



Data Profiling Functional Dependency Discovery

Thorsten Papenbrock
G-3.1.09, Campus III
Hasso Plattner Institut

Definition 2.1 (Functional dependency). 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 $\forall t_i, t_j \in r : t_i[X] = t_j[X] \Rightarrow t_i[A] = t_j[A]$.

FD discovery is in $\mathcal{O}(n^2 \cdot 2^m \cdot (\frac{m}{2})^2)$

Data Profiling

FD Discovery

Thorsten Papenbrock
Chart 2

- Motivation
- Approaches
- fdep
- HyFD
- Evaluation
- Metanome

Data Profiling

FD Discovery

Thorsten Papenbrock
Chart **3**

- **Motivation**
- Approaches
- fdep
- HyFD
- Evaluation
- Metanome



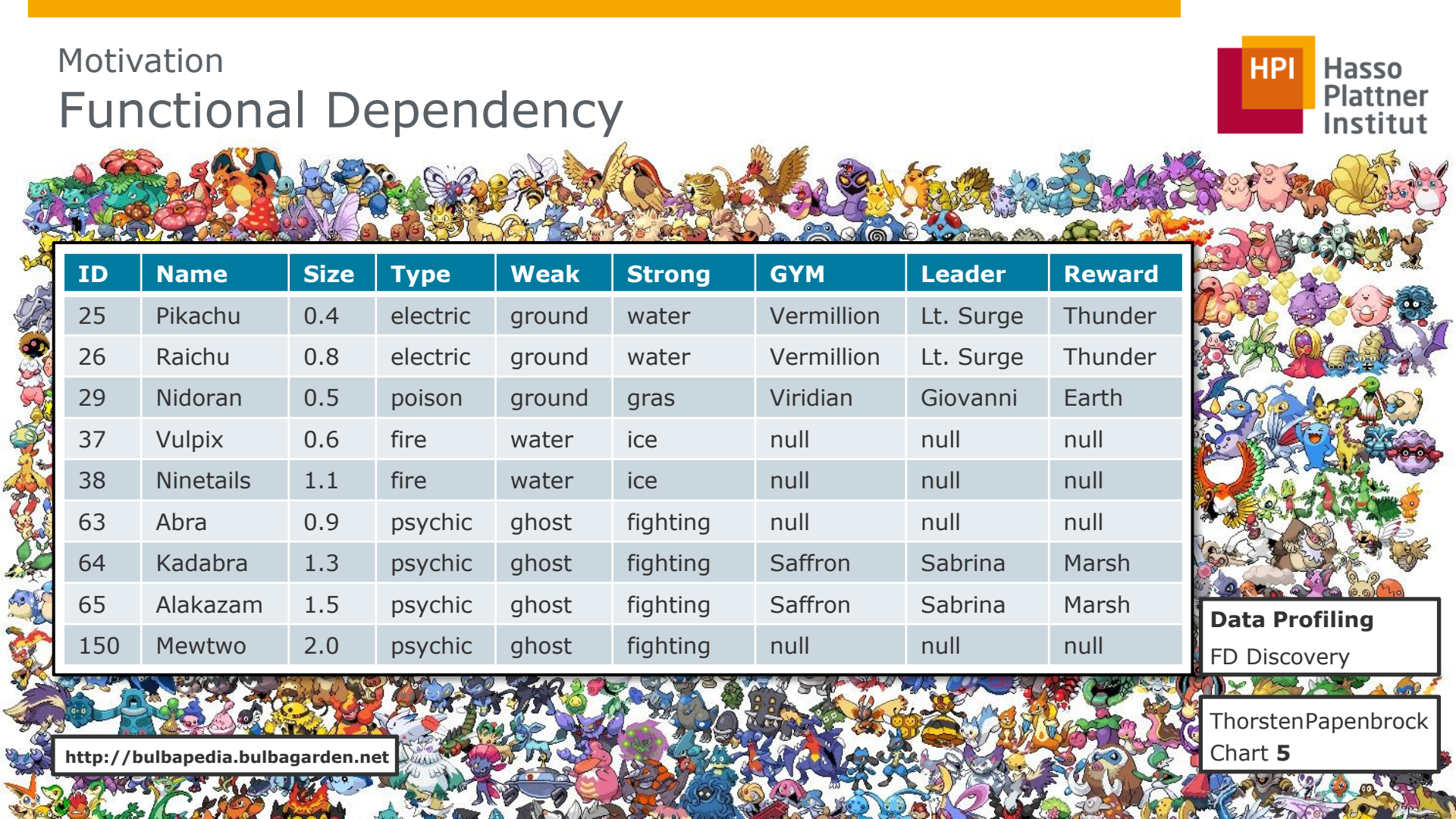
Data Profiling

FD Discovery

ThorstenPapenbrock
Chart 4

Motivation

Functional Dependency



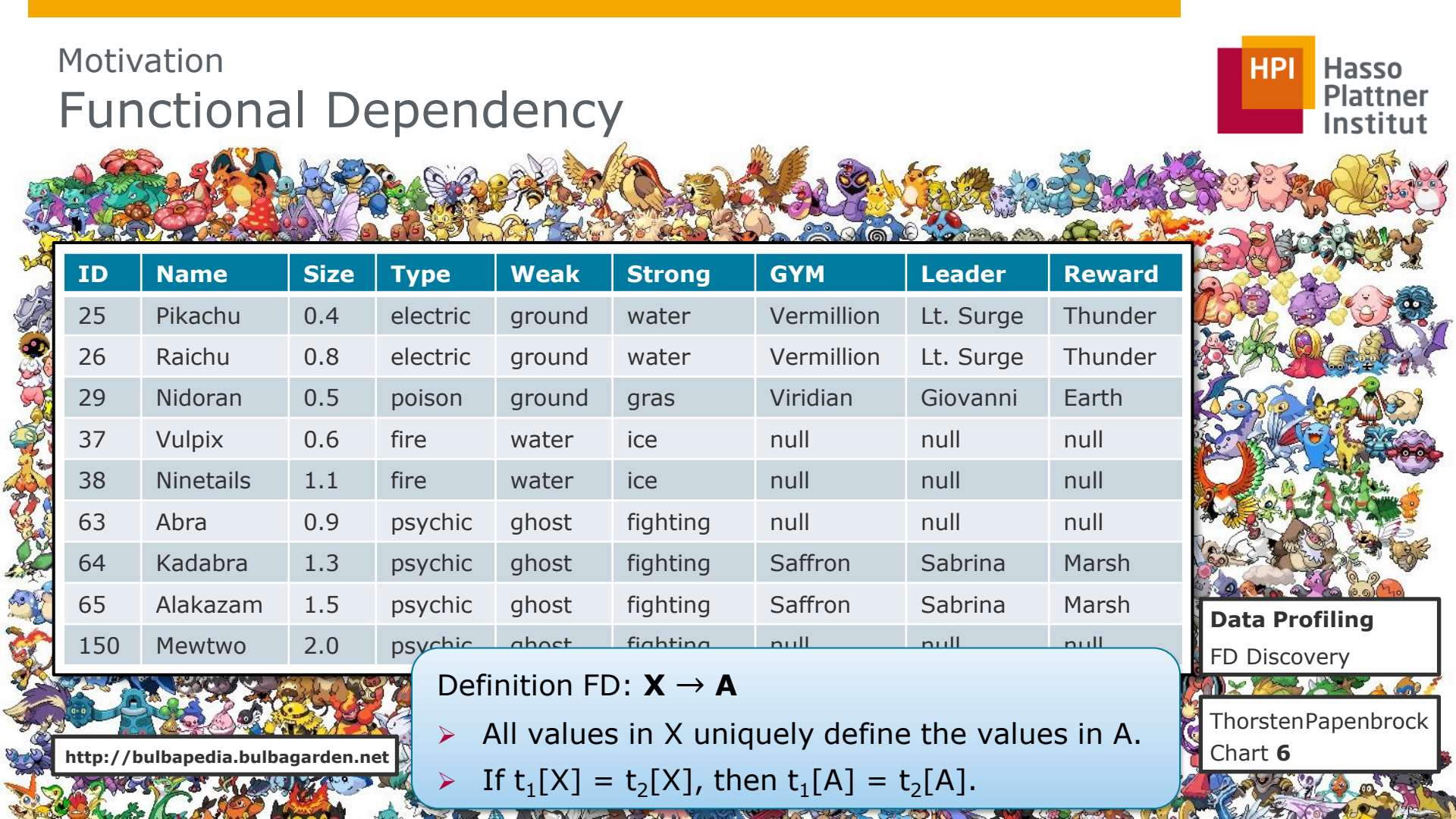
ID	Name	Size	Type	Weak	Strong	GYM	Leader	Reward
25	Pikachu	0.4	electric	ground	water	Vermillion	Lt. Surge	Thunder
26	Raichu	0.8	electric	ground	water	Vermillion	Lt. Surge	Thunder
29	Nidoran	0.5	poison	ground	grass	Viridian	Giovanni	Earth
37	Vulpix	0.6	fire	water	ice	null	null	null
38	Ninetails	1.1	fire	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

Data Profiling

FD Discovery

ThorstenPapenbrock
Chart 5

Functional Dependency



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	grass	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

<http://bulbapedia.bulbagarden.net>

Definition FD: $\mathbf{X} \rightarrow \mathbf{A}$

- All values in X uniquely define the values in A.
- If $t_1[X] = t_2[X]$, then $t_1[A] = t_2[A]$.

Data Profiling

FD Discovery

ThorstenPapenbrock
Chart 6


 Redundancy!

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

Data Profiling

FD Discovery


ThorstenPapenbrock
Chart 7

Functional Dependency



-25 cells!

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	gras
fire	water	ice
psychic	ghost	fighting



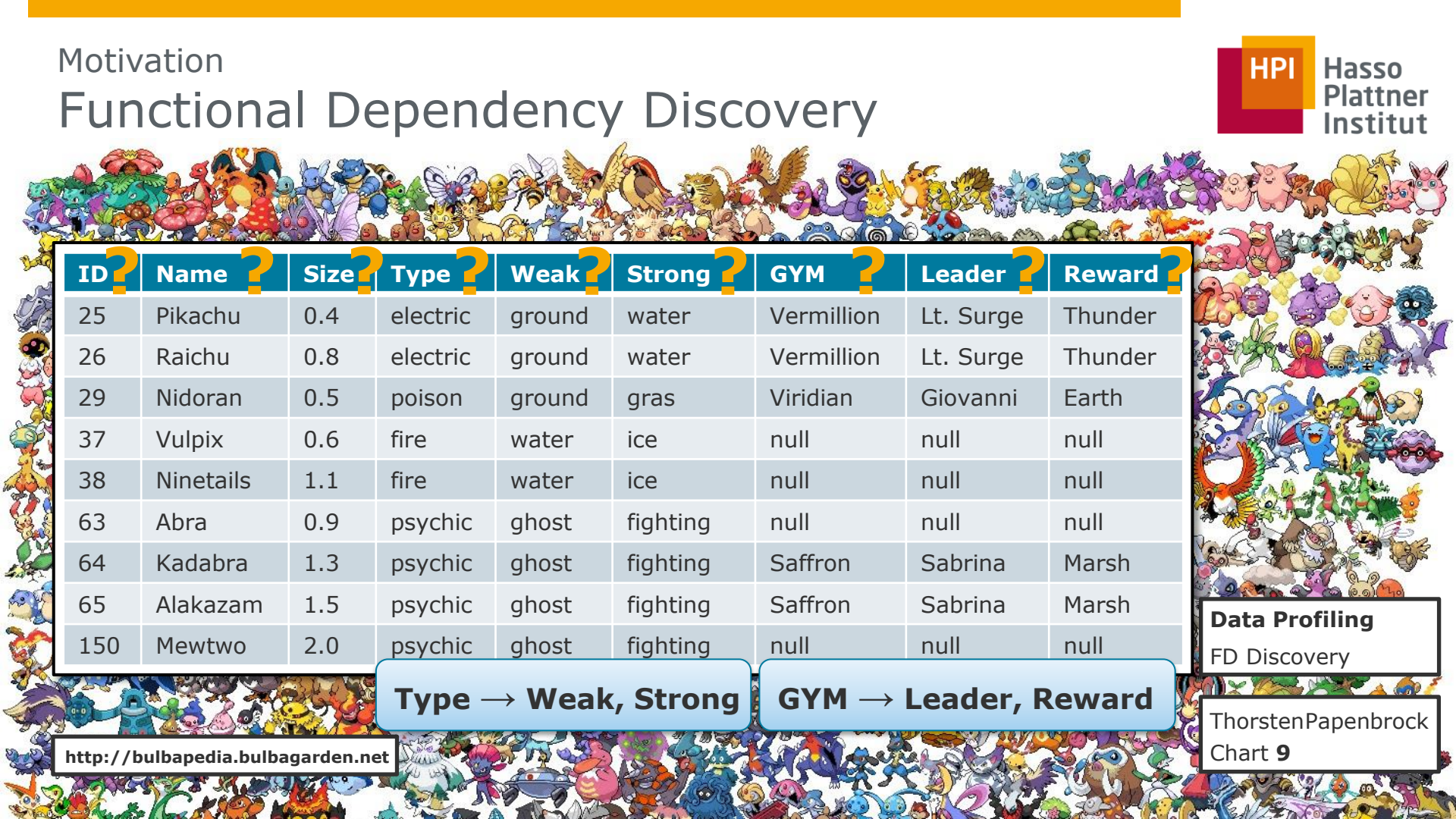
GYM	Leader	Reward
Vermillion	Lt. Surge	Thunder
Viridian	Giovanni	Earth
Saffron	Sabrina	Marsh

Data Profiling
FD Discovery

ThorstenPapenbrock
Chart 8

Motivation

Functional Dependency Discovery



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

Data Profiling
FD Discovery

ThorstenPapenbrock
Chart 9

<http://bulbapedia.bulbagarden.net>

Functional Dependency Discovery

- Motivation
- **Approaches**
 - Lattice Traversal
 - Dependency Induction
 - Difference- and Agree-Set
- fdep
- HyFD
- Evaluation
- Metanome

