# Scaling data pipelines @Telekom

Aleks & Georg

# Agenda

# About us

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

- meetup organizer and conference speaker

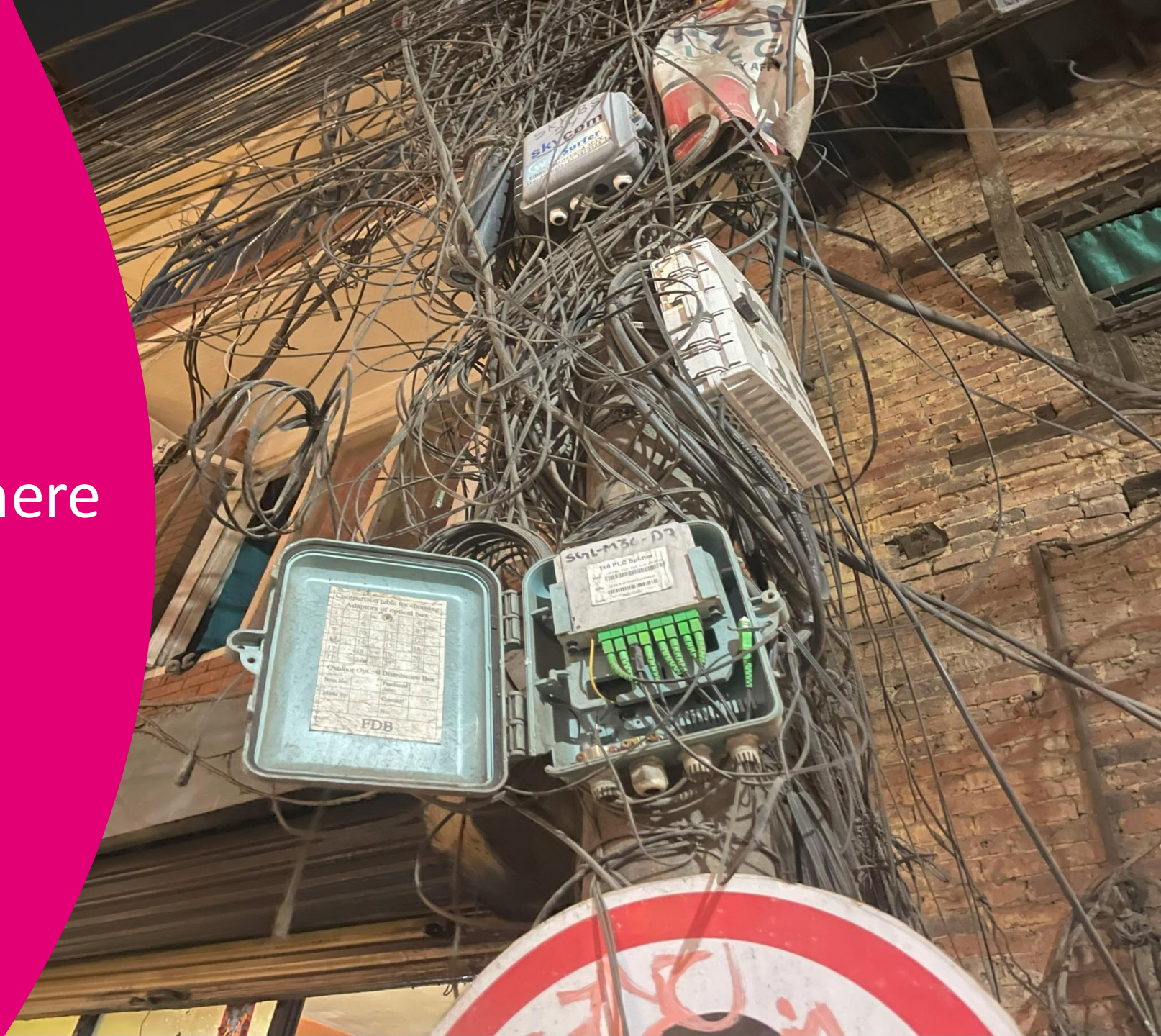- data architecture, multimodal and complex data challenges

🔗 **geoheil**

⬛ **geoheil**

🦋 **@geoheil.com**

**Data engineer** at Magenta Telecom

- database internals, data platform engineering

🔗 **milicevica23**
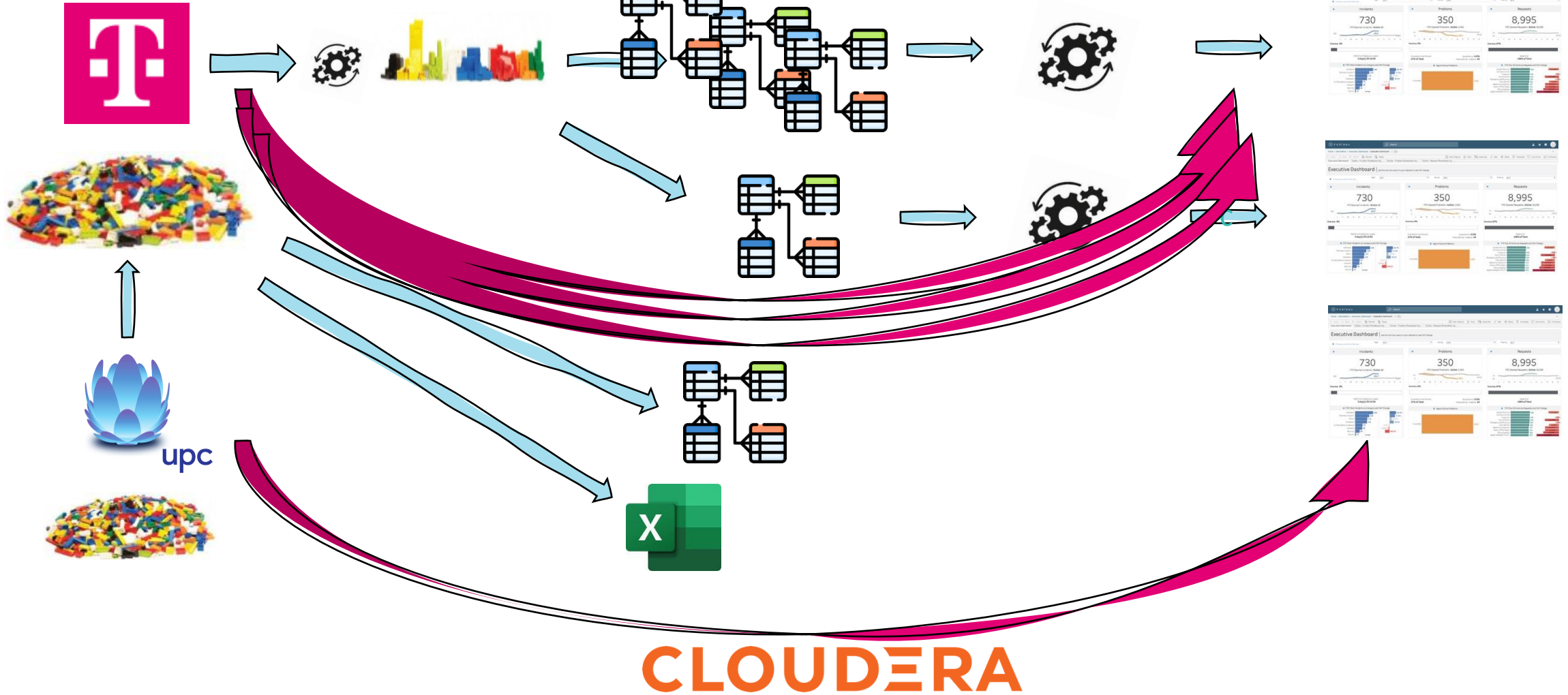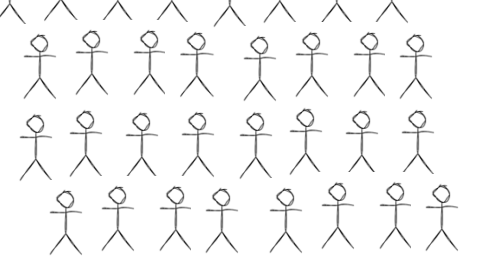
⬛ **milicevica23**

