



## Data Profiling with Metanome

Exercise: FDs

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 2

## Discovery of foreign key candidates

---



- All teams have passed the exercise:
  - 35 submissions
  - No duplicate algorithm names!
  - No incorrect results (on given datasets)!
  - No import errors in Metanome (apart from execution errors)!

### **Data Profiling with Metanome**

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

# Exercise 2

## Short presentations – Part 1

---



### **Data Profiling with Metanome**

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

# Exercise 2

## 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 1.8**

### **Data Profiling with Metanome**

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

AlexoFredInclusions  
curunir  
dennis\_marius.ind  
DJ\_IND  
dpdc-cnms-ind  
DreamteamInd  
FrohnOttoInc-preconfigured  
GansGans  
HorLehInDeBinder  
IncDep  
InclusionDetector  
ind\_grundke\_wiese  
IND\_Kirsten\_Zwerg  
IND\_schaeffer\_zoellner\_hashmap  
IND\_SPIRO  
IndBotheReissaus  
INDianaJonesRaiderOfTheLostInclusion  
INDIE  
INDJRFM  
IND-MoritzChristian  
IndPerchykSchmidt  
Jung  
JungRohloffInd  
Kankra  
klinger\_marten\_inclusion  
LucieKerstinIND  
MMINDs  
MuchDiscoVeryDisco  
MyInd  
PCInd  
PracticeIndJoerke  
RT\_BindBuck  
SBMMIND  
smart-data-cat-IND  
Tsun12Ind  
YuckInd

## WDC-planets dataset

### Data Profiling with Metanome

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

AlexoFredInclusions

**curunir**

dennis\_marius.ind

DJ\_IND

dpdc-cnms-ind

DreamteamInd

FrohnOttoInc-preconfigured

GansGans

HorLehInDeBinder

IncDep

InclusionDetector

ind\_grundke\_wiese

IND\_Kirsten\_Zwerg

IND\_schaeffer\_zoellner\_hashmap

IND\_SPIRO

IndBotheReissaus

INDianaJonesRaiderOfTheLostInclusion

INDIE

INDJRFM

IND-MoritzChristian

IndPerchykSchmidt

Jung

JungRohloffInd

Kankra

klinger\_marten\_inclusion

LucieKerstinIND

MMINDs

MuchDiscoVeryDisco

**MyInd**

PCInd

PracticeIndJoerke

RT\_BindBuck

SBMMIND

smart-data-cat-IND

Tsun12Ind

YuckInd

**Only One InputGenerator**

**WDC-planets dataset**

**No working parameterization found**

**Data Profiling with Metanome**

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

AlexoFredInclusions

dennis\_marius.ind

DJ\_IND

dpdc-cnms-ind

DreamteamInd

FrohnOttoInc-preconfigured

GansGans

HorLehInDeBinder

IncDep

InclusionDetector

ind\_grundke\_wiese

IND\_Kirsten\_Zwerg

IND\_schaeffer\_zoellner\_hashmap

IND\_SPIRO

IndBotheReissaus

INDianaJonesRaiderOfTheLostInclusion

INDIE

INDJRFM

IND-MoritzChristian

IndPerchykSchmidt

Jung

JungRohloffInd

Kankra

klinger\_marten\_inclusion

LucieKerstinIND

MMINDs

MuchDiscoVeryDisco

PCInd

PracticeIndJoerke

RT\_BindBuck

SBMMIND

smart-data-cat-IND

Tsun12Ind

YuckInd

## TPC-H dataset

(Customer, Supplier, Nation, Part)

