

Graph-based Virtual Data Integration

ANNA QUERALT, SERGI NADAL, OSCAR ROMERO
(FACULTAT D'INFORMÀTICA DE BARCELONA)

Data Integration

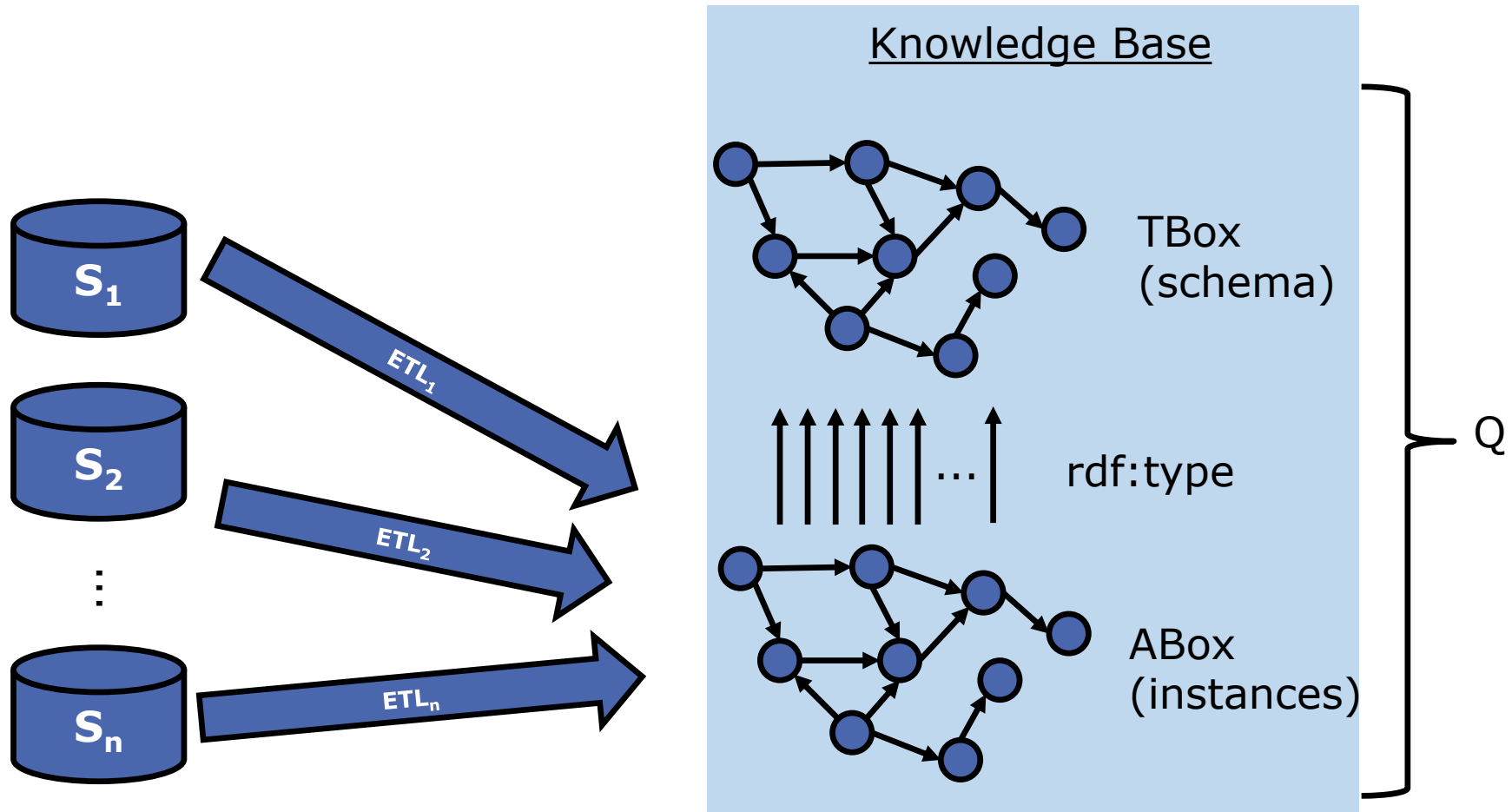
Data integration is an area of study within data management aimed at facilitating **transparent access** to a **variety of heterogeneous data sources**

- Two main options to perform data integration:
 - Physical data integration
 - Virtual data integration

In this course we will focus on using graphs to solve data integration

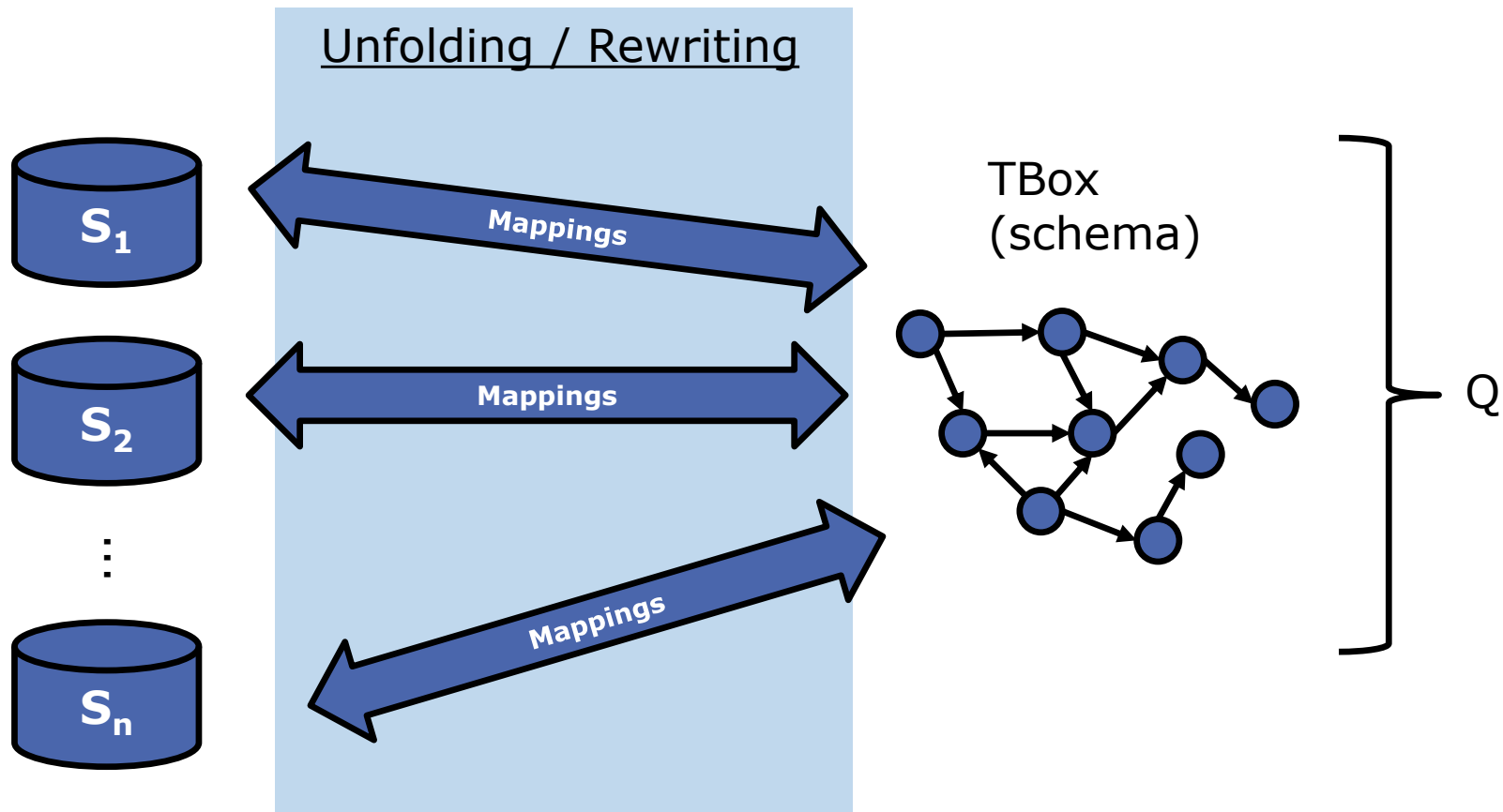
Graph-Based Data Integration at a Glance

Option 1: Physical Data Integration



Graph-Based Data Integration at a Glance

Option 2: Virtual Data Integration



Why Graph-Based Virtual Data Integration?

- ❑ When the data sources are not under your control and owners require a federation (e.g., data exchange between companies)
 - E.g., Data Portability
- ❑ When we do not want to move the data from where it resides
 - For example, key-based models are more performant than graph models for table scans
- ❑ When data freshness is crucial
 - ETLs run from time to time and the period between updates is called the *update window*
- ❑ To systematically trace all the data available within a company or organization (**data provenance**)
 - The construct created to perform data governance is typically called **data catalog** or **data fabric**

Why Graph-Based Virtual Data Integration?

- There is consensus, nowadays, that graph-based solutions are the way to go for data integration / data governance
 - Graph-based data models are richer than any other data model and can express any construct
 - Knowledge graphs are preferred since they facilitate linking TBOXes / ABOXes with little effort
 - Allow to represent all data integration constructs with the same formalism (target schema, source schemata and mappings)
- Most popular Big Data Integration systems follow this approach:
 - Data Tamer (<https://www.tamr.com/>)
 - The BigDAWG Polystore System (<https://dl.acm.org/citation.cfm?id=3226620>)
 - Ontop (<http://ontop.inf.unibz.it/>)
 - ODIN (<http://www.essi.upc.edu/~snadal/odin.html>)
 - Colibra (<https://www.colibra.com/us/en/platform/data-governance>)
 - Ataccama (<https://www.ataccama.com/>)
 - Company specific tools:
 - Amundsen (Lyft), DataHub (LinkedIn), Data Access Layer (Twitter), DataPortal (Airbnb), Databook (Uber), Metacat (Netflix), Lexikon (Spotify), Artifact (Shopify), Nemo (Meta), Atlas (Apache), Marquez (WeWork)... any Big Data Cloud provider also have their own

Two Main Approaches

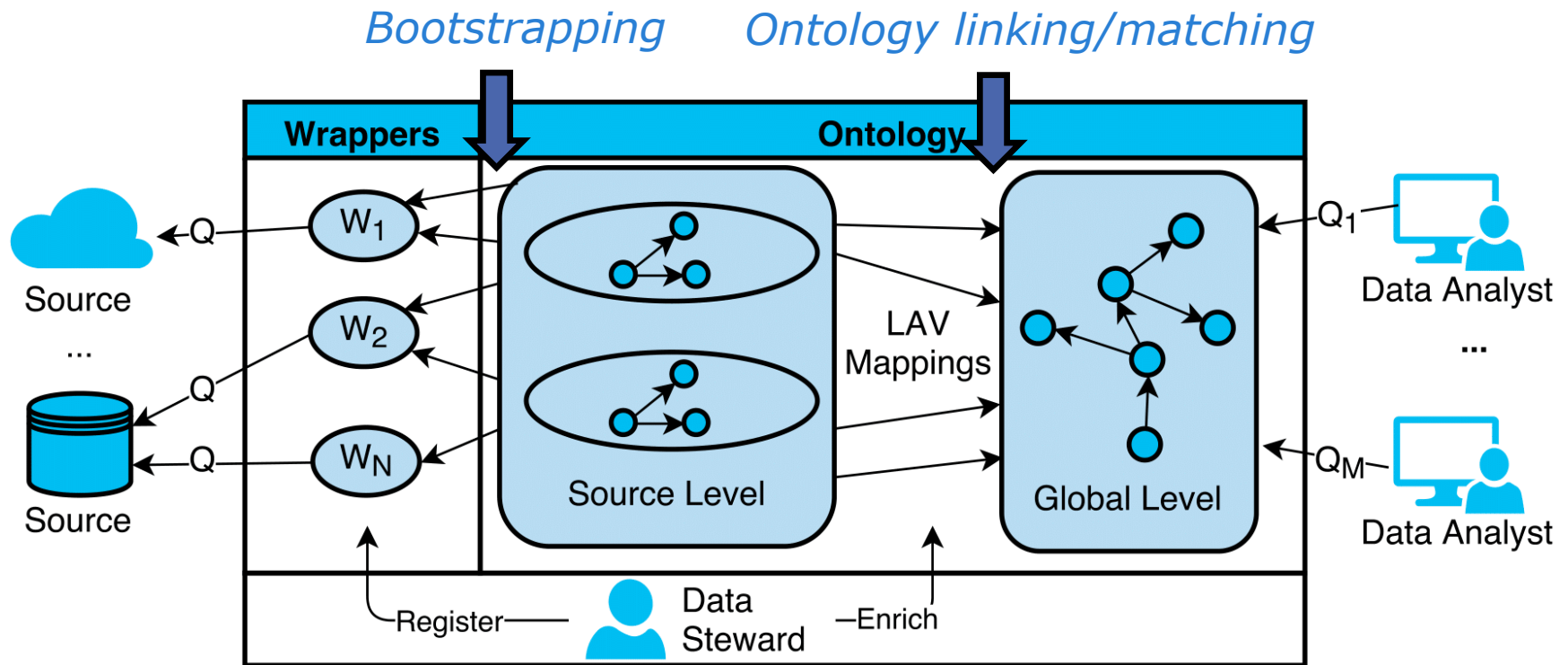
- Ontology-based data access
 - Monolithic approach
 - The TBox is directly related to the sources via mappings
- Ontology-mediated queries
 - Relies on the concept of wrapper
 - Thus, we can select a subset of the data source to be exposed to the whole integration System
 - Security
 - Modularity
 - It allows pay-as-you-go data Integration
 - The integrated schema is built incrementally as new data sources arrive

