes DataFusion’s architecture, detailing its execution model and key components. **Section 6** enumerates many of the standard query optimizations included in DataFusion. **Section 7** describes the APIs for extending DataFusion. **Section 8** evaluates DataFusion’s performance. We describe related work in **Section 9** and conclude in **Section 10**.

2 FOUNDATIONAL ECOSYSTEM

DataFusion is only possible due to the advent of several lower-level technologies: Apache Arrow’s in-memory columnar structure and compute kernels, Parquet’s efficient columnar storage, and the Rust ecosystem that enables a high-performance, yet comprehensible

implementation. Without these technologies, it is unlikely we could have built DataFusion with the relatively modest resources available. Additionally, using these technologies, systems built with DataFusion easily integrate with the broader ecosystem, directly sharing files and in-memory data streams without time-consuming and error-prone format transformations.

2.1 Apache Arrow

At its core, Apache Arrow [24] simply standardizes industrial best practices to represent data in memory using cache-efficient columnar layouts. By standardizing implementation details such as validity/null representations, endianness, variable length byte and character data, lists, and nested structures, systems built with Arrow benefit from well-known techniques and easy data interchange between applications. For example, while it is probably not critical for most systems if a NULL value is represented by a 0 or 1 in a bit mask, it is critical that all systems agree on the same convention for interoperability.

Originally, Arrow was designed as an in-memory interchange format and added compute-focused features such as `StringView` [21] and high-performance compute kernels over time. Arrow users can thus avoid re-implementing features that are well understood in academia and industry, but time-consuming to implement.

2.2 Apache Parquet

Apache Parquet [26] is an open source column-oriented data file format, originally designed for the Hadoop ecosystem and inspired by academic work on columnar storage [74]. It provides efficient data compression and encoding schemes, along with support for structured types via record shredding [56], embedded schema description, zone-map [57] like index structures and Bloom filters for fast data access.

Unlike Arrow, which is designed for fast random access and efficient in-memory processing, Parquet is optimized to store large amounts of data in a space-efficient manner. Like all formats, Parquet is not perfect, but it has become the de-facto standard for data storage and interchange in the analytic ecosystem. Its combination of open format, excellent compression across real-world data sets, broad ecosystem and library support, and embedded self-describing schema makes it a compelling choice for storing and exchanging compressed data. In addition to compression and compatibility, the file structure allows query engines to apply advanced projection and filter push-down techniques, such as late materialization, directly on files, yielding competitive performance compared with specialized formats [48].

2.3 Rust

Rust [76] is a relatively new system programming language, featuring a low-level, yet safe memory management approach and C-like performance. It incorporates an innovative memory ownership model that mitigates many of the worst memory and thread safety challenges prevalent in traditional C/C++ programming. Rust programs are easy to embed in other systems as they do not require a language run-time and have C ABI compatibility. Rust’s strong emphasis on zero-cost abstractions and its rich ecosystem of

performance-centric libraries, along with developer-friendly documentation and diagnostic tools, make it a compelling choice for implementing high-performance applications with a relatively lower engineering investment.

Unlike many C/C++ build systems, which can require substantial effort to just configure on a particular environment, Rust’s built-in Cargo Package Manager [12] and crate ecosystem makes adding DataFusion to most projects as simple as adding a single line to a configuration file.

3 USE CASES

A wide variety of commercial products and open source projects use DataFusion due to its combination of extensibility, feature set, fast query performance, and ecosystem compatibility. Projects built on DataFusion typically spend most time implementing value-adding features rather than replicating existing analytic engine technologies. While still early in adoption, DataFusion is already used in:

- (1) **Tailored database systems** for domain-specific use cases such as time series databases (e.g. InfluxDB 3.0 [41] and Coralogix [13]) and streaming SQL platforms (e.g. Synnada [75] and Arroyo [4]).
- (2) **Execution run-times** for specialized query front-ends, such as Comet for Apache Spark (Section 3.1), or Seafoam [69] for PostgreSQL [62] the Vega visualization language [55], and the InfluxQL [40] time series query language.
- (3) **SQL analysis tools** such as *dask-sql* [63] and SDF [68], which use DataFusion’s SQL parser, planner, and plan representation to analyze SQL queries.
- (4) **Table formats** such as the Rust implementations of Delta Lake [3], Apache Iceberg [30] and Lance [51], which use DataFusion expressions and query plans to fetch and decode remote data, implement predicate-based delete tombstones, push predicates to specialized secondary indexes, and compact files while retaining sort orders.

All these systems inherit the easy integration that comes with being Arrow-native. For example, Lance provides APIs for users to write Python functions that operate on RecordBatches using the existing *pyarrow* [28] library without any data conversion required.

3.1 Accelerating Apache Spark

A special use case is accelerating Apache Spark [77], an open-source analytic engine for large-scale data processing, widely adopted as a standard tool for data engineering, data science, and machine learning. Implemented primarily in JVM languages Scala and Java, its performance suffers from well-known JVM overheads.

Spark will likely remain a major component of data infrastructures in the near term given its high adoption and easy-to-use APIs. Fortunately, Spark’s design allows replacing *just* the execution engine with a specialized implementation such as Velox [60] (open-source) or Photon [7] (proprietary).

DataFusion is used by several native Spark runtimes, including Blaze [5] and Apache Arrow DataFusion Comet [11]. In these projects, Spark’s query front-ends, parsing, analysis, and optimization steps are used as is, while its execution plans are converted to DataFusion ExecutionPlans (Section 5.5) that execute through JNI interfaces with zero-copy data exchange via Apache Arrow.

In scenarios where Spark’s semantics differ from those offered by DataFusion, DataFusion’s extensible design (Section 7) permits projects to override and implement Spark specific expressions and operators (e.g., decimal related operations where Spark semantics deviate from ANSI SQL).

4 DECONSTRUCTED DATABASES

The rise of DataFusion and similar systems, such as Apache Calcite [6] and Velox [60], is part of a longer-term trend away from monolithic “one size fits all” general-purpose systems to “fit for purpose” specialized systems [73]. Given the expense of building the underlying technology, widespread proliferation of such specialized systems is only feasible when they can be assembled from reusable high-quality components, a trend which has been called the *Deconstructed Database* [44] [61].

The database systems literature offers a vast array of advanced and thoroughly studied techniques for most operations. However, due to economic and architectural constraints, these techniques have historically been confined to tightly integrated, often proprietary databases or analogous analytic systems. This tight integration limits reuse, leading to numerous costly reimplementations.

One classic example of a re-implementation is data science analysis tools, such as *pandas* [59]. The data science community innovated new APIs (DataFrame vs. SQL) and preferred a different deployment model (local files vs. networked servers), distinct from most contemporary database offerings. However, these tools initially performed poorly and did not incorporate many well-known techniques from database systems, such as query planning and optimization and parallel vectorized execution. In fact, Apache Arrow was initially born out of a desire to bring such well-studied database systems techniques to the data science ecosystem.

Another missed opportunity was the emergence of Map-Reduce [19] and its open source implementation, Hadoop, for parallel distributed processing. Although database researchers pointed out several ways that it was technically inferior [23], the lack of open reusable components inevitably led to the reimplementations of similar low-level analytical techniques.

4.1 Parallel with LLVM

The transition from very few monolithic implementations to many specialized systems sharing an open source foundation has occurred before. For example, system programming languages recently underwent a similar transformation, enabled by LLVM [52]. Compilers evolved from being tightly integrated with systems (IBM System / 390, Solaris, AIX HP-UX), to modular designs sharing the same LLVM backend (Rust [76], Swift [37], Zig [31], Julia [8]). Similarly, database systems evolved from tightly integrated systems directly managing hardware and software (Oracle [16], SQL Server [15], DB2 [14]) to fully modular designs sharing the same DataFusion backend (InfluxDB 3.0, GreptimeDB, and Coralogix).

