**CLOUDERA**

# Cloudera Navigator
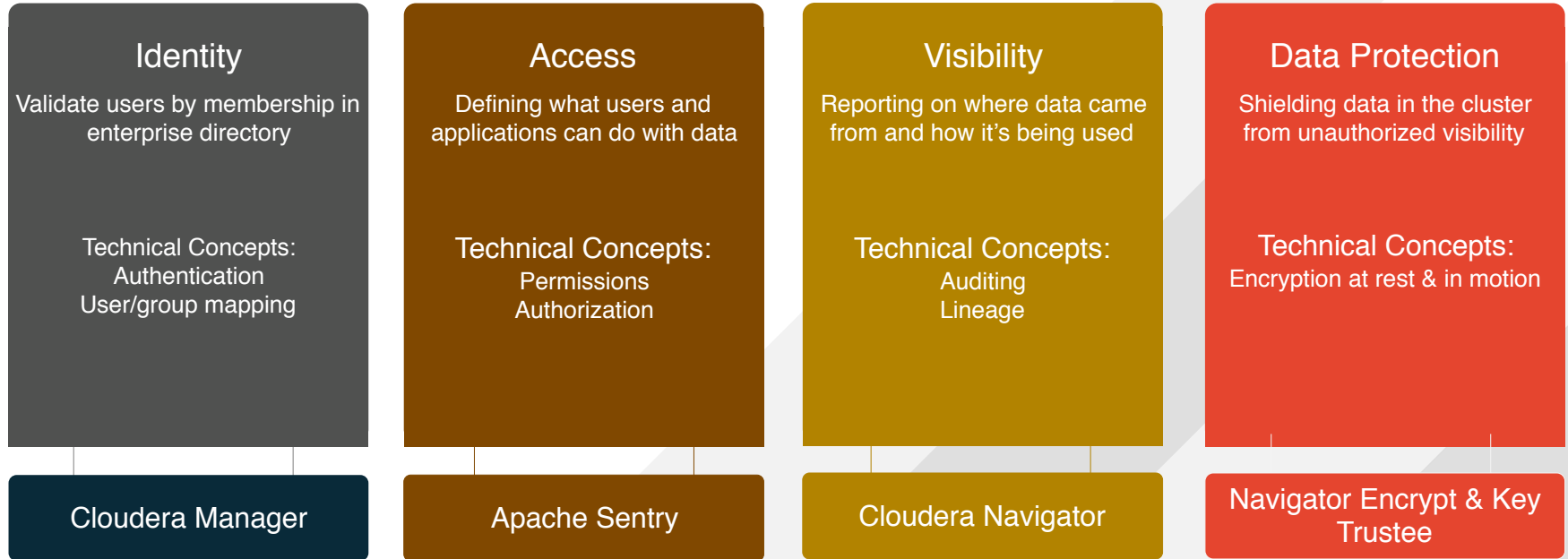
Enterprise Data Governance

Johannes Muselaers
Solutions Engineer Nordics
jmuselaers@cloudera.com
+46 72 588 1091

