

# Data Quality in Environmental and Sustainability Sciences

**Laure Berti-Equille**

IRD

[laure.berti@ird.fr](mailto:laure.berti@ird.fr)

# Outline

---

## I. **Data Quality Problems**

- Illustrative Examples
- Importance of Profiling Data Quality

## II. **Data Cleaning and Preparation**

- Generic data pipeline
- Automating data preparation

# 1. Examples of Data Quality Problems

Relational data quality problems

*Nobel Laureates in Chemistry*

Name	Institution	Institution_City	DoB
Skłodowska-Curie Marie	Institut Pasteur	Varsovie	07-11-1867
M. Curie	Pasteur Institute	Paris	1867-11-07
Melvin Calvin	UC Berkeley	Berkeley	1911-04-08
Marie Curien	Paris	Pasteur Institute	2007-11-07
Avram Hershko	NULL	Haifa	NULL
Ronald Hoffman		US	00000000

**Misfiled Value**

**Representation**

**Duplicates**

**Typos**

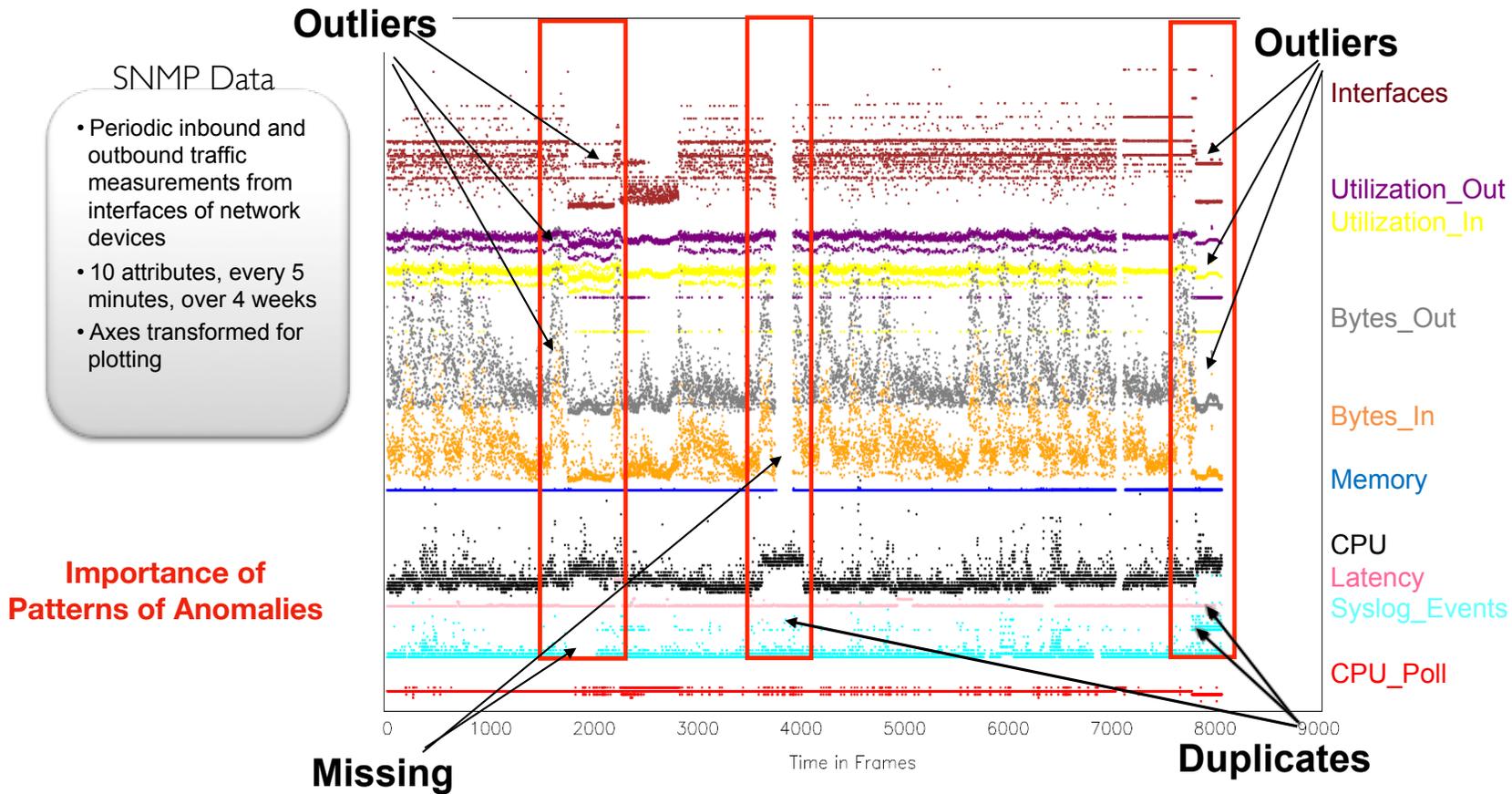
**Inconsistencies**

**Incorrect Value**

**Incorrect Values**

**Missing Values**

# 2. Examples of Sensor Data Quality Problems



# 3. PFAS Data Quality Problems (1/2)

---

## 1. Inconsistent and Incomplete Identification

- **Incomplete DB coverage**  
Current databases often cover only a fraction of existing PFAS compounds.
- **Lack of standardized naming conventions**  
Misidentification, duplication, and difficulty in combining data from different sources.
- **Unknown/Unreported PFAS**  
Many PFAS in use or present in the environment may be proprietary or simply unknown.

## 2. Variable Data Formats and Metadata: Lack of Standardization

- **Inconsistent Formats**  
PFAS Data stored in different formats (spreadsheets, scientific publications, proprietary databases) with various conventions: e.g., PFAS concentrations in parts per million (ppm), while another uses parts per billion (ppb), hindering integration for analysis.
- **Missing Metadata**  
Crucial details about samples (location, collection date, analytical method, etc.) are often incomplete or missing.

