# Data is the center of every digital experience.

In order to excel in digital, companies need to manage these 4 components through people, process & systems:

| **Infrastructure Monitoring** | + | **Data Monitoring** | + | **Model Monitoring** | + | **Application Monitoring** |
|---|---|---|---|---|---|---|
| IT / Infrastructure: "DevOps" | | Data Owners & Engineers: "DataOps" | | ML & Data Engineers: "MLOps" | | Software Development: "DevOps" |
| Platform: Data Dog | | Platform: Soda | | Platforms: Dataiku | | Platforms: New Relic |

# What could possibly go wrong?

Definitely. Many issues are actually silent (go unnoticed), until someone shouts from the rooftops.

- **Data ownership**

  Use of poorly architected systems
  Data contracts / SLAs not met
  Lack of process enforcement

- **Data engineering**

  Software/firmware changes
  (Partially) missing data
  Schema evolution

- **Data analytics**

  Errors in interpretation
  Changing data requirements
  Annotation distribution

⬇ Customer experience

⬆ Reputation risk

⬆ Operational inefficiencies

# Solution: Data Monitoring-as-a-Service

1  **Automatically instrument your data to create visibility & transparency at scale.**
   ML-driven analysis of data at rest and in-motion, whether it's in your data warehouse, lake or elsewhere. Analytics teams can layer-in their data requirements and track data usage SLAs.

2  **Analyse and collaborate on meaningful alerts using rich context.**
   Get alerts based on your role and inspect those using rich diagnostics and contextual info. Involve SMEs, find the appropriate resolution, and escalate problems where needed.

3  **Proactively share the status of datasets via the Data Quality Catalog.**
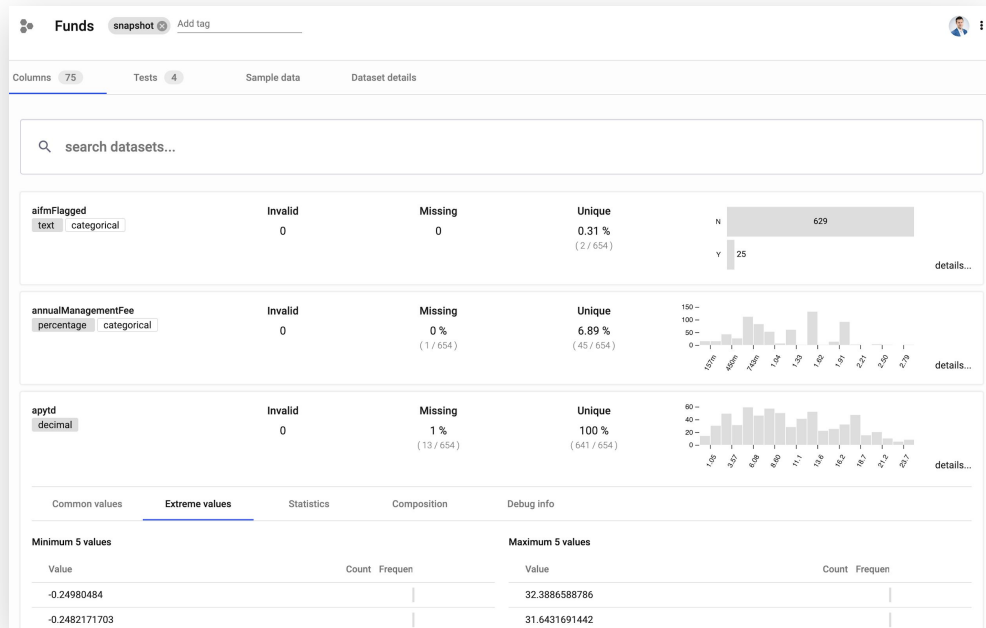   Give consumers of data products real-time insight into the quality of the data. Publish data quality indicators where data is searched and consumed to help increase trust in data.

