# Premium sponsors

data on
Power your data

twoday

exmon

Fellowmind

# Standard sponsors

devoteam
Data Driven

INVIXO

INSPARI
a valantic company

backstage

BIWISE
BASED ON INTELLIGENCE
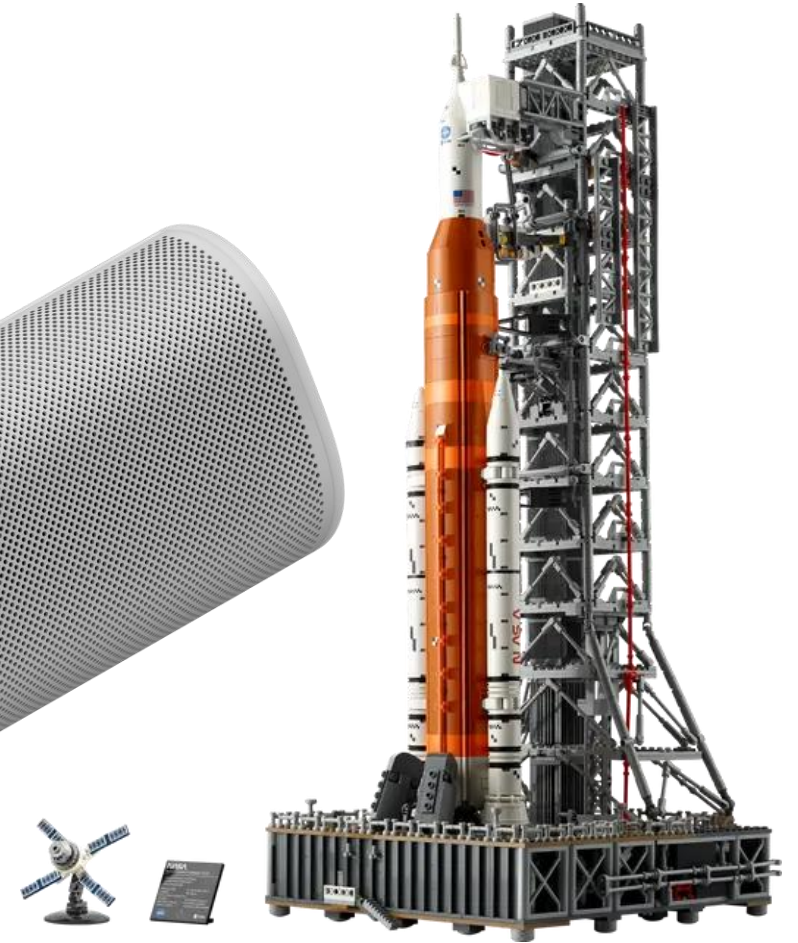
# Raffle Prizes

# Enhancing your Fabric Warehouse with dbt



Sean Douglas Thomsen

Business Intelligence Solution Architect @ LINAK

sdt@linak.com

# Agenda

Introduction to dbt

Tests

Contracts

Documentation

Deployment

# Introduction to dbt

- Data build tools
- Open-source framework for data transformations <u>within</u> data warehouses.
- Only does the T in ETL
- Focuses on transformations using <u>SQL</u>, automates DDL
- Brings software development best practices to data engineering
  - Version control
  - Testing
  - Code reusability
- dbt is written in Python ➔ works well in CI/CD
- Data never leaves our data warehouse!
- Huge community  >30.000 customers

