

# Beyond dbt: Modern SQL Transformation and Lineage

Tomas Peluritis





# Tomas Peluritis

Head of Data @ Mediatech  
AKA Uncle Data



<https://www.linkedin.com/in/tomaspeluritis/>



<https://podcasters.spotify.com/pod/show/duomenu-dede>



<https://uncledata.substack.com>



# Let's talk about elephant in the room



# SQL Transformation Timeline

Aug 2025



Slides polished and ready  
for November



# SQL Transformation Timeline

Aug 2025



Slides polished and ready  
for November

Sept 2025  
Acquisition



Fivetran acquires Tobiko.



# SQL Transformation Timeline





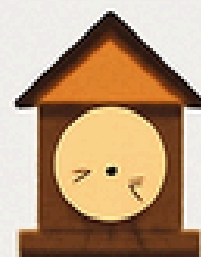
# Agenda

Data Lineage  
Market Overview

Quick intro  
to dbt

What's SQLGlot  
and SQLMesh

Common  
approaches  
to implement  
Data Lineage



# **PHASE 1      PHASE 2      PROFIT**

---

Collect  
Data Lineage



Profit



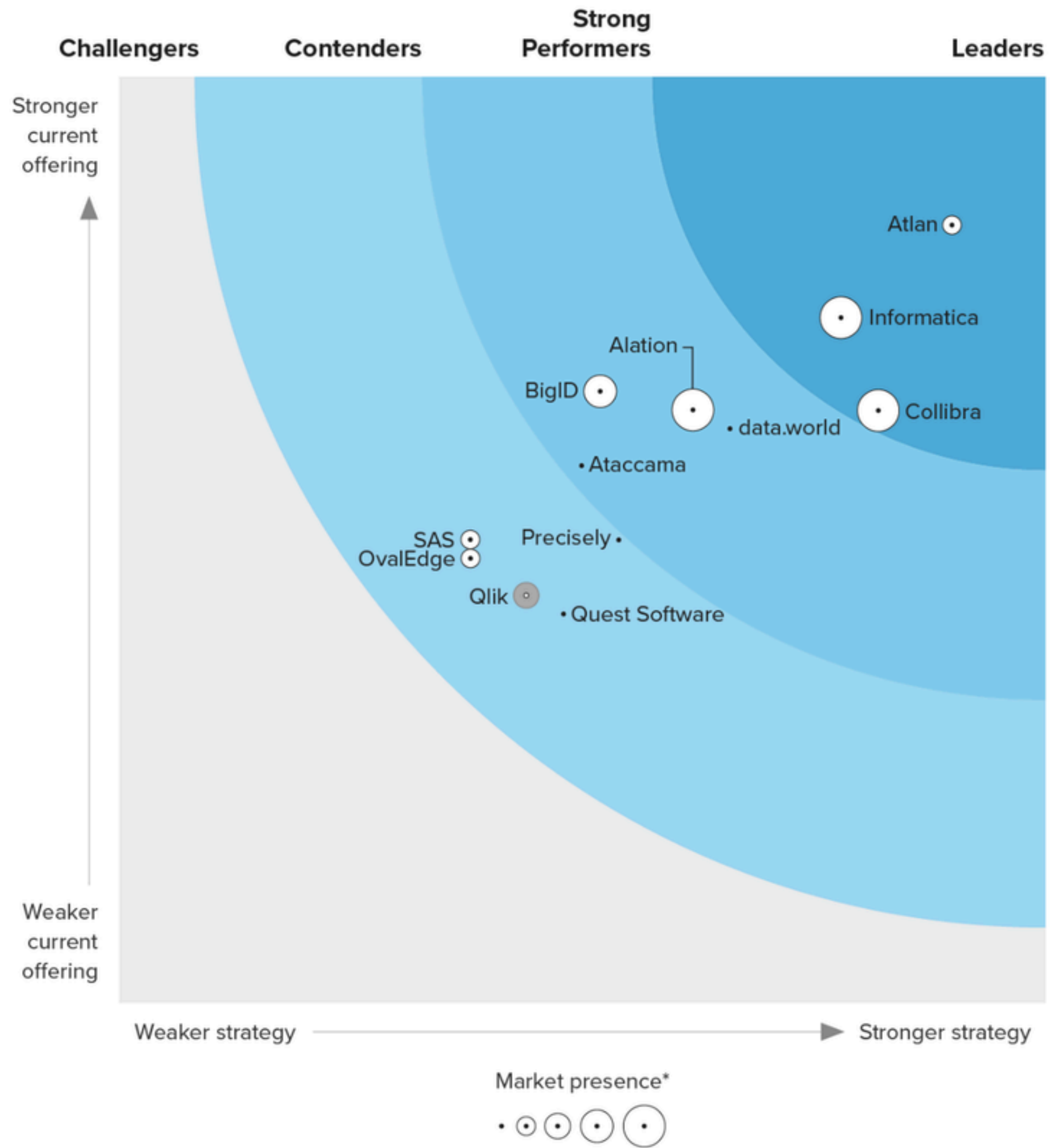


# \$1.5-2B

Lineage market size

The **growing demand** for data lineage solutions reflects the increasing complexity of data ecosystems, with businesses recognizing the need for improved **visibility and accountability** in data management.

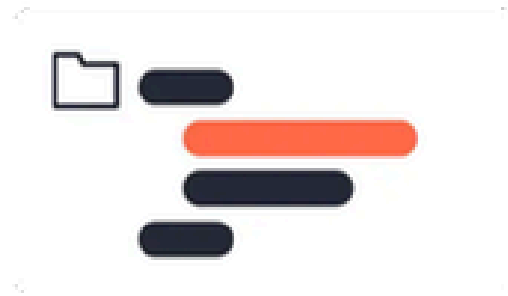
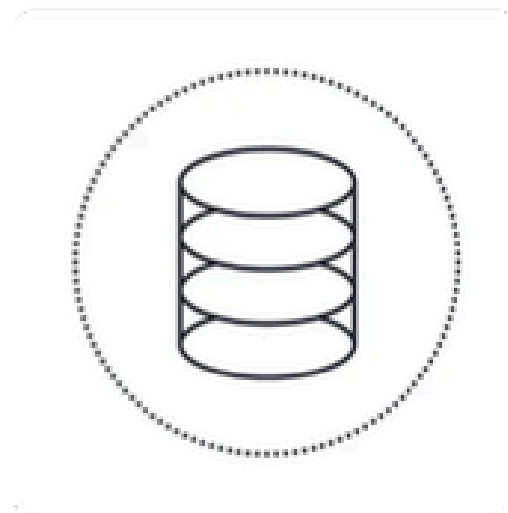
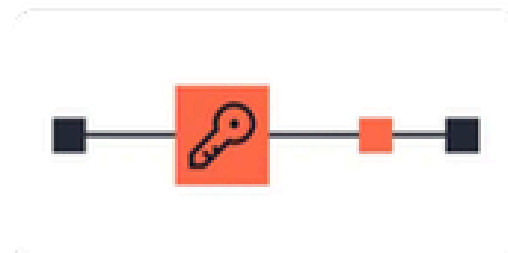
# Forrester Report



