# FROM FABRIC TO FANTASTIC

## HOW DBT MAKES LAKEHOUSES AND WAREHOUSES SHINE

Sam Debruyn

Fabric Global Online Conference
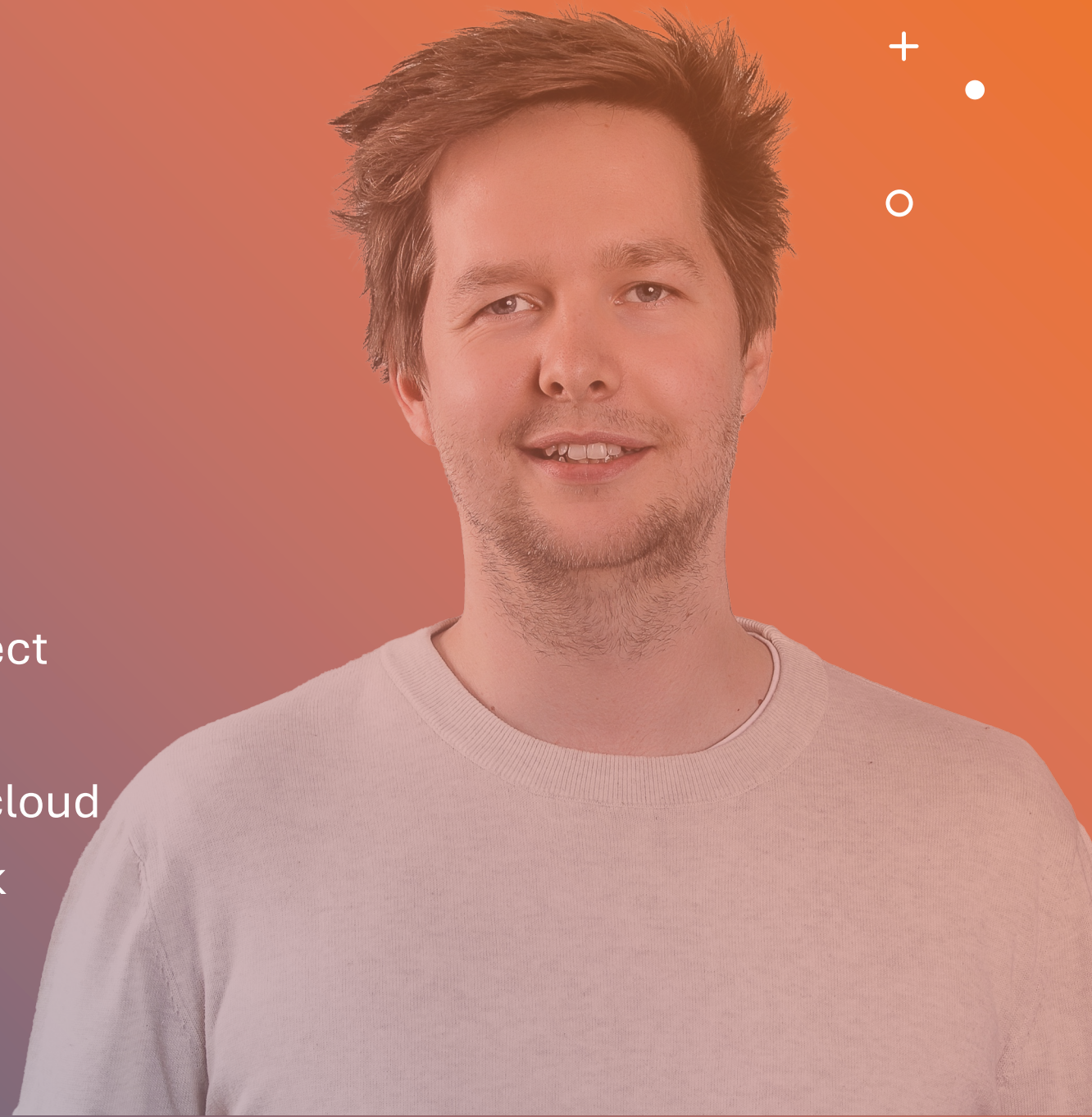September 2024

