# Introduction to using and scaling dagster

Aleksandar Milicevic
Georg Heiler

# Agenda

1 **Welcome and introduction to**
- What is a data platform
- What is the role of dagster

2 **Crash course on Dagster concepts**
- Asset based
- Metadata-created pipelines
- Resources and IO managers

3 **Hands-on lab** (local or GitHub Codespace)
- Spin up the environment
- Tour Dagster UI and run first asset
- Understand, run and extend the above Dagster concepts

4. **Dagster @Magenta**
- Open source implementation with local-data-stack

QA & Wrap-up

# About us



**Data expert** in academia and industry Magenta Telecom

- meetup organizer and conference speaker

- data architecture, multimodal and complex data challenges

in geoheil

geoheil

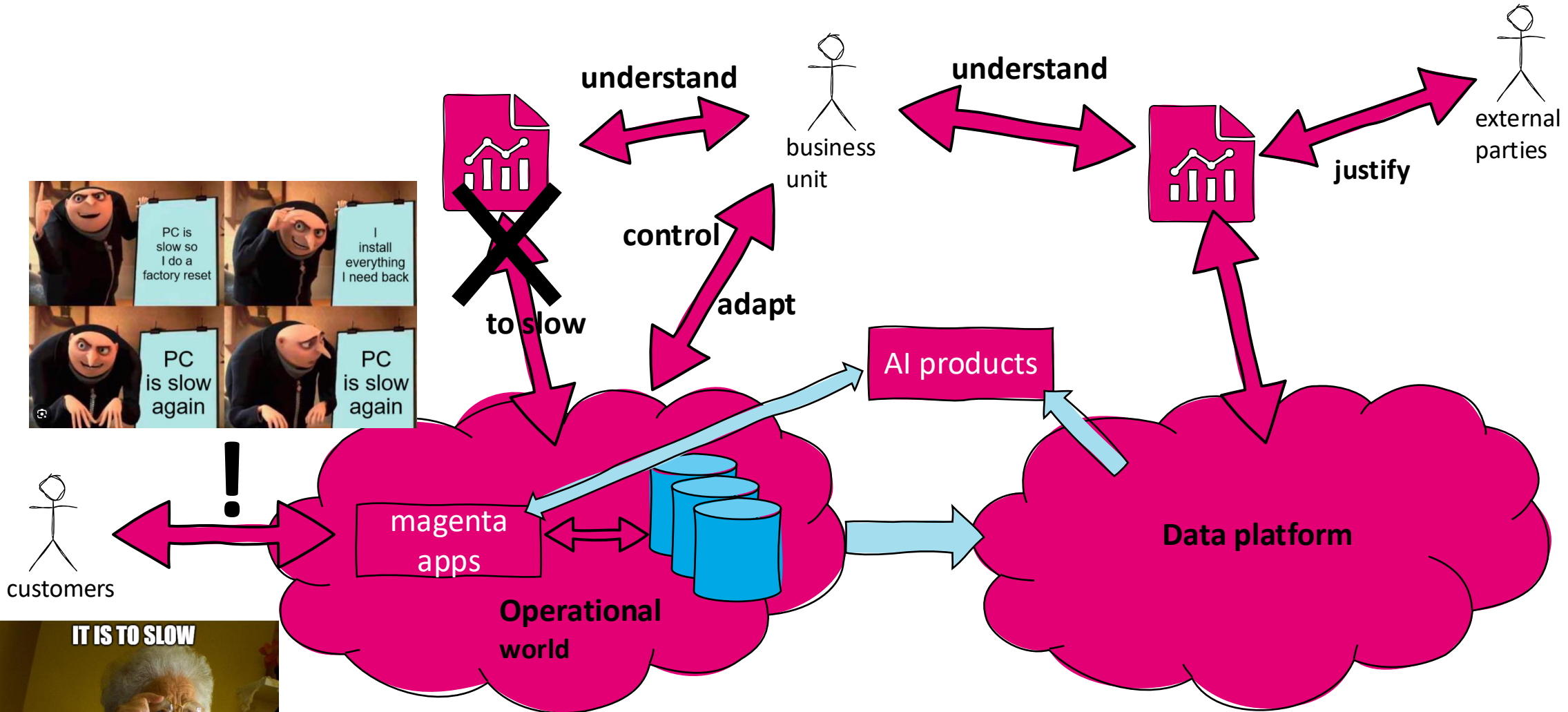@geoheil.com




**Data engineer** at Magenta Telecom

- database internals, data platform engineering

in milicevica23

milicevica23

@milicevica23.bsky.social

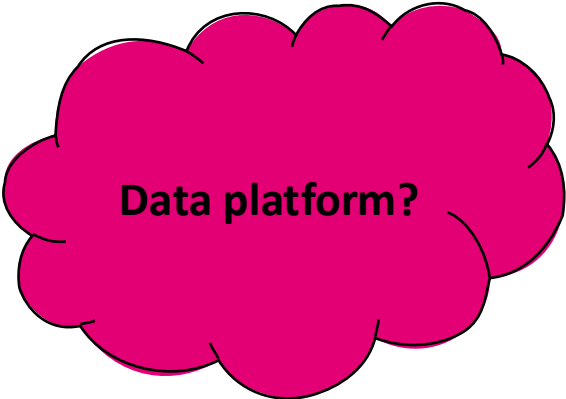# Data platform is there to help understand business process improve service



- customer facing performance is important – reporting is second priority (in a short run)
- different technology for different use case – oltp vs olap databases
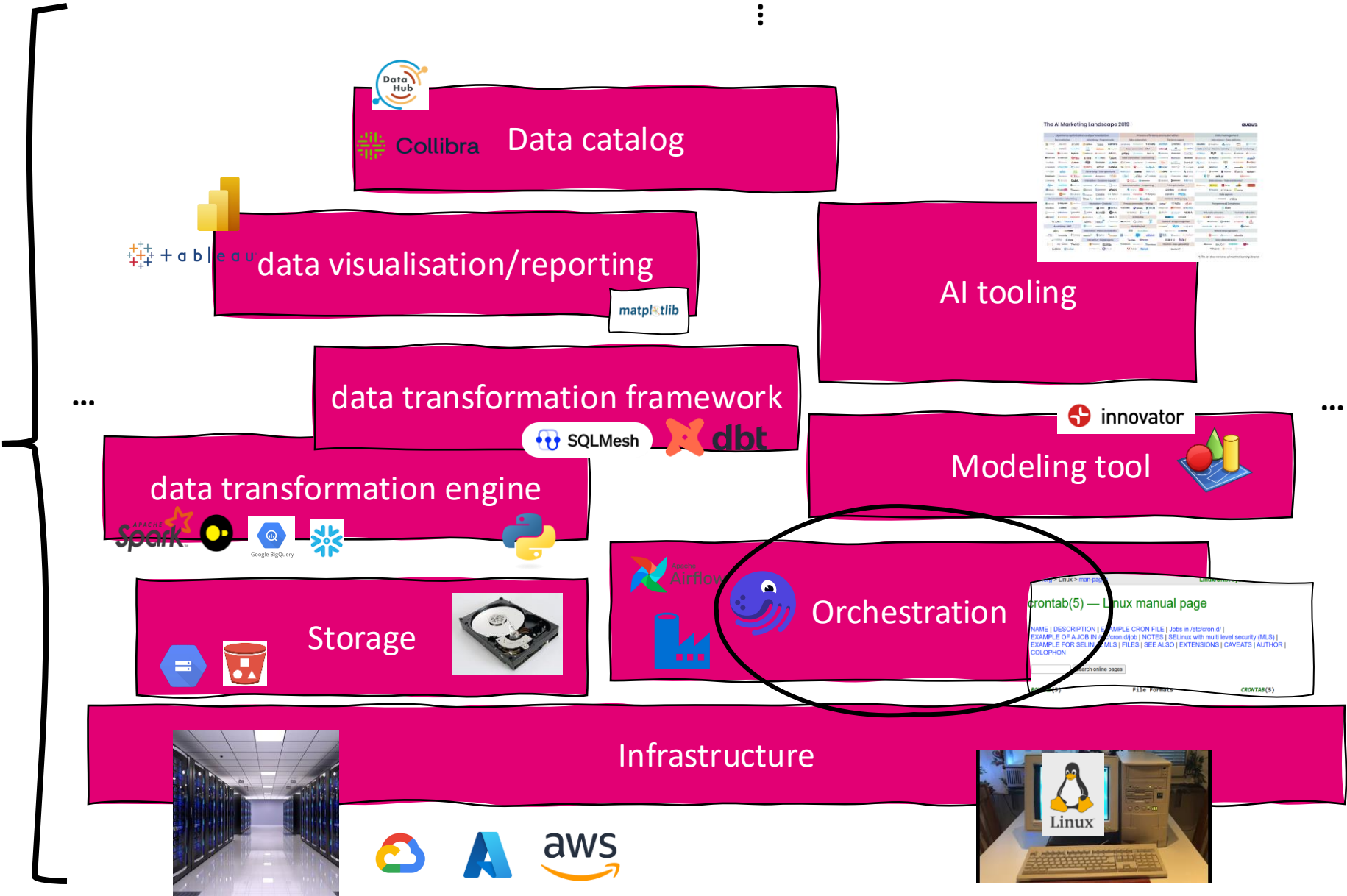- company data is an asset as car, phone or internet cable