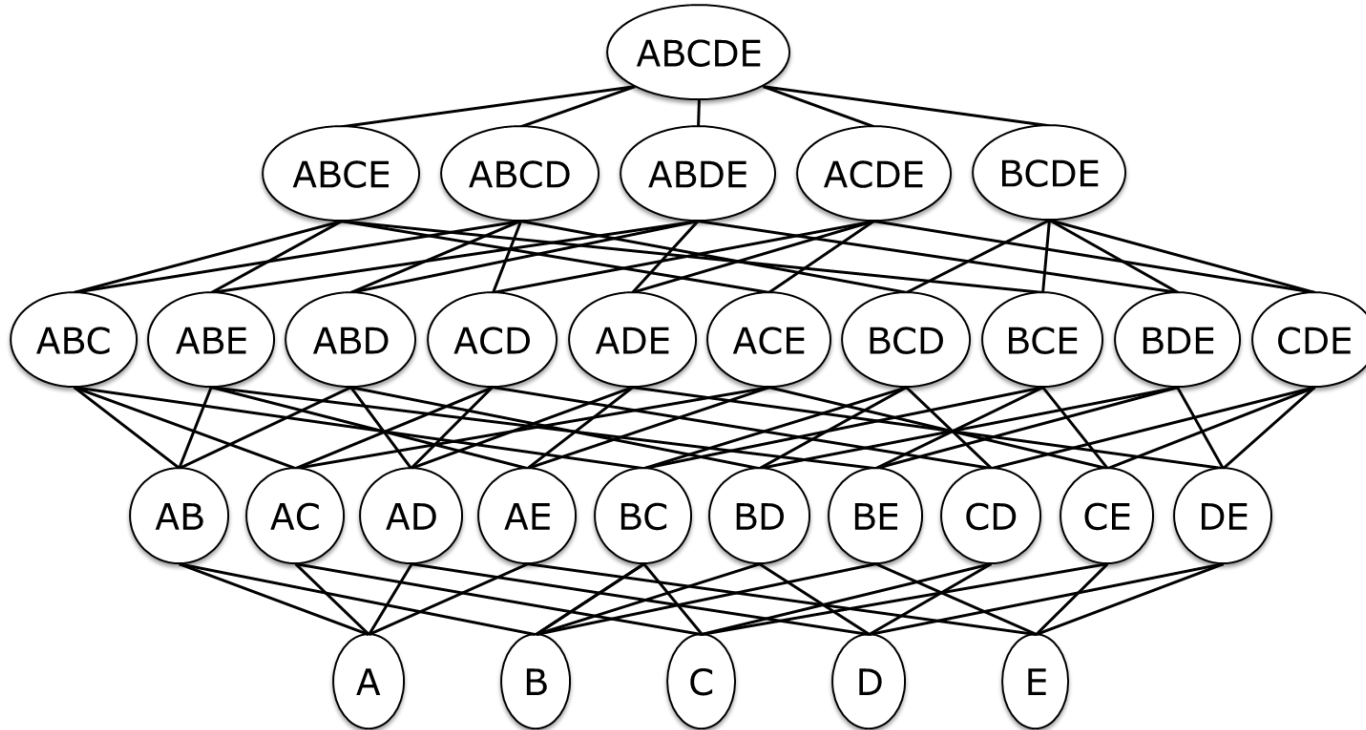


Data Profiling
Discovery

Thorsten Papenbrock
Chart 11

Approaches

Lattice Traversal



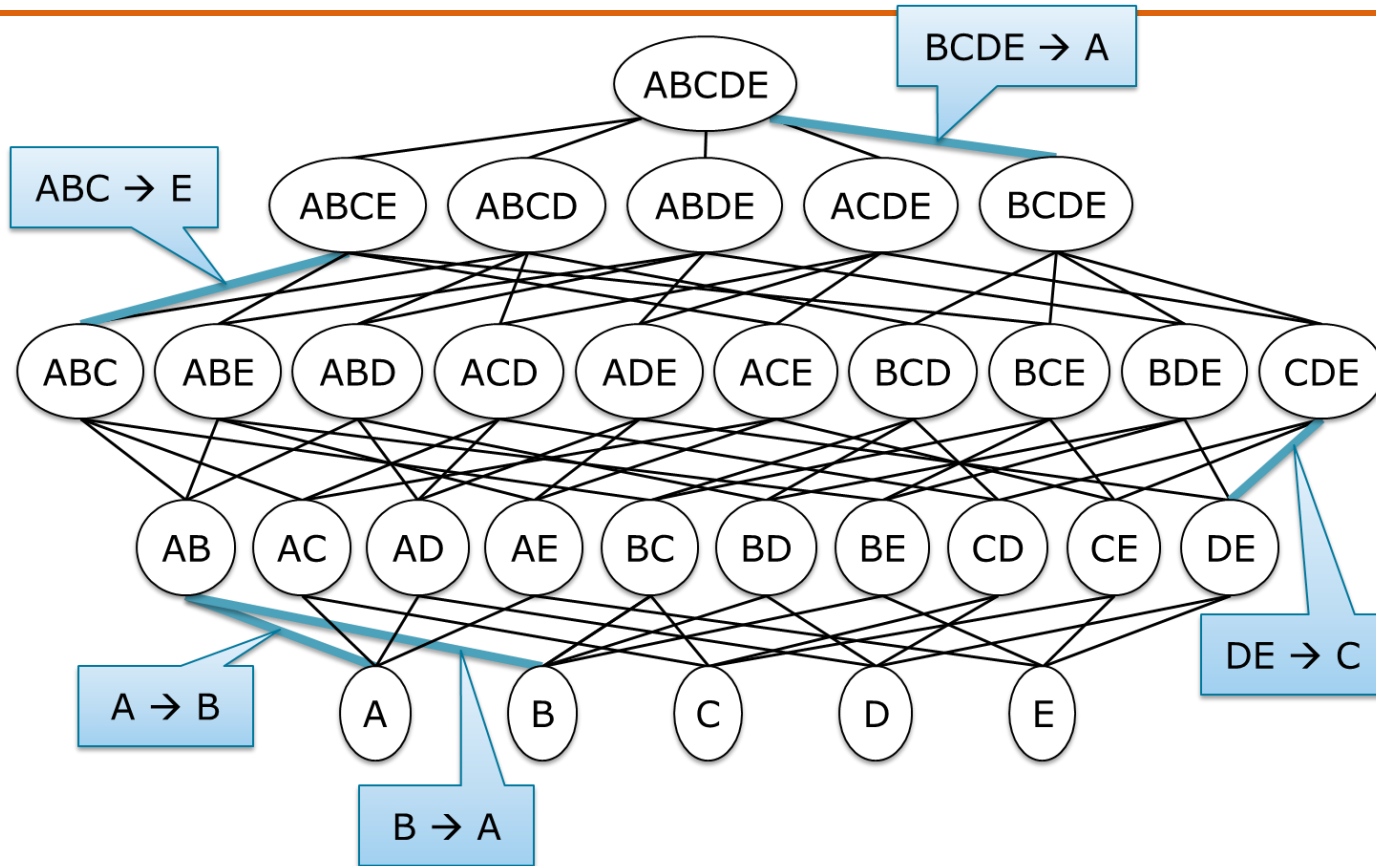
Data Profiling

FD Discovery

Thorsten Papenbrock
Chart 12

Approaches

Lattice Traversal



Data Profiling

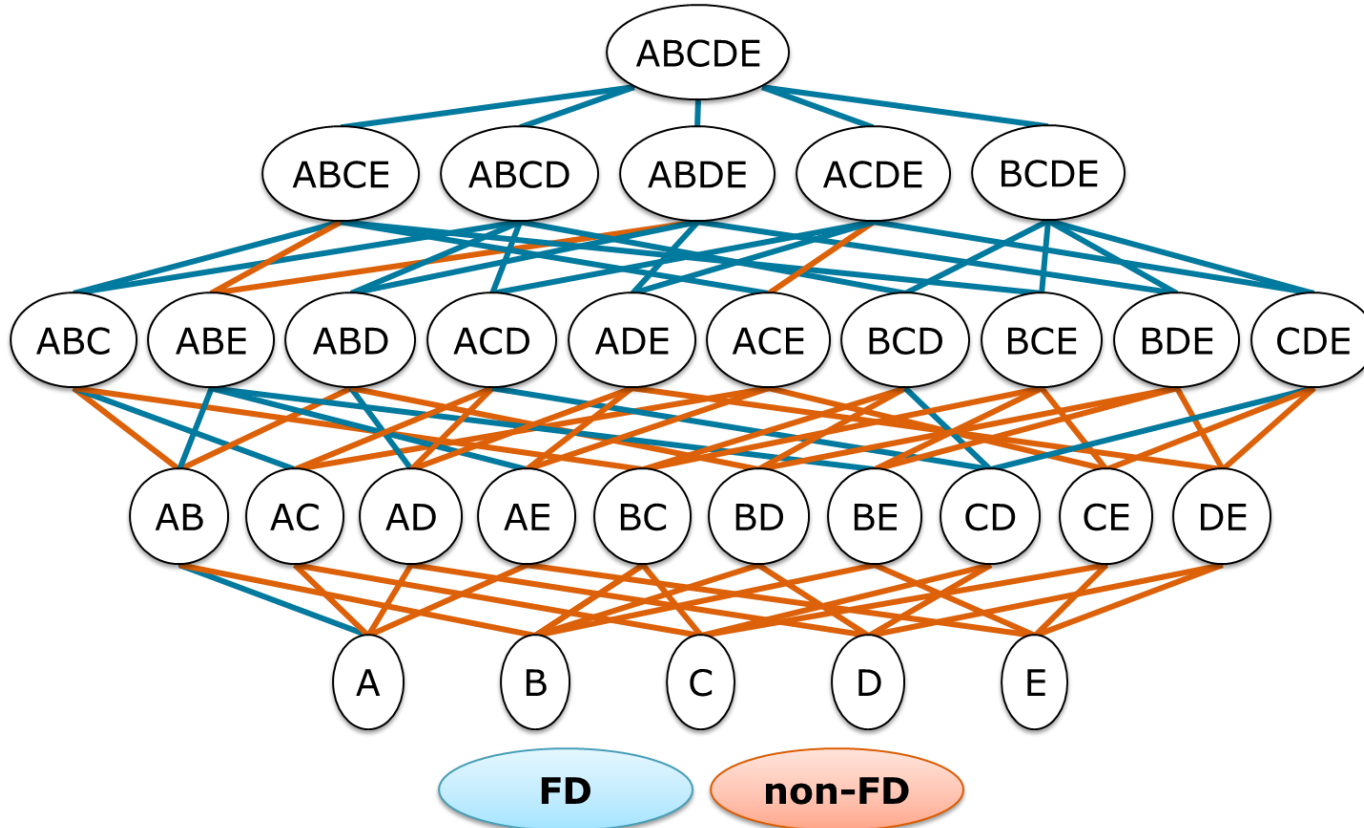
FD Discovery

ThorstenPapenbrock

Chart 13

Approaches

Lattice Traversal



Data Profiling

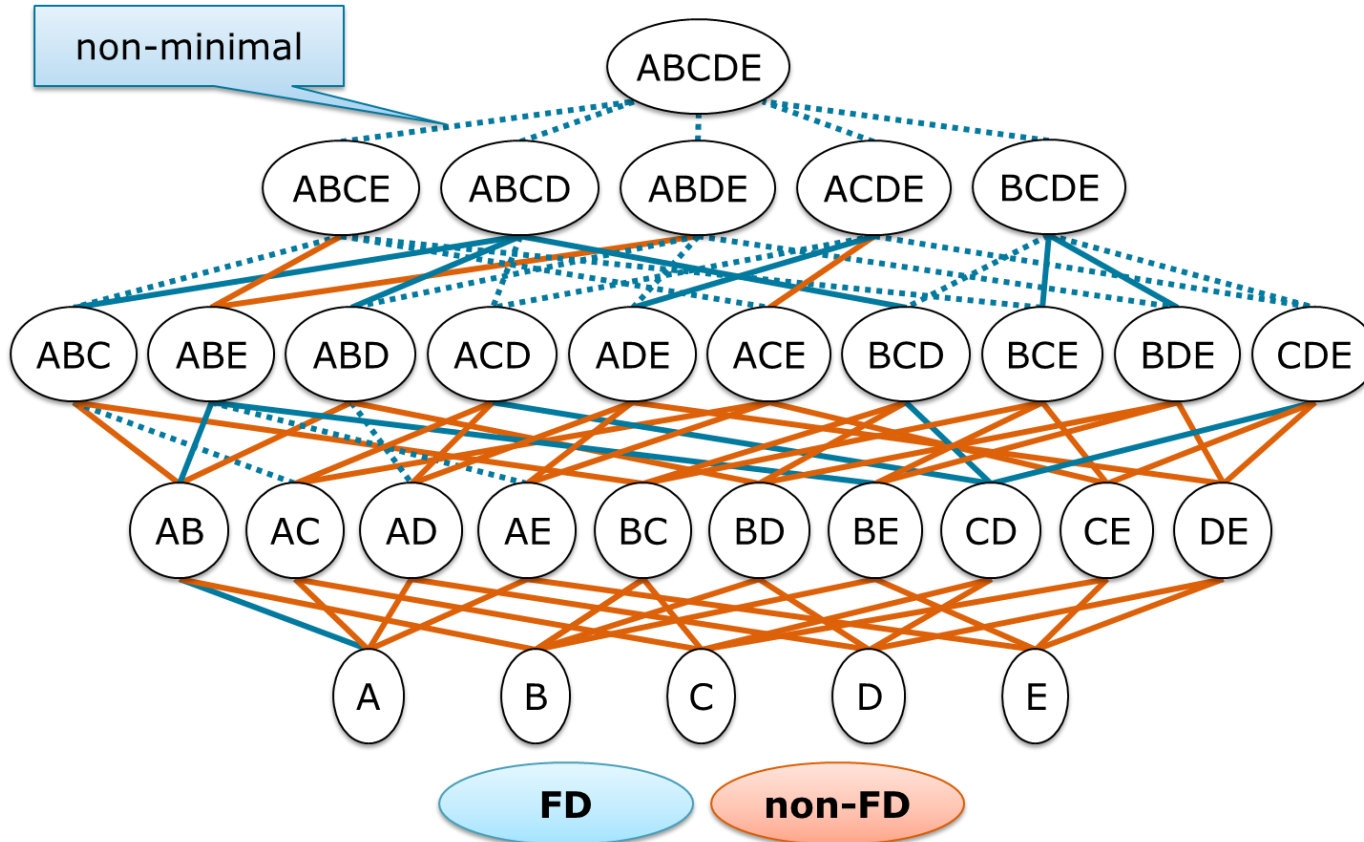
FD Discovery

Thorsten Papenbrock

Chart 14

Approaches

Lattice Traversal

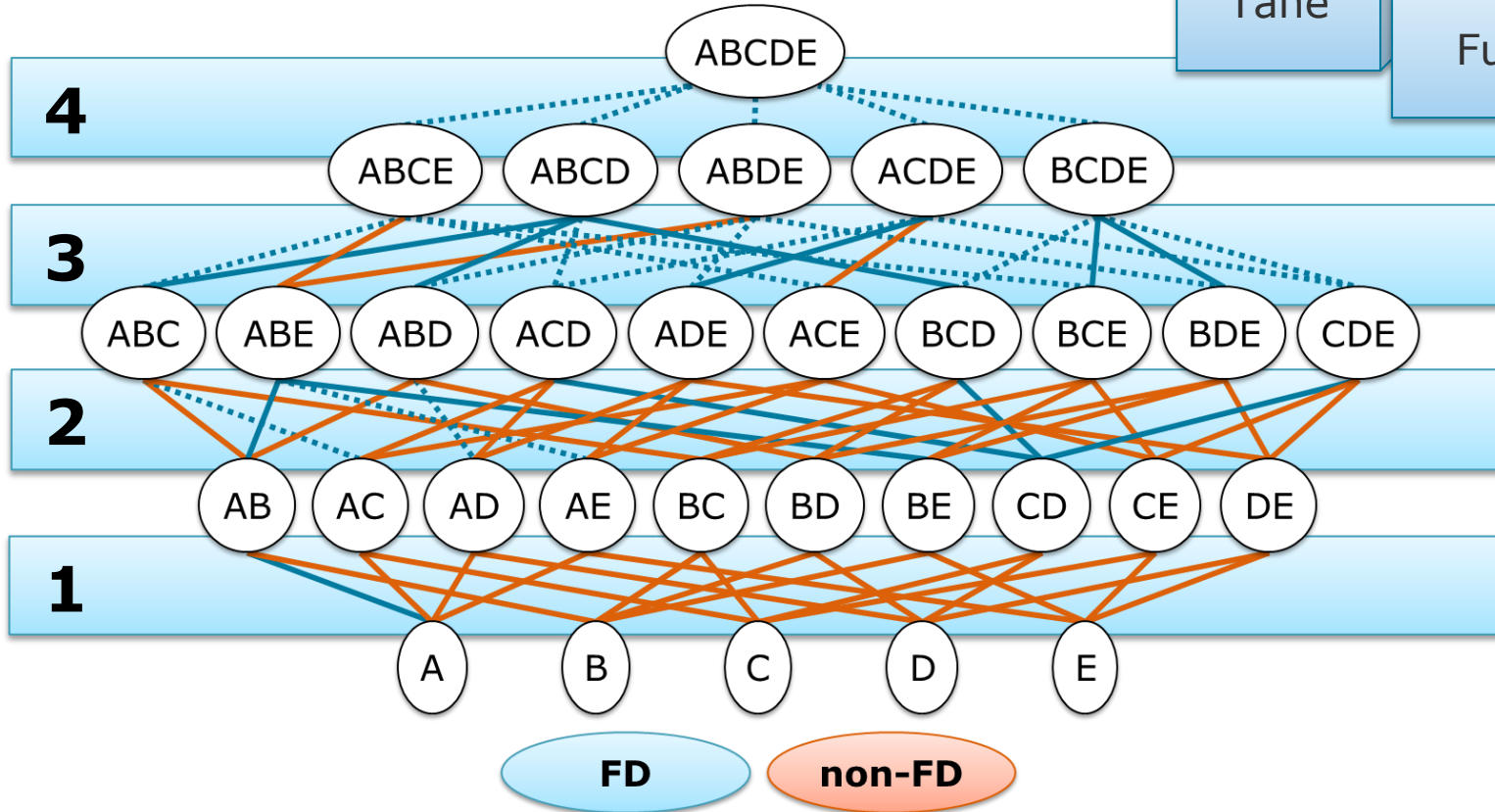
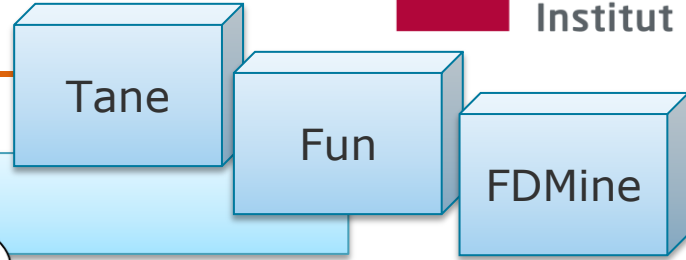


Data Profiling

FD Discovery

ThorstenPapenbrock
Chart 15

Lattice Traversal: Level-Wise



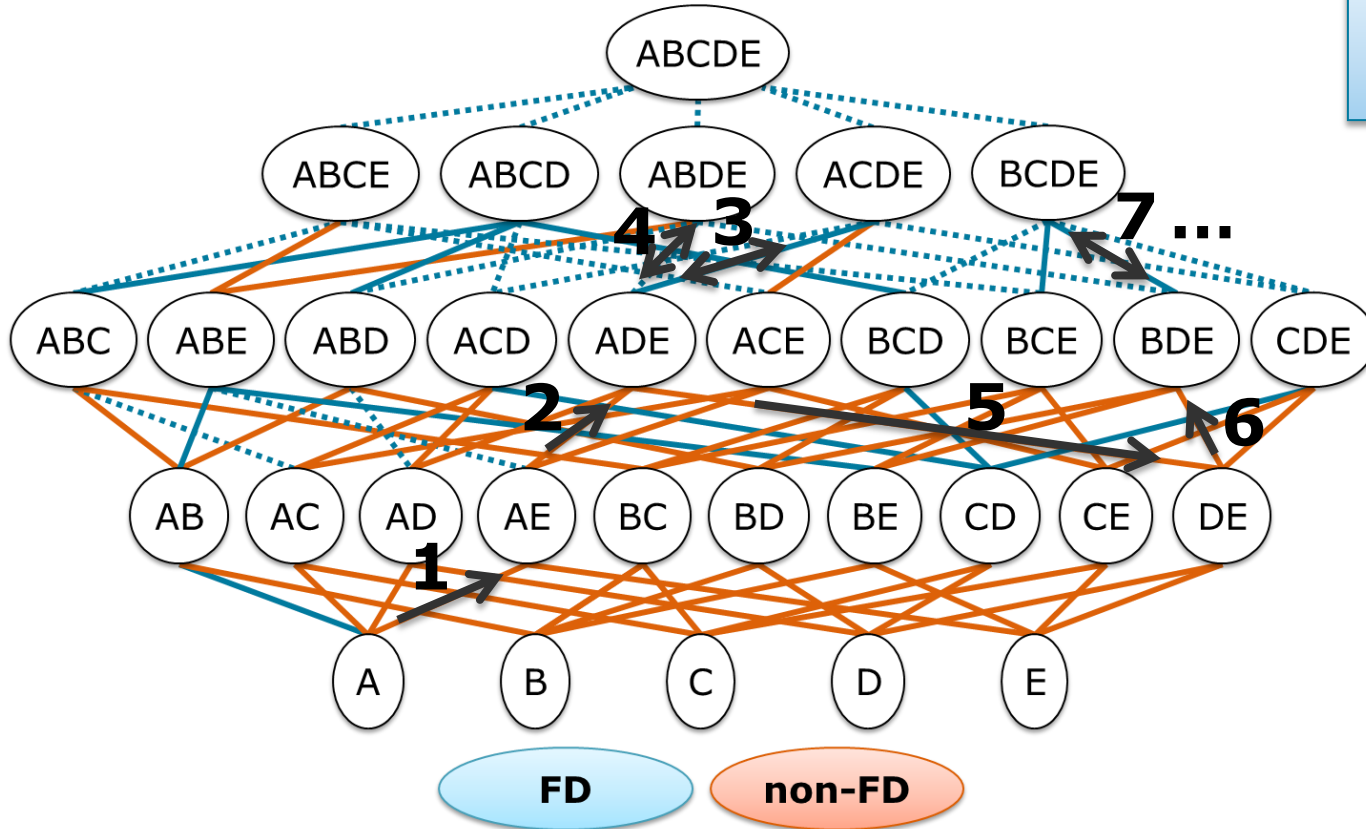
Data Profiling

FD Discovery

ThorstenPapenbrock

Chart 16

Lattice Traversal: Depth-First (Random Walk)



Data Profiling

FD Discovery

ThorstenPapenbrock
Chart 17



Comparison to the level-wise traversal:

Powerful pruning (up- and downwards)



Hole filling and minimization is expensive



Unreliable performance due to randomness

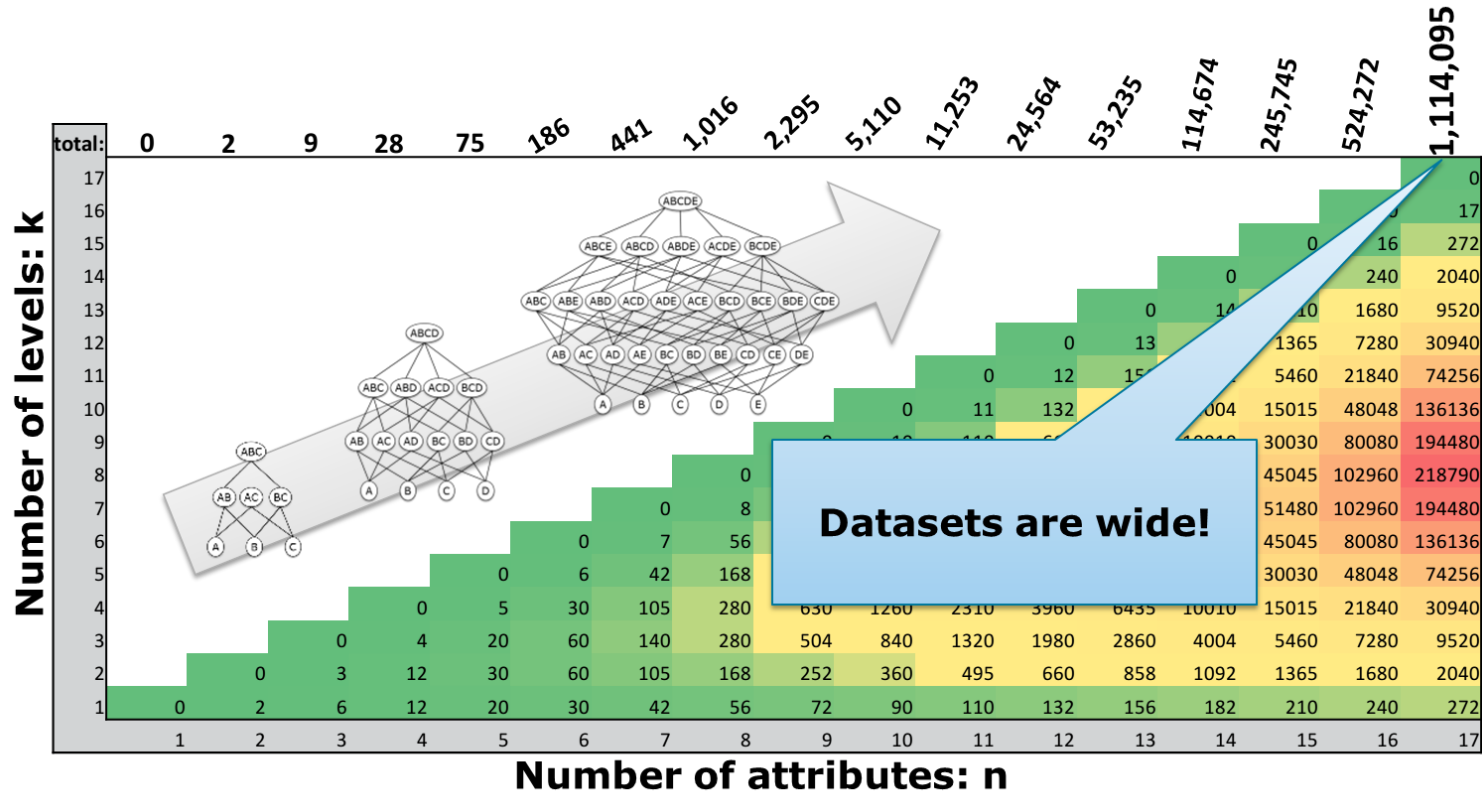


Data Profiling

FD Discovery

Approaches

Lattice Traversal



Data Profiling

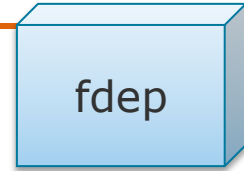
FD Discovery

Thorsten Papenbrock
Chart 19

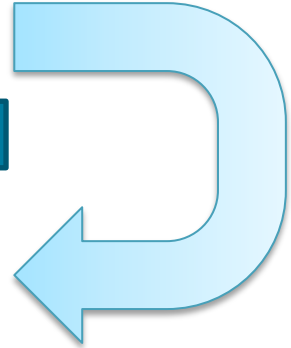
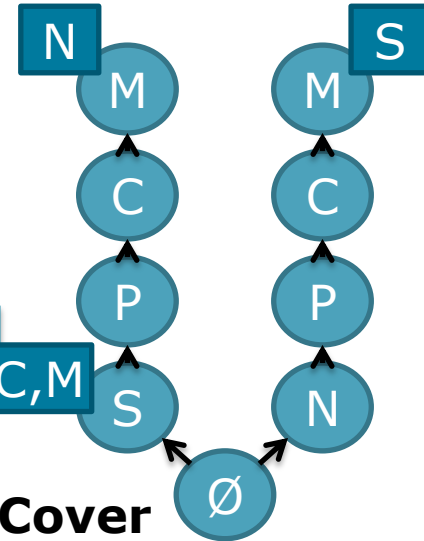
Lattice traversal algorithms: The algorithms TANE [Huhtala et al., 1999], FUN [Novelli and Cicchetti, 2001], FD_MINE [Yao et al., 2002], and DFD [Abedjan et al., 2014c] conceptually arrange all possible FD candidates in a powerset lattice of attribute combinations and then traverse this lattice. The first three algorithms search through the candidate lattice level-wise bottom-up using the *a priori-gen* candidate generation [Agrawal and Srikant, 1994], whereas DFD applies a depth-first random walk. Lattice traversal algorithms in general make intensive use of pruning rules and their candidate validation is based on position list indexes (see Section 2.4). They have been shown to perform well on long datasets, i.e., datasets with many records, but due to their candidate-driven search strategy, they scale poorly with the number of columns in the input dataset.

Dependency Induction

<u>Name</u>	<u>Surname</u>	<u>Postcode</u>	<u>City</u>	<u>Mayor</u>
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

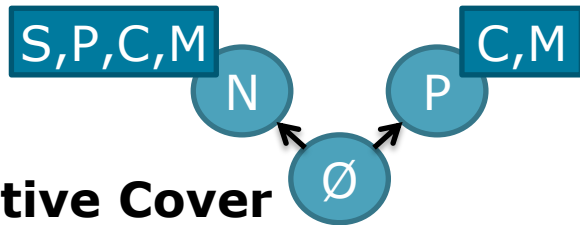


$N, P, C, M \Rightarrow S$



Data Profiling
FD Discovery

Thorsten Papenbrock
Chart 21



Negative Cover

Positive Cover

Approaches Dependency

ncvoter_1024001r_19c.csv - LibreOffice Calc

	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S
26964	924938	9882511	steehn	l	thynpops	558	10	10	108	breitengr pl	cafe	ncv	279311	ncv	01011982	2011-10			
26965	924938	9882512	steehn	l	thynpops	558	10	10	108	breitengr pl	cafe	ncv	279311	ncv	01011982	2011-10			
26966	924938	9882513	steehn	l	thynpops	558	10	10	108	breitengr pl	cafe	ncv	279311	ncv	01011982	2011-10			
26967	924938	9882514	steehn	l	thynpops	558	10	10	108	breitengr pl	cafe	ncv	279311	ncv	01011982	2011-10			
26968	924938	9882515	steehn	l	thynpops	558	10	10	108	breitengr pl	cafe	ncv	279311	ncv	01011982	2011-10			
26969	924938	9882516	steehn	l	thynpops	558	10	10	108	breitengr pl	cafe	ncv	279311	ncv	01011982	2011-10			
26970	924938	9882517	steehn	l	thynpops	558	10	10	108	breitengr pl	cafe	ncv	279311	ncv	01011982	2011-10			
26971	924938	9882518	steehn	l	thynpops	558	10	10	108	breitengr pl	cafe	ncv	279311	ncv	01011982	2011-10			
26972	924938	9882519	steehn	l	thynpops	558	10	10	108	breitengr pl	cafe	ncv	279311	ncv	01011982	2011-10			
26973	924938	9882520	steehn	l	thynpops	558	10	10	108	breitengr pl	cafe	ncv	279311	ncv	01011982	2011-10			
26974	924938	9882521	steehn	l	thynpops	558	10	10	108	breitengr pl	cafe	ncv	279311	ncv	01011982	2011-10			
26975	924938	9882522	steehn	l	thynpops	558	10	10	108	breitengr pl	cafe	ncv	279311	ncv	01011982	2011-10			
26976	924938	9882523	steehn	l	thynpops	558	10	10	108	breitengr pl	cafe	ncv	279311	ncv	01011982	2011-10			
26977	924938	9882524	steehn	l	thynpops	558	10	10	108	breitengr pl	cafe	ncv	279311	ncv	01011982	2011-10			
26978	924938	9882525	steehn	l	thynpops	558	10	10	108	breitengr pl	cafe	ncv	279311	ncv	01011982	2011-10			
26979	924938	9882526	steehn	l	thynpops	558	10	10	108	breitengr pl	cafe	ncv	279311	ncv	01011982	2011-10			
26980	924938	9882527	steehn	l	thynpops	558	10	10	108	breitengr pl	cafe	ncv	279311	ncv	01011982	2011-10			
26981	924938	9882528	steehn	l	thynpops	558	10	10	108	breitengr pl	cafe	ncv	279311	ncv	01011982	2011-10			
26982	924938	9882529	steehn	l	thynpops	558	10	10	108	breitengr pl	cafe	ncv	279311	ncv	01011982	2011-10			
26983	924938	9882530	steehn	l	thynpops	558	10	10	108	breitengr pl	cafe	ncv	279311	ncv	01011982	2011-10			
26984	924938	9882531	steehn	l	thynpops	558	10	10	108	breitengr pl	cafe	ncv	279311	ncv	01011982	2011-10			
26985	924938	9882532	steehn	l	thynpops	558	10	10	108	breitengr pl	cafe	ncv	279311	ncv	01011982	2011-10			
26986	924938	9882533	steehn	l	thynpops	558	10	10	108	breitengr pl	cafe	ncv	279311	ncv	01011982	2011-10			
26987	924938	9882534	steehn	l	thynpops	558	10	10	108	breitengr pl	cafe	ncv	279311	ncv	01011982	2011-10			
26988	924938	9882535	steehn	l	thynpops	558	10	10	108	breitengr pl	cafe	ncv	279311	ncv	01011982	2011-10			
26989	924938	9882536	steehn	l	thynpops	558	10	10	108	breitengr pl	cafe	ncv	279311	ncv	01011982	2011-10			
26990	924938	9882537	steehn	l	thynpops	558	10	10	108	breitengr pl	cafe	ncv	279311	ncv	01011982	2011-10			
26991	924938	9882538	steehn	l	thynpops	558	10	10	108	breitengr pl	cafe	ncv	279311	ncv	01011982	2011-10			
26992	924938	9882539	steehn	l	thynpops	558	10	10	108	breitengr pl	cafe	ncv	279311	ncv	01011982	2011-10			
26993	924938	9882540	steehn	l	thynpops	558	10	10	108	breitengr pl	cafe	ncv	279311	ncv	01011982	2011-10			
26994	924938	9882541	steehn	l	thynpops	558	10	10	108	breitengr pl	cafe	ncv	279311	ncv	01011982	2011-10			
26995	924938	9882542	steehn	l	thynpops	558	10	10	108	breitengr pl	cafe	ncv	279311	ncv	01011982	2011-10			
26996	924938	9882543	steehn	l	thynpops	558	10	10	108	breitengr pl	cafe	ncv	279311	ncv	01011982	2011-10			
26997	924938	9882544	steehn	l	thynpops	558	10	10	108	breitengr pl	cafe	ncv	279311	ncv	01011982	2011-10			
26998	924938	9882545	steehn	l	thynpops	558	10	10	108	breitengr pl	cafe	ncv	279311	ncv	01011982	2011-10			
26999	924938	9882546	steehn	l	thynpops	558	10	10	108	breitengr pl	cafe	ncv	279311	ncv	01011982	2011-10			
27000	924938	9882547	steehn	l	thynpops	558	10	10	108	breitengr pl	cafe	ncv	279311	ncv	01011982	2011-10			
27001	924938	9882548	steehn	l	thynpops	558	10	10	108	breitengr pl	cafe	ncv	279311	ncv	01011982	2011-10			
27002	924938	9882549	steehn	l	thynpops	558	10	10	108	breitengr pl	cafe	ncv	279311	ncv	01011982	2011-10			
27003	924938	9882550	steehn	l	thynpops	558	10	10	108	breitengr pl	cafe	ncv	279311	ncv	01011982	2011-10			
27004	924938	9882551	steehn	l	thynpops	558	10	10	108	breitengr pl	cafe	ncv	279311	ncv	01011982	2011-10			
27005	924938	9882552	steehn	l	thynpops	558	10	10	108	breitengr pl	cafe	ncv	279311	ncv	01011982	2011-10			

Datasets are long!

