# CLOUDERA NAVIGATOR

Philip Morch              Account Executive

Johannes Muselaers        Solution Engineer

Philippe Lanckvrind       Solution Engineer

# WHAT MAKES BIG DATA GOVERNANCE DIFFERENT?

Governing big data requires governing petabytes of diverse types of data

No one application will solve every big data governance problem

Applications are shifting to the cloud, and data governance must still be applied consistently

Self-service data discovery is mandatory for big data

# DATA MANAGEMENT IS THE FOUNDATION OF ADOPTION

## Governance & Compliance
Track, understand and protect access to data

- Am I prepared for an audit?
- Who's accessing sensitive data?
- What are they doing with the data?
- Is sensitive data governed and protected?

## Curation
Manage and organize data assets at Hadoop scale

- How can I identify and classify sensitive data in accordance with regulation?
- How can I organize and classify data for business users?
- How can I efficiently make data available to business users?

## Self-Service Discovery
Effortlessly find and trust the data that matters most

- How can I find explore data sets on my own?
- Can I trust what I find?
- How do I use what I find?
- How do I find and use related data sets?

## Stewardship
Boost user productivity and cluster performance

- How can I efficiently manage data lifecycle, from ingest to purge?
- How can I optimize my data models to support common access patterns?
- How can I migrate workloads to Hadoop risk-free?

### BIG DATA GOVERNANCE FOUNDATION

- Centralized audits
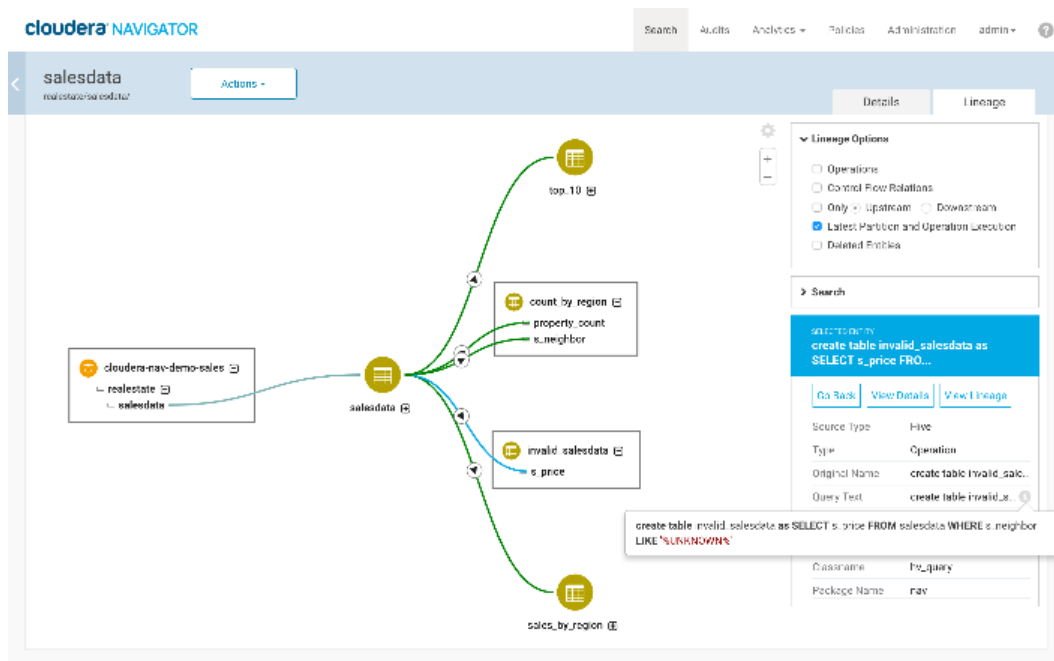- Unified data catalog
- Comprehensive lineage
- Data policies

**CLOUDERA**

# THE CURRENT STATE OF BIG DATA GOVERNANCE

**1**

**Governance & Compliance:** Raw governance artifact capture

**2**

**Curation:** Classification and tagging of data sets for compliance and discovery

**3**

**Self-Service Discovery:** Cataloging of data sets for end user self-service

**4**

**Stewardship:** Automation of lifecycle management, from ingest to purge

**5**

**Optimization & Refactoring:** Iterative, continuous improvement

**These are the most pressing big data governance challenges today**

# CLOUDERA NAVIGATOR

Governance, stewardship, and discovery for big data built on machine learning and analytics