# Feature Highlights

SODA

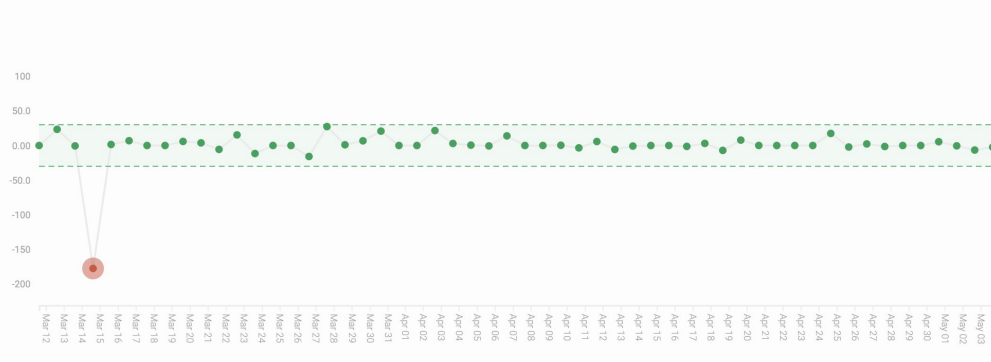# Soda automates data profiling for Data Quality

When new data comes in, Soda automatically re-runs profiling and analyzes the results historically. The approx. 50 out-of-the-box profiling metrics can be extended by adding custom metrics or test definitions.

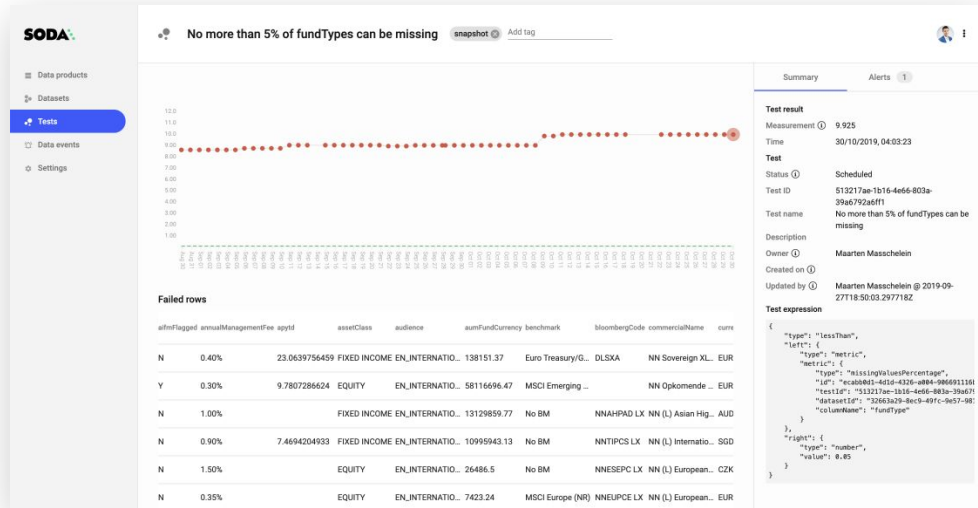# Soda finds and analyzes issues with your data.

Soda automatically measures key characteristics of your data over time in search of anomalies. Layer in business logic (via self-service) to get a robust operating model for data quality in place.



The average of all Fund NAV should not increase or decrease more than 30 day-over-day.

SODA:

# Route non-compliant records to the right people

Every metric Soda monitors contains details on what was calculated. If you're validating record-by-record, these details include all non-compliant records.

# Soda helps you prioritize & analyze alerts

Soda helps you understand which alerts are important to solve first by analyzing how many data products will be impacted and what the criticality of those are.

🔔 **Data events**

Reset to default

**Assigned to** ⓘ

👤

**Created** ⓘ

📅 dd/mm/yyyy — 📅 dd/mm/yyyy

**Criticality** ⓘ
- ☑ Critical
- ☑ Major
- ☑ Minor

**Dataset** ⓘ

Type a name of a dataset here...

**Columns** ⓘ

Type a column here...

| | | |
|---|---|---|
| 2 days ago | The average of all Fund NAV should not increase or decrease more than 25 day-over-day | |
| Critical | Assigned to Thijs Dhulster | |
| a week ago | The average of all Fund NAV should not increase or decrease more than 25 day-over-day | |
| Minor | Assigned to Tom Baeyens | |
| 3 weeks ago | The average of all Fund NAV should not increase or decrease more than 25 day-over-day | |
| Major | Assigned to Thijs Dhulster | |
| 3 weeks ago | The average of all Fund NAV should not increase or decrease more than 25 day-over-day | |
| Major | Assigned to Maarten Masschelein | |
| a month ago | No more than 5% of fundTypes can be missing | |
| Major | Assigned to Tom Baeyens | |
| a month ago | No more than 30% of WKN Codes can be missing | |
| Critical | Assigned to Tom Baeyens | |

