

# INF3580/4580 – Semantic Technologies – Spring 2017

## Lecture 14: Introduction to Ontology-based Data Access (OBDA)

Ernesto Jimenez Ruiz

15th May 2017



DEPARTMENT OF  
INFORMATICS



UNIVERSITY OF  
OSLO

## Today's Plan

- 1 Exposing data as RDF
- 2 Data Access in Statoil: limitations and solutions
- 3 OBDA Ingredients
  - Overview
  - Ontology
  - Mappings
  - Query rewriting
  - Bootstrapping
  - Visual Query Formulation
  - Optique

## Outline

- 1 Exposing data as RDF
- 2 Data Access in Statoil: limitations and solutions
- 3 OBDA Ingredients
  - Overview
  - Ontology
  - Mappings
  - Query rewriting
  - Bootstrapping
  - Visual Query Formulation
  - Optique

## RDF

- Why URIs?
  - URIs naturally have a “global” scope, unique throughout the web.
  - URIs are also addresses.
  - “A web of data.”
- Why triples?
  - Any information format can be transformed to triples.
  - Relationships are made explicit and are elements in their own right

## RDF on the web: Where is it?

- In files:
  - In some serialisation format: XML/RDF, Turtle, ...
  - Typically small RDF graphs, i.e., max. a few 100 triples, e.g.,
    - Vocabularies: <http://xmlns.com/foaf/spec/index.rdf>.
    - Tiny datasets: <http://folk.uio.no/martingi/foaf.rdf>.
- “Behind” *SPARQL endpoints*:
  - Data kept in a *triple store*, i.e., a database of triples.
  - RDF is served from endpoint as results of *SPARQL queries*.
  - Exposes data (in different ways)
    - with endpoint frontends, e.g., <http://dbpedia.org/resource/Norway>, or
    - by direct SPARQL query: <http://dbpedia.org/sparql>.

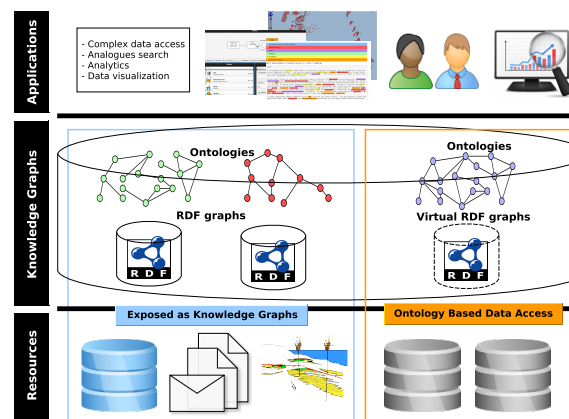
## RDF on the web: Where is it?

- “Behind” *OBDA repositories*:
  - OBDA: Ontology-based Data Access
  - Data kept in a traditional *relational database*
  - Access is transparent via SPARQL queries
  - SPARQL queries are “internally” transformed to SQL queries
  - An RDF representation of the relational database is “virtualized” to answer the SPARQL query.

## Exposing data as RDF

- **Virtual exposure of data (OBDA)**
  - ✓ End-users' friendly access to “unfriendly” relational data
  - ✓ Pay as you go data integration
  - ✗ Requires an ontology in OWL 2 QL
  - ✗ Data remains in old-fashioned databases
- **Data Export**
  - ✓ Easy to exchange data (over the Web)
  - ✓ Ontology not limited to OWL 2 QL
  - ✗ Data replication
  - ✗ Due to size or privacy it may not be possible to export the data

## Exposing data as RDF



## Outline

- 1 Exposing data as RDF
- 2 Data Access in Statoil: limitations and solutions
- 3 OBDA Ingredients
  - Overview
  - Ontology
  - Mappings
  - Query rewriting
  - Bootstrapping
  - Visual Query Formulation
  - Optique

## EU project Optique

- EU Project from 2012-2016
- Aimed at facilitating **scalable end-user access to big data** in the oil and gas industry.
- Advocated for an **OBDA approach**
  - ontology provides a virtual access to the data
  - mappings connect the ontology with the data source.
- Focused around two demanding use cases provided by the industry partners **Siemens** and **Statoil**
- Currently takes 30-70% of engineers' time (e.g., more than 250 MNOK annually).

