

Data Profiling

Effiziente Entdeckung Struktureller Abhängigkeiten

Dr. Thorsten Papenbrock
Information Systems Group, HPI

Art. 15 DSGVO Auskunftsrecht der betroffenen Person

Knowledge Discovery

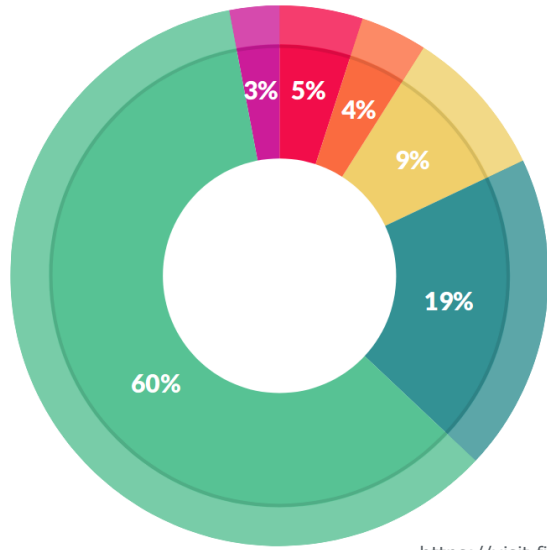
What data do you have?

- (1) Die betroffene Person hat das Recht, von dem Verantwortlichen eine Bestätigung darüber zu verlangen, ob sie betreffende personenbezogene Daten verarbeitet werden; ist dies der Fall, so hat sie ein Recht auf Auskunft über diese personenbezogenen Daten und auf folgende Informationen:
 - a) die Verarbeitungszwecke;
 - b) die Kategorien personenbezogener Daten, die verarbeitet werden;
 - c) die Empfänger oder Kategorien von Empfängern, gegenüber denen die personenbezogenen Daten offengelegt worden sind oder noch offengelegt werden, insbesondere bei Empfängern in Drittländern oder bei internationalen Organisationen;
 - d) falls möglich die geplante Dauer, für die die personenbezogenen Daten gespeichert werden, oder, falls dies nicht möglich ist, die Kriterien für die Festlegung dieser Dauer;
 - e) das Bestehen eines Rechts auf Berichtigung oder Löschung der sie betreffenden personenbezogenen Daten oder auf Einschränkung der Verarbeitung durch den Verantwortlichen oder eines Widerspruchsrechts gegen diese Verarbeitung;
 - f) das Bestehen eines Beschwerderechts bei einer Aufsichtsbehörde;
 - g) wenn die personenbezogenen Daten nicht bei der betroffenen Person erhoben werden, alle verfügbaren Informationen über die Herkunft der Daten;
 - h) das Bestehen einer automatisierten Entscheidungsfindung einschließlich Profiling gemäß [Artikel 22](#) Absätze 1 und 4 und – zumindest in diesen Fällen – aussagekräftige Informationen über die involvierte Logik sowie die Tragweite und die angestrebten Auswirkungen einer derartigen Verarbeitung für die betroffene Person.
- (2) Werden personenbezogene Daten an ein Drittland oder an eine internationale Organisation übermittelt, so hat die betroffene Person das Recht, über die geeigneten Garantien gemäß [Artikel 46](#) im Zusammenhang mit der Übermittlung unterrichtet zu werden.
- (3) ¹ Der Verantwortliche stellt eine Kopie der personenbezogenen Daten, die Gegenstand der Verarbeitung sind, zur Verfügung. ² Für alle weiteren Kopien, die die betroffene Person beantragt, kann der Verantwortliche ein angemessenes Entgelt auf der Grundlage der Verwaltungskosten verlangen. ³ Stellt die betroffene Person den Antrag elektronisch, so sind die Informationen in einem gängigen elektronischen Format zur Verfügung zu stellen, sofern sie nichts anderes angibt.
- (4) Das Recht auf Erhalt einer Kopie gemäß Absatz 3 darf die Rechte und Freiheiten anderer Personen nicht beeinträchtigen.

Many companies do not know what data they have!

- **Decentralized** storage and retrieval
- **Heterogeneous** data formats and systems
- **Unconnected** sources
- **Lack of metadata and integrity constraints**
- Different **access rights**
- Data **quality issues**
- Complicated **business processes**
- Data **backups** and **archives**
- Data **acquisition** and **sharing**
- ...

CrowdFlower Data Science Report 2016



What data scientists spend the most time doing

- Building training sets: 3%
- Cleaning and organizing data: 60%
- Collecting data sets; 19%
- Mining data for patterns: 9%
- Refining algorithms: 4%
- Other: 5%

~80% on data preparation!

https://visit.figure-eight.com/rs/416-ZBE-142/images/CrowdFlower_DataScienceReport_2016.pdf

Knowledge Discovery

Data Analytics

Data scientists spend most of their time on data preparation!

- Multiple, heterogeneous data sources
- Lack of metadata and documentation
- Data quality issues
- Data acquisition and sharing
- ...

Deep Visual-Semantic Alignments for Generating Image Descriptions

Andrej Karpathy and Li Fei-Fei, Stanford University, TPAMI, 2015



"a young boy is holding a baseball bat."



"a cat is sitting on a couch with a remote control."



"a woman holding a teddy bear in front of a mirror."



"a horse is standing in the middle of a road."

Knowledge Discovery

Data Analytics

AI Systems

AI systems learn what they see and understand

AI systems learn erroneous, non-interpretable behavior!

- Data quality issues
- Insufficient training data
- Heterogeneous data formats and systems
- Lack of metadata and documentation
- ...

Data Engineering for Data Science



Knowledge Discovery

Data Analytics

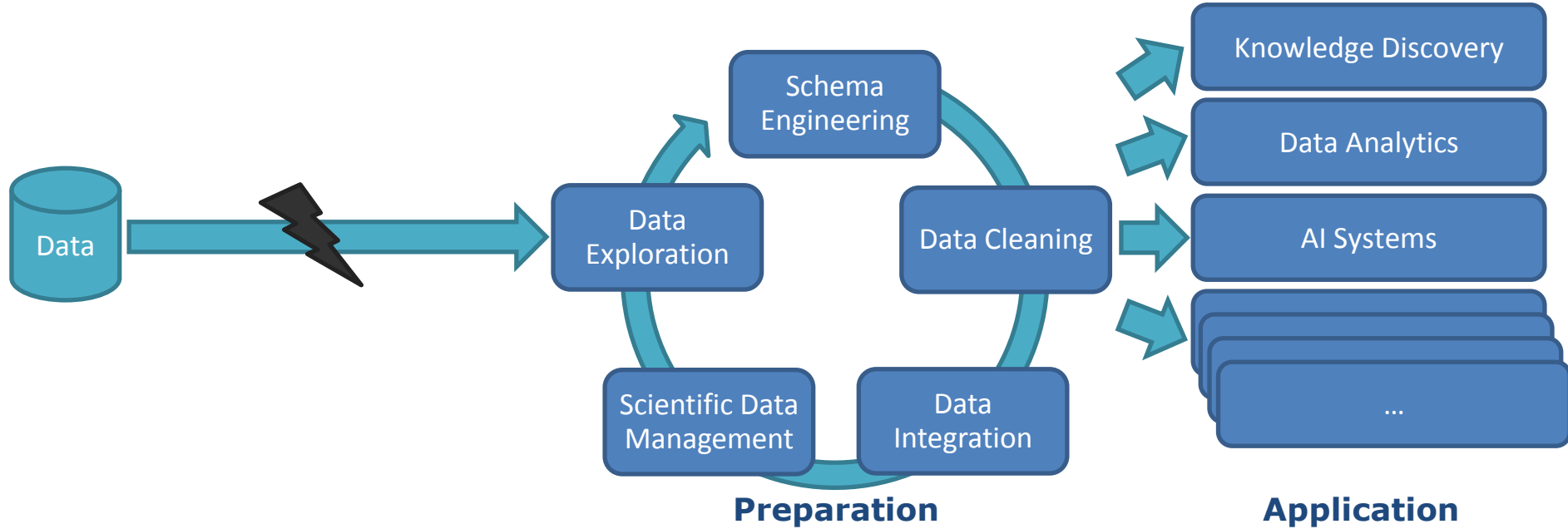
AI Systems

...