SODA

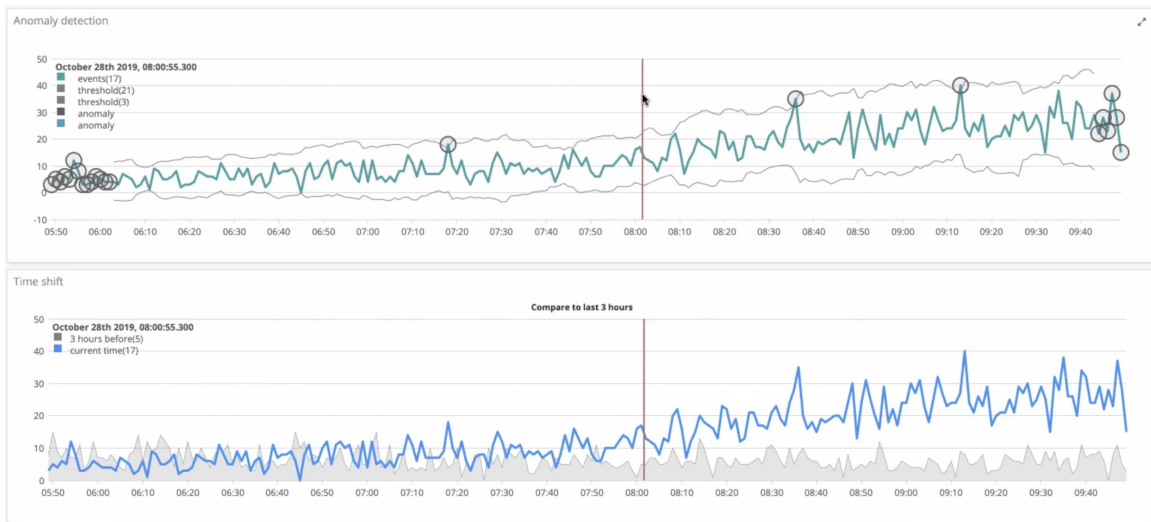# Soda integrates natively with data catalogs

Whether you use Alation, Colliba, Hive Metastore, or Amundsen, we natively integrate into your catalog to help provide operational context (real-time intelligence) to your datasets. This way, users can truly trust the data they use or produce.

# Soda uses ML based anomaly detection for DQ

Managing rules on data is unwieldy at scale. This is where Machine Learning comes in. Soda look at outliers in your data delivery process (when does data arrive, how does it look, ...).

# Architecture

SODA

01

# Soda Compute Engine

Big data engine calculates a wide variety of base metrics and allows business users to define & measure custom metrics that will be evaluated as new data comes in  (and where possible also historically).

E.g. "Percentage of valid values"
 = (countNulls(ds1.c1) + customValueCount("N/A") + outOfDomainValueCount(ds1.c1, ds2.c1)) / rowCount(ds1)

Soda supports the creation of metrics on filtered datasets as well.

**Soda Solution Architecture**

| Visualization, Annotation & Search |
| Alerting & Workflow<br>task assignment based on simple responsibility graph |
| Data Event System<br>using pub/sub & provides wide variety of event types |
| Monitoring & Timeseries Framework<br>e.g. Prometheus (incl. anomaly detection) |
| Soda DQ Engine<br>metrics calculation, schema evolution, ... |

Cloud Data Lake & DWH
(raw, batch, streaming, ...)

DQ-as-code
in Scala, Python, ...

