



Using Keywords to Find the Right Path through Relational Data

Roberto De Virgilio, [Antonio Maccioni](#) and Riccardo Torlone

SEBD 2012


Venice, 26-06-2012

Motivation Scenario



Number	LastName	FirstName	Position	Goals

Tabella di database



LastName	FirstName	Position

Tabella di recordset



The Problem

- **Keyword Search on Relational Databases**
 - **INPUT:** keyword query. Ex. **Q = {Database, Matters}**
 - **OUTPUT:** tuples matching keywords connected by foreign keys.

T_4 : **Author**

	<u>name</u>	<u>publication</u>
t_{13}	Lenzerini	Data Integration
t_{14}	Date	Foundation Matters

T_5 : **Publication**

	<u>title</u>	<u>year</u>
t_{15}	Data Integration	2002
t_{16}	Foundation Matters	2002

T_1 : **Person**

	<u>name</u>	<u>area</u>
t_1	Watson	Database
t_2	Lenzerini	Database
t_3	Date	Database
t_4	Hunt	Inf. Systems

T_2 : **Affiliated**

	<u>professor</u>	<u>department</u>
t_5	Watson	x123
t_6	Lenzerini	cs34
t_7	Date	cs34
t_8	Hunt	m111

T_3 : **Department**

	<u>id</u>	<u>dname</u>	<u>director</u>
t_9	x123	CS	Watson
t_{10}	cs34	IE	Hunt
t_{11}	ee67	EE	Date
t_{12}	m111	ME	Hunt

The Problem

- **Keyword Search on Relational Databases**
 - **INPUT:** keyword query. Ex. **Q = {Database, Matters}**
 - **OUTPUT:** tuples matching keywords connected by foreign keys.

T_4 : **Author**

	<u>name</u>	<u>publication</u>
t_{13}	Lenzerini	Data Integration
t_{14}	Date	Foundation Matters

T_5 : **Publication**

	<u>title</u>	<u>year</u>
t_{15}	Data Integration	2002
t_{16}	Foundation Matters	2002

T_1 : **Person**

	<u>name</u>	<u>area</u>
t_1	Watson	Database
t_2	Lenzerini	Database
t_3	Date	Database
t_4	Hunt	Inf. Systems

T_2 : **Affiliated**

	<u>professor</u>	<u>department</u>
t_5	Watson	x123
t_6	Lenzerini	cs34
t_7	Date	cs34
t_8	Hunt	m111

T_3 : **Department**

	<u>id</u>	<u>dname</u>	<u>director</u>
t_9	x123	CS	Watson
t_{10}	cs34	IE	Hunt
t_{11}	ee67	EE	Date
t_{12}	m111	ME	Hunt

The Problem

- **Keyword Search on Relational Databases**
 - **INPUT:** keyword query. Ex. **Q = {Database, Matters}**
 - **OUTPUT:** tuples matching keywords connected by foreign keys.

T_4 : Author

	<u>name</u>	<u>publication</u>
t_{13}	Lenzerini	Data Integration
t_1	Date	Foundation Matters

T_5 : Publication

	<u>title</u>	<u>year</u>
t_{15}	Data Integration	2002
t_1	Foundation Matters	2002

T_1 : Person

	<u>name</u>	<u>area</u>
t_1	Watson	Database
t_2	Lenzerini	Database
t_3	Date	Database
t_4	Hunt	Inf. Systems

T_2 : Affiliated

	<u>professor</u>	<u>department</u>
t_5	Watson	x123
t_6	Lenzerini	cs34
t_7	Date	cs34
t_8	Hunt	m111

T_3 : Department

	<u>id</u>	<u>dname</u>	<u>director</u>
t_9	x123	CS	Watson
t_{10}	cs34	IE	Hunt
t_{11}	ee67	EE	Date
t_{12}	m111	ME	Hunt

