



Корпоративная платформа хранения
и обработки больших данных

Что такое Data Mesh и примеры реализации

Интро

Александр Ермаков

- 17+ лет в DWH & Data Engineering
- 35+ проектов по миру
- Участник Apache Software Foundation
- Участник Open Source Initiative
- СТО в Arenadata

О чем поговорим?

- Совсем коротко про Arenadata
- Текущие архитектурные подходы обработки данных
- Архитектура следующего поколения платформы данных
- Реализация платформы данных для Data Mesh

Disclaimer

- Материал подготовлен на основе уже ранее описанных методологий и подходов реализации концепции и архитектуры Data Mesh
 - Data Mesh Principles and Logical Architecture (Zhamak Dehghani)
 - How to Move Beyond a Monolithic Data Lake to a Distributed Data Mesh (Zhamak Dehghani)
 - Data Mesh Architecture - Data Mesh From an Engineering Perspective (Jochen Christ, Dr. Larysa Visengeriyeva, Dr. Simon Harrer)
 - Evolution of Data Mesh Architecture Can Drive Significant Value in Modern Enterprise (JP Morgan)
 - Deploying Data Products at the speed of the business (Jon Cooke - CTO Dataception)

Совсем коротко про Arenadata



КТО МЫ

ARENADATA — российский разработчик платформы сбора, хранения и обработки данных

Коммитер ряда Open Source проектов



Решаемые задачи



Корпоративные хранилища больших данных

- Классические КХД и хранилища больших данных с обновлением в реальном и околореальном времени.
- Виртуализация и федерация данных, доступ ко всем данным через единое окно с помощью SQL-запросов.



Разгрузка дорогостоящих инженерных систем (Teradata, Exadata)

- Импортзамещающие ПАК на базе массово-параллельных СУБД с открытым исходным кодом Arenadata DB и Arenadata QuickMarts.
- Прозрачная методика перехода, минимум рисков и сохранение всех преимуществ.



Машинное обучение и искусственный интеллект

- Поддержка R, Python и Madlib.
- Distributed ML на базе Spark.
- Возможность эффективной эксплуатации моделей в SQL-среде с помощью встроенных функций Madlib.



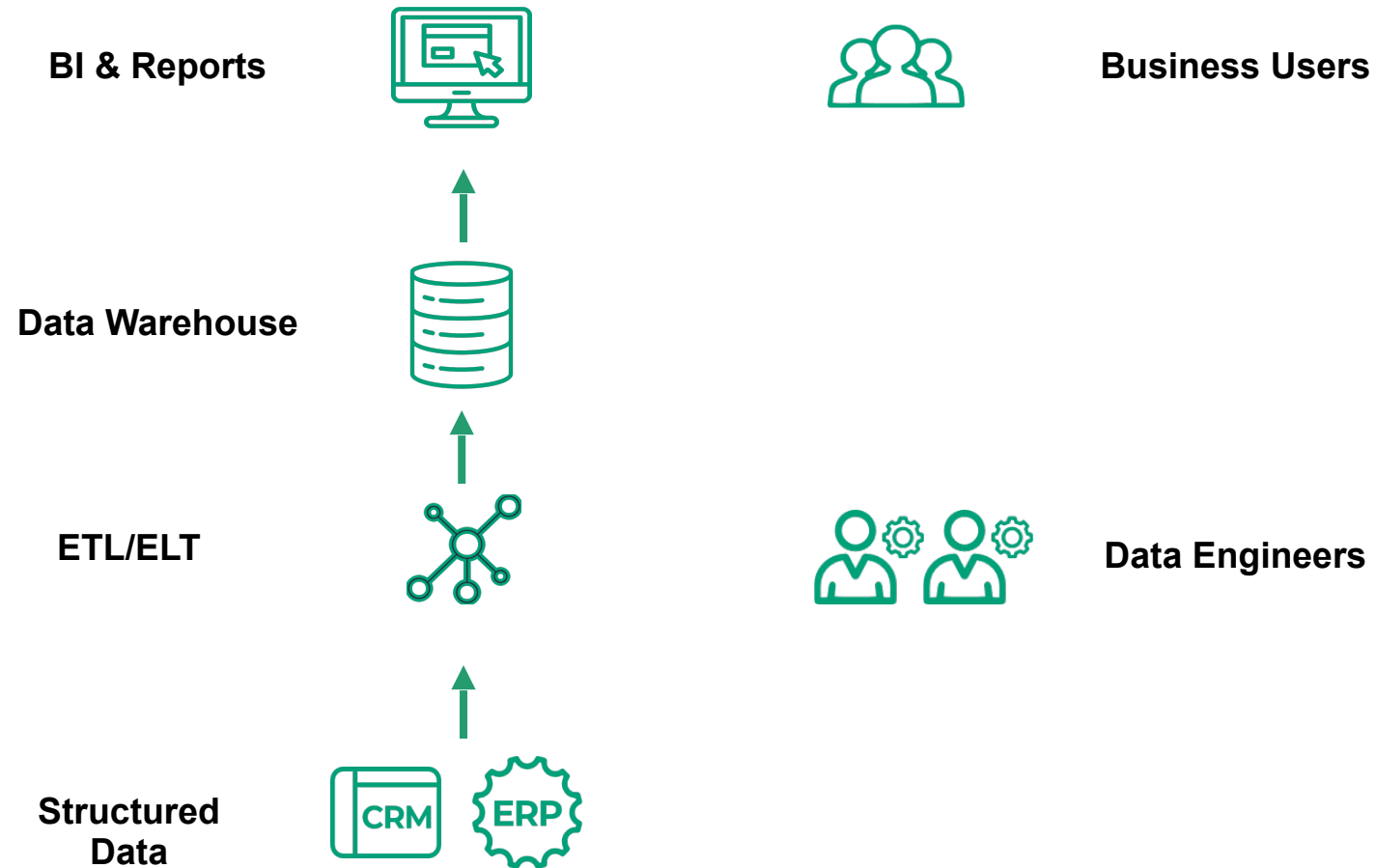
Озёра и фабрики данных

- Единый хаб всех данных компании.
- Возможность быстрого развёртывания и сворачивания зон «песочниц» для пилотных проектов и проверки статистических гипотез.
- Работа со всеми аналитическими инструментами в единой среде.

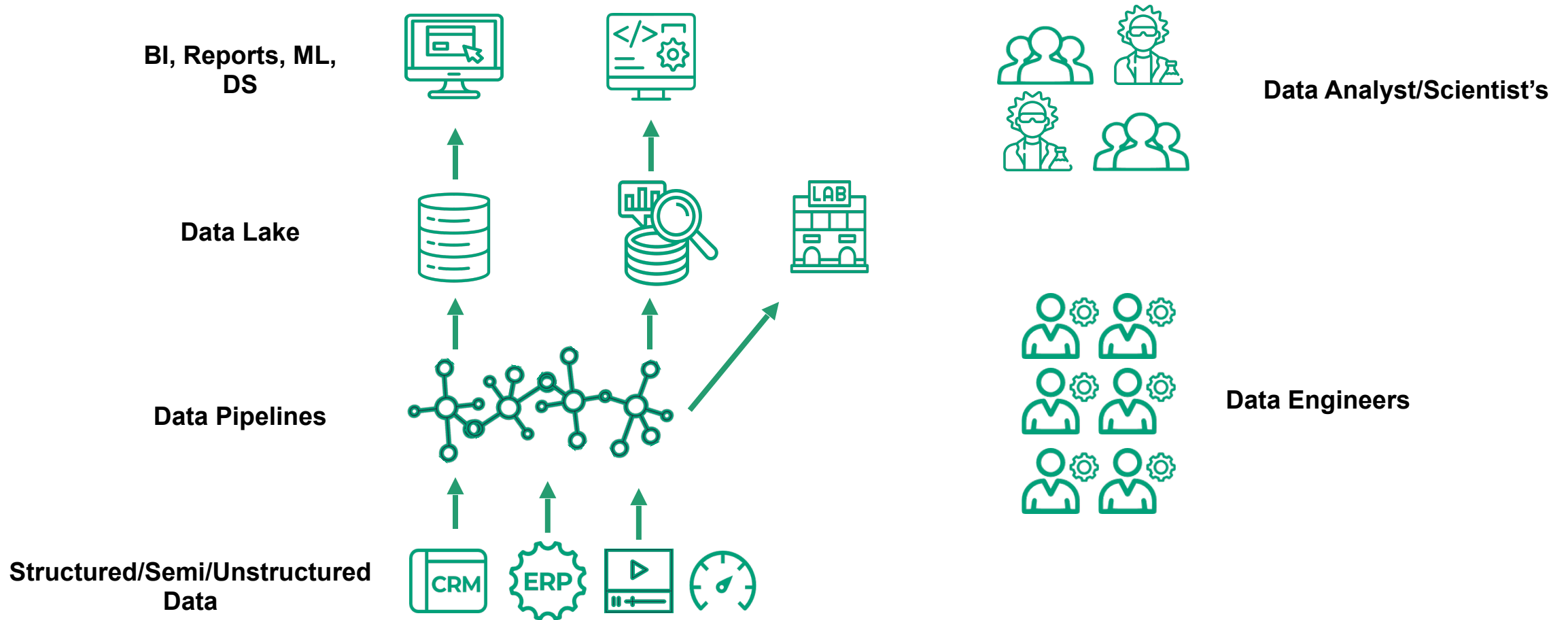
Текущие архитектурные подходы обработки данных