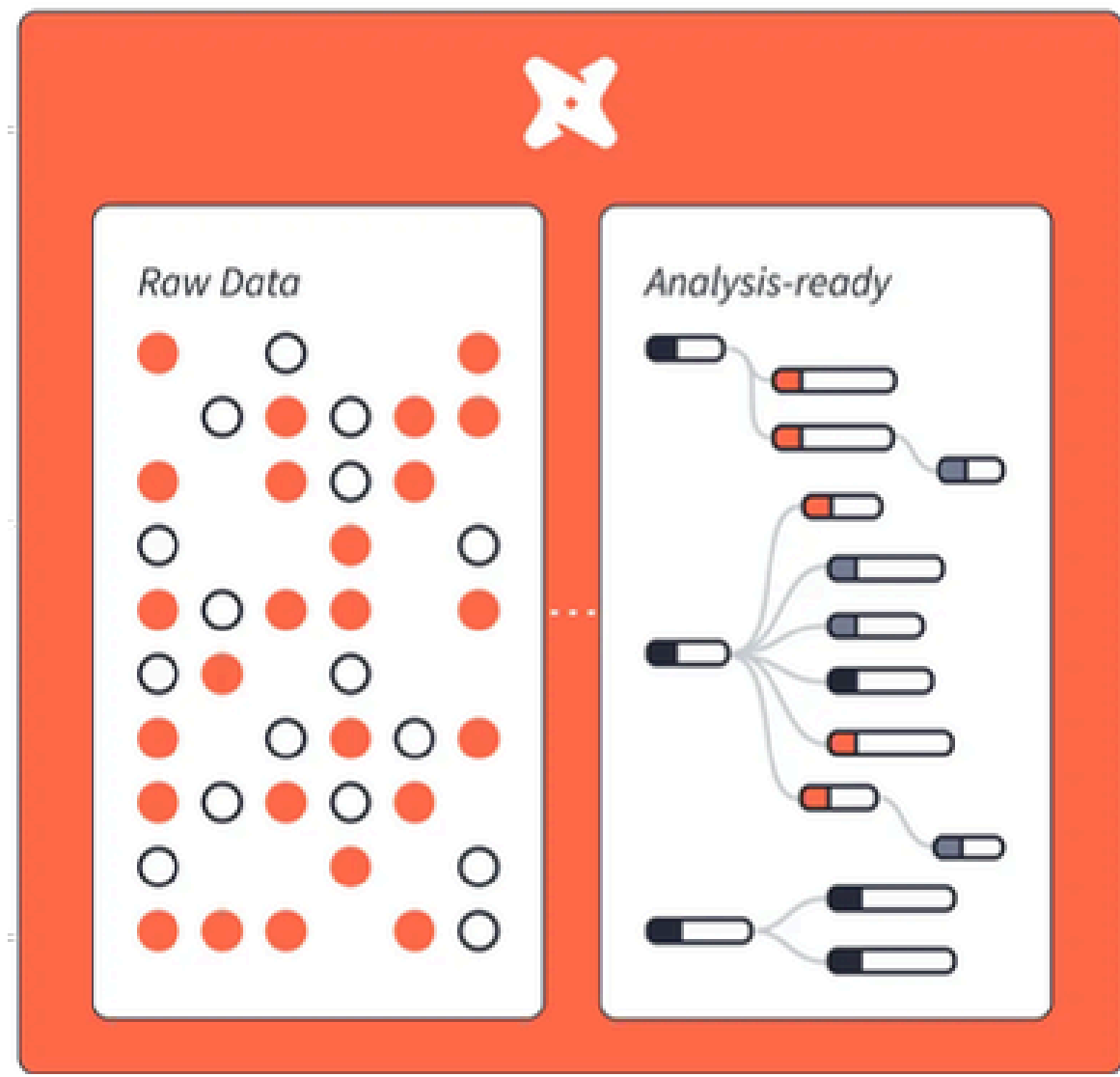
# Introduction to dbt



Data Sources



Data Platform



Analysis



# SQLGlott Deep Dive



# SQLGlot Deep Dive

SQL parser, transpiler, optimizer & engine

Zero dependencies

Supports 24+ dialects



# SQLGlot Deep Dive

SQL parser, transpiler, optimizer & engine

Zero dependencies

Supports 24+ dialects

Comprehensive SQL parsing with robust testing

Syntactically & semantically correct output

High performance despite pure Python

Format & translate SQL effortlessly



# SQLGlot Parser and AST



# Abstract Syntax Tree (AST)

## Understanding SQLGlot's Abstract Syntax Tree

The Abstract Syntax Tree (AST) visually represents SQL queries, showcasing the structure and relationships within the code, enabling easier parsing, optimization, and transformation for various SQL engines.



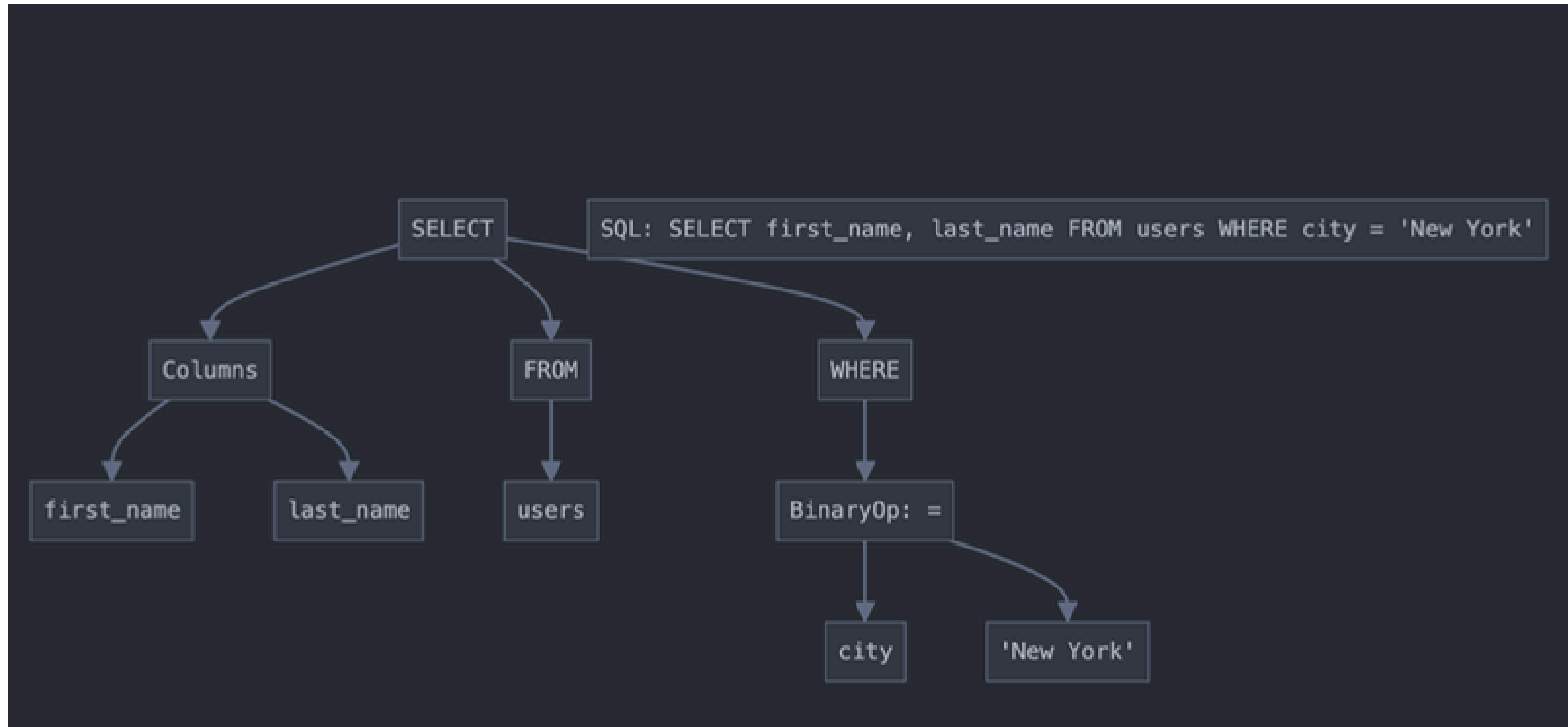
# Abstract Syntax Tree (AST)



```
select first_name, last_name from users where city = 'New York'
```



# Abstract Syntax Tree (AST)



# Abstract Syntax Tree (AST)

```
Select(  
  expressions=[  
    Column(  
      this=Identifier(this=first_name, quoted=False)),  
    Column(  
      this=Identifier(this=last_name, quoted=False))],  
  from=From(  
    this=Table(  
      this=Identifier(this=users, quoted=False))),  
  where=Where(  
    this=EQ(  
      this=Column(  
        this=Identifier(this=city, quoted=False)),  
      expression=Literal(this='New York', is_string=True)))
```

