# LAST CLASS

Snowflake Data Warehouse

# HISTORICAL CONTEXT

CWI researchers recognized that data scientists do not use the full query capabilities of DBMSs due to the overhead of setting up and accessing data.

In 2017 they created an embedded version of MonetDB called <u>MonetDBLite</u> to run inside of R applications.
→ Running in-process reduces the cost of transferring data back and forth between the DBMS and the application.

But MonetDB had too much legacy baggage…

# DUCKDB (2019)

Multi-threaded embedded (in-process, serverless) DBMS that executes SQL over disparate data files.
→ PostgreSQL-like dialect with quality-of-life enhancements.
→ ***"SQLite for Analytics"***

Provides zero-copy access to query results via Arrow to client code running in same process.

The core DBMS is nearly all custom C++ code with little to no third-party dependencies.
→ Relies on extensions ecosystem to expand capabilities.

DUCKDB: AN EMBEDDABLE
ANALYTICAL DATABASE
SIGMOD 2019

# DUCKDB

Shared-Everything

Push-based Vectorized Query Processing

Precompiled Primitives

Multi-Version Concurrency Control

Morsel Parallelism + Scheduling

PAX Columnar Storage

Sort-Merge + Hash Joins

Stratified Query Optimizer
→ Supports unnesting of arbitrary subqueries

# DUCKDB: PUSH-BASED PROCESSING

System originally used pull-based vectorized query processing but found it unwieldly to expand to support more complex parallelism.
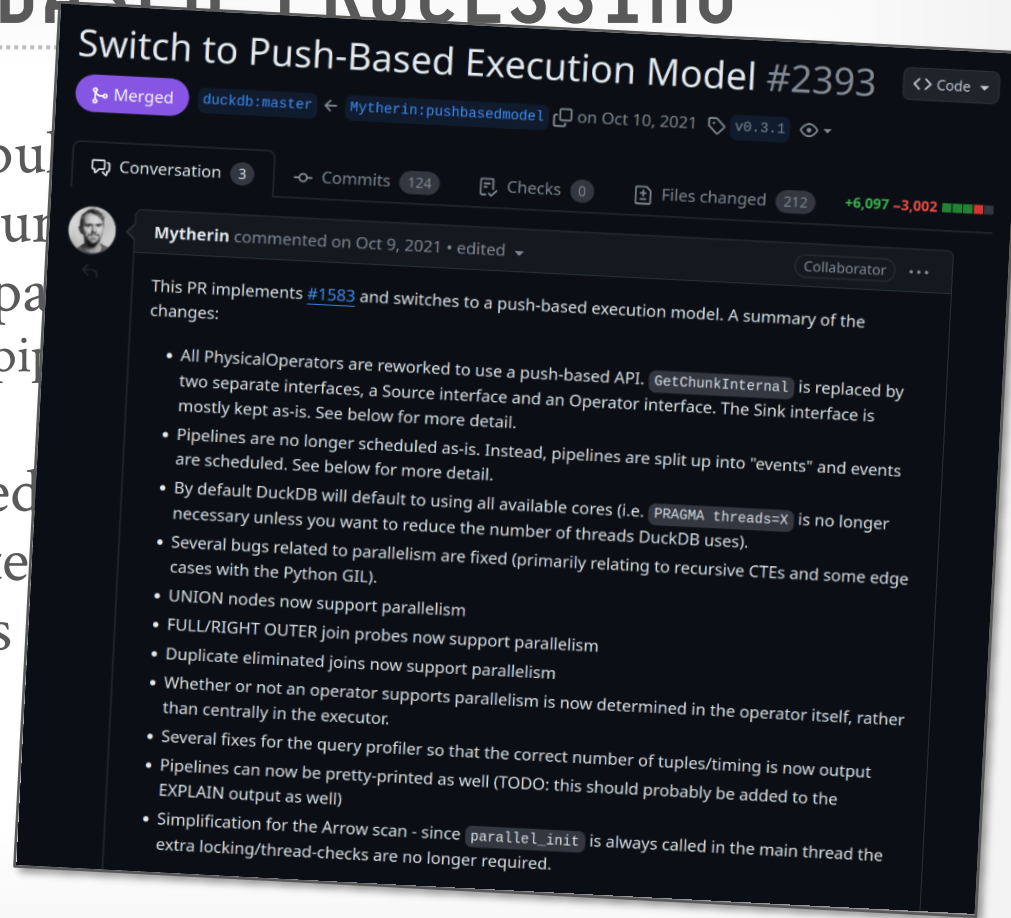→ Cannot invoke multiple pipelines simultaneously.

