

# Islands of tractability in ontology-based data access

Michael Zakharyashev

*Department of Computer Science and Information Systems,  
Birkbeck, University of London*

<http://www.dcs.bbk.ac.uk/~michael>

**EPSRC**

Engineering and Physical Sciences  
Research Council

supported by EPSRC grants ExODA EP/H05099X and iTract EP/M012670



## Data access in industry

(from Norwegian Petroleum Directorate's FactPages)

show me the wellbores completed before 2008 where Statoil as a drilling operator sampled less than 10 meters of cores



5 days later:

```
SELECT DISTINCT cores.wlbName, cores.lenghtM, wellbore.wlbDrillingOperator, wellbore.wlbCompletionYear
FROM
```

```
( (SELECT wlbName, wlbNpdidWellbore, (wlbTotalCoreLength * 0.3048) AS lenghtM
  FROM wellbore_core
  WHERE wlbCoreIntervalUom = '(ft)' )
```

```
UNION
```

```
(SELECT wlbName, wlbNpdidWellbore, wlbTotalCoreLength AS lenghtM
  FROM wellbore_core
  WHERE wlbCoreIntervalUom = '(m)' )
```

```
) as cores,
```

```
( (SELECT wlbNpdidWellbore, wlbDrillingOperator, wlbCompletionYear
  FROM wellbore_development_all
```

```
UNION
```

```
(SELECT wlbNpdidWellbore, wlbDrillingOperator, wlbCompletionYear
  FROM wellbore_exploration_all )
```

```
UNION
```

```
(SELECT wlbNpdidWellbore, wlbDrillingOperator, wlbCompletionYear
  FROM wellbore_shallow_all )
```

```
) as wellbore
```

```
WHERE wellbore.wlbNpdidWellbore = cores.wlbNpdidWellbore
```

```
...
```

