

# Scaling out the Discovery of Inclusion Dependencies

BTW 2015, Hamburg, Germany

Sebastian Kruse, Thorsten Papenbrock, Felix Naumann  
Research Assistant  
Hasso Plattner Institute, Potsdam, Germany

# Inclusion Dependencies Examples

## Customers

ID	Name	Address
1	Tanja Jager	Marseiller Str. 12
2	Sandra Möller	Flughafenstr. 63
3	Dennis Eberhart	Sonnenallee 19
4	Barbara Pabst	Ziegelstr. 76
5	Thorsten Mauer	Güntzelstr. 90

## Orders

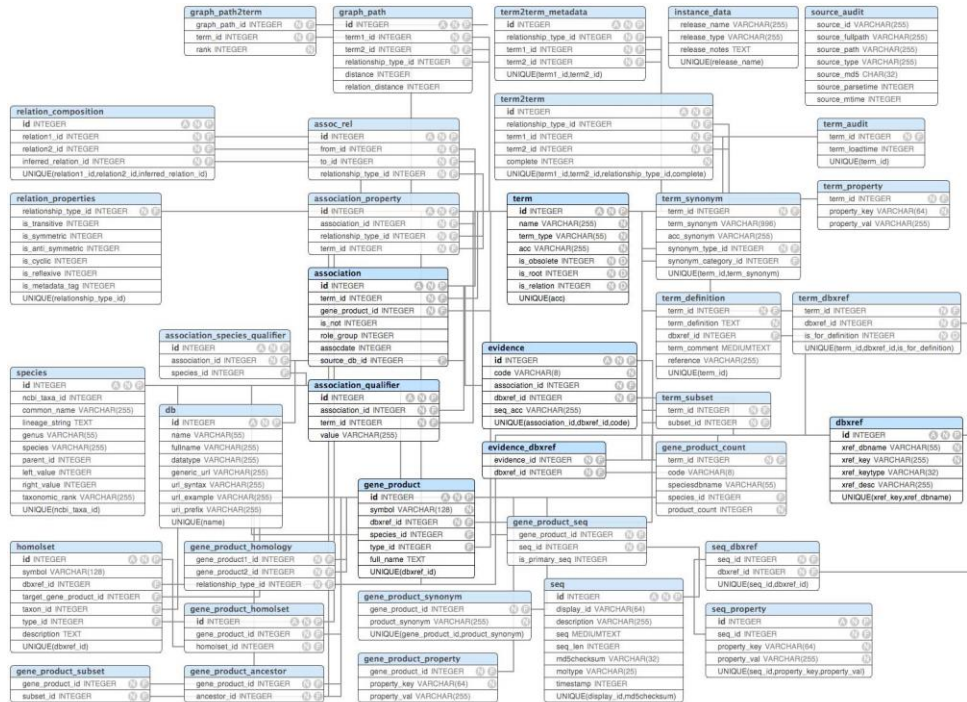
Customer	Item	Quantity
3	CK-242-1	1
3	DF-098-7	1
3	KE-883-6	1
1	LM-437-2	2
5	PE-383-5	1

$CustomerID \subseteq ID$

**Scaling out the  
Discovery of INDs**

Sebastian Kruse  
March 5, 2015

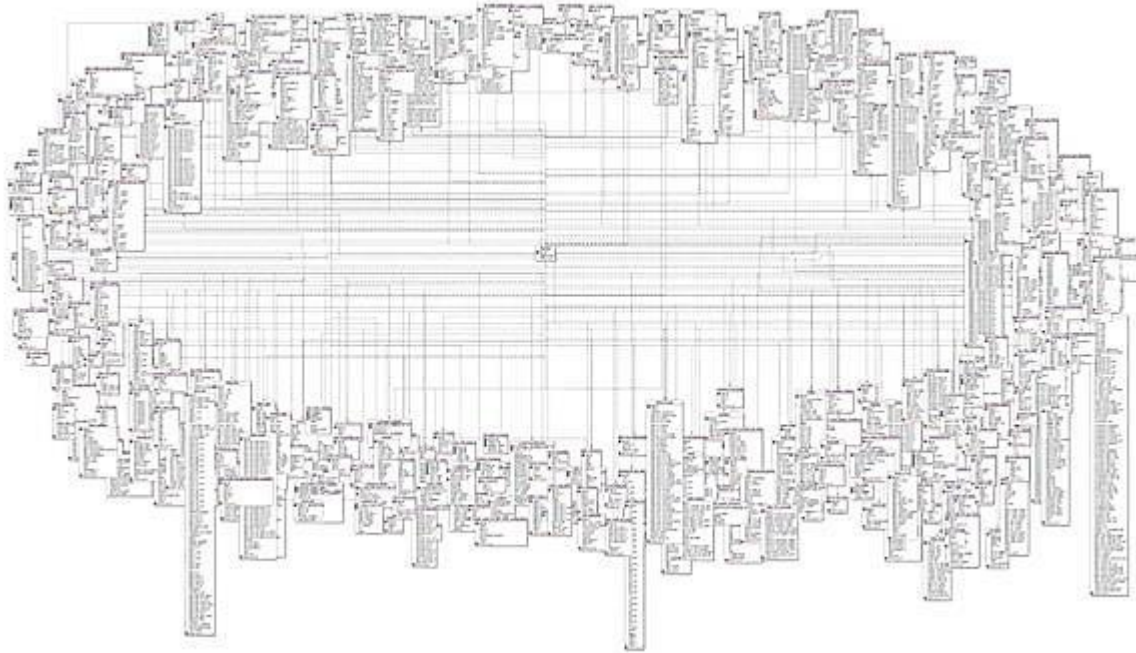
# Inclusion Dependencies Examples



Scaling out the Discovery of INDs

Sebastian Kruse  
March 5, 2015

# Inclusion Dependencies Example



**Scaling out the  
Discovery of INDs**

Sebastian Kruse  
March 5, 2015

# Scaling Out the Discovery of Inclusion Dependencies

## Agenda

---

- 1. Discovering Inclusion Dependencies**
2. Related Work
3. SINDY: A distributed discovery algorithm
4. Evaluation
5. Conclusions

# Scaling Out the Discovery of Inclusion Dependencies

## Agenda

---

1. Discovering Inclusion Dependencies
- 2. Related Work**
3. SINDY: A distributed discovery algorithm
4. Evaluation
5. Conclusions

# Related Work

## MIND

- Fabien De Marchi, Stéphan Lopes, and Jean-Marc Petit. Unary and n-ary inclusion dependency discovery in relational databases. *Journal of Intelligent Information Systems*, 32:53–73, 2009.

ID	Name	Address
1	Tanja Jager	Marseiller Str. 1
2	Sandra Möller	Flughafenstr. 63
3	Dennis Eberhart	Sonnenallee 19
4	Barbara Pabst	Ziegelstr. 76
5	Thorsten Mauer	Güntzelstr. 90

Customer	Item	Quantity
3	CK-242-1	1
3	DF-098-7	1
3	KE-883-6	1
1	LM-437-2	2
5	PE-383-5	1

### Scaling out the Discovery of INDs

Sebastian Kruse  
March 5, 2015

# Related Work

## MIND

Value	Attributes
1	ID, Customer, Quantity
Tanja Jager	Name
Marseiller Str. 12	Address
2	ID, Quantity
Sandra Möller	Name
Flughafenstr. 63	Address
...	...

Quantity  $\subseteq$  ID  
Quantity  $\subseteq$  Quantity

**Intersection**      **ID, Quantity**

**Scaling out the  
Discovery of INDs**

Sebastian Kruse  
March 5, 2015



## Related Work

# SPIDER

- Jana Bauckmann, Ulf Leser, and Felix Naumann. Efficiently Computing Inclusion Dependencies for Schema Discovery. In *ICDE Workshops*, 2006.