Local-as-view (LAV) Data Integration

ONTOLOGY-MEDIATED QUERIES

- It is a family of systems performing graph-based data integration with LAV
 - Conceptually, Global-as-view (GAV) is also possible
- **Based on the well-known wrapper-mediator architecture**
- To make the querying rewriting feasible, they adopt several measures:
 - Exact mappings (i.e., Closed-World assumption)
 - Very basic reasoning capabilities (taxonomies and domain / range inference)

Ontology-mediated Query



Virtual integration with LAV mappings

<https://www.semantic-web-journal.net/content/incremental-schema-integration-data-wrangling-knowledge-graphs-0>

Big Data Integration Ontology

- It revisits the Data Integration framework and construct an ontology as follows:
 - Global level (G) – Integrated view
 - Source levels (S) – Views on the data sources (wrappers)
 - Mappings (M) – LAV mappings between G and S

Wrappers

- They represent **a view** on the source
- You can think of a **named** query over the source. For example:

W1: SELECT a, b, c FROM T

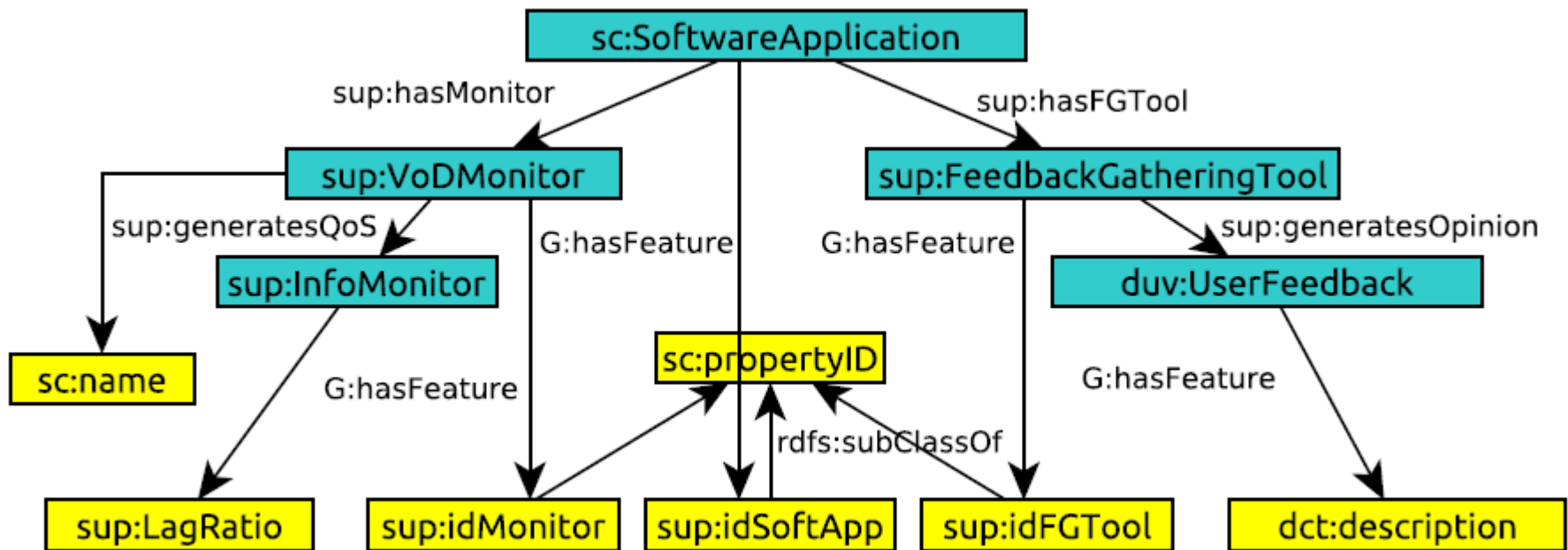
- Typical assumptions made by wrappers:
 - They expose the source in **tabular format** (1NF)
 - Thus, Cypher, SPARQL or MongoDB's aggregation framework would also meet the requirements
 - In general, most query languages produce tabular format
 - A data source may generate several wrappers
 - Typically, new versions of data are considered new wrappers

Example

- A company monitoring the status of their service providing streaming during the Olympics.
- They require to pose cross-domain queries on:
 - Monitored data on video players (e.g., lag ratio)
 - Tweets in English gathered through a feedback gathering tool

Global Level

- Green: concepts
- Yellow: attributes



Source Level

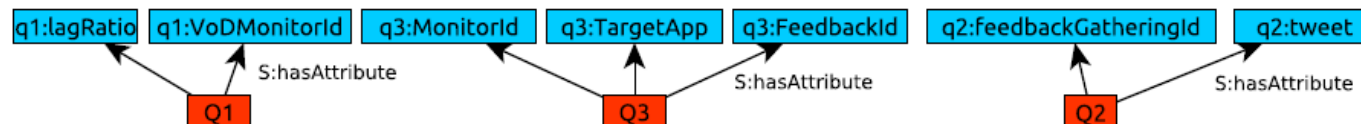
- Sources are exposed by means of wrappers
 - We automatically bootstrap the attributes projected by the wrappers

Q_1 : ID and compute the lag ratio

```
db.getCollection('vod').aggregate([  
  {$project: {"VoDmonitorId":true, "lagRatio": {$divide : ["$waitTime","$watchTime"]}}} ])
```

Q_2 - all attributes for tweets in english.

Q_3 - association target app \rightarrow monitor, feedback gathering tool

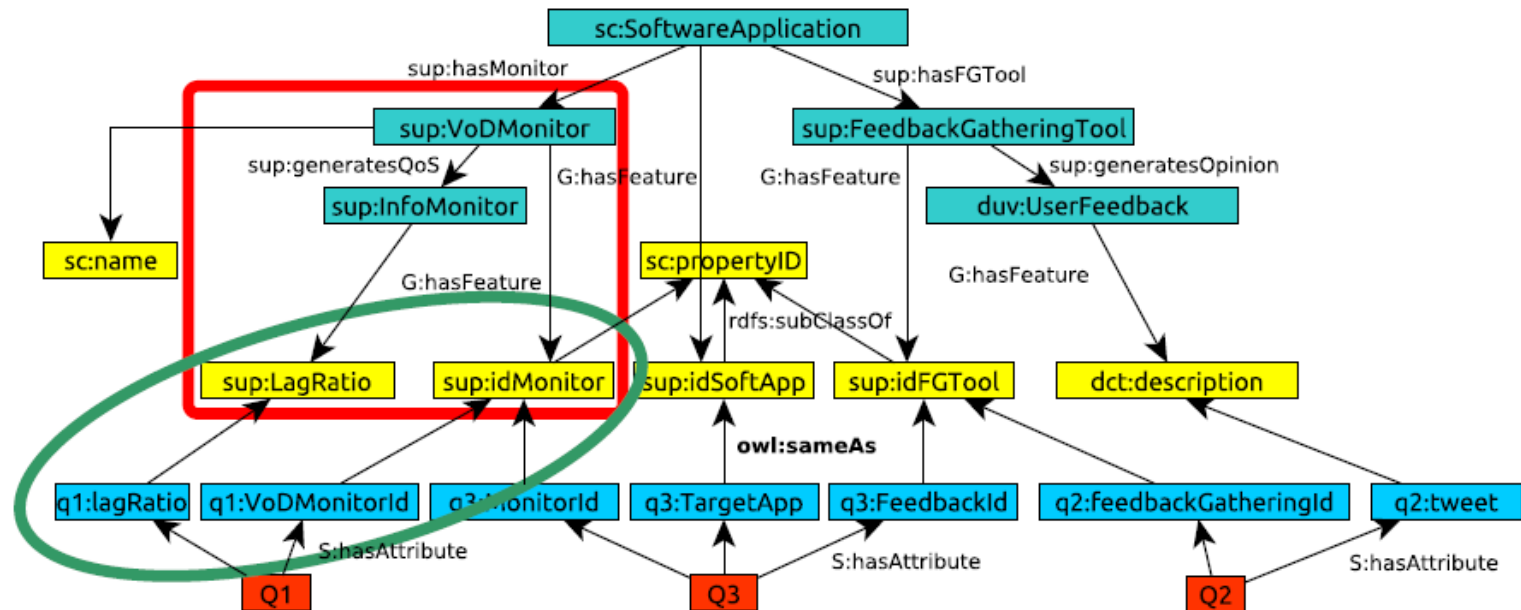


Red: Wrappers; Blue: Wrapper attributes

Mappings

- A Local-as-View (LAV) mapping for a wrapper Q is defined as: $M = \langle G, S \rangle$ where:
 - G is a named graph
 - S is a set of triples of the form:
 - $\langle x, \text{owl:sameAs}, y \rangle$, where
 - $\langle x, \text{rdf:type}, S:\text{Attribute} \rangle$ and
 - $\langle y, \text{rdf:type}, G:\text{Feature} \rangle$

LAV Mapping Example



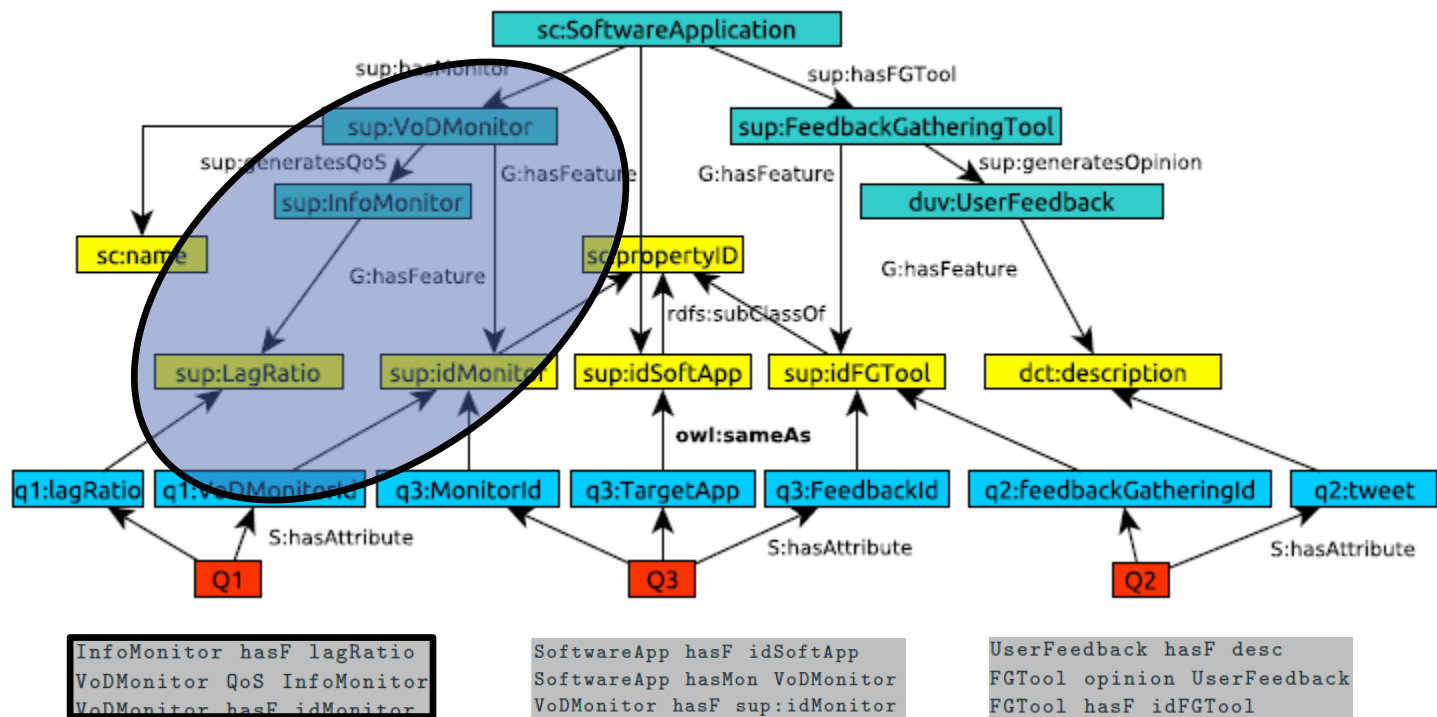
G (named graph):

```
Q1 S:provides { sup:InfoMonitor G:hasFeature sup:lagRatio .
sup:VoDMonitor sup:generatesQoS sup:InfoMonitor .
sup:VoDMonitor G:hasFeature sup:idMonitor }}
```