# Data platform components

**What do I need for a bare minimum E2E reporting?**

**How to make it last?**
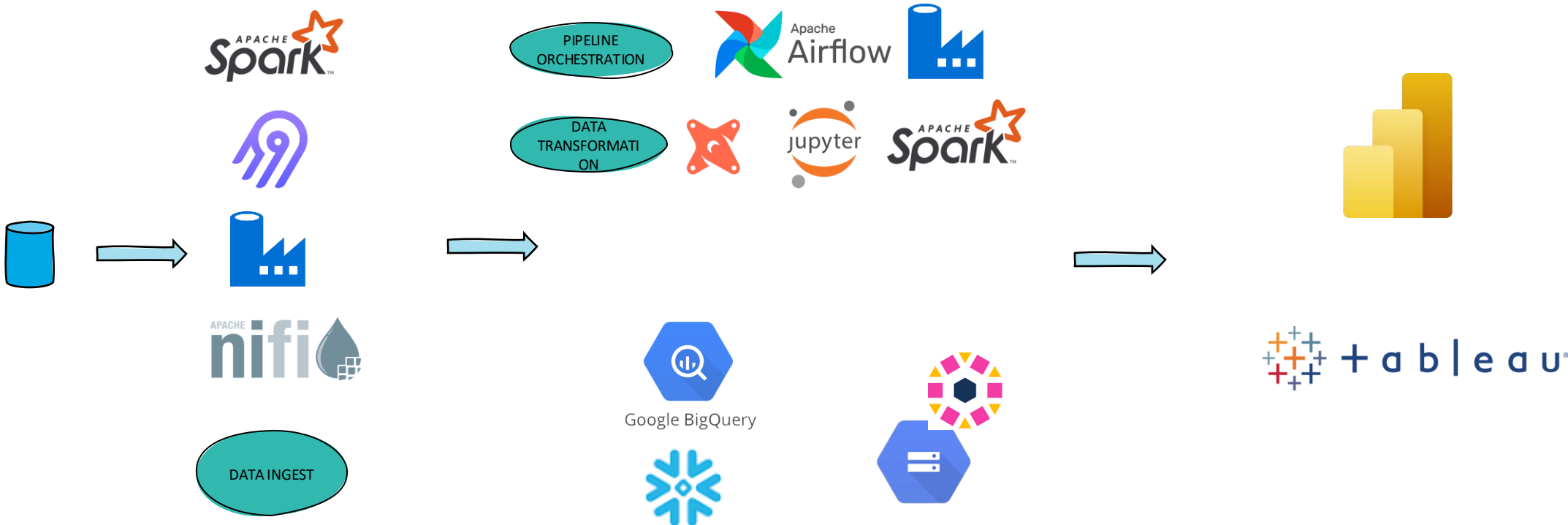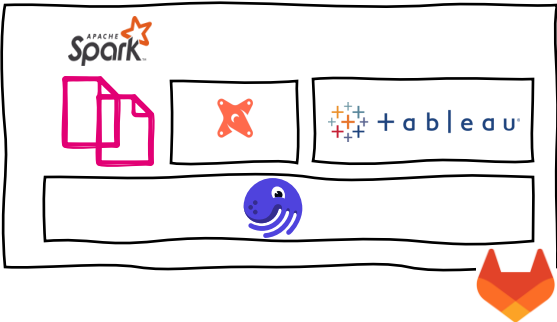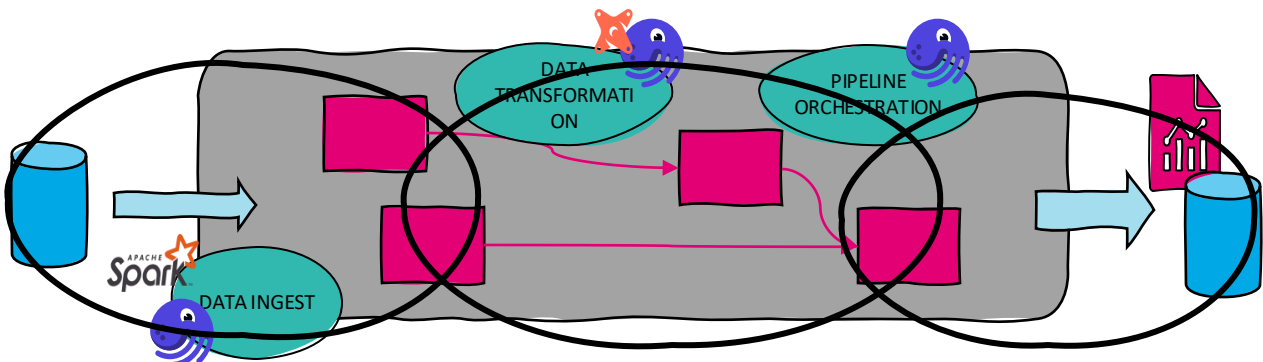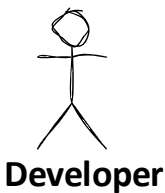**Modern tooling, governance, development…**

