

Apache DataFusion: Design Choices when Building Modern Analytic Systems

Boston University
Data System Seminar: October 28, 2024

Andrew Lamb
Staff Engineer @ InfluxData
Apache {DataFusion, Arrow} PMC



Andrew Lamb

Staff Engineer
InfluxData

> ~~20~~ 21 🤖 years in enterprise software development

Oracle: Database (2 years)

DataPower: XSLT compiler (2 years)

Vertica: DB / Query Optimizer (6 years)

Nutonian/DataRobot: ML Startups (7 years)

InfluxData: InfluxDB 3.0, Arrow, DataFusion (4 years)

Goal

Content: Case Study of Choices in Apache DataFusion Query Engine

- Exposure to important aspects of analytic query engines
- Introduction to currently important technologies

Convince you to build with / contribute to the building blocks:

- Rust
- Arrow
- Parquet
- DataFusion!
- ...

Outline

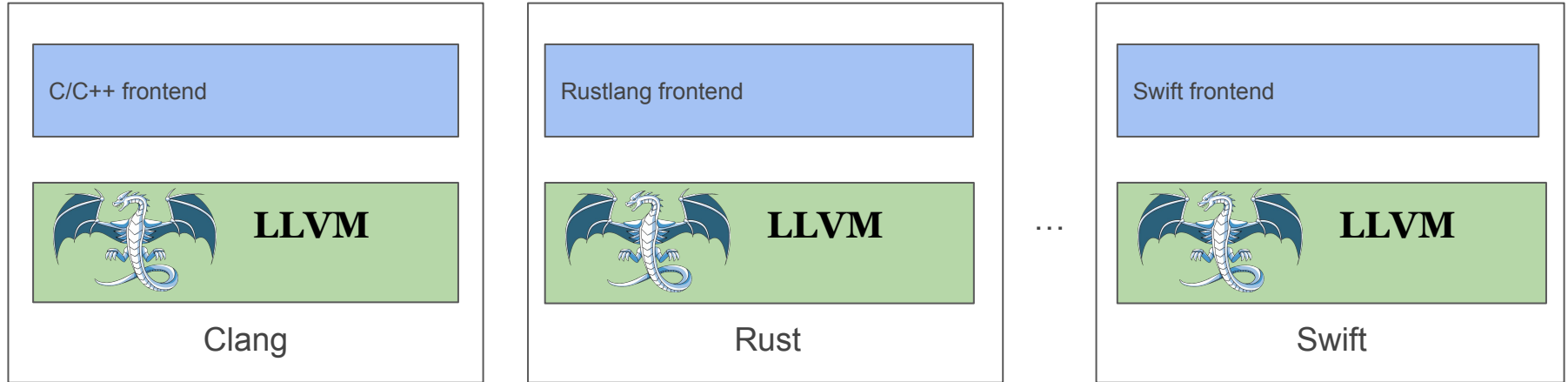
Brief Intro to DataFusion

Design Decisions

- Options
- What we chose and why
- Technical Overview
- Lessons learned

Intro to DataFusion

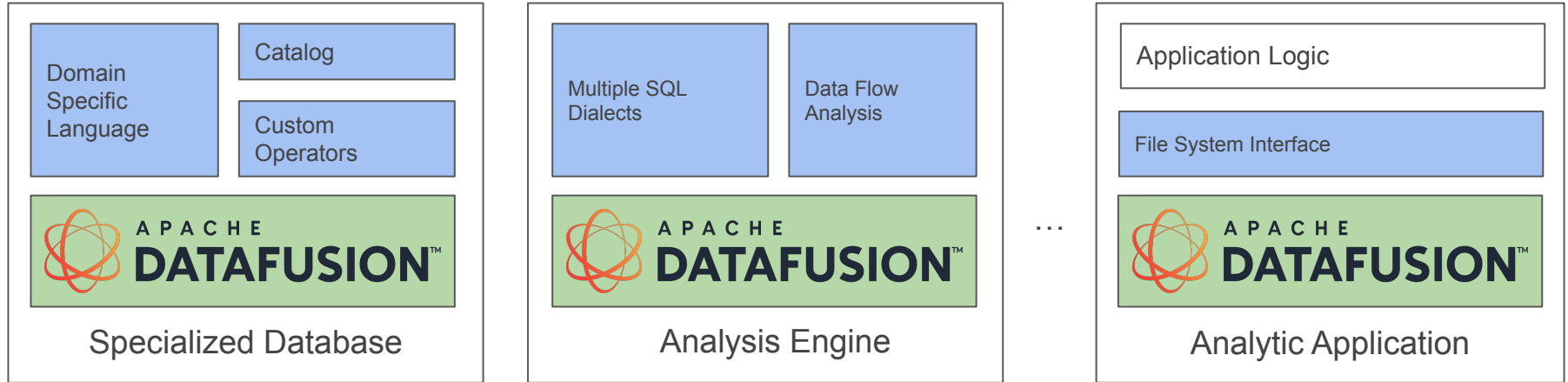
Analogy: DataFusion is LLVM for Databases



LLVM enabled innovation in programming languages:

- High quality reusable optimizer, code generator, debugger, lsp integration, etc.
- Focus on language design, ecosystem, libraries, etc

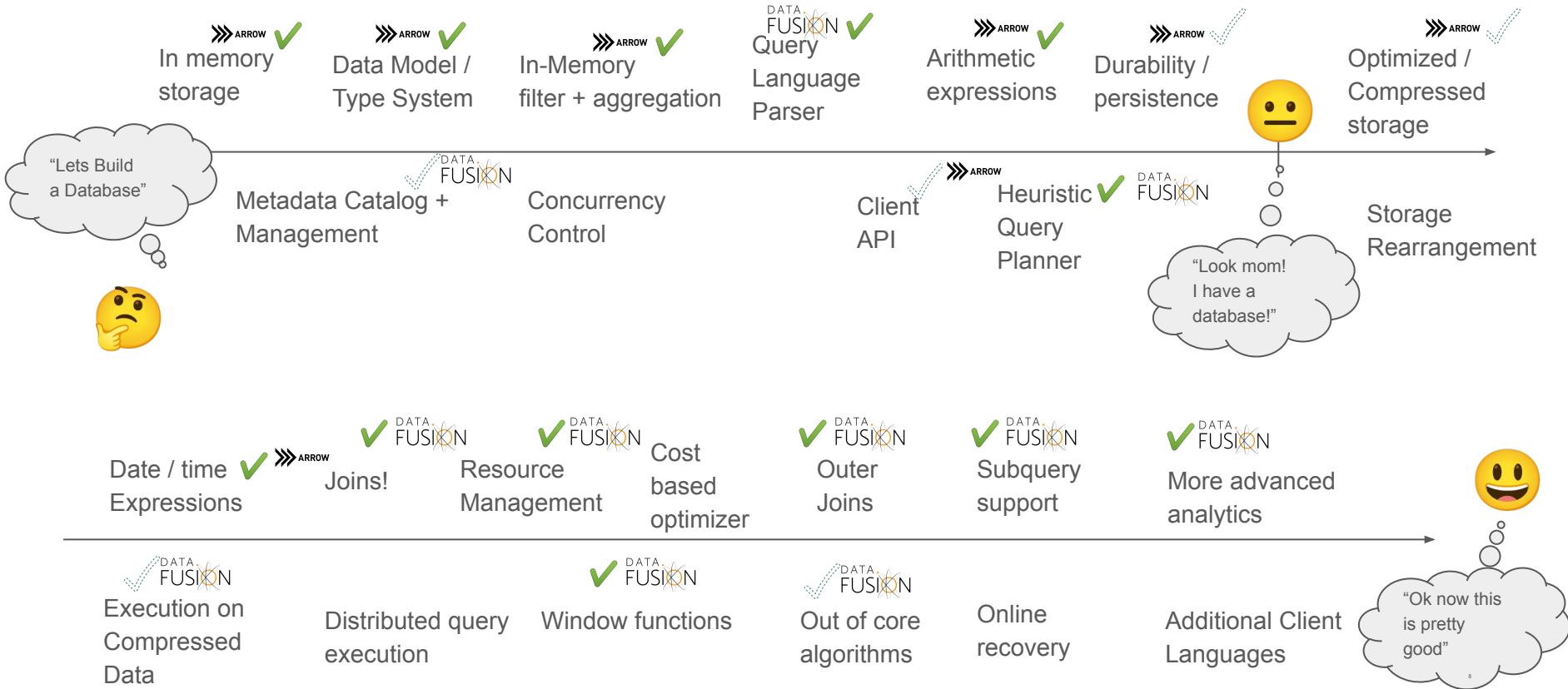
Analogy: DataFusion is LLVM for Databases



[DataFusion](#) enables innovation in data intensive systems

- High quality reusable SQL planner, optimizer, function library, vectorized operators, etc
- Focus on language design, data management, use case specific features

Implementation timeline for a new Database system



More Reading / Viewing / Background

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

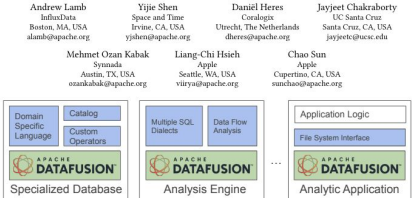


Figure 1: When building with DataFusion, system designers implement domain-specific features via extension APIs (blue), rather than re-implementing standard OLAP query engine technology (green).

ABSTRACT
Apache Arrow DataFusion [1] is a fast, embeddable, and extensible query engine written in Rust [2] that uses Apache Arrow [3] 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 experiments 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 [4].