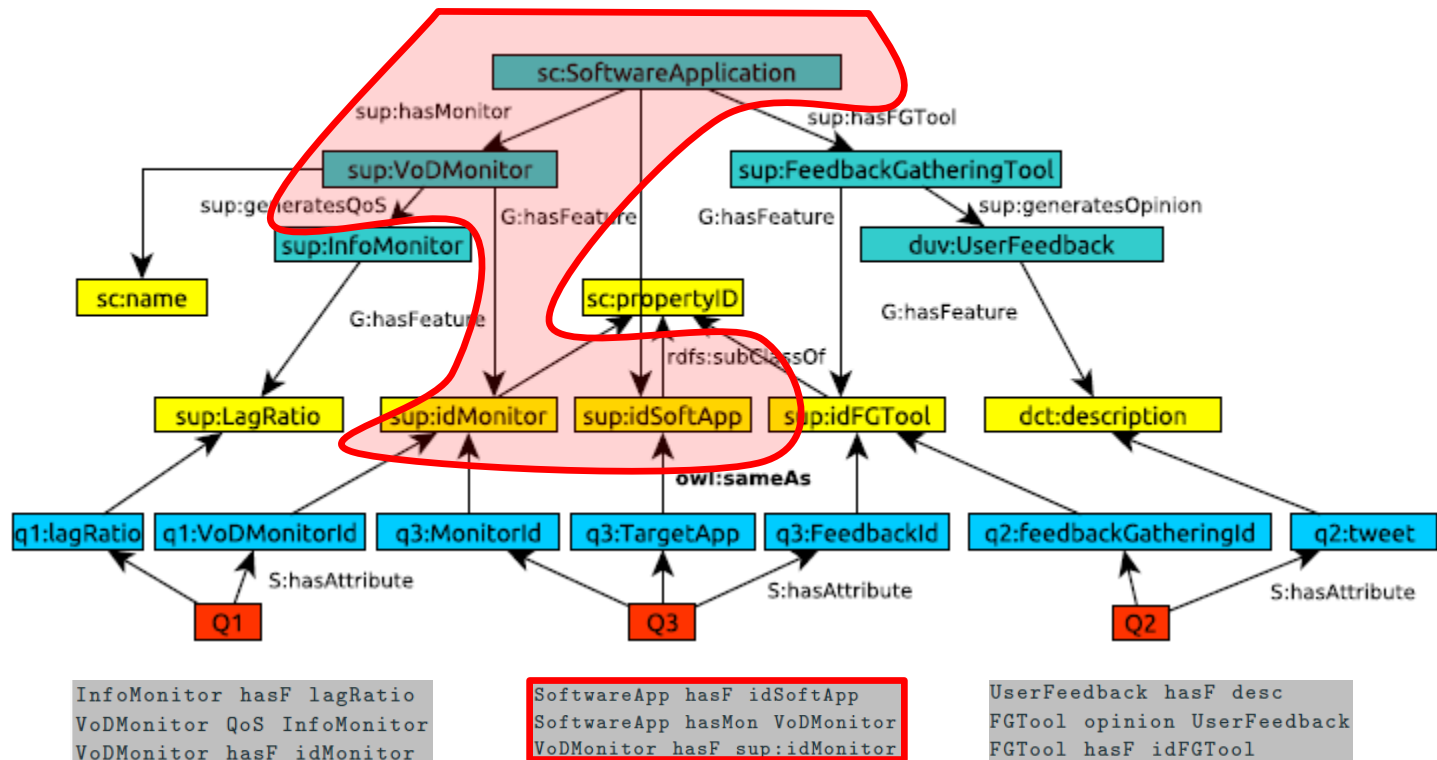
S ("same as" triples):

```
q1:lagRatio owl:sameAs sup:lagRatio
q1:VoDMonitorId owl:sameAs sup:idMonitor
```

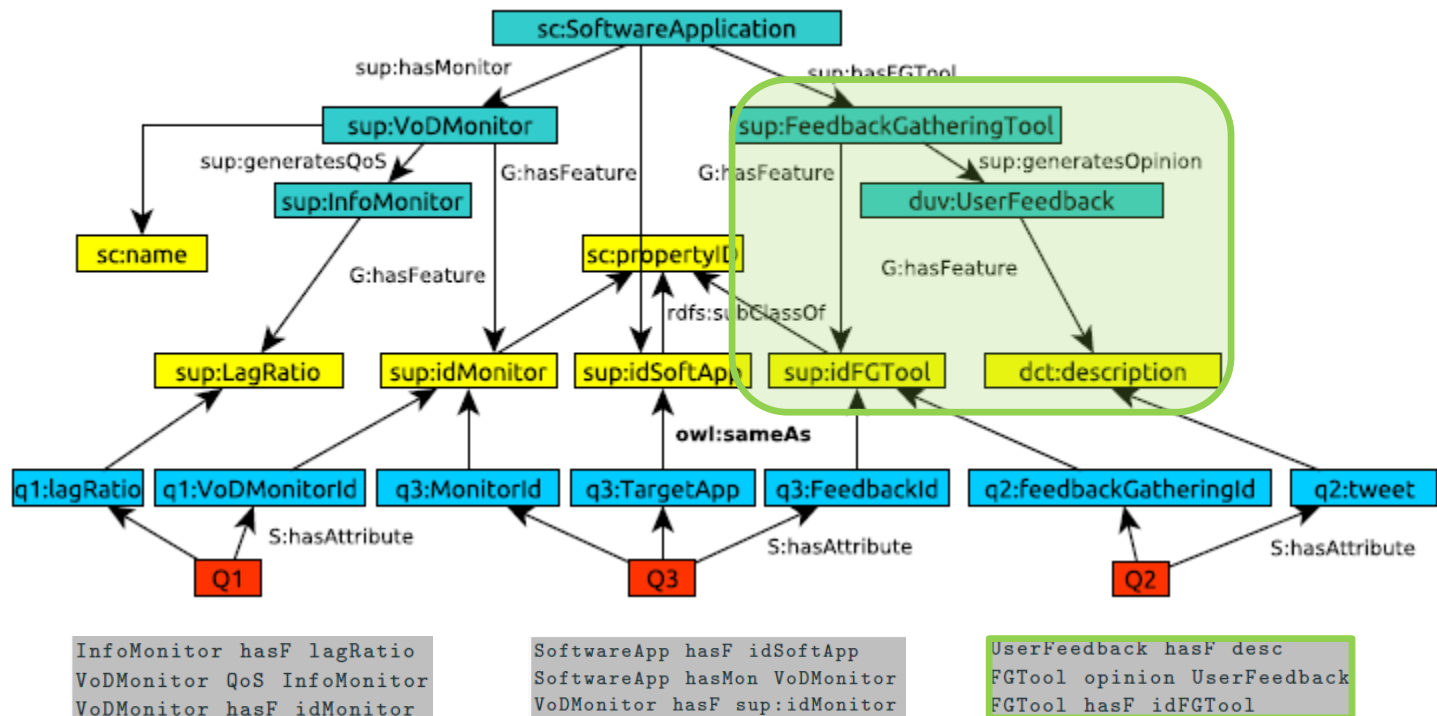
LAV Mappings (Q1)



LAV Mappings (Q3)



LAV Mappings (Q2)



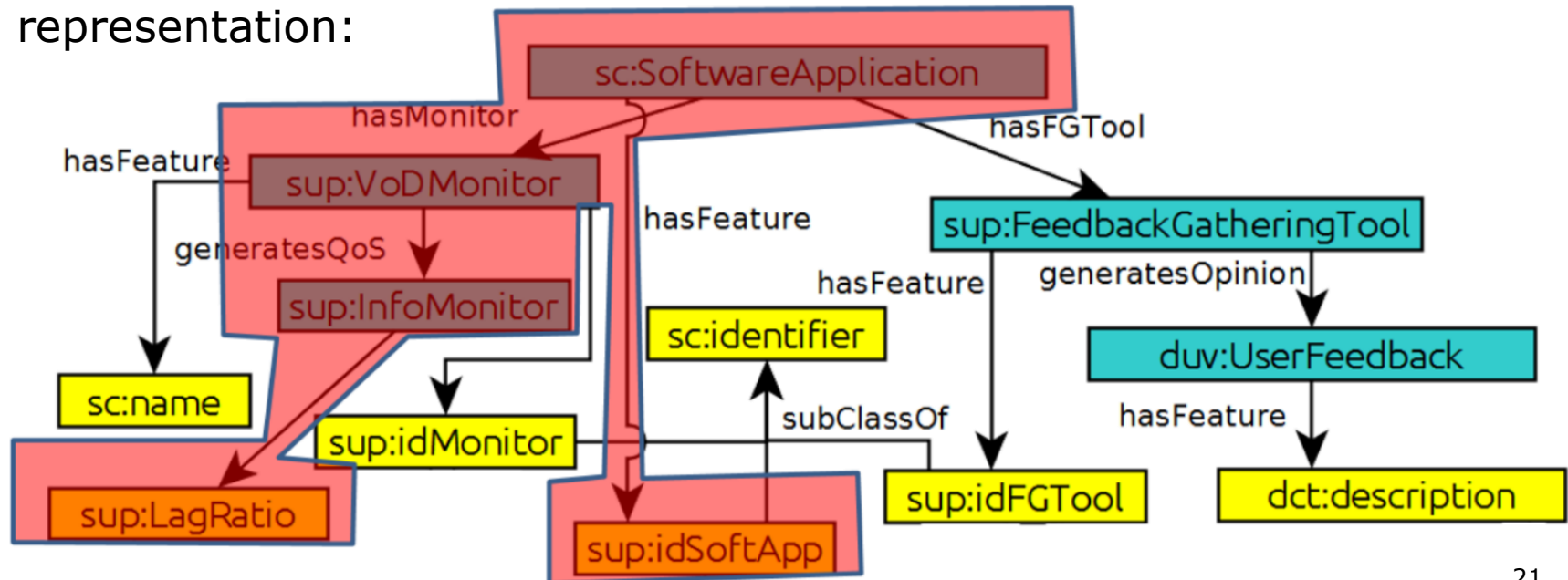
Query Answering – Rewriting Algorithm

- Any SPARQL query on the global graph must be rewritten as a query in terms of the wrappers
- Example of query over G :

SPARQL Query:

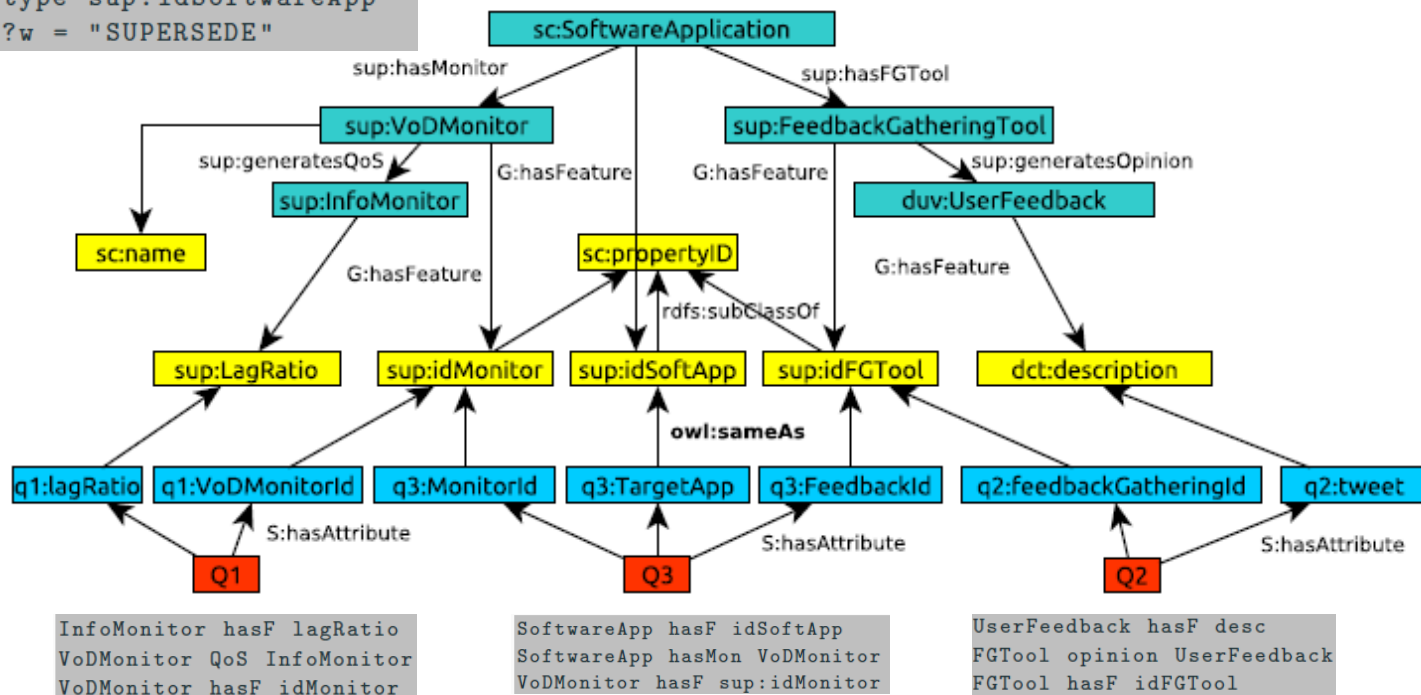
```
SELECT ?w,?t WHERE
  ?t rdf:type sup:lagRatio
  ?x G:hasFeature ?t
  ?x rdf:type sup:InfoMonitor
  ?y sup:generatesQoS ?x
  ?y rdf:type sup:VoDMonitor
  ?z sup:hasMonitor ?y
  ?z rdf:type sc:SoftwareApp
  ?z G:hasFeature ?w
  ?w rdf:type sup:idSoftwareApp
  FILTER ?w = "SUPERSEDE"
```