**In STATOIL:**

**1,000 TB of relational data**

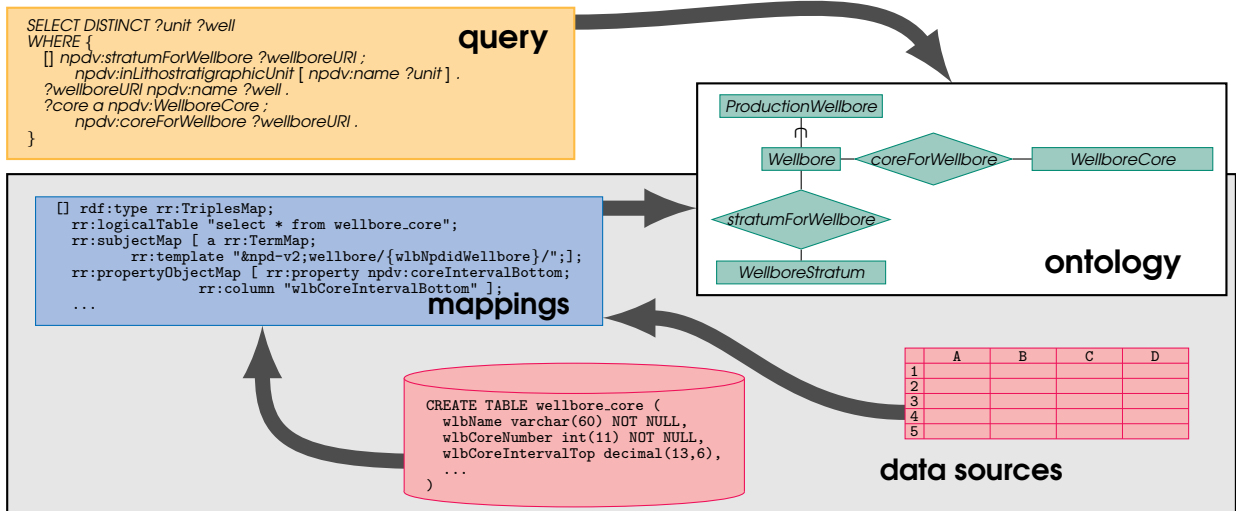
**2,000 tables**

**different schemas**

**30–70% of time on data gathering**

# Ontology-based data access (OBDA)

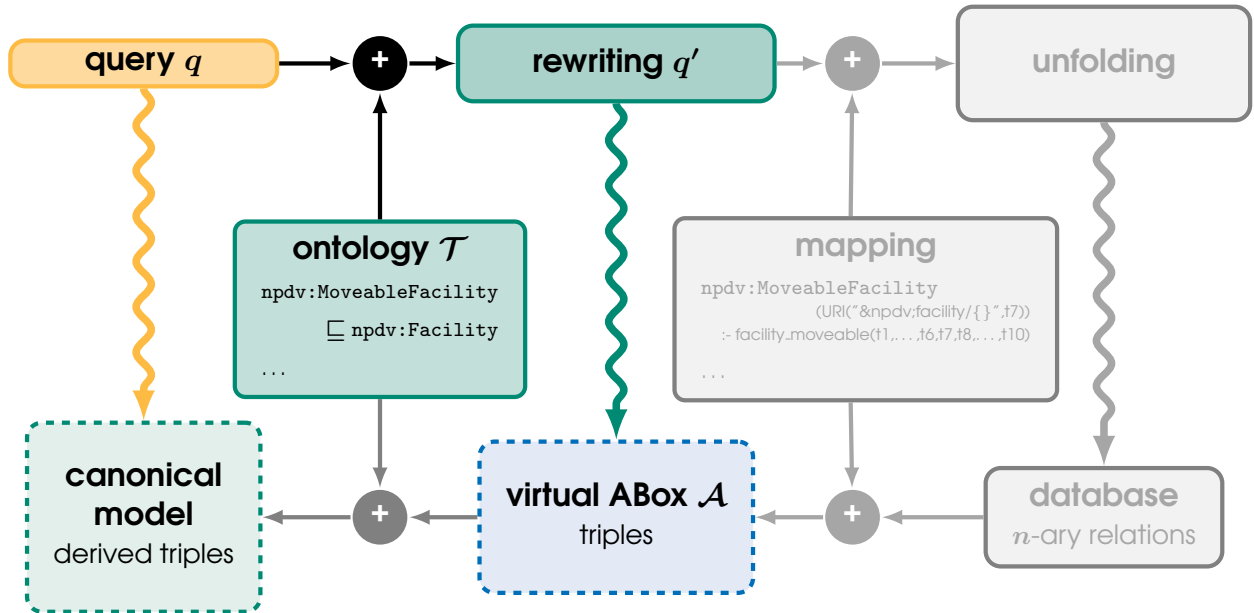
(the Romans  $\approx$  2007)



## Ontology

- gives a high-level conceptual view of the data
- provides a convenient & natural vocabulary for user queries
- enriches incomplete data with background knowledge

## OBDA via FO-rewriting



for all  $\mathcal{A}$  and  $\vec{a}$ ,  $\mathcal{T}, \mathcal{A} \models q(\vec{a}) \iff \mathcal{I}_{\mathcal{A}} \models q'(\vec{a})$

reduction to DB query evaluation

## OWL 2 QL profile of OWL 2 (W3C 2012)

### Roles

$$\varrho(x, y) ::= \top \mid P(x, y) \mid P(y, x)$$

$$R ::= \top \mid P \mid P^{-}$$

### Basic concepts

$$\tau(x) ::= \top \mid A(x) \mid \exists y \varrho(x, y)$$

$$B ::= \top \mid A \mid \exists R$$

### TBoxes

$$\forall x (\tau(x) \rightarrow \tau'(x))$$

$$B \sqsubseteq B'$$

$$\forall x, y (\varrho(x, y) \rightarrow \varrho'(x, y))$$

$$R \sqsubseteq R'$$

$$\forall x \varrho(x, x)$$

$R$  is reflexive

$$\forall x (\tau(x) \wedge \tau'(x) \rightarrow \perp)$$

$$B \sqcap B' \sqsubseteq \perp$$

$$\forall x, y (\varrho(x, y) \wedge \varrho'(x, y) \rightarrow \perp)$$

$$R \sqcap R' \sqsubseteq \perp$$

$$\forall x (\varrho(x, x) \rightarrow \perp)$$

$R$  is irreflexive

### Sugar

$$\forall x (\tau(x) \rightarrow \exists y (\varrho_1(x, y) \wedge \dots \wedge \varrho_k(x, y) \wedge \tau'(y)))$$

$$B \sqsubseteq \exists R.B'$$

(expressible via additional role inclusions)