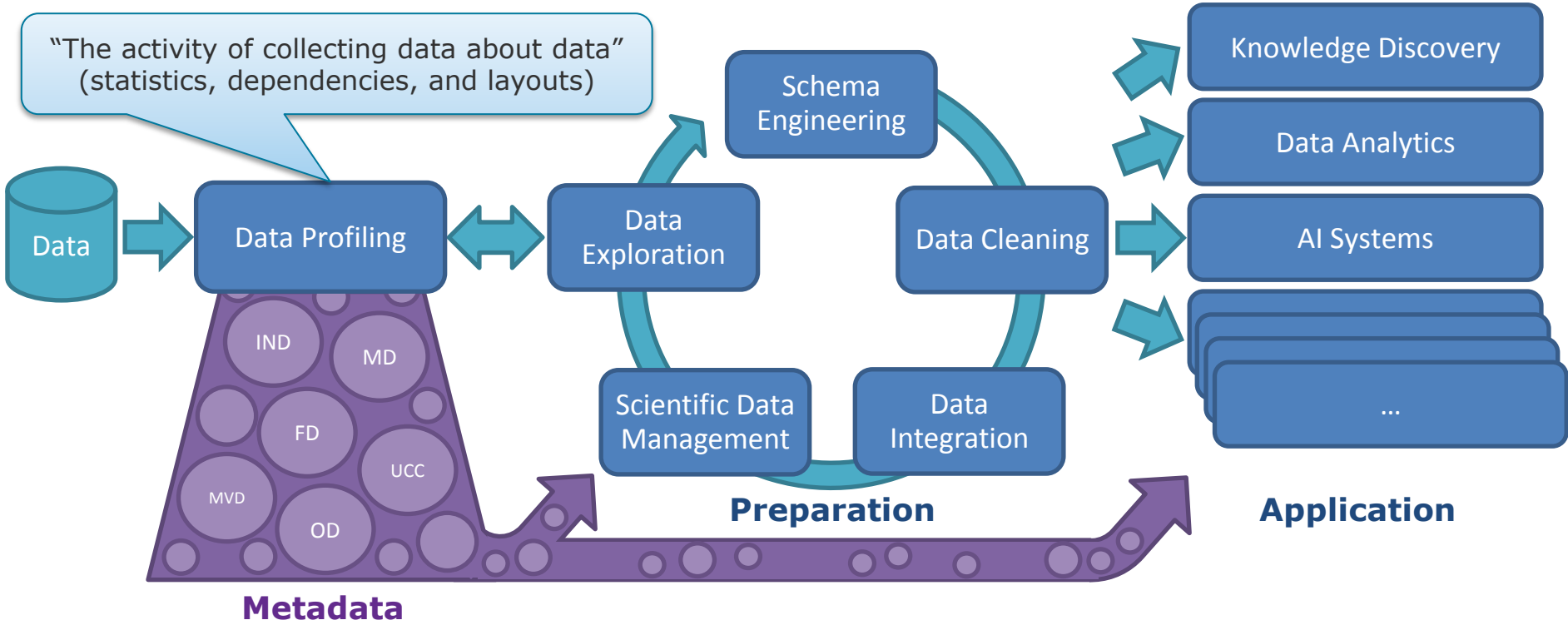
...

Application

Data Engineering for Data Science



Data Engineering for Data Science



Data Profiling



ID	Name	Evolution	Location	Sex	Weight	Size	Type	Weak	Strong	Special
25	Pikachu	Raichu	Viridian Forest	m/w	6.0	0.4	electric	ground	water	false
27	Sandshrew	Sandslash	Route 4	m/w	12.0	0.6	ground	gras	electric	false
29	Nidoran	Nidorino	Safari Zone	m	9.0	0.5	poison	ground	gras	false
32	Nidoran	Nidorina	Safari Zone	w	7.0	0.4	poison	ground	gras	false
37	Vulpix	Ninetails	Route 7	m/w	9.9	0.6	fire	water	ice	false
38	Ninetails	null	null	m/w	19.9	1.1	fire	water	ice	true
63	Abra	Kadabra	Route 24	m/w	19.5	0.9	psychic	ghost	fighting	false
64	Kadabra	Alakazam	Cerulean Cave	m/w	56.5	1.3	psychic	ghost	fighting	false
130	Gyarados	null	Fuchsia City	m/w	235.0	6.5	water	electric	fire	false
150	Mewtwo	null	Cerulean Cave	null	122.0	2.0	psychic	ghost	fighting	true



Data Profiling

format

density

#null = 3
%null = 30

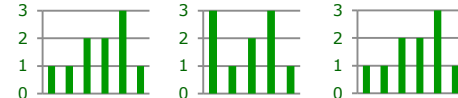
ranges

min = 0.4
max = 2.0

aggregations

sum = 14.3
avg = 1.43

distributions



size
= 10

ID	Name	Evolution	Location	Sex	Weight	Size	Type	Weak	Strong	Special
25	Pikachu	Raichu	Viridian Forest	m/w	6.0	0.4	electric	ground	water	false
27	Sandshrew	Sandslash	Route 4	m/w	12.0	0.6	ground	grass	electric	false
29	Nidoran	Nidorino	Safari Zone	m	9.0	0.5	poison	ground	grass	false
32	Nidoran	Nidorina	Safari Zone	w	7.0	0.4	poison	ground	grass	false
37	Vulpix	Ninetails	Route 7	m/w	9.9	0.6	fire	water	ice	false
38	Ninetails	null	null	m/w	19.9	1.1	fire	water	ice	true
63	Abra	Kadabra	Route 24	m/w	19.5	0.9	psychic	ghost	fighting	false
64	Kadabra	Alakazam	Cerulean Cave	m/w	56.5	1.3	psychic	ghost	fighting	false
130	Gyarados	null	Fuchsia City	m/w	235.0	6.5	water	electric	fire	false
150	Mewtwo	null	Cerulean Cave	null	122.0	2.0	psychic	ghost	fighting	true

INTEGER

CHAR(16)

CHAR(16)

CHAR(32)

CHAR(3)

FLOAT

FLOAT

CHAR(8)

CHAR(8)

CHAR(8)

BOOLEAN

data types

inclusion dependencies

Pokemon.Location \subseteq Location.Name

functional dependencies

Type \rightarrow Weak

ID	Name	Evolution	Location	Sex	Weight	Size	Type	Weak	Strong	Special
25	Pikachu	Raichu	Viridian Forest	m/w	6.0	0.4	electric	ground	water	false
27	Sandshrew	Sandslash	Route 4	m/w	12.0	0.6	ground	gras	electric	false
29	Nidoran	Nidorino	Safari Zone	m	9.0	0.5	poison	ground	gras	false
32	Nidoran	Nidorina	Safari Zone	w	7.0	0.4	poison	ground	gras	false
37	Vulpix	Ninetails	Route 7	m/w	9.9	0.6	fire	water	ice	false
38	Ninetails	null	null	m/w	19.9	1.1	fire	water	ice	true
63	Abra	Kadabra	Route 24	m/w	19.5	0.9	psychic	ghost	fighting	false
64	Kadabra	Alakazam	Cerulean Cave	m/w	56.5	1.3	psychic	ghost	fighting	false
130	Gyarados	null	Fuchsia City	m/w	235.0	6.5	water	electric	fire	false
150	Mewtwo	null	Cerulean Cave	null	122.0	2.0	psychic	ghost	fighting	true

