

Mapping Patterns for Virtual Knowledge Graphs

Diego Calvanese

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

Department of Computing Science
Umeå University, Sweden



International Conference on Conceptual Structures (ICCS 2021)
22 September 2021 – Virtual

Outline of the presentation

- 1 Challenges in Data Access
- 2 Knowledge Graphs for Data Access
- 3 Designing a (V)KG System
- 4 Mapping Patterns
- 5 The ADaMaP Algorithm
- 6 Conclusions

Challenges in the Big Data era

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

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



By 2016, it is projected there will be

18.9 BILLION NETWORK CONNECTIONS

— almost 2.5 connections per person on earth



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

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

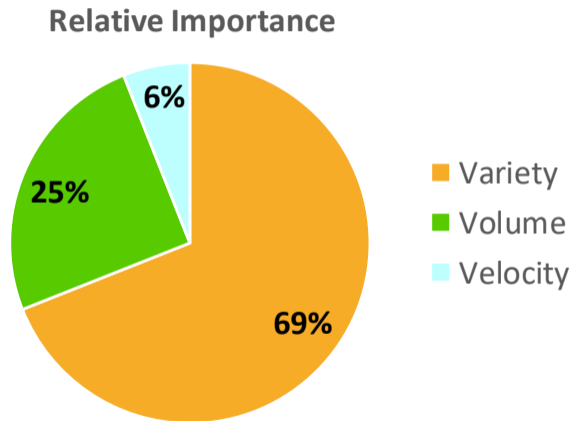
Poor data quality costs the US economy around

\$3.1 TRILLION A YEAR



Variety, not volume, is driving Big Data initiatives

MIT Sloan Management Review (28 March 2016)



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

The problem of accessing relevant data!

- Every day, **huge amounts of data** are produced by various actors.
- **Effective access to such data** is crucial to make value out of them.

Key issue

Queries are hard to formulate and navigating data through conventional tools is becoming an increasingly hard, and mostly ineffective, challenge.

Accessing heterogeneous data at Statoil

(Use case from FP7 IP Optique)

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



The humongous queries over *Slegge*

A typical query over *Slegge*, as formulated by the end-user in natural language

“Give me the names and chronostratigraphic zones found in the available wellbores as well as the top depths of the intervals they were found in; the values should be in the standard units and from standard reference points (depth in metres along the drill string). Also return all the lithostratigraphic zones from depths overlapping with the depths at which the chronostratigraphic zones were found.”

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

- Main table for wellbores has 38 columns (with cryptic names).
- The information about chronostratigraphic and lithostratigraphic zones and corresponding depths is stored in 8 additional tables.
- The resulting query contains a 24-table join.

The humongous queries over Slegge

```

SELECT
QVIEW1."IDENTIFIER" AS "wellbore",
QVIEW2."STRAT_UNIT_IDENTIFIER" AS "chronostrat_unit",
QVIEW15."STRAT_UNIT_IDENTIFIER" AS "lithostrat_unit",
QVIEW2."STRAT_ZONE_ENTRY_MD" AS "top_md_m"
FROM
"WELLBORE" QVIEW1,
"STRATIGRAPHIC_ZONE" QVIEW2,
"ROCK_FEATURE" QVIEW3,
"COMPONENT_MATERIAL" QVIEW4,
"DATA_COLLECTION" QVIEW5,
"DATA_COLLECTION_CONTENT" QVIEW6,
"ROCK_FEATURE" QVIEW7,
"MATERIAL_CLASS" QVIEW8,
"CLASSIFICATION_SYSTEM" QVIEW9,
"DATA_COLLECTION_CONTENT" QVIEW10,
"MATERIAL_CLASSIFICATION" QVIEW11,
"STRATIGRAPHIC_ZONE" QVIEW15,
"ROCK_FEATURE" QVIEW16,
"COMPONENT_MATERIAL" QVIEW17,
"DATA_COLLECTION" QVIEW18,
"DATA_COLLECTION_CONTENT" QVIEW19,
"ROCK_FEATURE" QVIEW20,
"MATERIAL_CLASS" QVIEW21,
"CLASSIFICATION_SYSTEM" QVIEW22,
"DATA_COLLECTION_CONTENT" QVIEW23,
"MATERIAL_CLASSIFICATION" QVIEW24
WHERE
QVIEW1."REF_EXISTENCE_KIND" = 'actual' AND
QVIEW1."IDENTIFIER" IS NOT NULL AND
QVIEW1."IDENTIFIER" = QVIEW2."WELLBORE" AND
QVIEW2."STRAT_ZONE_DEPTH_UOM" = 'm' AND
QVIEW2."STRAT_COLUMN_IDENTIFIER" IS NOT NULL AND
QVIEW2."STRAT_INTERP_VERSION" IS NOT NULL AND
QVIEW2."STRAT_ZONE_IDENTIFIER" IS NOT NULL AND
QVIEW2."STRAT_UNIT_IDENTIFIER" IS NOT NULL AND
QVIEW2."STRAT_UNIT_IDENTIFIER" = QVIEW3."DESCRIPTION" AND
QVIEW4."ENTITY_TYPE_NAME" = 'COMPONENT_MATERIAL' AND
QVIEW3."ROCK_FEATURE_S" = QVIEW4."INCORPORATE_S" AND
QVIEW3."ROCK_FEATURE_S" = QVIEW6."COLLECTION_PART_S" AND
QVIEW5."DATA_COLLECTION_S" = QVIEW6."PART_OF_S" AND
QVIEW2."STRAT_COLUMN_IDENTIFIER" = QVIEW5."NAME" AND
QVIEW5."REF_DATA_COLLECTION_TYPE" = 'stratigraphic hierarchy' AND
QVIEW9."KIND" = 'chronostratigraphy' AND
QVIEW8."CLASSIFICATION_SYSTEM" = QVIEW9."NAME" AND
QVIEW7."ROCK_FEATURE_S" = QVIEW10."COLLECTION_PART_S" AND
QVIEW5."DATA_COLLECTION_S" = QVIEW10."PART_OF_S" AND
QVIEW7."ROCK_FEATURE_S" = QVIEW11."MATERIAL_S" AND
QVIEW8."MATERIAL_CLASS_S" = QVIEW11."MATERIAL_CLASS_S" AND
QVIEW2."STRAT_ZONE_ENTRY_MD" IS NOT NULL AND
QVIEW1."IDENTIFIER" = QVIEW15."WELLBORE" AND
QVIEW15."STRAT_ZONE_DEPTH_UOM" = 'm' AND
(((QVIEW15."STRAT_ZONE_EXIT_MD" >=
QVIEW2."STRAT_ZONE_ENTRY_MD") AND
(QVIEW15."STRAT_ZONE_ENTRY_MD" <=
QVIEW2."STRAT_ZONE_EXIT_MD")) OR
((QVIEW2."STRAT_ZONE_EXIT_MD" >=
QVIEW15."STRAT_ZONE_ENTRY_MD") AND
(QVIEW2."STRAT_ZONE_ENTRY_MD" <=
QVIEW15."STRAT_ZONE_EXIT_MD"))
)) AND
QVIEW15."STRAT_COLUMN_IDENTIFIER" IS NOT NULL AND
QVIEW15."STRAT_INTERP_VERSION" IS NOT NULL AND
QVIEW15."STRAT_ZONE_IDENTIFIER" IS NOT NULL AND
QVIEW15."STRAT_UNIT_IDENTIFIER" IS NOT NULL AND
QVIEW15."STRAT_UNIT_IDENTIFIER" = QVIEW16."DESCRIPTION" AND
QVIEW17."ENTITY_TYPE_NAME" = 'COMPONENT_MATERIAL' AND
QVIEW16."ROCK_FEATURE_S" = QVIEW17."INCORPORATE_S" AND
QVIEW15."STRAT_COLUMN_IDENTIFIER" = QVIEW18."NAME" AND
QVIEW18."REF_DATA_COLLECTION_TYPE" = 'stratigraphic hierarchy' AND
QVIEW16."ROCK_FEATURE_S" = QVIEW19."COLLECTION_PART_S" AND
QVIEW18."DATA_COLLECTION_S" = QVIEW19."PART_OF_S" AND
QVIEW22."KIND" = 'lithostratigraphy' AND
QVIEW21."CLASSIFICATION_SYSTEM" = QVIEW22."NAME" AND
QVIEW20."ROCK_FEATURE_S" = QVIEW23."COLLECTION_PART_S" AND
QVIEW18."DATA_COLLECTION_S" = QVIEW23."PART_OF_S" AND
QVIEW20."ROCK_FEATURE_S" = QVIEW24."MATERIAL_S" AND
QVIEW21."MATERIAL_CLASS_S" = QVIEW24."MATERIAL_CLASS_S"
QVIEW21."MATERIAL_CLASS_S" = QVIEW24."MATERIAL_CLASS_S"

