

OSA CON 25



Scaling Data Pipelines @Magenta Telekom

Georg Heiler

Senior Data Expert
geoheil.com

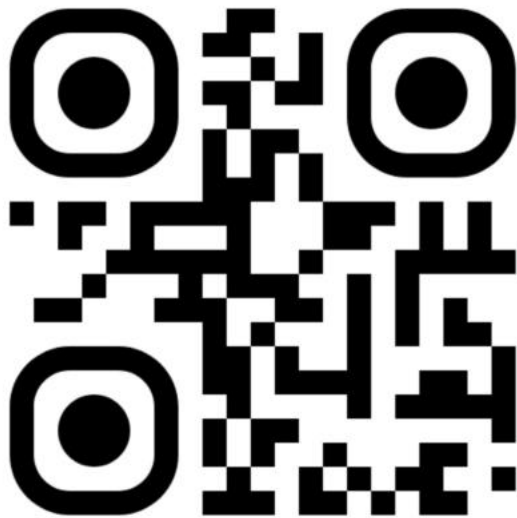
November 4-5, 2025

Georg Heiler

Solving challenges with data.

- Senior data expert @Magenta
- Research Software Engineer @ASCII & CSH
- Meetup organizer & frequent Speaker

geoheil.com, linkedin.com/in/geoheil



There is a chaos out there



How did we end here? Time!

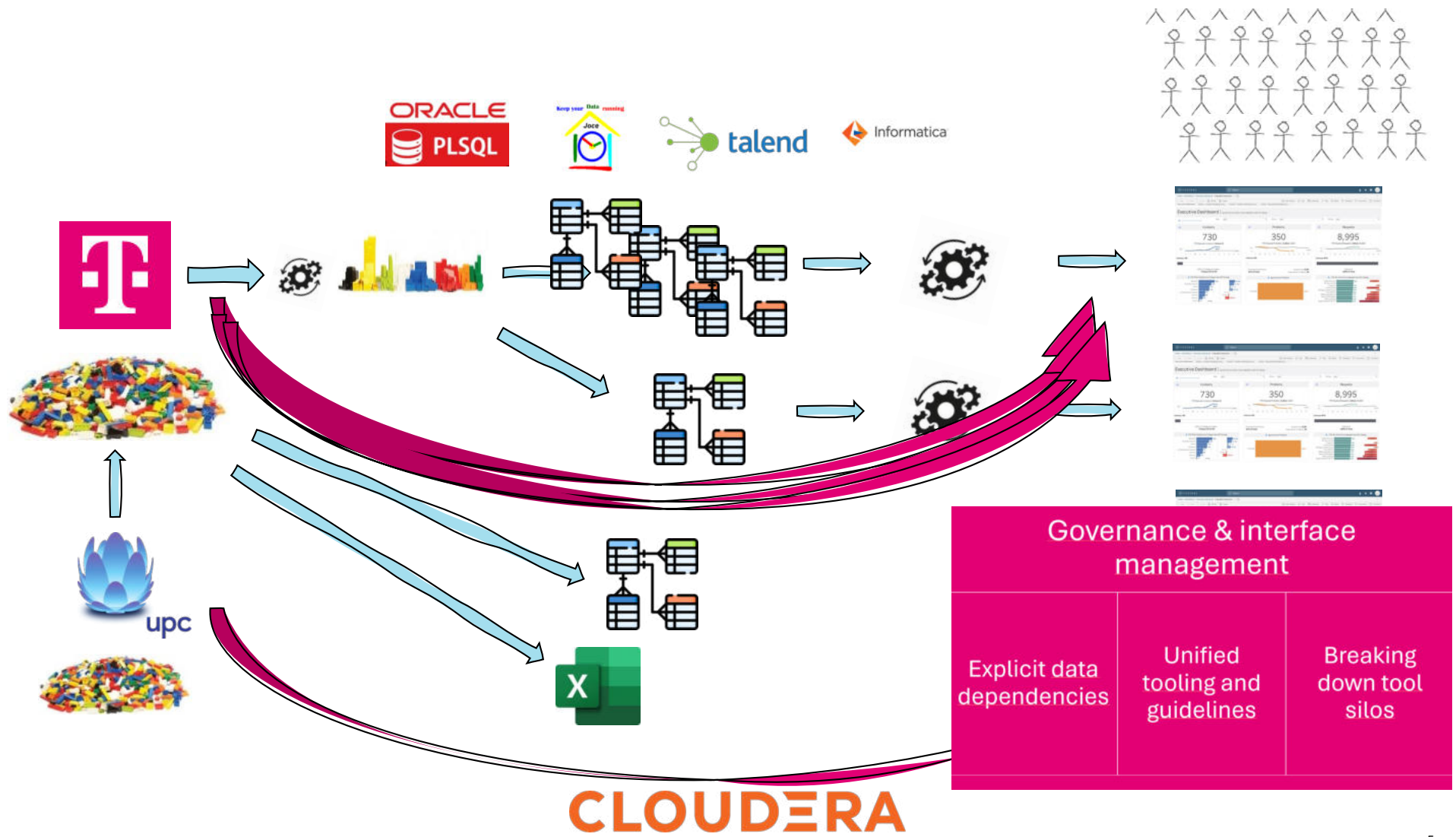
business grows (merger)
demand for data grows
methodology and tooling changes