Switched to a push-based query processing model in 2021. Each operator determines whether it will execute in parallel on its own instead of a centralized executor.

# DUCKDB: PUSH-BASED PROCESSING

System originally used pu[...]
processing but found it un[...]
support more complex pa[...]
→ Cannot invoke multiple pi[...]

Switched to a push-based [...]
2021. Each operator dete[...]
execute in parallel on its [...]
centralized executor.



Switch to Push-Based Execution Model #2393

**Merged**   duckdb:master ← Mytherin:pushbasedmodel   on Oct 10, 2021   v0.3.1

💬 Conversation 3    Commits 124    Checks 0    Files changed 212    +6,097 −3,002

**Mytherin** commented on Oct 9, 2021 • edited ▾     Collaborator

This PR implements #1583 and switches to a push-based execution model. A summary of the changes:

- All PhysicalOperators are reworked to use a push-based API. `GetChunkInternal` is replaced by two separate interfaces, a Source interface and an Operator interface. The Sink interface is mostly kept as-is. See below for more detail.
- Pipelines are no longer scheduled as-is. Instead, pipelines are split up into "events" and events are scheduled. See below for more detail.
- By default DuckDB will default to using all available cores (i.e. `PRAGMA threads=X` is no longer necessary unless you want to reduce the number of threads DuckDB uses).
- Several bugs related to parallelism are fixed (primarily relating to recursive CTEs and some edge cases with the Python GIL).
- UNION nodes now support parallelism
- FULL/RIGHT OUTER join probes now support parallelism
- Duplicate eliminated joins now support parallelism
- Whether or not an operator supports parallelism is now determined in the operator itself, rather than centrally in the executor.
- Several fixes for the query profiler so that the correct number of tuples/timing is now output
- Pipelines can now be pretty-printed as well (TODO: this should probably be added to the EXPLAIN output as well)
- Simplification for the Arrow scan - since `parallel_init` is always called in the main thread the extra locking/thread-checks are no longer required.

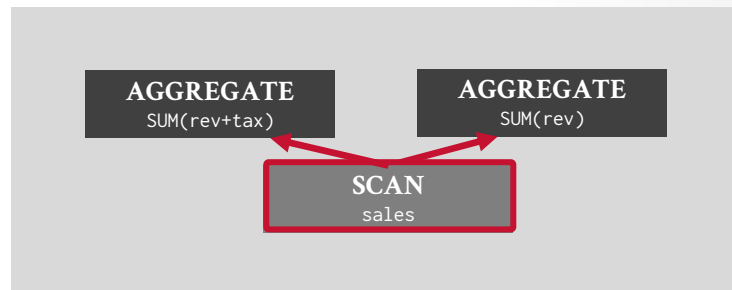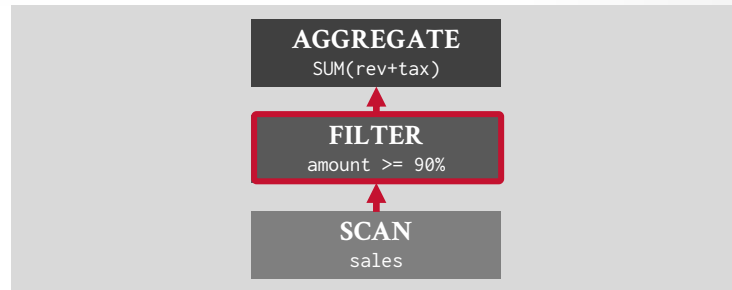# DUCKDB: FINE-GRAINED CONTROL

