

# exploiting Emergent Schemas to make RDF systems more efficient

Minh-Duc Pham

Peter Boncz

CWI & VU University, Amsterdam



WWW Conference, 2015

## Deriving an Emergent Relational Schema from RDF Data

Minh-Duc Pham<sup>△</sup>  
m.d.pham@vu.nl

Linnea Passing<sup>□</sup>  
passing@in.tum.de

Orri Erling<sup>◇</sup>  
oerling@openlinksw.com

Peter Boncz<sup>○</sup>  
boncz@cwi.nl

<sup>△</sup>Vrije Universiteit Amsterdam, The Netherlands

<sup>□</sup>Technische Universität München, Germany

<sup>◇</sup>OpenLink Software, UK

<sup>○</sup>CWI, The Netherlands

### ABSTRACT

We motivate and describe techniques that allow to detect an “emergent” *relational schema* from RDF data. We show that on a wide variety of datasets, the found structure explains well over 90% of the RDF triples. Further, we also describe technical solutions to the semantic challenge to give short names that humans find logical to these emergent tables, columns and relationships between tables. Our techniques can be exploited in many ways, e.g., to improve the efficiency of SPARQL systems, or to use existing SQL-based applications on top of any RDF dataset using a RDBMS.

0 (subject, property, object) columns<sup>1</sup>. SQL systems tend to be more efficient than triple stores, because the latter need query plans with many self-joins – one per SPARQL triple pattern. Not only are these extra joins expensive, but because the complexity of query optimization is exponential in the amount of joins, SPARQL query optimization is much more complex than SQL query optimization. As a result, large SPARQL queries often execute with a suboptimal plan, to much performance detriment. RDBMS's can further store data efficiently e.g. using advanced techniques such as column-wise compression, table partitioning, materialized views and multi-dimensional data clustering. These

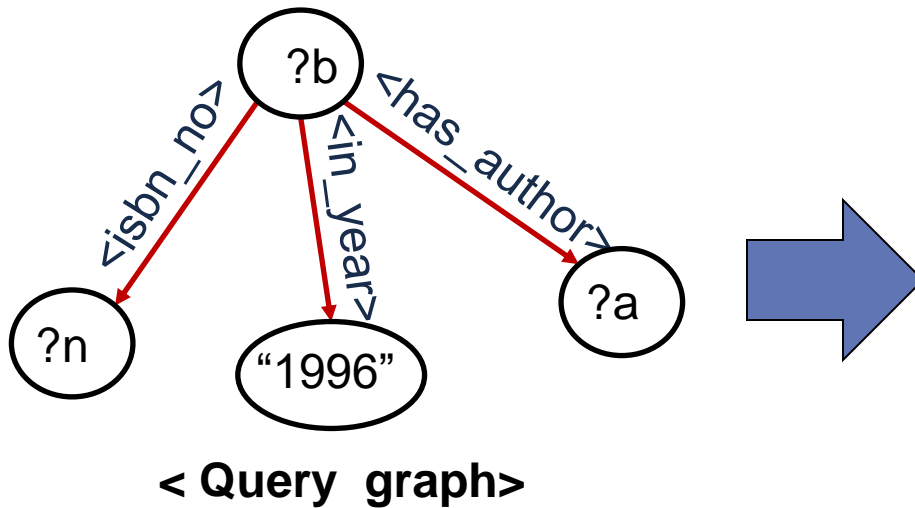
## Summary

# exploiting Emergent Schemas to make RDF systems more efficient

- **the 3 weaknesses of triple stores**
- “schema” confusion between DB and SW practitioners -- but both are right!
- Emergent Schemas to the rescue
- 3 technical details from the ISWC paper
  1. **structured storage is an efficient form of POS storage**
  2. **structure-aware SPARQL execution only as optimization (protects against ?P queries)**
  3. **relation plan does most work, new operator RDFscan adds the missing bindings**

