# From Software Engineering ...

**Java by Comparison**
Become a Java Craftsman in 70 Examples

The Pragmatic Programmers

Simon Harrer, Jörg Lenhard, Linus Dietz
Foreword by Venkat Subramaniam
Edited by Andrea Stewart

**Remote Mob Programming**

Jochen Christ
Simon Harrer
Martin Huber

At home, but not alone

INNOQ

Foreword by Mark Pearl

**GitOps**
Cloud-native Continuous Deployment

Florian Beetz · Anja Kammer · Simon Harrer

Quickstart with Kubernetes

INNOQ

**java.by-comparison.com**

**remotemobprogramming.org**

**gitops.tech**

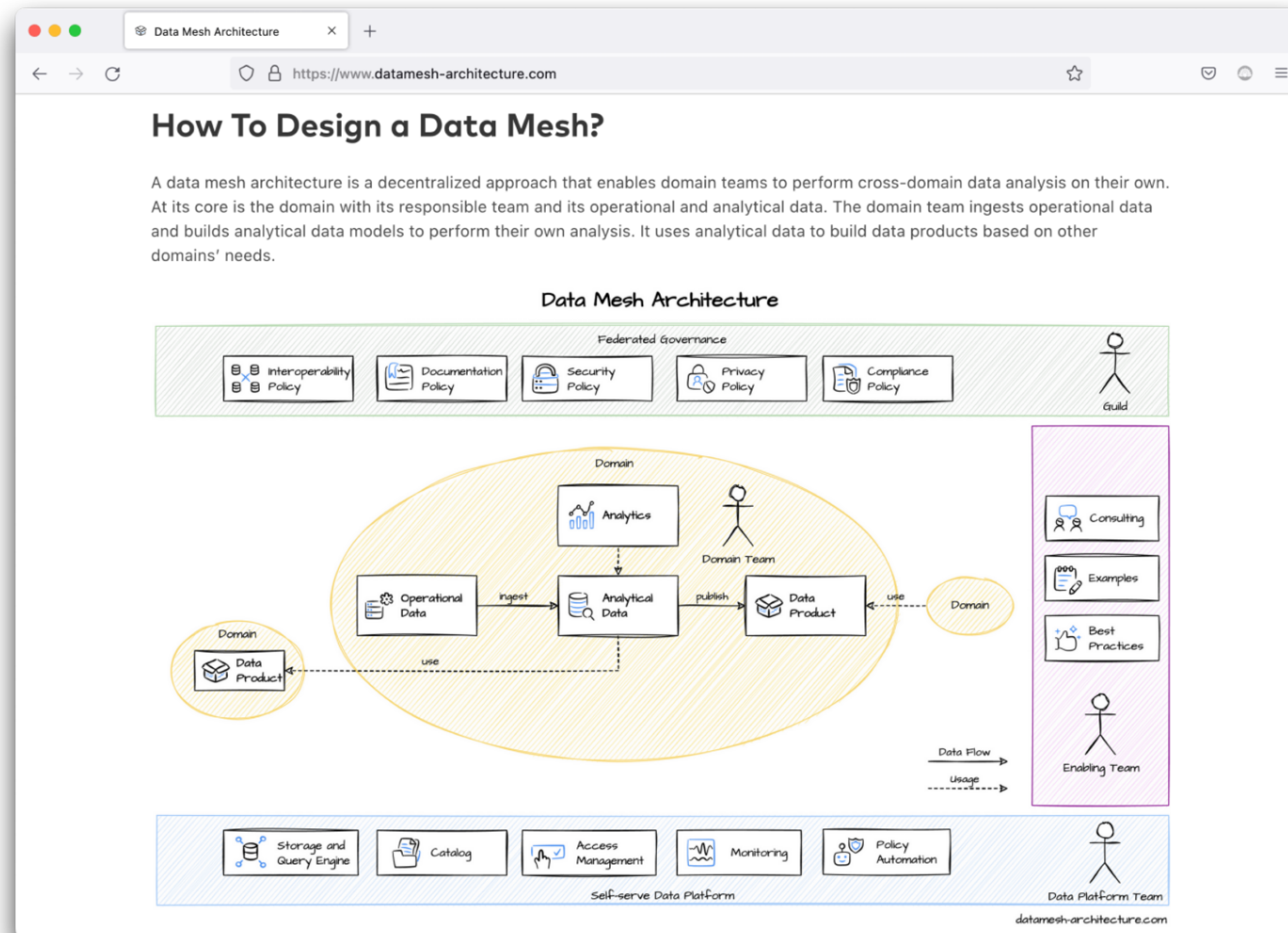**2021**

# ... to Data Mesh ...



**datamesh-architecture.com**



**oreilly.de/produkt/data-mesh**

# ... to Open Standards and Tools



Bitol | LF AI & DATA

Contributors to the data contract

- Data Engineers
- Data Scientists
- Data Product Owners
- Automation Tools

**Open Data Contract Standard v3**

| Fundamentals | Schema | Data Quality | Pricing |
|---|---|---|---|
| Name, Version, Descriptions... | Logical representation & physical implementation | Data quality rules, Data governance policies | Internal or external costs associated to usage |
| Team | Security | SLA | Infrastructure |
| History of team members | Roles | Latency, retention, frequency... | Servers, environment, and storage |
| Support | Business rules* | Tags & Custom Properties | |
| Support mechanisms for consumers | Data QoS applied to your business needs | Custom extensions & processing | |

**Enterprise-level consumers & contributors**
- Enterprise Data Governance
- Enterprise Security & Audit
- Enterprise Data Catalog
- Enterprise Ops

**Consumers of the data contract**
- Applications
- Observability
- Monitoring
- Other Tools
- Notification

* Future

## Data Contract CLI

SQL DDL · Avro · JSON Schema · Protobuf · BigQuery · Unity Catalog · AWS Data Catalog · ODCS

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
or odcs.yaml

diff

export →

SQL DDL · HTML · Avro · RDF · dbt · SodaCL · Terraform · ODCS

test ↓

AWS S3 · BigQuery · Azure · databricks · snowflake · Kafka

github.com/datacontract/cli

# ... to Open Standards and Tools



Star History

- bitol-io/open-data-contract-standard
- datacontract/datacontract-cli

star-history.com

I'll talk about data mesh,

data products,

and data contracts.


And, if time allows, using AI in this space.