# Who am I?

**Sam Debruyn**

📍 Heist-op-den-Berg, BE

💼 Consultant / Data & Cloud Architect

5️⃣ years in data

🔟 years in software / architecture / cloud

🤟 dbt, Microsoft, modern data stack
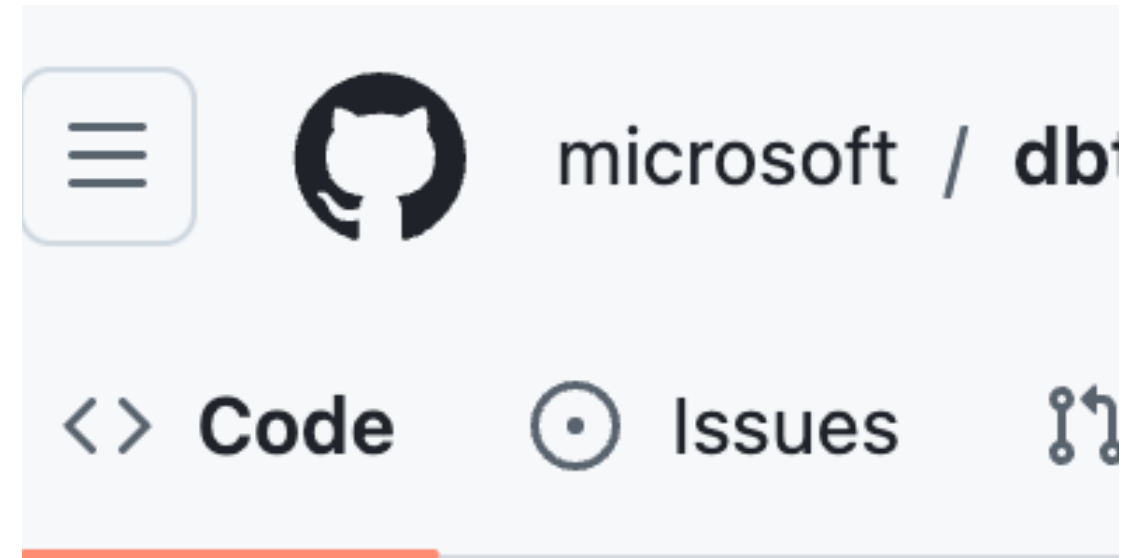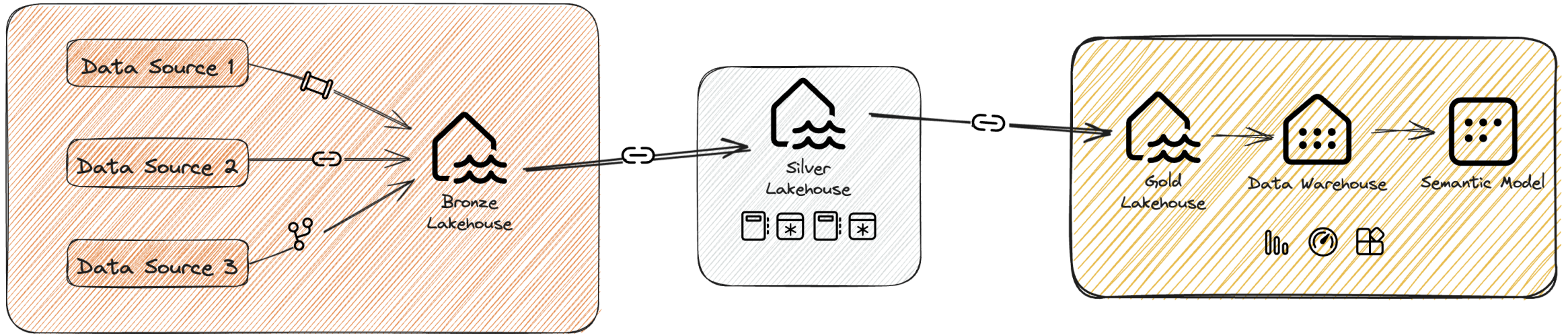
**MVP** **Microsoft®**
Most Valuable
Professional