State Of The Art

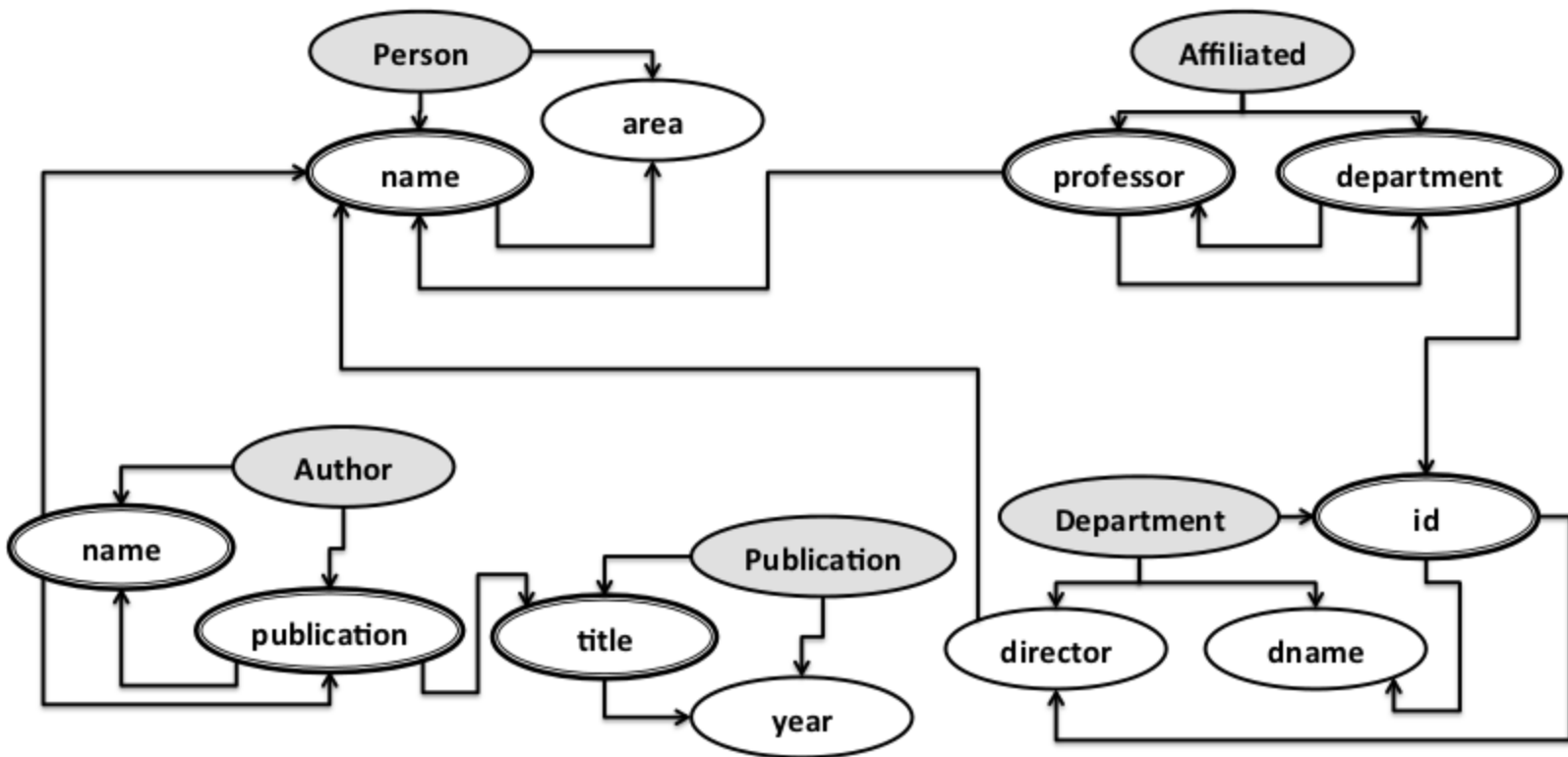
- **Schema-free**: graph based algorithms
[Kacholia et al., VLDB 2005] [He et al., SIGMOD 2007]
 - In-memory elaboration
 - Lot of I/O processing*RDBMS features are not exploited [Qin et al., SIGMOD 2009]*
- **Schema-based**: query rewriting
[Luo et al., SIGMOD 2007] [Bergamaschi et al., SIGMOD 2011]
 - Produce and Execute lot of SQL queries
 - Empty results*RDBMS features are not WELL exploited [Qin et al., SIGMOD 2009]*

The Idea

- Avoid *middleware approaches* and *graph navigations*.
- Build the solutions just **making use** of the **RDBMS**.
- Lean *exploration* guided by the use of **paths** resulting from the logical schema.

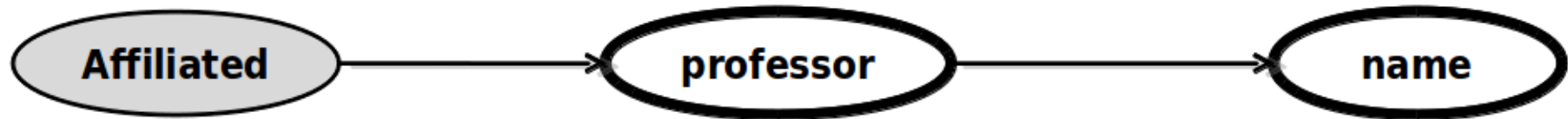
The Schema Graph

- Logical schema of the DB as a graph $SG = (V_s, E_s)$, where $V_s = (T, A)$ and $E_s \subseteq (T \times A) \cup (A \times A)$