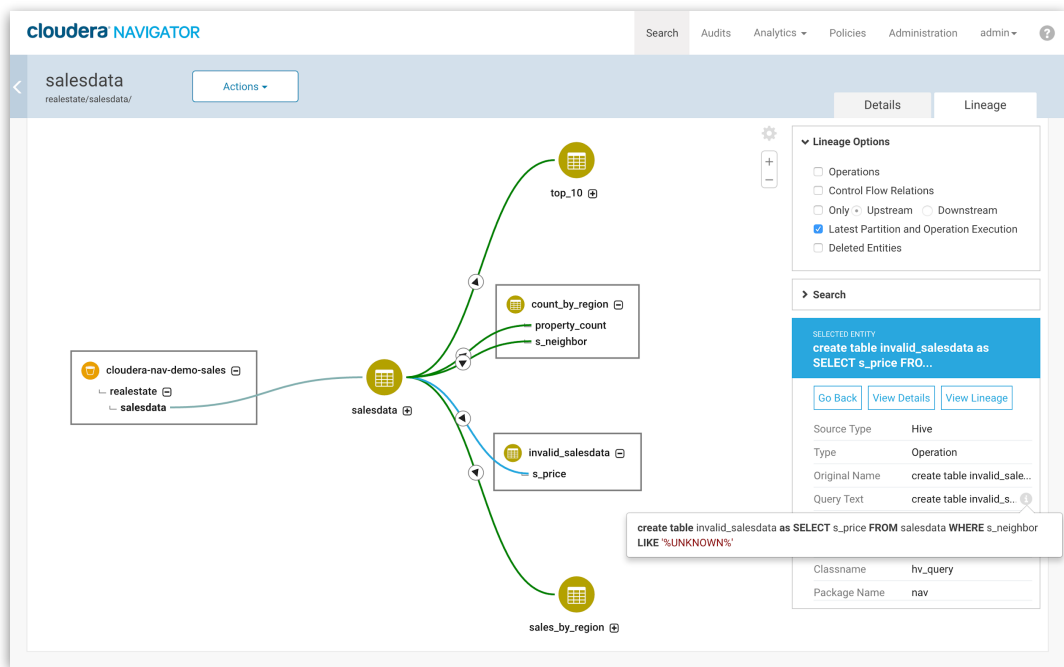# PDF of this presentation:

https://github.com/jmuselaers/Presentations/Cloudera Navigator.pdf

# Cloudera Enterprise-Grade Security and Governance

| Identity | Access | Visibility | Data Protection |
|---|---|---|---|
| Validate users by membership in enterprise directory | Defining what users and applications can do with data | Reporting on where data came from and how it's being used | Shielding data in the cluster from unauthorized visibility |
| Technical Concepts: Authentication User/group mapping | Technical Concepts: Permissions Authorization | Technical Concepts: Auditing Lineage | Technical Concepts: Encryption at rest & in motion |
| Cloudera Manager | Apache Sentry | Cloudera Navigator | Navigator Encrypt & Key Trustee |

# 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 five years in production

**Compliance-grade**

- The only Hadoop distribution to pass PCI audit

**Open and interoperable**

- Integrated with leading partner solutions

CLOUDERA

# Cloudera Navigator

- **Audit & Access Control**
  - Maintain full audit history across Hue, Cloudera Manager, HDFS, Impala, HBase & Solr
  - Ensuring appropriate permissions and reporting on data access for compliance
- **Discovery & Exploration**
  - Finding out what data is available and what it looks like
  - Tag datasets with metadata for pooling, knowledge sharing and common definitions
- **Lineage**
  - Tracing data back to its original source
  - Download and consume lineage data
- **Policy Engine**
  - automated classification on arrival
- **SDK** for integration with 3rd party tooling

# Auditing



- Who (IP), What (operation), When (timestamp) across the stack
- Any HDFS file access
- Denied logins to Hue, CM, Navigator
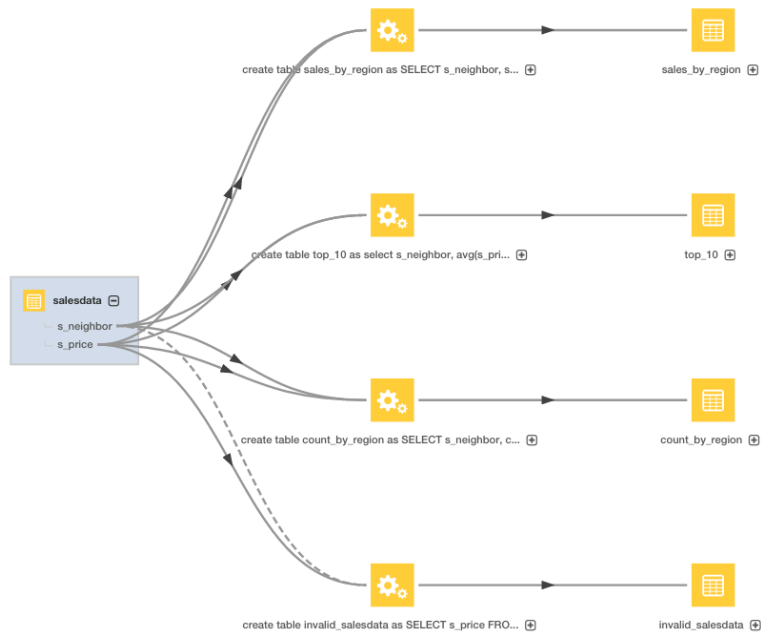- Query text from Impala, Hive, HBase, Solr, Sentry