## 02

# Monitoring & Timeseries Framework

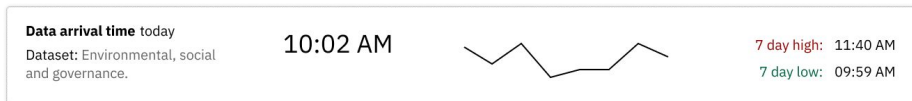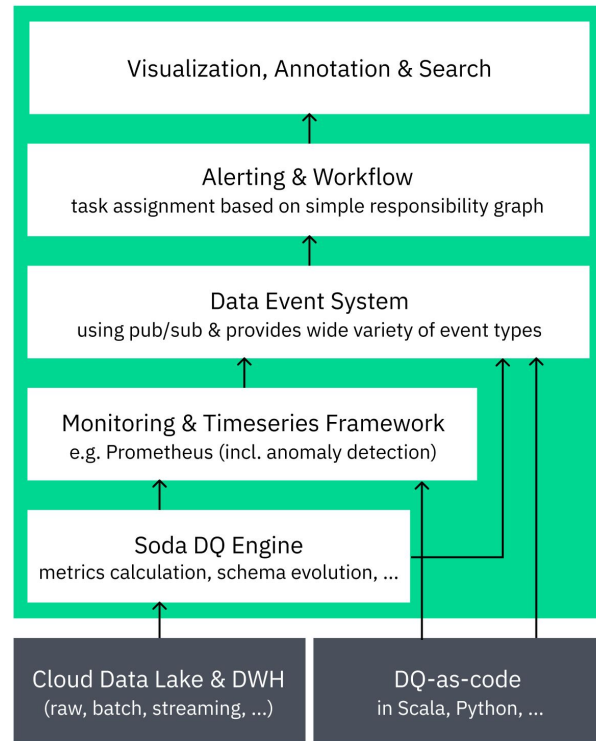Leverages open-source metrics framework like Prometheus & provides anomaly detection tools. Externally managed metrics can also be pushed in this framework (data quality as code). Add DQ specific functionality to show e.g. backfills, patched data, ...

**Data arrival time** today
Dataset: Environmental, social and governance.

10:02 AM

7 day high: 11:40 AM
7 day low: 09:59 AM

## Soda Solution Architecture

Visualization, Annotation & Search

Alerting & Workflow
task assignment based on simple responsibility graph

Data Event System
using pub/sub & provides wide variety of event types

Monitoring & Timeseries Framework
e.g. Prometheus (incl. anomaly detection)

Soda DQ Engine
metrics calculation, schema evolution, ...

Cloud Data Lake & DWH
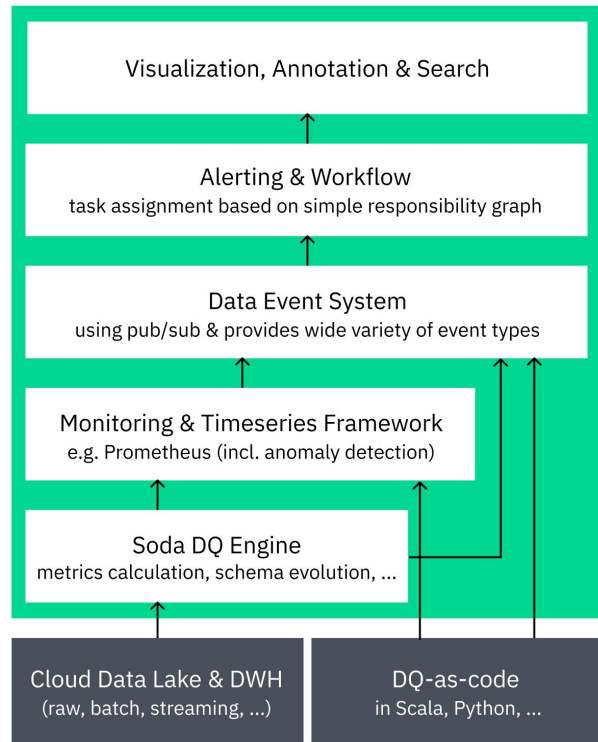(raw, batch, streaming, ...)

DQ-as-code
in Scala, Python, ...

# 03

# Data Event System

Events are created from (i) metrics that have a measurement that exceeds a thresholds, (ii) schema service or (iii) external code. Data consumers can subscribe to datasets to receive alerts. Events are displayed in an event stream that can be filtered and aggregated at different levels (e.g. dataset, user, ...).

## Soda Solution Architecture

Visualization, Annotation & Search

Alerting & Workflow
task assignment based on simple responsibility graph

Data Event System
using pub/sub & provides wide variety of event types

Monitoring & Timeseries Framework
e.g. Prometheus (incl. anomaly detection)

Soda DQ Engine
metrics calculation, schema evolution, ...

