

# Data Mesh

## Introduction



**JOCHEN CHRIST**  
@JOCHEN\_CHRIST



Hi,  
I am Jochen

## Jochen Christ

Data Mesh Consultant  
Product Manager Data Mesh Manager



*Java*



*Data Mesh*



*Data-driven Product Development*







Home > Data Products > Search Queries All

# Search Queries All

Show specification

Edit

Request Access

Search source-aligned active managed demo



## Info

Information about the data product

Data Catalog

Repository

Documentation

Name  
Search Queries All

ID  
urn:dataprodut:search:search-queries-all

Description  
All search queries and results since 2020.

## Metrics

Monitor business value, costs, and compliance

Consumers  
1

Costs  
\$6,560.00

Compliance  
0 / 1 policies

## Input Ports

The source of the data, source systems or other data products.

OpenSearch  
acme.search.clicks

## Output Ports

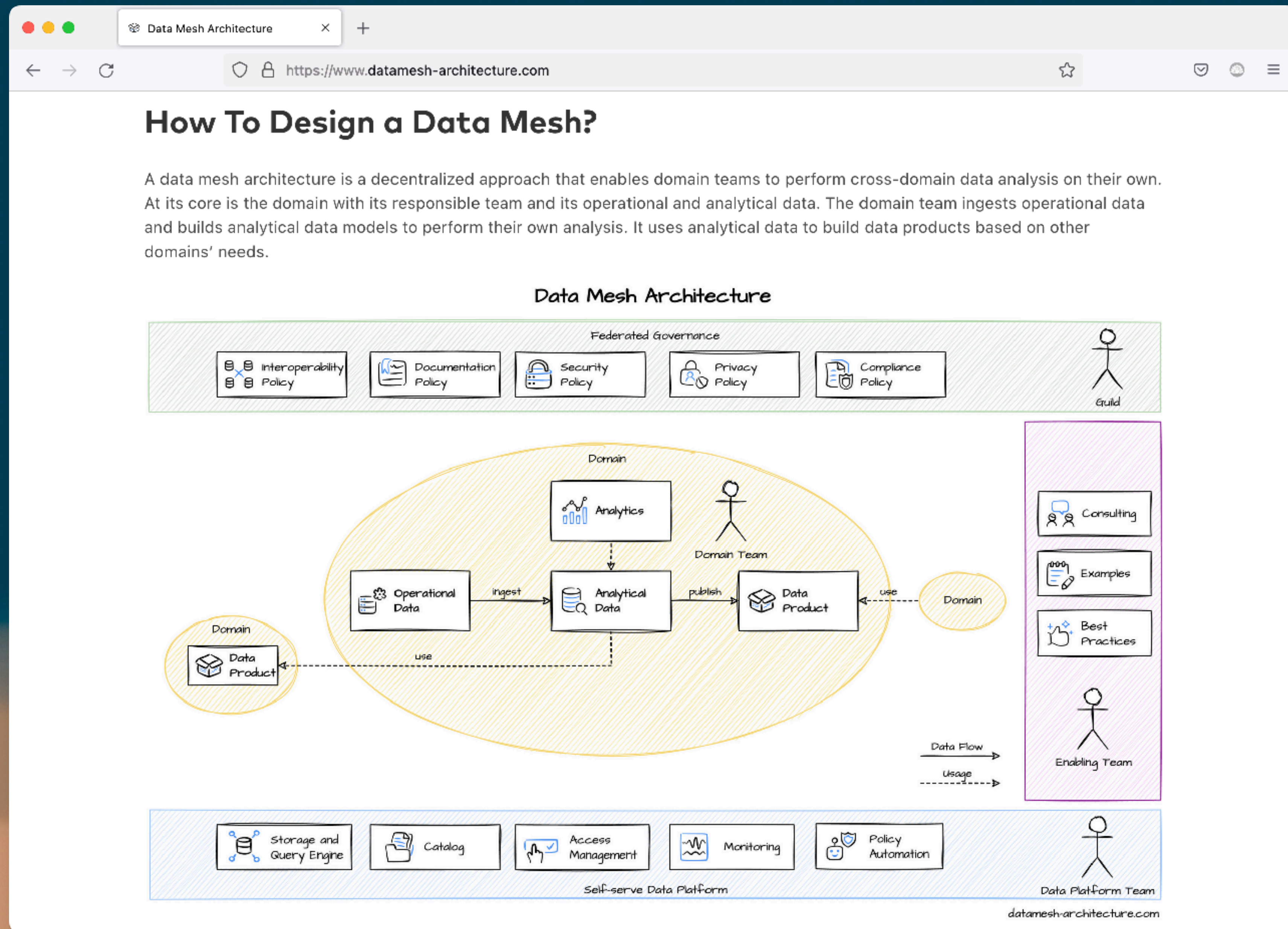
Technology, dataset, and version of provided data.

search\_queries\_all\_v1  
SEARCH\_DB.SEARCH\_QUERIES\_ALL\_NPII\_V1

Data Contract

1 usage





[datamesh-architecture.com](https://www.datamesh-architecture.com)



O'REILLY®

Deutsche  
Ausgabe

# Data Mesh

Eine dezentrale Datenarchitektur entwerfen



Zhamak Dehghani

Vorwort von Martin Fowler

Übersetzung von Jochen Christ und Simon Harrer

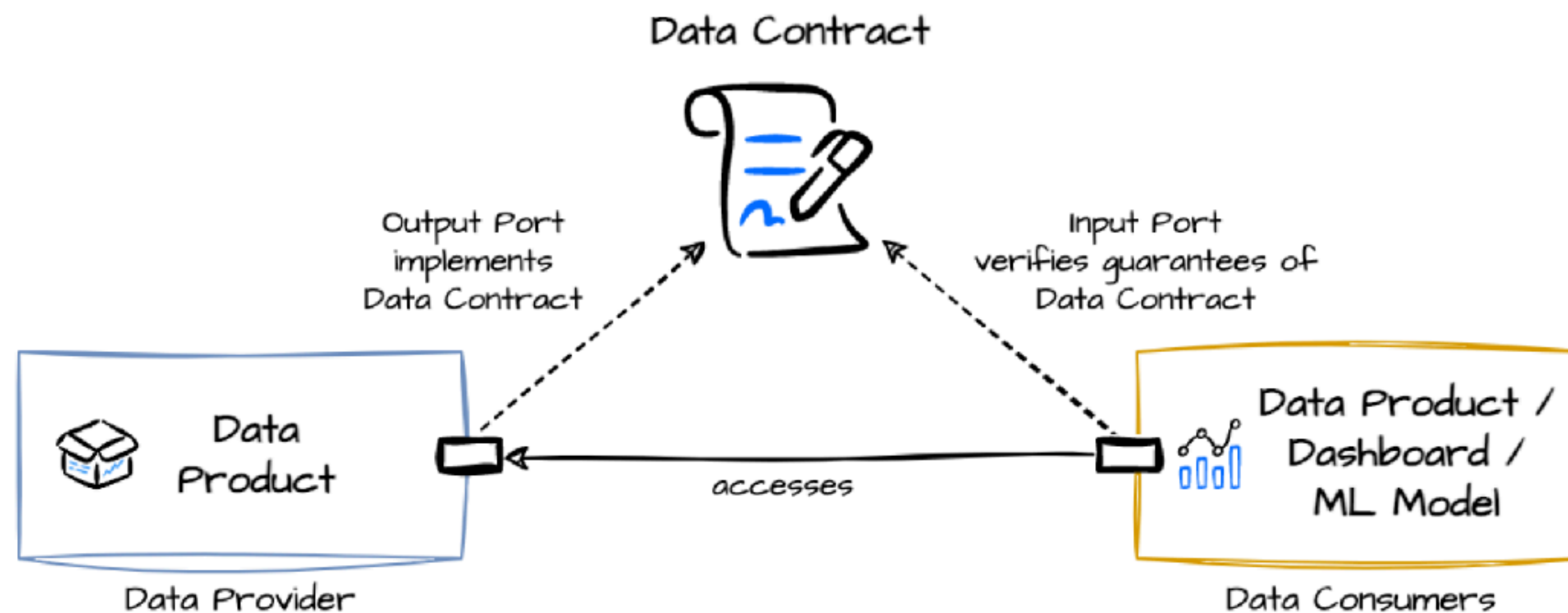




datacontract.com



# Data Contract Specification



Data contracts bring data providers and data consumers together.

A *data contract* is a document that defines the structure, format, semantics, quality, and terms of use for exchanging data between a data provider and their consumers. A data contract is implemented by a data product's output port or other data technologies. Data contracts can also be used for the input port to specify the expectations of data dependencies and verify given guarantees.

