

© Xeeldek (www.sketchport.com)

A Hybrid Approach to Functional Dependency Discovery

Thorsten Papenbrock and Felix Naumann

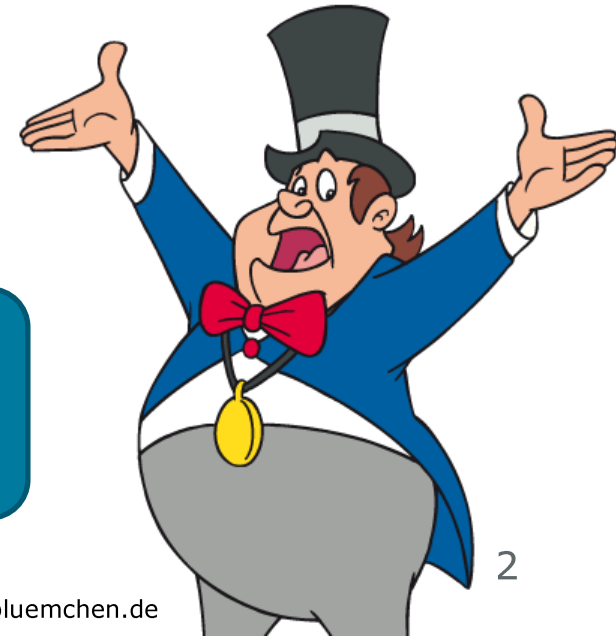


<u>Name</u>	<u>Surname</u>	<u>Postcode</u>	<u>City</u>	<u>Mayor</u>
Thomas	Miller	14482	Potsdam	Jakobs
Sarah	Miller	14482	Potsdam	Jakobs
Peter	Smith	60329	Frankfurt	Feldmann
Jasmine	Cone	01069	Dresden	Orosz
Thomas	Cone	14482	Jakobs	
Mike	Moore	60329	Frankfurt	Feldmann

Definition FD: **X → 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 → City
Postcode → Mayor





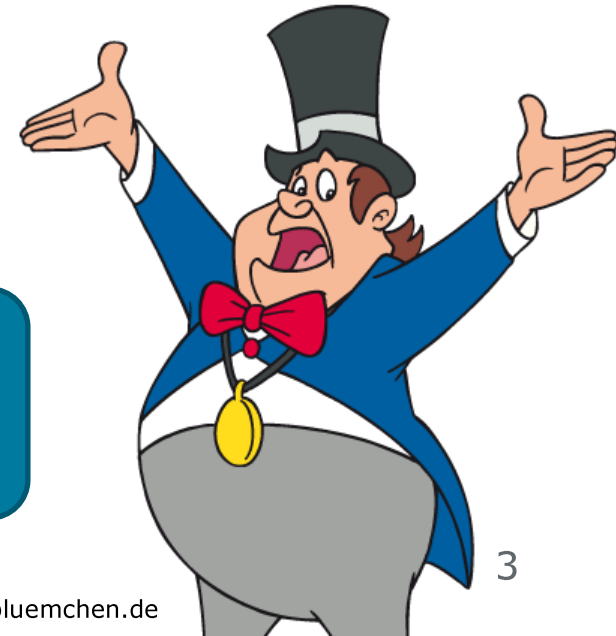
<u>Name</u>	<u>Surname</u>	<u>Postcode</u>
Thomas	Miller	14482
Sarah	Miller	14482
Peter	Smith	60329
Jasmine	Cone	01069
Thomas	Cone	14482
Mike	Moore	60329

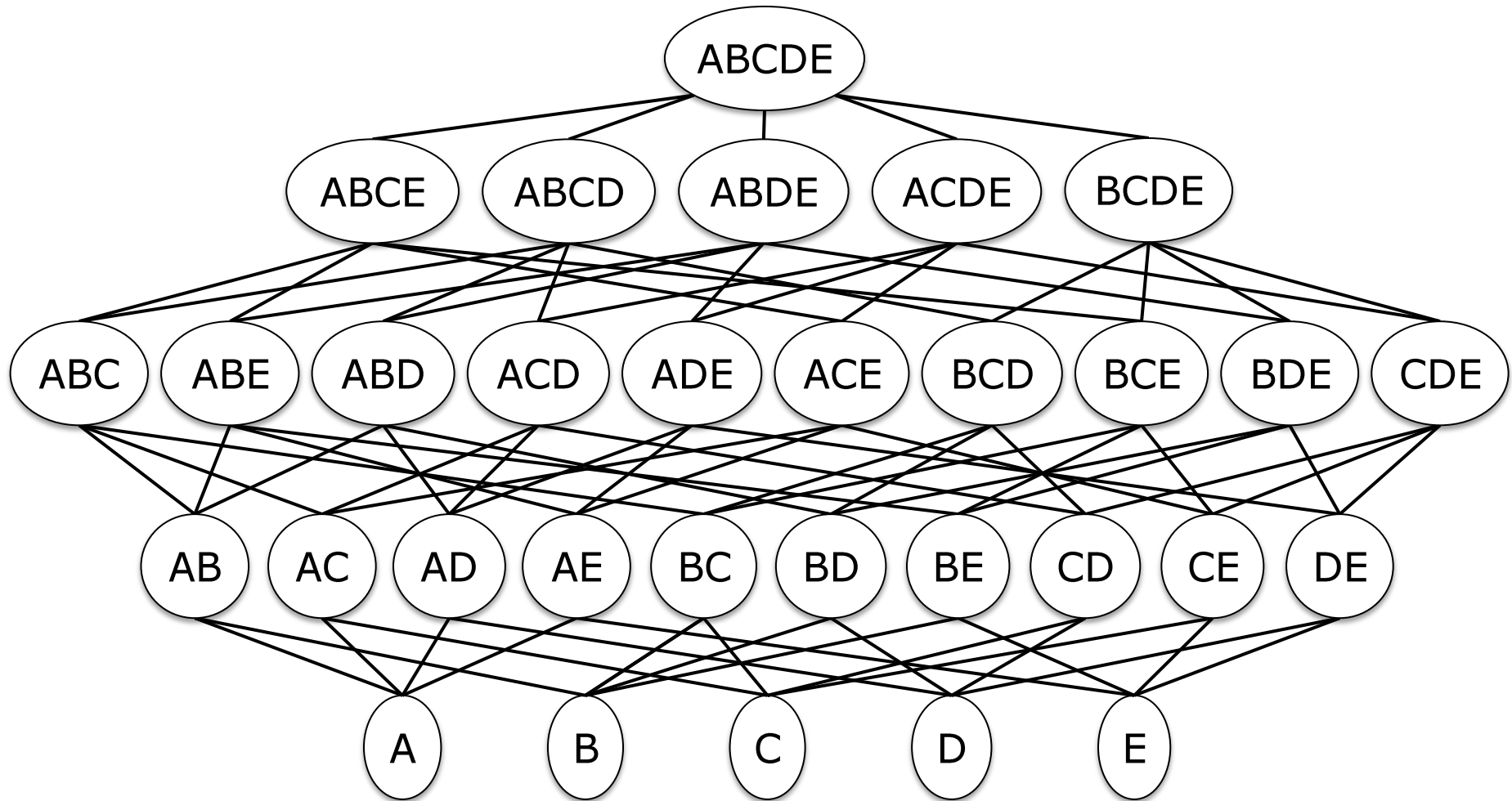
<u>Postcode</u>	<u>City</u>	<u>Mayor</u>
14482	Potsdam	Jakobs
60329	Frankfurt	Feldmann
01069	Dresden	Orosz