# Lineage

# 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

## Archive files older than 7 years

| | |
|---|---|
| Status: | ✔ Enabled |
| Search Query: | (sourceType:hdfs) AND (created:[* TO NOW-7YEAR]) AND -tags:archive |
| Policy Description: | |
| Last Run On: | Tuesday, July 14th 2015, 3:08 pm |
| Last Modified: | Friday, April 13th 2018, 10:53 pm |
| Last Modified By: | admin |
| Schedule: | Recurring |
| | Start Time: | Jun 30, 2017 5:00 AM |
| | End Time: | Dec 31, 2099 5:00 AM |
| | Interval: | 1 Week(s) |

**Metadata Assignments**

| | |
|---|---|
| Name: | |
| Description: | |
| Tags: | archive |
| Key-Value Pairs: | |

**Command Action**

| | |
|---|---|
| Type: | Move |
| Target Path: | /archive/ |

**JMS Notifications**

| | |
|---|---|
| Notification: | archive_file |
| Queue: | hdfs_archive_queue |

## Autoclassify incoming sales data

| | |
|---|---|
| Status: | ✔ Enabled |
| Search Query: | (sourceType:hive) AND (originalName:salesdata) AND (type:Table) |
| Policy Description: | |
| Last Run On: | Thursday, May 24th 2018, 10:10 am |
| Last Modified: | Saturday, May 26th 2018, 12:15 pm |
| Last Modified By: | admin |
| Schedule: | On Change |

**Metadata Assignments**

| | | |
|---|---|---|
| Name: | | |
| Description: | This is the real estate sales history from March 2014 in New York. Obtained by monthly MLS feed. | |
| Managed Metadata: | **Stewardship** | |
| | Steward | md@cloudera.com |
| | **Classification** | |
| | Department | Sales |
| | Keywords | nycmls |
| | PII | true |
| Tags: | mls,  realestate,  salesdata | |
| Key-Value Pairs: | month:  March | |
| | retain_until:2018-03-15 | |
| | source:  nycmls | |

cloudera NAVIGATOR

Search    Audits    Analytics    Policies    Administration    admin

# customers_kudu

/user/hive/warehouse/customers_kudu

Actions

Details    Lineage

View in Hue

Edit Metadata...

Move...

Move to Trash...

Directory
HDFS

Owner: impala
Parent Directory: ware
Group: hive
Permissions: rwxrwxr

1
**Inputs**

0
**Outputs**

May 2, 2018 4:13 PM
**Modified**

## Technical Metadata

| | |
|---|---|
| Source Type | HDFS |
| Parent Directory | warehouse |
| Path | /user/hive/warehouse/customers_kudu |
| Owner | impala |
| Group | hive |
| Permissions | rwxrwxrwt |
| Last Modified | May 2, 2018 4:13 PM |
| Created | May 2, 2018 4:13 PM |
| Source | HDFS-1 |
| Classname | HDFS Entity |
| Package Name | nav |

▼ Directory Contents (0)

No matches found

❯ Inputs (1)

❯ Outputs (0)

# Edit Metadata

**Name**

customers_kudu

**Description**

## Managed Metadata

No Managed Metadata properties available.