- Missing lineage
- Missing semantics
- Missing collaboration
- High lead times
- Limited quality

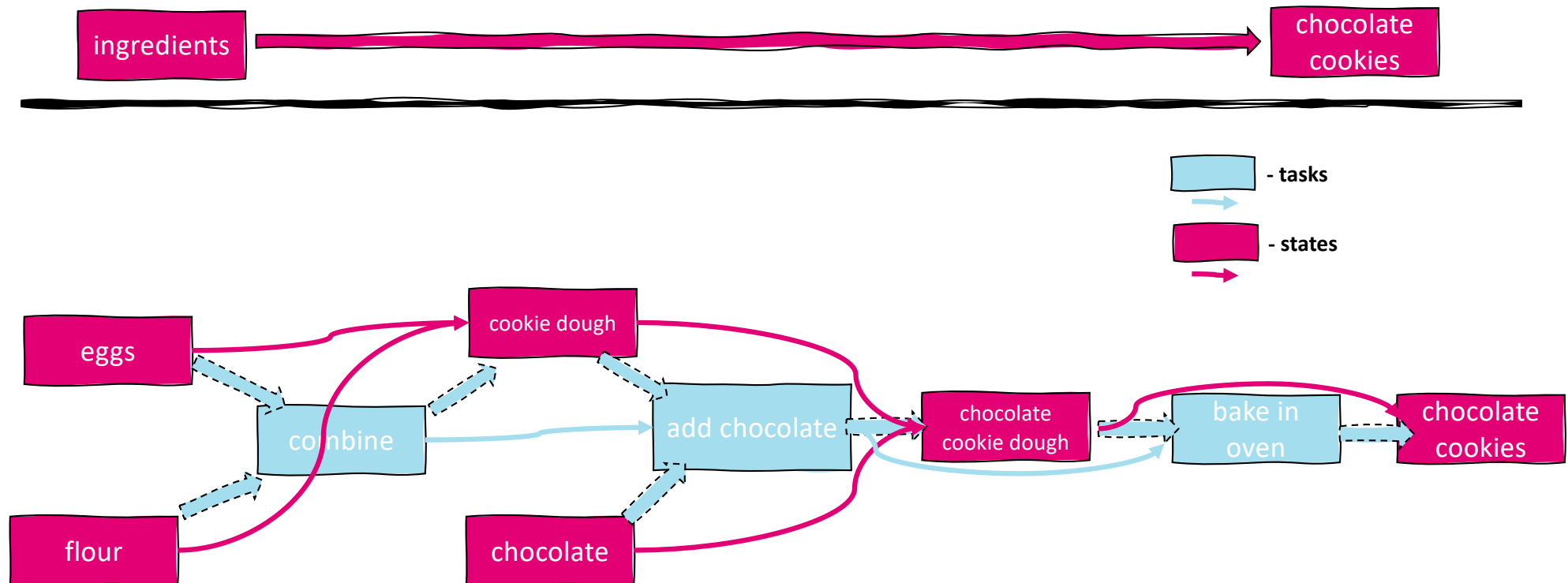
“[...] model behavior is not determined by architecture, hyperparameters, or optimizer choices. It's determined by your dataset, nothing else.”