Just as LLVM’s modular design catalyzed the development of system programming languages, DataFusion catalyzes the development of data systems. Authors of programming languages now focus on language-specific features and use LLVM for critical, yet commonplace, features such as IRs, auto-vectorization, and architecture-specific code generation. Authors of data systems can now focus on

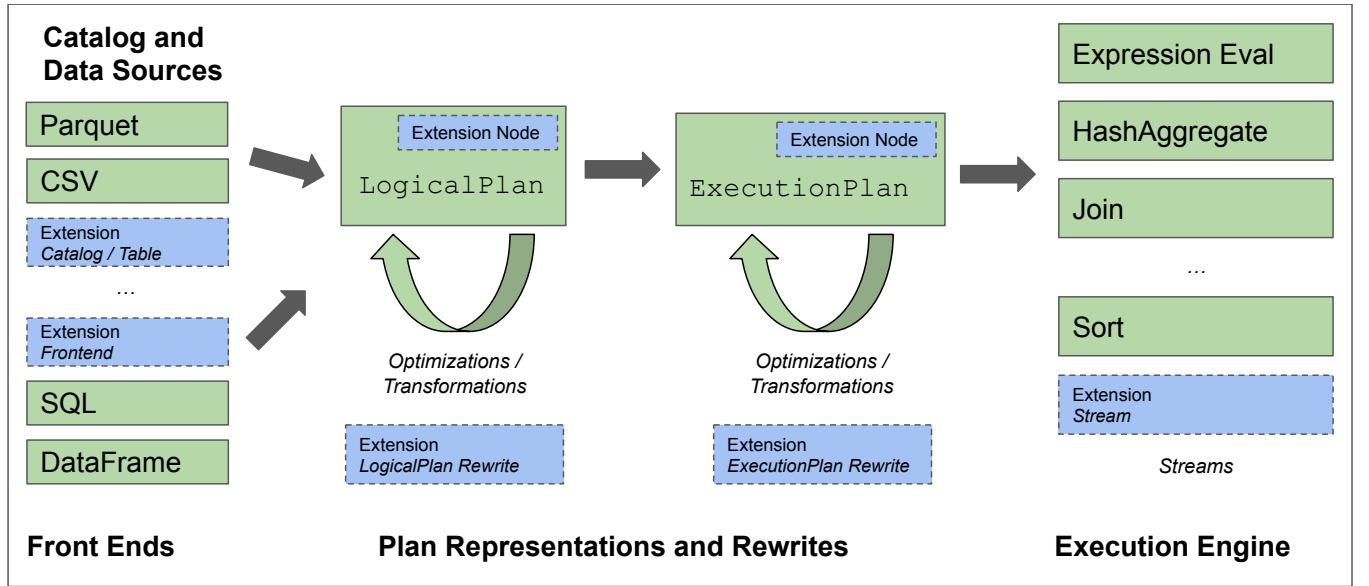


Figure 2: Architecture (Section 5). DataFusion’s standard query engine subsystems (green) run queries “out of the box”. Systems built on top of DataFusion customize behavior using extension APIs (blue).

value-added, domain-specific features and use DataFusion for features like SQL parsing, plan representations, optimizations, storage formats, and standard relational operators.

5 DATAFUSION FEATURES

5.1 Architecture Overview

DataFusion works “out of the box” while also providing extensive customization APIs, which we describe in Section 7. This architecture, shown in Figure 2, allows users to quickly start with a basic, high-performance engine and specialize the implementation over time to suit their needs and available engineering capacity. DataFusion’s implementation follows industrial best practices, informed by research literature, focusing on well-known patterns:

- (1) *Catalog and Data Sources* provide schema and data layout and location information.
- (2) A *Front End* creates a tree of relational operators and expressions, called a *LogicalPlan*.
- (3) *Optimizers*, using analysis APIs rewrite the *LogicalPlan* and expressions to a more optimal form.
- (4) The *LogicalPlan* is lowered to an *ExecutionPlan* that includes characteristics of the intermediate results such as sort order, and specific algorithm selections.
- (5) Additional optimizers rewrite the *ExecutionPlan* to match available physical resources and data layout.
- (6) The *ExecutionPlan* is executed by *Streams* that incrementally produce results.

5.2 Catalog and Data Sources

5.2.1 Catalog. In order to plan queries, like all query engines, DataFusion needs a *Catalog* to provide metadata such as which tables and columns exist, their data types, statistical information, and

storage details. DataFusion comes with a simple in-memory catalog and an Apache Hive [9]-style partitioned file/directory-based catalog. More complex catalog implementations are not included because catalog management is key system design decision and it is unlikely such implementations would be widely usable. Instead, most systems use the APIs described in Section 7.2 to supply catalog information (e.g., directly from a Hive metastore).

5.2.2 Data Sources. DataFusion includes *TableProviders* for commonly used file formats: Apache Parquet, Apache Avro, JSON, CSV, and Apache Arrow IPC files. The built in formats are implemented via the exact same API as user defined *TableProviders*, and support many features such as predicate, projection, and limit pushdown. The Parquet reader uses the native Arrow Rust implementation and implements predicate pushdown and late materialization (Section 6.8), bloom filters, and nested types. The CSV and JSON readers automatically infer schema, and the JSON reader fully supports nested types.

5.3 Front Ends

5.3.1 Data Types. DataFusion directly uses the Apache Arrow type system and inherits its broad range of supported types, including integral and floating-point numerics of various byte widths, fixed precision decimals, variable length character and binary strings, dates, times, timestamps, intervals, duration types, nested structs and lists. During execution, operators exchange data as either Arrow Arrays or scalar (single) values.

5.3.2 SQL Planner. DataFusion uses the *sqlparser-rs* [72] library to parse SQL text and generates a *LogicalPlan* from the parsed query representation. While it is likely that no SQL implementation should ever claim to be “complete” given the amorphous, ever-expanding SQL specification [42], DataFusion supports a large subset of SQL features including WHERE, GROUP BY, ORDER BY, LIMIT,

DISTINCT, WINDOW / OVER, UNION / INTERSECT, GROUPING SETS, and FULL / INNER / OUTER JOIN. It also supports more advanced functionality such ROWS / VALUES PRECEDING, FOLLOWING, UNBOUNDED window bounds, recursive CTEs, and GROUP BY with per-aggregate FILTER and ORDER BY.

5.3.3 DataFrame and LogicalPlanBuilder APIs. In addition to SQL, DataFusion also offers a DataFrame API, modeled after pandas [59], for building query plans in a procedural style. The DataFrame API generates the same underlying LogicalPlan representation (Section 5.4.1) as the SQL API, and is optimized and executed the same way. For more advanced uses, such as custom query languages, the LogicalPlanBuilder API offers a Rust builder-style interface for constructing plans directly.

5.4 LogicalPlan and Optimizer

5.4.1 Plans and Expressions. DataFusion’s API includes: (1) A full range of structures to represent and evaluate trees of expressions and relational operators, both at logical (Expr and LogicalPlan) and physical (PhysicalExpr and ExecutionPlan) levels, along with routines to create and manipulate them ergonomically; (2) Libraries to (de)serialize these structures from/to bytes suitable for network transport, both using Protocol Buffers as well as Substrait [22]; (3) Structures to describe statistics that may be known at planning time, such as row counts, null counts and minimum/maximum values.

5.4.2 Expression Analysis. In addition to basic expression evaluation, DataFusion provides libraries for simplification, interval analysis [58], and range propagation. Combined with statistics, these libraries provide predicate cardinality and selectivity estimates, and plan-time partition elimination (e.g. Parquet row group pruning, described in Section 6.8). These features are both usable directly by client systems and used to implement DataFusion’s built-in optimizations.

5.4.3 Function Library. DataFusion features a large library [29] of built-in scalar, window, and aggregate functions, including string operations, timestamp/date/time manipulations, interval arithmetic, and list/struct/map operations. These functions are implemented using the same API as user-defined functions by manipulating Arrow Arrays and can be invoked via both SQL or DataFrame APIs.

5.4.4 Rewrites. DataFusion includes an extensible plan rewriting framework, implemented as a series of LogicalPlan and ExecutionPlan rewrites. These passes handle details such as automatically coercing types to match available operator and function signatures, and introducing necessary sort and redistribution operations. The same framework is used for optimizations as well (Section 6.1).

5.5 Execution Engine

DataFusion uses a pull-based streaming execution model and distributes work across multiple cores using Volcano-style [33] exchange operators (viz. RepartitionExec).

5.5.1 Streaming Execution. Whenever possible, all operators produce output incrementally (Figure 3) as Arrow Arrays grouped

into RecordBatches, with a default size of 8192 rows. For pipeline-breaking operations such as a full sort, final aggregation, or a hash join, the operators buffer tuples, spilling to disk if necessary. Data flows between Streams as Arrow Arrays, which allows for seamless integration of user-defined operators (Section 7.7). Within each Stream, non-Arrow representations, such as the Row Format (Section 6.6) are used when necessary to increase performance.

```
impl Stream for MyOperator {
    ...
    // Pull next input (may yield at await)
    while let Some(batch) = stream.next().await {
        // Calculate, check if output is ready
        if Some(output) = self.process(&batch)? {
            // "Return" RecordBatch to output
            tx.send(batch).await
        }
    }
    ...
}
```

Figure 3: Streaming Execution. Each Stream (operator) implements the Rust Stream [32] trait, incrementally producing Apache Arrow RecordBatches that flow through the plan. Control flow is managed using Rust’s built-in await continuation generation, automatically marshaling the necessary state before yielding control. Each Stream attempts to output RecordBatches with a target number of tuples.

