# Why does this makes sense?

## The Goal

- **DuckDb + Unity Catalog Integration with dbt**

## Why use Unity catalog

- **Centralized Governance**
- **Access management**



## Why duckDB?

- **Fast**
- **Ease of Setup for analytics**



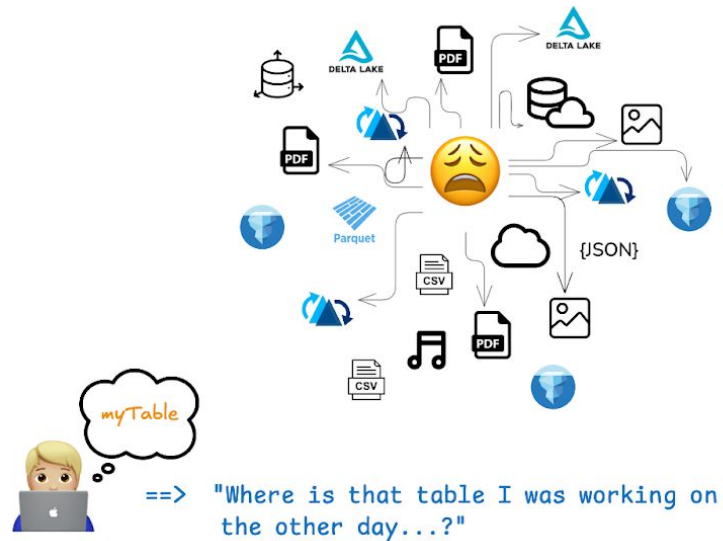## Why dbt?

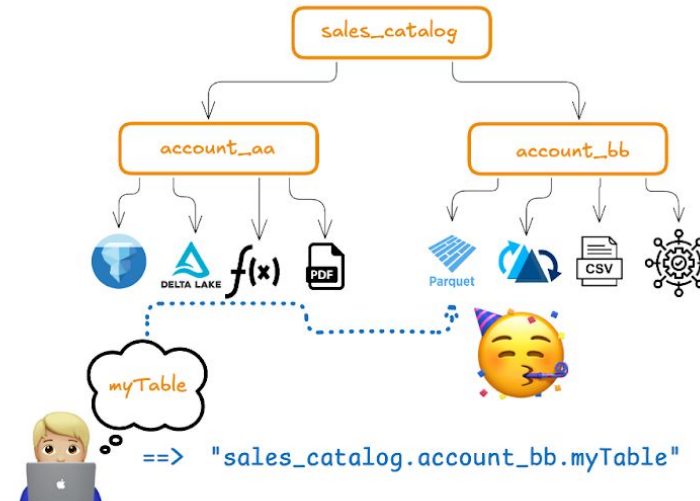- **To automate the the transformation and writing**

# What is Data Catalog ?

## Why might you need one ?

- **A Metadata layer on top of actual data**

- **A Centralized inventory of data assets**



without
a data catalog



==> "Where is that table I was working on the other day...?"

with
a data catalog

sales_catalog

account_aa          account_bb

myTable

==> "sales_catalog.account_bb.myTable"

What is Unity Catalog

# What is Unity Catalog ?

- **Open source version of Databricks' Unity Catalog. Firstly released in June 2024**

- **Has 2 components, CLI and UI.**

- Command line interface

- UI

# Unity Catalog - Structure

Unity Catalog stores all assets in a 3-level namespace:

- catalog
- schema
- assets like tables, volumes, functions, etc.

# What is Unity Catalog ?

# Unity Catalog - Auth

**External Identity Provider** (Google, Okta, Keycloak)

**Unity Catalog**

**User requests access**

① Authenticate User

③ Is user authorized?

② User authenticated

④ Use authorized

② User NOT authenticated

④ User NOT authorized

"Are you allowed to enter?"

② Authorization

"Are you who you say you are?"