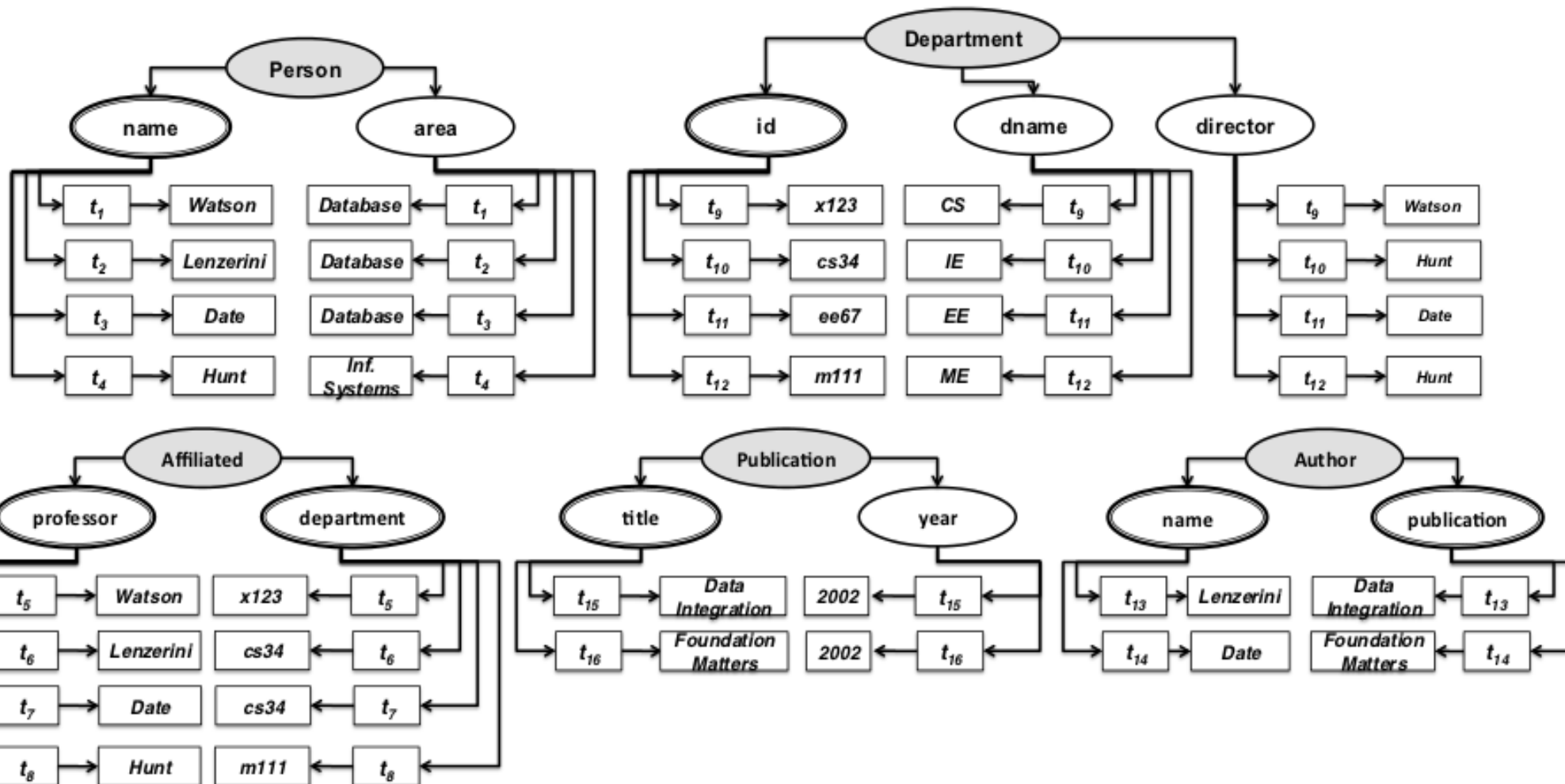


The Schema Paths

- A sequence $v_1 \rightarrow v_2 \rightarrow \dots \rightarrow v_f$ where
 - $v_1 \in T$
 - $(v_i, v_{i+1}) \in E_s$