# Part 1: Data Mesh

# Status Quo

# The Problem

**Great expectations of data**
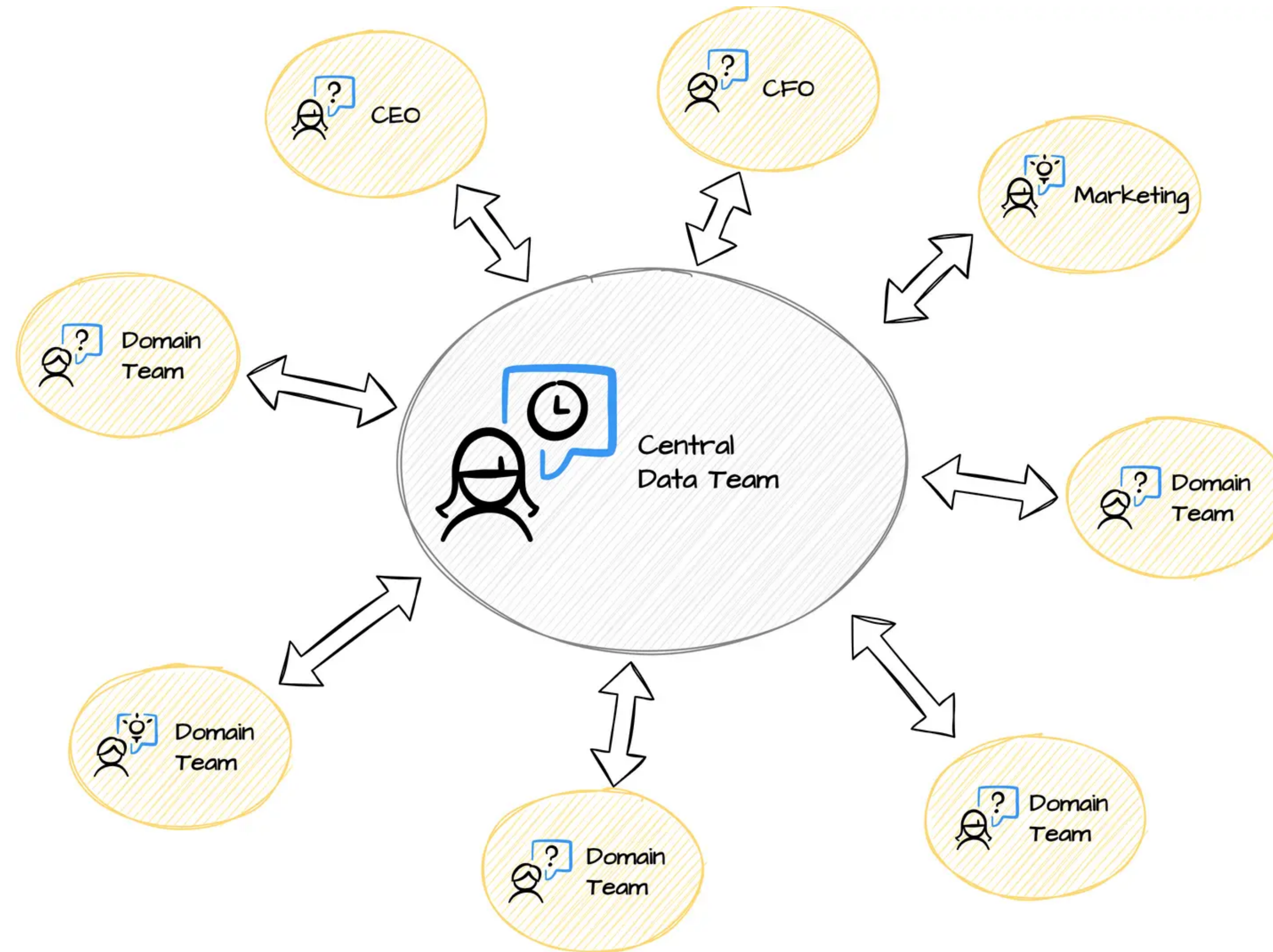Diverse and wide applications of ML and analytics

**Great divide of data**
Complexity risen from fragmentation of operational and analytical data

**Scale**
Large scale data source proliferation

**Business complexity & volatility**
Continuous change and growth of businesses

**Discord between data investments and returns**
Expensive data solutions lacking impact

**Organizational impact**
Agility in response to change
Get value from data

**Data mesh**
Promise to new heights

**Inflection point**
Approach to data management

No change in the approach

**Scale of complexity**
Business complexity and size
Data source proliferation
Data usage diversity

O'REILLY
Data Mesh
Eine dezentrale Datenarchitektur entwerfen

Zhamak Dehghani
Vorwort von Martin Fowler
Übersetzung von Jochen Christ und Simon Harrer

# The Solution



datamesh-architecture.com

# I brought a Data Mesh with me

# The Architecture Perspective



Data Mesh Architecture

Federated Governance
- Interoperability Policy
- Documentation Policy
- Security Policy
- Privacy Policy
- Compliance Policy

Governance Group

Domain
- Analytics
- Operational Data — ingest → Data Product — publish → Data Contract
- Domain Team

use

Domain
- Data Product

Domain — use

Consulting
Examples
Best Practices

Enabling Team

Data Flow
Usage

Self-serve Data Platform
- Storage and Query Engine
- Data Product Catalog
- Data Contract Management
- Monitoring
- Policy Automation

Data Platform Team

datamesh-architecture.com

In summary: the vast majority of larger companies try moving towards the principles of data mesh.

But be aware of the name data mesh.

# Part 2: Data Products

# What is a data product?

# My own practical definition:

A data product is a logical unit

that contains all components to process domain data

and provide data sets via output ports.

(The view of a software architect)

