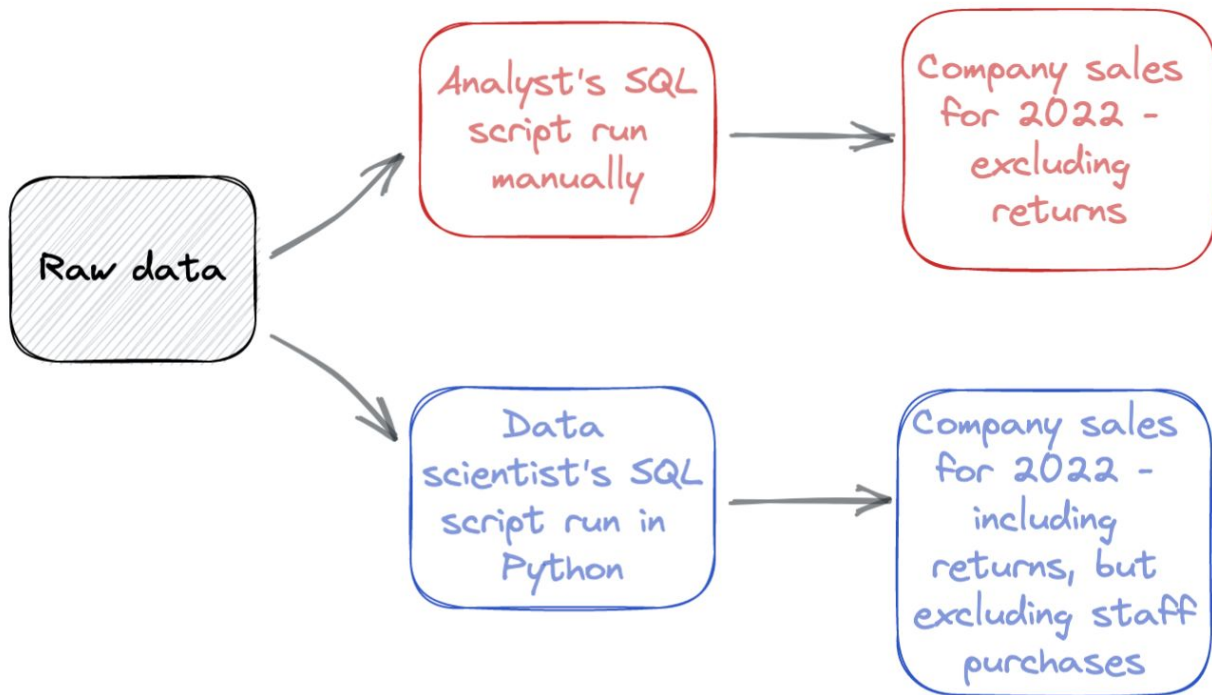


A brief history of data & dbt

How most data pipelines used to work

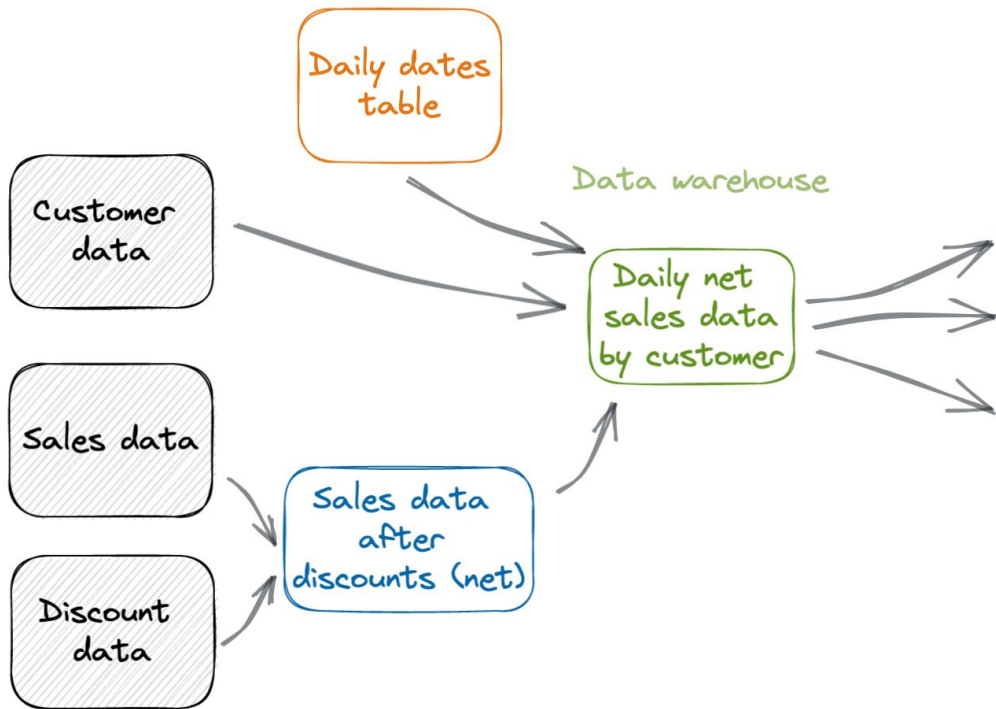
Lots of manual effort, lots of overlap, little consistency



Same metric,
different definitions,
run manually
and not reusable

Data warehousing

Centralised data models



The problems

- The people creating the data models were very separate from the people using them
- You have to manually specify what order to run your SQL scripts in
- Testing for data quality is a manual process (1 SQL file per check)
- Documenting data models is typically done in a different tool
- SQL has limitations compared to other tools like Python