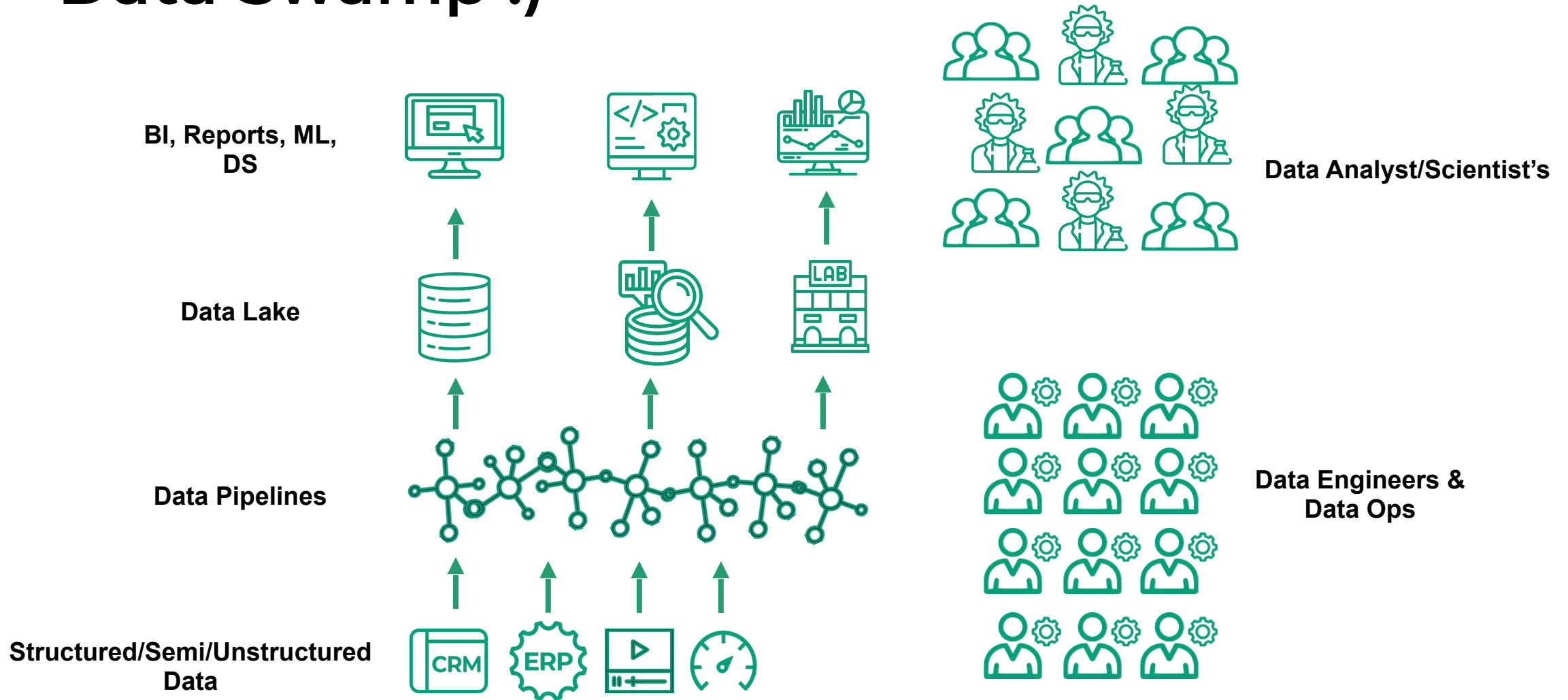
Data Warehouse



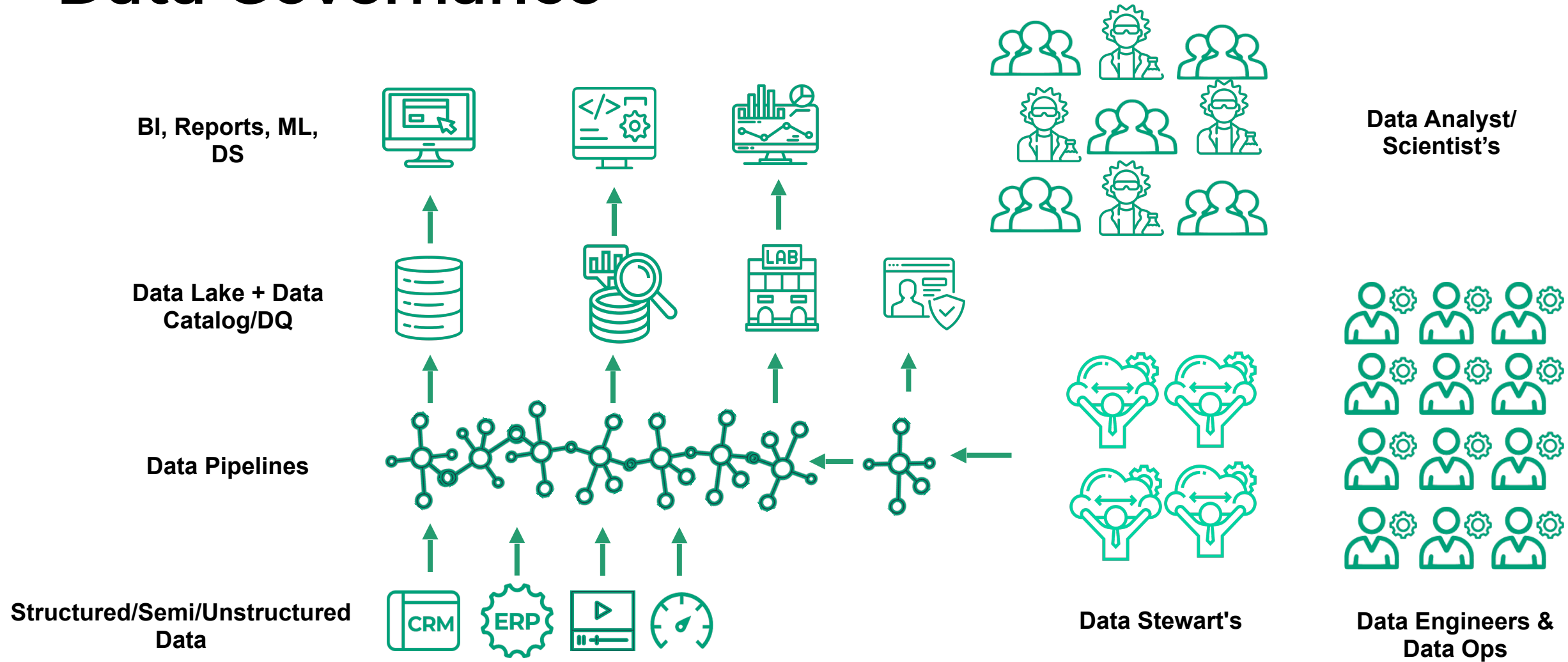
Data Lake



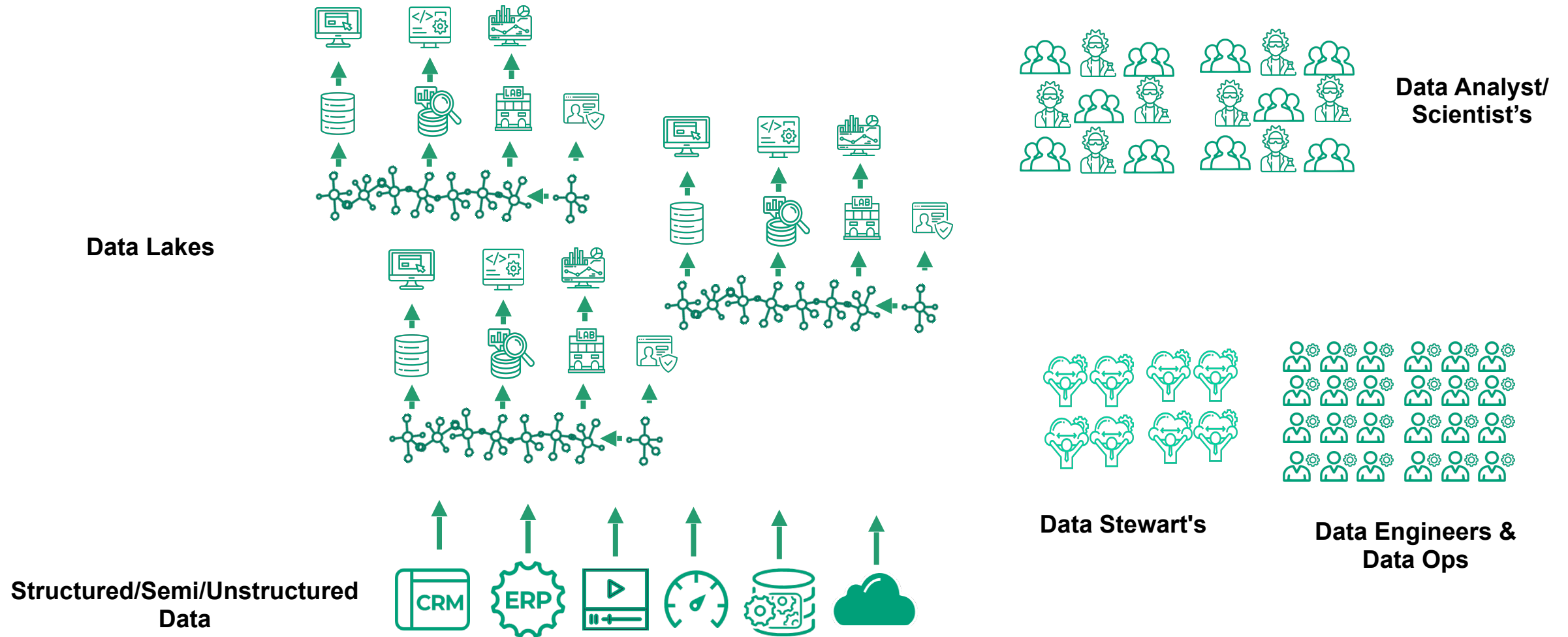
Data Swamp :)



