

COMBINING DATA LAKE AND DATA WRANGLING FOR ENSURING DATA QUALITY IN CRIS

Otmane Azeroual^{1*},

Joachim Schöpfel²,

Dragan Ivanovic³,

Anastasija Nikiforova^{4,5}

(1) German Centre for Higher Education Research and Science Studies (DZHW), Germany

(2) GERiCO-Labor, University of Lille, France

(3) University of NoviSad, Serbia

(4) University of Tartu, Institute of Computer Science, Estonia

(5) European Open Science Cloud Task Force «FAIR metrics and data quality», Belgium

15th International Conference on Current Research Information Systems (CRIS2022)

Dubrovnik, Croatia, May 12-14, 2022

BACKGROUND AND MOTIVATION

Today, billions of data sources continuously generate, collect, process, and exchange data. With the rapid increase in the number of devices and information systems in use, the amount and variety of data are increasing. This is also the case for the research / scientific domain.

Researchers as the end-users of RIS should be able to integrate ever-increasing volumes of data into their institutional database such as Current Research Information Systems (CRIS), regardless of the source, format or amount/ size of research information, where the data quality, flexibility and scalability in connecting and processing different data sources are crucial.



an effective mechanism should be employed to ensure faster value creation from these data

BACKGROUND AND MOTIVATION

an effective mechanism should be employed to ensure faster value creation from these data



DATA LAKE + DATA WRANGLING

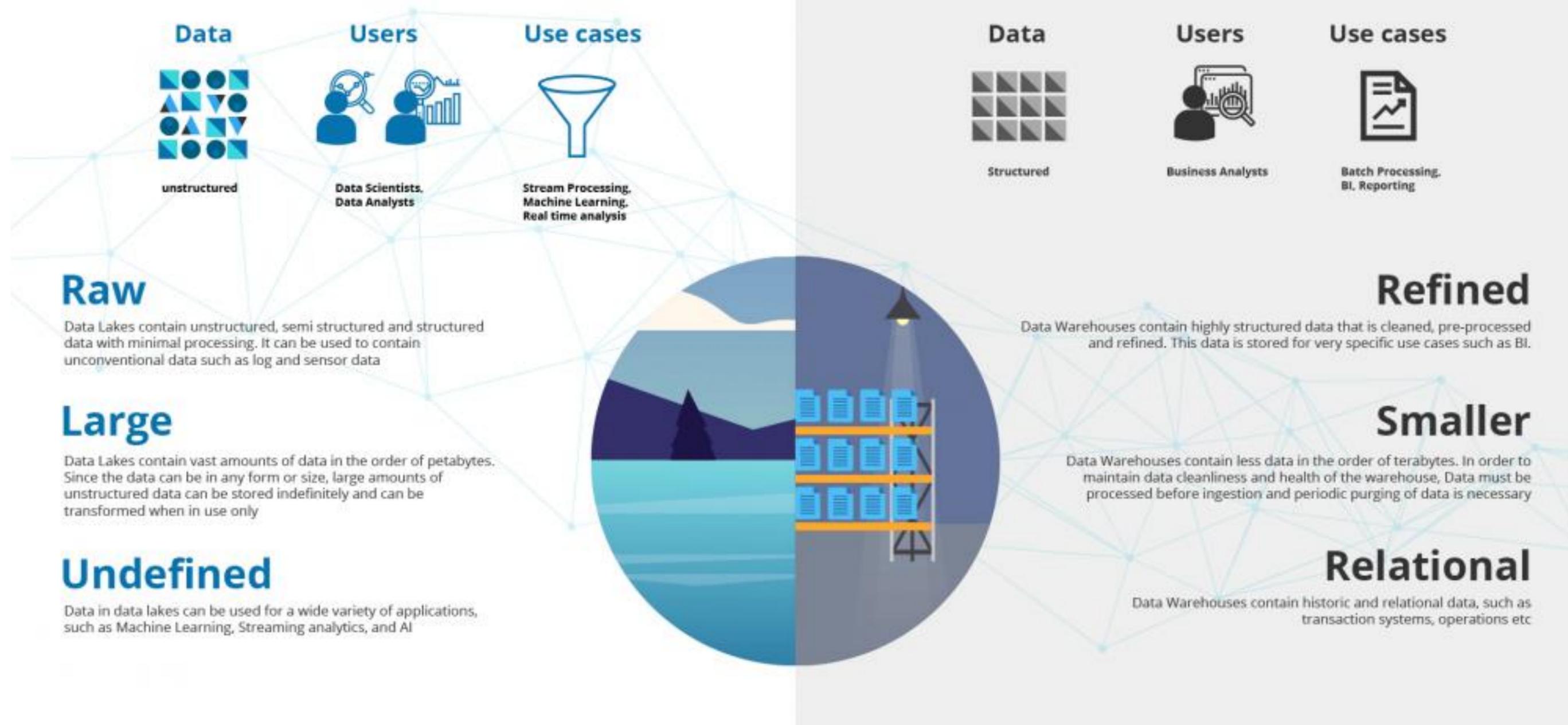
This study sets out the concept of a data lake with data wrangling process to be used in CRIS to clean up data from heterogeneous data sources as it is ingested and integrated.

DATA LAKE

DATA LAKE ✓

vs

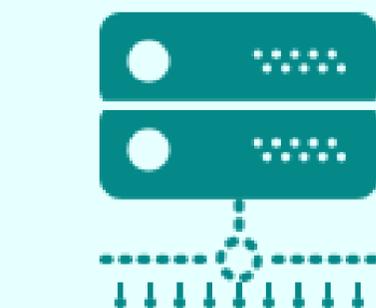
DATA WAREHOUSE



DATA WAREHOUSE

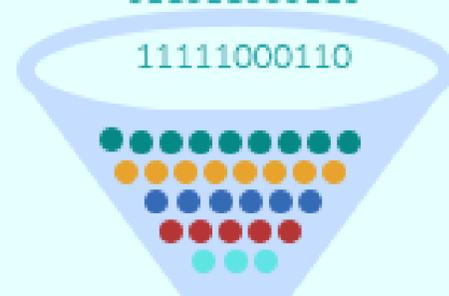
VS

DATA LAKE



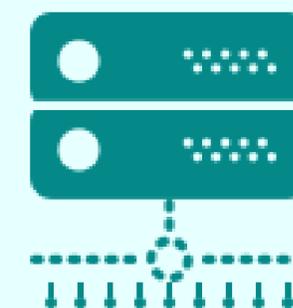
1110001101110
011011000110
11111000110

Data is processed and organized into a single schema before being put into the warehouse



