

RDF Analytics: Lenses over Semantic Graphs

Dario Colazzo, François Goasdoué, Ioana Manolescu, Alexandra Roatis.

23rd International World Wide Web Conference,
Apr 2014, Seoul, South Korea

Presented by Kim A. Jakobsen (kjakob09@cs.aau.dk)
Department of Computer Science, Aalborg University
Study Circle

November 21, 2014

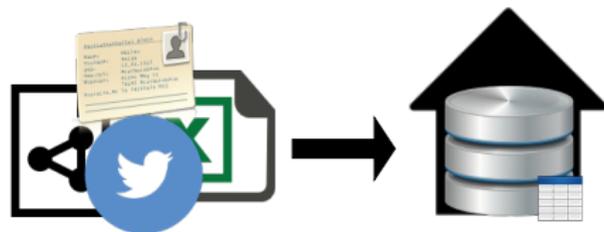
Motivation

Problem



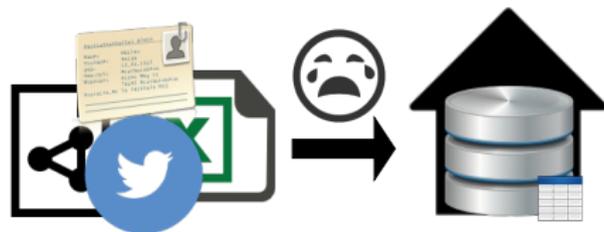
Motivation

Problem



Motivation

Problem



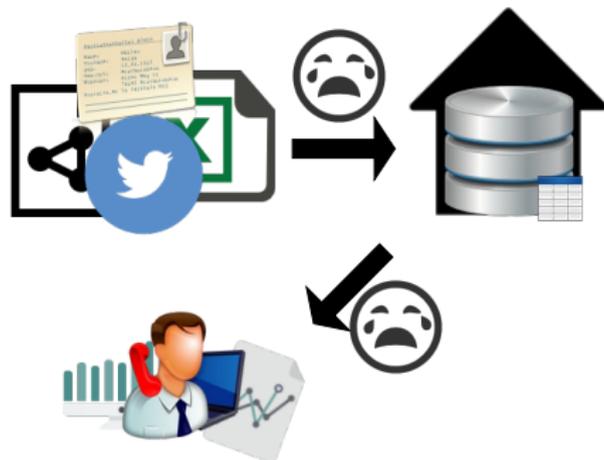
Motivation

Problem



Motivation

Problem



Motivation

WaRG: Warehousing RDF Graphs

- Full-RDF warehouse
- No single central concept
- Flexible choice of measures and classifiers
- Support cube analysis operations (slice, dice, roll-up ...)



Motivation

WaRG: Warehousing RDF Graphs

- Full-RDF warehouse
- No single central concept
- Flexible choice of measures and classifiers
- Support cube analysis operations (slice, dice, roll-up ...)



Outline

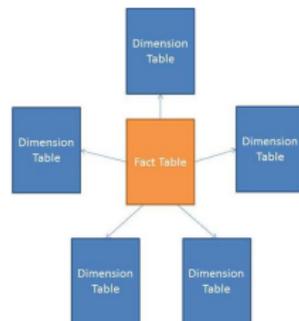
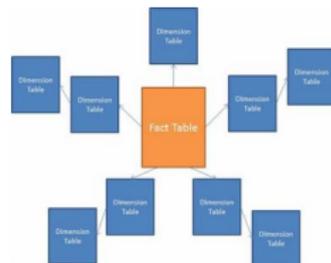
- 1 Motivation
- 2 Analytical Schema
- 3 Analytical Queries
- 4 Analytical Queries Answering
- 5 Experiments
- 6 Conclusion
- 7 Strengths and Weaknesses
- 8 Questions?

- 1 Motivation
- 2 Analytical Schema**
- 3 Analytical Queries
- 4 Analytical Queries Answering
- 5 Experiments
- 6 Conclusion
- 7 Strengths and Weaknesses
- 8 Questions?

