

# Cost-driven OBDA

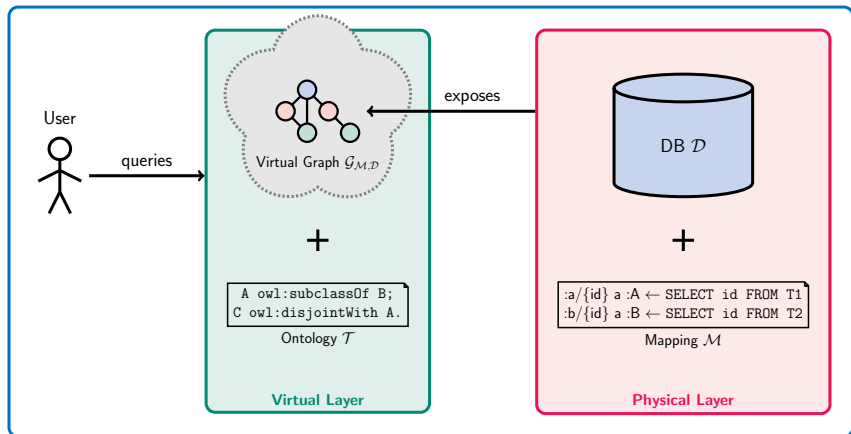
ISWC 2017

Vienna, Austria

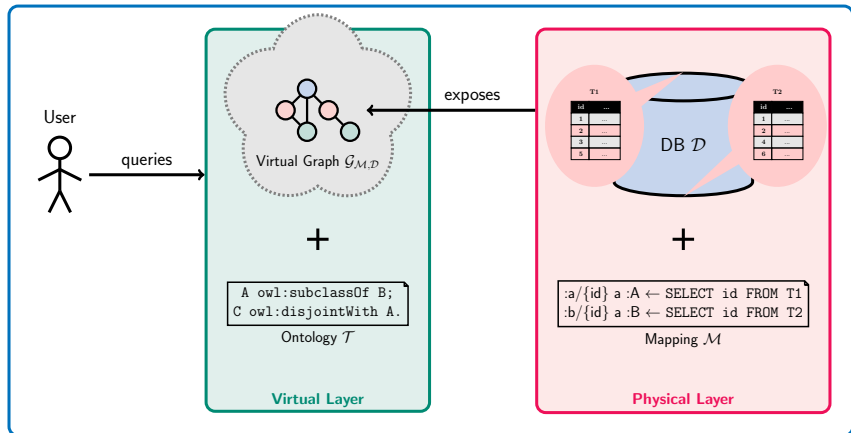
- ▶ **Speaker:** Davide Lanti
- ▶ **Joint Work With:** Guohui Xiao, Diego Calvanese
- ▶ **Institution:** Free University of Bozen-Bolzano

`dlanti@inf.unibz.it`

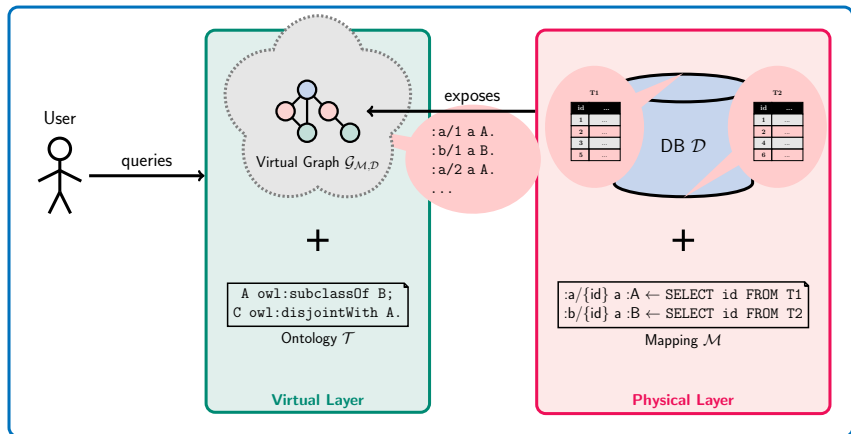
## An OBDA Scenario (Semantic Web View)



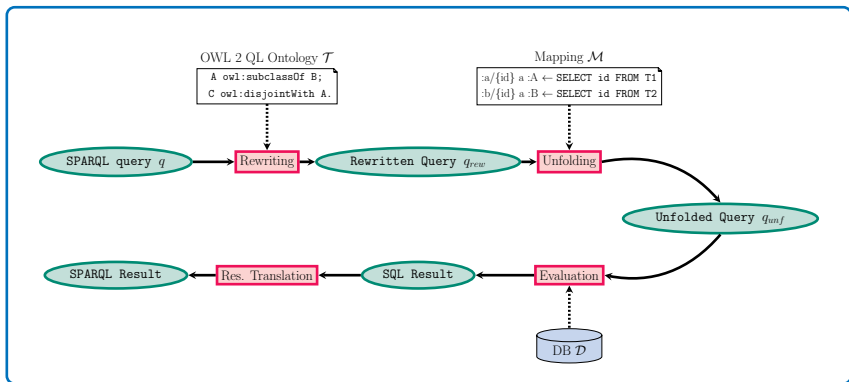
## An OBDA Scenario (Semantic Web View)