# Triplestore weakness #1

- superfluous joins explode query complexity



- query has unnecessary joins
  - in a relational DB, this is retrieving a record, not a join
  - problem #1: joins are **costly at query execution time**
  - problem #2: query **optimization complexity is  $O(3^N)$**

with **star patterns** size F, **exponentially worse** ( $3^F$ ) optimization space coverage

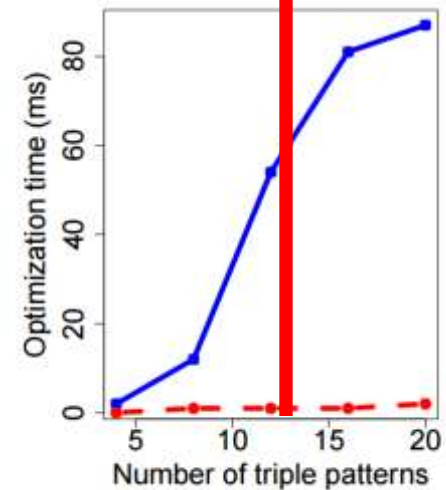
book query:

```
SELECT ?a ?n WHERE {
  ?b    <has_author> ?a .
  ?b    <in_year>    "1996" .
  ?b    <isbn_no>    ?n
}
```

Virtuoso starts cutting  
join order search space  
above 12 patterns: best  
plan may be missed

IdxScan  
<in\_year> "1996"

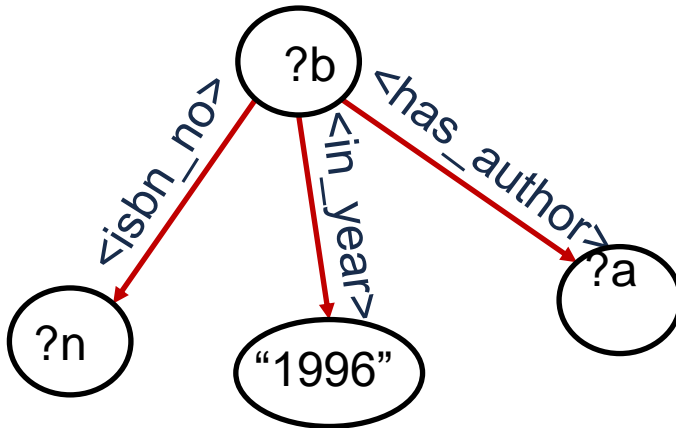
< Example q



(b) Virtuoso

# Triplestore weakness #2

- superfluous joins explode query complexity



< Query graph >

- structural correlations
  - if (?b has an <isbn\_no>) it's a book, and it has <in\_year> and <has\_author>
  - query optimizer estimates using the **independence assumption**
  - problem: result size estimates are wrong → wrong query plan chosen

book query:

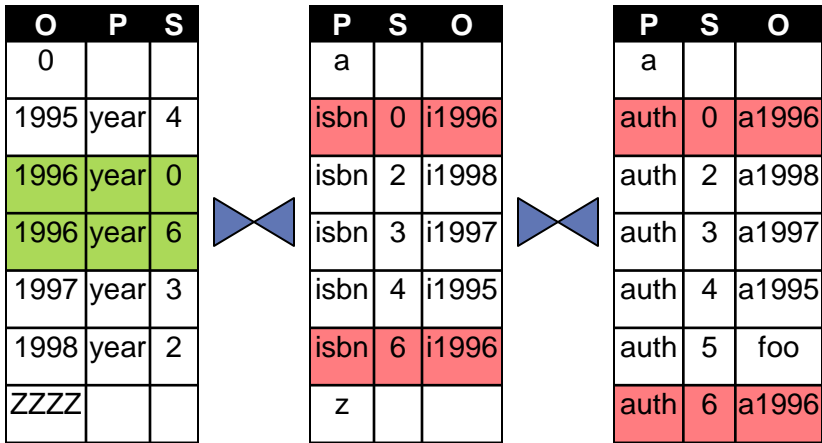
```

SELECT ?a ?n WHERE {
  ?b    <has_author> ?a .
  ?b    <in_year>    "1996" .
  ?b    <isbn_no>    ?n
}
  
```

**Query Optimization Quality of SPARQL is much worse than for SQL**

# Triplestore weakness #3

- RDF architect has no grip on data locality



book query:

```
SELECT ?a ?n WHERE {
  ?b    <has_author> ?a .
  ?b    <in_year>    "1996" .
  ?b    <isbn_no>    ?n
}
```

hexastorage (“indexing all orders”) does **not** save you from lack of locality!!

