



# Die OTTO BI

Wie man hunderte von Data  
Scientist\*innen mit Daten versorgt

# Anissa Ayoub

Product Owner @OTTO BI

Alumni FH Wedel Master E-Commerce (Abschluss  
2021)

<https://www.linkedin.com/in/Anissa-ayoub>

[anissa.ayoub@otto.de](mailto:anissa.ayoub@otto.de)

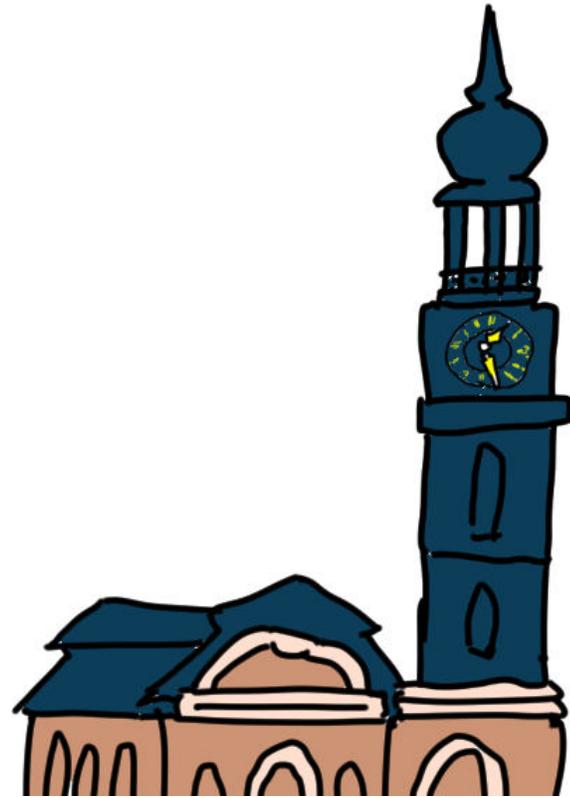
Bjørn Stachmann

<https://kapitel26.github.io>

@old\_stachi

BjoernArno.Stachmann@otto.de

**OTTO**



# business@OTTO

- 5,1 Milliarden € Umsatz
- 10 Millionen Produkte
- 19.000 Marken
- bis zu 10 Bestellungen pro Sekunde
- **otto.de**: 2,9 Millionen Visits pro Tag



# tech@OTTO

- 1000+ Techies
- 100+ Mio. Tech-Investitionen

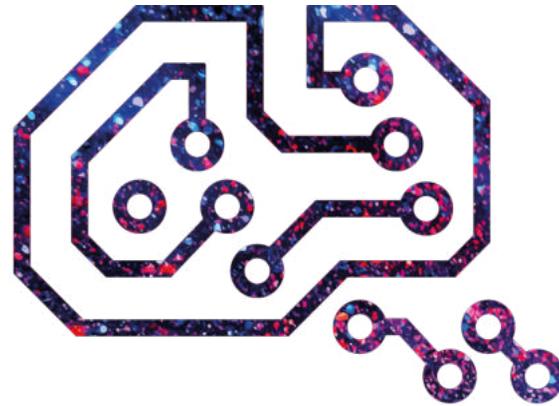
Und wir kommen aus der...



**OTTO BI**

INNOVATIONS AND EXCELLENT DECISIONS

Wir bauen...



**BRAIN**

**OTTO BI PLATTFORM**

Wir sind größer als einige Städte  
Deutschlands

**OTTO BI**

INNOVATIONS AND EXCELLENT DECISIONS

# Orga@BI

- 15 Abteilungen
- 400+ Kolleg\*innen
- 50+ Teams, agil organisiert

# Team MAMMUT

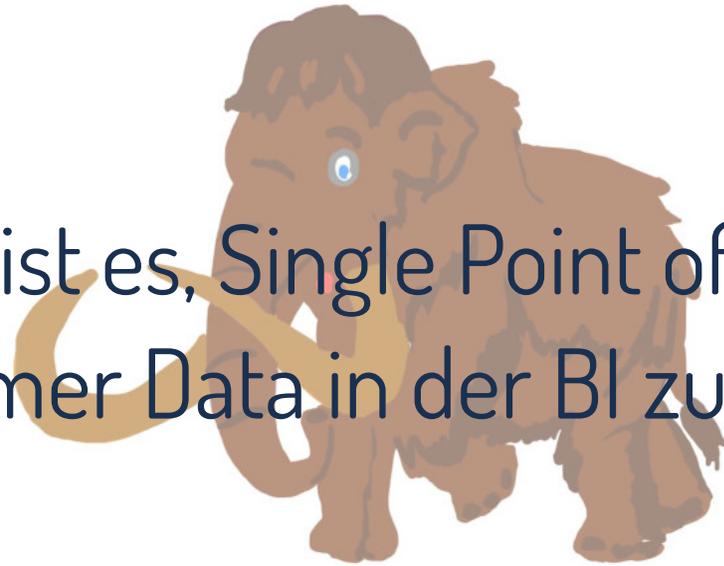


4 Data Eng.

1 Business Analyst

1 Product Owner

1 Agile Master



Unsere Vision ist es, Single Point of Truth für Core-  
Customer Data in der BI zu werden

REAL TIME  
ANALYTICS  
BIG DATA  
CLASSIC

REPORTING

Es fing alles ganz  
konventionell an.

REAL TIME

ANALYTIC

BIG DATA

CLASSIC

REPORTING

Aber die Zeiten begannen sich zu ändern.





Sie sind hungrig!

Sehr hungrig!