Data Profiling

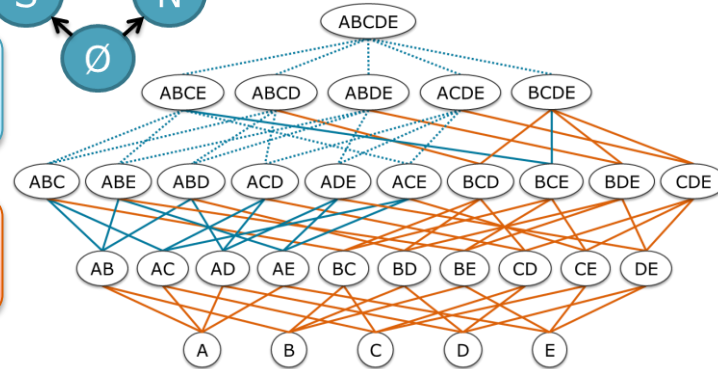
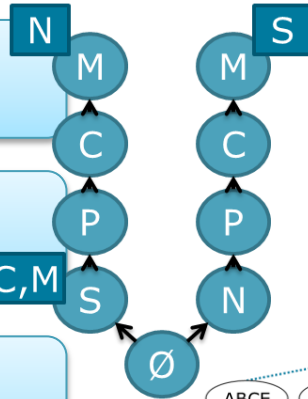
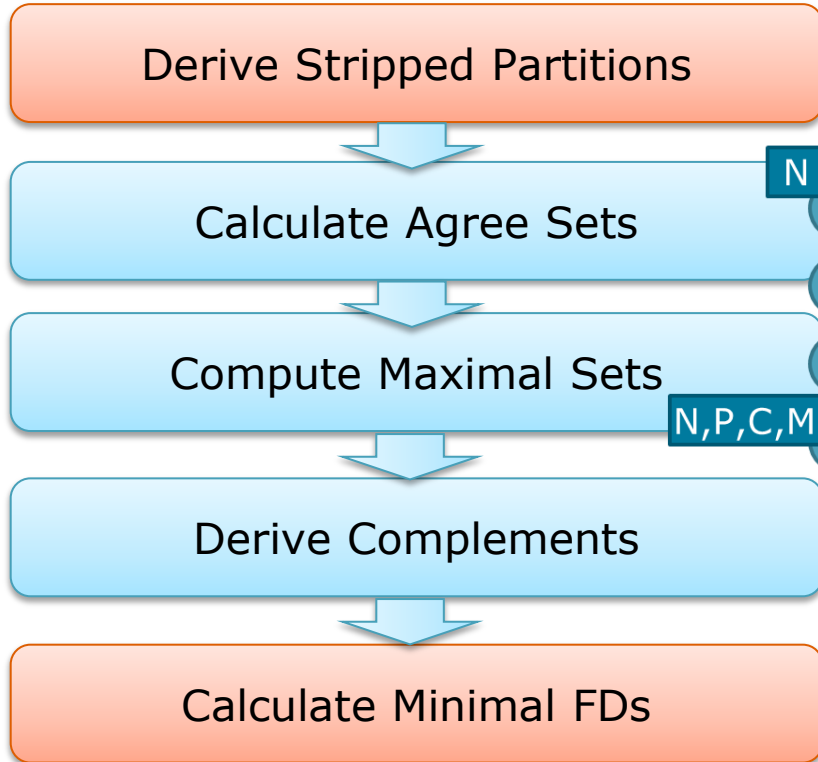
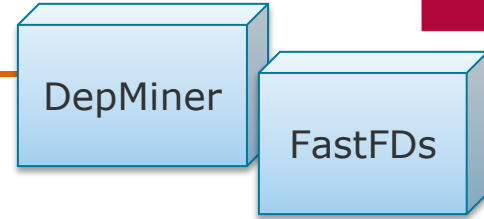
FD Discovery

ThorstenPapenbrock
Chart 22

Row 92694

Dependency induction algorithms: The FDEP [Flach and Savnik, 1999] algorithm also compares all records pair-wise to find all *invalid* functional dependencies. This set is called *negative cover* and is stored in a prefix tree. In contrast to DEP-MINER and FAST-FDS, FDEP translates this negative cover into the set of valid functional dependencies, i.e., the *positive cover*, not by forming complements but by successive specialization: The positive cover initially assumes that each attribute functionally determines all other attributes; these functional dependencies are then refined with every single non-FD in the negative cover. Apart from the fact that the pair-wise comparisons do not scale with the number of records in the input dataset, this discovery strategy has proven to scale well with the number of attributes.

Difference- and Agree-Set



Data Profiling
FD Discovery

ThorstenPapenbrock
Chart 24

Difference- and Agree-Set

Difference- and agree-set algorithms: The algorithms DEP-MINER [Lopes et al., 2000] and FASTFDS [Wyss et al., 2001] analyze a dataset for sets of attributes that agree on the values in certain tuple pairs. These so-called agree-sets are transformed into difference-sets from which all valid FDs can be derived. This discovery strategy scales better with the number of attributes than lattice traversal strategies, because FD candidates are generated only from concrete observations rather than being generated systematically. The required maximization of agree-sets or minimization of difference-sets respectively, however, reduces this advantage significantly. Furthermore, DEP-MINER and FASTFDS scale much worse than the previous algorithms with the number of records, because they need to compare all pairs of records.

Data Profiling

FD Discovery

ThorstenPapenbrock
Chart 25

- Motivation
- Approaches
- **fdep**
 - Peter A Flach and Iztok Sarnik. Database dependency discovery: a machine learning approach. *AI Communications*, 12(3):139–160, 1999.
- HyFD
- Evaluation
- Metanome



Data Profiling

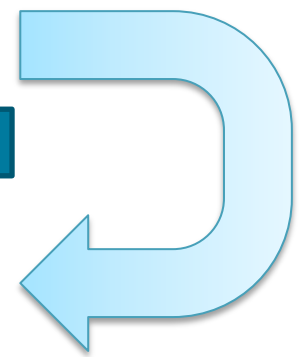
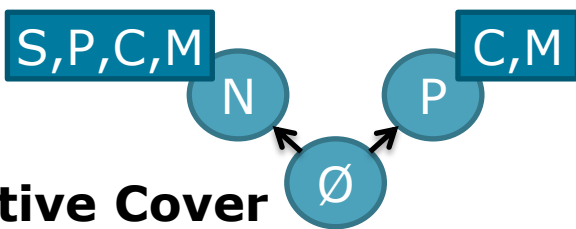
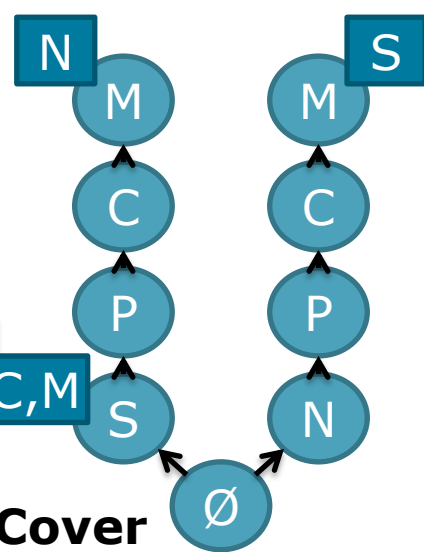
FD Discovery

Thorsten Papenbrock
Chart 26

fdep Algorithm Overview

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

N,P,C,M → S



Data Profiling
FD Discovery
ThorstenPapenbrock
Chart 27

Positive Cover

Negative Cover

fdep Example

<u>Name</u>	<u>Surname</u>	<u>Postcode</u>	<u>City</u>	<u>Mayor</u>
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 \nrightarrow Name
- Name, Postcode, City, Mayor \nrightarrow Surname
- Surname \nrightarrow Name, Postcode, City, Mayor

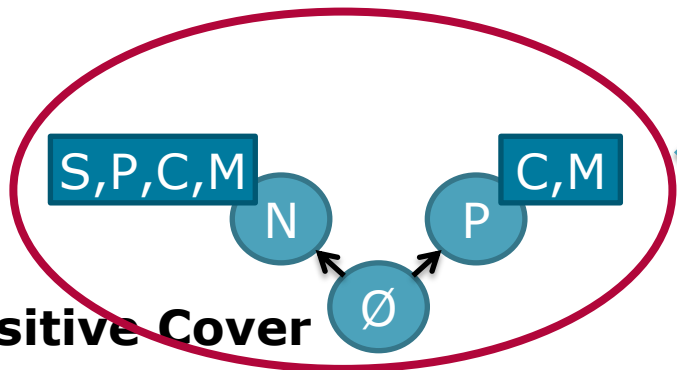


Postcode \rightarrow City
Postcode \rightarrow Mayor
Name \rightarrow Surname, ...

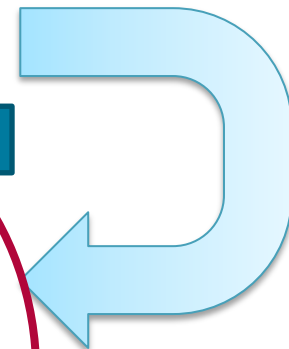
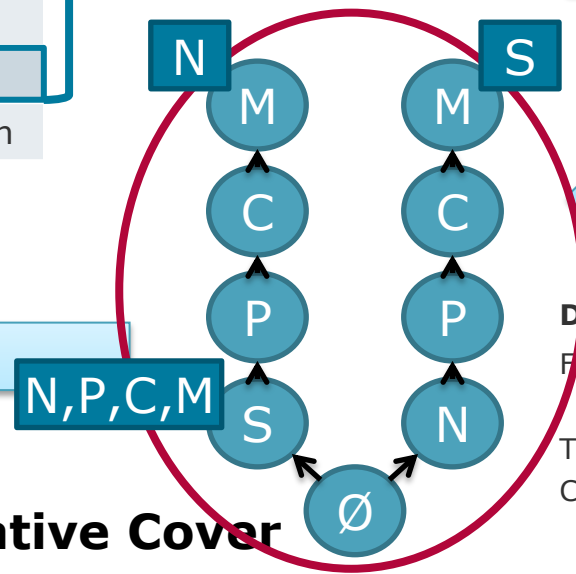
fdep Algorithm Overview

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

$N, P, C, M \rightarrow S$



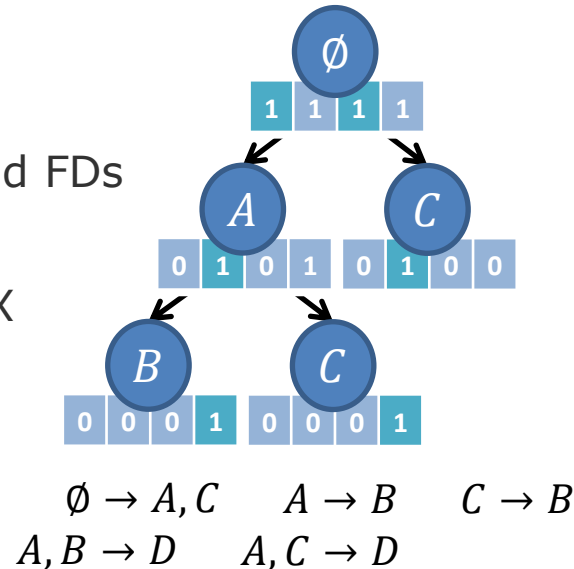
?



Data Profiling
FD Discovery

Thorsten Papenbrock
Chart 29

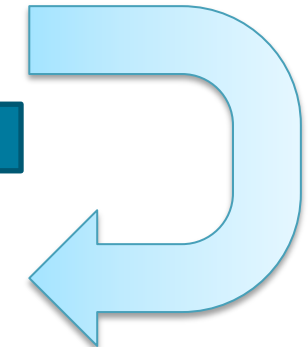
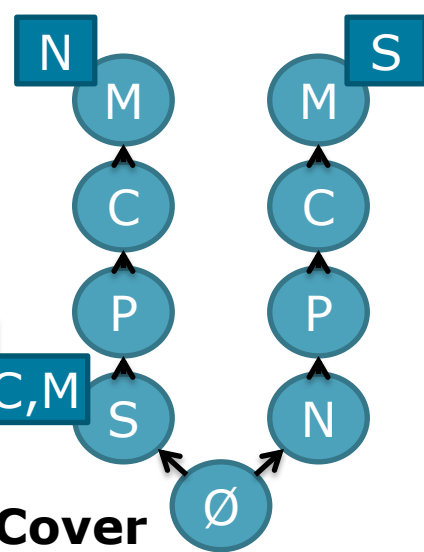
- Definition: An FDTree is a prefix tree (Trie) with annotations:
 - Nodes represent left-hand-side attributes
 - Paths represent left-hand-sides
 - Annotations (bitsets) represent right-hand-side attributes
 - Marked annotations represent valid FDs
- `getGeneralizations($X \rightarrow Y$)`
 - Generalization: $X' \rightarrow Y$ with $X' \subseteq X$
 - Simple depth-first search



fdep Algorithm Overview

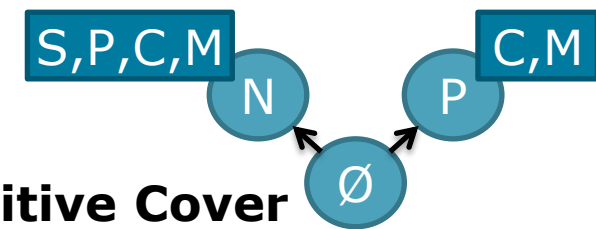
Name	Surname	Postcode	City	Mayor
Thomas	Miller	14482	Potsdam	Jakobs
Sarah	Miller	14482	Potsdam	Jakobs
Peter	Smith	60329	Frankfurt	Feldmann
Jasmine	Cone	01069	Dresden	Orosz
Thomas	Cone	14482	Potsdam	Jakobs
Mike	Moore	60329	Frankfurt	Feldmann

N,P,C,M \Rightarrow S



Data Profiling
FD Discovery

ThorstenPapenbrock
Chart 31



N,P,C,M

Negative Cover

Positive Cover

fdep

Cover Inversion

Algorithm 3: Functional Dependency Induction

Data: *nonFds*

Result: *fds*

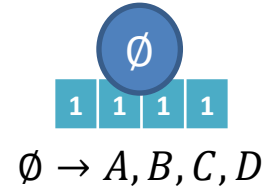
```

1 nonFds ← sort(nonFds, CARDINALITY_DESCENDING);
2 if fds = null then
3   | fds ← new FDTree;
4   | fds.add(∅ → {0, 1, ..., numAttributes});
5 for lhs ∈ nonFds do
6   | rhss ← lhs.clone().flip();
7   | for rhs ∈ rhss do
8     | | specialize(fds, lhs, rhs);
9 return fds;
```

```

function specialize(fds, lhs, rhs)
10 invalidLhss ← fds.getFdAndGenerals(lhs, rhs);
11 for invalidLhss ∈ invalidLhss do
12   | fds.remove(invalidLhss, rhs);
13   | for attr ∈ [0, numAttributes[ do
14     | if invalidLhss.get(attr) ∨
15       |   rhs = attr then
16       |   | continue;
17       |   newLhs ← invalidLhss ∪ attr;
18       |   if fds.findFdOrGeneral(newLhs, rhs) then
19       |   | continue;
20       |   | fds.add(newLhs, rhs);
```

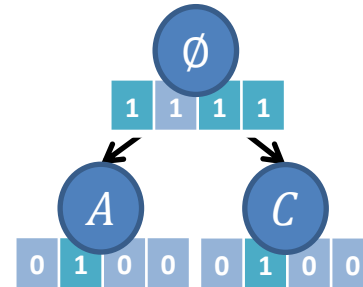
fds:



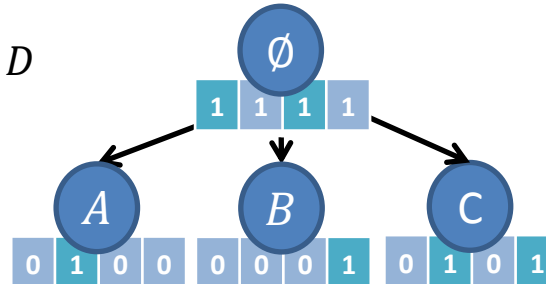
nonFds: $D \not\rightarrow B, \quad A \not\rightarrow D, \quad B \not\rightarrow D, \quad C \not\rightarrow D$

specialize():

$D \not\rightarrow B$



$A \not\rightarrow D$



fdep

Cover Inversion

Algorithm 3: Functional Dependency Induction

Data: $nonFds$

Result: fds

```

1  $nonFds \leftarrow sort(nonFds, CARDINALITY\_DESCENDING);$ 
2 if  $fds = null$  then
3    $fds \leftarrow new\ FDTree;$ 
4    $fds.add(\emptyset \rightarrow \{0, 1, \dots, numAttributes\});$ 
5 for  $lhs \in nonFds$  do
6    $rhss \leftarrow lhs.clone().flip();$ 
7   for  $rhs \in rhss$  do
8      $specialize(fds, lhs, rhs);$ 
9 return  $fds;$ 

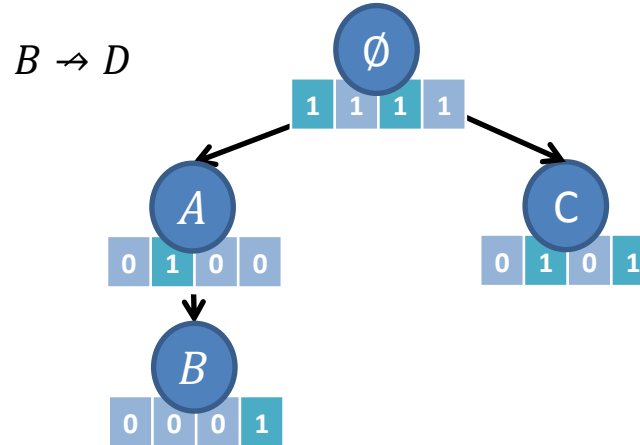
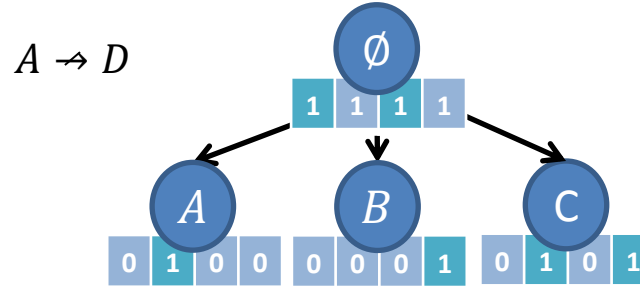
```

```

function  $specialize(fds, lhs, rhs)$ 
10  $invalidLhs \leftarrow fds.getFdAndGenerals(lhs, rhs);$ 
11 for  $invalidLhs \in invalidLhs$  do
12    $fds.remove(invalidLhs, rhs);$ 
13   for  $attr \in [0, numAttributes[$  do
14     if  $invalidLhs.get(attr) \vee$ 
15        $rhs = attr$  then
16        $continue;$ 
17      $newLhs \leftarrow invalidLhs \cup attr;$ 
18     if  $fds.findFdOrGeneral(newLhs, rhs)$  then
19        $continue;$ 
20      $fds.add(newLhs, rhs);$ 