### Data Profiling with Metanome

Thorsten Papenbrock,  
PhD Candidate,  
17<sup>th</sup> November, 2014  
Chart 8



**AlexoFredInclusions**

**OutOfMemory**

dennis\_marius.ind

**DJ\_IND**

**OutOfMemory**

dpdc-cnms-ind

DreamteamInd

FrohnOttoInc-preconfigured

GansGans

**HorLehInDeBinder**

**OutOfMemory**

IncDep

**InclusionDetector**

**OutOfMemory**

ind\_grundke\_wiese

**IND\_Kirsten\_Zwerg**

**>15 min**

**IND\_schaeffer\_zoellner\_hashmap**

**OutOfMemory**

IND\_SPIRO

IndBotheReissaus

INDianaJonesRaiderOfTheLostInclusion

**INDIE**

**OutOfMemory**

**INDJRFM**

**Wrong Result**

IND-MoritzChristian

**IndPerchykSchmidt**

**OutOfMemory**

**Jung**

**OutOfMemory**

**JungRohloffInd**

**OutOfMemory**

Kankra

**klinger\_marten\_inclusion**

**>15 min**

**LucieKerstinIND**

**OutOfMemory**

MMINDs

**MuchDiscoVeryDisco**

**OutOfMemory**

**PCInd**

**OutOfMemory**

**PracticeIndJoerke**

**>15 min**

RT\_BindBuck

**SBMMIND**

**OutOfMemory**

**smart-data-cat-IND**

**>15 min**

Tsun12Ind

**YuckInd**

**>15 min**

## TPC-H dataset

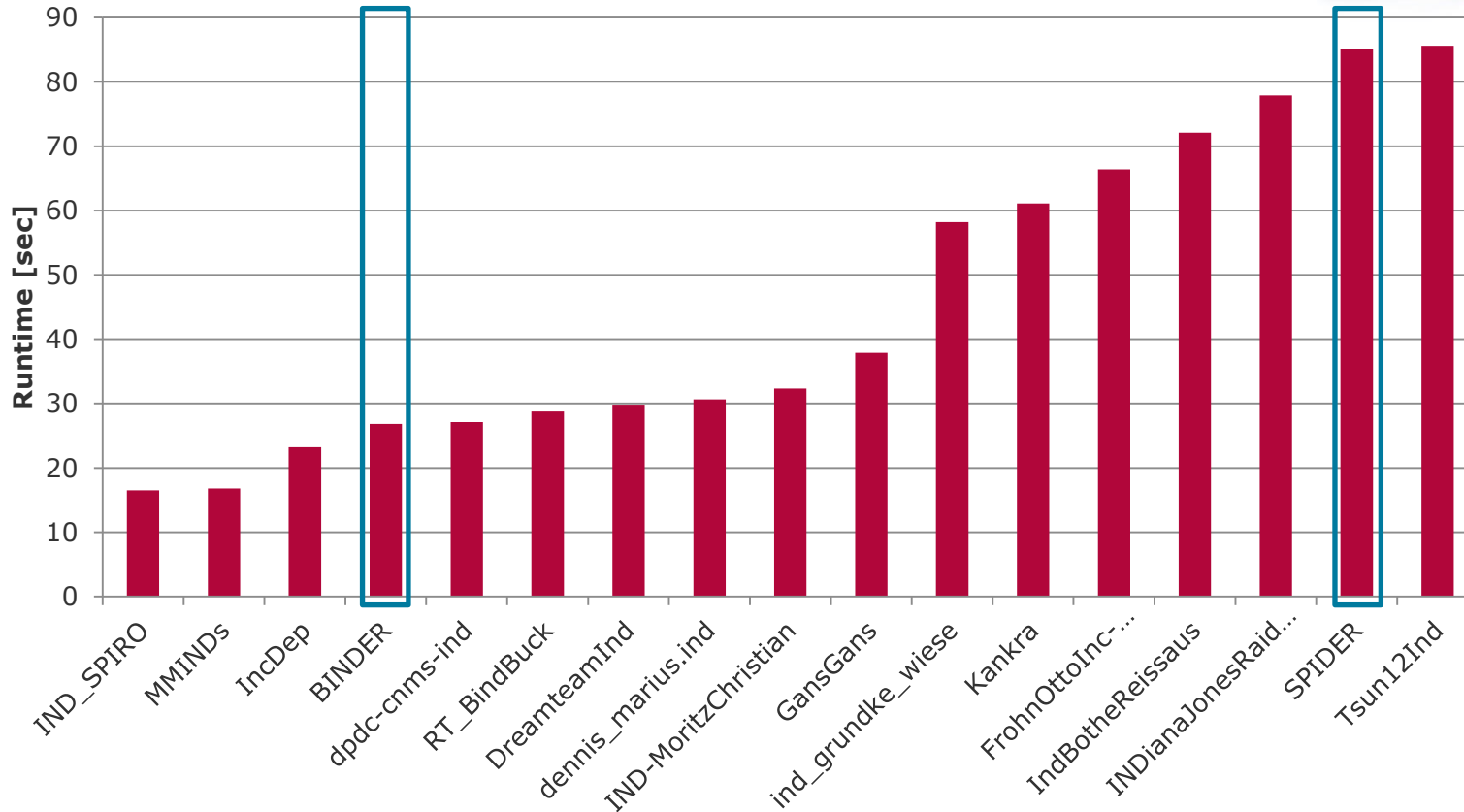
(Customer, Supplier, Nation, Part)

