



# Meetup 1.0

March 2024  
Austin, TX

Andrew Lamb, Staff Engineer, InfluxData



---

# Welcoming Remarks Andrew Lamb

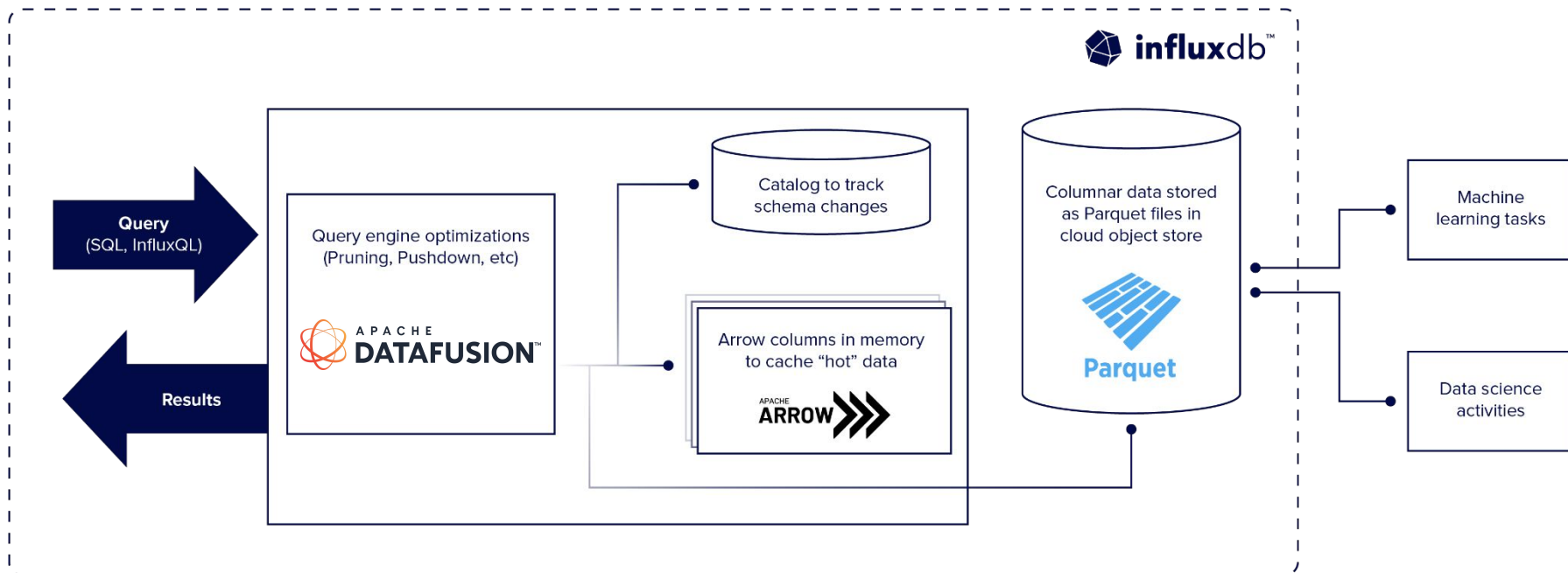
Thank you to our sponsor:  **influxdata**<sup>®</sup>

**Creators of InfluxDB** – the purpose built time series database

**InfluxDB 3.0** was completely rewritten using Apache Arrow, DataFusion, Parquet, and Rust

**Continuous open source commitment** – Long term believer and supporter of open source

# InfluxDB 3.0 Architecture



# Agenda

- Reflections
- Goals
- Logistics



---

**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: IOx, Arrow, DataFusion (4 years)

# Expectations

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.”

# Reality

## 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  
The Randstad, Netherlands  
dheres@apache.org

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

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

Chao Sun, Liang-Chi Hsieh  
Apple  
[Cupertino, Seattle], (CA, WA), USA  
[sunchao, viirya]@apache.org

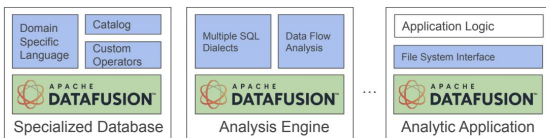


Figure 1: When building with DataFusion, system designers implement domain-specific features via extension APIs (blue), rather than re-implementing standard OLAP query engine technology (green).

