



# Functional Dependency Discovery: An Experimental Evaluation of Seven Algorithms

Thorsten Papenbrock, Jens Ehrlich, Jannik Marten, Tommy Neubert,  
Jan-Peer Rudolph, Martin Schönberg, Jakob Zwiener, Felix Naumann

# Motivation

## Functional Dependencies and Normalization

<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	Frankfurt	Jakobs
Mike	Moore	60329	Frankfurt	Feldmann

Definition FD:  $X \rightarrow A$

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

Postcode  $\rightarrow$  City  
Postcode  $\rightarrow$  Mayor

# Motivation

## Functional Dependencies and Normalization

<u>Name</u>	<u>Surname</u>	<u>Postcode</u>
Thomas	Miller	14482
Sarah	Miller	14482
Peter	Smith	60329
Jasmine	Cone	01069
Thomas	Cone	14482
Mike	Moore	60329

<u>Postcode</u>	<u>City</u>	<u>Mayor</u>
14482	Potsdam	Jakobs
60329	Frankfurt	Feldmann
01069	Dresden	Orosz

**Less Memory Consumption**  
**Less Anomaly Vulnerability**

Definition FD:  $X \rightarrow A$

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

Postcode  $\rightarrow$  City  
Postcode  $\rightarrow$  Mayor

# Motivation

## Very large Tables e.g. uniprot

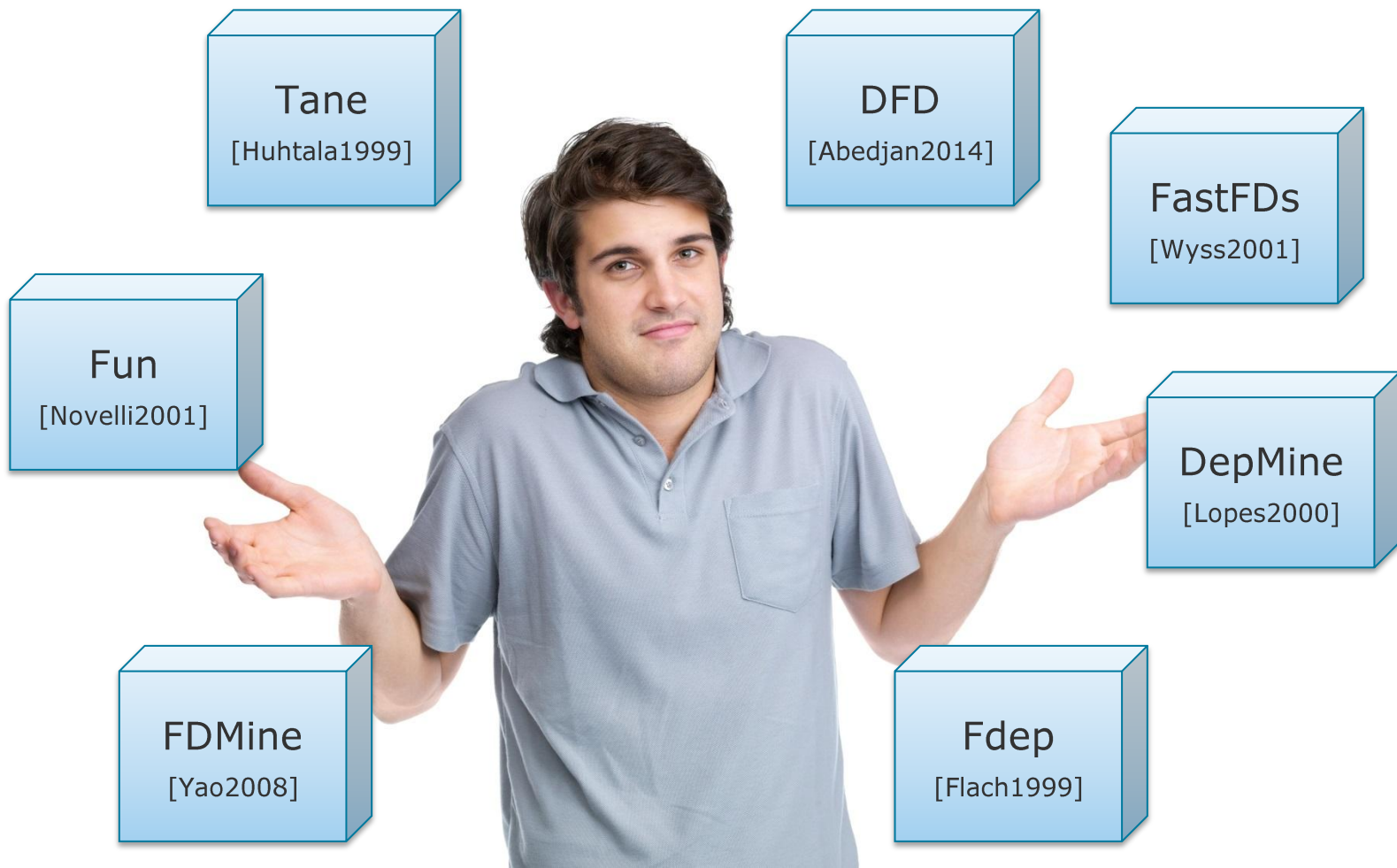
name	references pubmed	keywords	gene ontology	protein existence	virus host ncbi	sequence	sequence length	sequence mass	sequence checksum	sequence modified	taxonomic	taxonomic lineage	protein names	gene orf	cross reference	cross ref
017L_FRG3G	15165820	Complete proteome; Reference predicted			8295; 8316; 8404; 3034	METMSDYSKEVSE, 502	53469		A747EE6F952CBAD7	2004-07-19+01:00	654924	Viruses > dsDNA virus Uncharacterized p FV3-017L	AY548484; AAT0967e	YP_031595.1		
040R_FRG3G	15165820	Complete proteome; Reference predicted			8295; 8316; 8404; 3034	MIRALCTIVLIAAGV, 182	19577		FCF64F420039CF35	2004-07-19+01:00	654924	Viruses > dsDNA virus Uncharacterized p FV3-040R	AY548484; AAT0969e	YP_031618.1		
057L_IIV3	15165820	Complete proteome; Reference predicted			8295; 8316; 8404; 3034	MEWVDFE6MGGU, 182	19577		84854328631A0855	2004-07-19+01:00	654924	Viruses > dsDNA virus Uncharacterized p FV3-057L	AY548484; AAT0968e	YP_031617.1		
078L_FRG																
1001R_AS																
11014_AS																
1107L_AS																
124R_IIV3																
14331_CA																
14335_AR																
14338_VII																
1433T_BO																
1433_NEO																
14KL_BRU																
17KD_RIC																
1A11_ORY																
1A1D_BUI																
1A1D_PSE																
1A43_HUI																
1807_HUI																
1837_HUI																
1BFH_AEG																
242L_IIV6																
2A5D_HU																
2AAB_PIG																
2ABD_DA																
2B_PEBV																
2S51_BRA																
3001L_AS																
348R_IIV6																
3601S_AS																
3603L_AS																
385L_IIV6																
3BHS_HO																
3DHQ1_A																
3DHQ3_SC																
3HAO_CH																
3HGA_XA																
3MGA_BA																
3MGH_BC																
3MGH_CL																
3MGH_FR																
3MGH_MH																
3MGH_RH																
3MGH_ST																
3MG_HUN																
441R_IIV6																
4CL4_ORY																
4EBP3_HU																
5059S_AS																
5E5_RAT																
5HT1B_RA																
5HT2A_PI																
5HT5A_HU																
5NTSL_MOUSE	15489334; 1641072; 1946	Acetylation; Alternative splicing; evidence at transcript level				MKATVLMRQPGRI, 292	33578		3AE268CE06833085	2008-04-08+01:00	10090	Eukaryota > Metazoa > Cytosolic 5'-nucleotidase III-like BC015307; AAH1530	NP_080837.3			
60A_DROVI	8688461	Cytokine; Disulfide bond; Glyco inferred from homology				MTASLVLPSSLWLI, 436	49999		C74484AE58796692	1997-11-01Z	7244	Eukaryota > Metazoa > Protein 60A	U48595; AAC47262.1; Genomic_C			

