



Divide & Conquer-based Inclusion Dependency Discovery

Thorsten Papenbrock, Sebastian Kruse, Jorge-Arnulfo Quianè-Ruiz, Felix Naumann

Motivation

Unary Inclusion Dependencies (INDs)

Book

Title	Author	Price	Pages	Published
Database Systems	Ullman	214	1203	2007
Algorithms in Java	Sedgewick	130	768	2002
3D Computer Graphics	Watt	20	570	1999

Name \subseteq Title

Lending

ID	Name	Location	Student	Course
42	Database Systems	A-1.2	Miller	DBS 1
88	Database Systems	B-2.2	Miller	PT 1
73	Database Systems	A-1.2	Smith	DPDC
69	Algorithms in Java	C-E.1	Miller	PT 1
13	Algorithms in Java	C-E.1	Smith	DPDC

Motivation

N-ary Inclusion Dependencies (INDs)

Student

Name	Lecture	Credit	Semester	Verified
Miller	DBS 1	20	2	false
Miller	PT 1	15	2	false
Smith	DPDC	10	6	true

Student, Course \subseteq Name, Lecture

Lending

ID	Name	Location	Student	Course
42	Database Systems	A-1.2	Miller	DBS 1
88	Database Systems	B-2.2	Miller	PT 1
73	Database Systems	A-1.2	Smith	DPDC
69	Algorithms in Java	C-E.1	Miller	PT 1
13	Algorithms in Java	C-E.1	Smith	DPDC

Motivation

Necessity for IND Discovery

Independend Creation

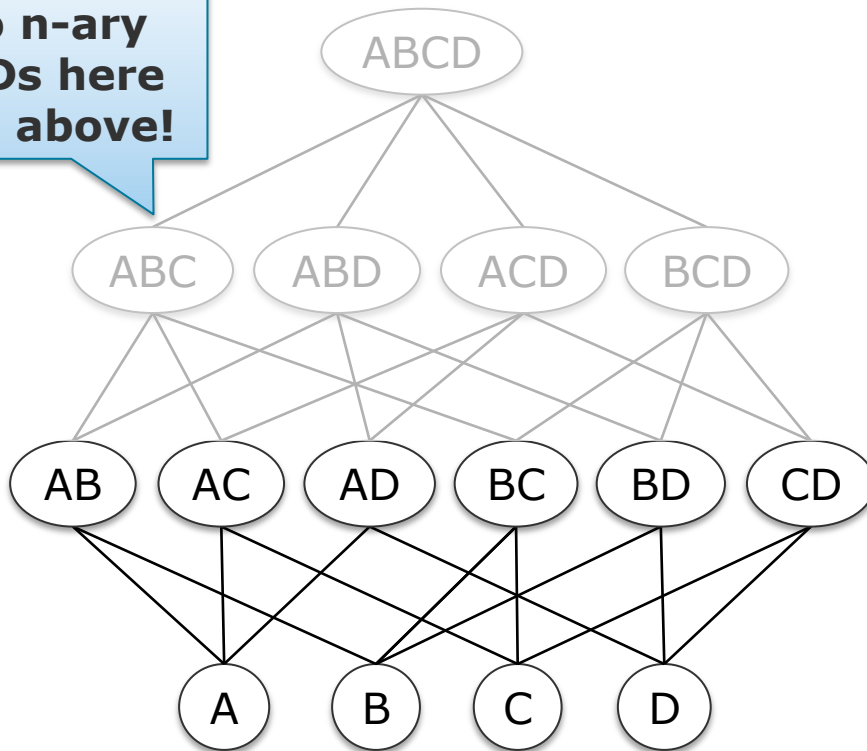


Technical Limitations



Motivation IND Detection Search Space

No n-ary
 INDs here
 and above!



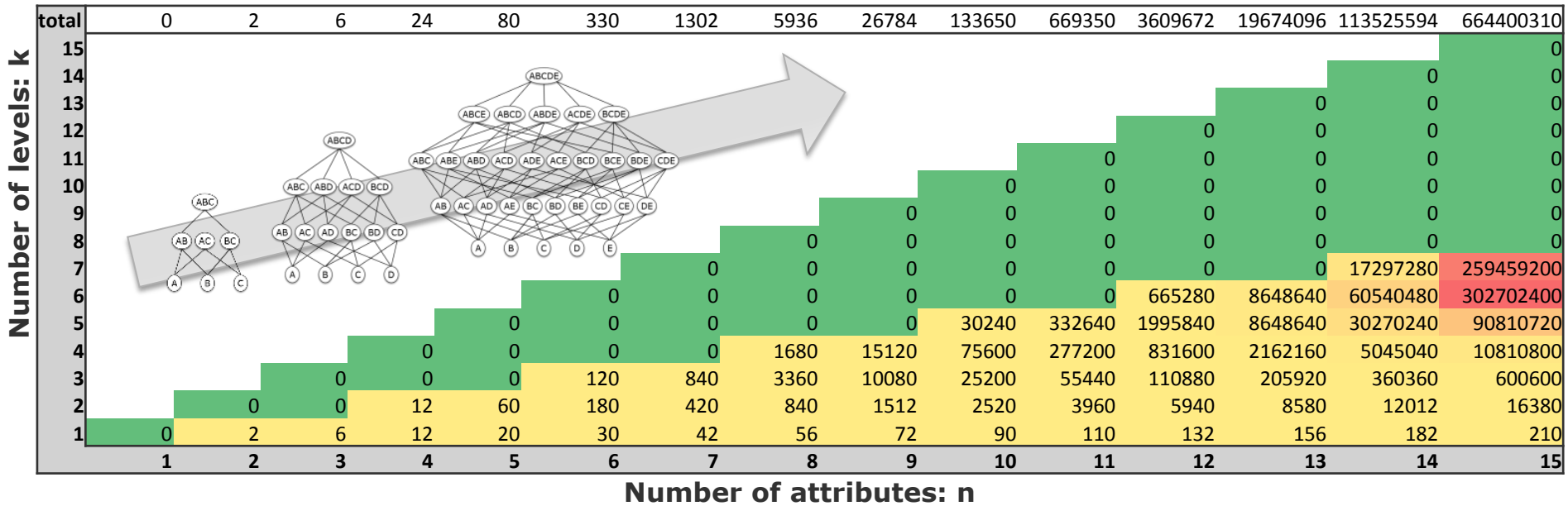
Test every combination
 with all other combina-
 tions of same size!