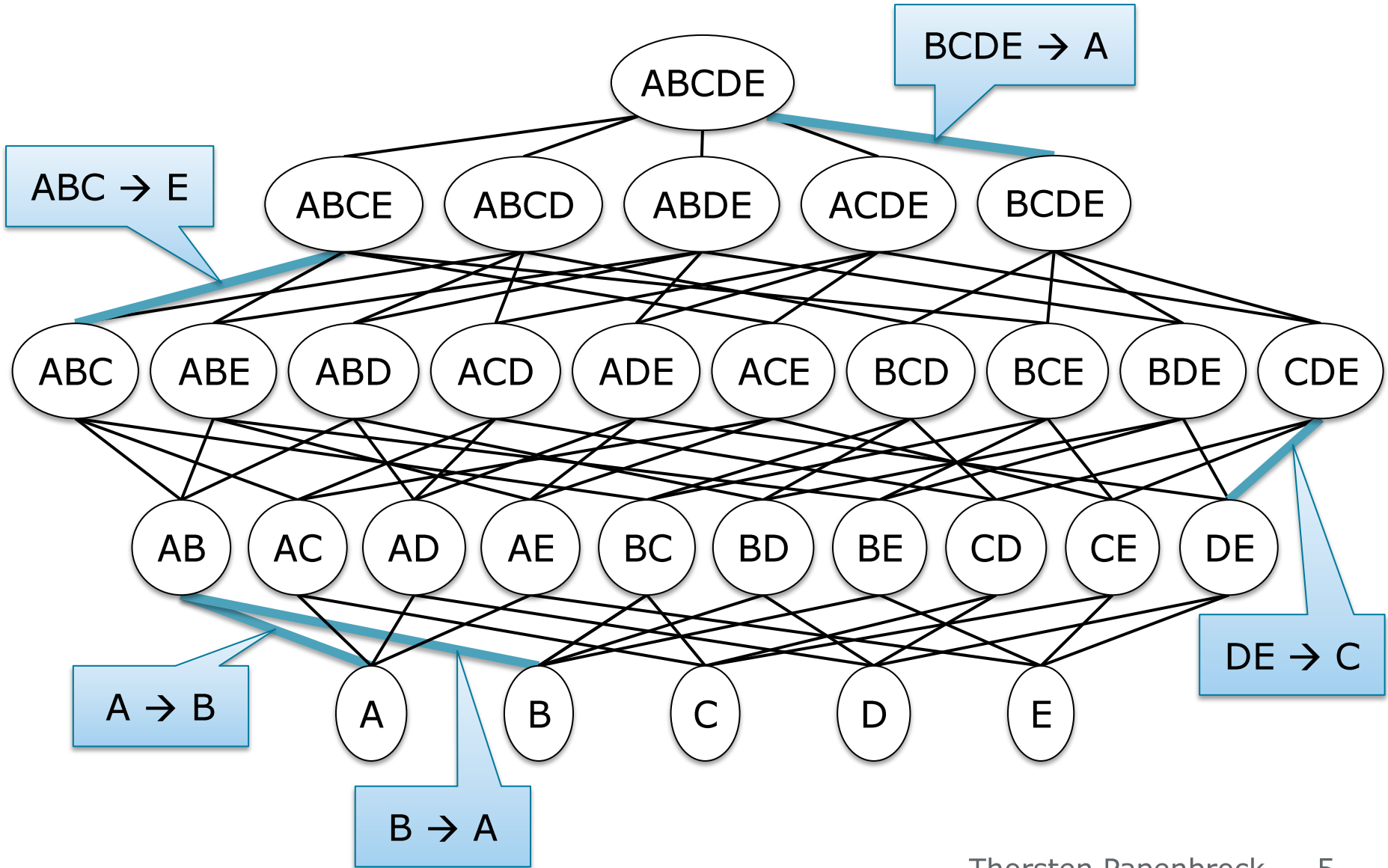
Definition FD: **X → A**

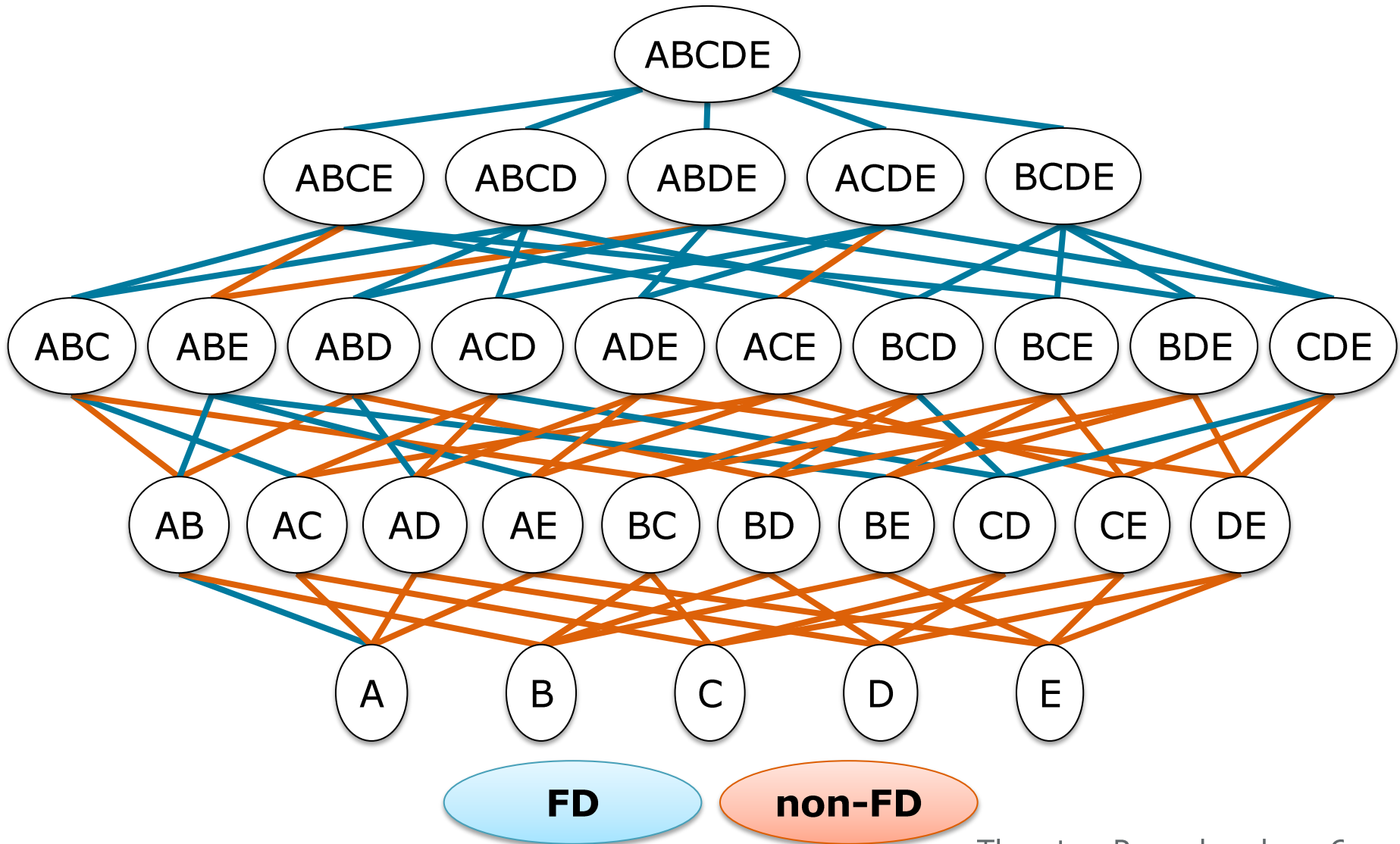
- All values in X uniquely define the values in A.
- If $t_1[X] = t_2[X]$, then $t_1[A] = t_2[A]$.