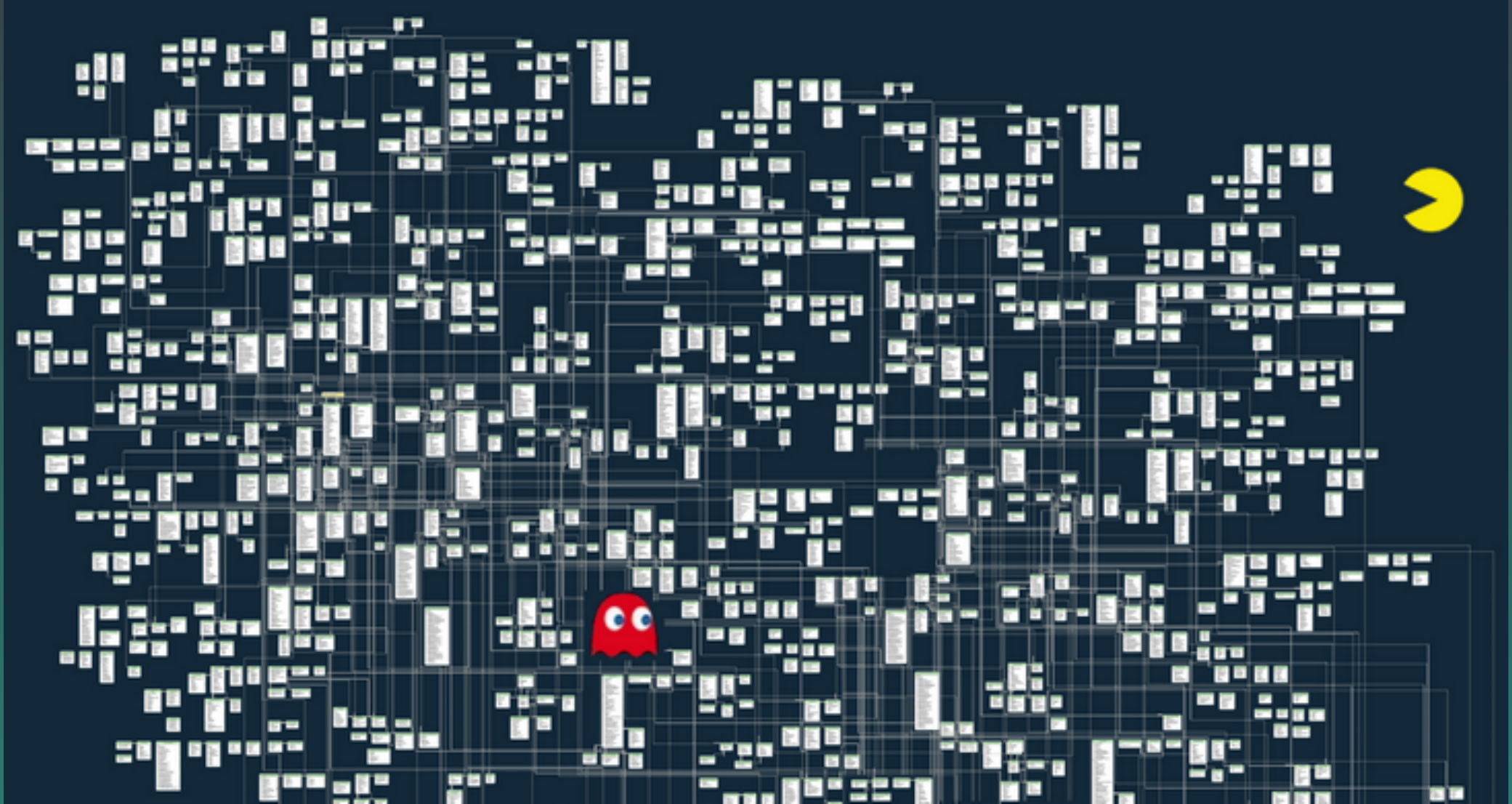
Data Governance



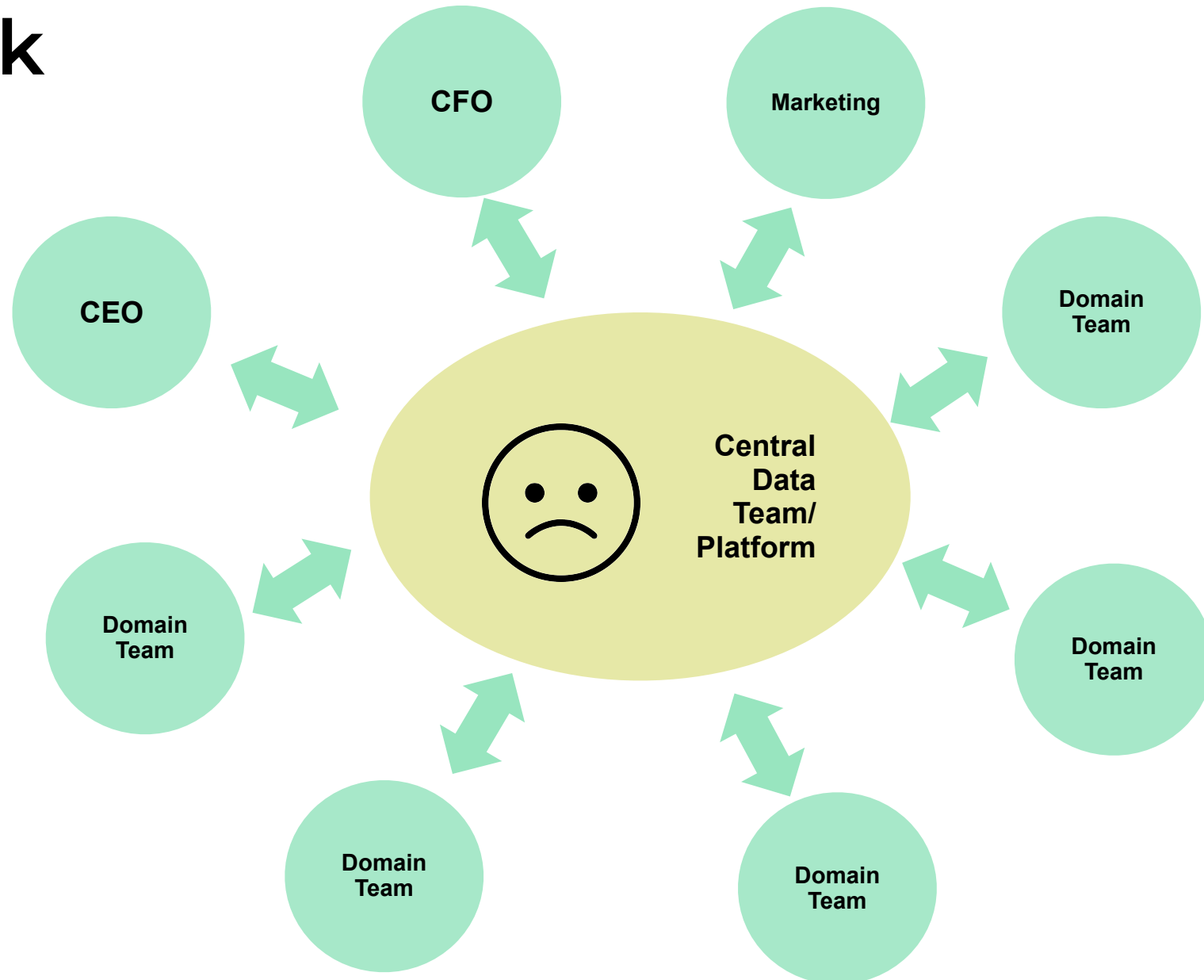
Data Oceans



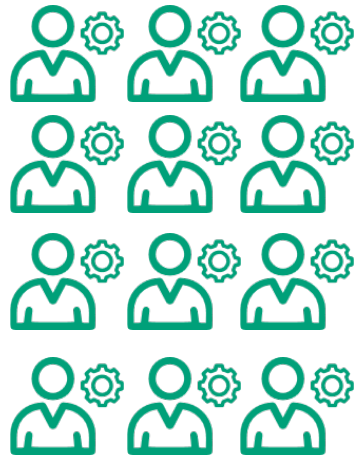
Модель данных



Bottleneck



Разрозненные и узко специализированные команды



Data & Ops Engineers



**Data Governance, Data Science &
ML Team**

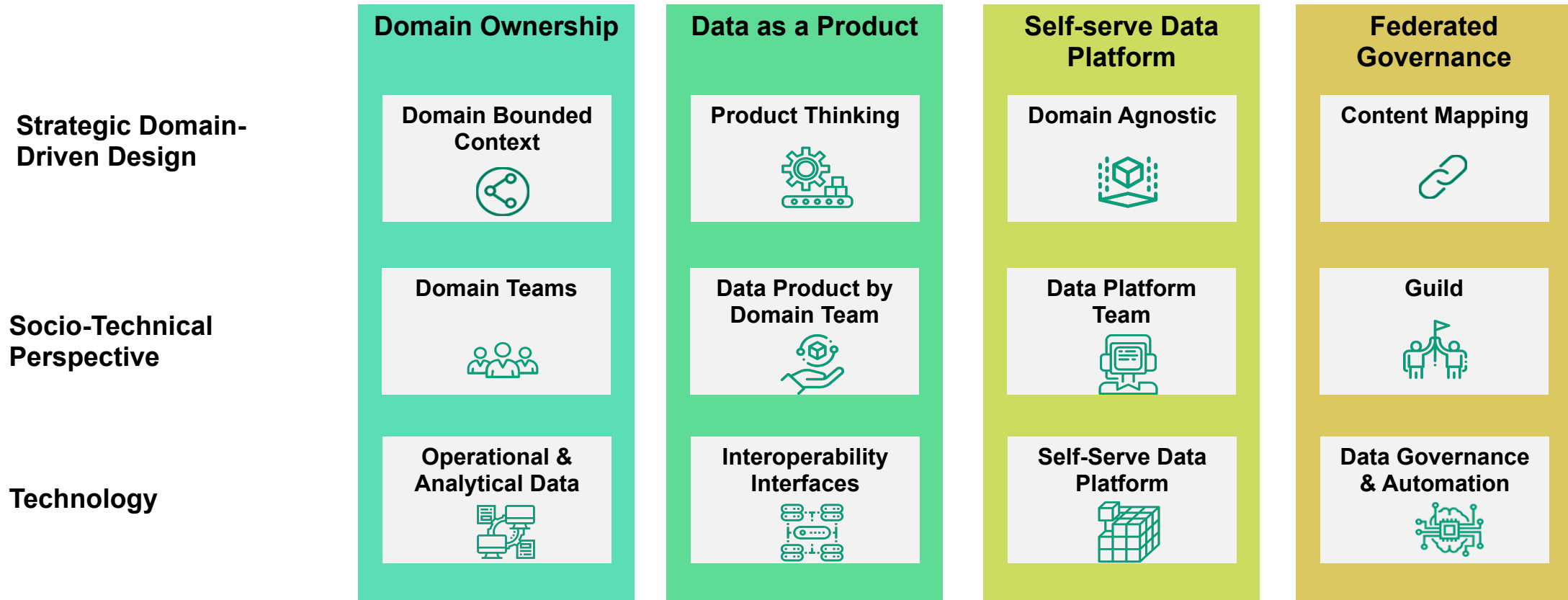


**Cross-Functional Domain
Oriented Customer/Consumer
Team**

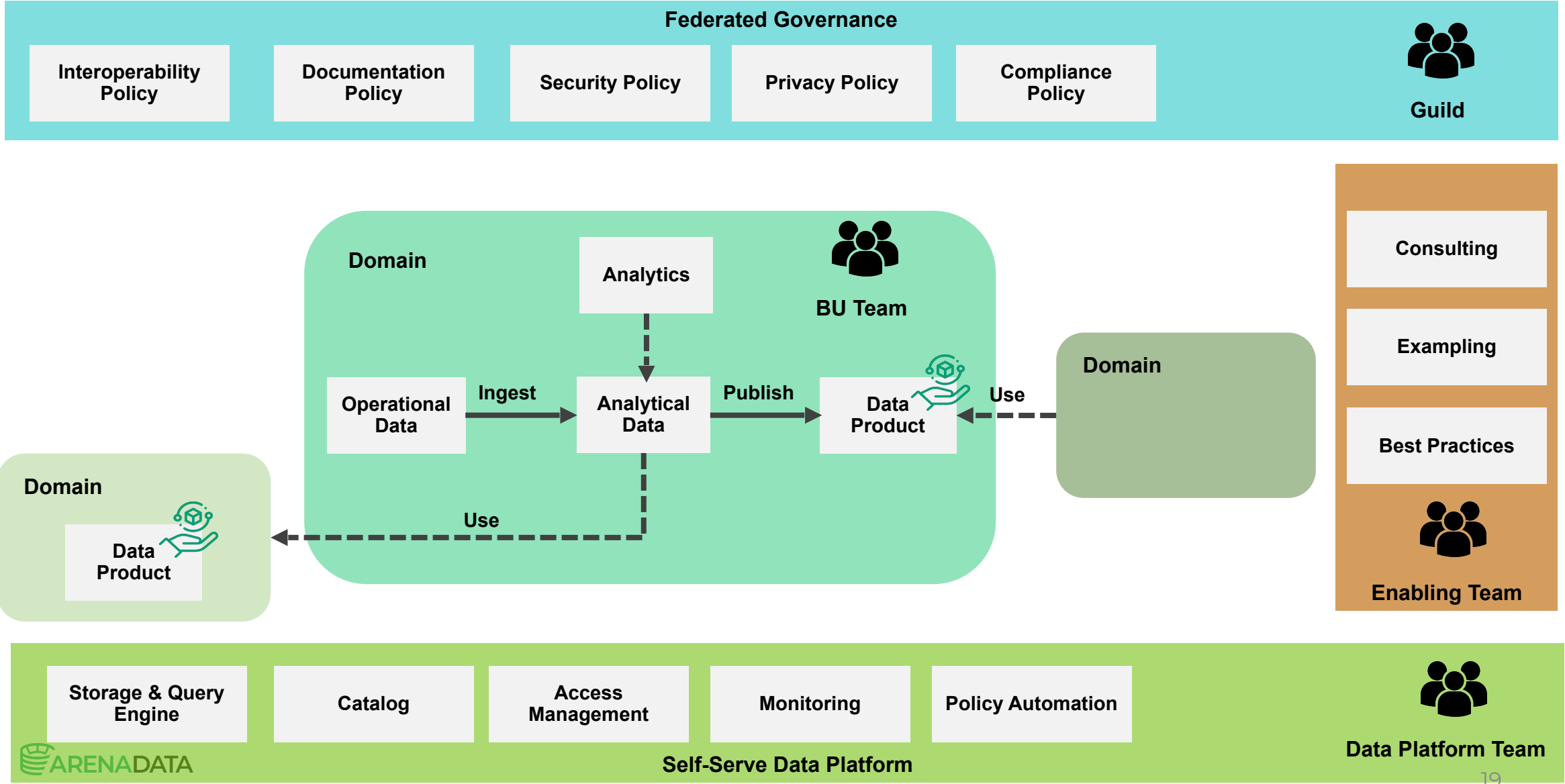
Архитектура следующего поколения платформы данных