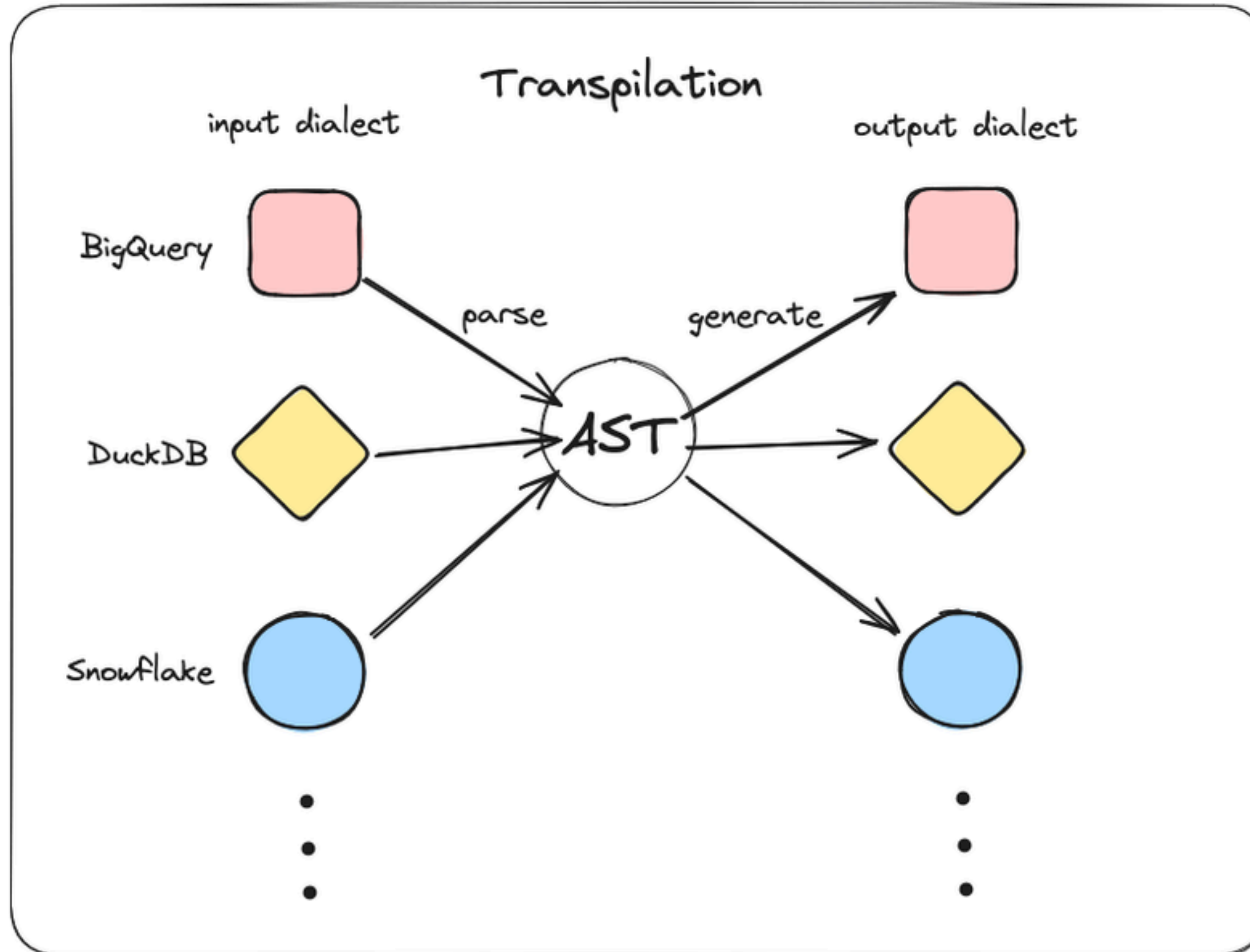


# Why ASTs Are Important



# Transpiration



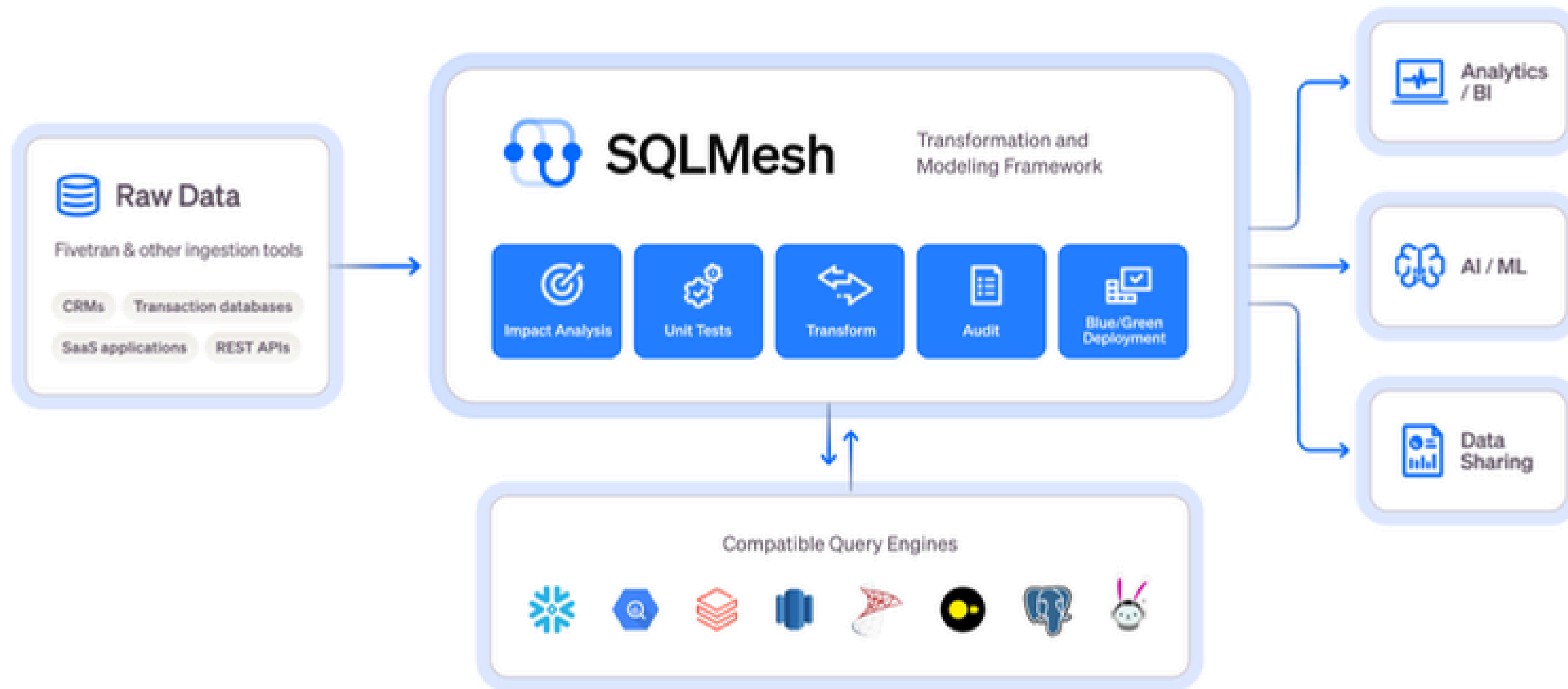


Source



# Introduction to SQLMesh





[sqlmesh.com](https://sqlmesh.com)



SQLMesh by Tobiko

Environment: dev | Changes: 1 | No Errors

Custom SQL: `stg_supplies.py`

```
1 from sqlmesh.core.model import model
2 import pandas as pd
3 import hashlib
4
5 @model
6 "jaffle_shop.staging.stg_supplies",
7 kind="FULL",
8 cross="pdply",
9 columns=[
10     "supply_uuid": "STRING",
11     "supply_id": "STRING",
12     "product_id": "STRING",
13     "supply_name": "STRING",
14     "supply_cost": "FLOAT",
15     "is_perishable_supply": "BOOLEAN"
16 ],
17 column_descriptions=[
18     "supply_uuid": "The unique key of our supplies per cost",
19     "supply_id": "Supply identifier",
20     "product_id": "Product identifier",
21     "supply_name": "Supply name",
22     "supply_cost": "Supply cost in dollars",
23     "is_perishable_supply": "Whether the supply is perishable"
24 ],
25 descriptions=[
26     "List of our supply expenses data with basic cleaning and transformation applied.
27     One row per supply cost, not per supply. As supply costs fluctuate they receive
28     a new row with a new UUID. Thus there can be multiple rows per supply_id."
29 ]
30
31 def stg_supplies(context, **kwargs):
32     # Get the raw supplies data
33     raw_supplies = context.fetcher.fetch("select * from {{context.resolve_table('jaffle_shop.raw.raw_supplies')}}")
34
35     # Apply transformations
36     df = pd.DataFrame(raw_supplies)
37
38     # Generate composite key directly instead using the macros
39     df["supply_uuid"] = df.apply(lambda row: hashlib.md5(f"{row['id']}-{row['sku']}".encode()).hexdigest(), axis=1)
40
41     # Convert cents to dollars
42     df["supply_cost"] = df["cost"] / 100
43
44     # Rename columns
45     df = df.rename(columns={
46         "id": "supply_id",
47         "sku": "product_id",
48         "name": "supply_name",
49         "perishable": "is_perishable_supply"
50     })
51
52     # Select only needed columns
53     df = df[["supply_uuid", "supply_id", "product_id", "supply_name", "supply_cost", "is_perishable_supply"]]
54
55     return df
```

Save | Language: Python | Select program: SQLMesh Type:

Linkage

jaffle\_shop.staging.stg\_supplies | 20:14 | 2024-04-11 | Upstream/Downstream: 1



# SQLMesh Capabilities



# SQLMesh State Database



# It has a DB for state

```
{% if is_incremental() %}  
  
    -- this filter will only be applied on an incremental run  
    -- (uses >= to include records whose timestamp occurred since the last run of this model)  
    -- (If event_time is NULL or the table is truncated, the condition will always be true and load all records.  
where event_time >= (select coalesce(max(event_time), '1900-01-01') from {{ this }} )  
  
{% endif %}
```



# It has a DB for state

```
{% if is_incremental() %}

-- this filter will only be applied to records that have occurred since the last run
-- (uses >= to include records whose event_time is equal to the last run of this model)
-- (If event_time is NULL or the table is empty, the condition will always be true and load all records)
where event_time >= (select coalesce(max(event_time), '2000-01-01') from {{ this }} )

{% endif %}
```

▼	sqlmesh
▼	Tables 7
>	_auto_restatements -1 rows, 16KB
>	_environment_statements -1 rows, 16KB
>	_environments -1 rows, 16KB
>	_intervals -1 rows, 32KB
>	_plan_dags -1 rows, 24KB
>	_snapshots -1 rows, 24KB
>	_versions -1 rows, 16KB



# It has a DB for state

Modified Directly

staging\_\_dev.stg\_orders Breaking Change

- └─ marts\_\_dev.orders
- └─ marts\_\_dev.customers
- └─ marts\_\_dev.order\_items