# Introduction to dbt

**dbt Labs**™

👍 This session only discusses dbt-core features

## dbt-core
- Open Source
- Free to use
- Command line only

## dbt-cloud
- Commercial Product
- Browser IDE
- Job scheduling
- CI/CD
- Documentation hosting

# Adapters

**AlloyDB**

✖ Set up in dbt Cloud
✖ Install with dbt Core

pypi package  1.7.2

**BigQuery**

✖ Set up in dbt Cloud
✖ Install with dbt Core

pypi package  1.7.2

**Databricks**

✖ Set up in dbt Cloud
✖ Install with dbt Core

pypi package  1.7.1

**Dremio**

✖ Install with dbt Core

pypi package  1.5.0

**Postgres**

✖ Set up in dbt Cloud
✖ Install with dbt Core

pypi package  1.7.2

**Redshift**

✖ Set up in dbt Cloud
✖ Install with dbt Core

pypi package  1.7.0

**Snowflake**

✖ Set up in dbt Cloud
✖ Install with dbt Core

pypi package  1.7.0

**Spark**

✖ Set up in dbt Cloud
✖ Install with dbt Core

pypi package  1.7.1

**Starburst/Trino**

✖ Set up in dbt Cloud
✖ Install with dbt Core

pypi package  1.7.0

**Microsoft Fabric**

✖ Set up in dbt Cloud
✖ Install with dbt Core

pypi package  1.7.0

**Azure Synapse**

✖ Install with dbt Core

pypi package  1.4.0

🚧 Verification in progress

**Teradata**

✖ Install with dbt Core

pypi package  1.6.7

# dbt is templating engine

We write SQL templates using jinja and configuration and metadata using YAML

SQL templates are compiled to runnable code during build

Compiled SQL is executed on target platform

# dbt terminology

Model: A select statement

Materialization: How should a model be created in target (table, view etc)

Source: External data

Exposure: Downstream use of models outside of dbt

# Models are defined in SQL



```
dbt > models > marts > 🗄 latest_inventory_movement.sql
1    WITH inventtrans
2    AS
3    (
4        SELECT *
5        FROM {{ ref('con_inventtrans') }}
6    ),
7    inventtable AS
8    (
9        SELECT *
10       FROM {{ ref('con_inventtable') }}
11   )
12   SELECT id AS itemid
13   , inventtable.primaryvendorid
14   , CAST(MAX(inventtrans.datephysical) as date) AS latest_datephysical
15   FROM inventtrans
16   INNER JOIN
17   inventtable ON inventtrans.itemid = inventtable.itemid
18            AND  inventtrans.dataareaid = inventtable.dataareaid
19   GROUP BY inventtable.id,
20   inventtable.primaryvendorid
```

# Macros – ref()

- Ref() is used to reference models in our project

- By using ref() we do not have to rewrite downstream models when we move target schemas

➔ Model names need to be unique in dbt

➔ Ref() used to build a DAG

```
dbt > models > marts > 🗄 latest_inventory_movement.sql
  1    with inventtrans
  2    as
  3    (
  4        SELECT *
  5        FROM {{ ref('con_inventtrans') }}
  6    ),
  7    inventtable as
  8    (
  9        SELECT *
 10        FROM {{ ref('con_inventtable') }}
 11    )
```

# We can add our own custom macros



```
        ,custinvoicetrans.qty
        ,custinvoicetrans.lineamount as salescur
        ,{{ safe_divider('custinvoicetrans.lineamount', 'custinvoicetrans.qty') }} as amount_per_unit
from   custinvoicetrans
inner join
custinvoicejour on custinvoicejour.dataareaid = custinvoicetrans.dataareaid
```

dbt > macros > 🛢 safe_divider.sql
```
1    {%- macro safe_divider(column_name, column_name2) -%}
2        CASE WHEN {{column_name2}}  <> 0 THEN {{column_name}} /  {{column_name2}} END
3    {%- endmacro -%}
```

