# Comparison of 2 Data Catalogues Tools

presented by Lukas Heubach and Leonhard Krause
DWH-TINF18D

# Agenda

- What is a Data Catalog?
- Why Data Catalogues?
- IBM InfoSphere IGC
- Lumada Data Catalog
- Comparison
- Sources

# What is a Data Catalog?

- digital inventory (directory)
- contains all company data
- Single source of trust - data inventory that is correct and can be relied upon
- Data catalog is filled with metadata of technical and business origin
- Data supply & demand
- provides functions for registering, retrieving, using, evaluating and analyzing data
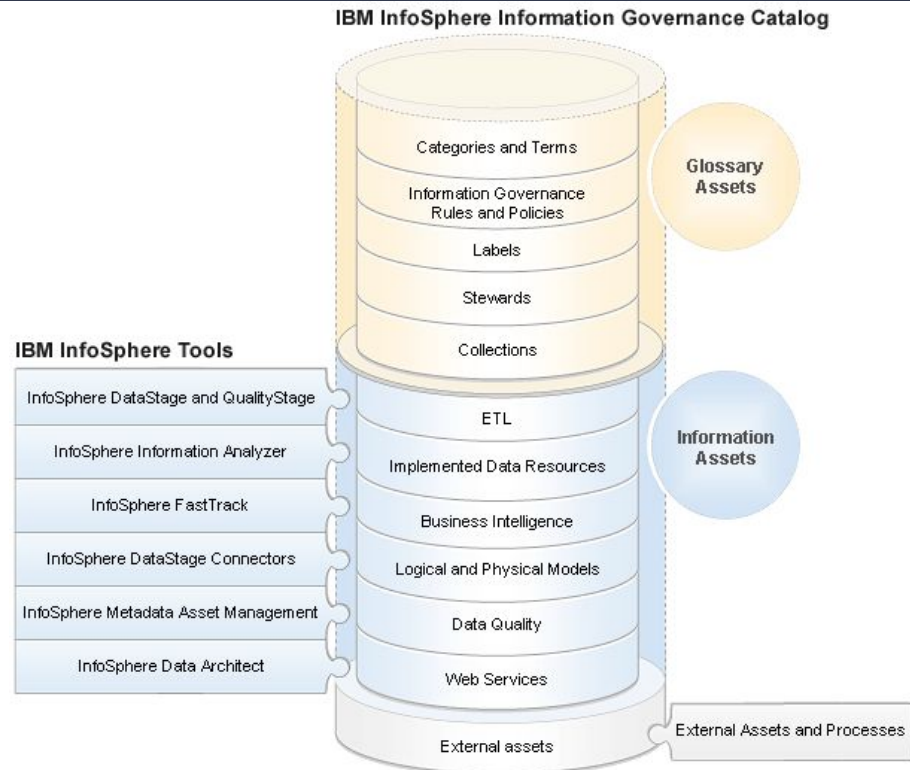
# Why Data Catalogues?

- Data is constantly accumulating, more and more and in new formats
- data sets should be transparently available in the company
- to organize company data
- Main objective: to promote collaboration within the company by making relevant data

# IBM InfoSphere IGC – General

InfoSphere Information Governance Catalog

- web-based tool
- Create, manage, share, use business knowledge
- Prices
-      individual offer
- Goals
  - Data -> reliable information
- Can be used in conjunction with other InfoSphere tools



**IBM InfoSphere Information Governance Catalog**

Glossary Assets
- Categories and Terms
- Information Governance Rules and Policies
- Labels
- Stewards
- Collections

**IBM InfoSphere Tools**
- InfoSphere DataStage and QualityStage
- InfoSphere Information Analyzer
- InfoSphere FastTrack
- InfoSphere DataStage Connectors
- InfoSphere Metadata Asset Management
- InfoSphere Data Architect

Information Assets
- ETL
- Implemented Data Resources
- Business Intelligence
- Logical and Physical Models
- Data Quality
- Web Services

External assets — External Assets and Processes

# IBM InfoSphere IGC – Functions

**Connection of data sources**

- different types of sources (asset types)
    - (AWS S3, IBM InfoSphere DB2, Oracle ...)
- Import e.g. via Metadata Asset Manager

# IBM InfoSphere IGC – Functions

**Glossary Assets**

- create/represent complex relationships between assets
- Categories
  - like a folder to structure Glossary Assets
- Terms
  - Word/phrase that describes a characteristic
- IG Rules
  - a natural language description of a criterion that determines whether an information asset meets a business objective.
- IG Policies
  - a natural language description of a subject area
- Labels