- relational **clustered index**

year	author	isbn
1975	a1995	i1995
1996	a1996	i1996
1996	a1996	i1996
1997	a1997	i1997
1998	a1998	i1998

- relational **partitioned table**

author	isbn
1995	a1995 i1995
1996	a1996 i1996
1997	a1997 i1997
1998	a1998 i1998

## Summary

# exploiting Emergent Schemas to make RDF systems more efficient

- the 3 weaknesses of triple stores
- **“schema” confusion between DB and SW practitioners -- but both are right!**
- Emergent Schemas to the rescue
- 3 technical details from the ISWC paper
  1. structured storage is an efficient form of POS storage
  2. structure-aware SPARQL execution only as optimization (protects against ?P queries)
  3. relation plan does most work, new operator RDFscan adds the missing bindings

# What does “schema” mean to you?

- Relational answer: the structure definition of a database (**tables + constraints**)
- Semantic web community answer: **Ontologies and Vocabularies**

# What does “schema” mean to you?

- Relational answer: the structure definition of a database (**tables + constraints**)
  - structure of one, particular, database, i.e. the structure of one dataset
  - not intended for reuse, or data integration 😞
  - give the query writer a **clear idea** of what the data looks like
  - must be declared **before** the data can be used → “schema first”
- Semantic web community answer: **Ontologies and Vocabularies**
  - model a knowledge universe so current and future users can denote concepts in a universally understood way in many different contexts → good for **data integration**
  - ontology classes are a poor descriptor of data structure:
    - **partial** use (DBpedia): **<30%** of ontology class attributes occur (on avg)
    - **mixed** use (Dbpedia): subjects combine **7** ontology classes (on avg)
  - symptom: SPARQL query comes back empty
  - **schema can evolve** by e.g. adding new property triples over time → “schema last”



# What does “schema” mean to you?

- Relational answer: the structure definition of a database (**tables + constraints**)
  - structure of one, particular, database, i.e. the structure of one dataset
  - not intended for reuse, or data integration ☹️
  - give the query writer a **clear idea** of what the data looks like
  - must be declared **before** the data can be used → “schema first”
- Semantic web community answer: **Ontologies and Vocabularies**

Both the DB and SW notions of “schema” are valuable!

Semantic Web applications can profit from a DB schema:

1. Users understand datasets better, no more empty SPARQL results
2. Systems can become more efficient (storage, qopt, execution) → ISWC2016

## Summary

# exploiting **Emergent Schemas** to make RDF systems more efficient

- the 3 weaknesses of triple stores
- “schema” confusion between DB and SW practitioners -- but both are right!
- **Emergent Schemas** to the rescue
- 3 technical details from the ISWC paper
  1. structured storage is an efficient form of POS storage
  2. structure-aware SPARQL execution only as optimization (protects against ?P queries)
  3. relation plan does most work, new operator RDFscan adds the missing bindings

# Emergent Schemas

- Detect the “DB-schema” in RDF data automatically!
  - “association rule mining” to find characteristic sets (CS) of co-occurring properties
  - merge similar CSs using multiple methods into **few tables**
  - derive human-friendly names for tables and columns (exploiting ontologies)

## WWW Conference, 2015. Pham, Passing, Erling, Boncz Deriving an Emergent Relational Schema from RDF

Datasets	Number of tables			Coverage – Metric $C$ (%)		
	before merging	after merging	remove small tables	remove small tables	prune infreq. prop.	final schema
LUBM	17	13	12	100	100	100.00
BSBM	49	8	8	100	100	100.00
SP2B	554	13	10	99.99	99.65	99.65
MusicBrainz	27	12	12	100	99.9	99.60
EuroStat	44	10	5	99.73	99.53	99.53
DBLP	249	9	6	100	99.68	99.60
PubMed	3340	14	12	100	99.75	99.73
WebData.	13354	3000	253	98.17	94.37	92.79
DBpedia	439629	542	234	99.12	96.68	95.82