# 3. PFAS Data Quality Problems (2/2)

---

## 3. Limited Accessibility, Interoperability, and Reusability

- **Fragmented Storage**

PFAS data is scattered across government agencies, academic institutions, industry, etc. There's no central, comprehensive repository.

- **Manual Curation**

Maintaining and updating data often requires manual curation, which can be time-consuming and prone to human error.

## 4. Measurement and Analytical Gaps

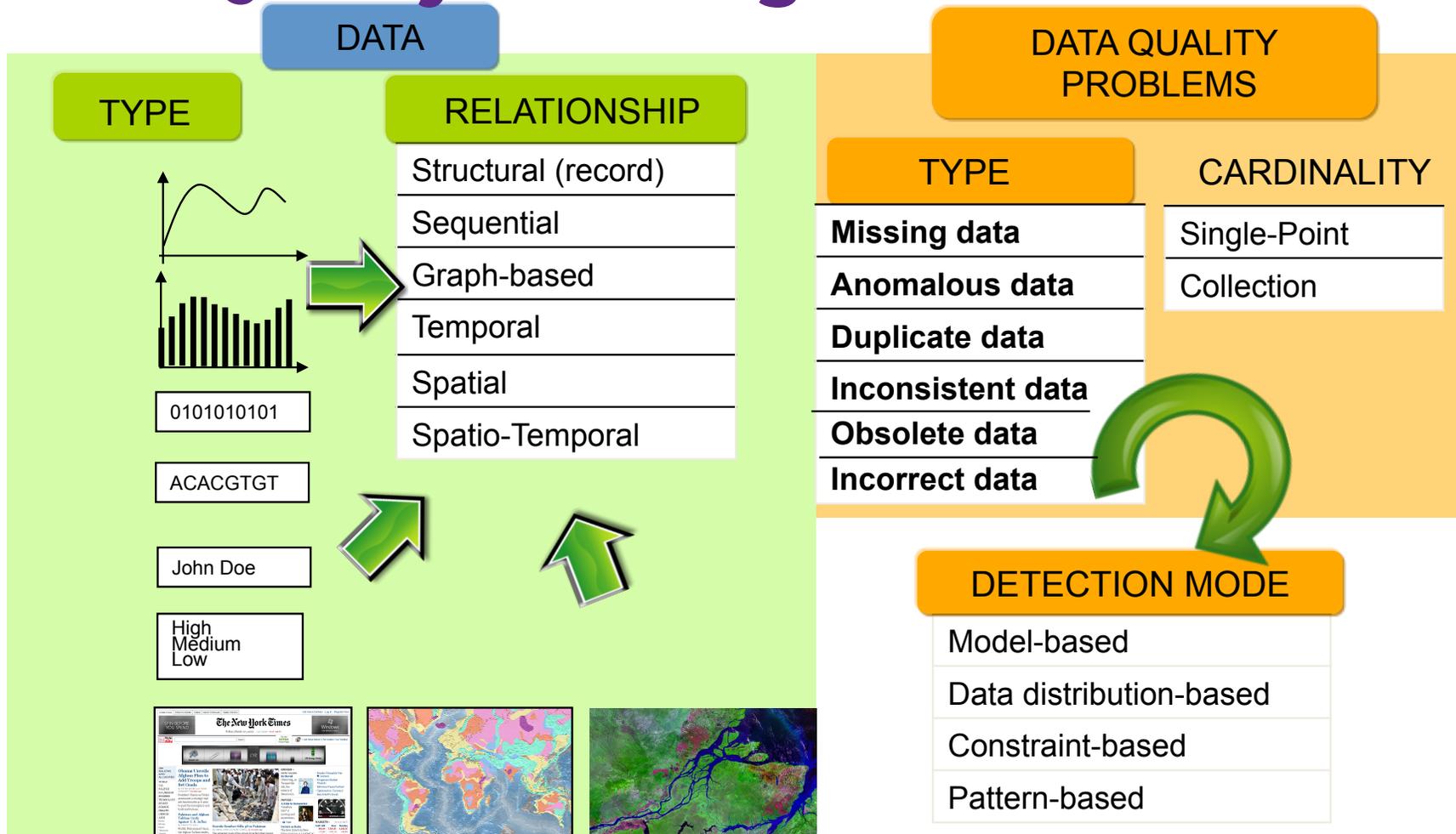
- **Detection Limits**

Analytical techniques for PFAS are constantly evolving. Older data may have high detection limits, meaning low levels of contamination weren't captured at the time.

- **Nontarget Analysis**

Targeted analysis looks for known PFAS. Nontarget analysis, which could identify novel PFAS, is less common and often less accessible.