🦋 **@milicevica23.bsky.social**

public live stream about the basics: georgheiler.com/event/magenta-pixi-25
https://yam-united.telekom.com/pages/data-magenta/apps/blog/updates/view/1496666c-c479-4e09-8ef5-0ebfe709e179

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
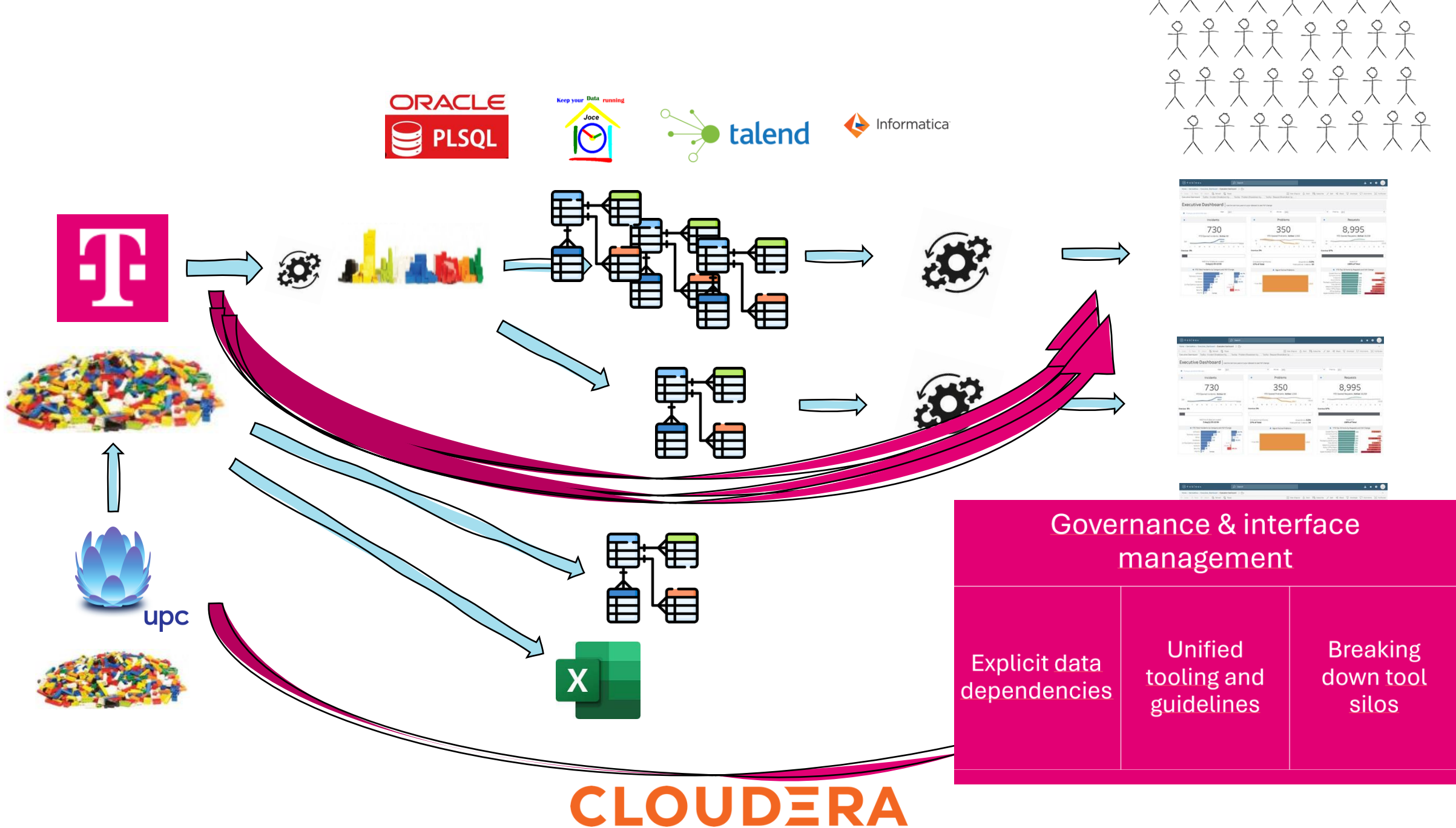- High lead times
- Limited quality
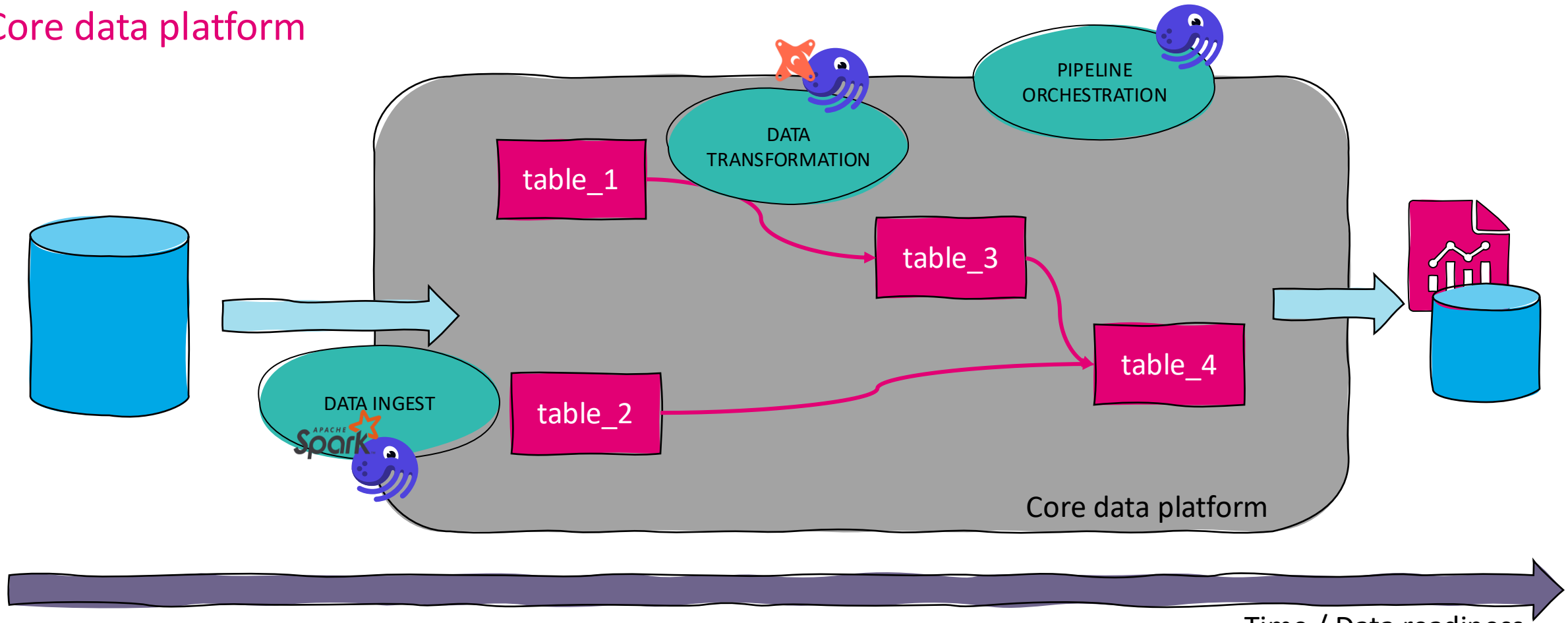
# Governance & interface management

| Explicit data dependencies | Unified tooling and guidelines | Breaking down tool silos |