# I see them as an architectural unit



datamesh-architecture.com

# That can form a large graph / mesh

# The Birth-Certificate

| Domain | Data Product | Date | |
|---|---|---|---|

## Data Product Architecture

### Processing
Streaming, batch

### Framework
dbt, Databricks, Java, …

### Costs
Budget and cost control

**Transformation Steps**
Ingestion, cleaning, conversions, aggregations, joining
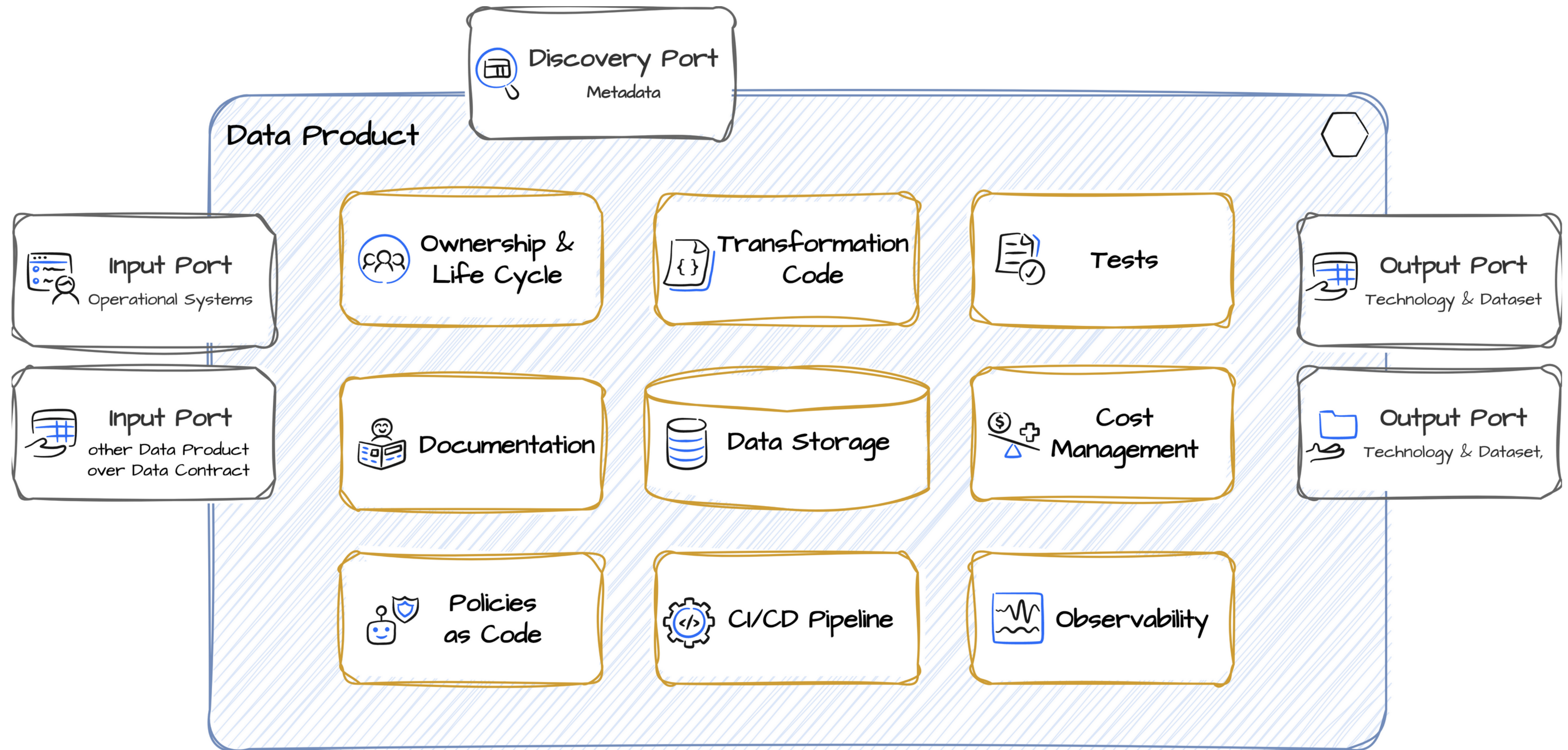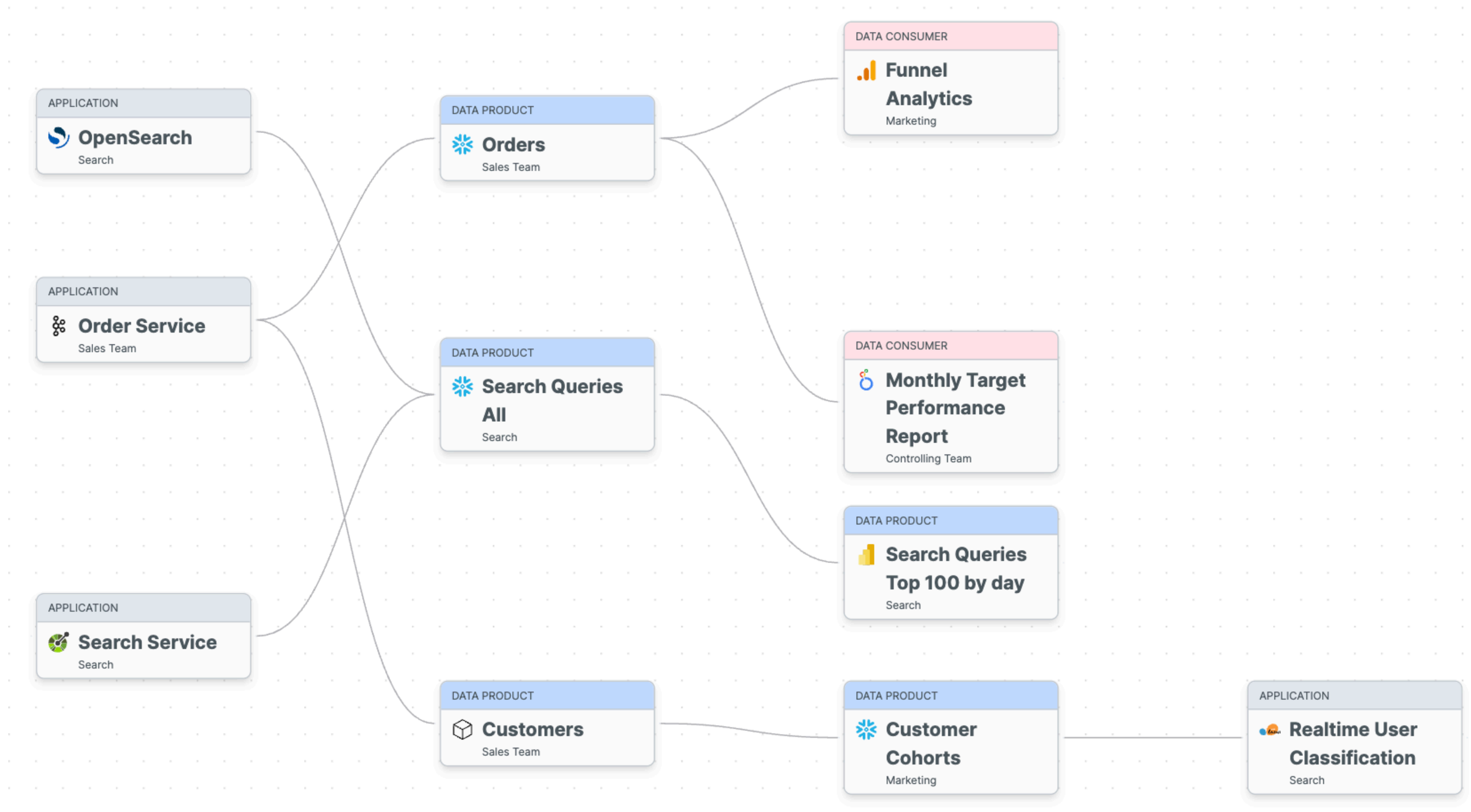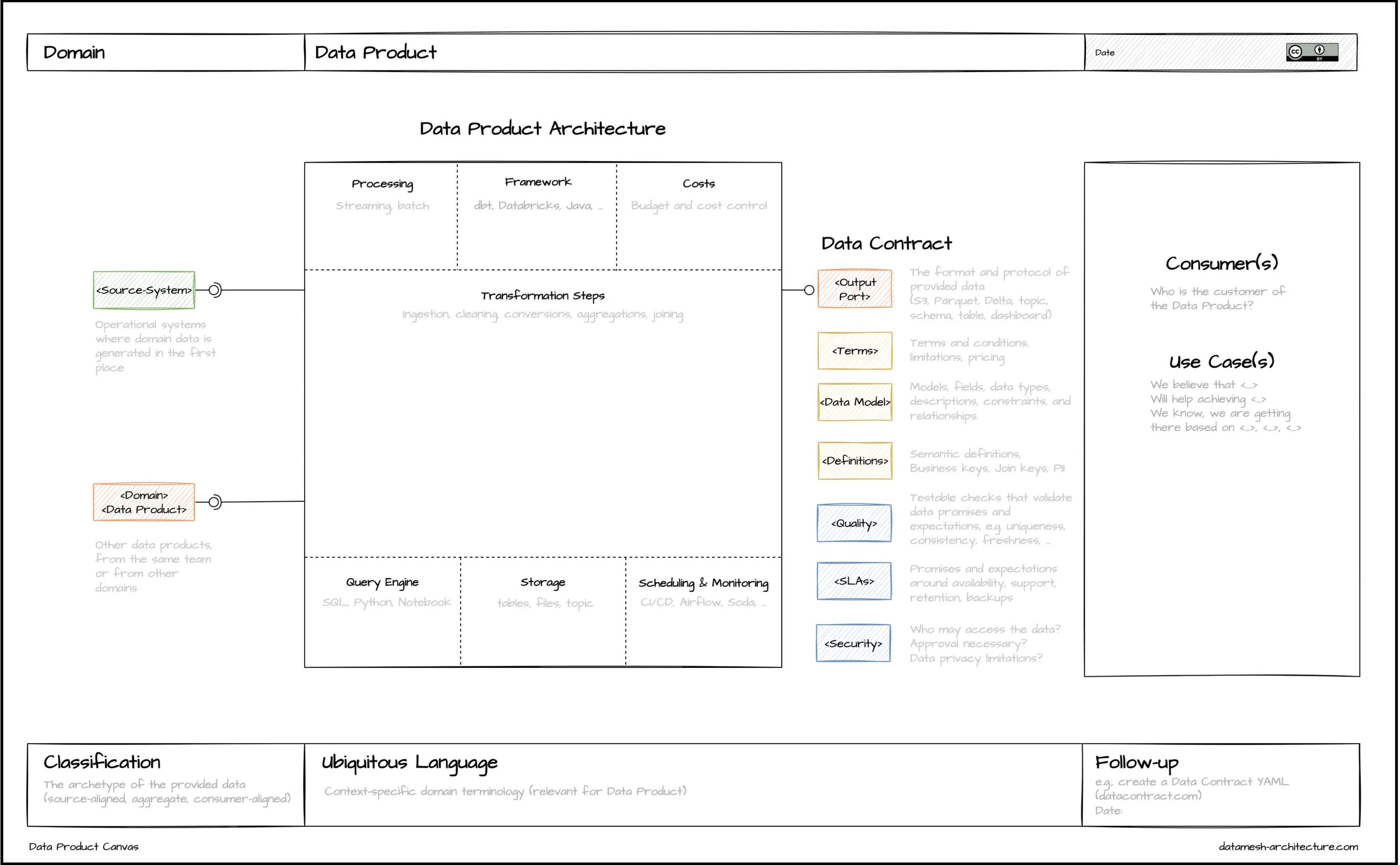
### Query Engine
SQL, Python, Notebook