Postcode → City
Postcode → Mayor

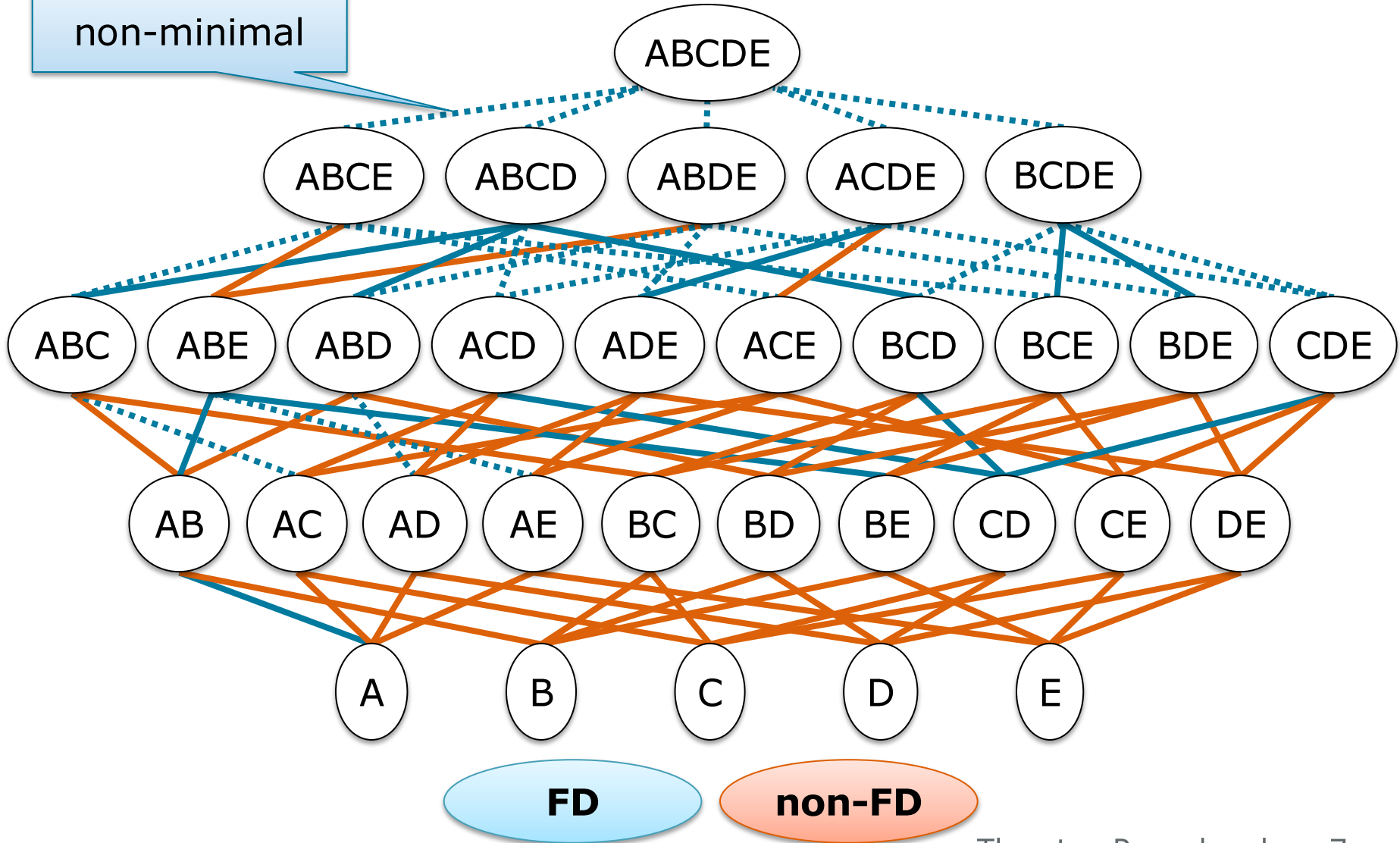


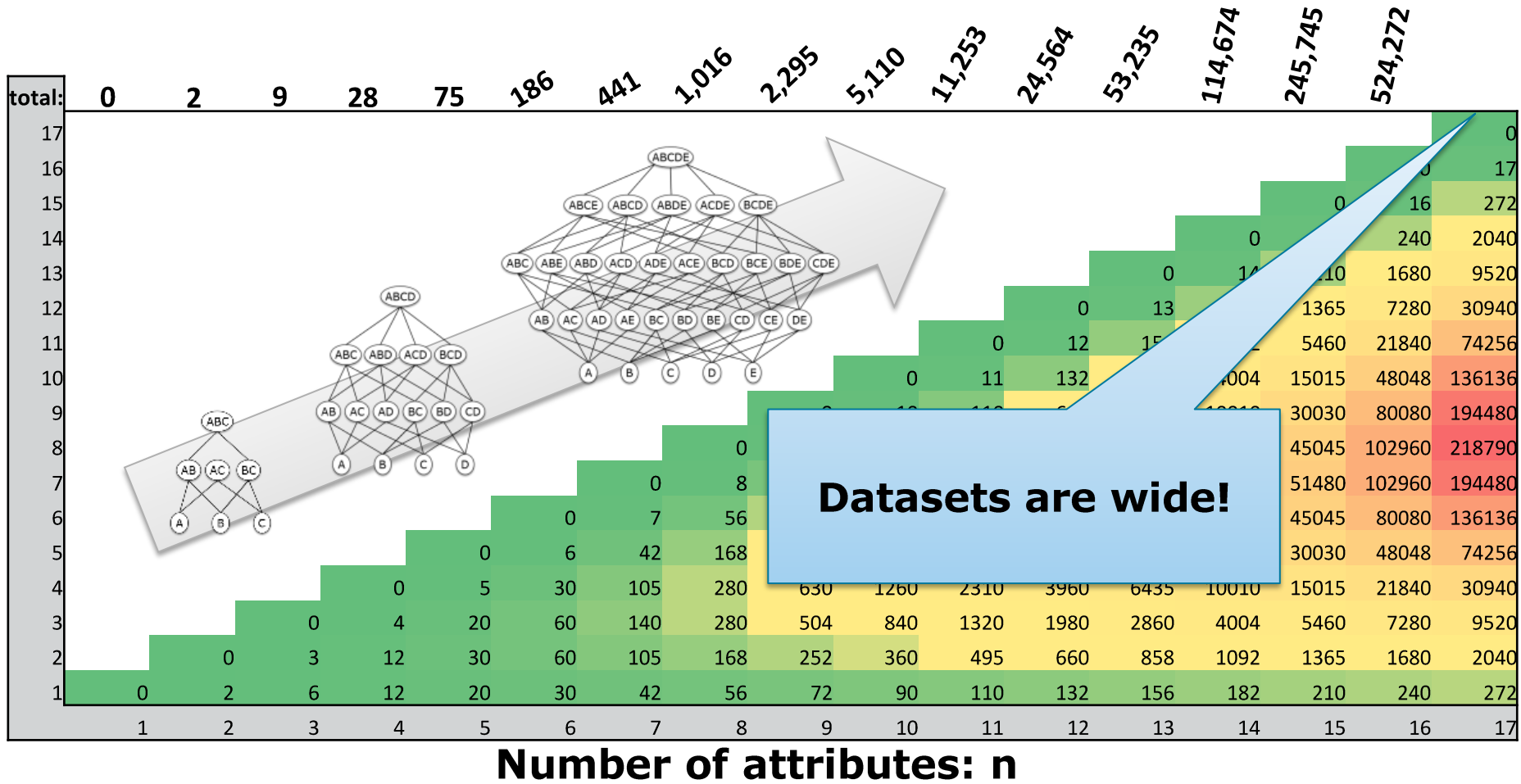






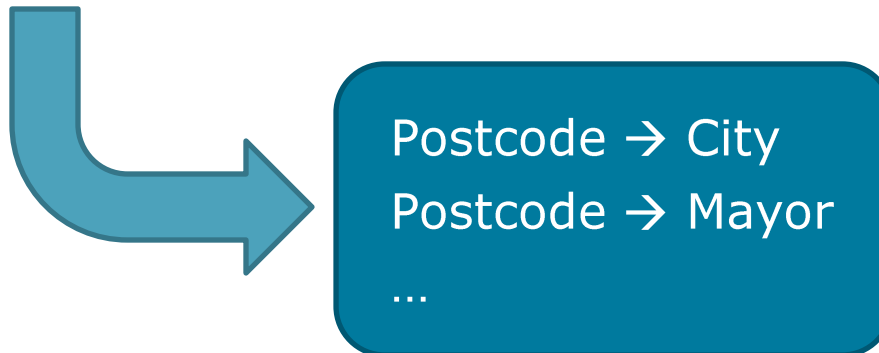
non-minimal



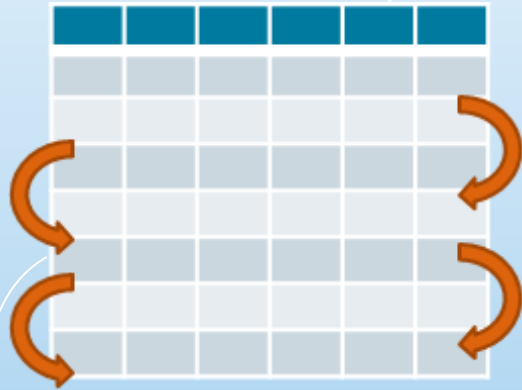
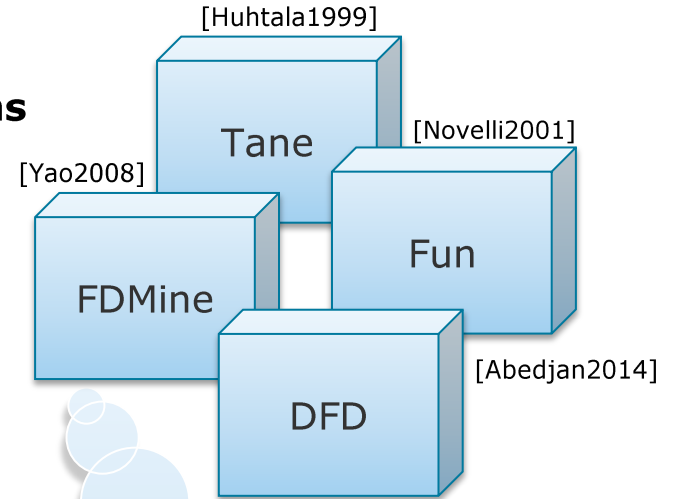
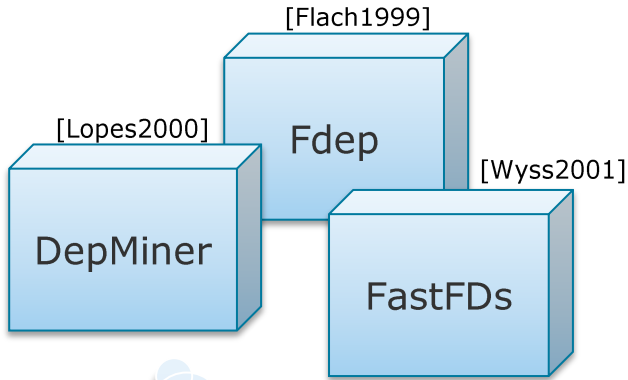


<u>Name</u>	<u>Surname</u>	<u>Postcode</u>	<u>City</u>	<u>Mayor</u>
Thomas	Miller	14482	Potsdam	Jakobs
Sarah	Miller	14482	Potsdam	Jakobs
Peter	Smith	60329	Frankfurt	Feldmann
Jasmine	Cone	01069	Dresden	Orosz
Thomas	Cone	14482	Potsdam	Jakobs
Mike	Moore	60329	Frankfurt	Feldmann

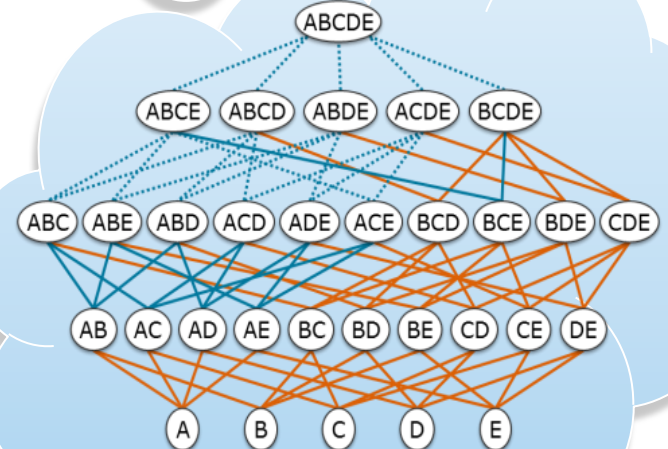
- Surname, Postcode, City, Mayor ↗ Name
- Name, Postcode, City, Mayor ↗ Surname
- Surname ↗ Name, Postcode, City, Mayor



FD Discovery Algorithms



column-efficient

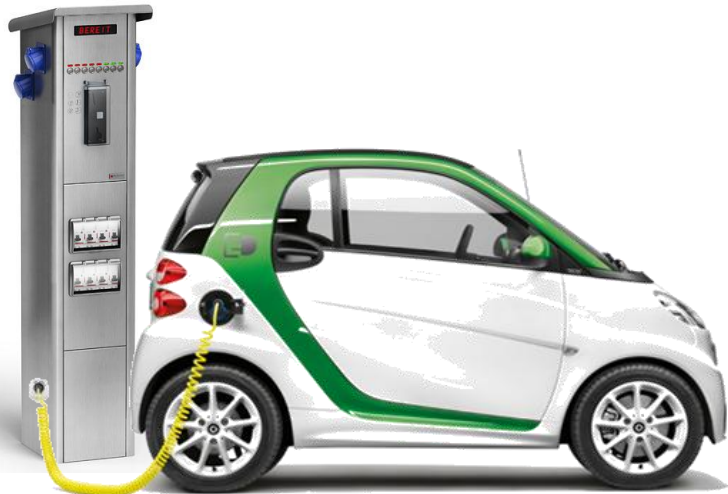
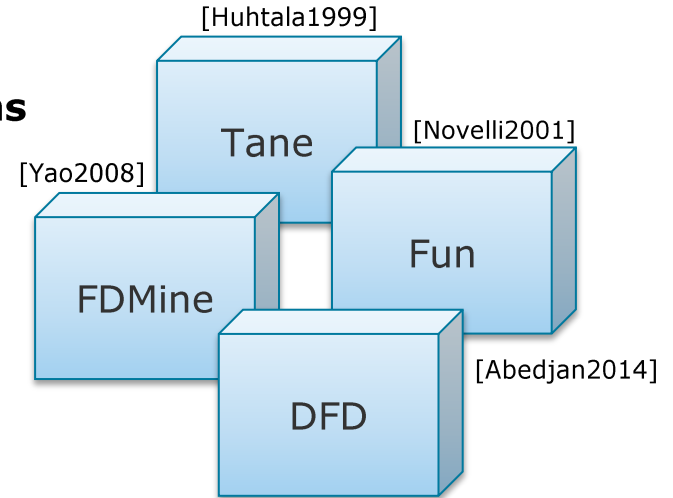
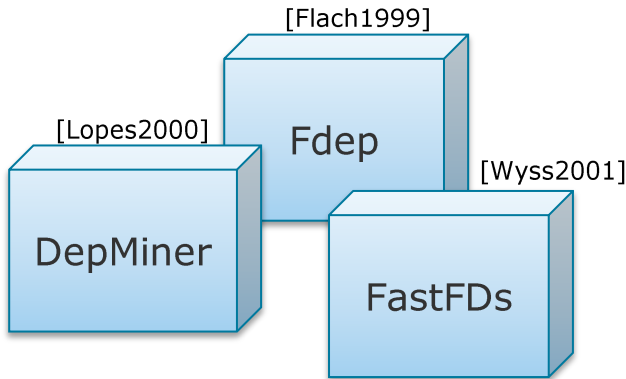


row-efficient



Idea Going Hybrid

FD Discovery Algorithms



© <http://4electric.de>

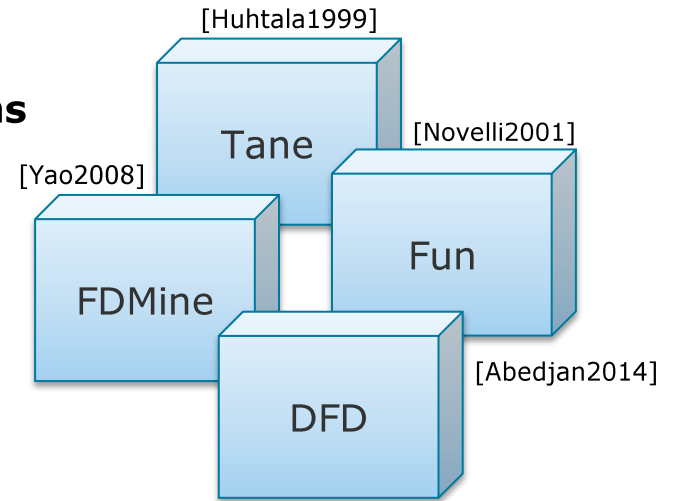
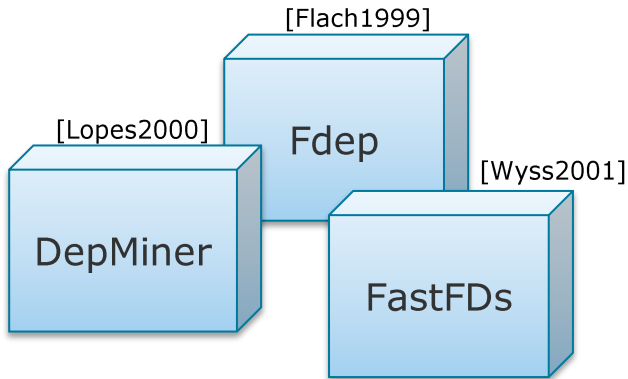


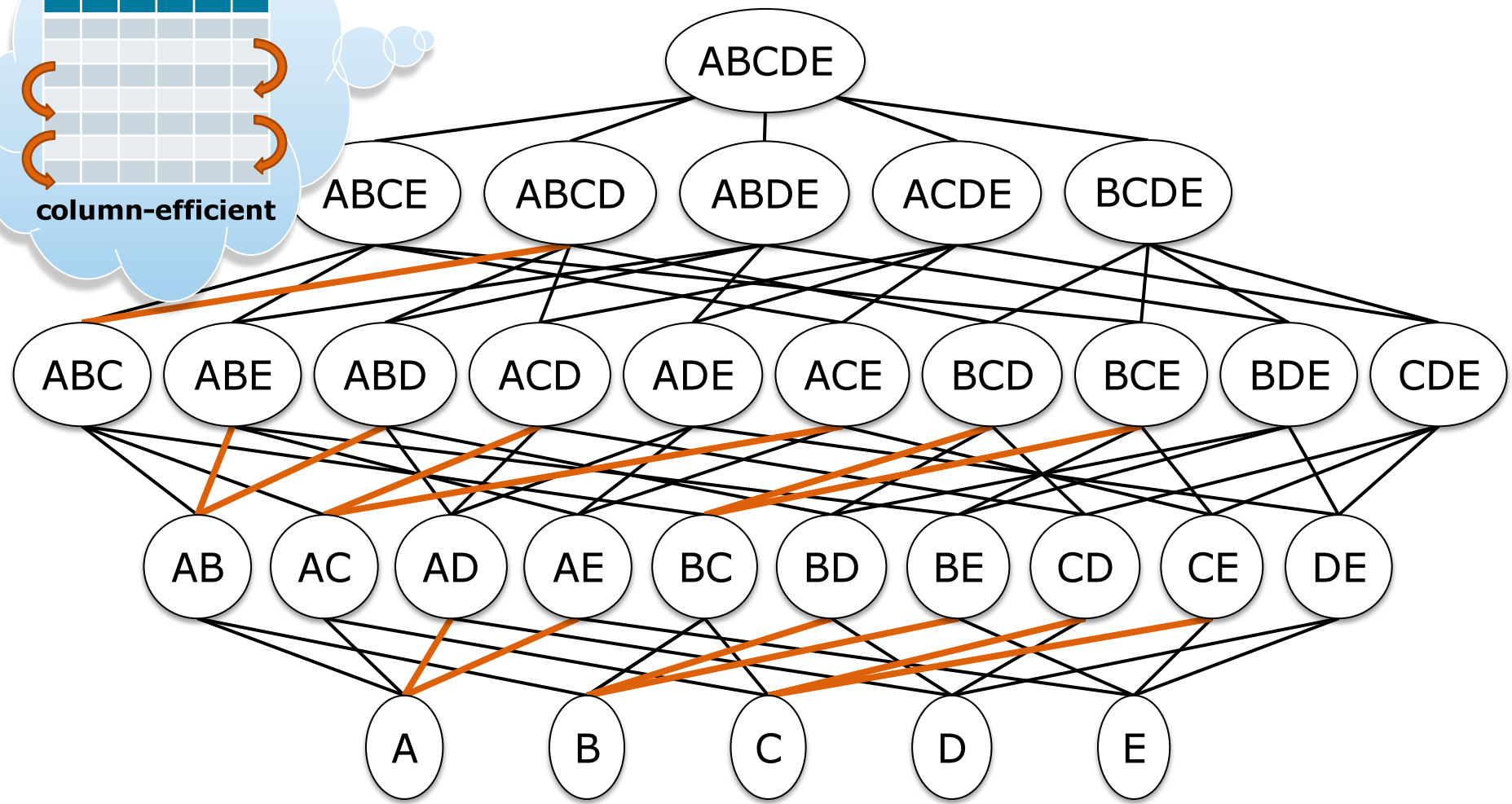
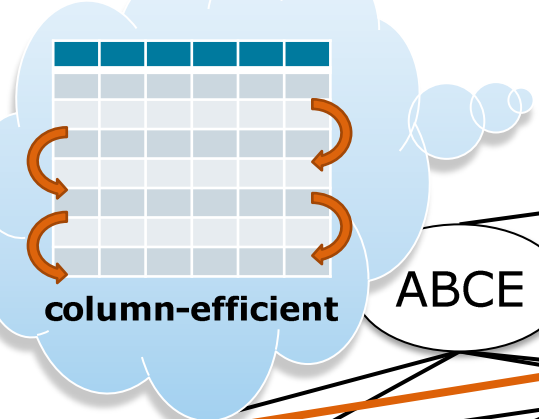
© <http://indiancarsbikes.in>



