# Spark Assignment

## Distributed Data Management

**Tobias Jordan & Nastassia Heumann - 7th July 2021**

# Our Approach
## Sindy Algorithm

- Our implementation is based on the Sindy slides from Data Profiling (WS20/21)

- We strongly make use of the functional methods `map(), reduce(), filter(), sort(), foreach()`

- What we found challenging:

  - Figuring out how to handle various Scala data types (DataFrame, WrappedArray, Dataset, …)

  - Learning Scala Syntax

- What surprised us is how little code is needed to detect unary INDs

# Our Approach
## Sindy Algorithm

1. Read the datasets

2. Map all cell values to the column

3. Join all columns with the same cell value

4. Based on the joined columns, aka attribute groups, compute all IND candidates

5. Consolidate IND candidates

6. Print INDs