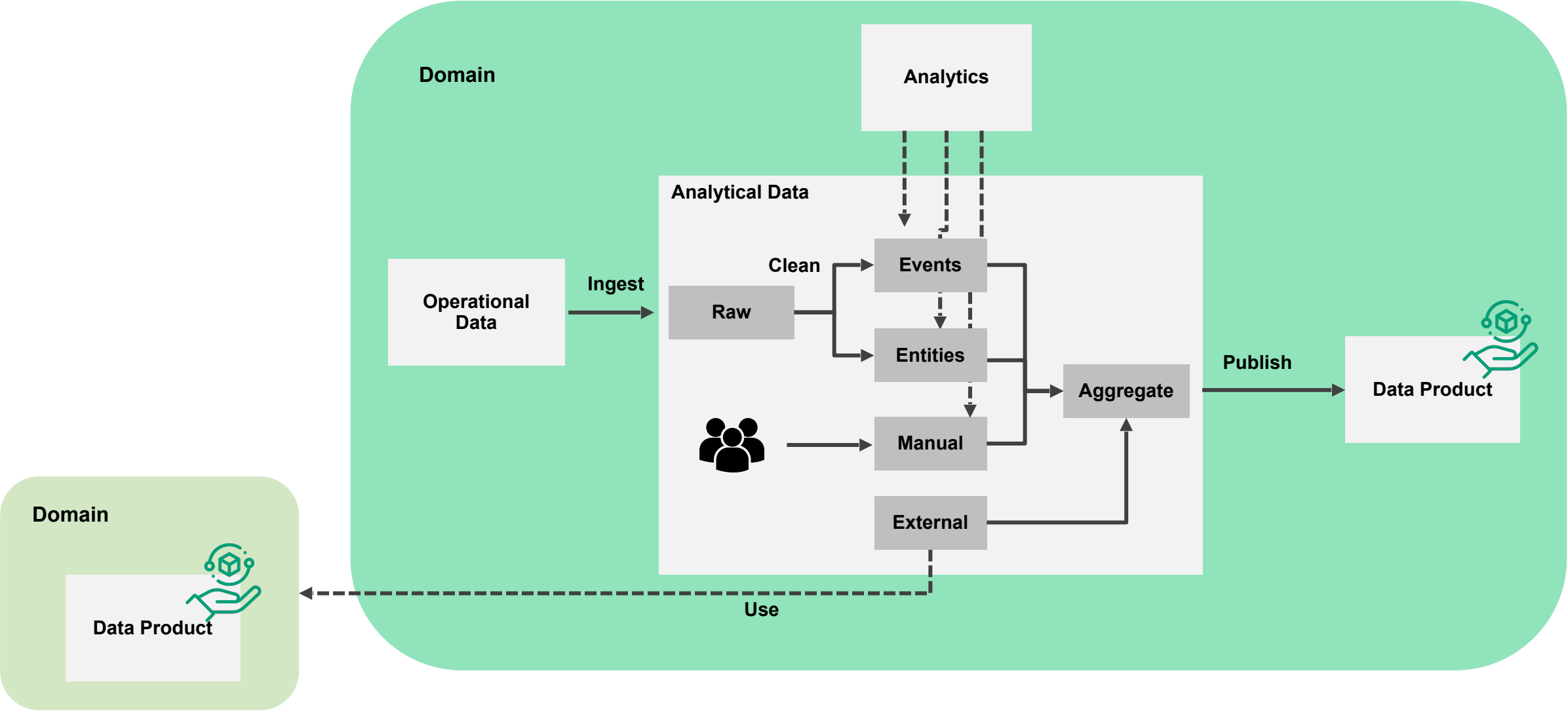
Data Mesh 101



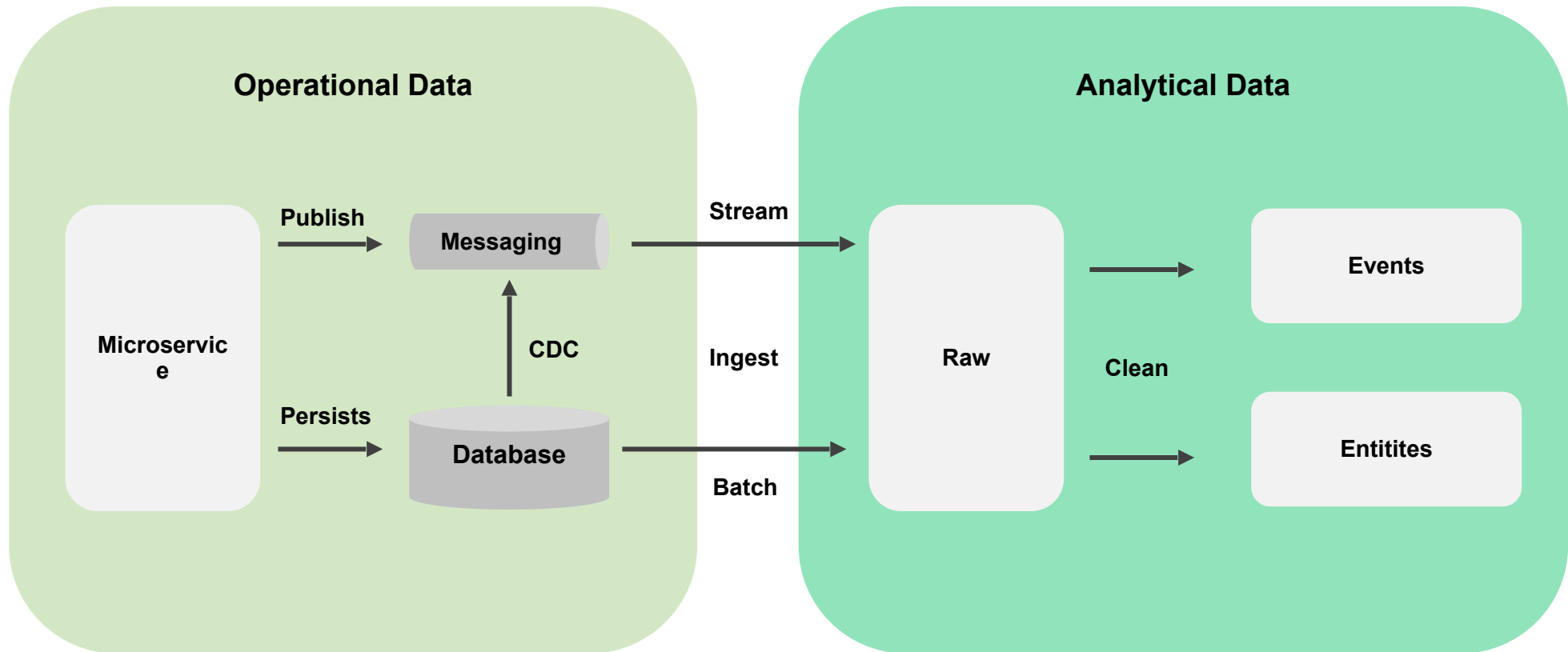
Data Mesh Concept Architecture



Analytical Data



Data Ingesting & Cleaning



Clean Data

- Структурирование
- Дедупликация
- Полнота
- Arteфакты

This feature helps you find groups of different cell values that might be alternative representations of the same thing. For example, the two strings "New York" and "new york" are very likely to refer to the same concept and just have capitalization differences, and "Gödel" and "Godel" probably refer to the same person. [Find out more...](#)

Method key collision Keying Function Fingerprint 7 clusters found

Cluster size	Row Count	Values in cluster	Merge?	New cell value
2	10	<ul style="list-style-type: none">PRINCESS FIELD (6 rows)PRINCESS FIELD (4 rows)	<input type="checkbox"/>	PRINCESS FIELD
2	22	<ul style="list-style-type: none">MOUNT ZION (18 rows)MOUNT ZION (4 rows)	<input type="checkbox"/>	MOUNT ZION
2	10	<ul style="list-style-type: none">SWABYS HOPE (5 rows)SWABYS HOPE (5 rows)	<input type="checkbox"/>	SWABYS HOPE
2	9	<ul style="list-style-type: none">BROGUE HILL (5 rows)BROGUE Hill (4 rows)	<input type="checkbox"/>	BROGUE HILL
2	9	<ul style="list-style-type: none">BALLARDS RIVER (7 rows)BALLARDS RIVER (2 rows)	<input type="checkbox"/>	BALLARDS RIVER
2	5	<ul style="list-style-type: none">SPRING MOUNTAIN (3 rows)MOUNTAIN SPRING (2 rows)	<input type="checkbox"/>	SPRING MOUNTAIN
2	19	<ul style="list-style-type: none">GRAVEL HILL (17 rows)GRAVEL HILL (2 rows)	<input type="checkbox"/>	GRAVEL HILL

Rows in cluster

5 — 22

Average length of choices

11 — 15

Length variance of choices

0 — 1

Select all Deselect all Export clusters Merge selected & re-cluster Merge selected & Close Close

Federated Governance

- Глобальная политика, которая определяет, как доменные команды должны создавать свои информационные продукты
- Политика интероперабельности
- Централизованный каталог метаданных
- Унифицированный способ безопасного доступа к фактическому продукту данных

Global Policy 1: Data Format

CSV-Files

Standard: RFC 4180

Delimiter: COMMA

Header-Record: Required

Global Policy 2: Location

S3

Bucket is managed by domain team

Global Policy 3: Data Discovery

A common table in Wiki that lists:

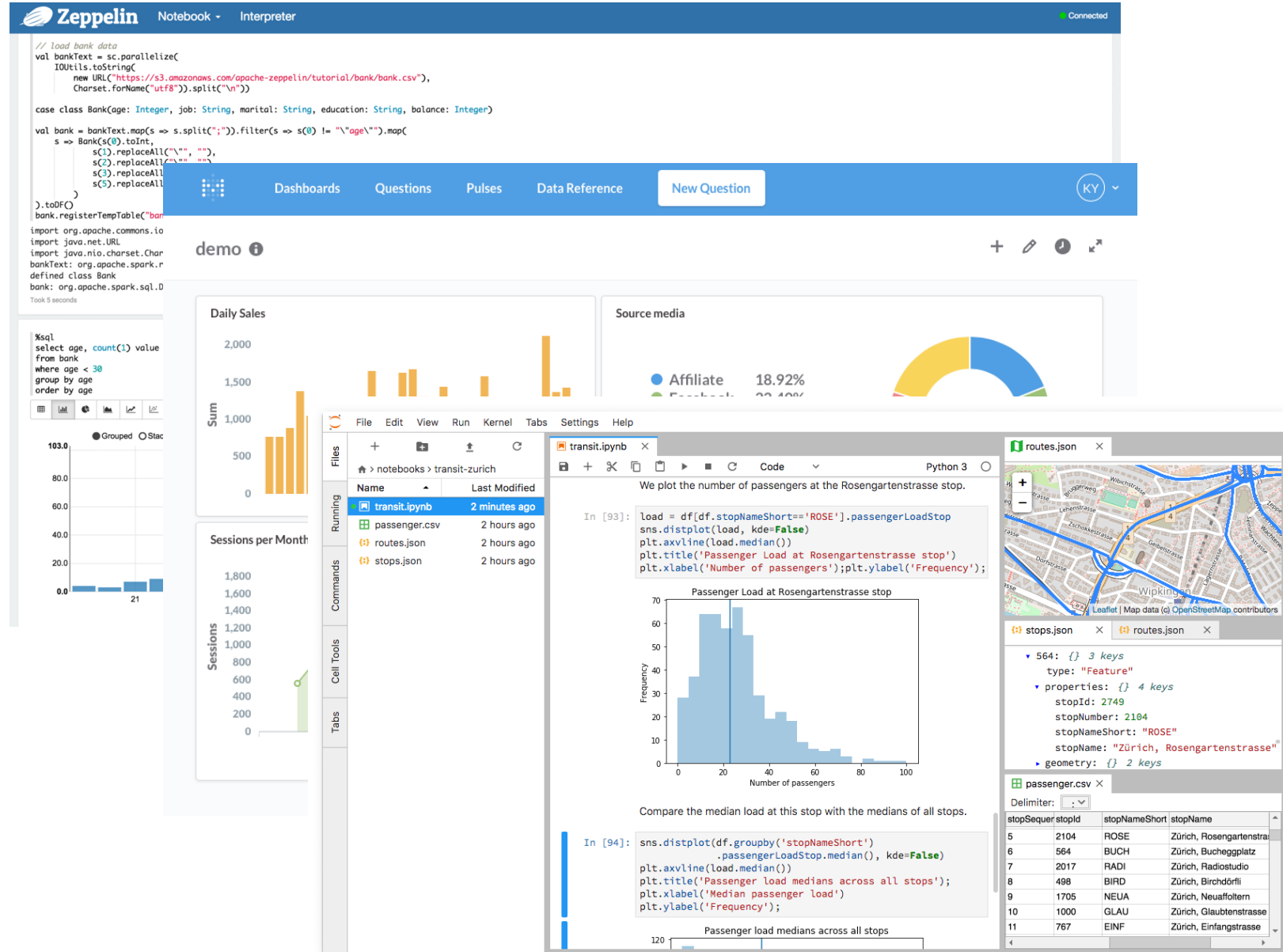
- Domain Team
- Data Product Name
- Responsible Product Owner
- S3-Location
- Data Freshness / Update frequency
- Tags
- Data Model

Global Policy 4: Access Control

AWS IAM Role based ACL

Analytics

- Поддержка интеграции с различными источниками
- Поддержка синтаксиса SQL
- Поддержка in-database и filter pushdown функционала



Data Product

- Продукт данных обычно представляет собой опубликованный набор данных, к которому могут получить доступ другие домены, аналогично API.
- Продукт данных описывается метаданными, включая информацию о владельце и контактную информацию, местоположение данных и доступ к ним, частоту обновления и спецификацию модели данных.
- Команда домена отвечает за работу продукта данных в течение всего его жизненного цикла.