Analytical Schema

- The lense through which the underlying data is observed
- Abstraction over the instance data
- Corresponds to the relation schema
- Flexible
- RDF Graph

- The lense through which the underlying data is observed
- Abstraction over the instance data
- Corresponds to the relation schema
- Flexible
- RDF Graph



Images borrowed from <http://social.technet.microsoft.com/>

Analytical Schema

Example

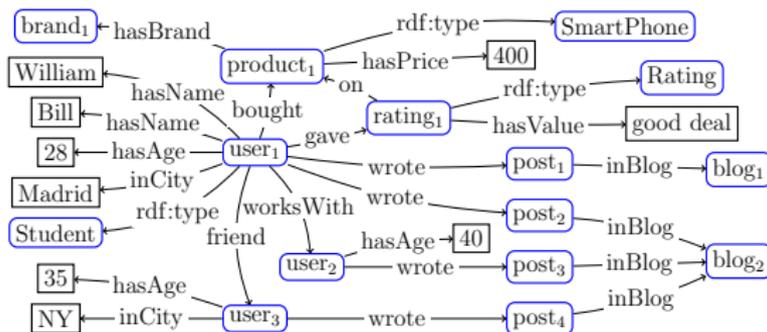


Figure 2: Running example: RDF graph.

Analytical Schema

Example

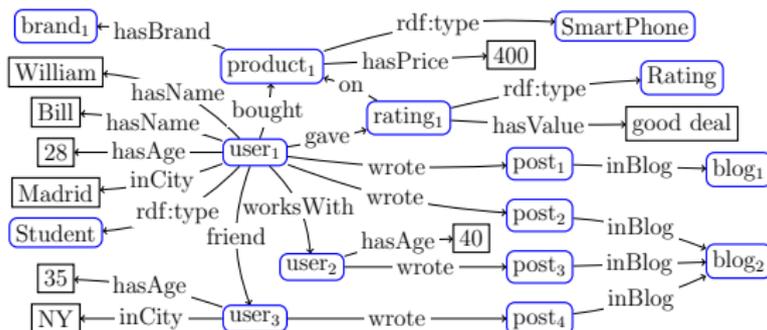


Figure 2: Running example: RDF graph.

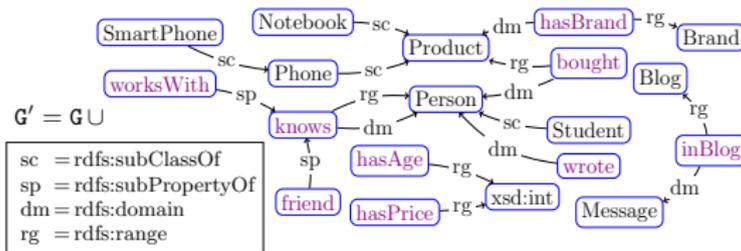
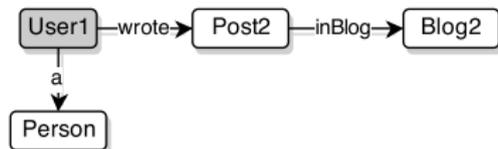


Figure 3: Running example: RDF Schema triples.

Analytical Schema

Nodes and Edges

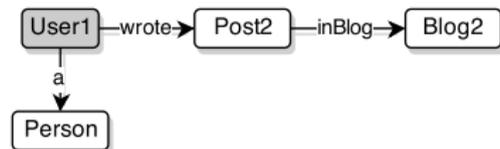
Blogger = $q(x):- x \text{ rdf:type } Person, x \text{ wrote } y, y \text{ inBlog } z$



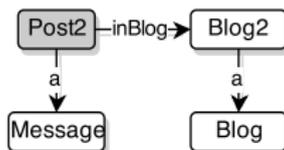
Analytical Schema

Nodes and Edges

Blogger = $q(x):- x \text{ rdf:type Person, } x \text{ wrote } y, y \text{ inBlog } z$



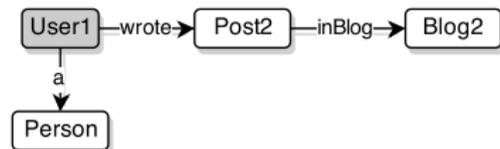
BlogPost = $q(x):- x \text{ rdf:type Message, } x \text{ inBlog } z, z \text{ rdf:type Blog}$



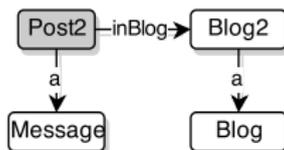
Analytical Schema

Nodes and Edges

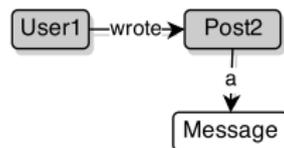
$\text{Blogger} = q(x):- x \text{ rdf:type Person}, x \text{ wrote } y, y \text{ inBlog } z$