The *data contract specification* defines a YAML format to describe attributes of provided data sets. It is data platform neutral, yet supports well-known formats to express schemas (e.g., dbt models, JSON Schema, Protobuf, SQL DDL) and quality tests (e.g., SodaCL, SQL queries) to avoid unnecessary abstractions. The data contract specification is an open initiative to define a common data contract format. Think of an [OpenAPI specification](#), but for data sets.

datacontract.com



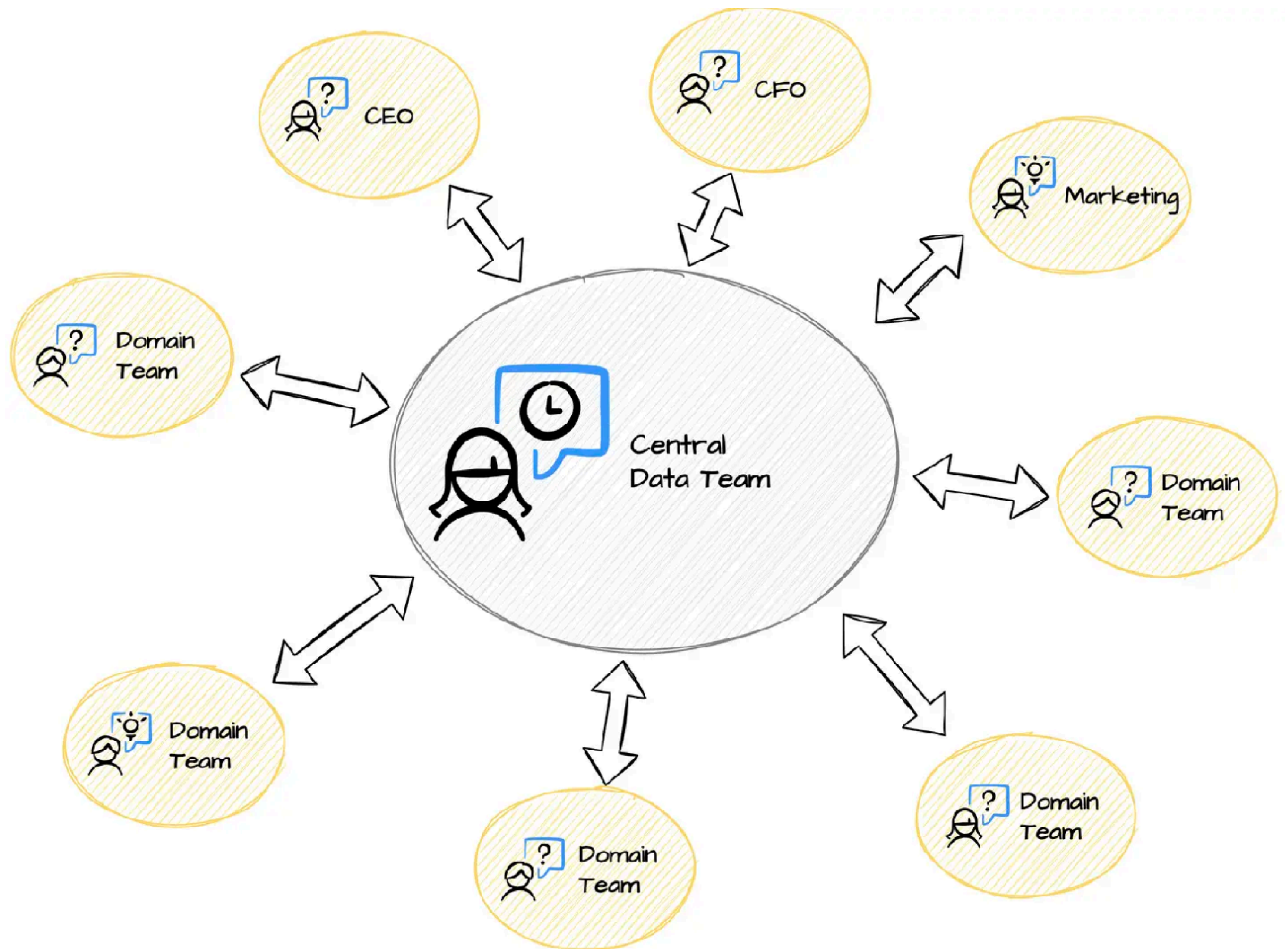
## Decentralized Data Architecture

# What?

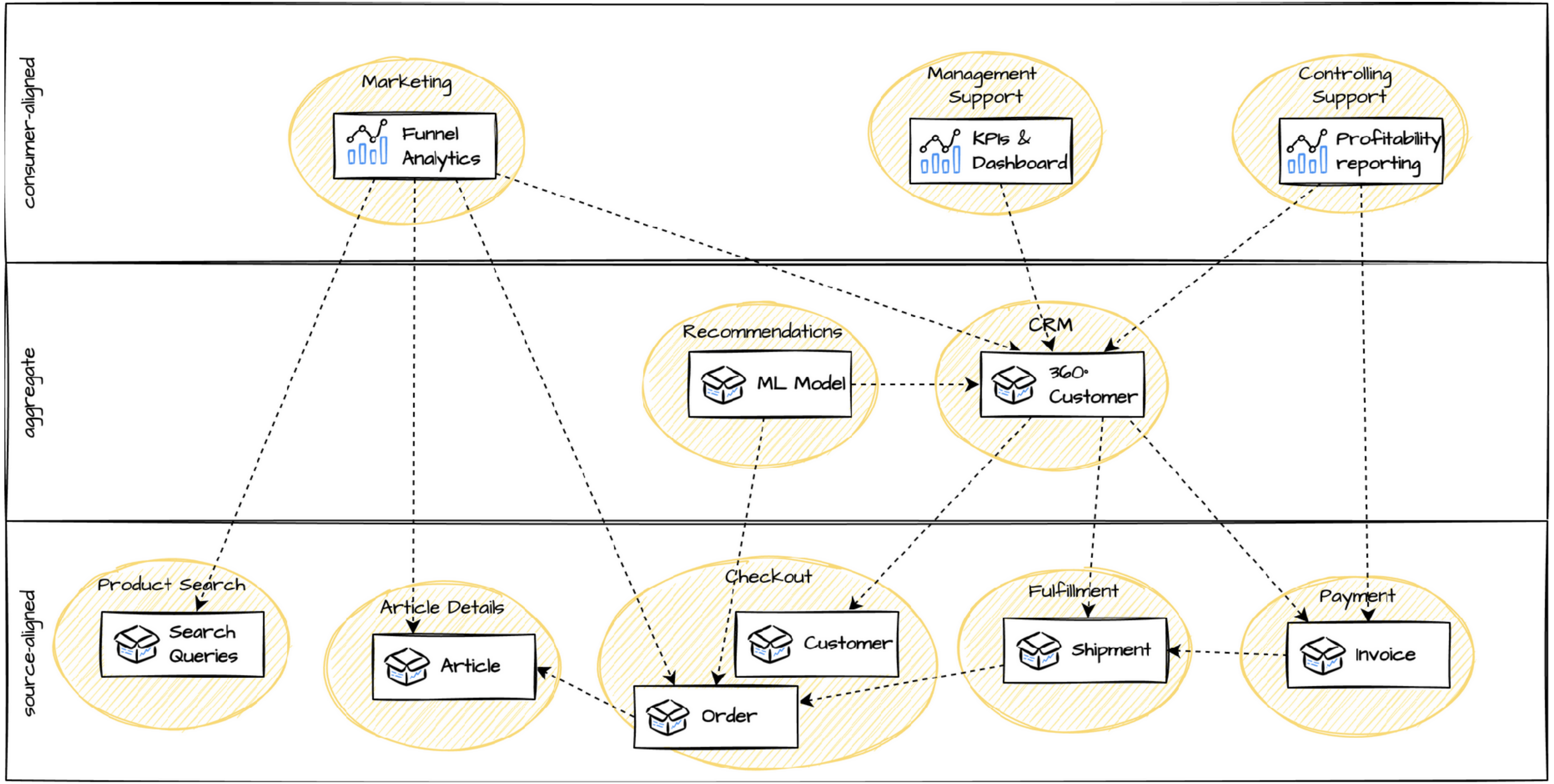
**A decentralized data architecture gives  
ownership and competence  
for (analytical) data  
to the teams that  
understand the business context.**

-- Jochen











## Decentralized Data Architecture

# Why?



### **Make qualified data-driven decisions in your domain**

Use data to better understand your users and system behavior. Derive features from insights, qualify value, and fast iterations. Also qualified rejection of unnecessary tasks.

Do the right things, purpose, motivation



### **Build innovative services in your domain**

Enhance your customer experience with data technologies, such as LLMs, visualizations, classifications, and ML models for predictions and recommendations.

Customer value through innovation



### **Provide data as business value for other domains**

Domain data is valuable for other business units as reference data and to aggregate. Needs managed, explained, high-quality and easy accessible data as products.

Company success



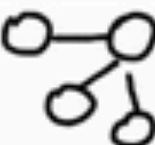
# What Is Data Mesh?