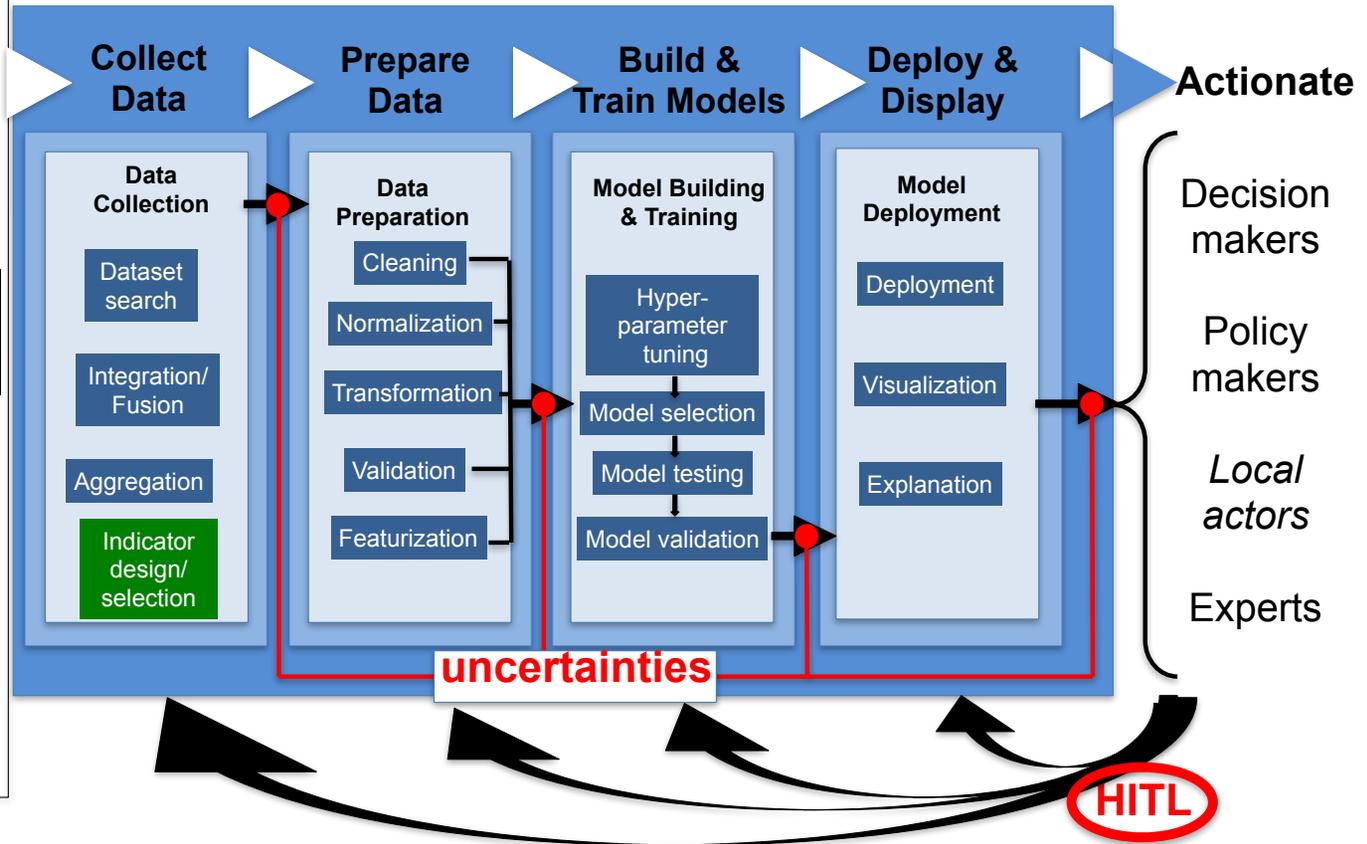
# Data Quality Profiling



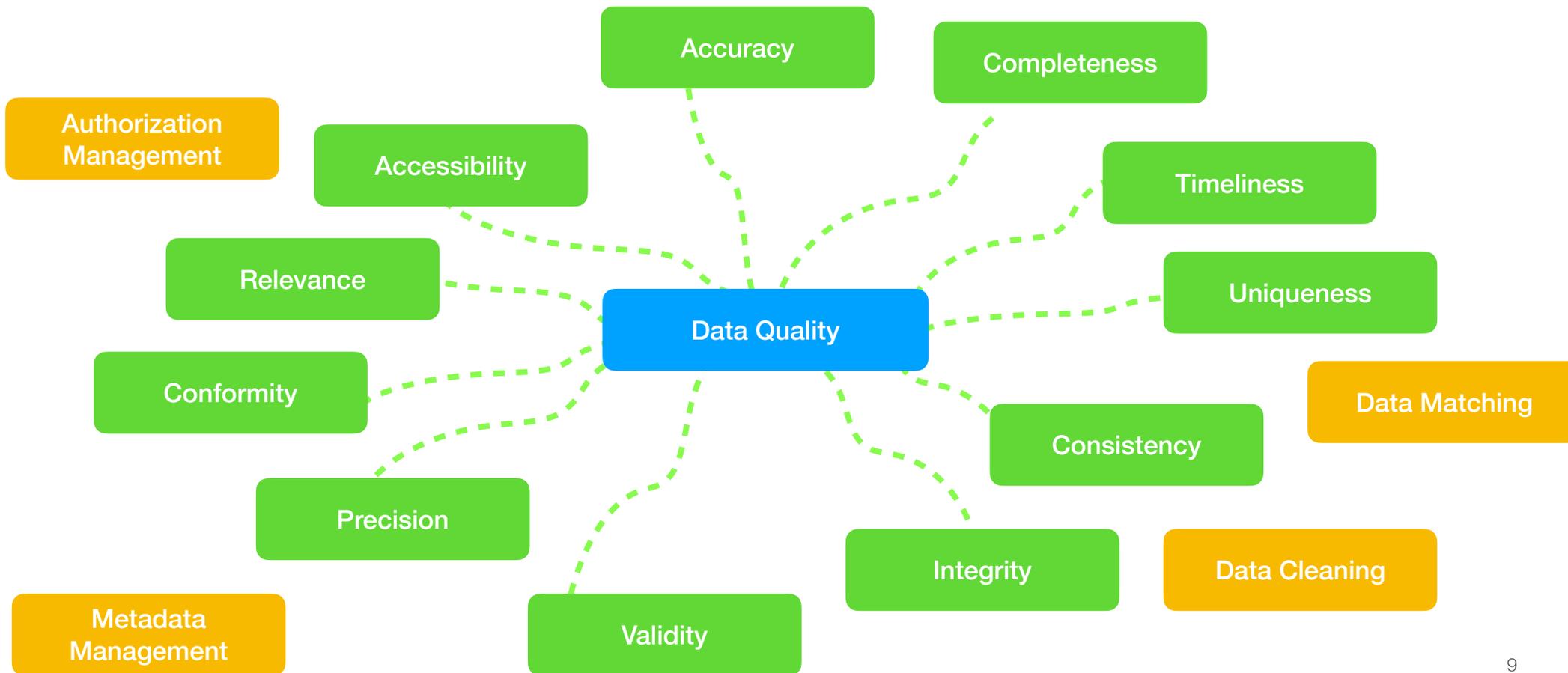
# Data Science Pipeline

*Data quality-awareness*

*Multiple datasets*



# Data Quality Profiling: Reporting Indicators



# Outline

---

## I. **Data Quality Problems**

- Illustrative Examples
- Profiling Data Quality

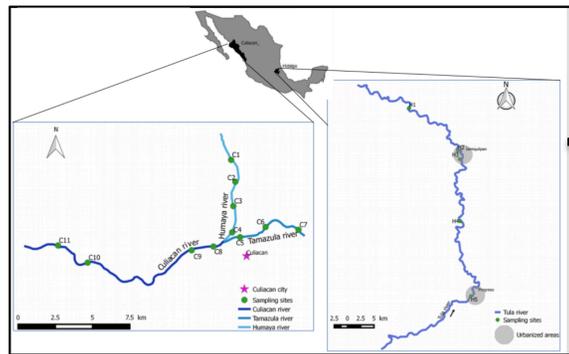
## II. **Data Cleaning and Preparation**