**Vector Cache:**
→ Buffer results between operators until it fills vector.

**Scan Sharing:**
→ Push results from one child operator to multiple parent operators (DAG plan).

**Backpressure / Async IO**
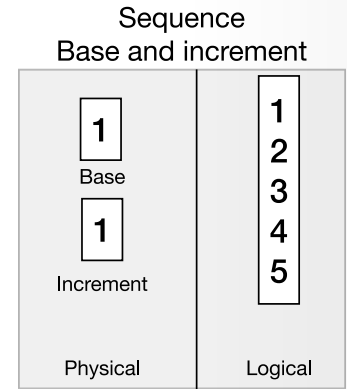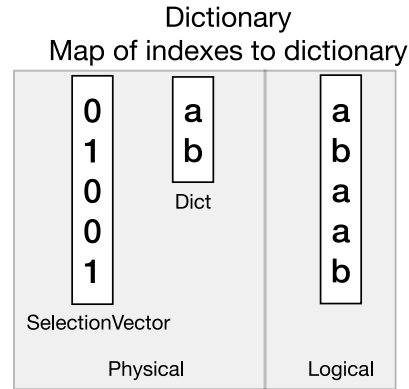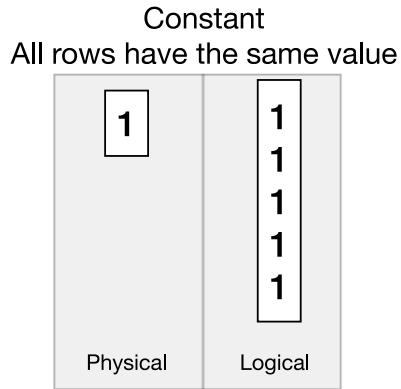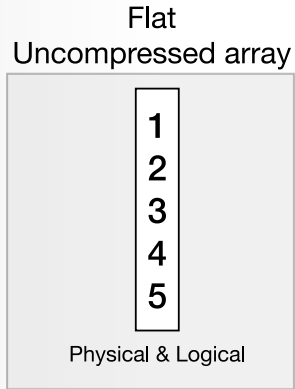→ Pause operator execution when buffers are full or when waiting for remote I/O.



Source: Mark Raasveldt

# DUCKDB: VECTORS

Custom internal vector layout for intermediate results that is compatible with Velox.

Supports multiple vector types:
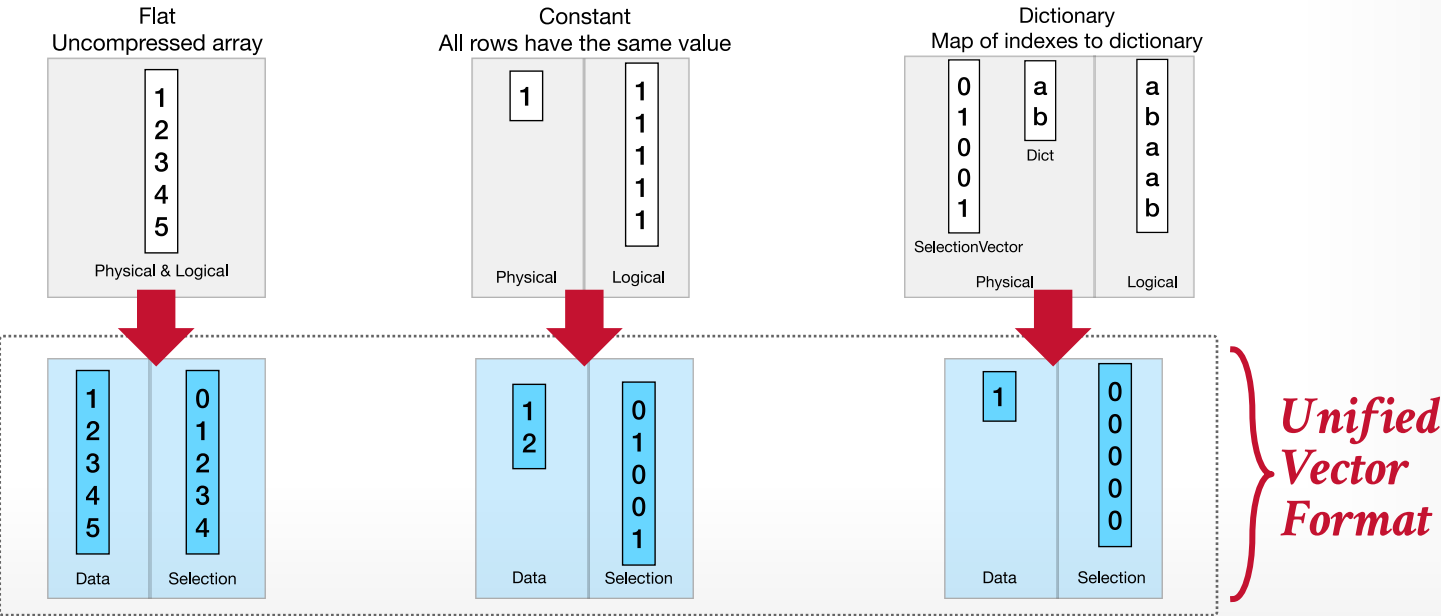


| Flat Uncompressed array | Constant All rows have the same value | Dictionary Map of indexes to dictionary | Sequence Base and increment |
|---|---|---|---|
| 1 2 3 4 5 | 1 / 1 1 1 1 1 | 0 1 0 0 1 / a b / a b a a b | 1 (Base) 1 (Increment) / 1 2 3 4 5 |
| Physical & Logical | Physical / Logical | SelectionVector / Dict / Physical / Logical | Physical / Logical |