Strategic  
Domain-driven  
Design

Socio-technical  
Perspective

Technology


Domain  
Ownership

Domain  
Bounded Context  


Domain Teams  


Operational &  
Analytical Data  


Data as a  
Product

Product Thinking  


Data Product by  
Domain Team  


Interoperability  
Interfaces  


Self-serve  
Data Platform

Domain-agnostic

Data Platform  
Team

Self-serve  
Data Platform

Federated  
Governance

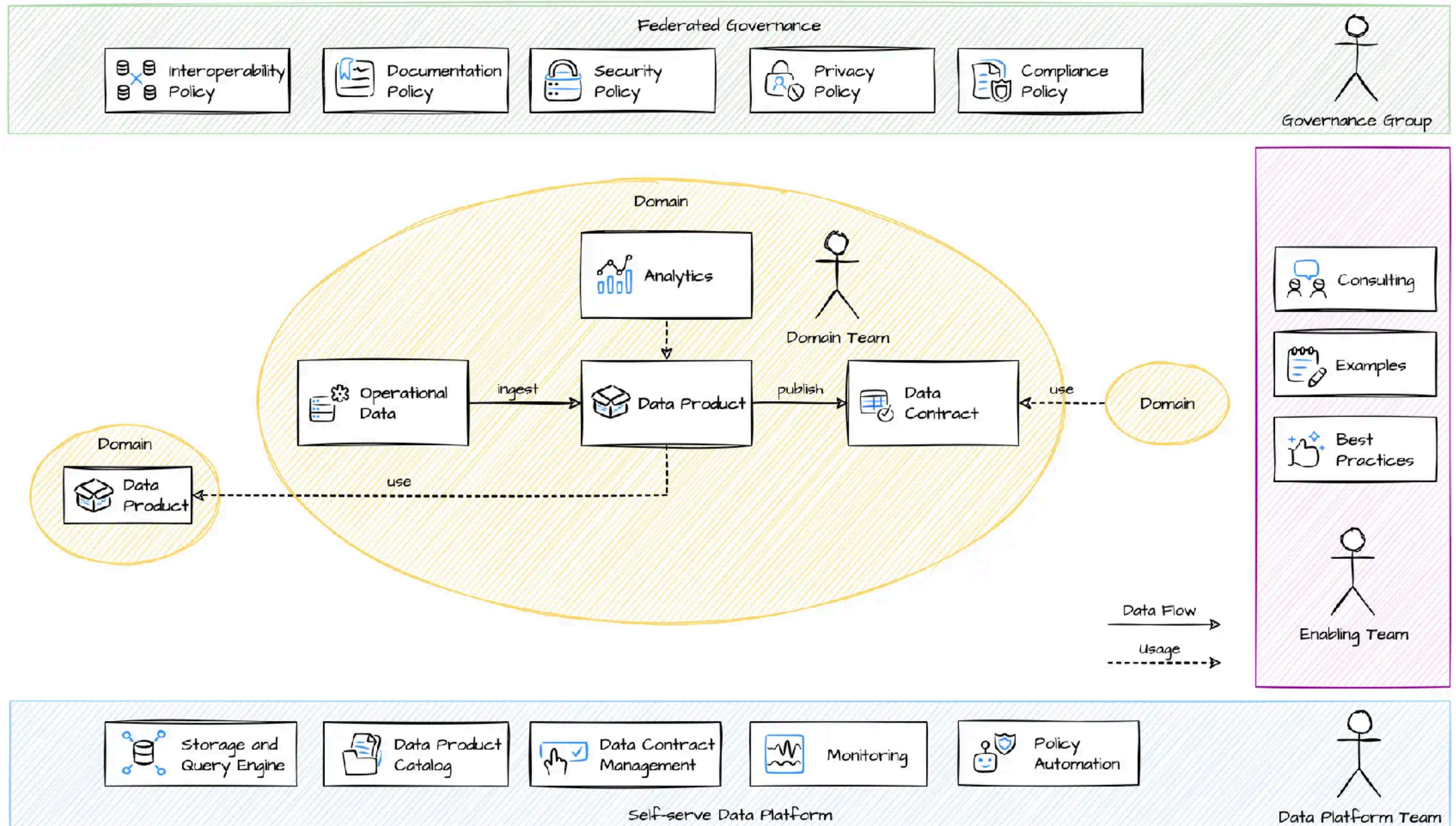
Context Mapping  


Guild  


Data Governance  
& Automation  

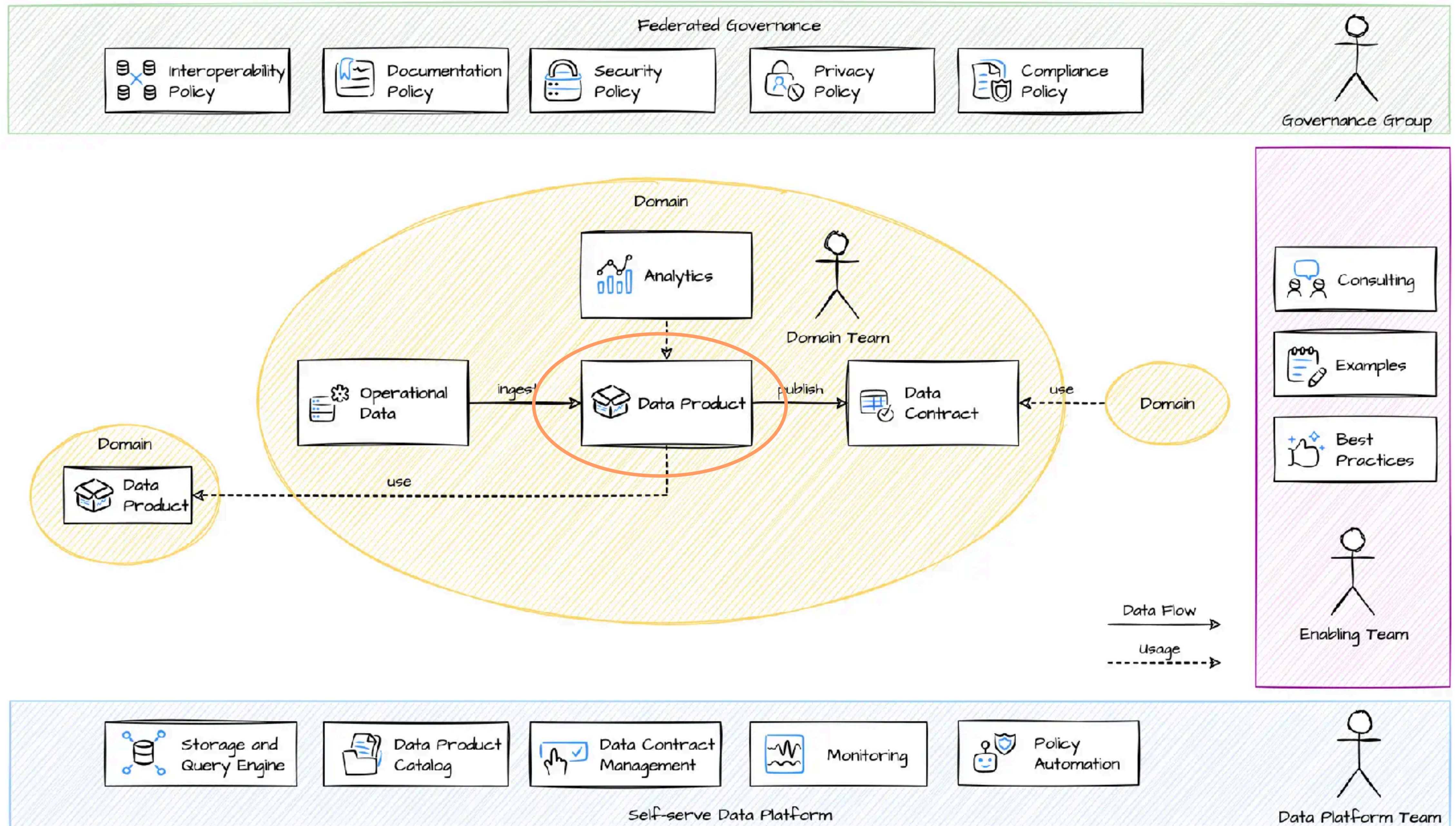



# Data Mesh Architecture



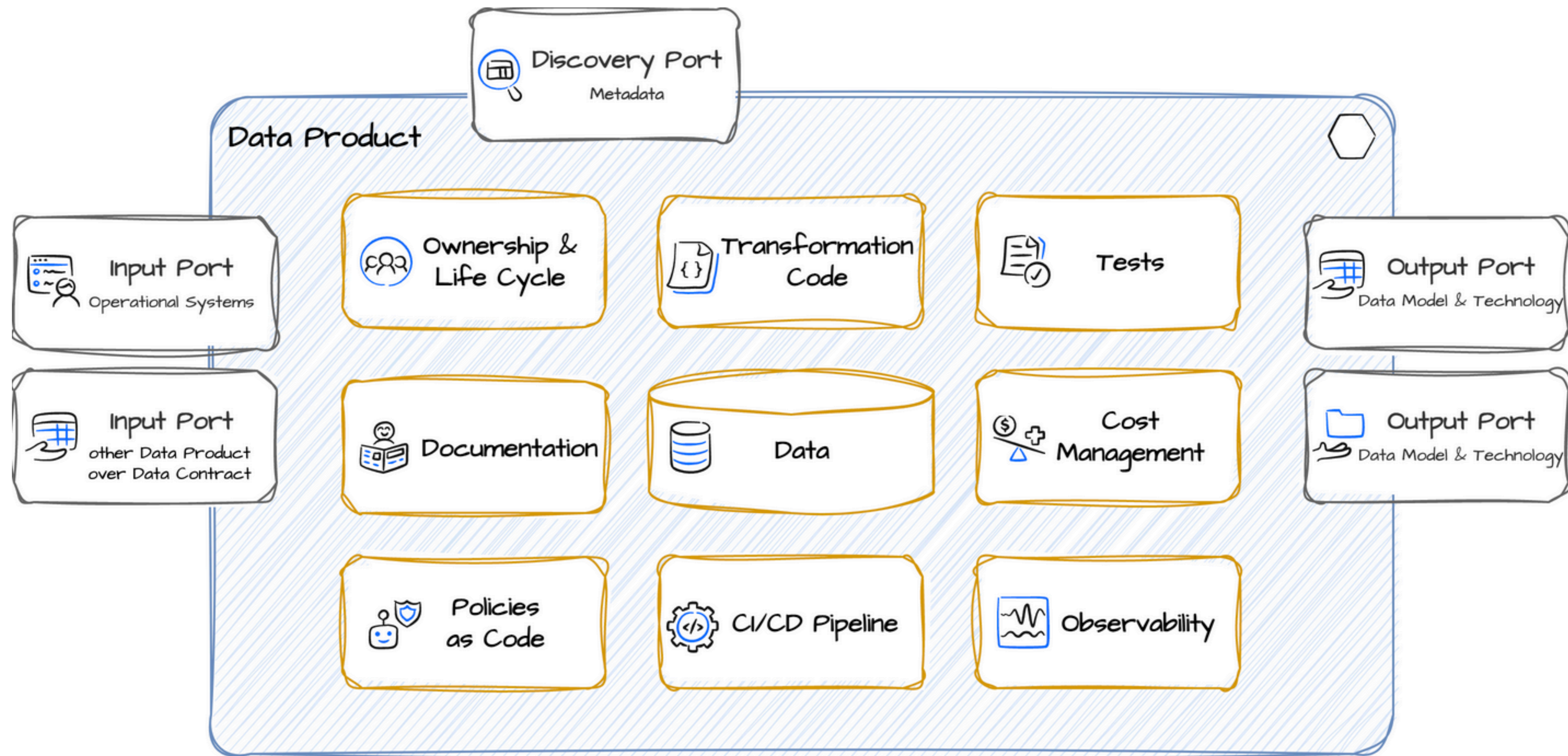


# Data Mesh Architecture





# Data Products are Modules





# Output Port Example

Row	sku	location	available	updated_at
1	9520010951145	20	0	2021-02-28 12:29:21 UTC
2	9520010951145	20	1	2021-03-02 09:07:21 UTC
3	9520010951145	20	0	2021-03-03 16:36:21 UTC
4	9520010951145	20	1	2021-03-04 13:03:21 UTC
5	9520010951145	20	2	2021-03-05 17:26:21 UTC
6	9520010951145	20	3	2021-03-06 03:35:21 UTC
7	9520010951145	20	2	2021-03-06 17:25:21 UTC
8	9520010951145	20	1	2021-03-07 18:10:21 UTC

- Technical Endpoint
- Hides implementation details
- Large data set
- Read-only
- Technology
  - Tables
  - Files in Bucket
  - Topic
- Data Model
  - With PII
  - Without PII
- Version



# Implementation depends on Tech Stack

Stack	Storage	Query Engine	Framework
