

# SPARQLytics: Multidimensional Analytics for RDF

Michael Rudolf

Database Technology Group, Technische Universität Dresden

March 8, 2017

# Agenda

MOTIVATION

RDF AND SPARQL

MULTIDIMENSIONAL ANALYTICS FOR RDF

# *Motivation*

## FOCUS MOVED FROM SINGLE ENTITY (OLTP)

- Bookkeeping
- Where is what?



## TO AGGREGATIONS OVER SETS OF ENTITIES OF THE SAME KIND (OLAP)

- Reporting
- What are the sales figures?



## TO CONNECTIONS BETWEEN ENTITIES

- Who likes what and why?
- What do the friends of your customers buy?



## SUPPLY CHAIN MANAGEMENT

- Transportation & logistics: routing, tendering, tracking, auditing, payment



<http://787updates.newairplane.com/787-Suppliers/World-Class-Supplier-Quality>

## SUPPLY CHAIN MANAGEMENT

- Transportation & logistics: routing, tendering, tracking, auditing, payment



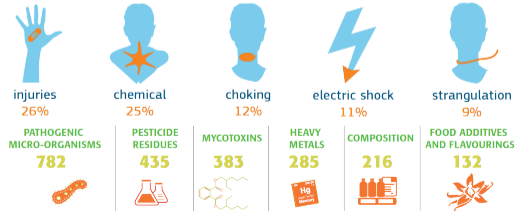
<http://787updates.newairplane.com/787-Suppliers/World-Class-Supplier-Quality>

## TRACK & TRACE

- Pinpoint product recalls
- Mandated by law for certain industries (e.g. pharmaceuticals, food, waste)

## EU Commission's Rapid Alert System

	non-food (RAPEX)	food & feed (RASFF)
2013	2364	3137
2014	2435	3157



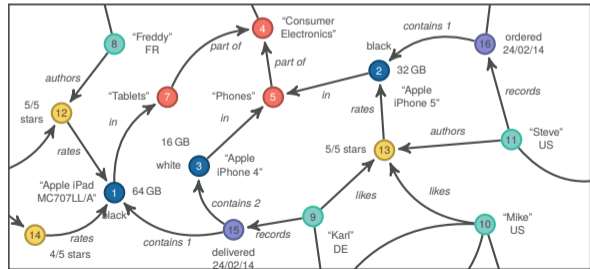
# *RDF and SPARQL*

- **Subjects** name an entity
- **Predicates** describe the relationship
- **Objects** can be literals or name

- no built-in schema
- can re-use vocabularies and ontologies
- suitable for inferencing facts

```
@prefix amazon: <http://www.amazon.com/#> .  
@prefix customer: <http://www.amazon.com/customer#> .  
@prefix product: <http://www.amazon.com/product#> .  
@prefix category: <http://www.amazon.com/category#> .
```

```
product:1 amazon:capacity "64 GB" .  
product:1 amazon:color "black" .  
product:1 amazon:in category:7 .  
category:7 amazon:name "Tablets" .  
category:7 amazon:partOf category:6 .  
category:6 amazon:name "Computers & Accessories" .  
user:8 amazon:country "FR" .  
user:8 amazon:rates product:1 .
```





- Built around pattern matching, produces pattern variable bindings
- Grouping and aggregation, CRUD operations
- No multidimensional concepts → complex and error-prone queries

```
PREFIX amazon: <http://www.amazon.com/#>
SELECT (AVG(?capacity) AS ?avgCap) (?name AS ?categoryName)
WHERE {
  ?product amazon:in      ?category .
  ?category amazon:name   ?name .
  ?category amazon:partOf+ category:6 .
  ?product amazon:capacity ?capacity
}
GROUP BY ?categoryName
```

# *Multidimensional Analytics for RDF*

## (BASE) FACTS

- Describe events and measurements
- Mostly numeric and continuous

## DIMENSIONS

- Provide context for facts
- If numeric, then often discrete
- Can embody structure

