

Data Lineage: Where 'It Depends' Finally Gets an Answer

2024-11-28



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>



Data Lineage Market Overview



Forrester Report



As of early 2024, estimates suggested the global data catalog market was valued around \$1.5-2 billion

Source: Claude.ai





Common Approaches



Just use dbt



The screenshot displays a code editor with a dark theme. On the left, there is a sidebar with 'Version Control' and 'File Explorer' sections. The 'File Explorer' shows a directory structure with folders like 'analysis', 'etc', 'logs', 'models', 'metrics', 'staging', and 'seeds'. The 'models' folder is expanded, showing files like 'docs.md', 'orders.sql', and 'orders.yml'. The main editor area shows the content of 'orders.sql', which is a SQL query using Common Table Expressions (CTEs). The query defines 'orders' as a CTE that selects from 'int_order_payments_pivoted', and 'customers' as a CTE that selects from 'int_customer_order_history_joined'. The final query selects from the 'orders' CTE, left joining with the 'customers' CTE on 'customer_id'. Below the code editor, there is a 'Lineage' tab showing a diagram of data dependencies. The diagram consists of nodes representing tables and metrics, connected by arrows indicating data flow. The nodes are: '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' node is highlighted in purple, and it is the source for 'average_order_amount', 'expenses', and 'revenue'. 'average_order_amount' and 'expenses' are highlighted in red, and 'revenue' is highlighted in blue. The 'revenue' node is also highlighted in blue, indicating a dependency on the 'revenue' node.

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

Source



but...



WHERE NON-SQL SOURCES?