Sie sind

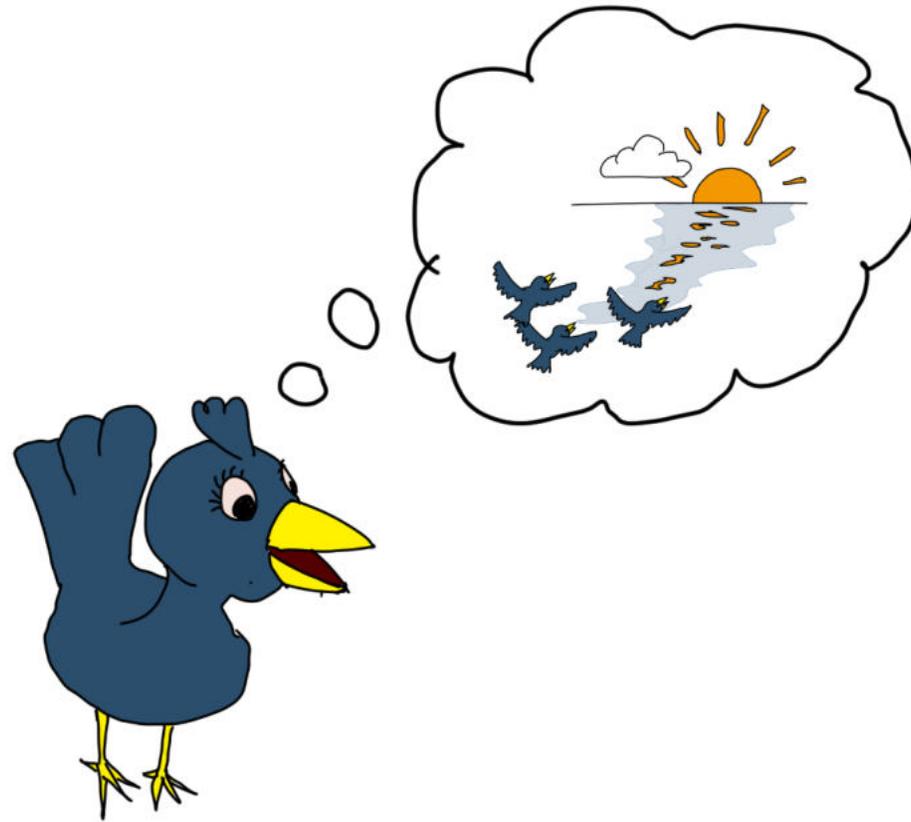
Data Scientisten!

Wer soll helfen?





Mama BI



- User experience
- Plattform-Partner-Management
- Ausbau der Sortimentsbreite
- Dynamic Pricing
- Logistik





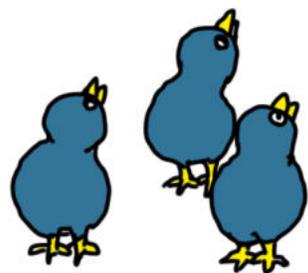
Deshalb wollen wir eine  
**Data Driven Company**  
werden.

Unsere Vision

Keiner holt aus Daten & KI mehr raus als OTTO

Daten & KI überall

Aber, so einfach ist das nicht:



Data Scientisten brauchen  
**mehr Zeit**  
für das **Abholen und Aufbereiten**  
von Daten  
als für Analysen.



Legacy DB

EXASOL

Kafka

adloop

ftp-Server



Und wie löst man das?



wir bauen ein

Big

Warehouse

Aber was genau ist ein  
Data Warehouse?

Moment ... mal eben googeln ...



**WIKIPEDIA**  
Die freie Enzyklopädie

[Hauptseite](#)

[Themenportale](#)

[Zufälliger Artikel](#)

[Mitmachen](#)

[Artikel verbessern](#)

 [Nicht angemeldet](#) [Diskussionsseite](#) [Beiträge](#)

[Benutzerkonto erstellen](#) [Anmelden](#)

Artikel

[Diskussion](#)

Weitere 



## Data Warehouse

Ein **Data Warehouse** (kurz **DWH** oder **DW**; wörtlich „Datenlager“, im Deutschen dominiert die englische Schreibweise, die Schreibweise *Datawarehouse* wird jedoch auch verwendet) ist eine für **Analysezwecke** optimierte **zentrale Datenbank**, die Daten **aus mehreren**, in der Regel **heterogenen Quellen zusammenführt**.<sup>[1]</sup>

# Business Intelligence Warehouse

## Alle Daten

- aus *verschiedenen Quellen*
- auf *einer Plattform*,
- *bereinigt, interpretiert, ggf. aggregiert*
- leicht *auffindbar, nutzbar* und *verknüpfbar*

Wie bringen wir Daten aus  
vielen Quellen in einen Topf?

# Our High-Performance Databases

## Exasol

## BigQuery

---

In Memory  
dutzende TB

Cloud Storage  
nahezy unbegrenzt

---

Super effiziente  
Joins und Views

Beliebig viele Clients  
Trennung von  
Storage & Compute

---

Skalierbar  
derzeit 40 Nodes

Automatische Skalierung  
fully managed

Wie erreichen wir Verknüpfbarkeit?

Wie schaffen wir es Dutzende von Teams zu koordinieren?

Und Dutzende von Fachbereichen zu versorgen?

# Warehouse

Core

Domain

Access

# Warehouse

Also alles gut, oder?

Core

Domain

Access

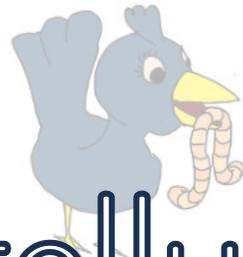
Aber ganz so einfach ist das auch wieder nicht ...



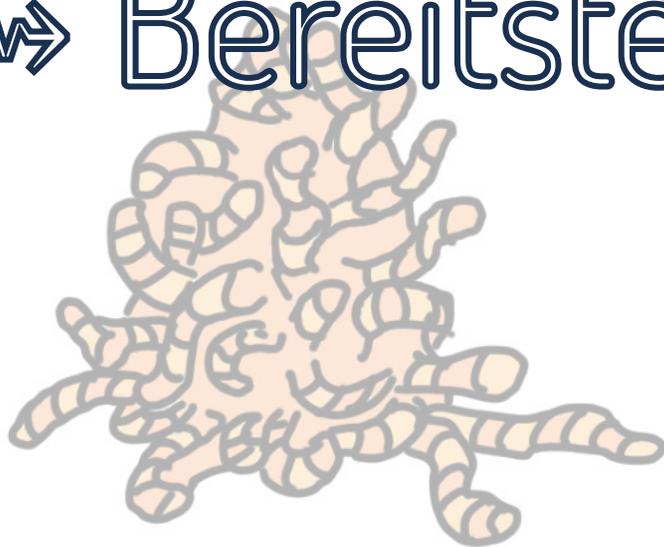
700+ Quellsystemen, zehntausende von Tabellen  
lädt man nicht "mal eben" in ein Warehouse.

Werden aber ASAP benötigt.