## MEASURES

- Are computed from grouped facts
- Are “arranged” in (hyper-)cubes

## (BASE) FACTS

- Describe events and measurements
- Mostly numeric and continuous



Slice

## DIMENSIONS

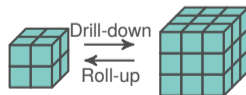
- Provide context for facts
- If numeric, then often discrete
- Can embody structure



Dice

## MEASURES

- Are computed from grouped facts
- Are “arranged” in (hyper-)cubes



## (BASE) FACTS

- Describe events and measurements
- Mostly numeric and continuous



Slice

## DIMENSIONS

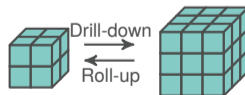
- Provide context for facts
- If numeric, then often discrete
- Can embody structure



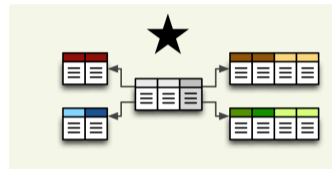
Dice

## MEASURES

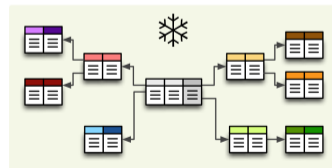
- Are computed from grouped facts
- Are “arranged” in (hyper-)cubes

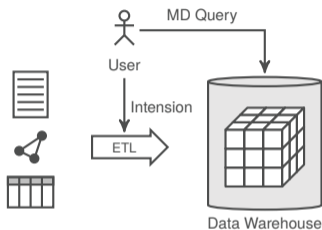


## Star schema



## Snowflake schema

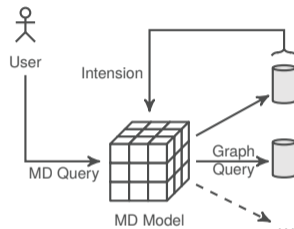
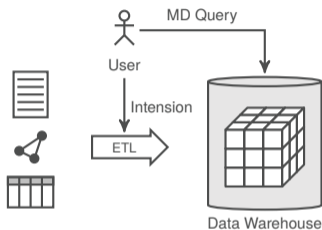




## DATA TRANSFORMATION

- Intension fixed by domain expert or metadata
- Import data using ETL process

# From Intensional to Extensional Analytics



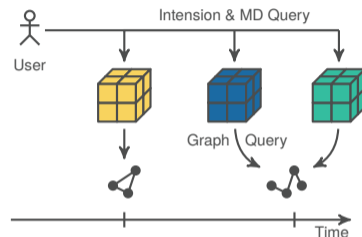
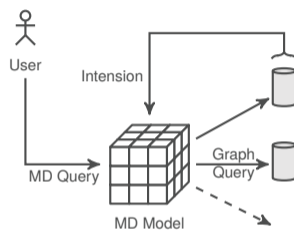
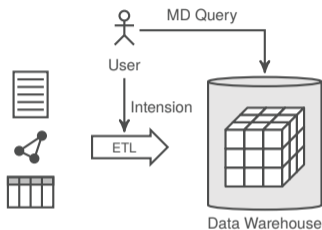
## DATA TRANSFORMATION

- Intension fixed by domain expert or metadata
- Import data using ETL process

## QUERY GENERATION

- Intension fixed by metadata
- Generate SPARQL queries from model

# From Intensional to Extensional Analytics



## DATA TRANSFORMATION

- Intension fixed by domain expert or metadata
- Import data using ETL process

## QUERY GENERATION

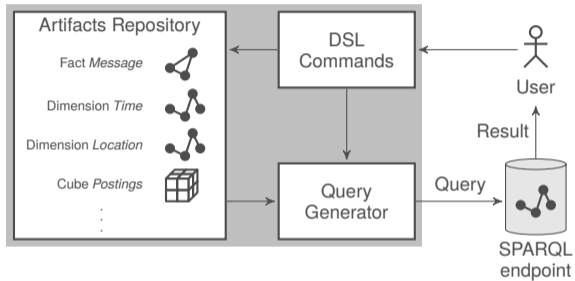
- Intension fixed by metadata
- Generate SPARQL queries from model

