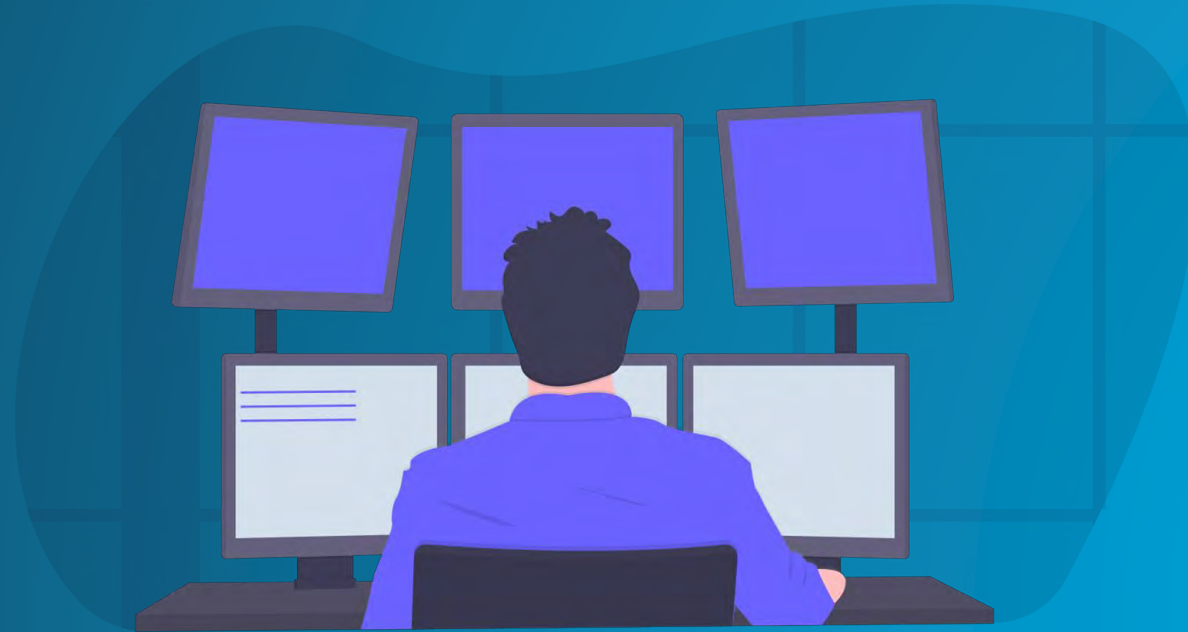




Data Mesh Infrastructure as a Service with Stackable Data Platform

Sönke Liebau



Stefan Igel

Stackable in a Nutshell

Founded

2020

 OpenCore

 b.telligent

IONOS

Stackable Data Platform

- Open Source
- Infrastructure as Code
- Cloud-native (Kubernetes)
- On-Premises, Cloud, Hybrid