### ABSTRACT

Apache Arrow DataFusion[28] is a fast, embeddable, and extensible query engine written in Rust[76] that uses Apache Arrow[27] 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 the features of a modern OLAP engine, and outline optimizations required for high performance. Next we describe DataFusion's architecture and extension APIs to illustrate the interfaces used in modular query engines to integrate with the systems built on 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[21, 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, and Chao Sun, Liang-Chi Hsieh. 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 9–11, 2024, Santiago, Chile. ACM, New York, NY, USA, 13 pages. <https://doi.org/10.1145/XXXXXX.XXXXXX>

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, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from [permissions.acm.org/](https://permissions.acm.org/).  
SIGMOD-Companion '24, June 9–11, 2024, Santiago, Chile.  
© 2024 Copyright held by the owner/authors. Publication rights licensed to ACM.  
ACM ISBN 979-8-4007-9421-1/24/06...\$15.00  
<https://doi.org/10.1145/XXXXXX.XXXXXX>

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

## To Appear: SIGMOD 2024

<https://github.com/apache/arrow-datafusion/issues/8373>



# Mind shift

---

## DataFusion is not

A product that we\* are providing

## DataFusion is

A project we are all working on together

\* we == Apache / the committers / myself

# Comparison with Enterprise Software Development

---

## Differences


- Don't need to fix all bugs (instead encourage others)
- Open contribution (not just [open source](#)): accept from everyone
- Less 🌶️🌶️🌶️ discussions
- Roadmap driven by who contributes

## Similarities

- Acknowledge humanity: thank for contributions (even if they don't get merged)
- Provide clear, explicit feedback on how to make progress
- Need to keep the code moving (Long PR queue still makes me anxious)
- Clear tickets get worked on

# Unreasonable 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 circled in yellow. A third comment from "alamb" says "Thank you @Lordworms". A large yellow text overlay reads "16 Minutes between file and claim", with arrows pointing 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

# Capitalism + Altruism

- I like getting paid
- InfluxData is a for profit company and responsible to investors
- DataFusion directly helps our bottom line
- ⇒ **Also** a greater good: drive up system wide efficiency
- More people can build / use better, faster, cheaper products



Non profit governance of open source communities

“Community over Code” - [The Apache Way](#)

# Apache: Benefits for DataFusion

- ⇒ Predictable Foundation
- **Stable License:** (ASL 20 years old) low risk of changes, (ahem OpenTofu)
- **Communication:** Predictable and open (if slow)
- **Multi-Vendor Participation:** Shared investment reduces individual risk
- **Long Term Maintenance:** Hedged against life changes, corporate strategy shifts, VC funding cycles
- ★★★★★: Works far better than could be reasonably expected

# New Top Level Project

- [\[DISCUSS\] Move Apache Arrow Datafusion to top level project](#)
- ⇒ Proposal Complete, Approved by Arrow Community
- ⇒ Expect final approval at April 2024 ASF board meeting
- Help: <https://github.com/apache/arrow-datafusion/issues/9691>
  - (see what I did there?)

# Conclusion

- Never imagined building such a cool project as open source
- Grateful to have a chance to work with you all
  - (and geek out about / obsess over performance)
- I hope 1,000 projects use DataFusion 🚀
- Thank you all 🙏



---

# Logistics

# Meetup Goals

1. Community building
2. Discuss common goals and challenges
3. Brainstorm ways to continue to grow and nurture the community.

**IMPORTANT:** much larger community than those in this room – please help make sure everything we say gets written down / communicated.

# Schedule

12:00 - 1:00: Lunch/Networking (food provided)

1:00 - 1:30: Welcome / Introductions: Andrew Lamb 

1:30 - 3:00: 6 x 10-15 minute talks about “How do you DataFusion”?

3:00 - 3:15: Coffee / Email Break

3:15 - 4:45: Breakouts / discussions

4:45 - 5:00: Close / Conclusion

5:00 - 7:00: DataFusion Social @ Rustic Tap

7:00 - 9:00: “Preparty” @ Speakeasy (part of datacouncil - [RSVP](#))

# DataFusion Social

**The Rustic Tap: 613 W 6th St, Austin, TX 78701**

14 min walk from this meetup

11 min from the Data Councils pre-party

They [serve food](#) and there is live music ⇒ Austin vibe

Thanks to Coordinators 🙌: Mattia Pavoni, Jacopo Tagliabue, Cheng She

# Introductions

Suggestions:

1. Name, Github handle
2. Optional affiliation
3. Optional reason you came to the meetup

---

How do you DataFusion?

# Schedule

12:00 - 1:00: Lunch/Networking (food provided)

1:00 - 1:30: Welcome / Introductions: Andrew Lamb

1:30 - 3:00: 6 x 10-15 minute talks about “How do you DataFusion”?