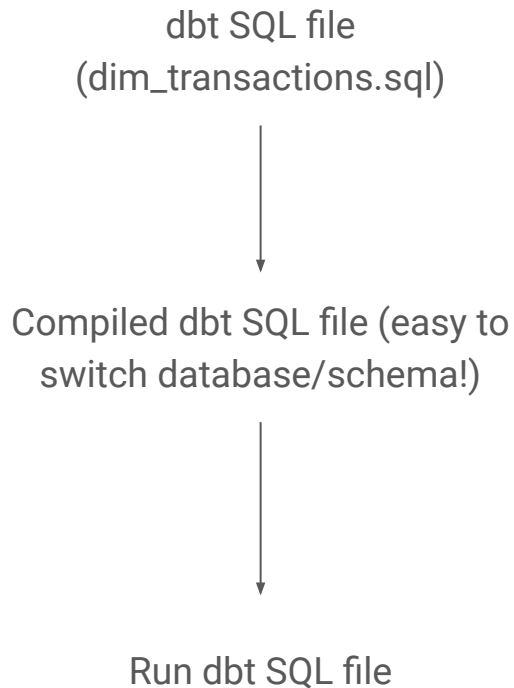
The main benefits of using dbt

**1. Inferring
dependencies**

**2. Documentation
& testing**

**3. Python-like
functionality**

1. Inferring dependencies



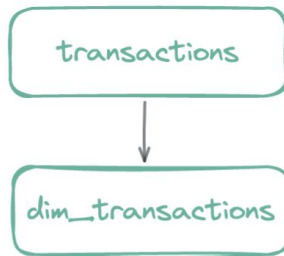
```
SELECT
  date,
  transaction_id,
  amount
FROM {{ ref('transactions') }}
```

```
SELECT
  date,
  transaction_id,
  amount
FROM `my_database.my_schema.transactions`
```

```
CREATE OR REPLACE TABLE `my_database.my_schema.dim_transactions` AS

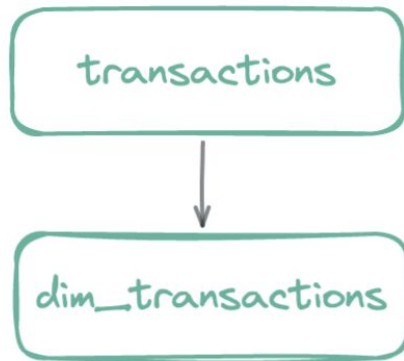
SELECT
  date,
  transaction_id,
  amount
FROM `my_database.my_schema.transactions`
```

Infers dependencies (DAG)



1a. What's a DAG?