Our Customers







IONOS

opencorporates

Our Team: 20 People

International
in Germany & Europe

Our Services

- Product Support
- Big Data Consulting
- Trainings

Network - Collaborations

 OSB Open Source
Business
ALLIANCE



KI BUNDESVERBAND

gaia-x



bitkom

eco



Stackable Data Platform



Stackable – An Open-Source Data Platform Ready for Data Mesh



Stackable

Data Visualisation



Analytics & AI



Federated Queries



Data Processing



Storage



Data Ingestion



Infrastructure
Orchestration



Security



Open Policy Agent

Monitoring

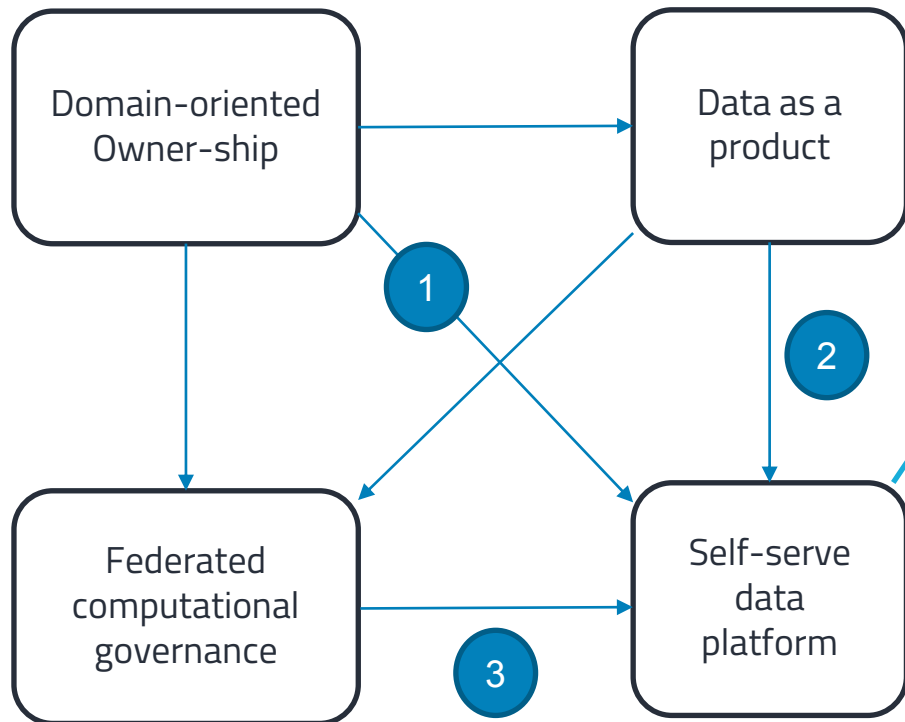


Beyond the Lake - Data Mesh and its Principles*

Data Mesh is a decentralized sociotechnical approach to share, access, and manage analytical data in complex and large-scale environments – within or across organizations.