```



The humongous queries over *Slegge*

```
SELECT
QVIEW1."IDENTIFIER" AS "wellbore",
QVIEW2."STRAT_UNIT_IDENTIFIER" AS "chronostrat_unit",
QVIEW15."STRAT_UNIT_IDENTIFIER" AS "lithostrat_unit",
QVIEW2."STRAT_ZONE_ENTRY_MD" AS "top_md_m"
FROM
"WELLBORE" QVIEW1,
"STRATIGRAPHIC_ZONE" QVIEW2,
```

```
QVIEW3."ROCK_FEATURE_S" = QVIEW6."COLLECTION_PART_S" AND
QVIEW5."DATA_COLLECTION_S" = QVIEW6."PART_OF_S" AND
QVIEW2."STRAT_COLUMN_IDENTIFIER" = QVIEW5."NAME" AND
QVIEW5."REF_DATA_COLLECTION_TYPE" = 'stratigraphic hierarchy' AND
QVIEW9."KIND" = 'chronostratigraphy' AND
QVIEW8."CLASSIFICATION_SYSTEM" = QVIEW9."NAME" AND
QVIEW7."ROCK_FEATURE_S" = QVIEW10."COLLECTION_PART_S" AND
QVIEW5."DATA_COLLECTION_S" = QVIEW10."PART_OF_S" AND
```

The formulation of such a query can require days of interaction between the IT personnel (database experts) and the end-users (domain experts).

```
"DATA_COLLECTION" QVIEW18,
"DATA_COLLECTION_CONTENT" QVIEW19,
"ROCK_FEATURE" QVIEW20,
"MATERIAL_CLASS" QVIEW21
```

```
QVIEW15."STRAT_ZONE_EXIT_MD")
)) AND
QVIEW15."STRAT_COLUMN_IDENTIFIER" IS NOT NULL AND
QVIEW15."STRAT_TYPERD_HERETON" IS NOT NULL AND
```

This is also very costly!
Equinor loses 50.000.000€ per year only due to this problem!!

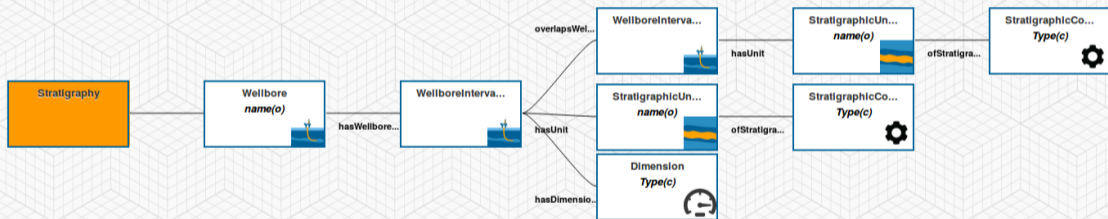
```
QVIEW4."ENTITY_TYPE_NAME" = 'COMPONENT_MATERIAL' AND
QVIEW3."ROCK_FEATURE_S" = QVIEW4."INCORPORATE_S" AND
```

```
QVIEW20."ROCK_FEATURE_S" = QVIEW24."MATERIAL_S" AND
QVIEW21."MATERIAL_CLASS_S" = QVIEW24."MATERIAL_CLASS_S"
QVIEW21."MATERIAL_CLASS_S" = QVIEW24."MATERIAL_CLASS_S"
```

A not-so humongous equivalent query

The FP7 IP Optique solution

```
SELECT ?wellbore ?chronostrat_unit ?top_md_m ?lithostrat_unit
WHERE {
  ?w a :Wellbore ;
     :name ?wellbore ;
```



```
?litho_intv :hasUnit ?lu .
?lu :name ?lithostrat_unit ;
     :ofStratigraphicColumn [ a :LithoStratigraphicColumn ] .
}
```

Uninformative names from SDSS¹

(Use case from H2020 RIA INODE)

Retrieve all stars that are White Dwarfs

```
SELECT s1.objID AS objid
FROM skyserverv3_correct.star AS S1
WHERE (S1.u - S1.g) < 0.4 AND (S1.g - S1.r) < 0.7 AND
      (S1.r - S1.i) > 0.4 AND (S1.i - S1.z) > 0.4
```

Retrieve all stars that are White Dwarfs

```
SELECT ?x WHERE { ?x a :WhiteDwarf }
```

¹The Sloan Digital Sky Survey, <https://www.sdss.org/>

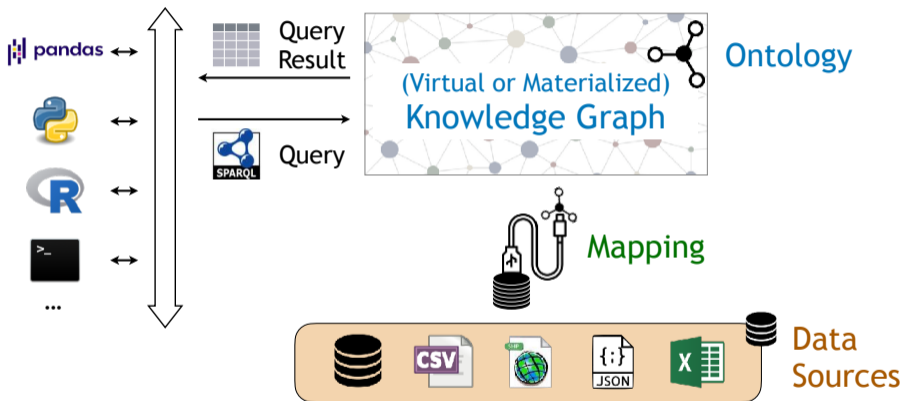
Outline

- 1 Challenges in Data Access
- 2 Knowledge Graphs for Data Access
- 3 Designing a (V)KG System
- 4 Mapping Patterns
- 5 The ADaMaP Algorithm
- 6 Conclusions

Outline

- 1 Challenges in Data Access
- 2 Knowledge Graphs for Data Access**
- 3 Designing a (V)KG System
- 4 Mapping Patterns
- 5 The ADaMaP Algorithm
- 6 Conclusions

(Virtual) Knowledge Graphs for data access



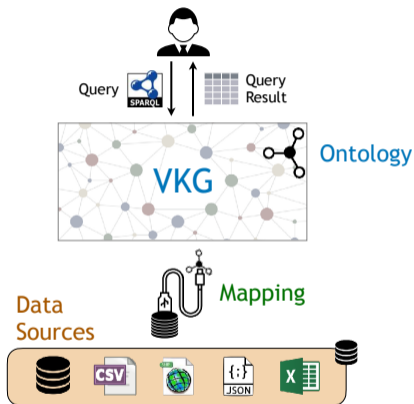
Components of the VKG framework

Which are the right languages for the components of the VKG framework?