- AB ⊆ CD AB ⊆ DC
 - AC ⊆ BD AC ⊆ DB
 - AD ⊆ BC AD ⊆ CB
 - BC ⊆ AD BC ⊆ DA
 - BD ⊆ AC BD ⊆ CA
 - CD ⊆ AB CD ⊆ BA
-
- A ⊆ B A ⊆ C A ⊆ D
 - B ⊆ A B ⊆ C B ⊆ D
 - C ⊆ A C ⊆ B C ⊆ D
 - D ⊆ A D ⊆ B D ⊆ C

Motivation

IND Detection Search Space

IND candidates



Unary IND detection:

$$O(n^2)$$

for n attributes

N-ary IND detection:

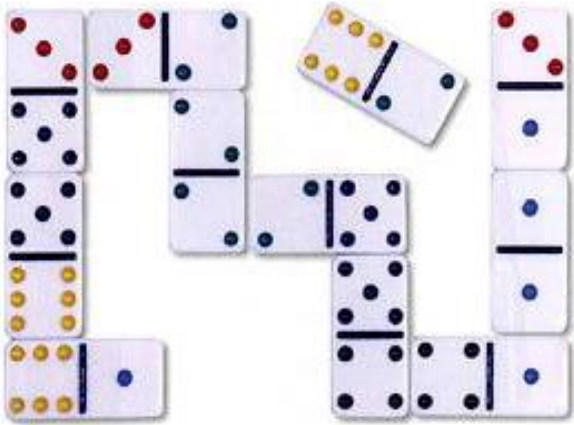
$$O(2^n \cdot n!)$$

for n attributes

Motivation Related Work

- **Bell:** S. Bell and P. Brockhausen. *Discovery of data dependencies in relational databases*. Technical report, Universität Dortmund, 1995
- **Koeller:** A. Koeller and E. A. Rundensteiner. *Discovery of high-dimensional inclusion dependencies*. In ICDE, 2002
- **Zigzag:** F. D. Marchi and J.-M. Petit. *Zigzag: A new algorithm for mining large inclusion dependencies in databases*. In ICDM, 2003
- **SPIDER:** J. Bauckmann, U. Leser, and F. Naumann. *Efficiently computing inclusion dependencies for schema discovery*. In ICDE Workshops, 2006 (best for **unary** INDs)
- **MIND:** F. D. Marchi, S. Lopes, and J.-M. Petit. *Unary and n-ary inclusion dependency discovery in relational databases*. JIIS, 32:53–73, 2009 (best for **n-ary** INDs)
- **Clim:** F. D. Marchi. *CLIM: Closed inclusion dependency mining in databases*. In ICDM Workshops, 2011