Data

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

Metadata

Tags (1)

Data Product

Display name	Value
Domain Team	Fulfillment
Update Frequency	Hourly
Environment	PRODUCTION

Schema

Filter

Enter property name or value

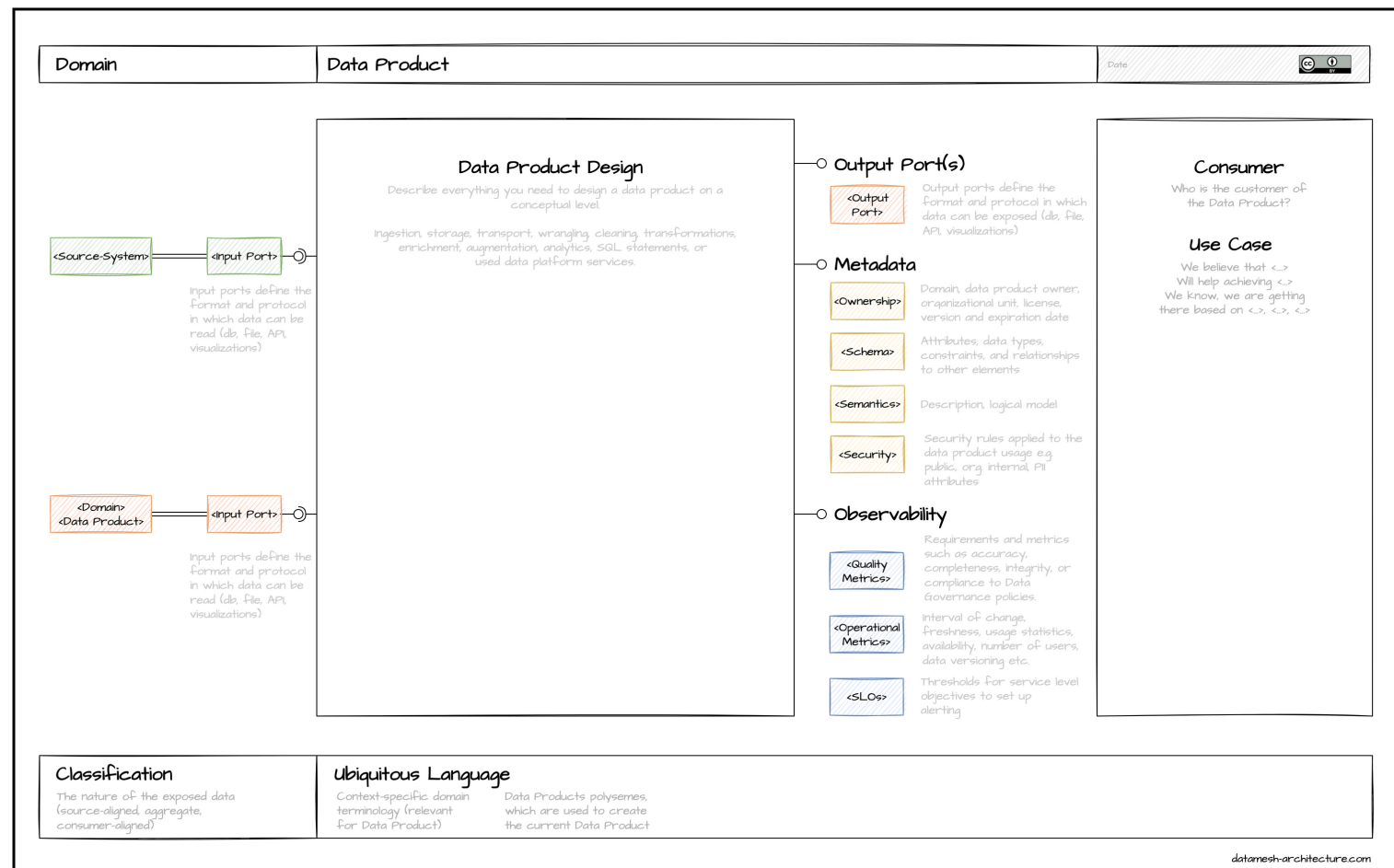
Field name	Type	Mode
sku	STRING	NULLABLE
location	STRING	NULLABLE
available	INT64	NULLABLE
updated_at	TIMESTAMP	NULLABLE

Operations

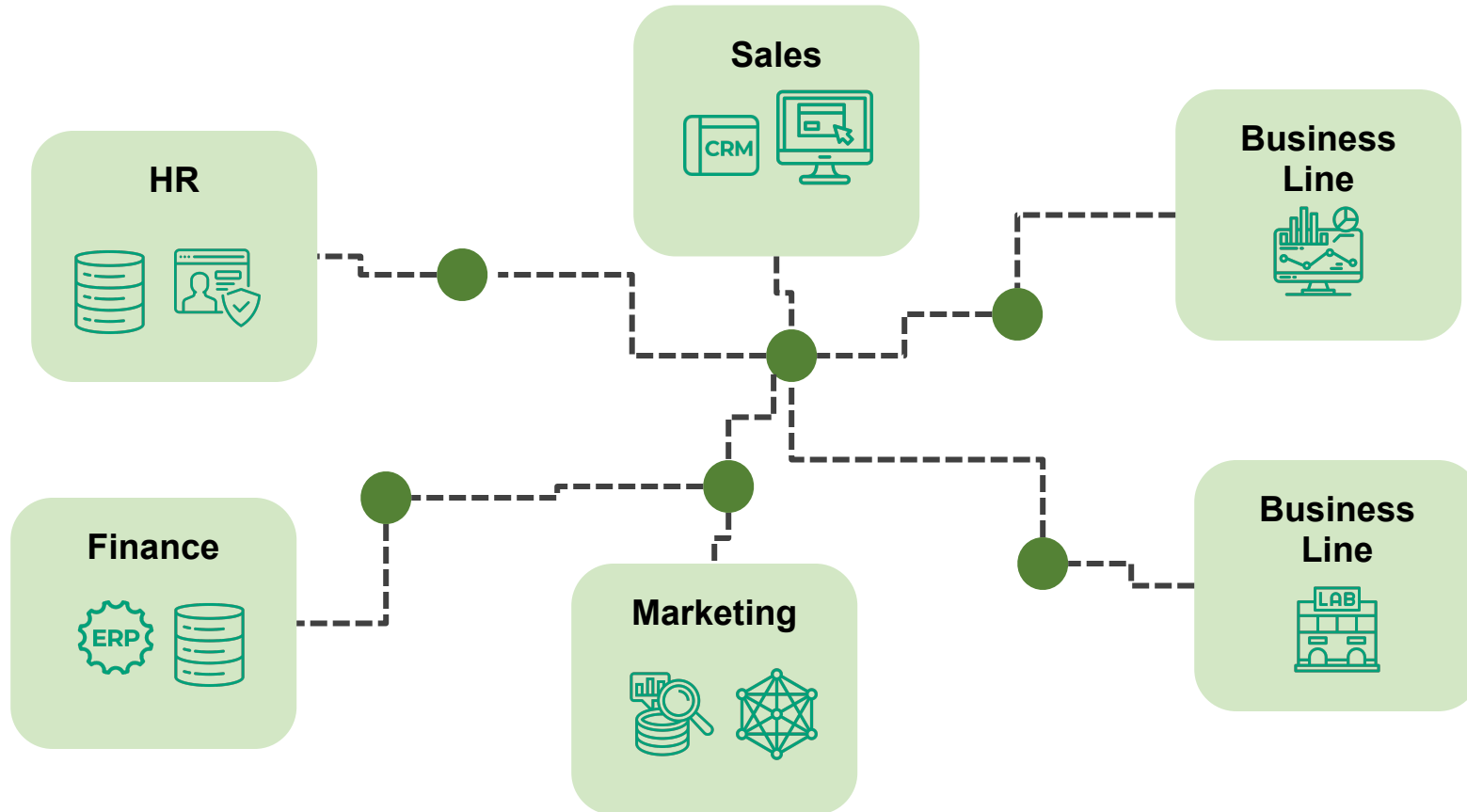
132M	0	5 min ago	42	\$1.4k
Rows	Duplicates	Freshness	24h Usages	Monthly Costs

Data Product Canvas

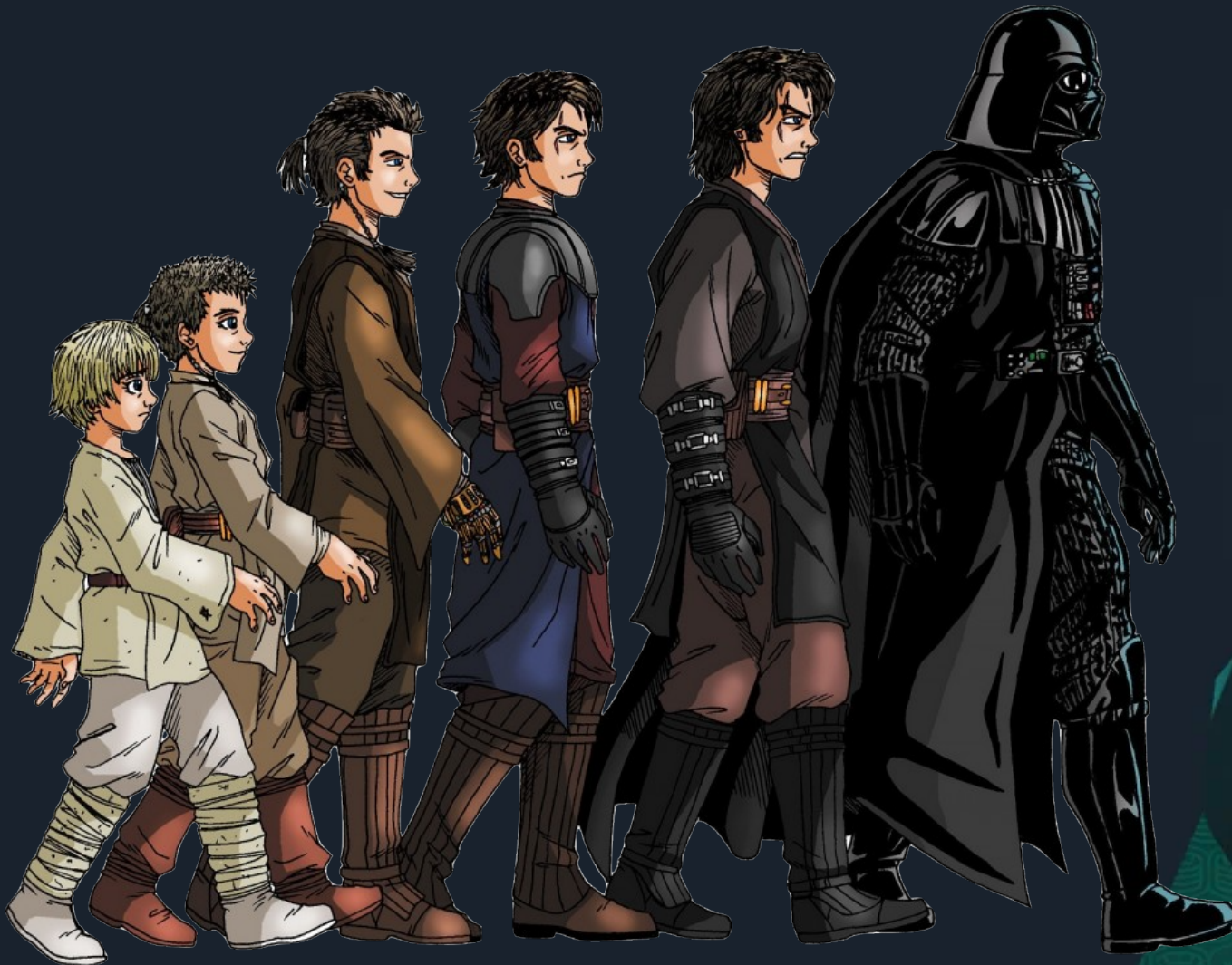
- Домен
- Название продукта
- Потребитель и пример использования
- Выходной порт
- Метаданные
- Входные порты
- Описание продукта данных
- Состояние продукта данных
 - Показатели качества
 - Операционные показатели
 - SLO
- Синтаксис взаимодействия
- Классификация продукта