We need to consider the **tradeoff between expressive power and efficiency**, where efficiency with respect to the data is the key aspect to consider.

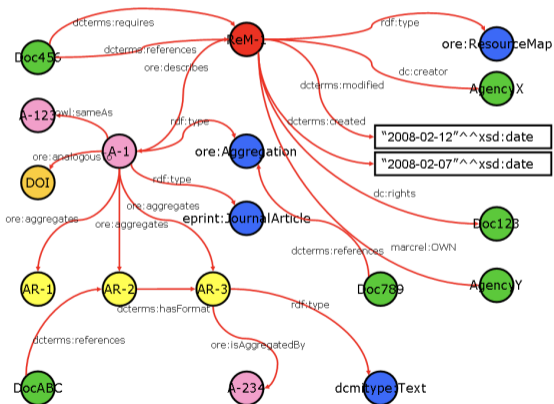
The W3C has standardized languages suitable for VKGs:

- 1 Knowledge graph: expressed in **RDF** [W3C Rec. 2014]
- 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]



RDF – Data is represented as a graph

The graph consists of a set of **subject-predicate-object triples** relating objects to other objects or values and to classes.



Class membership:

`<A-1> rdf:type :JournalArticle .`

Data property:

`<ReM-1> :created "2008-02-07" .`

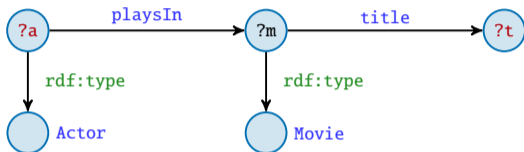
Object property:

`<A-1> ore:describes <ReM-1> .`

SPARQL query language

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

```
SELECT ?a ?t
WHERE { ?a rdf:type Actor .
        ?a playsIn ?m .
        ?m rdf:type Movie .
        ?m 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)
- ...

The OWL 2 QL ontology language

- **OWL 2 QL** is one of the three standard sub-languages of the very expressive standard ontology language OWL 2. [W3C Rec. 2012]
- It is considered a lightweight ontology language:
 - controlled expressive power
 - efficient inference
- Optimized for accessing large amounts of data
 - Queries over the ontology can be rewritten into SQL queries over the underlying relational database (**First-order rewritability**).
 - Logical consistency of ontology and data can also be checked by executing SQL queries over the underlying database.

Constructs of OWL 2 QL

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

- Subclass assertions $\text{MovieActor} \sqsubseteq \text{Actor}$
- Class disjointness $\text{Actor} \sqsubseteq \neg \text{Movie}$
- Domain of a property $\exists \text{actsIn} \sqsubseteq \text{MovieActor}$
- Range of a property $\exists \text{actsIn}^- \sqsubseteq \text{Movie}$
- Mandatory participation to a property $\text{MovieActor} \sqsubseteq \exists \text{actsIn}$
- Subproperty assertions $\text{actsIn} \sqsubseteq \text{playsIn}$
- Inverse properties $\text{actsIn} \equiv \text{hasActor}^-$

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 [Berardi, C. & De Giacomo 2005; Bergamaschi & Sartori 1992; Borgida 1995; C., Lenzerini & Nardi 1999; Lenzerini & Nobili 1990; Queralt et al. 2012].

MovieActor \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

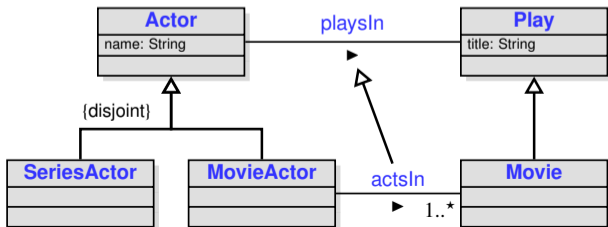
disjointness

domain

range

mandatory participation

sub-association



In fact, to visualize an OWL 2 QL ontology, we can use standard UML class diagrams.

Use of mappings

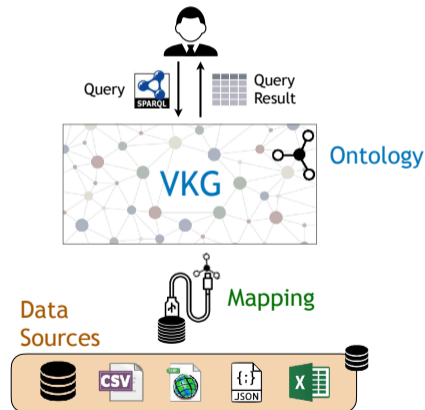
In the (V)KG framework, the **mapping** encodes how the **data in the sources** should be used to create the **(virtual) knowledge graph**, which is formulated in the vocabulary of the **ontology**.

(V)KG defined from the **mapping** and the **data**.

- Queries are answered with respect to the **ontology** and the data of the **(V)KG**.
- Such data is defined in terms of the **data sources** and the **mapping**.

In the virtual setting:

- The data of the **VKG** is not materialized.
- Instead, the information in the **ontology** and the **mapping** is used to translate queries over the **ontology** into queries formulated over the **sources**.
- The graph is **always up to date** wrt the data sources.



Mapping language

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

$$\begin{aligned}\Phi(\vec{x}) &\rightsquigarrow C(\mathbf{iri}(\vec{x})) \\ \Phi(\vec{x}) &\rightsquigarrow p(\mathbf{iri}_1(\vec{x}), \mathbf{iri}_2(\vec{x}))\end{aligned}$$

- $\Phi(\vec{x})$ is the **source query** in SQL.
- The **right hand side** is the **target**, consisting of a triple pattern involving an ontology class C or a (data or object) property p , and making use of the answer variables \vec{x} of the SQL query.

Intuition behind the mapping

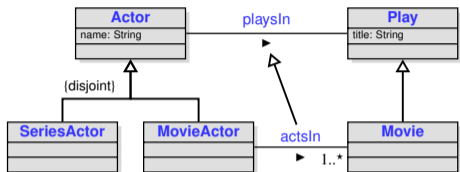
The **answers** returned by the **SQL query** in the left-hand side are used to create the **objects** (and values) that populate the **class / property** in the right-hand side.

Impedance mismatch: values in the DB vs. objects in the knowledge graph

In the **target**, we make use of **iri-templates** of the form $\mathbf{iri}(\vec{x})$, which transform database values into object identifiers (IRIs) or literals.

Mapping language – Example

Ontology \mathcal{O} :



Database \mathcal{D} :

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

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

Mapping \mathcal{M} :

m_1 : **SELECT** *mcode*, *mtitle* **FROM** MOVIE
WHERE *type* = "m"

\rightsquigarrow **Movie**($\text{iri}_m(\text{mcode})$)
title($\text{iri}_m(\text{mcode})$, *mtitle*)

m_2 : **SELECT** *M.mcode*, *A.acode* **FROM** MOVIE *M*, ACTOR *A*
WHERE *M.mcode* = *A.pcode* **AND** *M.type* = "m"

\rightsquigarrow **actsIn**($\text{iri}_a(\text{acode})$, $\text{iri}_m(\text{mcode})$)

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

Movie($\text{iri}_m(5118)$)

title($\text{iri}_m(5118)$, "The Matrix")

Movie($\text{iri}_m(2281)$)

title($\text{iri}_m(2281)$, "Blade Runner")

actsIn($\text{iri}_a(438)$, $\text{iri}_m(5118)$)

actsIn($\text{iri}_a(572)$, $\text{iri}_m(5118)$)

actsIn($\text{iri}_a(271)$, $\text{iri}_m(2281)$)

Outline

- 1 Challenges in Data Access
- 2 Knowledge Graphs for Data Access
- 3 Designing a (V)KG System**
- 4 Mapping Patterns
- 5 The ADaMaP Algorithm
- 6 Conclusions

Who provides the ontology?

- Designing an ontology is not an easy task.
- In many domains (e.g., the biomedical one) ontologies are developed independently by trained experts and already available to be re-used.
- Having “standardized ontologies” enables interoperability across different data sources.
- However, ontology design is a well investigated task, and methodologies and supporting tools are readily available. See, e.g.,
 - the series of *Workshops on Ontology Design Patterns* [<http://ontologydesignpatterns.org/>];
 - the OntoClean methodology for ontology analysis based on formal, domain-independent properties of classes [Guarino & Welty 2009].

Who provides the mappings?

(V)KG mappings:

- Map complex queries to complex queries – cf. GLAV relational mappings [Lenzerini 2002].
- Overcome the abstraction mismatch between relational data and target ontology.
- Are inherently more sophisticated than mappings for schema matching [Rahm & Bernstein 2001] and ontology matching [Euzenat & Shvaiko 2007].

As a consequence:

- Management of VKG mappings is an essentially manual effort that is **labor-intensive** and **error-prone**.
- Requires highly-skilled professionals [Spanos, Stavrou & Mitrou 2012].
- Writing mappings is challenging in terms of semantics, correctness, and performance.

Designing and managing mappings is the most critical bottleneck for the adoption of the VKG approach.

Who provides the mapping?

Writing mappings manually is a **time-consuming** and **error-prone** task.

Who provides the mapping?

Who provides the mapping?

The screenshot shows the Protege software interface with the following components:

- Top Bar:** cbio (http://purl.org/cbio/0.3) | [/home/tir/Seafle/Study/Papering/myPresentations/2021/INODE-Review-tutorial/protége-onto/oncmx_v0_3.owl]
- Menu Bar:** File Edit View Reasoner Tools Refactor Window Ontop Help
- Active Ontology:** cbio (http://purl.org/cbio/0.3)
- Class Hierarchy:** owl:Thing
- Mapping Editor:** Mapping Assistant - BETA
- Code Editor:**