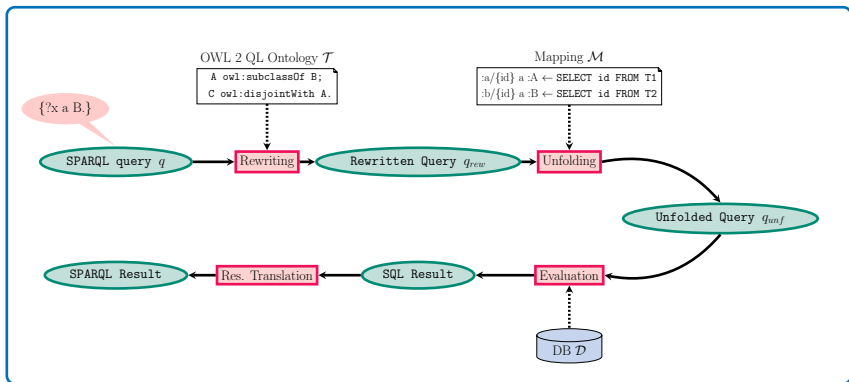
## An OBDA Scenario (Semantic Web View)



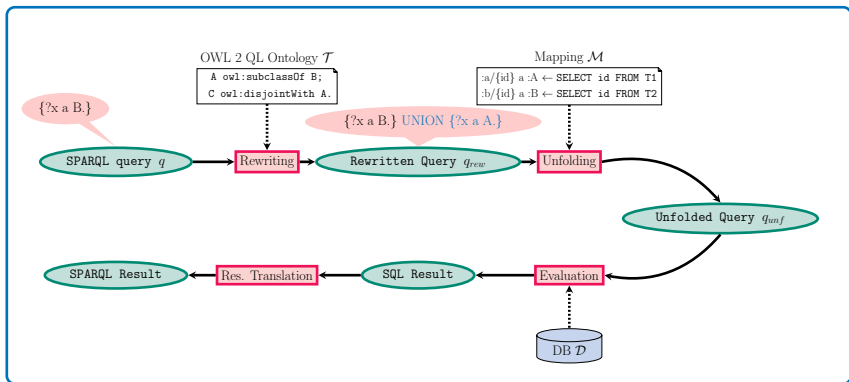
# Query Evaluation Process I



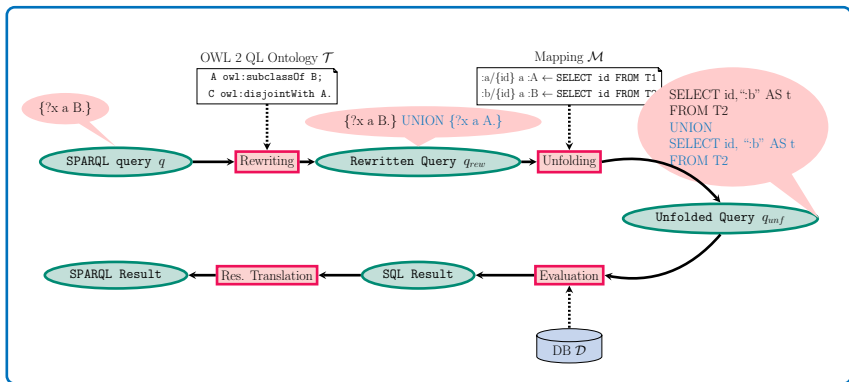
## Query Evaluation Process II



## Query Evaluation Process III

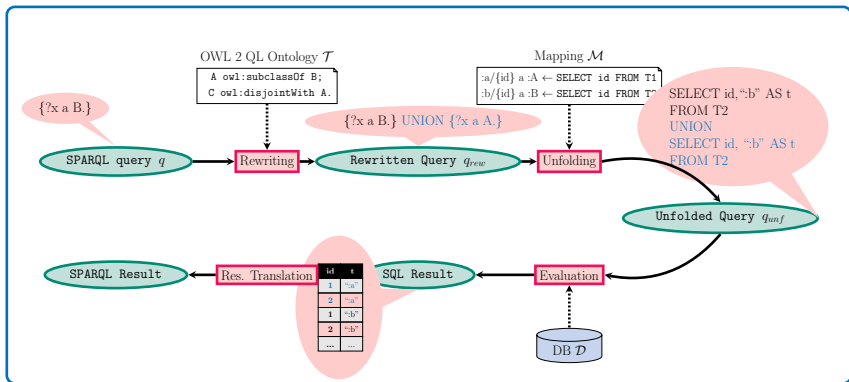


## Query Evaluation Process IV

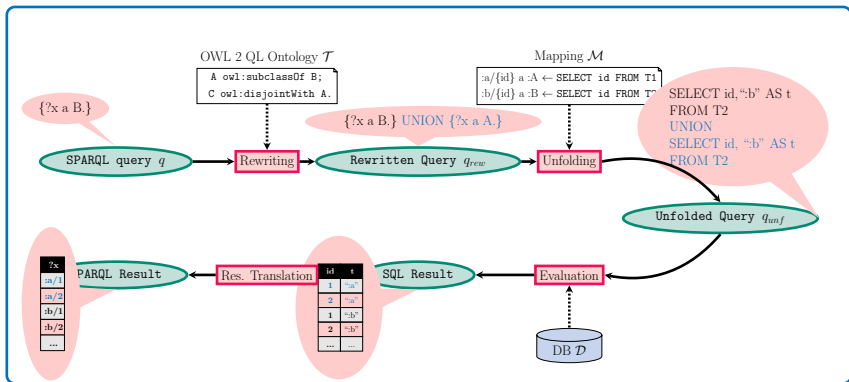




## Query Evaluation Process V



## Query Evaluation Process VI



## The Performance Problem

- ▶ **Problem:**
  - ▷ Performance is an issue in real-world applications of OBDA
- ▶ **Current Solutions:**
  - ▷ Data-independent Optimizations
    - ▶ Structural Optimizations
    - ▶ Semantic Query Optimizations

## The Performance Problem

- ▶ **Problem:**
  - ▷ Performance is an issue in real-world applications of OBDA
- ▶ **Current Solutions:**
  - ▷ Data-independent Optimizations
    - ▶ Structural Optimizations
    - ▶ Semantic Query Optimizations
- ▶ **Contribution of this work:**
  - ▷ Data-dependent Optimizations
    - ▶ Explore a space of alternatives for unfoldings
    - ▶ Choose the best alternative, given the instance