**223 attributes!**



# Motivation

## Functional Dependencies Discovery



## FD Algorithm Classes

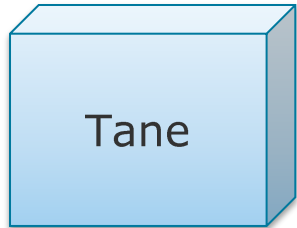


## Experimental Evaluation



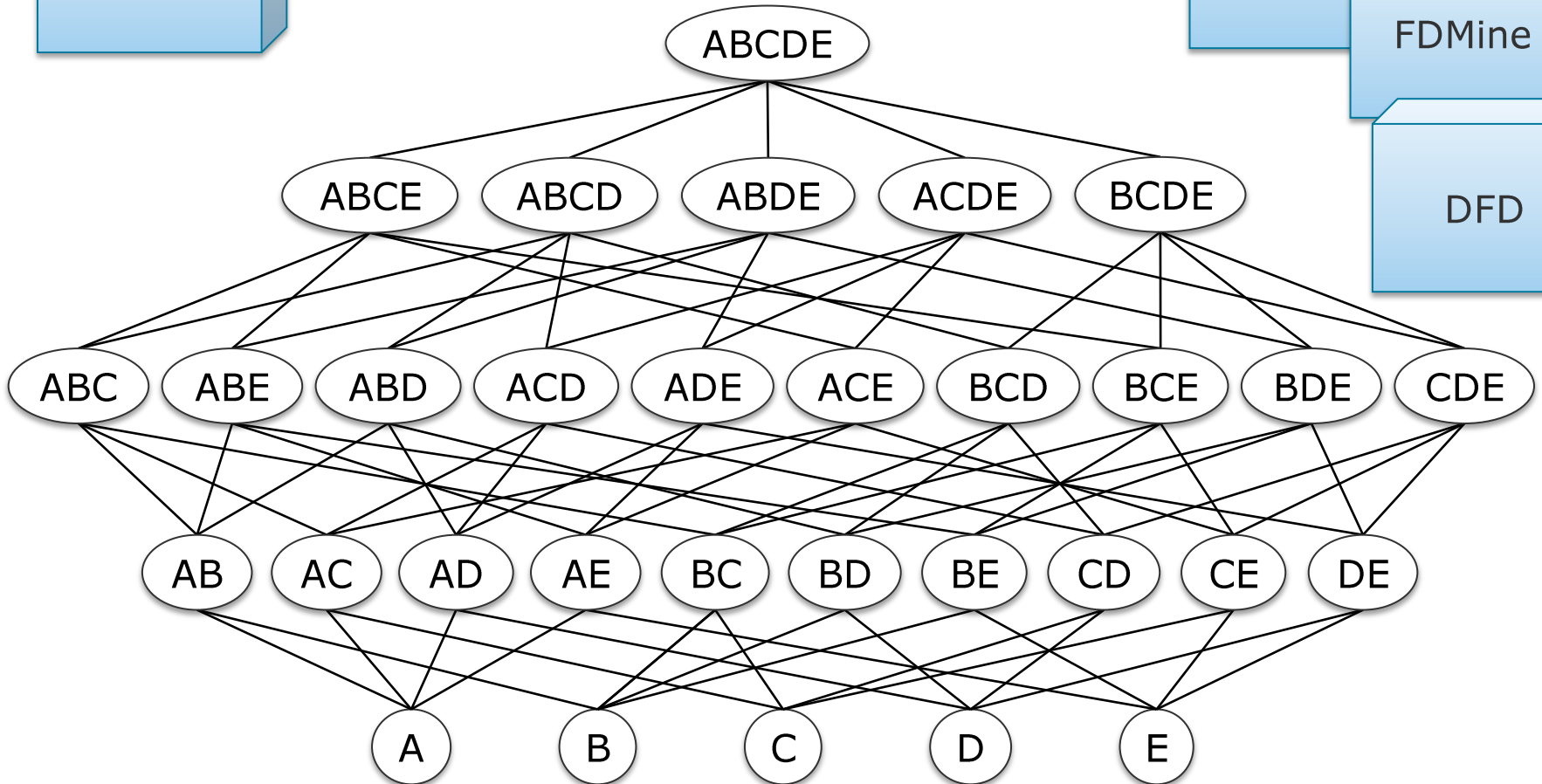
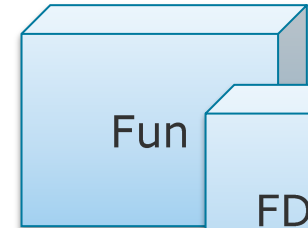
## Extrapolation of Results