{Name, Sex}

unique column combinations

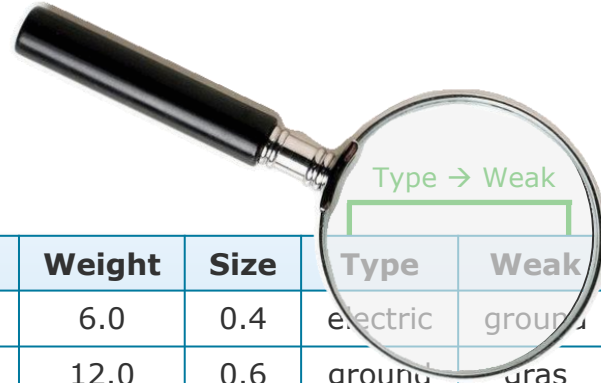
Weight \downarrow Size

order dependencies

Weak \neq Strong

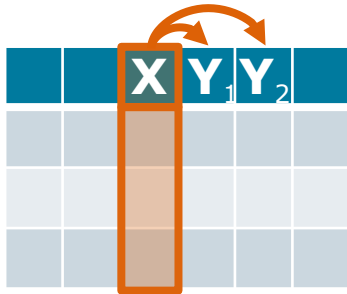
denial constraints

Data Profiling



ID	Name	Evolution	Location	Sex	Weight	Size	Type	Weak	Strong	Special
25	Pikachu	Raichu	Viridian Forest	m/w	6.0	0.4	electric	ground	water	false
27	Sandshrew	Sandslash	Route 4	m/w	12.0	0.6	ground	gras	electric	false
29	Nidoran	Nidorino	Safari Zone	m	9.0	0.5	poison	ground	gras	false
32	Nidoran	Nidorina	Safari Zone	w	7.0	0.4	poison	ground	gras	false
37	Vulpix	Ninetails	Route 7	m/w	9.9	0.6	fire	water	ice	false
38	Ninetails	null	null	m/w	19.9	1.1	fire	water	ice	true
63	Abra	Kadabra	Route 24	m/w	19.5	0.9	psychic	ghost	fighting	false
64	Kadabra	Alakazam	Cerulean Cave	m/w	56.5	1.3	psychic	ghost	fighting	false
130	Gyarados	null	Fuchsia City	m/w	235.0	6.5	water	electric	fire	false
150	Mewtwo	null	Cerulean Cave	null	122.0	2.0	psychic	ghost	fighting	true

Functional Dependencies



Definition: Given a relational instance r for a schema R . The **functional dependency** $X \rightarrow A$ with $X \subseteq R$ and $A \in R$ is valid in r , iff $\forall t_i, t_j \in r : t_i[X] = t_j[X] \Rightarrow t_i[A] = t_j[A]$.

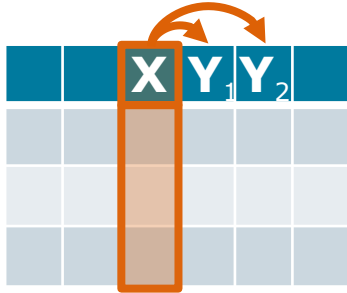
“The values in X functionally define the values in Y”

ID	Name	Size	Type	Weak	Strong	GYM	Leader	Reward
25	Pikachu	0.4	electric	ground	water	Vermillion	Lt. Surge	Thunder
26	Raichu	0.8	electric	ground	water	Vermillion	Lt. Surge	Thunder
29	Nidoran	0.5	poison	ground	gras	Viridian	Giovanni	Earth
37	Vulpix	0.6	fire	water	ice	null	null	null
38	Ninetails	1.1	fire	water	ice	null	null	null
63	Abra	0.9	psychic	ghost	fighting	null	null	null
64	Kadabra	1.3	psychic	ghost	fighting	Saffron	Sabrina	Marsh
65	Alakazam	1.5	psychic	ghost	fighting	Saffron	Sabrina	Marsh
150	Mewtwo	2.0	psychic	ghost	fighting	null	null	null

Type → Weak, Strong

GYM → Leader, Reward

Functional Dependencies




Definition: Given a relational instance r for a schema R . The **functional dependency** $X \rightarrow A$ with $X \subseteq R$ and $A \in R$ is valid in r , iff $\forall t_i, t_j \in r : t_i[X] = t_j[X] \Rightarrow t_i[A] = t_j[A]$.

“The values in X functionally define the values in Y”



ID	Name	Size	Type	GYM
25	Pikachu	0.4	electric	Vermillion
26	Raichu	0.8	electric	Vermillion
29	Nidoran	0.5	poison	Viridian
37	Vulpix	0.6	fire	null
38	Ninetails	1.1	fire	null
63	Abra	0.9	psychic	null
64	Kadabra	1.3	psychic	Saffron
65	Alakazam	1.5	psychic	Saffron
150	Mewtwo	2.0	psychic	null



Type	Weak	Strong
electric	ground	water
poison	ground	grass
fire	water	ice
psychic	ghost	fighting



GYM	Leader	Reward
Vermillion	Lt. Surge	Thunder
Viridian	Giovanni	Earth
Saffron	Sabrina	Marsh



2013

2014

2015

2016

2017

2018

TANE FUN FD_MINE DFD DEP-MINER FASTFDS FDEP

[TANE] *TANE: An efficient algorithm for discovering functional and approximate dependencies*,
Ykä Huhtala, Juha Kärkkäinen, Pasi Porkka and Hannu Toivonen, The Computer Journal, 1999.

[FUN] *FUN: An efficient algorithm for mining functional and embedded dependencies*,
Noël Novelli and Rosine Cicchetti, ICDT, 2001.

[FD_MINE] *FD Mine: discovering functional dependencies in a database using equivalences*,
Hong Yao, Howard J Hamilton and Cory J Butz, ICDM, 2002.

[DFD] *DFD: Efficient Functional Dependency Discovery*,
Ziawasch Abedjan, Patrick Schulze and Felix Naumann, CIKM, 2014.

[DEP-MINER] *Efficient discovery of functional dependencies and Armstrong relations*,
Stéphane Lopes, Jean-Marc Petit and Lotfi Lakhal, EDBT, 2000.

[FASTFDS] *FastFDS: A heuristic-driven, depth-first algorithm for mining functional dependencies from relation instances*,
Catharine Wyss, Chris Giannella and Edward Robertson, DaWaK, 2001.

[FDEP] *Database dependency discovery: a machine learning approach*,
Peter A Flach and Iztok Sarnik, AI Communications, 1999.