## The Performance Problem

- ▶ **Problem:**
  - ▷ Performance is an issue in real-world applications of OBDA
- ▶ **Current Solutions:**
  - ▷ Data-independent Optimizations
    - ▶ Structural Optimizations
    - ▶ Semantic Query Optimizations
- ▶ **Contribution of this work:**
  - ▷ Data-dependent Optimizations
    - ▶ Explore a space of alternatives for unfoldings
    - ▶ Choose the best alternative, given the instance

## Unfolding and Structural Optimizations I

### Mapping $\mathcal{M}$

#### % High-Level Mappings *target* $\leftarrow$ *source view*

$URI_1(\vec{x})$ <i>a</i> <b>A</b>	$\leftarrow$	$V_1(\vec{x})$
$URI_2(\vec{x})$ <i>a</i> <b>A</b>	$\leftarrow$	$V_2(\vec{x})$
$URI_1(\vec{x})$ <i>R</i> $URI_3(\vec{y})$	$\leftarrow$	$V_3(\vec{x}, \vec{y})$
$URI_1(\vec{x})$ <i>R</i> $URI_4(\vec{y})$	$\leftarrow$	$V_4(\vec{x}, \vec{y})$

#### % Low-Level Mappings *source view* $\leftarrow$ *source query*

$V_1(\vec{x})$	$\leftarrow$	$CQ_1(\vec{x})$
$V_2(\vec{x})$	$\leftarrow$	$CQ_2(\vec{x})$
$V_3(\vec{x}, \vec{y})$	$\leftarrow$	$CQ_3(\vec{x}, \vec{y})$
$V_4(\vec{x}, \vec{y})$	$\leftarrow$	$CQ_4(\vec{x}, \vec{y})$

## Unfolding and Structural Optimizations I

### Mapping $\mathcal{M}$

**% High-Level Mappings**      *target*  $\leftarrow$  *source view*

$URI_1(\vec{x})$ a A	$\leftarrow$	$V_1(\vec{x})$
$URI_2(\vec{x})$ a A	$\leftarrow$	$V_2(\vec{x})$
$URI_1(\vec{x})$ R $URI_3(\vec{y})$	$\leftarrow$	$V_3(\vec{x}, \vec{y})$
$URI_1(\vec{x})$ R $URI_4(\vec{y})$	$\leftarrow$	$V_4(\vec{x}, \vec{y})$

**% Low-Level Mappings**      *source view*  $\leftarrow$  *source query*

$V_1(\vec{x})$	$\leftarrow$	$CQ_1(\vec{x})$
$V_2(\vec{x})$	$\leftarrow$	$CQ_2(\vec{x})$
$V_3(\vec{x}, \vec{y})$	$\leftarrow$	$CQ_3(\vec{x}, \vec{y})$
$V_4(\vec{x}, \vec{y})$	$\leftarrow$	$CQ_4(\vec{x}, \vec{y})$

### Query $q$

**SELECT**  $?x$   $?y$  **WHERE** {  $?x$  a A.  $?x$  R  $?y$ . }

## Unfolding and Structural Optimizations II

Query  $q$

$SELECT ?x ?y WHERE \{ ?x \text{ a } A. ?x R ?y. \}$

### SELECT CONCAT

We introduce auxiliary notation:

$$\begin{aligned}
 SC_1(V_1) &= SELECT CONCAT('URI_1(', \vec{x}, ')') AS x FROM V_1 \\
 &\vdots \\
 SC_{1,4}(V_4) &= SELECT CONCAT('URI_1(', \vec{x}, ')') AS x, \\
 &\quad CONCAT('URI_4(', \vec{y}, ')') AS y FROM V_4
 \end{aligned}$$





## Unfolding and Structural Optimizations III

### *Structural Optimization 1: Transform to UCQ*

$SC_1(V_1) \bowtie SC_{1,3}(V_3)$

$SC_1(V_1) \bowtie SC_{1,4}(V_4)$

$SC_2(V_2) \bowtie SC_{1,3}(V_3)$

$SC_2(V_2) \bowtie SC_{1,4}(V_4)$

## Unfolding and Structural Optimizations III

### Structural Optimization 1: Transform to UCQ

$SC_1(V_1)$	$\bowtie$	$SC_{1,3}(V_3)$	
$SC_1(V_1)$	$\bowtie$	$SC_{1,4}(V_4)$	
$SC_2(V_2)$	$\bowtie$	$SC_{1,3}(V_3)$	} $URI_2 \neq URI_1$ !! (Incompatible URIs)
$SC_2(V_2)$	$\bowtie$	$SC_{1,4}(V_4)$	

## Unfolding and Structural Optimizations III

### Structural Optimization 1: Transform to UCQ

$SC_1(V_1)$	$\bowtie$	$SC_{1,3}(V_3)$	
$SC_1(V_1)$	$\bowtie$	$SC_{1,4}(V_4)$	
$SC_2(V_2)$	$\bowtie$	$SC_{1,3}(V_3)$	} $URI_2 \neq URI_1$ !! (Incompatible URIs)
$SC_2(V_2)$	$\bowtie$	$SC_{1,4}(V_4)$	