Breaking Change  
It will rebuild all models

Non-Breaking Change  
It will exclude all indirect models caused by this change

Forward-Only Change  
The change requires no rebuilding

```
---  
+++  
  
@@ -2,7 +2,13 @@  
  
name jaffle_shop.staging.stg_orders,  
start '2016-09-01',  
dialect postgres,  
- kind FULL,  
+ kind INCREMENTAL_BY_TIME_RANGE (  
+   time_column ("ordered_at", 'YYYY-MM-DD'),  
+   partition_by_time_column TRUE,  
+   forward_only FALSE,  
+   disable_restatement FALSE,  
+   on_destructive_change 'ERROR'  
+ ),  
columns (
```



# It has a DB for state

```
@cents_to_dollars(order_total) AS order_total,  
CAST(ordered_at AS DATE) AS ordered_at  
FROM jaffle_shop.raw.raw_orders  
+WHERE  
+ ordered_at BETWEEN @start_ds AND @end_ds
```

marts\_\_dev.orders Breaking Change

**Modified Indirectly**

- marts\_\_dev.customers
- marts\_\_dev.order\_items

**Backfills**

Models 4

- marts\_\_dev.customers
- marts\_\_dev.order\_items
- marts\_\_dev.orders
- staging\_\_dev.stg\_orders



# dbt vs SQLMesh Models



# Model Definition

```
1 select
2     id as customer_id,
3     name as customer_name
4 from {{ source('ecom', 'raw_customers') }}
5
```

```
1 models:
2   - name: stg_customers
3     description: Customer data with basic cleaning and transformation applied, one row per customer.
4     columns:
5       - name: customer_id
6         description: The unique key for each customer.
7         data_tests:
8           - not_null
9           - unique
10
```



# SQLMesh



```
1 MODEL (  
2   name jaffle_shop.staging.stg_customers,  
3   kind FULL,  
4   cron '@daily',  
5   description 'Customer data with basic cleaning and transformation applied, one row per customer',  
6   columns (  
7     customer_id STRING COMMENT 'The unique key for each customer',  
8     customer_name STRING COMMENT 'Customer full name'  
9   )  
10 );  
11  
12 SELECT  
13   id AS customer_id,  
14   name AS customer_name  
15 FROM  
16   jaffle_shop.raw.raw_customers
```



# Testing Approaches Comparison



# dbt

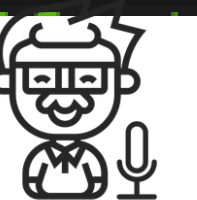
```
1 models:
2   - name: customers
3     description: Customer overview data mart, offering key details for each unique customer. One row per customer.
4     data_tests:
5       - dbt_utils.expression_is_true:
6         expression: "lifetime_spend_pretax + lifetime_tax_paid = lifetime_spend"
7     columns:
8       - name: customer_id
9         description: The unique key of the orders mart.
10        data_tests:
11          - not_null
12          - unique
13        - name: customer_name
14          description: Customers' full name.
15        - name: count_lifetime_orders
16          description: Total number of orders a customer has ever placed.
17        - name: first_ordered_at
18          description: The timestamp when a customer placed their first order.
19        - name: last_ordered_at
20          description: The timestamp of a customer's most recent order.
21        - name: lifetime_spend_pretax
22          description: The sum of all the pre-tax subtotals of every order a customer has placed.
23        - name: lifetime_tax_paid
24          description: The sum of all the tax portion of every order a customer has placed.
25        - name: lifetime_spend
26          description: The sum of all the order totals (including tax) that a customer has ever placed.
27        - name: customer_type
28          description: Options are 'new' or 'returning', indicating if a customer has ordered more than once or has only placed their first order to date.
29        data_tests:
30          - accepted_values:
31            values: ["new", "returning"]
```



# SQLMesh

```
1 audits (  
2   /* Primary key validation */  
3   not_null(columns := (customer_id)),  
4   unique_values(columns := (customer_id)),  
5  
6   /* Data integrity checks */  
7   forall(criteria := (  
8     ABS(lifetime_spend_pretax + lifetime_tax_paid - lifetime_spend) < 0.01  
9   )),  
10  
11  /* Categorical values validation */  
12  accepted_values(column := customer_type, is_in := ('new', 'returning')),  
13  
14  /* Logical consistency checks */  
15  forall(criteria := (  
16    (customer_type = 'returning' AND count_lifetime_orders > 1) OR  
17    (customer_type = 'new' AND count_lifetime_orders <= 1)  
18  )),  
19  
20  /* Temporal logic checks */  
21  forall(criteria := (  
22    (first_ordered_at IS NULL) OR (last_ordered_at IS NULL) OR (first_ordered_at <= last_ordered_at)  
23  )),  
24  
25  /* Non-negative values for monetary fields */  
26  forall(criteria := (  
27    lifetime_spend_pretax >= 0  
28  )),  
29  forall(criteria := (  
30    lifetime_tax_paid >= 0  
31  )),  
32  forall(criteria := (  
33    lifetime_spend >= 0  
34  )),  
35  
36  /* Check repeat buyer flag consistency */  
37  forall(criteria := (  
38    (is_repeat_buyer = TRUE AND count_lifetime_orders > 1) OR  
39    (is_repeat_buyer = FALSE AND count_lifetime_orders <= 1)  
40  ))  
41 )
```

```
Found 10 audit(s).  
not_null on model jaffle_shop.marts.customers ✓ PASS.  
unique_values on model jaffle_shop.marts.customers ✓ PASS.  
forall on model jaffle_shop.marts.customers ✓ PASS.  
accepted_values on model jaffle_shop.marts.customers ✓ PASS.  
forall on model jaffle_shop.marts.customers ✓ PASS.  
forall on model jaffle_shop.marts.customers ✓ PASS.  
forall on model jaffle_shop.marts.customers ✓ PASS.  
forall on model jaffle_shop.marts.customers ✓ PASS.  
forall on model jaffle_shop.marts.customers ✓ PASS.  
forall on model jaffle_shop.marts.customers ✓ PASS.  
Finished with 0 audit errors and 0 audits skipped.  
Done.
```



# Breaking Changes Comparison





# Breaking Changes

```
Models:
├── Directly Modified:
│   └── staging.stg_customers
├── Indirectly Modified:
│   └── marts.customers
└── ...

+++

@@ -9,6 +9,5 @@
)
)
)
SELECT
- id AS customer_id,
- name AS customer_name
+ id AS customer_id
FROM jaffle_shop.raw.raw_customers

Directly Modified: staging.stg_customers (Breaking)
├── Indirectly Modified Children:
│   └── marts.customers (Indirect Breaking)
Models needing backfill:
├── marts.customers: [full refresh]
├── staging.stg_customers: [full refresh]
Apply - Backfill Tables [y/n]: █
```

\*Bonus: suggests to backfill as well





# Detection Methods Comparison

## Manual vs. Automatic Processes

### dbt Manual Detection

Manual detection in dbt involves **extensive code review**, which can be time-consuming and prone to human error, leading to delays in identifying breaking changes during the development process.

### SQLMesh Automatic Detection

SQLMesh streamlines the detection process by offering **automatic identification** of breaking changes, which accelerates development cycles and enhances overall reliability in managing SQL transformation projects.