! dbt labs recommends not going overboard with macros. Favor readability over strict DRY adherence (Don't repeat yourself)

```
        ,custinvoicetrans.lineamount as salescur
        ,CASE WHEN custinvoicetrans.qty  <> 0 THEN custinvoicetrans.lineamount /  custinvoicetrans.qty END as amount_per_unit
from   custinvoicetrans
```
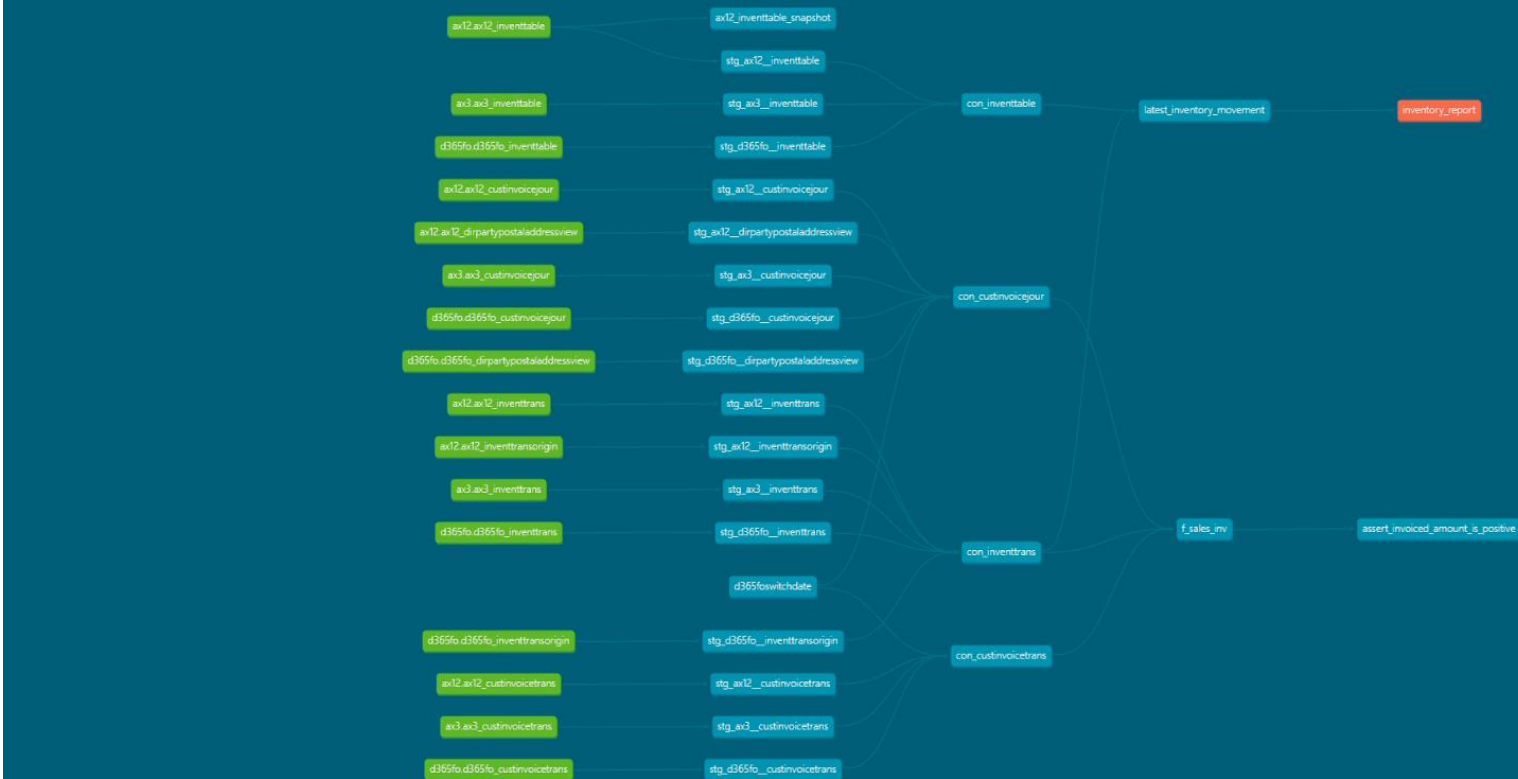
# DAG enables…