### Structural Optimization 2: Remove joins between incompatible URIs

$SC_1(V_1)$	$\bowtie$	$SC_{1,3}(V_3)$
$SC_1(V_1)$	$\bowtie$	$SC_{1,4}(V_4)$

## Unfolding and Structural Optimizations III

### Structural Optimization 1: Transform to UCQ

$$\begin{array}{l}
 \mathbf{SC}_1(V_1) \bowtie \mathbf{SC}_{1,3}(V_3) \\
 \mathbf{SC}_1(V_1) \bowtie \mathbf{SC}_{1,4}(V_4) \\
 \mathbf{SC}_2(V_2) \bowtie \mathbf{SC}_{1,3}(V_3) \\
 \mathbf{SC}_2(V_2) \bowtie \mathbf{SC}_{1,4}(V_4)
 \end{array}
 \left. \vphantom{\begin{array}{l} \mathbf{SC}_1(V_1) \bowtie \mathbf{SC}_{1,3}(V_3) \\ \mathbf{SC}_1(V_1) \bowtie \mathbf{SC}_{1,4}(V_4) \\ \mathbf{SC}_2(V_2) \bowtie \mathbf{SC}_{1,3}(V_3) \\ \mathbf{SC}_2(V_2) \bowtie \mathbf{SC}_{1,4}(V_4) \end{array}} \right\} \text{URI}_2 \neq \text{URI}_1 \text{ !! (Incompatible URIs)}$$

### Structural Optimization 2: Remove joins between incompatible URIs

$$\begin{array}{l}
 \mathbf{SC}_1(V_1) \bowtie \mathbf{SC}_{1,3}(V_3) \\
 \mathbf{SC}_1(V_1) \bowtie \mathbf{SC}_{1,4}(V_4)
 \end{array}$$

### Structural Optimization 3: Remove Joins over templates

$$\mathbf{SC}_{1,3}(V_1 \bowtie V_3) \cup \mathbf{SC}_{1,4}(V_1 \bowtie V_4)$$

## Unfolding and Structural Optimizations III

### Structural Optimization 1: Transform to UCQ

$$\begin{array}{l}
 SC_1(V_1) \bowtie SC_{1,3}(V_3) \\
 SC_1(V_1) \bowtie SC_{1,4}(V_4) \\
 SC_2(V_2) \bowtie SC_{1,3}(V_3) \\
 SC_2(V_2) \bowtie SC_{1,4}(V_4)
 \end{array}
 \left. \vphantom{\begin{array}{l} SC_1(V_1) \\ SC_1(V_1) \\ SC_2(V_2) \\ SC_2(V_2) \end{array}} \right\} URI_2 \neq URI_1 \text{ !! (Incompatible URIs)}$$

### Structural Optimization 2: Remove joins between incompatible URIs

$$\begin{array}{l}
 SC_1(V_1) \bowtie SC_{1,3}(V_3) \\
 SC_1(V_1) \bowtie SC_{1,4}(V_4)
 \end{array}$$

### Structural Optimization 3: Remove Joins over templates

$$SC_{1,3}(V_1 \bowtie V_3) \cup SC_{1,4}(V_1 \bowtie V_4)$$

We call the unfolding after structural optimizations **UCQ-unfolding**

## Trade-off Redundancy/Indexing

### Mapping $\mathcal{M}$

**% High-Level Mappings**     *target*  $\leftarrow$  *source*

$$\{ \text{URI}_1(\vec{x}) \text{ a } \mathbf{A} \leftarrow \text{V}_i(\vec{x}) \mid 1 \leq i \leq 6 \}$$

$$\{ \text{URI}_1(\vec{x}) \text{ P } \text{URI}_2(\vec{y}) \leftarrow \text{W}_i(\vec{x}, \vec{y}) \mid 1 \leq i \leq 6 \}$$

$$\{ \text{URI}_1(\vec{x}) \text{ R } \text{URI}_3(\vec{z}) \leftarrow \text{Z}_i(\vec{x}, \vec{z}) \mid 1 \leq i \leq 6 \}$$

## Trade-off Redundancy/Indexing

### Mapping $\mathcal{M}$

**% High-Level Mappings**     *target*  $\rightsquigarrow$  *source*

$\{ \text{URI}_1(\vec{x}) \text{ a } \mathbf{A} \rightsquigarrow \mathbf{V}_i(\vec{x}) \mid 1 \leq i \leq 6 \}$

$\{ \text{URI}_1(\vec{x}) \text{ P } \text{URI}_2(\vec{y}) \rightsquigarrow \mathbf{W}_i(\vec{x}, \vec{y}) \mid 1 \leq i \leq 6 \}$

$\{ \text{URI}_1(\vec{x}) \text{ R } \text{URI}_3(\vec{z}) \rightsquigarrow \mathbf{Z}_i(\vec{x}, \vec{z}) \mid 1 \leq i \leq 6 \}$

- ▶ SELECT DISTINCT \* WHERE { ?x rdf:type A. ?x P ?y. ?x R ?z }



## Trade-off Redundancy/Indexing

### Mapping $\mathcal{M}$

**% High-Level Mappings**     *target*  $\rightsquigarrow$  *source*

$\{ \text{URI}_1(\vec{x}) \text{ a A} \rightsquigarrow V_i(\vec{x}) \mid 1 \leq i \leq 6 \}$

$\{ \text{URI}_1(\vec{x}) \text{ P URI}_2(\vec{y}) \rightsquigarrow W_i(\vec{x}, \vec{y}) \mid 1 \leq i \leq 6 \}$

$\{ \text{URI}_1(\vec{x}) \text{ R URI}_3(\vec{z}) \rightsquigarrow Z_i(\vec{x}, \vec{z}) \mid 1 \leq i \leq 6 \}$

▶ SELECT DISTINCT \* WHERE { ?x rdf:type A. ?x P ?y. ?x R ?z }

▶ **UCQ-unfolding:**  $\bigcup_{i,j,k \in \{1, \dots, 6\}} \text{SC}(V_i \bowtie W_j \bowtie Z_k)$