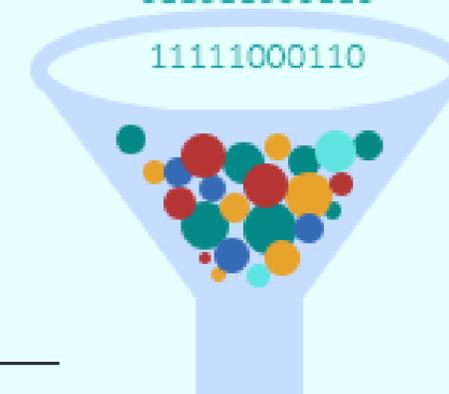
The analysis is done on the cleansed data in the warehouse

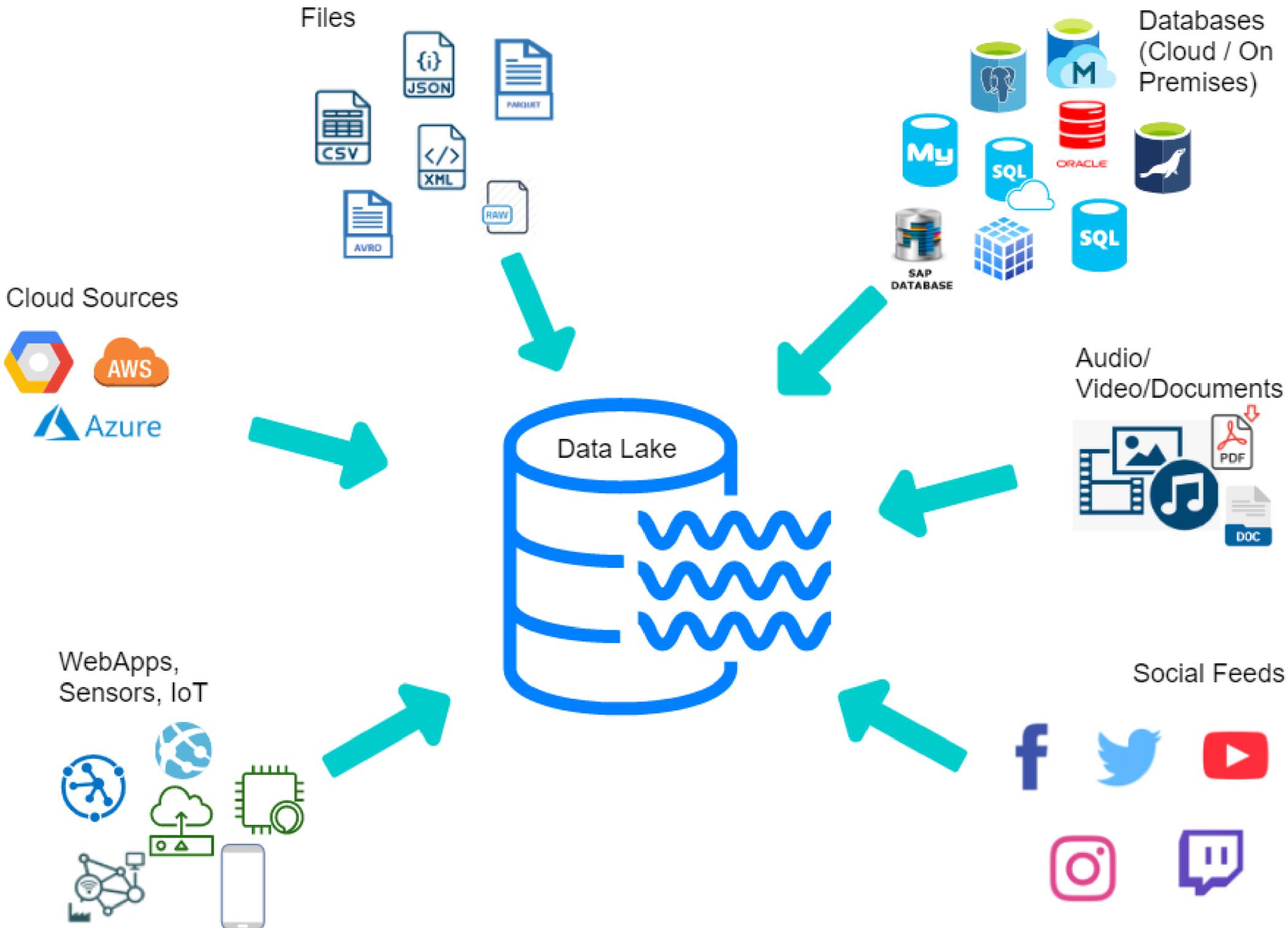
Raw and unstructured data goes into a data lake