### Storage
tables, files, topic

### Scheduling & Monitoring
CI/CD, Airflow, Soda, …

**<Source-System>**
Operational systems where domain data is generated in the first place

**<Domain> <Data Product>**
Other data products, from the same team or from other domains

## Data Contract

**<Output Port>**
The format and protocol of provided data (S3, Parquet, Delta, topic, schema, table, dashboard)

**<Terms>**
Terms and conditions, limitations, pricing

**<Data Model>**
Models, fields, data types, descriptions, constraints, and relationships.

**<Definitions>**
Semantic definitions, Business keys, Join keys, PII

**<Quality>**
Testable checks that validate data promises and expectations, e.g. uniqueness, consistency, freshness, …

**<SLAs>**
Promises and expectations around availability, support, retention, backups

**<Security>**
Who may access the data? Approval necessary? Data privacy limitations?

## Consumer(s)
Who is the customer of the Data Product?

## Use Case(s)
We believe that <…>
Will help achieving <…>
We know, we are getting there based on <…>, <…>, <…>

## Classification
The archetype of the provided data (source-aligned, aggregate, consumer-aligned)

## Ubiquitous Language
Context-specific domain terminology (relevant for Data Product)

## Follow-up
e.g. create a Data Contract YAML (datacontract.com)
Date:

Data Product Canvas

datamesh-architecture.com
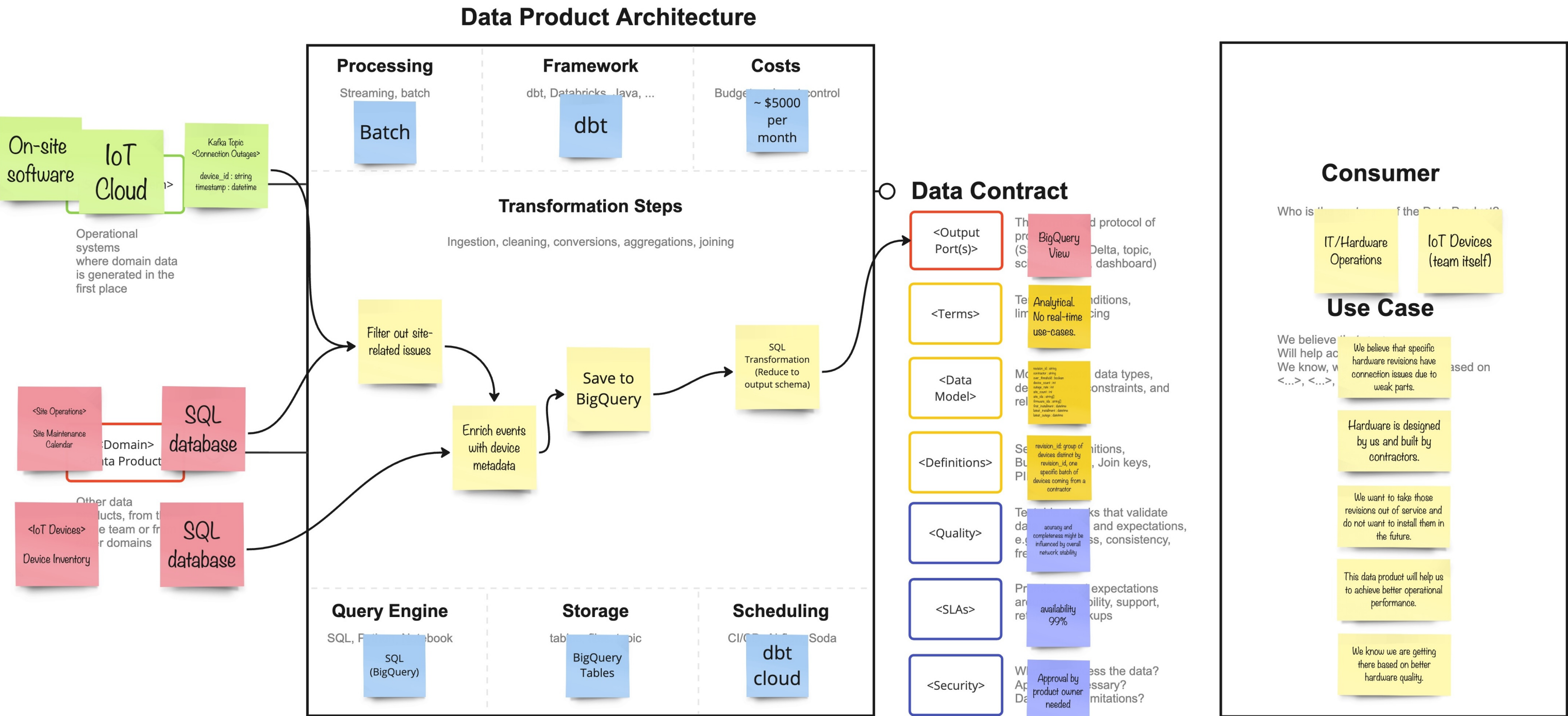
# The Birth-Certificate

In order to enhance operational performance, our IoT device company plans to develop a data product that identifies hardware revisions prone to connection issues caused by weak components. By analyzing data from our designed and contractor-built hardware, we can remove these problematic revisions from service and prevent their installation in the future, leading to improved hardware quality and overall operational efficiency.

| **Domain** *IoT Devices* | **Data Product** *Error prone device revisions* | **Date** *30.5.2023* |
| --- | --- | --- |

## Data Product Architecture

### Processing
Streaming, batch

Batch

### Framework
dbt, Databricks, Java, ...

dbt

### Costs
Budget, cost control

~ $5000 per month

**On-site software** · **IoT Cloud**

Kafka Topic <Connection Outages>
device_id : string
timestamp : datetime

Operational systems where domain data is generated in the first place

### Transformation Steps
Ingestion, cleaning, conversions, aggregations, joining

Filter out site-related issues

Enrich events with device metadata

Save to BigQuery

SQL Transformation (Reduce to output schema)

<Site Operations>
Site Maintenance Calendar

<Domain> Data Product

SQL database

Other data products, from the team or from other domains

<IoT Devices>
Device Inventory

SQL database

### Query Engine
SQL, Python, Notebook

SQL (BigQuery)

### Storage
table, file, topic

BigQuery Tables

### Scheduling
CI/CD, U-Flow, Soda

dbt cloud

## Data Contract

**<Output Port(s)>** — The standard protocol of providing data (SQL, files, Delta, topic, schema, dashboard)

BigQuery View

**<Terms>** — Terms and conditions, limitations, pricing

Analytical. No real-time use-cases.

**<Data Model>** — Model, data types, descriptions, constraints, relations

revision_id group of devices distinct by revision_id, one specific batch of devices coming from a contractor

**<Definitions>** — Semantic definitions, Business terms, Join keys, PII

**<Quality>** — Tests and checks that validate data quality and expectations, e.g. completeness, consistency, freshness

accuracy and completeness might be influenced by overall network stability

**<SLAs>** — Promised expectations around availability, support, response, backups

availability 99%

**<Security>** — Who can access the data? Approvals necessary? Data classification/limitations?

Approval by product owner needed

## Consumer
Who is the consumer of the Data Product?

IT/Hardware Operations

IoT Devices (team itself)

## Use Case
We believe that...
Will help ac...
We know, w... ...ed on
<...>, <...>...

We believe that specific hardware revisions have connection issues due to weak parts.

Hardware is designed by us and built by contractors.

We want to take those revisions out of service and do not want to install them in the future.

This data product will help us to achieve better operational performance.

We know we are getting there based on better hardware quality.

## Classification
The class of the exposed data ...med, ag... ...mer-al...

consumer-aligned

## Ubiquitous Language

device: unit of IoT hardware

revision: group of devices distinct by revision_id, one specific batch of devices coming from a contractor

firmware: software running on the iot device

site: a geographical address

installment: device was installed at site

outage: device was not able to communicate

## Follow-up
e.g., ...
YAML...
Date:...

5.5.2023
Let's create a data contract YAML

datamesh-architecture.com

# Specified in YAML

```yaml
dataProductSpecification: 0.0.1
id: shelf_warmers
info:
  title: Shelf Warmers
  description: Calculated shelfwarmers. Read about calculation in docs.
  status: active
  archetype: consumer-aligned
  owner: fulfillment
  domain: ecommerce
inputPorts: []
outputPorts:
  - id: glue_catalog_database_shelf_warmers_v1
    name: 'Glue Catalog Database: Shelf Warmers (v1)'
    description: All Shelf Warmers represented as a Glue Catalog table
    dataContractId: shelf_warmers_v1
    type: Glue
    status: active
    location: arn:aws:glue:eu-central-1:528115139298:table/fulfillment-shelf-warmers/shelf_warmers
    containsPii: false
    links:
      Athena Query Editor: https://eu-central-1.console.aws.amazon.com/athena/home?region=eu-central-1#
      Glue Table: https://eu-central-1.console.aws.amazon.com/glue/home?region=eu-central-1#/v2/data-ca
    custom:
      platform: aws
    tags:
      - glue
      - athena
```

Data Product Specification
(https://dataproduct-specification.com/)

Created by INNOQ

# Implementations

- Very different ways

- Depends highly on the data platform the company is using

- Typical: Group all code and YAMLs in a git repo per data product

- Examples:

  - Git repository with dbt that is scheduled in the CI/CD pipeline and runs queries in snowflake

  - Java Application sourcing data from a REST-API and pushing it on an AWS S3 bucket

  - Databricks Asset Bundle with pipelines written in Python

  - ...

All fine? Sadly, no.

# Some call them "Pure Data Products"



|  | Pure Data Products | | | Enhanced Data Products | | |
|---|---|---|---|---|---|---|
| **Types of Data Products** | Source-aligned Data | Merged Domain Data | Consumer-aligned Data | Decision Support | Applied Algorithms | Services/ Applications |
| **Examples** | CRM customer data | Customer „golden record" | Customer retention KPIs | Customer value dashboard | Customer classification | Customer app for smartphone |

©BARC | barc.com

If you ask 5 different people,
you get 6 different answers...

# Major Differences

deployment unit (like a Docker container)

logical unit (like a git repository)

data set + quality (like a database table with metadata on guarantees)

anything that heavily consumes data (like a report or application)

And we haven't talked about any details, like access-request-workflows and breaking-change-processes...