Domain Driven Design

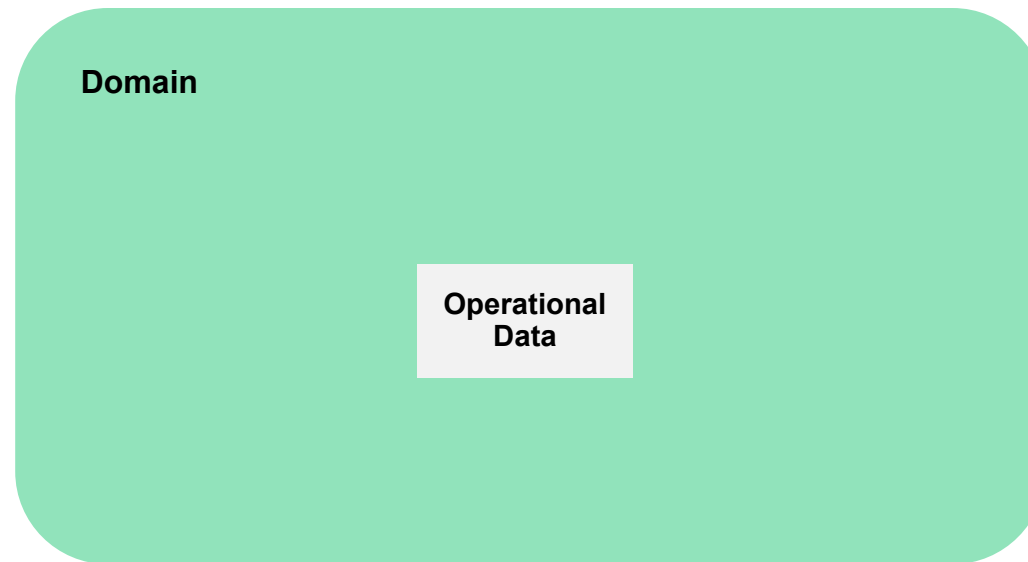


Enabling Team

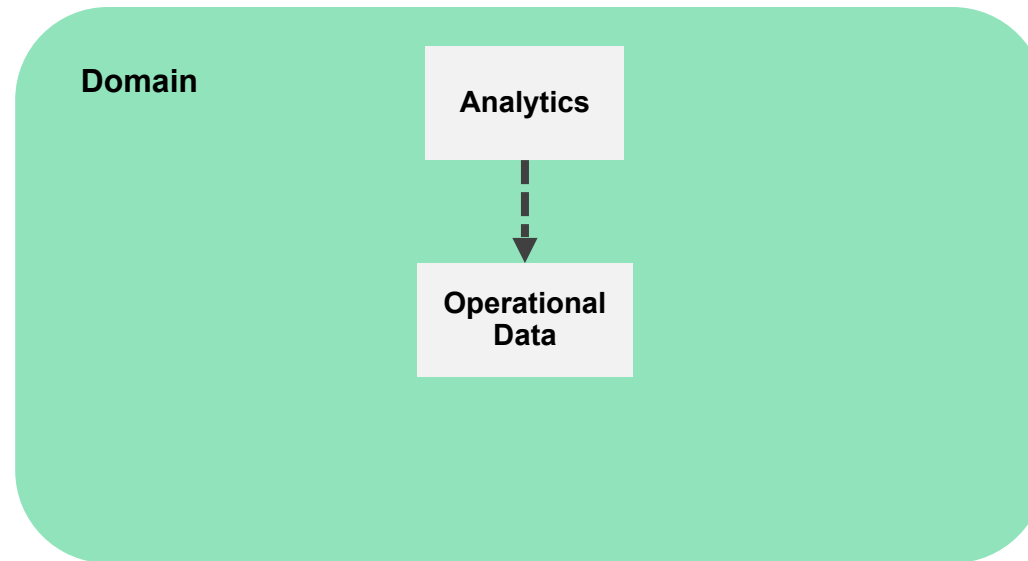


 ARENADATA

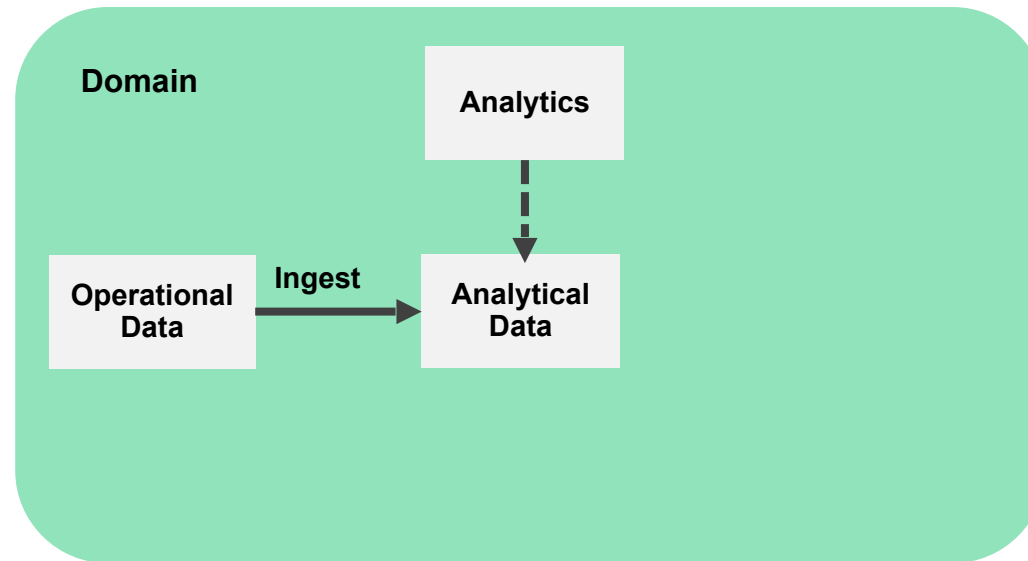
No Data Analytics (lvl 0)



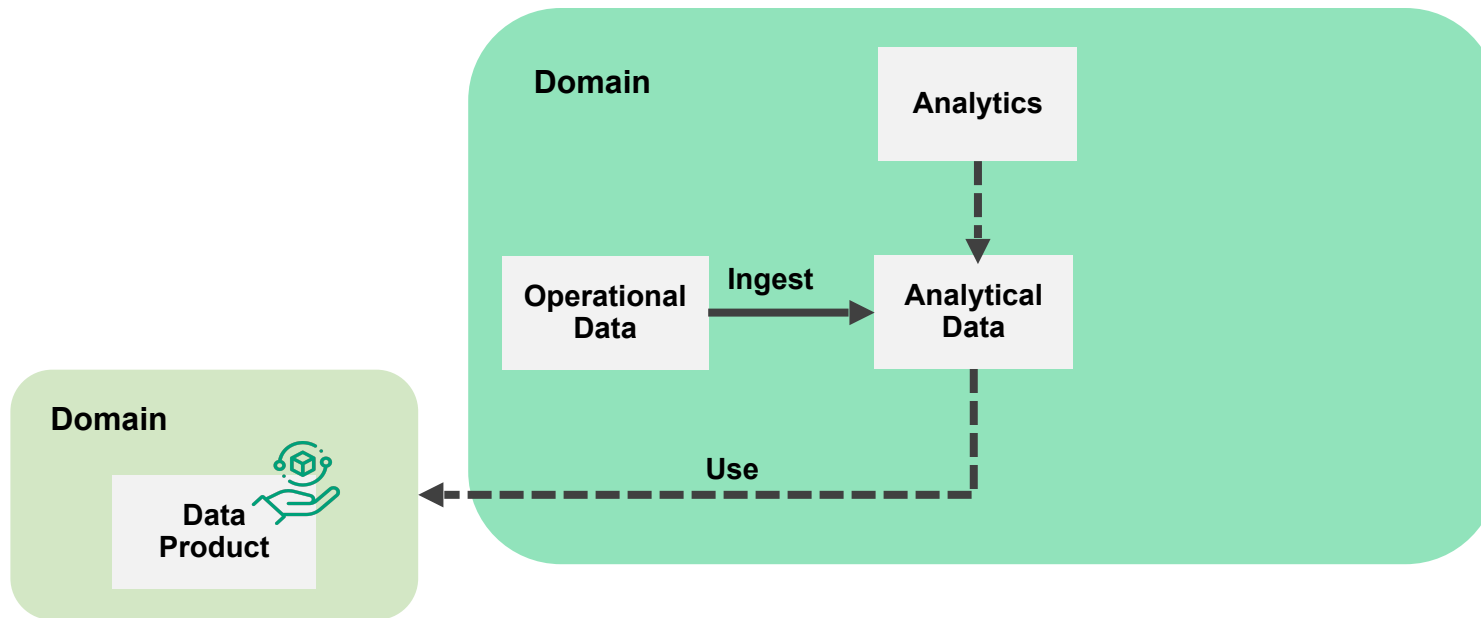
Operational Database Queries (lvl 1)



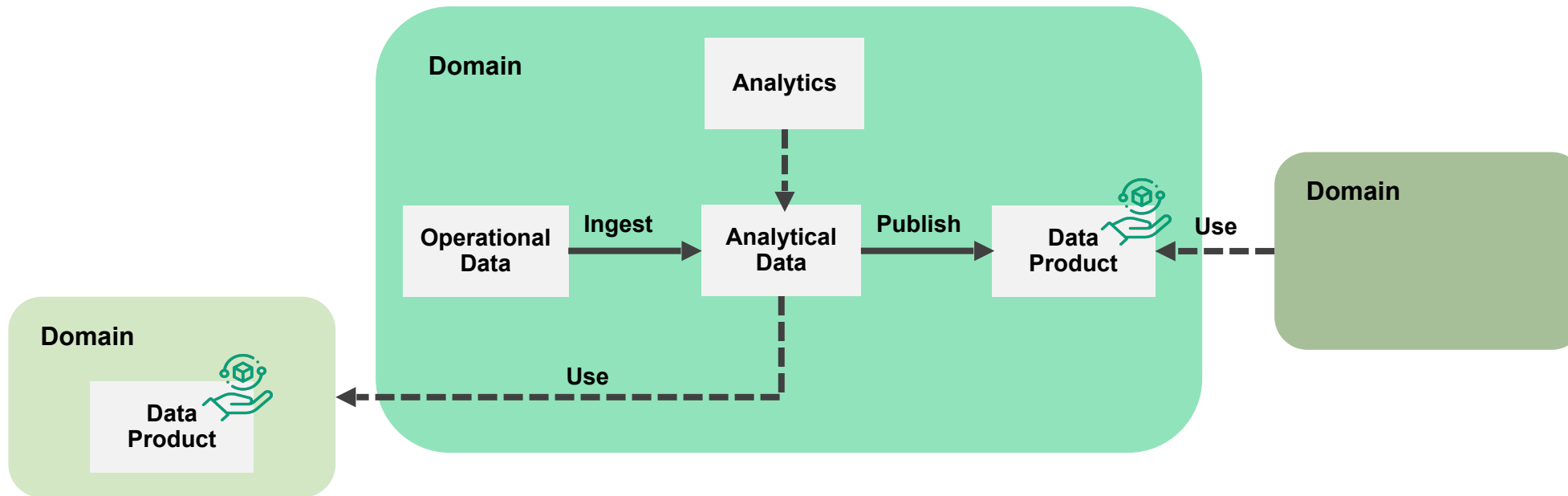
Analyze Own Data (lvl 2)