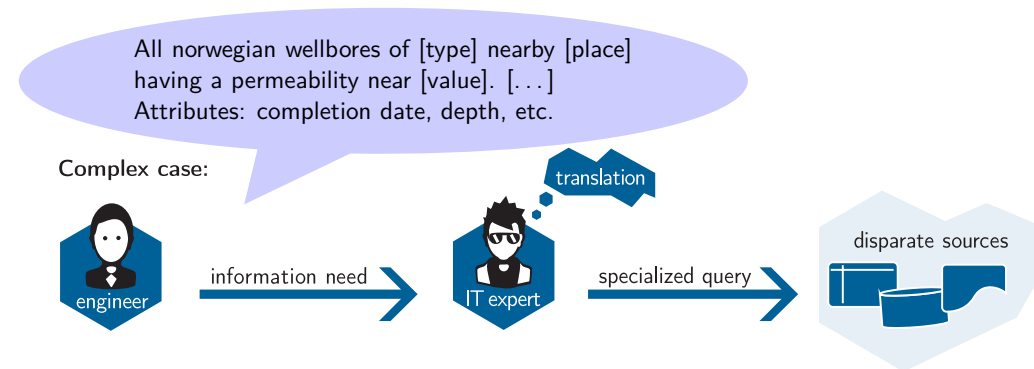
## Limitations

Simple case:



## Limitations

Complex case:



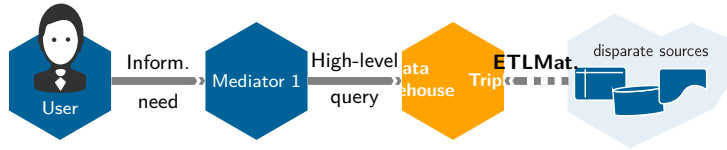
Takes 4 days in average



## Generating a new representation of the data

### 1. Extract Transform Load (ETL) / Materialization

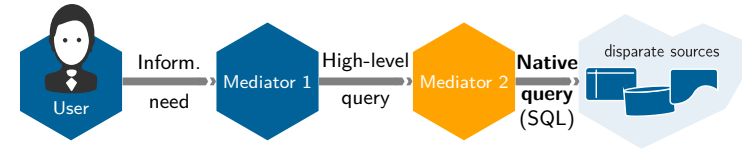
E.g., relational data warehouse, triplestore (RDF)



## Generating a new representation of the data

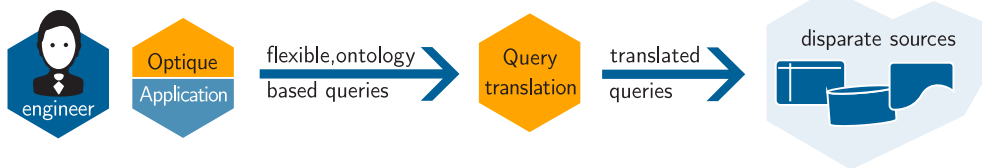
### 2. Virtual views

E.g., virtual databases (Teiid, Apache Drill, Exareme), **OBDA with Optique**



## Optique solution: Ontology-Based Data Access (OBDA)

### Optique solution



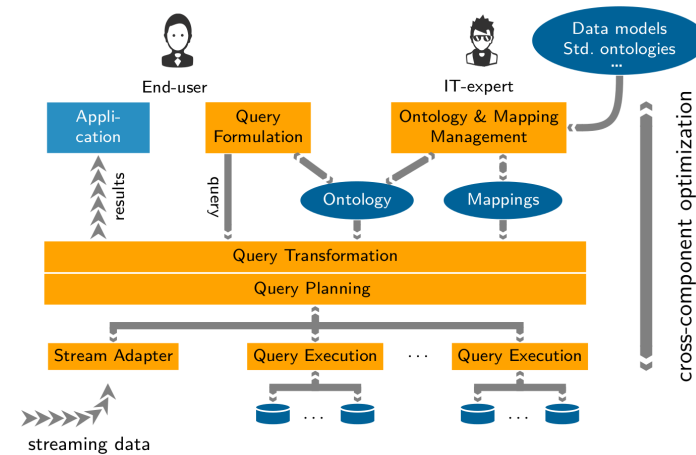
### Choice 1: Generating a new representation of the data

- 1 Virtual views

### Choice 2: Which data format for the virtual view

- 1 Resource Description Framework (RDF)

## Optique architecture



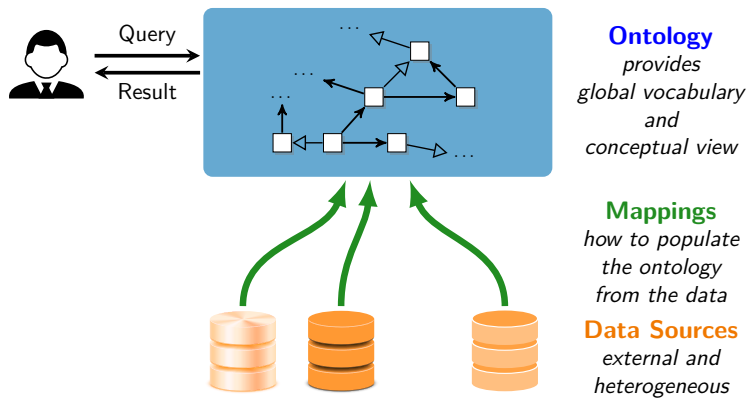
## Outline

- 1 Exposing data as RDF
- 2 Data Access in Statioil: limitations and solutions
- 3 **OBDA Ingredients**
  - Overview
  - Ontology
  - Mappings
  - Query rewriting
  - Bootstrapping
  - Visual Query Formulation
  - Optique

