June 8, 2021

# Data Governance FAQ's – Axon, EDC, IDQ, DPM

Steven Fleishman, Principal Consultant
Sumit Saraswat, Solution Architect
Informatica Professional Services

Informatica™

# Frequently Asked Questions from the Field

- What does CLAIRE exactly do ?
  - Are we leveraging it or how can we leverage it better ?
  - Does it learn from human actions of curation etc ?
- I see profiling in EDC, IDQ, and DPM. Are they the same? Which one is leveraged when?
- Where does AI/ML come into the picture in the Informatica DG Solution What does it do exactly?
- How can we make the information in Axon actionable ?
- What is the approach to classify information easily and efficiently in Axon ?
  - Similarly, how can I classify information in EDC if I do not have DPM?
- How can I easily generate a 360 graphical view of related and impacted assets in Axon?
- How to segment information based on LOBs/departments in Axon ?
  - Can I have a common/shared repository of assets and a department specific one.
  - Can I have local/private change management processes for my specific group ?
- How can I record and expose data dictionaries in the tools ?
- What are the best practices of scanning and cataloging now widely adopted solutions such as the data lake on S3 or Azure ?
- What does CLAIRE exactly do ?
  - Are we leveraging it or how can we leverage it better ?
  - Does it learn from human actions of curation etc ?
- I see profiling in EDC, IDQ, and DPM. Are they the same? Which one is leveraged when?
- Where does AI/ML come into the picture in the Informatica DG Solution What does it do exactly?
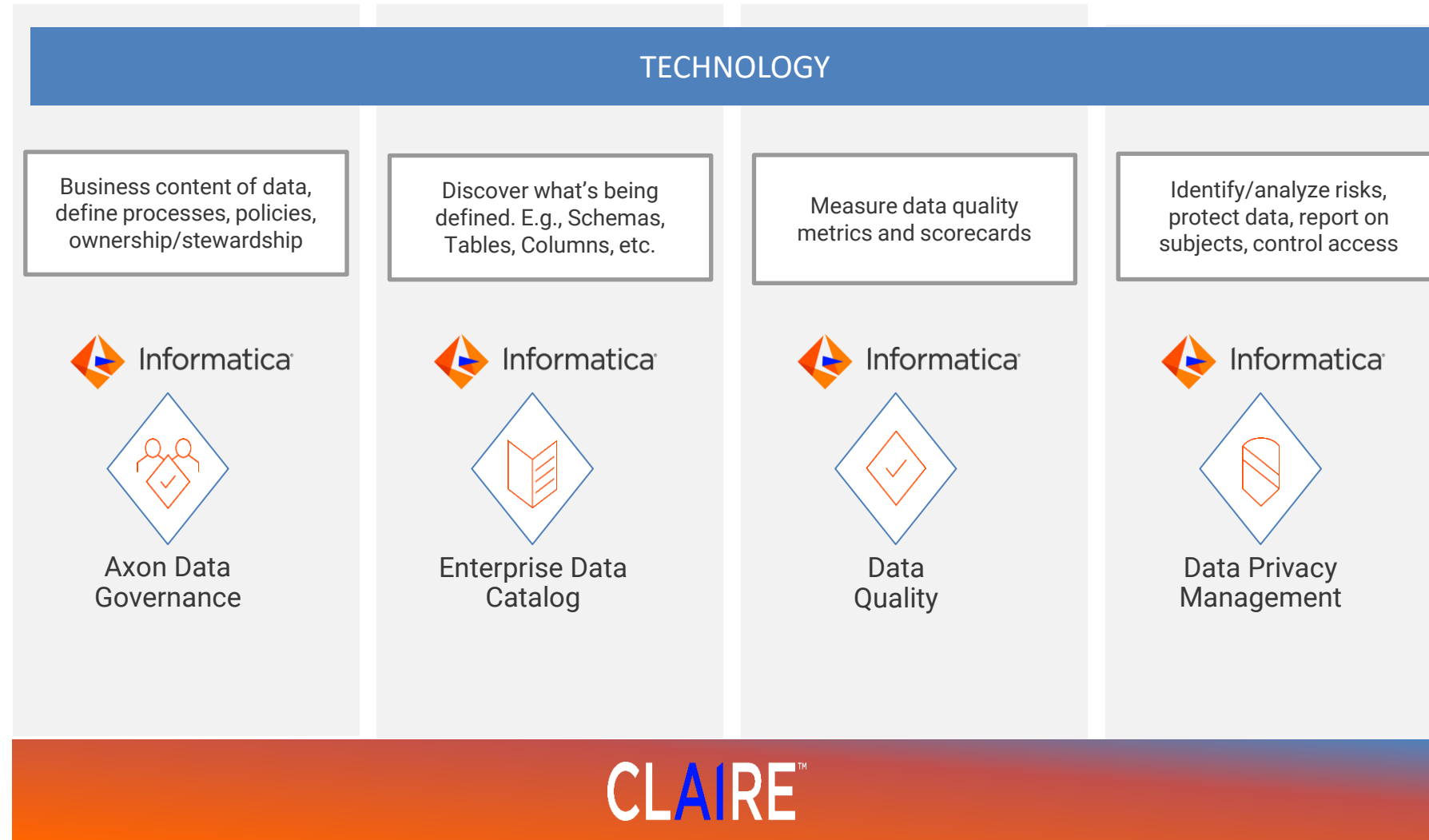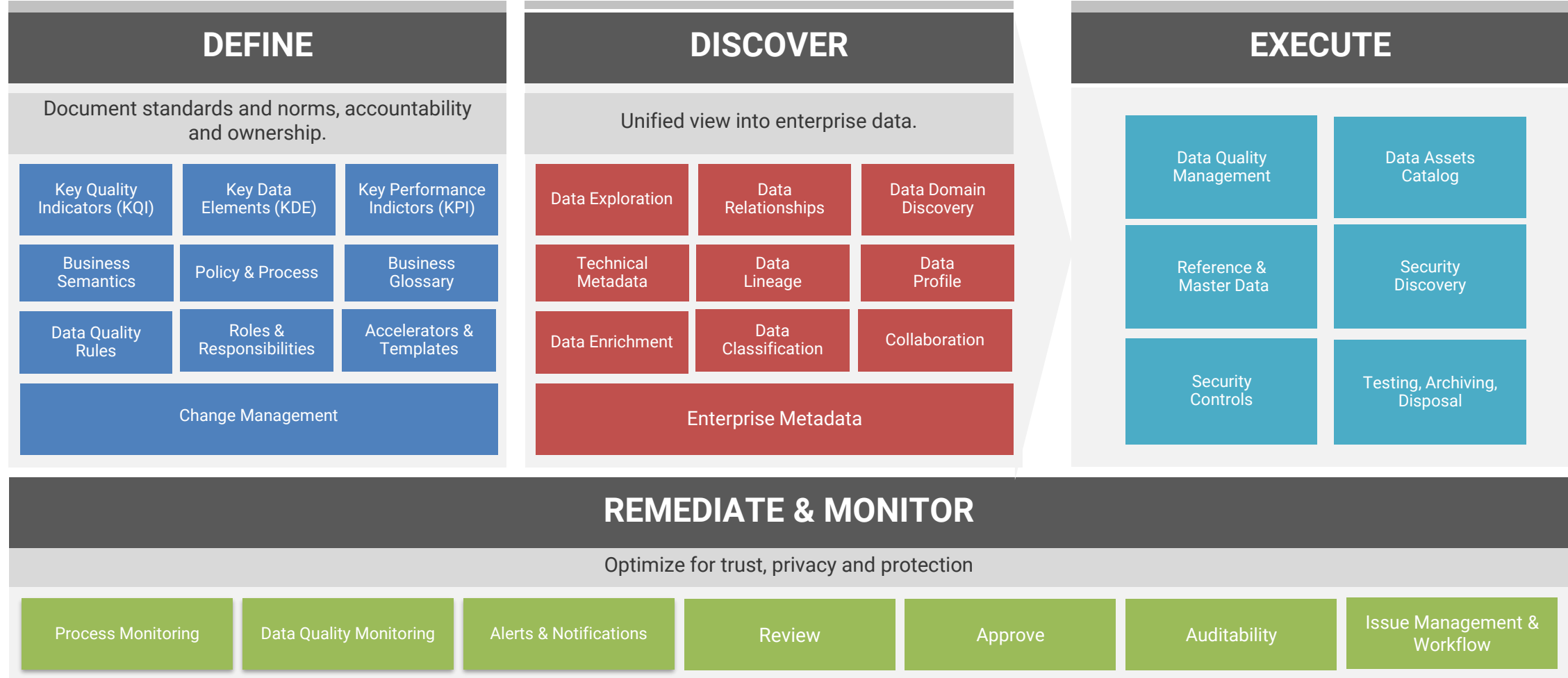
# Agenda

- Product Integration, Navigation, and Terminology

- Domains and Data Domain Types

- AI/ML - Claire

- Profiling functionality – EDC compared to IDQ

- Q&A

Informatica
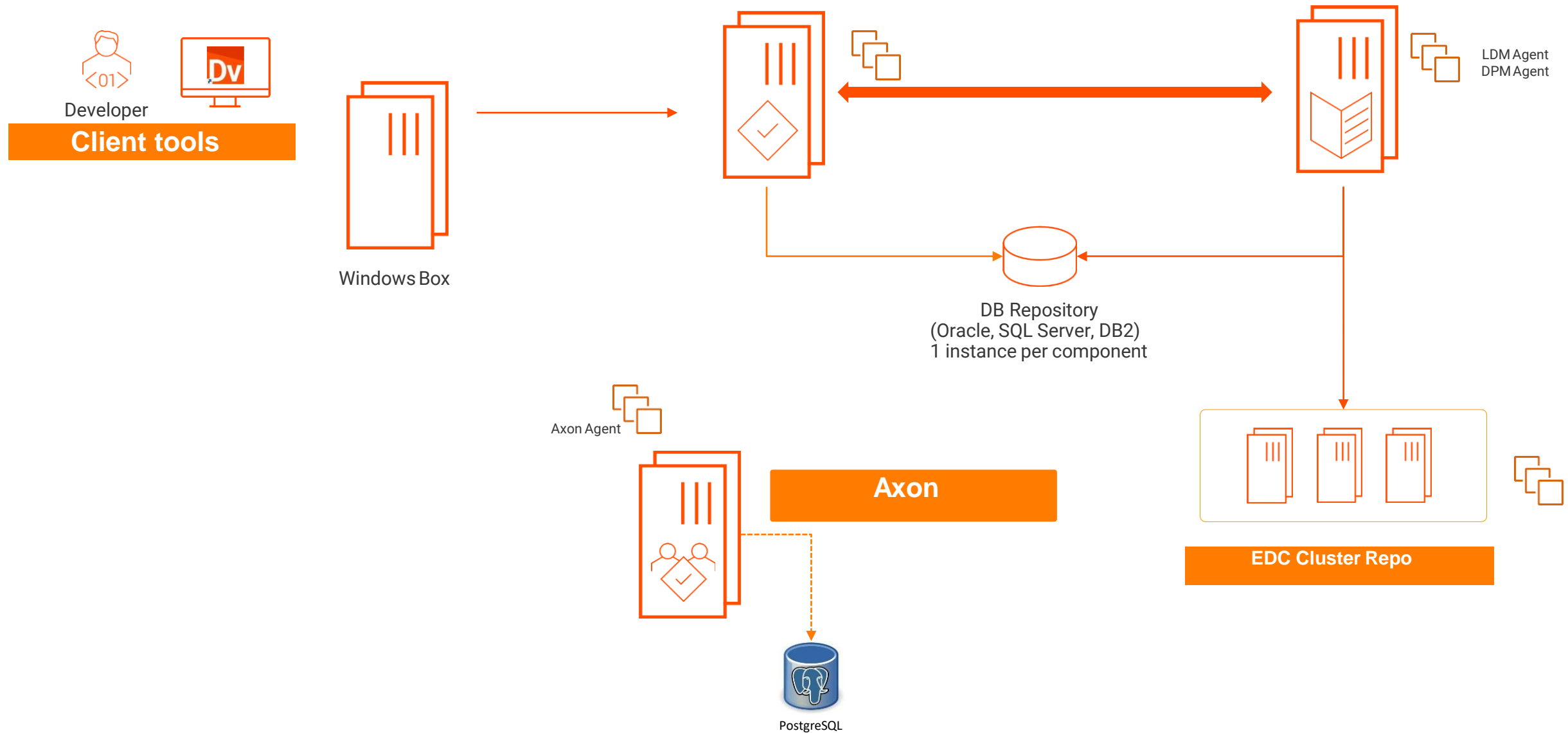
# Product Integration, Navigation, and Terminology

# Data Governance and Privacy Solutions

# Architecture to Support Data Governance Framework

## DEFINE

Document standards and norms, accountability and ownership.

| Key Quality Indicators (KQI) | Key Data Elements (KDE) | Key Performance Indictors (KPI) |
|---|---|---|
| Business Semantics | Policy & Process | Business Glossary |
| Data Quality Rules | Roles & Responsibilities | Accelerators & Templates |

Change Management

## DISCOVER

Unified view into enterprise data.

| Data Exploration | Data Relationships | Data Domain Discovery |
|---|---|---|
| Technical Metadata | Data Lineage | Data Profile |
| Data Enrichment | Data Classification | Collaboration |

Enterprise Metadata

## EXECUTE

| Data Quality Management | Data Assets Catalog |
|---|---|
| Reference & Master Data | Security Discovery |
| Security Controls | Testing, Archiving, Disposal |

## REMEDIATE & MONITOR

Optimize for trust, privacy and protection

| Process Monitoring | Data Quality Monitoring | Alerts & Notifications | Review | Approve | Auditability | Issue Management & Workflow |
|---|---|---|---|---|---|---|

Informatica™

# INFA Platform Architecture

**Client tools**

Developer

Windows Box

**IDQ**

**EDC+DPM**

LDM Agent
DPM Agent

DB Repository
(Oracle, SQL Server, DB2)
1 instance per component

Axon Agent

**Axon**
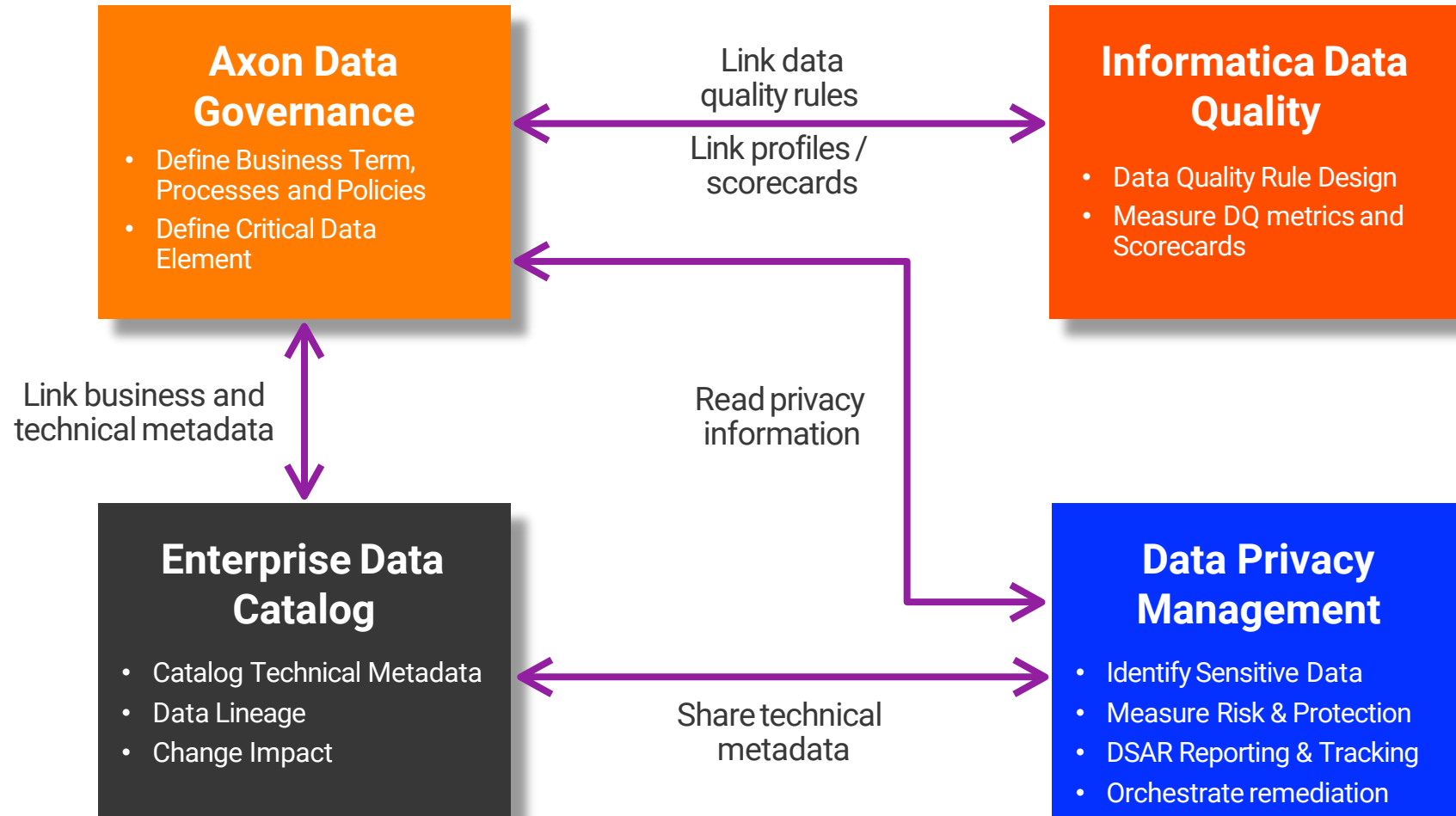
PostgreSQL

**EDC Cluster Repo**

# INFA Split Domain: EDC and IDQ

**Recommendation and Best Practice for EDC and/or DPM and IDQ to be installed in separate Domain, here are  pointers:**
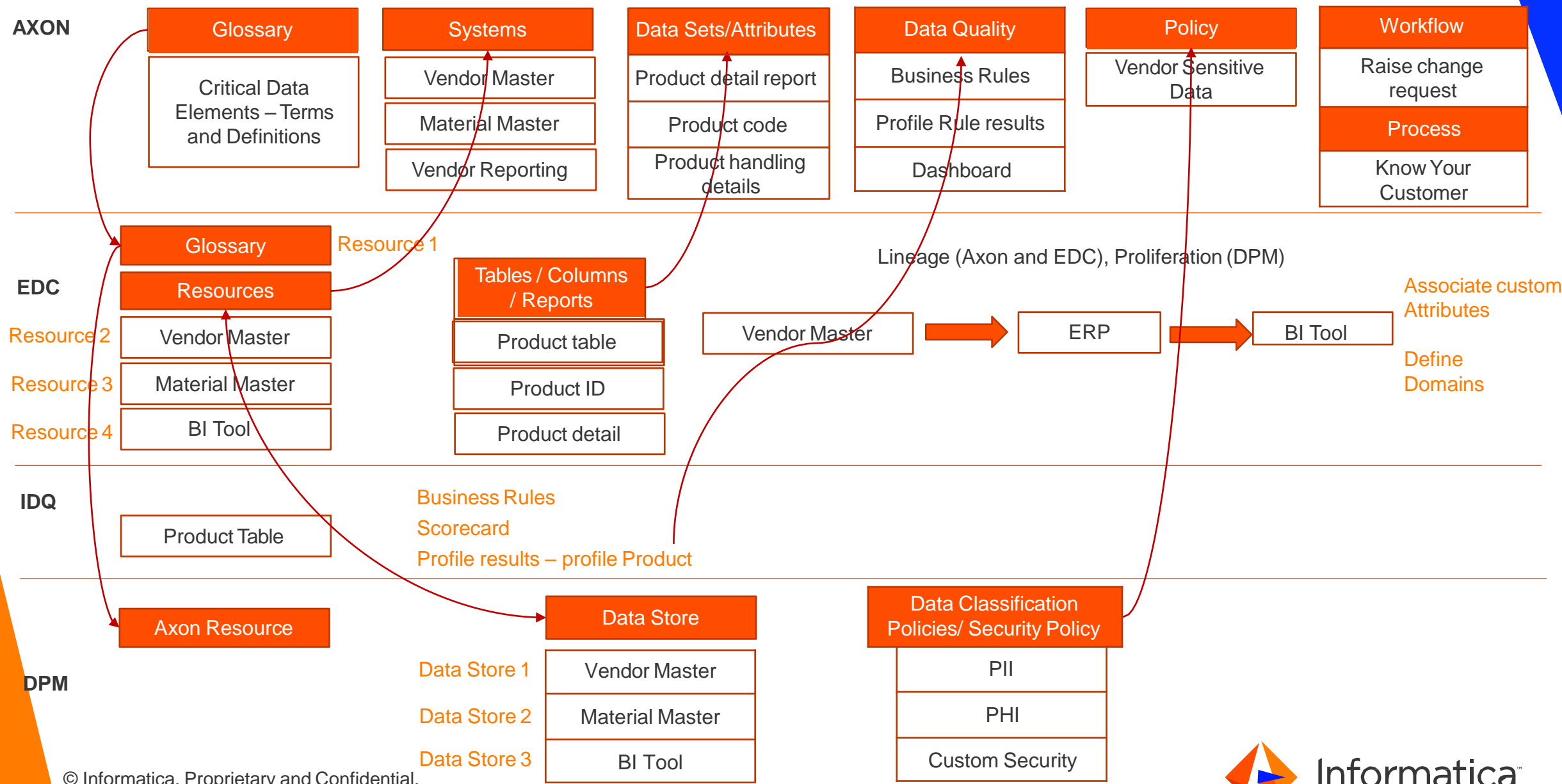
- Flexibility applying patches, fixes, upgrades for respective product
- IDQ is higher volume (longer running jobs-less jobs-more operational driven)
- EDC is Metadata (more jobs-less operational driven)
- IDQ licensing is based on number of cores in the machine, whereas EDC licensing is based on number of Resources
- Profiling: Context of Profiling in EDC is for Data Domain Discovery, Similarity Discovery, Unique Key Inference, CLAIRE on larger set of data, however context of Profiling on IDQ is to perform checks on Data Quality Rules, Scorecards focused on key sets of data.

Informatica

# Cross Product View



**Axon Data Governance**
- Define Business Term, Processes and Policies
- Define Critical Data Element

**Informatica Data Quality**
- Data Quality Rule Design
- Measure DQ metrics and Scorecards

Link data quality rules

Link profiles / scorecards

Link business and technical metadata

Read privacy information

**Enterprise Data Catalog**
- Catalog Technical Metadata
- Data Lineage
- Change Impact

**Data Privacy Management**
- Identify Sensitive Data
- Measure Risk & Protection
- DSAR Reporting & Tracking
- Orchestrate remediation

Share technical metadata

**Informatica**

# Data Governance Application Relationships

**AXON**

| Glossary | Systems | Data Sets/Attributes | Data Quality | Policy | Workflow |
|---|---|---|---|---|---|
| Critical Data Elements – Terms and Definitions | Vendor Master | Product detail report | Business Rules | Vendor Sensitive Data | Raise change request |
| | Material Master | Product code | Profile Rule results | | **Process** |
| | Vendor Reporting | Product handling details | Dashboard | | Know Your Customer |

**EDC**

Resource 1

Lineage (Axon and EDC), Proliferation (DPM)

| Glossary |
|---|
| **Resources** |
| Vendor Master |
| Material Master |
| BI Tool |

Resource 2
Resource 3
Resource 4

| Tables / Columns / Reports |
|---|
| Product table |
| Product ID |
| Product detail |

Vendor Master → ERP → BI Tool

Associate custom Attributes

Define Domains

**IDQ**

Business Rules
Scorecard
Profile results – profile Product

| Product Table |
|---|

**DPM**

| Axon Resource |
|---|

| Data Store |
|---|
| Vendor Master |
| Material Master |
| BI Tool |

Data Store 1
Data Store 2
Data Store 3

| Data Classification Policies/ Security Policy |
|---|
| PII |
| PHI |
| Custom Security |

Informatica™

# Data Domains and Domain Types

# "Domain" Usage in Data Governance and Privacy

- Informatica Domain
  - A collection of nodes and services that define the Informatica platform. You group nodes and services in a domain based on administration ownership

- Axon Domain
  - A glossary type, that's a way of classifying data
  - Describes a broad category of data concepts, for example, customer domain or transaction data domain
  - Specific to Axon and can be modified

- Data Domain
  - Predefined or user-defined Model repository
  - Based on the semantics of column data or a

**CLASSIFICATIONS**

BASIC CLASSIFICATIONS

Axon Status *

| Active | ⌄ |

Lifecycle *

| Approved | ⌄ |

Axon Viewing *

| Public | ⌄ |

Type *

| Domain | ⌄ |

Informatica™

# Types of Data Domains

- Rule-Based
  - Run against Metadata, Data or Both
  - 125+ predefined data domains
  - Regex - pattern
    - credit card, SSN, phone number
  - Reference – finite, non-overlapping
    - ISO country code, currency codes
  - Mapplet – Leverage Informatica Developer and Analyst for complex rules

Informatica™

# Types of Data Domains Continued

- ## Smart – Specific to EDC
  - Example based data domain
  - Data tagging and propagation

- ## Composite Data Domain

- ## Data Domain Group

- Collection of data domains or other composite data domains linked using rules
- Enables you to search for the required details of an entity across multiple schemas defined for the database



Informatica™

# Process to create custom data domains

```
Start
  │
  ▼
Create Data Domain
```

**UI used**

- Informatica Analyst
- Informatica Developer
- Catalog Administrator

```
Informatica Analyst ──► Use Domain Glossary for Domain Creation
Informatica Developer ──► Create Rule Specification ──► Generate Rule to save as Mapplet in Developer Tool
Create Mapplet ──► Validate Mapplet as Rule ──► Use Rule in domain glossary for domain creation
Catalog Administrator ──► Predefined rules ☑ / Use Reference Table ☑ / Use RegEx ☑
```

- Data Domain Created ──► Run DataDomain System Resource in Catalog Administrator ──► Stop

**Informatica**

# Out of the Box Data Domains

- The following data domains may create large number of false positives; Use with caution
  - Age
  - Salary
  - Weight
  - Height
  - Alphanumeric_specialCharacters
  - Date_allFormats
  - Admission_dates
  - JobPosition
  - Binary Value
  - Admission_date

- Avoid using "All" data domains

- Make a copy of the original data domain before modifying

Informatica™

# AI / ML - Claire

# CLAIRE

- CLAIRE stands for Cloud-Scale AI powered Real-Time Engine.
- Identifies all capabilities in Informatica products and services that use artificial intelligence (AI) and machine-learning techniques on enterprise-wide data and metadata to significantly boosts the productivity and experience of users of our technology.
- The only real way to discover velocity and diversity of data  manage this complexity is to increase automation and to significantly improve the productivity and effectiveness of the data management staff.
- This is where artificial intelligence and machine learning come in.

Enrichments

- Custom Attributes (22)
- Business terms
  - Business term ingestion
    - from Axon
    - from BG
  - Synonyms: Admin user input
  - Association: User input
    - Individual
    - Bulk
  - Intelligent Glossary association
    - Auto association
    - Curation
      - UI (1)
      - REST API
- Data Domains
  - Rule based data domain
    - Association: User input
    - data domain discovery
      - Auto association
      - Curation
        - UI (2)
        - REST API
  - Smart Domain
    - Creation: User input
    - Association: User input
    - Propagation
      - Auto association
      - Curation
  - Composite domain
- Relationship discovery (6)

CLAIRE
- Column similarity
  - Pattern Match
  - Name Match
  - Distinct Values Match
  - Data Overlap (Value frequency)
  - Similarity Confidence
- BT Propagation
  - Name match
  - Domain Conformance
- UK inference

Informatica

# Smart Data Domains

Process of discovering semantic meaning of data in the data sources

Smart domains

- Act as tags
- Learn by example and propagated by looking at column similarity.
- Exist as an object in the catalog and can be enriched as well.
- Requires access to the data

| (650) 385-5000 | 95008 | Darren | gp@gmail.com | Informatica |
|---|---|---|---|---|
| Phone Number | Zip | First Name | Email | Company Name |

Informatica

# Column similarity

- Identify clusters of columns that contain similar data within and across data sources.
- **Use:**
- Identifying data
- Detecting duplicates
- Combining individual data fields into business entities
- Propagating tags across data sets
- Recommending data sets to users

# Business term association through propagation



- When data domains are inferred against specific columns, the associated glossary terms are recommended for those columns.

- When data domains are accepted, associated glossary terms are also associated to the columns

---

- System propagates business glossary terms to similar column

- Similarity based on name match, unique value match and data match is used for business glossary propagation

# Business term association through Claire Match

- Match English phrases with technical names using sequence alignment

  - **Sequence Alignment / Delete-only Edit Distance :** The business term names that align well with asset names are sought. This approach can capture obvious abbreviations of business terms.

    - HEALTH PROGRAM CONSULTATION  (Business Term Title)
    - H- - LTH P - - G - - M C- NS - LT- T- -N  (Asset name)

  - **Synonym dictionary :** If available, user provides a dictionary of commonly used synonyms/abbreviations in technical asset names within the organization. This dictionary is used to improve glossary matching

- Additionally, prefix ignore options for discarding common technical prefixes(like TBL, VW etc for better matches)

Informatica

# How does profiling differ between IDQ and EDC?

# EDC − Broad Profiling Results − Table View

Business Title

Asset Certification

Data Owner

Data Domain

Business Terms

Custom Attributes

Basic Data Profile

# EDC − Broad Profiling Results − Column View

Basic Data Profile

Pattern Distribution

Data Domain

Value Frequencies

Column Similarity

**City**
Field3
Oracle_Hermes > SCALA11GR2 > HERMES > CUSTOMER_LOYALITY

Overview | Lineage and Impact | Relationships

**Description**
City

**Value Frequency**
Total Rows 179    Null | Distinct | Non Distinct
Min - Allen
Max - ZEBULON

| Value | Frequency | Percentage |
|---|---|---|
| Kansas City | 10 | 5.78 |
| Dallas | 7 | 4.05 |
| New York | 7 | 4.05 |

**Similar Columns**

.../DUPLName
Field3
**99%** Confidence
Distinct Values | Pattern | Data
Bank_Ro... | Address(... | +1

.../INDIA_CUSTOMERS
Field3
**99%** Confidence
Distinct Values | Pattern | Data
City | Bank_Ro... | +1

.../JAPAN_CUSTOMERS
Field3
**99%** Confidence
Distinct Values | Pattern | Data
City | Bank_Ro... | +1

**People**
Data Owner
Mary Smith
Data Steward
Peter Kiley

**Data Domains**
City | Address(58.6...

**Pattern**

| | rows | % |
|---|---|---|
| Others | 47 rows | 26.26% |
| X(8) | 23 rows | 12.85% |
| X(7) | 22 rows | 12.29% |
| X(5) | 13 rows | 7.26% |
| X(6)bX(4) | 13 rows | 7.26% |
| X(6) | 20 rows | 11.17% |
| X(10) | 20 rows | 11.17% |
| NULL | 6 rows | 3.35% |
| X(9) | 15 rows | 8.38% |

**Inferred Data Types**

| | rows | % |
|---|---|---|
| String(17) | 173 rows | 100.00% |

**Informatica**™

# IDQ – Broad **and** Deep Enterprise-Grade Data Management Solution

**Discovery, search & profiling**

**Role-based capabilities**
Enable business users to build and test logical business rules without relying on IT

**Rich set of transformations**
Manage and transform data with data standardization, validation, enrichment, de-duplication, and consolidation capabilities.

**Reusable rules & accelerators**
Apply pre-built business rules and accelerators and reuse common data quality rules to save time and resources.

**Exception management**
Allow business users to review, correct, and approve exceptions throughout the automated process.

**Informatica**™

# Select only columns to be profiled

IDQ



Select the columns you want to profile on.

# Compare Profile Results

Understand data quality trends through time by comparing historical profile results
Compare column and rule profile results between two profile runs
Detailed comparisons include changes in datatypes, patterns, nulls and distinct counts

# Resources

1. Configure Access Axon/IDQ: [Click Here](#)
2. Configure Access Axon/EDC: [Click Here](#)
3. Configure Access Axon/DPM: [Click Here](#)
4. Axon/EDC Automatic Onboarding Workflow: [Click Here](#)
5. Automate Data Quality Rules in Axon: [Click Here](#)
6. EDC Sizing Guide: [Click Here](#)
7. Profiling Sizing Guide: [Click Here](#)
8. Integrated Monitoring for Capacity Planning/Resource Utilization: [Click Here](#)
9. Product Availability Matrix (PAM): [Click Here](#)
10. AWS Informatica Marketplace Offerings: [Click Here](#)
11. Azure Informatica Marketplace Offerings: [Click Here](#)
12. Deploying DIS on GRID: [Click Here](#)
13. Informatica Axon Data Governance Playbook: [Click Here](#)
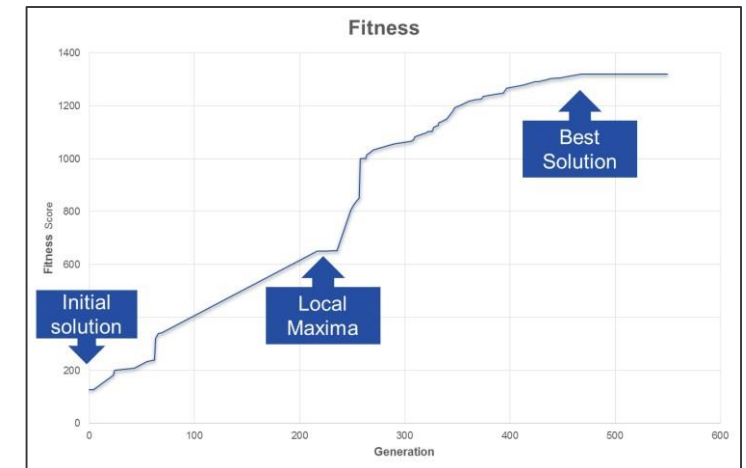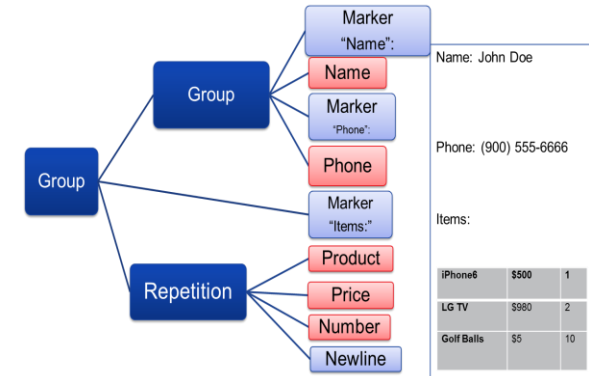
Thank You!

Questions?

# Appendix

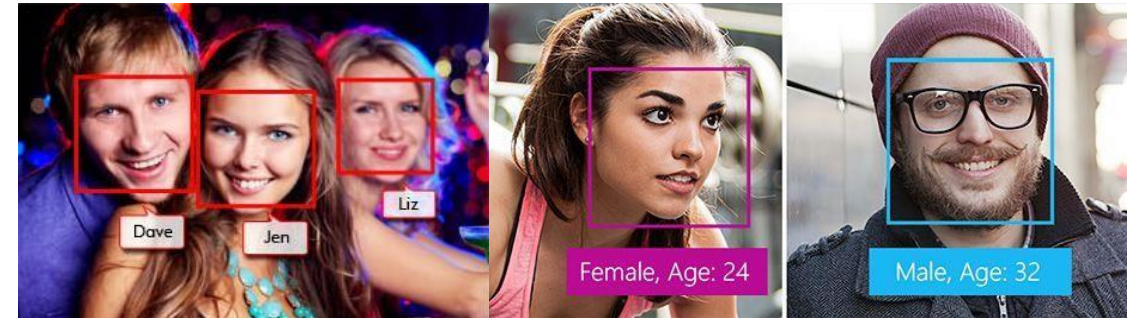# Additional Claire Details

# Decipher Data (schema extraction)

- High level analysis using A* based dynamic programming

- Genetic Algorithms to identify complex sub-structures

- Various NLP algorithms to modify model based on semantics
  - Identify text blocks that are not for parsing (comments, free text, etc)
  - Identifying patterns in the input
  - Element naming and semantics
  - Map between inputs and models

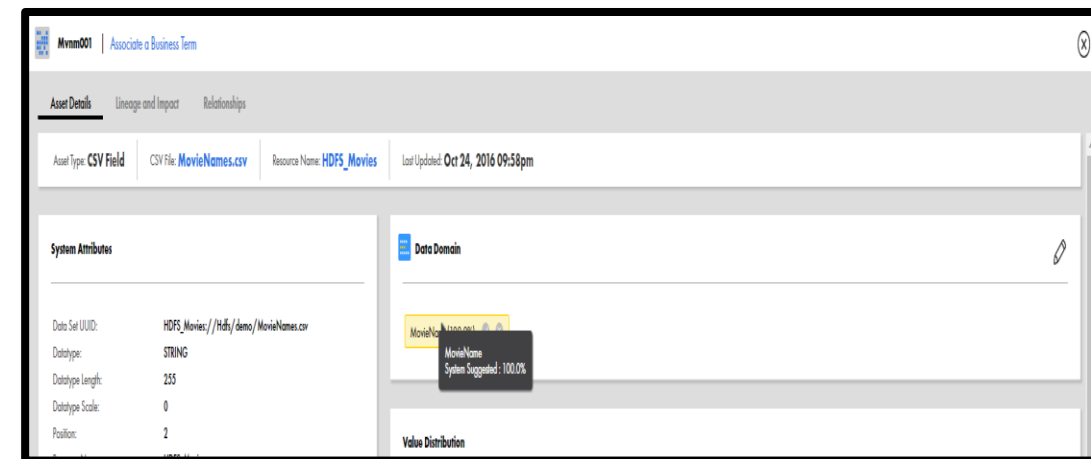Extendible with user and vertical specific types





© Informatica. Proprietary and Confidential.

Informatica

# Artificial Intelligence to Cluster Data

- Column Similarity based on Data Overlap

- Large Overlap of Distinct Values:

  - Jaccard distance = $1 - \dfrac{|S(X) \cap S(Y)|}{|S(X) \cup S(Y)|}$

- Similar Value Frequencies for overlapping columns

  - Bray Curtis Similarity: $\dfrac{\sigma_{i=1}^{n}}{n} \dfrac{|X_i - Y_i|}{X_j + Y_j}$

- Clustering based on Column metadata and Jaccard Coefficient and then computing Bray Curtis Similarity.



Like photo tagging
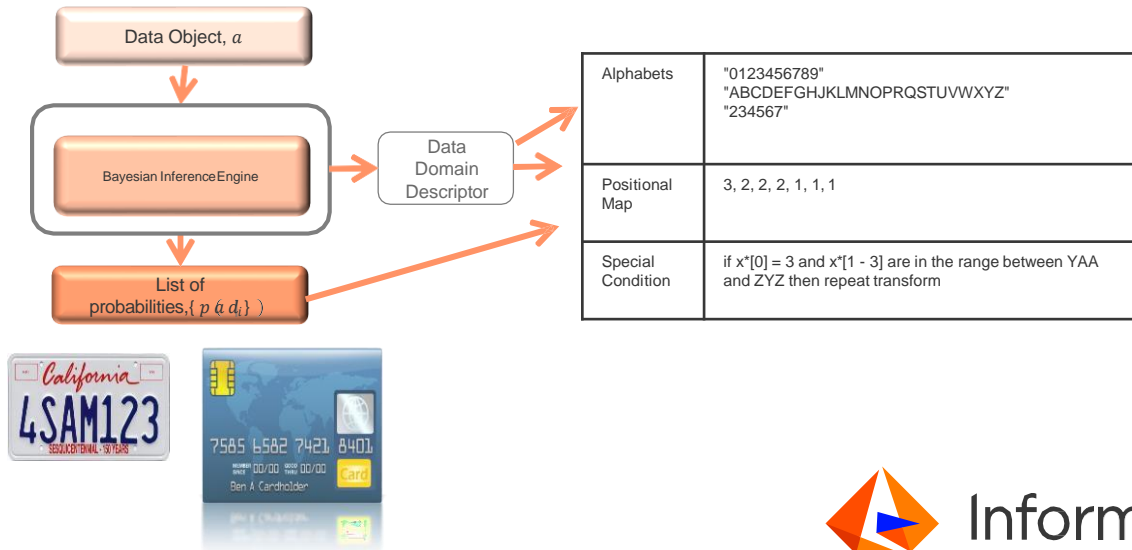CL*AI*RE for Columns

Informatica™

# Artificial Intelligence for Security Analytics

- Bayesian Inference for auto-morphism and format preserve masking

- UBA unsupervised machine learning combined with Principal component analysis to create multi-dimensional model of user activities

- BIRCH technique for unsupervised hierarchical clustering and to identify changes in user behavior

- Validation based on distance and density for outlier detection and Grubbs' test

The Grubbs' test statistic is defined as:

$$G = \frac{\max_{i=1,\dots,N} |Y_i - \bar{Y}|}{s}$$



| Alphabets | "0123456789" "ABCDEFGHJKLMNOPRQSTUVWXYZ" "234567" |
|---|---|
| Positional Map | 3, 2, 2, 2, 1, 1, 1 |
| Special Condition | if x*[0] = 3 and x*[1 - 3] are in the range between YAA and ZYZ then repeat transform |

Informatica™

# Artificial Intelligence to Extract Entities

- NLP techniques to identify and extract data entities from strings
  - Extract Product Code from product descriptions
  - Identify Organization vs. Person information
  - Extract entities from unstructured Data
- Use Classifier Transform (Mallet from UMASS) to categorize data based on a custom classification model
- Statistical algorithms identify common and uncommon elements of your data