- Table lineage
- Execution order
- Graph operators

👍 Select on: models, sources, exposures, tags, path, states, configs…

# dbt-core is a command line tool

# Side note: dbt vs Databricks Jobs

- No DAG

- Data products on different schedules required hand crafted jobs

# Models are described with YAML

```yaml
models:
  - name: latest_inventory_movement
    description: "Latest inventory movement by itemid used for inventory valuation"
    config:
      tags: ["inventory","hourly"]
      contract:
        enforced: true
      grants:
        select: ['sdt@linak.com']
    columns:
      - name: itemid
        description: "Item ID"
        data_type: varchar(30)
        constraints:
          - type: not_null
        tests:
          - unique
      - name: primaryvendorid
        description: "Primary Vendor ID"
        data_type: varchar(30)
        constraints:
          - type: not_null
      - name: latest_inventory_movement
        description: "Latest date of Physical Inventory"
        data_type: date
        tests:
          - date_is_not_in_future
```

# Data tests in dbt

Implemented as SQL queries

Fail when query returns rows

Tests can have different severities

Executed after model has been materialized

➔ Bad data can be in model accessible by users

# Generic Tests

- SQL Templates

- 4 generic tests are included
  - Unique
  - Not_null
  - Accepted_values
  - Relationships

- More can be added manually or by packages

dbt-core / core / dbt / include / global_project / macros / generic_test_sql / unique.sql

iknox-fa  Reformat core [CT-104 CT-105] (#4697)

| Code | Blame | 12 lines (9 loc) · 243 Bytes · |

```
1     {% macro default__test_unique(model, column_name) %}
2
3     select
4         {{ column_name }} as unique_field,
5         count(*) as n_records
6
7     from {{ model }}
8     where {{ column_name }} is not null
9     group by {{ column_name }}
10    having count(*) > 1
11
12    {% endmacro %}
```
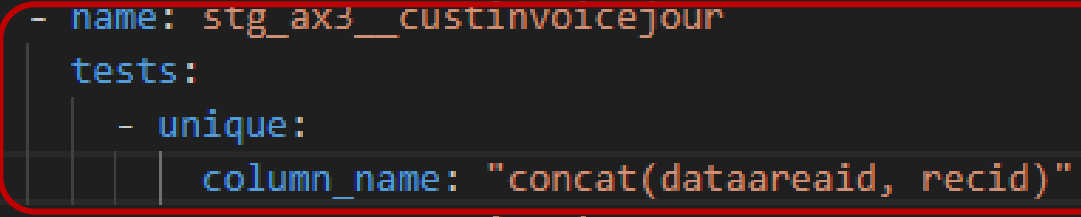
# Adding a generic test

```
models:
  - name: latest_inventory_movement
    description: "Latest inventory movement by itemid used for inventory valuation"
    config:
      tags: ["inventory","hourly"]
    columns:
      - name: itemid
        description: "Item ID"
        datatype: string
        tests:
          - unique
      - name: primaryvendorid
        description: "Primary Vendor ID"
        datatype: string
      - name: latest_inventory_movement
        description: "Latest date of Physical Inventory"
        datatype: timestamp
```

```sql
select
    itemid as unique_field,
    count(*) as n_records

from "msbid_warehouse"."dbo"."latest_inventory_movement"
where itemid is not null
group by itemid
having count(*) > 1
```

# Generic Test

```yaml
models:
  - name: stg_ax3__custinvoicejour
    tests:
      - unique:
          column_name: "concat(dataareaid, recid)"
  - name: stg_ax3__custinvoicetrans
  - name: stg_ax3__inventtable
  - name: stg_ax3__inventtrans
```

```sql
select
    concat(dataareaid, recid) as unique_field,
    count(*) as n_records

from "msbid_warehouse"."dbo"."stg_ax3__custinvoicejour"
where concat(dataareaid, recid) is not null
group by concat(dataareaid, recid)
having count(*) > 1
```

# Generic Tests from dbt_utils

👍 dbt_expectations contains dozens of tests inspired by Great Expectations

❗ Some templates use unsupported SQL syntax on Fabric

## Generic Tests

### equal_rowcount (source)

Asserts that two relations have the same number of rows.