2013

2014

2015

2016

2017

2018

HPI

Hasso
Plattner
Institut

Dataset	Columns [#]	Rows [#]	Size [KB]	FDs [#]	TANE	FUN	FD_MINE	DFD	DEP-MINER	FASTFDs	FDEP
iris	5	150	5	4							
balance-scale	5	625	7	1							
chess	7	28,056	519	1							
abalone	9	4,177	187	137							
nursery	9	12,960	1,024	1							
breast-cancer	11	699	20	46							
bridges	13	108	6	142							
echocardiogram	13	132	6	538							
adult	14	48,842	3,528	78							
letter	17	20,000	695	61							
ncvoter	19	1,000	151	758							
hepatitis	20	155	8	8,250							
horse	27	368	25	128,726							
fd-reduced-30	30	250,000	69,581	89,571							
plista	63	1,000	568	178,152							
flight	109	1,000	575	982,631							
uniprot	223	1,000	2,439	unknown							

2013

2014

2015

2016

2017

2018

Dataset	Columns [#]	Rows [#]	Size [KB]	FDs [#]	TANE	FUN	FD_MINE	DFD	DEP-MINER	FASTFDs	FDEP
iris	5	150	5	4	1.1	0.1	0.2	0.2	0.2	0.2	0.1
balance-scale	5	625	7	1	1.2	0.1	0.2	0.3	0.3	0.3	0.2
chess	7	28,056	519	1	2.9	1.1	3.8	1.0	174.6	164.2	125.5
abalone	9	4,177	187	137	2.1	0.6	1.8	1.1	3.0	2.9	3.8
nursery	9	12,960	1,024	1	4.1	1.8	7.1	0.9	121.2	118.9	46.8
breast-cancer	11	699	20	46	2.3	0.6	2.2	0.8	1.1	1.1	0.5
bridges	13	108	6	142	2.2	0.6	4.2	0.9	0.5	0.6	0.2
echocardiogram	13	132	6	538	1.6	0.4	69.9	1.2	0.5	0.5	0.2
adult	14	48,842	3,528	78	67.4	111.6	531.5	5.9	6039.2	6033.8	860.2
letter	17	20,000	695	61	260.0	529.0	7204.8	6.0	1090.0	1015.5	291.3
ncvoter	19	1,000	151	758	4.3	4.0	ML	5.1	11.4	1.9	1.1
hepatitis	20	155	8	8,250	12.2	175.9	ML	326.7	5576.5	9.5	0.8
horse	27	368	25	128,726	457.0	TL	ML	TL	TL	385.8	7.2
fd-reduced-30	30	250,000	69,581	89,571	41.1	77.7	ML	TL	377.2	382.4	TL
plista	63	1,000	568	178,152	ML	ML	ML	TL	TL	TL	26.9
flight	109	1,000	575	982,631	ML	ML	ML	TL	TL	TL	216.5
uniprot	223	1,000	2,439	unknown	ML	ML	ML	TL	TL	TL	ML

Results larger than 1,000 FDs are only counted

TL: time limit of 4 hours exceeded**ML:** memory limit of 100GB exceeded

Table 1: Runtimes in seconds for several real-world datasets

2013

2014

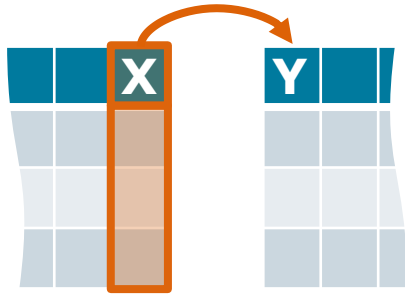
2015

2016

2017

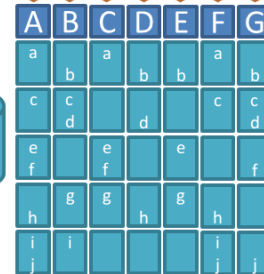
2018

Inclusion Dependencies

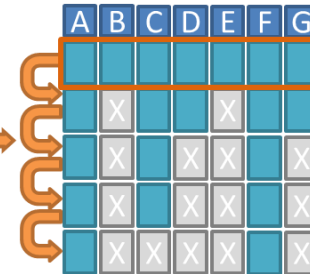


Divide

Rel. 1				Rel. 2		
A	B	C	D	E	F	G
a	b	g	d	e	g	c
b	c	e	b	h	i	b
c	a	g	b	i	c	a
a	j	b	b	e	j	b
j	i	c	d	b	c	f
e	f	g	a	g	e	c
f	g	a	f	e	f	d
h	i	e	d	g	h	d
i	j	f	d	h	i	c
a	a	a	b	g	j	d
a	a	a	b	g	j	d



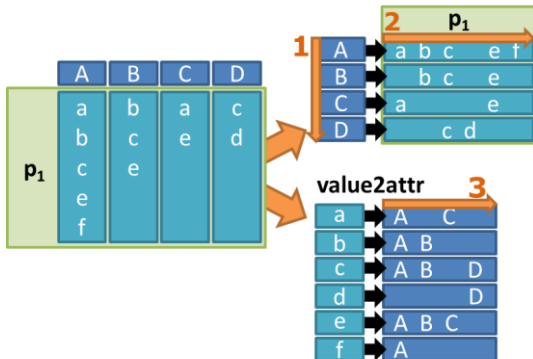
Conquer



validation?