## Trade-off Redundancy/Indexing

### Mapping $\mathcal{M}$

% High-Level Mappings      target  $\rightsquigarrow$  source

$\{ \text{URI}_1(\vec{x}) \text{ a A} \rightsquigarrow V_i(\vec{x}) \mid 1 \leq i \leq 6 \}$

$\{ \text{URI}_1(\vec{x}) \text{ P URI}_2(\vec{y}) \rightsquigarrow W_i(\vec{x}, \vec{y}) \mid 1 \leq i \leq 6 \}$

$\{ \text{URI}_1(\vec{x}) \text{ R URI}_3(\vec{z}) \rightsquigarrow Z_i(\vec{x}, \vec{z}) \mid 1 \leq i \leq 6 \}$

▶ SELECT DISTINCT \* WHERE { ?x rdf:type A. ?x P ?y. ?x R ?z }

▶ **UCQ-unfolding:**  $\bigcup_{i,j,k \in \{1, \dots, 6\}} \text{SC}(V_i \bowtie W_j \bowtie Z_k)$

▶  $6^3 = 216$  CQs

## Trade-off Redundancy/Indexing

### Mapping $\mathcal{M}$

% High-Level Mappings      target  $\rightsquigarrow$  source

$\{ \text{URI}_1(\vec{x}) \text{ a A} \rightsquigarrow V_i(\vec{x}) \mid 1 \leq i \leq 6 \}$

$\{ \text{URI}_1(\vec{x}) \text{ P URI}_2(\vec{y}) \rightsquigarrow W_i(\vec{x}, \vec{y}) \mid 1 \leq i \leq 6 \}$

$\{ \text{URI}_1(\vec{x}) \text{ R URI}_3(\vec{z}) \rightsquigarrow Z_i(\vec{x}, \vec{z}) \mid 1 \leq i \leq 6 \}$

▶ SELECT DISTINCT \* WHERE { ?x rdf:type A. ?x P ?y. ?x R ?z }

▶ **UCQ-unfolding:**  $\bigcup_{i,j,k \in \{1, \dots, 6\}} \text{SC}(V_i \bowtie W_j \bowtie Z_k)$

▶  $6^3 = 216$  CQs

▶ **Redundancy?**

## Trade-off Redundancy/Indexing

### Mapping $\mathcal{M}$

**% High-Level Mappings**      *target*  $\leftarrow$  *source*

$\{ \text{URI}_1(\vec{x}) \text{ a A} \leftarrow V_i(\vec{x}) \mid 1 \leq i \leq 6 \}$

$\{ \text{URI}_1(\vec{x}) \text{ P URI}_2(\vec{y}) \leftarrow W_i(\vec{x}, \vec{y}) \mid 1 \leq i \leq 6 \}$

$\{ \text{URI}_1(\vec{x}) \text{ R URI}_3(\vec{z}) \leftarrow Z_i(\vec{x}, \vec{z}) \mid 1 \leq i \leq 6 \}$

- ▶ SELECT DISTINCT \* WHERE { ?x rdf:type A. ?x P ?y. ?x R ?z }
- ▶ **UCQ-unfolding:**  $\bigcup_{i,j,k \in \{1, \dots, 6\}} \text{SC}(V_i \bowtie W_j \bowtie Z_k)$ 
  - ▶  $6^3 = 216$  CQs
  - ▶ **Redundancy?**
- ▶ SELECT DISTINCT \* WHERE { { ?x rdf:type A. ?x P ?y } { ?x rdf:type A. ?x R ?z } }

## Trade-off Redundancy/Indexing

### Mapping $\mathcal{M}$

% High-Level Mappings      *target*  $\rightsquigarrow$  *source*

$\{ \text{URI}_1(\vec{x}) \text{ a A} \rightsquigarrow V_i(\vec{x}) \mid 1 \leq i \leq 6 \}$

$\{ \text{URI}_1(\vec{x}) \text{ P URI}_2(\vec{y}) \rightsquigarrow W_i(\vec{x}, \vec{y}) \mid 1 \leq i \leq 6 \}$

$\{ \text{URI}_1(\vec{x}) \text{ R URI}_3(\vec{z}) \rightsquigarrow Z_i(\vec{x}, \vec{z}) \mid 1 \leq i \leq 6 \}$

▶ SELECT DISTINCT \* WHERE  $\{ ?x \text{ rdf:type A. } ?x \text{ P } ?y. ?x \text{ R } ?z \}$

▶ **UCQ-unfolding:**  $\bigcup_{i,j,k \in \{1, \dots, 6\}} \text{SC}(V_i \bowtie W_j \bowtie Z_k)$

▶  $6^3 = 216$  CQs

▶ **Redundancy?**

▶ SELECT DISTINCT \* WHERE  $\{ \{ ?x \text{ rdf:type A. } ?x \text{ P } ?y \} \{ ?x \text{ rdf:type A. } ?x \text{ R } ?z \} \}$

▶ **Join of UCQ-unfoldings:**  $(\bigcup_{i,j \in \{1, \dots, 6\}} \text{SC}(V_i \bowtie W_j)) \bowtie (\bigcup_{i,k \in \{1, \dots, 6\}} \text{SC}(V_i \bowtie Z_k))$

## Trade-off Redundancy/Indexing

### Mapping $\mathcal{M}$

% High-Level Mappings      *target*  $\rightsquigarrow$  *source*

$\{ \text{URI}_1(\vec{x}) \text{ A } \mathbf{A} \rightsquigarrow V_i(\vec{x}) \mid 1 \leq i \leq 6 \}$

$\{ \text{URI}_1(\vec{x}) \text{ P } \text{URI}_2(\vec{y}) \rightsquigarrow W_i(\vec{x}, \vec{y}) \mid 1 \leq i \leq 6 \}$

