

# Big Data on OSS @ SBB

Adrian Burri  
Zug, 14.05.2024





# Agenda.



OSS @ SBB



Train Delay Analysis



OSS for Big Data



Discussion

# Usage: Transition from Closed Source to Open Source

Dimension

until 2015

today

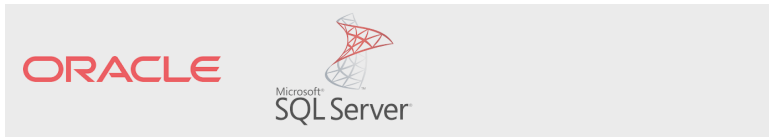
Application Runtime



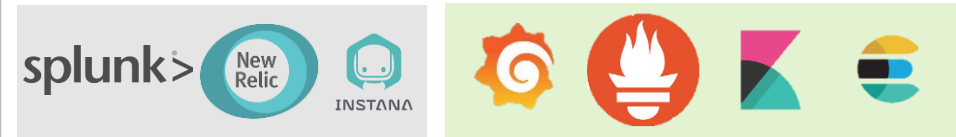
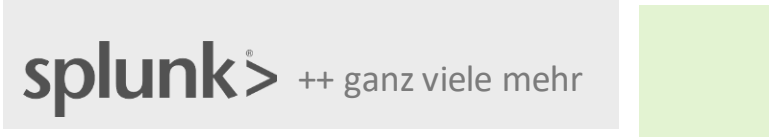
Messaging



Databases/Caches



Monitoring



DevOps Tools



Infrastructure



Open Source  
Closed Source

# OpenRail Association to Foster Collaboration of European Rail Companies



## OpenRail Association Governance

- **Neutral basis**
  - Open source licenses
  - Open source collaboration model
  - Stimulate active contribution
- **Project governance** on project level

## SBB Open Source Guide

- Strategic goals
- Principles for Use/Create
- Structured process





## Cooperation in Switzerland: Driver Advisory System

- **Cost Sharing:** SBB develops for itself and different other railway companies (BLS, SOB)
- **Open Source License**
  - Use
  - Enhancements
- **Transparency** for partners
- Built on open **standards**
- **International** opportunities





# Agenda.



OSS @ SBB



Train Delay Analysis



OSS for Big Data



Discussion

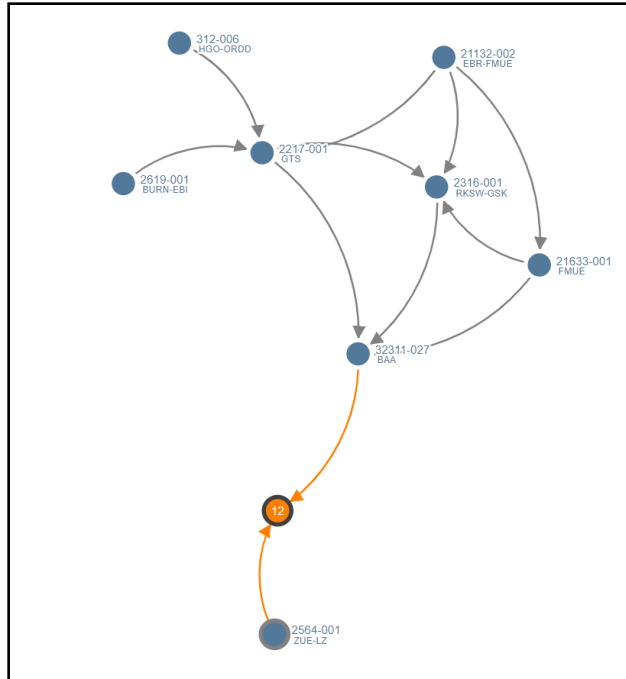
# Punctuality: Our DNA



## Definition:

Arrival time < 3 min.  
after schedule

# Train Delay Analysis



**Delay Analysis:**  
Root Cause Statistics,  
Deviation Graphs

**Planning:**  
Timetables, Connections,  
Rolling Stock, etc.

**Bern → Interlaken Ost**  
SBB CFF FS Richtung Interlaken Ost  
11 12 13 14 15 16 17 18 19 20 21 22 23 24 25 26 27 28 29 30 31 32 33 34 35 36 37 38 39 40 41 42 43 44 45 46 47 48 49 50 51 52 53 54 55 56 57 58 59 60 61 62 63 64 65 66 67 68 69 70 71 72 73 74 75 76 77 78 79 80 81 82 83 84 85 86 87 88 89 90 91 92 93 94 95 96 97 98 99 100