Source





**Column
level lineage
in dbt core**

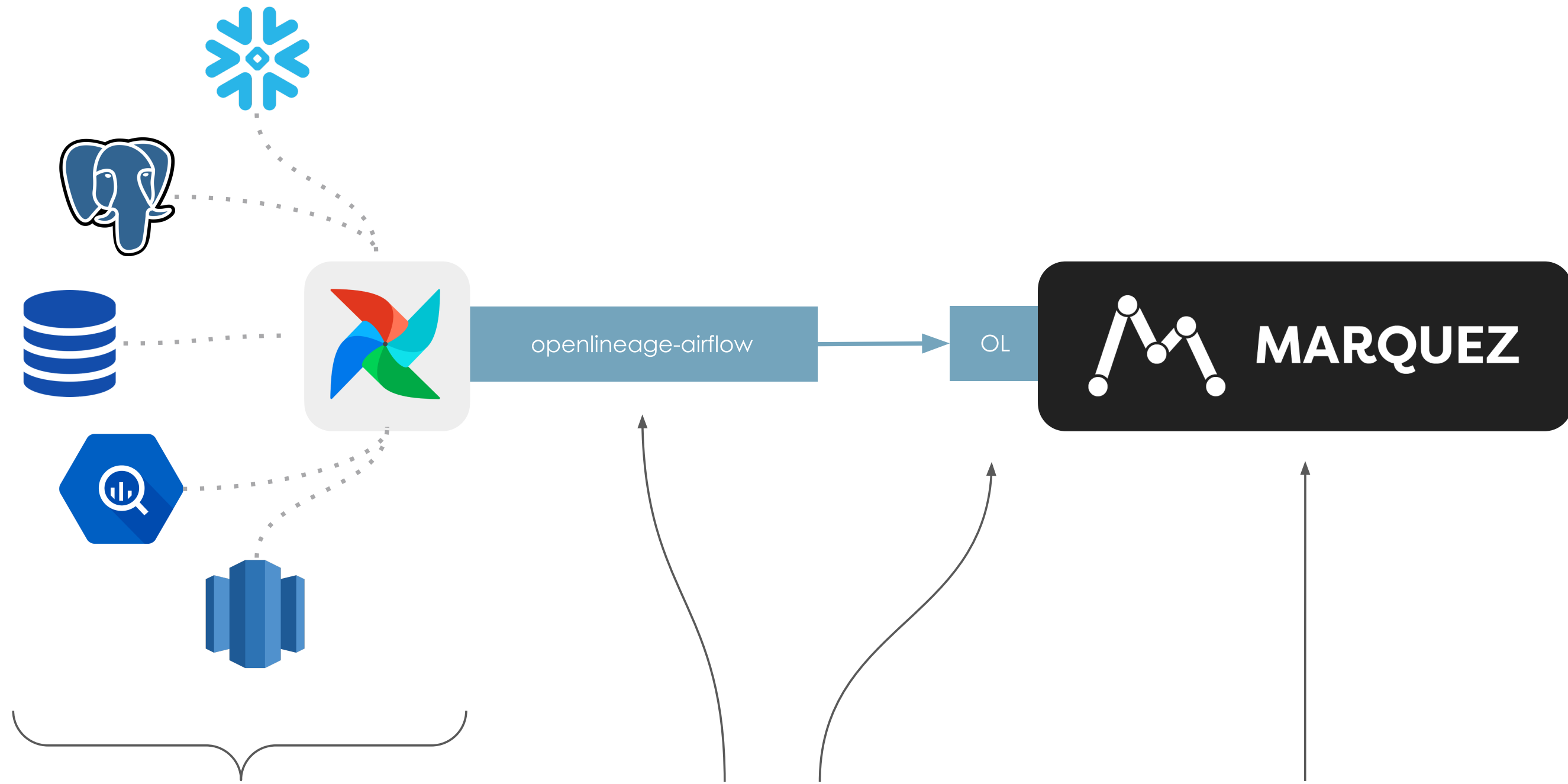


**Column-level
lineage in dbt
cloud (\$\$\$)**



Why not Airflow?





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**



Why not Dagster?