## EXTENSIONAL

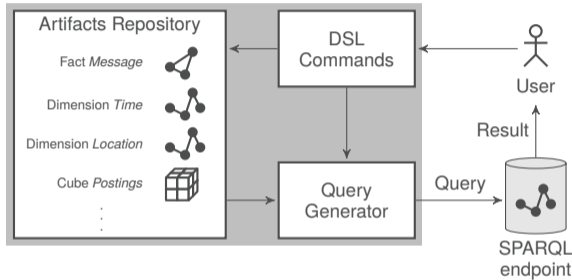
- Intension not fixed up-front
- Generate graph queries from user-specified intension



## SPARQLYTICS WORKFLOW



## SPARQLYTICS WORKFLOW

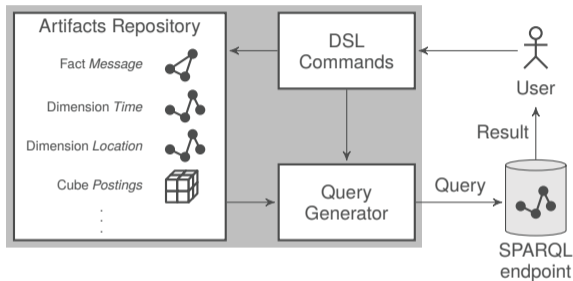


### 1. Create artifacts in repository

## EXAMPLE

```
USING REPOSITORY "myrepo";
SELECT FACTS {
  ?person rdf:type snvoc:Person ;
           snvoc:birthday ?birthday .
  FILTER (YEAR(NOW()) - YEAR(?birthday) >= 18)
};
DEFINE DIMENSION "Location" FROM (
  ?person snvoc:isLocatedIn ?city .
  ?city   snvoc:isPartOf ?country .
  ?country snvoc:isPartOf ?continent
) WITH (
  LEVEL "City" AS ?city,
  LEVEL "Country" AS ?country,
  LEVEL "Continent" AS ?continent
);
DEFINE MEASURE "Avg. No. Languages"
AS COUNT(DISTINCT ?language) WHERE (
  ?person snvoc:speaks ?language
) WITH "AVG";
CREATE CUBE "QB" FROM "Location", ...
WITH "Avg. No. Languages", ...;
```

## SPARQLYTICS WORKFLOW

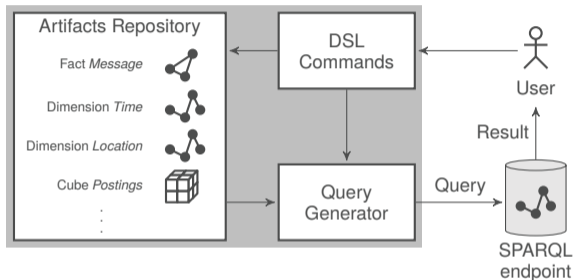


## EXAMPLE

```
USING CUBE "QB" OVER <http://localhost:3030/ds/sparql>;  
SLICE("Location", "Country", dbpedia:Italy);  
COMPUTE ("Avg. No. Languages");
```

1. Create artifacts in repository
2. Start session re-using artifacts

## SPARQLYTICS WORKFLOW



1. Create artifacts in repository
2. Start session re-using artifacts
3. Iteratively explore data, optionally create additional artifacts

## EXAMPLE

```
USING CUBE "QB" OVER <http://localhost:3030/ds/sparql>;  
SLICE("Location", "Country", dbpedia:Italy);  
COMPUTE ("Avg. No. Languages");
```

```
RESET FILTER("Location", "Country");  
ROLLUP("Location", 1);  
COMPUTE ("Avg. No. Languages");  
...
```

## BIG GRAPH DATA

- Not just social networks, also business scenarios
- Not enough data scientists, enable data enthusiasts

## RDF AND SPARQL

- Linked Open Data a rich source of information
- SPARQL does not expose multidimensional concepts

## SPARQLYTICS

- Re-use core SPARQL elements for defining multidimensional model
- Generate complex SPARQL queries from analytical session
- Stateful approach integrates well with data enthusiasts workflow

## *Additional Material & References*



Charu C. Aggarwal and Haixun Wang.

