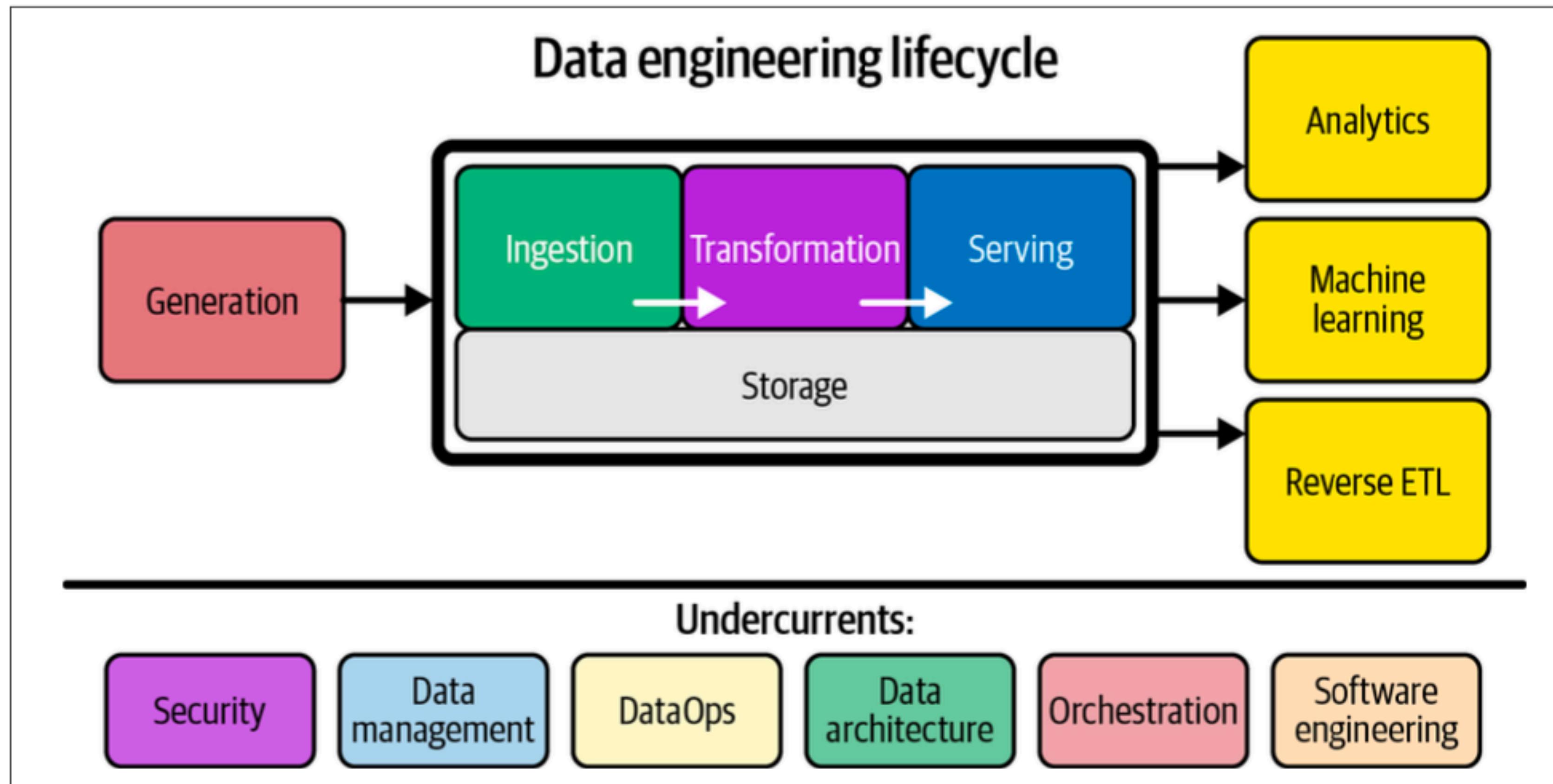


# Data Platforms and the Cloud

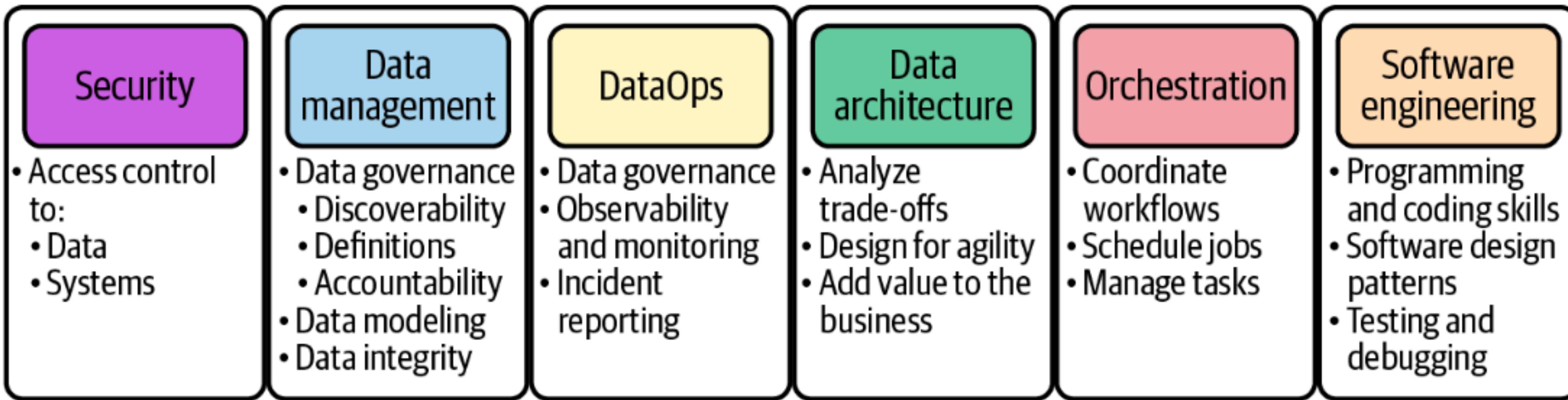
Rahul Dave(@rahuldave), Univ.Ai

# Data Engineering Lifecycle



*Figure 1-1. The data engineering lifecycle*

# Undercurrents

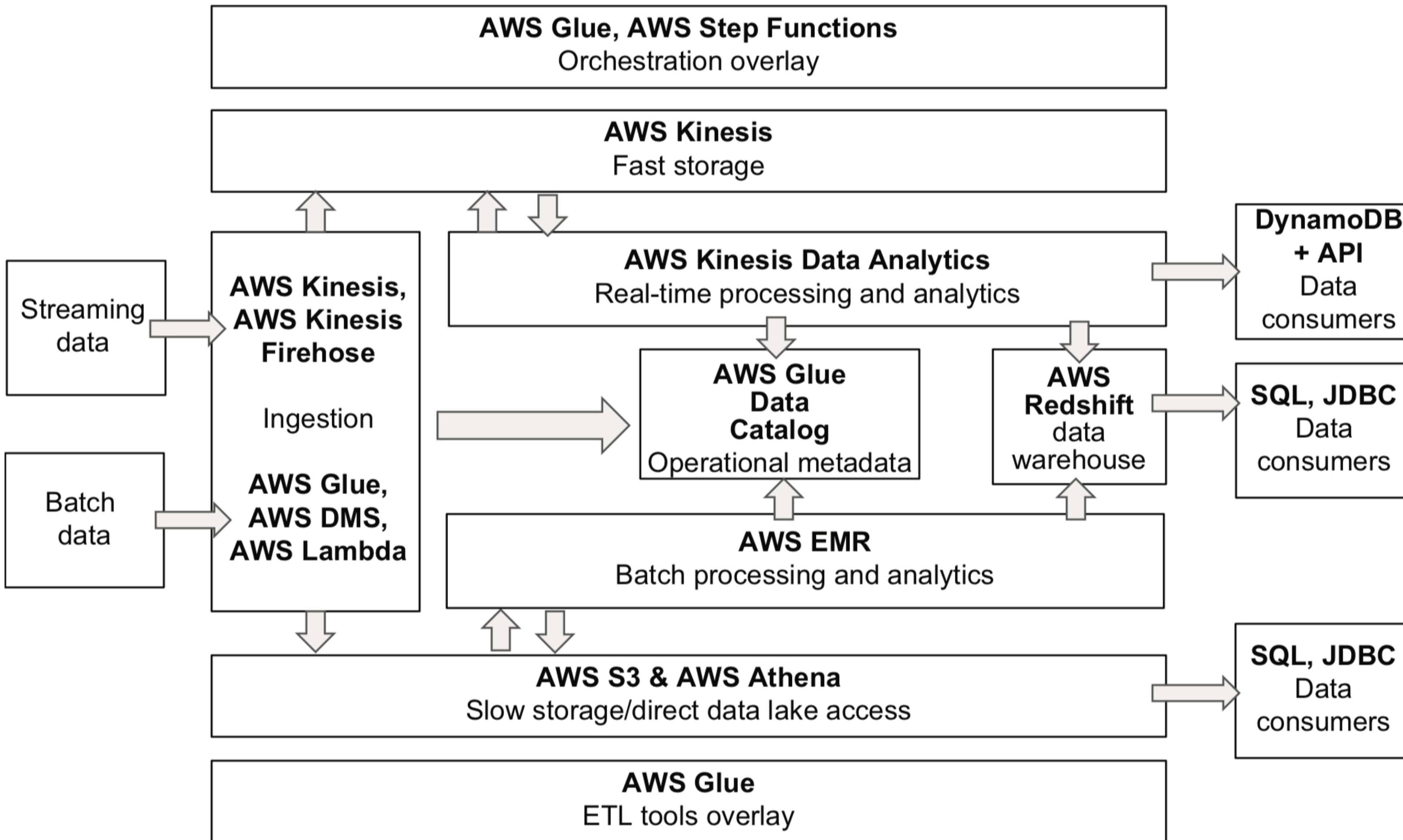


# What is left

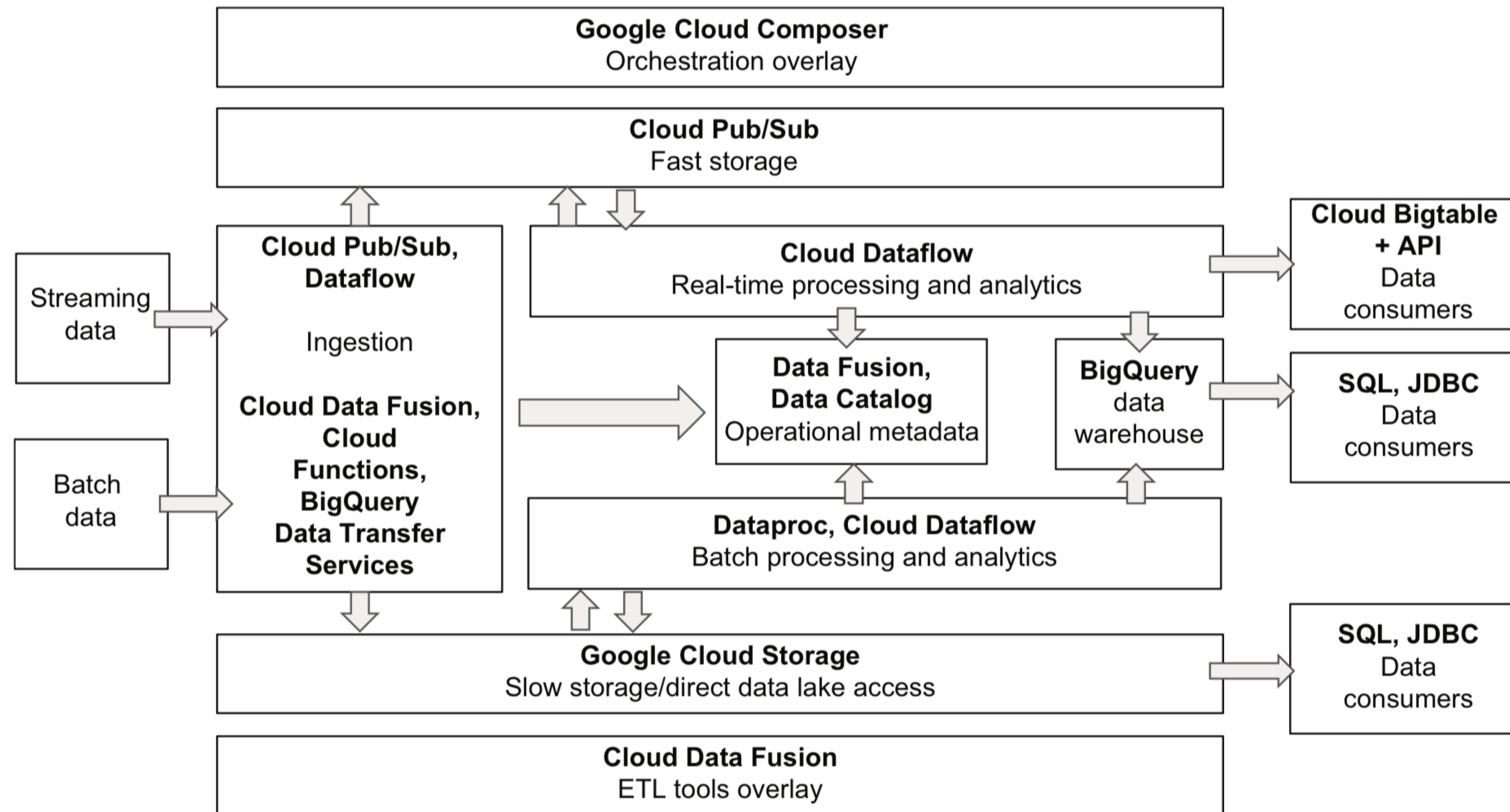
- We won't cover security and governance: these are huge topics in themselves
- we covered orchestration and data discoverability, and we'll cover it some more here
- we've talked about the metadata layer, and how we use it to track pipelines, discover data, and simplify data modeling (including schema on demand).
- we have covered data architecture in our discussion of the move from data lake and warehouse to data lakehouse, and really, the data platform: the combination of this all
- Now we'll focus on the stages of a data pipeline, and the best software development skills to organize and track data flow through them. We'll also focus on model integrity, testing.

