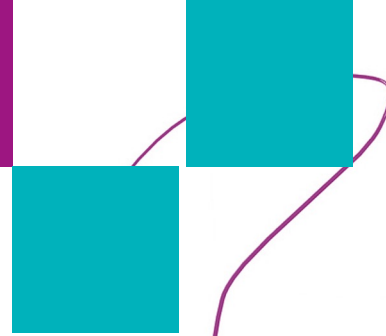


Datenzähmen leicht gemacht

Warum verlässliche
Data-Architekturen auf solidem
Software-Engineering basieren

Dr. Masud Afschar

iteratec



Agenda



Kurzvorstellung iteratec



Projekt und Technische Herausforderungen



Wichtige Erkenntnisse

A smiling man and woman are shown in a purple-tinted setting. The man is on the right, wearing a dark t-shirt, and the woman is on the left, wearing glasses and a dark t-shirt. They appear to be looking at something together. The background is a solid purple color. There are several colorful geometric shapes (squares and rectangles) scattered around the image: a yellow square in the top left, a white square and a cyan square below it, another cyan square to the right of the white square, a yellow square on the right side, a cyan square below it, a black square and a white square at the bottom right.

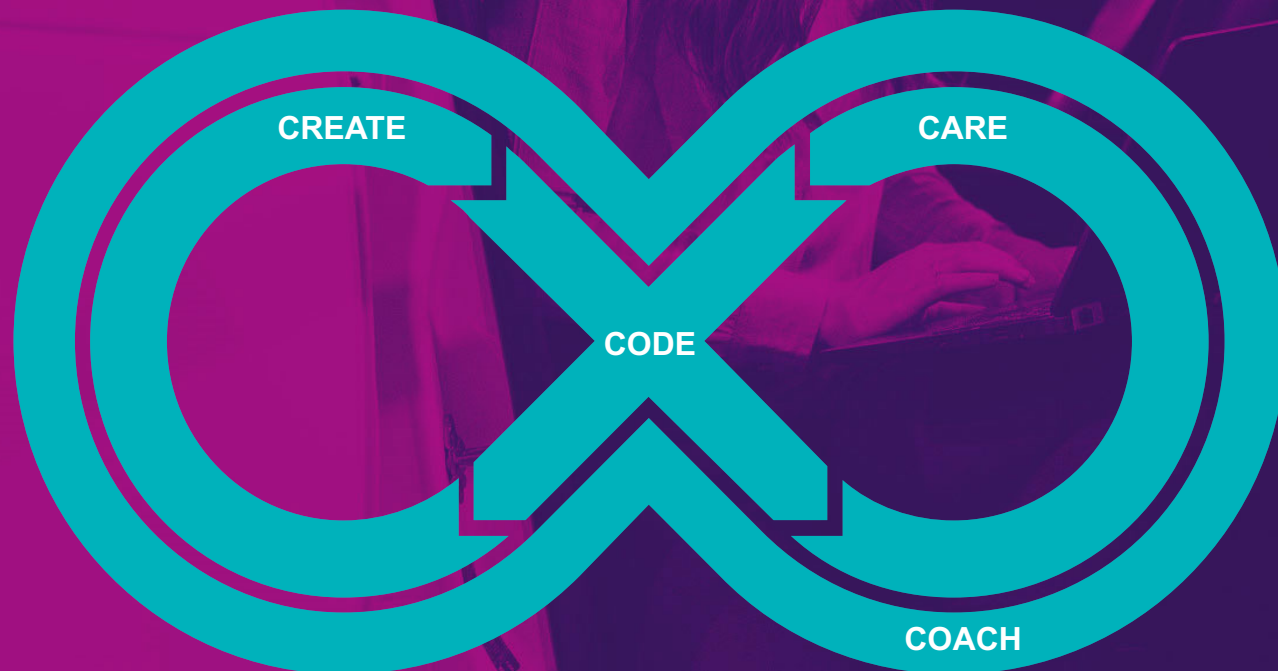
iteratec

Wie wir digitale Champions entwickeln

Innovationspotenziale erkunden, erproben und umsetzen

Digitale Lösungen & Leistungen entwickeln

Funktion, Performance und Sicherheit gewährleisten



Technologie und Wissen vermitteln, Zusammenarbeit und Strukturen verbessern



Gut aufgestellt

470

Kolleg*Innen

1000

Studierende

27

Jahre Erfahrung

64

Mio. Euro Umsatz

7

Standorte



Einstieg

Data Engineering und das Projekt



Was ist Data Engineering?

