# INTELLIGENT BUSINESS STRATEGIES

## *Accelerating Data Engineering Using a Data Catalogue*

Mike Ferguson
Managing Director
Intelligent Business Strategies
Quest for Knowledge
Amsterdam, June 2025

---

## About Intelligent Business Strategies

- A UK-based independent IT analyst and consulting firm founded 1992 specialising in data management and analytics
- Mike Ferguson is an independent IT Industry Analyst and consultant, Conference Chairman of Big Data LDN and a member of the EDM Council CDMC Executive Advisory Board
- Three main lines of business

### Research
- Market research
  - 4th Industrial Revolution Survey
- D&A product research
  - Data Catalogs
  - Data Governance
  - Data Fabric
  - Data Science Workbenches
  - Analytical Databases
  - The Agentic Enterprise

### Education
- Building a Data & AI Strategy for a Data Driven Enterprise
- Modern Data Architecture
- AI-Driven Active Data Governance
- Practical Guidelines for Implementing a Data Products
- Embedded Analytics, Intelligent Apps, AI Agents & AI Automation
- Data Warehouse Migration to the Cloud
- Data Warehouse Modernisation
- Public classes (anyone)
- On-site classes (single client)
  - Customers, vendors, systems integrators
- On-line (public & on-site classes)

### Consulting
- Customer consulting services
  - D&A Strategy, Data Architecture
  - D&A Technology selection
  - D&A Reviews, Data Governance
  - Project implementation advisory
- Vendor advisory services
  - Product strategy
  - Product positioning & go to market
  - Marketing support
    - Speaking at vendor events
    - White papers
    - Webinars
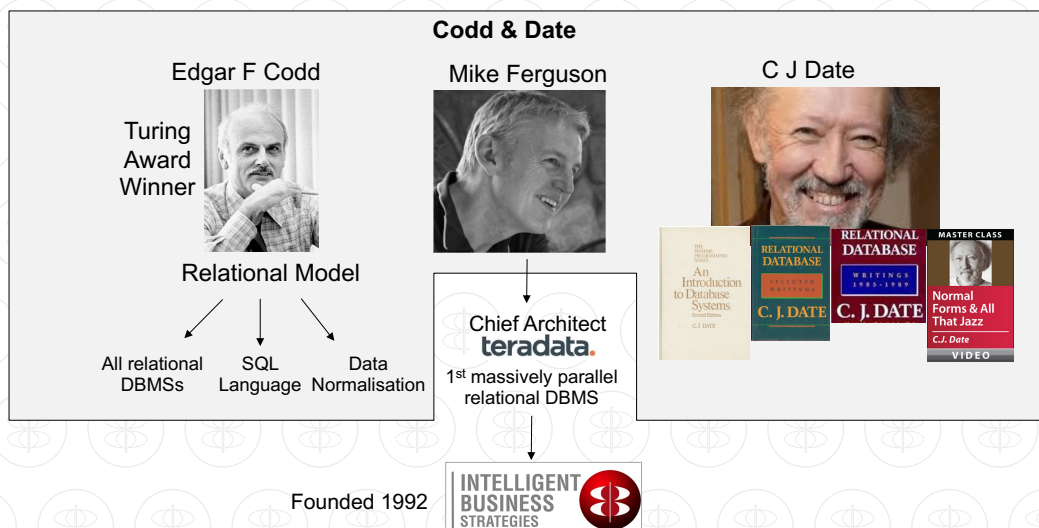- Venture Capitalists
  - Due-diligence, Asset advisory

www.intelligentbusiness.biz

2

## Who Is Mike Ferguson? – A Leading Analyst In Data Management & Analytics

**Codd & Date**

Edgar F Codd

Turing Award Winner

Relational Model

All relational DBMSs — SQL Language — Data Normalisation

Mike Ferguson

Chief Architect
**teradata.**
1st massively parallel relational DBMS

Founded 1992

**INTELLIGENT BUSINESS STRATEGIES**

C J Date

An Introduction to Database Systems, Second Edition — C.J. Date

RELATIONAL DATABASE SELECTED WRITINGS — C. J. DATE

RELATIONAL DATABASE WRITINGS 1985-1989 — C. J. DATE

MASTER CLASS — Normal Forms & All That Jazz — C.J. Date — VIDEO

3

## Mike Ferguson Is Europe's Leading Industry Analyst / Consultant In Data Management & Analytics And Conference Chairman Of Big Data LDN

Big Data LDN is the largest data & analytics conference in Europe
- 20000 delegates
- 200+ vendors
- 14 theatres
- 350+ speakers