```
SELECT distinct g.gene_symbol, g.ensembl_gene_id FROM xref_gene_ensembl as g WHERE g.appealid = '9606'  
  
disease_mutation  
oncmx:DM (id) a: DiseaseMutation ; terms:source {data_source}^^xsd:anyURI ; falds:reference obo:(chromosome_id) ; falds:location  
oncmx:LOCATION_CHROMOSOME(id)-(chromosome_pos) ; hasSequenceAlteration oncmx:SEQ_ALT(id)-(peptide_id), oncmx:SEQ_ALT(id)-(uniprotkb_ac) ;  
oncmx:SEQ_ALT(id)-(gene_symbol) ; hasTargetDisease obo:DOID({doid}) ; mutationFrequency (mutation_freq)^^xsd:nonNegativeInteger ;  
oncmx:LOCATION_CHROMOSOME(id)-(chromosome_pos) a falds:ExactPosition ; falds:position (chromosome_pos)^^xsd:integer ; gene:hasSequencedUnit  
<https://identifiers.org/hgnc:symbol:{gene_symbol}> ;  
  
SELECT id, CASE chromosome_id WHEN '1' THEN 'NCIT_C13204' WHEN '2' THEN 'NCIT_C13215' WHEN '3' THEN 'NCIT_C13219' WHEN '4' THEN 'NCIT_C13220' WHEN '5' THEN  
'NCIT_C13221' WHEN '6' THEN 'NCIT_C13222' WHEN '7' THEN 'NCIT_C13223' WHEN '8' THEN 'NCIT_C13224' WHEN '9' THEN 'NCIT_C13225' WHEN '10' THEN 'NCIT_C13205'  
WHEN '11' THEN 'NCIT_C13206' WHEN '12' THEN 'NCIT_C13207' WHEN '13' THEN 'NCIT_C13208' WHEN '14' THEN 'NCIT_C13209' WHEN '15' THEN 'NCIT_C13210' WHEN  
'16' THEN 'NCIT_C13211' WHEN '17' THEN 'NCIT_C13212' WHEN '18' THEN 'NCIT_C13213' WHEN '19' THEN 'NCIT_C13214' WHEN '20' THEN 'NCIT_C13216' WHEN '21'  
THEN 'NCIT_C13217' WHEN '22' THEN 'NCIT_C13218' WHEN 'X' THEN 'NCIT_C13285' WHEN 'Y' THEN 'NCIT_C13286' END AS chromosome_id, chromosome_pos, cds_pos,  
aa_pos,uniprotkb, mutation_freq, CASE data_source WHEN 'logs' THEN 'https://logs.org' WHEN 'cosmic' THEN 'https://cosmic.sanger.ac.uk/cosmic' WHEN 'toga' THEN  
'https://www.ncbi.nlm.nih.gov/clinvar' END AS data_source, doid, peptide_id, en_ensembl_transcript_id, np_uniprotkb_ac, gene_symbol FROM disease_mutation as dm  
join map_protein_disease_putation as mp on mp.ensembl_transcript_id = dm.ensembl_transcript_id left join xref_gene_uniprot as hugo on  
hugo.uniprotkb_ac=mp.uniprotkb_ac  
  
sequence_alteration  
oncmx:SEQ_ALT(id)-(peptide_id) a obo:SO_0001059 ; falds:reference <https://identifiers.org/ensembl:{peptide_id}> ; falds:location  
oncmx:LOCATION_PROT(id)-(peptide_pos) ; alteredFrom obo:(ref_aa), sio:(ref_aa_SIO), obo:(ref_aa_SIO) ; alteredTo obo:(alt_aa), sio:(alt_aa_SIO), obo:(alt_aa_X)  
; oncmx:LOCATION_PROT(id)-(peptide_pos) a falds:ExactPosition ; falds:position (peptide_pos)^^xsd:integer ; oncmx:SEQ_ALT(id)-(gene_symbol) a obo:SO_0001059 ;  
falds:reference <http://purl.uniprot.org/uniprot/{uniprotkb_ac}> ; falds:location oncmx:LOCATION_PROT(id)-{aa_pos_uniprotkb} ; alteredFrom obo:(ref_aa),  
sio:(ref_aa_SIO), obo:(ref_aa_X) ; alteredTo obo:(alt_aa), sio:(alt_aa_SIO), obo:(alt_aa_X) ; oncmx:LOCATION_PROT(id)-{aa_pos_uniprotkb} a falds:ExactPosition ;  
falds:position (aa_pos_uniprotkb)^^xsd:integer ; oncmx:SEQ_ALT(id)-(gene_symbol) a obo:SO_0001059 ; falds:reference <https://identifiers.org/hgnc:symbol:{gene_symbol}> ;  
<https://identifiers.org/hgnc:symbol:{gene_symbol}> ; falds:location oncmx:LOCATION_GENE(id)-(cds_pos) ; alteredTo obo:(ref_nt) ; alteredTo obo:(alt_nt),  
oncmx:LOCATION_GENE(id)-(gene_symbol) a falds:ExactPosition ; falds:position (cds_pos)^^xsd:integer  
  