### ABoxes

$$\{A(a), P(a, b), \dots\}$$

based on the 'DL-Lite family' designed by the Romans ( $\approx$  2005) and extended by Artale, Calvanese, Kontchakov & Z (2007–9)

## Example

### Staff ontology $\mathcal{T}$

$$\begin{aligned}\forall x \ (ProjectManager(x) \rightarrow \exists y \ (isAssistedBy(x, y) \wedge PA(y))) \\ \forall x \ (\exists y \ managesProject(x, y) \rightarrow ProjectManager(x)) \\ \forall x \ (ProjectManager(x) \rightarrow Staff(x)) \\ \forall x \ (PA(x) \rightarrow Secretary(x))\end{aligned}$$

**User query  $q$ :** find the staff assisted by secretaries

$$q(x) = \exists y \ (Staff(x) \wedge isAssistedBy(x, y) \wedge Secretary(y))$$

### PE-rewriting of ontology-mediated query $(\mathcal{T}, q)$

$$q'(x) = \exists y \ [Staff(x) \wedge isAssistedBy(x, y) \wedge (Secretary(y) \vee PA(y))] \vee \\ ProjectManager(x) \vee \exists z \ managesProject(x, z)$$

## Why are OWL 2 QL OMQs FO-rewritable?

✓ **Canonical model (chase)**  $\mathcal{C}_{\mathcal{T}, \mathcal{A}}$  of a given consistent  $(\mathcal{T}, \mathcal{A})$

homomorphically embeddable into every model of  $(\mathcal{T}, \mathcal{A})$

$$\longrightarrow \mathcal{T}, \mathcal{A} \models q \iff \mathcal{C}_{\mathcal{T}, \mathcal{A}} \models q \text{ for any CQ } q$$

**Example:**  $\mathcal{T} = \{A \sqsubseteq \exists R^-. \exists R. B, \quad B \sqsubseteq \exists S. B\} \quad \mathcal{A} = \{A(a)\}$



*all Horn DLs have canonical models* but OMQ  $(\{\exists R. A \sqsubseteq A\}, A(x))$  is not FO-rewritable  
(recursive datalog needed)

✓ **Bounded depth derivation property:** there is a function  $f$  such that  
 $\mathcal{T}, \mathcal{A} \models q \iff \mathcal{C}_{\mathcal{T}, \mathcal{A}}^N \models q$  with  $\mathcal{C}_{\mathcal{T}, \mathcal{A}}^N$  constructed in  $N = f(|\mathcal{T}|, |q|)$  steps

$\Leftrightarrow$  **FO-rewritability**

$f$  is **polynomial** for OWL 2 QL

## What is the price of OBDA?

- reduction to DB query evaluation could be too expensive  
→ OBDA would not be viable

