

Hierarchical Property Set Merging for SPARQL Query Optimization

Marios Meimaris, Athena Research Center, Greece



George Papastefanatos, Athena Research Center, Greece

Panos Vassiliadis, University of Ioannina, Greece

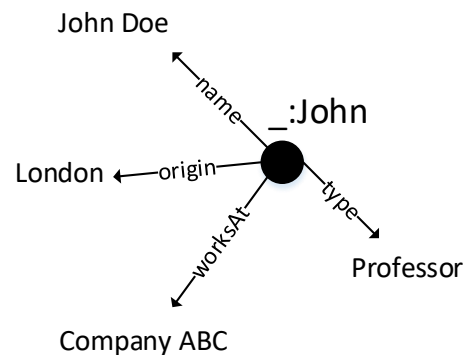


Preliminaries

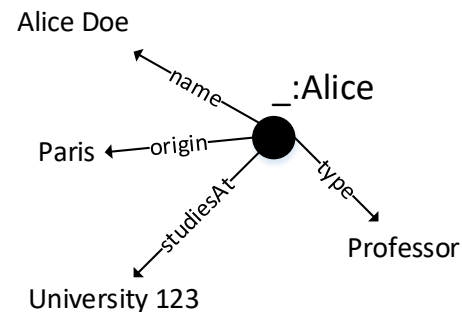
- RDF (Resource Description Framework)
 - Abstract Data model for Linked Data
 - Based on *Triples*: Subject-Predicate-Object
 - RDF datasets are *Directed Labelled Graphs*
- Characteristic Set (CS)
 - A CS is a set of properties with the same subject as source node
 - An RDF dataset can be described as a set of unique CSs
 - Each CS is an *implicit resource type*

Preliminaries

- › Use Characteristic Sets (CSs) and their links in order to store and index triples
- › **Characteristic Sets** (Neumann & Moerkotte, ICDE 2011)
 - › A *Characteristic Set* (CS) S_c of a node x is defined as the set of properties emitting from x (i.e., x as subject)



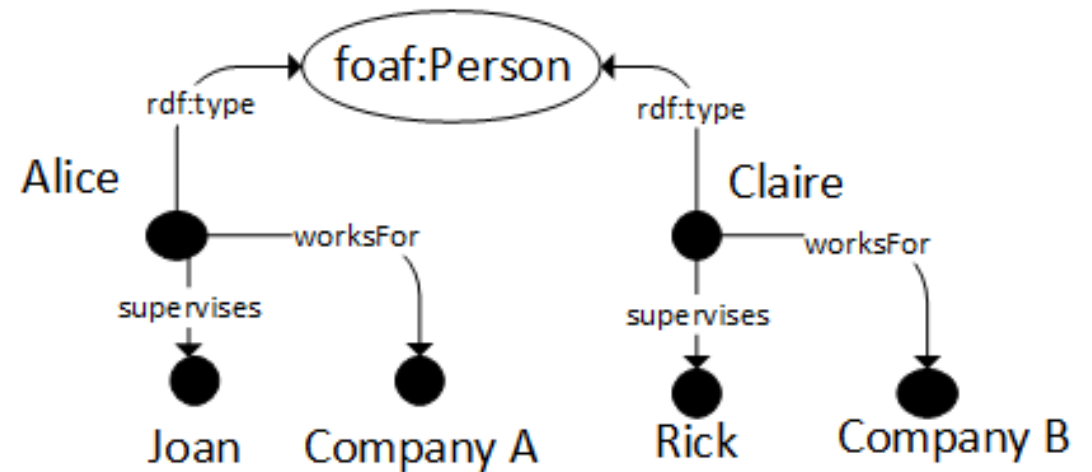
$$S_c(\text{John}) = \{\text{name, origin, worksAt, type}\}$$



$$S_c(\text{Alice}) = \{\text{name, origin, studiesAt, type}\}$$

Background

- › Derive a **relational representation** of an RDF dataset
- › Use CSs as tables and links between CSs as relationships
- › CS properties \rightarrow relation attributes