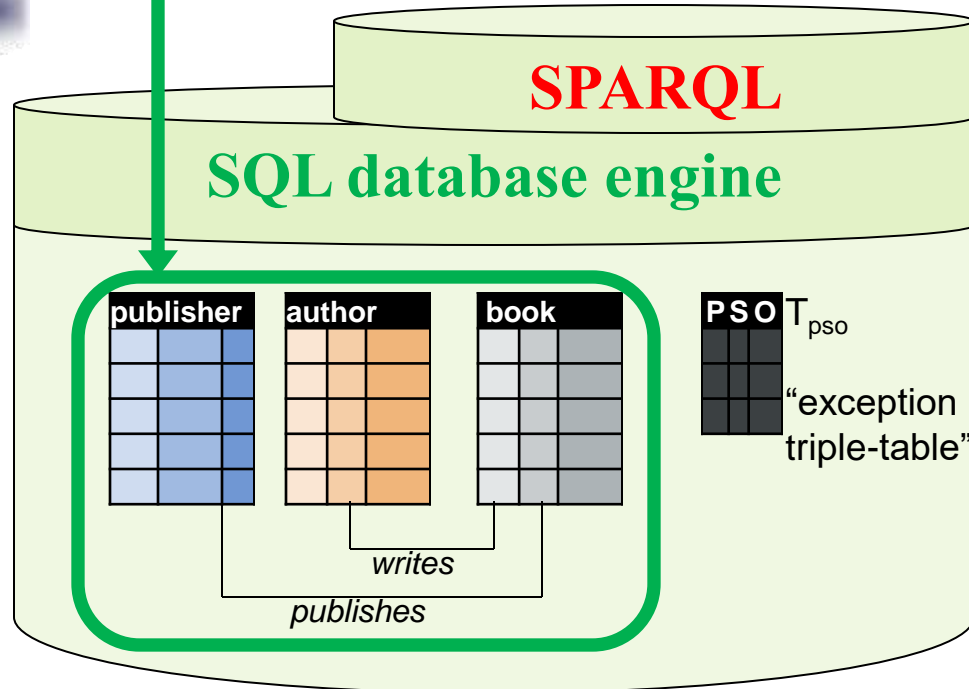
Table 4: #tables and metric  $C$  after merging & filtering

# Emergent Schemas

- Detect the “DB-schema” in RDF data automatically!
  - “association rule mining” to find characteristic sets (CS) of co-occurring properties
  - merge similar CSs using multiple methods into **few tables**
  - derive human-friendly names for tables and columns (exploiting ontologies)