# From Lakehouse to Platform

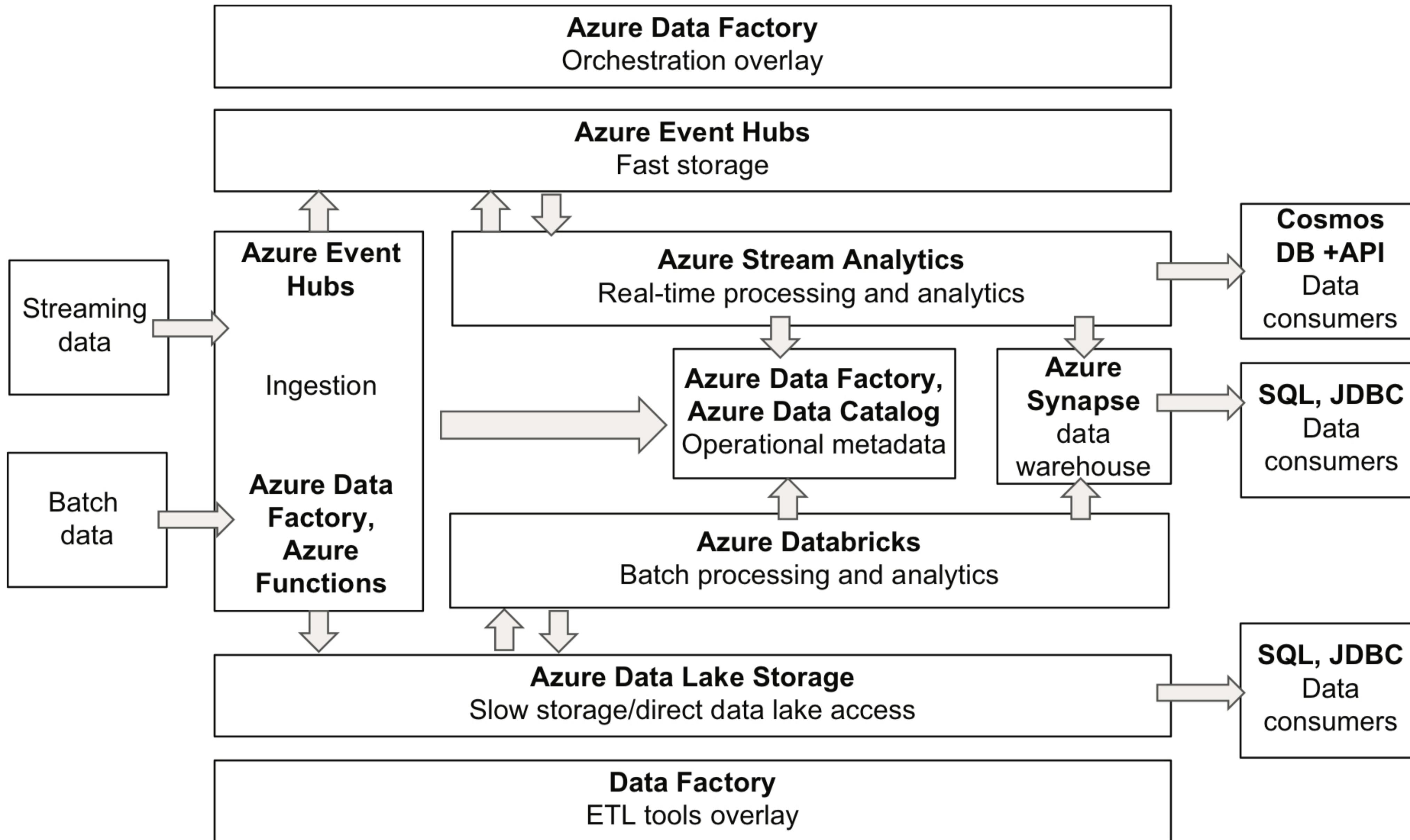
AWS



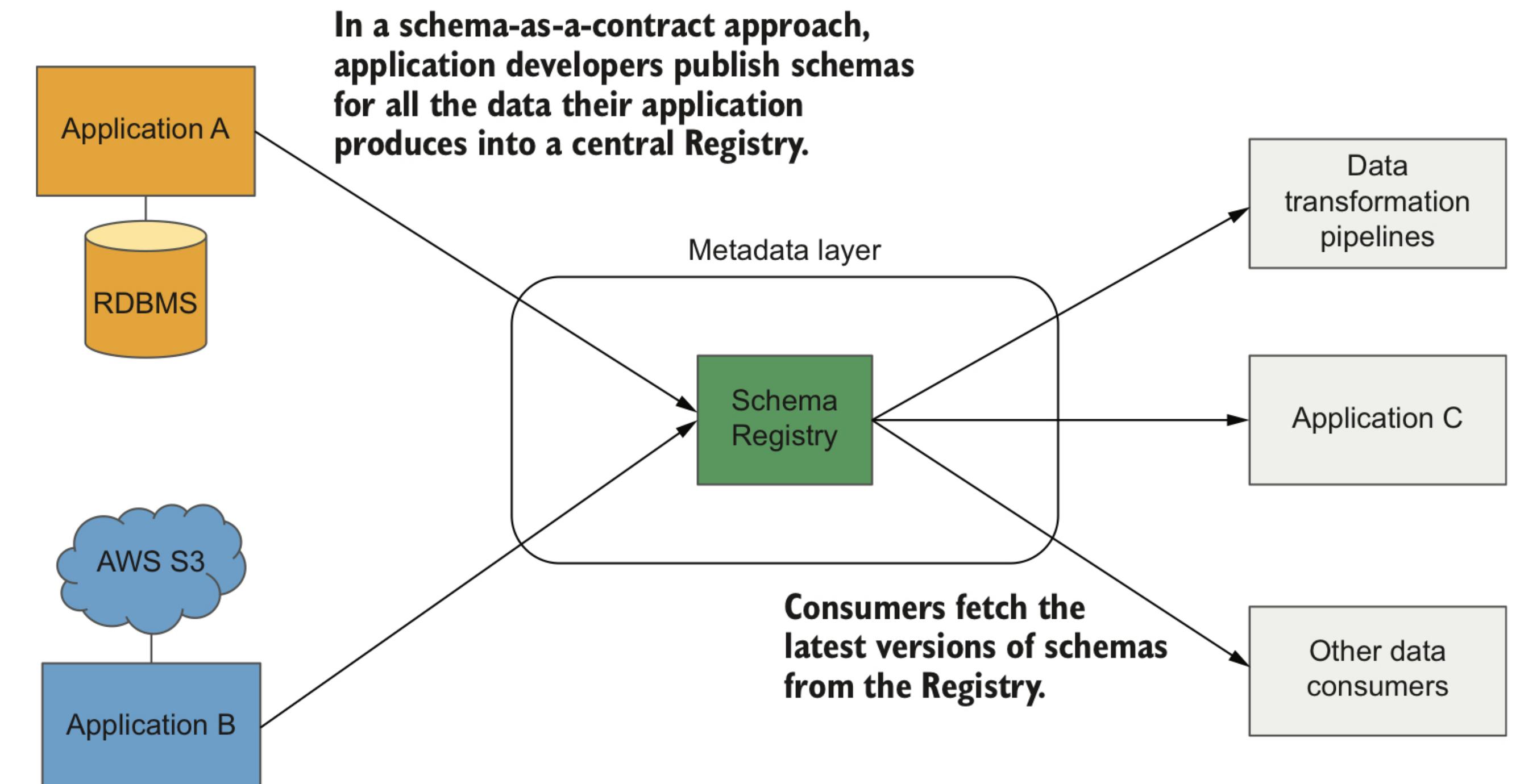
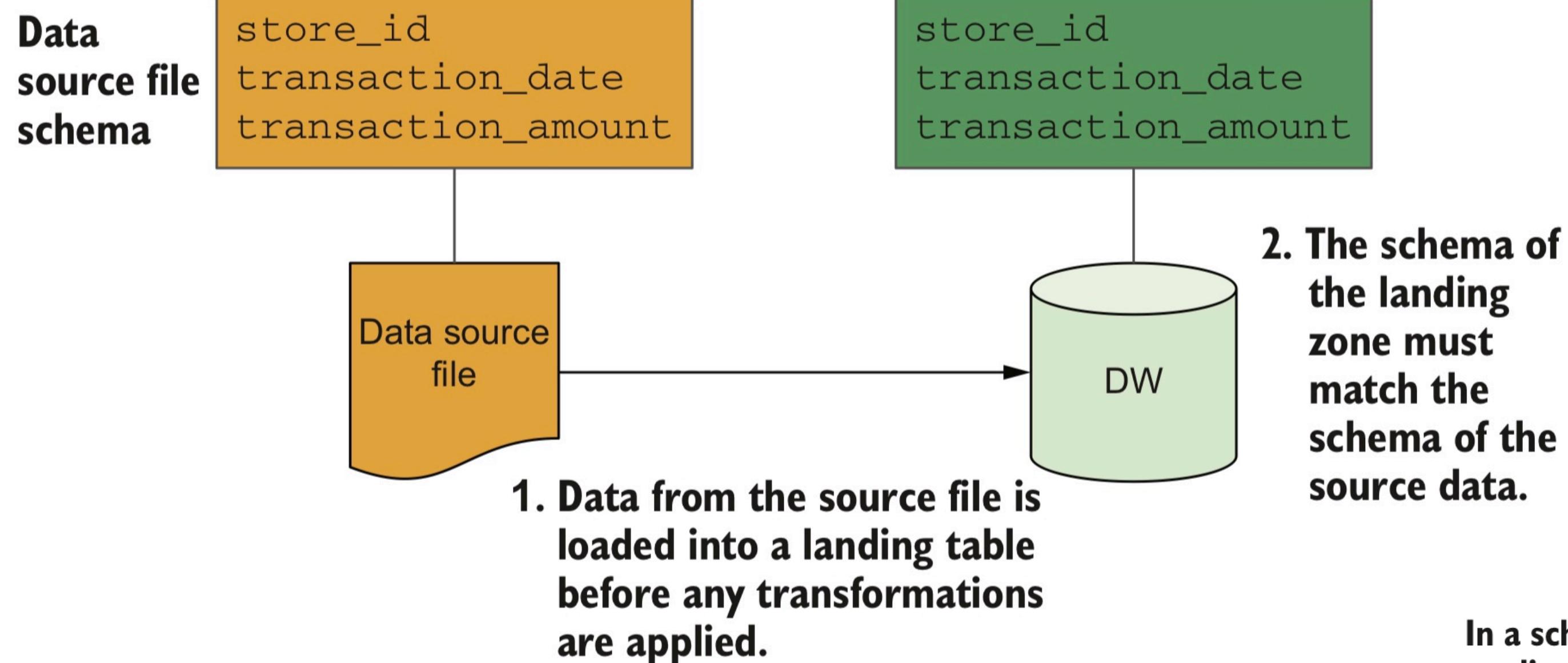
# Google Cloud