SELECT id, CASE ref_nt WHEN 'A' THEN 'CHEBI_16708' WHEN 'C' THEN 'CHEBI_16040' WHEN 'G' THEN 'CHEBI_16235' WHEN 'T' THEN 'CHEBI_17621' WHEN 'U' THEN  
'CHEBI_17568' END AS ref_nt, CASE alt_nt WHEN 'A' THEN 'CHEBI_16708' WHEN 'C' THEN 'CHEBI_16708' WHEN 'G' THEN 'CHEBI_16040' WHEN 'T' THEN 'CHEBI_16235' WHEN 'U'  
' THEN 'CHEBI_17568' END AS alt_nt, cds_pos, aa_pos,uniprotkb, CASE ref_aa WHEN 'A' THEN 'CHEBI_16449' WHEN 'C' THEN 'CHEBI_15356' WHEN 'D' THEN  
'CHEBI_22660' WHEN 'E' THEN 'CHEBI_18237' WHEN 'F' THEN 'CHEBI_28044' WHEN 'G' THEN 'CHEBI_15971' WHEN 'H' THEN 'CHEBI_15428' WHEN 'I' THEN 'CHEBI_17191'  
WHEN 'K' THEN 'CHEBI_25094' WHEN 'L' THEN 'CHEBI_25017' WHEN 'M' THEN 'CHEBI_16811' WHEN 'N' THEN 'CHEBI_22653' WHEN 'P' THEN 'CHEBI_17203' WHEN 'Q' THEN  
'CHEBI_18050' WHEN 'R' THEN 'CHEBI_16447' WHEN 'S' THEN 'CHEBI_16447' WHEN 'T' THEN 'CHEBI_16811' WHEN 'V' THEN 'CHEBI_14141' WHEN 'W' THEN  
'CHEBI_27897' WHEN 'Y' THEN 'CHEBI_18186' END AS ref_aa, CASE ref_aa WHEN 'A' THEN 'SIO_010448' END AS ref_aa_SIO, CASE ref_aa WHEN 'X' THEN 'ANY_OC000' END  
AS ref_aa_X, CASE alt_aa WHEN 'A' THEN 'CHEBI_16449' WHEN 'C' THEN 'CHEBI_15356' WHEN 'D' THEN 'CHEBI_22660' WHEN 'E' THEN 'CHEBI_18237' WHEN 'F' THEN  
'CHEBI_28044' WHEN 'G' THEN 'CHEBI_15428' WHEN 'H' THEN 'CHEBI_15971' WHEN 'I' THEN 'CHEBI_17191' WHEN 'K' THEN 'CHEBI_25094' WHEN 'L' THEN 'CHEBI_25017'  
WHEN 'M' THEN 'CHEBI_16811' WHEN 'N' THEN 'CHEBI_22653' WHEN 'P' THEN 'CHEBI_17203' WHEN 'Q' THEN 'CHEBI_18050' WHEN 'R' THEN 'CHEBI_16447' WHEN 'S' THEN  
'CHEBI_17115' WHEN 'T' THEN 'CHEBI_16957' WHEN 'V' THEN 'CHEBI_16414' WHEN 'W' THEN 'CHEBI_27897' WHEN 'Y' THEN 'CHEBI_18186' END AS alt_aa, CASE alt_aa  
WHEN 'U' THEN 'SIO_010448' END AS alt_aa_SIO, CASE alt_aa WHEN 'X' THEN 'ANY_OC000' END AS alt_aa_X, peptide_pos, mutation_freq, data_source, doid, peptide_id,  
en_ensembl_transcript_id, np_uniprotkb_ac, gene_symbol FROM disease_putation as dp join map_protein_disease_putation as mp on mp.ensembl_transcript_id =  
dp.ensembl_transcript_id left join xref_gene_uniprot as hugo on hugo.uniprotkb_ac=mp.uniprotkb_ac
```
- Left Panel:** Data property hierarchy, Object property hierarchy, Object property hierarchy: owl:topObjectProperty
- Bottom Bar:** Mapping size: 36 Search (any of):



Outline

- 1 Challenges in Data Access
- 2 Knowledge Graphs for Data Access
- 3 Designing a (V)KG System
- 4 Mapping Patterns**
- 5 The ADaMaP Algorithm
- 6 Conclusions

Mapping patterns

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

Several approaches and tools that deal with the problem of extracting a KG from a relational data source have been proposed, several of them based on mapping patterns.

However, to the best of our knowledge:

- There is no comprehensive approach for KG mapping patterns exploiting all of:
 - the relational schema with its constraints
 - extensional data stored in the DB
 - the domain knowledge that is encoded in ontology axioms
 - the conceptual schema at the basis of the relational schema
- only a few come with a systematic categorization of the mappings that they produce.
- None of them have drawn an explicit and precise connection between their outputs and conceptual modeling practices found in DB design.
- None of them attempts an analysis over real-world scenarios

Catalog of mapping patterns

We build 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

In specifying each pattern, we consider:

- 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.

For the moment, we do not fix what is given as input and what is produced as output, but we simply describe how the elements relate to each other, on a per-pattern basis.

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

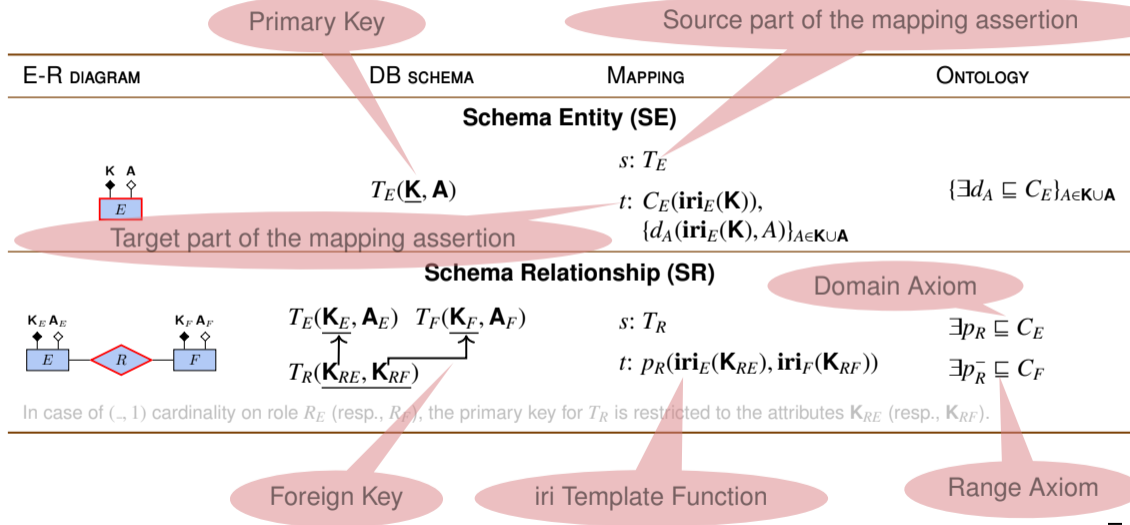
When defining the mapping patterns, we consider the traditional types of DB constraints:

- **Primary key constraint:** $T(\underline{\mathbf{K}}, \mathbf{A})$
- **Key constraint:** $\text{key}_T(\mathbf{K})$
- **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:

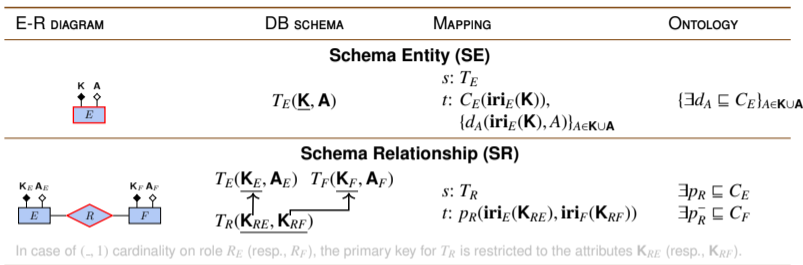


Note: We use normal font (e.g., A) for single attributes, and boldface for sets of attributes (e.g., \mathbf{A}).

Fragment of schema-driven patterns from [C., Gal, Lanti, et al. 2020]



Example 1: Schema Relationship Pattern



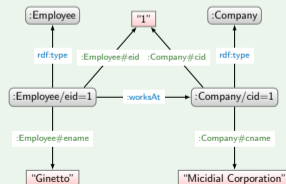
Conceptual



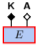

DB Schema



RDF (data only)



Example 2

E-R DIAGRAM	DB SCHEMA	MAPPING	ONTOLOGY
Schema Entity (SE)			
	$T_E(\underline{K}, \mathbf{A})$	$s: T_E$ $t: C_E(\text{iri}_E(\mathbf{K})),$ $\{d_A(\text{iri}_E(\mathbf{K}), A)\}_{A \in \mathbf{K} \cup \mathbf{A}}$	$\{\exists d_A \sqsubseteq C_E\}_{A \in \mathbf{K} \cup \mathbf{A}}$
