

The background of the slide is a photograph of a server rack. The server units are black with numerous yellow handles and ports. The perspective is from an angle, looking down into the rack. A semi-transparent orange banner is overlaid at the bottom of the image.

# Distributed Computing Data Profiling at Scale

Dr. Thorsten Papenbrock  
Information Systems Group, HPI



# Dr. Thorsten Papenbrock

## Education

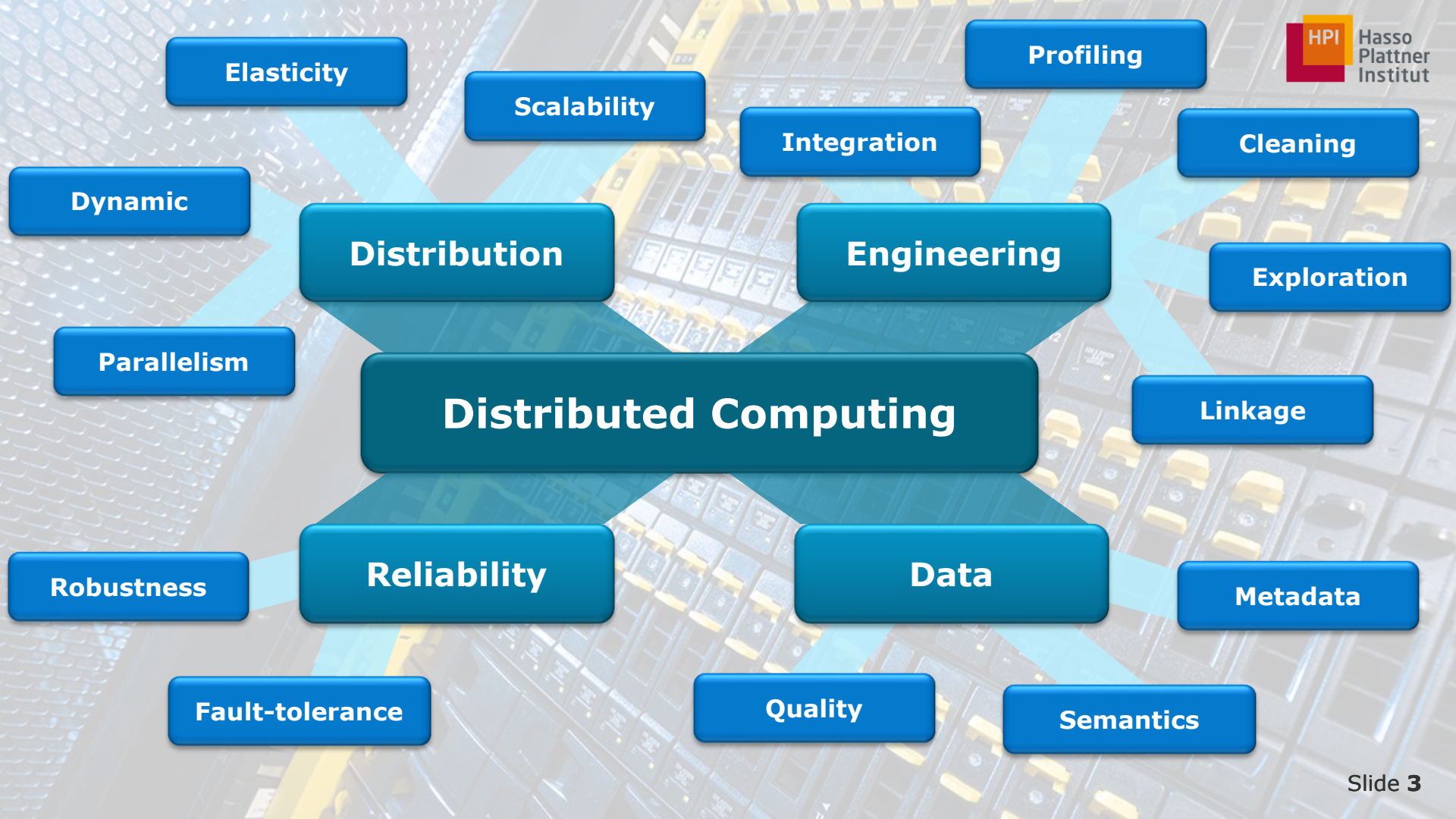
2017	PhD	HPI	Potsdam
2013	Master	HPI	Potsdam
2010	Bachelor	HPI	Potsdam
2007	Abitur	KvG	Mettingen

## Industry

2014	Intern	QCRI	Qatar
2011/12	Intern	SAP	Belfast
2010/11	Project	BBF	Munich
2009	Intern	SAP	Walldorf

## Research Interests

- Distributed Computing
- Data Profiling and Cleaning
- Information Systems and Databases