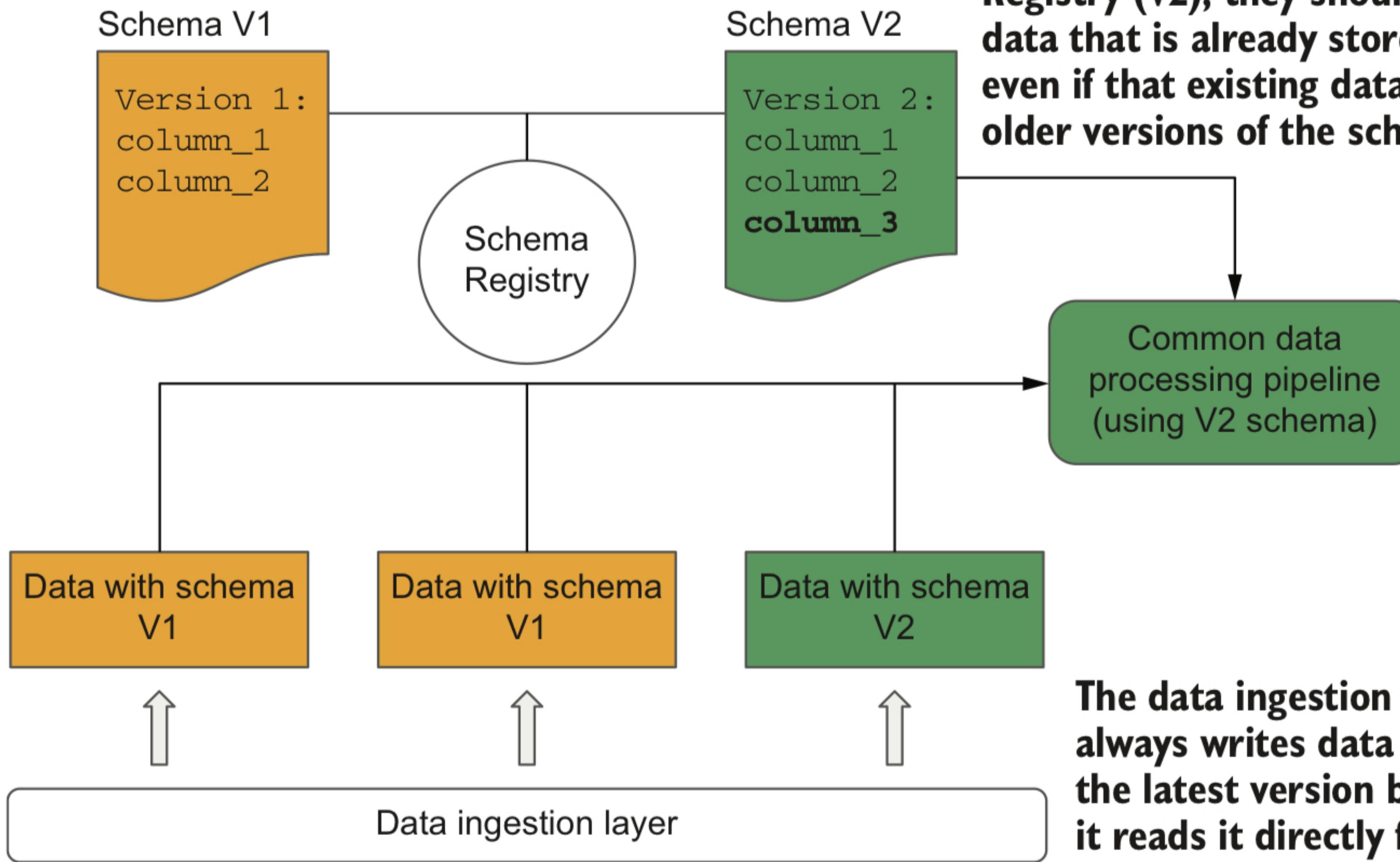


# Azure



# Schema Discovery

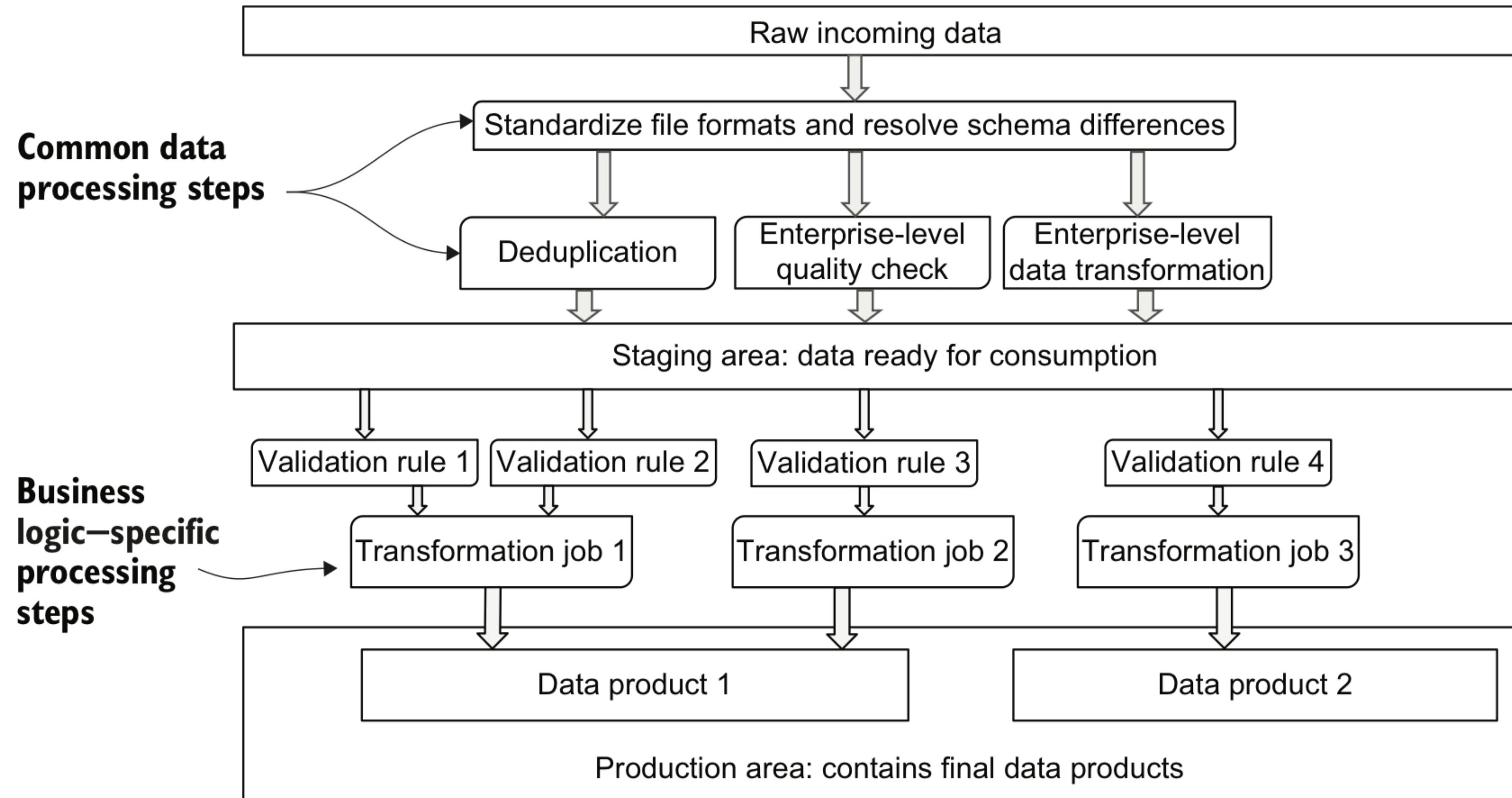




**Backward-compatible schema changes mean that if our data transformation pipelines use the latest version of the schema from the Registry (V2), they should be able to read all data that is already stored in the platform even if that existing data was written using older versions of the schema.**

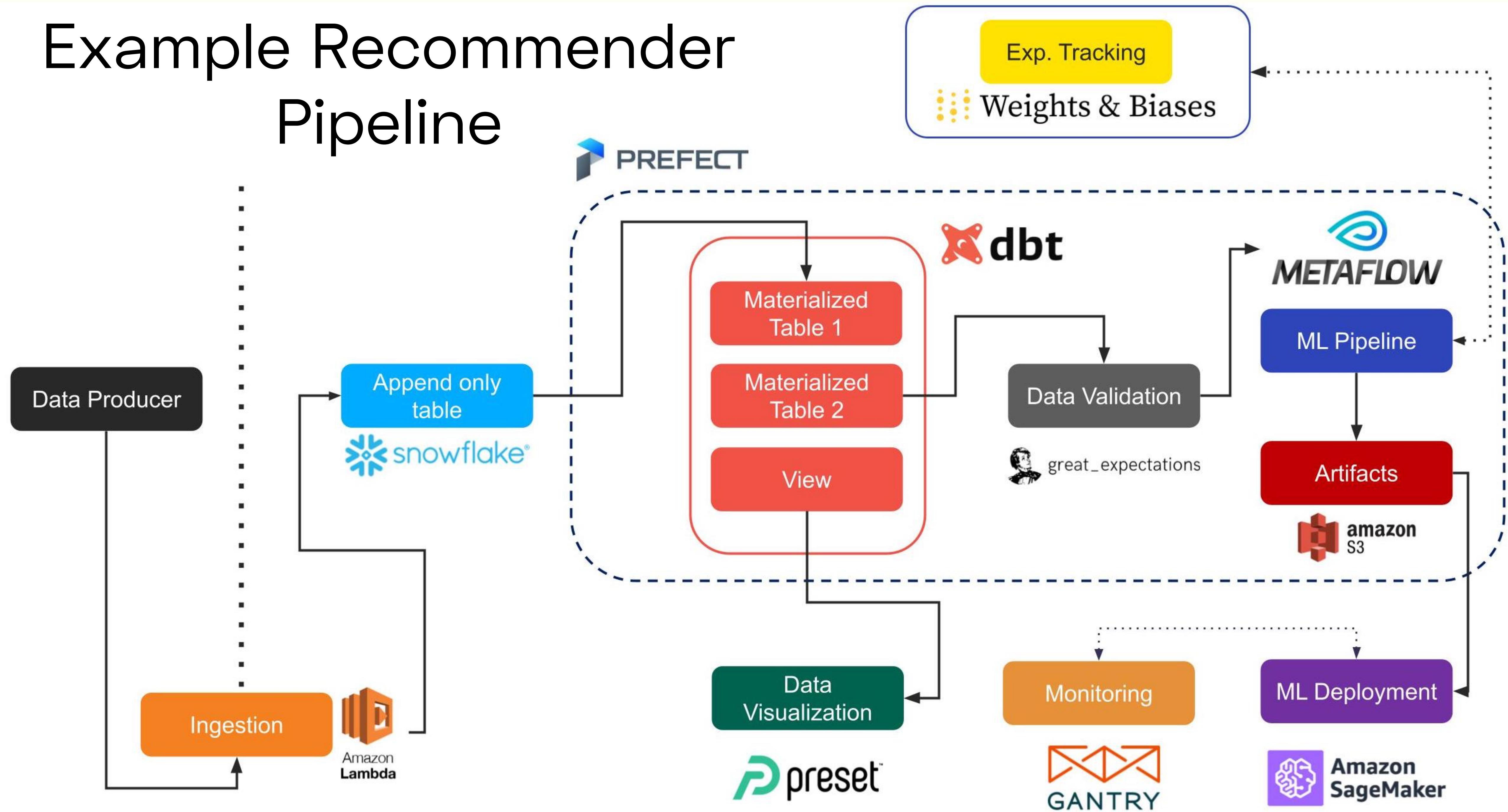
**The data ingestion layer always writes data with the latest version because it reads it directly from the data source.**

# How can we organize stuff like this?



- we need orchestration, which we have discussed
- we also need status reporting, which we have also discussed
- part of the status reporting can be data drift, and data statistics tests, which can lead to pipeline alerts
- we want to organize the code that's used to flow through this pipeline into reusable units, be able to see it all together, and to test it.