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

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

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

ACM Reference Format:
Andrew Lamb, Yijie Shen, Daniel Hertz, Coralogix, Jayjeet Chakraborty, Mehmet Ozan Kahak, Liang-Chi Hsieh, and Chao Sun. 2024. Apache Arrow DataFusion: A Fast, Embeddable, Modular Analytic Query Engine. In *Proceedings of the 2024 International Conference on Management of Data (SIGMOD/Conference '24)*, June 19–24, 2024, Seattle, WA, USA. <https://doi.org/10.1145/3632696.3633348>

Santiago
Chile



SIGMOD
PODS In SIGMOD 2024
2024

The screenshot shows the Apache DataFusion website. At the top, there's a navigation bar with the Apache DataFusion logo and the text "Apache DataFusion". Below the navigation bar, there's a search bar and a list of links categorized into "ASF LINKS", "LINKS", "USER GUIDE", and "LIBRARY USER GUIDE". The "ASF LINKS" section includes links to Apache Software Foundation, License, Donate, Thanks, Security, GitHub and Issue Tracker, crates.io, API Docs, Blog, and Code of conduct. The "LINKS" section includes links to GitHub and Issue Tracker, crates.io, API Docs, Blog, and Code of conduct. The "USER GUIDE" section includes links to Introduction, Example Usage, Core Configuration, DataFusion CLI, DataFusion API, Expression API, SQL Reference, Configuration Settings, Reading Explain Plans, and Frequently Asked Questions. The "LIBRARY USER GUIDE" section includes links to Introduction, Extensions List, Using the SQL API, Working with Eger's, Using the DataFusion API, Building Logical Plans, Catalog, Schemas, and Tables, and Links.

The screenshot shows a video player interface. At the top, there's the Apache DataFusion logo and the text "APACHE DATAFUSION". Below the logo, there's the URL "https://datafusion.apache.org/". The main content of the video is "CMU Database Seminar Series, Fall 2024 Database Building Blocks". The video is dated "September 23, 2024" and features "Andrew Lamb, InfluxData". The video player controls are visible at the bottom, showing a play button, a progress bar, and a timestamp of "0:47 / 1:06:51".

The graphic features a metal barrel on the left side. To the right of the barrel, the text reads "Carnegie Mellon University DATABASE BUILDING BLOCKS FALL 2024 SEMINAR SERIES". The text "DATABASE BUILDING BLOCKS" is in a large, bold, black font, and "FALL 2024 SEMINAR SERIES" is in a smaller, red, handwritten-style font.

<https://datafusion.apache.org/>

<https://db.cs.cmu.edu/seminar2024/>

Design Decisions

Programming Language?

Programming Language



Option 1: C/C++

Pros: Well understood, significant track record in databases

Cons: Hard to write correct code. The build system! Macros! etc.



Option 2: Rust

Pros: Memory and Thread safety, Hip language (attracts new developers), modern tooling (e.g. Cargo)

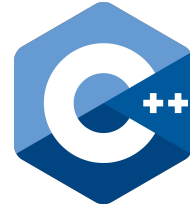
Cons: Not battle tested for Database implementations when choosing

Programming Language

What we did: Rust

Why:

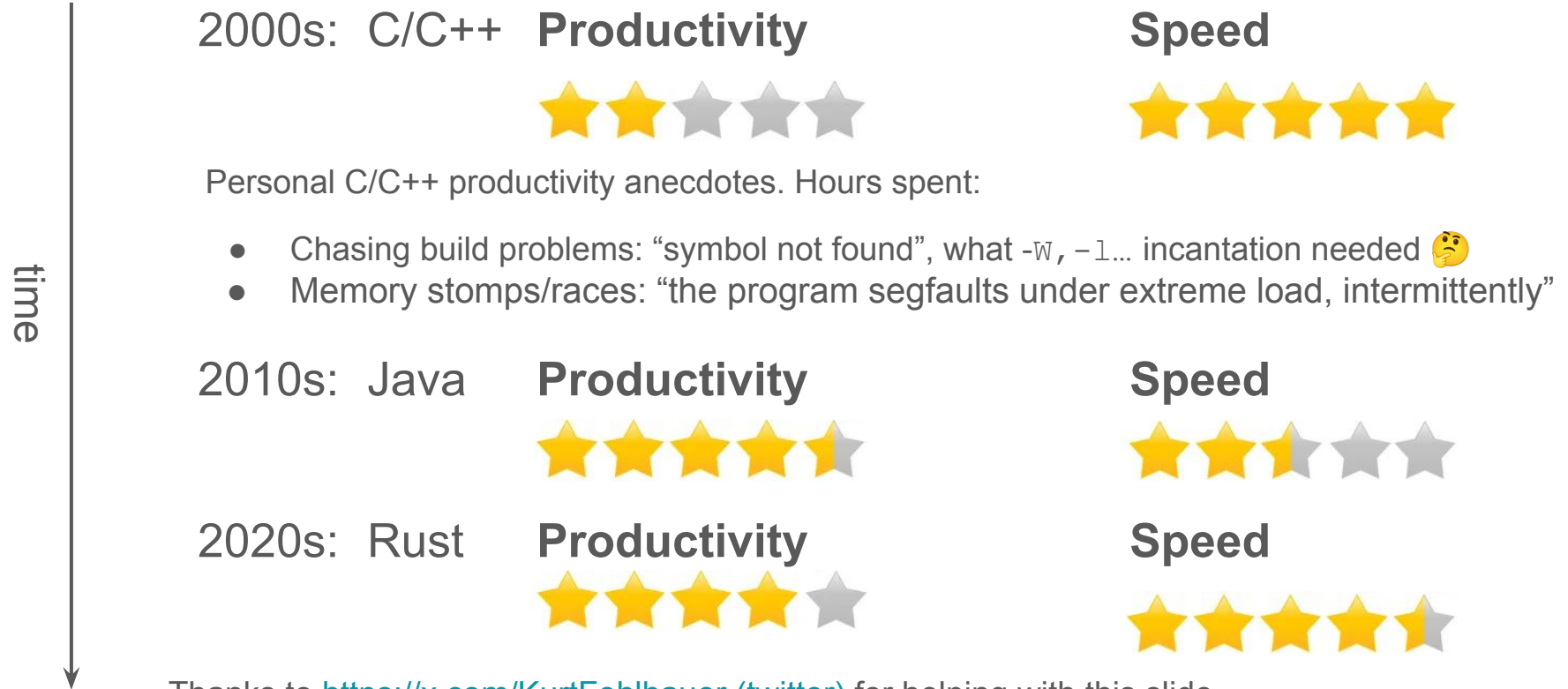
- Initially because Andy Grove believed in Rust*



* See <https://andygrove.io/2018/01/rust-is-for-big-data/>



Lamb Theory on Evolution of Systems Languages



Thanks to <https://x.com/KurtFehlhauer> (twitter) for helping with this slide

Quiz: does this program have undefined behavior?

```
std::vector<int> v { 10, 11 };
```

```
int *vptr = &v[1];
```

```
v.push_back(12);
```

```
std::cout << *vptr;
```



Use after free (maybe)

Points *into* v.

```
let mut v = vec![10, 11];
```

```
let vptr = &mut v[1];
```

```
v.push(12);
```

```
println!("{}", *vptr);
```



Compiler error

>_ codapi

C++ playground

```
1 #include <vector>
2 #include <iostream>
3
4 int main() {
5     std::vector<int> v { 10, 11 };
6     int *vptr = &v[1];
7     v.push_back(12);
8     std::cout << *vptr;
9 }
```




Run

☆ Share

⌵ Fullscreen

11

C++ is a general-purpose programming language that began as an extension of C, but later acquired object-oriented, generic, and functional features.

 [homepage](#) •  [tutorial](#) •  [community](#)

Example: Implement specialized group storage in Rust

```
fn insert_if_new_inner<MP, OP, B>(
    &mut self,
    values: &ArrayRef,
    mut make_payload_fn: MP,
    mut observe_payload_fn: OP,
) where
    MP: FnMut(Option<&[u8]>) -> V,
    OP: FnMut(V),
    B: ByteArrayType,
{
    // step 1: compute hashes
    let batch_hashes = &mut self.hashes_buffer;
    batch_hashes.clear();
    batch_hashes.resize(values.len(), 0);
    create_hashes(&[values.clone()], &self.random_state, batch_hashes)
        // hash is supported for all types and create_hashes only
        // returns errors for unsupported types
        .unwrap();

    // step 2: insert each value into the set, if not already present
    let values = values.as_bytes::<B>();

    // Ensure lengths are equivalent
    assert_eq!(values.len(), batch_hashes.len());
    ...
}
```

```

for (value, &hash) in values.iter().zip(batch_hashes.iter()) {
    // handle null value
    let Some(value) = value else {
        ... (handle nulls here)
    };
    observe_payload_fn(payload);
    continue;
};

// get the value as bytes
let value: &[u8] = value.as_ref();
let value_len = 0::usize_as(value.len());

// value is "small"
let payload = if value.len() <= SHORT_VALUE_LEN {
    let inline = value.iter().fold(0::usize, |acc, &x| {
        acc << 8 | x as usize
    });
    // is value is already present in the set?
    let entry = self.map.get_mut(hash, |header| {
        // compare value if hashes match
        if header.len != value_len {
            return false;
        }
        // value is stored inline so no need to consult buffer
        // (this is the "small string optimization")
        inline == header.offset_or_inline
    });
    if let Some(entry) = entry {
        entry.payload
    }
}

```

```

if let Some(entry) = entry {
    entry.payload
}
// if no existing entry, make a new one
else {
    // Put the small values into buffer
    self.buffer.append_slice(value);
    self.offsets.push(
        0::usize_as(self.buffer.len())
    );
    let payload = make_payload_fn(Some(value));
    let new_header = Entry {
        hash,
        len: value_len,
        offset_or_inline: inline,
        payload,
    };
    self.map.insert_accounted(
        new_header,
        |header| header.hash,
        &mut self.map_size,
    );
    payload
}
// value is not "small"
else {
    ...
}
// Check for overflow in offsets

```

[binary_map.rs](#) from DataFusion

```
// value is not "small"
else {
    // Check if the value is already present in the set
    let entry = self.map.get_mut(hash, |header| {
        // compare value if hashes match
        if header.len != value_len {
            return false;
        }
        // Need to compare the bytes in the buffer
        // SAFETY: buffer is only appended to, and we correctly inserted values and offsets
        let existing_value =
            < unsafe { self.buffer.as_slice().get_unchecked(header.range()) } >;
        value == existing_value
    });
```