Usage:

```
version: 2

models:
  - name: model_name
    tests:
      - dbt_utils.equal_rowcount:
          compare_model: ref('other_table_name')
```

This test supports the `group_by_columns` parameter; see Grouping in tests for details.

# Custom Generic Test

```
dbt > macros > tests > 🗄 test_date_is_not_in_future.sql
  1    {% macro test_date_is_not_in_future(model, column_name) %}
  2
  3    select
  4        *
  5    from {{ model }}
  6    where {{ column_name }} > getdate()
  7
  8    {% endmacro %}
```

```
select
    *
from "msbid_warehouse"."dbo"."latest_inventory_movement"
where latest_datephysical > getdate()
```

```
- name: latest_inventory_movement
  description: "Latest date of Physical Inventory"
  datatype: timestamp
  tests:
    - date_is_not_in_future
```

# Singular Tests – Work only on a specific model

```
dbt > tests > 🗃 assert_invoiced_amount_is_positive.sql
  1    select
  2        invoiceid,
  3        sum(salescur) as total_amount
  4    from {{ ref('f_sales_inv' )}}
  5    group by invoiceid
  6    having sum(salescur) < 0
```

# Generic Version of Singular Test

```
dbt > macros > tests > 🔲 test_aggregated_amount_by_columns_is_positive.sql
1    {% macro test_aggregated_amount_by_columns_is_positive(model, group_by_column_name, aggregation_column_name) %}
2
3    select
4        {{ group_by_column_name }},
5        sum( {{ aggregation_column_name }} ) as total_amount
6    from {{ model}}
7    group by {{ group_by_column_name }}
8    having sum( {{ aggregation_column_name }} ) < 0
9
10    {% endmacro %}
```

```
- name: f_sales_invoice
  description: "Sales Invoice"
  config:
    tags: ["sales"]
  tests:
    - aggregated_amount_by_columns_is_positive:
        group_by_column_name: invoiceid
        aggregation_column_name: salescur
  columns:
    - name: salesid
      description: "Sales ID"
    - name: custaccount
      description: "Customer Account"
    - name: invoicedate
      description: "Invoice Date"
```

# Unit Tests

- Brand new feature

- Instead of testing on data in warehouse we test on well defined data with a known result

- Are we sure that the logic to find the start of fiscal year always works?

```sql
with stg_sales_invoices as (
    select invoicedate
    from {{ ref('stg_sales_invoices') }}
)
SELECT  invoicedate
    ,CASE
            WHEN DATEPART(MONTH,invoicedate) < 7 -- Are we in H1 or H2 of current year.
                THEN DATEFROMPARTS(YEAR( DATEADD(YEAR,-1,invoicedate) ),7,1) -- Fiscal year start when in H1
                ELSE DATEFROMPARTS(YEAR(invoicedate),7,1) -- same H1
            END as FiscalYearStart
    FROM stg_sales_invoices
```

