

Incomplete Information in RDF

Charalampos Nikolaou and Manolis Koubarakis

charnik@di.uoa.gr

koubarak@di.uoa.gr

Department of Informatics and Telecommunications
National and Kapodistrian University of Athens



Web Reasoning and Rule Systems (RR) 2013
July 27, 2013

Outline

Motivation

Previous work

The RDFⁱ framework

SPARQL query evaluation over RDFⁱ databases

An algorithm for certain answer computation

Preliminary complexity results

Conclusions and future work



Motivation

- ▶ Incomplete information is an important issue in many research areas: relational databases, knowledge representation and the semantic web.
- ▶ Incomplete information arises in many practical settings (e.g., sensor data). RDF is often used to represent such data.
- ▶ Even if initial information is complete, incomplete information arises later on (e.g., relational view updates, data integration, data exchange).
- ▶ Although there is much work recently on incomplete information in XML, not much has been done for incomplete information in RDF.



Previous work

Relational

- ▶ Relations extended to tables with various models of incompleteness [Imielinski/Lipski '84]
- ▶ Complexity results for the associated decision problems [Abiteboul/Kanellakis/Grahne '91]
- ▶ Dependencies and updates [Grahne '91]



Previous work

Relational

- ▶ Relations extended to tables with various models of incompleteness [Imielinski/Lipski '84]
- ▶ Complexity results for the associated decision problems [Abiteboul/Kanellakis/Grahne '91]
- ▶ Dependencies and updates [Grahne '91]

XML

- ▶ Dynamic enrichment of incomplete information [Abiteboul/Segoufin/Vianu '01,'06]
- ▶ General models of incompleteness, query answering, and computational complexity [Barceló/Libkin/Poggi/Sirangelo '09,'10]



Previous work (cont'd)

RDF

- ▶ Blank nodes as existential variables in the RDF standard
- ▶ SPARQL query evaluation under certain answer semantics (Open World Assumption) [Arenas/Pérez '11]
- ▶ Anonymous timestamps in general temporal RDF graphs [Gutierrez/Hurtado/Vaisman '05]
- ▶ General temporal RDF graphs with temporal constraints [Hurtado/Vaisman '06]



Previous work (cont'd)

RDF

- ▶ Blank nodes as existential variables in the RDF standard
- ▶ SPARQL query evaluation under certain answer semantics (Open World Assumption) [Arenas/Pérez '11]
- ▶ Anonymous timestamps in general temporal RDF graphs [Gutierrez/Hurtado/Vaisman '05]
- ▶ General temporal RDF graphs with temporal constraints [Hurtado/Vaisman '06]

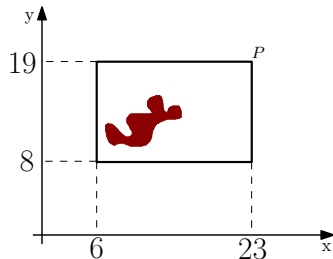
RDFⁱ: It captures incomplete information for property values using constraints. It is for RDF what the c-tables model is for the relational model.



RDFⁱ by example

Example

```
hotspot1    type      Hotspot .  
  fire1     type      Fire    .  
hotspot1   correspondsTo  fire1  .  
  fire1     occurredIn  R1    .
```

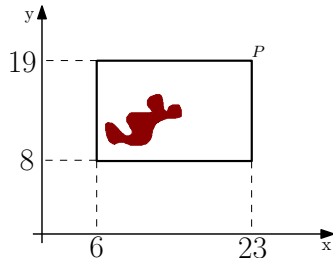


RDFⁱ by example

Example

```
hotspot1    type      Hotspot .  
  fire1     type      Fire    .  
hotspot1   correspondsTo  fire1  .  
  fire1     occurredIn  R1    .
```

R1 NTPP " $x \geq 6 \wedge x \leq 23 \wedge y \geq 8 \wedge y \leq 19$ "



RDFⁱ in a nutshell

- ▶ Extension of RDF for capturing **incomplete information** for property values that **exist** but are **unknown** or **partially known**
- ▶ Partial knowledge captured by **constraints** using an appropriate constraint language \mathcal{L} interpreted over a fixed structure $\mathbf{M}_{\mathcal{L}}$

Syntax

RDF graphs extended to RDFⁱ databases: pair (G, ϕ)

- ▶ G : RDF graph with a new kind of literals, called **e-literals**
- ▶ ϕ : quantifier-free formula of \mathcal{L}

Semantics

- ▶ Possible world semantics as in [Imielinski/Lipski '84] and [Grahne '91]



Constraint languages \mathcal{L}

Examples

ECL

- ▶ **Equality constraints**
interpreted over an infinite domain: $x \text{ EQ } y, x \text{ EQ } c$
- ▶ Blank nodes as existential variables



Constraint languages \mathcal{L}

Examples

ECL

- ▶ **Equality constraints** interpreted over an infinite domain: $x \text{ EQ } y, x \text{ EQ } c$
- ▶ Blank nodes as existential variables

diPCL/dePCL

- ▶ **Difference constraints** of the form $x - y \leq c$ interpreted over the integers or rationals
- ▶ Incomplete temporal information [Koubarakis '94]



Constraint languages \mathcal{L}

Examples

ECL

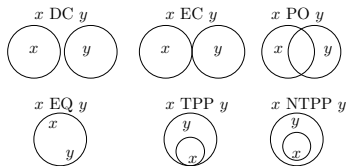
- ▶ **Equality constraints** interpreted over an infinite domain: $x \text{ EQ } y, x \text{ EQ } c$
- ▶ Blank nodes as existential variables

diPCL/dePCL

- ▶ **Difference constraints** of the form $x - y \leq c$ interpreted over the integers or rationals
- ▶ Incomplete temporal information [Koubarakis '94]