## Outline

- 1 Exposing data as RDF
- 2 Data Access in Statioil: limitations and solutions
- 3 **OBDA Ingredients**
  - Overview
  - Ontology
  - Mappings
  - Query rewriting
  - Bootstrapping
  - Visual Query Formulation
  - Optique

## OBDA framework



## OBDA framework

## Logical transparency in accessing data:

- does not know where and how the data is stored.
- can only see a **conceptual view** of the data.

## OBDA Ingredients

- Relies on . . .
  - **ontology** to provide a virtual access to the data
  - **mappings** to connect the ontology with the data
- Required infrastructure
  - **Query formulation system** to express the information needs in SPARQL (Mediator 1)
  - **Query transformation/rewriting system** to covert from SPARQL to (native) SQL (Mediator 2)
  - **Ontology and mapping bootstrapper**

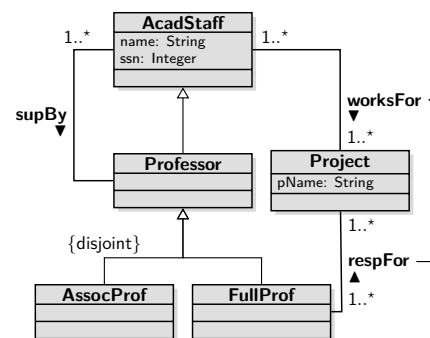
## Outline

- 1 Exposing data as RDF
- 2 Data Access in Statoil: limitations and solutions
- 3 OBDA Ingredients
  - Overview
  - **Ontology**
  - Mappings
  - Query rewriting
  - Bootstrapping
  - Visual Query Formulation
  - Optique

## The Ontology: OWL 2 QL profile

- OWL 2 QL is one of the three standard profiles of OWL 2.
- Is considered a lightweight ontology language:
  - controlled expressive power
  - efficient inference
- Optimized for accessing large amounts of data (i.e., for data complexity):
  - *First-order rewritability* of query answering: queries over the ontology can be rewritten into SQL queries over the underlying relational database.
  - Consistency checking is also first-order rewritable.
- The ontology data (ABox) in an OBDA setting is (usually) implicitly defined through the database and mappings.

## Capturing UML class diagrams/ER schemas in OWL 2 QL



Professor	⊆	AcadStaff
AssocProf	⊆	Professor
FullProf	⊆	Professor
AssocProf	⊆	¬FullProf
AcadStaff	⊆	∃ssn
∃ssn	⊆	AcadStaff
∃ssn <sup>-</sup>	⊆	Integer
∃worksFor	⊆	AcadStaff
∃worksFor <sup>-</sup>	⊆	Project
AcadStaff	⊆	∃worksFor <sup>-</sup>
Project	⊆	∃worksFor
respFor	⊆	worksFor

## Outline

- 1 Exposing data as RDF
- 2 Data Access in Statioil: limitations and solutions
- 3 **OBDA Ingredients**
  - Overview
  - Ontology
  - **Mappings**
  - Query rewriting
  - Bootstrapping
  - Visual Query Formulation
  - Optique

## OBDA Mappings

Global-As-View (GAV) mapping assertion  $\varphi \rightsquigarrow \psi$ 

- $\varphi$ : FO query (over DB predicates)
- $\psi$ : atom (over an RDF predicate)
- Open-World Assumption (by default)

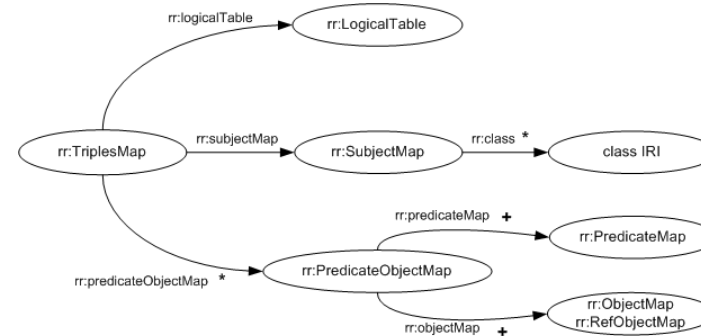
## Class instance (:Student)