# Unit Tests

```yaml
unit_tests:
  - name: test_FiscalYearStart
    description: "Check that the fiscal year starts correctly on the first of July"
    model: unit_test_demo
    given:
      - input: ref('stg_sales_invoices')
        rows:
          - {InvoiceDate: 2020-06-30}
          - {InvoiceDate: 2020-07-01}
          - {InvoiceDate: 2020-12-31}
          - {InvoiceDate: 2020-01-01}
    expect:
      rows:
        - {InvoiceDate: 2020-06-30, FiscalYearStart: 2019-07-01}
        - {InvoiceDate: 2020-07-01, FiscalYearStart: 2020-07-01}
        - {InvoiceDate: 2020-12-31, FiscalYearStart: 2020-07-01}
        - {InvoiceDate: 2020-01-01, FiscalYearStart: 2019-07-01}
```
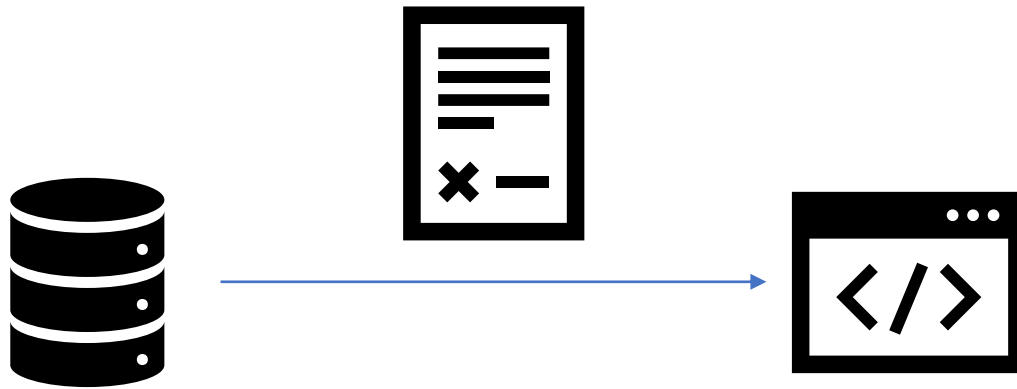
Known edge cases

Expected result

# Contracts

- Data platforms are used by down stream applications

- Any changes to data shapes may break integrations

- Contracts define and enforce a specified data shape

# Contracts are specified in YAML

```yaml
models:
  - name: latest_inventory_movement
    description: "Latest inventory movement by itemid used for inventory valuation"
    config:
      tags: ["inventory","hourly"]
      contract:
        enforced: true
      grants:
        select: ['sdt@linak.com']
    columns:
      - name: itemid
        description: "Item ID"
        data_type: varchar(30)
        constraints:
          - type: not_null
        tests:
          - unique
      - name: primaryvendorid
        description: "Primary Vendor ID"
        data_type: varchar(30)
        constraints:
          - type: not_null
      - name: latest_inventory_movement
        description: "Latest date of Physical Inventory"
        data_type: date
        tests:
          - date_is_not_in_future
```

# Contracts are enforced by platform

```sql
/* Table materialization without contract enabled */
CREATE TABLE [dbo].[latest_inventory_movement] AS
(
    SELECT *
    FROM [dbo].[latest_inventory_movement_temp_view]
)
```

```sql
/* With contract enabled */
CREATE TABLE [dbo].[latest_inventory_movement] (
    itemid VARCHAR(30) NOT NULL
    ,primaryvendorid VARCHAR(30) NOT NULL
    ,latest_datephysical DATE
    )


INSERT INTO [dbo].[latest_inventory_movement] (
    [itemid]
    ,[primaryvendorid]
    ,[latest_datephysical]
    )
SELECT [itemid]
    ,[primaryvendorid]
    ,[latest_datephysical]
FROM [dbo].[latest_inventory_movement_temp_view];
```

# Breaking contracts causes failure

# Tests compared to Contracts

|  | Contracts | Tests |
|---|---|---|
| What? | Data shape | Data quality |
| When? | During materialization. | After materialization. |
| How often? | On every materialization | As specified |
| Support? | Constraints depend on target platform | Works on every platform (if valid SQL) |

# Documentation

dbt docs generates a html page, containing all information and dependencies present in project

! dbt docs does not contain any data from your data warehouse.