**Elasticity**

**Scalability**

**Profiling**

**Integration**

**Cleaning**

**Dynamic**

**Distribution**

**Engineering**

**Exploration**

**Parallelism**

**Distributed Computing**

**Linkage**

**Robustness**

**Reliability**

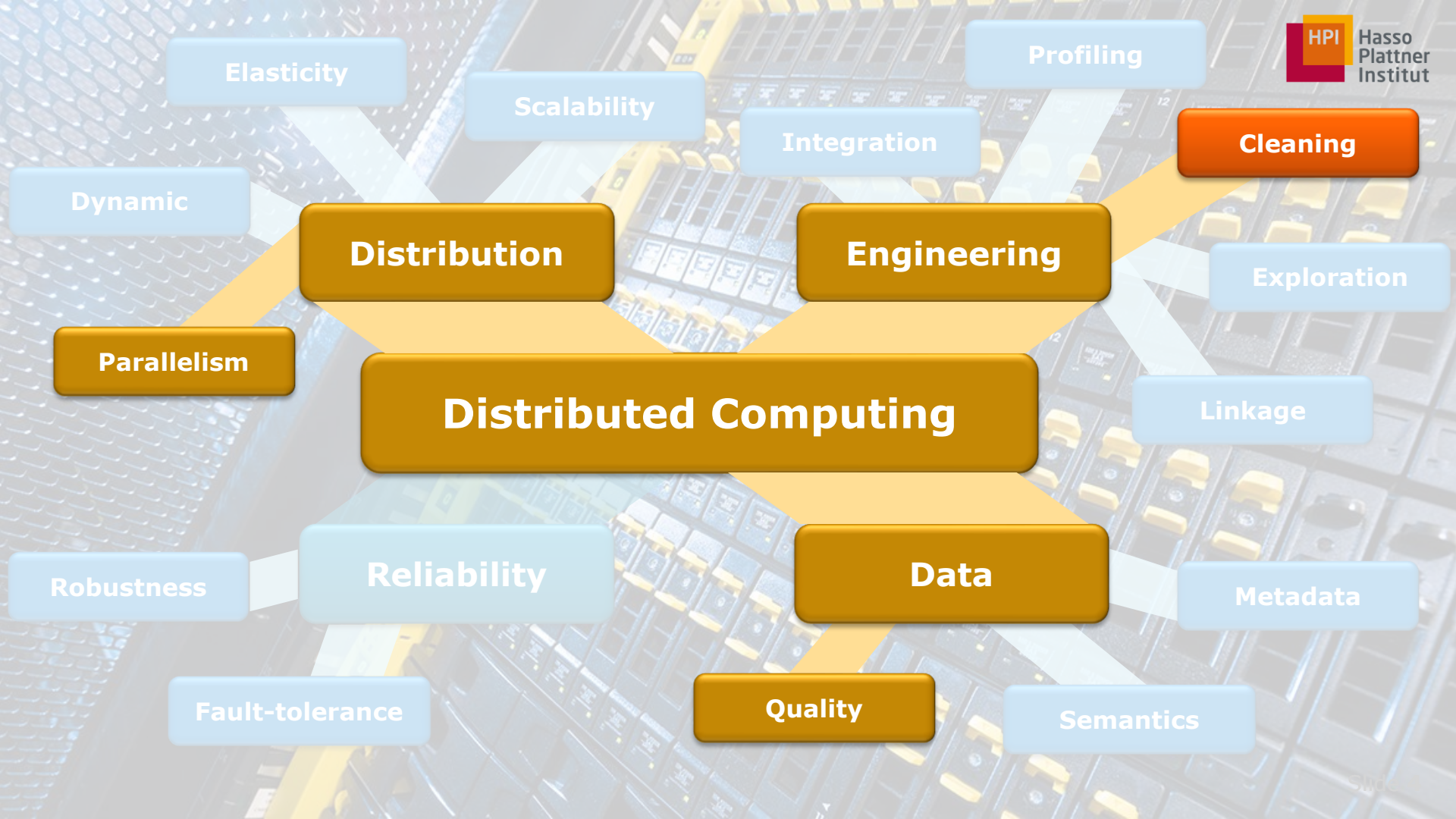
**Data**

**Metadata**

**Fault-tolerance**

**Quality**

**Semantics**



Elasticity

Scalability

Profiling

Cleaning

Dynamic

Distribution

Engineering

Exploration

Parallelism

Distributed Computing

Linkage

Robustness

Reliability

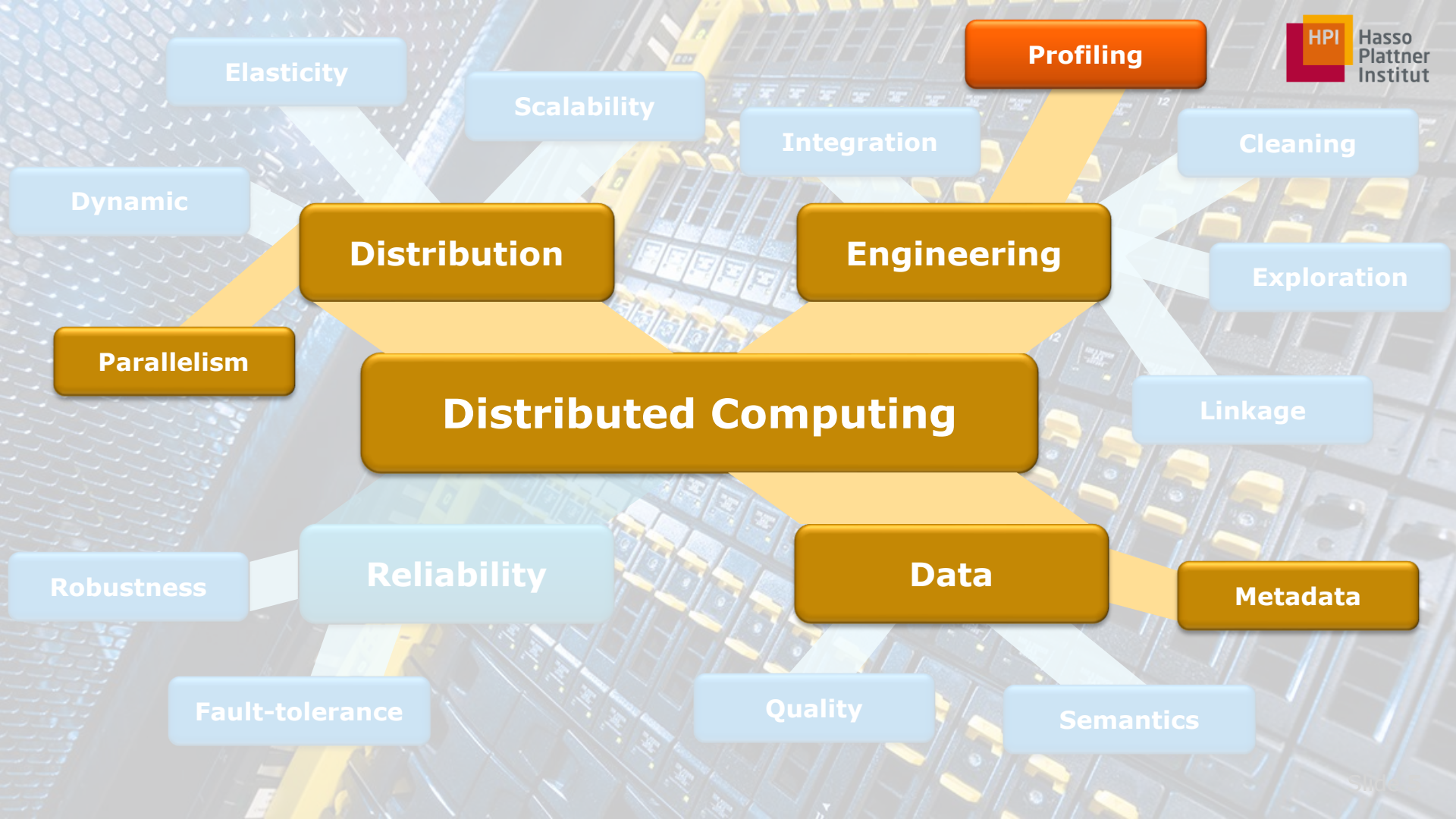
Data

Metadata

Fault-tolerance

Quality

Semantics



**Distributed Computing**

**Distribution**

**Engineering**

**Data**

**Reliability**

Elasticity

Scalability

Integration

Profiling

Cleaning

Exploration

Linkage

Metadata

Semantics

Quality

Fault-tolerance

Robustness

Dynamic

Parallelism

## Art. 15 DSGVO Auskunftsrecht der betroffenen Person

### Knowledge Discovery

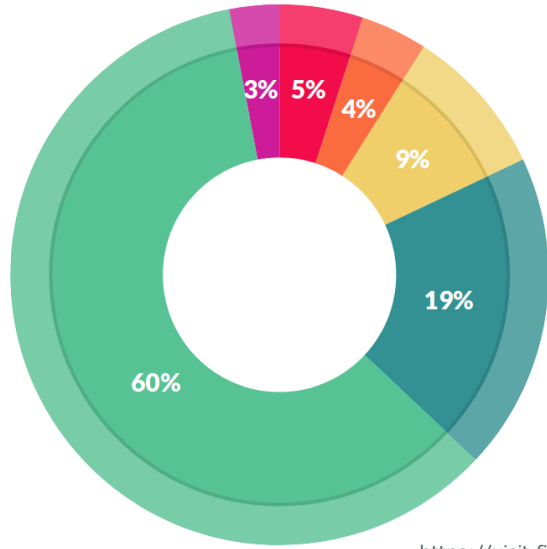
What data do you have?

- (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) <sup>1</sup> Der Verantwortliche stellt eine Kopie der personenbezogenen Daten, die Gegenstand der Verarbeitung sind, zur Verfügung. <sup>2</sup> Für alle weiteren Kopien, die die betroffene Person beantragt, kann der Verantwortliche ein angemessenes Entgelt auf der Grundlage der Verwaltungskosten verlangen. <sup>3</sup> 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!

Knowledge Discovery

Data Analytics

[https://visit.figure-eight.com/rs/416-ZBE-142/images/CrowdFlower\\_DataScienceReport\\_2016.pdf](https://visit.figure-eight.com/rs/416-ZBE-142/images/CrowdFlower_DataScienceReport_2016.pdf)

## 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
- ...

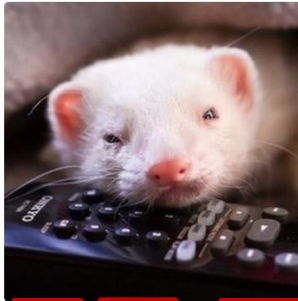


# Deep Visual-Semantic Alignments for Generating Image Descriptions

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

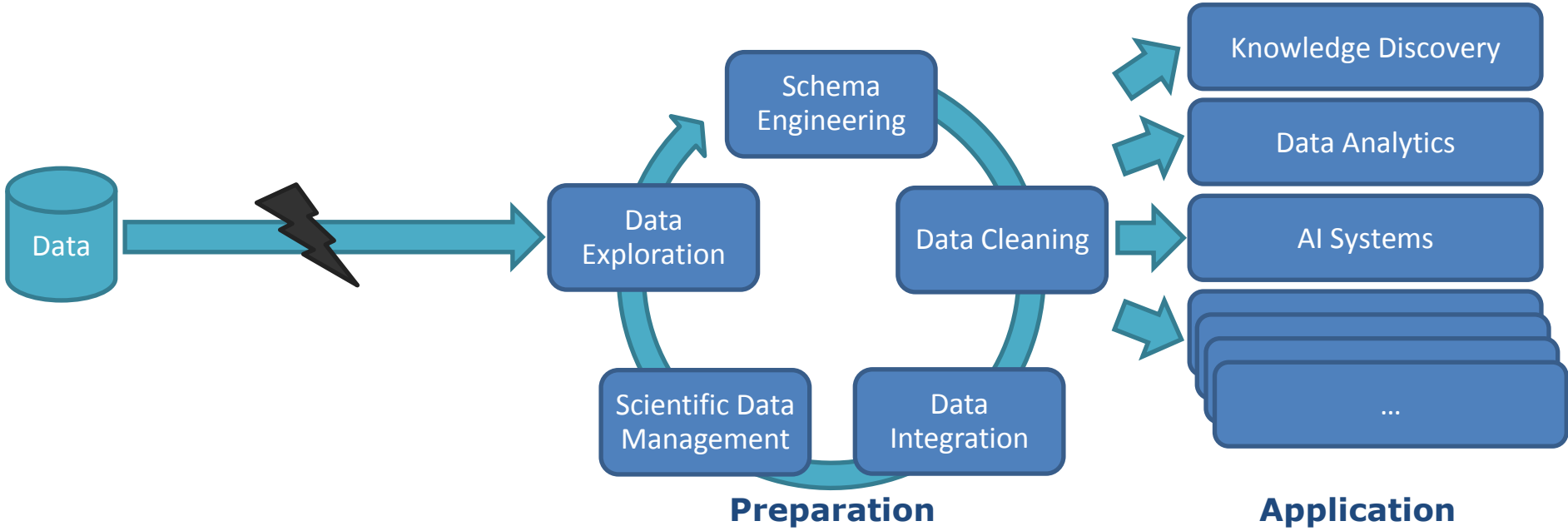
AI Systems

...

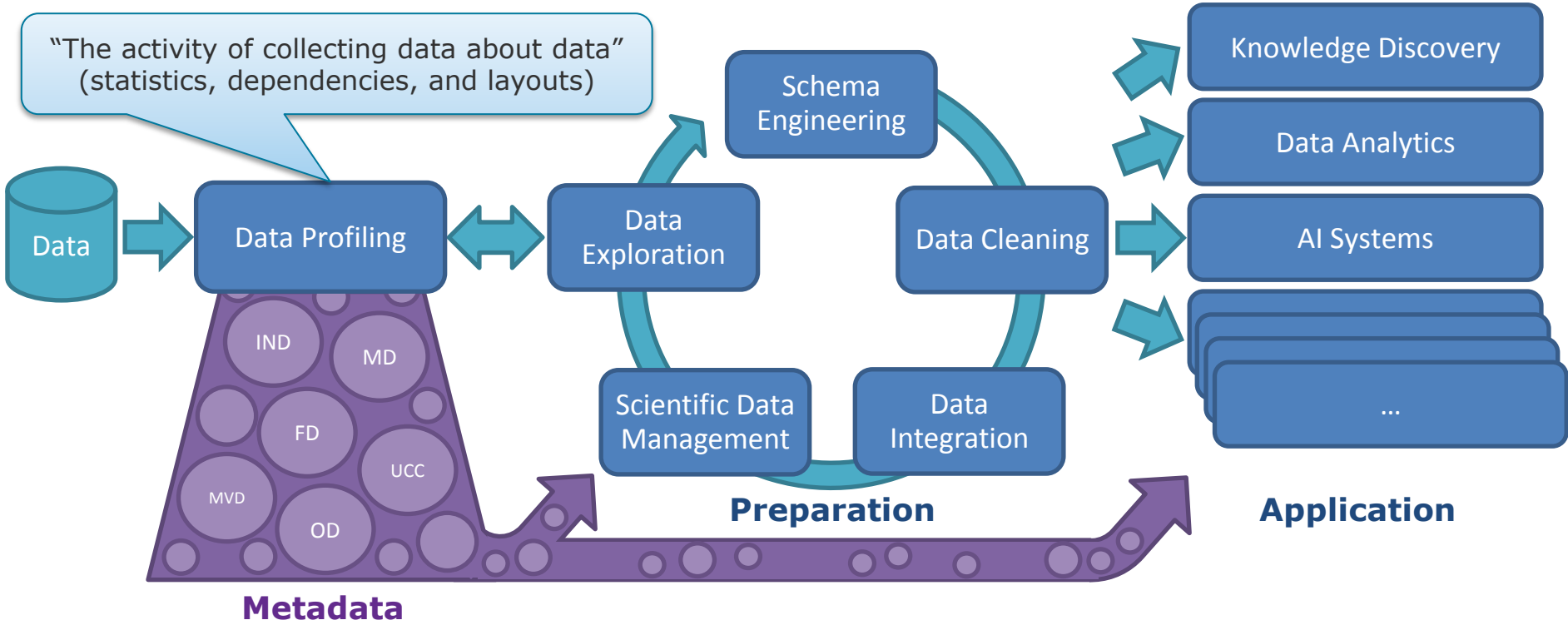
...

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	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

# Data Profiling

format

density

#null = 3  
%null = 30

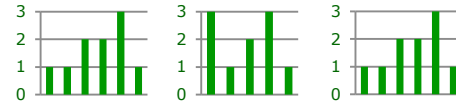
ranges

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

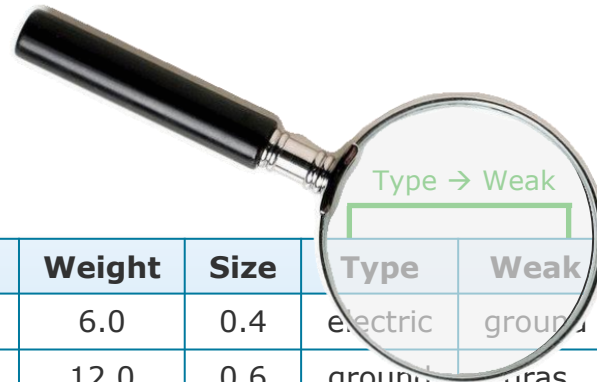
Weight  $\downarrow$  Size

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	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



A B C D E F G H I J K L M N O P Q R S T U V W X Y Z AA AB AC AD AE AF AG AH

county county des-voter rep-status-voter status reason voter status reason d-a-b-n-last name first name middle name name-res city desc res city desc zip q-mail addr1 mail-m-a-p-mail city mail-m-a-p-mail city mail-m-a-p-mail city num-birth-party cd gen-birth-party cd gen-birth-party cd

Table with columns for county, voter info, status, address, city, zip, and party. Includes rows for various individuals like EVELYN LARSEN, CHRISTINA CASTAGNA, CLAUDIA HAYDEN, etc.