```

$nonFds:$ $D \not\rightarrow B,$ $A \not\rightarrow D,$
 $B \not\rightarrow D,$ $C \not\rightarrow D$



fdep

Cover Inversion

Algorithm 3: Functional Dependency Induction

Data: $nonFds$

Result: fds

```

1  $nonFds \leftarrow sort(nonFds, CARDINALITY\_DESCENDING);$ 
2 if  $fds = null$  then
3    $fds \leftarrow new\ FDTree;$ 
4    $fds.add(\emptyset \rightarrow \{0, 1, \dots, numAttributes\});$ 
5 for  $lhs \in nonFds$  do
6    $rhss \leftarrow lhs.clone().flip();$ 
7   for  $rhs \in rhss$  do
8      $specialize(fds, lhs, rhs);$ 
9 return  $fds;$ 

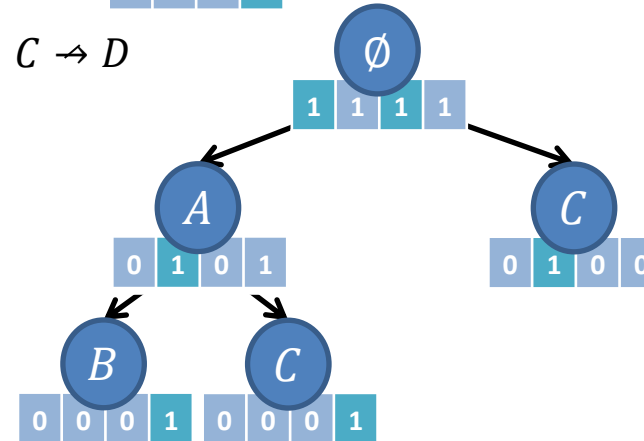
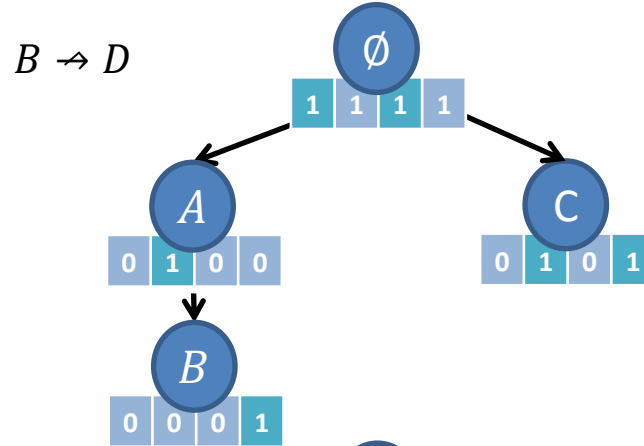
```

```

function  $specialize(fds, lhs, rhs)$ 
10  $invalidLhs \leftarrow fds.getFdAndGenerals(lhs, rhs);$ 
11 for  $invalidLhs \in invalidLhs$  do
12    $fds.remove(invalidLhs, rhs);$ 
13   for  $attr \in [0, numAttributes[$  do
14     if  $invalidLhs.get(attr) \vee$ 
15        $rhs = attr$  then
16        $continue;$ 
17      $newLhs \leftarrow invalidLhs \cup attr;$ 
18     if  $fds.findFdOrGeneral(newLhs, rhs)$  then
19        $continue;$ 
20      $fds.add(newLhs, rhs);$ 

```

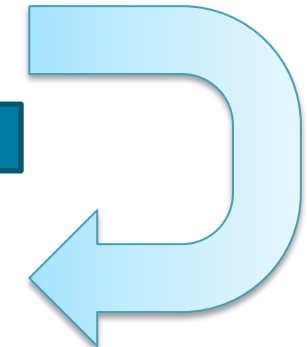
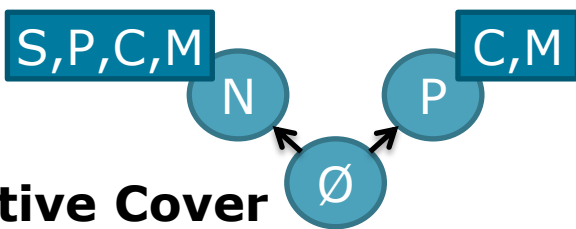
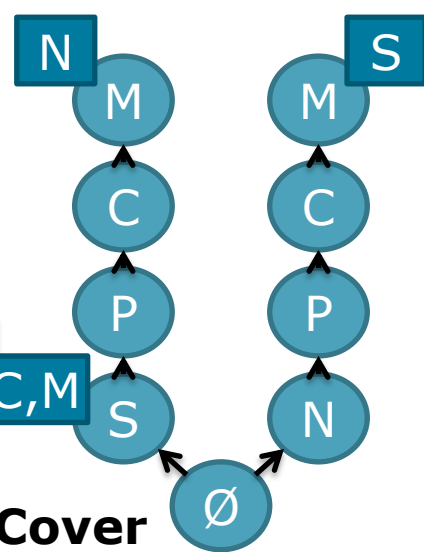
$nonFds:$ $D \rightarrow B, \quad A \rightarrow D,$
 $B \rightarrow D, \quad C \rightarrow D$



fdep Algorithm Overview

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

N,P,C,M \Rightarrow **S**



Data Profiling
FD Discovery

Positive Cover

Negative Cover

- Motivation
- Approaches
- fdep



■ HyFD

Thorsten Papenbrock and Felix Naumann. A Hybrid Approach to Functional Dependency Discovery. In *Proceedings of the International Conference on Management of Data (SIGMOD)*, pages 821–833, 2016.

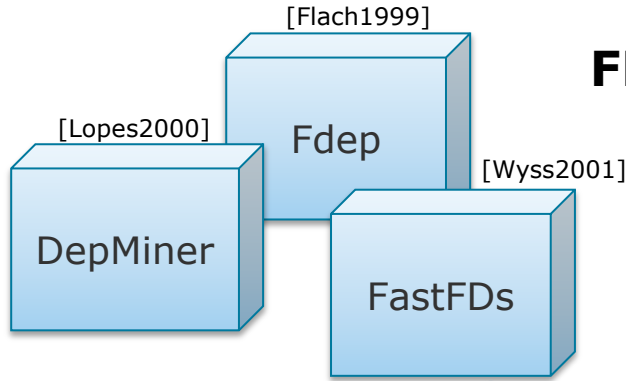
- Evaluation
- Metanome

Data Profiling

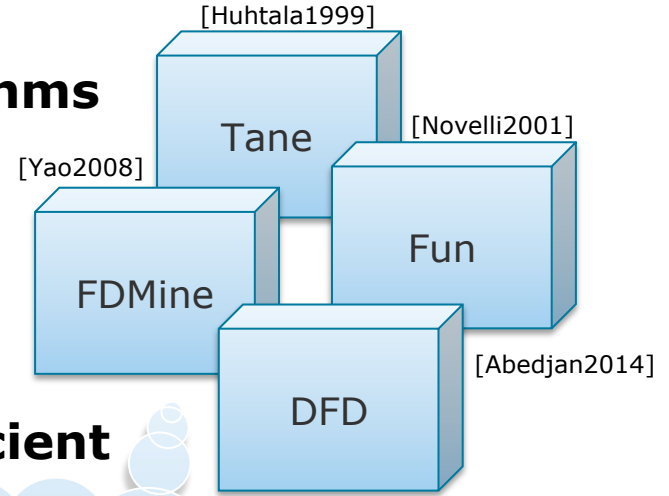
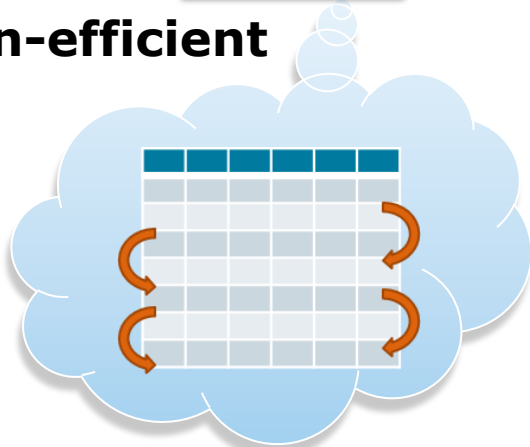
FD Discovery

ThorstenPapenbrock
Chart 36

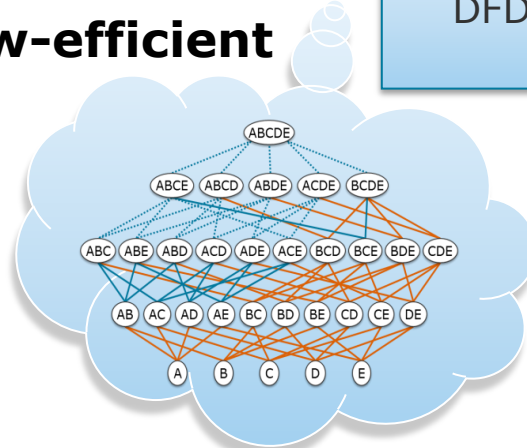
FD Discovery Algorithms



column-efficient



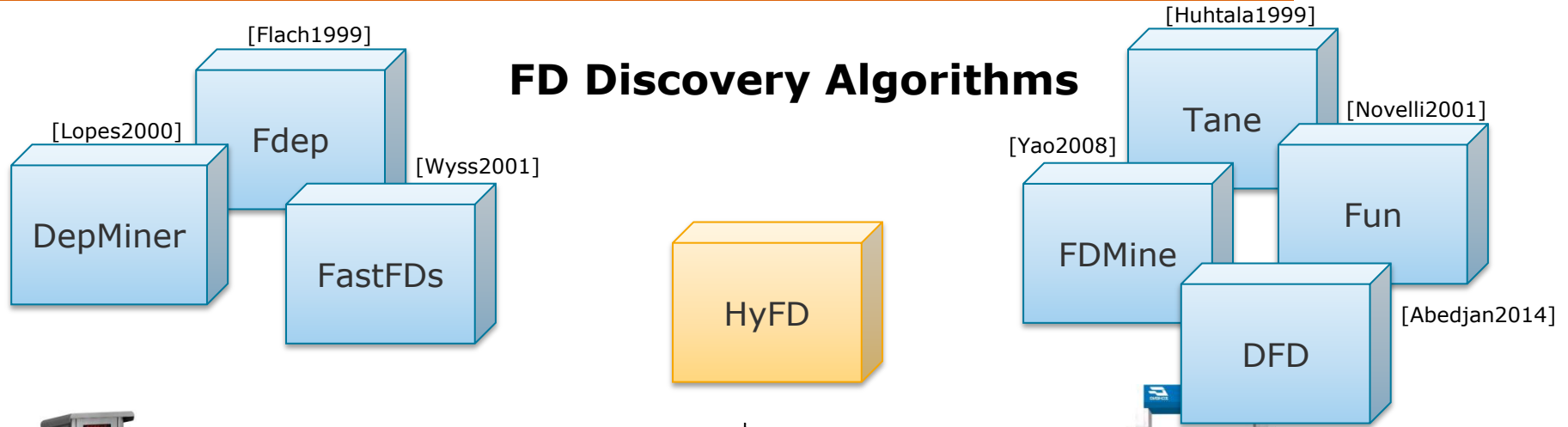
row-efficient



Data Profiling

FD Discovery

ThorstenPapenbrock
Chart 37

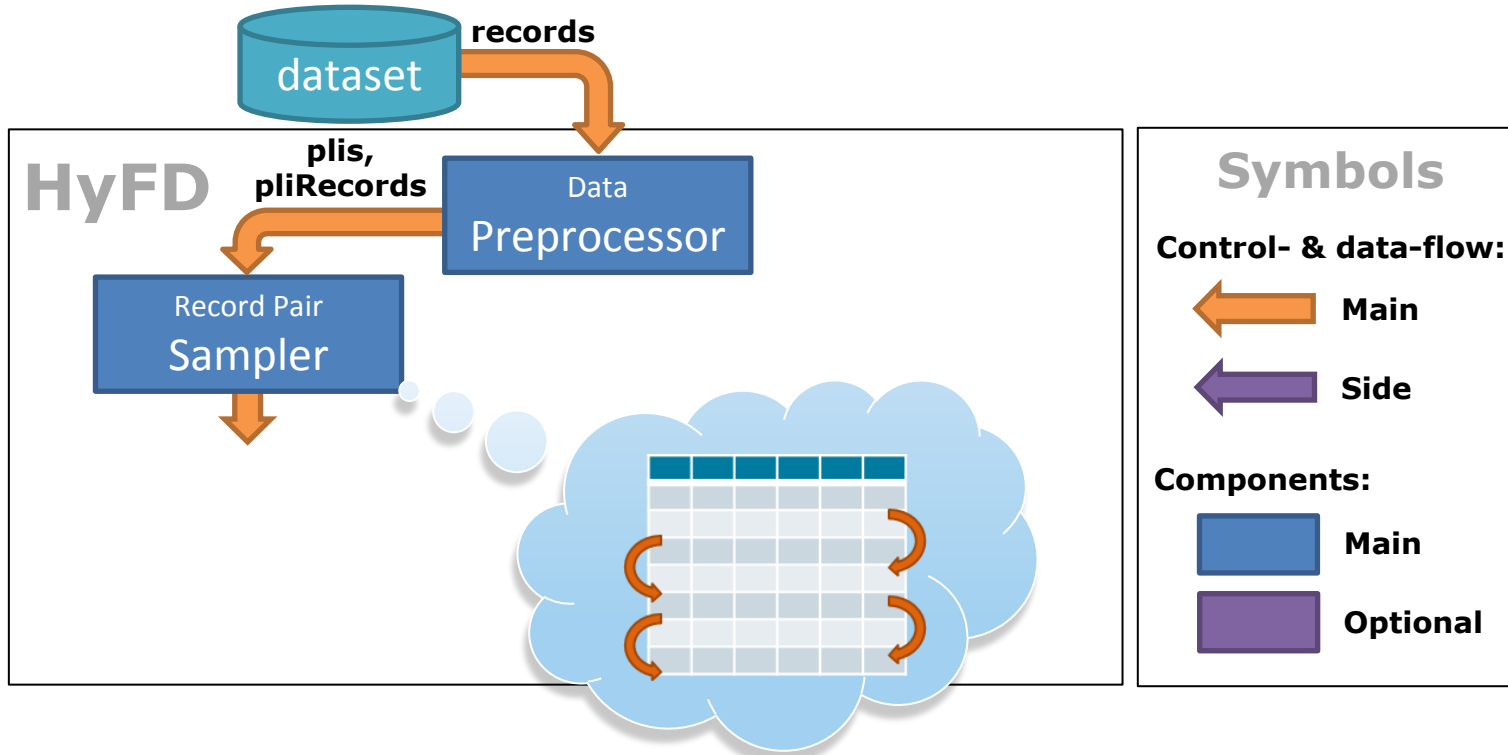


Data Profiling

FD Discovery

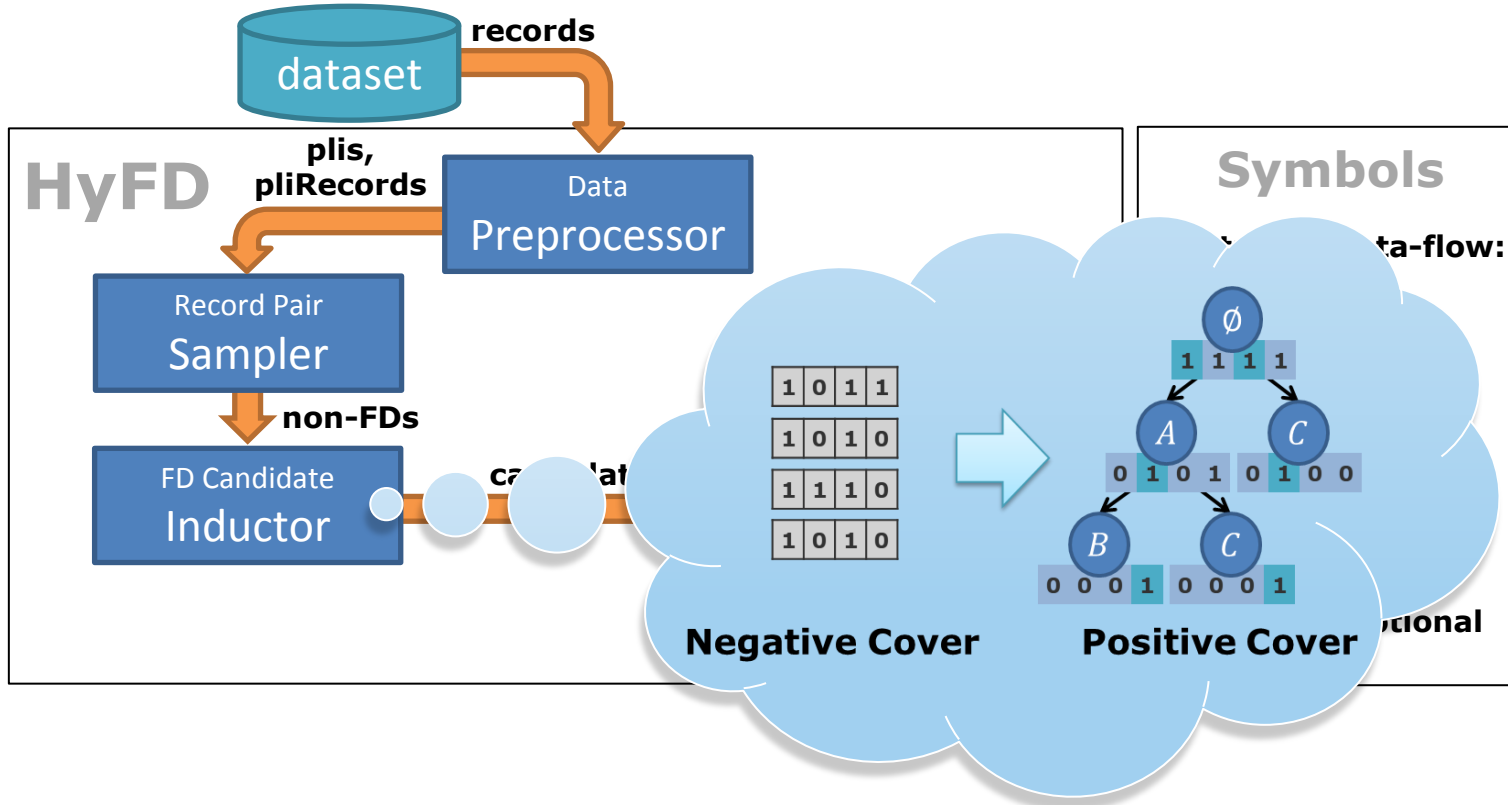
ThorstenPapenbrock
Chart 38

HyFD Hybrid FD Discovery



Data Profiling
FD Discovery

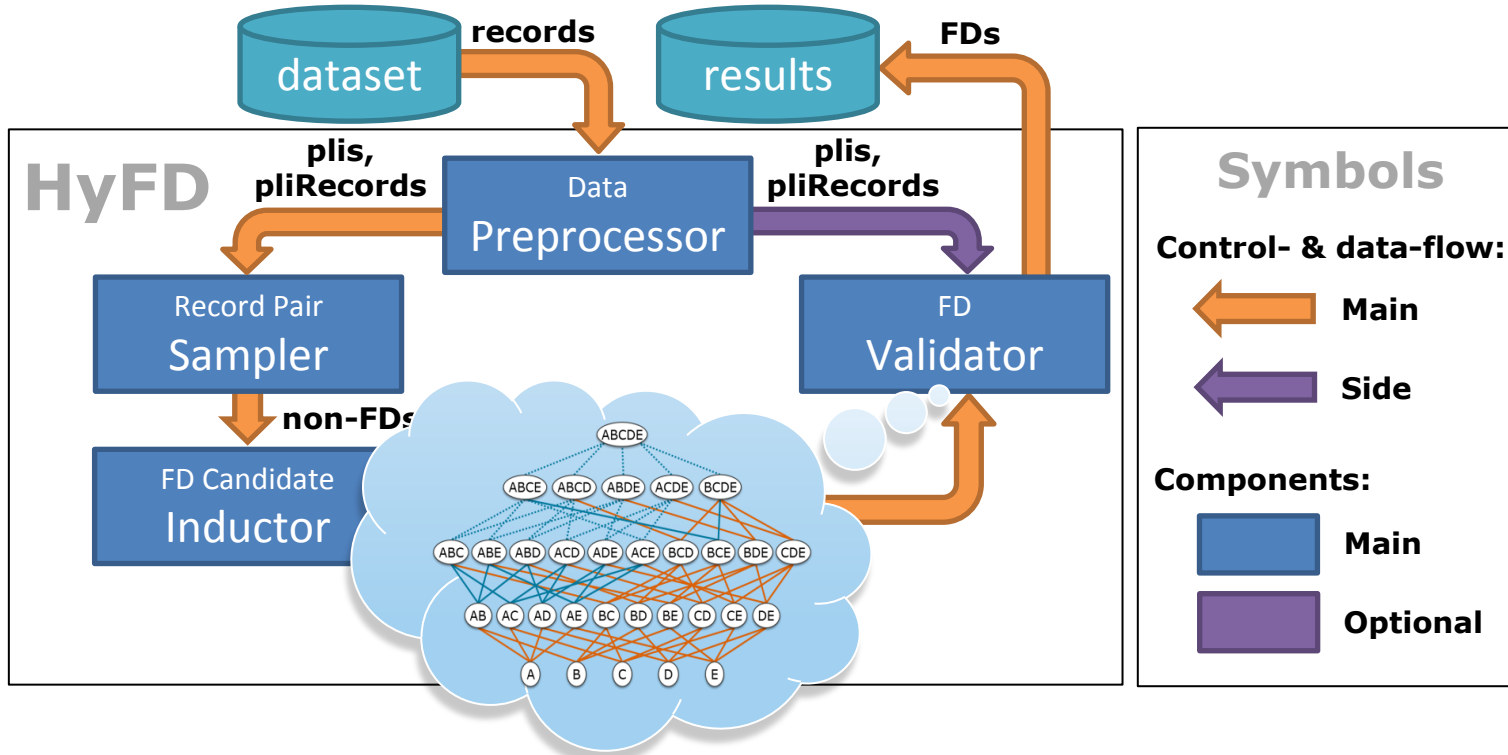
HyFD Hybrid FD Discovery



Data Profiling
FD Discovery

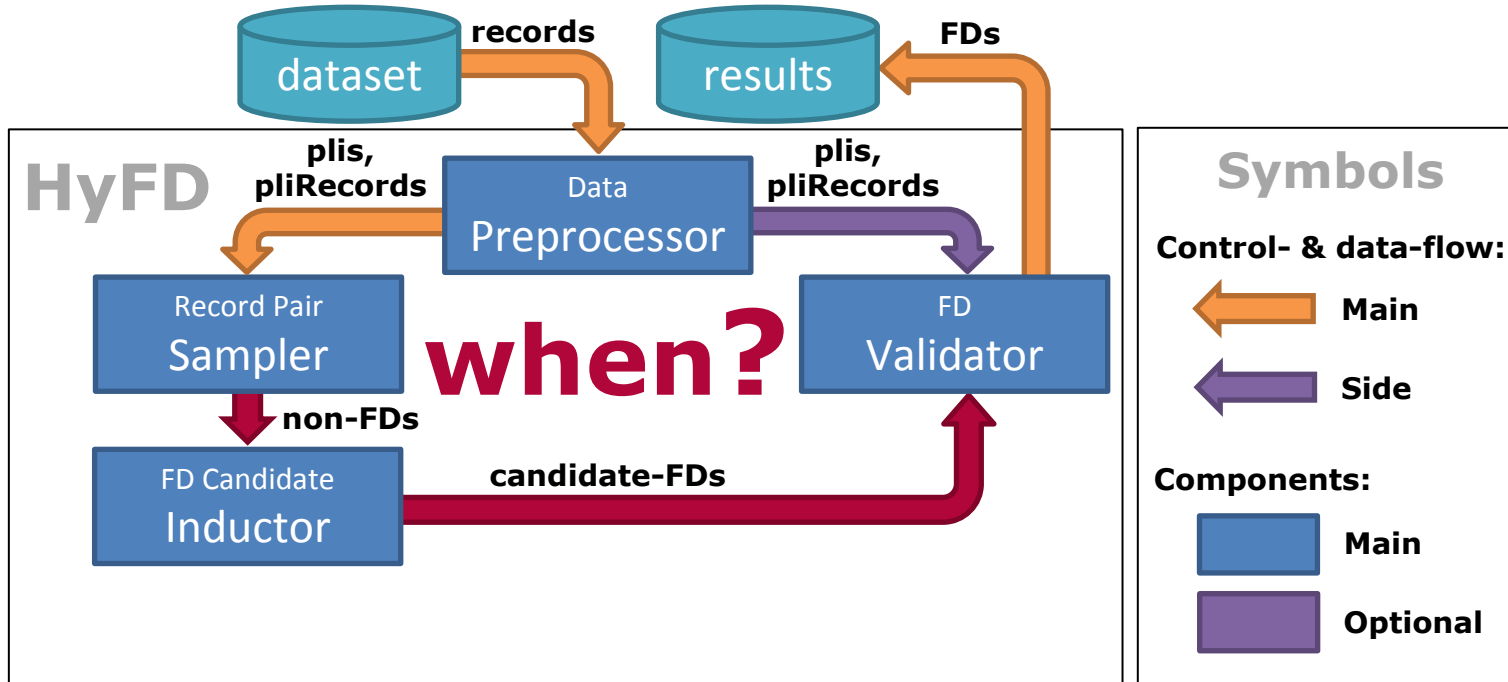
ThorstenPapenbrock
Chart 40

HyFD Hybrid FD Discovery



Data Profiling
FD Discovery

HyFD Hybrid FD Discovery



Data Profiling
FD Discovery

```

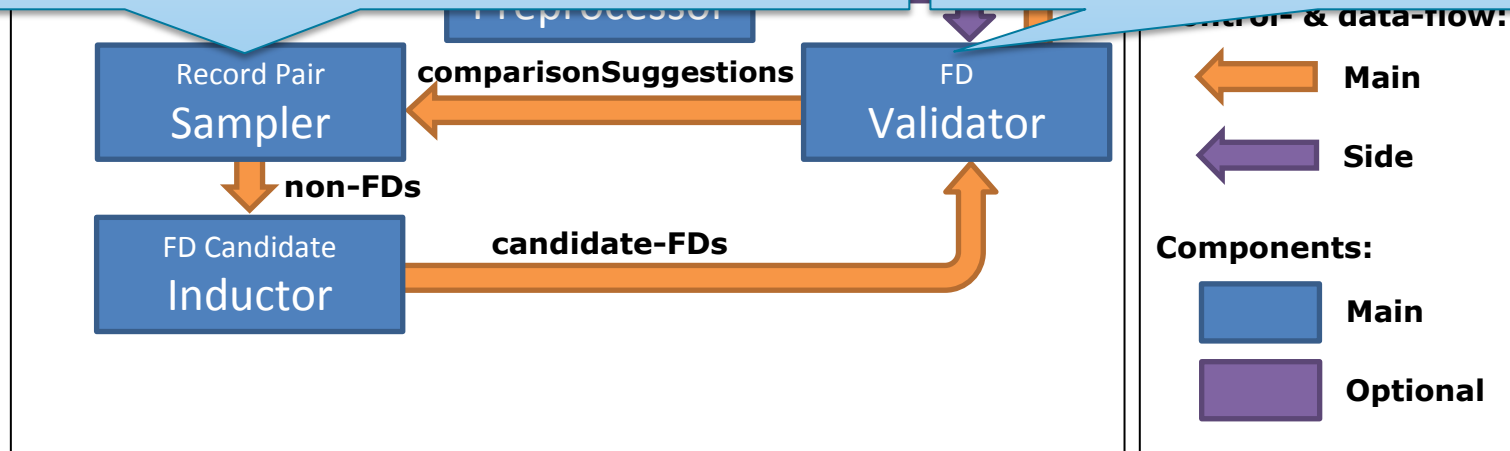
if (  $\frac{\#newFindings}{\#comparisons} < samplingThreshold$  ) {
    samplingThreshold --;
    exit();
}