$\text{BlogPost} = q(x):- x \text{ rdf:type Message}, x \text{ inBlog } z, z \text{ rdf:type Blog}$



$\text{wrotePost} = q(x, y):- x \text{ wrote } y, y \text{ rdf:type Message}$



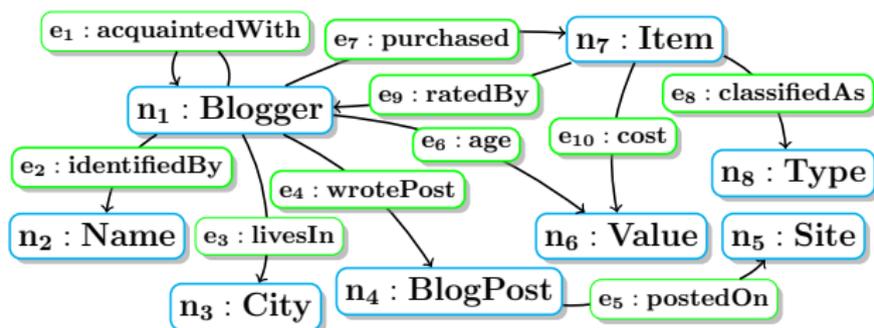


Figure 4: Sample Analytical Schema (*AnS*).

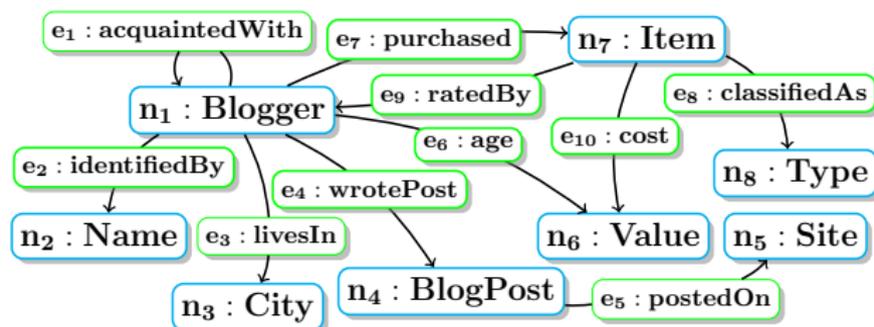


Figure 4: Sample Analytical Schema (*AnS*).

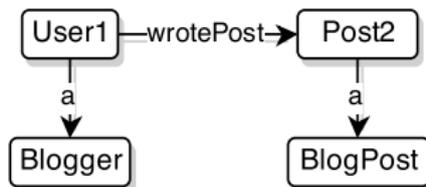


Figure: Analytical Schema Instance Data sample

Outline

- 1 Motivation
- 2 Analytical Schema
- 3 Analytical Queries**
- 4 Analytical Queries Answering
- 5 Experiments
- 6 Conclusion
- 7 Strengths and Weaknesses
- 8 Questions?

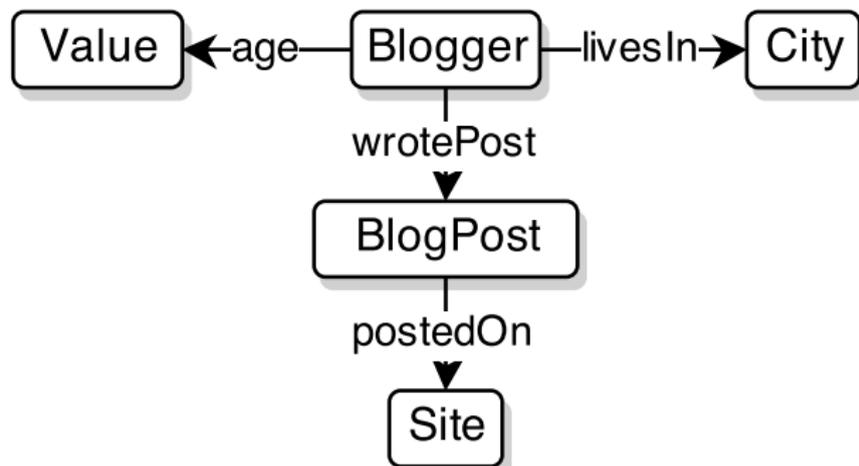
- Analytical Query (AnQ):
 - Classifiers / Dimensions
 - Measures
 - Aggregation function