# Data Profiling

~8 million records

94 attributes

085  
03S  
124  
124  
13  
35  
06S  
06N  
10N  
03C  
03S  
03W  
4  
7  
09S  
128  
128  
1210  
128  
03N  
127  
03C  
06W  
13  
4  
03N  
03N  
124  
06C  
09S  
64

# Agenda



Advances in Data Profiling



Challenges in Distributed Computing

2013

2014

2015

2016

2017

2018

Dataset	Columns [#]	Rows [#]	Size [KB]	FDs [#]	TANE [7]	FUN [14]	FD_MINE [21]	DFD [1]	DEP-MINER [12]	FASTFDs [20]	FDEP [6]
iris	5	150	5	4	1.1	<b>0.1</b>	0.2	0.2	0.2	0.2	<b>0.1</b>
balance-scale	5	625	7	1	1.2	<b>0.1</b>	0.2	0.3	0.3	0.3	0.2
chess	7	28,056	519	1	2.9	1.1	3.8	<b>1.0</b>	174.6	164.2	125.5
abalone	9	4,177	187	137	2.1	<b>0.6</b>	1.8	1.1	3.0	2.9	3.8
nursery	9	12,960	1,024	1	4.1	1.8	7.1	<b>0.9</b>	121.2	118.9	46.8
breast-cancer	11	699	20	46	2.3	0.6	2.2	0.8	1.1	1.1	<b>0.5</b>
bridges	13	108	6	142	2.2	0.6	4.2	0.9	0.5	0.6	<b>0.2</b>
echocardiogram	13	132	6	538	1.6	0.4	69.9	1.2	0.5	0.5	<b>0.2</b>
adult	14	48,842	3,528	78	67.4	111.6	531.5	<b>5.9</b>	6039.2	6033.8	860.2
letter	17	20,000	695	61	260.0	529.0	7204.8	<b>6.0</b>	1090.0	1015.5	291.3
ncvoter	19	1,000	151	758	4.3	4.0	ML	5.1	11.4	1.9	<b>1.1</b>
hepatitis	20	155	8	8,250	12.2	175.9	ML	326.7	5576.5	9.5	<b>0.8</b>
horse	27	368	25	128,726	457.0	TL	ML	TL	TL	385.8	<b>7.2</b>
fd-reduced-30	30	250,000	69,581	89,571	<b>41.1</b>	77.7	ML	TL	377.2	382.4	TL
plista	63	1,000	568	178,152	ML	ML	ML	TL	TL	TL	<b>26.9</b>
flight	109	1,000	575	982,631	ML	ML	ML	TL	TL	TL	<b>216.5</b>
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

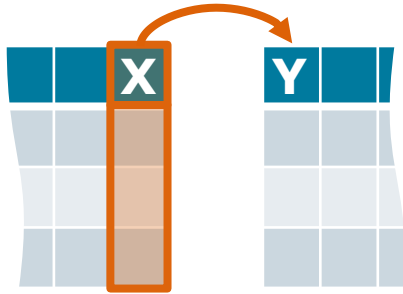
2015

2016

2017

2018

## Inclusion Dependencies



Foreign key relationships

**Definition:** Given two relational instances  $r_i$  and  $r_j$  for the schemata  $R_i$ , and  $R_j$ , respectively. The **inclusion dependency**  $R_i[X] \subseteq R_j[Y]$  (short  $X \subseteq Y$ ) with  $X \subseteq R_i$ ,  $Y \subseteq R_j$  and  $|X| = |Y|$  is valid, iff  $\forall t_i[X] \in r_i, \exists t_j[Y] \in r_j : t_i[X] = t_j[Y]$ .

“All values in X are also contained in Y”

2013

2014

2015

2016

2017

2018

Name	Type	Equatorial diameter	Mass	Orbital radius	Orbital period	Rotation period	Confirmed moons	Rings	Atmosphere
Mercury	Terrestrial	0.382	0.06	0.47	0.24	58.64	0	no	minimal
Venus	Terrestrial	0.949	0.82	0.72	0.62	-243.02	0	no	CO <sub>2</sub> , N <sub>2</sub>
Earth	Terrestrial	1.000	1.00	1.00	1.00	1.00	1	no	N <sub>2</sub> , O <sub>2</sub> , Ar
Mars	Terrestrial	0.532	0.11	1.52	1.88	1.03	2	no	CO <sub>2</sub> , N <sub>2</sub> , Ar
Jupiter	Giant	11.209	317.8	5.20	11.86	0.41	67	yes	H <sub>2</sub> , He
Saturn	Giant	9.449	95.2	9.54	29.46	0.43	62	yes	H <sub>2</sub> , He
Uranus	Giant	4.007	14.6	19.22	84.01	-0.72	27	yes	H <sub>2</sub> , He
Neptune	Giant	3.883	17.2	30.06	164.8	0.67	14	yes	H <sub>2</sub> , He

**Complexity:**  $O(n^2-n)$   
for  $n$  attributes

**Example:**  
10 attributes  $\sim$  90 checks  
1,000 attributes  $\sim$  999,000 checks

- Name  $\subseteq$  Type ?
- Name  $\subseteq$  Equatorial\_diameter ?
- Name  $\subseteq$  Mass ?
- Name  $\subseteq$  Orbital\_radius ?
- Name  $\subseteq$  Orbital\_period ?
- Name  $\subseteq$  Rotation\_period ?
- Name  $\subseteq$  Confirmed\_moons ?
- Name  $\subseteq$  Rings ?
- Name  $\subseteq$  Atmosphere ?
- Type  $\subseteq$  Name ?
- Type  $\subseteq$  Equatorial\_diameter ?
- Type  $\subseteq$  Mass ?
- Type  $\subseteq$  Orbital\_radius ?
- Type  $\subseteq$  Orbital\_period ?
- Type  $\subseteq$  Rotation\_period ?
- Type  $\subseteq$  Confirmed\_moons ?
- Type  $\subseteq$  Rings ?
- Type  $\subseteq$  Atmosphere ?
- Mass  $\subseteq$  Name ?
- Mass  $\subseteq$  Type ?
- Mass  $\subseteq$  Equatorial\_diameter ?
- ...

2013

2014

2015

2016

2017

2018

Name	Type	Equatorial diameter	Mass	Orbital radius	Orbital period	Rotation period	Confirmed moons	Rings	Atmosphere
Mercury	Terrestrial	0.382	0.06	0.47	0.24	58.64	0	no	minimal
Venus	Terrestrial	0.949	0.82	0.72	0.62	-243.02	0	no	CO <sub>2</sub> , N <sub>2</sub>
Earth	Terrestrial	1.000	1.00	1.00	1.00	1.00	1	no	N <sub>2</sub> , O <sub>2</sub> , Ar
Mars	Terrestrial	0.532	0.11	1.52	1.88	1.03	2	no	CO <sub>2</sub> , N <sub>2</sub> , Ar
Jupiter	Giant	11.209	317.8	5.20	11.86	0.41	67	yes	H <sub>2</sub> , He
Saturn	Giant	9.449	95.2	9.54	29.46	0.43	62	yes	H <sub>2</sub> , He
Uranus	Giant	4.007	14.6	19.22	84.01	-0.72	27	yes	H <sub>2</sub> , He
Neptune	Giant	3.883	17.2	30.06	164.8	0.67	14	yes	H <sub>2</sub> , He

Planet	Rotation Period	Revolution Period
Mercury	58.6 days	87.97 days
Venus	243 days	224.7 days
Earth	0.99 days	365.26 days
Mars	1.03 days	1.88 years
Jupiter	0.41 days	11.86 years
Saturn	0.45 days	29.46 years
Uranus	0.72 days	84.01 years
Neptune	0.67 days	164.79 years
Pluto	6.39 days	248.59 years

Symbol	Unicode	Glyph
Sun	U+2609	☉
Moon	U+263D	☾
Moon	U+263E	☾
Mercury	U+263F	☿
Venus	U+2640	♀
Earth	U+1F728	🌍
Mars	U+2642	♂
Jupiter	U+2643	♃
Saturn	U+2644	♄
Uranus	U+2645	♅
Uranus	U+26E2	♁
Neptune	U+2646	♆
Eris	≈ U+2641	♁
Eris	≈ U+29EC	♁
Pluto	U+2647	♇
Pluto	not present	--
Aries	U+2648	♈
Taurus	U+2649	♉
Gemini	U+264A	♊
Cancer	U+264B	♋
Leo	U+264C	♌
Virgo	U+264D	♍
Libra	U+264E	♎
Scorpio	U+264F	♏
Sagittarius	U+2650	♐
Capricorn	U+2651	♑
Capricorn	U+2651	♑
Aquarius	U+2652	♒
Pisces	U+2653	♓
Conjunction	U+260C	♆
...	...	...

Planet	Synodic period	Synodic period (mean)	Days in retrograde
Mercury	116	3.8	~21
Venus	584	19.2	41
Mars	780	25.6	72
Jupiter	399	13.1	121
Saturn	378	12.4	138
Uranus	370	12.15	151
Neptune	367	12.07	158

Planet	Mean distance	Relative mean distance
Mercury	57.91	1
Venus	108.21	1.86859
Earth	149.6	1.3825
Mars	227.92	1.52353
Ceres	413.79	1.81552
Jupiter	778.57	1.88154
Saturn	1,433.53	1.84123
Uranus	2,872.46	2.00377
Neptune	4,495.06	1.56488
Pluto	5,869.66	1.3058

Sign	House	Domicile	Detriment	Exaltation	Fall	Planetary Joy
Aries	1st House	Mars	Venus	Sun	Saturn	Mercury
Taurus	2nd House	Venus	Pluto	Moon	Uranus	Jupiter
Gemini	3rd House	Mercury	Jupiter	N/A	N/A	Saturn
Cancer	4th House	Moon	Saturn	Jupiter	Mars	Venus
Leo	5th House	Sun	Uranus	Neptune	Mercury	Mars
Virgo	6th House	Mercury	Neptune	Pluto, Mercury	Venus	Saturn
Libra	7th House	Venus	Mars	Saturn	Sun	Moon
Scorpio	8th House	Pluto	Venus	Uranus	Moon	Saturn
Sagittarius	9th House	Jupiter	Mercury	N/A	N/A	Sun
Capricorn	10th House	Saturn	Moon	Mars	Jupiter	Mercury
Aquarius	11th House	Uranus	Sun	Mercury	Neptune	Venus
Pisces	12th House	Neptune	Mercury	Venus	Pluto, Mercury	Moon

Planet	Calculated (in AU)	Observed (in AU)	Perfect octaves	Actual distance
Mercury	0.4	0.387	0	0
Venus	0.7	0.723	1	1.1
Earth	1	1	2	2
Mars	1.6	1.524	4	3.7
Asteroid belt	2.8	2.767	8	7.8
Jupiter	5.2	5.203	16	15.7
Saturn	10	9.539	32	29.9
Uranus	19.6	19.191	64	61.4
Neptune	38.8	30.061	96	-96.8
Pluto	77.2	39.529	128	127.7

2013

2014

2015

2016

2017

2018

 attributes     values  
 dataflow     ignored

**Divide****Conquer**

Rel. 1

Rel. 2

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

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

A	B	C	D	E	F	G
	X			X		
	X		X	X		X
	X		X	X		X
	X	X	X	X		X
	X	X	X	X		X

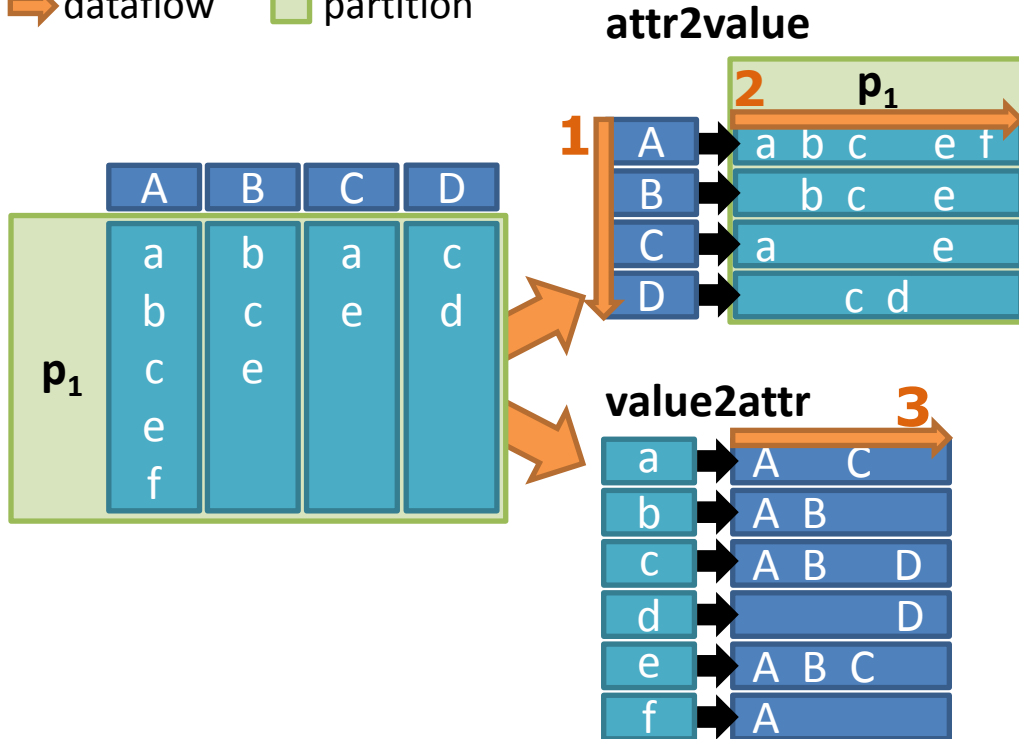
**validation?** $F \subseteq A$ 

### Dynamic Memory Handling:

Spill largest buckets to disk if memory is exhausted.