### Efficiency of Processes

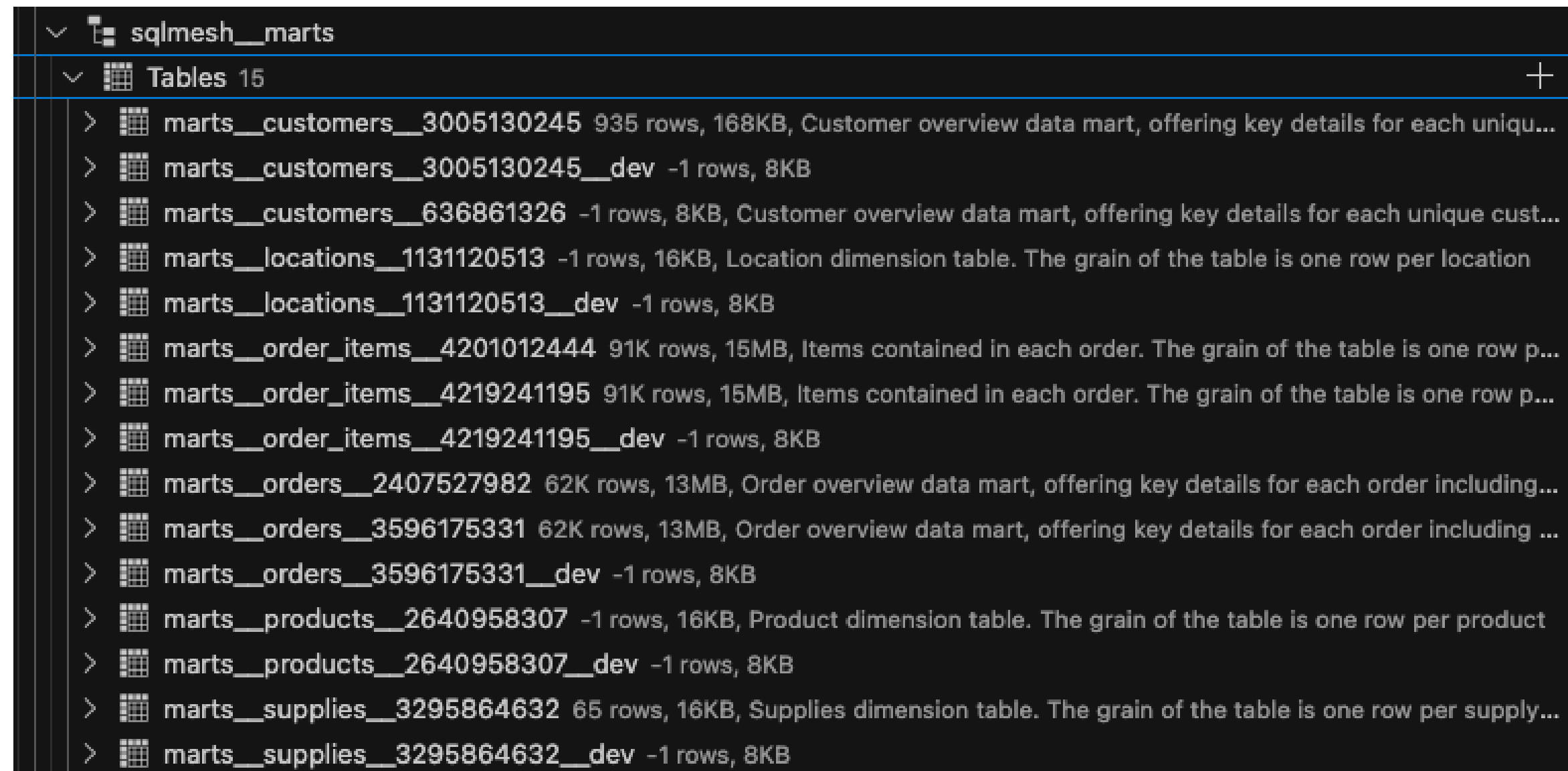
The shift from manual to automatic detection not only improves **efficiency** but also fosters a more collaborative environment, allowing teams to focus on innovation rather than burdening themselves with manual checks.

# Physical vs. Virtual Versioning



# Data Versioning

## Physical layer



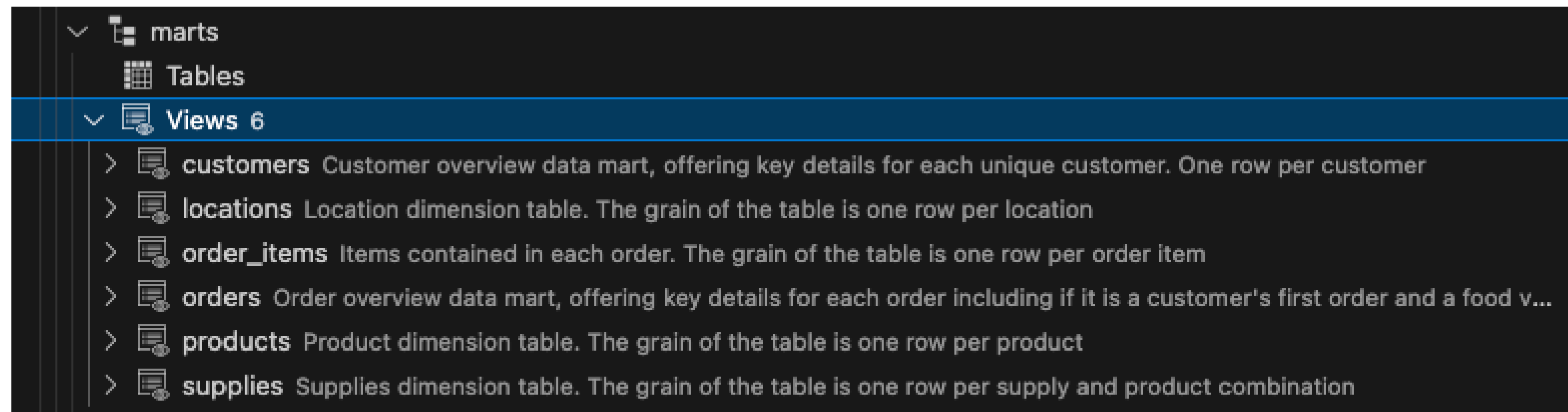
The screenshot displays a database interface showing a schema named 'sqlmesh\_\_marts'. Underneath, there is a folder labeled 'Tables 15' which contains a list of 15 tables. Each table entry includes its name, row count, size, and a brief description of its purpose and grain.

Table Name	Rows	Size	Description
marts__customers__3005130245	935	168KB	Customer overview data mart, offering key details for each unique...
marts__customers__3005130245__dev	-1	8KB	
marts__customers__636861326	-1	8KB	Customer overview data mart, offering key details for each unique cust...
marts__locations__1131120513	-1	16KB	Location dimension table. The grain of the table is one row per location
marts__locations__1131120513__dev	-1	8KB	
marts__order_items__4201012444	91K	15MB	Items contained in each order. The grain of the table is one row p...
marts__order_items__4219241195	91K	15MB	Items contained in each order. The grain of the table is one row p...
marts__order_items__4219241195__dev	-1	8KB	
marts__orders__2407527982	62K	13MB	Order overview data mart, offering key details for each order including...
marts__orders__3596175331	62K	13MB	Order overview data mart, offering key details for each order including ...
marts__orders__3596175331__dev	-1	8KB	
marts__products__2640958307	-1	16KB	Product dimension table. The grain of the table is one row per product
marts__products__2640958307__dev	-1	8KB	
marts__supplies__3295864632	65	16KB	Supplies dimension table. The grain of the table is one row per supply...
marts__supplies__3295864632__dev	-1	8KB	



# Data Versioning

## Virtual layer



\*Bonus adds metadata to DB if possible



# Data Versioning

## Virtual layer

```
1 CREATE OR REPLACE VIEW marts.customers
2 AS SELECT marts__customers__3005130245.customer_id,
3         marts__customers__3005130245.customer_name,
4         marts__customers__3005130245.count_lifetime_orders,
5         marts__customers__3005130245.is_repeat_buyer,
6         marts__customers__3005130245.first_ordered_at,
7         marts__customers__3005130245.last_ordered_at,
8         marts__customers__3005130245.lifetime_spend_pretax,
9         marts__customers__3005130245.lifetime_tax_paid,
10        marts__customers__3005130245.lifetime_spend,
11        marts__customers__3005130245.customer_type
12 FROM sqlmesh__marts.marts__customers__3005130245;
```



# Lineage Approaches

Lineage Approaches:  
Just Buy a Tool that fits  
your need

# Lineage Approaches: dbt Native

# Just Use dbt

The screenshot displays the dbt CLI interface. On the left, the 'File Explorer' shows a project structure with folders like 'analysis', 'dbt\_packages', 'etc', 'logs', 'models', 'intermediate', 'marts', 'metrics', 'staging', 'seeds', and 'tests'. The 'models' folder is expanded, showing 'orders.sql' selected. The main editor shows the following SQL code:

```
1 with orders as (  
2     select * from {{ ref('int_order_payments_pivoted') }}  
3 )  
4  
5 *  
6  
7 customers as (  
8     select * from {{ ref('int_customer_order_history_joined') }}  
9 )  
10  
11 *  
12  
13 final as (  
14     select  
15         *  
16         from orders  
17         left join customers using (customer_id)  
18     )  
19  
20 *  
21  
22 select * from final
```

Below the code editor, a lineage graph is visible, showing the flow of data from source tables to target metrics. The graph includes nodes for 'stg\_orders', 'stg\_customers', 'stg\_payments', 'int\_customer\_order\_history\_joined', 'int\_order\_payments\_pivoted', 'orders', 'average\_order\_amount', 'expenses', 'revenue', 'profit', and 'revenue'. The 'orders' model is highlighted in purple, and its lineage is shown in red.