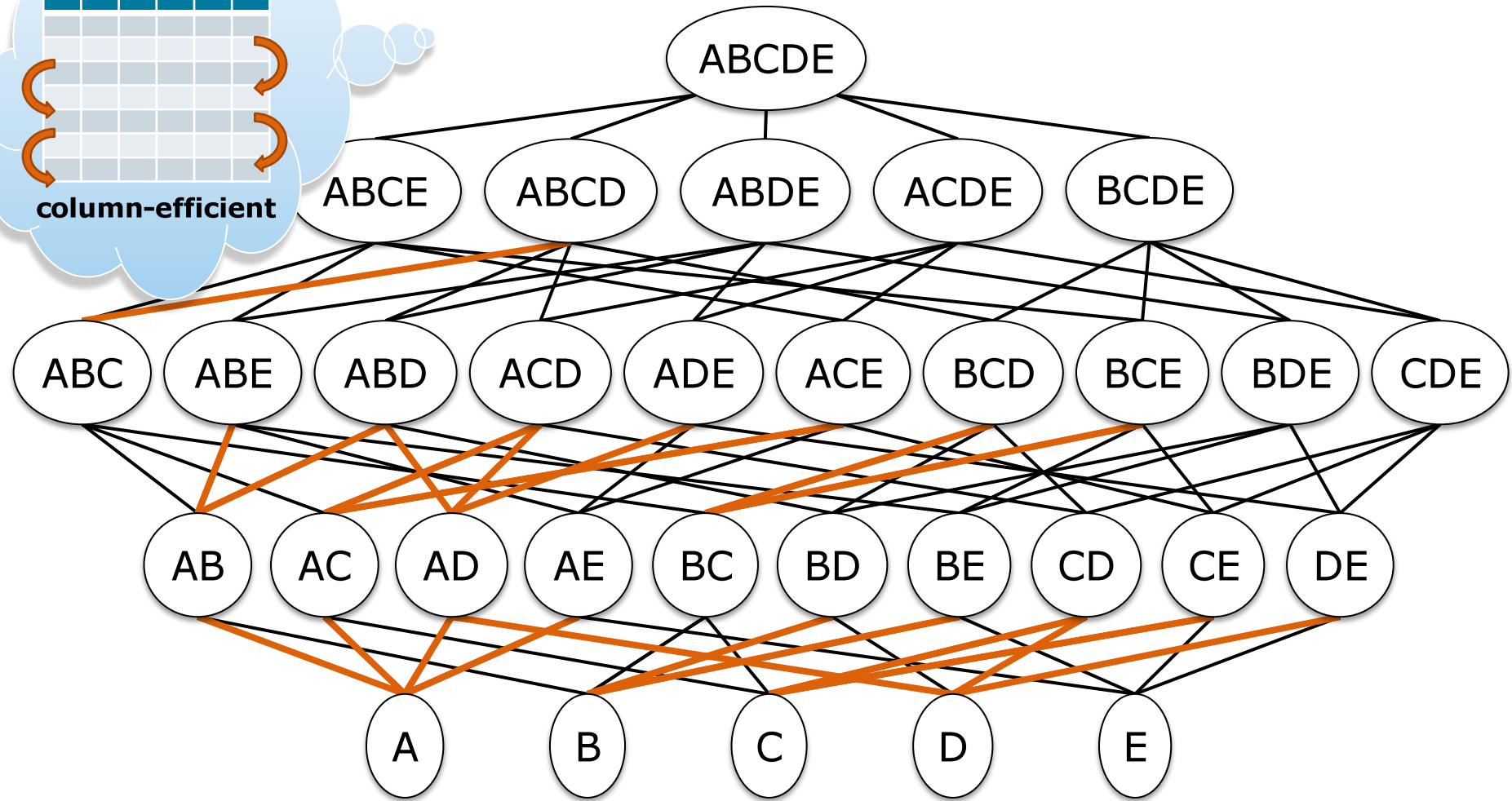
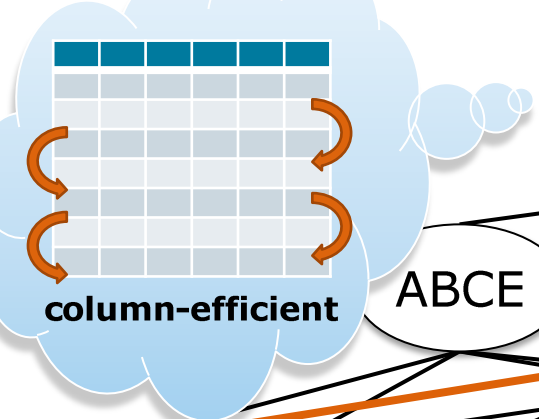
Idea Going Hybrid

FD Discovery Algorithms

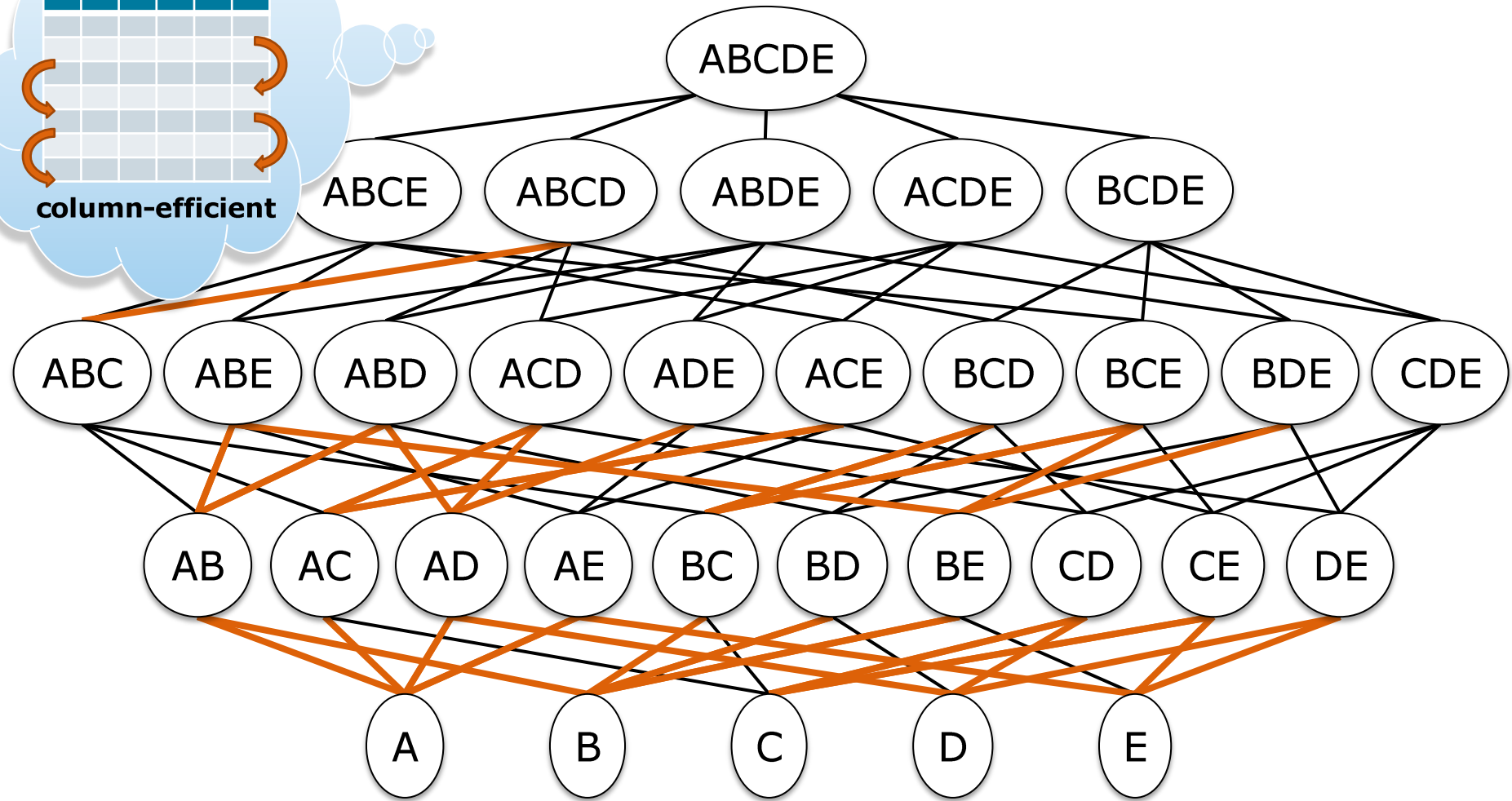
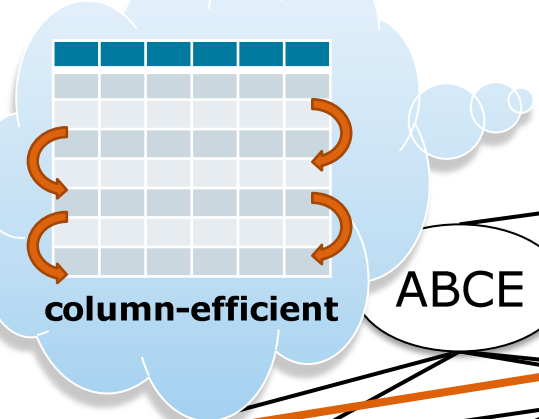