# FD Algorithm Classes

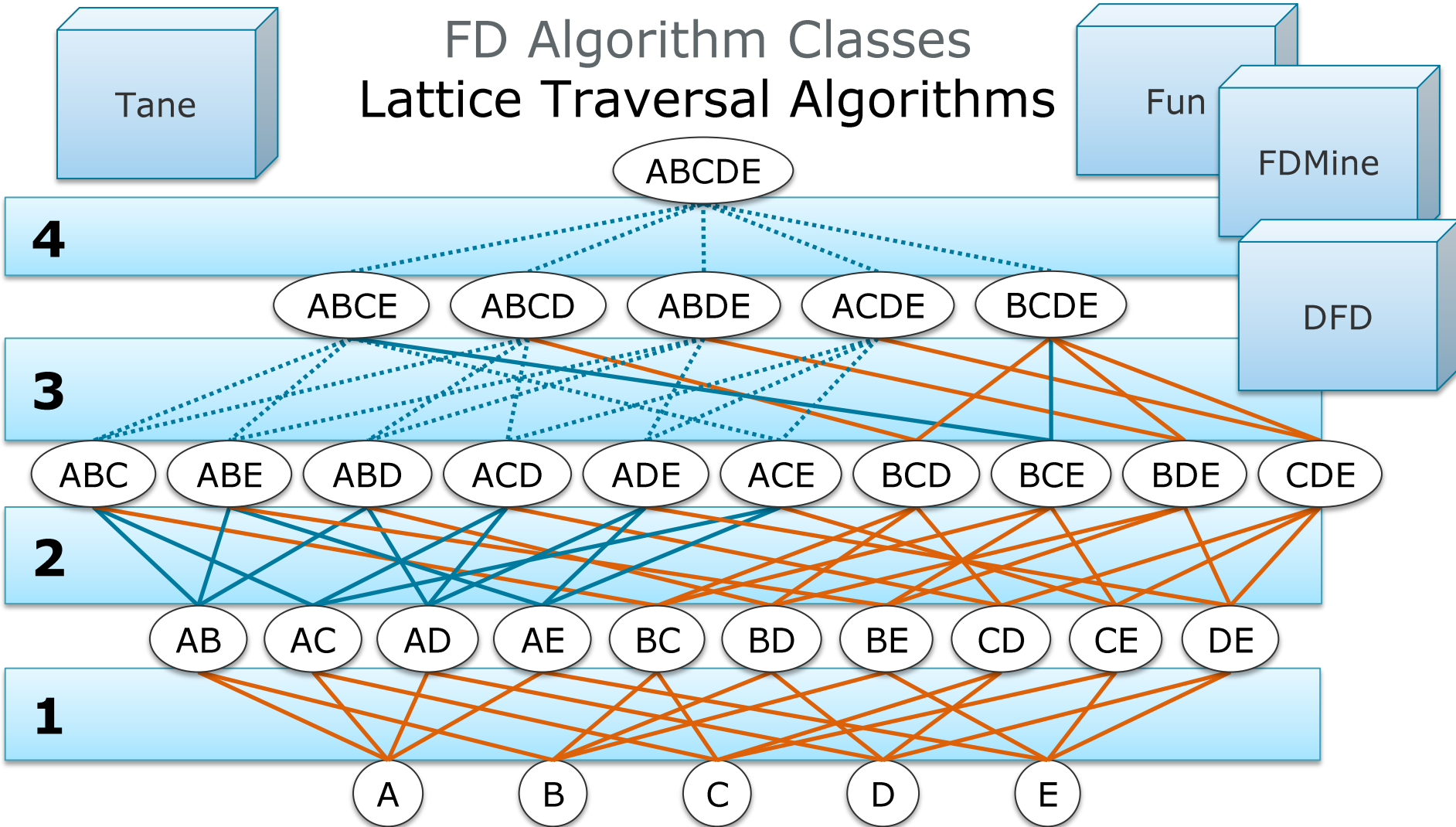
## Lattice Traversal Algorithms





# FD Algorithm Classes

## Lattice Traversal Algorithms



— invalid FD (non-FD)