Streckenverlauf von Romanshorn

15:58		
16:04	Bern	Gleis 6
	1.   2.	
16:23		
16:24	Thun	Gleis 1
	1.   2.	
16:33		
16:34	Spiez	Gleis 1
	1.   2.	
16:51		
16:53	Interlaken West	Gleis 1
	1.   2.	
16:58	Interlaken Ost	Gleis 7

**Operations:**  
Actual Times, Deviations



**Dispatching:**  
Incidents, Dispatcher  
Actions





# Example Code Assignment

## 91: Special Environmental Conditions





# Agenda.



OSS @ SBB



Train Delay Analysis

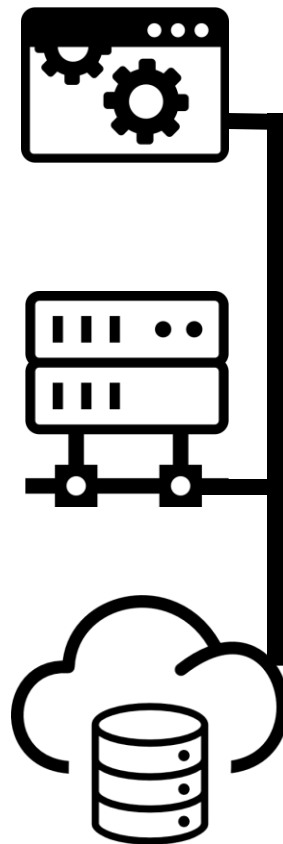


OSS for Big Data



Discussion

# Overview Tech Stack



**Workflow**

Apache Airflow

**Analysis**

jupyter

**Monitoring**

Prometheus Grafana

**CI/CD**

argo **TEKTON**

JFrog ARTIFACTORY Renovate

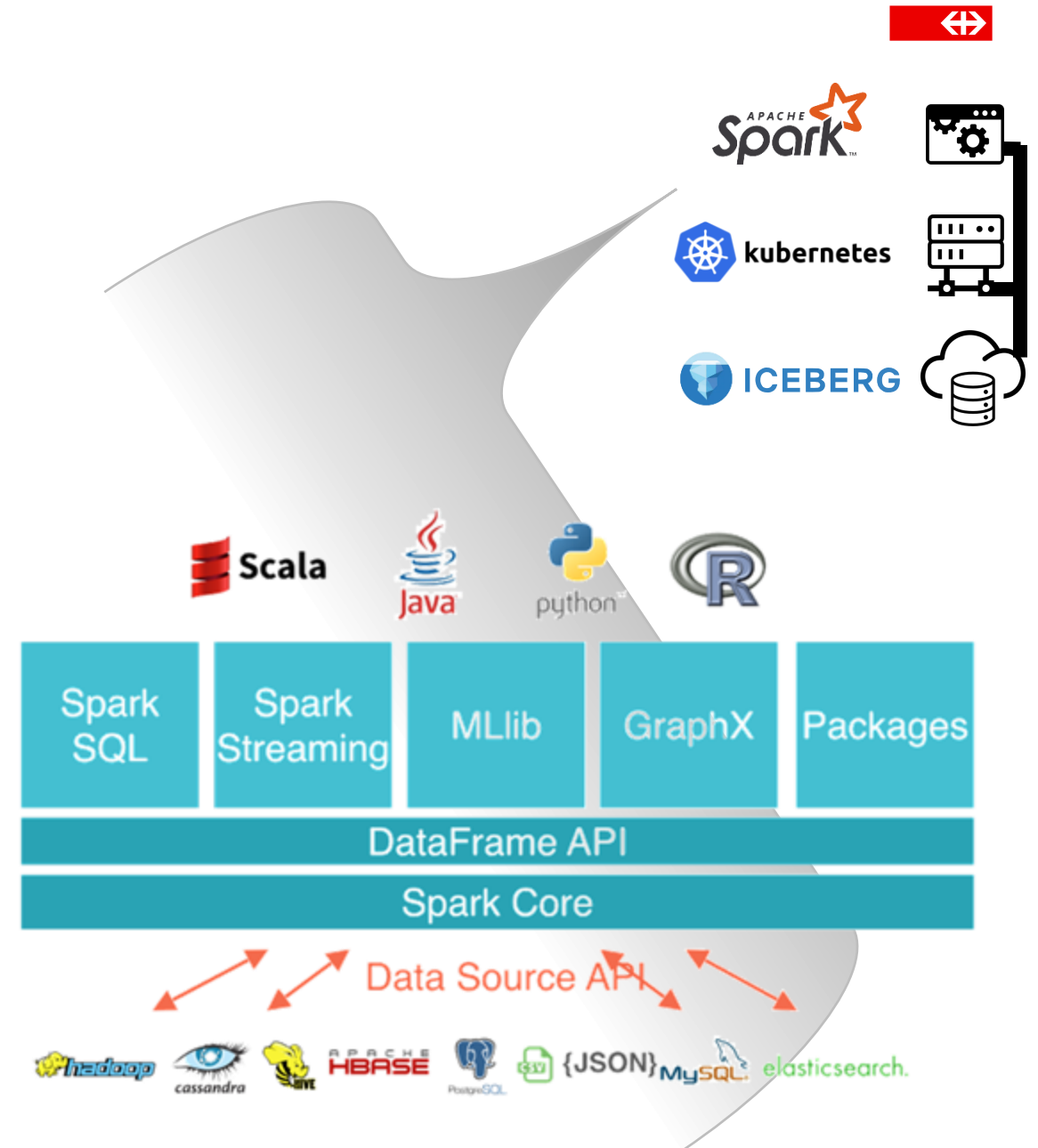