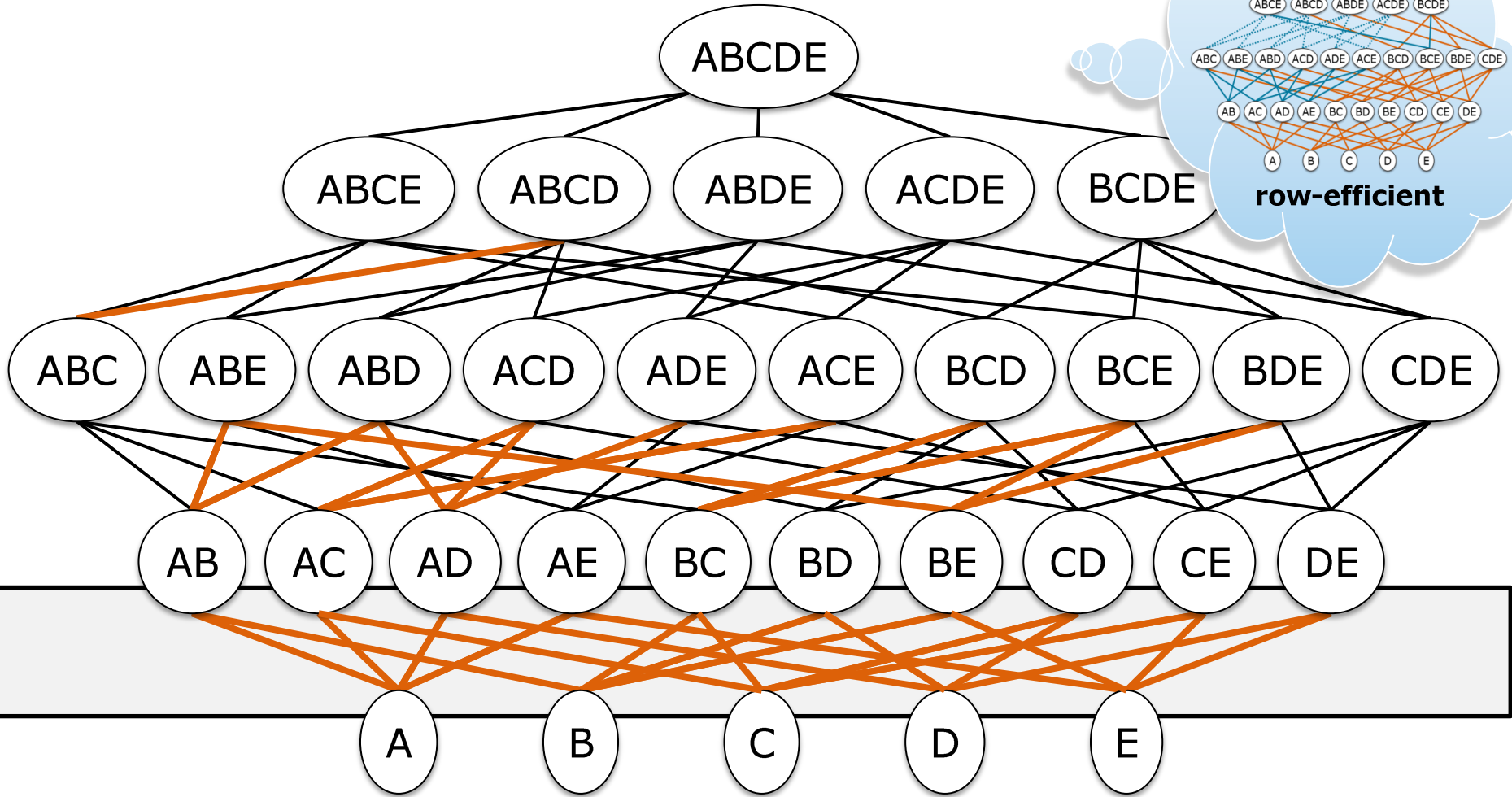
ABC ↗ DE

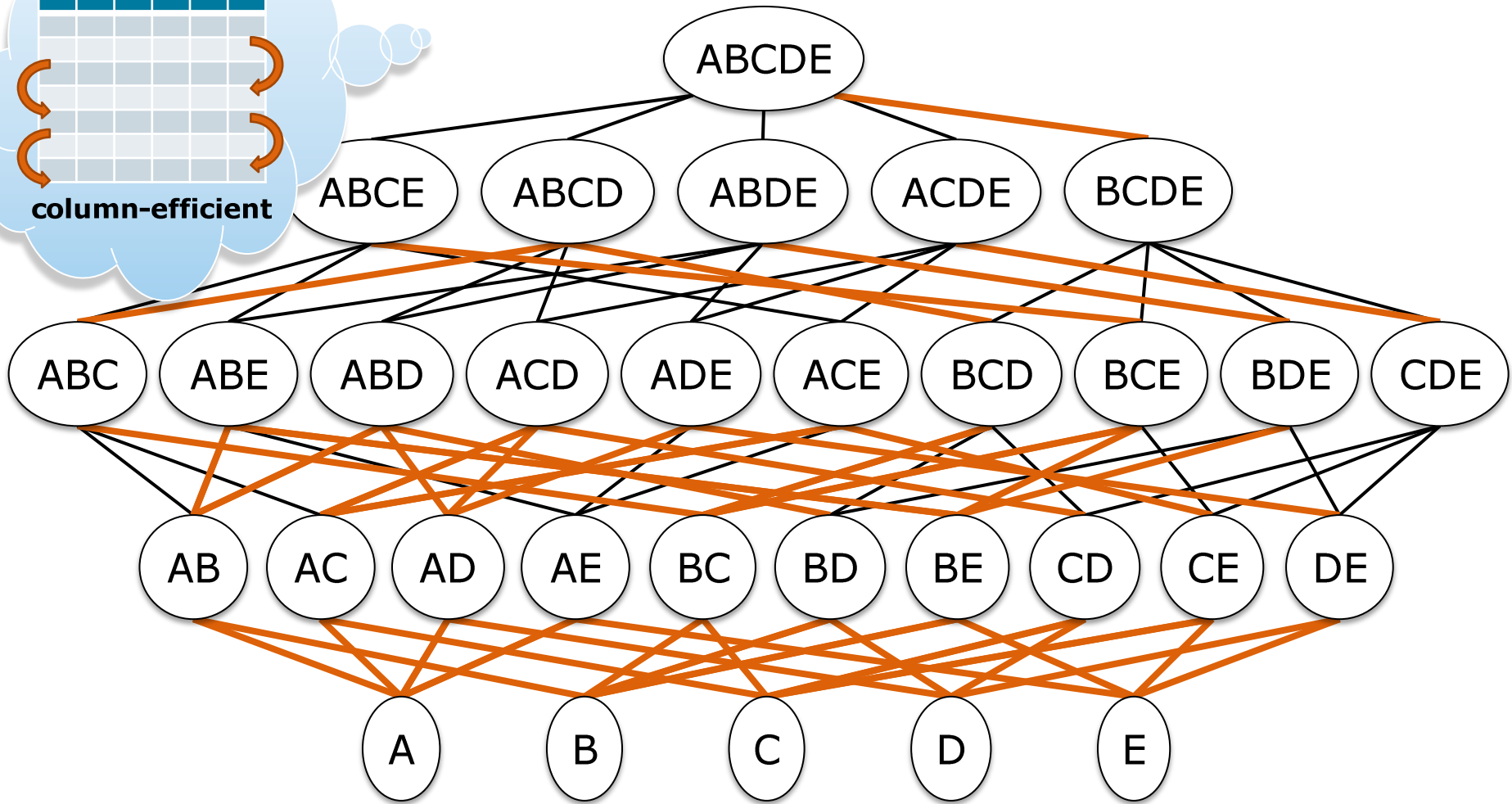
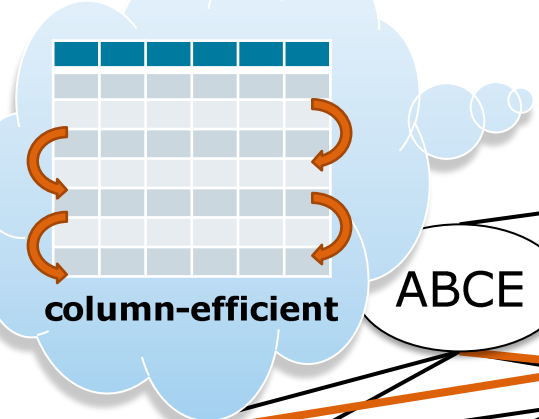


AD \rightarrow BCE

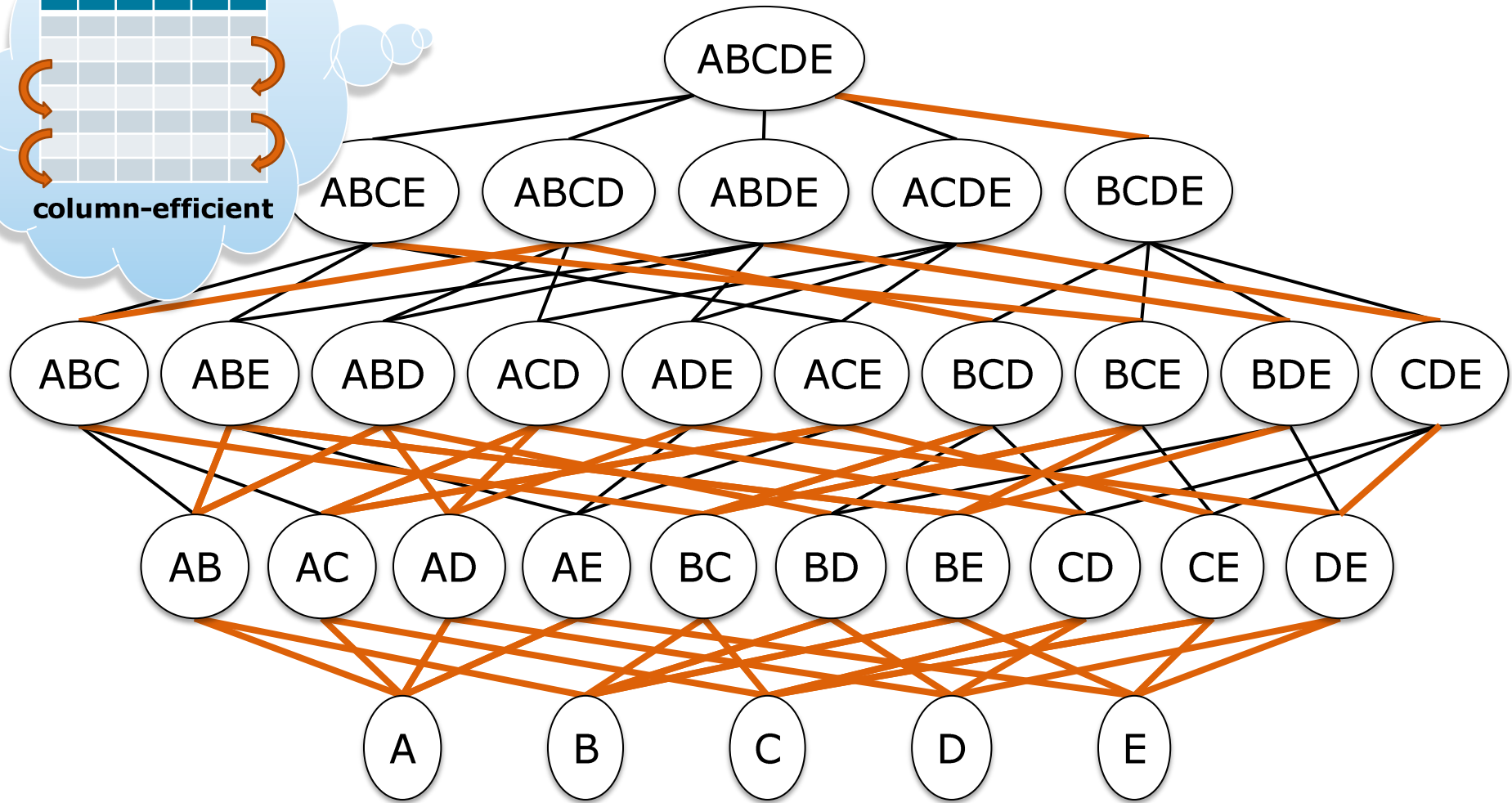
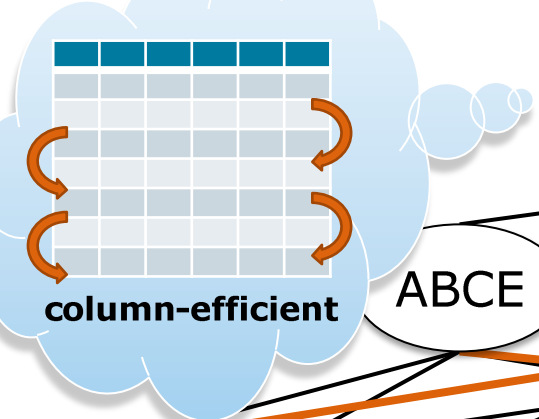