**Trusted for production**
- Deployed by hundreds of customers across multiple industries
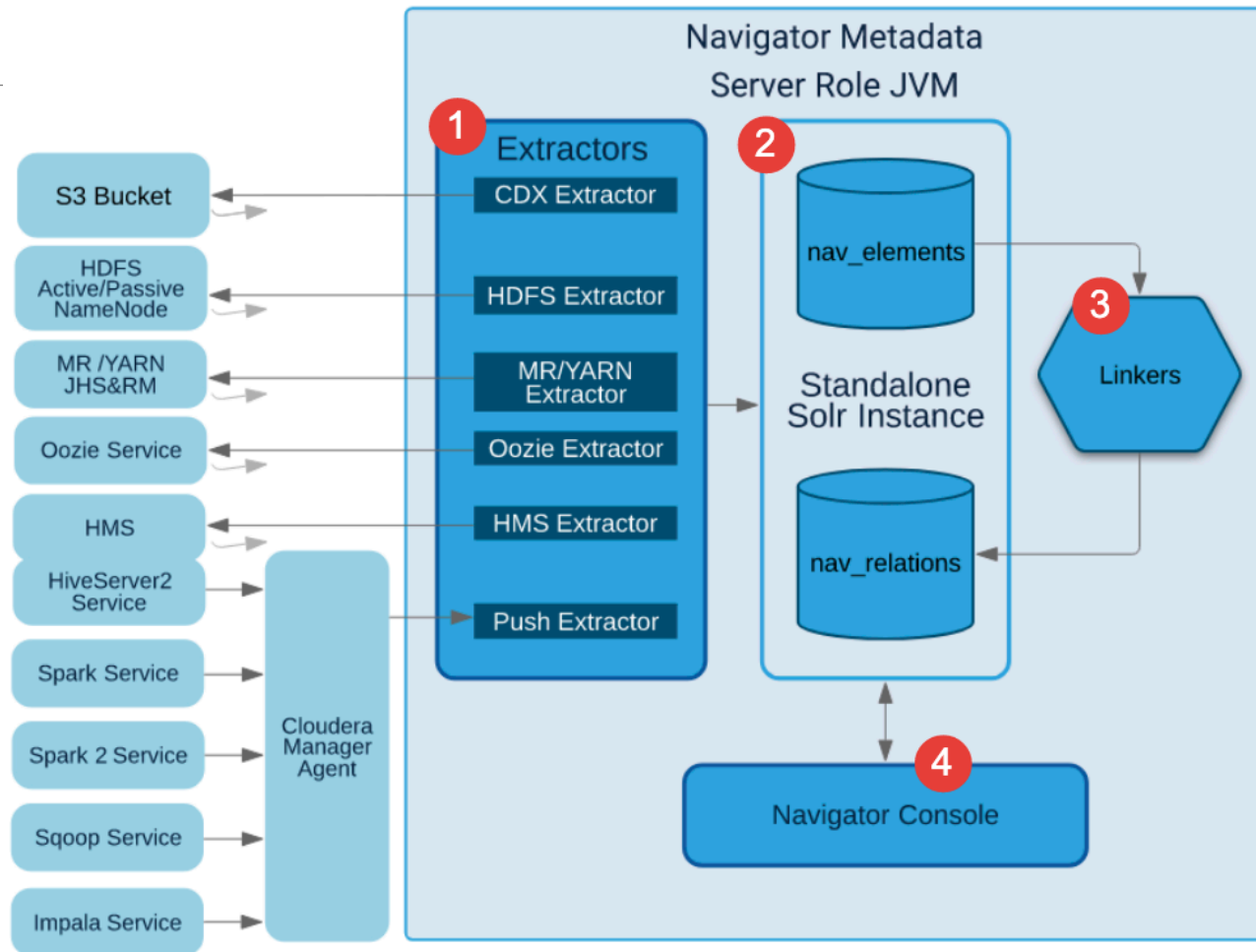- Over four years in production

**Compliance-grade**
- The only Hadoop distribution to pass PCI audit

**Open and interoperable**
- Integrated with leading partner solutions

# TURNKEY COMPLIANCE-GRADE **GOVERNANCE**

- Automatically collect column-level lineage and audit logs

- Effortlessly publish to enterprise governance frameworks, such as Informatica and IBM InfoSphere

# AUTOMATED DATA **CURATION**

- Set up policies that classify data sets automatically upon ingest
- Add business glossary definitions and profiling information for faster self-service
- Automatically trigger data preparation, profiling, and data quality activities

CLOUDERA

# OUT-OF-THE-BOX SELF-SERVICE **DISCOVERY**

- Let business users collaborate on and classify data sets while leveraging centrally-curated classifications

- Leverage deep analytics on historical usage to empower users to find, trust, and use data sets on their own

# STEWARDSHIP
BUILT ON MACHINE LEARNING AND ANALYTICS

At-a-glance stewardship metrics

CLOUDERA

Cloudera Navigator
Walk Through

# Use Cases & Best Practices

# USE CASES: COMPLIANCE

**Compliance**

Track, understand and protect access to data

- Am I prepared for an audit?
- Who's accessing sensitive data?
- What are they doing with the data?
- Is sensitive data governed and protected?