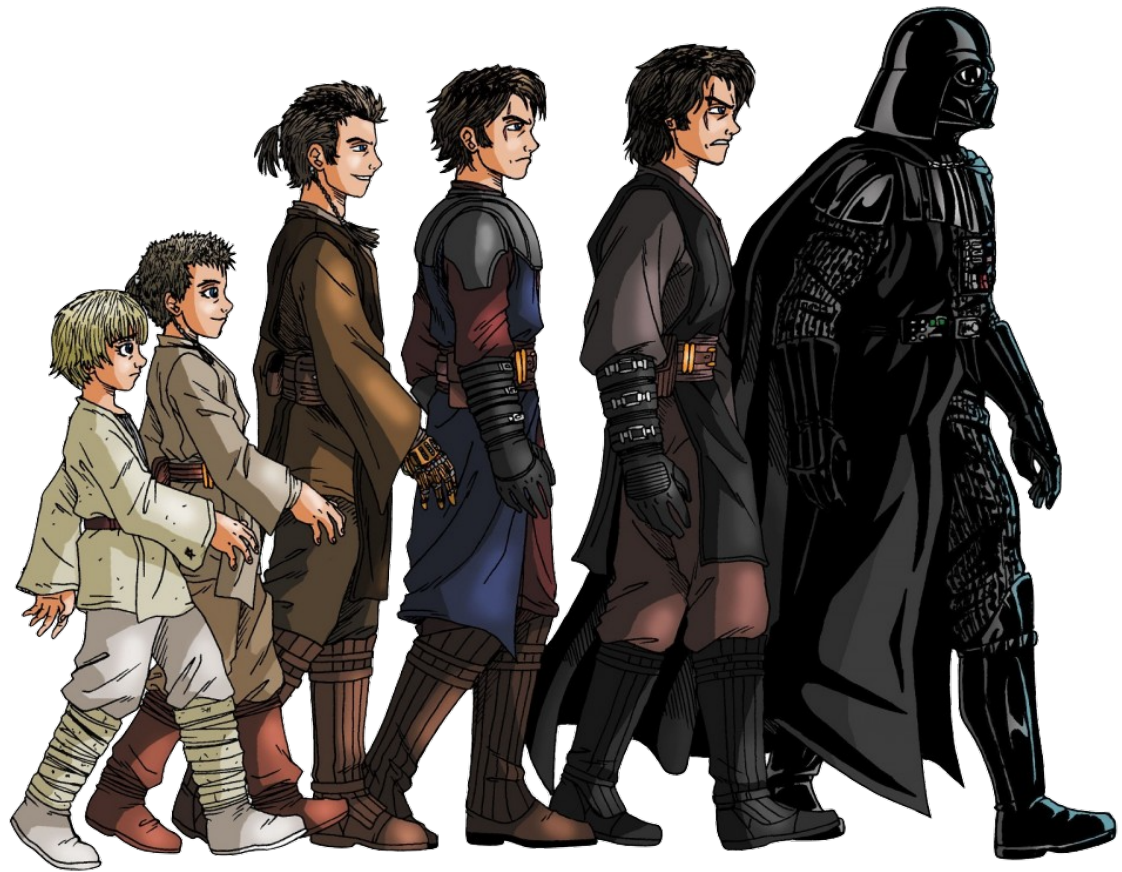
Analyze Cross Domain Data (lvl 3)



Publish Data as a Product (lvl 4)



Data Team Journey



Level 4

Publish Data as a Product

Level 3

Analyze Cross Domain Data

Level 2

Analyze Own Data

Level 1

Operational Database Queries

Level 0

No Data Analytics

Реализация платформы данных для Data Mesh



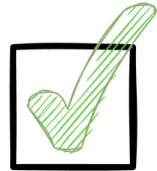
А надо ли?



Decentralize business
into domains



Decentralize engineering
into autonomous teams

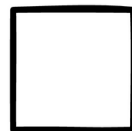


Decentralize monolith
into microservices



Decentralize operations
into DevOps-driven teams

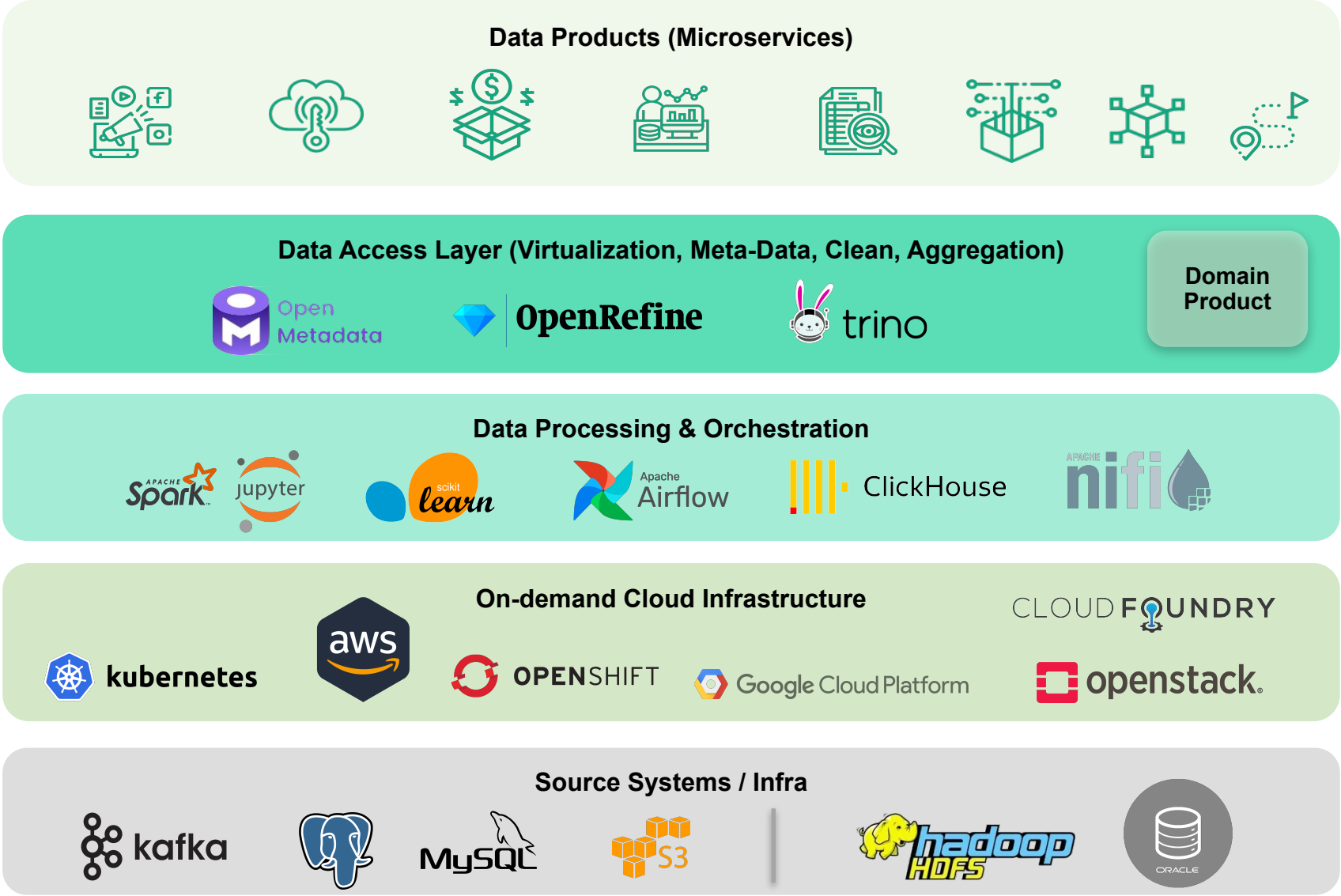
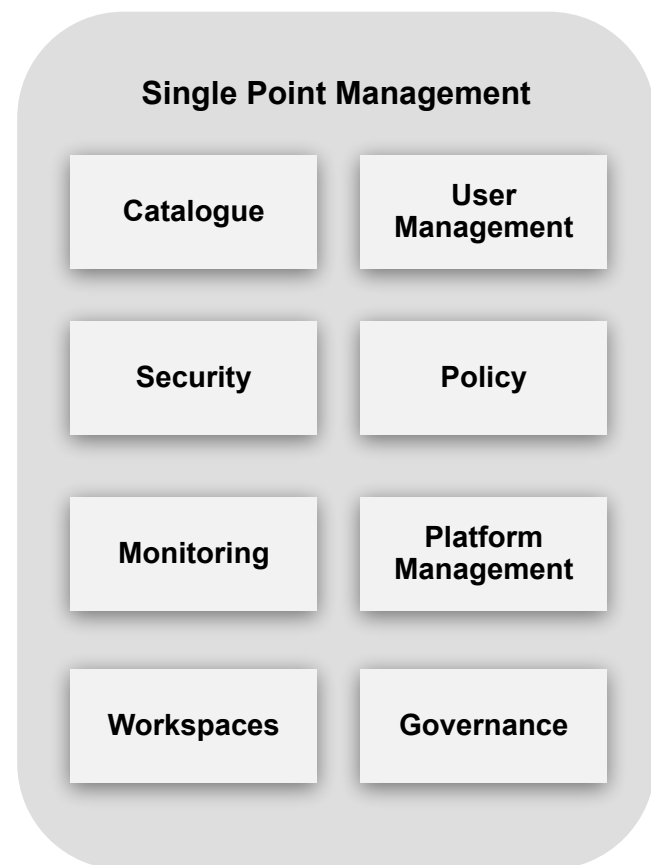
Scale up
software
development



