

# Ontology-based Data Access Made Practical

Diego Calvanese

KRDB Research Centre for Knowledge and Data  
Free University of Bozen-Bolzano, Italy



Department of Computing Science  
Umeå University, Sweden



Ontopic s.r.l.

**ONTOPIC**

EDBT-INTENDED Summer School  
4–9 July 2022 – Bordeaux (France)

# Challenges in data management

**40 ZETTABYTES**  
[ 43 TRILLION GIGABYTES ]  
of data will be created by 2020, an increase of 300 times from 2005

**6 BILLION PEOPLE**  
have cell phones



**Volume**  
SCALE OF DATA

2005

2020

It's estimated that **2.5 QUINTILLION BYTES**  
[ 2.3 TRILLION GIGABYTES ]  
of data are created each day



Most companies in the U.S. have at least **100 TERABYTES**  
[ 100,000 GIGABYTES ]  
of data stored

The New York Stock Exchange captures **1 TB OF TRADE INFORMATION** during each trading session



**Velocity**  
ANALYSIS OF  
STREAMING DATA

By 2016, it is projected there will be **18.9 BILLION NETWORK CONNECTIONS** – almost 2.5 connections per person on earth



Modern cars have close to **100 SENSORS** that monitor items such as fuel level and tire pressure



## The FOUR V's of Big Data

From traffic patterns and music downloads to web history and medical records, data is recorded, stored, and analyzed to enable the technology and services that the world relies on every day. But what exactly is big data, and how can these massive amounts of data be used?

As a leader in the sector, IBM data scientists break big data into four dimensions: **Volume, Velocity, Variety and Veracity**

Depending on the industry and organization, big data encompasses information from multiple internal and external sources such as transactions, social media, enterprise content, sensors and mobile devices. Companies can leverage data to adapt their products and services to better meet customer needs, optimize operations and infrastructure, and find new sources of revenue.

By 2015 **4.4 MILLION IT JOBS** will be created globally to support big data, with 1.9 million in the United States



As of 2011, the global size of data in healthcare was estimated to be

**150 EXABYTES**  
[ 161 BILLION GIGABYTES ]



**30 BILLION PIECES OF CONTENT** are shared on Facebook every month



**Variety**  
DIFFERENT  
FORMS OF DATA



By 2014, it's anticipated there will be **420 MILLION WEARABLE, WIRELESS HEALTH MONITORS**

**4 BILLION+ HOURS OF VIDEO** are watched on YouTube each month



**400 MILLION TWEETS** are sent per day by about 200 million monthly active users

**1 IN 3 BUSINESS LEADERS** don't trust the information they use to make decisions



**27% OF RESPONDENTS**

in one survey were unsure of how much of their data was inaccurate

**Veracity**  
UNCERTAINTY  
OF DATA

Poor data quality costs the US economy around **\$3.1 TRILLION A YEAR**

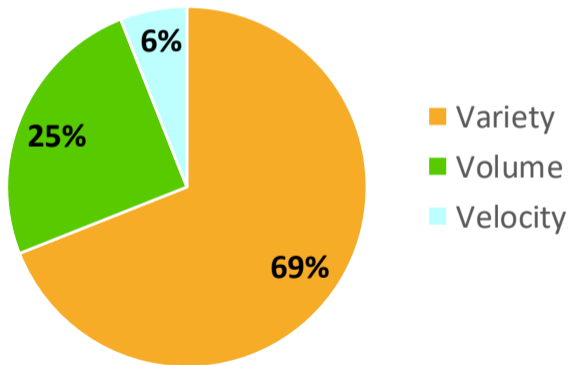


unibz

# Variety, not volume, is driving data management initiatives

MIT Sloan Management Review (28 March 2016)

## Relative Importance



<http://sloanreview.mit.edu/article/variety-not-volume-is-driving-big-data-initiatives/>

# How much time is spent searching for the right data?



Important problem: searching for data and establishing its quality

Example: in oil&gas, engineers spend 30–70% of their time on this problem  
(Crompton, 2008)

# Challenge: Accessing legacy data sources

## Statoil (now Equinor) Exploration

*Geologists at Statoil, prior to making decisions on drilling new wellbores, need to gather relevant information about previous drillings.*

*Slegge* relational database:

- Terabytes of relational data
- 1,545 tables and 1727 views
- each with dozens of attributes
- consulted by 900 geologists



# Problem: Translating information needs

## Information need expressed by geologists

*In my geographical area of interest, return all pressure data tagged with key stratigraphy information with understandable quality control attributes, and suitable for further filtering.*

To obtain the answer, this needs to be translated into SQL:

- Main table for wellbores has 38 columns (with cryptic names).
- To obtain pressure data requires a 4-table join with two additional filters.
- To obtain stratigraphic information requires a join with 5 more tables.

# Problem: Translating information needs

We would obtain the following SQL query:

```
SELECT WELLBORE.IDENTIFIER, PTY_PRESSURE.PTY_PRESSURE_S,  
       STRATIGRAPHIC_ZONE.STRAT_COLUMN_IDENTIFIER, STRATIGRAPHIC_ZONE.STRAT_UNIT_IDENTIFIER  
FROM WELLBORE,  
     PTY_PRESSURE,  
     ACTIVITY FP_DEPTH_DATA  
  LEFT JOIN (PTY_LOCATION_1D FP_DEPTH_PT1_LOC  
            INNER JOIN PICKED_STRATIGRAPHIC_ZONES ZS  
              ON ZS.STRAT_ZONE_ENTRY_MD <= FP_DEPTH_PT1_LOC.DATA_VALUE_1_0 AND  
                ZS.STRAT_ZONE_EXIT_MD >= FP_DEPTH_PT1_LOC.DATA_VALUE_1_0 AND  
                ZS.STRAT_ZONE_DEPTH_UOM = FP_DEPTH_PT1_LOC.DATA_VALUE_1_OU  
            INNER JOIN STRATIGRAPHIC_ZONE  
              ON  ZS.WELLBORE = STRATIGRAPHIC_ZONE.WELLBORE AND  
                ZS.STRAT_COLUMN_IDENTIFIER = STRATIGRAPHIC_ZONE.STRAT_COLUMN_IDENTIFIER AND  
                ZS.STRAT_INTERP_VERSION = STRATIGRAPHIC_ZONE.STRAT_INTERP_VERSION  AND  
                ZS.STRAT_ZONE_IDENTIFIER = STRATIGRAPHIC_ZONE.STRAT_ZONE_IDENTIFIER)  
  ON FP_DEPTH_DATA.FACILITY_S = ZS.WELLBORE AND  
     FP_DEPTH_DATA.ACTIVITY_S  = FP_DEPTH_PT1_LOC.ACTIVITY_S,  
     ACTIVITY_CLASS FORM_PRESSURE_CLASS  
WHERE WELLBORE.WELLBORE_S = FP_DEPTH_DATA.FACILITY_S AND  
     FP_DEPTH_DATA.ACTIVITY_S = PTY_PRESSURE.ACTIVITY_S AND  
     FP_DEPTH_DATA.KIND_S = FORM_PRESSURE_CLASS.ACTIVITY_CLASS_S AND  
     WELLBORE.REF_EXISTENCE_KIND = 'actual' AND  
     FORM_PRESSURE_CLASS.NAME = 'formation pressure depth data'
```

# Problem: Translating information needs

We would obtain the following SQL query:

```
SELECT WELLBORE_IDENTIFIED_BY_PRESSURE_DEPTH_DATA  
        STRATIGRAPHIC_ZONE  
FROM WELLBORE_ACTIVITIES  
        FORMATION_PRESSURE_DEPTH_DATA  
        LEAKAGE_DATA
```

**This can be very time consuming, and requires knowledge of the domain of interest, a deep understanding of the database structure, and general IT expertise.**

```
INNER JOIN STRATIGRAPHIC_ZONE  
        ON ZS.WELLBORE = STRATIGRAPHIC_ZONE.WELLBORE AND
```

```
ACTIVITIES  
WHERE WELLBORE_ACTIVITIES  
        FORMATION_PRESSURE_DEPTH_DATA  
        FP_DEPTH_DATA.KIND_S = FORMATION_PRESSURE_DEPTH_DATA.ACTIVITY_CLASS_S AND  
        WELLBORE.REF_EXISTENCE_KIND = 'actual' AND  
        FORMATION_PRESSURE_DEPTH_DATA.NAME = 'formation pressure depth data'
```

**This is also very costly!**

**Equinor loses 50.000.000€ per year only due to this problem!!**



# How to address data access (and integration) challenges?

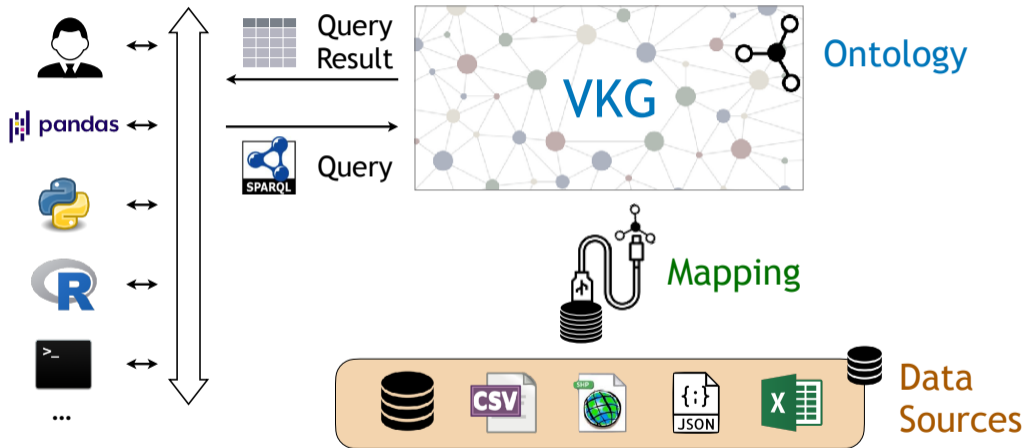
We combine three key ideas:

- 1 Use a global (or integrated) schema and **map the data sources to the global schema**.
- 2 Adopt a very flexible data model for the global schema  
     $\rightsquigarrow$  **Knowledge Graph** whose vocabulary is expressed in an **ontology**.
- 3 Exploit **virtualization**, i.e., the KG is not materialized, but kept virtual.

This gives rise to the **Virtual Knowledge Graph (VKG)** approach to data access/integration, also called **Ontology-based Data Access/Integration (OBDA)**.

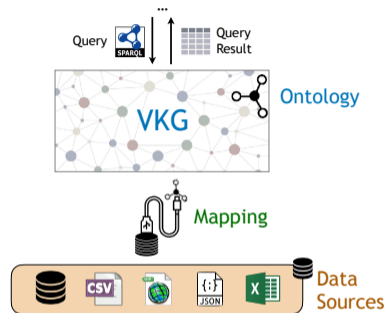
[Xiao, C., et al. 2018, IJCAI]

# Virtual Knowledge Graph (VKG) architecture



# Why an ontology?

An ontology is a structured formal representation of concepts and their relationships that are relevant for the domain of interest.



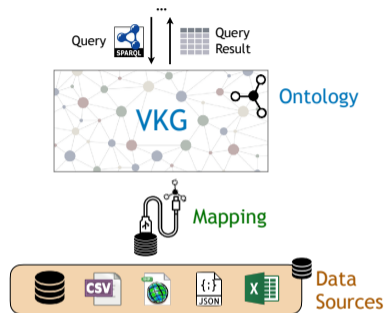
- In the VKG setting, the ontology has a twofold purpose:
  - It defines a **vocabulary of terms** to denote classes and properties that are familiar to the user.
  - It extends the data in the sources with **background knowledge about the domain of interest**, and this knowledge is machine processable.
- One can make use of **custom-built domain ontologies**.
- In addition, one can rely on **standard ontologies**, which are available for many domains.

# Why a Knowledge Graph for the global schema?

The traditional approach to data integration adopts a relational global schema.

A **Knowledge Graph**, instead:

- Does not require to commit early on to a specific structure.
- Can better accommodate heterogeneity.
- Can better deal with missing / incomplete information.
- Does not require complex restructuring operations to accommodate new information or new data sources.

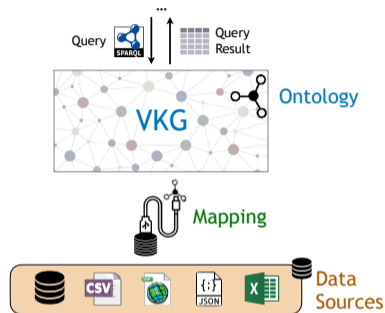


# Why mappings?

The traditional approach to data integration relies on mediators, which are specified through complex code.

**Mappings**, instead:

- Provide a **declarative specification**, and not code.
- Are **easier to understand**, and hence to design and to maintain.
- Support an **incremental approach** to integration.
- Are **machine processable**, hence are used in query answering and for query optimization.

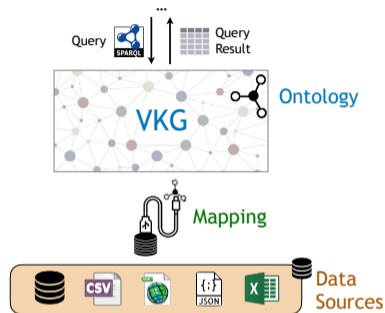


# Why virtualization?

Materialized data integration relies on extract-transform-load (ETL) operations, to load data from the sources into an integrated data store / data warehouse / materialized KG.

In the **virtual approach**, instead:

- The data stays in the sources and is only accessed at query time.
- No need to construct a large and potentially costly materialized data store and keep it up-to-date.
- Hence the data is always fresh wrt the latest updates at the sources.
- One can rely on the existing data infrastructure and expertise.
- There is better support for an incremental approach to integration.



# Outline

- 1 Challenges in Data Access
- 2 Virtual Knowledge Graphs for Data Access (and Integration)
- 3 Optimizing Query Answering in VKGs
- 4 Designing VKG Mappings
- 5 The Ontop System
- 6 Conclusions

# Outline

- 1 Challenges in Data Access
- 2 Virtual Knowledge Graphs for Data Access (and Integration)
  - Representing Data in RDF and RDFS
  - Representing Ontologies in OWL 2 QL
  - Query Language – SPARQL
  - Mapping an Ontology to a Relational Database
  - Formalizing the VKG Framework
- 3 Optimizing Query Answering in VKGs
- 4 Designing VKG Mappings
- 5 The Ontop System
- 6 Conclusions



# Incomplete information

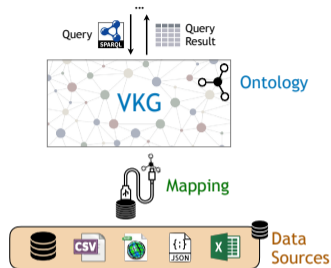
We are in a setting of **incomplete information!!!**

Incompleteness is introduced:

- by data sources, in general assumed to be incomplete;
- by domain constraints encoded in the ontology.

## Plus:

Ontologies are logical theories, and hence perfectly suited to deal with incomplete information!



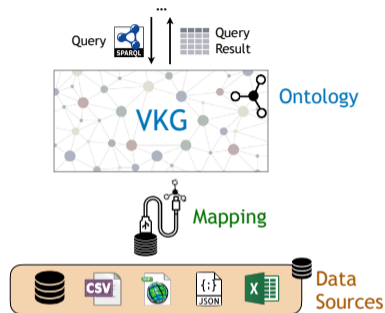
## Minus:

Query answering amounts to **logical inference**, and hence is significantly more challenging.

# Components of the VKG framework

We consider now the main components that make up the VKG framework, and the languages used to specify them.

In defining such languages, we need to consider the **tradeoff between expressive power and efficiency**, where the key point is efficiency with respect to the data.



The W3C has standardized languages that are suitable for VKGs:

- 1 Knowledge graph: expressed in **RDF** [W3C Rec. 2014] (v1.1)
- 2 Ontology  $\mathcal{O}$ : expressed in **OWL 2 QL** [W3C Rec. 2012]
- 3 Mapping  $\mathcal{M}$ : expressed in **R2RML** [W3C Rec. 2012]
- 4 Query: expressed in **SPARQL** [W3C Rec. 2013] (v1.1)

# Outline

- 1 Challenges in Data Access
- 2 Virtual Knowledge Graphs for Data Access (and Integration)**
  - Representing Data in RDF and RDFS
  - Representing Ontologies in OWL 2 QL
  - Query Language – SPARQL
  - Mapping an Ontology to a Relational Database
  - Formalizing the VKG Framework
- 3 Optimizing Query Answering in VKGs
- 4 Designing VKG Mappings
- 5 The Ontop System
- 6 Conclusions

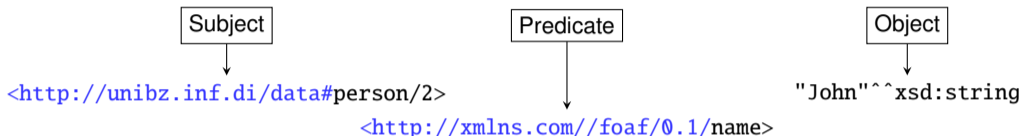
# Resource Description Framework (RDF)

- RDF is a language standardized by the W3C for representing information [W3C Rec. 2004] (v1.0) and [W3C Rec. 2014] (v1.1).
- RDF is a **graph-based data model**, where information is represented as (labeled) nodes connected by (labeled) edges.
- Nodes have three different forms:
  - literal: denotes a constant value, with an associated datatype;
  - IRI (for *internationalized resource identifier*): denotes a resource (i.e., an object), for which the IRI acts as an identifier;
  - blank node: represents an anonymous object.
- An IRI might also denote a **property**, connecting an object to a literal, or connecting two objects.

See also <https://www.w3.org/TR/rdf11-concepts/> for details.

# RDF triples

RDF provides a description of the domain of interest in terms of **triples**:



Triple elements: resources denoted by **global identifiers** (called **IRIs**)

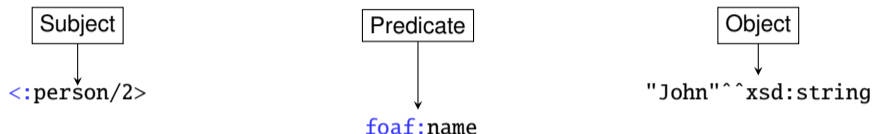
- 1 Subject: IRI of the described resource
- 2 Predicate: IRI of the property
- 3 Object: attribute value or IRI of another resource

**Prefixes**: useful abbreviations and/or references to external information

```
@prefix foaf: <http://xmlns.com/foaf/0.1/>
@prefix : <http://unibz.inf.di/data#>
@base <http://unibz.inf.di/>
```

# RDF triples

RDF provides a description of the domain of interest in terms of **triples**:



Triple elements: resources denoted by **global identifiers** (called **IRIs**)

- ① Subject: IRI of the described resource
- ② Predicate: IRI of the property
- ③ Object: attribute value or IRI of another resource

**Prefixes**: useful abbreviations and/or references to external information

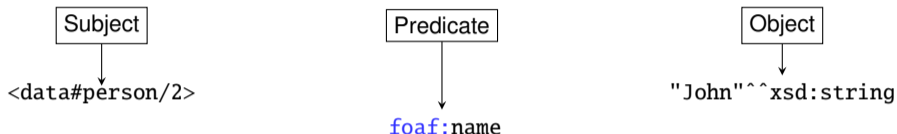
```
@prefix foaf: <http://xmlns.com/foaf/0.1/>
```

```
@prefix : <http://unibz.inf.di/data#>
```

```
@base <http://unibz.inf.di/>
```

# RDF triples

RDF provides a description of the domain of interest in terms of **triples**:



Triple elements: resources denoted by **global identifiers** (called **IRIs**)

- 1 Subject: IRI of the described resource
- 2 Predicate: IRI of the property
- 3 Object: attribute value or IRI of another resource

**Prefixes**: useful abbreviations and/or references to external information

@prefix foaf: <http://xmlns.com/foaf/0.1/>

@prefix : <http://unibz.inf.di/data#>

@base <http://unibz.inf.di/>

# RDF – Examples

## Class membership:

RDF triple	<code>&lt;uni2/p/25&gt; rdf:type :Professor</code>
Fact	<code>Professor(uni2/p/25)</code>

Note: This is typically abbreviated as

RDF triple	<code>&lt;uni2/p/25&gt; a :Professor</code>
------------	---

## Data property of an individual:

RDF triple	<code>&lt;uni2/p/25&gt; :lastName "Artale"</code>
Fact	<code>lastName(uni2/p/25, "Artale")</code>

## Object property of an individual:

RDF triple	<code>&lt;uni2/p/25&gt; :teaches &lt;uni2/c/7&gt;</code>
Fact	<code>teaches(uni2/p/25, uni2/c/7)</code>



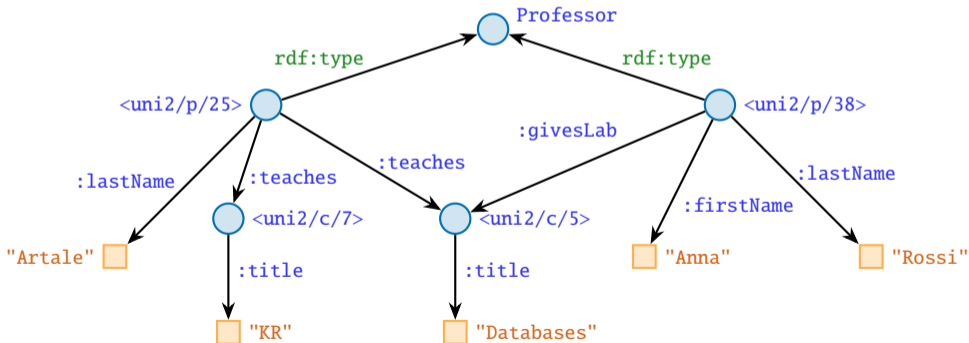
# RDF graph – Example

```

<uni2/p/25> rdf:type :Professor
<uni2/p/25> foaf:lastName "Artale"
<uni2/p/25> :teaches <uni2/c/5>
...