# Apache Spark

**Type:** Unified engine for data analytics  
**OSS since:** 2010

Spark is based on the concept of an RDD (Resilient Distributed Dataset)

## Key benefits:

- Serial programming model with parallel execution (incl. fault tolerance)
- In-memory processing
- Versatility
  - Many connectors (file, DB, messaging,...)
  - Many use cases (SQL, streaming, ML,...)



# Spark on Kubernetes

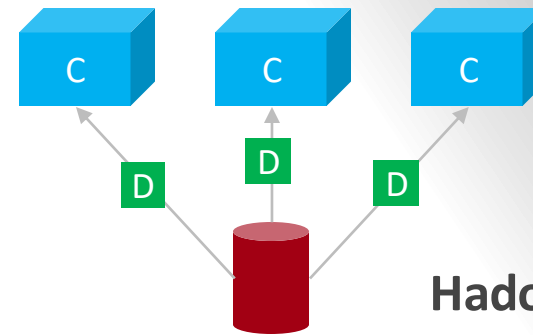
**Type:** Container orchestration  
**OSS since:** 2015

Runtime for transactional and analytical workloads in cloud environments.

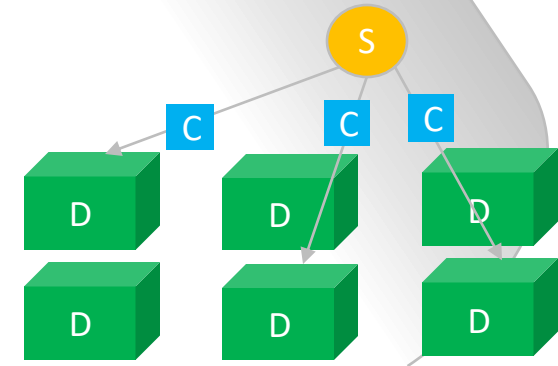
**Key benefits:**

- Scalability of workloads
- Portability between different cloud providers
- Rich ecosystem for DevOps