AWS	S3	Athena (SQL)	Lambda / Step Functions
Google Cloud	BigQuery	BigQuery (SQL)	dbt
Azure (MS Fabric)	OneLake	Spark	Fabric notebook
Databricks	Deltalake	Spark	Databricks Asset Bundles
Snowflake	Snowflake	Snowflake (SQL)	dbt
Trino	S3 compliant	Trino (SQL)	dbt
Dremio	S3 compliant	Dremio Sonar	dbt
Java	S3 compliant	Java	Spring Cloud Data Flow
DuckDB	S3 compliant	DuckDB (SQL)	dbt



# Transformation: SQL

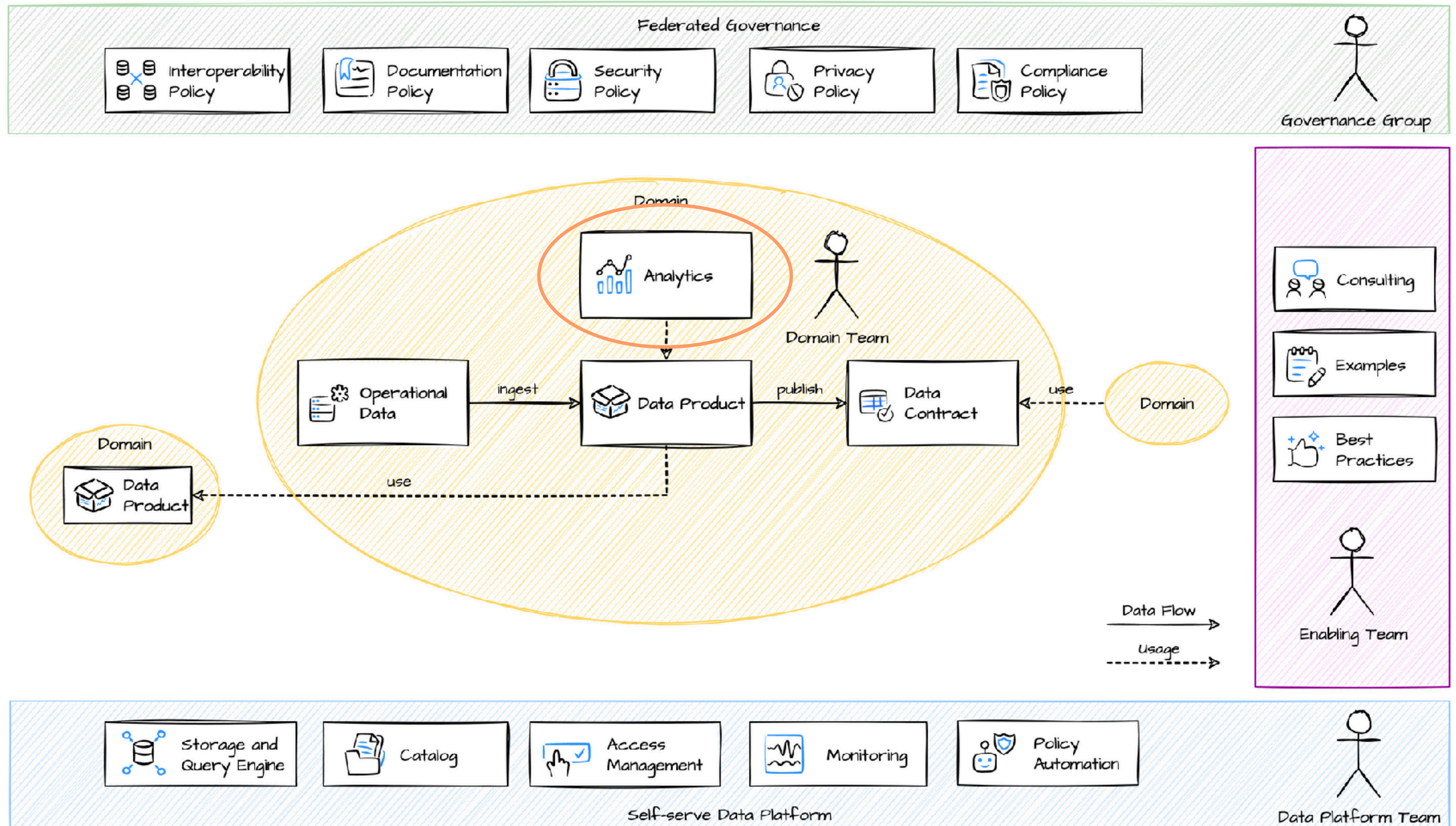
 entities\_\_inventory\_history.sql

Raw

```
1  -- Step 1: Deduplicate
2  WITH inventory_deduplicated AS (
3      SELECT *
4      EXCEPT (row_number)
5      FROM (
6          SELECT *,
7              ROW_NUMBER() OVER (PARTITION BY id ORDER BY time DESC) row_number
8          FROM `datameshexample-fulfillment.raw.inventory`)
9      WHERE row_number = 1
10 ),
11 -- Step 2: Parse JSON to columns
12 inventory_parsed AS (
13     SELECT
14         json_value(data, "$.sku") AS sku,
15         json_value(data, "$.location") AS location,
16         CAST(json_value(data, "$.available") AS int64) AS available,
17         CAST(json_value(data, "$.updated_at") AS timestamp) AS updated_at,
18     FROM inventory_deduplicated
19 )
20 -- Step 3: Actual Query
21 SELECT sku, location, available, updated_at
22 FROM inventory_parsed
23 ORDER BY sku, location, updated_at
```