### 1 what is the size of rewritings ?

- depending on the type of OMQs
- depending on the type of rewritings

new research (succinctness) problem

### 2 what is the combined complexity of OMQ answering ?

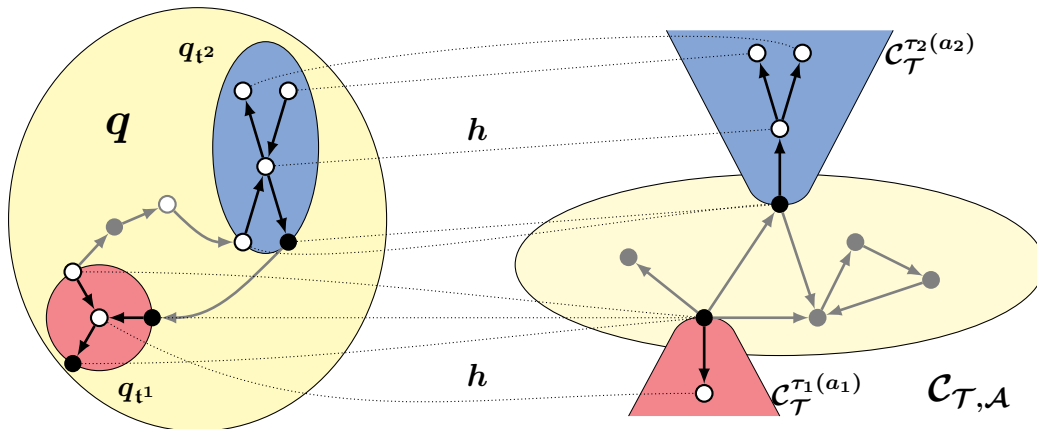
- depending on the type of OMQs

well-known problem in DB theory

it may turn out that reduction to DB query evaluation is not most optimal way of OMQ answering



## Tree-witness rewriting of OMQ $Q = (\mathcal{T}, q)$



$$q_{\text{tw}}(\vec{x}) = \bigvee_{\Theta \text{ independent set of tree witnesses}} \exists \vec{y} \left( \bigwedge_{S(\vec{z}) \in q \setminus q_{\Theta}} S(\vec{z}) \wedge \bigwedge_{t \in \Theta} \text{tw}_t \right)$$

$\Theta$  is **independent** if  $q_t \cap q_{t'} = \emptyset$ , for any distinct  $t, t' \in \Theta$

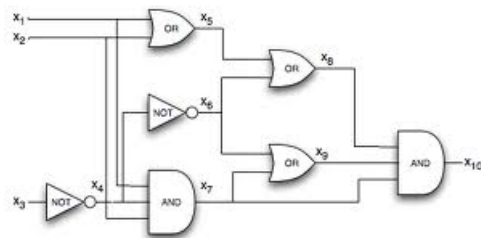


however, it can be simplified to a polynomial-size PE-rewriting:

can we always do this?

## Circuit complexity

**P/poly**: the class of problems decidable by  
**polynomial-size** circuit families



$$P \subseteq P/poly$$

if  $NP \not\subseteq P/poly$  then  $P \neq NP$

– almost all Boolean functions with  $n$  inputs require circuits of size  $\Theta(2^n/n)$

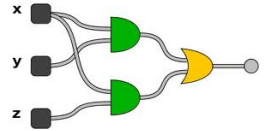
(Shannon 1949)

🤔 are there **complex** Boolean functions  $f_n$  in NP? (known lower bound:  $5n - o(n)$ )

nobody knows, but ...

# Monotone circuit complexity

(Razborov, Raz, *et al.* 1985)



Boolean variables  $e_{ij}$  give graph  $G = (V, E)$ :  $V = \{1, \dots, n\}$ ,  $E = \{\{i, j\} \mid e_{ij} = 1\}$

–  $\text{CLIQUE}_{n,k}(\vec{e}) = 1$  iff  $G$  contains a  **$k$ -clique** (e.g., for  $k \leq n^{1/4}$ )

monotone circuits: **exp** ( $2^{\varepsilon\sqrt{k}}$ )

monotone formulas: **exp**

formulas with  $\neg$ : **superpoly** unless  $\text{NP} \subseteq \text{P/poly}$

–  $\text{MATCHING}_n(\vec{e}) = 1$  iff the bipartite graph  $\vec{e}$  with  $n$  vertices in each part has a **perfect matching** (subset of edges containing every node once)

monotone formulas: **exp**

formulas with  $\neg$ : **poly**

## Tree-witness rewriting as a Boolean function

OMQ  $Q = (\mathcal{T}, q) \longrightarrow$  a **hypergraph**  $H_Q = (V, E)$  where  
vertices  $V$  = atoms of  $q$   
hyperedges  $E$  = tree witnesses  $q_t$

monotone Boolean **hypergraph function** for  $Q$  (or hypergraph  $H_Q$ )

$$f_Q = \bigvee_{E' \subseteq E \text{ independent}} \left( \bigwedge_{v \in V \setminus V_{E'}} p_v \wedge \bigwedge_{e \in E'} p_e \right)$$