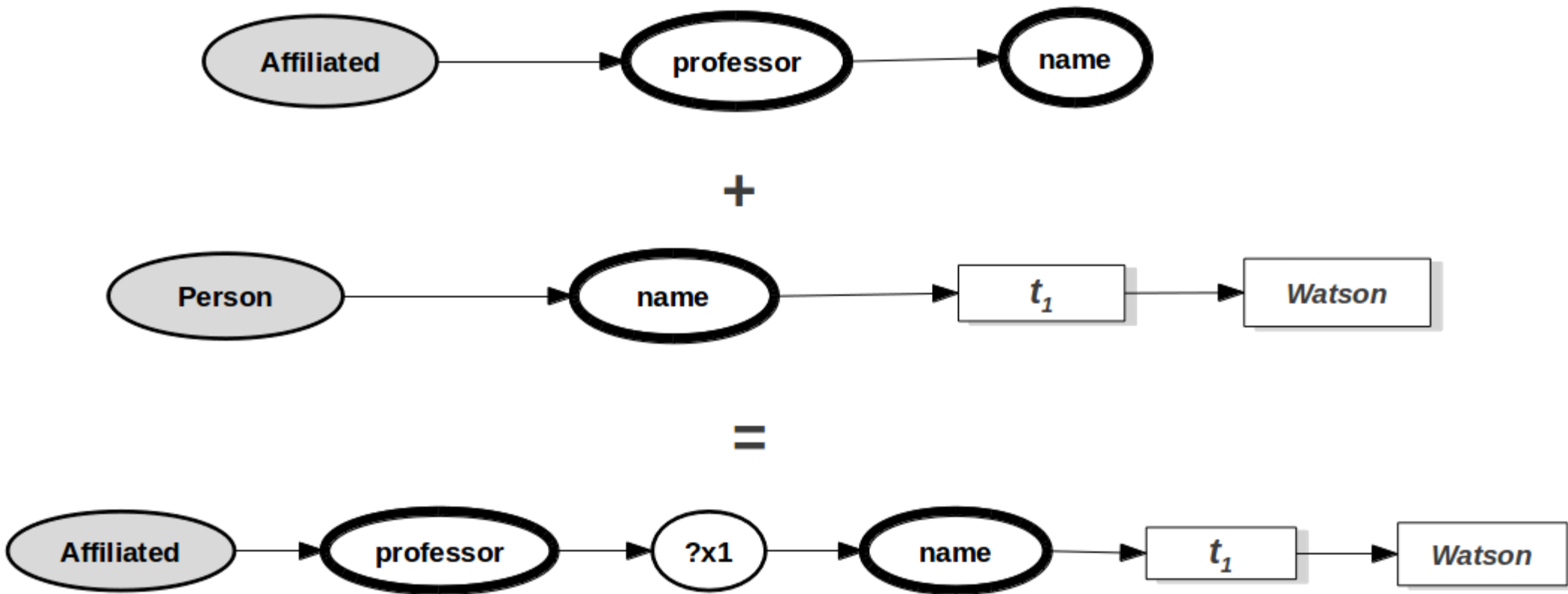


Data Graph And Data Paths



Expanded Data Paths

- Schema Path + Data Path = Expanded Data Path



The Approach

QUERYING

+

CLUSTERING

+

BUILDING

1. Retrieve the data paths
2. Expand the data paths

The Approach

QUERYING

+

CLUSTERING

+

BUILDING

1. Retrieve the data paths
2. Expand the data paths
3. Cluster expanded data paths
4. Sort the paths within a cluster

The Approach

QUERYING

+

CLUSTERING

+

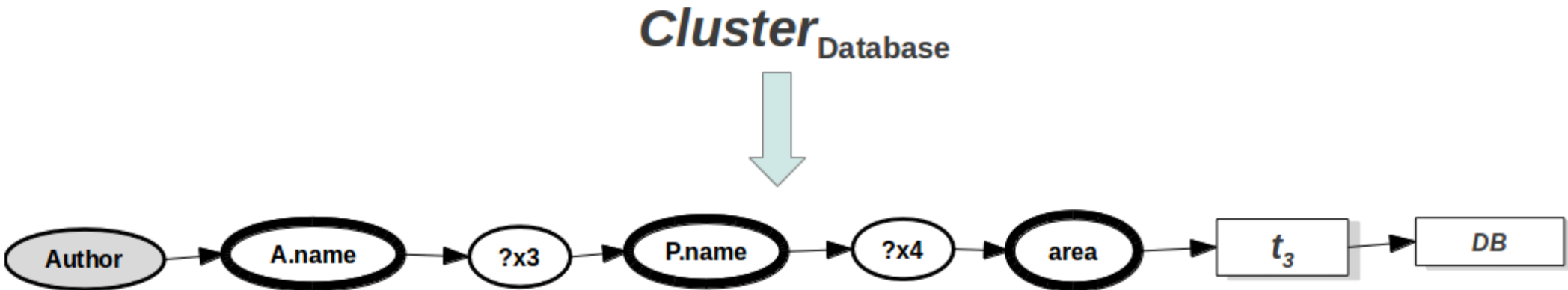
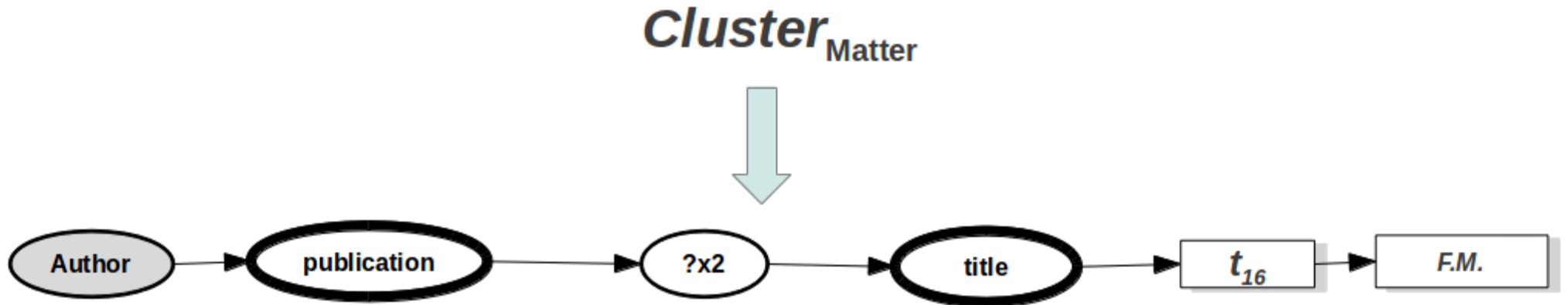
BUILDING

1. Retrieve the data paths
2. Expand the data paths
3. Cluster expanded data paths
4. Sort the paths within a cluster
5. Combine the paths
6. We explore the paths singularly