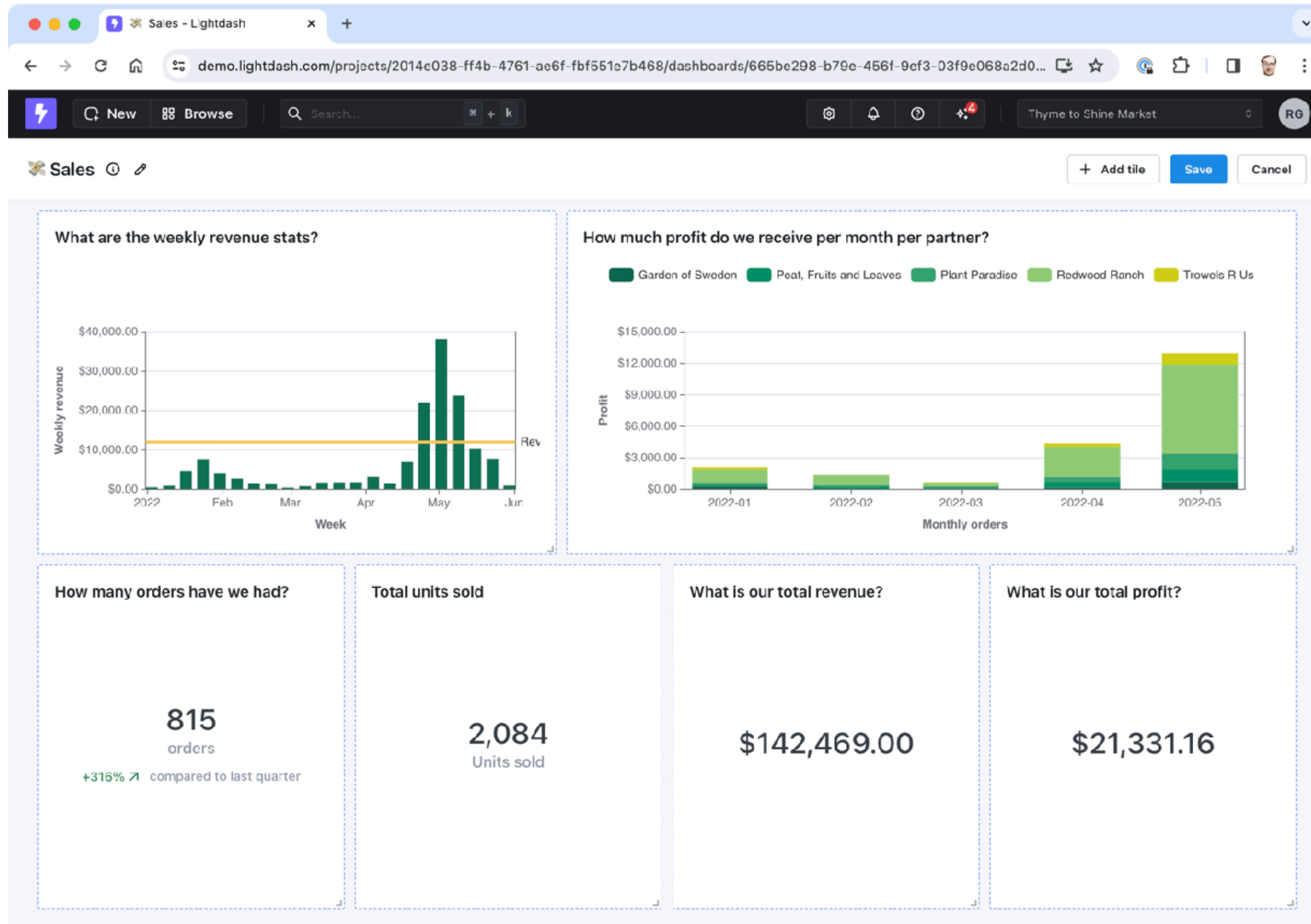


# Data Mesh Architecture



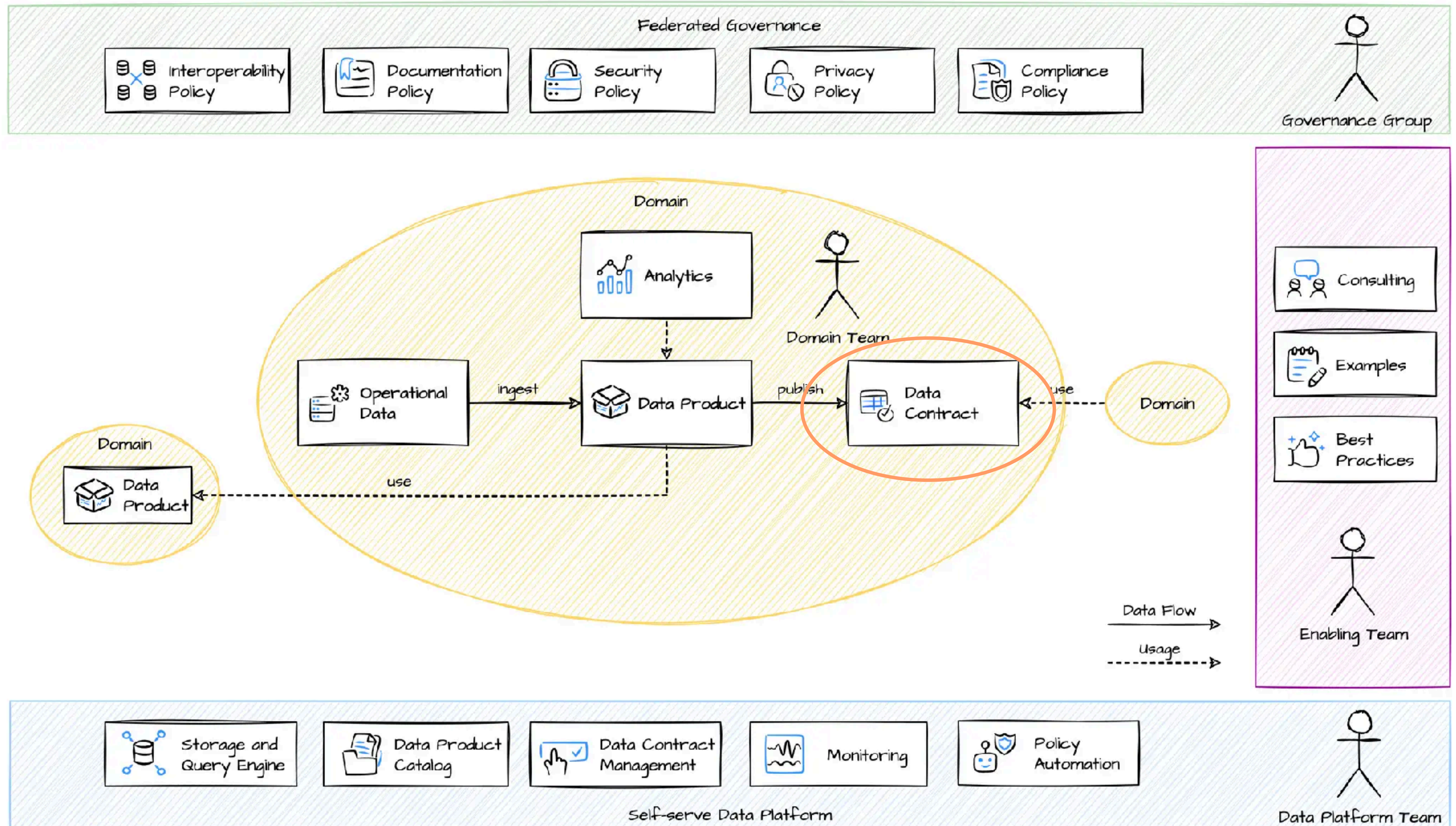


# Analytics: Enable Data Culture



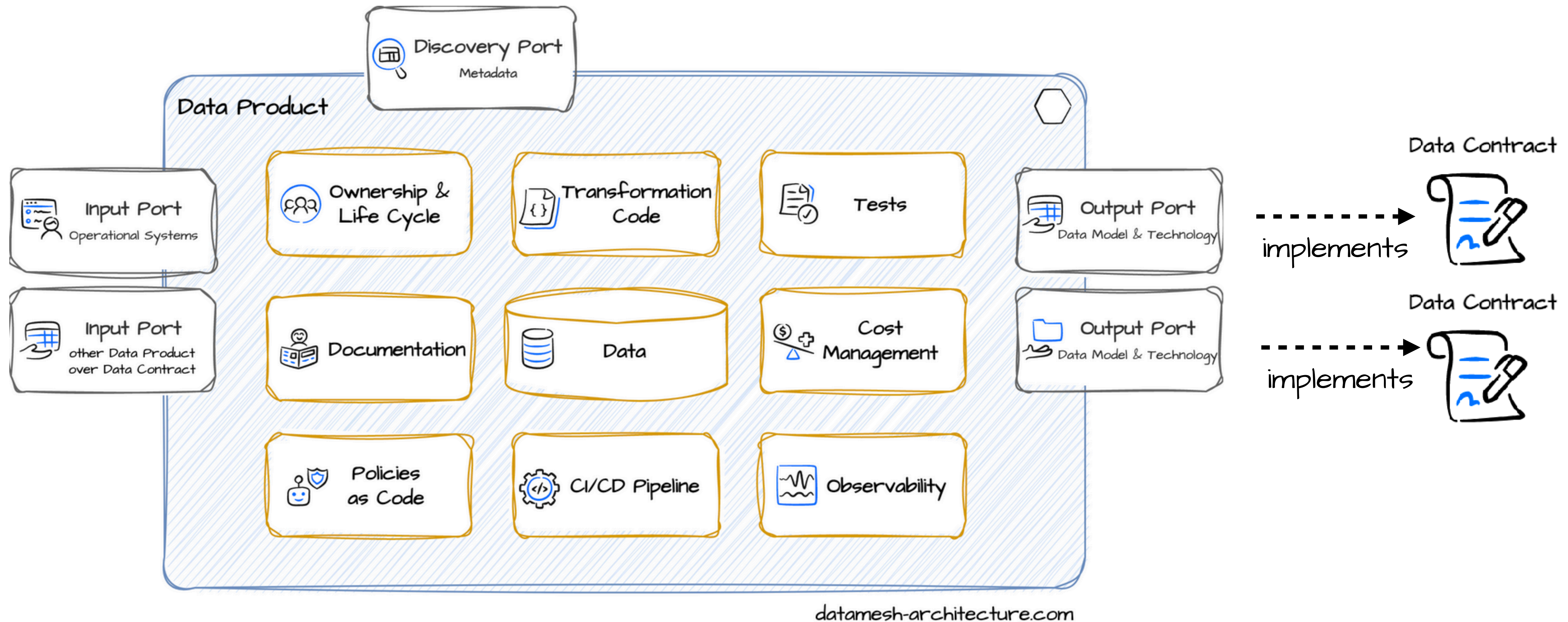


# Data Mesh Architecture





# Data Products Data Contracts





# APIs



# Messages



# Data





# Data Contract

```
dataContractSpecification: 0.9.1
id: web-orders-with-consent-v1
info:
  title: Web Orders With Consent V1
  version: 1.0.0
  description: "All orders made through the web channel.\r\nFiltered for orders where customers have expressed consent for analytical use."
  owner: checkout
  contact:
    url: https://teams.example.com/datacontracts/web-orders-with-consent-v1
terms:
  usage: "The data can be used for analytical and data science use cases, as the customer has expressed their consent."
  limitations: "As the dataset is filtered, these data set cannot be used to aggregate financial KPIs.\r\nNot suited for real-time use cases."
  billing: $1000 per month
  noticePeriod: P3M
models:
  orders:
    type: table
    description: A successful sale in the web shop
    fields:
      order_id:
        type: string
        description: Primary key of the order
      billing_customer_id:
        type: string
        description: Customer ID of the billing customer
      shipment_customer_id:
        type: string
        description: Customer ID of customer to ship the order to
      sold_timestamp:
        type: timestamp_tz
        description: The timestamp of the final confirmation step in the web form.
      total_amount:
        type: bigint
        description: The total order amount in the smallest unit of the currency (such as Eurocents)
```

**datacontract.com**  
**cli.datacontract.com**

- Interface Specification (like OpenAPI, but for data)
- YAML
- Define Requirements
- Make expectations explicit
- Make domain knowledge explicit
- Common language for data providers and consumers
- Owned by a team
- Contract-first
- Enforce Contract in CI/CD



dataContractSpecification: 0.9.1

id: web-orders-with-consent-v1

info:

title: Web Orders With Consent V1

version: 1.0.0

description: "All orders made through the web channel.\r\nFiltered for orders where customers have expressed consent for analytical use."

owner: checkout

contact:

url: <https://teams.example.com/datacontracts/web-orders-with-consent-v1>

terms:

usage: "The data can be used for analytical and data science use cases, as the customer has expressed their consent."

limitations: "As the dataset is filtered, these data set cannot be used to aggregate financial KPIs.\r\nNot suited for real-time use cases."

billing: \$1000 per month

noticePeriod: P3M

models:

orders:

type: table

description: A successful sale in the web shop

fields:

order\_id:

type: string

description: Primary key of the order

billing\_customer\_id:

type: string

description: Customer ID of the billing customer

shipment\_customer\_id:

type: string

description: Customer ID of customer to ship the order to

sold\_timestamp:

type: timestamp\_tz

description: The timestamp of the final confirmation step in the web form.

total\_amount:



# Contract Enforcement



order\_total:

description: Total amount the smallest monetary unit (e.g., cents)

type: long

required: true

customer\_id:

description: Unique identifier for the customer.

type: text

minLength: 10

maxLength: 20

customer\_email\_address:

description: The email address, as entered by the customer. The email

type: text

format: email

required: true



```
quality:
  type: SodaCL
  specification:
    checks for orders:
      - freshness(order_timestamp) < 24h
      - row_count > 500000
      - duplicate_count(order_id) = 0
    checks for line_items:
      - row_count > 500000
```



# Data Contract CLI

```
dataContractSpecification: 0.9.3
id: urn:datacontract:orders-latest
info:
  title: Orders Latest
  version: 1.0.0
models:
  orders:
    type: table
    fields:
      order_id:
        type: text
        format: uuid
```



*datacontract.yaml*

test



[github.com/datacontract/cli](https://github.com/datacontract/cli)



# **datacontract test**

```
$ datacontract test datacontract.yaml
```



# datacontract test

jochen — -zsh

jochen@Jochens-MacBook-Pro-2 ~ % datacontract test https://datacontract.com/examples/orders-latest/datacontract.yaml  
Testing <https://datacontract.com/examples/orders-latest/datacontract.yaml>

Result	Check	Field	Details
passed	Check that JSON has valid schema	orders	All JSON entries are valid.
passed	Check that JSON has valid schema	line_items	All JSON entries are valid.
passed	Check that field order_id is present	orders	
passed	Check that field order_timestamp is present	orders	
passed	Check that field order_total is present	orders	
passed	Check that field customer_id is present	orders	
passed	Check that field customer_email_address is present	orders	
passed	Check that field processed_timestamp is present	orders	
passed	row_count >= 5	orders	
passed	Check that required field order_id has no null values	orders.order_id	