— minimal FD

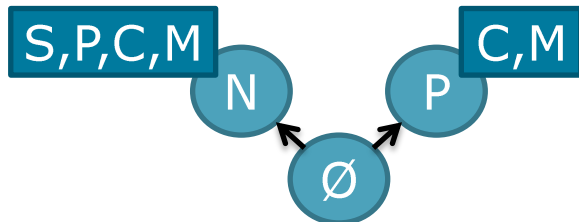
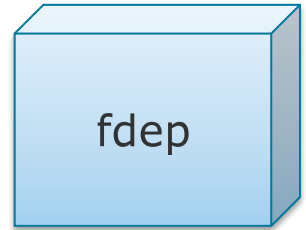
..... non-minimal FD

# FD Algorithm Classes

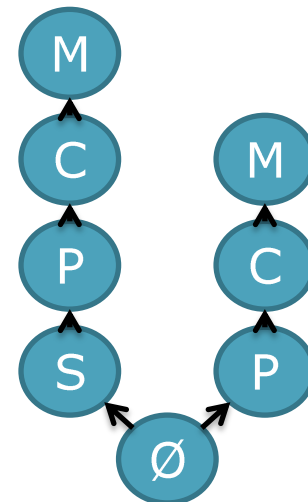
## Dependency Induction Algorithms

<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

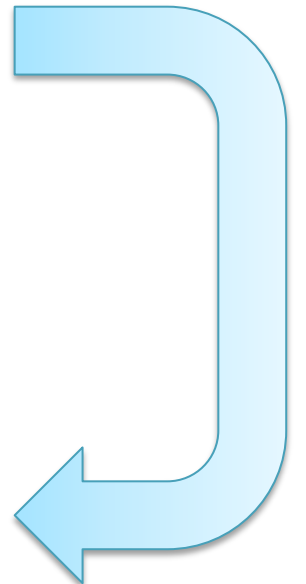
**N,P,C,M**  $\Rightarrow$  **S**



**Positive Cover**

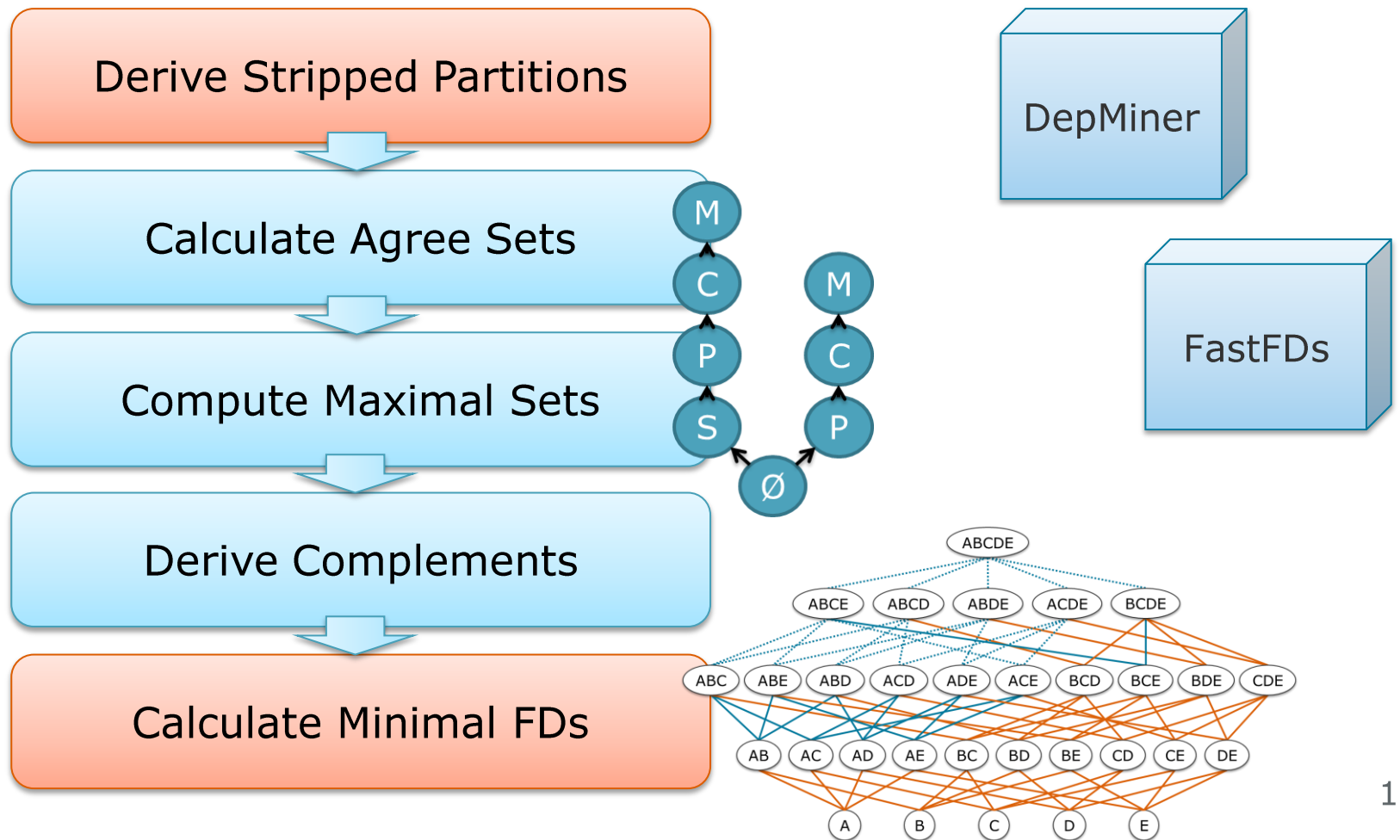


**Negative Cover**



# FD Algorithm Classes

## Difference- and Agree-Set