# Starburst Workshop:
## Python data pipelines & products

Starburst Academy
Data Universe 2024

# Workshop objectives

- Present the reference architecture spanning multiple zones
  - Land
  - Structure
  - Consume
- Build a Python-based pipeline that spans the reference architecture's zones
- Create and secure granular data products for your downstream consumers

Starburst

# For today...

- POIs
  - The lab guide has it ALL!
  - The S3 credentials will expire after the weekend
  - All code is in the Jupyter notebook – `Workshop4.ipynb`
- Approach
  - I'm going to perform the labs myself & you can...
    - do them along with me
    - or just watch (and maybe do later)
  - We ALL come from DIFFERENT experiences & backgrounds, so...
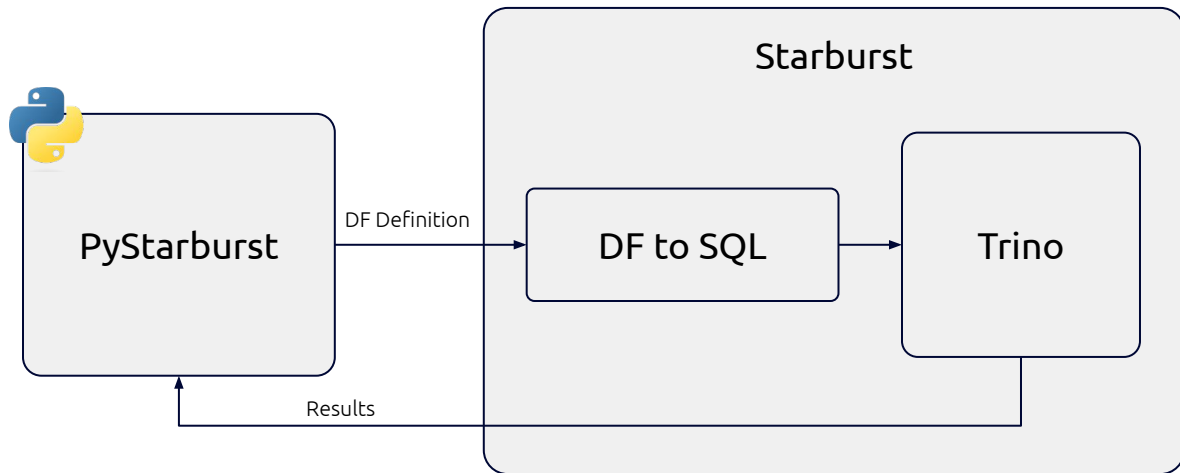    - **ASK QUESTIONS**

❄️ Starburst

# PyStarburst Overview

```
df_missions = df_missions.with_column("date", f.sql_expr("COALESCE(TRY(date_parse(\"date\", '%a %b %d, %Y %H:%i UTC')), NULL)"))

print(df_missions.schema)

df_missions = df_missions\
    .filter(col("date") > datetime(2000, 1, 1))\
    .sort(col("date"), ascending=True)

df_missions.show()
```

PyStarburst → DF Definition → Starburst

DF to SQL → Trino

Results

- PySpark-like Syntax
- Lazy execution
- Python gets converted to SQL
- Heavy lifting done by Trino

Links: API Docs & Example Code

Currently supported on Starburst Galaxy and Starburst Enterprise.

Starburst

# Data pipelines

## Reference architecture

Starburst

# Reference architecture

The reference architecture centers around the data lakehouse and how we classify our data assets into distinct zones.  Data pipelines populate the zones.

# Activities across the architecture



Data lakehouse

Multiple data sources → **Land** → **Structure** → **Consume** → Multiple BI / reporting engines & ad hoc queries

ML notebooks
Data modeling
Data applications

| Upstream data is created by apps, websites, tools, etc | Raw data is **landed** & kept | Data is cleansed, enriched & **structured** | Views & rollups are created for reporting | People & processes **consume** the data |

Starburst

# Tools & technologies

# Summary

# Data products

**Data consumers & producers**

Starburst

# Qualities of data products

**Discoverable**: end users and other domains need to be able to discover and access a given data product

**Addressable**: the data should have a straightforward and documented way of being programmatically accessed, e.g. via SQL

**Trustworthy**: end users should be able to understand the level of data quality and ideally view the provenance (lineage) of the data so they can be confident in any analyses using the data product

**Self-describing**: any end user outside the domain which produces the data product should have all of the information they require to use the data

**Interoperable**: governance should ensure that the data complies to any inter- or intra-domain standards or regulations, so the end user can confidently use the data without concern

**Secure**: data products should fold any authorization into the access control provided by the data mesh experience plane, which is where data product consumption occurs

Starburst

# Starburst Data Products

# Discover, create, publish, and manage
data products based on multiple datasets

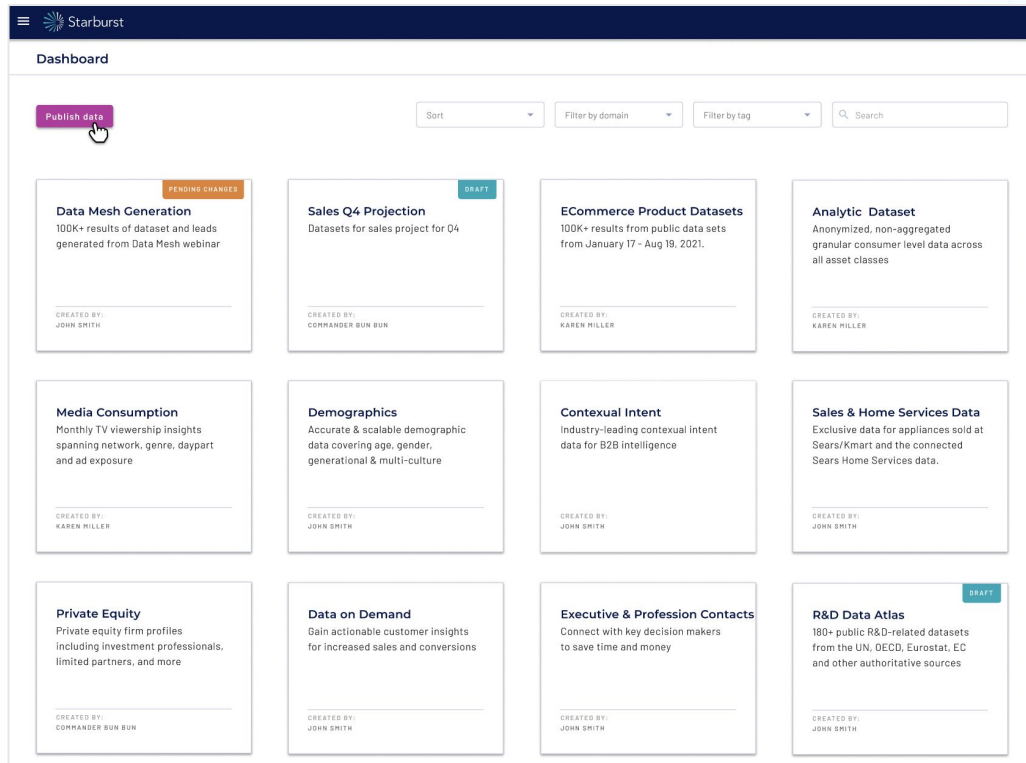Allow data owners & data engineers to **define relevant metadata** to data consumers

**Secure** your data products with **access control**, ensuring consistent **governance** from the source level

# Query data products

that are trusted and approved for frequent business use

# Share and rate your data products

internally and track usage