### Lazy Partition Refinement:

Split a partition if it does not fit into main memory.

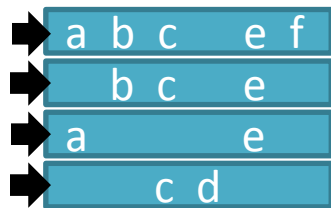


### Validation algorithm

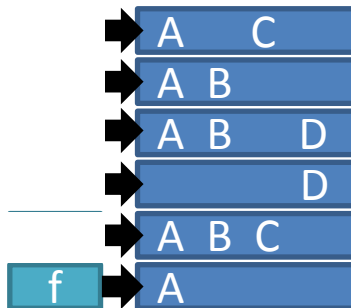
1. Iterate attributes
2. Iterate values
3. If value2attr entry exists
  - Intersect candidates with this list
  - Remove value2attr entry
  - If attribute removed from all candidates
  - Remove entry from attr2value



## attr2value



## value2attr

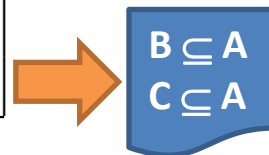


Never tested! →

f

1. Iterate attributes
2. Iterate values
3. If value2attr entry exists
  - Intersect candidates with this list
  - Remove value2attr entry
  - If attribute removed from all candidates
  - Remove entry from attr2value

	A	B	C	D
look up	B,C,D	A,C,D	A,B,D	A,B,C



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

Table 6: uIND performance on real-world datasets (minutes)

Datasets	Bell & Brockhausen	DeMarchi	Spider	Spider BF	S-indd	Binder	Size	Attributes	uInds	nInds
SCOP	0.14	<b>0.04</b>	0.08	0.07	0.09	0.07	16 MB	22	39	36
CATH	0.11	<b>0.02</b>	0.05	0.05	0.05	0.04	16 MB	25	51	81
CENSUS	1.05	0.97	0.15	0.18	0.17	<b>0.14</b>	112 MB	42	39	89
WIKIPEDIA	TL	<b>1.02</b>	1.47	1.45	1.56	1.22	540 MB	14	2	0
BIOSQL	4.98	<b>0.76</b>	1.30	1.43	1.48	1.41	560 MB	77	348	507
WIKIRANK	2.90	<b>0.73</b>	1.53	1.19	1.44	1.23	697 MB	29	99	103
LOD	0.34	<b>0.25</b>	0.45	0.43	0.41	0.30	830 MB	41	258	unknown
ENSEMBL	23.52	6.87	3.04	3.50	3.70	<b>2.39</b>	836 MB	130	364	100
TESMA	TL	4.53	3.30	<b>2.85</b>	4.75	4.27	1 GB	114	2	0
TPC-H 1	17.79	<b>1.96</b>	3.88	3.10	3.58	2.96	1 GB	61	96	8
TPC-H 10	TL	ML	44.43	35.06	36.54	<b>28.21</b>	10 GB	61	97	11
MUSICBRAINZ	TL	136.03	61.42	130.15	106.26	<b>45.69</b>	27 GB	1054	49829	unknown

TL: time limit of 4h exceeded

ML: memory limit of 83G exceeded

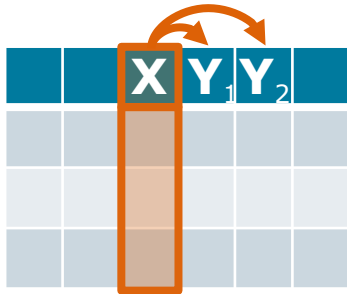
Table 7: nIND performance on real-world datasets in minutes

Datasets	Binder	Find2	Mind	Mind2	ZigZag
SCOP	<b>0.30</b>	0.37	0.36	1.85	0.50
CATH	4.18	3.24	<b>3.15</b>	29.16	3.23
CENSUS	<b>0.68</b>	2.28	2.01	N.A.	3.37
WIKIPEDIA	<b>1.40</b>	1.50	1.48	1.47	1.54
BIOSQL	<b>3.65</b>	5.51	5.36	TL	4.88
WIKIRANK	8.10	2.99	8.89	TL	<b>2.96</b>
ENSEMBL	<b>6.31</b>	8.33	7.13	TL	206.94
ghaIND	177.65	13.28	27.94	13.45	<b>13.13</b>
TESMA	4.35	3.39	<b>3.38</b>	8.79	3.39
TPC-H 1	9.97	7.49	<b>7.01</b>	TL	12.37
TPC-H 10	<b>121.09</b>	TL	TL	TL	TL

TL: time limit of 4h exceeded ML: memory limit of 83G exceeded

[Inclusion Dependency Discovery:  
 An Experimental Evaluation of Ten Algorithms,  
 F. Dürsch, A. Stebner, F. Windheuser, M. Fischer,  
 T. Friedrich, N. Strelow, T. Bleifuß, H. Harmouch,  
 L. Jiang, T. Papenbrock, F. Naumann,  
 submitted to VLDB]

## 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  $\forall t_i, t_j \in r : t_i[X] = t_j[X] \Rightarrow t_i[A] = t_j[A]$ .

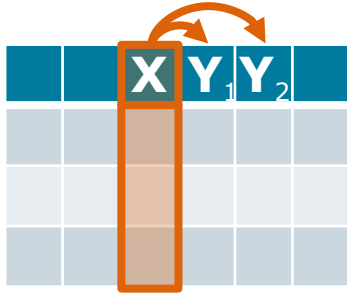
“The values in X functionally define the values in Y”