F ⊆ A

attr2value



look up

A → a → A,C
 A → b → A,B
 B → c → A,B,D
 B → e → A,B,C
 D → d → D

	A	B	C	D
B,C,D	A,C,D	A,C,D	A,B,D	A,B,C
C	A,C,D	A	A	A,B,C
B	-	A	A	A,B
B	-	A	A	A,B
D	-	A	A	-

2013

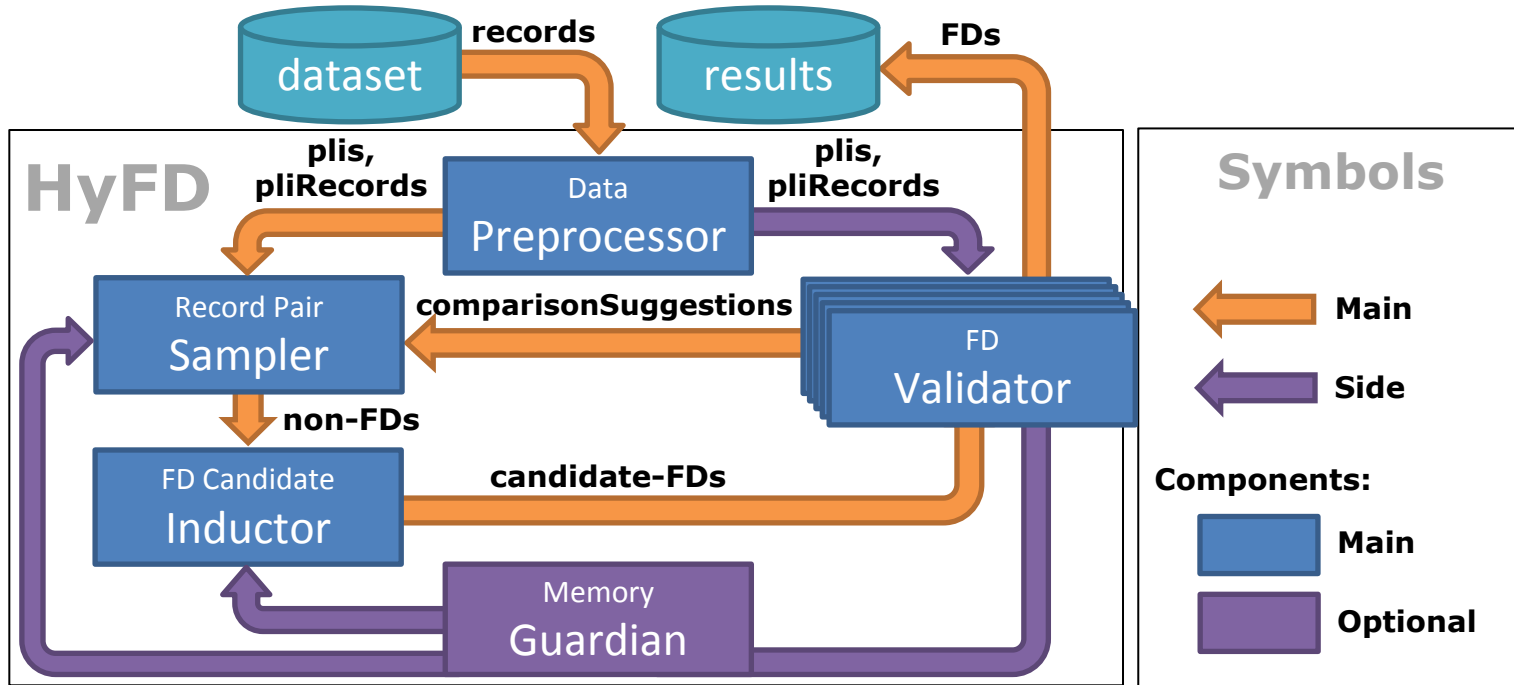
2014

2015

2016

2017

2018



2013

2014

2015

2016

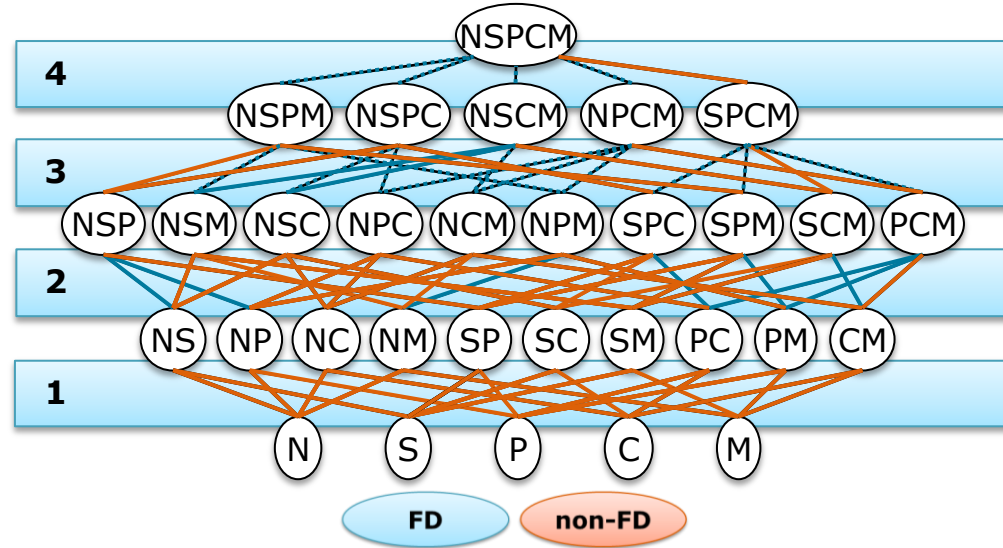
2017

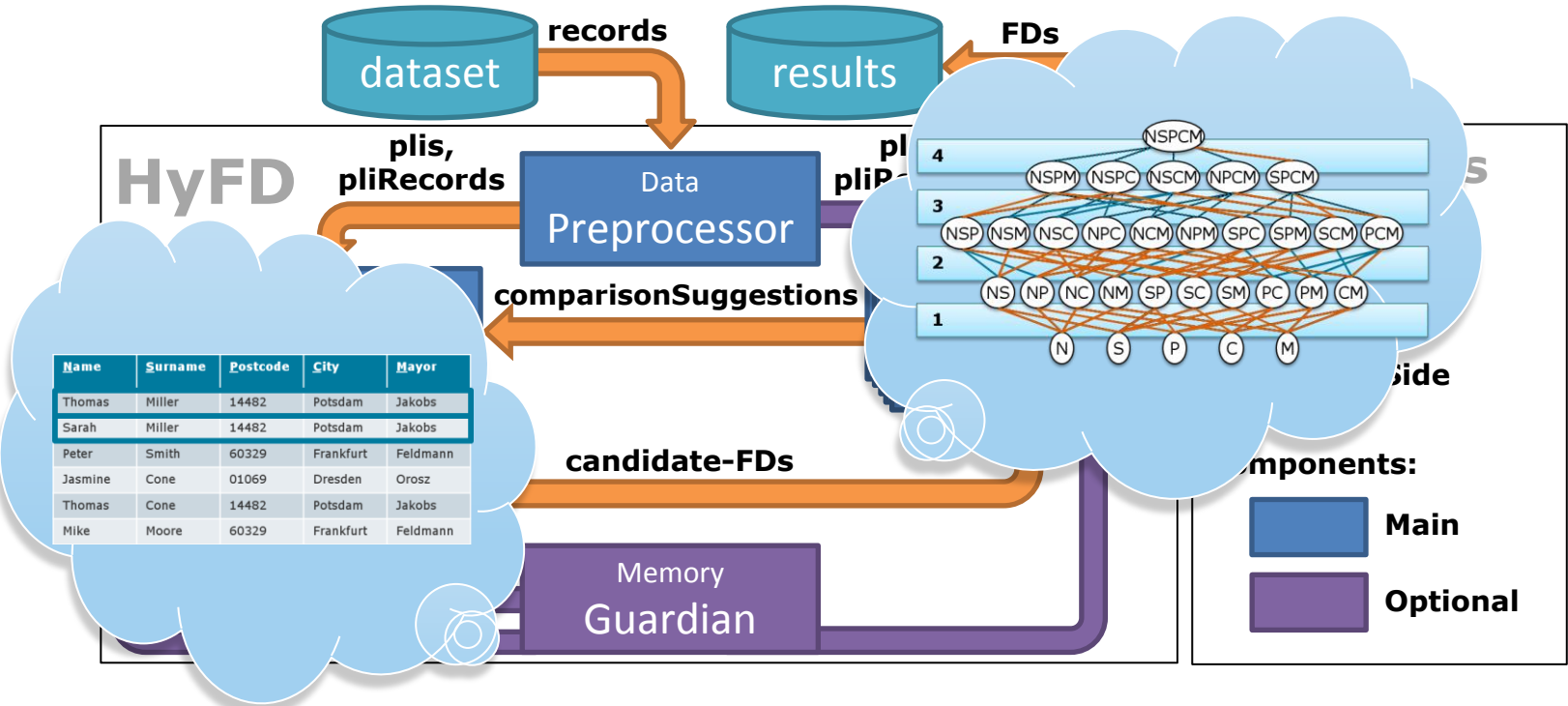
2018

Name	Surname	Postcode	City	Mayor
Thomas	Miller	14482	Potsdam	Jakobs
Sarah	Miller	14482	Potsdam	Jakobs
Peter	Smith	60329	Frankfurt	Feldmann
Jasmine	Cone	01069	Dresden	Orosz
Thomas	Cone	14482	Potsdam	Jakobs
Mike	Moore	60329	Frankfurt	Feldmann

- Surname, Postcode, City, Mayor $\not\rightarrow$ Name
- Name, Postcode, City, Mayor $\not\rightarrow$ Surname
- Surname $\not\rightarrow$ Name, Postcode, City, Mayor

Postcode \rightarrow City
 Postcode \rightarrow Mayor
 ...