5.5.2 Multi-Core Execution. Each ExecutionPlan is run using one or more Streams (i.e. operators) that execute in parallel. Most Streams coordinate only with their input(s), though some must coordinate with sibling Streams, such as a HashJoinExec when building a shared hash table or a RepartitionExec when redistributing data to different streams. The number of Streams created for each ExecutionPlan is called its *partitioning*, which is determined at plan time (Figure 4).

5.5.3 Thread Scheduling. DataFusion Streams are implemented as Rust async functions and run on a thread pool implemented as a Tokio [67] runtime. Tokio is one of the most widely used libraries in the Rust ecosystem and was initially designed for asynchronous network I/O. However, its combination of an efficient, work-stealing scheduler, first-class compiler support for automatic continuation generation, and exceptional performance makes it a compelling choice for CPU-intensive applications as well [46]. While some recent work [53] describes challenges with the Volcano model on NUMA architectures, DataFusion achieves similar scalability as systems that use alternate designs (Section 8.2).

5.5.4 Memory Management. DataFusion manages memory using a MemoryPool, which is shared between one or more concurrently running queries. Streams cooperatively record when their memory consumption changes substantially by calling grow and shrink APIs. Stream implementations use a pragmatic approach, accurately tracking the largest memory consumers (e.g., contents of the

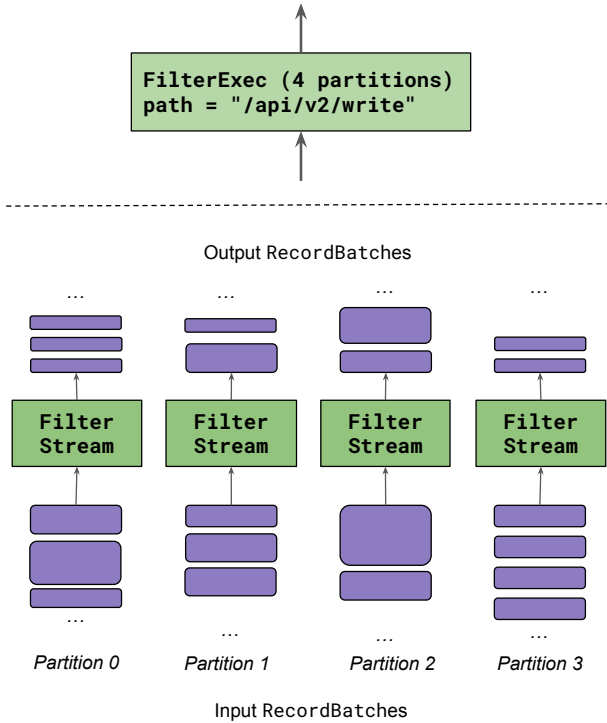


Figure 4: Partitioned Execution: Each `ExecutionPlan` is annotated with a number of partitions chosen by the planner, and a `Stream` (operator) is created for each partition. The `Streams` run independently on multiple threads. In this figure, the `FilterExec` `ExecutionPlan` (top) has 4 partitions. Thus, 4 distinct `FilterStream` operators are created during execution, and they run in parallel without coordination.

hash table for hash aggregation), but not small ephemeral allocations (e.g., memory for the current output batch).

DataFusion has two built-in memory pool implementations. `GreedyPool`, enforces per-process memory limits but does not attempt to distribute resources fairly across `Streams` in a query. `FairPool` distributes resources evenly among all pipeline-breaking `Streams`. DataFusion-based systems typically implement their own `MemoryPools` with domain-specific policies using the same API.

6 OPTIMIZATIONS

Query engines allow users to express desired results, and the engine handles the many details necessary to compute them efficiently. This section enumerates some of the techniques used by DataFusion to efficiently execute queries.

The techniques are not novel; each has been extensively studied and documented in research literature and implemented many times in commercial systems. DataFusion’s contribution is offering open-source, well-tested, extensible implementations which allow new systems to avoid the cost of re-implementing them (yet) again.

6.1 Query Rewrites

DataFusion includes a variety of query rewrites for both `LogicalPlans` and `ExecutionPlans`. `LogicalPlan` rewrites include projection pushdown, filter pushdown, limit pushdown, expression simplification, common subexpression elimination, join predicate extraction, correlated subquery flattening, and outer-to-inner join conversion. `ExecutionPlan` rewrites include eliminating unnecessary sorts, maximizing parallel execution, and determining specific algorithms such as Hash or Merge joins.

6.2 Sorting

Sorting, along with grouping and joining, is one of the most expensive operations in an analytic system and is well-studied in the literature. Most commercial analytic systems include heavily optimized multicolumn sorting implementations, and DataFusion is no exception. Broadly based on the techniques described in [34], it incorporates a tree of losers, normalized keys (`RowFormat` described in Section 6.6), the ability to spill to temporary disk files when memory is exhausted, and specialized implementations for `LIMIT` (aka “Top K”).

6.3 Grouping and Aggregation

Similarly to sorting, grouped aggregations are a core part of any analytic tool, as they create understandable summaries of large data volumes and thus are both well-studied and highly optimized in industrial systems. DataFusion contains a two-phase parallel partitioned hash grouping implementation [50], featuring vectorized execution, the ability to spill to disk when memory is exhausted, and special handling for no, partially ordered, and fully ordered group keys.

6.4 Joins

When joining multiple relations, DataFusion automatically identifies equality (equi-join) predicates, heuristically reorders joins based on statistics, pushes predicates through joins (subject to OUTER join restrictions), introduces transitive join predicates, and picks the optimal physical join algorithm. It includes parallel in-memory hash join, merge join, symmetric hash join, nested loops join and cross join implementations which each support Inner, Left, Right, Full, LeftSemi, RightSemi, LeftAnti, RightAnti joins, and are optimized for equality predicates. The in-memory hash join is implemented using vectorized hashing and collision checking similar to the algorithm described in [36]. While not implemented at the time of writing, we are working on additional join performance such as dynamically applying join filters during scans¹, a form of sideways information passing [70].

6.5 Window Functions

DataFusion supports SQL Window Functions (e.g. functions that have an `OVER` clause). Like most optimized window function implementations, DataFusion minimizes resorting by reusing existing sort orders, sorting only if necessary based on the `PARTITION BY` and `ORDER BY` clauses. It evaluates window functions incrementally [2], producing output once the required input window is present.

¹<https://github.com/apache/arrow-datafusion/issues/7955>

We have not yet found the need to implement newer, more sophisticated (and complex) schemes such as Physical Segment Trees [54] as the processing time of queries with window functions is typically dominated by other operations such as sorting.

6.6 Normalized Sort Keys / RowFormat

Columnar engines like DataFusion perform well on operations that naturally vectorize. However, query processing also requires efficient fundamentally row-based operations such as multi-column sorting and multi-column equality comparisons for grouping and joins, where the per row overhead can not be amortized by vectorization [43]. Within such operators, DataFusion uses a *RowFormat* [49], a form of normalized key [34] which 1) permits byte-wise comparisons using `memcmp` and 2) offers predictable memory access patterns. The *RowFormat* is densely packed, one column after another, with specialized encoding schemes for each data type, optionally adjusted for SQL sort options, such as ASC or DESC order and NULL placement. For example, unsigned and signed integers are encoded using their big-endian representation, whereas floating-point numbers are converted to a signed integer representation that incorporates the sign bit.

6.7 Leveraging Sort Order

DataFusion’s Optimizer is aware of, and takes advantage of, any pre-existing order in the input or intermediate results that flow from Stream to Stream. DataFusion 1) tracks multiple sort orders² and 2) includes Streams optimized for sorted or partially sorted input, such as Merge Join and partially ordered (streaming) Hash Aggregation.

Leveraging sort-order is important for at least two reasons:

- (1) **Physical Clustering:** Secondary indexes are often too expensive to maintain at the high ingest rates required for many OLAP systems, leaving the sort order of primary storage as the only available mechanism to cluster data.
- (2) **Memory Usage and Streaming Execution:** The sort order defines how the data that flows through Streams is partitioned during execution, defining where values may change and thus where intermediate results can be emitted.

6.8 Pushdown and Late Materialization

DataFusion pushes several operations down (towards the data sources): 1) projection (column selection) which elides unnecessary columns from intermediate results 2) LIMIT and OFFSET, which permits the plan to stop early when results are no longer needed and 3) predicates which moves filtering closer (or *in*) to data sources, minimizing the amount of data processed by the rest of the plan.