### Data Profiling with Metanome

Thorsten Papenbrock,  
PhD Candidate,  
17<sup>th</sup> November, 2014  
Chart 9

# Exercise 2

## Our evaluation – TPC-H (reduced)



### Data Profiling with Metanome

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

dennis\_marius.ind

dpdc-cnms-ind

DreamteamInd

FrohnOttoInc-preconfigured

GansGans

IncDep

ind\_grundke\_wiese

IND\_SPIRO

IndBotheReissaus

INDianaJonesRaiderOfTheLostInclusion

IND-MoritzChristian

Kankra

MMINDs

RT\_BindBuck

Tsun12Ind

**TPC-H dataset**

(complete)

**Data Profiling with  
Metanome**

ThorstenPapenbrock,  
PhD Candidate,  
17<sup>th</sup> November, 2014  
Chart **11**

<b>dennis_marius.ind</b>	<b>OutOfMemory</b>
<b>dpdc-cnms-ind</b> DreamteamInd FrohnOttoInc-preconfigured GansGans	<b>OutOfMemory</b>
IncDep	
ind_grundke_wiese	
IND_SPIRO IndBotheReissaus <b>INDianaJonesRaiderOfTheLostInclusion</b>	<b>Too many open files</b>
IND-MoritzChristian	
<b>Kankra</b>	<b>OutOfMemory</b>
MMINDs	
RT_BindBuck	
<b>Tsun12Ind</b>	<b>OutOfMemory</b>

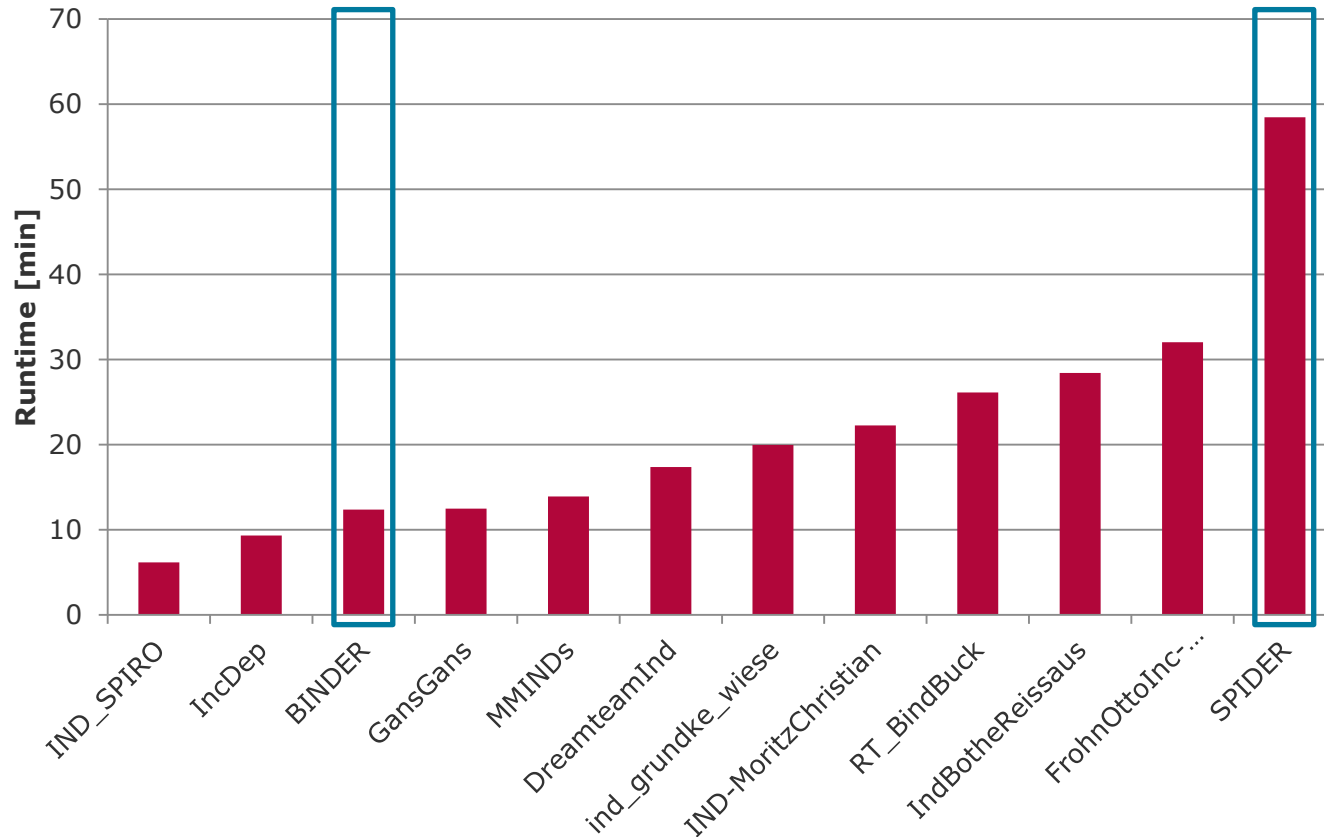
**TPC-H dataset**  
(complete)