```

We can represent such a set of facts graphically:



# Additional RDF features

RDF has additional features that we do not cover here:

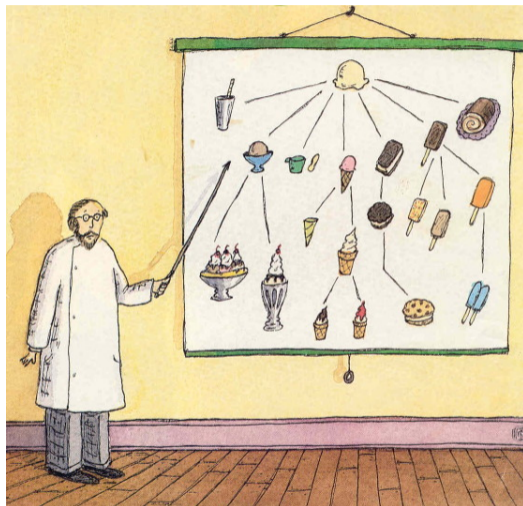
- datatypes
- blank nodes
- named graphs

# Outline

- 1 Challenges in Data Access
- 2 Virtual Knowledge Graphs for Data Access (and Integration)**
  - Representing Data in RDF and RDFS
  - Representing Ontologies in OWL 2 QL**
  - Query Language – SPARQL
  - Mapping an Ontology to a Relational Database
  - Formalizing the VKG Framework
- 3 Optimizing Query Answering in VKGs
- 4 Designing VKG Mappings
- 5 The Ontop System
- 6 Conclusions

# What is an ontology?

- An ontology conceptualizes a domain of interest in terms of **concepts/classes**, (binary) **relations**, and their **properties**.
- It typically organizes the concepts in a hierarchical structure.
- Ontologies are often represented as graphs.
- However, an ontology is actually a **logical theory**, expressed in a suitable fragment of first-order logic, or better, in **description logics**.





# What is an ontology?

- An ontology conceptualizes a domain of interest in terms of **concepts/classes**, (binary) **relations**, and their **properties**.
- It typically organizes the concepts in a hierarchical structure.
- Ontologies are often represented as graphs.
- However, an ontology is actually a **logical theory**, expressed in a suitable fragment of first-order logic, or better, in **description logics**.

$$\forall x. \text{Actor}(x) \rightarrow \text{Staff}(x)$$

$$\forall x. \text{SeriesActor}(x) \rightarrow \text{Actor}(x)$$

$$\forall x. \text{MovieActor}(x) \rightarrow \text{Actor}(x)$$

$$\forall x. \text{SeriesActor}(x) \rightarrow \neg \text{MovieActor}(x)$$

$$\forall x. \text{Staff}(x) \rightarrow \exists y. \text{ssn}(x, y)$$

$$\forall y. \exists x. \text{ssn}(x, y) \rightarrow \text{xsd:int}(y)$$

$$\forall x, y, y'. \text{ssn}(x, y) \wedge \text{ssn}(x, y') \rightarrow y = y'$$

$$\forall x. \exists y. \text{actsIn}(x, y) \rightarrow \text{MovieActor}(x)$$

$$\forall y. \exists x. \text{actsIn}(x, y) \rightarrow \text{Movie}(y)$$

$$\forall x. \text{MovieActor}(x) \rightarrow \exists y. \text{actsIn}(x, y)$$

$$\forall x. \text{Movie}(x) \rightarrow \exists y. \text{actsIn}(y, x)$$

$$\forall x, y. \text{actsIn}(x, y) \rightarrow \text{playsIn}(x, y)$$

...

# What is an ontology?

- An ontology conceptualizes a domain of interest in terms of **concepts/classes**, (binary) **relations**, and their **properties**.
- It typically organizes the concepts in a hierarchical structure.
- Ontologies are often represented as graphs.
- However, an ontology is actually a **logical theory**, expressed in a suitable fragment of first-order logic, or better, in **description logics**.

$\text{Actor} \sqsubseteq \text{Staff}$   
 $\text{SeriesActor} \sqsubseteq \text{Actor}$   
 $\text{MovieActor} \sqsubseteq \text{Actor}$   
 $\text{SeriesActor} \sqsubseteq \neg \text{MovieActor}$

$\text{Staff} \sqsubseteq \exists \text{ssn}$   
 $\exists \text{ssn}^- \sqsubseteq \text{xsd:int}$   
 (func ssn)

$\exists \text{actsIn} \sqsubseteq \text{MovieActor}$   
 $\exists \text{actsIn}^- \sqsubseteq \text{Movie}$   
 $\text{MovieActor} \sqsubseteq \exists \text{actsIn}$   
 $\text{Movie} \sqsubseteq \exists \text{actsIn}^-$   
 $\text{actsIn} \sqsubseteq \text{playsIn}$   
 ...

# The OWL 2 QL ontology language

- **OWL 2 QL** is one of the three standard profiles of OWL 2. [W3C Rec. 2012]
- Is derived from the *DL-Lite<sub>R</sub>* description logic (DL) of the *DL-Lite*-family [C., De Giacomo, et al. 2007].
- Is considered a lightweight ontology language:
  - controlled expressive power
  - efficient inference
- Optimized for accessing large amounts of data (i.e., for data complexity):
  - Queries over the ontology can be rewritten into SQL queries over the underlying relational database (**First-order rewritability** of query answering).
  - Consistency of ontology and data can also be checked by executing SQL queries (i.e., it is also first-order rewritable).



# Classes and properties in OWL 2 QL

All ontology languages based on OWL 2 (and hence also OWL 2 QL), provide three types of elements to construct an ontology:

- **Classes** (also called **concepts**), which allow one to structure the domain of interest, by grouping in a class objects with common properties.

**Examples:** [Movie](#), [Staff](#), [Actor](#), [SeriesActor](#), ...

- **Data properties** (also called **attributes**), which are binary relations that relate objects to values (or literals, in RDF terminology).


**Examples:**

- **title**, associating a string to a [Movie](#);
- **ssn**, associating an integer to a [Person](#).

- **Object properties** (also called **roles**), which are binary relations between objects.

**Examples:**

- **actsIn**, relating a [MovieActor](#) to a [Movie](#);
- **worksFor**, relating an [Employee](#) to a [Project](#).

In the following, to depict an OWL 2 QL ontology, we make use of a **graphical notation** inspired by  the one for UML class diagrams.

# OWL 2 QL knowledge bases

An **OWL 2 QL knowledge base (KB)** consists of two parts:

An **ontology**  $\mathcal{O}$  modeling the schema level information.

- Contains the declarations of the classes, data properties, and object properties of the ontology. This constitutes the **vocabulary** with which we can then query the ontology.
- Contains the **axioms** that capture the **domain knowledge**.
- These axioms express the conditions that must hold for the classes and properties in the ontology.

An **RDF graph**  $\mathcal{G}$ , modeling the extensional level information (i.e., facts).

The RDF graph  $\mathcal{G}$  consists of triples that express membership assertions of the following forms:

- An individual  $\langle a \rangle$  belongs to a class  $:C$ :  $\langle a \rangle \text{ rdf:type } :C$  .
- A pair individual  $\langle a \rangle$  and literal  $\langle l \rangle$  belongs to a data property  $:A$ :  $\langle a \rangle :A \langle l \rangle$  .
- A pair of individuals  $\langle a1 \rangle, \langle a2 \rangle$  belongs to an object property  $:P$ :  $\langle a1 \rangle :P \langle a2 \rangle$  .

*Note:* As we will see later, in the VKG setting, the RDF graph of a KB is not given explicitly, but is (usually) defined implicitly through the database(s) and the mappings.

# Axioms in an OWL 2 QL ontology

In an OWL 2 QL ontology, one can express knowledge about the classes and properties in the domain of interest by means of the following types of axioms.

Axiom type	OWL syntax	DL syntax
Class declaration	<code>:Actor rdf:type owl:Class</code>	<code>Actor</code>
Object property decl.	<code>:actsIn rdf:type owl:ObjectProperty</code>	<code>actsIn</code>
Data property declaration	<code>:title rdf:type owl:DatatypeProperty</code>	<code>title</code>
Subclass assertion	<code>:MovieActor rdfs:subClassOf :Actor</code>	<code>MovieActor <math>\sqsubseteq</math> Actor</code>
Class disjointness	<code>:Actor owl:disjointWith :Movie</code>	<code>Actor <math>\sqsubseteq</math> <math>\neg</math>Movie</code>
Domain of a property	<code>:actsIn rdfs:domain :MovieActor</code>	<code><math>\exists</math>actsIn <math>\sqsubseteq</math> MovieActor</code>
Range of a property	<code>:actsIn rdfs:range :Movie</code>	<code><math>\exists</math>actsIn<sup>-</sup> <math>\sqsubseteq</math> Movie</code>
Mandatory participation	<code>owl:someValuesFrom</code> in superclass expression	<code>MovieActor <math>\sqsubseteq</math> <math>\exists</math>actsIn</code>
Subproperty assertion	<code>:actsIn rdfs:subPropertyOf :playsIn</code>	<code>actsIn <math>\sqsubseteq</math> playsIn</code>
Inverse properties	<code>:actsIn owl:inverseOf :hasActor</code>	<code>actsIn <math>\equiv</math> hasActor<sup>-</sup></code>

# Syntax and semantics of OWL 2 QL axioms

Axiom type	OWL Syntax	DL Syntax	FOL
Class declaration	$C$ <code>rdf:type owl:Class</code>	$C$	$C(x)$
Object property declaration	$P$ <code>rdf:type owl:ObjectProperty</code>	$P$	$P(x, y)$
Data property declaration	$A$ <code>rdf:type owl:DatatypeProperty</code>	$A$	$A(x, y)$
Subclass assertion	$C_1$ <code>rdfs:subClassOf C_2</code>	$C_1 \sqsubseteq C_2$	$\forall x. C_1(x) \rightarrow C_2(x)$
Class disjointness	$C_1$ <code>owl:disjointWith C_2</code>	$C_1 \sqsubseteq \neg C_2$	$\forall x. C_1(x) \rightarrow \neg C_2(x)$
Domain of a property	$P$ <code>rdfs:domain C_1</code>	$\exists P \sqsubseteq C_1$	$\forall x. (\exists y. P(x, y)) \rightarrow C_1(x)$
Range of a property	$P$ <code>rdfs:range C_2</code>	$\exists P^- \sqsubseteq C_2$	$\forall y. (\exists x. P(x, y)) \rightarrow C_2(y)$
Mandatory participation	using <code>owl:someValuesFrom</code>	$C \sqsubseteq \exists R$	$\forall x. C(x) \rightarrow \exists y. R(x, y)$
Subproperty assertion	$P_1$ <code>rdfs:subPropertyOf R_2</code>	$P_1 \sqsubseteq R_2$	$\forall x, y. P_1(x, y) \rightarrow R_2(x, y)$
Inverse property	$P_2$ <code>owl:inverseOf P_1</code>	$P_1 \equiv P_2^-$	$\forall x, y. P_1(x, y) \leftrightarrow P_2(y, x)$

- We have used  $R$  to denote either an object property  $P$  or the inverse  $P^-$  of an object property.
- We have listed the axioms involving object properties, but OWL 2 QL allows for analogous axioms involving data properties, except that we might not use the inverse of a data property.

# Impact of disjointness and functionality on query answering

- Disjointness of classes and of properties cannot be expressed in RDFS, but can be expressed in OWL 2 QL.
- Functionality of properties cannot be expressed in OWL 2 QL, but can be expressed in OWL 2 (which is a much more powerful ontology language).
- However, both disjointness and functionality are supported by VKG systems such as *Ontop*.
- These constructs have an impact on consistency, i.e., they might be violated by the data and thus lead to an RDF graph that is inconsistent with the ontology.
- It turns out, however, that neither disjointness nor functionality affect query answering, as long as the ontology and the data are consistent. This means that they are actually ignored by the query evaluation algorithm.

# Representing OWL 2 QL ontologies as UML class diagrams/ER schemas

There is a close correspondence between OWL 2 QL and conceptual modeling formalisms, such as UML class diagrams and ER schemas. [Lenzerini & Nobili 1990; Bergamaschi & Sartori 1992; Borgida 1995; C., Lenzerini, et al. 1999; Borgida & Brachman 2003; Berardi et al. 2005; Queralt et al. 2012].

SeriesActor  $\sqsubseteq$  Actor

SeriesActor  $\sqsubseteq$   $\neg$ MovieActor

$\exists$ playsIn  $\sqsubseteq$  Actor

$\exists$ playsIn<sup>-</sup>  $\sqsubseteq$  Play

MovieActor  $\sqsubseteq$   $\exists$ actsIn

actsIn  $\sqsubseteq$  playsIn

...

rdfs:subClassOf

owl:disjointWith

rdfs:domain

rdfs:range

owl:someValuesFrom

rdfs:subPropertyOf

subclass

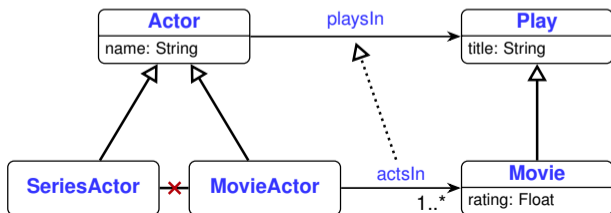
disjointness

domain

range

mandatory participation

sub-association



In fact, to visualize an OWL 2 QL ontology, we could have used standard UML class diagrams, instead of the specific graphical notation that we have introduced.

# Representing OWL 2 QL ontologies as UML class diagrams/ER schemas

There is a close correspondence between OWL 2 QL and conceptual modeling formalisms, such as UML class diagrams and ER schemas. [Lenzerini & Nobili 1990; Bergamaschi & Sartori 1992; Borgida 1995; C., Lenzerini, et al. 1999; Borgida & Brachman 2003; Berardi et al. 2005; Queralt et al. 2012].

SeriesActor  $\sqsubseteq$  Actor

SeriesActor  $\sqsubseteq$   $\neg$ MovieActor

$\exists$ playsIn  $\sqsubseteq$  Actor

$\exists$ playsIn $^-$   $\sqsubseteq$  Play

MovieActor  $\sqsubseteq$   $\exists$ actsIn

actsIn  $\sqsubseteq$  playsIn

...

rdfs:subClassOf

owl:disjointWith

rdfs:domain

rdfs:range

owl:someValuesFrom

rdfs:subPropertyOf

subclass

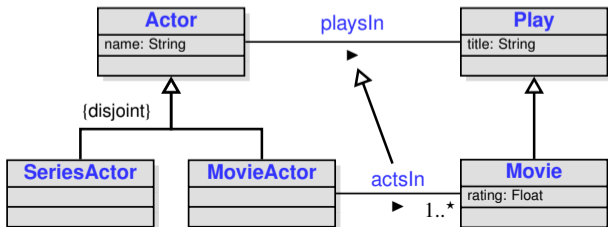
disjointness

domain

range

mandatory participation

sub-association



In fact, to visualize an OWL 2 QL ontology, we could have used standard UML class diagrams, instead of the specific graphical notation that we have introduced.

# Outline

- 1 Challenges in Data Access
- 2 Virtual Knowledge Graphs for Data Access (and Integration)**
  - Representing Data in RDF and RDFS
  - Representing Ontologies in OWL 2 QL
  - Query Language – SPARQL**
  - Mapping an Ontology to a Relational Database
  - Formalizing the VKG Framework
- 3 Optimizing Query Answering in VKGs
- 4 Designing VKG Mappings
- 5 The Ontop System
- 6 Conclusions



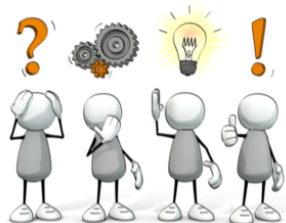
# Query answering – Which query language to use

## Querying under **incomplete information**

**Query answering** is not simply query evaluation, but **a form of logical inference**, and requires reasoning.

Two borderline cases for choosing the language for querying ontologies:

- 1 Use the **ontology language** as query language.
  - Ontology languages are tailored for capturing intensional relationships.
  - They are quite **poor as query languages**.
- 2 Use **Full SQL** (or equivalently, first-order logic).
  - Problem: in a setting with incomplete information, **query answering is undecidable** (FOL validity).



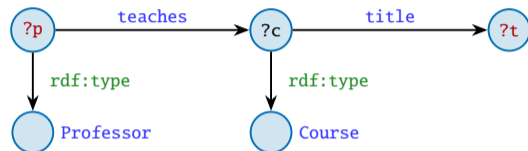
## Conjunctive queries – Are concretely represented in **SPARQL**

A good tradeoff is to use **conjunctive queries** (CQs) or unions of CQs (UCQs), corresponding to SQL/relational algebra **(union) select-project-join queries**.

# SPARQL query language

- Is the standard query language for RDF data. [W3C Rec. 2008, 2013]
- Core query mechanism is based on **graph matching**.

```
SELECT ?p ?t
WHERE {
  ?p rdf:type Professor .
  ?p teaches ?c .
  ?c rdf:type Course .
  ?c title ?t .
}
```



Additional language features (SPARQL 1.1):

- UNION: matches one of alternative graph patterns
- OPTIONAL: produces a match even when part of the pattern is missing
- complex FILTER conditions
- GROUP BY, to express aggregations
- MINUS, to remove possible solutions
- property paths (regular expressions)
- ...

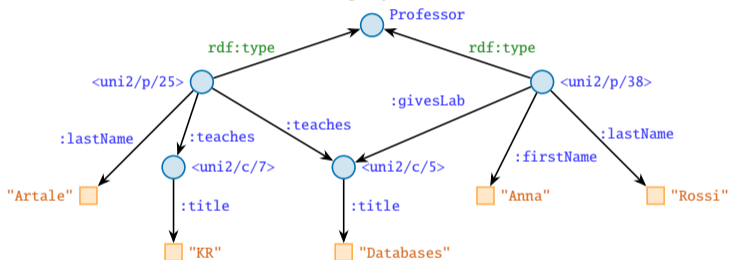
# SPARQL Basic Graph Patterns

**Basic Graph Pattern (BGP)** are the simplest form of SPARQL query, asking for a pattern in the RDF graph, made up of triple patterns.

## Example: BGP

```
SELECT ?p ?ln ?c ?t
WHERE {
  ?p :lastName ?ln .
  ?p :teaches ?c .
  ?c :title ?t .
}
```

When evaluated over the RDF graph



... the query returns:

p	ln	c	t
<uni2/p/25>	"Artale"	<uni2/c/5>	"Databases"
<uni2/p/25>	"Artale"	<uni2/c/7>	"KR"

# Abbreviated syntax for Basic Graph Patterns

We can use an abbreviated syntax for BGPs, that avoids repeating the subject of triple patterns.