It is 6 x size of Gartner's D&A conference

4

## Topics

- What is a data catalogue?
- The data catalogue marketplace
- Using a business glossary within a data catalogue to define data products
- Using a data catalogue to automatically discover data in multiple data sources
- Using a data catalogue to map raw data to common terms in a business glossary
- Using a data catalogue to automatically detect sensitive data in data sources
- Using a data catalogue to automatically profile data quality in your data sources and recommend fixes
- The power of metadata - Integrating data engineering tools and generative AI with a data catalogue to rapidly build data pipelines to produce data products
- Publishing data products in a data marketplace within the data catalog

5

## Topics

- ➢ What is a data catalogue?
- ➢ The data catalogue marketplace
- Using a business glossary within a data catalogue to define data products
- Using a data catalogue to automatically discover data in multiple data sources
- Using a data catalogue to map raw data to common terms in a business glossary
- Using a data catalogue to automatically detect sensitive data in data sources
- Using a data catalogue to automatically profile data quality in your data sources and recommend fixes
- The power of metadata - Integrating data engineering tools and generative AI with a data catalogue to rapidly build data pipelines to produce data products
- Publishing data products in a data marketplace within the data catalog

6

## What Is A Data Catalogue?

**Enterprise Data Catalogue**

Data catalogue software connects to data stores, IoT schema registries and tools in your organisation to automatically discover and classify the structured, semi-structured and unstructured data and data relationships that exist in your distributed data estate and also profile its quality.

It creates an inventory of all your data assets, labels data so you know how to govern it and discovers lineage to understand how data is processed.

A data catalogue is a metadata repository that contains a knowledge graph of your data and its relationships, to enable users to find and access information quickly and to enable you to create data governance policies in the catalogue to govern data in a consistently across your data estate.

7

## You Need A Data Catalog To Automatically Discover What Data Is Available, Its Quality, Sensitivity And Where It Is Across Your Data Estate to Govern It

Automatically discover, classify, data quality profile and catalog data

Enterprise Data Fabric Software

Data catalog

Scans can be scheduled to keep the catalog up to date

Automatic data discovery, classification & data quality profiling

| SaaS Applications | Cloud Based Applications | On-Premises Systems | Edge Devices |
|---|---|---|---|
| salesforce, workday, Marketo, ORACLE NETSUITE, slack, monday.com, Dropbox, Dynamics 365, Adobe Experience Cloud | OLTP systems, Analytical Systems, Files, Content — Multiple clouds (Azure, aws, Google Cloud) | Analytical Systems, OLTP Systems, Files, Content | IoT data (Sensor data) |

8

## Data Catalog Product Examples (A Very Crowded Market)

- Ab Initio data catalog
- Actian Zeenea
- Alation
- Alex Solutions
- Altair (Cambridge Semantics) Anzo Catalog
- Alteryx Connect
- Amazon Glue Catalog
- Apache Atlas (open source)
- Ataccama ONE Data Catalog
- Atlan Catalog
- BigID Data Catalog
- Boomi Atomsphere Data Platform Catalog
- Cloudera Data Platform SDX Catalog
- Collibra Catalog
- Databricks Unity Catalog
- Denodo Catalog
- Google Cloud Data Catalog
- Hitachi Vantara Pentaho Data Catalog

- IBM Knowledge Catalog
- Informatica IDMC Data Governance & Catalog
- Microsoft Purview
- Oracle Cloud Infrastructure Data Catalog
- Qlik
  - Enterprise Data Catalog
  - Qlik (Talend) Data Catalog
- Quest (formerly erwin) Data Catalog
- Rocket Software (formerly ASG) Intelligent Data Catalog
- SAP DataSphere Catalog (cloud)
- SAS Information Catalog
- Salesforce Tableau Catalog
- Servicenow data.world
- Snowflake Polaris
- TIBCO Cloud Metadata Catalog
- Top Quadrant TopBraid EDG Data Catalog
- Truist Zaloni Arena Data Catalog

9

## What Core Capabilities Should A Data Catalog Offer?

| Core Capability | Additional Comments |
|---|---|
| Business glossary management | • Create common business data names and definitions for common understanding |
| Automated data discovery | • On-premises, multiple clouds, SaaS Apps edge |
| Knowledge graph of your data estate | • Understand data relationships, who / what uses that data, dependencies & more |
| AI assisted search and faceted search to find data | • To help data producers quickly and easily find the data they need |
| Automated data governance classification | • Auto detection of sensitive data types, e.g., personal data, financial data<br>• Trainable classifiers to label data using user defined governance classification schemes<br>• Classification grouping, e.g., all personal data subject to GDPR |
| Automated data quality profiling | • Understand data anomalies, completeness |
| Automated extraction & derivation of lineage | • Understand how data has been processes, where it is used and who uses it |
| Data Governance of a distributed data estate | • Define policies based on data governance classifications & classification groups |
| Tag management | • Allow people to tag data and files for easy location and governance in the future |
| Data governance policy management | • Policies to govern data access security, privacy, retention, sharing and quality |
| Data usage behavioural observation | • Harvest database query logs, understand frequency of use, usage patterns, |
| Data marketplace | • Provide consumers with 'ready-made' , high quality data and govern data sharing |
| Lineage visualisation | • Understand data processing at a glance, understand reports & models using it |
| Integration with data fabric and 3rd party tools | • Access from 3rd party tools, launch of data preparation tools<br>• Governance of self-service data preparation jobs<br>• Governance of reports, dashboards, ML models…. |