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

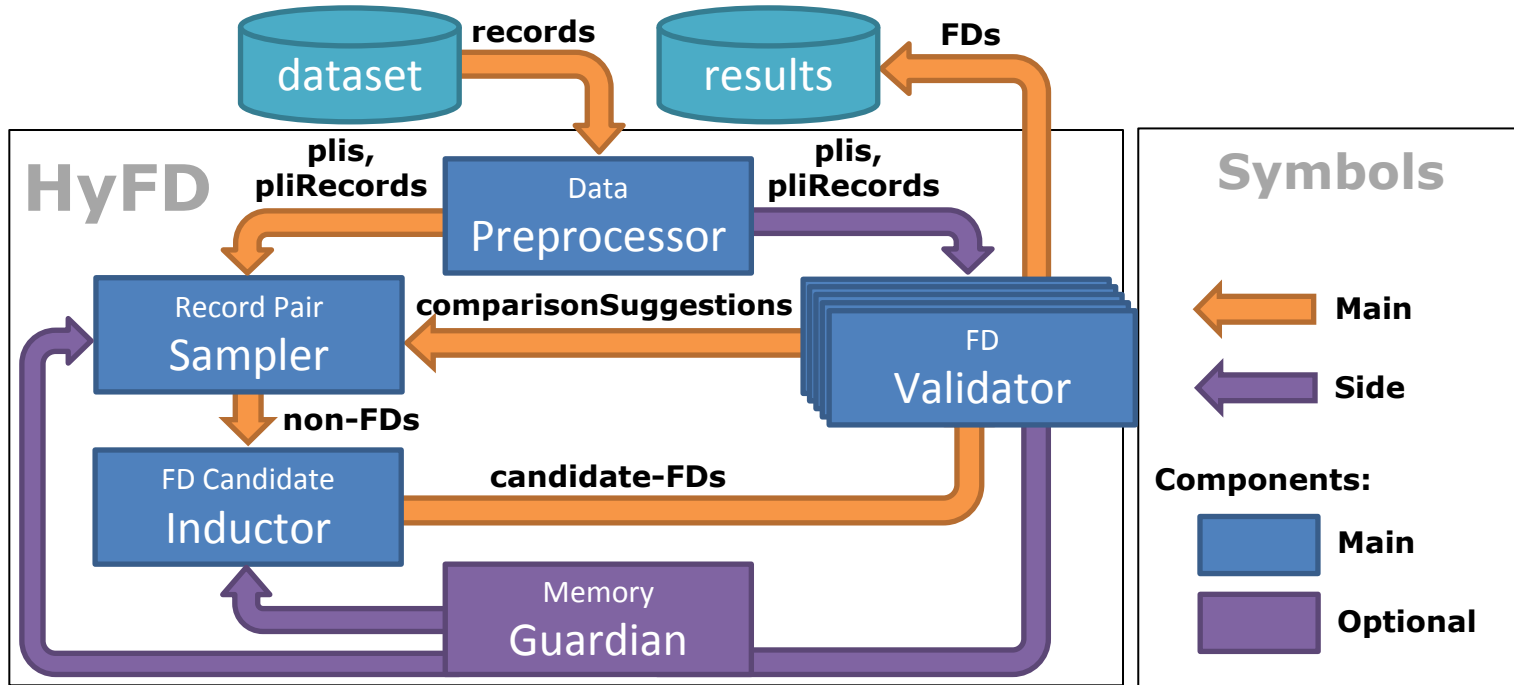
2014

2015

2016

2017

2018



2013

2014

2015

2016

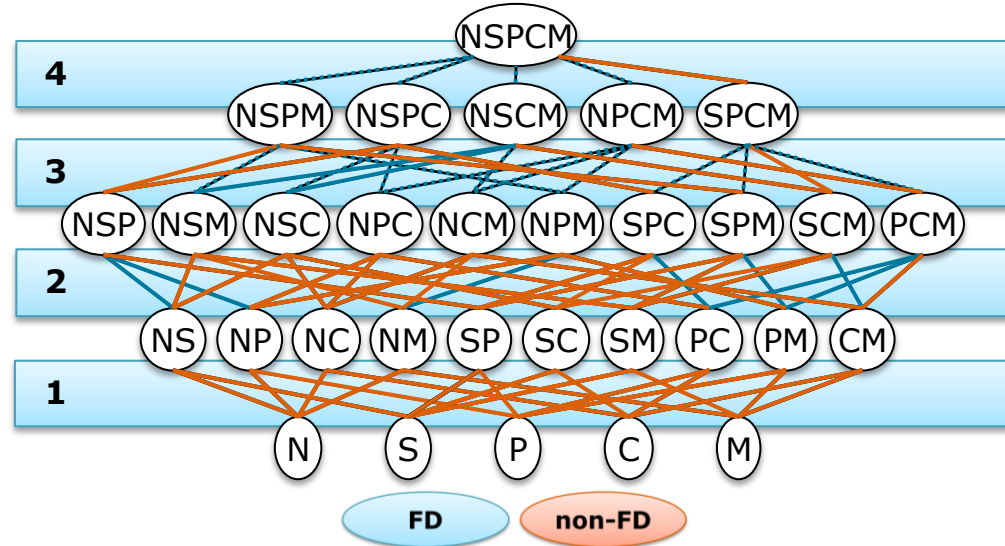
2017

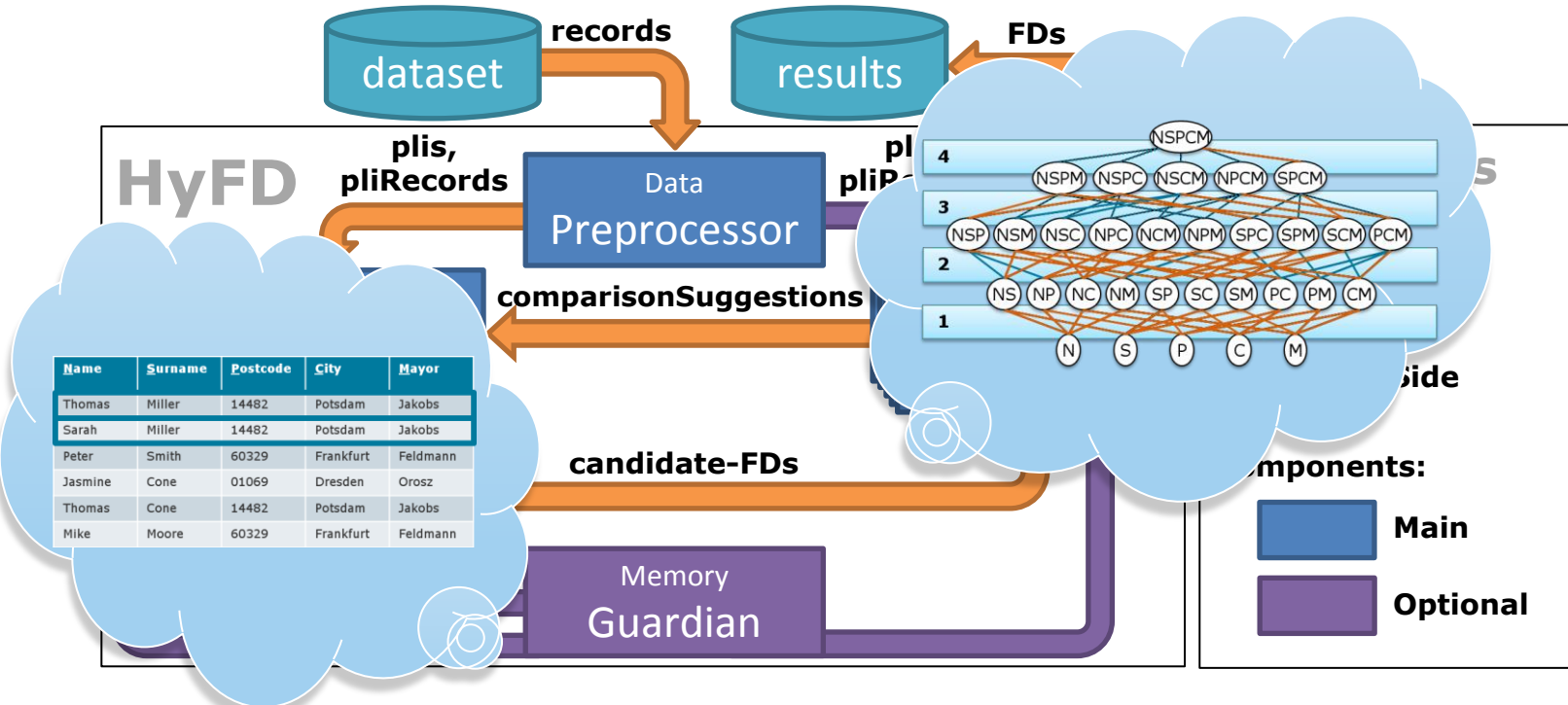
2018

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

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

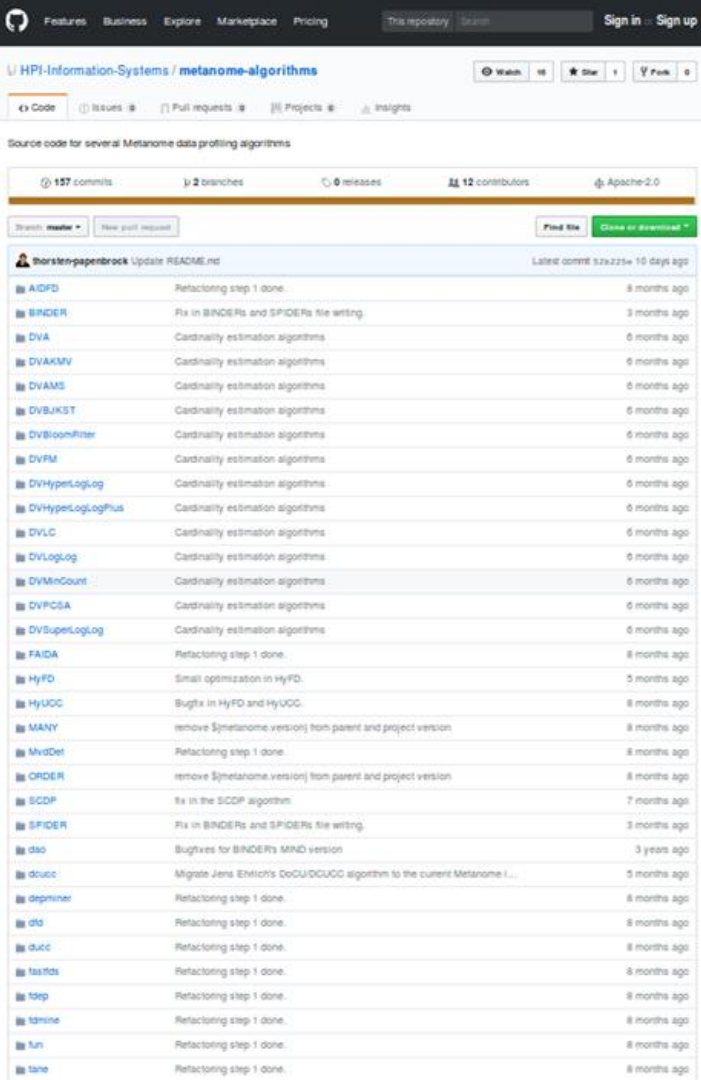
2015

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



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
	<b>21</b>	<b>17</b>

2013

2014

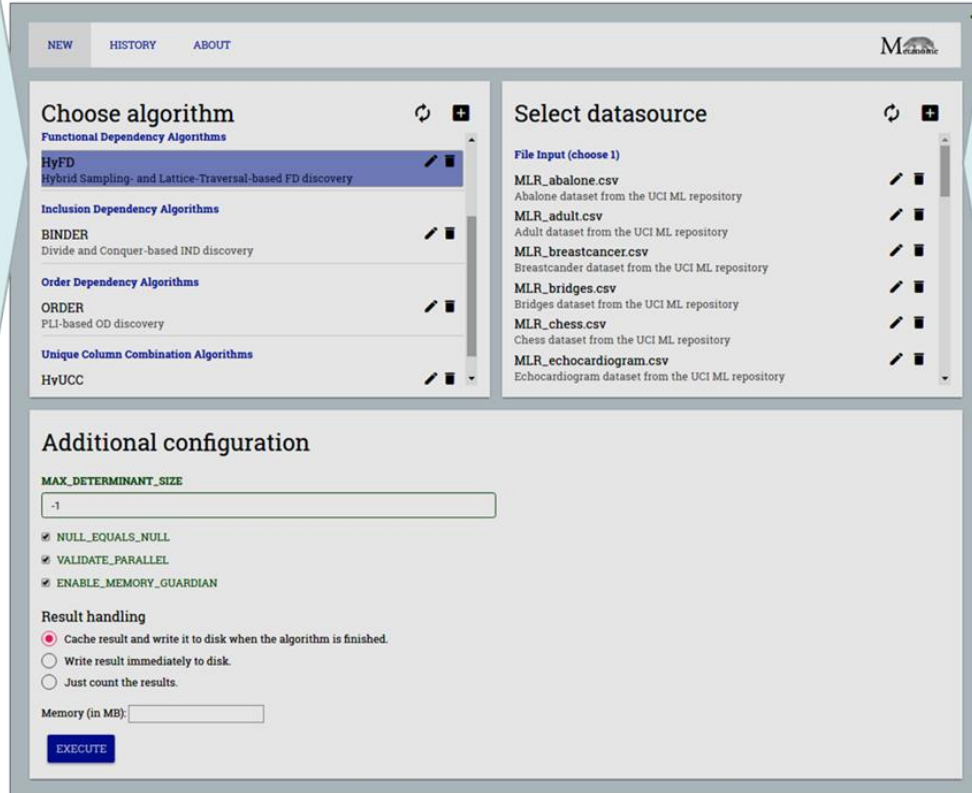
2015

2016

2017

2018

### Algorithms

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


NEW HISTORY ABOUT Metanome

#### 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

- 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**

### Datasets



2013

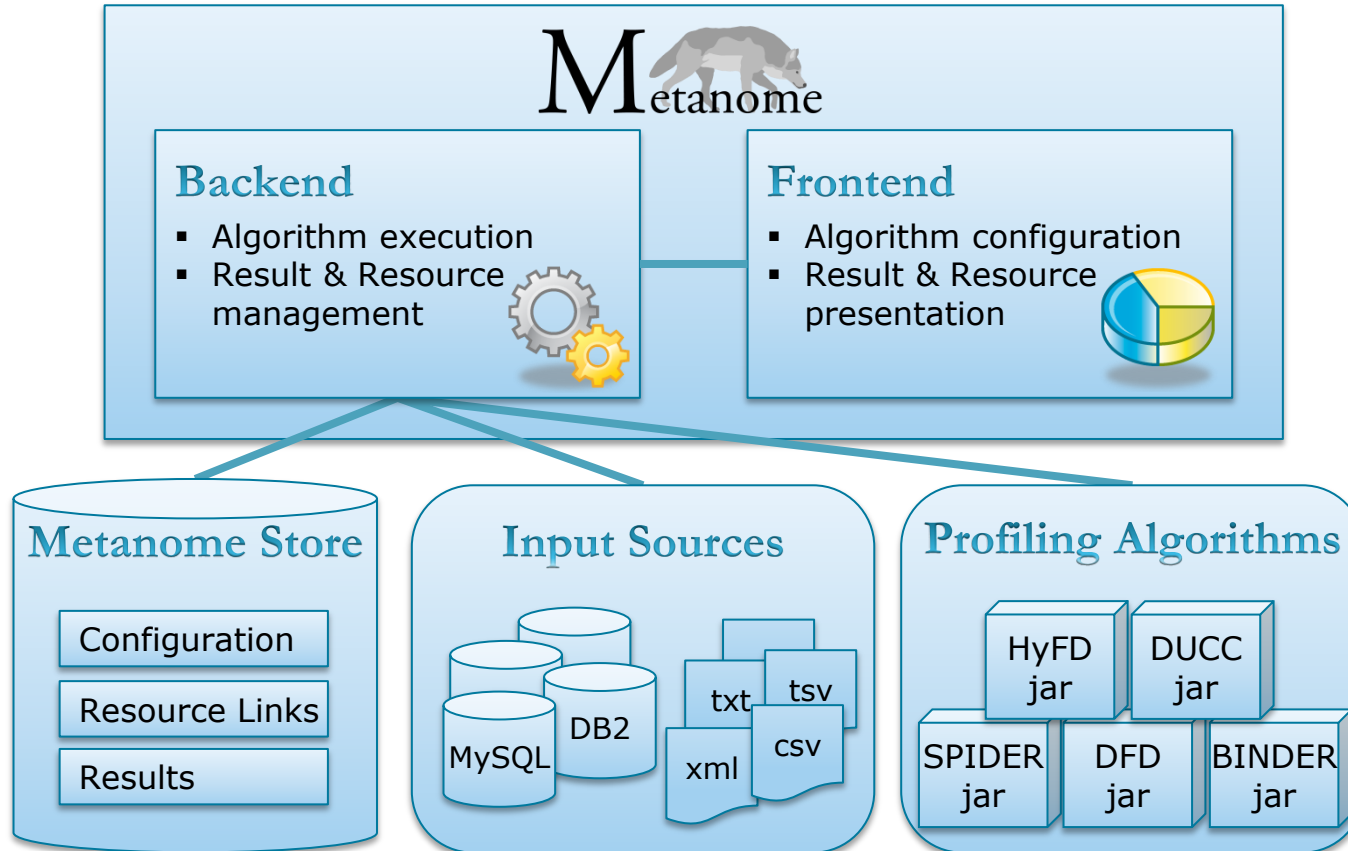
2014

2015