Unary IND Discovery



N-ary IND Discovery



Experimental Evaluation



Unary IND Discovery The Approach

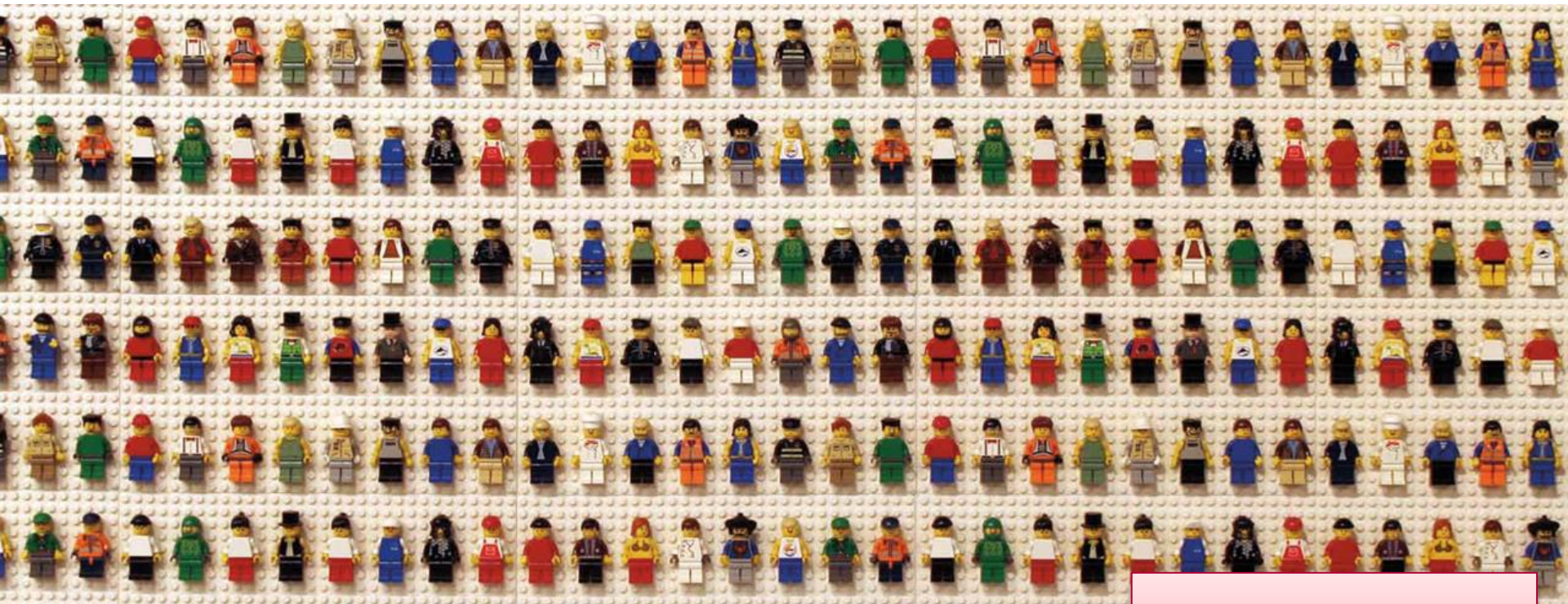


Unary IND Discovery The Approach



No preparation

Unary IND Discovery The Approach



Sorting

Unary IND Discovery The Approach

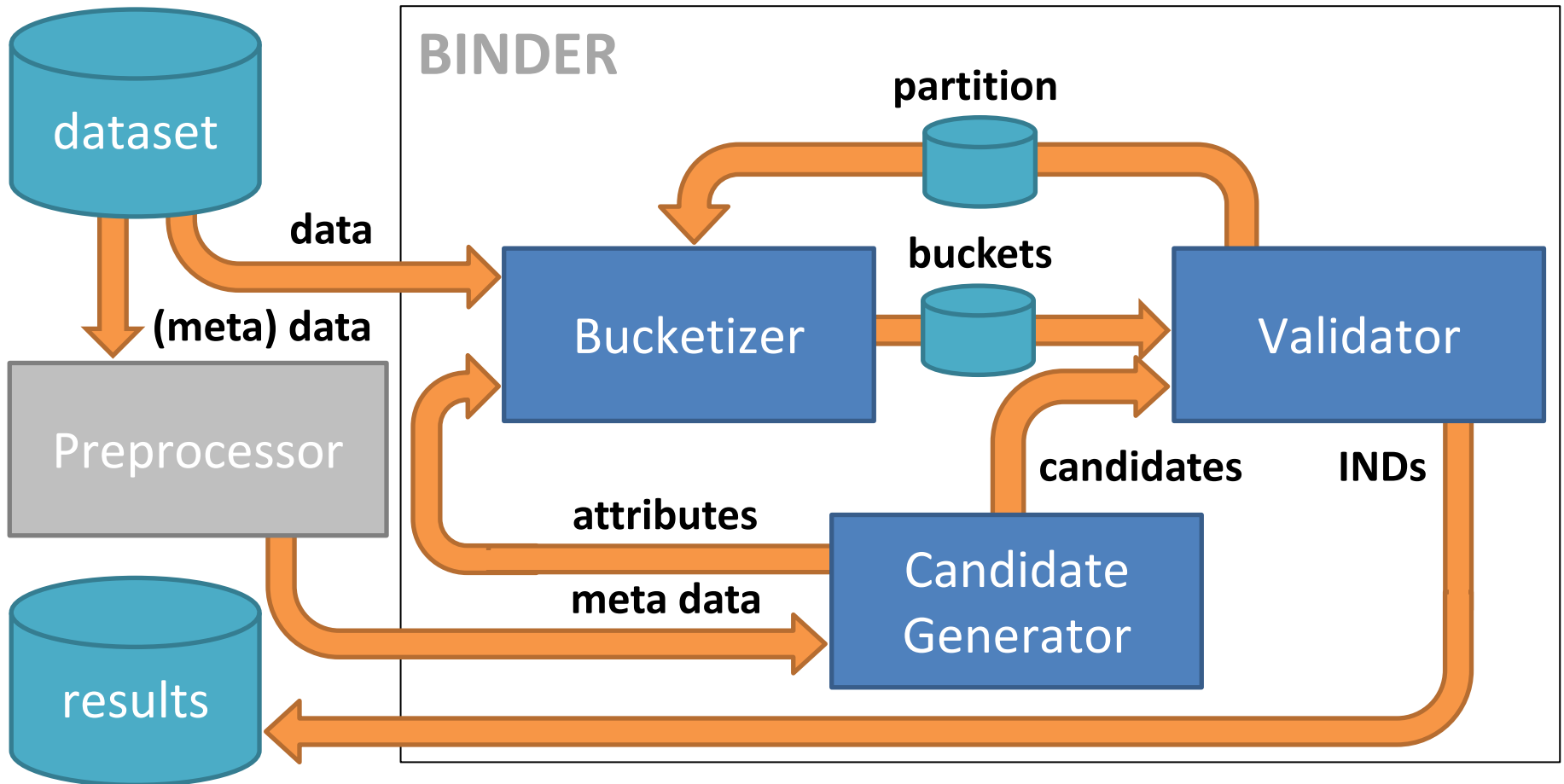


Bucketing

Unary IND Discovery The Approach



N-ary IND Discovery The BINDER Algorithm

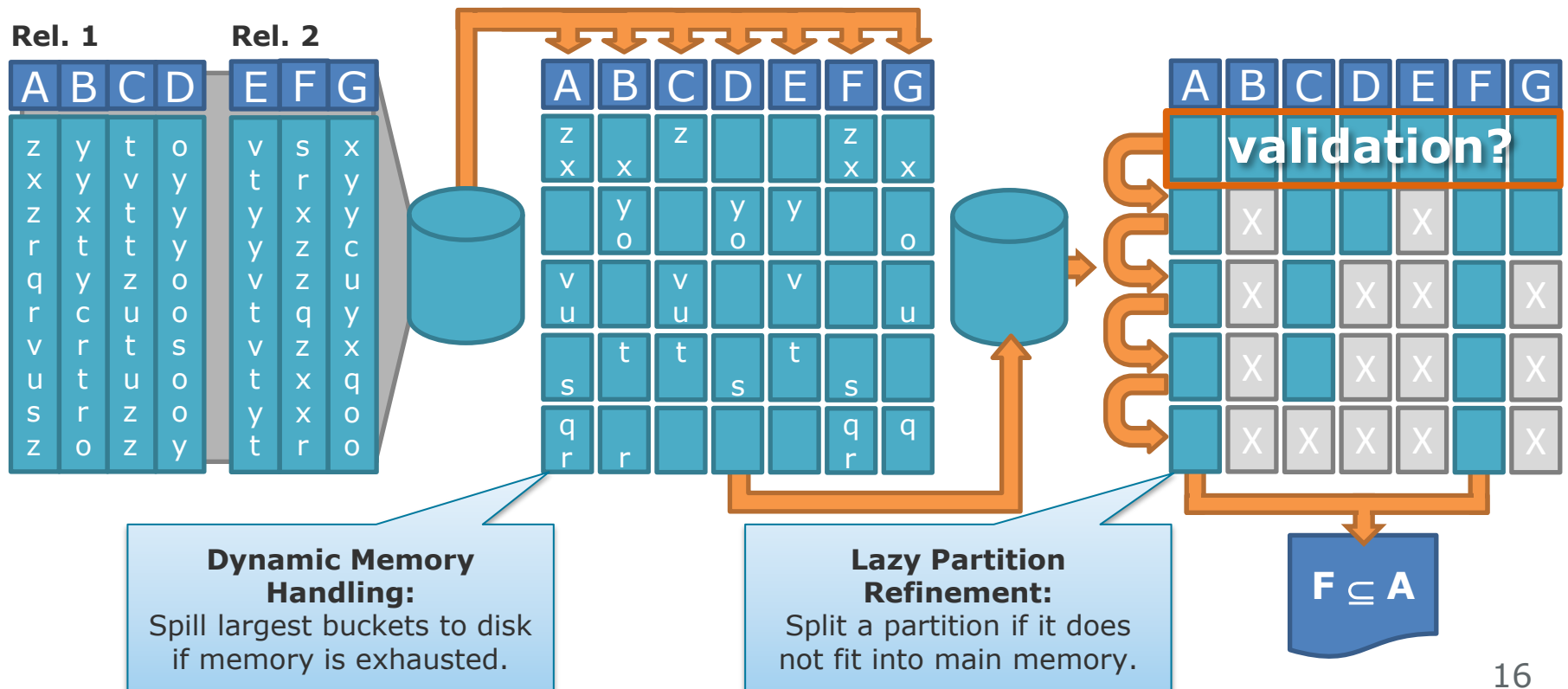


