

Revisiting Analytical Storage with DuckLake

A simpler take on the lakehouse

May 7, 2026

Table of contents

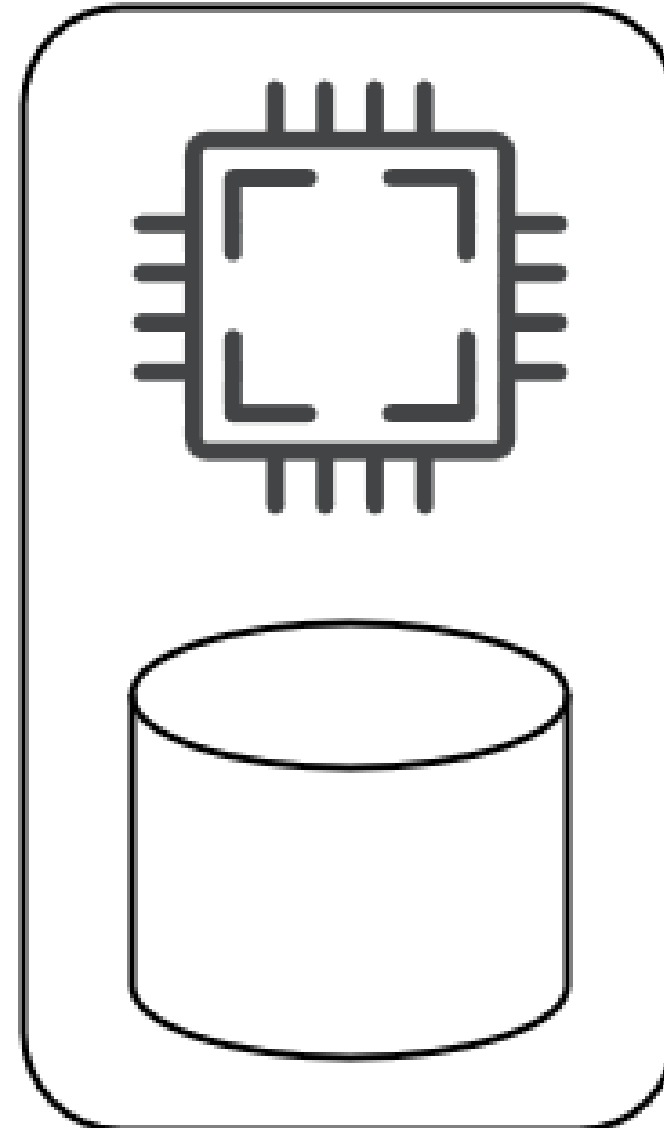
- **The historical evolution of analytical storage**
The evolution from database → warehouse → lake → lakehouse
- **DuckLake: Rethinking the Lakehouse Architecture**
Metadata as a database problem
- **DuckLake in practice**
Ecosystem integration and the early-adopter trade-off
- **Conclusion**

The historical evolution of analytical storage

From database → warehouse → lake → lakehouse

The transactional database (OLTP)