TCL

- ▶ **Topological constraints** of non-empty, regular closed subsets of topological space
- ▶ Six binary predicates: DC, EC, PO, EQ, TPP, NTPP



Constraint languages \mathcal{L}

Examples

ECL

- ▶ **Equality constraints** interpreted over an infinite domain: $x \text{ EQ } y, x \text{ EQ } c$
- ▶ Blank nodes as existential variables

diPCL/dePCL

- ▶ **Difference constraints** of the form $x - y \leq c$ interpreted over the integers or rationals
- ▶ Incomplete temporal information [Koubarakis '94]

TCL

- ▶ **Topological constraints** of non-empty, regular closed subsets of topological space
- ▶ Six binary predicates: DC, EC, PO, EQ, TPP, NTPP

PCL

- ▶ **TCL plus constant symbols** representing polygons in \mathbb{Q}^2
- ▶ e.g.,
 $r \text{ NTPP } "x - y \geq 0 \wedge x \leq 1 \wedge y \geq 0"$



RDFⁱ: Vocabulary

RDF	RDF ⁱ	\mathcal{L}
I (IRIs)	I	
B (blank nodes)	B	
L (literals)	L	
	C (literals)	constants
	U (e-literals)	variables
M (datatype map)	M	
	A (datatypes)	set of sorts



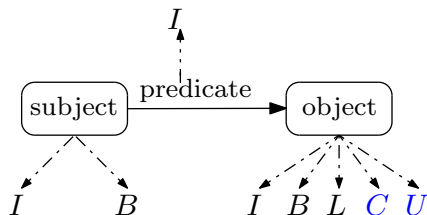
RDFⁱ: Vocabulary

RDF	RDF ⁱ	\mathcal{L}
I (IRIs)	I	
B (blank nodes)	B	
L (literals)	L	
	C (literals)	constants
	U (e-literals)	variables
M (datatype map)	M	
	A (datatypes)	set of sorts

$\mathbf{M}_{\mathcal{L}}$ interprets the constants of \mathcal{L} **in agreement with** function $L2V$ of M



RDFⁱ: Syntax



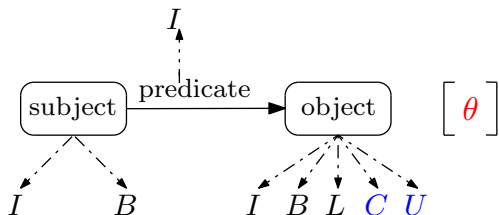
I : IRIs
 B : blank nodes
 L : literals
 C : constants of \mathcal{L}
 U : e-literals

Definition

- ▶ $(s, p, o) \in (I \cup B) \cup I \cup (I \cup B \cup L \cup C \cup U)$ is called an **e-triple**



RDFⁱ: Syntax



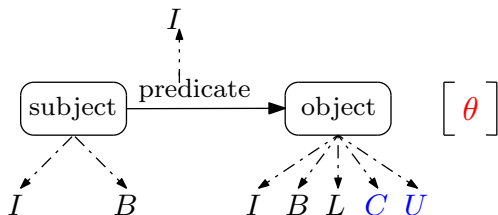
I : IRIs
 B : blank nodes
 L : literals
 C : constants of \mathcal{L}
 U : e-literals

Definition

- ▶ $(s, p, o) \in (I \cup B) \cup I \cup (I \cup B \cup L \cup C \cup U)$ is called an **e-triple**
- ▶ If t is an e-triple and θ a conjunction of \mathcal{L} -constraints, then the pair (t, θ) is called a **conditional triple**



RDFⁱ: Syntax



I : IRIs
 B : blank nodes
 L : literals
 C : constants of \mathcal{L}
 U : e-literals

Definition

- ▶ $(s, p, o) \in (I \cup B) \cup I \cup (I \cup B \cup L \cup C \cup U)$ is called an **e-triple**
- ▶ If t is an e-triple and θ a conjunction of \mathcal{L} -constraints, then the pair (t, θ) is called a **conditional triple**
- ▶ A set of conditional triples is called a **conditional graph**



RDFⁱ: Syntax (cont'd)

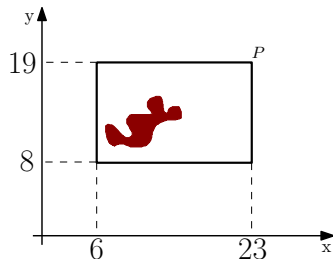
Definition

An RDFⁱ database D is a pair $D = (G, \phi)$ where G is a **conditional graph** and ϕ a Boolean combination of \mathcal{L} -constraints (**global constraint**)

Example

```
hotspot1    type      Hotspot .
  fire1     type      Fire   .
hotspot1 correspondsTo fire1 .
  fire1     occuredIn _R1   .

_R1 NTPP "x ≥ 6 ∧ x ≤ 23 ∧ y ≥ 8 ∧ y ≤ 19"
```



RDFⁱ: Semantics

RDFⁱ database

D

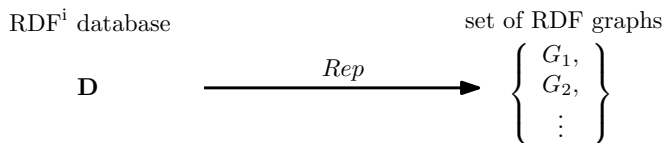
Rep

set of RDF graphs

$\left\{ \begin{array}{c} G_1, \\ G_2, \\ \vdots \end{array} \right\}$



RDFⁱ: Semantics



Definition

A **valuation** v is a function from U to C assigning to each e-literal from U a constant from C