# dbt-fabric: a quick lookback

OSS project (dbt adapter) created in 2019 to bring dbt support to SQL Server & Azure SQL
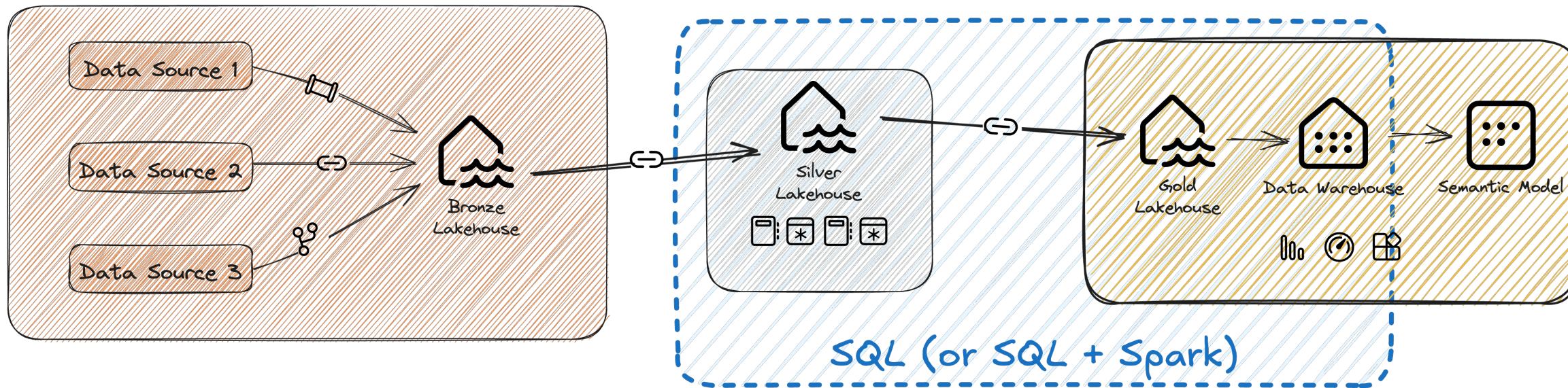
Led maintenance on adapters for SQL Server, Azure SQL, Synapse...

Worked with Microsoft to bring dbt to Fabric

# Lakehouse/Warehouse

Where does SQL fit in?

# Different ways to transform data

## Programming languages

Python and Scala. High learning curves and often creates a boundary between business users and specialized engineers. Very powerful and easy to maintain.

## Declarative languages

SQL, SAS, and the likes. Code is easy to write and understand but offers limited flexibility and can be hard to maintain (adopting software eng. best practices).
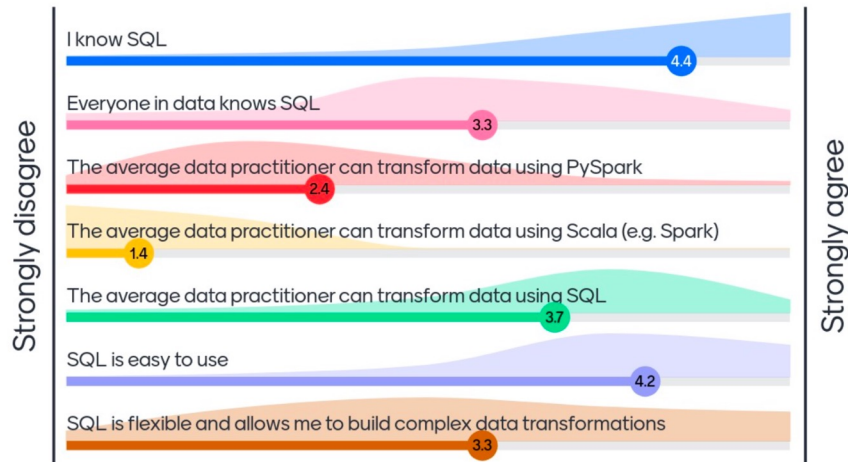
## Low-code / UI-based