Name	Surname	Postcode	City	Mayor
Thomas	Miller	14482	Potsdam	Jakobs
Sarah	Miller	14482	Potsdam	Jakobs
Peter	Smith	60329	Frankfurt	Feldmann
Jasmine	Cone	01069	Dresden	Orosz
Thomas	Cone	14482	Potsdam	Jakobs
Mike	Moore	60329	Frankfurt	Feldmann



2013 **2014** **2015** **2016** **2017** **2018**

Dataset	Cols [#]	Rows [#]	Size [KB]	FDs [#]	TANE	FUN	FD_MINE	DFD	DEP-MINER	FASTFDs	FDEP	HyFD
iris	5	150	5	4	1.1	0.1	0.2	0.2	0.2	0.2	0.1	0.1
balance-scale	5	625	7	1	1.2	0.1	0.2	0.3	0.3	0.3	0.2	0.1
chess	7	28,056	519	1	2.9	1.1	3.8	1.0	174.6	164.2	125.5	0.2
abalone	9	4,177	187	137	2.1	0.6	1.8	1.1	3.0	2.9	3.8	0.2
nursery	9	12,960	1,024	1	4.1	1.8	7.1	0.9	121.2	118.9	46.8	0.5
breast-cancer	11	699	20	46	2.3	0.6	2.2	0.8	1.1	1.1	0.5	0.2
bridges	13	108	6	142	2.2	0.6	4.2	0.9	0.5	0.6	0.2	0.1
echocardiogram	13	132	6	527	1.6	0.4	69.9	1.2	0.5	0.5	0.2	0.1
adult	14	48,842	3,528	78	67.4	111.6	531.5	5.9	6039.2	6033.8	860.2	1.1
letter	17	20,000	695	61	260.0	529.0	7204.8	6.0	1090.0	1015.5	291.3	3.4
ncvoter	19	1,000	151	758	4.3	4.0	ML	5.1	11.4	1.9	1.1	0.4
hepatitis	20	155	8	8,250	12.2	175.9	ML	326.7	5576.5	9.5	0.8	0.6
horse	27	368	25	128,727	457.0	TL	ML	TL	TL	385.8	7.2	7.1
fd-reduced-30	30	250,000	69,581	89,571	41.1	77.7	ML	TL	377.2	382.4	TL	513.0
plista	63	1,000	568	178,152	ML	ML	ML	TL	TL	TL	26.9	21.8
flight	109	1,000	575	982,631	ML	ML	ML	TL	TL	TL	216.5	53.4
uniprot	223	1,000	2,439	>2,437,556	ML	ML	ML	TL	TL	TL	ML	>5254.7

Results larger than 1,000 FDs are only counted

TL: time limit of 4 hours exceeded

ML: memory limit of 100 GB exceeded

2013

2014

2015

2016

2017

2018

Dataset	Cols	Rows	Size	FDs	HyFD
	[#]	[#]	[MB]	[#]	[s/m/h/d]
TPC-H.lineitem	16	6 m	1,051	4 k	4 m
PDB.POLY_SEQ	13	17 m	1,256	68	3 m
PDB.ATOM_SITE	31	27 m	5,042	10 k	64 m
SAP_R3.ZBC00DT	35	3 m	783	211	2 m
SAP_R3.ILOA	48	45 m	8,731	16 k	8 h
SAP_R3.CE4HI01	65	2 m	649	2 k	10 m
NCVoter.statewide	71	1 m	561	5 m	31 h
CD.cd	107	10 k	5	36 k	3 s

Features Business Explore Marketplace Pricing This repository Search Sign in Sign up

URI-Information-Systems / metanome-algorithms

Code Issues Pull requests Projects Insights

Source code for several Metanome data profiling algorithms

157 commits 2 branches 0 releases 12 contributors Apache-2.0

Search master New pull request Find file Clone or download

thorsten-papenbrock Update README.md Latest commit szazszs 10 days ago

- AIDFD Refactoring step 1 done 8 months ago
- BINDER Fix in BINDERs and SPIDERs file writing 3 months ago
- DVA Cardinality estimation algorithms 6 months ago
- DVAKMV Cardinality estimation algorithms 6 months ago
- DVAMS Cardinality estimation algorithms 6 months ago
- DVBJKST Cardinality estimation algorithms 6 months ago
- DVBloomFilter Cardinality estimation algorithms 6 months ago
- DVFM Cardinality estimation algorithms 6 months ago
- DVHyperLogLog Cardinality estimation algorithms 6 months ago
- DVHyperLogLogPlus Cardinality estimation algorithms 6 months ago
- DVLC Cardinality estimation algorithms 6 months ago
- DVLogLog Cardinality estimation algorithms 6 months ago
- DVMinCount Cardinality estimation algorithms 6 months ago
- DVPCSA Cardinality estimation algorithms 6 months ago
- DVSuperLogLog Cardinality estimation algorithms 6 months ago
- FAIDA Refactoring step 1 done 8 months ago
- HyFD Small optimization in HyFD. 5 months ago
- HyUCC Bugfix in HyFD and HyUCC 8 months ago
- MANY remove \$metanome version from parent and project version 8 months ago
- MvdDet Refactoring step 1 done 8 months ago
- ORDER remove \$metanome version from parent and project version 8 months ago
- SCDP fix in the SCDP algorithm. 7 months ago
- SPIDER Fix in BINDERs and SPIDERs file writing 3 months ago
- dao Bugfixes for BINDER's MIND version 3 years ago
- ducc Migrate Jens Ehrlich's DoCU/DCUCC algorithm to the current Metanome l... 5 months ago
- deprimer Refactoring step 1 done 8 months ago
- dfo Refactoring step 1 done 8 months ago
- ducc Refactoring step 1 done 8 months ago
- tsafds Refactoring step 1 done 8 months ago
- tstep Refactoring step 1 done 8 months ago
- tamine Refactoring step 1 done 8 months ago
- tun Refactoring step 1 done 8 months ago
- tane Refactoring step 1 done 8 months ago



Dependency	Algorithms (exact)	Algorithms (approximate)
Unique Column Combination (UCC)	2	0
Inclusion Dependency (IND)	5	2
Functional Dependency (FD)	8	2
Order Dependency (OD)	1	0
Matching Dependency (MD)	2	0
Multi-valued Dependency (MvD)	1	0
Denial Constraints (DC)	1	0
Statistics	1	13
	21	17