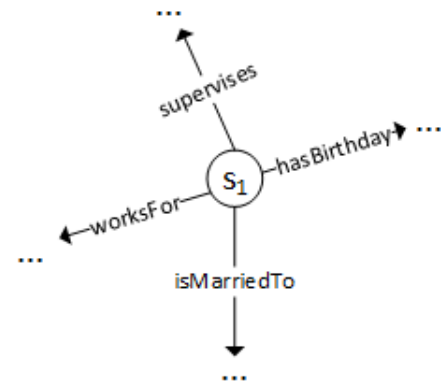


$c_1 = \{\text{rdf:type}, \text{worksFor}, \text{supervises}\}$

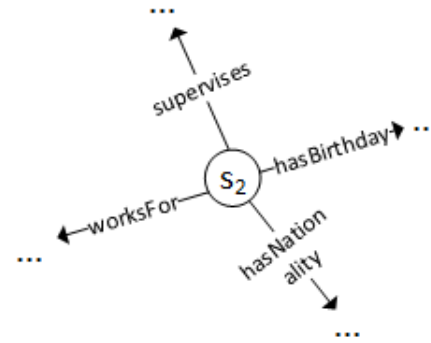
id	rdf:type	worksFor	supervises
Alice	foaf:Person	Company A	Joan
Claire	foaf:Person	Company B	Rick

Problem statement

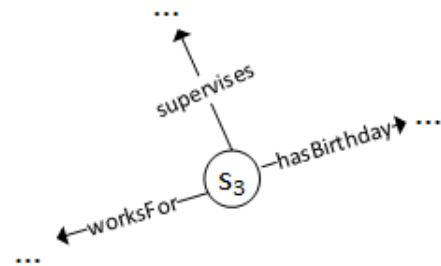
RDF structural looseness \rightarrow multiple CSs \rightarrow different representation strategies



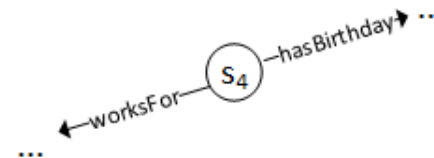
$c_1 = \{\text{worksFor, supervises, hasBirthday, isMarriedTo}\}$



$c_2 = \{\text{worksFor, supervises, hasBirthday, hasNationality}\}$

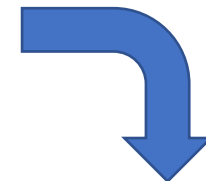
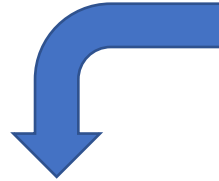
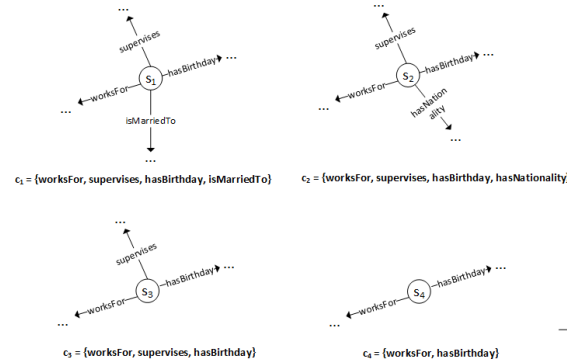


$c_3 = \{\text{worksFor, supervises, hasBirthday}\}$



$c_4 = \{\text{worksFor, hasBirthday}\}$

Trade-Off for creating a relational schema



A relational table for each different CS

A "universal" table for all CS's

id	supervises	worksFor	hasBirthday	isMarriedTo
S_1

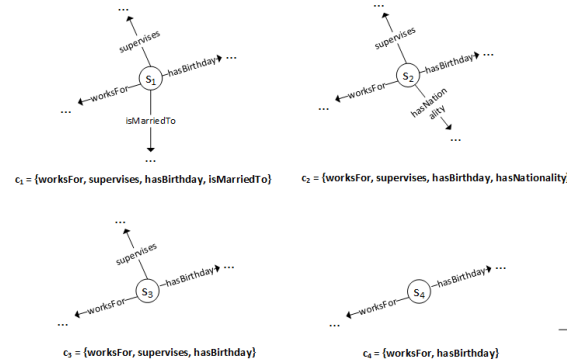
id	supervises	worksFor	hasBirthday	hasNationality
S_2

id	supervises	worksFor	hasBirthday
S_3

id	worksFor	hasBirthday
S_4

id	supervises	worksFor	hasBirthday	isMarriedTo	hasNationality
S_1	NULL
S_2	NULL	...
S_3	NULL	NULL
S_4	NULL	NULL	NULL
...