**Optimized for transactions,
not analytics**

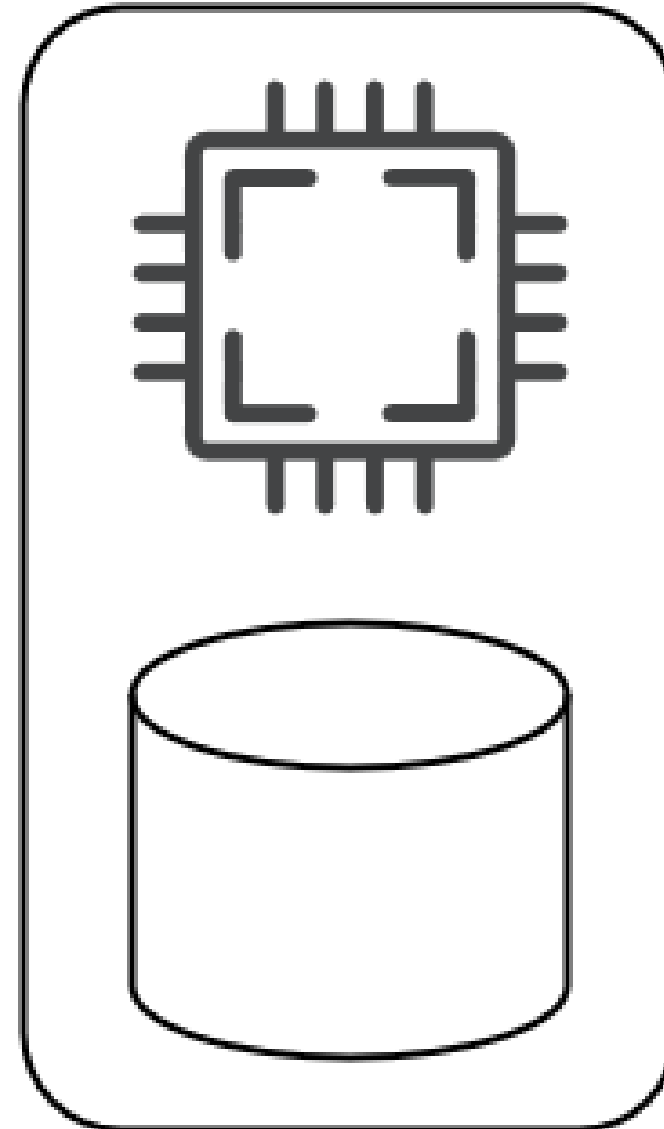


Transaction processing vs analytics

Property	Transaction processing systems (OLTP)	Analytic systems (OLAP)
Main read pattern	Small number of records per query, fetched by key	Aggregate over large number of records
Main write pattern	Random-access, low-latency writes from user input	Bulk import (ETL) or event stream
Primarily used by	End user/customer, via web application	Internal analyst, for decision support
What data represents	Latest state of data (current point in time)	History of events that happened over time
Dataset size	Gigabytes to terabytes	Terabytes to petabytes

The data warehouse (OLAP)

