# Mastering Multitenant Orchestration with dbt and Dagster

## Life of Data Engineers does not have to be that hard

Andrea Montes

DBT Bogotá Meetup

2024-06

# Audience questions

1. Airflow familiarity?
2. Crazy tools difficult to debug?

# What is needed?

- Daily updates to client dashboards. +250 different clients
- Custom reports per client
- Product and business questions

# Legacy product(s) - Data Warehouse
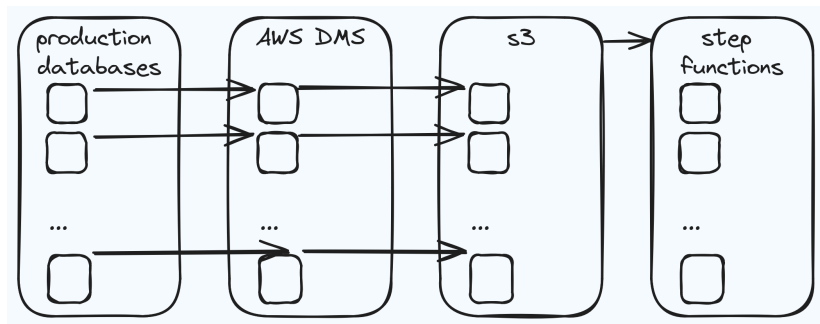


Figure: Legacy DWH product

- ▶ Daily refresh
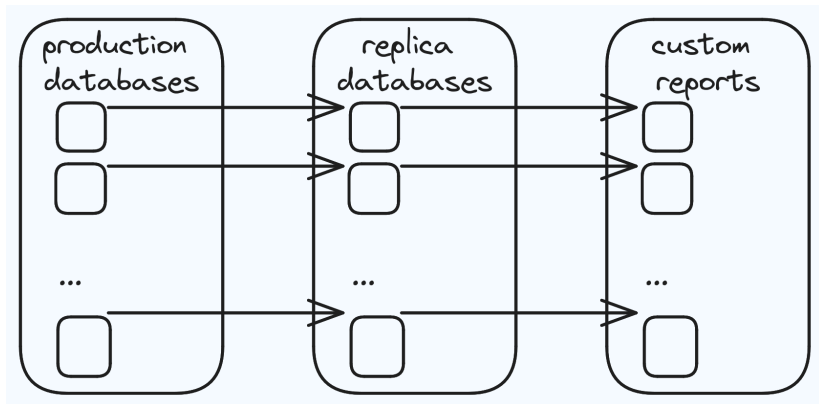- ▶ Failed 2 or 4 times a week for heavy clients

# Legacy product(s) - Reporting



Figure: Legacy reporting product

▶ Replicated once a week
▶ Raw sql
▶ Stored procedures

# We need to re-design

# Different users, different requirements

### Dashboards
Clients need their dashboards updated to know events statuses

### Business questions
What are the most used features?, how virtual vs in-person events attendance has changed after pandemic?

### Client questions
How long is taking a candidate to become an applicant? Do I have a diverse pipeline?
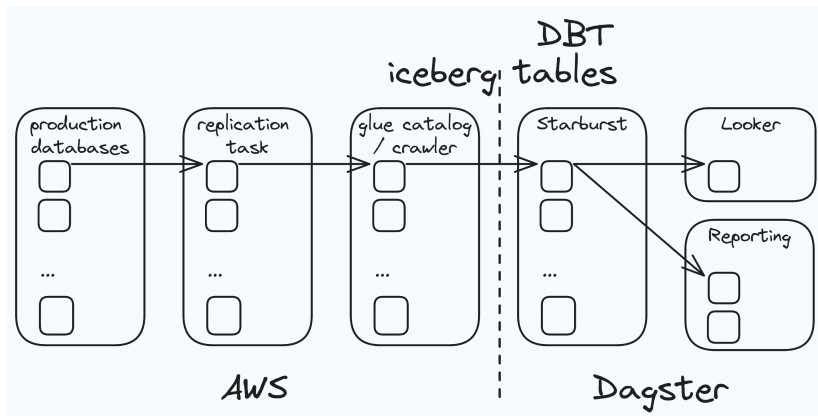
# New architecture