iteratec

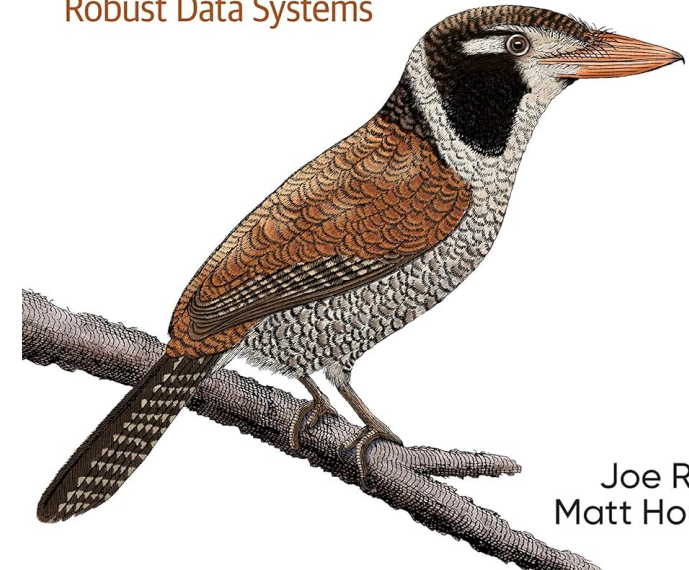
Und was macht ein Data Engineer?

- „Data engineering is the development, implementation, and maintenance of systems and processes that take in raw data and produce high-quality, consistent information that supports downstream uses, such as analysis and machine learning.“
- „A data engineer manages the data engineering lifecycle, beginning with getting data from source systems and ending with serving data for uses cases, such as analysis and machine learning“