(some tweaks required in case of exponentially-many tree witnesses)

- Boolean formula  $\varphi$  for  $f_Q$   $\longrightarrow$  FO-rewriting of size  $O(|\varphi| \cdot |Q|)$
- monotone Boolean formula  $\varphi$  for  $f_Q$   $\longrightarrow$  PE-rewriting \_\_\_\_\_
- monotone Boolean circuit  $\varphi$  for  $f_Q$   $\longrightarrow$  NDL-rewriting \_\_\_\_\_  
(nonrecursive datalog)

tool for obtaining upper succinctness and complexity bounds

using classical circuit complexity






## Tool for lower bounds

For any OMQ  $Q = (\mathcal{T}, q)$  and assignment  $\alpha: \text{predicates}(q) \rightarrow \{0, 1\}$ ,

$$\mathcal{A}_\alpha = \{A(a) \mid \alpha(A) = 1\} \cup \{P(a, a) \mid \alpha(P) = 1\}$$

ABox with **a single individual  $a$**

**Primitive evaluation function:**  $g_Q(\alpha) = 1 \Leftrightarrow \mathcal{T}, \mathcal{A}_\alpha \models q(\vec{a})$

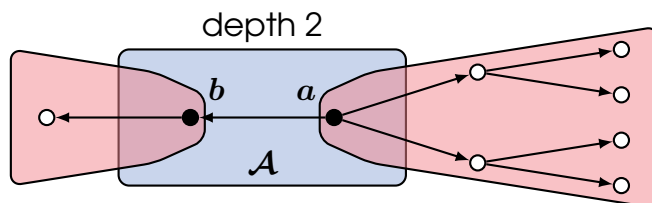
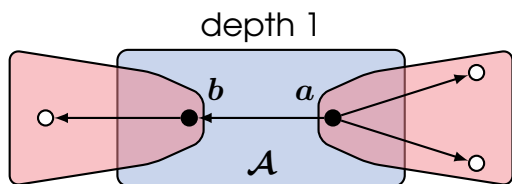
- FO-rewriting  $q'$  of  $Q$   Boolean formula for  $g_Q$  of size  $O(|q'|)$
- PE-rewriting  $q'$  of  $Q$   monotone Boolean formula for  $g_Q$  
- NDL-rewriting  $q'$  of  $Q$   monotone Boolean circuit for  $g_Q$    
(proof by quantifier elimination)

tool for obtaining lower succinctness bounds

using classical circuit complexity

## Case study: OMQs with ontologies of depth 1

no axioms such as  $A \sqsubseteq \exists P, \quad \exists P^- \sqsubseteq \exists R$



$Q = (\mathcal{T}, q)$  with  $\mathcal{T}$  of **depth 1**  $\longrightarrow$  hypergraph  $H_Q$  is of **degree  $\leq 2$**   
 each vertex belongs to  $\leq 2$  hyperedges

hypergraph  $H$  of degree  $\leq 2$   $\longrightarrow \exists$  OMQ  $Q_H$  with  $\mathcal{T}$  of depth 1 and  $H \cong H_{Q_H}$

What can hypergraph functions of degree 2 compute?