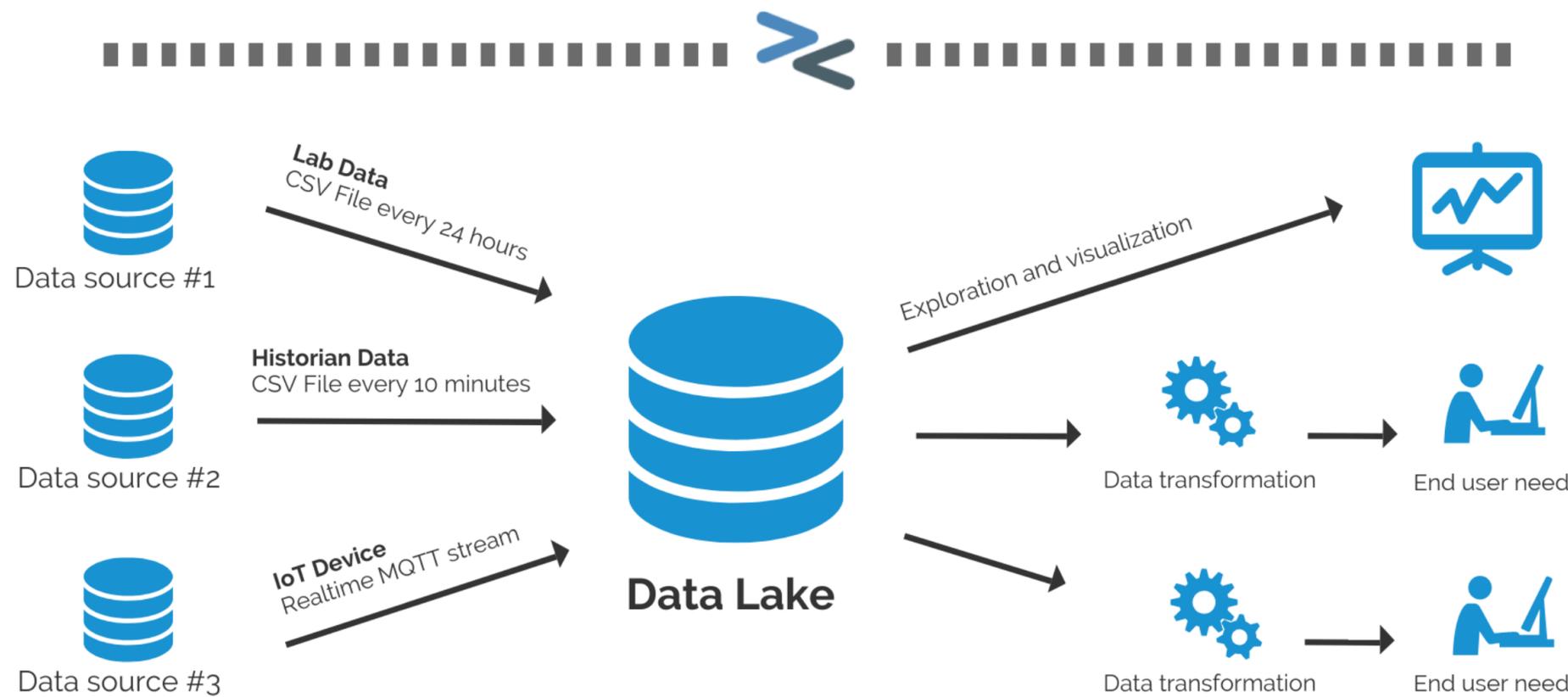
1110001101110
011011000110
11111000110

Data is selected and organized as and when needed





- **Data Lake provides a scalable platform for storing and processing large amounts of research data from various sources in their original raw format, regardless of their type, i.e., structured or unstructured data or text, numeric, images, video etc.**
- **The raw data are not cleaned, validated, or transformed → they are original data in their original format.**

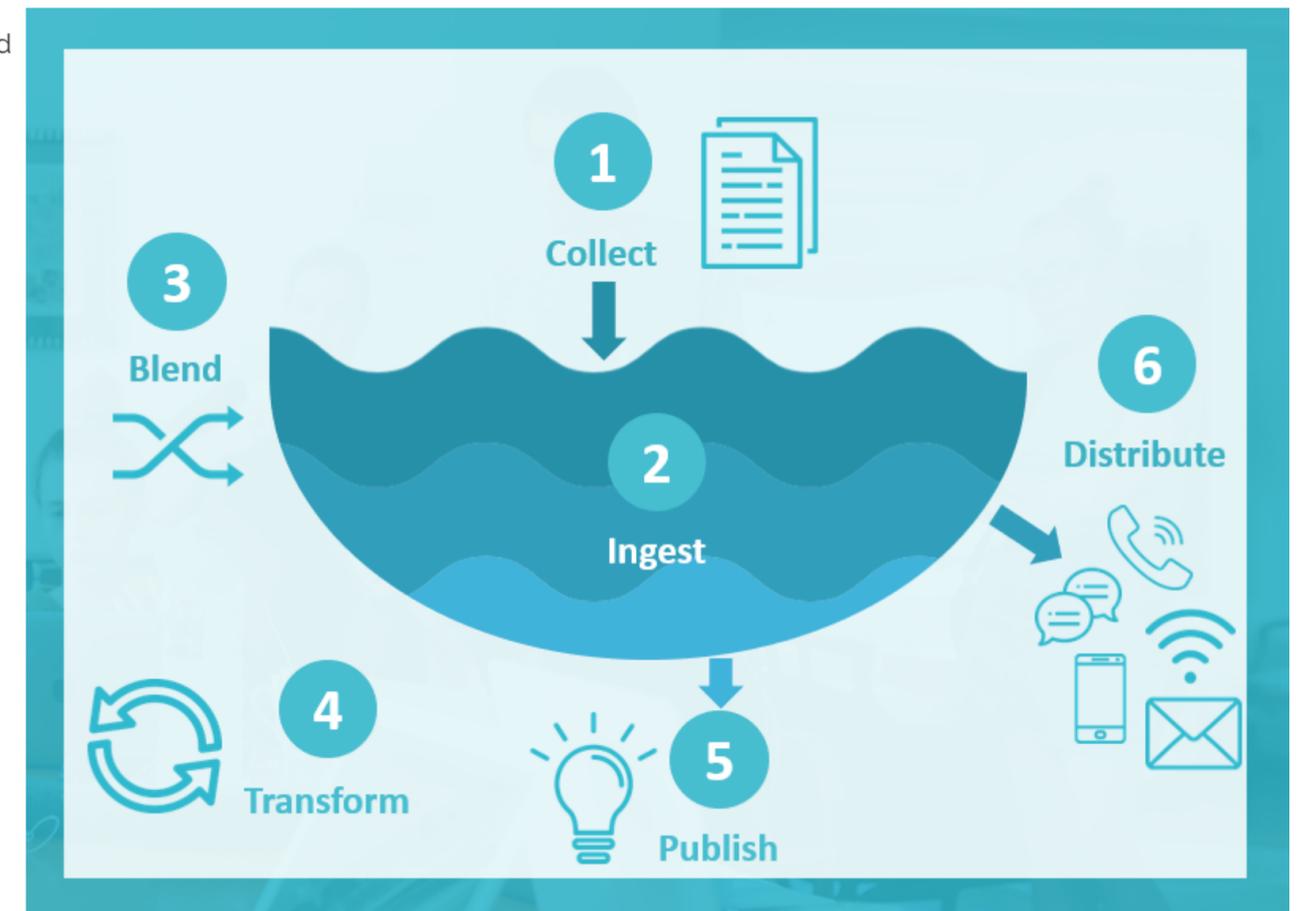


- When storing data / research information, the completeness of data and reduction of the cycle time between data generation and availability are important.
- The lack of pre-processing does not slow down data supply and does not lead to data loss.