Figure: New architecture to accomplish reporting and DWH requirements

# DMS

- Cheap...
- Not reliable as we would like
- Trade-off

# Glue catalog and Crawler

- ► Easy enough to implement
- ► Our compute engine support it

# Starburst - DBT

**Pros**:

- ► Starburst is a vendor option for Trino
- ► Cheaper than snowflake
- ► Trino has a good community

**Cons**:

- ► SQL ANSI
- ► It's a new product, random changes
- ► Compute throttling
- ► We were their first large user

# Minimum Example: DBT and multiple clients

Data used:

- ▶ Open data Colombia - Mobile Telephony subscribers by category
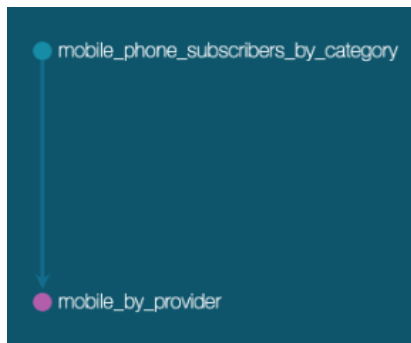- ▶ One CSV with information for all mobile telephony providers
- ▶ We need to split this data in a convenient way



Figure: Example: Example Lineage

# DBT and multiple clients

**Model: Mobile by provider**

```sql
SELECT
    ANO AS year
    , TRIMESTER AS trimester
    , PROVEEDOR AS provider
    , "LINEAS EN SERVICIO" AS lines_in_service
    , "LINEAS PREPAGO" AS prepaid_lines
    , "LINEAS POSPAGO" AS postpaid_lines
    , "LINEAS ACTIVADAS" AS enabled_lines
    , "LINEAS RETIRADAS" AS retired_lines
FROM {{ ref('mobile_phone_subscribers_by_category')}}
WHERE PROVEEDOR = '{{ var("client")}}'
```

# DBT and multiple clients

**DBT command:**

```
dbt run --profiles-dir dbt_project/.dbt --select
    mobile_by_provider --vars '{client:   AVANTEL S
    . A . S }'
```

**SQL compiled:**

```sql
CREATE TABLE
      "housing"."avantelsas"."mobile_by_provider"
    AS (
SELECT
    ANO AS year
    , TRIMESTER AS trimester
    , PROVEEDOR AS provider
    , "LINEAS EN SERVICIO" AS lines_in_service
    , "LINEAS PREPAGO" AS prepaid_lines
    , "LINEAS POSPAGO" AS postpaid_lines
    , "LINEAS ACTIVADAS" AS enabled_lines
    , "LINEAS RETIRADAS" AS retired_lines
FROM housing.telephony_sources.
    mobile_phone_subscribers_by_category
WHERE PROVEEDOR = 'AVANTEL S.A.S')
```

# DBT and multiple clients



Figure: Example: client = AVANTEL S.A.S



Figure: Example: client = UFF MOVIL

# Orchestration tool: Dagster



- ▶ We run the job daily starting at 2 am CST
- ▶ Each job client partition takes about 7 mins
- ▶ Job finishes before 7 am CST

Figure: Dagster: How do we run +250 transformations?

# Orchestration tool: Dagster



Figure: Dagster: How do we run +250 transformations?

# Takeaways

▶ Not getting used to what is poorly done

▶ Don't fall in love with a technology product

▶ Invest time learning new tools!

▶ Fail fast mindset

# Who am I?



Figure: Andrea Montes

**in** Andrea Montes - Senior Data Engineer

**○** mamontesp

**○** Code example: `https://github.com/mamontesp/mastering-orchestration-dbt-dagster`

**▪** My blog: `https://mamontesp.ghost.io/`

# QRs!



Figure: Andrea Montes in [in]



Figure: Code Example