Schema Relationship (SR)			
	$T_E(\underline{K}_E, \mathbf{A}_E)$ $T_F(\underline{K}_F, \mathbf{A}_F)$ $T_R(\underline{K}_{RE}, \mathbf{K}_{RF})$	$s: T_R$ $t: p_R(\text{iri}_E(\mathbf{K}_{RE}), \text{iri}_F(\mathbf{K}_{RF}))$	$\exists p_R \sqsubseteq C_E$ $\exists p_{\bar{R}} \sqsubseteq C_F$
<p>In case of $(_, 1)$ cardinality on role R_E (resp., R_F), the primary key for T_R is restricted to the attributes \mathbf{K}_{RE} (resp., \mathbf{K}_{RF}).</p>			

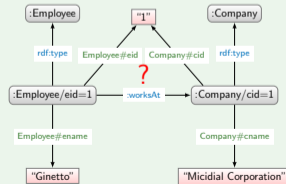
Conceptual



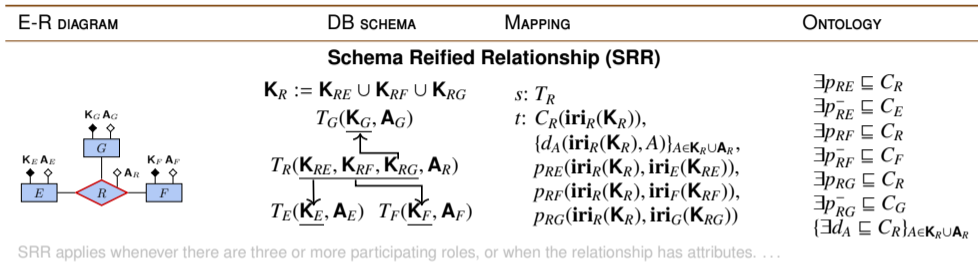
DB Schema



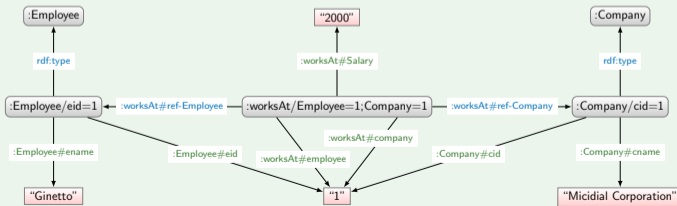
RDF (data only)



Example 2 – Revised



RDF (data only)



The other schema-driven patterns [C., Gal, Lanti, et al. 2020]

E-R DIAGRAM	DB SCHEMA	MAPPING	ONTOLOGY
Schema Relationship with Identifier Alignment (SRa)			
	$T_E(\underline{K_E}, \mathbf{A_E})$ $T_F(\underline{K_F}, \underline{U_F}, \mathbf{A_F})$ $T_R(\underline{K_{RE}}, \underline{U_{RF}}) \text{ key}_{R_F}(\underline{U_F})$	$s: T_R \bowtie_{U_{RF}=U_F} T_F$ $t: p_R(\text{iri}_E(K_{RE}), \text{iri}_F(K_F))$	$\exists p_R \sqsubseteq C_E$ $\exists p_R^- \sqsubseteq C_F$
In case of $(-, 1)$ cardinality on role R_E (resp., R_F), the primary key for T_R is restricted to the attributes K_{RE} (resp., U_{RF}).			
Schema Relationship with Merging (SRm)			
	$T_F(\underline{K_F}, \mathbf{A_F})$ $T_E(\underline{K_E}, \underline{K_{EF}}, \mathbf{A_E})$	$s: T_E$ $t: p_{EF}(\text{iri}_E(K_E), \text{iri}_F(K_{EF}))$	$\exists p_{EF} \sqsubseteq C_E$ $\exists p_{EF}^- \sqsubseteq C_F$
Schema Hierarchy (SH)			
	$T_E(\underline{K_E}, \mathbf{A_E})$ $T_F(\underline{K_{FE}}, \mathbf{A_F})$	$s: T_F$ $t: C_F(\text{iri}_E(K_{FE})),$ $\{d_A(\text{iri}_E(K_{FE}), A)\}_{A \in \mathbf{A}_F}$	$C_F \sqsubseteq C_E$ $\{\exists d_A \sqsubseteq C_F\}_{A \in \mathbf{A}_F}$
Schema Hierarchy with Identifier Alignment (SHa)			
	$T_E(\underline{K_E}, \mathbf{A_E}) \text{ key}_{T_F}(\underline{U_F})$ $T_F(\underline{K_F}, \underline{U_F}, \mathbf{A_F})$ <hr/> $T_E(\underline{K_E}, \mathbf{A_E}) \text{ key}_{V_F}(\underline{K_F})$ $V_F(\underline{K_F}, \underline{U_F}, \mathbf{A_F}) = T_F$	$s: T_F$ $t: C_F(\text{iri}_E(\underline{U_F})),$ $\{d_A(\text{iri}_E(\underline{U_F}), A)\}_{A \in \mathbf{K}_F \cup \mathbf{A}_F}$	$C_F \sqsubseteq C_E$ $\{\exists d_A \sqsubseteq C_F\}_{A \in \mathbf{K}_F \cup \mathbf{A}_F}$
In this pattern, the "alignment" is meant to align the primary identifier used in the child entity to the primary identifier used in the parent entity. ...			

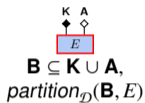
A “data”-driven pattern [C., Gal, Lanti, et al. 2020]

E-R DIAGRAM

DB SCHEMA

MAPPING

ONTOLOGY



Clustering Entity to Class (CE2C)

 $T_E(\mathbf{K}, \mathbf{A})$
 $unique_{T_E}(\mathbf{K})$
 $\mathbf{B} \subseteq \mathbf{K} \cup \mathbf{A}$
 $partition_D(\mathbf{B}, E)$
 $\{V_{E_v}(\mathbf{K}, \mathbf{A}) = \sigma_{\mathbf{B}=\mathbf{v}}(T_E)\}_{\mathbf{v} \in \pi_{\mathbf{B}}(T_E)}$
 $\{s : \sigma_{\mathbf{B}=\mathbf{v}}(T_E)$
 $t : C_E^{\mathbf{v}}(\mathbf{iri}_E(\mathbf{K}))\}_{\mathbf{v} \in \pi_{\mathbf{B}}(T_E)}$
 $\{C_E^{\mathbf{v}} \sqsubseteq C_E\}_{\mathbf{v} \in \pi_{\mathbf{B}}(T_E)}$

	eid	name	gender	
	1	Ginetto	M	
	2	Ginetta	F	
	3	Anna	F	
	4	Paolo	M	

Male

Female

Male

Employee

Design scenarios for VKG mapping patterns

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

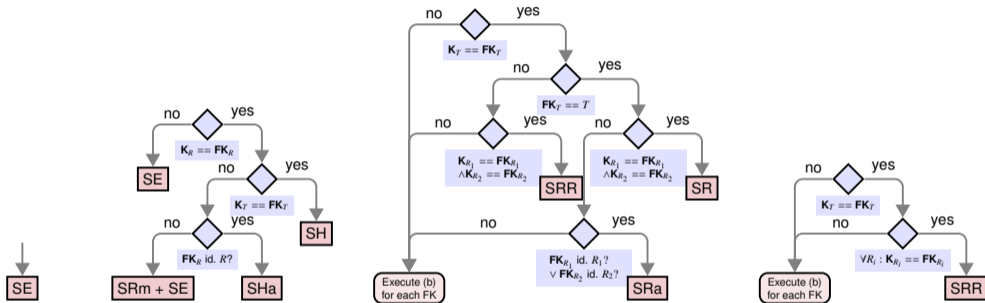
- 1 **Debugging of a VKG specification** that is already in place.
- 2 **Conceptual schema reverse engineering** for a DB that represents the domain of interest by using a given full VKG specification.
- 3 **Mapping bootstrapping** for a given DB and ontology that miss the mappings relating them.
- 4 **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.

Outline

- 1 Challenges in Data Access
- 2 Knowledge Graphs for Data Access
- 3 Designing a (V)KG System
- 4 Mapping Patterns
- 5 The ADaMaP Algorithm**
- 6 Conclusions

The ADaMaP algorithm [C., Gal, Haba, et al. 2021]

Idea: Classify each table according to a pattern



(a) $|fkeys| = 0$

(b) $|fkeys| = 1$

(c) $|fkeys| = 2$

(d) $|fkeys| \geq 3$

- **R:** Referenced Table
- **T:** Current Table
- **id:** "identifies"

- **K_T :** Key of current table
- **FK_T :** Considered foreign key of current table

Evaluation setting

To assess the applicability of our approach, we rely on two **non-trivial** and **real-world** scenarios:

- **NPD:**
 - Scenario built around the domain of oil and gas extraction.
 - Presents a **high number of mappings** (>1k).
 - **Most mappings** are **automatically generated** (Direct Mapping).
 - **Several** complex **manually-written** mappings as well.
