# Delivering Trusted Data: Collibra's Modern Approach To Data Quality

Henry Tram, Lead Sales Engineer – Data Quality
08/10/2022

# **Agenda**

- Trends in Data Quality

- Collibra Data Quality & Observability (User Journey)

- Product Demonstration

- Q&A

Collibra

# How would you manage this company's data?

**Consider the following features as we demonstrate them in action…**

### The data is created by your front office…

- Sales representatives record the transactions
- They have flexibility in their own software

**Human Error**

### The data is in multiple places…

- Your "controls" separated the data
- You have sales data and financial transactions

**Technical Error**

### Everyone up to CEO needs that data…

- What could possibly go wrong?
- How will you discover it?

**No Error**

Collibra

# Problem:  This leads to real costs across your organization

## $1.9B
projected data quality **spend** in 2021.

## 50-70%
average **time spent** on manual rule writing and management.

## 47%
of recently created data records have at least one **critical error**.

## 15-25%
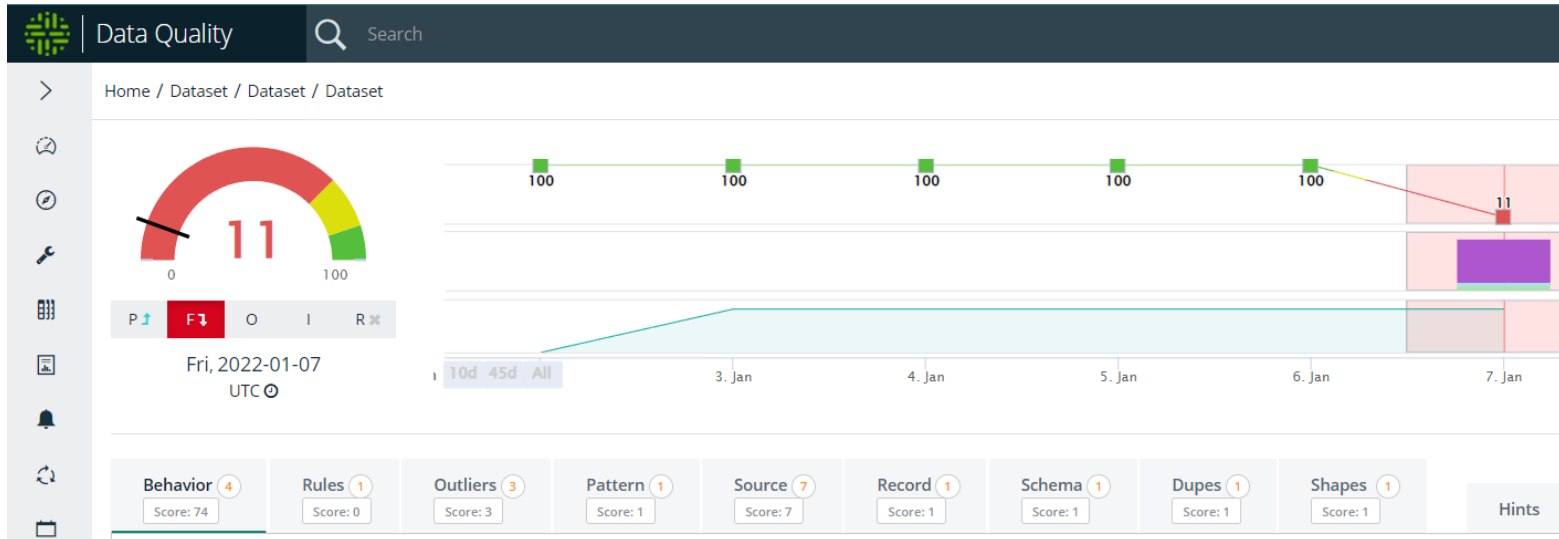**lost revenue** for most companies due to bad data.

# Collibra Data Quality & Observability

- **Discover data quality issues**
- **"Feels" statistical, not code**
- **Scales as easy a identify, scan, discover, train**

