

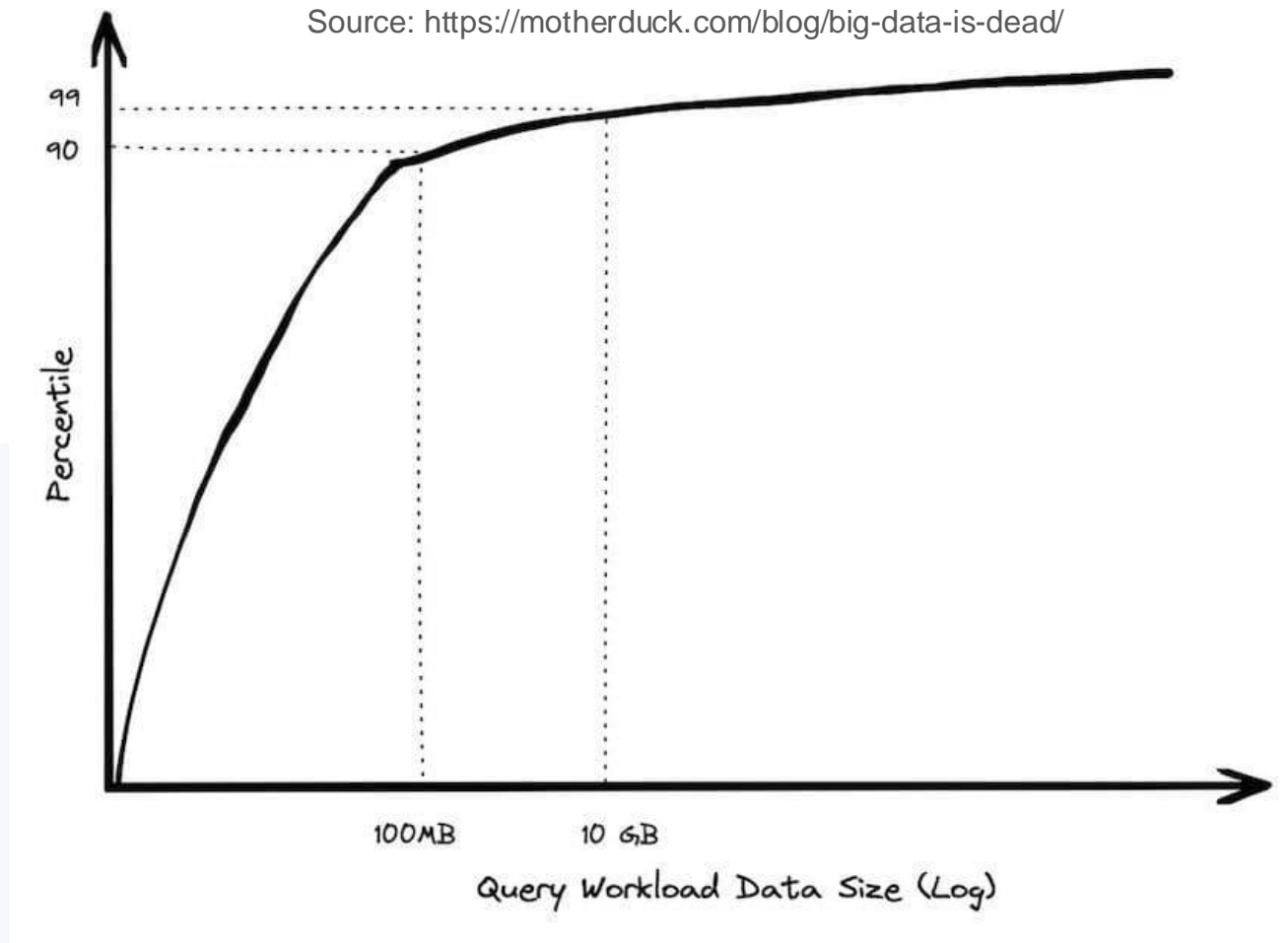


Ducklake

A journey to integrate DuckDB with Unity Catalog

Frank Mbonu, Diederik Greveling

Big data is dead?



in-process analytical database designed for fast query execution, especially suited for analytics workloads.