## Example: BGP

```
SELECT ?p ?ln ?c ?t ?r
WHERE {
  ?p :lastName ?ln .
  ?p :teaches ?c .
  ?c :title ?t .
  ?c :room ?r .
}
```

## Example: BGP with abbreviated syntax

```
SELECT ?p ?ln ?c ?t ?r
WHERE {
  ?p :lastName ?ln ;
  :teaches ?c .
  ?c :title ?t ;
  :room ?r .
}
```

When we end a triple pattern with a ';' (instead of '.'), the next triple pattern uses the same subject (which therefore is not repeated).

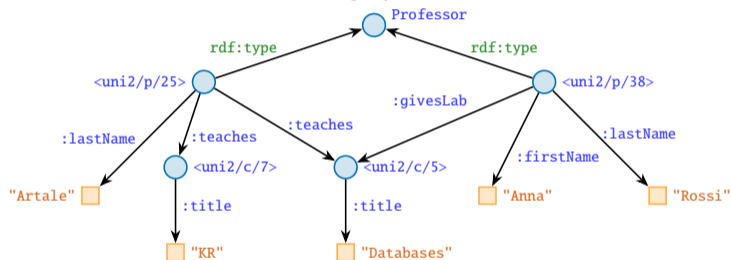
# Projecting out variables in a SPARQL query

A query may also return only a subset of the variables used in the BGP.

## Example: BGP with projection

```
SELECT ?ln ?t
WHERE {
  ?p :lastName ?ln .
  ?p :teaches ?c .
  ?c :title ?t .
}
```

When evaluated over the RDF graph



... the query returns:

ln	t
"Artale"	"Databases"
"Artale"	"KR"

# BGPs vs. conjunctive queries

We can write queries using the more compact and abstract syntax of conjunctive queries (CQs).

## Example: BGP

```
SELECT ?p ?ln ?c ?t
WHERE {
  ?p :lastName ?ln .
  ?p :teaches ?c .
  ?c :title ?t .
}
```

## vs. conjunctive query

$$q(p, ln, c, t) \leftarrow \text{lastName}(p, ln), \\ \text{teaches}(p, c), \\ \text{title}(c, t)$$

A **conjunctive query**  $q$  has the form  $q(\vec{x}) \leftarrow p_1(\vec{y}_1), \dots, p_k(\vec{y}_k)$  where

- $q(\vec{x})$  is called the **head** of  $q$ ,
- $p_1(\vec{y}_1), \dots, p_k(\vec{y}_k)$  is a conjunction of atoms called the **body** of  $q$ ,
- all variables  $\vec{x}$  in the head are among  $\vec{y}_1, \dots, \vec{y}_k$ , and
- the variables in  $\vec{y}_1, \dots, \vec{y}_k$  that are not among  $\vec{x}$  are **existentially quantified**.

# BGPs vs. conjunctive queries (cont.)

## Example: BGP with projection

```
SELECT ?ln ?t
WHERE {
  ?p :lastName ?ln .
  ?p :teaches ?c .
  ?c :title ?t .
}
```

## vs. conjunctive query with existential variables

$$q(ln, t) \leftarrow \text{lastName}(p, ln),$$
$$\text{teaches}(p, c),$$
$$\text{title}(c, t)$$

## But there is a difference in semantics when we have an ontology:

- In a SPARQL query, all variables, including those that are projected out, must match nodes of the RDF graph.
- In a conjunctive query, the existentially quantified variables can also match nodes that are existentially implied by the axioms of the ontology.

# BGPs vs. conjunctive queries – Example

Consider the KB  $\langle \mathcal{O}, \mathcal{A} \rangle$ , where the ontology is  $\mathcal{O} = \{ C \sqsubseteq \exists P \}$  and the RDF graph is  $\mathcal{A} = \{ C(a) \}$ . Consider further the following SPARQL BGP and the corresponding conjunctive query.

SPARQL query that has the form of a BGP

```
SELECT ?x WHERE { ?x rdf:type :C . ?x :P ?y . }
```

Conjunctive query

$$q(x) \leftarrow C(x), P(x, y)$$

Every model  $\mathcal{M}$  of  $\langle \mathcal{O}, \mathcal{A} \rangle$  contains the fact  $C(a)$  (recall that  $a^{\mathcal{M}} = a$ ), and since  $C \sqsubseteq \exists P \in \mathcal{O}$  also a fact  $P(a, o)$ , for some (existentially implied) object  $o$ . For example, the following are models of  $\langle \mathcal{O}, \mathcal{A} \rangle$ :

- $\mathcal{M}_1$ , with facts  $C(a)$  and  $P(a, o_1)$ ;
- $\mathcal{M}_2$ , with facts  $C(a)$  and  $P(a, o_2)$ ;
- $\mathcal{M}_3$ , with facts  $C(a)$ ,  $P(a, o_1)$ , and  $P(a, o_3)$ ;
- $\mathcal{M}_4$ , with facts  $C(a)$  and  $P(a, a)$ ;
- ...
- ...

Hence, for every model  $\mathcal{M}$  of  $\langle \mathcal{O}, \mathcal{A} \rangle$ , there is a homomorphism from the body of the conjunctive query  $q$  to  $\mathcal{M}$  that maps  $x$  to  $a$ . (Therefore, we have that  $a \in \text{cert}(q, \langle \mathcal{O}, \mathcal{A} \rangle)$ . – See later.)

Instead, even in the presence of an ontology, the SPARQL query must match on the RDF graph  $\mathcal{A}$  to produce an answer. Since  $\mathcal{A}$  contains only  $C(a)$ , the answer to the SPARQL query is empty.



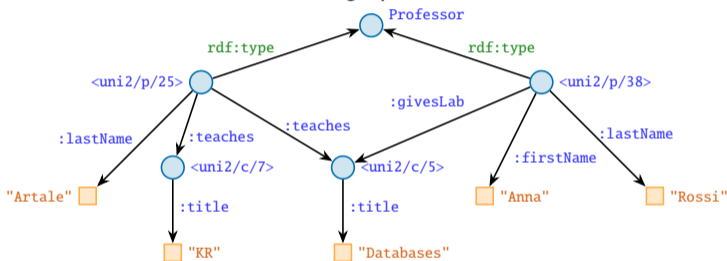
# Extending BGPs with OPTIONAL

We might want to add information when available, but **not reject** a solution **when some part of the query does not match**.

## Example: BGP with OPTIONAL

```
SELECT ?p ?fn ?ln
WHERE {
  ?p :lastName ?ln .
  OPTIONAL {
    ?p :firstName ?fn .
  }
}
```

When evaluated over the RDF graph



... the query returns:

<b>p</b>	<b>fn</b>	<b>ln</b>
<uni2/p/25>		"Artale"
<uni2/p/38>	"Anna"	"Rossi"

# SPARQL algebra

We have just seen the following features of the SPARQL algebra:

- Basic Graph Patterns
- OPTIONAL

The overall algebra has additional features:

- UNION
- ORDER BY, LIMIT, OFFSET
- FILTER conditions
- GROUP BY, to express aggregations and support aggregation operators
- MINUS, to remove possible solutions
- path expressions, corresponding to regular expressions

# Outline

- 1 Challenges in Data Access
- 2 Virtual Knowledge Graphs for Data Access (and Integration)**
  - Representing Data in RDF and RDFS
  - Representing Ontologies in OWL 2 QL
  - Query Language – SPARQL
  - Mapping an Ontology to a Relational Database**
  - Formalizing the VKG Framework
- 3 Optimizing Query Answering in VKGs
- 4 Designing VKG Mappings
- 5 The Ontop System
- 6 Conclusions

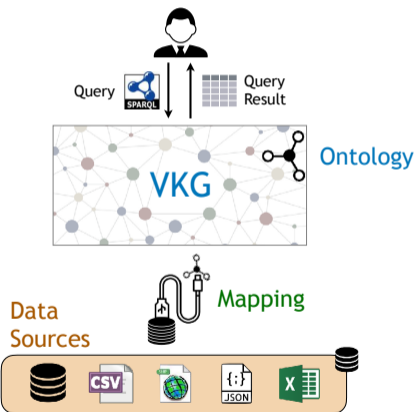
# Use of mappings

In the VKG framework, the **mapping** encodes how the **data in the sources** should be used to create the **Virtual Knowledge Graph**, which is formulated in the vocabulary of the **ontology**.

**VKG** defined from the **mapping** and the **data**.

- Queries are answered with respect to the **ontology** and the data of the **VKG**.
- The data of the **VKG** is not materialized (it is virtual!).
- Instead, the information in the **ontology** and the **mapping** is used to translate queries over the **ontology** into queries formulated over the **sources**.

Note: The graph is **always up to date** wrt the data sources.

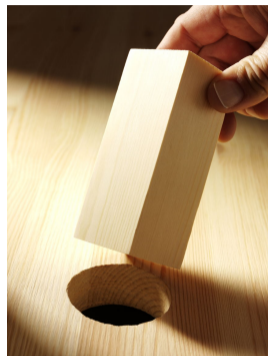


# Mismatch between data layer and ontology

## Impedance mismatch

- **Relational databases** store **values**.
- **Knowledge bases / ontologies** represent both **objects** and values.

We need to construct the ontology objects from the database values.



## Proposed solution

The specification of **how to construct the ontology objects** that populate the virtual knowledge graph from the database values **is embedded in the mapping** between the data sources and the ontology.

# VKG mapping

The **mapping** consists of a set of assertions of the form:

$$Q_{sql}(\vec{x}) \rightsquigarrow \Psi(\vec{t}, \vec{x})$$

- $Q_{sql}(\vec{x})$  is the **source query** expressed in SQL.
- $\Psi(\vec{t}, \vec{x})$  is the **target**, consisting of a set of **triple patterns** (i.e., atoms) that refer to the classes and properties of the ontology and make use of the answer variables  $\vec{x}$  of the SQL query.

To address the **impedance mismatch**, in the target query:

- we specify how to construct valid IRIs (that act as object identifiers), by concatenating database values and string constants;
- to refer to a database value, we use an answer variable of the source query;
- we call a term that constructs an IRI by referring to answer variables of the source query, an **IRI-template**.

# Triple patterns and IRI-templates

## Intuition behind the mapping

The **answers** returned by the **SQL query** in the source-part of the mapping are used to create, via the IRI-templates, the **objects** (and **values**) that populate the **classes / properties** in the target part.

More precisely:

- Each **triple pattern** in the target part has one of the forms:

$\text{iri}_1(\vec{x}_1)$  **rdf:type**  $C$

where  $C$  is a class of the ontology, or

$\text{iri}_1(\vec{x}_1)$  **prop**  $\text{iri}_2(\vec{x}_2)$

where **prop** is a (data or object) property of the ontology.

- For each answer tuple  $\vec{a}$  returned by the **source query**  $Q_{sql}(\vec{x})$  (when evaluated over the **database**), the **iri-template**  $\text{iri}_i(\vec{x}_i)$  generates an **object / value**  $\text{iri}_i(\vec{a}_i)$  of the **VKG**.
- Such objects / values are then used to populate the classes and properties of the ontology according to what specified in the **target** part of the mapping.

In this way we provide a solution to the **impedance mismatch** problem.

# A concrete mapping language

We describe the concrete mapping language adopted by the *Ontop* system.

In the *Ontop* mapping language, each mapping assertion is made up of three parts:

- A **mapping identifier**, which is convenient to refer to a specific mapping.
- The **source part**, which is a regular SQL query over the data source(s).
- The **target part**, which is a set of triple patterns that make use of IRI-templates.  
In the target part, the answer variables of the source part are enclosed in `{...}`.

## Mapping $m_1$

- Mapping identifier: `m1`
- Source part:
 

```
SELECT mcode, mtitle
FROM MOVIE
WHERE type = "m"
```
- Target part:
 

```
:m/{mcode} rdf:type :Movie .
:m/{mcode} :title {mtitle} .
```

## Mapping $m_2$

- Mapping identifier: `m2`
- Source part:
 

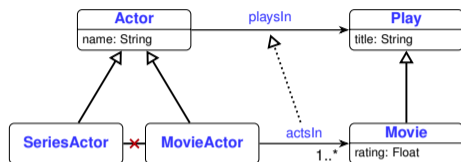
```
SELECT M.mcode, A.icode
FROM MOVIE M, ACTOR A
WHERE M.mcode = A.pcode
AND M.type = "m"
```
- Target part:
 

```
:a/{icode} :actsIn :m/{mcode} .
```



# Mapping language – Example

## Ontology $\mathcal{O}$ :



## Mapping $\mathcal{M}$ :

$m_1$ : **SELECT** mcode, mtitle **FROM** MOVIE

**WHERE** type = "m"

$\rightsquigarrow$  :m/{mcode} **rdf:type** :Movie .  
 :m/{mcode} :title {mtitle} .

$m_2$ : **SELECT** M.mcode, A.acode **FROM** MOVIE M, ACTOR A

**WHERE** M.mcode = A.pcode **AND** M.type = "m"

$\rightsquigarrow$  :a/{acode} :actsIn :m/{mcode} .

## Database $\mathcal{D}$ :

MOVIE				
mcode	mtitle	myear	type	...
5118	The Matrix	1999	m	...
8234	Altered Carbon	2018	s	...
2281	Blade Runner	1982	m	...

ACTOR			
pcode	acode	aname	...
5118	438	K. Reeves	...
5118	572	C.A. Moss	...
2281	271	H. Ford	...

The mapping  $\mathcal{M}$  applied to database  $\mathcal{D}$  generates the virtual knowledge graph  $\mathcal{M}(\mathcal{D})$ :

:m/5118 **rdf:type** :Movie .    :m/5118 :title "The Matrix" .  
 :m/2281 **rdf:type** :Movie .    :m/2281 :title "Blade Runner" .  
 :a/438 :actsIn :m/5118 .    :a/572 :actsIn :m/5118 .    :a/271 :actsIn :m/2281 .

# Standard mapping languages

Several proposals for concrete languages to map a relational DB to an ontology:

- They assume that the ontology is populated in terms of RDF triples.
- Some template mechanism is used to specify the triples to instantiate.

Examples: D2RQ<sup>1</sup>, SML<sup>2</sup>, Ontop<sup>3</sup>

## R2RML

- Most popular RDB to RDF mapping language
- W3C Recommendation 27 Sep. 2012, <http://www.w3.org/TR/r2rml/>
- R2RML mappings are themselves expressed as RDF graphs and written in Turtle syntax.

---

<sup>1</sup><http://d2rq.org/d2rq-language>

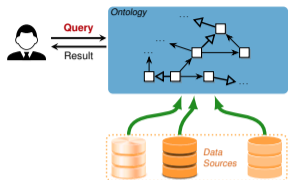
<sup>2</sup>[http://sparqlify.org/wiki/Sparqlification\\_mapping\\_language](http://sparqlify.org/wiki/Sparqlification_mapping_language)

<sup>3</sup>[https://github.com/ontop/ontop/wiki/ontopOBDAModel#Mapping\\_axioms](https://github.com/ontop/ontop/wiki/ontopOBDAModel#Mapping_axioms)

# Outline

- 1 Challenges in Data Access
- 2 Virtual Knowledge Graphs for Data Access (and Integration)**
  - Representing Data in RDF and RDFS
  - Representing Ontologies in OWL 2 QL
  - Query Language – SPARQL
  - Mapping an Ontology to a Relational Database
  - Formalizing the VKG Framework**
- 3 Optimizing Query Answering in VKGs
- 4 Designing VKG Mappings
- 5 The Ontop System
- 6 Conclusions

# VKGs: Formalization



To formalize VKGs, we distinguish between the intensional and the extensional level information.

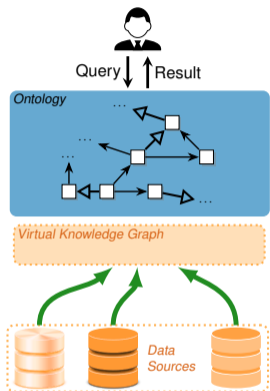
A **VKG specification** is a triple  $\mathcal{P} = \langle O, \mathcal{M}, \mathcal{S} \rangle$ , where:

- $O$  is an **ontology** (expressed in OWL 2 QL),
- $\mathcal{S}$  is a (possibly federated) **relational database schema** for the data sources, possibly with integrity constraints,
- $\mathcal{M}$  is a set of (R2RML) **mapping assertions** between  $O$  and  $\mathcal{S}$ .

A **VKG instance** is a pair  $\langle \mathcal{P}, \mathcal{D} \rangle$ , where

- $\mathcal{P} = \langle O, \mathcal{M}, \mathcal{S} \rangle$  is a VKG specification, and
- $\mathcal{D}$  is a (possibly federated) relational database compliant with  $\mathcal{S}$ .

# Semantics of VKGs



Remember:

- The mapping  $\mathcal{M}$  generates from the data  $\mathcal{D}$  in the sources a **virtual knowledge graph**  $\mathcal{V} = \mathcal{M}(\mathcal{D})$ .
- The set of constants that can appear in  $\mathcal{V}$  consists of:
  - values obtained directly from the database, and
  - IRIs, which are constructed by applying the `iri` function to string constants and database values.

We use  $\mathcal{C}_{\mathcal{V}}$ , i.e.,  $\mathcal{C}_{\mathcal{M}(\mathcal{D})}$ , to denote such set of constants.

A first-order interpretation  $\mathcal{I}$  of the ontology predicates and the constants in  $\mathcal{C}_{\mathcal{M}(\mathcal{D})}$  is a **model** of  $\langle \mathcal{P}, \mathcal{D} \rangle$  if

- it satisfies all axioms in  $\mathcal{O}$ , and
- contains all facts in  $\mathcal{M}(\mathcal{D})$ , i.e., retrieved from  $\mathcal{D}$  through  $\mathcal{M}$ .

Note:

- In general,  $\langle \mathcal{P}, \mathcal{D} \rangle$  has infinitely **many models**, and some of these might be infinite.
- However, for query answering, we do not need to compute such models.

# Query answering in VKGs – Certain answers

In VKGs, we want to answer queries formulated over the ontology, by using the data provided by the data sources through the mapping.

Consider our formalization of VKGs and a VKG instance  $\mathcal{J}$ .

## Certain answers $\text{cert}(q, \mathcal{J})$ – Intuition

Given a VKG instance  $\mathcal{J}$  and a query  $q$  over  $\mathcal{J}$ , the certain answers  $\text{cert}(q, \mathcal{J})$  to  $q$  over  $\mathcal{J}$  are those answers to  $q$  that hold in **every model** of  $\mathcal{J}$ .

## Certain answers $\text{cert}(q, \mathcal{J})$ – Formal definition

Given a VKG instance  $\mathcal{J} = \langle \mathcal{P}, \mathcal{D} \rangle$  and a query  $q$  over  $\mathcal{J}$ , a tuple  $\vec{c}$  of constants in  $\mathcal{C}_{\mathcal{M}(\mathcal{D})}$  is a **certain answer** to  $q$  over  $\mathcal{J}$ , i.e.,  $\vec{c} \in \text{cert}(q, \mathcal{J})$ , if for **every model**  $\mathcal{I}$  of  $\mathcal{J}$  we have that  $\vec{c} \in q(\mathcal{I})$ .

*Note:* Each certain answer  $\vec{c}$  is a tuple of constants in  $\mathcal{C}_{\mathcal{M}(\mathcal{D})}$ , but when we evaluate  $q$  over an interpretation  $\mathcal{I}$ , it returns tuples of elements of  $\Delta^{\mathcal{I}}$ . Therefore, we should actually require that  $\vec{c}^{\mathcal{I}} \in q(\mathcal{I})$ , and not that  $\vec{c} \in q(\mathcal{I})$ . However, due to the standard names assumption, we have that  $\vec{c}^{\mathcal{I}} = \vec{c}$ , so the two conditions are equivalent.

# First-order rewritability

To make computing certain answers viable in practice, the VKG setting relies on reducing it to evaluating SQL (i.e., first-order logic) queries over the data.

Consider a VKG specification  $\mathcal{P} = \langle \mathcal{O}, \mathcal{M}, \mathcal{S} \rangle$ .

## First-order rewritability

A query  $r(\vec{x})$  is a **first-order rewriting** of a query  $q(\vec{x})$  with respect to  $\mathcal{P}$  if, for every source DB  $\mathcal{D}$ , certain answers to  $q(\vec{x})$  over  $\langle \mathcal{P}, \mathcal{D} \rangle =$  answers to  $r(\vec{x})$  over  $\mathcal{D}$ .

For OWL 2 QL ontologies and R2RML mappings,  
(core) SPARQL queries are first-order rewritable.