Trade-Off for creating a relational schema



A relational table for each different CS

A "universal" table for all CS's

id	supervises	worksFor	hasBirthday	isMarriedTo
S ₁

id	supervises	worksFor	hasBirthday	hasNationality
S ₁	NULL

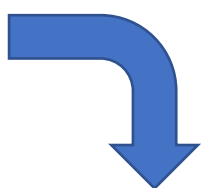
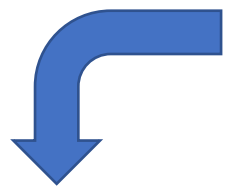
- Space efficient
- Large numbers of relational tables with few tuples in
- Too many joins to answer queries

id	supervises	worksFor	hasBirthday	isMarriedTo	hasNationality
S ₁	NULL

- Captures all CSs into a single table
- Too many NULL values
- Space inefficient

Trade-Off for answering a complex SPARQL query with many joins

```
SELECT ? x ?y ?z ?w
WHERE { ?x worksFor ?y .
        ?x supervises ?z .
        ?z hasBirthday '2011-02-24' .
        ?z isMarriedTo ?w .
        ?w hasNationality 'GR' }
```



A relational table for each different CS

A "universal" table for all CS's

id	supervises	worksFor	hasBirthday	isMarriedTo
S1

id	supervises	worksFor	hasBirthday	isMarriedTo	hasNationality
S1	NULL
S2	NULL	...
S3	NULL	NULL
S4	NULL	NULL	NULL

One self-join for each one of the *worksFor* , *supervises* and *isMarriedTo* query conditions

Additionally three joins between each CS table and all other CS tables in the database – i.e., 4 joins per table.

One self-join for each one of the *worksFor* , *supervises* and *isMarriedTo* query conditions

Problem to be solved

Context: Mapping heterogeneous RDF datasets to a relational schema with the aim to facilitate the processing of complex analytical SPARQL queries

Solution: automating the decision of which tables will be created for a set of CS, such that there are **no overly empty tables and extremely large numbers of joins.**

Observations

- › Based on previous findings:
 - › CS number is generally low but exhibits skewed distribution
 - › E.g., many CSs with very few (<10) subjects
 - › CS number affects number of joins
- › **Merging** closely related CSs helps storage & querying
 - › Less CSs means less joins
 - › Less CSs means less I/O costs in disk-based systems
 - › Compact schema easier to understand and maintain
- › CSs are hierarchical, i.e., their property sets can be super/subsets of each other
- › **Challenge:** *exploit the hierarchical structure in order to merge together closely related CSs*

Challenge

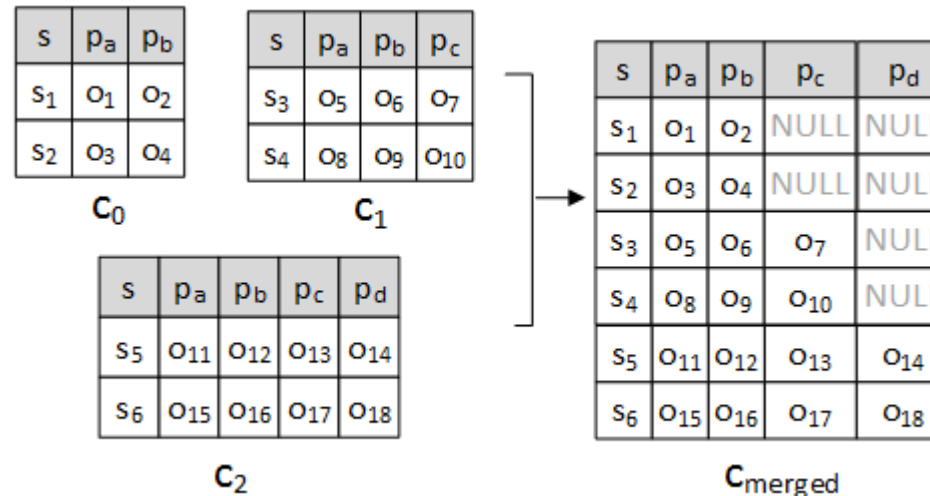
- › Each CS defines a relational table $(s, p_1, p_2, \dots, p_k)$
- › Merging of CS tables results in NULL values for non-shared attributes
- › Challenge: merge CSs and reduce NULL value effect

e.g.:

$c_0 = \{\text{name, age}\}$

$c_1 = \{\text{name, age, marriedTo}\}$

$c_2 = \{\text{name, age, marriedTo, worksAt}\}$



Approach

- › Use a **dense child** table and merge its parents into it
 - › **Why dense?** -> # of NULLs is proportional to # of records of table to be merged
 - › **Why child?** -> more specialized, thus will contain columns of parents
- › Identify dense CSs
 - › if $|c_i| > m \times |c_{max}|$ parameter $\Rightarrow c_i$ is dense
 - › Every resulting (merged) table will contain **exactly one dense node** (and several non-dense)
- › Find optimal merging of ancestors to dense child CSs