ORACLE PLSQL

Keep your Data running — Joce

talend

Informatica

Governance & interface management

| Explicit data dependencies | Unified tooling and guidelines | Breaking down tool silos |
|---|---|---|

Executive Dashboard

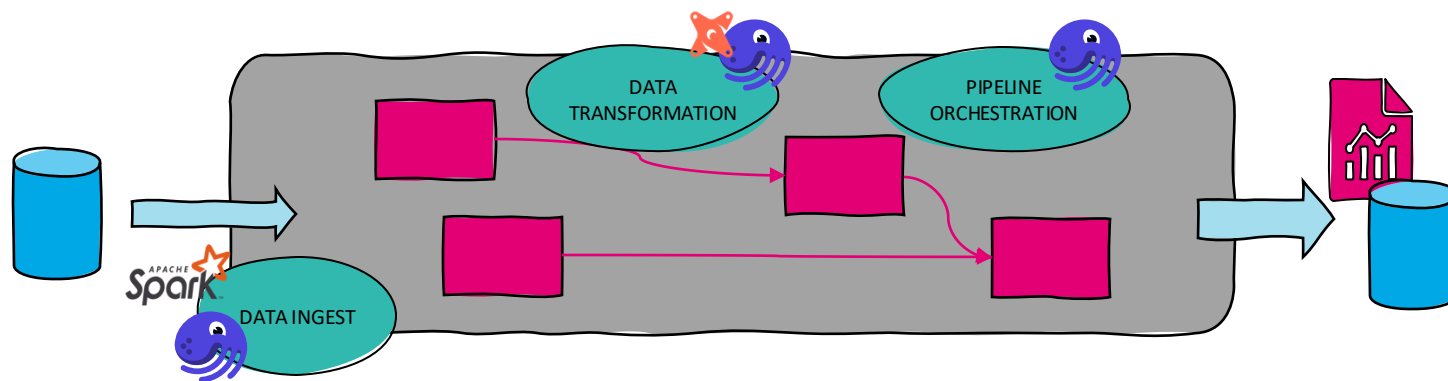| Incidents | Problems | Requests |
|---|---|---|
| 730 | 350 | 8,995 |

upc

CLOUDERA

# Core data platform



Source data:

- Kafka

- Files

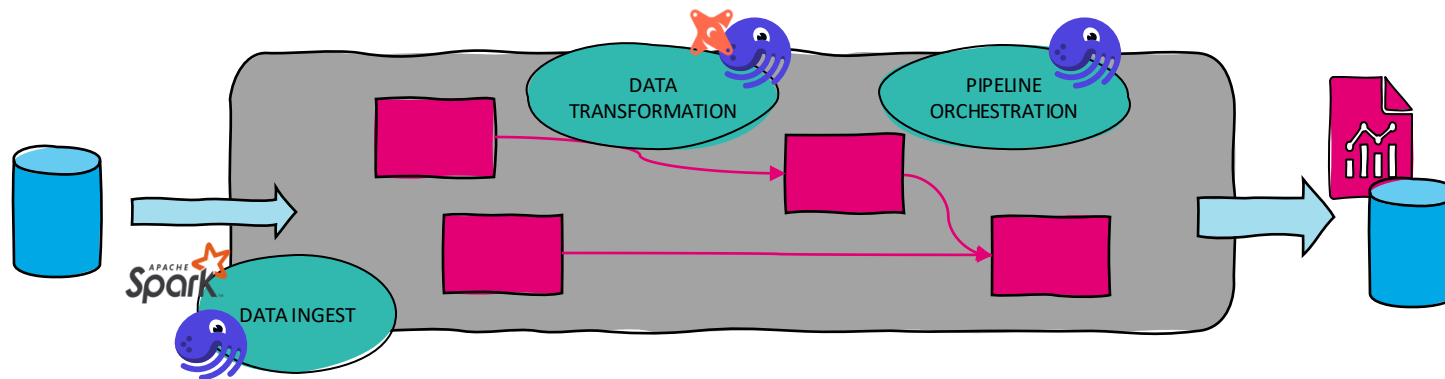- Database systems

Time / Data readiness

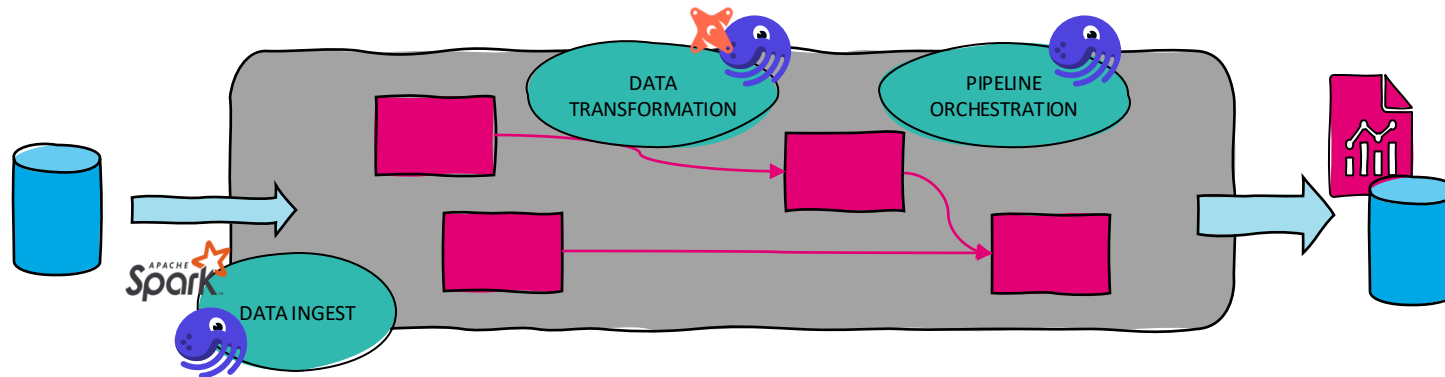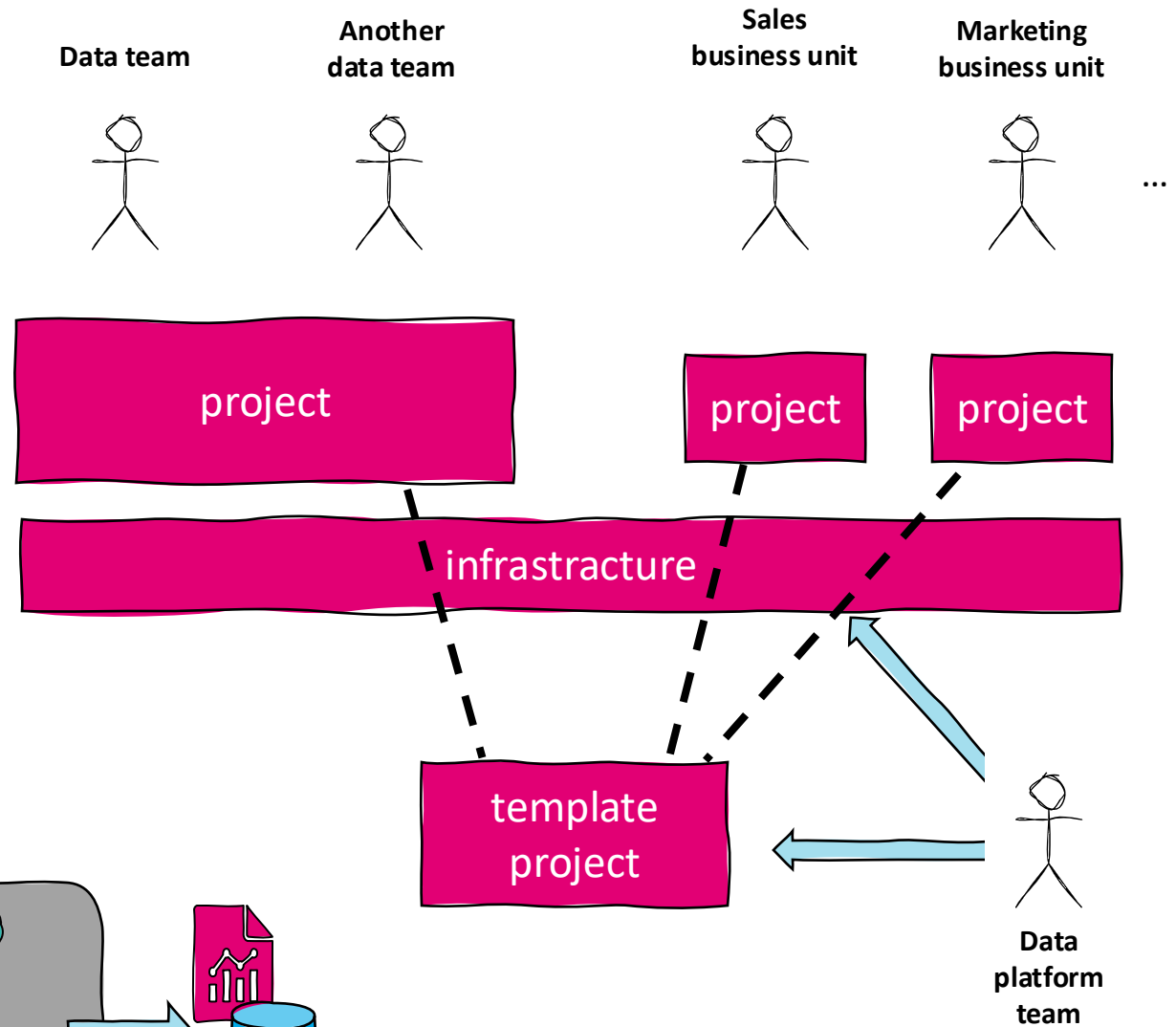# I fear that what we build is very hard to push into the business units

know easy
I ~~fear~~ that what we build is very ~~hard~~ to push into the business units

# Observation

- Process is straight forward: ingest, transform, use

- Everything we do - we do for business to provide better service

- Hard to scale across company

- Dividing people into **develop framework** and **use framework** groups

- Thinking in a **building block** structure

- Introduce modern tooling supporting software engineering practices: **dbt, dagster, pixi, docker**

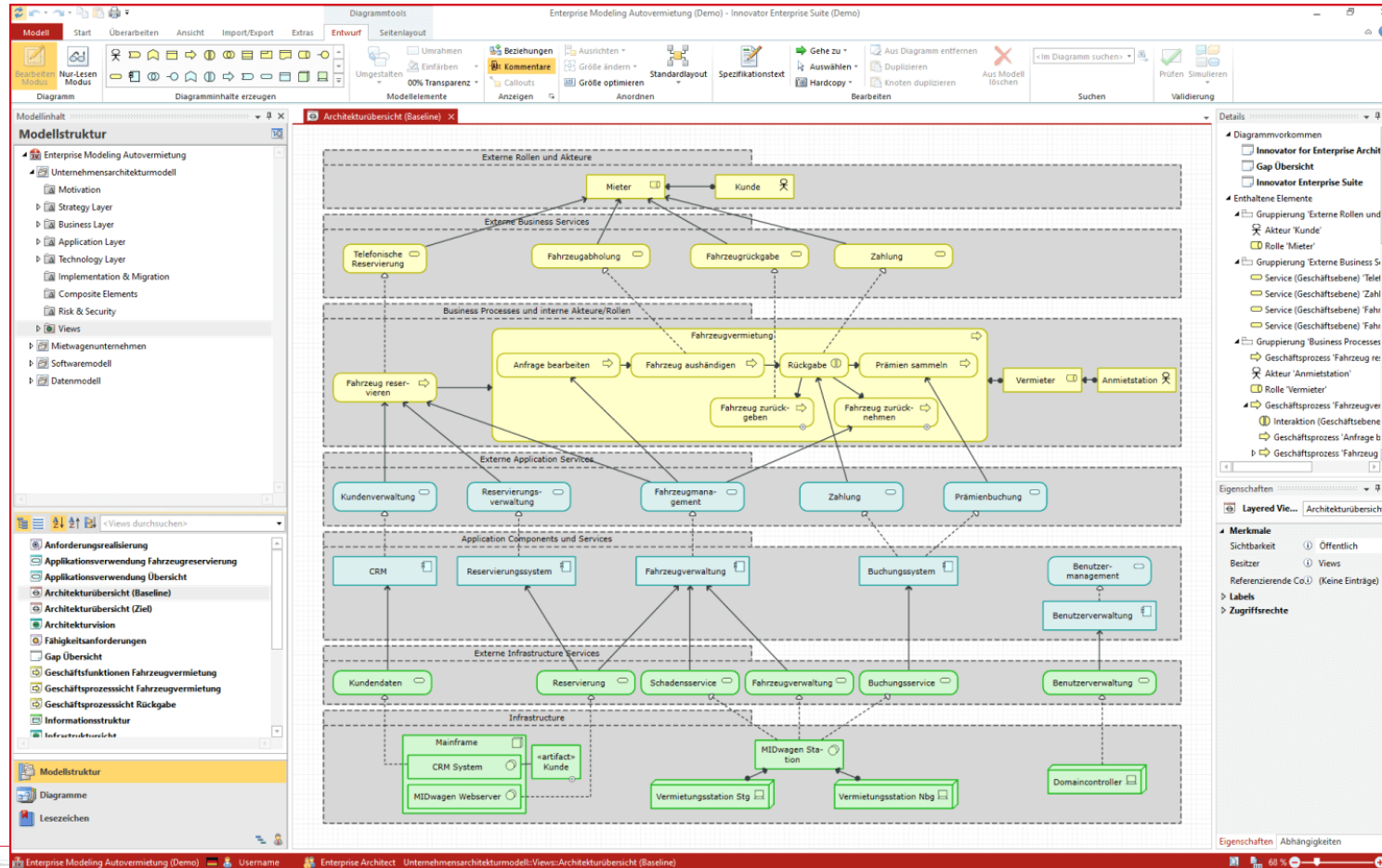- Introduce **new processes**, **modeling** and **metadata tooling** for better governance
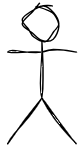
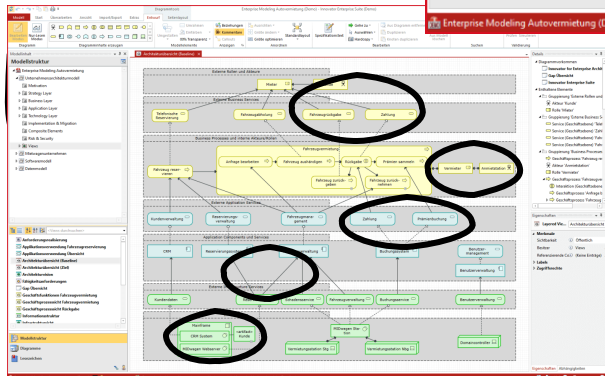# Having a data platform team doesn't mean that your platform scales