Unary IND Discovery The BINDER Algorithm



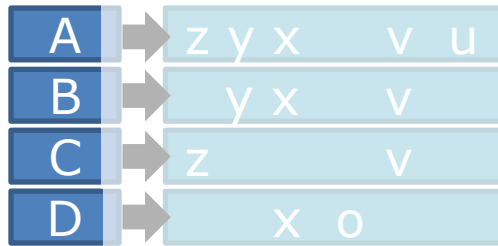
Divide

Conquer

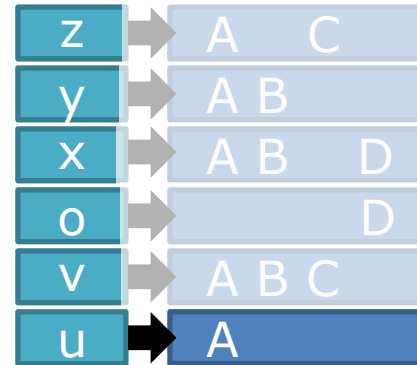


Unary IND Discovery The BINDER Algorithm

attr2value

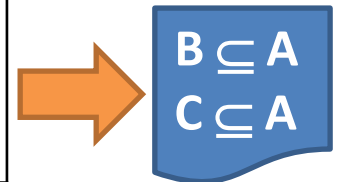


value2attr

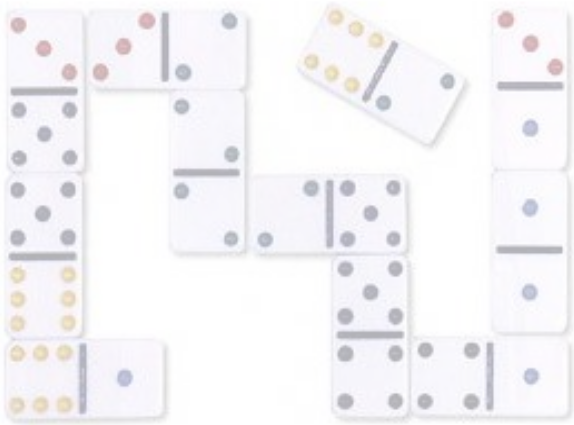


Never tested! →

	A	B	C	D
look up	B,C,D	A,C,D	A,B,D	A,B,C
A → z → A,C	C	A,C,D	A	A,B,C
A → y → A,B	-	A	A	A,B,C
B → x → A,B,D	-	A	A	A,B
B → v → A,B,C	-	A	A	A,B
D → o → D	-	A	A	-