**ENTERPRISE METADATA REPOSITORY**

informatica   Data Advantage Group   IBM   *adaptive*

**ENTERPRISE AUDITING & SECURITY**

splunk>   IMPERVA   IBM   RSA SECURITY

**HADOOP DATA GOVERNANCE & MANAGEMENT**

Unified metadata   Unified lineage   Unified auditing

**Common use cases:**
- Security breach detection
- Data access tracking for PCI compliance
- Audit defense

CLOUDERA

# USE CASES: STEWARDSHIP & CURATION

**Stewardship, Curation & Discovery**

Manage, classify, and use data assets at Hadoop scale

How can I efficiently manage data lifecycle, from ingest to purge?

How can I identify and classify sensitive data in accordance with regulation?

How can end users find, trust, and use data sets on their own?

**Define Business Metrics & Glossary**

**Deliver Visualizations, Analytics, Reporting Across Systems**

**Ingest & Prepare: Landing Area**

**Profiling, Collaboration and Tagging**

**Clean, Transform, Refine Data**

**HADOOP DATA GOVERNANCE & MANAGEMENT**

# USE CASES: STEWARDSHIP & CURATION
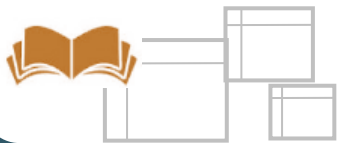
## Stewardship, Curation & Discovery

Manage, classify, and use data assets at Hadoop scale

- How can I efficiently manage data lifecycle, from ingest to purge?

- How can I identify and classify sensitive data in accordance with regulation?

- How can end users find, trust, and use data sets on their own?

### Define Business Metrics & Glossary
informatica
Data Advantage Group
Collibra
adaptive
IBM

### Deliver Visualizations, Analytics, Reporting Across Systems
tableau
SAP BusinessObjects
ZOOMDATA
platfora
SAS

### Ingest & Prepare: Landing Area
informatica
pentaho
syncsort
IBM
TRIFACTA
talend*
Paxata

### Profiling, Collaboration and Tagging
informatica
talend*
WATERLINE DATA SCIENCE
pentaho
TRIFACTA
IBM
Datameer
Paxata

### Clean, Transform, Refine Data
informatica
talend*
TRILLIUM SOFTWARE
pentaho
TRIFACTA
Datameer
platfora
IBM

**HADOOP DATA GOVERNANCE & MANAGEMENT**

# CLOUDERA NAVIGATOR'S VAST PARTNER ECOSYSTEM

Application

Platform

**collibra**

**podium data**

**informatica**

**talend**

Project management
Policy management
RACI
Stewardship workflows
ETL
Centralized curation
Centralized glossaries

Data quality
Uniqueness
Data valuation
Data profiling
Content enrichment

**Paxata**

**Waterline Data**

**TRIFACTA**

**ARCADIA DATA**

**Alation**

**tableau**

**Qlik Q**

Data wrangling
Data visualization
Query recommendations

**PRIVITAR**

**DATAGUISE**

Security profiling
Compliance: BCBS239,
GDPR

Unified technical metadata catalog
Extensible business metadata and glossary
Metadata policy engine
Comprehensive lineage
Unified audit/access logs
Dashboards and analytics
APIs for augmentation and consumption

End user collaboration
Crowdsourced metadata

**cloudera® NAVIGATOR**

Enterprise aggregation: metadata, lineage, SIEM, auditing

**informatica**

**Ab InITIO**

**IBM**

Centralized Stewardship

End User Discovery

**CLOUDERA**

15

# CURATION: ALIGN WITH KEY LIFECYCLE STAGES

## 1. Ingest Raw Data

## 2. Wrangle/Prepare Data

## 3. Publish Data

# ALIGN DATA CURATION WITH KEY LIFECYCLE STAGES

## 1. Ingest Raw Data

- Add source info (e.g., DB URL)
- Add retention information (e.g., retain for seven years)
- Add basic classifications (department, etc.)

## 2. Wrangle/Prepare Data

- Identify and classify sensitive data (e.g., PII, PHI)
- Add business glossary definitions
- Standardize field names (e.g., Zip and Zipcode)
- Integrate with DQ and profiling tools

## 3. Publish Data

- Collaborative metadata
- Crowdsourced metadata

# MANAGED METADATA VS CUSTOM METADATA

| | Managed Metadata | Custom Metadata |
|---|---|---|
| Intended usage | Centrally-curated metadata<br>Sentry ABAC<br>Metadata ABAC | End-user collaboration<br>Data set sharing |
| Assigned to specific entities (e.g., columns) | ✔ | |
| Typed and Validated (e.g., Boolean, Date, Enumeration) | ✔ | |
| Editable by data curators | ✔ | ✔ |
| Editable by end users | | ✔ |
| Viewable by data curators | ✔ | ✔ |
| Viewable by end users | ✔ | ✔ |