1.  Building block with governance, modeling and software-engineering principles
2.  Understanding data platform vendor war
3.  Understanding tool silos

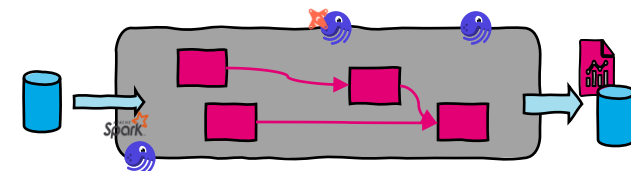# Break the silos with the building block
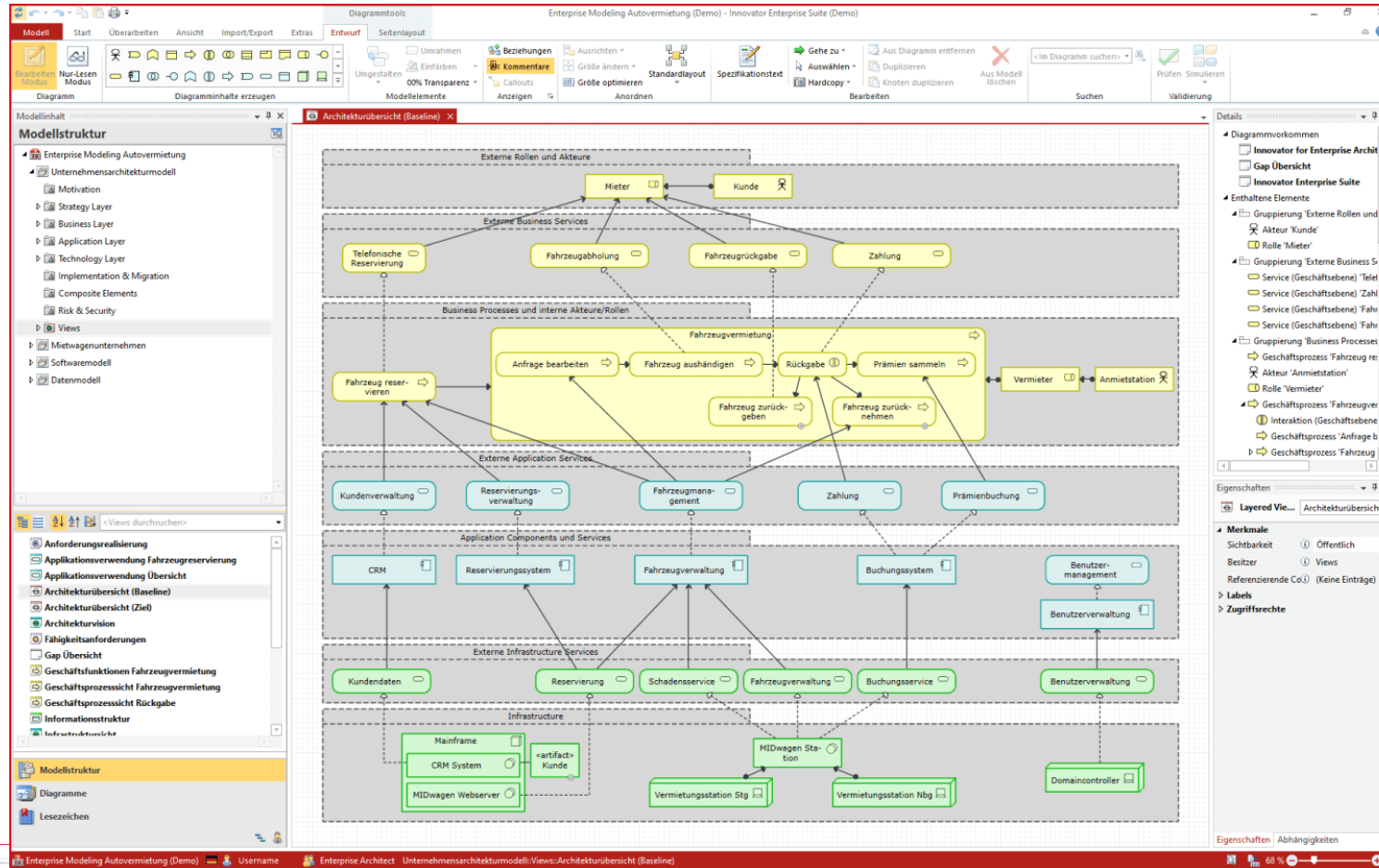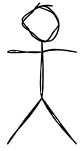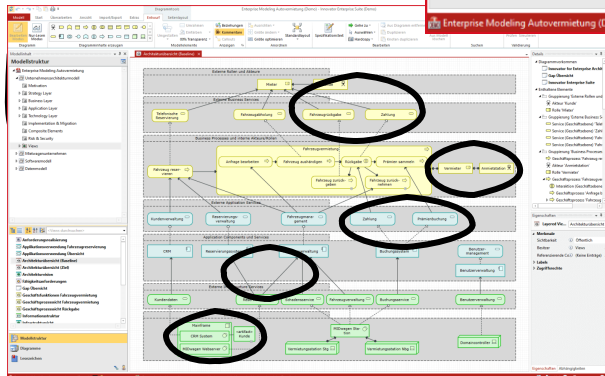
**Data team**