Easy to adopt, use, and achieve results. Very high vendor lock-in and limited flexibility and modularity.
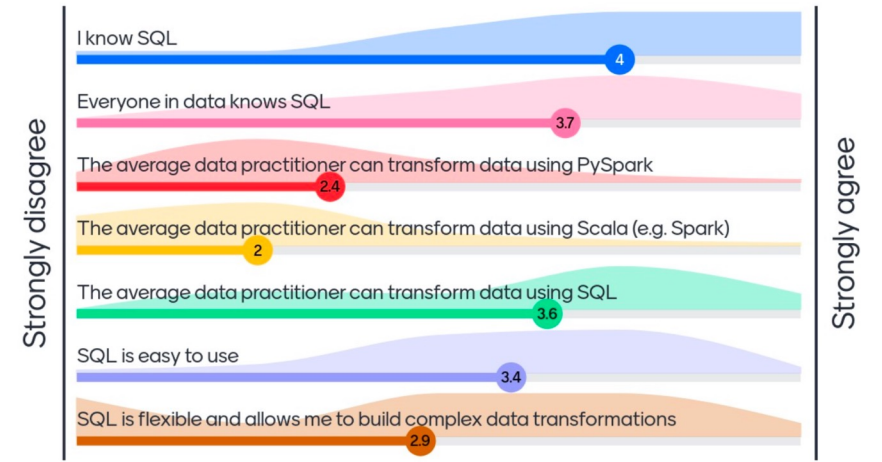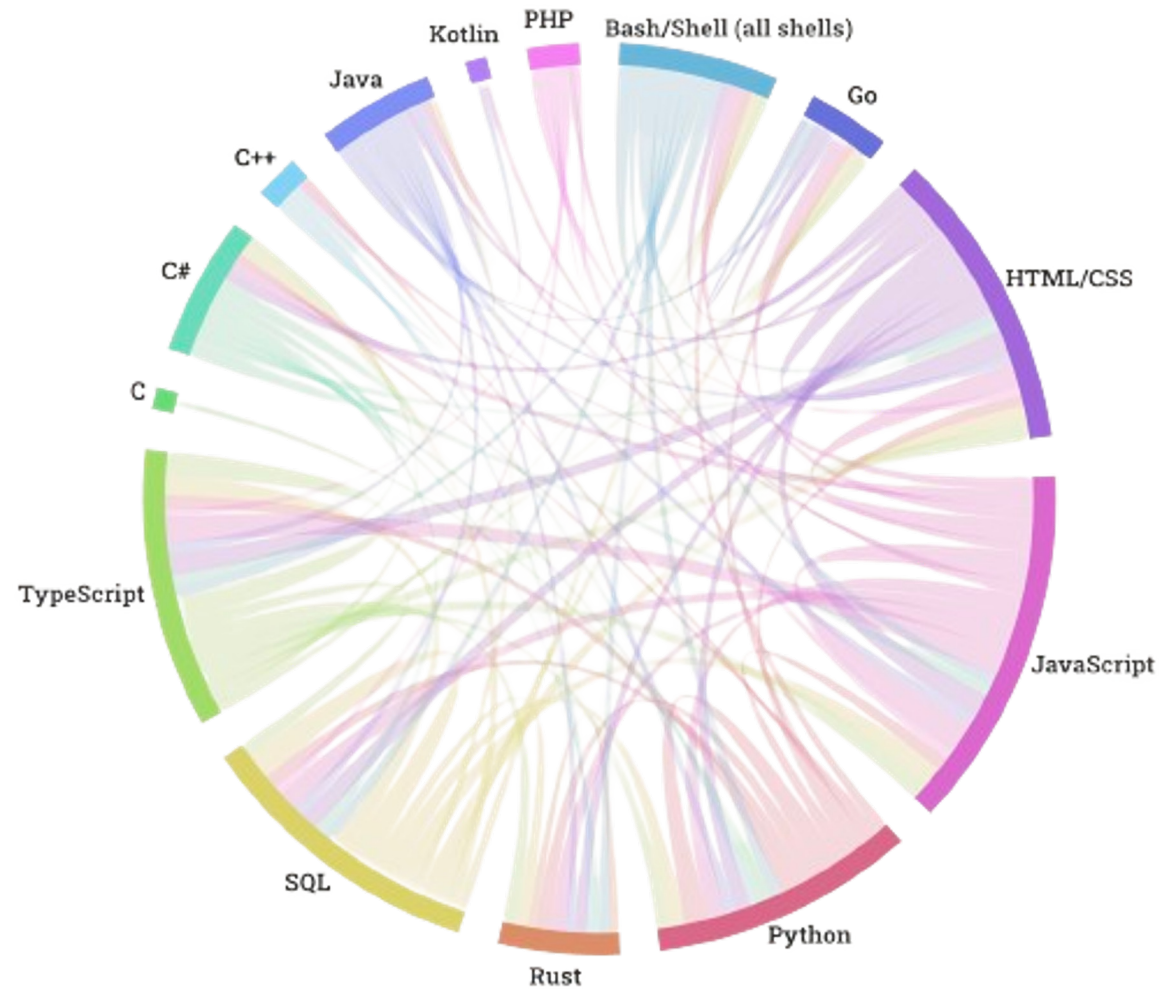
# A quick survey done at local meetups



How would you rate these statements?

| Statement | Rating |
|---|---|
| I know SQL | 4.4 |
| Everyone in data knows SQL | 3.3 |
| The average data practitioner can transform data using PySpark | 2.4 |
| The average data practitioner can transform data using Scala (e.g. Spark) | 1.4 |
| The average data practitioner can transform data using SQL | 3.7 |
| SQL is easy to use | 4.2 |
| SQL is flexible and allows me to build complex data transformations | 3.3 |



How would you rate these statements?

| Statement | Rating |
|---|---|
| I know SQL | 4 |
| Everyone in data knows SQL | 3.7 |
| The average data practitioner can transform data using PySpark | 2.4 |
| The average data practitioner can transform data using Scala (e.g. Spark) | 2 |
| The average data practitioner can transform data using SQL | 3.6 |
| SQL is easy to use | 3.4 |
| SQL is flexible and allows me to build complex data transformations | 2.9 |