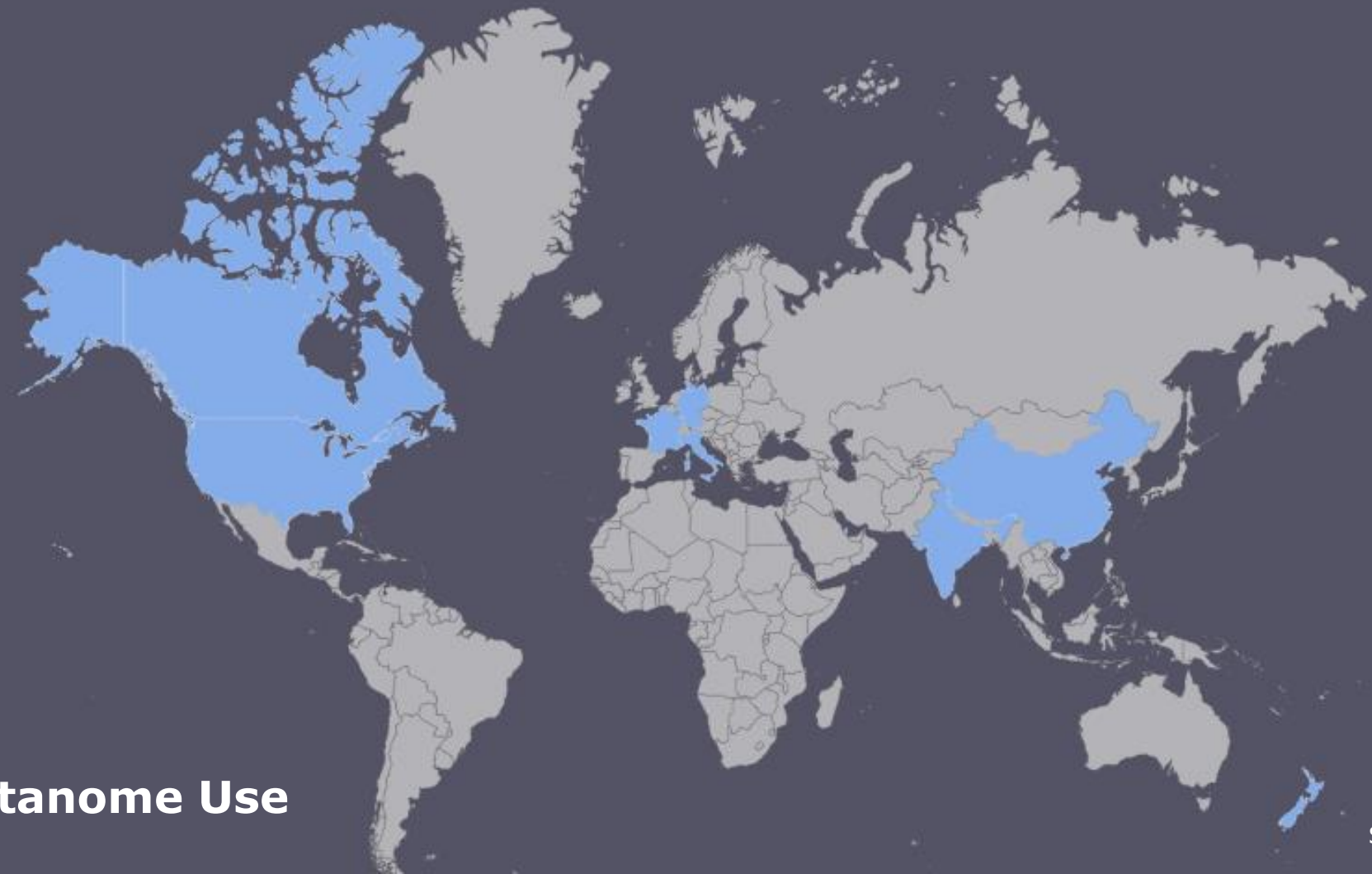
2016

2017

2018



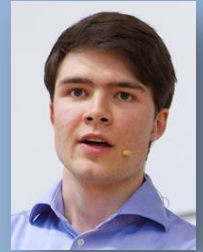
# Metanome Use



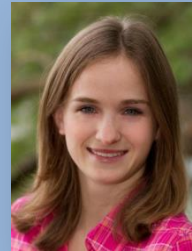
# Metanome Algorithm Research

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

# Agenda



Advances in Data Profiling



Challenges in Distributed Computing





## Query Optimization

(with Jan Kossmann, EPIC Chair)

## Data Cleaning

(with Ioannis Koumarelas, IS Chair)

Application



## Matching Dependency Discovery

(with Philipp Schirmer, Bakdata)

## Denial Constraint Discovery

(with Eduardo Pena, IS Chair)

Discovery



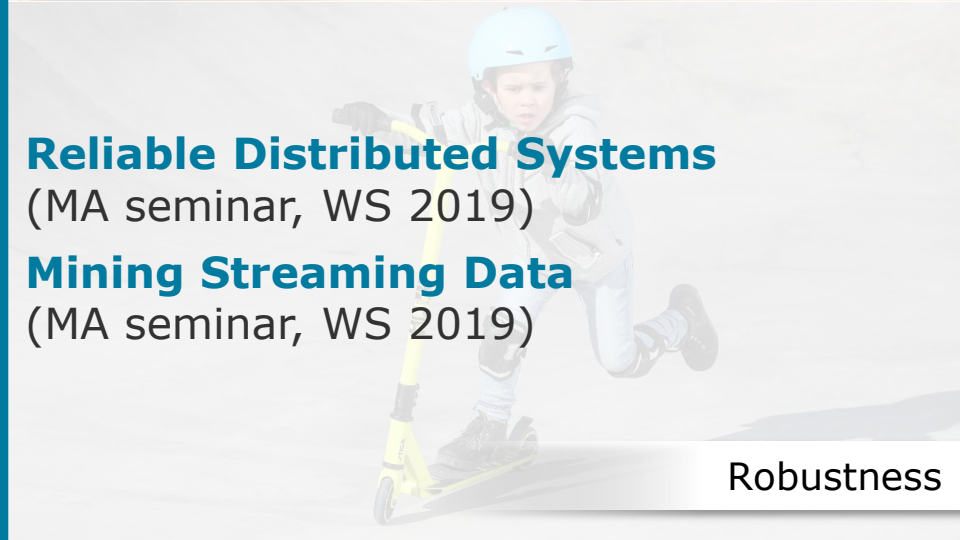
## Distributed UCC Discovery

(with Martin Schirneck, AE Chair)

## Distributed FD Discovery

(with Felix Naumann, IS Chair)

Distribution



## Reliable Distributed Systems

(MA seminar, WS 2019)

## Mining Streaming Data

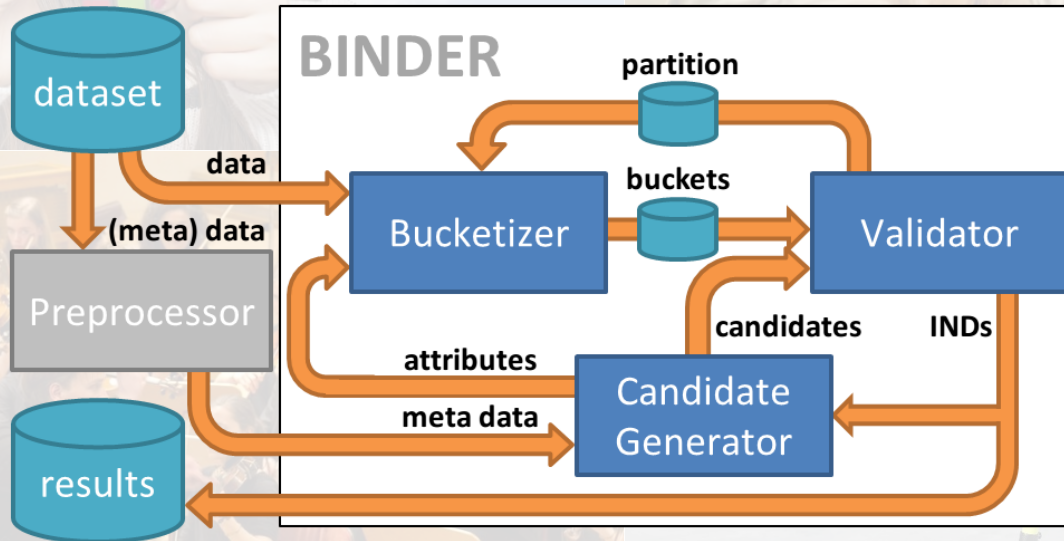
(MA seminar, WS 2019)

