



## Data Profiling with Metanome

Exercise: INDs

Thorsten Papenbrock

PhD Candidate

Hasso-Plattner-Institute

# Advanced Profiling

## Three important metadata



### Unique Column Combinations

#### Key candidates

| Name    | Type        | Equatorial diameter | Mass  |
|---------|-------------|---------------------|-------|
| Mercury | Terrestrial | 0.382               | 0.06  |
| Venus   | Terrestrial | 0.949               | 0.82  |
| Earth   | Terrestrial | 1.000               | 1.00  |
| Mars    | Terrestrial | 0.532               | 0.11  |
| Jupiter | Giant       | 11.209              | 317.8 |
| Saturn  | Giant       | 9.449               | 95.2  |
| Uranus  | Giant       | 4.007               | 14.6  |
| ...     | ...         | ...                 | ...   |

**|Name|**

### Inclusion Dependencies

#### Foreign key candidates

| Sign        | Domicile |
|-------------|----------|
| Aries       | Mars     |
| Taurus      | Venus    |
| Gemini      | Mercury  |
| Cancer      | Moon     |
| Leo         | Sun      |
| Virgo       | Mercury  |
| Libra       | Venus    |
| Scorpio     | Pluto    |
| Sagittarius | Jupiter  |
| Capricorn   | Saturn   |
| Aquarius    | Uranus   |
| ...         | ...      |

| Name    | Type        |
|---------|-------------|
| Mercury | Terrestrial |
| Venus   | Terrestrial |
| Earth   | Terrestrial |
| Mars    | Terrestrial |
| Jupiter | Giant       |
| Saturn  | Giant       |
| Uranus  | Giant       |
| ...     | ...         |

**Domicile  $\subseteq$  Name**

### Functional Dependencies

#### Normalization criterion

| Name    | Atmosphere                            | Rings |
|---------|---------------------------------------|-------|
| Mercury | minimal                               | no    |
| Venus   | CO <sub>2</sub> , N <sub>2</sub>      | no    |
| Earth   | N <sub>2</sub> , O <sub>2</sub> , Ar  | no    |
| Mars    | CO <sub>2</sub> , N <sub>2</sub> , Ar | no    |
| Jupiter | H <sub>2</sub> , He                   | yes   |
| Saturn  | H <sub>2</sub> , He                   | yes   |
| Uranus  | H <sub>2</sub> , He                   | yes   |
| ...     | ...                                   | ...   |

**Atmosphere  $\rightarrow$  Rings**

### Data Profiling with Metanome

Thorsten Papebrock,  
PhD Candidate,  
17<sup>th</sup> November, 2014  
Chart 2

# Exercise 1

## Discovery of key candidates

---



- All teams have passed the exercise:
  - 36 submissions
  - No duplicate algorithm names!
  - Publication usually not wanted
  
- Tips:
  - Use the mailing-list for help!
  - Check your algorithm before submission
  - Use the Metanome helper projects carefully

### **Data Profiling with Metanome**

Thorsten Papenbrock,  
PhD Candidate,  
17<sup>th</sup> November, 2014  
Chart **3**

# Exercise 1

## NULL values



### **NULL = NULL**

can cause ...

#### **larger PLIs**

(one additional large index-set for NULL values)

#### **larger UCCs**

(attributes with NULLs need other attributes to become unique)

### **NULL $\neq$ NULL**

can cause ...

#### **smaller PLIs**

(NULL values can directly be ignored during PLI creation)

#### **smaller UCCs**

(attributes with NULL values can be unique themselves)

**Data Profiling with Metanome**

Thorsten Papenbrock,  
PhD Candidate,  
17<sup>th</sup> November, 2014  
Chart 4

# Exercise 1

## NULL values



### NULL = NULL

can cause ...

### NULL ≠ NULL

can cause ...

#### more minimal UCCs

(attributes are unique in different combinations)

#### more minimal UCCs

(attributes with NULL values are unique itself)

| A | B | C | D | E | NULL = NULL                                    |
|---|---|---|---|---|--|
| - | - | - | 1 | 1 | {A,D}, {A,E},<br>{B,D}, {B,E},<br>{C,D}, {C,E} |
| - | - | - | 2 | 2 |  |
| 1 | 1 | 1 | - | - | NULL ≠ NULL<br>{A}, {B}, {C},<br>{D}, {E}      |
| 2 | 2 | 2 | - | - |  |

| A | B | C | NULL = NULL                  |
|---|---|---|------------------------------|
| - | 1 | 1 | {B}, {C}                     |
| - | 2 | 2 | NULL ≠ NULL<br>{A}, {B}, {C} |

### Data Profiling with Metanome

Thorsten Papebrock,  
PhD Candidate,  
17<sup>th</sup> November, 2014  
Chart 5

# Exercise 1

## Our evaluation

---



- DELL Optiplex 9010
  - CPU: Intel i5 3.2 GHz
  - RAM: 8 GB (2 GB for Metanome JVM)
  - OS: Debian 64-bit
  - JVM: Java 1.7

### **Data Profiling with Metanome**

Thorsten Papenbrock,  
PhD Candidate,  
17<sup>th</sup> November, 2014  
Chart **6**

AlexoFredUniques  
AWsome  
Buelow\_Flemming\_UCC  
CPUcc  
dennis\_marius.ucc  
DJ\_UCC  
dpdc-cnms-ucc  
DreamteamUcc  
DreamteamUcc-par  
EnteEnte  
FrohnOttoUcc  
HorLehUcc  
klinger\_marten\_ucc  
LucieKerstinUCC  
Metanomnomnom  
MgEsUcc  
MMUniqueColumnCombinations  
olorin  
PracticeUccJoerke  
RobertJPUniqueCol  
RohloffUcc  
SBMMUCC  
Tsun12Ucc  
ucc  
ucc\_bothe\_reissaus  
UCC\_Grundke\_Wiese  
Ucc\_Jung  
Ucc\_Kirsten\_Zwerg  
UCC\_schaeffer\_zoellner  
UCC\_SPIRO  
UCCDraegerSchueler  
UCCFMJR  
UccHeikoMax  
UccPerchykSchmidt  
UniColCom  
UniqueColumnCombinationDiscoverer  
YUCC



## Data Profiling with Metanome

ThorstenPapenbrock,  
PhD Candidate,  
17<sup>th</sup> November, 2014  
Chart **7**