# Programming languages 2023



source: Stack Overflow
worked with / wants to
work with

# The common language of data transformations is SQL

**Data architects**

ERWIN
Wherescape
**SQL**

**Data engineers**

Informatica
Matillion
**SQL**

**Analytics engineers**

Alteryx
Talend
**SQL**

**BI developers**

Tableau
Qlik
**SQL**

**Analysts**

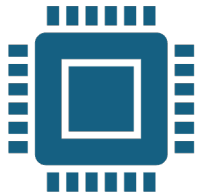Excel / Sheets
Power BI
**SQL**

# Introducing dbt

Open-source Python utility for building data transformations

Free/OSS version: dbt Core / version with all the bells & whistles included: dbt Cloud

The de facto default tool for analytics engineering

# 3 things to know

## No compute

dbt requires a data warehouse to function, it only sends SQL queries

## SQL with Jinja

dbt is built for SQL, in some cases you can also use Python
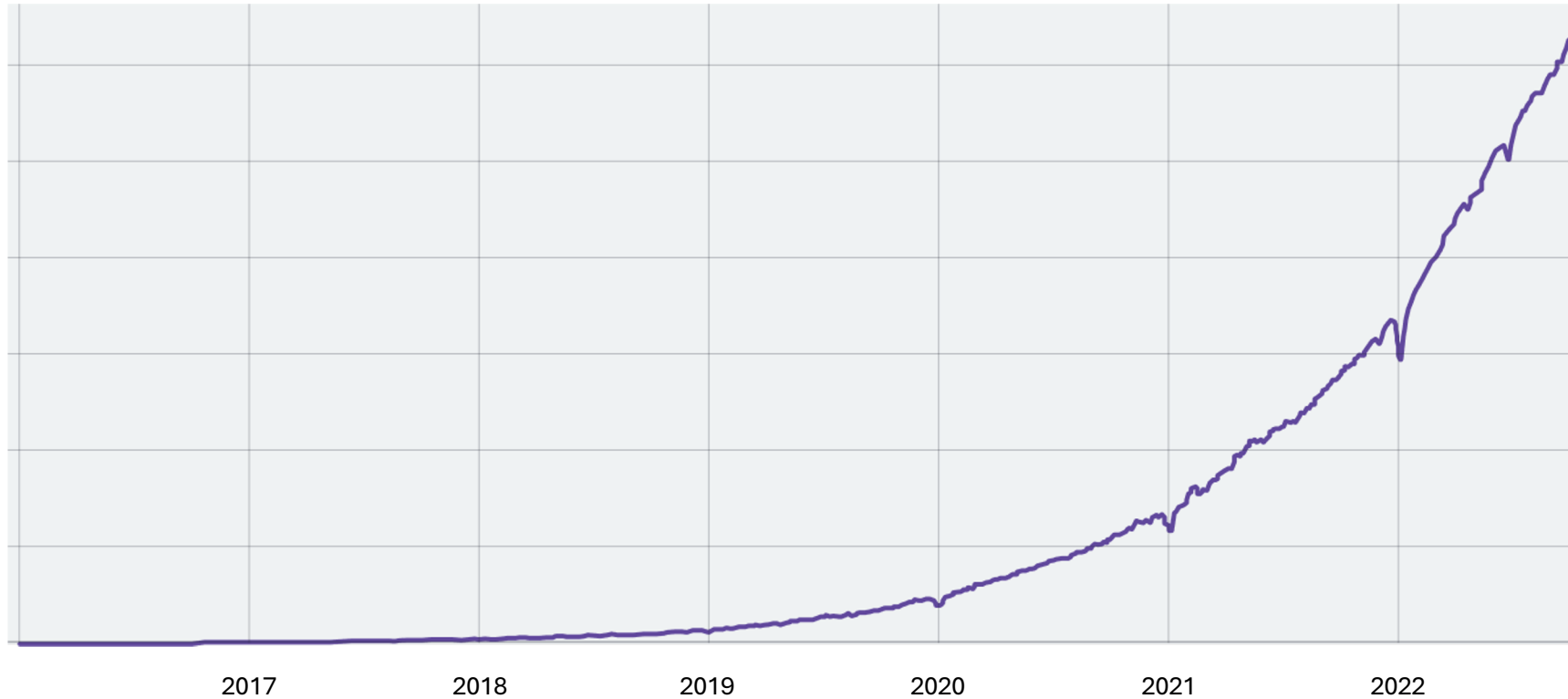
## Free/self-hosted or cloud

dbt Core is free but requires "plumbing" (e.g. an orchestrator)

dbt Cloud is paid, but will be cheaper than building everything around it manually