Definition

Let G be a conditional graph and v a valuation. Then $v(G)$ denotes the RDF graph

$$\{v(t) \mid (t, \theta) \in G \text{ and } \mathbf{M}_{\mathcal{L}} \models v(\theta)\}$$



RDFⁱ: Semantics (cont'd)

From RDFⁱ databases to sets of RDF graphs

An RDFⁱ database $D = (G, \phi)$ corresponds to the following set of RDF graphs:

$$\text{Rep}(D) = \left\{ H \mid \text{there exists valuation } v \text{ and RDF graph } H \right. \\ \left. \text{such that } \mathbf{M}_{\mathcal{L}} \models v(\phi) \text{ and } H \supseteq v(G) \right\}$$

- ▶ Relation \supseteq captures the OWA semantics
- ▶ An RDFⁱ database corresponds to an **infinite** number of RDF graphs



Question

How can we evaluate a query q over an RDFⁱ database D (compute $\llbracket q \rrbracket_D$)?



Question

How can we evaluate a query q over an RDFⁱ database D (compute $\llbracket q \rrbracket_D$)?

Semantic definition

$$\llbracket q \rrbracket_{Rep(D)} = \{ \llbracket q \rrbracket_G \mid G \in Rep(D) \}$$



Question

How can we evaluate a query q over an RDFⁱ database D (compute $\llbracket q \rrbracket_D$)?

Semantic definition

$$\llbracket q \rrbracket_{Rep(D)} = \{ \llbracket q \rrbracket_G \mid G \in Rep(D) \}$$



Question

How can we evaluate a query q over an RDFⁱ database D (compute $\llbracket q \rrbracket_D$)?

Semantic definition

$$\llbracket q \rrbracket_{Rep(D)} = \{ \llbracket q \rrbracket_G \mid G \in Rep(D) \}$$

In practice?

- ▶ Start with SPARQL algebra of [Pérez/Arenas/Gutierrez '06] with set semantics
- ▶ Define SPARQL query evaluation for RDFⁱ databases



From mappings to e-mappings...

$\{?F \rightarrow \text{fire1}, ?S \rightarrow "x \geq 1 \wedge x \leq 2 \wedge y \geq 1 \wedge y \leq 2"\}$



From mappings to e-mappings...

$\{?F \rightarrow \text{fire1}, ?S \rightarrow "x \geq 1 \wedge x \leq 2 \wedge y \geq 1 \wedge y \leq 2"\}$

$\{?F \rightarrow \text{fire1}, ?S \rightarrow \underline{R1}\}$



... to conditional mappings

$\{?F \rightarrow \text{fire1}, ?S \rightarrow "x \geq 1 \wedge x \leq 2 \wedge y \geq 1 \wedge y \leq 2" \}$



... to conditional mappings

$\left(\{ ?F \rightarrow \text{fire1}, ?S \rightarrow "x \geq 1 \wedge x \leq 2 \wedge y \geq 1 \wedge y \leq 2" \}, \text{true} \right)$



... to conditional mappings

$\left(\{ ?F \rightarrow \text{fire1}, ?S \rightarrow \text{R1} \}, \text{R1 EQ "x} \geq 1 \wedge x \leq 2 \wedge y \geq 1 \wedge y \leq 2" \right)$



From compatible mappings to possibly compatible mappings

Join of conditional mappings

$(\{?F \rightarrow fire1, ?S \rightarrow _R1\}, _R1 \text{ EQ } "x \geq 1 \wedge x \leq 2 \wedge y \geq 1 \wedge y \leq 2")$

$(\{ \quad \quad \quad ?S \rightarrow _R2\}, true)$



From compatible mappings to possibly compatible mappings

Join of conditional mappings

$(\{?F \rightarrow fire1, ?S \rightarrow _R1\}, _R1 \text{ EQ } "x \geq 1 \wedge x \leq 2 \wedge y \geq 1 \wedge y \leq 2")$

$(\{ ?S \rightarrow _R2\}, true)$



From compatible mappings to possibly compatible mappings

Join of conditional mappings

$$\left(\{ ?F \rightarrow \text{fire1}, ?S \rightarrow _R1 \}, _R1 \text{ EQ } "x \geq 1 \wedge x \leq 2 \wedge y \geq 1 \wedge y \leq 2" \right)$$

⊗

$$\left(\{ ?S \rightarrow _R2 \}, \text{true} \right)$$

=

$$\left(\{ ?F \rightarrow \text{fire1}, ?S \rightarrow _R1 \}, \text{true} \wedge _R1 \text{ EQ } _R2 \wedge _R1 \text{ EQ } "x \geq 1 \wedge x \leq 2 \wedge y \geq 1 \wedge y \leq 2" \right)$$


Operations on conditional mappings

Let Ω_1 and Ω_2 be sets of **conditional mappings**. We can define the operation of:

- ▶ Join ($\Omega_1 \bowtie \Omega_2$)
- ▶ Union ($\Omega_1 \cup \Omega_2$)
- ▶ Difference ($\Omega_1 \setminus \Omega_2$)
- ▶ Left-outer join ($\Omega_1 \bowtie\!\!\!\!\!\! \bowtie \Omega_2$)



Graph pattern evaluation

If D is an RDFⁱ database and P a graph pattern, the **evaluation** of P over D is defined recursively:



Graph pattern evaluation

If D is an RDFⁱ database and P a graph pattern, the **evaluation** of P over D is defined recursively:

base case:

P is the triple pattern t

recursion:

P is $(P_1 \text{ AND } P_2)$	\rightarrow	$\llbracket P_1 \rrbracket_D \bowtie \llbracket P_2 \rrbracket_D$
P is $(P_1 \text{ UNION } P_2)$	\rightarrow	$\llbracket P_1 \rrbracket_D \cup \llbracket P_2 \rrbracket_D$
P is $(P_1 \text{ OPT } P_2)$	\rightarrow	$\llbracket P_1 \rrbracket_D \bowtie \llbracket P_2 \rrbracket_D$



Graph pattern evaluation

If D is an RDFⁱ database and P a graph pattern, the **evaluation** of P over D is defined recursively:

base case:

P is the triple pattern t

recursion:

P is $(P_1 \text{ AND } P_2)$ \rightarrow $\llbracket P_1 \rrbracket_D \bowtie \llbracket P_2 \rrbracket_D$

P is $(P_1 \text{ UNION } P_2)$ \rightarrow $\llbracket P_1 \rrbracket_D \cup \llbracket P_2 \rrbracket_D$

P is $(P_1 \text{ OPT } P_2)$ \rightarrow $\llbracket P_1 \rrbracket_D \bowtie \llbracket P_2 \rrbracket_D$

P is $(P_1 \text{ FILTER } R)$

where R is a conjunction of \mathcal{L} -constraints



Graph pattern evaluation

If D is an RDFⁱ database and P a graph pattern, the **evaluation** of P over D is defined recursively:

base case:

P is the triple pattern t

recursion:

P is $(P_1 \text{ AND } P_2)$ \rightarrow $\llbracket P_1 \rrbracket_D \bowtie \llbracket P_2 \rrbracket_D$

P is $(P_1 \text{ UNION } P_2)$ \rightarrow $\llbracket P_1 \rrbracket_D \cup \llbracket P_2 \rrbracket_D$

P is $(P_1 \text{ OPT } P_2)$ \rightarrow $\llbracket P_1 \rrbracket_D \bowtie \llbracket P_2 \rrbracket_D$

P is $(P_1 \text{ FILTER } R)$

where R is a conjunction of \mathcal{L} -constraints



Triple pattern evaluation (case 1)

Example

Database D

fire1 occurredIn _R1 .

_R1 NTPP " $x \geq 6 \wedge x \leq 23 \wedge y \geq 8 \wedge y \leq 19$ "

Query q

?F occurredIn ?R



Triple pattern evaluation (case 1)

Example

Database D

`fire1` occurredIn `_R1` .

`_R1` NTPP " $x \geq 6 \wedge x \leq 23 \wedge y \geq 8 \wedge y \leq 19$ "

Query q

`?F` occurredIn `?R`

Answer (set of conditional mappings)

$$\llbracket q \rrbracket_D = \left\{ \left(\{ ?F \rightarrow \text{fire1}, ?R \rightarrow _R1 \}, \text{true} \right) \right\}$$



Triple pattern evaluation (case 2)

Example

Database D

fire1 occurredIn $_R1$.

$_R1$ NTPP " $x \geq 6 \wedge x \leq 23 \wedge y \geq 8 \wedge y \leq 19$ "

Query q

?F occurredIn

" $x \geq 1 \wedge x \leq 2 \wedge y \geq 1 \wedge y \leq 2$ "



Triple pattern evaluation (case 2)

Example

Database D

fire1 occurredIn $_R1$.

$_R1$ NTTPP " $x \geq 6 \wedge x \leq 23 \wedge y \geq 8 \wedge y \leq 19$ "

Query q

?F occurredIn

" $x \geq 1 \wedge x \leq 2 \wedge y \geq 1 \wedge y \leq 2$ "

Answer (set of conditional mappings)

$$\llbracket q \rrbracket_D = \left\{ \left(\{ ?F \rightarrow \text{fire1} \}, _R1 \text{ EQ } "x \geq 1 \wedge x \leq 2 \wedge y \geq 1 \wedge y \leq 2" \right) \right\}$$



Evaluation of FILTER graph patterns

Example

Database D

`fire1` occurredIn `_R1` .

`_R1` NTTPP " $x \geq 6 \wedge x \leq 23 \wedge y \geq 8 \wedge y \leq 19$ "

Query q

`?F` occurredIn `?R` .

FILTER (`?R` NTTPP

" $x \geq 1 \wedge x \leq 2 \wedge y \geq 1 \wedge y \leq 2$ ")



Evaluation of FILTER graph patterns

Example

Database D

fire1 occurredIn $_R1$.

$_R1$ NTPP "x \geq 6 \wedge x \leq 23 \wedge y \geq 8 \wedge y \leq 19"

Query q

?F occurredIn ?R .

FILTER (?R NTPP
"x \geq 1 \wedge x \leq 2 \wedge y \geq 1 \wedge y \leq 2")

Answer

$$\llbracket q \rrbracket_D = \left\{ \left(\{ ?F \rightarrow \text{fire1}, ?R \rightarrow _R1 \}, \right. \right. \\ \left. \left. _R1 \text{ NTPP "x} \geq 1 \wedge x \leq 2 \wedge y \geq 1 \wedge y \leq 2" \right) \right\}$$



SELECT queries

Example

Database D

fire1 occurredIn _R1 .

_R1 NTPP " $x \geq 6 \wedge x \leq 23 \wedge y \geq 8 \wedge y \leq 19$ "

Query q

SELECT ?F

WHERE {

?F occurredIn ?R .