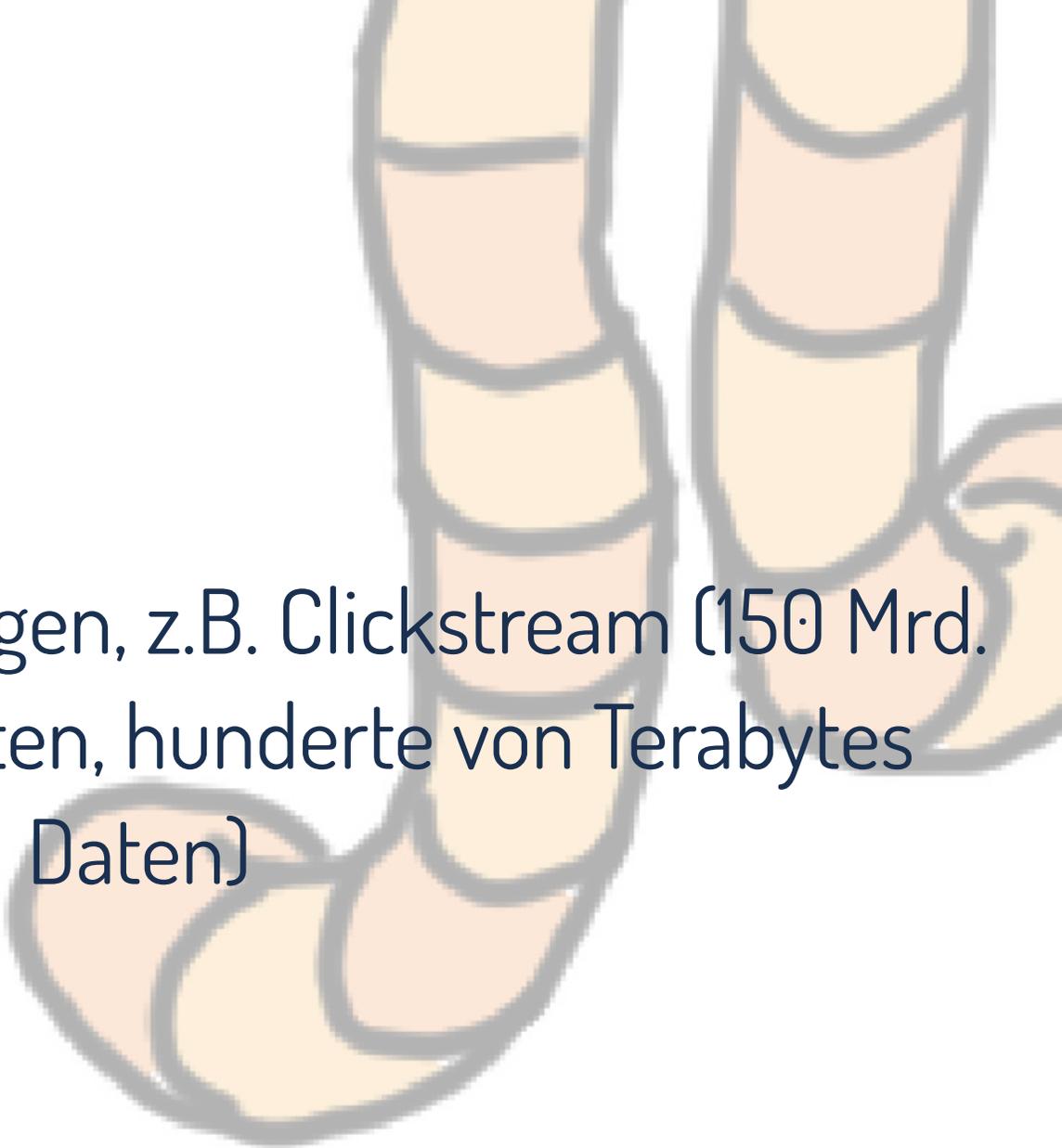




→ → Bereitstellung zu träge



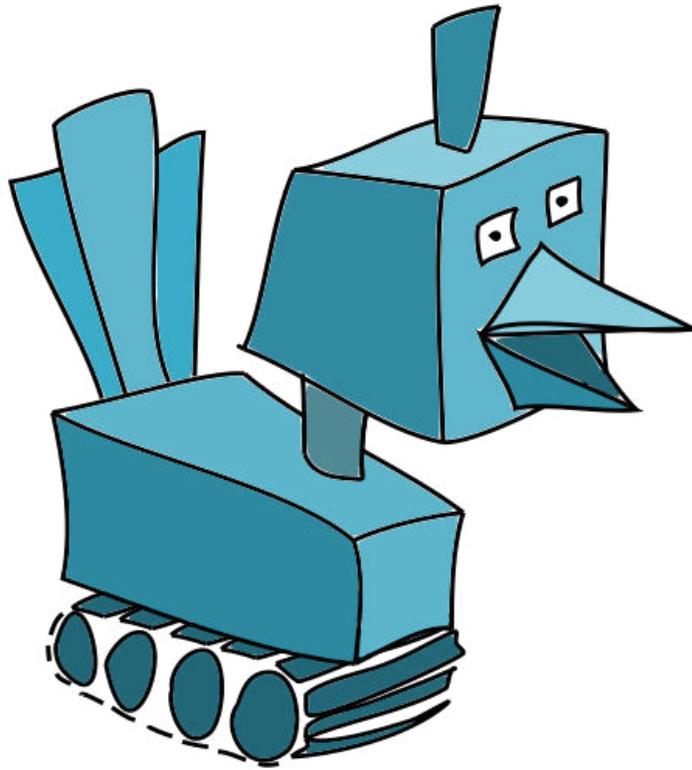


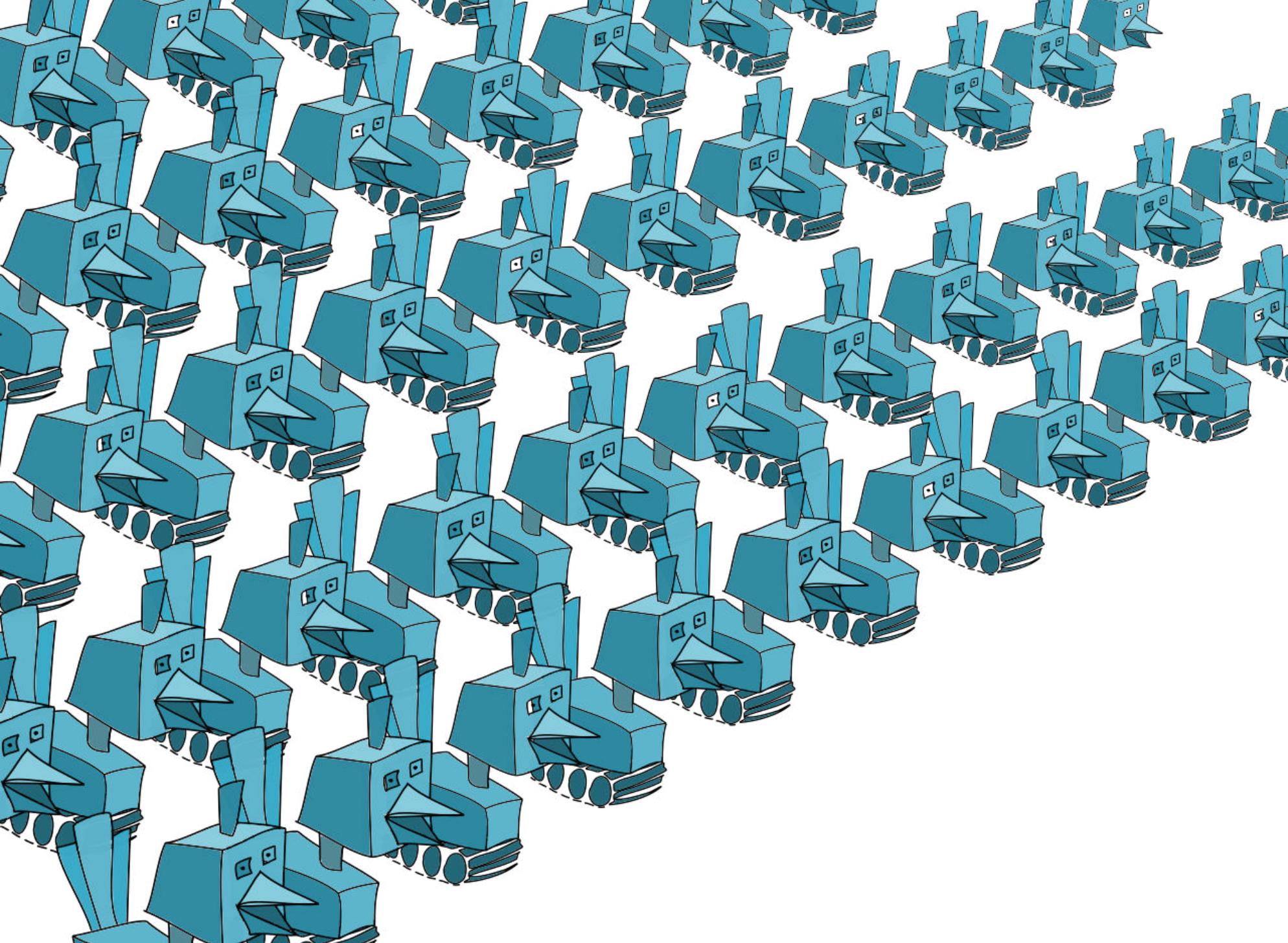


Gewaltige Datenmengen, z.B. Clickstream (150 Mrd.  
Zeilen, 2000+ Spalten, hunderte von Terabytes  