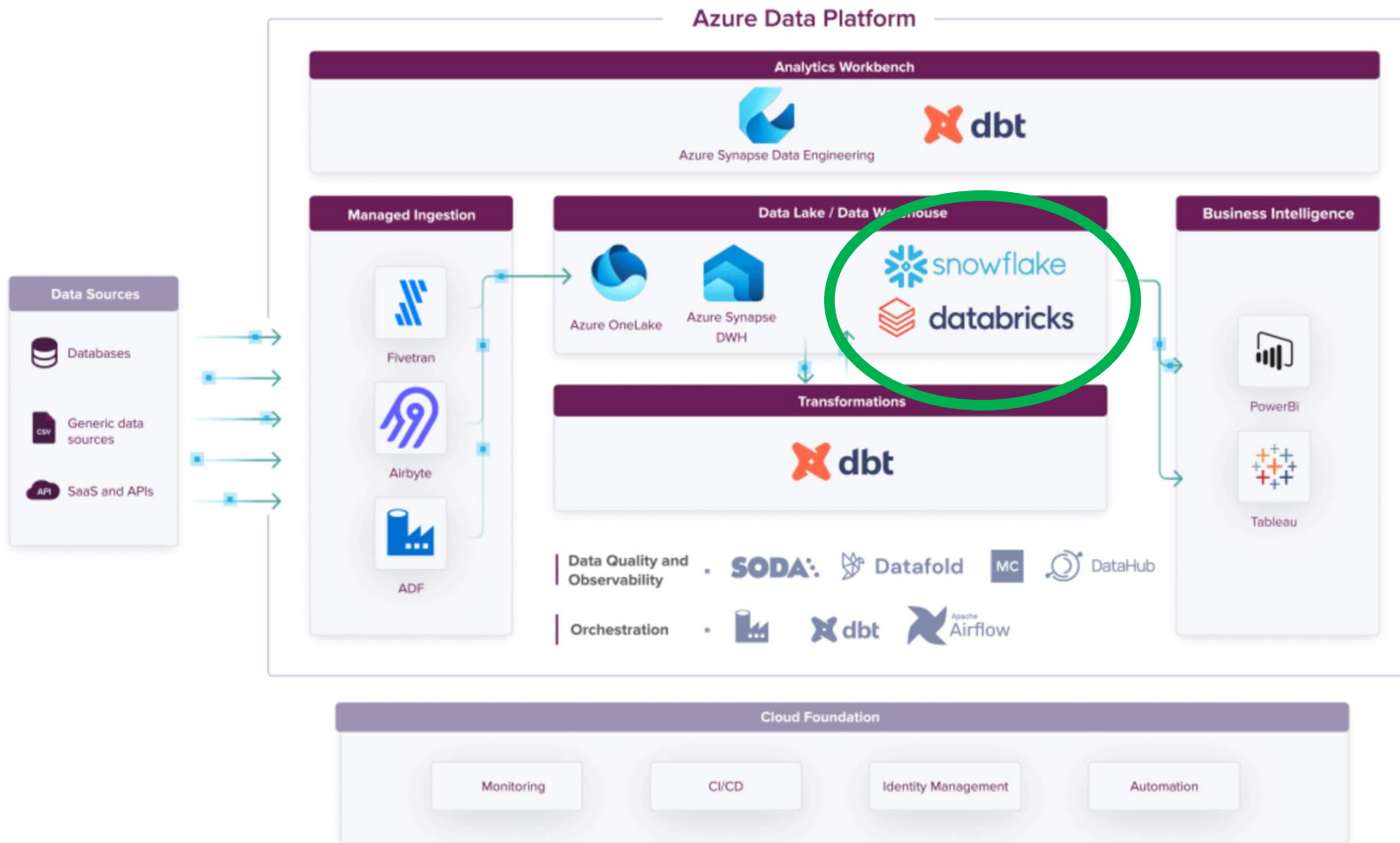
Big data is dead?