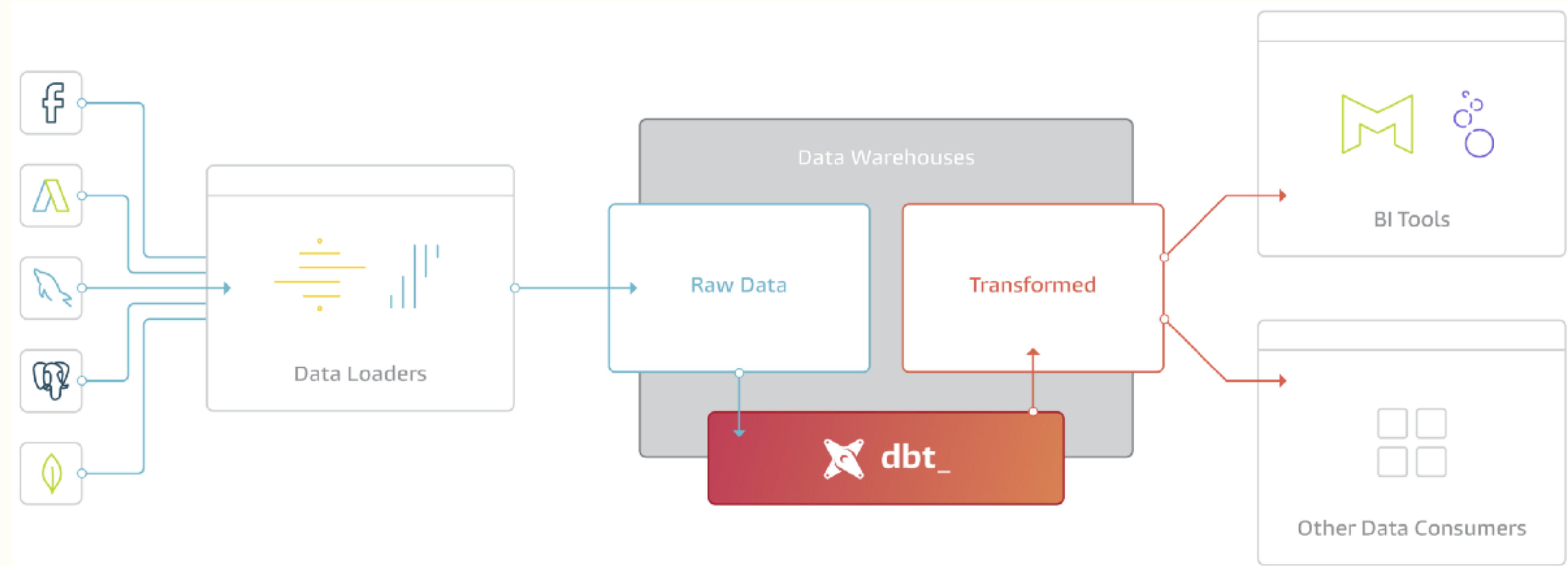
# Example Recommender Pipeline



# Orchestration, Code, and Transformation: What is dbt?

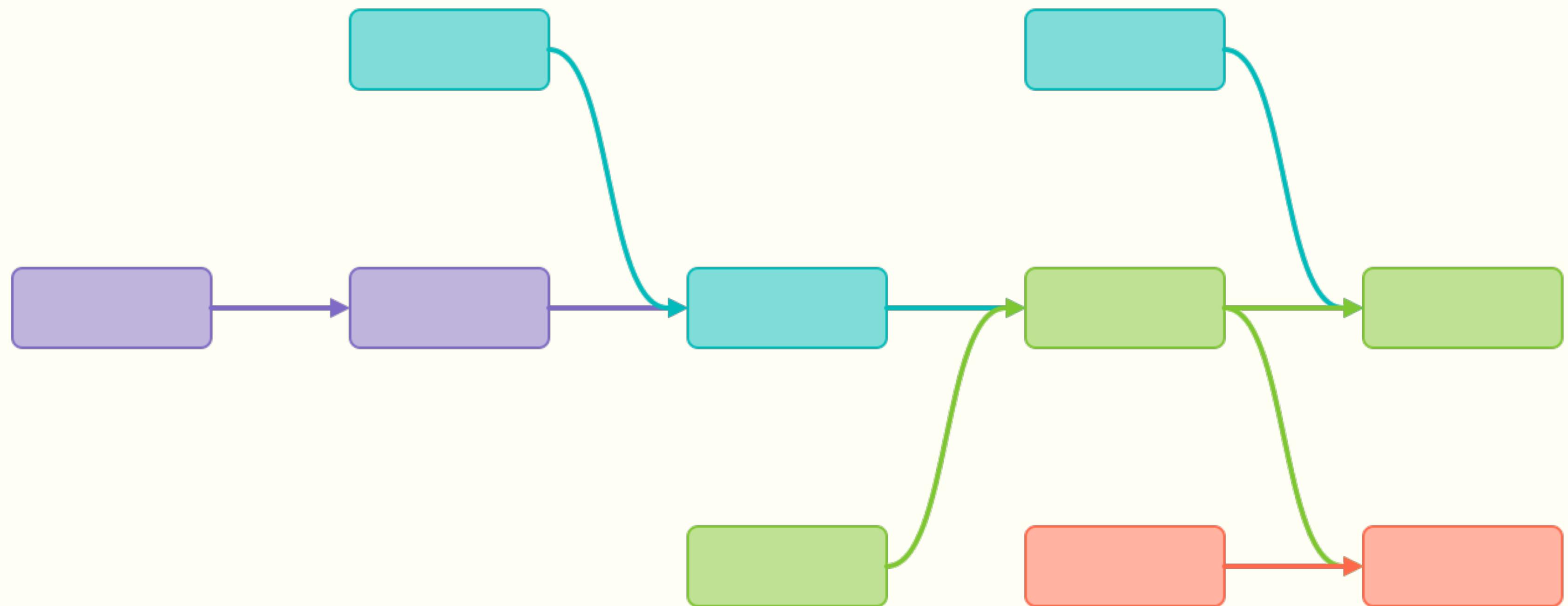
**dbt is a transformation workflow that helps you get more work done while producing higher quality results. You can use dbt to modularize and centralize your analytics code, while also providing your data team with guardrails typically found in software engineering workflows. Collaborate on data models, version them, and test and document your queries before safely deploying them to production, with monitoring and visibility.**

**dbt compiles and runs your analytics code against your data platform, enabling you and your team to collaborate on a single source of truth for metrics, insights, and business definitions. This single source of truth, combined with the ability to define tests for your data, reduces errors when logic changes, and alerts you when issues arise.**



# Its a DAG!

- This is a DAG for transformations only
- It will be embedded in a larger dag run by airflow or prefect or dagster or similar cloud software
- The lessons we learn from it apply everywhere though
- Such systematization is very new!



# Materialization

snapshots/orders.sql

```
{% snapshot orders_snapshot %}

{{ config(
    unique_key='id',
    strategy='timestamp',
    updated_at='updated_at',
)
}}
```

model\_with\_inline\_materialization.sql

```
{

    config(
        materialized='table',
        sort='timestamp',
        dist='user_id'
    )
}

select *
from ...
```

incremental\_model.sql

```
{

    config(
        materialized='incremental'
    )
}

select ...
```

# Tests

tests/assert\_payment\_amount\_is\_positive.sql

```
select
    order_id,
    sum(amount) as total_amount
from {{ ref('fct_payments') }}
group by 1
having not(total_amount >= 0)
```

schema.yml

```
version: 2
models:
  - name: events
    description: '{{ doc("table_events") }}'
    columns:
      - name: event_id
        description: This is a unique identifier for the event
    test:
      - unique
      - not_null
```

packages.yml

```
packages:
  - package: calogica/dbt_expectations
    version: 0.2.8
```

# Docs

exposures.yml

exposures:

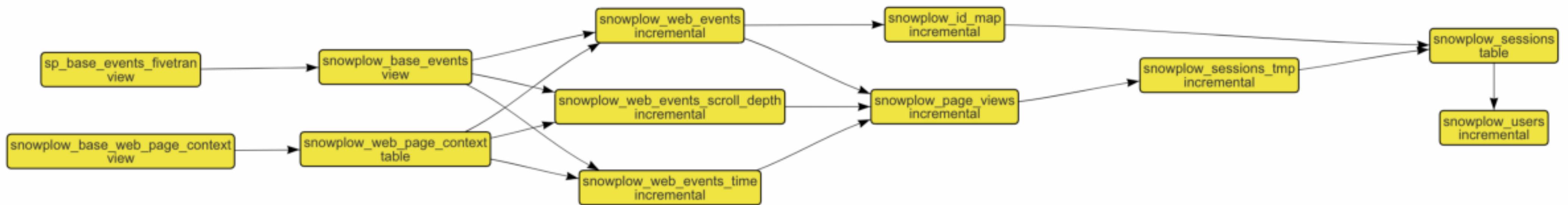
```
- name: weekly_jaffle_metrics
  type: dashboard
  url: https://bi.tool/dashboards/1
  depends_on:
    - ref('fct_orders')
    - ref('dim_customers')
  owner:
    name: Claire from Data
    email: data@jaffleshop.com
```

source\_freshness.yml

```
sources:
  - name: jaffle_shop
    database: raw
    freshness:
      warn_after: {count: 12, period: hour}
      error_after: {count: 24, period: hour}
      loaded_at_field: _etl_loaded_at
```

