# Data quality

**5ARB0: DATA ACQUISITION & ANALYSIS (2022 – 2023)**

**Uzay Kaymak, Jheronimus Academy of Data Science, u.kaymak@tue.nl**

Mastertrack: Artificial Intelligence & Engineering Systems

1

# Outline

- **Data governance and assurance**

- **Data quality**
  - Definition
  - Dimensions
  - Challenges
  - Tools

- **AHIMA data quality model (healthcare example)**
  - Data quality management functions
  - Characteristics of data quality

2

TU/e

## Data governance

Discipline of administering data and information assets across an organization through formal oversight of
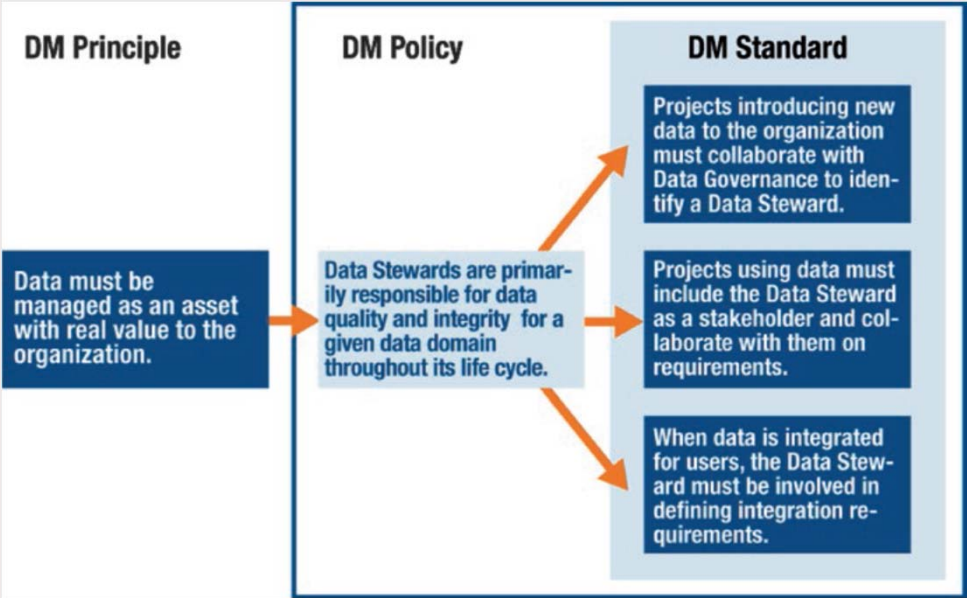
- the people,