Graph representation:



Notions on the Query Rewriting Alg.

```
SELECT ?w,?t WHERE
  ?t rdf:type sup:lagRatio
  ?x G:hasFeature ?t
  ?x rdf:type sup:InfoMonitor
  ?y sup:generatesQoS ?x
  ?y rdf:type sup:VoDMonitor
  ?z sup:hasMonitor ?y
  ?z rdf:type sc:SoftwareApp
  ?z G:hasFeature ?w
  ?w rdf:type sup:idSoftwareApp
  FILTER ?w = "SUPERSEDE"
```



Start from a Terminal Feature

```
SELECT ?w,?t WHERE
```

```
?t rdf:type sup:lagRatio
```

```
?x G:hasFeature ?t
```

```
?x rdf:type sup:InfoMonitor
```

```
?y sup:generatesQoS ?x
```

```
?y rdf:type sup:VoDMonitor
```

```
?z sup:hasMonitor ?y
```

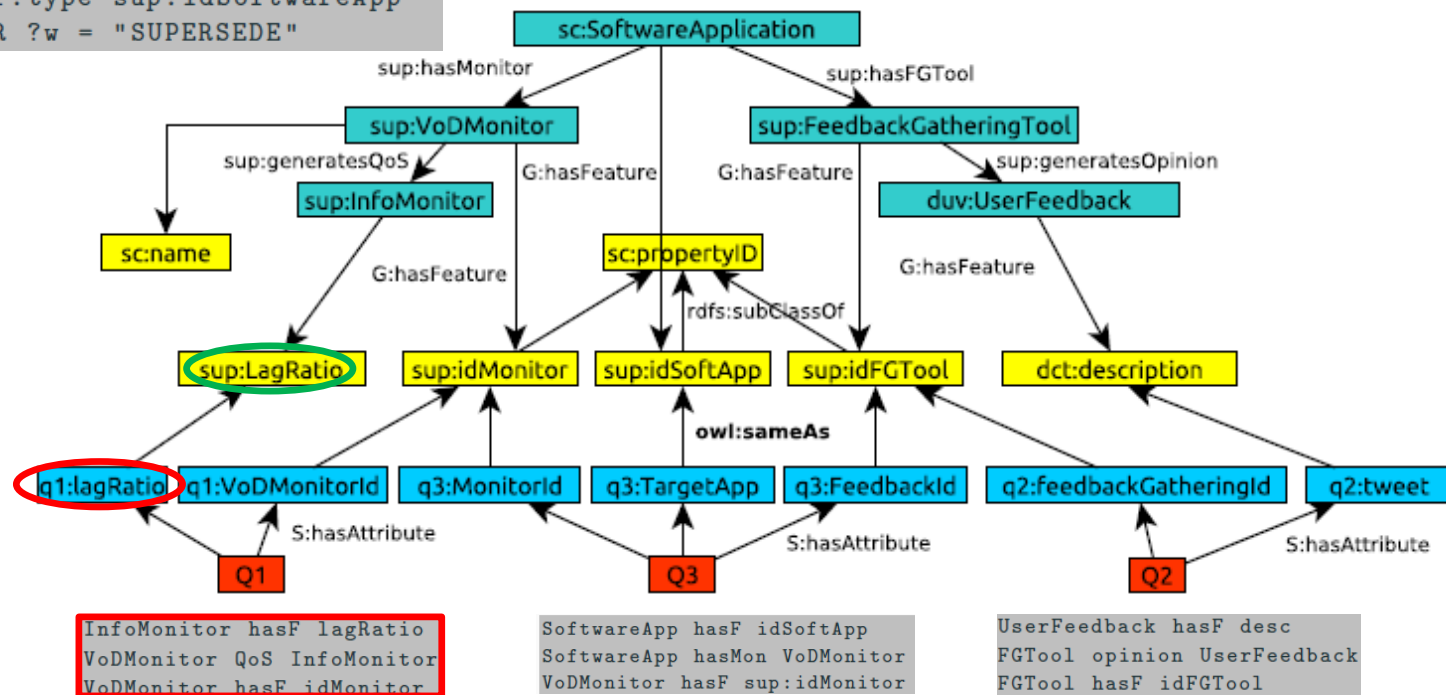
```
?z rdf:type sc:SoftwareApp
```

```
?z G:hasFeature ?w
```

```
?w rdf:type sup:idSoftwareApp
```

```
FILTER ?w = "SUPERSEDE"
```

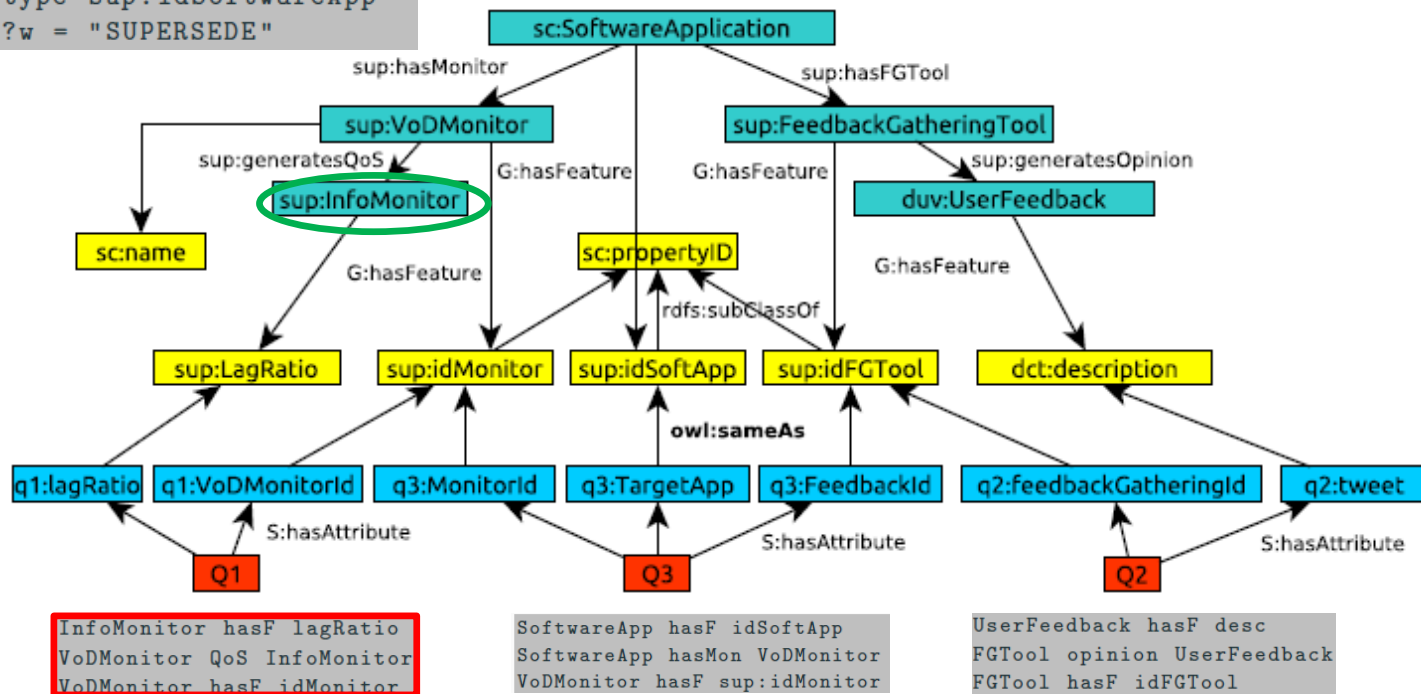
$$\sqcap t(\rho_{Q_1.lagRatio} \rightarrow t(Q_1))$$



Navigate G from the Feature

```
SELECT ?w,?t WHERE
  ?t rdf:type sup:lagRatio
  ?x G:hasFeature ?t
  ?x rdf:type sup:InfoMonitor
  ?y sup:generatesQoS ?x
  ?y rdf:type sup:VoDMonitor
  ?z sup:hasMonitor ?y
  ?z rdf:type sc:SoftwareApp
  ?z G:hasFeature ?w
  ?w rdf:type sup:idSoftwareApp
  FILTER ?w = "SUPERSEDE"
```

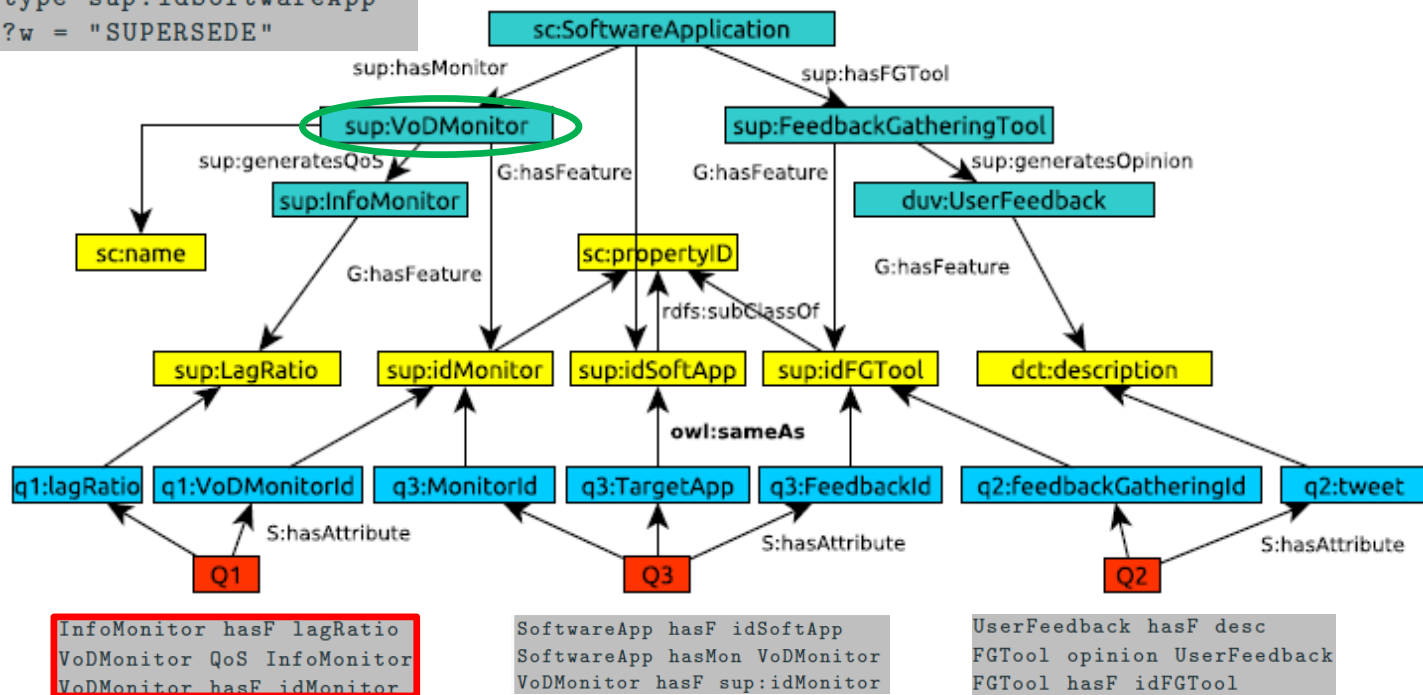
$$\sqcap t(\rho_{Q_1.lagRatio} \rightarrow t(Q_1))$$