## Intelligent

Auto-generate explainable and adaptive DQ rules

## Self-service

Empower users with a unified scoring system and flexible rule management

## Scalable

Scan large and diverse databases, files and streaming data

Collibra

# Collibra Data Quality & Observability Demonstration

## 1  Connect to Data Sources

Sources → Connector → Collibra DQ

.csv

snowflake

Amazon S3

ORACLE

JSON

SQL Server

MySQL

PostgreSQL

Parquet

Amazon Redshift

IBM DB2

HIVE

## 2  Create Dataset & View Findings

## 3  DQ Core Components

Collibra

# Collibra Data Quality & Observability Demonstration



## 1 Connect to Data Sources

**Sources** → **Connector** → **Collibra DQ**

- .CSV
- snowflake
- Amazon S3
- ORACLE
- JSON
- SQL Server
- MySQL
- PostgreSQL
- Parquet
- Amazon Redshift
- IBM DB2
- HIVE

## 2 Create Dataset & View Findings

**Create & Configure Datasets** → **Verify Data Profiles** → **View Findings (Validate Anomalies)**

### Dataset DQ Findings
- Behavior
- Outliers
- Data Rules
- Patterns
- Source to Target Validation
- Shapes
- Schema Changes
- Dupes Fuzzy Matching
- Data Profiling
- ML/AI

## 3 DQ Core Components

# Collibra Data Quality & Observability Demonstration

APACHE Spark

## 1 Connect to Data Sources

Sources → Connector → Collibra DQ

- .CSV
- snowflake
- Amazon S3
- ORACLE
- JSON
- SQL Server
- MySQL
- PostgreSQL
- Parquet
- Amazon Redshift
- IBM DB2
- HIVE

## 2 Create Dataset & View Findings

Create & Configure Datasets → Verify Data Profiles → View Findings (Validate Anomalies)

**Dataset DQ Findings**

- Behavior
- Outliers
- Data Rules
- Patterns
- Source to Target Validation
- Shapes
- Schema Changes
- Dupes Fuzzy Matching

Data Profiling | ML/AI

## 3 DQ Core Components

- Data Quality Rules
- Adaptive Rules
- Rule Discovery
- Alerts & Scheduling
- REST/Scala/Py4J API
- PII & Data Detection
- Reporting & Metastore