Pushing filters into data source enables `TableProviders` to apply filters *during* the scan, potentially avoiding significant work during execution. For example, DataFusion’s Parquet reader uses pushed-down predicates to 1) prune (skip) entire Row Groups and Data Pages based on metadata and Bloom filters and 2) apply predicates after decoding only a subset of column values, a form of late materialization [1] which can avoid the effort required to decode

values in other columns that will be filtered out. To illustrate, consider a query with the condition `A > 35 AND B = "F"`. DataFusion’s Parquet reader runs through these steps:

- (1) Prunes (skips) all Row Groups such that $A_{max} \leq 35$ or $B_{max} < "F"$ or $B_{min} > "F"$ using Row Group metadata.
- (2) Decodes column B, and evaluates `B = "F"`, capturing all rows which pass as a `RowSelection` (e.g. row indexes [100–244]).
- (3) Decodes only pages that contain the relevant rows from Column A, using the Page Index, and evaluates `A > 35` further refining the `RowSelection` (e.g. to row indexes [100–150]).
- (4) Decodes the pages containing the remaining `RowSelection` for any other selected columns (e.g. C).

Together, these techniques are very effective when predicate columns are cluster together such as when they appear early in the sort order of a sorted file [48].

7 EXTENSIBILITIES

This section describes the extension points for DataFusion, which are sufficiently flexible to support a wide variety of use cases (Section 3). This list of extension APIs offers a blueprint for other modular query engines and follows best practices for the internal boundaries of more tightly integrated systems.

Batches of data are represented as `ColumnarValues`, either a single scalar value or an Arrow Array. Because DataFusion itself uses Arrow throughout, extensions have the same performance as built-in functionality and they can use the same wide range of existing libraries, knowledge, and tools (e.g. computation kernels and Arrow Flight for network transfer).

7.1 Scalar, Aggregate, and Window Functions

Systems built on DataFusion often add use case specific functions that don’t belong in a general function library such as window functions that compute derivatives, calendar bucketing for time-series, and custom cryptography functions. Systems can register the following types of functions dynamically, which receive `ColumnarValues` as input and produce `ColumnarValues` as output:

- (1) **Scalar:** single output row for each input row
- (2) **Aggregate:** single output row for many input rows
- (3) **Window:** single output row for each input row; calculation has access to values in a surrounding window frame

DataFusion is not, of course, the first engine with user-defined function APIs. However, APIs in other systems are often limited in performance and functionality compared with built-in functions. Even when similarly performant APIs do exist, they are challenging to work with because they are tightly bound to engine specific data representations. Working with engine specific representations is especially challenging in column-oriented engines, due to the complexity of vectorization [43].

7.2 Catalog

Using a combination of the Catalog API and expression evaluation (Section 5.4.1), Catalogs built with DataFusion use file metadata (such as minimum and maximum values) to avoid reading entire files or parts of files (e.g. Parquet Row Groups). For example, the

²E.g. data is sorted by (A, B) and (A, C) via an order preserving join on B=C

Rust implementation of the Delta Lake table format [20] uses DataFusion to skip reading Parquet files based on filter predicates.

The Catalog API consists of 1) `TableProvider` for individual tables (Section 7.3), 2) `SchemaProvider`, a collection of `TableProviders`, and 2) `CatalogProvider`, a collection of `SchemaProviders`, sometimes referred to as a "catalog" or "database" in other systems. These APIs are async Rust functions, making it straightforward to integrate network I/O to access remote catalogs.

7.3 Data Sources

Using the DataFusion `TableProvider` trait, systems can query in-memory buffers of Arrow Arrays, stream data from remote servers (perhaps via Arrow Flight), or read from custom file formats. DataFusion's built-in providers (Section 5.2.2) are implemented with the same API exposed to users, and produce the same Rust async Stream of Arrow Arrays as `ExecutionPlans`. The `TableProvider` API supports 1) partitioned inputs, 2) pushdown of projection, filter, and limit, 3) parallel concurrent reads, and 4) pre-existing sort orders, 5) statistics and 6) updates.

Similarly to user-defined functions, it is challenging to offer user-defined data sources that perform as well as built-in formats in tightly integrated engines. Not only must the implementation produce data in the engine's native format, it must also interact with the engine's expression representation to implement predicate pushdown, and potentially perform asynchronous network I/O to implement incremental (streaming) output.

7.4 Execution Environment

Execution environments vary widely from system to system. For example, if fast local storage (e.g. NVMe) is available, caching meta-data in memory might make less sense than in environments where it is not available (e.g. some Kubernetes environments). Likewise, some systems run multiple queries concurrently, optimistically sharing resources, and others run a mix of queries with predefined resource budgets.

DataFusion can be customized for different environments using resource management APIs. `MemoryPool` is described in Section 5.5.4. `DiskManager` creates reference counted spill files if configured. `CacheManager` caches directory contents (e.g. expensive object store LIST operations) and per-file metadata such as statistics required for planning and pruning (Section 6.8). Caching such information is important when querying from disaggregated storage such as object stores. Similarly to other APIs, DataFusion comes with simple implementations, and users who require more tailored policies provide their own implementations (e.g. eviction policies or limiting temp space).

7.5 New Query / Language Frontends

Users extend the SQL supported by DataFusion by rewriting the AST prior to calling the DataFusion SQL planner (Section 5.3.2). For more substantial extensions or entirely different languages (e.g. PromQL or Vega), users implement their own frontends that create `LogicalPlans` using the structures described in Section 5.4.1.

7.6 Query Rewrites / Optimizer Passes

DataFusion users add domain-specific optimizations such as input re-ordering and macro expansions by implementing `OptimizerRules` and `PhysicalOptimizerRules`, which rewrite `LogicalPlan` and `ExecutionPlan` trees, respectively, with the same APIs as built-in rewrites (Section 6.1). Users can also specify the order in which rewrites are applied, both provided as well as their own.

7.7 Relational Operators

Domain-specific systems often require relational operations not found in SQL-only systems. For example, InfluxDB IOx [39] has specialized operators for timeseries gap filling, schema pivoting operations, and insert order resolution.

Users extend DataFusion by implementing the `ExecutionPlan` trait, the same as the nodes provided with DataFusion, such as join, filter, group by, and windowing. DataFusion does not distinguish between user-defined and built-in `ExecutionPlans` during optimization and execution. While other systems such as [71] offer similar functionality in the form of user defined table functions, those APIs typically restrict the syntax and placement of those operators, and are often unable to perform as well as built-in operators.

8 PERFORMANCE EVALUATION

To quantify the performance penalty of using open standards and a modular architecture, rather than a tightly integrated design, we compared DataFusion's performance to DuckDB [66], a system we think exemplifies a state-of-the-art, tightly integrated analytic engine at the time of writing. DataFusion performs similarly over a variety of real-world usecases. Although we acknowledge the challenges of benchmarking [65], different target use cases, and the rate of change in both engines, our experiments show that there is nothing fundamental about an open design that limits performance.

The most important performance metric for query engines is the end-to-end query execution time, which we measured for benchmarks that reflect commonly encountered data sizes and queries:

- (1) **ClickBench** [38] models large scale web analytic processing, with queries that filter and aggregate a large denormalized dataset. We used the unmodified 14 GB athena_partitioned dataset, which consists of 100 Parquet files, each approximately 140 MB in size.
- (2) **TPC-H** [17] models classic data warehouse analytics with 22 queries that join several tables in summary reports. We used the TPC-H data generator with Scale Factor=10 and converted the 8 resulting CSV files to parquet files, limiting row groups to 1M records, with total size of 2.5 GB.
- (3) **H2O-G** [35], models operations commonly found in data science workloads. We ran the groupby task queries on the G1_1e7_1e2_5_0.csv dataset, a single 488 MB comma separated value (CSV) file with 10M records.

We ran all benchmarks directly on the raw source data files. While transforming and loading into specialized per-database formats is common in previous generation systems, we believe that this approach is increasingly impractical as data flows become more fluid and dynamic. For our target systems, data is most commonly read and written using open formats using a diverse ecosystem of tools.