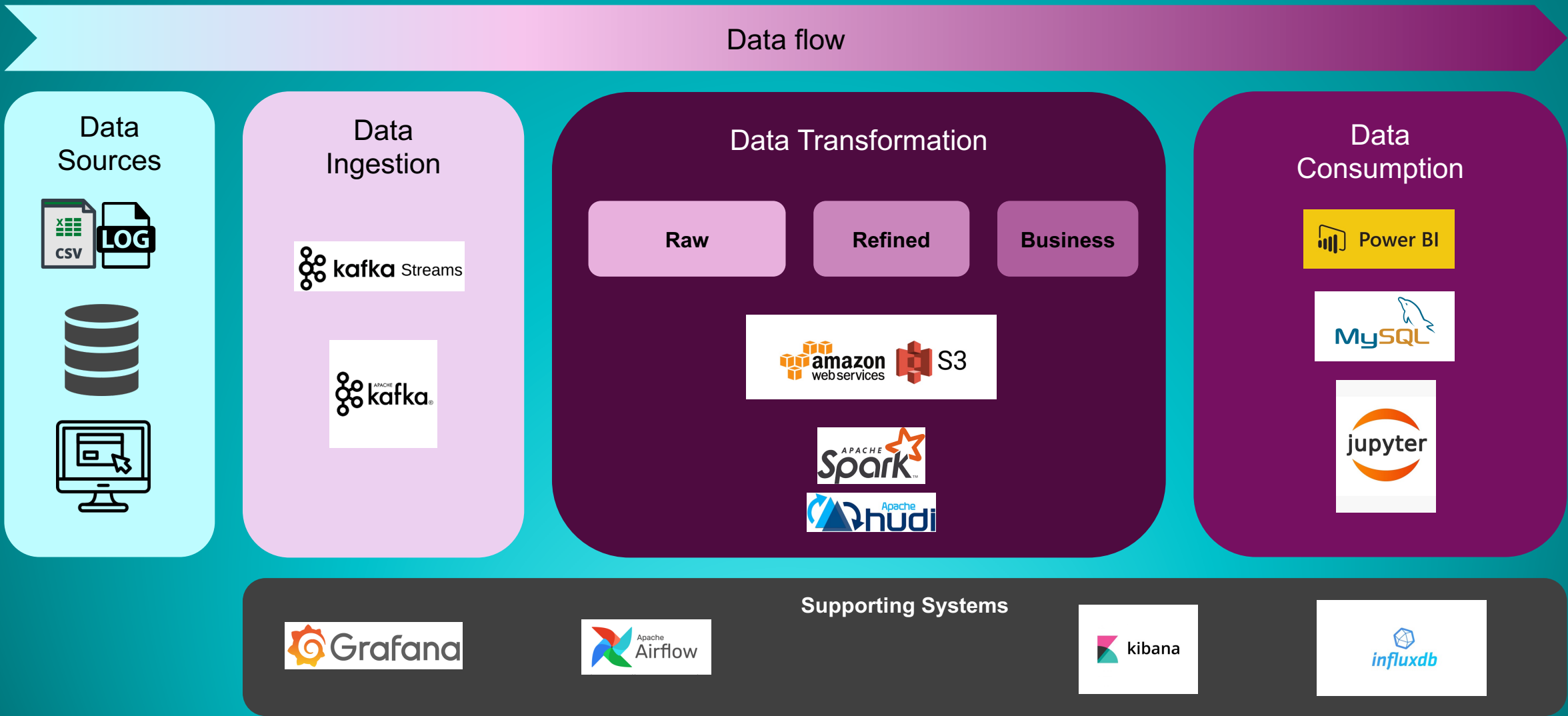
O'REILLY®

Fundamentals of Data Engineering

Plan and Build
Robust Data Systems



Joe Reis &
Matt Housley



Herausforderungen und Lösungsansätze

Pionierarbeit

- Erste Cloud-basierte Big Data-Plattform innerhalb des Unternehmens
- Seit mehr als 5 Jahren in Produktion

Datenquellen

Herausforderung: Extrahierung aus den operationalen Systemen
Lösung: Verringerung der Einstiegsbarrieren durch selbstentwickelte Libraries
Ergebnis: Alle relevanten Business-Domänen sind angebunden

Real-time Daten

Herausforderung: Anforderungen shiften in Richtung (near-)real-time
Lösung: Event-Driven Architektur von Anfang an
Ergebnis: Zukünftige near-real-time Anwendungsfälle können realisiert werden

Data-Ingestion

Herausforderung: Sicherstellung hoher Qualitätsstandards
Lösung: Einbindung von Domain Experts
Ergebnis: Hohe Datenqualität in der gesamten Transformationskette

Datenprozessierung

Herausforderung: Beherrschung komplexer Transformationsgraphen
Lösung: Einführung aufbauender Datenreifegrade und Smart-Scheduling
Ergebnis: bis zu 2.500 Transformationen pro Tag laufen unter kontrollierten Bedingungen

Data-Services

Herausforderung: Auslieferung der korrekten Daten an verschiedene Consumer-Rollen
Lösung: Individuelle Datenprodukte für Anwendungsfälle
Ergebnis: Diverse, aktive Datennutzerbasis im Unternehmen



iteratec

Herausforderung: Datenschutz



Datenschutz "at scale"