# DUCKDB: VECTORS

DuckDB uses a unified format to process all vector types without needing to decompress them first.
→ Reduce # of specialized primitives per vector type

# DUCKDB: DATAFRAMES

DuckDB supports DataFrame libraries to query databases without using SQL.
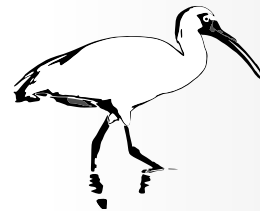→ **dpylr** (R-lang)
→ **Ibis** (Python)

Integration libraries generate DuckDB logical plans the DBMS converts into optimized physical plans.
→ Bypasses the SQL parser

Zero-copy result passing via Apache Arrow.

# DUCKDB: STORAGE FORMAT

DBMS's built-in storage format maintains a single PAX-oriented file per database.
→ Splits tables into row groups with 120k tuples.
→ On-disk encoding is different than in-memory representation.

Two phase compression scheme:
→ **Analyze**: Sample a small portion of a column to determine the best encoding scheme
→ **Compress:** Encode the values and write it to disk.

# DUCKDB: STORAGE FORMAT

DBMS's built-in s̶ format maintains a single
PAX-oriented file
→ Splits tables into r̶
→ On-disk encoding
representation.

Two phase comp
→ **Analyze**: Sample̶
the best encoding
→ **Compress:** Enco̶

| Version | Taxi | On-Time | Lineitem | Notes | Date |
|---|---|---|---|---|---|
| DuckDB v0.2.8 | 15.3GB | 1.73GB | 0.85GB | Uncompressed | July 2021 |
| DuckDB v0.2.9 | 11.2GB | 1.25GB | 0.79GB | RLE + Constant | September 2021 |
| DuckDB v0.3.2 | 10.8GB | 0.98GB | 0.56GB | Bitpacking | February 2022 |
| DuckDB v0.3.3 | 6.9GB | 0.23GB | 0.32GB | Dictionary | April 2022 |
| DuckDB v0.5.0 | 6.6GB | 0.21GB | 0.29GB | FOR | September 2022 |
| DuckDB v0.6.0 | 4.8GB | 0.21GB | 0.17GB | FSST + Chimp | October 2022 |
| CSV | 17.0GB | 1.11GB | 0.72GB | | |
| Parquet (Uncompressed) | 4.5GB | 0.12GB | 0.31GB | | |
| Parquet (Snappy) | 3.2GB | 0.11GB | 0.18GB | | |
| Parquet (ZSTD) | 2.6GB | 0.08GB | 0.15GB | | |

# DUCKDB: EXTERNAL TABLES

The DBMS can also access external data files via extensions.
→ Parquet, Arrow, SQLite, JSON,

Can also install extensions to retrieve files from remote filesystems (HTTP, S3)