Query	DataFusion	DuckDB	Delta
1	1.22	0.18	6.74x slower
2	0.36	0.81	2.25x faster
3	1.11	1.78	1.6x faster
4	1.09	1.5	1.38x faster
5	20.74	8.34	2.49x slower
6	17.81	11.98	1.49x slower
7	0.3	2.08	6.91x faster
8	0.37	0.83	2.24x faster
9	27.91	10.83	2.58x slower
10	25.84	14.11	1.83x slower
11	4.29	3.22	1.33x slower
12	4.67	8.69	1.86x faster
13	11.38	10.27	1.11x slower
14	26.96	14.61	1.84x slower
15	12.7	11.15	1.14x slower
16	13.31	9.12	1.46x slower
17	29.6	21.97	1.35x slower
18	29.09	21.23	1.37x slower
19	92.31	39.1	2.36x slower
20	0.8	1.33	1.65x faster
25	6.01	8.44	1.4x faster
26	5.02	6.11	1.22x faster
27	6.59	8.4	1.28x faster
28	23.62	23.85	1.01x faster
29	107.41	62.99	1.71x slower
30	5.91	69.08	11.7x faster
31	12.59	12.95	1.03x faster
32	14.85	15.93	1.07x faster
33	92.17	57.2	1.61x slower
36	27.89	11.48	2.43x slower
37	0.67	0.52	1.31x slower
38	0.34	0.38	1.12x faster
39	0.34	0.42	1.24x faster
40	2.05	0.83	2.46x slower
41	0.2	0.25	1.28x faster
42	0.17	0.24	1.43x faster
43	0.19	0.27	1.44x faster

Table 1: ClickBench performance on a single core, in seconds, processing a 14GB dataset partitioned into 100 parquet files.

We measured performance of DataFusion 32.0.0 and DuckDB 0.9.1, the most recently released versions as of the time of this writing, using their respective Python bindings. We evaluate per core efficiency in Section 8.1 and multithreaded scalability in Section 8.2. Our scripts are available online³.

8.1 Single Core Efficiency

To measure the CPU efficiency of each engine, we ran the benchmark queries using an Intel Broadwell CPU, 32 GB of RAM with 8 virtual cores on an e2-standard-8 VM instance on Google Cloud Platform. We measured using Ubuntu 22.04.3 LTS and the Linux

³<https://github.com/JayjeetAtGithub/datafusion-duckdb-benchmark>

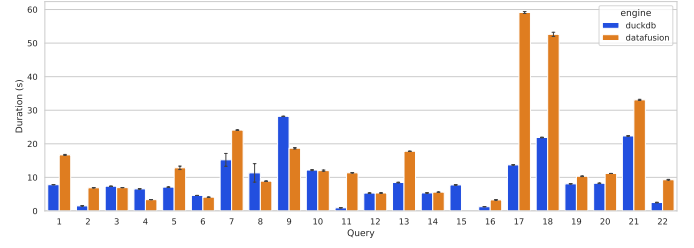


Figure 5: TPC-H SF=10 performance on a single core, one parquet file per table.

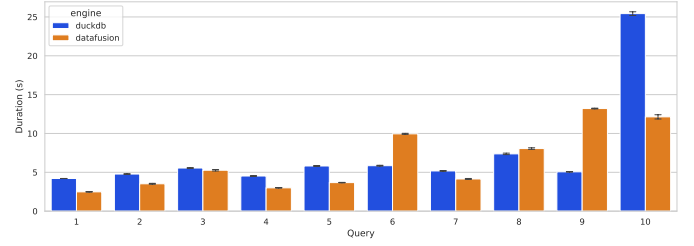


Figure 6: H2O-G (grouping) performance on single core with a single 488MB CSV file.

kernel version 6.2.0-1013-gcp. We limited each engine to a single core, in DataFusion by setting `target_partitions` to 1 and in DuckDB by setting the threads `PRAGMA` to 1.

ClickBench: Table 1 shows query execution time for the ClickBench queries. DataFusion performs better on queries with highly selective predicates such as Q2, Q8, and Q20 likely due to its ability to push predicates into the parquet scan to skip entire row groups. DataFusion also does well for queries with a single group such Q4 and Q7 and Q30, likely due to its vectorized aggregate updates. For queries with medium selectivity and medium group cardinality, such as Q15, Q31, Q32, Q41 and Q42 the engines have similar performance. For queries that have high group cardinality (10M groups or more) such as Q18, Q19, Q36, DuckDB performs better, likely due to its optimized parallel group by aggregation [45].

TPC-H: Figure 5 shows query execution time for TPC-H queries. Unlike ClickBench, most queries in TPC-H join multiple tables. DataFusion is faster for some queries such as Q4 and Q9, with highly selective predicates. There are some queries where performance is roughly equal such as Q3, Q6 and Q14. There are also several queries where DataFusion is well over 2x slower, such as Q11, Q17, Q18, and Q21. Much of this differences is due to a suboptimal join order⁴, and when we manually force a better join order, the performance of the two systems becomes similar.

H2O-G Figure 6 shows query execution time for the H2O-G queries. DataFusion has slightly better performance for most queries, though is significantly worse for Q9, due to an inefficient implementation of the `corr` aggregate function. The performance of all queries is largely dominated by the time spent parsing the CSV file, and DataFusion benefits from the optimized CSV parser included in the Rust implementation of Apache Arrow. Limiting to a single core may also unfairly penalize DuckDB, which seems to optimize

⁴<https://github.com/apache/arrow-datafusion/issues/7949>

multi-threaded parsing⁵, while a similar trade off doesn't exist for DataFusion.

Discussion Both engines perform similarly using a single core, with different strengths and weaknesses depending on attributes of the particular query. We conclude that there is nothing about using open standards that fundamentally limits DataFusion's performance. Our intuition and experience in implementing industrial systems is that the determining factor is instead available engineering investment. DataFusion's community already has projects underway to improve performance for query patterns where it lags DuckDB in these benchmarks, such as join ordering⁶ and high cardinality grouping⁷. Likewise, we expect that DuckDB's performance will improve in areas where it lags DataFusion such as low cardinality grouping, parquet predicate pushdown, and CSV parsing with additional investment.

8.2 Scalability

DataFusion is often used as a single-node engine, or the embedded engine in distributed systems, so its ability to scale "up" and use the resources of multiple cores in a single machine is important. Figure 7 shows how performance varies with increasing core count. We ran each ClickBench query 5 times, varying the number of cores from 1 to 192, plotting the final 3 runs to remove any caching or warm-up effects. We ran this experiment on the highest end CPU available to us on Google Cloud Platform, a c3-highcpu-176 instance with the Intel Sapphire Rapids micro-architecture, 176 virtual CPUs (cores), and 352 GB of memory. We ran all experiments using Ubuntu 22.04 with Linux kernel version 6.2.0-1016-gcp.

Relative performance The absolute value of the y axis is important. Some queries like Q10 take seconds to execute, while other queries like Q1 take less than a second. Thus, even while the relative performance difference between the two engines may appear substantial in some queries, such as Q1-Q4 or Q37-Q42, the absolute difference is 100s of milliseconds, while the absolute difference in queries such as Q19, Q32 and Q33 is an order of magnitude higher.

1, 2, 3, 8, 16, 32 cores Up to 32 cores, both DataFusion and DuckDB show excellent, near-linear decrease in execution times as the core counts increase.

64, 128, 192 cores At higher core counts, both engines show a mix of better and worse performance. In Q28 and Q29, performance continues to improve as the core count increases, close to the ideal curve. These queries contain low (6000) and medium (3M) cardinality grouping operations and require significant CPU effort to evaluate LIKE string matching predicates. In queries such as Q11, Q14, and Q32 both engines show a pronounced *increase* in query duration (they slow down) with more cores, likely because as the work done by each core decreases, the relative overhead of coordinating between the cores increases. In queries such as Q41, Q42 and Q43, the slowdown at high core count is more pronounced for DataFusion⁸, and in some queries such as Q25 and Q26 it is more pronounced for DuckDB.

⁵<https://github.com/duckdb/duckdb/issues/9136>

⁶<https://github.com/apache/arrow-datafusion/issues/7949>

⁷<https://github.com/apache/arrow-datafusion/issues/5546>

⁸Some of the slowdown in DataFusion is due to a poorly tuned hash table flushing strategy for high cardinalities <https://github.com/apache/arrow-datafusion/issues/6937>

Discussion DataFusion and DuckDB exhibit similar scaling behavior, and thus we conclude DataFusion's modular design and pull based scheduler do not preclude state of the art multi-core performance. The curves in Figure 7 for both engines are similar in shape, suggesting the differences are largely due to implementation details rather than fundamental differences in design.

9 RELATED WORK

The theme of more modular and composable architectures was observed at least as early as 2000 [10]. The term "Deconstructed Database" was initially popularized in 2018 [44], and there are recent calls to accelerate the adoption of modular designs [61].

Velox [60] and Apache Calcite [6] both provide components for assembling new databases and analytic systems. However, building a working end-to-end system requires substantial integration (e.g. bridging JVM and Native code and build systems), while using DataFusion requires a single configuration line. Modular designs allow swapping components based on use case, and the Photon [7] and Gluten [64] (based on Velox) projects replace just one module, the execution engine, of Apache Spark with a faster native implementation.