AlexoFredUniques  
AWsome  
Buelow\_Flemming\_UCC  
CPUcc  
dennis\_marius.ucc  
DJ\_UCC

**dpdc-cnms-ucc**

DreamteamUcc  
DreamteamUcc-par  
EnteEnte  
FrohnOttoUcc

**HorLehUcc**

klinger\_marten\_ucc

**LucieKerstinUCC**

Metanomnomnom

**MgEsUcc**

MMUniqueColumnCombinations

**olorin**

PracticeUccJoerke  
RobertJPUniqueCol  
RohloffUcc

**SBMUCC**

Tsun12Ucc  
ucc

ucc\_bothe\_reissaus  
UCC\_Grundke\_Wiese  
Ucc\_Jung

Ucc\_Kirsten\_Zwerg  
UCC\_schaeffer\_zoellner  
UCC\_SPIRO

UCCDraegerSchueler

**UCCFMJR**

UccHeikoMax

**UccPerchykSchmidt**

UniColCom  
UniqueColumnCombinationDiscoverer  
YUCC

**Incorrect Result**

**NullPointerException**

**NullPointerException**

**NullPointerException**

**NoSuchMethodException**

**NoSuchMethodError**

**NoSuchMethodError**

**NoSuchMethodError**



**Reasons:**

1. Java 1.8
2. Old Metanome version
3. Programming error

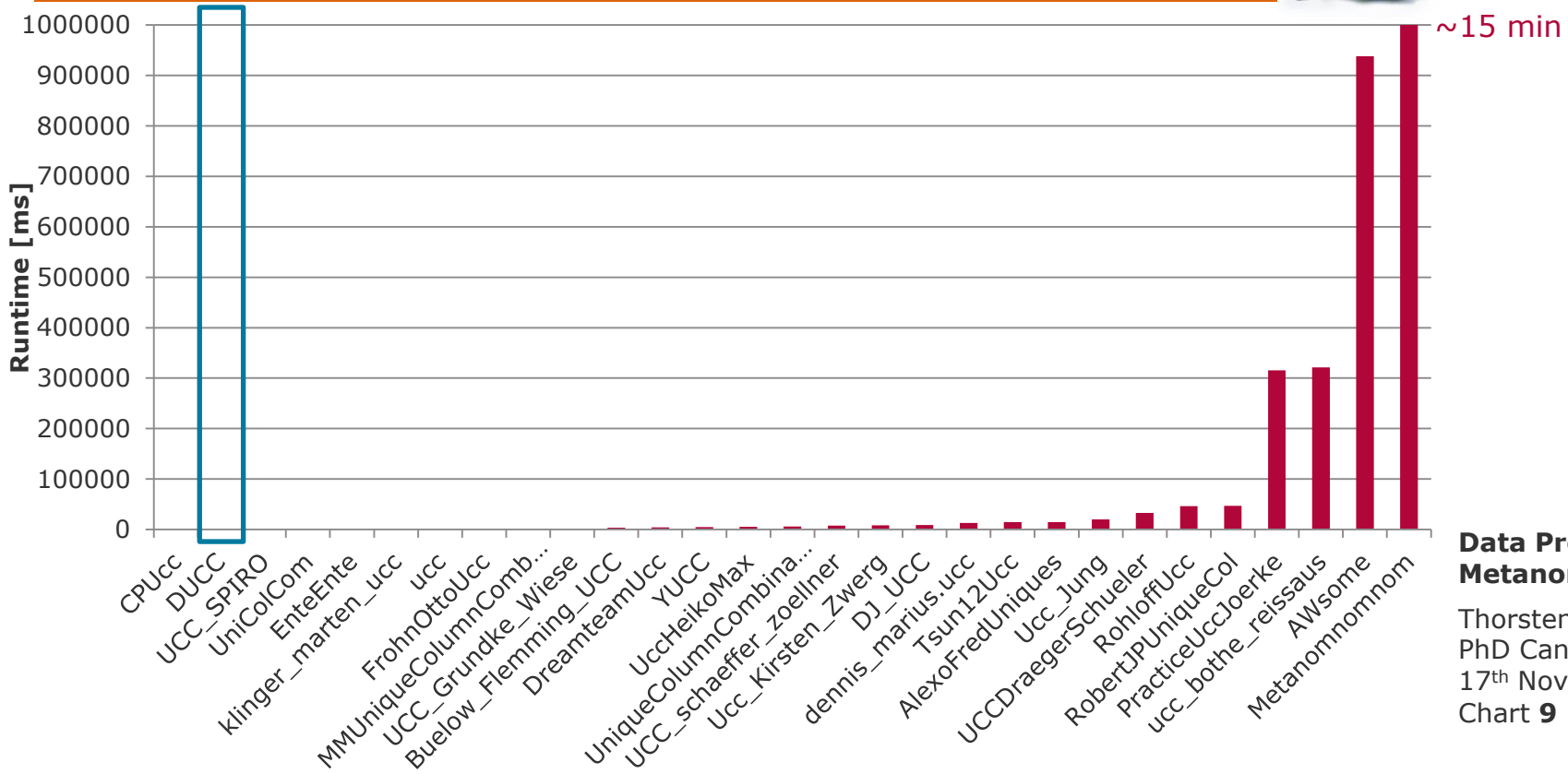
**Data Profiling with  
Metanome**

ThorstenPapenbrock,  
PhD Candidate,  
17<sup>th</sup> November, 2014  
Chart **8**