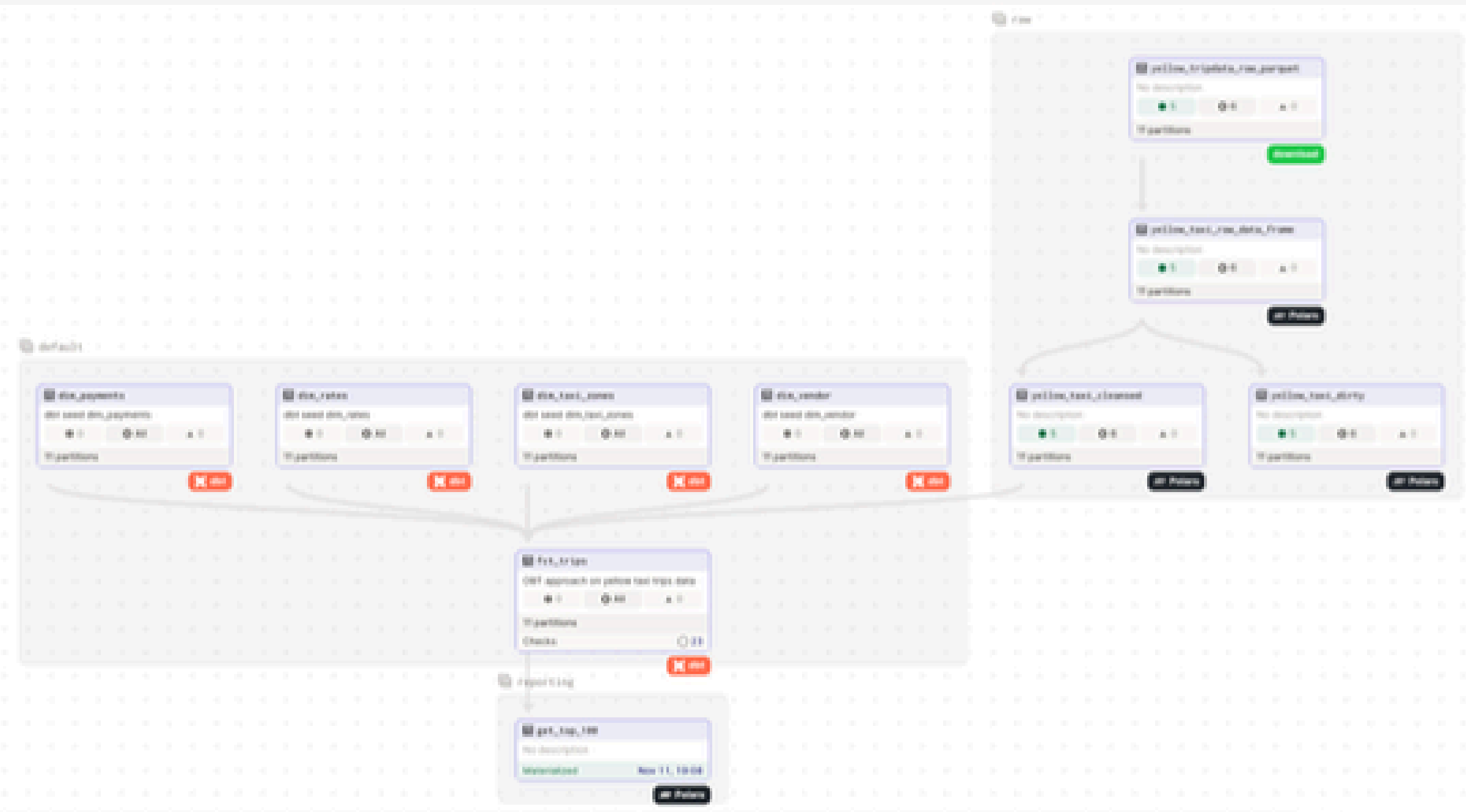


- yellow_taxi
- Jobs
- all_assets_job
- Asset groups
- dbt_seeds
- raw
- reporting
- yellow_taxi_model
- Resources
- dbt
- file_io_manager
- s3

Global Asset Lineage

Reload definitions

Filter Clear query Materialize all...



Search icon
 Zoom in icon
 Zoom out icon
 Refresh icon
 Close icon

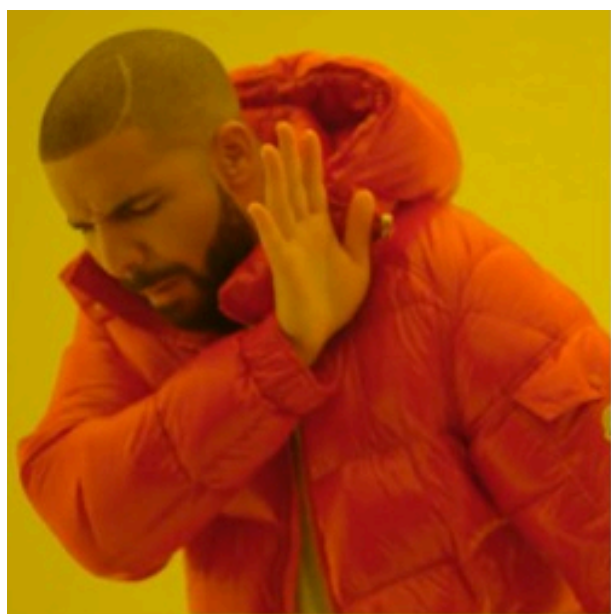


but...