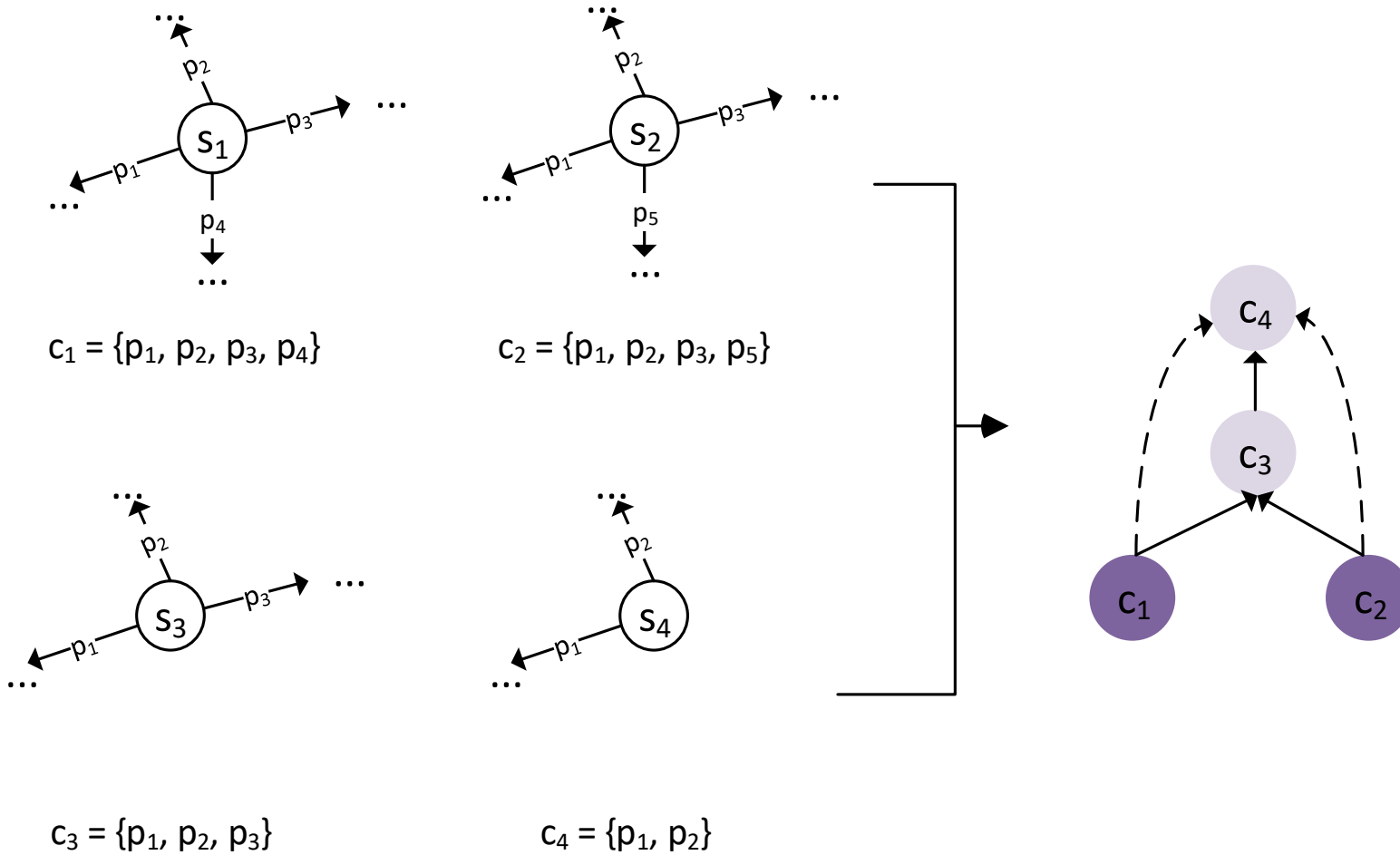
e.g.