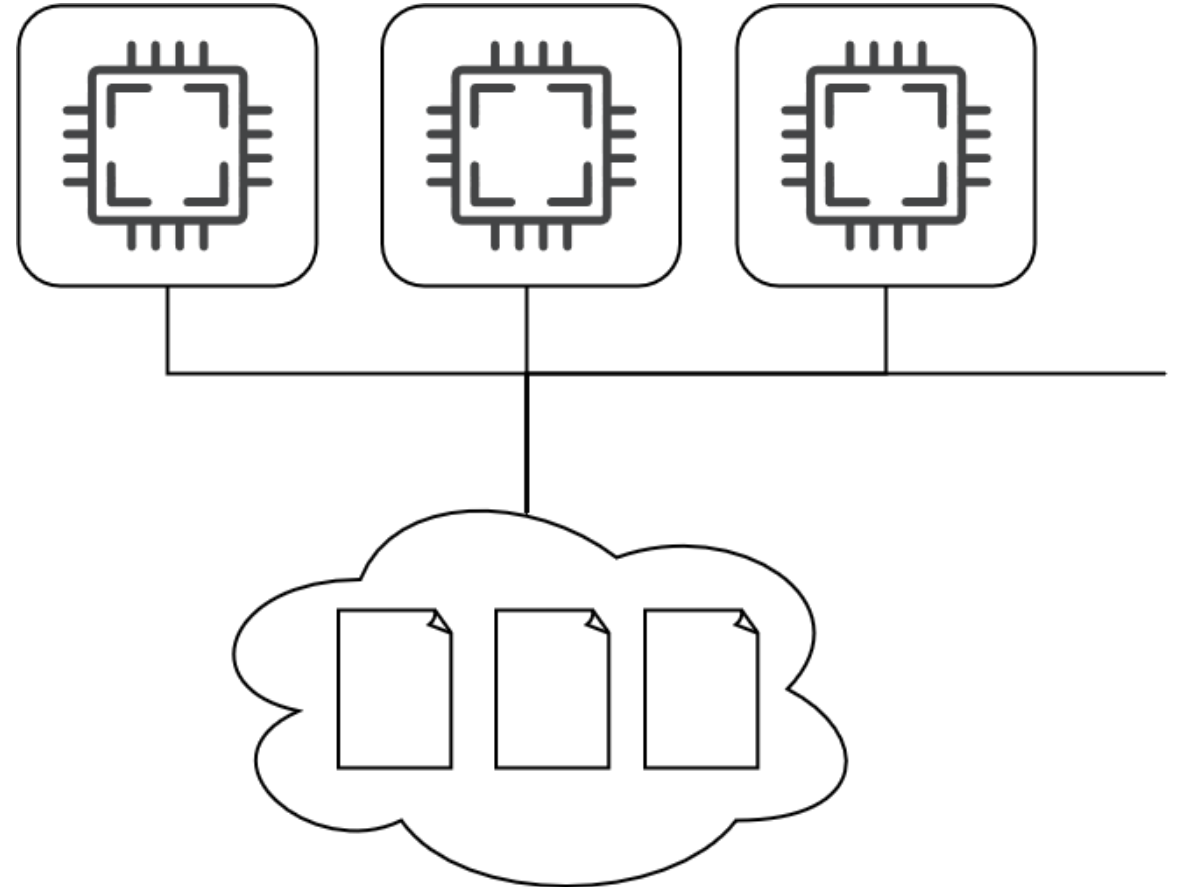
Storage & compute combined



The data lake

**Weak governance => can become a
“data swamp”**

No transactional guarantees

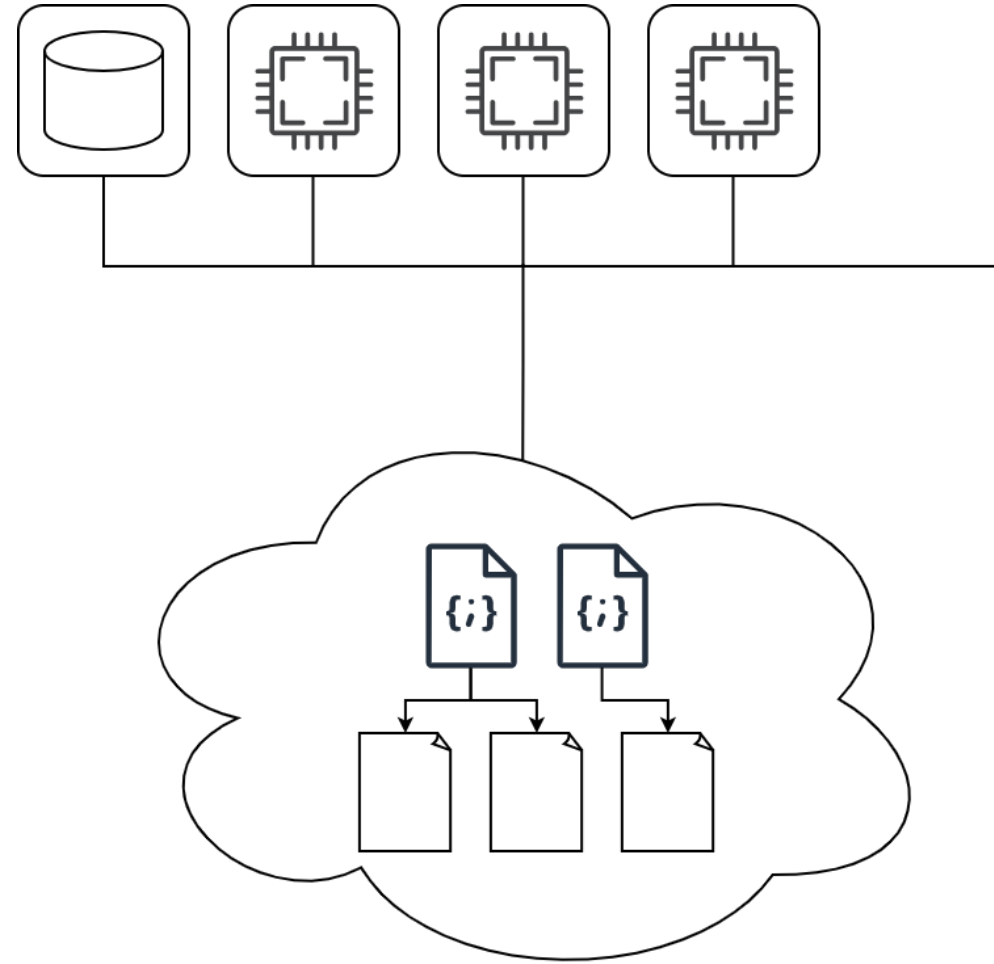


The data lakehouse

Example:

- Iceberg (+ Polaris)
- Delta Lake (+ Unity)
- DuckLake

**Bring order to the datalake
using metadata**

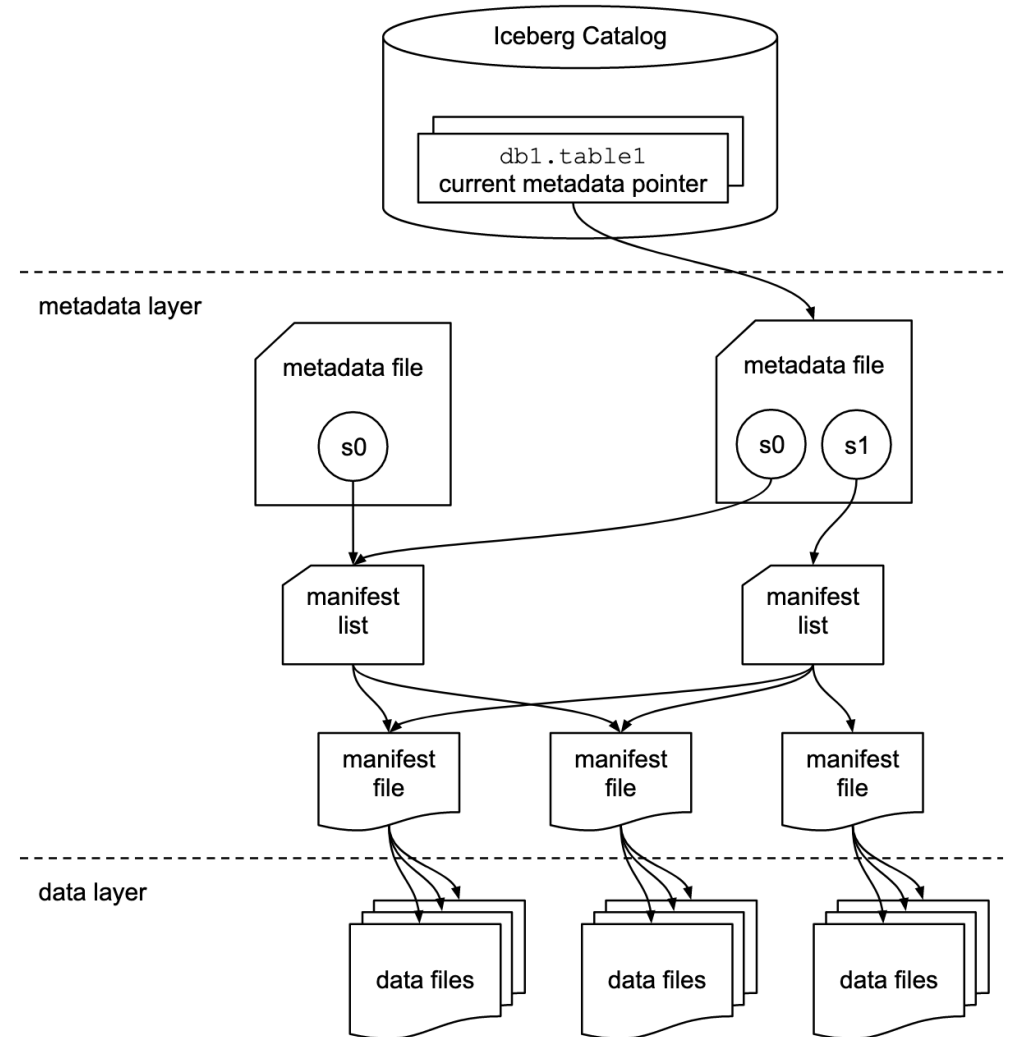


DuckLake: Rethinking the Lakehouse Architecture

Metadata as a database problem

The Iceberg / Delta model

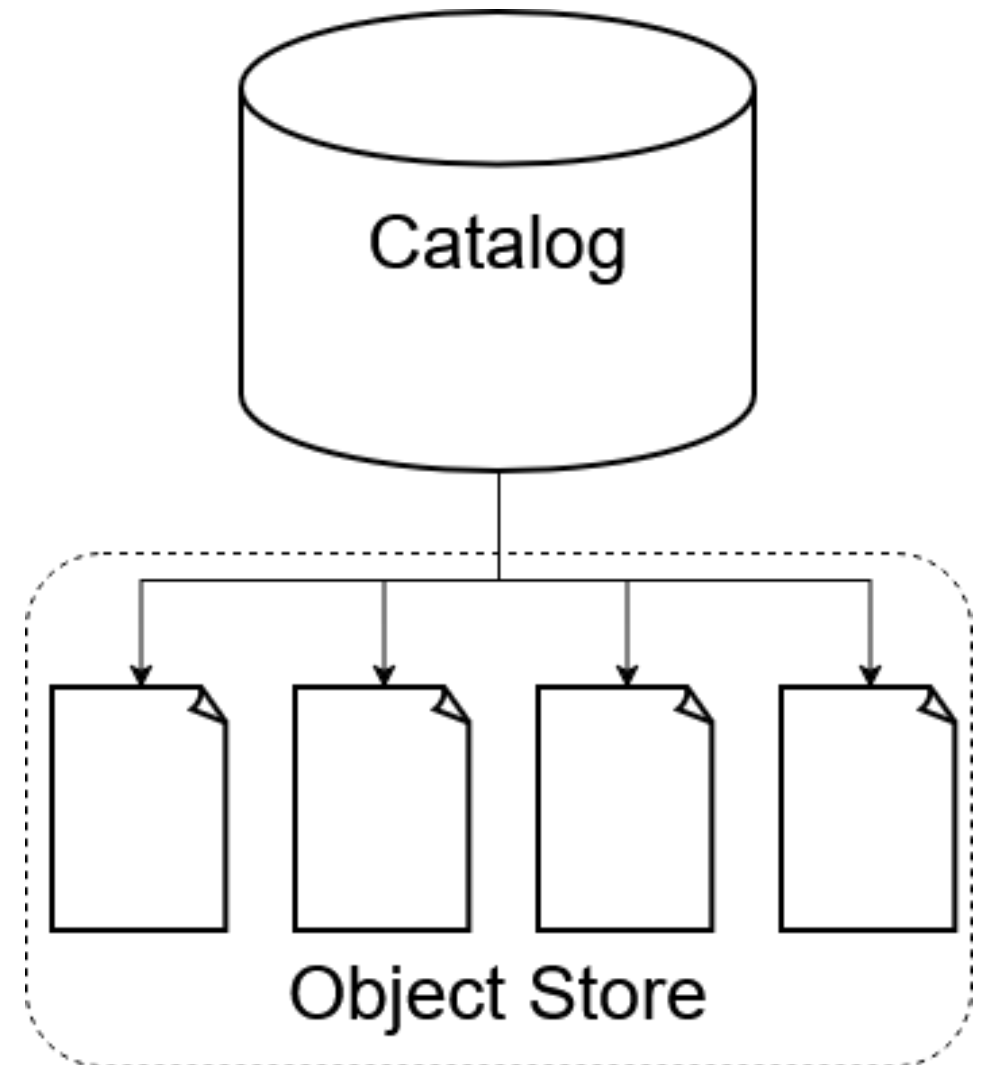
- Requires read / write of many small metadata files
- Increased overhead in query planning
- Increased operational complexity



The DuckLake model

- Leverages catalog DB, which is already part of the stack
- Reduced operational complexity
- Increased performance
- Like Snowflake / BigQuery