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





Column
level lineage
in dagster

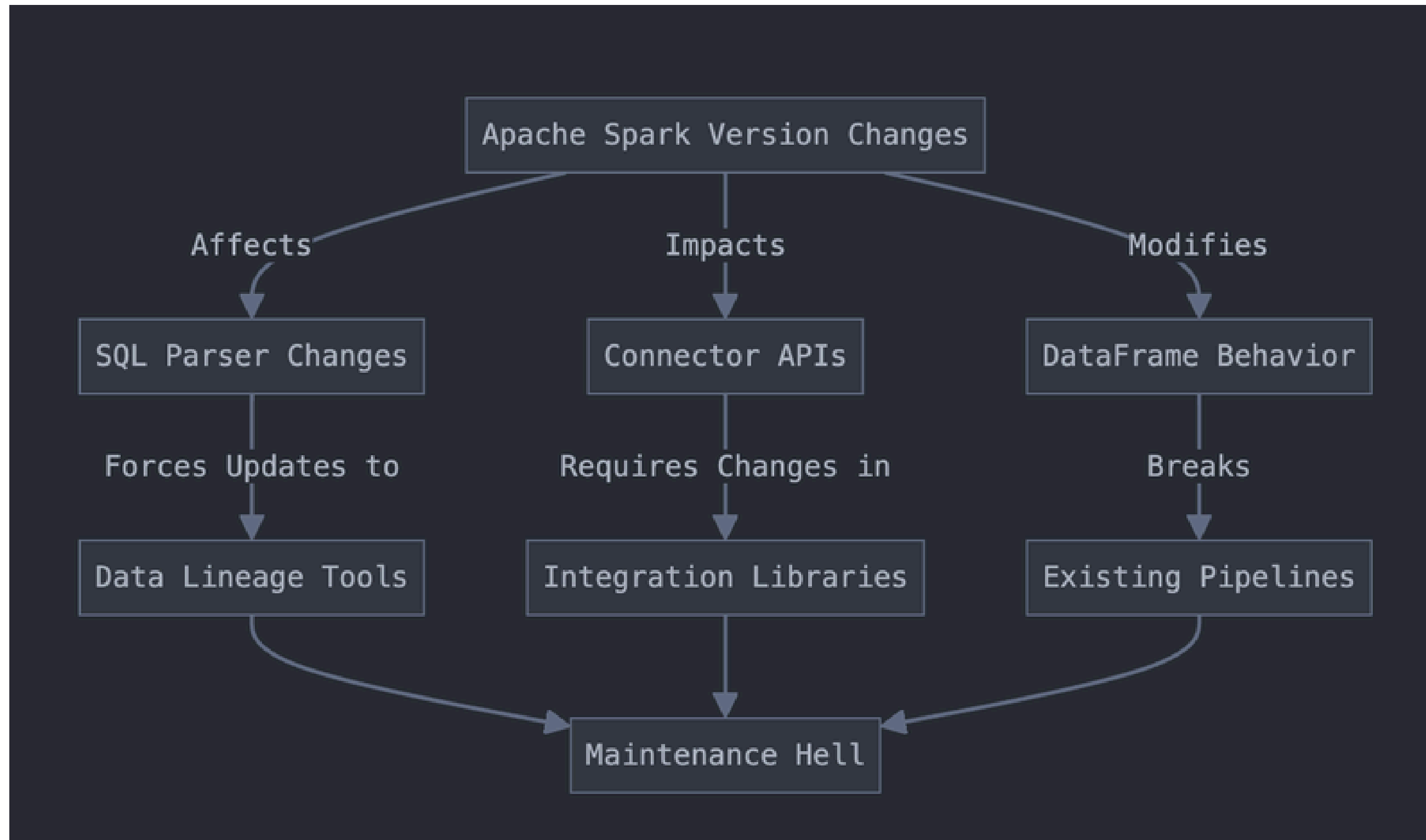


Column-level
lineage
in dagster
cloud (\$\$\$)



Elephant in the room





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

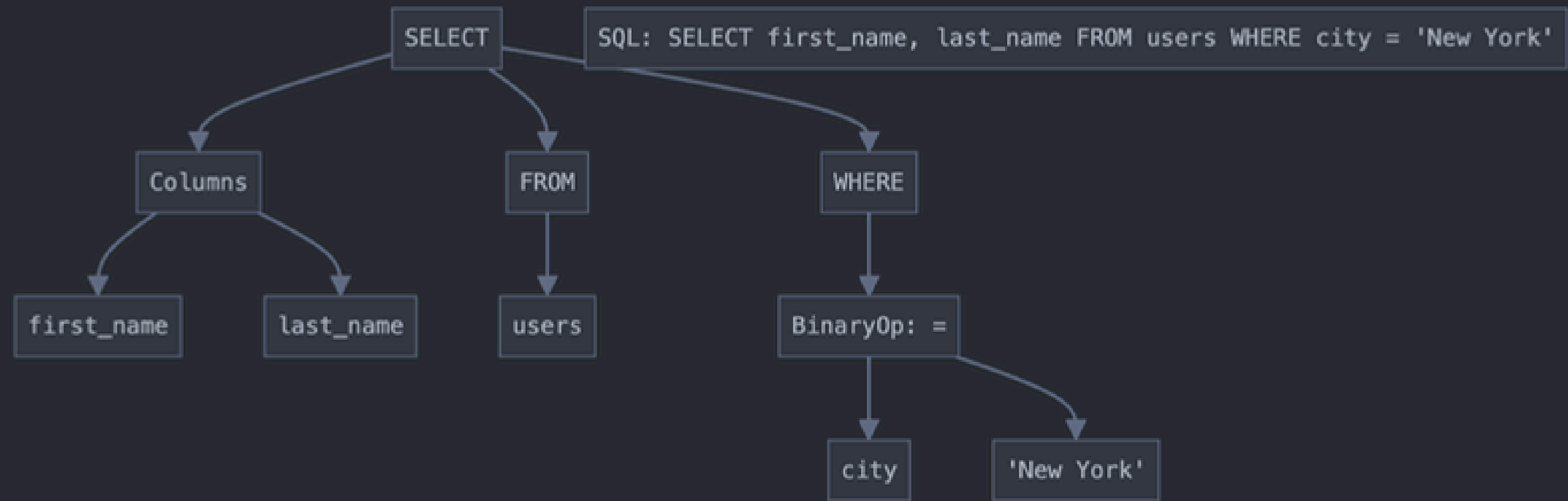


What's SQLGlot?