In other words, **in VKGs, we can compute the certain answers to a SPARQL query by evaluating over the sources its rewriting, which is a SQL query.**

# Computational complexity of query answering

**Theorem** [C., De Giacomo, et al. 2007; Poggi et al. 2008; Artale et al. 2009]

For OWL 2 QL (or *DL-Lite*) VKG instances  $\langle \mathcal{P}, \mathcal{D} \rangle$ , with  $\mathcal{P} = \langle \mathcal{O}, \mathcal{M}, \mathcal{S} \rangle$ , **query answering** for UCQs / SPARQL queries is:

- 1 Very efficiently tractable, i.e., in  $AC^0$ , in the size of the **database**  $\mathcal{D}$ .
- 2 Efficiently tractable, i.e., in **LOGSPACE**, in the size of the **ontology**  $\mathcal{O}$  and the **mapping**  $\mathcal{M}$ .
- 3 Exponential, more precisely **NP-complete**, in the size of the **query**.

In **theory this is not bad**, since this is the complexity of evaluating CQs in relational DBs.

**Note:** The  $AC^0$  result is a consequence of the fact that query answering in such a setting can be reduced to evaluating a SQL query over the relational database  $\mathcal{D}$ .

**Can we go beyond *DL-Lite* and maintain the same complexity results?**

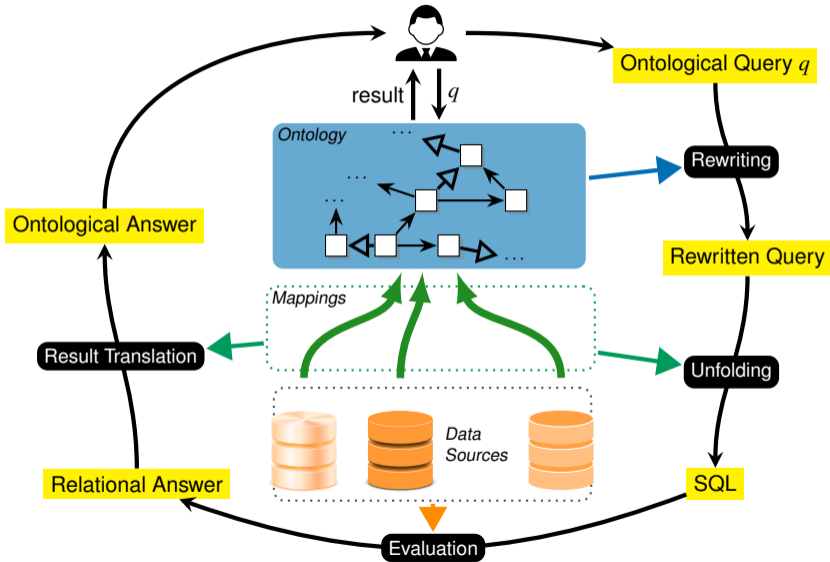
Essentially no! By adding essentially any additional constructs of OWL, we lose first-order rewritability and hence these nice computational properties. [C., De Giacomo, et al. 2006, 2013]



# Outline

- 1 Challenges in Data Access
- 2 Virtual Knowledge Graphs for Data Access (and Integration)
- 3 Optimizing Query Answering in VKGs**
  - Query rewriting wrt an OWL 2 QL ontology
  - Query unfolding wrt a mapping
  - Mapping saturation
  - Optimization of query reformulation
- 4 Designing VKG Mappings
- 5 The Ontop System
- 6 Conclusions

# Query answering via query reformulation – Conceptual framework



# Query answering via query reformulation – Optimizations needed

The above conceptual framework is realized as follows.

Computing certain answers to a SPARQL query  $q$  over a VKG instance  $\langle \mathcal{P}, \mathcal{D} \rangle$ , with  $\mathcal{P} = \langle \mathcal{O}, \mathcal{M}, \mathcal{S} \rangle$ :

- 1 Compute the perfect rewriting of  $q$  w.r.t.  $\mathcal{O}$ .
- 2 Unfold the perfect rewriting w.r.t. the mapping  $\mathcal{M}$ .
- 3 **Optimize** the unfolded query, using database constraints.
- 4 Evaluate the resulting SQL query over  $\mathcal{D}$ .

Steps 1–3 are collectively called **query reformulation**.

We analyze now these steps more in detail.

# Outline

- 1 Challenges in Data Access
- 2 Virtual Knowledge Graphs for Data Access (and Integration)
- 3 Optimizing Query Answering in VKGs**
  - Query rewriting wrt an OWL 2 QL ontology
  - Query unfolding wrt a mapping
  - Mapping saturation
  - Optimization of query reformulation
- 4 Designing VKG Mappings
- 5 The Ontop System
- 6 Conclusions

# Rewriting step

The rewriting Step ① deals with the knowledge encoded by the axioms of the ontology:

- hierarchies of classes and of properties;
- objects that are existentially implied by such axioms: existential reasoning.

We illustrate the need for dealing with these two aspects with two examples.

# Dealing with hierarchies

Suppose that every graduate student is a student, i.e.,

`GraduateStudent`  $\sqsubseteq$  `Student`

and john is a graduate student: `GraduateStudent(john)`.

What is the answer to the following query, asking for all students?

`q(x) ← Student(x)`

In SPARQL: `SELECT ?x WHERE { ?x a Student . }`

The answer should be `john`, since being a graduate student, he is also a student.

# Dealing with existential reasoning

Suppose that every student is supervised by some professor, i.e.,

$$\text{Student} \sqsubseteq \exists \text{isSupervisedBy}.\text{Professor}$$

and john is a student:  $\text{Student}(\text{john})$ .

What is the answer to the following query, asking for all individuals supervised by some professor?

$$q(x) \leftarrow \text{isSupervisedBy}(x, y), \text{Professor}(y)$$

In SPARQL:  $\text{SELECT ?x WHERE \{ ?x isSupervisedBy [ a Professor ] . \}}$

The answer should be  $\text{john}$ , even though we don't know who is John's supervisor (under existential reasoning).

# The query rewriting algorithm

The **query rewriting** algorithm takes into account hierarchies and existential reasoning, by “compiling” the axioms of the ontology into the query.

## Example

Consider the ontology axioms:

$$\text{Student} \sqsubseteq \exists \text{isSupervisedBy}.\text{Professor}$$
$$\text{GraduateStudent} \sqsubseteq \text{Student}$$

Using these axioms, the rewriting algorithm rewrites the query

$$q(x) \leftarrow \text{isSupervisedBy}(x, y), \text{Professor}(y)$$

into a union of conjunctive queries (or a SPARQL union query):

$$q(x) \leftarrow \text{isSupervisedBy}(x, y), \text{Professor}(y)$$
$$q(x) \leftarrow \text{Student}(x)$$
$$q(x) \leftarrow \text{GraduateStudent}(x)$$

Therefore, over the data  $\text{Student}(\text{john})$ , the rewritten query returns  $\text{john}$  as an answer.

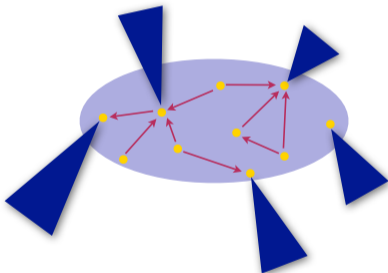
**Note:** In *Ontop*, existential reasoning needs to be switched on explicitly, since it affects performance. 



# Query rewriting and canonical model

## Canonical model

Every consistent *DL-Lite* / *OWL 2 QL* KB  $\mathcal{K} = \langle \mathcal{O}, \mathcal{A} \rangle$  has a **canonical model**  $\mathcal{I}_{\mathcal{K}}$ , which **gives the right answers to all CQs**, i.e.,  $\text{cert}(q, \mathcal{K}) = \text{ans}(q, \mathcal{I}_{\mathcal{K}})$



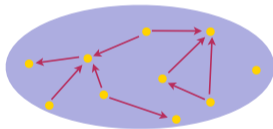
- The core part can be handled by  **saturating the mapping** .
- The anonymous part can be handled by  **tree-witness rewriting** .

# Query rewriting and canonical model

## Canonical model

Every consistent *DL-Lite* / *OWL 2 QL* KB  $\mathcal{K} = \langle \mathcal{O}, \mathcal{A} \rangle$  has a **canonical model**  $\mathcal{I}_{\mathcal{K}}$ , which **gives the right answers to all CQs**, i.e.,  $\text{cert}(q, \mathcal{K}) = \text{ans}(q, \mathcal{I}_{\mathcal{K}})$

**Core**  
individuals  
from  $\mathcal{A}$

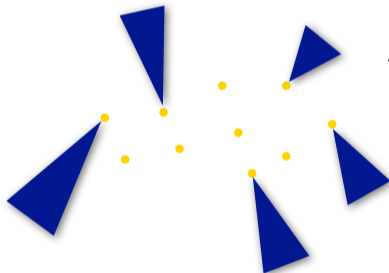


- The core part can be handled by  **saturating the mapping** .
- The anonymous part can be handled by  **tree-witness rewriting** .

# Query rewriting and canonical model

## Canonical model

Every consistent *DL-Lite* / *OWL 2 QL* KB  $\mathcal{K} = \langle \mathcal{O}, \mathcal{A} \rangle$  has a **canonical model**  $\mathcal{I}_{\mathcal{K}}$ , which **gives the right answers to all CQs**, i.e.,  $\text{cert}(q, \mathcal{K}) = \text{ans}(q, \mathcal{I}_{\mathcal{K}})$



### Anonymous part

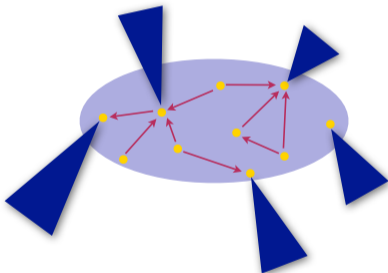
*trees rooted at individuals,  
using unnamed objects*

- The core part can be handled by  **saturating the mapping** .
- The anonymous part can be handled by  **tree-witness rewriting** .

# Query rewriting and canonical model

## Canonical model

Every consistent *DL-Lite* / *OWL 2 QL* KB  $\mathcal{K} = \langle \mathcal{O}, \mathcal{A} \rangle$  has a **canonical model**  $\mathcal{I}_{\mathcal{K}}$ , which **gives the right answers to all CQs**, i.e.,  $\text{cert}(q, \mathcal{K}) = \text{ans}(q, \mathcal{I}_{\mathcal{K}})$

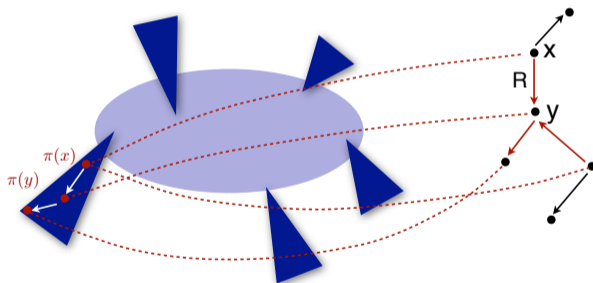


- The core part can be handled by  **saturating the mapping** .
- The anonymous part can be handled by  **tree-witness rewriting** .

# Query rewriting and canonical model

## Canonical model

Every consistent *DL-Lite* / *OWL 2 QL* KB  $\mathcal{K} = \langle \mathcal{O}, \mathcal{A} \rangle$  has a **canonical model**  $\mathcal{I}_{\mathcal{K}}$ , which **gives the right answers to all CQs**, i.e.,  $\text{cert}(q, \mathcal{K}) = \text{ans}(q, \mathcal{I}_{\mathcal{K}})$



- The core part can be handled by  **saturating the mapping** .
- The anonymous part can be handled by  **tree-witness rewriting** .

# The *PerfectRef* algorithm for query rewriting

We do not describe here the tree-witness rewriting algorithm, which is rather involved.

Instead, we describe *PerfectRef*, a simple query rewriting algorithm that maintains a set of queries and applies over them two types of transformations:

- rewriting steps that involve inclusion assertions of the ontology, and
- unification of query atoms.

These transformations are applied repeatedly until saturation, i.e., until the set of queries does not change anymore.

Given as input a (core) SPARQL query  $q$ , *PerfectRef* computes its **perfect rewriting**, which is still a SPARQL query (involving UNION).

*Note:* Disjointness assertions play a role in ontology satisfiability, but can be ignored during query rewriting. (This is called **separability**.)

# Query rewriting step: Basic idea

Intuition: an **inclusion assertion** corresponds to a **logic programming rule**.

## Basic rewriting step:

When an **atom in the query** unifies with the **head** of the rule, generate a new query by substituting the atom with the **body** of the rule.

We say that the inclusion assertion **applies to** the atom.

## Example

The inclusion assertion  $\text{Professor} \sqsubseteq \text{Teacher}$   
 corresponds to the logic programming rule  $\text{Teacher}(z) \leftarrow \text{Professor}(z)$ .

Consider the query  $q(x) \leftarrow \text{Teacher}(x)$ .

By applying the inclusion assertion to the atom  $\text{Teacher}(x)$ , we generate:

$q(x) \leftarrow \text{Professor}(x)$ .

This query is added to the input query, and contributes to the perfect rewriting.

# Query rewriting (cont'd)

## Example

Consider the query  $q(x) \leftarrow \text{teaches}(x, y), \text{Course}(y)$

and the inclusion assertion  $\exists \text{teaches}^- \sqsubseteq \text{Course}$

as a logic programming rule:  $\text{Course}(z_2) \leftarrow \text{teaches}(z_1, z_2)$ .

The inclusion applies to  $\text{Course}(y)$ , and we add to the rewriting the query

$$q(x) \leftarrow \text{teaches}(x, y), \text{teaches}(z_1, y).$$

## Example

Consider now the query  $q(x) \leftarrow \text{teaches}(x, y)$

and the inclusion assertion  $\text{Professor} \sqsubseteq \exists \text{teaches}$

as a logic programming rule:  $\text{teaches}(z, f(z)) \leftarrow \text{Professor}(z)$ .

The inclusion applies to  $\text{teaches}(x, y)$ , and we add to the rewriting the query

$$q(x) \leftarrow \text{Professor}(x).$$



# Query rewriting – Constants

## Example

Conversely, for the query  $q(x) \leftarrow \text{teaches}(x, \text{"databases"})$

and the same inclusion assertion as before  $\text{Professor} \sqsubseteq \exists \text{teaches}$

as a logic programming rule:  $\text{teaches}(z, f(z)) \leftarrow \text{Professor}(z)$

$\text{teaches}(x, \text{"databases"})$  does not unify with  $\text{teaches}(z, f(z))$ , since the **skolem term**  $f(z)$  in the head of the rule **does not unify** with the constant  $\text{"databases"}$ .

Remember: We adopt the **unique name assumption**.

We say that the **inclusion** does **not** apply to the atom  $\text{teaches}(x, \text{"databases"})$ .

## Example

The same holds for the following query, where  $y$  is **distinguished**, since unifying  $f(z)$  with  $y$  would correspond to returning a skolem term as answer to the query:

$$q(x, y) \leftarrow \text{teaches}(x, y).$$

# Query rewriting – Join variables

An analogous behavior to the one with constants and with distinguished variables holds when the atom contains **join variables** that would have to be unified with skolem terms.

## Example

Consider the query  $q(x) \leftarrow \text{teaches}(x, y), \text{Course}(y)$

and the inclusion assertion  $\text{Professor} \sqsubseteq \exists \text{teaches}$

as a logic programming rule:  $\text{teaches}(z, f(z)) \leftarrow \text{Professor}(z)$ .

The **inclusion assertion** above does **not** apply to the atom  $\text{teaches}(x, y)$ .

# Query rewriting – Reduce step

## Example

Consider now the query  $q(x) \leftarrow \text{teaches}(x, y), \text{teaches}(z, y)$

and the inclusion assertion  $\text{Professor} \sqsubseteq \exists \text{teaches}$

as a logic rule:  $\text{teaches}(z, f(z)) \leftarrow \text{Professor}(z)$ .

This inclusion assertion does not apply to  $\text{teaches}(x, y)$  or  $\text{teaches}(z, y)$ , since  $y$  is in join, and we would again introduce the skolem term in the rewritten query.

## Example

However, we can transform the above query by **unifying** the atoms  $\text{teaches}(x, y)$  and  $\text{teaches}(z, y)$ . This rewriting step is called **reduce**, and produces the query

$$q(x) \leftarrow \text{teaches}(x, y).$$

Now, we can apply the inclusion above, and add to the rewriting the query

$$q(x) \leftarrow \text{Professor}(x).$$

# Query rewriting – Summary

To compute the perfect rewriting of a query  $q$  with respect to an ontology  $O$ , start from  $q$ , iteratively get a CQ  $q'$  to be processed, and do one of the following:

- Apply to some atom of  $q'$  an inclusion assertion in  $O$  as follows:

$$\begin{array}{llll}
 C_1 \sqsubseteq C_2 & \dots, C_2(x), \dots & \rightsquigarrow & \dots, C_1(x), \dots \\
 \exists P \sqsubseteq C & \dots, C(x), \dots & \rightsquigarrow & \dots, P(x, \_), \dots \\
 \exists P^- \sqsubseteq C & \dots, C(x), \dots & \rightsquigarrow & \dots, P(\_, x), \dots \\
 C \sqsubseteq \exists P & \dots, P(x, \_), \dots & \rightsquigarrow & \dots, C(x), \dots \\
 C \sqsubseteq \exists P^- & \dots, P(\_, x), \dots & \rightsquigarrow & \dots, C(x), \dots \\
 \exists P_1 \sqsubseteq \exists P_2 & \dots, P_2(x, \_), \dots & \rightsquigarrow & \dots, P_1(x, \_), \dots \\
 P_1 \sqsubseteq P_2 & \dots, P_2(x, y), \dots & \rightsquigarrow & \dots, P_1(x, y), \dots \\
 P_1 \sqsubseteq P_2^- & \dots, P_2(x, y), \dots & \rightsquigarrow & \dots, P_1(y, x), \dots
 \end{array}$$

('\_' denotes a variable that appears only once)

- Choose two atoms of  $q'$  that unify, and apply the unifier to  $q'$ .

After each rewriting/unification step, the obtained query is added to the queries still to be processed.

**Note:** Unifying atoms can make rules applicable that were not so before, and is required for completeness of the method [C., De Giacomo, et al. 2007].

The UCQ resulting from this process is the **perfect rewriting**  $q_r$  of  $q$  w.r.t. the ontology  $O$ .

# Query rewriting algorithm

**Algorithm** *PerfectRef*( $Q, O_p$ )

**Input:** union of conjunctive queries  $Q$ , set  $O_p$  of *DL-Lite* / *OWL 2 QL* positive inclusion assertions

**Output:** union of conjunctive queries  $PR$

$PR := Q$ ;

**repeat**

$PR' := PR$ ;

**for each**  $q \in PR'$  **do**

**for each**  $g$  in  $q$  **do**

**for each** inclusion assertion  $I$  in  $O_p$  **do**

**if**  $I$  is applicable to  $g$  **then**  $PR := PR \cup \{ \text{ApplyPI}(q, g, I) \}$ ;

**for each**  $g_1, g_2$  in  $q$  **do**

**if**  $g_1$  and  $g_2$  unify **then**  $PR := PR \cup \{ \tau(\text{Reduce}(q, g_1, g_2)) \}$ ;

**until**  $PR' = PR$ ;

**return**  $PR$

*Observations:*

- Termination follows from having only finitely many different rewritings.
- Disjointness assertions and functionalities do not play any role in the rewriting of the query.

# Query answering in *DL-Lite* – Example

Ontology:

Professor  $\sqsubseteq$  Teacher

Teacher  $\sqsubseteq \exists \text{teaches}$

$\exists \text{teaches}^- \sqsubseteq \text{Course}$

Corresponding rules:

Teacher( $z$ )  $\leftarrow$  Professor( $z$ )

teaches( $z, f(z)$ )  $\leftarrow$  Teacher( $z$ )

Course( $z$ )  $\leftarrow$  teaches( $w, z$ )

Query:  $q(x) \leftarrow \text{teaches}(x, y), \text{Course}(y)$

Perfect rewriting:  $q(x) \leftarrow \text{teaches}(x, y), \text{Course}(y)$

$q(x) \leftarrow \text{teaches}(x, y), \text{teaches}(\_, y)$

$q(x) \leftarrow \text{teaches}(x, \_)$

$q(x) \leftarrow \text{Teacher}(x)$

$q(x) \leftarrow \text{Professor}(x)$

ABox: teaches(jim, databases)      Professor(jim)  
 teaches(julia, security)      Teacher(nicole)

Evaluating the perfect rewriting over the ABox (seen as a DB) produces as answer  
**{jim, julia, nicole}**.