Daten)



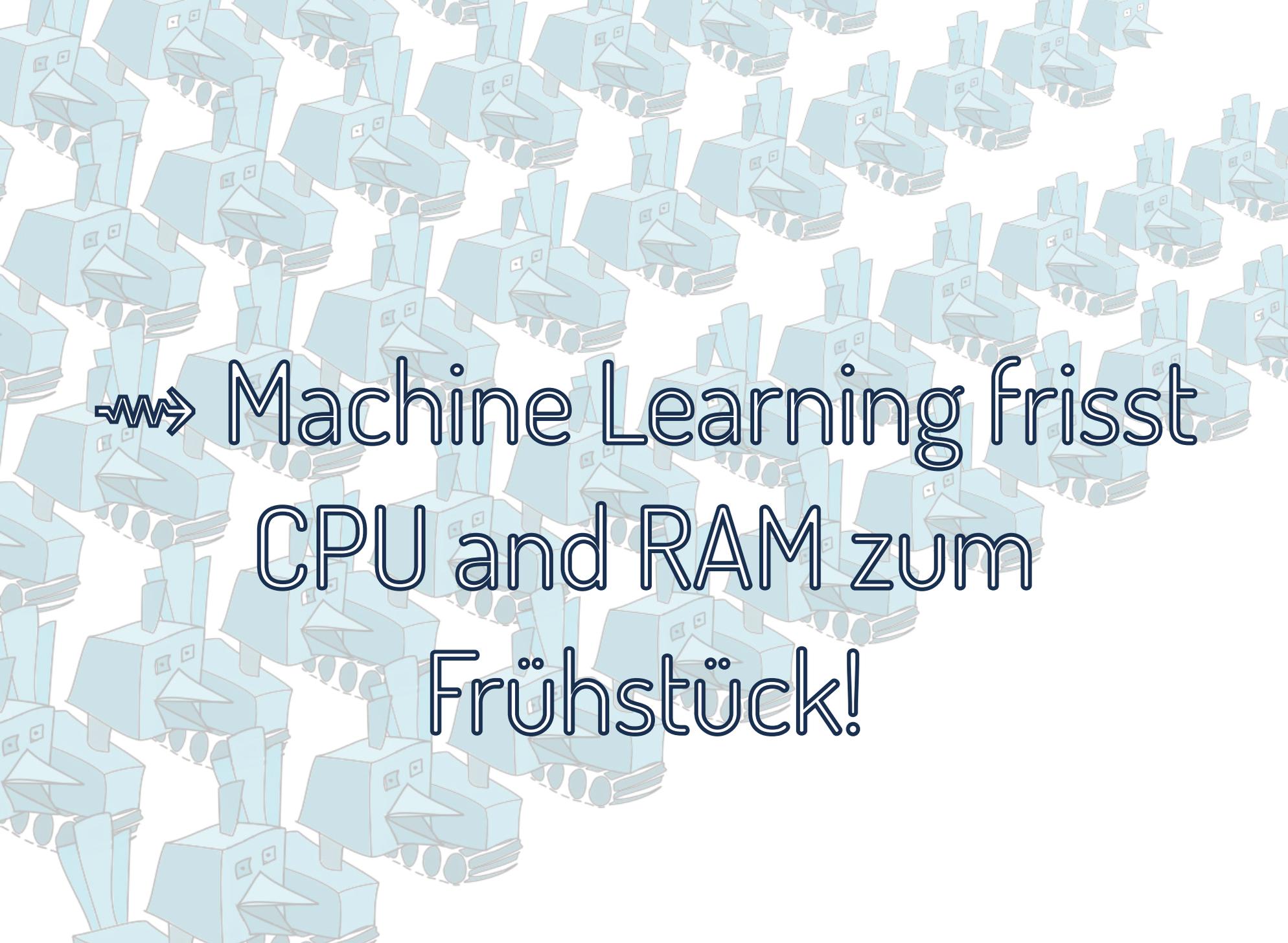
→ nicht 1:1 in der in-memory  
Exasol abbildbar



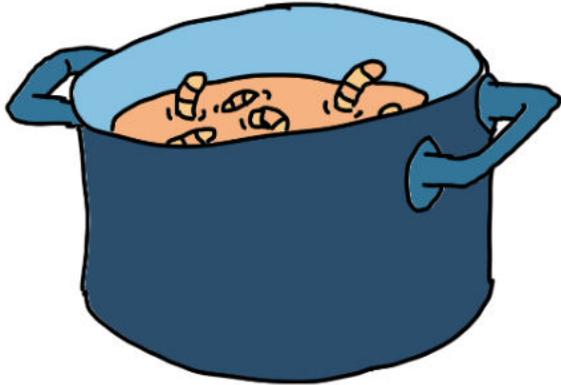




Machine Learning  
liebt Scans über  
alle Zeilen & alle Spalten.