$\{ \text{URI}_1(\vec{x}) \text{ R } \text{URI}_3(\vec{z}) \rightsquigarrow Z_i(\vec{x}, \vec{z}) \mid 1 \leq i \leq 6 \}$

▶ SELECT DISTINCT \* WHERE  $\{ ?x \text{ rdf:type } \mathbf{A}. ?x \text{ P } ?y. ?x \text{ R } ?z \}$

▶ **UCQ-unfolding:**  $\bigcup_{i,j,k \in \{1, \dots, 6\}} \text{SC}(V_i \bowtie W_j \bowtie Z_k)$

▶  $6^3 = 216$  CQs

▶ **Redundancy?**

▶ SELECT DISTINCT \* WHERE  $\{ \{ ?x \text{ rdf:type } \mathbf{A}. ?x \text{ P } ?y \} \{ ?x \text{ rdf:type } \mathbf{A}. ?x \text{ R } ?z \} \}$

▶ **Join of UCQ-unfoldings:**  $(\bigcup_{i,j \in \{1, \dots, 6\}} \text{SC}(V_i \bowtie W_j)) \bowtie (\bigcup_{i,k \in \{1, \dots, 6\}} \text{SC}(V_i \bowtie Z_k))$

▶  $6^2 + 6^2 = 64$  CQs

## Trade-off Redundancy/Indexing

### Mapping $\mathcal{M}$

% High-Level Mappings      *target*  $\leftarrow$  *source*

$\{ \text{URI}_1(\vec{x}) \text{ a } \mathbf{A} \leftarrow V_i(\vec{x}) \mid 1 \leq i \leq 6 \}$

$\{ \text{URI}_1(\vec{x}) \text{ P } \text{URI}_2(\vec{y}) \leftarrow W_i(\vec{x}, \vec{y}) \mid 1 \leq i \leq 6 \}$

$\{ \text{URI}_1(\vec{x}) \text{ R } \text{URI}_3(\vec{z}) \leftarrow Z_i(\vec{x}, \vec{z}) \mid 1 \leq i \leq 6 \}$

▶ SELECT DISTINCT \* WHERE  $\{ ?x \text{ rdf:type } \mathbf{A}. ?x \text{ P } ?y. ?x \text{ R } ?z \}$

▶ **UCQ-unfolding:**  $\bigcup_{i,j,k \in \{1, \dots, 6\}} \text{SC}(V_i \bowtie W_j \bowtie Z_k)$

▶  $6^3 = 216$  CQs

▶ **Redundancy?**

▶ SELECT DISTINCT \* WHERE  $\{ \{ ?x \text{ rdf:type } \mathbf{A}. ?x \text{ P } ?y \} \{ ?x \text{ rdf:type } \mathbf{A}. ?x \text{ R } ?z \} \}$

▶ **Join of UCQ-unfoldings:**  $(\bigcup_{i,j \in \{1, \dots, 6\}} \text{SC}(V_i \bowtie W_j)) \bowtie (\bigcup_{i,k \in \{1, \dots, 6\}} \text{SC}(V_i \bowtie Z_k))$

▶  $6^2 + 6^2 = 64$  CQs

▶ **Indexes?**

## Key Ideas

- ▶ The UCQ form for the unfolding is a “one-fits-all” solution
  - ▶ However, other forms could be more efficient over certain data instances
- ▶ Ideally, one would want to choose the best form for the given OBDA setting<sup>1</sup>
- ▶ This is typical in databases, where **cost models** are used to **select the best among different candidate query execution plans**
- ▶ [Bursztyn et al., 2015] adapted such cost-based solutions to the OBQA<sup>2</sup> context
- ▶ We here study how to adapt such techniques to the OBDA scenario

---

<sup>1</sup>OBDA setting: ontology, mappings, database schema and database instance.

<sup>2</sup>Ontology-based Query Answering. No mappings.



## Aspects to Consider

- ▶ **Adapting such techniques requires to consider several aspects:**
  - ▷ **Alternative unfoldings and structural optimizations**
  - ▷ **Additional information available in an OBDA input**
  - ▷ **More complex, structured, search space**
  - ▷ **Interaction of the mappings and data with the ontology**

## Aspects to Consider

- ▶ **Adapting such techniques requires to consider several aspects:**
  - ▷ **Alternative unfoldings and structural optimizations**
  - ▷ **Additional information available in an OBDA input**
  - ▷ **More complex, structured, search space**
  - ▷ **Interaction of the mappings and data with the ontology**

## Collecting Statistics (for Properties)

- ▶ Consider mappings  $\mathcal{M}$  and a data instance  $\mathcal{D}$
- ▶ Let  $URI_1(\vec{x}) \text{ P } URI_2(\vec{y}) \leftarrow V(\vec{x}, \vec{y})$  be a mapping in  $\mathcal{M}$
- ▶ We collect:
  - ▷  $|V^{\mathcal{D}}|$  (S1)
  - ▷  $|\pi_{\vec{x}}(V^{\mathcal{D}})|, |\pi_{\vec{y}}(V^{\mathcal{D}})|$  (S2)
- ▶ If  $URI_1(\vec{x}) \text{ R } URI_3(\vec{w}) \leftarrow W(\vec{x}, \vec{w})$  belongs to  $\mathcal{M}$ , then we also collect:
  - ▷  $|\pi_{\vec{x}}(V^{\mathcal{D}}) \cap \pi_{\vec{x}}(W^{\mathcal{D}})|$  (S3)

## Cardinality Estimation

- ▶ For a single join:

$$f_{\mathcal{D}}(V \bowtie_{V.\bar{x}=W.\bar{y}} W) = \left[ |(V \bowtie_{V.\bar{x}=W.\bar{y}} W)^{\mathcal{D}}| \cdot \frac{|V^{\mathcal{D}}|}{|\pi_{\bar{x}} V^{\mathcal{D}}|} \cdot \frac{|W^{\mathcal{D}}|}{|\pi_{\bar{y}} W^{\mathcal{D}}|} \right]$$