- Impact of Data Prep
- Generic Pipeline

# Data Preprocessing Impact (1/2)

Are the Mexican rivers of Humaya, Tula, Tamazula, and Culiacan polluted ?

[Serrano Balderas et al. 2015, 2017]



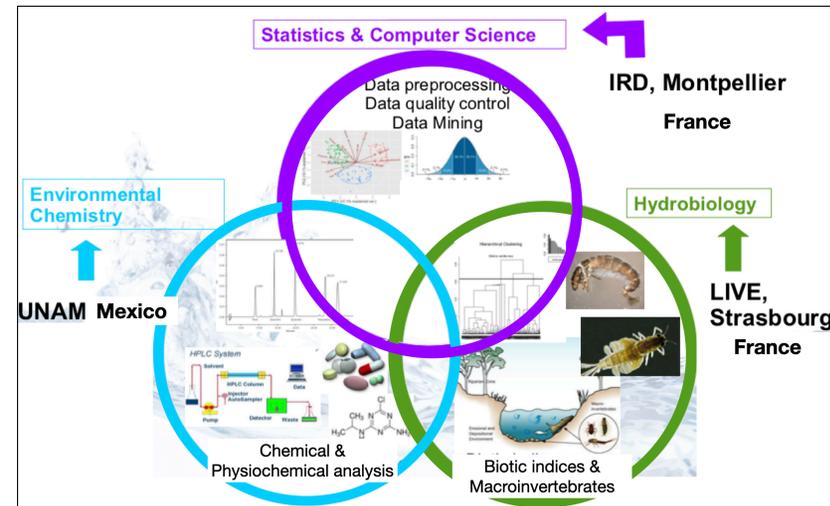
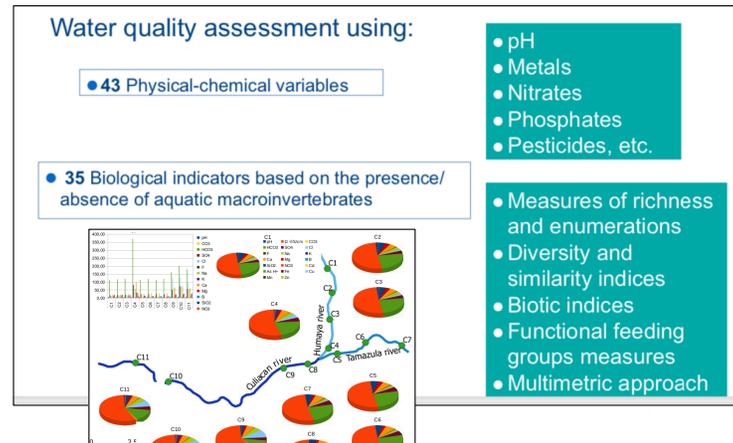
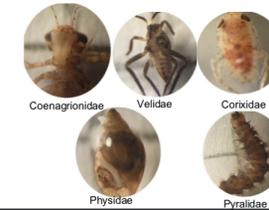
Sampling of liquid and biological samples