- processes,

- technologies, and

- lines of business

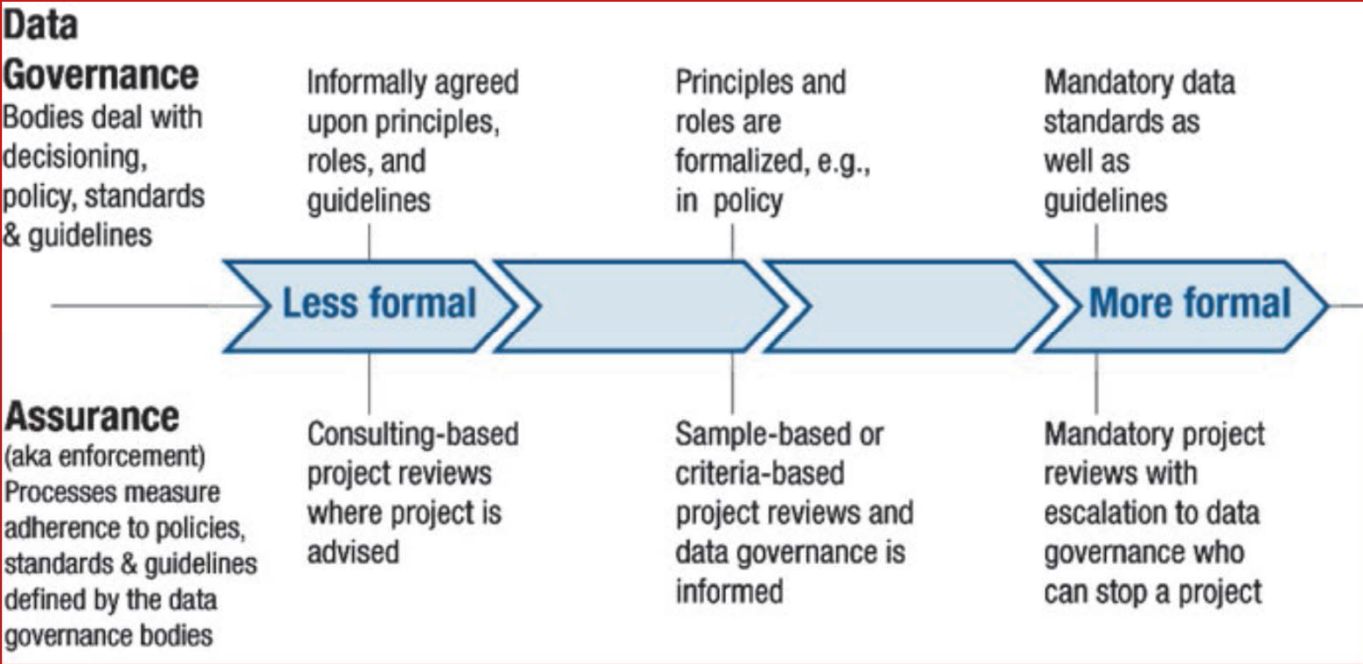that influence data and information outcomes to drive business performance

Orr (2011)