- **Directed Acyclic Graph**
 - **Directed:** arrows pointing from one node to another
 - **Acyclic:** no circular dependencies
 - **Graph:** points (nodes) and arrows (edges)
- Fancy way of saying a set of data models (SQL files) with arrows that show what order to run them, without any circular dependencies
- “What order do I need to run my SQL files?”



2. Document & test SQL data models

Using YAML (.yaml) files

```
version: 2

models:
  - name: dim_transactions
    description: One row per transaction_id, with an amount column
    columns:
      - name: date
        description: "Transaction date"
        tests:
          - not_null

      - name: transaction_id
        description: "Unique identifier for this transaction"
        tests:
          - not_null
          - unique
          - relationships:
              to: ref('transactions')
              field: transaction_id

      - name: amount
        description: "Transaction amount"
        tests:
          - not_null
          - dbt_utils.expression_is_true:
              expression: "> 0"
```

You can have:

- One .yaml file per model (e.g. dim_transactions.yaml),
- Or have multiple models in one file (e.g. core_models.yaml)

2. Document & test SQL data models

Using the YAML files to generate (and host) documents

The screenshot displays the dbt Docs web interface. On the left, a sidebar shows the project structure: 'acme' (data, models, mrr, utils). The main content area is titled 'mrr view' and includes tabs for 'Details', 'Description', 'Columns', and 'SQL'. The 'Details' tab is active, showing a table with model metadata:

TAGS	OWNER	TYPE	PACKAGE	RELATION
untagged	TRANSFORMER	view	acme	analytics.dbt_claire_playbook.mrr

Below this, the 'Description' section contains two paragraphs: 'This model represents one record per month, per account (months have been filled in to include any period...)' and 'This model classifies each month as one of: new, reactivation, upgrade, downgrade, or churn.' The 'Columns' section shows a table with model columns:

COLUMN	TYPE	DESCRIPTION
id	TEXT	
date_month	TIMESTAMP_NTZ	
customer_id	NUMBER	

On the right, a 'Lineage Graph' modal is open, showing a flow from 'customer_revenue_by_month' to 'customer_churn_month' to 'mrr'.

- Easy to share
- Shows model lineage

Source: <https://docs.getdbt.com/docs/collaborate/documentation>

3. Python like functionality - variables

Using `set` to create a list
which we then loop through



Compiled SQL

```
{% set categories = ["electronics", "clothes"] %}

SELECT
    date,
    {% for category in categories %}
    SUM(IF(item_category = '{{ category }}', amount, 0)) AS amount_sold_{{ category }},
    {% endfor %}
    SUM(amount) AS amount_sold
FROM {{ ref('transactions') }}
GROUP BY 1
```

```
SELECT
    date,
    SUM(IF(item_category = 'electronics', amount, 0)) AS amount_sold_electronics,
    SUM(IF(item_category = 'clothes', amount, 0)) AS amount_sold_clothes,
    SUM(amount) AS amount_sold
FROM `my_database.my_schema.transactions`
GROUP BY 1
```

3. Python like functionality - macros

Writing a macro, stored
in a separate SQL file

```
{% macro divide_by_100(column_name, precision=2) %}  
    ROUND({{ column_name }} / 100), {{ precision }}  
{% endmacro %}
```

Applying it in a SQL file

```
SELECT  
    {{ divide_by_100(amount_cents, precision=2) }} AS amount_usd  
FROM {{ ref('transactions') }}
```

The compiled SQL

```
SELECT  
    ROUND(amount_cents / 100, 2) AS amount_usd  
FROM `my_database.my_schema.transactions`
```