# IBM InfoSphere IGC – Functions

**Information Assets**

- Imported records
- Imported metadata
- Display of all included data sources
- Display of all data sets/metadata
- Assign to Glossary Assets

# IBM InfoSphere IGC – Functions

**Queries**

- create your own queries
- for Information Assets
- for Glossary Assets
- result: Table with information

# IBM InfoSphere IGC – Who is it for?

- Business Analysts
- Business experts
- Organizations that
  - want to manage a common enterprise vocabulary and governance practices
  - want to leverage the potential of integrated metadata
  - reduce the need for technical training

# Lumada Data Catalog (Waterline Data Catalog)

- Tool for managing data from diverse sources
- uses machine learning to build data inventory
- patented fingerprinting technology for data
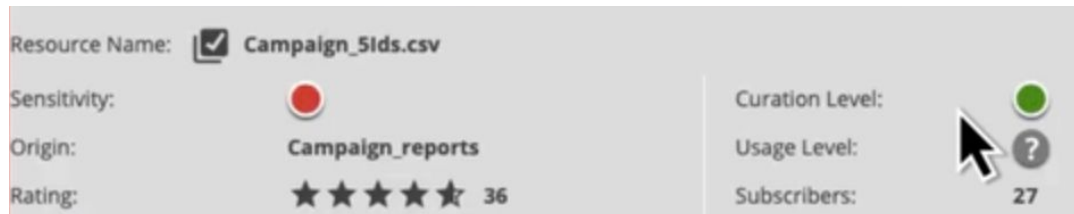- Management of data lakes
- Goal: analyze large amounts of data automatically
- Price and demo are only available on request

# Lumada Data Catalog – Functions

- Data recognition using metadata-based search
- Management of sensitive data
- Patented data fingerprinting
  - data processing based on machine learning
  - automatic recognition of sensitive data
  - enables Google-like search with corporate identifiers
- visualization of data, -origins and -relationships
- Detection of redundant data
- comment, subscription, rating function
- Recording of user data

| | | | | |
|---|---|---|---|---|
| Resource Name: | ☑ Campaign_5Ids.csv | | | |
| Sensitivity: | 🔴 | | Curation Level: | 🟢 |
| Origin: | Campaign_reports | | Usage Level: | ❓ |
| Rating: | ★★★★⯪ 36 | | Subscribers: | 27 |

# Lumada Data Catalog – Who is it for?

- Companies with large amounts of data
- Data analysts
- Fast and efficient compliance with data protection regulations
- Management of sensitive data

# Comparison

| Category | IBM InfoSphere IGC | Lumada Data Catalogue |
|---|---|---|
| **Similarities** | ● Data Catalog Tools<br>● individual price offer<br>● use of different data sources | |
| **User Interface** | Web Tool | no specification (possibly desktop application) |
| **Focus** | data management | data analysis |
| **AI** | No | AI-based data analysis and structuring |
| **Product Portfolio** | IBM InfoSphere Family | Lumada Data Services |
| **"Social-Media" Tools** | No | Yes (comments, subscription, rating) |

# Sources

- https://www.hitachivantara.com/de-de/products/data-management-analytics/lumada-data-services/lumada-data-catalog.html (retrieved on: 13.02.2021 at 14:32 Uhr)
- https://www.hitachivantara.com/en-us/products/data-management-analytics/lumada-data-catalog.html (retrieved on: 13.02.2021 at 14:32 Uhr)
- https://www.hitachivantara.com/en-us/pdfd/datasheet/lumada-data-lake-datasheet.pdf (abgerufen am: 13.02.2021 at 14:32 Uhr)
- https://www.hitachivantara.com/de-de/pdf/datasheet/lumada-data-catalog-datasheet-de.pdf (retrieved on: 13.02.2021 at 14:32 Uhr)
- https://www.talend.com/de/resources/what-is-data-catalog/(retrieved on: 13.02.2021 at 14:32 Uhr)
- https://www.cc-cdq.ch/data-catalogs (retrieved on: 13.02.2021 at 14:32 Uhr)
- https://www.ibm.com/de-de/marketplace/information-governance-catalog (retrieved on: 13.02.2021 at 14:32 Uhr)
- https://www.saracus.com/blog/der-ibm-infosphere-information-governance-catalog/ (retrieved on: 13.02.2021 at 14:32 Uhr)

# THANK YOU FOR YOUR ATTENTION!