No metadata files, only a catalog database



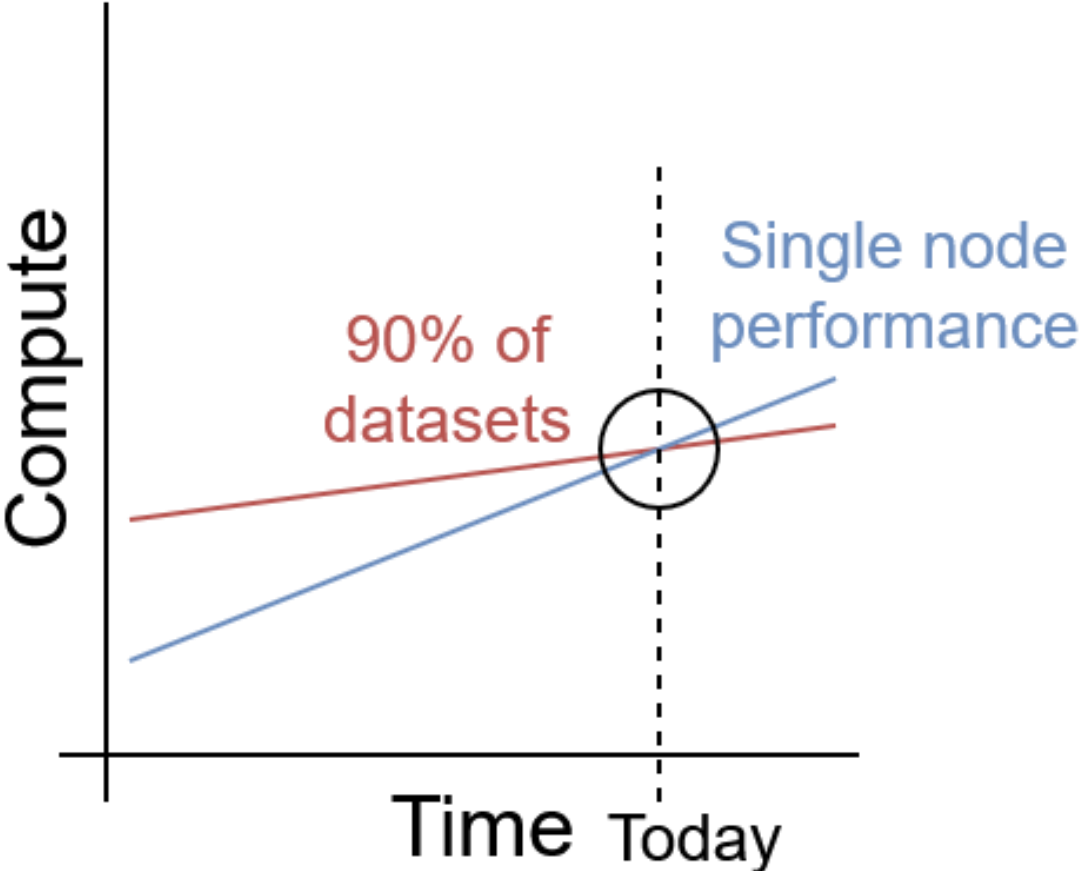
Distributed vs single-node compute



VS



DuckDB



DuckLake in practice

Ecosystem integration and the early-adopter trade-off

Easy setup, local & cloud

```
1 INSTALL ducklake;
2 INSTALL postgres;
3
4 CREATE SECRET pg_secret(
5   TYPE 'postgres',
6   HOST 'localhost',
7   DATABASE 'ducklake_catalog',
8   USER 'ducklake_user',
9   PASSWORD *****
10 );
11
12 CREATE SECRET (
13   TYPE 'ducklake',
14   METADATA_PATH '',
15   DATA_PATH '/path/to/data/dir',
16   METADATA_PARAMETERS MAP {'TYPE': 'postgres', 'SECRET':
17     'pg_secret'}
18 );
19 ATTACH 'ducklake:' AS my_ducklake;
20 USE my_ducklake;
```

```
1 INSTALL ducklake;
2 INSTALL postgres;
3 INSTALL aws;
4 INSTALL httpfs;
5
6 CREATE SECRET s3_secret (
7   TYPE s3,
8   PROVIDER credential_chain,
9   CHAIN config,
10  PROFILE development,
11  REGION 'eu-central-1'
12 );
13
14 CREATE SECRET pg_secret(
15   TYPE 'postgres',
16   HOST 'remote-host.rds.amazonaws.com',
17   DATABASE 'ducklake_catalog',
18   USER 'ducklake_user',
19   PASSWORD *****
20 );
21
22 CREATE SECRET (
23   TYPE 'ducklake',
24   METADATA_PATH '',
25   DATA_PATH 's3://path/to/data/on/s3',
26   METADATA_PARAMETERS MAP {'TYPE': 'postgres', 'SECRET':
27     'pg_secret'}
28 );
29 ATTACH 'ducklake:s3_secret' AS my_ducklake;
30 USE my_ducklake;
```

Ecosystem integration: dbt



Overall: **good**

Some rare minor workarounds required e.g.,

- SCD2 snapshots need dbt-duckdb macro edit + tmp schema

```
profiles.yml
1 my_dbt_profile:
2   target: dev
3   outputs:
4     local:
5       type: duckdb
6       extensions:
7         - httpfs
8         - aws
9         - ducklake
10        - postgres
11      secrets:
12        - type: postgres
13          name: catalog_credentials
14          host: "{{ env_var('DL_HOST') }}"
15          port: "{{ env_var('DL_PORT') }}"
16          user: "{{ env_var('DL_USER') }}"
17          password: "{{ env_var('DL_PASSWORD') }}"