# And many formats as well

- From vendors and standard bodies that are similar: Data Products Ontology (DPROD from OMG), Data Product Descriptor Specification (DPDS from Quantyca), Data Product Specification (DPS from witboost), Data Product Specification (DPS from INNOQ), Open Data Product Specification (ODPS from LF) …

- A popular one, the Open Data Product Specification, follows a different definition of what a data product is (data set + guarantees) … which makes everything even more confusing. (https://opendataproducts.org/)

- There is no clear leader yet

# Our Answer

# The Open Data Product Standard

- At the Linux Foundation, Data & AI, as part of Bitol

- Bitol: standards around data products, data contracts, data mesh, …

- Status: Work in Progress, release overdue ;-)



## Roadmap

Here is a quick look at the Bitol roadmap.

| | | | |
|---|---|---|---|
| **ODCS** | Open Data Contract Standard | v3.0.2 | March 2025 |
| **ODPS** | Open Data Product Standard | | Q2 2025 |
| **ODMS** | Open Data Mesh Standard | | 2025 |
| **OORS** | Open Observability Results Standard | | 2025 |
| **OOCS** | Open Orchestration & Control Standard | | 2025 |

In summary: data product is a catchy term
Everybody wants to use it for their own agenda

And now you have to clarify the term before every conversation...

# Part 3: Data Contracts

**APIs**

As every presentation needs an image generated with AI, here you go.

# APIs

- REST-API specified with OpenAPI



- Messages and Events specified with AsyncAPI



- What about sharing large datasets? How to specify these APIs?

  - Examples: JSON on AWS S3, SQL tables on BigQuery, Iceberg files on Azure One Lake, SQL views on Snowflake, Delta Live Tables on Databricks, CSV on sftp

# Existing Specs are Lacking

- data structure (string with length 5 is implemented with VARCHAR(5))

- data quality on columns (column is nullable, but only 3% null values max)

- terms and conditions (can I use the data for my use case?)

  - data classifications and PII anonymization

- Service-level agreements (freshness, latency, retention, ...)

- Semantics (what the column really is about)

**So we need something that fills this gap.**

**Terminology**

A **data contract** is not a **contract** as a mutual agreement, it is rather an **offer** to potential consumers.
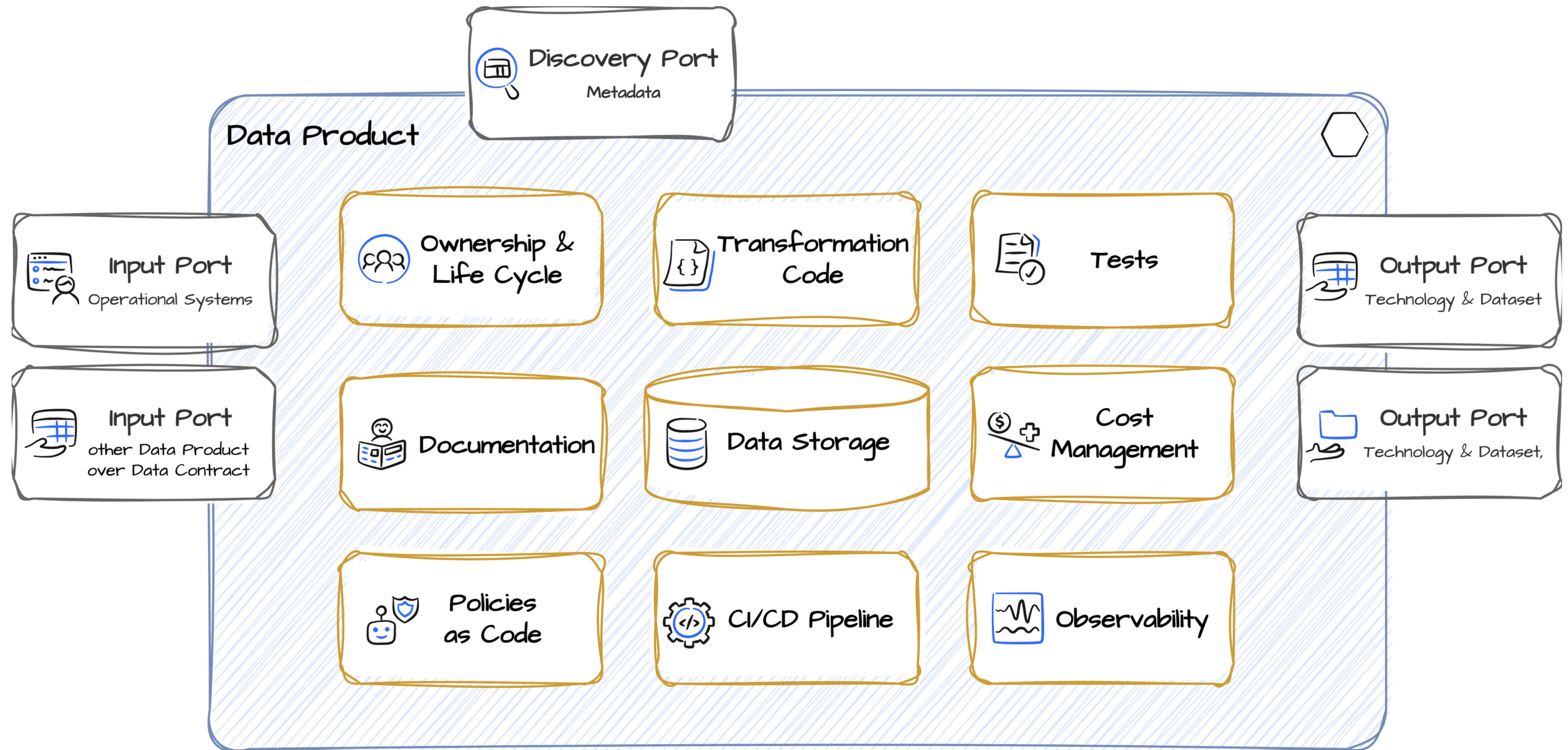
# Cardinality

A data contract is a **one-to-many** relationship.

**One** data provider to **many** data consumers.

# What is a data contract?

- Creates a link between a data producer and data consumers.

- Creates a link between business (logical representation of the data) and technology (its physical implementation).

- Describes meta meta data: rules, quality, and behavior.

- It's like an API Contract, but for data. And contains much more.


- **A data contract is the source of truth for your metadata.**

# Data Contracts Protect Output Ports



**Discovery Port**
Metadata

**Data Product**

**Input Port**
Operational Systems

**Input Port**
other Data Product
over Data Contract

Ownership &
Life Cycle

Transformation
Code

Tests

Documentation

Data Storage

Cost
Management

Policies
as Code

CI/CD Pipeline

Observability

**Output Port**
Technology & Dataset

**Output Port**
Technology & Dataset,

datamesh-architecture.com

# What are the problems it solves?

- Normalizing and keeping documentation relevant.

- Bringing quality data in AI workflows. Describing service-level expectations.

- Easing data & tools integration.

- Ending painful data discovery.

- Enabling data product thinking.

**Slide used by courtesy of Jeon-George Perrin**

Data contracts [..] are a bit like tax returns. You have to do them. Many people don't feel like doing them and some people even try to avoid them altogether. For us, data contracts form the foundation of a living data ownership culture.