ID	Name	Address
1	Tanja Jager	Marseiller Str. 1
2	Sandra Möller	Flughafenstr. 63
3	Dennis Eberhart	Sonnenallee 19
4	Barbara Pabst	Ziegelstr. 76
5	Thorsten Mauer	Güntzelstr. 90

Customer	Item	Quantity
3	CK-242-1	1
3	DF-098-7	1
3	KE-883-6	1
1	LM-437-2	2
5	PE-383-5	1

### Scaling out the Discovery of INDs

Sebastian Kruse  
March 5, 2015

## Related Work

# SPIDER

- Jana Bauckmann, Ulf Leser, and Felix Naumann. Efficiently Computing Inclusion Dependencies for Schema Discovery. In *ICDE Workshops*, 2006.

ID	Name	Customer	Item	Quantity
1	Barbara Pabst	1	CK-242-1	1
2	Dennis Eberhart	3	DF-098-7	2
3	Sandra Möller	5	KE-883-6	
4	Tanja Jager		LM-437-2	
5	Thorsten Mauer		PE-383-5	

### Scaling out the Discovery of INDs

Sebastian Kruse  
March 5, 2015

# Common proceeding

**Input Data**

ID	Name	Addr
1	T.J.	M.12
2	S.M.	F.63
Cus	Item	Qty
3	CK	1
4	DF	1
5	KE	1
3	LM	2
5	PE	1

**Full Outer Join**

ID	Name	Addr	Cus	Item	Qty
1			1		1
2					2
3			3		
4					
5			5		
	T.J.				
	S.M.				
...	...	...	...	...	...

**Inclusion Dependencies**

Quantity  $\subseteq$  ID  
Customer  $\subseteq$  ID

**Step 1:**  
Calculate full outer join of all attributes

**Step 2:**  
Extract inclusion dependencies

**Scaling out the Discovery of INDs**

Sebastian Kruse  
March 5, 2015

# Scaling Out the Discovery of Inclusion Dependencies

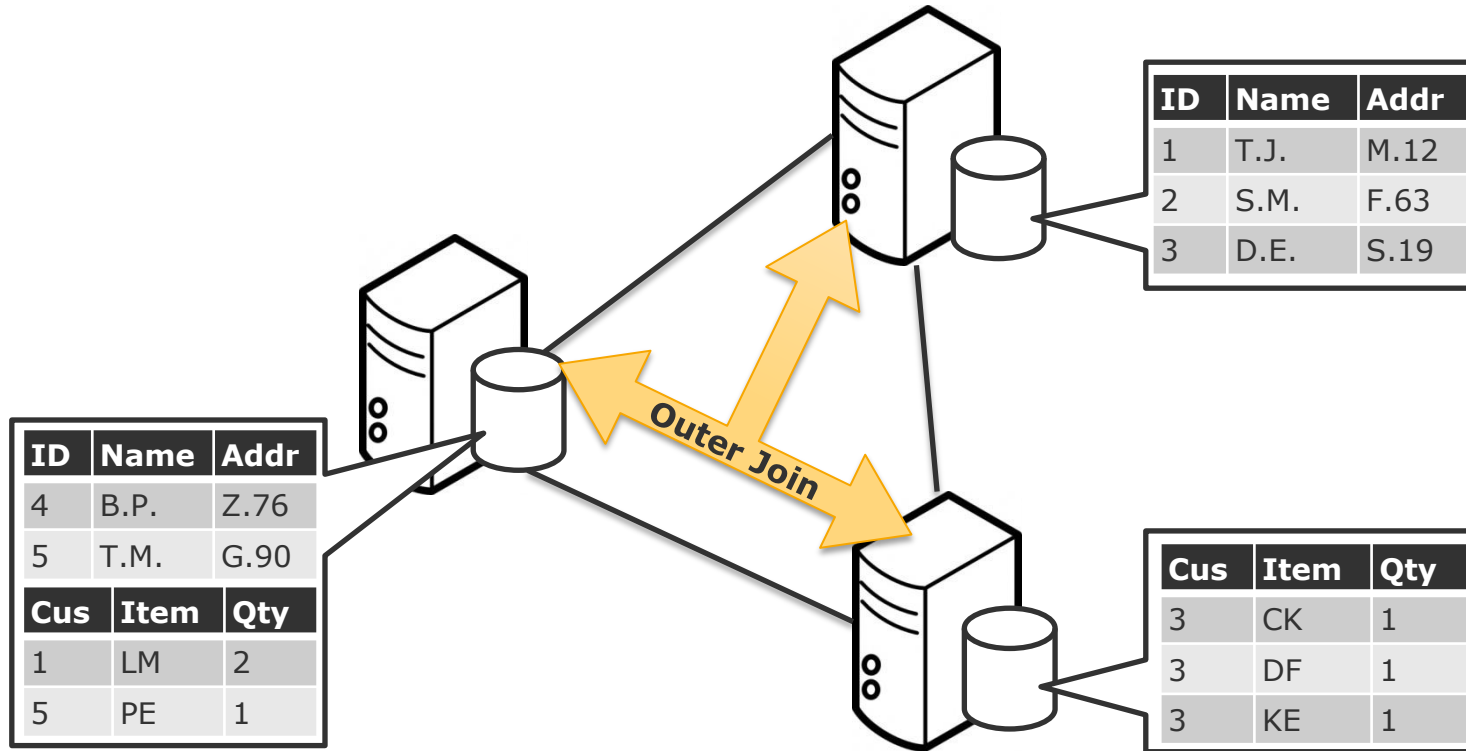
## Agenda

---

1. Discovering Inclusion Dependencies
2. Related Work
- 3. SINDY: A distributed discovery algorithm**
4. Evaluation
5. Conclusions