- The concept of data lake allows to store a variety of data within the memory

BUT

- there is a need to clean up dirty data and enrich them in a pre-processing process, where data wrangling is found to be suitable for these purposes.
- The goal is to convert complex data types and data formats into structured data without programming efforts → users should be able to prepare and transform their research information without the need of using the ETL tools or familiarity and use of programming languages, where these transformations should be automatically suggested after reading the data based on machine learning algorithms that greatly speeds up this process.



DATA WRANGLING

DATA WRANGLING VERSUS DATA CLEANING

DATA CLEANING

Process of detecting and removing corrupted or inaccurate records from a record set, table or database

Data cleansing is another name for data cleaning

Visit www.PEDIAA.com

DATA WRANGLING

Process of transforming and mapping data from one raw data form into another form with the intent of making it more appropriate and valuable for various tasks

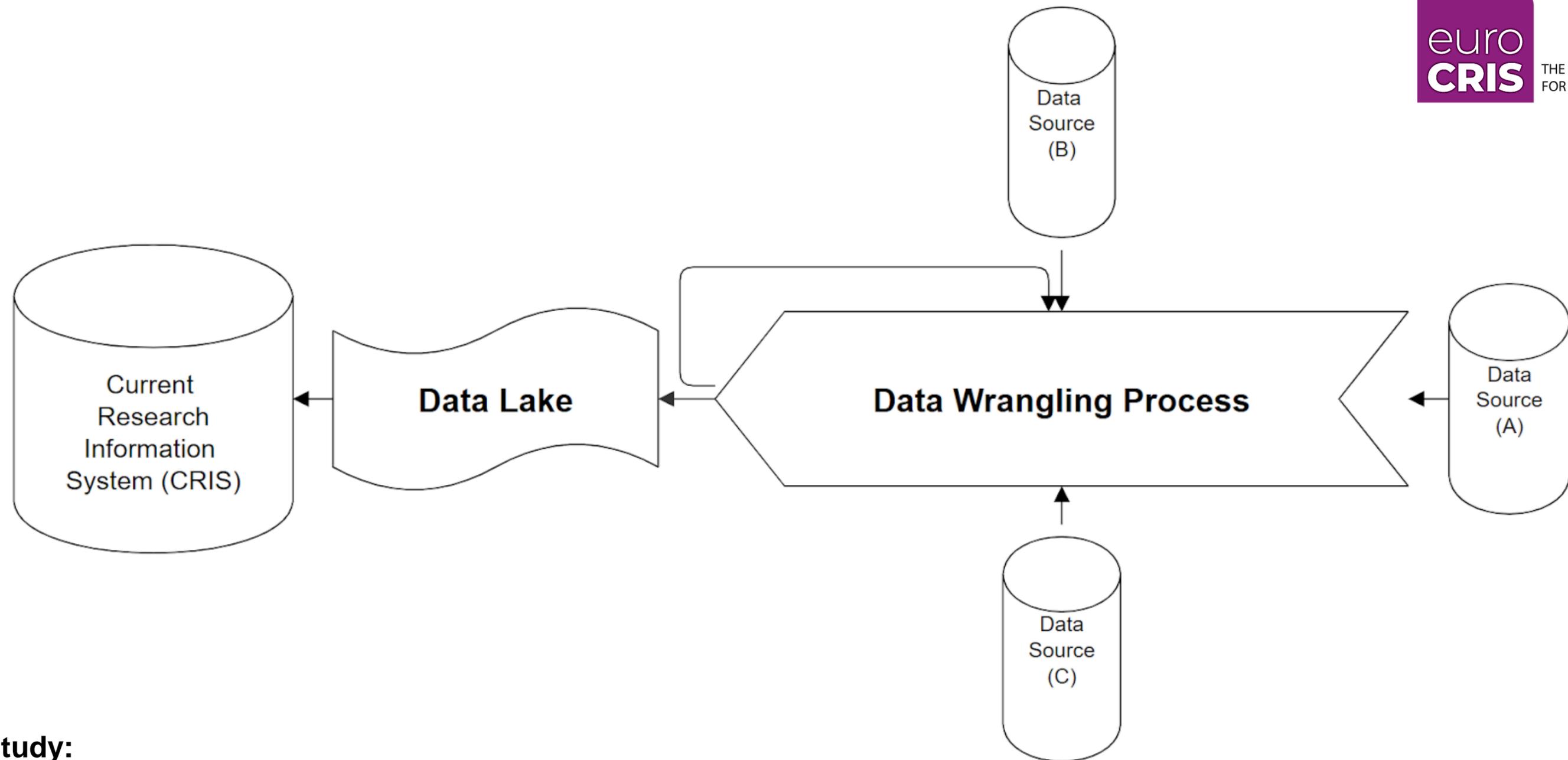
Data munging is another name for data wrangling



With the amount of data and data sources rapidly growing and expanding, it is getting increasingly essential for large amounts of available data to be organized for analysis.

Data wrangling is the process of cleaning and unifying messy and complex data sets for easy access and analysis.