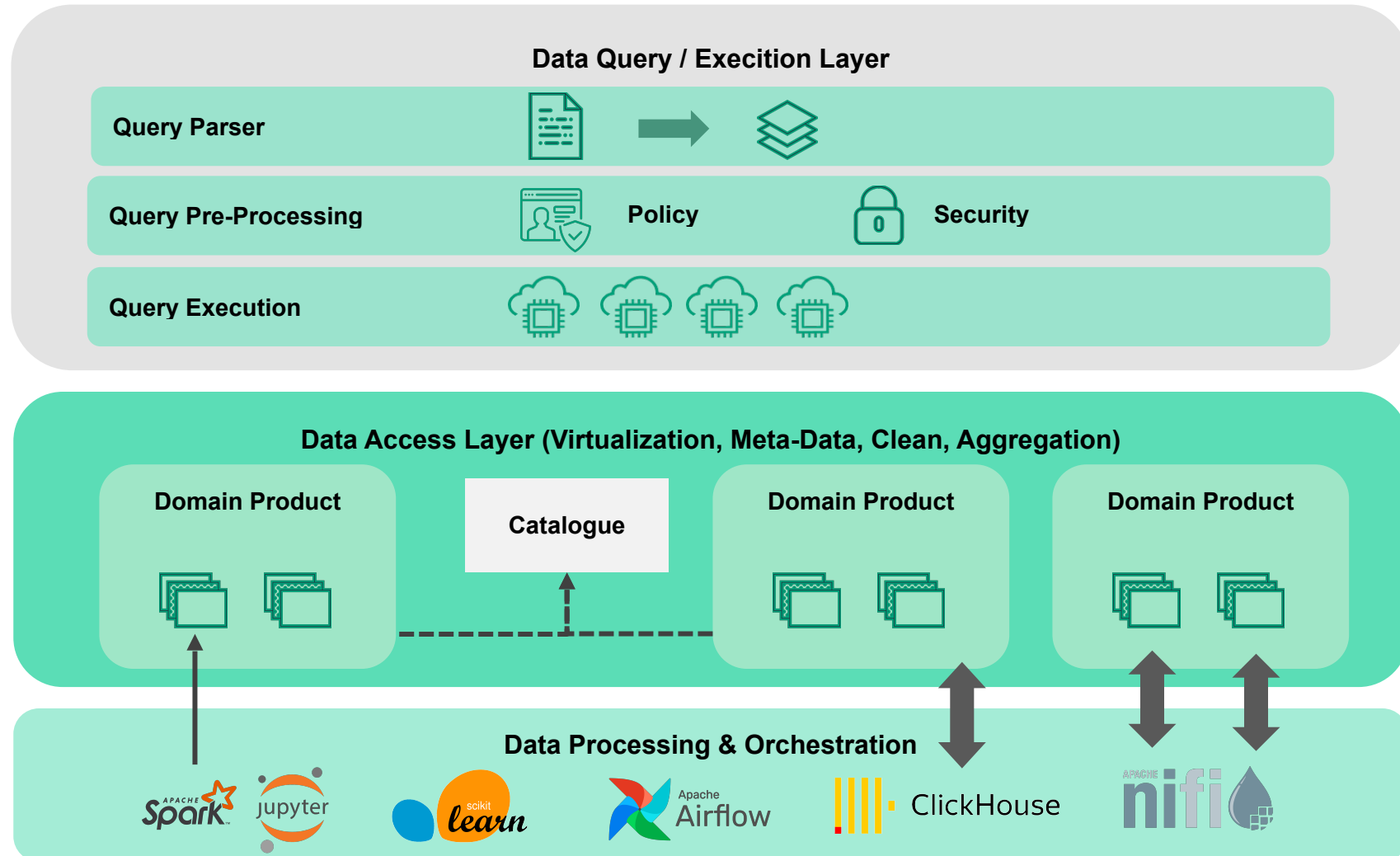
Decentralize data lake
into data mesh

Scale up
data analytics

Data Platform

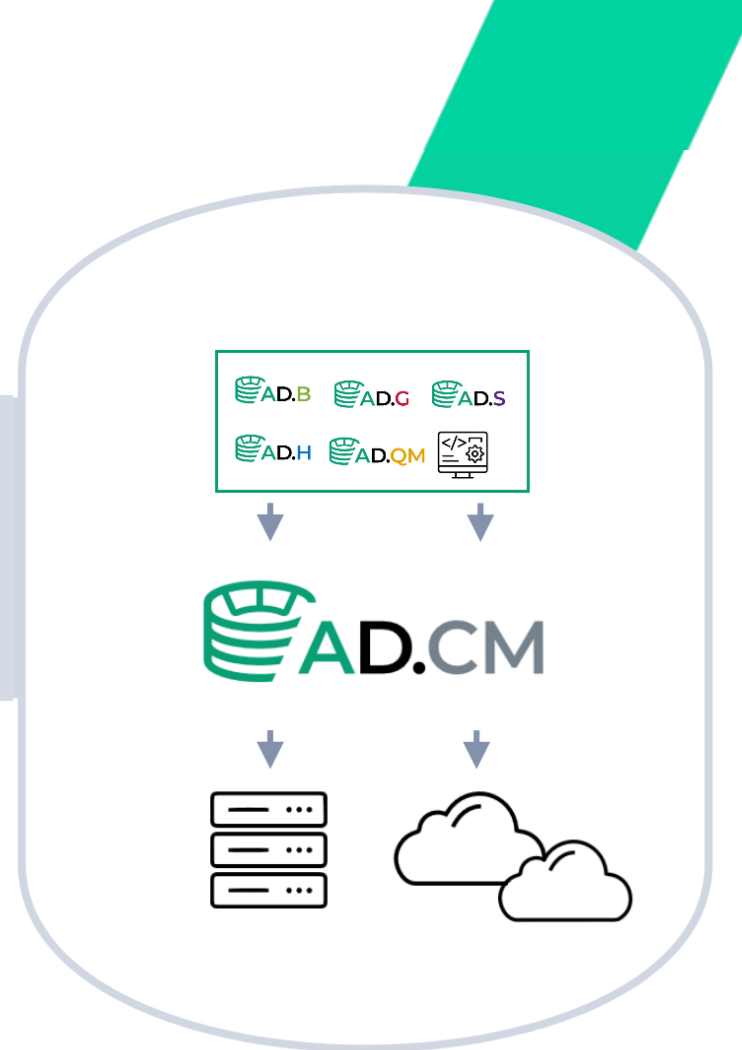


Data Virtualization

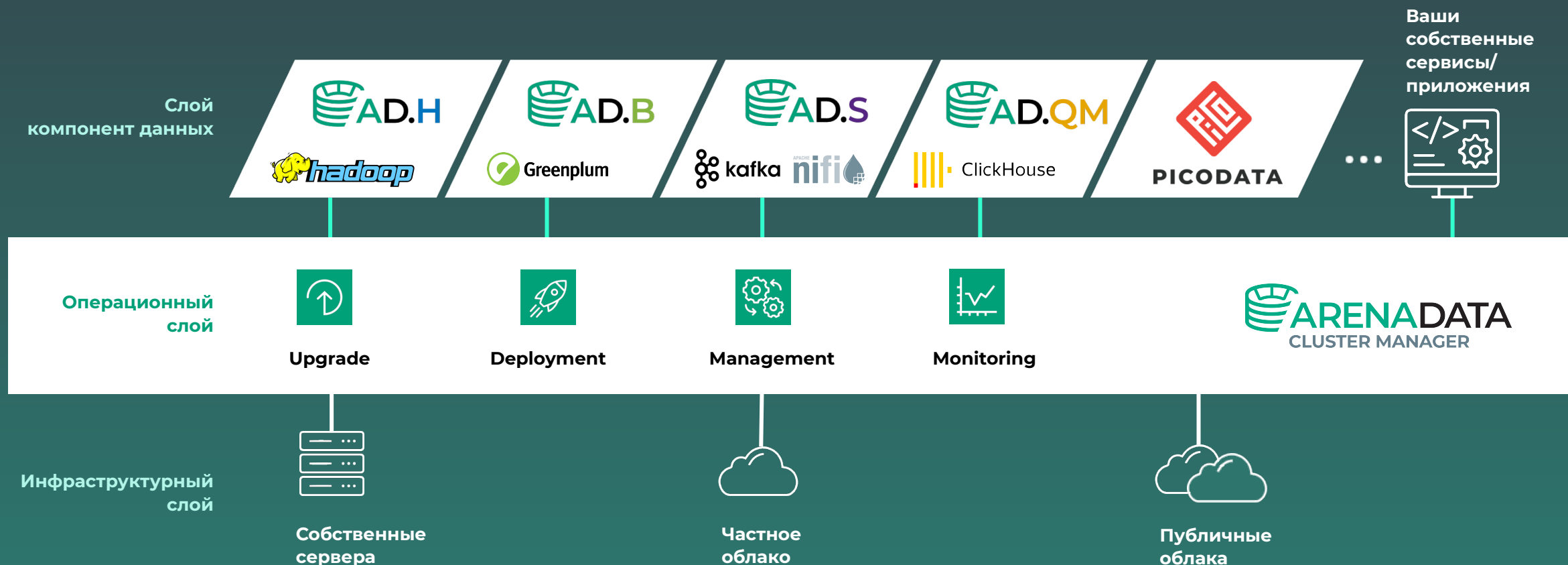


Arenadata Cluster Manager (ADCM) - универсальный оркестратор гибридного ландшафта. Позволяет быстро устанавливать, настраивать все data-сервисы компании и управлять ими независимо от инфраструктуры.

- В Arenadata Cluster Manager установка, настройка и обновление кластеров в оркестраторе производятся по нажатию кнопки в графическом интерфейсе или по запросу в API.
- При этом все настройки ОС, сервисов, сети, монтирование дисков происходят автоматически.
- Наиболее ярко преимущества ADCM раскрываются при работе с гетерогенной инфраструктурой, при которой появляется возможность размещать data-сервисы на различных типах инфраструктур: в облаке, on-premise или в качестве PaaS-сервисов.



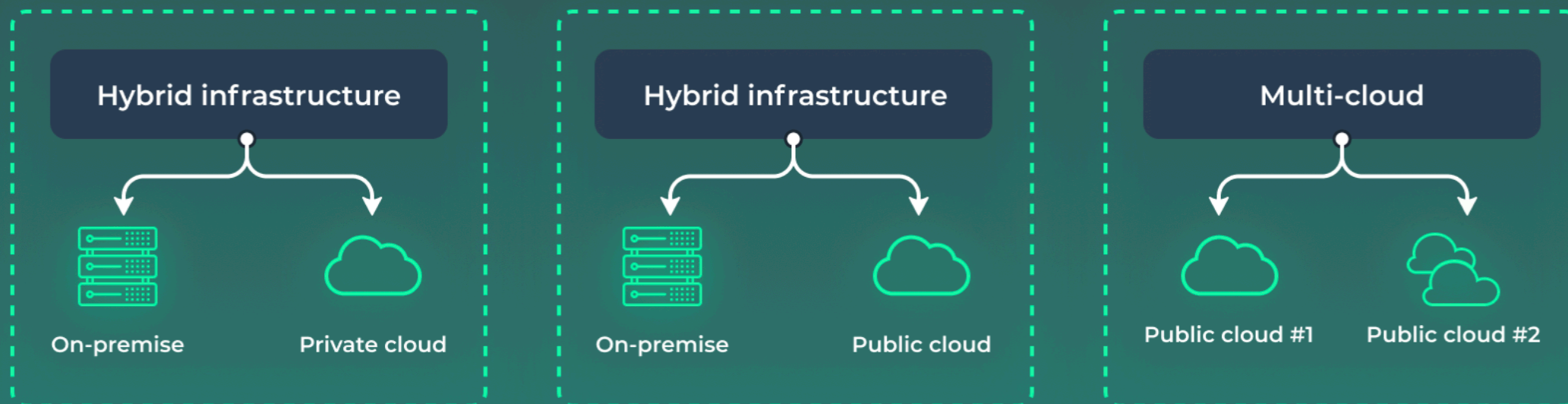
Платформа данных Arenadata



Объединение разных облачных провайдеров

Мы заложили в Arenadata Cluster Manager возможность поддерживать любые ИТ-инфраструктуры и виртуальные машины в самых разных облаках и на bare-metal оборудовании. В настоящий момент у нас есть коннекторы для Google, Яндекс.Облако, Mail.ru Cloud Solutions, VMware. При необходимости мы можем добавить коннекторы и к другим облакам.

Благодаря ADCM установка наших продуктов в гибридной ИТ-архитектуре происходит в несколько кликов из интерфейса оркестратора.





<https://github.com/arenadata/adcm>



`docker pull hub.arenadata.io/adcm/adcm:latest`



https://t.me/arenadata_cm



https://store.arenadata.io/#products/arenadata_cluster_manager





Корпоративная платформа хранения
и обработки больших данных

Оцените потенциал Big Data и Open Source вместе с Arenadata Enterprise Data Platform

- Узнайте больше на arenadata.tech
- Скачайте бесплатно на store.arenadata.io
- Получите консультацию по почте info@arenadata.io