```

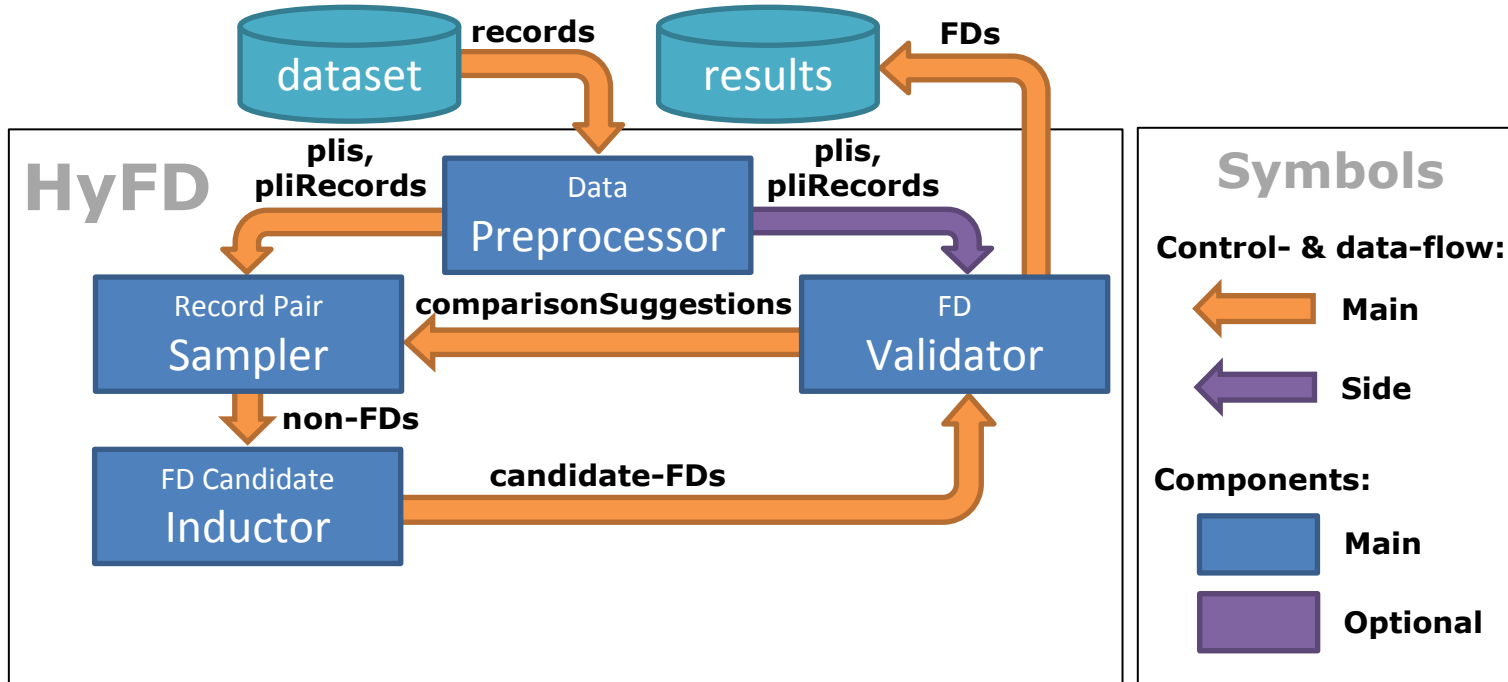
```

if (  $\frac{\#validFDs}{\#validations} < validationThreshold$  ) {
    exit();
}

```

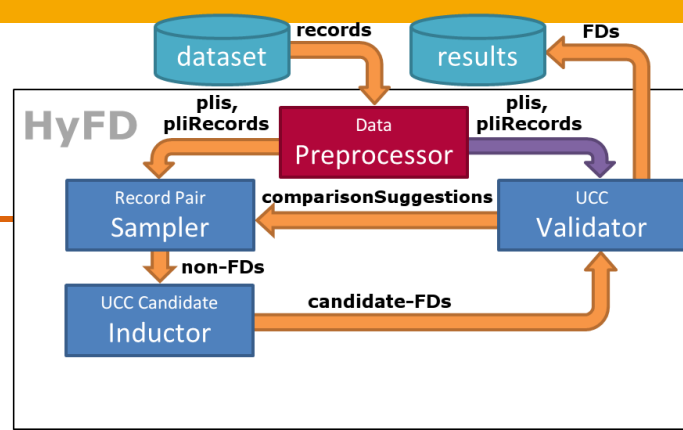


HyFD Hybrid FD Discovery

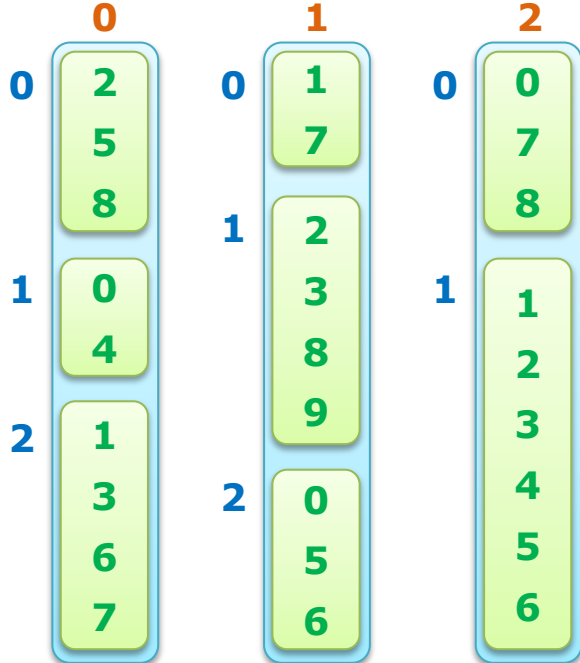


Data Profiling
FD Discovery

HyFD Hybrid FD Discovery



position list indexes
(plis)



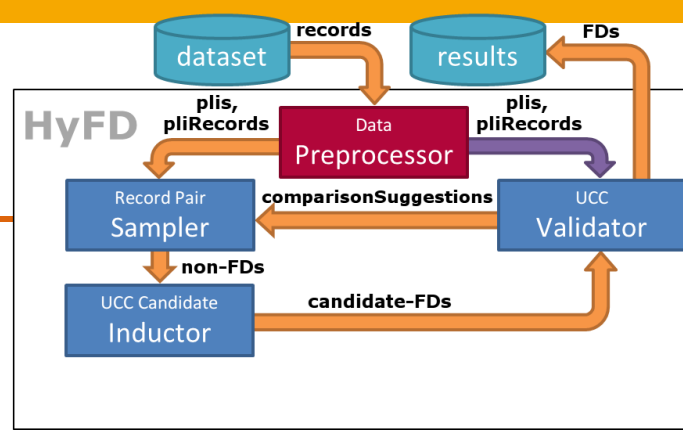
attributeID → clusterID → recordID

Implementation example:
List<List<Integer>> plis;

Data Profiling
FD Discovery

ThorstenPapenbrock
Chart 45

HyFD Hybrid FD Discovery



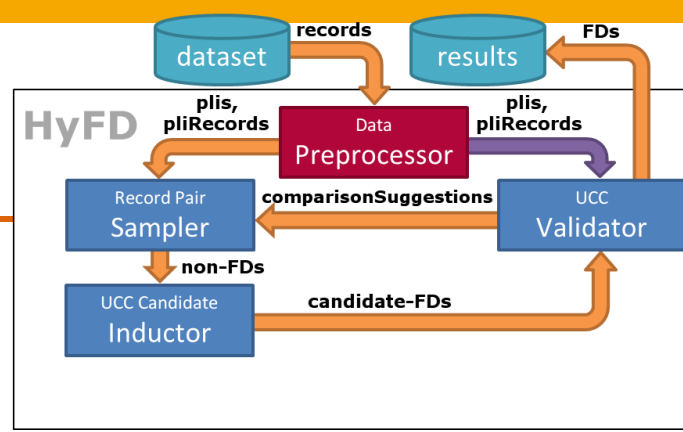
inverted position list indexes
(**invertedPlis**)

	0	1	2
0	1	2	0
1	2	0	1
2	0	1	1
3	2	1	1
4	1	-	1
5	0	2	1

attributeID → recordID → clusterID

Implementation example:
List<Integer[]> invertedPlis;

HyFD Hybrid FD Discovery



position list index records
(**pliRecords**)

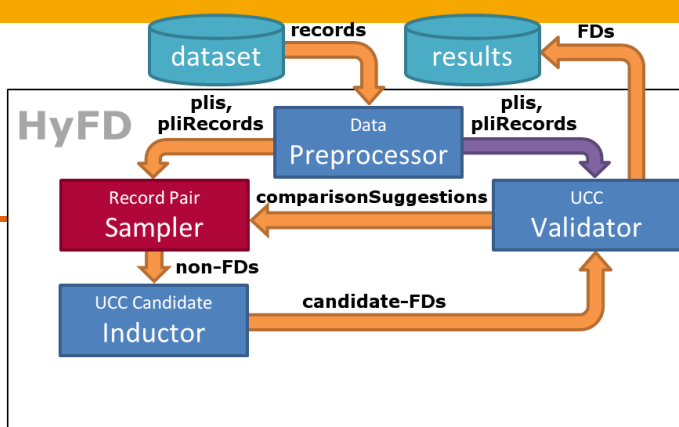
	0	1	2
0	1	2	0
1	2	0	1
2	0	1	1
3	2	1	1
4	1	-	1
5	0	2	1

recordID → attributeID → clusterID

Implementation example:
Integer[][] pliRecords;

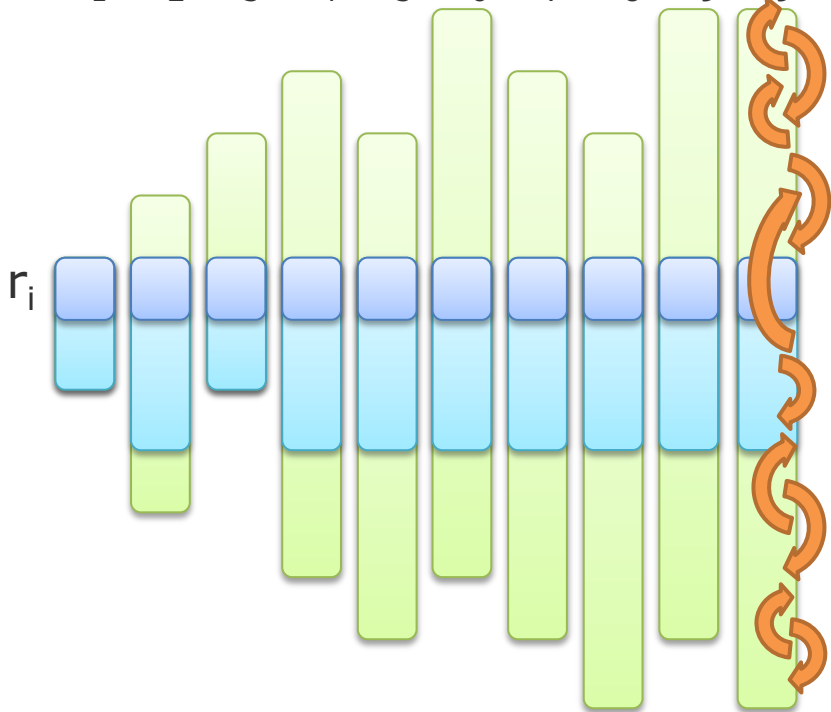
“Dictionary encoded records
using the plis”

HyFD Hybrid FD Discovery



plis:

A_1 A_2 A_3 A_4 A_5 A_6 A_7 A_8 A_9 A_9



DO NOT COMPARE EVERYTHING!

➤ Focused sampling:

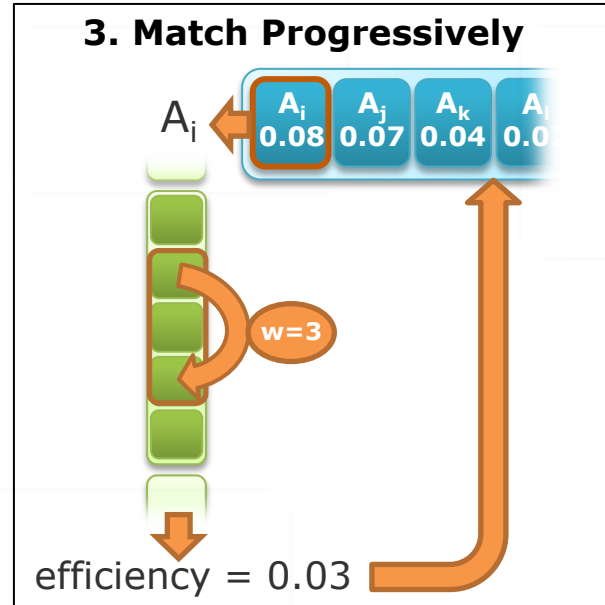
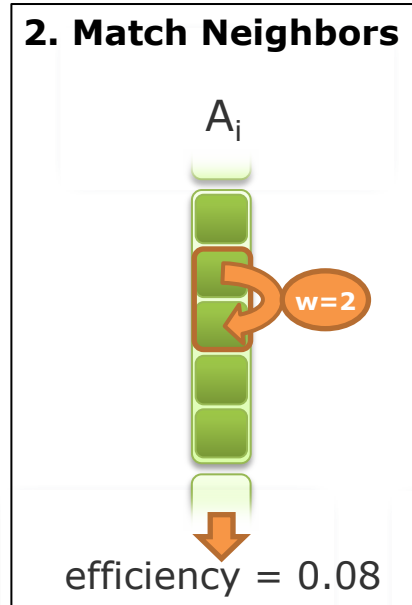
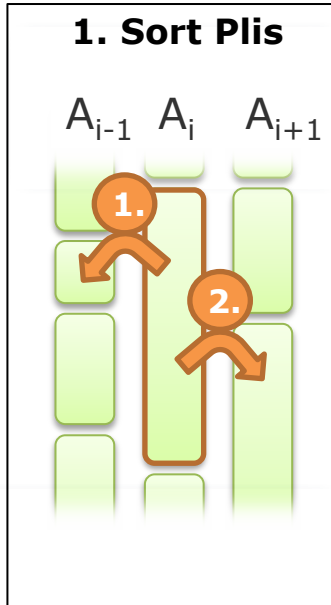
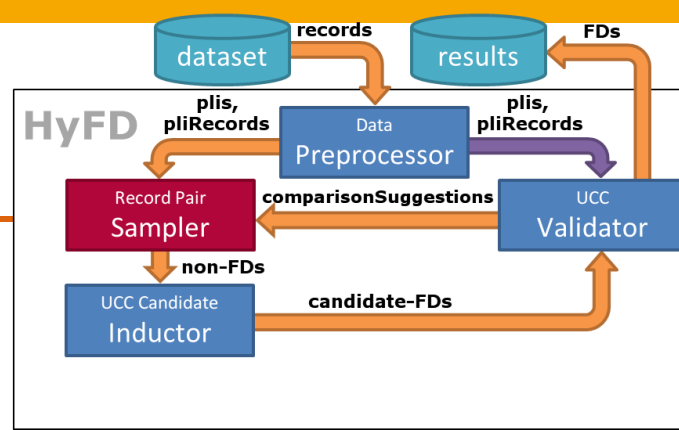
- Only compare records that share at least one cluster (otherwise they cannot violate any FD).
- Try to avoid duplicate comparisons.
- Records that share many clusters prune more effectively.

Data Profiling
FD Discovery

Thorsten Papenbrock
Chart 48

HyFD Hybrid FD Discovery

Progressive, focused sampling:

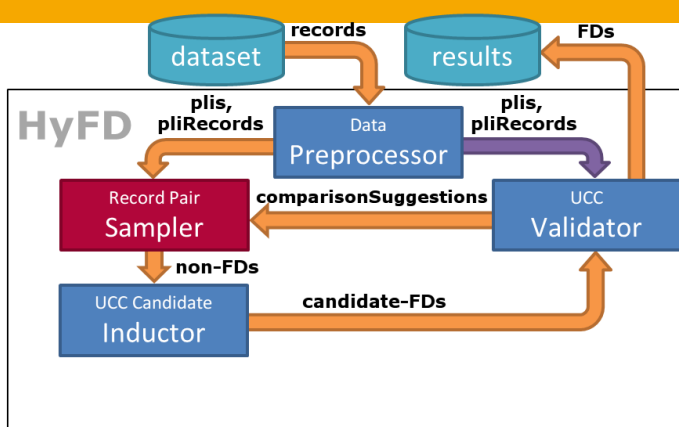
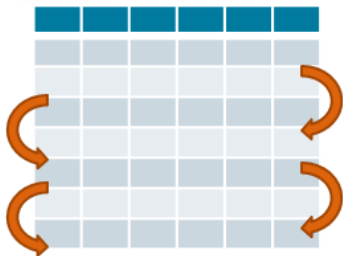
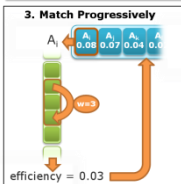
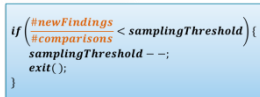
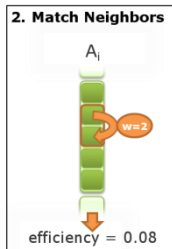
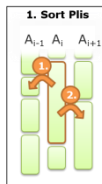


Algorithm 2: Record Pair Sampling

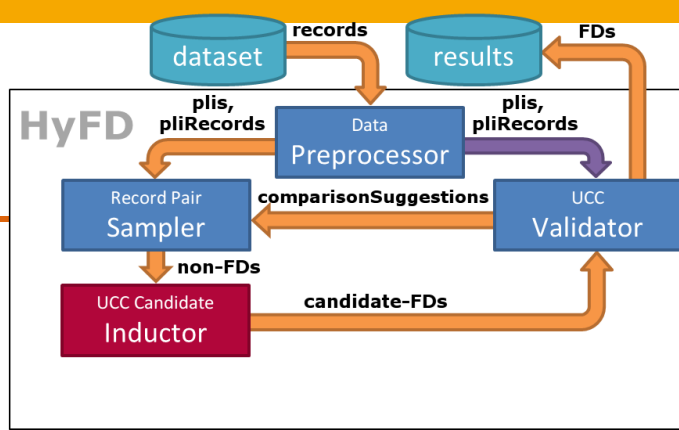
Data: *plis*, *pliRecords*, *comparisonSuggestions*

Result: *nonFds*

```
1 if efficiencyQueue =  $\emptyset$  then
2   for pli  $\in$  plis do
3     for cluster  $\in$  pli do
4       cluster  $\leftarrow$  sort(cluster, ATTR_LEFT_RIGHT);
5
6     nonFds  $\leftarrow$   $\emptyset$ ;
7     efficiencyThreshold  $\leftarrow$  0.01;
8     efficiencyQueue  $\leftarrow$  new PriorityQueue;
9     for attr  $\in$  [0, numAttributes] do
10      efficiency  $\leftarrow$  new Efficiency;
11      efficiency.attribute  $\leftarrow$  attr;
12      efficiency.window  $\leftarrow$  2;
13      efficiency.comps  $\leftarrow$  0;
14      efficiency.results  $\leftarrow$  0;
15      runWindow(efficiency, plis[attr], nonFds);
16      efficiencyQueue.append(efficiency);
17
18 else
19   efficiencyThreshold  $\leftarrow$  efficiencyThreshold / 2;
20   for sug  $\in$  comparisonSuggestions do
21     nonFds  $\leftarrow$  nonFds  $\cup$  match(sug[0], sug[1]);
22
23 while true do
24   bestEff  $\leftarrow$  efficiencyQueue.peak();
25   if bestEff.eval() < efficiencyThreshold then
26     break;
27   bestEff.window  $\leftarrow$  bestEff.window + 1;
28   runWindow(bestEff, plis[bestEff.attribute], nonFds);
29
30 return newFdsIn(nonFds);
31
32 function runWindow(efficiency, pli, nonFds)
33   prevNumNonFds  $\leftarrow$  |nonFds|;
34   for cluster  $\in$  pli do
35     for i  $\in$  [0, |cluster| - efficiency.window] do
36       pivot  $\leftarrow$  pliRecords[cluster[i]];
37       partner  $\leftarrow$  pliRecords[cluster[i + efficiency.window - 1]];
38       nonFds  $\leftarrow$  nonFds  $\cup$  match(pivot, partner);
39       efficiency.comps  $\leftarrow$  efficiency.comps + 1;
40
41   newResults  $\leftarrow$  |nonFds| - prevNumNonFds;
42   efficiency.results  $\leftarrow$  efficiency.results + newResults;
```

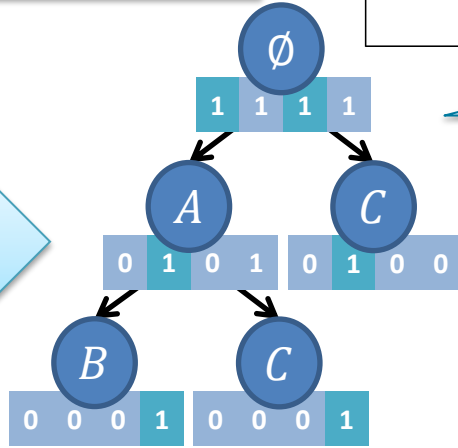
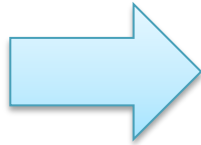


HyFD Hybrid FD Discovery



Bitmaps for negative cover

1	0	1	1
1	0	1	0
1	1	1	0
1	0	1	0



FDTree for positive cover

FDTree serves as candidate lattice
(and final result!)