Algorithms



## FD Algorithm Classes



## Experimental Evaluation

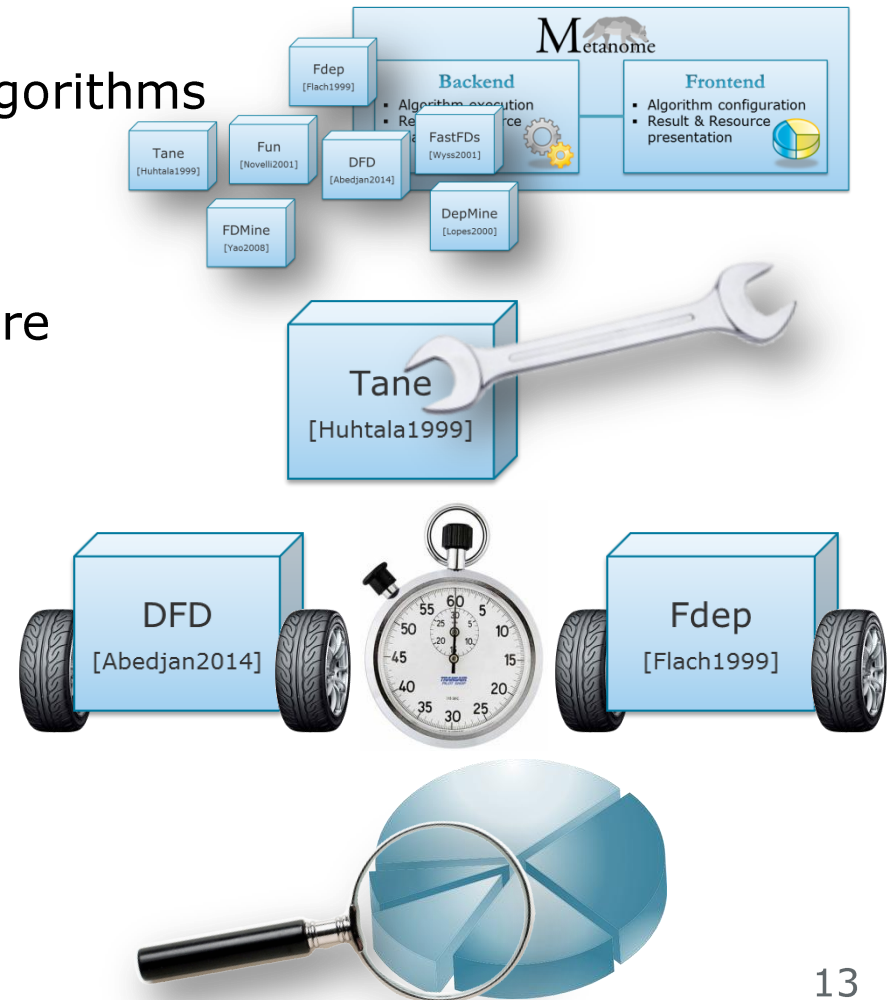


## Extrapolation of Results



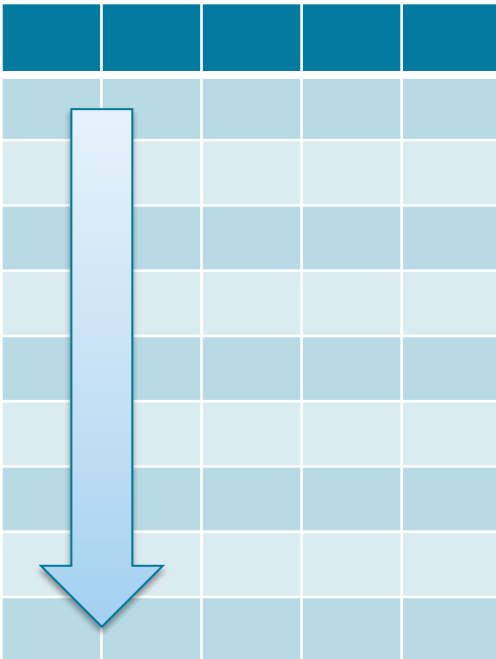
## Experimental Evaluation Contributions

- We **re-implemented** all seven algorithms for the Metanome profiling tool.
- We **amended** the algorithms where original descriptions were sparse.
- We **evaluated** the algorithms under comparable conditions.
- We **analyzed** and extrapolated the experimental results.