FILTER (?R NTPP

" $x \geq 1 \wedge x \leq 2 \wedge y \geq 1 \wedge y \leq 2$ ")}]



SELECT queries

Example

Database D

fire1 occurredIn $_R1$.

$_R1$ NTPP " $x \geq 6 \wedge x \leq 23 \wedge y \geq 8 \wedge y \leq 19$ "

Query q

SELECT ?F

WHERE {

?F occurredIn ?R .

FILTER (?R NTPP

" $x \geq 1 \wedge x \leq 2 \wedge y \geq 1 \wedge y \leq 2$ ")}

Answer (set of conditional mappings)

$$\llbracket q \rrbracket_D = \left\{ \left(\{ ?F \rightarrow \text{fire1} \}, \right. \right. \\ \left. \left. _R1 \text{ NTPP } "x \geq 1 \wedge x \leq 2 \wedge y \geq 1 \wedge y \leq 2" \right) \right\}$$



CONSTRUCT queries

Example

Database D

```
fire1 occuredIn _R1 .
```

```
_R1 NTPP "x ≥ 6 ∧ x ≤ 23 ∧ y ≥ 8 ∧ y ≤ 19"
```

Query q

```
CONSTRUCT { ?F type Fire }  
WHERE {  
    ?F occuredIn ?R  
}
```



CONSTRUCT queries

Example

Database D

```
fire1 occuredIn _R1 .
```

```
_R1 NTPP "x ≥ 6 ∧ x ≤ 23 ∧ y ≥ 8 ∧ y ≤ 19"
```

Query q

```
CONSTRUCT { ?F type Fire }  
WHERE {  
    ?F occuredIn ?R  
}
```

Answer (RDFⁱ database)

$$D' = (G', \phi)$$

```
fire1 type Fire .
```

```
_R1 NTPP "x ≥ 6 ∧ x ≤ 23 ∧ y ≥ 8 ∧ y ≤ 19"
```



CONSTRUCT queries

Example

Database D

```
fire1 occuredIn _R1 .
```

```
_R1 NTPP "x ≥ 6 ∧ x ≤ 23 ∧ y ≥ 8 ∧ y ≤ 19"
```

Query q

```
CONSTRUCT { ?F type Fire }  
WHERE {  
    ?F occuredIn ?R  
}
```

Answer (RDFⁱ database)

$D' = (G', \phi)$

fire1 type Fire .

_R1 NTPP "x ≥ 6 ∧ x ≤ 23 ∧ y ≥ 8 ∧ y ≤ 19"

Closure property



Correctness of SPARQL query evaluation for RDFⁱ

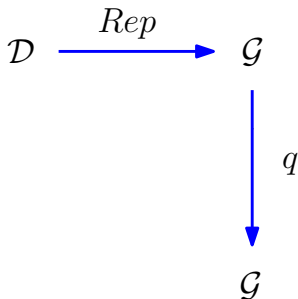
Does query evaluation compute the **correct answer**
(the answer agrees with the semantic definition)?



Correctness of SPARQL query evaluation for RDFⁱ

Does query evaluation compute the **correct answer**
(the answer agrees with the semantic definition)?

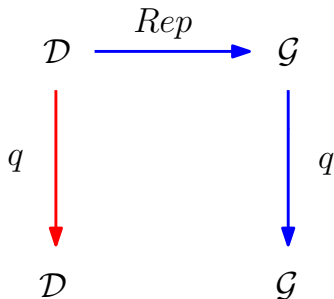
The following diagram should **commute**



Correctness of SPARQL query evaluation for RDFⁱ

Does query evaluation compute the **correct answer**
(the answer agrees with the semantic definition)?

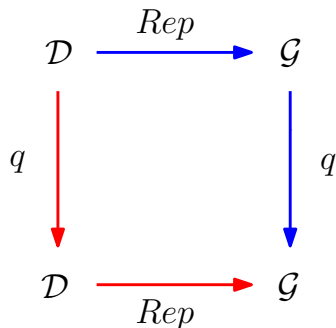
The following diagram should **commute**



Correctness of SPARQL query evaluation for RDFⁱ

Does query evaluation compute the **correct answer**
(the answer agrees with the semantic definition)?

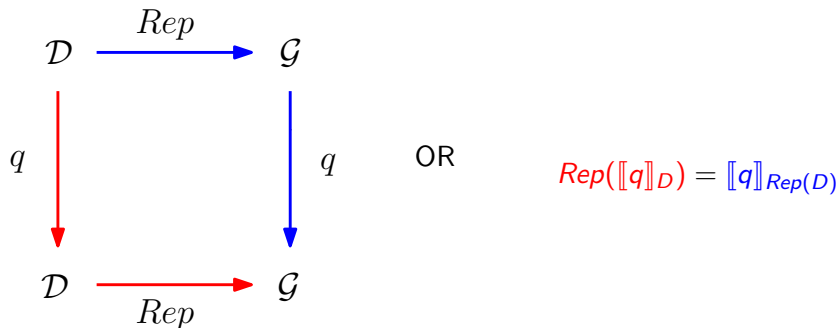
The following diagram should **commute**



Correctness of SPARQL query evaluation for RDFⁱ

Does query evaluation compute the **correct answer**
(the answer agrees with the semantic definition)?