Source	$q(s) \leftarrow \text{uni1-student}(s, f, l)$ <pre>SELECT s_id FROM uni1.student</pre>
Target	<pre>Student(URI<sub>1</sub>(s)) ex:uni1/student/{s_id} a :Student .</pre>

## OBDA Mappings: R2RML

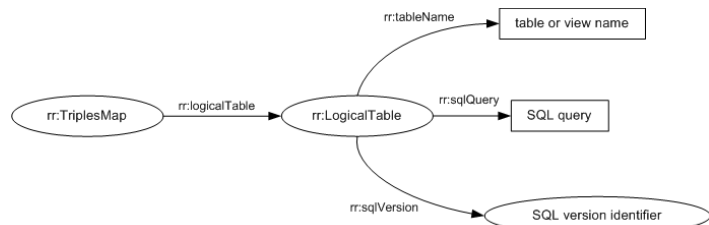
- R2RML is a W3C recommended RDB-to-RDF mapping language
  - <https://www.w3.org/TR/r2rml/>
- Generates RDF triples from a relational database based on specific mappings
- The mappings are specified in Turtle syntax
- The R2RML mapping is an RDF graph consisting of several `rr:TriplesMaps`
  - how to map a logical table in the input relational database into RDF

## OBDA Mappings: R2RML





## OBDA Mappings: R2RML



## OBDA Mappings: R2RML example

## Triples map to populate Student class"

```

<TriplesMap1>      a rr:TriplesMapClass;
  rr:logicalTable [rr:SQLQuery "Select s_id, name, from STUDENT"];
  rr:subjectMap
  [
    rr:template "http://example.com/uni1/student/{s_id}";
    rr:class ex:Student;
  ] ;
  rr:predicateObjectMap
  [
    rr:predicate foaf:name;
    rr:objectMap [ rr:column "name" ]
  ] .
  
```

## OBDA Mappings: R2RML example

## Table STUDENT:

s_id	name
1	Ernesto
2	Martin
3	Leif

## Triples:

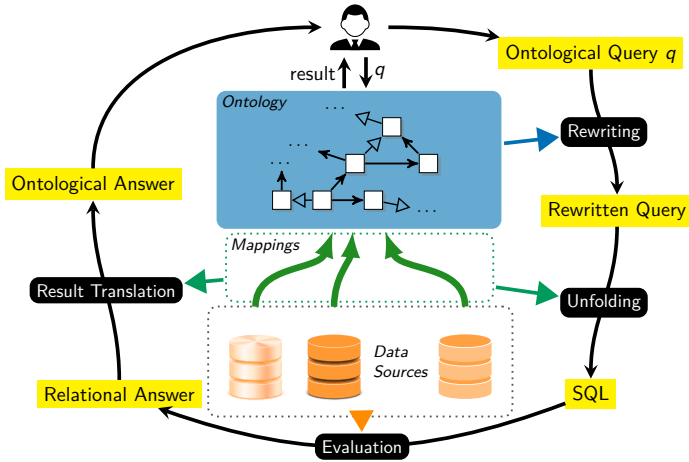
```

http://example.com/uni1/student/1  foaf:name  Ernesto
http://example.com/uni1/student/2  foaf:name  Martin
http://example.com/uni1/student/3  foaf:name  Leif
  
```

## Outline

- 1 Exposing data as RDF
- 2 Data Access in Statioil: limitations and solutions
- 3 OBDA Ingredients
  - Overview
  - Ontology
  - Mappings
  - **Query rewriting**
  - Bootstrapping
  - Visual Query Formulation
  - Optique

### Query answering by rewriting



### Query answering by rewriting (example)

Database:

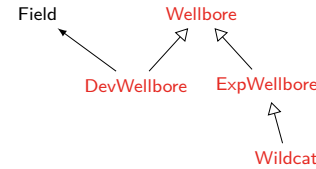
```
wlb_dev(name, ...)
wlb_exp(name, purpose, ...)
```

Query: List all wellbores.