# dbt adoption past 6 years



October 2023: 30000+ weekly active projects

```sql
{% set item_types=["food", "drink"] %}

compute_booleans as (

    select
        orders.*,

        order_items_summary.order_cost,
        order_items_summary.order_items_subtotal,
        order_items_summary.count_food_items,
        order_items_summary.count_drink_items,
        order_items_summary.count_order_items,
        {% for type in item_types %}
        case
            when order_items_summary.count_{{ type }}_ite
            else 0
        end as is_{{ type }}_order
        {% if not loop.last %},{% endif %}
        {% endfor %}

    from orders

    left join
        order_items_summary
        on orders.order_id = order_items_summary.order_id

),
```

# Modular development

Write transformations in separate version-controlled files

SQL on steroids with Jinja: control logic, loops

Customize and parametrize with variables

Reusable code blocks with macros

Easy to follow DRY principles

# Sources

```yaml
__sources.yml                    ×

models / staging / __sources.yml

1    version: 2
2
3    sources:
4      - name: ecom
5        schema: raw
6        description: E-commerce data for the Jaffle Shop
7        freshness:
8          warn_after:
9            count: 24
10           period: hour
11       tables:
         Generate model
12         - name: raw_customers
13           identifier: customers
14           description: One record per person who has purchased one or more items
         Generate model
15         - name: raw_orders
16           identifier: orders
17           description: One record per order (consisting of one or more order items)
18           loaded_at_field: ordered_at
         Generate model
19         - name: raw_items
20           identifier: items
21           description: Items included in an order
         Generate model
22         - name: raw_stores
23           identifier: stores
24           loaded_at_field: opened_at
```

Manage data sources and monitor data freshness

# Sources

```
stg_customers.sql  ⦿

models / staging / stg_customers.sql  🧭

1   with source as (
2       select *
3       from {{ source('ecom', 'raw_customers') }}
4   ),
5
6   renamed as (
7       select
8           id as customer_id,
9           name as customer_name
10      from source
11  )
12
13  select * from renamed
14  |
```
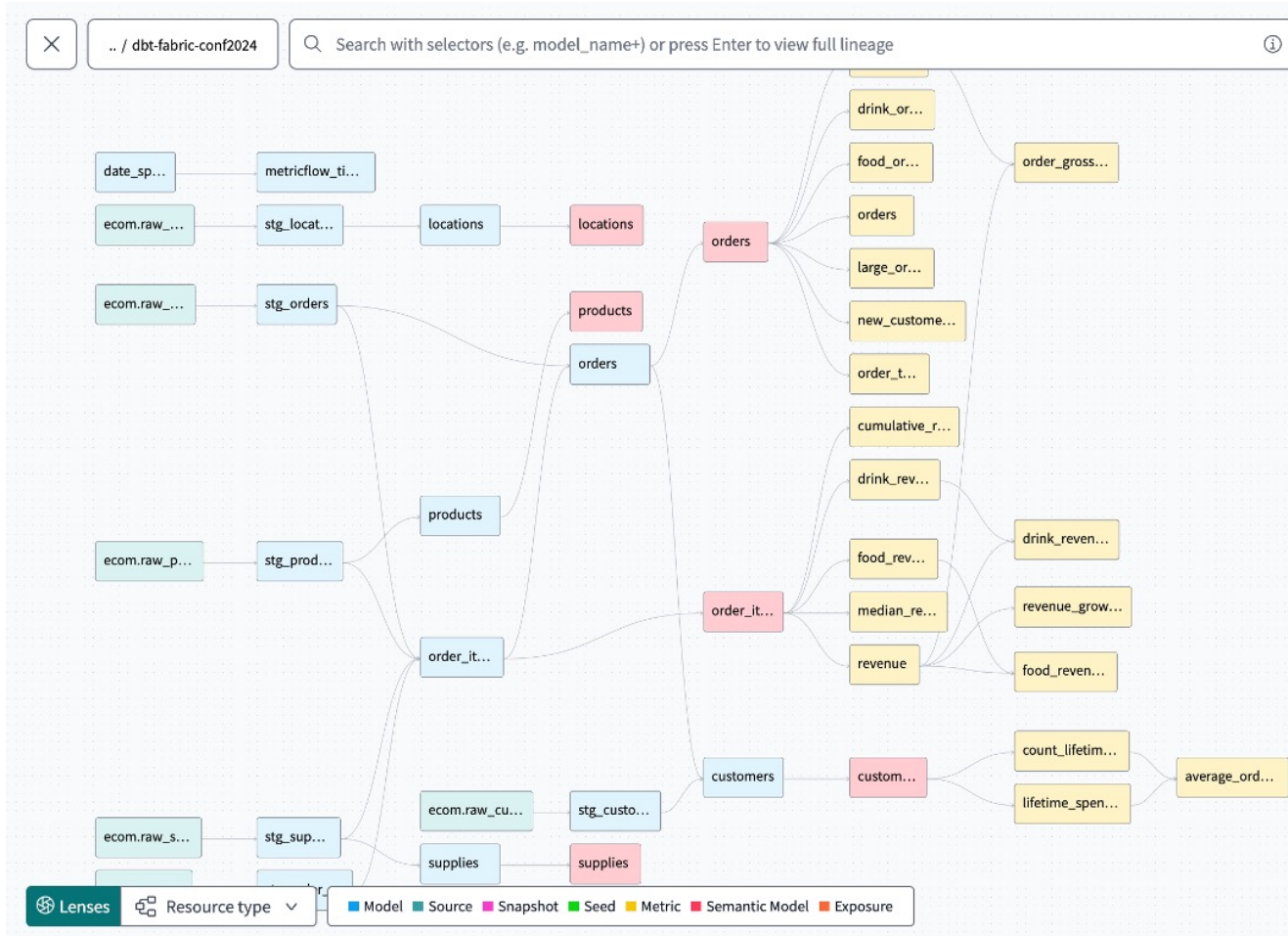
Dynamic schema selection