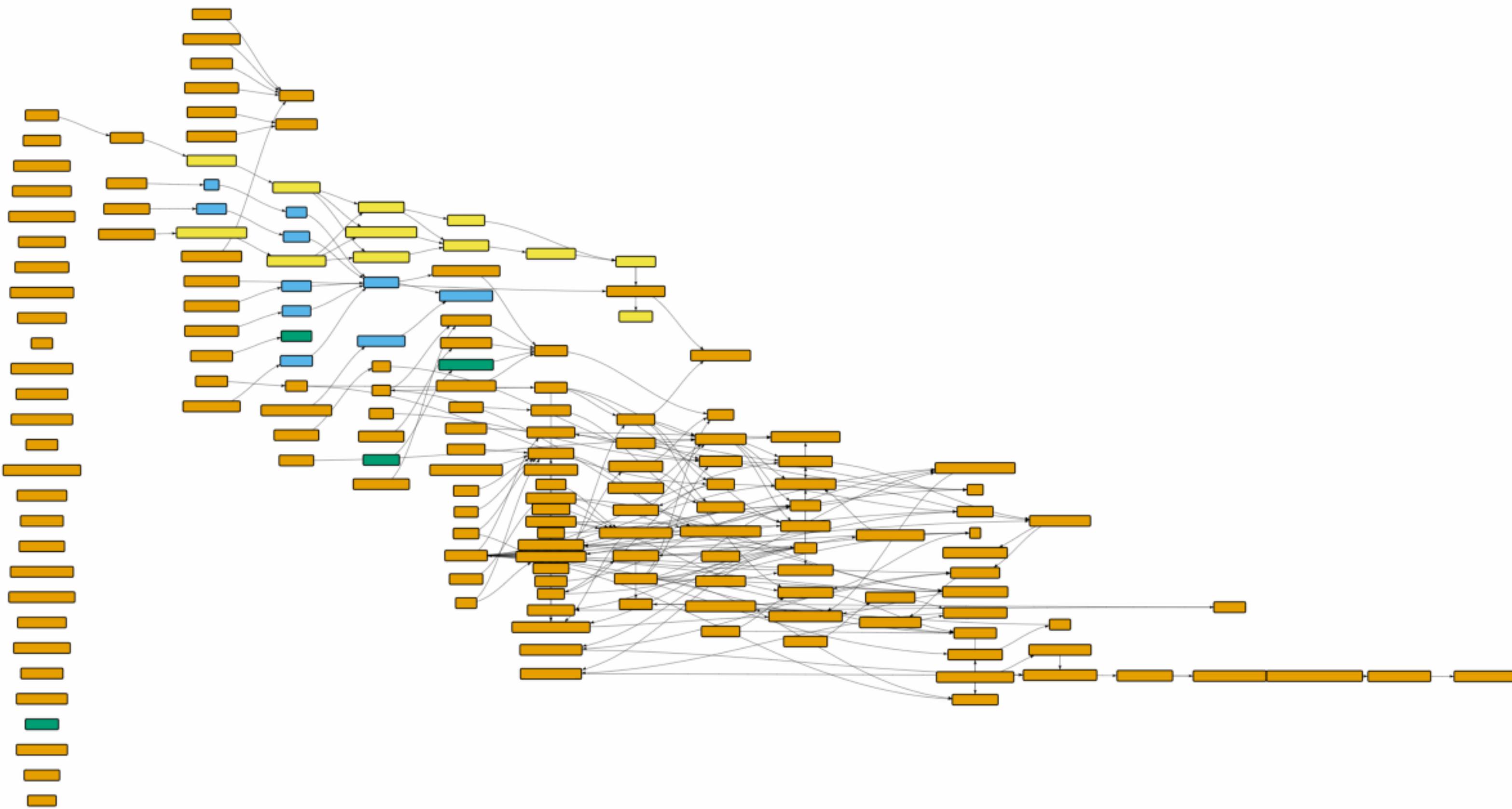
# Useful Features

- dbt constructs views/materialized views/tables automatically as models from sql files
- you can add integrity statements and tests
- using the ref macro, you can refer to a model abstractly so that the implementation is resolved at run time. This enables different development and production flows, even using different databases. Ditto the source macro
- For example you might test against duckdb locally, and move to spark/trino/dremio on s3 or iceberg or just bigquery, snowflake, or redshift.
- dbt operates “inside the warehouse, for the transformation layer”, but using a multi-sql tool like dremio or spark, you can use it everywhere. It is not the E or the L layer, but in this way it can be a multi database T layer.



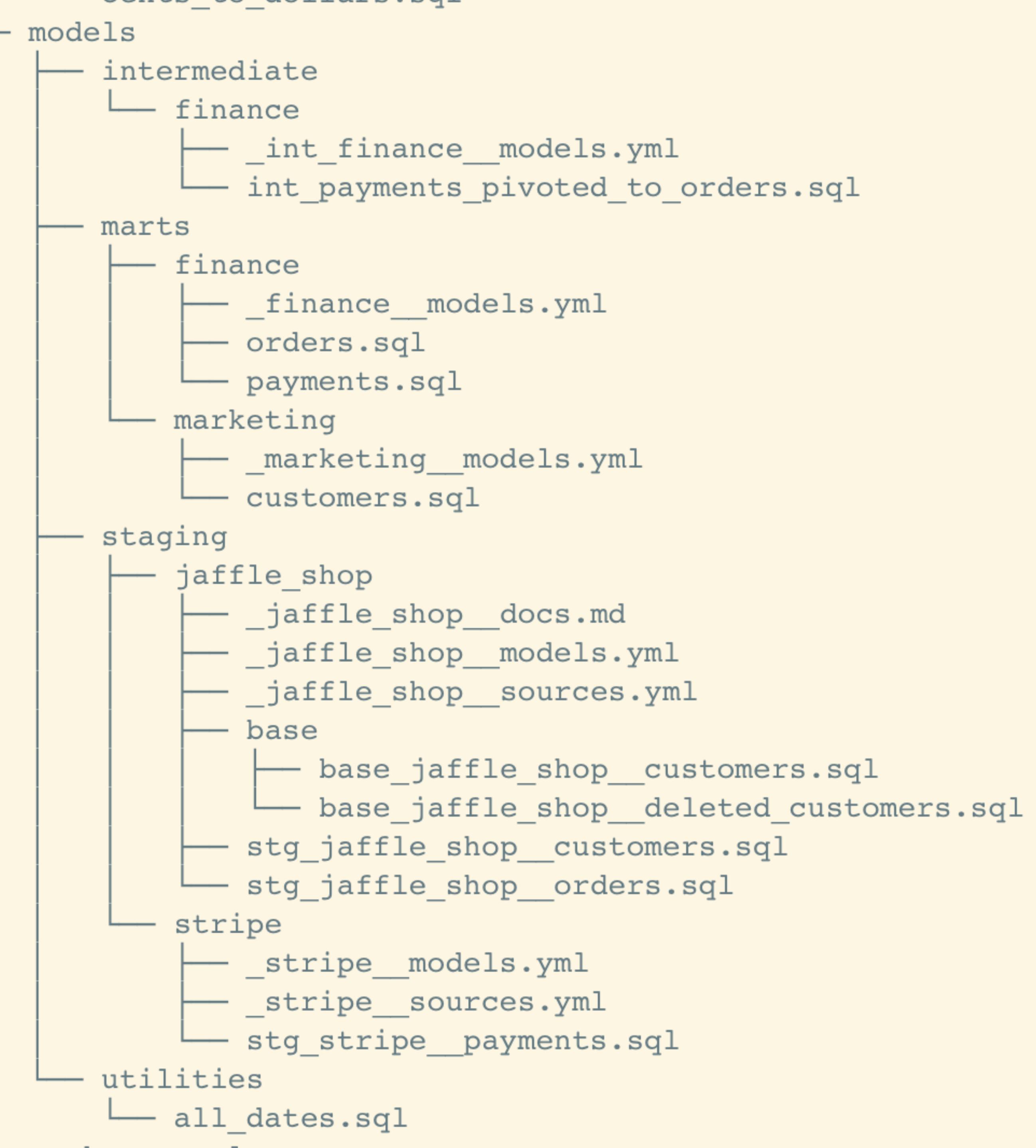
**Snowplow incrementally sessionized DAG**

**Ecommerce DAG**



# Example Layout

- We start with the staging, where we minimally change stuff coming in from other systems (over something like meltano/airbyte/spark/etc)
- We then make intermediate transformations
- finally groups can (or you can for them) make your own transformations in data marts: you will remember this architecture from before.
- inputs or outputs could be in warehouse, lake, or lakehouse
- read: <https://docs.getdbt.com/guides/best-practices/how-we-structure/1-guide-overview>