# SINDY: A distributed discovery algorithm

## Distributed setting

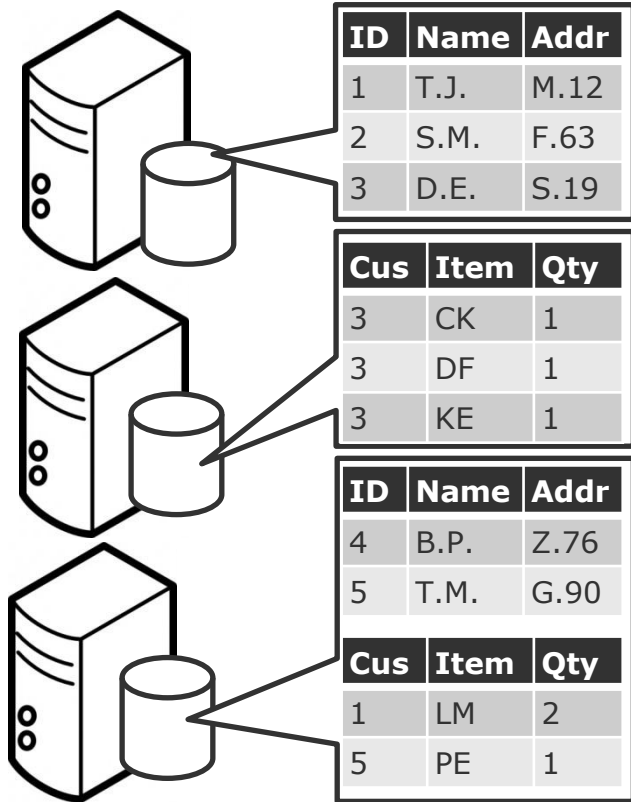


Scaling out the Discovery of INDs

Sebastian Kruse  
March 5, 2015

# SINDY: A distributed discovery algorithm

## Perform full outer join



1	ID	T.J.	Name	M.12	Addr
2	ID	S.M.	Name	F.63	Addr
3	ID	D.E.	Name	S.19	Addr

3	Cus	CK	Item	1	Qty
3	Cus	DF	Item	1	Qty
3	Cus	KE	Item	1	Qty

4	ID	B.P.	Name	Z.76	Addr
5	ID	T.M.	Name	G.90	Addr
1	Cus	LM	Item	2	Qty
5	Cus	PE	Item	1	Qty

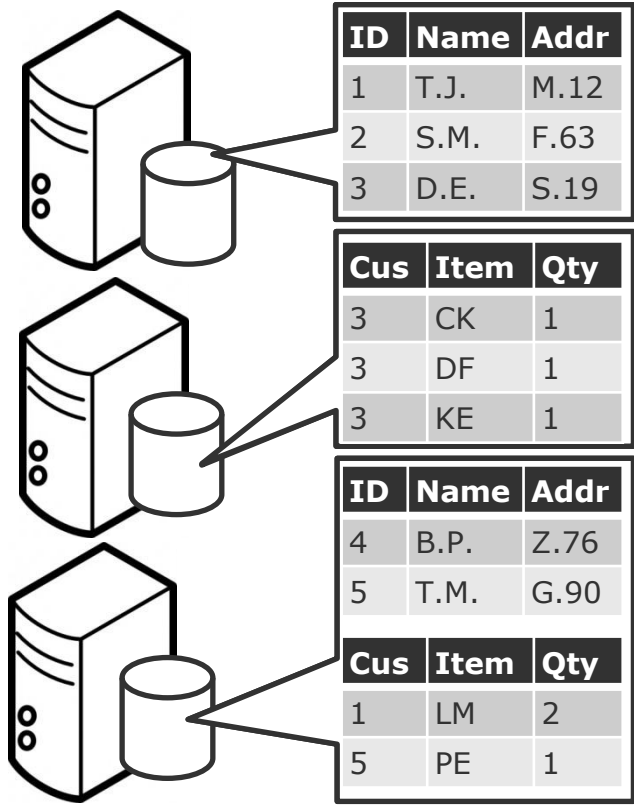
1	Cus, Qty
5	Cus, ID

**Scaling out the Discovery of INDs**

Sebastian Kruse  
March 5, 2015

# SINDY: A distributed discovery algorithm

## Perform full outer join

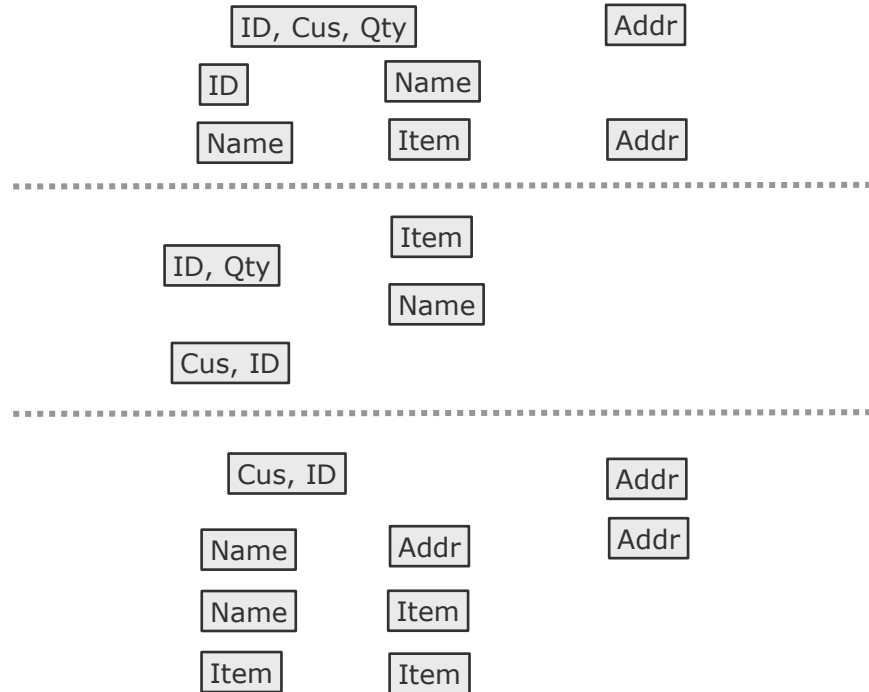
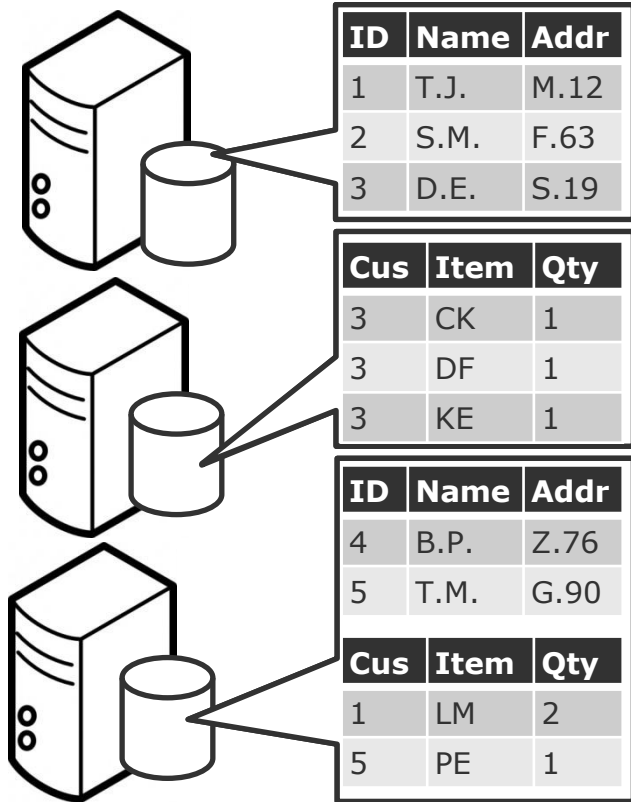


1	1	ID, Cus, Qty	me	M.12	Addr
2	ID	S.M.	Name	F.63	Addr
3	ID	D.E.	Name	S.19	Addr
-----					
5	Cus	CK	Item	1	Qty
2	ID, Qty	DF	Item		
		KE	Item		
-----					
4	3	ID, Cus	Name	Z.76	Addr
		T.M.	Name	G.90	Addr
		LM	Item	2	Qty
		PE	Item		