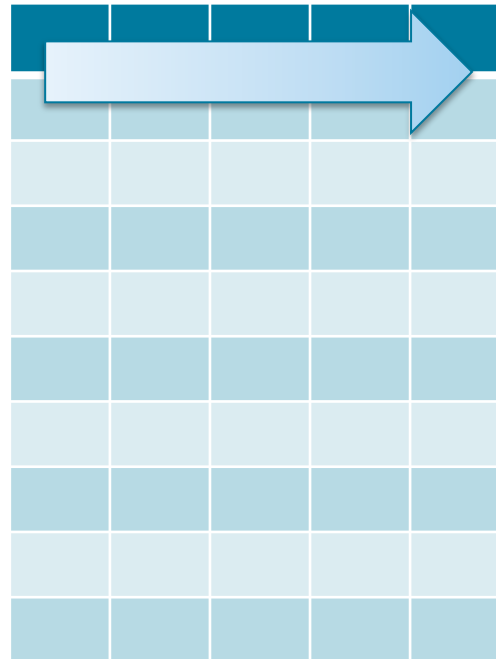


# Experimental Evaluation Evaluation Strategy

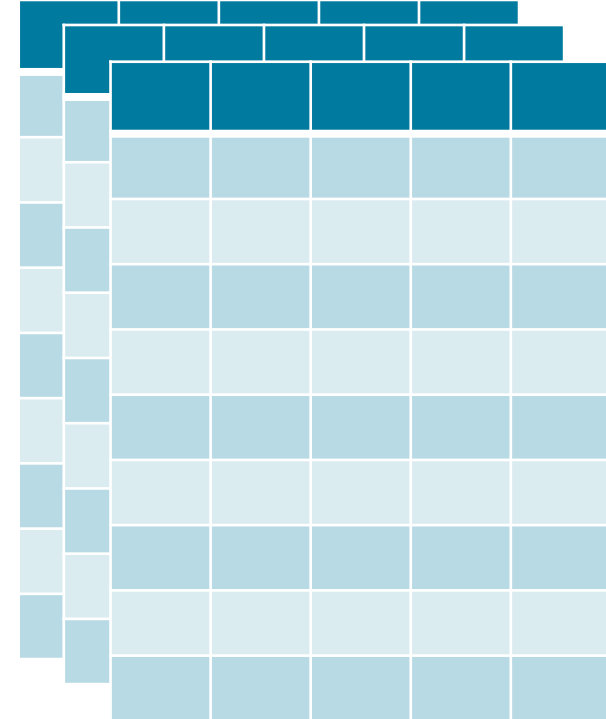
**Row Scalability**



**Column Scalability**

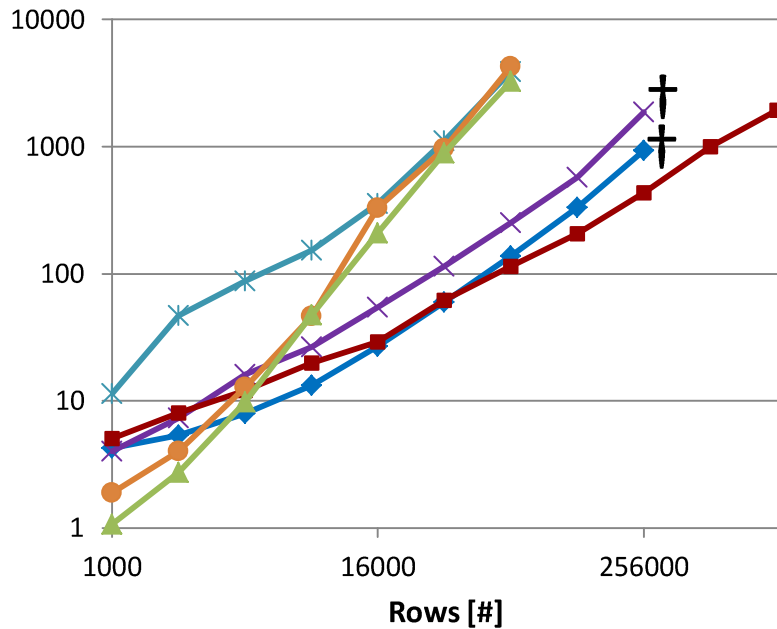


**Different Datasets**

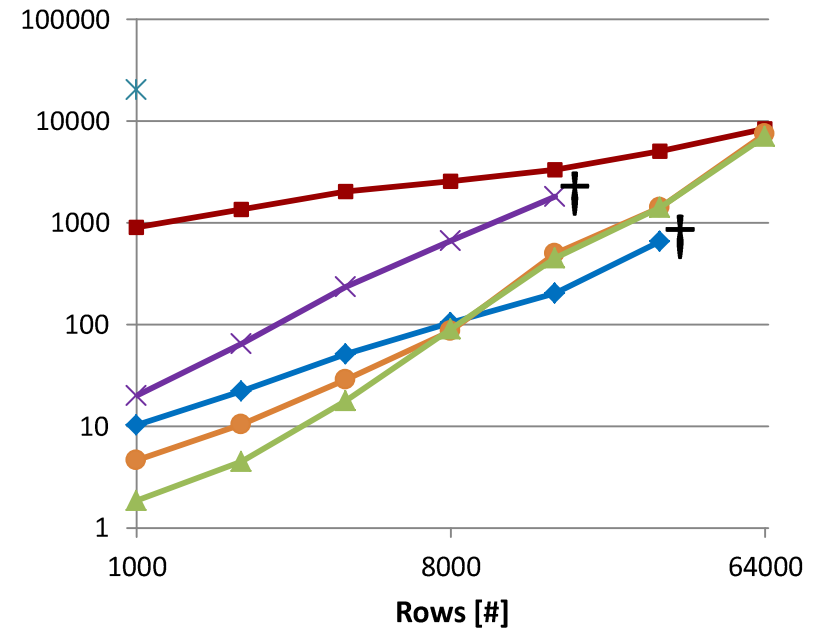


# Experimental Evaluation Row-Scalability

**ncvoter dataset (19 columns)**



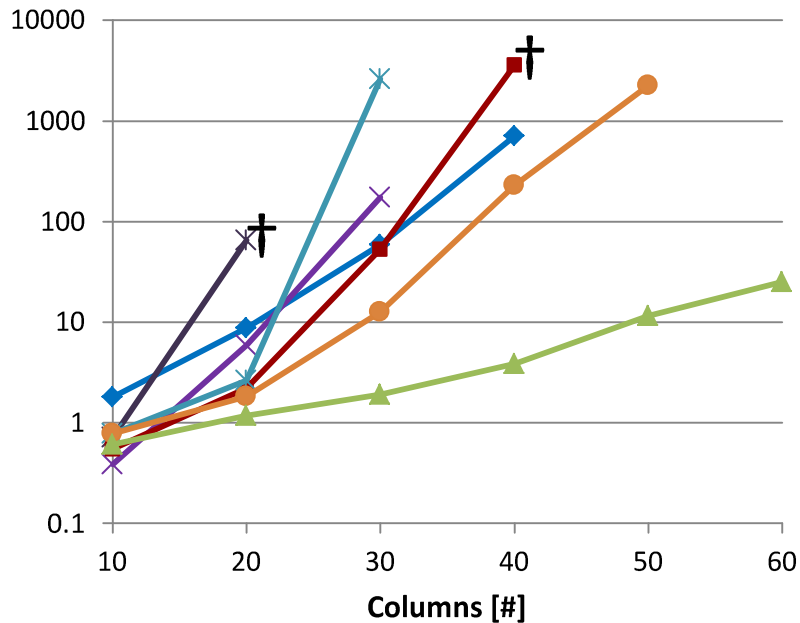
**uniprot dataset (30 columns)**



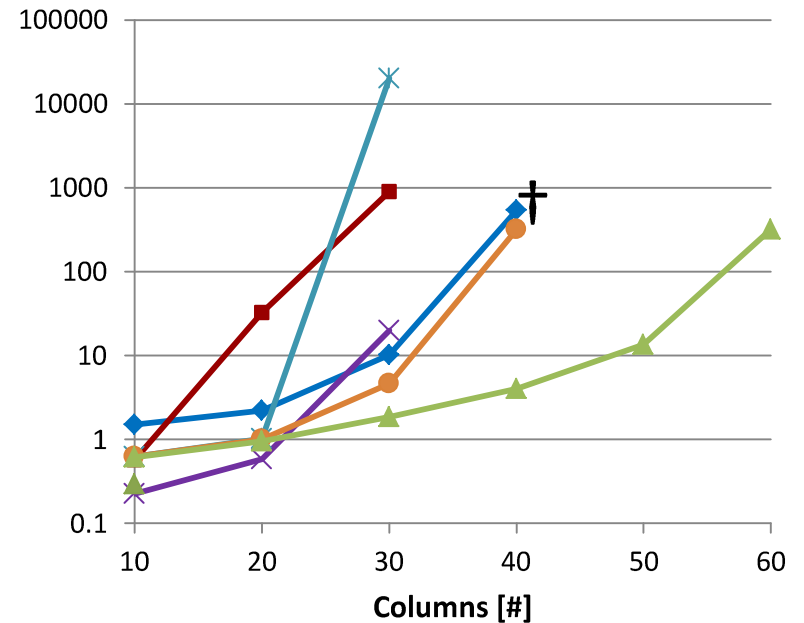
† 100 GB memory exhausted

# Experimental Evaluation Column-Scalability

**plista dataset (1000 rows)**



**uniprot dataset (1000 rows)**



† 100 GB memory exhausted



# Experimental Evaluation Different Datasets

Dataset	Columns [#]	Rows [#]	Size [KB]	FDs [#]	Runtime [sec]						
					TANE [7]	FUN [15]	FD_MINE [21]	DFD [1]	DEP-MINER [12]	FASTFDs [20]	FDEP [6]
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