# Query answering in *DL-Lite* – An interesting example

TBox:  $\text{Person} \sqsubseteq \exists \text{hasFather}$   
 $\exists \text{hasFather}^- \sqsubseteq \text{Person}$

ABox:  $\text{Person}(\text{john})$

Query:  $q(x) \leftarrow \text{Person}(x), \text{hasFather}(x, y_1), \text{hasFather}(y_1, y_2), \text{hasFather}(y_2, y_3)$

$q(x) \leftarrow \text{Person}(x), \text{hasFather}(x, y_1), \text{hasFather}(y_1, y_2), \text{hasFather}(y_2, -)$

$\Downarrow$  **Apply**  $\text{Person} \sqsubseteq \exists \text{hasFather}$  to the atom  $\text{hasFather}(y_2, -)$

$q(x) \leftarrow \text{Person}(x), \text{hasFather}(x, y_1), \text{hasFather}(y_1, y_2), \text{Person}(y_2)$

$\Downarrow$  **Apply**  $\exists \text{hasFather}^- \sqsubseteq \text{Person}$  to the atom  $\text{Person}(y_2)$

$q(x) \leftarrow \text{Person}(x), \text{hasFather}(x, y_1), \text{hasFather}(y_1, y_2), \text{hasFather}(-, y_2)$

$\Downarrow$  **Unify** atoms  $\text{hasFather}(y_1, y_2)$  and  $\text{hasFather}(-, y_2)$

$q(x) \leftarrow \text{Person}(x), \text{hasFather}(x, y_1), \text{hasFather}(y_1, y_2)$

$\Downarrow$

...

$q(x) \leftarrow \text{Person}(x), \text{hasFather}(x, -)$

$\Downarrow$  **Apply**  $\text{Person} \sqsubseteq \exists \text{hasFather}$  to the atom  $\text{hasFather}(x, -)$

$q(x) \leftarrow \text{Person}(x)$

# Exponential blowup in the rewriting

Even with a flat hierarchy of classes in which a single inclusion assertion can be applied to each atom of a query  $q$ , the rewriting may contain an number of CQs that is exponential in the length of  $q$ .

Consider a query:  $q(x) \leftarrow C_1^1(x), C_2^1(x), \dots, C_n^1(x)$

and the ontology:  $O = \{ C_1^2 \sqsubseteq C_1^1, C_2^2 \sqsubseteq C_2^1, \dots, C_n^2 \sqsubseteq C_n^1 \}$

Each atom  $C_i^1(x)$  in  $q$  can either stay as is, or we can apply to it the inclusion assertion  $C_i^2 \sqsubseteq C_i^1$ , and generate a new CQ in which  $C_i^1(x)$  is replaced by  $C_i^2(x)$ .

Hence, in the rewriting we have **one CQ**

$$q(x) \leftarrow C_1^{j_1}(x), C_2^{j_2}(x), \dots, C_n^{j_n}(x)$$

for each possible combination of  $j_i \in \{1, 2\}$ , for  $i \in \{1, \dots, n\}$ .

Hence, the rewriting of  $q$  with respect to  $O$  contains  **$2^n$  CQs**.



# Outline

- 1 Challenges in Data Access
- 2 Virtual Knowledge Graphs for Data Access (and Integration)
- 3 Optimizing Query Answering in VKGs**
  - Query rewriting wrt an OWL 2 QL ontology
  - Query unfolding wrt a mapping**
  - Mapping saturation
  - Optimization of query reformulation
- 4 Designing VKG Mappings
- 5 The Ontop System
- 6 Conclusions

# Query unfolding

We consider now Step ② of reformulation, i.e., the unfolding w.r.t. the mapping  $\mathcal{M}$ .

In principle, we have two approaches to exploit the mapping:

- bottom-up approach: simpler, but typically less efficient
- top-down approach: more sophisticated, but also more efficient

Both approaches require to first **split** the set of atoms in the target queries of the mapping assertions into the constituent atoms.

*Note:* In the following, to make notation more compact, we represent an IRI-template of the form

**$:xxx/\{v_1\}/\{v_2\}/\dots/\{v_n\}$**

more compactly as

**$xxx(v_1, \dots, v_n)$** .

# Splitting of mappings

A mapping assertion  $\Phi \rightsquigarrow \Psi$ , where the target query  $\Psi$  is constituted by the atoms  $X_1, \dots, X_k$ , can be split into  $k$  mapping assertions:

$$\Phi \rightsquigarrow X_1 \quad \dots \quad \Phi \rightsquigarrow X_k$$

This is possible, since  $\Psi$  does not contain non-distinguished variables.

## Example

$m_1$ : `SELECT pcode, acode, aname FROM ACTOR`  $\rightsquigarrow$  `Play(pl(pcode)),`  
`Actor(act(acode)),`  
`name(act(acode), aname),`  
`actsIn(act(acode), pl(pcode))`

is split into

$m_1^1$ : `SELECT pcode, acode, aname FROM ACTOR`  $\rightsquigarrow$  `Play(pl(pcode))`  
 $m_1^2$ : `SELECT pcode, acode, aname FROM ACTOR`  $\rightsquigarrow$  `Actor(act(acode))`  
 $m_1^3$ : `SELECT pcode, acode, aname FROM ACTOR`  $\rightsquigarrow$  `name(act(acode), aname)`  
 $m_1^4$ : `SELECT pcode, acode, aname FROM ACTOR`  $\rightsquigarrow$  `actsIn(act(acode), pl(pcode))`

# Bottom-up approach to deal with mappings: Materialization

Consists in a straightforward application of the mappings to the data:

- 1 Propagate the data from  $\mathcal{D}$  through  $\mathcal{M}$ , **materializing** the RDF graph  $\mathcal{V} = \mathcal{M}(\mathcal{D})$  (the constants in such an RDF graph are values and object terms obtained from the database values).
- 2 Apply to  $\mathcal{V}$  and to the ontology  $\mathcal{O}$ , the query answering algorithm (based on query rewriting) developed for *DL-Lite* / OWL 2 QL.

This approach has several drawbacks:

- The technique is no more  $AC^0$  in the size of the data, since the RDF graph  $\mathcal{V}$  to materialize is in general polynomial in the size of the data.
- $\mathcal{V}$  may be very large, and thus it may be infeasible to actually materialize it.
- Freshness of  $\mathcal{V}$  with respect to the underlying data source(s) may be an issue, and one would need to propagate source updates (cf. Data Warehousing).

# Top-down approach to deal with mappings: Unfolding

The top-down approach is realized by computing from the (rewritten) query  $q_r$  a new query  $q_{\text{unf}}$ , by **unfolding**  $q_r$  using (the split version of) the mappings  $\mathcal{M}$ .

Consider the mapping assertions  $\Phi_i \rightsquigarrow \Psi_i$ .

- Essentially, each atom in  $q_r$  that unifies with an atom in some  $\Psi_i$  is substituted with the corresponding query  $\Phi_i$  over the database.
- The unfolded query  $q_{\text{unf}}$  is such that for each database  $\mathcal{D}$  we have that:

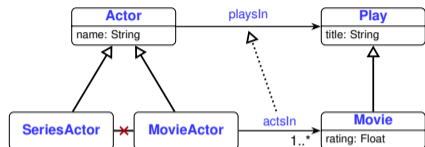
$$q_{\text{unf}}(\mathcal{D}) = \text{Eval}_{\text{CWA}}(q_r, \mathcal{M}(\mathcal{D})).$$

# Unfolding

To unfold a query  $q_r$  with respect to a set  $\mathcal{M}$  of mapping assertions:

- 1 For each non-split mapping assertion  $\Phi_i(\vec{x}) \rightsquigarrow \Psi_i(\vec{t}, \vec{y})$ :
  - 1 Introduce a **view symbol**  $Aux_i$  of arity equal to that of  $\Phi_i$ .
  - 2 Add a **view definition**  $Aux_i(\vec{x}) \leftarrow \Phi_i(\vec{x})$ .
  
- 2 For each split version  $\Phi_i(\vec{x}) \rightsquigarrow X_i^j(\vec{t}, \vec{y})$  of a mapping assertion, introduce a **clause**  $X_i^j(\vec{t}, \vec{y}) \leftarrow Aux_i(\vec{x})$ .
  
- 3 Obtain from  $q_r$  in all possible ways queries  $q_{aux}$  defined over the view symbols  $Aux_i$  as follows:
  - 1 Find a most general unifier  $\vartheta$  that unifies each atom  $X(\vec{z})$  in the body of  $q_r$  with the head of a clause  $X(\vec{t}, \vec{y}) \leftarrow Aux_i(\vec{x})$ .
  - 2 Substitute each atom  $X(\vec{z})$  with  $\vartheta(Aux_i(\vec{x}))$ , i.e., with the body the unified clause to which the unifier  $\vartheta$  is applied.
  
- 4 The unfolded query  $q_{unf}$  is the **union** of all queries  $q_{aux}$ , together with the view definitions for the predicates  $Aux_i$  appearing in  $q_{aux}$ .

# Unfolding – Example



$m_1$ : `SELECT pcode, acode, aname`  $\rightsquigarrow$  `Play(pl(pcode)), Actor(act(acode), name(act(acode), aname), actsIn(act(acode), pl(pcode)))`  
`FROM ACTOR`

$m_2$ : `SELECT mcode, acode, mtitle`  $\rightsquigarrow$  `Movie(pl(mcode), playsIn(act(acode), pl(mcode)), title(pl(mcode), mtitle))`  
`FROM MOVIE M, ACTOR A`  
`WHERE M.mcode = A.pcode`  
`AND M.type = "m"`

We define a view  $Aux_i$  for the source query of each mapping  $m_i$ .

For each (split) mapping assertion, we introduce a clause:

```

Play(pl(pcode)) ← Aux1(pcode, -, -)
Actor(act(acode)) ← Aux1(-, acode, -)
name(act(acode), aname) ← Aux1(-, acode, aname)
actsIn(act(acode), pl(pcode)) ← Aux1(pcode, acode, -)
Movie(pl(mcode)) ← Aux2(mcode, -, -)
playsIn(act(acode), pl(mcode)) ← Aux2(mcode, acode, -)
title(pl(mcode), mtitle) ← Aux2(mcode, -, mtitle)
  
```

# Unfolding – Example (cont'd)

Query over the ontology: Actors with their name who act in a movie whose title is "The Matrix":

$q(a, n) \leftarrow \text{Actor}(a), \text{name}(a, n), \text{actsIn}(a, p), \text{Movie}(p), \text{title}(p, \text{"The Matrix"})$

A unifier  $\vartheta$  between the atoms in  $q$  and the clause heads is:

$\vartheta(a) = \text{act}(acode) \quad \vartheta(n) = \text{aname}$

$\vartheta(p) = \text{pl}(pcode) \quad \vartheta(mcode) = pcode \quad \vartheta(mtitle) = \text{"The Matrix"}$

$\text{Actor}(\text{act}(acode))$	$\leftarrow$	$\text{Aux}_1(-, acode, -)$
$\text{name}(\text{act}(acode), \text{aname})$	$\leftarrow$	$\text{Aux}_1(-, acode, \text{aname})$
$\text{actsIn}(\text{act}(acode), \text{pl}(pcode))$	$\leftarrow$	$\text{Aux}_1(pcode, acode, -)$
$\text{Movie}(\text{pl}(mcode))$	$\leftarrow$	$\text{Aux}_2(mcode, -, -)$
$\text{title}(\text{pl}(mcode), \text{mtitle})$	$\leftarrow$	$\text{Aux}_2(mcode, -, \text{mtitle})$

After applying  $\vartheta$  to  $q$ , we obtain:

$q(\text{act}(acode), \text{aname}) \leftarrow \text{Actor}(\text{act}(acode)), \text{name}(\text{act}(acode), \text{aname}), \text{actsIn}(\text{act}(acode), \text{pl}(pcode)),$   
 $\text{Movie}(\text{pl}(pcode)), \text{title}(\text{pl}(pcode), \text{"The Matrix"})$

Substituting the atoms with the bodies of the clauses (after having applied the unifier), we obtain:

$q(\text{act}(acode), \text{aname}) \leftarrow \text{Aux}_1(-, acode, -), \text{Aux}_1(-, acode, \text{aname}), \text{Aux}_1(pcode, acode, -),$   
 $\text{Aux}_2(pcode, -, -), \text{Aux}_2(pcode, -, \text{"The Matrix"})$



# Exponential blowup in the unfolding

When there are multiple mapping assertions for each atom, the unfolded query may be exponential in the original one.

Consider a query:  $q(y) \leftarrow C_1(y), C_2(y), \dots, C_n(y)$

and the mappings:  $m_i^1: \Phi_i^1(x) \rightsquigarrow C_i(\mathbf{iri}(x))$  (for  $i \in \{1, \dots, n\}$ )  
 $m_i^2: \Phi_i^2(x) \rightsquigarrow C_i(\mathbf{iri}(x))$

We add the view definitions:  $\text{Aux}_i^j(x) \leftarrow \Phi_i^j(x)$

and introduce the clauses:  $C_i(\mathbf{iri}(x)) \leftarrow \text{Aux}_i^j(x)$  (for  $i \in \{1, \dots, n\}, j \in \{1, 2\}$ ).

There is a single unifier, namely  $\vartheta(y) = \mathbf{iri}(x)$ , but each atom  $C_i(y)$  in the query unifies with the head of two clauses.

Hence, we obtain **one unfolded query**

$$q(\mathbf{iri}(x)) \leftarrow \text{Aux}_1^{j_1}(x), \text{Aux}_2^{j_2}(x), \dots, \text{Aux}_n^{j_n}(x)$$

for each possible combination of  $j_i \in \{1, 2\}$ , for  $i \in \{1, \dots, n\}$ .

Hence, we obtain  **$2^n$  unfolded queries**.

# Implementation of top-down approach to query answering

To implement the top-down approach, we need to generate an SQL query.

We can follow different strategies:

- 1 Substitute each view predicate in the unfolded queries with the corresponding SQL query over the source:
  - + joins are performed on the DB attributes, hence can be done efficiently, e.g., by exploiting indexes;
  - + does not generate doubly nested queries;
  - the number of unfolded queries may be exponential.
- 2 Construct for each atom in the original query a new view. This view takes the union of all SQL queries corresponding to the view predicates, and constructs also the IRIs based on the IRI templates:
  - + avoids exponential blow-up of the resulting query, since the union (of the queries coming from multiple mappings) is done before the joins;
  - joins are performed on IRIs, i.e., on terms built using string concatenation, hence are highly inefficient;
  - generates doubly nested queries, which per se the database has difficulty in optimizing.

Which method is better, depends on various parameters, and there is no definitive answer.

In general, one needs a mixed approach that applies different strategies to different parts of the query.

# Outline

- 1 Challenges in Data Access
- 2 Virtual Knowledge Graphs for Data Access (and Integration)
- 3 Optimizing Query Answering in VKGs**
  - Query rewriting wrt an OWL 2 QL ontology
  - Query unfolding wrt a mapping
  - Mapping saturation**
  - Optimization of query reformulation
- 4 Designing VKG Mappings
- 5 The Ontop System
- 6 Conclusions

# Contributions of rewriting and unfolding

- We are interested in computing certain answers to SPARQL queries over a VKG instance  $\langle \mathcal{P}, \mathcal{D} \rangle$ , with  $\mathcal{P} = \langle \mathcal{O}, \mathcal{M}, \mathcal{S} \rangle$ .
- In practice, by computing the rewriting  $q_r$  of  $q$  w.r.t.  $\mathcal{O}$  and its unfolding w.r.t.  $\mathcal{M}$ , the resulting query  $q_{\text{unf}}$  might become very large, and costly to execute over  $\mathcal{D}$ .

Let us consider the contributions of rewriting and unfolding to the query answers:

- In principle, evaluating the unfolding  $q_{\text{unf}}$  (of  $q_r$  w.r.t.  $\mathcal{M}$ ) over  $\mathcal{D}$ , gives the same result as evaluating  $q_r$  over the RDF graph  $\mathcal{V} = \mathcal{M}(\mathcal{D})$  extracted through the mapping  $\mathcal{M}$  from the data  $\mathcal{D}$ .
- Instead, the impact of the rewriting on the query answers consists of two components:
  - ① the rewriting w.r.t. class and property hierarchies, i.e.,  $C_1 \sqsubseteq C_2$ ,  $P_1 \sqsubseteq P_2$ ;
  - ② the rewriting taking into account existential reasoning, i.e.,  $C \sqsubseteq \exists R$ ,  $C_1 \sqsubseteq \exists R.C_2$ .

*Note:* Component ① corresponds to computing the saturation  $\mathcal{V}_{\text{sat}}$  of  $\mathcal{V}$  w.r.t. class and property hierarchies, while component ② can be handled only through rewriting.

# Tree-witness rewriting and saturated mapping

We want to avoid materializing  $\mathcal{V}$  and  $\mathcal{V}_{\text{sat}}$ , but also want to avoid computing the query rewriting w.r.t. class and property hierarchies.

Therefore we proceed as follows:

- 1 We rewrite  $q$  only w.r.t. the inclusion assertions that cause existential reasoning (i.e.,  $C \sqsubseteq \exists R.C_1$ ,  $C_1 \sqsubseteq \exists R.C_2$ ).  
 $\rightsquigarrow$  **tree-witness rewriting**  $q_{\text{tw}}$  [Kikot et al. 2012]
- 2 We use instead class and property hierarchies (i.e.,  $C_1 \sqsubseteq C_2$ ,  $P_1 \sqsubseteq P_2$ ) to enrich the mapping  $\mathcal{M}$ .  
 $\rightsquigarrow$  **saturated mapping**  $\mathcal{M}_{\text{sat}}$  [Rodríguez-Muro et al. 2013; Kontchakov, Rezk, et al. 2014]
- 3 We unfold the tree-witness rewriting  $q_{\text{tw}}$  w.r.t. the saturated mapping  $\mathcal{M}_{\text{sat}}$ .

It is possible to show that the resulting query is equivalent to the perfect rewriting  $q_r$  (as obtained, e.g., through ordinary rewriting w.r.t.  $\mathcal{O}$  and unfolding w.r.t.  $\mathcal{M}$ ).

For more details, we refer also to [Kontchakov & Zakharyashev 2014].

# Saturated mapping

Intuitively, the **saturated mapping**  $\mathcal{M}_{\text{sat}}$  is obtained as the composition of  $\mathcal{M}$  and the ontology  $\mathcal{O}$ .

For each mapping assertion in $\mathcal{M}$	and each TBox assertion in $\mathcal{O}$	we add a mapping assertion to $\mathcal{M}_{\text{sat}}$
$\Phi(x) \rightsquigarrow C_1(\text{iri}(x))$	$C_1 \sqsubseteq C_2$	$\Phi(x) \rightsquigarrow C_2(\text{iri}(x))$
$\Phi(x, y) \rightsquigarrow P(\text{iri}_1(x), \text{iri}_2(y))$	$\exists P \sqsubseteq C_1$	$\Phi(x, y) \rightsquigarrow C_1(\text{iri}_1(x))$
$\Phi(x, y) \rightsquigarrow P(\text{iri}_1(x), \text{iri}_2(y))$	$\exists P^- \sqsubseteq C_2$	$\Phi(x, y) \rightsquigarrow C_2(\text{iri}_2(y))$
$\Phi(x, y) \rightsquigarrow P_1(\text{iri}_1(x), \text{iri}_2(y))$	$P_1 \sqsubseteq P_2$	$\Phi(x, y) \rightsquigarrow P_2(\text{iri}_1(x), \text{iri}_2(y))$

Due to saturation,  $\mathcal{M}_{\text{sat}}$  will contain at most  $|\mathcal{O}| \cdot |\mathcal{M}|$  many mappings.

*Note:* The saturated mapping has also been called **T-mapping** in the literature.