passed	Check that unique field order_id has no duplicate values	orders.order_id	
passed	duplicate_count(order_id) = 0	orders.order_id	
passed	Check that required field order_timestamp has no null values	orders.order_timestamp	
passed	Check that required field order_total has no null values	orders.order_total	
passed	Check that required field customer_email_address has no null values	orders.customer_email_address	
passed	Check that required field processed_timestamp has no null values	orders.processed_timestamp	
passed	Check that field lines_item_id is present	line_items	
passed	Check that field order_id is present	line_items	
passed	Check that field sku is present	line_items	
passed	values in (order_id) must exist in orders (order_id)	line_items.order_id	
passed	row_count >= 5	line_items	
passed	Check that required field lines_item_id has no null values	line_items.lines_item_id	
passed	Check that unique field lines_item_id has no duplicate values	line_items.lines_item_id	

data contract is valid. Run 23 checks. Took 6.776398 seconds.

jochen@Jochens-MacBook-Pro-2 ~ %



Discover: Real User Interaction

Change column name · dataco

github.com/datacontract/cli-examples/actions/runs/6423146219/job/17441032658?pr=2

Code

Issues

Pull requests

1

Actions

Projects

Wiki

Security

Insights

Settings

← Back to pull request #2

Change column name #11

Re-run jobs

Summary

Jobs

checkBreakingChanges

Run details

Usage

Workflow file

checkBreakingChanges

failed 5 days ago in 7s

Search logs

> Set up job

1s

> Run actions/checkout@v4

1s

> Get CLI

0s

> Check backwards compatibility

0s

1 ▶ Run ./datacontract breaking --with https://raw.githubusercontent.com/datacontract/cli-examples/main/datacontract.yaml

4 Found 1 differences between the data contracts!

5

6 Difference 1:

7 Description: field 'my\_table.my\_column' was removed

8 Type: field-removed

9 Severity: breaking

10 Level: field

11 Model: my\_table

12 Field: my\_column

13 Exiting application with error: found breaking differences between the data contracts

14 Error: Process completed with exit code 1.



# Data Contract CLI



import

```
dataContractSpecification: 0.9.3
id: urn:datacontract:orders-latest
info:
  title: Orders Latest
  version: 1.0.0
models:
  orders:
    type: table
    fields:
      order_id:
        type: text
        format: uuid
```

*datacontract.yaml*

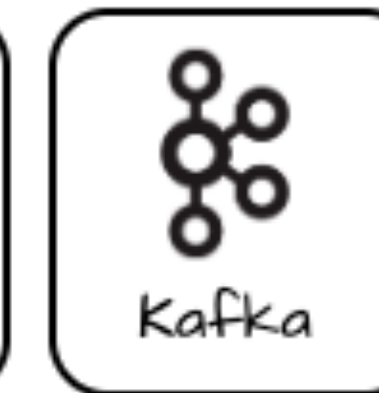


diff

export



test





# Data Marketplace





# Data Contracts

Definitions

+ Add Data Contract

Owner

Data Product

Sort

<div><div>Customer Cohorts</div><div><div>Marketing</div><div>Customer Cohorts</div><div>1</div></div><div><div></div><div></div></div></div> <div>A table with customer cohorts and their properties</div>	<div><div>s3_customers_history_npii_v1</div><div><div>Checkout</div><div>Customers</div></div><div><div></div><div></div></div></div> <div>All customer states, updated on every modifying event. PII removed.</div>	<div><div>s3_customers_history_pii_v1</div><div><div>Checkout</div><div>Customers</div></div><div><div></div><div></div></div></div> <div>All customer states, updated on every modifying event. PII included.</div>
<div><div>search_queries_all_v1</div><div><div>Search</div><div>Search Queries All</div><div>1</div></div><div><div></div><div></div></div></div> <div>All search queries and result sets with PII removed.</div>	<div><div>snowflake_articles_history</div><div><div>Products</div><div>Articles history</div><div>2</div></div><div><div></div><div></div></div></div> <div>All article snapshots since 2020</div>	<div><div>snowflake_articles_latest</div><div><div>Products</div><div>Articles latest</div><div>1</div></div><div><div></div><div></div></div></div> <div>Current state of all articles</div>
<div><div>snowflake_customers_latest_npii_v1</div><div><div>Checkout</div><div>Customers</div></div><div><div></div><div></div></div></div> <div>All customers in their latest state, PII removed.</div>	<div><div>snowflake_customers_latest_pii_v1</div><div><div>Checkout</div><div>Customers</div><div>1</div></div><div><div></div><div></div></div></div> <div>All customers in their latest state, PII included.</div>	<div><div>snowflake_fulfillment_shelf_warmers</div><div><div>Fulfillment</div><div>Shelf Warmers</div></div><div><div></div><div></div></div></div> <div>A list of articles with no sales in last 6 months</div>
<div><div>snowflake_fulfillment_stock_update_events</div><div><div>Fulfillment</div><div>Stock Update Events</div><div>1</div></div><div><div></div><div></div></div></div> <div>All stock updates since 2020</div>	<div><div>snowflake_orders_npii_v2</div><div><div>Checkout</div><div>Orders</div><div>1</div></div><div><div></div><div></div></div></div> <div>All order-created events, PII removed.</div>	<div><div>snowflake_orders_pii_v2</div><div><div>Checkout</div><div>Orders</div><div>1</div></div><div><div></div><div></div></div></div> <div>All order-created events, with PII.</div>