① Authentication

# Unity Catalog - AI

# Unity Catalog Demo

**Using CLI**

**Using Python SDK**

**Unity Catalog UI**

**AI**

# Unity Catalog - Difficulties with Unity catalog

- Writing to unity catalog

- Querying

- Giving permissions to users

Unity Catalog Project(dbt+azure
containers +duckdb)

# Unity Catalog Project - Initial Logic (Docker)

# Unity Catalog Project - Current Setup (Azure)

# Unity Catalog Project Demo

Configurations

dbt

Functions

Azure resources

DuckDB

Final Thoughts and limitations

# Unity Catalog vs other Catalogs

- **Direct enforcement of access control on the data assets (Gatekeeper)**
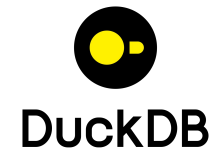
- **Unified Governance for Data & AI with Direct Data Lake Control**

- **Structured Namespace for already known from Databricks**

## Data Catalog Comparison Metascore

# Unity Catalog - What is missing

**THE GOOD**

- Unified Governance (tables, volumes, functions, ML models) and access management

- Open Source and free

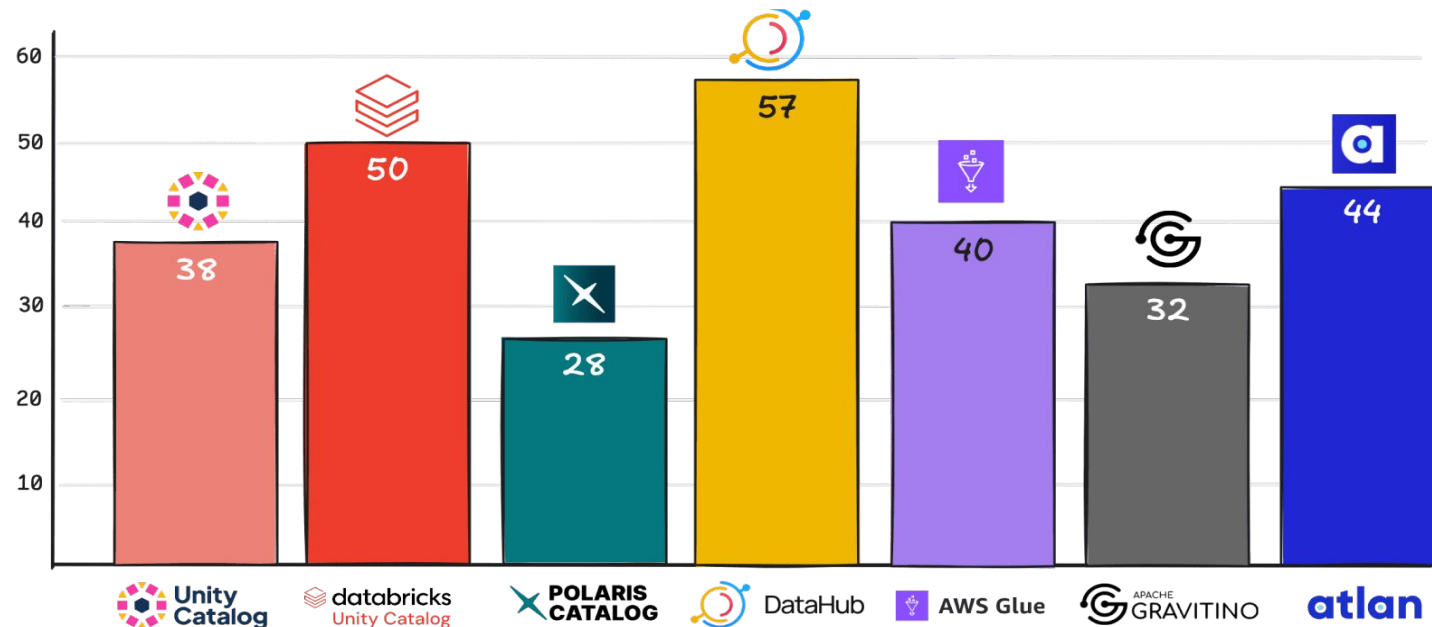- Lots of support and integrations(duckdb, spark)

- Lots of room to grow as it is only on 0.2 version

**THE BAD**

- Limited Data Lineage Capabilities (in OSS):

- No user group permission

- Operational overhead as there is extra layer to manage/debug

- Potential Performance Bottlenecks & Scalability

**THE UGLY**

- (Currently)Cloud Credential Vending Challenges

- Currently there are some bugs and some things not working as intended

- Maturity Gap Compared to Commercial Offerings:

# Additional Use cases

## This can be something for you?

Unity Catalog has integrations with :

- Celerdata
- Daft
- DuckDB
- PuppyGraph
- Spark
- SpiceAI
- Trino
- XTable

Also some ai integrations

- LangChain
- LlamaIndex
- OpenAI
- Anthropic
- CrewAI
- AutoGen
- LiteLLM
- Gemini

dataroots
a TAIAN company

Thank you and questions