Source: <https://motherduck.com/blog/big-data-is-dead/>

Key takeaways

- Most Organizations Don't Have "Big Data"
 - "The majority of companies, even large enterprises, typically have data warehouses smaller than a terabyte"
- Workloads Are Smaller Than Data Sizes
 - "Most analytical queries process only a small fraction of the total data stored. For example, 90% of BigQuery queries process less than 100 MB of data"
- The Big Data Frontier is Shrinking
 - "As hardware improves, fewer workloads require distributed systems. A single machine today can handle what required thousands of nodes a decade ago"

For whom is this interesting?



Distributed compute (e.g. Spark) makes a lot of sense for large datasets.

But for smaller ones (< 1Tb) it might not be the best fit

Introducing Ducklake



DuckDB

in-process analytical database designed for fast query execution, especially suited for analytics workloads.

+



centralized governance, features like access controls and data lineage. Open source since the summer of 2024.

=

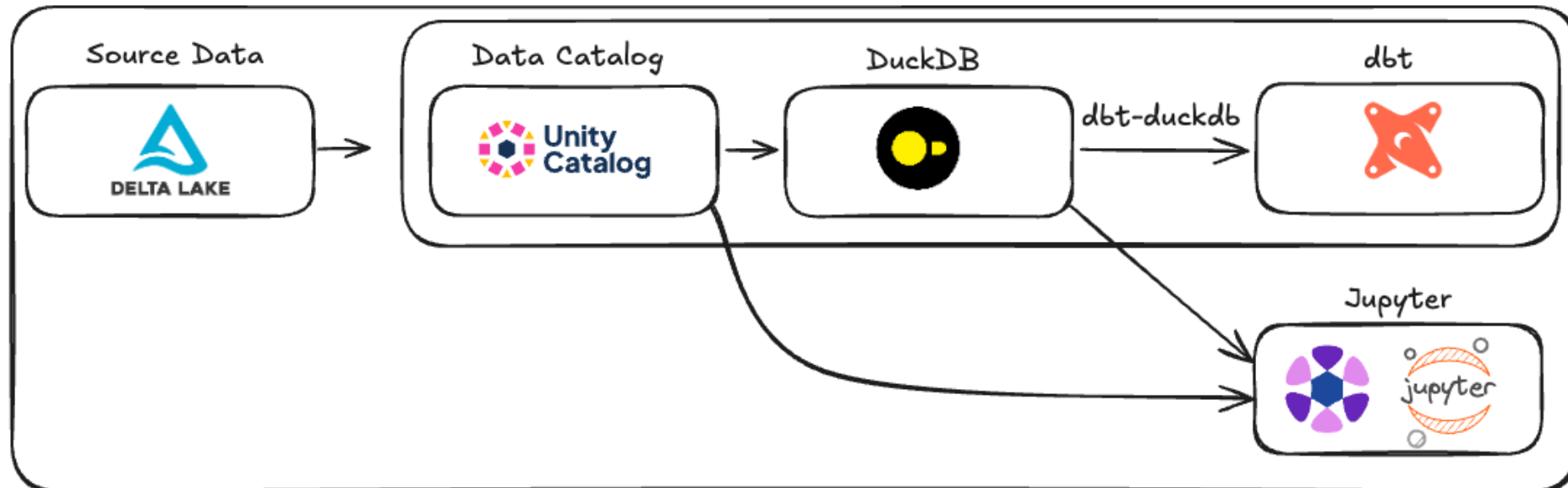


Ducklake: Let's combine DuckDB with the Unity Catalog through the DuckDB Unity Catalog extension

Ducklake high-level design

Requirements:

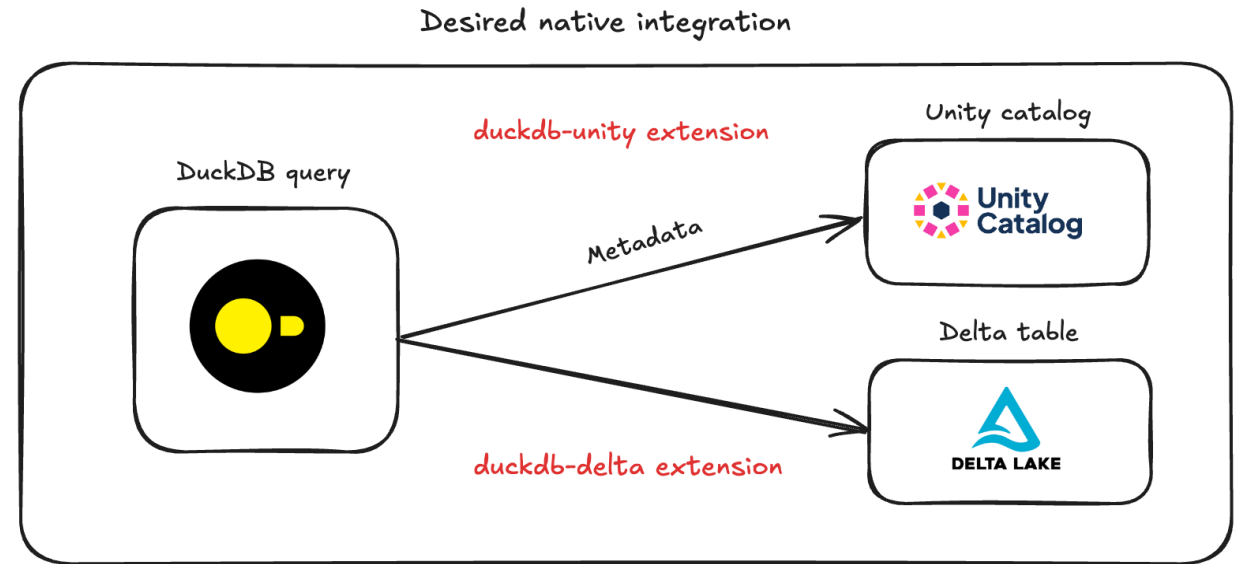
- Access Control
- Support for ACID transactions
 - To ensure data is complete, correct, conflict-free
- Read-write integration with an open data format
- Notebook interface for interactive development
- Storage decoupled from compute
- Full dbt integration



Challenges encountered

Lack of full native compatibility

- The DuckDB delta extension depends on the **delta-kernel-rs***, which currently only supports reads (and blind appends)
- The DuckDB unity extension currently does not support all CRUD operations (e.g. table create, schema create)



*The **Delta Kernel (Java or Rust (C and C++ bindings))** is a set of libraries that provides high-level APIs for interacting with Delta Lake tables designed to make Delta Lake integration easier and more efficient.

*DuckDB is written in C++ hence the delta-kernel-rs is used for building the delta connector