**Scaling out the Discovery of INDs**  
 Sebastian Kruse  
 March 5, 2015

# SINDY: A distributed discovery algorithm

## Perform full outer join



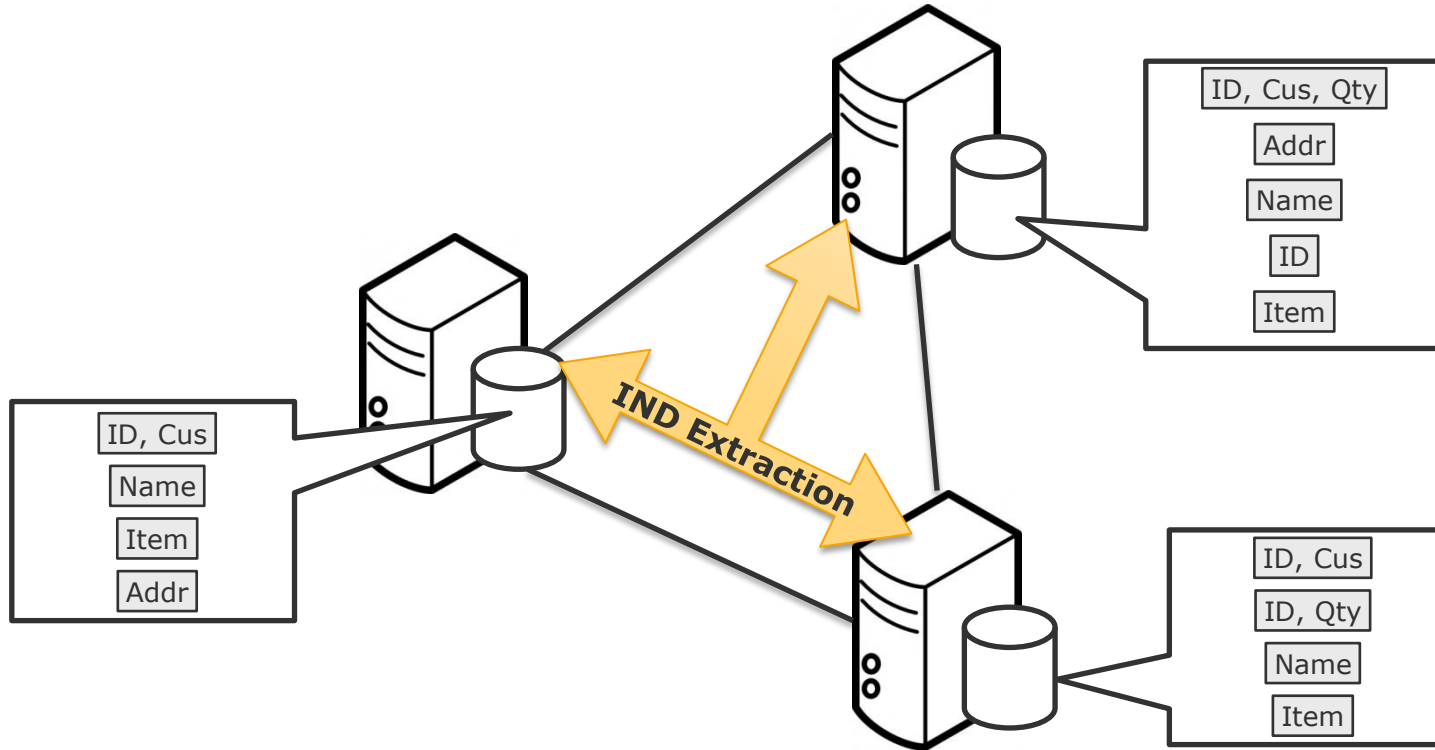
**Scaling out the Discovery of INDs**

Sebastian Kruse  
March 5, 2015



# SINDY: A distributed discovery algorithm

## Distributed join product

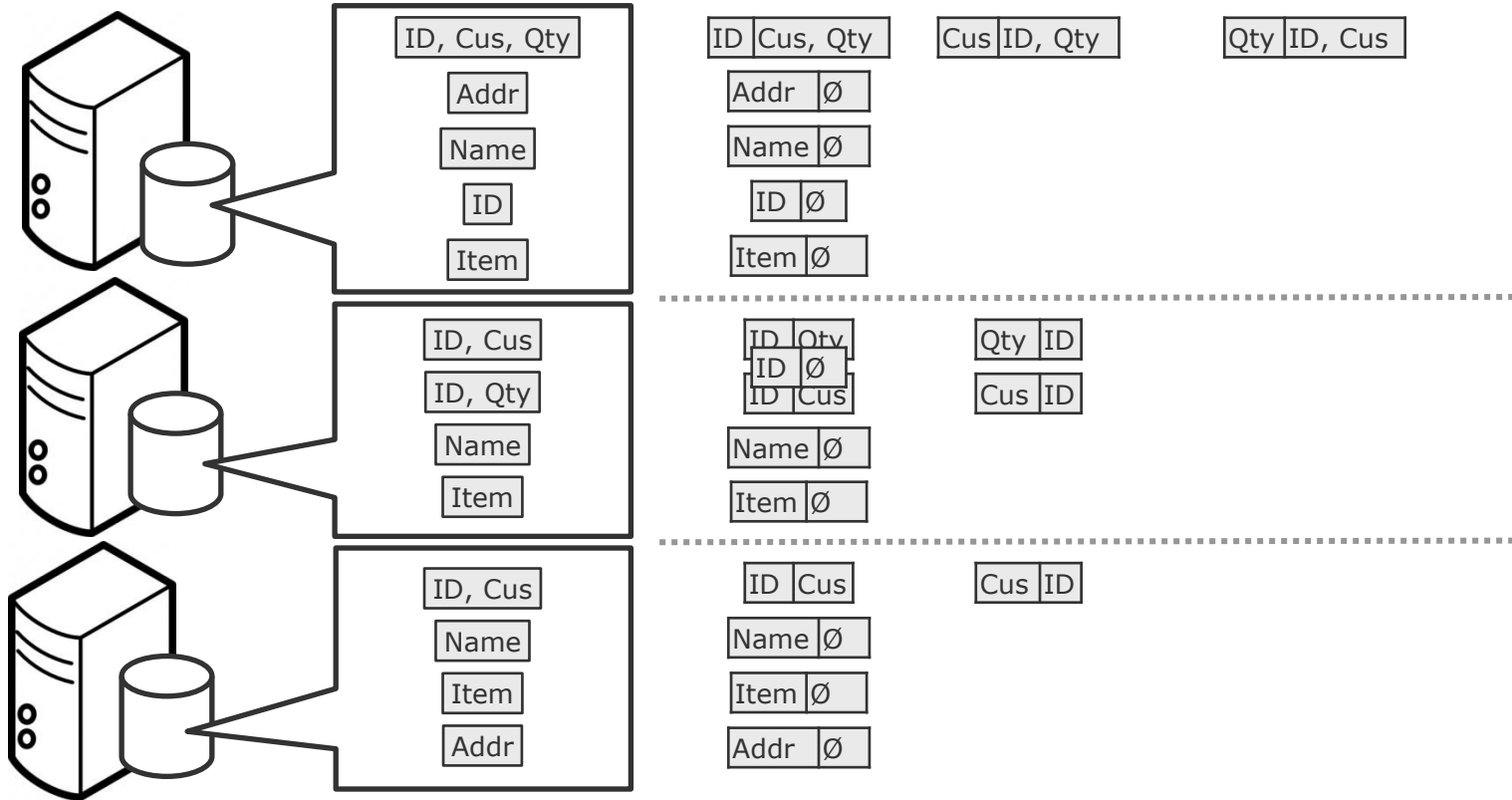


### Scaling out the Discovery of INDs

Sebastian Kruse  
March 5, 2015

# SINDY: A distributed discovery algorithm

## Evaluate full outer join

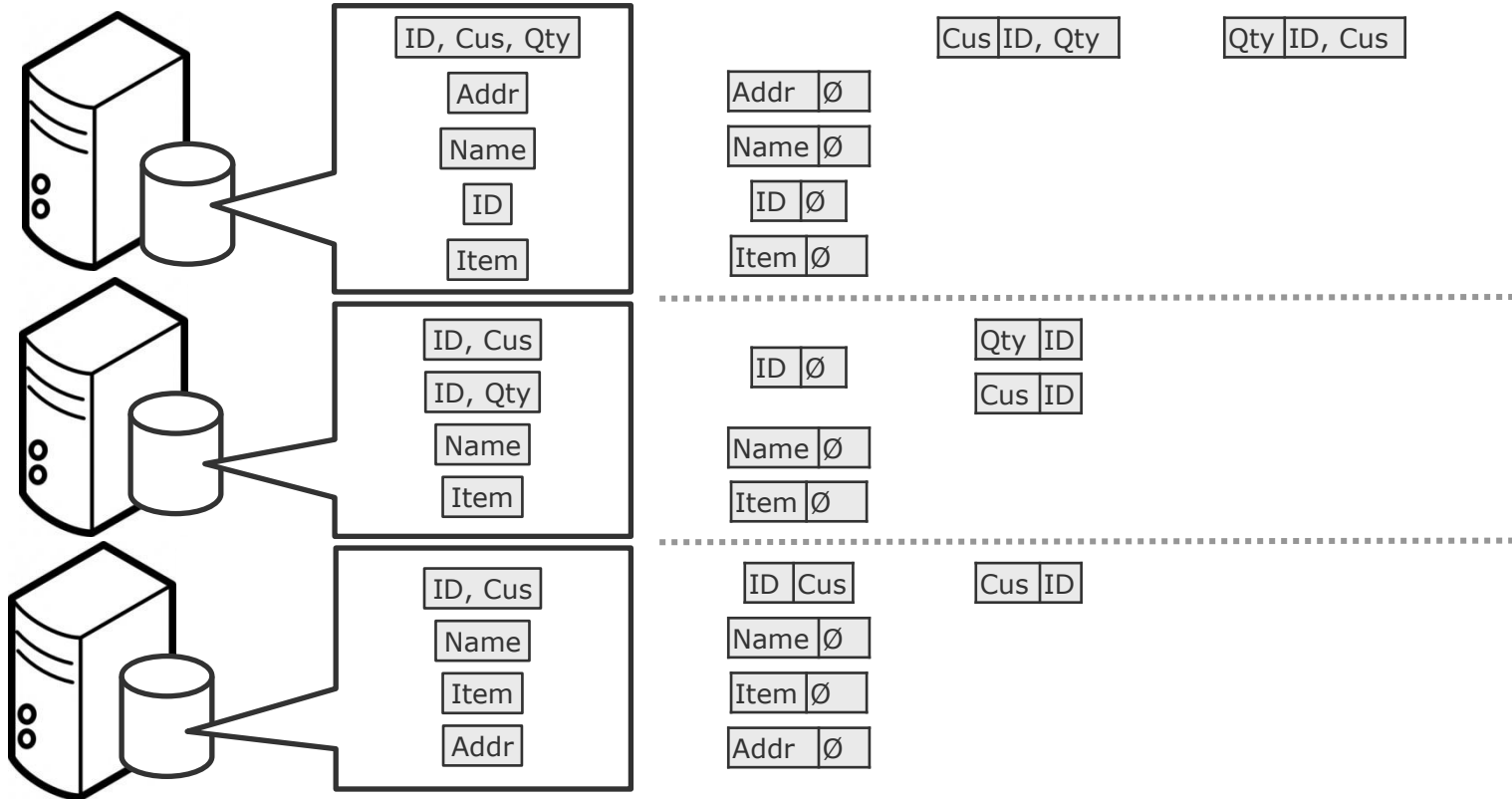


**Scaling out the Discovery of INDs**

Sebastian Kruse  
March 5, 2015

# SINDY: A distributed discovery algorithm

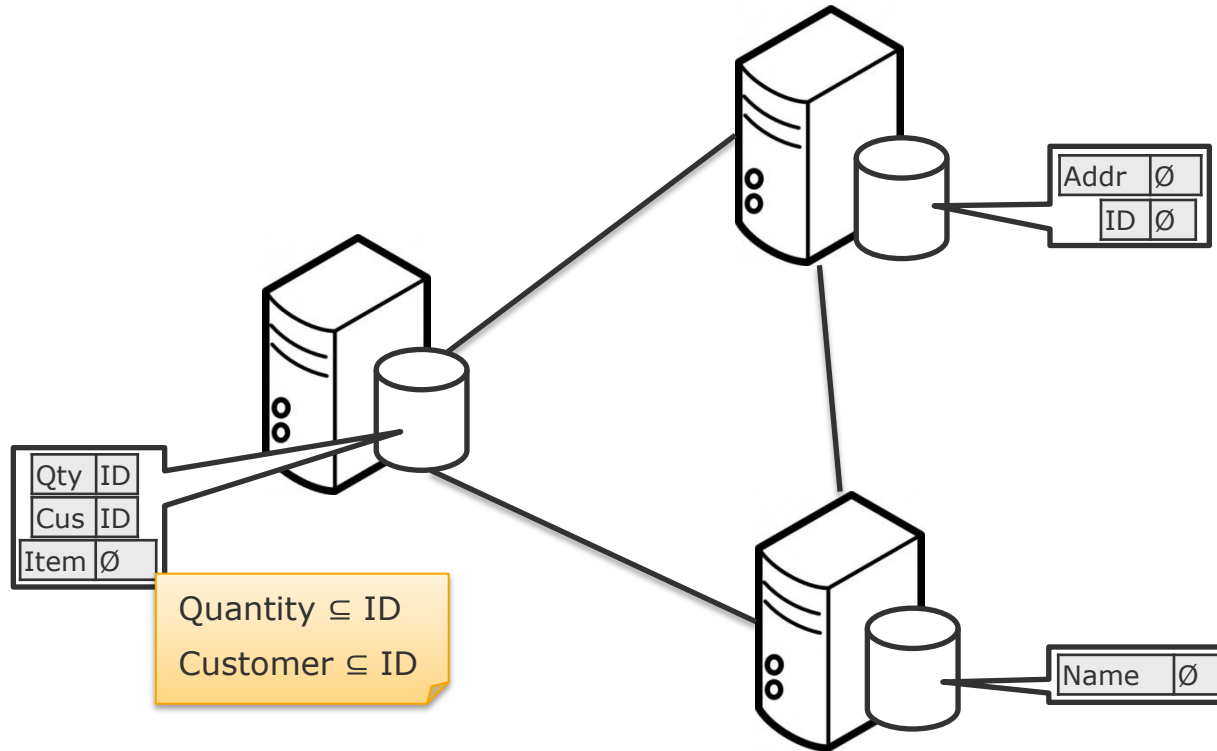
## Evaluate full outer join



**Scaling out the Discovery of INDs**

Sebastian Kruse  
March 5, 2015

# Distributed inclusion dependencies



**Scaling out the  
Discovery of INDs**

Sebastian Kruse  
March 5, 2015

- Inclusion dependencies on combinations of columns (aka n-ary INDs)
  - Adaption: Create cells for combinations of values
  - Powerful in combination with apriori-like proceeding
  
- Partial inclusion dependencies
  - Adaption: aggregate IND candidates with multiset union instead of intersection
  - Compare with number of distinct values of dependent column

**Scaling out the  
Discovery of INDs**

Sebastian Kruse  
March 5, 2015

# Scaling Out the Discovery of Inclusion Dependencies

## Agenda

---

1. Discovering Inclusion Dependencies
2. Related Work
3. SINDY: A distributed discovery algorithm
- 4. Evaluation**
5. Conclusions

### ■ Cluster Setup

- 1 master node (Intel Xeon @ 2x2.67 GHz, 8 GiB RAM)
- 10 worker nodes (Intel Core 2 Duo @ 2x2.6 GHz, 8 GiB RAM)
- Apache HDFS 2.2, Apache Flink 0.6.2