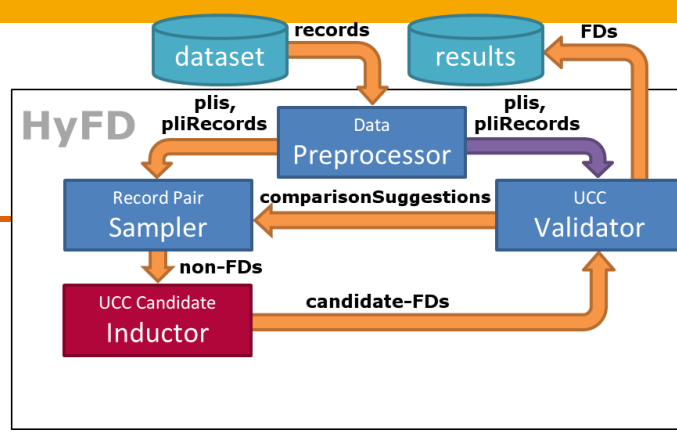
Negative Cover
(**non-FDs**)

Positive Cover
(**candidate-FDs**)

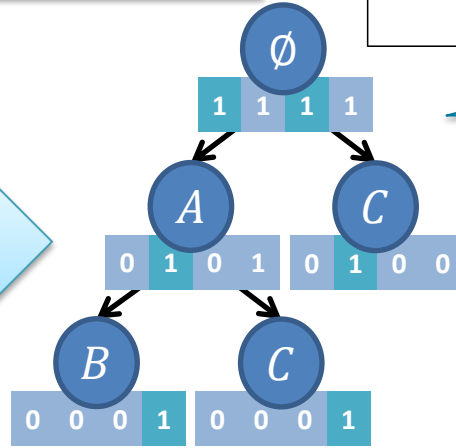
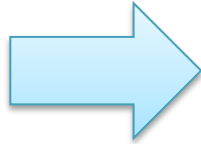
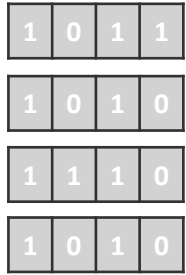
Data Profiling
FD Discovery

ThorstenPapenbrock
Chart **51**

HyFD Hybrid FD Discovery



Bitmaps for negative cover



FDTree for positive cover

Negative Cover
(**non-FDs**)

Positive Cover
(**candidate-FDs**)

Algorithm 3: Functional Dependency Induction

Data: *nonFds*

Result: *fds*

```

1 nonFds ← sort(nonFds, CARDINALITY_DESCENDING);
2 if fds = null then
3   fds ← new FDTree;
4   fds.add(∅ → {0, 1, ..., numAttributes});
5 for lhs ∈ nonFds do
6   rhss ← lhs.clone().flip();

```

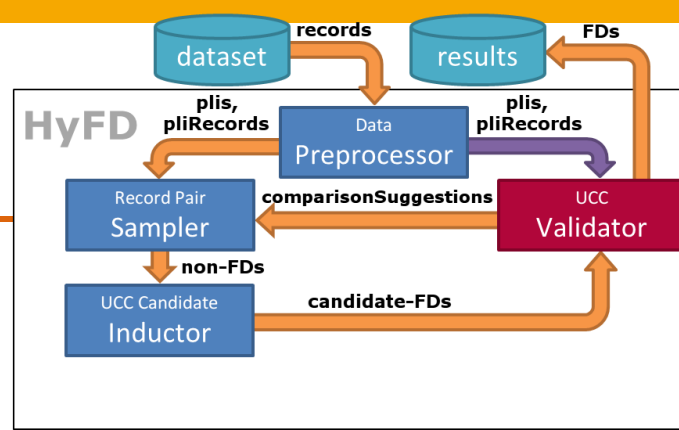
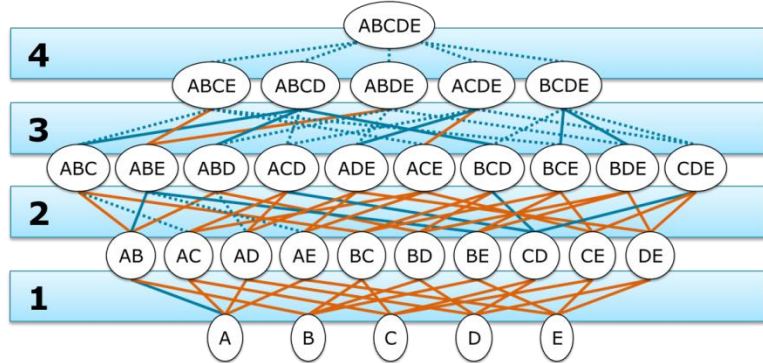
Induction = Cover Inversion in fdep

Use non-FDs with greatest pruning impact first!

HyFD

Hybrid FD Discovery

Bottom-up lattice traversal:



Optimizations:

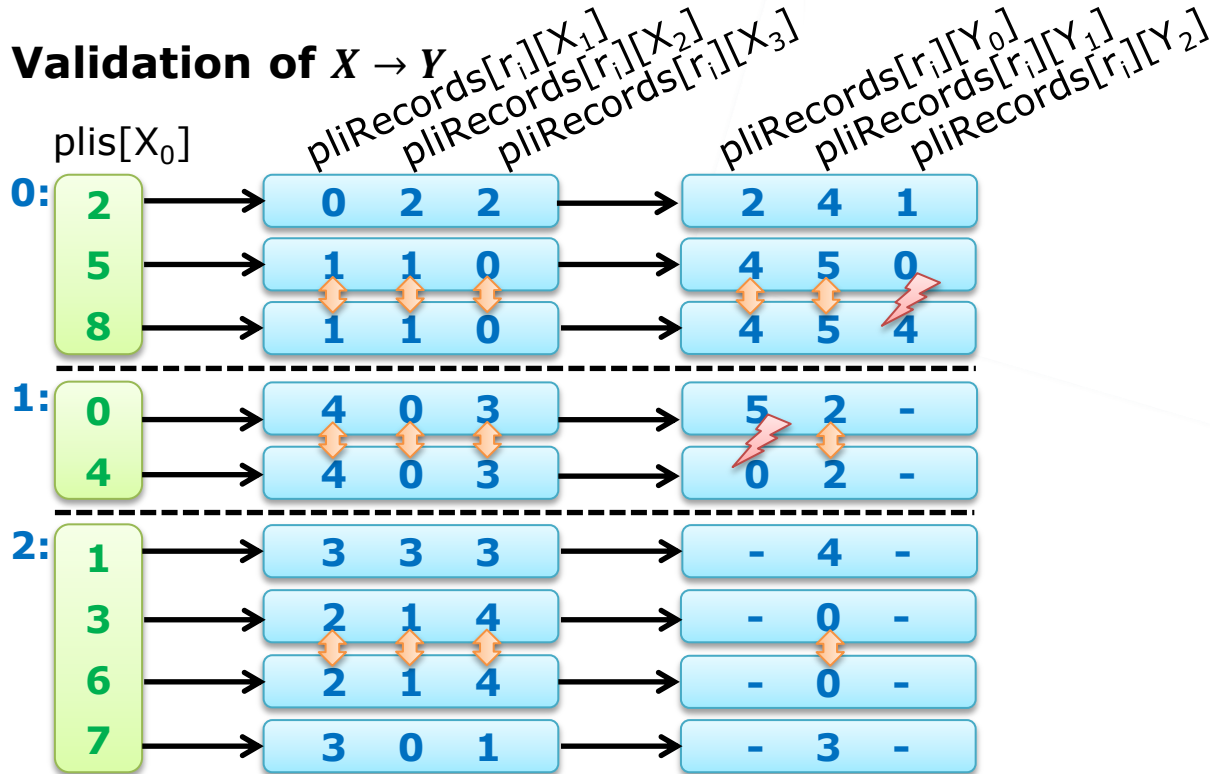
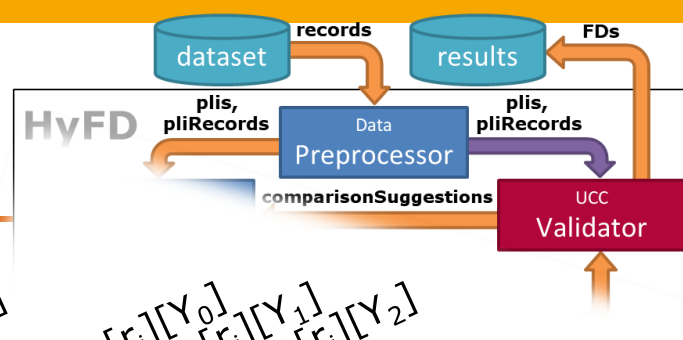
- **Additional Pruning:** Use the results from the sampling.
- **No PLI-caching:** Use only the single column plis and pliRecords for the validations (do not remember plis of higher arity).
- **Early termination:** Stop PLI intersects if candidates are invalid.
- **Bulk validations:** Check all rhss for one given lhs jointly.

Data Profiling

FD Discovery

ThorstenPapenbrock
Chart 53

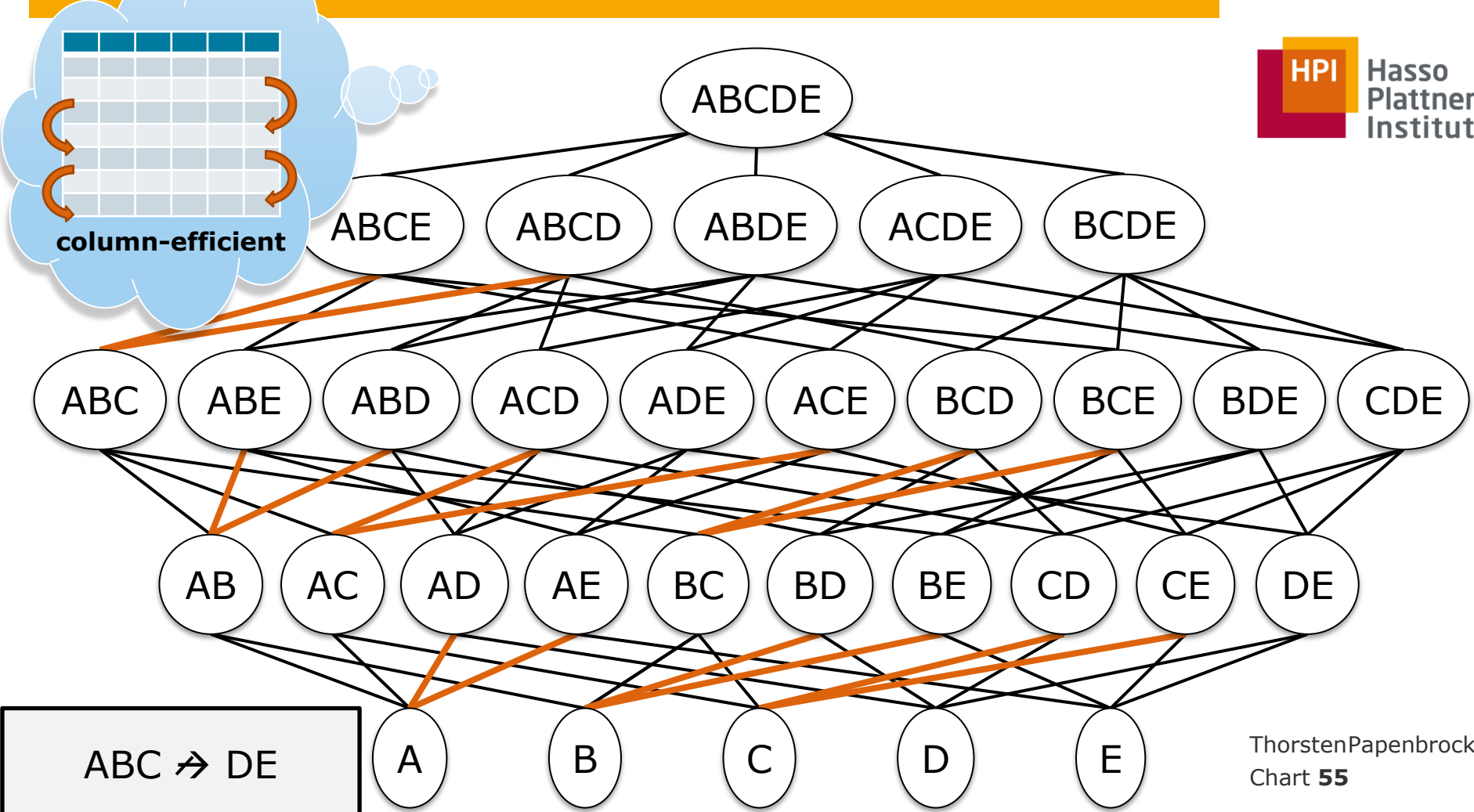
HyFD Hybrid FD Discovery



Data Profiling

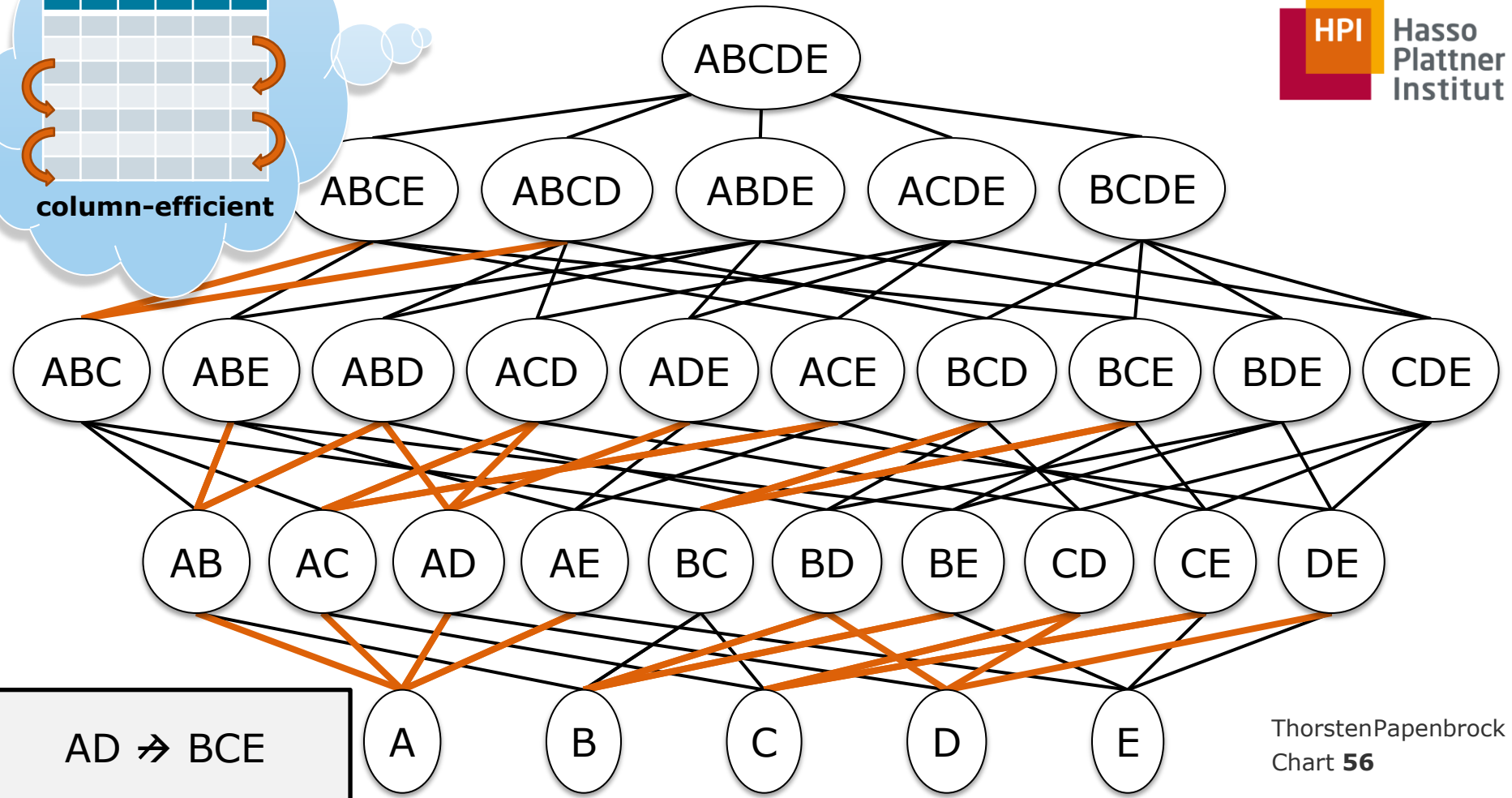
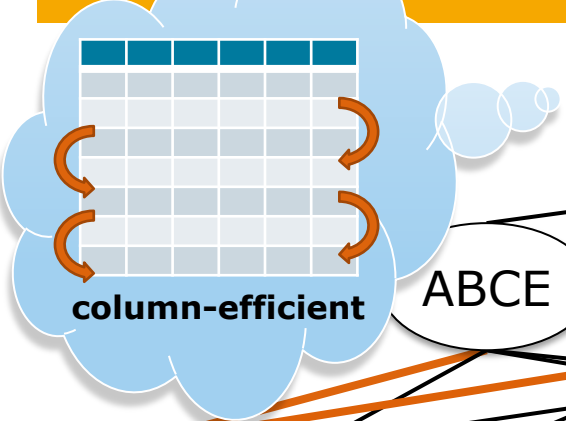
FD Discovery

ThorstenPapenbrock
Chart 54

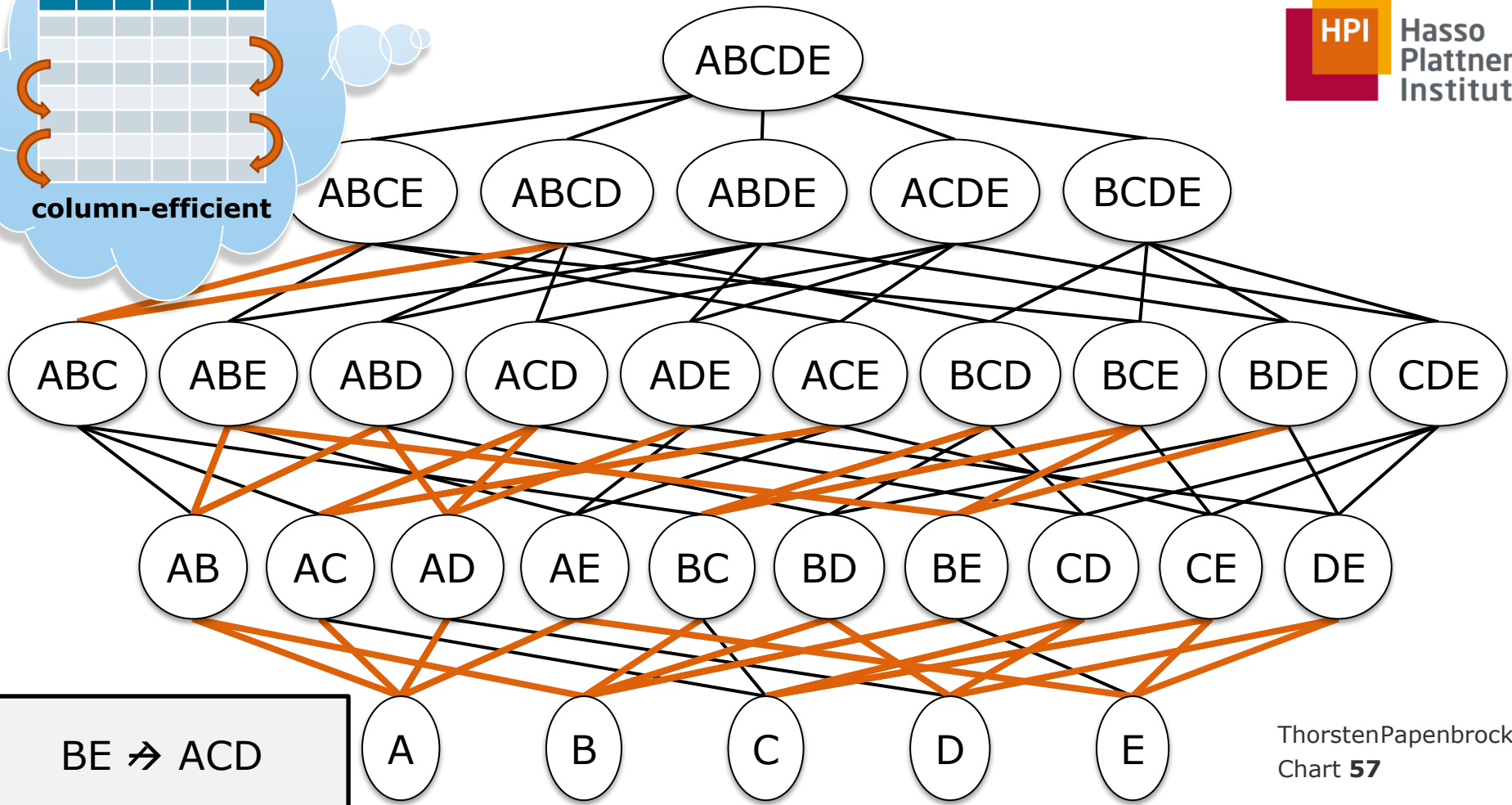
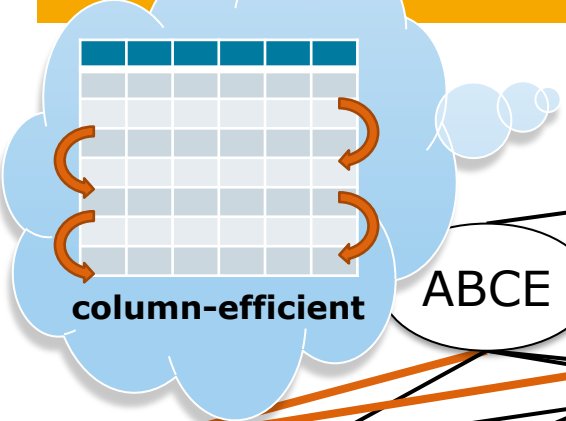