The workaround

Creating our own dbt-duckdb* plugin



Delta-rs

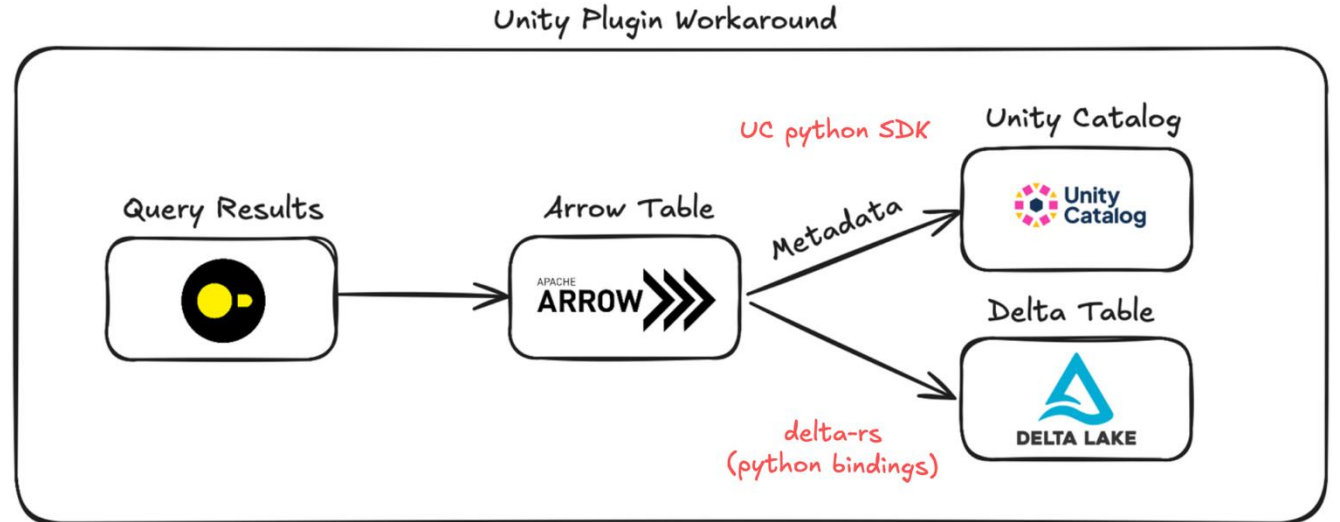
- Writing the result in-memory arrow table to a delta table



UC SDK

- Create schemas
- Create tables
- Get temporary storage credentials (AWS)

**dbt-duckdb is a dbt adapter for DuckDB*



A dbt build will

1. **Compile** our dbt models into executable SQL
2. **Execute** the compiled code against DuckDB
3. **Convert** the query results to a PyArrow table
4. **Create** the unity schema if it doesn't exist
5. **Create** the unity table if it doesn't exist
6. **Write** the PyArrow table to a Delta table

Ducklake in action

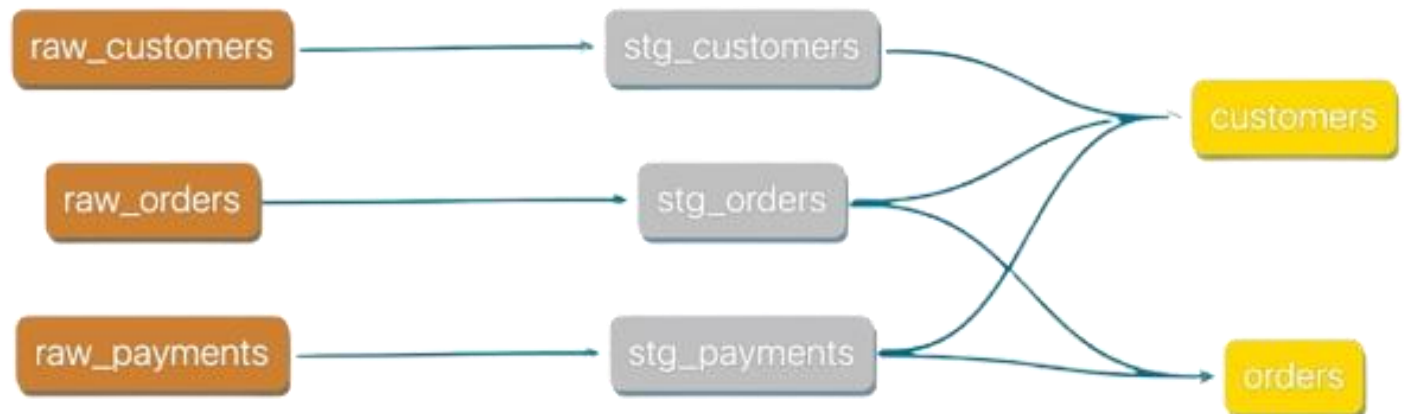
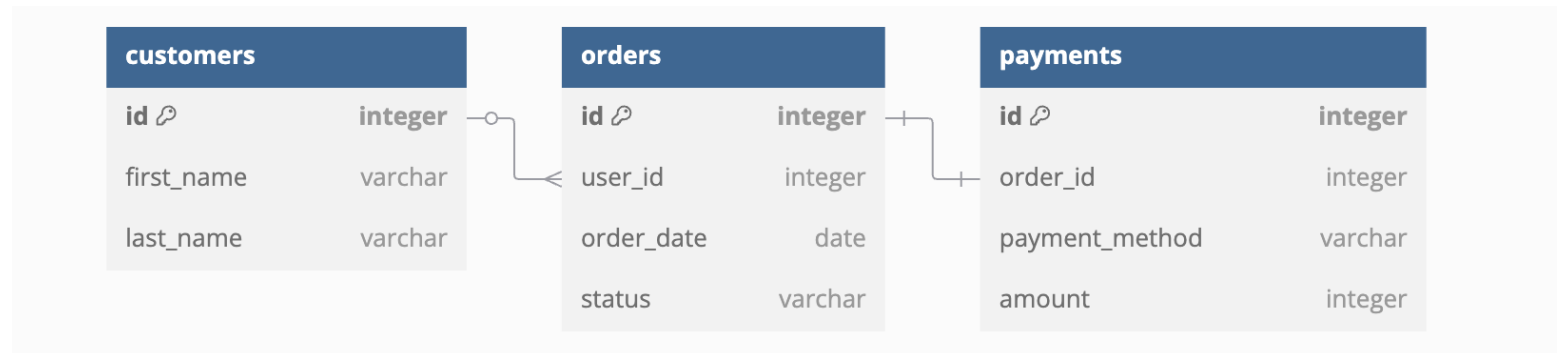
Now that we have created an integrated setup with DuckDB and Unity Catalog