### ■ Single node (for SPIDER)

- Intel Xeon @ 8x2GHz, 128 GiB RAM, RAID-1

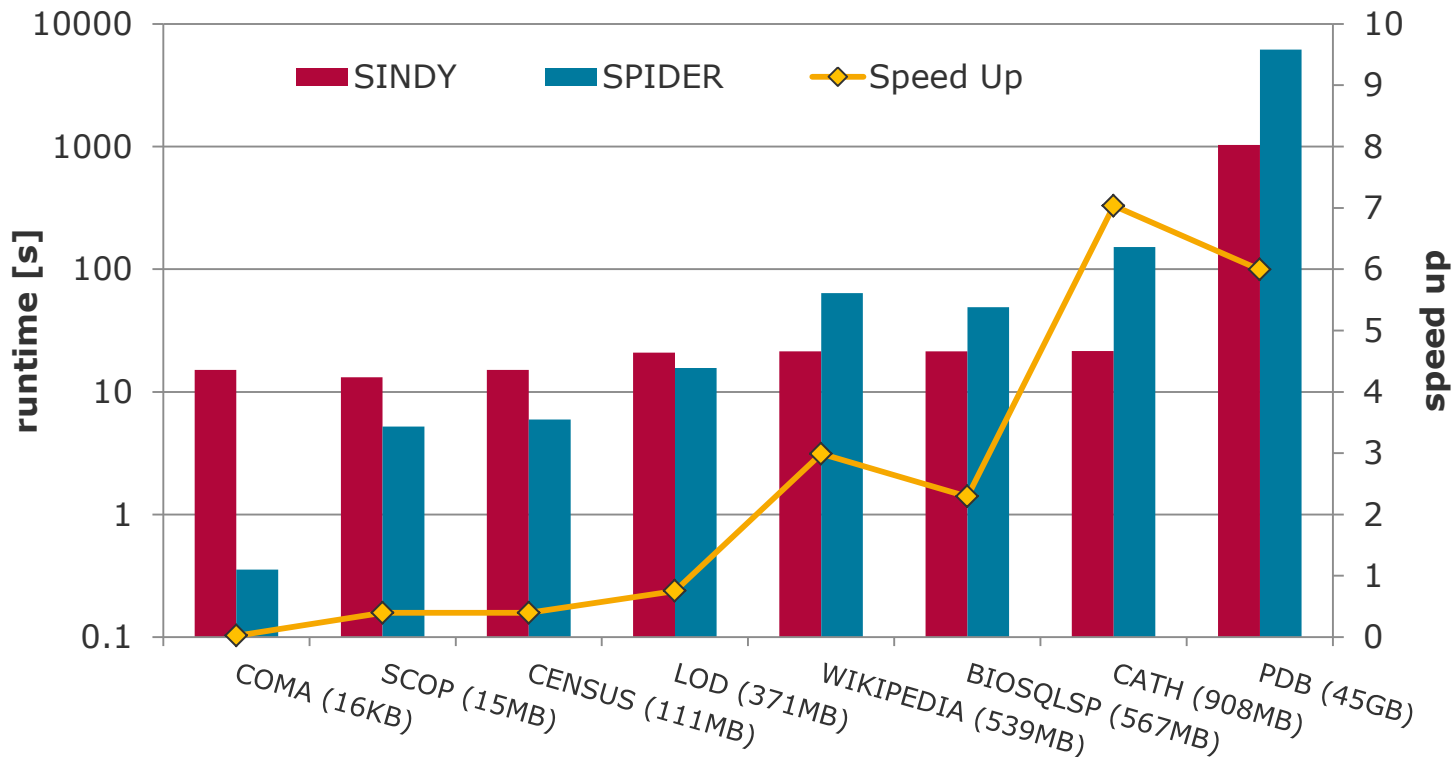
### ■ Datasets

- Relational datasets from different domains
- 16 KB to 44.9 GB

**Scaling out the  
Discovery of INDs**

Sebastian Kruse  
March 5, 2015

# Performance comparison with SPIDER



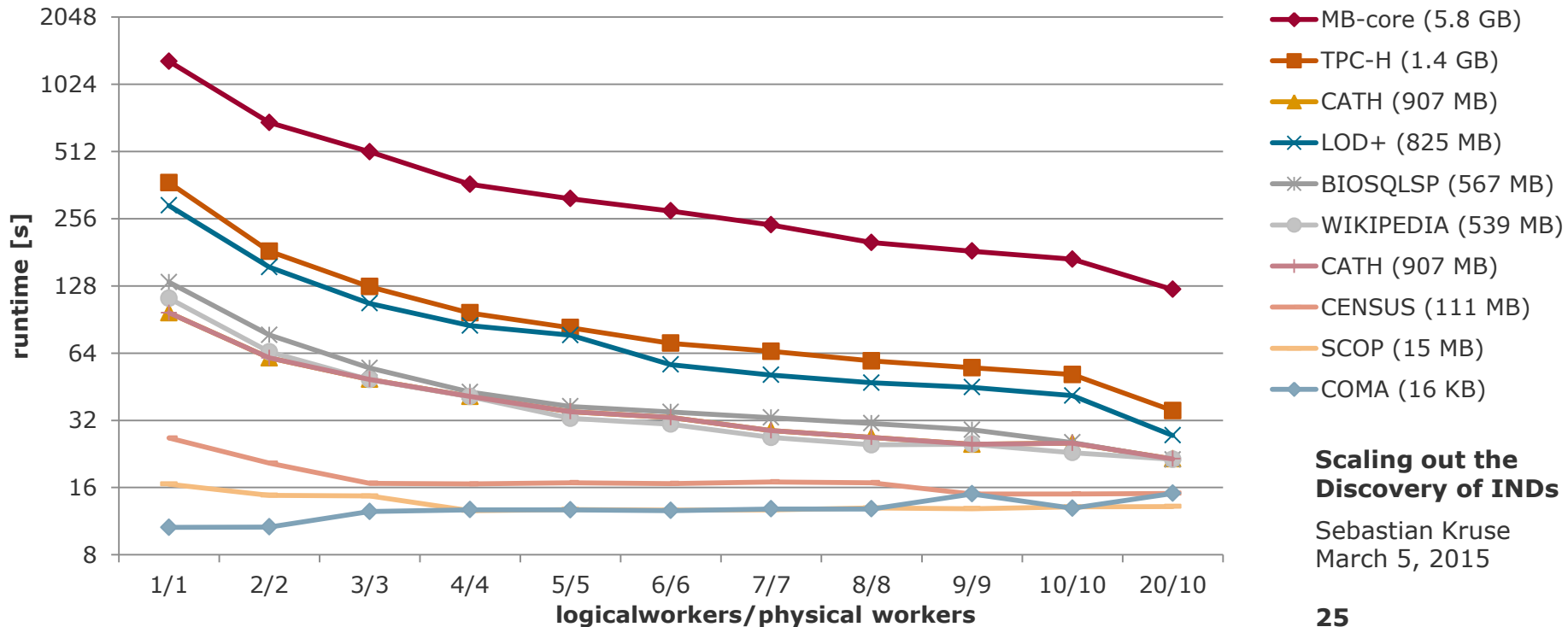
**Scaling out the Discovery of INDs**

Sebastian Kruse  
March 5, 2015



# Evaluation

## Scale-Out Behavior



**Scaling out the Discovery of INDs**

Sebastian Kruse  
March 5, 2015

# Scaling Out the Discovery of Inclusion Dependencies

## Agenda

---

1. Discovering Inclusion Dependencies
2. Related Work
3. SINDY: A distributed discovery algorithm
4. Evaluation
- 5. Conclusions**

# Scaling Out the Discovery of Inclusion Dependencies

## Conclusions

---

- Presented new distributed IND discovery algorithm
  - Applicable for unary, n-ary, and partial inclusion dependencies
  - Consists of full outer join calculation and extraction phase
  - Scales well on large datasets
  
- Open questions
  - How can one continuously maintain inclusion dependencies?
  - How can the algorithm be applied to similar problems, e.g., RDF data?