Harmonisierung von Consumer-Anforderungen und Data-Governance



iteratec



Data-Konnektoren

- Standardisierte API
- Einfache Integration



Data-Realms

- Aufteilung gemäß Anforderungen
- Einstellbares Routing



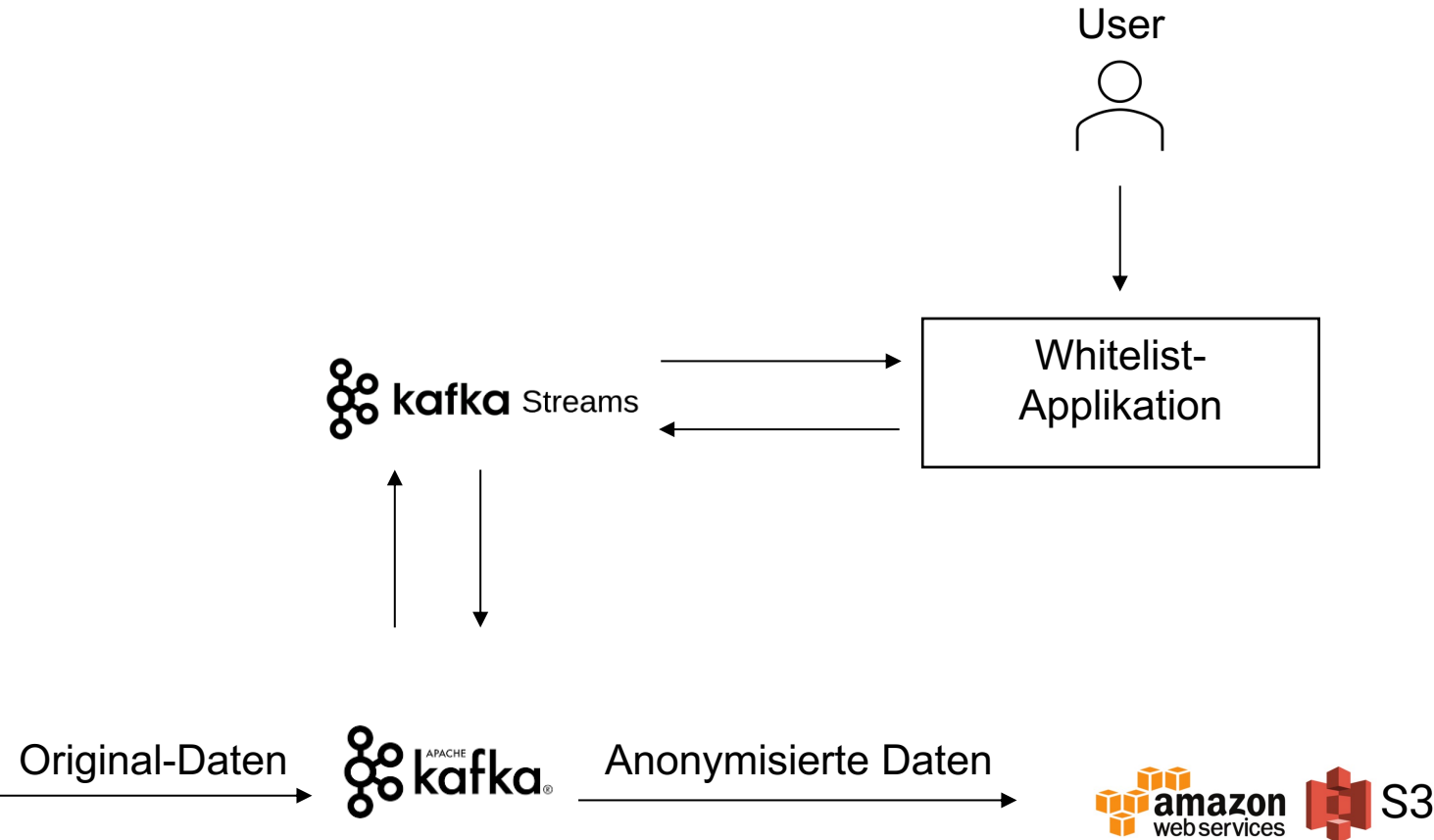
Quality-Gates

- Individuelle UI
- Verschiedene Modi

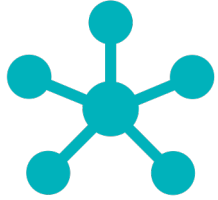
Datenanonymisierung

Technische Lösung

iteratec



- Daten werden von Domänen an Kafka geschickt
- Separate Kafka Streams-Applikation greift Daten sofort ab
- Erhält aus eigens entwickeltem System sog. Whitelist
- User können Whitelist über UI pflegen und Anonymisierung konfigurieren



iteratec

**Herausforderung:
Komplexität**



Datentransformation

Komplexität beherrschen



iteratec

- Über 2.000 Transformationen pro Tag
- Unterschiedliche Granularität und Load
- Abhängigkeiten bilden komplexe Verarbeitungsgraphen

Datentransformation

Komplexität beherrschen



iteratec

- Über 2.000 Transformationen pro Tag
- Unterschiedliche Granularität und Load
- Abhängigkeiten bilden komplexe Verarbeitungsgraphen

- Smartes Scheduling und Monitoring
- Eigenentwickeltes Test-Framework für anwendungsfallorientierte End-To-End-Tests