Zhamak Dehghani

Principals



Self-serve data platform

Data Infrastructure optimized for infrastructure utilization and performance

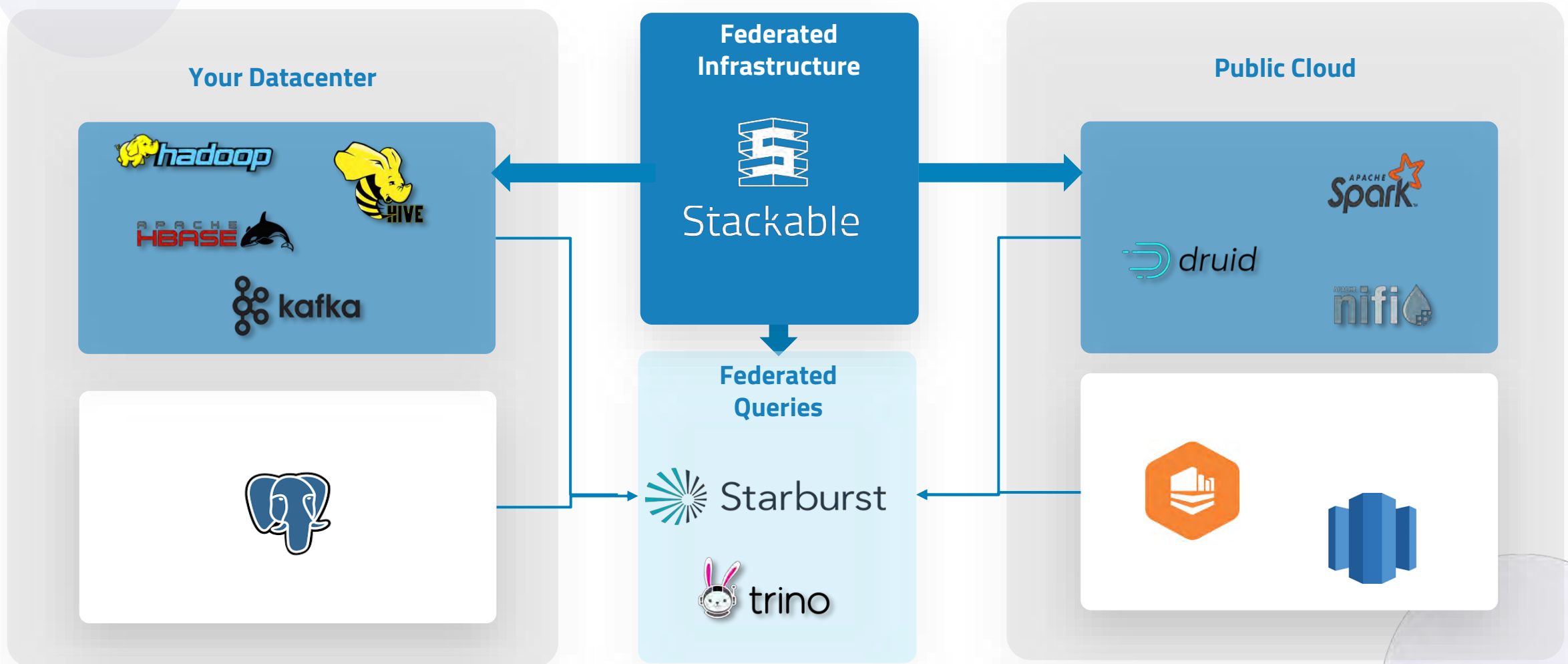


Purpose

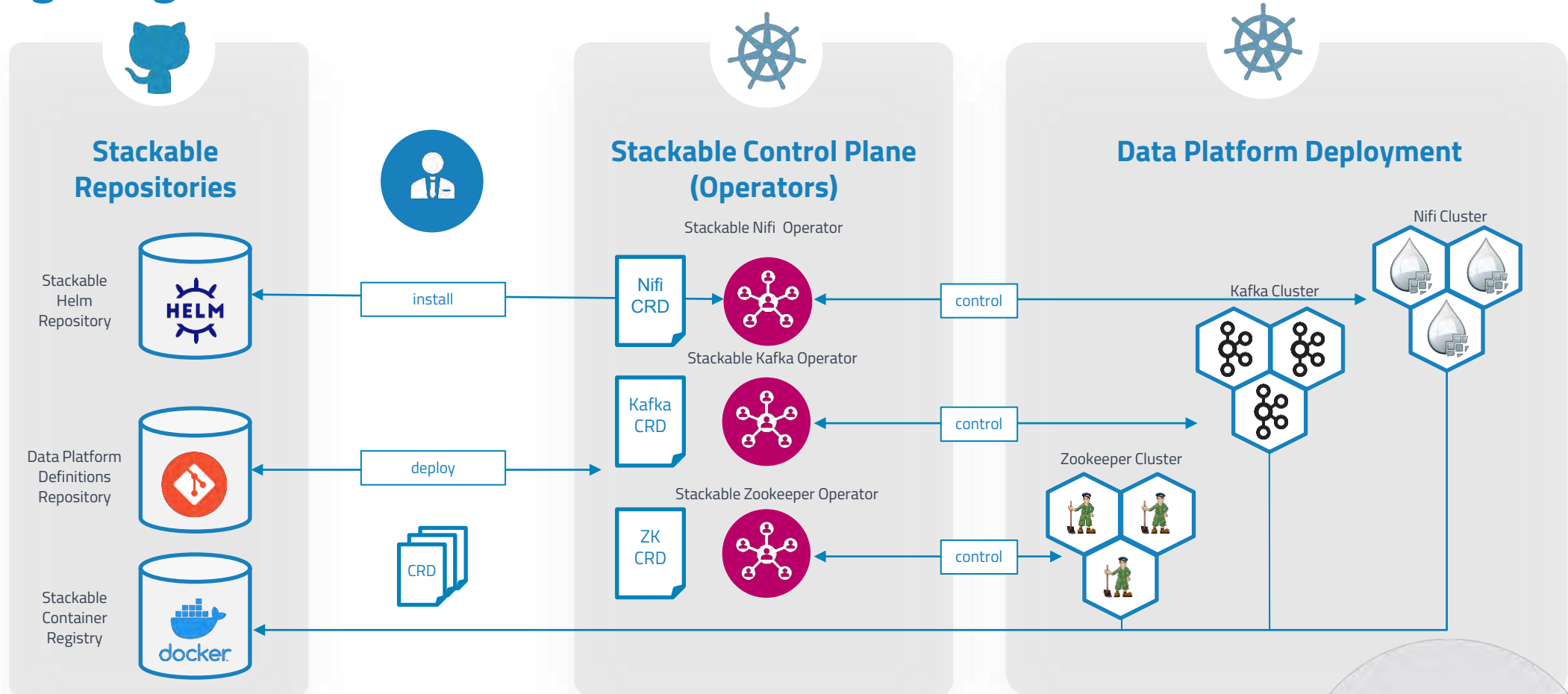
1. Empower Domain Teams
2. Reduce data product cost of ownership
3. Mesh-level consistent and reliable Policy enforcement