**Scaling out the  
Discovery of INDs**

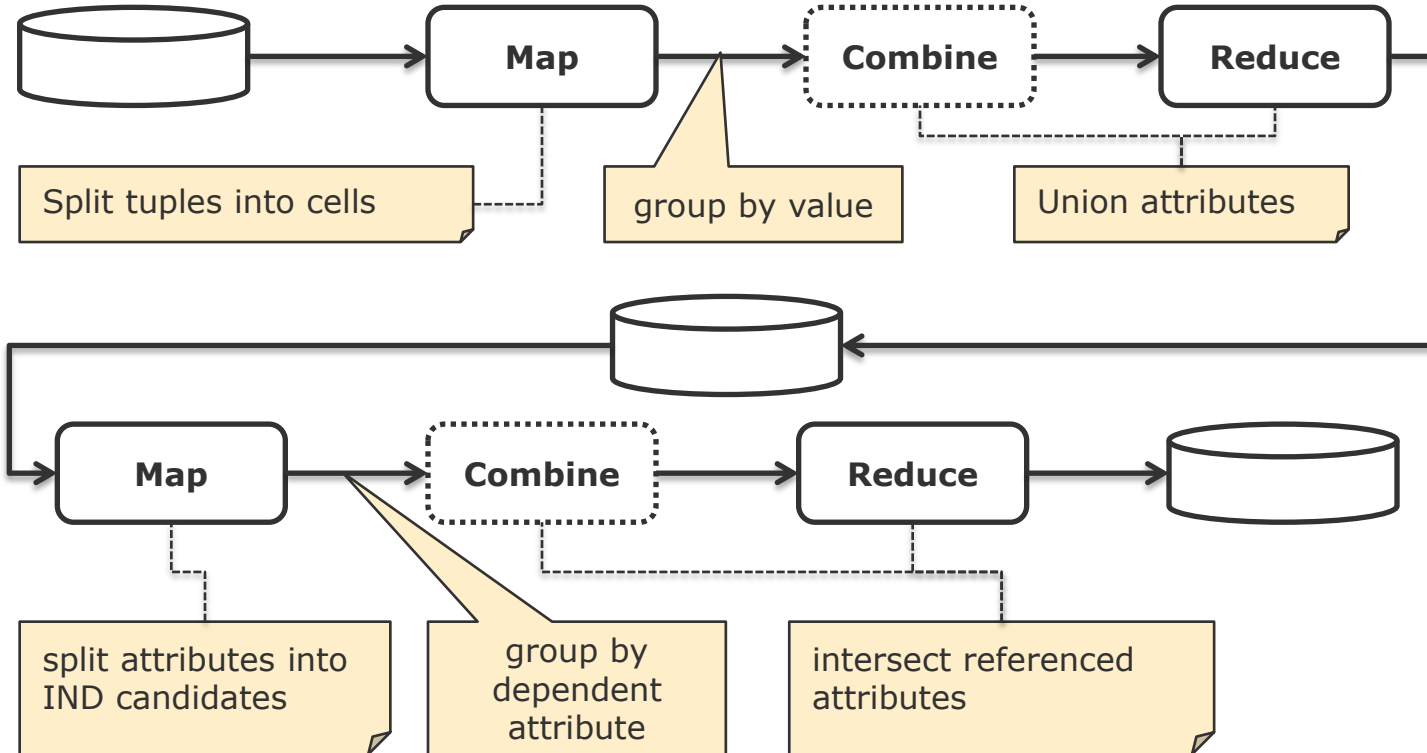
Sebastian Kruse  
March 5, 2015

# Scaling out the Discovery of Inclusion Dependencies

## BTW 2015, Hamburg, Germany

Sebastian Kruse, Thorsten Papenbrock, Felix Naumann  
Research Assistant  
Hasso Plattner Institute, Potsdam, Germany

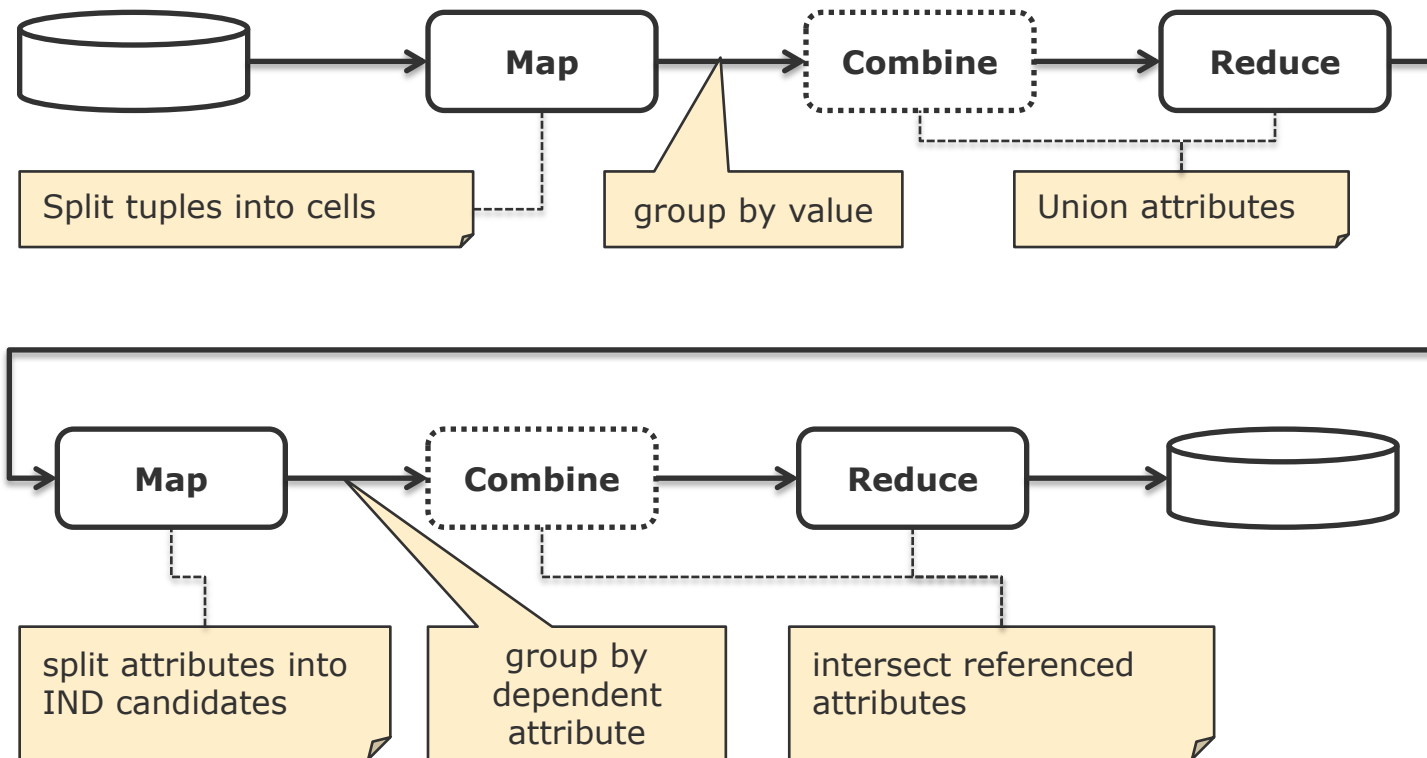
# Apache Hadoop Implementation



**Scaling out the Discovery of INDs**

Sebastian Kruse  
March 5, 2015

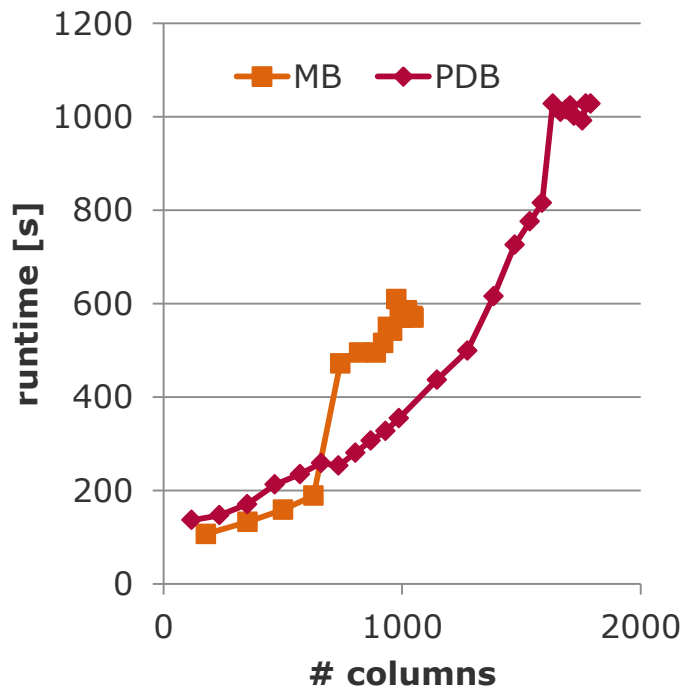
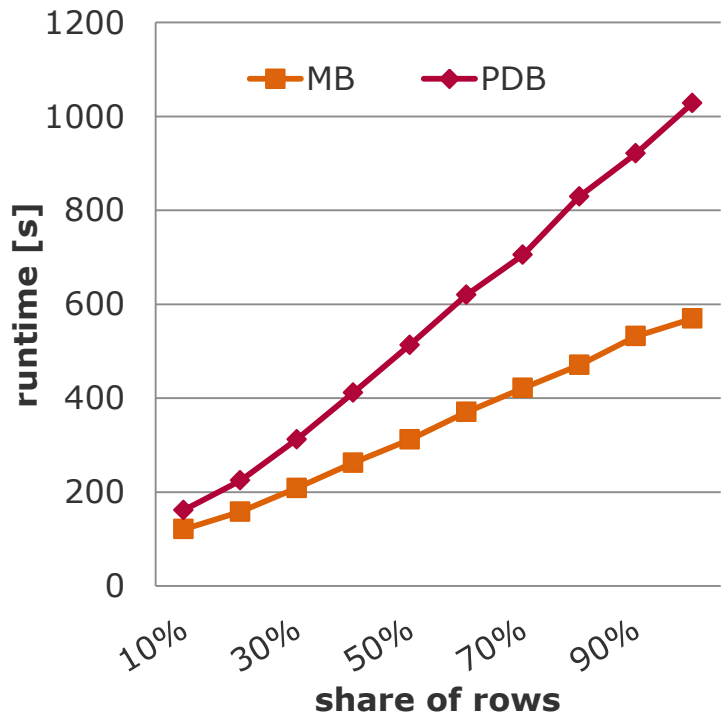
# Apache Flink/Spark Implementation