$c_1: \{\text{name, age, address}\}, c_2: \{\text{name, age}\}: c_1$ child of c_2

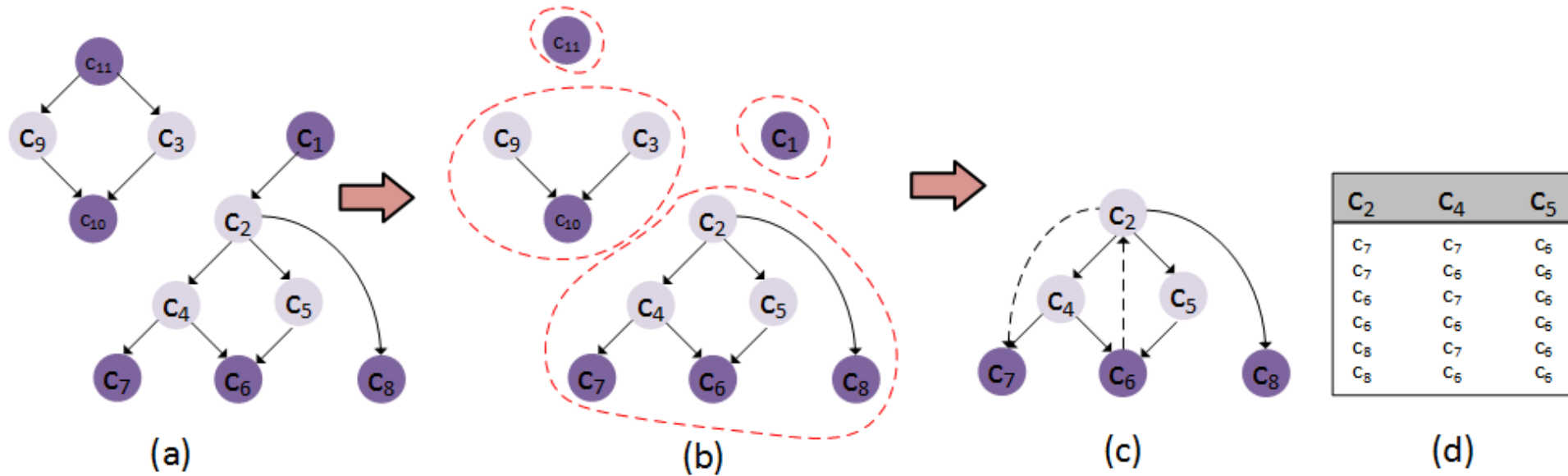
$\text{hier_merge}(c_1, c_2) = c_{12}$

$c_{12}: \{\text{name, age, address}\}$

CS Graph Example



Approach - Example



Approach – Loading and Merging

- › Finding the optimal solution is equivalent to enumerating all possible sub-graphs -> *exponential*
- › Greedy approximation
 - › At each step, merge parent CS and dense child CS that minimize objective cost function
 - › Cost function minimizes the number of NULL values introduced by the merge
- › Tuning of m parameter

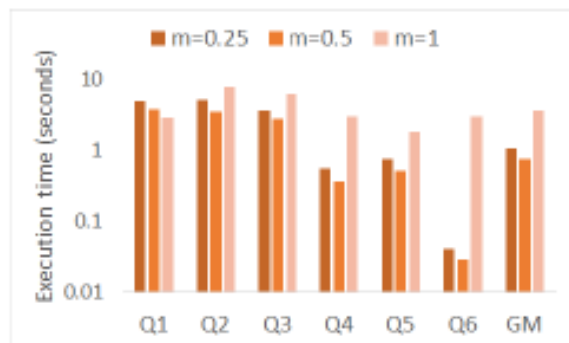
Approach – Querying

- › Parse incoming SPARQL queries
 - › Identify query CSs that match merged CSs in the dataset
 - › Rewrite query as an SQL statement with UNIONS between matched CSs
 - › In case of SO/OS joins, prune off CSs that are not linked
- › Pass final query to relational optimizer
- › Build and output results

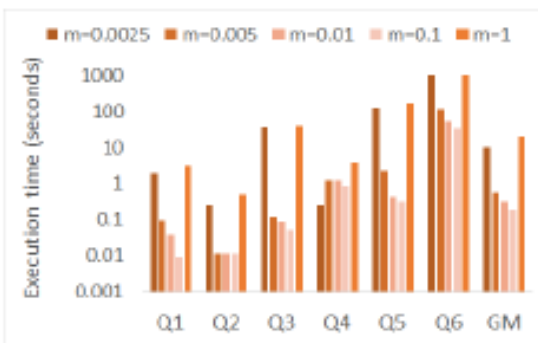
Implementation & Evaluation (Loading)