Cloud Data Lake & DWH
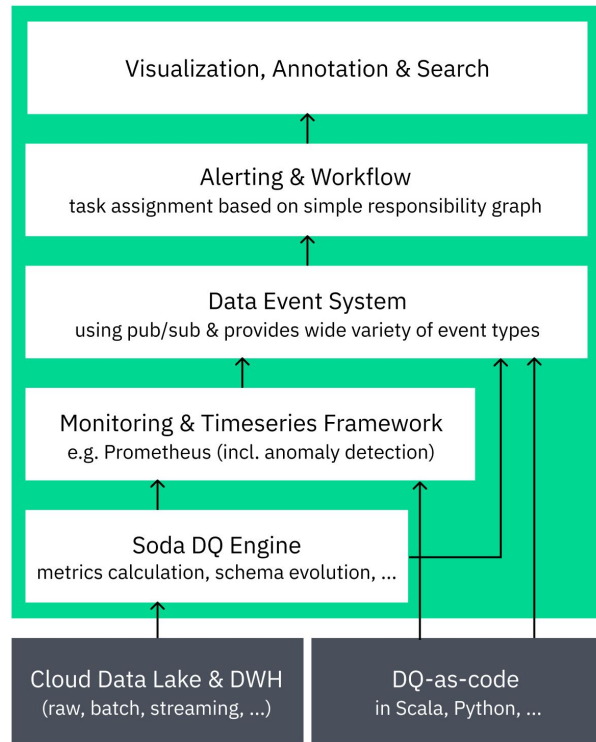(raw, batch, streaming, ...)

DQ-as-code
in Scala, Python, ...

04

# Alerting & Workflow

Provide, easy yet flexible notification & task mgmt.



Traverse the graph to add the users you want

**Soda Solution Architecture**



Visualization, Annotation & Search

Alerting & Workflow
task assignment based on simple responsibility graph

Data Event System
using pub/sub & provides wide variety of event types

Monitoring & Timeseries Framework
e.g. Prometheus (incl. anomaly detection)

Soda DQ Engine
metrics calculation, schema evolution, ...

Cloud Data Lake & DWH
(raw, batch, streaming, ...)

DQ-as-code
in Scala, Python, ...
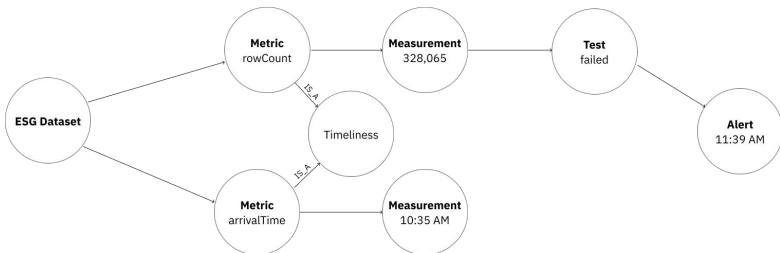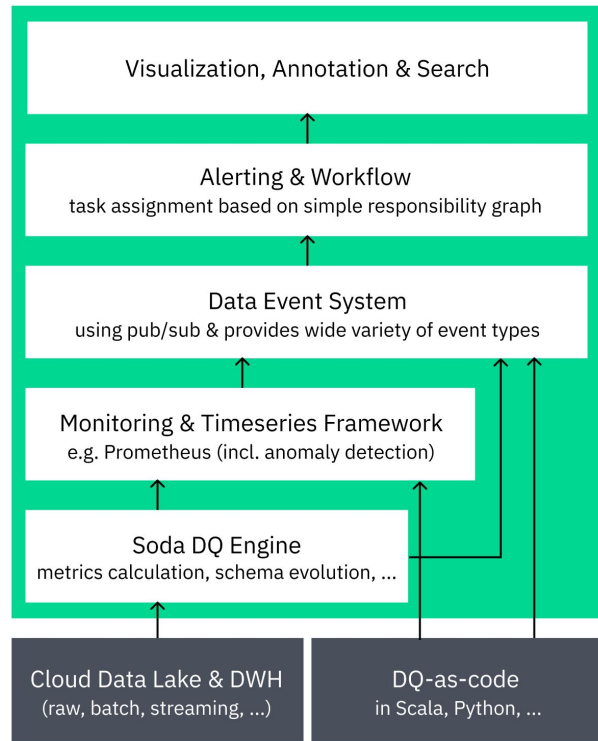
# 05

# Visualization, Annotation & Search

Search & correlate events to find the root cause. Provide out-of-the-box aggregation dashboards per persona or connect to third party tools (Tableau, ...).



**Soda Solution Architecture**

# Thank you.

SODA

Maarten Masschelein
Founder & CEO
(929) 920-1414
maarten@sodadata.io