$q$ : Wellbore( $x$ )

$q_0$ : Wellbore( $x$ )  $\cup$  DevWellbore( $x$ )  $\cup$  ExpWellbore( $x$ )  $\cup$  Wildcat( $x$ )

Ontology:



$q_{SQL}$ : SELECT name FROM wlb\_dev UNION  
SELECT name FROM wlb\_exp  
UNION  
SELECT name FROM wlb\_exp WHERE  
purpose = 'WILDCAT'

Mappings:

```
DevWellbore(name)  $\mapsto$  SELECT name FROM wlb_dev
ExpWellbore(name)  $\mapsto$  SELECT name FROM wlb_exp
Wildcat(name)  $\mapsto$  SELECT name FROM wlb_exp
WHERE purpose = 'WILDCAT'
```

$q'_{SQL}$ : SELECT name FROM wlb\_dev UNION  
SELECT name FROM wlb\_exp

### Query answering by rewriting (example)

Database:

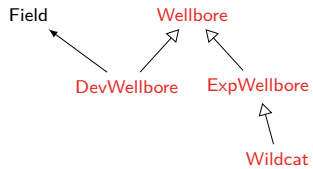
```
wlb_dev(name, ...)
wlb_exp(name, purpose, ...)
```

Query: List all wellbores.

$q$ : Wellbore( $x$ )

$q_0$ : Wellbore( $x$ )  $\cup$  DevWellbore( $x$ )  $\cup$  ExpWellbore( $x$ )  $\cup$  Wildcat( $x$ )

Ontology:



$q_{SQL}$ : SELECT name FROM wlb\_dev UNION  
SELECT name FROM wlb\_exp  
UNION  
SELECT name FROM wlb\_exp WHERE  
purpose = 'WILDCAT'

Mappings:

```
DevWellbore(name)  $\mapsto$  SELECT name FROM wlb_dev
ExpWellbore(name)  $\mapsto$  SELECT name FROM wlb_exp
Wildcat(name)  $\mapsto$  SELECT name FROM wlb_exp
WHERE purpose = 'WILDCAT'
```

$q'_{SQL}$ : SELECT name FROM wlb\_dev UNION  
SELECT name FROM wlb\_exp

### Query answering by rewriting (example)

Database:

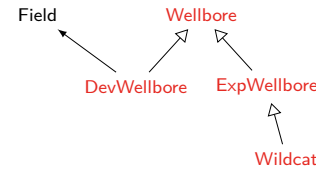
```
wlb_dev(name, ...)
wlb_exp(name, purpose, ...)
```

Query: List all wellbores.

$q$ : Wellbore( $x$ )

$q_0$ : Wellbore( $x$ )  $\cup$  DevWellbore( $x$ )  $\cup$  ExpWellbore( $x$ )  $\cup$  Wildcat( $x$ )

Ontology:



$q_{SQL}$ : SELECT name FROM wlb\_dev UNION  
SELECT name FROM wlb\_exp  
UNION  
SELECT name FROM wlb\_exp WHERE  
purpose = 'WILDCAT'

Mappings:

```
DevWellbore(name)  $\mapsto$  SELECT name FROM wlb_dev
ExpWellbore(name)  $\mapsto$  SELECT name FROM wlb_exp
Wildcat(name)  $\mapsto$  SELECT name FROM wlb_exp
WHERE purpose = 'WILDCAT'
```

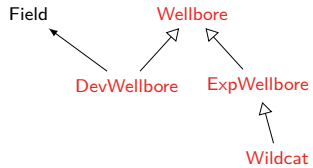
$q'_{SQL}$ : SELECT name FROM wlb\_dev UNION  
SELECT name FROM wlb\_exp

## Query answering by rewriting (example)

## Database:

```
wlb_dev(name, ...)
wlb_exp(name, purpose, ...)
```

## Ontology:



## Mappings:

```
DevWellbore(name) → SELECT name FROM wlb_dev
ExpWellbore(name) → SELECT name FROM wlb_exp
Wildcat(name) → SELECT name FROM wlb_exp
WHERE purpose = 'WILDCAT'
```

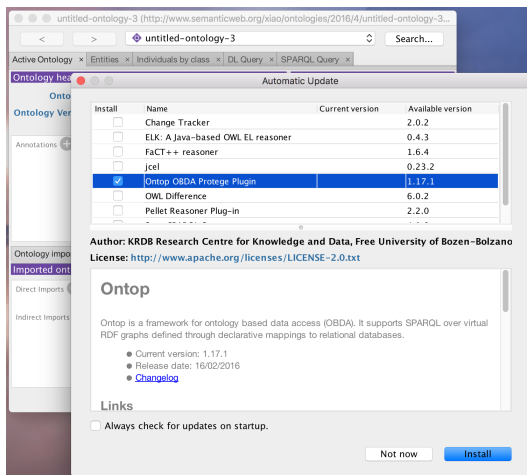
## Query: List all wellbores.

```
q: Wellbore(x)
qO: Wellbore(x) ∪ DevWellbore(x) ∪
ExpWellbore(x) ∪ Wildcat(x)
qSQL: SELECT name FROM wlb_dev
UNION
SELECT name FROM wlb_exp
UNION
SELECT name FROM wlb_exp WHERE
purpose = 'WILDCAT'
q'SQL: SELECT name FROM wlb_dev UNION
SELECT name FROM wlb_exp
```