2013

2014

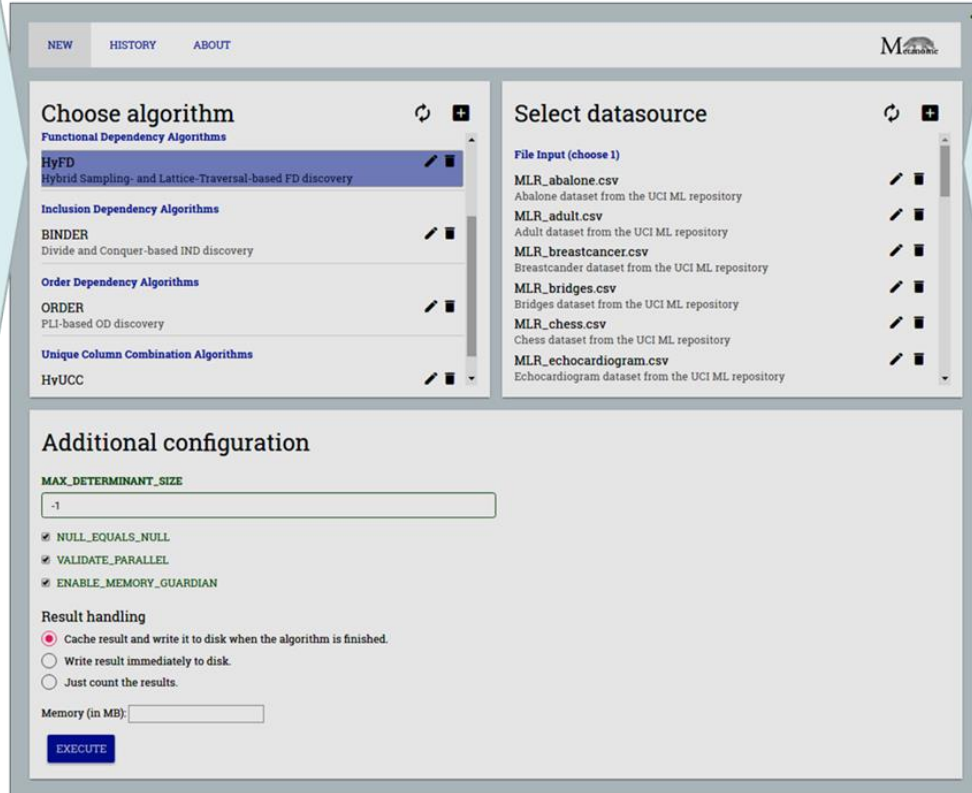
2015

2016

2017

2018

Algorithms

 $A \rightarrow B$
 $A \subset B$
 $A \gg B$
 $A \subset B \subset C \subset D$


The screenshot shows the Metanome web interface with the following sections:

- Choose algorithm:**
 - Functional Dependency Algorithms:** HyFD (Hybrid Sampling- and Lattice-Traversal-based FD discovery)
 - Inclusion Dependency Algorithms:** BINDER (Divide and Conquer-based IND discovery)
 - Order Dependency Algorithms:** ORDER (PLI-based OD discovery)
 - Unique Column Combination Algorithms:** HyUCC
- Select datasource:**
 - File Input (choose 1)
 - MLR_abalone.csv (Abalone dataset from the UCI ML repository)
 - MLR_adult.csv (Adult dataset from the UCI ML repository)
 - MLR_breastcancer.csv (Breastcancer dataset from the UCI ML repository)
 - MLR_bridges.csv (Bridges dataset from the UCI ML repository)
 - MLR_chess.csv (Chess dataset from the UCI ML repository)
 - MLR_echocardiogram.csv (Echocardiogram dataset from the UCI ML repository)
- Additional configuration:**
 - MAX_DETERMINANT_SIZE: -1
 - NULL_EQUALS_NULL
 - VALIDATE_PARALLEL
 - ENABLE_MEMORY_GUARDIAN
 - Result handling:**
 - Cache result and write it to disk when the algorithm is finished.
 - Write result immediately to disk.
 - Just count the results.
 - Memory (in MB):
 - EXECUTE** button