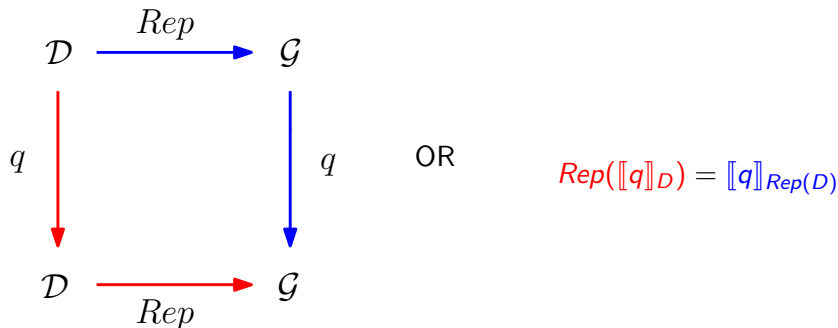
The following diagram should **commute**



Correctness of SPARQL query evaluation for RDFⁱ

Does query evaluation compute the **correct answer**
(the answer agrees with the semantic definition)?

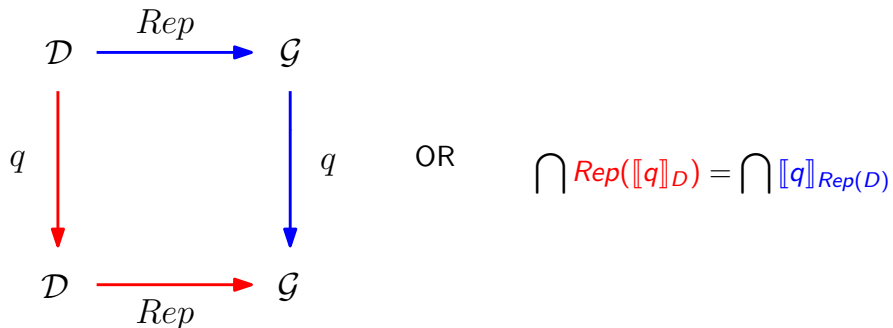
The following diagram should **commute**. **Does it?**



Correctness of SPARQL query evaluation for RDFⁱ

Does query evaluation compute the **correct answer**
(the answer agrees with the semantic definition)?

The following diagram should **commute**. **Does it?**



Certain answer to the rescue

Definition

The **certain answer** to query q over a set of RDF graphs \mathcal{G} is set

$$\bigcap \{ \llbracket q \rrbracket_G \mid G \in \mathcal{G} \}$$



Certain answer to the rescue

Definition

The **certain answer** to query q over a set of RDF graphs \mathcal{G} is set

$$\bigcap \{ [q]_G \mid G \in \mathcal{G} \}$$

Using the notion of certain answer we can **relax the earlier equality requirement** to one that uses **Q-equivalence**.



Certain answer to the rescue

Definition

The **certain answer** to query q over a set of RDF graphs \mathcal{G} is set

$$\bigcap \{ \llbracket q \rrbracket_G \mid G \in \mathcal{G} \}$$

Using the notion of certain answer we can **relax the earlier equality requirement** to one that uses **\mathcal{Q} -equivalence**.

Definition

Let \mathcal{Q} be a fragment of SPARQL. Two sets of RDF graphs \mathcal{G}, \mathcal{H} will be **\mathcal{Q} -equivalent** (denoted by $\mathcal{G} \equiv_{\mathcal{Q}} \mathcal{H}$) if they give the same **certain answer** to every query $q \in \mathcal{Q}$

$$\bigcap \{ \llbracket q \rrbracket_G \mid G \in \mathcal{G} \} = \bigcap \{ \llbracket q \rrbracket_H \mid H \in \mathcal{H} \}$$



Representation system

Let

- ▶ \mathcal{D} be the set of all RDFⁱ databases
- ▶ \mathcal{G} be the set of all RDF graphs
- ▶ $Rep : \mathcal{D} \rightarrow \mathcal{G}$ be a function determining the set of possible RDF graphs corresponding to an RDFⁱ database, and
- ▶ Q be a fragment of SPARQL

$\langle \mathcal{D}, Rep, Q \rangle$ is a **representation system** if for all $D \in \mathcal{D}$ and all $q \in Q$, there exists an RDFⁱ database $\llbracket q \rrbracket_D$ such that

$$Rep(\llbracket q \rrbracket_D) \equiv_Q \llbracket q \rrbracket_{Rep(D)}$$



Representation system

Let

- ▶ \mathcal{D} be the set of all RDFⁱ databases
- ▶ \mathcal{G} be the set of all RDF graphs
- ▶ $Rep : \mathcal{D} \rightarrow \mathcal{G}$ be a function determining the set of possible RDF graphs corresponding to an RDFⁱ database, and
- ▶ \mathcal{Q} be a fragment of SPARQL

$\langle \mathcal{D}, Rep, \mathcal{Q} \rangle$ is a **representation system** if for all $D \in \mathcal{D}$ and all $q \in \mathcal{Q}$, there exists an RDFⁱ database $\llbracket q \rrbracket_D$ such that

$$Rep(\llbracket q \rrbracket_D) \equiv_{\mathcal{Q}} \llbracket q \rrbracket_{Rep(D)}$$

Are there interesting fragments \mathcal{Q} of SPARQL that lead to a representation system?



Representation systems for RDFⁱ

Theorem

The following fragments of SPARQL can give us representation systems for RDFⁱ (with D and Rep as defined):

- ▶ Q_{AUF}^C : CONSTRUCT queries using only *AND*, *UNION*, and *FILTER* graph patterns, and *without blank nodes* in their templates
- ▶ Q_{WD}^C : CONSTRUCT queries using only *well-designed* graph patterns, and *without blank nodes* in their templates

Well-designed graph patterns [Pérez/Arenas/Gutierrez '06]

- ▶ AND, FILTER, OPT fragment
- ▶ P FILTER R : **safe**
- ▶ P_1 OPT P_2 : variables in P_2 are **properly scoped**



Representation systems for RDFⁱ (cont'd)

Monotonicity

Definition

A fragment \mathcal{Q} of SPARQL is **monotone** if for every $q \in \mathcal{Q}$ and RDF graphs G and H such that $G \subseteq H$, it is $\llbracket q \rrbracket_G \subseteq \llbracket q \rrbracket_H$.