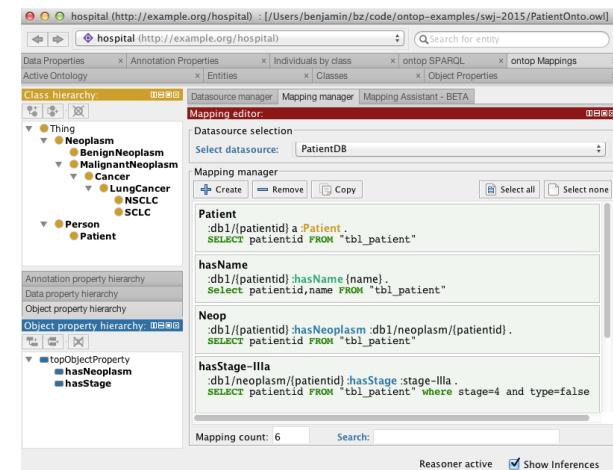
## Query answering by rewriting (tool support)

- **Ontop**: state-of-the-art OBDA system. <http://ontop.inf.unibz.it/>
- Compliant with the RDFS, OWL 2 QL, R2RML, and SPARQL standards.
- Supports all major relational DBs
  - Oracle, DB2, MS SQL Server, Postgres, MySQL, Teiid, Exareme, etc.
- **Open-source** and released under Apache 2 license
- Development of -ontop-:
  - Development started in 2009
  - -ontop- supports (essentially) all features of SPARQL 1.0 and the OWL 2 QL entailment regime of SPARQL 1.1.
  - Other features of SPARQL 1.1 (e.g., aggregates, property path queries, negation) are work in progress.

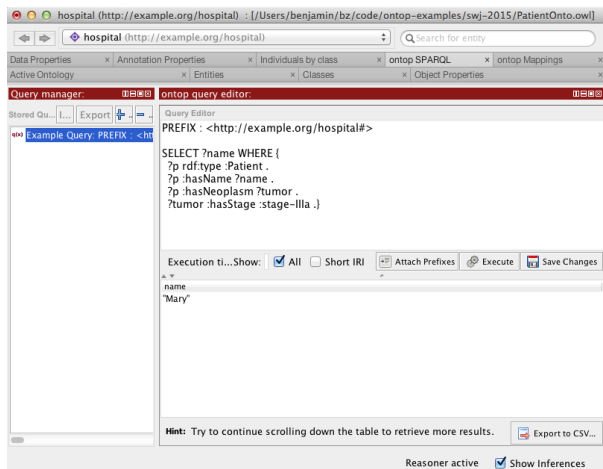
## Ontop plugin available from Protégé



## Ontop: Mapping editor in Protégé



# Ontop: SPARQL query answering in Protégé



# Outline

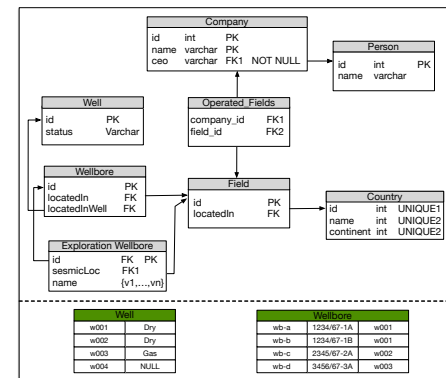
- 1 Exposing data as RDF
- 2 Data Access in Statioil: limitations and solutions
- 3 **OBDA Ingredients**
  - Overview
  - Ontology
  - Mappings
  - Query rewriting
  - **Bootstrapping**
  - Visual Query Formulation
  - Optique

# Bootstrapping overview

- Given a relational database (semi)automatically extracts ontological vocabulary, ontology and R2RML mappings
- Bootstrappers may also accept as input an ontology
  - R2RML mappings will link the given ontoloy to the database, or
  - The given ontology will be *aligned* with the bootstrapped ontology
- Type of mappings:
  - W3C direct mapping specification (*schema driven*)
    - To a given or bootstrapped ontology vocabulary
  - Mappings beyond direct ones (*data driven*)
    - Clusters of tuples
    - Joinable tuples

# Bootstrapping: Vocabulary Generation

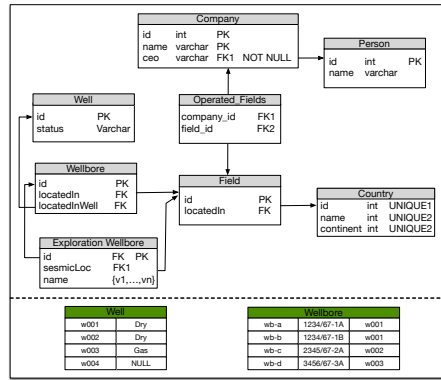
- **W3C directives**
  - Tables → classes
  - Foreign Keys → object properties
  - Data columns → data properties
  - Binary tables → fresh object properties
- **Attribute naming schema:**
  - Unique names (e.g. Person.name)
  - Reusable names (e.g. name)