**DATA LAKE + DATA WRANGLING
=
DATA QUALITY IN CRIS**



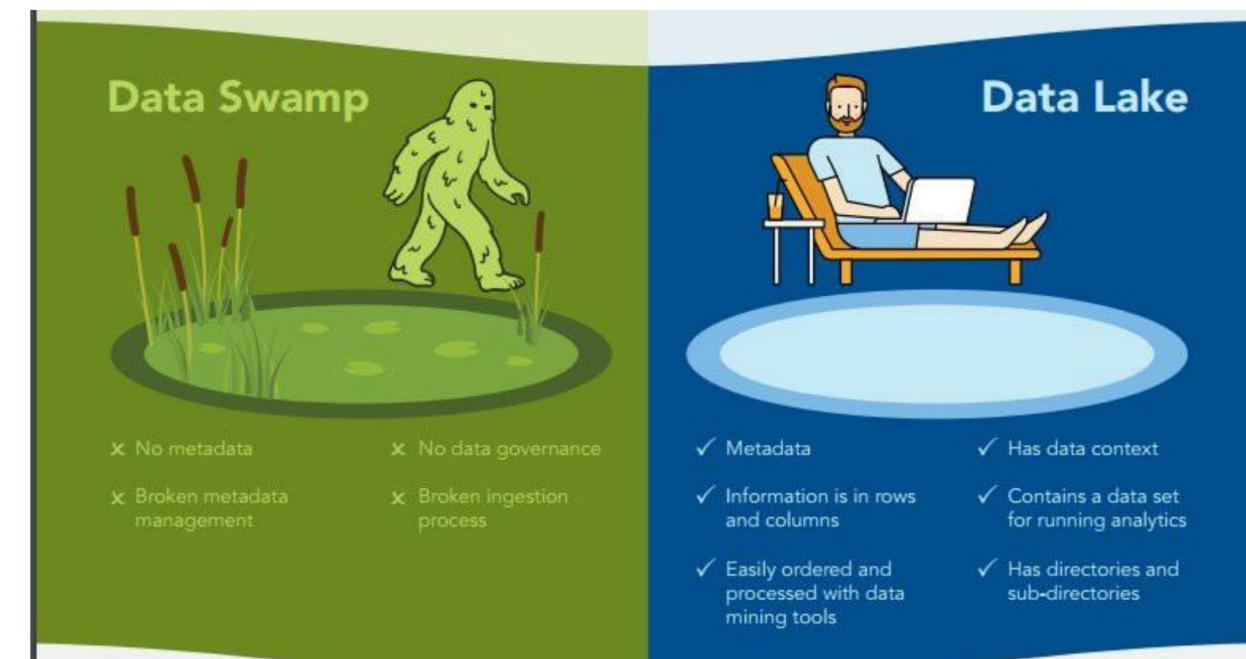
In this study:

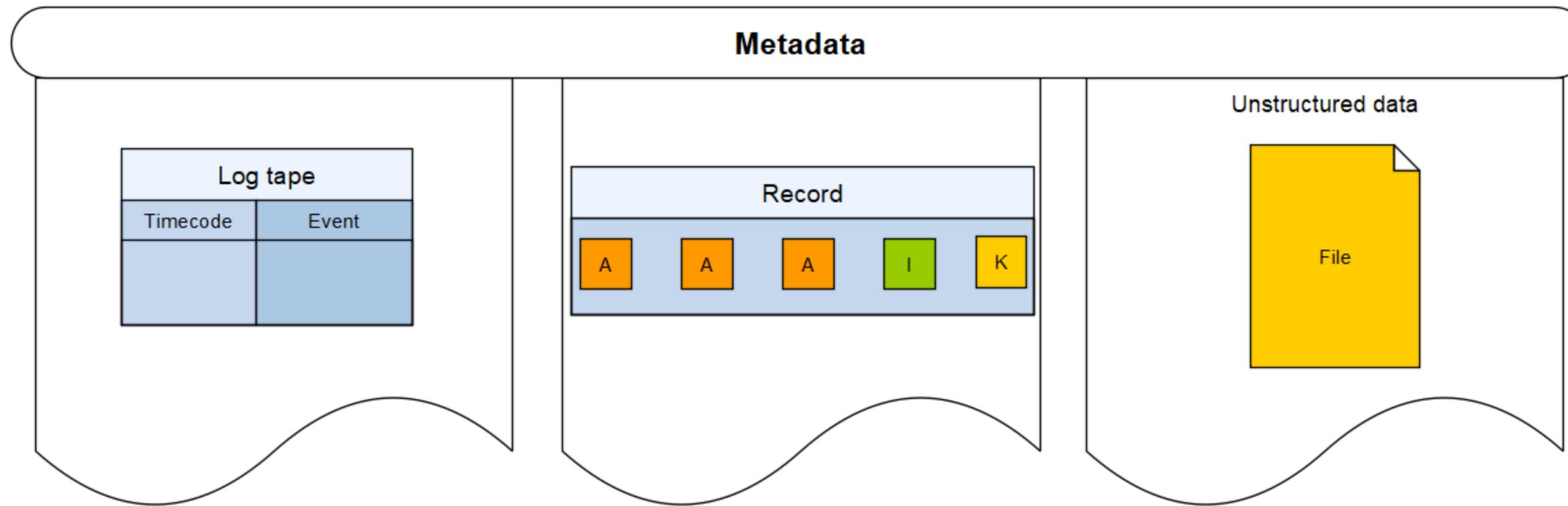
- an architectural model is first designed and specified, which analyzes the research information, adjusts it and transforms it into CRIS;
- a data lake makes both structured and unstructured data available in a reliable, trustworthy, secure and controlled way;
- the data wrangling process is used to verify and improve the quality of data, which also protects data from misuse → data are properly updated, retained, and eventually deleted according to the stage of its lifecycle.

The data wrangling process consists of several sequential steps. Depending on the IS and the desired or required target quality*, these individual steps should be carried out several times → data wrangling is a continuous process that repeats itself repeatedly at regular intervals.

SEVERAL ASPECTS AFFECTING A DATA LAKE

Aspect	Description
Metadata	describes a dataset in more detail containing <u>data about the origin, structure and content of the data</u> + <u>sorting, filtering or categorizing properties</u> + are used for system management and administration
Data mapping	describes the <u>context of the data</u> → integration map - a detailed specification of which application data from which data sources are linked / associated with which characteristics (mostly metadata)
Data lake context	describes the higher-level use case on which the data lake is based → the selection of the required data sources is more targeted. !!! This avoids the misuse of the data lake as a data swamp!!!
Data context	the individual datasets and their context so that they can be better classified for analysis purposes. The context for records can be data origin, categorization, or other contextual feature in the metadata
Processing logging	refers to the raw data processing that takes place in the data lake. The data record and its metadata are manipulated in the process → is of particular interest to data analysts to analyze data lake usage, data set and use case