# MOTHERDUCK

Cloud-based service that provides automatic execution of DuckDB queries on serverless compute nodes.
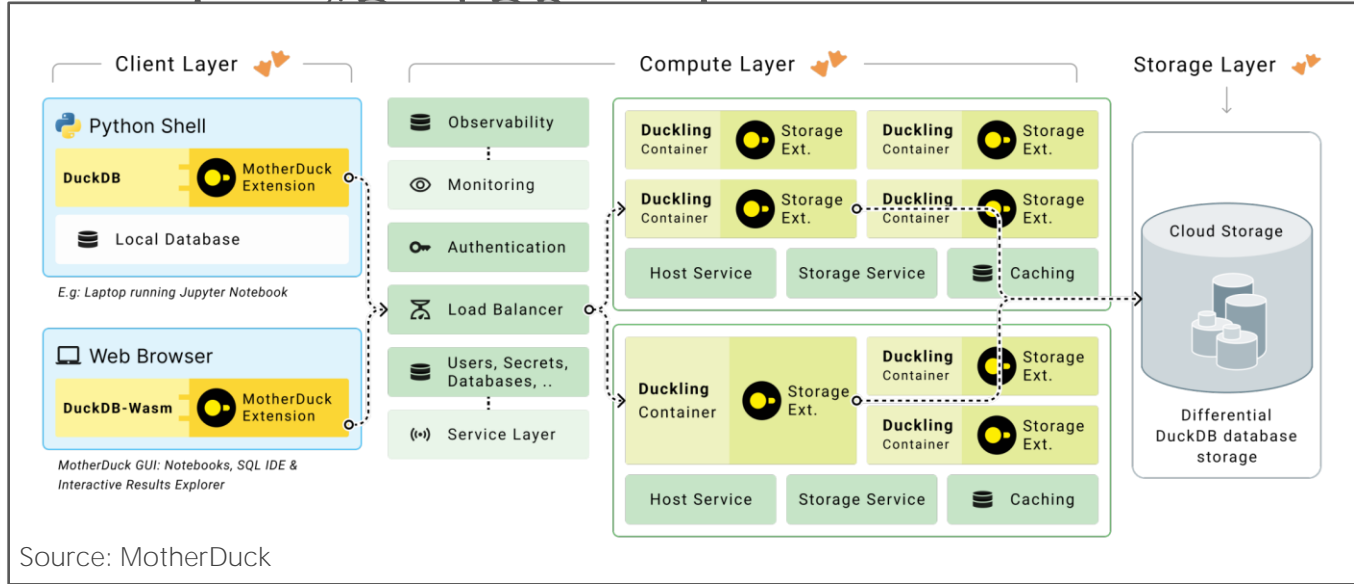
→ Remote nodes are DuckDB instances running inside of containers and connected to object stores.

→ Exposes remote catalog to local instance.

The latest versions of DuckDB already include extension to connect to MotherDuck.

# MOTHERDUCK

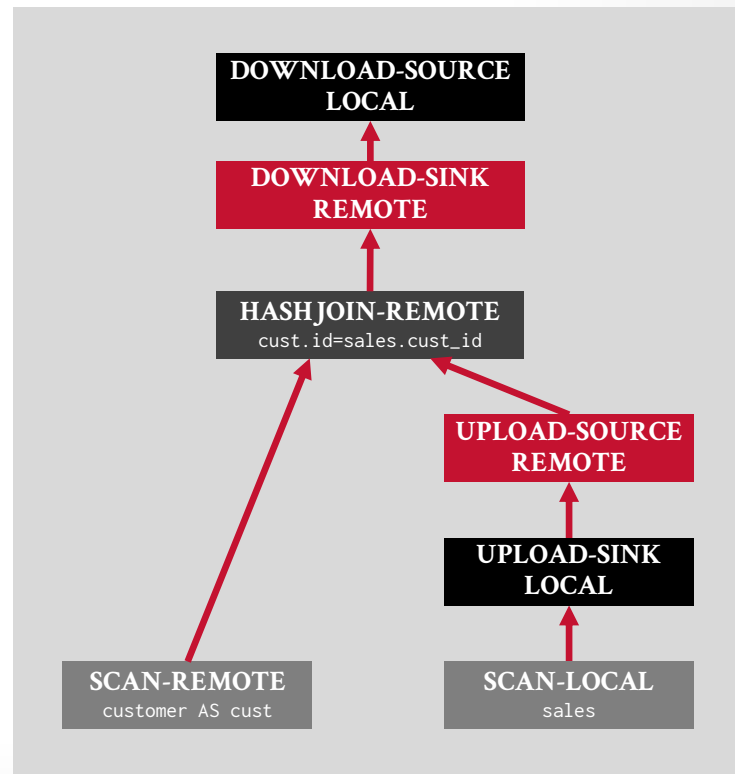Cloud-based service that provides automatic



Source: MotherDuck

# MOTHERDUCK: HYBRID QUERY PROCESSING

Introduces a new "bridge" operators that passes tuple streams between local and remote DuckDB instances.
→ Leverages operator pausing feature that DuckDB added from switching to push-based execution.

Query optimization occurs on the local instance as normal and then uses cost-based rules to decide what to run locally vs. remote.

# PARTING THOUGHTS

DuckDB is brilliant and its adoption is enviable.
→ Right place. Right time. Right problem.

Andy bet his earlier research agenda wrongly on in-memory DBMSs.

This is what HyPer/Umbra could have become if they were open-source…

# PARTING THOUGHTS

DuckDB is brilliant and its adoption is enviable.
→ Right place. Right time. Right problem.

Andy bet his earlier research agenda wrongly on in-memory DBMSs.

This is what HyPer/Umbra could have become if they were open-source…

# NEXT CLASS

Yellowbrick