Datasets



2013

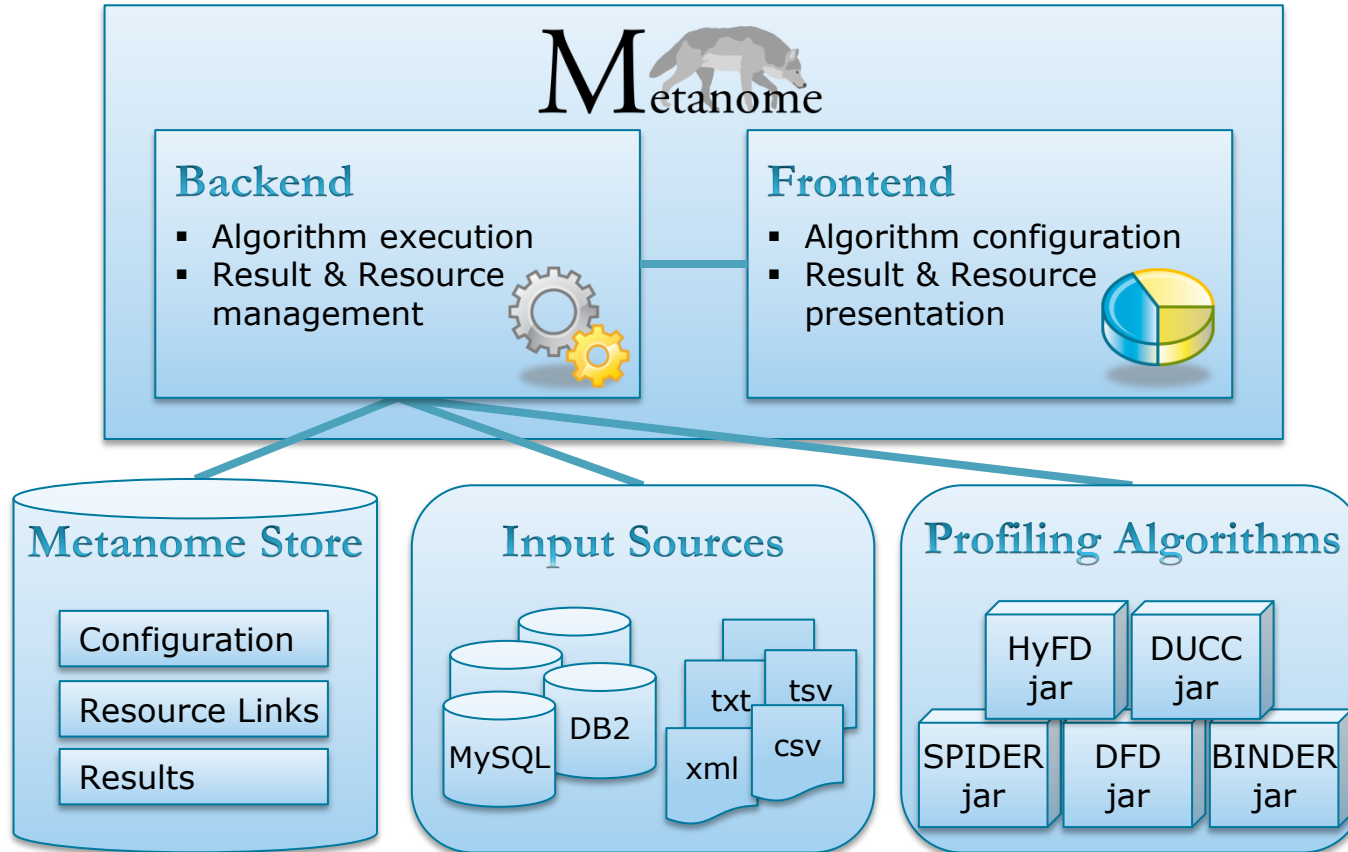
2014

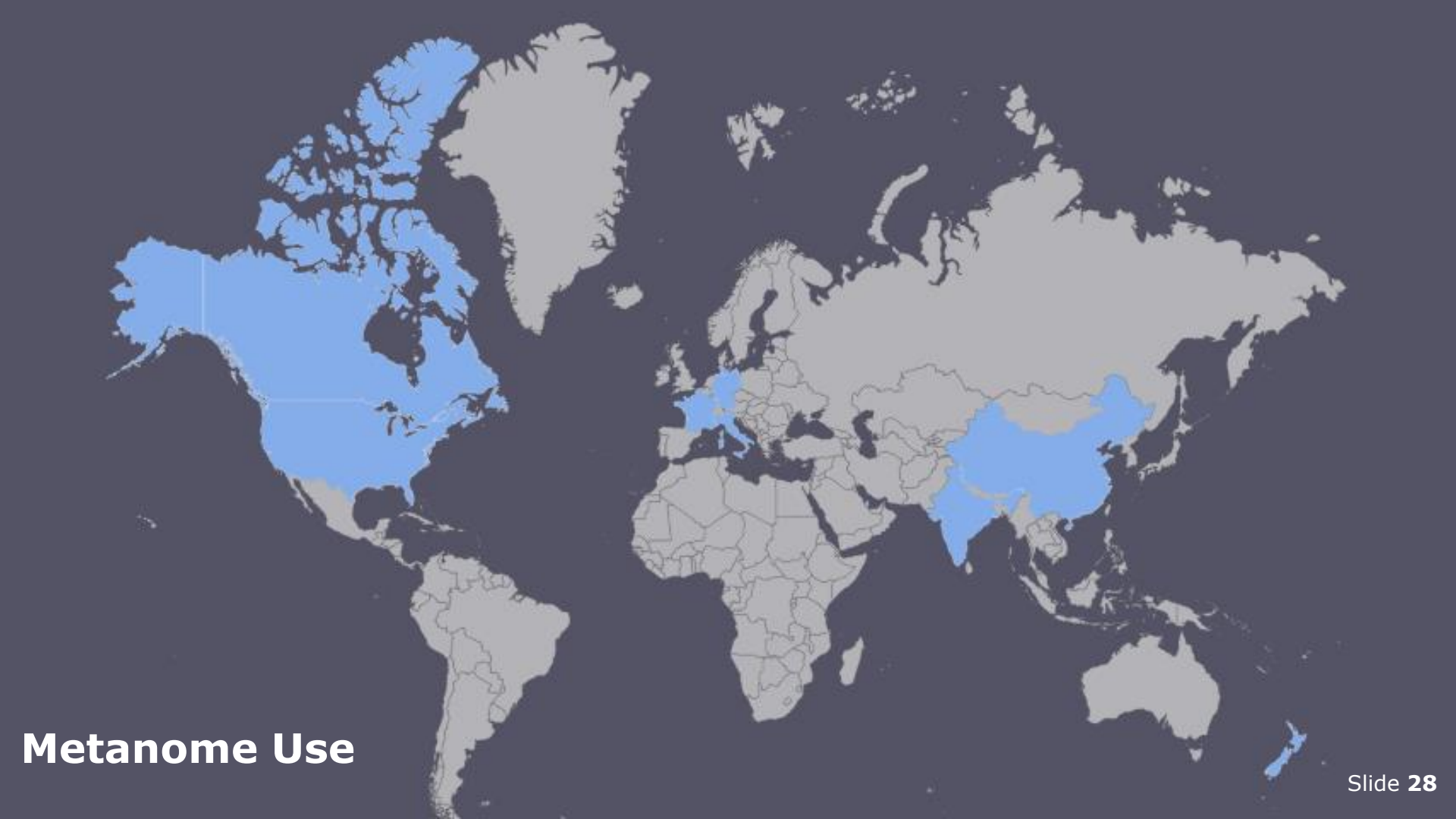
2015

2016

2017

2018



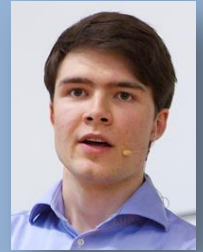


Metanome Use

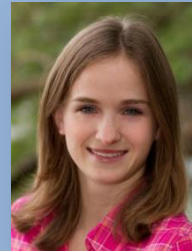
Metanome Algorithm Research

- Anja Jentzsch (RDF)
- Arvid Heise (UCC)
- Fabian Tschirschnitz (IND)
- Felix Naumann (Research lead)
- Hazar Harmouch (Single Column Profiling)
- Jens Ehrlich (Conditional UCC)
- Jorge-Arnulfo (UCC, IND)
- Maximilian Grundke (Conditional FD)
- Moritz Finke (Approximate FD/IND)
- Philipp Langer (OD)
- Philipp Schirmer (MVD)
- Sebastian Kruse (IND, partial FD; Metadata Store)
- **Thorsten Papenbrock** (IND, UCC, FD, ...; Metanome)
- Tim Draeger (MVD)
- Tobias Bleifuß (DC)
- Ziawasch Abedjan (UCC)

Tool Development



Jakob Zwiener
(Backend & Frontend)



Claudia Exeler
(Frontend)



Tanja Bergmann
(Backend & Frontend)



Moritz Finke
(Backend)



Carl Ambroselli
(Frontend)



Maxi Fischer
(Backend & Frontend)



Vincent Schwarzer
(Backend)

2019	DynFD: Functional Dependency Discovery in Dynamic Datasets P. Schirmer, T. Papenbrock, S. Kruse, D. Hempfing, T. Mayer, D. Neuschäfer-Rube, F. Naumann An Actor Database System for Akka S. Schmidl, F. Schneider, T. Papenbrock	(EDBT) (BTW)	
2018	Data Profiling – Synthesis Lectures on Data Management Z. Abedjan, L. Golab, F. Naumann, T. Papenbrock	(Morgan & Claypool)	
2017	Detecting Inclusion Dependencies on Very Many Tables F. Tschirschnitz, T. Papenbrock, F. Naumann Data-driven Schema Normalization T. Papenbrock, F. Naumann A Hybrid Approach for Efficient Unique Column Combination Discovery T. Papenbrock, F. Naumann Fast Approximate Discovery of Inclusion Dependencies S. Kruse, T. Papenbrock, C. Dullweber, M. Finke, M. Hegner, M. Zabel, C. Zöllner, F. Naumann	(TODS) (EDBT) (BTW) (BTW)	
2016	A Hybrid Approach to Functional Dependency Discovery T. Papenbrock, F. Naumann Data Anamnesis: Admitting Raw Data into an Organization S. Kruse, T. Papenbrock, H. Harmouch, F. Naumann Holistic Data Profiling: Simultaneous Discovery of Various Metadata J. Ehrlich, M. Roick, L. Schulze, J. Zwiener, T. Papenbrock, F. Naumann RDFind: Scalable Conditional Inclusion Dependency Discovery in RDF Datasets S. Kruse, A. Jentzsch, T. Papenbrock, Z. Kaoudi, J. Quiané-Ruiz, F. Naumann Approximate Discovery of Functional Dependencies for Large Datasets T. Bleifuß, S. Bülow, J. Frohnhofen, J. Risch, G. Wiese, S. Kruse, T. Papenbrock, F. Naumann	(SIGMOD) (IEEE Data Engineering Bulletin) (EDBT) (SIGMOD)	
2015	Functional Dependency Discovery: An Experimental Evaluation of Seven Algorithms T. Papenbrock, J. Ehrlich, J. Marten, T. Neubert, J. Rudolph, M. Schönberg, J. Zwiener, F. Naumann Data Profiling with Metanome T. Papenbrock, T. Bergmann, M. Finke, J. Zwiener, F. Naumann Divide & Conquer-based Inclusion Dependency Discovery T. Papenbrock, S. Kruse, J. Quiané-Ruiz, F. Naumann Scaling Out the Discovery of Inclusion Dependencies S. Kruse, T. Papenbrock, F. Naumann Progressive Duplicate Detection T. Papenbrock, A. Heise, F. Naumann	(CIKM) (VLDB) (VLDB) (VLDB) (BTW) (TKDE)	
2013	Ein Datenbankkurs mit 6000 Teilnehmern F. Naumann, M. Jenders, T. Papenbrock Duplicate Detection on GPUs B. Forchhammer, T. Papenbrock, T. Stening, S. Viehmeier, U. Draisbach, F. Naumann	(Informatik-Spektrum) (BTW)	
2011	BlackSwan: Augmenting Statistics with Event Data J. Lorey, F. Naumann, B. Forchhammer, A. Mascher, P. Retzlaff, A. Zamani Farahani, S. Discher, C. Fähnrich, S. Lemme, T. Papenbrock, R. C. Peschel, S. Richter, T. Stening, S. Viehmeier	(CIKM)	



Discovery



Application



Distribution



Robustness