- Analytical Query (AnQ):
 - Classifiers / Dimensions
 - Measures
 - Aggregation function
- Number of sites where each blogger posts, classified by the bloggers's age and city
- $\langle c(x, y_1, y_2), m(x, z), \oplus \rangle$
 - $c(x, y_1, y_2)$:- x age y_1 , x livesIn y_2
 - $m(x, z)$:- x wrotePost y , y postedOn z
 - \oplus :- *count*

Analytical Queries / Cubes

Example

- Number of sites where each blogger posts, classified by the bloggers's age and city



Outline

- 1 Motivation
- 2 Analytical Schema
- 3 Analytical Queries
- 4 Analytical Queries Answering**
- 5 Experiments
- 6 Conclusion
- 7 Strengths and Weaknesses
- 8 Questions?

Analytical Queries Answering

A Query

- Number of sites where each blogger posts, classified by the bloggers's age and city

```
SELECT ?Blogger ?Age ?City COUNT(*)  
WHERE {  
    ?Blogger a Blogger .  
    ?Blogger wrotePost ?BlogPost .  
    ?BlogPost postedOn ?Site .  
    ?Blogger age ?Age .  
    ?Blogger livesIn ?City .  
}  
GROUP BY ?Site
```

Analytical Queries Answering

A Query

- Number of sites where each blogger posts, classified by the bloggers's age and city

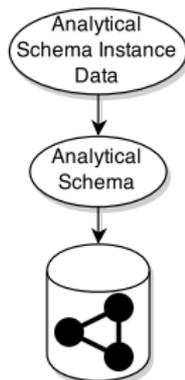
```
SELECT ?Blogger ?Age ?City COUNT(*)  
WHERE {  
    ?Blogger a Blogger .  
    ?Blogger wrotePost ?BlogPost .  
    ?BlogPost postedOn ?Site .  
    ?Blogger age ?Age .  
    ?Blogger livesIn ?City .  
}  
GROUP BY ?Site
```

- They use the language q not SPARQL

Analytical Queries Answering

Materialization and Query Reformulation

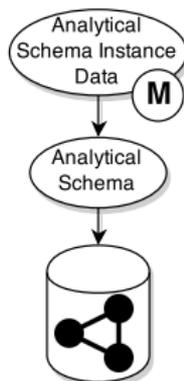
```
SELECT ?Blogger ?Age ?City COUNT(*)  
WHERE {  
  ?Blogger a Blogger .  
  ?Blogger wrotePost ?BlogPost .  
  ?BlogPost postedOn ?Site .  
  ?Blogger age ?Age .  
  ?Blogger livesIn ?City .  
}  
GROUP BY ?Site
```



Analytical Queries Answering

Materialization and Query Reformulation

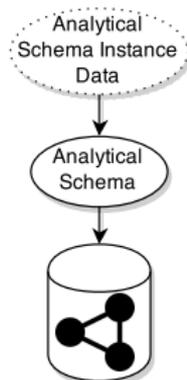
```
SELECT ?Blogger ?Age ?City COUNT(*)  
WHERE {  
  ?Blogger a Blogger .  
  ?Blogger wrotePost ?BlogPost .  
  ?BlogPost postedOn ?Site .  
  ?Blogger age ?Age .  
  ?Blogger livesIn ?City .  
}  
GROUP BY ?Site
```



Analytical Queries Answering

Materialization and Query Reformulation

```
SELECT ?Blogger ?Age ?City COUNT(*)  
WHERE {  
  ?Blogger a Blogger .  
  ?Blogger wrotePost ?BlogPost .  
  ?BlogPost postedOn ?Site .  
  ?Blogger age ?Age .  
  ?Blogger livesIn ?City .  
}  
GROUP BY ?Site
```



Analytical Queries Answering

Query Reformulation

```
SELECT ?Blogger ?Age ?City COUNT(*)  
WHERE {  
    ?Blogger a Blogger .  
    ?Blogger wrotePost ?BlogPost .  
    ?BlogPost postedOn ?Site .  
    ?Blogger age ?Age .  
    ?Blogger livesIn ?City .  
}  
GROUP BY ?Site
```