Start tracking lineage from the source

# Data lineage



Understand the flow of data

Impact of modifying a transformation

How a dimension/fact is constructed

# Data lineage



Spot and detect bad data model design

# Data tests & unit tests

Automated testing for your code, as well as for your data

Tests can be integrated in other tooling to get a good view on your data quality

Simple YAML- or SQL-based syntax to define tests

# Documentation and tests

**customers.yml**

models / marts / customers.yml

```yaml
1   models:
2     - name: customers
3       description: Customer overview data mart, offering key details for each unique customer. One row per customer.
4       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 customers mart.
10          tests:
11            - not_null
12            - unique
13        - name: customer_type
14          description: Options are 'new' or 'returning', indicating if a customer has ordered more than once or has only placed their first order to date.
15          tests:
16            - accepted_values:
17                values: ["new", "returning"]
18
```

# dbt docs



Clear convention-based data documentation

Good step-up to a data catalog

# dbt packages: don't reinvent the wheel

Similar to libraries in software development

Benefit from global knowledge by using pre-built common data transformations and data modelling techniques

Share publicly or privately within your organization

Can contain models (transformations), macros, tests, …

# Date dimension in 1 line

models / marts / date_dimension.sql 🧭     💾 Save

```
1  {% set sql_stmt %}
2      select {{ dateadd(datepart="year", interval=1, from_date_or_timestamp=current_timestamp()) }} as val
3  {% endset %}
4  {{ dbt_date.get_date_dimension('2017-01-01', dbt_utils.get_single_value(sql_stmt)) }}
```

▦ Preview    </> Compile    🔧 Build ⌄    Format     Results    Code quality    Compiled code    Lineage

15.2s | Results limited to 500 rows. ⓘ    Change row display        Download CSV

| date_day | prior_date_day | next_date_day | prior_year_date_day | prior_year_over_year... | day_of_week | day_of_week_iso | day_of_week_ |
|---|---|---|---|---|---|---|---|
| 2017-01-01 | 2016-12-31 | 2017-01-02 | 2016-01-01 | 2016-01-03 | 1 | 7 | Sunday |
| 2017-01-02 | 2017-01-01 | 2017-01-03 | 2016-01-02 | 2016-01-04 | 2 | 1 | Monday |
| 2017-01-03 | 2017-01-02 | 2017-01-04 | 2016-01-03 | 2016-01-05 | 3 | 2 | Tuesday |
| 2017-01-04 | 2017-01-03 | 2017-01-05 | 2016-01-04 | 2016-01-06 | 4 | 3 | Wednesday |
| 2017-01-05 | 2017-01-04 | 2017-01-06 | 2016-01-05 | 2016-01-07 | 5 | 4 | Thursday |
| 2017-01-06 | 2017-01-05 | 2017-01-07 | 2016-01-06 | 2016-01-08 | 6 | 5 | Friday |

# There is more

Implement SCD with *snapshots*

*Incremental* loads

*Hooks* & *operations*

Run *Python models* through Spark (coming soon on Fabric)

Manage access with *grants*
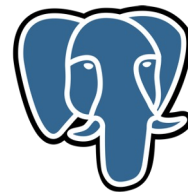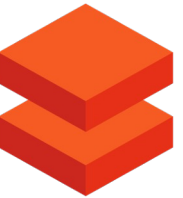
Track dataset usage in BI & ML with *exposures*