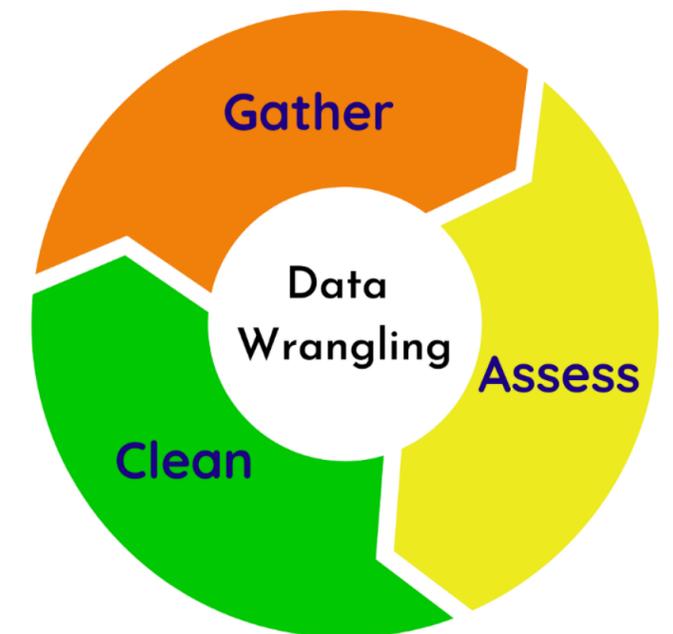


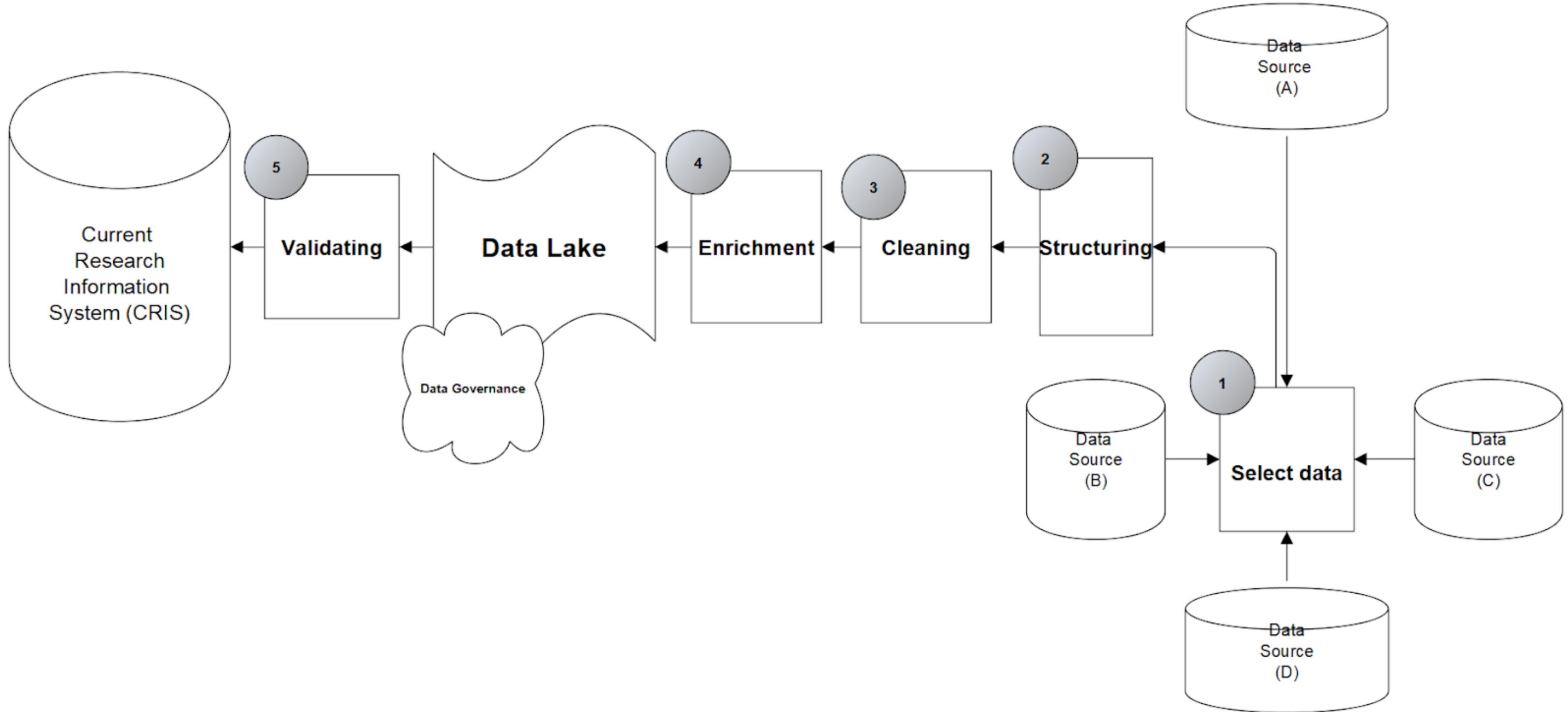
- **Analog data**: data sources automatically generate data in a specific predefined and therefore known data format. Due to the automatic generation, they accumulate in a very large amount and are mostly repeated / duplicated. For this reason, they are usually stored in tabular form in so-called "log tapes".
- **Application data**: also have a known structure but are significantly different from analog data in their origin - analog data typically represents physical measurement data, application data arises during the operation, transactions of an application (e.g. transmitted system data or analysis data). So-called "records" are used as a common storage solution for these data characterized by their uniform / homogeneous structure.

A data record usually consists of a key attribute K, an index attribute I, and other predefined attributes A. Depending on the data origin and data type of the application data, the predefined attributes may differ from each other. This application data structure is based on DBMS.
- **Text-based data** that are also closely related to the application, but are stored as separate files with metadata. A transformation is required to be carried out for further processing of this data. The process of converting them into analytically processable data is called textual disambiguation.

DATA WRANGLING

- In the context of research information, data wrangling refers to the process of identifying, extracting, preparing and integrating data into a database system such as CRIS.
- The use of data wrangling eliminates low-quality data, i.e. redundant, incomplete, inaccurate or incorrect data, etc., in order to preserve only high-quality research information from which the reliable and value-adding knowledge can be obtained.
- This adjusted research information is then entered into the appropriate target CRIS system to be used in further phases of the analysis (e.g. by analytical applications and protected from unauthorized access by access control).
- This should minimize the effort of analysis and enriching large volumes of data and metadata and achieve far-reaching added value in the procurement of information staff, developers and end users of the CRIS.