→ Machine Learning frisst  
CPU and RAM zum  
Frühstück!





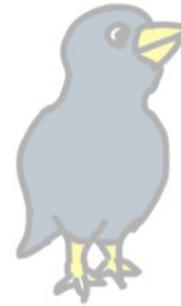
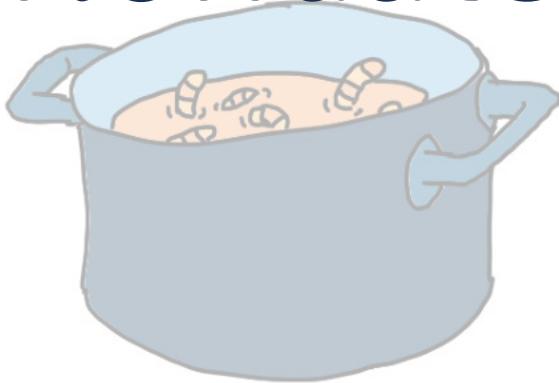
Manche mögen's roh!

Insbesondere Data Scientisten





→ Rohdaten erforderlich



# Warum kann man nicht einfach **alles** in das große **Warehouse** tun?

- Bereitstellung zu träge
- nicht 1:1 im Warehouse abbildbar
- ML frisst CPU and RAM
- Rohdaten erforderlich
- Archiv & Back-Up

Wir haben gelernt,  
**Data Lake** und **Warehouse**  
zu kombinieren.

## Data Lake

quellnah

roh

heterogen

semi-structured, binär

big data

redundanz

Archiv & Backup

## Data Warehouse

einheitlich modelliert

bereinigt

verknüpfbar

columnar, SQL

universell, verknüpfbar

single point of truth

# Lake

structured

unstructured

# Warehouse

Core

Domain

Access

# Domain

Assortment

Pricing

Finance



# Catalog

... und dadurch können wir

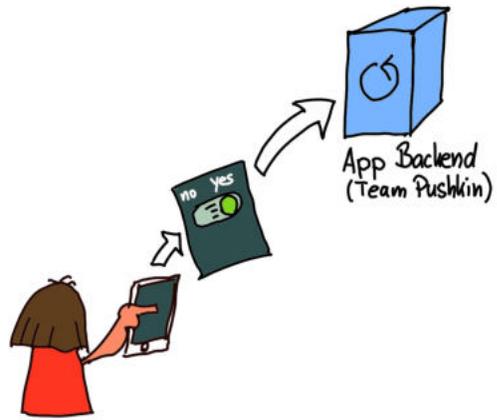
- interactive Queries
- Reporting
- Batch Processing
- Big Data
- Machine Learning

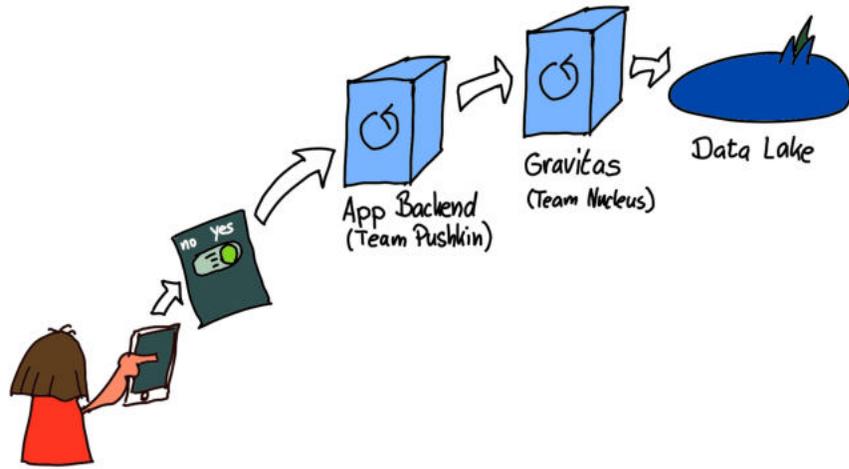
unterstützen.

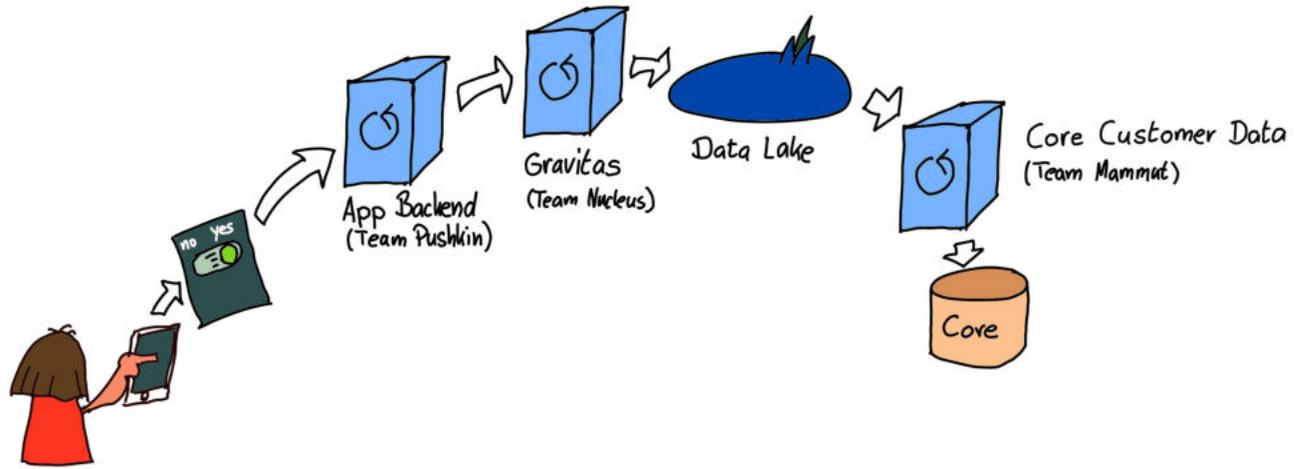
... und was heißt das jetzt in  
der Praxis?