*Data contracts*

…

Compatibility

# Accomplish great things

**Version controlled and reproducible**

↗ Collaboration within the team & other teams
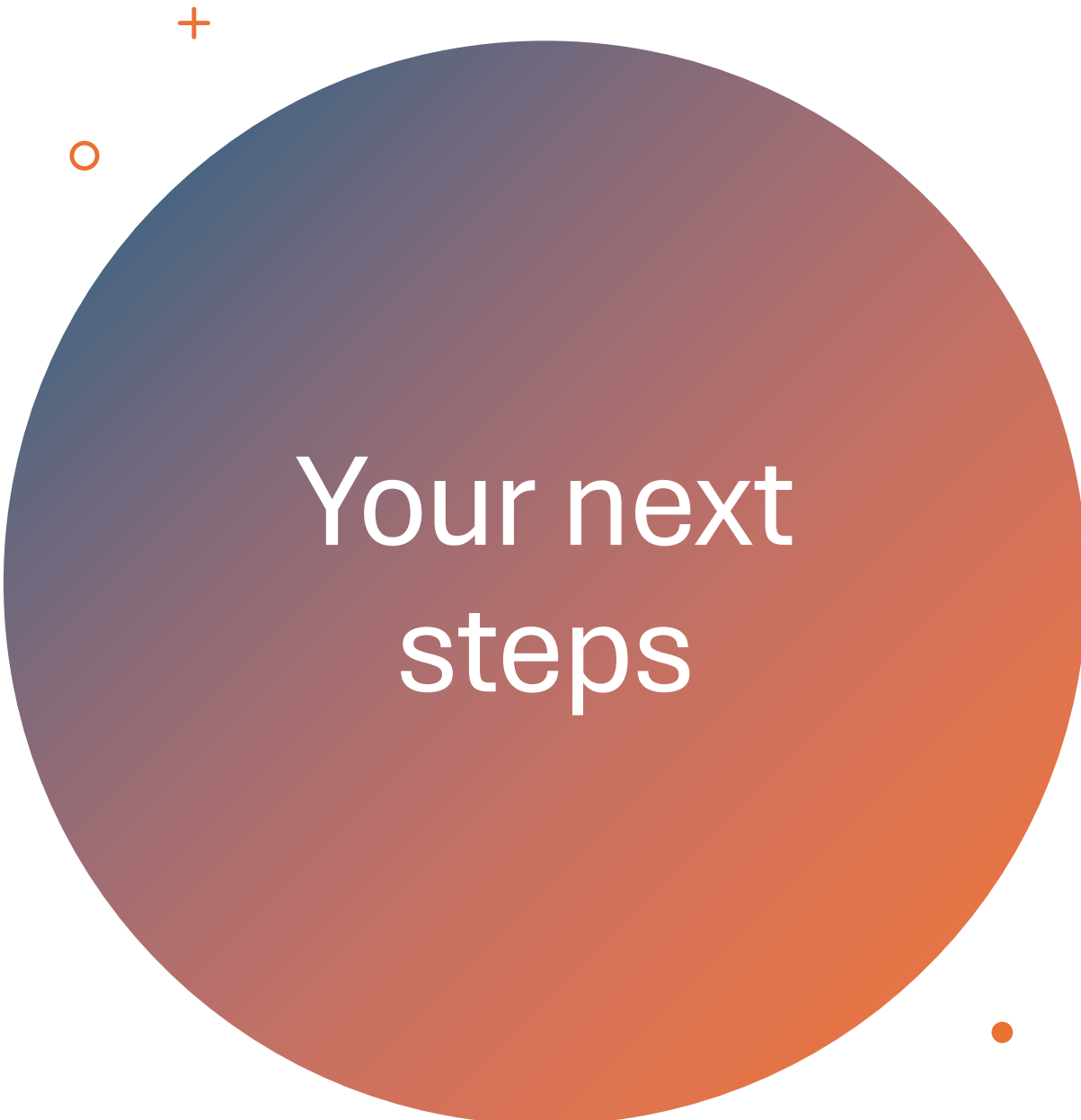
**Built-in docs & lineage**

↗ Know and understand your data

**Test code & data**

↗ Deploy & run with confidence

**Modular & easy to use**

↗ Easy to extend and maintain

# Your next steps

dbt Community: over 100K members

Active and helpful Slack channels

A lot of development in the open-source space

Local meetups all over the world

learn.getdbt.com: free courses to get started with dbt

# Questions?

sam@debruyn.dev

https://debruyn.dev