Herausforderung: Austauschbar- und Erweiterbarkeit



Data Services

Flexible Datenarchitekturen für neue Anwendungsfälle

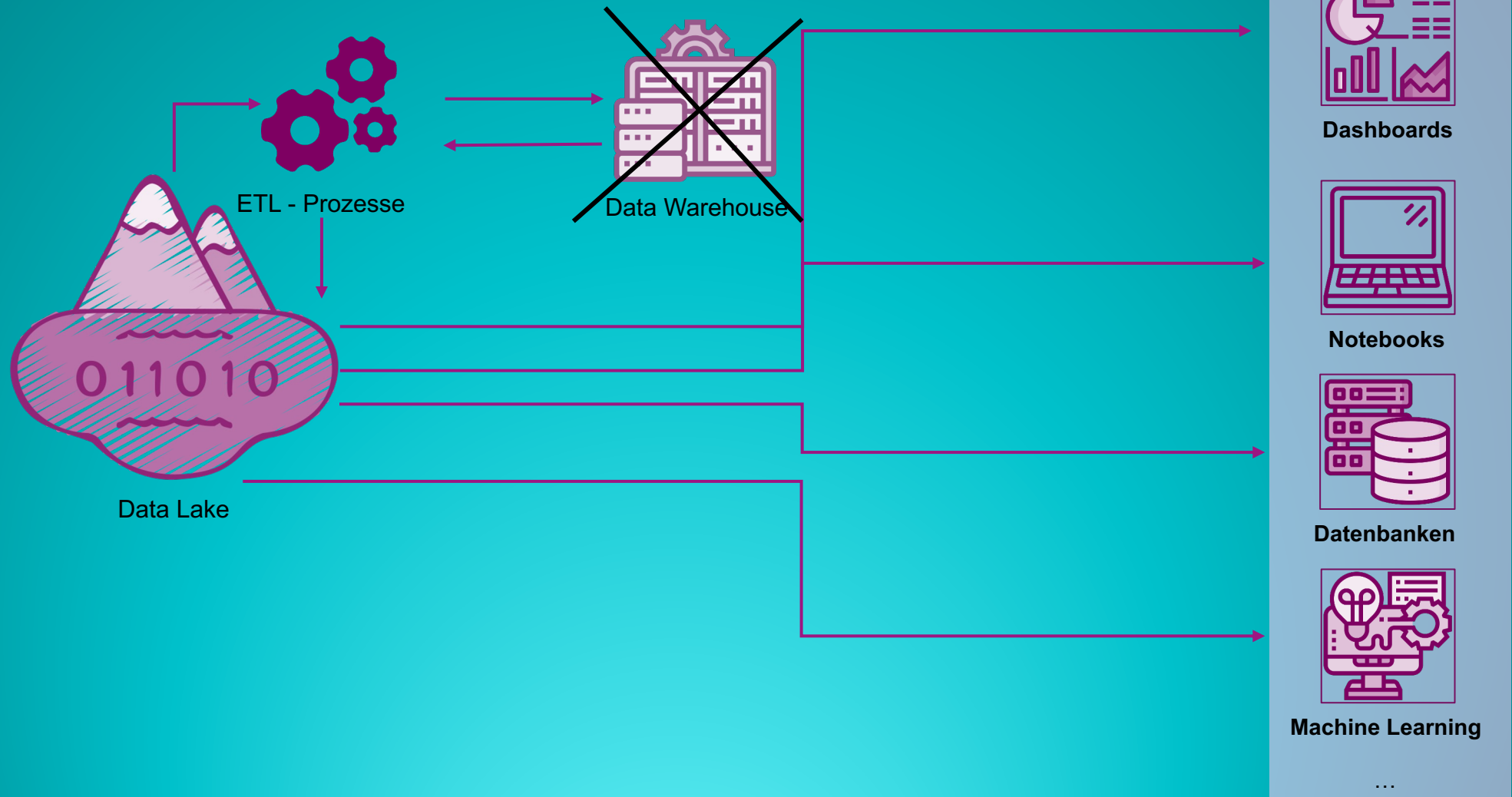
iteratec



Data Services

Flexible Datenarchitekturen für neue Anwendungsfälle

iteratec



iteratec

Abschluss



Warum gutes Data-Engineering gutes Software-Engineering voraussetzt

- Datenarchitekturen sind weiter auf individuelle Software-Lösungen angewiesen
 - Vorhandene Tools müssen ggf. ergänzt, erweitert oder verbaut werden
 - Lösungen haben nur marginal mit „Data“ zu tun
- Datenarchitekturen haben wie jedes andere Software-System hohe Anforderungen an Security, DevOps, Logging, Modularität, etc.
 - Data-Engineers müssen Systeme hoch-verfügbar, robust, resilient, fehlertolerant, sicher halten
- Integratives Testen im Data-Bereich benötigt bis heute viel „Handarbeit“
- Cutting-Edge-Technologien brauchen zur Integration umfangreiches Engineering-Wissen
 - Je neuer desto weniger Erfahrung, Dokumentation, ...



Dr. Masud Afschar
Senior AI & Data Engineer

iteratec GmbH
Westhafenplatz 1
60327 Frankfurt
+49 170 3748731
masud.afschar@iteratec.com