Unary IND Discovery



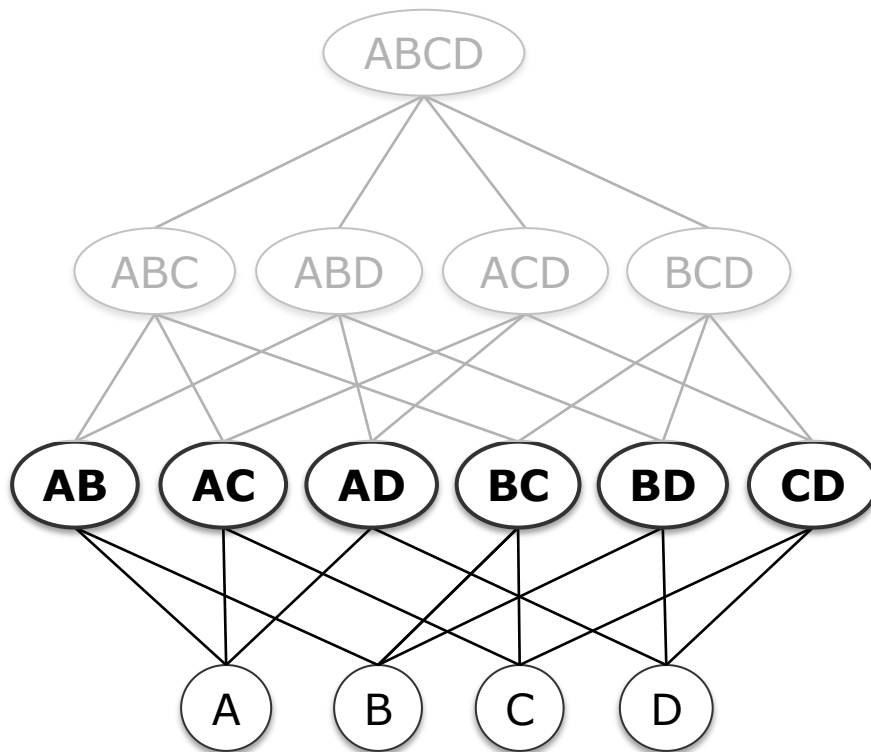
N-ary IND Discovery



Experimental Evaluation



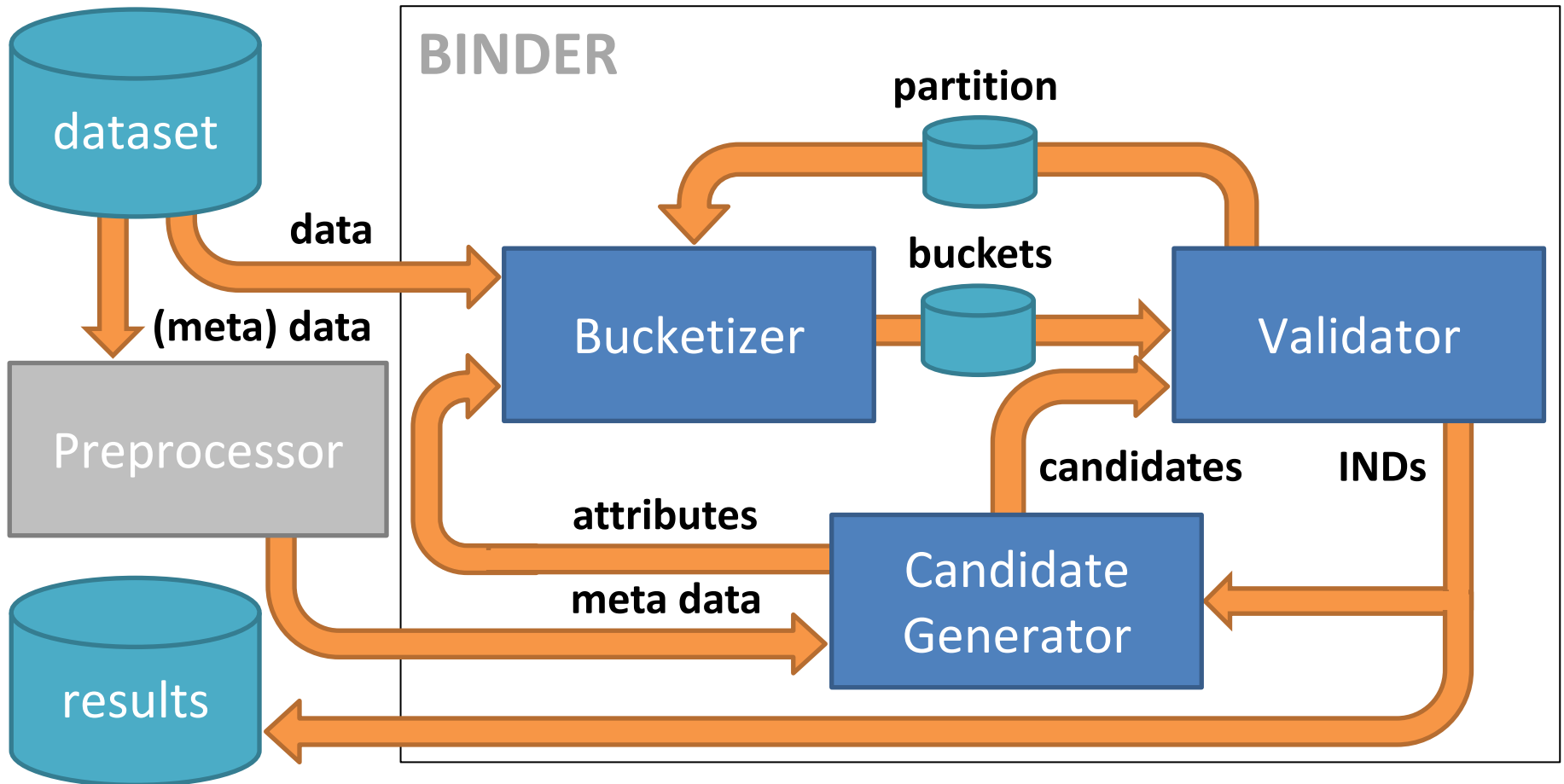
N-ary IND Discovery The BINDER Algorithm