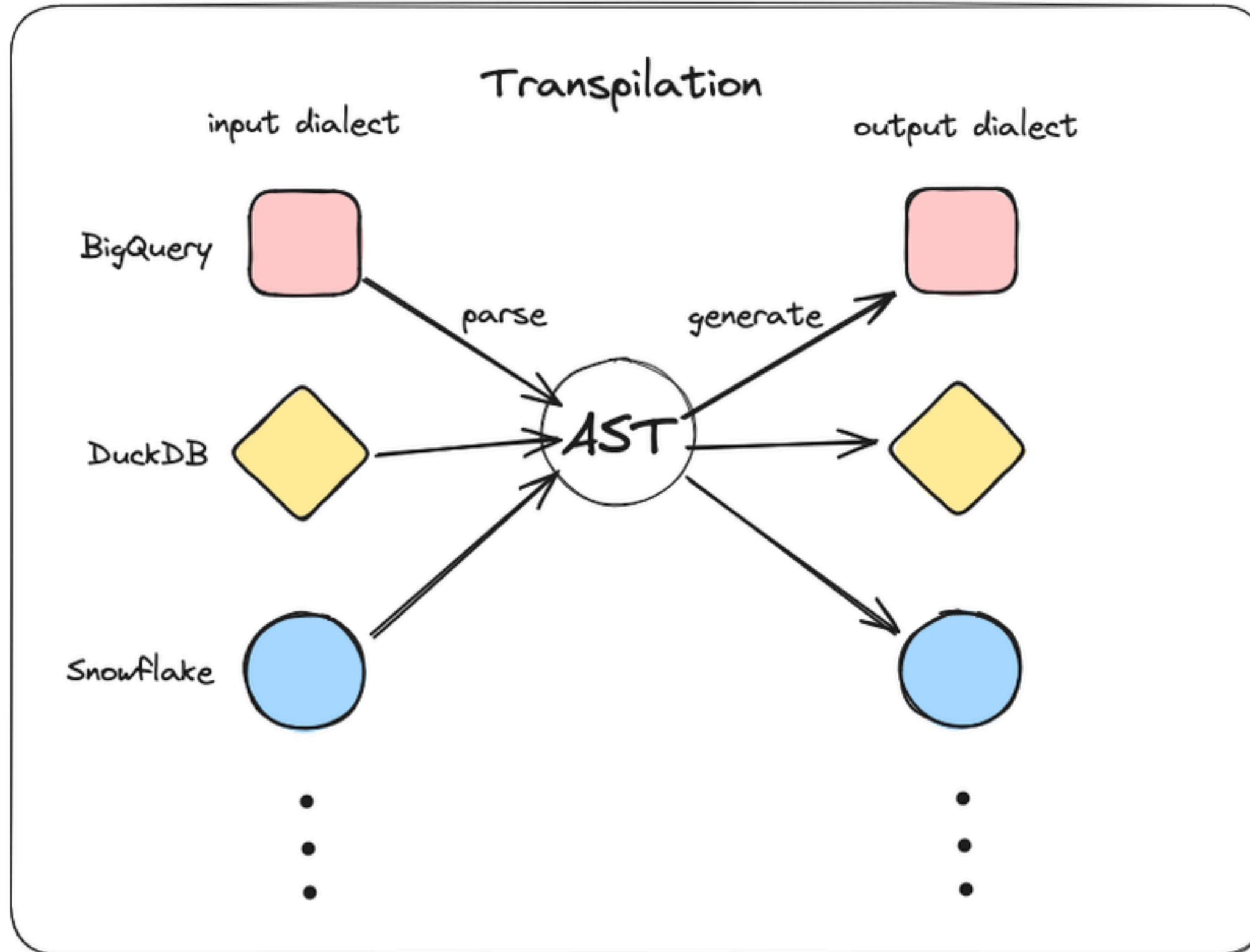
Abstract Syntax Tree (AST)





Transpilation





Source



Demo

Questions?