3. Python like functionality - tests

Writing a column level
test, stored in a
separate SQL file



Applying it to a column
in a YAML file

```
{% test all_values_equal_to(model, column_name, equal_to) %}

-- Returns any records where the column doesn't match a given value
SELECT
*
FROM {{ model }}
WHERE {{ column_name }} != {{ equal_to }}

{% endtest %}
```

```
- name: store_id
  description: "Unique identifier for the store"
  tests:
    - test_all_values_equal_to:
        equal_to: 1234
```

The problems (recap)

- The people creating the data models were very separate from the people using them
- You have to manually specify what order to run your SQL scripts in
- Testing for data quality is a manual process (1 SQL file per check)
- Documenting data models is typically done in a different tool
- SQL has limitations compared to other tools like Python

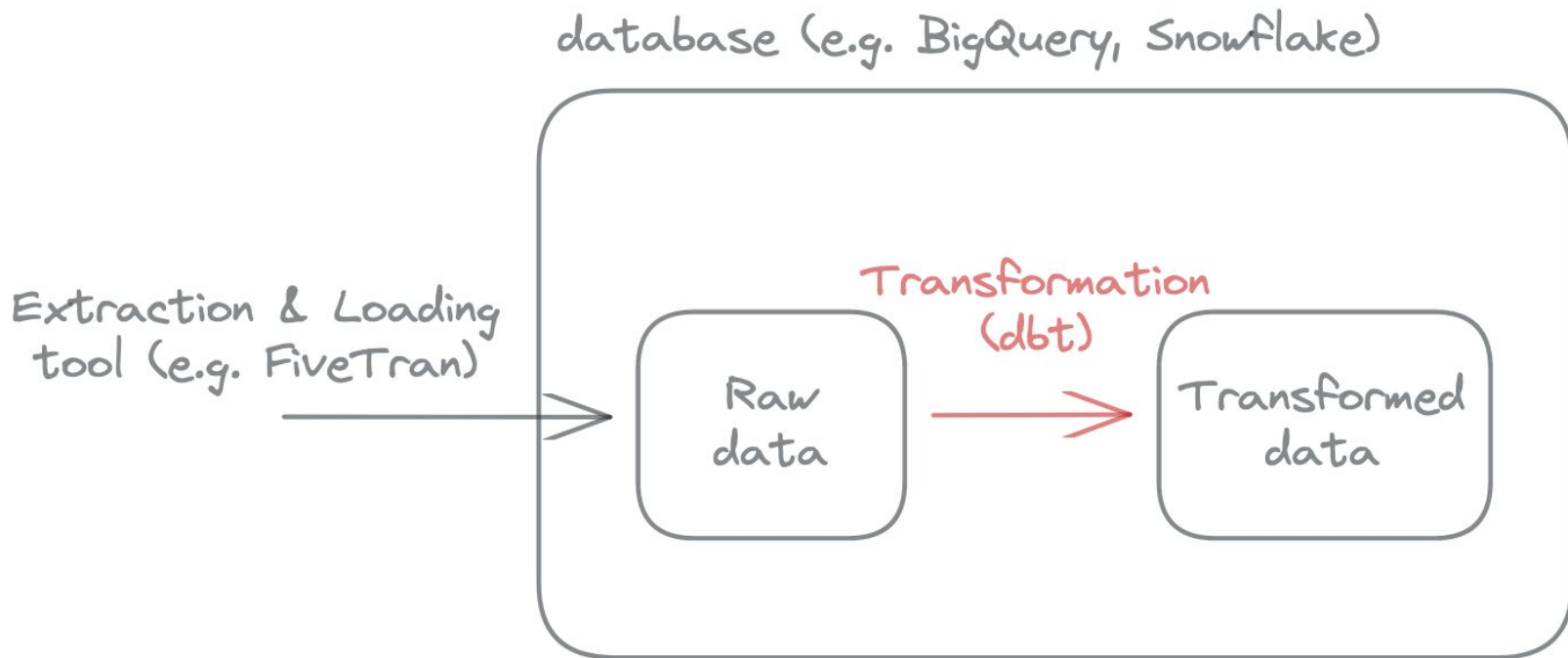
How dbt solved common data problems

- Lowered the technical barrier for creating data pipelines & models
- Automatically infers what order to run SQL models
- Documenting and testing data models is very simple
- Jinja allows for functionality not found in native SQL