3:00 - 3:15: Coffee / Email Break

3:15 - 4:45: Breakouts / discussions

4:45 - 5:00: Close / Conclusion

5:00 - 7:00: DataFusion Social @ Rustic Tap

7:00 - 9:00: “Preparty” @ Speakeasy (part of datacouncil - [RSVP](#))

# Speaker Lineup

- Devin D'Angelo: [datafusion-meetup-devind.pptx](#)
- Lukas Schulte & Bo Lin: [SDF\\_Datafusion\\_DC2024Meetup.pdf](#)
- Bruce Ritchie: [Replacing Spark with DataFusion.pptx](#)
- Dan Harris:
- Andy Grove: [Andy Grove @ DataFusion Meetup](#)
- Jackson Newhouse (Arroyo): [Arroyo And DataFusion](#)
- QP Hou (Neuralink):



---

Begin 

---

Break 😄 (15 minutes)

# Schedule

12:00 - 1:00: Lunch/Networking (food provided)

1:00 - 1:30: Welcome / Introductions: Andrew Lamb

1:30 - 3:00: 6 x 10-15 minute talks about “How do you DataFusion”?

3:00 - 3:15: Coffee / Email Break

3:15 - 4:45: Breakouts / discussions

4:45 - 5:00: Close / Conclusion

5:00 - 7:00: DataFusion Social @ Rustic Tap

7:00 - 9:00: “Preparty” @ Speakeasy (part of datacouncil - [RSVP](#))

---

# Breakout Discussions

# Schedule

12:00 - 1:00: Lunch/Networking (food provided)

1:00 - 1:30: Welcome / Introductions: Andrew Lamb

1:30 - 3:00: 6 x 10-15 minute talks about “How do you DataFusion”?

3:00 - 3:15: Coffee / Email Break

3:15 - 4:45: Breakouts / discussions

4:45 - 5:00: Close / Conclusion

5:00 - 7:00: DataFusion Social @ Rustic Tap

7:00 - 9:00: “Preparty” @ Speakeasy (part of datacouncil - RSVP)

# Breakouts / Discussions

- Break into small groups for 20ish minute sessions
- I will ring a gong every 20 minutes  $\Rightarrow$  find a new group
- 🙏 🙏 🙏 write up insights, and post them as tickets / discussions
- We will regroup around 4:45

# Possible Discussion Topics



- Harnessing the community (e.g. good-first-issue tickets). Parable of [6195](#).
- Symbiotic relationship: maintainers organize work and community get it done?
- How do we get better automated testing? sql fuzz tests (e.g. different output target batch sizes, different partition counts ), automated performance regression tests
- Can we get community contributions? Do we need to pay people for this?
- How do we make it easier for people to build systems on DataFusion?
- How do we make it easier for people to contribute back to the project?
- What features would you like to see / not see?

---

# Closing



# Schedule

12:00 - 1:00: Lunch/Networking (food provided)

1:00 - 1:30: Welcome / Introductions: Andrew Lamb

1:30 - 3:00: 6 x 10-15 minute talks about “How do you DataFusion”?

3:00 - 3:15: Coffee / Email Break

3:15 - 4:45: Breakouts / discussions

4:45 - 5:00: Close / Conclusion

5:00 - 7:00: DataFusion Social @ Rustic Tap

7:00 - 9:00: “Preparty” @ Speakeasy (part of datacouncil - [RSVP](#))

# DataFusion Historical Trivia Quiz

**Rules:** 30-60s per Question, No internet, talk with neighbors.

Q: Original author of DataFusion?

**A: Andy Grove**

Q: Book is DataFusion based on?

**A: [How Query Engines Work](#)**

Q: # github stars?

**A: 4.8k**

Q: # distinct contributors?

**A: [538](#) / 350\***

Q: First version in ClickBench?

**A: [10.0.0](#)**

Q: When did PR/Issue 10,000 happen?

**A: in 1 months (est) #9786**

# Onwards!

1,000+ projects!

# Thank you again

- DataFusion is amazing and could not have been done without you
- Thank you for attending, hope you found face to face time valuable
- Treat DataFusion like it is your own project (because it is)
  - Review PRs, submit issues, etc

# Reminder

Help *write* down everything interesting you found here

- for the community not able to attend physically
- for ourselves in 4 months

I will posts slides / videos in the next day or two

# DataFusion Social

**The Rustic Tap: 613 W 6th St, Austin, TX 78701**

14 min walk from this meetup

11 min from the Data Councils pre-party

They [serve food](#) and there is live music ⇒ Austin vibe

Thanks to Coordinators 🙌: Mattia Pavoni, Jacopo Tagliabue, Cheng She



---

THANK YOU