Let's put our Ducklake to test...

First, we need some data to work with. Since we're using dbt in our setup, the jaffle shop seems appropriate to use.

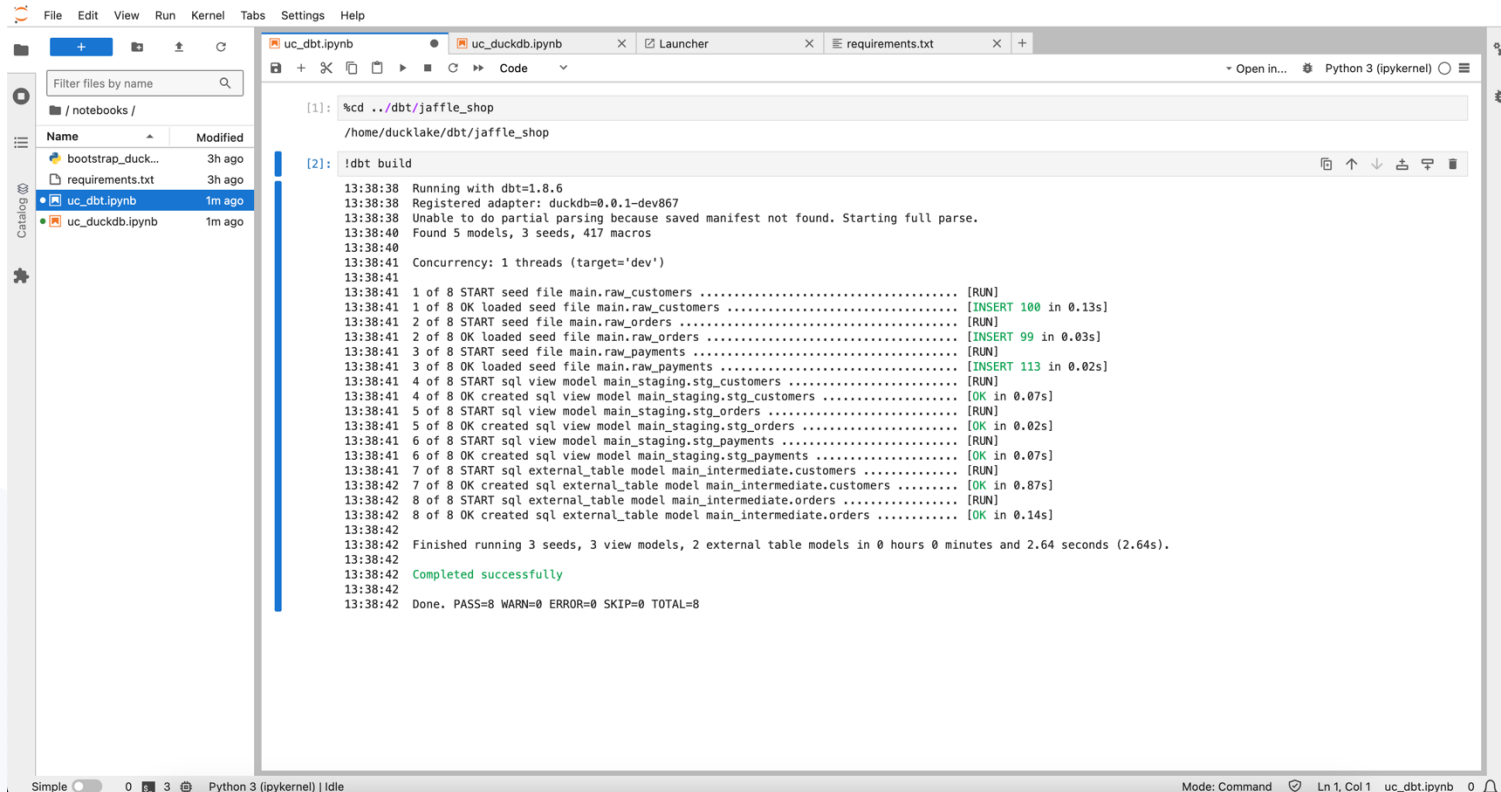
The Jaffle shop is a fictional e-commerce store often used for dbt demos. It transforms raw csv data into customer and order models.