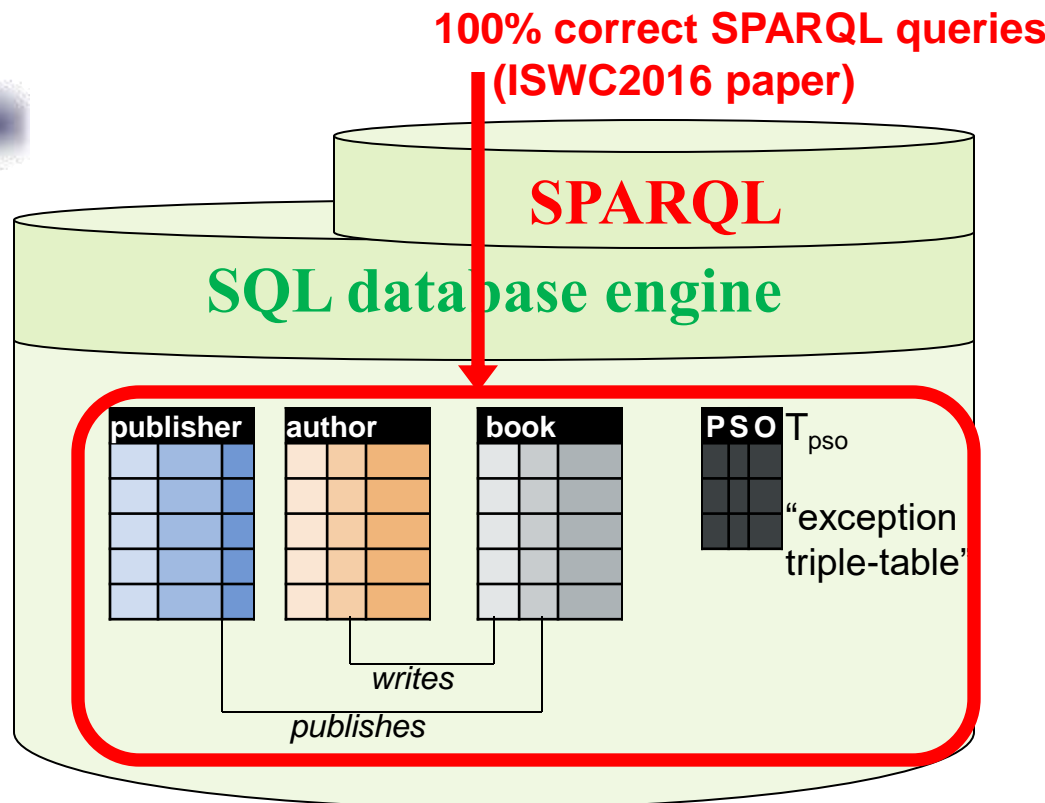


SQL queries on >95% of the RDF  
(WWW2015 paper)



# Emergent Schemas

- Detect the “DB-schema” in RDF data automatically!
  - “association rule mining” to find characteristic sets (CS) of co-occurring properties
  - merge similar CSs using multiple methods into **few tables**
  - derive human-friendly names for tables and columns (exploiting ontologies)



## Summary

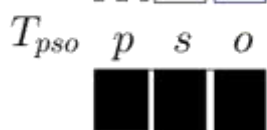
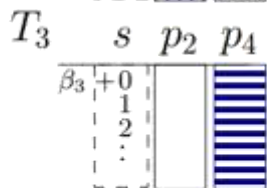
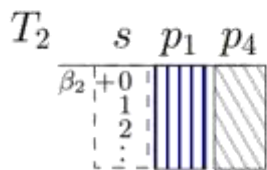
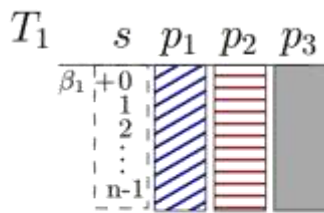
# exploiting Emergent Schemas to make RDF systems more efficient

- the 3 weaknesses of triple stores
- “schema” confusion between DB and SW practitioners -- but both are right!
- the semantic web is mostly a bunch of tables (is it really very graphy?)
- Emergent Schemas to the rescue
- 3 technical details from **the ISWC paper**
  1. structured storage is an efficient form of POS storage
  2. structure-aware SPARQL execution only as optimization (protects against ?P queries)
  3. relation plan does most work, new operator RDFscan adds the missing bindings

# Emergent Table – Smart Column Storage

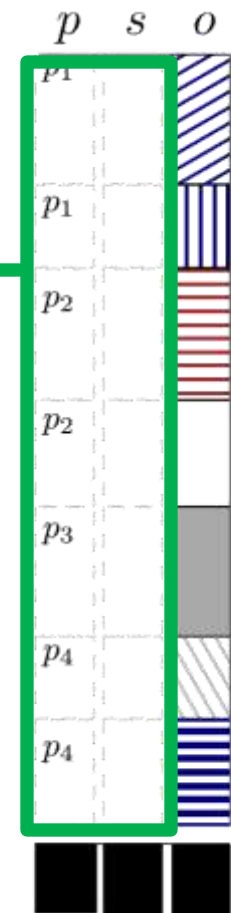
- URI → integer OID mapping (standard technique)
  - we manipulate this dictionary during emergent schema detection
  - **subject URIs** belonging to table  $i$  get **dense** numbers starting at  $\beta_i$
  - $\beta_i = i * 2^{40}$  (really large number  
so OIDs from different tables do not overlap)