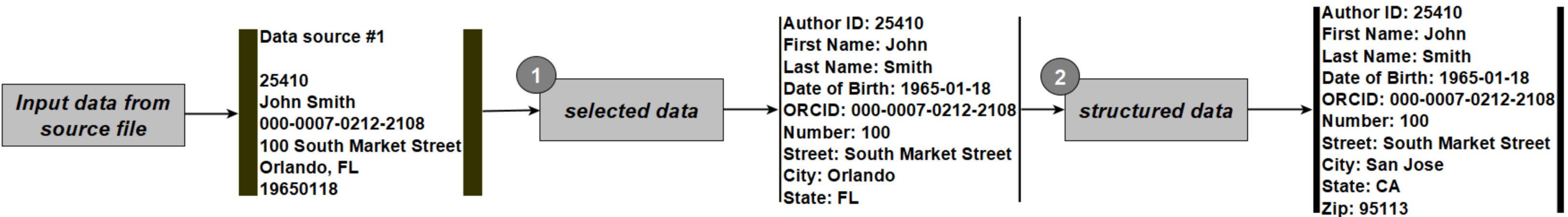


The data wrangling process (steps of the process are indicated by numbers) to prepare research information and integrate it into CRIS

Step	Description
Select data	The required data records are identified in different data sources. When selecting data, a record is evaluated by its value → if there is added value, the availability and terms of use of the data and subsequent data from this data source are checked
Structuring	In most cases, there is little or no structure in the data → change the structure of the data for easier accessibility.
Cleaning	Almost every dataset contains some outliers that can skew the analysis results → the data are extensively cleaned for better analysis (<i>processing of null values, removing duplicates and special characters, and standardization of the formatting to improve data consistency</i>)
Enrichment	<p>The data needs to be enriched - an inventory of the data set and a strategy for improving it by adding additional data should be carried out. The data set is enriched with various metadata:</p> <ul style="list-style-type: none"> ✓ Schematic metadata provide basic information about the processing and ingestion of data → the data wrangler analyzes / parses data records according to an existing schema. ✓ Conversation metadata are exchanged between accessing instances with the idea to document information obtained during the processing or analysis of these data for subsequent users. <p>The recognized peculiarities/ features of a data set can be saved.</p>
<i>*Data lake</i>	<p>The physical transfer of data in the data lake. Although data are prepared using metadata, the record is not pre-processed.</p> <p>The goal is to avoid a data swamp → estimate the value of the data and decide on their lifespan depending on the data quality and its interconnectedness with the rest of the DB.</p> <p>Analyzes are not performed directly in the data lake, but only on the relevant data. To be able to use the data, the requester needs the appropriate access rights → Data Wrangler performs data extraction, however, general viewing and exploration of the data should be possible directly in the data lake.</p>
<i>*Data governance</i>	The contents of the data lake, technologies and hardware used are subject to change → an audit is required to take care of the care and maintenance of the data lake. The main principles / guidelines and measures that regulates data maintenance coordinating all processes in the data lake and responsibilities are defined
Validating	<p>the data are checked one more time before they are integrated into the target CRIS to identify problems with the data quality and consistency of the data, or to confirm that the transformation has been successful.</p> <p>Verify that the values of the attribute are correct and conform to the syntactic and distribution constraints, thus ensuring high data quality AND document every change so that older versions can be restored, or history of changes can be viewed. If new data are generated during data analysis in CRIS, it can be re-included in Data Lake**</p> <p><i>**New data go through the data wrangling process, starting with the step 2 of data validating and structuring the data.</i></p>

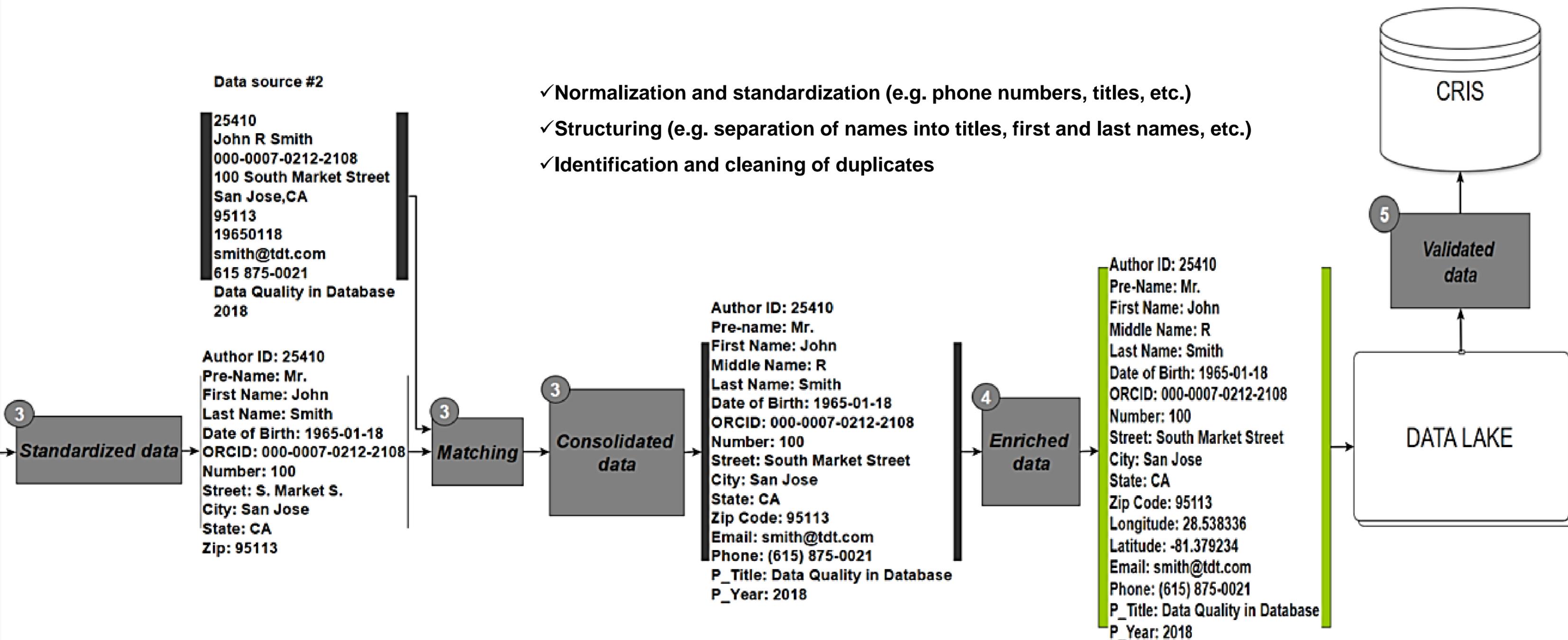
At the end of this process, research information can be used by analytical applications and protected from unauthorized access by access control