### Data Profiling with Metanome

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

# Exercise 2

## Our evaluation – TPC-H (complete)



### Data Profiling with Metanome

Thorsten Papebrock,  
PhD Candidate,  
17<sup>th</sup> November, 2014  
Chart 13

DreamteamInd  
FrohnOttoInc-preconfigured  
GansGans

IncDep

ind\_grundke\_wiese

IND\_SPIRO  
IndBotheReissaus

IND-MoritzChristian

MMINDs

RT\_BindBuck

**TPC-H and PLISTA dataset**  
(complete and sample; 1GB RAM)

## Data Profiling with Metanome

ThorstenPapenbrock,  
PhD Candidate,  
17<sup>th</sup> November, 2014  
Chart **14**

DreamteamInd  
FrohnOttoInc-preconfigured  
GansGans

**OutOfMemory**  
>12 h

Also happens to SPIDER and BINDER:  
 $\text{size}(\text{candidates}) > \text{size}(\text{RAM})$

IncDep

**OutOfMemory**

ind\_grundke\_wiese

IND\_SPIRO  
IndBotheReissaus

**OutOfMemory**  
**OutOfMemory**

**TPC-H and PLISTA dataset**  
(complete and sample; 1GB RAM)

IND-MoritzChristian

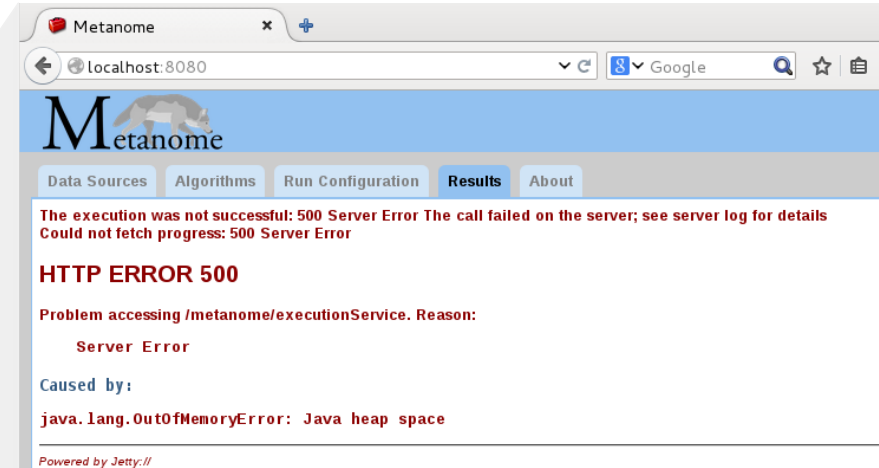
NumberFormatException

MMINDs

**OutOfMemory**

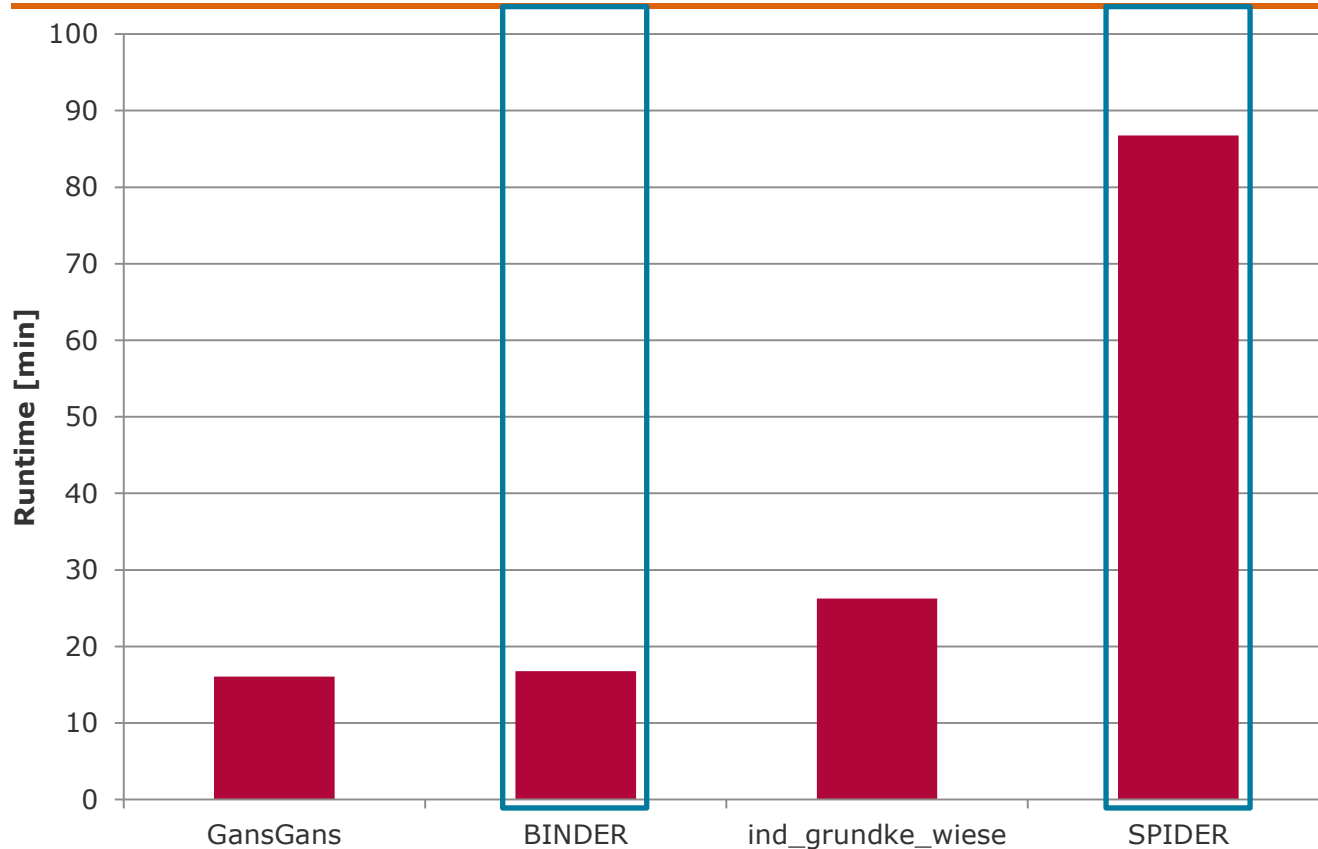
RT\_BindBuck

**OutOfMemory**



# Exercise 2

## Our evaluation – TPC-H and PLISTA



### Data Profiling with Metanome

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



# Exercise 2

## Short presentations – Part 2

---



### **Data Profiling with Metanome**

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

# 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 18

# 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

Metanome

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 19

# Exercise 3

## Discovery of functional dependencies

---

### Exercise 3

### Functional Dependencies

- Deadline: **Monday, 05.01.15**
- 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 **20**