Navigate G from the Feature

```
SELECT ?w,?t WHERE
  ?t rdf:type sup:lagRatio
  ?x G:hasFeature ?t
  ?x rdf:type sup:InfoMonitor
  ?y sup:generatesQoS ?x
  ?y rdf:type sup:VoDMonitor
  ?z sup:hasMonitor ?y
  ?z rdf:type sc:SoftwareApp
  ?z G:hasFeature ?w
  ?w rdf:type sup:idSoftwareApp
FILTER ?w = "SUPERSEDE"
```

$$\Pi_{t(\rho_{Q_1.lagRatio} \rightarrow t(Q_1))}$$



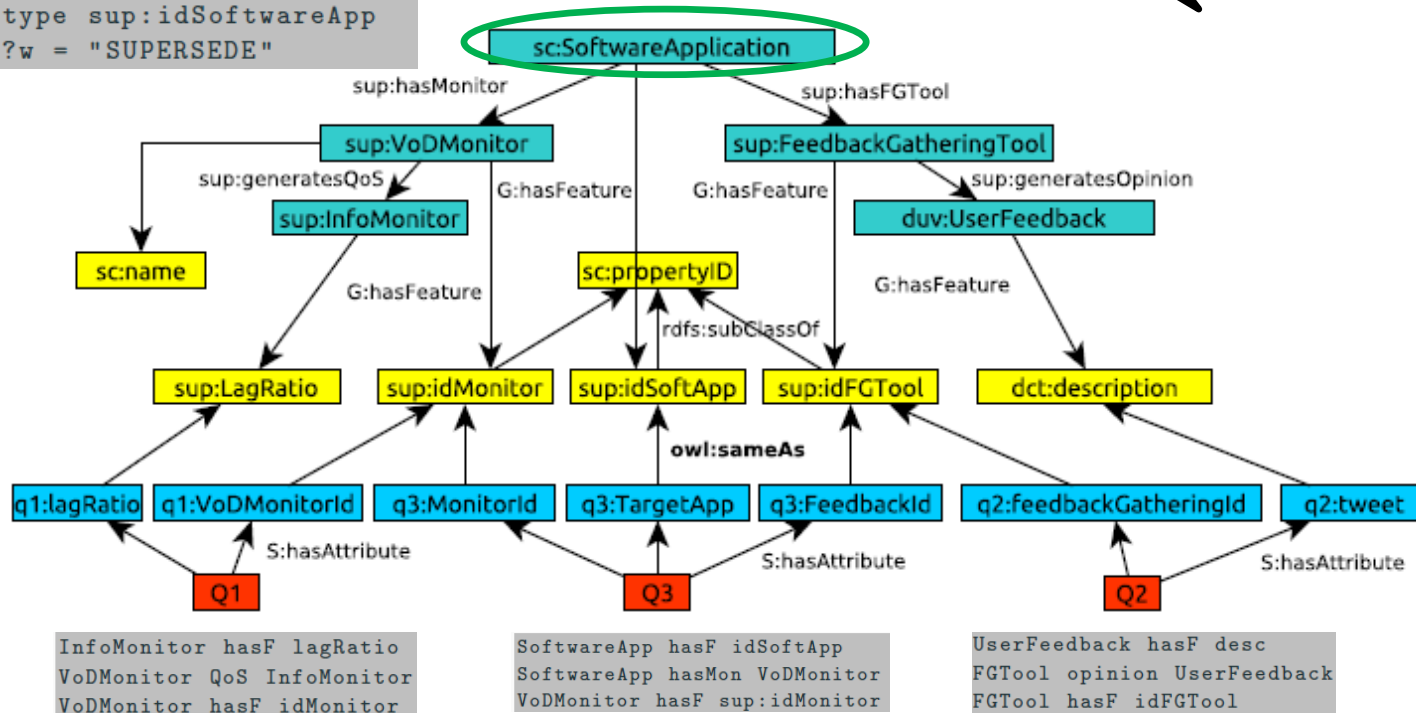
Navigate G from the Feature

```
SELECT ?w,?t WHERE
  ?t rdf:type sup:lagRatio
  ?x G:hasFeature ?t
  ?x rdf:type sup:InfoMonitor
  ?y sup:generatesQoS ?x
  ?y rdf:type sup:VoDMonitor
  ?z sup:hasMonitor ?y
  ?z rdf:type sc:SoftwareApp
  ?z G:hasFeature ?w
  ?w rdf:type sup:idSoftwareApp
  FILTER ?w = "SUPERSEDE"
```

$$\sqcap \quad t(\rho_{Q_1.lagRatio} \rightarrow t(Q_1))$$



Subpath not
contained in Q1!!



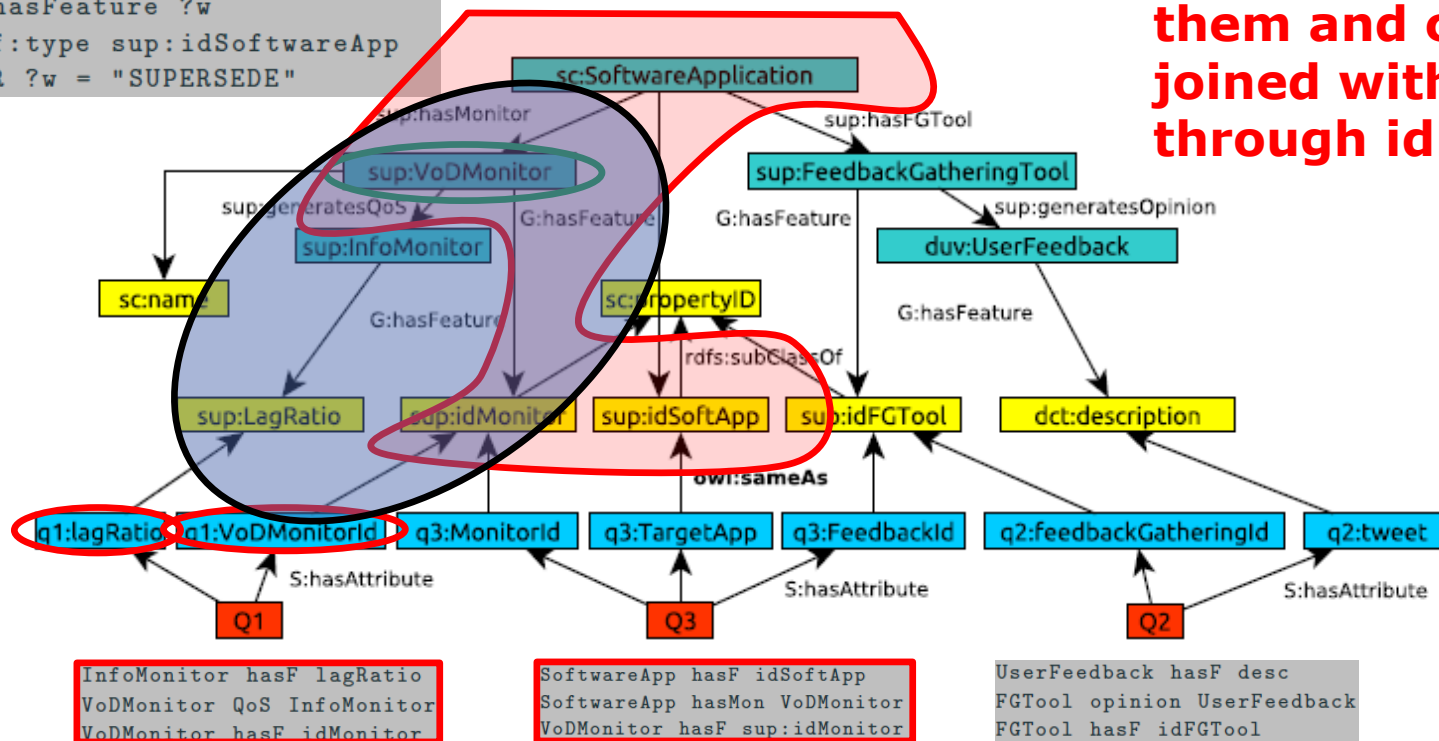
Explore Join Candidates

```
SELECT ?w,?t WHERE
  ?t rdf:type sup:lagRatio
  ?x G:hasFeature ?t
  ?x rdf:type sup:InfoMonitor
  ?y sup:generatesQoS ?x
  ?y rdf:type sup:VoDMonitor
  ?z sup:hasMonitor ?y
  ?z rdf:type sc:SoftwareApp
  ?z G:hasFeature ?w
  ?w rdf:type sup:idSoftwareApp
  FILTER ?w = "SUPERSEDE"
```

$$\Pi_{t(\rho_{Q_1}.lagRatio \rightarrow t(Q_1))}$$

Any other wrapper contains these triples?

Yes! Q3 covers them and can be joined with Q1 through idMonitor!

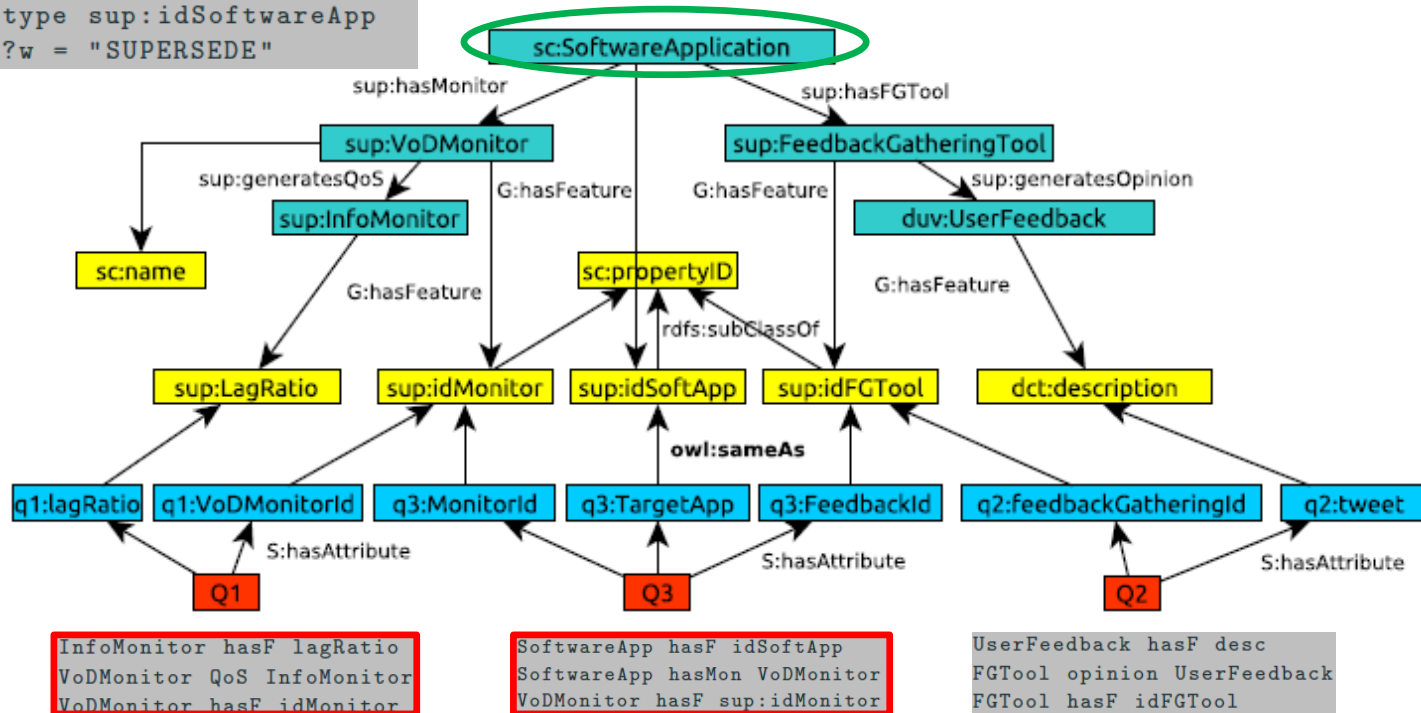


Join to an Alternative Wrapper

```
SELECT ?w,?t WHERE
  ?t rdf:type sup:lagRatio
  ?x G:hasFeature ?t
  ?x rdf:type sup:InfoMonitor
  ?y sup:generatesQoS ?x
  ?y rdf:type sup:VoDMonitor
  ?z sup:hasMonitor ?y
  ?z rdf:type sc:SoftwareApp
  ?z G:hasFeature ?w
  ?w rdf:type sup:idSoftwareApp
  FILTER ?w = "SUPERSEDE"
```

$$\Pi \quad t(\rho_{Q_1}.lagRatio \rightarrow t$$

$$\sigma_{Q_1.VoDMonitorId = Q_3.MonitorId}$$

$$(Q_1 \times Q_3))$$


Continue Navigating G

```
SELECT ?w,?t WHERE
  ?t rdf:type sup:lagRatio
  ?x G:hasFeature ?t
  ?x rdf:type sup:InfoMonitor
  ?y sup:generatesQoS ?x
  ?y rdf:type sup:VoDMonitor
  ?z sup:hasMonitor ?y
  ?z rdf:type sc:SoftwareApp
  ?z G:hasFeature ?w
  ?w rdf:type sup:idSoftwareApp
  FILTER ?w = "SUPERSEDE"
```