Clustering

$$\begin{aligned} cl_{Matters} : & \left(\begin{array}{l} dp'_1 : \text{Author} \rightarrow \text{Author.publication} \rightarrow t_{14} \rightarrow F.M. \\ dp'_2 : \text{Publication} \rightarrow \text{Publication.title.} \rightarrow t_{16} \rightarrow F.M. \\ dp'_3 : \text{Author} \rightarrow \text{Author.name} \rightarrow ?x_1 \rightarrow \text{Author.publication} \rightarrow t_{14} \rightarrow F.M. \\ dp'_4 : \text{Author} \rightarrow \text{Author.publication} \rightarrow ?x_2 \rightarrow \text{Publication.title} \rightarrow t_{16} \rightarrow F.M. \\ \dots \end{array} \right) \\ \\ cl_{Database} : & \left(\begin{array}{l} dp'_5 : \text{Person} \rightarrow \text{Person.area} \rightarrow t_1 \rightarrow Db \\ dp'_6 : \text{Person} \rightarrow \text{Person.area} \rightarrow t_2 \rightarrow Db \\ dp'_7 : \text{Person} \rightarrow \text{Person.area} \rightarrow t_3 \rightarrow Db \\ \dots \\ dp'_8 : \text{Author} \rightarrow \text{Author.name} \rightarrow ?x_3 \rightarrow \text{Person.name} \rightarrow ?x_4 \rightarrow \text{Person.area} \rightarrow t_3 \rightarrow Db \\ \dots \end{array} \right) \end{aligned}$$

Building



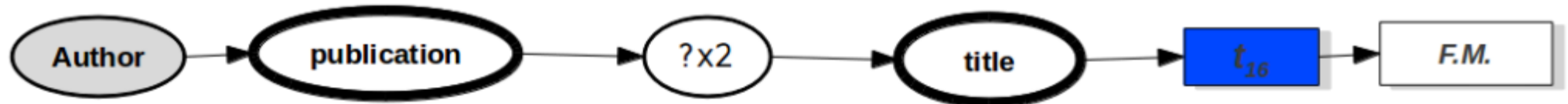
Backward Exploration

T_4 : Author

	<u>name</u>	<u>publication</u>
t_{13}	Lenzerini	Data Integration
t_{14}	Date	Foundation Matters

T_5 : Publication

	<u>title</u>	<u>year</u>
t_{15}	Data Integration	2002
t_{16}	Foundation Matters	2002



$$S = \{t_{16}\}$$

Backward Exploration

T_4 : Author

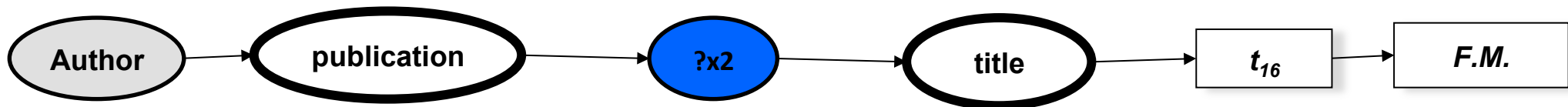
	<u>name</u>	<u>publication</u>
t_{13}	Lenzerini	Data Integration
t_{14}	Date	Foundation Matters

T_5 : Publication

	<u>title</u>	<u>year</u>
t_{15}	Data Integration	2002
t_{16}	Foundation Matters	2002