$\text{Blogger} = q(x):- x \text{ rdf:type Person, } x \text{ wrote } y, y \text{ inBlog } z$

$\text{BlogPost} = q(x):- x \text{ rdf:type Message, } x \text{ inBlog } z, z \text{ rdf:type Blog}$

$\text{wrotePost} = q(x,y):- x \text{ wrote } y, y \text{ rdf:type Message}$

Analytical Queries Answering

Materialization vs Query Reformulation

- Materialization
 - Fast query times
- Query Reformulation
 - Low Storage Cost
 - No maintenance at update

Outline

- 1 Motivation
- 2 Analytical Schema
- 3 Analytical Queries
- 4 Analytical Queries Answering
- 5 Experiments**
- 6 Conclusion
- 7 Strengths and Weaknesses
- 8 Questions?

- Database
 - kdb+ v3.0
 - Commercial software
 - Query language q
 - In-memory column store
 - Light weight
- Hardware
 - 8-core 2.13 GHz
 - 16 GB RAM
 - Linux Kernel 2.6.31.14
- Data
 - Ontology and Ontology Infobox from DBpedia
 - 34,000,000 triples (4.4 GB)
 - Scale by replicate data

Experiments

Materialization (1)

$\text{Blogger}(3) = q(x):-x \text{ rdf:type Person, } x \text{ wrote } y, y \text{ inBlog } z$

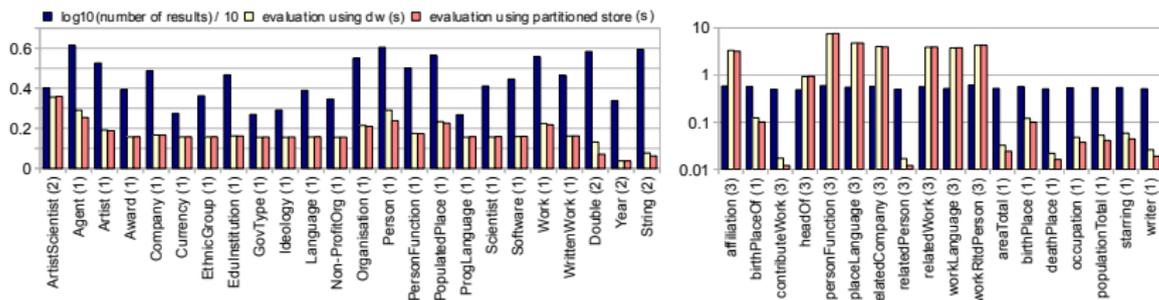


Figure 6: Evaluation time (s) and number of results for AnS node queries (left) and edge queries (right).

Blogger(3) = $q(x):-x \text{ rdf:type Person, } x \text{ wrote } y, y \text{ inBlog } z$

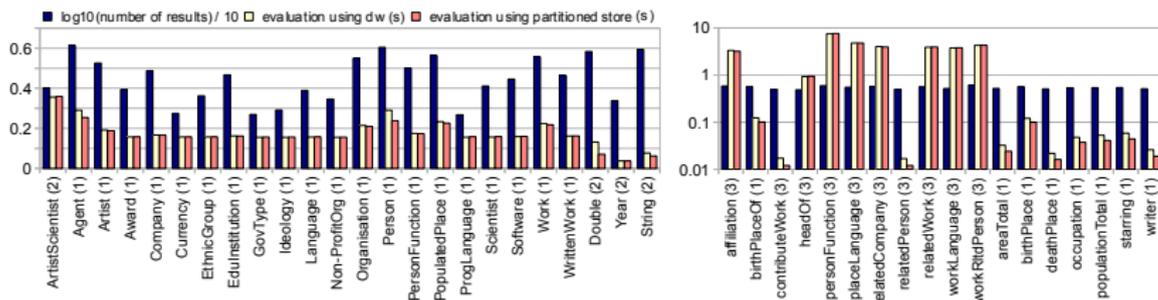


Figure 6: Evaluation time (s) and number of results for *AnS* node queries (left) and edge queries (right).