# Exercise 1

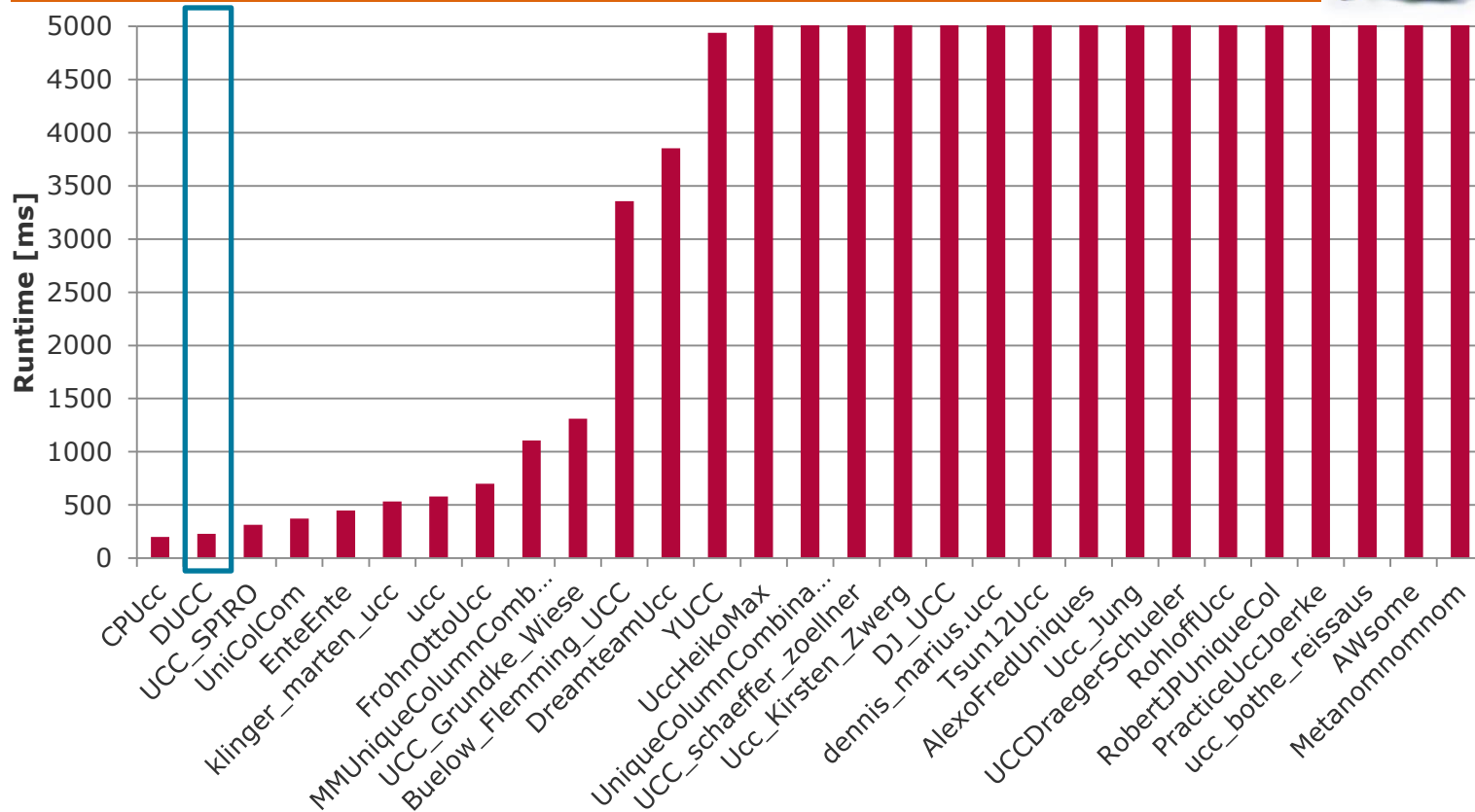
## Our evaluation – ncvoter-1k



**Data Profiling with Metanome**  
Thorsten Papenbrock,  
PhD Candidate,  
17<sup>th</sup> November, 2014  
Chart 9

# Exercise 1

## Our evaluation – ncvoter-1k (<5s)

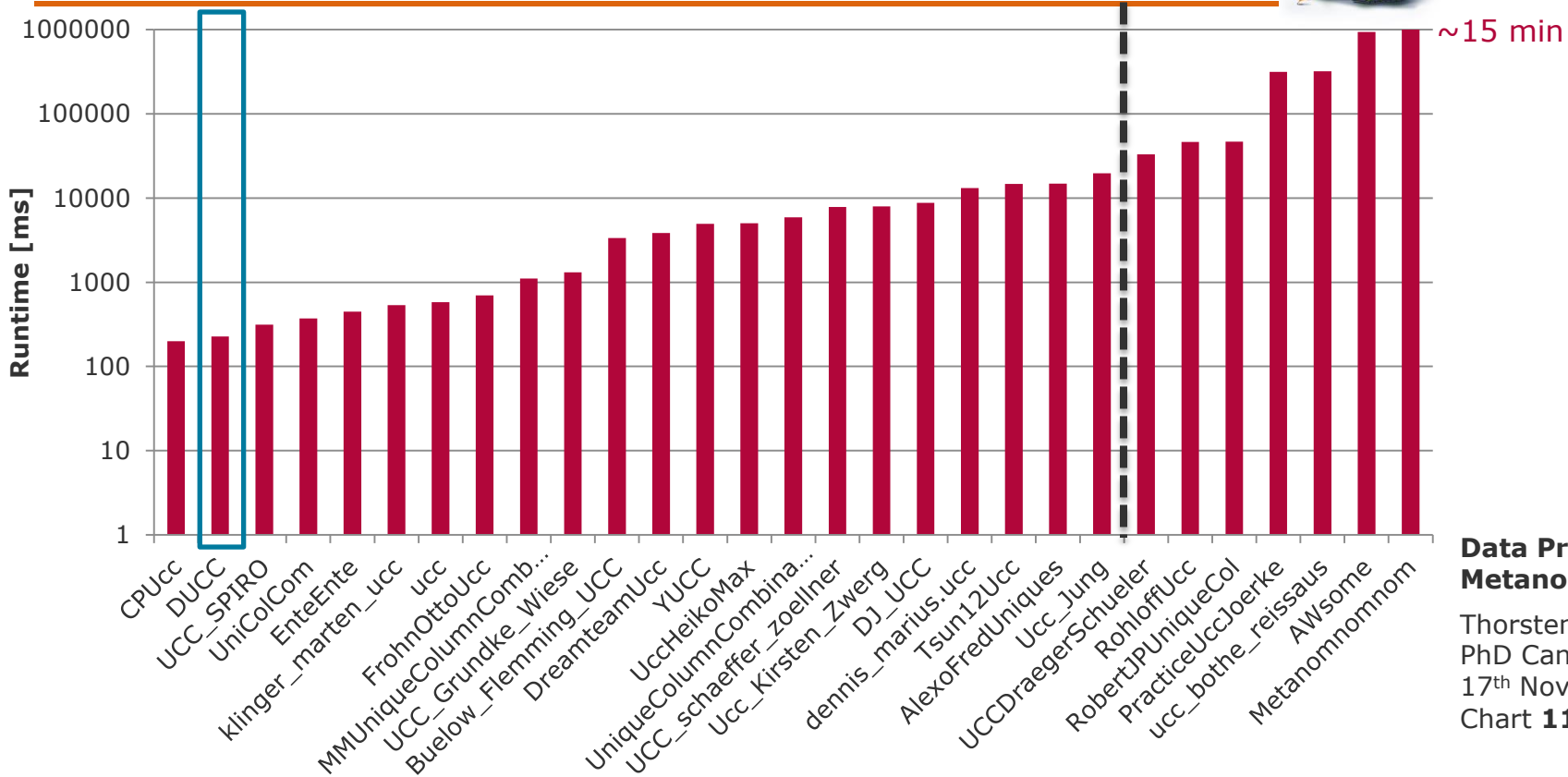


### Data Profiling with Metanome

Thorsten Papanbrock,  
PhD Candidate,  
17<sup>th</sup> November, 2014  
Chart 10