Collibra

Collibra Data Quality & Observability Demonstration

# Real data, real use case… Sales Data

| Id | Date | Time | Daily_ID | First_Name | Last_Name | Email | Vendor_Type | Cost | Cost_Code | Cost_Description | Sales Rep | Sale_State | State_Tax |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | 2022-01-03 | 8:00 AM | 1 | Waly | Measor | wmeasoro@wordpress.org | Electrician | 1000.53 | 16-000 | General Electrical | Kal | DE | 0.000 |
| 1001 | 2022-01-04 | 8:00 AM | 1 | Waly | Measor | wmeasoro@wordpress.org | Electrician | 1275.12 | 16-000 | General Electrical | John | DE | 0.000 |
| 2001 | 2022-01-05 | 8:00 AM | 1 | Waly | Measor | wmeasoro@wordpress.org | Electrician | 1378.11 | 16-000 | General Electrical | Sara | DE | 0.000 |
| 3001 | 2022-01-06 | 8:00 AM | 1 | Waly | Measor | wmeasoro@wordpress.org | Electrician | 1100.98 | 16-000 | General Electrical | John | DE | 0.000 |
| 4001 | 2022-01-07 | 8:00 AM | 1 | George | Malarachy | gerogegotyou@qq.com | Electrician | 1280.12 | 16-000 | General Electrical | John | DE | 0.000 |
| 4980 | 2022-01-07 | 3:50 PM | 3D4 | Marshal | Morrice | mmorriceic@comcast.net | Subcontractor | 5921.10 | 01-515 | Temporary Lighting | John | PA | 0.085 |
| 4981 | 2022-01-07 | 3:50 PM | 3D5 | Elane | Blas | eblasqy@blogtalkradio.com | Supervisor | 7221.23 | 01-517 | Temporary Telephone | John | CT | 0.110 |
| 4982 | 2022-01-07 | 3:51 PM | 3D6 | Barron | Danilovich | bdanilovich43@squarespace.com | Estimator | 7052.15 | 01-510 | Temporary Utilities | John | NY | 0.115 |
| 4983 | 2022-01-07 | 3:52 PM | 3D7 | Rogerio | Sappell | rsappell1w@miibeian.gov.cn | Electrician | 2199.64 | 07-200 | Thermal Protection-Insulation | John | NY | 0.115 |
| 4984 | 2022-01-07 | 3:52 PM | 3D9 | Holly | Sephton | hsephtonnp@hao123.com | Construction Expeditor | 4188.61 | 04-800 | Masonry Assemblies | John | PA | 0.085 |
| 4985 | 2022-01-07 | 3:53 PM | 3D0 | Dannel | Vannozzii | dvannozzii2h@bing.com | Subcontractor | 5146.13 | 01-630 | Product Substitution Procedures | John | DE | 0.000 |
| 4986 | 2022-01-07 | 3:53 PM | 3DB | Juliette | Boughen | jboughen4c@outlook.com | Electrician | 1954.03 | 04-500 | General Refractories | John | DE | 0.115 |
| 4988 | 2022-01-07 | 3:53 PM | 3DD | Francine | Todarini | ftodarini1o@aol.com | Engineer | 3807.93 | 02-200 | Site Preparation | John | NY | 0.115 |
| 4989 | 2022-01-07 | 3:54 PM | 3DE | Francine | Todarini | ftodarini1@aol.com | Engineer | 3807.93 | 02-200 | Site Preparation | John | NY | 0.115 |
| 4990 | 2022-01-07 | 3:54 PM | 3DF | Dorothea | Daltrey | ddaltreypd@1und1.de | Surveyor | 5745.87 | 1.630 | Product Substitution Procedures | John | PA | 0.085 |
| 4991 | 2022-01-07 | 3:55 PM | 3E0 | Lorelei | Walkey | lwalkeybw@paginegialle.it | Engineer | 1323.83 | 01-530 | Temporary Construction | John | NJ | 0.090 |
| 4994 | 2022-01-07 | 3:57 PM | 3E3 | Flinn | Skocroft | fskocroftr8@icio.us | DOCTOR | 8548.72 | 02-311 | Final Grading | John | NY | 0.115 |
| 4996 | 2022-01-07 | 3:58 PM | 3E5 | Katerine | McTaggart | kmctaggarteq@apple.com | Subcontractor | 6225.50 | 13-500 | Recording Instrumentation | John | CT | 0.110 |
| 4997 | 2022-01-07 | 3:59 PM | 3E6 | Desiree | Sylett | dsylett4k | Construction Worker | 2618.29 | 02-822 | Ornamental Metal Fences and Gates | John | NYY | 0.115 |
| 4998 | 2022-01-07 | 3:59 PM | 3E7 | Lorelei | Walkey | lwalkeybw@paginegialle.it | Supervisor | 239100 | 15-400 | Plumbing Fixtures and Equipment | John | NY | 0.115 |
| 4999 | 2022-01-07 | 3:59 PM | 3E8 | Cletis | Pennington | cpenningtonnv@deliciousdays.com | Surveyor | 0 | 13-120 | Pre-Engineered Structures | John | PA | 0.085 |

Collibra

13

# Use Cases For Large Data Processing

**150 Million rows, 31 columns,** *in 45 minutes*
- 18 Spark Executors, 12GB RAM, 4 CPUs per Executor and 4GB RAM per Executor
- Success!

**44 Million rows, 495 columns,** *in 79 minutes*
- 46 Spark Executors, 24GB RAM, 8 CPUs per Executor and 4GB RAM per Executor
- Success!

**65 Million rows, 270 columns,** *in 4 hours*
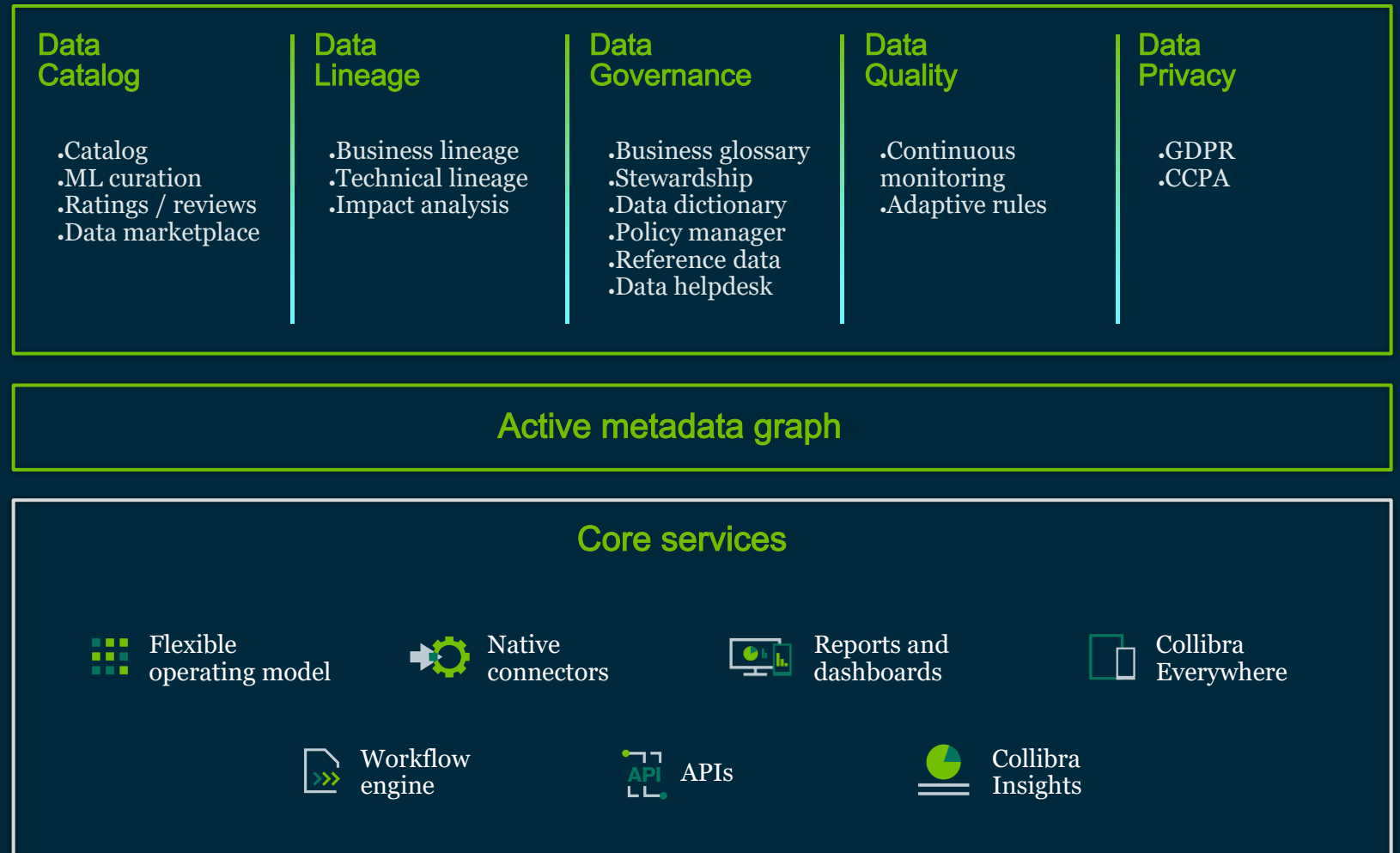- 10 Spark Executors, 25GB RAM, 8 CPUs per Executor and 4GB RAM per Executor
- Success!

**65 Million rows, 100 columns,** *in 1.5 hours*
- AWS EMR, 5 nodes, 20 Executors, 4 CPUs per Executor and 9GB RAM per Executor
- Success!

Collibra

# Data quality is core to data intelligence

## Govern. Trust. Access.

**Collibra Data Intelligence Cloud**

| **Data Catalog** | **Data Lineage** | **Data Governance** | **Data Quality** | **Data Privacy** |
|---|---|---|---|---|
| •Catalog<br>•ML curation<br>•Ratings / reviews<br>•Data marketplace | •Business lineage<br>•Technical lineage<br>•Impact analysis | •Business glossary<br>•Stewardship<br>•Data dictionary<br>•Policy manager<br>•Reference data<br>•Data helpdesk | •Continuous monitoring<br>•Adaptive rules | •GDPR<br>•CCPA |

### Active metadata graph

### Core services

- Flexible operating model
- Native connectors
- Reports and dashboards
- Collibra Everywhere
- Workflow engine
- APIs
- Collibra Insights

# Q&A

Remember to submit your questions!

# Thank you!

Collibra