BE ↗ ACD

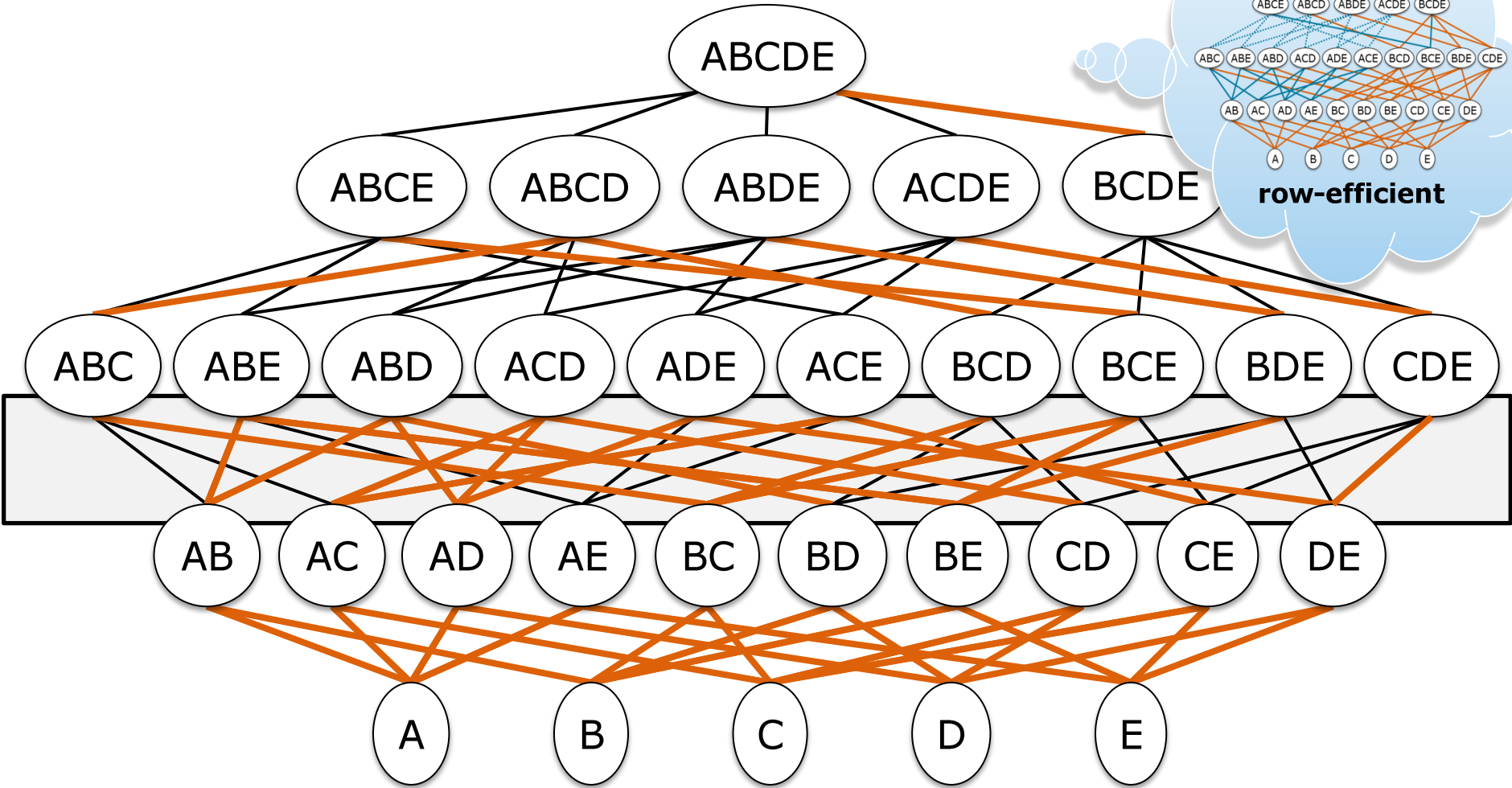


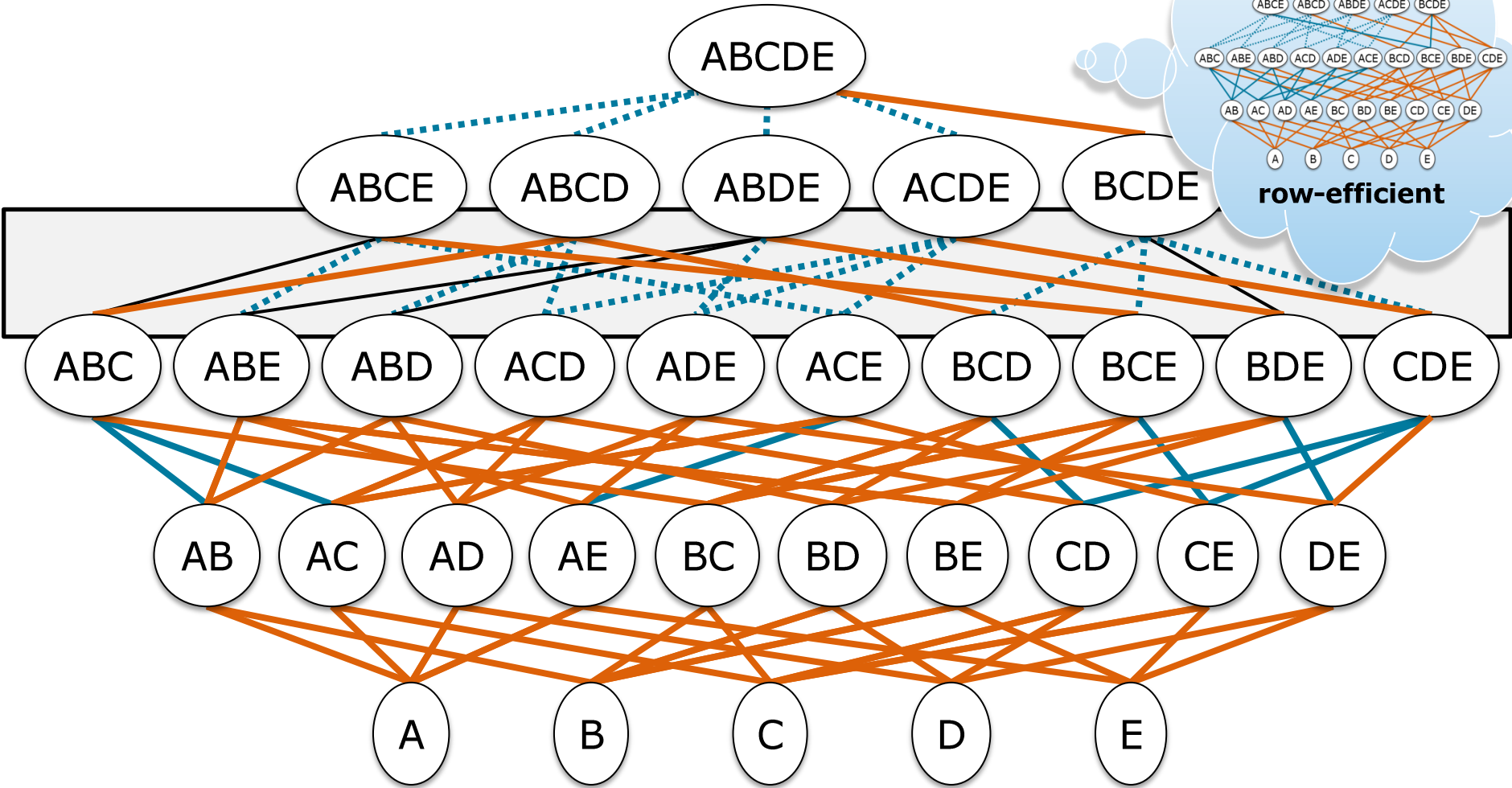


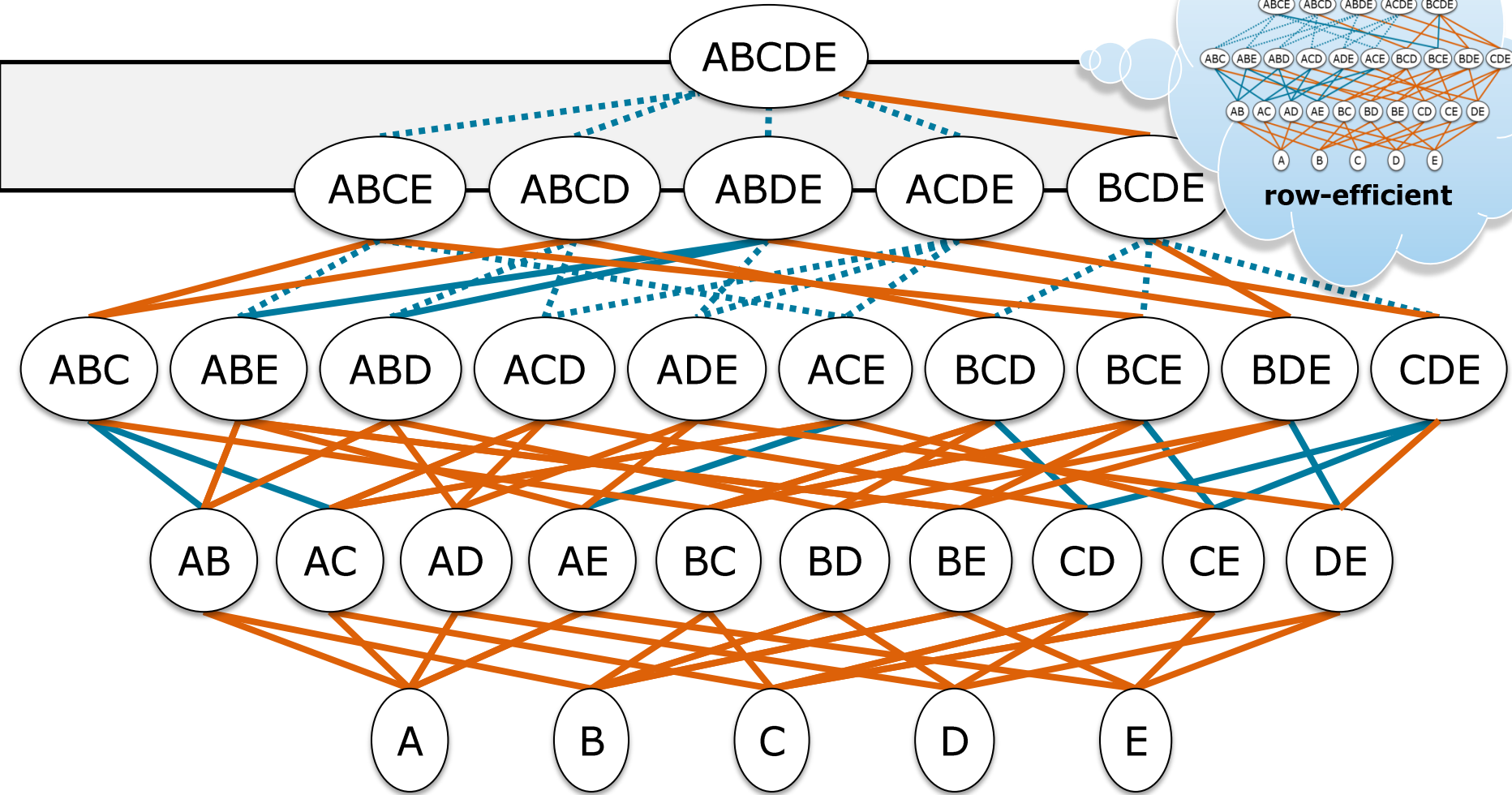
BCDE \rightarrow A

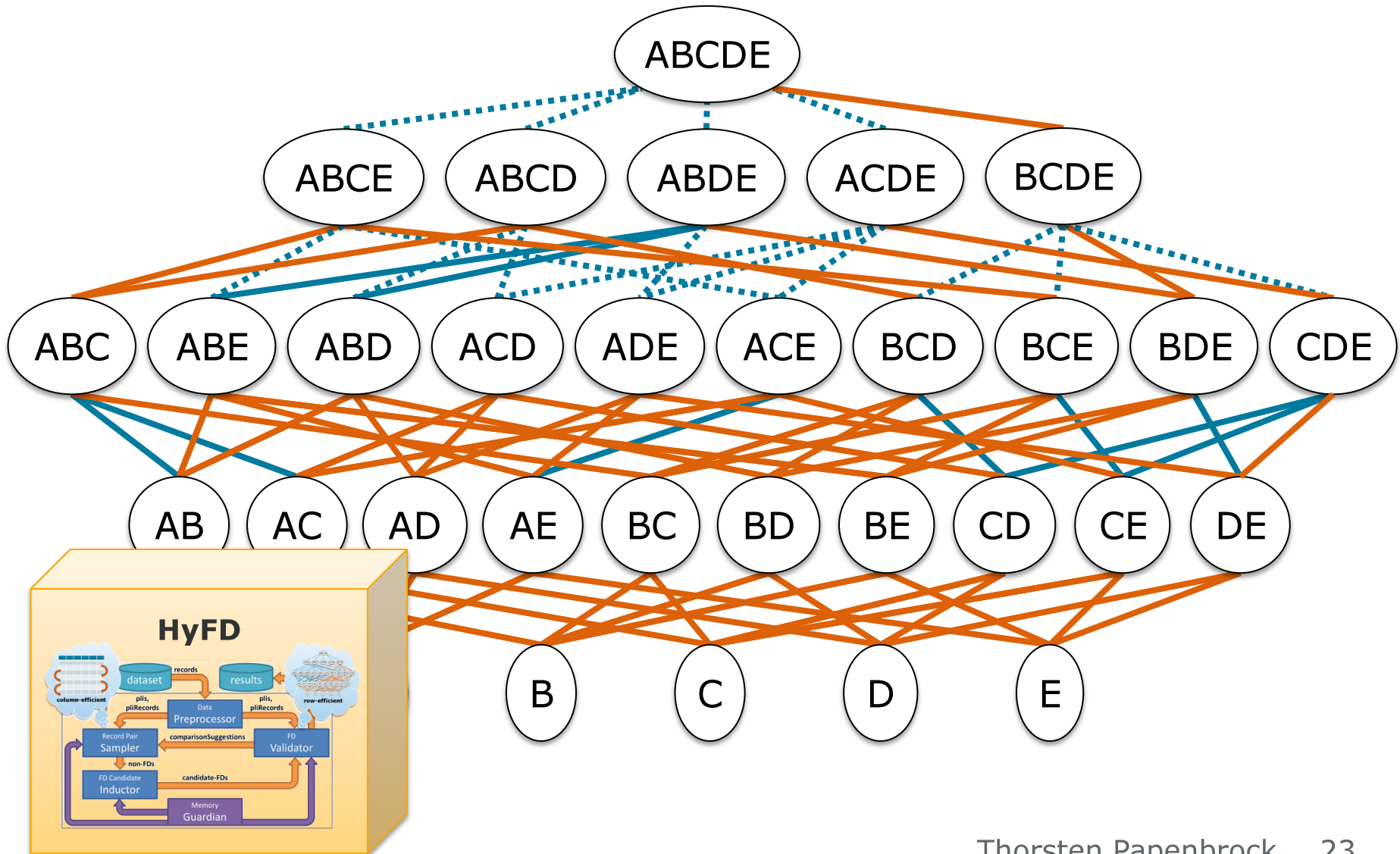


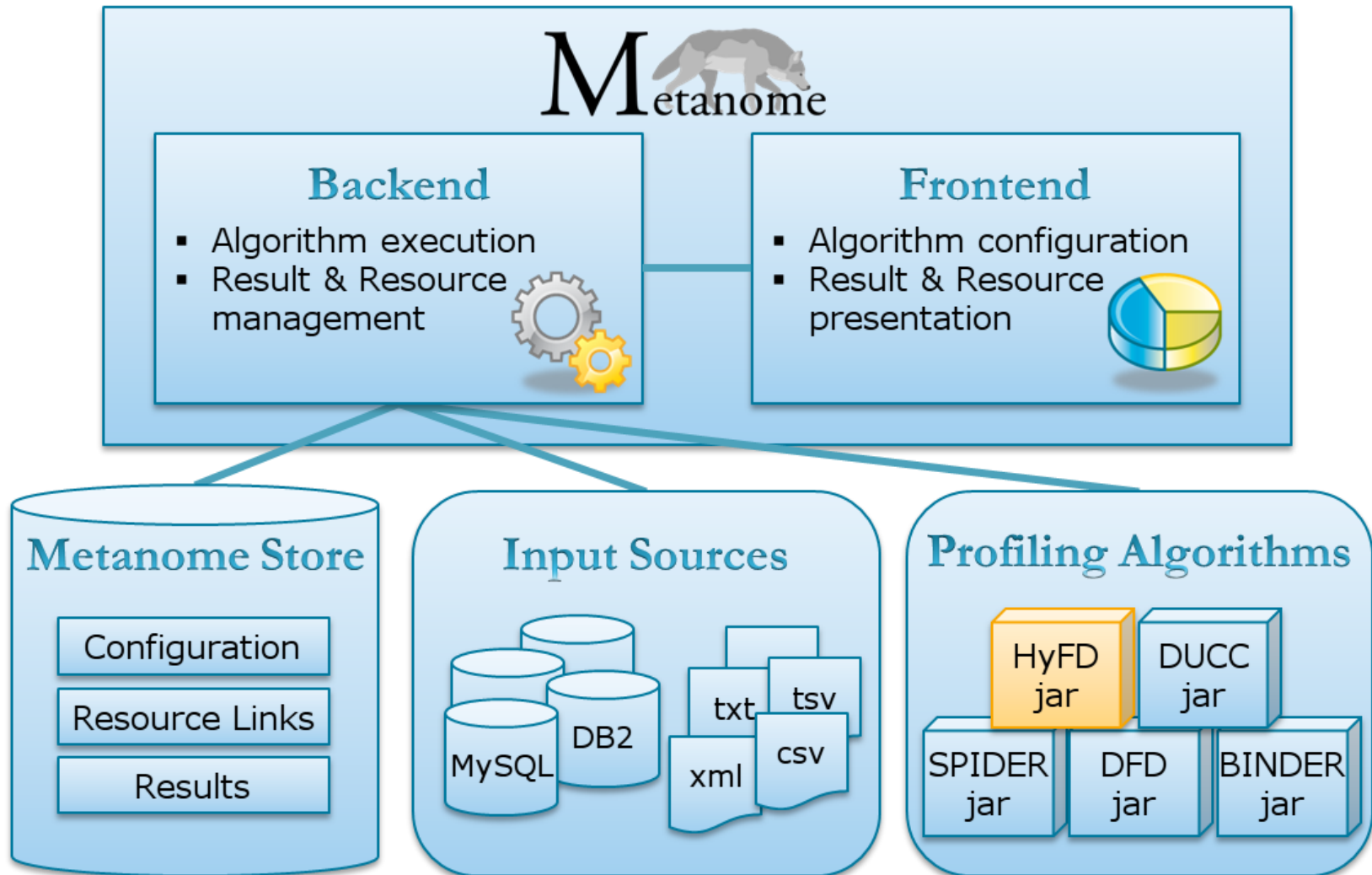
DE \rightarrow C



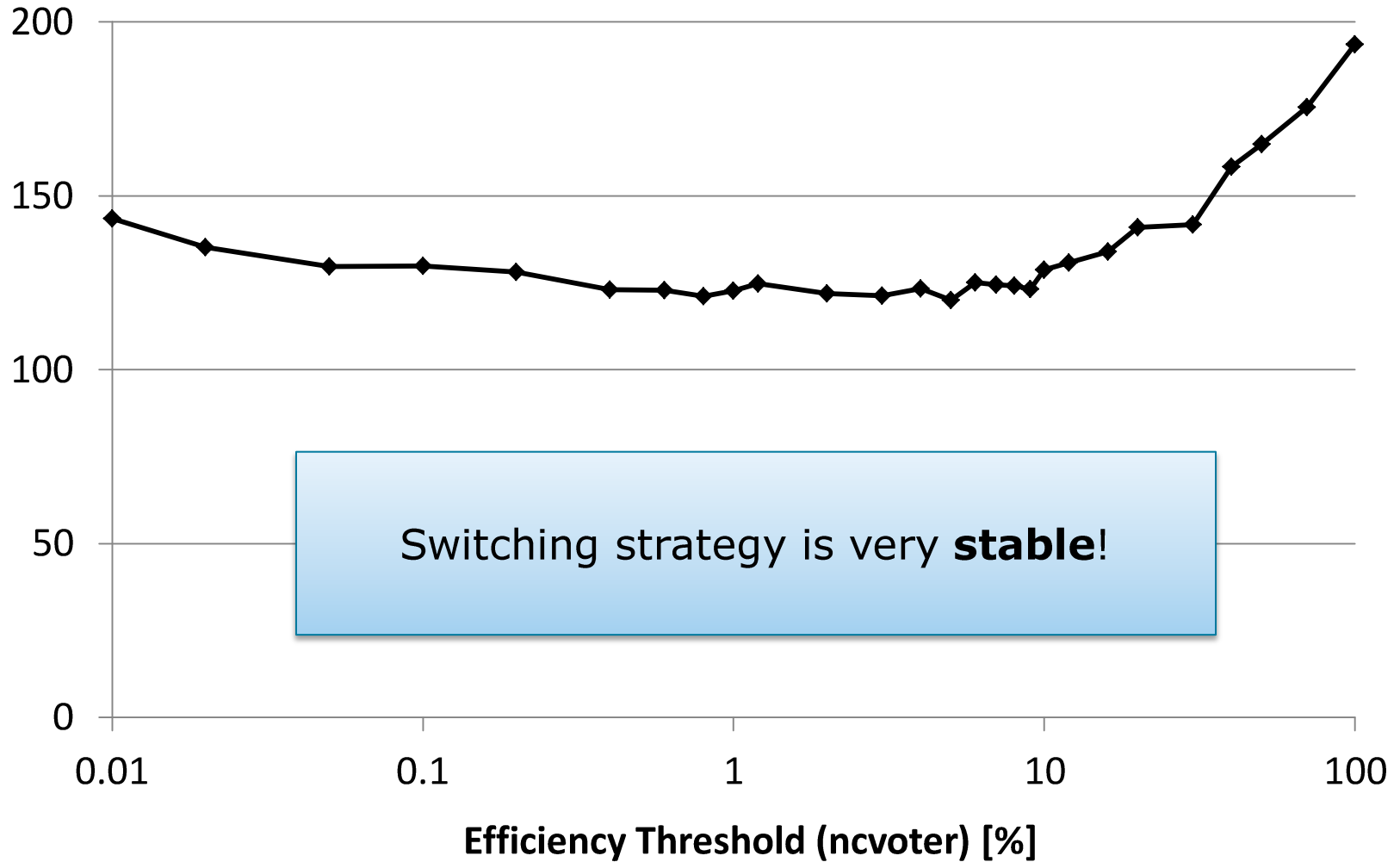








www.metanome.de



◆ HyFD

Dataset	Columns [#]	Rows [#]	Size [KB]	FDs [#]	TANE [7]	FUN [15]	FD_MINE [21]	DFD [1]	DEP-MINER [12]	FASTFDs [20]	FDEP [6]	HyFD
iris	5	150	5	4	1.1	0.1	0.2	0.2	0.2	0.2	0.1	0.1s
balance-scale	5	625	7	1	1.2	0.1	0.2	0.3	0.3	0.3	0.2	0.1s
chess	7	28,056	519	1	2.9	1.1	3.8	1.0	174.6	164.2	125.5	0.2s
abalone	9	4,177	187	137	2.1	0.6	1.8	1.1	3.0	2.9	3.8	0.2s
nursery	9	1								8.9	46.8	0.5s
breast-cancer	11									1.1	0.5	0.2s
bridges	13									0.6	0.2	0.1s
echocardiogram	13									0.5	0.2	0.1s
adult	14	4								3.8	860.2	1.1s
letter	17	2								5.5	291.3	3.4s
ncvoter	19									1.9	1.1	0.4s
hepatitis	20	155	8	8,250	12.2	175.9	ML	326.7	5576.5	9.5	0.8	0.6s
horse	27	368	25	128,726	457.0	TL	ML	TL	TL	385.8	7.2	7.1s
fd-reduced-30	30	250,000	69,581	89,571	41.1	77.7	ML	TL	377.2	382.4	TL	513.0s
plista	63	1,000	568	178,152	ML	ML	ML	TL	TL	TL	26.9	21.8s
flight	109	1,000	575	982,631	ML	ML	ML	TL	TL	TL	216.5	53.4s
uniprot	223	1,000	2,439	unknown	ML	ML	ML	TL	TL	TL	ML	>5254.7s

Faster on all real-world datasets!