# Staging

```
-- stg_stripe__payments.sql

with
source as (
    select * from {{ source('stripe','payment') }}
),
renamed as (
    select
        -- ids
        id as payment_id,
        orderid as order_id,
        -- strings
        paymentmethod as payment_method,
        case
            when payment_method in ('stripe', 'paypal', 'credit_card', 'gift_card') then 'credit'
            else 'cash'
        end as payment_type,
        status,
        -- numerics
        amount as amount_cents,
        amount / 100.0 as amount,
        -- booleans
        case
            when status = 'successful' then true
            else false
        end as is_completed_payment,
        -- dates
        date_trunc('day', created) as created_date,
        -- timestamps
        created::timestamp_ltz as created_at
    from source
)
select * from renamed
```

```
-- base_jaffle_shop__customers.sql

with
source as (
    select * from {{ source('jaffle_shop','customers') }}
),
customers as (
    select
        id as customer_id,
        first_name,
        last_name
    from source
)
select * from customers
```

```
-- stg_jaffle_shop__customers.sql

with
customers as (
    select * from {{ ref('base_jaffle_shop__customers') }}
),
deleted_customers as (
    select * from {{ ref('base_jaffle_shop__deleted_customers') }}
),
join_and_mark_deleted_customers as (
    select
        customers.*,
        case
            when deleted_customers.deleted_at is not null then true
            else false
        end as is_deleted
    from customers
    left join deleted_customers on customers.customer_id = deleted_customers.customer_id
)
select * from join_and_mark_deleted_customers
```



```
-- int_payments_pivoted_to_orders.sql

{%- set payment_methods = ['bank_transfer','credit_card','coupon','gift_card'] -%}
with
payments as (
    select * from {{ ref('stg_stripe_payments') }}
),
pivot_and_aggregate_payments_to_order_grain as (
    select
        order_id,
        {% for payment_method in payment_methods -%}
            sum(
                case
                    when payment_method = '{{ payment_method }}' and
                        status = 'success'
                    then amount
                    else 0
                end
            ) as {{ payment_method }}_amount,
        {%- endfor %}
        sum(case when status = 'success' then amount end) as total_amount
    from payments
    group by 1
)
select * from pivot_and_aggregate_payments_to_order_grain
```

# Intermediate Transforms

- create tables or views that are useful to everyone in different groups, start with views, then move to tables, and then move to incremental tables
- great place to do course grained kimball tables
- Using templating can make your SQL easier to read and much more DRY. You could use this idea in your E and T stages as well.

```
-- customers.sql

with
customers as (
    select * from {{ ref('stg_jaffle_shop__customers')}}
),
orders as (
    select * from {{ ref('orders')}}
),
customer_orders as (
    select
        customer_id,
        min(order_date) as first_order_date,
        max(order_date) as most_recent_order_date,
        count(order_id) as number_of_orders,
        sum(amount) as lifetime_value
    from orders
    group by 1
),
customers_and_customer_orders_joined as (
    select
        customers.customer_id,
        customers.first_name,
        customers.last_name,
        customer_orders.first_order_date,
        customer_orders.most_recent_order_date,
        coalesce(customer_orders.number_of_orders, 0) as number_of_orders,
        customer_orders.lifetime_value
    from customers
    left join customer_orders on customers.customer_id = customer_orders.customer_id
)
select * from customers_and_customer_orders_joined
```

# Vertical Marts

- materialize as tables or incremental models
- unlike kimball, this is a great place to keep your stuff as wide denormalized tables as it will help your BI software
- After marts, what? depends on what you are doing with the data. You might use BI to make a report
- Consider adding a metrics or semantic layer so that more stuff becomes as self-serve as possible

# Also check out sqlmesh

The screenshot shows the SQLMesh interface. On the left, there's a file tree for a project named "/sqlmesh-example". The "models" folder contains "example\_full\_model.sql" and "example\_incremental\_model.sql". The "tests" folder contains "test\_example\_full...". There are also "config.yaml" and "db.db" files. The main area shows the content of "example\_full\_model.sql". The code defines a model named "example\_full\_model" with a cron schedule of "@daily", an audit rule for positive order IDs, and a query that counts distinct item IDs from the "example\_incremental\_model" table. To the right, there's a sidebar with configuration details: Model Name (sqlmesh\_example.example\_full), Start Date (2023-04-06), End Date (2023-04-07), Latest Date (2023-04-06), and Limit (1000). At the bottom, there are "Validate" and "Evaluate" buttons, and tabs for "Table", "Query", and "Console". The "Table" tab is selected, showing a preview of the data:

| item_id | num_orders |
|---------|------------|
| 1       | 6          |
| 2       | 1          |
| 3       | 1          |

# Entire Process Considerations

- You need to decide what is your platform: perhaps warehouse for reports/BI, lakehouse for lower latency BI, data for ML, perhaps lake for raw data and even some ML products, a catalog for a transactional lakehouse, schema discovery, and a model registry/dataset registry
- How do you want to move data through your system: full refreshes, log structured append only tables, incremental ingestion, etc. See <https://tobikodata.com/correctly-loading-incremental-data-at-scale.html> and <https://discourse.getdbt.com/t/on-the-limits-of-incrementality/303> for some considerations
- Your project should have an organized structure and be well tested. This is not just in the dbt parts, but the E and L parts as well.
- We'll revisit all of these in the lab, using a complete open-source stack.

# Post warehouse/mart: metrics

The screenshot shows a development environment with three main panes:

- EXPLORER**: Shows a tree view of project files under "TRY-MALLOY [GITHUB]". Files include ".vscode", "airports.csv", "airports.malloy" (selected), "LICENSE", and "README.md".
- CODE**: An editor pane titled "airports.malloy" containing Malloy schema code. It includes sections for "faa\_region", "by\_state", "by\_facility\_type", and "by\_elevation". A "Run" button is present after each section.
- TERMINAL**: A terminal pane titled "airports->airports\_by\_region\_dashboard" showing query results. It includes tabs for "HTML", "JSON", and "SQL", with "HTML" selected. Results are displayed in three cards:
  - faa\_region**: A bar chart for "AGL" with a value of 4,437.
  - by\_state**: A choropleth map of the United States where states are shaded by airport count. A color scale on the right ranges from 180 (light) to 890 (dark blue).
  - by\_facility\_type**: A table showing the count of airports by facility type:

| fac_type      | airport_count |
|---------------|---------------|
| AIRPORT       | 3,443         |
| HELIPORT      | 826           |
| SEAPLANE BASE | 119           |
| ULTRALIGHT    | 30            |
| STOLPORT      | 11            |
| BALLOONPORT   | 4             |
| GLIDERPORT    | 4             |

**Metrics layer:**  
**malloy, ibis,**  
**dbt semantic**  
**layer**

# Post warehouse/mart: reports/dashboard

## Presentation layer: preset/superset, evidence.dev

Superset Dashboards Charts Datasets SQL

### Quarterly Sales (By Product Line)

Added to 1 dashboard Not available 3 days ago

Chart Source: public.cleaned\_sales\_data

DATA CUSTOMIZE

TIME-SERIES AREA CHART View all charts

Metrics: COUNT(\*)

Time: order\_date

TIME GRAIN: Quarter

TIME RANGE: No filter

Query

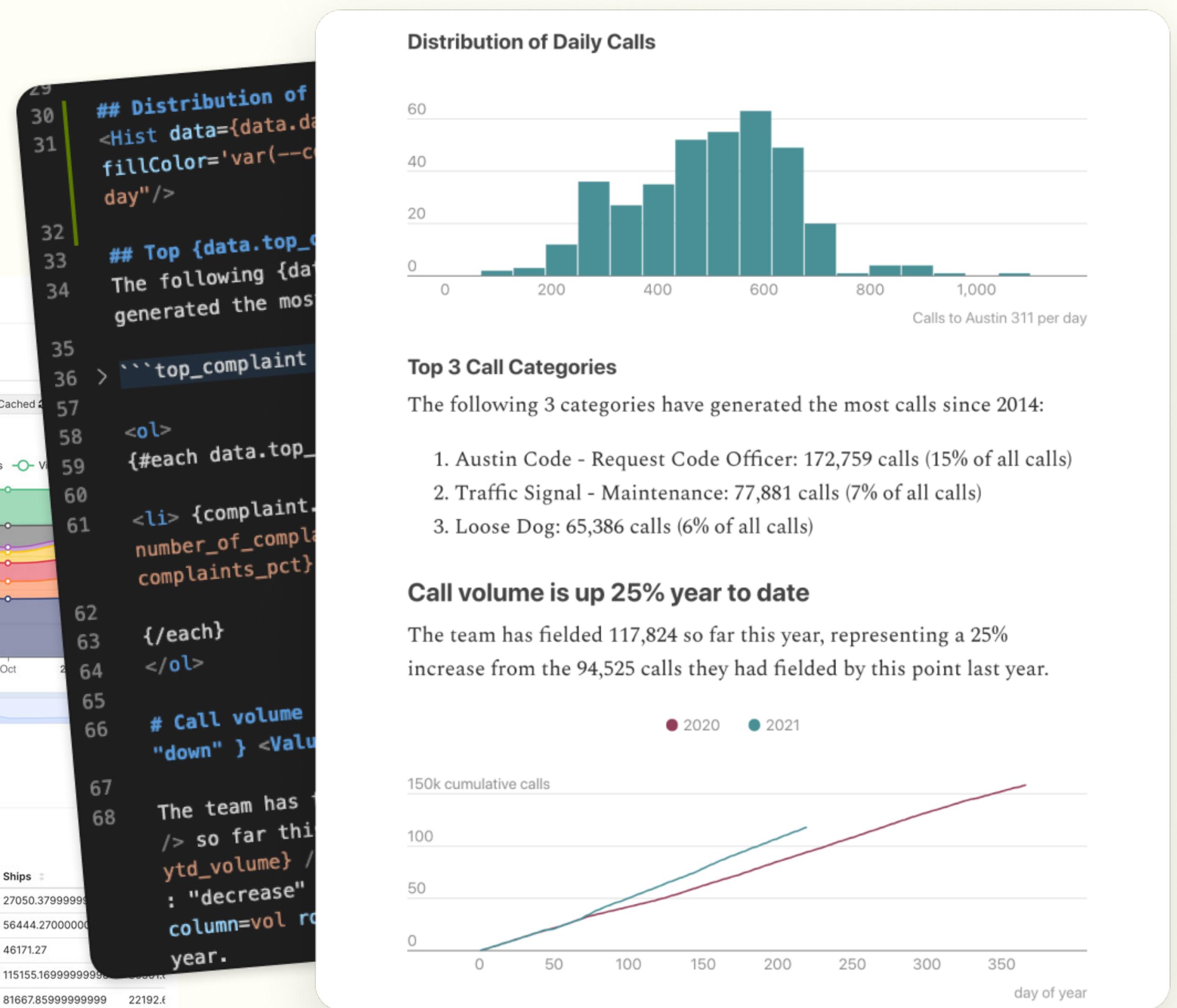
METRICS: SUM(sales)