# Exercise 1

## Our evaluation – ncvoter-1k (log)

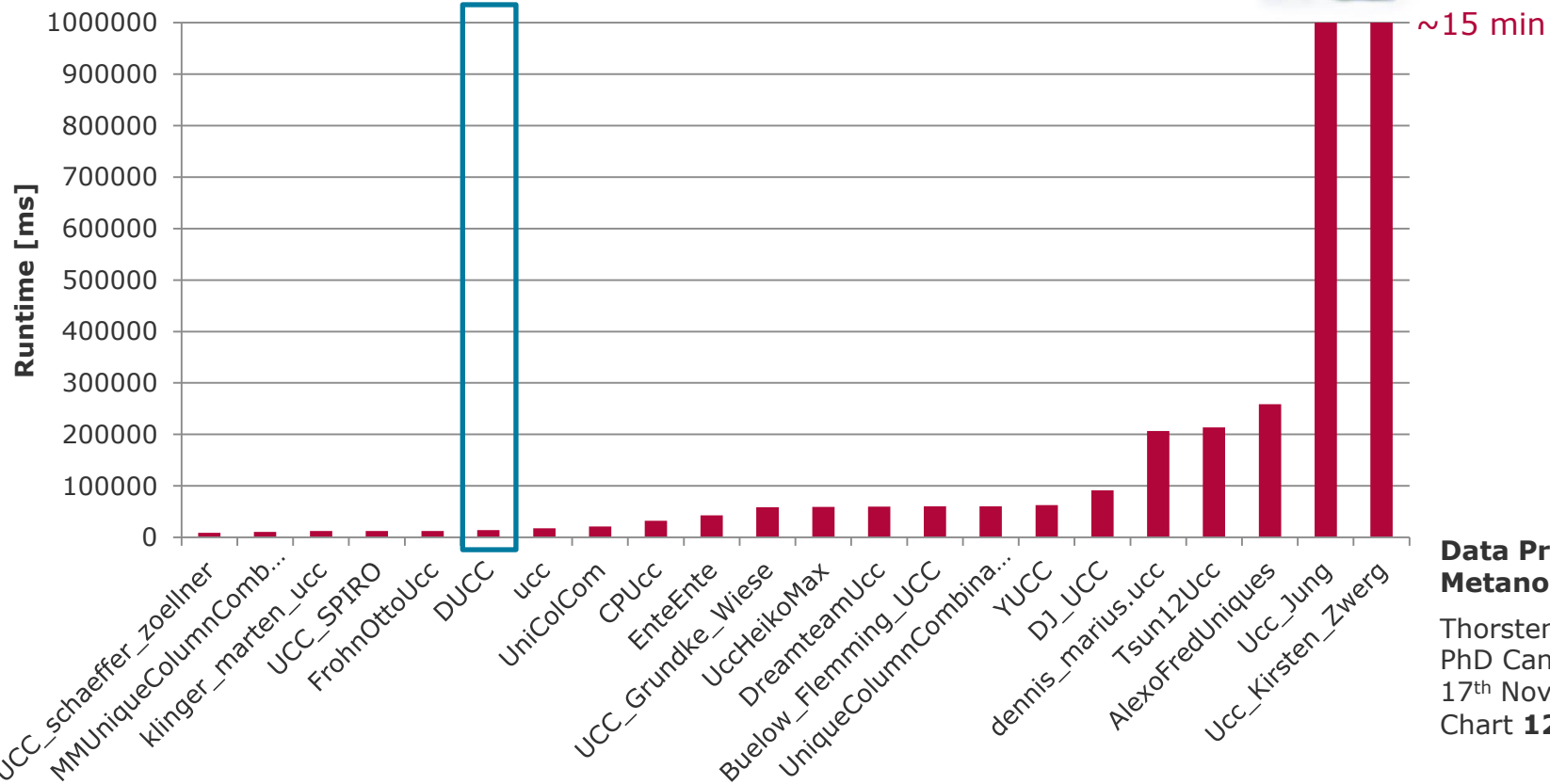


**Data Profiling with Metanome**

Thorsten Papanbrock,  
PhD Candidate,  
17<sup>th</sup> November, 2014  
Chart 11