*nach Zhamak Dehghani: Data Mesh, O'Reilly 2022

Stackable & Starburst - a Match made in Heaven



Streaming & Big Data Infrastructure as Code on Kubernetes



How K8s and Operators Support Data Mesh Features

```
kubectl apply -f - <<EOF
---
apiVersion: kafka.stackable.tech/v1alpha1
kind: KafkaCluster
metadata:
  name: simple-kafka
spec:
  version: 2.8.1
  zookeeperConfigMapName: simple-kafka-znode
  brokers:
    roleGroups:
      brokers:
        replicas: 1
        selector:
          matchLabels:
            node: quickstart-1
---
apiVersion: zookeeper.stackable.tech/v1alpha1
kind: ZookeeperZnode
metadata:
  name: simple-kafka-znode
spec:
  clusterRef:
    name: simple-zk
    namespace: default
EOF
```

Example: Custom Resource Definition (CRD)

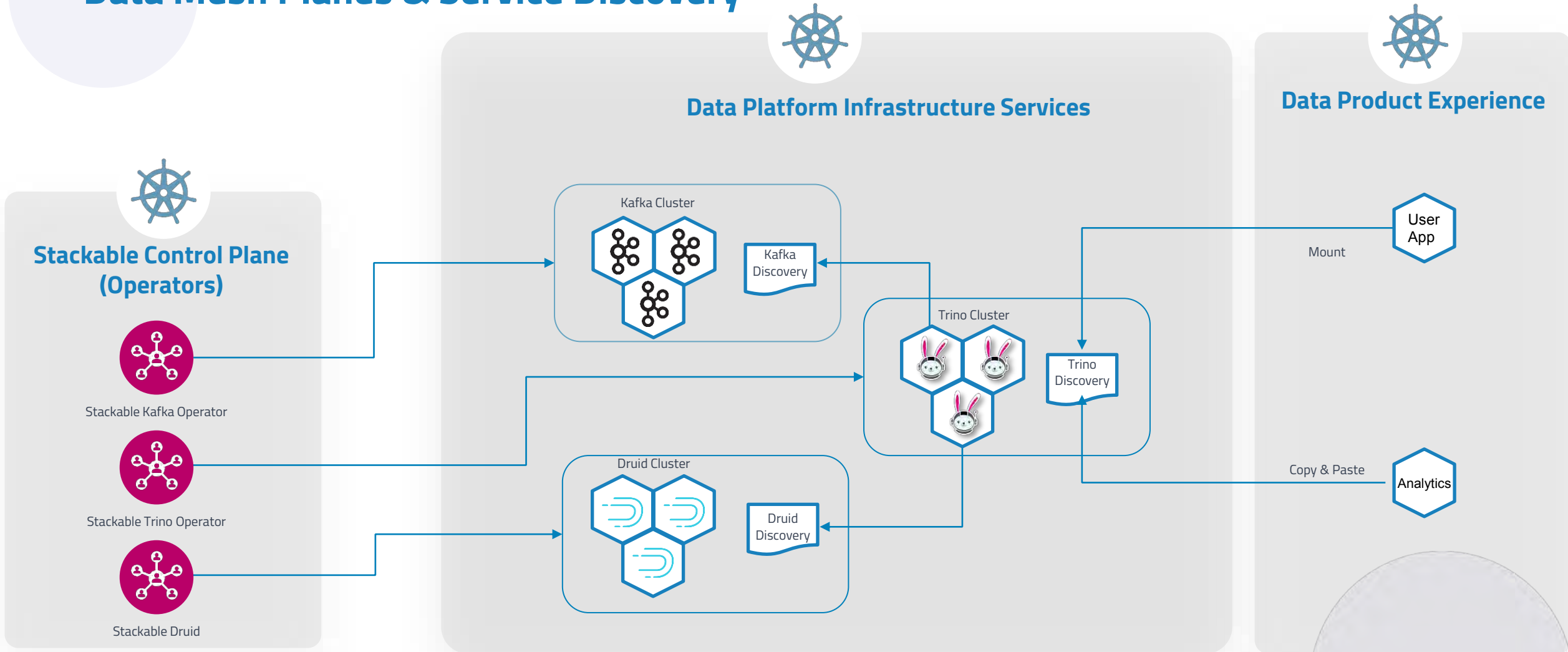
Operators are software extensions to Kubernetes that make use of custom resources to manage applications and their components.*

- Scalability of compute resources is managed by K8S
- Ship platform components as containers managed by operators
- Storage: S3 and HDFS Operators or external
- Portable, reduces vendor lock-in