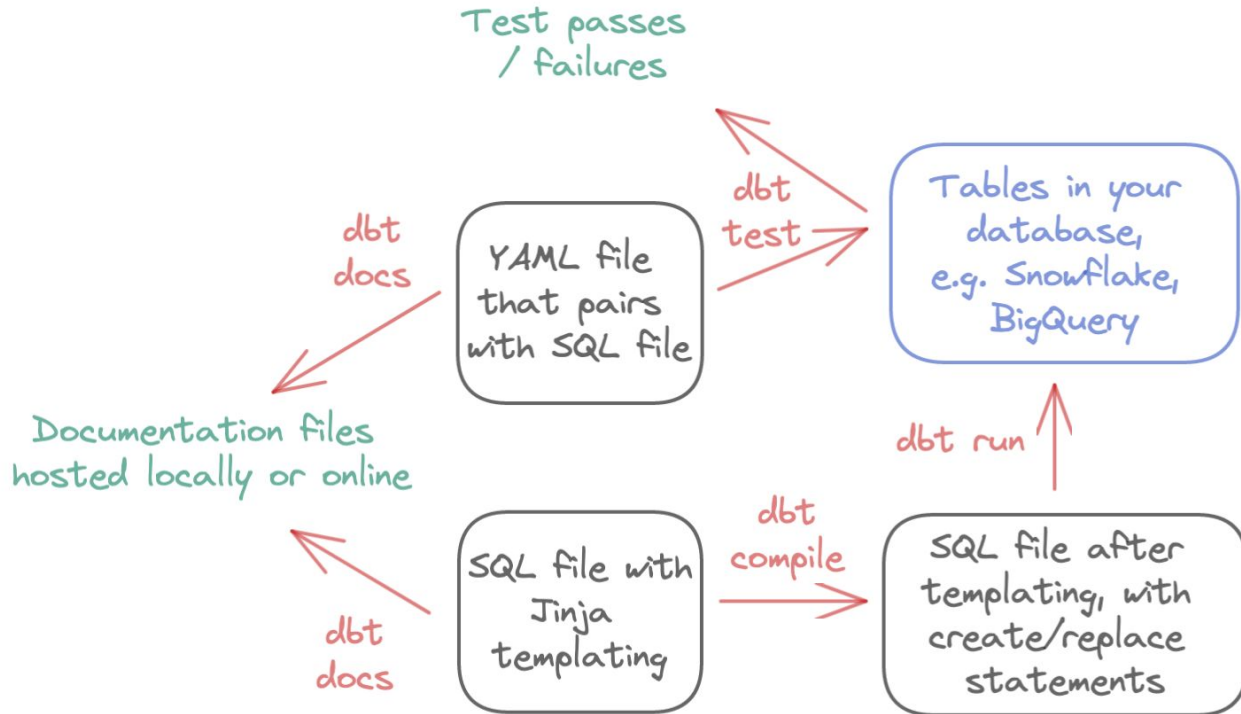
How dbt fits in

Dbt is the “T” in “ELT” (Extract, Load, Transform)


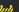


It transforms data within your database




dbt takes care of transforming, testing, and documenting your data



What databases can you use dbt with?

Supported Data Platforms	
Verified Adapters	
Data Platform (click to view setup guide)	latest verified version
AlloyDB	(same as <code>dbt-postgres</code>)
Azure Synapse	1.3.0 
BigQuery	1.2.0
Databricks	1.3.0 
Dremio	1.3.0 
Postgres	1.2.0
Redshift	1.2.0
Snowflake	1.2.0
Spark	1.2.0
Starburst & Trino	1.2.0 

: Verification in progress

Adapter - thing that plugs dbt into a specific type of database, e.g. `dbt-bigquery`

Source: <https://docs.getdbt.com/docs/supported-data-platforms>