# Exercise 1

## Our evaluation – fd\_reduced\_15

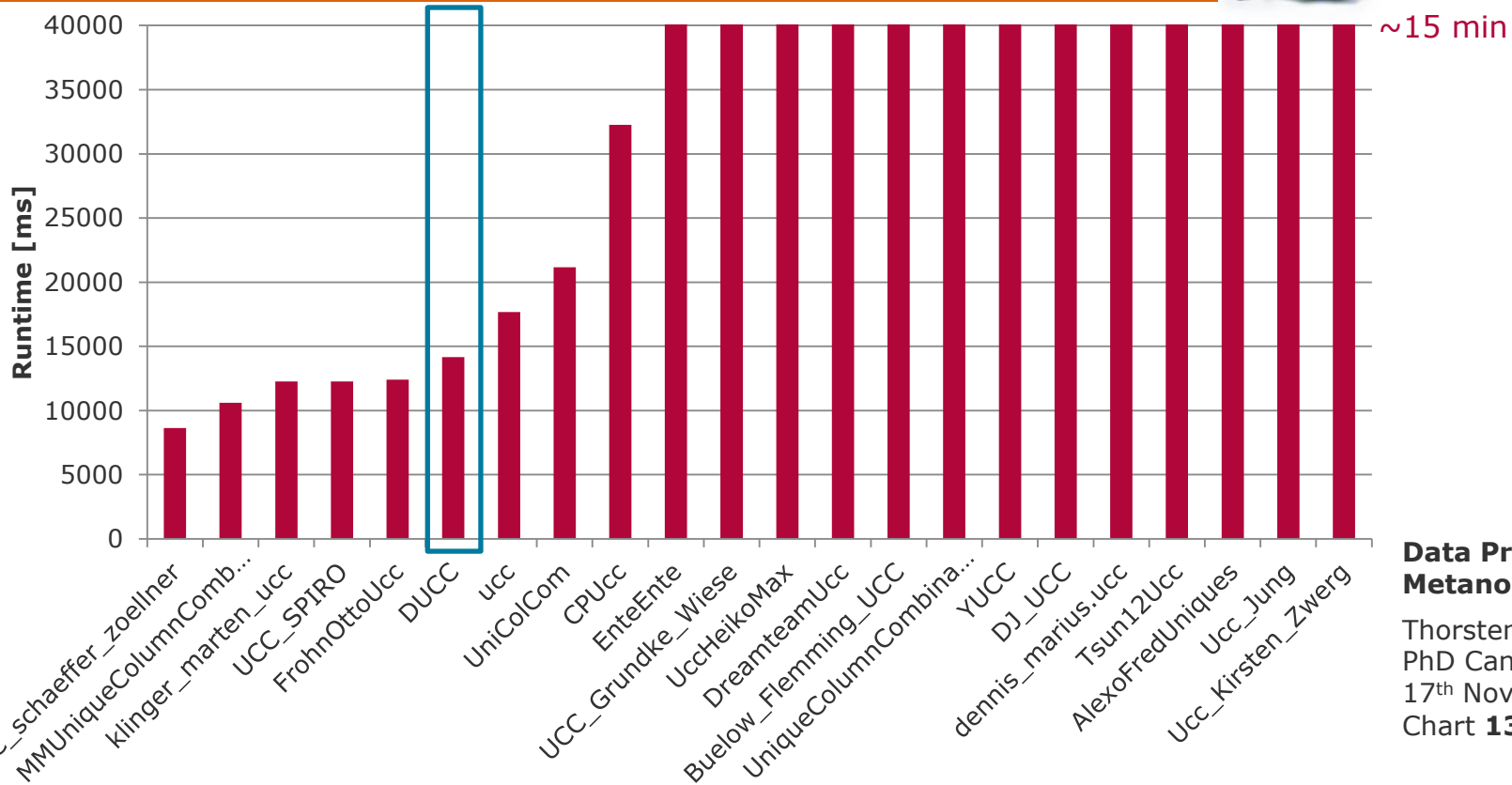


### Data Profiling with Metanome

Thorsten Papenbrock,  
PhD Candidate,  
17<sup>th</sup> November, 2014  
Chart 12