- Infrastructure-as-code via CRDs
- Service Discovery
- Central secret management (certificates) by Secret Operator
- Flexible authorization (as code) through Open Policy Agent Operator
- Unified telemetry (Monitoring, Logging, Alerting) configurable via CRDs

*<https://kubernetes.io/docs/concepts/extend-kubernetes/operator/>

Data Mesh Planes & Service Discovery



Security Policies – Authorization as Code



Open Policy Agent

Source: <https://cncf-branding.netlify.app/projects/opa/>

```
1 allow {
2     # Find grants for the user.
3     some grant
4     user_is_granted[grant]
5
6     # Check if the grant permits the action.
7     input.action == grant.action
8     input.type == grant.type
9 }
10
11 user_is_granted[grant] {
12     some role in data.user_roles[input.user]
13     some grant in data.role_grants[role]
14 }
```

- Policies-as-Code (“Rego Rules”)
- Authorization plugins added to the components where possible
 - Trino
 - Apache Druid
 - Apache Kafka
- Group lookup done once!
 - We’re adding a dedicated way to look up groups
 - No more configuring a dozen tools with the same settings

How K8s and Operators Support Data Mesh Features

Challenges to address

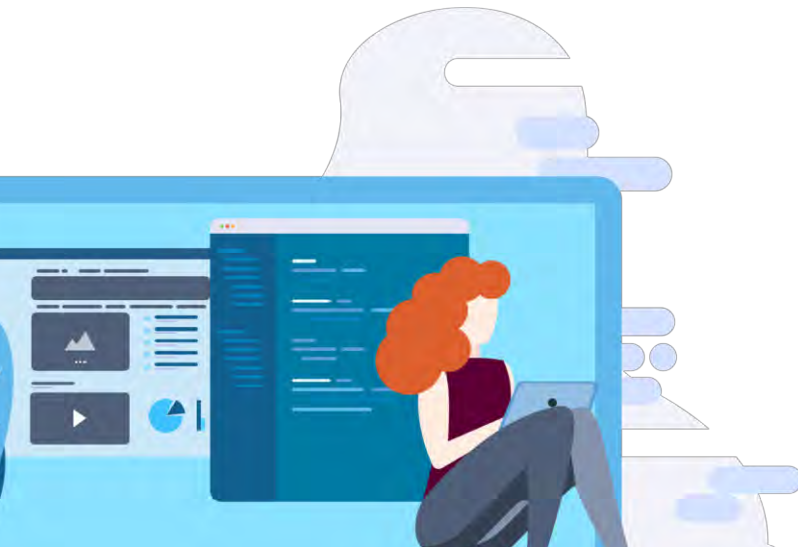
Self-Serve Data Platform

Distributed data architecture will lead to

- duplication of efforts in each domain
- Increased cost of operation
- Inconsistencies and incompatibilities across domains

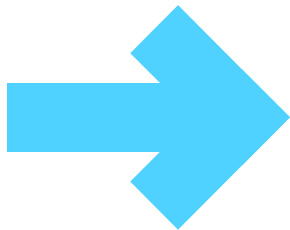
Benefits

- Standardized Logging, Monitoring, Auditing
- Similar Operators (same look and feel)
- More providing standards and examples than running a central platform
- Easy to run many instances at the same time
- Easy to define entire stacks and deploy them multiple times
 - Every team has its own stack to run
 - Can be easily shared with other teams





**Thank
you**



Contact

Dr. Stefan Igel
stefan.igel@stackable.de
+49 (160) 6171731
[linkedin.com/in/stefan-igel/](https://www.linkedin.com/in/stefan-igel/)