- ▶ **Statistics Used:** (S1), (S2), (S3)

- ▶ **Assumptions:**

- ▶ **Uniformity:**  $P(C = v_1) = P(C = v_2)$ , for all  $v_1, v_2 \in C$
- ▶ **Independence:**  $P(C_1 = v_1 | C_2 = v_2) = P(C_1 = v_1)$ , for all  $v_1 \in C_1$  and  $v_2 \in C_2$

## Cardinality Estimation for Basic CQs

Given a basic CQ  $E'$ ,  $f_{\mathcal{D}}(E')$  estimates the number  $|E'^{\mathcal{D}}|$  of distinct results in the evaluation of  $E'$  over  $\mathcal{D}$ . We define it as:

$$f_{\mathcal{D}}(E \bowtie_{V_{[p]}. \vec{x} = W_{[q]}. \vec{y}} W_{[q]}) = \begin{cases} \left[ \frac{k_{\mathcal{D}}(V_{[p]} \bowtie_{V_{[p]}. \vec{x} = W_{[q]}. \vec{y}} W_{[q]}) \cdot |V^{\mathcal{D}}| \cdot |W^{\mathcal{D}}|}{\text{dist}_{\mathcal{D}}(V, V_{[p]}. \vec{x}) \cdot \text{dist}_{\mathcal{D}}(W, W_{[q]}. \vec{y})} \right], & \text{if } E = V \\ \left[ \frac{k_{\mathcal{D}}(E \bowtie_{V_{[p]}. \vec{x} = W_{[q]}. \vec{y}} W_{[q]}) \cdot f_{\mathcal{D}}(E) \cdot |W^{\mathcal{D}}|}{\text{dist}_{\mathcal{D}}(E, V_{[p]}. \vec{x}) \cdot \text{dist}_{\mathcal{D}}(W, W_{[q]}. \vec{y})} \right], & \text{otherwise.} \end{cases}$$

## Cost Model

### Cost for the Unfolding of a UCQ

$$c(\mathbf{q}^{ucq}) = \sum_i c(\mathbf{q}_i^{cq}) + c_u(\mathbf{q}^{ucq})$$

- ▶  $c(\mathbf{q}_i^{cq})$  is the cost of evaluating each  $\mathbf{q}_i^{cq}$  in  $\mathbf{q}^{ucq}$
- ▶  $c_u(\mathbf{q}^{ucq})$  is the cost of removing duplicate results

### Cost for the Unfolding of a Join of UCQs (JUCQ)

$$c(\mathbf{q}^{jucq}) = \sum_i c(\mathbf{q}_i^{ucq}) + \sum_{i \neq k} c_{mat}(\mathbf{q}_i^{ucq}) + c_{mj}(\mathbf{q}^{jucq}) + c_u(\mathbf{q}^{jucq})$$

- ▶  $c(\mathbf{q}_i^{ucq})$  is the cost of evaluating each UCQ component  $\mathbf{q}_i^{ucq}$
- ▶  $\sum_{i \neq k} c_{mat}(\mathbf{q}_i^{ucq})$  is the cost of materializing the intermediate results from  $\mathbf{q}_i^{ucq}$
- ▶  $c_{mj}(\mathbf{q}^{jucq})$  is the cost of a merge join over the materialized intermediate results
- ▶  $c_u(\mathbf{q}^{jucq})$  is the cost of removing duplicate results

## Empirical Evaluation

- ▶ **Test suites:** System testing (OBDA extension of the Wisconsin Benchmark) and application-driven (NPD Benchmark)
- ▶ **Hardware:** HP Proliant 24 cores@3.47GHz; 106 GB of RAM
- ▶ **RDBMS:** PostgreSQL 9.6

## Wisconsin Experiment

► **Comparison between:**

▷ `SELECT DISTINCT * WHERE { ?x :mrP1 ?y. ?x :jmrP2 ?z. ?x :jmrP3 ?w }` ( $q_{ucq}$ )

▷ `SELECT DISTINCT ?x ?y ?z ?w WHERE  
 { { ?x :mrP1 ?y. ?x :jmrP2 ?z } { ?x :mrP1 ?y1. ?x :jmrP3 ?w } }` ( $q_{jucq}$ )

► **Parameters:**

▷  **$m$ : number of mappings**, in  $\{1, \dots, 6\}$

▷  **$r$ : number of redundant mappings<sup>3</sup>**, in  $\{1, \dots, m - 1\}$

▷  **$j$ : join-selectivity<sup>4</sup>**, in  $\{5, 10, 15, 20\}$

► **All intra-BGP joins are over indexed values**

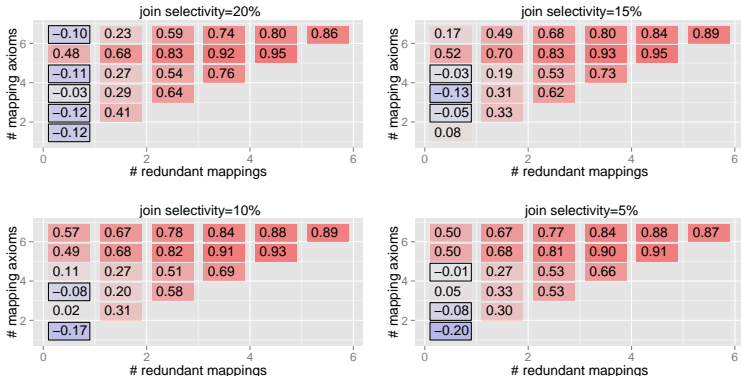
<sup>3</sup>Mappings which do not contribute to new assertions