**Data platform?**

**How to make it appealing?**
**AI…**

Data catalog

Collibra

data visualisation/reporting

AI tooling

data transformation framework

SQLMesh    dbt

Modeling tool

innovator

data transformation engine

Spark    Google BigQuery

Storage

Airflow

Orchestration

crontab(5) — Linux manual page

Infrastructure

aws

# Understanding tool silos

**What should i do to get E2E reporting use case done?**

**Developer**

DATA INGEST

DATA TRANSFORMATION

PIPELINE ORCHESTRATION

---

APACHE **Spark**

PIPELINE ORCHESTRATION

Apache **Airflow**

DATA TRANSFORMATION

jupyter

APACHE **Spark**

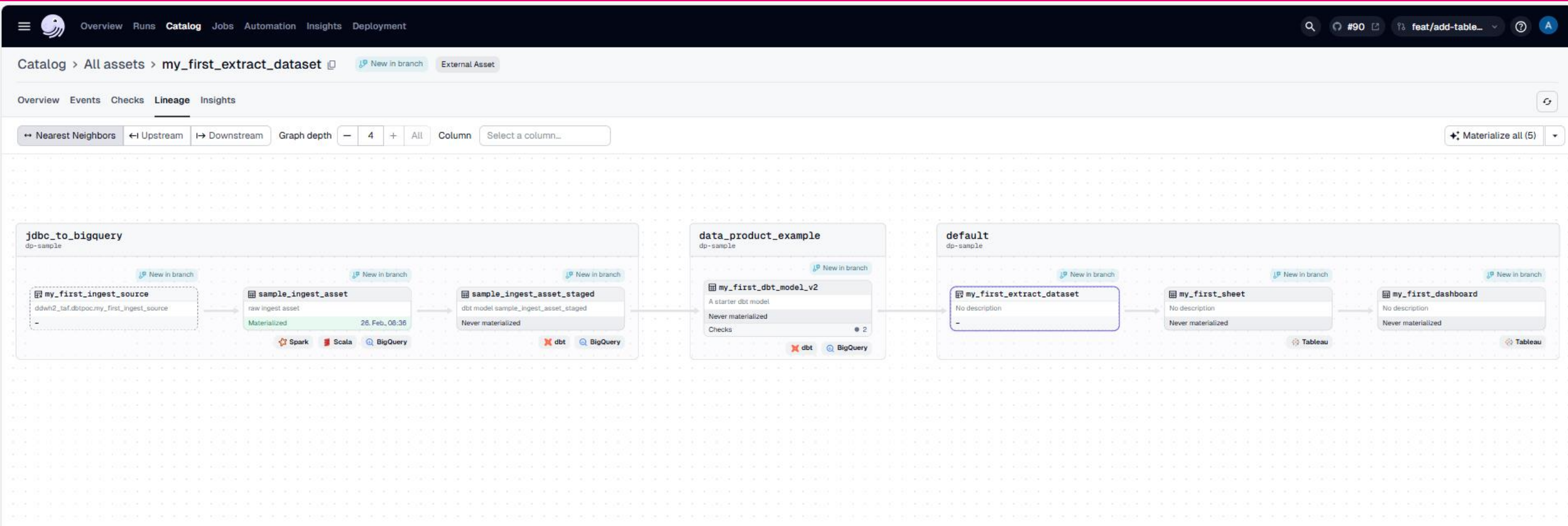APACHE **nifi**

Google BigQuery

DATA INGEST

tableau

# Dagster as the core of the platform



- at Magenta we decided to build around the orchestrator

  - hybrid deployment – controlplane SaaS – runtime in our k8s

  - software engineering best practices for project development and deployment