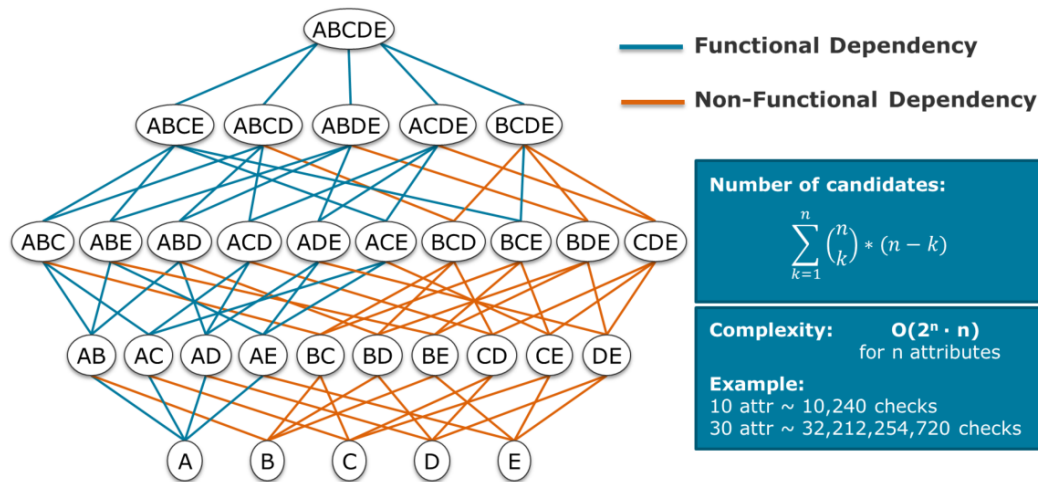
# Exercise 3

## Discovery of functional dependencies

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

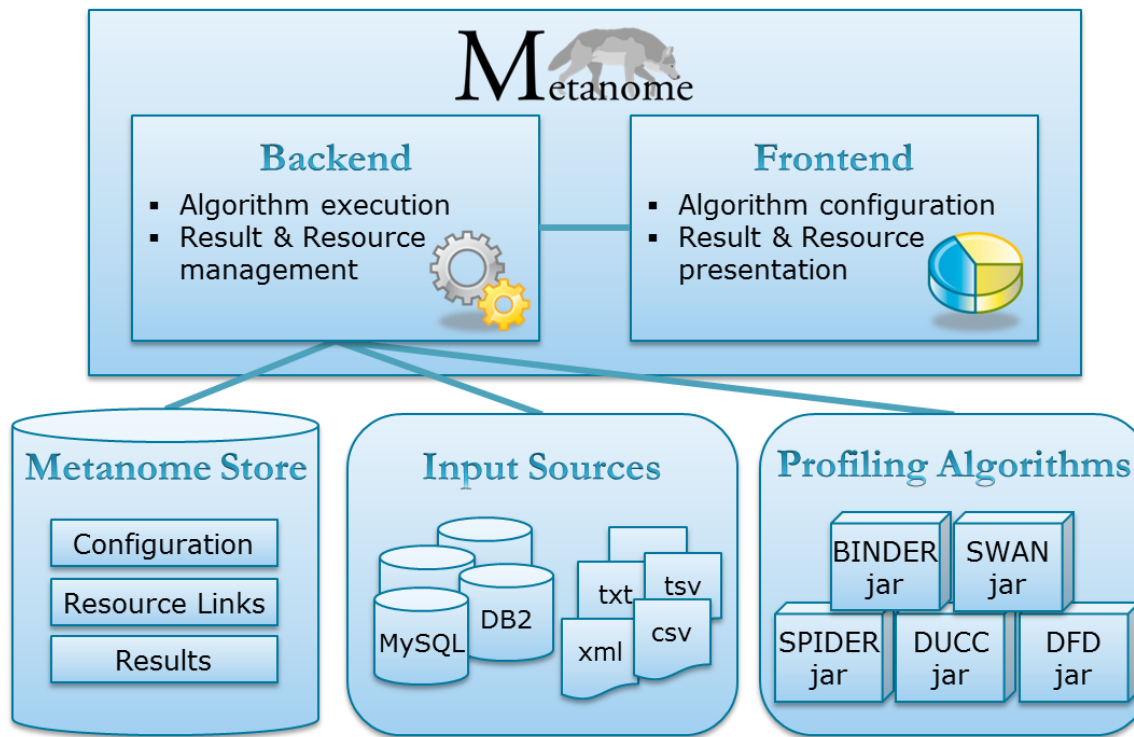
Write an algorithm that discovers *all unary and n-ary* functional dependencies on the given datasets. The rules for your implementation are as follows:

- 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.
- For the NULL semantic, assume  $\text{NULL} = \text{NULL}$ .



### Data Profiling with Metanome

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



## Data Profiling with Metanome

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

Thorsten Papenbrock

PhD Candidate

Hasso-Plattner-Institute