James Betker
Research Engineer, Open AI

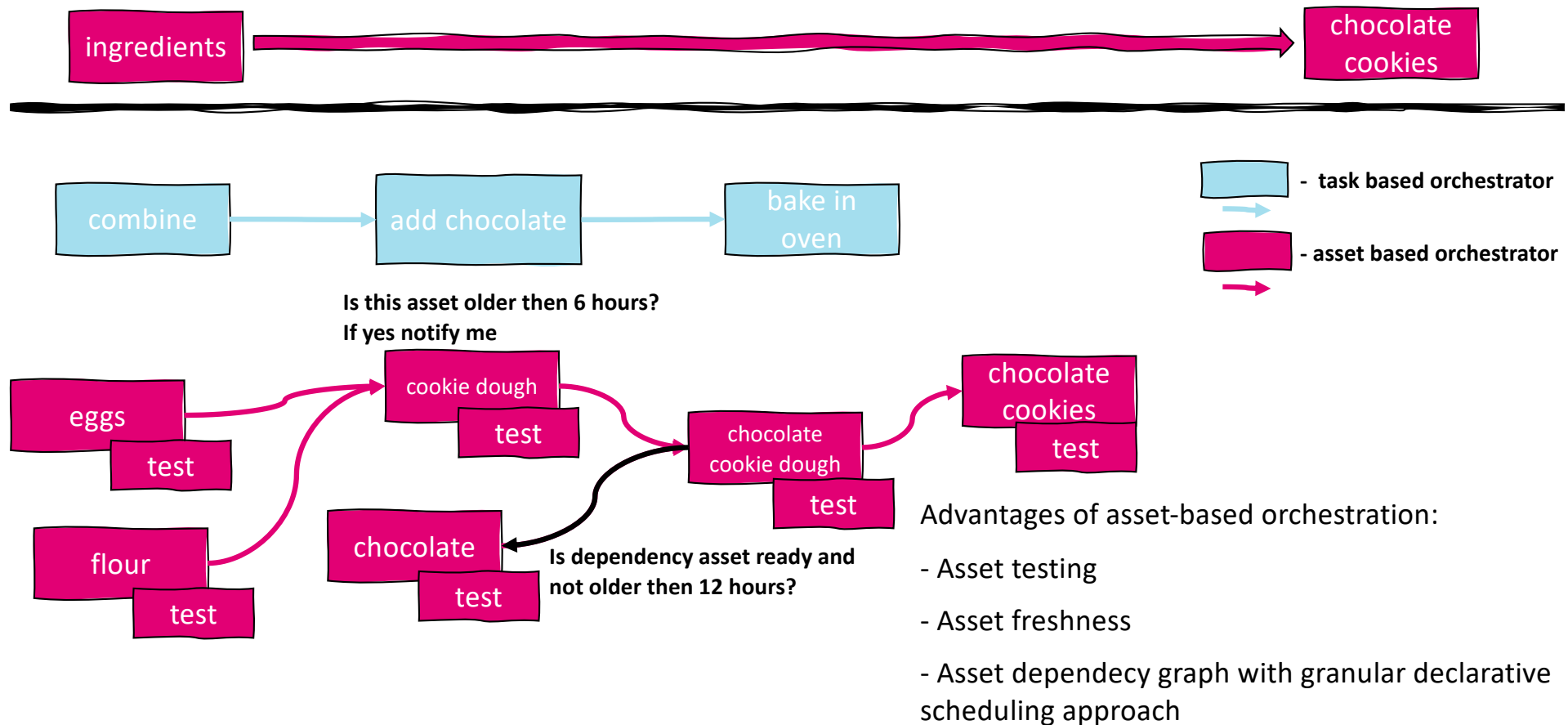
<https://www.youtube.com/watch?v=lvhtTu9CTAU>