3

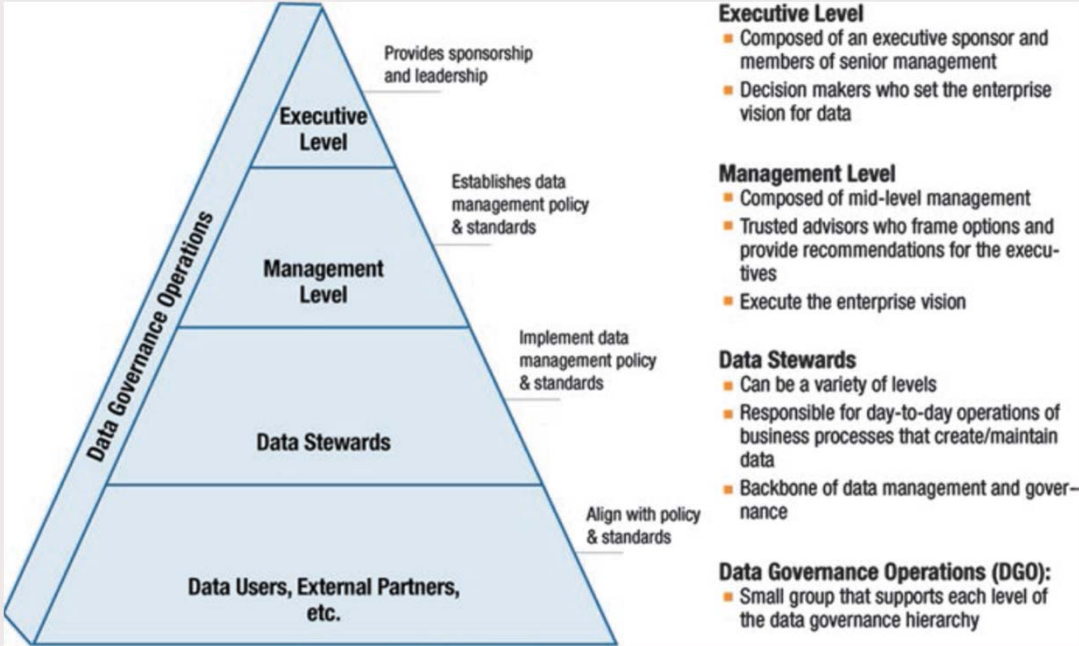**TU/e**

# Organizational principles, policies and standards

Data governance is
a management function



| DM Principle | DM Policy | DM Standard |
|---|---|---|
| Data must be managed as an asset with real value to the organization. | Data Stewards are primarily responsible for data quality and integrity for a given data domain throughout its life cycle. | Projects introducing new data to the organization must collaborate with Data Governance to identify a Data Steward. |
| | | Projects using data must include the Data Steward as a stakeholder and collaborate with them on requirements. |
| | | When data is integrated for users, the Data Steward must be involved in defining integration requirements. |

4

**TU/e**

# Data governance and assurance

**Data Governance**
Bodies deal with decisioning, policy, standards & guidelines

Informally agreed upon principles, roles, and guidelines

Principles and roles are formalized, e.g., in policy

Mandatory data standards as well as guidelines

Less formal → → → More formal

**Assurance**
(aka enforcement) Processes measure adherence to policies, standards & guidelines defined by the data governance bodies

Consulting-based project reviews where project is advised

Sample-based or criteria-based project reviews and data governance is informed

Mandatory project reviews with escalation to data governance who can stop a project