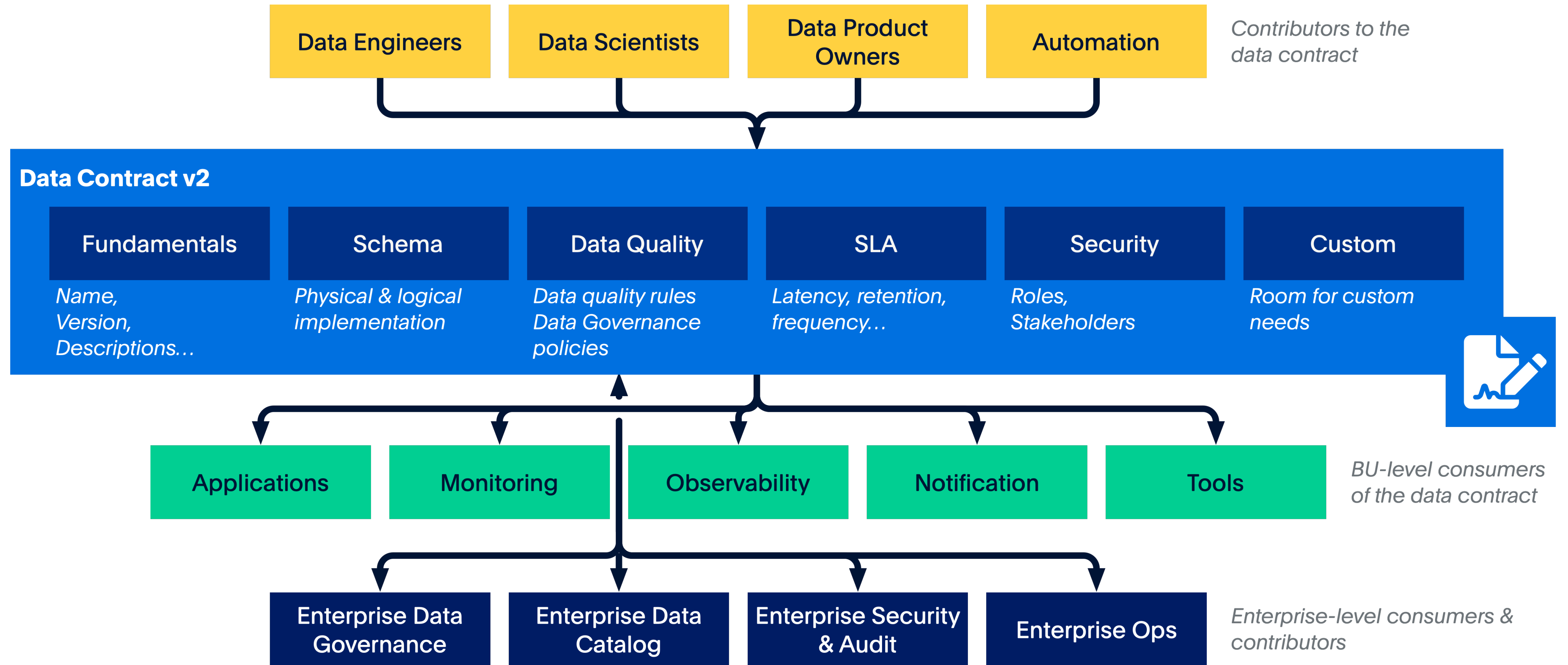*– David Brandstädter, Director Data Enablement*
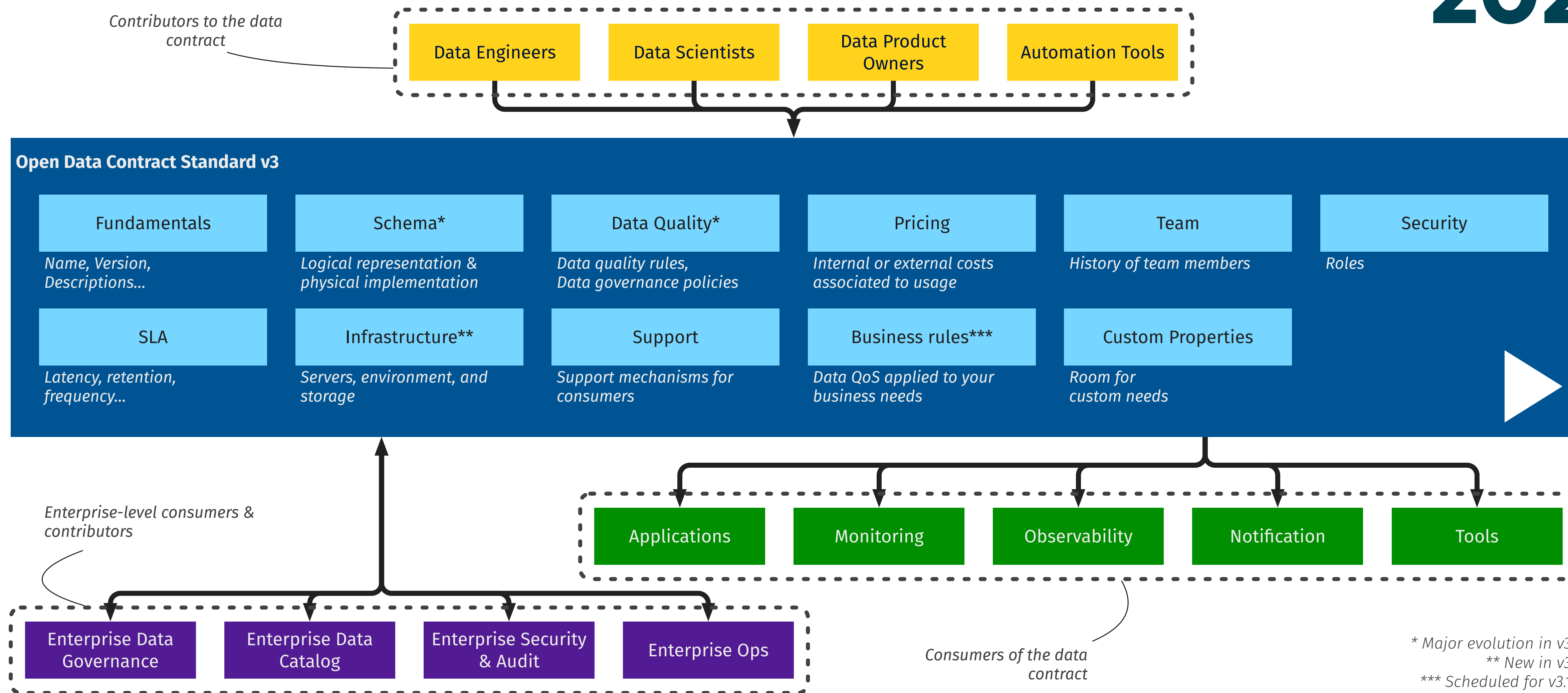*Lidl eCommerce*

# What happened ...

... with the rise of data products, the need for data contracts grew tremendously ...

... and many many company created their own data contract format ...

... and every vendor did the same ...

# Open Data Contract Standard
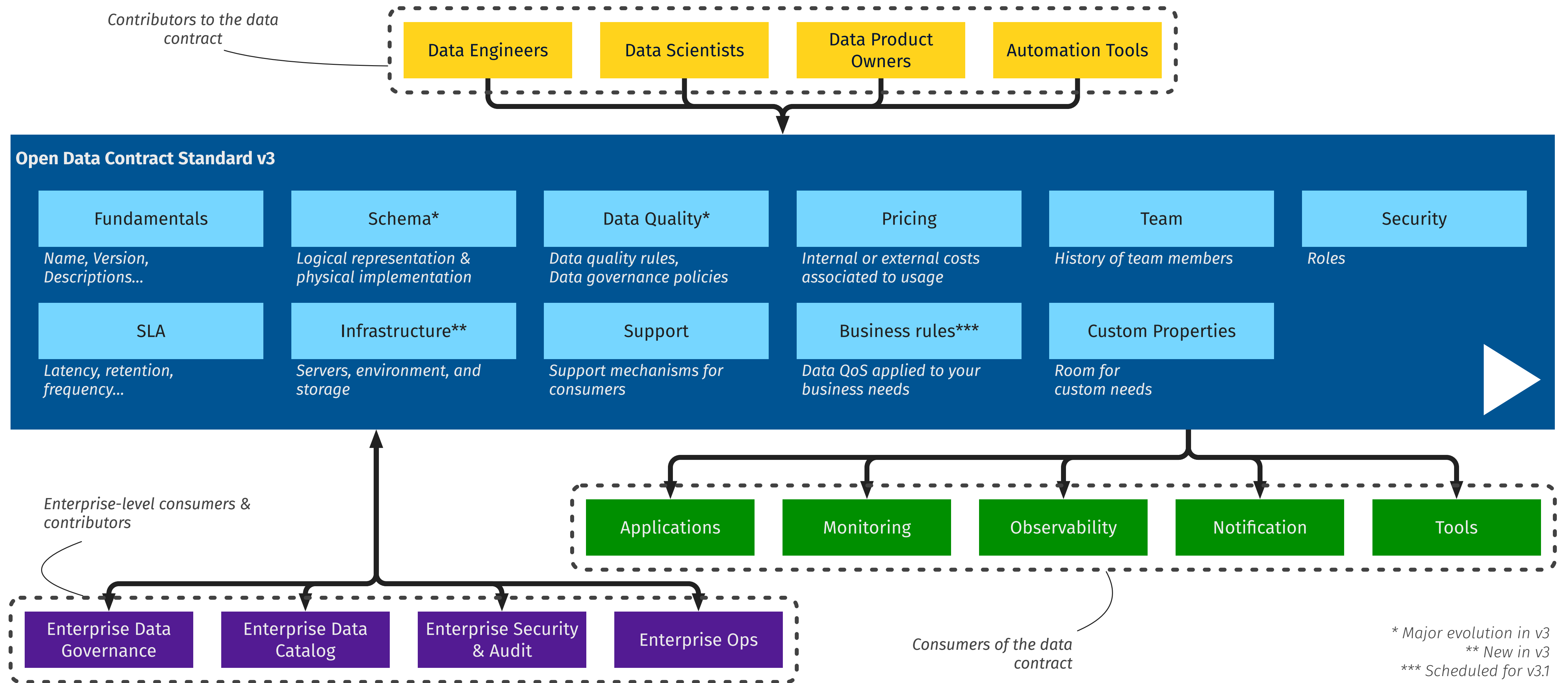## to the rescue (this time just-in-time)

2023

**Data Engineers** · **Data Scientists** · **Data Product Owners** · **Automation** — *Contributors to the data contract*