Proposition [Arenas/Pérez '11]

- ▶ The fragment of SPARQL corresponding to AND, UNION, and FILTER graph patterns is **monotone**.
- ▶ The fragment of SPARQL corresponding to well-designed graph patterns is **weakly-monotone** (\sqsubseteq).

Proposition

Fragments \mathcal{Q}_{AUF}^C and \mathcal{Q}_{WD}^C are **monotone**.



Computing certain answers

- ▶ Representation systems guarantee correctness of query evaluation for RDFⁱ and SPARQL
- ▶ Query evaluation computes an RDFⁱ database

$$\llbracket q \rrbracket_D = D' = (G', \phi)$$

- ▶ How could we compute the certain answer?

$$\bigcap \text{Rep}(\llbracket q \rrbracket_D)$$

- ▶ $\text{Rep}(\llbracket q \rrbracket_D)$ is **infinite!**



Computing certain answers (cont'd)

Theorem

For $D = (G, \phi)$ and q from \mathcal{Q}_{AUF}^C or \mathcal{Q}_{WD}^C , the *certain answer* of q over D can be computed as follows:

- i) compute $\llbracket q \rrbracket_D = D_q = (G_q, \phi)$,
- ii) compute the RDFⁱ database $(H_q, \phi) = ((D_q)^{EQ})^*$, and
- iii) return the set of *RDF triples*

$$\{(s, p, o) \mid ((s, p, o), \theta) \in H_q \text{ such that } \phi \models \theta \text{ and } o \notin U\}$$



The certainty problem

$CERT(q, H, D)$

Input

An RDF graph H , a CONSTRUCT query q , and an RDFⁱ database D

Question

Does H belong to the certain answer of q over D ?

$$H \subseteq \bigcap \llbracket q \rrbracket_{Rep(D)}?$$



The certainty problem

$$CERT(q, H, D)$$

Input

An RDF graph H , a CONSTRUCT query q , and an RDFⁱ database D

Question

Does H belong to the certain answer of q over D ?

$$H \subseteq \bigcap \llbracket q \rrbracket_{Rep(D)}?$$

We study the **data complexity** of $CERT(q, H, D)$

- ▶ H and D are part of the input
- ▶ q is fixed



Deciding the certainty problem

Theorem

$CERT(q, H, D)$ is *equivalent* to deciding whether formula

$$\bigwedge_{t \in H} (\forall \mathbf{l})(\phi(\mathbf{l}) \supset \Theta(t, q, D, \mathbf{l}))$$

is *true*

- ▶ \mathbf{l} is the vector of all e-literals in D
- ▶ $\Theta(t, q, D, \mathbf{l})$ is of the form $\theta_1 \vee \dots \vee \theta_k$, where θ_i is a conjunction of \mathcal{L} -constraints



Computational complexity

Problem	\mathcal{L}	data complexity
$CERT(q, H, D)$	ECL/diPCL/dePCL/RCL	coNP-complete
	TCL/PCL (RCC-5)	EXPTIME



Computational complexity

Problem	\mathcal{L}	data complexity
$CERT(q, H, D)$	ECL/diPCL/dePCL/RCL	coNP-complete
	TCL/PCL (RCC-5)	EXPTIME

Problem	combined complexity	data complexity
SPARQL	PSPACE-complete	LOGSPACE
SPARQL _{AUF}	NP-complete	
SPARQL _{WD}	coNP-complete	



Conclusions

RDFⁱ framework

- ▶ Modeling of incomplete information for **property values**
- ▶ Formal semantics through **possible worlds** semantics
- ▶ SPARQL query evaluation and **certain answer** semantics
- ▶ Two **representation** systems for RDFⁱ and SPARQL
- ▶ **Algorithm** for certain answer computation
- ▶ Preliminary **complexity analysis**



Future work

- ▶ More general models of incomplete information (subject, predicate)
- ▶ More refined complexity results
- ▶ Scalable implementation when \mathcal{L} expresses topological constraints with/without constants (TCL/PCL)
- ▶ Connection with query processing for the topology vocabulary extension of GeoSPARQL
- ▶ Probabilistic extension to RDFⁱ
- ▶ Data integration theory for linked data (only practice exists so far)
- ▶ Connection to geospatial OBDA using DL logics



Thank you

Constraint languages \mathcal{L}

Properties of \mathcal{L}

- ▶ Many-sorted first-order language
- ▶ Interpreted over a fixed (intended) structure $\mathbf{M}_{\mathcal{L}}$
- ▶ EQ: distinguished equality predicate
- ▶ \mathcal{L} -constraints: quantifier-free formulae of \mathcal{L}
- ▶ Weakly closed under negation: the negation of every atomic \mathcal{L} -constraint is equivalent to a disjunction of \mathcal{L} -constraints

Correctness of SPARQL query evaluation for RDFⁱ

An easy negative example

Example (classical RDF - OWA)

D

s p o .

q

```
CONSTRUCT { s ?p ?o }  
WHERE { s ?p ?o }
```

Correctness of SPARQL query evaluation for RDFⁱ

An easy negative example

Example (classical RDF - OWA)

D

s p o .