## Custom Metadata

**Tags**

sensitive ×

**Key-Value Pairs**

+

Cancel   Save

**cloudera** NAVIGATOR

Search    Audits    Analytics ▾    Policies    Administration    admin ▾    ❓

# Search

🔍 sensitive ⟵                                              ✕          Actions ▾

**Filters**          Add Filters    Clear All Filters

**5** results                                          Show full query

∨ SOURCE TYPE

☐ ● Hive          4        🗄 Hive **s_neighbor**                                      View in Hue
                              Data Type **string**   Parent Path **/default/top_10**   Source **HIVE-1**
☐ ● HDFS          1

                              🗄 Hive **s_neighbor**                                      View in Hue
∨ TYPE                         Data Type **string**   Parent Path **/default/salesdata**   Source **HIVE-1**

☐ Field          4        🗄 Hive **s_neighbor**                                      View in Hue
                              Data Type **string**   Parent Path **/default/sales_by_region**   Source **HIVE-1**
☐ Directory      1

                              🗄 Hive **s_neighbor**                                      View in Hue
∨ OWNER                        Data Type **string**   Parent Path **/default/count_by_region**   Source **HIVE-1**

☐ impala         1
                              🗂 HDFS **customers_kudu**                                  View in Hue
Add New Value                  Path **/user/hive/warehouse/customers_kudu**   Owner **impala**   Group **hive**   Permissions **rwxrwxrwt**
                              Last Modified **May 2, 2018 4:13 PM**   Created **May 2, 2018 4:13 PM**   Source **HDFS-1**

∨ CLUSTER TEMPLATE

☐ Cluster 1      1

Add New Value

# 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
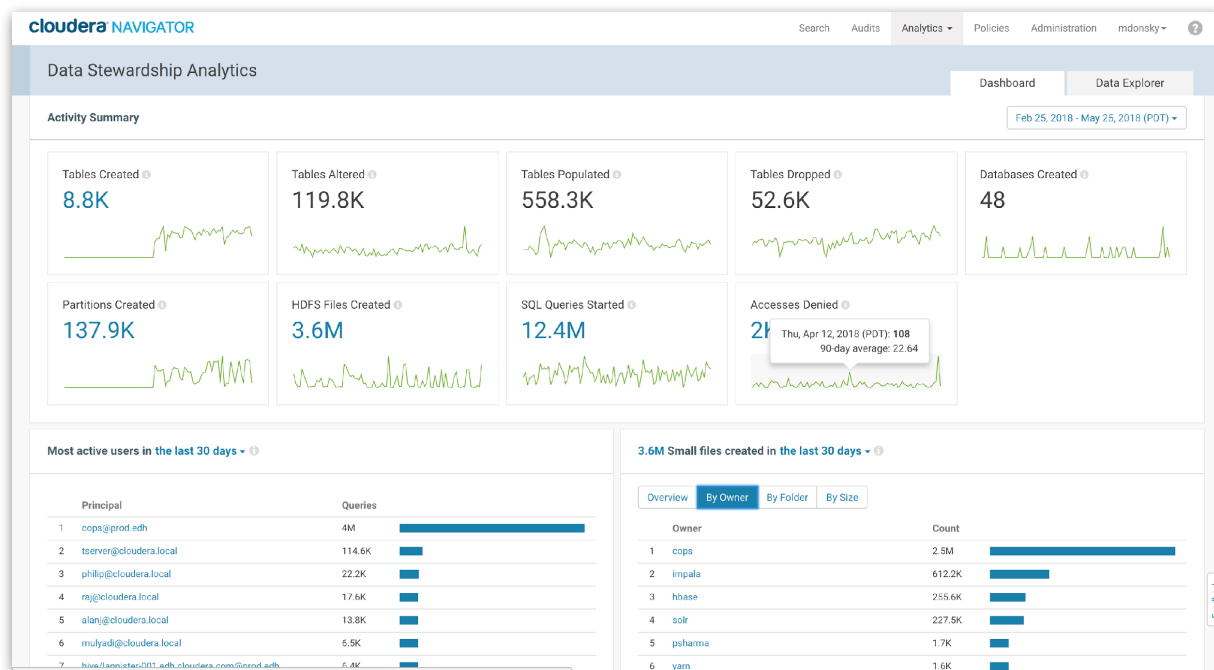
# STEWARDSHIP
## BUILT ON MACHINE LEARNING AND ANALYTICS