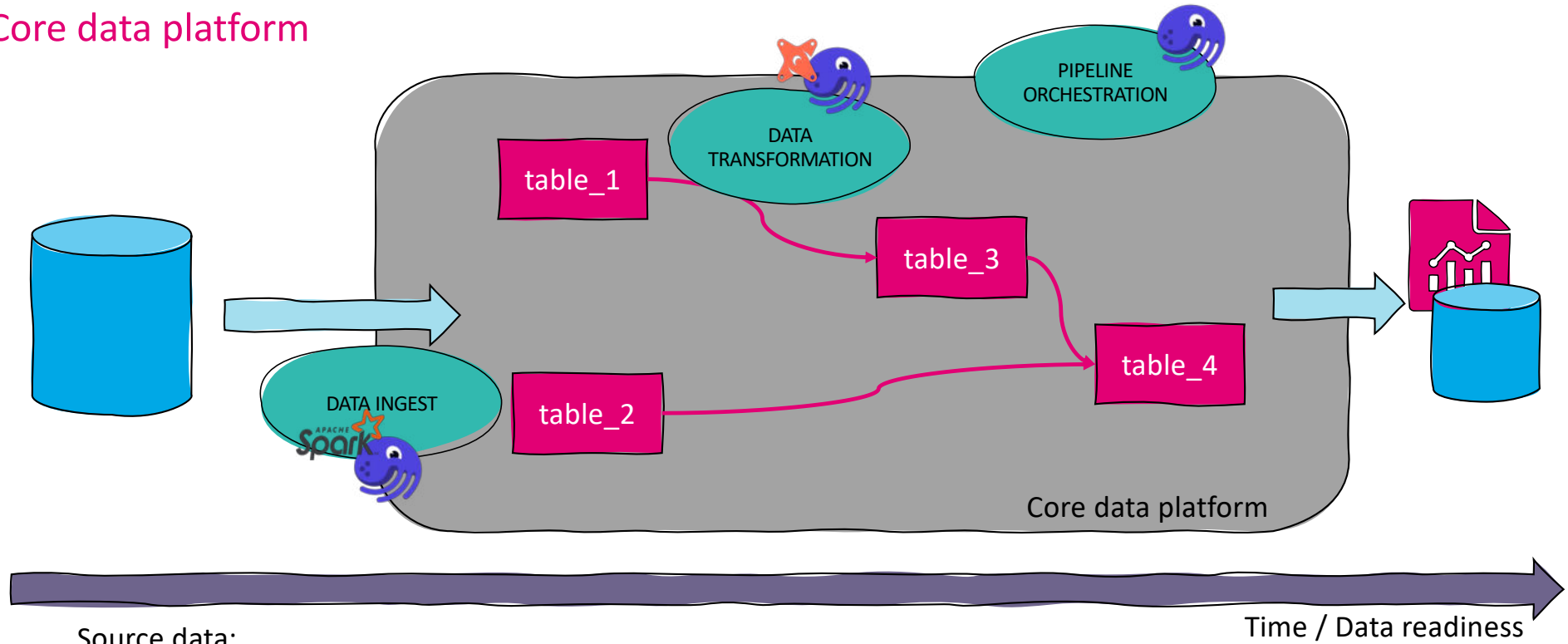
Asset and Task based orchestration: Chocolate cookie example