## Hypergraph programs (HGPs)

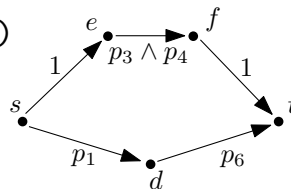
An **HGP** is a hypergraph  $H = (V, E)$  with every vertex labelled by  $0$ ,  $1$ ,  $p_i$  or  $\neg p_i$

**computes**  $f$ :  $f(\vec{\alpha}) = 1 \Leftrightarrow$  there is an independent  $E' \subseteq E$  **covering all zeros**  
(contains all vertices whose labels evaluate to  $0$  under  $\vec{\alpha}$ )

**monotone** if no  $\neg p_i$  among the labels

Any monotone HGP based on  $H$  computes a sub-function of  $f_H$

✓ HGPs based on hypergraphs of **degree  $\leq 2$**  are polynomially equivalent to  
**nondeterministic branching programs (NBPs)**



**HGP<sup>2</sup> = NBP = NL/poly**

functions computable by NLogSpace TMs with polynomial advice functions  
(non-uniform analogue of NLOGSPACE)



## Rewritings for OMQs with ontologies of depth 1



$\mathbf{HGP}^2 = \mathbf{NL}/\mathbf{poly} \subseteq \mathbf{P}/\mathbf{poly}$  (for monotone functions)



polynomial-size **NDL**-rewritings



there is a monotone  $f$  computable by a polynomial-size NBP, but  
any monotone Boolean formula computing  $f$  is of size  $n^{\Omega(\log n)}$



$\exists$  OMQ with superpolynomial **PE**-rewritings only



all OMQs have polynomial **FO**-rewritings



$\mathbf{NC}^1 = \mathbf{NL}/\mathbf{poly}$



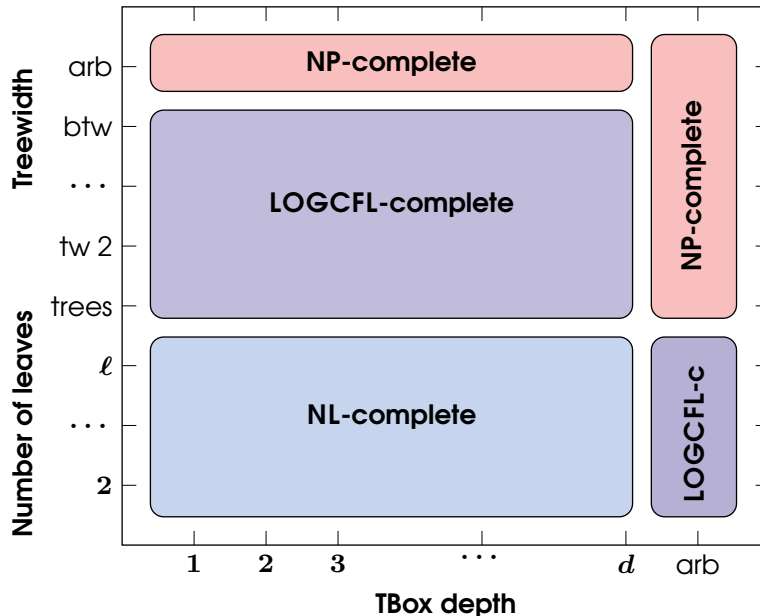
all OMQs with CQs of bounded treewidth have polynomial PE-rewritings  
all tree-shaped OMQs have polynomial-size  $\Pi_4$ -rewritings ( $\wedge \vee \wedge \vee$ )

(SPARQL queries under OWL 2 QL entailment regime)