  - asset-based mindset for data flows (graph like a calculator for data dependencies)
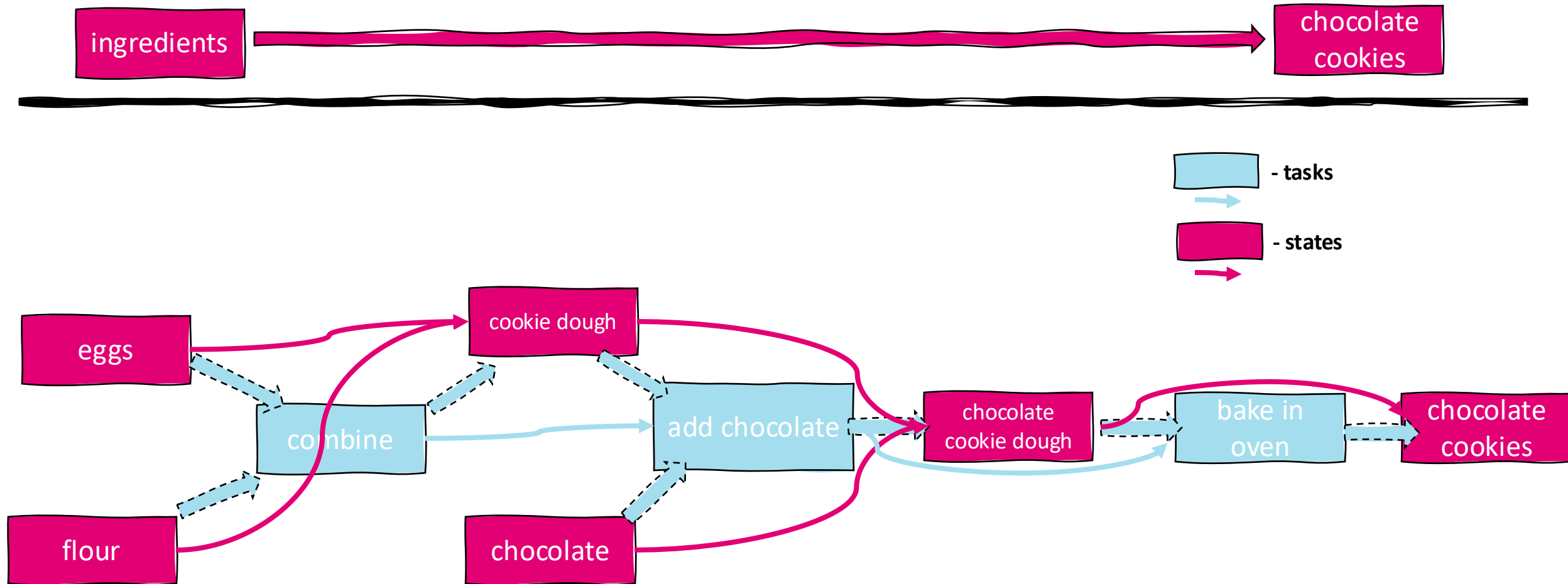
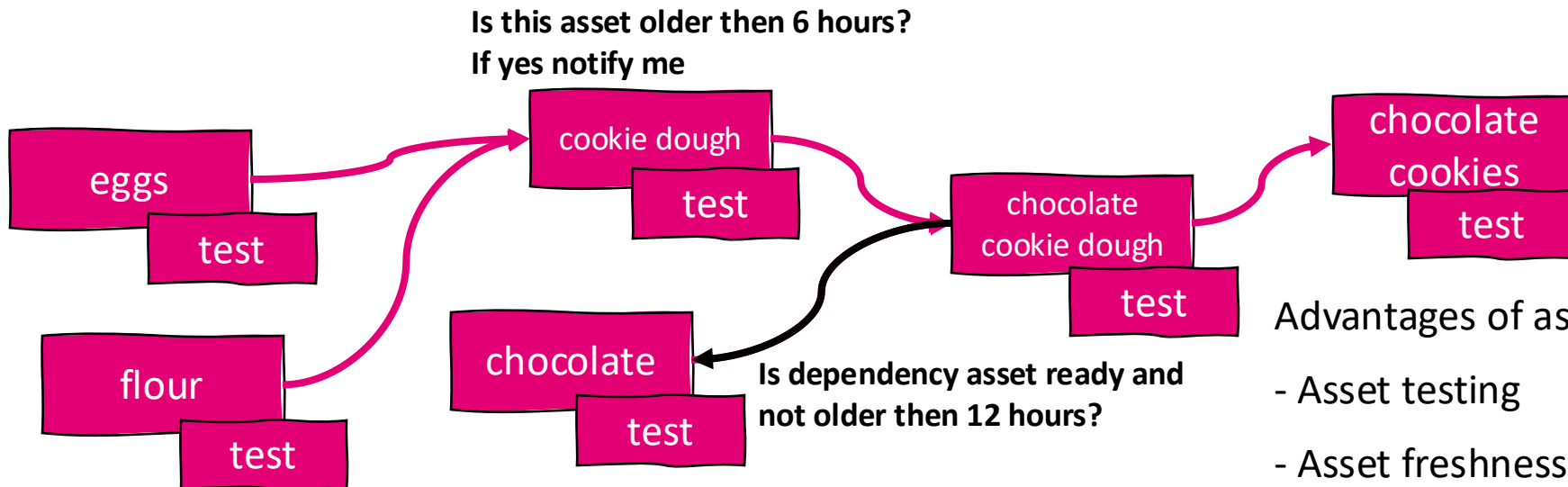- new concepts in orchestration...
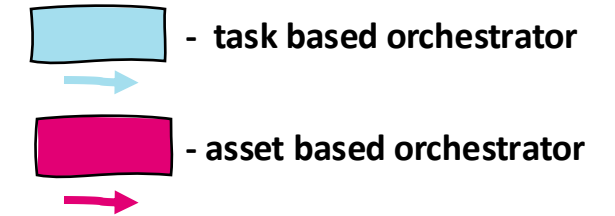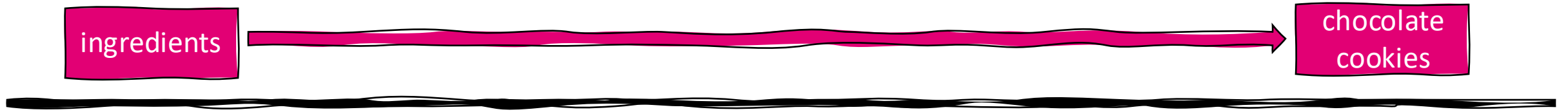
# New enabled concepts

- Asset based graph
- Metadata driven pipeline creation
- Reusable objects (resource, io manager)

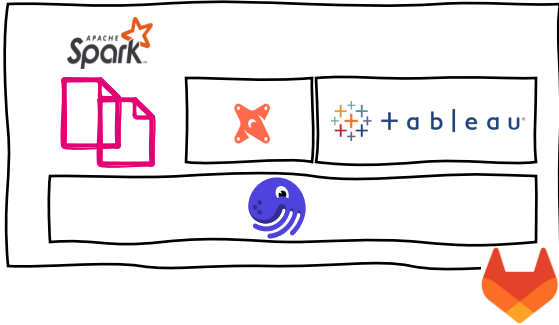# Asset and Task based orchestration: Chocolate cookie example