[binary_map.rs](#) from DataFusion

Rust: Lessons Learned



- Rust lived up to the hype
 - In 4 years, we had ~1 memory issue, and no multi-threaded bugs / race conditions
- Learning curve is quite steep
 - Be prepared to curse the compiler for a while
- Ecosystem / package manager is amazingly productive
 - `cargo new my_project`
 - `cd my_project`
 - `cargo add datafusion`

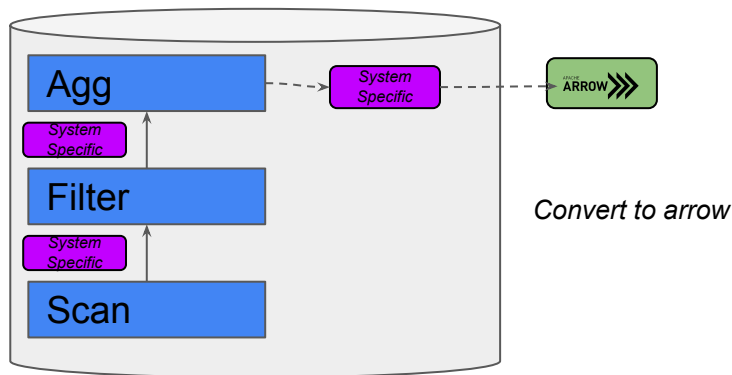
Rust: Lessons Learned



- Improves Open Source Project Velocity
 - Compiler enforces memory safety rather than relying on code reviews
 - Reviewer bandwidth is the most limited resource we have
- Improved Quality of Contributions
 - Non trivial dedication to learn Rust → filtering effect increases contribution quality

Memory Format?

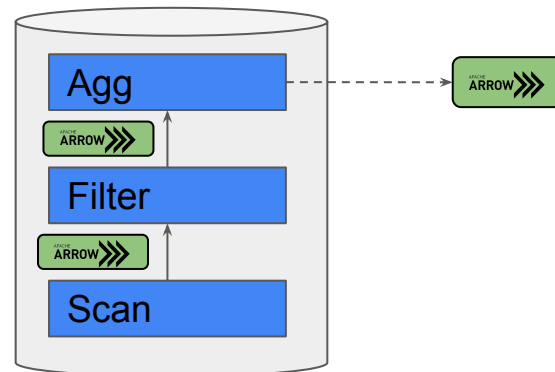
Memory Format



Option 1: Use specialized structures internally, convert to Arrow at edges (Spark, Velox, DuckDB, ...)

Pros: Can use specialized structures

Cons: Maintain specialized code



Option 2: Use Arrow Internally (pola.rs, Acero)

Pros: Fast interchange, reuse Arrow libraries, UDF* become trivial

Cons: Constrained(*) to Arrow

Memory Format

What we did:

- Used Apache Arrow

Why:

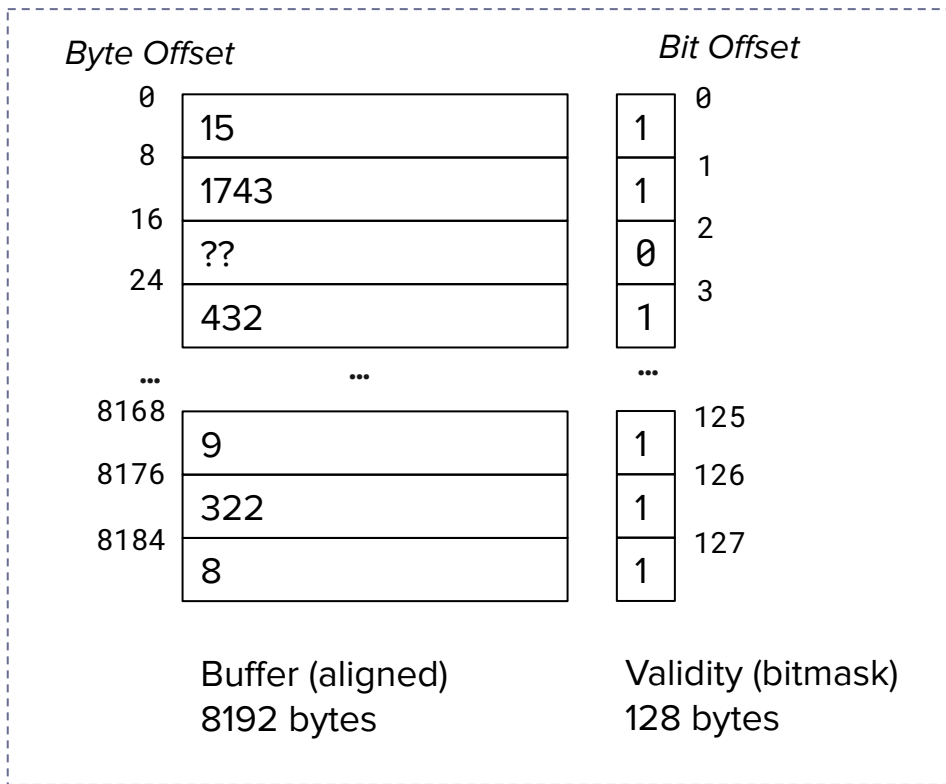
- Theory: Using Arrow is “good enough” compared to specialized structures



Arrow Array: Int64Array

| |
|------|
| 15 |
| 1743 |
| NULL |
| 432 |
| .. |
| 9 |
| 322 |
| 8 |

Logically 1024
8 byte integers



Arrow Array

Pretty much
what you will find
in every
vectorized
column store
engine

Compute Kernels



| | | | |
|----|---|----|---|
| 10 | > | 2 | 1 |
| 20 | | 33 | 0 |
| 17 | | 2 | 1 |
| 5 | | 1 | 1 |
| 23 | | 6 | 1 |
| 5 | | 7 | 0 |
| 9 | > | 8 | 1 |
| 12 | | 2 | 1 |
| 4 | | 7 | 0 |
| 5 | | 2 | 1 |
| 76 | > | 5 | 1 |
| 2 | | 6 | 0 |
| 3 | | 7 | 0 |
| 5 | > | 8 | 0 |

left right output

```
let output = gt(  
  &left,  
  &right  
);
```

The *gt* (greater than) kernel computes an output `BooleanArray` where each element is `left > right`

Kernels handle nulls (validity masks), optimizations for different data sizes, etc.

~50 different kernels, full list: [docs.rs page](#)



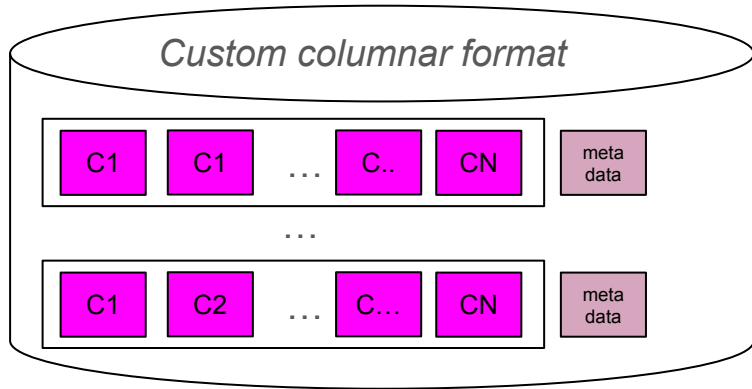
Memory Format - Lessons Learned

- Upstream wasn't quite ready \Rightarrow needed lots of help
- arrow-rs optimized kernels were as important as layout
- Missing Features: * Selection Vectors / "String Views" / RLE encoding
- Single constant value (ScalarValue) should have been in Rust Arrow
- Awkward that Arrow `DataTypes` both logical and physical (`DictionaryArray`)

* RLE + StringView were added later

Storage/File Format?

File Format

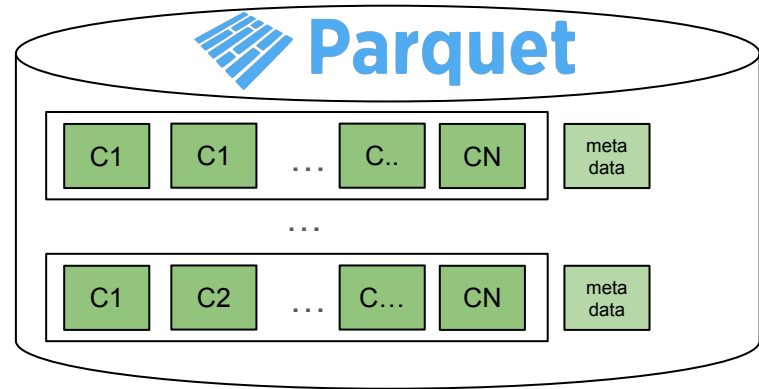


PAX style, encodings, statistics in Zone Maps, etc

Option 1: Custom format (e.g. [DuckDB](#), Snowflake, Vertica, ...)

Pros: Can use specialized structures, encodings, control the format

Cons: Maintain specialized code, pay to copy data in / out of this format



Option 2: Existing Format (e.g. Parquet)

Pros: Well understood, ecosystem interoperability

Cons: Constrained(*) to formats + existing implementations

File Format

What we did:

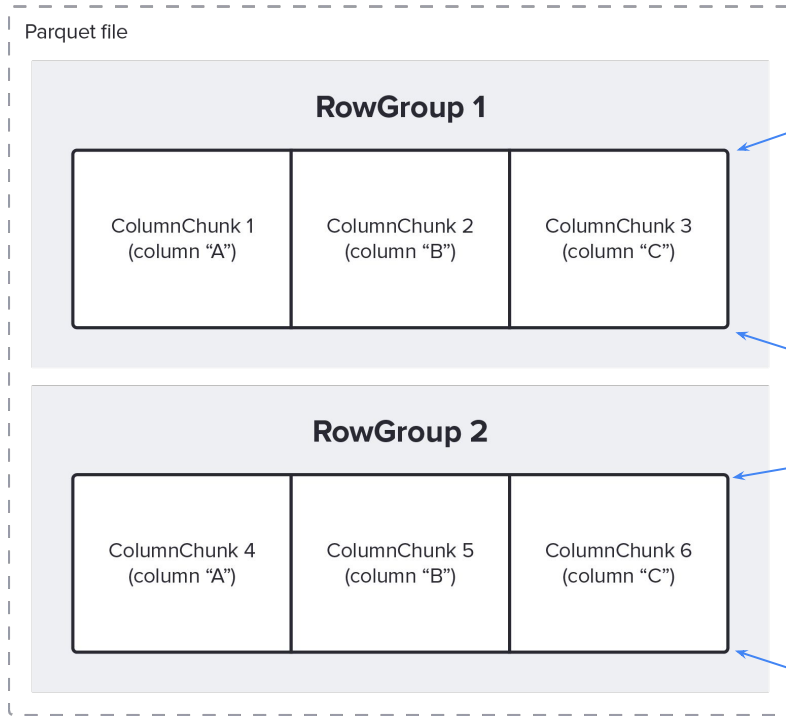
- Use Parquet, Avro, Json, CSV, Arrow
- Extension APIs for others

Why?