Asset based orchestration



Core data platform



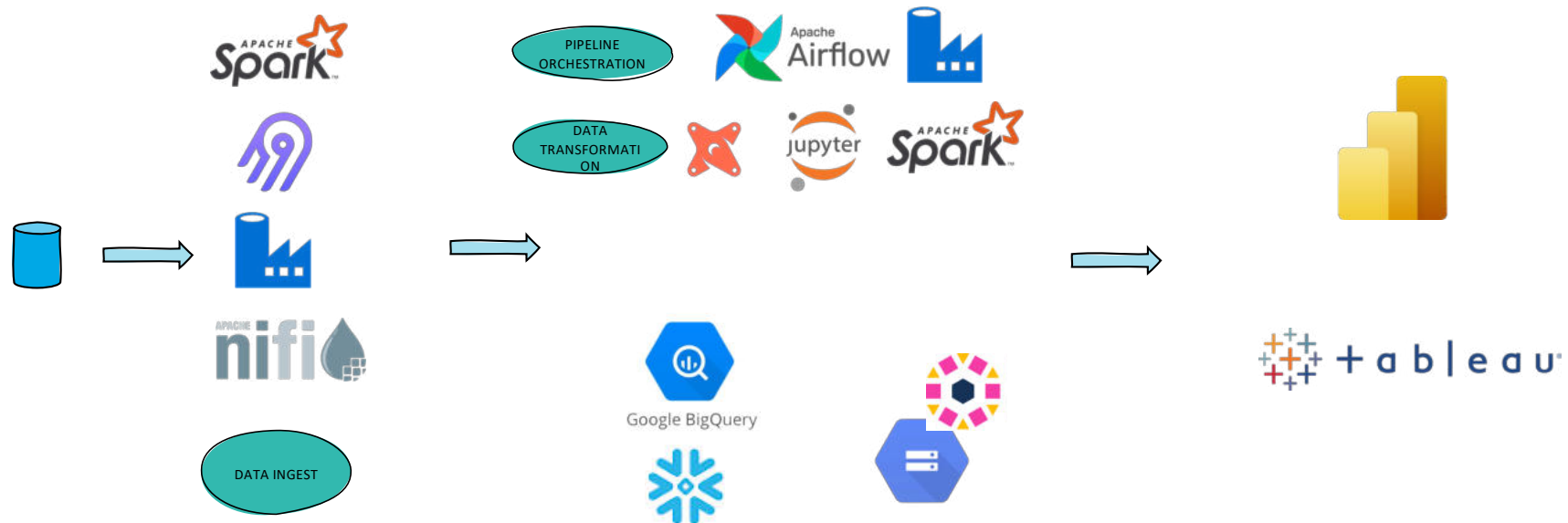
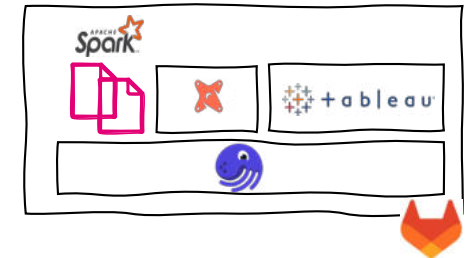
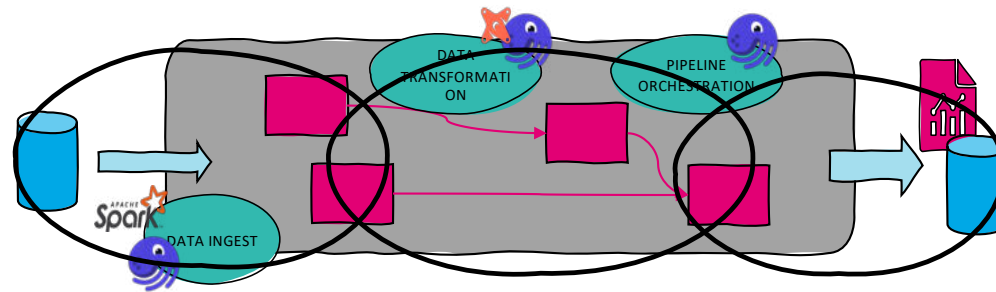
Source data:

- Kafka
- Files
- Database systems

Understanding tool silos

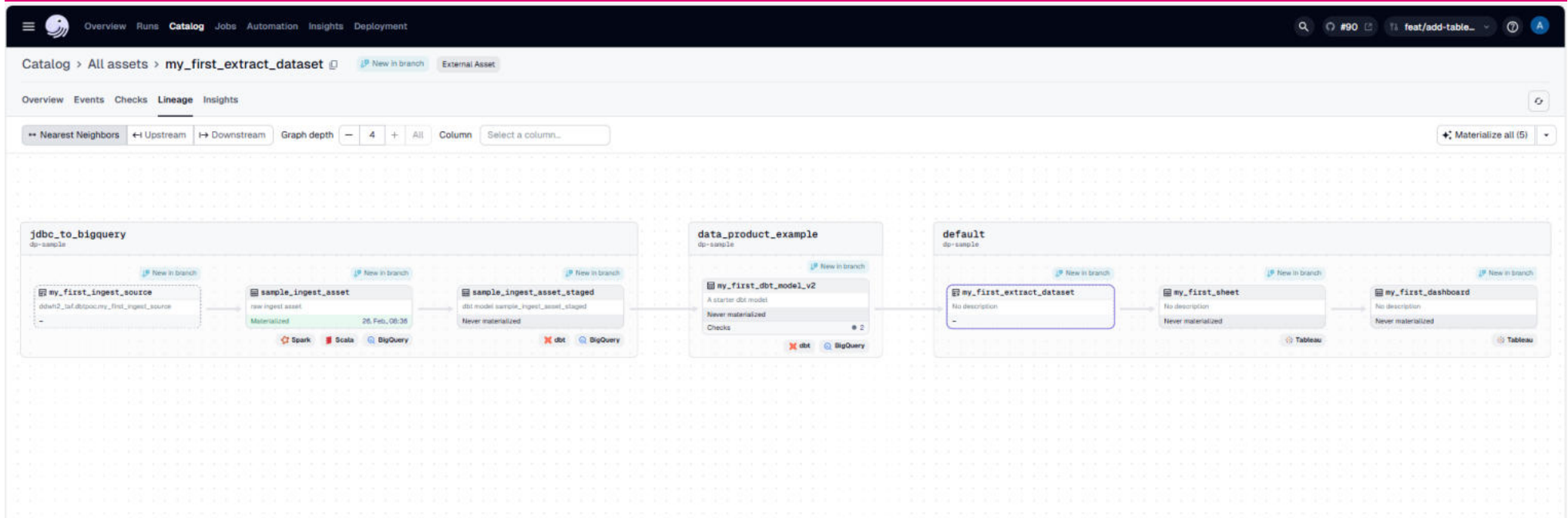
What should i do to get E2E reporting use case done?