- 2 node queries and 57 edge queries are omitted because they are below 0.01 seconds

Experiments

Materialization (2)

- Instance table (Subject, Predicate, Object)
- Partitioned table Node(Subject) Edge(Subject, Object)

Experiments

Materialization (2)

- Instance table (Subject, Predicate, Object)
- Partitioned table Node(Subject) Edge(Subject, Object)

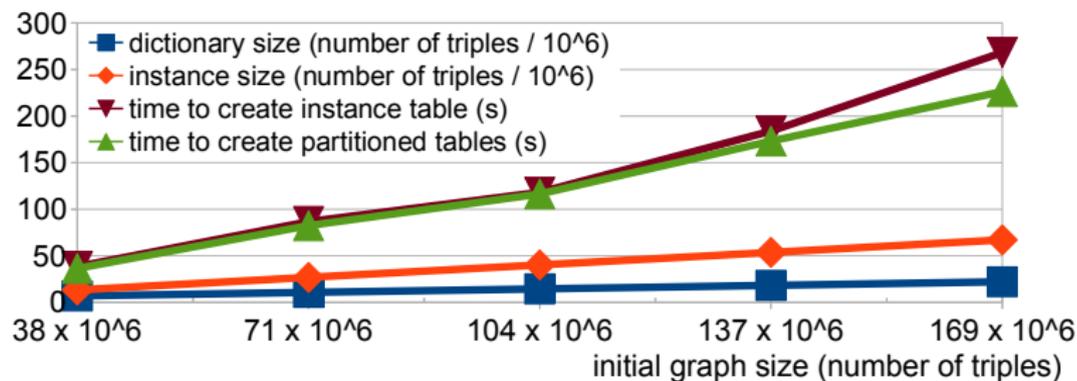


Figure 7: \mathcal{I} materialization time vs. \mathcal{I} size.

Experiments

Materialized Analytical Queries / Cubes

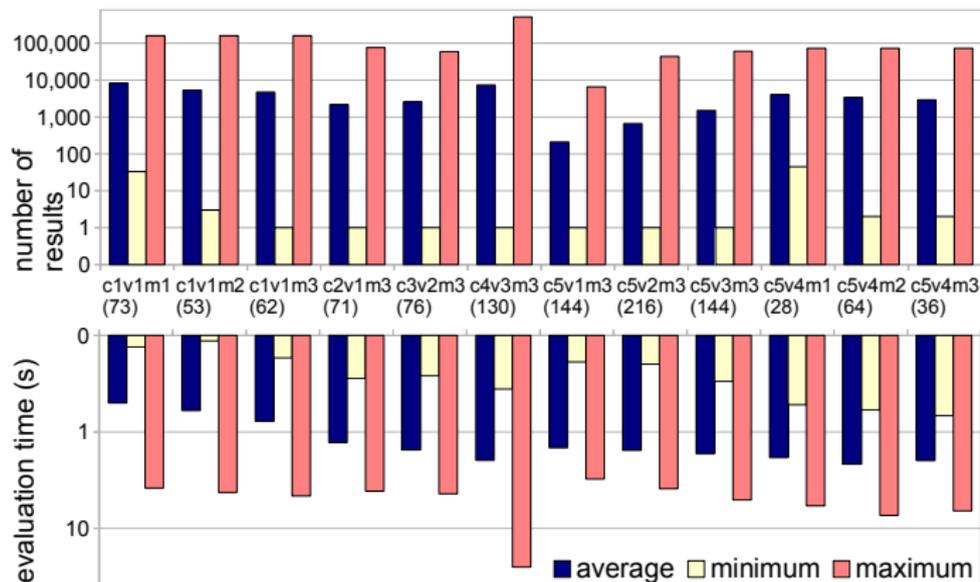


Figure 8: *AnQ* statistics for query patterns.

Experiments

Materialized Analytical Queries Scaling

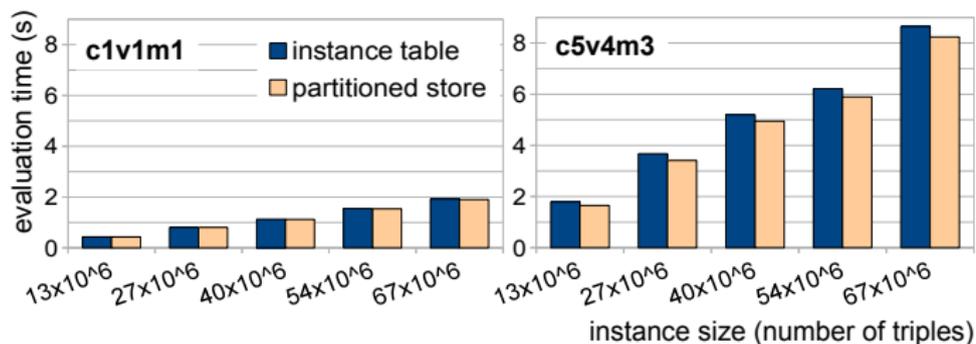


Figure 9: *AnQ* evaluation time over large datasets.