- **Cordis:**
 - Domain of competitive research projects, provided by SIRIS Academic S.L., a consultancy company specializing in higher education and research.
 - The mappings were **manually-written**, and they amount to 120.

Coverage Analysis

In this analysis, we check how many mappings in the analyzed scenarios can be explained through mapping patterns.

Pattern	#usages	#mappings
SE	13	60
SR	3	3
SRm	3	3
SRR	1	16

Covered Mappings: 89 (out of 120)

CORDIS Coverage

Pattern	#usages	#mappings
SE	61	454
SRm	74	74
SRR	1	12
SH	3	132

Covered Mappings: 672 (out of 1173)

NPD Coverage

Mismatches analysis

- We compare the classification returned by ADaMAP to a classification manually verified by a human expert.
- M : the output of ADaMAP
- M^* : the manually verified classification

Precision

$$P_{M^*}(M) = \frac{|M \cap M^*|}{|M|}$$

Recall

$$R_{M^*}(M) = \frac{|M \cap M^*|}{|M^*|}$$

F1-measure

The harmonic mean of $P_{M^*}(M)$ and $R_{M^*}(M)$

Mismatch analysis (Cordis)

topics	CP						
subject_areas	CP						
projects	CP			CP			
project_topics		CP					
project_subject_areas		CP					
project_programmes		CP					
project_members	DP			DP	DN		
project_member_roles	CP						
project_erc_panels		DN		DP		DP	
programmes	CP			CP			
people	CP						
institutions	CP			CP			
funding_schemes	CP						
eu_territorial_units	DP						
erc_research_domains	CP						
erc_panels	CP			CP			
ec_framework_programs	CP						
countries	CP						
activity_types	CP						
	SE	SR	SRa	SRm	SRR	SH	SHa

- **DP/DN**: Discoordinated Positive/Negative
- **CP/CN**: Coordinated Positive/Negative
- **SE**: Schema Entity
- **SR/SRa/SRm/SRR**: Schema Relationship/with alignment/with merging/with reification
- **SH/SHa**: Schema Hierarchy/with alignment

$$P_{M^*}(M) = R_{M^*}(M) = F_{M^*}(M) = 0.8$$

- Overall, ADaMAP and the manual extraction have 20% of disagreements.
- All but one disagreement stem from the fact that multiple conceptual schemata can correspond to the same database schema.
- The algorithm cannot determine which of these equally valid choices is actually the one that was adopted by the human designer.

Mismatch analysis (NPD)

wbpoint	DN			DN		CP	
wellbore_shallow_all						CP	
wellbore_oil_sample	CP			CP			
wellbore_npdid_overview	CP						
wellbore_mud	CP			CP			
wellbore_formation_top	CP			CP			
wellbore_exploration_all				DP		DN	DP
wellbore_dst	CP			CP			
wellbore_document	CP			DP			
wellbore_development_all	DN			CP		DN	DP
wellbore_core_photo	CP			CP			
wellbore_core	CP			DP			
wellbore_coordinates	DN			DN		CP	
wellbore_casing_and_lot	CP			CP			
luf_petreg_message	CP			CP			
luf_petreg_licence_oper	DN			DP		DP	
luf_petreg_licence_licence					CP		
luf_petreg_licence	CP						
luf_owner_hst	CP			CP			
luf_operator_hst	CP			CP			
strat_litho_wellbore_core	CP			CP			
strat_litho_wellbore	CP			CP			
seis_acquisition_progress	CP			DN			
seis_acquisition_coordinates_inc_turnarea	CP			DN			
seis_acquisition	CP			DN			
seamultiline	CP			DN		DN	
seaaarea	CP			DN			
priareasplitbyblock	CP			CP			
priarea	CP			CP			
pipeline	CP			CP			
licence_transfer_hst	CP			CP			
licence_task	CP			CP			
licence_phase_hst	CP			CP			
licence_petreg_message	CP			CP			
licence_petreg_licence_oper	DN			CP		CP	
licence_petreg_licence_licence					CP		
licence_petreg_licence						CP	

licence_oper_hst	CP				CP		
licence_licencee_hst	CP				CP		
licence_area_poly_hst	CP				CP		
licence	CP						
fldarea	CP				CP		
field_reserves	DN				DN		CP
field_production_yearly	CP						
field_production_totait_ncs_year	DP						
field_production_totait_ncs_month	CP						
field_production_monthly	CP						
field_owner_hst	CP				CP		
field_operator_hst	CP				CP		
field_licencee_hst	CP				CP		
field_investment_yearly	CP				CP		
field_description	DP				DP		
field_activity_status_hst	CP				CP		
field	CP				DP		
fcipoint	DN				CP		DP
facility_moveable	CP				CP		
facility_fixed	CP						DN
dscarea	CP				CP		
discovery_reserves	CP				CP		
discovery	CP				CP		
company_reserves						DP	
company	DP				DN		
bsns_arr_area_transfer_hst	CP				CP		
bsns_arr_area_operator	DN				CP		CP
bsns_arr_area_licencee_hst	CP				CP		
bsns_arr_area_area_poly_hst	CP				CP		
bsns_arr_area	CP						
baaarea	CP				CP		
apaareanet	CP						
apaareagross	CP						
	SE	SR	SRa	SRm	SRR	SH	SHa

The algorithm and the manual analysis disagree on 35 instances, with 14 discoordinated positives and 21 discoordinated negatives.

In terms of precision (P), recall (R), and F1-measure (F), ADaMAP obtains the following results:

$$P_{M^*}(M) = 0.88, R_{M^*}(M) = 0.82, F_{M^*}(M) = 0.85$$

Outline

- 1 Challenges in Data Access
- 2 Knowledge Graphs for Data Access
- 3 Designing a (V)KG System
- 4 Mapping Patterns