At-a-glance stewardship metrics

# USE CASES I SELF-SERVICE DISCOVERY

## Find, trust, and use data sets

- **Find** data sets based on business context

- **Trust** data sets with insight based on ironclad governance

- **Use** data sets with analytics on on historical usage patterns

# Navigator Encrypt

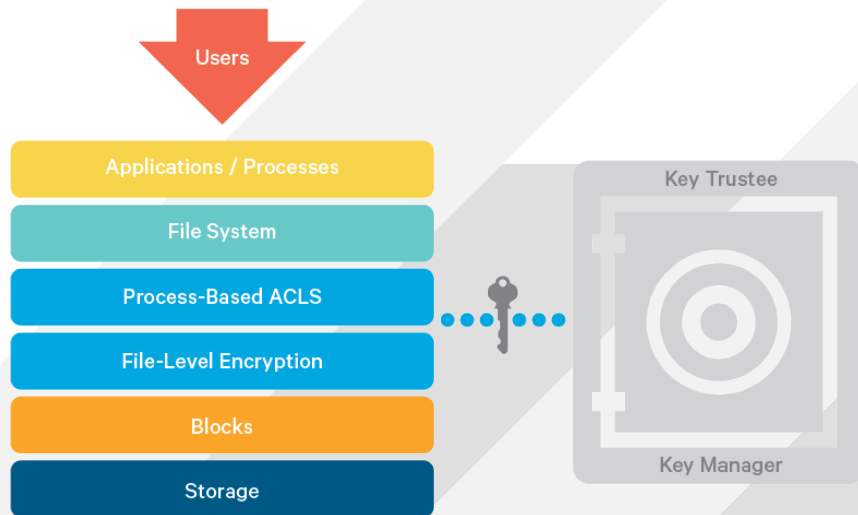Transparent layer between application and file system

Compliance-Ready

Massively Scalable

High Performance: Optimized for Intel

Separation of Duties

Key Management with Navigator Key Trustee

# Spark limitations

## Spark Lineage Limitations and Requirements

Spark lineage diagrams are supported in the following releases:

- Spark 2.3 as of Cloudera Manager 5.14.0 (Navigator 2.13.0)

- Spark 1.6 as of Cloudera Manager 5.11.0 (Navigator 2.10.0)

Lineage is not available for Spark when Cloudera Manager is running in single user mode. In addition to these requirements, Spark lineage has the following limitations:

- Lineage is produced only for data that is read/written and processed using the Dataframe and SparkSQL APIs. Lineage is not available for data that is read/written or processed using Spark's RDD APIs.

- Lineage information is not produced for calls to aggregation functions such as `groupBy()`.

- The default lineage directory for Spark on Yarn is `/var/log/spark/lineage`. No process or user should write files to this directory—doing so can cause agent failures. In addition, changing the Spark on Yarn lineage directory has no effect: the default remains `/var/log/spark/lineage`.

# Latest releases

https://www.cloudera.com/documentation/enterprise/release-notes/topics/cn_rn_new_features.html

# Roadmap

# THANK YOU

**CLOUDERA**

[http://sdx-nav-demo-1.gce.cloudera.com:7187/login.html](http://sdx-nav-demo-1.gce.cloudera.com:7187/login.html)