- Parquet has enough (Pax, Bloom Filters, Zone Maps)
- Huge amount of Parquet already out there (table stakes!)



Parquet Organization

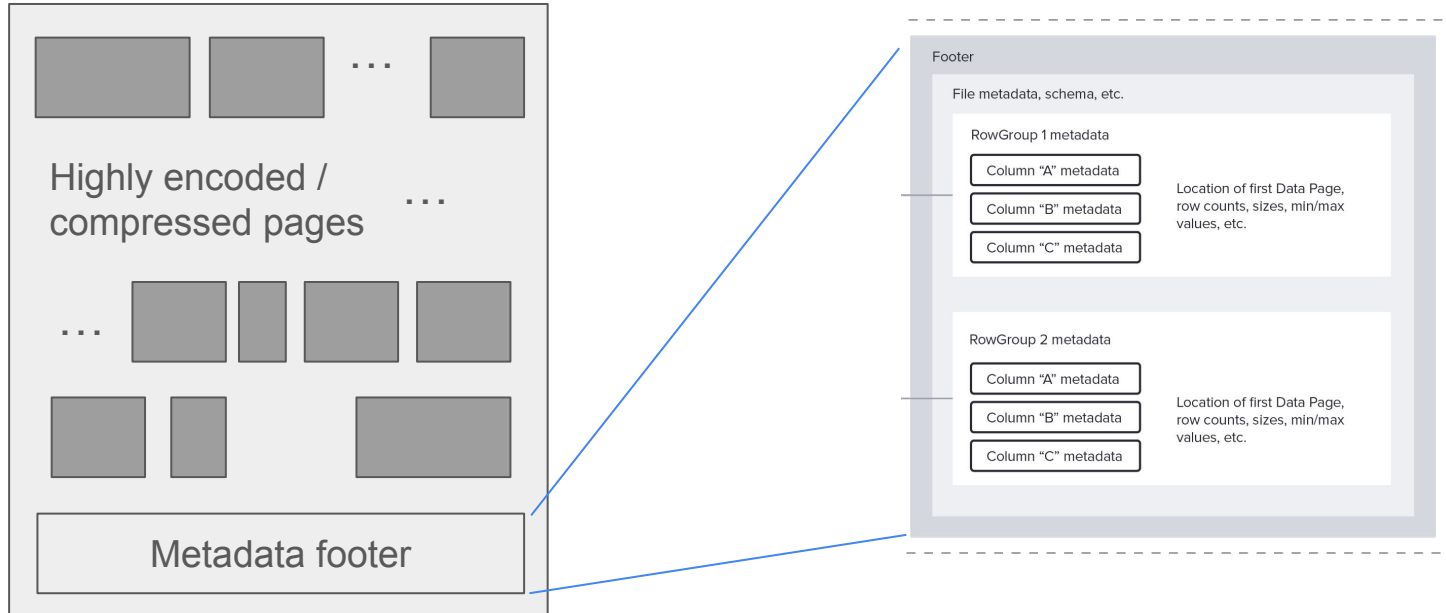


| A | B | C |
|-----|-------|------------|
| 15 | Foo | 1/1/2023 |
| ... | | |
| 11 | Bar | ..1/5/2023 |
| 50 | Baz | 1/1/2023 |
| ... | | |
| 32 | Blarg | 1/6/2023 |

Source: <https://www.influxdata.com/blog/querying-parquet-millisecond-latency/>

(“PAX” in DB literature)

Parquet Structure + Metadata



Footer contains location of pages, and statistics such as min/max/count/nullcount.

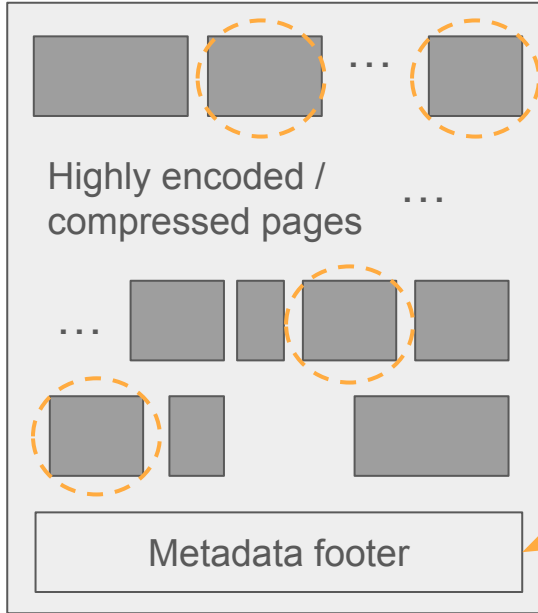
("Zone Maps", "Small Materialized Aggregates" in DB literature)

Source:

<https://www.influxdata.com/blog/querying-parquet-millisecond-latency/>

Parquet Projection + Filter Pushdown

2. Only read/decode necessary pages



Metadata + query to prune (skip) pages that aren't needed

```
SELECT A  
...  
WHERE C > 25
```

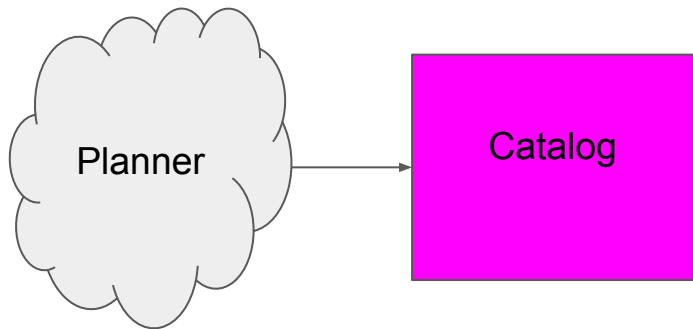
1. Consult metadata

File Format - Lessons Learned

- Standard formats \Rightarrow large community interested, able and willing to help
- Parquet has many optimization opportunities
 - Rust Parquet implementation is now really good
 - Statistics Pruning (file, row group, page index)
 - Filter pushdown / Late Materialization / IO Interleaving
- Leveraging existing format and invest heavily in the software implementation
 - Work with community to evolve rather than replace Parquet

Catalog?

Catalog Format

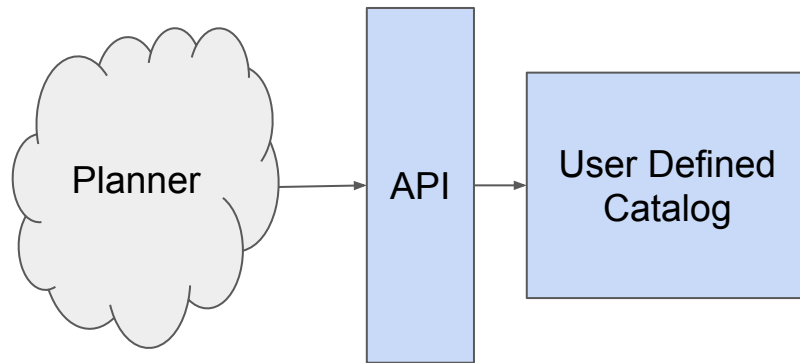


Concrete catalog implementation

Option 1: Provide catalog (e.g. `.sqlite`, `.duckdb`, `Iceberg`, etc)

Pros: Fast to start using

Cons: catalog implementation bound to usecase (e.g. local files vs remote service), planner may be more coupled to catalog



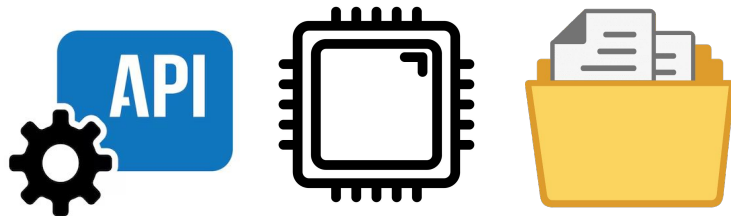
User provides the catalog implementation

Option 2: Provide an API (e.g. Calcite)

Pros: can tailor the catalog to needs, planner not coupled to catalog

Cons: higher startup cost (have to implement catalog)

Catalog Format



What we did:

- 2 simple built in catalogs (memory + file based / Hive-style partitioning)
- APIs for others

Why:

- Catalog format is very usecase / system dependent
- Nothing we could have built in would likely work well

Listing Table Catalog

```
SELECT ...  
FROM table2  
WHERE date='2024-06-02'
```

SQL Query

Tables:

- Directories of files
- “Standard” hive-style directory partitioning

```
table1  
├── file1.parquet  
└── file2.parquet
```

```
table2  
├── date=2024-06-01  
│   ├── file3.parquet  
│   └── file4.parquet  
└── date=2024-06-02  
    └── file5.parquet
```

Filesystem directory structure



Catalog Format - Lessons Learned

- Providing basic “get started” implementation and extension APIs worked great
 - Let people start quickly, but customize as needed
- ListingTable (directory Filesystem):
 - More complicated than expected
 - Should have had a cleaner separation from the start

Planning?