**Scaling out the Discovery of INDs**

Sebastian Kruse  
March 5, 2015

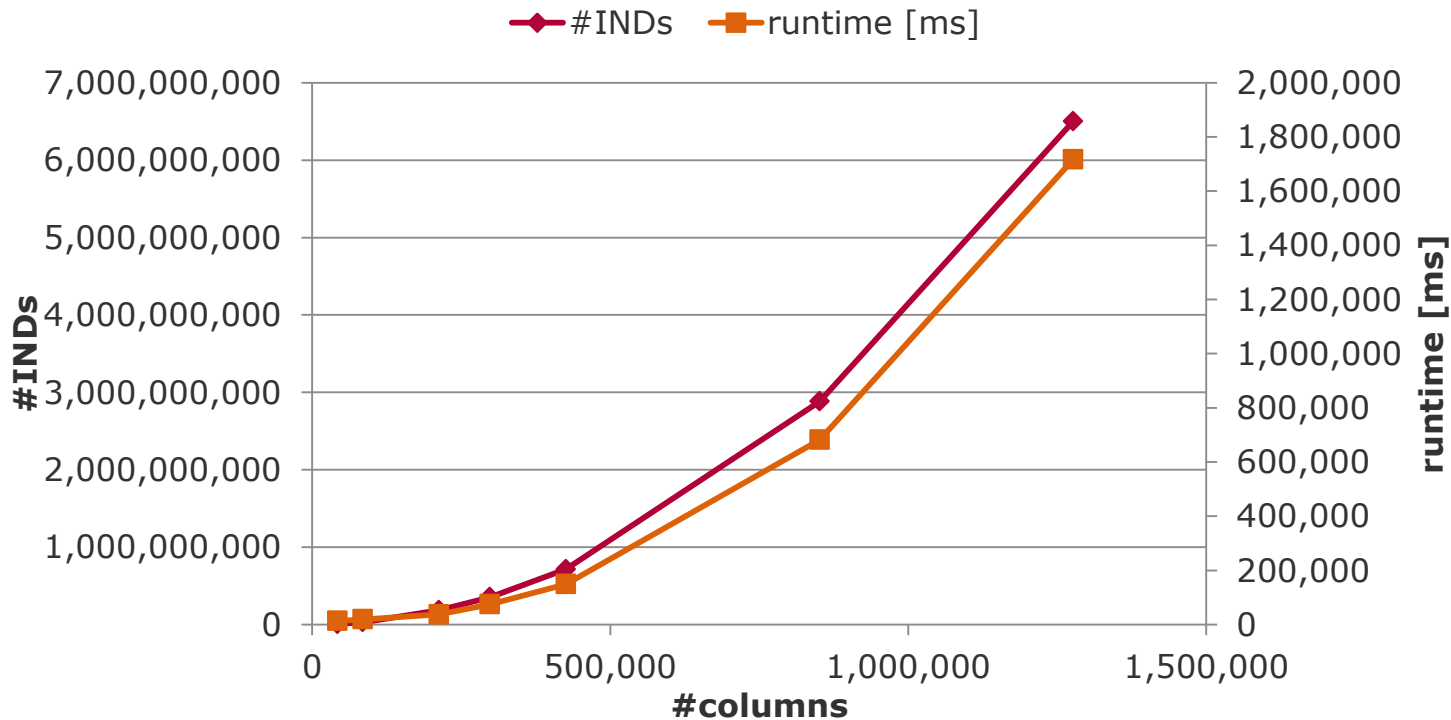
# Column and row scaling behavior



**Scaling out the  
Discovery of INDs**

Sebastian Kruse  
March 5, 2015

# Evaluation: Wikitables

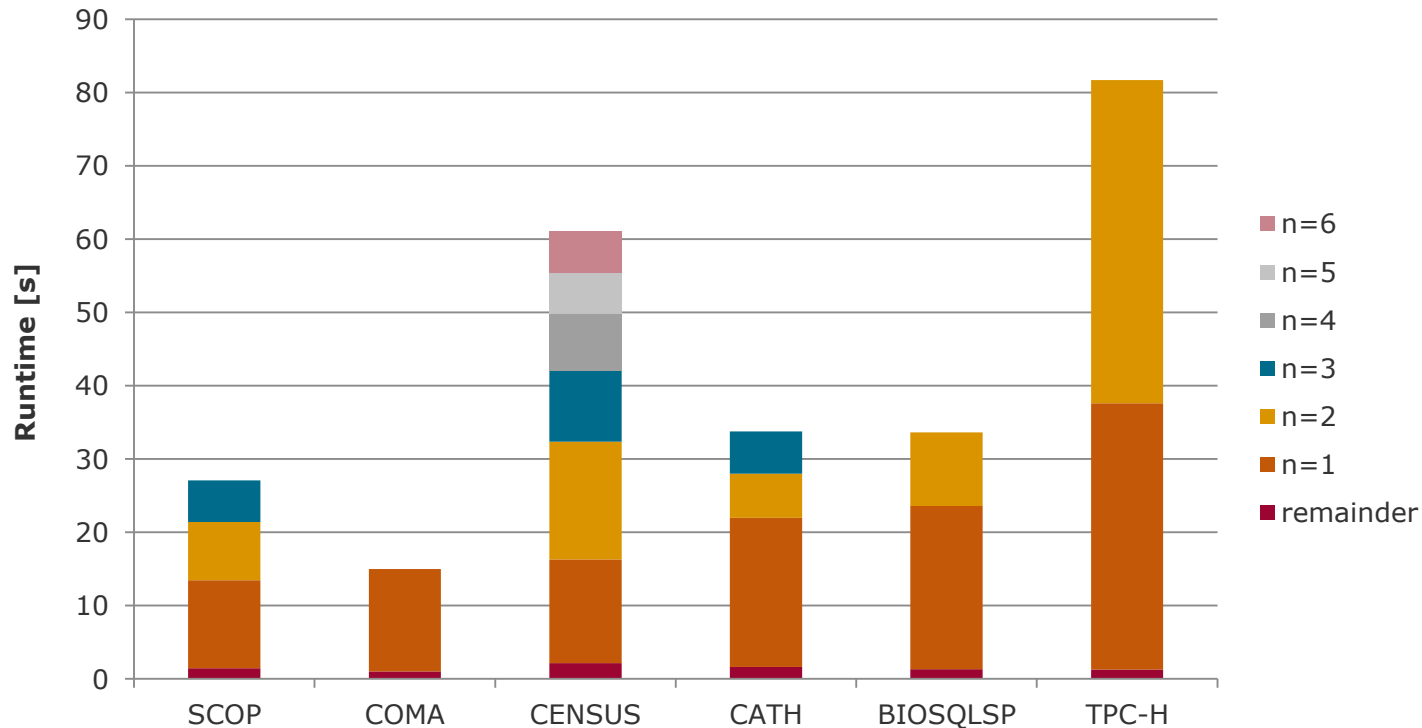


**Scaling out the Discovery of INDs**

Sebastian Kruse  
March 5, 2015



# Evaluation: n-ary INDs



**Scaling out the Discovery of INDs**

Sebastian Kruse  
March 5, 2015