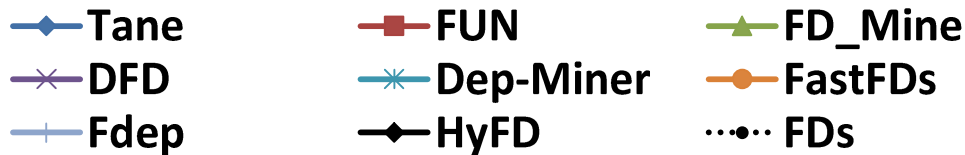
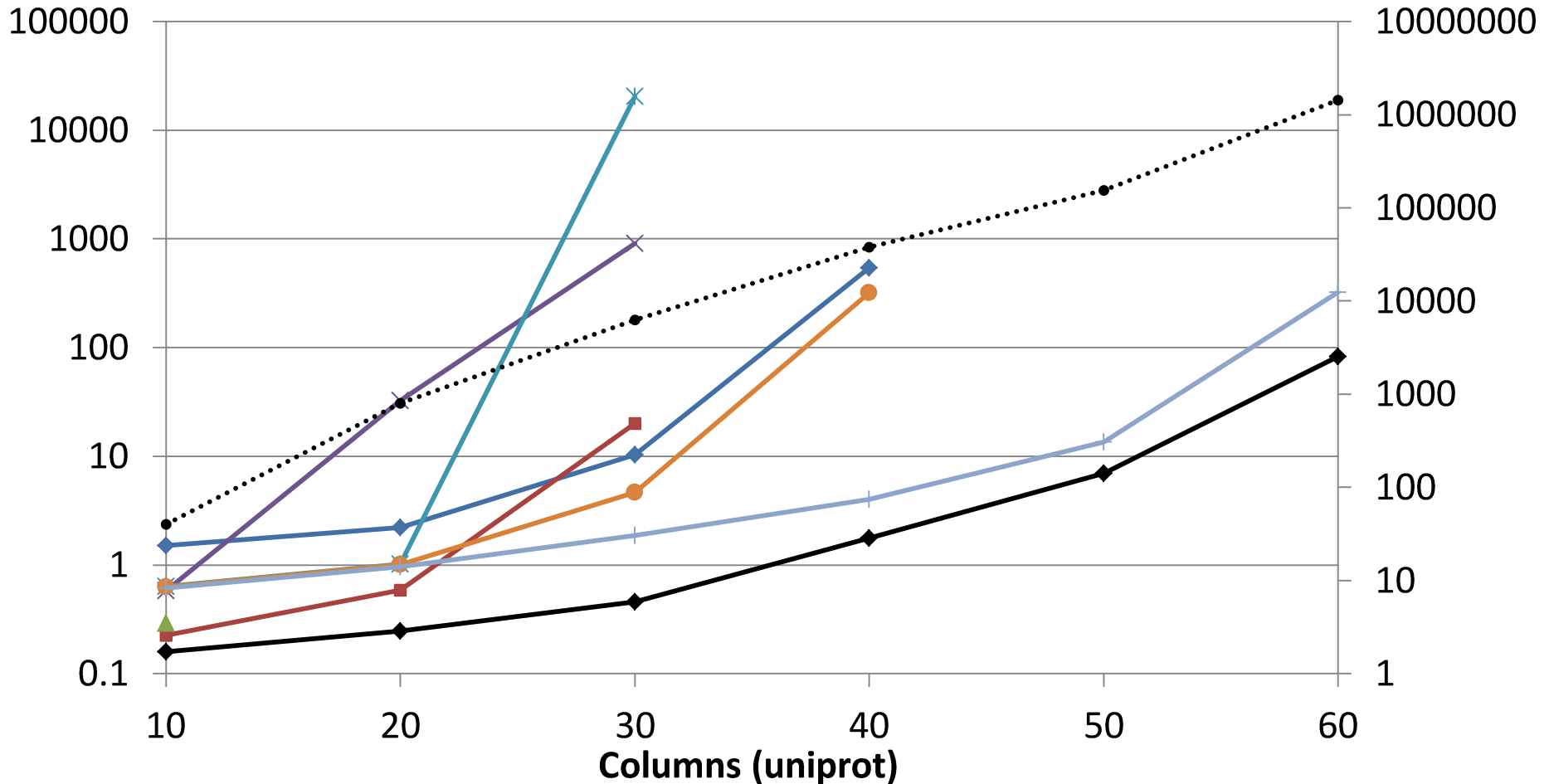
Results larger than 1,000 FDs are only counted

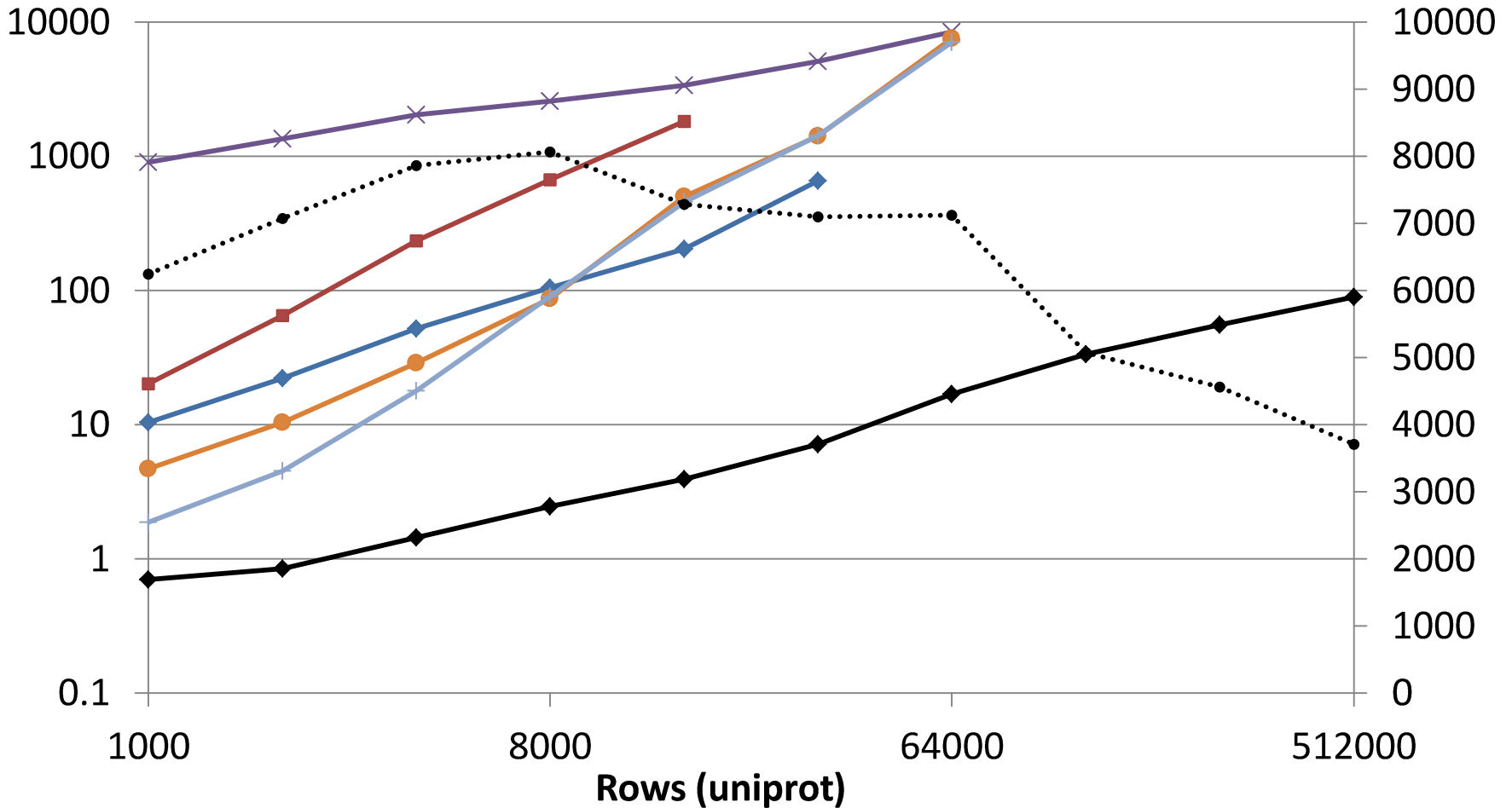
TL: time limit of 4 hours exceeded

ML: memory limit of 100GB exceeded

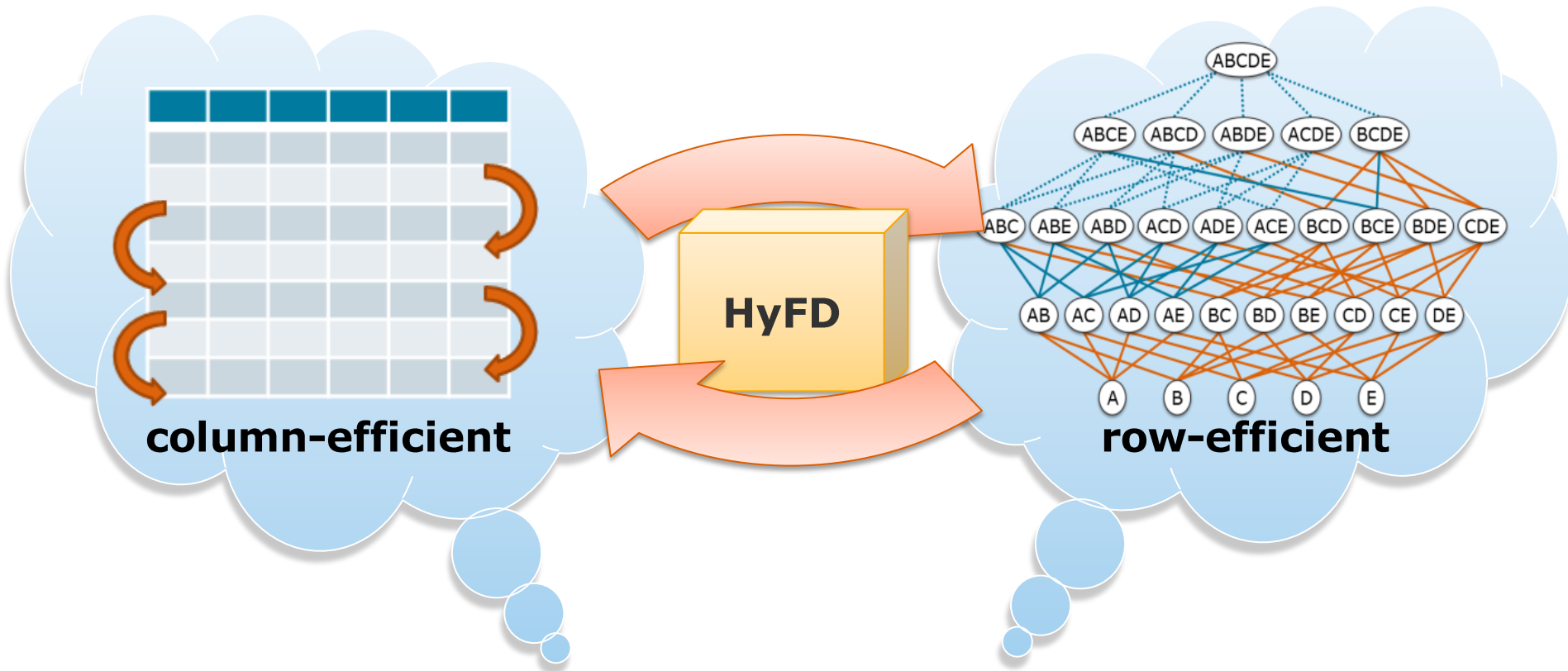
Dataset	Cols [#]	Rows [#]	Size [MB]	FDs [#]	HyFD [s/m/h/d]
TPC-H.lineitem	16	6 m	1.051	4 k	39 m 4 m
PDB.POLY_SEQ					
PDB.ATOM_SITE					
SAP_R3.ZBC00DT					
SAP_R3.ILOA					
SAP_R3.CE4HI01					
NCVoter.statewide					
CD.cd	107	10 k	5	36 k	5 s 3 m

Larger datasets than ever before!





- ◆ Tane
- ◆ FUN
- ◆ FD_Mine
- ◆ DFD
- ◆ Dep-Miner
- ◆ FastFDs
- ◆ Fdep
- ◆ HyFD
- ◆ FDs



A Hybrid Approach to Functional Dependency Discovery

Thorsten Papenbrock and Felix Naumann