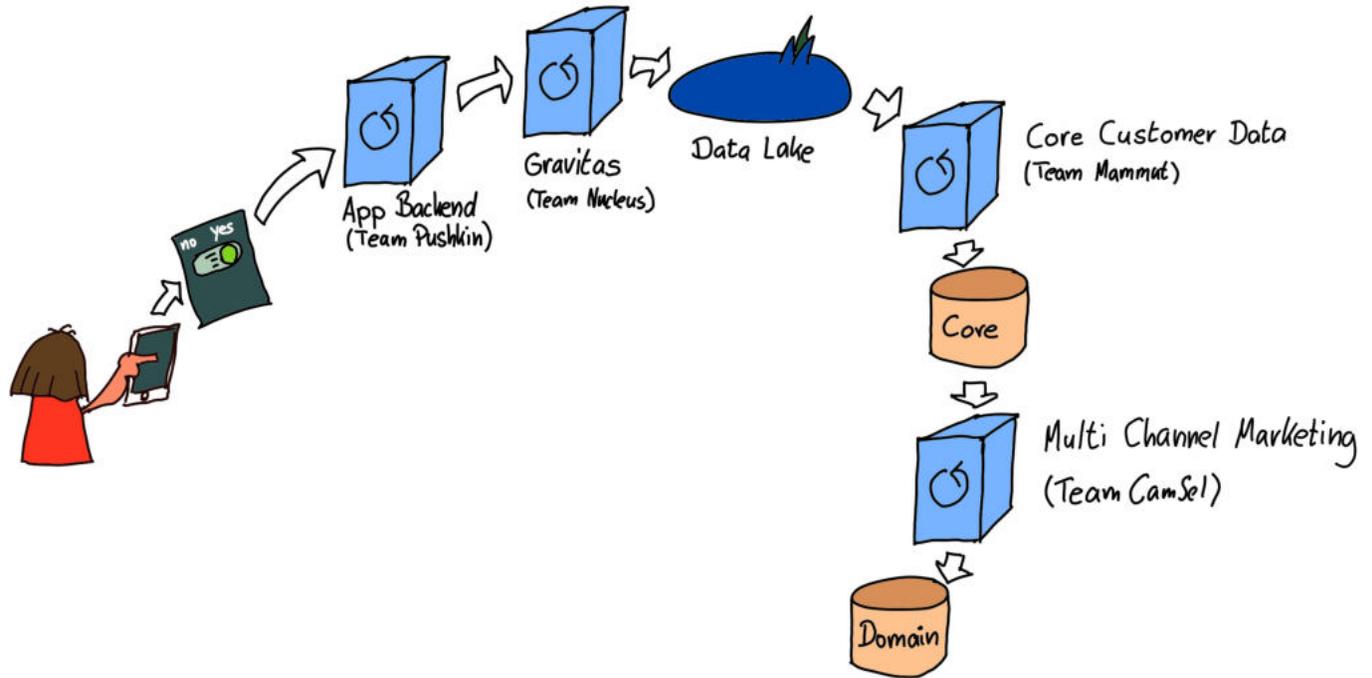


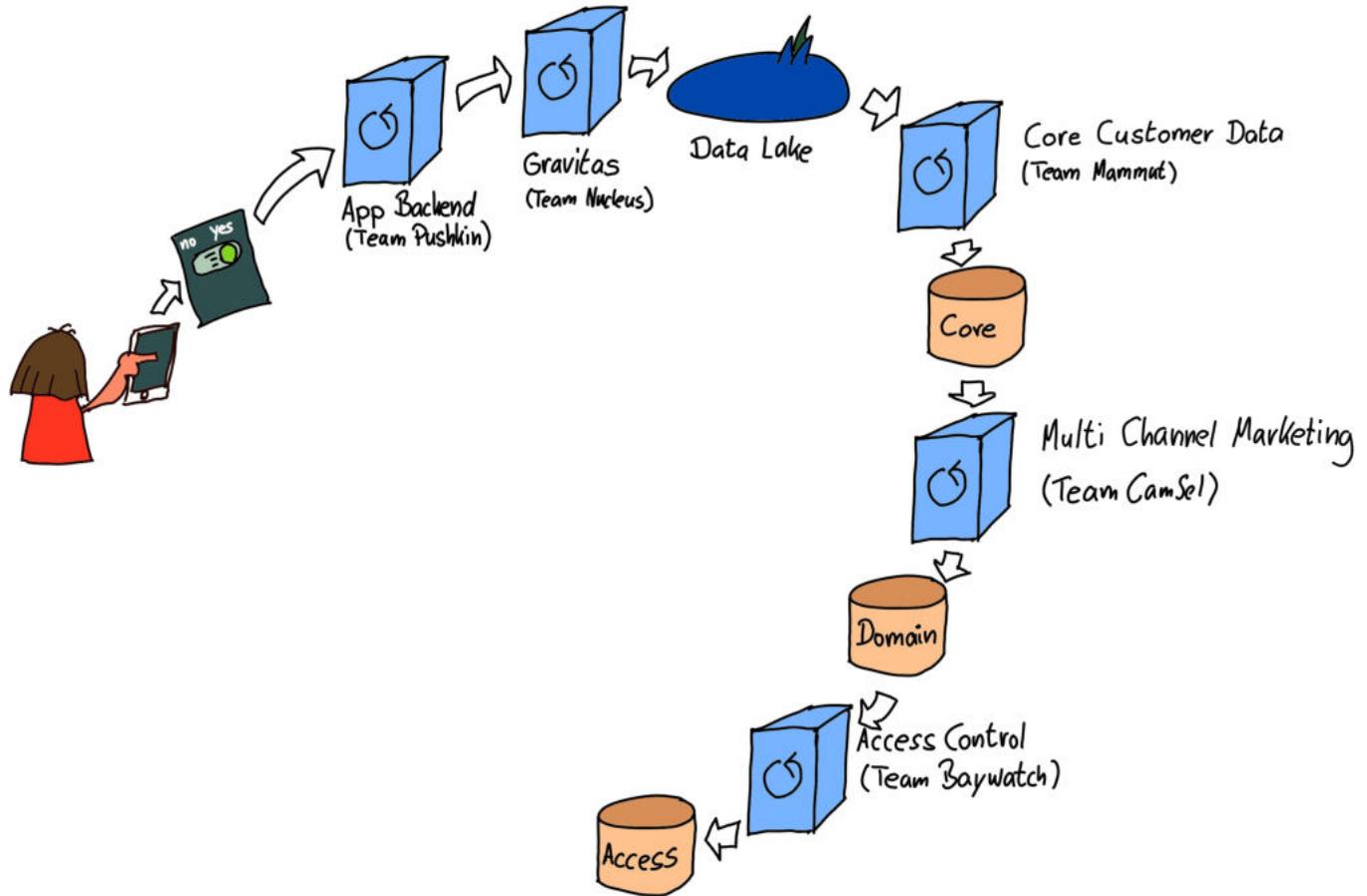


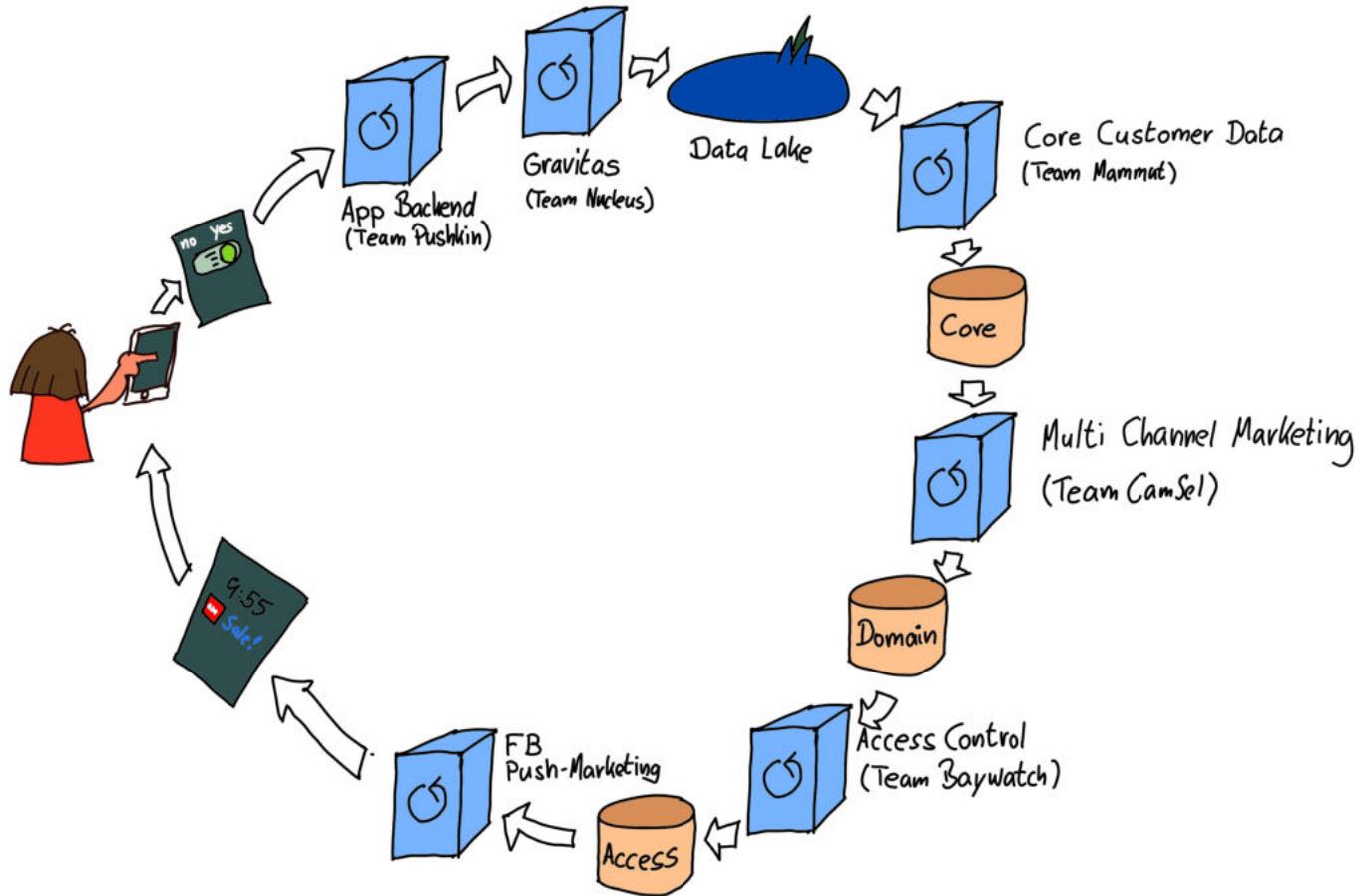


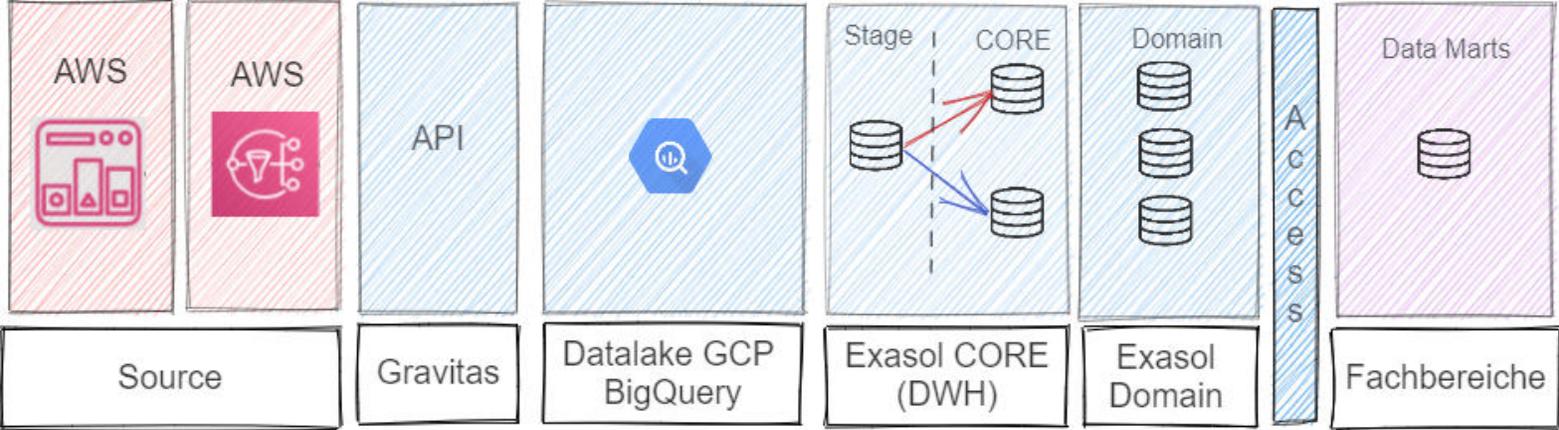


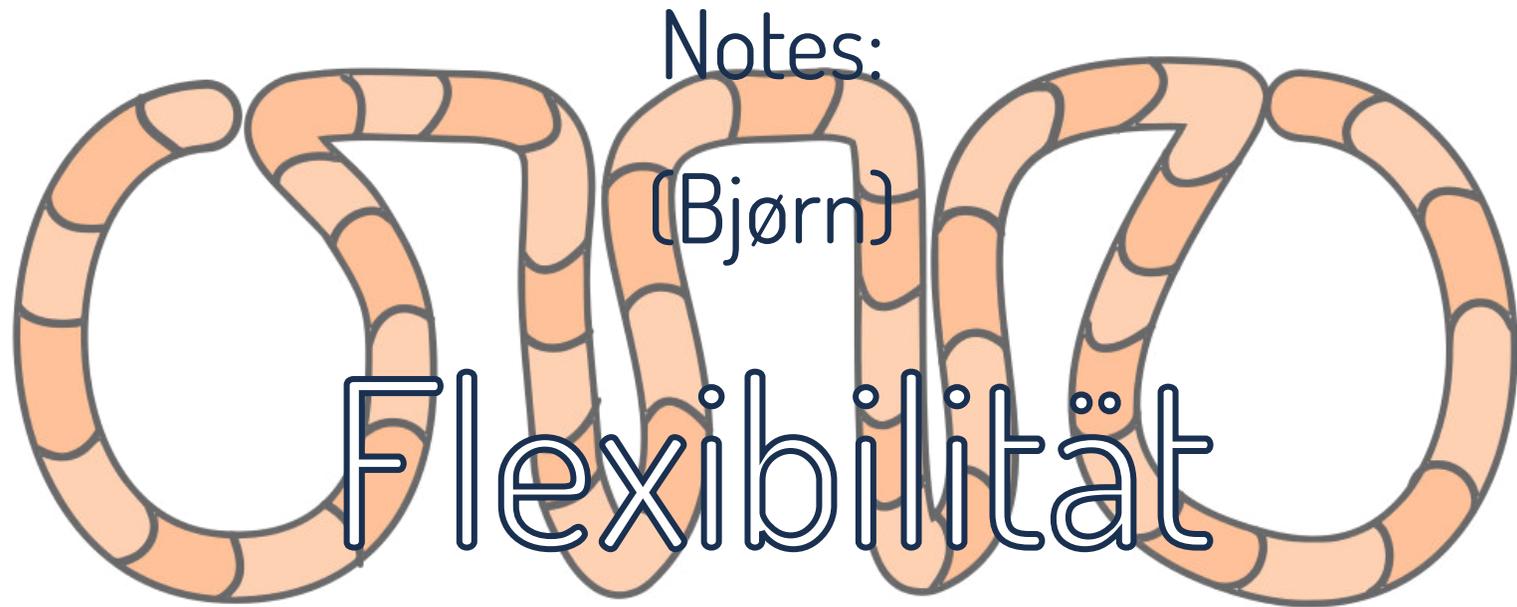