q

```
CONSTRUCT { s ?p ?o }  
WHERE { s ?p ?o }
```

Then,

$$\llbracket q \rrbracket_D = D$$

Correctness of SPARQL query evaluation for RDFⁱ (cont'd)

An easy negative example

Example

Let us compare the the set of graphs represented by $\llbracket q \rrbracket_D$ with $\llbracket q \rrbracket_{Rep(D)}$

Correctness of SPARQL query evaluation for RDFⁱ (cont'd)

An easy negative example

Example

Let us compare the the set of graphs represented by $\llbracket q \rrbracket_D$ with $\llbracket q \rrbracket_{Rep(D)}$

$$Rep(\llbracket q \rrbracket_D) = \left\{ \left\{ (s, p, o) \right\}, \left\{ \begin{array}{l} (s, p, o) \\ (c, d, e) \end{array} \right\}, \left\{ \begin{array}{l} (s, p, o) \\ (s, b, c) \end{array} \right\}, \dots \right\}$$

Correctness of SPARQL query evaluation for RDFⁱ (cont'd)

An easy negative example

Example

Let us compare the the set of graphs represented by $\llbracket q \rrbracket_D$ with $\llbracket q \rrbracket_{Rep(D)}$

$$Rep(\llbracket q \rrbracket_D) = \left\{ \left\{ \begin{array}{c} (s, p, o) \\ (c, d, e) \end{array} \right\}, \left\{ \begin{array}{c} (s, p, o) \\ (s, b, c) \end{array} \right\}, \dots \right\}$$

$$\llbracket q \rrbracket_{Rep(D)} = \left\{ \left\{ \begin{array}{c} (s, p, o) \\ (s, b, c) \end{array} \right\}, \dots \right\}$$

Correctness of SPARQL query evaluation for RDFⁱ (cont'd)

An easy negative example

Example

Let us compare the the set of graphs represented by $\llbracket q \rrbracket_D$ with $\llbracket q \rrbracket_{Rep(D)}$

$$Rep(\llbracket q \rrbracket_D) = \left\{ \left\{ (s, p, o) \right\}, \left\{ \begin{matrix} (s, p, o) \\ (c, d, e) \end{matrix} \right\}, \left\{ \begin{matrix} (s, p, o) \\ (s, b, c) \end{matrix} \right\}, \dots \right\}$$

$$\llbracket q \rrbracket_{Rep(D)} = \left\{ \left\{ \begin{matrix} (s, p, o) \\ (s, b, c) \end{matrix} \right\}, \dots \right\}$$

There is no $g \in \llbracket q \rrbracket_{Rep(D)}$ containing the triple (c, d, e) !

Correctness of SPARQL query evaluation for RDFⁱ (cont'd)

An easy negative example

Example

Let us compare the the set of graphs represented by $\llbracket q \rrbracket_D$ with $\llbracket q \rrbracket_{Rep(D)}$

$$Rep(\llbracket q \rrbracket_D) = \left\{ \left\{ (s, p, o) \right\}, \left\{ (s, p, o), (c, d, e) \right\}, \left\{ (s, p, o), (s, b, c) \right\}, \dots \right\}$$

$$\llbracket q \rrbracket_{Rep(D)} = \left\{ \left\{ (s, p, o) \right\}, \left\{ (s, p, o), (s, b, c) \right\}, \dots \right\}$$

There is no $g \in \llbracket q \rrbracket_{Rep(D)}$ containing the triple (c, d, e) !

- This would work if RDF made the CWA

Correctness of SPARQL query evaluation for RDFⁱ (cont'd)

An easy negative example

Example

Let us compare the the set of graphs represented by $\llbracket q \rrbracket_D$ with $\llbracket q \rrbracket_{Rep(D)}$

$$Rep(\llbracket q \rrbracket_D) = \left\{ \left\{ \begin{array}{l} (s, p, o) \\ (c, d, e) \end{array} \right\}, \left\{ \begin{array}{l} (s, p, o) \\ (s, b, c) \end{array} \right\}, \dots \right\}$$

$$\llbracket q \rrbracket_{Rep(D)} = \left\{ \left\{ \begin{array}{l} (s, p, o) \\ (s, b, c) \end{array} \right\}, \dots \right\}$$

There is no $g \in \llbracket q \rrbracket_{Rep(D)}$ containing the triple (c, d, e) !

- ▶ This would work if RDF made the CWA
- ▶ We know this already from the relational case [Imielinski/Lipski '84]

Computing certain answers

Definitions

Definition (EQ-completion)

The EQ-completed form of $D = (G, \phi)$, denoted by $D^{EQ} = (G^{EQ}, \phi)$, is taken from D by replacing all e-literals $_l \in U$ appearing in G by the constant $c \in C$ such that $\phi \models _l \text{EQ } c$

Computing certain answers

Definitions

Definition (EQ-completion)

The EQ-completed form of $D = (G, \phi)$, denoted by $D^{EQ} = (G^{EQ}, \phi)$, is taken from D by replacing all e-literals $l \in U$ appearing in G by the constant $c \in C$ such that $\phi \models l \text{ EQ } c$

Definition (Normalization)

The normalized form of D is the RDFⁱ database $D^* = (G^*, \phi)$ where G^* is the set

$$\{(t, \theta) \mid (t, \theta_i) \in G \text{ for all } i = 1 \dots n, \text{ and } \theta \text{ is } \bigvee_i \theta_i\}$$

$$G = \{(t, \theta_1), (t, \theta_2), (t', \theta')\}$$

$$G^* = \{(t, \theta_1 \vee \theta_2), (t', \theta')\}$$