PSO union view



**dense sequences compress away**  
 otherwise the s-column would be ordered integers  
**dense sequences allow array lookup**  
 otherwise mergejoins would be needed

- no B-tree or HashMap needed
- makes new RDFscan operator (up next..) cheap

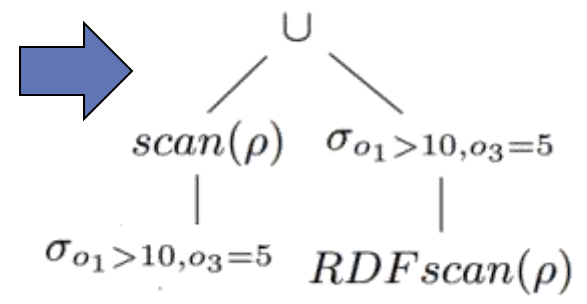


- 1. structured storage is an efficient form of POS storage
- 2. structure-aware SPARQL execution only as optimization (protects against ?P queries)

# The RDFscan operator

- the strategy is to replace SPARQL **star patterns** with **table scan(s)**
  - re-using scan is good: RDBMS pushes down selections
    - it triggers index selections, partition pruning,...(goodness)
  - this would only return 95% of the data. How to generate the rest?
- the new RDFscan operator

```
SELECT ?s ?o1 ?o2 WHERE { ?s p1 ?o1 .
                          ?s p2 ?o2 .
                          ?s p3 5 . FILTER (?o1 > 10) }
```



$T_1$				$T_2$				$T_{pso}$			Result		
s	p1	p2	p3	s	p1	p3	p4	p	s	o	s	o1	o2
100	11	2	5	200	11	7	1	p1	0	20	100	11	2
101	13	4	6	201		5	2	p1	1	9	104	15	8
102	14		5	202	13	9	3	p1	201	15	0	20	8
103	9	6						p2	0	8	102	14	6
104	15	8	5					p2	102	6	201	15	4
								p2	201	4			
								p3	0	5			
								p6	6	7			

**RDFscan**

Fig. 8: Example RDF data and expected query result.



## Summary

# exploiting Emergent Schemas to make RDF systems more efficient

- the 3 weaknesses of triple stores
  1. superfluous joins explode query complexity
  2. structural correlations destroy estimates
  3. architect has no grip data locality (no data clustering, no table partitioning, ...)
- SCHEMA: confusion between DB and SW practitioners -- but both are right!
- Emergent Schemas to the rescue
  - neutralizes the 3 triple-store weaknesses!
  - = DB notion of schema for RDF data
- 3 technical details from the ISWC paper
  1. structured storage is an efficient form of POS storage
  2. structure-aware SPARQL execution only as optimization (protects against ?P queries)
  3. relation plan does most work, new operator RDFscan adds the missing bindings

## Summary

# exploiting Emergent Schemas to make RDF systems more efficient

- the 3 weaknesses of triple stores
  1. superfluous joins explode query complexity
  2. structural correlations destroy estimates
  3. architect has no grip data locality (no data clustering, no table partitioning, ...)
- SCHEMA: confusion between DB and SW practitioners -- but both are right!

• Both the DB and SW notions of “schema” are valuable!

- n  
= l

Semantic Web applications can profit from a DB schema:

- 3
  1. Users understand datasets better, no more empty SPARQL results
  2. Systems can become more efficient (storage, query, execution) → ISWC2016

**DB practitioners can also profit from a SW schema!**

1. Tables, columns, constraints should have URIs, link to SW-schemas.
2. Keys could be URIs. Extend SQL! Query SQL+RDF datasets in SQL.