# Management of Scientific Data

ChemBioSys

Matthias Bruhns, Tobias Giesemann, Lucas Schneider, Tabitha Uphoff
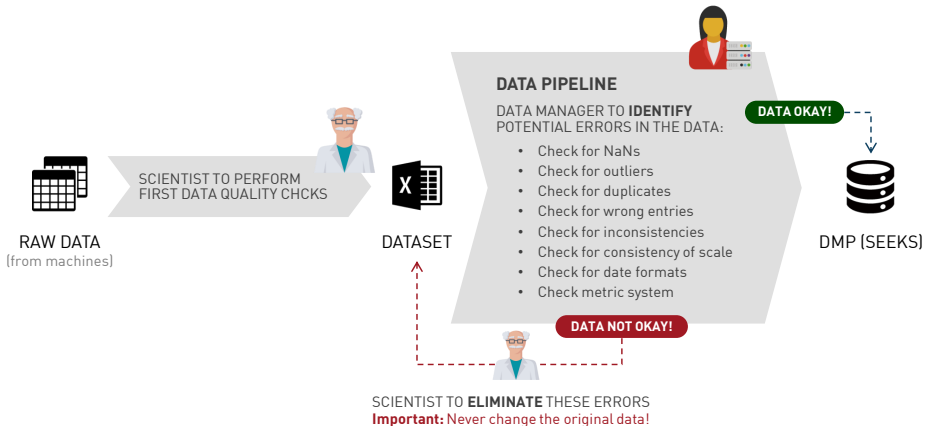13.05.2019

Friedrich-Schiller-Universität Jena

# General principles

**RAW DATA**
(from machines)

SCIENTIST TO PERFORM
FIRST DATA QUALITY CHCKS

**DATASET**

**DATA PIPELINE**

DATA MANAGER TO **IDENTIFY**
POTENTIAL ERRORS IN THE DATA:

- Check for NaNs
- Check for outliers
- Check for duplicates
- Check for wrong entries
- Check for inconsistencies
- Check for consistency of scale
- Check for date formats
- Check metric system

**DATA OKAY!**

**DATA NOT OKAY!**

**DMP (SEEKS)**

SCIENTIST TO **ELIMINATE** THESE ERRORS
**Important:** Never change the original data!

# Inconsistency, wrong entries and duplicates

# Data Quality

| | | |
|---|---|---|
| 🦴 | **Inconsistencies** | Formats: Time Stamps, Date, Numericals, … <br> Scales: (kg, g) (m, mm) <br> System: (metric, imperial) |

| | | |
|---|---|---|
| ✖ | **Wrong Entries:** | Impossible Values -> set to NULL for analysis <br> Outliers -> exclude from analysis <br> NEVER overwrite original set |

| | | |
|---|---|---|
| 🔑 | **Duplicates:** | In DBS: Primary Keys <br> Consistent IDs should be used |

# Data Quality @ ChemBioSys

| | | |
|---|---|---|
| 🦴 | Inconsistencies | Formats: Machine-created Data<br>Scales: Machine-created Data e.g. NMR-Machine, Mass-Spectrometry, …<br>System: metric |
| ✖ | Wrong Entries: | Impossible values: Machine Errors<br>Outlier Detection:<br>Checkup Routine done by individual Scientists |
| 🔑 | Duplicates: | Data Management Platform (SEEK)<br>ID-consistency in each spreadsheet<br>Overlooked by each individual scientist |

# NaNs and Outliers

### NaNs

- NaNs can appear when the machine could not measure a value
- Add comments in Meta Data file
    - text based: not computer-readable
    - might lead to misunderstandings
- Replace NaNs with zero
    - might lead to wrong analysis
- Delete NaNs from data file
    - if more than 20% of the data are NaNs, data should not be used or pulished
- Indicate "% of missed values" in meta data file

### Outliers

- Outliers can appear when irregularities during the experiment happen
- Use statistical tests to identify "real" outliers $\rightarrow$ expensive
- Exclude entries above or below a certain threshold

# Usage of machines/computers

### Pro

- Most machines provide meta data file for each experiment
- Machines can perform first data quality check internally
- Errors can only be caused by external factors (assuming good software)

### Con

- Proprietrary software might become a problem
    - File formats: Updates can render old formats useless
    - Bugs: No way to check the quality of the internal software
  - $\rightarrow$ Some formats might be convertable using scripts

Questions?