# Saturated mapping – Exercise

## Ontology $\mathcal{O}$

Student  $\sqsubseteq$  Person  
 PostDoc  $\sqsubseteq$  Faculty  
 Professor  $\sqsubseteq$  Faculty  
 $\exists$ teaches  $\sqsubseteq$  Faculty  
 Faculty  $\sqsubseteq$  Person

## User-defined mapping assertions $\mathcal{M}$

$\text{student}(scode, fn, ln) \rightsquigarrow \text{Student}(\text{iri1}(scode))$  (1)  
 $\text{academic}(acode, fn, ln, pos), pos = 9 \rightsquigarrow \text{PostDoc}(\text{iri2}(acode))$  (2)  
 $\text{academic}(acode, fn, ln, pos), pos = 2 \rightsquigarrow \text{Professor}(\text{iri2}(acode))$  (3)  
 $\text{teaching}(course, acode) \rightsquigarrow \text{teaches}(\text{iri2}(acode), \text{iri3}(course))$  (4)  
 $\text{academic}(acode, fn, ln, pos) \rightsquigarrow \text{Faculty}(\text{iri2}(acode))$  (5)

By **saturation the mapping**, we obtain  $\mathcal{M}_{\text{sat}}$ , containing additional mapping assertions for the classes **Faculty** and **Person**.

$\text{student}(scode, fn, ln) \rightsquigarrow \text{Person}(\text{iri1}(scode))$  (6)  
 $\text{academic}(acode, fn, ln, pos), pos = 9 \rightsquigarrow \text{Faculty}(\text{iri2}(acode))$  (7)  
 $\text{academic}(acode, fn, ln, pos), pos = 9 \rightsquigarrow \text{Person}(\text{iri2}(acode))$  (8)  
 $\text{academic}(acode, fn, ln, pos), pos = 2 \rightsquigarrow \text{Faculty}(\text{iri2}(acode))$  (9)  
 $\text{academic}(acode, fn, ln, pos), pos = 2 \rightsquigarrow \text{Person}(\text{iri2}(acode))$  (10)  
 $\text{academic}(acode, fn, ln, pos) \rightsquigarrow \text{Person}(\text{iri2}(acode))$  (11)  
 $\text{teaching}(course, acode) \rightsquigarrow \text{Faculty}(\text{iri2}(acode))$  (12)  
 $\text{teaching}(course, acode) \rightsquigarrow \text{Person}(\text{iri2}(acode))$  (13)

# Properties of saturated mappings

## H-complete RDF graph

An RDF graph  $\mathcal{G}$  is **H-complete** w.r.t. an ontology  $\mathcal{O}$ , if, for every RDF triple  $(s, p, o)$ , we have:

$$\langle \mathcal{O}, \mathcal{G} \rangle \models (s, p, o) \quad \text{iff} \quad (s, p, o) \in \mathcal{G}$$

The **saturation**  $\mathcal{G}_{\text{sat}}$  of  $\mathcal{G}$  w.r.t.  $\mathcal{O}$  is the smallest RDF graph that contains  $\mathcal{G}$  and is H-complete w.r.t.  $\mathcal{O}$ .

Intuitively,  $\mathcal{G}_{\text{sat}}$  is obtained from  $\mathcal{G}$  by applying the class and property inclusions of  $\mathcal{O}$ , but without introducing new nodes.

## Relationship between the saturated mapping $\mathcal{M}_{\text{sat}}$ and the saturation of $\mathcal{M}(\mathcal{D})$

- We have that  $\mathcal{M}_{\text{sat}}(\mathcal{D}) = (\mathcal{M}(\mathcal{D}))_{\text{sat}}$  (hence, it is an H-complete RDF graph).
- $\mathcal{M}_{\text{sat}}$  does not depend on the SPARQL query  $q$ , hence it can be pre-computed.
- It can be optimized (by exploiting query containment).



# Mapping optimization – Exercise

## Saturated mapping assertions $\mathcal{M}_{\text{sat}}$

...

<code>academic(acode,fn,ln,pos)</code>	$\rightsquigarrow$ Faculty( <b>iri2</b> (acode))	(5)
<code>student(scode,fn,ln)</code>	$\rightsquigarrow$ Person( <b>iri1</b> (scode))	(6)
<code>academic(acode,fn,ln,pos), pos = 9</code>	$\rightsquigarrow$ Faculty( <b>iri2</b> (acode))	(7)
<code>academic(acode,fn,ln,pos), pos = 9</code>	$\rightsquigarrow$ Person( <b>iri2</b> (acode))	(8)
<code>academic(acode,fn,ln,pos), pos = 2</code>	$\rightsquigarrow$ Faculty( <b>iri2</b> (acode))	(9)
<code>academic(acode,fn,ln,pos), pos = 2</code>	$\rightsquigarrow$ Person( <b>iri2</b> (acode))	(10)
<code>academic(acode,fn,ln,pos)</code>	$\rightsquigarrow$ Person( <b>iri2</b> (acode))	(11)
<code>teaching(course,acode)</code>	$\rightsquigarrow$ Faculty( <b>iri2</b> (acode))	(12)
<code>teaching(course,acode)</code>	$\rightsquigarrow$ Person( <b>iri2</b> (acode))	(13)

Consider also a foreign key over the database relations

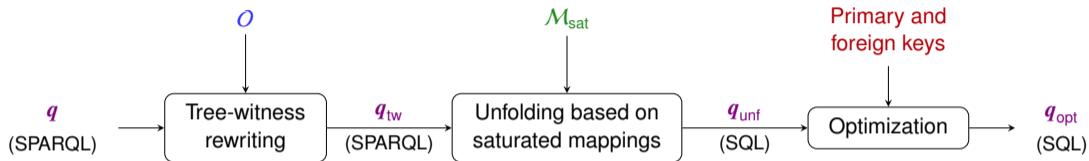
FK:  $\exists y_1. \text{teaching}(y_1, x) \rightarrow \exists y_2 y_3 y_4. \text{academic}(x, y_2, y_3, y_4)$

We can **optimize the mapping** using query containment and the FK. This removes mapping assertions 7, 8, 9, 10, 12, and 13.

...

<code>academic(acode,fn,ln,pos)</code>	$\rightsquigarrow$ Faculty( <b>iri2</b> (acode))	(5)
<code>student(scode,fn,ln)</code>	$\rightsquigarrow$ Person( <b>iri1</b> (scode))	(6)
<code>academic(acode,fn,ln,pos)</code>	$\rightsquigarrow$ Person( <b>iri2</b> (acode))	(11)

# Query reformulation as implemented by the Ontop system



	Step	Input	Output
1.	Tree-witness rewriting	$q$ and $O$	$q_{tw}$
2.	Unfolding	$q_{tw}$ and $M_{sat}$	$q_{unf}$
3.	Optimization	$q_{unf}$ , primary and foreign keys	$q_{opt}$

Let us now consider the optimization step.

# Outline

- 1 Challenges in Data Access
- 2 Virtual Knowledge Graphs for Data Access (and Integration)
- 3 Optimizing Query Answering in VKGs**
  - Query rewriting wrt an OWL 2 QL ontology
  - Query unfolding wrt a mapping
  - Mapping saturation
  - Optimization of query reformulation**
- 4 Designing VKG Mappings
- 5 The Ontop System
- 6 Conclusions

# SQL query optimization

Objective : produce SQL queries that are ...

- similar to manually written ones
- adapted to existing query planners

## Structural optimization

- From join-of-unions to union-of-joins
- IRI decomposition to improve performance of joins

## Semantic optimization

- Redundant join elimination
- Redundant union elimination
- Using functional constraints

## Integrity constraints

- Primary and foreign keys, uniqueness constraints
- Sometimes implicit
- **Vital for query reformulation!**

# Reformulation example – 1. Unfolding

## Saturated mapping

$\text{academic}(acode, fn, ln, pos), pos \in [1..8]$   
 $\rightsquigarrow \text{Teacher}(\text{iri2}(acode))$   
 $\text{teaching}(course, acode) \rightsquigarrow \text{Teacher}(\text{iri2}(acode))$   
 $\text{student}(scode, fn, ln) \rightsquigarrow \text{firstName}(\text{iri1}(scode), fn)$   
 $\text{academic}(acode, fn, ln, pos) \rightsquigarrow \text{firstName}(\text{iri2}(acode), fn)$   
 $\text{student}(scode, fn, ln) \rightsquigarrow \text{lastName}(\text{iri1}(scode), ln)$   
 $\text{academic}(acode, fn, ln, pos) \rightsquigarrow \text{lastName}(\text{iri2}(acode), ln)$

Query (we assume that the ontology is empty, hence  $q_r = q$ )

$q(x, y, z) \leftarrow \text{Teacher}(x), \text{firstName}(x, y), \text{lastName}(x, z)$

We apply **query unfolding**, and then **normalization** to make the join conditions explicit.

$$q_{\text{norm}}(x, y, z) \leftarrow q1_{\text{unf}}(x), q2_{\text{unf}}(x1, y), q3_{\text{unf}}(x2, z), x = x1, x = x2$$

$$q1_{\text{unf}}(\text{iri2}(acode)) \leftarrow \text{academic}(acode, fn, ln, pos), pos \in [1..8]$$

$$q1_{\text{unf}}(\text{iri2}(acode)) \leftarrow \text{teaching}(course, acode)$$

$$q2_{\text{unf}}(\text{iri1}(scode), fn) \leftarrow \text{student}(scode, fn, ln)$$

$$q2_{\text{unf}}(\text{iri2}(acode), fn) \leftarrow \text{academic}(acode, fn, ln, pos)$$

$$q3_{\text{unf}}(\text{iri1}(scode), ln) \leftarrow \text{student}(scode, fn, ln)$$

$$q3_{\text{unf}}(\text{iri2}(acode), ln) \leftarrow \text{academic}(acode, fn, ln, pos)$$

# Reformulation example – 2. Structural optimization

## Unfolded normalized query

$$q_{\text{norm}}(x, y, z) \leftarrow q1_{\text{unf}}(x), q2_{\text{unf}}(x_1, y), \\ q3_{\text{unf}}(x_2, z), \\ x = x_1, x = x_2$$

$$q1_{\text{unf}}(\text{iri2}(a)) \leftarrow \text{academic}(a, f, l, p), \\ p \in [1..8]$$

$$q1_{\text{unf}}(\text{iri2}(a)) \leftarrow \text{teaching}(c, a)$$

$$q2_{\text{unf}}(\text{iri1}(s), f) \leftarrow \text{student}(s, f, l)$$

$$q2_{\text{unf}}(\text{iri2}(a), f) \leftarrow \text{academic}(a, f, l, p)$$

$$q3_{\text{unf}}(\text{iri1}(s), l) \leftarrow \text{student}(s, f, l)$$

$$q3_{\text{unf}}(\text{iri2}(a), l) \leftarrow \text{academic}(a, f, l, p)$$

- While flattening, we can avoid to generate those queries that contain in their body an equality between two terms with incompatible IRI templates.
- This might avoid a potential exponential blowup.

## Flattening (URI template lifting) – Part 1/2

$$q_{\text{liff}}(\text{iri2}(a), y, z) \leftarrow \text{academic}(a, f_1, l_1, p_1), \\ \text{student}(s, f_2, l_2), \\ \text{student}(s_1, f_3, l_3), \\ \text{iri2}(a) = \text{iri1}(s), \\ \text{iri2}(a) = \text{iri1}(s_1), \\ p_1 \in [1..8]$$

$$q_{\text{liff}}(\text{iri2}(a), y, z) \leftarrow \text{academic}(a, f_1, l_1, p_1), \\ \text{student}(s, f_2, l_2), \\ \text{academic}(a_2, f_3, z, p_3), \\ \text{iri2}(a) = \text{iri1}(s), \\ \text{iri2}(a) = \text{iri2}(a_2), \\ p_1 \in [1..8]$$

(One sub-query not shown)

$$q_{\text{liff}}(\text{iri2}(a), y, z) \leftarrow \text{academic}(a, f_1, l_1, p_1), \\ \text{academic}(a_1, y, l_2, p_2), \\ \text{academic}(a_2, f_3, z, p_3), \\ \text{iri2}(a) = \text{iri2}(a_1), \\ \text{iri2}(a) = \text{iri2}(a_2), \\ p_1 \in [1..8]$$

# Reformulation example – 2. Structural optimization

## Unfolded normalized query

$$q_{\text{norm}}(x, y, z) \leftarrow q1_{\text{unf}}(x), q2_{\text{unf}}(x_1, y), \\ q3_{\text{unf}}(x_2, z), \\ x = x_1, x = x_2$$

$$q1_{\text{unf}}(\text{iri2}(a)) \leftarrow \text{academic}(a, f, l, p), \\ p \in [1..8]$$

$$q1_{\text{unf}}(\text{iri2}(a)) \leftarrow \text{teaching}(c, a)$$

$$q2_{\text{unf}}(\text{iri1}(s), f) \leftarrow \text{student}(s, f, l)$$

$$q2_{\text{unf}}(\text{iri2}(a), f) \leftarrow \text{academic}(a, f, l, p)$$

$$q3_{\text{unf}}(\text{iri1}(s), l) \leftarrow \text{student}(s, f, l)$$

$$q3_{\text{unf}}(\text{iri2}(a), l) \leftarrow \text{academic}(a, f, l, p)$$

- While flattening, we can avoid to generate those queries that contain in their body an equality between two terms with incompatible IRI templates.
- This might avoid a potential exponential blowup.

## Flattening (URI template lifting) – Part 2/2

$$q_{\text{lifft}}(\text{iri2}(a), y, z) \leftarrow \text{teaching}(c, a), \\ \text{student}(s, f_2, l_2), \\ \text{student}(s_1, f_3, l_3), \\ \text{iri2}(a) = \text{iri1}(s), \\ \text{iri2}(a) = \text{iri1}(s_1)$$

$$q_{\text{lifft}}(\text{iri2}(a), y, z) \leftarrow \text{teaching}(c, a), \\ \text{student}(s, f_2, l_2), \\ \text{academic}(a_2, f_3, z, p_3), \\ \text{iri2}(a) = \text{iri1}(s), \\ \text{iri2}(a) = \text{iri2}(a_2)$$

(One sub-query not shown)

$$q_{\text{lifft}}(\text{iri2}(a), y, z) \leftarrow \text{teaching}(c, a), \\ \text{academic}(a_1, y, l_2, p_2), \\ \text{academic}(a_2, f_3, z, p_3), \\ \text{iri2}(a) = \text{iri2}(a_1), \\ \text{iri2}(a) = \text{iri2}(a_2)$$

## Reformulation example – 3. Semantic optimization

We are left with just two queries, which we can simplify by eliminating equalities

$$q_{\text{struct}}(\mathbf{iri2}(a), y, z) \leftarrow \text{academic}(a, f_1, l_1, p_1), p_1 \in [1..8], \\ \text{academic}(a, y, l_2, p_2), \\ \text{academic}(a, f_3, z, p_3)$$

$$q_{\text{struct}}(\mathbf{iri2}(a), y, z) \leftarrow \text{teaching}(c, a), \\ \text{academic}(a, y, l_2, p_2), \\ \text{academic}(a, f_3, z, p_3)$$

We can then exploit database constraints (e.g., primary keys) for semantic optimization of the query.

### Self-join elimination (semantic optimization)

$$\text{PK: } \text{academic}(a, f, l, p) \wedge \text{academic}(a, f', l', p') \rightarrow (f = f') \wedge (l = l') \wedge (p = p')$$

$$q_{\text{opt}}(\mathbf{iri2}(a), y, z) \leftarrow \text{academic}(a, y, z, p_1), p_1 \in [1..8]$$

$$q_{\text{opt}}(\mathbf{iri2}(a), y, z) \leftarrow \text{teaching}(c, a), \text{academic}(a, y, z, p_2)$$



# Outline

- 1 Challenges in Data Access
- 2 Virtual Knowledge Graphs for Data Access (and Integration)
- 3 Optimizing Query Answering in VKGs
- 4 Designing VKG Mappings**
- 5 The Ontop System
- 6 Conclusions

# Designing VKG mappings

The form of the mapping in VKGs is critical in ensuring that:

- the resulting VKG specification captures correctly the domain semantics, and
- queries posed over a VKG instance can be answered efficiently.

In designing the mapping assertions, we can rely on some **simple observations**:

- For each **atom in the target part**, the **source query** should be the **simplest SQL query** that retrieves the data that is necessary to populate that atom.
- In particular, we should **avoid unnecessary joins** in the **source query**.
- We should **combine two (or more) atoms** in a single mapping assertion only if they require the **same source query**.
- We need to pay attention to the form of the IRI-templates, to ensure that the “same” ontology object retrieved through multiple mappings is constructed with the same IRI-template.

**However, these observations in general are not sufficient to ensure a good mapping design.**

# Patterns in data sources

- In order to simplify the task of mapping design, it is convenient to identify whether the data source satisfies certain common patterns.
- Each such data pattern can be captured in a sort of “standard” way through a specific form of mapping assertions, combined with some specific form of ontology axiom.
- The presence of a pattern in a data source, and hence the applicability of the corresponding standard encoding into mapping (and ontology axioms), is signaled by the presence of some (combination of) constraints that hold over the relational tables.
- Notice that such constraints might hold:
  - either because they are explicitly declared in the database, and hence enforced by the DBMS,
  - or because they are implied by the semantics of the domain, even though they might not be declared explicitly in the database.

# Looking at database design principles

In relational database design, **well-established conceptual modeling principles** and **methodologies** are usually employed.

- The resulting schema should suitably reflect the application domain at hand.
- This design phase relies on semantically-rich representations such as ER diagrams.
- However, these representations, typically:
  - get lost during deployment, since they are not conveyed together with the database itself, or
  - quickly get outdated due to continuous adjustments triggered by changing requirements.

## Key Observation

While the relational model may be semantically-poor with respect to ontological models, the original semantically-rich design of the application domain **leaves recognizable footprints** that can be converted into ontological mapping patterns.

# VKG mapping patterns

Therefore, in designing VKG mapping patterns, we draw an explicit and precise connection with conceptual modeling practices found in DB design, while exploiting all of:

- the relational schema with its constraints
- the conceptual schema at the basis of the relational schema
- extensional data stored in the DB (when available)
- the domain knowledge that is encoded in ontology axioms

# Catalog of mapping patterns

To come up with a catalog of mapping patterns, we can rely on well-established methodologies and patterns studied in:

- data management – e.g., W3C Direct Mapping Specification [Arenas et al. 2012] and extensions
- data analysis – e.g., algorithms for discovering dependencies, and
- conceptual modeling

The specification of each pattern includes:

- the three components of a VKG specification: DB schema, ontology, mapping between the two;
- the conceptual schema of the domain of interest;
- underlying data, when available.

Note that the patterns do not fix what is given as input and what is produced as output, but simply describe how the different elements relate to each other.

# Two major groups of mapping patterns

## Schema-driven patterns

Are shaped by the structure of the DB schema and its explicit constraints.

## Data-driven patterns

- Consider also constraints emerging from specific configurations of the data in the DB.
  - For each schema-driven pattern, we identify a data-driven version:  
The constraints over the schema are not explicitly specified, but hold in the data.
  - We provide also data-driven patterns that do not have a schema-driven counterpart.
- 
- We use also additional semantic information from the ontology  $\rightsquigarrow$  **Pattern modifiers**
  - Some patterns come with **views over the DB-schema**:
    - Views reveal structures over the DB-schema, when the pattern is applied.
    - Views can be used to identify the applicability of further patterns.

# Constraints on the data

When defining the mapping patterns, we consider the traditional types of DB constraints:

- **Primary key constraint:**  $T(\underline{\mathbf{K}}, \mathbf{A})$
- **Key constraint:**  $key_T(\mathbf{K})$
- **Foreign key constraint:**  $T_1[\mathbf{A}] \subseteq T_2[\mathbf{K}]$ , where  $\mathbf{K}$  is a (typically primary) key of relation  $T_2$ .  
We use the notation:



*Note:* We use normal font (e.g.,  $A$ ) for single attributes, and boldface for sets of attributes (e.g.,  $\mathbf{A}$ ).



# Types of mapping patterns

In the following, we discuss the following mapping patterns:

- Entity (MpE)
- Relationship (MpR)
- Relationship with Identifier Alignment (MpRa)
- Relationship with Merging (MpRm)
- 1-1 Relationship with Merging (MpR11m)
- Reified Relationship (MpRR)
- Hierarchy (MpH)
- Hierarchy with Identifier Alignment (MpHa)
- Clustering Entity to Class / Data Property / Object Property (MpCE2X)

We present each mapping pattern by specifying the following four components:

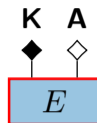
- 1 The constraints over the relational schema/data that make the patterns applicable.
- 2 A possible conceptual schema (specified as an Entity-Relationship diagram) that corresponds to such constraints. The elements that are directly affected by the pattern and that give rise to the mapping assertions are outlined in **red**.
- 3 The source and target part of the resulting mapping assertion(s).
- 4 The ontology axioms that should hold.