Developer



New enabled concepts

- Asset based graph
- Metadata driven pipeline creation
- Reusable Components
-



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 conf
2 list
3 for
4
5
6
7
8
9
10
11
12
13
14 Defi
```

```
1 m
2 l
3 f
4
5
6
7
8
9
10
11
12
13
14
15 D
```

```
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)
```

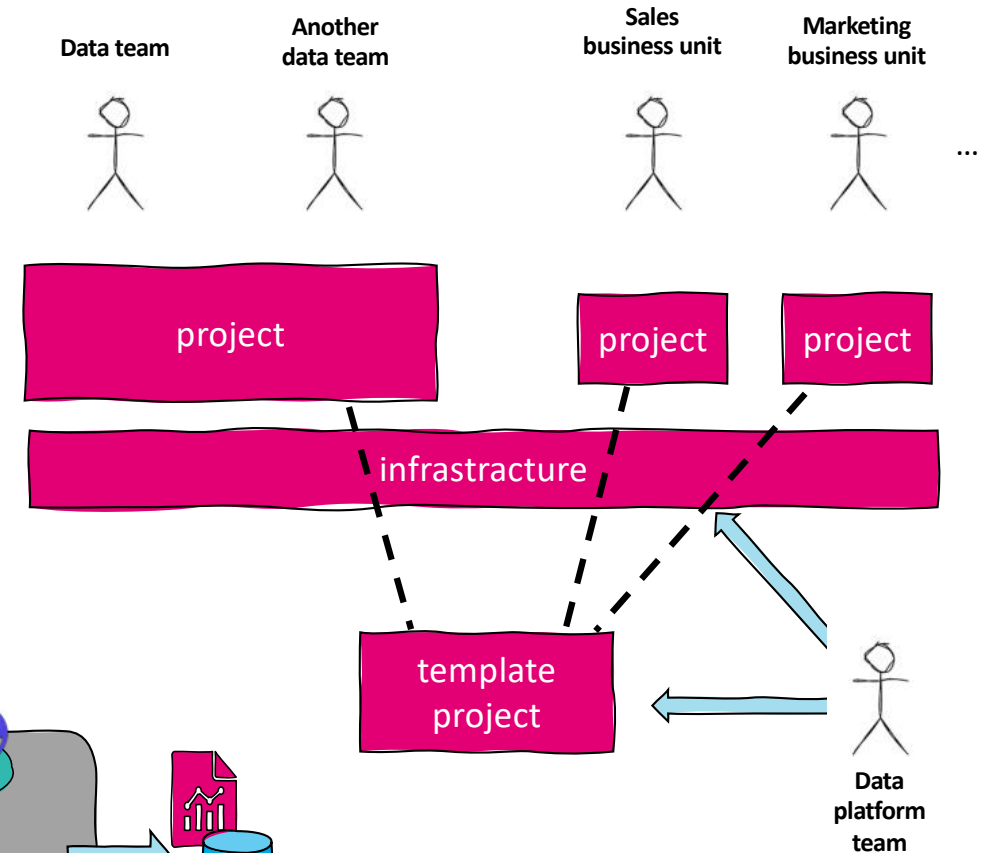
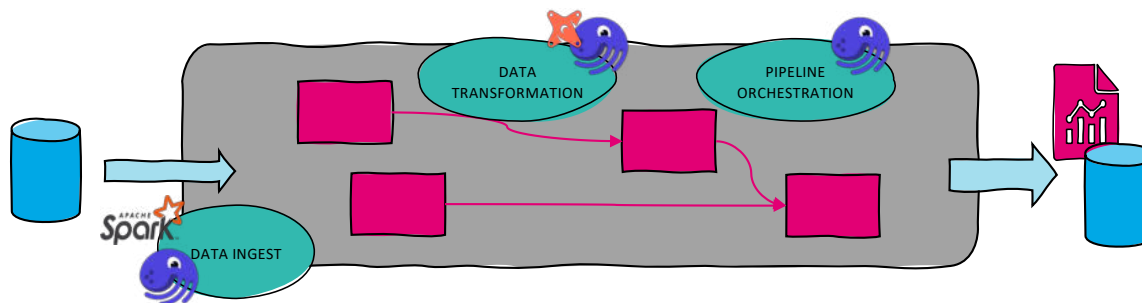
Reusable components

- Define once, test & reuse
- 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