# Exercise 1

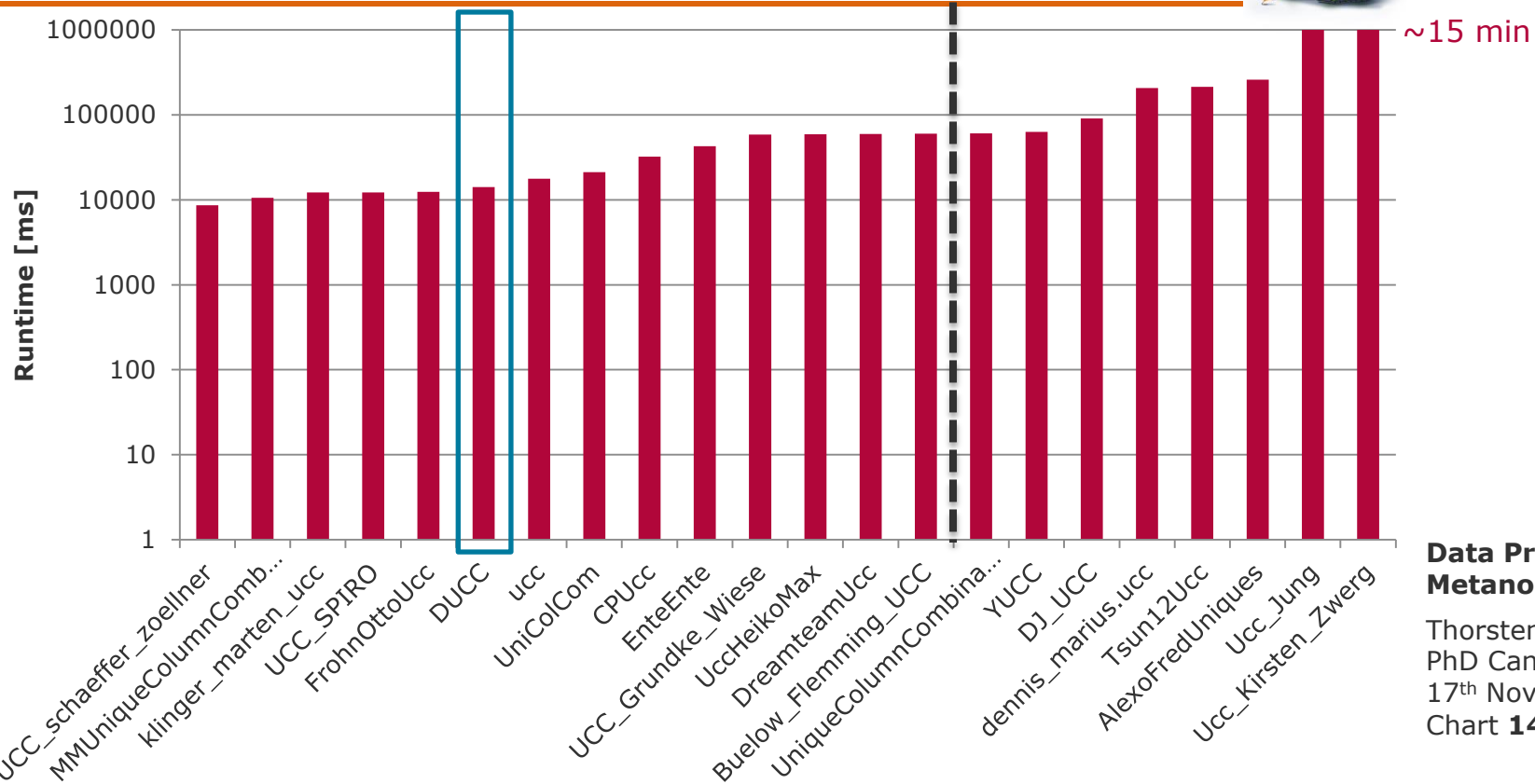
## Our evaluation – fd\_reduced\_15



**Data Profiling with Metanome**  
Thorsten Papenbrock,  
PhD Candidate,  
17<sup>th</sup> November, 2014  
Chart 13

# Exercise 1

## Our evaluation – fd\_reduced\_15

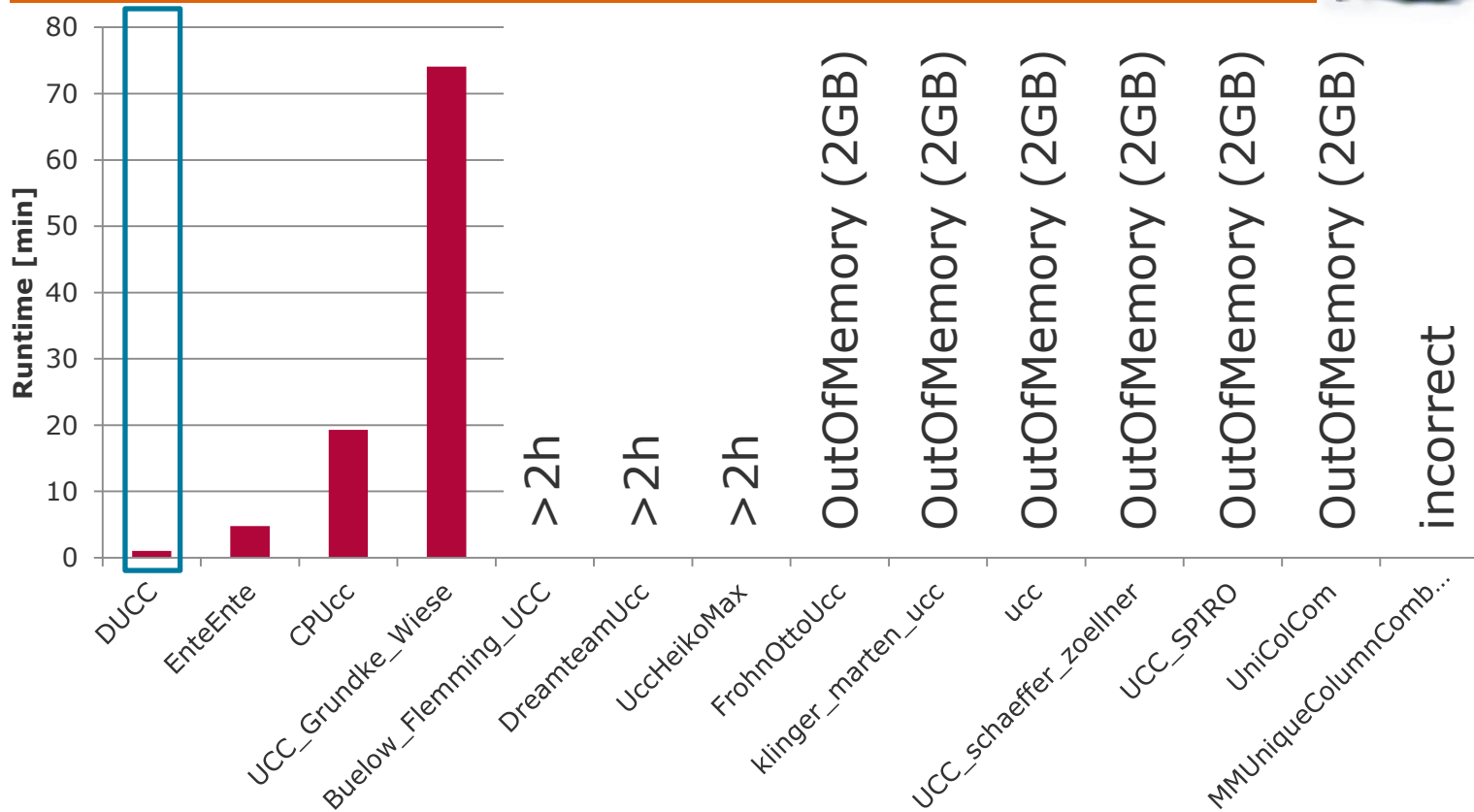


### Data Profiling with Metanome

Thorsten Papenbrock,  
PhD Candidate,  
17<sup>th</sup> November, 2014  
Chart 14

# Exercise 1

## Our evaluation – ncvoter-1m



### Data Profiling with Metanome

Thorsten Papanbrock,  
PhD Candidate,  
17<sup>th</sup> November, 2014  
Chart 15

# Exercise 1

## Short presentations

---



### **Data Profiling with Metanome**

Thorsten Papenbrock,  
PhD Candidate,  
17<sup>th</sup> November, 2014  
Chart **16**



# Advanced Profiling

## Three important metadata



### Unique Column Combinations

#### Key candidates

| Name    | Type        | Equatorial diameter | Mass  |
|---------|-------------|---------------------|-------|
| Mercury | Terrestrial | 0.382               | 0.06  |
| Venus   | Terrestrial | 0.949               | 0.82  |
| Earth   | Terrestrial | 1.000               | 1.00  |
| Mars    | Terrestrial | 0.532               | 0.11  |
| Jupiter | Giant       | 11.209              | 317.8 |
| Saturn  | Giant       | 9.449               | 95.2  |
| Uranus  | Giant       | 4.007               | 14.6  |
| ...     | ...         | ...                 | ...   |

|Name|

### Inclusion Dependencies

#### Foreign key candidates

| Sign        | Domicile |
|-------------|----------|
| Aries       | Mars     |
| Taurus      | Venus    |
| Gemini      | Mercury  |
| Cancer      | Moon     |
| Leo         | Sun      |
| Virgo       | Mercury  |
| Libra       | Venus    |
| Scorpio     | Pluto    |
| Sagittarius | Jupiter  |
| Capricorn   | Saturn   |
| Aquarius    | Uranus   |
| ...         | ...      |

| Name    | Type        |
|---------|-------------|
| Mercury | Terrestrial |
| Venus   | Terrestrial |
| Earth   | Terrestrial |
| Mars    | Terrestrial |
| Jupiter | Giant       |
| Saturn  | Giant       |
| Uranus  | Giant       |
| ...     | ...         |