# Looking Forward
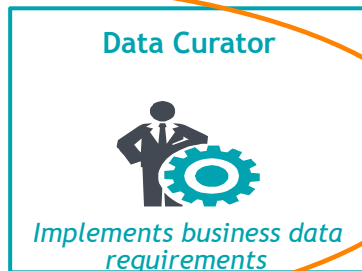
# MERGING OF TECHNOLOGIES

Bringing new capabilities

Authorization

Governance

Data fabric

# DATA GOVERNANCE: IT'S A TEAM SPORT!

| Sponsor | Data Owner | Data Council |
|---------|-----------|--------------|
| *Champions data governance across enterprise* | *Accountable for all data generated by an agency* | *Coordinate cross-agency data management activities* |

| Data Steward | Data Curator | Business Data SME |
|--------------|-------------|-------------------|
| *Manages business requirements for data sharing* | *Implements business data requirements* | *Supports the Data Steward in data related activities* |

**CLOUDERA**

# Discover With Data Steward Studio

**DETECT**
Find where important data assets are located

**REPORT**
Create and view multiple dashboards, reports, and summarizations of data

**INVENTORY**
Locate and catalog all data globally

**ENRICH**
Add classifications and annotations

**SECURE**
Protect data assets and monitor access and usage

**VERIFY**
Understand sources and complete chain of custody for all data (lineage and impact)

**COLLABORATE**
Crowdsource and leverage knowledge across the enterprise

**ORGANIZE**
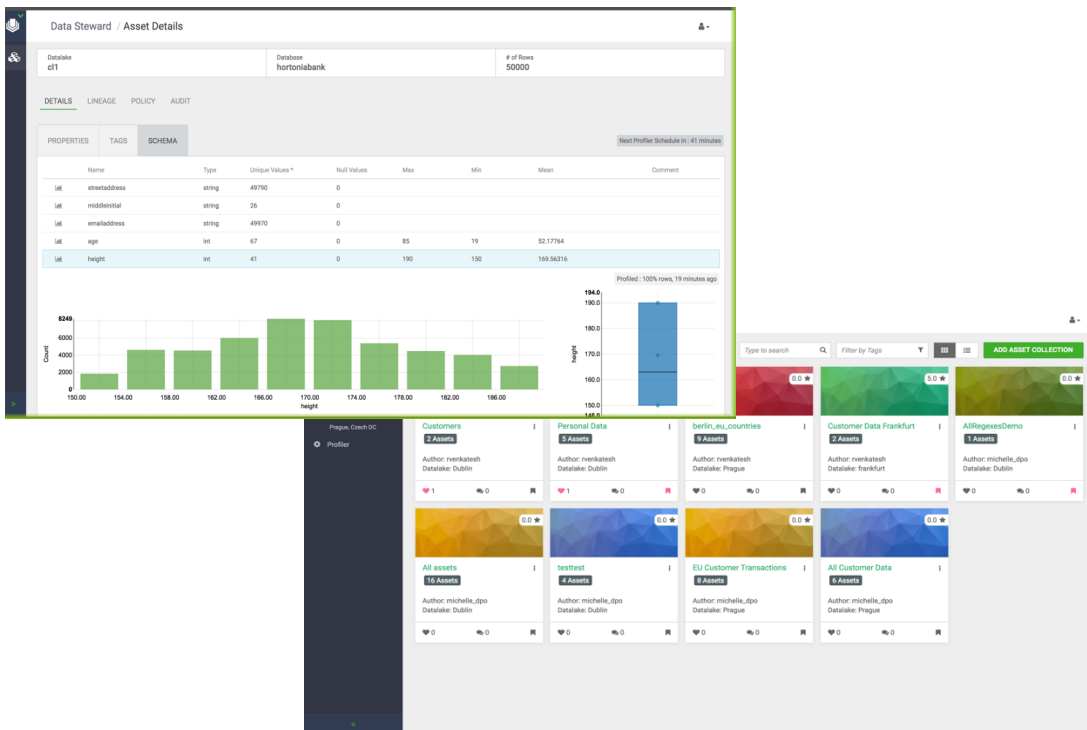Curate and group data based on different characteristics

CLOUDERA

# DATA STEWARD STUDIO

DSS provides the "tooling" part of the People, Processes, and Technology required for Hybrid Data Lake Governance



- **Profile Data** for understanding shape and structure
- Organize and **curate data** for e.g. by domains they belong to or data usage
- Identify **sensitive** data
- **Collaborate** with broader teams on how data needs to be used and by who and provide **community ratings** for crowdsourcing knowledge
- **Monitor** ongoing usage, **visualize** chain of custody and trustworthiness for longer term use, understand data protection

# THANK YOU

**CLOUDERA**