10

## Data Catalogs Can Support Multiple Personas

- Personas
  - Chief data officers
  - Data stewards
  - Data engineers (data producers)
  - Data scientists
  - Business analysts
  - IT developers (including catalog access via APIs)
- Data catalogs also integrate with many other tools
  - Data science workbenches
  - ETL tools
  - BI tools
  - Data modelling tools
  - API management tools
  - Data automation tools
  - Data observability tools

11

## Topics – Where Are We?

- What is a data catalogue?
- The data catalogue marketplace
- ➢ Using a business glossary within a data catalogue to define data products
- Using a data catalogue to automatically discover data in multiple data sources
- Using a data catalogue to map raw data to common terms in a business glossary
- Using a data catalogue to automatically detect sensitive data in data sources
- Using a data catalogue to automatically profile data quality in your data sources and recommend fixes
- The power of metadata - Integrating data engineering tools and generative AI with a data catalogue to rapidly build data pipelines to produce data products
- Publishing data products in a data marketplace within the data catalog

12

## Creating Reusable Data Products Is An Approach Fast Gaining Momentum As A Way of Incrementally Building Up A Secure And Compliant Data Foundation



SaaS Applications
- salesforce
- workday
- Marketo
- ORACLE NETSUITE
- slack
- monday.com
- Dropbox
- Dynamics 365
- Adobe Experience Cloud

Cloud Based Applications

OLTP systems
Analytical Systems
Files
Content

Multiple clouds

Analytical Systems
OLTP Systems
Files
Content

Automatic data discovery, classification & data quality profiling

Data catalog

IoT data

High Quality, Compliant Data Products

Data Fabric

**13**

## Key Requirements – Organisational Implications If You Are To Become Data Driven – We Need Data Product Producers And Consumers



Data producers

Business domain

Embedded IT professional data architect

raw data

discover classify, profile

clean integrate & analyse

trusted data product

publish

data catalog

data catalog

Citizen Data engineers

Data marketplace

data catalog

Ready made data products

data consumers

business analysts / data scientists

search
find
shop
order

provision / consume

BI tool or notebook

- ▪ Need to make use of
  - • A business glossary and data catalog
  - • A collaborative approach to producing data and analytical products
  - • A catalog / data marketplace to quickly find reusable trusted data products to drive business value

**14**

Federated Operating Model – Co-ordinated Business Strategy Aligned Teams Producing Data Products Using Common Data Fabric and a Data Catalog



A Modern Data Architecture for Converged Analytical Workloads With Access to Shared Data Products in Stored in Lakehouse Open Tables That Are Available to Multiple Analytical Engines

A Best Practice Continuous Improvement **Top-Down** Methodology For Creating Structured Data Products That Starts With Defining Your Business Classary

**17**

There is Also a **Bottom-Up** Variation Of The Methodology To Create Structured Data Products By Discovering What Data Exists and Auto Generate Your Business Glossary Terms

**18**

## Topics – Where Are We?

- What is a data catalogue?
- The data catalogue marketplace
- ➢ Using a business glossary within a data catalogue to define data products
- Using a data catalogue to automatically discover data in multiple data sources
- Using a data catalogue to map raw data to common terms in a business glossary
- Using a data catalogue to automatically detect sensitive data in data sources
- Using a data catalogue to automatically profile data quality in your data sources and recommend fixes
- The power of metadata - Integrating data engineering tools and generative AI with a data catalogue to rapidly build data pipelines to produce data products
- Publishing data products in a data marketplace within the data catalog

19

## The Foundation For Data Management And Governance Is A Common Vocabulary And Lineage



The SBV acts as the foundation for sharing data across systems irrespective of whether those systems are on-premises or in the cloud - It is fundamental to getting rid of complexity

20

## A Shared Business Vocabulary Is Defined In A Business Glossary Which Is Typically Part Of A Data Catalog

Business Glossary

Shared business vocabulary (SBV) → Documented in → [Business Glossary screenshot] → Contained within → Data catalog

21

## Options For Defining Your Own Shared Business Vocabulary

1. Adopt pre-built vocabularies as a standard

2. Leverage existing enterprise data models
   • If these models exist in your organisation, they are probably managed by data architects and may be widely adopted

3. Use a business data concept model to start with and incrementally define your own vocabulary from scratch
   • Use a data catalog to automatically discover data to help build out all attributes and to help automate the mapping of physical data assets to the terms in your glossary