Dataset	Size (MB)	Time	# Tables (CSs)	# of ECSs	Dense CS Coverage
Reactome Simple	781	3min	112	346	100%
Reactome (m=0.05)	675	4min	35	252	97%
Reactome (m=0.25)	865	4min	14	73	77%
Geonames Simple	4991	69min	851	12136	100%
Geonames (m=0.0025)	4999	70min	82	2455	97%
Geonames (m=0.05)	5093	91min	19	76	87%
Geonames (m=0.1)	5104	92min	6	28	83%
LUBM Simple	591	3min	14	68	100%
LUBM (m=0.25)	610	3min	6	21	90%
LUBM (m=0.5)	620	3min	3	6	58%
WatDiv Simple	4910	97min	5667	802	100%
WatDiv (m=0.01)	5094	75min	67	99	77%
WatDiv (m=0.1)	5250	75min	25	23	63%
WatDiv (m=0.5)	5250	77min	16	19	55%

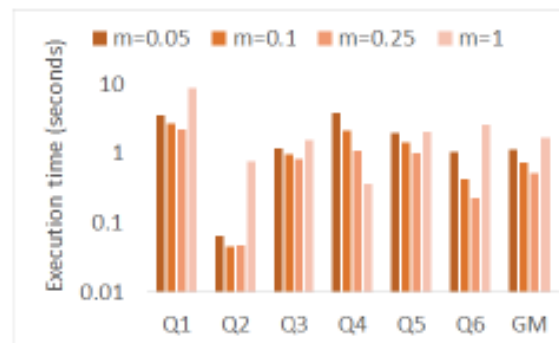
Implementation & Evaluation (Querying)



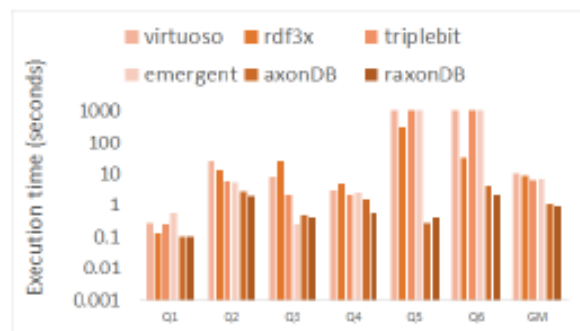
(a) Execution time (seconds) for LUBM



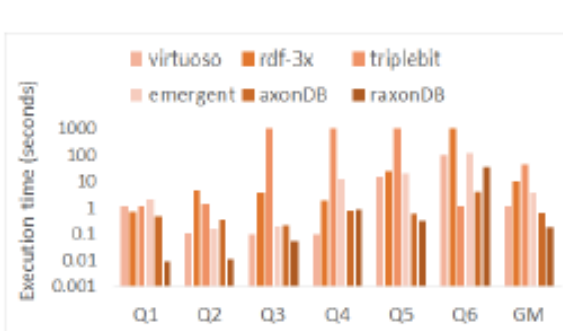
(b) Execution time (seconds) for Geonames



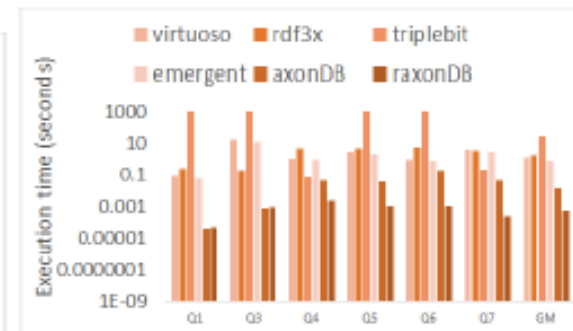
(c) Execution time (seconds) for Reactome



(a) Execution time (seconds) for LUBM2000



(b) Execution time (seconds) for Geonames



(c) Execution time (seconds) for Reactome

Future Work

- › Distributed version of *raxonDB*
 - › CS-based partitioning scheme
 - › Distributed query processing
- › Refined cost function
- › Different ways of defining density

Thank you

{m.meimaris, gpapas}@athenarc.gr, pvassil@cs.uoi.gr

<https://github.com/mmeimaris/raxonDB>

<https://visualfacts.imsi.athenarc.gr/>



This research is funded by the project [VisualFacts](#) (#1614) - 1st Call of the Hellenic Foundation for Research and Innovation Research Projects for the support of post-doctoral researchers.