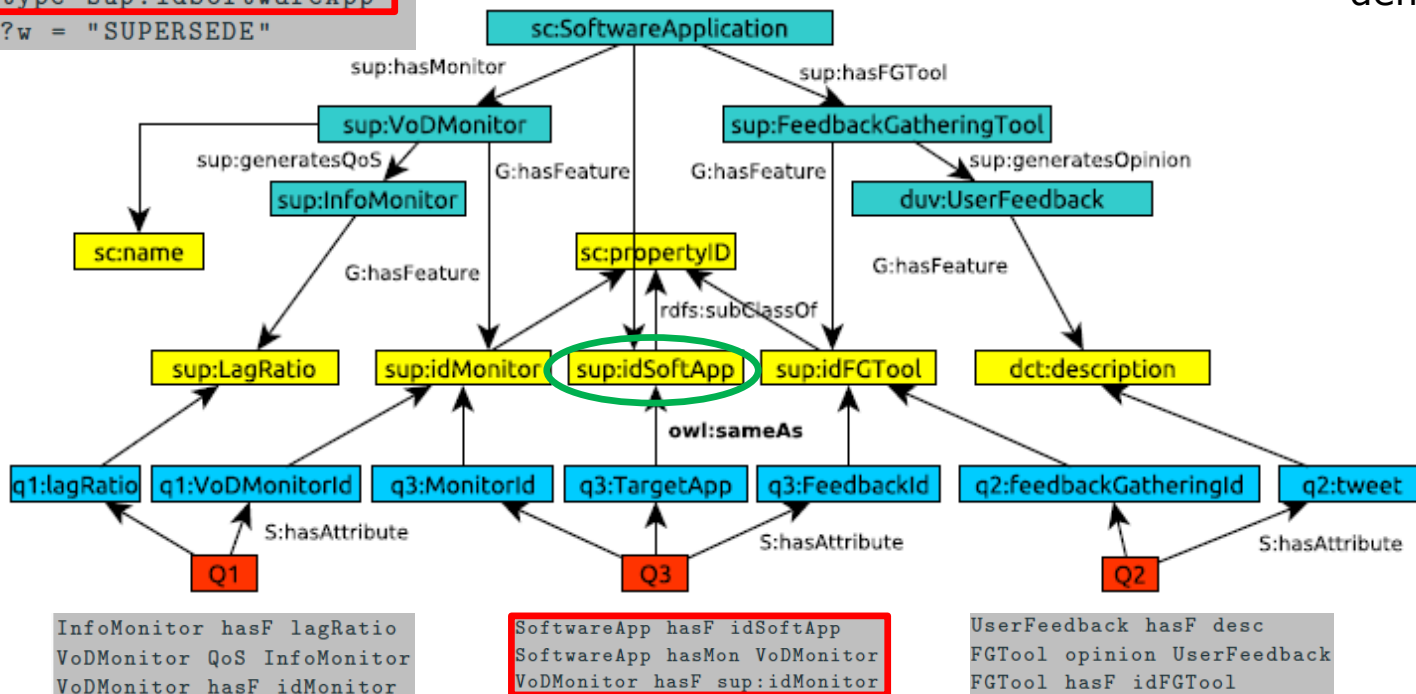
$$\Pi_{w,t}(\rho_{Q_1.lagRatio \rightarrow t}$$

$$\rho_{Q_3.TargetApp \rightarrow w}$$

$$\sigma_{Q_1.VoDMonitorId=Q_3.MonitorId}$$

$$(Q_1 \times Q_3))$$

This is a query over the wrappers! Now, each wrapper name must be replaced by its definition query



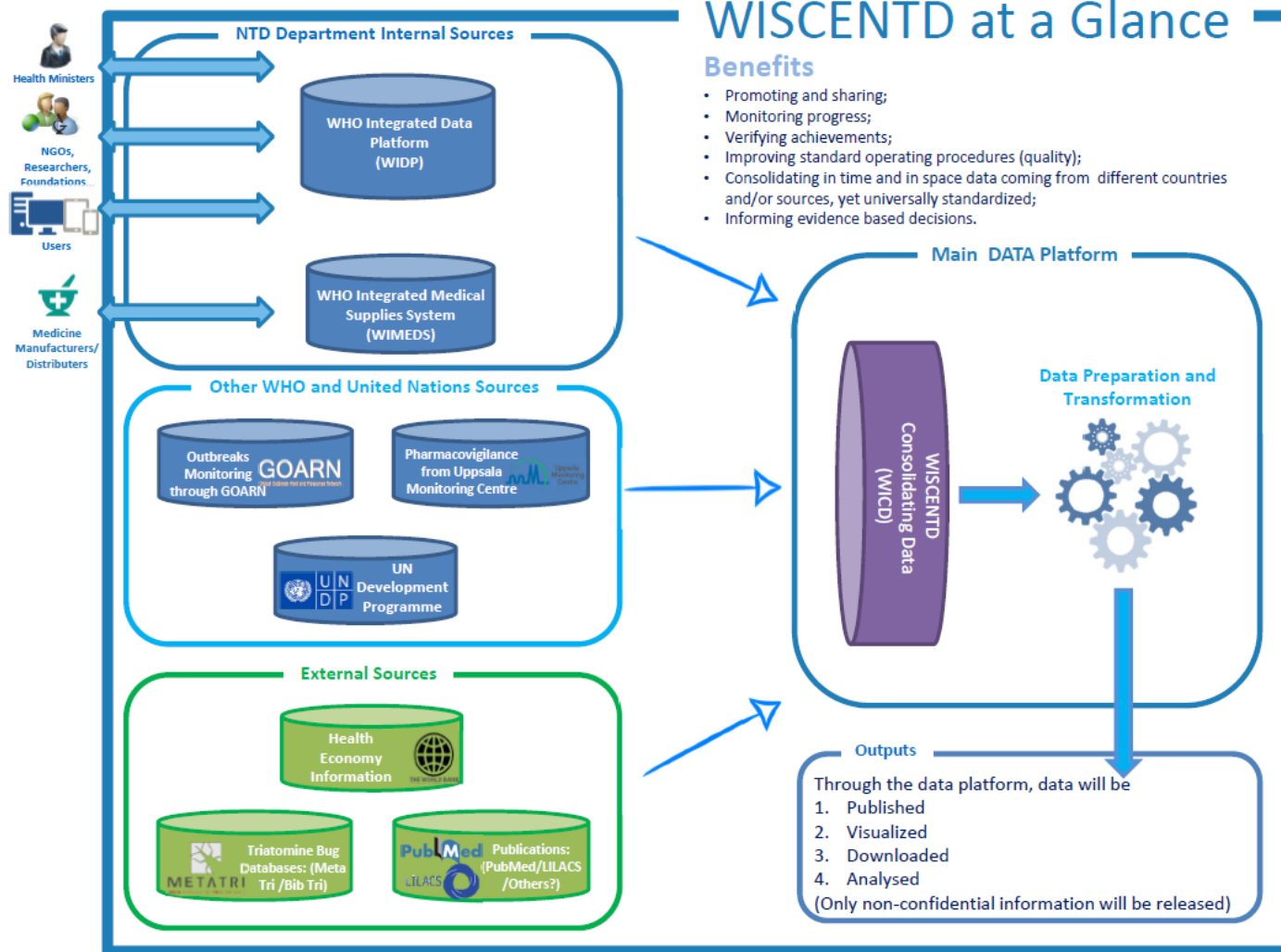
Computational Complexity

- This query rewriting algorithm is:
 - Linear in the size of the subgraph of G to navigate
 - Linear in the size of the wrappers mappings
 - Exponential in the number of wrappers that may join
 - Our experiments show that typically Big Data sources have few join points and therefore this exponential complexity is affordable in real cases



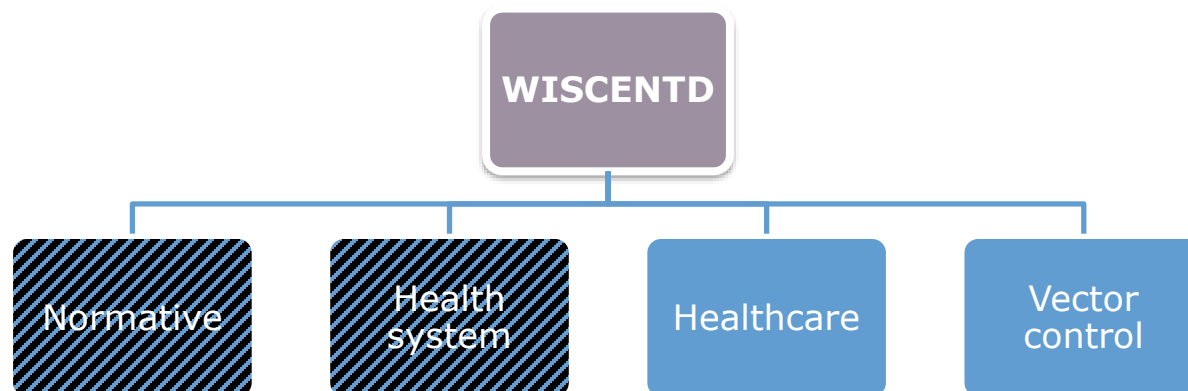
Example of application: The World Health Organisation

WISCENTD



Standardisation

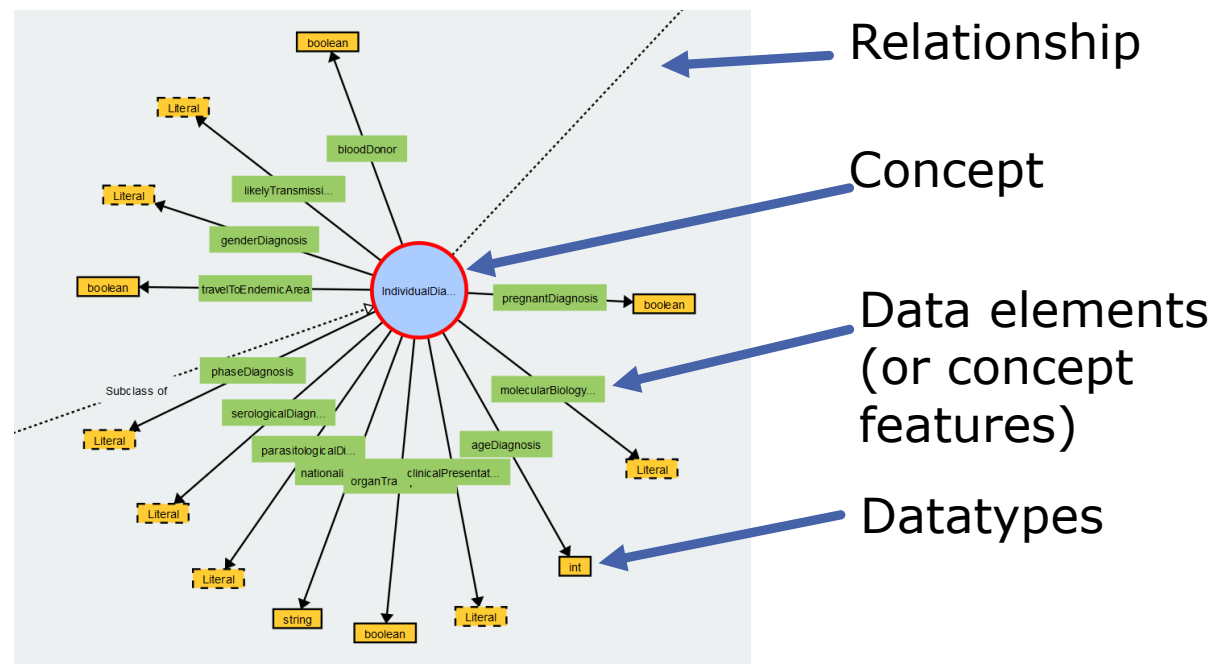
- Data has been organized into 4 packages
 - **Healthcare**: to collect patient data
 - **Vector control**: to collect data on vector control activities
 - **Health system**: to collect general information on how NTDs are included in the national health systems
 - **Normative**: to collect information about regulations implemented in to country in order to control and eliminate NTDs