**Data Contract v2**

| Fundamentals | Schema | Data Quality | SLA | Security | Custom |
|---|---|---|---|---|---|
| *Name, Version, Descriptions…* | *Physical & logical implementation* | *Data quality rules Data Governance policies* | *Latency, retention, frequency…* | *Roles, Stakeholders* | *Room for custom needs* |

**Applications** · **Monitoring** · **Observability** · **Notification** · **Tools** — *BU-level consumers of the data contract*

**Enterprise Data Governance** · **Enterprise Data Catalog** · **Enterprise Security & Audit** · **Enterprise Ops** — *Enterprise-level consumers & contributors*

**https://github.com/paypal/data-contract-template**

PayPal

**2025**

Contributors to the data contract

Data Engineers | Data Scientists | Data Product Owners | Automation Tools

**Open Data Contract Standard v3**

| Fundamentals | Schema* | Data Quality* | Pricing | Team | Security |
|---|---|---|---|---|---|
| Name, Version, Descriptions... | Logical representation & physical implementation | Data quality rules, Data governance policies | Internal or external costs associated to usage | History of team members | Roles |

| SLA | Infrastructure** | Support | Business rules*** | Custom Properties |
|---|---|---|---|---|
| Latency, retention, frequency... | Servers, environment, and storage | Support mechanisms for consumers | Data QoS applied to your business needs | Room for custom needs |

Enterprise-level consumers & contributors

Enterprise Data Governance | Enterprise Data Catalog | Enterprise Security & Audit | Enterprise Ops

Applications | Monitoring | Observability | Notification | Tools

Consumers of the data contract

* Major evolution in v3
** New in v3
*** Scheduled for v3.1

**https://github.com/bitol-io/open-data-contract-standard**

Bitol | LF AI & DATA

**ODCS (Open Data Contract Standard)** can be found & used at 30+ companies:

- **End users:** peer-to-peer payment leader, major cable company, major retailers, SMB to Fortune 500.
- Software **Vendors**: several data-oriented startups and vendors in Europe, NA, and APAC.
- **Service providers** in Europe and NA.
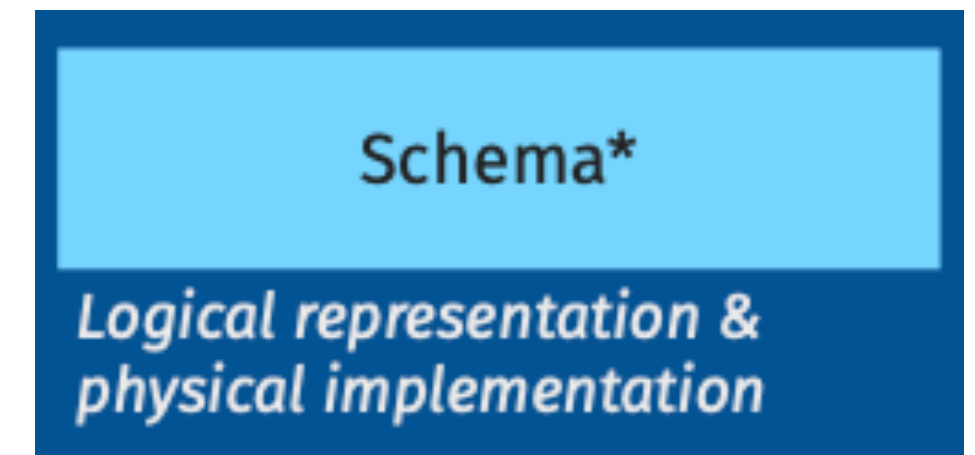- A **going community** behind the standard.
- **Academia** (JADS, HTW).

**Industry-led governance**.

Contributors to the data contract

**Data Engineers**  **Data Scientists**  **Data Product Owners**  **Automation Tools**

**Open Data Contract Standard v3**

| Fundamentals | Schema* | Data Quality* | Pricing | Team | Security |
|---|---|---|---|---|---|
| Name, Version, Descriptions... | Logical representation & physical implementation | Data quality rules, Data governance policies | Internal or external costs associated to usage | History of team members | Roles |

| SLA | Infrastructure** | Support | Business rules*** | Custom Properties |
|---|---|---|---|---|
| Latency, retention, frequency... | Servers, environment, and storage | Support mechanisms for consumers | Data QoS applied to your business needs | Room for custom needs |

Enterprise-level consumers & contributors

**Applications**  **Monitoring**  **Observability**  **Notification**  **Tools**

**Enterprise Data Governance**  **Enterprise Data Catalog**  **Enterprise Security & Audit**  **Enterprise Ops**

Consumers of the data contract

* Major evolution in v3
** New in v3
*** Scheduled for v3.1

**https://github.com/bitol-io/open-data-contract-standard**

Bitol | LF AI & DATA

```yaml
apiVersion: v3.0.0
kind: DataContract
id: urn:datacontract:checkout:orders-latest
name: Orders Latest
version: 2.0.0
status: active
description:
  usage: |
    Data can be used for reports, analytics and machine learning use cases.
    Order may be linked and joined by other tables'
  limitations: |
    Not suitable for real-time use cases.
    Data may not be used to identify individual customers.
    Max data processing per day: 10 TiB'
```

```yaml
schema:
- name: orders
  physicalName: orders
  logicalType: object
  physicalType: table
  description: One record per order. Includes cancelled and deleted orders.
  properties:
  - name: order_id
    businessName: Order ID
    logicalType: string
    physicalType: text
    description: An internal ID that identifies an order in the online shop.
    isNullable: false
    isUnique: true
    classification: restricted
    examples:
    - 243c25e5-a081-43a9-aeab-6d5d5b6cb5e2
    primaryKey: true
    primaryKeyPosition: 1
    customProperties:
    - property: pii
      value: true
    tags:
    - orders
```



Schema*

Logical representation &
physical implementation

elements

objects

properties

**Orders**
id
date
customer_id
delivery_ctry
...

logicalType: object
physicalType: table
name: Orders
physicalName: order

logicalType: integer
physicalType: serial
name: Id
physicalName: id

**OrderLines**
id
order_id
qty
product_id
...

logicalType: object
physicalType: table
name: OrdersLines
physicalName: or_line

logicalType: number
physicalType: float(10)
name: Quantity
physicalName: qty

Bitol | LF AI & DATA

```yaml
quality:
- type: sql
  description: The maximum duration between two orders should be less that 3600
    seconds
  query: |
    SELECT MAX(duration) AS max_duration FROM (SELECT EXTRACT(EPOCH FROM (order_timestamp
    - LAG(order_timestamp) OVER (ORDER BY order_timestamp))) AS duration FROM orders)
  mustBeLessThan: 3600
- type: sql
  description: Row Count
  query: 'SELECT count(*) as row_count FROM orders'
  mustBeGreaterThan: 5
```

```yaml
slaProperties:
- property: generalAvailability
  value: The server is available during support hours
- property: retention
  value: P1Y
support:
- channel: other
  url: https://teams.microsoft.com/l/channel/example/checkout
servers:
- server: production
  type: s3
  environment: prod
  format: json
  delimiter: new_line
  location: s3://datacontract-example-orders-latest/v2/{model}/*.json
  roles:
  - name: analyst_us
    description: Access to the data for US region
  - name: analyst_cn
    description: Access to the data for China region
customProperties:
- property: owner
  value: Checkout Team
```

**SLA**

*Latency, retention, frequency...*

**Infrastructure\*\***

*Servers, environment, and storage*

**Custom Properties**

*Room for custom needs*

## Code Generation

- Java
- Python in Pydantic
- dbt Models and Sources
- SQL DDL and Queries

## Test

- Compare contract with real data
- Breaking data detection in PR
- Breaking metadata detection
- Continuous Monitoring
- Consumer-driven Contract Testing

## Metadata Distribution

- Metastores: Hive, ...
- Data Catalogs: Colibra, ...
- Data Contract Catalog: Data Mesh Manager, ...
- Software Catalogs: LeanIX, ...