Domicile  $\subseteq$  Name

### Functional Dependencies

#### Normalization criterion

| Name    | Atmosphere                            | Rings |
|---------|---------------------------------------|-------|
| Mercury | minimal                               | no    |
| Venus   | CO <sub>2</sub> , N <sub>2</sub>      | no    |
| Earth   | N <sub>2</sub> , O <sub>2</sub> , Ar  | no    |
| Mars    | CO <sub>2</sub> , N <sub>2</sub> , Ar | no    |
| Jupiter | H <sub>2</sub> , He                   | yes   |
| Saturn  | H <sub>2</sub> , He                   | yes   |
| Uranus  | H <sub>2</sub> , He                   | yes   |
| ...     | ...                                   | ...   |

Atmosphere  $\rightarrow$  Rings

#### Data Profiling with Metanome

Thorsten Papebrock,  
PhD Candidate,  
17<sup>th</sup> November, 2014  
Chart 17

# Metanome

## The exercise: Build your own profiling tool

The screenshot shows the Metanome web application running on localhost:8080. The interface includes a navigation menu with 'Data Sources', 'Algorithms', 'Run Configuration', 'Results', and 'About'. The main content area lists several algorithm categories: 'Unique Column Combinations', 'Conditional Unique Column Combinations', 'Functional Dependencies', 'Inclusion Dependencies', and 'Basic Statistics'. Below these is a form titled 'Add A New Algorithm' with fields for 'File Name', 'Algorithm Name', 'Author', and 'Description', along with 'Refresh' and 'Save' buttons. Four orange callout boxes are overlaid on the image, pointing to specific features: 'Exercise 1: UCC algorithm' points to 'Unique Column Combinations'; 'Exercise 2: IND algorithm' points to 'Inclusion Dependencies'; 'Exercise 3: FD algorithm' points to 'Functional Dependencies'; and 'Exercise 4: Duplicate Detection' points to the 'Add A New Algorithm' form.

localhost:8080

**M**etanome

Data Sources Algorithms Run Configuration Results About

Unique Column Combinations

Conditional Unique Column Combinations

Functional Dependencies

Inclusion Dependencies

Basic Statistics

Add A New Algorithm

File Name -- Refresh

Algorithm Name

Author

Description

Save

**Exercise 1: UCC algorithm**

**Exercise 2: IND algorithm**

**Exercise 3: FD algorithm**

**Exercise 4: Duplicate Detection**

### Data Profiling with Metanome

Thorsten Papebrock,  
PhD Candidate,  
17<sup>th</sup> November, 2014  
Chart 18

# Exercise 2

## Discovery of foreign-key candidates

---

### Exercise 2

#### Inclusion Dependencies

- Deadline: **Monday, 01.12.14**
- The admission to the exam requires *all* exercises to be solved.
- The exercises should be solved in teams of two students.
- The Metanome project is available at GitHub:  
<https://github.com/HPI-Information-Systems/Metanome>
- The datasets and supplemental material can be found at network drive S:  
`\\fs3\bbs\DPDC`
- The submission system can be found at:  
<https://www.dcl.hpi.uni-potsdam.de/submit/>
- To solve an exercise, please submit a zip file containing the following items:
  - **<algorithm\_name>.jar**: An executable Metanome algorithm.
  - **<algorithm\_name>.zip**: The algorithm's source code (maven project).
  - **<algorithm\_name>.docu.pdf**: Short documentation of the algorithm.
  - **<algorithm\_name>\_pres.pptx/ppt/pdf**: Two slides presentation of the algorithm.

#### Data Profiling with Metanome

Thorsten Papenbrock,  
PhD Candidate,  
17<sup>th</sup> November, 2014  
Chart **19**

# Exercise 2

## Discovery of foreign-key candidates

### Task 1: Inclusion Dependencies - A discovery algorithm

Write an algorithm that discovers *all unary* inclusion dependencies on the given datasets. The rules for your implementation are as follows:

- The algorithm discovers *exact* results, so no approximate or fuzzy results are allowed.
- The algorithm is not allowed to use parallelization.
- The algorithm implements the Metanome interface and is compatible with Metanome.
- The algorithm takes an arbitrary number of tables as input.
- The algorithm ignores NULL values.

Note that in contrast UCC and FD discovery, IND discovery uses *multiple tables*. INDs should therefore be discovered within and between tables. You can re-implement an existing algorithm from literature or find your own algorithm. To test and evaluate the algorithm, use the datasets provided on the network share. Your algorithm should at least be able to process the WDC dataset! Before submitting your algorithm, check that it correctly executes within Metanome!

*BONUS TASK: If you like to dive deeper, you can try to find n-ary inclusion dependencies as well. Can you think of pruning rules for the n-ary discovery? You can also try parallelization or approximate strategies on your algorithm. If you made changes to your main algorithm, please submit them as a separate algorithm, e.g. <algorithm\_name>\_nary.jar*

### Data Profiling with Metanome

Thorsten Papebrock,  
PhD Candidate,  
17<sup>th</sup> November, 2014  
Chart 20

# Exercise 2

## Discovery of foreign-key candidates

```
public class RelationalMyIndDetector
    extends MyIndDetector
    implements InclusionDependencyAlgorithm, RelationalInputParameterAlgorithm {

    public enum Identifier {INPUT_GENERATOR}

    @Override
    public ArrayList<ConfigurationRequirement> getConfigurationRequirements() {
        ArrayList<ConfigurationRequirement> configs = new ArrayList<ConfigurationRequirement>(1);
        configs.add(new ConfigurationRequirementRelationalInput(
            RelationalMyIndDetector.Identifier.INPUT_GENERATOR.name(), ConfigurationRequirement.ARBITRARY_NUMBER_OF_VALUES));
        return configs;
    }
    @Override
    public void setRelationalInputConfigurationValue(String identifier, RelationalInputGenerator... values)
        throws AlgorithmConfigurationException {
        if (RelationalMyIndDetector.Identifier.INPUT_GENERATOR.name().equals(identifier)) {
            this.relationalInputGenerators = values;
            // Useful for relation naming: this.relationalInputGenerators[0].generateNewCopy().relationName()
        }
        else
            this.handleUnknownConfiguration(identifier, CollectionUtils.concat(values, ","));
    }
    @Override
    public void setResultReceiver(InclusionDependencyResultReceiver resultReceiver) {
        this.resultReceiver = resultReceiver;
    }
    protected void handleUnknownConfiguration(String identifier, String value)
        throws AlgorithmConfigurationException {
        throw new AlgorithmConfigurationException("Unknown configuration: " + identifier + " -> " + value);
    }
    @Override
    public void execute() throws AlgorithmExecutionException {
        super.execute();
    }
}
```