Source



...but



...messy UI



Source



...extra features = extra money



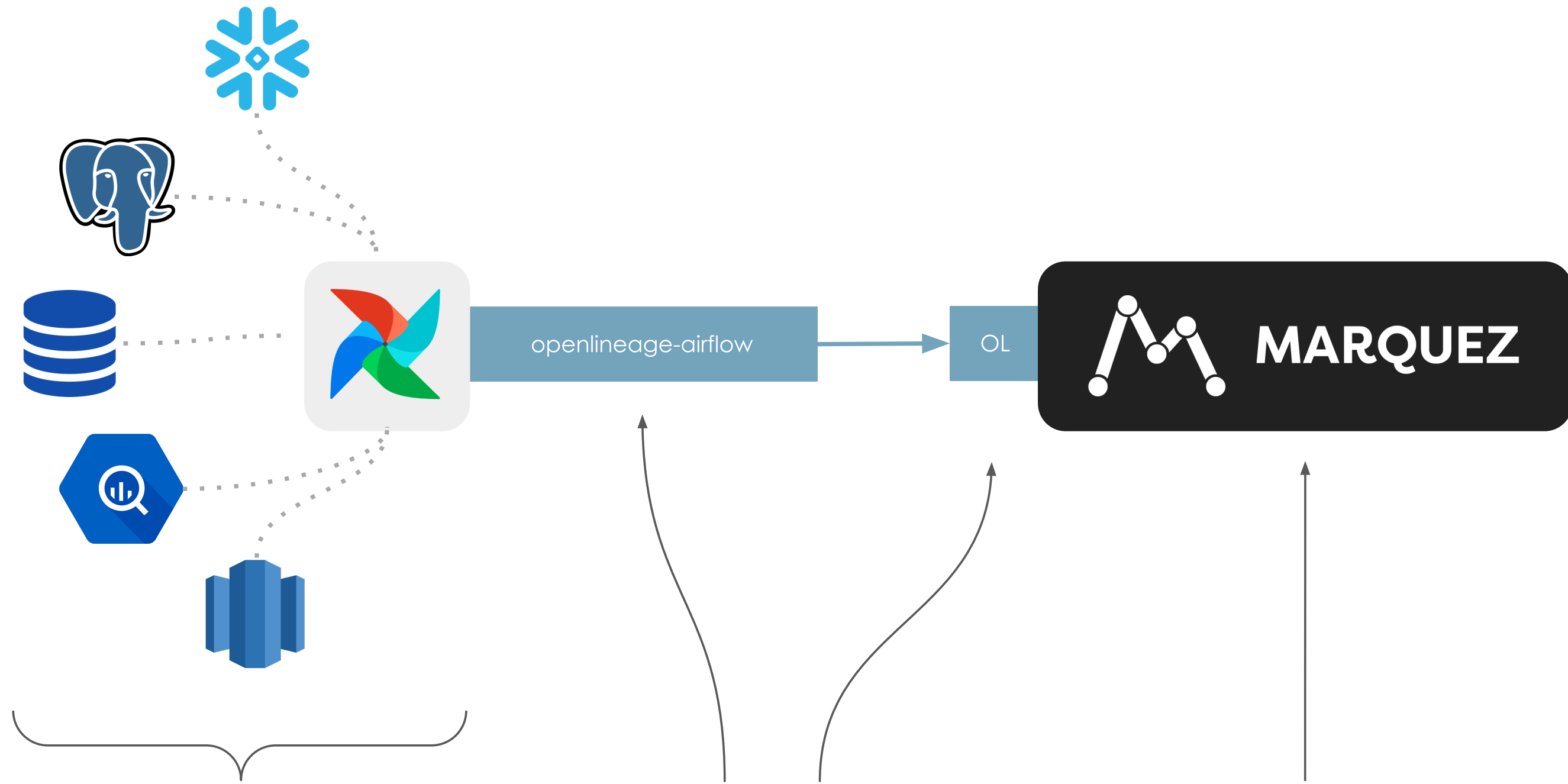
**Column  
level lineage  
in dbt core**



**Column-level  
lineage in dbt  
cloud (\$\$\$)**



Lineage Approaches:  
Airflow OpenLineage



Airflow executes a set of DAGs that operate on various data sources (e.g., Snowflake, Redshift, Bigquery, PostgreSQL)

The openlineage-airflow Python module automatically captures metadata and calls the OpenLineage Collection API

Metadata is stored in an OpenLineage-compatible repository like Marquez, where it can be analyzed



# ...but

- Not all operators have OpenLineage integrations
- Custom Operators require manual OL implementation
- SQL Parsing doesn't work well on custom cases
- Column level lineage is not ready yet
- Dynamic Tasks + SQL



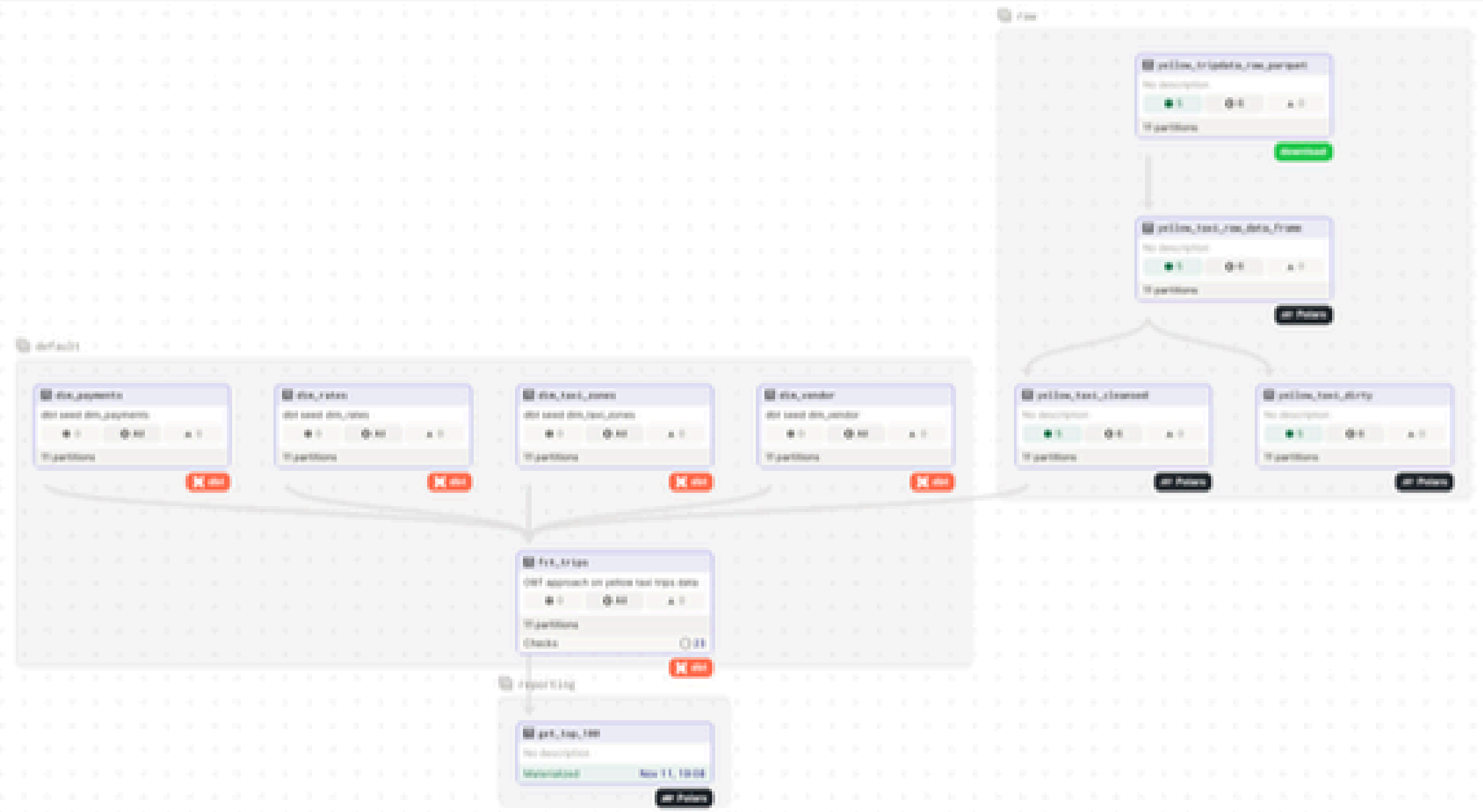
# Lineage Approaches: Dagster

- yellow\_taxi
- Jobs
  - all\_assets\_job
- Asset groups
  - dbt\_seeds
  - raw
  - reporting
  - yellow\_taxi\_model
- Resources
  - dbt
  - file\_io\_manager
  - s3

### Global Asset Lineage

Reload definitions

Filter Type an asset subset... (ex: yellow\_taxi\_raw/yellow\_taxi\_raw\_data\_frame+) Clear query Materialize all...



...but

- Relatively small community
- Requires mindset shift
- and...