5

**TU/e**

# Data governance hierarchy

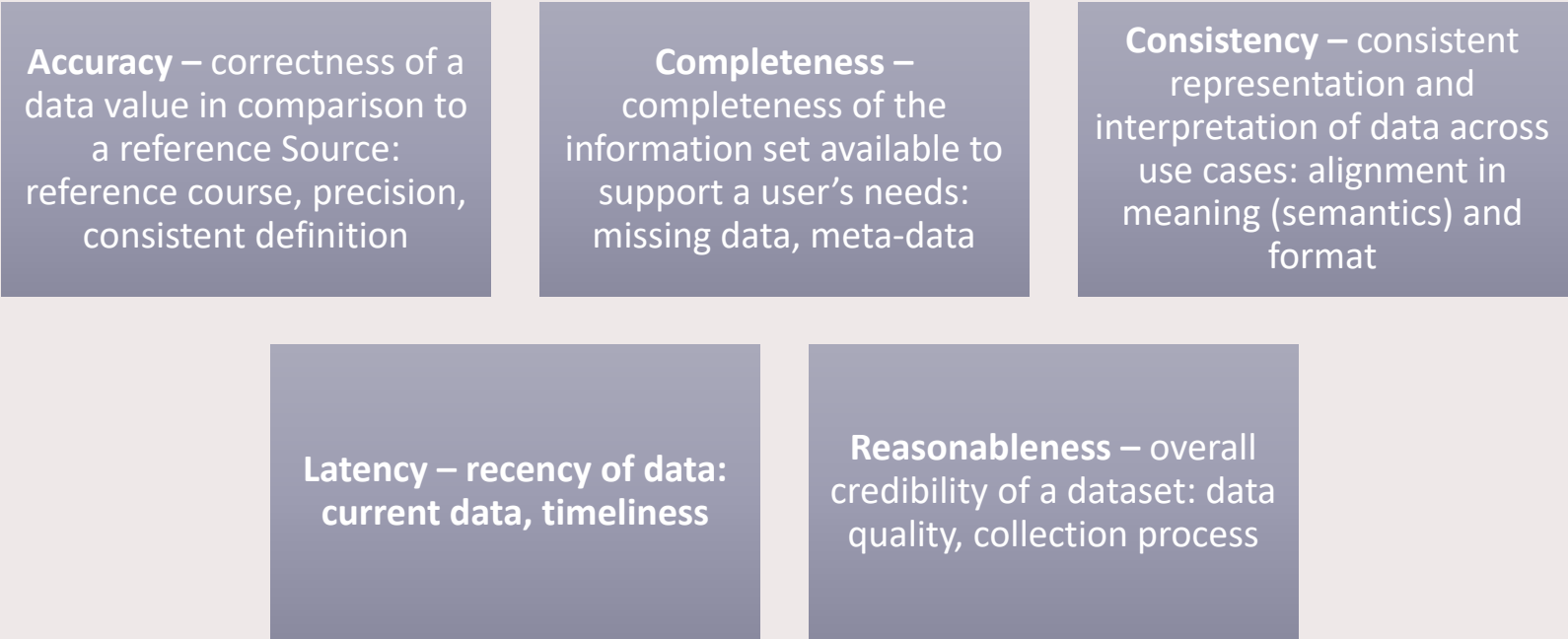U. Kaymak                    6

# Data quality

Data are considered high quality if "they are fit for their intended uses in operations, decision making and planning" for the business

Data quality management activities enable the availability of data that is "fit for use" for the business user
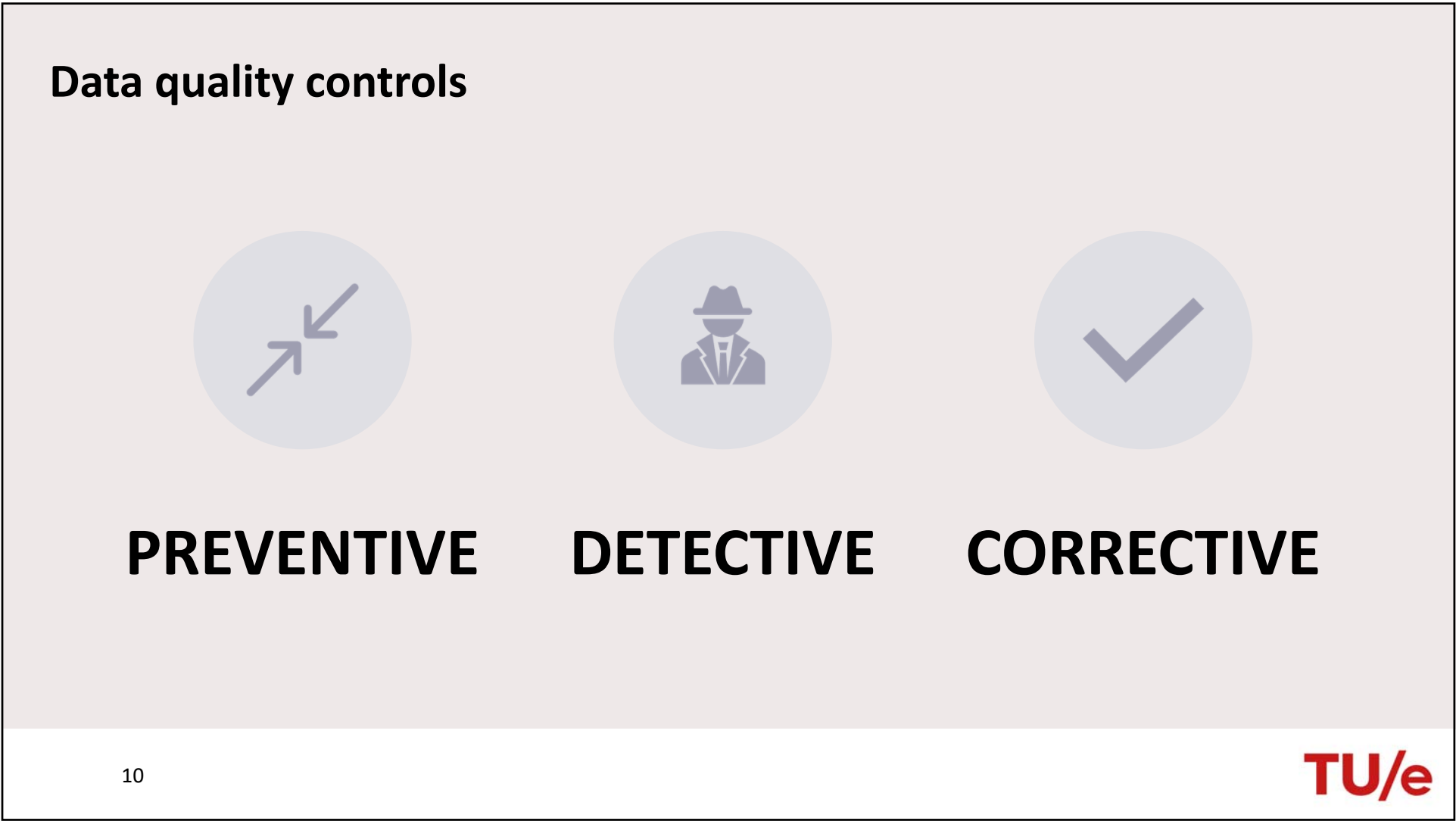
7