$\sigma_{\text{publication}='Foundation Matters'}(\text{Author})$



$$S = \{t_{16}\}$$

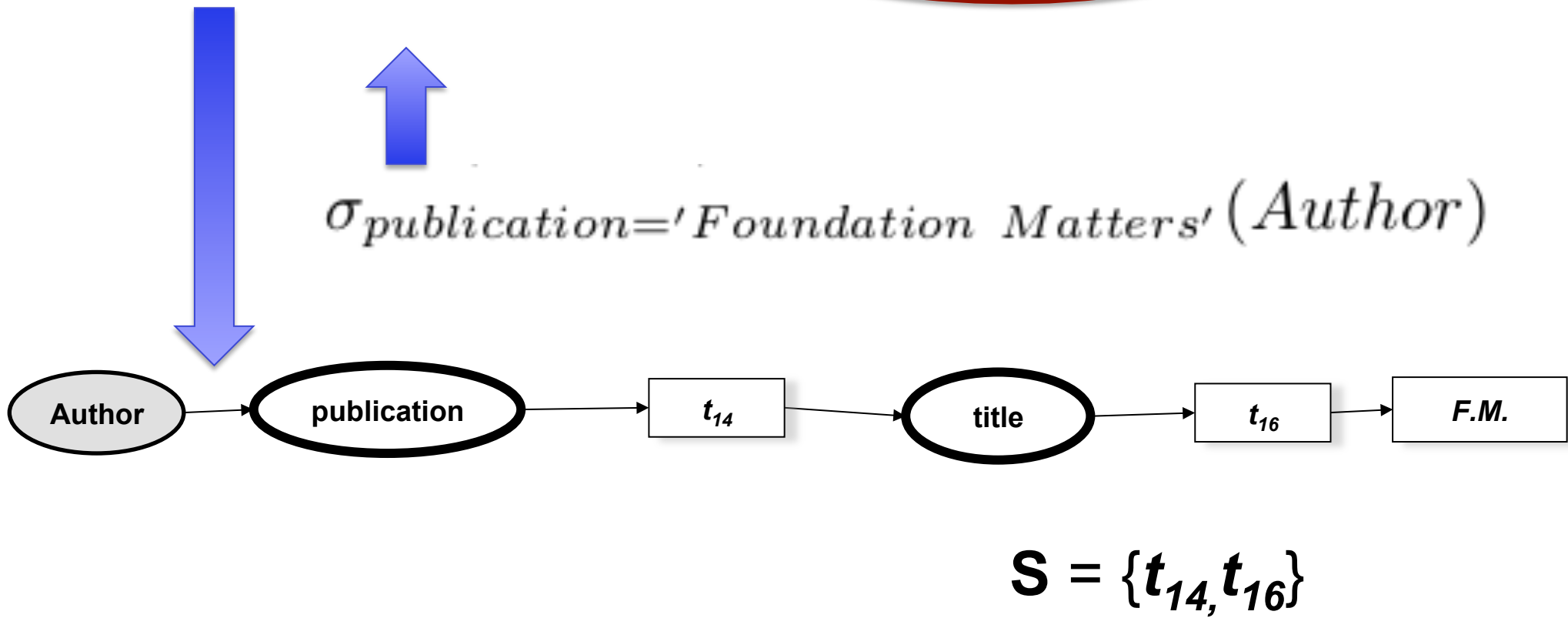
Backward Exploration

T_4 : Author

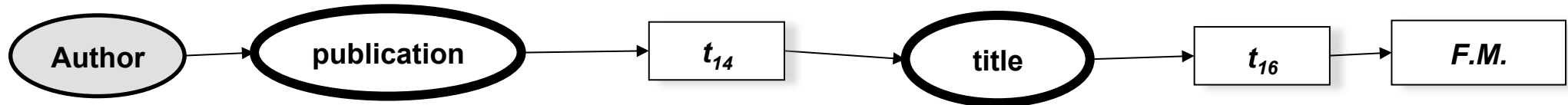
	<u>name</u>	<u>publication</u>
t_{13}	Lenzerini	Data Integration
t_{14}	Date	Foundation Matters

T_5 : Publication

	<u>title</u>	<u>year</u>
t_{15}	Data Integration	2002
t_{16}	Foundation Matters	2002



Backward Exploration

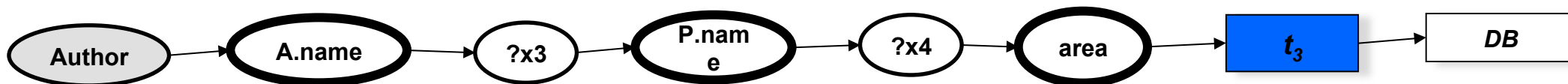


T_4 : Author

	<u>name</u>	<u>publication</u>
t_{13}	Lenzerini	Data Integration
t_{14}	Date	Foundation Matters

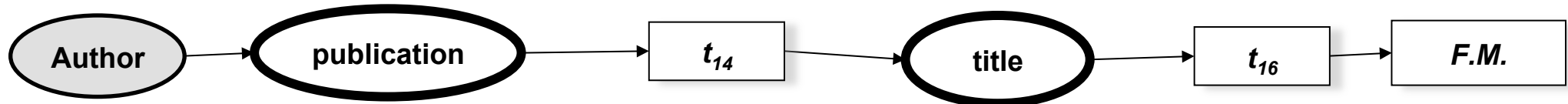
T_1 : Person

	<u>name</u>	<u>area</u>
t_1	Watson	Database
t_2	Lenzerini	Database
t_3	Date	Database
t_4	Hunt	Inf. Systems



$$S = \{t_3, t_{14}, t_{16}\}$$

Backward Exploration



T_4 : Author

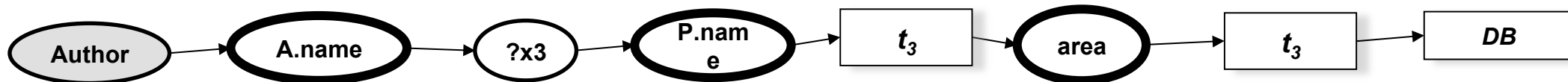
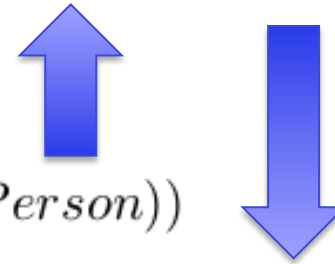
	<u>name</u>	<u>publication</u>
t_{13}	Lenzerini	Data Integration
t_{14}	Date	Foundation Matters

T_1 : Person

	<u>name</u>	<u>area</u>
t_1	Watson	Database
t_2	Lenzerini	Database
t_3	Date	Database
t_4	Hunt	Inf. Systems

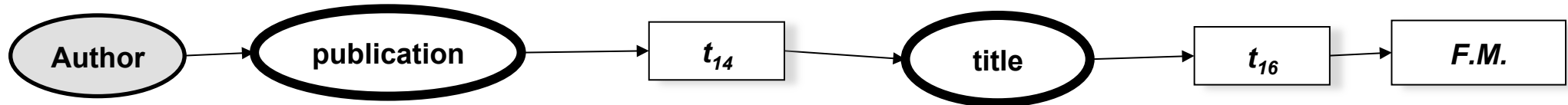


$\pi_{name}(\sigma_{tid=t_3}(Person))$



$$S = \{t_3, t_{14}, t_{16}\}$$

Backward Exploration



T_4 : Author

	<u>name</u>	<u>publication</u>
t_{13}	Lenzerini	Data Integration
t_{14}	Date	Foundation Matters