DIMENSIONS: abc\_product\_line

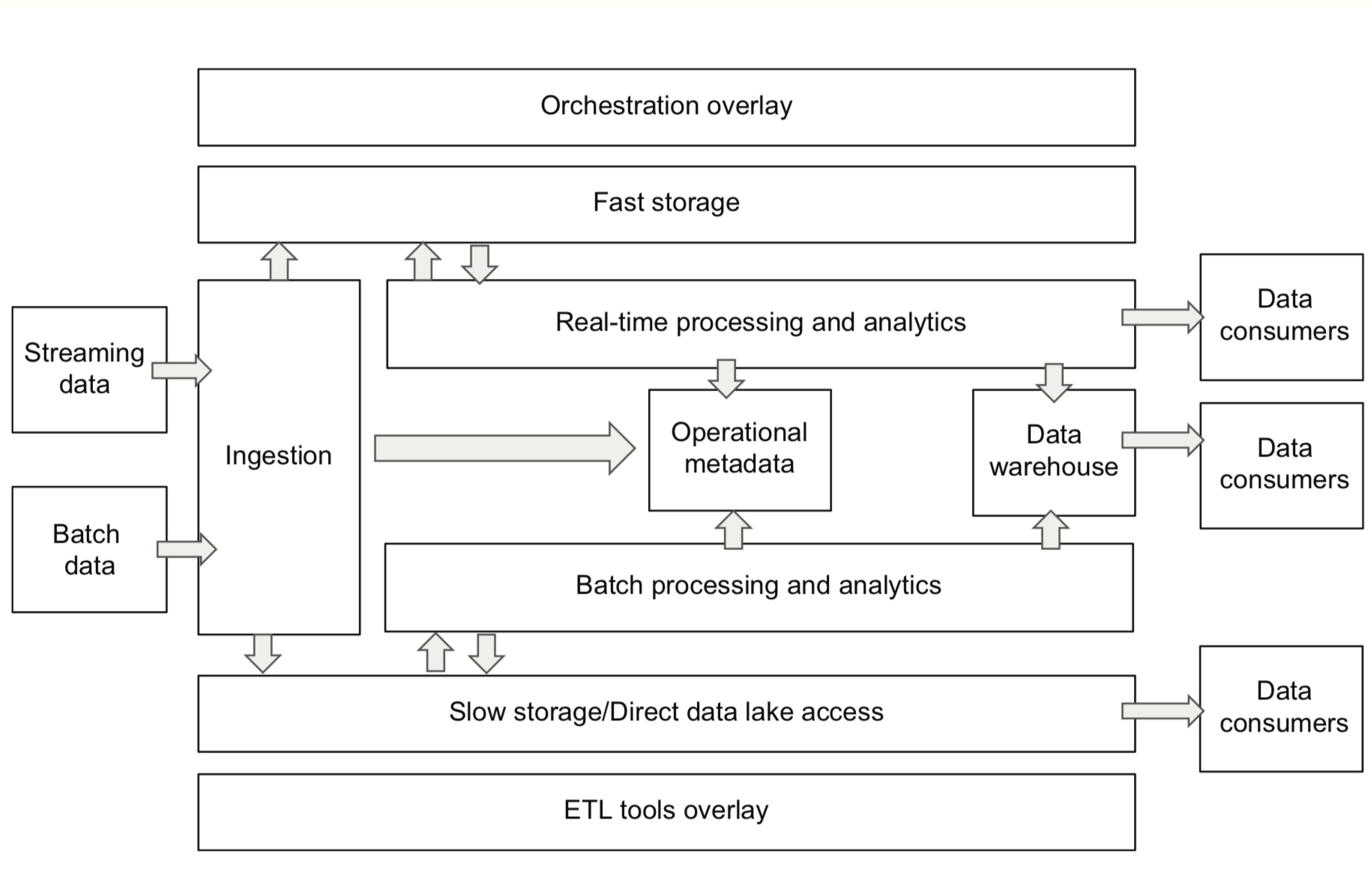
RESULTS SAMPLES

| abc_product_line    | Classic Cars       | Motorcycles       | Planes             | Ships              |
|---------------------|--------------------|-------------------|--------------------|--------------------|
| 2003-01-01 00:00:00 | 166682.87000000002 | 38422.91          | 39205.31000000005  | 27050.37999999999  |
| 2003-04-01 00:00:00 | 208309.87          | 48214.92000000006 | 68678.99           | 56444.27000000001  |
| 2003-07-01 00:00:00 | 280129.06          | 85244.71999999997 | 33938.33           | 46171.27           |
| 2003-10-01 00:00:00 | 829663.49          | 199013.03         | 130434.96999999999 | 115155.16999999998 |
| 2004-01-01 00:00:00 | 336970.3699999999  | 90267.02          | 73113.46999999999  | 81667.85999999999  |

UPDATE CHART



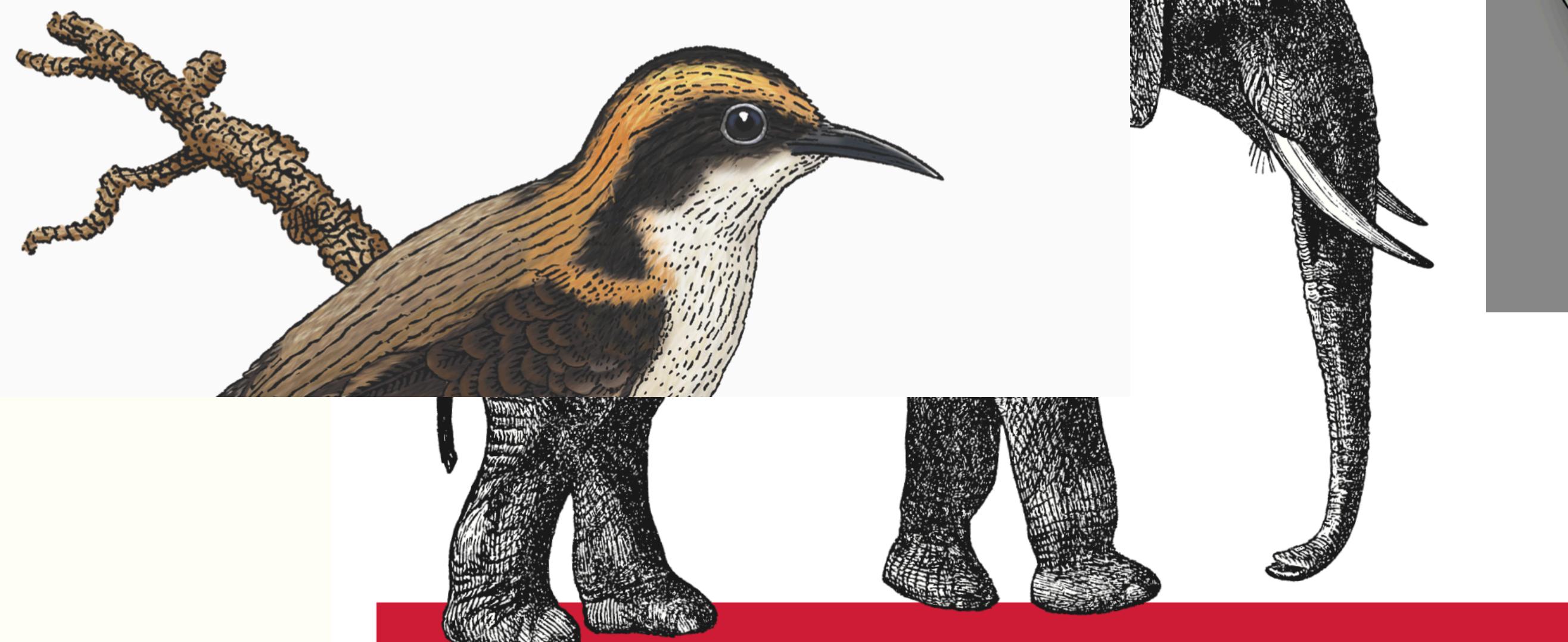
# Infra for the holy grail



O'REILLY®

# Data Algorithms with Spark

Recipes and Design Patterns for Scaling Up  
Using PySpark



Hadoop

Revised



Ergest Xheblati

Minimum  
Viable  
**SQL**  
**Patterns**

*Hands on Design Patterns and  
Best Practices with SQL*

# Designing Cloud Data Platforms

Danil Zburivsky  
Lynda Partner

Univ. AI

Fin (not really)