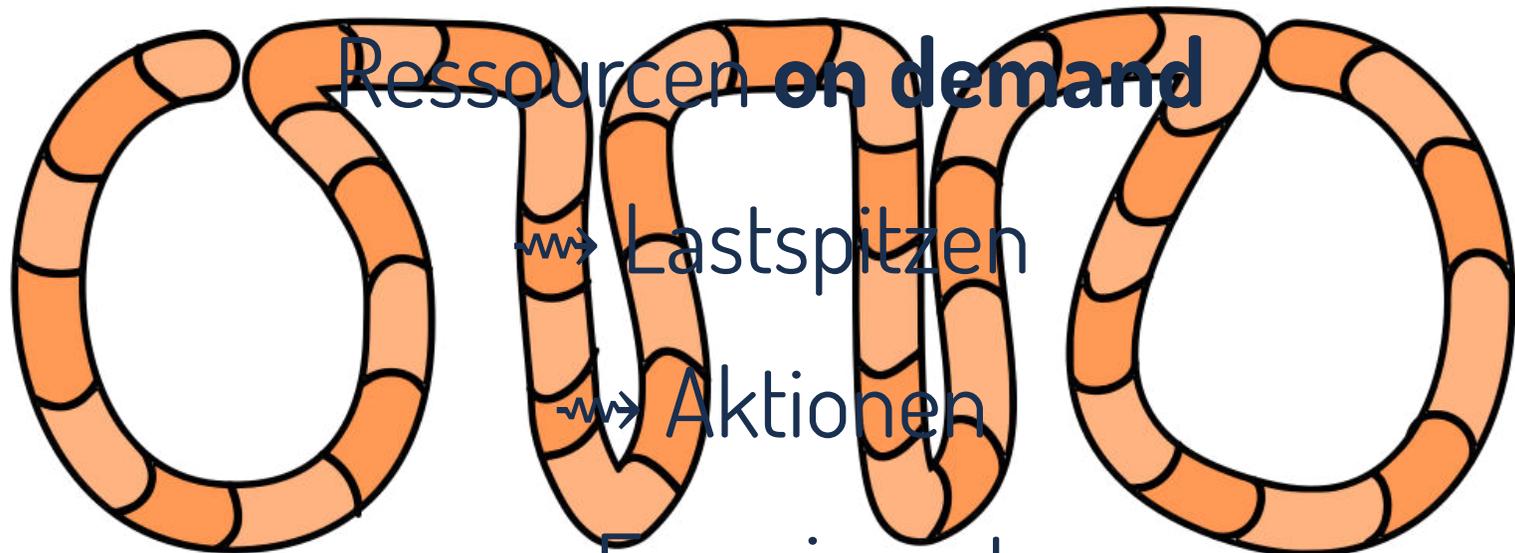
Notes:

(Bjørn)

Flexibilität

Notes: (Bjørn)

# Flexibilität, Unabhängigkeit



Ressourcen **on demand**

↳ Lastspitzen

↳ Aktionen

↳ Experimente

↳ Konkurrierende Lösungen

# Empowerment

Es muss für jeden bei OTTO leichter werden, Daten zu nutzen.



# Empowerment

- großes Ökosystem mit einsatzbereiten Produkten
- Managed Services  $\rightarrow$  einfach
- On-Demand, Kostentransparenz