SQL Dialects

```
SELECT
  age,
  sum(civility) AS total_civility
FROM star_wars_universe
GROUP BY ALL
ORDER BY ALL;
```

Friendlier SQL with DuckDB

Option 1: Implement your own

Pros: have exactly the semantics you want, friendlier language

Cons: It is a lot of work (implementation *and* education)

```
SELECT
  age,
  sum(civility) AS
total_civility
FROM star_wars_universe
GROUP BY age
ORDER BY age, total_civility;
```

“Standard” SQL

Option 2: Use existing dialect

Pros: Avoid having to define semantics

Cons: Bug for bug compatible. Have to pick one dialect

SQL Dialects



What we did:

- Emulate postgres semantics default
- Extension APIs

Why:

- Well understood
- It is time consuming / expensive to invent semantics

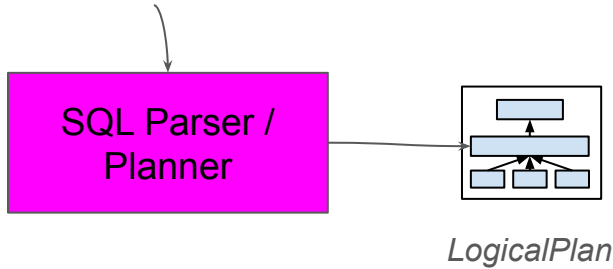
SQL Dialects - Lessons Learned

- Dialect syntax only part of the story. Others:
 - Function Library and semantics (e.g `null` or error on invalid args?)
 - DataTypes (`VARCHAR2`? `CHAR` / `VARCHAR`?)
 - Type Coercion Rules
- No one ideal choice: Spark Dialect vs Postgres (A: UDFs!)
- Postgres + Arrow timestamp representation impedance mismatch
- On the whole this was still a good idea

SQL Planner

```
SELECT status, COUNT(1)
FROM http_api_requests_total
WHERE ...
```

SQL Text



Option 1: Implement sql parser / planner

Pros: Minimize dependencies, native integration into plan structures / exprs

Cons: Much more work

```
SELECT status, COUNT(1)
FROM http_api_requests_total
WHERE ...
```

SQL Text



Option 2: Calcite

Pros: Mature

Cons: Java (dependencies), bridge plan representation to internal representation

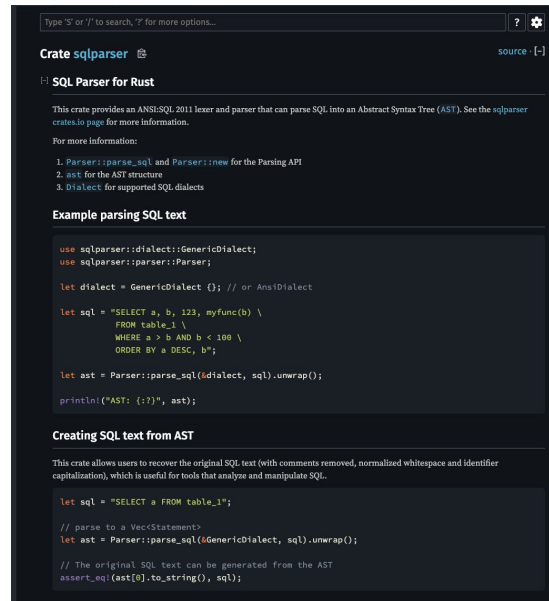
SQL Planner

What we did

- Implemented own sqlparser-rs and planner in Rust
- + extensions

Why:

- Avoid Java dependencies / have a pure Rust stack



The screenshot shows the documentation for the `sqlparser` crate. At the top, there is a search bar and a settings icon. Below that, the crate name `Crate sqlparser` is displayed. The main heading is `SQL Parser for Rust`. A brief description states: "This crate provides an ANSI-SQL 2011 lexer and parser that can parse SQL into an Abstract Syntax Tree (AST). See the `sqlparser` `crates.io` page for more information." Below this, there are three numbered items for more information: 1. `Parser::parse_sql` and `Parser::new` for the Parsing API, 2. `ast` for the AST structure, and 3. `Dialect` for supported SQL dialects. A section titled "Example parsing SQL text" contains a Rust code snippet that demonstrates parsing a SQL query. The code defines a `GenericDialect`, creates a `Parser`, and uses `Parser::parse_sql` to parse a query. The output is printed as an AST. Below this, a section titled "Creating SQL text from AST" shows a code snippet that demonstrates how to generate SQL text from an AST node.

```
Type 'S' or '/' to search, '?' for more options...
Crate sqlparser
SQL Parser for Rust

This crate provides an ANSI-SQL 2011 lexer and parser that can parse SQL into an Abstract Syntax Tree (AST). See the sqlparser crates.io page for more information.

For more information:
1. Parser::parse_sql and Parser::new for the Parsing API
2. ast for the AST structure
3. Dialect for supported SQL dialects

Example parsing SQL text

use sqlparser::dialect::GenericDialect;
use sqlparser::parser::Parser;

let dialect = GenericDialect {}; // or AnsiDialect

let sql = "SELECT a, b, 123, myfunc(b) \
FROM table_1 \
WHERE a > b AND b < 100 \
ORDER BY a DESC, b";

let ast = Parser::parse_sql(&dialect, sql).unwrap();

println!("AST: {:?}", ast);

Creating SQL text from AST

This crate allows users to recover the original SQL text (with comments removed, normalized whitespace and identifier capitalization), which is useful for tools that analyze and manipulate SQL.

let sql = "SELECT a FROM table_1";

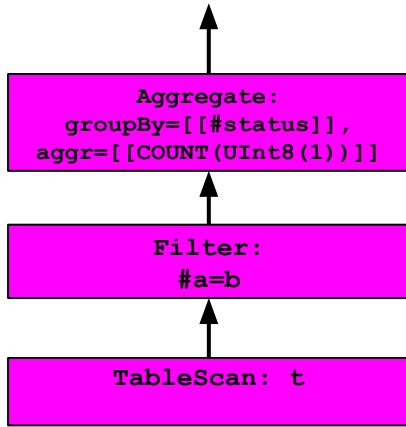
// parse to a Vec<Statement>
let ast = Parser::parse_sql(&GenericDialect, sql).unwrap();

// The original SQL text can be generated from the AST
assert_eq!(ast[0].to_string(), sql);
```

SQL Planner - Lessons Learned

- Implementing a SQL planner is a LOT of work
 - SQL is a crazy creole language
- Modularity helped
 - sqlparser-rs has a clean split from DataFusion
 - Means it is used by many projects, and thus benefits from larger community
- Would recommend avoiding this if you can

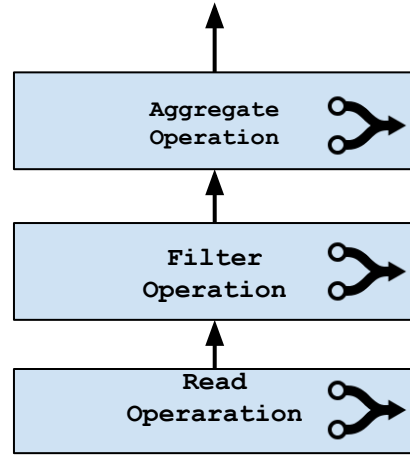
Plan Representation / API



Option 1: Custom Structures

Pros: Native APIs, make it ergonomic

Cons: (very) large API surface area + code to maintain, have to define semantics precisely



Option 2: Use existing library

Pros: Code is simple, planning is predictable

Cons: Limited to whatever is available

Plan Representation / API

What we did:

- Custom structures and API
- + Extension APIs

Why:

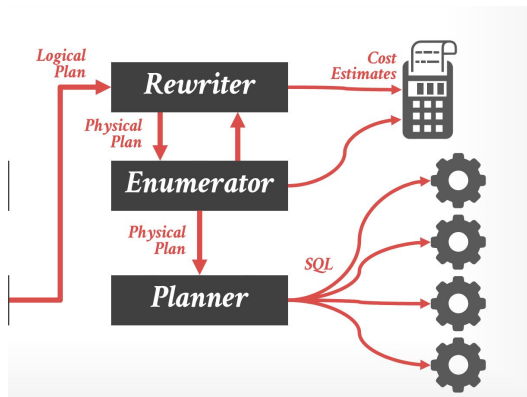
- No compelling alternative



Plan Representation / API - Lessons Learned

- It is a lot of code and API design
- TreeNode API is quite cool (unified Expr/Plan walking)
- Custom serialization takes takes lots of time
- Would / should have used [substrate](#) if it was ready
- Should have used rewrite in place APIs for performance reasons

Cost Based / Stratified / Unified / Join Order Optimizer

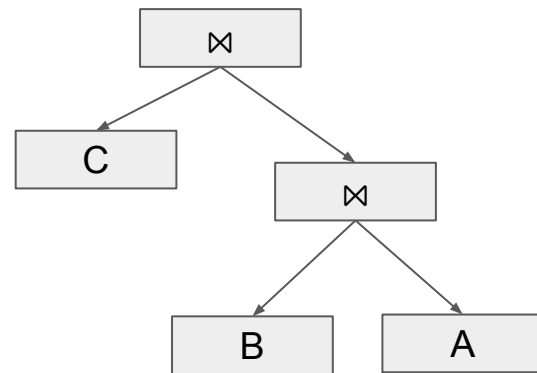


Option 1: CBO + Heuristics

Pros: Well understood pattern in DBs

Cons: Known hard problem: cost estimates, cardinality estimates, correlations, performance cliffs, etc

```
SELECT ...  
FROM  
  A  
  JOIN B ON ...  
  JOIN C ON ...
```



Option 2: “Syntactic” optimizer (whatever order user tells you)

Pros: Code is simple, planning is predictable

Cons: Complex join orders ⇒ 😞

Cost Based / Stratified / Unified / Join Order Optimizer

What we did:

- Syntactic optimizer
- + Just enough to avoid TPCH disasters

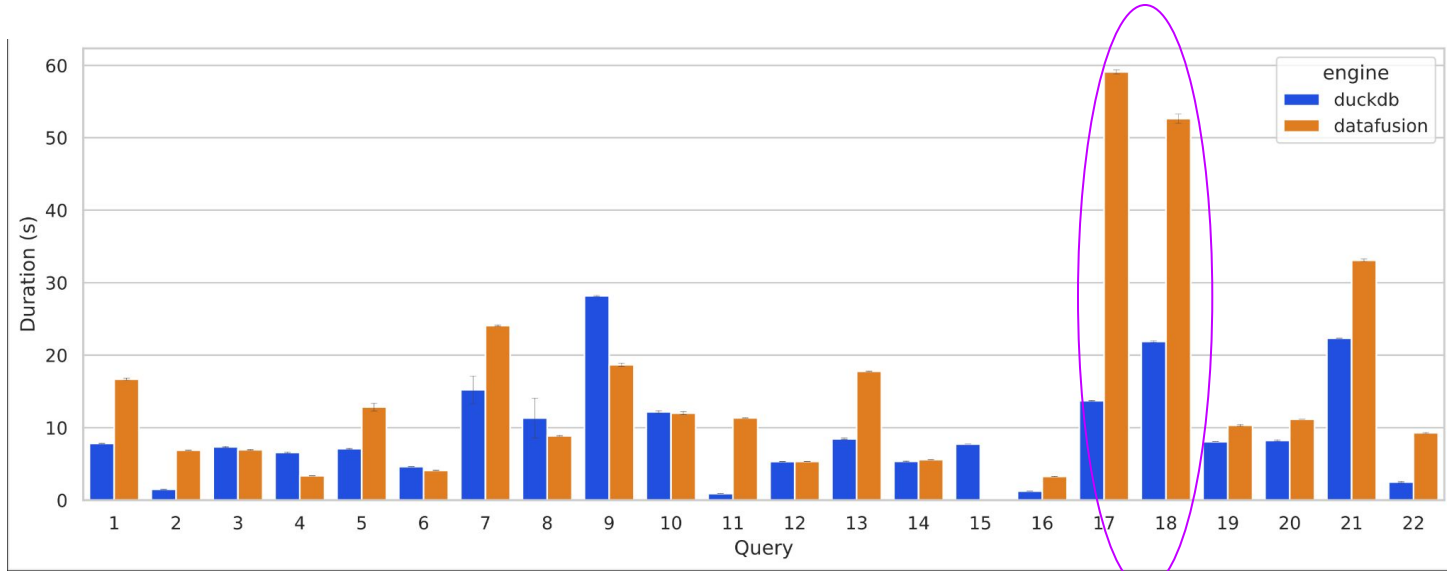
Why:

- Handling complex join reordering is hard (both theoretically and practically):
Depends a lot on cost model + accurate statistics (also hard)
- Denormalized tables very common in olap workloads (so join order relatively less important)
- Users can implement more sophisticated strategies as rewrites



Cost Based Optimizer: Lessons Learned

- Worked OK (Some TPCCH embarrassments)



Join order disaster (subquery cardinality estimation, fixed in #7949)

Cost Based / Stratified / Unified / Join Order Optimizer

- Example of implement Join Ordering as user defined rule

“Currently, optd is integrated into Apache Arrow Datafusion as a physical optimizer. It receives the logical plan from Datafusion, implements various physical optimizations (e.g., determining the join order), and subsequently converts it back into the Datafusion physical plan for execution.”

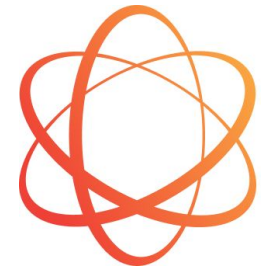
<https://github.com/cmu-db/optd>

Conclusion 😄💧

Conclusion and Takeaways

- Analytic Systems take a lot of work
- Rust and Apache {Arrow, Parquet, DataFusion} are awesome
- Reusing open building blocks saved lots of effort
 - Not free: contributed a lot back to help make them better
- Basic implementations + Extension APIs: kept core “simple”
 - Same APIs for Built in and User Defined
 - Forces the API to be complete / no special casing (e.g. UDFs)

Thank you



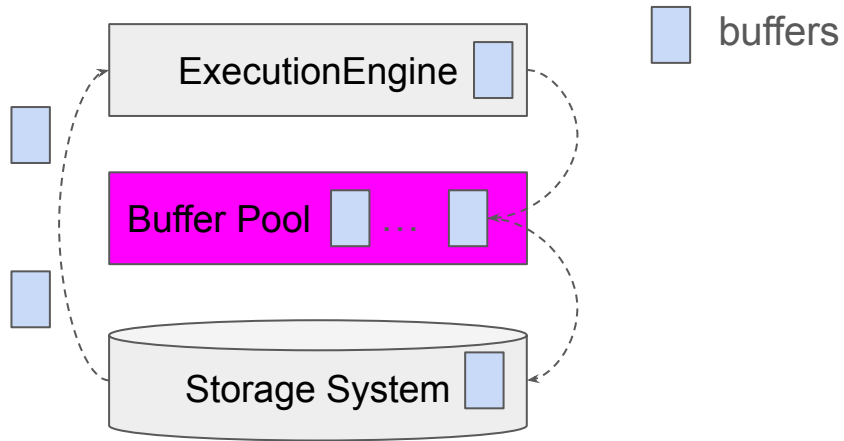
A P A C H E

DATAFUSIONTM

Come join us!

<https://datafusion.apache.org/>

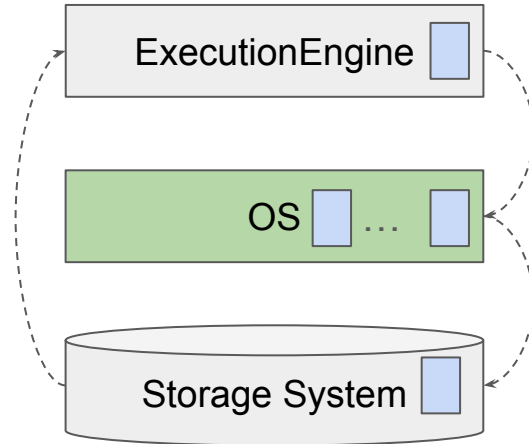
Buffer Pool



Option 1: Custom Bufferpool

Pros: Well understood pattern in DBs

Cons: Significant system complexity, have to manage memory distribution between I/O, execution, etc. Tune pool to workload



Option 2: OS Alloc + Page Cache

Pros: No code to manage

Cons: Beholden to OS, Potential MMAP 🤩 situations

Buffer Pool

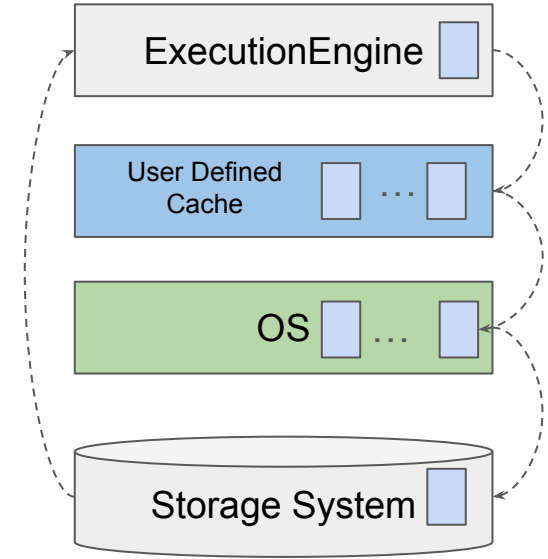
What we did: Use OS + User defined cache

Why: Simple

- Optimal Caching strategy almost always highly dependent on system / environment
- Can implement caching strategies (aka buffer pool) via extensions

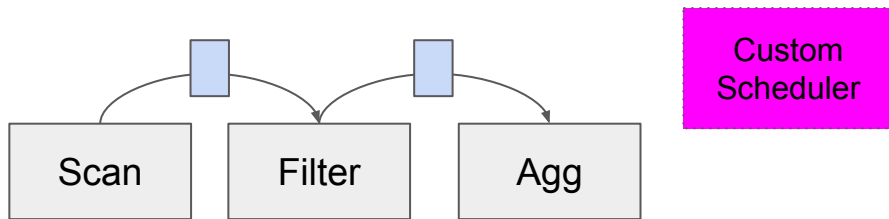
Lessons Learned:

- This has worked very well – basic implementation is easy to understand and very predictable



Backup Content

Execution Engine Scheduler

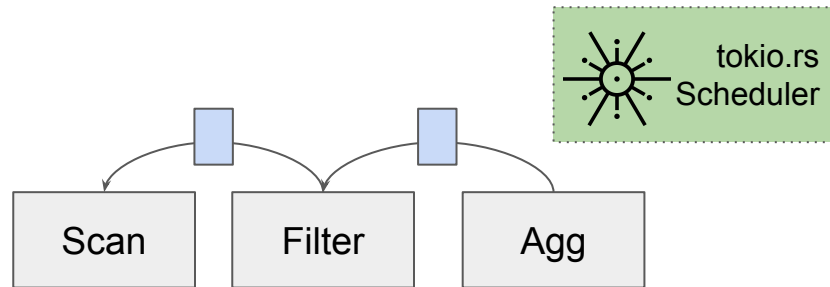


Scheduler "pushes" blocks from scan through plan

Option 1: Write own (push based) scheduler

Pros: Tight control over behavior, prioritization, etc

Cons: Very hard to write correctly and tune well, especially under load, network backpressure, etc



Scheduler "pulls" blocks from scan through plan

Option 2: Use tokio scheduler (pull)

Pros: Someone else writes scheduler and tools, integrated IO + CPU (tokio) patterns, already present in many rust apps

Cons: Less control

Execution Engine Scheduler



What we did: Used Tokio + Futures + Rust async continuations

Why:

- Super well tested, built in compiler language support and tools
- Didn't really have budget to make our own

Lessons Learned:

- I would do it again (though I will admit others are not convinced)
- Performance turned out to be no different than custom scheduler
- TUM papers are convincing (esp with DuckDB's story)
- Easy to run IO and CPU on the same pool, so care is warranted (though I would argue it was good we didn't need to worry about this until it actually mattered for our scale)

More details in <https://thenewstack.io/using-rustlangs-async-tokio-runtime-for-cpu-bound-tasks/>

Execution Engine Scheduler: Morsels?