4. Use Generative AI to automatically generated business glossary terms for the data it finds during automatic data discovery within a data catalog

22

## Pre-Built Common Business Vocabularies Example - IBM Knowledge Accelerators - E.g. IBM Knowledge Accelerator for Financial Services



Source: IBM

23

## Pre-Built Vocabularies Example - Microsoft Common Data Model



Source: https://learn.microsoft.com/en-us/common-data-model/

24

## A Good Top-Down Approach To Starting A Business Glossary Is To Create A Business Data Concept Model

- Identify the data concepts, properties and relationships and construct a data concept model
- It is good practice to highlight all master and transaction data concepts in your data concept model

Business Data Concept Model Example



Some vendors provide complete definitions for business data entities and all their attributes to get you started quickly, e.g., Collibra uses Schema.org, Microsoft has a CDM

**25**

## Steps To Creating A Common Vocabulary
## – From Data Concept Model To Data Marketplace

1. Identify the data concepts, properties and relationships and construct a data concept model

2. The data concepts become the 'skeleton data entities' in the common vocabulary

3. Each data concept and its attributes should be defined in the business glossary as a data entity with a data owner

4. Use the catalog to discover the data for each data entity in underlying data stores across the data landscape

5. Design DataOps component-based pipelines to create the data products with common vocabulary data names

6. Publish all data products in a data marketplace

Data Concept Model Example

**26**

## Several Communities Are Likely To Be Involved In Defining Common Data Names And Definitions For Data In A Business Glossary

Communities may include a mix of IT and subject matter expert (SME) business users

**27**

## Data Catalogs Can Also Use AI to Auto Generate Business Data Names for Your Business Glossary That Can Be Accepted or Rejected - E.g. User Confirmed Business Terms in Alation

This is a bottom-up approach to creating a business glossary



Source: Alation

Green 'robot heads' indicate that the AI Generated business term has been confirmed by a business user

**28**

## Data Catalogs Can Also Recommend Business Terms That You Can Accept or Decline - E.g., Informatica IDMC Data Governance and Catalog

**29**

## Using Generative AI For Business Glossary Metadata Enrichment - Auto-Generation of Business Term Descriptions in Atlan Business Glossary



Business Glossary      Data catalog

Source: Atlan

Add a Readme, pick from a template and auto generate the Readme text using a generative LLM

**30**

## Topics – Where Are We?

- What is a data catalogue?
- The data catalogue marketplace
- Using a business glossary within a data catalogue to define data products
- ➢ Using a data catalogue to automatically discover data in multiple data sources
- ➢ Using a data catalogue to map raw data to common terms in a business glossary
- Using a data catalogue to automatically detect sensitive data in data sources
- Using a data catalogue to automatically profile data quality in your data sources and recommend fixes
- The power of metadata - Integrating data engineering tools and generative AI with a data catalogue to rapidly build data pipelines to produce data products
- Publishing data products in a data marketplace within the data catalog

31

## Automatic Data Discovery Is About Discovering What Data Exists In Your Data Stores – The Data Catalog Is Populated With Physical Data Names



Image Source: IBM

32

## Automated Data Discovery Using Data Catalog Crawlers (Scanners) - Scan The Data Stores, Populate The Data Catalog And Build Search Indexes



Scan time will vary depending on data volumes, scope of the scan, variety of the data and whether or not an incremental or full scan has been selected

automated data discovery

**Issues**
- Could result in millions of physical data names
- Multiple different names for same data in different data sources
- Still a challenge to find data
  - Not all data has meaningful data names
  - Searching may not reveal all unless relationships are understood

Oracle    Azure SQL    S3    MongoDB    Hadoop

33

## Registering Data Sources In A Data Catalog For Automated Discovery – Product Example: Microsoft Purview



Data sources can also be grouped into collections for easy management

Source: Microsoft

34

## Some Data Catalog Products Offer Incremental Data Discover Scans
## - E.g., Informatica IDMC Data Governance and Catalog



Source: Informatica

35

## The Mapping Of Discovered Data In A Data Store To Business Metadata In A Business Glossary Enables Us To Understand The Meaning Of Data



Source: IBM

36

## Automated Mapping of Raw Data to Business Terms Seeks to Determine the Meaning of Discovered Data and to Identify Complete Data Entities From Data Across Your Data Estate



- Automated determination and grouping of data across multiple systems into business entities

37

## How Do You Map Physical Data Names To Business Terms In A Business Glossary?

- Manual mapping done by business domain subject matter experts and data stewards
  - For occasional use and overriding incorrect automated mappings
  - Very time consuming and not practical for millions of files and thousands of tables and columns
  - Highly unlikely that any employee will understand data assets in SaaS applications

- **Automatically mapping** using several techniques
  - Pre-built machine learning models (AI-driven algorithmic scoring)
  - AI assisted mapping
    - Auto data discovery
    - Auto data clustering of similar data
    - Observability of manual mapping to glossary terms
    - Auto labelling of similar data based on observations
    - Training machine learning models to do it automatically
  - User defined rules, e.g., regular expressions
  - Reference data

38

## Slide 39

**Automated** Data Discovery, **AND Mapping To Common Business Terms In The Business Glossary** Using Data Catalog Crawlers, Rules and Machine Learning

Mapping to a business glossary allows the meaning of data items to be understood

**Data Catalog**

Business Glossary

| Customer | Product | Order |
|---|---|---|
| Customer_ID | | |
| Customer_First_Name | | |
| Customer_Surname | | |
| Customer_DOB | | |
| ... | | |

+ Physical data names

+ Search indexes

👍 Automatic mapping physical names to business terms

Oracle   Azure SQL   S3   MongoDB   Hadoop

**39**

## Slide 40

Enabling Glossary Business Term Association In Informatica IDMC Data Governance and Catalog To Map Physical Data Names To Business Terms in a Business Glossary

Informatica, can automatically map physical data names to business terms in a business glossary. It does this by using an algorithm (based on accepted business terms on data domains, column similarity, and name match between a column and business term) to calculate a confidence score that it can map it correctly. You can specify a threshold for the confidence score above which it will automatically assign a business term to a physical data asset.

Source: Informatica

**40**

## Auto Discovery of Complete Data Entities Such as Customer, Address, Product is Also Possible in Informatica Data Governance and Catalog Using Entity Classifications



Source: Informatica

Copyright © Intelligent Business Strategies 1992-2025

41

## Identifying Complete Data Entities - Informatica IDMC Data Governance and Catalog Can Classify Data Elements and Complete Data Entities



Source: Informatica

Copyright © Intelligent Business Strategies 1992-2025

42

## You Can Also Add Your Own Rules to Extend Data Catalogs - Using Regex in Pentaho Data Catalog to Identify Social Security Number Data and Map It Your Business Glossary Terms



Source: Hitachi Vantara

43

## Topics – Where Are We?

- What is a data catalogue?
- The data catalogue marketplace
- Using a business glossary within a data catalogue to define data products
- Using a data catalogue to automatically discover data in multiple data sources
- Using a data catalogue to map raw data to common terms in a business glossary
- ➢ Using a data catalogue to automatically detect sensitive data in data sources
- ➢ Using a data catalogue to automatically profile data quality in your data sources and recommend fixes
- The power of metadata - Integrating data engineering tools and generative AI with a data catalogue to rapidly build data pipelines to produce data products
- Publishing data products in a data marketplace within the data catalog

44

## Several Vendors Support Predefined Classifiers To Classify Sensitive Data Types – E.g., AWS Glue Studio Can Detect And Process Sensitive Data

○ Detect PII in each cell
Scan the entire data set, and act on each occurance individually.

○ Detect fields containing PII
To reduce costs and improve performance, sample only a portion of the data and act on fields across all records.

**Select entities to detect**                    ×

Available entities (19)   ↻   Select all   Clear all   Create new ↗   Manage ↗

🔍 Find entities                    All categories ▼    ‹ 1 ›

| Entity name | ▽ | Category | ▲ |
|---|---|---|---|
| ☐ Person's name | | Universal, HIPAA | |
| ☐ Email (General) | | Universal | |
| ☐ Credit Card | | Universal | |
| ☐ IP Address | | Networking | |
| ☐ MAC Address | | Networking | |
| ☐ US Phone | | United States, HIPAA | |
| ☐ US Passport | | United States | |
| ☐ Social Security Number (SSN) | | United States, HIPAA | |
| ☐ US Individual Taxpayer Identification Number (ITIN) | | United States, HIPAA | |
| ☐ US/Canada bank account | | United States, HIPAA | |
| ☐ US driving license | | HIPAA | |
| ☐ Healthcare Common Procedure Coding System (HCPCS) code | | HIPAA | |
| ☐ National Drug Code (NDC) | | HIPAA | |
| ☐ National Provider Identifier (NPI) | | HIPAA | |
| ☐ Drug Enforcement Agency (DEA) Registration Number | | HIPAA | |
| ☐ Health Insurance Claim Number (HICN) | | HIPAA | |
| ☐ Medicare Beneficiary Identifier | | HIPAA | |

- Universal (e.g., Email, Credit Card)
- HIPAA (e.g., US Driving License, Healthcare Common Procedure Coding System (HCPCS) code)
- Networking (e.g., IP Address, MAC Address)
- United States (e.g., US Phone, US Passport)
- United Kingdom (e.g., UK Bank Account, UK VAT)
- Japan (e.g., Japan My Number, Japan Passport)

Choose what you want to do with identified PII data

**Actions**
Choose actions to take on detected entities.

○ Enrich data with detection results
Create a new column that will contain any entity type detected in that row.

○ Redact detected text
Replace detected entity with a string you choose.

○ Apply cryptographic hash
Apply a SHA-256 cryptographic hash function to the input string.

Source: AWS

45

## Informatica Metadata Command Center Let's You Import Over 200 Predefined Classifications That You Can Use To Identify Country Specific Personal Data In Sources



Informatica  Metadata Command Center   monitoringdev2.fromselfservice ∨

**Explore**                                   Import Predefined Content

Data Classifications ▼

**Data Classifications (23)**          ↻  ↑↓  ▽  Find  ×

| Name ^ | Type | Description | Updated On | Updated By |
|---|---|---|---|---|
| 👤 ABA Routing No | Data Element | This identifies ABA r | 10 Nov 2021, 03:33 | monitoringdev2 |
| 👤 All Match | Data Element | | 10 Nov 2021, 03:29 | monitoringdev2 |
| 👤 Amex | Data Element | This identifies Amer | 10 Nov 2021, 03:33 | monitoringdev2 |
| 👤 Birth date | Data Element | It matches a birthda | 10 Nov 2021, 03:33 | monitoringdev2 |
| 👤 Company Ticker USA | Data Element | This identifies the Sy | 10 Nov 2021, 03:33 | monitoringdev2 |
| 👤 Credit Card | Data Element | This validates VISA, | 10 Nov 2021, 03:33 | monitoringdev2 |

Source: Informatica

Note that Informatica predefined classifications cannot be modified

46

## Several Vendors Support Predefined Classifiers To Classify Sensitive Data Types – E.g. 200+ Pre-Defined Classifiers Used In Microsoft Purview

### Examples

- People's names
- Phone numbers
- Email addresses
- Postal codes
- Passport numbers
- Driving license numbers
- Social security numbers
- Credit card numbers
- Bank account numbers
....

**Sensitive info types**

These are examples of very specific sensitive data classifications

You can also add your own sensitive data types

**Data**

Files RDBMS Hadoop NoSQL Cloud storage

Scan & classify

You need to identify where this sensitive data is in all data stores and content

Scan & classify

**Content**

Copyright © Intelligent Business Strategies 1992-2025

47

---

## Microsoft Purview Data Catalog – Table Schema And Automated Classification Of A Discovered Data Asset



| Column name | Data type | Classifications | Glossary terms | Description |
|---|---|---|---|---|
| CustomerID | int | CustomerID | Customer | Unique Customer Identifier |
| CompanyName | int | | | Customer Company Name |
| EmailAddress | int | EmailAddress | | Customer Email Address |
| Phone | nvarchar | | | Customer Phone Number |
| FirstName | nvarchar | Person's Name | | Customer Contact First Name |
| LastName | nvarchar | Person's Name | | Customer Contact Last Name |
| SalesInvoiceNumber | int | InvoiceNumber | | Unique Sales Invoice Number |
| OrderQuantity | int | | SalesOrder | Sales Order Quantity |
| ProductCode | nvarchar | ProductCode | | Standardized Product Code |
| DiscountAmount | nvarchar | | | Discounted Amount |
| BillingID | int | | Finance | Unique Billing Identifier |
| CostCenterCode | nvarchar | CostCenterCode | | Unique Code for Accounting Cost Center |
| SalesAmount | int | | | Sales Amount |
| TaxAmount | int | | Tax | Tax Amount |

Source: Microsoft

Copyright © Intelligent Business Strategies 1992-2025

48

---

**Microsoft Purview Data Catalog – Classification Counts Show How Many Data Sources, Files And Tables Have That Data Of That Classification Type**

This allows you to see the proliferation of sensitive data across your data estate so that you can create policies to protect it

Source: Microsoft

49



**Generative AI in Data Governance**
**– Defining Data Quality Data Validation Rules Using Collibra AI Assistant**

Source: Collibra

50

## Some Data Catalogs Also Offer AI-Assisted Auto-generation of Data Quality Rules – e.g. Microsoft AI-Suggested Data Quality Rules



Source: Microsoft

Source: Microsoft

51

## Data Catalogs Can Perform Automatic Data Quality Profiling and Generate Metadata In the Data Catalog – E.g. IBM Knowledge Catalog



Quickly profile the data source for a quick overview of the data quality using data sampling – a good option for large data sets

This is much longer running as it does a full analysis of your data to provide detailed data quality insights

Source: IBM

Automated full and incremental data quality profiling is also supported by other vendors e.g. Informatica EDC

52

Automated Data Quality Profiling – Cloud Pak For Data and IBM Knowledge Catalog

Publish puts the DQ metadata into IBM Knowledge catalog for all to see

Source: IBM



IBM Cloud Pak For Data Automated Data Quality Profiling Highlights Data Quality Issues, E.g., Data Violates The Matching Rules Set Up In A Data Class

Source: IBM

## IBM Knowledge Catalog Allows You To Add Poor Quality Data To A Project And Immediately Start Cleaning It Using Self-Service Data Preparation



Source: IBM

**55**

## Data Observability Can Also Monitor Executing Data Pipelines, Cause Alerts and Pass Issues to Data Catalogs

**56**

## Real-Time AI-Driven Active Data Governance In Data Quality Can Come From Integration Between Data Catalogs, Data Observability and DBMSs – E.g. Alation DQ Alerts



Source: Alation

Integrates with SODA, Monte Carlo, Snowflake DMFs, Databricks etc.
Alation has DQ processes for Snowflake and Databricks to help govern DQ in these platforms

57

## In Some Data Catalogs, Data Quality Can Also Be Overlayed on Top of Data Asset Lineage to See Data Quality Issues at a Glance – Alation Data Quality Overlay



Source: Alation

58

## Topics – Where Are We?

- What is a data catalogue?
- The data catalogue marketplace
- Using a business glossary within a data catalogue to define data products
- Using a data catalogue to automatically discover data in multiple data sources
- Using a data catalogue to map raw data to common terms in a business glossary
- Using a data catalogue to automatically detect sensitive data in data sources
- Using a data catalogue to automatically profile data quality in your data sources and recommend fixes
- ➤ The power of metadata - Integrating data engineering tools and generative AI with a data catalogue to rapidly build data pipelines to produce data products
- Publishing data products in a data marketplace within the data catalog

59

## Data Producers Should Check If Data Is Classified As Sensitive In The Data Catalog And Mask Sensitive Data In The Pipeline When Creating Data Products

60

## Data Masking In A Pipeline That Produces a Data Product – E.g. IBM StreamSets



Source: StreamSets

Masking capability is available in several data fabric software offerings

61

## Use Metadata, Gen AI Data Automation To Generate Pipelines To Produce A Data Product For Each Entity Defined In The Business Glossary Within A Data Catalog
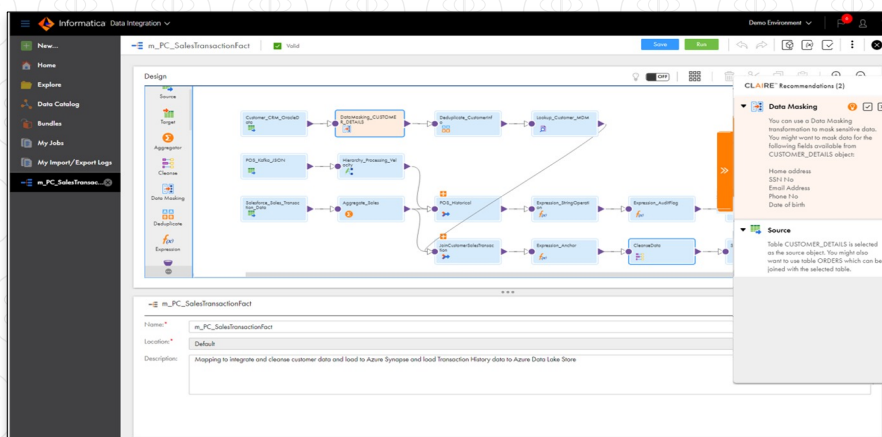
62

31

## AI Assisted Data Fabric Software Can Now Automatically Generate Pipelines – Informatica IDMC CLAIRE Co-Pilot AI-Driven Automated Pipeline Generation



Source: Informatica

- Fully automated data pipeline generation
- Support for unstructured, semi-structured and hierarchical data sources
- Automated creation data masking, cleansing and standardization rules
- Support for join, union, normalization and denormalization of data sources

- Simplifies and accelerates development
- Improves productivity

63

## Generative AI In Data Engineering Opens Up Data Product Development To Citizen Data Engineers – SnapLogic SnapGPT Prompt Based Data Engineering



Generated pipeline preview

Import the pipeline

64

## Generative AI Prompt Based Data Engineering
– SnapLogic SnapGPT Pipeline Configuration Wizard



Source: SnapLogic

65

## Topics – Where Are We?

- What is a data catalogue?
- The data catalogue marketplace
- Using a business glossary within a data catalogue to define data products
- Using a data catalogue to automatically discover data in multiple data sources
- Using a data catalogue to map raw data to common terms in a business glossary
- Using a data catalogue to automatically detect sensitive data in data sources
- Using a data catalogue to automatically profile data quality in your data sources and recommend fixes
- The power of metadata - Integrating data engineering tools and generative AI with a data catalogue to rapidly build data pipelines to produce data products
➢ Publishing data products in a data marketplace within the data catalog

66

## Where Are We? – Publishing Trusted Data Products In A Marketplace For Consumption



- Need to make use of
  - A business glossary and data catalog
  - A collaborative approach to producing data and analytical products
  - A catalog / data marketplace to quickly find reusable trusted data products to drive business value

67

## What Is An Enterprise Data Marketplace?

### Enterprise Data Marketplace

A catalog application that governs the sharing of published ready-made, trusted, data and analytical products that are available as services with common data names documented in a business glossary, full metadata lineage and that are tagged and organised to make them easy to find, access, share and reuse across the enterprise

Product examples:
- Alation Marketplaces
- Amazon Datazone
- Collibra Data Marketplace
- Databricks Data Marketplace
- CData (Data Virtuality) Data Shop
- Harbr Data
- IBM Data Product Hub
- Informatica IDMC Data Marketplace
- One Data Data Products Marketplace
- Microsoft Purview
- Quest Data Marketplace
- Qlik Data Products Catalog
- Snowflake Internal Data Marketplace
- Snowflake External Data Marketplace
- SAP DataSphere Data Marketplace

68

## Key Elements In Enterprise Data Sharing – Governing The Publishing Of Data Products

▪ Quality assurance and governance of data products is needed before they are published in a marketplace for consumption

Business Glossary +
Data Lineage
+ Data Freshness
+ Data Owners

| Data Quality score | ✓ |
|---|---|
| Common data names (business glossary) | ✓ |
| Full metadata lineage | ✓ |
| Access security policy defined | ✓ |
| Privacy policy defined | ✓ |
| Retention policy defined | ✓ |
| Version management | ✓ |
| **Data sharing agreement defined (terms & conditions)** | ✓ |
| Tagged by business objective | ✓ |

Data product quality assurance

Data marketplace

Data Catalog

Data product

69

## Collibra Data Marketplace Includes Recommended Data



70

## Information Consumers Shopping For Data In An Enterprise Data Marketplace – E.g. Informatica IDMC Data Marketplace Data Categories



Audit history, reporting, analytics on data consumption patterns

Categories contain data collections

Data collections contain data assets

Data assets can contain one or more data elements

Source: Informatica

71

## Informatica IDMC Data Marketplace – Includes Data Products And Machine Learning Models



AI model plus its training datasets and tables from different sources

Source: Informatica

72

## Informatica IDMC Data Marketplace – Data Product Terms Of Use (Data Contract)



Terms of use indicate that this data collection is controlled because it contains PII data

Source: Informatica

Copyright © Intelligent Business Strategies 1992-2025

73

## Some Data Catalogs Will Now Let You View Data Quality Scores By Business Domain and By Data Products – E.g. Microsoft Purview

In this case there is a hierarchy showing Corporate Functions and then business domains



Source: Microsoft

Copyright © Intelligent Business Strategies 1992-2025

74

Microsoft Purview
- You Can Also Monitor Health Actions Associated With Data Products



Copyright © Intelligent Business Strategies 1992-2025

75

Microsoft Purview - Generative AI Can Suggest New Assets to Add To a Data Products



Copyright © Intelligent Business Strategies 1992-2025

76

Microsoft Purview - Data Product Owners Can Also Set Policies Including the Duration of the Data Access Policy

Copyright © Intelligent Business Strategies 1992-2025

77



The Last Mile - BI, ML And AI Are 'Wired Into' Every Process At All Levels To Enable Timely, Guided And Automated Decisions To Maximise Value – Create A Closed Loop

Embedded analytics in applications

Intelligent applications

Decision automation

AI Agents & AI-automated tasks

Access to commonly understood, data, BI reports, ML models, AI services across the business is critical to success

Copyright © Intelligent Business Strategies 1992-2025

78

## About Mike Ferguson

www.intelligentbusiness.biz

mferguson@intelligentbusiness.biz

@mikeferguson1

(+44) 1625 520700

Mike Ferguson is CEO of Intelligent Business Strategies Limited a UK boutique research and consulting firm. He has been an independent consultant and IT industry analyst for over 30 years and is a world authority in data management and AI. Mike has over 40 years of IT experience and has consulted worldwide for dozens of companies on data strategy, data architecture, data management technology selection, data governance, AI and decision automation. Mike is also conference chairman of Big Data LDN, the largest data and AI conference in Europe and is a member of the EDM Council CDMC Executive Advisory Board. He has consulted and spoken at events all over the world and has written countless articles. Formerly he was a co-founder and principal of Codd and Date with Turing Award winner Dr E.F. Codd – the inventor of the Relational Model which caused the birth of relational databases and the SQL language. He was also formerly Chief Architect at Teradata helping pioneer the first massively parallel DBMS. He teaches popular master classes in 9 countries in Data & AI Strategy, Modern Data Architecture, Practical Guidelines for Implementing Data Products, AI Driven Active Data Governance, and Embedded Analytics, Intelligent Apps, AI Agents and AI Automation

**79**