# Asset based orchestration

ingredients $\longrightarrow$ chocolate cookies

---

combine $\rightarrow$ add chocolate $\rightarrow$ bake in oven

- task based orchestrator

- asset based orchestrator

**Is this asset older then 6 hours?**
**If yes notify me**

eggs

test

cookie dough

test

flour

test

chocolate

test

**Is dependency asset ready and not older then 12 hours?**

chocolate cookie dough

test

chocolate cookies

test

Advantages of asset-based orchestration:

- Asset testing

- Asset freshness

- Asset dependecy graph with granular declarative scheduling approach

# Machine-readable metadata pipeline generation



```
1   metadata = collect_metadata()
2
3   list_of_assets = []
4   for each_metadata_node in metadata:
5       #.. use_metadata
6       @asset
7       def asset():
8           #..use_metadata
9
10      list_of_assets.append(asset)
11
12  Definitions(assets = list_of_assets)
```

```
1   tableau_server = TableauServer(creds)
2   #send rest api call to tableau server and get information
3   tableau_metadata = tableau_server.get_metadata()
4   list_of_assets = []
5   for each_tableau_object in tableau_metadata:
6       tableau_object_deps = each_tableau_object.deps
7       tableau_object_name = each_tableau_object.deps
8       if each_tableau_object.type == extract_datasource:
9           @asset(
10              name = tableau_object_name,
11              deps = tableau_object_deps
12          )
13          def refresh_extract(tableau_server):
14              #send api call to refresh object
15              tableau_server.refresh_extract(tableau_object_name)
16          list_of_assets.append(run_dbt_asset)
17      else:
18          @asset(name = tableau_object_deps)
19          def asset():
20              pass
21          list_of_assets.append(asset)
22
23  Definitions(assets = list_of_assets)
```
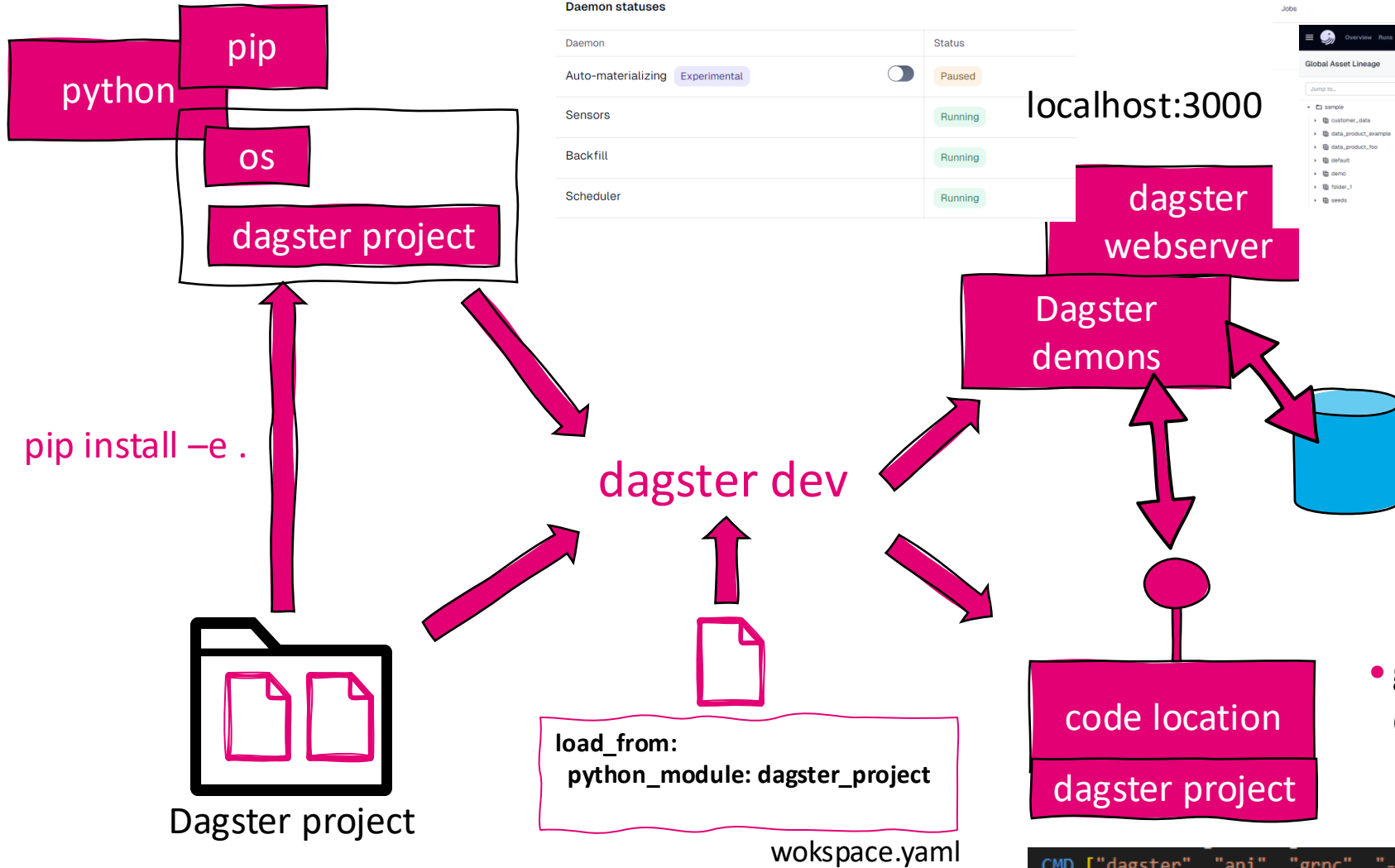
11