We actually tried to implement a push based (morsel driven) scheduler

“the rayon-based [push based] scheduler in its current incomplete incarnation provides minimal benefits over the current tokio-based approach, and is a non-trivial amount of fairly complex code.” - [removal PR](#)

More details in the [Epic](#)

Morsel-Driven Parallelism: A NUMA-Aware Query Evaluation Framework for the Many-Core Age

Viktor Leis^{*} Peter Boncz[†] Alfons Kemper^{*} Thomas Neumann[†]

^{*} Technische Universität München [†] CWI

^{*}{leis,kemper,neumann}@in.tum.de [†]p.boncz@cw.nl

ABSTRACT

With modern computer architecture evolving, two problems conspire against the state-of-the-art approaches in parallel query execution: (i) to take advantage of many-cores, all query work must be distributed evenly among (soon) hundreds of threads in order to achieve good speedup, yet (ii) dividing the work evenly is difficult even with accurate data statistics due to the complexity of modern out-of-order cores. As a result, the existing approaches for “plan-driven” parallelism run into load balancing and context-switching bottlenecks, and therefore no longer scale. A third problem faced by many-core architectures is the decentralization of memory controllers, which leads to Non-Uniform Memory Access (NUMA). In response, we present the “morsel-driven” query execution framework, where scheduling becomes a fine-grained run-time task that is NUMA-aware. Morsel-driven query processing takes small fragments of input data (“morsels”) and schedules these to worker threads that run entire operator pipelines until the next pipeline breaker. The degree of parallelism is not baked into the plan but can elastically change during query execution, so the dispatcher can react to execution speed of different morsels but also adjust resources dynamically in response to newly arriving queries in the workload. Further, the dispatcher is aware of data locality of the NUMA-local morsels and operator state, such that the great majority of executions takes place on NUMA-local memory. Our evaluation on the TPC-H and SSB benchmarks shows extremely high absolute performance and an average speedup of over 30 with 32 cores.

In response, we present the “morsel-driven” query execution framework, where scheduling becomes a fine-grained run-time task that is NUMA-aware. Morsel-driven query processing takes small fragments of input data (“morsels”) and schedules these to worker threads that run entire operator pipelines until the next pipeline breaker.

The degree of parallelism is not baked into the plan but can elastically change during query execution, so the dispatcher can react to execution speed of different morsels but also adjust resources dynamically in response to newly arriving queries in the workload. Further, the dispatcher is aware of data locality of the NUMA-local morsels and operator state, such that the great majority of executions takes place on NUMA-local memory. Our evaluation on the TPC-H and SSB benchmarks shows extremely high absolute performance and an average speedup of over 30 with 32 cores.

Categories and Subject Descriptors

H.2.4 [Systems]: Query processing

Keywords

Morsel-driven parallelism, NUMA-awareness

1. INTRODUCTION

The main impetus of hardware performance improvement nowadays comes from increasing multi-core parallelism rather than from speeding up single-threaded performance [2]. By SIGMOD 2014

Permission to make digital or hard copies of part or all 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. Copyright for third-party components of this work must be honored. For all other uses, contact the owner(s). Copyright is held by the author(s).

SIGMOD '14, June 22–27, 2014, Snowbird, UT, USA.
SCM 978-1-4503-2316-5/14/06.
http://dx.doi.org/10.1145/258555.2610507.

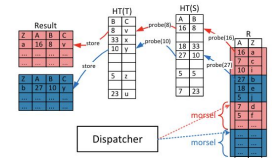


Figure 1: Idea of morsel-driven parallelism: $R \bowtie_A S \bowtie_B T$

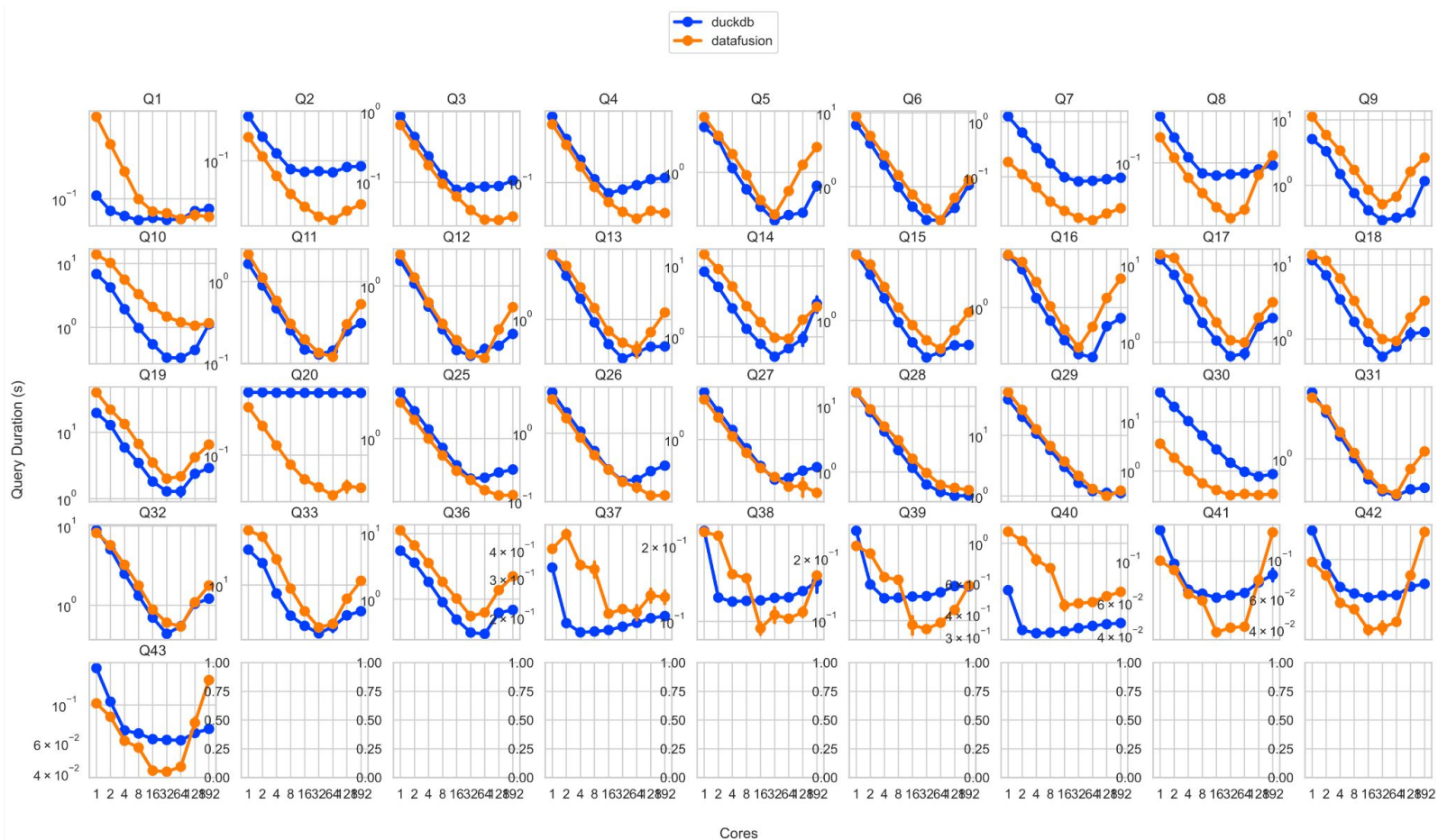
Intel’s forthcoming mainstream server Ivy Bridge EX, which can run 120 concurrent threads, will be available. We use the term *many-core* for such architectures with tens or hundreds of cores.

At the same time, increasing main memory capacities of up to several TB per server have led to the development of main-memory database systems. In these systems query processing is no longer I/O bound, and the huge parallel compute resources of many-cores can be truly exploited. Unfortunately, the trend to move memory controllers into the chip and hence the decentralization of memory access, which was needed to scale throughput to huge memories, leads to non-uniform memory access (NUMA). In essence, the computer has become a network in itself as the access costs of data items varies depending on which chip the data and the accessing thread are located. Therefore, many-core parallelization needs to take RAM and cache hierarchies into account. In particular, the NUMA division of the RAM has to be considered carefully to ensure that threads work (mostly) on NUMA-local data.

Abundant research in the 1990s into parallel processing led the majority of database systems to adopt a form of parallelism inspired by the Volcano [12] model, where operators are kept largely unaware of parallelism. Parallelism is encapsulated by so-called “exchange” operators that route tuple streams between multiple threads each executing identical pipelined segments of the query plan. Such implementations of the Volcano model can be called *plan-driven*: the optimizer statically determines at query compile-time how many threads should run, instantiates one query operator plan for each thread, and connects these with exchange operators.

In this paper we present the adaptive *morsel-driven* query execution framework, which we designed for our main-memory database system Hyper [16]. Our approach is sketched in Figure 1 for the three-way-join query $R \bowtie_A S \bowtie_B T$. Parallelism is achieved

Paper: [Morsel-Driven Parallelism: A NUMA-Aware Query Evaluation Framework for the Many-Core Age](#)



ClickBench Scaling of DataFusion (orange) vs DuckDB (blue) 1 - 172 cores \Rightarrow shapes are very similar

More details in our SIGMOD paper: <https://s.apache.org/datafusion-sigmod-2024>

Governance



[“Benevolent Dictator for Life”](#), “Lone Maintainer”, etc

Option 1: Custom Governance

Pros: Very flexible, custom tailored
doesn't slow down

Cons: Can be non trivial to setup,
unclear governance long term
sustanaibility



Foundation

Option 2: Open, existing system

Pros: Predictable well understood
governance

Cons: Often slower / not ideal process, may
be hard to join

Governance



What we did: Open Governance via Apache Software Foundation

Why:

- Allows contributors from many companies
- Makes governance clear and predictable, especially growing capacity

Lessons Learned

- Totally open / public community takes getting used to
- Works shockingly well
- Not the fastest way of making decisions

* Full Disclosure: I am an ASF foundation member (via work on Arrow + Datafusion)

Architecture / Roadmap



Cathedral

Option 1: Tight Control

Pros: Unified architecture, coherent APIs, decided by a small group of people

Cons: scaling up, limits contributors

Eric Raymond: https://en.wikipedia.org/wiki/The_Cathedral_and_the_Bazaar



Bazaar

Option 2: Accept anything

Pros: Can draw from many people

Cons: Can end up with frankenstein APIs, half finished features

Architecture / Roadmap

What we did: default accept features

- Datafusion-contrib: outlet for people to contribute without as many constraints / still have community

Why:

- Encourages growth and contribution as people feel control over direction

Lessons Learned

- Shockingly effective
- Insist on tests for all PRs: likely more important than the actual code
- Code quality can vary (which is ok to clean it up over time, with tests)
- Need to invest in tracking follow on work with tickets

Architecture / Roadmap

Example contributor driven features

- Predicate pushdown into TableProvider - Thanks @returnString
- Window functions – Thanks @jimexist
- COPY statement, parallelized parquet writer – Thanks @devinjdangelo
- Parquet bloom filter / data page pruning – Thanks @Ted-Jiang

Examples of half finished things

- Recursive CTEs (then @jonahgao picked it up and brought it over the line)
- Predicate selectivity estimation

But quality?