Summary

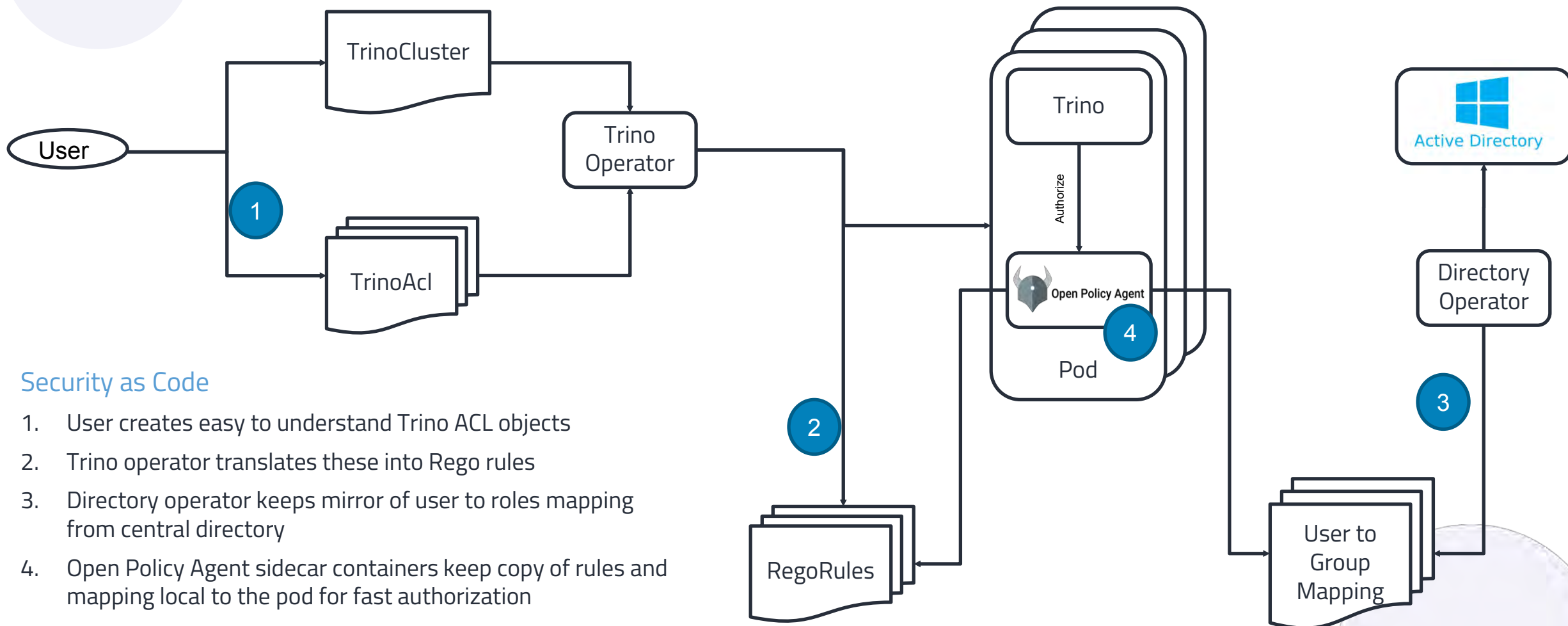
- Data platform architectures have evolved over time together with the enabling technologies
- Kubernetes has been a great paradigm shift
 - K8S provides a scalable compute platform for data workloads
 - Operators allow enforcement of standards
 - Containers and K8S facilitate data lakes and meshes
 - Some tweaks necessary to enable data lake gen 1 technologies
- Self-Serve Data platform relies heavily on X-as-Code
- Modern data platforms can be setup vendor-independent by open-source tools

Data Lake gen 1 with K8s - challenges

- Due to data locality paradigm a lot of effort was put into running calculations where the data is - K8S is not really interested in this
- A lot of the early data (or big data) tools are from a different era of computing
 - Stable Network
 - Bare Metal access
 - "Simple" DNS
 - Predictable Restarts
 - ...
- Complexity from "back then" is not gone, it is just hidden - until ...