Similarly to DataFusion, DuckDB [66] is an open-source SQL system that does not require a separate server. DuckDB is targeted at users who run SQL, while DataFusion is targeted at people building new systems (that may run SQL as well as other types of processing). DuckDB has a more limited extension API and its own custom in-memory representation, storage format, Parquet implementation, and thread scheduler.

9.1 Future Research

Evaluating additional lines of code or engineering hours required to build a system from scratch without DataFusion, a study of how DataFusion's extensibility APIs are used in practice, and a systematic evaluation of end to end performance of DataFusion based applications would all help explore the trade offs encountered building systems using modular query engines. There is also a need for modular systems like DataFusion to accelerate other areas of database implementation, such as transaction processing and distributed key/value stores. First-class support, either as bindings to DataFusion or separate implementations, for other systems languages like C/C++ and Swift, are also needed.

10 CONCLUSION

Since the introduction of LLVM, the need to build compilers from scratch has been significantly reduced. With the emergence of technologies like DataFusion, the need to build database systems from the scratch should become similarly rare. Of course, with sufficient engineering investment, a tightly integrated engine can theoretically outperform a modular one. However, as the effort to reach state-of-the-art functionality and performance increases ever more, we believe that widely used, modular engines such as DataFusion will attract the required investment from large open-source communities and provide a richer feature set and better performance than all but the most well-resourced tightly integrated designs.

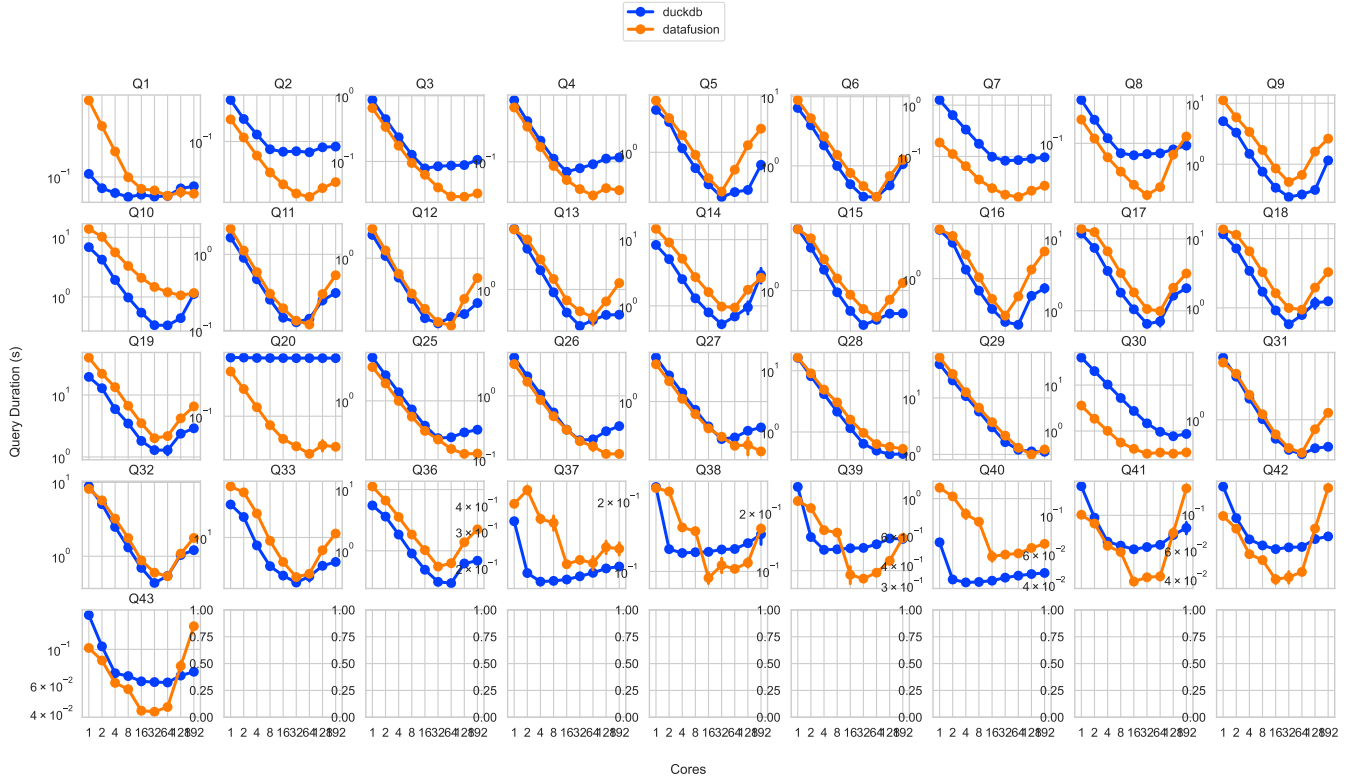


Figure 7: Query duration for ClickBench queries using 1, 2, 4, 8, 16, 32, 64, 128, or 192 cores, respectively.

Modular designs are by no means the only strategy for building systems, and we continue to see new tightly integrated systems emerge. However, as awareness of systems such as DataFusion increases, we predict adoption will accelerate and an explosion of new analytic systems will emerge that would previously not have been possible.

11 ACKNOWLEDGMENTS

DataFusion is a community-driven project encompassing a diverse array of individuals over a considerable span of time. At the time of writing, DataFusion had over 5500 proposed contributions (Pull Requests) from over 350 distinct individuals. The authors thank everyone who contributed ideas, reviews, bug reports, code, and tests over the years that made DataFusion possible. Finally, we thank the anonymous SIGMOD reviewers for their comments and suggestions on how to improve this paper.

REFERENCES

- [1] Daniel J. Abadi, Daniel S. Myers, David J. DeWitt, and Samuel Madden. 2007. Materialization Strategies in a Column-Oriented DBMS. In *Proceedings of the 23rd International Conference on Data Engineering, ICDE 2007, The Marmara Hotel, Istanbul, Turkey, April 15-20, 2007*, Rada Chirkova, Asuman Dogac, M. Tamer Özsu, and Timos K. Sellis (Eds.). IEEE Computer Society, 466–475. <https://doi.org/10.1109/ICDE.2007.367892>
- [2] Mustafa Akur and Mehmet Ozan Kabak. 2023. *Running Windowing Queries in Stream Processing*. <https://www.synnada.ai/blog/running-window-query-in-stream-processing>
- [3] Michael Armbrust, Tathagata Das, Sameer Paranjpye, Reynold Xin, Shixiong Zhu, Ali Ghodsi, Burak Yavuz, Mukul Murthy, Joseph Torres, Liwen Sun, Peter A. Boncz, Mostafa Mokhtar, Herman Van Hovell, Adrian Ionescu, Alicja Luszczak,

- Michal Switakowski, Takuya Ueshin, Xiao Li, Michal Szafranski, Pieter Senster, and Matei Zaharia. 2020. Delta Lake: High-Performance ACID Table Storage over Cloud Object Stores. *Proc. VLDB Endow.* 13, 12 (2020), 3411–3424. <https://doi.org/10.14778/3415478.3415560>
- [4] Arroyo. 2023. *Arroyo - Serverless Stream Processing*. <https://www.arroyo.dev/>
- [5] The Blaze Authors. 2023. *The Blaze accelerator for Apache Spark*. <https://github.com/blaze-init/blaze>
- [6] Edmon Begoli, Jesús Camacho-Rodríguez, Julian Hyde, Michael J. Mior, and Daniel Lemire. 2018. Apache Calcite: A Foundational Framework for Optimized Query Processing Over Heterogeneous Data Sources. In *Proceedings of the 2018 International Conference on Management of Data, SIGMOD Conference 2018, Houston, TX, USA, June 10-15, 2018*, Gautam Das, Christopher M. Jermaine, and Philip A. Bernstein (Eds.). ACM, 221–230. <https://doi.org/10.1145/3183713.3190662>
- [7] Alexander Behm, Shoumik Palkar, Utkarsh Agarwal, Timothy Armstrong, David Cashman, Ankur Dave, Todd Greenstein, Shant Hovsepian, Ryan Johnson, Arvind Sai Krishnan, Paul Leventis, Ala Luszczak, Prashanth Menon, Mostafa Mokhtar, Gene Pang, Sameer Paranjpye, Greg Rahn, Bart Samwel, Tom van Bussel, Herman Van Hovell, Maryann Xue, Reynold Xin, and Matei Zaharia. 2022. Photon: A Fast Query Engine for Lakehouse Systems. In *SIGMOD '22: International Conference on Management of Data, Philadelphia, PA, USA, June 12 - 17, 2022*, Zachary G. Ives, Angela Bonifati, and Amr El Abbadi (Eds.). ACM, 2326–2339. <https://doi.org/10.1145/3514221.3526054>
- [8] Tyler A. Cabutto, Sean P. Heeney, Shaun V. Ault, Guifen Mao, and Jin Wang. 2018. An Overview of the Julia Programming Language. In *Proceedings of the 2018 International Conference on Computing and Big Data (Charleston, SC, USA) (ICCBD '18)*, Association for Computing Machinery, New York, NY, USA, 87–91. <https://doi.org/10.1145/3277104.3277119>
- [9] Jesús Camacho-Rodríguez, Ashutosh Chauhan, Alan Gates, Eugene Koifman, Owen O'Malley, Vineet Garg, Zoltan Haindrich, Sergey Shelukhin, Prasanth Jayachandran, Siddharth Seth, Deepak Jaiswal, Slim Bouguerra, Nishant Bangarwa, Sankar Hariappan, Anishek Agarwal, Jason Dere, Daniel Dai, Thejas Nair, Nita Dembla, Gopal Vijayaraghavan, and Günther Hagleitner. 2019. Apache Hive: From MapReduce to Enterprise-grade Big Data Warehousing. In *Proceedings of the 2019 International Conference on Management of Data, SIGMOD Conference 2019, Amsterdam, The Netherlands, June 30 - July 5, 2019*, Peter A. Boncz, Stefan Manegold, Anastasia Ailamaki, Amol Deshpande, and Tim Kraska (Eds.). ACM, 1773–1786. <https://doi.org/10.1145/3299869.3314045>