**Sales business unit**

infrastracture

template project

# Break the silos with the building block

**Data team**



**Sales business unit**
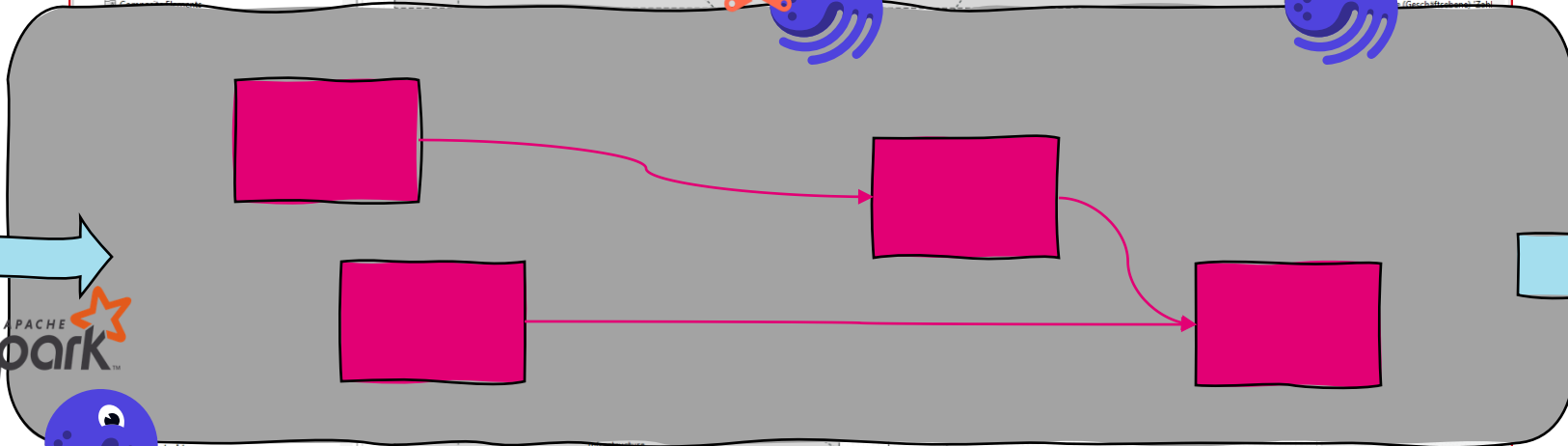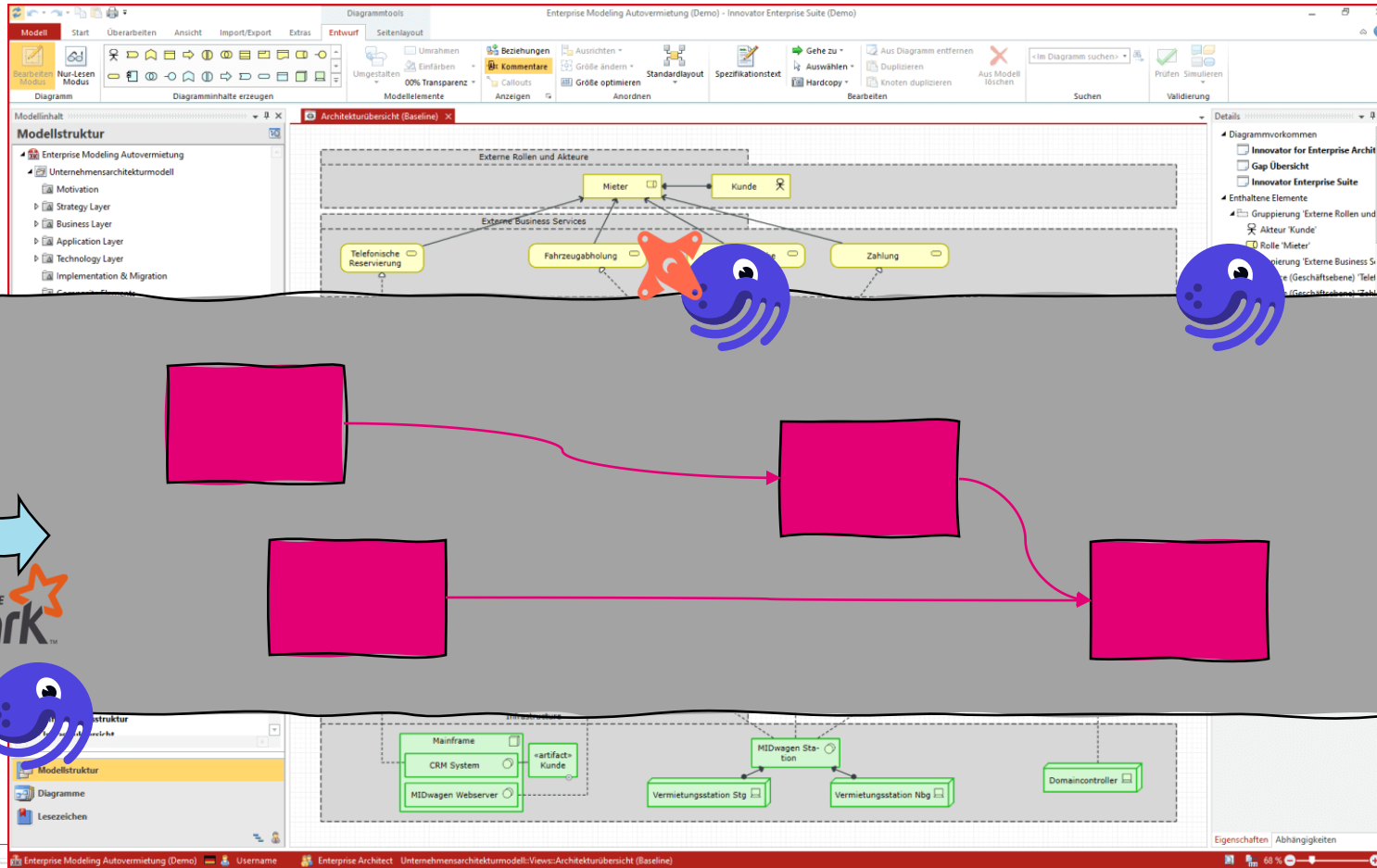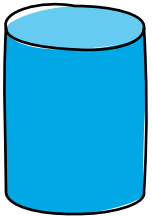
infrastracture

template project

# Break the silos with the building block



**Data team**

**Sales business unit**

infrastracture

template project

16

# Break the silos with the building block
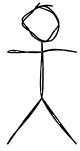
**Data team**
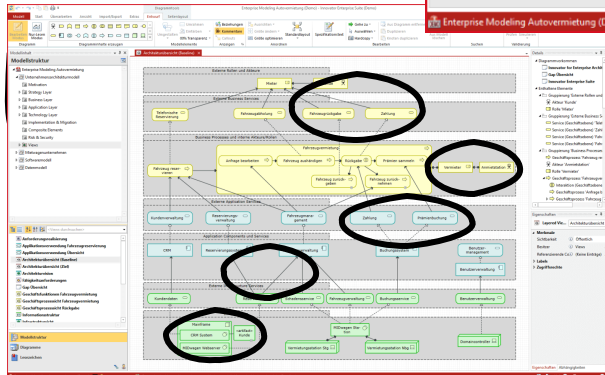
**Sales business unit**

infrastracture

template project

# Break the silos with the building block

**Data team**

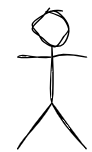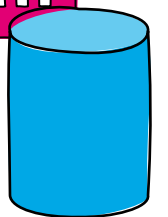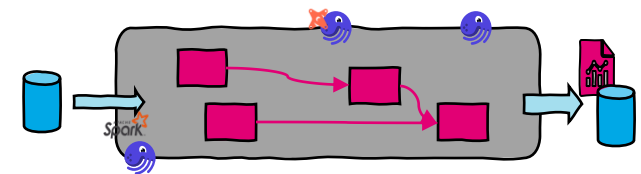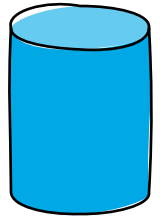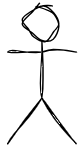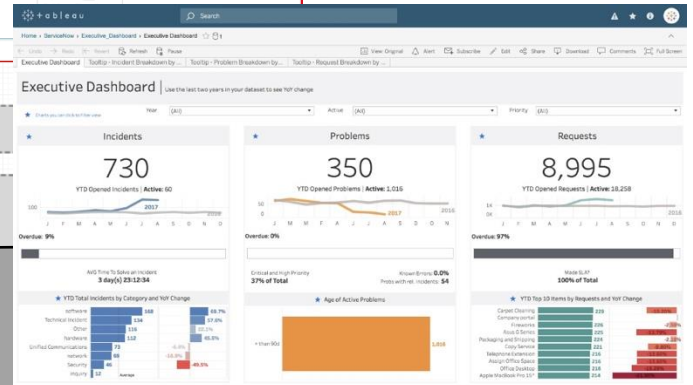**Sales business unit**