...extra features = extra money



Column  
level lineage  
in dagster

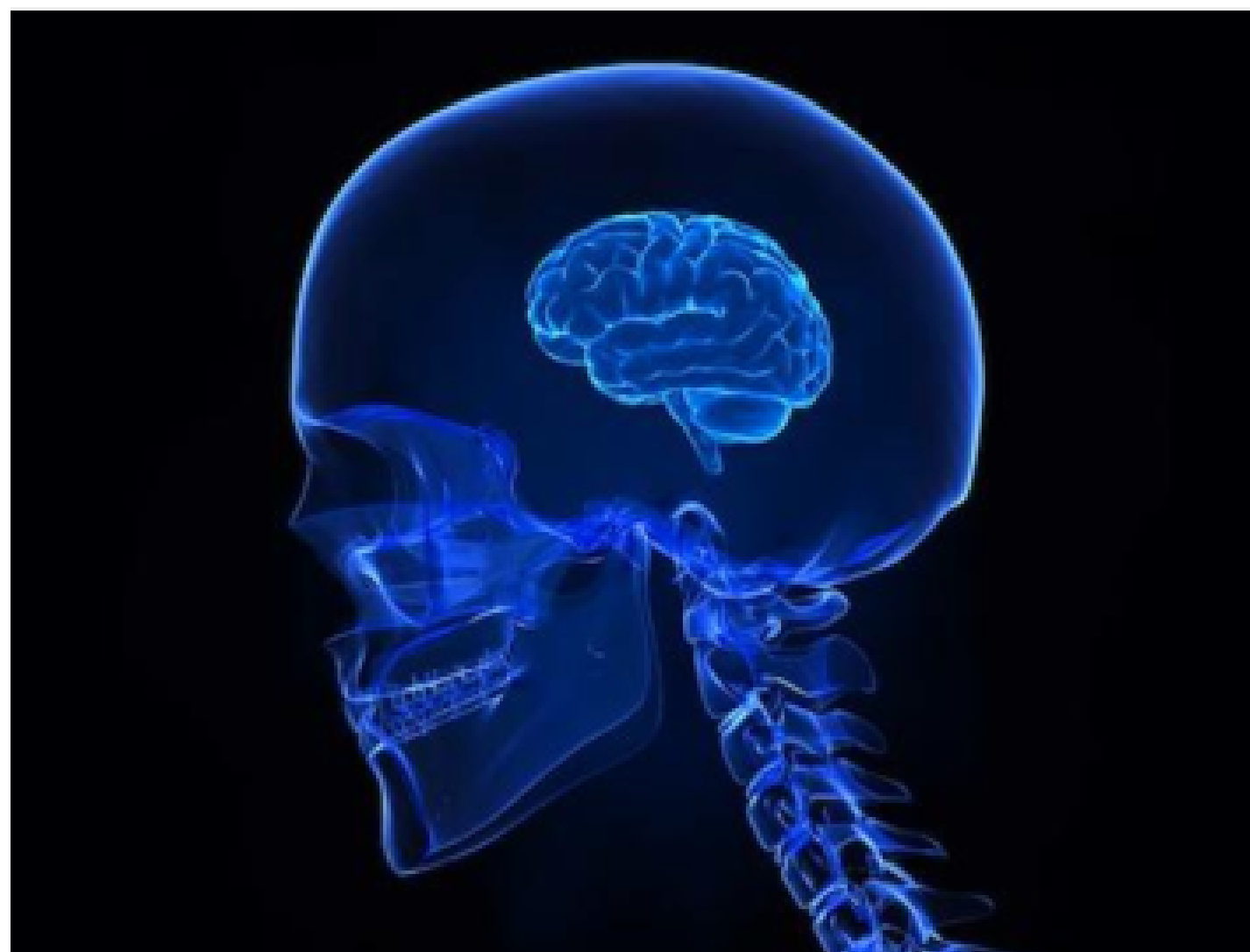


Column-level  
lineage  
in dagster  
cloud (\$\$\$)



I'm brave ~~stupid~~  
enough to DIY





Let's start simple



```
from sqlglot import parse_one
sql = """
CREATE TABLE mart.customer_orders AS
SELECT
    c.customer_name,
    c.email,
    o.order_id,
    o.order_date,
    o.amount,
    o.price,
    o.amount * o.price AS total_price
FROM staging.orders o
JOIN dim.customers c ON o.customer_id = c.customer_id"""

parse_one(sql)
```



```

Create(
  this=Table(
    this=Identifier(this=customer_orders, quoted=False),
    db=Identifier(this=mart, quoted=False)),
  kind=TABLE,
  expression=Select(
    expressions=[
      Column(
        this=Identifier(this=customer_name, quoted=False),
        table=Identifier(this=c, quoted=False)),
      Column(
        this=Identifier(this=email, quoted=False),
        table=Identifier(this=c, quoted=False)),
      Column(
        this=Identifier(this=order_id, quoted=False),
        table=Identifier(this=o, quoted=False)),
      Column(
        this=Identifier(this=order_date, quoted=False),
        table=Identifier(this=o, quoted=False)),
      Column(
        this=Identifier(this=amount, quoted=False),
        table=Identifier(this=o, quoted=False)),
      Column(
        this=Identifier(this=price, quoted=False),
        table=Identifier(this=o, quoted=False)),
      Alias(
        this=Mul(
          this=Column(
            this=Identifier(this=amount, quoted=False),
            table=Identifier(this=o, quoted=False)),
          expression=Column(
            this=Identifier(this=price, quoted=False),
            table=Identifier(this=o, quoted=False))),
        alias=Identifier(this=total_price, quoted=False))),
    from=From(
      this=Table(
        this=Identifier(this=orders, quoted=False),
        db=Identifier(this=staging, quoted=False),
        alias=TableAlias(
          this=Identifier(this=o, quoted=False))),
      joins=[
        Join(
          this=Table(
            this=Identifier(this=customers, quoted=False),
            db=Identifier(this=dim, quoted=False),
            alias=TableAlias(
              this=Identifier(this=c, quoted=False))),
          on=EQ(
            this=Column(
              this=Identifier(this=customer_id, quoted=False),
              table=Identifier(this=o, quoted=False)),
            expression=Column(
              this=Identifier(this=customer_id, quoted=False),
              table=Identifier(this=c, quoted=False))))))

```





Let's increase complexity a bit



```
-----  
sql_queries = {  
    "staging": """  
    CREATE TABLE staging.orders AS  
    SELECT order_id, customer_id, order_date, amount, price  
    FROM raw_data.orders  
    """,  
  
    "mart": """  
    CREATE TABLE mart.customer_orders AS  
    SELECT  
        c.customer_name,  
        c.email,  
        o.order_id,  
        o.order_date,  
        o.amount,  
        o.price,  
        o.amount * o.price AS total_price  
    FROM staging.orders o  
    JOIN dim.customers c ON o.customer_id = c.customer_id  
    """  
}
```





```
Analyzing staging query:  
Source tables: {'raw_data.orders'}  
Target tables: {'staging.orders'}
```

```
Column lineage:  
order_id <- ['order_id']  
customer_id <- ['customer_id']  
order_date <- ['order_date']  
amount <- ['amount']  
price <- ['price']
```

```
Analyzing mart query:  
Source tables: {'dim.customers AS c', 'staging.orders AS o'}  
Target tables: {'mart.customer_orders'}
```

```
Column lineage:  
customer_name <- ['c.customer_name']  
email <- ['c.email']  
order_id <- ['o.order_id']  
order_date <- ['o.order_date']  
amount <- ['o.amount']  
price <- ['o.price']  
total_price <- ['o.amount', 'o.price']
```





more...



```

sql_queries = {
  "staging": """
CREATE TABLE staging.enriched_orders AS
WITH daily_rates AS (
  SELECT date, currency, exchange_rate
  FROM raw_data.exchange_rates
  WHERE currency = 'USD'
),
order_metrics AS (
  SELECT
    o.order_id,
    o.customer_id,
    o.order_date,
    o.amount * COALESCE(r.exchange_rate, 1.0) as usd_amount,
    COUNT(*) OVER (PARTITION BY o.customer_id) as customer_order_count
  FROM raw_data.orders o
  LEFT JOIN daily_rates r ON o.order_date = r.date
)
SELECT
  om.*,
  c.customer_segment,
  c.country,
  CASE
    WHEN usd_amount > 1000 AND customer_order_count > 5 THEN 'VIP'
    WHEN usd_amount > 500 THEN 'Premium'
    ELSE 'Standard'
  END as customer_tier
FROM order_metrics om
JOIN dim.customers c ON om.customer_id = c.customer_id
""",
  "mart": """
CREATE TABLE mart.customer_analytics AS
SELECT
  c.customer_segment,
  c.country,
  DATE_TRUNC('month', o.order_date) as month,
  COUNT(DISTINCT o.customer_id) as unique_customers,
  COUNT(*) as total_orders,
  SUM(o.usd_amount) as total_revenue,
  AVG(CASE WHEN o.customer_tier = 'VIP' THEN o.usd_amount ELSE 0 END) as avg_vip_order_value
FROM staging.enriched_orders o
JOIN dim.customers c ON o.customer_id = c.customer_id
GROUP BY 1, 2, 3
HAVING COUNT(*) > 10
""",
}

```



Analyzing staging query:  
Source tables: {'raw\_data.orders AS o', 'order\_metrics AS om', 'dim.customers AS c', 'raw\_data.exchange\_rates', 'daily\_rates AS r'}  
Target tables: {'staging.enriched\_orders'}

Column lineage:

Table: staging.enriched\_orders

Target column: \*  
Source columns: [('om', '\*')]

Target column: customer\_segment  
Source columns: [('c', 'customer\_segment')]

Target column: country  
Source columns: [('c', 'country')]

Target column: customer\_tier  
Source columns: [(None, 'usd\_amount'), (None, 'usd\_amount'), (None, 'customer\_order\_count')]

Analyzing mart query:  
Source tables: {'staging.enriched\_orders AS o', 'dim.customers AS c'}  
Target tables: {'mart.customer\_analytics'}

Column lineage:

Table: mart.customer\_analytics

Target column: customer\_segment  
Source columns: [('c', 'customer\_segment')]

Target column: country  
Source columns: [('c', 'country')]

Target column: month  
Source columns: [('o', 'order\_date')]