- [10] Surajit Chaudhuri and Gerhard Weikum. 2000. Rethinking Database System Architecture: Towards a Self-Tuning RISC-Style Database System. In *VLDB 2000, Proceedings of 26th International Conference on Very Large Data Bases, September 10-14, 2000, Cairo, Egypt*, Amr El Abbadi, Michael L. Brodie, Sharma Chakravarthy, Umeshwar Dayal, Nabil Kamel, Gunter Schlegel, and Kyu-Young Whang (Eds.). Morgan Kaufmann, 1–10. <http://www.vldb.org/conf/2000/P001.pdf>
- [11] The Apache Arrow DataFusion Comet. 2024. *The Comet accelerator for Apache Spark*. <https://github.com/apache/arrow-datafusion-comet>
- [12] The Rust community. 2023. *Cargo: Rust's built-in package manager*. <https://crates.io/>
- [13] Coralogix. 2023. *Coralogix - Full-Stack Observability Platform with In-Stream Data Analytics*. <https://coralogix.com>
- [14] IBM Corporation. 2023. *IBM DB2*. <https://www.ibm.com/products/db2>
- [15] Microsoft Corporation. 2023. *Microsoft SQL Server*. <https://www.microsoft.com/en-us/sql-server>
- [16] Oracle Corporation. 2023. *The Oracle Database Server*. <https://www.oracle.com/database/>
- [17] The Transaction Processing Council. 2023. *The TPC-H Benchmark*. <https://www.tpc.org/tpch/>
- [18] Voltron Data. 2023. *The Composable Codex*. <https://voltrondata.com/codex>
- [19] Jeffrey Dean and Sanjay Ghenawat. 2004. MapReduce: Simplified Data Processing on Large Clusters. (2004), 137–150. <http://www.usenix.org/events/osdi04/tech/dean.html>
- [20] Delta-rs. 2024. *A native Rust library for Delta Lake*. <https://github.com/delta-io/delta-rs>
- [21] Arrow developers. 2023. *Mailing list: [DISCUSS][Format] Starting the draft implementation of the ArrayView array format*. <https://lists.apache.org/thread/r28rw5n39jwvtn080lj09d4q2c1ysvb>
- [22] Substrait Developers. 2023. *Substrait: Cross-Language Serialization for Relational Algebra*. <https://substrait.io/>
- [23] David J. DeWitt and Michael Stonebraker. 2008. *MapReduce: A major step backwards*. https://homes.cs.washington.edu/~billhowe/mapreduce_a_major_step_backwards.html
- [24] Apache Software Foundation. 2023. *Apache Arrow*. <https://arrow.apache.org>
- [25] Apache Software Foundation. 2023. *Apache Arrow DataFusion*. <https://arrow.apache.org/datafusion/>
- [26] Apache Software Foundation. 2023. *Apache Parquet*. <https://parquet.apache.org>
- [27] Apache Software Foundation. 2023. *A Primer on ASF Governance*. <https://www.apache.org/foundation/governance/>
- [28] Apache Software Foundation. 2023. *PyArrow - Apache Arrow Python bindings*. <https://arrow.apache.org/docs/python/index.html>
- [29] The Apache Software Foundation. 2023. *Apache DataFusion SQL reference*. <https://arrow.apache.org/datafusion/user-guide/sql/index.html>
- [30] The Apache Software Foundation. 2024. *Apache Iceberg: The open table format for analytic datasets*. <https://iceberg.apache.org/>
- [31] The Zig Software Foundation. 2023. *The Zig programming language*. <https://ziglang.org/>
- [32] Rust futures crate. 2023. *Stream trait*. <https://docs.rs/futures/0.3.28/futures/prelude/stream/trait.Stream.html>
- [33] Goetz Graefe. 1990. Encapsulation of Parallelism in the Volcano Query Processing System. In *Proceedings of the 1990 ACM SIGMOD International Conference on Management of Data, Atlantic City, NJ, USA, May 23-25, 1990*, Hector Garcia-Molina and H. V. Jagadish (Eds.). ACM Press, 102–111. <https://doi.org/10.1145/93597.98720>
- [34] Goetz Graefe. 2006. Implementing sorting in database systems. *ACM Comput. Surv.* 38, 3 (2006), 10. <https://doi.org/10.1145/1132960.1132964>
- [35] H2O.ai. 2023. *Database-like ops benchmark*. <https://h2oai.github.io/db-benchmark/>
- [36] Stratos Idreos, Fabian Groffen, Niels Nes, Stefan Manegold, K. Sjoerd Mullender, and Martin L. Kersten. 2012. MonetDB: Two Decades of Research in Column-oriented Database Architectures. *IEEE Data Eng. Bull.* 35, 1 (2012), 40–45. <http://sites.computer.org/debull/A12mar/monetdb.pdf>
- [37] Apple Inc. 2023. *The Swift programming language*. <https://developer.apple.com/swift/>
- [38] ClickHouse Inc. 2023. *ClickBench — a Benchmark For Analytical DBMS*. <https://benchmark.clickhouse.com/>
- [39] InfluxData Inc. 2023. *Announcing InfluxDB IOx - The Future Core of InfluxDB Built with Rust and Arrow*. <https://www.influxdata.com/blog/announcing-influxdb-iox/>
- [40] InfluxData Inc. 2023. *The Influx Query Language Specification*. <https://github.com/influxdata/influxql>
- [41] InfluxData Inc. 2023. *InfluxDB — open source time series, metrics, and analytics database*. <https://influxdata.com/>
- [42] ISO/IEC 9075:2023 2023. *Information technology - Database languages - SQL Standard*. International Organization for Standardization.
- [43] Timo Kersten, Viktor Leis, Alfons Kemper, Thomas Neumann, Andrew Pavlo, and Peter A. Boncz. 2018. Everything You Always Wanted to Know About Compiled and Vectorized Queries But Were Afraid to Ask. *Proc. VLDB Endow.* 11, 13 (2018), 2209–2222. <https://doi.org/10.14778/3275366.3275370>
- [44] Amandeep Khurana and Julien Le Dem. 2018. The Modern Data Architecture: The Deconstructed Database. *login Usenix Mag.* 43, 4 (2018). <https://www.usenix.org/publications/login/winter-2018-vol-43-no-4/khurana>
- [45] DuckDB Labs. 2023. *Parallel Grouped Aggregation in DuckDB*. <https://duckdb.org/2022/03/07/aggregate-hashtable.html>
- [46] Andrew Lamb. 2022. *Using Rustlang's Async Tokio Runtime for CPU-Bound Tasks*. <https://thenewstack.io/using-rustlangs-async-tokio-runtime-for-cpu-bound-tasks/>
- [47] Andrew Lamb, Matt Fuller, Ramakrishna Varadarajan, Nga Tran, Ben Vandiver, Lyric Doshi, and Chuck Bear. 2012. The Vertica Analytic Database: C-Store 7 Years Later. *Proc. VLDB Endow.* 5, 12 (2012), 1790–1801. <https://doi.org/10.14778/2367502.2367518>
- [48] Andrew Lamb and Raphael Taylor-Davies. 2022. *Querying Parquet with Millisecond Latency*. <https://arrow.apache.org/blog/2022/12/26/querying-parquet-with-millisecond-latency/>
- [49] Andrew Lamb and Raphael Taylor-Davies. 2023. *Fast and Memory Efficient Multi-Column Sorts in Apache Arrow Rust*. <https://arrow.apache.org/blog/2022/11/07/multi-column-sorts-in-arrow-rust-part-1/>
- [50] Andrew Lamb, Raphael Taylor-Davies, and Daniël Heres. 2023. *Aggregating Millions of Groups Fast in Apache Arrow DataFusion*. <https://www.influxdata.com/blog/aggregating-millions-groups-fast-apache-arrow-datafusion>
- [51] Lance. 2023. *Lance: modern columnar data format for ML*. <https://lancedb.github.io/lance/>
- [52] Chris Lattner and Vikram S. Adve. 2004. LLVM: A Compilation Framework for Lifelong Program Analysis & Transformation. In *2nd IEEE / ACM International Symposium on Code Generation and Optimization (CGO 2004), 20-24 March 2004, San Jose, CA, USA*, IEEE Computer Society, 75–88. <https://doi.org/10.1109/CGO.2004.1281665>
- [53] Viktor Leis, Peter A. Boncz, Alfons Kemper, and Thomas Neumann. 2014. Morsel-driven parallelism: a NUMA-aware query evaluation framework for the many-core age. In *International Conference on Management of Data, SIGMOD 2014, Snowbird, UT, USA, June 22-27, 2014*, Curtis E. Dyreson, Feifei Li, and M. Tamer Özsu (Eds.). ACM, 743–754. <https://doi.org/10.1145/2588555.2610507>
- [54] Viktor Leis, Kan Kundhikanjana, Alfons Kemper, and Thomas Neumann. 2015. Efficient Processing of Window Functions in Analytical SQL Queries. *Proc. VLDB Endow.* 8, 10 (2015), 1058–1069. <https://doi.org/10.14778/2794367.2794375>
- [55] Jon Mease. 2023. *VegaFusion: Server side scaling for the Vega visualization library*. <https://vegafusion.io/>
- [56] Sergey Melnik, Andrey Gubarev, Jing Jing Long, Geoffrey Romer, Shiva Shivakumar, Matt Tolton, and Theo Vassilakis. 2010. Dremel: Interactive Analysis of Web-Scale Datasets. *Proc. VLDB Endow.* 3, 1 (2010), 330–339. <https://doi.org/10.14778/1920841.1920886>
- [57] Guido Moerkotte. 1998. Small Materialized Aggregates: A Light Weight Index Structure for Data Warehousing. In *VLDB'98, Proceedings of 24th International Conference on Very Large Data Bases, August 24-27, 1998, New York City, New York, USA*, Ashish Gupta, Oded Shmueli, and Jennifer Widom (Eds.). Morgan Kaufmann, 476–487. <http://www.vldb.org/conf/1998/p476.pdf>
- [58] Ramon E. Moore. 1966. *Interval analysis*. Vol. 4. Prentice-Hall Englewood Cliffs.
- [59] The pandas development team. 2020. *pandas-dev/pandas: Pandas*. <https://doi.org/10.5281/zenodo.3509134>
- [60] Pedro Pedreira, Orri Erling, Maria Basmanova, Kevin Wilfong, Laith S. Sakka, Krishna Pai, Wei He, and Biswapesh Chattopadhyay. 2022. Velox: Meta's Unified Execution Engine. *Proc. VLDB Endow.* 15, 12 (2022), 3372–3384. <https://doi.org/10.14778/3554821.3554829>
- [61] Pedro Pedreira, Orri Erling, Konstantinos Karanasos, Scott Schneider, Wes McKinney, Sathyanarayana R. Valluri, Mohamed Zait, and Jacques Nadeau. 2023. The Composable Data Management System Manifesto. *Proc. VLDB Endow.* 16, 10 (2023), 2679–2685. <https://doi.org/10.14778/3603581.3603604>
- [62] PostgreSQL. 2024. *The PostgreSQL Relational Database*. <https://www.postgresql.org/>
- [63] The Dask project. 2023. *The dask-sql project*. <https://dask-sql.readthedocs.io/en/latest/>
- [64] The OAP project. 2023. *Gluten: Plugin to Double SparkSQL's Performance*. <https://h2oai.github.io/db-benchmark/>
- [65] Mark Raasveldt, Pedro Holanda, Tim Gubner, and Hannes Mühleisen. 2018. Fair Benchmarking Considered Difficult: Common Pitfalls In Database Performance Testing. In *Proceedings of the 7th International Workshop on Testing Database Systems, DBTest@SIGMOD 2018, Houston, TX, USA, June 15, 2018*, Alexander Böhm and Tilmann Rabl (Eds.). ACM, 2:1–2:6. <https://doi.org/10.1145/3209950.3209955>
- [66] Mark Raasveldt and Hannes Mühleisen. 2019. DuckDB: an Embeddable Analytical Database. In *Proceedings of the 2019 International Conference on Management of Data, SIGMOD Conference 2019, Amsterdam, The Netherlands, June 30 - July 5, 2019*, Peter A. Boncz, Stefan Manegold, Anastasia Ailamaki, Amol Deshpande, and Tim Kraska (Eds.). ACM, 1981–1984. <https://doi.org/10.1145/3299869.3320212>
- [67] Tokio rs Developers. 2023. *Tokio: As asynchronous Rust runtime*. <https://tokio.rs/>
- [68] SDF. 2023. *SDF*. <https://www.sdf.com/engine>
- [69] Seafowl. 2024. *Seafowl Postgres Accelerator*. <https://seafowl.io/>