**Note:** In the following, we make use of IRI-templates of the form “:E/{K}”, where we assume that “:E/” is a prefix that is specific for the instances of a class  $C_E$ .

# Mapping pattern: Entity (MpE)

Relational schema and constraints:

$T_E(\underline{\mathbf{K}}, \mathbf{A})$



Mapping assertion:

$s : T_E$

$t : :E/\{\mathbf{K}\} \text{ rdf:type } C_E .$

$\{ :E/\{\mathbf{K}\} d_A \{A\} . \}_{A \in \mathbf{KUA}}$

Ontology axioms:

$\{ \exists d_A \sqsubseteq C_E \}_{A \in \mathbf{KUA}}$

For the application of the mapping pattern, we observe the following:

- This fundamental pattern considers a single table  $T_E$  with primary key  $\underline{\mathbf{K}}$  and other relevant attributes  $\mathbf{A}$ .
- The pattern captures how  $T_E$  is mapped into a corresponding class  $C_E$ .
- The primary key  $\underline{\mathbf{K}}$  of  $T_E$  is used to construct the objects that are instances of  $C_E$ , using a template  $:E/\{\mathbf{K}\}$  specific for  $C_E$ .
- Each relevant attribute of  $T_E$  is mapped to a data property of  $C_E$ .

# Mapping pattern: Entity (MpE) – Example

Consider a **TClient** table containing **ssns** of clients, together with **name**, **dateOfBirth**, and **hobbies** as additional attributes.

```
TClient(ssn, name, dateOfBirth, hobbies)
```

**Mapping:** **TClient** is mapped to a **Client** class using the attributes **ssn** to construct its objects. In addition, the **ssn**, **name**, and **dateOfBirth** are used to populate in the object position the three data properties **ssn**, **name**, and **dob**, respectively. The attribute **hobbies** is ignored.

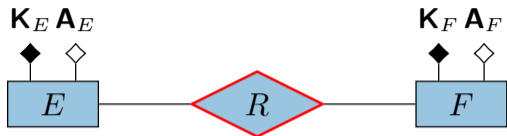
```
mappingId MClient
source     SELECT ssn, name, dateOfBirth FROM TClient
target     :C/{ssn} rdf:type :Client ;
           :ssn {ssn} ;
           :name {name} ;
           :dob {dateOfBirth} .
```

# Mapping pattern: Relationship (MpR)

Relational schema and constraints:

$$T_E(\mathbf{K}_E, \mathbf{A}_E) \quad T_F(\mathbf{K}_F, \mathbf{A}_F)$$

$$T_R(\mathbf{K}_{RE}, \mathbf{K}_{RF})$$



Mapping assertion:

$$s : T_R$$

$$t : :E/\{\mathbf{K}_{RE}\} \quad p_R \quad :F/\{\mathbf{K}_{RF}\} .$$

Ontology axioms:

$$\exists p_R \sqsubseteq C_E$$

$$\exists p_R^- \sqsubseteq C_F$$

For the application of the mapping pattern, we observe the following:

- This pattern considers three tables  $T_R$ ,  $T_E$ , and  $T_F$ .
- The primary key of  $T_R$  is partitioned into two parts  $\mathbf{K}_{RE}$  and  $\mathbf{K}_{RF}$  that are foreign keys to  $T_E$  and  $T_F$ , respectively.
- $T_R$  has no additional (relevant) attributes.
- The pattern captures how  $T_R$  is mapped to an object property  $p_R$ , using the two parts  $\mathbf{K}_{RE}$  and  $\mathbf{K}_{RF}$  of the primary key to construct respectively the subject and the object of the triples in  $p_R$ .

# Mapping pattern: Relationship (MpR) – Example

An additional **TAddress** table in the client registry stores the addresses at which each client can be reached, and such table has a foreign key to a table **TLocation** storing locations using attributes **city** and **street**.

```
TClient(ssn, name, dateOfBirth, hobbies)
```

```
TLocation(city, street)
```

```
TAddress(client, locCity, locStreet)
```

```
FK: TAddress[client] -> Tclient[ssn]
```

```
FK: TAddress[locCity, locStreet] -> TLocation[city, street]
```

**Mapping:** The **TAddress** table is mapped to an **address** object property, for which the ontology asserts that the domain is the class **Person** and the range an additional class **Location**, corresponding to the **TLocation** table.

```
mappingId MAddress
source     SELECT client, locCity, locStreet FROM TAddress
target     :C/{client} :address :L/{locCity}/{locStreet} .
```

# Mapping pattern: Relationship with Identifier Alignment (MpRa)

Relational schema and constraints:

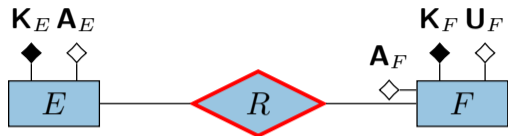
$$T_E(\underline{K_E}, A_E) \quad T_F(\underline{K_F}, U_F, A_F)$$

$$T_R(\underline{K_{RE}}, U_{RF}) \quad \text{key}_{R_F}(U_F)$$

Mapping assertion:

$$s : T_R \bowtie_{U_{RF}=U_F} T_F$$

$$t : :E/\{K_{RE}\} \ p_R \ :F/\{K_F\}.$$



Ontology axioms:

$$\exists p_R \sqsubseteq C_E$$

$$\exists p_R^- \sqsubseteq C_F$$

For the application of the mapping pattern, we observe the following:

- Such pattern is a variation of pattern **MpR**, in which the foreign key in  $T_R$  does not point to the primary key  $K_F$  of  $T_F$ , but to an additional key  $U_F$ .
- Since the instances of class  $C_F$  corresponding to  $T_F$  are constructed using the primary key  $K_F$  of  $T_F$  (cf. pattern **MpE**), also the pairs that populate  $p_R$  should refer in their object position to  $K_F$ .
- Note that  $K_F$  can only be retrieved by a join between  $T_R$  and  $T_F$  on the additional key  $U_F$ .

# Mapping pattern: Rel. with Identifier Alignment (MpRa) – Example

The primary key of the table **TLocationCoord** is now not given by the **city** and **street**, which are used in the table **TAddress** that relates clients to their addresses, but is given by the **latitude** and **longitude** of locations.

```
TClient(ssn, name, dateOfBirth, hobbies)
```

```
TLocationCoord(latitude, longitude, city, street)    key[TLocation]: city, street
```

```
TAddress(client, locCity, locStreet)
```

```
FK: TAddress[client] -> TClient[ssn]
```

```
FK: TAddress[locCity, locStreet] -> TLocationCoord[city, street]
```

**Mapping:** The **Address** table is mapped to an **address** object property, for which the ontology asserts that the domain is the class **Person** and the range an additional class **Location**, corresponding to the **Location** table.

```
mappingId MAddressCoord
source     SELECT client, latitude, longitude
          FROM TAddress JOIN TLocationCoord ON locCity = city AND locStreet = street
target    :C/{client} :address :LC/{latitude}/{longitude} .
```

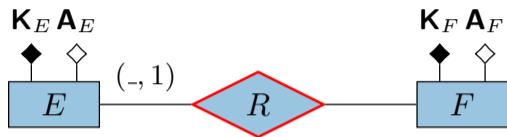
# Mapping pattern: Relationship with Merging (MpRm)

Relational schema and constraints:

$$T_F(\underline{\mathbf{K}}_F, \mathbf{A}_F)$$

$$\uparrow$$

$$T_E(\underline{\mathbf{K}}_E, \mathbf{K}_{EF}, \mathbf{A}_E)$$



Mapping assertion:

$$s : T_E$$

$$t : :E/\{\mathbf{K}_E\} p_R :F/\{\mathbf{K}_{EF}\} .$$

Ontology axioms:

$$\exists p_R \sqsubseteq C_E$$

$$\exists p_R^- \sqsubseteq C_F$$

For the application of the mapping pattern, we observe the following:

- Such pattern is characterized by a table  $T_E$  in which the foreign key  $\mathbf{K}_{EF}$  to a table  $T_F$  is disjoint from its primary key  $\mathbf{K}_E$ .
- The table  $T_E$  is mapped to an object property  $p_{EF}$ , whose subject and object are derived respectively from  $\mathbf{K}_E$  and  $\mathbf{K}_{EF}$ .



# Mapping pattern: Relationship with Merging (MpRm) – Example

The relationship between a client and its unique billing address has been merged into the **TClient** table. The ontology defines a **billingAddress** object property, whose domain is the **Client** class and whose range is the **Location** class.

```
TLocation(city, street)
```

```
TClient(ssn, name, dateOfBirth, billCity, billStreet, hobbies)
```

```
FK: TClient[billCity, billStreet] -> TLocation[city, street]
```

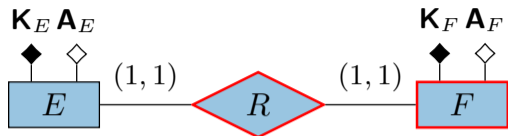
**Mapping:** The billing address information is extracted by a mapping from the **TClient** table to **billingAddress**.

```
mappingId MBillingAddress
source     SELECT ssn, billCity, billStreet FROM TClient
target     :C/{ssn} :billingAddress :L/{billCity}/{billStreet} .
```

# Mapping pattern: 1-1 Relationship with Merging (MpR11m)

Relational schema and constraints:

$$\begin{array}{c}
 T_E(\underline{\mathbf{K}}_E, \mathbf{A}_E, \mathbf{K}_F, \mathbf{A}_F) \quad \text{key}_{T_E}(\mathbf{K}_F) \\
 \hline
 V_E(\underline{\mathbf{K}}_E, \mathbf{A}_E) = \pi_{\mathbf{K}_E, \mathbf{A}_E}(T_E) \quad V_F(\mathbf{K}_F, \mathbf{A}_F) = \pi_{\mathbf{K}_F, \mathbf{A}_F}(T_E) \\
 \underbrace{\qquad\qquad\qquad}_{V_R(\underline{\mathbf{K}}_E, \mathbf{K}_F) = \pi_{\mathbf{K}_E, \mathbf{K}_F}(T_E)}
 \end{array}$$



Mapping assertion:

$s : T_E$

$t : :F/\{\mathbf{K}_F\} \text{ rdf:type } C_F.$

$\{ :F/\{\mathbf{K}_F\} d_A \{A\} . \}_{A \in \mathbf{K}_F \cup \mathbf{A}_F}$

$:E/\{\mathbf{K}_E\} p_R :F/\{\mathbf{K}_F\}.$

Ontology axioms:

$\{\exists d_A \sqsubseteq C_F\}_{A \in \mathbf{K}_F \cup \mathbf{A}_F}$

$\exists p_R \sqsubseteq C_E$

$\exists p_R^- \sqsubseteq C_F$

For the application of the mapping pattern, we observe the following:

- The pattern could be applied when a table  $T_E$  has a primary key  $\mathbf{K}_E$  and an additional key  $\mathbf{K}_F$ .
- Moreover, domain knowledge of the ontology indicates that objects with IRI  $:F/\{\mathbf{K}_F\}$  are relevant in the domain, and that they have data properties that correspond to the attributes  $\mathbf{A}_F$  of  $T_E$ .
- When this pattern is applied, the key  $\mathbf{K}_F$  and the attributes  $\mathbf{A}_F$ , can be projected out from  $T_E$ , resulting in a view  $V_E$  to which further patterns can be applied, including **SR11m** itself.

# Mapping pattern: 1-1 Relationship with Merging (MpR11m) – Example

A single table **TUniversity**, containing the information about universities, contains also information about their rector. The given ontology contains both a **University** and a **Rector** class.

```
TUniversity(uname, numfaculties, recssn, recname, recdob, salary)
  key[TUniversity]: recssn
```

**Mapping:** The attribute **recssn** in **TUniversity**, identifying the rector, is used to form the IRIs for the instances of **Rector**, and the attributes **recname** and **recdob**, intuitively belonging to the rector, are mapped to data properties that have as domain **Rector** (as opposed to **University**).

```
mappingId MUniversity
source    SELECT uname, numfaculties FROM TUniversity
target    :U/{uname} rdf:type :University ; :numfac {numfaculties} .
```

```
mappingId MRector
source    SELECT recssn, recname, recdob FROM TUniversity
target    :P/{recssn} rdf:type :Rector ;
           :ssn {recssn} ; :name {recname} ; :dob {recdob} .
```

```
mappingId MhasRector
source    SELECT uname, recssn FROM TUniversity
target    :U/{uname} :hasRector :P/{recssn} .
```

# Mapping pattern: 1-1 Relationship with Merging (MpR11m) – Notes

- Notice that to apply pattern **MpR11m**, domain knowledge is inherently required to determine to which class the attributes should be associated.
- For example, assume that the table **TUniversity** contains an attribute for the **salary** of the rector. Then, we have two possibilities:
  - the salary is considered a property of the rector, e.g., if the salary is negotiated individually by the rector.
  - the salary is considered a property of the university, e.g., if the salary of the rector is determined by some regulation of the university.

Distinguishing which of these two possibilities is the correct one, requires in-depth knowledge about the domain.

- The necessary domain knowledge may also come from the ontology, e.g., if the data properties corresponding to the attributes are already present in the ontology, and their domain has been declared.

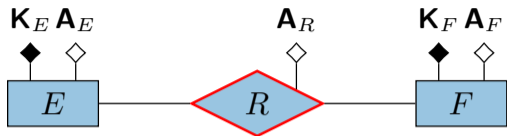
# Mapping pattern: Reified Relationship (MpRR) – Attribute case

Relational schema and constraints:



Mapping assertion:

$s : T_R$   
 $t : :R/\{K_{RE}\}/\{K_{RF}\} \text{ rdf:type } C_R .$   
 $\{ :R/\{K_{RE}\}/\{K_{RF}\} d_A \{A\} . \}_{A \in K_{RE} \cup K_{RF} \cup A_R}$   
 $:R/\{K_{RE}\}/\{K_{RF}\} p_{RE} :E/\{K_{RE}\} .$   
 $:R/\{K_{RE}\}/\{K_{RF}\} p_{RF} :F/\{K_{RF}\} .$



Ontology axioms:

$\exists p_{RE} \sqsubseteq C_R$        $\exists p_{RF} \sqsubseteq C_R$   
 $\exists p_{RE}^- \sqsubseteq C_E$        $\exists p_{RF}^- \sqsubseteq C_F$   
 $\{ \exists d_A^- \sqsubseteq C_R \}_{A \in K_{RE} \cup K_{RF} \cup A}$

For the application of the mapping pattern, we observe the following:

- The pattern applies to a table  $T_R$  whose primary key is partitioned in (at least) two parts  $K_{RE}$  and  $K_{RF}$  that are foreign keys to additional tables, and there are additional attributes  $A_R$  in  $T_R$ .
- Since  $T_R$  corresponds to a conceptual element that has itself properties (corresponding to  $A_R$ ), to represent it in the ontology we require a class  $C_R$  whose instances have an IRI  $:R/\{K_{RE}\}/\{K_{RF}\}$ .
- The mapping ensures that each components of the relationship is represented by an object property ( $p_{RE}$ ,  $p_{RF}$ ), and that the tuples instantiating them can all be derived from  $T_R$  alone.

# Mapping pattern: Reified Relationship (MpRR) – $n$ -ary relationship case

Relational schema and constraints:



Mapping assertion:  $\mathbf{K}_R := \mathbf{K}_{RE} \cup \mathbf{K}_{RF} \cup \mathbf{K}_{RG}$

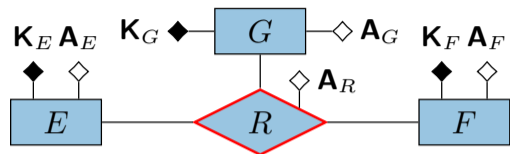
$s : T_R$

$t : :R/\{\mathbf{K}_R\} \text{ rdf:type } C_R .$

$\{ :R/\{\mathbf{K}_R\} d_A \{A\} . \}_{A \in \mathbf{K}_R \cup \mathbf{A}_R}$

$:R/\{\mathbf{K}_R\} p_{RE} :E/\{\mathbf{K}_{RE}\} .$

$:R/\{\mathbf{K}_R\} p_{RF} :F/\{\mathbf{K}_{RF}\} . \quad :R/\{\mathbf{K}_R\} p_{RG} :G/\{\mathbf{K}_{RG}\} .$



Ontology axioms:

$\exists p_{RE} \sqsubseteq C_R$

$\exists p_{RF} \sqsubseteq C_R$

$\exists p_{RG} \sqsubseteq C_R$

$\exists p_{RE}^- \sqsubseteq C_E$

$\exists p_{RF}^- \sqsubseteq C_F$

$\exists p_{RG}^- \sqsubseteq C_G$

$\{ \exists d_A^- \sqsubseteq C_R \}_{A \in \mathbf{K}_R \cup \mathbf{A}_R}$

For the application of the mapping pattern, we observe the following:

- The pattern applies to a table  $T_R$  whose primary key is partitioned in at least three parts  $\mathbf{K}_{RE}$ ,  $\mathbf{K}_{RF}$ , and  $\mathbf{K}_{RG}$ , that are foreign keys to three additional tables.
- Additional attributes  $\mathbf{A}_R$  might also be present in  $T_R$ .
- Apart from the arity of the relationship, the pattern behaves analogously to **MpRR** for the attribute case.

# Mapping pattern: Reified Relationship (MpRR) – Example

Consider a table **TExam** containing information about university exams, (which involve a student, a course, and a professor teaching that course), that has foreign keys towards three tables, namely **TStudent**, **TCourse**, and **TProfessor**.

**TExam**(student, course, professor, grade)

**TStudent**(ssn, sname)

**TCourse**(cid, cname, credits)

**TProfessor**(ssn, pname, level)

FK: **TExam**[student] -> **TStudent**[ssn]

FK: **TExam**[course] -> **TCourse**[cid]

FK: **TExam**[professor] -> **TProfessor**[ssn]

**Mapping:** This information is represented by a relationship that is inherently ternary. The ontology should contain a class **Exam** corresponding to the reified relationship, connected via three object properties to the classes **Student**, **Course**, and **Professor**. The mapping ensures that the class **Exam** is instantiated with objects whose IRI is constructed from the identifiers of the component classes.

```
mappingId MExam
source      SELECT student, course, professor, grade FROM TExam
target      :E/{student}/{course}/{professor} rdf:type :Exam ;
                                                    :examOf :P/{student} ;
                                                    :examFor :C/{course} ;
                                                    :examBy  :P/{professor} ;
                                                    :examGrade {grade} .
```

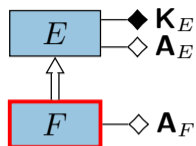
# Mapping pattern: Hierarchy (MpH)

Relational schema and constraints:

$$T_E(\underline{\mathbf{K}}_E, \mathbf{A}_E)$$

$$\uparrow$$

$$T_F(\underline{\mathbf{K}}_{FE}, \mathbf{A}_F)$$



Mapping assertions:

$$s : T_F$$

$$t : :E/\{\mathbf{K}_{FE}\} \text{ rdf:type } C_F .$$

$$\{ :E/\{\mathbf{K}_{FE}\} d_A \{A\} . \}_{A \in \mathbf{A}_F}$$

Ontology axioms:

$$C_F \sqsubseteq C_E$$

$$\{ \exists d_A^- \sqsubseteq C_F \}_{A \in \mathbf{A}_F}$$

For the application of the mapping pattern, we observe the following:

- The pattern considers a table  $T_F$  whose primary key is a foreign key to a table  $T_E$ .
- Then,  $T_F$  is mapped to a class  $C_F$  in the ontology that is a sub-class of the class  $C_E$  to which  $T_E$  is mapped.
- Hence,  $C_F$  “inherits” the template  $:E/\{\cdot\}$  of  $C_E$ , so that the instances of the two classes are “compatible”.