$AB \subseteq CD$	$AB \subseteq DC$
$AC \subseteq BD$	$AC \subseteq DB$
$AD \subseteq BC$	$AD \subseteq CB$
$BC \subseteq AD$	$BC \subseteq DA$
$BD \subseteq AC$	$BD \subseteq CA$
$CD \subseteq AB$	$CD \subseteq BA$

$A \subseteq B$	$A \subseteq C$	$A \subseteq D$
$B \subseteq A$	$B \subseteq C$	$B \subseteq D$
$C \subseteq A$	$C \subseteq B$	$C \subseteq D$
$D \subseteq A$	$D \subseteq B$	$D \subseteq C$

N-ary IND Discovery The BINDER Algorithm



N-ary IND Discovery The BINDER Algorithm

New Candidate?



1. Both valid:

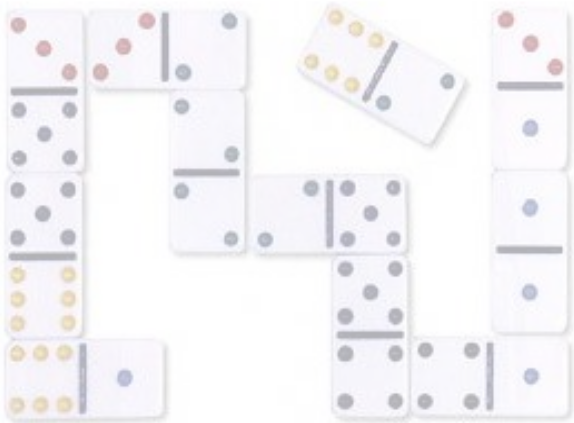
$$A \subseteq C$$

$$B \subseteq D$$

2. Non-trivial:

$$A \neq B \neq C \neq D$$

Unary IND Discovery



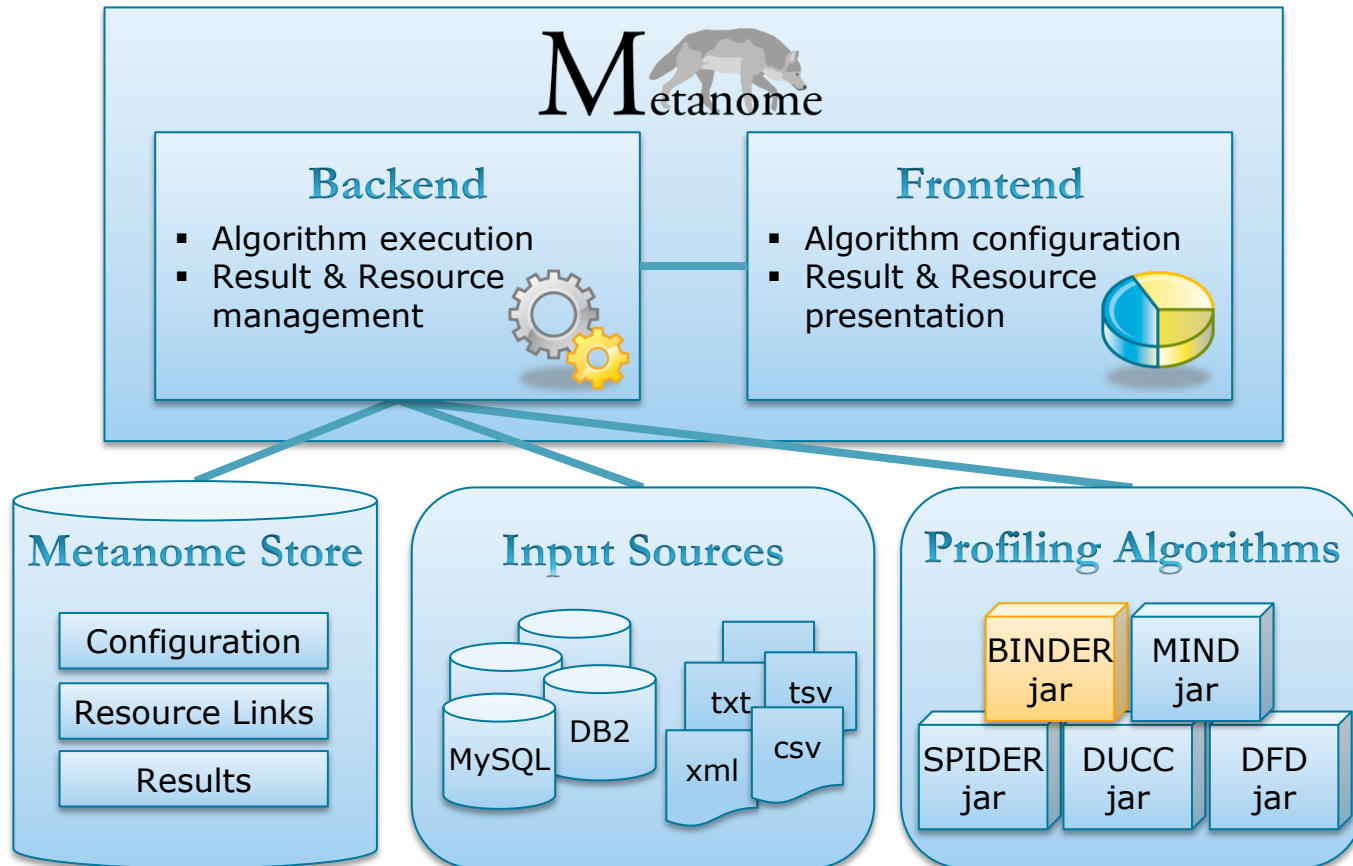
N-ary IND Discovery



Experimental Evaluation



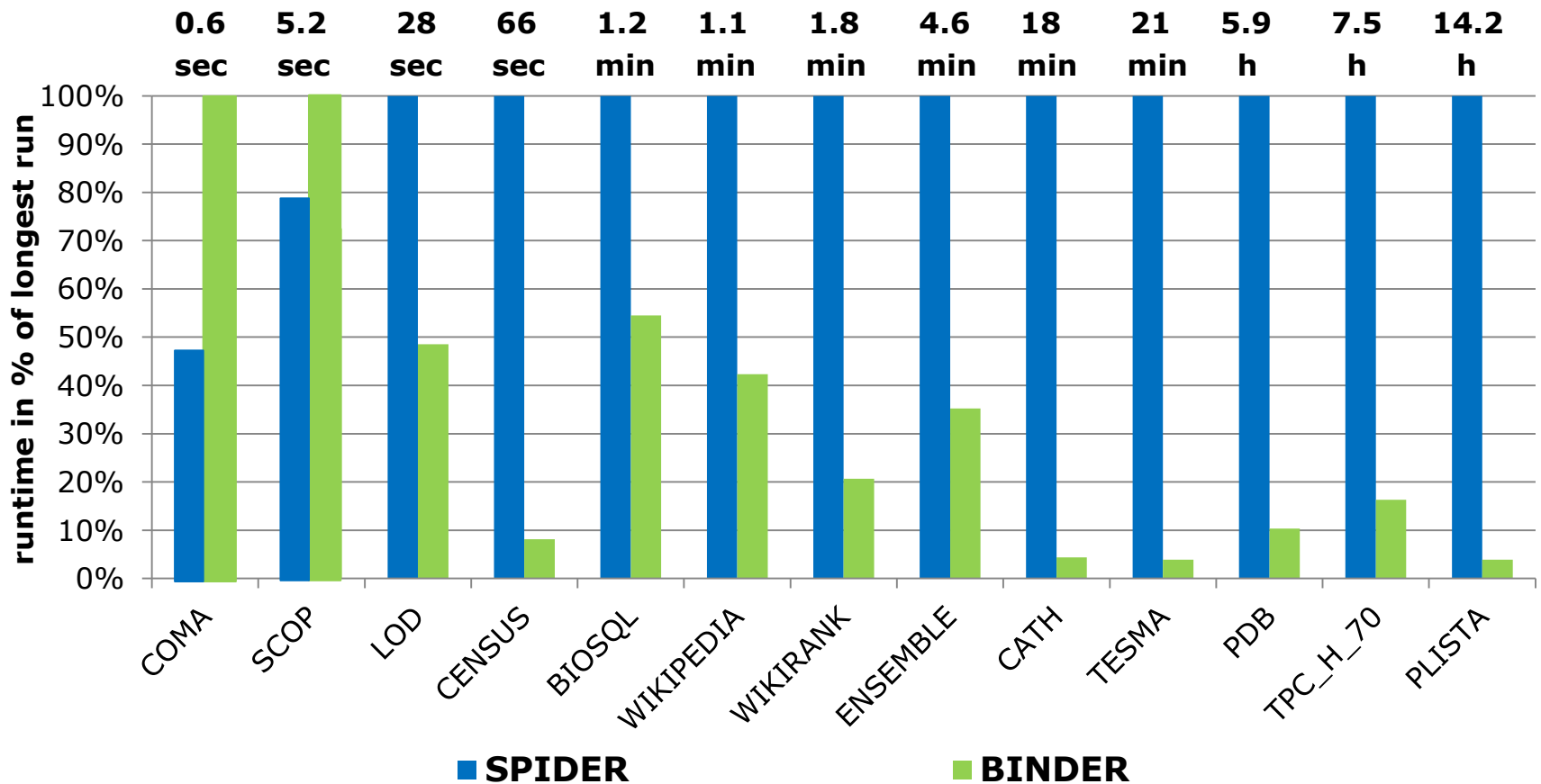
Experimental Evaluation The Metanome Profiling Tool