Robustness

# The Standard Batch Approach

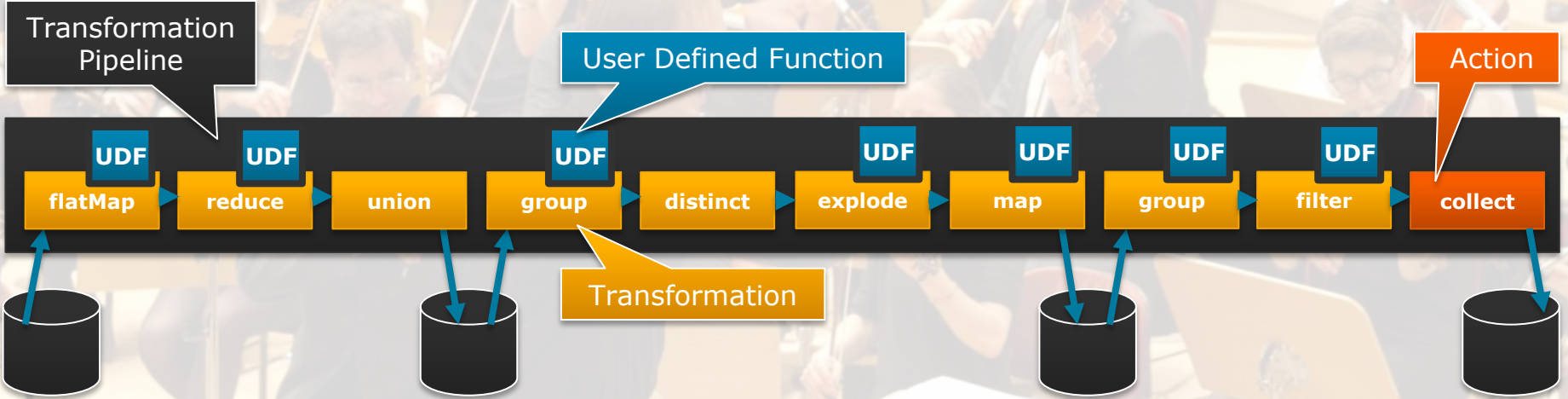


Apache Flink



Distribution

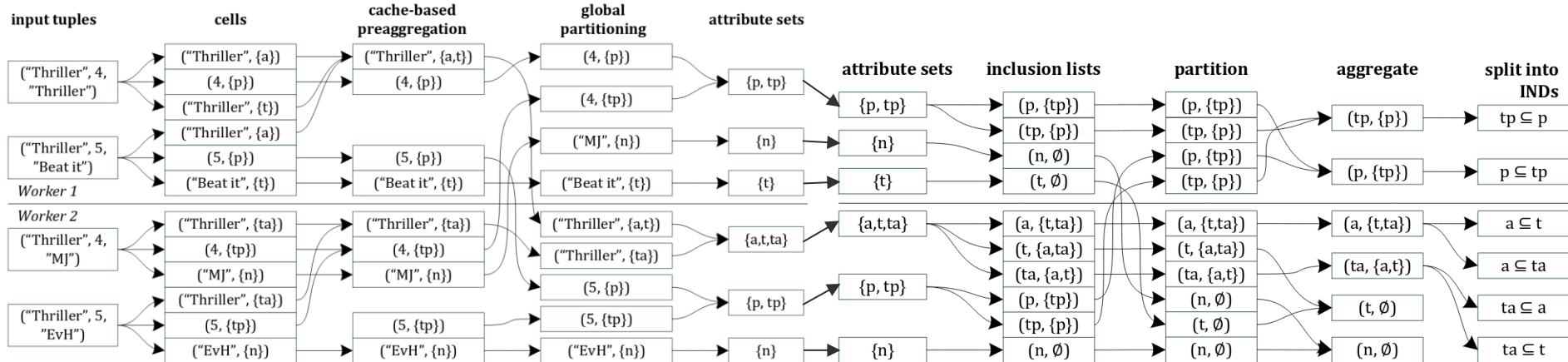
# The Standard Batch Approach



# The Standard Batch Approach



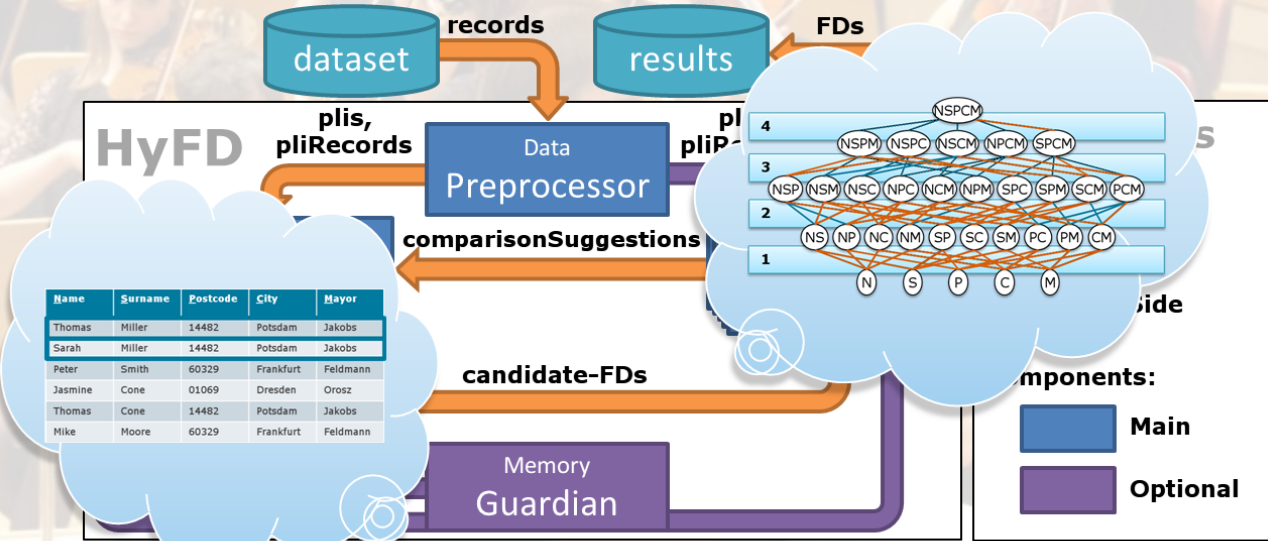
Apache Flink



# The Standard Batch Approach



Apache Flink

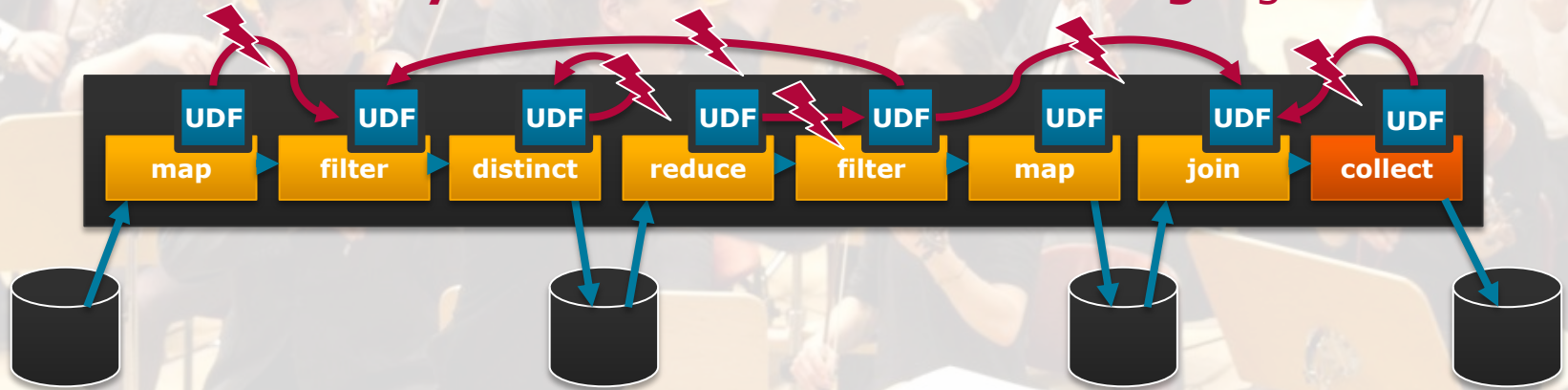


Name	Surname	Postcode	City	Mayor
Thomas	Miller	14482	Potsdam	Jakobs
Sarah	Miller	14482	Potsdam	Jakobs
Peter	Smith	60329	Frankfurt	Feldmann
Jasmine	Cone	01069	Dresden	Orosz
Thomas	Cone	14482	Potsdam	Jakobs
Mike	Moore	60329	Frankfurt	Feldmann

# The Standard Batch Approach



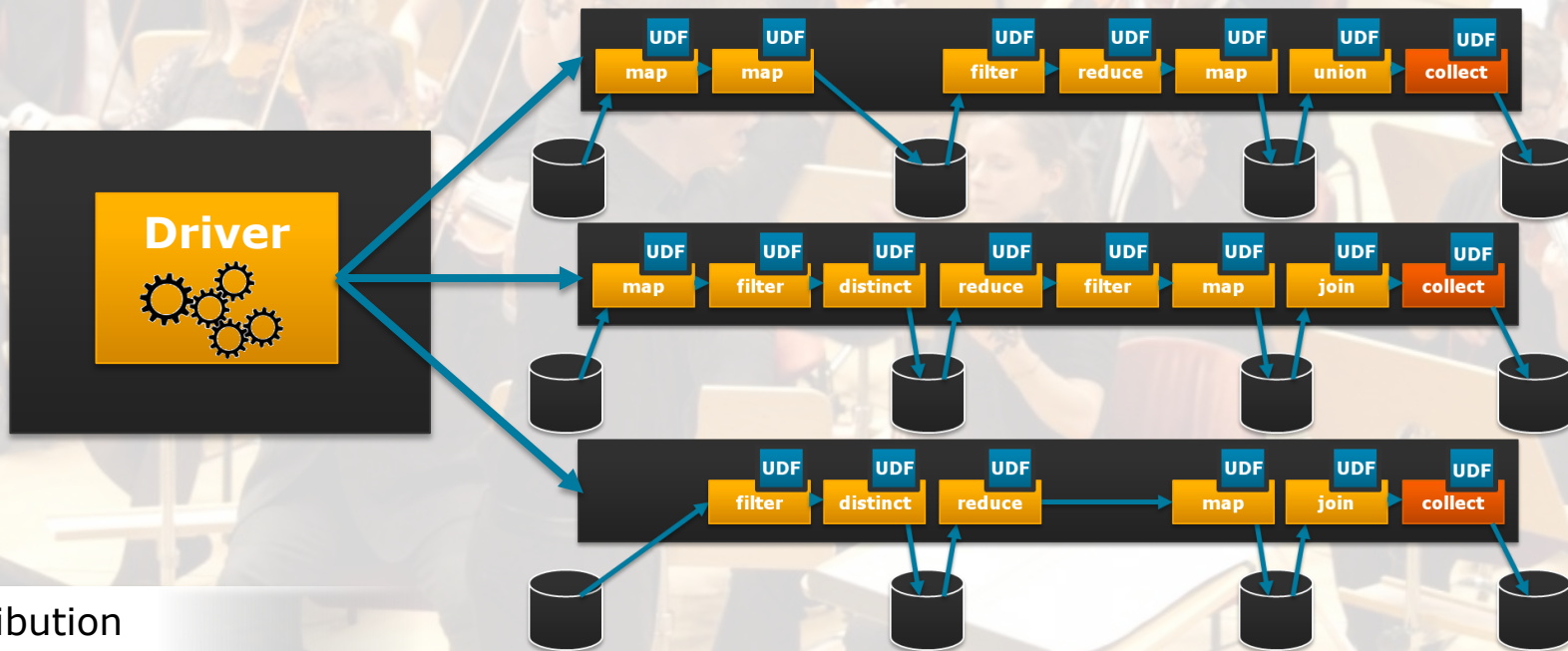
**dynamic** behavior and **branching** logic