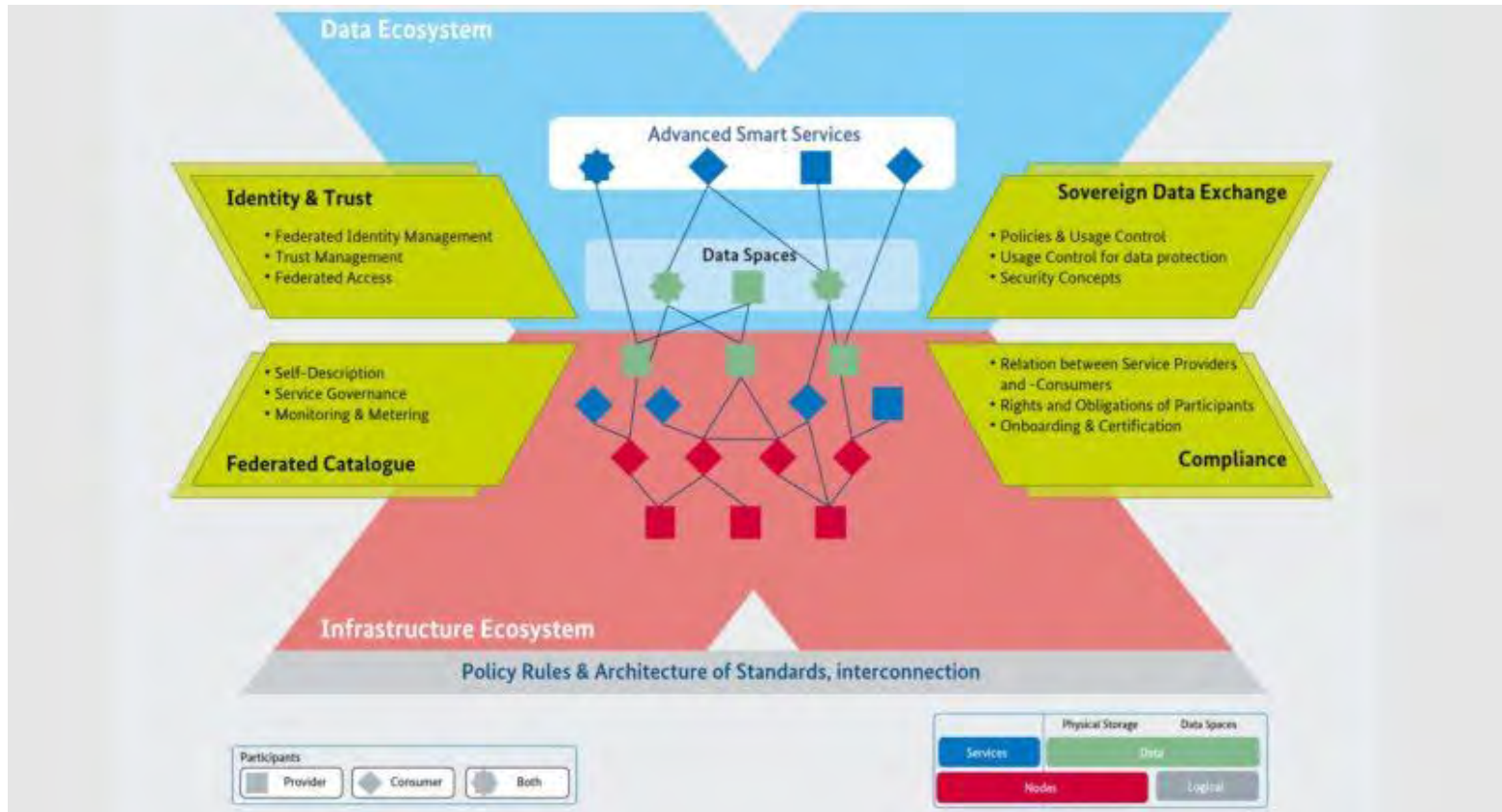
Security as Code – Putting it all together



Security as Code

1. User creates easy to understand Trino ACL objects
2. Trino operator translates these into Rego rules
3. Directory operator keeps mirror of user to roles mapping from central directory
4. Open Policy Agent sidecar containers keep copy of rules and mapping local to the pod for fast authorization

What's next? Gaia-X Sovereign Data Spaces



conceptual

Cross-organization data mesh

Federation Services

Data Sovereignty

Data Products

Governance

Technical Architecture

API based

Revival of Compute-to-Data

K8S part of the reference architecture
(SCS stack)





marispace-x Building a maritime Data Space and Connect the dots.

GAIA-X Lighthouse project

- Drive the digitization of the ocean
- Facilitate digital collaboration in marine research
- Develop a smart maritime dataspace including Cloud-, Fog- and Edge-Computing
- Leveraging GAIA-X Federation Services for data sovereignty
- Funded by BMWK

Stackable Role

- Dataspace Platform Service Layer
- Data Storage & Compute
- Data Security & Governance
- GAIA-X Interoperability



Consortium

- Universität Kiel
- Universität Rostock
- GEOMAR Helmholtz Zentrum
- Fraunhofer IGD
- EGEOS GmbH
- TrueOcean GmbH
- MacArtney Germany
- IONOS SE
- Stackable GmbH

Use Cases

- Internet of Underwater Things (IoUT)
- Offshore Wind – renewable energy
- Marine protection – ammunition in the sea
- Bio climate protection - decarbonization

Learn more



 Stackable

<https://marispacex.com/>