T_1 : Person

	<u>name</u>	<u>area</u>
t_1	Watson	Database
t_2	Lenzerini	Database
t_3	Date	Database
t_4	Hunt	Inf. Systems

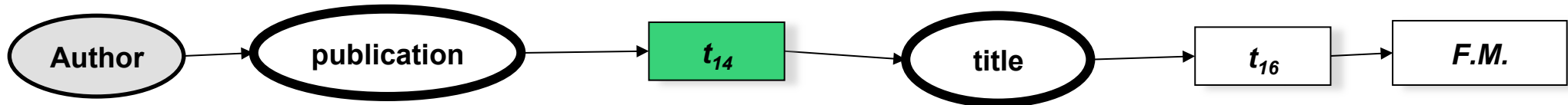


$\sigma_{name='Date'}(Author)$



$$S = \{t_3, t_{14}, t_{16}\}$$

Backward Exploration



T_4 : Author

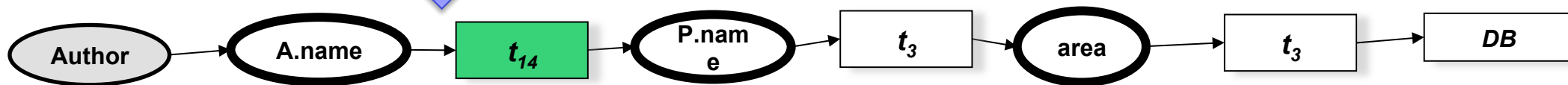
	<u>name</u>	<u>publication</u>
t_{13}	Lenzerini	Data Integration
t_{14}	Date	Foundation Matters

T_1 : Person

	<u>name</u>	<u>area</u>
t_1	Watson	Database
t_2	Lenzerini	Database
t_3	Date	Database
t_4	Hunt	Inf. Systems

$\sigma_{name='Date'}(Author)$

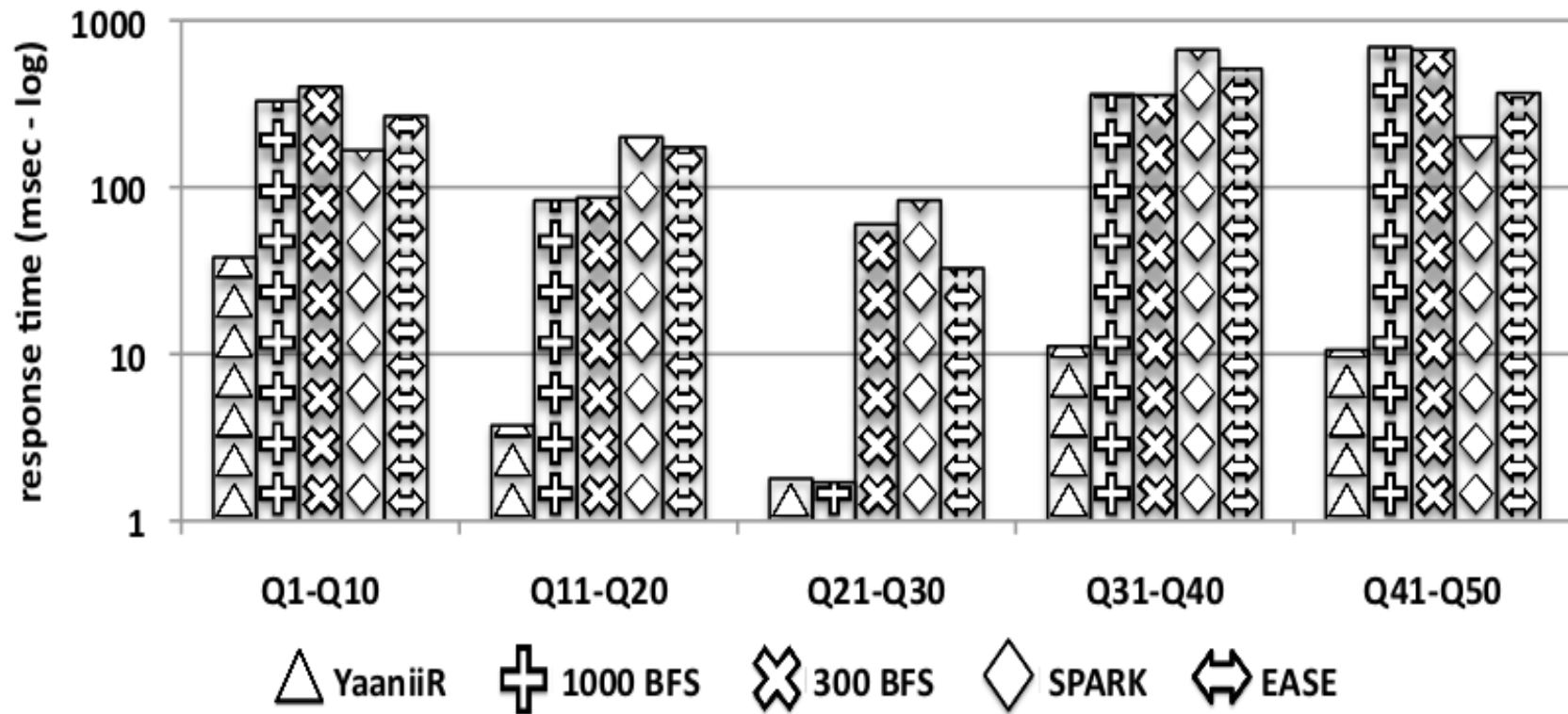
$S = \{t_3, t_{14}, t_{16}\}$



Implementation & Experiments

Implementation in **PostgreSQL 9.1** (*PL/pgSQL*)