column-efficient

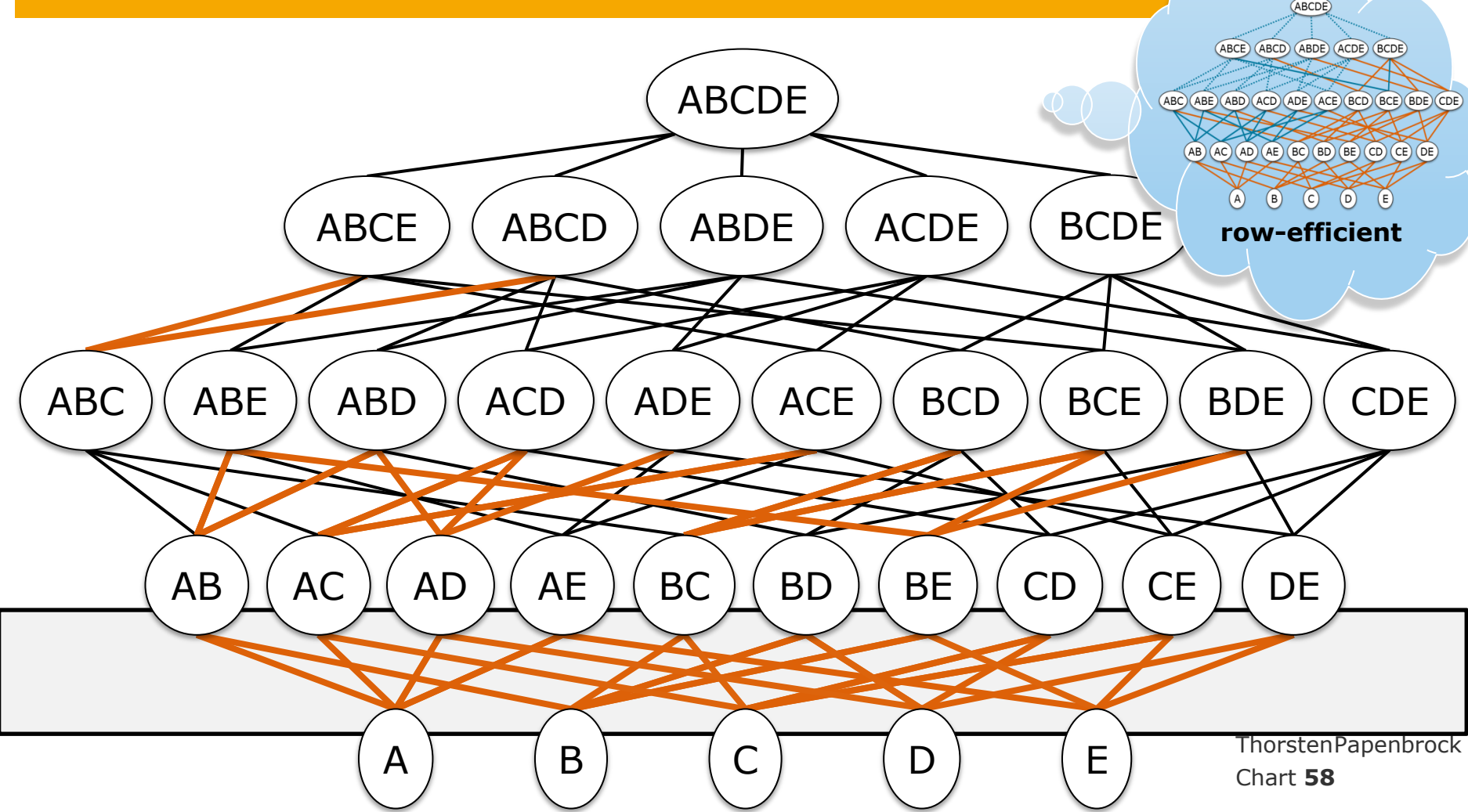
$ABC \rightarrow DE$

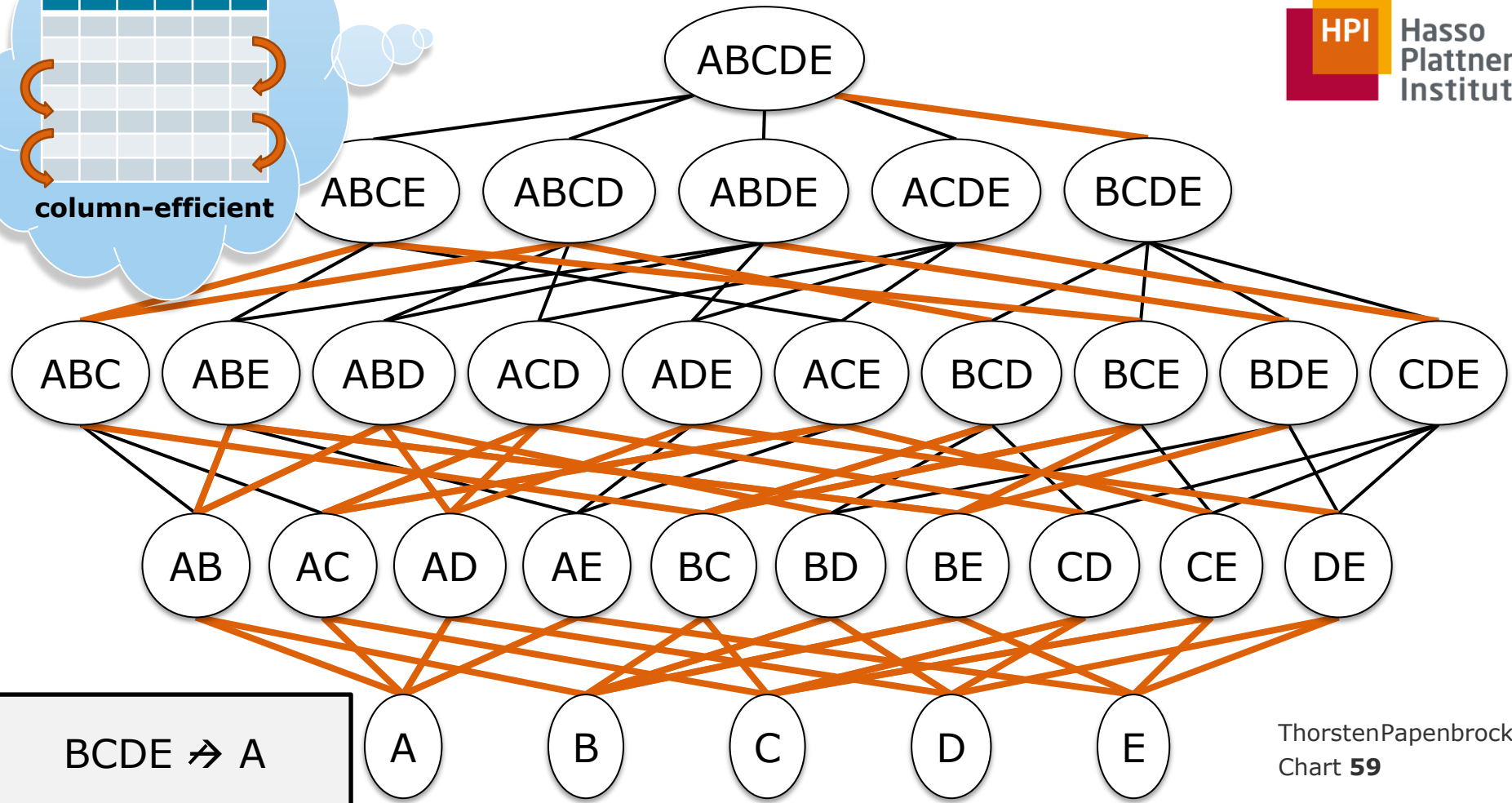
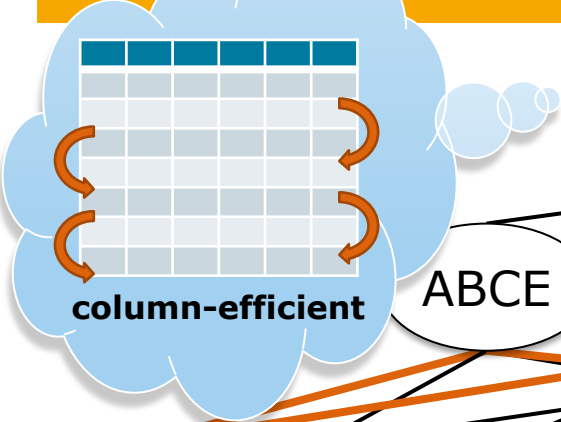