infrastracture

template project

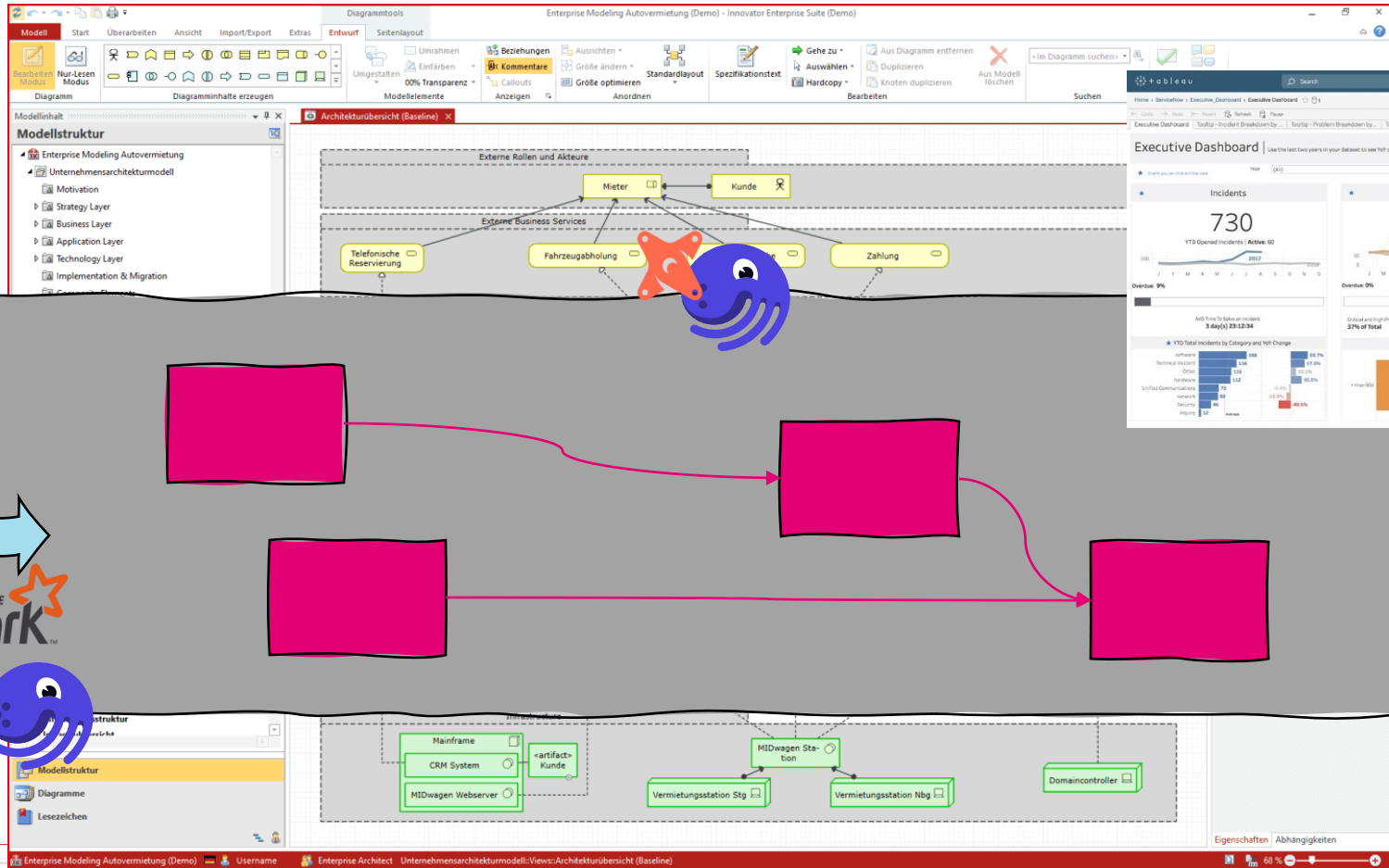# Break the silos with the building block

**Data team**



**Sales business unit**

infrastracture

template project

19

# Multi-project setup is a challenge!

1. Dedicated team maintaning infrastructure and template project
2. Governance and modeling tools

# Modeling and governance are keys for success

**Development phase**

**Exploration phase**

me, you

infrastructure

# Understanding data platform vendor war



THE 2024 MAD (MACHINE LEARNING, ARTIFICIAL INTELLIGENCE & DATA) LANDSCAPE

# execution engine

Cloud:
Engine:
DevEnv: jupyter | jupyter | jupyter
Orc:

Google BigQuery

- **Full stack offered by one vendor**
- **E2E integration lock in**
- **Deployment is always a workaround**
- **No SWE, no local development**
- **Orchestrator is second class citizen and always task based**
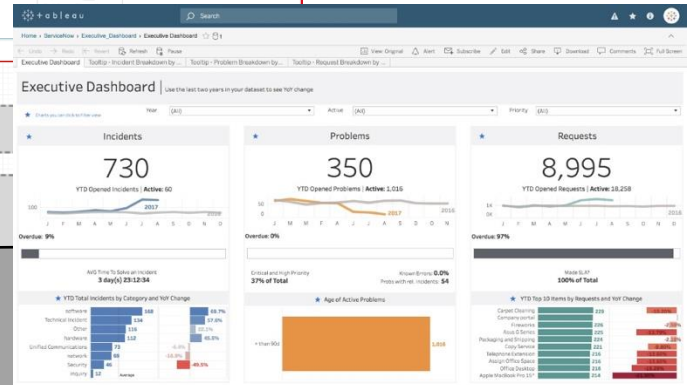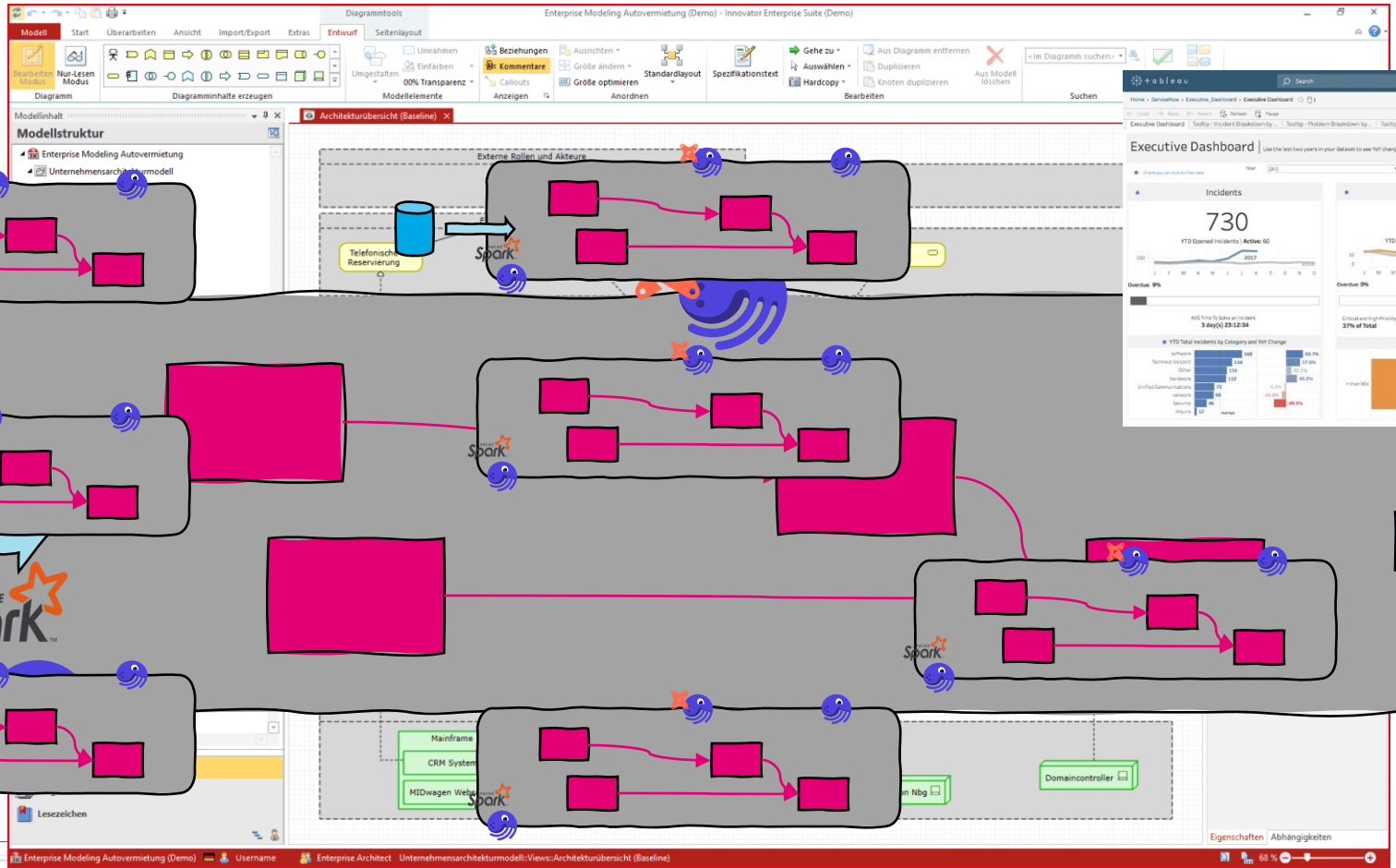
- **+ Frontier in the lakehouse approach**
- **+ Notebooks environment is very convenient**
- **+ Everything on one place**