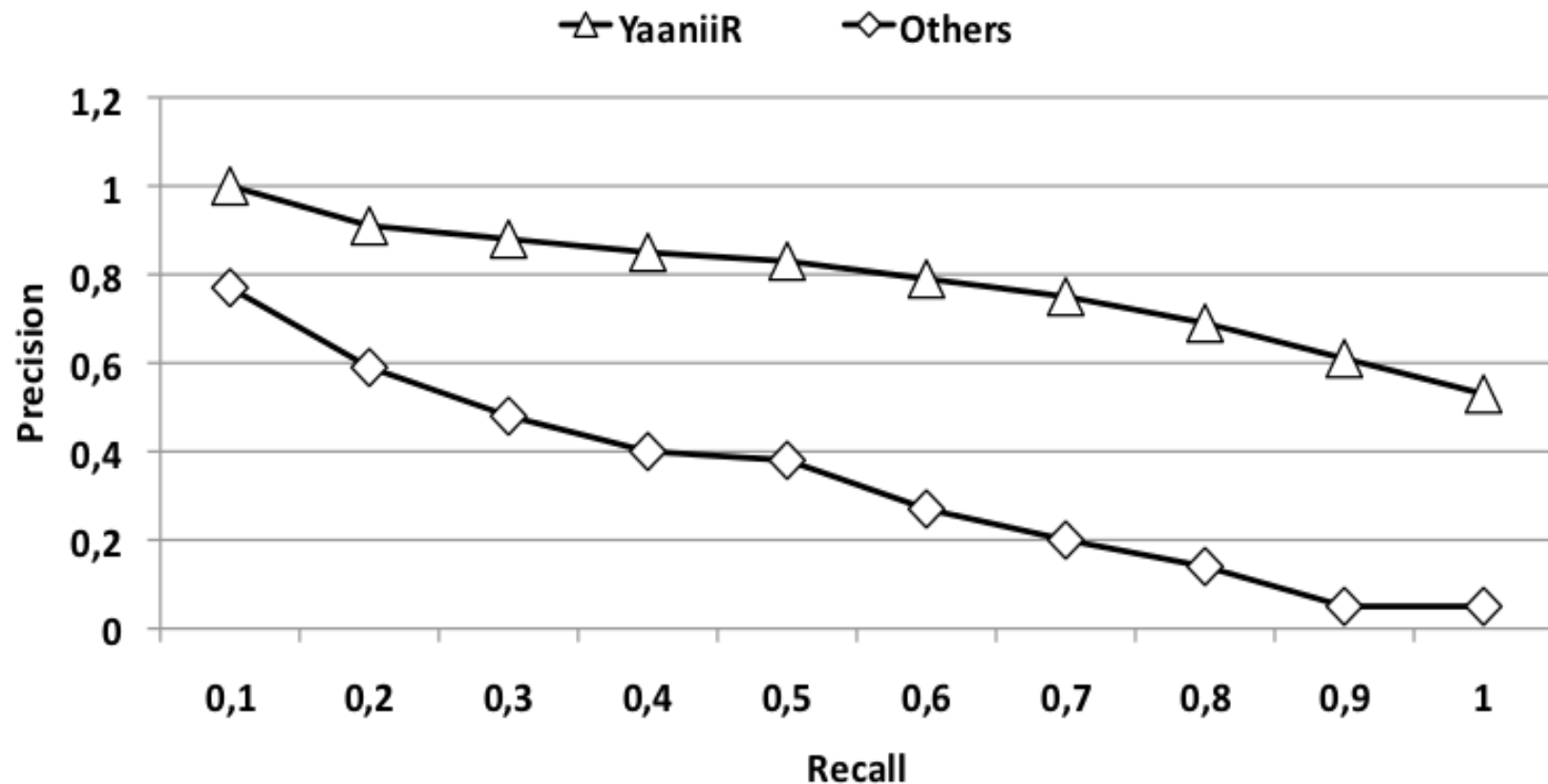
Tested with a common benchmark [*Coffman et al., CIKM 2010*]



Implementation & Experiments

Implementation in **PostgreSQL 9.1** (*PL/pgSQL*)

Tested with a common benchmark [*Coffman et al., CIKM 2010*]



Conclusion

- A strategy to **Structured Keyword Search** that:
 - Works exploiting the **RDBMS** engine
 - Retrieves solutions in order
 - Uses only operations of ***selection*** and ***projection***

Future Works

- *Forward exploration*.
- *Query Execution Plan* comparison with query rewriting approaches.
- Exploit *IR functionalities* on top of the approach.

Thanks for the attention...