## Bootstrapping: Axiom Generation

- OWL 2 expressiveness

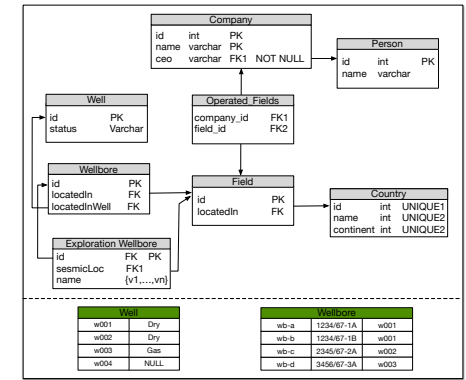
- OWL 2 QL (e.g. OBDA/Optique)
- OWL 2 EL (e.g. EOLO)
- OWL 2 RL (e.g. RDFox)
- OWL 2 (e.g. PAGOdA, HermiT)



## Bootstrapping: Axiom Generation

- Unique constraints

- Person **HasKey**: id (OWL 2 RL/EL)
- Global onto. constraints
  - **Functional**: Person.name (OWL 2 RL)
  - name **Domain**: Person (all profiles)
  - Person.name **Range**: xsd:string (all profiles)
- Local onto. constraints
  - Person **subclassOf**: name **some** xsd:string (OWL 2 QL/EL)
  - Person **subclassOf**: name **only** xsd:string (OWL 2 RL)
  - Person **subclassOf**: name **exactly 1** xsd:string (OWL 2)



## Bootstrapping: Datatypes

- Clear mapping between SQL and XML schema datatypes
- Not all XML Schema datatypes are included in OWL 2
  - xsd:date not in OWL 2
  - xsd:boolean and xsd:double not in OWL 2 QL/EL
- Value spaces of primitive datatypes are disjoint (e.g. xsd:double and xsd:decimal)
- ✗ Problems when materializing the data or using the ontology for virtual access to data
- ✓ Solution: rdfs:Literal

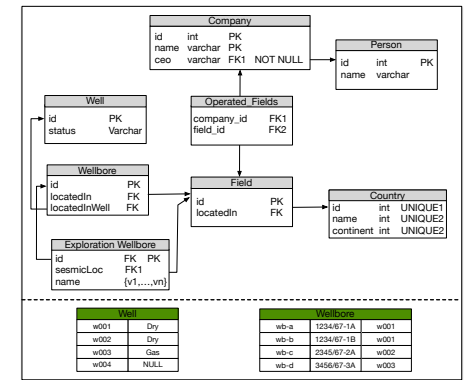
## Bootstrapping: Taxonomy Generation

- Data driven

- Clusters of tuples
- Joinable tuples
- e.g. Well\_Dry **SubclassOf**: Well (w001, w002)
- e.g. Well\_with\_Wellbore **SubclassOf**: Well (w001, w002, w003)

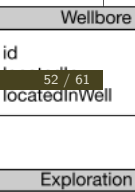
- Schema driven

- A single-column Foreign Key and



- e.g. Exploration\_Wellbore **SubclassOf**: Wellbore

Wellbore		
wb-a	1234/67-1A	w001



## Tool support and Lessons learnt

- **BootOX**: <http://www.cs.ox.ac.uk/isg/tools/BootOX/>
- **Feedback from use cases and evaluation...**
  - ✓ Good as a first approximation of the ontology and mappings
  - ✓ Competitive results in (academic) benchmarks
  - ✗ For the largest Statioil datasources, the solution is far from perfect
  - ✗ Ontology close to the original database
  - ✗ Large amount of ontology entities and R2RML mappings
  - ✓ Implementation of an incremental/interactive bootstrapping
- **Future work...**
  - New ways to exploit the data and schema (e.g. views)

## Outline

- 1 Exposing data as RDF
- 2 Data Access in Statioil: limitations and solutions
- 3 **OBDA Ingredients**
  - Overview
  - Ontology
  - Mappings
  - Query rewriting
  - Bootstrapping
  - **Visual Query Formulation**
  - Optique

## Visual query formulation (OptiqueVQS)

The screenshot shows the OptiqueVQS interface. At the top, a visual query graph is displayed with nodes for 'Wellbore name(o)', 'Company name(o)', and 'drillingOper...'. The graph includes labels for 'kernel', 'inverted', and 'pivot'. Below the graph, there are controls for 'Delete Node', 'Undo', 'Redo', 'New Query', 'Save Query', 'Stored Queries', 'Q-Config', 'SPARQL Query', and 'Run Query'. The bottom panel shows a list of ontology classes: 'ProductionLicenceAreaPerBlock (inv) country', 'Company wellOperatedBy', 'License wellLicense', and 'WellboreOperatedByPerBlock'. A search bar is visible above the list.