Standardisation

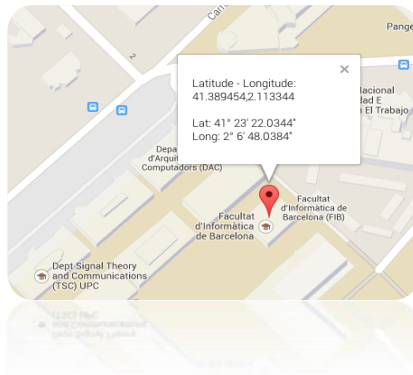
- WISCENTD provides a single **standardised** view of the whole domain

- WISCENTD provides a graph-based metaphor representing the domain:
- Concepts
 - Data elements of each concept (and their datatypes)
 - Relationships between concepts



Master Data: Geographic and Temporal Components

Coordinates



Timestamps



Polygons



Time periods

August 2013							September 2013						
Su	Mo	Tu	We	Th	Fr	Sa	Su	Mo	Tu	We	Th	Fr	Sa
28	29	30	31	1	2	3	25	26	27	28	29	30	31
4	5	6	7	8	9	10	1	2	3	4	5	6	7
11	12	13	14	15	16	17	8	9	10	11	12	13	14
18	19	20	21	22	23	24	15	16	17	18	19	20	21
25	26	27	28	29	30	31	22	23	24	25	26	27	28
1	2	3	4	5	6	7	29	30	1	2	3	4	5
Tuesday, August 13, 2013							Tuesday, August 13, 2013						

LOCATION

TIME

Data Analysis

MDM

Global Graph ▾ Data Sources ▾ Wrappers ▾ LAV Mappings ▾ Ontology-Mediated Queries ▾

Signed in as sergi ▾

Ontology querying

Legend:

Concept

Feature ID

Feature

Show features:

On

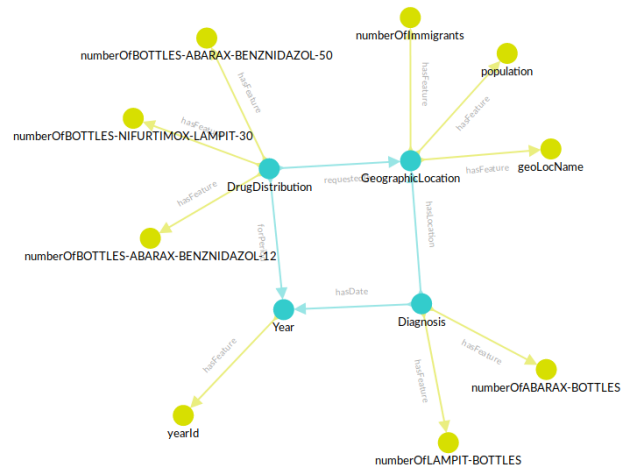
WHO

Clear Query

Get features

Execute query

Features:

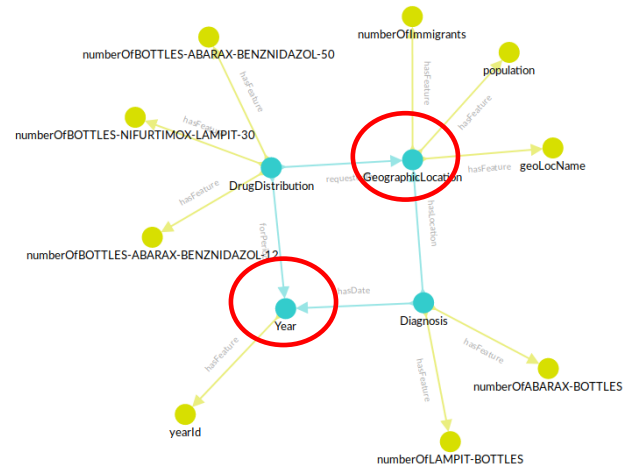


Data Analysis

MDM Global Graph ▾ Data Sources ▾ Wrappers ▾ LAV Mappings ▾ Ontology-Mediated Queries ▾ Signed in as sergi ▾

Ontology querying Legend: Concept Feature ID Feature Show features: On WHO Clear Query

Get features Execute query Features:



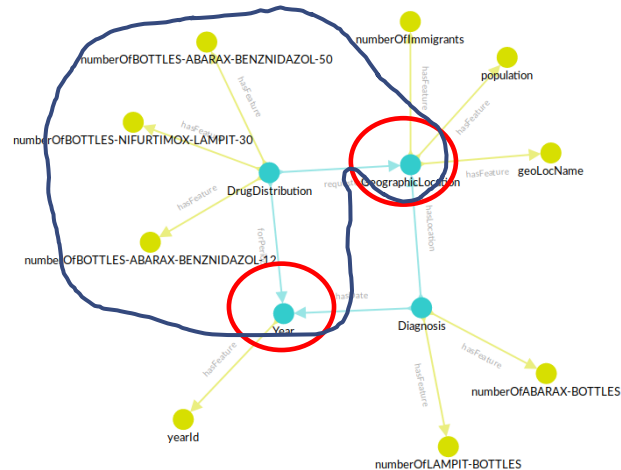
Master data:
geographical
and temporal
components

Data Analysis

MDM Global Graph ▾ Data Sources ▾ Wrappers ▾ LAV Mappings ▾ Ontology-Mediated Queries ▾ Signed in as sergi ▾

Ontology querying Legend: Concept Feature ID Feature Show features: On WHO Clear Query

Get features Execute query Features:



Master data:
geographical
and temporal
components

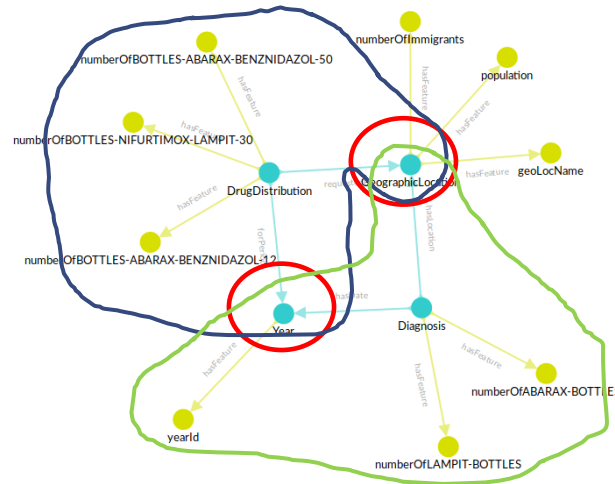
WIMEDS:
medicament
request and
distribution

Data Analysis

MDM Global Graph ▾ Data Sources ▾ Wrappers ▾ LAV Mappings ▾ Ontology-Mediated Queries ▾ Signed in as sergi ▾

Ontology querying Legend: Concept Feature ID Feature Show features: On WHO Clear Query

Get features Execute query Features:

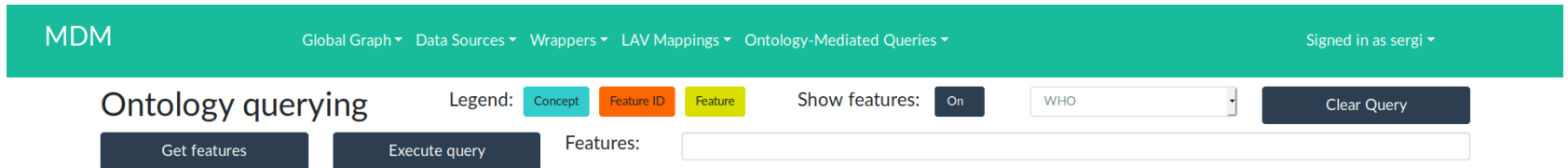


Master data:
geographical
and temporal
components

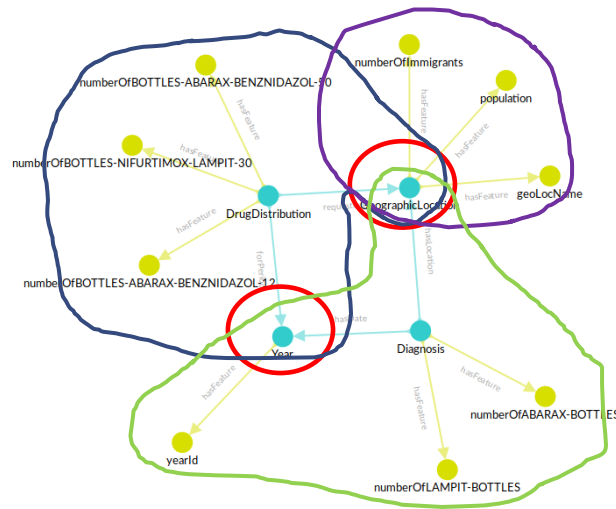
WIMEDS:
medicament
request and
distribution

WIDP:
diagnosis and
treatment

Data Analysis



UN Data:
Health
economics
(two sources:
population and
immigration
data)



Master data:
geographical
and temporal
components

WIMEDS:
medicament
request and
distribution

WIDP:
diagnosis and
treatment

Data Analysis

"I would like to correlate the number of treatments with the population and number of immigrants of a specific geographical area per year"

Data Analysis

MDM

Global Graph ▾ Data Sources ▾ Wrappers ▾ LAV Mappings ▾ Ontology-Mediated Queries ▾

Signed in as sergi ▾

Ontology querying

Legend:

Concept

Feature ID

Feature

Show features:

Off

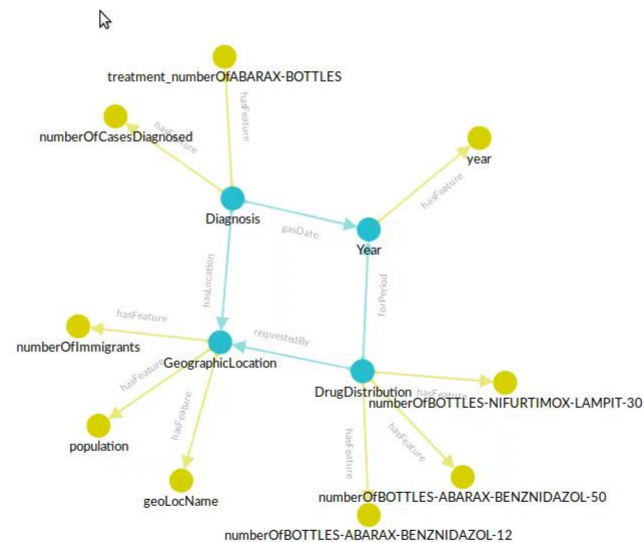
who

Clear Query

Get features

Execute query

Features:



Find this video in Learn-SQL (video 1)

Data Analysis

"I would like to correlate the number of treatments with the number of medicines distributed in a specific geographical area per year.

This information should also include the population and the number of immigrants of that area"

Data Analysis

MDM

Global Graph ▾ Data Sources ▾ Wrappers ▾ LAV Mappings ▾ Ontology-Mediated Queries ▾

Signed in as sergi ▾

Ontology querying

Legend:

Concept

Feature ID

Feature

Show features: Off

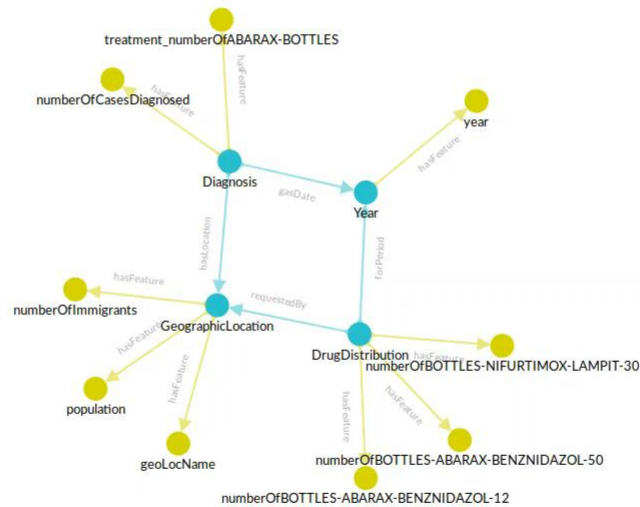
who

Clear Query

Get features

Execute query

Features:



Find this video in Learn-SQL (video 2)

Management: Extending the Ontology

- *My new data source contains data elements not covered by the attributes of the standardised model. How do I extend it?*

Management: Registering a Source

- ❑ *My data source contains data that is not covered in the attributes of the standardised model. How do I extend it?*
- ❑ *Great! Now, I want to register a new source providing such data*

MDM: Governing Evolution in Big Data Ecosystems

On-demand integration of multiple data sources is a critical requirement in many Big Data settings. This has been coined as the data variety challenge, which refers to the complexity of dealing with an heterogeneous set of data sources to enable their integrated analysis. In Big Data settings, data sources are commonly represented by external REST APIs, which provide data in their original format and continuously apply changes in their structure (i.e., schema). Thus, data analysts face the challenge to integrate such multiple sources, and then continuously adapt their analytical processes to changes in the schema. To address this challenges we present the Metadata Management System, shortly MDM, a tool that supports data stewards and analysts to manage the integration and analysis of multiple heterogeneous sources under schema evolution. MDM adopts a vocabulary-based integration-oriented ontology to conceptualize the domain of interest and relies on local-as-view mappings to link it with the sources. MDM provides user-friendly mechanisms to manage the ontology and mappings. Finally, a query rewriting algorithm ensures that queries posed to the ontology are correctly resolved to the sources in the presence of multiple schema versions, a transparent process to data analysts.