# Reusable components

- Reusable objects

  - Resources → Encapsulate complex logic to interact with external systems
  - IO manager → Make complex IO interactions substitutable & testable

- Benefits

  - Dependency injection
  - Day 1 productivity: Scale the data pipeline down to a single laptop
  - Increase self-service: Business/DS focus not required to handle complex IO

```python
@asset(
    io_manager_key="bigquery_io_manager",
)
def awesome_ml_model(context, reference_addresses: pd.DataFrame, bigquery: BigQueryResource) -> pd.DataFrame:
    # simple normal python code here
    # IO is abstracted
    context.log.info(f"from source: \n{reference_addresses.head()}")
    # auth & complexity (imagine web API) is abstracted
    with bigquery.get_client() as client:
        job = client.query("select * from example.upstream")
        query_result = job.result().to_dataframe()
        context.log.info(f"direct query: \n{query_result.head()}")
    return pd.DataFrame({"foo": [1,2,3]})
```

Lab - Demo
- basic asset and UI overview
- resources
- metadata pipeline creation
- components

# dagster - pip install dagster

pip

python

os

dagster project

**Daemon statuses**

| Daemon | | Status |
|---|---|---|
| Auto-materializing | Experimental | Paused |
| Sensors | | Running |
| Backfill | | Running |
| Scheduler | | Running |

localhost:3000



pip install –e .

dagster dev

Dagster project

load_from:
  python_module: dagster_project

wokspace.yaml

dagster webserver

Dagster demons

code location

dagster project

- storage for dagster runs and metadata
  - sqlite/postgres

- grpc server serving the assets to the dagster deamon
- needs definition object from the module

CMD ["dagster", "api", "grpc", "-h", "0.0.0.0", "-p", "4000", "-f", "sample/__init__.py"]

# Takeaways

- Integrated asset-based graph is key (from ingest, transformation, reporting, tests – to AI)
  - Event driven connection
  - Better collaboration (scaling)
- Software engineering principles enable business self service
  - Blueprint
  - Automate all the things: CI/CD (stateful & stateless)
  - DRY: build tested foundation – dependency injection
  - Make business departments part of the key processes and pipelines
- Executable specification (metadata, contracts)
  - Interface Mangement
  - Preserve semantics
  - Preserve compliance (security classification, PII, retention)

Q&A

Dagster @Magenta
Local Modern Data Stack

# Open-source contributions: L-MDS + Tableau integration + VertexAI integration



# Building a data platform is team work

# Further material

- [public live stream about the basics: georgheiler.com/event/magenta-pixi-25](https://georgheiler.com/event/magenta-pixi-25)

- n-depth Magenta implementation [georgheiler.com/event/magenta-data-architecture-25](https://georgheiler.com/event/magenta-data-architecture-25)

- [https://courses.dagster.io/](https://courses.dagster.io/)