### Data Profiling with Metanome

Thorsten Papebrock,  
PhD Candidate,  
17<sup>th</sup> November, 2014  
Chart 21

# Exercise 2

## Discovery of foreign-key candidates

### Task 2: Documentation

Write a short (max one A4 page) documentation for your algorithm describing the algorithm that you implemented:

- a) Describe the algorithm's basic idea. How does the algorithm cope with the complexity of the given task?
- b) If you used an algorithm from literature, provide a reference to the according publication.
- c) If you came up with an own approach, provide one or two arguments why it is or could be better than related algorithms.
- d) If your algorithm implements an adaption or optimization of existing approaches, describe these briefly.
- e) State if your algorithm (jar) should be published as Metanome algorithm providing you as the authors.
- f) *If you solved a bonus task, please discuss your findings here as well.*

In the same document, answer the following questions:

- a) How many inclusion dependencies did your algorithm find on the provided datasets?
- b) How long did the discovery take on each dataset and what machine did you use?
- c) Did you discover any limitations of your approach (e.g. runtime or memory consumption) that made computing a certain dataset impossible?
- d) Why is it a good idea to ignore NULL values for IND discovery?

### Data Profiling with Metanome

Thorsten Papebrock,  
PhD Candidate,  
17<sup>th</sup> November, 2014  
Chart **22**

# Exercise 2

## Discovery of foreign-key candidates

---

### Task 3: Presentation

Prepare two slides for a short, 5 min presentation of your algorithm in the lecture. One slide about the algorithm and one slide about its performance. Here are some ideas for these slides:

- a) Explain the idea of your approach and why you think it works well in comparison to others.
- b) How did your algorithm perform on the given datasets?
- c) Did you make any interesting observations?
- d) What lets your algorithm crash?
- e) How hard was it to implement your algorithm?
- f) *This is also a good time to introduce an optimization or adaption of your algorithm.*

Note that each team will present its work once!

### Data Profiling with Metanome

Thorsten Papenbrock,  
PhD Candidate,  
17<sup>th</sup> November, 2014  
Chart **23**

# Exercise 2

## Hinweise zum Referenzalgorithmus BINDER

---

- **INPUT\_GENERATPOR: WDC\_planets, WDC\_astronomical, ...**

- <die Tabellen, die durchsucht werden sollen; nur angegebene Tabellen werden genutzt>

- **INPUT\_RO\_LIMIT: -1**

- <der Parameter erlaubt die künstliche Verkürzung aller Tabellen auf X Zeilen; -1 nimmt alle Tupel>

- **TEMP\_FOLDER\_PATH: temp\**

- <Verzeichnis für temporäre Daten; relativ zum Verzeichnis der Metanome-Instanz>

- **CLEAN\_TEMP: true**

- <sollen die temporären Daten gelöscht werden nach der Ausführung?>

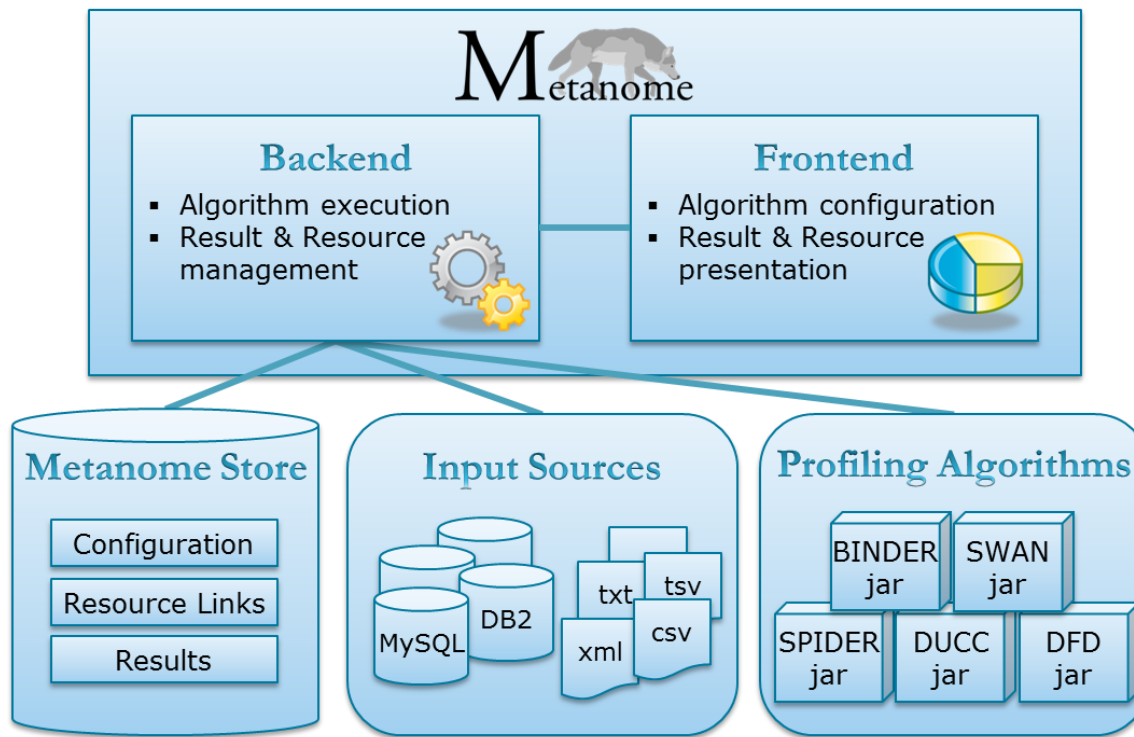
- **DETECT\_NARY: false**

- <BINDER kann auch n-ary INDs finden; für eure Übung braucht ihr aber nur unary INDs>

### Data Profiling with Metanome

Thorsten Papenbrock,  
PhD Candidate,  
17<sup>th</sup> November, 2014  
Chart 24





## Data Profiling with Metanome

[www.metanome.de](http://www.metanome.de)

Thorsten Papenbrock

PhD Candidate

Hasso-Plattner-Institute