**TU/e**

## Data quality dimensions

**Accuracy –** correctness of a data value in comparison to a reference Source: reference course, precision, consistent definition

**Completeness –** completeness of the information set available to support a user's needs: missing data, meta-data

**Consistency –** consistent representation and interpretation of data across use cases: alignment in meaning (semantics) and format

**Latency – recency of data: current data, timeliness**

**Reasonableness –** overall credibility of a dataset: data quality, collection process

8

**TU/e**

## Data quality challenges

- **Inadequate controls at the point of origin**

- **Volume, variety, velocity**

- **Environment complexity**

- **Extensive proliferation and duplication**

- **Poor metadata, unclear definitions, multiple interpretations**



9

**TU/e**

# Data quality controls



## PREVENTIVE        DETECTIVE        CORRECTIVE

10

TU/e

10

# Common functions in data quality tools

**Profiling and metrics**

**Standardization and normalization**

**Identity resolution**

**Monitoring**

**Issue resolution and workflow**

**Cleaning and enhancement**

11

TU/e

# AHIMA Data Quality Model



FIGURE 2.2.    AHIMA Data Quality Management Model

**Characteristics of Data Quality**

- Accessibility
- Consistency
- Currency
- Granularity
- Precision
- Accuracy
- Comprehensiveness
- Definition
- Relevancy
- Timeliness

**Application** – The purpose for which the data are collected.

**Collection** – The processes by which data elements are accumulated.

**Warehousing** – Processes and systems used to archive data and data journals.

**Analysis** – The process of translating data into information utilized for an application.

*Source*: AHIMA, Data Quality Management Task Force, 1998.

# AHIMA DQM – data quality management functions

**Application:** The purpose for the data collection

**Collection:** The processes by which data elements are accumulated

**Warehousing:** Processes and systems used to archive data

**Analysis:** The process of translating data into meaningful information

13

## AHIMA DQM – characteristics of data quality (1)

**Accuracy:** The extent to which the data are free of identifiable errors

**Accessibility:** The level of ease and efficiency at which data are legally obtainable, within a well protected and controlled environment

**Comprehensiveness:** The extent to which all required data within the entire scope are collected, documenting intended exclusions

**Consistency:** The extent to which the healthcare data are reliable, identical, and reproducible by different users across applications

**Currency:** The extent to which data are up-to-date; a datum value is up-to-date if it is current for a specific point in time, and it is outdated if it was current at a preceding time but incorrect at a later time

**TU/e**

14

# AHIMA DQM – characteristics of data quality (2)

**Definition:** The specific meaning of a healthcare-related data element

**Granularity:** The level of detail at which the attributes and characteristics of data quality in healthcare data are defined

**Precision:** The degree to which measures support their purpose, and/or the closeness of two or more measures to each other

**Relevancy:** The extent to which healthcare-related data are useful for the purposes for which they were collected

**Timeliness:** The availability of up-to-date data within the useful, operative, or indicated time

15

**TU/e**

# Exercise

16