# Automate all the things!

## Infrastructure Provisioning

- Output Port (S3 Bucket, ...)
- Input Port (dbt sources.yml)
- Transformations (anonymisation)
- Access Control (IAM permissions)

## Collaboration

- Contract-First (in workshop)
- Data-First (import from ...)
- Semantics

## Governance

- Policies (naming conventions, ...)
- Schema Evolution (Notice period)
- Usage Agreements
- Approval Workflows

# Data Contract CLI

# 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
or odcs.yaml

diff

import

export

test

SQL DDL · Avro · JSON Schema · Protobuf · BigQuery · Unity Catalog · AWS Data Catalog · ODCS

SQL DDL · HTML · Avro · RDF · dbt · SodaCL · Terraform · ODCS

AWS S3 · BigQuery · Azure · databricks · snowflake · Kafka

github.com/datacontract/cli

**Usage:** `datacontract [OPTIONS] COMMAND [ARGS]...`

The datacontract CLI is an open source command-line tool for working with Data Contracts (https://datacontract.com).
It uses data contract YAML files to lint the data contract, connect to data sources and execute schema and quality tests, detect breaking changes, and export to different formats.

```
┌─ Options ─────────────────────────────────────────────────────────────────────┐
│  --version        Prints the current version.                                  │
│  --help           Show this message and exit.                                  │
└────────────────────────────────────────────────────────────────────────────────┘
```

```
┌─ Commands ────────────────────────────────────────────────────────────────────┐
│  init       Create an empty data contract.                                      │
│  lint       Validate that the datacontract.yaml is correctly formatted.         │
│  test       Run schema and quality tests on configured servers.                 │
│  export     Convert data contract to a specific format. Saves to file specified by `output` option if present, otherwise prints to stdout. │
│  import     Create a data contract from the given source location. Saves to file specified by `output` option if present, otherwise prints to stdout. │
│  publish    Publish the data contract to the Data Mesh Manager.                 │
│  catalog    Create a html catalog of data contracts.                            │
│  breaking   Identifies breaking changes between data contracts. Prints to stdout. │
│  changelog  Generate a changelog between data contracts. Prints to stdout.      │
│  diff       PLACEHOLDER. Currently works as 'changelog' does.                   │
│  api        Start the datacontract CLI as server application with REST API.     │
└────────────────────────────────────────────────────────────────────────────────┘
```

DEMO

# Modern Data Governance

## Responsibility

### Data Product Owner

Data is owned decentralized by business & IT experts where data is generated
Product owners are responsible for what happens with their data

## Concepts & Tools

### Data Contracts

Define the syntax, semantics, quality, and terms of use as YAML

### Data Marketplace

Data discovery with a self-service access request workflow

### Global Policies

The conventions and rules of play for data on the data platform

## Automation

### Contract Enforcement

Test that data products correctly implement the data contract

### Automated Permission Granting

Give table access based on access request approvals
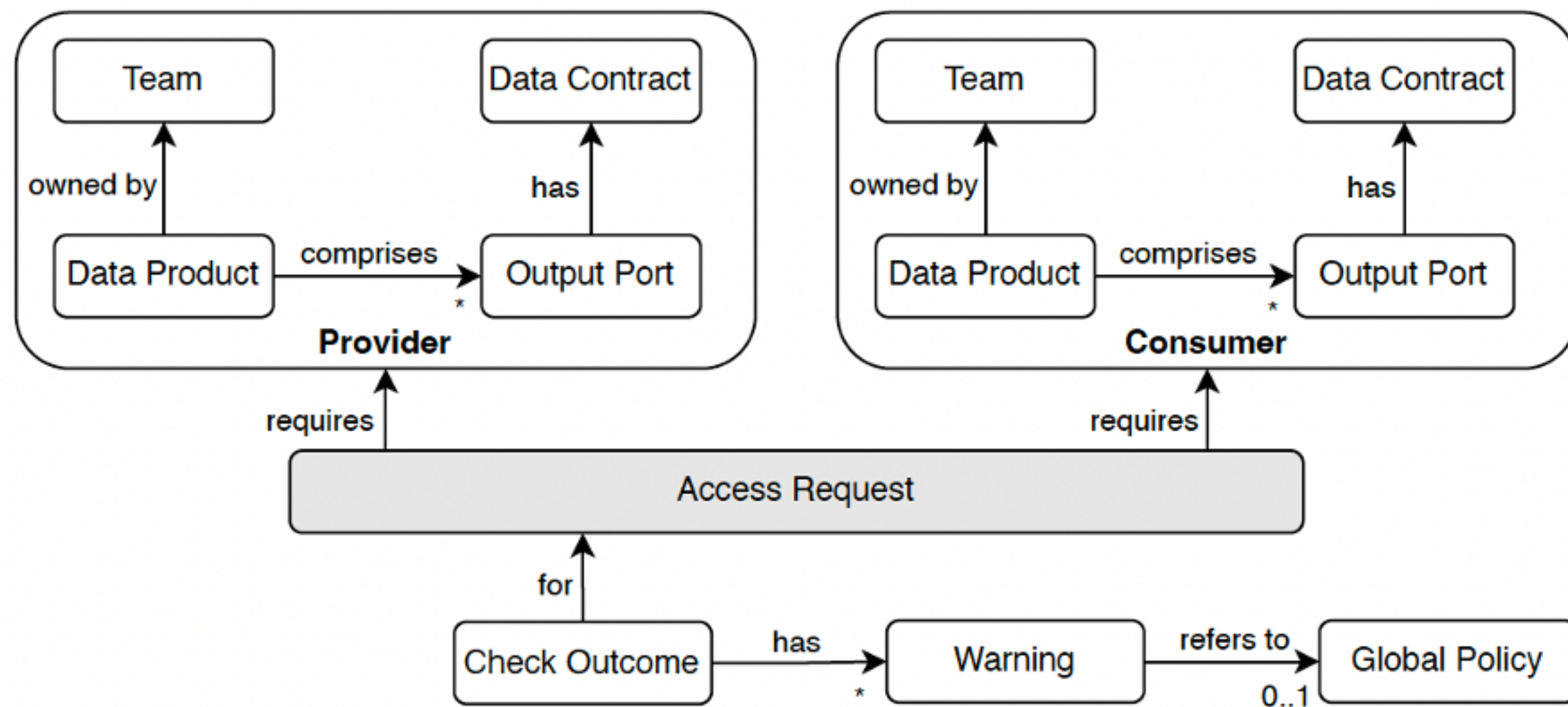
### AI-based Policy Checking

Check that policies are correctly adopted by data product owners

datamesh-manager.com

In summary: data contracts bring API thinking and specification into the data world.

Swagger was first released in 20211 and is everywhere today. ODCS was first released in 2025 and will be everywhere in ???

My prediction: This will be huge!

# Bonus-Part: Data Governance with AI (with Arif)

# Modern Data Governance

## Responsibility

**Data Product Owner**

Data is owned decentralized by business & IT experts where data is generated
Product owners are responsible for what happens with their data

## Concepts & Tools

| Data Contracts | Data Marketplace | Global Policies |
|---|---|---|
| Define the syntax, semantics, quality, and terms of use as YAML | Data discovery with a self-service access request workflow | The conventions and rules of play for data on the data platform |

## Automation

| Contract Enforcement | Automated Permission Granting | AI-based Policy Checking |
|---|---|---|
| Test that data products correctly implement the data contract | Give table access based on access request approvals | Check that policies are correctly adopted by data product owners |

datamesh-manager.com

# Use Case: Data Access Requests

# LLMs to the rescue!



Requester

Data product with
sensitive data

Privacy policies etc.

Let the
robot do
tedious
reading
work

**System prompt**

1) **Task.** We describe the main task: analyzing access requests.
2) **Persona.** We asked the model to adopt the persona of a Data Governance Expert.
3) **Steps.** We describe the six steps how to analyze an access request.

**User Prompt**

1) **Access Request** that needs to be analyzed (YAML)
2) **Provider** side of the access request, including the providing data product, the relevant output port, the data contract, and the providing team (YAML)
3) **Consumer** side of the access request, including the consuming data product, all output ports, data contracts, and the consuming team (YAML)
4) **Global Policies** governing the data mesh (text).
5) **Detailed Instructions** about the task, the requirements, and additional constraints.
6) **Required Elements** of the output with an explanation. The structure of the required elements was enforced using the "Structured Outputs JSON mode," cf. https://platform.openai.com/docs/guides/structured-outputs

In summary: Data Goverance perfectly fits to the strengths of LLMs and can help in federation data management.

# Thank You! Want to learn more?

- Talk to me at this event, or later via LinkedIn  /in/simonharrer

- Try out cli.datacontract.com with an ODCS data contract

- And, of course, give us a star on Github at

  github.com/datacontract/datacontract-cli
  github.com/bitol-io/open-data-contract-standard