**Kubernetes: Data-to-code**



**Hadoop: Code-to-data**



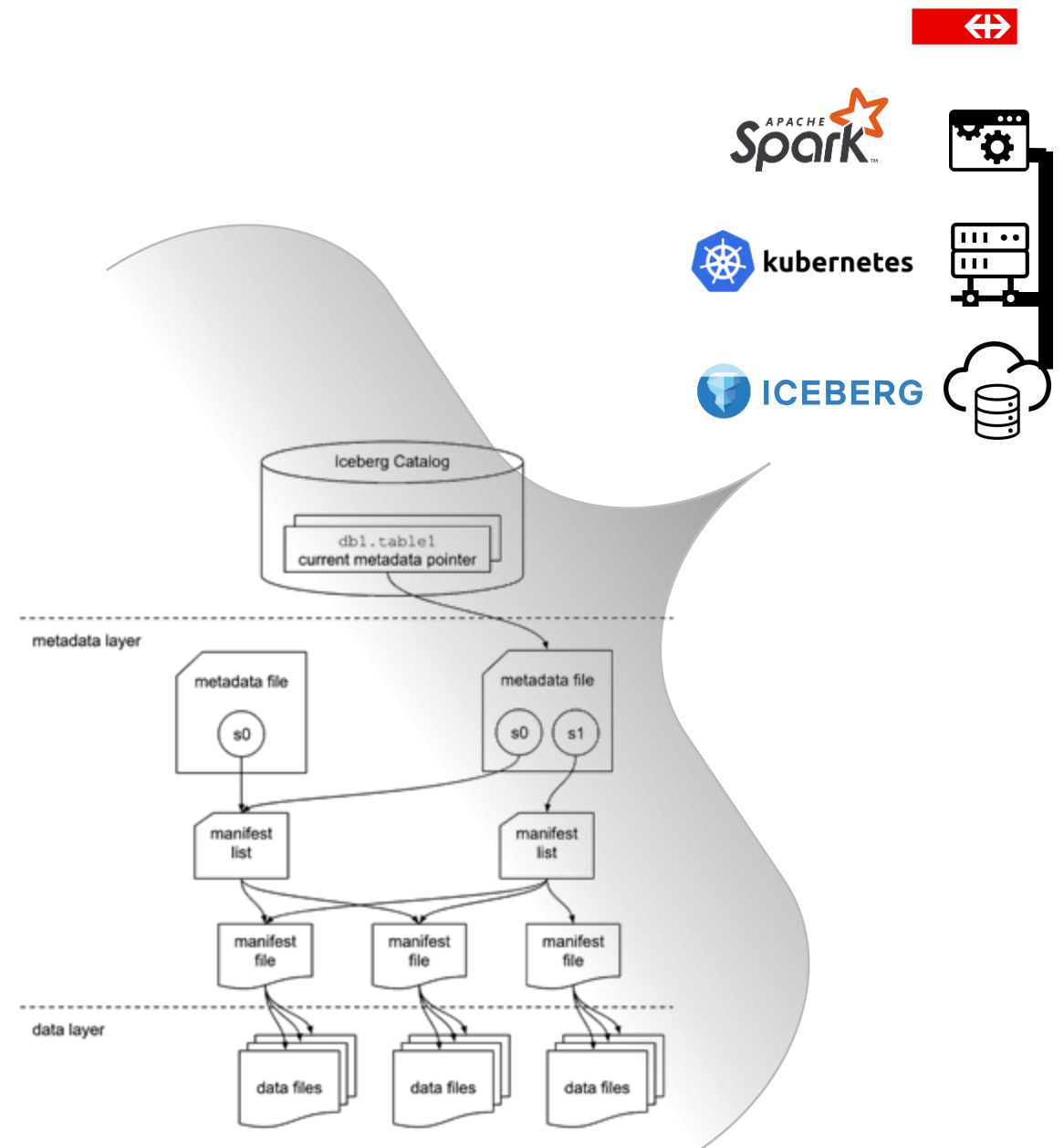
# Apache Iceberg

**Type:** Open table format for analytics  
**OSS since:** 2018

Iceberg provides a table abstraction on unstructured blob storage through hierarchical metadata.

## Key benefits:

- Correctness through ACID transactions
- Efficient query planning
- Easy schema evolution (incl. partition change)
- Timetravel





# Agenda.



OSS @ SBB



Train Delay Analysis



OSS for Big Data



Discussion

# Key Takeaways

Analytics workloads can be processed in a cost-efficient on a purely OSS tech stack.