In our example the gold models are **materialized as tables** in Unity Catalog



Ducklake in action

First, let's populate our unity catalog with the jaffle shop materialized tables using `dbt build`



```
[1]: %cd ../dbt/jaffle_shop
/home/ducklake/dbt/jaffle_shop

[2]: !dbt build

13:38:38 Running with dbt=1.8.6
13:38:38 Registered adapter: duckdb=0.0.1-dev867
13:38:38 Unable to do partial parsing because saved manifest not found. Starting full parse.
13:38:40 Found 5 models, 3 seeds, 417 macros
13:38:40
13:38:41 Concurrency: 1 threads (target='dev')
13:38:41
13:38:41 1 of 8 START seed file main.raw_customers ..... [RUN]
13:38:41 1 of 8 OK loaded seed file main.raw_customers ..... [INSERT 100 in 0.13s]
13:38:41 2 of 8 START seed file main.raw_orders ..... [RUN]
13:38:41 2 of 8 OK loaded seed file main.raw_orders ..... [INSERT 99 in 0.03s]
13:38:41 3 of 8 START seed file main.raw_payments ..... [RUN]
13:38:41 3 of 8 OK loaded seed file main.raw_payments ..... [INSERT 113 in 0.02s]
13:38:41 4 of 8 START sql view model main_staging.stg_customers ..... [RUN]
13:38:41 4 of 8 OK created sql view model main_staging.stg_customers ..... [OK in 0.07s]
13:38:41 5 of 8 START sql view model main_staging.stg_orders ..... [RUN]
13:38:41 5 of 8 OK created sql view model main_staging.stg_orders ..... [OK in 0.02s]
13:38:41 6 of 8 START sql view model main_staging.stg_payments ..... [RUN]
13:38:41 6 of 8 OK created sql view model main_staging.stg_payments ..... [OK in 0.07s]
13:38:41 7 of 8 START sql external_table model main_intermediate.customers ..... [RUN]
13:38:42 7 of 8 OK created sql external_table model main_intermediate.customers ..... [OK in 0.07s]
13:38:42 8 of 8 START sql external_table model main_intermediate.orders ..... [RUN]
13:38:42 8 of 8 OK created sql external_table model main_intermediate.orders ..... [OK in 0.14s]
13:38:42
13:38:42 Finished running 3 seeds, 3 view models, 2 external table models in 0 hours 0 minutes and 2.64 seconds (2.64s).
13:38:42
13:38:42 Completed successfully
13:38:42
13:38:42 Done. PASS=8 WARN=0 ERROR=0 SKIP=0 TOTAL=8
```