Analysis of samples

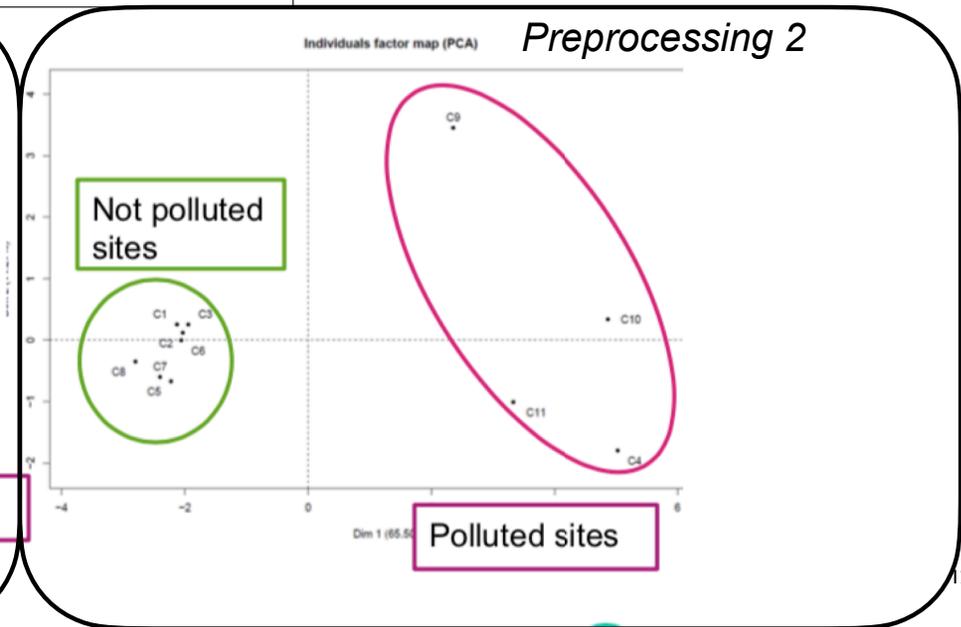
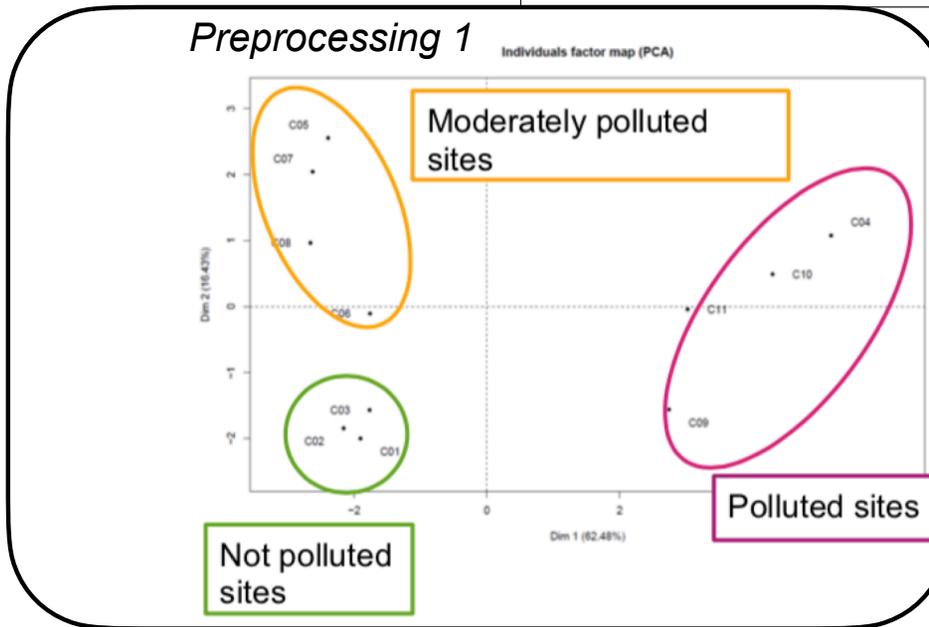
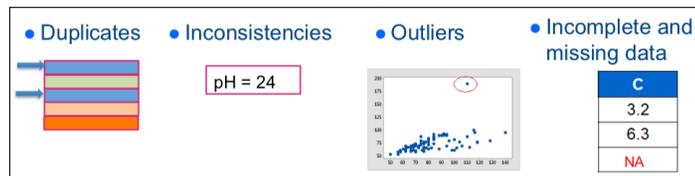


Solid Phase Extraction (SPE)



# Data Preprocessing Impact (2/2)

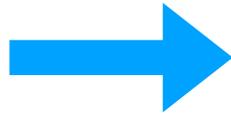
Various data cleaning/preparation strategies lead to different and misleading conclusions