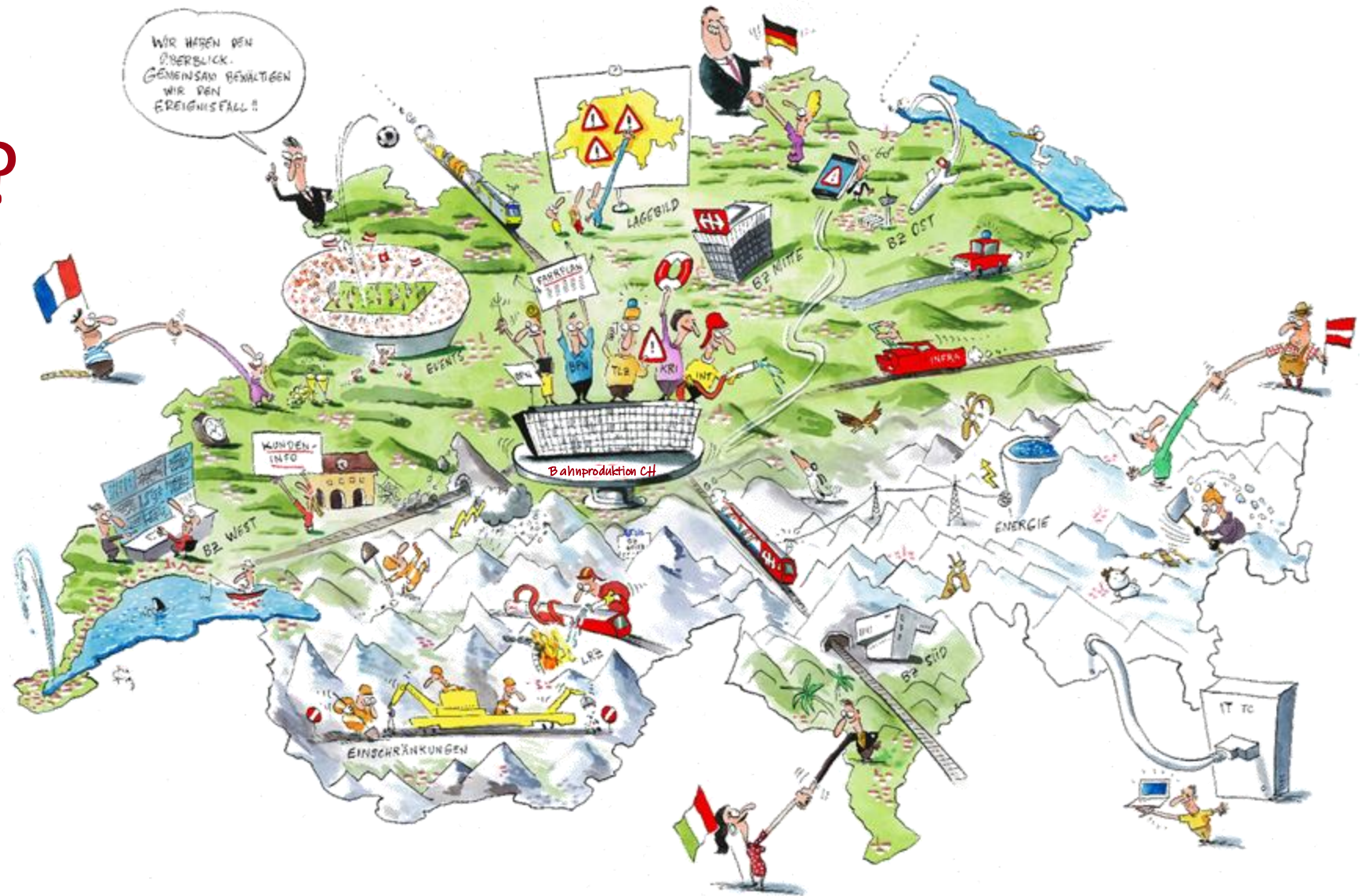
## Key ingredients:

- Spark as analytics engine offers a simple and efficient computing model for various Big Data use cases
- In a cloud environment, the same runtime environment (Kubernetes) as for transactional workloads can be re-used
- Iceberg (or other table formats) allow for cheap table-like storage in unstructured cloud storage





# Questions?



A group of people are sitting around a table, looking at a smartphone held by one of them. The scene is a close-up, focusing on the hands and the phone. The background is blurred, showing the upper bodies and arms of several people. The lighting is warm and natural.

Danke, merci, grazie  
and thank you.