<sup>4</sup>% of the number of retrieved rows for each mapping defining the property (each mapping retrieves **200k** rows)



## Analysis of the Trade-Off Redundancy/Indexing

- ▶ **Blue color:**  $q^{ucq}$  is better
- ▶ **Red color:**  $q^{jucq}$  is better
- ▶ **Numbers in cells:**  $1 - (time(q^{jucq})/time(q^{ucq}))$



**Figure:** Performance gain of JUCQ compared with UCQ

# Cost

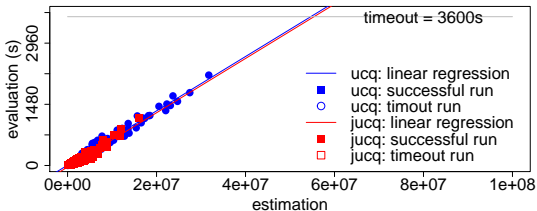


Figure: Our cost model vs. evaluation

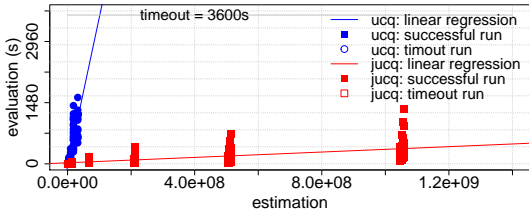


Figure: PostgreSQL cost model vs. evaluation

# Evaluation of PostgreSQL Cardinality Estimation (UCQs)

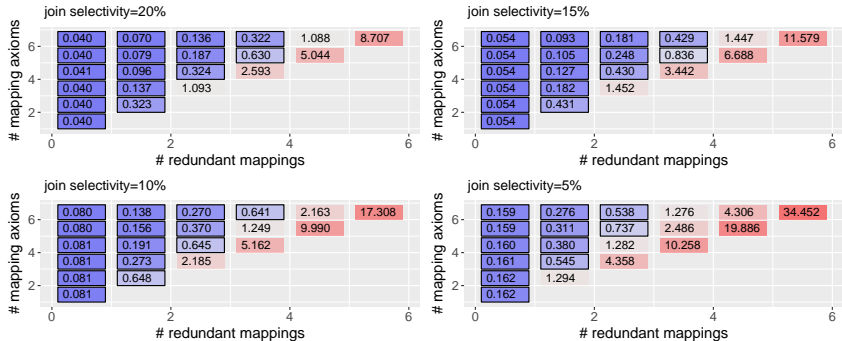
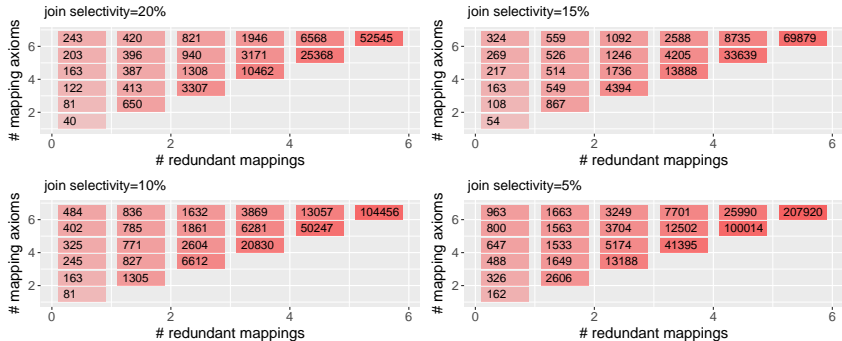


Figure: UCQs: (PostgreSQL estimated cardinality) / (real cardinality)

# Evaluation of PostgreSQL Cardinality Estimation (JUCQs)



**Figure:** JUCQs: (PostgreSQL estimated cardinality) / (real cardinality)

## NPD Benchmark Experiment

- ▶ The NPD benchmark is an OBDA-specific benchmark aimed at reproducing the use of an OBDA system in an industrial real world scenario

**Table:** Evaluation over the NPD benchmark

SPARQL Query		Unfolding for UCQs		Unfolding for JUCQs		
name	# triple patterns	time (s)	# CQs	time (s)	# Frags	# CQs
<i>q<sub>6</sub></i>	7	2.18	48	1.20	2	14
<i>q<sub>11</sub></i>	8	3.39	24	0.40	2	12
<i>q<sub>12</sub></i>	10	6.67	48	0.47	2	14
<i>q<sub>31</sub></i>	10	54.27	3840	1.58	2	327

## Conclusion and Future Work

### ▶ **Contribution:**

- ▶ We studied how to find the most efficient unfolding of a user query
  - ▶▶ We introduced a space of unfolding alternatives
  - ▶▶ We defined structural optimizations for the new alternatives
  - ▶▶ We devised a cost model to search for the best alternative
- ▶ We performed two empirical evaluations, and observed that:
  - ▶▶ Alternative unfoldings can be orders of magnitude more efficient than traditional UCQ unfoldings
  - ▶▶ The devised cost model is well suited to select the best unfolding

### ▶ **Future work:**

- ▶ Integrate cost estimation with semantic optimizations
- ▶ Explore more unfolding alternatives
- ▶ Relax estimation assumptions by integrating our model with existing techniques based on histograms

## References I

- [Bursztyn et al., 2015] **Bursztyn, D., Goasdoué, F., and Manolescu, I. (2015).**  
**Reformulation-based query answering in RDF: alternatives and performance.**  
*Proc. of VLDB*, 8(12):1888–1891.
- [Chebotko et al., 2009] **Chebotko, A., Lu, S., and Fotouhi, F. (2009).**  
**Semantics preserving sparql-to-sql translation.**  
*Data & Knowledge Engineering*, 68(10):973 – 1000.