# The Standard Batch Approach

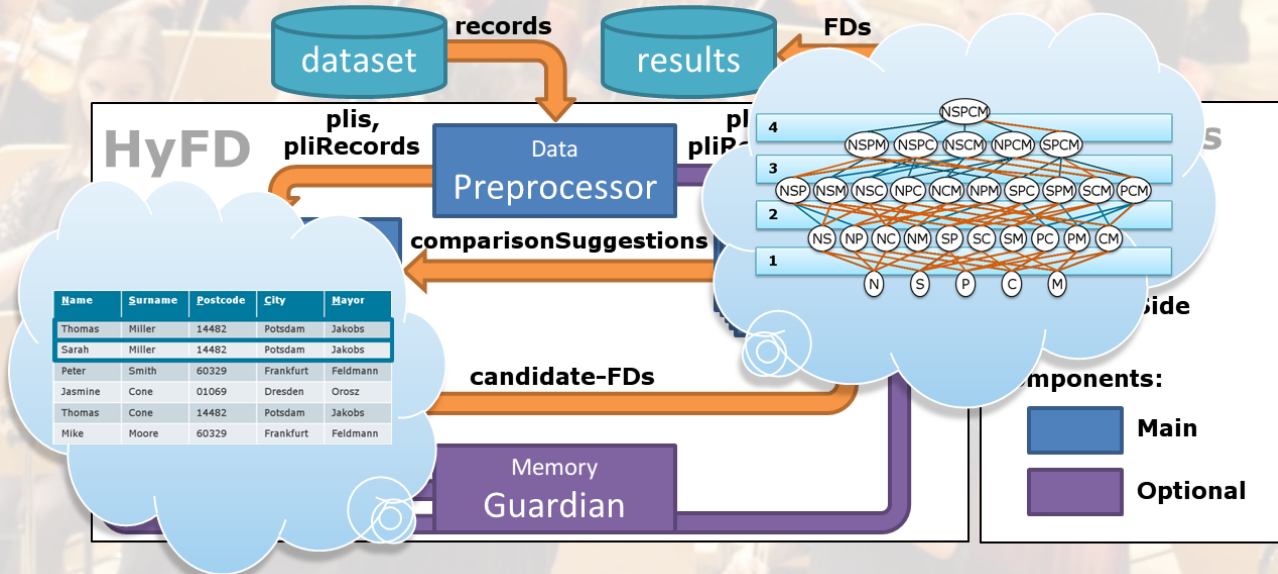


Apache Flink



Distribution

# The Low-Level Message-Passing Approach

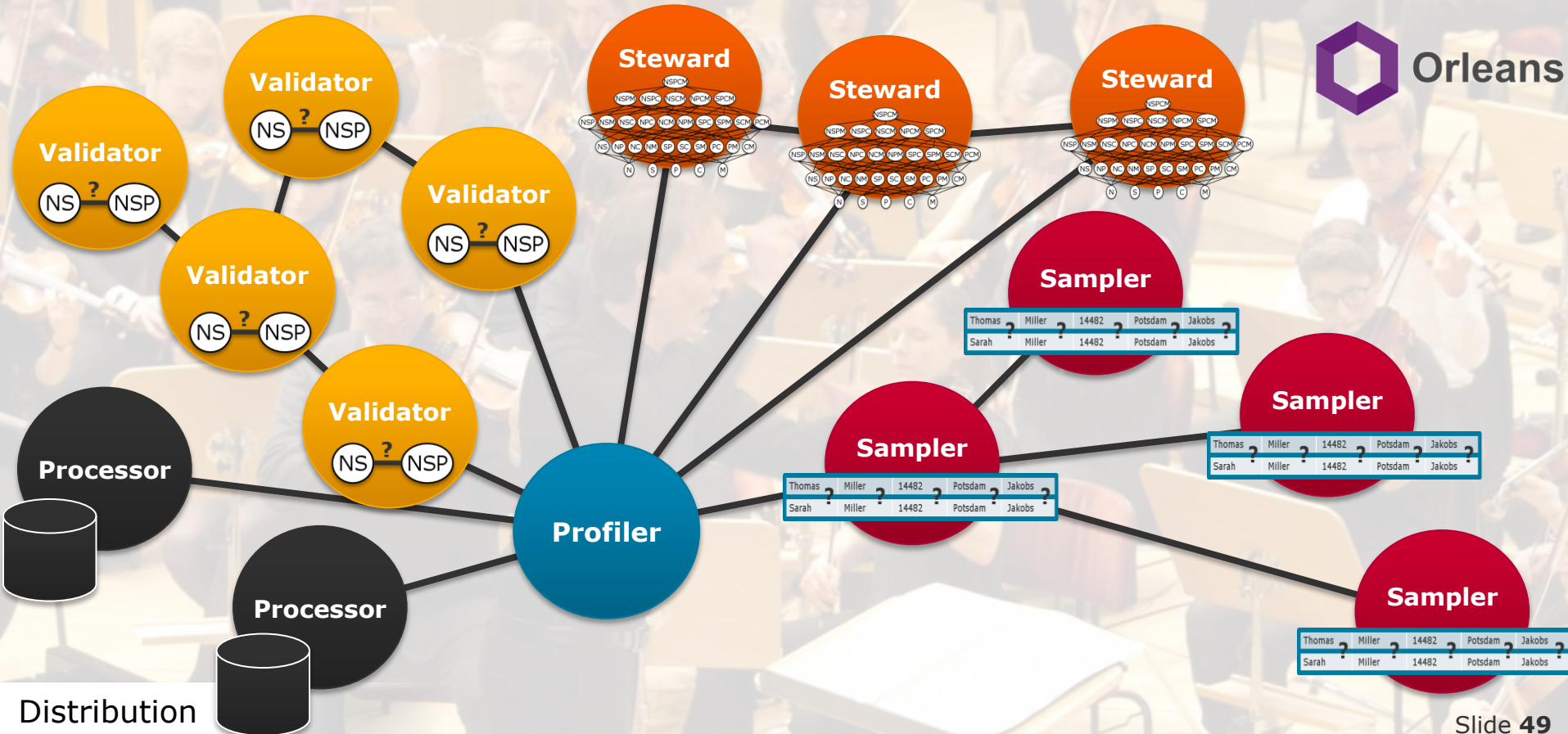


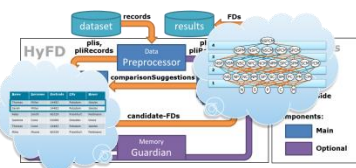
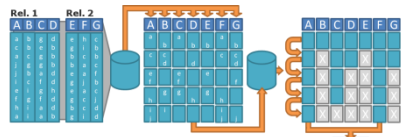
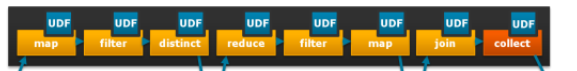
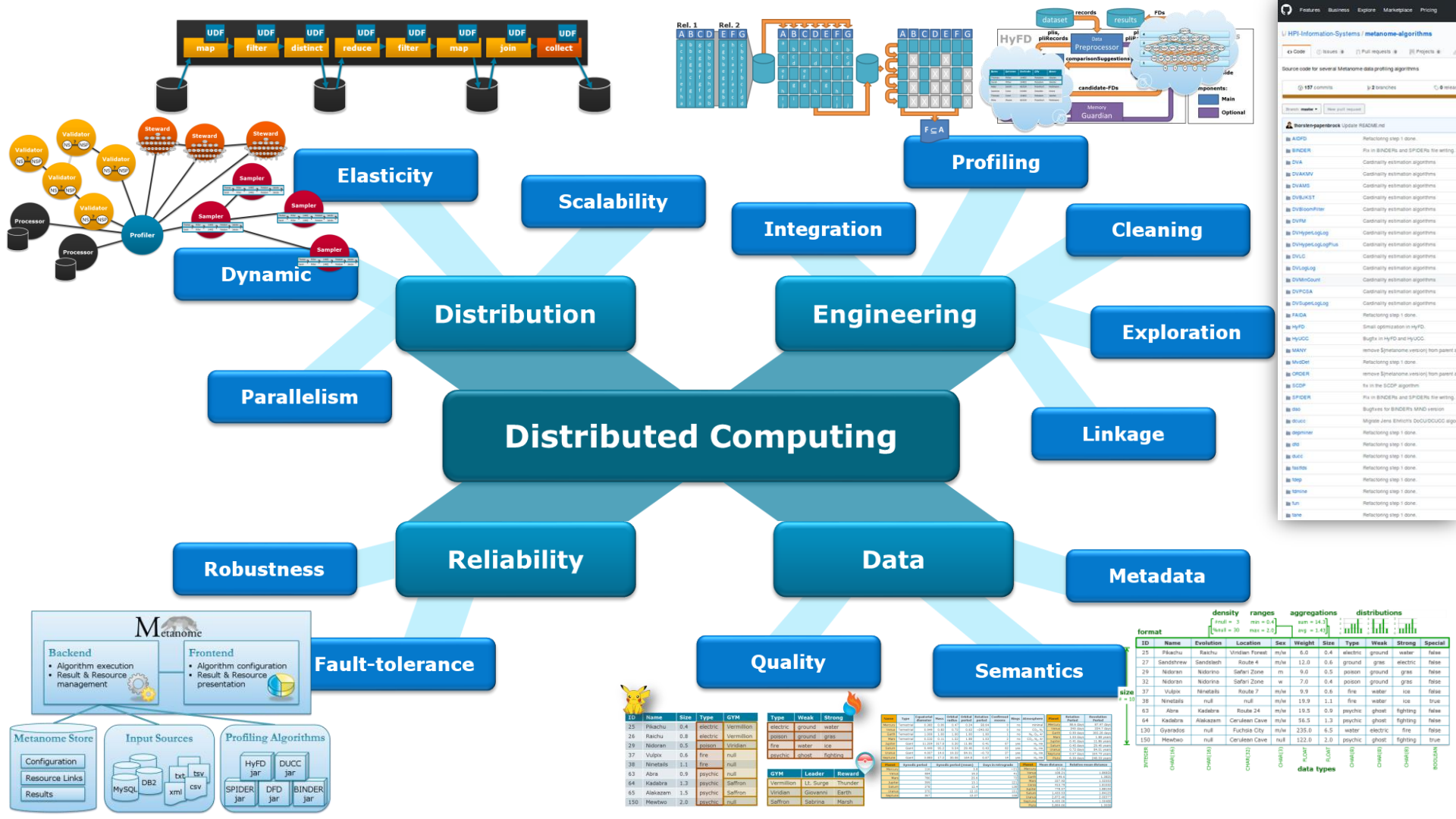
Name	Surname	Postcode	City	Mayor
Thomas	Miller	14482	Potsdam	Jakobs
Sarah	Miller	14482	Potsdam	Jakobs
Peter	Smith	60329	Frankfurt	Feldmann
Jasmine	Cone	01069	Dresden	Orosz
Thomas	Cone	14482	Potsdam	Jakobs
Mike	Moore	60329	Frankfurt	Feldmann

- Components:**
- Main**
  - Optional**



# The Low-Level Message-Passing Approach





UHP: Information Systems / metanome-algorithms

Code Issues Pull requests Projects

Source code for several Metanome data profiling algorithms

157 commits 2 branches 0 releases

metanome metanome

**metanome-preprocessor** Update README.md

- A-DFD Refactoring step 1 done
- BINDER Fix in BINDERs and SPIDERs file access
- DVA Centrality estimation algorithms
- DVAKMV Centrality estimation algorithms
- DVAMS Centrality estimation algorithms
- DVBKST Centrality estimation algorithms
- DVbloomfilter Centrality estimation algorithms
- DVFM Centrality estimation algorithms
- DVHyperLogLog Centrality estimation algorithms
- DVHyperLogLogPlus Centrality estimation algorithms
- DVLC Centrality estimation algorithms
- DVLogLog Centrality estimation algorithms
- DVMinCount Centrality estimation algorithms
- DVPCSA Centrality estimation algorithms
- DVSuperLogLog Centrality estimation algorithms
- FADA Refactoring step 1 done
- HyFD Small optimization in HyFD
- HyFOCC Bugfix in HyFD and HyFOCC
- HyFOV remove SpineSearcher, vectors, toString() and refactoring step 1 done
- HyKST remove SpineSearcher, vectors, toString() and refactoring step 1 done
- OCER remove SpineSearcher, vectors, toString() and refactoring step 1 done
- SCOP fix in the SCOP algorithm
- SPIDER fix in SPIDERs and SPIDERs file access
- SHO Bugfixes to BINDERs SHD version
- SHO1Migrate Java Elements ZooKeeper/SHO1Migrate step 1 done
- SHO2 Refactoring step 1 done
- SHO3 Refactoring step 1 done
- SHO4 Refactoring step 1 done
- SHO5 Refactoring step 1 done
- SHO6 Refactoring step 1 done
- SHO7 Refactoring step 1 done
- SHO8 Refactoring step 1 done
- SHO9 Refactoring step 1 done
- SHO10 Refactoring step 1 done



ID	Name	Size	Type	GYM	Type	Weak	Strong
25	Pikachu	0.4	electric	Vermillion	electric	ground	water
26	Raichu	0.8	electric	Vermillion	poison	ground	grass
28	Nidoran	0.6	poison	Ninetales	fire	water	ice
37	Vulpix	0.6	fire	null	psychic	ghost	fighting
38	Ninetales	1.1	fire	null	fire	ice	false
63	Abra	0.9	psychic	null	psychic	ghost	fighting
64	Kabra	1.2	psychic	Saffron	psychic	ghost	fighting
65	Alakazam	1.5	psychic	Saffron	psychic	ghost	fighting
150	Mewtwo	2.0	psychic	null	Saffron	Sabrina	Marsh

Year	Population	Area	Population Density	Population Growth	Population Change	Population Change Rate	Population Change Rate Rate
2000	1000000	10000000000	1000	10000000000	10000000000	10000000000	10000000000
2001	1000000	10000000000	1000	10000000000	10000000000	10000000000	10000000000
2002	1000000	10000000000	1000	10000000000	10000000000	10000000000	10000000000
2003	1000000	10000000000	1000	10000000000	10000000000	10000000000	10000000000
2004	1000000	10000000000	1000	10000000000	10000000000	10000000000	10000000000
2005	1000000	10000000000	1000	10000000000	10000000000	10000000000	10000000000
2006	1000000	10000000000	1000	10000000000	10000000000	10000000000	10000000000
2007	1000000	10000000000	1000	10000000000	10000000000	10000000000	10000000000
2008	1000000	10000000000	1000	10000000000	10000000000	10000000000	10000000000
2009	1000000	10000000000	1000	10000000000	10000000000	10000000000	10000000000
2010	1000000	10000000000	1000	10000000000	10000000000	10000000000	10000000000
2011	1000000	10000000000	1000	10000000000	10000000000	10000000000	10000000000
2012	1000000	10000000000	1000	10000000000	10000000000	10000000000	10000000000
2013	1000000	10000000000	1000	10000000000	10000000000	10000000000	10000000000
2014	1000000	10000000000	1000	10000000000	10000000000	10000000000	10000000000
2015	1000000	10000000000	1000	10000000000	10000000000	10000000000	10000000000
2016	1000000	10000000000	1000	10000000000	10000000000	10000000000	10000000000
2017	1000000	10000000000	1000	10000000000	10000000000	10000000000	10000000000
2018	1000000	10000000000	1000	10000000000	10000000000	10000000000	10000000000
2019	1000000	10000000000	1000	10000000000	10000000000	10000000000	10000000000
2020	1000000	10000000000	1000	10000000000	10000000000	10000000000	10000000000

format

density ranges aggregations distributions

format = 2 (first = 2 min = 0.4 (last = 30 max = 2.0) sum = 34.1 avg = 1.4)

ID	Name	Evolution	Location	Sex	Weight	Size	Type	Weak	Strong	Special
25	Pikachu	Raichu	Vindian Forest	m/f	6.0	0.4	electric	ground	water	false
27	Sandshrew	Sandslash	Route 4	m/f	12.0	0.6	ground	grass	electric	false
29	Nidoran	Nidorina	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	Ninetales	Route 7	m/f	9.9	0.6	fire	water	ice	false
38	Ninetales	null	null	m/f	19.9	1.1	fire	water	ice	true
63	Abra	Kadabra	Route 24	m/f	18.5	0.9	psychic	ghost	fighting	false
64	Kadabra	Alakazam	Cerulean Cave	m/f	16.5	1.3	psychic	ghost	fighting	false
130	Overdrive	null	Fuchsia City	m/f	235.0	6.5	water	electric	fire	false
150	Mewtwo	null	Cerulean Cave	null	122.0	2.0	psychic	ghost	fighting	true

size = 10

data types