# Request Access

You are requesting access to the data product [Articles history](#) on output port [snowflake\\_articles\\_history](#) .  
The system will create a data usage agreement for the team [Products](#) to approve.

Consumer

Required

Article Profitability Analysis (Team Controlling)



Select your data product that wants to access and use the provided data.

Purpose

Required

Profitability analysis.



Why do you want access and what do you want to do with the data?

Cancel

Customize

Request access



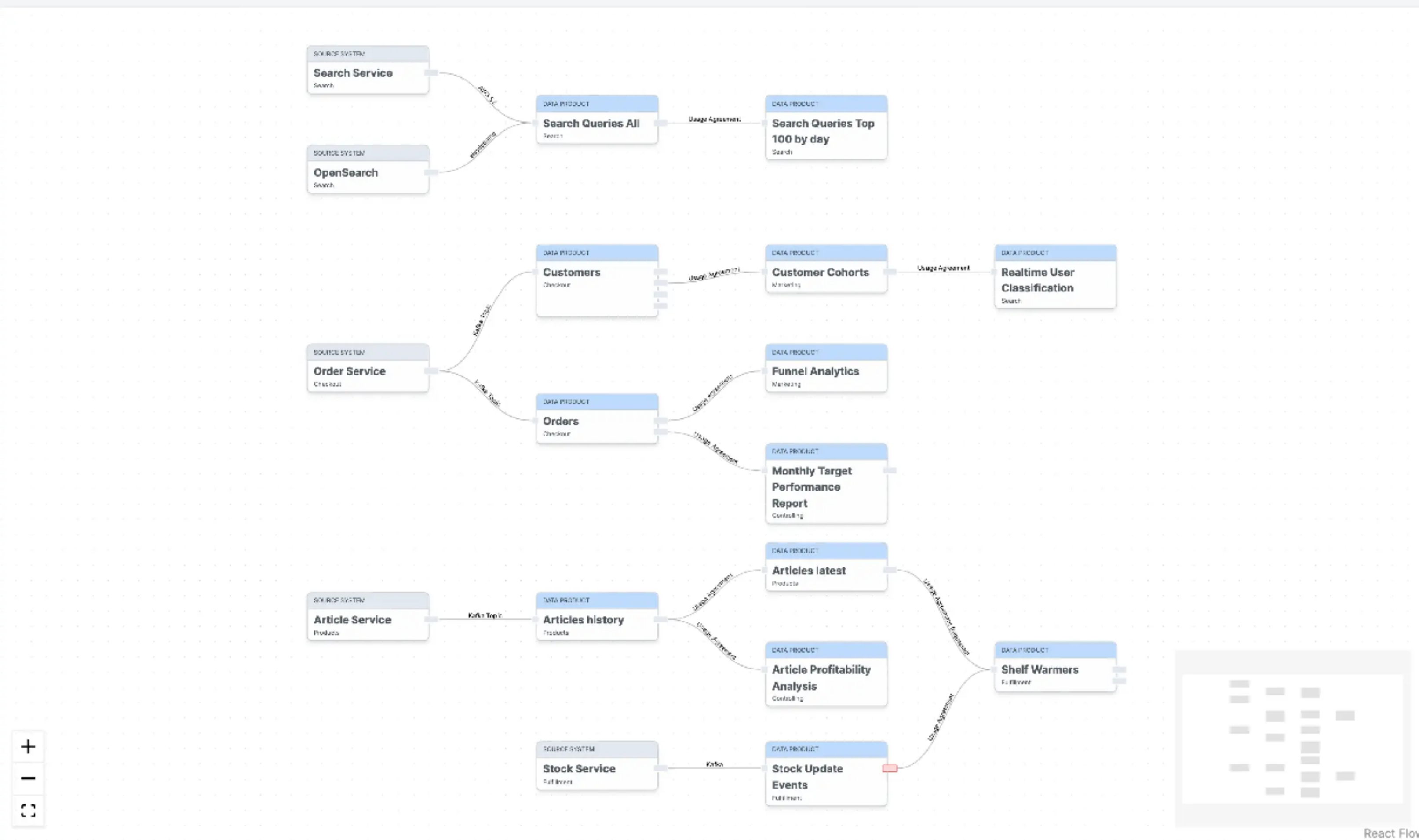


## Data Map

[Data Products](#)[Source Systems](#)[Domain Teams](#)[+ Add Data Product](#)

Owner ▾

Status ▾





## Decentralized Data Architecture

# Start using your own data



### **Make qualified data-driven decisions** in your domain

Use data to better understand your users and system behavior. Derive features from insights, qualify value, and fast iterations. Also qualified rejection of unnecessary tasks.

Do the right things, purpose, motivation



### **Build innovative services** in your domain

Enhance your customer experience with data technologies, such as LLMs, visualizations, classifications, and ML models for predictions and recommendations.

Customer value through innovation



## Decentralized Data Architecture

# And share when there is demand



### **Make qualified data-driven decisions** in your domain

Use data to better understand your users and system behavior. Derive features from insights, qualify value, and fast iterations. Also qualified rejection of unnecessary tasks.

Do the right things, purpose, motivation



### **Build innovative services** in your domain

Enhance your customer experience with data technologies, such as LLMs, visualizations, classifications, and ML models for predictions and recommendations.

Customer value through innovation



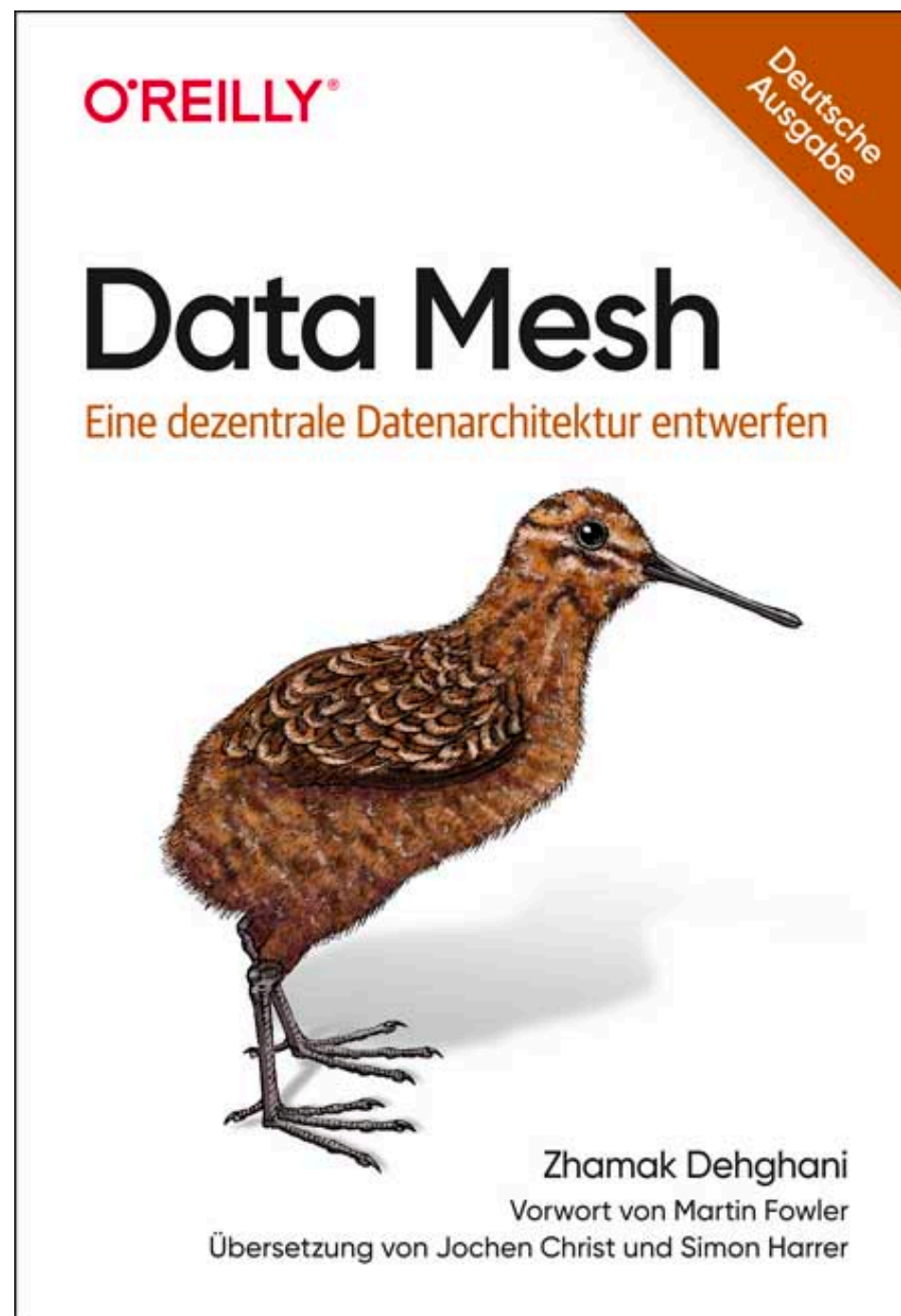
### **Provide data as business value** for other domains

Domain data is valuable for other business units as reference data and to aggregate. Needs managed, explained, high-quality and easy accessible data as products.