# sql transformation framework

dbt | SQLMesh

- **Orchestration is just for DWH/SQL part of the platform**

- **+ Frontier in the SQL development**
- **+ SWE for DWH development**

# orchestration engine

dbt

MAGE

- **More technical knowledge needed to setup and use correctly**

- **+ Integration is important**
- **+ full SWE and local development**

# Understanding tool silos

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

**Developer**

DATA INGEST

DATA TRANSFORMATION

PIPELINE ORCHESTRATION

PIPELINE ORCHESTRATION — Apache Airflow

DATA TRANSFORMATION — Jupyter — Apache Spark

DATA INGEST

Apache Spark

Apache nifi

Google BigQuery

tableau

# Dagster as the core of the platform



- at Magenta we decided to build around the orchestrator and not around a execution engine
  - 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 Components

- ....

# Asset and Task based orchestration: Chocolate cookie example

# Asset based orchestration



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

Is dependency asset ready and
not older then 12 hours?

- task based orchestrator
- asset based orchestrator

Advantages of asset-based orchestration:

- Asset testing

- Asset freshness

- Asset dependecy graph with granular declarative scheduling approach

# New enabled concepts

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

- ....

# Machine-readable metadata pipeline generation

# Machine-readable metadata pipeline generation

# 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   configuration_files = read_ingest_configuration_folder(path)
2   list_of_assets = []
3   for each_config_file in configuration_files:
4       config = parse_config(each_config_file)
5
6       @asset(
7           name = config.name
8       )
9       def ingest_asset():
10          df = spark.read(config.source)
11          df.write(config.source)
12      list_of_assets.append(ingest_asset)
13
14  Definitions(assets = list_of_assets)
```

# 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   manifest = read_dbt_manifest(path)
2   list_of_assets = []
3   for each_node in manifest:
4       model_name = each_node.name
5       model_deps = each_node.deps
6       @asset(
7           name = model_name,
8           deps = model_deps
9       )
10      def run_dbt_asset(dbt: DbtClient)
11          dbt.run(f"--select {model_name}")
12
13      list_of_assets.append(run_dbt_asset)
14
15  Definitions(assets = list_of_assets)
```

# 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)
```
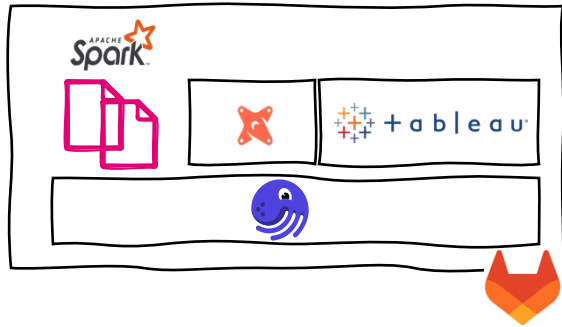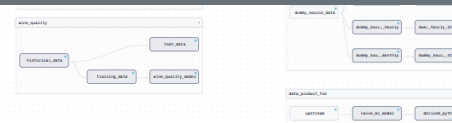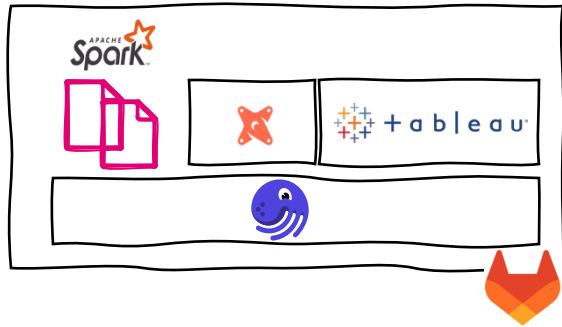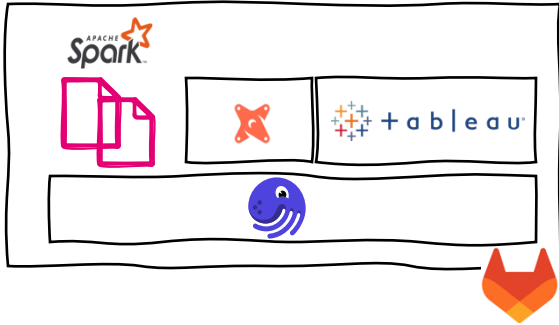
```
1   tableau_server = TableauServer(creds)
2   #send rest api call to tableau server and get information
3   tableau_metadata = tableau_server.get_metadata()
4   list_of_assets = []
5   for each_tableau_object in tableau_metadata:
6       tableau_object_deps = each_tableau_object.deps
7       tableau_object_name = 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)
```

# 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

```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]})
```

# 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)

# Last things



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

# Building block

- Dagster is the core
- Need for java because of custom spark code
- Development tooling for testing and code quality
- Pixi is holding everything together as the environment manager

# What should I do now?



developer

How to setup development environment?

ingest config file

dbt Core.

development environment

What is and how to use gitlab?

clone project

push changes

gitlab server

Development process?

project deployment

Dagster+

How to configure ingest configuration files?

How to start and debug jobs in local dagster?

What is dbt and how to use dbt?

Big query engine? cost model?

How to start and debug jobs in dagster cloud?

BigQuery

table_1

table_2

table_3

table_4

data ingest

data transformation