poly FO iff NL/poly  $\subseteq$  NC<sup>1</sup>



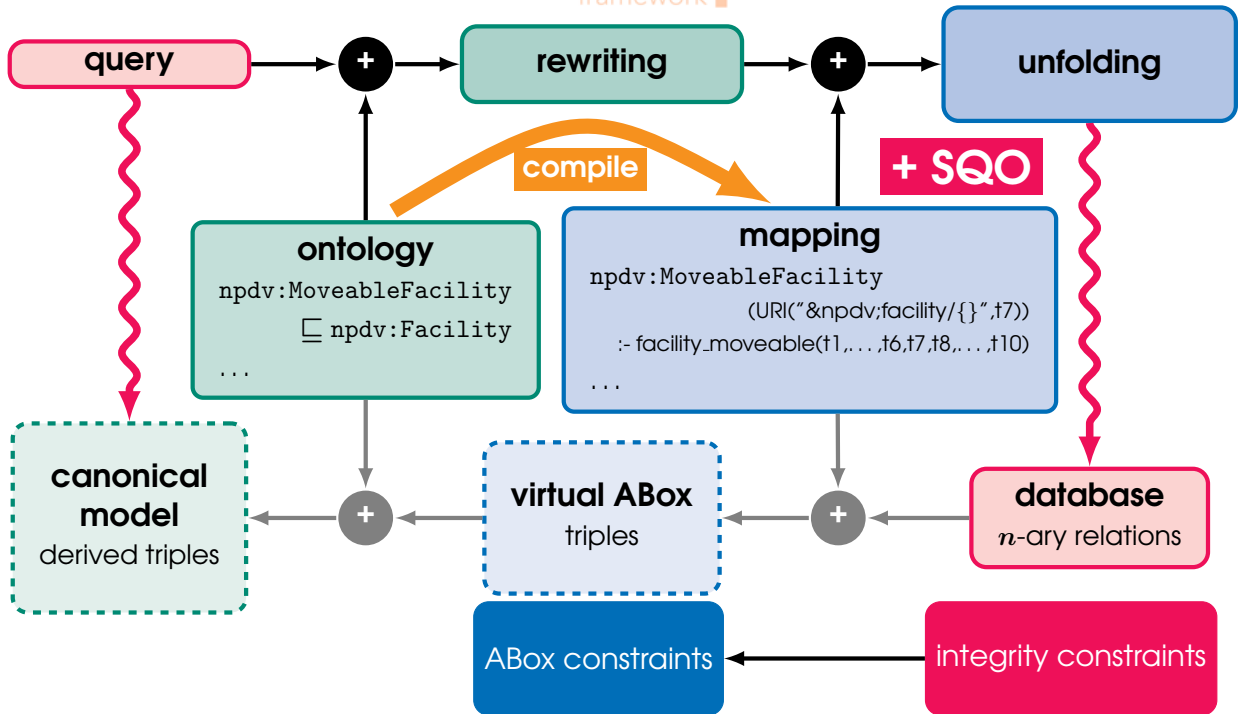
## Combined complexity landscape



- CQ evaluation over databases is **NP-complete**  $L \subseteq NL \subseteq LOGCFL \subseteq NC^2 \subseteq P \subseteq NP$
- bounded treewidth CQ evaluation is **LOGCFL-complete** (logspace reducible to a CFL)

Gottlob et al. 2001

# ontop framework



Rodriguez-Muro, Calvanese, Kontchakov, Rezk, Xiao, Z 2010–15

## Ontop in practice

- ✓ T-mappings compile (big parts of) OWL 2 QL ontologies into mappings (domain and range constraints, concept and role hierarchies)
- ✓ can be optimised **offline**
- ✓ few tree witnesses in real-world OBDA → polynomial-size rewritings
- ✓ database constraints and SQO significantly simplify T-mappings

→ efficient SQL queries over the data

✗ some important conceptual modelling constructs are missing in OWL 2 QL

$A \sqsubseteq B \sqcup C$   $\exists R.A \sqsubseteq B$  owl:sameAs

? **islands** of OMQ rewritability & succinctness for expressive languages

# iTract: Islands of Tractability in Ontology-Based Data Access

## EPSRC UK project:

- (i) establish a novel, OMQ-centric approach to OBDA aiming to identify **islands of tractable** OMQs in rich ontology and query languages
- (ii) develop uniformly efficient OMQ answering techniques for the identified islands
- (iii) implement and test these techniques in practice, using state-of-the-art OBDA systems

## Team:

- London: MZ (PI), S Kikot (RA), R Kontchakov (co-I), I Razgon (co-I)
- Liverpool: F Wolter (PI), F Papacchini (RA), A Hernich, B Konev

## Project partners:

- University of Bozen-Bolzano (Diego Calvanese)
- University of Bremen (Carsten Lutz)
- University of Oslo (Arild Waaler)
- IBM Watson, New York (Mariano Rodriguez-Muro)