Outline

- 1 Motivation
- 2 Analytical Schema
- 3 Analytical Queries
- 4 Analytical Queries Answering
- 5 Experiments
- 6 Conclusion**
- 7 Strengths and Weaknesses
- 8 Questions?

Conclusion

- RDF data warehouse
- Dynamic Analytical Schema to model data
- Make Analytical Queries (cubes)
- Linear scale of Analytical Queries (cubes)

Outline

- 1 Motivation
- 2 Analytical Schema
- 3 Analytical Queries
- 4 Analytical Queries Answering
- 5 Experiments
- 6 Conclusion
- 7 Strengths and Weaknesses**
- 8 Questions?

- The story about Alice (Context) :)

Strengths

- The story about Alice (Context) :)
- Well written (Art of complicated writing)

Strengths

- The story about Alice (Context) :)
- Well written (Art of complicated writing)
- Relational concepts get new Semantic Web names
 - Analytical Schema = (Snowflake/star) Schema
 - Analytical Queries = Cube
 - OLAP (Specified)
 - Dimensions are define by a classifier query

Strengths

- The story about Alice (Context) :)
- Well written (Art of complicated writing)
- Relational concepts get new Semantic Web names
 - Analytical Schema = (Snowflake/star) Schema
 - Analytical Queries = Cube
 - OLAP (Specified)
 - Dimensions are define by a classifier query
- Precise formulations

Strengths

- The story about Alice (Context) :)
- Well written (Art of complicated writing)
- Relational concepts get new Semantic Web names
 - Analytical Schema = (Snowflake/star) Schema
 - Analytical Queries = Cube
 - OLAP (Specified)
 - Dimensions are define by a classifier query
- Precise formulations
- Running example throughout the paper

Strengths

- The story about Alice (Context) :)
- Well written (Art of complicated writing)
- Relational concepts get new Semantic Web names
 - Analytical Schema = (Snowflake/star) Schema
 - Analytical Queries = Cube
 - OLAP (Specified)
 - Dimensions are define by a classifier query
- Precise formulations
- Running example throughout the paper
- Two materializations strategies
 - Instance table
 - Partitioned tables

- The story about Alice (Context) :)
- Well written (Art of complicated writing)
- Relational concepts get new Semantic Web names
 - Analytical Schema = (Snowflake/star) Schema
 - Analytical Queries = Cube
 - OLAP (Specified)
 - Dimensions are define by a classifier query
- Precise formulations
- Running example throughout the paper
- Two materializations strategies
 - Instance table
 - Partitioned tables
- Basic OLAP support

- No query reformulation results, this is the hard part :(
 - In another paper (informal French conference)

- No query reformulation results, this is the hard part :(
 - In another paper (informal French conference)
- What kind of index do they use on their data?
 - PSOC?

- No query reformulation results, this is the hard part :(
 - In another paper (informal French conference)
- What kind of index do they use on their data?
 - PSOC?
- Analytical Queries (Cubes) can only have one type of aggregation

- No query reformulation results, this is the hard part :(
 - In another paper (informal French conference)
- What kind of index do they use on their data?
 - PSOC?
- Analytical Queries (Cubes) can only have one type of aggregation
- Simple cubes in experiments (not complex)
 - No “nextLevel” in the dataset
 - No aggregation
 - No unbalanced hierarchies

- No query reformulation results, this is the hard part :(
 - In another paper (informal French conference)
- What kind of index do they use on their data?
 - PSOC?
- Analytical Queries (Cubes) can only have one type of aggregation
- Simple cubes in experiments (not complex)
 - No “nextLevel” in the dataset
 - No aggregation
 - No unbalanced hierarchies
- No BI queries :S They only have “setup” times!
 - Time to get Analytical Schema Nodes and Edges
 - Creation time of Analytical Queries (cubes) and how this scales

Questions?