Target column: unique\_customers  
Source columns: [('o', 'customer\_id')]

Target column: total\_orders  
Source columns: []

Target column: total\_revenue  
Source columns: [('o', 'usd\_amount')]

Target column: avg\_vip\_order\_value  
Source columns: [('o', 'usd\_amount'), ('o', 'customer\_tier')]





MORE

```

sql_queries = {
  "staging": """
CREATE TABLE staging.enriched_orders AS
WITH daily_rates AS (
  SELECT date, currency, exchange_rate
  FROM raw_data.exchange_rates
  WHERE currency = 'USD'
),
order_metrics AS (
  SELECT
    o.order_id,
    o.customer_id,
    o.order_date,
    o.amount * COALESCE(r.exchange_rate, 1.0) as usd_amount,
    COUNT(*) OVER (PARTITION BY o.customer_id) as customer_order_count
  FROM raw_data.orders o
  LEFT JOIN daily_rates r ON o.order_date = r.date
)
SELECT
  om.*,
  c.customer_segment,
  c.country,
  CASE
    WHEN usd_amount > 1000 AND customer_order_count > 5 THEN 'VIP'
    WHEN usd_amount > 500 THEN 'Premium'
    ELSE 'Standard'
  END as customer_tier
FROM order_metrics om
JOIN dim.customers c ON om.customer_id = c.customer_id
""",

  "mart": """
CREATE TABLE mart.customer_analytics AS
SELECT
  c.customer_segment,
  c.country,
  DATE_TRUNC('month', o.order_date) as month,
  COUNT(DISTINCT o.customer_id) as unique_customers,
  COUNT(*) as total_orders,
  SUM(o.usd_amount) as total_revenue,
  AVG(CASE WHEN o.customer_tier = 'VIP' THEN o.usd_amount ELSE 0 END) as avg_vip_order_value
FROM staging.enriched_orders o
JOIN dim.customers c ON o.customer_id = c.customer_id
GROUP BY 1, 2, 3
HAVING COUNT(*) > 10
""",
}

```

```
Analyzing staging query:
Source tables: ("dim"."customers" AS "c", "raw_data"."exchange_rates" AS "exchange_rates", "order_metrics" AS "om", "raw_data"."orders" AS "o", "daily_rates" AS "r")
Target tables: ("staging"."enriched_orders")

Column lineage:

Table: "staging"."enriched_orders"

Target column: order_id
Source columns: (('om', 'order_id'))

Target column: customer_id
Source columns: (('om', 'customer_id'))

Target column: order_date
Source columns: (('om', 'order_date'))

Target column: usd_amount
Source columns: (('om', 'usd_amount'))

Target column: customer_order_count
Source columns: (('om', 'customer_order_count'))

Target column: customer_segment
Source columns: (('c', 'customer_segment'))

Target column: country
Source columns: (('c', 'country'))

Target column: customer_tier
Source columns: (('om', 'usd_amount'), ('om', 'usd_amount'), ('om', 'customer_order_count'))

Analyzing mart query:
Source tables: ("staging"."enriched_orders" AS "o", "dim"."customers" AS "c")
Target tables: ("mart"."customer_analytics")

Column lineage:

Table: "mart"."customer_analytics"

Target column: customer_segment
Source columns: (('c', 'customer_segment'))

Target column: country
Source columns: (('c', 'country'))

Target column: month
Source columns: (('o', 'order_date'))

Target column: unique_customers
Source columns: (('o', 'customer_id'))

Target column: total_orders
Source columns: ()

Target column: total_revenue
Source columns: (('o', 'usd_amount'))

Target column: avg_vip_order_value
Source columns: (('o', 'usd_amount'), ('o', 'customer_tier'))
```

A large, dark metal cauldron with a handle is filled with glowing golden coins. The background is a dark, fiery, and textured environment, possibly a cave or a forge, with a dragon's head visible on the left side. The overall lighting is warm and orange-red, creating a dramatic and intense atmosphere.

People who use SELECT \*  
in production

```
▶ sql = "SELECT orders.* FROM orders JOIN customers ON orders.customer_id = customers.id"  
print(parse_one(sql).sql(pretty=True))
```

```
⇨ SELECT  
  orders.*  
FROM orders  
JOIN customers  
  ON orders.customer_id = customers.id
```

```
▶ schema = {"orders": {"id": "INT", "order_date": "DATE", "amount": "FLOAT", "price": "FLOAT", "customer_id": "int"},
            "customers": {"id": "INT", "name": "varchar", "email": "varchar"},
            "exchange_rates": {"date": "DATE", "currency": "STRING", "exchange_rate": "FLOAT"}}
expression = parse_one(sql)
ast = qualify(expression, schema=schema, expand_stars=True).sql(pretty=True)
print(ast)
```

```
↳ SELECT
  "orders"."id" AS "id",
  "orders"."order_date" AS "order_date",
  "orders"."amount" AS "amount",
  "orders"."price" AS "price",
  "orders"."customer_id" AS "customer_id"
FROM "orders" AS "orders"
JOIN "customers" AS "customers"
  ON "orders"."customer_id" = "customers"."id"
```

What if I use dbt?

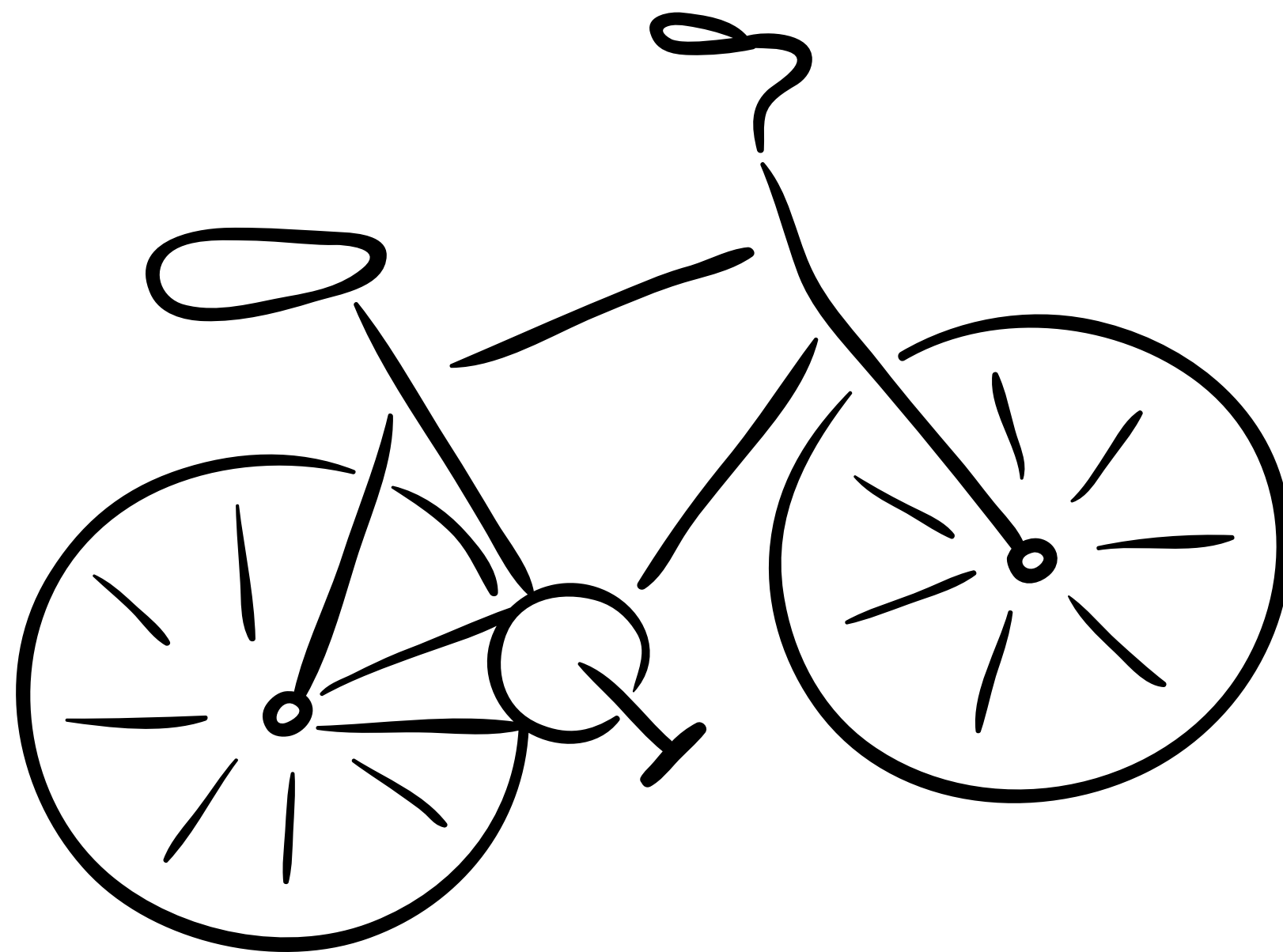
Compile your files

Run parsing over them

What's next?

Process Flow: Extract SQL → Parse → Push to Database → Add API on top

# Congratulations!



Q&A