AD \rightarrow BCE

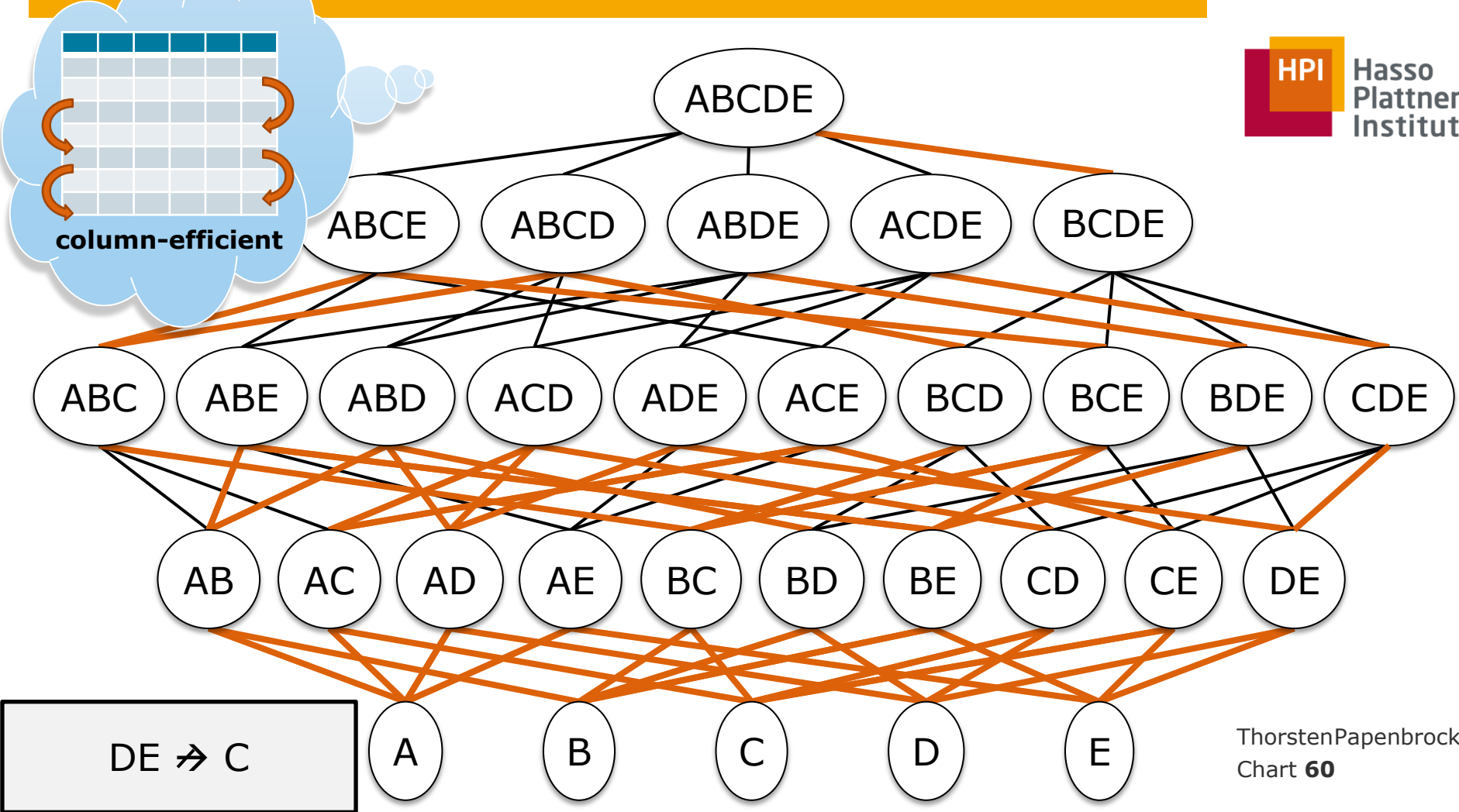


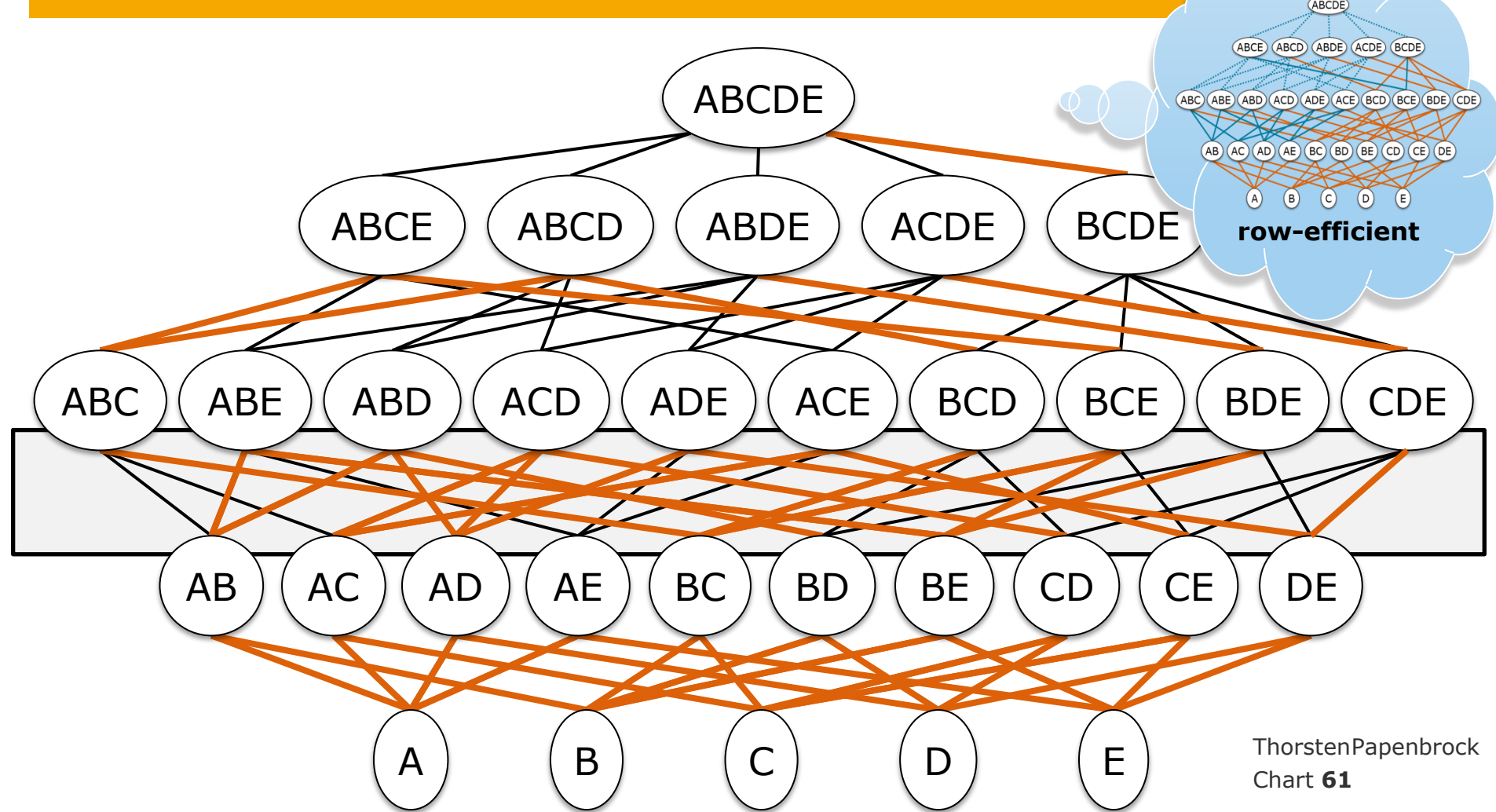
BE \rightarrow ACD

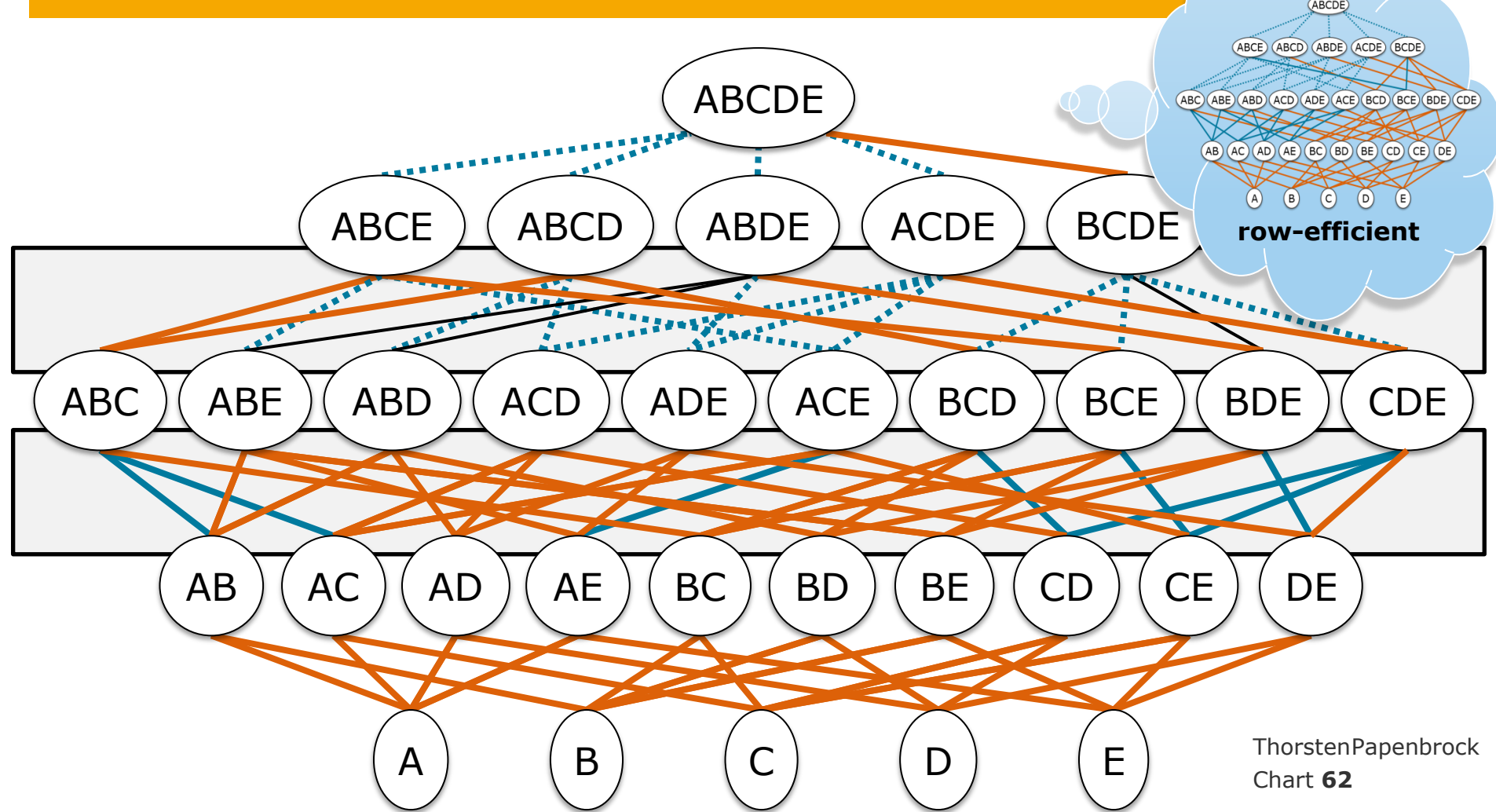


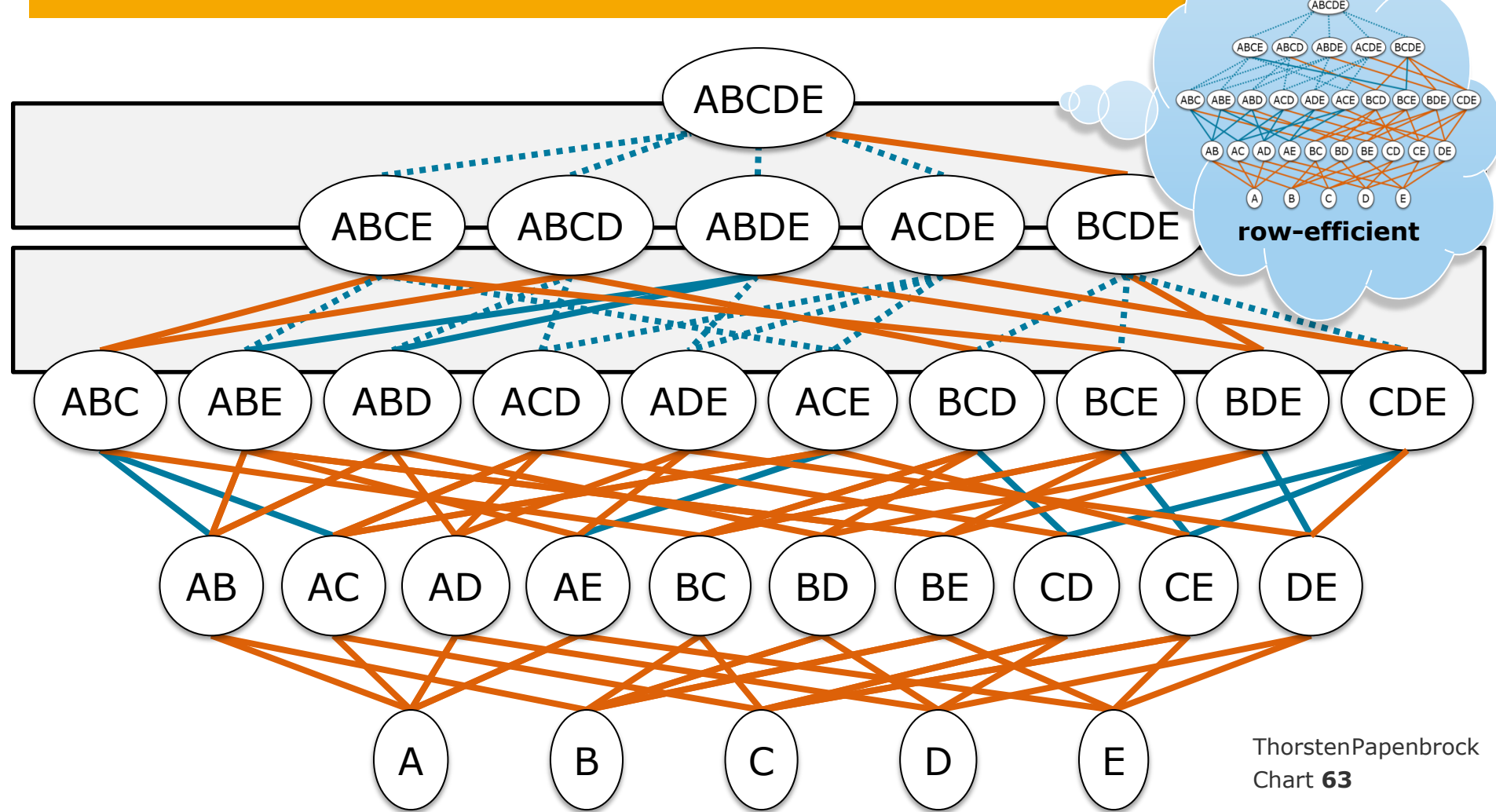


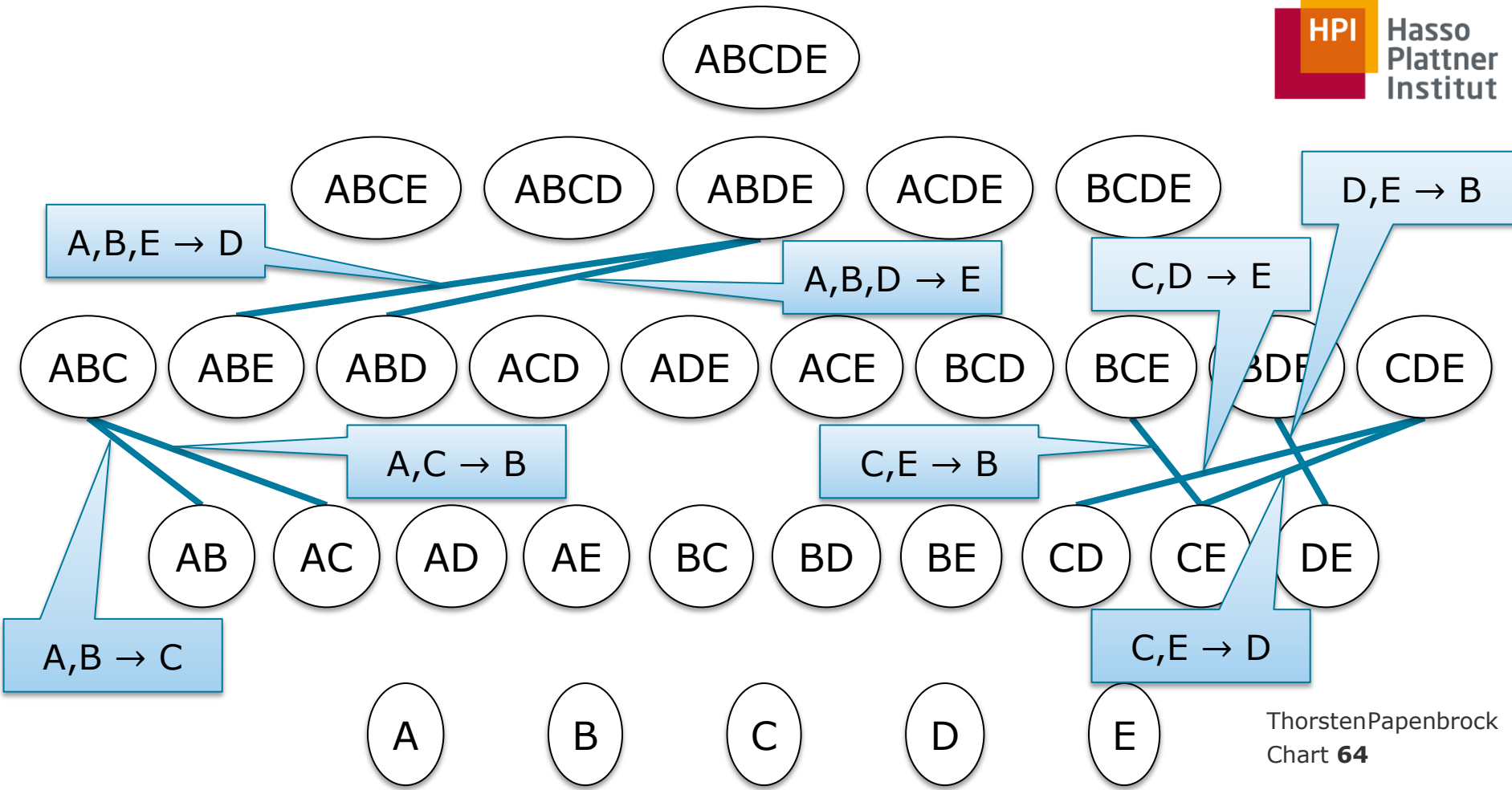
BCDE \rightarrow A



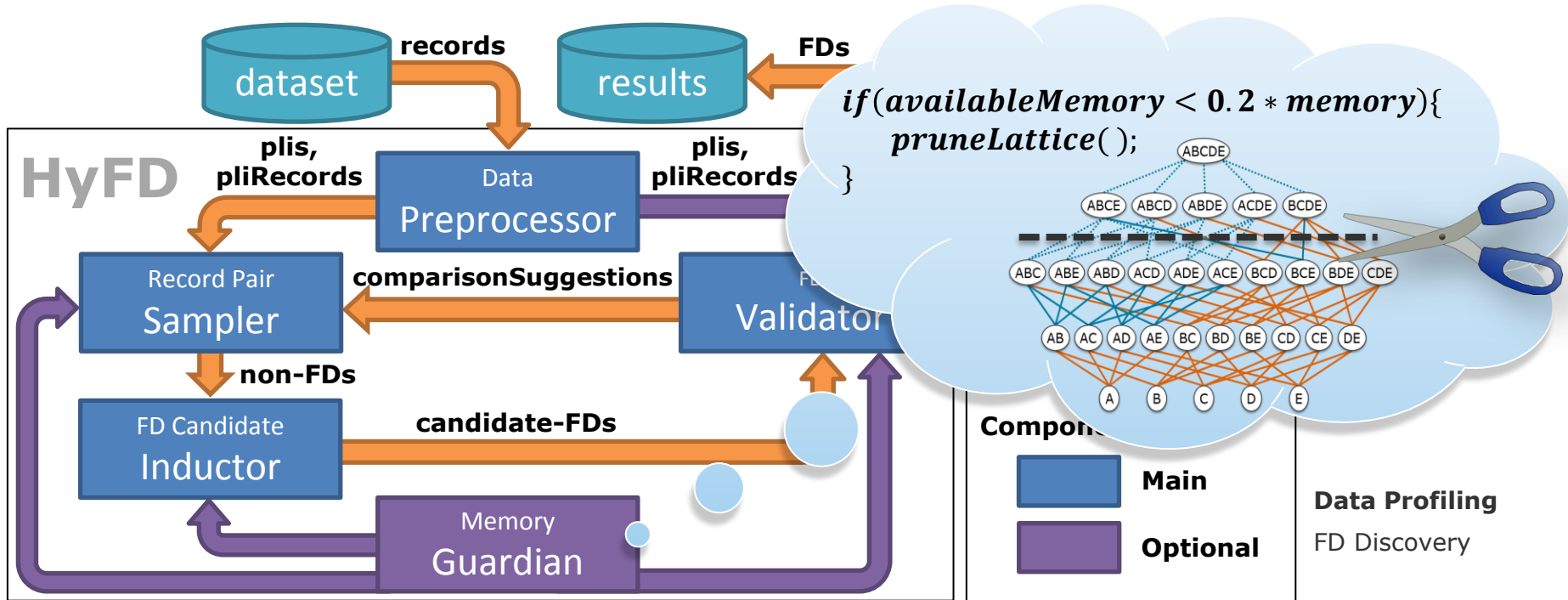




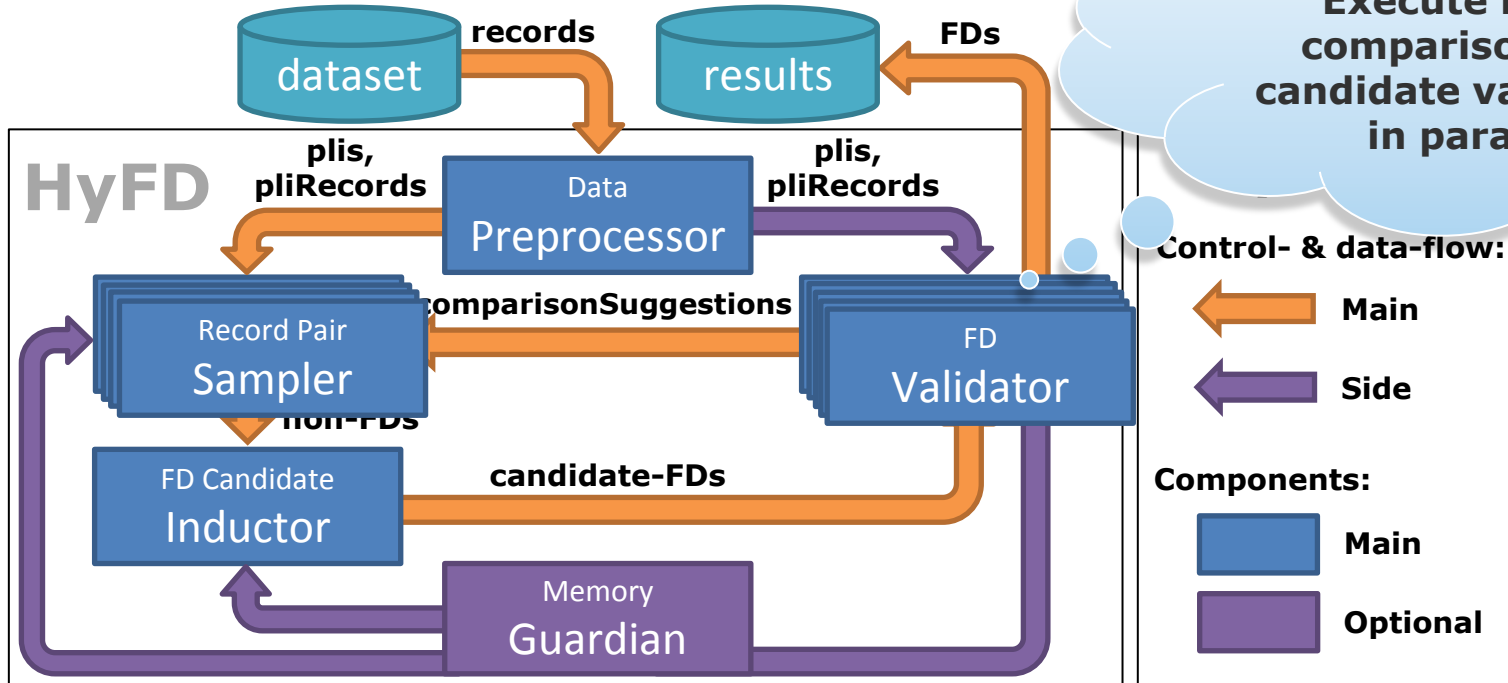




HyFD Hybrid FD Discovery



HyFD Hybrid FD Discovery



Execute record comparisons and candidate validations in parallel.

Data Profiling
FD Discovery

- Motivation
- Approaches
- fdep
- HyFD
- **Evaluation**
- Metanome



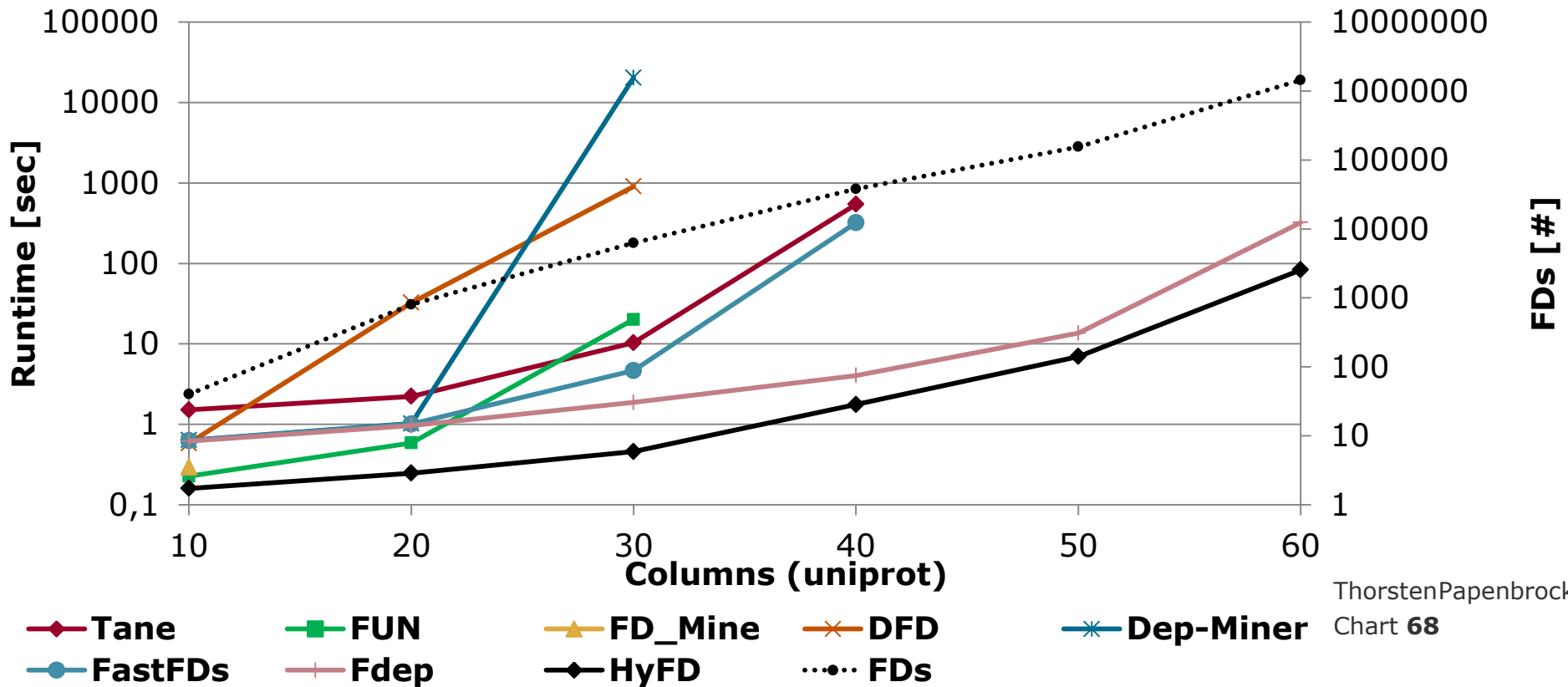
Data Profiling

FD Discovery

ThorstenPapenbrock
Chart **67**

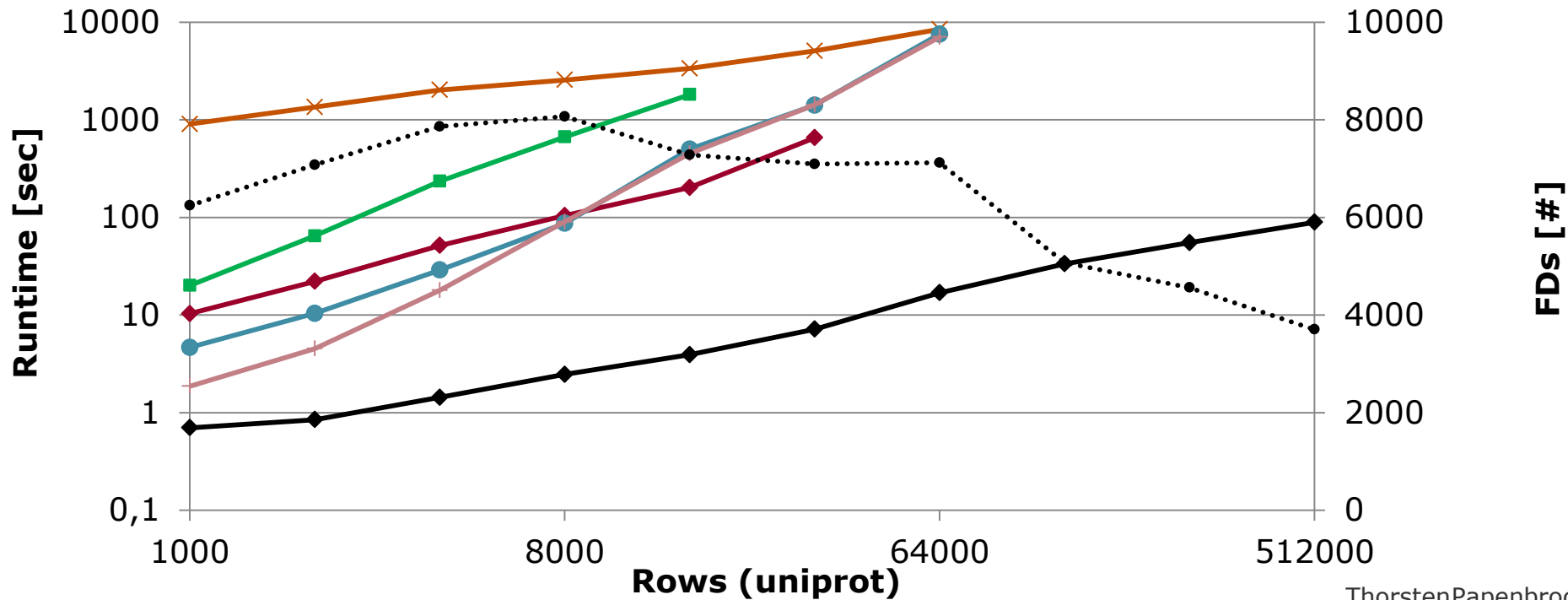
Evaluation

Column Scalability



Evaluation

Row Scalability



Evaluation

Various Datasets

Dataset	Cols [#]	Rows [#]	Size [KB]	FDs [#]
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	527
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	128,727
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	>2,437,556

Results larger than 1,000 FDs are only counted

TL: time limit of 4 hours exceeded

ML: memory limit of 100 GB exceeded

Table 3.1: Runtimes in seconds for several real-world datasets.

TANE [Huhtala et al., 1999]
DEP-MINER [Lopes et al., 2000]

FUN [Novelli and Cicchetti, 2001]
FASTFDs [Wyss et al., 2001]

FD_MINE [Yao et al., 2002]
FDEP [Flach and Savnik, 1999]

DFD [Abedjan et al., 2014c]
HYFD [Papenbrock and Naumann, 2016]

Data Profiling

FD Discovery

Thorsten Papenbrock

Chart 70

Evaluation

Various Datasets

Dataset	Cols	Rows	Size	FDs	HyFD	
	[#]	[#]	[MB]	[#]	[s/m/h/d]	
TPC-H.lineitem	16	6 m	1,051	4 k	39 m	4 m
PDB.POLY_SEQ	13	17 m	1,256	68	4 m	3 m
PDB.ATOM_SITE	31	27 m	5,042	10 k	12 h	64 m
SAP_R3.ZBC00DT	35	3 m	783	211	4 m	2 m
SAP_R3.ILOA	48	45 m	8,731	16 k	35 h	8 h
SAP_R3.CE4HI01	65	2 m	649	2 k	17 m	10 m
NCVoter.statewide	71	1 m	561	5 m	10 d	31 h
CD.cd	107	10 k	5	36 k	5 s	3 s

Data Profiling

FD Discovery

ThorstenPapenbrock
Chart **71**

sequential

32 x parallel

- Motivation
- Approaches
- fdep
- HyFD
- Evaluation
- **Metanome**



www.metanome.de

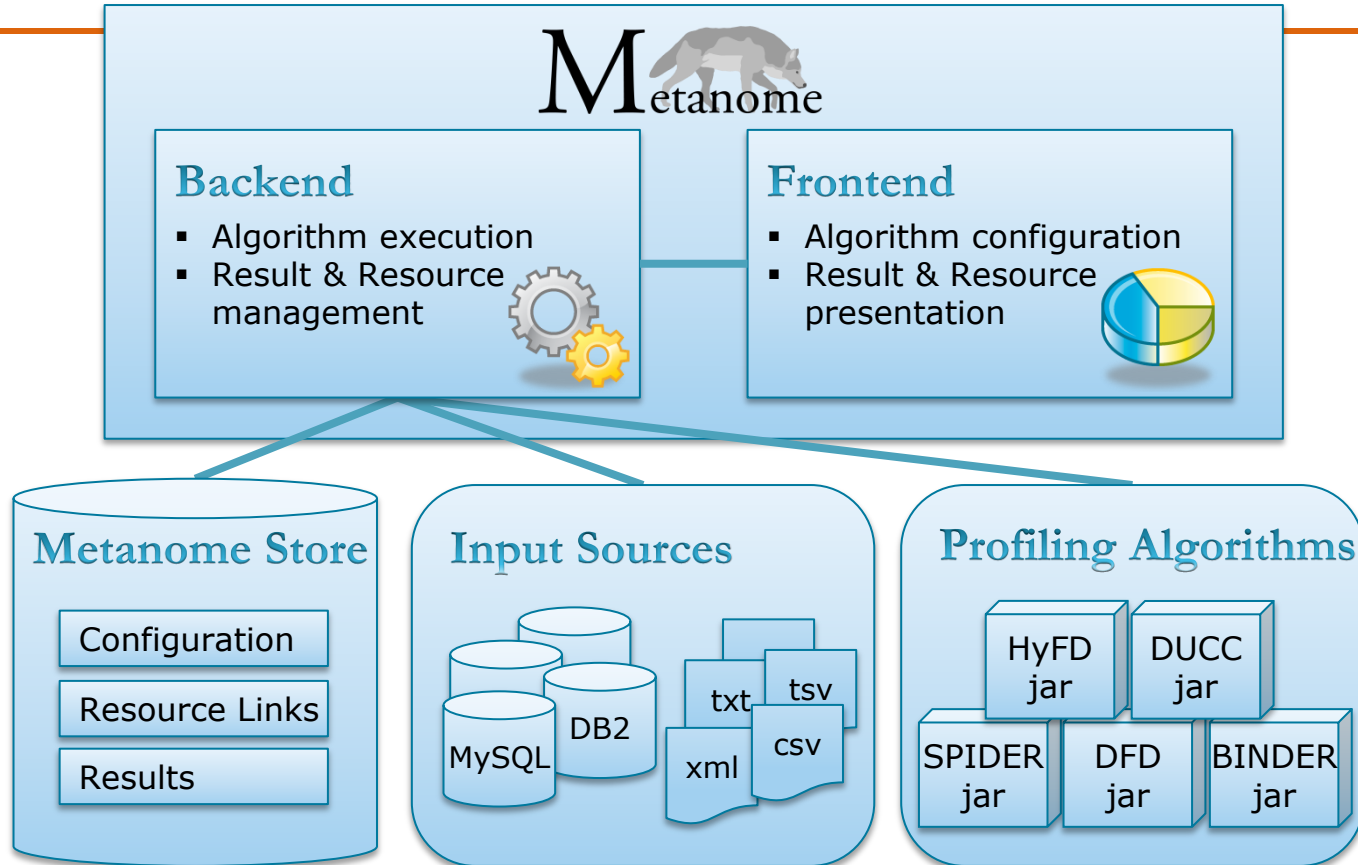
Data Profiling

FD Discovery

ThorstenPapenbrock

Chart 72

Metanome Architecture



Data Profiling

FD Discovery

ThorstenPapenbrock
Chart 73

Algorithms

$A \rightarrow B$

$A \subset B$

$A \gg B$

$A \underline{B} C D$

Datasets



NEW HISTORY ABOUT Melanotic

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

Functional Dependency

Determinant	Dependant
[WDC_planets.csv.Name]	WDC_planets.csv.EquatorialDiameter
[WDC_planets.csv.Name]	WDC_planets.csv.Mass
[WDC_planets.csv.Name]	WDC_planets.csv.OrbitalRadius
[WDC_planets.csv.Name]	WDC_planets.csv.OrbitalPeriod
[WDC_planets.csv.Name]	WDC_planets.csv.RotationPeriod
[WDC_planets.csv.Name]	WDC_planets.csv.ConfirmedMoons
[WDC_planets.csv.Name]	WDC_planets.csv.Atmosphere
[WDC_planets.csv.Name]	WDC_planets.csv.Type
[WDC_planets.csv.Name]	WDC_planets.csv.Rings
[WDC_planets.csv.EquatorialDiameter]	WDC_planets.csv.Name

Rows per page: 10 1 - 10 of 66

Data Profiling

FD Discovery

ThorstenPapenbrock
Chart 74

Metanome Developers



Jakob Zwiener



Claudia Exeler



Tanja Bergmann



Moritz Finke



Fabian Tschirschnitz



Carl Ambroselli



Maxi Fischer



Vincent Schwarzer



?



?

Data Profiling

FD Discovery

Thorsten Papenbrock
Chart **75**



Data Profiling Functional Dependency Discovery

Thorsten Papenbrock
G-3.1.09, Campus III
Hasso Plattner Institut