# Deploying dbt to Fabric

- There is no Fabric native way to run dbt
- We need python environment
- LINAK Requirement: Run needs to be scheduled by Azure Data Factory

# Deploying to Azure Function



- My colleague Allan created a docker container with dbt installed
- Dbt runs can be started using API
- Solution is cheap, reliable and fast
- Requires knowledge on Azure Functions and docker containers

https://www.linkedin.com/pulse/deploy-dbt-core-workloads-azure-using-durable-allan-rasmussen-6oflf/

# Deploying to Fabric Notebook

```
1   DBT_SCHEMA = "xcu_test"
2   DBT_COMMAND = "run --select  +ax_assettrans"
3   DBT_RETRY_COUNT = 3
✓ 2 min 45 sec - Apache Spark session ready in 2 min 43 sec 258 ms. Command executed in 2 sec 130 ms by Sean Dou
```

- pip install dbt-fabric (or use an environment)

- az copy or git clone to bring in project files

- Use dbt programmatic interface

- Simple, but notebook execution can only be triggered inside Fabric due to API limitations

```
1   import os
2   os.environ['SERVICE_PRINCIPAL_CLIENTID'] = mssparkutils.credentials. \
3               getSecret("<akv-URL>","<secret-name>")
4   os.environ['SERVICE_PRINCIPAL_SECRET'] = mssparkutils.credentials. \
5               getSecret("<akv-URL>","<secret-name>")
6   os.environ['DBT_SCHEMA'] = DBT_SCHEMA
7   os.environ['DBT_PROJECT_DIR'] = '/mnt/repo'
8   os.environ['DBT_PROFILES_DIR'] = '/mnt/repo'
✓ 1 sec - Command executed in 909 ms by Sean Douglas Thomsen on 12:44:56 PM, 6/08/24
```

```
1   from dbt.cli.main import dbtRunner, dbtRunnerResult
2   import re
3   import time
4
5   retriable_commands = {
6           "build",
7           "compile",
8           "seed",
9           "snapshot",
10          "test",
11          "run",
12          "run-operation",
13      }
14
15  # initialize
16  dbt = dbtRunner()
17
18  # create CLI args as a list of strings. Unnessary white spaces are removed.
19  cli_args = re.sub(r"\s+", " ", DBT_COMMAND).strip().split(sep=' ')
20
21  # run the command
22  for attempt in range(0,1+DBT_RETRY_COUNT):
23      if attempt==0:
24          res: dbtRunnerResult = dbt.invoke(cli_args)
25      if attempt >0 and cli_args[0] in retriable_commands:
26          res: dbtRunnerResult = dbt.invoke(["retry"])
27      if res.success:
28          break
29      time.sleep(15)
✓ 1 min 4 sec - Command executed in 1 min 3 sec 835 ms by Sean Douglas Thomsen on 12:46:00 PM, 6/08/24
```

# New Announcements

## Data pipeline support for DBT CLI

Estimated release timeline: Q3 2024

Release Type: Public preview

## Fabric Core REST APIs support Service Principal

Estimated release timeline: Q3 2024

Release Type: General availability

## Data workflows: Build data pipelines powered by Apache Airflow

Estimated release timeline: Q2 2024

Release Type: Public preview

Data workflows are powered by Apache Airflow and offer an integrated Apache Airflow runtime environment, enabling you to author, execute, and schedule Python DAGs with ease.

# dbt is much more

- dbt seeds – Define manual data in your projects
- Logging – have dbt log run executions to your warehouse
- Snapshots – use dbt to snapshot your source data or create SCD
- Packages – include open-source packages in your project
- Hooks
- Grant – use dbt to manage access to your models
- Incremental models
- Slim ci
- …

# Key Take-aways

☺ dbt is in essence a templating engine
☺ Well suited for teams that prefer SQL
☺ Encourages reusing code, by using templates and macros
☺ Less boiler plate code, more focus on transformations
☺ Provides a data quality testing framework
☺ Can enforce data contracts
☺ Makes documentation a part of development workflow
☹ Deployment to Fabric could be easier
☹ Many popular packages do not support dbt-fabric yet

# Thank you for your time!