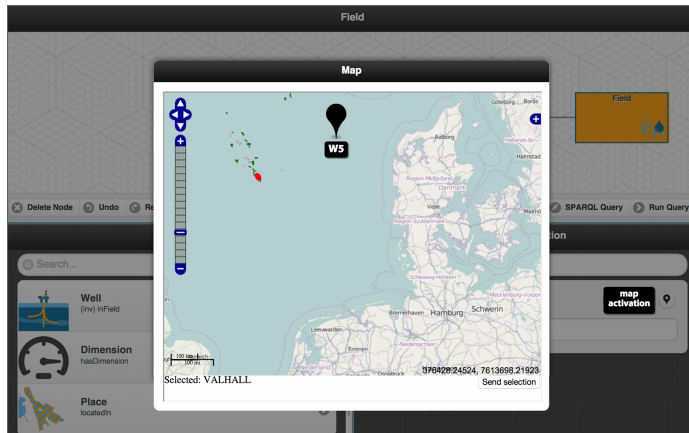
## Visual query formulation (OptiqueVQS)

The screenshot shows the OptiqueVQS interface with a SPARQL query entered in the top panel. The query is:
 

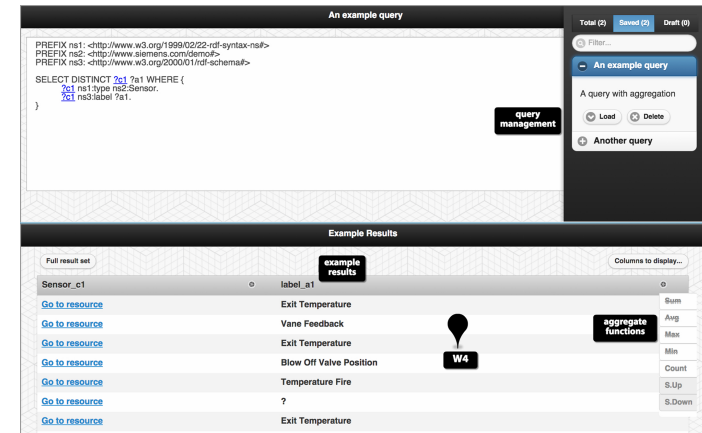
```
PREFIX ns3: <http://www.optique-project.eu/well-ontology/>
SELECT ?c1 ?a3 ?a2 ?a1 ?c3 ?a4 WHERE {
?c1 ns1:type ns2:Wellbore.
?a2 ns1:type ns3:Well.
?c3 ns1:type ns2:Company.
?c1 ns3:hasWellbore ?a2.
?c1 ns2:drillingOperatorCompany ?c3.
?c1 ns2:name ?a3.
?a2 ns2:dateSyncNPD ?a1.
?a2 ns2:wellType ?a2.
?a3 ns2:name ?a4.
FILTER(regex(?a2, "DEVELOPMENT", "i"))}
```

 The bottom panel shows the same list of ontology classes as in the previous screenshot, with a search bar above it.

## Visual query formulation (OptiqueVQS)



## Visual query formulation (OptiqueVQS)



## Outline

- 1 Exposing data as RDF
- 2 Data Access in Statoil: limitations and solutions
- 3 OBDA Ingredients
  - Overview
  - Ontology
  - Mappings
  - Query rewriting
  - Bootstrapping
  - Visual Query Formulation
  - Optique

## Optique infrastructure

- Training material:
  - <http://optique-northwind.fluidops.net> (demo/demo)
- OptiqueVQS can be tested online
- Local installation possible (academic license):
  - <https://appcenter.fluidops.com/resource/Search?search=optique>
- What is next?
  - SIRIUS: <http://sirius-labs.no/>

## Questions?

Ernesto Jiménez-Ruiz (ernestoj@ifi.uio.no)  
Office hours: from 9:00 to 16:00 at OJD 8165

### Acknowledgements

- -ontop- team (Diego Calvanese, Benjamin Cogrel, Guohui Xiao, etc.) and Dag Hovland for sharing their wonderful slides
  - <http://ontop.inf.unibz.it/ekaw-2016-tutorial/>
- Ahmet Soylu: main person behind OptiqueVQS
- Optique partners