# Mapping pattern: Hierarchy (MpH) – Example

Consider a table **TPerson** containing information about persons, and a table **TStudent** containing information about students, which has a foreign key towards **TPerson**.

```
TPerson(ssn, name, dateOfBirth)
```

```
TStudent(ssn, sid, credits)
```

```
FK: TStudent[ssn] -> TPerson[ssn]
```

**Mapping:** The two tables **TPerson** and **TStudent** are mapped to two classes **Person** and **Student**, respectively, each with data properties corresponding to the attributes of the table. Moreover, the ontology will contain an axiom stating that **Student** is a sub-class of **Person**.

```
mappingId MPerson
source      SELECT ssn, name, dob FROM TPerson
target      :P/{ssn} rdf:type :Person ;
            :name {name} ;
            :dob {dateOfBirth} .
```

```
mappingId MStudent
source      SELECT ssn, sid FROM TStudent
target      :P/{ssn} rdf:type :Student ;
            :studentId {sid} ;
            :credits {credits} .
```

# Mapping pattern: Hierarchy with Identifier Alignment (MpHa)

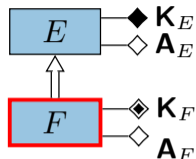
Relational schema and constraints:

$$T_E(\underline{\mathbf{K}}_E, \mathbf{A}_E) \text{ key}_{T_F}(\mathbf{U}_F)$$

$$T_F(\underline{\mathbf{K}}_F, \mathbf{U}_F, \mathbf{A}_F)$$

$$\frac{T_E(\underline{\mathbf{K}}_E, \mathbf{A}_E) \text{ key}_{V_F}(\mathbf{K}_F)}{V_F(\underline{\mathbf{K}}_F, \mathbf{U}_F, \mathbf{A}_F) = T_F}$$

$$V_F(\underline{\mathbf{K}}_F, \mathbf{U}_F, \mathbf{A}_F) = T_F$$



Mapping assertions:

$$s : T_F$$

$$t : :E/\{\mathbf{U}_F\} \text{ rdf:type } C_F.$$

$$\{ :E/\{\mathbf{U}_F\} d_A \{A\} . \}_{A \in \mathbf{K}_F \cup \mathbf{A}_F}$$

Ontology axioms:

$$C_F \sqsubseteq C_E$$

$$\{ \exists d_A^- \sqsubseteq C_F \}_{A \in \mathbf{K}_F \cup \mathbf{A}_F}$$

For the application of the mapping pattern, we observe the following:

- Such pattern is like **MpH**, but the foreign key in  $T_F$  is over a key  $\mathbf{U}_F$  that is not primary.
- The objects for  $C_F$  have to be built out of  $\mathbf{U}_F$ , rather than out of its primary key  $\mathbf{K}_F$ .
- For this purpose, the pattern creates a view  $V_F$  in which  $\mathbf{U}_F$  is the primary key, and the foreign key relations are preserved.

# Mapping pattern: Hierarchy with Identifier Alignment (MpHa) – Example

Consider the tables **TPerson** and **TStudent** of the previous example, but assume now that the primary key of **TStudent** is **sid**. Consider also an additional table **TEnrolled**, recording course enrollments.

**TPerson**(ssn, name, dateOfBirth)

**TStudent**(sid, ssn, credits)

**TEnrolled**(student, course)

FK: TStudent[ssn] -> TPerson[ssn]      key[TStudent]: ssn

FK: TEnrolled[student] -> TStudent[sid]

**Mapping:** By applying pattern **MpHa**, we identify the instances of **Student** by their **ssn**, and we create a view **VStudent**(sid, ssn, credits). But now, considering this view instead of **TStudent**, in order to map **TEnrolled** into an object property **enrolledIn**, we need to apply pattern **MpRa** rather than **MpR**.

```

mappingId MPerson
source    SELECT ssn, name, dob FROM TPerson
target    :P/{ssn} rdf:type :Person ;      :name {name} ;      :dob {dateOfBirth} .

mappingId MStudent
source    SELECT sid, ssn, credits FROM TStudent
target    :P/{ssn} rdf:type :Student ;     :studentId {sid} ;     :credits {credits} .

mappingId MEnrolled
source    SELECT ssn, course FROM TEnrolled JOIN TStudent ON student = sid
target    :P/{ssn} :enrolledIn :C/{course} .
  
```

# Mapping pattern: Clustering Entity to Class (MpCE2C)

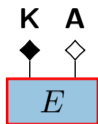
Relational schema and constraints:

$T_E(\underline{\mathbf{K}}, \mathbf{A})$ ,

$\mathbf{B} \subseteq \mathbf{K} \cup \mathbf{A}$  such that  $partition_{\mathcal{D}}(\mathbf{B}, T_E)$

---

$\{ V_{E_v}(\mathbf{K}, \mathbf{A}) = \sigma_{\mathbf{B}=v}(T_E) \}_{v \in \pi_{\mathbf{B}}(T_E)}$



$\mathbf{B} \subseteq \mathbf{K} \cup \mathbf{A}$   
such that  $partition_{\mathcal{D}}(\mathbf{B}, E)$

Mapping assertions:

$\{ s : \sigma_{\mathbf{B}=v} T_E$

$t : :E/\{\mathbf{K}\} \text{ rdf:type } C_E^v . \}_{v \in \pi_{\mathbf{B}}(T_E)}$

Ontology axioms:

$\{ C_E^v \sqsubseteq C_E \}_{v \in \pi_{\mathbf{B}}(T_E)}$

For the application of the mapping pattern, we observe the following:

- This pattern is characterized by a table  $T_E$  corresponding to a class  $C_E$ , and a **derivation rule** defining sub-classes of  $C_E$  according to the values for attributes  $\mathbf{B}$  in  $T_E$ .
- Accordingly, instances in  $T_E$  can be mapped to ontology objects in the sub-classes  $C_E^v$  of  $C_E$ .
- As for other patterns, this pattern produces views according to the possible values  $v$  of  $\mathbf{B}$ .

# Mapping pattern: Clustering Entity to Class (MpCE2C) – Example

Consider a table **TPerson** containing persons with an attribute defining their gender and ranging over 'F' or 'M'.

```
TPerson(ssn, name, dob, gender)
```

**Mapping:** The ontology defines a class **Person** with two subclasses **Female** and **Male**. Pattern **MpCE2C** clusters the table according to the gender attribute, and instantiates the classes **Female** and **Male** accordingly.

```
mappingId MPerson
source    SELECT ssn, name, dob FROM TPerson
target    :P/{ssn} rdf:type :Person ;      :name {name} ;      :dob {dateOfBirth} .

mappingId MFemale
source    SELECT ssn FROM TPerson WHERE gender = 'F'
target    :P/{ssn} rdf:type :Female .

mappingId MMale
source    SELECT ssn FROM TPerson WHERE gender = 'M'
target    :P/{ssn} rdf:type :Male .
```

# Further mapping patterns

- Similarly to the previous pattern, which clusters instances of a class into different subclasses, we can consider patterns that generate a cluster of data properties, or a cluster of object properties, according to different criteria that can be applied to the source data.
- In order to understand when such patterns can be applied, and then define the corresponding mapping assertions and the expected underlying ontology axioms, we can proceed in a way similar to the case of a cluster of (sub)classes.
- More in general, we might conceive also additional patterns that involve more complex operations or queries over the data.
- Also, in any (sufficiently complex) real-world integration scenario, many cases will occur for which none of the specified pattern applies.
- Therefore, based on (the knowledge that the designer has about) the domain semantics, and the constructs that are available in the ontology, in general also ad-hoc mappings need to be defined.

# Additional considerations on IRI-templates

- As we have seen, it is a good practice to include in the IRI-template a prefix that depends on the kind of object (i.e., the class).
- In the case of **ISA hierarchies**, one has to pay attention on whether to use **the same or different templates** for the various classes in the hierarchy:
  - Using the same template allows for specifying joins across the various classes of the hierarchy.
  - Using different templates allows for differentiating the different classes and for applying stricter pruning of queries (as we have seen).
- One has also to consider whether to include info about the **data source** as part of the IRI-template or not:
  - In general, this is not done, which makes the data sources transparent to the user who queries.
  - By including the data source in the IRI-template, such information is recorded in the created objects.

# Design scenarios for VKG mapping patterns

Depending on what information is available, we can consider different design scenarios where the patterns can be applied:

- 1 **Debugging of a VKG specification** that is already in place.
- 2 **Conceptual schema reverse engineering** for a DB that represents the domain of interest by using a given full VKG specification.
- 3 **Mapping bootstrapping** for a given DB and ontology that miss the mappings relating them.
- 4 **Ontology + mapping bootstrapping** from a given DB with constraints, and possibly a conceptual schema.
- 5 **VKG bootstrapping**, where the goal is to set up a full VKG specification from a conceptual schema of the domain.



# Automating the mapping design process

- In a complex real-world scenario, understanding the domain semantics, the semantics of the data sources, and how the sources have to be related to the global schema/ontology can be rather resource intensive and therefore costly.
- Currently, there are no tools that completely automate this process, and it is unlikely that a completely automated solution is possible at all.
- However, there are tools that provide automated support for the (already difficult) task of understanding which elements in one schema (e.g., a source) can correspond to which elements of another schema (e.g., the global schema). This task is called **schema matching**.
- Based on a proposed match between elements, mapping patterns can provide valuable indications on how to convert the match into an actual mapping, i.e., how to define the (SQL) queries that correctly relate the semantics of the sources to that of the ontology.
- Also, mapping patterns can be automatically discovered, either by considering the constraints on the data sources, or, more interestingly, derive the constraints from the actual data, even when they are not defined over the sources at the schema level.

# Outline

- 1 Challenges in Data Access
- 2 Virtual Knowledge Graphs for Data Access (and Integration)
- 3 Optimizing Query Answering in VKGs
- 4 Designing VKG Mappings
- 5 The Ontop System**
- 6 Conclusions

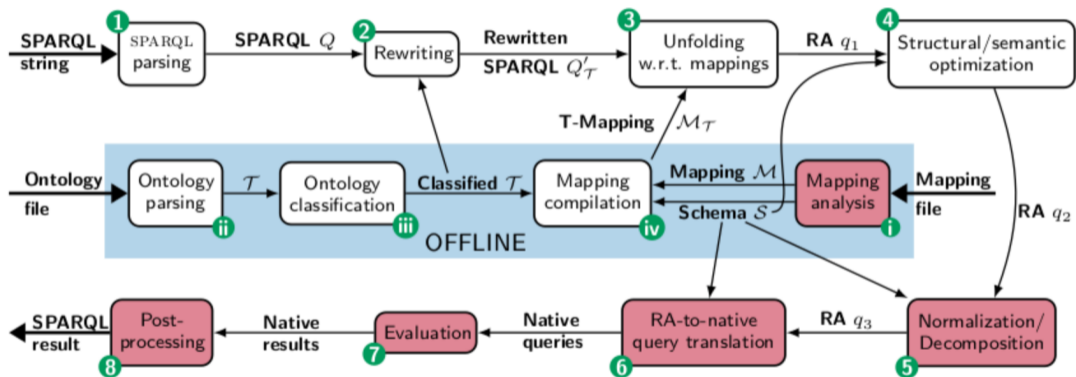
# The *Ontop* system [C., Cogrel, et al. 2017, Semantic Web J.], [Xiao, Lanti, et al. 2020, ISWC]



<https://ontop-vkg.org/>

- State-of-the-art VKG system.
- Implements the presented techniques for query answering and optimization.
- Addresses the key challenges of scalability and performance.
- Compliant with all relevant Semantic Web standards:  
RDF, RDFS, OWL 2 QL, R2RML, SPARQL, and GeoSPARQL.
- Supports all major relational DBMSs:  
Oracle, DB2, MS SQL Server, Postgres, MySQL, Teiid, Dremio, Denodo, etc.
- **Open-source** and released under Apache 2 license.

# Query answering in *Ontop*



# The *Ontopic* spinoff of unibz

# ONTOPIC

<https://ontopic.ai/>

Funded in April 2019 as the first spin-off of the Free University of Bozen-Bolzano.

- **Ontopic Studio** ready to be released
  - Ensures scalability, reliability, and cost-efficiency at design and runtime of VKG solutions.
  - Strong focus on usability.
- **Technical services**
  - Technical support for Ontop and Ontopic Studio.
  - Customized developments.
- **Consulting** on adoption of VKG-based solutions for data access and integration.

# Outline

- 1 Challenges in Data Access
- 2 Virtual Knowledge Graphs for Data Access (and Integration)
- 3 Optimizing Query Answering in VKGs
- 4 Designing VKG Mappings
- 5 The Ontop System
- 6 Conclusions**

# Conclusions

- VKGs are by now a mature technology to address the challenges related to data access and integration.
- It has been well-investigated and applied in many different scenarios mostly for the case of relational data sources.
- The technology is general purpose, and it can be tailored towards specific domains, relying also on standard ontologies.
- Performance and scalability w.r.t. larger datasets (**volume**), larger and more complex ontologies (**variety**, **veracity**), and multiple heterogeneous data sources (**variety**, **volume**) is a challenge.
- Currently, VKGs are being investigated for alternative types of data, such as **temporal data**, **graph data**, **tree structured data**, **linked open data**, and **geo-spatial data**.
- Performance and scalability are even more critical for these more complex domains.

Thank you!



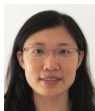
# A great thank you to all my collaborators



Elena  
Botoeva



Julien  
Corman



Linfang  
Ding



Elem  
Güzel



Davide  
Lanti



Marco  
Montali



Alessandro  
Mosca



Mariano  
Rodriguez  
Muro



Guohui  
Xiao

Technion  
Haifa



Avigdor  
Gal



Roei  
Shraga

Birkbeck  
College  
London



Roman  
Kontchakov



Vladislav  
Ryzhikov

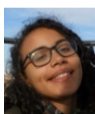


Michael  
Zakharyashev

Ontopic  
s.r.l.



Benjamin  
Cogrel



Sarah  
Komla Ebri

U. Roma  
"La  
Sapienza"



Giuseppe  
De Giacomo



Domenico  
Lembo



Maurizio  
Lenzerini



Antonella  
Poggi



Riccardo  
Rosati

# References I

- [1] Guohui Xiao, Diego C., Roman Kontchakov, Domenico Lembo, Antonella Poggi, Riccardo Rosati & Michael Zakharyashev. “Ontology-Based Data Access: A Survey”. In: *Proc. of the 27th Int. Joint Conf. on Artificial Intelligence (IJCAI)*. IJCAI Org., 2018, pp. 5511–5519. doi: 10.24963/ijcai.2018/777.
- [2] Boris Motik, Bernardo Cuenca Grau, Ian Horrocks, Zhe Wu, Achille Fokoue & Carsten Lutz. *OWL 2 Web Ontology Language Profiles (Second Edition)*. W3C Recommendation. Available at <http://www.w3.org/TR/owl2-profiles/>. World Wide Web Consortium, Dec. 2012.
- [3] Diego C., Giuseppe De Giacomo, Domenico Lembo, Maurizio Lenzerini & Riccardo Rosati. “Tractable Reasoning and Efficient Query Answering in Description Logics: The *DL-Lite* Family”. In: *J. of Automated Reasoning* 39.3 (2007), pp. 385–429. doi: 10.1007/s10817-007-9078-x.
- [4] Maurizio Lenzerini & Paolo Nibili. “On the Satisfiability of Dependency Constraints in Entity-Relationship Schemata”. In: *Information Systems* 15.4 (1990), pp. 453–461.
- [5] Sonia Bergamaschi & Claudio Sartori. “On Taxonomic Reasoning in Conceptual Design”. In: *ACM Trans. on Database Systems* 17.3 (1992), pp. 385–422.

## References II

- [6] Alexander Borgida. “Description Logics in Data Management”. In: *IEEE Trans. on Knowledge and Data Engineering* 7.5 (1995), pp. 671–682.
- [7] Diego C., Maurizio Lenzerini & Daniele Nardi. “Unifying Class-Based Representation Formalisms”. In: *J. of Artificial Intelligence Research* 11 (1999), pp. 199–240.
- [8] Alexander Borgida & Ronald J. Brachman. “Conceptual Modeling with Description Logics”. In: *The Description Logic Handbook: Theory, Implementation and Applications*. Ed. by Franz Baader, Diego C., Deborah McGuinness, Daniele Nardi & Peter F. Patel-Schneider. Cambridge University Press, 2003. Chap. 10, pp. 349–372.
- [9] Daniela Berardi, Diego C. & Giuseppe De Giacomo. “Reasoning on UML Class Diagrams”. In: *Artificial Intelligence* 168.1–2 (2005), pp. 70–118.
- [10] Anna Queralt, Alessandro Artale, Diego C. & Ernest Teniente. “OCL-Lite: Finite Reasoning on UML/OCL Conceptual Schemas”. In: *Data and Knowledge Engineering* 73 (2012), pp. 1–22. doi: 10.1016/j.datak.2011.09.004.

# References III

- [11] Antonella Poggi, Domenico Lembo, Diego C., Giuseppe De Giacomo, Maurizio Lenzerini & Riccardo Rosati. “Linking Data to Ontologies”. In: *J. on Data Semantics* 10 (2008), pp. 133–173. doi: 10.1007/978-3-540-77688-8\_5.
- [12] Alessandro Artale, Diego C., Roman Kontchakov & Michael Zakharyashev. *The DL-Lite Family and Relations*. Tech. rep. BBKCS-09-03. Available at <http://www.dcs.bbk.ac.uk/research/techreps/2009/bbkcs-09-03.pdf>. London: School of Computer Science and Information Systems, Birbeck College, 2009.
- [13] Diego C., Giuseppe De Giacomo, Domenico Lembo, Maurizio Lenzerini & Riccardo Rosati. “Data Complexity of Query Answering in Description Logics”. In: *Proc. of the 10th Int. Conf. on Principles of Knowledge Representation and Reasoning (KR)*. 2006, pp. 260–270.
- [14] Diego C., Giuseppe De Giacomo, Domenico Lembo, Maurizio Lenzerini & Riccardo Rosati. “Data Complexity of Query Answering in Description Logics”. In: *Artificial Intelligence* 195 (2013), pp. 335–360. doi: 10.1016/j.artint.2012.10.003.

# References IV

- [15] Stanislav Kikot, Roman Kontchakov & Michael Zakharyashev. “Conjunctive Query Answering with OWL 2 QL”. In: *Proc. of the 13th Int. Conf. on Principles of Knowledge Representation and Reasoning (KR)*. 2012, pp. 275–285.
- [16] Mariano Rodriguez-Muro, Roman Kontchakov & Michael Zakharyashev. “Ontology-Based Data Access: Ontop of Databases”. In: *Proc. of the 12th Int. Semantic Web Conf. (ISWC)*. Vol. 8218. Lecture Notes in Computer Science. Springer, 2013, pp. 558–573. doi: 10.1007/978-3-642-41335-3\_35.
- [17] Roman Kontchakov, Martin Rezk, Mariano Rodriguez-Muro, Guohui Xiao & Michael Zakharyashev. “Answering SPARQL Queries over Databases under OWL 2 QL Entailment Regime”. In: *Proc. of the 13th Int. Semantic Web Conf. (ISWC)*. Vol. 8796. Lecture Notes in Computer Science. Springer, 2014, pp. 552–567. doi: 10.1007/978-3-319-11964-9\_35.

# References V

- [18] Roman Kontchakov & Michael Zakharyashev. “An Introduction to Description Logics and Query Rewriting”. In: *Reasoning Web: Reasoning on the Web in the Big Data Era – 10th Int. Summer School Tutorial Lectures (RW)*. Vol. 8714. Lecture Notes in Computer Science. Springer, 2014, pp. 195–244. doi: 10.1007/978-3-319-10587-1\_5.
- [19] Marcelo Arenas, Alexandre Bertails, Eric Prud’hommeaux & Juan Sequeda. *A Direct Mapping of Relational Data to RDF*. W3C Recommendation. Available at <http://www.w3.org/TR/rdb-direct-mapping/>. World Wide Web Consortium, Sept. 2012.
- [20] Diego C., Benjamin Cogrel, Sarah Komla-Ebri, Roman Kontchakov, Davide Lanti, Martin Rezk, Mariano Rodriguez-Muro & Guohui Xiao. “Ontop: Answering SPARQL Queries over Relational Databases”. In: *Semantic Web J.* 8.3 (2017), pp. 471–487. doi: 10.3233/SW-160217.
- [21] Guohui Xiao, Davide Lanti, Roman Kontchakov, Sarah Komla-Ebri, Elem Güzel-Kalayci, Linfang Ding, Julien Corman, Benjamin Cogrel, Diego C. & Elena Botoeva. “The Virtual Knowledge Graph System Ontop”. In: *Proc. of the 19th Int. Semantic Web Conf. (ISWC)*. Vol. 12507. Lecture Notes in Computer Science. Springer, 2020, pp. 259–277. doi: 10.1007/978-3-030-62466-8\_17.