- 5 The ADaMaP Algorithm
- 6 Conclusions**

Conclusions

- The design of mappings is a complex task, that currently is the major bottleneck in the wider adoption of the VKG paradigm for data access and integration.
- Mapping patterns are a promising approach for simplifying this complex task.
- In our work, we have identified a catalog of mapping patterns for the VKG framework.
- We have introduced ADaMAP, an algorithmic technique that extracts semantics from a relational data source, by automatically identifying how ontology mapping patterns are applied to fragments of its schema
- ADaMAP can be used to support the automatic generation of ontologies and mappings
- The patterns identified by ADaMAP provide a solid basis that can be manually improved by human experts
- The validation of ADaMAP in two significant real-world case studies confirms that the identified patterns by-and-large agree with those detected by a human expert

Future Work

ADaMAP comes with some limitations that should be tackled:

- For a given relational schema there are in general many possible combinations of mapping patterns that are, in principle, equally valid . . .
 - while ADaMAP only returns the “most typical” one.
- ADaMAP ignores data, however “data”-driven patterns² are also important . . .
 - especially in those scenarios where the DB schema is poorly structured or denormalized.

²As in C., Gal, Lanti, et al. 2020

Thank you!

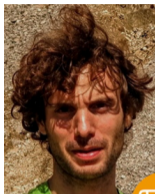
- E: calvanese@inf.unibz.it
- H: <http://www.inf.unibz.it/~calvanese/>

The logo for 'ontop' features the word in a lowercase, rounded, orange font. A thin horizontal line passes through the middle of the letters, creating a sense of depth or a shadow effect.The logo for 'ONTOPIC' is in a bold, uppercase, black font. The letters are stylized with a blocky, geometric appearance, featuring some internal cutouts or a stencil-like effect.

- *Ontop* website: <https://ontop-vkg.org/>
- Github: <http://github.com/ontop/ontop/>
- Facebook: <https://www.facebook.com/obdaontop/>
- Twitter: @ontop4obda
- *Ontopic* website: <https://ontopic.biz/>

A great thank you to all collaborators

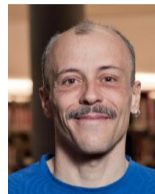
Unibz



Davide Lanti



Marco Montali



Alessandro Mosca

Technion
Haifa



Avigdor Gal



Roe Shraga

References I

- [1] 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.
- [2] Daniela Berardi, Diego C. & Giuseppe De Giacomo. “Reasoning on UML Class Diagrams”. In: *Artificial Intelligence* 168.1–2 (2005), pp. 70–118.
- [3] Sonia Bergamaschi & Claudio Sartori. “On Taxonomic Reasoning in Conceptual Design”. In: *ACM Trans. on Database Systems* 17.3 (1992), pp. 385–422.
- [4] Alexander Borgida. “Description Logics in Data Management”. In: *IEEE Trans. on Knowledge and Data Engineering* 7.5 (1995), pp. 671–682.
- [5] Diego C., Avigdor Gal, Naor Haba, Davide Lanti, Marco Montali, Alessandro Mosca & Roei Shraga. “ADaMaP: Automatic Alignment of Data Sources using Mapping Patterns”. In: *Proc. of the 33rd Int. Conf. on Advanced Information Systems Engineering (CAiSE 2021)*. Vol. 12751. Lecture Notes in Computer Science. Springer, 2021, pp. 193–209. doi: 10.1007/978-3-030-79382-1_12.

References II

- [6] Diego C., Avigdor Gal, Davide Lanti, Marco Montali, Alessandro Mosca & Roei Shraga. *Mapping Patterns for Virtual Knowledge Graphs*. CoRR Technical Report arXiv:2012.01917. arXiv.org e-Print archive, 2020.
- [7] Diego C., Maurizio Lenzerini & Daniele Nardi. “Unifying Class-Based Representation Formalisms”. In: *J. of Artificial Intelligence Research* 11 (1999), pp. 199–240.
- [8] Jérôme Euzenat & Pavel Shvaiko. *Ontology Matching*. Springer, 2007.
- [9] Nicola Guarino & Christopher A. Welty. “An Overview of OntoClean”. In: *Handbook on Ontologies*. Ed. by Steffen Staab & Rudi Studer. International Handbooks on Information Systems. Springer, 2009, pp. 201–220. doi: 10.1007/978-3-540-92673-3_9.
- [10] Maurizio Lenzerini. “Data Integration: A Theoretical Perspective.”. In: *Proc. of the 21st ACM Symp. on Principles of Database Systems (PODS)*. 2002, pp. 233–246. doi: 10.1145/543613.543644.
- [11] Maurizio Lenzerini & Paolo Nobili. “On the Satisfiability of Dependency Constraints in Entity-Relationship Schemata”. In: *Information Systems* 15.4 (1990), pp. 453–461.

References III

- [12] 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.
- [13] Erhard Rahm & Philip A. Bernstein. “A Survey of Approaches to Automatic Schema Matching”. In: *Very Large Database J.* 10.4 (2001), pp. 334–350.
- [14] Dimitrios-Emmanuel Spanos, Periklis Stavrou & Nikolas Mitrou. “Bringing relational databases into the Semantic Web: A survey”. In: *Semantic Web J.* 3.2 (2012), pp. 169–209.