# Data Preprocessing pipeline

---

Real-world  
Dirty Data



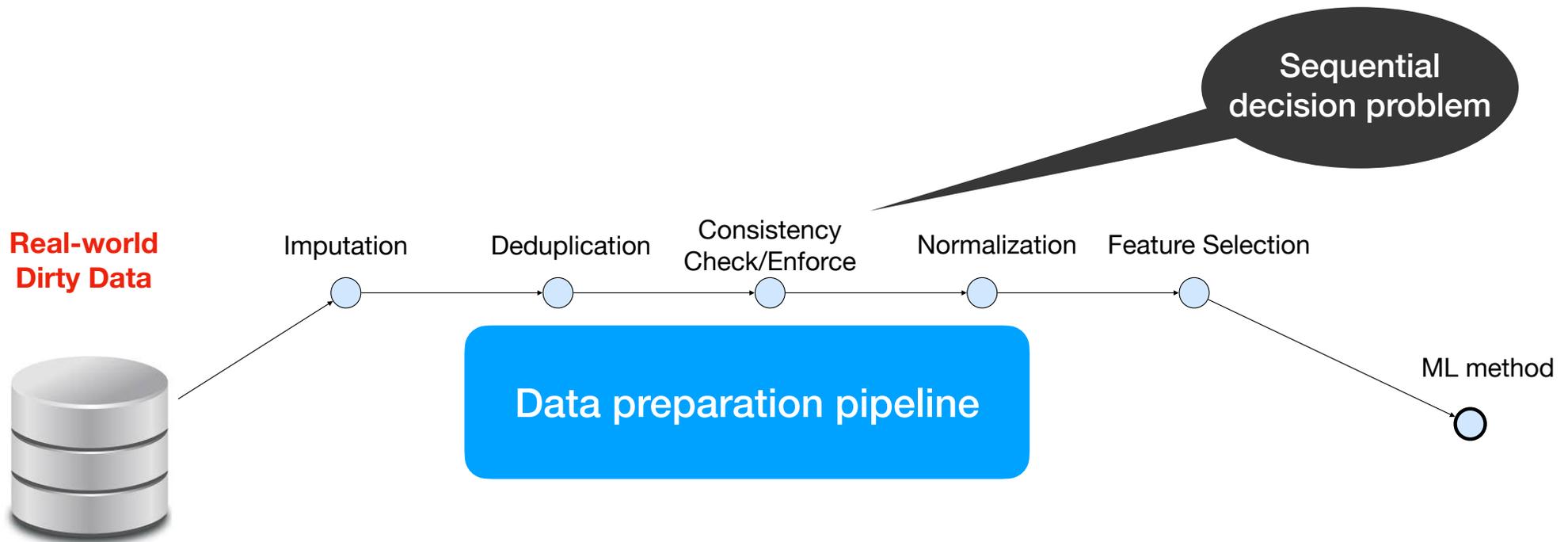
Data preparation pipeline



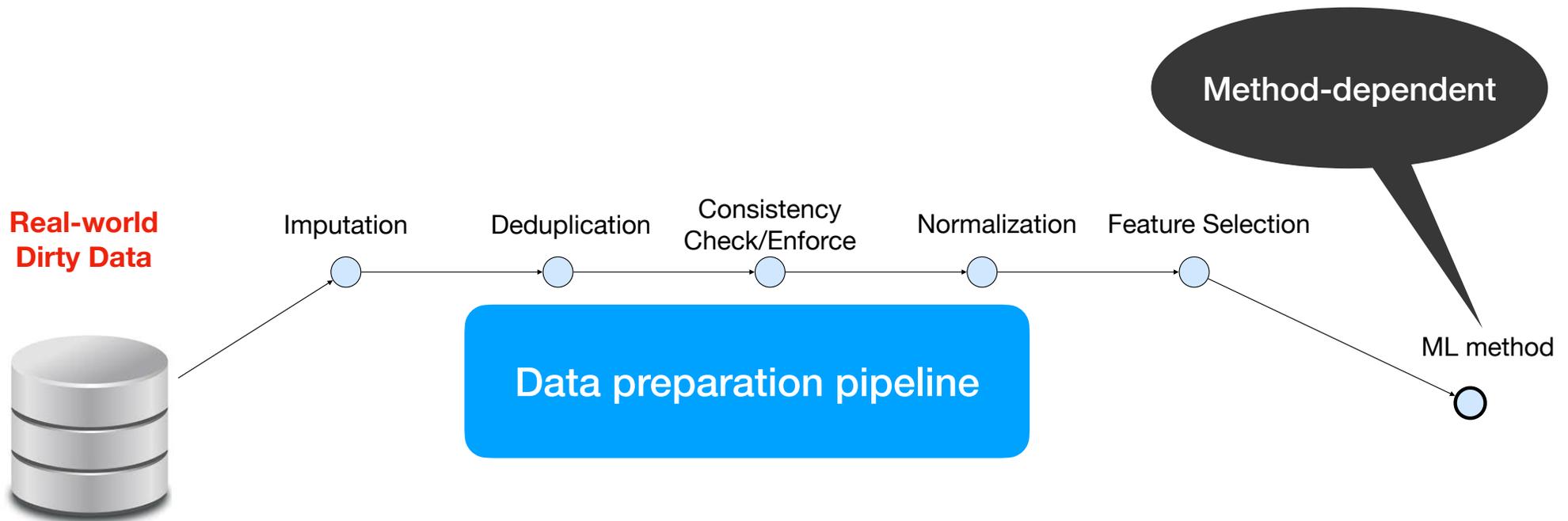
ML method



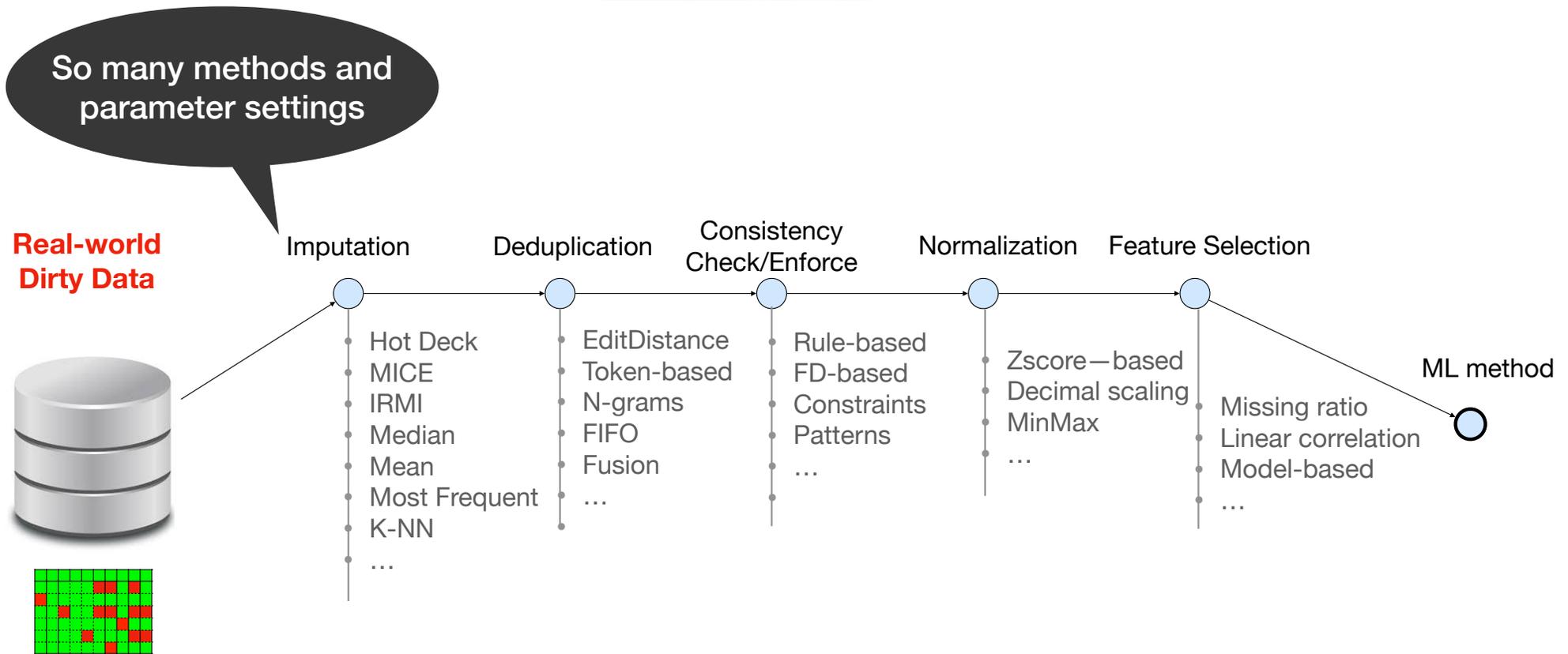
# Data Preprocessing pipeline



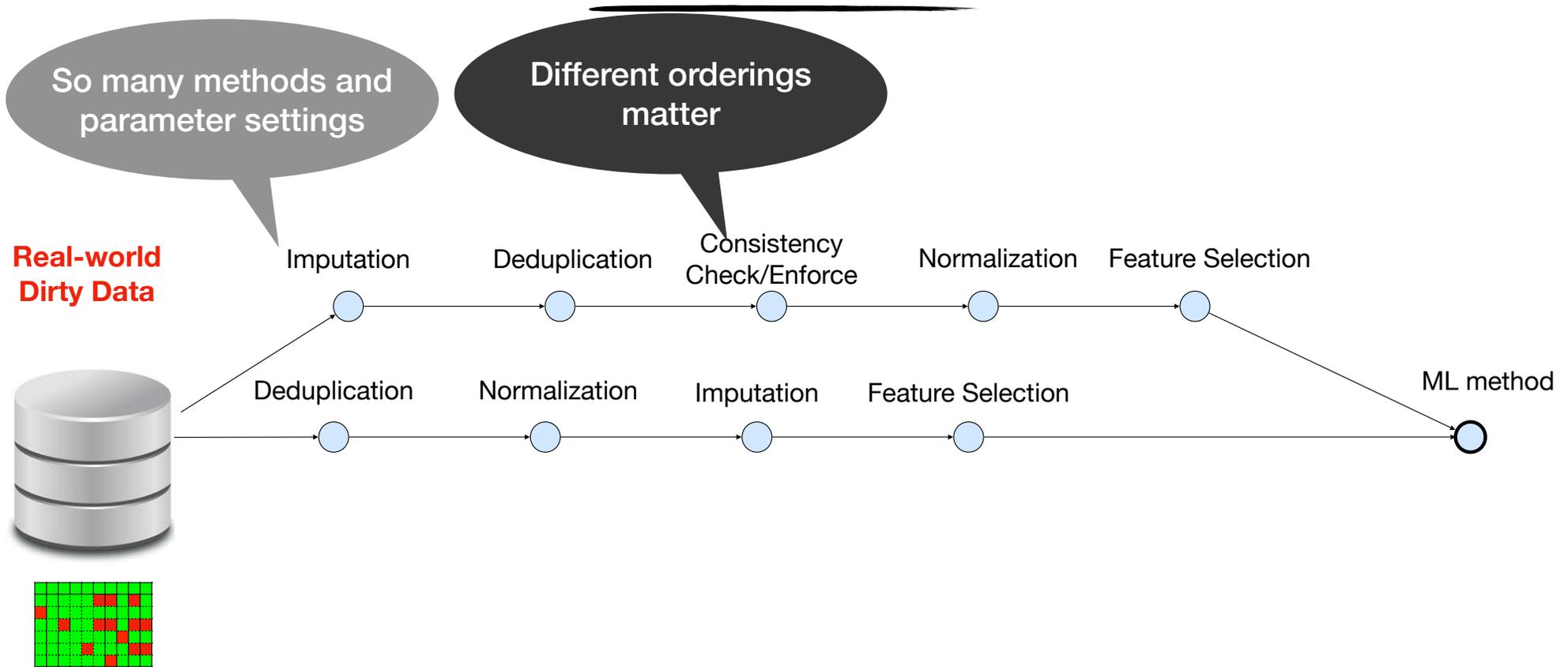
# Data Preprocessing pipeline



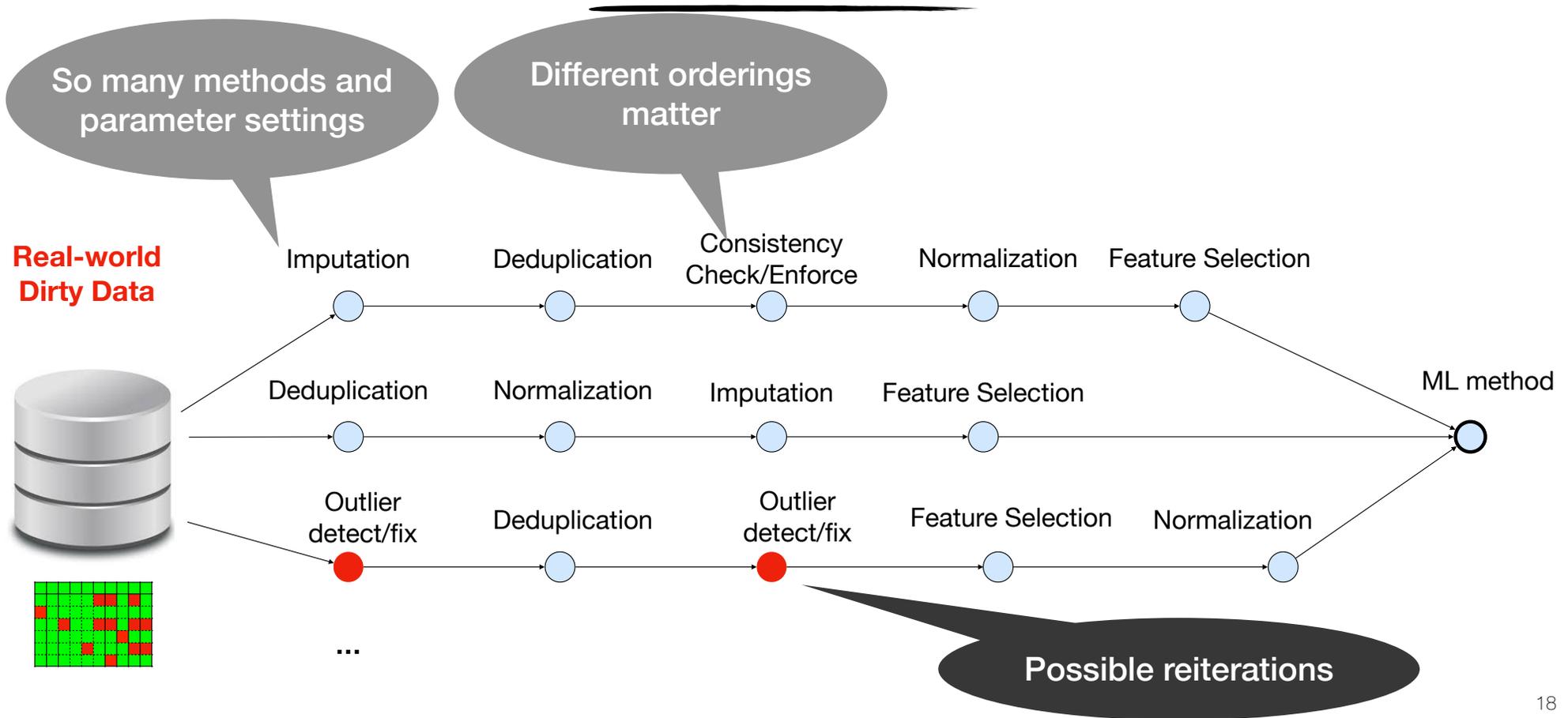