A Survey of Clustering Algorithms for Graph Data.

In Charu C. Aggarwal and Haixun Wang, editors, *Managing and Mining Graph Data*, volume 40 of *Advances in Database Systems*, chapter 9, pages 275–301. Springer US, 2010.



Seyed-Mehdi-Reza Beheshti, Boualem Benatallah, Hamid Reza Motahari-Nezhad, and Mohammad Allahbakhsh.

A framework and a language for on-line analytical processing on graphs.

In *Proceedings of the 13<sup>th</sup> International Conference on Web Information Systems Engineering (WISE)*, volume 7651 of *Lecture Notes in Computer Science*, pages 213–227. Springer, 2012.



Peter Boncz.

LDBC: Benchmarks for Graph and RDF Data Management.

In *Proc. IDEAS*, pages 1–2. ACM, 2013.



Fabio Crestani.

Application of spreading activation techniques in information retrieval.

*Artificial Intelligence Review*, 11(6):453–482, December 1997.



Chen Chen, Xifeng Yan, Feida Zhu, Jiawei Han, and Philip S. Yu.

Graph OLAP: Towards Online Analytical Processing on Graphs.

In *Proceedings of the 8<sup>th</sup> International Conference on Data Mining*, pages 103–112. IEEE, December 2008.



Hartmut Ehrig, Gregor Engels, Hans-Jörg Kreowski, and Grzegorz Rozenberg, editors.

*Handbook of Graph Grammars and Computing by Graph Transformation: Applications, Languages and Tools*, volume 2. World Scientific, 1997.



Steven Harris and Andy Seaborne.

SPARQL 1.1 query language.  
W3C recommendation, W3C, March 2013.



Dirk Koschützki, Katharina Anna Lehmann, Leon Peeters, Stefan Richter, Dagmar Tenfelde-Podehl, and Oliver Zlotowski.

*Centrality Indices*, volume 3418 of *Lecture Notes in Computer Science*, chapter 3, pages 16–61.  
Springer, 2005.



Sven Kosub.

*Local Density*, volume 3418 of *Lecture Notes in Computer Science*, chapter 6, pages 112–142.  
Springer, 2005.



Ralph Kimball and Margy Ross.

*The Data Warehouse Toolkit: The Definitive Guide to Dimensional Modeling*.  
Wiley, 3<sup>rd</sup> edition, 2013.



Kristen LeFevre and Evimaria Terzi.

Grass: Graph structure summarization.  
In *Proc. SDM*, pages 454–465. SIAM, 2010.



Saket Navlakha, Rajeev Rastogi, and Nisheeth Shrivastava.

Graph summarization with bounded error.  
In *Proc. SIGMOD*, pages 419–432. ACM, 2008.





Satu Elisa Schaeffer.

Graph clustering.

*Computer Science Review*, 1(1):27–64, August 2007.



Yuanyuan Tian and Jignesh M. Patel.

TALE: A Tool for Approximate Large Graph Matching.

In *2008 IEEE 24<sup>th</sup> International Conference on Data Engineering*, pages 963–972. IEEE, April 2008.



David Wood, Markus Lanthaler, and Richard Cyganiak.

RDF 1.1 concepts and abstract syntax.

W3C recommendation, W3C, February 2014.

<http://www.w3.org/TR/2014/REC-rdf11-concepts-20140225/>.



Peixiang Zhao, Xiaolei Li, Dong Xin, and Jiawei Han.

Graph Cube: On Warehousing and OLAP Multidimensional Networks.

In *Proceedings of the ACM SIGMOD International Conference on Management of Data*, pages 853–864. ACM, 2011.



Ning Zhang, Yuanyuan Tian, and Jignesh M. Patel.

Discovery-Driven Graph Summarization.

In *Proceedings of the 26<sup>th</sup> International Conference on Data Engineering*, pages 880–891. IEEE, 2010.