USE-CASE



- ✓ Data formatting
- ✓ Correction of incorrect data (e.g. address data)

USE-CASE



USE-CASE: TRIFACTA FOR DATA WRANGLING

The screenshot displays the Trifacta data wrangling interface. At the top, the file path is 'PUBLICATIONS DATA.XLSX / Publikationsdaten aus WoS.xlsx/Publikationsliste - 2'. The main area shows a data table with columns: #, AUTHORID, RBC, FIRSTNAME, RBC, LASTNAME, GENDER, DATE OF BIRTH, RBC, ORCID, and RBC. Above the table, there are histograms for each column. The 'Recipe' panel on the right lists various transformation options:

- Scale to min max: Scale a column to a specific min max range
- One hot encode: Create a column for each unique value indicating its presence or absence
- Scale to mean: Scale a column to zero mean and unit variance
- Bin column: Bin values into ranges of equal or custom size
- New formula: Create a new column from the result of a formula
- Edit with formula: Set one or more columns to the result of a formula
- Window: Perform calculations across multiple ordered rows
- Schema: Change column type (Change the data type of a column)
- Delete columns: Delete one or more columns
- Move columns: Move one or more columns before or after another column
- Rename columns: Rename one or more columns
- Rename with pattern: Rename columns using a pattern
- Rename with prefix: Rename with prefix

At the bottom of the interface, it shows '9 Columns 550 Rows 5 Data Types'.

CONCLUSIONS

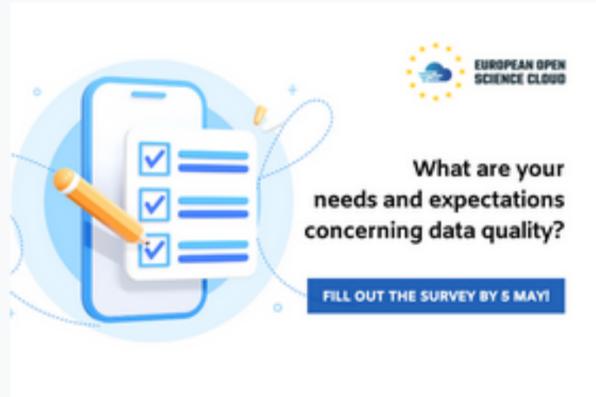
- ✓ **As the volume of research information and data sources increases, the prerequisite for data to be complete, findable, comprehensively accessible, interoperable, reusable (compliant with FAIR principles), but also securely stored, structured, and networked in order to be useful remain critical but at the same time become more difficult to fulfill → data wrangling can be seen a valuable asset in ensuring this.**
- ✓ **The goal is to counteract the growing number of data silos that isolate research information from different areas of the organization. Once successfully implemented, data can be retrieved, managed and made available and accessible to everyone within the entity.**
- ✓ **A data lake and data wrangling can be implemented to improve and simplify IT infrastructure and architecture, governance and compliance. They provide valuable support for predictive analytics and self-service analysis by making it easier and faster to access large amount of data from multiple sources.**
- ✓ **The proper organization of the data lake makes it easier to find the research information the user needs. Managing the data that have already been pre-processed results in an increased efficiency and cost saving, as preparing data for their further use is the most resource-consuming part of data analysis. By providing pre-processed research information, users with limited or no experience in data preparation (low level of data literacy) can be supported and analyzes can be carried out faster and more accurately.**

**TO BE
CONTINUED** →



Complete the Data Quality survey!

12 April 2022



The Data Quality subgroup of the EOSC Association's [Task Force Fair metrics and Data Quality](#) created a survey intended to capture needs and expectations for data quality from a wide range of stakeholders. Your responses will help to identify common approaches and formulate recommendations to support data quality within EOSC (European Open Science Cloud).

[Fill in the survey](#), it will take around **10 minutes** of your time and will be open for you to respond till **Thursday 5th May**.

By completing all answers, we are happy to acknowledge your contribution in the forthcoming EOSC data quality recommendations.

Please email us back (carlo.lacagnina@bsc.es) if you would like your name to be included in the formal acknowledgement. For any further feedback do not hesitate to contact us for an individual interview.

Your opinion is very important for shaping the EOSC!

[Fill in the survey](#)

EOSC - Data Quality Task Force subgroup

This survey is intended to capture needs and expectations for data quality from a wide range of stakeholders. Your responses will help us to identify common approaches and formulate recommendations to support data quality within [EOSC](#) (European Open Science Cloud).

Your opinion is very important for shaping the [EOSC](#)!

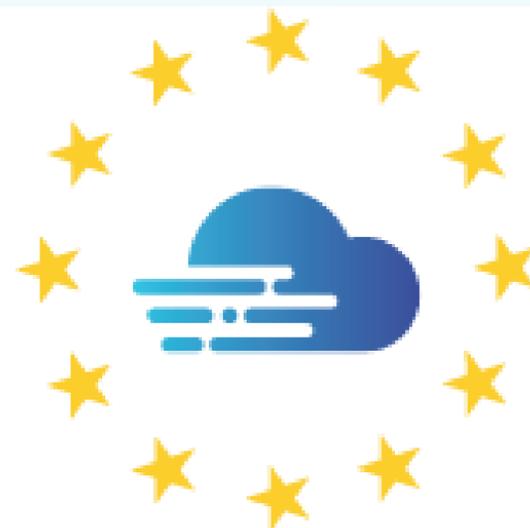
The survey is split into three thematic blocks, the expected time to fill it in is 10 minutes and it will be open for you to respond till Thursday 5th May.



EUROPEAN OPEN SCIENCE CLOUD

[Continue](#) press Enter ↵

[EOSC Data Quality survey](#)



EUROPEAN OPEN SCIENCE CLOUD

**THANK YOU FOR
ATTENTION!
QUESTIONS?**

*For questions or any other queries,
contact us via email -*

Nikiforova.Anastasija@gmail.com,

azeroual@dzhw.eu