The New York Times

Boaty McBoatface: What You Get When You Let the Internet Decide

Share full article



A computer image of the research vessel, which is still being designed and is scheduled to set sail in 2019. The Natural Environment Research Council

“Originally suggested by former BBC radio presenter James Hand, by the end of the poll on 16 April Boaty McBoatface had garnered 124,109 votes and 33% of the total vote.”

⇒ “Science Minister Jo Johnson said there were “more suitable” names.”

But quality?

- Quality is quite high
- Bugs and regressions do get introduced
- But fixed almost as soon as they are filed (often by the original author)
- I believe quality is better than many of the enterprise software systems I have worked on

Example: Shocking effectiveness of tickets

The screenshot shows a GitHub issue titled "Add API to read a `Vec<RecordBatch>` from `SessionContext` #9157". The issue is marked as "Closed" and "Fixed by #9197". It was opened by user "alamb" on Feb 7. The main comment from "alamb" asks, "Is your feature request related to a problem or challenge?" and provides context about APIs in DataFusion. A second comment from "Lordworms" says "I can do this one" and is timestamped "Feb 7, 2024, 9:25 PM EST". A third comment from "alamb" says "Thank you @Lordworms". A large yellow text overlay reads "16 Minutes between file and claim". Two red arrows point from this text to the timestamps of the first and second comments.

🤖: File ticket and often there is a PR up within 24 hours to solve it

⇒ “file clear requests” and people want to help

It does take non-coding effort

Make the engine Distributed Engine

What we did instead: Focused on single core, features to build distributed engine

Plan Serialization + use arrow flight

Why:

Distribution strategy varies greatly on usecase

All can be built with a single core

Had other projects (like Ballista) that added distribution

Lessons Learned:

Optimized Comparisons

Why: optimize multi column sort

Add RowFormat to Arrow (see Blog)

Used in merge, multi-column grouping and join comparisons

Lessons Learned:

- Dictionary interning was a bad idea (memory exploded)
-

Manually vectorized kernels

What we did instead: Use autovectoruixation in Rust

Why: Partly Rust's SIMD intrincs story wasn't ideal, but we found we could often get better results with careful rust code than with custom kernels

Lessons Learned:



Optimized Hash Aggregates

Why: aggregation key

What: Two phase hash aggregate (todo blog post link)

Special column based key storage, type specialized, etc

Lessons Learned:

- Need special case code for different types
- Rust auto-vectorization is quite nice

Function Library

Why: Turns out a lot of SQL functionality

What: Massive library (TODO quantify) of functions, all use the same API as user functions

Lessons Learned:

- Had a split initially between built in functions and udfs
- Should have had easier UDFs earlier on
- That way behavior can be customized more easily rather than the core ever growing

Pruning Predicate

Why: Key building block for many optimizations, including fast parquet reader

What: given expression and min/max/nullcounts (and bloom filter info) tell “can this possibly have rows that evaluate to true)?

Lessons Learned:

- API needs to be vectorized (quickly prune large numbers of files)
- Pruning is crazy tricky (because you are trying to prove the null) – testing

Listing Table: Hive style partitioned data

Why: lots of data is in this format (predates table formats)

What: TODO show example directory

Lessons Learned:

- Should have had this outside the core earlier on (as people almost immediately wanted to plug in other file types in there)

Foundations: Apache Parquet

- Open column-oriented data file format,
- Provides efficient data compression and encoding schemes, along with support for structured types via record shredding [56], embedded schema descriptions, zone-map [57] like index structures and Bloom filters for fast data access.

Differences

- Arrow: fast random access and efficient in-memory processing
- Parquet: store large amounts of data in a space-efficient manner.

Why Parquet and not a specialized format?

- De-facto standard for data storage and interchange in the analytic ecosystem.
- Open format, excellent compression across real-world data sets, broad ecosystem and library support, and embedded self-describing schema
- Structure allows query engines (like DataFusion) to apply projection and filter pushdown techniques, such as late materialization, directly on files, yielding competitive performance compared with specialized formats [48].



Parquet

Foundations: Apache Arrow

Apache Arrow

- Standardizes industrial best practices to represent data in memory using cache-efficient columnar layouts.
- Users avoid re-implementing features that are well understood in academia and industry, but time-consuming to implement.

Specifies

- Validity/null representations; Endianness
- Variable length byte and character data; lists, and nested structures,

Benefits

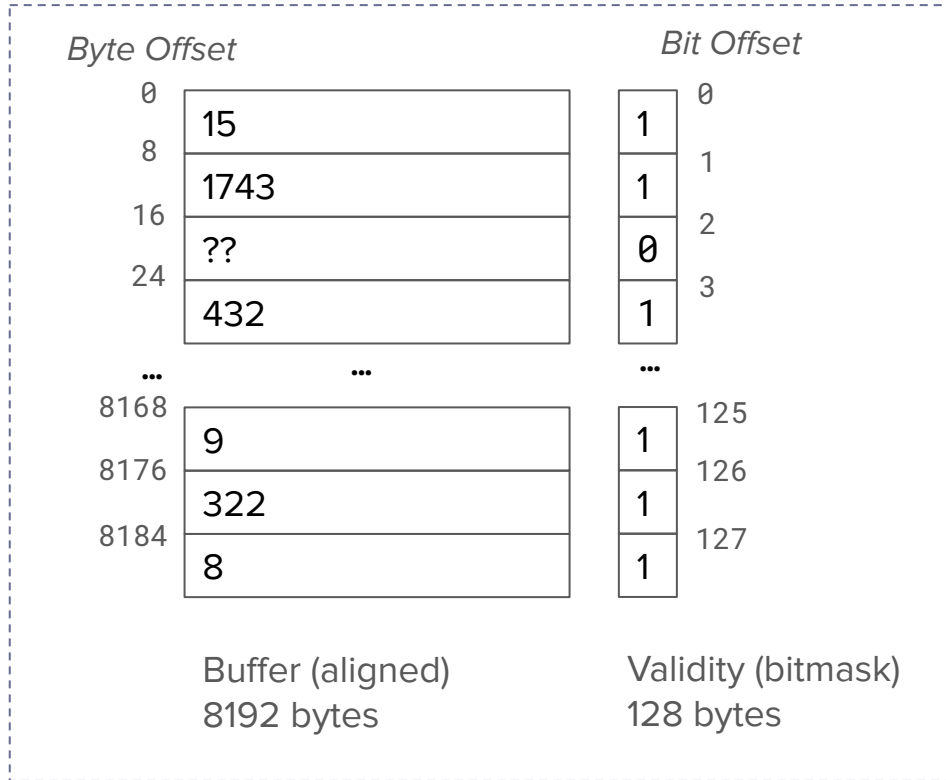
- Well-known techniques (e.g. vectorized compute kernels, special case nulls/no-nulls)
- Easy + zero cost data interchange (e.g. is a NULL value is represented by a 0 or 1 in a bit mask?).
- Arrow evolves over time (e.g. StringView [21] and high-performance compute kernels)



Arrow Array: Int64Array

| |
|------|
| 15 |
| 1743 |
| NULL |
| 432 |
| .. |
| 9 |
| 322 |
| 8 |

Logically 1024
8 byte integers



Arrow Array



Pretty much what
you will find in
every vectorized
column store
engine

Why Arrow Internally (and not just at interface)?

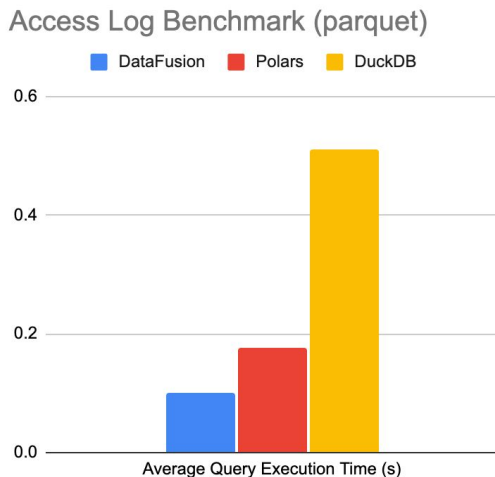
Theory: Using Arrow is “good enough” compared to specialized structures

Pooled open source development → invest heavily in optimized parquet reader

So far results
are encouraging

Good: Sorting, Filtering,
Projection, Parquet

Could Improve:
Grouping, Joining



<https://github.com/tustvoid/access-log-bench>

Foundations: Rust and its Ecosystem

Rust [76]: new system programming language (LLVM based)

1. **Excellent Performance:** similar to C/C++
2. **Easy to Embed:** no language run-time and have C ABI compatibility.
3. **Resource Efficient:** Low level (but safe!) memory management
4. **Productive Ecosystem:** Cargo Package Manager[12] and crate ecosystem make adding DataFusion to most projects as simple as adding a single line to a configuration file.

Physical Segment Trees

What we did instead: Classic sort based optimizations (sorted input by PARTITION BY / ORDER BY clause)

Why: Simple, well understood

Lessons Learned: