# Thank you to our Fabric February Friends!

twoday

bouvet

sopra steria

DATAmasterminds

ADVANCING ANALYTICS

Evidi

Profisee
Master Data Management

Tabular Editor

KURANT

Fraktal

CluedIn

Dufrain
THE DATA COMPANY

#FabricFebruary

fabric FEBRUARY

# Wolfgang Strasser

Data Juggler

**ACP CUBIDO**



/in/wolfgang-strasser

workingondata.wordpress.com

@WolfgangStrasser

Hi, my name is Wolfgang

# A Day in the Life of a Data Worker..

# Sounds about right...

**infoCleanse**

| Year | Number of Bad Records | Number of Clean Records |
|------|----------------------|------------------------|
| 2015 | 20,000 | 80,000 |
| 2016 | 28,000 | 112,000 |
| 2017 | 39,200 | 1,56,000 |

■ **Number of Bad Records**

■ **Number of Clean Records**

Keeping the **$100 Per Dirty Record** in mind, and using a **100,000-Record Database**, here ist the astronomical cost of bad data over three years.

On average, corporate data grows at **40%** per year.

Approximately **20%** of the average database is dirty

**COST OF BAD RECORDS**

$2,000,000 — 2015

$2,800,000 — 2016

$3,920,000 — 2017

# $ 3.1T

est. cost of bad data for US businesses

# 12%

average loss in revenue due to bad data

# 60%

of organizations do not measure the financial cost of poor data quality

https://hbr.org/2016/09/bad-data-costs-the-u-s-3-trillion-per-year
https://www.gartner.com/smarterwithgartner/how-to-stop-data-quality-undermining-your-business

# Some Risks of Bad Data (Quality)

- Poor decision making

- Business inefficiencies

- Missed opportunities

- Siloed information – Redundancy

- Mistrust

- Lost Reputation

- Lost Revenue

# What is Data Quality?

# Fit-for-Use Purpose

# Fit-for-Use purpose

Data quality is often defined as the extent to which data is "fit for use" in a specific context. In this view, high-quality data is that which meets the needs of its users in operations, decision-making, and planning. This definition typically emphasizes that data must be:

- **Accurate** (reflecting the true values),

- **Complete** (containing all required information),

- **Consistent** (uniformly formatted across systems),

- **Timely** (up-to-date and available when needed), and

- **Relevant** (useful for the specific task or decision).

For example, a business using customer data to target marketing campaigns would expect the data to be correct, complete, and current to ensure that its campaigns reach the right audience.

# Representation and Reliability Perspective

# Definition 2: Representation and Reliability Perspective

Another common definition of data quality focuses on the degree to which data accurately represents real-world entities or events and is free from errors. Under this definition, data quality means that the information is not only correct but also reliable for analysis. It includes aspects such as:

- **Correctness:** Data values reflect actual conditions,

- **Duplication-free:** Each entity is represented only once (eliminating redundant records), and

- **Validity:** Data conforms to defined formats or rules.

In a business context, consider a logistics company whose routing decisions depend on location data. If the data is not an accurate reflection of actual geographic positions—or if duplicate or conflicting entries exist—the company might misroute deliveries, leading to increased costs and customer dissatisfaction.

Data Quality Dimensions

DQ

# Data Quality Dimensions

Accuracy

Completeness

Consistency

Timeliness

Validity

Uniqueness

# DQ Process(es)

DQ Planning & Definition

Data Profiling & Assessment

Data Cleansing & Transformation

Data Validation & Verification

Monitoring, Governance & Continuous Improvement

#FabricFebruary

# Data Quality in Microsoft Fabric

# Data Preparation in Microsoft Fabric

# Data Prep & DQ in Microsoft Fabric

# Demo Time

DQ in Microsoft Fabric

# Data Quality in Microsoft Fabric

Data Profiling — Data Wrangler, ydata (https://docs.profiling.ydata.ai/latest/)

Data Cleansing — pyspark, Copilot

Data Quality Checks — Great Expectations (LH, Model)

https://greatexpectations.io/


(Data Lineage)

(Endorsements)

# Data Quality with Microsoft Purview

# Microsoft Purview

## Data Security

**+**

## Data Governance

**+**

## Data Compliance

Secure your data
along its lifecycle

*Data Loss Prevention*
*Insider Risk Management*
*Information Protection*
*Adaptive Protection*

Unlock value creation
from your data

*Discovery and Access*
*Catalog Management*
*Health Management*
*Master Data Management\**

Manage critical risks and
regulatory requirements

*Compliance Manager*
*eDiscovery and Audit*
*Communication Copliance*
*Data Lifecycle Management*
*Records Management*

### Shared platform capabilities

*\* Through*
*3rd Party Integration*

### Structured, Semi-structured and Unstructured Data

### Along the Data Lifecycle

# Data Quality in Microsoft Purview



#FabricFebruary

# Data Quality in Microsoft Purview (2025-02)

https://learn.microsoft.com/en-us/purview/data-quality-overview

https://learn.microsoft.com/en-us/purview/data-quality-for-fabric-data-estate

## Supported multicloud data sources

- Azure Data Lake Storage (ADLS Gen2)
  - File Types: Delta Parquet and Parquet
- Azure SQL Database
- Fabric data estate in OneLake including shortcut and mirroring data estate. Data quality scanning is supported only for Lakehouse delta tables and parquet files.
  - Mirroring data estate: Cosmos DB, Snowflake, Azure SQL
  - Shortcut data estate: AWS S3, GCS, AdlsG2, and dataverse
- Azure Synapse serverless and data warehouse
- Azure Databricks Unity Catalog
- Snowflake
- Google Big Query (Private Preview)
- Iceberg data in ADLS Gen2, Microsoft Fabric Lakehouse, AWS S3, and GCP GCS

#FabricFebruary

# Data Quality Rules

# Data Quality Scans



Business Domain
Data Product
Rules per Asset

# Health Management



#FabricFebruary

# Data Quality Reporting



#FabricFebruary

# Think about Data Quality along the life-cycle

Where are you with DQ today?
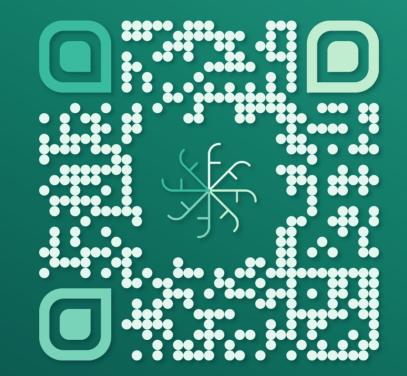
Get-to-know your data first

Get-to-know your data people (data owners)

Incrementally add more rules

Track data quality success

DQ Planning &
Definition

Data Profiling &
Assessment

Monitoring, Governance
& Continuous
Improvement

Data Validation &
Verification

Data Cleansing &
Transformation

Share **your thoughts** and help our speakers!

**fabfeb.app/feedback**

#FabricFebruary

# Thank you!

fabric
**FEBRUARY**