# Data Preprocessing pipeline



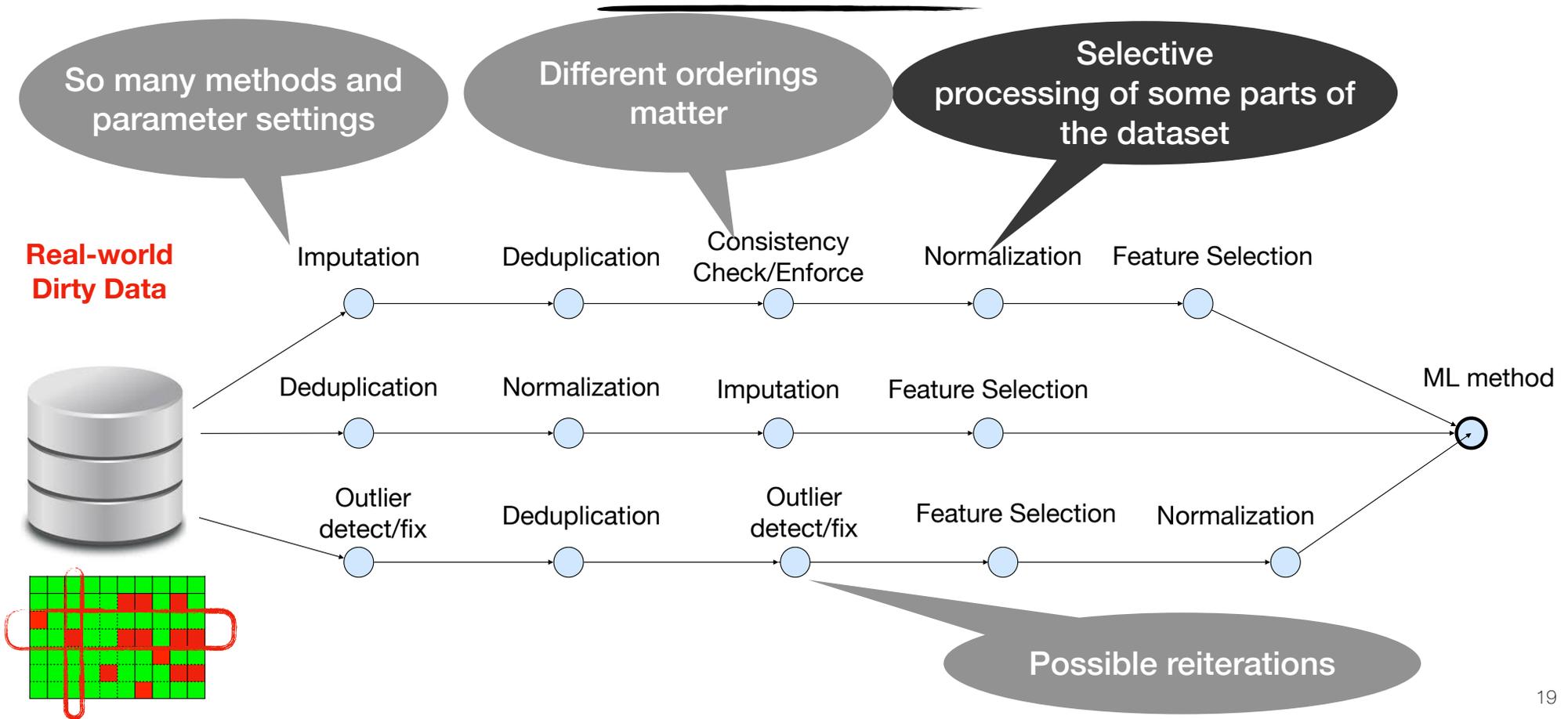
# Data Preprocessing pipeline



# Data Preprocessing pipeline



# Data Preprocessing pipeline

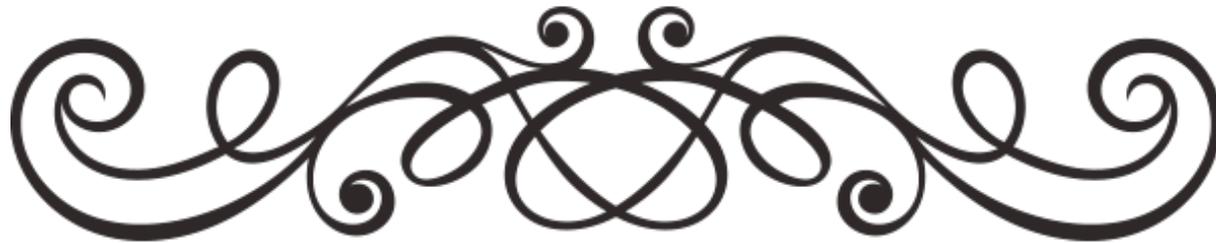




# Conclusions

- PFAS Data standardization and data quality profiling are mandatory starting points.
- Principled data cleaning and data preparation are essential and have a great impact on the analysis results.
- Many data preprocessing tasks require seamless integration of automated tasks and Human-in-the-Loop expertise.
- Many tools and R&D solutions exist for quantifying & monitoring DQ

# Thanks!



**Laure Berti-Equille**

contact: [laure.berti@ird.fr](mailto:laure.berti@ird.fr)

<https://laureberti.github.io/website/>