- [70] Lakshmikanth Shrinivas, Sreenath Bodagala, Ramakrishna Varadarajan, Ariel Cary, Vivek Bharathan, and Chuck Bear. 2013. Materialization strategies in the Vertica analytic database: Lessons learned. In *29th IEEE International Conference on Data Engineering, ICDE 2013, Brisbane, Australia, April 8-12, 2013*, Christian S. Jensen, Christopher M. Jermaine, and Xiaofang Zhou (Eds.). IEEE Computer Society, 1196–1207. <https://doi.org/10.1109/ICDE.2013.6544909>
- [71] Moritz Sichert and Thomas Neumann. 2022. User-Defined Operators: Efficiently Integrating Custom Algorithms into Modern Databases. *Proc. VLDB Endow.* 15, 5 (2022), 1119–1131. <https://doi.org/10.14778/3510397.3510408>
- [72] The sqlparser-rs authors. 2023. *sqlparser-rs: Extensible SQL Lexer and Parser for Rust*. <https://github.com/sqlparser-rs/sqlparser-rs>
- [73] Michael Stonebraker. 2008. Technical perspective - One size fits all: an idea whose time has come and gone. *Commun. ACM* 51, 12 (2008), 76. <https://doi.org/10.1145/1409360.1409379>
- [74] Michael Stonebraker, Daniel J. Abadi, Adam Batkin, Xuedong Chen, Mitch Cherniack, Miguel Ferreira, Edmond Lau, Amerson Lin, Samuel Madden, Elizabeth J. O’Neil, Patrick E. O’Neil, Alex Rasin, Nga Tran, and Stanley B. Zdonik. 2005. C-Store: A Column-oriented DBMS. In *Proceedings of the 31st International Conference on Very Large Data Bases, Trondheim, Norway, August 30 - September 2, 2005*, Klemens Böhm, Christian S. Jensen, Laura M. Haas, Martin L. Kersten, Per-Ake Larson, and Beng Chin Ooi (Eds.). ACM, 553–564. <http://www.vldb.org/archives/website/2005/program/paper/thu/p553-stonebraker.pdf>
- [75] Synnada. 2023. *Synnada realtime data platform*. <https://www.synnada.ai/>
- [76] The Rust team. 2023. *The Rust programming language*. <https://www.rust-lang.org/>
- [77] Matei Zaharia, Reynold S. Xin, Patrick Wendell, Tathagata Das, Michael Armbrust, Ankur Dave, Xiangrui Meng, Josh Rosen, Shivaram Venkataraman, Michael J. Franklin, Ali Ghodsi, Joseph Gonzalez, Scott Shenker, and Ion Stoica. 2016. Apache Spark: a unified engine for big data processing. *Commun. ACM* 59, 11 (2016), 56–65. <https://doi.org/10.1145/2934664>