Find this video in Learn-SQL (video 4)

Management: Querying a New Source

- ❑ *My data source contains data that is not covered in the attributes of the standardised model. How do I extend it?*
- ❑ *Great! Now, I want to register a new source providing such data*
- ❑ *Good! Now "I would like to see the medicine distribution per medicine sender, per year and geographical area"*

Ontology querying

Legend:

Concept

Feature ID

Feature

Show features:

On

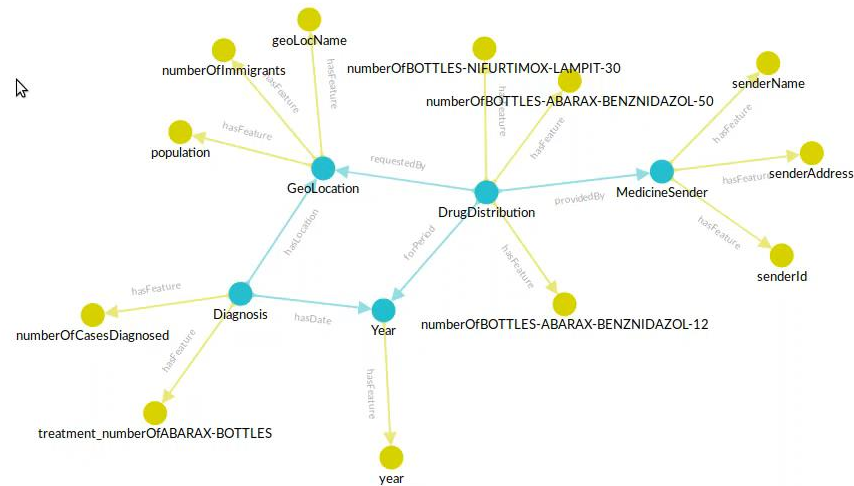
WHO

Clear Query

Get features

Execute query

Features:



Find this video in Learn-SQL (video 5)

A Tool for OMQ

□ ODIN:

<http://www.essi.upc.edu/~snadal/odin.html>

- RDFS / OWL (but limited reasoning)
- LAV mappings
- Pay-as-you-go data integration



Summary

- Graph-based Virtual Data Integration
- Ontology-based Data Access
 - DL-Lite
 - GAV mappings
 - Linking Data to Ontologies
- Ontology-mediated Queries
 - RDFS
 - LAV mappings
 - Sources exposed as wrappers

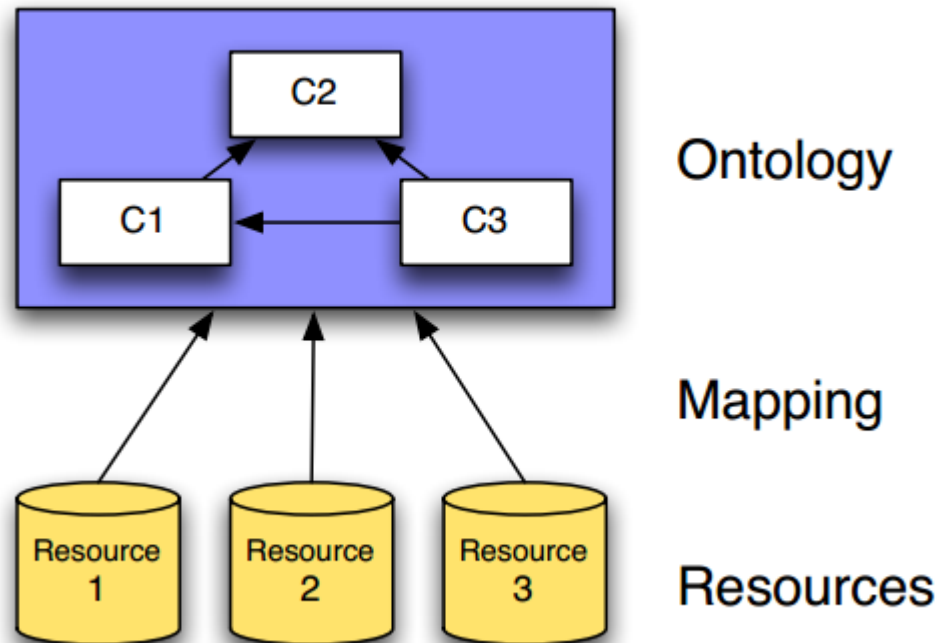
APPENDIX

GAV Data Integration

ONTOLOGY-BASED DATA ACCESS

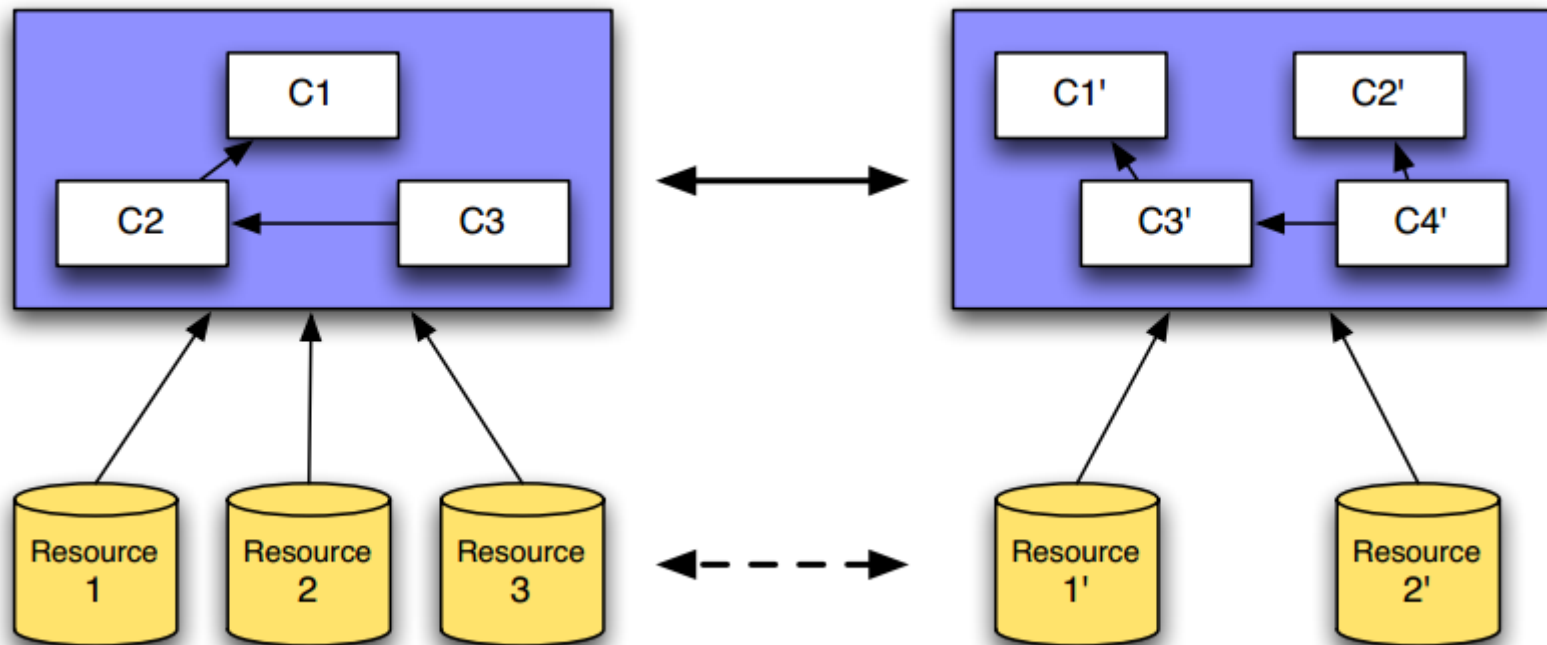
Ontology-Based Data Access

□ Ontology-mediated data access



Ontology-Based Data Access

□ An Approach for Data Integration



The DL-Lite Family

- Right trade-off between expressiveness and data complexity query answering
 - PTime in the size of the TBOX
 - LogSpace in the size of the ABOX
- Two maximal DLs satisfying this trade-off
 - DL-Lite_F
 - DL-Lite_R

Remember the DL-Lite family maps to **OWL 2 QL**

DL-Lite_F

TBox assertions:

- Concept inclusion assertions: $Cl \sqsubseteq Cr$, with:

$$\begin{array}{lcl} Cl & \longrightarrow & A \mid \exists Q \\ Cr & \longrightarrow & A \mid \exists Q \mid \neg A \mid \neg \exists Q \\ Q & \longrightarrow & P \mid P^- \end{array}$$

- Functionality assertions: (**funct** Q)

ABox assertions: $A(c)$, $P(c_1, c_2)$, with c_1, c_2 constants

Observations:

- Captures all the basic constructs of UML Class Diagrams and ER
- Notable exception: covering constraints in generalizations.

Semantics of DL-Lite

- It basically captures the expressivity of a UML class diagram

ISA between classes	$A_1 \sqsubseteq A_2$
Disjointness between classes	$A_1 \sqsubseteq \neg A_2$
Domain and range of relations	$\exists P \sqsubseteq A_1 \quad \exists P^- \sqsubseteq A_2$
Mandatory participation	$A_1 \sqsubseteq \exists P \quad A_2 \sqsubseteq \exists P^-$
Functionality of relations (in $DL-Lite_{\mathcal{F}}$)	$(\mathbf{funct} \ P) \quad (\mathbf{funct} \ P^-)$
ISA between relations (in $DL-Lite_{\mathcal{R}}$)	$Q_1 \sqsubseteq Q_2$
Disjointness between relations (in $DL-Lite_{\mathcal{R}}$)	$Q \sqsubseteq \neg Q$

Semantics of DL-Lite

Construct	Syntax	Example	Semantics
atomic conc.	A	Doctor	$A^I \subseteq \Delta^I$
exist. restr.	$\exists Q$	$\exists \text{child}^-$	$\{d \mid \exists e. (d, e) \in Q^I\}$
at. conc. neg.	$\neg A$	$\neg \text{Doctor}$	$\Delta^I \setminus A^I$
conc. neg.	$\neg \exists Q$	$\neg \exists \text{child}$	$\Delta^I \setminus (\exists Q)^I$
atomic role	P	child	$P^I \subseteq \Delta^I \times \Delta^I$
inverse role	P^-	child^-	$\{(o, o') \mid (o', o) \in P^I\}$
role negation	$\neg Q$	$\neg \text{manages}$	$(\Delta_O^I \times \Delta_O^I) \setminus Q^I$
conc. incl.	$Cl \sqsubseteq Cr$	$\text{Father} \sqsubseteq \exists \text{child}$	$Cl^I \subseteq Cr^I$
role incl.	$Q \sqsubseteq R$	$\text{hasFather} \sqsubseteq \text{child}^-$	$Q^I \subseteq R^I$
funct. asser.	(funct Q)	(funct succ)	$\forall d, e, e'. (d, e) \in Q^I \wedge (d, e') \in Q^I \rightarrow e = e'$
mem. asser.	$A(c)$	$\text{Father}(\text{bob})$	$c^I \in A^I$
mem. asser.	$P(c_1, c_2)$	$\text{child}(\text{bob}, \text{ann})$	$(c_1^I, c_2^I) \in P^I$

Linking Data to Ontologies

- The ABOX is stored in a relational database
 - OBDA has recently been extended to other kind of sources, like document-stores
- Direct mappings between the TBOX and DB
 - Query answering is reformulated in terms of the TBOX, a set of mappings and a RDBMS

Theorem

Query answering in a $DL-Lite_A$ OBDM system $\mathcal{O} = \langle \mathcal{T}, \mathcal{M}, \mathcal{D} \rangle$ is

- ① NP-complete in the size of the query.
- ② PTIME in the size of the TBox \mathcal{T} and the mappings \mathcal{M} .
- ③ LOGSPACE in the size of the database \mathcal{D} .

Mappings

- OBDA works with GAV mappings
- Typically, they use RDF-based mapping languages to express them
 - R2RML (a language to express mappings from global concepts to relational databases)

mappingId	Actor
target	imdb:name/{person_id} a dbpedia:Actor .
source	select person_id from cast_info where cast_info.role_id = 1

- RML is a generalisation to map to any kind of source (<http://rml.io/>)

A Tool for OBDA

- Ontop: <http://ontop.inf.unibz.it/>
 - OWL 2 QL
 - RDFS



- Code, examples and more:
<https://github.com/ontop/ontop>