TL: time limit of 4 hours exceeded

ML: memory limit of 100GB exceeded

1. There are **many FDs** to discover!

2. Algorithms are either **row or column efficient!**

3. No algorithm can handle datasets of **real-world size!**

## FD Algorithm Classes



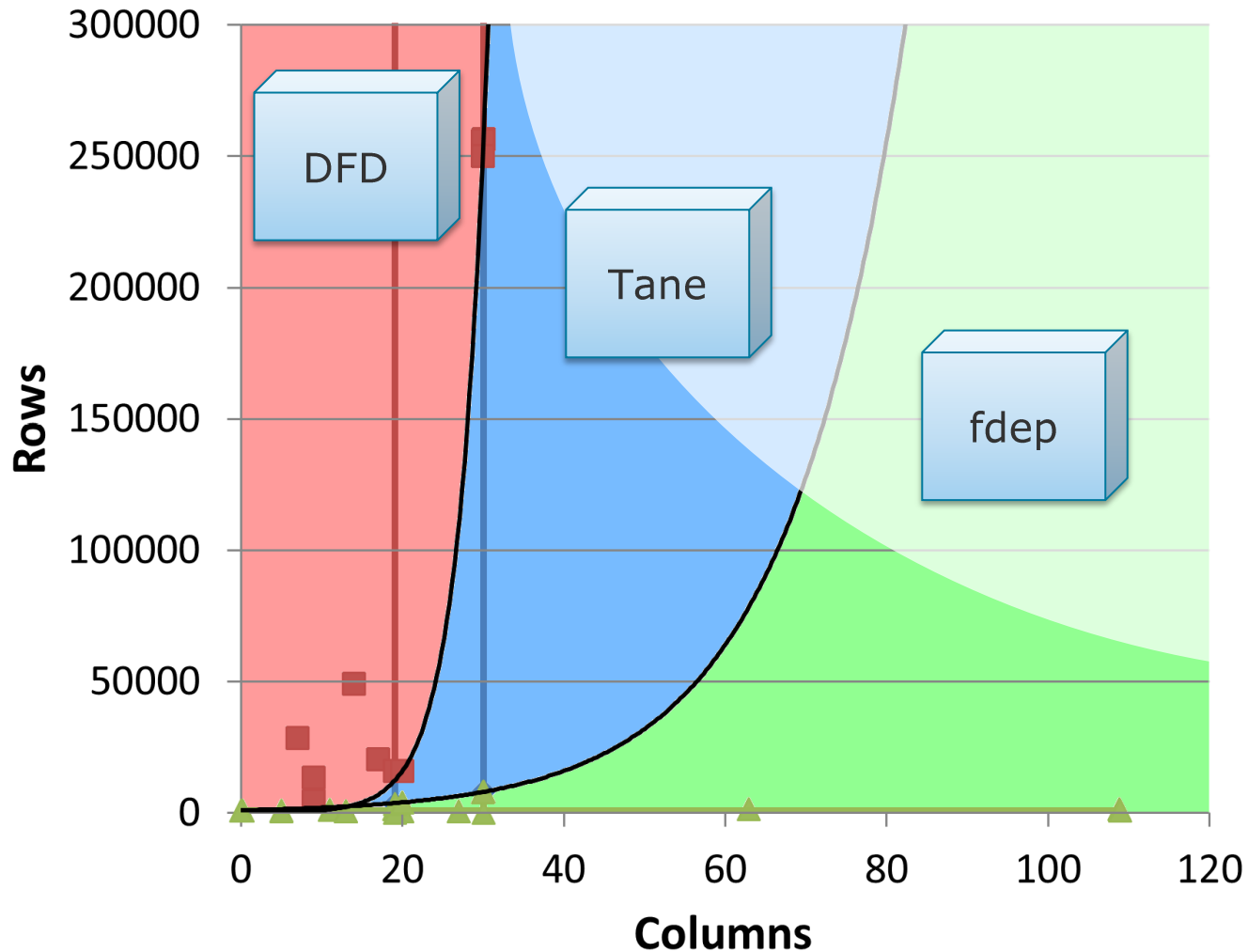
## Experimental Evaluation



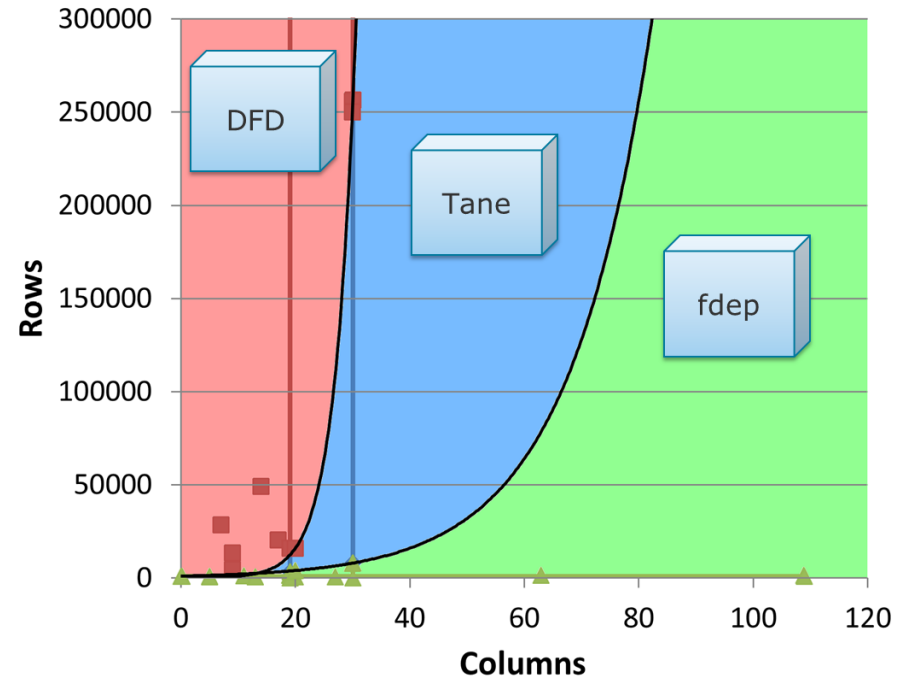
## Extrapolation of Results



## Extrapolation of Results Decision Chart



TANE [7]	FUN [15]	FD_MINE [21]	DFD [1]	DEP-MINER [12]	FASTFDS [20]	FDEP [6]
1.1	<b>0.1</b>	0.2	0.2	0.2	0.2	<b>0.1</b>
1.2	<b>0.1</b>	0.2	0.3	0.3	0.3	0.2
2.9	1.1	3.8	<b>1.0</b>	174.6	164.2	125.5
2.1	<b>0.6</b>	1.8	1.1	3.0	2.9	3.8
4.1	1.8	7.1	<b>0.9</b>	121.2	118.9	46.8
2.3	0.6	2.2	0.8	1.1	1.1	<b>0.5</b>
2.2	0.6	4.2	0.9	0.5	0.6	<b>0.2</b>
1.6	0.4	69.9	1.2	0.5	0.5	<b>0.2</b>
67.4	111.6	531.5	<b>5.9</b>	6039.2	6033.8	860.2
260.0	529.0	7204.8	<b>6.0</b>	1090.0	1015.5	291.3
4.3	4.0	ML	5.1	11.4	1.9	<b>1.1</b>
12.2	175.9	ML	326.7	5576.5	9.5	<b>0.8</b>
457.0	TL	ML	TL	TL	385.8	<b>7.2</b>
<b>41.1</b>	77.7	ML	TL	377.2	382.4	TL
ML	ML	ML	TL	TL	TL	<b>26.9</b>
ML	ML	ML	TL	TL	TL	<b>216.5</b>
ML	ML	ML	TL	TL	TL	ML



## Functional Dependency Discovery: An Experimental Evaluation of Seven Algorithms

Thorsten Papenbrock, Jens Ehrlich, Jannik Marten, Tommy Neubert,  
Jan-Peer Rudolph, Martin Schönberg, Jakob Zwiener, Felix Naumann