18          database: "{{ env_var('DL_DATABASE') }}"
19        - type: s3
20          provider: credential_chain
21          validation: "exists"
22          refresh: auto
23      attach:
24        - path: "ducklake:postgres:"
25          alias: catalog
26          options:
27            meta_secret: catalog_credentials
28            data_path: "{{ env_var('DL_DATA_PATH') }}"
29      database: catalog
30      schema: main
31      threads: 16
```

Ecosystem integration: dagster



Overall: **very good**

Some rare minor workarounds required e.g.,

- Subclassing of DuckDBResource and connection definition

```
profiles.yml

1 import dagster as dg
2 from dagster_duckdb import DuckDBResource
3 from duckdb import DuckDBPyConnection, connect
4
5
6 class DuckLakeResource(DuckDBResource):
7     """Get a python connection to a DuckLake."""
8
9     @classmethod
10    @contextmanager
11    def get_connection(cls) -> Generator[DuckDBPyConnection]:
12        ducklake = connect()
13        ducklake.sql("""
14            ATTACH 'ducklake:postgres:dbname=ducklake_catalog
15            host=localhost' AS my_ducklake
16            (DATA_PATH 'data_files/');
17            USE my_ducklake;
18        """)
19
20        try:
21            yield ducklake
22        finally:
23            ducklake.close()
```

Ducklake's biggest problem

Early-adopter tax

Works where DuckDB works, but often not without adjustments

- DBT
- Dagster
- Metabase



DuckLake

Ducklake inconveniences

RBAC

- No DuckLake built-in abstraction
- Is delegated to catalog database & object store
- Fix planned for v2.0

Git-like branching

- Currently missing (planned for v2.0)



DuckLake

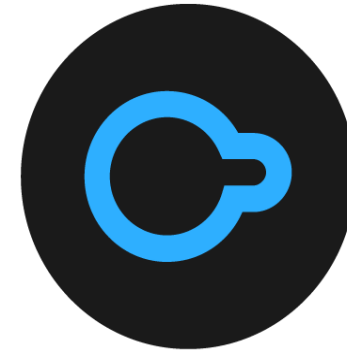
Conclusion

The current state of DuckLake

The current state of DuckLake

Promising project in its early stages

- FOSS
- “Multiplayer” DuckDB
- Simple Lakehouse
- For prod & quick local testing
- More hands-on than fully managed service



DuckLake

Thank you for listening