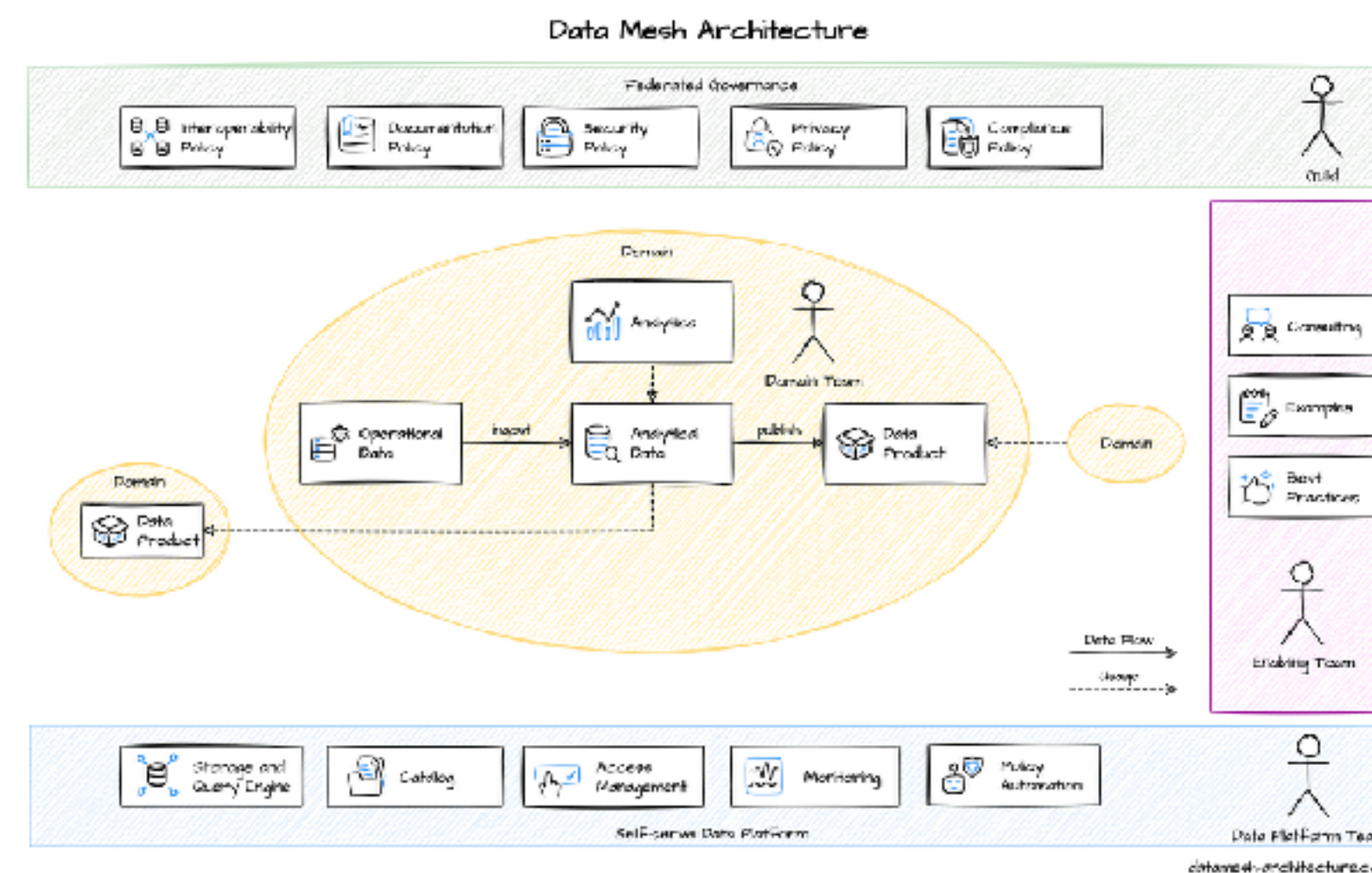
Company success



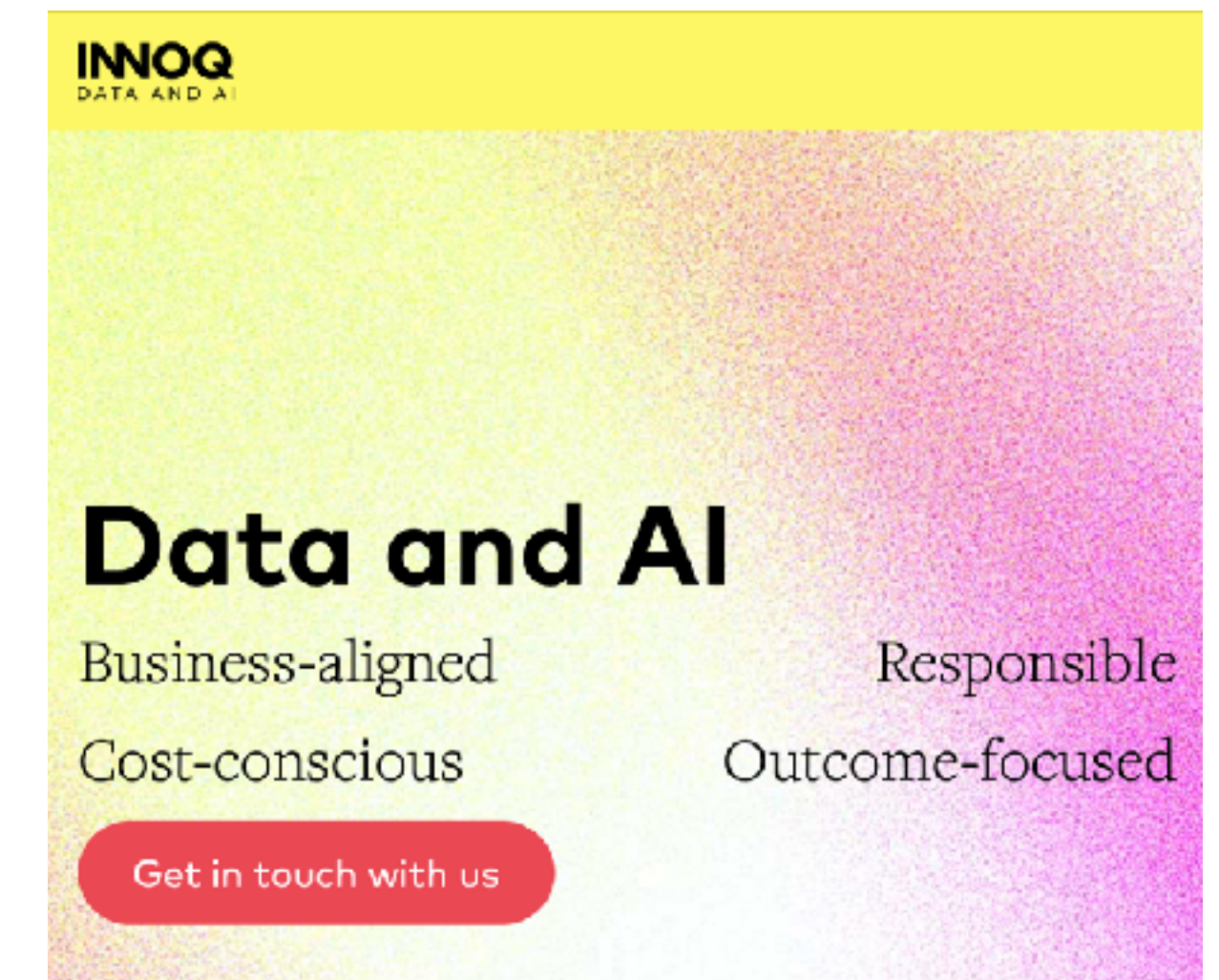
# Learn more



[oreilly.de/produkt/data-mesh](https://oreilly.de/produkt/data-mesh)



[datamesh-architecture.com](https://datamesh-architecture.com)



[INNOQ.ai](https://INNOQ.ai)  
Data Mesh Consulting, Trainings,  
Data Product Engineering



# Data Mesh

## Introduction



**JOCHEN CHRIST**  
@JOCHEN\_CHRIST