```
2 @asset(  
3     io_manager_key="bigquery_io_manager",  
4 )  
5 def awesome_ml_model(context, reference_addresses: pd.DataFrame, bigquery: BigQueryResource) -> pd.DataFrame:  
6     # simple normal python code here  
7     # IO is abstracted  
8     context.log.info(f"from source: \n{reference_addresses.head()}")  
9     # auth & complexity (imagine web API) is abstracted  
10    with bigquery.get_client() as client:  
11        job = client.query("select * from example.upstream")  
12        query_result = job.result().to_dataframe()  
13        context.log.info(f"direct query: \n{query_result.head()}")  
14    return pd.DataFrame({"foo": [1,2,3]})
```


Observation

- Process is straight forward: ingest, transform, use
- Everything we do - we do for business to provide better service
- Hard to scale across company
- Divide people into **develop framework** and **use framework** groups
- Thinking in **building blocks**
- Tooling supporting software engineering practices: **dbt**, **dagster**, **pixi**, **docker**
- Introduction of **new processes, modeling** and **metadata tooling** for better governance



Takeaways

- Integrated asset-based graph is key (from ingest, transformation, reporting, tests – to AI)
 - Event driven connection
 - Better collaboration (scaling)
 - Can enable execution environment re-targeting in advanced cases
- 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)

Explore for yourself!

The collage consists of three overlapping images. On the left is a screenshot of the GitHub repository 'local-data-stack' by 'l-mds'. It shows the file tree with folders like 'doc', 'img', and 'cookiecutter.project_slug', and files like '.gitignore', 'README.md', 'cookiecutter.json', 'pixi.lock', 'pyproject.toml', and 'yamllintconfig.yaml'. In the center is a screenshot of a GitHub pull request titled '[dagster-tableau] Exploring embedded data sources #27218' by 'maximearmstrong'. It shows a 'Summary & Motivation' section with two points and a 'How I Tested These Changes' section. On the right is a screenshot of a YouTube video titled 'Pixi powering Telekom data cloud' by 'Georg Heiler'. The video features a cartoon character holding a star and the text 'Pixi powering Telekom data cloud'.

[Scaling data pipelines @Telekom](#)

[Pixi powering Telekom data cloud](#)

[Declarative Execution](#)

[In-depth explanation](#)

[dagster-ducklake](#)











[dagster-vertexai](#)

Data platform is team work and
we are very proud and excited about the journey ahead

OLAP ELO Ranking

- Relative
- Continuously improving
- Robust

georgheiler.com/post/elo-data-challenge

1	 *StarRocks	2017	
2	 *Kinetica	2015	
3	 *Apache Hudi	2012	-2
4	 *Databricks SQL	2005	-5
	 *DuckDB	2004	-20
	 *Clickhouse	2004	+19
	 *Delta Lake	1985	-2
8	 *TrinoDB	1985	
9	 *Snowflake	1981	-5
10	 *Apache Spark	1975	