Experimental Evaluation Unary IND Discovery

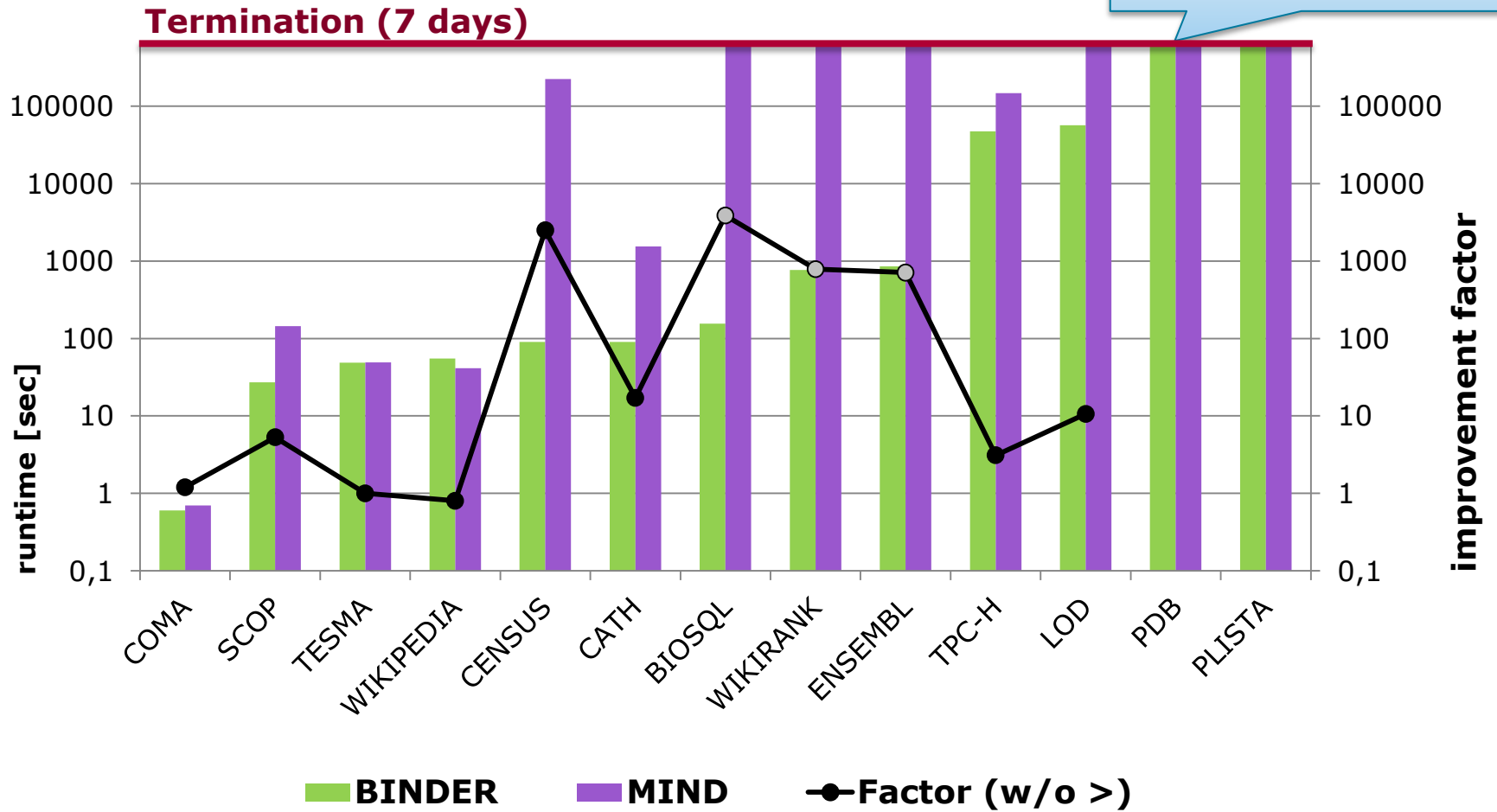
Longest run:

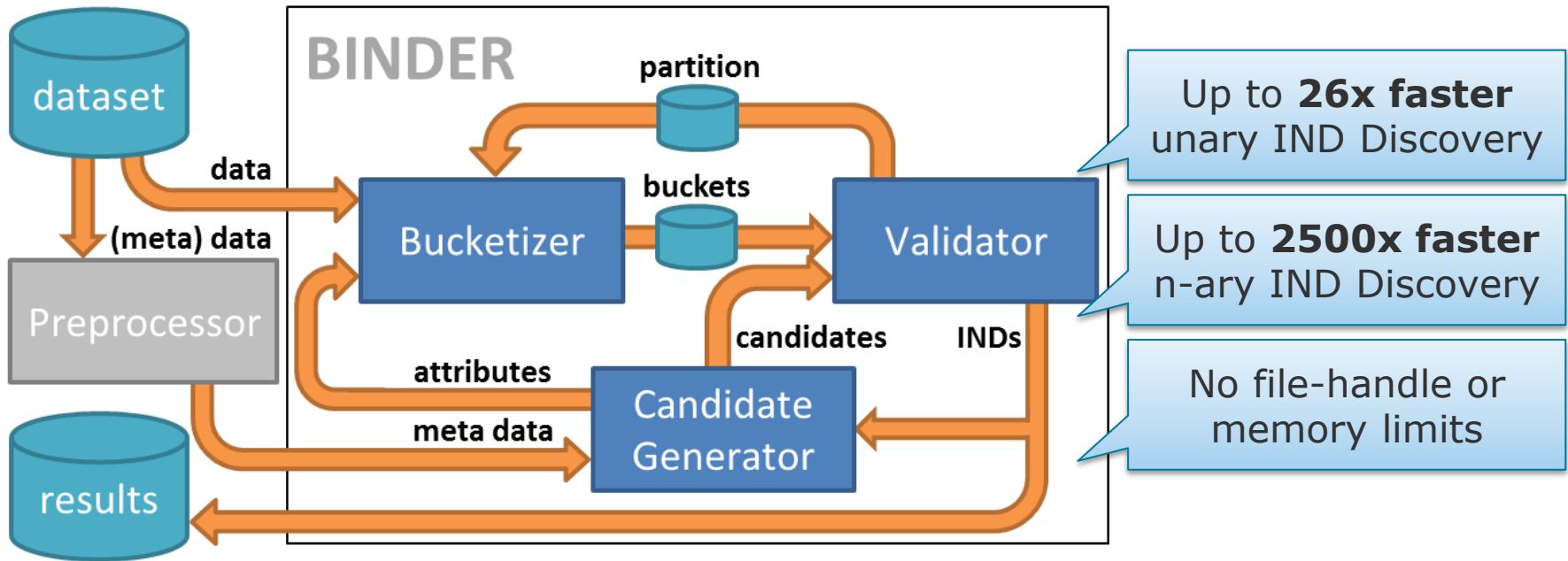
128 GB RAM



Experimental Evaluation N-ary IND Discovery

> 1 Billion INDs





Divide & Conquer-based Inclusion Dependency Discovery

Thorsten Papenbrock, Sebastian Kruse, Jorge-Arnulfo Quianè-Ruiz, Felix Naumann