This will

1. **Compile** our dbt models into executable SQL
2. **Execute** the compiled code against DuckDB
3. **Convert** the query results to a PyArrow table
4. **Create** the unity schema if it doesn't exist
5. **Create** the unity table if it doesn't exist
6. **Write** the PyArrow table to a Delta table

Ducklake in action

Now that our Unity Catalog is populated, let's try to query our materialized tables

[3]: `ATTACH 'unity' AS unity (TYPE UC_CATALOG);` ← ATTACH our Unity Catalog to DuckDB so that we can read/write from/to UC tables
Database 'unity' attached successfully.

[4]: `SHOW ALL TABLES;`







database	schema	name	column_names	column_types	temporary
unity	default	marksheet	['id' 'name' 'marks']	['INTEGER' 'VARCHAR' 'INTEGER']	False
unity	default	marksheet_uniform	['id' 'name' 'marks']	['INTEGER' 'VARCHAR' 'INTEGER']	False
unity	default	numbers	['as_int' 'as_double']	['INTEGER' 'DOUBLE']	False
unity	intermediate	customers	['customer_id' 'first_name' 'last_name' 'first_order' 'most_recent_order' 'number_of_orders' 'customer_lifetime_value']	['INTEGER' 'VARCHAR' 'VARCHAR' 'DATE' 'DATE' 'BIGINT' 'DOUBLE']	False
unity	intermediate	orders	['order_id' 'customer_id' 'order_date' 'status' 'credit_card_amount' 'coupon_amount' 'bank_transfer_amount' 'gift_card_amount' 'amount']	['INTEGER' 'INTEGER' 'DATE' 'VARCHAR' 'DOUBLE' 'DOUBLE' 'DOUBLE' 'DOUBLE' 'DOUBLE']	False

[5]: `SELECT * FROM unity.intermediate.customers LIMIT 10;`

customer_id	first_name	last_name	first_order	most_recent_order	number_of_orders	customer_lifetime_value
1	Michael	P.	2018-01-01 00:00:00	2018-02-10 00:00:00	2	33
2	Shawn	M.	2018-01-11 00:00:00	2018-01-11 00:00:00	1	23
3	Kathleen	P.	2018-01-02 00:00:00	2018-03-11 00:00:00	3	65
6	Sarah	R.	2018-02-19 00:00:00	2018-02-19 00:00:00	1	8
7	Martin	M.	2018-01-14 00:00:00	2018-01-14 00:00:00	1	26
8	Frank	R.	2018-01-29 00:00:00	2018-03-12 00:00:00	2	45
9	Jennifer	F.	2018-03-17 00:00:00	2018-03-17 00:00:00	1	30
11	Fred	S.	2018-03-23 00:00:00	2018-03-23 00:00:00	1	3
12	Amy	D.	2018-03-03 00:00:00	2018-03-03 00:00:00	1	4
13	Kathleen	M.	2018-03-07 00:00:00	2018-03-07 00:00:00	1	26

DuckLake works—no quacks, just results 

Next steps

-  Implement RBAC (UC v0.2.0)
-  Support schema evolution
-  Deploy Ducklake in the cloud (AWS)
-  Improve UI
-  Make Ducklake pluggable
 - Choose your engine
 - Choose your catalog
-  Write more blogs about our journey

Read our blog post if you haven't already!



Xebia



xebia.com